Volume 4 (2025) ISSN 2817-1705 Editors-in-Chief: Khaled El Emam, Bradley Malin, PhD

Contents

Reviews

Identification of Use Cases, Target Groups, and Motivations Around Adopting Smart Speakers for Health Care and Social Care Settings: Scoping Review (e55673)	
Sebastian Merkel, Sabrina Schorr.	4
Survey on Pain Detection Using Machine Learning Models: Narrative Review (e53026)	
Ruijie Fang, Elahe Hosseini, Ruoyu Zhang, Chongzhou Fang, Setareh Rafatirad, Houman Homayoun	18
Studying the Potential Effects of Artificial Intelligence on Physician Autonomy: Scoping Review (e59295)	
John Grosser, Juliane Düvel, Lena Hasemann, Emilia Schneider, Wolfgang Greiner	50
Disease Prediction Using Machine Learning on Smartphone-Based Eye, Skin, and Voice Data: Scoping Review (e59094)	
Research Dawadi, Mai Inoue, Jie Tay, Agustin Martin-Morales, Thien Vu, Michihiro Araki.	330

Viewpoint

Striking a Balance: Innovation, Equity, and Consistency in AI Health Technologies (e57421)	
Eric Perakslis, Kimberly Nolen, Ethan Fricklas, Tracy Tubb	63

Research Letter

Evaluating ChatGPT's Efficacy in Pediatric Pneumonia Detection From Chest X-Rays: Comparative Analysis	
of Specialized AI Models (e67621)	
Nitin Chetla, Mihir Tandon, Joseph Chang, Kunal Sukhija, Romil Patel, Ramon Sanchez.	74

Corrigenda and Addenda

JMIR AI 2025 | vol. 4 | p.1

XSL•FO RenderX

Original Papers

Estimation of Static Lung Volumes and Capacities From Spirometry Using Machine Learning: Algorithm Development and Validation (e65456)	
Scott Helgeson, Zachary Quicksall, Patrick Johnson, Kaiser Lim, Rickey Carter, Augustine Lee.	79
Limitations of Binary Classification for Long-Horizon Diagnosis Prediction and Advantages of a Discrete-Time Time-to-Event Approach: Empirical Analysis (e62985)	
De Loh, Elliot Hill, Nan Liu, Geraldine Dawson, Matthew Engelhard.	92
Using AI to Translate and Simplify Spanish Orthopedic Medical Text: Instrument Validation Study (e70222) Saman Andalib, Aidin Spina, Bryce Picton, Sean Solomon, John Scolaro, Ariana Nelson.	108
Deep Learning Models to Predict Diagnostic and Billing Codes Following Visits to a Family Medicine Practice: Development and Validation Study (e64279)	
Akshay Rajaram, Michael Judd, David Barber.	120
Improving the Robustness and Clinical Applicability of Automatic Respiratory Sound Classification Using Deep Learning–Based Audio Enhancement: Algorithm Development and Validation (e67239)	
Jing-Tong Tzeng, Jeng-Lin Li, Huan-Yu Chen, Chu-Hsiang Huang, Chi-Hsin Chen, Cheng-Yi Fan, Edward Huang, Chi-Chun Lee	134
A Conversational Agent Using Natural Language Processing for Postpartum Care for New Mothers: Development and Engagement Analysis (e58454)	
Kirstin Leitner, Clare Cutri-French, Abigail Mandel, Lori Christ, Nathaneal Koelper, Meaghan McCabe, Emily Seltzer, Laura Scalise, James Colbert, Anuja Dokras, Roy Rosin, Lisa Levine.	154
The Diagnostic Performance of Large Language Models and Oral Medicine Consultants for Identifying Oral Lesions in Text-Based Clinical Scenarios: Prospective Comparative Study (e70566)	
Sarah AlFarabi Ali, Hebah AlDehlawi, Ahoud Jazzar, Heba Ashi, Nihal Esam Abuzinadah, Mohammad AlOtaibi, Abdulrahman Algarni, Hazzaa Alqahtani, Sara Akeel, Soulafa Almazrooa.	164
Comparison of Deep Learning Approaches Using Chest Radiographs for Predicting Clinical Deterioration: Retrospective Observational Study (e67144)	
Mahmudur Rahman, Jifan Gao, Kyle Carey, Dana Edelson, Askar Afshar, John Garrett, Guanhua Chen, Majid Afshar, Matthew Churpek 1 7 1	
Large Language Models for Thematic Summarization in Qualitative Health Care Research: Comparative Analysis of Model and Human Performance (e64447)	
Arturo Castellanos, Haoqiang Jiang, Paulo Gomes, Debra Vander Meer, Alfred Castillo	185
Evaluation of ChatGPT Performance on Emergency Medicine Board Examination Questions: Observational Study (e67696)	
Mila Pastrak, Sten Kajitani, Anthony Goodings, Austin Drewek, Andrew LaFree, Adrian Murphy	197
Enhancing Interpretable, Transparent, and Unobtrusive Detection of Acute Marijuana Intoxication in Natural Environments: Harnessing Smart Devices and Explainable AI to Empower Just-In-Time Adaptive Interventions: Longitudinal Observational Study (e52270)	
Sang Bae, Tammy Chung, Tongze Zhang, Anind Dey, Rahul Islam.	205
Investigating the Classification of Living Kidney Donation Experiences on Reddit and Understanding the Sensitivity of ChatGPT to Prompt Engineering: Content Analysis (e57319)	
Joshua Nielsen, Xiaoyu Chen, LaShara Davis, Amy Waterman, Monica Gentili	225

XSL•FO RenderX

Urgency Prediction for Medical Laboratory Tests Through Optimal Sparse Decision Tree: Case Study With Echocardiograms (e64188)	242
Advancing Privacy-Preserving Health Care Analytics and Implementation of the Personal Health Train:	
Ananya Choudhury, Leroy Volmer, Frank Martin, Rianne Fijten, Leonard Wee, Andre Dekker, Johan Soest.	258
Predicting Satisfaction With Chat-Counseling at a 24/7 Chat Hotline for the Youth: Natural Language Processing Study (e63701)	
Silvan Hornstein, Ulrike Lueken, Richard Wundrack, Kevin Hilbert	274
Leveraging Medical Knowledge Graphs Into Large Language Models for Diagnosis Prediction: Design and Application Study (e58670)	
Yanjun Gao, Ruizhe Li, Emma Croxford, John Caskey, Brian Patterson, Matthew Churpek, Timothy Miller, Dmitriy Dligach, Majid Afshar 2	
GPT-4 as a Clinical Decision Support Tool in Ischemic Stroke Management: Evaluation Study (e60391)	
Amit Shmilovitch, Mark Katson, Michal Cohen-Shelly, Shlomi Peretz, Dvir Aran, Shahar Shelly.	304
Creation of Scientific Response Documents for Addressing Product Medical Information Inquiries: Mixed Method Approach Using Artificial Intelligence (e55277)	
Jerry Lau, Shivani Bisht, Robert Horton, Annamaria Crisan, John Jones, Sandeep Gantotti, Evelyn Hermes-DeSantis.	315
Prompt Engineering an Informational Chatbot for Education on Mental Health Using a Multiagent Approach for Enhanced Compliance With Prompt Instructions: Algorithm Development and Validation (e69820)	
Per Waaler, Musarrat Hussain, Igor Molchanov, Lars Bongo, Brita Elvevåg	351
Generative Large Language Model—Powered Conversational AI App for Personalized Risk Assessment: Case Study in COVID-19 (e67363)	
Mohammad Roshani, Xiangyu Zhou, Yao Qiang, Srinivasan Suresh, Steven Hicks, Usha Sethuraman, Dongxiao Zhu.	368
Insights on the Side Effects of Female Contraceptive Products From Online Drug Reviews: Natural Language Processing–Based Content Analysis (e68809)	
Nicole Groene, Audrey Nickel, Amanda Rohn.	385
Trust, Anxious Attachment, and Conversational AI Adoption Intentions in Digital Counseling: A Preliminary Cross-Sectional Questionnaire Study (e68960)	
Xiaoli Wu, Kongmeng Liew, Martin Dorahy.	403
Harnessing Moderate-Sized Language Models for Reliable Patient Data Deidentification in Emergency Department Records: Algorithm Development, Validation, and Implementation Study (e57828)	
Océane Dorémus, Dylan Russon, Benjamin Contrand, Ariel Guerra-Adames, Marta Avalos-Fernandez, Cédric Gil-Jardiné, Emmanuel Lagarde.	
Utility-based Analysis of Statistical Approaches and Deep Learning Models for Synthetic Data Generation With Focus on Correlation Structures: Algorithm Development and Validation (e65729)	
Marko Miletic, Murat Sariyar.	434

XSL•FO-RenderX JMIR AI 2025 | vol. 4 | p.3

Identification of Use Cases, Target Groups, and Motivations Around Adopting Smart Speakers for Health Care and Social Care Settings: Scoping Review

Sebastian Merkel^{1*}, PhD; Sabrina Schorr^{1*}, MA

Faculty of Social Science, Ruhr University Bochum, Bochum, Germany ^{*}all authors contributed equally

Corresponding Author: Sebastian Merkel, PhD Faculty of Social Science Ruhr University Bochum GD E1/ 155 Universitätsstraße 150 Bochum, 44801 Germany Phone: 49 0234 32 25411 Email: sebastian.merkel@ruhr-uni-bochum.de

Abstract

Background: Conversational agents (CAs) are finding increasing application in health and social care, not least due to their growing use in the home. Recent developments in artificial intelligence, machine learning, and natural language processing have enabled a variety of new uses for CAs. One type of CA that has received increasing attention recently is smart speakers.

Objective: The aim of our study was to identify the use cases, user groups, and settings of smart speakers in health and social care. We also wanted to identify the key motivations for developers and designers to use this particular type of technology.

Methods: We conducted a scoping review to provide an overview of the literature on smart speakers in health and social care. The literature search was conducted between February 2023 and March 2023 and included 3 databases (PubMed, Scopus, and Sociological Abstracts), supplemented by Google Scholar. Several keywords were used, including technology (eg, voice assistant), product name (eg, Amazon Alexa), and setting (health care or social care). Publications were included if they met the predefined inclusion criteria: (1) published after 2015 and (2) used a smart speaker in a health care or social care setting. Publications were excluded if they met one of the following criteria: (1) did not report on the specific devices used, (2) did not focus specifically on smart speakers, (3) were systematic reviews and other forms of literature-based publications, and (4) were not published in English. Two reviewers collected, reviewed, abstracted, and analyzed the data using qualitative content analysis.

Results: A total of 27 articles were included in the final review. These articles covered a wide range of use cases in different settings, such as private homes, hospitals, long-term care facilities, and outpatient services. The main target group was patients, especially older users, followed by doctors and other medical staff members.

Conclusions: The results show that smart speakers have diverse applications in health and social care, addressing different contexts and audiences. Their affordability and easy-to-use interfaces make them attractive to various stakeholders. It seems likely that, due to technical advances in artificial intelligence and the market power of the companies behind the devices, there will be more use cases for smart speakers in the near future.

(JMIR AI 2025;4:e55673) doi:10.2196/55673

KEYWORDS

conversational agents; smart speaker; health care; social care; digitalization; scoping review; mobile phone

Introduction

Background

In the context of ongoing public debates on artificial intelligence (AI), dialogue systems or conversational agents (CAs) are receiving increasing attention. Their potential applications are being discussed in various fields, including health care [1,2] and social care [3]. CAs have been used in both fields for several years, but recent developments in AI have fueled the scientific discourse [4,5]. The developments in the field of machine learning and natural language processing (NLP), as well as the success of commercially available CAs, such as Amazon's Alexa or Apple's Siri, have been particularly decisive in this regard.

The use of CAs is not limited to a single context; rather, they are used in a variety of settings, including those pertaining to the acquisition of information related to health [6]. CAs using NLP offer a number of features that can be implemented in a variety of health care and social care settings. The field of AI has witnessed considerable progress in recent years, with speech recognition (SR) and NLP advancing significantly. This has enabled the processing of medical terminology in various settings [7]. Although SR in health care has a long tradition dating back to the 1980s, when initial attempts were made to dictate doctor's letters [8], CAs offer multiple additional features. In the context of hands-free interaction, CAs have been used for the purposes of medication reminders [9], symptom management [10], documentation [11], or communication between patients and nurses or doctors, covering multiple medical fields. These include diabetes care [12], monitoring of pregnant women [13], children with special health care needs [11], hearing tests [14], cardiovascular disease [15], and the support of persons with dementia, to name a few [16].

The Rise of Smart Speakers

The term "CA" is not clearly defined, and within the literature, multiple synonyms are used interchangeably. These include "virtual assistants," "AI-driven digital assistants," "voice-based assistants," "voice-controlled intelligent personal assistants," and others. In the study by Laranjo et al [1], the term "CA" is defined as encompassing a range of technologies, including chatbots, embodied CA, which involves a computer-generated character such as an avatar, and smart conversational interfaces, such as Apple's Siri or Amazon's Alexa. In order to characterize CAs, the authors propose that it is necessary to differentiate between the type of technology in question (eg, if the software application is delivered through a mobile device or the telephone), the type of dialogue management (finite-state, frame-based, or agent-based), the actors with control over the dialogue initiative (the user, the system, or a combination of both), the input or output modality (spoken or written, or visual in the case of the output), and whether the system is task-oriented or not [1].

This paper is particularly interested in the use of CAs that are embodied in a physical stationary artifact, which is referred to as a smart speaker. Examples of such devices include Amazon's Echo and Apple's HomePod. Smart speakers are typically confined to a specific location and serve as a platform for a

```
https://ai.jmir.org/2025/1/e55673
```

smart conversational interface or AI-driven digital assistant that can be operated through voice input. In the case of the Echo, this is "Alexa", while in the HomePod, it is "Siri". Such assistants are capable of fulfilling a range of tasks, including answering simple questions, switching on lights in conjunction with a smart home system, and playing music. The devices are equipped with one or multiple microphones and software that is capable of analyzing and generating spoken language. In order to operate the devices, the user must utter a designated wake word, such as "Alexa" or "Computer" in the case of Amazon's Echo [17].

The diffusion of smart speakers has been observed to be high in private households in Europe and North America. Amazon launched the first smart speaker in the United States in 2015. As of 2022, approximately 35% of the total US population had used smart speakers [18]. In comparison to the figures from 2019, this represents an increase of 11.1% [19]. A number of studies conducted by market research companies in other countries have reached similar conclusions. For instance, these studies have found that 33% of internet households in the United States, 34% in the United Kingdom [20], and approximately 12%-33% of all households in Germany own at least one smart speaker [21,22].

A recent study by Gaspar and Neus [23] of smart speaker users in the United States, United Kingdom, and Germany shows that Amazon is still the current market leader (United States: 58%; United Kingdom: 71%; and Germany: 68%) followed by Google (United States: 34%; United Kingdom: 22%; and Germany: 25%) and other brands (United States, United Kingdom, and Germany: 7%). It was also found that in all countries, at least 40% (United States: 46%; United Kingdom: 40%; and Germany: 44%) of respondents use smart speakers several times a day. Participants were also asked about the attractiveness of certain application scenarios, including medical diagnosis. Here, participants gave high ratings: United States (19% very attractive and 36% attractive), United Kingdom (12% very attractive and 34% attractive), and Germany (13% very attractive and 35% attractive).

In light of the commercial success of smart speakers and the aforementioned technological advantages in SR and NLP, there has been a growing body of literature on smart speakers in different health care and social care settings [1,24-27]. Commercial devices, such as Amazon's Echo, offer a multitude of features. These devices can be used without any direct contact, are relatively inexpensive and easy to operate, and can be customized and personalized by installing new applications and features [28]. These factors have played a pivotal role in the dissemination of the technology. Finally, the widespread adoption of the technology was driven by the pandemic and the subsequent shift in clinical practices toward greater reliance on digital technologies [29]. Nevertheless, the pervasive use of these devices has also given rise to a multitude of issues and concerns, most notably data collection, storage, and protection [8].

Hence, the devices have attracted increasing attention, with several reviews on CAs in health care settings having been published recently. Each of these reviews has a specific focus:

XSL•FO RenderX

these include, for instance, design and evaluation challenges [30], effectiveness and usability [31], or chronic conditions [32,33]. To the best of our knowledge, no review has been conducted to date that specifically examines the use of smart speakers within health care and social care settings.

As evidenced by the current state of research, smart speakers are becoming increasingly prevalent in the field of health care and social care. However, there is currently no systematic review available that specifically investigates use cases, settings in which the devices are used, or target groups. To address this gap, our main research question is as follows: What are the scenarios of the use of smart speakers in health care and social care? To address this research question, the main aim of this paper is to present a review of the current research on the use of smart speakers in health care.

Methods

Overview

In order to provide an overview of the existing literature on smart speakers in health care and social care, we conducted a scoping review. The main aim of this approach is to observe, synthesize, and understand current trends [34]. In contrast to a systematic review, which is more suitable for the presentation of a specific clinical question or the presentation of evidence for practice, a scoping review is particularly suitable for identifying features and concepts. Furthermore, it does not aim to provide a synthesizing result for a specific question but rather to provide an overview of a specific topic [34,35]. Thus, the scoping review is a particularly suitable instrument for analyzing the research interest. This encompasses the identification of the nature of the literature, the collation of information on key topics, and the identification of knowledge gaps [35]. Its methodological framework was first published by Arksey and

O'Malley [36] and later adapted by Levac, Colquhoun, and O'Brien [37]. Contrary to a systematic review, search terms can be adjusted along the process of a scoping review [36,38]. For the conduction of the present review, the guidelines of Peters et al [39], the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) [40] and its extension for Scoping Reviews (PRISMA-ScR) [41] were followed. The results were presented according to the PRISMA checklist (Multimedia Appendix 1).

Search Strategy and Selection Criteria

The literature search was conducted between February 2023 and March 2023. This included a systematic literature search of 3 databases (PubMed, Scopus, and Sociological Abstracts) and a cross-search of the first 20 pages of Google Scholar. This was supplemented by tracing reference lists for further relevant studies. We used the program Citavi 6 for literature management. The review protocol is available on request from the authors. The following keywords were applied in varying combinations and spellings for the systematic search (Table 1):

- 1. Technology: Here, several terms described above that are found in the literature on CA were used. As the focus of this review is on smart speakers, the search was restricted to this specific type of CA.
- 2. Product name: As smart speakers were introduced to the market by major American information technology companies, which often use the product names as synonyms for the product, we also included the product or brand names in our search. Globally, Amazon, Google, and Apple are the 3 leading manufacturers; therefore, we included the names of their brands in our search [42].
- 3. Setting: In order to ensure the most comprehensive search results, we elected to limit our search to the 2 domains of health care and social care without imposing any further restrictions.

Table 1. Keywords used in the literature review.

Technology	Vendor, brand, and product	Setting
Smart speaker	Amazon Alexa	Health care
Voice assistant	Amazon Echo	Social care
Voice-based assistant	Apple HomePod	Care
Voice-controlled assistant	Apple Siri	Nursing
Artificial intelligence-driven digital assistant	Google Home	a
Conversational agent	Google Nest	_
Virtual assistant	_	_

^aNot applicable.

The terms were linked using Boolean operators. Multiple combinations of the search terms were used using different operators (Multimedia Appendix 2).

To select studies relevant to our research interest, we defined the following inclusion criteria for the full-text screening: (1) publications that were released after 2015, as this was the year in which the first commercial smart speaker was introduced to the market, and (2) the use of a smart speaker in health care and

https://ai.jmir.org/2025/1/e55673

social care settings. No restrictions were placed on the specific setting, including hospitals or long-term care facilities. Furthermore, articles were included in which the devices were not implemented in real settings but were developed for specific settings. Studies were excluded if they met one of the following exclusion criteria: (1) papers that do not report on the specific devices that were used (for instance, in some cases, the authors described the use of a personal assistant without explicitly indicating the specific device on which the assistant was

operational), (2) studies that did not specifically focus on smart speakers (this encompasses the development of voice-operated applications for use on smartphones or tablets), (3) systematic reviews and other forms of literature-based publications, and (4) articles not published in the English language.

Process of Study Selection and Data Extraction

We first screened the titles and abstracts for relevance by both authors. No exclusion criteria were applied to the type of publication during the title and abstract search. Should the title or abstract screening indicate the use of a smart speaker in a health care or social care context, the articles were deemed eligible for full-text screening. For the title and abstract screening, as well as the full-text screening, the same 2 authors reviewed each article independently in order to decide on its inclusion or exclusion. In the event of conflicting decisions regarding inclusion or exclusion, the authors attempted to reach a consensus through discussion. As there was no disagreement, there was no need to involve a third party. The data extraction table contains the following information about each article: (1) authors, (2) year of publication, and (3) country of publication. Furthermore, data were collected on the product and the use case. Furthermore, the following aspects were considered: the settings, the target groups, the motivation for using smart speakers, and the limitations of using such a device. As the primary focus was not on methodological aspects, and due to the heterogeneity of the included literature (some described only technical development while others also included user testing and the often-limited reporting of methods), no such information was collected. The articles included were subjected to qualitative thematic analysis in accordance with the

methodology outlined in [43]. Using Kuckartz's [43] approach to qualitative thematic text analysis, researchers identify codes through analysis based on the data gathered. During the process, these codes are then refined. Researchers then identify themes or categories that represent the main findings of the analysis. Identifying themes is a process of examining patterns and similarities between codes and then relating the themes to each other. Consequently, all papers included were read and re-read by both authors, with initial codes being identified. The codes were then compared by the authors, discussed, and grouped into themes. In particular, this included an analysis of the motivation for using the devices and the limitations encountered during the research and development process.

Ethical Considerations

Given the nature of the study, there were no direct interactions with human participants, and thus, no participants to recruit or consent, and no institutional ethical approval was required.

Results

Overview

In total, our search yielded 1975 articles. After removing 316 duplicates, 1659 titles and abstracts were screened by the 2 reviewers. The screening of titles and abstracts resulted in the exclusion of 1571 records, leaving 88 full texts to be assessed for eligibility. Of these, 61 articles were excluded, resulting in a final pool of 27 articles for analysis (Figure 1). The data extraction table for the articles included can be found in Multimedia Appendix 3 [3,9,13-15,44-65].

Figure 1. PRISMA (Preferred Reporting Items for Systematic Reviews and Meta-Analyses) flowchart of the search process.



Year and Country of Publication

The majority of articles included in the analysis were published in the United States (n=15 [9, 15, 39, 45, 46, 49, 51, 53-57, 61, 63, 65]), followed by the United Kingdom (n=4 [3,13,43,62]), North Macedonia (n=2 [58,59]), and Australia (n=2 [52,64]). All articles were published between 2018 and 2022, with 2021 being the year with the highest number of publications, with 11 articles.

Technology

There was a clear preference for the devices used: Amazon products were used in 23 of the articles, followed by Google (5). A total of 3 papers used a prototype. It should be noted that some articles used devices from several companies. We found 2 types of articles: Those that use the devices, including the infrastructure (eg, frameworks) provided by the developers, and those that mainly use the hardware (eg, for heart rhythm monitoring; Multimedia Appendix 3 [3,9,13-15,44-65]).

The devices were found to be used In 3 main ways: (1) as standard smart speakers without any further modification, for example, to communicate with patients or to support people living alone (for instance, [44,47]); (2) to develop a skill for a specific use case or multiple use cases (for instance, [48]); and (3) to use the smart speaker and, in some cases, the skill to feed information into another system or as a communication device for other systems (for instance, [15]).

Settings and Target Groups

Given the diverse range of health care and social care settings, we have defined the following categories (Textbox 1). It should be noted that not all articles reported the testing of smart speakers in real health care and social care settings. In some cases, applications were tested in laboratory environments. In the event that this was the case, the intended setting was coded.

Textbox 1. We used the following settings within the domains of health care and special care.

Private homes

• The private living environment includes a person's own home.

Hospitals

• This setting covers acute care hospitals as well as urgent care centers.

Long-term care facilities

• This category includes all settings in which long-term care is provided, for example, nursing homes or rehabilitation centers.

Outpatient services

• This category covers specialized outpatient services, for example, dental or pain management clinics.

Other

• In case the device was tested in a setting not matching the definition of the ones listed above, we categorized it as "other." For instance, this could be in a car.

Furthermore, 4 target groups were identified. It should be noted that an article can have several target groups, including (1) patients, (2) medical staff members such as physicians, (3) nurses and professional caregivers, and (4) informal caregivers who provide unpaid help to a friend or family member. Moreover, category (5), "other," was defined for all target groups not matching any of the aforementioned. It should be noted that multiple target groups were covered in one article. Only those who directly interact with the device were included. For instance, Domínguez et al [50] developed a system to support assisted reproduction treatment. Although physicians are involved, only the patients interact with a smart speaker and hence were included.

The most prevalent setting mentioned in the studies included was home care (n=20), followed by hospitals (n=6). Outpatient care (n=3) was less frequently observed (Multimedia Appendix 3 [3,9,13-15,44-65]). In one instance, the setting was not specified [14]. However, it is best classified under home care.

Among the target groups, patients are the most frequent users mentioned in 23 of the articles (Multimedia Appendix 3

```
https://ai.jmir.org/2025/1/e55673
```

[3,9,13-15,44-65]). Older adults, in particular, were often seen as a promising target group, and we found that 11 of the included publications focus on this target group [66] (Multimedia Appendix 3 [3,9,13-15,44-65]). While some articles included descriptions of the development and testing of skills specifically designed for older adults [51,52], others explored the general acceptance and potential of the technology for older adults. For instance, Lee et al [51] developed multiple skills aimed at older persons, including a reminder to take medication, a diet tracking system, and a skill alerting caregivers in case of a fall. Nallam et al [49] simulated a CA to answer health-related questions asked by older persons. O'Brien et al [47] used off-the-shelf devices without any form of modification to investigate the effects on home-bound older adults with social isolation. The participants used the devices for a variety of purposes, including monitoring their health and well-being, as well as for emergency communication. Some authors report that older adults constitute the largest group of first adopters of smart speakers. In addition, smart speakers allow easy contact with caregivers [12] or low-threshold access to health information [13]. Older adults as potential users of CA have been the focus before [39,67,68].

The second most frequent target group was physicians (n=11), followed by other health professionals (eg, nurses; n=9) and informal caregivers (n=1; Multimedia Appendix 3 [3,9,13-15,44-65]). These results demonstrate that the majority of articles focus on supporting nonresidential care.

Table 2 provides an overview of all settings and target groups. It is important to note that a single paper can include multiple settings and target groups.

Table 2. Settings and user groups.

While previous statements covered the number of papers included in the review, Table 2 combines user groups and settings across the studies covered. It shows that patients are the most common target group, while home care is the most common setting.

	Patients	Physicians	Older adults	Nurses and so on	Informal caregivers	Other	Total
Home care	19	5	11	7	1	1	44
Hospitals	4	5	0	2	0	0	11
Outpatient care	2	2	1	1	0	0	6
Total	25	12	12	10	1	1	

Use Cases

We found several use cases covering, among others, hearing tests [14], cardiovascular diseases [15,46], pregnancy companion [13], cancer management [eg, 58,59], or medication reminders [69]. It must be noted that several articles reported that smart speakers were used in multiple use cases. For example, Wright [70] describes that a local authority was involved in developing applications, including "a Skill that prompted users to take their medicine; a Skill that helped to record and manage care tasks; a Skill to facilitate communication with caregivers by recording messages; and a Skill to connect users to a trusted LA directory of services" [44]. Jadczyk et al [71], who developed a voice-enabled automated platform for the collection of medical data from patients with cardiovascular disease, describe 5 use cases within their study: (1) education, (2) process optimization, (3) patient support, and (4) data collection, and (5) medical device grade solutions (eg, diagnose and treatment). The devices were used to open patient files and images, initiate conference calls, or record images and videos [4].

While most of the identified use cases were found in the domain of health care, social care played a subordinate role. Still, we found several articles reporting on the use of smart speakers in this domain. Within this field, elderly care was the most relevant area. For instance, O'Brien et al [47] use a smart speaker to reduce loneliness and social isolation among older adults living at home. Palumbo et al [72] developed personalized coaching for older individuals to increase their well-being by aiming at the areas of physical activity, nutrition, cognition, and social relationships. In the domain of social care, older adults living at home or care home residents were the main user group (eg, [3]).

Motivation for Use

The reasons for using smart speakers in health care are framed with various arguments. Besides their low acquisition costs [51], this also includes aspects applying to digital technologies in health care and social care in general, such as the possibility to deliver care remotely without restrictions in time and space (eg, Sadavarte et al [13]). Another motivation is the fact that smart speakers are already widely accepted as a consumer

https://ai.jmir.org/2025/1/e55673

technology [45,52]. Hence, users already know how to operate the devices and are also familiar with their limitations. Other aspects cover potentially increased productivity across the use cases that we identified. For instance, Bhatt et al [45] used a voice-based assistant to access and update an electronic health record. They see advantages in terms of efficiency (less time spent on data input) and accuracy, as speech-to-text might result in fewer errors. Ultimately, this might also benefit patients as waiting time is reduced [45]. Jadczyk et al [71] highlighted the main potential in the possibility of automating traditional telehealth services: "Voice chatbots can support routine care through automatic at-home monitoring, triaging, screening, providing medical recommendations and guidelines, and improving operational workflow" [15,71].

Another advantage is the user interface, which is easy to navigate [11]. Cheng et al [55] argue that the main advantage of the technology is that it: "eliminate[s] the struggles that are associated with strictly tactile screens." (2018); or that human-like verbal communication that feels more natural and intuitive and particularly that the devices can be used hands-free [55]. Jansons et al [52] drive on the research of Foehr and Germelmann [73] and argue that the devices "may enhance adherence to remotely-delivered exercise interventions [...], because the human-like attributes associated with these technologies may elicit a sense of familiarity, social presence, and human engagement" [52]. Moreover, the authors see this as an advantage for older users [53] who support this viewpoint and argue that "digital non-natives" might be especially benefitting from this technology. For instance, Kim [4] tested the experiences of older adults who used the devices for the first time and found that due to the simple interaction, health-related questions were a typical use case.

The form of smart speakers and their design were mentioned in some publications. Gouda et al [74] saw the fact that smart speakers are "non-invasive" technology as a main advantage. As the devices can be placed nearly anywhere in the room and can be operated without the need to see them, it allows for new ways of interaction. Luo et al [56] also see a benefit in the fact that the immobility of the devices is as helpful as this helps, in contrast to mobile phones, in establishing habits and routines.

Wright [44] describes the use of smart speakers in trials run by local authorities in England. Drawing on interviews with managers from 8 English local authorities, benefits are seen in the low-cost supplement or alternative to telecare. Or, as one of his interview partners put it: "have the advantages of being sophisticated and powerful, relatively cheap, already widely used and familiar, designed with a degree of accessibility and intuitive use in mind, and a growing level of interoperability with other networked digital devices aided by an open development framework" [44]. One of the results of the study is that local authorities chose Amazon's Echo because of "councils facing depleted funds, a lack of expert guidance on care technologies, and an increasingly complex and fragmented care technology marketplace" [44].

Limitations of Smart Speakers

In addition, various limitations of the technology were addressed in the included articles. Here, most technical limitations were named (1) insufficient hearing comprehension [57], speech recognition [51], or emotion recognition [54]; (2) that there is no interruption of the recording during slow speeches allowed [14]; (3) difficult functioning in the natural living environments due to interfering noises [3]; (4) that the correctness of the answer is not always accurate [51]; and (5) that the devices allow longer conversations [49]. Internet access must also be provided [48,75]. Besides these technical aspects, there were also social aspects mentioned. This covered the (lack of) user acceptance, particularly among older users and professional caregivers [45,76], but also their lack of basic digital skills [75]. These supposedly low digital skills might lead to challenges in interacting with the devices. Users might forget the wake word, there may be timing issues when communicating with the devices, or they might have difficulties in setting up the devices [47,53]. Another issue that was mentioned regularly was data protection. Here, the misuse of sensitive data is particularly pointed out. For example, if security measures are inadequate, it would be possible to manipulate the medication and thus actively harm the patient [12]. Cheng et al [55] also argue for multimodal solutions as people might feel uncomfortable talking to devices in front of other people.

Discussion

Principal Findings

Our aim was to identify use cases and scenarios in which smart speakers can be used within health care and social care. The results show that smart speakers are used in various contexts and for multiple reasons. The main features used are NLP and hands-free interaction. Moreover, the fact that the technology is widely used in private homes and hence many persons are used to interact with the devices are important aspects. In addition to offering relatively inexpensive hardware, smart speakers and the companies behind them provide software frameworks and infrastructure, such as Amazon's skill, which assists developers in the design and marketing of their products.

It is important to note that there is no clear definition of smart speakers. One challenge of this study was the varying definitions of the technology, with the term often being used interchangeably with personal assistants such as Siri or Cortana.

```
https://ai.jmir.org/2025/1/e55673
```

These assistants play an important role in the use of smart speakers, which arguably only serve as a shell equipped with microphones and loudspeakers for them. However, we argue that smart speakers should be considered a distinct technology. Based on this review, we understand smart speakers as a type of CA bound to a fixed location. Within the field of health care and social care, the technology can be used in various settings and use cases such as communication, documentation, or diagnosis and therapy of diseases hands-free. Smart speakers are equipped with microphones and loudspeakers and connected to the internet. They usually come with an integrated digital assistant, but even without such an assistant, they offer multiple features that can be used across various settings. Smart speakers can be customized using either skills or apps that can be installed on the devices.

The results show that all publications were published between 2018 and 2021. Furthermore, the majority were published in the United States. The following explanations can be given for these 2 results. Alexa was the first voice assistant that was compliant with the Health Insurance Portability and Accountability Act (HIPAA), allowing it to be the access example of clinical records. In England, the National Health Service contracted with Amazon to enable Alexa in 2019 to answer health-related questions, raising questions about privacy and how health care data would be used [44,45]. The HIPAA compliance and the fact that the National Health Service contracted with Amazon explains why most studies have been carried out in the United States and the United Kingdom. Arguably, European countries are not as present due to more strict data protection regulations. Moreover, the use of smart speakers is significantly higher in the United States than in other countries, which in turn could also be related to data protection regulations [77]. Interestingly, Asian countries have, with few exceptions, also not been represented in the included articles. This seems counterintuitive as, in terms of market sales, smart speaker technology by Asian technology companies is more and more successful [42].

It also became clear that the devices were clearly dominant in the publications. This should be criticized from a scientific point of view. We were able to identify the following explanations for this result.

Since Amazon entered the market in 2015 and continuously updates its product line, off-the-shelf devices have recently increased in terms of market penetration, making them more popular for research and development. That Amazon's Echo was used in the vast majority of articles included comes as no surprise, and Amazon's market dominance is based on several factors. First, the company was the first to release a smart speaker to consumers. Second, Amazon's voice assistant, Alexa, has been embedded in a broad range of devices, including wall clocks, by third-party manufacturers. Third, Amazon sells products of the Echo family at comparably low prices, starting at around US \$20. Fourth, Amazon offers an infrastructure through its Skill Store and several frameworks for developers. Fifth, in the United States, the Echo is HIPAA-compliant.

The dominance of Amazon's smart speaker in the included papers poses several risks depending on the use case, some of

XSL•FO RenderX

which are discussed in the papers themselves. In terms of the devices themselves in their off-the-shelf version, the interaction is limited. For example, Nallam et al [49] used a smart speaker prototype as they argue that developed solutions often do not support conversational interactions and explore scenarios that are not yet supported.

The articles included in this publication address a diverse range of use cases across various settings, thereby demonstrating the versatility of smart speakers and the technology of NLP and AI incorporated in them. This technology can be used in a multitude of contexts within the domains of health care and social care. Overall, 2 general use cases can be distinguished: (1) supporting patients and their relatives in their private living environments and (2) supporting professional health care workers in clinical settings. As the devices were originally developed for private home environments and primarily for entertainment and e-commerce applications, it is unsurprising that this setting was the dominant one across the papers included in this review. This could be seen as an indicator of the restructuring of health care services, with an increased focus on the private living environment. Several clinical use cases supported by smart speakers could be automated and not be restricted to clinical settings (eg, [14,48]). Only in a few cases does the paper focus on clinical use cases and professional personnel (eg, [4,45,71]).

That patients, and particularly older adults, were the main target group supports this conclusion. Moreover, this also underlines that the role of patients and practices of health and care change against the background of digitalization and the use of AI [78]. While some of the use cases identified were exclusively designed for clinical settings, the majority can, in theory, be implemented in multiple settings. This could support patient empowerment, as smart speakers can be used to support the household as a central place of health care. An argument supporting the fit of the devices for older adults is that smart speakers do not require "reasonable levels of vision and manual dexterity" [79,80].

A key rationale for using the devices is not only their competitive pricing but also the potential to reduce expenditure by enhancing the efficiency of staff members and care processes, for instance, through enhanced documentation or facilitating straightforward communication with patients, colleagues, or clients. Although the majority of the papers reviewed argue that smart speakers could provide such benefits, these potential benefits depend on several circumstances. The first is whether the devices can be installed as they are or whether new skills or, more complexly, additional hardware or modifications are required. This depends on the use case and also the target group. Although many people are used to interacting with the devices, older adults might not have any experience and could need training.

The majority of the papers in our sample can be classified as exploratory in nature. The research designs used are predominantly qualitative, with sample sizes that are relatively small and no long-term studies conducted in real-world scenarios. This underscores the fact that the technology itself is still relatively new, particularly within the context of health care and social care. In addition, researchers and developers are

```
https://ai.jmir.org/2025/1/e55673
```

XSL•FO

still exploring the technology's potential applications in health care and social care, which may have become more apparent in the context of the pandemic. Both sectors are currently experiencing financial strain due to rising expenditure and a shortage of qualified personnel [81]. New technologies are frequently viewed as a potential solution to these challenges [70].

Smart speakers and digital voice assistants like Alexa are quite limited in terms of their initial dialogue management, which can be seen as an important motivator to using the systems as they are easier to develop and control. This finding is in line with a systematic review of CA in health care carried out by Laranjo et al [1]. The authors could identify 17 articles using 14 different CA. Most papers covered by the review evaluated task-oriented CA that aims at supporting patients and clinicians. Systems allowing the management of complex dialogues were only identified in 1 case. Even though conversational systems have proven to be beneficial for health-related purposes, most assistants allow only constrained user input (eg, multiple-choice answers) [1,82]. Clark et al [83] argue that users interact in "clearly delineated task-based conversations" and "fall short of reflexive and adaptive interactivity." According to the authors, the term conversation is "a poor description of the current interaction experience" with an AI using common smart speakers [83]. Hence, they suggest testing "human-agent interaction as a new genre of conversation, with its own rules, norms and expectations" [83]. The devices have only a limited capability to actually be able to engage in a conversational dialogue. Conversations are task-oriented instead of offering interactions initiated by the user and not by the device. While this might be true, it seems to be only a matter of time before future updates might be used to allow more natural dialogues, as is already the case with generative AI such as ChatGPT.

The analysis showed that change in existing practices and routines is an important aspect. Drawing on Sezgin et al [84], Capasso and Umbrello [85] argue that the novelty of CAs is that they act as "intermediaries between the health care system as a whole and the public," changing practices in health care and social care. Here, several studies follow the normative aim to implement innovative technologies in order to improve processes and outcomes. The use of smart speakers—or CAs in general—follows a technology-driven approach. Already existing technologies are transferred to the domains of health and social care. Due to the exploratory design of most studies, the emphasis is put on the technology and not on the context, like organizational or social factors. The logic of a "fitting" technology seems to be a main driver of many studies, neglecting the analysis of potentially changing social practices.

The dominance of Amazon in our sample has to be seen from a critical perspective. The company itself began offering the service Alexa Together and was able to emulate existing approaches and leverage its financial and market clout to challenge competitors. Moreover, developers depend on the technology, that is, the hardware and also the software frameworks of one company. As a consequence, the dominant position of Amazon might increase due to research using the company's products. If only one product from a particular company is examined, the capabilities of other products are not

taken into account, as they may perform better, for example, and might be used to copy promising applications.

Limitations

This paper has several limitations. First, the number of databases searched. To address this limitation, a cross-search was performed in Google Scholar to rule out the possibility that important articles were not found. In addition, to broaden the search strategy, other forms of literature, such as trial reports, could be included in future studies. For instance, a few trials using smart speakers are registered on clinicaltrials.gov. However, we decided not to include these as they did not provide all the information we wanted to obtain (eg, motivations for using the devices). Second, we restricted our search to the English language only. Few papers were found from the Asian region, probably due to the language limitation of the search. This limitation was mitigated by using brand names as search terms focusing on the brands with the highest market share. However, as recent market research shows, there is a shift toward products developed in Asian countries, and future studies should include a wider range of brands and products. Another limitation is that we only looked at smart speakers, which excludes other voice assistants that use essentially the same technology (such as digital assistants on smartphones and tablets). We deliberately excluded these as this review focused specifically on smart speakers as a form of CA, and we argue that the technology of smart speakers needs to be seen as a technology in its own right.

Conclusion

In this paper, a scoping review was conducted on the use of smart speakers in health care and social care settings. The analysis showed that—due to the widespread use of devices like Amazon's Echo—smart speaker technology has been tested and implemented in various settings and use cases in the health and social care sectors. The main setting was the private home environment, and the main user group was patients. There are, however, also approaches to making use of the technology in other settings, such as hospitals. It seems likely that due to technical progress in the field of AI and the market power of the companies behind the devices, there will be more use cases of smart speakers in the (near) future.

Acknowledgments

This study was supported by the Federal Ministry of Education and Research (grant number 16SV8791).

Data Availability

The datasets generated during and/or analyzed during this study are available from the corresponding author upon reasonable request.

Conflicts of Interest

None declared.

Multimedia Appendix 1 Preferred Reporting Items for Systematic reviews and Meta-Analyses extension for Scoping Reviews (PRISMA-ScR) checklist. [DOCX File, 22 KB - ai v4i1e55673 app1.docx]

Multimedia Appendix 2 Database search details. [DOCX File , 18 KB - ai v4i1e55673 app2.docx]

Multimedia Appendix 3 Data extraction table. [DOCX File, 38 KB - ai v4i1e55673 app3.docx]

References

- Laranjo L, Dunn AG, Tong HL, Kocaballi AB, Chen J, Bashir R, et al. Conversational agents in healthcare: a systematic review. J Am Med Inform Assoc 2018;25(9):1248-1258 [FREE Full text] [doi: 10.1093/jamia/ocy072] [Medline: 30010941]
- Provoost S, Lau HM, Ruwaard J, Riper H. Embodied conversational agents in clinical psychology: a scoping review. J Med Internet Res 2017;19(5):e151 [FREE Full text] [doi: 10.2196/jmir.6553] [Medline: 28487267]
- 3. Edwards KJ, Jones RB, Shenton D, Page T, Maramba I, Warren A, et al. The use of smart speakers in care home residents: implementation study. J Med Internet Res 2021;23(12):e26767 [FREE Full text] [doi: 10.2196/26767] [Medline: 34932010]
- 4. Kim S. Exploring how older adults use a smart speaker-based voice assistant in their first interactions: qualitative study. JMIR mHealth uHealth 2021;9(1):e20427 [FREE Full text] [doi: 10.2196/20427] [Medline: 33439130]

https://ai.jmir.org/2025/1/e55673

- Tudor Car L, Dhinagaran DA, Kyaw BM, Kowatsch T, Joty S, Theng Y, et al. Conversational agents in health care: scoping review and conceptual analysis. J Med Internet Res 2020;22(8):e17158 [FREE Full text] [doi: <u>10.2196/17158</u>] [Medline: <u>32763886</u>]
- 6. Ermolina A, Tiberius V. Voice-controlled intelligent personal assistants in health care: international Delphi study. J Med Internet Res 2021;23(4):e25312 [FREE Full text] [doi: 10.2196/25312] [Medline: 33835032]
- 7. Chung AE, Griffin AC, Selezneva D, Gotz D. Health and fitness apps for hands-free voice-activated assistants: content analysis. JMIR mHealth uHealth 2018;6(9):e174 [FREE Full text] [doi: 10.2196/mhealth.9705] [Medline: 30249581]
- Kumah-Crystal YA, Pirtle CJ, Whyte HM, Goode ES, Anders SH, Lehmann CU. Electronic health record interactions through voice: a review. Appl Clin Inform 2018;9(3):541-552 [FREE Full text] [doi: 10.1055/s-0038-1666844] [Medline: 30040113]
- Corbett CF, Combs EM, Chandarana PS, Stringfellow I, Worthy K, Nguyen T, et al. Medication adherence reminder system for virtual home assistants: mixed methods evaluation study. JMIR Form Res 2021;5(7):e27327 [FREE Full text] [doi: 10.2196/27327] [Medline: 34255669]
- 10. Vuppalapati JS, Kedari S, Ilapakurti A, Kedari S, Gudivada M, Vuppalapati C. The role of voice service technologies in creating the next generation outpatient data driven electronic health record (EHR). : IEEE; 2017 Presented at: 2017 Intelligent Systems Conference (IntelliSys); September 7-8, 2017; London, United Kingdom. [doi: 10.1109/intellisys.2017.8324289]
- Sezgin E, Noritz G, Elek A, Conkol K, Rust S, Bailey M, et al. Capturing at-home health and care information for children with medical complexity using voice interactive technologies: multi-stakeholder viewpoint. J Med Internet Res 2020;22(2):e14202 [FREE Full text] [doi: 10.2196/14202] [Medline: 32053114]
- 12. Basatneh R, Najafi B, Armstrong DG. Health sensors, smart home devices, and the internet of medical things: an opportunity for dramatic improvement in care for the lower extremity complications of diabetes. J Diabetes Sci Technol 2018;12(3):577-586 [FREE Full text] [doi: 10.1177/1932296818768618] [Medline: 29635931]
- Sadavarte SS, Bodanese E. Pregnancy companion chatbot using Alexa and Amazon Web Services. : IEEE; 2019 Presented at: 2019 IEEE Pune Section International Conference (PuneCon); December 18-20, 2019; Pune, India. [doi: 10.1109/punecon46936.2019.9105762]
- 14. Ooster J, Moreta PNP, Bach JH, Holube I, Meyer BT. 'Computer, Test My Hearing': accurate speech audiometry with smart speakers. 2019 Presented at: Interspeech 2019; 2019 September 15-19; Graz, Austria. [doi: <u>10.21437/interspeech.2019-2118</u>]
- 15. Jadczyk T, Kiwic O, Khandwalla RM, Grabowski K, Rudawski S, Magaczewski P, et al. Feasibility of a voice-enabled automated platform for medical data collection: cardioCube. Int J Med Inform 2019;129:388-393. [doi: 10.1016/j.ijmedinf.2019.07.001] [Medline: 31445282]
- Rampioni M, Stara V, Felici E, Rossi L, Paolini S. Embodied conversational agents for patients with dementia: thematic literature analysis. JMIR mHealth uHealth 2021;9(7):e25381 [FREE Full text] [doi: 10.2196/25381] [Medline: 34269686]
- 17. Waldhör K. Smarte objekte wie smart speaker und smarthome die medizinische und pflegerische versorgung zu hause unterstützen werden [Book in German]. In: Digitale Transformation von Dienstleistungen im Gesundheitswesen VI. Wiesbaden: Springer Fachmedien Wiesbaden; 2019:389-406.
- 18. Smart speakers statistics: report 2022. Speakergy. 2022. URL: <u>https://speakergy.com/smart-speakers-statistics/</u> #:~:text=The%20United%20States%20Smart%20Speaker,a%206%25%20increase%20from%202020 [accessed 2022-09-30]
- 19. Petrock V. Voice assistant and smart speaker users 2020: more time at home means more time to talk. 2020. URL: <u>https://www.emarketer.com/content/voice-assistant-and-smart-speaker-users-2020</u> [accessed 2021-12-02]
- 20. INSIGHTS 2020: device usage 2020. AudienceProject. 2020. URL: <u>https://www.audienceproject.com/wp-content/uploads/</u> audienceproject study device usage 2020.pdf [accessed 2021-12-02]
- 21. Initiative D21 e. V. D21-Digital-Index 2021/2022 [Website in German]. Jährliches Lagebild zur Digitalen Gesellschaft. 2022. URL: <u>https://initiatived21.de/publikationen/d21-digital-index/2021-2022</u> [accessed 2024-12-10]
- 22. Welcome to 'The Age of Voice 3.0': OMD Germany. OMD. 2021. URL: <u>https://www.omd.com/news/</u> welcome-to-the-age-of-voice-3-0/ [accessed 2023-02-18]
- 23. Gaspar C, Neus A. Smart-speaker-report 2023: erfahrungen, bewertungen und wünsche der nutzer in Deutschland, UK und Den USA [Article in German]. Nürnberg Institut für Marktentscheidungen e.V. 2023. URL: <u>https://www.nim.org/</u> publikationen/detail/smart-speaker-report-2023 [accessed 2024-12-10]
- 24. Baertsch MA, Decker S, Probst L, Joneleit S, Salwender H, Frommann F, et al. Convenient access to expert-reviewed health information via an alexa voice assistant skill for patients with multiple myeloma: development study. JMIR Cancer 2022;8(2):e35500 [FREE Full text] [doi: 10.2196/35500] [Medline: 35679096]
- 25. Beaman J, Lawson L, Keener A, Mathews ML. Within clinic reliability and usability of a voice-based Amazon Alexa administration of the patient health questionnaire 9 (PHQ 9). J Med Syst 2022;46(6):38 [FREE Full text] [doi: 10.1007/s10916-022-01816-0] [Medline: 35536347]
- 26. Brewer RN. 'If Alexa knew the state I was in, it would cry': older adults' perspectives of voice assistants for health. 2022 Presented at: CHI Conference on Human Factors in Computing Systems Extended Abstracts; April 29, 2022; New Orleans, LA, USA p. 1-8. [doi: 10.1145/3491101.3519642]
- 27. Sunshine J. Smart speakers: the next frontier in mHealth. JMIR mHealth uHealth 2022;10(2):e28686. [doi: 10.2196/28686] [Medline: 35188467]

- Kocaballi AB, Berkovsky S, Quiroz JC, Laranjo L, Tong HL, Rezazadegan D, et al. The personalization of conversational agents in health care: systematic review. J Med Internet Res 2019;21(11):e15360 [FREE Full text] [doi: 10.2196/15360] [Medline: 31697237]
- 29. Kraus S, Schiavone F, Pluzhnikova A, Invernizzi AC. Digital transformation in healthcare: analyzing the current state-of-research. J Bus Res 2021;123:557-567. [doi: 10.1016/j.jbusres.2020.10.030]
- Kocaballi AB, Sezgin E, Clark L, Carroll JM, Huang Y, Huh-Yoo J, et al. Design and evaluation challenges of conversational agents in health care and well-being: selective review study. J Med Internet Res 2022;24(11):e38525 [FREE Full text] [doi: 10.2196/38525] [Medline: 36378515]
- Milne-Ives M, de Cock C, Lim E, Shehadeh MH, de Pennington N, Mole G, et al. The effectiveness of artificial intelligence conversational agents in health care: systematic review. J Med Internet Res 2020;22(10):e20346 [FREE Full text] [doi: 10.2196/20346] [Medline: <u>33090118</u>]
- Bin Sawad A, Narayan B, Alnefaie A, Maqbool A, Mckie I, Smith J, et al. A systematic review on healthcare artificial intelligent conversational agents for chronic conditions. Sensors (Basel) 2022;22(7):2625 [FREE Full text] [doi: 10.3390/s22072625] [Medline: 35408238]
- Schachner T, Keller R, V Wangenheim F. Artificial intelligence-based conversational agents for chronic conditions: systematic literature review. J Med Internet Res 2020;22(9):e20701 [FREE Full text] [doi: <u>10.2196/20701</u>] [Medline: <u>32924957</u>]
- 34. Jahan N, Naveed S, Zeshan M, Tahir MA. How to conduct a systematic review: a narrative literature review. Cureus 2016;8(11):e864 [FREE Full text] [doi: 10.7759/cureus.864] [Medline: 27924252]
- Munn Z, Peters MDJ, Stern C, Tufanaru C, McArthur A, Aromataris E. Systematic review or scoping review? Guidance for authors when choosing between a systematic or scoping review approach. BMC Med Res Methodol 2018;18(1):143
 [FREE Full text] [doi: 10.1186/s12874-018-0611-x] [Medline: 30453902]
- 36. Arksey H, O'Malley L. Scoping studies: towards a methodological framework. Int J Soc Res 2005;8(1):19-32. [doi: 10.1080/1364557032000119616]
- 37. Levac D, Colquhoun H, O'Brien KK. Scoping studies: advancing the methodology. Implement Sci 2010;5:69 [FREE Full text] [doi: 10.1186/1748-5908-5-69] [Medline: 20854677]
- 38. Peters MDJ, Godfrey CM, Khalil H, McInerney P, Parker D, Soares CB. Guidance for conducting systematic scoping reviews. Int J Evid Based Healthc 2015;13(3):141-146. [doi: 10.1097/XEB.0000000000000050] [Medline: 26134548]
- O'Brien K, Liggett A, Ramirez-Zohfeld V, Sunkara P, Lindquist LA. Voice-controlled intelligent personal assistants to support aging in place. J Am Geriatr Soc 2020;68(1):176-179. [doi: <u>10.1111/jgs.16217</u>] [Medline: <u>31617581</u>]
- 40. Page MJ, McKenzie JE, Bossuyt PM, Boutron I, Hoffmann TC, Mulrow CD, et al. The PRISMA 2020 statement: an updated guideline for reporting systematic reviews. BMJ 2021;372:n71 [FREE Full text] [doi: 10.1136/bmj.n71] [Medline: 33782057]
- Tricco AC, Lillie E, Zarin W, O'Brien KK, Colquhoun H, Levac D, et al. PRISMA extension for scoping reviews (PRISMA-ScR): checklist and explanation. Ann Intern Med 2018;169(7):467-473 [FREE Full text] [doi: 10.7326/M18-0850] [Medline: 30178033]
- 42. Strategy analytics: global smart speaker shipments declined 5% in 1Q22 amid disruption from war and a resurgent COVID virus. Businesswire. 2022. URL: <u>https://www.businesswire.com/news/home/20220606005136/en/</u> <u>Strategy-Analytics-Global-Smart-Speaker-Shipments-Declined-5-in-1Q22-Amid-Disruption-from-War-and-a-Resurgent-COVID-Virus</u> [accessed 2022-09-30]
- 43. Kuckartz U. Qualitative Text Analysis: A Systematic Approach. In: Compendium for Early Career Researchers in Mathematics Education. Cham: Springer Nature; 2019:181-197.
- 44. Wright J. The Alexafication of adult social care: virtual assistants and the changing role of local government in England. Int J Environ Res Public Health 2021;18(2):812 [FREE Full text] [doi: 10.3390/ijerph18020812] [Medline: 33477872]
- Bhatt V, Li J, Maharjan B. DocPal: a voice-based EHR assistant for health practitioners. : IEEE; 2021 Presented at: 2020 IEEE International Conference on E-health Networking, Application & Services (HEALTHCOM); 2021 March 01-02; Shenzhen, China. [doi: 10.1109/healthcom49281.2021.9399013]
- 46. Wang A, Nguyen D, Sridhar AR, Gollakota S. Using smart speakers to contactlessly monitor heart rhythms. Commun Biol 2021;4(1):319 [FREE Full text] [doi: 10.1038/s42003-021-01824-9] [Medline: 33750897]
- 47. O'Brien K, Light SW, Bradley S, Lindquist L. Optimizing voice-controlled intelligent personal assistants for use by home-bound older adults. J Am Geriatr Soc 2022;70(5):1504-1509 [FREE Full text] [doi: 10.1111/jgs.17625] [Medline: 35029296]
- Sharma A, Oulousian E, Ni J, Lopes R, Cheng MP, Label J, et al. Voice-based screening for SARS-CoV-2 exposure in cardiovascular clinics. Eur Heart J Digit Health 2021;2(3):521-527 [FREE Full text] [doi: 10.1093/ehjdh/ztab055] [Medline: 36713601]
- Nallam P, Bhandari S, Sanders J, Martin-Hammond A. A question of access: exploring the perceived benefits and barriers of intelligent voice assistants for improving access to consumer health resources among low-income older adults. Gerontol Geriatr Med 2020;6:2333721420985975 [FREE Full text] [doi: 10.1177/2333721420985975] [Medline: 33457459]

```
https://ai.jmir.org/2025/1/e55673
```

- Domínguez D, Morales L, Sánchez N. IoMT-Driven eHealth: a technological innovation proposal based on smart speakers. In: Rojas I, Valenzuela O, Rojas F, editors. Bioinformatics and Biomedical Engineering. Cham: Springer International Publishing; 2020:378-386.
- 51. Lee E, Vesonder G, Wendel E. Eldercare robotics Alexa. : IEEE; 2020 Presented at: 2020 11th IEEE Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON); October 28-31, 2020; New York, NY, USA p. 820-825. [doi: 10.1109/uemcon51285.2020.9298147]
- Jansons P, Fyfe J, Via JD, Daly RM, Gvozdenko E, Scott D. Barriers and enablers for older adults participating in a home-based pragmatic exercise program delivered and monitored by Amazon Alexa: a qualitative study. BMC Geriatr 2022;22(1):248 [FREE Full text] [doi: 10.1186/s12877-022-02963-2] [Medline: 35337284]
- 53. Qiu L, Kanski B, Doerksen S, Winkels RM, Schmitz K, Abdullah S. Nurse AMIE: using smart speakers to provide supportive care intervention for women with metastatic breast cancer. : ACM; 2021 Presented at: CHI '21: CHI Conference on Human Factors in Computing Systems; 2021 May 8-13; Yokohama Japan. [doi: 10.1145/3411763.3451827]
- 54. Thomas G. Patient and clinician-centric healthcare enhancement through speech recognition: a research proposal. 2019 Presented at: 7th Annual International Conference on Architecture and Civil Engineering (ACE 2019) GSTF 2019; May 27-28, 2019; Singapore URL: <u>https://dl4.globalstf.org/products-page/books/</u> patient-and-clinician-centric-healthcare-enhancement-through-speech-recognition/ [doi: 10.5176/2301-394X_ACE19.581]
- 55. Cheng A, Raghavaraju V, Kanugo J, Handrianto YP, Shang Y. Development and evaluation of a healthy coping voice interface application using the Google home for elderly patients with type 2 diabetes. : IEEE; 2018 Presented at: 2018 15th IEEE Annual Consumer Communications & Networking Conference (CCNC); January 12-15, 2018; Las Vegas, NV, USA p. 1-5. [doi: 10.1109/ccnc.2018.8319283]
- 56. Luo Y, Lee B, Choe E. TandemTrack: shaping consistent exercise experience by complementing a mobile app with a smart speaker. : ACM; 2020 Presented at: Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems; 2020 April 25-30; Honolulu HI USA p. 1-13. [doi: 10.1145/3313831.3376616]
- 57. Arem H, Scott R, Greenberg D, Kaltman R, Lieberman D, Lewin D. Assessing breast cancer survivors' perceptions of using voice-activated technology to address insomnia: feasibility study featuring focus groups and in-depth interviews. JMIR Cancer 2020;6(1):e15859 [FREE Full text] [doi: 10.2196/15859] [Medline: 32348274]
- 58. Dojchinovski D, Ilievski A, Gusev M. Interactive home healthcare system with integrated voice assistant. 2019 Presented at: 42nd International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO); May 20-24, 2019; Opatija, Croatia URL: <u>https://ieeexplore.ieee.org/document/8756983</u> [doi: 10.23919/MIPRO.2019.8756983]
- Ilievski A, Dojchinovski D, Gusev M. Interactive voice assisted home healthcare systems. New York, NY: Association for Computing Machinery; 2019 Presented at: BCI'19: 9th Balkan Conference in Informatics; September 26-28, 2019; Sofia, Bulgaria. [doi: 10.1145/3351556.3351572]
- Yoo TK, Oh E, Kim H, Ryu IH, Lee IS, Kim JS, et al. Deep learning-based smart speaker to confirm surgical sites for cataract surgeries: a pilot study. PLoS One 2020;15(4):e0231322 [FREE Full text] [doi: 10.1371/journal.pone.0231322] [Medline: 32271836]
- Ismail HO, Moses AR, Tadrus M, Mohamed EA, Jones LS. Feasibility of use of a smart speaker to administer Snellen visual acuity examinations in a clinical setting. JAMA Netw Open 2020 Aug 03;3(8):e2013908 [FREE Full text] [doi: 10.1001/jamanetworkopen.2020.13908] [Medline: 32822489]
- 62. Chambers R, Beaney P. The potential of placing a digital assistant in patients' homes. Br J Gen Pract 2020 Jan;70(690):8-9 [FREE Full text] [doi: 10.3399/bjgp20X707273] [Medline: 31879289]
- 63. Kim JH, Um R, Liu J, Patel J, Curry E, Aghabaglou F, et al. Development of a smart hospital assistant: integrating artificial intelligence and a voice-user interface for improved surgical outcomes. Proc SPIE Int Soc Opt Eng 2021 Feb;11601:116010U [FREE Full text] [doi: 10.1117/12.2580995] [Medline: 35341075]
- 64. Jansons P, Dalla Via J, Daly RM, Fyfe JJ, Gvozdenko E, Scott D. Delivery of home-based exercise interventions in older adults facilitated by Amazon Alexa: a 12-week feasibility trial. J Nutr Health Aging 2022;26(1):96-102 [FREE Full text] [doi: 10.1007/s12603-021-1717-0] [Medline: 35067710]
- 65. Apergi LA, Bjarnadottir MV, Baras JS, Golden BL, Anderson KM, Chou J, et al. Voice interface technology adoption by patients with heart failure: pilot comparison study. JMIR mHealth uHealth 2021 Apr 01;9(4):e24646 [FREE Full text] [doi: 10.2196/24646] [Medline: 33792556]
- Martin-Hammond A, Vemireddy S, Rao K. Exploring older adults' beliefs about the use of intelligent assistants for consumer health information management: a participatory design study. JMIR Aging 2019;2(2):e15381 [FREE Full text] [doi: 10.2196/15381] [Medline: <u>31825322</u>]
- 67. Bickmore TW, Caruso L, Clough-Gorr K. Acceptance and usability of a relational agent interface by urban older adults. 2005 Presented at: Extended Abstracts Proceedings of the 2005 Conference on Human Factors in Computing Systems, CHI 2005; 2005 April 2-7; Portland, Oregon, USA p. 1212-1215. [doi: 10.1145/1056808.1056879]
- 68. Vardoulakis LP, Ring L, Barry B, Sidner CL, Bickmore T. Designing Relational Agents as Long Term Social Companions for Older Adults. In: Hutchison D, Kanade T, Kittler J, editors. Intelligent Virtual Agents. Berlin, Heidelberg: Springer Berlin Heidelberg; 2012:289-302.

```
https://ai.jmir.org/2025/1/e55673
```

- 69. Jesús-Azabal M, Medina-Rodríguez J, Durán-García J, García-Pérez D. Remembranza Pills: Using Alexa to Remind the Daily Medicine Doses to Elderly. In: García-Alonso J, Fonseca C, editors. Gerontechnology. Cham: Springer International Publishing; 2020:151-159.
- Henwood F, Marent B. Understanding digital health: productive tensions at the intersection of sociology of health and science and technology studies. Sociol Health Illn 2019;41 Suppl 1:1-15. [doi: 10.1111/1467-9566.12898] [Medline: 31599984]
- 71. Jadczyk T, Wojakowski W, Tendera M, Henry TD, Egnaczyk G, Shreenivas S. Artificial intelligence can improve patient management at the time of a pandemic: the role of voice technology. J Med Internet Res 2021;23(5):e22959 [FREE Full text] [doi: 10.2196/22959] [Medline: 33999834]
- 72. Palumbo F, Crivello A, Furfari F, Girolami M, Mastropietro A, Manferdelli G, et al. 'Hi This Is NESTORE, Your Personal Assistant': design of an integrated IoT system for a personalized coach for healthy aging. Front Digit Health 2020;2:545949 [FREE Full text] [doi: 10.3389/fdgth.2020.545949] [Medline: 34713033]
- 73. Foehr J, Germelmann CC. Alexa, can I trust you? Exploring consumer paths to trust in smart voice-interaction technologies. J Assoc for Consum Res 2020;5(2):181-205. [doi: 10.1086/707731]
- 74. Gouda P, Ganni E, Chung P, Randhawa VK, Marquis-Gravel G, Avram R, et al. Feasibility of incorporating voice technology and virtual assistants in cardiovascular care and clinical trials. Curr Cardiovasc Risk Rep 2021;15(8):13 [FREE Full text] [doi: 10.1007/s12170-021-00673-9] [Medline: 34178205]
- 75. Sheon AR, Bolen SD, Callahan B, Shick S, Perzynski AT. Addressing disparities in diabetes management through novel approaches to encourage technology adoption and use. JMIR Diabetes 2017;2(2):e16 [FREE Full text] [doi: 10.2196/diabetes.6751] [Medline: 30291090]
- 76. Kowalska M, Gładyś A, Kalańska-Łukasik B, Gruz-Kwapisz M, Wojakowski W, Jadczyk T. Readiness for voice technology in patients with cardiovascular diseases: cross-sectional study. J Med Internet Res 2020;22(12):e20456 [FREE Full text] [doi: 10.2196/20456] [Medline: 33331824]
- 77. Coyne M, Franzese C. The Promise of Voice: Connecting Drug Delivery Through Voice-Activated Technology. East Sussex, United Kingdom: Frederick Furness Publishing Ltd; 2017.
- 78. Marent B, Henwood F. Digital health: a sociomaterial approach. Sociol Health Illn 2023;45(1):37-53 [FREE Full text] [doi: 10.1111/1467-9566.13538] [Medline: 36031756]
- Ho DKH. Voice-controlled virtual assistants for the older people with visual impairment. Eye (Lond) 2018;32(1):53-54 [FREE Full text] [doi: 10.1038/eye.2017.165] [Medline: 28776586]
- Even C, Hammann T, Heyl V, Rietz C, Wahl H, Zentel P, et al. Benefits and challenges of conversational agents in older adults : a scoping review. Z Gerontol Geriatr 2022;55(5):381-387. [doi: <u>10.1007/s00391-022-02085-9</u>] [Medline: <u>35852588</u>]
- Marjanovic S, Altenhofer M, Hocking L, Chataway J, Ling T. Innovating for improved healthcare: sociotechnical and innovation systems perspectives and lessons from the NHS. Science and Public Policy 2020;47(2):1-15. [doi: 10.1093/scipol/scaa005]
- Anastasiadou U, Alexiadis A, Polychronidou E, Votis K, Tzovaras D. A prototype educational virtual assistant for diabetes management. : IEEE; 2020 Presented at: 2020 IEEE 20th International Conference on Bioinformatics and Bioengineering (BIBE); October 26-28, 2020; Cincinnati, OH, USA p. 999-1004. [doi: 10.1109/bibe50027.2020.00169]
- Clark L, Pantidi N, Cooney O, Doyle P, Garaialde D, Edwards J, et al. What makes a good conversation? Challenges in designing truly conversational agents. : ACM; 2019 Presented at: Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems; May 4-9, 2019; Glasgow, Scotland, UK p. 1-12. [doi: <u>10.1145/3290605.3300705</u>]
- 84. Sezgin E, Huang Y, Ramtekkar U. Readiness for voice assistants to support healthcare delivery during a health crisis and pandemic. NPJ Digit Med 2020;3(1):122 [FREE Full text] [doi: 10.1038/s41746-020-00332-0]
- 85. Capasso M, Umbrello S. Responsible nudging for social good: new healthcare skills for AI-driven digital personal assistants. Med Health Care Philos 2022;25(1):11-22 [FREE Full text] [doi: 10.1007/s11019-021-10062-z] [Medline: 34822096]

Abbreviations

AI: artificial intelligence
CA: conversational agent
HIPAA: Health Insurance Portability and Accountability Act
NLP: natural language processing
PRISMA: Preferred Reporting Items for Systematic Reviews and Meta-Analyses
PRISMA-ScR: Preferred Reporting Items for Systematic Reviews and Meta-Analyses Extension for Scoping Reviews
SR: speech recognition



Edited by JL Raisaro; submitted 20.12.23; peer-reviewed by M Chatzimina, H Younes, H Huang; comments to author 18.04.24; revised version received 13.06.24; accepted 24.11.24; published 13.01.25. <u>Please cite as:</u> Merkel S, Schorr S Identification of Use Cases, Target Groups, and Motivations Around Adopting Smart Speakers for Health Care and Social Care Settings: Scoping Review JMIR AI 2025;4:e55673 URL: https://ai.jmir.org/2025/1/e55673 doi:10.2196/55673 PMID:<u>39804689</u>

©Sebastian Merkel, Sabrina Schorr. Originally published in JMIR AI (https://ai.jmir.org), 13.01.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Survey on Pain Detection Using Machine Learning Models: Narrative Review

Ruijie Fang¹, BEng; Elahe Hosseini¹, MS; Ruoyu Zhang¹, MS; Chongzhou Fang¹, BE; Setareh Rafatirad², PhD; Houman Homayoun¹, PhD

¹Department of Electrical and Computer Engineering, University of California, Davis, CA, United States ²Department of Computer Science, University of California, Davis, CA, United States

Corresponding Author:

Ruijie Fang, BEng Department of Electrical and Computer Engineering University of California One Shields Avenue Davis, CA, 95616 United States Phone: 1 5308676009 Email: rjfang@ucdavis.edu

Abstract

Background: Pain, a leading reason people seek medical care, has become a social issue. Automated pain assessment has seen notable advancements over recent decades, addressing a critical need in both clinical and everyday settings.

Objective: The objective of this survey was to provide a comprehensive overview of pain and its mechanisms, to explore existing research on automated pain recognition modalities, and to identify key challenges and future directions in this field.

Methods: A literature review was conducted, analyzing studies focused on various modalities for automated pain recognition. The modalities reviewed include facial expressions, physiological signals, audio cues, and pupil dilation, with a focus on their efficacy and application in pain assessment.

Results: The survey found that each modality offers unique contributions to automated pain recognition, with facial expressions and physiological signals showing particular promise. However, the reliability and accuracy of these modalities vary, often depending on factors such as individual variability and environmental conditions.

Conclusions: While automated pain recognition has progressed considerably, challenges remain in achieving consistent accuracy across diverse populations and contexts. Future research directions are suggested to address these challenges, enhancing the reliability and applicability of automated pain assessment in clinical practice.

(JMIR AI 2025;4:e53026) doi:10.2196/53026

KEYWORDS

pain; pain assessment; machine learning; survey; mobile phone

Introduction

Pain is "an unpleasant sensory and emotional experience associated with actual or potential tissue damage, or described in terms of such damage," according to the International Association for the Study of Pain [1]. However, the discussion on the most precise definition of pain is still ongoing, and the advances in the understanding of pain instantiate the biopsychosocial perspective on pain to capture evidence-based understanding and the evolution of pain [2]. On the basis of the pain origin, it is categorized as nociceptive (due to stimulation of sensory nerve fibers), neuropathic (due to impaired somatosensory nervous system), or psychogenic pain (caused, increased, or prolonged by mental, emotional, or behavioral factors). On the basis of the time duration of the pain, it may be categorized as acute (short duration) or chronic (long duration, may last >3 months).

Approximately 20% of adults have chronic pain in the United States, and chronic pain is the most common reason adults seek medical care. For society, chronic pain contributes to an estimated US \$560 million each year in medical expenses, lost productivity, and disability caused by types of pain such as low back pain, arthritis, and joint pain [3,4]. These negative impacts make chronic pain a persistent public health concern.

Inappropriate pain management can lead to very deleterious physical, psychological, social, and financial consequences for patients. Untreated pain can lead to chronic pain syndrome, which is often accompanied by decreased mobility, impaired immunity, decreased concentration, anorexia, and sleep disturbances. More importantly, the use of prescription opioids for the treatment of chronic noncancer pain is associated with a substantial risk for abuse, dependence, and overdose [5].

As the first step of pain management, pain assessment holds an essential role [6]. Unprecise pain assessment can lead to severe consequences. Undertreatment of pain not only causes psychological consequences but also physiological consequences, for example, increased blood pressure and heart rate. By contrast, overtreatment of pain may result in nausea, vomiting, or constipation immediately and drug addiction in the long term. Traditionally, pain assessment is conducted through self-reports or observational scales. Self-report refers to the conscious communication of pain-related information by the person in pain, typically using spoken or written language or gestures. Various pain rating scales have been developed to capture patients' self-report of pain intensity. Traditional approaches used to play an important role in pain assessment, including the Verbal Rating Scale [7], the Visual Analog Scale

[8], the Numerical Rating Scale [9], and the Wong-Baker FACES Scale [10].

However, such scoring methods are not feasible for certain patients, such as such as those who are unconscious. For this, different observational pain scales, such as the Behavioral Pain Scale [11], Pain Assessment in Advanced Dementia [12], or Neonatal Infant Pain Scale [13], are used in clinical settings. Most scales consider facial expressions, vocalizations, and body language, while some include vital parameters. It is difficult to assess and compare the validity of the various scales because studies differ a lot in design, methodology, participants, and conceptualization of the pain phenomenon. Pain assessment through observation is very challenging and is affected by the subjective biases and errors in beliefs of the observer [14].

To solve these challenges, it is necessary to develop an objective, accurate, continuous pain assessment method, as shown in Figure 1. In the last decades, multiple studies have been conducted to evaluate the feasibility of automated pain assessment using multimodality and machine learning (ML) techniques. This paper surveys and reviews the recent advances in the field in terms of datasets, modalities, and ML models. Finally, we present the challenges remaining in the field and propose future directions.

Figure 1. Typical pipeline of automated pain assessment. FN: false negative; FP: false positive; PR: precision-recall; RNN: recurrent neural network; ROC: receiver operating characteristic; SVM: support vector machine.; TN: true negative; TP: true positive.



Pain Mechanism

The pain mechanism is not completely understood because of its complexity and diversity [15]. Pain, created by the brain, is a psychological state rather than a physical one [16]. Unlike pain, nociception refers to the response of the peripheral and central nervous systems to internal or external stimuli, triggered

Figure 2. Pain mechanism.

by the activation of nociceptors [17]. The noxious stimulus damages the tissue or potentially activates the nociceptors in the peripheral structure. Then, the information is transmitted to the spinal cord dorsal horn or the nucleus caudalis. From there, the information continues to the cerebral cortex via the brainstem in the brain, and the perception of pain is generated. Thus, no brain, no pain [18]. Figure 2 presents the mechanism of pain.



Usually, pain is regarded as chronic or acute according to its duration. Acute pain is a type of sudden pain. The mechanism of momentary pain is well understood [19]. The nociceptors generate the nociception, and the information is transmitted to the brain, where the perception of pain is caused. There are 2 major types of nociceptors responding to different stimuli: C-fibers, associated with unmyelinated axons, and A-delta fibers, associated with thinly myelinated axons [20]. C-fibers generate slow, diffuse pain, while A-delta fibers are related to sharp, pricking pain. Silent nociceptors typically respond to endogenous chemical mediators related to tissue injury [19].

Chronic pain, lasting >3 months, does not have a useful biological function and is challenging to treat due to its varied etiologies [21-23]. According to the *International Classification of Diseases, Eleventh Revision*, chronic pain can be categorized into musculoskeletal, neuropathic, visceral, and cancer pain [21].

Psychological distress refers to a diffuse subjective experience as an internal response to noxious stimuli. Many patients argue that psychological pain is more severe than intense physical pain [24]. Chronic pain can lead to psychological pain and depression, while depression can exacerbate chronic pain [25,26]. Psychogenic pain is physical pain caused or increased by mental and emotional factors [27]. Treatments such as transcutaneous electrical nerve stimulation or psychotherapy are often more effective for reducing psychogenic pain compared to traditional painkillers [28,29].

The body responds to pain via multiple physiological processes: the sympathetic nervous system (SNS), neuroendocrine system, immune system, as well as emotions [30]. The SNS, known for the fight or flight response, increases heart rate and blood pressure via hormones such as catecholamines, epinephrine, and norepinephrine when activated [31]. The SNS also activates sweat glands via acetylcholine, reflecting the active level of

https://ai.jmir.org/2025/1/e53026

SNS through the volume of secreted sweat within a time range [32].

Pain Datasets

Data that are representative are crucial in the creation of a pain recognition system and the demonstration of its efficacy. Crucially, the system should perform optimally within the intended medical context, a fact that must be validated through clinical studies involving patients. In the early stages of development, experimental pain research with healthy volunteers could be useful. This approach allows for strictly controlled conditions, larger participant pools, and the repeated application of pain stimuli. These data are foundational to the development of ML models for automated pain detection.

For studying pain in healthy adults, an external stimulus is needed. Common methods include heat applied via contact (eg, heated objects and electrical heaters) or radiant sources (eg, infrared light). Table 1 summarizes the publicly available datasets that were used for pain recognition research. The UNBC-McMaster Shoulder Pain Expression Archive Database [33] includes 200 video sequences that capture the facial expressions of 25 participants experiencing shoulder pain. Each video sequence includes individuals performing a series of active and passive range-of-motion tests to provoke visible responses to pain, providing a unique dataset rich in both the variety and volume of pain expressions. The dataset includes self-reported and observer assessments of pain intensity at the video level, along with Facial Action Coding System (FACS) coding at the frame level. The BioVid Heat Pain Database [34] is a collection of physiological data and videos from 90 healthy adults subjected to controlled heat stimuli. BioVid consists of several sections: A, B, and C, which focus on pain stimulation, along with sections D and E, which are dedicated to posed expressions and emotion elicitation, respectively. The MIntPAIN

XSL•FO RenderX

database [35] collected color, depth, and thermal videos from 20 healthy adults who were subjected to approximately 1600 instances of electrical pain stimuli at 4 different intensity levels. EmoPain [36], SenseEmotion [37], X-ITE Pain [38], BP4D-Spontaneous [39], and BP4D+ [40] datasets are substantially resources for pain and emotion studies. EmoPain contains video, audio, motion, and a surface electromyogram (sEMG) for lower back pain. SenseEmotion and X-ITE Pain

include audio and physiological data from healthy adults subjected to experimental pain stimuli, while X-ITE provides thermal videos, body movement data, and electromyography measurements. BP4D-Spontaneous and BP4D+ offer facial video recordings from individuals undergoing the cold presser task, with BP4D+ further providing 3D and thermal videos, along with physiological signals.

Table 1. Pain databases.

Database	Participants	Modalities	Annotation
Database with adults			
UNBC-McMaster [33]	25 adults with shoulder pain	Video of the face (RGB ^a)	FACS ^b , VAS ^c , and OPI ^d
BioVid [34]	87 healthy adults	Video of face (RGB), EDA ^e , electrocardiogram, and electromyo- graphy	Stimulus (calibrated per person)
MIntPAIN [35]	20 healthy adults	Video of face (RGB, depth, and thermal)	Stimulus (calibrated per person)
EmoPain [36]	22 adults with chronic back pain	Video, audio, electromyography, and motion capture	Self-report and naive OPI
SenseEmotion [37]	45 healthy adults	Video of face, audio, EDA, electrocardiogram, and electromyography	Stimulus (calibrated per person)
X-ITE [38]	134 healthy adults	Video of face, video of body, audio, EDA, electrocardiogram, and electromyography	Stimulus (calibrated per person)
BP4D-spontaneous [39]	41 healthy adults	Video of face (RGB and 3D)	Stimulus and FACS
BP4D+ [40]	140 healthy adults	Video of face (RGB, 3D, and thermal), heart rate, respiration rate, blood pressure, and EDA	Stimulus and FACS
Database with neonates			
iCOPE [41]	26 healthy neonates	204 RGB photographs of face	Category (pain, rest, cry, air puff, and friction)
YouTube [42]	142 infants	Video and audio	FLACC ^f
APN-db [43]	112 healthy neonates	Video of face (RGB)	NFLAPS ^g , NIPS ^h , and NFCS ⁱ
NPAD-ID [44]	36 healthy neonates and 9 neonates who underwent surgery	Video of face and body (RGB)	NIPS and N-PASS
iCOPEvid [45]	49 neonates	Video of face (grayscale)	Category (pain and no pain)
USF-MNPAD-I [46]	36 neonates	Video of face (RGB), audio, heart rate, blood pressure, SpO_2^{j} , deoxyhemoglobin (HbH), oxyhemoglobin (HbO ₂)	NIPS and N-PASS ^k

^aRGB: Red, green, blue color model.

^bFACS: Facial Action Coding System.

^cVAS: Visual Analog Scale.

^dOPI: Observed Pain Intensity.

^eEDA: electrodermal activity.

^fFLACC: Face, Legs, Activity, Cry, Consolability Scale.

^gNFLAPS: Neonatal Face and Limb Acute Pain Scale

^hNIPS: Neonatal Infant Pain Scale.

¹NFCS: Neonatal Facial Coding System.

^JSpO₂: saturation of peripheral oxygen.

^kN-PASS: Neonatal Pain, Agitation and Sedation Scale.

In the field of infant pain research, the iCOPE [41], YouTube [42], APN-db [43], iCOPEvid [45], and USF-MNPAD-I [46]

https://ai.jmir.org/2025/1/e53026

RenderX

databases are the publicly available datasets. The iCOPE consists of 204 static photographs that capture 26 neonates during various

procedures. The images provide valuable insights into the facial expressions associated with infant pain experiences. The YouTube dataset offers 142 videos accompanied by audio, showcasing the reactions of different infants undergoing immunizations. The APN-db is a dataset that includes >200 videos of infants undergoing various procedures, and it features unique annotations, such as Neonatal Face and Limb Acute Pain intensity. The USF-MNPAD-I dataset collects video, audio, and physiological data from 58 neonates during their hospitalization in the neonatal intensive care unit (ICU) and is annotated using the Neonatal Infant Pain Scale and N-PASS scales.

Postoperative Pain

Although automated pain assessment in controlled settings is well studied, postoperative pain has not been extensively researched due to the difficulty of data collection. Postoperative pain results from tissue injury following surgery and is critical to manage, as inadequate treatment can lead to serious physiological and psychological outcomes. Postoperative pain datasets often exhibit imbalanced distributions and may contain missing labels due to variability in patient experiences and and clinical settings, further complicating accurate comprehensive pain assessment. The NPAD-IA database [44] captures video, audio, and physiological data from 40 infants undergoing procedural (heel lancing and immunization) and postoperative (gastrostomy tube) pain. Notably, it includes postoperative pain data, addressing the complexity and variability of pain levels in real-world clinical settings, thereby enhancing the ecological validity of the assessment. Salekin et al [47] present a novel fully automated deep learning framework to assess neonatal postoperative pain. It uses a bilinear convolutional neural network (B-CNN) to extract facial features and a recurrent neural network (RNN) to model the temporal patterns of postoperative pain. The study uses a dataset of >600 minutes of visual, vocal, and physiological data from neonates, demonstrating the feasibility and efficiency of combining B-CNN and RNN for continuous and accurate assessment of postoperative pain intensity in clinical settings. Salekin et al [46] introduce an automated system for assessing neonatal postoperative pain by integrating visual, vocal, and physiological data. The study also uses a B-CNN for spatial feature extraction but uses a long short-term memory (LSTM) network for capturing temporal patterns, demonstrating that the multimodal spatial-temporal approach significantly outperforms unimodal methods, achieving an area under the curve (AUC) of 0.87 and accuracy of 79%. Automated postoperative pain assessment is still in its nascent stages, primarily hindered by a lack of comprehensive datasets and consistent research efforts. The current methods, often unimodal and focused on short-term procedural pain, fail to capture the complex and prolonged nature of postoperative pain. There is a pressing need for more extensive and diverse datasets to improve the accuracy and

reliability of these systems. Despite these challenges, the potential benefits of automated pain assessment are immense, offering more consistent and objective pain management that can significantly enhance patient outcomes and reduce the burden on health care providers.

Automatic Pain Assessment

Overview

Automated tools for pain assessment have great promise. Because pain results in different physiological and behavioral responses, signals that capture these may be used to detect the presence of pain. However, prior research work has been limited, and automated approaches have not yet become widely used in clinical practice. In this section, we briefly outline the different approaches relevant to the development of automated pain assessment methods described in the research literature. Specifically, we review their system architecture (inputs and outputs) and describe the data sources available for the research and development of ML-based automated pain assessment tools, together with an overview of system validation challenges. This section summarizes the results of the survey of automatic pain detection approaches.

The Use of Modalities

The selection of sensors is a critical aspect of automated pain assessment, as different sensors can convey varying levels of information and have different discriminative abilities. Modalities commonly used in this field can be broadly classified into 3 categories: video, audio, and physiological signals, as shown in Table 2. Functional magnetic resonance imaging (fMRI) was found to be the most prevalent sensor in pain studies, with a prevalence score of 95.9. Electroencephalogram and electrocardiogram were also frequently used, with prevalence scores of 69.6 and 39.1, respectively. In contrast, functional near-infrared spectroscopy (fNIRS) and photoplethysmography had much lower prevalence scores of <10. Moreover, Multimedia Appendix 1 also includes information on modalities used in studies (including brain activity, cardiovascular activity, electrodermal activity (EDA), respiration activity, and pupil size). In terms of physiological signals, activity can be measured brain using electroencephalograms, fMRI, and fNIRS. Cardiovascular activity can be measured using an electrocardiogram or photoplethysmography, while EDA is often measured by skin conductance level or sEMG. To gain insight into the prevalence of each modality, we conducted a search for "Modality AND Pain AND Machine learning" (eg, "EEG AND Pain AND Machine learning") on PubMed and Scopus, limiting the search to the period from January 1, 2010, to August 1, 2023. We then recorded the number of results and normalized them to the range of (0-100) for each database. The prevalence scores were then calculated as the average of the normalized results from PubMed and Scopus.



Table 2. Summary of the commonly used modalities.

Ca	tegory and name	Description	Prevalence ^a	References	
Vio	leo				
	Video analysis	Analyzes facial expressions and body movements to assess pain levels [48].	100	[33,35]	
Au	dio				
	Audio analysis	Analyzes vocal characteristics and speech patterns to assess pain [49].	48.2	[49]	
Pu	pil size				
	Pupil size measurement	Measures changes in pupil diameter as an indicator of pain [50].	12.7	[51,52]	
Bra	ain activity				
	Electroencephalogram	It is a test that detects tiny electrical charges that result from the activity of brain cells [53].	69.6	[54-56]	
	Functional magnetic resonance imaging	It uses magnetic resonance imaging to measure the changes in hemodynamics caused by neuronal activity [57].	95.9	[58-60]	
	Functional near-infrared spectroscopy	It uses scattering arising from the main components of blood upon exposure to near-infrared light (600 nm to 900 nm) to measure changes in oxyhemoglobin and deoxyhemoglobin during brain activity [50].	7.9	[61,62]	
Ca	rdiovascular activity				
	Electrocardiogram	It is a test that measures the electrical activity of the heartbeat [63].	39.1	[64-66]	
	Photoplethysmograph	It is an optical technique that can be used to detect blood volume changes in the microvascular bed of tissue [58].	9.4	[65,67]	
Ele	ectrodermal activity				
	Skin conductance level	It is the measurement of the electrical conductivity of the skin [60].	25.9	[65,66,68]	
	Surface electromyogram	It is a technique to measure muscle activity noninvasively using surface electrodes placed on the skin overlying the muscle [61].	25.6	[66,69,70]	
Re	spiration				
	Respiration	Respiration refers to a person's breathing and the movement of air into and out of the lungs [66].	17.5	[69,71]	

^aPrevalence is measured by the weighted search results from Scopus and PubMed, covering the period from 2010 to 2023, using the keywords "Name" AND "Pain" AND "Machine learning" as of August 1, 2023; the results are standardized on a scale of 0 to 100.

As shown in Table 2, video was found to be the most prevalent sensor in pain studies, with a prevalence score of 100. fMRI, electroencephalogram, and electrocardiogram were also frequently used, with prevalence scores of 95.9, 69.6, and 39.1, respectively. In contrast, fNIRS and photoplethysmography had much lower prevalence scores of <10.

Convenience and feasibility should also be considered when selecting sensors. For example, some sensors such as electroencephalograms and fMRI are nonwearable and can be invasive, which may limit their utility in certain settings. Moreover, complex signals require more sophisticated processing techniques and computing resources, which may not be practical in some situations, such as those involving microprocessors.

Facial Expression

Overview

Facial expression during the experience of pain is not unspecific grimacing but conveys pain-specific information. Studies investigating facial expressions of pain have most often used FACS [48], the gold standard for facial expression research. FACS is a fine-grained, objective, and anatomically based coding system that differentiates between 44 facial movements known as action units (AUs). Coders are trained to apply specific operational criteria to determine the onset and offset as well as the intensity of the AUs. Using FACS, it was shown that facial expressions of pain are composed of a small subset of facial activities, namely, lowering the brows (AU4), cheek raise or lid tightening (AUs 6 and 7), nose wrinkling or raising the upper lip (AUs 9 and 10), and eye closure for >0.5 seconds (AU 43). Prkachin and Solomon [72] developed the Prkachin and Solomon Pain Intensity metric based on this observation, which is a 16-level scale based on the contribution of the individual intensity of pain-related AUs and is defined as follows:

Pain=AU4+(AU6,AU7)+(AU9+AU10)+AU43

Figure 3 shows samples of different PSPI levels from UNBC-McMaster pain dataset. The list of pain-related AUs has been further expanded in more extensive research [73] to include lip corner puller (AU12), lip stretch (AU20), lips part (AU25), jaw drop (AU26), and mouth stretch (AU27).

Fang et al

MCL 2.7.1

Figure 3. Image frame samples of the UNBC-McMaster shoulder pain database. PSPI: Prkachin and Solomon Pain Intensity.



PSPI=0



PSPI=2



PSPI=6



PSPI=10



PSPI=12



PSPI=14

Facial activities during experimental and clinical pain are largely inborn but not uniform across individuals. People display different parts or combinations of facial activities. Cluster analyses identified four distinct facial activity patterns: (1) narrowed eyes with raised upper lip or nose wrinkling and furrowed brows, (2) narrowed eyes with furrowed brows, (3) narrowed eyes with mouth opening, and (4) raised eyebrows, which are less frequent and stable, often indicating novelty or surprise in response to pain. Recognizing these patterns improves pain detection more than focusing on a single expression. Thus, acknowledging variability in facial expressions can enhance pain communication.

Facial expression analysis uses spatial and spatiotemporal features. Spatial features capture static details of the face, such as the geometric and textural characteristics of the eyes, eyebrows, nose, lips, and facial contours, using techniques such as facial landmark detection, geometric feature extraction, Gabor filters, local binary patterns (LBPs), and histogram of oriented gradients (HOG). Spatiotemporal features capture dynamic changes in expressions over time using techniques such as optical flow or differences between consecutive frames. Advanced methods may involve 3D facial modeling or LSTM networks to identify temporal dependencies. Combining spatial and spatiotemporal features provides a comprehensive analysis of facial expressions.

Vision-Based Spatial Features

In the research conducted by Ashraf et al [74] and Lucey et al [75], features derived from the Active Appearance Model were input into support vector machine (SVM) classifiers for the purpose of frame-level pain recognition. In addition, they implemented pain detection at the sequence level by averaging the frame-level predictions. Gholami et al [76] used a Bayesian extension of SVM, known as the relevance vector machine, to

```
https://ai.jmir.org/2025/1/e53026
```

XSL•FO

differentiate between instances of pain and no pain in neonates. They also used this methodology to assess varying pain intensity levels. Meanwhile, Hammal et al [77] identified 4 levels of pain intensity through the use of log-normal filter-based features and an SVM classifier. Kaltwang et al [78] conducted a comparative study involving 3 separate methodologies. They used facial landmarks, discrete cosine transform, and LBP features to train 3 distinct relevance vector regression (RVR) models for estimating Prkachin and Solomon Pain Intensity. The best results were achieved by training an additional RVR model that consolidated the predictions from the 3 previously trained RVR models. The system [79] used a pyramid HOG for shape information and a pyramid LBP for appearance information, offering a more automated and objective approach to pain monitoring.

Pedersen [80] implementation used a 4-layer contractive autoencoder, along with SVM, which resulted in an effective pain detection system at the frame level. Egede et al [81] extracted features using both deep learning models and handcrafted methodologies. Facial landmarks, HOG, and deep vectors drawn from pretrained VGG-16 [82] and ResNet-50 [83] models were used. Rudovic et al [84] introduced a personalized federated deep learning technique for pain estimation derived from facial images. This approach involved using a compact convolutional neural network (CNN) architecture across various clients without the need to share their facial images. Contrary to the full sharing of model parameters, the personalized federated deep learning technique keeps the last layer localized. Hosseini et al [85] used a pretrained ResNet-18 model on the large emotion recognition dataset FER+ [86] and used transfer learning techniques to improve accuracy and performance. Huang et al [87] proposed a pain-awareness multistream CNN approach for feature extraction, focusing on specific regions most relevant to pain

expression instead of entire face images. Semwal and Londhe [88] proposed an Ensemble of Compact CNNs using 3 compact CNNs (variants of VGG, MobileNet, and GoogleNet) and integrating their predictions using the average ensemble rule. Kharghanian et al [89,90] developed a 4-layer convolutional deep belief network, trained as convolutional restricted Boltzmann machines to extract features. Semwal et al [91] introduced a novel fusion method for pain severity assessment in unconstrained environments using a decision-level fusion of 3 distinct features: data-driven red, green, blue color model (RGB) features, entropy-based texture features, and complementary features from both RGB and texture data. Using 3 CNNs (VGG-TL, ETNet, and DSCNN) with transfer learning, entropy texture network, and dual stream CNN, the model and various data augmentation techniques avoid overfitting and improve performance. The system demonstrates a 94% F_1 -score on a self-generated dataset from an unconstrained hospital setting.

Alghamdi and Alaghband [92] presented a facial expressions-based automatic pain assessment system using 2 concurrent subsystems that analyze both the full face and upper half of the face through pretrained CNNs, such as VGG16, InceptionV3, ResNet50, or ResNeXt50. Dai et al [93] developed a real-time pain detection system by mixing pain and emotion datasets for optimal real-time performance and conducting a cross-corpus test. The study experiments with both AU-based and non-AU-based methods, ultimately implementing the method on a robot for frozen shoulder therapy, thus emphasizing the need for balanced and ecologically valid pain datasets and the importance of real-world application and testing. Karamitsos et al [94] use the Haarcascade frontal face detector (OpenCV) for face detection; then, faces undergo gray scaling, histogram equalization, cropping, mean filtering, and normalization. The CNN is built upon a modified VGG16 architecture, achieving an impressive 92.5% accuracy. Barua et al [95] used a shutter blinds-based model inspired by spontaneous facial expressions and patch-based learning to achieve >95% accuracy in pain detection from facial images, leveraging transfer learning for efficient deep feature extraction. The model uniquely uses horizontal dynamic-sized patches, or "shutter blinds," to mine hidden facial signatures. Semwal et al [91] assess pain severity in unconstrained hospital environments using a decision-level fusion of 3 distinct types of features: data-driven RGB, entropy-based texture, and complementary features. They used 3 CNNs (VGG-CNN with transfer learning, entropy texture network, dual stream CNN) and various data augmentation techniques to avoid overfitting. The system demonstrates a 94.0% F_1 -score on a self-generated dataset from an unconstrained hospital setting.

Li et al [53] introduced a video-based infant monitoring system to analyze infant pain using 3 databases: Train-Data, Data-Clinic, and Data-YouTube. Using Fast Region-Based Convolutional Neural Network with object tracking and a hidden Markov model, the system precisely detects infant expressions and states. With a significant dataset from varied sources, including >16,000 images and real-world clinical videos, the approach offers enhanced accuracy and reliability in infant pain detection. Zamzmi et al [96] introduced a neonatal CNN that

```
https://ai.jmir.org/2025/1/e53026
```

uses a cascaded architecture with 3 convolutional branches. This design merges image-specific and general information for pain detection. The neonatal CNN demonstrated 91% accuracy and 0.93 AUC on the Neonatal Pain Assessment Dataset and 84.5% accuracy on the Infant Classification of Pain Expression dataset. Witherow et al [97] developed Facial Expressions Fusing Betamix Selected Landmark Features (FACE-BE-SELF), a novel deep adaptive method for adult-child facial expression classification. It fuses facial landmark data with deep feature representations, achieving domain-invariant classification. Using a unique mixture of beta distributions, facial features are selected based on expression, domain, and identity correlations. The FACE-BE-SELF method stands out by concurrently adapting adult-child domains, providing a unified expression representation for both groups. Compared to standard approaches, it surpasses in aligning latent representations of expressions across age groups.

Vision-Based Spatiotemporal Features

Bargshady et al [98] present an ensemble deep learning model that combines a 3-stream hybrid neural network with CNNs to extract facial features and classify pain levels. The VGG-Face, integrated with principal component analysis (PCA), is used for early feature extraction, while a 3-layer hybrid of CNN and bidirectional LSTM is developed for late fusion classification. This approach, tested on multiple pain databases, surpasses competing models with an accuracy of >89%. Sparse Autoencoders for Facial Expressions-Based Pain Assessment [57] reconstructs the upper part of the face from input images and then feeds both the original and reconstructed images into 2 concurrent and coupled InceptionV3 using Sparse Autoencoders. This dual-input approach emphasizes the upper facial features, essential for pain detection. By eliminating the need for conventional preprocessing steps such as face detection and adeptly handling varying head poses, Sparse Autoencoders for Facial Expressions-Based Pain Assessment offers enhanced performance and accuracy across multiple datasets, even in challenging profile views. Karamitsos et al [94] modified temporal convolutional network algorithm and processed facial features extracted from fine-tuned VGG-Face and PCA combined with hue, saturation, and value color spaces. The temporal convolutional network-based approach showcases faster performance and higher efficiency, achieving an accuracy of 92.44% and an AUC of 85%. Bargshady et al [99] propose an enhanced joint hybrid CNN-Bidirectional LSTM network model by leveraging a fine-tuned VGG-Face for feature extraction and apply PCA to focus on the most significant features, improving computational efficiency. These features are then classified by a CNN-Bidirectional LSTM network hybrid network into 4 levels of pain intensity.

The 3D CNNs have gained attention in several studies. Tavakolian and Hadid [100,101] created a 3D CNN that captures dynamic facial representations from videos and emphasizes the typical use of a fixed temporal kernel depth in research, which often misses capturing different time ranges. In the study by Huang et al [102], a hybrid network by combining 3D, 2D, and 1D CNNs has been introduced to extract spatiotemporal, spatial, and geometric features from image sequences. Wang et al [103] used the convolutional 3D network for pain expression

XSL•FO RenderX

recognition, which primarily uses a $3 \times 3 \times 3$ convolutional layer. However, this method often fails to capture the full spectrum of facial expression variations. To address this, they combined 3 distinct features: 3D CNN, HOG, and geometric features using support vector regression for pain estimation. They integrated the convolutional 3D network for spatiotemporal facial feature extraction and used the HOG in 2D images for geometric information to discern pain levels in facial expressions. De et al [104] present a deep learning architecture, the Decomposed Multiscale Spatiotemporal Network (DMSN). It uses 3 innovative blocks, DMSN-A, DMSN-B, and DMSN-C, to efficiently capture varied facial dynamics across conditions such as depression and pain. DMSN-A block focuses on pain, which might vary rapidly. It uses a sequence of $3 \times 1 \times 1$ temporal convolutions, capturing short to long temporal ranges. The studies by Granger and Cardinal [105] and Praveen et al [106] implemented weak-supervised domain adaptation, focusing on a shift from general affective expressions to specific pain expressions. Their framework used an inflated 3D CNN [107] with 3 convolutional layers and 3 inception modules, extracting both spatial and temporal data from videos.

Physiological Signals

Overview

While facial expressions are commonly used to identify pain, physiological signals are also a valuable modality for automatic pain detection. As detailed in the Pain Mechanism section, pain triggers changes in physiological signals, such as increased heart rate and skin conductivity, due to the activation of the SNS and peripheral nervous system [108]. Conversely, changes in physiological signals can indicate the presence of pain. However, extracting discriminative information from physiological signals is challenging. By contrast, they are objective indicators of pain because they cannot be artificially controlled [109], while exterior signals, such as facial expressions and gestures, may be unreliable, as individuals can deliberately disguise their behaviors. It makes physiological signals more reliable than exterior signals. In addition, physiological signals can be measured during daily life, while video and hand gestures can only be measured in laboratory settings. Thus, researchers have invested significant effort in exploring the feasibility of using physiological signals for pain assessment. Recent advances in sensor technology, signal processing, feature extraction, and ML algorithms are essential to the success of physiological signal-based automatic pain assessment.

This section provides a comprehensive review of the latest developments in pain detection approaches based on physiological signals. Four key components are exploited: (1) the use of modalities, (2) measurement devices, (3) feature extraction methods, and (4) ML models. The use of modalities refers to the type of physiological signals used for pain detection, including electroencephalogram, fMRI, electrocardiogram, and EDA. Measurement devices include both wearable and nonwearable devices, encompassing cardiac monitors, skin conductivity sensors, temperature sensors, accelerometers, and more. Feature extraction methods are techniques used to extract informative features from physiological signals, such as

```
https://ai.jmir.org/2025/1/e53026
```

XSL•F() RenderX time-domain features, frequency-domain features, and time-frequency features. Finally, ML models, such as SVM, artificial neural networks, and random forest (RF), are used to classify pain based on the extracted features.

Electroencephalogram as a Pain Indicator

Electroencephalography is a noninvasive technique widely used in the automatic detection of pain. The electrodes detect electrical activity and amplify it, producing a graphical of the representation brain activity over time. Electroencephalogram recordings typically show a series of waveforms or oscillations that are grouped into different frequency bands, such as delta, theta, alpha, beta, and gamma. These frequency bands have been associated with different mental states and cognitive functions. Various studies have shown the potential of electroencephalogram-based pain detection, and different approaches have been proposed to extract discriminative features from electroencephalogram signals for pain classification. For instance, Panavaranan et al [110] extracted the power spectral density of an electroencephalogram using fast Fourier transform and used SVM to classify thermal pain. Hadjileontiadis et al [54] proposed a novel approach that analyzes wavelet higher-order spectral features of an electroencephalogram to predict tonic cold pain. Vijayakumar et al [111] extracted time-frequency wavelet representations of independent components from electroencephalogram data and trained a RF model to classify pain levels, achieving an intrasubject accuracy of 93.26%.

The use of electroencephalogram techniques for pain detection has great potential to provide objective measures of pain, as these methods directly measure brain activity related to pain perception. However, these techniques also have limitations, including high cost, limited availability, and the need for specialized expertise for data analysis.

fMRI as a Pain Indicator

fMRI is a powerful neuroimaging tool that measures changes in blood flow within the brain as a proxy for neural activity. By measuring changes in the blood oxygen level–dependent signal, fMRI can indirectly map changes in neural activity in response to a specific stimulus, such as a painful stimulus.

The fMRI technique has been widely used in pain research, revealing a network of brain regions that are activated by painful stimuli. These regions include the primary and secondary somatosensory cortex, thalamus, insular cortex, and anterior cingulate cortex, among others. The activation of these regions is believed to be involved in the sensory and affective components of pain processing.

Activation of these regions is thought to be involved in the sensory discrimination aspects of pain processing. Thus, neuroimaging techniques allow us to visualize and quantify brain activities and then quantify pain. It is frequently used in the research of automatic pain assessment. Wager et al [112] used the least absolute shrinkage and selection operator ML regression algorithm to recognize induced heat pain by assessing the fMRI activity patterns. Shen et al [60] derived primary, dorsal, and ventral visual networks from blood oxygen level–dependent fMRI scans by using independent component

analysis and used a ML algorithm SVM to distinguish between patients with chronic low back pain and healthy volunteers and achieved an accuracy of 79.3%. Tu et al [59] proposed a novel sliced inverse regression-based fMRI decoding method to reduce the fMRI data dimension and showed overperformance compared to traditional regularization-based decoding analyses (principal component analysis and discriminant analysis, partial least squares-discriminant analysis, and least absolute shrinkage and selection operator). Robinson et al [58] scanned fMRI and applied ML algorithms to classify patients with fibromyalgia and healthy volunteers.

Electrocardiogram as a Pain Indicator

An electrocardiogram is a widely used technique to measure the electrical activity of the heart and its changes during each cardiac cycle. The electrocardiogram waveform consists of several characteristic waves and intervals that correspond to the different phases of the cardiac cycle, including the P wave, QRS complex, and T-wave. By analyzing the size, shape, and timing of these waves and intervals, a wide range of cardiac conditions, such as arrhythmias, heart attacks, and heart failure, can be diagnosed. The use of electrocardiograms in pain detection assumes that pain can cause a physiological stress response, leading to cardiovascular changes that are related to the pain stimuli. The autonomic nervous system responds to pain by increasing sympathetic tension and decreasing parasympathetic tension, leading to an increase in heart rate and blood pressure. By analyzing the electrocardiogram signal, features that reflect the autonomic nervous system status, such as heart rate variability (HRV), can be extracted and used to detect pain.

Several studies have shown the potential of electrocardiograms for pain detection. Walter et al [34] collected electrocardiogram data from 90 subjects using heat as pain stimuli and created the BioVid dataset, which also included skin conductance level, sEMG, and video data. Adjei et al [56] performed spectral analysis on electrocardiogram data and extracted HRV features, such as the low-frequency (LF) component and high-frequency (HF) component, which were significantly correlated with pain level. Jiang et al [64] extracted time-domain and frequency-domain HRV features from electrocardiogram data to classify pain level and obtained an AUC of 0.82 in the receiver operating characteristic curve.

However, there are also studies that suggest a lack of correlation between HRV and pain level. Meeuse et al [113] found no significant correlation between HRV features and heat pain level in their study. It is important to note that an electrocardiogram alone may not be sufficient to accurately detect pain, and other physiological signals, such as skin conductance and electromyography, may need to be considered as well. Furthermore, individual differences in pain perception and the variability of pain stimuli may affect the reliability of pain detection using an electrocardiogram.

EDA as a Pain Indicator

EDA, also referred to as galvanic skin response, is a physiological gauge of the skin's electrical conductance. This conductance changes according to the functioning of sweat glands within the skin [114]. The measurement of EDA is a noninvasive process involving the placement of 2 electrodes, often on the fingers or palms. Activation of the SNS, triggered by situations such as stress or pain, leads to increased sweat gland activity, causing a rise in the skin's electrical conductance.

Within the context of automated pain recognition, EDA serves as a valuable indicator due to its reflection of SNS activity [115], which is closely linked to the body's response to pain. Numerous research studies have highlighted EDA's potential in pain detection. For instance, in the BioVid dataset developed by Walter et al [34], EDA was used as one of the methods, revealing a correlation between EDA features and the intensity of pain.

sEMG is another important tool for measuring EDA in automatic pain detection. sEMG can measure the electrical activity of muscles and has been used to measure facial expression [116] or muscle movement of specific body parts, such as the back muscles [117]. These measures can provide additional information about the pain experience and may be used in combination with other modalities for better pain detection accuracy [118].

Devices

Data collection is indeed crucial in research, especially in statistical and ML-based studies. It is essential to ensure that the data collected are accurate, informative, and clean. However, selecting the right measurement devices is crucial for obtaining high-quality data.

Table 3 is a summary of previously used measurement devices in pain assessment studies. Figure 4 [115-117] presents 3 typical types of devices used in physiological signal-based pain assessment: wristband, headset, and chest band. The importance of wearable devices in this context cannot be overstated; they enable ubiquitous, real-time data collection [119,120], especially with the rise of body sensor networks. This technological advancement allows for extensive data gathering in wearable and remote settings, making continuous monitoring both feasible and affordable.



 Table 3. Physiological signal measurement devices used in pain assessment studies.

Device	Physiological signals	Connectivity	Туре	FDA ^a -cleared	Reference
Bioharness 3	Electrocardiogram	Bluetooth	Chest band	Yes	[64,69]
Affectiva Q sensor	EDA ^b	Bluetooth	Wristband	Yes	[68]
Procomp+	EDA and heart rate	Wired	Measurement hub	Yes	[121]
Emotive EPOC 14-channel elec- troencephalogram wireless record- ing headset	Electroencephalogram	Bluetooth	Headset	No	[54]
RespiBan	Respiration rate	Bluetooth	Chest band	No	[71]
Empatica E4	EDA, BVP ^c , and respiration rate	Wired	Wired sensor	Yes	[71]
Infiniti 3000A platform with Flex and Pro sensors	BVP, electrocardiogram, and EDA	Wired	Sensorhub	Yes	[65,67]
Polar RS800CX	HRV ^d	Wired	Watch	No	[122]

^aFDA: Food and Drug Administration.

^bEDA: electrodermal activity.

^cBVP: blood volume pulse.

^dHRV: heart rate variability.

Figure 4. Devices used in physiological signal-based pain assessment: WeBe band.



XSL•FO RenderX

There are several studies that have evaluated the usability and reliability of different measurement devices. Researchers can refer to these studies when choosing measurement devices for their own research. Ajayi et al [123] evaluated the Empatica E4 by comparing the results with nurse-recorded data and pooling questionnaires from participants. Nazari et al [124] tested the reliability of Bioharness and Fitbit measures of heart rate and activity at rest status. Rawstorn et al [125] evaluated the BioHarness by testing it on volunteers with both sinus rhythm and atrial fibrillation during simulated daily activities as well as low-, moderate-, and high-intensity exercises. Loberg et al [126] evaluated 4 different respiratory effort sensors and compared them with a respiratory sensor from NOX Medical as the golden reference device.

Feature Extraction

Overview

In the field of ML, pattern recognition, and image processing, feature extraction is a crucial step that involves transforming raw data into informative and nonredundant features to facilitate subsequent learning and generalization. Physiological signals typically carry implicit information that needs to be revealed through appropriate feature extraction techniques. While deep learning methods often generate features automatically, traditional ML methods require manual feature extraction.

For physiological signals, time window segmentation is commonly used to extract features. This involves segmenting the signals into chunks of equal time intervals and generating a row vector for each segment with 1 feature value for each feature, for example, the mean value of the segmentation. Physiological signal features can be classified into 4 categories: time-domain, frequency-domain, time-frequency-domain, and space-domain features.

Time-domain features describe the statistical and morphological properties of physiological signals, such as maximum value, SD, entropy, and mean R-R interval in electrocardiogram signals. Frequency-domain features characterize the spectral properties of signals, such as LF band power and low-high frequency ratio. Time-frequency-domain features consider both time-domain and frequency-domain properties simultaneously to account for the short duration and changing nature of physiological signals. Space-domain features, such as multispectral imaging and topography, are used to represent topographic characteristics of brain activity features, including electroencephalograms, fMRI, and fNIRS.

The complexity of physiological signals can guide feature selection. Signals with high stochastic stationarity and low signal-to-noise ratio, such as photoplethysmography and EDA, are considered low in complexity and can be represented by 1 or 2 feature domains. Signals with low stochastic stationarity and high signal-to-noise ratio, such as electrocardiogram, electroencephalogram, and fMRI, are high in complexity and require 3 to 4 feature domains to capture all relevant information. Nowadays, numerous Python libraries are available that facilitate the rapid extraction of features in physiological signals [127,128], electroencephalograms [129], video [130], and audio [131] domains. A summary of the commonly used features is presented in Table 4.



Fang et al

Table 4. Summary of the commonly used physiological signal features in pain assessment studies.

Cat	egory, feature, and description	Reference
HR	V ^a time-domain measures	[132]
	SD of NN ^b intervals	
	SD of RR ^c intervals	
	STD^{d} of the average NN intervals for each 5 min segment of a 24-hour HRV recording	
	Mean of the STD of all the NN intervals for 5-min segment of a 24-hour HRV recording	
	Percentage of successive RR intervals that differ by >50 ms	
	Average difference between the highest and lowest heart rates during each respiratory cycle	
	Root mean square of successive RR interval differences	
	Integral of the density of the RR interval histogram divided by its height	
	Baseline width of the RR interval histogram	
HR	V frequency-domain measures	[132]
	Absolute power of the ultra LF^e band (≤ 0.003 Hz)	
	Absolute power of the very-LF band (0.0033-0.04 Hz)	
	Peak frequency of the LF band (0.04-0.15 Hz)	
	Absolute power of the LF band (0.04-0.15 Hz)	
	Relative power of the LF band (0.04-0.15 Hz) in normal units	
	Relative power of the LF band (0.04-0.15 Hz)	
	Peak frequency of the HF ^f band (0.15-0.4 Hz)	
	Absolute power of the HF band (0.15-0.4 Hz)	
	Relative power of the HF band (0.15-0.4 Hz) in normal units	
	Relative power of the HF band (0.15-0.4 Hz)	
	Ratio of LF to HF power	
HR	V nonlinear measures	[132]
	Area of the ellipse that represents the total HRV	
	Poincare plot SD perpendicular to the line of identity	
	Poincare plot SD along the line of identity	
	Ratio of SD1 to SD2	
	Detrended fluctuation analysis, which describes short-term fluctuations	
	Detrended fluctuation analysis, which describes long-term fluctuations	
	Correlation dimension, which estimates the minimum number of variables required to construct a model of system dynamics	
Am	plitude	
	Peak amplitude	[133]
	Peak to peak amplitude	[133]
	Root mean square	[134]
	Mean absolute value	[134]
	Mean relative time of the peaks	[135]
	Mean relative time of the valleys	[135]
Vai	iability	
	IQR	[135]
	Range	[133]
	SD	[133]

https://ai.jmir.org/2025/1/e53026

XSL•FO RenderX JMIR AI 2025 | vol. 4 | e53026 | p.30 (page number not for citation purposes)

Category, feature, and description	Reference
Variance	[134]
Mean resting rate	[132]
Slope resting rate	[132]
Stationarity	
Integral degree of stationarity	[136]
Modified integral degree of stationarity	[136]
Modified mean degree of stationarity	[136]
Median	[133]
SD of SD vector	[133]
Entropy	
Approximate entropy	[137]
Fuzzy entropy	[138]
Sample entropy	[139]
Shannon entropy	[140]
Spectral entropy	[141]
Linearity	[133]
Lag dependence function	[136]
Population lag dependence function	[136]
Similarity	
Correlation coefficient	[142]
Median coherence	[143]
Mean coherence	[143]
Modified mean coherence	[143]
Modified integral of coherence	[143]
Mutual information	[144]
Frequency	
Bandwidth	[133]
Center frequency	[133]
Median frequency	[134]
Mean frequency	[134]
Mode frequency	[133]
Zero crossings	[134]

^aHRV: heart rate variability.
^bNN: neural network.
^cRR: 2 consecutive R waves.
^dSTD: SD.
^eLF: low-frequency.
^fHF: high-frequency.

Brain Activity Features

Physiological signals, including electroencephalograms, fMRI, and fNIRS, have unique characteristics that require specific feature extraction techniques. Electroencephalogram signals, for example, have high topological complexity as multiple channels are measuring simultaneously. They can be divided

```
https://ai.jmir.org/2025/1/e53026
```

XSL•FO RenderX into different frequency bands, such as delta, theta, alpha 1, alpha 2, beta 1, beta 2, gamma 1, and gamma 2. To assess pain, Panavaranan et al [110] used power spectral density features calculated using fast Fourier transform. Hadjileontiadis et al [54] combined continuous wavelet transform with higher-order statistics and spectra to create a new feature space for electroencephalograms. Rissacher et al [55] found temporal

Fang et al

JMIR AI

parietal alpha of electroencephalograms to be a useful feature for pain assessment.

In fMRI, Tu et al [59] proposed a novel dimension reduction method by incorporating singular value decomposition into sliced inverse regression to overcome the limitations of sliced inverse regression when dealing with high-dimensional data. This method was used to assess pain, achieving 77.61% binary classification accuracy.

There are various feature extraction approaches for electroencephalogram signals, as summarized by Behzadfar et al [145]. For brain activity signals in general, van der Miesen et al [146] outlined the state and progress in pain detection using these signals.

Electrocardiogram Features

Unlike general statistical feature extraction methods, electrocardiogram feature extraction involves more human experience on electrocardiograms and is more interpretable. Shaffer et al [132] provided an overview of HRV features, covering time-domain, frequency-domain, and non-linear measures. Time-domain and frequency-domain features are widely used in pain assessment studies. On the BioVid dataset, Werner et al [147] derived mean resting rate, root mean square of successive differences, and slope resting rate from the electrocardiogram signal. Gruss et al [148], Campbell et al [149], and Kachele et al [150] used the same 3 features in their studies. Kachele et al [150] also applied 4-level wavelet decomposition on detected R peaks to extract the mean alpha 1 coefficients. Jiang et al [64] extracted time-domain features, such as average interval between normal heart beats, SD of normal heart beat intervals, root mean square of successive differences, and percentage of successive RR intervals that differ by more than 20 ms, and frequency-domain features, such as LF, HF, and LF or HF, from an electrocardiogram and attained an AUC of 0.82 for induced electrical pain and an AUC of 0.75 for induced thermal pain.

Apart from HRV, other features have been used for various purposes. For instance, some studies have used morphological features, such as QRS complex duration and amplitude, T-wave amplitude, and ST-segment changes, for diagnosing cardiac abnormalities [150].

EDA and Electromyography Features

EDA and electromyography are critical tools in pain detection because they measure physiological responses that are directly linked to the autonomic nervous system's reactions, which vary significantly with pain perception [114,151]. Walter et al [133] systematically gathered and summarized feature extraction methods for EDA or electromyography signals from previous research and categorized them into mathematical groups of (1) amplitude, (2) frequency [152], (3) stationarity [136], (4) entropy [153], (5) linearity [144], and (6) variability. In total, 33 different features were listed, and their efficiency in pain assessment on the BioVid dataset was proved. Then, Gruss et al [148] deployed the feature table and derived it to 39 features. Campbell et al [149] also developed a feature list based on the study by Walter et al [133]. They also proposed a ML-based feature selection approach that deploys univariate feature selection and sequential forward selection for 100 epochs, with cross-validation as the metric to explore the optimal feature set. From their results, a relationship table between features and pain was displayed, illustrating the discriminative strength of features. In addition, amplitude, power, and unique functional features of electromyography signals are noted as useful in all different feature sets. Table 4 summarized the features used in previous studies.

Models

Overview

In the field of ML, the "no free lunch" theorem has been referred to often when talking about model selection [154]. This theorem illustrates that "any two optimization algorithms are equivalent when their performance is averaged across all possible problems," which implies that no single algorithm always has the best performance for all ML tasks. Thus, appropriate model selection is necessary for the success of ML-based pain assessment. In this section, we compare different ML algorithms by illustrating their advantages and disadvantages and their applicable scenarios. Table 5 provides a summary of the prevalent ML algorithms used in pain assessment.

 Table 5. Summary of the prevalent machine learning algorithms used in pain assessment studies.

Model	Advantages	Disadvantages		Reference
Support vector machine	Suitable for small datasetsTakes advantage of kernel functions	•	Low performance in multiclass tasks	[64,71]
Decision tree	Easily interpretableComputation friendly	•	High risk of overfitting Discards correlations between features	[155]
Random forest	 Applicable on large datasets Fixes the overfitting problem of decision tree Easy to parallelize 	•	Low performance on low-dimensional datasets Time consuming	[156,157]
Neural networks	High performance with large amounts of dataFlexible with layer configurations	•	Uninterpretable Computation consuming	[158,159]

SVM for Pain Classification

The first commonly used ML model in physiological signal-based automatic pain detection is SVM [64,71]. SVM is a type of generalized linear classifier that classifies data in a supervised learning way [160]. Its decision boundary is the maximum margin hyperplane for learning samples. SVM also includes kernel tricks, which makes it a substantially nonlinear classifier. The final decision of SVM only depends on the support vectors, which makes it suitable for small sample learning. On the contrary, SVM lacks the ability to provide restoration of variables to the formation of derived predictors [161], which is important in some areas such as financial prediction and health applications. In addition, SVM requires delicate preprocessing and tuning to acquire the best performance. Panavaranan et al [110] applied polynomial kernel SVM on electroencephalogram data and obtained an accuracy of 96.97%. Gruss et al [148] used SVM on the BioVid dataset and gained 90.94% accuracy on pain tolerance classification. In addition, Jiang et al [64] obtained an AUC of 0.82 with the use of SVM. More recently, Badura et al [71] achieved 94% accuracy using Gaussian kernel SVM.

Decision Tree for Pain Classification

Unlike SVM, decision tree is known for its interpretable characteristic. The decision tree algorithm is a method of approximating the value of a discrete function [162,163]. It is a typical classification method that uses an induction algorithm to generate readable rules and decision trees and then uses decision-making to analyze the new data. Essentially, a decision tree is a process of classifying data through a series of rules. Because of their inherent interpretability, tree-based algorithms help ML processes move beyond the "black box" model [164]. By contrast, due to the simple structure of tree-based models, overfitting easily happened on tree-based models [165]. Besides, they lack the ability to deal with missing data due to the continuity of tree structure.

RF for Pain Classification

RF is an algorithm that integrates multiple trees through the idea of ensemble learning. Its basic unit is a decision tree, and essentially, it belongs to a large branch of the ML "ensemble learning" method. Intuitively, each decision tree acts as a classifier, so for a given input sample, N decision trees will produce N classification results. RF integrates all classification voting results and designates the category with the most votes as the final output, which is a "bagging" idea. With the tree base and bagging theory RF holds, it has advantages such as preventing overfitting, easy to parallelize, and friendly with high-dimensional data [166]. In contrast, RFs require more time for training and prediction compared to decision trees. Vijayakumar et al [111] applied RF on 25 subjects' electroencephalogram data and obtained 89.45% accuracy. Naeini et al [167] used RF on the BioVid dataset and achieved an accuracy of 79%. Werner et al [168] used RF on their new "X- ITE" dataset and achieved 94.3% accuracy for phasic electrical pain classification.

Neural Networks for Pain Classification

NN have also been used by scholars for automatic pain detection [158,159]. NN abstracts the human brain neuron network from

```
https://ai.jmir.org/2025/1/e53026
```

the perspective of information processing, establishes a certain simple model, and composes different networks according to different connection structures. Thanks to the development of the digital society, the amount of data available for ML has grown substantially. NN, which can go deep in its layer structure, can reveal implicit information from data. Therefore, as the amount of data grows, the performance of NN keeps increasing, while traditional algorithms, such as SVM and RF, are limited. Nevertheless, NN has the defect of "black box" characteristic. Such uninterpretability keeps NN from blooming in certain fields, such as text and code analysis [169], judicial decision, and artificial intelligence medicine, because such fields understandable, interpretable require а clear, and decision-making process. Martinez et al [170] used NN on the BioVid dataset and obtained 82.75% accuracy on multitask classification. Jiang et al [69] applied an artificial neural network on 30 subjects and gained an average accuracy of 83.3%. The deviation of neural networks is widely used in automated pain assessment, such as CNN [156], RNN [171], and LSTM neural network [172].

Audio Analysis

Infant crying is a common sign of discomfort, hunger, or pain. It conveys information that helps caregivers assess the infant's emotional state and react appropriately. Crying analysis can be divided into two main stages: (1) the signal processing stage, which includes preprocessing the signal and extracting representative features; and (2) the classification stage. We classified the existing methods of signal processing stage into (1) time-domain methods; (2) frequency-domain methods; and (3) cepstral-domain methods.

Time-domain analysis is the analysis of a signal with respect to time (ie, the variation of a signal's amplitude over time). Linear prediction coding is one of the most common time-domain methods for analyzing sounds. The main concept behind linear prediction coding is the use of a linear combination of the past time-domain samples to predict the current time-domain sample. Other time-domain features that are commonly used for infants' sound analysis are energy, amplitude, and pause duration. Vempada et al [49] presented a time-domain method to detect discomfort-relevant cries. The proposed method was evaluated on a dataset consisting of 120 cry corpuses collected during pain (30 corpuses), hunger (60 corpuses), and wet diaper (30 corpuses). We want to note that the paper does not provide information about the stimulus that triggered the pain state or the data collection procedure. The infants' age ranges from 12 to 40 weeks. All corpses were recorded using a Sony digital recorder with a sampling rate of 44.1 kHz. In the feature extraction stage, two features were calculated: (1) short-time energy, which is the average of the square of the sample values in a suitable window; and (2) pause duration within the crying segment. Part of these features were used to build SVM, and the remaining features were used to evaluate its performance. The recognition performance of pain cry, hunger cry, and wet diaper cry were 83.33%, 27.78%, and 61.11%, respectively. The average recognition rate was 57.41%.

XSL•FO

Pupil Size

The measurement of changes in pupil size has been shown to be a promising physiological indicator of pain intensity. Pupil size can be used to monitor the effects of painful stimuli in the brain. The pupil dilates in response to pain due to the activation of the sympathetic branch, which releases norepinephrine, and the inhibition of the parasympathetic branch, which is responsible for constriction of the pupil. This section discusses the mechanism of using pupil dilation as a pain indicator and literature reviews of using pupil dilation for automated pain assessment.

The pupil dilation is a complex physiological response regulated automatically by 2 muscles in the eye, the sphincter pupillae and the dilator pupillae. The sphincter pupillae is controlled by the parasympathetic system to contract the pupil, while the dilator pupillae is dominated by the sympathetic system to dilate the pupil [50].

Höfle et al [51] investigated the influence of different luminance conditions on pupillometry for pain detection and found that the baseline pupil size values significantly differed under different luminance conditions, while the peak dilation remained the same. Bertrand et al [173] explored the influence of gender and anxiety on pupil dilation for pain detection and concluded that pupil dilation changes similarly in both men and women and are exacerbated in the presence of anxiety. Connelly et al Fang et al

[52] conducted an experiment on 30 children undergoing elective surgical correction of pectus excavatum and found that maximum pupil size, percent change in pupil size, and maximum constriction velocity were the most related features to pain intensity. Chapman et al [174] reported a delay of 1.25 seconds in 20 adult volunteers under noxious stimulation, while Eisenacha et al [175] reported a peak in pupil size with a lag of 4.25 seconds after the onset of heat pain on 28 adult volunteers. Wang et al [176] found that the pupillary response together with ML algorithms could be a promising method of objective pain level assessment by measuring pupillary response during induced cold pain on 32 subjects.

Multimodal Pain Detection

Including more modalities can possibly increase information density, which leads to increased accuracy. Thus, researchers have been increasingly turning to multimodal approaches to enhance the accuracy and reliability of automated pain assessment systems. These approaches combine information from multiple modalities, such as biomedical signals and facial expressions, to provide a more comprehensive understanding of the patient's pain experience. Furthermore, a multimodal approach can capture a more nuanced and diverse range of pain responses, which is particularly important given the wide variation in pain perception among individuals with different characteristics and cultural backgrounds. Figure 5 presents a typical flow of multimodal pain assessment.

Figure 5. Multimodal pain assessment.



Fusion strategies commonly used in multimodal pain assessment can be categorized into early fusion and late fusion. Early fusion involves the combination of features from different modalities before the training of a classifier, while late or decision fusion combines the predictions of individual classifiers after training. Common methods of combining predictions include fixed methods such as taking the mean or product and trainable methods such as using a pseudoinverse. Figure 6 illustrates the early and late fusion strategies. Some research has explored combining early and decision fusion by merging specific features at the feature level and then fusing those with other features at the decision level [46].

Figure 6. Fusion strategies.



(B) Late fusion

The first study to combine video and physiological signals for automated pain detection was conducted by Werner et al [147], who used an early fusion strategy to concatenate features from both modalities. The optimal fusion set is found to be the combination of all video and physiological signals, achieving accuracies of 80.6% and 77.8% for person-specific and generic classifiers, respectively, in detecting baseline and highest tolerable pain using a RF ensemble–based classifier. Kachele et al [177] applied both early and late fusion strategies using SVM with linear kernel and RF for recognizing baseline and

highest tolerable pain, achieving accuracies of 68.2% and 76.6% for early and late fusion, respectively.

Continuing the BioVid dataset, Kachele et al [178] applies early and late fusion techniques with new features included, achieving slightly better results with late fusion (83.1%) than early fusion (82.7%). Thiam et al [179] proposed a hierarchical fusion architecture that divides multimodal data into 3 subsets. These subsets are used for the first layer of RF training, followed by pseudo-inverse mapping, multilayer perceptron mapping, and a final layer that combines both pseudo-inverse and multilayer perceptron fusion mapping. Kessler et al [180] took advantage of the fusion strategy proposed by Thiam et al [179] and applied it to remote photoplethysmography.

Other studies focus on incorporating additional modalities, such as audio. Velana et al [37] published the SenseEmotion database, which captures video, physiological signals, and audio for the first time. Thiam et al [181] merged features from video, physiological signal, and audio data on the SenseEmotion dataset, exploring different data fusion strategies, including early fusion, group late fusion, and individual late fusion. Results show that individual late fusion outperforms other strategies slightly on leave-subject-out experiment, while group late fusion slightly outperforms on user-specific task. There is also a dataset for neonatal pain assessment that includes video, audio, and physiological signals [46,171].

Recent studies have explored new fusion approaches. Bellmann et al [182] proposed a dominant channel fusion approach that identifies the most relevant input channel and combines it with the remaining channels to create an ensemble of classifiers. Bellman et al [183] proposed a novel late fusion approach that combines a mixture of experts and stacked generalization approaches and is assessed on different datasets involving the biophysiological modalities electromyography, electrocardiogram, and EDA. Thiam et al [159] proposed an information theoretic approach that uses a deep denoising convolutional autoencoder to learn and aggregate latent representations based on each input channel.

However, it is evident that late fusion, using multiple models as part of an ensemble learning approach, requires significantly more computational power and storage space compared to early fusion methods. As pain assessment is an emerging field, the current focus is predominantly on enhancing predictive accuracy rather than on resource use, and discussions on model complexity are relatively scarce. However, with the advent of Tiny ML and the rise of edge computing [184], running large models on microprocessors becomes challenging. Consequently, early fusion might gain popularity on edge devices, where the ability to run simpler, more compact models efficiently is crucial. This shift could make early and lightweight fusion approaches more viable and preferred in scenarios where computational resources are limited. In addition, with the increasing inclusion of multimodal data, we can envisage future fusion methods potentially incorporating recently developed self-attention algorithms [185].

Discussion

The pain assessment field is faced with several challenges and opportunities for future development. This section will focus on 3 areas of concern—data, ML techniques, and ethical considerations—and then propose future research directions.

Data

Automatic pain assessment is challenged by the limited availability of clinical pain data, as most studies have focused on experimental or induced pain. Widely used datasets such as BioVid, BP4D+, and X-ITE are collected from healthy volunteers and use external thermal or electrical pain. These studies are conducted under consistent experimental conditions that differ from real-world scenarios. Furthermore, induced pain has different mechanisms than disease pain, which encompasses different types of pain, such as nociceptive and central pain. Therefore, it is important to test models trained on experimental data using clinical pain data. In addition, more clinical pain data should be collected to facilitate the development of automatic pain assessment models and enable their use in clinical trials.

Pupil dilation has been identified as a promising indicator of brain activity and pain levels. However, in previous studies, pain was often used as the stimuli for measuring brain activity, rather than the focus of the study. Consequently, only a few studies have directly correlated pupil dilation with pain levels. A potential research direction is to include pupil dilation in the automatic pain assessment modality family. Pupil dilation has been shown to be effective in affective computing, with datasets such as the MAHNOB-HCI and SEED containing eye-tracking data that demonstrate the contribution of pupil data to arousal detection. As pain can also be regarded as physiological arousal, transferring pupil dilation to automatic pain assessment studies is a worthwhile area of research.

Personalization of Pain Responses

In the following subsection, we explore personalized pain detection, focusing on the considerable differences in pain experiences among individuals. Pain perception varies widely due to a mix of biological factors and social-psychological influences. These differences are shaped by demographics such as gender, age, and ethnicity, which are linked to varying rates of chronic pain. In addition, factors such as genetic predispositions and psychological processes also significantly impact pain responses, whether in clinical settings or experimental scenarios. Importantly, these elements interact in complex ways, crafting the unique pain experiences of everyone. Research has highlighted that genetic markers associated with pain can differ across genders and ethnicities and interact with psychological aspects such as stress, affecting pain perception. These myriads of interacting factors culminates in a distinctive set of influences for each person's experience of pain [186].

Jiang et al [187] introduced a method that enhances pain assessment by incorporating personalized features. They used ML to analyze individual pain data, enabling the model to tailor its predictions to each patient's unique physiological and psychological characteristics. This approach improves the accuracy of pain management by adapting to personal pain

```
https://ai.jmir.org/2025/1/e53026
```
profiles. Casti et al [188] developed a platform to improve pain diagnosis by leveraging personalized data. Using a combination of visual, speech, and physiological indicators, they used ML techniques to tailor assessments to individual patient profiles, enhancing the precision and effectiveness of pain management strategies. Martinez et al [189] proposed a method to refine pain estimation by integrating personalized features. They used ML to analyze individual facial expressions, allowing the model to adjust its predictions based on each person's unique facial expressiveness score. This approach enhances the accuracy of Visual Analog Scale estimations by adapting to individual pain profiles [189].

Most papers on personalized pain assessment claim personalization at the model level, focusing on enhancing ML models to suit individualized approaches or using ML techniques to delve deeper into databases for extracting personalized information to improve predictions. The predominant reliance on public databases for research is evident, as most researchers use these readily available datasets. This reliance restricts personalization efforts to the data provided by these databases, making highly tailored training challenging. In addition, most pain-related datasets globally are derived from experiments involving artificially induced pain, which must pass rigorous ethical or clinical trial reviews, further limiting the quantity of available data. Looking to the future, personalization will undoubtedly be a crucial focus. It is foreseeable that researchers will collect more personalized data during experiments, including variables such as personality traits and ethnicity. This will likely lead to the generation of more nuanced datasets that include varied physiological responses to different pain stimuli, enhancing the granularity and effectiveness of personalized pain management solutions.

Real-Time Pain Detection

Building on our earlier discussion about the personalization of pain responses, it is essential to delve into another critically relevant clinical application: real-time monitoring [190]. The goal of such monitoring is not just to detect pain but to enable timely and effective interventions that can significantly enhance patient outcomes. Real-time monitoring of pain becomes particularly crucial in postoperative care, where accurately gauging a patient's pain levels is vital for adjusting analgesic dosages. This not only helps in managing the pain effectively but also minimizes the risk of both undermedication and overmedication, which can lead to complications such as opioid dependency or inadequate pain relief. In ICUs, the stakes are even higher. Many patients in ICUs are unable to communicate due to their conditions or sedation, making verbal reports of pain unreliable. Here, real-time monitoring systems can play a transformative role by continuously tracking pain indicators through physiological signals such as heart rate, blood pressure, and facial expressions. These data can then be analyzed to provide a dynamic, real-time assessment of pain, informing caregivers when an intervention is necessary. Moreover, real-time monitoring integrates seamlessly with the concept of personalized pain management. By continuously collecting and analyzing data specific to each patient, health care providers can tailor their interventions more precisely to the individual's pain profile and response to treatment. This approach not only

XSL•FO

improves the quality of care but also enhances patient comfort and satisfaction. As technology advances, the potential for real-time pain monitoring grows. Innovations in wearable technology, ML algorithms, and data integration are paving the way for even more accurate and responsive pain management systems. These systems promise to transform how pain is managed in health care settings, making care more proactive, patient centered, and effective.

In the academic sphere, the development of real-time pain monitoring is primarily concentrated on 2 aspects: improving model efficiency to enable fast judgments suitable for real-time applications and developing practical tools such as wearable and mobile apps to facilitate devices widespread implementation. Enhancing the processing speed of models involves not only maintaining accuracy but also integrating advanced ML technologies, such as deep learning. Meanwhile, the development of tools such as wearables and mobile apps allows for the noninvasive collection of physiological data and real-time analysis, helping patients and health care providers to promptly assess pain levels and treatment effectiveness. This combination of improved models and practical tools is driving pain management toward more precise, personalized, and proactive solutions. Kong et al [191] introduced a smartphone app that enhances real-time pain detection using EDA signals collected from a wrist-worn device. They tested the app with thermal grill and electrical pulse data, demonstrating high accuracy in pain detection with a RF model. This approach offers a practical solution for objective, near-real-time pain assessment in everyday settings. Dai et al [93] address automatic pain detection using a mix of pain and emotion datasets to enhance model robustness, achieving 88.4% accuracy. They criticize CNNs for overfitting on biased data and validate their method through experiments on a humanoid robot in physiotherapy, emphasizing the importance of real-time, real-world testing and assessing the system's practical utility and accuracy.

In summary, the advancement of real-time pain monitoring represents a significant enhancement in health care, enabling precise and timely interventions that are tailored to the unique needs of each patient. This technology not only improves the accuracy of pain assessments but also enriches the quality of care by integrating cutting-edge ML models and wearable technologies. As this field continues to evolve, it holds the promise of transforming pain management into a more responsive, personalized, and patient-centered practice.

ML Techniques

Although deep learning has revolutionized computer vision and physiological signal analysis, traditional ML algorithms still dominate the field of physiological signal–based automatic pain assessment. One possible reason for this is that deep learning requires extensive data, which is time consuming and resource intensive to collect. Therefore, studies often include only a small number of participants, typically in the tens, making it difficult to gather comprehensive datasets.

In this context, transfer learning, a prominent topic in artificial intelligence, offers a promising alternative solution. Transfer learning involves applying knowledge gained from a source

domain to a new target domain, which can be particularly useful in scenarios where data collection is challenging. Differing data distributions between the source and target domains can lead to performance degradation if models are applied directly. Transfer learning helps bridge this gap, ensuring better model performance across different settings [192].

Kächele et al [193] proposed an adaptive confidence learning method for personalizing pain intensity estimation systems, demonstrating the efficacy of transfer learning in this field. Feature extraction involved specific preprocessing steps for each signal type, such as bandpass filtering and artifact correction for electromyography. A multistage ensemble classifier was applied to learn the confidence of a regression system. This method involved selecting confident samples from unlabeled data of the test participants to iteratively adapt the model. Their experiments showed that the adaptive learning approach significantly improved the performance of pain intensity estimation.

Chen et al [194] implemented "TrAdaboost," a transfer learning algorithm, to improve facial expression recognition, including pain expressions. They used the PAINFUL database, which contains video sequences of 25 patients with shoulder injuries, encompassing 48,398 frames of spontaneous pain expressions. The primary challenge addressed was the variability in pain expressions across different individuals. They proposed an inductive transfer learning algorithm to develop person-specific models. This algorithm first trains a set of weak classifiers on source data from multiple subjects and then selects the most relevant classifiers for the target subject. Experimental results showed that inductive transfer learning significantly improved pain detection accuracy. For example, the AUC for pain detection increased from 0.769 to 0.782 with just 10 target samples and reached 0.891 with 100 samples. Furthermore, this approach drastically reduced training time compared to traditional methods, making it feasible for rapid retraining in clinical settings.

While traditional ML remains prevalent in automatic pain assessment due to data constraints, transfer learning presents a viable alternative. It addresses the challenges associated with varying data distributions and limited dataset sizes, enhancing model robustness and performance. Future research should explore the potential of transfer learning algorithms further, integrating them into clinical practice to improve pain management outcomes.

Ethical Considerations

Automatic pain assessment raises several ethical concerns that need to be addressed. One primary concern is the privacy and security of patients' health data. The use of physiological signals, such as facial expressions, speech patterns, and pupil dilation, to assess pain levels can lead to the collection of sensitive health data. Therefore, it is essential to ensure that the data collected are secure and protected from unauthorized access.

Another ethical consideration is the potential for bias in automatic pain assessment models. ML models are only as good as the data they are trained on, and if the training data are biased,

```
https://ai.jmir.org/2025/1/e53026
```

the model will be biased too. Bias can result in inaccurate pain assessment, leading to inadequate pain management and, in some cases, even harm to patients. Therefore, it is crucial to ensure that the data used to train the models are representative and unbiased.

Future Directions

Automated pain assessment has made significant strides in recent years, leveraging technological advancements and data-driven approaches to enhance the accuracy and efficiency of pain detection. However, several promising directions for future research remain unexplored. Addressing these areas could lead to the development of more sophisticated and reliable automated pain assessment systems.

First, integrating data from various sources, such as pupil dilation, voice analysis, and body movement, could offer a more comprehensive understanding of pain. This requires a more comprehensive, clinical, and clean database to be released. Second, exploring novel deep learning architectures, including transformer-based models and generative adversarial networks, may yield improved performance in pain assessment tasks. These architectures could capture intricate patterns and dependencies within pain-related data, leading to enhanced predictive capabilities. Third, collaboration with health care professionals is crucial to validate the effectiveness and reliability of automated pain assessment systems in real-world clinical settings. Integrating these systems into clinical workflows could provide valuable insights and assist health care providers in making informed decisions. Finally, using transfer learning can provide new insights. In scenarios where large, annotated datasets are scarce, exploring transfer learning techniques and methods to adapt models to smaller datasets could prove beneficial. These approaches could enable the development of accurate pain assessment models even with limited training data.

Conclusions

This survey reviewed the current advancements in automated pain assessment using ML techniques. Traditional pain assessment methods, reliant on self-reports and observational scales, face significant limitations, particularly for patients who are noncommunicative. We explored various modalities for automated pain detection, including facial expressions, physiological signals, audio, and pupil dilation. While each modality has its strengths, combining multiple modalities can enhance accuracy but also introduces challenges in data fusion and model complexity. Despite progress, challenges remain, such as the scarcity of diverse clinical pain datasets and ethical concerns regarding patient privacy. Personalized pain assessment models are also necessary due to variability in pain perception across populations. Future research should focus on developing more robust algorithms and leveraging deep learning and transfer learning. Collaborative efforts to create comprehensive pain datasets are crucial, as is integrating real-time pain monitoring into clinical practice. In summary, automated pain assessment has the potential to transform pain management. Continued interdisciplinary research and collaboration are key to overcoming current challenges and fully realizing these technologies' benefits.

XSL•FO RenderX

Acknowledgments

RF was responsible for writing the Abstract and Introduction sections on physiological signals and pupil size, the multimodal study, the Discussion and Conclusions sections, and organizing and formatting the paper. EH was responsible for writing the Facial Expression section. RZ was responsible for writing the Pain Mechanism and Electrodermal Activity sections. SR was responsible for collecting information, reviewing, and final editing. HH was responsible for reviewing and funding acquisition.

Conflicts of Interest

None declared.

Multimedia Appendix 1 Summary of studies table. [PDF File (Adobe PDF File), 139 KB - ai v4i1e53026 app1.pdf]

References

- 1. Merskey H. The definition of pain. Eur Psychiatr 2020 Apr 16;6(4):153-159. [doi: 10.1017/s092493380000256x]
- Williams AC, Craig KD. Updating the definition of pain. Pain 2016 Nov 18;157(11):2420-2423. [doi: 10.1097/j.pain.0000000000613] [Medline: 27200490]
- 3. Yong RJ, Mullins PM, Bhattacharyya N. Prevalence of chronic pain among adults in the United States. Pain 2022 Feb 01;163(2):e328-e332. [doi: 10.1097/j.pain.0000000002291] [Medline: 33990113]
- 4. Gaskin DJ, Richard P. The economic costs of pain in the United States. J Pain 2012 Aug;13(8):715-724 [FREE Full text] [doi: 10.1016/j.jpain.2012.03.009] [Medline: 22607834]
- 5. Manchikanti L, Helm S, Fellows B, Janata JW, Pampati V, Grider JS, et al. Opioid epidemic in the United States. Pain Physician 2012 Jul;15(3 Suppl):ES9-E38 [FREE Full text] [doi: 10.36076/ppj.2012/15/es9] [Medline: 22786464]
- 6. Fink R. Pain assessment: the cornerstone to optimal pain management. Proc (Bayl Univ Med Cent) 2000 Jul 11;13(3):236-239 [FREE Full text] [doi: 10.1080/08998280.2000.11927681] [Medline: 16389388]
- Gracely RH, McGrath P, Dubner R. Ratio scales of sensory and affective verbal pain descriptors. Pain 1978 Jun;5(1):5-18. [doi: <u>10.1016/0304-3959(78)90020-9</u>] [Medline: <u>673440</u>]
- McCormack HM, Horne DJ, Sheather S. Clinical applications of visual analogue scales: a critical review. Psychol Med 1988 Nov 09;18(4):1007-1019. [doi: <u>10.1017/s0033291700009934</u>] [Medline: <u>3078045</u>]
- 9. Downie WW, Leatham PA, Rhind VM, Wright V, Branco JA, Anderson JA. Studies with pain rating scales. Ann Rheum Dis 1978 Aug 01;37(4):378-381 [FREE Full text] [doi: 10.1136/ard.37.4.378] [Medline: <u>686873</u>]
- 10. Wong DL, Baker CM. Smiling faces as anchor for pain intensity scales. Pain 2001 Jan;89(2-3):295-300. [doi: 10.1016/s0304-3959(00)00375-4] [Medline: 11291631]
- Dehghani H, Tavangar H, Ghandehari A. Validity and reliability of behavioral pain scale in patients with low level of consciousness due to head trauma hospitalized in intensive care unit. Arch Trauma Res 2014 Mar 30;3(1):e18608 [FREE Full text] [doi: 10.5812/atr.18608] [Medline: 25032173]
- 12. Warden V, Hurley AC, Volicer L. Development and psychometric evaluation of the Pain Assessment in Advanced Dementia (PAINAD) scale. J Am Med Dir Assoc 2003;4(1):9-15. [doi: 10.1097/01.JAM.0000043422.31640.F7] [Medline: 12807591]
- Lawrence J, Alcock D, McGrath P, Kay J, MacMurray SB, Dulberg C. The development of a tool to assess neonatal pain. Neonatal Netw 1993 Sep;12(6):59-66. [Medline: <u>8413140</u>]
- 14. Kappesser J, de C Williams AC. Pain estimation: asking the right questions. Pain 2010 Feb;148(2):184-187. [doi: 10.1016/j.pain.2009.10.007] [Medline: 19880252]
- Merskey H. The taxonomy of pain. Med Clin North Am 2007 Jan;91(1):13-20, vii. [doi: <u>10.1016/j.mcna.2006.10.009</u>] [Medline: <u>17164101</u>]
- Gorczyca R, Filip R, Walczak E. Psychological aspects of pain. Ann Agric Environ Med 2013;Spec no. 1:23-27 [FREE Full text] [Medline: 25000837]
- 17. Garland EL. Pain processing in the human nervous system: a selective review of nociceptive and biobehavioral pathways. Prim Care 2012 Sep;39(3):561-571 [FREE Full text] [doi: 10.1016/j.pop.2012.06.013] [Medline: 22958566]
- Council NR, Criado A. Recognition and alleviation of pain in laboratory animals. Lab Anim 2010 Oct 01;44(4):380. [doi: 10.1258/LA.2010.201003]
- 19. Kandel ER, Schwartz JH, Jessell TM. Principles Of Neural Science. Volume 4. New York, NY: McGrawhill; 2000.
- 20. Julius D, Basbaum AI. Molecular mechanisms of nociception. Nature 2001 Sep 13;413(6852):203-210. [doi: 10.1038/35093019] [Medline: 11557989]
- Treede RD, Rief W, Barke A, Aziz Q, Bennett MI, Benoliel R, et al. A classification of chronic pain for ICD-11. Pain 2015 Jun;156(6):1003-1007 [FREE Full text] [doi: 10.1097/j.pain.000000000000160] [Medline: 25844555]

- 22. Markenson JA. Mechanisms of chronic pain. Am J Med 1996 Jul 31;101(1A):6S-18S [FREE Full text] [doi: 10.1016/s0002-9343(96)00133-7] [Medline: 8764755]
- 23. Borsook D. A future without chronic pain: neuroscience and clinical research. Cerebrum 2012 May;2012:7 [FREE Full text] [Medline: 23447793]
- 24. Mee S, Bunney BG, Reist C, Potkin SG, Bunney WE. Psychological pain: a review of evidence. J Psychiatr Res 2006 Dec;40(8):680-690. [doi: 10.1016/j.jpsychires.2006.03.003] [Medline: 16725157]
- 25. Bair MJ, Robinson RL, Katon W, Kroenke K. Depression and pain comorbidity: a literature review. Arch Intern Med 2003 Nov 10;163(20):2433-2445. [doi: 10.1001/archinte.163.20.2433] [Medline: 14609780]
- Von Korff M, Simon G. The relationship between pain and depression. Br J Psychiatry Suppl 1996 Jun;1688(30):101-108. [doi: <u>10.1192/s0007125000298474</u>] [Medline: <u>8864155</u>]
- 27. Engel GL. Psychogenic pain and the pain-prone patient. Am J Med 1959 Jun;26(6):899-918. [doi: 10.1016/0002-9343(59)90212-8] [Medline: 13649716]
- 28. Bassler M, Krauthauser H, Hoffmann SO. Inpatient psychotherapy with chronic psychogenic pain patients. Psychother Psychosom Med Psychol 1994;44(9-10):299-307. [Medline: <u>7972647</u>]
- 29. Paxton SL. Clinical uses of TENS. A survey of physical therapists. Phys Ther 1980 Jan;60(1):38-44. [doi: 10.1093/ptj/60.1.38] [Medline: 6965323]
- Ziemssen T, Kern S. Psychoneuroimmunology--cross-talk between the immune and nervous systems. J Neurol 2007 May;254 Suppl 2(S2):II8-II1. [doi: <u>10.1007/s00415-007-2003-8</u>] [Medline: <u>17503136</u>]
- 31. Teff KL. Visceral nerves: vagal and sympathetic innervation. JPEN J Parenter Enteral Nutr 2008 Sep;32(5):569-571. [doi: 10.1177/0148607108321705] [Medline: 18753395]
- 32. Singaram S, Ramakrishnan K, Selvam J, Senthil M, Narayanamurthy V. Sweat gland morphology and physiology in diabetes, neuropathy, and nephropathy: a review. Arch Physiol Biochem 2024 Aug 05;130(4):437-451. [doi: 10.1080/13813455.2022.2114499] [Medline: 36063413]
- Lucey P, Cohn JF, Prkachin KM, Solomon PE, Matthews I. Painful data: the UNBC-McMaster shoulder pain expression archive database. In: Proceedings of the 2011 IEEE International Conference on Automatic Face & Gesture Recognition. 2011 Presented at: FG '11; March 21-25, 2011; Santa Barbara, CA p. 57-64 URL: <u>https://ieeexplore.ieee.org/document/5771462</u> [doi: 10.1109/fg.2011.5771462]
- 34. Walter S, Gruss S, Ehleiter H, Tan J, Traue HC, Werner P, et al. The biovid heat pain database data for the advancement and systematic validation of an automated pain recognition system. In: Proceedings of the 2013 IEEE International Conference on Cybernetics. 2013 Presented at: CYBCO '13; June 13-15, 2013; Lausanne, Switzerland p. 128-131 URL: <u>https://ieeexplore.ieee.org/document/6617456</u> [doi: 10.1109/cybconf.2013.6617456]
- 35. Haque MA, Bautista RB, Noroozi F, Kulkarni K, Laursen CB, Irani R, et al. Deep multimodal pain recognition: a database and comparison of spatio-temporal visual modalities. In: Proceedings of the 13th IEEE International Conference on Automatic Face & Gesture Recognition. 2018 Presented at: FG '18; May 15-19, 2018; Xi'an, China p. 250-257 URL: <u>https://ieeexplore.ieee.org/document/8373837</u> [doi: 10.1109/fg.2018.00044]
- 36. Aung MS, Kaltwang S, Romera-Paredes B, Martinez B, Singh A, Cella M, et al. The automatic detection of chronic pain-related expression: requirements, challenges and the multimodal emopain dataset. IEEE Trans Affective Comput 2016 Oct 1;7(4):435-451. [doi: 10.1109/taffc.2015.2462830]
- 37. Velana M, Gruss S, Layher G, Thiam P, Zhang Y, Schork D, et al. The SenseEmotion database: a multimodal database for the development and systematic validation of an automatic pain- and emotion-recognition system. In: Proceedings of the 4th IAPR TC 9 Workshop on Multimodal Pattern Recognition of Social Signals in Human-Computer-Interaction. 2016 Presented at: MPRSS '16; December 4, 2016; Cancun, Mexico p. 127-139 URL: <u>https://link.springer.com/chapter/10.1007/</u> 978-3-319-59259-6_11 [doi: 10.1007/978-3-319-59259-6_11]
- 38. Gruss S, Geiger M, Werner P, Wilhelm O, Traue HC, Al-Hamadi A, et al. Multi-modal signals for analyzing pain responses to thermal and electrical stimuli. J Vis Exp 2019 Apr 05(146). [doi: <u>10.3791/59057</u>] [Medline: <u>31009005</u>]
- 39. Zhang X, Yin L, Cohn JF, Canavan S, Reale M, Horowitz A, et al. BP4D-Spontaneous: a high-resolution spontaneous 3D dynamic facial expression database. Image Vis Comput 2014 Oct;32(10):692-706. [doi: 10.1016/j.imavis.2014.06.002]
- Zhang Z, Girard JM, Wu Y, Zhang X, Liu P, Ciftci U. Multimodal spontaneous emotion corpus for human behavior analysis. In: Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition. 2016 Presented at: CVPR '16; June 27-30, 2016; Las Vegas, NV p. 3438-3446 URL: <u>https://ieeexplore.ieee.org/abstract/document/7780743</u> [doi: 10.1109/cvpr.2016.374]
- 41. Brahnam S, Chuang CF, Shih FY, Slack MR. SVM classification of neonatal facial images of pain. In: Proceedings of the 6th International Workshop on Fuzzy Logic and Applications. 2005 Presented at: WILF '05; September 15-17, 2005; Crema, Italy p. 128 URL: <u>https://link.springer.com/chapter/10.1007/11676935_15</u> [doi: <u>10.1007/11676935_15</u>]
- 42. Harrison D, Sampson M, Reszel J, Abdulla K, Barrowman N, Cumber J, et al. Too many crying babies: a systematic review of pain management practices during immunizations on YouTube. BMC Pediatr 2014 May 29;14(1):134 [FREE Full text] [doi: 10.1186/1471-2431-14-134] [Medline: 24885559]
- 43. Egede J, Valstar M, Torres MT, Sharkey D. Automatic neonatal pain estimation: an acute pain in Neonates database. In: Proceedings of the 8th International Conference on Affective Computing and Intelligent Interaction. 2019 Presented at:

ACII '19; September 3-6, 2019; Cambridge, UK p. 1-7 URL: <u>https://ieeexplore.ieee.org/document/8925480</u> [doi: 10.1109/acii.2019.8925480]

- 44. Zamzmi G, Pai CY, Goldgof D, Kasturi R, Ashmeade T, Sun Y. A comprehensive and context-sensitive neonatal pain assessment using computer vision. IEEE Trans Affective Comput 2022 Jan 1;13(1):28-45. [doi: <u>10.1109/taffc.2019.2926710</u>]
- 45. Brahnam S, Nanni L, McMurtrey S, Lumini A, Brattin R, Slack M, et al. Neonatal pain detection in videos using the iCOPEvid dataset and an ensemble of descriptors extracted from gaussian of local descriptors. Appl Comput Inform 2020 Jul 17;19(1/2):122-143 [FREE Full text] [doi: 10.1016/j.aci.2019.05.003]
- 46. Salekin MS, Zamzmi G, Hausmann J, Goldgof D, Kasturi R, Kneusel M, et al. Multimodal neonatal procedural and postoperative pain assessment dataset. Data Brief 2021 Apr;35:106796 [FREE Full text] [doi: 10.1016/j.dib.2021.106796] [Medline: <u>33644268</u>]
- 47. Salekin MS, Zamzmi G, Goldgof D, Kasturi R, Hoppe T, Sun Y. First investigation into the use of deep learning for continuous assessment of neonatal postoperative pain. In: Proceedings of the 15th IEEE International Conference on Automatic Face and Gesture Recognition. 2020 Presented at: FG '20; November 16-20, 2020; Buenos Aires, Argentina p. 415-419 URL: <u>https://ieeexplore.ieee.org/document/9320233</u> [doi: 10.1109/fg47880.2020.00082]
- 48. Ekman P, Friesen WV. Facial Action Coding System: Investigator's Guide. Palo Alto, CA: Consulting Psychologists Press; 1978.
- 49. Rao KS, Koolagudi SG, Vempada RR. Emotion recognition from speech using global and local prosodic features. Int J Speech Technol 2012 Aug 4;16(2):143-160. [doi: 10.1007/s10772-012-9172-2]
- Zambach SA, Cai C, Helms HC, Hald BO, Dong Y, Fordsmann JC, et al. Precapillary sphincters and pericytes at first-order capillaries as key regulators for brain capillary perfusion. Proc Natl Acad Sci U S A 2021 Jun 29;118(26):e2023749118 [FREE Full text] [doi: 10.1073/pnas.2023749118] [Medline: 34155102]
- 51. Höfle M, Kenntner-Mabiala R, Pauli P, Alpers GW. You can see pain in the eye: pupillometry as an index of pain intensity under different luminance conditions. Int J Psychophysiol 2008 Dec;70(3):171-175. [doi: 10.1016/j.ijpsycho.2008.06.008] [Medline: 18644409]
- Connelly MA, Brown JT, Kearns GL, Anderson RA, St Peter SD, Neville KA. Pupillometry: a non-invasive technique for pain assessment in paediatric patients. Arch Dis Child 2014 Dec 03;99(12):1125-1131 [FREE Full text] [doi: 10.1136/archdischild-2014-306286] [Medline: 25187497]
- Li C, Pourtaherian A, van Onzenoort L, Ten WE, de With PH. Infant facial expression analysis: towards a real-time video monitoring system using R-CNN and HMM. IEEE J Biomed Health Inform 2021 May;25(5):1429-1440. [doi: 10.1109/JBHI.2020.3037031] [Medline: 33170787]
- 54. Hadjileontiadis LJ. EEG-based tonic cold pain characterization using wavelet higher order spectral features. IEEE Trans Biomed Eng 2015 Aug;62(8):1981-1991. [doi: 10.1109/TBME.2015.2409133] [Medline: 25769141]
- 55. Rissacher D, Dowman R, Schuckers SA. Identifying frequency-domain features for an EEG-based pain measurement system. In: Proceedings of the 33rd Annual Northeast Bioengineering Conference. 2007 Presented at: NEBC '07; March 10-11, 2007; Stony Brook, NY p. 114-115 URL: <u>https://ieeexplore.ieee.org/document/4413305</u> [doi: 10.1109/nebc.2007.4413305]
- 56. Adjei T, Von Rosenberg W, Goverdovsky V, Powezka K, Jaffer U, Mandic DP. Pain prediction from ECG in vascular surgery. IEEE J Transl Eng Health Med 2017;5:2800310 [FREE Full text] [doi: 10.1109/JTEHM.2017.2734647] [Medline: 29026686]
- 57. Alghamdi T, Alaghband G. SAFEPA: an expandable multi-pose facial expressions pain assessment method. Applied Sciences 2023 Jun 16;13(12):7206. [doi: <u>10.3390/app13127206</u>]
- Robinson ME, O'Shea AM, Craggs JG, Price DD, Letzen JE, Staud R. Comparison of machine classification algorithms for fibromyalgia: neuroimages versus self-report. J Pain 2015 May;16(5):472-477 [FREE Full text] [doi: 10.1016/j.jpain.2015.02.002] [Medline: 25704840]
- 59. Tu Y, Fu Z, Tan A, Huang G, Hu L, Hung Y, et al. A novel and effective fMRI decoding approach based on sliced inverse regression and its application to pain prediction. Neurocomputing 2018 Jan;273:373-384. [doi: 10.1016/j.neucom.2017.07.045]
- 60. Shen W, Tu Y, Gollub RL, Ortiz A, Napadow V, Yu S, et al. Visual network alterations in brain functional connectivity in chronic low back pain: a resting state functional connectivity and machine learning study. Neuroimage Clin 2019;22:101775 [FREE Full text] [doi: 10.1016/j.nicl.2019.101775] [Medline: 30927604]
- 61. Karunakaran KD, Peng K, Berry D, Green S, Labadie R, Kussman B, et al. NIRS measures in pain and analgesia: fundamentals, features, and function. Neurosci Biobehav Rev 2021 Jan;120:335-353. [doi: <u>10.1016/j.neubiorev.2020.10.023</u>] [Medline: <u>33159918</u>]
- 62. Fernandez Rojas R, Huang X, Ou KL. A machine learning approach for the identification of a biomarker of human pain using fNIRS. Sci Rep 2019 Apr 04;9(1):5645 [FREE Full text] [doi: 10.1038/s41598-019-42098-w] [Medline: 30948760]
- 63. Electroencephalogram (EEG). Johns Hopkins Medicine. URL: <u>https://www.hopkinsmedicine.org/health/</u> <u>treatment-tests-and-therapies/</u> <u>electroencephalogram-eeg#:~:text=An%20EEG%20is%20a%20test,activity%20of%20your%20brain%20cells</u> [accessed 2024-04-29]

- 64. Jiang M, Mieronkoski R, Rahmani AM, Hagelberg N, Salanterä S, Liljeberg P. Ultra-short-term analysis of heart rate variability for real-time acute pain monitoring with wearable electronics. In: Proceedings of the 2017 IEEE International Conference on Bioinformatics and Biomedicine. 2017 Presented at: BIBM '17; November 13-16, 2017; Kansas City, MO p. 1025-1032 URL: <u>https://ieeexplore.ieee.org/document/8217798</u> [doi: 10.1109/bibm.2017.8217798]
- 65. Chu Y, Zhao X, Yao J, Zhao Y, Wu Z. Physiological signals based quantitative evaluation method of the pain. IFAC Proc Vol 2014;47(3):2981-2986. [doi: 10.3182/20140824-6-za-1003.01420]
- 66. Werner P, Al-Hamadi A, Niese R, Walter S, Gruss S, Traue HC. Towards pain monitoring: facial expression, head pose, a new database, an automatic system and remaining challenges. In: Proceedings of the 2013 Conference on British Machine Vision. 2013 Presented at: BMVC '13; September 9-13, 2013; Bristol, UK p. 1-13 URL: <u>https://citeseerx.ist.psu.edu/</u> document?repid=rep1&type=pdf&doi=03f075e95638bc66e687badd97a58c5de67e58e6 [doi: <u>10.5244/c.27.119</u>]
- Chu Y, Zhao X, Han J, Su Y. Physiological signal-based method for measurement of pain intensity. Front Neurosci 2017 May 26;11:279 [FREE Full text] [doi: 10.3389/fnins.2017.00279] [Medline: 28603478]
- 68. Susam BT, Akcakaya M, Nezamfar H, Diaz D, Xu XL, de Sa VR, et al. Automated pain assessment using electrodermal activity data and machine learning. Annu Int Conf IEEE Eng Med Biol Soc 2018 Jul;2018:372-375 [FREE Full text] [doi: 10.1109/EMBC.2018.8512389] [Medline: 30440413]
- 69. Jiang M, Mieronkoski R, Syrjälä E, Anzanpour A, Terävä V, Rahmani AM, et al. Acute pain intensity monitoring with the classification of multiple physiological parameters. J Clin Monit Comput 2019 Jun 26;33(3):493-507 [FREE Full text] [doi: 10.1007/s10877-018-0174-8] [Medline: 29946994]
- 70. Mark JN, Hu Y, Luk K. ICA-based ECG removal from surface electromyography and its effect on low back pain assessment. In: Proceedings of the 3rd International IEEE/EMBS Conference on Neural Engineering. 2007 Presented at: CNE '07; May 2-5, 2007; Kohala Coast, HI p. 646-649 URL: <u>https://ieeexplore.ieee.org/document/4227361</u> [doi: <u>10.1109/cne.2007.369756</u>]
- Badura A, Masłowska A, Myśliwiec A, Piętka E. Multimodal signal analysis for pain recognition in physiotherapy using wavelet scattering transform. Sensors (Basel) 2021 Feb 12;21(4):1311 [FREE Full text] [doi: 10.3390/s21041311] [Medline: 33673097]
- 72. Prkachin KM, Solomon PE. The structure, reliability and validity of pain expression: evidence from patients with shoulder pain. Pain 2008 Oct 15;139(2):267-274. [doi: 10.1016/j.pain.2008.04.010] [Medline: 18502049]
- 73. Williams AC. Facial expression of pain: an evolutionary account. Behav Brain Sci 2002 Aug 11;25(4):439-455. [doi: 10.1017/s0140525x02000080] [Medline: 12879700]
- 74. Ashraf AB, Lucey S, Cohn JF, Chen T, Ambadar Z, Prkachin KM, et al. The painful face pain expression recognition using active appearance models. Image Vis Comput 2009 Oct;27(12):1788-1796 [FREE Full text] [doi: 10.1016/j.imavis.2009.05.007] [Medline: 22837587]
- 75. Lucey P, Cohn JF, Matthews I, Lucey S, Sridharan S, Howlett J, et al. Automatically detecting pain in video through facial action units. IEEE Trans Syst Man Cybern B Cybern 2011 Jun;41(3):664-674 [FREE Full text] [doi: 10.1109/TSMCB.2010.2082525] [Medline: 21097382]
- 76. Gholami B, Haddad WM, Tannenbaum AR. Relevance vector machine learning for neonate pain intensity assessment using digital imaging. IEEE Trans Biomed Eng 2010 Jun;57(6):1457-1466 [FREE Full text] [doi: 10.1109/TBME.2009.2039214] [Medline: 20172803]
- 77. Hammal Z, Cohn JF. Automatic detection of pain intensity. Proc ACM Int Conf Multimodal Interact 2012 Oct;2012:47-52 [FREE Full text] [doi: 10.1145/2388676.2388688] [Medline: 32724903]
- 78. Kaltwang S, Rudovic O, Pantic M. Continuous pain intensity estimation from facial expressions. In: Proceedings of the 8th International Symposium Conference on Advances in Visual Computing. 2012 Presented at: ISVC '12; July 16-18, 2012; Crete, Greece p. 368-377 URL: <u>https://link.springer.com/chapter/10.1007/978-3-642-33191-6_36</u> [doi: 10.1007/978-3-642-33191-6_36]
- 79. Khan RA, Meyer A, Konik H, Bouakaz S. Pain detection through shape and appearance features. In: Proceedings of the 2013 IEEE International Conference on Multimedia and Expo. 2013 Presented at: ICME '13; July 15-19, 2013; San Jose, CA p. 1-6 URL: <u>https://ieeexplore.ieee.org/document/6607608</u> [doi: <u>10.1109/icme.2013.6607608</u>]
- Pedersen H. Learning appearance features for pain detection using the UNBC-McMaster shoulder pain expression archive database. In: Proceedings of the 10th International Conference on Computer Vision Systems. 2015 Presented at: ICVS '15; July 6-9, 2015; Copenhagen, Denmark p. 10-36 URL: <u>https://dl.acm.org/doi/10.1007/978-3-319-20904-3_12</u> [doi: 10.1007/978-3-319-20904-3_12]
- 81. Egede JO, Song S, Olugbade TA, Wang C, Williams AC, Meng G, et al. EMOPAIN challenge 2020: multimodal pain evaluation from facial and bodily expressions. In: Proceedings of the 15th IEEE International Conference on Automatic Face and Gesture Recognition. 2020 Presented at: FG' 20; November 16-20, 2020; Buenos Aires, Argentina p. 849-856 URL: <u>https://dl.acm.org/doi/10.1109/FG47880.2020.00078</u> [doi: <u>10.1109/fg47880.2020.00078</u>]
- 82. Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition. arXiv Preprint posted online September 4, 2014 [FREE Full text]
- He K, Zhang X, Rennke S, Sun J. Deep residual learning for image recognition. In: Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition. 2016 Presented at: CVPR '16; June 27-30, 2016; Las Vegas, NV p. 770-778 URL: <u>https://ieeexplore.ieee.org/document/7780459</u> [doi: <u>10.1109/cvpr.2016.90</u>]

```
https://ai.jmir.org/2025/1/e53026
```

- 84. Rudovic O, Tobis N, Kaltwang S, Schuller B, Rueckert D, Cohn JF, et al. Personalized federated deep learning for pain estimation from face images. arXiv Preprint posted online January 12, 2021 [FREE Full text]
- Hosseini E, Fang R, Zhang R, Chuah CN, Orooji M, Rafatirad S, et al. Convolution neural network for pain intensity assessment from facial expression. Annu Int Conf IEEE Eng Med Biol Soc 2022 Jul;2022:2697-2702. [doi: 10.1109/EMBC48229.2022.9871770] [Medline: 36085712]
- 86. Barsoum E, Zhang C, Ferrer CC, Zhang Z. Training deep networks for facial expression recognition with crowd-sourced label distribution. In: Proceedings of the 18th ACM International Conference on Multimodal Interaction. 2016 Presented at: ICMI '16; November 12-16, 2016; Tokyo, Japan p. 278-283 URL: <u>https://dl.acm.org/doi/10.1145/2993148.2993165</u> [doi: 10.1145/2993148.2993165]
- 87. Huang D, Xia Z, Li L, Wang K, Feng X. Pain-awareness multistream convolutional neural network for pain estimation. J Electron Imag 2019 Jul 1;28(04):1. [doi: 10.1117/1.jei.28.4.043008]
- 88. Semwal A, Londhe ND. ECCNet: an ensemble of compact convolution neural network for pain severity assessment from face images. In: Proceedings of the 11th International Conference on Cloud Computing, Data Science & Engineering. 2021 Presented at: Confluence '21; January 28-29, 2021; Noida, India p. 761-766 URL: <u>https://ieeexplore.ieee.org/document/</u> <u>9377197</u> [doi: 10.1109/confluence51648.2021.9377197]
- 89. Kharghanian R, Peiravi A, Moradi F. Pain detection from facial images using unsupervised feature learning approach. In: Proceedings of the 38th Annual International Conference of the IEEE Engineering in Medicine and Biology Society. 2016 Presented at: EMBC '16; August 16-20, 2016; Orlando, FL p. 419-422 URL: <u>https://ieeexplore.ieee.org/document/7590729</u> [doi: 10.1109/embc.2016.7590729]
- 90. Kharghanian R, Peiravi A, Moradi F, Iosifidis A. Pain detection using batch normalized discriminant restricted Boltzmann machine layers. J Vis Commun Image Represen 2021 Apr;76:103062. [doi: <u>10.1016/j.jvcir.2021.103062</u>]
- 91. Semwal A, Londhe ND. MVFNet: a multi-view fusion network for pain intensity assessment in unconstrained environment. Biomed Signal Process Control 2021 May;67:102537. [doi: <u>10.1016/j.bspc.2021.102537</u>]
- 92. Alghamdi T, Alaghband G. Facial expressions based automatic pain assessment system. Appl Sci 2022 Jun 24;12(13):6423. [doi: 10.3390/app12136423]
- 93. Dai L, Broekens J, Truong KP. Real-time pain detection in facial expressions for health robotics. In: Proceedings of the 8th International Conference on Affective Computing and Intelligent Interaction Workshops and Demos. 2019 Presented at: ACIIW '19; September 3-6, 2019; Cambridge, UK p. 277-283 URL: <u>https://ieeexplore.ieee.org/document/8925192</u> [doi: 10.1109/aciiw.2019.8925192]
- 94. Karamitsos I, Seladji I, Modak S. A modified CNN network for automatic pain identification using facial expressions. J Softw Eng Appl 2021;14(08):400-417. [doi: 10.4236/jsea.2021.148024]
- 95. Barua PD, Baygin N, Dogan S, Baygin M, Arunkumar N, Fujita H, et al. Automated detection of pain levels using deep feature extraction from shutter blinds-based dynamic-sized horizontal patches with facial images. Sci Rep 2022 Oct 14;12(1):17297 [FREE Full text] [doi: 10.1038/s41598-022-21380-4] [Medline: 36241674]
- 96. Zamzmi G, Paul R, Goldgof D, Kasturi R, Sun Y. Pain assessment from facial expression: neonatal convolutional neural network (N-CNN). In: Proceedings of the 2019 International Joint Conference on Neural Networks. 2019 Presented at: IJCNN '19; July 14-19, 2019; Budapest, Hungary p. 1-7 URL: <u>https://ieeexplore.ieee.org/document/8851879</u> [doi: 10.1109/ijcnn.2019.8851879]
- 97. Witherow MA, Samad MD, Diawara N, Bar HY, Iftekharuddin KM. Deep adaptation of adult-child facial expressions by fusing landmark features. IEEE Trans Affective Comput 2024 Jul;15(3):847-858. [doi: <u>10.1109/taffc.2023.3297075</u>]
- 98. Bargshady G, Zhou X, Deo RC, Soar J, Whittaker F, Wang H. Ensemble neural network approach detecting pain intensity from facial expressions. Artif Intell Med 2020 Sep;109:101954. [doi: <u>10.1016/j.artmed.2020.101954</u>] [Medline: <u>34756219</u>]
- 99. Bargshady G, Zhou X, Deo RC, Soar J, Whittaker F, Wang H. Enhanced deep learning algorithm development to detect pain intensity from facial expression images. Expert Syst Appl 2020 Jul;149:113305. [doi: 10.1016/j.eswa.2020.113305]
- 100. Tavakolian M, Hadid A. Deep spatiotemporal representation of the face for automatic pain intensity estimation. In: Proceedings of the 24th International Conference on Pattern Recognition. 2018 Presented at: ICPR '18; August 20-24, 2018; Beijing, China p. 350-354 URL: <u>https://ieeexplore.ieee.org/document/8545324</u> [doi: <u>10.1109/icpr.2018.8545324</u>]
- Tavakolian M, Hadid A. A spatiotemporal convolutional neural network for automatic pain intensity estimation from facial dynamics. Int J Comput Vis 2019 Jun 25;127(10):1413-1425. [doi: <u>10.1007/s11263-019-01191-3</u>]
- 102. Huang Y, Qing L, Xu S, Wang L, Peng Y. HybNet: a hybrid network structure for pain intensity estimation. Vis Comput 2021 Feb 04;38(3):871-882. [doi: <u>10.1007/s00371-021-02056-y</u>]
- Wang J, Sun H. Pain intensity estimation using deep spatiotemporal and handcrafted features. IEICE Trans Inf Syst 2018;E101.D(6):1572-1580. [doi: <u>10.1587/transinf.2017edp7318</u>]
- 104. de Melo WC, Granger E, Lopez MB. Facial expression analysis using decomposed multiscale spatiotemporal networks. Expert Syst Appl 2024 Feb;236:121276. [doi: <u>10.1016/j.eswa.2023.121276</u>]
- 105. Granger E, Cardinal P, Praveen RG. Deep domain adaptation for ordinal regression of pain intensity estimation using weakly-labelled videos. arXiv Preprint posted online August 13, 2020 [FREE Full text]
- 106. Praveen RG, Granger E, Cardinal P. Deep weakly supervised domain adaptation for pain localization in videos. In: Proceedings of the 15th IEEE International Conference on Automatic Face and Gesture Recognition. 2020 Presented at:

FG '20; November 16-20, 2020; Buenos Aires, Argentina p. 473-480 URL: <u>https://ieeexplore.ieee.org/document/9320216</u> [doi: <u>10.1109/fg47880.2020.00139</u>]

- 107. Carreira J, Zisserman A. Quo vadis, action recognition? A new model and the kinetics dataset. In: Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition. 2017 Presented at: CVPR '17; July 21-26, 2017; Honolulu, HI p. 4724-4733 URL: <u>https://ieeexplore.ieee.org/document/8099985</u> [doi: <u>10.1109/cvpr.2017.502</u>]
- 108. Shu L, Xie J, Yang M, Li Z, Li Z, Liao D, et al. A review of emotion recognition using physiological signals. Sensors (Basel) 2018 Jun 28;18(7):2074 [FREE Full text] [doi: 10.3390/s18072074] [Medline: 29958457]
- 109. Li W, Zhang Z, Song A. Physiological-signal-based emotion recognition: an odyssey from methodology to philosophy. Measurement 2021 Feb;172:108747. [doi: 10.1016/j.measurement.2020.108747]
- 110. Panavaranan P, Wongsawat Y. EEG-based pain estimation via fuzzy logic and polynomial kernel support vector machine. In: Proceedings of the 2013 Biomedical Engineering International Conference. 2013 Presented at: BMEiCon '13; October 23-25, 2013; Amphur Muang, Thailand p. 1-4 URL: <u>https://ieeexplore.ieee.org/document/6687668</u> [doi: 10.1109/bmeicon.2013.6687668]
- 111. Vijayakumar V, Case M, Shirinpour S, He B. Quantifying and characterizing tonic thermal pain across subjects from EEG data using random forest models. IEEE Trans Biomed Eng 2017 Dec;64(12):2988-2996 [FREE Full text] [doi: 10.1109/TBME.2017.2756870] [Medline: 28952933]
- 112. Wager TD, Atlas LY, Lindquist MA, Roy M, Woo C, Kross E. An fMRI-based neurologic signature of physical pain. N Engl J Med 2013 Apr 11;368(15):1388-1397 [FREE Full text] [doi: <u>10.1056/NEJMoa1204471</u>] [Medline: <u>23574118</u>]
- 113. Meeuse JJ, Löwik MS, Löwik SA, Aarden E, van Roon AM, Gans RO, et al. Heart rate variability parameters do not correlate with pain intensity in healthy volunteers. Pain Med 2013 Aug 01;14(8):1192-1201. [doi: <u>10.1111/pme.12133</u>] [Medline: <u>23659489</u>]
- 114. Hosseini E, Fang R, Zhang R, Rafatirad S, Homayoun H. Emotion and stress recognition utilizing galvanic skin response and wearable technology: a real-time approach for mental health care. In: Proceedings of the 2023 IEEE International Conference on Bioinformatics and Biomedicine. 2023 Presented at: BIBM '23; December 5-8, 2023; Istanbul, Turkey p. 1125-1131 URL: <u>https://www.computer.org/csdl/proceedings-article/bibm/2023/10386049/1TObUqDKemQ</u> [doi: 10.1109/bibm58861.2023.10386049]
- 115. Hosseini E, Fang R, Zhang R, Parenteau A, Hang S, Rafatirad S. A low cost EDA-based stress detection using machine learning. In: Proceedings of the 2022 IEEE International Conference on Bioinformatics and Biomedicine. 2022 Presented at: BIBM '22; December 6-8, 2022; Las Vegas, NV p. 2619-2623 URL: <u>https://ieeexplore.ieee.org/document/9995093</u> [doi: 10.1109/bibm55620.2022.9995093]
- 116. Merletti R, Farina D. Surface Electromyography: Physiology, Engineering, and Applications. Hoboken, NJ: John Wiley & Sons; 2016.
- 117. Srinivasan J, Balasubramanian V. Low back pain and muscle fatigue due to road cycling—an sEMG study. J Bodyw Mov Ther 2007 Jul;11(3):260-266. [doi: 10.1016/j.jbmt.2006.08.009]
- 118. Jiang M, Rahmani AM, Westerlund T, Liljeberg P, Tenhunen H. Facial expression recognition with sEMG method. In: Proceedings of the 2015 IEEE International Conference on Computer and Information Technology; Ubiquitous Computing and Communications; Dependable, Autonomic and Secure Computing; Pervasive Intelligence and Computing, 2015 Presented at: IUCC '15; October 26-28, 2015; Liverpool, UK p. 981-988 URL: <u>https://ieeexplore.ieee.org/document/7363189</u> [doi: 10.1109/cit/iucc/dasc/picom.2015.148]
- 119. Zhang Z, Zhang R, Chang CW, Guo Y, Chi YW, Pan T. iWRAP: a theranostic wearable device with real-time vital monitoring and auto-adjustable compression level for venous thromboembolism. IEEE Trans Biomed Eng 2021 Sep;68(9):2776-2786. [doi: 10.1109/TBME.2021.3054335] [Medline: <u>33493109</u>]
- 120. Zhang R, Fang C, Homayoun H, Berk GG. Privee: a wearable for real-time bladder monitoring system. In: Proceedings of the Adjunct Proceedings of the 2023 ACM International Joint Conference on Pervasive and Ubiquitous Computing & the 2023 ACM International Symposium on Wearable Computing. 2023 Presented at: UbiComp/ISWC '23; October 8-12, 2023; Cancun, Mexico p. 291-295 URL: <u>https://dl.acm.org/doi/10.1145/3594739.3610782</u> [doi: <u>10.1145/3594739.3610782</u>]
- 121. Loggia ML, Juneau M, Bushnell CM. Autonomic responses to heat pain: heart rate, skin conductance, and their relation to verbal ratings and stimulus intensity. Pain 2011 Mar;152(3):592-598. [doi: 10.1016/j.pain.2010.11.032] [Medline: 21215519]
- 122. Hautala AJ, Karppinen J, Seppanen T. Short-term assessment of autonomic nervous system as a potential tool to quantify pain experience. Annu Int Conf IEEE Eng Med Biol Soc 2016 Aug;2016:2684-2687. [doi: <u>10.1109/EMBC.2016.7591283</u>] [Medline: <u>28268874</u>]
- 123. Ajayi TA, Salongo L, Zang Y, Wineinger N, Steinhubl S. Mobile health-collected biophysical markers in children with serious illness-related pain. J Palliat Med 2021 Apr 01;24(4):580-588 [FREE Full text] [doi: 10.1089/jpm.2020.0234] [Medline: <u>33351729</u>]
- 124. Nazari G, MacDermid JC, Sinden KE, Richardson J, Tang A. Reliability of Zephyr bioharness and Fitbit charge measures of heart rate and activity at rest, during the modified Canadian aerobic fitness test, and recovery. J Strength Cond Res 2019 Feb;33(2):559-571. [doi: 10.1519/JSC.000000000001842] [Medline: 30689619]

- 125. Rawstorn JC, Gant N, Warren I, Doughty RN, Lever N, Poppe KK, et al. Measurement and data transmission validity of a multi-biosensor system for real-time remote exercise monitoring among cardiac patients. JMIR Rehabil Assist Technol 2015 Mar 20;2(1):e2 [FREE Full text] [doi: 10.2196/rehab.3633] [Medline: 28582235]
- 126. Løberg F, Goebel V, Plagemann T. Quantifying the signal quality of low-cost respiratory effort sensors for sleep apnea monitoring. In: Proceedings of the 3rd International Workshop on Multimedia for Personal Health and Health Care. 2018 Presented at: HealthMedia '18; October 22, 2018; Seoul, Republic of Korea p. 3-11 URL: <u>https://dl.acm.org/doi/10.1145/ 3264996.3264998</u> [doi: <u>10.1145/3264996.3264998</u>]
- 127. Fang R, Zhang R, Hosseini E, Fang C, Rafatirad S, Homayoun H. Introducing an open-source Python toolkit for machine learning research in physiological signal based affective computing. In: Proceedings of the 2023 IEEE International Conference on Bioinformatics and Biomedicine. 2023 Presented at: BIBM '23; December 5-8, 2023; Istanbul, Turkiye p. 1890-1894 URL: <u>https://ieeexplore.ieee.org/document/10385965</u> [doi: <u>10.1109/bibm58861.2023.10385965</u>]
- 128. Makowski D, Pham T, Lau ZJ, Brammer JC, Lespinasse F, Pham H, et al. NeuroKit2: a Python toolbox for neurophysiological signal processing. Behav Res Methods 2021 Aug;53(4):1689-1696. [doi: 10.3758/s13428-020-01516-y] [Medline: 33528817]
- 129. Cabañero-Gomez L, Hervas R, Gonzalez I, Rodriguez-Benitez L. eeglib: a Python module for EEG feature extraction. SoftwareX 2021 Jul;15:100745. [doi: 10.1016/j.softx.2021.100745]
- 130. Iashin V, Korbar B, Georgievski B, Hoppe J. v-iashin / video_features. GitHub. URL: <u>https://github.com/v-iashin/video_features</u> [accessed 2024-04-29]
- Lenain R, Weston J, Shivkumar A, Fristed E. Surfboard: audio feature extraction for modern machine learning. arXiv Preprint posted online May 18, 2020 [FREE Full text] [doi: 10.21437/interspeech.2020-2879]
- Shaffer F, Ginsberg JP. An overview of heart rate variability metrics and norms. Front Public Health 2017 Sep 28;5:258
 [FREE Full text] [doi: 10.3389/fpubh.2017.00258] [Medline: 29034226]
- 133. Walter S, Gruss S, Limbrecht-Ecklundt K, Traue HC, Werner P, Al-Hamadi A, et al. Automatic pain quantification using autonomic parameters. Psychol Neurosci 2014;7(3):363-380. [doi: 10.3922/j.psns.2014.041]
- 134. Phinyomark A, Phukpattaranont P, Limsakul C. Feature reduction and selection for EMG signal classification. Expert Syst Appl 2012 Jun;39(8):7420-7431. [doi: 10.1016/j.eswa.2012.01.102]
- 135. Phinyomark A, Scheme E. An investigation of temporally inspired time domain features for electromyographic pattern recognition. In: Proceedings of the 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society. 2018 Presented at: EMBC '18; July 18-21, 2018; Honolulu, HI p. 5236-5240 URL: <u>https://ieeexplore.ieee.org/ document/8513427</u> [doi: 10.1109/embc.2018.8513427]
- 136. Cao C, Slobounov S. Application of a novel measure of EEG non-stationarity as 'Shannon- entropy of the peak frequency shifting' for detecting residual abnormalities in concussed individuals. Clin Neurophysiol 2011 Jul;122(7):1314-1321 [FREE Full text] [doi: 10.1016/j.clinph.2010.12.042] [Medline: 21216191]
- Pincus SM. Approximate entropy as a measure of system complexity. Proc Natl Acad Sci U S A 1991 Mar 15;88(6):2297-2301 [FREE Full text] [doi: 10.1073/pnas.88.6.2297] [Medline: 11607165]
- 138. Kosko B. Fuzzy entropy and conditioning. Inf Sci 1986 Dec;40(2):165-174. [doi: 10.1016/0020-0255(86)90006-X]
- Richman JS, Moorman JR. Physiological time-series analysis using approximate entropy and sample entropy. Am J Physiol Heart Circ Physiol 2000 Jun;278(6):H2039-H2049 [FREE Full text] [doi: 10.1152/ajpheart.2000.278.6.H2039] [Medline: 10843903]
- 140. Lin J. Divergence measures based on the Shannon entropy. IEEE Trans Inform Theory 1991;37(1):145-151. [doi: 10.1109/18.61115]
- 141. Zhang A, Yang B, Huang L. Feature extraction of EEG signals using power spectral entropy. In: Proceedings of the 2008 International Conference on BioMedical Engineering and Informatics. 2008 Presented at: BMEI '08; May 27-30, 2008; Sanya, China p. 435-439 URL: <u>https://ieeexplore.ieee.org/document/4549210</u> [doi: <u>10.1109/bmei.2008.254</u>]
- 142. Kennedy HL. A new statistical measure of signal similarity. In: Proceedings of the 2007 Conference on Information, Decision and Control, Adelaide. 2007 Presented at: IDC '07; February 12-14, 2007; Adelaide, Australia p. 112-117 URL: <u>https://ieeexplore.ieee.org/document/4252487</u> [doi: 10.1109/idc.2007.374535]
- 143. Dukic S, Iyer PM, Mohr K, Hardiman O, Lalor EC, Nasseroleslami B. Estimation of coherence using the median is robust against EEG artefacts. Annu Int Conf IEEE Eng Med Biol Soc 2017 Jul;2017:3949-3952. [doi: <u>10.1109/EMBC.2017.8037720</u>] [Medline: <u>29060761</u>]
- 144. Chen HM, Varshney PK, Arora MK. Performance of mutual information similarity measure for registration of multitemporal remote sensing images. IEEE Trans Geosci Remote Sensing 2003 Nov;41(11):2445-2454. [doi: 10.1109/tgrs.2003.817664]
- 145. Behzadfar N. A brief overview on analysis and feature extraction of electroencephalogram signals. Signal Process Renew Energy 2022;6(1):39-64 [FREE Full text]
- 146. van der Miesen MM, Lindquist MA, Wager TD. Neuroimaging-based biomarkers for pain: state of the field and current directions. Pain Rep 2019;4(4):e751 [FREE Full text] [doi: 10.1097/PR9.000000000000751] [Medline: 31579847]
- 147. Werner P, Al-Hamadi A, Niese R, Gruss S, Traue HC. Automatic pain recognition from video and biomedical signals. In: Proceedings of the 22nd International Conference on Pattern Recognition. 2014 Presented at: ICPR '14; August 24-28, 2014; Stockholm, Sweden p. 4582-4587 URL: <u>https://ieeexplore.ieee.org/document/6977497</u> [doi: <u>10.1109/icpr.2014.784</u>]

- 148. Gruss S, Treister R, Werner P, Traue HC, Crawcour S, Andrade A, et al. Pain intensity recognition rates via biopotential feature patterns with support vector machines. PLoS One 2015 Oct 16;10(10):e0140330 [FREE Full text] [doi: 10.1371/journal.pone.0140330] [Medline: 26474183]
- Campbell E, Phinyomark A, Scheme E. Feature extraction and selection for pain recognition using peripheral physiological signals. Front Neurosci 2019 May 7;13:437 [FREE Full text] [doi: 10.3389/fnins.2019.00437] [Medline: 31133782]
- 150. Kachele M, Thiam P, Amirian M, Schwenker F, Palm G. Methods for person-centered continuous pain intensity assessment from bio-physiological channels. IEEE J Sel Top Signal Process 2016 Aug;10(5):854-864. [doi: 10.1109/jstsp.2016.2535962]
- 151. Fang R, Zhang R, Hosseini SM, Faghih M, Rafatirad S, Rafatirad S, et al. Pain level modeling of intensive care unit patients with machine learning methods: an effective congeneric clustering-based approach. In: Proceedings of the 4th International Conference on Intelligent Medicine and Image Processing. 2022 Presented at: IMIP '22; March 18-21, 2022; Tianjin, China p. 89-95 URL: <u>https://dl.acm.org/doi/pdf/10.1145/3524086.3524100</u> [doi: <u>10.1145/3524086.3524100</u>]
- 152. Nakano K, Ota Y, Ukai H, Nakamura K, Fujita H. Frequency detection method based on recursive DFT algorithm. In: Proceedings of the 14th International Conference on Power Systems Computation. 2002 Presented at: PSCC '02; June 24-28, 2002; Seville, Spain p. 1-7 URL: <u>https://www.researchgate.net/publication/</u> 255601650 Frequency detection method based on recursive DFT algorithm
- 153. Chen W, Zhuang J, Yu W, Wang Z. Measuring complexity using FuzzyEn, ApEn, and SampEn. Med Eng Phys 2009 Jan;31(1):61-68. [doi: 10.1016/j.medengphy.2008.04.005] [Medline: 18538625]
- 154. Wolpert DH, Macready WG. No free lunch theorems for optimization. IEEE Trans Evol Computat 1997;1(1):67-82. [doi: 10.1109/4235.585893]
- 155. Bellmann P, Thiam P, Kestler HA, Schwenker F. Machine learning-based pain intensity estimation: where pattern recognition meets chaos theory—an example based on the Biovid heat pain database. IEEE Access 2022;10:102770-102777. [doi: 10.1109/access.2022.3208905]
- 156. Gouverneur P, Li F, Adamczyk WM, Szikszay TM, Luedtke K, Grzegorzek M. Comparison of feature extraction methods for physiological signals for heat-based pain recognition. Sensors (Basel) 2021 Jul 15;21(14):4838 [FREE Full text] [doi: 10.3390/s21144838] [Medline: 34300578]
- 157. Othman E, Werner P, Saxen F, Fiedler MA, Al-Hamadi A. An automatic system for continuous pain intensity monitoring based on analyzing data from Uni-, Bi-, and multi-modality. Sensors (Basel) 2022 Jul 01;22(13):4992 [FREE Full text] [doi: 10.3390/s22134992] [Medline: 35808487]
- 158. Pouromran F, Lin Y, Kamarthi S. Personalized deep Bi-LSTM RNN based model for pain intensity classification using EDA signal. Sensors 2022 Oct 22;22(21):8087. [doi: 10.3390/s22218087]
- 159. Thiam P, Hihn H, Braun DA, Kestler HA, Schwenker F. Multi-modal pain intensity assessment based on physiological signals: a deep learning perspective. Front Physiol 2021 Sep 1;12:720464 [FREE Full text] [doi: 10.3389/fphys.2021.720464] [Medline: 34539444]
- 160. Cortes C, Vapnik V. Support-vector networks. Mach Learn 1995;20:297 [FREE Full text]
- 161. Verikas A, Gelzinis A, Bacauskiene M. Mining data with random forests: a survey and results of new tests. Pattern Recognit 2011 Feb;44(2):330-349. [doi: 10.1016/j.patcog.2010.08.011]
- 162. Breiman L, Friedman JH, Olshen RA, Stone CJ. Classification and Regression Trees. New York, NY: Routledge; 2017.
- 163. Breiman L. Random forests. Mach Learn 2001;45(1):5-32 [FREE Full text]
- 164. Gilpin LH, Bau D, Yuan BZ, Bajwa A, Specter M, Kagal L. Explaining explanations: an overview of interpretability of machine learning. In: Proceedings of the 5th International Conference on Data Science and Advanced Analytics. 2018 Presented at: DSAA '18; October 1-3, 2018; Turin, Italy p. 80-89 URL: <u>https://ieeexplore.ieee.org/document/8631448</u> [doi: 10.1109/dsaa.2018.00018]
- 165. Pal M, Mather PM. An assessment of the effectiveness of decision tree methods for land cover classification. Remote Sens Environ 2003 Aug;86(4):554-565. [doi: <u>10.1016/s0034-4257(03)00132-9</u>]
- 166. Fang R, Zhang R, Hosseini E, Parenteau AM, Hang S, Rafatirad S. Prevent over-fitting and redundancy in physiological signal analyses for stress detection. In: Proceedings of the 2022 IEEE International Conference on Bioinformatics and Biomedicine. 2022 Presented at: BIBM '22; December 6-8, 2022; Las Vegas, NV p. 2585-2588 URL: <u>https://ieeexplore.ieee.org/document/9995121</u> [doi: 10.1109/bibm55620.2022.9995121]
- 167. Naeini EK, Shahhosseini S, Subramanian A, Yin T, Rahmani AM, Dutt N. An edge-assisted and smart system for real-time pain monitoring. In: Proceedings of the 2019 IEEE/ACM International Conference on Connected Health: Applications, Systems and Engineering Technologies. 2019 Presented at: CHASE '19; September 25-27, 2019; Arlington, VA p. 47-52 URL: <u>https://ieeexplore.ieee.org/document/8908653</u> [doi: 10.1109/chase48038.2019.00023]
- 168. Werner P, Al-Hamadi A, Gruss S, Walter S. Twofold-multimodal pain recognition with the X-ITE pain database. In: Proceedings of the 8th International Conference on Affective Computing and Intelligent Interaction Workshops and Demos. 2019 Presented at: ACIIW '19; September 3-6, 2019; Cambridge, UK p. 290-296 URL: <u>https://ieeexplore.ieee.org/document/ 8925061</u> [doi: 10.1109/aciiw.2019.8925061]
- 169. Fang C, Miao N, Srivastav S, Liu J, Zhang R, Fang R, Asmita, et al. Large language models for code analysis: do LLMS really do their job? arXiv Preprint posted online October 18, 2023 [FREE Full text]

- 170. Lopez-Martinez D, Picard R. Multi-task neural networks for personalized pain recognition from physiological signals. In: Proceedings of the 7th International Conference on Affective Computing and Intelligent Interaction Workshops and Demos. 2017 Presented at: ACIIW '17; October 23-26, 2017; San Antonio, TX p. 181-184 URL: <u>https://www.computer.org/csdl/proceedings-article/aciiw/2017/08272611/12OmNAZfxKZ</u> [doi: 10.1109/aciiw.2017.8272611]
- 171. Salekin MS, Zamzmi G, Goldgof D, Kasturi R, Ho T, Sun Y. Multimodal spatio-temporal deep learning approach for neonatal postoperative pain assessment. Comput Biol Med 2021 Feb;129:104150 [FREE Full text] [doi: 10.1016/j.compbiomed.2020.104150] [Medline: <u>33348218</u>]
- 172. Pinzon-Arenas JO, Kong Y, Chon KH, Posada-Quintero HF. Design and evaluation of deep learning models for continuous acute pain detection based on phasic electrodermal activity. IEEE J Biomed Health Inform 2023 Sep;27(9):4250-4260. [doi: 10.1109/JBHI.2023.3291955] [Medline: <u>37399159</u>]
- 173. Bertrand AL, Garcia JB, Viera EB, Santos AM, Bertrand RH. Pupillometry: the influence of gender and anxiety on the pain response. Pain Physician 2013;16(3):E257-E266 [FREE Full text] [doi: 10.36076/ppj.2013/16/e257] [Medline: 23703424]
- 174. Chapman CR, Oka S, Bradshaw DH, Jacobson RC, Donaldson GW. Phasic pupil dilation response to noxious stimulation in normal volunteers: relationship to brain evoked potentials and pain report. Psychophysiology 1999 Jan 20;36(1):44-52. [doi: 10.1017/s0048577299970373] [Medline: 10098379]
- 175. Eisenach JC, Curry R, Aschenbrenner CA, Coghill RC, Houle TT. Pupil responses and pain ratings to heat stimuli: Reliability and effects of expectations and a conditioning pain stimulus. J Neurosci Methods 2017 Mar 01;279:52-59 [FREE Full text] [doi: 10.1016/j.jneumeth.2017.01.005] [Medline: 28089758]
- 176. Wang L, Guo Y, Dalip B, Xiao Y, Urman RD, Lin Y. An experimental study of objective pain measurement using pupillary response based on genetic algorithm and artificial neural network. Appl Intell 2021 May 17;52(2):1145-1156. [doi: 10.1007/s10489-021-02458-4]
- 177. Kächele M, Werner P, Al-Hamadi A, Palm G, Walter S, Schwenker F. Bio-visual fusion for person-independent recognition of pain intensity. In: Proceedings of the 12th International Workshop on Multiple Classifier Systems. 2015 Presented at: MCS '15; June 29-July 1, 2015; Günzburg, Germany p. 220-230 URL: <u>https://link.springer.com/chapter/10.1007/</u> <u>978-3-319-20248-8_19</u> [doi: <u>10.1007/978-3-319-20248-8_19</u>]
- 178. Kächele M, Thiam P, Amirian M, Werner P, Walter S, Schwenker F, et al. Multimodal data fusion for person-independent, continuous estimation of pain intensity. In: Proceedings of the 16th International Conference on Engineering Applications of Neural Networks. 2015 Presented at: EANN '15; September 25-28, 2015; Rhodes, Greece p. 275-285 URL: <u>https://link.springer.com/chapter/10.1007/978-3-319-23983-5_26</u> [doi: 10.1007/978-3-319-23983-5_26]
- 179. Thiam P, Kessler V, Schwenker F. Hierarchical combination of video features for personalised pain level recognition. In: Proceedings of the 2017 Conference on European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning. 2017 Presented at: ESANN '17; April 26-28, 2017; Bruges, Belgium p. 465-470 URL: <u>https://www. esann.org/sites/default/files/proceedings/legacy/es2017-104.pdf</u>
- 180. Kessler V, Thiam P, Amirian M, Schwenker F. Multimodal fusion including camera photoplethysmography for pain recognition. In: Proceedings of the 2017 International Conference on Companion Technology. 2017 Presented at: ICCT '17; September 11-13, 2017; Ulm, Germany p. 1-4 URL: <u>https://ieeexplore.ieee.org/document/8287083</u> [doi: <u>10.1109/companion.2017.8287083</u>]
- 181. Thiam P, Schwenker F. Multi-modal data fusion for pain intensity assessment and classification. In: Proceedings of the 7th International Conference on Image Processing Theory, Tools and Applications. 2017 Presented at: IPTA '17; November 28-December 1, 2017; Montreal, QC p. 1-6 URL: <u>https://ieeexplore.ieee.org/document/8310115</u> [doi: <u>10.1109/ipta.2017.8310115</u>]
- 182. Bellmann P, Thiam P, Schwenker F. Dominant channel fusion architectures-an intelligent late fusion approach. In: Proceedings of the 2020 International Joint Conference on Neural Networks. 2020 Presented at: IJCNN '20; July 19-24, 2020; Glasgow, Scotland p. 1-8 URL: <u>https://ieeexplore.ieee.org/document/9206814</u> [doi: <u>10.1109/ijcnn48605.2020.9206814</u>]
- 183. Bellmann P, Thiam P, Schwenker F. Using meta labels for the training of weighting models in a sample-specific late fusion classification architecture. In: Proceedings of the 25th International Conference on Pattern Recognition. 2021 Presented at: ICPR '21; January 10-15, 2021; Milan, Italy p. 2604-2611 URL: <u>https://ieeexplore.ieee.org/document/9412509</u> [doi: 10.1109/icpr48806.2021.9412509]
- 184. Oliveira F, Costa DG, Assis F, Silva I. Internet of intelligent things: a convergence of embedded systems, edge computing and machine learning. Internet Things 2024 Jul;26:101153. [doi: 10.1016/j.iot.2024.101153]
- 185. Xiong Y, Zeng Z, Chakraborty R, Tan M, Fung G, Li Y, et al. Nyströmformer: a Nyström-based algorithm for approximating self-attention. In: Proceedings of the 35th AAAI Conference on Artificial Intelligence. 2021 May 18 Presented at: AAAI '21; February 2-9, 2021; Vancouver, BC p. 14138-14148 URL: https://tinyurl.com/yc3epb39 [doi: 10.1609/aaai.v35i16.17664]
- 186. Nielsen CS, Staud R, Price DD. Individual differences in pain sensitivity: measurement, causation, and consequences. J Pain 2009 Mar;10(3):231-237 [FREE Full text] [doi: 10.1016/j.jpain.2008.09.010] [Medline: 19185545]
- 187. Jiang M, Rosio R, Salanterä S, Rahmani AM, Liljeberg P, da Silva DS, et al. Personalized and adaptive neural networks for pain detection from multi-modal physiological features. Expert Syst Appl 2024 Jan;235:121082. [doi: 10.1016/j.eswa.2023.121082]

```
https://ai.jmir.org/2025/1/e53026
```

- 188. Casti P, Mencattini A, Filippi J, D'Orazio M, Comes MC, Giuseppe DD. A personalized assessment platform for non-invasive monitoring of pain. In: Proceedings of the 2020 IEEE International Symposium on Medical Measurements and Applications. 2020 Presented at: MeMeA '20; June 1-4, 2020; Bari, Italy p. 1-5 URL: <u>https://ieeexplore.ieee.org/document/9137138</u> [doi: 10.1109/memea49120.2020.9137138]
- 189. Lopez Martinez D, Rudovic O, Picard R. Personalized automatic estimation of self-reported pain intensity from facial expressions. In: Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops. 2017 Presented at: CVPRW '17; July 21-26, 2017; Honolulu, HI p. 2318-2327 URL: <u>https://ieeexplore.ieee.org/document/8015020</u> [doi: 10.1109/cvprw.2017.286]
- 190. Zhang R, Fang R, Zhang Z, Hosseini E, Orooji M, Homayoun H. Short: real-time bladder monitoring by bio-impedance analysis to aid urinary incontinence. In: Proceedings of the 2023 IEEE/ACM Conference on Connected Health: Applications, Systems and Engineering Technologies. 2023 Presented at: CHASE '23; June 21-23, 2023; Orlando, FL p. 138-142 URL: https://ieeexplore.ieee.org/document/10183756 [doi: 10.1145/3580252.3586985]
- 191. Kong Y, Posada-Quintero HF, Chon KH. Real-time high-level acute pain detection using a smartphone and a wrist-worn electrodermal activity sensor. Sensors (Basel) 2021 Jun 08;21(12):3956 [FREE Full text] [doi: 10.3390/s21123956] [Medline: 34201268]
- 192. Fang R, Zhang R, Hosseini E, Parenteau AM, Hang S, Rafatirad S. Towards generalized ML model in automated physiological arousal computing: a transfer learning-based domain generalization approach. In: Proceedings of the 2022 IEEE International Conference on Bioinformatics and Biomedicine. 2022 Presented at: BIBM '22; December 6-8, 2022; Las Vegas, NV p. 2577-2584 URL: <u>https://ieeexplore.ieee.org/document/9995340</u> [doi: 10.1109/bibm55620.2022.9995340]
- 193. Kächele M, Amirian M, Thiam P, Werner P, Walter S, Palm G, et al. Adaptive confidence learning for the personalization of pain intensity estimation systems. Evol Syst 2016 Jul 16;8(1):71-83. [doi: 10.1007/s12530-016-9158-4]
- Chen J, Liu X, Tu P, Aragones A. Learning person-specific models for facial expression and action unit recognition. Pattern Recognit Lett 2013 Nov;34(15):1964-1970. [doi: <u>10.1016/j.patrec.2013.02.002</u>]

Abbreviations

AU: action unit AUC: area under the curve B-CNN: bilinear convolutional neural network CNN: convolutional neural network **DMSN:** Decomposed Multiscale Spatiotemporal Network **EDA:** electrodermal activity FACE-BE-SELF: Facial Expressions Fusing Betamix Selected Landmark Features FACS: Facial Action Coding System fMRI: functional magnetic resonance imaging fNIRS: functional near-infrared spectroscopy HF: high-frequency HOG: histogram of oriented gradients **HRV:** heart rate variability **ICU:** intensive care unit LBP: local binary pattern LF: low-frequency **LSTM:** long short-term memory ML: machine learning PCA: principal component analysis RF: random forest **RGB:** red, green, blue color model **RNN:** recurrent neural network **RVR:** relevance vector regression **sEMG:** surface electromyogram **SNS:** sympathetic nervous system SVM: support vector machine



Edited by JL Raisaro; submitted 22.09.23; peer-reviewed by A Naser, S Kisvarday, A Subramanian, P Lakshman, A Mazumder; comments to author 11.04.24; revised version received 06.06.24; accepted 23.07.24; published 24.02.25. <u>Please cite as:</u> Fang R, Hosseini E, Zhang R, Fang C, Rafatirad S, Homayoun H Survey on Pain Detection Using Machine Learning Models: Narrative Review JMIR AI 2025;4:e53026 URL: https://ai.jmir.org/2025/1/e53026 doi:10.2196/53026 PMID:

©Ruijie Fang, Elahe Hosseini, Ruoyu Zhang, Chongzhou Fang, Setareh Rafatirad, Houman Homayoun. Originally published in JMIR AI (https://ai.jmir.org), 24.02.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Studying the Potential Effects of Artificial Intelligence on Physician Autonomy: Scoping Review

John Grosser¹, MA, MSc; Juliane Düvel², MSc; Lena Hasemann¹, MSc; Emilia Schneider¹; Wolfgang Greiner¹, Prof Dr

¹Department of Health Economics and Health Care Management, School of Public Health, Bielefeld University, Bielefeld, Germany ²Centre for Electronic Public Health Research (CePHR), School of Public Health, Bielefeld University, Bielefeld, Germany

Corresponding Author: John Grosser, MA, MSc Department of Health Economics and Health Care Management School of Public Health Bielefeld University Universitätsstraße 25 Bielefeld, 33615 Germany Phone: 49 52110686319 Email: john.grosser@uni-bielefeld.de

Abstract

Background: Physician autonomy has been found to play a role in physician acceptance and adoption of artificial intelligence (AI) in medicine. However, there is still no consensus in the literature on how to define and assess physician autonomy. Furthermore, there is a lack of research focusing specifically on the potential effects of AI on physician autonomy.

Objective: This scoping review addresses the following research questions: (1) How do qualitative studies conceptualize and assess physician autonomy? (2) Which aspects of physician autonomy are addressed by these studies? (3) What are the potential benefits and harms of AI for physician autonomy identified by these studies?

Methods: We performed a scoping review of qualitative studies on AI and physician autonomy published before November 6, 2023, by searching MEDLINE and Web of Science. To answer research question 1, we determined whether the included studies explicitly include physician autonomy as a research focus and whether their interview, survey, and focus group questions explicitly name or implicitly include aspects of physician autonomy. To answer research question 2, we extracted the qualitative results of the studies, categorizing them into the 7 components of physician autonomy introduced by Schulz and Harrison. We then inductively formed subcomponents based on the results of the included studies in each component. To answer research question 3, we summarized the potentially harmful and beneficial effects of AI on physician autonomy in each of the inductively formed subcomponents.

Results: The search yielded 369 studies after duplicates were removed. Of these, 27 studies remained after titles and abstracts were screened. After full texts were screened, we included a total of 7 qualitative studies. Most studies did not explicitly name physician autonomy as a research focus or explicitly address physician autonomy in their interview, survey, and focus group questions. No studies addressed a complete set of components of physician autonomy; while 3 components were addressed by all included studies, 2 components were addressed by none. We identified a total of 11 subcomponents for the 5 components of physician autonomy that were addressed by at least 1 study. For most of these subcomponents, studies reported both potential harms and potential benefits of AI for physician autonomy.

Conclusions: Little research to date has explicitly addressed the potential effects of AI on physician autonomy and existing results on these potential effects are mixed. Further qualitative and quantitative research is needed that focuses explicitly on physician autonomy and addresses all relevant components of physician autonomy.

(JMIR AI 2025;4:e59295) doi:10.2196/59295

KEYWORDS

RenderX

autonomy, professional autonomy; physician autonomy; ethics; artificial intelligence; clinical decision support systems; CDSS; ethics of artificial intelligence; AI ethics; AI; scoping review; physician; acceptance; adoption

Introduction

The use of artificial intelligence (AI) systems in medicine has increased significantly in recent years. AI in medicine can take a number of forms and fulfill a number of tasks, ranging from risk prediction or diagnosis and screening to AI-powered clinical decision support systems (CDSS) [1]. AI systems have also been introduced across a range of medical specialties, including oncology, pulmonology, and radiology [2].

Physician autonomy has been found to play a role in physician acceptance and adoption of medical technologies [3], and in particular, AI [1]. Although physician autonomy has become an increasingly important concept in recent decades [4-7], there is still no consensus definition in the literature. However, physician autonomy is generally seen as including both clinical freedoms, as well as social and economic freedoms [6,7]. The former concerns physician autonomy in clinical practice, including their control over the diagnosis and treatment of patients and over evaluations of their care. The latter concerns the autonomy of physicians as professionals, including their choice of specialty and control over the nature and volume of their tasks [5]. A number of recent reviews have found that the feared loss of physician autonomy represents a barrier to the acceptance of AI [1,8-10]. However, although these reviews (partially) address physician autonomy as a barrier to acceptance, there is little research so far focusing primarily on the effects of AI on physician autonomy. Furthermore, such reviews rarely systematically address both clinical, social, and economic freedoms.

Our aim is to begin to fill this gap by performing a scoping review of qualitative studies on AI and physician autonomy. In particular, this review addresses the following research questions: (1) How do these studies conceptualize and assess physician autonomy? (2) Which aspects of physician autonomy are addressed by these studies? (3) What are the potential benefits and harms of AI for physician autonomy identified by

Textbox 1. Inclusion and exclusion criteria.

Inclusion criteria

- Empirical, qualitative, or mixed methods study
- Focus on artificial intelligence (AI) in clinical care
- Physician autonomy addressed in the study
- The study population includes physicians
- English or German language

Exclusion criteria

- Nonempirical or purely quantitative study
- No focus on AI
- Focus on AI in veterinary medicine or public health
- Physician autonomy not addressed in the study
- The study population does not include physicians
- Language other than English or German

these studies? To address research question 1, we investigate whether and how the studies include physician autonomy as a research focus in their interview, survey, and focus group questions. To answer research question 2, we identify the components of physician autonomy addressed by the studies based on the 7-component model proposed by Schulz and Harrison [5]. For each of these components, we then inductively form subcomponents based on the results of the included studies. To answer research question 3, we summarize the potential benefits and harms of AI for physician autonomy reported by the included studies in each subcomponent. These questions lend themselves to a scoping review approach, rather than a systematic review since we aim to answer broader conceptual and methodological questions, rather than perform a risk of bias assessment or meta-analysis [11].

Methods

Search Strategy

We performed a scoping review of qualitative studies on AI and physician autonomy and drafted the paper according to the PRISMA (Preferred Reporting Items for Systematic Reviews and Meta-Analyses) checklist (Multimedia Appendix 1) [11]. We searched MEDLINE and Web of Science using a search string based on the following combination of concepts: "Physician" AND "Artificial Intelligence" AND "Autonomy" AND "Qualitative Research." The complete search terms for both databases (including Medical Subject Headings terms and keywords) can be found in Multimedia Appendix 2. The cutoff date for the search was November 6, 2023.

Screening

After removing duplicates, the titles and abstracts of the remaining studies were screened by 2 authors (JD and LH) according to predefined inclusion and exclusion criteria (Textbox 1). This was followed by a screening of the remaining full texts. Disagreements and concerns regarding the results were resolved in consultation with a third researcher (JG).

Data Extraction and Synthesis

For each included study, we first extracted relevant study characteristics, including country, design, and study population, as well as the AI system under consideration. We also ascertained whether the included studies explicitly include physician autonomy as a research focus and reviewed supplemental material, where available, to determine whether their interview, survey, and focus group questions explicitly name physician autonomy or implicitly include aspects of physician autonomy. We then extracted the qualitative results of the studies, categorizing them into 7 components of physician autonomy introduced by Schulz and Harrison [5]. This categorization contains 3 social and economic freedoms (Textbox 2) and 4 clinical freedoms (Textbox 3).

Textbox 2. Social and economic components of physician autonomy [5].

Choice of specialty and practice location

• Potential limitations on autonomy include market restrictions, bureaucratic restrictions, and educational restrictions

Control over earnings

• Potential limitations on autonomy include workload controls, fee schedules, reimbursement rates, salaried status, and control over permitted earnings

Control over the nature and volume of medical tasks

• Potential limitations on autonomy include hierarchical management, contractual obligations, and the need to share scarce resources

Textbox 3. Clinical components of physician autonomy [5].

Acceptance of patients

• Potential limitations on autonomy include compelling physicians to accept or reject certain patients based on geography, medical specialty, or insurance status

Control over diagnosis and treatment

• Potential limitations on autonomy include individual and aggregate constraints on tests or prescription costs, preset budgets, enforcement of clinical protocols, and gatekeeping

Control over evaluation of care

• Potential limitations on autonomy include peer review, medical audit systems, and comparative information on care outcomes

Control over other professionals

• Potential limitations on autonomy include limitations on physicians' ability to directly manage other health professionals and include precise instructions in referrals for diagnosis or therapy

To paint a more detailed picture of the effect of AI on physician autonomy, we inductively formed subcomponents from the results in each component. To avoid overgeneralizing based on individual participants and studies, we only considered subcomponents that were addressed by at least 2 included studies. Finally, we summarized the potentially harmful and beneficial effects of AI on physician autonomy in each of the inductively formed subcomponents.

Results

Selection of Sources of Evidence

The search yielded 369 studies after duplicates were removed (Figure 1). Of these, 27 studies remained after titles and abstracts were screened. After full texts were screened, we included a total of 7 qualitative studies [12-18].







Study Characteristics

All 7 included studies had a cross-sectional design; most studies (n=5) used (qualitative) semistructured interviews, which 1 study [13] combined with a focus group. The remaining studies used co-design workshops [16] and a mixed methods survey consisting of both quantitative and qualitative items [15] (although we focus only on the qualitative results). More than half of the studies (n=4) were conducted in Europe; 2 studies were conducted in Asia and one in Australia (Table 1). Radiologists [13,17] and general practitioners (GPs) or primary care physicians [16,18] were the focus of 2 studies each, while

the remaining studies recruited participants across multiple specialties. Some studies also included further groups, such as patients or family members [12,18], medical students [15], and radiographers [13], in addition to physicians. The most common form of (medical) AI investigated was CDSS (n=3). Digital disease surveillance systems and documentation assistants were investigated by 1 study each. The remaining 2 studies investigated various forms of AI in medicine. However, only 1 study [17] explicitly recruited participants who had experience with medical AI systems; the remaining studies merely provided participants with vignettes or videos of possible AI systems.



Table 1. Study characteristics of the included studies.

Study	Country	Study period	Participants	AI ^a system
Amann et al (2023) [12]	Germany and Switzerland	2019-2020	14 health care professionals, 14 stroke survivors, and 6 family members of stroke survivors	CDSS ^b
Chen et al (2021) [13]	United Kingdom	2018-2020	12 physicians (radiologists) and 6 radiographers	Various
Huang et al (2023) [14]	Singapore and In- dia	2022	45 physicians	CDSS
Jussupow et al (2022) [15]	Germany	2017-2019	164 medical students and 42 medical professionals	CDSS
Kocaballi et al (2020) [16]	Australia	NR	16 physicians (GPs ^c)	DA ^d
Lombi and Rossero (2023) [17]	Italy	2021	12 physicians (radiologists)	Various
Wong et al (2023) [18]	China	2021	16 physicians (PCPs ^e) and 24 patients	DDS ^f

^aAI: artificial intelligence.

^bCDSS: clinical decision support systems.

^cGP: general practitioner.

^dDA: documentation assistant.

^ePCP: primary care physician.

^fDDS: digital disease surveillance.

Conceptualizing and Assessing Physician Autonomy

The studies differed significantly in how they conceptualized physician autonomy and to what extent physician autonomy was the focus of their research. In particular, only 1 study [17]

explicitly named (the effect of AI on) physician autonomy as a research focus (Table 2). The remaining studies focused on expectations and acceptance of or views and attitudes toward AI.

Table 2. The role of physician autonomy in the included studies.

	[12] ^a	[13]	[14]	[15]	[16]	[17]	[18]
Physician autonomy is an explicit focus of the study						1	
Questions explicitly include physician autonomy			1			\checkmark	
Questions implicitly include physician autonomy	✓		1		1	1	

^aThe interview questions reference "autonomy," but not explicitly physician autonomy.

Only 2 of 7 included studies [14,17] explicitly included physician autonomy in their interview, survey, or focus group questions, and of these, only one study [17] uses a concrete theoretical framework for physician autonomy. Nevertheless, more than half of the studies (implicitly) included at least some aspects of physician autonomy in their interview questions, even if they did not explicitly relate them to physician autonomy. The remaining studies did not include physician autonomy in their interview questions but did identify aspects of physician autonomy in their participants' responses. Therefore, although most studies did not explicitly name physician autonomy as a research focus or in their interview questions, the qualitative results of all studies include a number of themes related to physician autonomy. We categorized these results into the 7 components of physician autonomy proposed by Schulz and Harrison [5] and formed 2-3 subcomponents for each component, described in the following sections.

Social and Economic Subcomponents of Physician Autonomy

For the choice of specialty and practice location, we identified two subcomponents: (1) AI replacing physicians and (2) AI replacing specialties. Three studies [12,15,16] reported that

```
https://ai.jmir.org/2025/1/e59295
```

RenderX

physicians feared becoming redundant or being replaced by AI. This represents an (indirect) threat to physician autonomy in choosing their specialty and practice location, as this choice will not be available to physicians who have been replaced by AI. In contrast, however, participants in 2 studies [12,16] argued that AI cannot or will not replace physicians, either because fully autonomous medical AI was seen as unrealistic (at least in the near future) or because AI was seen as unable to perform core tasks of (human) physicians, such as empathy and human warmth or communication.

A number of studies also addressed the risk of certain physician specialties, such as GPs [16] and radiologists [13,17], being replaced by or becoming mere assistants of AI—a direct threat to physician autonomy in choosing specialty and practice location. However, 2 studies [13,17] also found that radiologists were seen as less vulnerable to replacement by AI since their roles encompass a wide range of challenging activities (including complex diagnoses and patient relationships), which AI cannot replace as easily as routine reporting activities.

For control over the nature and volume of medical tasks, we identified three subcomponents: (1) the effect of AI on workflow and efficiency, (2) the ability of physicians to personalize and

customize AI tools, and (3) involving physicians in AI design and creation. Participants in all 7 studies [12-18] believed that AI could increase efficiency by redefining workflows, taking over mundane and repetitive administrative tasks, and allowing faster decision-making. This would help address workforce shortages and free up more time for physicians to pursue other, more preferred tasks, such as research or treating complex cases. In this way, AI could enhance physician autonomy over the nature and volume of their tasks. However, participants in 3 of these studies [14,16,17] also expressed hesitation about the time-saving potential of AI, noting that additional time and effort may be required to input required data, fix errors, and train both physicians and AI systems.

Two studies [14,16] addressed further subcomponents relevant to physician control over the nature and volume of medical tasks. At the micro level, these studies addressed the ability of physicians to personalize and customize AI systems. In particular, AI systems may also enhance physician autonomy over the nature and volume of their work through personalized and adaptive features [16], although physicians in 1 study did not find AI customizability necessary [14]. At the macro level, both studies [14,16] addressed the importance of involving physicians in the design and creation of AI systems. While not every physician can be involved in the cocreation of AI, this would nevertheless increase the control of physicians as a group over the AI systems they will be working with. Table 3 shows the distribution of the components or subcomponents for social and economic freedoms among the included studies. Note that none of the included studies addressed control over earnings.

Table 3. Social and economic components or subcomponents of physician au	utonomy
--	---------

Component or subcomponent	Number of studies	Studies			
Choice of specialty and practice location					
AI ^a replacing physicians	3	[12,15,16]			
AI replacing specialties	3	[13,16,17]			
Total	5	[12,13,15-17]			
Control over earnings					
Total	0	b			
Control over the nature and volume of medical tasks					
AI and workflow or efficiency	7	[12-18]			
AI customization or personalization	2	[14,16]			
Involving physicians in AI design or creation	2	[14,16]			
Total	7	[12-18]			

^aAI: artificial intelligence.

^bNot applicable.

Clinical Subcomponents of Physician Autonomy

For control over diagnosis and treatment, we identified two subcomponents: (1) the (direct) effect of AI on clinical decision-making and (2) the effect of AI on physicians' expertise and skills. Five studies [12-14,16,18] reported concerns that AI may negatively affect physicians' clinical decision-making autonomy; participants in most of these studies [12-14] agreed that physicians should remain the final authority in clinical decision-making. Participants in other studies were less concerned about this risk, arguing that AI systems will not negatively affect physician autonomy when their adoption is voluntary [14] or when they are used as only one of many criteria informing physicians' clinical decisions [17].

In contrast, 4 studies [12,14-16] reported that AI systems may enhance physician autonomy in clinical decision-making, particularly for less experienced physicians, by affirming their decisions and increasing decision certainty, providing inspiration and offering new possibilities of care, or helping clinicians adhere to guidelines (note that while Amann et al [12] describe better adherence to guidelines as a positive effect of AI, a close reading of Schulz and Harrison [5] suggests that strict adherence

```
https://ai.jmir.org/2025/1/e59295
```

to guidelines may, in fact, decrease physician control over diagnosis and treatment).

All but 1 study [12,14-18] addressed the risk of automation bias, or the overreliance of physicians on AI systems, particularly when the use of such systems is mandated [14]. In addition to diagnostic errors [17], this overreliance may lead to deskilling and loss of expertise, especially in younger generations of physicians [12,14], indirectly reducing physicians' control over diagnosis and treatment by making some courses of action unavailable. Participants in 2 studies [13,17], however, were less concerned about this risk. For example, radiologists in 1 study [13] argued that their wide array of high-level tasks made them less vulnerable to deskilling by AI.

Conversely, 4 studies [12,13,15,16] found that AI systems may enhance the expertise and skills of physicians, thereby increasing rather than decreasing their control over diagnosis and treatment. For example, AI may assist physicians who are struggling to be empathetic by suggesting empathetic statements [16] or providing relevant and up-to-date information, especially for novice physicians [15].

Concerning control over the evaluation of care, we identified two subcomponents: (1) the effect of AI on the risk of medicolegal consequences for physicians and (2) the effect of AI on evaluations of care by patients. All but 1 study [12-17] addressed the risk of medicolegal consequences resulting from the use of AI systems. On the one hand, physicians feared the liability issues that may arise from disagreeing with AI decisions or recommendations [15,16], particularly in light of potential data biases in AI systems. On the other hand, they feared that AI systems may be used as auditing tools [16], retrospectively assessing physician's consultation and treatment records for potential errors in diagnosis or treatment. While many study participants agreed that the responsibility-and liability-for medical decisions involving AI rests with physicians as the final decision makers [12,14,17], a number of participants suggested that other actors, such as developers [12], host units [13], or hospitals [14], could share this responsibility (in full or in part).

Five studies [12,14-16,18] addressed the effects of AI on patient evaluations of care. On the one hand, participants in most of these studies feared that patients would negatively react to the use of AI because dependence on AI may undermine patients' faith in the competence of physicians and their recommendations [15,16], because intransparency about AI's use of patient data may threaten patient trust in physicians [18] or because patients may simply prefer human physicians [14]. On the other hand, some studies suggested that patients may approve of the use of AI as an evidence-based approach that can lead to improved care outcomes [14,15], and while Amann et al [12] found that patients should have a say when it comes to the use of AI, Huang et al [14] found that many physicians felt it unnecessary to discuss AI use with patients. Finally, we identified two subcomponents for control over other professionals: (1) indirect control and (2) direct control, which were addressed by two studies each. Indirect control refers to the status and prestige of physicians (individually and as a profession) in relation to other professionals, including other physicians. While Jussupow et al [15] found that AI systems were seen as leading to a loss in status and prestige for physicians in general, Lombi and Rossero [17] suggested that the advent of AI may present an opportunity for radiologists to reconfigure their professional identity and actually increase their status and prestige by becoming proficient in these technologies.

Direct control refers to the ability of physicians to directly influence or exercise authority over other professionals, including other physicians. While 2 studies [14,17] addressed this component, they conceptualized the effect of AI on professional control in different ways and no overarching themes emerged between them. On the one hand, Huang et al [14] found that senior physicians would encourage junior physicians to use AI and that physicians would, in fact, be influenced by colleagues to adopt AI. On the other hand, Lombi and Rossero [17] found that AI may transform and expand radiologists' interprofessional collaboration (including with nonclinical professionals). AI was seen as threatening professional boundaries and risking a loss of radiologist authority to other clinical professionals but was not seen as challenging radiologists' professional boundaries or authority concerning nonclinical professionals [17]. Table 4 shows the distribution of the components or subcomponents for clinical freedoms among the included studies. Note that none of the included studies addressed the acceptance of patients.

Table 4. Clinical components or subcomponents of physician autonor	ny.
--	-----

Component or subcomponent	Number of studies	Studies
Acceptance of patients		
Total	0	a
Control over diagnosis and treatment		
AI ^b and clinical decision-making	7	[12-18]
AI and physician expertise or skills	7	[12-18]
Total	7	[12-18]
Control over the evaluation of care		
AI and medicolegal consequences	6	[12-17]
AI and patient evaluations of care	5	[12,14-16,18]
Total	7	[12-18]
Control over other professionals		
AI and indirect control over other professionals	2	[15,17]
AI and direct control over other professionals	2	[14,17]
Total	3	[14,15,17]

^aNot applicable.

^bAI: artificial intelligence.

Potential Benefits and Harms of AI for Physician Autonomy

The main results of the included studies in each subcomponent are summarized in Textboxes 4 (for social and economic freedoms) and 5 (for clinical freedoms). For 6 of 11 subcomponents, we found mixed results concerning the potential benefits and harms of AI for physician autonomy. In particular, studies disagreed on whether AI will increase or decrease workflow efficiency, enhance or impede clinical decision-making, improve or worsen physician skills and expertise, lead to patient approval or disapproval, and increase or decrease physician status or prestige. Studies were also split on how AI will affect physicians' direct control over other professionals.

Textbox 4. Potential benefits and harms of artificial intelligence (AI) for social and economic freedoms, indicated by (+) and (-), respectively. Circles indicate relevant findings that are neither harms nor benefits.

Choice of specialty and practice location

AI replacing physicians (n=3)

- (+) AI (currently) lacks the capabilities, such as empathy, necessary to replace physicians
- (-) AI may replace physicians in the future

AI replacing specialties (n=3)

- (+) Radiologists are less vulnerable to AI replacement due to their wide range of challenging activities
- (-) AI may replace radiologists in the future
- (-) AI may replace general practitioners in the future

Control over the nature and volume of medical tasks

AI and workflow or efficiency (n=7)

- (+) AI can increase efficiency by handling mundane activities, freeing up time for other tasks
- (-) AI may decrease efficiency due to the time and effort required for data input, error correction and training

AI customization or personalization (n=2)

• (+) AI may support physicians through personalized and adaptive features

Involving physicians in AI design or creation (n=2)

• (o) Physicians should be involved in AI design or creation

For 2 subcomponents (AI replacing physicians and AI replacing specialties), we found mixed to negative results. On the one hand, the studies that addressed these 2 components found that physicians and some specialties (radiologists and GPs or primary care physicians) may be at risk of replacement by AI. On the

other hand, the studies gave a number of reasons why physicians and some specialties may be less vulnerable to such replacement, at least in the near future. However, while these results are not fully negative, we did not find any results indicating that AI may improve physician autonomy in these subcomponents.



Textbox 5. Potential benefits and harms of artificial intelligence (AI) for clinical freedoms, indicated by (+) and (-), respectively. Circles indicate relevant findings that are neither harms nor benefits.

Control over diagnosis and treatment

AI and clinical decision-making (n=7)

- (+) AI may enhance clinical autonomy by increasing decision certainty and providing inspiration
- (-) AI may harm clinical decision-making autonomy
- (o) Physicians should remain the final authority in clinical decision-making

AI and physician expertise or skills (n=7)

- (+) AI may enhance physicians' expertise
- (-) AI may lead to loss of expertise through overreliance and automation bias

Control over evaluation of care

AI and medicolegal consequences (n=6)

- (-) AI decisions and recommendations may lead to liability issues for physicians
- (-) AI systems may be used as post hoc auditing tools
- (o) Developers, hospitals, or other actors should (partially) share responsibility for medical decisions involving AI

AI and patient evaluations of care (n=5)

- (+) Patients may approve of AI use (eg, due to improved outcomes)
- (-) AI may lead to patient disapproval or mistrust
- (-) AI may undermine patients' faith in physicians' care

Control over other professionals

AI and indirect control over other professionals (n=2)

- (+) AI may offer radiologists an opportunity to increase their status and prestige
- (-) AI systems may lead to a loss in status and prestige for physicians in general

AI and direct control over other professionals (n=2)

- (+) AI may expand radiologists' interprofessional collaboration with nonclinical professionals
- (-) AI may threaten radiologists' authority over other clinical professionals
- (-) Physicians may be influenced by peers and superiors to adopt AI

In contrast, we found general agreement between the included studies for the remaining 3 subcomponents. For AI customization or personalization, this consensus was positive: both studies addressing this component found that customizable AI systems would support physician autonomy. Furthermore, there was agreement between studies that AI represented potential harms (but not benefits) to physician autonomy in the AI and medicolegal consequences component. Finally, both studies that addressed involving physicians in AI design or creation found that such involvement should take place (although this more accurately represents a recommendation or demand rather than a potential benefit or harm).

Discussion

Principal Results

These results show that research on the potential effects of AI on physician autonomy is still in its nascency. In particular, there is no consensus definition or operationalization of

```
https://ai.jmir.org/2025/1/e59295
```

physician autonomy in qualitative research. Most studies did not name physician autonomy as a focus of their research or explicitly include physician autonomy in their interview, survey, or focus group questions. In fact, only 1 study [17] specified a clear theoretical framework for physician autonomy. These results align with existing research on the professional autonomy of nurses, which has been found to face challenges due to inconsistent definitions and inappropriate measures of nurse autonomy [19] and the confounding of the clinical and nonclinical aspects of nurse autonomy [20].

No studies addressed a complete set of components of physician autonomy (as defined by Schulz and Harrison [5]). Furthermore, coverage between components varies significantly: while all 7 studies addressed control over the nature and volume of medical tasks, control over diagnosis and treatment, and control over the evaluation of care, none of the included studies addressed control over earnings and acceptance of patients.

We identified a total of 11 subcomponents for the 5 components of physician autonomy that were addressed by at least 1 study. For most of these subcomponents, studies reported mixed results concerning the potential harms and benefits of AI for physician autonomy. A notable exception addressed by most studies was AI and medicolegal consequences, with studies reporting only potential harms for this subcomponent. AI customization or personalization was the only subcomponent in which only potential benefits were reported, although this subcomponent was only addressed by 2 studies. Overall, there is a need for further research that focuses specifically on physician autonomy and includes a full conception of its components and subcomponents.

Some of the results within subcomponents align with recent reviews of the academic literature, which have found positive effects of AI on clinical and administrative workflow or efficiency or patient-physician trust [21,22]. A recent review of the "grey literature" also found that clinical and administrative AI applications impact physician job autonomy, skills, and professional relationships [23]. However, not all of these results are reported by the reviews as components of physician autonomy.

Limitations

However, the methodological limitations of our scoping review should be considered when interpreting our results. In particular, we identified only 7 studies that fit the inclusion criteria. Furthermore, although 4 of 7 studies [12,14,17,18] were published in 2023, only 1 study [14] specified a data collection period later than 2021 and 3 studies completed their data collection before the end of 2020. Considering the rapid evolution of AI in medicine, such as the recent introduction of large language models such as ChatGPT [24,25], there is a clear need for additional, up-to-date research on physician autonomy and new AI systems.

Furthermore, we included only qualitative studies in this review. In our view, expanding our scope to include a full systematic review of quantitative studies on AI and physician autonomy would have been premature, as the field is comparatively new and because we were focused particularly on how physician autonomy is defined and conceptualized by researchers and participants. However, the subcategories we have identified provide a useful roadmap for future systematic reviews of quantitative studies on physician autonomy and AI, and such reviews should be conducted.

Our review may also have missed further studies that were not included in the databanks we searched or that did not explicitly mention (physician) autonomy. However, these studies may still be relevant: while we assigned study results to components of physician autonomy in order to form inductive subcomponents, most of the included studies do not conceptualize physician autonomy as covering each of these components. For example, subcomponents such as AI and workflow or efficiency, AI and physician expertise or skills, or AI and patient evaluations of care were addressed by a number of studies, but usually not explicitly related to physician autonomy. This indicates that there may be further studies that address relevant components without explicitly mentioning

```
https://ai.jmir.org/2025/1/e59295
```

autonomy. This should also be considered when conducting future systematic reviews of quantitative studies on physician autonomy and AI. In particular, search terms related to specific subcomponents (but not physician autonomy) may lead to the inclusion of additional relevant studies.

Future research should also explicitly include the 2 components that were not addressed by any of the studies in our review: control over earnings and acceptance of patients. In particular, one should not conclude from our review that AI will have no effect on physician autonomy for these components. Such a conclusion seems implausible since examples of possible effects are easily constructed. For example, if AI systems were to take on the role of gatekeepers and play some part in deciding which patients can be seen by which physicians, this would represent harm to physician autonomy. Instead, the absence of these components from our review should be taken to indicate that respondents (or researchers) did not conceive of control over earnings and acceptance of patients as (relevant) aspects of physician autonomy.

Studies also differed in their definition of AI, which complicates evidence comparison and synthesis. While some studies considered AI-based CDSS, others considered different AI systems or AI innovations more broadly, and while 1 study [17] recruited participants who had actual working experience with AI systems, most merely presented participants with vignettes describing possible AI systems. This means that most studies report only the potential harms and benefits of AI (as feared or hoped for by participants), not actual harms and benefits. As a systematic comparison of the effects of different types of AI systems on physician autonomy was not possible with only 7 included studies, our scoping review is further limited to a broader discussion of the potential effects of AI in general. However, further research should analyze these differences in effect, based (where possible) on evaluations of actual AI systems, rather than vignettes.

Initial evidence also suggests that participants in different regions or cultures perceive different potential harms and benefits of AI for physician autonomy. For example, Huang et al [14] found that views on (the effects of AI on) some aspects of physician autonomy differed between physicians in Singapore and India, while Wong et al [18] discuss the fragility of doctor-patient trust specifically in China. While we were unable to analyze these differences due to the limited number of studies, future research should more thoroughly investigate such cultural and geographic differences in attitudes toward both AI and physician autonomy.

Overall, our results are based on a limited number of studies and should be seen as opening, rather than closing, lines of inquiry into the effects of AI on physician autonomy. Fully understanding these effects will require an ambitious research program. First, there is a need for further qualitative studies focusing explicitly on physician autonomy. Second, a definitive understanding of AI and physician autonomy will require quantitative studies using validated and reliable instruments designed for this purpose. Finally, the current literature focuses almost exclusively on self-reported physician autonomy. However, it may also be possible to measure the effect of AI

XSL•FO RenderX

on physician autonomy using objective quantitative indicators, such as the number of alerts and reviews triggered by AI systems or test results from experimental studies of physician expertise. Future research should consider if and when the use of such indicators in addition to self-reported assessments of physician autonomy is appropriate.

Conclusions

Little research to date has addressed the potential effects of AI on physician autonomy. Existing results on AI and physician autonomy are mostly secondary findings or merely part of larger analyses into physicians' attitudes toward and acceptance of AI. Most studies addressed physician autonomy only indirectly in their research focus and interview, survey, or focus group questions.

While 3 of the components of physician autonomy proposed by Schulz and Harrison [5] were addressed by all included studies, 2 components were not addressed by any studies. In eleven (inductively formed) subcomponents, the included studies reported a number of potential effects of AI on physician autonomy. However, results were mixed, with studies reporting both potential harms and benefits of AI for physician autonomy in most subcomponents.

In conclusion, further qualitative and quantitative research is needed that focuses explicitly on physician autonomy and addresses all relevant components of physician autonomy. Where possible, research on the effects of AI on physician autonomy should be based on real experience with AI systems, rather than vignettes, and consider the differences between different AI systems and between physicians in different cultural and geographic settings.

Acknowledgments

All authors contributed to the study's conception and design. JD and LH devised the search strategy and performed the screening. JG was consulted to resolve disagreements. JG, JD, and ES performed the data extraction and synthesis. JG and LH drafted the manuscript, which was edited, discussed, and approved by all authors. No funding was received to assist with the preparation of this manuscript. We acknowledge support for this publication by the DFG, Deutsche Forschungsgemeinschaft, and the Open Access Publication Fund of Bielefeld University.

Data Availability

The data sets generated during or analyzed during this study are available from the corresponding author upon reasonable request.

Conflicts of Interest

None declared.

Multimedia Appendix 1

PRISMA-ScR (Preferred Reporting Items for Systematic reviews and Meta-Analyses extension for Scoping Reviews) checklist. [DOC File, 120 KB - ai v4i1e59295 app1.doc]

Multimedia Appendix 2

Search terms for PubMed/MEDLINE and Web of Science. [DOC File , 37 KB - ai v4i1e59295 app2.doc]

References

- 1. Tang L, Li J, Fantus S. Medical artificial intelligence ethics: a systematic review of empirical studies. Digital Health 2023;9:20552076231186064 [FREE Full text] [doi: 10.1177/20552076231186064] [Medline: 37434728]
- Bitkina OV, Park J, Kim HK. Application of artificial intelligence in medical technologies: a systematic review of main trends. Digital Health 2023;9:20552076231189331 [FREE Full text] [doi: 10.1177/20552076231189331] [Medline: 37485326]
- 3. Walter Z, Lopez MS. Physician acceptance of information technologies: role of perceived threat to professional autonomy. Decis Support Syst 2008;46(1):206-215. [doi: 10.1016/j.dss.2008.06.004]
- 4. Harrison S, Ahmad WIU. Medical autonomy and the UK State 1975 to 2025. Sociology 2025;34(1):129-146. [doi: 10.1017/s0038038500000092]
- 5. Schulz R, Harrison S. Physician autonomy in the federal republic of Germany, Great Britain and the United States. Int J Health Plann Manage 1986;1(5):335-355. [doi: 10.1002/hpm.4740010504] [Medline: 10281783]
- 6. Marjoribanks T, Lewis JM. Reform and autonomy: perceptions of the Australian general practice community. Soc Sci Med 2003;56(10):2229-2239. [doi: 10.1016/s0277-9536(02)00239-3] [Medline: 12697211]
- Salvatore D, Numerato D, Fattore G. Physicians' professional autonomy and their organizational identification with their hospital. BMC Health Serv Res 2018;18(1):775 [FREE Full text] [doi: 10.1186/s12913-018-3582-z] [Medline: 30314481]

- Lambert SI, Madi M, Sopka S, Lenes A, Stange H, Buszello C, et al. An integrative review on the acceptance of artificial intelligence among healthcare professionals in hospitals. NPJ Digital Med 2023;6(1):111 [FREE Full text] [doi: 10.1038/s41746-023-00852-5] [Medline: 37301946]
- Eltawil FA, Atalla M, Boulos E, Amirabadi A, Tyrrell PN. Analyzing barriers and enablers for the acceptance of artificial intelligence innovations into radiology practice: a scoping review. Tomography 2023;9(4):1443-1455 [FREE Full text] [doi: 10.3390/tomography9040115] [Medline: <u>37624108</u>]
- 10. Vo V, Chen G, Aquino YSJ, Carter SM, Do QN, Woode ME. Multi-stakeholder preferences for the use of artificial intelligence in healthcare: a systematic review and thematic analysis. Soc Sci Med 2023;338:116357 [FREE Full text] [doi: 10.1016/j.socscimed.2023.116357] [Medline: <u>37949020</u>]
- Tricco AC, Lillie E, Zarin W, O'Brien KK, Colquhoun H, Levac D, et al. PRISMA extension for scoping reviews (PRISMA-ScR): checklist and explanation. Ann Intern Med 2018;169(7):467-473 [FREE Full text] [doi: 10.7326/M18-0850] [Medline: 30178033]
- Amann J, Vayena E, Ormond KE, Frey D, Madai VI, Blasimme A. Expectations and attitudes towards medical artificial intelligence: a qualitative study in the field of stroke. PLoS One 2023;18(1):e0279088 [FREE Full text] [doi: 10.1371/journal.pone.0279088] [Medline: 36630325]
- 13. Chen Y, Stavropoulou C, Narasinkan R, Baker A, Scarbrough H. Professionals' responses to the introduction of AI innovations in radiology and their implications for future adoption: a qualitative study. BMC Health Serv Res 2021;21(1):813 [FREE Full text] [doi: 10.1186/s12913-021-06861-y] [Medline: 34389014]
- Huang Z, George MM, Tan YR, Natarajan K, Devasagayam E, Tay E, et al. Are physicians ready for precision antibiotic prescribing? A qualitative analysis of the acceptance of artificial intelligence-enabled clinical decision support systems in India and Singapore. J Global Antimicrob Resist 2023;35:76-85 [FREE Full text] [doi: 10.1016/j.jgar.2023.08.016] [Medline: 37640155]
- 15. Jussupow E, Spohrer K, Heinzl A. Identity threats as a reason for resistance to artificial intelligence: survey study with medical students and professionals. JMIR Form Res 2022;6(3):e28750 [FREE Full text] [doi: 10.2196/28750] [Medline: 35319465]
- 16. Kocaballi AB, Ijaz K, Laranjo L, Quiroz JC, Rezazadegan D, Tong HL, et al. Envisioning an artificial intelligence documentation assistant for future primary care consultations: a co-design study with general practitioners. J Am Med Inform Assoc 2020;27(11):1695-1704 [FREE Full text] [doi: 10.1093/jamia/ocaa131] [Medline: 32845984]
- 17. Lombi L, Rossero E. How artificial intelligence is reshaping the autonomy and boundary work of radiologists. a qualitative study. Sociol Health Illn 2024;46(2):200-218. [doi: 10.1111/1467-9566.13702] [Medline: 37573551]
- 18. Wong WCW, Zhao IY, Ma YX, Dong WN, Liu J, Pang Q, et al. Primary care physicians' and patients' perspectives on equity and health security of infectious disease digital surveillance. Ann Fam Med 2023;21(1):33-39 [FREE Full text] [doi: 10.1370/afm.2895] [Medline: 36635084]
- 19. Varjus SL, Leino-Kilpi H, Suominen T. Professional autonomy of nurses in hospital settings—a review of the literature. Scand J Caring Sci 2011;25(1):201-207. [doi: <u>10.1111/j.1471-6712.2010.00819.x</u>] [Medline: <u>20707857</u>]
- Pursio K, Kankkunen P, Sanner-Stiehr E, Kvist T. Professional autonomy in nursing: an integrative review. J Nurs Manag 2021;29(6):1565-1577. [doi: <u>10.1111/jonm.13282</u>] [Medline: <u>33548098</u>]
- 21. Ali O, Abdelbaki W, Shrestha A, Elbasi E, Alryalat MAA, Dwivedi YK. A systematic literature review of artificial intelligence in the healthcare sector: benefits, challenges, methodologies, and functionalities. J Innovation Knowl 2023;8(1):100333. [doi: 10.1016/j.jik.2023.100333]
- 22. Alowais SA, Alghamdi SS, Alsuhebany N, Alqahtani T, Alshaya AI, Almohareb SN, et al. Revolutionizing healthcare: the role of artificial intelligence in clinical practice. BMC Med Educ 2023;23(1):689 [FREE Full text] [doi: 10.1186/s12909-023-04698-z] [Medline: 37740191]
- 23. Tursunbayeva A, Renkema M. Artificial intelligence in health care: implications for the job design of healthcare professionals. Asia Pac J Human Res 2022;61(4):845-887. [doi: 10.1111/1744-7941.12325]
- 24. Li J, Dada A, Puladi B, Kleesiek J, Egger J. ChatGPT in healthcare: a taxonomy and systematic review. Comput Methods Programs Biomed 2024;245:108013 [FREE Full text] [doi: 10.1016/j.cmpb.2024.108013] [Medline: 38262126]
- 25. Thirunavukarasu AJ, Hassan R, Mahmood S, Sanghera R, Barzangi K, El Mukashfi M, et al. Trialling a large language model (ChatGPT) in general practice with the applied knowledge test: observational study demonstrating opportunities and limitations in primary care. JMIR Med Educ 2023;9:e46599 [FREE Full text] [doi: 10.2196/46599] [Medline: 37083633]

Abbreviations

AI: artificial intelligence
CDSS: clinical decision support system
GP: general practitioner
PRISMA: Preferred Reporting Items for Systematic Reviews and Meta-Analyses



```
https://ai.jmir.org/2025/1/e59295
```

Edited by D Manuel; submitted 08.04.24; peer-reviewed by E Rossero, B Mesko; comments to author 24.04.24; revised version received 15.05.24; accepted 31.12.24; published 13.03.25. <u>Please cite as:</u> Grosser J, Düvel J, Hasemann L, Schneider E, Greiner W Studying the Potential Effects of Artificial Intelligence on Physician Autonomy: Scoping Review JMIR AI 2025;4:e59295 URL: https://ai.jmir.org/2025/1/e59295 doi:10.2196/59295 PMID:

©John Grosser, Juliane Düvel, Lena Hasemann, Emilia Schneider, Wolfgang Greiner. Originally published in JMIR AI (https://ai.jmir.org), 13.03.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Striking a Balance: Innovation, Equity, and Consistency in AI Health Technologies

Eric Perakslis^{1,2}, PhD; Kimberly Nolen³, BS, PharmD; Ethan Fricklas¹, MSE; Tracy Tubb¹, RN, MSN

¹Duke Clinical Research Institute, Duke University School of Medicine, Durham, NC, United States ²Pluto Health, Durham, NC, United States ³Pfizer Inc, New York, NY, United States

Corresponding Author:

Ethan Fricklas, MSE Duke Clinical Research Institute, Duke University School of Medicine, Durham, NC, United States

Abstract

With the explosion of innovation driven by generative and traditional artificial intelligence (AI), comes the necessity to understand and regulate products that often defy current regulatory classification. Tradition, and lack of regulatory expediency, imposes the notion of force-fitting novel innovations into pre-existing product classifications or into the essentially unregulated domains of wellness or consumer electronics. Further, regulatory requirements, levels of risk tolerance, and capabilities vary greatly across the spectrum of technology innovators. For example, currently unregulated information and consumer electronic suppliers set their own editorial and communication standards without extensive federal regulation. However, industries like biopharma companies are held to a higher standard in the same space, given current direct-to-consumer regulations like the Sunshine Act (also known as Open Payments), the federal Anti-Kickback Statute, the federal False Claims Act, and others. Clear and well-defined regulations not only reduce ambiguity but facilitate scale, showcasing the importance of regulatory clarity in fostering innovation and growth. To avoid highly regulated industries like health care and biopharma from being discouraged from developing AI to improve patient care, there is a need for a specialized framework to establish regulatory evidence for AI-based medical solutions. In this paper, we review the current regulatory environment considering current innovations but also pre-existing legal and regulatory responsibilities of the biopharma industry and propose a novel, hybridized approach for the assessment of novel AI-based patient solutions. Further, we will elaborate the proposed concepts via case studies. This paper explores the challenges posed by the current regulatory environment, emphasizing the need for a specialized framework for AI medical devices. By reviewing existing regulations and proposing a hybridized approach, we aim to ensure that the potential of AI in biopharmaceutical innovation is not hindered by uneven regulatory landscapes.

(JMIR AI 2025;4:e57421) doi:10.2196/57421

KEYWORDS

artificial intelligence; algorithm; regulatory landscape; predictive model; predictive analytics; predictive system; practical model; machine learning; large language model; natural language processing; deep learning; digital health; regulatory; health technology

Introduction

Background

The convergence of algorithms, artificial intelligence (AI), big data, and digital health technologies (DHTs) is a sea change not seen since the "dot.com" era, which has significantly changed the way we work, play, and learn [1,2]. However, the lack of comprehensive regulatory guidance has led to the force-fitting of novel innovations into existing categories, leading to ambiguous boundaries between medical devices and consumer electronics. This results in added ambiguity for innovators seeking to share valuable product concepts. What is lacking is a comprehensive approach to evaluating medical benefits, risks, and evidence that can be universally applied across different product categories, product types, and regulatory regimes [3]. Such an approach would be flexible, allowing for the distinctions

RenderX

between various products to be properly addressed. Clear regulations play an important role in enabling easier scaling, highlighting the mutually beneficial relationship between regulatory clarity and the acceleration of innovation. Moreover, the pharmaceutical industry's comprehensive understanding of scaling and marketing extends beyond the confines of drug development, presenting a valuable paradigm for other sectors in the AI landscape. This paper addresses the ambiguity faced by innovators and proposes models for evidence strategies, particularly focusing on the distinct regulatory challenges faced by the biopharma industry.

At the time of this writing, Gartner [4] has placed the increasingly popular generative AI technology at the peak of inflated expectations for emerging technologies in 2024. Whether the transformation that generative or other forms of AI and DHTs bring to health care occurs gradually or rapidly,

there is widespread anticipation of both advancement and potential challenges [5,6]. The many use cases ranging across diagnostics, logistics, clerical improvements, and new treatment modalities in general medicine and across medical subspecialties have been described in detail within the literature [7,8]. Similarly, the use of these technologies holds equally compelling promise within biomedical product development [9,10].

Current Challenges

AI and DHTs represent a wide range of intricate and interconnected technologies, and the vast array of applications are equally diverse and complex. For example, most DHTs rely on proprietary algorithms trained from and across a mixture of private and public data sources. As multiple technologies are integrated into a tool, more information may be shared, and the capabilities and risks aggregate [11]. This additive complexity challenges traditional domain-based regulatory regimes. The US Food and Drug Administration (FDA) regulates medical devices, but not all algorithms and apps meet the regulatory definition of a medical device as defined in Section 201(h) of the Food, Drug, and Cosmetic Act. These apps often access and transmit data across the internet but do not fit neatly into the codified remit of the Federal Communications Commission or the Federal Trade Commission (FTC) [12]. Even within regulatory regimes, there are qualified gaps. The Office for Civil Rights enforces health care privacy, but only for covered entities, leaving a loophole and resulting in gaps in protection that are

systematically being exploited [13]. Further, these overlapping, complex, and intricate interregulatory and intraregulatory regimes create confusion and inequities, hindering progress.

Objectives

While there have been efforts to establish standardized approaches to the regulatory assessment of DHT and AI medical products, many of these frameworks take the approach of a single regulator and regulatory regime versus approaches that inform regulatory decision-making across the spectrum of relevant regulators [14-16]. Further, these frameworks incorrectly assume that all innovators are alike. Consumer electronic companies and health technology startups, providing solutions that may overlap or compete with offerings from traditional medical device and pharmaceutical companies, often navigate a regulatory landscape that offers them more adaptability in their operations, which differs from the established health care regulations governing other sectors. The coexistence of unregulated and highly regulated makers in the same market can lead to various challenges, including issues related to safety, quality, and fair competition. Balancing innovation and oversight is crucial in this context. We need solutions that promote fair competition while maintaining a high standard of safety and product effectiveness, without creating a disparity between the heavily and lightly regulated entities. It is also helpful to level set on terminology. Textbox 1 provides definitions for common terms used in this space.

Textbox 1. Key terms and definitions.

Digital health technologies: technologies consisting of hardware (eg, sensors or transmitters) or software (eg, connectivity software, algorithms, or artificial intelligence) components that are used for health care–related purposes.

Medical device software: term primarily used in the European Union to define software with a medical purpose that can be used either alone or in combination with a regulated medical device. This is not interchangeable with software as a medical device [17].

Software as a medical device: software that is used for medical purposes and may do so independently of a hardware medical device as well as not being a required component of a hardware medical device [18].

Mobile medical apps: mobile apps serving as medical devices, which integrate software functionality that aligns with the Food and Drug Administration definition of a device, as outlined in Section 201(h) of the Food, Drug, and Cosmetic Act. These apps may function as accessories to regulated medical devices or convert a mobile platform into a regulated medical device [19].

Digital therapeutics: software-based interventions intended to prevent, manage, or treat medical conditions based on evidence of a demonstrable positive therapeutic impact on a patient's health [20].

Direct to consumer: marketing products or services directly to consumers without the involvement of a health care provider.

Regulatory Regimes and Industries in DHT

The key is that not all makers are subject to regulations in the same manner nor do they exhibit the same affinity for risk. For example, while direct-to-consumer advertising is highly regulated for pharmaceutical products, the oversight is less consistent for over-the-counter (OTC) "devices," such as some medical tests not regulated by the FDA or FTC [21,22]. This not only results in a less than comprehensive regulatory coverage of AI medical devices and DHTs but also involuntarily creates an ecosystem where makers develop and market their products around the varying gaps in regulatory coverage. Certain products, such as OTC medical device algorithms to detect sleep apnea, may be subject to less regulation and thus have an

advantage over established health care products. OTC sleep apnea devices represent a category of products that can fall in the "interstitial spaces" of regulatory oversight, as they are not always subject to the same level of scrutiny as prescription devices. These products often include wearable sensors, smartphone apps, or other consumer-grade devices that purport to detect symptoms of sleep apnea, such as disrupted breathing or low oxygen levels during sleep. Many OTC sleep apnea devices may fall into class I or II and thus may not require premarket approval, which is the most stringent type of device marketing application required by the FDA. Instead, they may only need to meet the requirements for 510(k) clearance, which is a less rigorous process and does not require clinical trials. However, there are also some OTC sleep apnea devices that do not fall under any FDA regulation because they are marketed as "wellness" or "lifestyle" products rather than medical devices.

They do not detect signs of sleep apnea and are not marketed as a medical device but as a sleep improvement system; therefore, they do not fall under FDA regulation.

This lack of consistent regulation can create opportunities for companies to market products with less oversight and potentially greater profit margins, but it can also lead to consumer confusion and potential safety risks if the products do not perform as advertised. It is a clear example of how the existing regulatory framework may struggle to keep up with the rapid pace of innovation in DHT. Clear and well-defined regulations play a pivotal role, especially during the transition from exploratory phases to scaling products, enabling smoother, efficient scaling processes.

In many ways, the opposite situation exists for larger established health care organizations. A highly regulated pharmaceutical company that is already subject to the many previously discussed compliance regimes and other complex corporate regulatory obligations may find it too difficult or risky to attempt digital innovation, as the burden of reporting, evidence, and oversight are all greatly heightened compared to niche innovators. This scenario must be discouraged, as these organizations have deep expertise within the disease areas where they have successfully delivered drugs and devices. This expertise could lend itself to success beyond many health technology startups, which often fail due to a lack of market fit of their products [23]. Encouraging a balance between regulatory clarity and flexibility is paramount to fostering innovation across diverse players in the digital health landscape. Indeed, there is a cost to regulatory compliance, which is more readily absorbed by well-resourced companies. Smaller startups may not have sufficient funding to run the optimal size and scale validation study. They may have funding constrained by the need to showing promise to investors in order to survive to their next round of funding.

Evidence Requirements and Claims

While one side of the matrix is the nature of the products being developed and the types of makers, the evidence supporting these products is equally if not more diverse. Companies with very different sizes and areas of expertise may be competing openly within a range of product categories and evidence strategies with clinical development plans that seem lacking. Much of this may be attributed to the relative lack of maturation of the AI medical device and DHT spaces, which has led to a wide range of interpretation of the guidelines. This can pose challenges for companies with more rigid regulatory boundaries, which wish to participate in this evolving experimental domain and have substantial evidence strategies to support product development but are uncomfortable as the space is not mature. For most, the first step is to determine the type of product being developed. However, when dealing with products designed for medical purposes or functions, it is essential to ensure that it addresses an unmet medical need. In the United States, the product type can vary from a device software or algorithm that may be classified as mobile medical apps, software functions that are not medical devices, clinical decision support software, or software as a medical device (SaMD) [24-26]. Each of these product types needs different types and levels of evidence to support them in the market and may need regulatory approval. While the FDA offers guidance on how to determine the product type, significant judgment is required due to similarities within the categories as well as the severity of the disease indication. This necessitates an iterative thought process, considering multiple regulatory guidance alongside the evidence strategy and clinical development plan [27].

To simplify this process, we developed the graphical consolidated regulatory decision framework shown in Figure 1 [28-30]. This framework builds on the approach in the FDA Guidance, Software as a Medical Device (SAMD): Clinical Evaluation [31]. Additional details supporting this framework are available in Multimedia Appendix 1. A precursor to the workflow is determining whether the clinical association is well-established or novel. This can be nontrivial, as many SaMD products lack clinically established standards due to the novelty of the product. When there is a well-established clinical association, these SaMD have outputs with well-documented association as identified in sources such as clinical guidelines, clinical studies in peer-reviewed journals, consensus for the use of the SaMD, international reference materials, or other similar well-established comparators of previously marketed devices. When the clinical association is novel, these SaMD may involve new inputs, algorithms, outputs, new intended target population, or new intended use. An example may include the combination of nonstandard inputs (eg, mood or pollen count), with standard inputs (eg, blood pressure or other physiological signals), that uses novel algorithms to detect deterioration of health or diagnosis of a disease.



Figure 1. Consolidated regulatory classification decision framework (for the reader's convenience, Multimedia Appendix 1 gives a set of figures referenced in Figure 1). CDS: clinical decision support; IVD: in vitro diagnostic; SaMD: software as a medical device.



The importance of objective consideration of these questions cannot be overstated. Frequently, innovators have embedded biases within their assumptions that cloud judgment in this assessment.

Any one or combination of these biases can threaten or derail the development of novel technology products including product integrity and patient safety. For example, the well-publicized case of the FDA warning letter that caused Owlet to cease selling their Smart Sock and copackaged products includes several of these biases [32]. The FDA determined that the product was a medical device but the maker had not reached the same determination [33]. According to the FDA, the apnea alarms had inadequate clinical evidence, and parents could potentially seek emergency care due to product alarms that had inadequately established clinical association (prestep), specifically dips in oxygen saturation as determined by pulse oximetry in infants during sleep [34]. This incident exemplifies how the misclassification of a product type and inadequate clinical evidence highlight the challenges in navigating ambiguous regulatory guidelines. Clear regulations not only mitigate the risk of inadequate clinical evidence, reducing the likelihood of triggering unnecessary care in this case, but also highlight the significance of aligning evidence strategies with robust clinical development plans to avoid such pitfalls.

In contrast, Apple's approach to irregular heartbeat notification serves as an example that a technology company that operates

https://ai.jmir.org/2025/1/e57421

RenderX

outside the traditional health care sphere can diligently address product-type classification and evidence. Apple developed these FDA-cleared features via rigorous and traditional approaches using randomized controlled trials [35,36]. The resulting product label is considered FDA-regulated; therefore, the claims made about the product must not exceed the evidence produced, the specifics listed within the clearance letter, or the resulting label [37]. However, Apple is an exception within the technology industry in size, scope, and resourcing. The average medical device maker is much smaller and must build their product strategies around regulatory regimes in order to get their product to market. This is not solely about the avoidance of regulation. Many innovations pass the bar for regulatory clearance but not the bar for reimbursement, which can be a difficult and unprofitable market situation. Alternatively, developing a product that does not meet the bar for FDA regulation can result in a highly profitable "health" or "fitness" application that, while not regulated or reimbursed by health insurance, can be highly profitable by volume of sales even at very low-price points.

If the maker determines that the software product is not a medical device, the next step is to determine whether the product is a decision support tool, and this can be accomplished with the aid of the earlier-referenced FDA guidance document.

Continuing with the framework, if the software is a medical device, the next step is to determine whether the software is a SaMD. If not, the user is directed toward traditional medical

device pathways. If so, they are guided through the SaMD categorization process. This is complicated by the requirement for a deep understanding of the medical indication and all possible outcomes from any chosen route. Determining whether a SaMD treats or diagnoses, whether it drives clinical management, or simply informs clinical care yields a level of interpretation that often varies by stakeholders. In addition, it must be simultaneously determined whether these actions occur in a critical, serious, or nonserious medical situation. The FDA has published guidance on when and whether independent review of these decisions should be included.

Once the SaMD is categorized, the next step is to determine whether the product is an in vitro diagnostic (IVD) or non-IVD SaMD. The rubric directs the user toward evidence generation in either case. When the SaMD is non-IVD and novel clinically, the evidence generation process requires the establishment of a valid clinical association, analytical validation, and clinical validation of the product. These steps have been elaborated in detail within the previously referenced work by Goldsack et al [16]. When the SaMD is an IVD, the evidence generation process is analogous to the non-IVD case and can follow the stages of the clinical evidence assessment process as outlined in the Global Harmonization Task Force's Clinical Evidence for IVD Medical Devices [38].

Product Labeling Standards as a Guide

One improbable solution to creating clear and concise regulations would be the reorganization of current regulatory regimes to produce a new agency focused directly on the regulation of health care technologies. However, we can gain some insights from the nutrition industry.

Today, there are 3 agencies responsible for the regulation of food and nutrition information, the FDA, the FTC, and the Food Safety and Inspection Service (FSIS) of the US Department of Agriculture [39]. This may appear logical, assuming that each agency shares part of an overall mission. However, the reality is that handoffs, overlapping, and gray areas decrease the regulatory effectiveness or, at minimum, create confusion of roles and responsibilities. Continuing the example of food labeling, the FTC regulates food advertising, while the other 2 agencies share responsibility for regulating labels: FSIS regulates meat, poultry, and egg labeling and FDA regulates labeling for all other foods and nonspecified red meat (game). The Nutrition Labeling and Education Act addressed FDA-regulated packages and FSIS-issued parallel regulations. As an example of a gap, there are no provisions in the regulatory authorities defined by Congress that allow the FDA to approve dietary supplements for safety before they reach the consumer [40]. This results in fragmented safety data and little ability to forecast or prevent harmful products from reaching consumers [41].

There is a great deal that the digital health space can learn from nutrition and food labeling. Specifically, the FDA or Nutrition Labeling and Education Act and the FSIS oversee 3 elements of food package labeling: nutrient content, nutrient content claims, and "disease" claims. Further, the FDA has restricted health claims to a small number of permitted claims [42]. This type of strategic and comprehensive approach to labeling is a model that if applied to AI medical devices and DHTs would improve transparency and clarify their benefits and risks. Elements are starting to appear in relevant subdomains of digital health such as cybersecurity, computing hardware, clinical decision support, and medical devices [43-45]. Examples extending the nutrition-based health claims into the digital domain are shown in Figure 2. This figure shows health risk areas linked to nutrients having FDA-approved health claims. The figure shows those same health risk areas linked to digital health elements that could also have an impact on risk. Standardized AI medical device and DHT product labeling across prescription and nonprescription products could address the diversity of innovators and makers just as nutrition labeling levels the field between small farms and large, industrial food production corporations, as all are held to the same standards.

Different aspects of the creation, oversight, and enforcement of such labeling regulations would likely fall within the purview of existing regulatory bodies assessing medical devices and algorithm-driven DHTs. As Table 1 suggests, this is a current patchwork of different agencies without necessarily one central authority. The table indicates how the different agencies could each play a role in the oversight of AI medical devices and DHTs and the relevant product development.



Figure 2. Health impacts and their associated nutrients and digital elements. FDA: Food and Drug Administration.





Table . Regulatory authority across components of medical devices and algorithm-driven digital health technologies (DHTs).

Federal agency	Relevant aspects of current role	Possible enhanced role
US Food and Drug Administration	Regulation of medical devices, MMAs ^a , SaMDs ^b	Develop specific guidelines for labeling AI ^c medical devices, including continuous learning algorithms, and establish a clear pathway for marketing authorizations and DTC ^d
Federal Trade Commission	Enforcement of federal laws that prevent fraud, deception, and unfair business practices (as in advertising)	Regulate use of labels in DTC advertising of AI medical devices and DHTs
Centers for Medicare & Medicaid Services	Administers the Open Payments program, enhanc- ing transparency by collecting and publicly dis- closing information about financial relationships between health care providers and pharmaceutical or medical device manufacturers	Outline required transparency elements in the label, capturing health care provider financial interactions with companies in the design, testing, and use of DHTs
US Department of Health and Human Services, Office for Civil Rights	Ensures the privacy and security of protected health information	Ensure that the AI medical device label provides transparency in how AI devices use patient data and enforce penalties for noncompliance
National Institute of Standards and Technology	Develops and maintains technical standards	Establish the benchmarks cited on the label for AI performance, safety, and interoperability with other medical systems
US Consumer Product Safety Commission	Protects the public from the risk of injury or death associated with consumer products	Issue product recalls for all other DHTs or AI- enabled health care devices not under the Food and Drug Administration's purview
_		

^aMMA: mobile medical app.

^bSaMD: software as a medical device.

^cAI: artificial intelligence.

^dDTC: direct to consumer.

Al-Based DHT Package Labeling

Regardless of the type of maker, a significant unmet need that can inform innovators' strategies is standardized package labeling. Based upon the known potential harms and known potential limitations in AI-based DHTs, the minimum content for analogous product labeling would include the type of algorithm, the framework used for evidence generation, qualification and quantification of reproducibility, the ethical framework used, a description of the data used to train the model including how it was collected and consented, a statement on how bias was minimized or quantified, the risk management framework used to aggregate these various elements, and performance metrics [46-52]. This would enable a DHT packaging label, similar to the example shown in Table 2.

Table . Example of a permitted claims approach to artificial intelligence (AI)-based digital health technologies (DHTs).

DHT or algorithm design element	Permitted health claim	Example
Type of evidence basis	Primary benefit or utility	Randomized controlled trials, real-world evi- dence, etc
Ethical framework	Population benefit-risk	Inclusion and exclusion criteria as well as inter- pretability and explainability
Reproducibility	Statistically quantified and qualified claims	Specific indications and efficacy, data lineage, model versioning
Training data description	Applicability and specificity to populations	Rationale for included and excluded populations, training and testing data split
Disclosure of bias	Limitations of use and contraindications	Phenotypic traits such as skin tone
Risk management framework	Product integrity	Cybersecurity resilience, prevention of AI poi- soning, measures for protecting user data (such as differential privacy)
Performance	Primary benefit or utility	Sensitivity, specificity, negative predictive value, positive predictive value

Labeling would directly counter the real and perceived black box problem and inform clinicians, researchers, patients, and caregivers in a manner that is equivalent to how they study, learn, and use new prescription and OTC drugs, diagnostics, and medical devices today [53,54].

New regulatory frameworks often face pushback from those they regulate, and labeling regulations are likely no exception. This resistance can be mitigated somewhat with guidance documents that incorporate feedback from various manufacturers, educational initiatives, and avenues for direct interaction between companies and regulatory bodies.

Conclusions

AI medical products and DHTs hold immense promise, but the diverse regulatory constraints among product makers necessitate a standardized approach. This is especially critical for smaller AI developers who operate in a different landscape than health care or larger industries. To create consistency, adopting minimum product labeling requirements, understanding claims, and having substantial evidence plans become essential. As innovation accelerates, ensuring equity in the ecosystem will allow both emerging and mature technology innovators to contribute meaningfully without being hindered by not only regulatory disparities but also ambiguities, such as the uncertain classification of certain AI applications and the lack of clear communication standards. Future research exploring the various reimbursement strategies and ethical implications of AI across product makers would be valuable to providing a more complete picture of this space. The perspectives from a wide range of digital health ecosystem stakeholders should be included to ensure that their diverse needs and expectations are being addressed.

Acknowledgments

This research was a collaboration between Duke University Clinical Research Institute and Pfizer Inc. This research was funded by Pfizer Inc.

Authors' Contributions

EP, KN, EF, and TT contributed to the conceptualization, writing of the original draft, and the reviewing and editing of the manuscript. EP, EF, and TT prepared and finalized all figures. All authors have read and approved the final manuscript.

Conflicts of Interest

KN is an employee and shareholder of Pfizer Inc. EP was affiliated with Duke Clinical Research Institute at the time of this initiative and is currently affiliated with Pluto Health.

Multimedia Appendix 1

Set of linked figures referenced in Figure 1. [PPTX File, 10310 KB - <u>ai_v4i1e57421_app1.pptx</u>]

References

- 1. Geier B. What did we learn from the dotcom stock bubble of 2000? Time. 2015 Apr 14. URL: <u>https://time.com/3741681/</u> 2000-dotcom-stock-bust [accessed 2024-01-07]
- 2. Dewey C. 36 ways the web has changed us. The Washington Post. 2014 Mar 12. URL: <u>https://www.washingtonpost.com/</u> news/arts-and-entertainment/wp/2014/03/12/36-ways-the-web-has-changed-us [accessed 2024-01-07]
- 3. What is digital health? US Food and Drug Administration. 2020. URL: <u>https://www.fda.gov/medical-devices/</u> <u>digital-health-center-excellence/what-digital-health</u> [accessed 2024-01-07]
- 4. What's driving the Hype Cycle for generative AI, 2024. Gartner. 2024 Nov 14. URL: <u>https://www.gartner.com/en/articles/</u> <u>hype-cycle-for-genai</u> [accessed 2024-01-07]
- 5. Yu KH, Beam AL, Kohane IS. Artificial intelligence in healthcare. Nat Biomed Eng 2018 Oct;2(10):719-731. [doi: 10.1038/s41551-018-0305-z] [Medline: 31015651]
- Haug CJ, Drazen JM. Artificial intelligence and machine learning in clinical medicine, 2023. N Engl J Med 2023 Mar 30;388(13):1201-1208. [doi: <u>10.1056/NEJMra2302038</u>] [Medline: <u>36988595</u>]
- Rajpurkar P, Chen E, Banerjee O, Topol EJ. AI in health and medicine. Nat Med 2022 Jan;28(1):31-38. [doi: 10.1038/s41591-021-01614-0] [Medline: 35058619]
- Jiang F, Jiang Y, Zhi H, et al. Artificial intelligence in healthcare: past, present and future. Stroke Vasc Neurol 2017 Dec;2(4):230-243. [doi: <u>10.1136/svn-2017-000101</u>] [Medline: <u>29507784</u>]
- 9. Paul D, Sanap G, Shenoy S, Kalyane D, Kalia K, Tekade RK. Artificial intelligence in drug discovery and development. Drug Discov Today 2021 Jan;26(1):80-93. [doi: 10.1016/j.drudis.2020.10.010] [Medline: 33099022]
- 10. Smith GF. Artificial intelligence in drug safety and metabolism. Methods Mol Biol 2022;2390:483-501. [doi: 10.1007/978-1-0716-1787-8_22] [Medline: 34731484]
- 11. Sivakumar CLV, Mone V, Abdumukhtor R. Addressing privacy concerns with wearable health monitoring technology. WIREs Data Min Knowl 2024 May;14(3):e1535. [doi: <u>10.1002/widm.1535</u>]
- 12. McKinnon J, Kendall B. Is the FTC up to the task of internet regulation? The Wall Street Journal. 2017 Dec 15. URL: https://www.wsj.com/articles/is-ftc-up-to-the-task-of-internet-regulation-1513349967 [accessed 2024-01-07]

- Mandl KD, Perakslis ED. HIPAA and the leak of "deidentified" EHR data. N Engl J Med 2021 Jun 10;384(23):2171-2173. [doi: <u>10.1056/NEJMp2102616</u>] [Medline: <u>34110112</u>]
- 14. Silberman J, Wicks P, Patel S, et al. Rigorous and rapid evidence assessment in digital health with the evidence DEFINED framework. NPJ Digit Med 2023 May 31;6(1):101. [doi: 10.1038/s41746-023-00836-5] [Medline: 37258851]
- Coravos A, Doerr M, Goldsack J, et al. Modernizing and designing evaluation frameworks for connected sensor technologies in medicine. NPJ Digit Med 2020;3:37. [doi: <u>10.1038/s41746-020-0237-3</u>] [Medline: <u>32195372</u>]
- Goldsack JC, Coravos A, Bakker JP, et al. Verification, analytical validation, and clinical validation (V3): the foundation of determining fit-for-purpose for Biometric Monitoring Technologies (BioMeTs). NPJ Digit Med 2020;3:55. [doi: <u>10.1038/s41746-020-0260-4]</u> [Medline: <u>32337371</u>]
- 17. Software as a medical device: possible framework for risk categorization and corresponding considerations. International Medical Device Regulators Forum. 2014 Sep 18. URL: <u>http://www.imdrf.org/docs/imdrf/final/technical/</u> imdrf-tech-140918-samd-framework-risk-categorization-141013.pdf [accessed 2024-01-07]
- Guidance on qualification and classification of software in regulation (EU) 2017/745—MDR and regulation (EU) 2017/746—IVDR. Medical Device Coordination Group. URL: <u>https://health.ec.europa.eu/system/files/2020-09/</u>md mdcg 2019 11 guidance qualification classification software en 0.pdf [accessed 2025-03-24]
- 19. Device software functions, including mobile medical applications. US Food and Drug Administration. URL: <u>https://www.fda.gov/medical-devices/digital-health-center-excellence/device-software-functions-including-mobile-medical-applications</u> [accessed 2024-01-07]
- 20. Digital therapeutics (DTx). European Data Protection Supervisor. URL: <u>https://edps.europa.eu/press-publications/publications/</u> techsonar/

digital-therapeutics-dtx_en#:~:text=Digital%20Therapeutics%20(DTx)%20are%20evidence,a%20medical%20disorder%20or%20disease [accessed 2024-01-07]

- 21. Refuse to accept policy for 510(k)s guidance for industry and Food and Drug Administration staff. US Food and Drug Administration. 2022. URL: <u>https://www.fda.gov/regulatory-information/search-fda-guidance-documents/</u> refuse-accept-policy-510ks [accessed 2024-01-07]
- 22. From the manufacturers' mouth to your ears: direct to consumer advertising. US Food and Drug Administration. 2015. URL: <u>https://www.fda.gov/drugs/special-features/manufacturers-mouth-your-ears-direct-consumer-advertising</u> [accessed 2024-01-07]
- 23. Rigg K. The top 3 reasons health care startups fail. High Tech World. 2022. URL: <u>https://www.htworld.co.uk/news/</u> <u>the-top-3-reasons-health-tech-startups-fail</u> [accessed 2025-03-24]
- 24. Examples of software functions that are not medical devices. US Food and Drug Administration. 2022. URL: <u>https://www.fda.gov/medical-devices/device-software-functions-including-mobile-medical-applications/</u> examples-software-functions-are-not-medical-devices#~:text=Software%20functions%20that%20are%20intended,are%20not%20intended%20for%20use [accessed 2024-01-07]
- 25. Clinical Decision Support Software. US Food and Drug Administration. 2022. URL: <u>https://www.fda.gov/</u> regulatory-information/search-fda-guidance-documents/clinical-decision-support-software [accessed 2024-01-07]
- 26. Software as a Medical Device (SaMD). US Food and Drug Administration. 2018. URL: <u>https://www.fda.gov/medical-devices/</u> <u>digital-health-center-excellence/software-medical-device-samd</u> [accessed 2024-01-07]
- 27. Classify your medical device. US Food and Drug Administration. 2020. URL: <u>https://www.fda.gov/medical-devices/</u> overview-device-regulation/classify-your-medical-device [accessed 2024-01-07]
- 28. Guidance—MDCG endorsed documents and other guidance. European Commission—Public Health. URL: <u>https://health.</u> <u>ec.europa.eu/medical-devices-sector/new-regulations/guidance-mdcg-endorsed-documents-and-other-guidance_en</u> [accessed 2024-01-07]
- 29. Clinical Decision Support Software Frequently Asked Questions (FAQs). US Food and Drug Administration. 2022. URL: <u>https://www.fda.gov/medical-devices/software-medical-device-samd/</u> <u>your-clinical-decision-support-software-it-medical-device</u> [accessed 2024-01-07]
- 30. Is your software a medical device? European Commission. 2021 Mar 23. URL: <u>https://health.ec.europa.eu/system/files/</u>2021-03/md_mdcg_2021_mdsw_en_0.pdf [accessed 2025-03-24]
- 31. Software as a medical device (SAMD): clinical evaluation—guidance for industry and Food and Drug Administration staff. US Food and Drug Administration. 2017 Dec 8. URL: <u>https://www.fda.gov/regulatory-information/</u> search-fda-guidance-documents/software-medical-device-samd-clinical-evaluation [accessed 2024-01-07]
- Marchante M. Owlet Smart Sock stops U.S. sales after FDA warning. Miami Herald. 2021 Dec. URL: <u>https://www.miamiherald.com/news/recalls/article256291067.html</u> [accessed 2025-04-01]
- 33. Warning letter to Owlet Baby Care, Inc. US Food and Drug Administration. URL: <u>https://www.fda.gov/inspections-compliance-enforcement-and-criminal-investigations/warning-letters/owlet-baby-care-inc-616354-10052021</u> [accessed 2025-03-24]
- 34. Bonafide CP, Jamison DT, Foglia EE. The emerging market of smartphone-integrated infant physiologic monitors. JAMA 2017 Jan 24;317(4):353-354. [doi: 10.1001/jama.2016.19137] [Medline: 28118463]

- 35. Direct-to-consumer tests. US Food and Drug Administration. 2018. URL: <u>https://www.fda.gov/medical-devices/</u> in-vitro-diagnostics/direct-consumer-tests [accessed 2025-03-24]
- Perez MV, Mahaffey KW, Hedlin H, et al. Large-scale assessment of a smartwatch to identify atrial fibrillation. N Engl J Med 2019 Nov 14;381(20):1909-1917. [doi: <u>10.1056/NEJMoa1901183</u>] [Medline: <u>31722151</u>]
- 37. Statement from FDA Commissioner Scott Gottlieb, M.D., and Center for Devices and Radiological Health Director Jeff Shuren, M.D., J.D., on agency efforts to work with tech industry to spur innovation in digital health. US Food and Drug Administration. 2018. URL: <u>https://www.fda.gov/news-events/press-announcements/</u>
- statement-fda-commissioner-scott-gottlieb-md-and-center-devices-and-radiological-health-director [accessed 2024-01-07]
 38. Clinical evidence for IVD medical devices. International Medical Device Regulators Forum. 2012 Nov. URL: <u>https://www.imdrf.org/sites/default/files/docs/ghtf/final/sg5/technical-docs/</u>
- <u>ghtf-sg5-n6-2012-clinical-evidence-ivd-medical-devices-121102.pdf</u> [accessed 2024-01-07]
 39. Consumer use of information: implications for food policy. US Department of Agriculture, Economic Research Service. URL: <u>https://www.ers.usda.gov/webdocs/publications/41905/51665_ah715c.pdf?v=0</u> [accessed 2025-04-01]
- 40. Questions and answers on dietary supplements. US Food and Drug Administration. 2022. URL: <u>https://www.fda.gov/food/</u> information-consumers-using-dietary-supplements/questions-and-answers-dietary-supplements [accessed 2024-01-07]
- 41. Beach C. Report finds enormous increase in number of food items recalled in 2022. Food Safety News. 2023 Mar 15. URL: https://www.foodsafetynews.com/2023/03/report-finds-enormous-increase-in-number-of-food-items-recalled-in-2022/ [accessed 2024-01-07]
- 42. Label claims for conventional foods and dietary supplements. US Food and Drug Administration. 2022. URL: <u>https://www.fda.gov/food/food-labeling-nutrition/label-claims-conventional-foods-and-dietary-supplements</u> [accessed 2024-01-07]
- 43. Lemos R. Cybersecurity 'nutrition' labels still a work in progress. Dark Reading. 2022 Nov 11. URL: <u>https://www.darkreading.com/dr-tech/cybersecurity-nutrition-labels-still-a-work-in-progress</u> [accessed 2024-01-07]
- 44. Sendak MP, Gao M, Brajer N, Balu S. Presenting machine learning model information to clinical end users with model facts labels. NPJ Digit Med 2020;3:41. [doi: <u>10.1038/s41746-020-0253-3</u>] [Medline: <u>32219182</u>]
- 45. A new role for the FDA in medical device regulation. The Center for Growth and Opportunity. 2019 Jun 17. URL: <u>https://www.thecgo.org/research/a-new-role-for-the-fda-in-medical-device-regulation/</u> [accessed 2024-01-07]
- Meskó B, Topol EJ. The imperative for regulatory oversight of large language models (or generative AI) in healthcare. NPJ Digit Med 2023 Jul 6;6(1):120. [doi: <u>10.1038/s41746-023-00873-0</u>] [Medline: <u>37414860</u>]
- 47. Cruz Rivera S, Liu X, Chan AW, Denniston AK, Calvert MJ, SPIRIT-AI and CONSORT-AI Working Group. Guidelines for clinical trial protocols for interventions involving artificial intelligence: the SPIRIT-AI extension. Lancet Digit Health 2020 Oct;2(10):e549-e560. [doi: 10.1016/S2589-7500(20)30219-3] [Medline: 33328049]
- 48. Beam AL, Manrai AK, Ghassemi M. Challenges to the reproducibility of machine learning models in health care. JAMA 2020 Jan 28;323(4):305-306. [doi: <u>10.1001/jama.2019.20866</u>] [Medline: <u>31904799</u>]
- 49. Murphy K, Di Ruggiero E, Upshur R, et al. Artificial intelligence for good health: a scoping review of the ethics literature. BMC Med Ethics 2021 Feb 15;22(1):14. [doi: 10.1186/s12910-021-00577-8] [Medline: 33588803]
- Frenkel S, Thompson S. 'Not for machines to harvest': data revolts break out against A.I. The New York Times. 2023 Jul 15. URL: <u>https://www.nytimes.com/2023/07/15/technology/artificial-intelligence-models-chat-data.html</u> [accessed 2024-01-07]
- 51. Igoe K, Harvard T. Algorithmic bias in health care—how to prevent it. Chan School of Public Health. 2012. URL: <u>https://www.hsph.harvard.edu/ecpe/how-to-prevent-algorithmic-bias-in-health-care/</u> [accessed 2024-01-07]
- 52. AI risk management framework. National Institute of Standards and Technology. URL: <u>https://www.nist.gov/itl/</u> <u>ai-risk-management-framework</u> [accessed 2024-01-07]
- 53. Blouin L. AI's mysterious 'black box' problem, explained. Dearborn. 2023. URL: <u>https://umdearborn.edu/news/</u> <u>ais-mysterious-black-box-problem-explained</u> [accessed 2024-01-07]
- 54. Wadden JJ. Defining the undefinable: the black box problem in healthcare artificial intelligence. J Med Ethics 2021 Jul 21:medethics-2021-107529. [doi: 10.1136/medethics-2021-107529] [Medline: 34290113]

Abbreviations

AI: artificial intelligence
DHT: digital health technology
FDA: Food and Drug Administration
FSIS: Food Safety and Inspection Service
FTC: Federal Trade Commission
IVD: in vitro diagnostic
OTC: over-the-counter
SaMD: software as a medical device


Edited by B Malin, KE Emam; submitted 16.02.24; peer-reviewed by G Berntsen, M Lotfinia, R Jung, U Lokala; revised version received 10.01.25; accepted 23.02.25; published 07.04.25. <u>Please cite as:</u> Perakslis E, Nolen K, Fricklas E, Tubb T Striking a Balance: Innovation, Equity, and Consistency in AI Health Technologies JMIR AI 2025;4:e57421 URL: https://ai.jmir.org/2025/1/e57421 doi:10.2196/57421

© Eric Perakslis, Kimberly Nolen, Ethan Fricklas, Tracy Tubb. Originally published in JMIR AI (https://ai.jmir.org), 7.4.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Research Letter

Evaluating ChatGPT's Efficacy in Pediatric Pneumonia Detection From Chest X-Rays: Comparative Analysis of Specialized AI Models

Nitin Chetla¹, BS; Mihir Tandon², BA; Joseph Chang³, BS; Kunal Sukhija⁴, MD; Romil Patel¹, BS; Ramon Sanchez⁵, MD

¹Department of Radiology, University of Virginia School of Medicine, Charlottesville, VA, United States

²Department of Orthopaedics, Albany Medical College, Albany, NY, United States

³Department of Radiology, University of Passau, Passau, Germany

⁴Department of Emergency Medicine, Kaweah Health Medical Center, Visalia, CA, United States

⁵Department of Radiology, Children's National Hospital, Washington, DC, United States

Corresponding Author:

Mihir Tandon, BA Department of Orthopaedics Albany Medical College 43 New Scotland Ave Albany, NY, 12208 United States Phone: 1 3322488708 Email: tandonm@amc.edu

(JMIR AI 2025;4:e67621) doi:10.2196/67621

KEYWORDS

artificial intelligence; ChatGPT; pneumonia; chest x-ray; pediatric; radiology; large language models; machine learning; pneumonia detection; diagnosis; pediatric pneumonia

Introduction

Recent studies have demonstrated the versatility of ChatGPT in health care [1]. In contrast, convolutional neural networks (CNNs) have an established history in medical imaging, particularly in identifying pneumonia from chest x-rays. CNNs are a class of deep learning algorithms that recognize patterns in images, making them invaluable tools in radiology and other imaging-based diagnostics [2]. Numerous studies demonstrate CNNs' effectiveness in medical imaging [3].

With advancements and developments in artificial intelligence (AI) technology, this research aims to evaluate the effectiveness of using ChatGPT-4 to detect pneumonia on x-ray images and compare its performance with specialized CNNs. These technologies could address radiologist shortages.

Community-acquired pneumonia incidence has reached 450 million cases worldwide annually [4]. In diagnosing pneumonia, a clinical history, physical examination, and laboratory tests are required, but clinical guidelines consider chest x-ray as the gold standard for distinguishing pneumonia from other respiratory tract infections [5]. However, interobserver agreement has been poor in chest radiographs of pediatric pneumonia [6].

RenderX

Technological improvements such as ChatGPT and AI can help detect and diagnose pediatric pneumonia.

Methods

This study used a dataset of chest x-rays from the Kaggle dataset "Chest X-Ray Images (Pneumonia)," originally sourced from the Guangzhou Women and Children's Medical Center [3,7]. The dataset consists of 5863 pneumonia and normal chest x-ray images. The images were selected from retrospective cohorts of pediatric patients, aged 1-5 years, who underwent anterior-posterior chest x-rays as part of their workup. For quality assurance, the diagnoses associated with the images were graded by three expert physicians. The dataset includes bacterial and viral pneumonia cases but does not specify the type of pneumonia or distinguish between simple and complicated pneumonia.

The study used a subset of this dataset, consisting of 500 x-rays with pneumonia and 500 without pneumonia. Each image is stored in a subfolder labeled "Pneumonia" or "Normal," enabling straightforward categorization and access. ChatGPT-4 was then prompted with "Based on the image, does the patient have A) pneumonia or B) no pneumonia? Only output the answer as A or B." The results were analyzed.

Results

(Table 1 and Figure 1). The substantial bias affects the statistical measures used. ChatGPT-40 performs slightly better overall, except in sensitivity and specificity.

ChatGPT-4 Turbo was biased toward the answer nonpneumonia

Figure 1. Confusion matrix of ChatGPT-4 Turbo.



Table 1. Statistical overview table of results of ChatGPT-4 Turbo and GPT-4o.

Statistic	ChatGPT-4 Turbo	ChatGPT-40
Accuracy (95% CI)	0.541 (0.511-0.571)	0.612 (0.582-0.642)
Precision (95% CI)	0.579 (0.548-0.607)	0.576 (0.545-0.607)
Specificity (95% CI)	0.780 (0.754-0.806)	0.839 (0.816-0.861)
Sensitivity (95% CI)	0.302 (0.274-0.333)	0.850 (0.828-0.872)
F_1 -score (95% CI)	0.397 (0.367-0.427)	0.685 (0.656-0.714)

Discussion

Although ChatGPT-4 Turbo demonstrated a slight ability to differentiate between pneumonia and nonpneumonia cases, this accuracy was overshadowed by the model's strong bias, making its distinction between the two classes unreliable for clinical use. ChatGPT-40 is equally unreliable for clinical use.

Compared with Kermany et al [3], our ChatGPT results are subpar. ChatGPT's best accuracy was 61.2% (ChatGPT-4o) in this study, compared to 92.8%. ChatGPT-4o's sensitivity and specificity were also lower in this study: 85% and 38% compared to 93.2% and 90.1%, respectively. Noticeably, ChatGPT-4o's specificity was very low comparatively. ChatGPT-4 Turbo's sensitivity and specificity results were nearly reversed compared to its successor, indicating a substantial shift in predictive behavior. Our experiment only involved 1000 testing samples in total, while Kermany et al [3] trained with 5232 samples and tested another 624 samples.

Several challenges exist in using ChatGPT-4 Turbo for diagnosing pneumonia from chest x-ray radiographs. The model's strong bias toward classifying images as nonpneumonia significantly affected the accuracy and other measures used to evaluate the model's performance. The high number of false negatives could lead to delayed or missed diagnoses in a clinical setting.

A limitation of this study is that the lack of complex pattern recognition of pediatric pneumonia by ChatGPT may be anticipated as the program has likely not been fine-tuned to assess these types of patterns. However, numerous studies have mentioned that programs like ChatGPT may replace radiologists,

but studies are needed to improve these programs, and radiologists will continue to be vital to health care [8]. By providing empirical evidence of the limitations of generalist AI models, this study underscores the need for task-specific fine-tuning and integration with computer vision models, which can help further develop these programs.

ChatGPT-4 has limitations when diagnosing pneumonia from chest x-ray radiographs as shown by this research. The model's

Conflicts of Interest

None declared.

References

- 1. Tan S, Xin X, Wu D. ChatGPT in medicine: prospects and challenges: a review article. Int J Surg 2024 Jun 01;110(6):3701-3706 [FREE Full text] [doi: 10.1097/JS9.00000000001312] [Medline: 38502861]
- Li M, Jiang Y, Zhang Y, Zhu H. Medical image analysis using deep learning algorithms. Front Public Health 2023;11:1273253 [FREE Full text] [doi: 10.3389/fpubh.2023.1273253] [Medline: 38026291]
- Kermany DS, Goldbaum M, Cai W, Valentim CC, Liang H, Baxter SL, et al. Identifying medical diagnoses and treatable diseases by image-based deep learning. Cell 2018 Feb 22;172(5):1122-1131.e9 [FREE Full text] [doi: 10.1016/j.cell.2018.02.010] [Medline: 29474911]
- 4. Sattar S, Nguyen A, Sharma S. Bacterial pneumonia. In: StatPearls. Treasure Island, FL: StatPearls Publishing; 2024.
- 5. Htun TP, Sun Y, Chua HL, Pang J. Clinical features for diagnosis of pneumonia among adults in primary care setting: a systematic and meta-review. Sci Rep 2019 May 20;9(1):7600. [doi: 10.1038/s41598-019-44145-y] [Medline: 31110214]
- Voigt GM, Thiele D, Wetzke M, Weidemann J, Parpatt P, Welte T, et al. Interobserver agreement in interpretation of chest radiographs for pediatric community acquired pneumonia: findings of the pedCAPNETZ-cohort. Pediatr Pulmonol 2021 Aug;56(8):2676-2685 [FREE Full text] [doi: 10.1002/ppul.25528] [Medline: 34076967]
- 7. Mooney P. Chest x-ray images (pneumonia). Kaggle. URL: <u>https://www.kaggle.com/datasets/paultimothymooney/</u> <u>chest-xray-pneumonia</u> [accessed 2024-12-18]
- Lecler A, Duron L, Soyer P. Revolutionizing radiology with GPT-based models: current applications, future possibilities and limitations of ChatGPT. Diagn Interv Imaging 2023 Jun;104(6):269-274 [FREE Full text] [doi: 10.1016/j.diii.2023.02.003] [Medline: 36858933]

Abbreviations

AI: artificial intelligence CNN: convolutional neural network

Edited by Y Huo; submitted 16.10.24; peer-reviewed by CH Chan; comments to author 23.11.24; revised version received 24.11.24; accepted 04.12.24; published 10.01.25.

<u>Please cite as:</u>

Chetla N, Tandon M, Chang J, Sukhija K, Patel R, Sanchez R Evaluating ChatGPT's Efficacy in Pediatric Pneumonia Detection From Chest X-Rays: Comparative Analysis of Specialized AI Models JMIR AI 2025;4:e67621 URL: https://ai.jmir.org/2025/1/e67621 doi:10.2196/67621 PMID:

©Nitin Chetla, Mihir Tandon, Joseph Chang, Kunal Sukhija, Romil Patel, Ramon Sanchez. Originally published in JMIR AI (https://ai.jmir.org), 10.01.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.

Corrigenda and Addenda

Correction: "Prompt Engineering an Informational Chatbot for Education on Mental Health Using a Multiagent Approach for Enhanced Compliance With Prompt Instructions: Algorithm Development and Validation"

Per Niklas Waaler¹, MS; Musarrat Hussain¹, PhD; Igor Molchanov¹, MS; Lars Ailo Bongo¹, PhD; Brita Elvevåg², PhD

¹Department of Computer Science, UiT The Arctic University of Norway, Tromsø, Norway ²Department of Clinical Medicine, UiT The Arctic University of Norway, Tromsø, Norway

Corresponding Author:

Per Niklas Waaler, MS Department of Computer Science UiT The Arctic University of Norway Hansine Hansens vei 54 Tromsø, N-9037 Norway Phone: 47 776 44056 Email: <u>pwa011@uit.no</u>

Related Article:

Correction of: https://ai.jmir.org/2025/1/e69820

(JMIR AI 2025;4:e75191) doi:10.2196/75191

In "Prompt Engineering an Informational Chatbot for Education on Mental Health Using a Multiagent Approach for Enhanced Compliance With Prompt Instructions: Algorithm Development and Validation" (JMIR Res Protoc 2025;4:e69820) the authors noted two errors.

The affiliation of Per Niklas Waaler was changed from:

Department of Computer Science, UiT The Arctic University of Norway, Lund, Sweden

to:

Department of Computer Science, UiT The Arctic University of Norway, Tromsø, Norway

The contact information for the corresponding author was also changed from:

Corresponding Author: Per Niklas Waaler, MS Department of Computer Science UiT The Arctic University of Norway Backgatan 35, Södra Sandby Lund, 24731 Sweden Phone: 46 944 44096 Email: pwa011@uit.no

to:

Corresponding Author: Per Niklas Waaler, MS Department of Computer Science UiT The Arctic University of Norway Hansine Hansens vei 54 Tromsø, N-9037 Norway Phone: 47 776 44056 Email: pwa011@uit.no

The correction will appear in the online version of the paper on the JMIR Publications website together with the publication of this correction notice. Because this was made after submission to PubMed, PubMed Central, and other full-text repositories, the corrected article has also been resubmitted to those repositories.



Submitted 29.03.25; this is a non-peer-reviewed article; accepted 01.04.25; published 10.04.25. <u>Please cite as:</u> Waaler PN, Hussain M, Molchanov I, Bongo LA, Elvevåg B Correction: "Prompt Engineering an Informational Chatbot for Education on Mental Health Using a Multiagent Approach for Enhanced Compliance With Prompt Instructions: Algorithm Development and Validation" JMIR AI 2025;4:e75191 URL: https://ai.jmir.org/2025/1/e75191 doi:10.2196/75191 PMID:

©Per Niklas Waaler, Musarrat Hussain, Igor Molchanov, Lars Ailo Bongo, Brita Elvevåg. Originally published in JMIR AI (https://ai.jmir.org), 10.04.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Estimation of Static Lung Volumes and Capacities From Spirometry Using Machine Learning: Algorithm Development and Validation

Scott A Helgeson¹, MS, MD; Zachary S Quicksall², MS; Patrick W Johnson², MS; Kaiser G Lim³, MD; Rickey E Carter², PhD; Augustine S Lee¹, MD

¹Division of Pulmonary and Critical Care Medicine, Mayo Clinic, 4500 San Pablo Road S, Jacksonville, FL, United States ²Digital Innovation Laboratory, Department of Quantitative Health Sciences, Mayo Clinic, Jacksonville, FL, United States ³Division of Pulmonary and Critical Care Medicine, Mayo Clinic, Rochester, MN, United States

Corresponding Author:

Scott A Helgeson, MS, MD Division of Pulmonary and Critical Care Medicine, Mayo Clinic, 4500 San Pablo Road S, Jacksonville, FL, United States

Abstract

Background: Spirometry can be performed in an office setting or remotely using portable spirometers. Although basic spirometry is used for diagnosis of obstructive lung disease, clinically relevant information such as restriction, hyperinflation, and air trapping require additional testing, such as body plethysmography, which is not as readily available. We hypothesize that spirometry data contains information that can allow estimation of static lung volumes in certain circumstances by leveraging machine learning techniques.

Objective: The aim of the study was to develop artificial intelligence-based algorithms for estimating lung volumes and capacities using spirometry measures.

Methods: This study obtained spirometry and lung volume measurements from the Mayo Clinic pulmonary function test database for patient visits between February 19, 2001, and December 16, 2022. Preprocessing was performed, and various machine learning algorithms were applied, including a generalized linear model with regularization, random forests, extremely randomized trees, gradient-boosted trees, and XGBoost for both classification and regression cohorts.

Results: A total of 121,498 pulmonary function tests were used in this study, with 85,017 allotted for exploratory data analysis and model development (ie, training dataset) and 36,481 tests reserved for model evaluation (ie, testing dataset). The median age of the cohort was 64.7 years (IQR 18 - 119.6), with a balanced distribution between genders, consisting 48.2% (n=58,607) female and 51.8% (n=62,889) male patients. The classification models showed a robust performance overall, with relatively low root mean square error and mean absolute error values observed across all predicted lung volumes. Across all lung volume categories, the models demonstrated strong discriminatory capacity, as indicated by the high area under the receiver operating characteristic curve values ranging from 0.85 to 0.99 in the training set and 0.81 to 0.98 in the testing set.

Conclusions: Overall, the models demonstrate robust performance across lung volume measurements, underscoring their potential utility in clinical practice for accurate diagnosis and prognosis of respiratory conditions, particularly in settings where access to body plethysmography or other lung volume measurement modalities is limited.

(JMIR AI 2025;4:e65456) doi:10.2196/65456

KEYWORDS

artificial intelligence; machine learning; pulmonary function test; spirometry; total lung capacity; AI; ML; lung; lung volume; lung capacity; spirometer; lung disease; database; respiratory; pulmonary

Introduction

Pulmonary function testing (PFT) provides physiological measurements of the respiratory system across multiple dimensions, typically classified into (1) spirometry, which measures air flow, lung volumes, and capacities during a expiratory forced vital capacity (FVC) maneuver; (2) static lung volumes; and (3) gas exchange parameters such as the diffusing

https://ai.jmir.org/2025/1/e65456

capacity for carbon monoxide and oxygen saturations [1]. PFTs are critical for the diagnosis and prognostication of respiratory disorders, and provide a noninvasive method for measuring and monitoring the degree of respiratory impairment [2]. They are recommended for the initial evaluation of patients with chronic dyspnea and other respiratory symptoms, as well as for individuals at risk of respiratory complications due to transplant or surgery [3,4].

Basic spirometry remains the most widely used component of PFT, largely due to its size and portability, allowing it to be performed in clinic office settings or remotely at home with adequate training. However, spirometry, by definition is an expiratory FVC maneuver that focuses on assessing airflow limitations and does not directly measure static lung volumes, which can be integral to understanding many respiratory conditions [4]. Accurate determination of static lung volumes traditionally necessitates more complex and resource-intensive techniques such as body plethysmography or gas dilution methods, with body plethysmography serving as the current gold standard [3,5,6]. However, these methods, while precise, may not always be readily accessible, cost-effective, or suitable for routine clinical practice outside a specialized pulmonary function laboratory.

Advancements in artificial intelligence (AI) techniques have introduced new avenues in health care, offering the potential to derive comprehensive insights from existing data, including patterns not easily recognizable through human interpretation or standard statistical modeling. A prior study by Beverin et al [7] examined the prediction of total lung capacity from spirometry using three tree-based machine learning (ML) models, achieving a mean squared error of 560.1 mL. They further developed models to classify restrictive ventilatory impairment, achieving a sensitivity and specificity of 83% and 92%, respectively. However, they did not explore prediction of the complete lung volume assessments. Predicting functional residual capacity status, for example, could facilitate the prevention of atelectasis during anesthesia [8]. Another study by Evankovich et al [9] developed a regression model in patients with chronic obstructive pulmonary disease (COPD) to predict residual volume and its elevation status, achieving an area under the receiver operating characteristic curve (ROC) of 0.95 for predicting residual volume above 175%. However, these models lack applicability beyond the COPD cohort [9]. Given this context, we hypothesized that ML models could predict static lung volumes using spirometry alone across a diverse cohort of lung conditions. Such an approach could reduce the need for identifying those who would benefit most from formal lung volume assessments. In this study, we applied ML approaches to develop and validate an algorithm for estimating lung volumes and capacities from standard spirometry. We further examined the model performance among subsets of physiologic derangements such as obstructive and restrictive ventilatory disorders.

Methods

Cohort Selection

This study was approved by the Institutional Review Board (20 - 009821) with a waiver of consent. The dataset curated for this study was obtained from the Mayo Clinic PFT database, which houses PFT data from two distinct US regions (Midwest and Southeast), with records from February 19, 2001, to December 16, 2022. The PFTs performed on the same day—with paired spirometry and lung volume data, without the use of methacholine or a bronchodilator—were identified. Individuals under 18 years of age and patients who opted out

```
https://ai.jmir.org/2025/1/e65456
```

of authorizing their data for research use were excluded from the analysis. All lung volume measurements were performed using body plethysmography. For models trained to classify normal versus abnormal lung volume measures, an additional requirement was applied to ensure nonmissing demographics within the boundaries of the Global Lung Initiative GLI2021 lung volume estimation equations [10]. If an individual underwent multiple PFTs, only their most recent PFT measurement comprising both lung volumes and spirometry was used. The following lung volume measures were selected for prediction: expiratory reserve volume (ERV), functional residual capacity (FRC), residual volume (RV), total lung capacity (TLC), the ratio of RV to TLC as a percentage (RV/TLC), and vital capacity (VC).

Preprocessing

Following the initial database query, the dataset was augmented with reference lung function measures for both spirometry and lung volume measures, including the lower limit of normal function (LLN), the upper limit of normal function (ULN), and the expected volume. These values were generated using a custom package built according to the Global Lung Initiative pulmonary function testing reference equation publications [1,11,12]. The LLN and ULN values were used to assign "normal" (within the LLN/ULN range) or "abnormal" (below LLN or above ULN) status to reformulate the lung volume regression problem into a classification task.

Both the regression and classification data sets were split into independent training and testing subsets using a randomized 70/30 split before any downstream exploratory analysis or model development. Features provided to the models included forced expiratory volume in the first second of exhalation (FEV1), forced vital capacity (FVC), the ratio of FEV1 and FVC (FEV1/FVC), peak expiratory flow, estimated maximum vital capacity, age, gender, height, weight, and race (White, African American, Northeast Asian, Southeast Asian, and Other).

Model Selection and Evaluation

A randomized grid search was performed using various ML algorithms, including a generalized linear model with regularization, distributed random forests, extremely randomized trees, gradient-boosted trees, and XGBoost. Models were tuned using appropriate parameter grids via five-fold cross-validation on the training dataset to provide estimates of performance summarized using applicable metrics, including root mean squared error (RMSE) for regression and area under the receiver operating characteristic curve (ROC-AUC) for classification [13]. Final tuning parameters were selected from the candidate model with the highest cross-validation performance (lowest RMSE for regression, highest ROC-AUC for classification), which was ranked highest among all explored configurations. The model was then refitted to the full training data set using the chosen hyperparameters before evaluation on the testing dataset (Multimedia Appendix 1). For the classification models, the probability threshold was selected to maximize the Youden index on the training data set.

The regression model performance was evaluated visually using prediction scatter plots and summary metrics, including RMSE,

mean absolute error (MAE), mean signed difference, mean percentage error (MPE), mean absolute percentage error (MAPE), and the correlation-based coefficient of determination [14]. The classification model was evaluated with the area under the receiver-operating-characteristic curve (AUC), accuracy, sensitivity (SENS), specificity, positive predictive value, negative predictive value (NPV), precision, recall, positive likelihood ratio (LRT+), negative likelihood ratio (LRT-), odds ratio, and F1-score. All modeling was performed using the H2O AutoML cluster (version 3.44.0.3) [15]. Further details regarding the grid search process, parameter tuning, and model implementation are available in the H2O official documentation [15] (Multimedia Appendix 2).

In the cohort summary tables, categorical data were displayed as counts and percentages, while continuous data were displayed as medians and ranges. Standardized mean differences were computed to identify significant differences in variables between the training and testing datasets, with insignificant differences defined as a value <0.1. The regression and classification models were applied to the specific PFT patterns (normal, obstructed, restricted, and mixed pattern) defined by the American Thoracic Society (ATS) [10]. All analyses were performed using R software (version 4.2.2; R Foundation for Statistical Computing) on a Google Cloud Platform virtual machine.

Ethical Considerations

This study was approved by the Mayo Clinic Institutional Review board (22-009471) and was determined to be exempt (45 CFR 46.104d, Category 4). All data was deidentified for this study, and no compensation was provided to the participants

Results

A total of 121,498 PFTs were used in this study, with 85,017 allocated for exploratory data analysis and model development and 36,481 tests reserved for model evaluation. The median age across the cohort was 64.7 years (IQR 18 - 119.6), with a nearly balanced gender distribution between genders, with 48.2% (n=58,607) female patients and 51.8% (n=62,889) male patients. The cohort was predominantly White (n= 114,388, 94.1%), followed by African American patients (n=4,656, 3.8%). Of particular importance, the distribution of baseline PFT measures-both spirometry and lung volumes-showed no differences between the training and testing datasets. Standardized mean differences, indicating the degree of difference between the training and testing sets, were minimal across all variables, suggesting a well-balanced model development and testing cohorts. A complete breakdown is provided in Table 1.

Table . Cohort summary.

•				
Variables	Training dataset (n=85,015)	Testing dataset (n=36,481)	Total (N=121,496)	Standardized difference
Age (years), median (IQR)	64.7 (18.0-119.6)	64.7 (18.0-101.0)	64.7 (18.0-119.6)	.005
Gender, n (%)				.004
Female	40,964 (48.2)	17,643 (48.4)	58,607 (48.2)	
Male	44,051 (51.8)	18,838 (51.6)	62,889 (51.8)	
Race, n (%)				.01
White	80,048 (94.2)	34,340 (94.1)	114,388 (94.1)	
African American	3223 (3.8)	1433 (3.9)	4656 (3.8)	
Southeast Asian	508 (0.6)	213 (0.6)	721 (0.6)	
Northeast Asian	64 (0.1)	27 (0.1)	91 (0.1)	
Other	1172 (1.4)	468 (1.3)	1640 (1.3)	
Height (m), median (IQR)	1.7 (0.5-2.2)	1.7 (0.2-2.0)	1.7 (0.2-2.2)	.001
Weight (kg), median (IQR)	82.8 (7.8-253.4)	82.9 (12.9-400.0)	82.8 (7.8, 400.0)	.001
ATS ^a Pattern, n (%)				.007
Normal	33,150 (41.2)	14,346 (41.6)	47,496 (41.3)	
Obstruction	16,810 (20.9)	7173 (20.8)	23,983 (20.9)	
Restriction	19,856 (24.7)	8482 (24.6)	28,338 (24.7)	
Mixed defect	10,611 (13.2)	4512 (13.1)	15,123 (13.2)	
PFT ^b measures, median (IQ	R)			
FEV1 ^c	2.0 (0.2-6.8)	2.0 (0.2-6.1)	2.0 (0.2-6.8)	.005
FVC ^d	2.9 (0.3-8.8)	2.9 (0.5-8.3)	2.9 (0.3-8.8)	.004
FEV1/FVC ^e	71.6 (16.2-100.0)	71.5 (16.2-100.0)	71.6 (16.2-100.0)	.002
PEF ^f	6.1 (0.7-18.8)	6.2 (0.6-17.5)	6.2 (0.6-18.8)	.001
VC (Spiro) ^g	2.9 (0.3-8.8)	2.9 (0.5-8.3)	2.9 (0.3-8.8)	.004
RV ^h	2.3 (0.0-11.8)	2.3 (0.1-10.4)	2.3 (0.0-11.8)	.003
TLC ⁱ	5.5 (0.9-13.9)	5.5 (1.3-13.1)	5.5 (0.9-13.9)	.004
RV/TLC ^j	43.6 (1.2-90.7)	43.6 (3.4-89.7)	43.6 (1.2-90.7)	.002
FRC ^k	3.2 (0.5-12.3)	3.2 (0.4-10.8)	3.2 (0.4-12.3)	.004
ERV ¹	0.8 (0.0-4.4)	0.8 (0.0-4.1)	0.8 (0.0-4.4)	.003
VC (Pleth) ^m	3.0 (0.3-8.8)	3.0 (0.5-8.4)	3.0 (0.3-8.8)	.003

^aATS: American Thoracic Society.
^bPulmonary function test.
^cFEV1: Forced expiratory volume in the first second.
^dFVC: Forced vital capacity.
^eFEV/FVC: Ratio of FEV1 to FVC (as a percentage).

^fPEF: Peak expiratory flow.

^gVC (Spiro): Vital capacity measured via spirometry.

^hRV: Residual volume.

ⁱTLC: Total lung capacity.

^jRV/TLC: Ratio of RV to TLC (as a percentage).

^kFRC: Functional residual capacity.

^lERV: Expiratory reserve volume.

https://ai.jmir.org/2025/1/e65456



^mVC (Pleth): Vital capacity measured via body plethysmography.

Multimedia Appendix 3 stratifies the same cohort according to the ATS classification criteria for pulmonary function patterns (ie, normal, obstructive, restrictive, and mixed pattern). This stratification highlights differences in demographics and pulmonary function measures between individuals with normal, obstructive, restrictive, or mixed patterns assigned using spirometry. Predictably, spirometry measures—including FEV1, FVC, and the FEV1/FVC ratio—significantly differed between groups (*P* values<.001), as did all phenotype-related parameters presented in the table.

Lung Volume Regression

The final models chosen for evaluation were selected based on the lowest RMSE values and varied minimally in type across the lung volumes of interest. XGBoost models were identified as the best approach for predicting all lung volumes except TLC, for which traditional gradient-boosted trees showed superior performance. Model metrics were similar between the training and testing cohorts, suggesting a reasonable trade-off between overfitting and underfitting during model training (Table 2). Findings showed a strong performance overall, with relatively low RMSE and MAE values observed across all predicted lung volumes. MPE showed a negative skew across all lung volumes. However, quantile-quantile plot analyses showed that predicted values closely followed a theoretical normal distribution, with slight underprediction and overprediction of high and low values at the extremes, respectively. Paired with mean signed differences of zero-also known as the mean bias error-these evaluations suggest no global bias in the direction of model predictions. Instead, these skewed MPE values were the result of extreme values at the tails of the distribution. A complete breakdown of model performance metrics is presented in Table 2, with complementary prediction scatter plots in Figure 1. Further subgroup analysis with different ATS patterns showed relatively similar results overall and across all categories in Multimedia Appendix 2).

Table . Regression model performance metrics.

Variables	Training d	lataset					Testing dataset						
Volume	RMSE (L) ^a	MAE ^b	MSD (L) ^c	MPE (%) ^d	MAPE (%) ^e	RSQ ^f	RMSE (L)	MAE	MSD (L)	MPE(%)	MAPE (%)	RSQ	
Expirato- ry Re- serve Volume (ERV)	0.31	0.24	0	-40.12	60.28	0.64	0.33	0.25	0.00	-39.10	59.95	0.61	
Function- al Residu- al Capaci- ty (FRC)	0.56	0.42	0	-2.83	12.93	0.78	0.59	0.44	0.00	-2.91	13.51	0.75	
Residual Volume (RV)	0.54	0.40	0	-4.86	17.29	0.73	0.56	0.41	0.00	-4.92	17.80	0.71	
RV / TLC	5.07	3.93	0	-1.61	9.55	0.82	5.20	4.03	0.03	-1.58	9.83	0.81	
Total Lung Ca- pacity (TLC)	0.55	0.41	0	-1.07	7.57	0.87	0.58	0.43	0.00	-1.10	7.92	0.85	
Vital Ca- pacity (VC)	0.15	0.11	0	-0.27	3.73	0.98	0.15	0.11	0.00	-0.33	3.91	0.98	

^aRoot mean squared error.

^bMean absolute error.

^cMean signed deviation.

^dMean percent error.

^eMean absolute percent error.

^fR-Squared.





Figure 1. Regression scatter plots of predicted versus true lung volume measures.

Lung Volume Classification

Due to limitations in demographic information (ie, age and race) required for the calculation of LLN and ULN boundaries, a total of 114,377 PFTs from the regression cohort were successfully recharacterized for the development of classification models, with 34,314 PFTs reserved for model evaluation. A comparison of demographics, spirometry, and lung volumes between the training and testing data sets can be seen in Multimedia Appendices 5 and 6. These tables mirror the factors presented in Table 1, except for the lung volume classes (normal vs abnormal), which are unique to this subset.

Similar to the regression tasks, the final classification models selected for downstream evaluation varied minimally in type across lung volumes and were selected based on the largest ROC-AUC values. Traditional gradient-boosted trees ranked best for classifying lung volume status for FRC and vital capacity. XGBoost models ranked at the top for all other lung volume classifications. Across all lung volume categories, the models demonstrated strong discriminatory capacity, as indicated by high AUC values ranging from 0.85 to 0.99 in the training dataset and 0.81 to 0.98 in the testing dataset. High accuracy scores, ranging from 0.74 to 0.93, illustrate the ability of each model to correctly classify instances overall, with sensitivity scores ranging from 0.73 to 0.93 in the testing data set, indicating the effectiveness in identifying positive cases (ie, lung volume measurements outside the expected normal range). The high NPVs (ranging from 0.84 to 0.94) highlight each model's ability to correctly identify normal lung volumes. The greater variation in positive predictive value across the lung volume classes (ranging from 0.35 - 0.94) suggests that some models may struggle to identify positive cases correctly, relative to the larger population of normal test findings. Classification performance metrics can be found in Table 3, with complementary ROC curves in Figure 2.



Helgeson et al

Table . Classification model performance metrics.

Vol- ume	/ol- Training dataset me								Testing dataset											
	AUCa	ACC ^b	SENS	SPEC ^d	PPV ^e	NPV ^f	LRT+ ^g	LRT- ^h	OR ⁱ	F1 ^j	AUC	ACC	SENS	SPEC	PPV	NPV	LRT+	LRT-	OR	F1
Expi- rato- ry re- serve vol- ume (ERV)	0.85	0.76	0.78	0.76	0.38	0.95	3.24	0.29	11.23	0.51	0.81	0.74	0.73	0.75	0.35	0.94	2.87	0.36	7.95	0.47
Func- tion- al resid- ual ca- paci- ty (FRC)	0.88	0.80	0.79	0.80	0.58	0.92	3.99	0.26	15.16	0.67	0.84	0.78	0.75	0.78	0.55	0.90	3.48	0.32	10.90	0.63
Resid- ual vol- ume (RV)	0.90	0.82	0.80	0.83	0.60	0.93	4.70	0.24	19.89	0.69	0.87	0.80	0.76	0.81	0.56	0.91	4.01	0.30	13.40	0.65
RMIC (%)	0.91	0.82	0.82	0.83	0.78	0.86	4.77	0.22	21.60	0.80	0.90	0.81	0.80	0.82	0.78	0.84	4.43	0.24	18.52	0.79
To- tal lung ca- paci- ty (ILC)	0.93	0.85	0.84	0.85	0.73	0.92	5.71	0.19	30.77	0.78	0.89	0.82	0.79	0.83	0.69	0.89	4.70	0.25	18.86	0.74
Vital ca- paci- ty (VC)	0.99	0.95	0.95	0.94	0.95	0.94	16.59	0.05	30954	0.95	0.98	0.93	0.93	0.92	0.94	0.91	12.13	0.08	160.18	0.93

^aAUC: area under the receiver operating curve.

^bACC: accuracy.

^cSENS: sensitivity.

^dSPEC: specificity.

^ePPV: positive predictive value.

^fNPV: negative predictive value.

^gLRT+: likelihood ratio test+.

^hLRT-: likelihood ratio test-.

ⁱOR: odds ratio.

^jF1: F1-score.



Expiratory Reserve Volume (ERV)

Figure 2. Classification receiver operating characteristic (ROC) curves.





Functional Residual Capacity (FRC

1.00

When stratified by PFT patterns, unique strengths, and weaknesses were observed across subgroups (Multimedia Appendix 7). These variations can be attributed to the limitations of the training data, feature space, and models, while others were driven by the rarity of certain lung volume abnormalities in specific spirometry-defined patterns. For instance, in classifying ERV status-arguably the most challenging lung volume explored in this study-the model showed consistently high NPVs across all spirometry pattern types, highlighting general confidence in predicting normal lung volume status. However, it achieved notably better sensitivity in the "restriction" and "mixed pattern" subsets (0.91 and 0.75). Comparing these sensitivities and other metrics to those in the "normal" and "obstruction" subgroups, the model seems to struggle to detect positive cases in patients with normal or obstructive spirometry findings.

Discussion

The development of ML models to predict lung volume status (normal vs abnormal findings) from spirometry in over 110,000 patients has yielded highly encouraging results, displaying remarkable discriminatory power with high AUC values (0.81 - 0.95) across measured lung volumes. Estimates of FRC, TLC, RV, and the RV/TLC ratio status show strong sensitivity and specificity. These metrics remain largely consistent across spirometry-defined pattern subgroups, with a few exceptions that can generally be attributed to the rarity of abnormal lung volume measures in certain spirometry patterns. The ability to predict lung volume measures without having to perform extensive testing represents a promising innovation for improving the diagnosis and management of dyspnea and chronic respiratory diseases, particularly in the primary care

```
https://ai.jmir.org/2025/1/e65456
```

RenderX

setting [16]. The strong predictive performance of lung volume measurement underscores the potential of these models as a transformative tool in respiratory medicine, offering substantial clinical implications and opportunities for enhancing patient care.

The performance of the regression models showed a high correlation between the training and testing datasets, suggesting that the models were able to effectively capture the relationship between spirometry-derived features and measured lung volumes and capacities derived from body plethysmography. The effectiveness of the models was evident in their ability to closely approximate lung volumes with minimal deviation from true values on average. The RMSE and MAE values are low relative to their respective lung volume ranges. For instance, the median TLC measure in the cohort was 5.5 L, with the model attaining an MAE of 0.43 L and an MAPE of 7.92%. The ability to accurately estimate the RV/TLC ratio further highlights the potential of these models in capturing the dynamic interplay between these volumes, which is particularly relevant in differentiating between common lung conditions such as COPD, asthma, and restrictive lung diseases [17-20]. The high R-squared values observed for TLC (0.87 in the training set and 0.85 in the testing set) underscore the model's capacity to capture a significant portion of the variance in TLC measurement. Similarly, the robust estimation of RV (R-squared of 0.73 in the training set and 0.71 in the testing set) and FRC (R-squared of 0.78 in the training set and 0.75 in the testing set) further validates model reliability in estimating lung volumes crucial for the evaluation of respiratory function. The model demonstrated a high correlation for vital capacity (R^2 =0.98). However, this finding is misleading, as spirometry already provides an accurate estimate of vital capacity, making it trivial

0.00

0.00

to map to a similar value obtained via body plethysmography, assuming minimal measurement error and consistent effort on the part of the patient when executing breathing maneuvers. A significant change in TLC has been reported to be 10% over one year, whereas this model was able to predict TLC within 7.5% and 550 mL [10]. No significant changes were reported in FRC or RV over time. Considering the performance metrics as a whole, the potential of these models to augment clinical practice is encouraging, with R-squared values exceeding 0.7 for all volumes except ERV, which seems to be the most challenging volume to predict accurately. Estimation of TLC, RV, and their ratio (RV/TLC) is particularly promising, as the accurate estimation of the RV/TLC ratio facilitates the identification of air trapping and hyperinflation, which are key factors in many patients' symptomatology [3,17-20]. Moreover, the reasonable estimation of FRC suggests its potential utility as an indicator for restrictive lung disease diagnosis and treatment. This is particularly important as body plethysmography directly measures only FRC, which is then used to calculate the other variables.

Focusing on the estimation of ERV, the notably high MAPE indicates a relatively subpar overall performance. Given that ERV has the narrowest range of measured values (ie, median 0.8 L, (IQR 0-44) L and a large RMSE of 0.31 relative to the ERV range, this elevated MAPE may be partially influenced by the smaller margin for error [21]. ERV measures the volume of air that an individual can exhale after completing a normal tidal breath. Pairing this with spirometry, individuals with a higher ERV may experience more difficulty with exhalation or exhibit an obstructive pattern on spirometry with a lower FEV1 measure [22,23]. A higher ERV could be a sign of lung hyperinflation, while other factors like obesity, pregnancy, and significant ascites can decrease ERV [22,24]. Lung hyperinflation in obstructed patients, which is defined as elevated FRC, RV, RV/TLC, or occasionally ERV, is highly variable in patients and occurs inconsistently over time [23,25]. This inconsistency, combined with ERV's narrow range, makes it challenging to predict.

Highlighting a more robust model, predictions for the RV/TLC ratio are strong overall, with AUC values ranging from 0.8 to 0.86 across all patterns and 0.91 in the full cohort. Except for normal pattern PFTs, the model consistently achieved sensitivities >0.84, but it struggled to identify positive cases in normal spirometry tests. While spirometry alone does not directly measure RV or TLC, FEV1 and FVC can indirectly reflect changes in lung volumes. In obstructive lung diseases, a reduction in FEV1/FVC ratio combined with an increase in the RV/TLC ratio often indicates air trapping [22-25]. In restrictive diseases, such as pulmonary fibrosis, spirometry may show decreased FVC with a preserved or decreased RV/TLC ratio, suggesting reduced air trapping [22-25]. Given the absence of abnormal FEV1 and FVC values, normal spirometry patterns would not usually suggest the existence of an abnormal RV/TLC ratio, potentially explaining the reduced sensitivity to predicting abnormal RV/TLC in normal spirometry.

A previous study used a CatBoost model to predict the TLC from spirometry, yielding good results [7]. The study reports an MSE of 560.1 mL for TLC and a positive predictive value

```
https://ai.jmir.org/2025/1/e65456
```

for reduced TLC of 8% or 67%, depending on the model parameters. However, this study only focused on TLC and did not assess other pulmonary physiologic parameters obtained through lung volume measurements, such as FRC and RV. These parameters are necessary as they are crucial for assessing prognosis in various respiratory diseases [26-30].

Several studies have highlighted the importance of lung volume assessments for the diagnosis and prognosis of respiratory diseases [31]. In routine practice, it can aid in the early detection, diagnosis, and monitoring of respiratory conditions such as COPD, restrictive lung diseases, and neuromuscular disorders affecting respiratory function [10,32,33]. For instance, lung volume measurements (specifically, FRC and TLC) strongly correlate with mortality risk among patients with idiopathic pulmonary fibrosis [27,28,30]. This illustrates that the prediction of lung volumes from traditional spirometry holds substantial promise in clinical scenarios where lung volume measurements cannot be directly performed, such as primary care offices, or health care facilities in rural areas where the equipment for measuring lung volumes is not readily accessible. Another scenario is when a patient is not capable of physically performing lung volume measurements, which could involve physical conditions that prevent them or any number of other limitations that could potentially limit them. Additionally, it may facilitate personalized treatment plans by providing a more nuanced understanding of a patient's lung capacities, as lung volume measurements are typically performed only after a patient is determined to have an abnormal spirometry, unless in specialized centers.

Accurate assessment of lung volumes is pivotal in diagnosing and monitoring various respiratory conditions, including COPD, interstitial lung diseases, neuromuscular disorders, and restrictive lung diseases [4,32]. If lung volume measurements are not performed, vital capacity is often used as a surrogate [34,35]. However, there is a significant error in the application of this method, as a reduced vital capacity can be seen in restrictive lung disease and obstructive lung disease with increased residual volume [36]. A restrictive defect on lung volume measurements has rarely been seen occurring with normal vital capacity, and approximately 58% of the time with low vital capacity measurements [36]. Another study showed that when forced vital capacity >100% predicted in males or >85% predicted in females ruled out a restrictive pattern on lung volumes [37]. The use of direct lung volume prediction models, such as those developed in this study, have a significantly better performance than those used in these prior studies and could reduce the frequency of clinical scenarios where lung volumes are unknown.

The AI model's ability to estimate lung volumes from readily available spirometry data streamlines these diagnostic procedures. A typical spirometry test may take approximately 30 - 45 minutes, while lung volume measurements add another 15 - 30 minutes [38,39]. Replacing or complementing traditional, more resource-intensive lung volume measurement techniques with the AI model's predictions from spirometry data offers cost-effective alternatives. The physician fee for spirometry ranges from \$29.62 to \$150.68, depending upon the medications used, while measuring lung volumes adds another

XSL•FO RenderX

\$59.98 to the cost [40]. This approach optimizes healthcare resources, reduces patient burden associated with additional tests, and potentially increases the efficiency of healthcare delivery.

The accessibility of spirometry in various healthcare settings, coupled with the estimation of both lung volumes via the developed models, opens avenues for telemedicine applications. Remote monitoring and assessment of spirometry are already being performed and could be facilitated and enhanced with automated decision support systems utilizing models such as those developed in this study [41-43]. Such strategies could enable the continuous monitoring of patients with chronic respiratory conditions that affect lung volumes [41-43]. This aligns with the evolving landscape of telemedicine, emphasizing its potential in respiratory care.

Despite the remarkable performance of the predictive models, certain limitations warrant consideration. Model training and testing relied on datasets with potential biases in demographic variables, including a majority-White population (91%) of older adults (median age 64.7) years. These factors potentially limit the generalizability to diverse populations, although this model was developed with patients of all ages from two distinct regions of the United States (Midwest and Southeast). Further validation across broader demographic groups from various clinical settings is essential to establish widespread applicability and reliability.

Moreover, continuous refinement and validation of the models using larger datasets encompassing a broader spectrum of respiratory conditions and disease severities is imperative. This iterative process would enhance model performance while preventing model drift, ensuring its efficacy in diverse clinical scenarios even as standard clinical practices are updated or changed.

In conclusion, the development of AI models for predicting lung volumes from spirometry represents an advancement in pulmonary function assessment. The remarkable sensitivity and specificity offered by the classification models affect a transformative approach to complement traditional lung volume measurement techniques. While the regression models may not attain the same level of performance, the continuous nature of their estimates provides a unique addition to supplement and contextualize binary classifications, potentially elucidating new insights into the remote monitoring of pulmonary function. If integrated into clinical practice, these models hold the promise revolutionizing respiratory care, enabling of more comprehensive and accessible assessments of lung function, and ultimately improving patient outcomes. Overall, the models demonstrate robust performance across lung volume measurements, underscoring their potential utility in clinical practice for accurate diagnosis and prognosis of respiratory conditions in locations where access to body plethysmography or other lung volume measurement modalities is challenging ...

Acknowledgments

This publication was made possible through the support of the Walter and Leonare Annenberg Career Development Award in Pulmonary Medicine (2 of 2).

Conflicts of Interest

None declared.

Multimedia Appendix 1 Classification model parameters. [DOCX File, 18 KB - ai v4i1e65456 app1.docx]

Multimedia Appendix 2 Regression model parameters. [DOCX File, 18 KB - ai v4i1e65456 app2.docx]

Multimedia Appendix 3 Regression model cohort summary. [DOCX File, 21 KB - ai v4i1e65456 app3.docx]

Multimedia Appendix 4 Classification model cohort summary. [DOCX File, 19 KB - ai v4i1e65456 app4.docx]

Multimedia Appendix 5 Classification model cohort summary by American Thoracic Society patterns. [DOCX File, 22 KB - ai_v4i1e65456_app5.docx]

Multimedia Appendix 6

https://ai.jmir.org/2025/1/e65456

Regression model performance metrics. [DOCX File, 36 KB - ai v4i1e65456 app6.docx]

Multimedia Appendix 7 Classification model performance metrics. [DOCX File, 27 KB - ai_v4i1e65456_app7.docx]

References

- Hall GL, Filipow N, Ruppel G, et al. Official ERS technical standard: Global Lung Function Initiative reference values for static lung volumes in individuals of European ancestry. Eur Respir J 2021 Mar;57(3):2000289. [doi: 10.1183/13993003.00289-2020] [Medline: 33707167]
- 2. Crapo RO. Pulmonary-function testing. N Engl J Med 1994 Jul 7;331(1):25-30. [doi: <u>10.1056/NEJM199407073310107</u>] [Medline: <u>8202099</u>]
- O'Donnell DE, Milne KM, Vincent SG, Neder JA. Unraveling the causes of unexplained dyspnea: the value of exercise testing. Clin Chest Med 2019 Jun;40(2):471-499. [doi: 10.1016/j.ccm.2019.02.014] [Medline: 31078223]
- 4. Ruppel GL. What is the clinical value of lung volumes? Respir Care 2012 Jan;57(1):26-35. [doi: <u>10.4187/respcare.01374</u>] [Medline: <u>22222123</u>]
- 5. Ip A, Asamoah-Barnieh R, Bischak DP, Davidson WJ, Flemons WW, Pendharkar SR. Using operational analysis to improve access to pulmonary function testing. Can Respir J 2016;2016:5269374. [doi: 10.1155/2016/5269374] [Medline: 27445545]
- 6. Sassi-Dambron DE, Eakin EG, Ries AL, Kaplan RM. Treatment of dyspnea in COPD. A controlled clinical trial of dyspnea management strategies. Chest 1995 Mar;107(3):724-729. [doi: 10.1378/chest.107.3.724] [Medline: 7874944]
- Beverin L, Topalovic M, Halilovic A, Desbordes P, Janssens W, De Vos M. Predicting total lung capacity from spirometry: a machine learning approach. Front Med (Lausanne) 2023;10:1174631. [doi: <u>10.3389/fmed.2023.1174631</u>] [Medline: <u>37275373</u>]
- 8. Hedenstierna G, Rothen HU. Atelectasis formation during anesthesia: causes and measures to prevent it. J Clin Monit Comput 2000;16(5-6):329-335. [doi: 10.1023/a:1011491231934] [Medline: 12580216]
- Evankovich JW, Nouraie SM, Sciurba FC. A model to predict residual volume from forced spirometry measurements in chronic obstructive pulmonary disease. Chronic Obstr Pulm Dis 2023 Jan 25;10(1):55-63. [doi: <u>10.15326/jcopdf.2022.0354</u>] [Medline: <u>36563054</u>]
- 10. Stanojevic S, Kaminsky DA, Miller MR, et al. ERS/ATS technical standard on interpretive strategies for routine lung function tests. Eur Respir J 2022 Jul;60(1):2101499. [doi: <u>10.1183/13993003.01499-2021</u>] [Medline: <u>34949706</u>]
- 11. Quanjer PH, Stanojevic S, Cole TJ, et al. Multi-ethnic reference values for spirometry for the 3-95-yr age range: the global lung function 2012 equations. Eur Respir J 2012 Dec;40(6):1324-1343. [doi: 10.1183/09031936.00080312] [Medline: 22743675]
- Stanojevic S, Graham BL, Cooper BG, et al. Official ERS technical standards: Global Lung Function Initiative reference values for the carbon monoxide transfer factor for Caucasians. Eur Respir J 2017 Sep;50(3):1700010. [doi: 10.1183/13993003.00010-2017] [Medline: 28893868]
- Chen T, Guestrin C. XGBoost: A scalable tree boosting system. Presented at: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining; Aug 13-17, 2016; San Francisco California USA p. 785-794.
- 14. Kuhn MVD, Hvitfeldt E. Yardstick: tidy characterizations of model performance. R package version 1.3.1 2024. URL: https://yardstick.tidymodels.org [accessed 2025-03-12]
- 15. LeDell E, Poirier S. H2O automl: scalable automatic machine learning. 2020 Jul 18 Presented at: Proceedings of the AutoML Workshop at ICML URL: <u>https://api.semanticscholar.org/CorpusID:221338558</u> [accessed 2025-03-12]
- 16. Budhwar N, Syed Z. Chronic dyspnea: diagnosis and evaluation. Am Fam Physician 2020 May 1;101(9):542-548. [Medline: 32352727]
- Casanova C, Cote C, de Torres JP, et al. Inspiratory-to-total lung capacity ratio predicts mortality in patients with chronic obstructive pulmonary disease. Am J Respir Crit Care Med 2005 Mar 15;171(6):591-597. [doi: <u>10.1164/rccm.200407-867OC</u>] [Medline: <u>15591470</u>]
- Marin JM, Carrizo SJ, Gascon M, Sanchez A, Gallego B, Celli BR. Inspiratory capacity, dynamic hyperinflation, breathlessness, and exercise performance during the 6-minute-walk test in chronic obstructive pulmonary disease. Am J Respir Crit Care Med 2001 May;163(6):1395-1399. [doi: 10.1164/ajrccm.163.6.2003172] [Medline: 11371407]
- O'Donnell DE, Webb KA. Exertional breathlessness in patients with chronic airflow limitation. The role of lung hyperinflation. Am Rev Respir Dis 1993 Nov;148(5):1351-1357. [doi: <u>10.1164/ajrccm/148.5.1351</u>] [Medline: <u>8239175</u>]
- Shin TR, Oh YM, Park JH, et al. The prognostic value of residual volume/total lung capacity in patients with chronic obstructive pulmonary disease. J Korean Med Sci 2015 Oct;30(10):1459-1465. [doi: <u>10.3346/jkms.2015.30.10.1459</u>] [Medline: <u>26425043</u>]

- 21. Makridakis S. Accuracy measures: theoretical and practical concerns. Int J Forecast 1993 Dec;9(4):527-529. [doi: 10.1016/0169-2070(93)90079-3]
- 22. Quanjer PH, Tammeling GJ, Cotes JE, Pedersen OF, Peslin R, Yernault JC. Lung volumes and forced ventilatory flows. Eur Respir J 1993 Mar 1;6(Suppl 16):5-40. [doi: 10.1183/09041950.005s1693]
- Papandrinopoulou D, Tzouda V, Tsoukalas G. Lung compliance and chronic obstructive pulmonary disease. Pulm Med 2012;2012:542769. [doi: <u>10.1155/2012/542769</u>] [Medline: <u>23150821</u>]
- 24. O'Donnell DE, Laveneziana P. Physiology and consequences of lung hyperinflation in COPD. Eur Respir Rev 2006 Dec;15(100):61-67. [doi: 10.1183/09059180.00010002]
- 25. Leith DE, Brown R. Human lung volumes and the mechanisms that set them. Eur Respir J 1999 Feb;13(2):468-472. [doi: 10.1183/09031936.99.13246899] [Medline: 10065702]
- 26. Budweiser S, Harlacher M, Pfeifer M, Jörres RA. Co-morbidities and hyperinflation are independent risk factors of all-cause mortality in very severe COPD. COPD 2014 Aug;11(4):388-400. [doi: 10.3109/15412555.2013.836174] [Medline: 24111878]
- 27. Erbes R, Schaberg T, Loddenkemper R. Lung function tests in patients with idiopathic pulmonary fibrosis. Are they helpful for predicting outcome? Chest 1997 Jan;111(1):51-57. [doi: <u>10.1378/chest.111.1.51</u>] [Medline: <u>8995992</u>]
- Kishaba T, Maeda A, Yamazato S, Nabeya D, Yamashiro S, Nagano H. Radiological and physiological predictors of IPF mortality. Medicina (Kaunas) 2021 Oct 18;57(10):1121. [doi: <u>10.3390/medicina57101121</u>] [Medline: <u>34684158</u>]
- 29. Nishimura K, Izumi T, Tsukino M, Oga T. Dyspnea is a better predictor of 5-year survival than airway obstruction in patients with COPD. Chest 2002 May;121(5):1434-1440. [doi: 10.1378/chest.121.5.1434] [Medline: 12006425]
- 30. King TE, Tooze JA, Schwarz MI, Brown KR, Cherniack RM. Predicting survival in idiopathic pulmonary fibrosis. Am J Respir Crit Care Med 2001 Oct 1;164(7):1171-1181. [doi: <u>10.1164/ajrccm.164.7.2003140</u>]
- 31. Lutfi MF. The physiological basis and clinical significance of lung volume measurements. Multidiscip Respir Med 2017;12:3. [doi: 10.1186/s40248-017-0084-5] [Medline: 28194273]
- 32. Agustí A, Celli BR, Criner GJ, et al. Global initiative for chronic obstructive lung disease 2023 report: GOLD executive summary. Eur Respir J 2023 Apr;61(4):2300239. [doi: 10.1183/13993003.00239-2023]
- Chiang J, Mehta K, Amin R. Respiratory diagnostic tools in neuromuscular disease. Children (Basel) 2018 Jun 15;5(6):78. [doi: <u>10.3390/children5060078</u>] [Medline: <u>29914128</u>]
- 34. Mehrparvar AH, Sakhvidi MJZ, Mostaghaci M, Davari MH, Hashemi SH, Zare Z. Spirometry values for detecting a restrictive pattern in occupational health settings. Tanaffos 2014;13(2):27-34. [Medline: 25506373]
- 35. Pellegrino R, Viegi G, Brusasco V, et al. Interpretative strategies for lung function tests. Eur Respir J 2005 Nov;26(5):948-968. [doi: 10.1183/09031936.05.00035205] [Medline: 16264058]
- 36. Dykstra BJ, Scanlon PD, Kester MM, Beck KC, Enright PL. Lung volumes in 4,774 patients with obstructive lung disease. Chest 1999 Jan;115(1):68-74. [doi: 10.1378/chest.115.1.68] [Medline: 9925064]
- Vandevoorde J, Verbanck S, Schuermans D, et al. Forced vital capacity and forced expiratory volume in six seconds as predictors of reduced total lung capacity. Eur Respir J 2008 Feb;31(2):391-395. [doi: <u>10.1183/09031936.00032307</u>] [Medline: <u>17928313</u>]
- 38. What is spirometry and why it is done. American Lung Association. 2023. URL: <u>https://www.lung.org/lung-health-diseases/</u> <u>lung-procedures-and-tests/spirometry</u> [accessed 2024-07-20]
- 39. Pulmonary function tests. National Heart Lung, and Blood Institute. URL: <u>https://www.nhlbi.nih.gov/science/pulmonary-function-lab/tests</u> [accessed 2024-07-20]
- 40. Physician fee schedule. Centers for Medicare and Medicaid Services. 2024. URL: <u>https://www.cms.gov/medicare/payment/</u><u>fee-schedules/physician?redirect=/PhysicianFeeSched</u> [accessed 2024-08-05]
- 41. Burgos F, Disdier C, de Santamaria EL, et al. Telemedicine enhances quality of forced spirometry in primary care. Eur Respir J 2012 Jun;39(6):1313-1318. [doi: 10.1183/09031936.00168010] [Medline: 22075488]
- 42. Congrete S, Metersky ML. Telemedicine and remote monitoring as an adjunct to medical management of bronchiectasis. Life (Basel) 2021 Nov 6;11(11):1196. [doi: 10.3390/life11111196] [Medline: 34833072]
- 43. Liao CA, Young TH, Cheng CT, et al. The feasibility and efficiency of remote spirometry system on the pulmonary function for multiple ribs fracture patients. J Pers Med 2021 Oct 23;11(11):1067. [doi: <u>10.3390/jpm11111067</u>] [Medline: <u>34834419</u>]

Abbreviations:

AI: artificial intelligence
AUC: area under the receiver-operating-characteristic curve
COPD: chronic obstructive pulmonary disease
ERV: expiratory reserve volume
FEV1: forced expiratory volume in the first second of exhalation
FEV1/FVC: ratio of FEV1 and FVC
FRC: functional residual volume
FVC: forced vital capacity
LLN: lower limit of normal

```
https://ai.jmir.org/2025/1/e65456
```

LRT+: positive likelihood ratio LRT-: negative likelihood ratio MAE: mean absolute error MAPE: mean absolute percentage error ML: machine learning MPE: mean percentage error **NPV:** negative predictive value **PFT:** pulmonary function test **PPV:** positive predictive value RMSE: root mean squared error **RV:** residual volume RV/TLC: ratio of residual volume to total lung capacity **SPEC:** specificity TLC: total lung capacity ULN: upper limit of normal VC: vital capacity

Edited by KE Emam; submitted 01.10.24; peer-reviewed by K Singh, S Liu; revised version received 18.12.24; accepted 09.02.25; published 24.03.25.

<u>Please cite as:</u>

Helgeson SA, Quicksall ZS, Johnson PW, Lim KG, Carter RE, Lee AS Estimation of Static Lung Volumes and Capacities From Spirometry Using Machine Learning: Algorithm Development and Validation JMIR AI 2025;4:e65456 URL: <u>https://ai.jmir.org/2025/1/e65456</u> doi:<u>10.2196/65456</u>

© Scott A Helgeson, Zachary S Quicksall, Patrick W Johnson, Kaiser G Lim, Rickey E Carter, Augustine S Lee. Originally published in JMIR AI (https://ai.jmir.org), 24.3.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Limitations of Binary Classification for Long-Horizon Diagnosis Prediction and Advantages of a Discrete-Time Time-to-Event Approach: Empirical Analysis

De Rong Loh^{1,2}, BSc; Elliot D Hill², MS; Nan Liu¹, PhD; Geraldine Dawson³, PhD; Matthew M Engelhard², MD, PhD

¹Duke-NUS Medical School, 8 College Road, Singapore, Singapore

²Department of Biostatistics and Bioinformatics, Duke University School of Medicine, Durham, NC, United States

³Department of Psychiatry and Behavioral Sciences, Duke University School of Medicine, Durham, NC, United States

Corresponding Author: De Rong Loh, BSc Duke-NUS Medical School, 8 College Road, Singapore, Singapore

Abstract

Background: A major challenge in using electronic health records (EHR) is the inconsistency of patient follow-up, resulting in right-censored outcomes. This becomes particularly problematic in long-horizon event predictions, such as autism and attention-deficit/hyperactivity disorder (ADHD) diagnoses, where a significant number of patients are lost to follow-up before the outcome can be observed. Consequently, fully supervised methods such as binary classification (BC), which are trained to predict observed diagnoses, are substantially affected by the probability of sufficient follow-up, leading to biased results.

Objective: This empirical analysis aims to characterize BC's inherent limitations for long-horizon diagnosis prediction from EHR; and quantify the benefits of a specific time-to-event (TTE) approach, the discrete-time neural network (DTNN).

Methods: Records within the Duke University Health System EHR were analyzed, extracting features such as *ICD-10* (*International Classification of Diseases, Tenth Revision*) diagnosis codes, medications, laboratories, and procedures. We compared a DTNN to 3 BC approaches and a deep Cox proportional hazards model across 4 clinical conditions to examine distributional patterns across various subgroups. Time-varying area under the receiving operating characteristic curve (AUC_t) and time-varying average precision (AP_t) were our primary evaluation metrics.

Results: TTE models consistently had comparable or higher AUC_t and AP_t than BC for all conditions. At clinically relevant operating time points, the area under the receiving operating characteristic curve (AUC) values for DTNN_{YOB≤2020} (year-of-birth) and DCPH_{YOB≤2020} (deep Cox proportional hazard) were 0.70 (95% CI 0.66 - 0.77) and 0.72 (95% CI 0.66 - 0.78) at *t*=5 for autism, 0.72 (95% CI 0.65 - 0.76) and 0.68 (95% CI 0.62 - 0.74) at *t*=7 for ADHD, 0.72 (95% CI 0.70 - 0.75) and 0.71 (95% CI 0.69 - 0.74) at *t*=1 for recurrent otitis media, and 0.74 (95% CI 0.68 - 0.82) and 0.71 (95% CI 0.63 - 0.77) at *t*=1 for food allergy, compared to 0.6 (95% CI 0.55 - 0.66), 0.47 (95% CI 0.40 - 0.54), 0.73 (95% CI 0.70 - 0.75), and 0.77 (95% CI 0.71 - 0.82) for BC_{YOB≤2020}, respectively. The probabilities predicted by BC models were positively correlated with censoring times, particularly for autism and ADHD prediction. Filtering strategies based on YOB or length of follow-up only partially corrected these biases. In subgroup analyses, only DTNN predicted diagnosis probabilities that accurately reflect actual clinical prevalence and temporal trends.

Conclusions: BC models substantially underpredicted diagnosis likelihood and inappropriately assigned lower probability scores to individuals with earlier censoring. Common filtering strategies did not adequately address this limitation. TTE approaches, particularly DTNN, effectively mitigated bias from the censoring distribution, resulting in superior discrimination and calibration performance and more accurate prediction of clinical prevalence. Machine learning practitioners should recognize the limitations of BC for long-horizon diagnosis prediction and adopt TTE approaches. The DTNN in particular is well-suited to mitigate the effects of right-censoring and maximize prediction performance in this setting.

(JMIR AI 2025;4:e62985) doi:10.2196/62985

KEYWORDS

machine learning; artificial intelligence; deep learning; predictive models; practical models; early detection; electronic health records; right-censoring; survival analysis; distributional shifts



Introduction

Electronic health records (EHR) are a rich source of data that can be used to develop effective clinical prediction models to improve patient care [1]. However, a major challenge is that patients have inconsistent follow-ups, leading to right-censored outcomes, and follow-up length typically depends on observed covariates. This challenge is exacerbated in long-horizon event prediction, such as prediction of an autism and attention-deficit/hyperactivity disorder (ADHD) diagnosis early in life, because many patients are lost to follow-up before the outcome can be observed. Consequently, the probability of observing a diagnosis depends not only on the probability of diagnosis but also on the probability of sufficient follow-up (ie, the probability that diagnosis occurs before censoring). As a result, binary classification (BC) models trained to predict observed diagnoses are substantially affected by the probability of sufficient follow-up unless filtering strategies are carefully applied [2].

A common filtering strategy to mitigate this effect is to exclude all individuals with insufficient follow-up. However, this is not feasible for many long-term prediction tasks. For example, sufficient follow-up for ADHD would extend into adolescence and adulthood; therefore, this criterion would preclude the development of early ADHD prediction models. Even in cases where such a criterion is feasible, it can significantly reduce the sample size available for learning and introduce systematic biases [3], as it tends to exclude subpopulations with shorter follow-up, including disadvantaged groups.

Time-to-event (TTE; ie, survival analysis) methods are the natural alternative, as they are designed for right-censored outcomes. Various versions of classification trees and random forests [4,5], Bayesian networks [6,7], Cox proportional hazards regression [8] and neural networks [9,10] have been applied to survival data with mixed success, and have been adapted to the EHR setting [11]. Deep learning [12] models such as DeepSurv [13] or deep Cox proportional hazards (DCPHs), which follow the Cox proportional hazards framework but uses a neural network to predict the log-hazard ratio, have become popular for EHR prediction tasks. Neural network-based TTE approaches are advantageous because they can efficiently process large, unstructured, high-dimensional inputs and capture complex nonlinear relationships between features and outcomes.

However, common TTE approaches also have limitations relevant to long-horizon diagnosis prediction. Unlike in survival analysis, the event of interest never occurs in most patients, and typically we are more concerned with predicting diagnosis probability than predicting diagnosis timing. Consequently, approaches that predict the probability of diagnosis separately from its timing [14] are well-suited for long-horizon diagnosis prediction, whereas DCPH and other approaches that assume relative likelihood does not change over time are less appropriate. These considerations motivate our current work to use a discrete-time neural network (DTNN), which combines the benefits of BC and TTE approaches.

First, the DTNN offers significant flexibility. Specifically, it does not assume a particular parametric form for the event time

```
https://ai.jmir.org/2025/1/e62985
```

density, and in particular, allows the effect of covariates on risk to vary across the time horizon. Second, the DTNN predicts the probability of no-event within the time horizon, which is useful in diagnosis prediction where the event of interest may often not occur. For these reasons, we have found DTNN to be advantageous in our work.

In this paper, we examine the advantages of the DTNN approach compared to BC and DCPH across 4 long-horizon, EHR-based event prediction tasks. We hypothesize that the DTNN approach will yield higher discrimination performance and more accurate likelihood predictions compared to BC even after common filtering strategies are applied due to the inability of BC to disentangle the probability of diagnosis from that of insufficient follow-up. We further hypothesize that DTNN performance will be higher than DCPH, and DTNN predictions will better reflect real-world clinical prevalence and patterns. The code for our work is available online [15].

Methods

Ethical Considerations

All study procedures were approved by the Duke Health Institutional Review Board (Pro00111224) and comply with institutional policies and federal regulations. A waiver of participant consent was approved due to the minimal risk posed by study procedures and the infeasibility of obtaining consent in a large retrospective cohort. No compensation was provided to the participants. Identifiers were omitted during analysis, which was executed within the Duke PACE (Protected Analytics Computing Environment), a highly secure virtual network space designed for protected health information.

Cohort Identification

Analyses were based on inpatient and outpatient encounters within the Duke University Health System (DUHS), a large academic medical center based in Durham, NC. DUHS provides care to approximately 85% of children in Durham and surrounding Durham County, which has a diverse population with varying demographic and socioeconomic status [16]. Records were extracted from the current (2014 - 2023) DUHS EHR, which is based on the platform developed by Epic.

Study inclusion criteria were the following: (1) date of birth between January 1, 2014 and October 29, 2022; and (2) \geq 1 visit within the DUHS before aging 30 days. DUHS encounters between January 1, 2014 and June 2, 2023 were extracted for individuals meeting these criteria. See Figure S1 in Multimedia Appendix 1 for the distribution of year of birth for this identified cohort.

Diagnosis Identification

We focused on 4 clinical diagnoses: autism spectrum disorder (autism), ADHD, recurrent otitis media (ROM), and food allergy (FA). We used computable phenotypes previously established within DUHS [17] or formulated in consultation with clinicians. The classification criteria are provided in Tables S1 and S2 in Multimedia Appendix 1.

XSL•FO RenderX

Experimental Setup

BC models predicting observed diagnoses are significantly influenced by adequate follow-up probabilities, requiring meticulous filtering strategies. We first conducted baseline experiments to establish the performance of BC models with and without exclusion criteria based on year-of-birth (YOB) or follow-up length. Correspondingly, we have 3 models trained on different cohort subsets, which are denoted as BC_{YOB≤2020}, BC_{YOB≤2018}, and BC_{t≥5} (where t denotes follow-up length). The upper limit of the dataset for the prediction tasks was capped at 2020 due to the rarity of autism and ADHD diagnoses before the age of 2 years (Figure 1). For subset YOB ≤ 2018 , we excluded all children who were age younger than 5 years at the end of our observation window to limit effects of early censoring on model predictions. For subset $t\geq 5$, we excluded all children with <5 years of follow-up as a more aggressive measure; note that this subset overlaps the subset YOB ≤ 2018 . Next, we introduced 2 TTE models, namely DTNN_{YOB ≤ 2020} and DCPH_{YOB ≤ 2020}, and evaluated their performance against the 3 BC approaches. To summarize, we explored the effect of each setup when training the corresponding model to predict each of the 4 conditions, yielding 20 models in total.

Figure 1. Distribution of observed diagnosis ages in years (upper panel) and months (lower panel). Children with diagnoses before respective diagnosis age cutoffs (marked by the red line) were excluded. Note that there were 2 ADHD diagnoses before the age cutoff of 3 years. ADHD: attention-deficit/hyperactivity disorder; FA: food allergy; ROM: recurrent otitis media.



Our features were based on encounters taking place before the following predefined, condition-specific prediction ages: 15 months, 3 years, 4 months, and 3 months for autism, ADHD, ROM, and FA, respectively (Figure 1). These ages were chosen to be clinically useful prediction times that were earlier than most observed diagnoses. Individuals diagnosed or censored before these cutoffs were excluded from the analysis. To prevent temporal data leakage, the events used for prediction were limited to those taking place before the first diagnosis code (*ICD-10* [*International Classification of Diseases, Tenth Revision*]) associated with the outcome of interest. The distribution of censoring ages can be found in Figure S2 in Multimedia Appendix 1.

The use of predefined diagnosis age cutoffs was a deliberate design decision. First, we aimed to demonstrate the predictive value of detection models based solely on EHR data collected from early ages [17]. Second, using fixed age-offs standardizes the data collection period for all individuals, which simplifies analysis and ensures consistency across the dataset. This approach allows us to focus on understanding model



performance across various clinical conditions without the additional complexity of time-dependent updates.

For each diagnosis, the dataset was partitioned randomly, allocating 60% for training, 20% for validation, and 20% for testing.

Model Development

Overview

Each observation was represented by the triplet {X,T,S}, where X \subseteq Rd is a *d*-dimensional feature vector, T \in (0,Emax] is an observed event or censoring time over a finite time horizon, and S \in {0,1} indicates whether *T* is a right-censoring time (*S*=0) or an event time (*S*=1). The observed time *T* is the minimum of the event time *E* and the right-censoring time *C*, that is, T=min(E,C).

The model selection process began with experimenting with different combinations of fully connected layers and transformer architectures. See Figure 2 for the final model architectures.



Figure 2. Model architectures of DTNN, DCPH, and BC. BC: binary classification; DCPH: deep Cox proportional hazard; DTNN: discrete-time neural network; FC: fully connected; MLP: multilayer perceptron; ReLU: rectified linear unit.



Pretraining Medical Concept Embeddings

Patient histories were represented as timestamped sequences of DUHS EHR events, including *ICD-10* diagnosis codes, medications (RxNorm [18] codes), procedures (Current Procedural Terminology [19] codes), and laboratories (Logical Observation Identifiers Names and Codes [20] codes). Events were mapped to corresponding Word2Vec embeddings, which were learned by training the model on these event sequences to capture contextual relationships between codes. The model used a Continuous Bag of Words approach with negative sampling, producing embeddings of size 256. Padding and out-of-vocabulary indices were also included and mapped to a vector of zeroes. Table S3 in Multimedia Appendix 1 details the hyperparameters used during the training process.

Encoder Architecture

The BC and TTE models all shared a common underlying encoder architecture comprised of (1) an embedding layer, (2) a fully connected layer with rectified linear unit activation applied in parallel to each embedding, (3) a global mean pooling layer, and (4) a fully connected layer with rectified linear unit activation. The embedding layer was initialized with frozen pretrained weights from the Word2Vec model. The sequence length was fixed at 512. Shorter sequences were padded, while longer sequences were truncated by selecting the most recent events preceding the age cutoff for a given model. The mean pooling layer was applied across the sequence dimension, resulting in a single fixed-length vector with dimension equal to that of the embeddings.

Prediction Head

In DTNN, the prediction head was a single fully connected hidden layer with Softmax activation, producing a probability distribution across multiple bins. The bin boundaries can be found in Table S4 in Multimedia Appendix 1. Under the common assumption of noninformative right-censoring, we may ignore the censoring density and optimize the likelihood $P(t, s | x; \theta)$ over the observed data $D=\{xi,ti,si\}i=1N$ by minimizing the following loss:

LMLE(θ)=-(silog p θ (ti|xi)+(1-si)log P θ (ti|xi))

where P θ is the survival function associated with p θ and *T* has been discretized such that each ti indicates which interval contains min(E,C).

In BC and DCPH, the prediction head was a fully connected hidden layer predicting the log-odds and log-hazard ratio, respectively, with corresponding binary cross entropy or cox negative partial log-likelihood [21] loss. Whereas BC directly predicts the probability that diagnosis will be observed (by applying the logistic function to the predicted log-odds), with DCPH this probability may be derived from the predicted log-hazard ratio and baseline hazard function. Note that for BC, we assumed a constant predicted probability irrespective of the time point.

Hyperparameter Tuning

The hyperparameters, consisting of learning rate and weight decay, were then chosen through a grid search to minimize loss on the validation set (Table S5 in Multimedia Appendix 1). These optimized models were subsequently used for evaluation on the test set.

Model Evaluation

Calibration Curves

The BC models were evaluated using the probability calibration module from the *scikit-learn* library [22], while the TTE models were evaluated by comparing the observed probabilities (ie, estimated survival probabilities of the Kaplan-Meier estimator) and the predicted probabilities at selected time intervals [23].

Performance Metrics

Our primary evaluation metrics were the time-varying area under the receiving operating characteristic curve (AUC_t) and time-varying average precision (AP_t) [24], which quantify the model's ability to discriminate between individuals diagnosed before the age t (positives; S=1, t≤t) and individuals remaining event-free beyond age t (negatives; t>t). This time-dependent approach is necessary due to censoring, which prevents many diagnoses from being observed. In contrast, the standard area under the receiving operating characteristic curve (AUC) and average precision (AP) do not differentiate between nondiagnosed individuals with short versus long follow-up,



making them unsuitable for evaluating predicted diagnosis probabilities.

Harrell concordance index [25] was also used to quantify the agreement between likelihood predictions and event times. This metric quantifies the model's ability to discriminate between individuals diagnosed earlier and those diagnosed later or not at all.

For each metric, we computed the 95% CI of the distribution over performance obtained from 100 bootstrap samples in the test set.

As we were unable to directly assess the accuracy of the predicted probabilities because diagnoses were not fully observed in the dataset, we instead contextualized them and reasoned about their correctness by analyzing the corresponding published trends.

Subgroup Analysis

To explore possible differential effects of each model setup on specific demographics, we analyzed model predictions and performance in subgroups defined by YOB, follow-up length (ie, age at censoring), sex, race, and insurance. Biological sex was classified as male or female. Race was categorized into the following groups: Asian, Black or African American, White, unavailable, and other. Insurance status was separated into public, private, and other categories.

To assess the performance of our models on out-of-distribution (OOD) data, we extended the evaluation to include children born after 2018 and individuals with a follow-up duration of <5 years for the YOB and follow-up length plots, respectively. For the YOB plots, 2019 and 2020 were designated as OOD years for BC_{YOB<2018}. Since BC_{t≥5} also fulfilled the YOB<2018 criteria, the same years were, by extension, considered OOD. Similarly, for the follow-up length plots, individuals with a

follow-up duration of ≥ 5 years were categorized as in-distribution, while those with <5 years were classified as OOD.

Semisynthetic ROM Dataset

To further explore the effect of early censoring on each method's ability to predict diagnosis probability, we simulated early censoring for ROM cases. Unlike ADHD, most ROM diagnoses were observed rather than censored due to the earlier age of diagnosis. Leveraging prior knowledge of true ROM labels, we introduced artificial censoring by scaling the true censoring distribution such that the maximum age is at 1.2 years to mimic the ADHD scenario. Generating a semisynthetic ROM dataset served 2 purposes: reproducing earlier findings on BC limitations with censored data and demonstrating DTNN model performance under such conditions. Additional DTNN and BC models were trained on this semisynthetic train dataset and subsequently evaluated on the original test dataset.

This study follows the Consolidated Reporting of Machine Learning Studies guidelines (Checklist 1) [26].

Results

Patient Characteristics

Records for 57,701 unique patients meeting study criteria were initially extracted. After excluding children born after 2020, the evaluation dataset comprised 43,536 patients (Table 1). Based on the respective diagnosis age cutoffs (Figure 1), we further excluded 1 individual with autism as an outlier due to a diagnosis within the first month of birth, along with 2 individuals with ADHD, 25 individuals with ROM, and 70 with FA. Additionally, individuals with censoring ages preceding the age cutoffs were excluded: 9332 from the autism dataset, 17,691 from the ADHD dataset, 6171 from the ROM dataset, and 5847 from the FA dataset.



Table . Patient demographics.

Variable and category or value	All	Autism	ADHD ^a	ROM ^b	FA ^c
Total, n (%)	43,536 (100)	749 (1.7)	618 (1.4)	5201 (11.9)	916 (2.1)
Sex					
Male, n (%)	22,583 (51.9)	590 (78.8)	432 (69.9)	2951 (56.7)	544 (59.4)
Female, n (%)	20,953 (48.1)	159 (21.2)	186 (30.1)	2250 (43.3)	372 (40.6)
Chi-square (<i>df</i>)	N/A ^d	221.9 (1)	79.7 (1)	58.8 (1)	21.5 (1)
<i>P</i> value	N/A	<.001	<.001	<.001	<.001
Race, n (%)					
Asian	1835 (4.2)	23 (3.1)	8 (1.3)	145 (2.8)	63 (6.9)
Black or African American	13,132 (30.2)	272 (36.3)	206 (33.3)	1226 (23.6)	278 (30.3)
White	18,681 (42.9)	266 (35.5)	326 (52.8)	2936 (56.5)	418 (45.6)
Unavailable	3874 (8.9)	57 (7.6)	29 (4.7)	390 (7.5)	45 (4.9)
Other	6014 (13.8)	131 (17.5)	49 (7.9)	504 (9.7)	112 (12.2)
Chi-square (<i>df</i>)	N/A	22.1 (4)	55 (4)	521.9 (4)	44.7 (4)
<i>P</i> value	N/A	<.001	<.001	<.001	<.001
Insurance, n (%)					
Public	23,262 (53.4)	431 (57.5)	326 (52.8)	2011 (38.7)	319 (34.8)
Private	20,127 (46.2)	316 (42.2)	288 (46.6)	3178 (61.1)	596 (65.1)
Other	147 (0.3)	2 (0.3)	4 (0.6)	12 (0.2)	1 (0.1)
Chi-square (<i>df</i>)	N/A	4.3 (2)	0.7 (2)	571.1 (2)	141.7 (2)
<i>P</i> value	N/A	.12	.69	<.001	<.001

^aADHD: attention-deficit/hyperactivity disorder.

^bROM: recurrent otitis media.

^cFA: food allergy.

^dN/A: not applicable.

Male-to-female ratios were 3.7 for autism, 2.3 for ADHD, 1.3 for ROM, and 1.5 for FA. All diagnoses were associated with sex (P<.001) and racial status (P<.001). ROM and FA were associated with insurance status (P<.001), but autism and ADHD were not (P=.12 and P=.69, respectively). Private insurance rates were 3178/5201 (61.1%) and 596/916 (65.1%) in the ROM and FA groups, respectively, compared to 316/749 (42.2%) and 288/618 (46.6%) in the autism and ADHD groups, respectively.

The mean age at diagnosis for autism and ADHD was 3.75 years and 6.22 years, respectively, higher than that for ROM and FA, which were 1.57 years and 2.01 years, respectively (Figure 1).

Analysis of Performance Metrics

In general, the TTE models consistently matched or outperformed BC models with higher AUC_t values across all conditions (Figure 3 and Table S6 in Multimedia Appendix 1). At clinically relevant operating time points, the AUC values for $DTNN_{YOB\leq 2020}$ and $DCPH_{YOB\leq 2020}$ were 0.70 (95% CI

0.66 - 0.77) and 0.72 (95% CI 0.66 - 0.78) at t=5 for autism, 0.72 (95% CI 0.65 - 0.76) and 0.68 (95% CI 0.62 - 0.74) at t=7 for ADHD, 0.72 (95% CI 0.70 - 0.75) and 0.71 (95% CI 0.69 - 0.74) at t=1 for ROM, and 0.74 (95% CI 0.68 - 0.82) and 0.71 (95% CI 0.63 - 0.77) at t=1 for FA, compared to 0.60 (95% CI 0.55 - 0.66), 0.47 (95% CI 0.40 - 0.54), 0.73 (95% CI 0.70 - 0.75), and 0.77 (95% CI 0.71 - 0.82) for BC_{YOB≤2020}, respectively.

Conversely, the regular AUC values for $BC_{YOB \le 2020}$ were consistently higher than those for $DTNN_{YOB \le 2020}$ and $DCPH_{YOB \le 2020}$. Notably, a statistically significant difference (*P*<.05) was observed in the ADHD prediction task (BCYOB \le 2020ADHD: AUC 0.75, 95% CI 0.71 - 0.80; DTNNYOB \le 2020ADHD: AUC 0.64, 95% CI 0.59 - 0.69; DCPHYOB \le 2020ADHD: AUC 0.64, 95% CI 0.60 - 0.69). With filtering, $BC_{YOB \le 2020}$ and $BC_{t \ge 5}$ exhibited decreased regular AUC, with the latter experiencing a larger decline.



Figure 3. Comparison of AUC_t (solid lines) and regular AUC (bar graphs). ADHD: attention-deficit/hyperactivity disorder; AUC: area under the receiving operating characteristic curve; AUC_t : time-varying area under the receiving operating characteristic curve; BC: binary classification; DCPH: deep Cox proportional hazard; DTNN: discrete-time neural network; FA: food allergy; ROM: recurrent otitis media; t: t denotes follow-up length; YOB: year-of-birth.



ROM FA 0.75 $DTNN^{ROM}_{YOB \, \leq \, 2020}$ DTNN^{FA} YOB ≤ 2020 0.8 DCPH^{ROM} YOB ≤ 2020 DCPH^{FA} YOB ≤ 2020 0.70 BC_{YOB} ≤ 2020 BCFA YOB ≤ 2020 0.65 0.7 AUCt AUCt BC ROM BCFA 0.60 BCFA 0.6 0.55 0.50 0.5 2 6 8 2 8 Years (t) Т 0.6 0.6 Regular AUC Regular AUC 0.4 0.4 0.2 0.2 108 5100 108 5200 5000 DCH108 5200 52000 off silver to solve silver 5²⁾⁰¹⁻ 8^{C108}52⁰¹⁸6⁰⁰⁴3 0.0 0.0 8C108 22018 CF25 DTHN\$ 605 2020 DTHN 908 52020

The regular AP and AP_t exhibited similar trends as described above, with higher AP_t but lower regular AP for TTE models (Figure S3 and Table S7 in Multimedia Appendix 1). However, direct comparison and interpretation are difficult due to the variation in test prevalence across different datasets. The concordance index, comparing ordered predicted event probabilities with observed event times, further demonstrates that the TTE models consistently performed as well as or better than the BC models (Table S8 in Multimedia Appendix 1). In particular, DTNN_{YOB<2020} and DCPH_{YOB<2020} achieved 0.656 and 0.667 for autism, 0.682 and 0.657 for ADHD, as compared to 0.629 and 0.558 for BC_{YOB<2020}, respectively.

The predicted probabilities for all models closely align with the observed estimates for in-distribution years, demonstrating overall good calibration, while OOD curves (ie, years 2019 and 2020) for BC_{YOB≤2018} and BC_{t≥5} show poor calibration (Figure 4).



Figure 4. Calibration analysis. The predicted probabilities were compared with observed event rates across different probability bins, using Kaplan-Meier estimates for the TTE models and true binary outcomes for the BC models. OOD curves (ie, years 2019 and 2020) were also added for $BC_{YOB \le 2018}$ and $BC_{t \ge 5}$. ADHD: attention-deficit/hyperactivity disorder; BC: binary classification; DCPH: deep Cox proportional hazard; DTNN: discrete-time neural network; FA: food allergy; OOD: out-of-distribution; ROM: recurrent otitis media; t: t denotes follow-up length; TTE: time-to-event.



Semisynthetic Censoring Experiment Results

The DTNNYOB ≤ 2020 ROM, ss performance remained comparable to DTNNYOB ≤ 2020 ROM and BCYOB ≤ 2020 ROM, exhibited good calibration, AUC_t and regular AUC values. However, BCYOB ≤ 2020 ROM, ss displayed worse calibration

due to underprediction, and had lower AUC_t and regular AUC values (Figure 5). Note that comparing performances beyond 1.2 years would be unfair, as those observed times were not available for model learning during training in the semisynthetic setup.



Figure 5. Comparison of performance metrics evaluated on the original test set between BC and DTNN models trained on original and semisynthetic ROM train datasets. AUC: area under the receiving operating characteristic curve; AUC_t: time-varying area under the receiving operating characteristic curve; BC: binary classification; DTNN: discrete-time neural network; ROM: recurrent otitis media; YOB: year-of-birth.



Subgroup Analyses

Probabilities predicted by $BC_{YOB \le 2020}$ decreased over time across all conditions. This trend was less pronounced for $BC_{YOB \le 2018}$ and $BC_{t \ge 5}$ (Figure 6). In contrast, the probabilities predicted by

DTNN_{YOB} ≤ 2020 for autism and ADHD showed a consistent yearly increase. For ROM, predicted probabilities declined from 2014 to 2017, then increased from 2018 onward. For FA, predicted probabilities modestly increased from 2014 to 2015, then stabilized at approximately 3.4% - 3.5% in subsequent years. The results for DCPH_{YOB} ≤ 2020 were heterogeneous.

Figure 6. Grouped analysis of predicted probability distributions by year-of-birth. ADHD: attention-deficit/hyperactivity disorder; BC: binary classification; DCPH: deep Cox proportional hazard; DTNN: discrete-time neural network; FA: food allergy; ROM: recurrent otitis media; t: t denotes follow-up length; YOB: year-of-birth.



We expanded our YOB subgroup analysis to include 2019 and 2020 to evaluate BC model behaviours during these OOD years (Figure 6). BC_{t≥5} exhibited a modest decrease in predicted probabilities across all the conditions, more pronounced in 2020 than in 2019, while BC_{YOB<2018} remained relatively stable.

There was a positive correlation observed between the predicted probability and follow-up length in all BC models, albeit to a lesser extent in BC_{YOB≤2020} and BC_{t≥5} (Figure 7). A similar trend was apparent in the analysis of the concordance between predicted nonevent probabilities with the observed censoring times (Table 2), with BCYOB≤2020ADHD showing the highest concordance index of 0.734. BC predictions appeared to align with the test prevalence (Figures S7-S9 in Multimedia Appendix 1), whereas DTNN and DCPH predictions did not (Figures S5 nd S6 in Multimedia Appendix 1).

Figure 7. Grouped analysis of predicted probability distributions by follow-up length in years. ADHD: attention-deficit/hyperactivity disorder; BC: binary classification; DCPH: deep Cox proportional hazard; DTNN: discrete-time neural network; FA: food allergy; ROM: recurrent otitis media; t: t denotes follow-up length; YOB: year-of-birth.



Table. Concordance index by comparing ordered predicted nonevent probabilities of BC^a models with observed censoring times.

	Autism	ADHD ^b	ROM ^c	FA ^d
BC _{YOB≤2020} ^e	0.581	0.734	0.605	0.558
$BC_{YOB \le 2018}$	0.533	0.625	0.605	0.535
$BC_{t \ge 5}f$	0.5	0.605	0.576	0.491

^aBC: binary classification.

^bADHD: attention-deficit/hyperactivity disorder.

^cROM: recurrent otitis media.

^dFA: food allergy.

^eYOB: year-of-birth.

^ft denotes follow-up length.

In all 4 conditions, DTNN predicted a greater likelihood of diagnosis for males. Among the racial groups, Asians had the highest predicted probability for autism and FA, while White individuals displayed the highest predicted probability for ADHD and ROM. Regarding insurance status, individuals with private insurance were more likely to be diagnosed with ROM

and FA; however, findings for autism and ADHD were equivocal (Figure 8).

The individual results of the subgroup analysis by demographics for each model setup are available in Figures S10-S12 in Multimedia Appendix 1.





Figure 8. Demographics analysis of probability distributions by $DTNN_{YOB \le 2020}$. The subgroups are sex, race, and insurance status. ADHD: attention-deficit/hyperactivity disorder; DTNN: discrete-time neural network; FA: food allergy; ROM: recurrent otitis media; YOB: year-of-birth.



Discussion

Principal Findings

Our study contributes to the understanding of how right-censoring influences model performance and predicted probabilities over time using EHR data. We highlight inherent limitations of BC in such contexts, even with filtering strategies. Furthermore, our results reinforce the potential of TTE approaches, particularly DTNN, in mitigating bias from the censoring distribution, leading to superior discrimination, calibration, and clinical prevalence prediction.

Principal Results

First, we demonstrated that BC cannot disentangle the probability of diagnosis and early censoring, even with filtering. The BC models displayed poor AUC_t performance, despite achieving high regular AUC scores (Figure 3 and Table S6 in Multimedia Appendix 1). This discrepancy arises because AUC_t calculation excludes individuals censored before prediction time *t* whereas regular AUC calculation does not. Thus, the AUC is artificially inflated by "correctly" predicting diagnosed individuals in this subgroup of individuals who were censored early as negative cases. With filtering, BC_{YOB≤2018} and BC_{t≥5} benefitted less, resulting in lower regular AUC scores because more true cases with later diagnoses were excluded.

Spurious positive correlations between the predicted probability and follow-up length imply that BC models were unduly benefitting from early censoring (Figure 7), along with increased concordance between predicted nonevent probabilities and observed censoring times (Table 2). Similarly, these differences

RenderX

were less prominent in $BC_{YOB \le 2020}$ and even less in $BC_{t \ge 5}$, but not completely absent.

This contrast was exacerbated in long-horizon prediction tasks such as ADHD, with the degree of variation corresponding with the tail end of the diagnosis age distributions (Figure 1). ADHD showed the highest proportion of later diagnoses, followed by autism and FA, and the lowest in ROM. These results corroborate observations associating censoring with biased improved outcomes, where hazard ratios fall below 1 compared to complete follow-up and correlate inversely with the proportion of censored cases [27].

Second, we found that TTE models outperformed BC models on all datasets. In diagnoses with longer time horizons, heavy right-censoring leads to many individuals having unknown status, while shorter prediction time horizons tend to have better follow-up. DTNN_{YOB≤2020} and DCPH_{YOB≤2020} achieved comparable or higher AUC_t scores in predicting ROM and FA (Figure 3 and Table S6 in Multimedia Appendix 1), suggesting that TTE models matched or surpassed BC models on datasets with less censoring. This superiority is particularly pronounced in autism and ADHD datasets, which experience heavier censoring. The main insight is that TTE models are well-suited to predict clinical outcomes, especially those with prolonged time horizons.

In our semisynthetic ROM censoring experiment, we reproduced the limitations of BC as evidenced by the deterioration in AUC_t and regular AUC performance of BCYOB \leq 2020ROM, ss when evaluated on the original dataset (Figure 5 and Table S6 in Multimedia Appendix 1). This result supports our earlier claim that the BC models were underpredicting diagnosed individuals with early censoring. We also demonstrated that

DTNNYOB ≤ 2020 ROM, ss remained well-calibrated and maintained comparable AUC_t performance as DTNNYOB ≤ 2020 ROM (Figure 5), demonstrating the applicability of our TTE approach in situations with partially observed information.

We also examined the impact of BC filtering strategies on OOD years. Specifically, we extended the evaluation to include 2019 and 2020 (Figure 6). Notably, a discernible decline in predicted probabilities was observed for BC_{t≥5} across all clinical conditions, with a slightly more pronounced drop in 2020 compared to 2019. In contrast, predicted probabilities by BC_{YOB≤2018} remained relatively stable during the same OOD years. This suggests that the inclusion of older individuals (ie, born before 2018) with shorter follow-up (ie, <5 years) makes predictions more stable on OOD years. However, including these individuals results in declining predicted probabilities due to early censoring on in-distribution years, as we have previously demonstrated. Moreover, BC_{YOB≤2018} and BC_{t≥5} showed poor calibration for all diagnoses on OOD years (Figure 4), rendering them unsuitable for clinical deployment.

Temporal and demographics trends were poorly represented in BC and DCPH. The probability of diagnosis should remain stable or increase over time due to improved awareness and tools unless specific interventions are implemented. However, $BC_{YOB\leq2020}$ exhibited declining predicted probability for all diagnoses because the models assigned lower probability scores to individuals born later, despite the absence of temporal information during learning. Inadvertently, BC predictions follow test prevalence, which also contributes to its poor performance in the demographics subgroup analysis.

The unclear patterns in DCPH models likely result from a violation of the proportional hazards assumption, which is common in practice. For example, varying severity levels in autism and ADHD diagnoses can lead to nonproportionality, where low-likelihood groups initially exhibit delays in hazard before catching up with the high-likelihood groups [28]. By assuming constant hazard rates over time, DCPH models may not fully leverage the complexity of likelihood representations and time-dependent covariate impacts. While excelling in providing generalized representations at a population level (Figure 3 and Figure S4 in Multimedia Appendix 1), our findings suggest inconsistent or inaccurate outcomes in subgroup analyses (Figures 6 and 7, and Figures S10-S12 in Multimedia Appendix 1). DTNN, however, does not assume proportional hazards, enabling better capture of time-dependent covariate influences on survival.

In contrast to the BC and DCPH models, the diagnosis probabilities predicted by the DTNN models (Figures 6 and 8) are in keeping with actual prevalence, reflecting both temporal and demographic trends. For example, autism prevalence increased from 2.24% in 2014 to 2.79% in 2019 [29], with higher rates among males and Black individuals [30]. Our demographics analysis for ADHD also concurs with trends toward increased prevalence in males and White individuals [31]. Note that the reported prevalence in DUHS may exceed

nationwide estimates, given its status as a regional hub for neurodevelopmental diagnosis.

Interestingly, for ROM, our DTNN models appear consistent with distinctive temporal patterns including (1) declining prevalence from 2014 to 2017 associated with the availability of postpneumococcal conjugate vaccines [32] and (2) increasing prevalence from 2018 to 2020 amid the COVID-19 pandemic [33]. The DTNN models also accurately predict increased likelihood associated with male sex, White race, lower socioeconomic status [32,34], and private insurance, which reflect health care use disparities [35,36].

Our models suggest stable FA prevalence (~3.4% - 3.5%), adding to mixed data that challenge whether rates have increased (range: 4.8% - 8%) [37]. This discrepancy may arise due to difficulties in estimating true prevalence [38,39] or our stricter diagnostic criteria (*ICD-10* code+IgE-based laboratory test) compared to other studies using surrogate laboratory tests or self-report, which tend to overestimate rates of clinical disease [40-42]. Demographically, our findings corroborate higher FA prevalence among males [43] and Asian and non-Hispanic Black individuals compared to non-Hispanic White individuals [44]. Additionally, our models corroborated the lower FA prevalence reported among children with public insurance [45].

Our findings suggest that TTE models, particularly the DTNN, should be preferred in clinical settings dealing with right censored outcomes. First, the DTNN models outperformed BC models, yielding clinically meaningful discriminatory performance with AUC_t≥0.7 at early ages across all 4 clinical conditions, supporting earlier diagnoses and timely interventions. Second, the DTNN approach addresses label bias that may lead to underprediction, as evidenced by its superior discrimination, calibration and ability to reflect clinical prevalence. While the modelling approach is arguably more challenging, it avoids the need for complex and often opaque filtering procedures.

Limitations

Our study has important limitations. First, it is confined to data from DUHS only, which primarily serves a population with a high representation of Black and White individuals. This demographic makeup may limit the generalizability of the results to other health systems with different patient demographics. Second, computable phenotypes are imperfect, as the identification and timing of diagnosis can vary in practice. Third, not all information, including vital signs and laboratory values, was used during the training process. Fourth, we do not include every possible filtering strategy and competing model, which may contribute to the breadth of our findings. Fifth, sex bias may also influence diagnosis trends, with males being more likely to be diagnosed with autism in practice. To the extent that sex affects the distribution of event times, the discrete-time approach can help mitigate this bias, because it does not conflate diagnosis probability with timing unlike BC and DCPH approaches. However, to the extent that sex also influences the probability of diagnosis at any given point, this is not a bias that we can overcome by choice of model alone and will require efforts to change assessment practices. Finally, the constrained size of our dataset prevents us from conducting finer subgroup

analyses. For example, we could not explore temporal trends among different demographics, such as instances where autism rates among Black children surpassed those among White children [46]. To address these limitations, we recommend incorporating data from diverse health systems, including a broader range of clinically relevant EHR data, exploring additional filtering strategies, and expanding dataset size to enable more detailed subgroup analyses.

Conclusion

Machine learning practitioners should acknowledge the inherent limitations of BC on right-censored outcomes and consider TTE approaches, particularly DTNN, in the clinical context. Our study paves the way for future research to identify and optimize models to improve patient outcomes.

Acknowledgments

This work was supported by the National Institute of Mental Health (K01-MH127309; principal investigator ME) and Eunice Kennedy Shriver National Institute of Child Health and Human Development (NICHD P50HD093074; principal investigator GD).

Data Availability

The datasets generated or analyzed during this study are not publicly available due to privacy regulations and ethical considerations related to electronic health record and cannot be shared.

Authors' Contributions

DRL completed all analyses, drafted the initial paper, and revised this paper. ME conceptualized this study, provided feedback on analyses, and reviewed and revised this paper. EH, NL, and GD provided feedback on analyses, and reviewed and revised this paper. All authors contributed to the study design and concept. All authors approved the final paper as submitted and agree to be accountable for all aspects of the work.

Conflicts of Interest

GD is on the Scientific Advisory Board of Tris Pharma, Inc, and the Nonverbal Learning Disability Project and received book royalties from Guilford Press and Springer Nature Press.

Multimedia Appendix 1 Additional figures and tables. [DOCX File, 2326 KB - ai v4i1e62985 app1.docx]

Checklist 1

CREMLS checklist. CREMLS: Consolidated Reporting of Machine Learning Studies. [DOCX File, 22 KB - ai v4i1e62985 app2.docx]

References

- 1. Rajkomar A, Oren E, Chen K, et al. Scalable and accurate deep learning with electronic health records. NPJ Digit Med 2018;1(1):18. [doi: 10.1038/s41746-018-0029-1] [Medline: 31304302]
- 2. Stajduhar I, Dalbelo-Basić B, Bogunović N. Impact of censoring on learning Bayesian networks in survival modelling. Artif Intell Med 2009 Nov;47(3):199-217. [doi: 10.1016/j.artmed.2009.08.001] [Medline: 19833488]
- 3. Weber GM, Adams WG, Bernstam EV, et al. Biases introduced by filtering electronic health records for patients with "complete data". J Am Med Inform Assoc 2017 Nov 1;24(6):1134-1141. [doi: <u>10.1093/jamia/ocx071</u>] [Medline: <u>29016972</u>]
- 4. Ishwaran H, Kogalur UB, Blackstone EH, Lauer MS. Random survival forests. Ann Appl Stat 2008;2(3):841-860. [doi: 10.1214/08-AOAS169]
- 5. Ibrahim N, Kudus A, Daud I, Bakar M. Decision tree for competing risks survival probability in breast cancer study. Int J Biol Med Sci 2008;3:25-29. [doi: <u>10.5281/zenodo.1078975</u>]
- Bandyopadhyay S, Wolfson J, Vock DM, et al. Data mining for censored time-to-event data: a Bayesian network model for predicting cardiovascular risk from electronic health record data. Data Min Knowl Disc 2015 Jul;29(4):1033-1069. [doi: 10.1007/s10618-014-0386-6]
- Brownstein NC, Bunn V, Castro LM, Sinha D. Bayesian analysis of survival data with missing censoring indicators. Biometrics 2021 Mar;77(1):305-315. [doi: <u>10.1111/biom.13280</u>] [Medline: <u>32282929</u>]
- 8. Cox DR. Regression models and life-tables. J R Stat Soc Ser B 1972 Jan 1;34(2):187-202. [doi: 10.1111/j.2517-6161.1972.tb00899.x]

- Biganzoli E, Boracchi P, Mariani L, Marubini E. Feed forward neural networks for the analysis of censored survival data: a partial logistic regression approach. Stat Med 1998 May 30;17(10):1169-1186. [doi: 10.1002/(sici)1097-0258(19980530)17:10<1169::aid-sim796>3.0.co;2-d] [Medline: <u>9618776</u>]
- Gensheimer MF, Narasimhan B. A scalable discrete-time survival model for neural networks. PeerJ 2019;7:e6257. [doi: 10.7717/peerj.6257] [Medline: 30701130]
- Vock DM, Wolfson J, Bandyopadhyay S, et al. Adapting machine learning techniques to censored time-to-event health record data: a general-purpose approach using inverse probability of censoring weighting. J Biomed Inform 2016 Jun;61:119-131. [doi: 10.1016/j.jbi.2016.03.009] [Medline: 26992568]
- 12. Solares JRA, Raimondi FED, Zhu Y, et al. Deep learning for electronic health records: a comparative review of multiple deep neural architectures. J Biomed Inform 2020 Jan;101:103337. [doi: 10.1016/j.jbi.2019.103337] [Medline: 31916973]
- Katzman JL, Shaham U, Cloninger A, Bates J, Jiang T, Kluger Y. DeepSurv: personalized treatment recommender system using a Cox proportional hazards deep neural network. BMC Med Res Methodol 2018 Feb 26;18(1):24. [doi: 10.1186/s12874-018-0482-1] [Medline: 29482517]
- Engelhard M, Henao R. Disentangling whether from when in a neural mixture cure model for failure time data. In: Camps-Valls G, Ruiz FJR, Valera I, editors. Proceedings of The 25th Int Conf Artif Intell Stat, PMLR 2022;151:9571-9581 [FREE Full text] [Medline: <u>35937033</u>]
- 15. LongHorizonDiagnosis. GitHub. URL: <u>https://github.com/engelhard-lab/LongHorizonDiagnosis</u> [accessed 2025-03-12]
- 16. Stolte A, Merli MG, Hurst JH, Liu Y, Wood CT, Goldstein BA. Using electronic health records to understand the population of local children captured in a large health system in Durham County, NC, USA, and implications for population health research. Soc Sci Med 2022 Mar;296:114759. [doi: 10.1016/j.socscimed.2022.114759] [Medline: 35180593]
- Engelhard MM, Henao R, Berchuck SI, et al. Predictive value of early autism detection models based on electronic health record data collected before age 1 year. JAMA Netw Open 2023 Feb 1;6(2):e2254303. [doi: 10.1001/jamanetworkopen.2022.54303] [Medline: 36729455]
- RxNorm. US National Library of Medicine. URL: <u>https://www.nlm.nih.gov/research/umls/rxnorm/index.html</u> [accessed 2025-03-12]
- 19. CPT. American Medical Association. 2024 Aug 23. URL: <u>https://www.ama-assn.org/practice-management/cpt</u> [accessed 2025-03-12]
- 20. LOINC. Regenstrief institute. URL: <u>https://loinc.org/</u> [accessed 2025-03-12]
- 21. Kvamme H, Hart B, Pati S, Sellereite N. pycox. GitHub. 2022 Jan. URL: <u>https://github.com/havakv/pycox.git</u> [accessed 2025-03-12]
- 22. Pedregosa F, Varoquaux G, Gramfort A, et al. Scikit-learn: machine learning in python. J Mach Learn Res 2011;12:2825-2830 [FREE Full text]
- 23. D'Agostino RB, Nam BH. Evaluation of the performance of survival analysis models: discrimination and calibration measures. In: Handbook of Statistics: Elsevier, Vol. 2003:1-25. [doi: 10.1016/S0169-7161(03)23001-7]
- 24. Antolini L, Boracchi P, Biganzoli E. A time-dependent discrimination index for survival data. Stat Med 2005 Dec 30;24(24):3927-3944. [doi: 10.1002/sim.2427] [Medline: 16320281]
- 25. Longato E, Vettoretti M, Di Camillo B. A practical perspective on the concordance index for the evaluation and selection of prognostic time-to-event models. J Biomed Inform 2020 Aug;108:103496. [doi: 10.1016/j.jbi.2020.103496] [Medline: 32652236]
- El Emam K, Leung TI, Malin B, Klement W, Eysenbach G. Consolidated Reporting Guidelines for Prognostic and Diagnostic Machine Learning Models (CREMLS). J Med Internet Res 2024 May 2;26:e52508. [doi: <u>10.2196/52508</u>] [Medline: <u>38696776</u>]
- 27. Barrajon E, Barrajon L. Effect of right censoring bias on survival analysis. JCO 2019 May 20;37(15_suppl):e18188. [doi: 10.1200/JCO.2019.37.15_suppl.e18188]
- 28. Aalen OO, Gjessing HK. Understanding the shape of the hazard rate: a process point of view (with comments and a rejoinder by the authors). Statist Sci 2001;16(1):1-22. [doi: 10.1214/ss/998929473]
- 29. Yuan J, Li M, Lu ZK. Racial/ethnic disparities in the prevalence and trends of autism spectrum disorder in US children and adolescents. JAMA Netw Open 2021 Mar 1;4(3):e210771. [doi: 10.1001/jamanetworkopen.2021.0771] [Medline: 33666658]
- 30. Maenner MJ, Warren Z, Williams AR. Prevalence and characteristics of autism spectrum disorder among children aged 8 years autism and developmental disabilities monitoring network, 11 sites, United States, 2020. MMWR Surveill Summ 2020;72(2):1-14. [doi: 10.15585/mmwr.ss7202a1]
- Xu G, Strathearn L, Liu B, Yang B, Bao W. Twenty-year trends in diagnosed attention-deficit/hyperactivity disorder among us children and adolescents, 1997-2016. JAMA Netw Open 2018 Aug 3;1(4):e181471. [doi: <u>10.1001/jamanetworkopen.2018.1471</u>] [Medline: <u>30646132</u>]
- 32. Kaur R, Morris M, Pichichero ME. Epidemiology of acute otitis media in the postpneumococcal conjugate vaccine era. Pediatrics 2017 Sep;140(3):e20170181. [doi: 10.1542/peds.2017-0181] [Medline: 28784702]
- 33. Allen DZ, Challapalli S, McKee S, et al. Impact of COVID-19 on nationwide pediatric otolaryngology: otitis media and myringotomy tube trends. Am J Otolaryngol 2022;43(2):103369. [doi: 10.1016/j.amjoto.2021.103369] [Medline: 35033925]

- 34. Smith DF, Boss EF. Racial/ethnic and socioeconomic disparities in the prevalence and treatment of otitis media in children in the United States. Laryngoscope 2010 Nov;120(11):2306-2312. [doi: 10.1002/lary.21090] [Medline: 20939071]
- 35. Patel S, Schroeder JW. Disparities in children with otitis media: the effect of insurance status. Otolaryngol Neck Surg 2011;144(1):73-77. [doi: 10.1177/0194599810391428]
- Nieman CL, Tunkel DE, Boss EF. Do race/ethnicity or socioeconomic status affect why we place ear tubes in children? Int J Pediatr Otorhinolaryngol 2016 Sep;88:98-103. [doi: <u>10.1016/j.ijporl.2016.06.029</u>] [Medline: <u>27497394</u>]
- Dunlop JH, Keet CA. Epidemiology of food allergy. Immunol Allergy Clin North Am 2018 Feb;38(1):13-25. [doi: 10.1016/j.iac.2017.09.002] [Medline: 29132669]
- Tang MLK, Mullins RJ. Food allergy: is prevalence increasing? Intern Med J 2017 Mar;47(3):256-261. [doi: 10.1111/imj.13362] [Medline: 28260260]
- Keet CA, Savage JH, Seopaul S, Peng RD, Wood RA, Matsui EC. Temporal trends and racial/ethnic disparity in self-reported pediatric food allergy in the United States. Ann Allergy Asthma Immunol 2014 Mar;112(3):222-229. [doi: 10.1016/j.anai.2013.12.007] [Medline: 24428971]
- 40. Bock SA. Prospective appraisal of complaints of adverse reactions to foods in children during the first 3 years of life. Pediatrics 1987 May;79(5):683-688. [doi: 10.1542/peds.79.5.683] [Medline: 3575022]
- Eggesbø M, Botten G, Halvorsen R, Magnus P. The prevalence of CMA/CMPI in young children: the validity of parentally perceived reactions in a population-based study. Allergy 2001 May;56(5):393-402. [doi: 10.1034/j.1398-9995.2001.056005393.x] [Medline: 11350302]
- 42. Osborne NJ, Koplin JJ, Martin PE, et al. Prevalence of challenge-proven IgE-mediated food allergy using population-based sampling and predetermined challenge criteria in infants. J Allergy Clin Immunol 2011 Mar;127(3):668-676. [doi: 10.1016/j.jaci.2011.01.039] [Medline: 21377036]
- 43. Pali-Schöll I, Jensen-Jarolim E. Gender aspects in food allergy. Curr Opin Allergy Clin Immunol 2019 Jun;19(3):249-255. [doi: <u>10.1097/ACI.000000000000529</u>] [Medline: <u>30893085</u>]
- 44. Jiang J, Warren CM, Brewer A, Soffer G, Gupta RS. Racial, ethnic, and socioeconomic differences in food allergies in the US. JAMA Netw Open 2023 Jun 1;6(6):e2318162. [doi: 10.1001/jamanetworkopen.2023.18162] [Medline: 37314805]
- 45. Bilaver LA, Kanaley MK, Fierstein JL, Gupta RS. Prevalence and correlates of food allergy among Medicaid-enrolled United States children. Acad Pediatr 2021;21(1):84-92. [doi: <u>10.1016/j.acap.2020.03.005</u>] [Medline: <u>32200110</u>]
- 46. Nevison C, Zahorodny W. Race/ethnicity-resolved time trends in United States ASD prevalence estimates from IDEA and ADDM. J Autism Dev Disord 2019 Dec;49(12):4721-4730. [doi: <u>10.1007/s10803-019-04188-6</u>] [Medline: <u>31435818</u>]

Abbreviations

ADHD: attention-deficit/hyperactivity disorder **AP:** average precision **AP_t:** time-varying average precision AUC: area under the receiving operating characteristic curve AUC_t: time-varying area under the receiving operating characteristic curve BC: binary classification **DCPH:** deep Cox proportional hazard **DTNN:** discrete-time neural network **DUHS:** Duke University Health System **EHR:** electronic health record **FA:** food allergy ICD-10: International Classification of Diseases, Tenth Revision **OOD:** out-of-distribution **PACE:** Protected Analytics Computing Environment ROM: recurrent otitis media t: t denotes follow-up length TTE: time-to-event YOB: year-of-birth



Edited by B Malin, KE Emam; submitted 07.06.24; peer-reviewed by M Aria, S Sengupta, X Ruan; revised version received 23.02.25; accepted 23.02.25; published 27.03.25. <u>Please cite as:</u> Loh DR, Hill ED, Liu N, Dawson G, Engelhard MM Limitations of Binary Classification for Long-Horizon Diagnosis Prediction and Advantages of a Discrete-Time Time-to-Event Approach: Empirical Analysis JMIR AI 2025;4:e62985 URL: https://ai.jmir.org/2025/1/e62985 doi:10.2196/62985

© De Rong Loh, Elliot D Hill, Nan Liu, Geraldine Dawson, Matthew M Engelhard. Originally published in JMIR AI (https://ai.jmir.org), 27.3.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Using AI to Translate and Simplify Spanish Orthopedic Medical Text: Instrument Validation Study

Saman Andalib^{1*}, BS; Aidin Spina^{1*}, BS; Bryce Picton¹, BS; Sean S Solomon¹, BS; John A Scolaro², MD; Ariana M Nelson³, MD

¹UCI School of Medicine, University of California, 1001 Health Sciences Rd, Irvine, CA, United States

²Department of Orthopaedic Surgery, UC Irvine Health, Orange, United States

³Department of Anesthesiology, UC Irvine Health, Orange, United States

*these authors contributed equally

Corresponding Author:

Aidin Spina, BS UCI School of Medicine, University of California, 1001 Health Sciences Rd, Irvine, CA, United States

Abstract

Background: Language barriers contribute significantly to health care disparities in the United States, where a sizable proportion of patients are exclusively Spanish speakers. In orthopedic surgery, such barriers impact both patients' comprehension of and patients' engagement with available resources. Studies have explored the utility of large language models (LLMs) for medical translation but have yet to robustly evaluate artificial intelligence (AI)–driven translation and simplification of orthopedic materials for Spanish speakers.

Objective: This study used the bilingual evaluation understudy (BLEU) method to assess translation quality and investigated the ability of AI to simplify patient education materials (PEMs) in Spanish.

Methods: PEMs (n=78) from the American Academy of Orthopaedic Surgery were translated from English to Spanish, using 2 LLMs (GPT-4 and Google Translate). The BLEU methodology was applied to compare AI translations with professionally human-translated PEMs. The Friedman test and Dunn multiple comparisons test were used to statistically quantify differences in translation quality. A readability analysis and feature analysis were subsequently performed to evaluate text simplification success and the impact of English text features on BLEU scores. The capability of an LLM to simplify medical language written in Spanish was also assessed.

Results: As measured by BLEU scores, GPT-4 showed moderate success in translating PEMs into Spanish but was less successful than Google Translate. Simplified PEMs demonstrated improved readability when compared to original versions (P<.001) but were unable to reach the targeted grade level for simplification. The feature analysis revealed that the total number of syllables and average number of syllables per sentence had the highest impact on BLEU scores. GPT-4 was able to significantly reduce the complexity of medical text written in Spanish (P<.001).

Conclusions: Although Google Translate outperformed GPT-4 in translation accuracy, LLMs, such as GPT-4, may provide significant utility in translating medical texts into Spanish and simplifying such texts. We recommend considering a dual approach—using Google Translate for translation and GPT-4 for simplification—to improve medical information accessibility and orthopedic surgery education among Spanish-speaking patients.

(JMIR AI 2025;4:e70222) doi:<u>10.2196/70222</u>

KEYWORDS

large language models; LLM; patient education; translation; bilingual evaluation understudy; GPT-4; Google Translate

Introduction

It has been well documented that racial and ethnic minority patient groups in the United States endure substantial limitations in patient care [1]. Specifically, significant disparities in health care outcomes between White populations and Hispanic populations persist in several overarching domains of medicine, including but not limited to rates of diabetes, hypertension, and insurance status [2]. Moreover, previous research suggests that language barriers may be associated with larger lapses in perioperative process-of-care outcomes [3], and patient populations who experience language barriers also face increased predisposition to hospital readmission and emergency department visits, further highlighting their susceptibility to undesired health care outcomes [4].

In the field of orthopedic surgery, these disparities are broadly evident [5-7]. From initial access to orthopedic care to postoperative outcomes, Spanish-speaking patients contend

```
https://ai.jmir.org/2025/1/e70222
```
with significant barriers in accessing high-quality care [6,7]. Hispanic populations often have limitations in their ability to schedule appointments for orthopedic concerns and often do not pursue revision surgery in cases of nonoptimal outcomes after surgical intervention [7,8]. During orthopedic clinic visits, more than half of Spanish-speaking patients have been asked to rely on nonqualified or ad hoc interpreters rather than professional services, indicating that this patient group faces limitations in access to clear and accurate information about orthopedic procedures and services [9]. These disparities may interact and thereby have implications on patient-reported outcome measures (PROMs) for Spanish-speaking populations. Additionally, recent work has evaluated the suitableness of PROMs for Spanish-speaking populations [10]. Commonly used PROMs for Spanish-speaking patient groups were shown to be written at a reading level above the recommended complexity for patient populations in the United States. Technological advancements can provide avenues to address these concerns if they are implemented in a manner that is tailored to their intended patient populations [11,12]. Thus, given the widespread documentation of disparities in orthopedic care that Spanish-speaking patients endure, further evaluation of how emerging technologies can address these lapses is extremely important.

Artificial intelligence (AI) has provided unique solutions to problems in health care, including those related to graduate medical education and patients' comprehension of medical text [13-17]. Recent work has turned to using publicly available large language models (LLMs) to translate patient discharge summaries and frequently asked questions. The utility of these tools in translating medical text has been illustrated in qualitative textual evaluations conducted via human grading [18,19]. However, studies have yet to evaluate AI-enabled textual translation through robust quantitative analysis involving bilingual evaluation understudy (BLEU) analysis [20]. This methodology quantitatively rates machine-translated text against human translation and has been used in clinical studies [21-23]. Additionally, no study has evaluated AI-driven simplification of Spanish medical text, although AI-driven simplification is a functionality that our group previously quantitatively evaluated for English medical text [16,24,25].

The goals of this study were twofold. First, we aimed to conduct a robust quantitative evaluation of machine translations of medical text by using BLEU analysis, and second, we aimed to assess whether AI platforms can be used to simplify orthopedic medical text written in Spanish.

Methods

Study Design

A total of 78 patient education materials (PEMs) from the American Academy of Orthopaedic Surgery (AAOS) were translated from English into Spanish, using 4 different GPT-4 input prompts via the application programming interface (prompts 1 - 4; Multimedia Appendix 1) [26] and Google Translate via the googletrans package (SuHun Han). Each machine-generated translation was compared to the professionally human-translated reference from the AAOS,

```
https://ai.jmir.org/2025/1/e70222
```

using BLEU analysis via the Natural Language Toolkit (NLTK) [27]; BLEU scores range from 0 to 1, with scores of ≥ 0.5 indicating high similarity to a designated reference text. A Friedman test, followed by a Dunn multiple comparisons test, was performed for each BLEU score to quantify differences in translation quality. Unigram, bigram, trigram, and fourgram precision analyses were conducted to further assess the translation quality. A Friedman test was followed by Dunn multiple comparisons for each precision metric.

To assess the simplification of the PEMs, we compared the readability of translations generated by GPT-4's prompt 1 and that of the original AAOS Spanish versions before and after simplification. Spanish text was simplified by using a standardized prompt that was validated for medical use cases [16]. Text complexity was analyzed by counting sentences, words, and syllables with custom functions and the NLTK library [27]. Readability was evaluated by using the Fernández-Huerta readability formula (FH = $206.84 - [0.60 \times$ P] – $[1.02 \times F]$; FH: reading ease score; P: average number of syllables per 100 words; F: average number of sentences per 100 words) [28] and the INFLESZ readability formula $(INFLESZ = 206.835 - [62.3 \times S/P] - [P/F]; S: total number$ of syllables; P: total number of words; F: total number of sentences) [29]. The Wilcoxon matched-pairs signed rank test was applied to compare the original and simplified versions, and the Spearman correlation coefficient was used to measure the strength of the association between the simplification process and improved readability.

To assess the impact of original English text features on translation quality, a feature analysis was performed. Random forest regression was completed, using 4 input features (number of words, average number of words per sentence, total number of syllables, and average number of syllables per sentence) of the original English PEM, to predict 20 distinct BLEU scores. These scores encompassed 4 BLEU scoring methods for Google Translate and 4 different GPT-4 input prompts. A 5-fold cross-validation was used to minimize overfitting of the data and to ensure robust feature importance calculations. Average importance scores across all folds were calculated to assess the contribution of each feature for translation performance.

Ethical Considerations

No application was submitted for review board assessment because no human or animal participants participated directly or indirectly in this study. The University of California, Irvine Institutional Review Board does not require assessment of studies that do not directly or indirectly involve human or animal participants. This study consisted solely of a quantitative evaluation of machine translations and was hence exempt from any institutional review.

Results

BLEU Analysis

BLEU 1 scores (Figure 1A) revealed a statistically significant difference between Google Translate and each prompt (prompt 1: rank sum difference=63.00; *P*=.01; prompt 2: rank sum difference=81.00; *P*<.001; prompt 3: rank sum difference=65.00;

XSL•FO RenderX

P=.01; prompt 4: rank sum difference=71.00; P=.003). No significant differences were observed among the 4 GPT prompts (all *P* values were >.05). For BLEU 1, Google Translate had the highest rank sum (290.0), while prompt 2 had the lowest (209.0). Prompt 1 had a rank sum of 227.0, while prompts 3 and 4 had rank sums of 225.0 and 219.0, respectively.

For BLEU 2 scores (Figure 1B), a similar trend was observed, with significant differences between Google Translate and prompts 1, 2, 3, and 4. The rank sum difference was 76.00 between Google Translate and prompt 1 (P<.001), 79.00 between prompt 2 and Google Translate (P<.001), 73.00 between prompt 3 and Google Translate (P=.002), and 77.00 between prompt 4 and Google Translate (P<.001). Again, no statistically significant differences were found between the 4 GPT prompts (all P values were >.05). The rank sum for Google Translate was the highest (295.0), followed by those for prompt 3 (222.0), prompt 1 (219.0), and prompt 4 (218.0). Prompt 2 had the lowest rank sum (216.0).

For the BLEU 3 scores (Figure 1C), the Dunn test also showed significant differences between Google Translate and each prompt (prompt 1: rank sum difference=72.00; P=.003; prompt 2: rank sum difference=85.00; P<.001; prompt 3: rank sum difference=76.00; P=.001; prompt 4: rank sum difference=82.00; P<.001). No significant differences were found between the 4 GPT prompts (all P values were >.05). The rank sums were as follows: 297.0 for Google Translate, 225.0 for prompt 1, 212.0 for prompt 2, 221.0 for prompt 3, and 215.0 for prompt 4.

Finally, BLEU 4 scores (Figure 1D) followed the same pattern as the BLEU scores in all 3 prior BLEU analyses, as the Dunn test revealed significant differences between Google Translate and each prompt (prompt 1: rank sum difference=74.00; P=.002; prompt 2: rank sum difference=77.00; P<.001; prompt 3: rank sum difference=82.00; P<.001). Google Translate had the highest rank sum (295.0), followed by prompt 3 (223.0), prompt 1 (221.0), and prompt 2 (218.0). Prompt 4 had the lowest rank sum (213.0).

Figure 1. BLEU scores for Google Translate and 4 GPT-4 input prompts (prompts 1-4). Box plots display the BLEU 1 (A), BLEU 2 (B), BLEU 3 (C), and BLEU 4 (D) scores for translations generated by Google Translate and the 4 different GPT-4 input prompts. BLEU: bilingual evaluation understudy.



N-Gram Precision Analysis

The unigram precision analysis (Figure 2A) revealed significant differences between Google Translate and prompts 1, 2, 3, and 4. The rank sum difference was 71.50 between Google Translate and prompt 1 (P=.003), 64.00 between prompt 2 and Google

```
https://ai.jmir.org/2025/1/e70222
```

XSL•FO RenderX Translate (P=.01), 55.50 between prompt 3 and Google Translate

(P=.05), and 74.00 between prompt 4 and Google Translate

(P=.002). Google Translate had the highest rank sum (287.0),

followed by prompt 3 (231.5), prompt 2 (223.0), and prompt 1

(215.5). Prompt 4 had the lowest rank sum (213.0).

The bigram precision analysis (Figure 2B) also revealed significant rank sum differences between Google Translate and each prompt (prompt 1: rank sum difference=93.00; P<.001; prompt 2: rank sum difference=88.50; P<.001; prompt 3: rank sum difference=79.50; P<.001; prompt 4: rank sum difference=99.00; P<.001). Google Translate had the highest rank sum (306.0), followed by prompt 3 (226.5). Prompt 2 followed with a rank sum of 217.5, and prompts 1 and 4 had a rank sum of 213.0 and 207.0, respectively.

For the trigram precision analysis (Figure 2C), the Dunn test revealed a pattern that was slightly different from the previously established pattern, with significant differences between Google Translate and prompt 1 (rank sum difference=80.00; P<.001), between Google Translate and prompt 2 (rank sum difference=73.00; P=.002), and between Google Translate and prompt 4 (rank sum difference=74.00; P=.002). There was no significant difference in trigram precision between Google Translate and prompt 3 (P=.07). Google Translate had the

highest rank sum (290.0), followed by prompt 3 (237.0). Prompt 2 had a rank sum of 217.0, while prompt 4 had a rank sum of 216.0. The lowest rank sum for trigram precision was recorded for prompt 1 (210.0).

The fourgram precision analysis (Figure 2D) showed the same pattern of significance as that in the trigram analysis, with significant differences between Google Translate and GPT prompts 1, 2, and 4. The rank sum difference between Google Translate and prompt 1 was 71.00 (P=.003). The rank sum differences between Google Translate and prompt 2 and between Google Translate and prompt 4 were 72.00 (P=.003) and 78.00 (P<.001), respectively. Fourgram precision showed no statistically significant difference between Google Translate and prompt 3 (P=.06). Google Translate had the highest rank sum (289.0), while prompt 3 ranked second with a rank sum of 235.0. Prompt 1 had a rank sum of 218.0, and prompt 2 closely followed with a rank sum of 217.0. Prompt 4 had the lowest rank sum (211.0).



Figure 2. N-gram precision for Google Translate and 4 GPT-4 input prompts (prompts 1-4). Box plots display unigram (A), bigram (B), trigram (C), and fourgram (D) precision scores for translations generated by Google Translate and the 4 different GPT-4 input prompts.



Simplification Analysis

As measured by the Fernández-Huerta scores, the simplified prompt 1 PEM translations and simplified AAOS Spanish PEMs demonstrated significant improvements in readability when

https://ai.jmir.org/2025/1/e70222

XSL•FO RenderX

(P<.001); the median difference was 7.846, and the Spearman

(W) test for prompt 1 showed a significant difference between

the original and simplified translations, with a W value of 3059

correlation coefficient was 0.6459 (P<.001). For the AAOS Spanish version, the Wilcoxon test revealed a significant improvement after simplification, with a W value of 3055 (P<.001) and a median difference of 5.807; the Spearman correlation coefficient was 0.6731 (P<.001).

For the INFLESZ scores, similar results were observed. For prompt 1, the Wilcoxon matched-pairs signed rank test indicated

a significant difference between the original and simplified translations, with a W value of 3058 (P<.001); the median difference was 7.830, and the Spearman correlation coefficient was 0.6591 (P<.001). For the AAOS Spanish PEMs, the Wilcoxon test showed a significant improvement after simplification, with a W value of 3045 (P<.001) and a median difference of 5.887; the Spearman correlation coefficient was 0.6926 (P<.001).

Figure 3. Fernández-Huerta and INFLESZ scores for the original translations by prompt 1 and the AAOS and for their simplified versions. Box plots display the Fernández-Huerta readability scores (A and B) and INFLESZ readability scores (C and D) for the original and simplified versions of the PEMs generated by GPT-4's prompt 1 (A and C) and for the original and simplified AAOS translations (B and D). AAOS: American Academy of Orthopaedic Surgery; PEM: patient education material.



Feature Analysis

The feature importance analysis of the original English text features revealed that the total number of syllables was the most influential predictor of BLEU scores across Google Translate and GPT-4 prompts, serving as the most important feature (ie, input variable) in every iteration, with scores ranging from 0.27

https://ai.jmir.org/2025/1/e70222

XSL•FO RenderX of words was 0.2 to 0.23, that for the average number of words per sentence was 0.19 to 0.27, and that for the average number of syllables per sentence was 0.22 to 0.27. Overall, syllable-based features, particularly the total number of syllables, served as the highest-importance features in determining BLEU scores across all translation methods.

to 0.35 (Figure 4). The feature importance range for the number

Andalib et al

Figure 4. Feature importance scores of English text characteristics for predicting BLEU scores. The heat map shows the relative importance of 4 input features—number of words, average number of words per sentence, total number of syllables, and average number of syllables per sentence—in predicting BLEU scores across the 4 BLEU analyses for each of the 5 translation methods. Darker colors represent higher feature importance. avg: average; BLEU: bilingual evaluation understudy; num: number.



Discussion

Context

Disparities in communication with Spanish-speaking populations can negatively affect patient education and subsequent outcomes in the field of orthopedic surgery [5-7]. Accurate translation of medical text is one component of properly educating Spanish-speaking patient populations about orthopedic conditions. For orthopedic surgeons, it is vital to ensure that Spanish-speaking patients are properly informed about their conditions and opportunities for surgery, given their increased propensity for hospital readmission, complications, and negative outlooks on surgical intervention [6-8]. Previous work provided a foundation for quantitatively evaluating AI-based medical text translation; however, no study has used BLEU methodology to provide a robust, machine learning-based evaluation of translation success. Additionally, no study has evaluated the AI-enabled simplification of Spanish text. Given the recently outlined need for simplified Spanish text among Spanish-speaking patient populations, this is a pressing need in the field [10]. Our study used a robust corpus of patient-facing orthopedic medical text that included language from across various subspecialties and topics of orthopedic surgery, including the spine, hip, knee, and upper extremities, among others. Through analyzing the success of openly accessible LLMs in translating such text, we aimed to comprehensively assess the translation options available for orthopedic practice.

https://ai.jmir.org/2025/1/e70222

Translation Success

This study demonstrated that LLMs, such as ChatGPT, can translate orthopedic PEMs with moderate success, as quantified through BLEU analysis. By experimenting with 4 different model prompts, we explored whether prompt optimization could enhance translation effectiveness. Our findings suggest that while prompt optimization can improve translation outcomes, Google Translate generally provides superior translation quality when compared to human-translated benchmarks. This superior performance highlights the potential of Google Translate for rapid translation tasks, such as translating patient directives in discharge summaries and other patient-facing documents. However, despite its prevalent use, Google Translate's limitations underscore the need for alternative translation solutions [19,30,31]. The feature analysis conducted within our study also revealed that the syllable complexity of the original English text is a critical predictor of successful translation for both Google Translate and ChatGPT, indicating areas for further refinement in translation approaches. An example AI translation, along with the original English and Spanish versions of the same PEM, can be found in Multimedia Appendix 1.

Simplification Success

We also assessed the capability of ChatGPT in simplifying medical texts written in Spanish, using a standardized simplification prompting structure that was previously evaluated by our group. Although the platform was able to simplify the

text, it did not achieve the targeted grade level specified in our prompts. This limitation aligns with prior studies that highlighted challenges in simplifying English medical texts [16]. However, despite existing challenges with the precision of AI-simplified text in meeting prespecified grade levels, the ability of ChatGPT to simplify texts could greatly benefit Spanish-speaking patients, given that no alternative exists to aid patient comprehension in this way. This is of great importance, considering the complexity of the PROMs and other tools used to assess the operative success of orthopedic procedures in this patient group [10]. Further studies should elucidate ways to best optimize the simplification of Spanish texts via AI platforms.

Recommendations

Based on our results, we offer several recommendations for orthopedic surgeons. Although Google Translate remains a superior tool for translating English to Spanish due to its adherence to human translation quality, LLMs, such as ChatGPT, also show moderate success and can be considered for specific use cases. Importantly, ChatGPT's ability to simplify Spanish texts makes it a valuable tool for enhancing patient comprehension and engagement, particularly when translation by a native Spanish speaker is not feasible. We recommend using ChatGPT as an adjunct tool for both translating and simplifying medical texts. Surgeons should continue to use Google Translate for straightforward translations, but they should also consider leveraging ChatGPT's simplification capabilities to improve the accessibility of medical information. Further research into simplification methodologies is essential for optimizing PROMs and ultimately enhancing patient satisfaction following surgical care. We believe that this technology, once it is fully optimized and vetted, will have the potential to be incorporated into the electronic health record to aid in medical record management through textual translation of records for patients.

Limitations

This study, while providing insights into the potential of LLMs for translating and simplifying medical texts, has several limitations. First, this study assessed existing models, only tested English-to-Spanish translations, and used a relatively small amount of content, thereby limiting the generalizability of our findings. Second, the BLEU metric, which we used to evaluate translation accuracy, primarily measures literal translation and may not fully capture semantic equivalence, which is critical in medical contexts. Future research could benefit from incorporating additional evaluations that involve human assessment to provide a more nuanced analysis. Third, this study's focus was on technical performance; we did not directly measure the impact on patient outcomes, such as comprehension, adherence, and satisfaction. Future studies should aim to link the quality of translations and simplifications to specific patient-centered outcomes. Clinical studies would provide valuable insights into the way that Spanish-speaking patient populations interact with and subsequently benefit from AI-enhanced PEMs, such as those analyzed in this study. Lastly, although the corpus of 78 PEMs covered a broad scope of orthopedic literature from all subspecialties, this means that the results of this study only reflect the language used in standard orthopedic practice. Future studies should aim to replicate our results in other medical specialties to provide a broad understanding of the capabilities of AI in translation and simplification.

Conclusions

This study highlights the utility and limitations of AI-driven tools in translating and simplifying medical texts for Spanish-speaking orthopedic patients. Our findings indicate that while Google Translate provides superior accuracy in translating medical texts, LLMs, such as ChatGPT, demonstrate moderate success and offer significant benefits in simplifying complex medical information into more comprehensible formats. Our recommended dual approach-leveraging Google Translate for accuracy and ChatGPT for simplification-presents a practical solution for enhancing patient education and engagement. Such advancements underscore the potential of AI to bridge the language gap in health care and thereby improve treatment outcomes. Future research should continue to refine these AI tools and enhance their precision and accessibility to meet the diverse needs of patient populations, thereby ensuring that all patients receive care that is both understandable and culturally competent.

Conflicts of Interest

None declared.

Multimedia Appendix 1

Example artificial intelligence–translated patient education material (PEM) with original English and original Spanish PEMs. [DOCX File, 31 KB - ai v4i1e70222 app1.docx]

References

- Woloshin S, Bickell NA, Schwartz LM, Gany F, Welch HG. Language barriers in medicine in the United States. JAMA 1995 Mar 1;273(9):724-728. [Medline: <u>7853631</u>]
- Odlum M, Moise N, Kronish IM, et al. Trends in poor health indicators among Black and Hispanic middle-aged and older adults in the United States, 1999-2018. JAMA Netw Open 2020 Nov 2;3(11):e2025134. [doi: 10.1001/jamanetworkopen.2020.25134] [Medline: 33175177]

- 3. Joo H, Fernández A, Wick EC, Moreno Lepe G, Manuel SP. Association of language barriers with perioperative and surgical outcomes: a systematic review. JAMA Netw Open 2023 Jul 3;6(7):e2322743. [doi: <u>10.1001/jamanetworkopen.2023.22743</u>] [Medline: <u>37432686</u>]
- 4. Chu JN, Wong J, Bardach NS, et al. Association between language discordance and unplanned hospital readmissions or emergency department revisits: a systematic review and meta-analysis. BMJ Qual Saf 2024 Jun 19;33(7):456-469. [doi: 10.1136/bmjqs-2023-016295] [Medline: <u>38160059</u>]
- 5. Busigo Torres R, Yendluri A, Stern BZ, et al. Is limited English proficiency associated with differences in care processes and treatment outcomes in patients undergoing orthopaedic surgery? A systematic review. Clin Orthop Relat Res 2024 Aug 1;482(8):1374-1390. [doi: 10.1097/CORR.0000000003034] [Medline: 39031039]
- 6. Azua E, Fortier LM, Carroll M, et al. Spanish-speaking patients have limited access scheduling outpatient orthopaedic appointments compared with English-speaking patients across the United States. Arthrosc Sports Med Rehabil 2023 Feb 26;5(2):e465-e471. [doi: 10.1016/j.asmr.2023.01.015] [Medline: 37101862]
- Aggarwal A, Naylor JM, Adie S, Liu VK, Harris IA. Preoperative factors and patient-reported outcomes after total hip arthroplasty: multivariable prediction modeling. J Arthroplasty 2022 Apr;37(4):714-720.e4. [doi: <u>10.1016/j.arth.2021.12.036</u>] [Medline: <u>34990754</u>]
- Nguyen KH, Suarez P, Sales C, Fernandez A, Ward DT, Manuel SP. Patients who have limited English proficiency have decreased utilization of revision surgeries after hip and knee arthroplasty. J Arthroplasty 2023 Aug;38(8):1429-1433. [doi: 10.1016/j.arth.2023.02.024] [Medline: <u>36805120</u>]
- Greene NE, Fuentes-Juárez BN, Sabatini CS. Access to orthopaedic care for Spanish-speaking patients in California. J Bone Joint Surg Am 2019 Sep 18;101(18):e95. [doi: <u>10.2106/JBJS.18.01080</u>] [Medline: <u>31567810</u>]
- Garavito JA, Rodarte P, Navarro RA. Readability analysis of Spanish-language patient-reported outcome measures in orthopaedic surgery. J Bone Joint Surg Am 2024 Oct 16;106(20):1934-1942. [doi: <u>10.2106/JBJS.23.01367</u>] [Medline: <u>38781322</u>]
- Cook DJ, Moradkhani A, Douglas KSV, Prinsen SK, Fischer EN, Schroeder DR. Patient education self-management during surgical recovery: combining mobile (iPad) and a content management system. Telemed J E Health 2014 Apr;20(4):312-317. [doi: 10.1089/tmj.2013.0219] [Medline: 24443928]
- Cohen SM, Baimas-George M, Ponce C, et al. Is a picture worth a thousand words? A scoping review of the impact of visual aids on patients undergoing surgery. J Surg Educ 2024 Sep;81(9):1276-1292. [doi: <u>10.1016/j.jsurg.2024.06.002</u>] [Medline: <u>38955659</u>]
- 13. Liu J, Wang C, Liu S. Utility of ChatGPT in clinical practice. J Med Internet Res 2023 Jun 28;25:e48568. [doi: 10.2196/48568] [Medline: 37379067]
- Massey PA, Montgomery C, Zhang AS. Comparison of ChatGPT-3.5, ChatGPT-4, and orthopaedic resident performance on orthopaedic assessment examinations. J Am Acad Orthop Surg 2023 Dec 1;31(23):1173-1179. [doi: <u>10.5435/JAAOS-D-23-00396</u>] [Medline: <u>37671415</u>]
- 15. Scherr R, Halaseh FF, Spina A, Andalib S, Rivera R. ChatGPT interactive medical simulations for early clinical education: case study. JMIR Med Educ 2023 Nov 10;9:e49877. [doi: <u>10.2196/49877</u>] [Medline: <u>37948112</u>]
- Spina A, Andalib S, Flores D, Vermani R, Halaseh FF, Nelson AM. Evaluation of generative language models in personalizing medical information: instrument validation study. JMIR AI 2024 Aug 13;3:e54371. [doi: <u>10.2196/54371</u>] [Medline: <u>39137416</u>]
- 17. Picton B, Andalib S, Spina A, et al. Assessing AI simplification of medical texts: readability and content fidelity. Int J Med Inform 2025 Mar;195:105743. [doi: 10.1016/j.ijmedinf.2024.105743] [Medline: 39667051]
- Garcia Valencia OA, Thongprayoon C, Jadlowiec CC, et al. AI-driven translations for kidney transplant equity in Hispanic populations. Sci Rep 2024 Apr 12;14(1):8511. [doi: <u>10.1038/s41598-024-59237-7</u>] [Medline: <u>38609476</u>]
- Brewster RCL, Gonzalez P, Khazanchi R, et al. Performance of ChatGPT and Google Translate for pediatric discharge instruction translation. Pediatrics 2024 Jul 1;154(1):e2023065573. [doi: <u>10.1542/peds.2023-065573</u>] [Medline: <u>38860299</u>]
- Papineni K, Roukos S, Ward T, Zhu WJ. BLEU: a method for automatic evaluation of machine translation. Presented at: 40th Annual Meeting of the Association for Computational Linguistics; Jul 7-12, 2002; Philadelphia, Pennsylvania. [doi: 10.3115/1073083.1073135]
- 21. Hasani AM, Singh S, Zahergivar A, et al. Evaluating the performance of Generative Pre-trained Transformer-4 (GPT-4) in standardizing radiology reports. Eur Radiol 2024 Jun;34(6):3566-3574. [doi: <u>10.1007/s00330-023-10384-x</u>] [Medline: <u>37938381</u>]
- 22. Nicolson A, Dowling J, Koopman B. Improving chest x-ray report generation by leveraging warm starting. Artif Intell Med 2023 Oct;144:102633. [doi: 10.1016/j.artmed.2023.102633] [Medline: 37783533]
- Perea-Trigo M, Botella-López C, Martínez-Del-Amor M, Álvarez-García JA, Soria-Morillo LM, Vegas-Olmos JJ. Synthetic corpus generation for deep learning-based translation of Spanish sign language. Sensors (Basel) 2024 Feb 24;24(5):1472. [doi: <u>10.3390/s24051472</u>] [Medline: <u>38475008</u>]
- 24. Andalib S, Solomon SS, Picton BG, Spina AC, Scolaro JA, Nelson AM. Source characteristics influence AI-enabled orthopaedic text simplification: recommendations for the future. JB JS Open Access 2025 Jan 8;10(1):e24.00007. [doi: 10.2106/JBJS.OA.24.00007] [Medline: 39781102]

```
https://ai.jmir.org/2025/1/e70222
```

RenderX

- Spina AC, Fereydouni P, Tang JN, Andalib S, Picton BG, Fox AR. Tailoring glaucoma education using large language models: addressing health disparities in patient comprehension. Medicine (Baltimore) 2025 Jan 10;104(2):e41059. [doi: 10.1097/MD.000000000041059] [Medline: <u>39792725</u>]
- 26. Overview OpenAI API. OpenAI. URL: <u>https://platform.openai.com</u> [accessed 2025-03-03]
- 27. Bird S, Klein E, Loper E. Natural Language Processing with Python, 1st edition: O'Reilly Media Inc; 2009.
- 28. Fernández-Huerta J. Medidas sencillas de lecturabilidad [Article in Spanish]. Consigna 1959;214:29-32.
- 29. Barrio-Cantalejo IM, Simón-Lorda P, Melguizo M, Escalona I, Marijuán MI, Hernando P. Validación de la Escala INFLESZ para evaluar la legibilidad de los textos dirigidos a pacientes [Article in Spanish]. Anales Sis San Navarra 2008;31(2):135-152. [doi: 10.4321/S1137-66272008000300004] [Medline: 18953362]
- 30. Taira BR, Kreger V, Orue A, Diamond LC. A pragmatic assessment of Google Translate for emergency department instructions. J Gen Intern Med 2021 Nov;36(11):3361-3365. [doi: 10.1007/s11606-021-06666-z] [Medline: 33674922]
- 31. Patil S, Davies P. Use of Google Translate in medical communication: evaluation of accuracy. BMJ 2014 Dec 15;349:g7392. [doi: 10.1136/bmj.g7392] [Medline: 25512386]

Abbreviations

AAOS: American Academy of Orthopaedic Surgery AI: artificial intelligence BLEU: bilingual evaluation understudy LLM: large language model NLTK: Natural Language Toolkit PEM: patient education material PROM: patient-reported outcome measure

Edited by S Gardezi, Z Yin; submitted 17.12.24; peer-reviewed by C Zickler, Y Xie; revised version received 06.02.25; accepted 12.02.25; published 21.03.25.

<u>Please cite as:</u> Andalib S, Spina A, Picton B, Solomon SS, Scolaro JA, Nelson AM Using AI to Translate and Simplify Spanish Orthopedic Medical Text: Instrument Validation Study JMIR AI 2025;4:e70222 URL: <u>https://ai.jmir.org/2025/1/e70222</u> doi:<u>10.2196/70222</u>

© Saman Andalib, Aidin Spina, Bryce Picton, Sean S Solomon, John A Scolaro, Ariana M Nelson. Originally published in JMIR AI (https://ai.jmir.org), 21.3.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Deep Learning Models to Predict Diagnostic and Billing Codes Following Visits to a Family Medicine Practice: Development and Validation Study

Akshay Rajaram^{1,2}, MD, MMI; Michael Judd¹, BSc; David Barber¹, MD

¹Department of Family Medicine, Queen's University, 220 Bagot Street, Kingston, ON, Canada ²Department of Emergency Medicine, Queen's University, Kingston, ON, Canada

Corresponding Author: Akshay Rajaram, MD, MMI Department of Family Medicine, Queen's University, 220 Bagot Street, Kingston, ON, Canada

Abstract

Background: Despite significant time spent on billing, family physicians routinely make errors and miss billing opportunities. In other disciplines, machine learning models have predicted Current Procedural Terminology codes with high accuracy.

Objective: Our objective was to derive machine learning models capable of predicting diagnostic and billing codes from notes recorded in the electronic medical record.

Methods: We conducted a retrospective algorithm development and validation study involving an academic family medicine practice. Visits between July 1, 2015, and June 30, 2020, containing a physician-authored note and an invoice in the electronic medical record were eligible for inclusion. We trained 2 deep learning models and compared their predictions to codes submitted for reimbursement. We calculated accuracy, recall, precision, F_1 -score, and area under the receiver operating characteristic curve.

Results: Of the 245,045 visits eligible for inclusion, 198,802 (81%) were included in model development. Accuracy was 99.8% and 99.5% for the diagnostic and billing code models, respectively. Recall was 49.4% and 70.3% for the diagnostic and billing code models, respectively. Precision was 55.3% and 76.7% for the diagnostic and billing code models, respectively. The area under the receiver operating characteristic curve was 0.983 for the diagnostic code model and 0.993 for the billing code model.

Conclusions: We developed models capable of predicting diagnostic and billing codes from electronic notes following visits to a family medicine practice. The billing code model outperformed the diagnostic code model in terms of recall and precision, likely due to fewer codes being predicted. Work is underway to further enhance model performance and assess the generalizability of these models to other family medicine practices.

(JMIR AI 2025;4:e64279) doi:10.2196/64279

KEYWORDS

machine learning; ML; artificial intelligence; algorithm; predictive model; predictive analytics; predictive system; family medicine; primary care; family doctor; family physician; income; billing code; electronic notes; electronic health record; electronic medical record; EMR; patient record; health record; personal health record

Introduction

Previous research has revealed that family physicians spend nearly 50% of their day on electronic medical records (EMRs) and that most of this time is spent on administrative tasks, including documentation of notes and billing [1]. Physicians in the United States and Canada spend an average of 3.4 hours and 2.2 hours per week, respectively, writing, reviewing, submitting, and disputing claims with significant financial losses [2,3]. Tseng et al [4] estimated total professional billing costs for a typical primary care physician at nearly US \$100,000 using time-driven activity-based costing. In addition to billing costs, attending and resident family physicians routinely make significant errors and miss opportunities in the context of billing [5,6].

While reasons for these errors and missed opportunities are multifactorial, experts have focused on a lack of education as a primary driver [7,8]. However, the literature demonstrates that even when robust practice management curricula are introduced, billing performance does not improve significantly [9]. Moreover, experienced attending family physicians report challenges with complex billing tasks, suggesting that accumulated experience does not enhance comfort [10].

Given limitations in education and training as quality improvement interventions, other system-focused strategies are warranted [11]. One potential solution is the use of artificial intelligence to predict diagnostic and billing codes from notes.



RenderX

Kim et al [12] demonstrated 87% accuracy of their machine learning model to predict Current Procedural Terminology (CPT) codes for spine surgery from operative dictations. Another study demonstrated 98% accuracy of a neural network in assigning CPT codes to pathology reports [13].

Little is known about whether similar approaches would work in family medicine, where presenting problems and assessments are highly diverse. Our primary objective was to assess the accuracy of machine learning models in predicting diagnostic and billing codes from the notes recorded in EMRs for visits to family physicians. Based on similar studies, we hypothesized that both the diagnostic and billing code models would generate predictions with at least 90% accuracy [12-14].

Methods

Design and Setting

We conducted a retrospective model development and validation study at a large academic Family Health Team (FHT) in Ontario, Canada, with approximately 50,000 visits per year. The FHT is in a more urban setting with a patient census of approximately 21,000 rostered to 26 attending physicians. Approximately 55-60 first-year resident physicians rotate through annually.

Faculty physicians at this site are primarily compensated through capitation payments but also submit invoices for individual visits as part of the province's Family Health Organization funding model. A single-payer system predominates, with most invoices submitted to the provincial health insurance plan for reimbursement. A minority of invoices are submitted to other insurance plans, including the Workplace Safety and Insurance Board or a third party (eg, Blue Cross) or directly to patients. In addition to faculty and residents, locum physicians provide clinical coverage and submit invoices for individual visits.

Following a patient visit, physicians document their note in an EMR often in the SOAP (subjective, objective, assessment, plan) format. To submit an invoice, physicians must select 1 or more diagnostic codes and 1 or more billing codes. Invoices are compiled electronically in the EMR, reviewed by FHT billing personnel, and subsequently submitted to the provincial health insurance plan for payment every month.

Oscar is the EMR used in this study, and it contains a combination of structured and unstructured data organized into modules. Structured fields include demographics, billing (invoice number, diagnostic codes, billing codes, and billing history), preventative interventions, disease registry, laboratory results, measurements, consultations, allergies, medications, risk factors, and family history. Unstructured fields include social history, medical history, and free text chart notes.

Ethical Considerations

This study received local research ethics board approval (FMED-6780 - 20) from Queen's University Health Sciences Research Ethics Board. The approval covered secondary analyses of these data without additional consent. Physicians were given an opportunity to censor specific patients or opt out of participation. Following the opt-out process, data of the included patients were exported as a flat file and stored on a

secure server meeting local privacy requirements. Data were subsequently anonymized and deidentified during the preprocessing stage.

Participants and Sampling

Between July 1, 2015, and June 30, 2020, 245,045 visits containing a documented note and an invoice submitted to the provincial health insurance plan for payment were eligible for inclusion. The included data comprised invoices containing diagnostic and billing codes and information about the status of reimbursement, corresponding visit information including the length of appointment, the date of birth of the patient, the patient's gender, and the physician's free text note for the visit. We excluded visits that had invoices that were not paid or were deleted.

Data Preprocessing

We first transformed data into a Pandas Dataframe for additional preprocessing, including deidentification, linkage of appointments with relevant features, feature scaling, and clinical text processing.

Deidentification

Data were initially in an identifiable form but were anonymized using an automated PERL-based deidentification software package designed for free-text medical records [15]. The software uses a combination of lexical look-up tables, regular expressions, and simple heuristics to locate traditional personal health information, including common names and date variations [15]. This information was then tokenized and removed.

Linking of Appointments With Relevant Features

In Oscar, appointments are associated with both billing and diagnostic codes and contain the length of time for the visit. We linked appointments as an entity with the following data:

- 1. Demographic data for the patient, including age at the time of the appointment and gender.
- 2. Free text chart notes from the relevant table: Oscar does not relate a single note entity to an appointment. Notes were linked with their corresponding appointment by an exact match of dates. The signed and verified note by the attending physician was matched in cases of multiple notes from 1 session.
- 3. Historical diagnostic codes listed 6 months preceding the appointment date: these codes were recorded, and the frequency of the codes was summed.

Feature Scaling for Structured Data

To facilitate the use of neural networks with a gradient descent approach, we scaled our data to achieve values between 0 and 1. We used different feature scaling for different fields: (1) *MinMax scaler* from Scikit-learn for age and appointment duration [16]; (2) binary encoding for male and female; and (3) *MultiLabelBinarizer* for one-hot encoding of historical diagnostic codes [16].

Clinical Text Processing

We applied the following preprocessing steps to overcome common challenges encountered with clinical text, including

domain-specific language, spelling mistakes, and redundant phrases [17]:

- 1. Stop words: we removed stop words (eg, "a," "the," "is") from the text using the list contained in the NLTK package in Python [18].
- 2. Oscar-specific domain language: clinical notes signed by physicians include a phrase "SIGNED AND VERIFIED BY," so *regex* was applied to remove this phrase from the text.
- 3. Deidentification tokens: the deidentification tool replaces all personally identifiable information with specific tokens. We removed these tokens from the text.
- 4. Spelling mistakes: we corrected potential spelling errors by applying the Symmetric Delete spelling correction

algorithm (SymSpell) with the MEDLINE unigram dictionary, which includes over 28 million unique terms.

- 5. Punctuation: we removed punctuation from the text.
- 6. Vectorization: we vectorized the text into a sequence of numbers in the *term frequency-inverse document frequency* format [19].

Model Training and Testing

We used Tensorflow and Keras to construct one model each for the prediction of diagnostic codes and billing codes. Each model uses the same model architecture with the following layers. A graphical representation of the model architecture is presented in Figure 1.



Figure 1. Graphical representation of model architecture. ReLU: rectified linear unit.



One input layer for the vectorized note and 1 input layer are assigned for each structured data feature including age, gender, previous diagnostic codes, and appointment duration. For text classification, we used a submodel architecture called *fasttext* [20]. For structured data classification, we used a simple, fully connected, single-level Dense layer followed by a Dropout layer [21]. Weights were randomly set in the inputs. We then concatenated the text classification output layer and each structured data output layer and applied multiple layers of a Dense network followed by a Dropout layer with a rectified linear unit (ReLU) activation function. The final output layer contains a sigmoid activation function and returns multilabel outputs.

Analysis

We divided data for model development into training, testing, and validation sets, using 70% (139,161/198,802) of notes for training and 30% (59,641/198,802) for testing and validation.

XSL•FO RenderX

In the testing set, the diagnostic code model assigned 1 of 459 unique diagnostic codes while the billing code model assigned 1 of 157 unique billing codes. These codes are based on the Ontario Health Insurance Plan Schedule of Benefits for family medicine [22]. Each model initially returned a prediction score for each code ranging from 0 to 1. The prediction threshold to transform scores into labels (ie, the most likely diagnostic and billing code for the note) was selected by optimizing for the F_1 -score. The diagnostic and billing codes predicted by the deep learning models were compared to the codes selected by the clinician or updated by the FHT's billing personnel that were ultimately billed to the health insurance plan.

Given the size of both datasets, we were unable to manually review and validate the diagnostic and billing codes of notes. However, the family medicine practice in our study benefits from having dedicated administrative staff who review invoices monthly and correct errors prior to submission for reimbursement.

Several metrics of model performance, including accuracy (correct predictions divided by total predictions), recall or sensitivity (true positives/[true positives+true negatives]), precision or positive predictive value (true positives/[true positives+false positives]), F_1 -score (2*true positives/[2*true positives+false positives+false negatives]), and area under the receiver operating characteristic curve, were calculated after testing using bootstrapping. We report 95% confidence intervals. Given the multiclass nature of diagnostic and billing code prediction and anticipated class imbalances, we report microaverages as a default unless otherwise specified. We generated performance metrics using *sklearn* in Python.

Results

Of the 245,045 visits eligible for inclusion, 198,802 (81%) were included in model derivation, representing 32,425 unique patients. Three physicians opted out of participation in the study. Collectively, there were 448 unique note authors (faculty, physicians, resident physicians, or nurses). For training, 139,161 notes were used, while 29,820 and 29,821 notes were used for testing and validation, respectively. The mean length of notes was 195 (SD 102) words in the training, testing, and validation sets are compared in Table 1.

Table . (Comparison	of the training,	testing,	and validation	datasets in	model	development.
-----------	------------	------------------	----------	----------------	-------------	-------	--------------

	Training (n=139,161)	Testing (n=29,820)	Validation (n=29,821)
Ages, n (%)			·
Patients aged 0-17 years	76,539 (55)	16,341 (54.8)	16,431 (55.1)
Patients aged 18-65 years	40,078 (28.8)	8707 (29.2)	8678 (29.1)
Patients aged >65 years	22,405 (16.1)	4771 (16)	4771 (16)
Sex, n (%)			
Male patients	85,027 (61.1)	18,160 (60.9)	18,370 (61.6)
Female patients	54,134 (38.9)	11,660 (39.1)	11,451 (38.4)
Notes, mean (SD)			
Note length (number of words)	194.7 (102.2)	195.0 (102.0)	194.7 (101.4)
Number of diagnostic codes per appointment	1.3 (0.6)	1.3 (0.6)	1.3 (0.6)
Number of billing codes per appointment	1.0 (0.1)	1.0 (0.1)	1.0 (0.1)
Codes, n (%)			
799	16,268 (11.7)	3426 (11.5)	3477 (11.7)
300	7779 (5.6)	1706 (5.7)	1706 (5.7)
916	6708 (4.8)	1440 (4.8)	1428 (4.8)
250	6666 (4.8)	1381 (4.6)	1425 (4.8)
401	5747 (4.1)	1223 (4.1)	1217 (4.1)
A007A	90,803 (65.2)	19,601 (65.7)	19,470 (65.3)
A001A	7139 (5.1)	1521 (5.1)	1563 (5.2)
G590A	6596 (4.7)	1378 (4.6)	1396 (4.7)
K005A	5887 (4.2)	1279 (4.3)	1235 (4.1)
G010A	4745 (3.4)	972 (3.3)	1041 (3.5)

The overall accuracy of the diagnostic and billing code models were 99.8% (95% CI 99.79% - 99.80%) and 99.5% (95% CI 99.57% - 99.60%), respectively. The recall (sensitivity) was

https://ai.jmir.org/2025/1/e64279

RenderX

49.4% (95% CI 49.07% - 51.77%) for the diagnostic code model

and 70.3% (95% CI 68.68% - 72.17%) for the billing code

model. The precision (positive predictive value) was 55.3%

(95% CI 54.31% - 55.79%) for the diagnostic code model and 76.7% (95% CI 72.29% - 74.58%) for the billing code model. The F_1 -scores were 52.2% (95% CI 51.56% - 52.16%) and 73.4% (95% CI 72.29% - 74.58%) for the diagnostic and billing

code models, respectively. Measures of model performance are reported in Table 2. The area under the receiver operating characteristic curves for the diagnostic and billing code models are shown in Figures 2 and 3, respectively. The precision-recall curves are shown in Figures 4 and 5, respectively.

	Diagnostic code model (95% CI)	Billing code model (95% CI)			
Accuracy, %	99.8 (99.79 - 99.80)	99.5 (99.5 - 99.60)			
Recall, %	49.4 (49.07 - 51.77)	70.3 (68.68 - 72.17)			
Precision, %	55.3 (54.31 - 55.79)	76.7 (72.29 - 74.58)			
F_1 -score, %	52.2 (51.56 - 52.16)	73.4 (72.29 - 74.58)			
AUC ^a	0.983 (0.9833 - 0.9863)	0.993 (0.9921 - 0.9943)			

^aAUC: area under the receiver operating characteristic curve.







Figure 3. Area under the ROC curve for the billing code model. ROC: receiver operating characteristic.







XSL•FO RenderX

Figure 5. Precision-recall (PR) curve for the billing code model.



In the testing set, code 799 ("symptoms, signs and ill-defined conditions") was the most commonly appearing diagnostic code (n=3425) followed by code 300 ("mental disorders – neuroses and personality disorders"; n=1707) and then code 916 ("well baby care"; n=1439). Code A007 ("intermediate assessment or well baby care") was the most billed code (n=19,601). Code

A001 ("minor assessment") was the second most billed code (n=1520), followed by code G590A ("immunization – influenza agent"; n=1783). The top 10 most common diagnostic and billing codes and corresponding model performances are listed in Table 3.



Table . Prevalence and model prediction performance for the top 10 diagnostic and billing codes in the testing set.

	Description	Support, n	Precision, %	Recall, %	F_1 -score, %
Diagnostic code					
799	Symptoms, signs and ill-defined conditions	3425	78.3	63.5	70.1
300	Mental disorders – neuroses and personali- ty disorders	1707	59.2	70.2	64.3
916	Well baby care	1439	83.9	92.2	87.8
250	Diabetes mellitus in- cluding complications	1382	73.7	82.8	78.0
401	Hypertension, essential	1222	62.4	68.2	65.2
650	50 Delivery – normal; pregnancy – uncompli- cated; complications of pregnancy, childbirth and the puerperium – normal pregnancy		86.2	92.8	89.4
847	Neck strain/sprain	856	51.1	57.5	54.1
311	Depressive or other non-psychotic disorder (not classified else- where)	790	53.6	53.4	53.5
844	Strains, sprains, and other trauma – knee, leg	685	51.4	65.7	57.7
787	Abdominal – pain, masses	639	45.4	47.0	46.2
Billing code					
A007A	Intermediate assess- ment or well baby care	19,601	85.7	89.6	87.6
A001A	Minor assessment	1520	45.1	46.5	45.8
G590A	Immunization – influen- za agent	1378	91.1	63.9	75.1
K005A	Primary mental health care – individual care	1278	49.2	71.5	58.3
G010A	One or more parts of above without mi- croscopy	972	58.5	63.2	60.8
K030A	Diabetic management assessment	920	66.8	84.4	74.6
P004A	Minor prenatal assess- ment	810	80.9	93.0	86.5
E430A	Pap (Papanicolaou) smear tray fee when performed outside of hospital	681	75.2	85.9	80.2
Q015A	Newborn care episodic fee	609	65.4	74.2	69.5
G365A	Pap (Papanicolaou) smear - periodic	583	69.9	90.1	78.7

XSL•FO RenderX

Discussion

Principal Results

To our knowledge, this study is the first to report the development and internal validation of machine learning models for the prediction of diagnostic and billing codes in family medicine. While the models were highly accurate in terms of predictions, their recall and precision were much lower. These differences in performance are characteristic of multiclassification problems where high rates of overall accuracy are driven by higher classification of true negatives than identification of true positives. In the context of diagnostic and billing codes, however, correctly generating the relevant codes is much more useful than excluding irrelevant or inappropriate codes.

Unsurprisingly, the billing code model outperformed the diagnostic code model likely due to fewer codes being predicted. The lower precision and F_1 -score of the diagnostic code model suggest that the model struggles to correctly identify and classify true positive cases. There are a few possible explanations for this finding. First, the dataset was imbalanced with most diagnostic labels relating to ill-defined conditions (code 799), mental disorders (code 300), well baby care (code 916), and diabetes mellitus (code 250). Performance for these codes was noticeably better than for the overall dataset with recall ranging from 63% - 92% and precision ranging from 59% - 84%. Second, misclassification was also possible. Patients of the academic FHT where the study was conducted are known to be medically comorbid and socially complex. Consequently, encounter notes may yield several diagnostic labels; however, only 1 code may be selected for the visit.

Part of the challenge in selecting a diagnostic label for these encounters is observed among the top performing diagnostic codes. Although code 799 ("symptoms, signs and ill-defined conditions") was the most frequent code in the dataset, recall was higher for several other codes, including codes 650 ("delivery – normal; pregnancy – uncomplicated; complications of pregnancy, childbirth and the puerperium – normal pregnancy"), 916 ("well baby care") and 250 ("diabetes mellitus including complications"). These differences in performance are likely due to challenges in making sense of nonspecific symptoms in the case of code 799 as opposed to pregnancy (code 650) for a patient seeking antenatal care or a patient following up for diabetes (code 250).

We anticipated that the billing code model would perform better at predicting codes that were more frequently selected. The highest recall was for P004A, the billing code for minor prenatal assessment. Patients are seen several times during their pregnancy leading to the accumulation of these codes in historical invoices. Along with straightforward visit documentation, we suspect the model was able to predict the P004A code more fluently.

Limitations

While our study is the first to derive and validate models to predict diagnostic and billing codes in family medicine, our results should be interpreted with caution. Our data were drawn

```
https://ai.jmir.org/2025/1/e64279
```

from 1 academic FHT located in a single province and our models have not yet been externally validated. As a result, our findings may not be generalizable to other family medicine settings (eg, community or nonacademic) or other jurisdictions.

We observed heterogeneity in the performance of the model in classifying diagnostic and billing codes. Due to the size of the dataset, limited resources, and administrative constraints, we were unable to perform more detailed analyses relating to the interpretability and explainability for the diagnostic and billing code predictions. Such analyses may have uncovered factors influencing the model's performance for each code and remain an important target for future work.

One factor that likely influenced performance is clinical note quality [23]. Generally, longer notes provide more information with the corollary being that more information tends to yield better predictions. However, longer notes may also contain more copied information, which may negatively impact natural language processing performance [23]. Similarly, previous work has shown differences in the documentation practices of trainee and attending physicians [24]. The notes of trainee physicians tend to be longer and more complete while attending physicians are most interested in the assessment and plan section of notes [24-26]. Critically, quality of documentation is challenging to assess, especially in family medicine settings where no validated tools exist.

Comparison With Prior Work

Our findings are generally consistent with the results of previous studies. Using the open-source Medical Information Mart for Intensive Care III (MIMIC-III) database, various groups have developed machine learning models for the prediction of diagnostic (*International Classification of Diseases, Ninth Revision* [*ICD-9*]) codes from discharge summaries achieving micro F_1 -scores between 57.5 - 58.9 [27]. Performance discrepancies between our diagnostic code model and the models in these studies may be attributed to differences between encounter notes and discharge summaries. The latter tend to be more comprehensive in capturing details regarding a patient's initial presentation, their course and management in the hospital, and follow-up plans after discharge. These sections provide ample substrate on which to base predictions.

In the context of billing, Ye [13] developed a 3-layer neural network to predict CPT codes based on the diagnosis header and diagnosis recorded in pathology reports and achieved accuracy of 97.5%. However, their model only predicted 5 codes using text with a median length of 12 words. In contrast, Burns et al [14] developed a neural network to predict 232 CPT codes from procedural text with a mean word count of 10 words per text and achieved 82.1% accuracy. On average, notes in our study were approximately 10 times larger than those in the study by Burns et al, with a comparable number of billing codes and much higher accuracy [14].

Implications

Despite the challenges associated with billing, including missed revenue opportunities and errors, the performance of our models suggest that more work is needed before machine-learned solutions for diagnostic and billing code prediction can be

XSL•FO RenderX

deployed in practice. Such work includes external validation with other academic and community family medicine clinics, prospective validation to compare performance with physicians, and the testing of generative pretrained transformer architectures.

Once completed, there are different ways these models could be embedded within existing billing workflows. Models could be integrated with existing EMRs providing diagnostic and billing code predictions to end-users in real-time. Physicians could review predictions before finalizing codes for submission. Alternatively, physicians could bill visits as they currently do with the model surfacing its predictions for encounters for which a code was missed or an error was made. Additionally, the model could be combined with rule-based approaches to reduce common errors.

Conclusions

Our study is the first to describe the development and validation of machine learning models for the prediction of diagnostic and billing codes in family medicine. Model performance was heterogeneous and requires further analysis to uncover the factors associated with the prediction of specific diagnostic and billing codes. In addition to addressing model explainability, future work will incorporate additional structured data, consider the impacts of note characteristics and authorship on model performance, and explore validation in other family medicine settings.

Acknowledgments

We would like to thank Dr Angela Coderre-Ball for her time in reviewing and providing feedback on the manuscript.

Conflicts of Interest

AR and MJ cofounded 12676362 Canada Inc doing business as Caddie Health. Both AR and MJ hold an equity stake in the company. DB previously served as an advisor to Caddie Health and held an equity stake in the company. Caddie Health had previously licensed the models described in this work for commercialization. At the time of writing, the company is not active commercially and has no sources of revenue.

References

- Arndt BG, Beasley JW, Watkinson MD, et al. Tethered to the EHR: primary care physician workload assessment using EHR event log data and time-motion observations. Ann Fam Med 2017 Sep;15(5):419-426. [doi: <u>10.1370/afm.2121</u>] [Medline: <u>28893811</u>]
- Morra D, Nicholson S, Levinson W, Gans DN, Hammons T, Casalino LP. US physician practices versus Canadians: spending nearly four times as much money interacting with payers. Health Aff (Millwood) 2011 Aug;30(8):1443-1450. [doi: 10.1377/hlthaff.2010.0893] [Medline: 21813866]
- Dunn A, Gottlieb JD, Shapiro AH, Sonnenstuhl DJ, Tebaldi P. A denial a day keeps the doctor away. : National Bureau of Economic Research; 2021 URL: <u>https://www.nber.org/system/files/working_papers/w29010/w29010.pdf</u> [accessed 2025-02-21]
- Tseng P, Kaplan RS, Richman BD, Shah MA, Schulman KA. Administrative costs associated with physician billing and insurance-related activities at an academic health care system. JAMA 2018 Feb 20;319(7):691-697. [doi: 10.1001/jama.2017.19148] [Medline: 29466590]
- 5. Evans DV, Cawse-Lucas J, Ruiz DR, Allcut EA, Andrilla CHA, Norris T. Family medicine resident billing and lost revenue: a regional cross-sectional study. Fam Med 2015 Mar;47(3):175-181. [Medline: <u>25853527</u>]
- Al Achkar M, Kengeri-Srikantiah S, Yamane BM, Villasmil J, Busha ME, Gebke KB. Billing by residents and attending physicians in family medicine: the effects of the provider, patient, and visit factors. BMC Med Educ 2018 Jun 13;18(1):136. [doi: 10.1186/s12909-018-1246-7] [Medline: 29895287]
- Faux M, Adams J, Wardle J. Educational needs of medical practitioners about medical billing: a scoping review of the literature. Hum Resour Health 2021 Jul 15;19(1):84. [doi: <u>10.1186/s12960-021-00631-x</u>] [Medline: <u>34266457</u>]
- 8. Burks K, Shields J, Evans J, Plumley J, Gerlach J, Flesher S. A systematic review of outpatient billing practices. SAGE Open Med 2022;10:20503121221099021. [doi: 10.1177/20503121221099021] [Medline: 35646364]
- 9. Nguyen D, O'Mara H, Powell R. Improving coding accuracy in an academic practice. US Army Med Dep J 2017(2-17):95-98. [Medline: <u>28853126</u>]
- 10. Chin S, Li A, Boulet M, Howse K, Rajaram A. Resident and family physician perspectives on billing: an exploratory study. Perspect Health Inf Manag 2022;19(4):1g. [Medline: <u>36348730</u>]
- Soong C, Shojania KG. Education as a low-value improvement intervention: often necessary but rarely sufficient. BMJ Qual Saf 2020 May;29(5):353-357. [doi: 10.1136/bmjqs-2019-010411] [Medline: 31843878]
- Kim JS, Vivas A, Arvind V, et al. Can natural language processing and artificial intelligence automate the generation of billing codes from operative note dictations? Global Spine J 2023 Sep;13(7):1946-1955. [doi: <u>10.1177/21925682211062831</u>] [Medline: <u>35225694</u>]
- 13. Ye JJ. Construction and utilization of a neural network model to predict current procedural terminology codes from pathology report texts. J Pathol Inform 2019;10:13. [doi: 10.4103/jpi.jpi 3 19] [Medline: 31057982]

RenderX

- Burns ML, Mathis MR, Vandervest J, et al. Classification of current procedural terminology codes from electronic health record data using machine learning. Anesthesiology 2020 Apr;132(4):738-749. [doi: <u>10.1097/ALN.00000000003150</u>] [Medline: <u>32028374</u>]
- Neamatullah I, Douglass MM, Lehman LWH, et al. Automated de-identification of free-text medical records. BMC Med Inform Decis Mak 2008 Jul 24;8:32. [doi: <u>10.1186/1472-6947-8-32</u>] [Medline: <u>18652655</u>]
- 16. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: machine learning in Python. J Mach Learn Res 2011;12(85):2825-2830 [FREE Full text]
- Cohen R, Elhadad M, Elhadad N. Redundancy in electronic health record corpora: analysis, impact on text mining performance and mitigation strategies. BMC Bioinformatics 2013 Jan 16;14:10. [doi: <u>10.1186/1471-2105-14-10</u>] [Medline: <u>23323800</u>]
- 18. nltk package. NLTK. 2023. URL: https://www.nltk.org/api/nltk.html [accessed 2025-02-21]
- 19. Abadi M, Agarwal A, Barham P, Brevdo E, Chen Z, Citro C, et al. TensorFlow: large-scale machine learning on heterogeneous distributed systems. arXiv. Preprint posted online on Mar 14, 2016. [doi: 10.48550/arXiv.1603.04467]
- 20. Bojanowski P, Grave E, Joulin A, Mikolov T. Enriching word vectors with subword information. arXiv. Preprint posted online on Jul 15, 2016. [doi: 10.48550/arXiv.1607.04606]
- 21. Srivastava N, Hinton GE, Krizhevsky A, Sutskever I, Salakhutdinov R. Dropout: a simple way to prevent neural networks from overfitting. J Mach Learn Res 2014;15(1):1929-1958 [FREE Full text]
- 22. Schedule of benefits: physician services under the health insurance act. Government of Ontario. 2024. URL: <u>https://www.ontario.ca/files/2024-08/moh-schedule-benefit-2024-08-30.pdf</u> [accessed 2025-02-21]
- 23. Liu J, Capurro D, Nguyen A, Verspoor K. "Note Bloat" impacts deep learning-based NLP models for clinical prediction tasks. J Biomed Inform 2022 Sep;133:104149. [doi: 10.1016/j.jbi.2022.104149] [Medline: 35878821]
- Lai FW, Kant JA, Dombagolla MH, Hendarto A, Ugoni A, Taylor DM. Variables associated with completeness of medical record documentation in the emergency department. Emerg Med Australas 2019 Aug;31(4):632-638. [doi: 10.1111/1742-6723.13229] [Medline: 30690885]
- 25. Koopman RJ, Steege LMB, Moore JL, et al. Physician information needs and electronic health records (EHRs): time to reengineer the clinic note. J Am Board Fam Med 2015;28(3):316-323. [doi: <u>10.3122/jabfm.2015.03.140244</u>] [Medline: <u>25957364</u>]
- Rajaram A, Patel N, Hickey Z, Wolfrom B, Newbigging J. Perspectives of undergraduate and graduate medical trainees on documenting clinical notes: implications for medical education and informatics. Health Informatics J 2022;28(2):14604582221093498. [doi: 10.1177/14604582221093498] [Medline: 35593170]
- 27. Medical code prediction on MIMIC-III. Papers With Code. 2022. URL: <u>https://paperswithcode.com/sota/</u> medical-code-prediction-on-mimic-iii [accessed 2025-02-21]

Abbreviations

CPT: Current Procedural Terminology EMR: electronic medical record FHT: Family Health Team *ICD-9: International Classification of Diseases, Ninth Revision* MIMIC-III: Medical Information Mart for Intensive Care III ReLU: rectified linear unit SOAP: subjective, objective, assessment, plan

Edited by G Luo; submitted 13.07.24; peer-reviewed by D Nuryunarsih, Q Dong; revised version received 19.01.25; accepted 08.02.25; published 07.03.25.

<u>Please cite as:</u> Rajaram A, Judd M, Barber D Deep Learning Models to Predict Diagnostic and Billing Codes Following Visits to a Family Medicine Practice: Development and Validation Study JMIR AI 2025;4:e64279 URL: <u>https://ai.jmir.org/2025/1/e64279</u> doi:<u>10.2196/64279</u>

© Akshay Rajaram, Michael Judd, David Barber. Originally published in JMIR AI (https://ai.jmir.org), 7.3.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium,

provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.

Improving the Robustness and Clinical Applicability of Automatic Respiratory Sound Classification Using Deep Learning–Based Audio Enhancement: Algorithm Development and Validation

Jing-Tong Tzeng¹, BSc; Jeng-Lin Li², PhD; Huan-Yu Chen², PhD; Chu-Hsiang Huang³, MD; Chi-Hsin Chen³, MD; Cheng-Yi Fan³, MD; Edward Pei-Chuan Huang^{3*}, MD; Chi-Chun Lee^{2*}, PhD

¹College of Semiconductor Research, National Tsing Hua University, Hsinchu, Taiwan

²Department of Electrical Engineering, National Tsing Hua University, Hsinchu, Taiwan

³Department of Emergency Medicine, National Taiwan University Hsin-Chu Hospital, Hsinchu, Taiwan

*these authors contributed equally

Corresponding Author:

Chi-Chun Lee, PhD Department of Electrical Engineering National Tsing Hua University 101, Section 2, Kuang-Fu Road Hsinchu, 300 Taiwan Phone: 886 35162439 Email: cclee@ee.nthu.edu.tw

Abstract

Background: Deep learning techniques have shown promising results in the automatic classification of respiratory sounds. However, accurately distinguishing these sounds in real-world noisy conditions poses challenges for clinical deployment. In addition, predicting signals with only background noise could undermine user trust in the system.

Objective: This study aimed to investigate the feasibility and effectiveness of incorporating a deep learning–based audio enhancement preprocessing step into automatic respiratory sound classification systems to improve robustness and clinical applicability.

Methods: We conducted extensive experiments using various audio enhancement model architectures, including time-domain and time-frequency–domain approaches, in combination with multiple classification models to evaluate the effectiveness of the audio enhancement module in an automatic respiratory sound classification system. The classification performance was compared against the baseline noise injection data augmentation method. These experiments were carried out on 2 datasets: the International Conference in Biomedical and Health Informatics (ICBHI) respiratory sound dataset, which contains 5.5 hours of recordings, and the Formosa Archive of Breath Sound dataset, which comprises 14.6 hours of recordings. Furthermore, a physician validation study involving 7 senior physicians was conducted to assess the clinical utility of the system.

Results: The integration of the audio enhancement module resulted in a 21.88% increase with P<.001 in the ICBHI classification score on the ICBHI dataset and a 4.1% improvement with P<.001 on the Formosa Archive of Breath Sound dataset in multi-class noisy scenarios. Quantitative analysis from the physician validation study revealed improvements in efficiency, diagnostic confidence, and trust during model-assisted diagnosis, with workflows that integrated enhanced audio leading to an 11.61% increase in diagnostic sensitivity and facilitating high-confidence diagnoses.

Conclusions: Incorporating an audio enhancement algorithm significantly enhances the robustness and clinical utility of automatic respiratory sound classification systems, improving performance in noisy environments and fostering greater trust among medical professionals.

(JMIR AI 2025;4:e67239) doi:10.2196/67239

KEYWORDS

RenderX

respiratory sound; lung sound; audio enhancement; noise robustness; clinical applicability; artificial intelligence; AI

Introduction

Background

Respiratory sounds play a crucial role in pulmonary pathology. They provide insights into the condition of the lungs noninvasively and assist disease diagnosis through specific sound patterns and characteristics [1,2]. For instance, wheezing is a continuous high-frequency sound that often indicates typical symptoms of chronic obstructive pulmonary disease and asthma [3]; crackling, on the other hand, is an intermittent low-frequency sound with a shorter duration that is a common respiratory sound feature among patients with lung infections [4]. The advancement of machine learning algorithms and medical devices enables researchers to investigate approaches for developing automated respiratory sound classification systems, reducing the reliance on manual inputs from physicians and medical professionals.

In earlier studies, researchers have engineered handcrafted audio features for respiratory sound classification [5]. Recently, neural network-based methods have become the de facto methods for lung sound classification. For example, Kim et al [6] fine-tuned the pretrained VGG16 algorithm, outperforming the conventional support vector machine (SVM) classifier. Wanasinghe et al [7] incorporated mel spectrograms, mel-frequency cepstral coefficients, and chroma features to expand the feature set input to a convolutional neural network (CNN), demonstrating promising results in the identification of pulmonary diseases. Pessoa et al [8] proposed a hybrid CNN model architecture that integrates time-domain information with spectrogram-based features, delivering a satisfactory performance. Moreover, various advanced architectures have been proposed to extract both long-term and short-term information from respiratory sounds based on the characteristics of crackle and wheeze sounds and have shown enhanced performance [9-13]. Recent works have used advanced contrastive learning strategies to enhance intraclass compactness and interclass separability for further improvements [14-17]. These advancements in neural network structures have shown increasing promise in achieving reliable respiratory sound classification.

Despite these advancements, significant challenges remain for the clinical deployment of automatic respiratory sound classification systems due to complex real-world noisy conditions [6,18]. Augmentation techniques, such as time shifting, speed tuning, and noise injection, have been key strategies to effectively improve the noise robustness and generalizability of a machine learning model [9,14,16,19-23]. While these approaches have shown promising results in respiratory sound classification tasks, their practical utility as modules for building clinical decision support systems remains in doubt. This is primarily attributed to their inability to provide clinicians with intelligible raw audio to listen to facilitate decision-making, thus making the current augmentation-based approach seem black box and hindering acceptance and adoption by medical professionals.

In fact, given the blooming use of artificial intelligence (AI) in health care, the issue of liability has been the focus. The prevailing public opinion suggests that physicians are the ones to bear responsibility for errors attributed to AI [24]. Hence, when these systems are opaque and inaccessible to physicians, it becomes challenging to have them assume responsibility without a clear understanding of the decision-making process. This difficulty is particularly pronounced for seasoned and senior physicians, who hesitate to endorse AI recommendations without transparent rationale. The resulting lack of trust contributes to conflicts in clinical applications. Therefore, elucidating the decision-making process is crucial to establishing the trust of physicians [25]. Moreover, exceptions are frequent in the field of medicine. For instance, in cases in which bronchioles undergo significant constriction, the wheezing sound may diminish to near silence, a phenomenon referred to as silent wheezing. This intricacy could confound AI systems, necessitating human intervention (ie, listening directly to the recorded audio) [26].

To address these challenges, we propose an approach that involves integrating an audio enhancement module into the respiratory sound classification system, as shown in Figure 1. This module aims to achieve noise-robust respiratory sound classification performance while providing clean audio recordings on file to support physicians' decision-making. By enhancing the audio quality and preserving critical information, our system aimed to facilitate more accurate assessments and foster trust among medical professionals. Specifically, we devised 2 major experiments to evaluate this approach in this study. First, we compared the performance of our noise-robust system through audio enhancement to the conventional method of noise augmentation (noise injection) under various clinical noise conditions and signal-to-noise ratios (SNRs). Second, we conducted a physician validation study to assess confidence and reliability when listening to our cleaned audio for respiratory sound class identification. To the best of our knowledge, this is the first study showing that deep learning enhancement architecture can effectively remove noise while preserving discriminative information for respiratory sound classification algorithms and physicians. Importantly, this study validates the clinical potential and practicality of our proposed audio enhancement front-end module, contributing to more robust respiratory sound classification systems and aiding physicians in making accurate and reliable assessments.



Figure 1. An overview of our proposed noise-robust respiratory sound classification system with audio enhancement. CNN: convolutional neural network; CNN14: 14-layer CNN; conformer: convolution-augmented transformer; ISTFT: inverse short-time Fourier transform; STFT: short-time Fourier transform; TS: 2 stage.



Related Work

Audio Enhancement

Audio enhancement is a technique that has been widely used in the speech domain, where it is referred to as speech enhancement. These techniques are primarily used in the front-end stage of automatic speech recognition systems to improve intelligibility [27-29]. Within speech enhancement, deep neural network approaches can be categorized into 2 main domains: time-frequency–domain approaches and time-domain approaches.

Time-frequency-domain approaches are used to estimate clean audio from the short-time Fourier transform (STFT) spectrogram, which provides both time and frequency information. Kumar and Florencio [30] leveraged noise-aware training [31] with psychoacoustic models, which decided the importance of frequency for speech enhancement. The result demonstrated the potential of deep neural network-based speech enhancement in complex multiple-noise conditions, such as real-world environments. In the research by Yin et al [32], they designed a 2-stream architecture that predicts amplitude and phase separately and further improves the performance. However, various research studies [33-35] have indicated that the conventional loss functions used in regression models (eg, L_1 and L_2) do not strongly correlate with speech quality, intelligibility, and word error rate. To address the issue of discriminator evaluation mismatch, Fu et al [36] introduced MetricGAN. This approach tackles the problem of metrics that are not entirely aligned with the discriminator's way of distinguishing between real and fake samples. They used perceptual evaluation of speech quality (PESQ) [37] and short-time objective intelligibility (STOI) [38] as evaluation functions, which are commonly used for assessing speech quality and intelligibility, as labels for the discriminator. Furthermore, the performance of MetricGAN can be enhanced by adding a

```
https://ai.jmir.org/2025/1/e67239
```

RenderX

learnable sigmoid function for mask estimation, including noisy recording for discriminator training, and using a replay buffer to increase sample size [39]. Recently, convolution-augmented transformers (conformers) have been widely used in automatic speech recognition and speech separation tasks due to their capacity in long-range and local contexts [40-42]. Cao et al [43] introduced a conformer-based metric generative adversarial network (CMGAN), which leverages the conformer structure along with MetricGAN for speech enhancement. In the CMGAN model, multiple 2-stage conformers are used to aggregate magnitude and complex spectrogram information in the encoder. In the decoder, the prediction of the magnitude and complex spectrogram are decoupled and then jointly incorporated to reconstruct the enhanced recordings. Furthermore, CMGAN achieved state-of-the-art results on the VoiceBank+DEMAND dataset [44,45].

On the other hand, time-domain approaches directly estimate the clean audio from the raw signal, encompassing both the magnitude and phase information, enabling them to enhance noisy speech in both domains jointly. Macartney and Weyde [46] leveraged Wave-U-Net, proposed in the study by Thiemann et al [44], to use the U-Net structure in a 1D time domain and demonstrated promising results in audio source separation for speech enhancement. Wave-U-Net uses a series of downsampling and upsampling blocks with skip connections to make predictions. However, its effectiveness in representing long signal sequences is limited due to its restricted receptive field. To overcome this limitation, the approaches presented in the studies by Pandey and Wang [47] and Wang et al [48] divided the signals into small chunks and repeatedly processed local and global information to expand the receptive field. This dual-path structure successfully improved the efficiency in capturing long sequential features. However, dual-path structures are not memory efficient as they require retaining the entire long signal during training. To address the memory efficiency issue, Park et al [49] proposed a multi-view attention network.

They used residual conformer blocks to enrich channel affectin representation and introduced multi-view attention blocks consisting of channel, global, and local attention mechanisms, enabling the extraction of features that reflect both local and not spec

global information. This approach also demonstrated state-of-the-art performance on the VoiceBank+DEMAND dataset [44,45].

Both approaches have made significant progress in performance improvements in recent years. However, their suitability for enhancing respiratory sounds collected through stethoscopes remains unclear. Therefore, for this study, we applied these 2 branches of enhancement models and compared their effectiveness in enhancing respiratory sounds in real-world noisy hospital settings [32,43,46,49].

Respiratory Sound Classification

In recent years, automatic respiratory sound classification systems have become an active research area. Several studies have explored the use of pretrained weights from deep learning models, showing promising results. Kim et al [6] demonstrated improved performance over SVMs by fine-tuning the pretrained VGG16 algorithm. Gairola et al [22] used effective preprocessing methods, data augmentation techniques, and transfer learning from ImageNet [50] pretrained weights to address data scarcity and further enhance performance.

As large-scale audio datasets [51,52] become more accessible, pretrained audio models are gaining traction, exhibiting promising performance in various audio tasks [53-55]. Studies have explored leveraging these pretrained audio models for respiratory sound classification. Moummad and Farrugia [17] incorporated supervised contrastive loss on metadata with the pretrained 6-layer CNN architecture [53] to improve the quality of learned features from the encoder. Chang et al [56] introduced a novel gamma patch-wise correction augmentation technique, which they applied to the fine-tuned 14-layer CNN (CNN14) architecture [53], achieving state-of-the-art performance. Bae et al [16] used the pretrained Audio Spectrogram Transformer (AST) [54] with a Patch-Mix strategy to prevent overfitting and improve performance. Kim et al [57] proposed a representation-level augmentation technique to effectively leverage different pretrained models with various input types, demonstrating promising results on the pretrained ResNet, EfficientNet, 6-layer CNN, and AST.

However, few of these studies have explicitly addressed the challenge of noise robustness in clinical settings. To improve noise robustness, data augmentation techniques such as adding white noise, time shifting, stretching, and pitch shifting have been commonly used [9,14]. These augmentations enable networks to learn efficient features under diverse recording conditions. Nonetheless, the augmented recordings may not accurately represent the conditions in clinical settings, potentially introducing artifacts and limiting performance improvement. In contrast to the aforementioned works, Kochetov et al [18] proposed a noise-masking recurrent neural network to filter out noisy frames during classification. They concatenated a binary noise classifier and an anomaly classifier with a mask layer to suppress the noisy preventing noises from

affecting the classification. However, the International Conference in Biomedical and Health Informatics (ICBHI) database lacks noise labels in the metadata, and the paper did not specify how these labels were obtained, rendering the results nonreproducible. Emmanouilidou et al [58] used multiple noise suppression techniques to address various noise sources, including ambient noise, signal artifacts, heart sounds, and crying, using a soft-margin nonlinear SVM classifier with handcrafted features. Similarly, our work uses a pipeline for noise enhancement and respiratory sound classification. However, we advanced this approach by using deep learning models for both tasks, enabling our system to handle diverse noise types and levels without the need for bespoke strategies for each noise source. Furthermore, we validated our system's practical utility through experiments across 2 respiratory sound databases and a physician validation study, demonstrating its improved performance and clinical relevance.

Methods

Datasets

This section presents 2 respiratory sound datasets and 1 clinical noise dataset used in this study.

ICBHI 2017 Dataset

The ICBHI 2017 database is one of the largest publicly accessible datasets for respiratory sounds, comprising a total of 5.5 hours of recorded audio [59]. These recordings were independently collected by 2 research teams in Portugal and Greece from 126 participants of all ages (79 adults, 46 children, and 1 unknown). The data acquisition process involved heterogeneous equipment and included recordings from both clinical and nonclinical environments. The duration of the recorded audio varies from 10 to 90 seconds. Within this database, 6898 respiratory cycles result in 920 annotated audio samples. Among these samples, 1864 contain crackles, 886 contain wheezes, and 506 include both crackles and wheezes, whereas the remaining cycles are categorized as normal.

Formosa Archive of Breath Sound

The Formosa Archive of Breath Sound (FABS) database comprises 14.6 hours of respiratory sound recordings collected from 1985 participants. Our team collected these recordings at the emergency department of the Hsin-Chu Branch at the National Taiwan University Hospital (NTUH). We used the CaRDIaRT DS101 electronic stethoscope, where each recording is 10 seconds long.

To ensure the accuracy of the annotations, a team of 7 senior physicians meticulously annotated the audio samples. The annotations focused on identifying coarse crackles, wheezes, or normal respiratory sounds. Unlike the ICBHI 2017 database, our annotation process treated each audio sample in its entirety rather than splitting it into respiratory cycles. This approach reduces the need for extensive segmentation procedures and aligns with regular clinical practice. To enhance the quality of the annotations, we implemented an annotation validation flow called "cross-annotator model validation." This involved training multiple models based on each annotator's data and validating the models on data from other annotators. Any data with

incongruent predictions were initially identified. These data then underwent additional annotation by 3 senior physicians randomly selected from the original annotation team for each sample to achieve the final consensus label. The FABS database encompasses 5238 annotated recordings, with 715 containing coarse crackles, 234 containing wheezes, and 4289 labeled as normal respiratory sound recordings. The detailed comparison between the ICBHI 2017 dataset and the FABS database is shown in Table 1.

Table 1. Comparison between the International Conference in Biomedical and Health Informatics (ICBHI) and Formosa Archive of Breath Sound (FABS) datasets.

	ICBHI (n=126 patients)	FABS (n=1985 patients)
Age (y), mean (SD)	42.99 (32.08)	66.04 (17.64)
BMI (kg/m ²), mean (SD)	27.19 (5.34)	23.95 (4.72)
Sex, n (%)		
Male	79 (62.7)	974 (49.1)
Female	46 (36.5)	841 (42.4)
Unknown	1 (0.8)	170 (8.6)
Sampling rate (kHz)	4-44.1	16
Duration (hours)	5.5	14.6
Label	Crackle and wheeze, crackle, wheeze, and normal	Coarse crackle, wheeze, and normal
Equipment	AKG C417L microphone, Littmann Classic II SE stethoscope, Littmann 3200 electronic stethoscope, and Welch Allyn Meditron electronic stethoscope	CaRDIaRT DS101 electronic stethoscope

NTUH Clinical Noise Dataset

The noise dataset used in this study was sourced from the NTUH Hsin-Chu Branch. To replicate the noise sounds that physicians typically encounter in real-world clinical settings, we used the CaRDIaRT DS101 electronic stethoscope for collecting the noise samples. The NTUH clinical noise dataset consists of 3 different types of clinical noises: 8 friction noises produced by the stethoscope moving on different fabric materials; 18 environment noises recorded at various locations within the hospital; and 12 patient noises generated by patients during auscultation through conversations, coughing, and snoring.

Proposed Methods

As shown in Figure 1, our proposed noise-robust respiratory sound classification system includes two main components: (1) audio enhancement and (2) respiratory sound classifier.

Audio Enhancement Module

Audio enhancement is usually approached as a supervised learning problem [30,31,33-36,39,43], where the goal is to map noisy respiratory sound inputs to their clean counterparts. Mathematically, this task can be represented as learning a function f, mapping X_{noisy} to X_{clean} , where X_{noisy} represents the input noisy sound and X_{clean} denotes the corresponding clean sound. The enhanced output, X'_{clean} , is obtained as $X'_{\text{clean}}=f(X_{\text{noisy}})$ (1), where f is the audio enhancement model optimized during training.

To achieve high-quality enhancement, it is crucial to carefully select reference clean recordings from the respiratory sound database to generate high-quality paired noisy-clean sound data. To address this, we used an "audio-tagging filter" approach. This approach leverages a large pretrained audio-tagging model

```
https://ai.jmir.org/2025/1/e67239
```

to identify clean samples and exclude recordings with irrelevant tags from the database. Specifically, we used the CNN14 pretrained audio neural network [53] that was trained on AudioSet [51], a comprehensive audio dataset containing 2,063,839 training audio clips sourced from YouTube covering 527 sound classes. Audio samples with the following audio event labels were filtered out: "music," "speech," "fire," "animal," "cat," and "domestic animals, pets." These labels were chosen as they were among the top commonest predictions of the audio-tagging model, indicating a higher likelihood of significant irrelevant noise in the recordings. By excluding these labels, we could ensure that the selected recordings could be used as reference clean audio. To validate the effectiveness of the filtering process, we manually checked the filtered recordings. The results showed that the tagging precision was 92.5%, indicating that this method is efficient and trustworthy. Moreover, as it is fully automatic, it is easy to reproduce the results.

In the ICBHI 2017 database, 889 clean audio samples were retained after filtering, consisting of 1812 cycles with crackling sounds, 822 cycles with wheezing sounds, 447 cycles with both crackling and wheezing sounds, and 3538 cycles with normal respiratory sounds. Alternatively, the filtered FABS clean samples comprised 699 recordings of coarse crackle respiratory sounds, 225 recordings of wheeze respiratory sounds, and 4238 recordings of normal respiratory sounds.

In this study, we used Wave-U-Net [46], Phase-and-Harmonics–Aware Speech Enhancement Network (PHASEN) [32], Multi-View Attention Network for Noise Erasure [49], and CMGAN [43] to compare the effectiveness of different model structures in enhancing respiratory sounds.

Respiratory Sound Classification

Training a classification model from scratch using a limited dataset may lead to suboptimal performance or overfitting. Therefore, we selected the CNN14 model proposed in the study by Kong et al [53], which had been pretrained on AudioSet [51], as our main classification backbone, and we further fine-tuned it on our respiratory datasets. We used log-mel spectrograms as the input feature, similar to previous works in respiratory sound classifications [6,9-11,14]. As the dataset is highly imbalanced, we used the balanced batch-learning strategy. To further improve model generalizability and performance, we incorporated data augmentation techniques, including Mixup [60] and SpecAugment [61], along with triplet loss [15,62] to enhance feature separability.

Mathematically, the classification task is formulated as a multi-class classification problem. The goal is to learn a mapping function, $g: Z \rightarrow Y(2)$, where Z represents the extracted features and Y denotes the target class labels. To obtain Z, input-enhanced audio signals X'_{clean} are transformed using the STFT to generate a spectrogram, followed by mel-filter banks to convert the frequency scale to the mel scale: $Z=log-mel(STFT[X'_{clean}])$ (3).

During training, the total loss function L_c combines cross-entropy loss and triplet loss: $L_c = L_{CE} + \lambda L_{triplet}$ (4).

Through grid search, λ =0.01 leads to the best performance.

Physician Validation Study

To further evaluate the effectiveness of audio enhancement for respiratory sound, we conducted a physician validation study

Textbox 1. Methods for various levels of noise intensity.

Clean

The respiratory sound classification models were only trained on clean data and tested on clean data. This approach served to establish the upper-bound performance for the overall comparison.

Noisy

The respiratory sound classification models were trained on clean data but tested on noisy data. As the models were not optimized for noise robustness, a significant drop in performance was expected.

Noise injection

The respiratory sound classification models were trained on synthesized noisy data and tested on noisy data. This approach represents the conventional method to enhance the noise robustness of the model.

Audio enhancement

The audio enhancement model functions as a front-end preprocessing step for the classification model. To achieve this, we first optimized the audio enhancement model to achieve a satisfactory enhancement performance. Subsequently, the respiratory sound classification model was trained on the enhanced data and tested on the enhanced data.

Experiment Setup

To evaluate the efficiency of our proposed method, we followed a similar setup as that in prior work [6,11,14] to have an 80%-20% train-test split on the database. Furthermore, the training set was mixed with the noise recordings from the NTUH clinical noise dataset with 4 SNRs (15, 10, 5, and 0 dB) with random time shifting. The test set was mixed with unseen noise data with 4 SNRs (17.5, 12.5, 7.5, and 2.5 dB), also subjected to random time shifting. For evaluation, we used the metrics of

accuracy, sensitivity, specificity, and ICBHI score. Sensitivity is defined as the recall of abnormal respiratory sounds. Specificity refers to the recall of normal respiratory sounds. The ICBHI score, calculated as the average of sensitivity and specificity, provides a balanced measure of the model's classification performance.

Tzeng et al

using the clean, noisy, and enhanced recordings from a randomly selected 25% of the testing set on the ICBHI 2017 database. In this study, we invited 7 senior physicians to independently annotate these recordings without access to any noise level or respiratory sound class label. We instructed the physicians to label the respiratory class with a confidence score ranging from 1 to 5. The objective was to demonstrate that our proposed method not only enhances the performance of the classification model but also improves the accuracy of the respiratory sound classification and increases the confidence in manual judgment done by physicians. The physician validation study was a critical step in validating the clinical practicality and effectiveness of our proposed audio enhancement preprocessing technique in clinical settings.

Ethical Considerations

This study was approved by the institutional review board of the NTUH Hsin-Chu Branch (109-129-E) and complies with ethical guidelines for human research. It involved both prospective and retrospective data collection, with retrospective data fully deidentified to protect participant privacy. All prospective participants provided informed consent before data collection. No financial compensation was provided to participants, ensuring voluntary and unbiased participation.

Results

Overview

To assess the noise robustness of our proposed method, we conducted a comparative analysis using methods across various levels of noise intensity, as outlined in Textbox 1.

Implementation Details

Technical Setup

The models were implemented using PyTorch (version 1.12; Meta AI) with the CUDA Toolkit (version 11.3; NVIDIA Corporation) for graphics processing unit acceleration. Training was conducted on an NVIDIA A100 graphics processing unit with 80 GB of memory. For clarity and reproducibility, the detailed implementation and computational setup is provided in Multimedia Appendix 1.

Preprocessing

We first resampled all recordings to 16 kHz. Next, each respiratory cycle was partitioned into 10-second audio segments before proceeding with feature extraction. In cases in which cycles were shorter in duration, we replicated and concatenated them to form 10-second clips in the ICBHI dataset. As the recordings in the FABS dataset are initially labeled per recording, there was no requirement for a segmentation process. Subsequently, these audio clips were mixed with the NTUH clinical noise dataset, generating pairs of noisy and clean data for further processing.

Enhancement Model Training

For enhancement model training, the 10-second audio clips were divided into 4-second segments. When implementing Wave-U-Net [43], the channel size was set to 24, the batch size was set to 4, and the number of layers of convolution upsampling and downsampling was set to 8. The model was trained using the Adam optimizer with a learning rate of 10^{-5} for 40 epochs when training using pretrained weights and 10^{-4} for 30 epochs when training from scratch. For the Multi-View Attention Network for Noise Erasure model [49], the channel size was set to 60, the batch size was set to 4, and the number of layers of up and down convolution was set to 4. The model was trained using the Adam optimizer with a learning rate of 10^{-6} for 10 epochs when training using pretrained weights and a learning rate of 10^{-5} for 10 epochs when training from scratch. When implementing PHASEN [32], which is trained in the time-frequency domain, we followed the original setup using a Hamming window of 25 ms in length and a hop size of 10 ms to generate STFT spectrograms. The number of 2-stream blocks was set to 3, the batch size was set to 4, the channel number for the amplitude stream was set to 24, and the channel number for the phase stream was set to 12. The model was trained using the Adam optimizer with a learning rate of 5×10^{-5} for 20 epochs when training using pretrained weights and a learning rate of 5×10^{-4} for 30 epochs when training from scratch. For CMGAN [43], we followed the original setting using a Hamming window of 25 ms in length and a hop size of 6.25 ms to generate STFT spectrograms. The number of 2-stage conformer blocks was set to 4, the batch size was set to 4, and the channel number in the generator was set to 64. The channel numbers in the discriminator were set to 16, 32, 64, and 128. The model was trained using the Adam optimizer with a learning

rate of 5×10^{-5} for 20 epochs when training using pretrained weights and a learning rate of 5×10^{-4} for 30 epochs when training from scratch. These hyperparameters are also listed in Multimedia Appendix 2.

The pretrained weights for these models were trained on the VoiceBank+DEMAND dataset [44,45], which is commonly used in speech enhancement research.

Classification Model Training

For the classification model, the 4-second enhanced segments were concatenated back into 10-second audio clips. To generate the log-mel spectrogram, the waveform was transformed using STFT with a Hamming window size of 512 and a hop size of 160 samples. The STFT spectrogram was then processed through 64 mel filter banks to generate the log-mel spectrogram. In the training stage, we set the batch size to 32 and used the Adam optimizer with a learning rate of 10^{-4} for 14,000 iterations using pretrained weights from the model trained on the 16-kHz AudioSet dataset [51]. These hyperparameters are also listed in Multimedia Appendix 2.

Evaluation Outcomes

In this study, we compared the classification performance of conventional noisy data augmentation with our proposed audio-enhanced preprocessing. The test set was split into 2 groups, and each classification model was trained 10 times, yielding 20 values for statistical analysis. We conducted a 1-tailed t test to assess whether models trained on CMGAN-enhanced audio using pretrained weights showed significant improvements over other models. In addition, we reported speech quality metrics for various audio enhancement models and analyzed their correlation with classification performance.

The experiment results, as shown in Table 2, highlight the effectiveness of our proposed audio enhancement preprocessing strategy for noise-robust performances. In the case of the ICBHI 2017 database, the model trained solely on clean data experienced a 33.95% drop in the ICBHI score when evaluated on the synthesized noisy dataset. Noise injection improved the score by 19.73%, but fine-tuning PHASEN achieved the highest score, outperforming noise injection by 2.28%. Regarding the FABS database, using the classification model trained on clean recordings on the noisy recordings led to a 12.48% drop in the ICBHI score. Noise injection improved performance by 1.31%, but fine-tuning CMGAN outperformed noise injection by 2.79%. Across both datasets, the audio enhancement preprocessing method consistently improved performance compared to the noise injection augmentation technique. Furthermore, it showed improved sensitivity for all enhancement model structures, with the most significant improvement being 6.31% for the ICBHI database and 13.54% for the FABS database. This indicates that the audio enhancement preprocessing method enhanced the classification model's ability to distinguish abnormal respiratory sounds, which is crucial for the early detection of potential illnesses in clinical use.

RenderX

Tzeng et al

 Table 2. Comparison of classification performance on both the International Conference in Biomedical and Health Informatics (ICBHI) and Formosa Archive of Breath Sound (FABS) datasets.

Method	Enhancement model	Accuracy, mean (SD)	P value	Sensitivity, mean (SD)	P value	Specificity, mean (SD)	P value	ICBHI score, mean (SD)	P value
ІСВНІ			·	·	·	·			
Clean	a	79.90 (0.01)	>.99	71.43 (0.02)	>.99	87.27 (0.01)	>.99	79.35 (0.01)	>.99
Noisy	_	45.70 (0.03)	<.001	40.99 (0.04)	<.001	49.80 (0.08)	<.001	45.40 (0.03)	<.001
Noise injec- tion	_	65.85 (0.01)	<.001	54.89 (0.04)	<.001	75.37 (0.04)	.98	65.13 (0.01)	<.001
AE ^b	Wave-U-Net	60.86 (0.02)	<.001	55.35 (0.04)	<.001	65.66 (0.05)	<.001	60.50 (0.02)	<.001
AE	Wave-U-Net ^c	61.29 (0.02)	<.001	55.04 (0.02)	<.001	66.72 (0.04)	<.001	60.88 (0.02)	<.001
AE	PHASEN ^d	66.81 (0.01)	.02	57.61 (0.03)	.001	74.81 (0.04)	.91	66.21 (0.01)	.005
AE	PHASEN ^c	68.09 ^e (0.01)	.84	57.71 ^f (0.03)	.004	77.12 ^f (0.04)	>.99	67.41 ^e (0.01)	.64
AE	MANNER ^g	67.62 (0.01)	.39	53.09 (0.03)	<.001	80.26 ^e (0.04)	>.99	66.67 (0.01)	.03
AE	MANNER ^c	60.36 (0.02)	<.001	57.67 (0.02)	<.001	62.70 (0.04)	<.001	60.19 (0.02)	<.001
AE	CMGAN ^h	64.75 (0.01)	<.001	55.84 (0.03)	<.001	72.50 (0.02)	.17	64.17 (0.01)	<.001
AE	CMGAN ^c	67.70 ^f (0.01)	_	61.20 ^e (0.03)	_	73.35 (0.02)	_	67.28 ^f (0.01)	_
FABS									
Clean		85.02 (0.01)	>.99	62.07 (0.04)	>.99	90.01 (0.02)	<.001	76.04 (0.02)	>.99
Noisy	_	81.02 (0.02)	<.001	36.41 (0.04)	<.001	90.71 (0.02)	.004	63.56 (0.02)	<.001
Noise injec- tion	_	84.53 (0.01)	>.99	34.29 (0.05)	<.001	95.44 (0.01)	>.99	64.87 (0.02)	<.001
AE	Wave-U-Net	85.97 ^e (0.01)	>.99	36.74 (0.03)	<.001	96.66 ^f (0.01)	>.99	66.70 (0.01)	.04
AE	Wave-U-Net ^c	85.88 ^f (0.01)	>.99	29.08 (0.05)	<.001	98.22 ^e (0.01)	>.99	63.65 (0.02)	<.001
AE	PHASEN	85.29 (0.004)	>.99	33.64 (0.02)	<.001	96.51 (0.01)	>.99	65.07 (0.01)	<.001
AE	PHASEN ^c	85.33 (0.01)	>.99	35.82 (0.03)	<.001	96.09 (0.01)	>.99	65.95 (0.02)	<.001
AE	MANNER	83.01 (0.01)	.05	37.50 (0.08)	.01	92.89 (0.03)	.67	65.20 (0.03)	.004
AE	MANNER ^c	79 (0.03)	<.001	47.83 ^e (0.06)	>.99	85.77 (0.05)	<.001	66.80 ^f (0.02)	.08
AE	CMGAN	82.47 (0.01)	<.001	37.61 (0.05)	<.001	92.22 (0.01)	.19	64.91 (0.02)	<.001
AE	CMGAN ^c	83.67 (0.01)	—	42.77 ^f (0.03)	—	92.55 (0.01)	_	67.66 ^e (0.01)	—

^aWithout any audio enhancement module.

https://ai.jmir.org/2025/1/e67239



^bAE: audio enhancement.

^cThe model is fine-tuned from the pretrained weight.

^dPHASEN: Phase-and-Harmonics–Aware Speech Enhancement Network.

^eBest performance across all methods for this metric.

^fSecond-best performance across all methods for this metric.

^gMANNER: Multi-View Attention Network for Noise Erasure.

^hCMGAN: convolution-augmented transformer–based metric generative adversarial network.

Comparing the 2 types of enhancement approaches, the time-frequency domain models (PHASEN and CMGAN) exhibited better performance in terms of ICBHI scores. In addition, CMGAN consistently showed high sensitivity across both datasets, indicating its potential for preserving respiratory sound features during audio enhancement. The spectrogram of the audio enhanced using CMGAN also revealed that it preserves more high-frequency information across all respiratory sound classes, as illustrated in Figure 2. In contrast, audio enhanced using other models either lost high-frequency

information or retained too much noise, leading to misclassification as normal, resulting in higher specificity for those models. Moreover, we observed that, while our focus was on training a respiratory sound enhancement model, using pretrained weights from models trained on the VoiceBank+DEMAND dataset, which were originally designed for speech, still significantly improved classification performance in most cases. This highlights the cross-domain effectiveness of pretrained weights from the speech domain in respiratory sound tasks.

Figure 2. The log-mel spectrograms of 4 different types of respiratory sounds on the International Conference in Biomedical and Health Informatics 2017 database. Each subfigure contains clean audio, noisy audio, and 4 types of enhanced audio from different audio enhancement approaches. CMGAN: convolution-augmented transformer–based metric generative adversarial network; MANNER: Multi-View Attention Network for Noise Erasure; PHASEN: Phase-and-Harmonics–Aware Speech Enhancement Network.



https://ai.jmir.org/2025/1/e67239

RenderX

To evaluate whether speech quality metrics, originally designed for speech, are effective for respiratory sounds, we analyzed their correlation with the ICBHI score and sensitivity. As shown in Table 3, the mean opinion score (MOS) of background noise intrusiveness (CBAK) and segmental SNR (SSNR) exhibited relatively higher correlations than other metrics, such as PESQ, STOI, the MOS of signal distortion, and the MOS of overall quality. Unlike these other metrics, which are primarily designed to assess speech intelligibility and quality, CBAK and SSNR focus on background noise intrusiveness and the SNR between recordings. This distinction explains why CBAK and SSNR show stronger correlations with classification performance, highlighting their potential applicability for respiratory sound analysis.

We evaluated the inference times of 4 audio enhancement models. Wave-U-Net generates 1 second of enhanced audio in just 1.5 ms, PHASEN does so in 3.9 ms, and MANNER does so in 11.7 ms. In contrast, CMGAN processes 1 second of audio in 26 ms—a longer time that is offset by its superior classification performance.

To further analyze the effectiveness of our proposed audio enhancement preprocessing method in handling different types of noise, we compared its performance using the noise injection method across various SNR levels. On the basis of the consistently outstanding performance of CMGAN across both datasets, we selected it for further analysis.

On the ICBHI database, as illustrated in Figure 3, the noise injection method performed better with environmental noises at SNR values of 2.5 and 12.5 dB. However, the front-end audio enhancement consistently performed better for patient and friction noises across almost all noise levels.

Regarding the FABS dataset, as shown in Figure 4, the noise injection method performed better with environmental and friction noises at an SNR value of 17.5 dB and patient noises at an SNR value of 2.5 and 7.5 dB. In all other situations, the audio enhancement preprocessing method demonstrated superior ICBHI scores.

These results suggest that our proposed strategy effectively mitigates the effects of various noise types while maintaining strong classification performance. This highlights the robustness and reliability of our approach in handling diverse noise scenarios and intensities, showcasing its potential for practical applications in clinical settings.



Tzeng et al

Table 3. Comparison of audio enhancement (AE) performance on both the International Conference in Biomedical and Health Informatics (ICBHI) and Formosa Archive of Breath Sound (FABS) datasets.

Method	Enhancement model	Parameters (mil- lions)	PESQ ^{a,b}	CSIG ^{c,d}	CBAK ^{e,f}	COVL ^{g,h}	SSNR ^{i,j}	STOI ^{k,1}
ICBHI			·					
Noisy	m	_	0.58	2.98	2.83	2.13	14.10	0.50
AE	Wave-U-Net	3.3	0.56	3.07	3.25	2.18	20.30	0.49
AE	Wave-U-Net ⁿ	3.3	0.57	3.10	3.25	2.20	20.20	0.50
AE	PHASEN ^o	7.7	0.57	3.07	3.34	2.19	21.41	0.52
AE	PHASEN ⁿ	7.7	0.56	3.04	3.32	2.17	21.26	0.51
AE	MANNER ^p	24	0.59	3.23	3.24	2.27	19.85	0.55
AE	MANNER ⁿ	24	0.66	3.38 ^q	3.24	2.39 ^r	19.17	0.60 ^r
AE	CMGAN ^s	1.8	0.75 ^q	3.31 ^r	3.46 ^r	2.40 ^q	22.06 ^r	0.61 ^q
AE	CMGAN ⁿ	1.8	0.74 ^r	3.29	3.47 ^q	2.38	22.31 ^q	0.61 ^q
FABS								
Noisy	_	_	2.10	3.80 ^q	3.41	3.03 ^q	12.99	0.62 ^r
AE	Wave-U-Net	3.3	1.78	1.96	3.16	1.90	10.97	0.52
AE	Wave-U-Net ⁿ	3.3	1.75	1.89	3.13	1.86	10.74	0.50
AE	PHASEN	7.7	1.93	2.34	3.26	2.19	11.54	0.58
AE	PHASEN ⁿ	7.7	1.84	2.11	3.20	2.03	11.27	0.57
AE	MANNER	24	2.14 ^r	3.35	3.44 ^r	2.81	12.87	0.61
AE	MANNER ⁿ	24	2.18 ^q	3.57 ^r	3.44 ^r	2.95 ^r	12.57	0.63 ^q
AE	CMGAN	1.8	2.01	1.79	3.42	1.96	13.59 ^r	0.59
AE	CMGAN ⁿ	1.8	2.06	1.68	3.48 ^q	1.91	13.98 ^q	0.59

^aPESQ: perceptual evaluation of speech quality.

^bICBHI: sensitivity correlation coefficient=0.36 and ICBHI score correlation coefficient=0.23; FABS: sensitivity correlation coefficient=0.72 and ICBHI score correlation coefficient=0.16.

^cCSIG: mean opinion score (MOS) of signal distortion.

^dICBHI: sensitivity correlation coefficient=0.51 and ICBHI score correlation coefficient=0.40; FABS: sensitivity correlation coefficient=0.34 and ICBHI score correlation coefficient=-0.25.

^eCBAK: MOS of background noise intrusiveness.

^fICBHI: sensitivity correlation coefficient=0.92 and ICBHI score correlation coefficient=0.90; FABS: sensitivity correlation coefficient=0.71 and ICBHI score correlation coefficient=0.23.

^gCVOL: MOS of overall quality.

^hICBHI: sensitivity correlation coefficient=0.52 and ICBHI score correlation coefficient=0.39; FABS: sensitivity correlation coefficient=0.42 and ICBHI score correlation coefficient=-0.20.

ⁱSSNR: segmental signal-to-noise ratio.

^jICBHI: sensitivity correlation coefficient=0.92 and ICBHI score correlation coefficient=0.93; FABS: sensitivity correlation coefficient=0.59 and ICBHI score correlation coefficient=0.22.

^kSTOI: short-time objective intelligibility.

^lICBHI: sensitivity correlation coefficient=0.45 and ICBHI score correlation coefficient=0.36; FABS: sensitivity correlation coefficient=0.68 and ICBHI score correlation coefficient=0.13.

^mWithout any audio enhancement module.

ⁿThe model is fine-tuned from the pretrained weight.

^oPHASEN: Phase-and-Harmonics-Aware Speech Enhancement Network.

^pMANNER: Multi-View Attention Network for Noise Erasure.

^qBest performance across all methods for this metric.

XSL•F() RenderX
^rSecond-best performance across all methods for this metric.

^sCMGAN: convolution-augmented transformer-based metric generative adversarial network.

Figure 3. Performance comparison of different approaches for each noise type with various signal-to-noise ratio (SNR) values on the International Conference in Biomedical and Health Informatics (ICBHI) 2017 database.



Figure 4. Performance comparison of different approaches for each noise type with various signal-to-noise ratio (SNR) values on the Formosa Archive of Breath Sound database. ICBHI: International Conference in Biomedical and Health Informatics.



https://ai.jmir.org/2025/1/e67239

Physician Validation Study

To assess the practical utility of our proposed approach in clinical settings, we conducted a physician validation study using the ICBHI dataset. This study involved comparing the annotation results provided by 7 senior physicians under 3 different conditions: clean, noisy, and enhanced recordings. By evaluating physician assessments across these conditions, we aimed to determine the effectiveness of our enhancement approach in improving diagnostic accuracy and confidence.

As shown in Table 4, the presence of noise in the recordings had a noticeable impact on the physicians' ability to conduct a reliable judgment, reducing accuracy by 1.81% and sensitivity by 6.46% compared to the clean recordings. However, the recordings with audio enhancement exhibited notable

improvement, with a 3.92% increase in accuracy and an 11.61% increase in sensitivity compared to the noisy recordings. The enhanced audio successfully preserved sound characteristics crucial for physicians in classifying respiratory sounds, leading to higher true positive rates in distinguishing adventitious sounds.

The enhanced audio recordings also received higher annotation confidence scores than the noisy recordings, as indicated in Figure 5 and Table 4. Moreover, the speech quality metrics PESQ, MOS of signal distortion, CBAK, MOS of overall quality, SSNR, and STOI positively correlated with the physicians' annotation confidence, as shown in Figure 6. These results underscore the potential of audio enhancement preprocessing techniques for practical application in real-world clinical settings.

Table 4. Annotation results from physicians on different types of recordings.

Type of recording	Accuracy (%)	Sensitivity (%)	Specificity (%)	ICBHI ^a score (%)	Confidence mean (SD)
Clean	49.4	23.23	72.32	47.77	2.88 (1.50)
Noisy	47.59	16.77	74.58	45.68	2.32 (1.29)
Enhanced	51.51	28.38	71.75	50.07	2.65 (1.36)

^aICBHI: International Conference in Biomedical and Health Informatics.

Figure 5. Physicians' annotation confidence score comparison among clean, noisy, and enhanced recordings.





Figure 6. Relationship between physicians' annotation confidence score and speech quality metrics. CBAK: mean opinion score (MOS) of background noise intrusiveness; CSIG: MOS of signal distortion; CVOL: MOS of overall quality; PESQ: perceptual evaluation of speech quality; SSNR: segmental signal-to-noise ratio; STOI: short-time objective intelligibility.



Ablation Study

Other Classification Model

To assess the effectiveness of our proposed speech enhancement preprocessing technique with different classification models, we conducted an ablation study. The hyperparameters used in this study are detailed in Multimedia Appendix 2. We used the fine-tuned CMGAN as the speech enhancement module as it showed consistently outstanding performance in previous experiments, as shown in Table 2.

For the ICBHI dataset, the speech enhancement preprocessing technique increased the sensitivity by 11.71% and the ICBHI score by 1.4% when using the AST model [54]. Similarly, when using the AST model with the Patch-Mix strategy [16], the speech enhancement preprocessing technique increased the

sensitivity by 17.08% and the ICBHI score by 1.6%, as shown in Tables 5 and 6.

4

5

Regarding the FABS dataset, the speech enhancement preprocessing technique increased the sensitivity by 18.48% and the ICBHI score by 5.46% when fine-tuning the AST model [54]. When fine-tuning the AST model using the Patch-Mix strategy [16], the speech enhancement preprocessing technique increased the sensitivity by 13.04% and the ICBHI score by 0.68%, as shown in Tables 7 and 8.

These results demonstrate that the speech enhancement preprocessing technique effectively improves the performance of various respiratory sound classification models, including fine-tuning the AST and AST using the Patch-Mix strategy, on both the ICBHI and FABS datasets.

Table 5. Comparison of the classification performance on the International Conference in Biomedical and Health Informatics (ICBHI) database by fine-tuning the Audio Spectrogram Transformer [54].

	Accuracy (%)	Sensitivity (%)	Specificity (%)	ICBHI score (%)
Clean	70.65	64.88	75.67	70.27
Noisy	24.13	30.41	18.67	24.54
Noise injection	53.78	35.28	69.87	52.58
Audio enhancement	54.46	46.99	60.96	53.98

Tzeng et al

	Accuracy (%)	Sensitivity (%)	Specificity (%)	ICBHI score (%)
Clean	70.73	61.79	78.5	70.14
Noisy	29.05	35.45	23.48	29.46
Noise injection	58.02	23.9	87.69	55.8
Audio enhancement	58.55	40.98	73.83	57.4

Table 6. Comparison of the classification performance on the International Conference in Biomedical and Health Informatics (ICBHI) database using the Patch-Mix training strategy from the Audio Spectrogram Transformer pretrained weight [16].

 Table 7. Comparison of the classification performance on the Formosa Archive of Breath Sound database by fine-tuning the Audio Spectrogram Transformer [54].

	Accuracy (%)	Sensitivity (%)	Specificity (%)	ICBHI ^a score (%)
Clean	85.74	46.74	94.21	70.48
Noisy	83.03	36.96	93.03	65
Noise injection	83.8	31.52	95.16	63.34
Audio enhancement	80.89	50	87.6	68.8

^aICBHI: International Conference in Biomedical and Health Informatics.

Table 8. Comparison of the classification performance on the Formosa Archive of Breath Sound database using the Patch-Mix training strategy from the Audio Spectrogram Transformer pretrained weight [16].

	Accuracy (%)	Sensitivity (%)	Specificity (%)	ICBHI ^a score (%)
Clean	86.13	42.39	95.63	69.01
Noisy	82.15	29.35	93.62	61.49
Noise injection	82.44	44.57	90.67	67.62
Audio enhancement	75.17	57.61	78.98	68.3

^aICBHI: International Conference in Biomedical and Health Informatics.

Metric Discriminator

Given that the metric discriminator optimizes PESQ, a metric primarily used in the speech domain for speech quality, a potential mismatch problem may arise when applied to respiratory sound tasks. To explore this issue, we conducted ablation studies on CMGAN's discriminator, examining the conformer generator-only model, the conformer generative adversarial network without PESQ estimation discriminator (with normal discriminator), and the complete setup (with metric discriminator). As shown in Table 9, the addition of a metric discriminator improved overall accuracy, sensitivity, and ICBHI score. This outcome indicates a positive contribution of the metric discriminator on PESQ to respiratory sound classification.

 Table 9. Classification results of the convolution-augmented transformer-based metric generative adversarial network with different discriminator setups on the International Conference in Biomedical and Health Informatics (ICBHI) 2017 database.

Setup	Accuracy (%)	Sensitivity (%)	Specificity (%)	ICBHI score (%)
Generator only	65.81	58.21	72.42	65.32
With normal discriminator	66.19	55.61	75.39	65.5
With metric discriminator	66.72	62.28	70.58	66.43

Discussion

Principal Findings

This paper proposes a deep learning audio enhancement preprocessing pipeline for respiratory sound classification tasks. We also introduced a collection of clinical noise and a real-world respiratory sound database from the emergency department of the Hsin-Chu Branch at the NTUH. Our noise-robust method enhances model performance in noisy environments and

https://ai.jmir.org/2025/1/e67239

RenderX

provides physicians with improved audio recordings for manual assessment even under heavy noise conditions.

The experimental results indicated that audio enhancement significantly improved performance across all 3 types of noise commonly encountered during auscultation. Specifically, our approach achieved a 2.15% improvement (P<.001) over the conventional noise injection method on the ICBHI dataset and outperformed it by 2.79% (P<.001) on the FABS dataset. Moreover, time-frequency-domain enhancement techniques demonstrated superior performance for this task. Analyzing the

correlation between classification performance and speech quality metrics, we observed that CBAK and SSNR exhibited higher correlations with ICBHI scores. These metrics are strongly influenced by background noise but are unrelated to speech intelligibility, aligning with the experimental settings. In the physician validation study, enhanced recordings showed an 11.61% increase in sensitivity and a 14.22% improvement in classification confidence. A positive correlation was also observed between speech quality metrics and diagnostic confidence, highlighting the effectiveness of enhanced recordings in aiding physicians in detecting abnormal respiratory sounds. Our ablation study on various classification model structures revealed that audio enhancement preprocessing consistently improved performance. The findings showed enhanced sensitivity and higher ICBHI scores across both databases when tested with 2 state-of-the-art respiratory sound classification models. Furthermore, incorporating the metric discriminator PESQ was found to enhance downstream classification performance.

These findings validate the feasibility and effectiveness of integrating deep learning–based audio enhancement techniques into respiratory sound classification systems, addressing the critical challenge of noise robustness and paving the way for the development of reliable clinical decision support tools.

Limitations and Future Work

Despite the encouraging findings in this study, there is a need to explore the co-optimization of front-end audio enhancement and classification models. As most audio enhancement tasks primarily focus on speech, the evaluation metrics are not highly correlated with respiratory sounds, potentially leading to inefficient optimization. Addressing this issue is crucial for achieving better performance in respiratory sound classification in future work. Furthermore, future studies should incorporate other types of noise and more complex noise mixture strategies to enable the development of a more noise-robust respiratory sound classification model for real-world clinical use. By considering a diverse range of noise scenarios, the model can be better prepared to handle the variability and challenges encountered in actual clinical settings. In addition, we have to speed up the model inference by simplifying the model to make it suitable for real-time applications. At the same time, we must ensure that enhancement quality is maintained and critical respiratory sound characteristics are preserved. In our long-term future work, we aim to deploy this model in real clinical environments by integrating it into electronic stethoscopes. To ensure the method's generalizability, we plan to collect cross-site respiratory sound recordings from 100 patients across various clinical environments. Of these recordings, data from 80 patients will be used for training, whereas data from the remaining 20 patients will be reserved for testing as part of a validation process aligned with Food and Drug Administration requirements. This approach will help validate the model's performance and facilitate its adoption for practical use in clinical settings.

Conclusions

In this study, we investigated the impact of incorporating a deep learning–based audio enhancement module into automatic respiratory sound classification systems. Our results demonstrated that this approach significantly improved the system's robustness and clinical applicability, particularly in noisy environments. The enhanced audio not only improved classification performance on the ICBHI and FABS datasets but also increased diagnostic sensitivity and confidence among physicians. This study highlights the potential of audio enhancement as a critical component in developing reliable and trustworthy clinical decision support systems for respiratory sound analysis.

Acknowledgments

This research is funded by the National Science and Technology Council of Taiwan under grant 112-2320-B-002-044-MY3.

Conflicts of Interest

None declared.

Multimedia Appendix 1 Details of the technical setup used in this study. [DOCX File, 17 KB - ai v4i1e67239 app1.docx]

Multimedia Appendix 2 Hyperparameters for training enhancement and classification models. [DOCX File, 20 KB - ai_v4i1e67239_app2.docx]

References

- Bohadana A, Izbicki G, Kraman SS. Fundamentals of lung auscultation. N Engl J Med 2014 Feb 20;370(8):744-751. [doi: 10.1056/nejmra1302901]
- Arts L, Lim EH, van de Ven PM, Heunks L, Tuinman PR. The diagnostic accuracy of lung auscultation in adult patients with acute pulmonary pathologies: a meta-analysis. Sci Rep 2020 Apr 30;10(1):7347 [FREE Full text] [doi: 10.1038/s41598-020-64405-6] [Medline: 32355210]

- 3. Huang W, Tsai Y, Wei Y, Kuo P, Tao C, Cheng S, et al. Wheezing, a significant clinical phenotype of COPD: experience from the Taiwan Obstructive Lung Disease Study. Int J Chronic Obstr Pulm Dis 2015 Oct;10(1):2121-2126. [doi: 10.2147/copd.s92062]
- 4. Piirila P, Sovijarvi AR. Crackles: recording, analysis and clinical significance. Eur Respir J 1995 Dec 01;8(12):2139-2148. [doi: 10.1183/09031936.95.08122139]
- Chambres G, Hanna P, Desainte-Catherine M. Automatic detection of patient with respiratory diseases using lung sound analysis. In: Proceedings of the International Conference on Content-Based Multimedia Indexing. 2018 Presented at: CBMI 2018; September 4-6, 2018; La Rochelle, France. [doi: 10.1109/cbmi.2018.8516489]
- Kim Y, Hyon Y, Jung SS, Lee S, Yoo G, Chung C, et al. Respiratory sound classification for crackles, wheezes, and rhonchi in the clinical field using deep learning. Sci Rep 2021 Aug 25;11(1):17186 [FREE Full text] [doi: 10.1038/s41598-021-96724-7] [Medline: 34433880]
- 7. Wanasinghe T, Bandara S, Madusanka S, Meedeniya D, Bandara M, Díez ID. Lung sound classification with multi-feature integration utilizing lightweight CNN model. IEEE Access 2024;12:21262-21276. [doi: 10.1109/access.2024.3361943]
- Pessoa D, Petmezas G, Papageorgiou VE, Rocha BM, Stefanopoulos L, Kilintzis V. Pediatric respiratory sound classification using a dual input deep learning architecture. In: Proceedings of the IEEE Biomedical Circuits and Systems Conference. 2023 Presented at: BioCAS 2023; October 19-21, 2023; Toronto, ON. [doi: 10.1109/biocas58349.2023.10388733]
- Acharya J, Basu A. Deep neural network for respiratory sound classification in wearable devices enabled by patient specific model tuning. IEEE Trans Biomed Circuits Syst 2020 Jun;14(3):535-544. [doi: <u>10.1109/TBCAS.2020.2981172</u>] [Medline: <u>32191898</u>]
- Yu S, Ding Y, Qian K, Hu B, Li W, Schuller BW. A glance-and-gaze network for respiratory sound classification. In: Proceedings of the 2022 IEEE International Conference on Acoustics, Speech and Signal Processing. 2022 Presented at: ICASSP 2022; May 23-27, 2022; Singapore, Singapore. [doi: <u>10.1109/icassp43922.2022.9746053</u>]
- Zhao Z, Gong Z, Niu M, Ma J, Wang H, Zhang Z. Automatic respiratory sound classification via multi-branch temporal convolutional network. In: 2022 IEEE International Conference on Acoustics, Speech and Signal Processing. 2022 Presented at: ICASSP 2022; May 23-27, 2022; Singapore, Singapore. [doi: 10.1109/icassp43922.2022.9746182]
- 12. He W, Yan Y, Ren J, Bai R, Jiang X. Multi-view spectrogram transformer for respiratory sound classification. In: Proceedings of the 2024 IEEE International Conference on Acoustics, Speech and Signal Processing. 2024 Presented at: ICASSP 2024; April 14-19, 2024; Seoul, Republic of Korea. [doi: 10.1109/icassp48485.2024.10445825]
- 13. Zhang Y, Huang Q, Sun W, Chen F, Lin D, Chen F. Research on lung sound classification model based on dual-channel CNN-LSTM algorithm. Biomed Signal Process Control 2024 Aug;94:106257. [doi: 10.1016/j.bspc.2024.106257]
- Song W, Han J, Song H. Contrastive embeddind learning method for respiratory sound classification. In: Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing. 2021 Presented at: ICASSP 2021; June 6-11, 2021; Toronto, ON. [doi: 10.1109/icassp39728.2021.9414385]
- 15. Roy A, Satija U. AsthmaSCELNet: a lightweight supervised contrastive embedding learning framework for asthma classification using lung sounds. In: Proceedings of the 24th INTERSPEECH Conference. 2023 Presented at: INTERSPEECH 2023; August 20-24, 2023; Dublin, Ireland. [doi: 10.21437/interspeech.2023-428]
- Bae S, Kim JW, Cho WY, Baek H, Son S, Lee B, et al. Patch-mix contrastive learning with audio spectrogram transformer on respiratory sound classification. arXiv Preprint posted online on May 23, 2023 [FREE Full text] [doi: 10.21437/interspeech.2023-1426]
- 17. Moummad I, Farrugia N. Pretraining respiratory sound representations using metadata and contrastive learning. In: Proceedings of the IEEE Workshop on Applications of Signal Processing to Audio and Acoustics. 2023 Presented at: WASPAA 2023; October 22-25, 2023; New Paltz, NY. [doi: <u>10.1109/waspaa58266.2023.10248130</u>]
- Kochetov K, Putin E, Balashov M, Filchenkov A, Shalyto A. Noise masking recurrent neural network for respiratory sound classification. In: Proceedings of the 27th International Conference on Artificial Neural Networks and Machine Learning. 2018 Presented at: ICANN 2018; October 4-7, 2018; Rhodes, Greece. [doi: 10.1007/978-3-030-01424-7_21]
- 19. Ma Y, Xu X, Li Y. LungRN+NL: an improved adventitious lung sound classification using non-local block ResNet neural network with mixup data augmentation. In: Proceedings of the INTERSPEECH 2020. 2020 Presented at: INTERSPEECH 2020; October 25-29, 2020; Virtual Event, China. [doi: 10.21437/interspeech.2020-2487]
- 20. Wang Z, Wang Z. A domain transfer based data augmentation method for automated respiratory classification. In: Proceedings of the 2022 IEEE International Conference on Acoustics, Speech and Signal Processing. 2022 Presented at: ICASSP 2022; May 23-27, 2022; Singapore, Singapore. [doi: 10.1109/icassp43922.2022.9746941]
- 21. Nguyen T, Pernkopf F. Lung sound classification using co-tuning and stochastic normalization. IEEE Trans Biomed Eng 2022 Sep;69(9):2872-2882. [doi: 10.1109/tbme.2022.3156293]
- 22. Gairola S, Tom F, Kwatra N, Jain M. RespireNet: a deep neural network for accurately detecting abnormal lung sounds in limited data setting. In: Proceedings of the 43rd Annual International Conference of the IEEE Engineering in Medicine & Biology Society. 2021 Presented at: EMBC 2021; November 1-5, 2021; Mexico City, Mexico. [doi: 10.1109/embc46164.2021.9630091]

- Zhao X, Shao Y, Mai J, Yin A, Xu S. Respiratory sound classification based on BiGRU-attention network with XGBoost. In: Proceedings of the IEEE International Conference on Bioinformatics and Biomedicine. 2020 Presented at: BIBM 2020; December 16-19, 2020; Seoul, Republic of Korea. [doi: 10.1109/bibm49941.2020.9313506]
- Khullar D, Casalino LP, Qian Y, Lu Y, Chang E, Aneja S. Public vs physician views of liability for artificial intelligence in health care. J Am Med Inform Assoc 2021 Jul 14;28(7):1574-1577 [FREE Full text] [doi: 10.1093/jamia/ocab055] [Medline: <u>33871009</u>]
- 25. Diprose WK, Buist N, Hua N, Thurier Q, Shand G, Robinson R. Physician understanding, explainability, and trust in a hypothetical machine learning risk calculator. J Am Med Inform Assoc 2020 Apr 01;27(4):592-600 [FREE Full text] [doi: 10.1093/jamia/ocz229] [Medline: 32106285]
- Shim CS, Williams MHJ. Relationship of wheezing to the severity of obstruction in asthma. Arch Intern Med 1983 May;143(5):890-892. [Medline: <u>6679232</u>]
- 27. Kinoshita K, Ochiai T, Delcroix M, Nakatani T. Improving noise robust automatic speech recognition with single-channel time-domain enhancement network. In: Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing. 2020 Presented at: ICASSP 2020; May 4-8, 2020; Barcelona, Spain. [doi: 10.1109/icassp40776.2020.9053266]
- Pandey A, Liu C, Wang Y, Saraf Y. Dual application of speech enhancement for automatic speech recognition. In: Proceedings of the IEEE Spoken Language Technology Workshop. 2021 Presented at: SLT 2021; January 19-22, 2021; Shenzhen, China. [doi: 10.1109/slt48900.2021.9383624]
- 29. Lu YJ, Chang X, Li C, Zhang W, Cornell S, Ni Z, et al. ESPnet-SE++: speech enhancement for robust speech recognition, translation, and understanding. arXiv Preprint posted online on July 19, 2022 [FREE Full text] [doi: 10.21437/interspeech.2022-10727]
- 30. Kumar A, Florencio D. Speech enhancement in multiple-noise conditions using deep neural networks. arXiv Preprint posted online on May 9, 2016 [FREE Full text] [doi: 10.21437/interspeech.2016-88]
- Seltzer ML, Yu D, Wang Y. An investigation of deep neural networks for noise robust speech recognition. In: Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing. 2013 Presented at: ICASSP 2013; May 26-31, 2013; Vancouver, BC. [doi: 10.1109/icassp.2013.6639100]
- 32. Yin D, Luo C, Xiong Z, Zeng W. PHASEN: a phase-and-harmonics-aware speech enhancement network. Proc AAAI Conf Artif Intell 2020;34(05):9458-9465. [doi: 10.1609/aaai.v34i05.6489]
- Bagchi D, Plantinga P, Stiff A, Fosler-Lussier E. Spectral feature mapping with MIMIC loss for robust speech recognition. In: Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing. 2018 Presented at: ICASSP 2018; April 15-20, 2018; Calgary, AB. [doi: 10.1109/icassp.2018.8462622]
- Fu SW, Wang TW, Tsao Y, Lu X, Kawai H. End-to-end waveform utterance enhancement for direct evaluation metrics optimization by fully convolutional neural networks. IEEE/ACM Trans Audio Speech Lang Process 2018 Sep;26(9):1570-1584. [doi: 10.1109/taslp.2018.2821903]
- Koizumi Y, Niwa K, Hioka Y, Kobayashi K, Haneda Y. DNN-based source enhancement to increase objective sound quality assessment score. IEEE/ACM Trans Audio Speech Lang Process 2018 Oct;26(10):1780-1792. [doi: 10.1109/taslp.2018.2842156]
- 36. Fu SW, Liao CF, Tsao Y, Lin SD. MetricGAN: generative adversarial networks based black-box metric scores optimization for speech enhancement. arXiv Preprint posted online on May 13, 2019 [FREE Full text]
- 37. Rix AW, Beerends JG, Hollier MP, Hekstra AP. Perceptual evaluation of speech quality (PESQ)-a new method for speech quality assessment of telephone networks and codecs. In: Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing. 2001 Presented at: ICASSP 2001; May 07-11, 2001; Salt Lake City, UT. [doi: 10.1109/icassp.2001.941023]
- 38. Taal CH, Hendriks RC, Heusdens R, Jensen J. A short-time objective intelligibility measure for time-frequency weighted noisy speech. In: Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing. 2010 Presented at: ICASSP 2010; March 14-19, 2010; Dallas, TX. [doi: <u>10.1109/icassp.2010.5495701</u>]
- 39. Fu SW, Yu C, Hsieh TA, Plantinga P, Ravanelli M, Lu X, et al. MetricGAN+: an improved version of MetricGAN for speech enhancement. arXiv Preprint posted online on April 8, 2021 [FREE Full text] [doi: 10.21437/interspeech.2021-599]
- 40. Chen S, Wu Y, Chen Z, Wu J, Li J, Yoshioka T. Continuous speech separation with conformer. In: Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing. 2021 Presented at: ICASSP 2021; June 6-11, 2021; Toronto, ON. [doi: 10.1109/icassp39728.2021.9413423]
- 41. Gulati A, Qin J, Chiu CC, Parmar N, Zhang Y, Yu J, et al. Conformer: convolution-augmented transformer for speech recognition. arXiv Preprint posted online on May 16, 2020 [FREE Full text] [doi: 10.21437/interspeech.2020-3015]
- Zeineldeen M, Xu J, Lüscher C, Michel W, Gerstenberger A, Schlüter R. Conformer-based hybrid ASR system for switchboard dataset. In: Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). 2022 Presented at: ICASSP 2022; May 23-27, 2022; Singapore, Singapore. [doi: 10.1109/icassp43922.2022.9746377]
- 43. Cao R, Abdulatif S, Yang B. CMGAN: conformer-based metric GAN for speech enhancement. arXiv Preprint posted online on March 28, 2022 [FREE Full text] [doi: 10.21437/interspeech.2022-517]

- 44. Thiemann J, Ito N, Vincent E. The Diverse Environments Multi-channel Acoustic Noise Database (DEMAND): a database of multichannel environmental noise recordings. Proc Mtgs Acoust 2013 May 14;19:035081. [doi: 10.1121/1.4799597]
- 45. Valentini-Botinhao C. Noisy speech database for training speech enhancement algorithms and TTS models. University of Edinburgh. 2017. URL: <u>https://datashare.ed.ac.uk/handle/10283/2791</u> [accessed 2025-02-28]
- 46. Macartney C, Weyde T. Improved speech enhancement with the Wave-U-Net. arXiv Preprint posted online on November 27, 2018 [FREE Full text]
- 47. Pandey A, Wang D. Dual-path self-attention RNN for real-time speech enhancement. arXiv Preprint posted online on October 23, 2020 [FREE Full text]
- 48. Wang K, He B, Zhu WP. TSTNN: two-stage transformer based neural network for speech enhancement in the time domain. In: Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing. 2021 Presented at: ICASSP 2021; June 06-11, 2021; Toronto, ON. [doi: <u>10.1109/icassp39728.2021.9413740</u>]
- 49. Park HJ, Kang BH, Shin W, Kim JS, Han SW. MANNER: multi-view attention network for noise erasure. In: Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing. 2022 Presented at: ICASSP 2022; May 23-27, 2022; Singapore, Singapore. [doi: 10.1109/icassp43922.2022.9747120]
- 50. Deng J, Dong W, Socher R, Li LJ, Kai L, Fei-Fei L. ImageNet: a large-scale hierarchical image database. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2009 Presented at: CVPR 2009; June 20-25, 2009; Miami, FL. [doi: 10.1109/cvprw.2009.5206848]
- 51. Gemmeke JF, Ellis DP, Freedman D, Jansen A, Lawrence W, Moore RC. Audio set: an ontology and human-labeled dataset for audio events. In: Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing. 2017 Presented at: ICASSP 2017; March 5-9, 2017; New Orleans, LA. [doi: 10.1109/icassp.2017.7952261]
- 52. Piczak KJ. ESC: dataset for environmental sound classification. In: Proceedings of the 23rd ACM International Conference on Multimedia. 2015 Presented at: MM '15; October 26-30, 2015; Brisbane, Australia. [doi: <u>10.1145/2733373.2806390</u>]
- Kong Q, Cao Y, Iqbal T, Wang Y, Wang W, Plumbley MD. PANNs: large-scale pretrained audio neural networks for audio pattern recognition. IEEE/ACM Trans Audio Speech Lang Process 2020;28:2880-2894. [doi: <u>10.1109/taslp.2020.3030497</u>]
- 54. Gong Y, Chung YA, Glass J. AST: audio spectrogram transformer. arXiv Preprint posted online on April 5, 2021 [FREE Full text] [doi: 10.21437/interspeech.2021-698]
- 55. Gong Y, Lai CI, Chung YA, Glass J. SSAST: self-supervised audio spectrogram transformer. arXiv Preprint posted online on October 19, 2021 [FREE Full text] [doi: 10.1609/aaai.v36i10.21315]
- 56. Chang AY, Tzeng JT, Chen HY, Sung CW, Huang CH, Huang EP, et al. GaP-Aug: gamma patch-wise correction augmentation method for respiratory sound classification. In: Proceedings of the 2024 IEEE International Conference on Acoustics, Speech and Signal Processing. 2024 Presented at: ICASSP 2024; April 14-19, 2024; Seoul, Republic of Korea. [doi: 10.1109/icassp48485.2024.10447967]
- Kim JW, Toikkanen M, Bae S, Kim M, Jung HY. RepAugment: input-agnostic representation-level augmentation for respiratory sound classification. arXiv Preprint posted online on May 5, 2024 [FREE Full text] [doi: 10.1109/embc53108.2024.10782363]
- Emmanouilidou D, McCollum ED, Park DE, Elhilali M. Computerized lung sound screening for pediatric auscultation in noisy field environments. IEEE Trans Biomed Eng 2018 Jul;65(7):1564-1574 [FREE Full text] [doi: 10.1109/TBME.2017.2717280] [Medline: 28641244]
- 59. Rocha BM, Filos D, Mendes L, Vogiatzis I, Perantoni E, Kaimakamis E, et al. A respiratory sound database for the development of automated classification. In: Proceedings of the International Conference on Biomedical and Health Informatics. 2018 Presented at: ICBHI 2017; November 18-21, 2017; Thessaloniki, Greece. [doi: 10.1007/978-981-10-7419-6_6]
- 60. Zhang H, Cisse M, Dauphin YN, Lopez-Paz D. mixup: beyond empirical risk minimization. In: Proceedings of the International Conference on Learning Representations. 2018 Presented at: ICLR 2018; April 30-May 3, 2018; Vancouver, BC. [doi: 10.1007/978-981-19-9711-2_6]
- 61. Park DS, Chan W, Zhang Y, Chiu CC, Zoph B, Cubuk ED, et al. SpecAugment: a simple data augmentation method for automatic speech recognition. arXiv Preprint posted online on April 18, 2019 [FREE Full text] [doi: 10.21437/interspeech.2019-2680]
- Dong X, Shen J. Triplet loss in Siamese network for object tracking. In: Proceedings of the 15th European Conference on Computer Vision. 2018 Presented at: ECCV 2018; September 8-14, 2018; Munich, Germany. [doi: 10.1007/978-3-030-01261-8_28]

Abbreviations

AI: artificial intelligence
AST: Audio Spectrogram Transformer
CBAK: mean opinion score of background noise intrusiveness
CMGAN: convolution-augmented transformer–based metric generative adversarial network
CNN: convolutional neural network

https://ai.jmir.org/2025/1/e67239

CNN14: 14-layer convolutional neural network Conformer: convolution-augmented transformer FABS: Formosa Archive of Breath Sound ICBHI: International Conference in Biomedical and Health Informatics MOS: mean opinion score NTUH: National Taiwan University Hospital PESQ: perceptual evaluation of speech quality PHASEN: Phase-and-Harmonics–Aware Speech Enhancement Network SNR: signal-to-noise ratio SSNR: segmental signal-to-noise ratio STFT: short-time Fourier transform STOI: short-time objective intelligibility SVM: support vector machine

Edited by G Luo; submitted 06.10.24; peer-reviewed by T Abd El-Hafeez, D Meedeniya; comments to author 03.12.24; revised version received 26.01.25; accepted 27.01.25; published 13.03.25.

Please cite as:

Tzeng JT, Li JL, Chen HY, Huang CH, Chen CH, Fan CY, Huang EPC, Lee CC Improving the Robustness and Clinical Applicability of Automatic Respiratory Sound Classification Using Deep Learning–Based Audio Enhancement: Algorithm Development and Validation JMIR AI 2025;4:e67239 URL: https://ai.jmir.org/2025/1/e67239 doi:<u>10.2196/67239</u> PMID:

©Jing-Tong Tzeng, Jeng-Lin Li, Huan-Yu Chen, Chu-Hsiang Huang, Chi-Hsin Chen, Cheng-Yi Fan, Edward Pei-Chuan Huang, Chi-Chun Lee. Originally published in JMIR AI (https://ai.jmir.org), 13.03.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



A Conversational Agent Using Natural Language Processing for Postpartum Care for New Mothers: Development and Engagement Analysis

Kirstin Leitner¹, MD; Clare Cutri-French², MD; Abigail Mandel³, BA; Lori Christ⁴, MD; Nathaneal Koelper⁵, MPH; Meaghan McCabe⁶, MPH; Emily Seltzer⁷, MPH; Laura Scalise², MSN, BSN; James A Colbert⁸, MD, MBA; Anuja Dokras⁹, MD, MHCI, PhD; Roy Rosin¹⁰, MBA; Lisa Levine¹¹, MD

¹Department of Obstetrics and Gynecology, University of Pennsylvania, 3701 Market Street, 3rd Floor, Philadelphia, PA, United States ¹⁰Penn Medicine, Philadelphia, PA, United States

²Hospital of the University of Pennsylvania, Philadelphia, PA, United States

³Medical University of South Carolina, Charleston, SC, United States

⁴Intensive Care Nursery, Department of Pediatrics, University of Pennsylvania, Philadelphia, PA, United States

⁵School of Medicine, University of Pennsylvania, Philadelphia, PA, United States

⁶Maternal Fetal Medicine Research Center, School of Medicine, University of Pennsylvania, Philadelphia, PA, United States

⁷Penn Medicine Center for Health Care Transformation and Innovation, Philadelphia, PA, United States

⁸Memora Health, San Francisco, CA, United States

⁹Division of Reproductive Endocrinology and Infertility, Department of Obstetrics and Gynecology, University of Pennsylvania, Philadelphia, PA, United States

Corresponding Author:

Kirstin Leitner, MD

Department of Obstetrics and Gynecology, University of Pennsylvania, 3701 Market Street, 3rd Floor, Philadelphia, PA, United States

Abstract

Background: The "fourth trimester," or postpartum time period, remains a critical phase of pregnancy that significantly impacts parents and newborns. Care poses challenges due to complex individual needs as well as low attendance rates at routine appointments. A comprehensive technological solution could provide a holistic and equitable solution to meet care goals.

Objective: This paper describes the development of patient engagement data with a novel postpartum conversational agent that uses natural language processing to support patients post partum.

Methods: We report on the development of a postpartum conversational agent from concept to usable product as well as the patient engagement with this technology. Content for the program was developed using patient- and provider-based input and clinical algorithms. Our program offered 2-way communication to patients and details on physical recovery, lactation support, infant care, and warning signs for problems. This was iterated upon by our core clinical team and an external expert clinical panel before being tested on patients. Patients eligible for discharge around 24 hours after delivery who had delivered a singleton full-term infant vaginally were offered use of the program. Patient demographics, accuracy, and patient engagement were collected over the first 6 months of use.

Results: A total of 290 patients used our conversational agent over the first 6 months, of which 112 (38.6%) were first time parents and 162 (56%) were Black. In total, 286 (98.6%) patients interacted with the platform at least once, 271 patients (93.4%) completed at least one survey, and 151 (52%) patients asked a question. First time parents and those breastfeeding their infants had higher rates of engagement overall. Black patients were more likely to promote the program than White patients (P=.047). The overall accuracy of the conversational agent during the first 6 months was 77%.

Conclusions: It is possible to develop a comprehensive, automated postpartum conversational agent. The use of such a technology to support patients postdischarge appears to be acceptable with very high engagement and patient satisfaction.

(JMIR AI 2025;4:e58454) doi:10.2196/58454

¹¹Division of Maternal Fetal Medicine, Department of Obstetrics and Gynecology, University of Pennsylvania, Philadelphia, PA, United States

KEYWORDS

Leitner et al

conversational agent; postpartum care; text messaging; postpartum; natural language processing; pregnancy; parents; newborns; development; patient engagement; physical recovery; infant; infant care; survey; breastfeeding; support; patient support; patient satisfaction

Introduction

The "fourth trimester," or postpartum time period, is often a forgotten "trimester" of pregnancy, yet plays a critical role in parental and newborn well-being. While undergoing numerous physiologic and emotional changes following birth, patients are also susceptible to complications such as infection, thrombosis, and hypertensive disorders as well as the new onset or exacerbation of mental health disorders [1,2]. The potential for medical complications post partum is of particular concern as over one-half of pregnancy related deaths occur after the birth of the infant [3,4]. These deaths also disproportionately affect Black women with maternal mortality rates nearly 3-times that of non-Hispanic White women [5]. The American College of Obstetricians and Gynecologists (ACOG) recommends that care during the postpartum period should be an "ongoing process" rather than the traditional 1-time postpartum visit [6]. A study evaluating the clinical features of postpartum presentation for emergency care indicated that while rates are overall low around 5%, most visits occur within the first 2 weeks post partum and are more likely to occur in Black patients [7]. Yet, even when follow up is recommended, nearly 50% of patients in the United States do not attend their routine postpartum appointment and adding additional clinic visits to increase access is impractical and impossible for both patients and clinicians [8].

This gap between patient needs, clinical recommendations and reality of health care access presents a significant challenge to patients and practicing providers. Innovative methods of identifying needs and providing ongoing care for the postpartum patient are needed without added burden to already overextended providers. A wide range of SMS text messaging health care interventions have been developed and trialed with varied success [9]. Within the realm of postpartum care these innovations have largely focused on specific individual conditions regarding postpartum recovery such as breastfeeding [10-12], blood pressure monitoring [13], and weight loss [14-17]. While many of these interventions have shown great promise in improving compliance with care and reducing health care disparities [13], there are limited comprehensive technologic interventions to support patients holistically during the fourth trimester. A technology-based solution has the potential to meet ACOG's goals of continued contact and comprehensive postpartum care for patients. In this manuscript we describe the development of a novel comprehensive postpartum conversational agent, which uses natural language processing (NLP) to provide anticipatory guidance and respond to patients' questions in real time. We also describe patient engagement and satisfaction with this novel technology.

Methods

Program Design and Content Development

We sought to create a comprehensive technology-based postpartum support program, "Healing at Home," which would provide 24/7 support to individuals through the use of SMS text messages for 6 weeks post partum. Content included anticipatory guidance regarding physical recovery, infant care and feeding, clinical algorithms to respond to urgent needs and postpartum depression screening through the Edinburgh Postnatal Depression Screen (EPDS). The EPDS is a clinically validated 10-question survey that is considered the standard for screening patients for postpartum depression. We postulated that a 24/7 SMS text message-based holistic support would result in increased engagement of patients and allow providers to identify symptoms before they resulted in complications. Automation of messaging and responses, alongside the ability to focus attention efficiently on patients with demonstrated higher needs, could also minimize care team workload. Patients could be quickly connected to their care team and receive in-the-moment answers to their concerns.

We used a multipronged approach to optimize discharge planning and maintain postpartum connection for patients delivering at the Hospital of the University of Pennsylvania (HUP), described in detail by Gaulton et al [18]. We called this program of optimized discharge planning and increased postpartum support "Healing at Home." Pertinent to the innovation described here, this preintervention pilot leveraged a "fake back end" SMS text message-based support during business hours (8 AM-5 PM) for patients for the first 6 weeks post partum. During this preintervention phase described by Gaulton et al [18], 90 patients were enrolled and encouraged to text their questions to the team. Text messages were monitored by nonclinical as well as clinical staff viewing and responding to patients. The team used a clinical reference guide, which was elaborated on throughout the pilot, outlining responses to frequently asked questions. While this method was effective at connecting with patients, it required significant time monitoring messages and responding to patients. Over 2000 text messages were exchanged with this cohort of 90 patients. In addition, we identified highly complex and individual needs ranging from inquiries about physical recovery specific to delivery mode (vaginal vs cesarean) to care of newborns (diapering and umbilical cord care) and infant feeding difficulties (pain with breastfeeding, difficulty pumping, and preparing formula). This complexity led us to conclude that a "simple" algorithmic approach was unlikely to be successful in providing this population with the holistic support required.

Conversational agents are designed to simulate conversation with human users and have become nearly ubiquitous in business, but their development within health care has been slow. Given the complexity and individualized needs of the

```
XSL•FO
RenderX
```

postpartum patient we postulated that a conversational agent using NLP might be a good solution and be acceptable to this population. We envisioned a 24/7 available SMS text message—based support program that interpreted patients' postpartum concerns, responded in real sentences and could also alert clinicians in real time when appropriate.

Conversational Agent Development

We partnered with Memora Health to undertake a 4-step process to develop the conversational agent using NLP to interact with patients in a HIPAA (Health Insurance Portability and Accountability Act)-compliant manner. Unlike a basic chatbot which uses rigid decision trees to respond to people, this type of conversational agent leverages NLP to understand and interpret patient messages, providing appropriate responses, leading to a conversational experience. First, a frequently asked question bank was used to generate accurate mapping of questions to the appropriate responses. Second, surveys (standardized conversation templates designed to collect patient data) were created by patients' clinical characteristics (ie, breastmilk vs formula fed, Figure 1A). Third, creation of anticipatory guidance specific to patient clinical characteristic was planned. Finally, algorithms for potentially acute clinical concerns were designed and layered onto the program. Throughout this process we incorporated personal touches into responses, such as patients or infants' names and worked to develop a consistent and empathetic tone.

Figure 1. Layering of patient clinical characteristics (1A), example of clinical algorithm with symptom triage for lower extremity edema (1B), Memora Health patient dashboard (1C), comprehensive list of clinical symptoms for which algorithms were developed (1D). BP: blood pressure; EPDS: Edinburgh Postnatal Depression Screen; HTN: hypertension.



The frequently asked questions were generated from both patients (through our 90-patient preintervention pilot) and clinicians (obstetrics, neonatology, lactation, and social workers on our mother-baby unit). Clinicians were encouraged to "think like a patient" and ask questions they had either received or conceived as important. An example question might be whether nipple pain with feeding is normal. Topics from both patients and providers were categorized (obstetrics, neonatology, or lactation), reviewed by our team for accurate clinical content and then made available in a frequently asked question bank.

Surveys, that is, structured questions designed to collect patient data, were used and incorporated into this program including validated clinical questionnaires such as the EPDS and net promoter score (NPS). The NPS is a customer satisfaction and loyalty metric used to measure the likelihood that customers will recommend a product, service, or experience to others. People are asked to provide a score from 0 to 10. Promoters are those who score a program 9 or 10, passives score of 7 or 8, and detractors 6 or less. The NPS is calculated by subtracting the percentage of promotors from the percentage of detractors. Scores of 50 or greater indicates exceptional loyalty. The NPS was collected from patients during week two of the program.

https://ai.jmir.org/2025/1/e58454

We also developed multiple surveys with branching logic to dynamically respond to patients around topics such as infant feeding and the importance of attending scheduled appointments. Surveys were added to the program at scheduled times according to clinical needs.

Next, structured anticipatory guidance customized to patient characteristics (Figure 1) was generated by our clinical team such that patients with certain characteristics received appropriate educational materials at the right time (when they needed it and not before). Examples of customized anticipatory guidance include information on the volume of feeds by feeding method (breast vs formula).

Finally, we created a series of algorithms designed to address specific clinical scenarios outside of the conversational agent. For example, when asking about lower extremity edema, it cannot be assumed that this is normal swelling post partum, so triage regarding possible signs of venous thromboembolism is essential for providing safety (Figure 1). We also developed a "latch" algorithm, which was designed to address some common concerns in early lactation and difficulty with breastfeeding.

Leitner et al

JMIR AI

All content, including ad hoc content, surveys, anticipatory guidance, and clinical algorithms, was reviewed not only by our internal clinical leads but also by an expert clinical panel. This external panel, composed of a group of clinician leaders at our institution uninvolved with the design or development, provided us with insight and perspective on the content, our approach, and identification of perceived patient risks in development.

While we aimed to minimize unnecessary escalation, we erred on the side of caution with standard language and responses. For example, if the conversational agent is unable to answer a patient's concern about their infant, they would receive the following response: "It sounds like this is something that infant name's doctor can help with. Please call their office at 555-555-555." We also instructed patients to use the phrase "TEXT ME" to indicate that their concern was not addressed, alerting the team to review the conversation and intervene. This was primarily accomplished through messaging patients directly in the Memora Health dashboard. Best practices when interacting with the conversational agent, such as using single-sentence questions and rephrasing questions to improve responses, were shared with the patients through flyers and a short educational video at the time of the program start.

Conversational Agent Testing

Testing of the conversational agent required multiple phases, including internal and external testing. First, we tested this program internally by asking our own team members to ask questions that they would imagine patients might ask, attempting to cover a wide range of common clinical scenarios. Suggested scenarios included concerns around the color of infant stools, pain management, etc. Accuracy of the responses was also improved through a rapid-fire test using Mechanical Turk as described in detail by Lin et al [19]. Once these tests had been completed, we tested the program with providers external to our team and a small set of patients.

First, we recruited 23 providers from obstetrics, neonatology, nursing, and lactation who had not been involved in the design of the program to use the conversational agent as if they were a patient (they were assigned a patient characteristic such as feeding method for testing). In the second phase, we recruited 37 patients from the HUP postpartum unit to use the conversational agent in their own recovery. These patients were selected to be representative of our population with examples of demographics including race, parity, marital status, insurance, and age. During this initial patient testing phase, monitoring of the platform occurred once daily at a minimum by our clinical team.

Chatbot Enrollment and Clinical Monitoring

Clinical criteria for patient participation in Healing at Home was determined by our clinical team and expert panel at the start of this program. Program participants were patients who were planned for discharge around 24 hours of after birth who had an uncomplicated vaginal delivery at term of a singleton infant; full exclusion criteria are outlined in Textbox 1. These clinical criteria were selected as clinicians felt most comfortable with a new technology being used by a group of patients less likely to experience postpartum complications. The data presented here regarding patient engagement includes 290 patients who met these clinical characteristics. Upon discharge, patients were enrolled in the texting platform by the postpartum nurse and verbally consented. Our first nontest patient was enrolled March 6, 2020. Once the program was live, we took a data-driven approach to improve the patient experience and the conversational agent itself. For example, when we discovered that many patients were asking about their infant's umbilical cord care during the first week, we programmed a message to be proactively sent at that time.



Textbox 1. Clinical exclusion criteria to enrollment in the "Healing at Home" program (planned discharge around 24 h post birth).

Maternal exclusion criteria:

- Age<18
- Cesarean delivery
- Gestational age<37 weeks and 0 days
- Multiple gestation
- Blood loss>1000cc
- 3rd or 4th degree perineal laceration
- Preexisting diabetes mellitus or gestational diabetes on medication
- Preeclampsia with severe features
- Chronic hypertension on medication

Infant exclusion criteria:

- Intensive care nursery admission
- Birth weight<2500 g
- Direct antiglobulin test positive
- 24-hour glucose<50 mg/dL
- 24-hour bilirubin>6 mg/dL
- No void at 24 hours of life
- Elevated sepsis risk score (≥ 0.7)
- Weight loss >7% at 24 hours of life
- <6 feeds in first 24 hours

Other exclusion criteria:

- No access to texting
- Non-English speaking primary language
- Adoption case
- Department of Human Services involvement
- Patient opt out
- Provider opt out
- Latch score ≤6

While interactions are designed to be automated, it was assumed that unanswered questions or clinical concerns would occur. Clinical teams were assigned to respond to these escalations which were assigned an acuity level by our team as appropriate. Some alerts were received through email, while more critical alerts were sent by SMS text messages if deemed to be emergent. While patients interacted with the conversational agent by text message, monitoring of the program occurred through the Memora Health dashboard, where patients could be viewed, chat history seen, and patient characteristics could be edited as appropriate (for example, changing the feeding method from breast milk to formula). On the dashboard, clinicians could directly message with patients in addition to traditional methods of patient contact by phone (Figure 1C).

Collection of Patient Engagement Data

A complete summary of clinical outcomes regarding users of this platform is beyond the scope of this paper, but we present users. This study was approved as a Quality Improvement Project by the University of Pennsylvania Institutional Review Board. Demographic data including age, parity, race, feeding method and insurance were collected from our electronic medical record. Patient engagement metrics and chatbot accuracy were extracted from individual patient text messages reviewed manually by our investigators. Classical descriptive statistics were generated using mean and SD for continuous variables and frequency and percentage for categorical variables. To measure the differences between demographically different groups, the chi-square test was used for categorical variables, and *t* test and ANOVA for continuous variables. Spearman rank correlations were used to describe relationships between continuous variables. All statistical analysis was performed using Stata (StataCorp).

demographic and engagement data from the first 6 months of

https://ai.jmir.org/2025/1/e58454

Engagement with the chatbot was measured by the number of total texts, number of questions asked, and survey response rates. Questions were classified by content category (maternal, baby, lactation, and social work) and by whether the question was prompted or unprompted. Prompted questions were defined as questions related to a previous message (ie, "Remember you can ask me questions about your health or baby") or directly following another interaction. Messages that were unrelated or temporally distant (>3 h since last message) were defined as unprompted. Binary data (asked vs did not ask question) and the total number of questions were recorded. Reworded questions did not count towards the total number of questions asked. Interactions between patients by a pleasantry (emoji, "ok," "thanks!") were recorded in a binary fashion. Patient satisfaction was collected through the NPS. Chatbot accuracy was measured by the percentage of correct answers per patient, excluding ignored interactions and no content situations. No qualitative interviews were conducted.

Ethical Considerations

This study was approved as a Quality Improvement Project by the University of Pennsylvania Institutional Review Board.

Results

A total of 290 patients used our chatbot over the first 6 months of use from March to August 2020. The average patient age was 28.8 (SD 5.47) years, 112 out of 290 (38.6%) patients were first time parents, 134 (46%) had private insurance, and 163 (56%) were Black (Table 1). This distribution is representative of the population at our large urban academic medical center. Of these 290 patients, 286 (98.6%) responded to the platform at least once, with 271 (93.4%) completing at least one survey, 151 (52%) asking a question (prompted or unprompted), and 162 (55.9%) interacting by a pleasantry. All patients were sent the EPDS at least 3 times over 6 weeks with 128 (44%) patients completing at least one EPDS. In addition, 93 (32%) of patients completed an NPS with an overall NPS score of 34.

Demographic characteristics		n (%)	
Age			
	<20	14 (4.8)	
	20 - 29	139 (48.1)	
	30 - 39	130 (44.7)	
	≥40	7 (2.4)	
Parity			
	0	112 (38.6)	
	1	94 (32.4)	
	2	51 (17.6)	
	≥3	33 (11.7)	
Race and ethnicity			
	Asian	13 (4.5)	
	Black	163 (56)	
	East Indian	3 (1)	
	Hispanic Latino/Black	3 (1)	
	Hispanic Latino/White	10 (3.4)	
	Other	11 (3.8)	
	Patient declined	2 (0.7)	
	Unknown	2 (0.7)	
	White	83 (28.6)	
Insurance			
	Private	134 (46)	
	Medicaid	156 (54)	
Feeding type			
	Breast	194 (66.9)	
	Formula	43 (14.8)	
	Both	53 (18.3)	

Black patients were statistically more likely to promote the program (score 9 or 10 on a scale of 0 - 10; P=.047) with an NPS score of 53 compared a NPS score of 18 for White patients. Engagement through survey completion and questions asked is shown in Table 2. White patients completed more surveys than Black patients (10.64 vs 6.64; P<.001), but there was no significant difference in the number of questions asked. Patients with private insurance completed more surveys than those with Medicaid (9.83 vs 6.43; P<.001), however, again no difference in questions asked. Patients feeding their infant breastmilk were more likely to ask questions (8.92 vs 4.65; P<.001) and complete

surveys (10 vs 4; P<.001). There were a total of 32 "super users" (patients who asked more than 4 questions) of which 25 (78%) were non-White, with 19 (59.3%) of these "super users" being Black and exclusively breastfeeding (although only 87 out of the 290 patients in this cohort (30%) were both Black and exclusively breastfeeding; see Figure 2A). Patients with lower parity, that is, patients who had experienced their first birth, asked more questions (P<.001) and completed more surveys (P<.001) than patients who had already birthed 1 or more children. Each unit increase in parity decreased the total number of questions by 0.36 (Figure 2B).

Table . Engagement data by patient demographics.

Demographic	Black race	White race	Other race	P value	Private insurance	Medicaid insurance	P value	Breast- milk on- ly	Formula only	Both	P value
Median total ques- tions (IQR)	0 (0 - 2)	1 (0 - 2)	1 (0 - 2)	.64	1 (0 - 2)	0 (0 - 2)	.26	1 (0 - 2)	0 (0 - 0)	1 (0 - 2)	<.001
Total questions ^a , n (%)				.89			.26				<.001
0	82 (50)	39 (47)	18 (41)		60 (45)	79 (51)		83 (43)	33 (77)	23 (43)	
1	37 (23)	20 (24)	11 (25)		32 (24)	36 (23)		46 (24)	6 (14)	16 (30)	
2	18 (11)	7 (8)	6 (14)		12 (9)	19 (12)		42 (10)	1 (2)	10 (19)	
3+	26 (16)	17 (20)	9 (20)		30 (22)	22 (14)		45 (23)	3 (7)	4 (8)	
Total questions ^b , n (%)				.31			.94				.03
<4	141 (87)	76 (92)	41 (93)		119 (89)	139 (89)		166 (86)	42 (98)	50 (94)	
4+	22 (13)	7 (8)	3 (7)		15 (11)	17 (11)		28 (14)	1 (2)	3 (6)	
Completed surveys, median (IQR)	6 (3-10)	12 (7-15)	8 (4-12)	<.001	10 (6-14)	6 (3-10)	<.001	10 (5-13)	4 (1-7)	7 (4-9)	<.001

^aTotal number of patients with 0,1,2, and 3+ questions.

^bTotal number of patients with <4 or 4+ questions.

Figure 2. Parity versus number of questions asked (2A) and completed surveys (2B). Dots represent individual patients and were jittered to minimize overplotting.





conversational agent with an overall chatbot accuracy of 77%

(correctly answered questions/correctly answered plus

incorrectly answered questions) with no difference in accuracy

by parity, race, or insurance status. The additional 97 (27%)

questions were not answered as they occurred concurrently with

a survey (58/97) or had no developed content (39/97), these are

A total of 422 questions were asked by patients, 177 (42%) were prompted and 244 (58%) were unprompted. In addition, 211 (50%) questions of patient questions were related to infant concerns, 135 (32%) to maternal health, 72 (17%) to lactation concerns, and 4 (1%) to social work concerns. Approximately, 325 (73%) of all patient questions could be answered by the

https://ai.jmir.org/2025/1/e58454

excluded from the overall accuracy rate reported here. As this was a fluid development process, responses were created to questions that were missing content for future users.

Discussion

Here we report on the development of a comprehensive postpartum conversational agent that leverages NLP to support patients during the "fourth trimester." Satisfaction from patients using this texting program was the highest among Black patients with high rates of engagement by all users regardless of race. The maternal health crisis is real in the United States and impacts Black patients at significantly higher rates [5,20]. Contrary to the current design of prenatal care where emphasis is placed on the pregnant patient and not the postpartum patient, we aimed to design a scalable approach to support patients during the fourth trimester by SMS text messages using augmented intelligence and NLP through a novel postpartum conversational agent. Its holistic rather than problem-based design gives this technology the potential for scalability beyond what previous models or interventions have been able to achieve. We have shown here that patient engagement is high (>98% interaction rate and >93% survey completion rate) and that patient satisfaction in Black patients is high, with Black patients were more likely to promote this program than White patients (P=.047). As we look to solutions for the maternal health crisis, we must keep a critical eye on the impact that racism has on health and find solutions that specifically target these disproportionately impacted populations.

A confounding aspect to the engagement data presented here is the time during which we collected data: March-August 2020. Our go-live date for the program coincided very closely to the start of the shut-down related to the COVID-19 pandemic with local restrictions going into effect in the second week of patient use with this platform. The influence of this may have had a significant impact on patients' experience with this platform and health care in general, especially in this cohort of patients who all completed the program by the end of 2020. Yet, we have continued to use this technology at our institution and will be able to determine whether and how moving out of the global pandemic impacts user engagement and patient satisfaction. Given the iterative nature of development, additional limitations include that significant improvements that were made over time (NLP improvement, mapping improvement) may not be reflected here (such as accuracy). Engagement by feeding method is confounded by race, with Black patients less likely to be exclusively breastfeeding, and NPS results are confounded by low response rates (<30%). An additional limitation of our engagement data presented here is that no qualitative interview of patients was performed. Without qualitative feedback from patients it is hard to draw any conclusions about the reason for NPS score disparity by race.

There are several outstanding questions that we have and plan to address in future work with this technology. First, we plan to report on the clinical outcomes on a larger cohort of patients with a specific focus on health care use and postpartum health goals such as visit attendance rates, ED and readmission rates as well as breastfeeding and contraception acceptance. In terms of health care use, we hope to gather data on number of phone calls to the office as well as amount of time per patient needed to manage concerns. In addition to these clinical and health care use outcomes, a key component to successful and broad implementation of such a program is intentional learning from the patients and providers who use this program. This allows for continued improvements and iterations on the program. We hope that future qualitative work with both providers and patients will help to elucidate barriers and facilitators to such a program. Within the framework of our layered program design, we have purposefully designed for flexibility in who, for instance, is asked to manage alerts or what content to include in the program to fit different hospital systems and teams.

We continue to use and expand upon this program at our own institution with to date over 1800 patients using this postpartum SMS text message-based support. Beyond our own application, we very much hope that the framework for the development of a comprehensive health care conversational agent (Figure 3) can help other clinical teams in their development, regardless of the specific clinical need addressed.



Figure 3. Conceptual framework for development of health care conversational agent. NPS: net promotor score.



Conflicts of Interest

AM is a former employee of Memora Health. JC is an employee of Memora Health. The remaining authors declare no conflict of interest.

References

- Hamilton N, Stevens N, Lillis T, Adams N. The fourth trimester: toward improved postpartum health and healthcare of mothers and their families in the United States. J Behav Med 2018 Oct;41(5):571-576. [doi: <u>10.1007/s10865-018-9969-9</u>] [Medline: <u>30302656</u>]
- 2. Paladine HL, Blenning CE, Strangas Y. Postpartum care: an approach to the fourth trimester. Am Fam Physician 2019 Oct 15;100(8):485-491. [Medline: <u>31613576</u>]
- Kassebaum NJ, Bertozzi-Villa A, Coggeshall MS, et al. Global, regional, and national levels and causes of maternal mortality during 1990-2013: a systematic analysis for the Global Burden of Disease Study 2013. Lancet 2014 Sep 13;384(9947):980-1004. [doi: 10.1016/S0140-6736(14)60696-6] [Medline: 24797575]
- Petersen EE, Davis NL, Goodman D, et al. Vital signs: Pregnancy-related deaths, United States, 2011-2015, and strategies for prevention, 13 States, 2013-2017. MMWR Morb Mortal Wkly Rep 2019 May 10;68(18):423-429. [doi: 10.15585/mmwr.mm6818e1] [Medline: <u>31071074</u>]
- 5. Hoyert DL. Maternal mortality rates in the United States, 2021. : National Center for Health Statistics; 2023 URL: <u>https://www.cdc.gov/nchs/data/hestat/maternal-mortality/2021/maternal-mortality-rates-2021.htm</u> [accessed 2025-03-07]
- 6. McKinney J, Keyser L, Clinton S, Pagliano C. ACOG Committee opinion no. 736: optimizing postpartum care. Obstet Gynecol 2018 Sep;132(3):784-785. [doi: 10.1097/AOG.0000000002849] [Medline: 30134408]
- 7. Rodriguez AN, Patel S, Macias D, Morgan J, Kraus A, Spong CY. Timing of emergency postpartum hospital visits in the fourth trimester. Am J Perinatol 2021 Mar;38(4):319-325. [doi: 10.1055/s-0040-1716842] [Medline: 32992354]
- 8. Bennett WL, Chang HY, Levine DM, et al. Utilization of primary and obstetric care after medically complicated pregnancies: an analysis of medical claims data. J Gen Intern Med 2014 Apr;29(4):636-645. [doi: 10.1007/s11606-013-2744-2] [Medline: 24474651]
- 9. Milne-Ives M, de Cock C, Lim E, et al. The effectiveness of artificial intelligence conversational agents in health care: systematic review. J Med Internet Res 2020 Oct 22;22(10):e20346. [doi: <u>10.2196/20346</u>] [Medline: <u>33090118</u>]
- Jiang H, Li M, Wen LM, et al. Effect of short message service on infant feeding practice: findings from a community-based study in Shanghai, China. JAMA Pediatr 2014 May;168(5):471-478. [doi: <u>10.1001/jamapediatrics.2014.58</u>] [Medline: <u>24639004</u>]
- Ahmed AH, Roumani AM, Szucs K, Zhang L, King D. The effect of interactive web-based monitoring on breastfeeding exclusivity, intensity, and duration in healthy, term infants after hospital discharge. J Obstet Gynecol Neonatal Nurs 2016;45(2):143-154. [doi: 10.1016/j.jogn.2015.12.001] [Medline: 26779838]

- 12. Gallegos D, Russell-Bennett R, Previte J, Parkinson J. Can a text message a week improve breastfeeding? BMC Pregnancy Childbirth 2014 Nov 6;14:374. [doi: 10.1186/s12884-014-0374-2] [Medline: 25369808]
- Hirshberg A, Downes K, Srinivas S. Comparing standard office-based follow-up with text-based remote monitoring in the management of postpartum hypertension: a randomised clinical trial. BMJ Qual Saf 2018 Nov;27(11):871-877. [doi: <u>10.1136/bmjqs-2018-007837</u>] [Medline: <u>29703800</u>]
- 14. Gilmore LA, Klempel MC, Martin CK, et al. Personalized mobile health intervention for health and weight loss in postpartum women receiving women, infants, and children benefit: a randomized controlled pilot study. J Womens Health (Larchmt) 2017 Jul;26(7):719-727. [doi: 10.1089/jwh.2016.5947] [Medline: 28338403]
- 15. Herring SJ, Cruice JF, Bennett GG, et al. Intervening during and after pregnancy to prevent weight retention among African American women. Prev Med Rep 2017 Sep;7:119-123. [doi: 10.1016/j.pmedr.2017.05.015] [Medline: 28660118]
- 16. van der Pligt P, Ball K, Hesketh KD, et al. A pilot intervention to reduce postpartum weight retention and central adiposity in first-time mothers: results from the mums OnLiNE (Online, Lifestyle, Nutrition & Exercise) study. J Hum Nutr Diet 2018 Jun;31(3):314-328. [doi: 10.1111/jhn.12521] [Medline: 29034545]
- 17. Phelan S, Hagobian T, Brannen A, et al. Effect of an internet-based program on weight loss for low-income postpartum women: a randomized clinical trial. JAMA 2017 Jun 20;317(23):2381-2391. [doi: 10.1001/jama.2017.7119] [Medline: 28632867]
- 18. Gaulton JS, Leitner K, Hahn L, et al. Healing at home: applying innovation principles to redesign and optimise postpartum care. BMJ Innov 2022 Jan;8(1):37-41. [doi: 10.1136/bmjinnov-2021-000791]
- 19. Lin J, Joseph T, Parga-Belinkie JJ, et al. Development of a practical training method for a healthcare artificial intelligence (AI) chatbot. BMJ Innov 2021 Apr;7(2):441-444. [doi: <u>10.1136/bmjinnov-2020-000530</u>]
- 20. Callaghan WM. Overview of maternal mortality in the United States. Semin Perinatol 2012 Feb;36(1):2-6. [doi: 10.1053/j.semperi.2011.09.002] [Medline: 22280858]

Abbreviations

ACOG: American College of Obstetricians and Gynecologists EPDS: Edinburgh Postnatal Depression Screen HIPAA: Health Insurance Portability and Accountability Act HUP: Hospital of the University of Pennsylvania NLP: natural language processing NPS: net promotor score

Edited by Z Yin; submitted 01.04.24; peer-reviewed by A Lakshmanan, L Narbey; revised version received 01.10.24; accepted 24.10.24; published 22.04.25.

<u>Please cite as:</u> Leitner K, Cutri-French C, Mandel A, Christ L, Koelper N, McCabe M, Seltzer E, Scalise L, Colbert JA, Dokras A, Rosin R, Levine L A Conversational Agent Using Natural Language Processing for Postpartum Care for New Mothers: Development and Engagement Analysis JMIR AI 2025;4:e58454 URL: <u>https://ai.jmir.org/2025/1/e58454</u> doi:<u>10.2196/58454</u>

© Kirstin Leitner, Clare Cutri-French, Abigail Mandel, Lori Christ, Nathaneal Koelper, Meaghan McCabe, Emily Seltzer, Laura Scalise, James A Colbert, Anuja Dokras, Roy Rosin, Lisa Levine. Originally published in JMIR AI (https://ai.jmir.org), 22.4.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



The Diagnostic Performance of Large Language Models and Oral Medicine Consultants for Identifying Oral Lesions in Text-Based Clinical Scenarios: Prospective Comparative Study

Sarah AlFarabi Ali^{1*}, PhD; Hebah AlDehlawi^{1*}, PhD; Ahoud Jazzar^{1*}, PhD; Heba Ashi^{2*}, PhD; Nihal Esam Abuzinadah^{3*}, PhD; Mohammad AlOtaibi⁴, BDS; Abdulrahman Algarni^{4*}, BDS; Hazzaa Alqahtani^{4*}, BDS; Sara Akeel^{1*}, PhD; Soulafa Almazrooa¹, DMSc

¹Department of Oral Diagnostic Sciences, Faculty of Dentistry, King Abdulaziz University, AlSulaimaniya, Jeddah, Saudi Arabia

²Department of Public Health, Faculty of Dentistry, King Abdulaziz University, Jeddah, Saudi Arabia

³Department of Computer Science, Faculty of Computing and Information Technology, King Abdulaziz University, Jeddah, Saudi Arabia

⁴Faculty of Dentistry, King Abdulaziz University, Jeddah, Saudi Arabia

*these authors contributed equally

Corresponding Author:

Sarah AlFarabi Ali, PhD Department of Oral Diagnostic Sciences, Faculty of Dentistry, King Abdulaziz University, AlSulaimaniya, Jeddah, Saudi Arabia

Abstract

Background: The use of artificial intelligence (AI), especially large language models (LLMs), is increasing in health care, including in dentistry. There has yet to be an assessment of the diagnostic performance of LLMs in oral medicine.

Objective: We aimed to compare the effectiveness of ChatGPT (OpenAI) and Microsoft Copilot (integrated within the Microsoft 365 suite) with oral medicine consultants in formulating accurate differential and final diagnoses for oral lesions from written clinical scenarios.

Methods: Fifty comprehensive clinical case scenarios including patient age, presenting complaint, history of the presenting complaint, medical history, allergies, intra- and extraoral findings, lesion description, and any additional information including laboratory investigations and specific clinical features were given to three oral medicine consultants, who were asked to formulate a differential diagnosis and a final diagnosis. Specific prompts for the same 50 cases were designed and input into ChatGPT and Copilot to formulate both differential and final diagnoses. The diagnostic accuracy was compared between the LLMs and oral medicine consultants.

Results: ChatGPT exhibited the highest accuracy, providing the correct differential diagnoses in 37 of 50 cases (74%). There were no significant differences in the accuracy of providing the correct differential diagnoses between AI models and oral medicine consultants. ChatGPT was as accurate as consultants in making the final diagnoses, but Copilot was significantly less accurate than ChatGPT (P=.015) and one of the oral medicine consultants (P<.001) in providing the correct final diagnosis.

Conclusions: ChatGPT and Copilot show promising performance for diagnosing oral medicine pathology in clinical case scenarios to assist dental practitioners. ChatGPT-4 and Copilot are still evolving, but even now, they might provide a significant advantage in the clinical setting as tools to help dental practitioners in their daily practice.

(JMIR AI 2025;4:e70566) doi:10.2196/70566

KEYWORDS

artificial intelligence; ChatGPT; Copilot; diagnosis; oral medicine; diagnostic performance; large language model; lesion; oral lesion

Introduction

Creating models that accurately replicate the complexity of the human brain and thinking has been a longstanding challenge for the scientific community [1]. The term "artificial intelligence" (AI) was first coined by John McCarthy in 1956, and this evolving scientific and engineering challenge focuses on computationally understanding intelligent behavior and

https://ai.jmir.org/2025/1/e70566

RenderX

creating applications that demonstrate such behavior [2]. AI has also emerged as a promising avenue for enhancing the precision and efficiency of diagnosing oral lesions. The diagnosis of pathological conditions within the oral cavity has traditionally relied on visual examination, histopathological analysis, and clinical expertise [3]. However, AI algorithms have the potential to analyze various data sources, including clinical images, patient records, and radiographs, to provide valuable insights

and suggestions for clinicians to facilitate the diagnosis of oral lesions [4].

ChatGPT is a recently introduced AI tool developed by OpenAI. ChatGPT is a large language model (LLM) trained with extensive data and capable of understanding and generating human-like responses accurately and consistently. ChatGPT currently operates on the GPT-4 architecture, allowing it to understand and respond to complex queries in a conversational manner [5]. ChatGPT can be used in medicine by rapidly providing appropriate answers to queries (or "prompts"), for instance, by assisting in decision-making based on up-to-date research and guidelines. There are high expectations for ChatGPT in the health sciences, including for education, research, and practice across different medical disciplines [6], and it can be embedded in various platforms.

Microsoft Copilot is another AI-driven assistant that can be accessed via a web interface or through seamless integration within the Microsoft 365 suite [5]. Leveraging LLMs and insights from Microsoft Graph, Microsoft Copilot delivers tailored support, enhancing the users' experience across Microsoft 365 applications such as Word, Excel, and PowerPoint. Copilot offers real-time suggestions and completions based on the context of the existing request. Powered by GPT-4 Turbo, it also has access to information, enhancing its utility for up-to-date coding tasks.

ChatGPT has also been used in several areas of medicine. For example, ChatGPT provided excellent responses on basic knowledge, lifestyle advice, and treatment for cirrhosis and hepatocellular carcinoma but performed less well for diagnosis and prevention [7]. In an analysis of ChatGPT responses to 284 medical questions, the results were highly accurate but incomplete [8]. AI has also been applied to dentistry [9-11]. In endodontics, AI models have been used to explore the anatomy of the root canal system, predict the health of dental pulp stem cells, detect root fractures and periapical lesions, and predict the success of retreatment procedures [12,13]. In oral medicine, ChatGPT was used to address questions about oral potentially malignant disorders. Guidelines on oral potentially malignant disorders from scientific societies were used to create questions for input into ChatGPT, which showed moderate knowledge about oral potentially malignant disorders as assessed by specialist reviewers [6]. AI also shows promise for scheduling, patient management, managing drug interactions, predictive tasks, and even robotic endodontic surgery [14], although the cost-effectiveness, reliability, and practicality of implementation still need to be assessed before widespread adoption [11].

To the best of our knowledge, no study has examined the use of AI-powered tools (ChatGPT and Copilot) in oral medicine, especially with respect to the diagnosis of oral lesions. To address this gap, herein, we compared the accuracy of ChatGPT and Copilot with oral medicine consultants in providing differential and final diagnoses from text-based clinical case scenarios.

Methods

Study Design

This was a comparative analytical study conducted at the King Abdulaziz University Faculty of Dentistry in Jeddah, Saudi Arabia. The primary objective was to assess and compare the accuracy of ChatGPT and Copilot with oral medicine consultants for diagnosing oral lesions from written clinical scenarios.

Ethical Considerations

The Research Ethics Committee-Faculty of Dentistry, King Abdulaziz University granted ethical approval (no. 209-11-23).

Data Collection

Sixty clinical case scenarios were collected from the Oral Medicine and Oral Pathology Division of the Oral Diagnostic Sciences Department. The final diagnosis was determined on the basis of the results of laboratory investigations, radiographs, and histopathological examination. Ten cases were excluded by an external reviewer, as they were deemed to be poorly written. The remaining 50 cases included patient age, chief complaint, history of the chief complaint, medical history, allergies, intra- and extraoral findings, a description of the lesions, and any additional information, including laboratory investigations and specific clinical features. An example clinical scenario is shown in Multimedia Appendix 1. The LLMs and oral medicine consultants were not provided with the histopathological features.

The cases were given to 3 oral medicine consultants (with clinical experiences of 7 years, 10 years, and 5 years for consultants 1, 2, and 3, respectively), who were asked to formulate differential and final diagnoses. Two specific prompts were designed for entry into ChatGPT and Copilot to formulate differential and final diagnoses (Figure 1): the first prompt enquired about the differential diagnoses for each clinical scenario ("As an oral medicine consultant, what is your differential diagnosis of the case?"), and the second enquired about the final diagnosis ("What is your final diagnosis based on the provided information?").



Figure 1. Schematic of the study design, describing the distribution of the clinical case scenarios to the oral medicine consultants and artificial intelligence (AI)-powered tools.



Responses were reviewed and evaluated independently by two reviewers who specialized in oral pathology and medicine and who were involved in case selection. Any discrepancies were resolved by a third reviewer. Each response was assessed for accuracy and assigned a score based on the following criteria: the differential diagnoses responses by ChatGPT, Copilot, and consultants were assigned a score of 2 (correctly identified all the listed differential diagnoses), 1 (correctly identified all the listed differential diagnoses). For the final diagnosis, responses were categorized as 1 (correct) or 0 (incorrect).

Statical Analysis

The performances of ChatGPT, Copilot, and the oral medicine consultants in providing differential and final diagnoses for oral lesions in the clinical scenarios are presented as frequency tables. The χ^2 or Fisher exact test was used to compare the performance distributions between the AI tools and consultants. A *P* value of .05 was considered significant. All statistical analyses were performed using IBM SPSS Statistics version 29.0.0 (IBM Statistics).

Results

Comparison of Differential Diagnoses Between AI Tools and Oral Medicine Consultants

ChatGPT exhibited the highest accuracy, correctly diagnosing 74% (37/50) of the cases, partially diagnosing 24% (12/50) of the cases correctly, and making completely incorrect diagnoses in only 2% (1/50) of the cases. In contrast, Copilot provided all correct differential diagnoses for 60% (30/50) of the cases, only one correct diagnosis in 34% (17/50) of the cases, and all wrong diagnoses in 6% (3/50) of the cases. There was no significant difference in the accuracy between the two models (P=.32).

In comparison to these AI models, oral medicine consultant 1 correctly diagnosed 60% (30/50), partially diagnosed 34% (17/50), and incorrectly diagnosed 6% (3/50) of the cases (P=.32 vs ChatGPT and P≥.99 vs Copilot). Oral medicine consultant 2 correctly diagnosed 72% (36/50) of the cases, partially diagnosed 22% (11/50) of the cases, and incorrectly diagnosed 6% (3/50) of the cases (P=.75 vs ChatGPT and P=.41 vs Copilot). Lastly, oral medicine consultant 3 accurately diagnosed 54% (27/50) of the cases, partially diagnosed 38% (19/50) of the cases, and incorrectly diagnosed 54% (27/50) of the cases, partially diagnosed 8% (4/50) of the cases (P=.10 vs ChatGPT and P=.82 vs Copilot). The AI models had similar accuracy in providing the differential diagnoses as oral medicine consultants, as shown in Table 1.



Table .	Comparison of	accuracy of ar	tificial intelligence	(AI)	versus oral	medicine	consultants f	for differential	diagnoses.
---------	---------------	----------------	-----------------------	------	-------------	----------	---------------	------------------	------------

AI model or consultant	Differential diagnosis, r	Differential diagnosis, n (%)						
	All wrong	One correct	All correct	P value ^a	<i>P</i> value ^b			
Oral medicine consul- tant 1	3 (6)	17 (34)	30 (60)	.31	≥.99			
Oral medicine consul- tant 2	3 (6)	11 (22)	36 (72)	.74	.42			
Oral medicine consul- tant 3	4 (8)	19 (38)	27 (54)	.11	.82			
ChatGPT	1 (2)	12 (24)	37 (74)	_c	.32			
Copilot	3 (6)	17 (34)	30 (60)	.32	_			

^a*P* values in comparison to ChatGPT.

^b*P* values in comparison to Copilot.

^c"–": not applicable.

Comparison of Final Diagnoses Between AI Tools and Oral Medicine Consultants

With respect to the definitive diagnoses, ChatGPT again showed the highest accuracy: 70% (35/50) correct diagnoses and 30% (15/50) incorrect diagnoses. Copilot performed less well, providing 46% (23/50) correct diagnoses and 40% (27/50) incorrect diagnoses.

Oral medicine consultant 1 correctly diagnosed 66% (33/50) of the cases and incorrectly diagnosed 34% (17/50) of the cases (P=.66 vs ChatGPT and P=.04 vs Copilot). Oral medicine consultant 2 had the highest diagnostic accuracy, diagnosing 80% (40/50) of the cases correctly and 20% (10/50) incorrectly (P=.25 vs ChatGPT and P<.001 vs Copilot). Oral medicine consultant 3 correctly diagnosed 64% (32/50) of the cases and incorrectly diagnosed 36% (18/50) of the cases (P=.52 vs ChatGPT and P=.07 vs Copilot); the data are shown in Table 2.

Table .	Comparison	of accuracy	of artificial	intelligence	(AI) versus	oral medicine	consultants for fir	nal diagnoses.
---------	------------	-------------	---------------	--------------	-------------	---------------	---------------------	----------------

AI model or consultant	Final diagnosis, n (%)			
	Wrong	Correct	<i>P</i> value ^a	<i>P</i> value ^b
Oral medicine consultant 1	17 (34)	33 (66)	.67	.04
Oral medicine consultant 2	10 (20)	40 (80)	.25	<.001
Oral medicine consultant 3	18 (36)	32 (64)	.52	.07
ChatGPT	15 (30)	35 (70)	_ ^c	.02
Copilot	27 (54)	23 (46)	.02	_

^aP values in comparison to ChatGPT.

^b*P* values in comparison to Copilot.

^c"–": not applicable.

Discussion

In this study, we compared the diagnostic accuracy of AI language models (ChatGPT-4 and Copilot) with three oral medicine consultants in providing differential and final diagnoses for oral lesions from text-based clinical scenarios. We found that the diagnostic accuracy of the LLMs and oral medicine consultants for providing accurate differential diagnoses was similar. However, Copilot was significantly less accurate than ChatGPT (P=.015) and one of the oral medicine consultants (P<.001) in providing the correct final diagnoses. Our results suggest that advanced language models, especially ChatGPT, can provide comparable diagnostic insights to human experts in the context of oral lesion diagnosis. ChatGPT-4 and Copilot are still evolving, but even now, they might provide a

RenderX

significant advantage in the clinical setting as tools to help dental practitioners in their daily practice. Copilot may have underperformed in making the final diagnoses compared to ChatGPT and consultants due to differences in training, dataset variations, and algorithmic constraints. ChatGPT is exposed to a broader range of medical and dental literature, whereas Copilot is optimized for general productivity, affecting its diagnostic precision. Additionally, Copilot's customization for enterprise applications may limit its ability to provide accurate clinical diagnoses [15].

Our findings are consistent with those obtained by Altamimi et al [16], who concluded that AI tools can be useful in clinical settings to provide diagnoses for certain conditions. Friederichs et al [17] evaluated the performance of ChatGPT using 400 multiple-choice questions from the progress test administered

in German-speaking countries, reporting that ChatGPT surpassed most first- to third-year medical students by correctly answering two-thirds of the multiple-choice questions, with proficiency equivalent to the level required for the German state licensing examination in Progress Test Medicine. Several studies have reported similar accuracy and efficacy of ChatGPT. A recent study from India demonstrated that ChatGPT was a reliable tool for addressing complex problems that involved higher-level cognitive skills such as interpretation, analysis, evaluation, and evidence-based opinion or prediction, correctly answering 100 complex questions in pathology [18]. Das et al [19] reported that ChatGPT could be considered a tool for answering direct inquiries regarding microbiology, showing 80% accuracy in its responses. Furthermore, Johnson et al [8] found that ChatGPT consistently provided accurate and comprehensive responses to a variety of questions in the medical field.

Copilot showed promising performance in providing differential diagnoses compared with oral medicine consultants, albeit with higher rates of all wrong differential diagnoses. Kaftan et al [20] recently examined the accuracy of AI-powered tools for interpreting biochemical data, reporting the highest accuracy for Copilot compared with ChatGPT-3.5 and Gemini. However, ChatGPT-3.5 had fewer capabilities than ChatGPT-4, which we used here. While Copilot is based on GPT-4, as noted above, its outputs differ due to Microsoft-specific customizations, including specialized training for productivity tasks, integration with enterprise tools, and compliance filters, perhaps explaining the difference in results between the two LLMs [20]. Tepe and Emekli [21] similarly observed significant variability between LLMs for answering prompts related to breast imaging. ChatGPT-4 showed high accuracy in responding to these questions, outperforming Gemini and Copilot. Moreover, AI-powered tools tended to give more differential diagnoses for each clinical scenario, regardless of whether the answers were all correct or not, with only two answers needed for analysis. Accordingly, expert judgment, knowledge, and experience are required to evaluate these answers to construct specific differential diagnoses for each case.

Diniz-Freitas et al [6] reported that integrating ChatGPT into oral medicine could significantly accelerate decision-making for patient diagnosis, treatment, and care. We found that oral medicine consultants outperformed Copilot with respect to the final diagnosis. However, one of the oral medicine consultants outperformed Copilot in providing an accurate final diagnosis, and this was the consultant with the most experience. Clinicians accumulate subject-specific knowledge and experience. Consequently, AI tools like ChatGPT, when paired with health care practitioners' expertise, could yield even more dependable and efficient outcomes for patients requiring oral medicine treatment [6].

AI tools obtain their datasets from different sources and have different training, which influences their applications and affects their performance. Training AI tools with medical or dental datasets reviewed by specialists might be expected to improve results and transform diagnostic health care services. The clinician's experience, which is influenced by solid knowledge and experience and unaffected by dataset variability, plays a major role in their superiority over AI tools [22].

As the training and refinement of filtered datasets improve AI tools, LLMs are expected to be integrated into clinical workflows, especially in areas without access to specialized consultants in the field. During the implementation of such technologies, ethical concerns should be considered and governed. The privacy and safety of patient data are major concerns in the use of AI in health care, requiring adherence to regulations like the HIPAA (Health Insurance Portability and Accountability Act) to prevent unauthorized access [23,24]. There are also medico-legal concerns, as AI-related errors could lead to liability issues, necessitating clear regulatory frameworks. Ethically, AI should serve as an assistive tool rather than a replacement for clinical expertise to maintain fairness and reliability. Clinician reliance on AI must be balanced to ensure that decision-making remains informed by human judgment, supported by proper training.

This study has some limitations. It was a pilot study that focused solely on evaluating the application of AI-powered tools in diagnosing text-based clinical scenarios specific to oral medicine. Therefore, the findings and conclusions may not be applicable or generalizable to other subjects or domains. Depending on text-based clinical scenarios makes it more difficult to provide both differential and definitive diagnoses. Using clinical images and histopathological findings greatly improves the accuracy of diagnostics, which were not provided in this study. Moreover, we only studied a limited number of cases (50 questions), and 10 cases were excluded by an external reviewer, which may have introduced bias. The formulation of the input "prompts" when interacting with language models can greatly impact the quality and nature of the generated responses. Consequently, further studies are needed to examine the optimal prompts that provide the best and most accurate responses. Moreover, it remains uncertain whether LLMs consistently produce identical or similar responses to the same query at different times. In this study, each question was submitted only once to the AI-powered tools, which may have limited the assessment of response consistency. Additional studies are needed to overcome these limitations and explore the real-world potential of using AI-powered tools in oral medicine.

In conclusion, LLMs such as ChatGPT and Copilot showed promising performance in making diagnoses in oral medicine clinical case scenarios. ChatGPT-4 and Copilot are still evolving, but even now might provide a significant advantage in the clinical setting as tools to help dental practitioners in their daily practice. Such technologies could particularly benefit dentists in rural areas or areas with no access to oral medicine consultants, who—provided the technology is further validated—could collect medical histories, perform extra- and intraoral examinations, and provide these data to LLMs systems to provide a set of relevant differential diagnoses to help with decision-making regarding further testing, referral, or simple management.

```
XSL•FO
```

None declared.

Multimedia Appendix 1 Example clinical scenario. [DOCX File, 14 KB - ai_v4i1e70566_app1.docx]

References

- 1. Alexander B, John S. Artificial intelligence in dentistry: current concepts and a peep into the future. IJAR 2018;6(12):1105-1108. [doi: 10.21474/IJAR01/8242]
- 2. Shapiro SC. Encyclopedia of Artificial Intelligence, 2nd edition: A Wiley Interscience Publication; 1992. URL: <u>https://cse.</u> <u>buffalo.edu/~rapaport/Papers/belrepsys92.encyai2.pdf</u> [accessed 2025-04-21]
- Aubreville M, Knipfer C, Oetter N, et al. Automatic classification of cancerous tissue in laserendomicroscopy images of the oral cavity using deep learning. Sci Rep 2017 Sep 20;7(1):11979. [doi: <u>10.1038/s41598-017-12320-8</u>] [Medline: <u>28931888</u>]
- 4. Khanagar SB, Al-Ehaideb A, Maganur PC, et al. Developments, application, and performance of artificial intelligence in dentistry a systematic review. J Dent Sci 2021 Jan;16(1):508-522. [doi: 10.1016/j.jds.2020.06.019] [Medline: 33384840]
- 5. Kongas K. GitHub copilot and chatgpt comparison in improving software development productivity [Master's thesis]. : LUT University; 2024.
- Diniz-Freitas M, Rivas-Mundiña B, García-Iglesias JR, García-Mato E, Diz-Dios P. How ChatGPT performs in oral medicine: the case of oral potentially malignant disorders. Oral Dis 2024 May;30(4):1912-1918. [doi: <u>10.1111/odi.14750</u>] [Medline: <u>37794649</u>]
- 7. Yeo YH, Samaan JS, Ng WH, et al. Assessing the performance of ChatGPT in answering questions regarding cirrhosis and hepatocellular carcinoma. Clin Mol Hepatol 2023 Jul;29(3):721-732. [doi: 10.3350/cmh.2023.0089] [Medline: 36946005]
- Johnson D, Goodman R, Patrinely J, et al. Assessing the accuracy and reliability of AI-generated medical responses: an evaluation of the Chat-GPT model. Res Sq 2023 Feb 28:rs.3.rs-2566942. [doi: <u>10.21203/rs.3.rs-2566942/v1</u>] [Medline: <u>36909565</u>]
- 9. Alhaidry HM, Fatani B, Alrayes JO, Almana AM, Alfhaed NK. ChatGPT in dentistry: a comprehensive review. Cureus 2023 Apr;15(4):e38317. [doi: 10.7759/cureus.38317] [Medline: 37266053]
- 10. Asiri AF, Altuwalah AS. The role of neural artificial intelligence for diagnosis and treatment planning in endodontics: a qualitative review. Saudi Dent J 2022 May;34(4):270-281. [doi: 10.1016/j.sdentj.2022.04.004] [Medline: 35692236]
- 11. Agrawal P, Nikhade P. Artificial intelligence in dentistry: past, present, and future. Cureus 2022 Jul;14(7):e27405. [doi: 10.7759/cureus.27405] [Medline: 36046326]
- 12. Aminoshariae A, Kulild J, Nagendrababu V. Artificial intelligence in endodontics: current applications and future directions. J Endod 2021 Sep;47(9):1352-1357. [doi: 10.1016/j.joen.2021.06.003] [Medline: 34119562]
- 13. Boreak N. Effectiveness of artificial intelligence applications designed for endodontic diagnosis, decision-making, and prediction of prognosis: a systematic review. J Contemp Dent Pract 2020 Aug 1;21(8):926-934. [Medline: <u>33568617</u>]
- 14. Meghil MM, Rajpurohit P, Awad ME, McKee J, Shahoumi LA, Ghaly M. Artificial intelligence in dentistry. Dentistry Review 2022 Mar;2(1):100009. [doi: 10.1016/j.dentre.2021.100009]
- 15. Daza J, Bezerra LS, Santamaría L, et al. Evaluation of four chatbots in autoimmune liver disease: a comparative analysis. Ann Hepatol 2025 Jan;30(1):101537. [doi: 10.1016/j.aohep.2024.101537]
- 16. Altamimi A, Aldughaim A, Alotaibi S, Alrehaili J, Bakir M, Almuhainy A. Evaluating the precision of ChatGPT artificial intelligence in emergency differential diagnosis. JMLPH 2024;4(1):327-337. [doi: <u>10.52609/jmlph.v4i1.113</u>]
- 17. Friederichs H, Friederichs WJ, März M. ChatGPT in medical school: how successful is AI in progress testing? Med Educ Online 2023 Dec;28(1):2220920. [doi: 10.1080/10872981.2023.2220920] [Medline: 37307503]
- Sinha RK, Deb Roy A, Kumar N, Mondal H. Applicability of ChatGPT in assisting to solve higher order problems in pathology. Cureus 2023 Feb;15(2):e35237. [doi: <u>10.7759/cureus.35237</u>] [Medline: <u>36968864</u>]
- Das D, Kumar N, Longjam LA, et al. Assessing the capability of ChatGPT in answering first- and second-order knowledge questions on microbiology as per competency-based medical education curriculum. Cureus 2023 Mar;15(3):e36034. [doi: <u>10.7759/cureus.36034</u>] [Medline: <u>37056538</u>]
- 20. Kaftan AN, Hussain MK, Naser FH. Response accuracy of ChatGPT 3.5 Copilot and Gemini in interpreting biochemical laboratory data a pilot study. Sci Rep 2024 Apr 8;14(1):8233. [doi: 10.1038/s41598-024-58964-1] [Medline: 38589613]
- Tepe M, Emekli E. Assessing the responses of large language models (ChatGPT-4, Gemini, and Microsoft Copilot) to frequently asked questions in breast imaging: a study on readability and accuracy. Cureus 2024 May;16(5):e59960. [doi: <u>10.7759/cureus.59960</u>] [Medline: <u>38726360</u>]
- 22. Yau JYS, Saadat S, Hsu E, et al. Accuracy of prospective assessments of 4 large language model chatbot responses to patient questions about emergency care: experimental comparative study. J Med Internet Res 2024 Nov 4;26:e60291. [doi: 10.2196/60291] [Medline: 39496149]

- Giannakopoulos K, Kavadella A, Aaqel Salim A, Stamatopoulos V, Kaklamanos EG. Evaluation of the performance of generative ai large language models ChatGPT, Google Bard, and Microsoft Bing Chat in supporting evidence-based dentistry: comparative mixed methods study. J Med Internet Res 2023 Dec 28;25:e51580. [doi: <u>10.2196/51580</u>] [Medline: <u>38009003</u>]
- MacIntyre MR, Cockerill RG, Mirza OF, Appel JM. Ethical considerations for the use of artificial intelligence in medical decision-making capacity assessments. Psychiatry Res 2023 Oct;328:115466. [doi: <u>10.1016/j.psychres.2023.115466</u>] [Medline: <u>37717548</u>]

ABBREVIATIONS

AI: artificial intelligence HIPAA: Health Insurance Portability and Accountability Act LLM: large language model

Edited by F Dankar, S Gardezi; submitted 26.12.24; peer-reviewed by K Sunil, Z Ehtesham; revised version received 17.03.25; accepted 18.03.25; published 24.04.25.

Please cite as:

AlFarabi Ali S, AlDehlawi H, Jazzar A, Ashi H, Esam Abuzinadah N, AlOtaibi M, Algarni A, Alqahtani H, Akeel S, Almazrooa S The Diagnostic Performance of Large Language Models and Oral Medicine Consultants for Identifying Oral Lesions in Text-Based Clinical Scenarios: Prospective Comparative Study JMIR AI 2025;4:e70566 URL: https://ai.jmir.org/2025/1/e70566 doi:10.2196/70566

© Sarah AlFarabi Ali, Hebah AlDehlawi, Ahoud Jazzar, Heba Ashi, Nihal Esam Abuzinadah, Mohammad AlOtaibi, Abdulrahman Algarni, Hazzaa Alqahtani, Sara Akeel, Soulafa Almazrooa. Originally published in JMIR AI (https://ai.jmir.org), 24.4.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Comparison of Deep Learning Approaches Using Chest Radiographs for Predicting Clinical Deterioration: Retrospective Observational Study

Mahmudur Rahman¹, PhD; Jifan Gao², MS; Kyle A Carey³, MPH; Dana P Edelson³, MD, MS; Askar Afshar¹, MS; John W Garrett^{2,4}, PhD; Guanhua Chen², PhD; Majid Afshar^{1,2}, MD, MSCR; Matthew M Churpek^{1,2}, MD, MPH, PhD

¹Department of Medicine, University of Wisconsin-Madison, 610 Walnut St, Madison, WI, United States

²Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison, Madison, WI, United States

³Department of Medicine, University of Chicago, Chicago, IL, United States

⁴Department of Radiology, University of Wisconsin-Madison, Madison, WI, United States

Corresponding Author:

Matthew M Churpek, MD, MPH, PhD Department of Medicine, University of Wisconsin-Madison, 610 Walnut St, Madison, WI, United States

Abstract

Background: The early detection of clinical deterioration and timely intervention for hospitalized patients can improve patient outcomes. The currently existing early warning systems rely on variables from structured data, such as vital signs and laboratory values, and do not incorporate other potentially predictive data modalities. Because respiratory failure is a common cause of deterioration, chest radiographs are often acquired in patients with clinical deterioration, which may be informative for predicting their risk of intensive care unit (ICU) transfer.

Objective: This study aimed to compare and validate different computer vision models and data augmentation approaches with chest radiographs for predicting clinical deterioration.

Methods: This retrospective observational study included adult patients hospitalized at the University of Wisconsin Health System between 2009 and 2020 with an elevated electronic cardiac arrest risk triage (eCART) score, a validated clinical deterioration early warning score, on the medical-surgical wards. Patients with a chest radiograph obtained within 48 hours prior to the elevated score were included in this study. Five computer vision model architectures (VGG16, DenseNet121, Vision Transformer, ResNet50, and Inception V3) and four data augmentation methods (histogram normalization, random flip, random Gaussian noise, and random rotate) were compared using the area under the receiver operating characteristic curve (AUROC) and the area under the precision-recall curve (AUPRC) for predicting clinical deterioration (ie, ICU transfer or ward death in the following 24 hours).

Results: The study included 21,817 patient admissions, of which 1655 (7.6%) experienced clinical deterioration. The DenseNet121 model pretrained on chest radiograph datasets with histogram normalization and random Gaussian noise augmentation had the highest discrimination (AUROC 0.734 and AUPRC 0.414), while the vision transformer having 24 transformer blocks with random rotate augmentation had the lowest discrimination (AUROC 0.598).

Conclusions: The study shows the potential of chest radiographs in deep learning models for predicting clinical deterioration. The DenseNet121 architecture pretrained with chest radiographs performed better than other architectures in most experiments, and the addition of histogram normalization with random Gaussian noise data augmentation may enhance the performance of DenseNet121 and pretrained VGG16 architectures.

(JMIR AI 2025;4:e67144) doi:<u>10.2196/67144</u>

KEYWORDS

chest X-ray; critical care; deep learning; chest radiographs; radiographs; clinical deterioration; predictive; deterioration; retrospective; data; dataset; artificial intelligence; AI; chest; patient; hospitalized

Introduction

Clinical deterioration is common in hospitalized patients and can lead to adverse outcomes, including increased morbidity and mortality if not identified and managed properly [1]. The early detection of patient deterioration and timely intervention

https://ai.jmir.org/2025/1/e67144

RenderX

can improve patient outcomes [2]. Various early warning scores (EWS) have been developed to identify the deterioration risk by monitoring different clinical variables, and the implementation of machine-learning EWS, such as the electronic cardiac arrest risk triage (eCART) score, has been associated with improved mortality [3-6]. Current EWS rely on structured

data, such as vital signs and laboratory values, to predict clinical deterioration and ignore other data modalities that could potentially enhance prediction accuracy [7]. This results in lower detection and higher false-positive rates for these scores that could be mitigated by incorporating additional modalities [8].

Because respiratory failure is a common cause of clinical deterioration, the use of computer vision models with chest radiographs is a promising direction for improving EWS performance [9]. Although traditional computer vision models have historically been used to analyze chest radiographs, prior work on chest radiographs is limited to identifying specific diagnoses [10-12]. In some recent studies, chest radiographs are used to detect lung disease [[13,14]], acute respiratory distress syndrome [15], pneumonia [16,17], tuberculosis [18,19], and COVID-19 [20]. However, to facilitate other tasks with comprehensive machine understanding, chest X-ray interpretation models are being more commonly used with the help of computer vision and transformer-based natural language processing models [21,22]. The advancements in predictive analytics with deep learning methods have led to increased capabilities to extract meaningful information from medical images, including chest radiographs [23]. However, deep learning models have never been trained with chest radiographs to predict clinical deterioration outside the intensive care unit (ICU). There are numerous deep learning architectures for chest radiograph prediction models, such as VGG16, ResNet50, DenseNet121, and Vision Transformer, and the performance of these models is unknown for this specific task. Additionally, there are different data augmentation techniques available to further enhance the performance of a vision model by improving model generalization, but it is unknown whether these data augmentation techniques would improve the performance of the prediction model for this task.

To address these knowledge gaps, the objective of this study was to compare different computer vision architectures and augmentation methods with chest radiographs for predicting clinical deterioration. Our training pipeline incorporates extensive hyperparameter tuning through Bayesian optimization and validates the generalizability of models in a separate hold-out test set. The findings of our experiments have important implications for researchers developing computer vision deep learning models for clinical applications with chest radiographs.

Methods

Ethical Considerations

The study protocol was reviewed and approved by the University of Wisconsin Institutional Review Board (approval #2019 - 1258). This study was a secondary analysis of limited HIPAA data from hospital electronic health records. The study was approved with a waiver of informed consent.

All direct identifiers of patients whose data were used in this study were de-identified prior to analysis to ensure participants privacy and confidentiality. Minimal necessary identifiable information was accessed or stored during the study beyond possible HIPAA data in clinical notes, radiological images, and real dates.

Participants did not receive any compensation for this data analysis, as no new data were collected and no direct contact with participants occurred.

Study Population and Data Collection

All adult patients (age \geq 18 years) hospitalized at the University of Wisconsin Health System (UW Health) between 2009 and 2020 with an elevated eCART score ≥93 (which is the threshold used in clinical practice at UW Health) on the medical-surgical wards were eligible for inclusion in this retrospective cohort study. The eCART score [3] is a validated EWS currently in clinical practice and cleared by the Food and Drug Administration that combines demographics, vital signs, and laboratory results in a gradient-boosted machine model to predict future clinical deterioration. The rationale for only including patients with an elevated score is based on creating an enriched cohort where chest radiograph models can enhance the prediction and mitigate the false-positive alerts from these scores. Furthermore, this simplifies the prediction task to a single time point, making it more feasible to compare multiple models and augmentation strategies. Patients with a chest radiograph within 48 hours before the first elevated eCART score were included in the study. Available anterior-posterior or posterior-anterior views were included in the study cohort. In addition to chest radiographs, additional study variables that were collected included patient demographics, admission time, vital signs, laboratory values, patient location, and discharge disposition, which were all collected via the clinical research data warehouse. Figure 1 shows the patient encounter flow chart for inclusion into the analytic cohort.





Outcome

The study outcome of clinical deterioration was defined as a direct ward-to-ICU transfer or ward death within 24 hours of the time of the patient's first elevated eCART score.

Data Preprocessing

The chest radiograph closest to (but before) the time of the elevated eCART score was used to predict the corresponding deterioration outcome. To address variations in image acquisition and processing protocols, all radiographs were rescaled to a uniform size of 224×224 pixels using nearest neighbor interpolation. Additionally, to address the variabilities in imaging exposure levels, pixel intensity values were normalized to a range of [0, 1] by applying min-max scaling. The clinical deterioration outcome (ie, ICU transfer or mortality within 24 hours from the prediction time point) was encoded as binary labels, with one-hot encoding used for the binary prediction task. These preprocessing steps ensured the creation of a high-quality robust dataset for training deep learning models to predict clinical deterioration from chest radiographs.

```
https://ai.jmir.org/2025/1/e67144
```

RenderX

Model Development

For the prediction task, computer vision deep learning models were trained and optimized with the dataset created from the cohort. Five publicly available computer vision models were compared for our task: (1) VGG16 [24], (2) DenseNet121 [25], (3) Vision Transformer [26], (4) ResNet50 [27], and (5) Inception V3 [28]. DenseNet121 is a convolutional neural network notable for its dense connections between layers, improving efficiency and reducing risk of overfitting, and VGG16 is known for its simplicity using a series of convolutional layers with small filters followed by max pooling layers. The Vision Transformer model is based on the transformer architecture and uses the self-attention mechanism to process the images. The main rationale of adopting these computer vision models for clinical deterioration tasks is that they are widely used in other chest radiograph detection tasks in clinical setups [29-31]. In addition, these models are easy to implement, and various pretrained weights are readily available. As clinical tasks require fine-grained image understanding for different tasks, these models provide that performance with a

manageable model size. However, the main shortcoming of using these models is they do not provide any generalized image understanding for explainability.

We used two different versions of the VGG16 architecture, one using randomly initialized weights (without pretraining) and the other using model weights pretrained on ImageNet [32]. Two different versions of the DenseNet121 architecture were also used: one with model weights pretrained on Imagenet [32] and one pretrained on publicly available radiograph datasets [33]. Specifically, the radiograph datasets used for pretraining consisted of the following datasets: NIH aka Chest X-ray14 [34], PC aka PadChest [35], CheX aka CheXpert [36], MIMIC-CXR [37], OpenI [38], Google [39], and RSNA Pneumonia Detection Challenge [40]. For the Vision Transformer model, we trained two models without any pretrained weights of two different sizes, one with 12 transformer blocks and another with 24 transformer blocks. We employed batch normalization layers after every block to ensure the stability of the optimization process during the model training. Figure 2 presents the overall structure of this study.

For each of the above architectures, we compared them with and without different preprocessing and data augmentation approaches. These included histogram normalization, random rotation (± 15 degrees), horizontal flipping, and the addition of random Gaussian noise. Briefly, histogram normalization addresses the regional discrepancy of exposure levels in the case of some images. Additionally, given the presence of noise and artifacts during the acquisition of the radiographs, random Gaussian noise, which was implemented as 0.1 probability with 0 mean and 0.1 standard deviation, may make the models more robust to noise in the input image samples. Figure 3 shows the examples of all the augmentation methods we have used in this work.

Figure 2. Overall structure of this study. *VGG16, DenseNet121, ResNet50, and Inception V3 models were trained from randomly initialized weights and pretrained weights. Other models were trained with randomly initialized weights only.





Figure 3. Examples of different image augmentation methods we have utilized. HN: histogram normalization; RGN: random Gaussian noise.



We used the Bayesian optimization algorithm to find the optimal hyperparameters that maximize the area under the receiver operating characteristic curve (AUROC). Details of the hyperparameters are presented in Multimedia Appendix 1. To make the training procedure faster, we used Ray Tune [41] to parallelize the hyperparameter search process in a multi-GPU environment. We trained the model with a randomly selected 60% of the encounters in the dataset and validated it with the development set consisting of 20% of the encounters to optimize hyperparameters and determine the final settings. The remaining 20% of the encounters were completely separated for independent final model evaluation of the optimized models as a test set. We trained the models for 20 epochs and decreased the learning rate by a factor of 0.5 in every epoch. During the training, early stopping was used if the validation AUROC failed to improve in three consecutive epochs. We used Adam mini-batch gradient descent optimization with a batch size from the search space of 32, 64, and 128.

Model Evaluation

All combinations of image augmentations and deep learning computer vision architectures for the clinical deterioration task were evaluated using the test dataset. Predicted probabilities for the deterioration outcome were calculated for every encounter during the evaluation. Model discrimination was assessed using the AUROC and its 95% CI, calculated via the DeLong method [42] as the primary metric and the area under the precision-recall curve (AUPRC) as the secondary metric.

Table . Population characteristics of the study cohort (N=21,817).

The p-values of the AUROC scores are presented in Multimedia Appendix 1. As *P*<.001 in all cases, our AUROC scores are statistically significant.

Data cleaning and cohort selection with descriptive analysis were conducted using Stata version 16.1 (StataCorp). We used Python version 3.8.10, along with the Monai framework version 1.2.0 (NVIDIA) and Pytorch version 2.0.0 (Facebook) to develop the deep learning models. Additionally, the AUROC score and its 95% CI were calculated using FastDeLong implementation from VMAF (Video Multimethod Assessment Fusion; Netflix) [43].

Results

Cohort Characteristics

A total of 258,621 admissions occurred during the study period, and 92,845 had an elevated eCART score. Of these, for 21,817 admissions, a chest radiograph was obtained within 48 hours of the time of the elevated score and was included in the analysis (Figure 1). The characteristics of the final cohort are presented in Table 1. The patients in the final cohort had a median age of 63 (IQR 52-74) years, with a higher likelihood of being male (56.1%, 12,249/21,817); 5.7% were black (1252/21,187). The median time to eCART score elevation from admission was 21.8 (7.1-47.6) hours and the median time to eCART score elevation from the last radiograph was 9 (7.1-47.6) hours. About 7.5% (1655/21,817) of the encounters had an outcome event, including 4.1% (893/21,817) cases of in-hospital death.

Variable	Value
Age, years, median (IQR)	63 (52-74)
Female, n (%)	9568 (43.9)
Black race, n (%)	1252 (5.7)
Elevated eCART ^a score, n (%)	1655 (7.59)
Time to the elevated eCART score from admission, hours, median (IQR)	21.8 (7.1-47.6)
Time to the elevated eCART score from the last radiograph, hours, median (\mbox{IQR})	9.0 (7.1-47.6)
In-hospital mortality, n (%)	893 (4.1)

^aeCART: electronic cardiac arrest risk triage

Model Discrimination

The model performance AUROC and AUPRC values for all models across all image augmentation methods are presented

https://ai.jmir.org/2025/1/e67144

and AUPRC are presented in Multimedia Appendix 1. Additionally, receiver operating characteristic (ROC curves and

in Tables 2 and 3, respectively, and the 95% CI of the AUROC

precision-recall curves are shown in Figures 4 and 5, respectively.

Table . Model performance area under the receiver operating characteristic curve (AUROC) with the validation dataset across different model architectures, pretrained weights, and image augmentation methods.

Model	Pretrained weights	No transforma- tion	Histogram normalization (HN)	Random flip	Random Gaus- sian noise (RGN)	Random rotate	HN + RGN	Average AU- ROC score ^a
VGG16	Random init	0.694	0.723	0.698	0.701	0.674	0.712	0.700
VGG16	ImageNet	0.712	0.717	0.692	0.710	0.689	0.719	0.707
DenseNet121	ImageNet	0.683	0.701	0.672	0.700	0.678	0.716	0.692
DenseNet121	Radiographs	0.723	0.716	0.713	0.696	0.701	0.734	0.714
ResNet50	Random init	0.588	0.684	0.629	0.678	0.638	0.651	0.645
ResNet50	ImageNet	0.715	0.707	0.694	0.694	0.669	0.712	0.700
Inception V3	Random init	0.691	0.672	0.671	0.661	0.703	0.690	0.681
Inception V3	ImageNet	0.714	0.712	0.710	0.706	0.686	0.713	0.707
Vision Trans- former (12 Blocks)	Random init	0.661	0.648	0.617	0.652	0.623	0.652	0.642
Vision Trans- former (24 Blocks)	Random init	0.654	0.663	0.609	0.651	0.598	0.662	0.640
Average Score over models ^b	c	0.684	0.694	0.671	0.685	0.666	0.696	—
Average Im- provement ^d	—	—	0.010	-0.013	0.001	-0.028	0.012	—

^aThe average AUROC score is for a particular model over different augmentation methods

^bThe "Average score over models" row presents the average AUROC score of a particular augmentation method over different models.

^c"—" indicates not applicable.

^dThe "Average improvement" row shows the average AUROC improvement of an augmentation method over the baseline score without any transformation.



Rahman et al

Table .	Model performance area under the precision-recall curve (AUPRC) score	s with the	validation	dataset acro	oss different	model	architectures,
pretrain	ed weights, and image augmentation methods.						

Model	Pretrained weights	No transforma- tion	Histogram normalization (HN)	Random flip	Random Gaus- sian noise (RGN)	Random rotate	HN + RGN	Average AUPRC score ^a
VGG16	Random init	0.346	0.398	0.329	0.349	0.320	0.378	0.353
VGG16	ImageNet	0.371	0.403	0.306	0.343	0.311	0.389	0.354
DenseNet121	ImageNet	0.321	0.373	0.360	0.355	0.365	0.379	0.359
DenseNet121	Radiographs	0.395	0.326	0.338	0.360	0.358	0.414	0.365
ResNet50	Random init	0.135	0.229	0.147	0.243	0.215	0.174	0.191
ResNet50	ImageNet	0.405	0.378	0.357	0.320	0.288	0.344	0.349
Inception V3	Random init	0.319	0.247	0.247	0.304	0.343	0.339	0.300
Inception V3	ImageNet	0.440	0.340	0.421	0.399	0.361	0.369	0.388
Vision Trans- former (12 Blocks)	Random init	0.205	0.189	0.143	0.209	0.139	0.204	0.182
Vision Trans- former (24 Blocks)	Random init	0.187	0.219	0.121	0.177	0.118	0.196	0.170
Average score over models ^b	c	0.313	0.310	0.277	0.306	0.282	0.319	—
Average im- provement ^d	—	—	-0.003	-0.036	-0.007	-0.031	0.006	—

^aThe average AUPRC score is for a particular model over different augmentation methods

^bThe "Average score over models" row presents the average AUROC score of a particular augmentation method over different models.

c"___" indicates not applicable.

^dThe "Average improvement" row shows the average AUROC improvement of an augmentation method over the baseline score without any transformation.





Figure 4. Receiver operating characteristic (ROC) curve of the best-performing models in every network architecture. Actual AUROC values are included in the corresponding label. HN: histogram normalization; RGN: random Gaussian noise; AUROC: area under the receiver operating characteristic curve.

Figure 5. Precision-recall curves of the best-performing models in every network architecture. Actual AUPRC values are included in the corresponding label. Best viewed in color. HN: histogram normalization; RGN: random Gaussian noise; AUPRC: area under the precision-recall curve.



https://ai.jmir.org/2025/1/e67144

Across all architectures and augmentation combinations, the DenseNet121 model pretrained with chest radiographs and augmented with histogram normalization and Gaussian noise had the highest AUROC (0.734) across all the models. Similarly, when averaged across all augmentation methods, the DenseNet121 models pretrained with chest radiographs had a higher average discrimination than any other architecture in terms of the AUROC (0.714). The vision transformer architectures (12 and 24 transformer blocks) performed similarly to each other on average and had worse average AUROC than other models (0.642 and 0.640 for 12 and 24 transformer blocks, respectively). In terms of the AUPRC, DenseNet121 pretrained with chest radiographs and augmented with histogram normalization and Gaussian noise also had the highest performance (0.414). Accordingly, compared with other models, Inception V3 pretrained with ImageNet had the highest AUPRC (0.388) on average.

In terms of the image augmentation methods, the histogram normalization with random Gaussian noise image augmentations had the best mean AUROC (0.696) when averaged across all architectures, followed by histogram normalization augmentation alone (0.694). The random rotate augmentation had the worst average performance in terms of the AUROC (0.666). In terms of the AUPRC, histogram normalization with random Gaussian noise image augmentations also had the highest average AUPRC (0.319) across the models, and the models with no transformation alone had the next highest average AUPRC of 0.310. Unlike the AUROC results, the random flip augmentation had the worst AUPRC among all the four other augmentation methods.

Discussion

Principal Findings and Comparison With Previous Works

In this retrospective study with over 20,000 hospital admissions, we compared three deep learning computer vision architectures and four image augmentation methods for the early detection of clinical deterioration. We found that the DenseNet121 model pretrained on different publicly available chest radiographs had better discrimination than the VGG16 and Vision Transformer models based on the average AUROC metric. Among different image augmentation methods, a combination of histogram normalization and random Gaussian noise augmentations achieved higher AUROCs and AUPRCs on average than random flip and random rotate transformation. In all of the cases, we found that random flip and random rotate transformation lowered the discrimination compared to the baseline model in terms of both AUROC and AUPRC metrics. To the best of our knowledge, this is the first study to compare different computer vision models and image augmentation methods for predicting clinical deterioration outside the ICU. These findings have important implications in the field of using deep learning models to correctly identify patients showing clinical deterioration and to improve existing EWS applications in health systems.

Although DenseNet121 pretrained on chest radiographs achieved the maximum discrimination with histogram normalization and random Gaussian noise data augmentation, our investigation

```
https://ai.jmir.org/2025/1/e67144
```

found multiple models exhibiting competitive performance across different data augmentation methods considering the AUROC. This may be due to our extensive hyperparameter search with Bayesian optimization that enables all models to achieve similar performances. Overall, the pretrained models performed better with respect to the models trained from scratch. This is consistent with the existing literature, as pretrained models already learned the fundamental building blocks of features (eg, lines and shades) from large number of images of the pretrained dataset [44,45]. However, as the VGG16 model was pretrained on the ImageNet [32] dataset, which is a collection of thousands of general-purpose images, and our dataset only contains chest radiographs, there may be a domain gap present in this scenario that prohibits the maximum benefits of the pretraining network. To analyze and mitigate that domain gap, we compared the performance of the DenseNet121 network pretrained on ImageNet and on a collection of the radiograph dataset. In almost all of the cases, DenseNet121 pretrained on radiographs outperformed the DenseNet121 model pretrained on ImageNet in terms of both the AUROC and AUPRC metrics. These experimental results proved our hypothesis and provided important insights into the use of pretrained networks with chest radiograph datasets. A prior study involving the classification of chest radiographs also found DenseNet networks achieving superior performance [10], which aligns with our findings. For example, Alhudhaif et al found that DenseNet201 achieved the highest discrimination in determining COVID-19 pneumonia with chest radiographs [10]. However, another work by Sitaula et al found that the VGG-16 model performed better than the DenseNet121 model for the classification of COVID-19 chest radiographs [11]. This discrepancy may be explained by differences in hyperparameter settings and the use of pretrained weight initialization. They tuned the hyperparameters manually, whereas we tuned the hyperparameters automatically with Bayesian optimization. As the DenseNet121 network is deeper than VGG-16 in terms of the number of layers, better hyperparameter tuning may enable DenseNet121 to learn more complex relationships without overfitting, hence achieving better performance than the VGG-16 network. Although DenseNet121 has more layers than VGG-16, DenseNet121 has fewer parameters than VGG-16 (7.98M vs 138.36M parameters). This parameter efficiency may reduce the risk of overfitting, which is important in medical imaging applications where datasets are often small. We also found that the Vision Transformer model underperformed in almost all the cases compared to other CNN-based models in the clinical deterioration prediction task. This finding contrasts with the recent success of Vision Transformer in general computer vision tasks [46]. However, in the case of classification tasks with radiographs, the lack of pretraining may harm the performance of the Vision Transformer models [47]. For the networks where we compared performance with random initialization and a pretrained model, in most of the cases, the pretrained model performed better than the randomly initialized one. This could be the main cause for the underperformance of the Vision Transformer models in our work, as we trained it from scratch.

In this study, we found that the models trained with histogram normalization combined with random Gaussian noise among different image augmentation methods achieved better

XSL•FO RenderX
performance, exhibiting the highest AUROC four times and the highest AUPRC three times for different architectures with different combinations of pretraining methods. However, the other two augmentation methods, random flip and random rotate, actually worsened the performance. Our findings align with the existing literature presenting performance improvements with histogram normalization and Gaussian random noise. Gielczyk et al showed that the combination of histogram normalization and Gaussian random noise achieved higher performances than the baseline method in detecting COVID-19 and pneumonia with chest radiographs [48]. However, this can be task dependent involving the useful features of that particular task. Lakhani et al presented a deep convolutional neural network for determining the presence and position of endotracheal tubes where random rotation and random flip augmentation achieved higher performances over the baseline values [12]. As that task was geometry-dependent, regularization introduced by random rotate and random flip augmentation might improve the performance. In contrast, our task of predicting clinical deterioration is not geometry dependent and hence did not benefit from geometric transformations like random rotate and random flips. These insights might be helpful in selecting appropriate image augmentation techniques in models involving chest radiographs.

Strengths

Our study has several important strengths. First, our study cohort consisted of elevated-risk patients with an eCART score ≥93. Predicting deterioration in these patients is more challenging due to their rapid and unpredictable progressions compared to lower-risk patients. Second, we compared multiple deep learning architectures to evaluate their efficacy in predicting clinical deterioration. This comparative approach allows for a more robust understanding of a model's performance in this context. Furthermore, by testing different data augmentation methods, the study explores ways to improve model performance. This

aspect is crucial for enhancing the generalizability and robustness of the models. Incorporating Bayesian optimization with a large search space provides the models to achieve the most optimal performance.

Limitations

Our study also has some limitations. First, we only considered the latest radiograph for our models to avoid bias and complexity. Although we reasoned that the latest radiograph conveys the most updated features of patients, prior radiographs and trends over time might carry important features for the model to predict clinical deterioration. Second, we focused on a few popular deep learning architectures with four different augmentation methods. Although recent studies have introduced numerous computer vision architectures, a more comprehensive study would be difficult considering our study's dataset size. Third, in the deterioration prediction model, we only considered the features on chest radiographs. Incorporating other modalities, such as structured data and clinical notes, could improve the accuracy and robustness of our models and will be an interesting future work. Finally, even though our study is the largest of its kind, this was a single-center study, and future studies in other centers are needed to evaluate the external validity of our models.

Conclusion

Our study demonstrates that the DenseNet121 model pretrained on chest radiographs often outperforms VGG16 and the Vision Transformer model with chest radiographs for the prediction of clinical deterioration. Furthermore, we found that model performance improves with histogram normalization along with random Gaussian noise augmentation in most models in terms of both the AUROC and AUPRC metrics. These results show that accurate prediction of patient clinical deterioration is feasible by utilizing chest radiographs while offering valuable insights into the use of computer vision-aided risk prediction.

Acknowledgments

This work was supported by the National Institute of Health (NIH) National Heart, Lung, and Blood Institute grant number R01HL157262 (MMC), and NIH National Library of Medicine grant number 1R01LM013151 (JWG).

Data Availability

The data used in this study were acquired from The University of Wisconsin health systems following approval from the Institutional Review Board. The data use agreements prohibit sharing data due to regulatory and legal constraints, and therefore, the data cannot be shared publicly.

Authors' Contributions

MMC, MA, DPE, and GC conceptualized the study. MMC, MA, and JWG managed data collection and Institutional Review Board approvals. KAC, AA, and MR led data preprocessing and cleaning. MR and JG were responsible for data analysis and methodology. MR prepared the original draft, and all authors participated in revising and editing the article.

Conflicts of Interest

MCM and DPE have a patent issued (#11,410,777) for risk stratification algorithms for hospitalized patients. DPE is employed by and has an equity interest in AgileMD, San Francisco, CA. The other authors have declared no potential conflicts of interest.

Multimedia Appendix 1



Supplementary tables and figures. [DOCX File, 42 KB - ai v4i1e67144 app1.docx]

References

- Padilla RM, Mayo AM. Clinical deterioration: a concept analysis. J Clin Nurs 2018 Apr;27(7-8):1360-1368. [doi: 10.1111/jocn.14238] [Medline: 29266536]
- 2. Vincent JL, Einav S, Pearse R, et al. Improving detection of patient deterioration in the general hospital ward environment. Eur J Anaesthesiol 2018 May;35(5):325-333. [doi: 10.1097/EJA.000000000000798] [Medline: 29474347]
- 3. U.S. food and drug administration. ECARTv5 clinical deterioration suite. URL: <u>https://www.accessdata.fda.gov/cdrh_docs/pdf23/K233253.pdf</u> [accessed 2025-03-25]
- 4. Subbe CP, Kruger M, Rutherford P, Gemmel L. Validation of a modified Early Warning Score in medical admissions. QJM 2001 Oct;94(10):521-526. [doi: 10.1093/qjmed/94.10.521] [Medline: 11588210]
- 5. Prytherch DR, Smith GB, Schmidt PE, Featherstone PI. ViEWS--Towards a national early warning score for detecting adult inpatient deterioration. Resuscitation 2010 Aug;81(8):932-937. [doi: 10.1016/j.resuscitation.2010.04.014] [Medline: 20637974]
- 6. Smith GB, Prytherch DR, Meredith P, Schmidt PE, Featherstone PI. The ability of the National Early Warning Score (NEWS) to discriminate patients at risk of early cardiac arrest, unanticipated intensive care unit admission, and death. Resuscitation 2013 Apr;84(4):465-470. [doi: 10.1016/j.resuscitation.2012.12.016] [Medline: 23295778]
- 7. Fang AHS, Lim WT, Balakrishnan T. Early warning score validation methodologies and performance metrics: a systematic review. BMC Med Inform Decis Mak 2020 Jun 18;20(1):111. [doi: 10.1186/s12911-020-01144-8] [Medline: 32552702]
- Xia C, et al. A multi-modality network for cardiomyopathy death risk prediction with CMR images and clinical information. Medical Image Computing and Computer Assisted Intervention – MICCAI 2019 2019;11765:577-585. [doi: 10.1007/978-3-030-32245-8_64]
- Rawal K, Sethi G, Walia GK. Impact of machine learning and deep learning in medical image analysis. In: Rabie K, Karthik C, Chowdhury S, Dutta PK, editors. Deep Learning in Medical Image Processing and Analysis 2023:187-199. [doi: 10.1049/PBHE059E_ch11]
- Alhudhaif A, Polat K, Karaman O. Determination of COVID-19 pneumonia based on generalized convolutional neural network model from chest X-ray images. Expert Syst Appl 2021 Oct 15;180:115141. [doi: <u>10.1016/j.eswa.2021.115141</u>] [Medline: <u>33967405</u>]
- 11. Sitaula C, Hossain MB. Attention-based VGG-16 model for COVID-19 chest X-ray image classification. Appl Intell (Dordr) 2021;51(5):2850-2863. [doi: 10.1007/s10489-020-02055-x] [Medline: 34764568]
- 12. Lakhani P. Deep convolutional neural networks for endotracheal tube position and X-ray image classification: challenges and opportunities. J Digit Imaging 2017 Aug;30(4):460-468. [doi: 10.1007/s10278-017-9980-7] [Medline: 28600640]
- 13. Al-qaness MAA, Zhu J, AL-Alimi D, et al. Chest X-ray images for lung disease detection using deep learning techniques: a comprehensive survey. Arch Computat Methods Eng 2024 Aug;31(6):3267-3301. [doi: 10.1007/s11831-024-10081-y]
- Alshmrani GMM, Ni Q, Jiang R, Pervaiz H, Elshennawy NM. A deep learning architecture for multi-class lung diseases classification using chest X-ray (CXR) images. Alexandria Engineering Journal 2023 Feb;64:923-935. [doi: 10.1016/j.aej.2022.10.053]
- Sjoding MW, Taylor D, Motyka J, et al. Deep learning to detect acute respiratory distress syndrome on chest radiographs: a retrospective study with external validation. Lancet Digit Health 2021 Jun;3(6):e340-e348. [doi: <u>10.1016/S2589-7500(21)00056-X]</u> [Medline: <u>33893070</u>]
- 16. Sharma S, Guleria K. A systematic literature review on deep learning approaches for pneumonia detection using chest X-ray images. Multimed Tools Appl 2023 Aug;83(8):24101-24151. [doi: 10.1007/s11042-023-16419-1]
- 17. Ibrahim AU, Ozsoz M, Serte S, Al-Turjman F, Yakoi PS. Pneumonia classification using deep learning from chest X-ray images during COVID-19. Cognit Comput 2021 Jan 4;16(4):1-13. [doi: 10.1007/s12559-020-09787-5] [Medline: 33425044]
- 18. Vats S, Sharma V, Singh K, et al. Incremental learning-based cascaded model for detection and localization of tuberculosis from chest x-ray images. Expert Syst Appl 2024 Mar;238:122129. [doi: 10.1016/j.eswa.2023.122129]
- 19. Sharma V, Gupta SK, Shukla KK. Deep learning models for tuberculosis detection and infected region visualization in chest X-ray images. Intelligent Medicine 2024 May;4(2):104-113. [doi: <u>10.1016/j.imed.2023.06.001</u>]
- 20. Sunnetci KM, Alkan A. Biphasic majority voting-based comparative COVID-19 diagnosis using chest X-ray images. Expert Syst Appl 2023 Apr 15;216:119430. [doi: 10.1016/j.eswa.2022.119430] [Medline: 36570382]
- 21. Chen Z, et al. A vision-language foundation model to enhance efficiency of chest X-ray interpretation. arXiv. Preprint posted online on 2024 URL: <u>http://arxiv.org/abs/2401.12208</u> [accessed 2024-09-25] [doi: <u>10.48550/arXiv.2401.12208</u>]
- Cid YD, Macpherson M, Gervais-Andre L, et al. Development and validation of open-source deep neural networks for comprehensive chest x-ray reading: a retrospective, multicentre study. Lancet Digit Health 2024 Jan;6(1):e44-e57. [doi: 10.1016/S2589-7500(23)00218-2] [Medline: 38071118]
- 23. Meedeniya D, Kumarasinghe H, Kolonne S, Fernando C, Díez ILT, Marques G. Chest X-ray analysis empowered with deep learning: a systematic review. Appl Soft Comput 2022 Sep;126:109319. [doi: 10.1016/j.asoc.2022.109319] [Medline: 36034154]

- 24. Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition. arXiv. Preprint posted online on 2015 URL: <u>http://arxiv.org/abs/1409.1556</u> [accessed 2024-09-24] [doi: <u>10.48550/arXiv.1409.1556</u>]
- Huang G, Liu Z, Van Der Maaten L, Weinberger KQ. Densely connected convolutional networks. Presented at: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR); Jul 21-26, 2017; Honolulu, HI p. 2261-2269. [doi: 10.1109/CVPR.2017.243]
- 26. Dosovitskiy A, et al. An image is worth 16x16 words: transformers for image recognition at scale. arXiv. Preprint posted online on 2021 URL: <u>http://arxiv.org/abs/2010.11929</u> [accessed 2025-03-25] [doi: <u>10.48550/arXiv.2010.11929</u>]
- 27. He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. arXiv. Preprint posted online on 2015 URL: https://arxiv.org/abs/1512.03385 [accessed 2025-03-25]
- 28. Szegedy C, Vanhoucke V, Ioffe S, Shlens J, Wojna Z. Rethinking the inception architecture for computer vision. 2016 Presented at: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR); Jun 26 to Jul 1, 2016; Las Vegas, NV, USA. [doi: <u>10.1109/CVPR.2016.308</u>]
- 29. Ahmed F, Abbas S, Athar A, et al. Identification of kidney stones in KUB X-ray images using VGG16 empowered with explainable artificial intelligence. Sci Rep 2024 Mar 14;14(1):6173. [doi: <u>10.1038/s41598-024-56478-4</u>] [Medline: <u>38486010</u>]
- Islam M, Hannan T, Sarker L, Ahmed Z. COVID-densenet: A deep learning architecture to detect COVID-19 from chest radiology images. In: Saraswat M, Chowdhury C, Mandal CK, Gandomi AH, editors. Presented at: Proceedings of International Conference on Data Science and Applications; Feb 7, 2023 p. 397-415. [doi: 10.1007/978-981-19-6634-7_28]
- Duong LT, Le NH, Tran TB, Ngo VM, Nguyen PT. Detection of tuberculosis from chest X-ray images: Boosting the performance with vision transformer and transfer learning. Expert Syst Appl 2021 Dec;184:115519. [doi: 10.1016/j.eswa.2021.115519]
- 32. Deng J, Dong W, Socher R, Li LJ. ImageNet: A large-scale hierarchical image database. Presented at: 2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops (CVPR Workshops); Jun 20-25, 2005; Miami, FL p. 248-255. [doi: 10.1109/CVPR.2009.5206848]
- Cohen JP, et al. TorchXRayVision: A library of chest X-ray datasets and models. arXiv. Preprint posted online on 2021. [doi: <u>10.48550/ARXIV.2111.00595</u>]
- Wang X, Peng Y, Lu L, Lu Z, Bagheri M, Summers RM. ChestX-ray8: hospital-scale chest X-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. Presented at: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR); Jul 21-26, 2017; Honolulu, HI p. 3462-3471. [doi: 10.1109/CVPR.2017.369]
- 35. Bustos A, Pertusa A, Salinas JM, de la Iglesia-Vayá M. PadChest: A large chest x-ray image dataset with multi-label annotated reports. Med Image Anal 2020 Dec;66:101797. [doi: 10.1016/j.media.2020.101797] [Medline: 32877839]
- 36. Irvin J, Rajpurkar P, Ko M, et al. CheXpert: a large chest radiograph dataset with uncertainty labels and expert comparison. . 2019(1) p. 590-597. [doi: <u>10.1609/aaai.v33i01.3301590</u>]
- 37. Johnson AEW, Pollard TJ, Berkowitz SJ, et al. MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports. Sci Data 2019 Dec 12;6(1):317. [doi: 10.1038/s41597-019-0322-0] [Medline: 31831740]
- 38. Demner-Fushman D, Kohli MD, Rosenman MB, et al. Preparing a collection of radiology examinations for distribution and retrieval. J Am Med Inform Assoc 2016 Mar;23(2):304-310. [doi: <u>10.1093/jamia/ocv080</u>] [Medline: <u>26133894</u>]
- Majkowska A, Mittal S, Steiner DF, et al. Chest radiograph interpretation with deep learning models: assessment with radiologist-adjudicated reference standards and population-adjusted evaluation. Radiology 2020 Feb;294(2):421-431. [doi: <u>10.1148/radiol.2019191293</u>] [Medline: <u>31793848</u>]
- 40. RSNA pneumonia detection challenge. URL: <u>https://www.kaggle.com/competitions/rsna-pneumonia-detection-challenge</u> [accessed 2025-03-25]
- 41. Liaw R, Liang E, Nishihara R, Moritz P, Gonzalez JE, Stoica I. Tune: A research platform for distributed model selection and training. arXiv. Preprint posted online on 2018. [doi: <u>10.48550/arXiv.1807.05118</u>]
- 42. DeLong ER, DeLong DM, Clarke-Pearson DL. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. Biometrics 1988 Sep;44(3):837-845. [doi: <u>10.2307/2531595</u>] [Medline: <u>3203132</u>]
- 43. Netflix. GitHub Netflix/vmaf: Perceptual video quality assessment based on multi-method fusion. URL: <u>https://github.</u> <u>com/Netflix/vmaf/</u> [accessed 2025-03-25]
- 44. He K, Girshick R, Dollar P. Rethinking imagenet pre-training. Presented at: 2019 IEEE/CVF International Conference on Computer Vision (ICCV); Oct 27 to Nov 2, 2019; Seoul, Korea (South) p. 4917-4926 URL: <u>https://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=8972782</u> [doi: <u>10.1109/ICCV.2019.00502</u>]
- 45. Lee J, Lee EJ. Self-supervised pre-training improves fundus image classification for diabetic retinopathy. In: Kehtarnavaz N, Carlsohn MF, editors. Presented at: Real-Time Image Processing and Deep Learning 2022; Apr 3-7, 2022; Orlando, United States. [doi: 10.1117/12.2632901]
- 46. Khan S, Naseer M, Hayat M, Zamir SW, Khan FS, Shah M. Transformers in vision: a survey. ACM Comput Surv 2022 Jan 31;54(10s):1-41. [doi: 10.1145/3505244]
- 47. Jain A, Bhardwaj A, Murali K, Surani I. A comparative study of CNN, ResNet, and Vision Transformers for multi-classification of chest diseases. arXiv 2024. [doi: 10.48550/arXiv.2406.00237]

48. Giełczyk A, Marciniak A, Tarczewska M, Lutowski Z. Pre-processing methods in chest X-ray image classification. PLoS ONE 2022;17(4):e0265949. [doi: 10.1371/journal.pone.0265949] [Medline: 35381050]

Abbreviations

AUPRC: area under the precision-recall curve AUROC: area under the receiver operating characteristic curve eCART: electronic cardiac arrest risk triage EWS: early warning scores ICU: intensive care unit UW Health: University of Wisconsin Health System VMAF: Video Multimethod Assessment Fusion

Edited by Y Huo; submitted 03.10.24; peer-reviewed by S Wang, S Lu; revised version received 08.03.25; accepted 10.03.25; published 10.04.25.

<u>Please cite as:</u>

Rahman M, Gao J, Carey KA, Edelson DP, Afshar A, Garrett JW, Chen G, Afshar M, Churpek MM Comparison of Deep Learning Approaches Using Chest Radiographs for Predicting Clinical Deterioration: Retrospective Observational Study JMIR AI 2025;4:e67144 URL: https://ai.jmir.org/2025/1/e67144 doi:10.2196/67144

© Mahmudur Rahman, Jifan Gao, Kyle A Carey, Dana P Edelson, Askar Afshar, John W Garrett, Guanhua Chen, Majid Afshar, Matthew M Churpek. Originally published in JMIR AI (https://ai.jmir.org), 10.4.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Large Language Models for Thematic Summarization in Qualitative Health Care Research: Comparative Analysis of Model and Human Performance

Arturo Castellanos^{1*}, PhD; Haoqiang Jiang^{2*}, PhD; Paulo Gomes^{3*}, PhD; Debra Vander Meer^{3*}, PhD; Alfred Castillo³, PhD

¹Mason School of Business, William & Mary, Williamsburg, VA, United States

²College of Informatics, Northern Kentucky University, Highland Heights, KY, United States

³Information Systems and Business Analytics Department, College of Business, Florida International University, 11200 SW 8th Street, Miami, FL, United States

*these authors contributed equally

Corresponding Author:

Paulo Gomes, PhD

Information Systems and Business Analytics Department, College of Business, Florida International University, 11200 SW 8th Street, Miami, FL, United States

Abstract

Background: The application of large language models (LLMs) in analyzing expert textual online data is a topic of growing importance in computational linguistics and qualitative research within health care settings.

Objective: The objective of this study was to understand how LLMs can help analyze expert textual data. Topic modeling enables scaling the thematic analysis of content of a large corpus of data, but it still requires interpretation. We investigate the use of LLMs to help researchers scale this interpretation.

Methods: The primary methodological phases of this project were (1) collecting data representing posts to an online nurse forum, as well as cleaning and preprocessing the data; (2) using latent Dirichlet allocation (LDA) to derive topics; (3) using human categorization for topic modeling; and (4) using LLMs to complement and scale the interpretation of thematic analysis. The purpose is to compare the outcomes of human interpretation with those derived from LLMs.

Results: There is substantial agreement (247/310, 80%) between LLM and human interpretation. For two-thirds of the topics, human evaluation and LLMs agree on alignment and convergence of themes. Furthermore, LLM subthemes offer depth of analysis within LDA topics, providing detailed explanations that align with and build upon established human themes. Nonetheless, LLMs identify coherence and complementarity where human evaluation does not.

Conclusions: LLMs enable the automation of the interpretation task in qualitative research. There are challenges in the use of LLMs for evaluation of the resulting themes.

(JMIR AI 2025;4:e64447) doi:10.2196/64447

KEYWORDS

artificial intelligence; generative AI; large language models; ChatGPT; machine learning; health care

Introduction

Background

Qualitative studies in health care shed light on the perceptions, narratives, and discourses that underlie human behavior. This approach enhances understanding of both clinicians and patients' experiences and expectations, thereby informing decision-making for health policy [1]. Traditionally, these studies involved data collection through face-to-face interviews, observation or artifact analysis, transcription, and manual human coding for sense-making. Recent online advances, such as social media interactions, online reviews, news articles, and in-depth

```
https://ai.jmir.org/2025/1/e64447
```

forum discussions, allow researchers and policy makers to collect larger data samples at lower time costs compared with direct interviews [2]. The advent of text mining tools, which allow researchers to cluster text samples into groups based on statistical similarity, has enabled partial automation of the sense-making step. For instance, the use of natural language processing (NLP) to identify risk factors from unstructured free-text clinical notes [3]. Yet, these tools provide only the groupings, leaving the human to apply thematic interpretation [4,5].

Recent advances in generative artificial intelligence (AI) provide valuable tools for researchers conducting qualitative studies,

offering support in both data analysis and interpretation. In particular, large language models (LLMs), which are statistical models built using internet-scale datasets, can generate human-style writing in response to natural-language prompts, and assist in analyzing textual data to identify patterns, themes, and underlying meanings [6]. LLMs can aid researchers in conducting thematic analysis by identifying recurrent themes, concepts, or ideas across a dataset supporting the automation of thematic interpretation.

Previous Work

Topic modeling is a popular approach to uncovering insights in text mining. It identifies patterns in word usage and clusters words into topics, making it a popular method for exploring large, unstructured text datasets. Latent Dirichlet allocation (LDA) is a widely applied method for topic modeling. Previous work has used LDA modeling to analyze social media data and derive insights on key topics [4,7,8]. Despite the new perspectives LDA approaches offer for scientific research [9], using LDA for topic modeling presents challenges [10], notably the significant role of human interpretative contribution in the process [11], which limits scalability. In addition, there is a noted lack of user-friendly tools that support the entire workflow, necessitating a human-in-the-loop to interpret the derived topics. In this paper, we argue that LLMs can help resolve some of the challenges of LDA analysis, specifically in interpreting and labeling topics.

LLMs are emerging as an increasingly reliable and effective tool for interpretative qualitative research, combining the scale

Table . Challenges of large language models.

that computational techniques allow for with the human's qualitative logic [12,13]. Previous studies show that ChatGPT (OpenAI) yields comparable results to manual coding with substantial time savings [14]. These studies compare emergent themes in human and AI-generated qualitative analyses, revealing similarities and differences. For instance, some themes are recognized by human coders but missed by ChatGPT, and vice versa [15]. LLMs can highlight novel connections within the data that are not apparent to human coders. In both deductive and inductive thematic analysis, ChatGPT extended the researchers' views of the themes present in the data [12].

There are challenges associated with the use of LLMs. In the previously cited study [14], ChatGPT was able to recreate the themes originally identified through more traditional methods. However, it was less successful at identifying subtle, interpretive themes, and more successful with concrete, descriptive themes. LLMs may miss themes that require a deep understanding of context or specific domain knowledge. For example, themes related to niche cultural practices or specific professional areas may not be accurately identified by AI without targeted training.

LLMs can also reflect biases present in its training data, potentially overlooking or misinterpreting themes that deviate from its learned patterns. On the other hand, LLM analyses can identify patterns and themes that might be overlooked by human coders due to their preconceived notions or cognitive biases. Further challenges associated with the use of LLMs are shown in Table 1.

Challenge	Description	Citations
Ambiguity resolution	LLMs ^a might struggle to disambiguate certain terms or topics, leading to unclear topic catego-rization.	[16,17]
Overfitting	LLMs can sometimes focus too much on com- mon or popular topics, missing out on niche or less frequently discussed topics.	[18,19]
Lack of context	Without external knowledge or the ability to track long-term context, LLMs might misinter- pret or miss certain topic nuances.	[20]
Bias	LLMs are trained on vast amounts of data, which may contain biases. This can affect topic analysis results.	[21,22]
Overgeneralization	LLMs might overly generalize topics, missing out on specific subtopics or nuances.	[23]
Sensitivity to input	Small changes in input phrasing can sometimes lead to different topic interpretations by the LLM.	[24]
Memory limitations	Due to token limits, LLMs might not capture very long or detailed discussions effectively for topic analysis.	[25]
Interactivity limitations	While LLMs can process static text effectively, they might struggle with dynamic topic analysis, where user feedback or real-time adjustments are required.	[26]

^aLLM: large language model.



Given these challenges, some studies suggest that the most effective qualitative analyses may involve a combination of human and AI insights, as human coders often recognize nuanced themes related to context, emotions, and cultural subtleties that AI may miss. For example, a study demonstrates the feasibility of AI as a research assistant, presenting a successful case-study of human-AI collaboration in research by merging the efficiency and accuracy of ChatGPT with human flexibility and context awareness [27]. In addition, the usefulness of ChatGPT in qualitative analysis may depend on the researcher's ability to ask appropriate questions (prompts), with the output evaluated and supplemented by a human researcher before the final report and publication [28].

There is little guidance in the literature about how LLMs can be integrated into thematic analysis. Challenges associated with the use of LLMs, including overgeneralization and overfitting, need to be investigated in the context of using LLMs for interpreting the relevance of identified topics. Our focus in this work considers inductive thematic analysis, where themes are derived from data without preconceived frameworks, and semantic analysis, in which themes are identified within the explicit content of the data [29]. We plan to consider a hybrid inductive and deductive approach in future work [30].

Study Objectives

This study considers the possibility of enhancing human productivity by applying LLMs in the interpretation and labeling stage of topic modeling. We present a case study in which data were gathered from an online forum and grouped using text mining tools, and then interpreted for themes in parallel: (1) by human coders and (2) by providing text samples from each classification group to an LLM and prompting the LLM for thematic summarization.

We compared the human- and LLM-generated themes along 4 qualitative dimensions: alignment, convergence, coherence, and complementarity. Based on this analysis, we demonstrate the feasibility of using an LLM to support human thematic interpretation for qualitative research and offer insights into where researchers may find benefit in using LLMs to support thematic interpretation, and where they should exercise caution.

Methods

Overview

The proposed methodology is based on three phases: (1) construction of a dataset and topic modeling using LDA, (2) labeling identified groups into topics through human interpretation and through use of LLM, and (3) comparison of identified topics.

Data Collection and Preprocessing

The data comprises discussions from a publicly accessible Nurse Forum [4]. Data come from posts aggregated over 28 2-week periods from March 2020 to April 2021. Our preprocessing approach ensures that the data is clean, standardized, and focused on the most relevant linguistic features, allowing for a clearer identification of the key aspects discussed in the nurse forum over time. Texts were tokenized using the Python library

```
https://ai.jmir.org/2025/1/e64447
```

Gensim [31]. Preprocessing included lowercasing and removing punctuation to ensure uniformity and reduce noise in the text. Stop words, including domain-specific terms like "covid" and "covid 19," were removed, in addition to those in the Natural Language Toolkit (NLTK) library, to focus on meaningful content. Bigram and trigrams were added to the corpora to identify common multiword expressions, which enhances the detection of contextually significant phrases. Finally, texts were lemmatized using SpaCy (Explosion) [32], retaining only nouns, adjectives, verbs, and adverbs, to normalize words to their base forms and reduce dimensionality.

Topic Modeling

Topic modeling was conducted using LDA to identify underlying themes in the text data. The LDA algorithm began with random assignments of topics to documents and words to topics. Through iterative optimization, it adjusts these assignments based on the likelihood of word-topic and topic-document distributions. We experimented with different numbers of topics and adjusted hyperparameters, to find the optimal model configuration. Coherence scores, which measure the semantic similarity of words within a topic, were computed for each run. Higher coherence scores indicate more meaningful and interpretable topics. The model with the highest coherence score was selected [33].

This optimal model is then used to extract the top keywords for each topic, summarizing the themes present in the data. The distribution of topics across the corpus was visualized to interpret their prevalence in individual documents and the entire dataset, providing insights into the prominent themes discussed in the nurse forum during the specified period.

Identification of Topics Through Human Interpretation

Thematic analysis was conducted by 2 coders working independently to familiarize themselves with the data by exhaustively reading the top 10 posts within each topic (ranked based on coherence scores) generated by the topic models [34]. The selected theme names for the labeled topics were compared, which achieved an initial interannotator agreement of 68% (210/310), and 94% (292/310) after a subsequent round. For the remaining 6% (18/310), the underlying posts were examined together to resolve the disagreements, which left no unresolved annotations. The interpretation analysis resulted in 16.5% (15/310) of the identified themes being categorized as having low coherence.

Theme Derivation Using Large Language Models

Following topic modeling, an LLM was used to derive themes from the identified topics. We created a custom function that takes a system message and a list of user-assistant message pairs, ensuring proper formatting and role assignment. We use the GPT-3.5 based model, specifying the structured messages, temperature, and seed for reproducibility [35]. The system prompt is embedded ensuring consistency in use of the associated set of instructions. The model was chosen for its advanced NLP capabilities, including context-awareness and adaptability to specific thematic contexts, and accessibility to the research team. The prompt instructs ChatGPT to generate

themes and subthemes for more nuanced theme identification, addressing the issue of overly broad categorizations observed in initial experiments. An overview of the modeling steps is provided in Multimedia Appendix 1.

Comparison of Identified Topics

The reliability of coding textual data can be challenging as the goal in content analysis is to attain a "scientific" analysis characterized by reliability, which implies stability in the phenomenon being studied and explicit analytic procedures to ensure that any reasonably qualified person would yield identical results [36]. Intercoder agreement emerges as a key tool in achieving a reliable coding scheme, assessing the extent to which coders assign identical codes to the same set of data [34]. A 5-item ordinal scale typically measures this agreement, with the anchors of "Perfect Agreement," representing where coders completely agree on codes or categories assigned to data, and "Slight Agreement," representing very little consensus, or

significant disagreement, among the coders in how they code the data. This agreement scale is provided in Multimedia Appendix 1.

A novel 7-point scale was developed following a pilot test conducted by two of the authors to address the complexities of comparing codes generated by humans and ChatGPT. This scale, presented in the first column of Table 2, focuses on exploring the complementary and divergent insights between human-generated and ChatGPT-generated codes. It emphasizes the value of examining differences, especially in cases of low coherence among human-coded data, which allows researchers to uncover nuanced perspectives and understandings contained in ChatGPT-generated themes and within subthemes variability, with the possibility of revealing new and meaningful insights. It serves as a dynamic tool that stresses the importance of learning from intercoding differences rather than seeking strict agreement and validation, as is valued among qualitative researchers [37].

Table . Agreement between large language model (LLM) and human coding.

Agreement scale	Number of topics, n	Rate of agreement, %
ChatGPT and human coding themes are aligned, coders largely interpret and code the data in a consistent manner.	95	30.6
Substantial agreement: ChatGPT's subthemes are aligned with human coding, some subthemes provide complementary perspectives or unique insights.	101	32.6
Substantial agreement: ChatGPT's themes are divergent, human coding classified as low coherence.	51	16.5
Moderate agreement: there is a reasonable level of consensus between ChatGPT and human coding, but there are significant differences in interpretation or coding for some subthemes.	15	0
Fair agreement: ChatGPT's themes are considered too broad, there are substantial discrepancies between ChatGPT's subthemes compared with human coding.	30	0
Poor agreement: ChatGPT's theme specific, yet divergent from human coding.	4	0
Poor agreement: ChatGPT's theme specific, yet low coherence in human coding.	14	0
Grand total	310	79.7

We then use GPT-4 for topic comparison, accessing the ChatGPT engine through an application programming interface (API) for programmatic purpose. Each prompt included the human-coded themes and the LLM-generated themes, requesting the LLM to assess the agreement based on 4 criteria: alignment, convergence, coherence, and complementarity between the themes. A detailed overview of the prompts is provided in Multimedia Appendix 1.

Alignment assesses the correspondence between ChatGPT and human themes in terms of contextual agreement between the themes [38], rather than lexical agreement. Convergence provides a similar comparison at the level of specific "ChatGPT

https://ai.jmir.org/2025/1/e64447

Subthemes" with reference to the "Human Theme." Coherence evaluates the logical consistency within the "ChatGPT Theme" and its subthemes, emphasizing the cohesion in both logic and meaning [39]. Complementarity looks at whether the ChatGPT subthemes offer valuable additional insights or perspectives that enhance the human theme by providing detailed mechanistic explanations that align with and build upon the established human theme without contradicting it [40,41],

LLM outputs were parsed to extract values for alignment, coherence, convergence, and complementarity. Human coders then compared the remaining results of reliability analysis with the LLM-generated comparison.

Ethical Considerations

This study does not involve human subjects, identifiable private information, or direct interactions with individuals. Instead, it relies exclusively on publicly available, anonymized social media posts. Consequently, institutional review board approval was deemed unnecessary.

Results

Analysis of Reliability

The LDA analysis identified 310 topics. In thematic analysis, the team considered the topics identified, groups of words, and representative blog post samples in each topic and categorized the 310 topics into 58 subthemes.

A total of 2 authors independently classified the level of agreement on each topic against the themes and subthemes generated by ChatGPT, using the 7-point agreement scale in Table 2. The authors then met to compare assessments and resolve disagreements. The overall reliability is estimated at 79.7% (247/310), which represents substantial agreement according to the intercoder reliability benchmark [36].

Table 2 provides a breakdown of agreement along the comparison scale, with 30.6% (95/310) reflecting taxonomic

agreement in themes identified by the human coder and ChatGPT. For example, in one case the human-coder's theme is "PPE resource availability and control" and the ChatGPT theme is "Mask Availability and Usage in Healthcare Settings."

In 32.6% (101/310) of the themes the agreement is at the subtheme level. For example, in one instance the human-coded theme is "Testing policies in different settings," while the ChatGPT theme is "Challenges and Controversies Surrounding COVID-19 Testing," which was not considered at the same level of specificity of the human coder's theme. The ChatGPT subthemes are "Allocation of Testing Resources," "Flaws in Testing Systems," and "Impact on Public Health and Society." In the first subtheme the discussion revolves around whether COVID-19 tests should be prioritized for hospitalized patients or health care workers, matching the theme identified by humans. Adding to the reliability of the method, we have the agreement on the lack of coherence of the posts included in the LLM topic, representing 16.5% (51/310) of the topics.

Alignment and Convergence

LLM provided results on alignment and convergence that we compare with the human evaluation of agreement. The results are displayed in Table 3.



Table . Analysis of alignment (theme level) and convergence (subtheme level).

		Alignment: Compare the "human theme" and the "ChatGPT theme"		Convergence: Compare the specifics in "ChatGPT Subtheme" with the "human theme."		in "ChatGPT	
Agreement scale	Total, n	Aligned, n	Misaligned, n	Meets expecta- tion, %	Convergent, n	Divergent, n	Meets expecta- tion, %
ChatGPT and human coding themes are aligned, coders largely interpret and code the da- ta in a consistent manner.	95	86 ^a	9	91	86 ^b	9	91
Substantial agreement: ChatGPT's sub- themes are aligned with hu- man coding, some subthemes provide comple- mentary perspec- tives or unique insights.	101	90 ^a	11	89	91 ^b	10	90
Substantial agreement: ChatGPT's themes are diver- gent, human coding classified as low coher- ence.	51	6	45 ^c	88	5	36 ^d	71
Moderate agree- ment: there is a reasonable level of consensus be- tween ChatGPT and human cod- ing, but there are significant differ- ences in interpre- tation or coding for some sub- themes.	15	10 ^a	5	67	11 ^b	4	73
Fair agreement: ChatGPT's themes are con- sidered too broad, there are substantial dis- crepancies be- tween ChatGPT subthemes com- pared with hu- man coding.	30	18	12	0	17	13 ^d	43
Poor agreement: ChatGPT's theme specific, yet divergent from human coding.	4	1	3 [°]	75	1	2	0



		Alignment: Compare the "human theme" and the "ChatGPT theme"		Convergence: Compare the specifics in "ChatGPT Subtheme" with the "human theme."		in "ChatGPT	
Agreement scale	Total, n	Aligned, n	Misaligned, n	Meets expecta- tion, %	Convergent, n	Divergent, n	Meets expecta- tion, %
Poor agreement: ChatGPT's theme specific, yet low coher- ence in human coding.	14	2	12 ^c	86	1	8 ^d	57

^aExpectation is "aligned" for items in the agreement scale.

^bExpectation is "convergent" in the agreement scale.

^cExpectation is "misaligned" in the agreement scale.

^dExpectation is "divergent" in the agreement scale.

We found high level of alignment and convergence for themes classified as high on agreement by human coder. For scale item 1 there was 91% (86/95) alignment and 91% (86/95) convergence, and for scale item 2, there was 89% (90/101) alignment and 90% (91/101) convergence. As expected, we find misalignment for scale items 3 and 7.

The results for scale item 5 (ChatGPT's themes are considered too broad, there are substantial discrepancies between ChatGPT subthemes compared with human coding) reveal specific nuances of the LLM comparison. Although we expect subthemes to be divergent based on human classification, only 43% (13/30) were classified as divergent by the LLM. For example, a topic labeled by human-coders as "Knowledge about virus," due to posts in general discuss the nature of COVID-19, was labeled by LLM as "COVID-19 and its implications for healthcare workers," which is considered much broader although aligned. However, the first subtheme, "Understanding the nature of coronaviruses and COVID-19" is both aligned and convergent

with human-generated theme while the other two subthemes, "Importance of proper PPE and testing for healthcare workers" and "Concerns and challenges in healthcare settings and home care," are clearly divergent from the narrow scope defined by human-theme. Although some subthemes may be tangential, the LLM still classifies them as convergent within a broader framework of idea similarity.

Coherence

Coherence evaluates the logical consistency within the "ChatGPT Theme" and its subthemes. The results are displayed in Table 4. Coherence was high for items 1,2, and 4 in the agreement scale, meeting expectations. We expected coherence to be low for scale item 5. However, contrary to our expectations, ChatGPT identified 97% (29/30) of cases as coherent. Although human interpretation viewed the LLM theme as broad and the subthemes as tangential, the LLM found logical consistency among these items within the broader scope of the theme.



Table . Analysis of coherence and complementarity.

	Analysis of coherence			Analysis of complementarity		
	Coherent, n	Low coherence, n	Meets expectation, %	Complementary, n	Meets expectation, %	
ChatGPT and human coding themes are aligned, coders largely interpret and code the data in a consistent manner.	94 ^a	1	99	93 ^b	98	
Substantial agreement: ChatGPT's subthemes are aligned with human coding, some sub- themes provide comple- mentary perspectives or unique insights.	101 ^a	0	100	97 ^b	96	
Substantial agreement: ChatGPT's themes are divergent, human cod- ing classified as low coherence.	49	2 ^c	4	8	0	
Moderate agreement: there is a reasonable level of consensus be- tween ChatGPT and human coding, but there are significant differences in interpre- tation or coding for some subthemes.	15 ^a	0	100	15 ^b	100	
Fair agreement: ChatG- PT's themes are consid- ered too broad, there are substantial discrep- ancies between ChatG- PT subthemes com- pared with human cod- ing.	29	1 ^c	3	26	87	
Poor agreement: Chat- GPT's theme specific, yet divergent from hu- man coding	3	1	0	1	0	
Poor agreement: Chat- GPT's theme specific, yet low coherence in human coding	14	0	0	4	0	

^aExpectation is "coherent" in the agreement scale.

^bExpectation is "complementary" in the agreement scale.

^cExpectation is "low coherence" in the agreement scale.

Another unexpected result concerns scale item 3, where 96% (49/51) of the topics were marked as coherent despite being rated as "low coherence" by human coders. Contrary to expectations, 49 out of 51 cases were classified as coherent. LLM relies on single posts to generate subthemes with logical consistency. We illustrate this finding with 2 examples.

One topic that ChatGPT themed as "Nurses' Safety and Well-being" with the subthemes of "Personal sacrifices and concerns for personal safety," and "Need for better protection and compensation." However, the second subtheme was

https://ai.jmir.org/2025/1/e64447

RenderX

generated based on a single post that mentions hazard pay: "It would be nice if hospitals offered hazard pay, but I'm sure they're also hurting financially given all of the new measures they're having to put into place. [...]; many are losing a lot of anticipated revenue because they've canceled their non-emergency surgeries." There is insufficient evidence to support the inclusion of this theme.

Another topic ChatGPT themed as "Challenges and Considerations in Nursing and Healthcare" with the subthemes of "Trust and Distrust in Healthcare" and "Disparities in

Healthcare." Although these are considered consistent with theme, the first subtheme is based on a post highlighting the impact of past negative experiences on trust, and the second subtheme is described by ChatGPT as emphasizing the importance of recognizing and addressing disparities that affect various groups, such as gender, age, ethnicity, and socioeconomic status; however, it is based on the following post: "There are disparities... People we love and care about. Yes, I think it's important to identify areas that are of particular concern and groups that are especially vulnerable. We need to learn and use that knowledge to try to improve our collective future."

A total of 2 topics were classified as low coherence, which agreed with the corresponding "low coherence" human theme designation. ChatGPT themed 1 topic as "Medications and Health Concerns" with the subthemes of "Medication Switch and COVID-19," "Casual Conversations and Expressions," and "Concern and Well-Wishes for Health," yet recognized as low coherence. The second topic, ChatGPT themed as "Controversial Issues in Healthcare" and has subthemes of "Use of Hydroxychloroquine for COVID-19 Treatment," "Systemic Racism and Police Brutality," and "Challenges in Ensuring Compliance with Infection Control Measures."

Complementarity

The analysis of complementarity is also provided in Table 4. For scale items 1, 2, and 4 the expectation was that the subthemes provide complementarity to the human-generated theme and the results meet expectations (98% (93/94), 96% (97/101), and 100% (15/15), respectively). For example, one topic with the human-generated theme of "Testing policies in different settings" was associated with the ChatGPT subthemes of "Allocation of Testing Resources," "Flaws in Testing Systems," and "Impact on Public Health and Society." The first subtheme is about whether testing availability should be prioritized for hospitalized patients or health care workers, but the second subtheme highlights significant complementary issues with regards to flaws in the CDC's COVID-19 testing protocols, delays in fixing the tests, and the impact on the ability to detect and track the spread of the virus. The third theme expanded further into the social implications of the impact of inadequate testing resources, limited testing on the perception of the virus's severity, and the potential spread of the virus due to lack of testing and preventive measures.

Conversely, the expectation for the agreement scale item 5 was that complementarity would be low, yet ChatGPT found 87% complementarity. For instance, in the example mentioned above, the topic labeled by human-coders as "Knowledge about virus," the subthemes are considered divergent ("Importance of proper PPE and testing for healthcare workers") and too broad ("Concerns and challenges in healthcare settings and home care") when compared with the scope defined by "knowledge about the virus." The posts on these themes cover diverse topics such as the importance of proper personal protective equipment (PPE), concerns about testing and returning to work, the potential risks involved in home care, questions about Health Insurance Portability and Accountability Act (HIPAA) regulations, and the need for research on treatment options. The complementarity of themes only exists in a very broad sense and can be considered as "out of context."

Discussion

Principal Findings

Our study offers several significant insights into the use of ChatGPT for the augmentation of topic models. A key finding is the importance of considering different levels of abstraction in theme analysis. The division into themes and subthemes is crucial for uncovering specific nuances, addressing the risk of overgeneralization inherent in LLMs.

Furthermore, our exploration of subthemes reveals that LLMs, in general, can resolve ambiguity, aiding in the clear categorization of topics, even from a limited dataset. The effective handling of "low-coherence" topics such as "health disparities" and the complementary insights provided on subthemes of "Testing policies in different settings" demonstrate the LLM's proficiency in navigating and categorizing complex subject matter at the subtheme level.

In terms of overall reliability, our study estimates a 79.7% (247/310) agreement level, positioning it at the high end of substantial agreement (60%-80%) and the low end of almost perfect agreement (80%-100%) on the intercoder reliability benchmark scale. This suggests a robust level of agreement between human coders and the LLM, indicating a reliable consistency in the classification of topics.

However, the examination of alignment and convergence reveals a nuanced aspect of LLM performance. While LLMs exhibit high accuracy in identifying alignment and convergence for topics classified by human analysis as aligned, a notable challenge arises when classifying divergent subthemes. The LLM tends to classify divergent subthemes as convergent, particularly when one of the subthemes converges in similar ideas, leading to a potential misrepresentation of thematic divergence.

The evaluation of coherence, yields an unexpected result, highlighting the issue of "overfitting." Specifically, topics classified as coherent by the LLM contradict human coders' assessments of low coherence. This suggests a potential challenge where ChatGPT may force-fit solutions that match specific data points (posts) but are "too good to be true" from a pattern standpoint, lacking the broader pattern consistency expected in thematic coherence. ChatGPT may be construing the theme based on the wealth of data at its disposal.

The analysis of complementarity confirms that LLMs identify subthemes that provide additional insights to themes in human researchers' findings. LLMs can successfully identify niche topics, showcasing their potential to uncover unique thematic elements.

Our study emphasizes the critical importance of providing adequate contextual framing to ChatGPT-based classification. The challenge of lack of context becomes apparent, as LLMs may misinterpret or overlook certain topic nuances without external knowledge or the ability to track long-term context.

Limitations

The study is limited by (1) our focus on a single social media source and (2) the LLM used. First, we focus on data from a single nurse forum, but future inclusion of additional social media sites, including those used in other countries and by users who speak other languages, may enhance the results reported here. Furthermore, while we used the OpenAI chat completion API (GPT-3.5 and GPT-4) for thematic analysis due to its accessibility to the research team, other language models have since emerged. These newer models should be tested to determine if they perform better in different contexts. Furthermore, we kept the LLM prompts as simple as possible to demonstrate that even using a simple approach the generative AI could produce solid results. Further work can apply fine tuning to prompting and design approaches to enhance the analysis capabilities of LLMs, thematic such as retrieval-augmented generation (RAG). Finally, we focus on inductive thematic analysis and short form content data. We recognize that long-form text data may pose distinct challenges in applying LLMs.

Implications

For the LLM challenges found in this study, such as overgeneralization and overfitting, future study may apply different guardrails, such as implement algorithms that detect and mitigate biases during both training and generation phases. These guardrails monitor and filter the outputs of LLMs addressing different requirements such as hallucinations in LLM outputs [42].

Future research could investigate the potential of feeding raw transcripts into ChatGPT and incorporating AI-generated themes into triangulation discussions. By contributing to triangulation, this approach promises to unveil potential oversights, present alternative perspectives, and highlight inherent researchers' personal biases. By seamlessly incorporating AI into the discourse analysis process, researchers may uncover a richer understanding of the subject matter, fostering a more comprehensive and nuanced exploration of diverse perspectives. This integration not only enhances the depth of analysis but also provides a valuable tool for refining methodologies and mitigating potential biases, ultimately contributing to the advancement of research methodologies in the burgeoning field of AI-driven discourse analysis.

Conclusions

Overall, this study underscores the multifaceted nature of using ChatGPT for thematic analysis, acknowledging both its strengths and challenges. The insights gained contribute to a more nuanced understanding of the capabilities and limitations of LLMs in handling complex topical data in the healthcare field, offering valuable considerations for future research in the intersection of artificial intelligence and discourse analysis.

Conflicts of Interest

None declared.

Multimedia Appendix 1

Large language model's use for thematic analysis and classifying agreements. [DOCX File, 20 KB - ai v4i1e64447 app1.docx]

References

- Hussain MI, Figueiredo MC, Tran BD, et al. A scoping review of qualitative research in JAMIA: past contributions and opportunities for future work. J Am Med Inform Assoc 2021 Feb 15;28(2):402-413. [doi: <u>10.1093/jamia/ocaa179</u>] [Medline: <u>33225361</u>]
- Ranade-Kharkar P, Weir C, Norlin C, et al. Information needs of physicians, care coordinators, and families to support care coordination of children and youth with special health care needs (CYSHCN). J Am Med Inform Assoc 2017 Sep 1;24(5):933-941. [doi: 10.1093/jamia/ocx023] [Medline: 28371887]
- 3. Scharp D, Hobensack M, Davoudi A, Topaz M. Natural language processing applied to clinical documentation in post-acute care settings: a scoping review. J Am Med Dir Assoc 2024 Jan;25(1):69-83. [doi: 10.1016/j.jamda.2023.09.006] [Medline: 37838000]
- Jiang H, Castellanos A, Castillo A, Gomes PJ, Li J, VanderMeer D. Nurses' work concerns and disenchantment during the COVID-19 pandemic: machine learning analysis of web-based discussions. JMIR Nurs 2023 Feb 6;6:e40676. [doi: 10.2196/40676] [Medline: <u>36608261</u>]
- Hobensack M, Ojo M, Barrón Y, et al. Documentation of hospitalization risk factors in electronic health records (EHRs): a qualitative study with home healthcare clinicians. J Am Med Inform Assoc 2022 Apr 13;29(5):805-812. [doi: 10.1093/jamia/ocac023] [Medline: 35196369]
- 6. Zhao WX, Zhou K, Li J, et al. A survey of large language models. arXiv. Preprint posted online on Mar 31, 2023. [doi: 10.48550/arXiv.2303.18223]
- Älgå A, Eriksson O, Nordberg M. Analysis of scientific publications during the early phase of the COVID-19 pandemic: topic modeling study. J Med Internet Res 2020 Nov 10;22(11):e21559. [doi: <u>10.2196/21559</u>] [Medline: <u>33031049</u>]
- Chandrasekaran R, Mehta V, Valkunde T, Moustakas E. Topics, trends, and sentiments of tweets about the COVID-19 pandemic: temporal infoveillance study. J Med Internet Res 2020 Oct 23;22(10):e22624. [doi: <u>10.2196/22624</u>] [Medline: <u>33006937</u>]

- 9. Kavvadias S, Drosatos G, Kaldoudi E. Supporting topic modeling and trends analysis in biomedical literature. J Biomed Inform 2020 Oct;110:103574. [doi: 10.1016/j.jbi.2020.103574] [Medline: 32971274]
- Jelodar H, Wang Y, Yuan C, et al. Latent Dirichlet allocation (LDA) and topic modeling: models, applications, a survey. Multimed Tools Appl 2019 Jun;78(11):15169-15211. [doi: <u>10.1007/s11042-018-6894-4</u>]
- Buenano-Fernandez D, Gonzalez M, Gil D, Lujan-Mora S. Text mining of open-ended questions in self-assessment of university teachers: an LDA topic modeling approach. IEEE Access 2020;8:35318-35330. [doi: 10.1109/ACCESS.2020.2974983]
- 12. Ibrahim EI, Voyer A. The augmented qualitative researcher: using generative AI in qualitative text analysis. SocArXiv. Preprint posted online on Jan 22, 2024. [doi: 10.31235/osf.io/gkc8w]
- 13. Gilardi F, Alizadeh M, Kubli M. ChatGPT outperforms crowd workers for text-annotation tasks. Proc Natl Acad Sci U S A 2023 Jul 25;120(30):e2305016120. [doi: 10.1073/pnas.2305016120] [Medline: 37463210]
- 14. Morgan DL. Exploring the use of artificial intelligence for qualitative data analysis: the case of ChatGPT. Int J Qual Methods 2023 Oct;22:16094069231211248. [doi: 10.1177/16094069231211248]
- 15. Hamilton L, Elliott D, Quick A, Smith S, Choplin V. Exploring the use of AI in qualitative analysis: a comparative study of guaranteed income data. Int J Qual Methods 2023 Oct;22. [doi: 10.1177/16094069231201504]
- 16. Devlin J, Chang MW, Lee K, Toutanova K. BERT: pre-training of deep bidirectional transformers for language understanding. Preprint posted online on Oct 11, 2018. [doi: <u>10.48550/arXiv.1810.04805</u>]
- 17. Ji Z, Lee N, Frieske R, et al. Survey of hallucination in natural language generation. ACM Comput Surv 2023 Dec 31;55(12):1-38. [doi: 10.1145/3571730]
- 18. Srivastava N, Hinton G, Krizhevsky A, Sutskever I, Salakhutdinov R. Dropout: a simple way to prevent neural networks from overfitting. J Mach Learn Res 2014;15:1929-1958. [doi: <u>10.5555/2627435.2670313</u>]
- 19. Xue F, Fu Y, Zhou W, Zheng Z, You Y. To repeat or not to repeat: insights from scaling LLM under token-crisis. In: Oh A, Naumann T, Globerson A, Saenko K, Hardt M, Levine S, editors. NIPS '23: Proceedings of the 37th International Conference on Neural Information Processing System: Curran Associates; 2023, Vol. 36:59304-59322.
- 20. Peters M, Neumann M, Iyyer M, et al. Deep contextualized word representations. In: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers): Association for Computational Linguistics; 2018. [doi: 10.18653/v1/N18-1202]
- 21. Bolukbasi T, Chang KW, Zou JY, Saligrama V, Kalai AT. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. In: Lee D, Sugiyama M, Luxburg U, Guyon I, Garnett R, editors. NIPS'16: Proceedings of the 30th International Conference on Neural Information Processing Systems: Curran Associates; 2016:4356-4364. [doi: 10.5555/3157382.3157584]
- 22. Lin L, Wang L, Guo J, Wong KF. Investigating bias in LLM-based bias detection disparities between llms and human perception. arXiv. Preprint posted online on Mar 22, 2024. [doi: <u>10.48550/arXiv.2403.14896</u>]
- 23. Griffiths T, Jordan M, Tenenbaum J, Blei D. Hierarchical topic models and the nested chinese restaurant process. In: Thrun S, Saul L, Schölkopf B, editors. Advances in Neural Information Processing Systems 2003, Vol. 16. URL: <u>https://proceedings.neurips.cc/paper_files/paper/2003/file/7b41bfa5085806dfa24b8c9de0ce567f-Paper.pdf</u> [accessed 2025-03-27]
- 24. Hajikhani A, Cole C. A critical review of large language models: sensitivity, bias, and the path toward specialized AI. Quant Sci Stud 2024 Aug 1;5(3):736-756. [doi: 10.1162/qss_a_00310]
- 25. Brown T, Mann B, Ryder N, et al. Language models are few-shot learners. In: Larochelle H, Ranzato M, Hadsell R, Balcan MF, Lin H, editors. Advances in Neural Information Processing Systems: Curran Associates; 2020, Vol. 33.
- 26. Amershi S, Cakmak M, Knox WB, Kulesza T. Power to the people: the role of humans in interactive machine learning. AI Mag 2014 Dec;35(4):105-120. [doi: 10.1609/aimag.v35i4.2513]
- 27. Koch MA. Turning chaos into meaning: a ChatGPT-assisted exploration of COVID-19 narratives [Master's thesis]. : University of Twente; 2023 URL: <u>https://purl.utwente.nl/essays/96885</u> [accessed 2025-03-27]
- 28. Mesec B. The language model of artificial inteligence ChatGPT a tool of qualitative analysis of texts. Authorea. Preprint posted online on Apr 18, 2023. [doi: 10.22541/au.168182047.70243364/v1]
- 29. Braun V, Clarke V. Thematic analysis: a practical guide. In: Adv Neural Inf Process Syst: MIT Press; 2021.
- 30. Proudfoot K. Inductive/deductive hybrid thematic analysis in mixed methods research. J Mix Methods Res 2023 Jul;17(3):308-326. [doi: 10.1177/15586898221126816]
- 31. Rehurek R, Sojka P. Gensim–Python framework for vector space modelling. : NLP Centre, Faculty of Informatics, Masaryk University; 2011 URL: <u>https://www.fi.muni.cz/usr/sojka/posters/rehurek-sojka-scipy2011.pdf</u> [accessed 2025-03-27]
- 32. Srinivasa-Desikan B. Natural Language Processing and Computational Linguistics: A Practical Guide to Text Analysis with Python, Gensim, Spacy, and Keras: Packt Publishing Ltd; 2018:978.
- Stevens K, Kegelmeyer P, Andrzejewski D, Buttler D. Exploring topic coherence over many models and many topics. Presented at: 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning; Jul 12-14, 2012; Jeju Island, Korea.
- 34. Landis JR, Koch GG. The measurement of observer agreement for categorical data. Biometrics 1977 Mar;33(1):159-174. [doi: <u>10.2307/2529310</u>] [Medline: <u>843571</u>]

- 35. Xu FF, Alon U, Neubig G, Hellendoorn VJ. A systematic evaluation of large language models of code. Presented at: MAPS '22: 6th ACM SIGPLAN International Symposium on Machine Programming; Jun 13, 2022; San Diego, CA, United States. [doi: 10.1145/3520312.3534862]
- 36. Neuendorf KA. The Content Analysis Guidebook, 2nd edition: SAGE; 2017. [doi: 10.4135/9781071802878]
- 37. Saldaña J. The Coding Manual for Qualitative Researchers, 4th edition: SAGE; 2021.
- Pickering MJ, Garrod S. Toward a mechanistic psychology of dialogue. Behav Brain Sci 2004 Apr;27(2):169-190. [doi: 10.1017/s0140525x04000056] [Medline: 15595235]
- 39. van Dijk TA, Kintsch W. Strategies of Discourse Comprehension: Academic Press; 1983.
- 40. Clark HH, Brennan SE. Grounding in communication. In: Resnick LB, Levine JM, Teasley SD, editors. Perspectives on Socially Shared Cognition: American Psychological Association; 1991:127-149. [doi: <u>10.1037/10096-006</u>]
- 41. Gernsbacher MA, Givón T, editors. Coherence in Spontaneous Text Papers Presented at the Symposium on Coherence in Spontaneous Text: University of Oregon; 1992.
- 42. Dong Y, Mu R, Jin G, et al. Building guardrails for large language models. arXiv. Preprint posted online on Feb 2, 2024. [doi: <u>10.48550/arXiv.2402.01822</u>]

Abbreviations

AI: artificial intelligence
API: application programming interface
HIPAA: Health Insurance Portability and Accountability Act
LDA: latent Dirichlet allocation
LLM: large language model
NLP: natural language processing
NLTK: Natural Language Toolkit
PPE: personal protective equipment
RAG: retrieval-augmented generation

Edited by F Dankar; submitted 17.07.24; peer-reviewed by C Wang, E Ash, X Zhang, Y Feng; revised version received 02.12.24; accepted 27.02.25; published 04.04.25.

Please cite as:

Castellanos A, Jiang H, Gomes P, Vander Meer D, Castillo A Large Language Models for Thematic Summarization in Qualitative Health Care Research: Comparative Analysis of Model and Human Performance JMIR AI 2025;4:e64447 URL: <u>https://ai.jmir.org/2025/1/e64447</u> doi:<u>10.2196/64447</u>

© Arturo Castellanos, Haoqiang Jiang, Paulo Gomes, Debra Vander Meer, Alfred Castillo. Originally published in JMIR AI (https://ai.jmir.org), 4.4.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Evaluation of ChatGPT Performance on Emergency Medicine Board Examination Questions: Observational Study

Mila Pastrak^{1*}, BSc; Sten Kajitani^{1*}, BSc; Anthony James Goodings¹, DEC; Austin Drewek², MD; Andrew LaFree³, MD; Adrian Murphy^{1,4}, MB, BCh, BAO, PhD

¹School of Medicine, University College Cork, Cork, Ireland

²Department of Emergency Medicine, Johns Hopkins University, Baltimore, MD, United States

³Department of Emergency Medicine, University of California, San Diego, 200 W. Arbor Dr. #8676, San Diego, CA, United States

⁴Department of Emergency Medicine, Cork University Hospital, Cork, Ireland

^{*}these authors contributed equally

Corresponding Author:

Andrew LaFree, MD

Department of Emergency Medicine, University of California, San Diego, 200 W. Arbor Dr. #8676, San Diego, CA, United States

Abstract

Background: The ever-evolving field of medicine has highlighted the potential for ChatGPT as an assistive platform. However, its use in medical board examination preparation and completion remains unclear.

Objective: This study aimed to evaluate the performance of a custom-modified version of ChatGPT-4, tailored with emergency medicine board examination preparatory materials (Anki flashcard deck), compared to its default version and previous iteration (3.5). The goal was to assess the accuracy of ChatGPT-4 answering board-style questions and its suitability as a tool to aid students and trainees in standardized examination preparation.

Methods: A comparative analysis was conducted using a random selection of 598 questions from the Rosh In-Training Examination Question Bank. The subjects of the study included three versions of ChatGPT: the Default, a Custom, and ChatGPT-3.5. The accuracy, response length, medical discipline subgroups, and underlying causes of error were analyzed.

Results: The Custom version did not demonstrate a significant improvement in accuracy over the Default version (P=.61), although both significantly outperformed ChatGPT-3.5 (P<.001). The Default version produced significantly longer responses than the Custom version, with the mean (SD) values being 1371 (444) and 929 (408), respectively (P<.001). Subgroup analysis revealed no significant difference in the performance across different medical subdisciplines between the versions (P>.05 in all cases). Both the versions of ChatGPT-4 had similar underlying error types (P>.05 in all cases) and had a 99% predicted probability of passing while ChatGPT-3.5 had an 85% probability.

Conclusions: The findings suggest that while newer versions of ChatGPT exhibit improved performance in emergency medicine board examination preparation, specific enhancement with a comprehensive Anki flashcard deck on the topic does not significantly impact accuracy. The study highlights the potential of ChatGPT-4 as a tool for medical education, capable of providing accurate support across a wide range of topics in emergency medicine in its default form.

(JMIR AI 2025;4:e67696) doi:10.2196/67696

KEYWORDS

artificial intelligence; ChatGPT-4; medical education; emergency medicine; examination; examination preparation

Introduction

Background

The integration of artificial intelligence (AI) into medical education represents a frontier with the potential to significantly enhance learning outcomes and examination preparation strategies [1-5]. This advancement comes at a crucial time when the medical field faces the dual challenges of rapidly evolving knowledge bases and the increasing complexity of patient care.

https://ai.jmir.org/2025/1/e67696

generate human-like text based on a vast database of information has sparked interest in its application for medical board examination preparation. Previous studies have shown mixed results regarding the effectiveness of A Lin medical education, with certain limitations

effectiveness of AI in medical education, with certain limitations identified in AI's ability to replicate the depth of understanding needed to answer questions correctly in high-stakes examinations [7-12]. Building upon this background, our study

Among the AI tools making strides in educational contexts,

ChatGPT has emerged as a notable platform [6]. Its ability to

seeks to determine whether a targeted enhancement of ChatGPT-4 can increase the accuracy of the model in answering board examination questions, particularly for the American Board of Emergency Medicine (ABEM) Examinations.

ChatGPT provides relatively accurate responses to questions in examinations such as the USMLE (United States Medical Licensing Examination) [13,14] and the ABFM (American Board of Family Medicine) examination [5]. This may instill the confidence in takers of these examinations to use ChatGPT as an additional tool to aid in preparation. For instance, when reviewing a question set, the trainee may use ChatGPT to provide the rationale for a correct answer or help them understand the questions that they responded incorrectly to. This provides the potential to streamline the preparation process by reducing the need to consult textbooks or internet-based resources, as retaining interaction with multiple sources, such as a validated question bank, flashcards, and ChatGPT, is likely to bolster confidence in the overall educational outcome [15]. Additionally, the functionality of ChatGPT enables the user to ask follow-up questions or for further clarification if the initial response is insufficient.

In the pursuit of enhancing the capabilities of ChatGPT-4 for emergency medicine board examination preparation, a comprehensive Anki deck was utilized as a resource for custom modification [16,17]. The specific Anki deck chosen, "The Emergency Medicine Residents' Deck," also called "Rob's Emergency Medicine Deck" [18], is a collection of emergency medicine knowledge, aggregating content from various premade decks and covering a wide array of topics pertinent to the field.

The information within this deck is sourced from a variety of educational resources and study aids [18]. The deck's development and maintenance are overseen by medical professionals, with a commitment to regular updates and improvements based on the latest research, peer-reviewed consensus, and user feedback.

Rationale

Medical learners seem to generally have a positive view on generative AI [19-21]. Incorporating its potential with another popular and effective resource [22,23], Anki, could be useful to this population. The hypothesis driving this study posits that a ChatGPT-4 model, when enhanced with the comprehensive knowledge contained in this Anki deck, would outperform its standard counterpart in emergency medicine board examination preparation. This assumption is grounded in the belief that the deck's content could significantly bolster the AI's understanding and response accuracy to examination-relevant questions. Moreover, a positive outcome from this hypothesis could suggest that medical students who use this Anki deck for preparation could potentially be equipped with all the knowledge to excel in the board examination.

The Anki deck was chosen as it is designed to be a comprehensive resource. Additionally, Anki has become one of the most popular study methods among trainees and medical students. The approach of spaced repetition is particularly useful in helping people recall information. While an AI model would not engage in spaced repetition, the content of the decks can be

```
https://ai.jmir.org/2025/1/e67696
```

used to train the AI. By using this method, it can allow us to evaluate the performance of ChatGPT when provided with a widely used, evidence-based resource. Relative to other resources such as textbooks, an Anki user endeavors to recall every piece of information in the deck, while a textbook is generally not used in the same way.

Aims and Objectives

This study aimed to explore the efficacy of ChatGPT-4, specifically a custom-modified version tailored with specialized preparatory materials, in the context of emergency medicine board examination preparation. The objectives of this work were to: (1) evaluate the accuracy of ChatGPT-3.5 (released in 2022) in answering board examination style questions, (2) assess the baseline capabilities of the standard ChatGPT-4 model (released in 2023) in answering board examination questions accurately and consistently, and (3) evaluate whether a version custom-trained with a comprehensive flashcard resource exhibits superior performance. This comparison aimed to shed light on the potential of AI as a tool for medical education and identify pathways for its optimization in this domain.

Methods

Resources and Procedure

We used the Rosh In-Training Examination Question Bank, comprising 2000 questions, as the primary resource for questions. In order to customize ChatGPT-4 and transform it into a more specialized emergency medicine language model, "Rob's Emergency Medicine Deck," a comprehensive Anki deck for the ABEM Examinations, was converted to a TXT file and used to train the modified ChatGPT-4 model named "Emergency Medicine Residency Board Examination Expert."

Questions were selected from the question bank by selecting the "unused questions" option during the creation of individual practice examination question sets to ensure random selection and no overlapping questions.

Statistical Analysis

Sample Size

To examine if the sample size of 598 questions that were evaluated out of 2000 questions from the Rosh Review database is sufficient to make a conclusion about the performance of the two language models being equal, the following statistical assessment of the proportion of correct answers in each database was performed: the two-proportion z test was implemented to determine if there is a significant difference in error rates between the two language models; the alpha level of 0.05 was set to test the null hypothesis. The power was set at 0.80. The CIs for the difference between the two proportions were calculated; for the 5% significance level, a CI of 95% that included 0 would imply no significant difference between the error rates of the two language models.

The analysis showed that the two-proportion z score of approximately -0.073 corresponded to a P value of 0.942. Therefore, no statistically significant difference between the error rates indicates equal performance of the two language models. The z score close to 0 is also within the range of typical

sampling variation. In addition, the CIs for the proportions of correct answers using the Wilson Score Interval were approximately 77.3% to 83.6% for Custom ChatGPT-4 versus 77.1% to 83.5% for Default ChatGPT-4. The CI for the differences between the two proportions ranged between -4.7% and 4.3%. This narrow difference between the two proportions included 0, further showing no significant difference in the performance of the two language models.

Hence, a sample size of 598 questions that represent 29.9% of the Rosh Review database is sufficient to reliably assess the performance of the two language models.

Comparative Analysis

The performance of both the default and enhanced ChatGPT-4 models was compared based on the number of correct and incorrect answers. The incorrect responses were categorized according to the reason for error (logical error, informational error, or other), an approach used in previous studies [5,24], and analyzed for patterns.

A logical error is when the response successfully identified the relevant information but failed to effectively transform it into an answer. For example, the model identifies that a patient is struggling with the consistent use of topical acne medications due to a busy schedule and yet selects the answer that is a daily treatment over a less frequent regimen.

An informational error is when ChatGPT missed a crucial detail, either contained within the question or from external sources that should be part of its expected knowledge base. For example, a young woman is seeking birth control with a history of deep vein thrombosis, yet it recommends the oral contraceptive pill when deep vein thrombosis is a contraindication.

All remaining errors that are not related to the nonadequate connection to information, had insufficient consideration of all elements of the information, or had an arithmetic mistake were classified as "other". For example, the model identifies that a patient has cardiac failure yet inaccurately classifies the patient per the New York Heart Association Classification.

Incorrect Response Analysis and Question Type Assessment

For each incorrect response, the explanation provided by ChatGPT-4 was quantified (as response length in characters without spaces). Incorrect questions were classified by type (cardiac emergencies, neurological emergencies, respiratory emergencies, etc) to identify specific areas of weakness.

Statistical Analysis and Data Manipulation

A combination of statistical tests and data manipulation techniques were employed, facilitated by Python. The data were managed and manipulated using Pandas [25], a Python library offering data structures and tools designed for efficient data manipulation and analysis. Tasks such as filtering data,

computing descriptive statistics, and organizing data into contingency tables for further statistical testing were conducted.

For statistical analyses, several methods were employed to assess differences in performance between versions of ChatGPT. The McNemar test was carried out using the SciPy library [26] to compare paired nominal data across different subgroups. Additionally, for comparisons involving proportions, the proportions_*z* test function from the Statsmodels library [27], which provides comprehensive classes and functions for estimating different statistical models and performing statistical tests, was used.

Furthermore, the Wilcoxon signed-rank test, through the SciPy library, was applied for the analysis of paired proportions with nonparametric methods to assess the statistical significance of differences between the versions without assuming the normal distribution of the data. CIs for proportions were estimated using a normal approximation method, underlining the assumptions made regarding the distribution of the sample proportions.

Ethical Considerations

As an observational study involving an AI system, there were no human or animal subjects, thus minimizing ethical concerns. Ethical approval was not required for this study in accordance with the criteria of the Clinical Research Ethics Committee of the Cork Teaching Hospitals, University College Cork.

Results

Data Collection

All results were collected from February 24, 2024 to March 13, 2024. The default ChatGPT-4 model was tested by manually entering a randomized selection of 598 questions from the Rosh In-Training Examination Question Bank. The ChatGPT-3.5 model was tested using a randomized selection of 269 questions from the same set of questions presented to the default ChatGPT-4 model.

Comparison of Models

Percent of Questions Correct

Table 1 shows the performance of Custom ChatGPT-4, Default ChatGPT-4, and Default ChatGPT-3.5 on the randomized 598 question Rosh Review bank. Custom ChatGPT-4 and Default ChatGPT-4 answered 481 questions (80.4%, 95% CI 77.3% to 83.6%) and 480 questions (80.3%, 95% CI 77.1% to 83.5%) correct, respectively, with *P*=.61. These results indicate that the overall performance for correctly answering is similar between the two versions, with overlapping CIs, suggesting no significant difference in their ability. However, Custom ChatGPT-4 significantly outperformed ChatGPT-3.5 by 17.6% while Default ChatGPT-4 significantly outperformed Default ChatGPT-3.5 by 17.5% (*P*<.001 and *P*<.001, respectively).

Table. The performance of three language models on the American Board of Emergency Medicine examination using the Rosh Review question bank.

	Custom ChatGPT-4 (n=598)	Default ChatGPT-4 (n=598)	Default ChatGPT-3.5 (n=269)
Number of Correct Questions	481	480	169
Correct (%)	80.4	80.3	62.8

https://ai.jmir.org/2025/1/e67696

Length of Responses

The Custom ChatGPT-4 had significantly shorter response lengths, 929 (SD=408) characters without spaces versus 1371 (SD=444) characters without spaces for the Default ChatGPT-4 (P<.001). This suggests that Default ChatGPT-4 provided either more comprehensive or verbose responses.

Responses by Discipline

In Table 2, we conducted a subgroup analysis to explore the performance of the Custom ChatGPT-4 and Default ChatGPT-4 versions across 15 different disciplines within emergency medicine. There were no statistically significant differences in the number of correct questions per discipline between Custom ChatGPT-4 and Default ChatGPT-4 in the 15 groups: 12/15 of the subgroups had P=1.0, except ear, nose, and throat (P=.23); obstetrics and gynecology (P=.50); and other (P=.77).

 Table . Comparison of custom ChatGPT-4 and default ChatGPT-4 correct performance in Rosh Review subgroup analysis.

Subgroup	Custom ChatGPT-4 (n, %)	Default ChatGPT-4 (n, %)
Cardiology	81 (72.8)	81 (71.6)
Respirology	48 (70.8)	48 (73.5)
Neurology	33 (87.9)	33 (84.9)
Infectious Diseases	72 (84.7)	72 (83.1)
Gastrointestinal	51 (80.4)	51 (82.4)
Renal	15 (80.0)	15 (86.7)
Reproductive	9 (88.9)	9 (88.9)
Endocrine	23 (78.3)	23 (78.3)
Musculoskeletal	37 (73.0)	37 (73.0)
Ear, Nose, and Throat	26 (80.8)	26 (92.3)
Dermatology	16 (81.3)	16 (81.3)
Ophthalmology	20 (90.0)	20 (85.0)
Obstetrics and Gynecology	24 (87.5)	24 (79.2)
Oncology and Hematology	30 (86.2)	30 (90.0)
Other (Environmental)	113 (82.5)	113 (80.5)

Error Type Analysis

In Table 3, the type of error made by the Custom ChatGPT-4 and Default ChatGPT-4 was evaluated. There was no significant

difference between Custom ChatGPT-4 and Default ChatGPT-4 for logical error (75.2% vs 80.5%), informational error (12.0% vs 13.6%), or other (12.8% vs 5.9%), with P=.41, P=.87, and P=.11, respectively.

 $\ensuremath{\textbf{Table}}$. Assessment of the type of error conducted in two language models.

Error type	Custom ChatGPT-4 (n, %)	Default ChatGPT-4 (n, %)
Logical error	88 (75.2)	95 (80.5)
Informational error	14 (12.0)	16 (13.6)
Other	15 (12.8)	7 (5.9)
Total	117 (100)	118 (100)

Probability of Achieving a Passing Score

The passing probability of each ChatGPT model as predicted by the Rosh Review according to the individual ChatGPT performance was compared to the true performance of emergency medicine residents who wrote the ABEM in 2023. The newest ChatGPT models, ChatGPT-4 had a 99% chance of passing in both the Custom and Default versions. These were higher than the 85% probability of the default ChatGPT-3.5 version to pass and the 88% overall pass rate for the human counterparts. Notably, the human counterparts outperformed the ChatGPT-3.5 model.

```
https://ai.jmir.org/2025/1/e67696
```

RenderX

Discussion

Principal Findings

A prominent characteristic highlighted through the development of ChatGPT is its capacity to grasp the context and key details that are pertinent to the discussed subject. Our study demonstrates that this capability is also applicable within the medical field by evaluating three versions of ChatGPT with the same data set. We found that both the custom and default models are highly likely capable of passing the ABEM written examination. This is supported by the Rosh Review [28], which

had a predictive measure of passing the examination with the probability of passing at 98.8% accuracy; the Rosh Review found that both models had a 99% probability of passing. However, ChatGPT-3.5 had an 85% probability of passing. This prediction suggests that the enhancements made for the custom-modified version did not significantly improve accuracy over the default version of ChatGPT-4 and also shows that advancements made between ChatGPT versions have potential applications in the medical field. These findings imply that the core capabilities of ChatGPT-4 are already sufficiently advanced for tasks such as aiding in emergency medicine board examination preparation. Furthermore, the recorded national average pass rate for first-time test takers is 91%, with the 2023 pass rate being 88% [29], suggesting that ChatGPT-4 has an improved performance while ChatGPT-3.5 is less equipped compared to humans.

In addition, our results illustrate that both models had consistent performance across various medical disciplines and highlight the versatility of ChatGPT as an educational tool. This versatility is particularly relevant in the context of emergency medicine, where a broad spectrum of knowledge is required, and suggests that AI can offer comprehensive support across diverse subject areas. Additionally, the integration of an Anki deck into a ChatGPT-4 model could help identify the specific flashcards and topics that the learners should focus on, an area for future research.

Comparison of Error Types and Response Length

The custom and default models had a similar level of drawing incorrect conclusions and omitting important components of questions, both of which hint at areas for improvement in both models. The high percentage of logical errors, compared to the other two errors, indicates that language models may not be particularly well suited in deductive reasoning [30]. It may be possible to address this by careful prompt engineering [31], for instance, instructing the model to follow a hierarchy of information sources to deliver the most reliable answers consistently. This is an area that could be the subject of further research.

Additionally, the response length analysis revealed that longer responses do not necessarily correlate with increased accuracy. Prompt engineering could be used to enhance the ease of learning by outlining a preferred explanation format. This finding has practical implications for the design of AI-driven study tools, suggesting that brevity, combined with accuracy, could enhance the efficiency of study sessions and information retention for learners. In contrast, it could be argued that longer responses reflect more comprehensive explanations. Future studies and particularly a qualitative analysis could be done to interrogate these hypotheses.

Effect of Custom Training on Performance

The results underscore the rapid advancements in AI technology, particularly in natural language processing and knowledge retrieval, which have significant implications for medical education. The observed improvements from version 3.5 to the more recent iterations of ChatGPT reflect a trajectory in AI development that could increasingly support complex learning

XSI•FC

needs. This evolution underscores the potential of AI to become an increasingly valuable asset in educational settings [6,19], offering up-to-date knowledge and adaptive learning paths on balance with a general cautious optimism among medical professionals [32]. Despite the lack of observed benefit from custom modifications in this context, the findings highlight the critical role of up-to-date AI models in enhancing learning outcomes. Furthermore, the results illustrate that the untrained ChatGPT-4 has a higher likelihood of passing compared to human test takers, who extensively prepared for the board examinations, suggesting that, even without custom modifications, ChatGPT-4 has sufficient accuracy to serve as a customizable tutor.

Overall, while the investigation revealed no significant difference in performance accuracy between the custom-modified and default versions of ChatGPT-4, both showed considerable improvement over the older 3.5 version. These findings prompt a re-evaluation of the presumed advantage of tailoring AI through specific educational content, suggesting that the core capabilities of advanced AI models might already be sufficiently robust for some less highly subspecialized educational applications. Additionally, these findings promote investigation into future upcoming ChatGPT models to evaluate if their advancements have accelerated benefit in the medical field.

When evaluating the reason for the Custom model not being significantly better than the Default model, we must consider that the Default version has already been trained on sufficiently similar data that the information provided did not contribute anything new to the knowledge base. The need for AI to be trained on up-to-date data is well established [33]. A previous study has hypothesized that training the model on static knowledge could potentially be a limiting factor [5], the reason for this being that online resources can be constantly updated with the latest guidelines and treatments. Basing training on a well-maintained dynamic knowledge source such as UpToDate® (Wolters Kluwer) could potentially provide more useful outcomes. It seems that general medicine knowledge has been well incorporated into the training material for the ChatGPT-4 model, and this can explain the similar performance between the two versions of ChatGPT-4 we tested. However, for more niche and subspecialized fields, there may exist a more pronounced benefit, and this is something future works could explore.

This study reaffirmed the potential of AI, particularly ChatGPT-4, as a powerful tool in medical education [6,34], capable of supporting learners in high-stakes examination preparation without the need for specialized enhancements. It highlighted the importance of leveraging the inherent capabilities of advanced AI models and provided a foundation for further research into effective integration strategies in educational settings. As AI continues to evolve, its role in education is likely to expand, offering opportunities to enrich learning experiences and access to knowledge.

Limitations

While this study provides valuable insights, it is not without limitations. The scope was restricted to emergency medicine,

limiting the generalizability of the findings to other fields of medicine or education. Future research could explore the application of AI in different specialties to assess its versatility and effectiveness further.

Additionally, the study's design focused on the efficacy of AI in answering board examination questions, which may not fully capture the nuances of applying that knowledge in clinical practice [35]. Further studies could investigate the impact of AI-assisted learning on clinical skills and decision-making processes [36,37]. The results of this study are not generalizable to the use of AI in contexts of medical education beyond the use case described for examination preparation.

The study's limitations suggest caution in generalizing the findings to other disciplines or educational objectives. Future research could broaden the scope to include diverse medical specialties and different types of educational content to verify the applicability of these results more widely.

Conclusion

This study reaffirmed the potential of AI, particularly ChatGPT-4, as a powerful tool in medical education, capable of supporting learners in high-stakes examination preparation without the need for specialized enhancements. It highlighted the importance of leveraging the inherent capabilities of advanced AI models and provided a foundation for further research into effective integration strategies in educational settings. This could be accomplished by determining if linking ChatGPT to a dynamic and reliable data source provides benefits, focusing in on highly subspecialized fields with static information sources, and ultimately comparing evaluation and management plans generated by AI to physician counterparts. As AI continues to evolve, its role in education and potentially clinical practice is likely to expand, offering opportunities to enrich learning experiences and access to knowledge.

Authors' Contributions

MP and SHK contributed to the conceptual design, data collection, data analysis, and drafting of the manuscript. MP and SHK are equal contributors. AJG contributed to the conceptual design, data analysis, and drafting of the manuscript. AD and AL provided critical feedback conceptual design and contributed to editing and revision of the manuscript. AM provided critical feedback conceptual design, contributed to editing and revision of the manuscript, and supervised the project.

Conflicts of Interest

None declared.

References

- 1. Krive J, Isola M, Chang L, Patel T, Anderson M, Sreedhar R. Grounded in reality: artificial intelligence in medical education. JAMIA Open 2023 Jul;6(2):00ad037. [doi: 10.1093/jamiaopen/00ad037] [Medline: 37273962]
- 2. Civaner MM, Uncu Y, Bulut F, Chalil EG, Tatli A. Artificial intelligence in medical education: a cross-sectional needs assessment. BMC Med Educ 2022 Nov 9;22(1):772. [doi: 10.1186/s12909-022-03852-3] [Medline: 36352431]
- 3. Nagi F, Salih R, Alzubaidi M, et al. Applications of artificial intelligence (AI) in medical education: a scoping review. Stud Health Technol Inform 2023 Jun 29;305:648-651. [doi: 10.3233/SHTI230581] [Medline: 37387115]
- 4. Eysenbach G. The role of ChatGPT, generative language models, and artificial intelligence in medical education: a conversation with ChatGPT and a call for papers. JMIR Med Educ 2023 Mar 6;9:e46885. [doi: 10.2196/46885] [Medline: 36863937]
- Goodings AJ, Kajitani S, Chhor A, et al. Assessment of ChatGPT-4 in family medicine board examinations using advanced ai learning and analytical methods: observational study. JMIR Med Educ 2024 Oct 8;10:e56128. [doi: <u>10.2196/56128</u>] [Medline: <u>39378442</u>]
- 6. Abd-Alrazaq A, AlSaad R, Alhuwail D, et al. Large language models in medical education: opportunities, challenges, and future directions. JMIR Med Educ 2023 Jun 1;9:e48291. [doi: <u>10.2196/48291</u>] [Medline: <u>37261894</u>]
- 7. Joly-Chevrier M, Nguyen AXL, Lesko-Krleza M, Lefrançois P. Performance of ChatGPT on a practice dermatology board certification examination. J Cutan Med Surg 2023;27(4):407-409. [doi: <u>10.1177/12034754231188437</u>] [Medline: <u>37489920</u>]
- 8. Lai UH, Wu KS, Hsu TY, Kan JKC. Evaluating the performance of ChatGPT-4 on the United Kingdom medical licensing assessment. Front Med (Lausanne) 2023;10:1240915. [doi: 10.3389/fmed.2023.1240915] [Medline: 37795422]
- 9. Mihalache A, Popovic MM, Muni RH. Performance of an artificial intelligence chatbot in ophthalmic knowledge assessment. JAMA Ophthalmol 2023 Jun 1;141(6):589-597. [doi: <u>10.1001/jamaophthalmol.2023.1144</u>] [Medline: <u>37103928</u>]
- Mihalache A, Huang RS, Popovic MM, Muni RH. ChatGPT-4: An assessment of an upgraded artificial intelligence chatbot in the United States medical licensing examination. Med Teach 2024 Mar;46(3):366-372. [doi: <u>10.1080/0142159X.2023.2249588</u>] [Medline: <u>37839017</u>]
- 11. Ali R, Tang OY, Connolly ID, et al. Performance of ChatGPT and GPT-4 on neurosurgery written board examinations. Neurosurgery 2023 Dec 1;93(6):1353-1365. [doi: 10.1227/neu.00000000002632] [Medline: 37581444]
- 12. Barbour AB, Barbour TA. A radiation oncology board exam of ChatGPT. Cureus 2023 Sep;15(9):e44541. [doi: 10.7759/cureus.44541] [Medline: 37790062]

- Gilson A, Safranek CW, Huang T, et al. How Does ChatGPT perform on the United States medical licensing examination (USMLE)? The implications of large language models for medical education and knowledge assessment. JMIR Med Educ 2023 Feb 8;9:e45312. [doi: 10.2196/45312] [Medline: 36753318]
- 14. Brin D, Sorin V, Vaid A, et al. Comparing ChatGPT and GPT-4 performance in USMLE soft skill assessments. Sci Rep 2023 Oct 1;13(1):16492. [doi: 10.1038/s41598-023-43436-9] [Medline: 37779171]
- Hu JM, Liu FC, Chu CM, Chang YT. Health care trainees' and professionals' perceptions of ChatGPT in improving medical knowledge training: rapid survey study. J Med Internet Res 2023 Oct 18;25:e49385. [doi: <u>10.2196/49385</u>] [Medline: <u>37851495</u>]
- 16. Anki powerful, intelligent flashcards [Internet]. 2025 Jan 25. URL: https://apps.ankiweb.net/ [accessed 2025-03-04]
- 17. What is anki? [internet]. Am Med Assoc. 2023 Jan 25. URL: <u>https://www.ama-assn.org/medical-students/usmle-step-1-2/</u> what-anki [accessed 2025-03-04]
- 18. Rob's emergency medicine deck ankiweb [internet]. 2025 Jan 25. URL: <u>https://ankiweb.net/shared/info/790760070</u> [accessed 2025-03-04]
- 19. Chan CKY, Hu W. Students' voices on generative AI: perceptions, benefits, and challenges in higher education. Int J Educ Technol High Educ 2023;20(1):43. [doi: 10.1186/s41239-023-00411-8]
- 20. Bisdas S, Topriceanu CC, Zakrzewska Z, et al. Artificial intelligence in medicine: a multinational multi-center survey on the medical and dental students' perception. Front Public Health 2021;9:795284. [doi: 10.3389/fpubh.2021.795284] [Medline: 35004598]
- Ooi SKG, Makmur A, Soon AYQ, et al. Attitudes toward artificial intelligence in radiology with learner needs assessment within radiology residency programmes: a national multi-programme survey. Singapore Med J 2021 Mar;62(3):126-134. [doi: 10.11622/smedj.2019141] [Medline: 31680181]
- Gilbert MM, Frommeyer TC, Brittain GV, et al. A cohort study assessing the impact of Anki as a spaced repetition tool on academic performance in medical school. Med Sci Educ 2023 Aug;33(4):955-962. [doi: <u>10.1007/s40670-023-01826-8</u>] [Medline: <u>37546209</u>]
- 23. French BN, Marxen TO, Akhnoukh S, et al. A call for spaced repetition in medical education. Clin Teach 2024 Feb;21(1):e13669. [doi: 10.1111/tct.13669] [Medline: <u>37787460</u>]
- 24. Gupta R, Herzog I, Park JB, et al. Performance of ChatGPT on the plastic surgery inservice training examination. Aesthet Surg J 2023 Nov 16;43(12):NP1078-NP1082. [doi: <u>10.1093/asj/sjad128</u>] [Medline: <u>37128784</u>]
- 25. pandas Python Data Analysis Library [Internet]. 2025 Jan 25. URL: https://pandas.pydata.org/ [accessed 2025-03-04]
- 26. SciPy [Internet]. 2025 Jan 25. URL: <u>https://scipy.org/</u> [accessed 2025-03-04]
- 27. Perktold J, Seabold S, Sheppard K, et al. Statsmodels/statsmodels: release 0.14.2 [internet]. Zenodo. 2025 Jan 25. URL: https://zenodo.org/doi/10.5281/zenodo.593847 [accessed 2025-03-04]
- 28. Michael SS. Rosh review as a predictive instrument for ABEM concerttm exam performance. West J Emerg Med Integrating Emerg Care Popul Health [Internet]. 2014 Jan 25. URL: <u>https://escholarship.org/uc/item/1kh68596</u> [accessed 2025-03-04]
- 29. ABEM | exam & certification statistics [internet]. ABEM. 2025 Jan 25. URL: <u>https://www.abem.org/resources/</u> <u>exam-and-certification-statistics/</u> [accessed 2025-03-04]
- 30. Mondorf P, Plank B. Comparing inferential strategies of humans and large language models in deductive reasoning [internet]. arXiv. Preprint posted online on Jan 25, 2025 URL: <u>http://arxiv.org/abs/2402.14856</u> [accessed 2025-03-04]
- 31. Meskó B. Prompt engineering as an important emerging skill for medical professionals: tutorial. J Med Internet Res 2023 Oct 4;25:e50638. [doi: 10.2196/50638] [Medline: 37792434]
- 32. Sallam M. ChatGPT utility in healthcare education, research, and practice: systematic review on the promising perspectives and valid concerns. Healthcare (Basel) 2023 Mar 19;11(6):887. [doi: <u>10.3390/healthcare11060887</u>] [Medline: <u>36981544</u>]
- 33. Fatima A, Shafique MA, Alam K, Fadlalla Ahmed TK, Mustafa MS. ChatGPT in medicine: A cross-disciplinary systematic review of ChatGPT's (artificial intelligence) role in research, clinical practice, education, and patient interaction. Medicine (Baltimore) 2024 Aug 9;103(32):e39250. [doi: 10.1097/MD.00000000039250] [Medline: 39121303]
- 34. Lee H. The rise of ChatGPT: Exploring its potential in medical education. Anat Sci Educ 2024;17(5):926-931. [doi: 10.1002/ase.2270] [Medline: 36916887]
- 35. Hiller K, Franzen D, Heitz C, Emery M, Poznanski S. Correlation of the national board of medical examiners emergency medicine advanced clinical examination given in July to intern American board of emergency medicine in-training examination scores: a predictor of performance? West J Emerg Med 2015 Nov;16(6):957-960. [doi: 10.5811/westjem.2015.9.27303] [Medline: 26594299]
- 36. Joo H, Mathis MR, Tam M, et al. Applying AI and guidelines to assist medical students in recognizing patients with heart failure: protocol for a randomized trial. JMIR Res Protoc 2023 Oct 24;12:e49842. [doi: 10.2196/49842] [Medline: 37874618]
- Scaioli G, Lo Moro G, Conrado F, Rosset L, Bert F, Siliquini R. Exploring the potential of ChatGPT for clinical reasoning and decision-making: a cross-sectional study on the Italian Medical Residency Exam. Ann Ist Super Sanita 2023;59(4):267-270. [doi: <u>10.4415/ANN_23_04_05</u>] [Medline: <u>38088393</u>]



Abbreviations

ABEM: American Board of Emergency MedicineABFM: American Board of Family MedicineAI: artificial intelligenceUSMLE: United States Medical Licensing Examination

Edited by Z Yin; submitted 18.10.24; peer-reviewed by D Li, E Bai, J Krive; revised version received 12.02.25; accepted 12.02.25; published 12.03.25.

<u>Please cite as:</u> Pastrak M, Kajitani S, Goodings AJ, Drewek A, LaFree A, Murphy A Evaluation of ChatGPT Performance on Emergency Medicine Board Examination Questions: Observational Study JMIR AI 2025;4:e67696 URL: <u>https://ai.jmir.org/2025/1/e67696</u> doi:<u>10.2196/67696</u>

© Mila Pastrak, Sten Kajitani, Anthony James Goodings, Austin Drewek, Andrew LaFree, Adrian Murphy. Originally published in JMIR AI (https://ai.jmir.org), 12.3.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Original Paper

Enhancing Interpretable, Transparent, and Unobtrusive Detection of Acute Marijuana Intoxication in Natural Environments: Harnessing Smart Devices and Explainable AI to Empower Just-In-Time Adaptive Interventions: Longitudinal Observational Study

Sang Won Bae¹, PhD; Tammy Chung², PhD; Tongze Zhang¹, MSc; Anind K Dey³, PhD; Rahul Islam¹, BSc

¹Human-Computer Interaction and Human-Centered AI Systems Lab, AI for Healthcare Lab, Charles V. Schaefer, Jr. School of Engineering and Science, Stevens Institute of Technology, Hoboken, NJ, United States

²Institute for Health, Healthcare Policy and Aging Research, Rutgers University, Newark, NJ, United States

³Information School, University of Washington, Seattle, WA, United States

Corresponding Author:

Sang Won Bae, PhD Human-Computer Interaction and Human-Centered AI Systems Lab AI for Healthcare Lab, Charles V. Schaefer, Jr. School of Engineering and Science Stevens Institute of Technology 1 Castle Point Terrace Hoboken, NJ, 07030-5906 United States Phone: 1 4122658616 Email: sbae4@stevens.edu

Abstract

Background: Acute marijuana intoxication can impair motor skills and cognitive functions such as attention and information processing. However, traditional tests, like blood, urine, and saliva, fail to accurately detect acute marijuana intoxication in real time.

Objective: This study aims to explore whether integrating smartphone-based sensors with readily accessible wearable activity trackers, like Fitbit, can enhance the detection of acute marijuana intoxication in naturalistic settings. No previous research has investigated the effectiveness of passive sensing technologies for enhancing algorithm accuracy or enhancing the interpretability of digital phenotyping through explainable artificial intelligence in real-life scenarios. This approach aims to provide insights into how individuals interact with digital devices during algorithmic decision-making, particularly for detecting moderate to intensive marijuana intoxication in real-world contexts.

Methods: Sensor data from smartphones and Fitbits, along with self-reported marijuana use, were collected from 33 young adults over a 30-day period using the experience sampling method. Participants rated their level of intoxication on a scale from 1 to 10 within 15 minutes of consuming marijuana and during 3 daily semirandom prompts. The ratings were categorized as not intoxicated (0), low (1-3), and moderate to intense intoxication (4-10). The study analyzed the performance of models using mobile phone data only, Fitbit data only, and a combination of both (MobiFit) in detecting acute marijuana intoxication.

Results: The eXtreme Gradient Boosting Machine classifier showed that the MobiFit model, which combines mobile phone and wearable device data, achieved 99% accuracy (area under the curve=0.99; F_1 -score=0.85) in detecting acute marijuana intoxication in natural environments. The F_1 -score indicated significant improvements in sensitivity and specificity for the combined MobiFit model compared to using mobile or Fitbit data alone. Explainable artificial intelligence revealed that moderate to intense self-reported marijuana intoxication was associated with specific smartphone and Fitbit metrics, including elevated minimum heart rate, reduced macromovement, and increased noise energy around participants.

Conclusions: This study demonstrates the potential of using smartphone sensors and wearable devices for interpretable, transparent, and unobtrusive monitoring of acute marijuana intoxication in daily life. Advanced algorithmic decision-making provides valuable insight into behavioral, physiological, and environmental factors that could support timely interventions to

reduce marijuana-related harm. Future real-world applications of these algorithms should be evaluated in collaboration with clinical experts to enhance their practicality and effectiveness.

(JMIR AI 2025;4:e52270) doi:10.2196/52270

KEYWORDS

digital phenotyping; smart devices; intoxication; smartphone-based sensors; wearables; mHealth; marijuana; cannabis; data collection; passive sensing; Fitbit; machine learning; eXtreme Gradient Boosting Machine classifier; XGBoost; algorithmic decision-making process; explainable artificial intelligence; XAI; artificial intelligence; JITAI; decision support; just-in-time adaptive interventions; experience sampling

Introduction

Background

Acute effects of marijuana use impair motor skills and cognitive functions, such as attention and information processing [1-3], leading to adverse outcomes like poor academic and work performance, as well as an increased risk of motor vehicle crashes and fatal collisions [2,4]. Delta-9 tetrahydrocannabinol (THC), the principal psychoactive constituent of marijuana, binds to brain receptors, inducing a feeling of "euphoria" or being "high" [5]. Given the risks associated with THC-induced impairment, there is a critical need to detect episodes of marijuana intoxication in real time in the natural environment.

Several studies have explored the use of phone sensors or wearable devices to detect acute marijuana consumption. For example, a laboratory study with 10 participants used smartphone sensors (accelerometer, gyroscope) to detect acute marijuana use (3% or 7% THC vs placebo) and found that gait analysis with a support vector machine model achieved 92% accuracy (F_1 -score=0.93) [6]. Another study (n=1) developed an electrochemical biosensor ring that detected salivary THC (minimum of 0.5 µM) and blood alcohol levels (minimum of 0.2 mM) within three minutes [7]. However, these studies were conducted in controlled environments, highlighting the need for research on using smartphone and wearable sensors to detect acute marijuana use in nonlaboratory, natural settings.

Detecting marijuana use in daily life could enable Just-In-Time interventions to reduce harm, such as avoiding driving while intoxicated [8]. However, challenges exist in detecting acute marijuana-related intoxication [9]. THC could be detected in an individual's blood or urine for several days after consumption depending on factors such as recency, frequency, and chronicity of use [10]. Thus, a person who tests positive for THC might not be intoxicated or impaired at the time of testing [10]. Existing testing methods (eg, blood, urine, saliva, and breath) are not suitable for real-time detection, as THC can remain detectable in the body for days after consumption, which does not necessarily indicate current impairment [10].

To address these limitations, our recent study [11] used passive sensing via smartphones, coupled with self-reported intoxication, to detect marijuana use with 90% accuracy, using sensor-derived data from mobile phones alongside temporal variables, including time of day and day of week. Building on these findings [11], this study explores the use of wearable devices (eg, Fitbit) to enhance detection capabilities by incorporating physiological indicators, thereby improving the accuracy and immediacy of identifying marijuana effects in natural environments.

Wearable device-reported heart rate (HR) was examined as a potential physiological indicator of acute marijuana intoxication, based on laboratory studies, showing a dose-dependent increase in resting HR shortly after smoking or vaping marijuana [12-14]. Specifically, laboratory research reports that within 2-3 minutes of smoking marijuana, there is an acute increase (20%-60%) dose-dependent) in resting HR [13], which might represent a "physiological signal" of the onset of a marijuana smoking episode. HR peaks 10-15 minutes after reaching maximum THC levels, followed by a rapid decline [12-14]. While tolerance to this effect may develop (eg, from a mean increase of 44.6 to 6.6 beats per minute (bpm) after 18-20 days of use) with chronic use, [12-14]. The acute HR increases have been validated in laboratory settings but have remained unexplored in real-world contexts. This study examines using off-the-shelf wearable devices, such as Fitbit, to detect acute HR increases as a physiological signal potentially correlated with self-reported marijuana intoxication.

Research Objectives and Contributions

While laboratory studies have established the link between HR changes and marijuana intoxication [12-14], its applicability in real-world scenarios is unexplored. To address this gap, we propose that combining wearable device data with smartphone sensors could improve algorithms for detecting marijuana intoxication in real-life settings. To enhance the interpretability of our algorithms and provide insights for just-in-time adaptive interventions, we incorporated explainable artificial intelligence (XAI) into our machine-learning pipeline. XAI helps clarify the role of digital biomarkers associated with self-reported marijuana intoxication in natural environments.

This study aims to determine whether data from smartphones (eg, accelerometer and GPS) and wearable devices (eg, Fitbit) can detect self-reported marijuana intoxication ("feeling high") in the natural environment, a topic not previously investigated. Two hypotheses drive this research: (1) the novel MobiFit model, which combines smartphones and Fitbit data will outperform models that use only one data source in detecting self-reported intoxication; (2) HR and daily behavioral data (eg, step count) from Fitbit are important features for detecting self-reported marijuana intoxication. If either hypothesis is validated, it indicates the value of integrating wearable device data into daily life monitoring.

This study evaluates the performance of sensor-based models using (1) only smartphone sensors, (2) only Fitbit data, and (3)

XSL•FO RenderX

the combined MobiFit model. We also used XAI to enhance understanding of key digital features from both smartphone sensors and Fitbit data associated with self-reported marijuana intoxication. Identifying smartphone-based sensors and Fitbit features that accurately detect self-reported marijuana intoxication in natural environments could ultimately trigger just-in-time interventions.

This study presents a comprehensive approach toward using mobile and wearable technology for detecting self-reported acute marijuana intoxication in real-life settings, emphasizing interpretability and transparency through XAI. This study demonstrates the potential of integrating smart devices with advanced analytical techniques to improve detection accuracy and support timely interventions based on detected intoxication levels.

Methods

Recruitment and Participants

A total of 57 participants aged 18-24 years were recruited through flyers, advertisements, and local communities. Eligibility criteria were (1) using marijuana at least twice a week, (2) owning a personal mobile phone, (3) not currently seeking treatment for substance abuse, (4) no self-reported history of psychosis, and (5) not taking any medication or using any medical device (eg, pacemaker) that could affect HR. Of the 57 participants, 24 participants were excluded from the analysis due to missing data (eg, no HR data and no mobile sensor data).

The final analysis focused on 33 participants aged 18-24 years, with an average age of 19.64 (SD 1.77) years. Among these, 23 participants identified as White, 4 participants as Black, and 6 participants as other race or ethnicity. The average age of first marijuana use was 16.48 (SD 1.84, range 13-22) years, and the average age of regular marijuana use was 17.03 (SD 1.72) years. In this subset, 24% (n=8) reported daily marijuana use, 9% (n=3) reported using it 5-6 times per week, and 67% (n=22) reported using it 2-4 times per week. Notably, 97% (n=32) of participants primarily used iOS smartphones, with only 3% (n=1) using Android devices.

Ethical Considerations

This naturalistic, observational follow-along study was approved by the university's institutional review board (Stevens 2020-008 [23-COAS3], Rutgers Pro2019002365). In line with similar Institutional Review Board–approved observational studies [15], all participants were informed about local medical and mental health resources. The study obtained a National Institutes of Health Certificate of Confidentiality. Written consent was obtained from participants, who were informed about privacy protections and the voluntary nature of their participation [16]. The research staff explained the types of data to be collected, the duration of data collection, and the purpose of the study.

Study Design

Participants completed a baseline laboratory assessment including interviews, questionnaires, and cognitive testing. They downloaded study apps from the App Store or Google Play

```
https://ai.jmir.org/2025/1/e52270
```

Store to their smartphones. Research staff trained participants on how to use the apps and the study provided Fitbit Charge 2 for data collection. The AWARE mobile app [17] delivered experience sampling method (ESM) questions on marijuana use. Participants wore the Fitbit Charge 2 wristband to collect data on HR, physical activity (eg, step count), and sleep (eg, time, duration, and quality; see Table S2 in Multimedia Appendix 1 for Fitbit variables). The study collected continuous sensor data from smartphones and Fitbit devices, along with self-reported data on marijuana intoxication, for up to 30 days. A 30-day period was chosen to ensure sufficient data, given the study's inclusion criteria of frequent marijuana use. At the end of the study, participants completed a debriefing interview about their experience.

Participants were compensated for their time and effort, receiving US \$75 for completing the baseline assessment, and US \$25 for the debriefing interview. They earned US \$10 for each day on which they completed more than 75% of data collection (eg, Fitbit and ESM).

Mobile Sensing Framework and Applications for Data Collection

AWARE App

AWARE is a mobile sensing framework [17] that passively and continuously collects data from smartphone sensors. This data can be used to infer human behavior patterns using various sensors: location (eg, distance traveled and circadian rhythm), physical movements (eg, acceleration and activity), device usage (eg, unlock, charge, keypress, and app usage), social patterns (eg, communication and conversations), and environmental context (eg, Wi-Fi, Bluetooth, sound or ambient noise, and light). The app, developed to track participants' natural behaviors in real-life settings, runs in the background 24/7 and collects sensor data with associated metadata, such as time stamps and communication logs. The data is transferred to a secure MySQL database owned and operated by the research team.

ESM

The mobile app also captured self-reports of marijuana use by participants. Two types of surveys were used [18]. Participants manually reported marijuana use within 15 minutes of consumption, detailing the amount used, mode of consumption, and the people whom the participant consumed marijuana with. They also rated their subjective intoxication on a scale from 0 (none) to 10 (a lot) [19]. Two hours later, the app prompted participants to complete an end-session survey indicating when intoxication symptoms subsided. In addition, fixed-time surveys were delivered daily at 10 AM, 3 PM, and 8 PM to collect information on the participants' daily lives, including time since last marijuana use, cravings, mood, and feelings (eg, relaxed, anxious, and sad), and other substance use (eg, alcohol and tobacco). Survey response windows were open for 5 hours to accommodate participants' schedules.

Fitbit Charge 2

Participants were provided with Fitbit Charge 2 devices and asked to wear them as much as possible. Fitbit collected

physiological data (eg, HR), activity data (eg, step count), and sleep. The study hypothesized that HR and behavioral data could signal episodes of acute marijuana intoxication. Fitbit data were retrieved from the Fitbit server at the end of the study using the Fitbit application programming interface.

Preparing Self-Report and Fitbit Data for Analysis

An episode of self-reported subjective marijuana intoxication was defined based on the ESM item: "How high are you feeling right now?" rated from 0 to 10 (0=not high to 10=a lot) [18,19]. To include episodes in the analysis, both start and end times had to be reported to calculate duration and label the sensor data. To capture behaviors without marijuana use, 1556 reports where participants answered "no" to the question "Did you smoke marijuana since the last report?" during afternoon

Figure 1. Flowchart of participants and the data included in the analyses.

(n=1151) and evening (n=950) surveys were labeled as "0" for the subjective rating of marijuana intoxication.

From all participants, we received 641 self-reports (mean 9.86, SD 8.49; median 7, IQR 4-13) and 1556 with no marijuana use reports (Figure 1). Out of 641 reports, 168 reports had a subjective intoxication rating of 0 and 10, and 6 reports had no rating. After excluding 6 reports without ratings and 108 duplicate reports, 527 samples remained. Reports with missing start and end times, or implausible episode durations (eg, longer than 3 hours) were excluded based on laboratory research indicating that smoked or vaped marijuana effects last less than 3 hours [20]. A total of 136 self-reports were excluded for exceeding this duration, leaving 1556 reports where no marijuana use was recorded [20].



For model building, episodes without mobile sensor data (n=72) were excluded, leaving 221 marijuana self-reports. Furthermore,

episodes without Fitbit sensor data (n=17) were excluded, leaving 50 participants. These participants provided 132

https://ai.jmir.org/2025/1/e52270

marijuana use self-reports and 909 "no marijuana use" reports. We analyzed reports from each participant, excluding those who only reported not using marijuana or had a rating of 0 for subjective intoxication, leaving a total of 642 with no marijuana use report or who reported 0 subjective intoxications when using marijuana and 34 people. Finally, to prevent participants from using Fitbit incorrectly, we excluded users without HR data, leaving a total of 33 people, who provided a total of 769 events: 640 "no marijuana use" reports and 129 marijuana use self-reports.

Extracting Smartphone and Fitbit Sensor Features

Following previous studies, we extracted audio features to detect social interactions [21,22] potentially associated with marijuana use. Audio features were extracted using the conversation plug-in, which detects whether a person was engaged in a conversation. Raw audio signals are converted to amplitude using the Euclidean norm [23], which categorizes ambient levels into silence, noise, voice, and unknown [24]. We also computed device use features, such as smartphone unlock minutes and the duration of device interaction sessions. In addition to audio features, we extracted GPS features to examine movement patterns related to marijuana use [25-28]. These included the radius of gyration, time at a location cluster, total distance traveled, number of clusters within a 5-minute window, acceleration, and phone angles. Environmental features, such as the number of Bluetooth devices detected, the most frequently contacted Wi-Fi access point, and light features (eg, average [avg], and maximum [max] lux) were also extracted. For most features, we calculated the minimum (min), max, avg, median (med), and SD. Further details on smartphone features can be found in Multimedia Appendix 1.

We used a 5-minute time window for extracting sensor feature statistics, as laboratory studies show a dose-dependent acute in resting HR within 2-3 minutes of marijuana use. Using larger time intervals could include data not related to marijuana use, given the average reported marijuana session duration is 75 (SD 46.2) minutes.

Raw data for HR, sleep, and steps were extracted from Fitbit. We first obtained per-minute HR and step count data using the Fitbit application programming interface. To exclude outliers, we refined data selection to omit instances where HR was below 40 bpm, as recommended by the American Heart Association [29,30]. We extracted feature statistics such as avg, SD, min, med, and max HR within a 5-minute window to explore the relationship between HR and marijuana intoxication levels ("moderate-intensive," "low," and "none"). Resting HR was

defined as HR data collected when the participant was sedentary (ie, no steps taken) for more than 5 minutes. To further analyze HR patterns related to marijuana intoxication, we examined the degree of peakedness (kurtosis) and asymmetry (skewness) in HR data, as these features may reveal physiological changes associated with marijuana intoxication [31]. For more details, refer to Table S2 in Multimedia Appendix 2.

Ground Truth and Labeling Sensor Data

To accurately label the collected sensor data, we defined the duration of marijuana use episodes as those equal to or less than 3 hours, based on reported start and end times. We excluded 3 hours of sensor data following the reported end time to account for the continued effects of marijuana, even when participants reported a subjective intoxication level of 0. For example, if marijuana use was reported from 6 PM to 6:30 PM, data from 6:30 PM to 9:30 PM were excluded to account for residual effects. We also excluded data from 30 minutes before the reported start time to account for potential delays in self-reporting, based on pilot study findings that delays could range from 5 to 15 minutes. To collect nonmarijuana data, we randomly sampled sensor data from days when participants did not use marijuana (ie, nonmarijuana days). These samples were labeled using morning, afternoon, and evening surveys in which participants reported "no" to the ESM item "Did you smoke marijuana since the last report?" and indicated that the last use was more than 5 hours before the ESM time stamp (Figure 2).

We aimed to capture acute intoxication versus nonuse, classifying intoxication levels into three categories: 0 as "not intoxicated," 1-3 as "low intoxication," and 4-10 as "moderate-intensive intoxication" (MI). In total, we labeled 32,722 sensor stream samples (5-minute windows) as "not intoxicated" (154 from self-initiated survey coded as 0 high, and 32,586 from time-based self-reports), 423 samples as "low intoxication" (ratings between 1 and 3) and 772 samples as "moderate-intensive" (ratings between 4 and 10, with 10 indicating "a lot").

Data from smartphones and Fitbit resulted in two datasets of different sizes. To ensure consistency, we down-sampled the smartphone dataset to include only samples overlapping with Fitbit data during the same time frames. This resulted in three datasets: (1) eXtreme Gradient Boosting (XGBoost)-Mobile: mobile phone only, (2) XGBoost-Fitbit: Fitbit-only, and (3) XGBoost-MobiFit: combined mobile and Fitbit data. The rationale for choosing Machine Learning (ML) models is detailed in Multimedia Appendix 3 and model comparison with different classifiers can be found in Multimedia Appendix 4.



Figure 2. Marijuana use episodes and labeling principle.



ML Pipeline

Feature Selection

We began data analysis by randomly partitioning the labeled sensor data into training (80%) and test (20% holdout) datasets. As shown in Figure 3, we first calculated Pearson correlation coefficients between features in the training dataset to identify highly covariant feature pairs (correlation coefficients >0.9) [32]. We then systematically removed one feature from each pair to reduce redundancy and improve model performance by retaining the most relevant and independent features. Next, we selected statistically significant features with a Gini coefficient importance [33] greater than 0.005. Details can be found in Multimedia Appendix 2.

Figure 3. Study overview. AI: artificial intelligence; HR: heart rate; SHAP: Shapley Additive exPlanations; SMOTE: Synthetic Minority Over-Sampling Technique; XGBoost: eXtreme Gradient Boosting Machine.



XSL•FO

Hyper-Parameter Tuning and Cross-Validation

As shown in Figure 3, during hyper-parameter tuning in the training dataset, we used cross-validation to randomly leave 10% of the samples out, training the model on the remaining 90% and testing on the withheld 10%. We used the Synthetic Minority Over-Sampling Technique [34] to ensure equal representation across all classes. We further optimized model performance with a Bayesian-optimization-driven method called Optuna [35] to select the best combination of hyperparameters and 10-fold cross-validation on models with Optuna-optimized hyperparameters.

For the final model evaluation, we used the reserved test data (20% unseen data, as shown in Figure 3). The model was evaluated on predictions made on the test data. Finally, as shown in Figure 3 (right column), we conducted an XAI analysis to better understand the decision-making process of our final predictive model. We generated SHapley Additive exPlanations (SHAP) on the unseen test data to ensure our findings were explainable for data the model had not seen.

Model Evaluation Metrics

We evaluated model performance using F_1 -score, recall, and precision, and selecting the best model based on the F_1 -score [36]. Low precision indicates too many false positives (ie, detecting intoxication when there is none), here we would mistakenly intervene or notify the participant. Low recall indicates too many false negatives (ie, not detecting intoxication when it occurs), potentially leading to unsafe behaviors such as impaired driving. Therefore, while we prioritize the F_1 -score, we also consider precision and recall.

Given our imbalanced samples, we used the area under the curve (AUC) metric, which provides a robust evaluation across all classification thresholds and is resilient to class imbalance.

XAI: Interpretation Approaches for Black-Box ML Models

To enhance algorithmic transparency, we used SHAP, a widely used interpretability method for ML models [37,38]. SHAP explains how specific data features influence model predictions, providing insights into the model's decision-making process. We identified the top 30 most significant features associated with marijuana intoxication reports, including their importance scores and visual summaries calculated by SHAP (see "Key Features Contributing to Model Performance" under the Results section). XGboost was selected due to its superior performance compared to other classifiers. The use of tree SHAP in this context reduces the computation time for SHAP values from exponential to polynomial [37].

Results

Timing, Duration, and Rating of Subjective Marijuana Intoxication

During the 30-day period, participants averaged 14 (SD 8.59) days of active participation. A total of 129 ESM self-initiated reports of marijuana use met the criteria for inclusion in the analysis: 101 reports of subjective marijuana intoxication (feeling high rated 1-10 out of 10) and 28 reports of feeling not high (0). Events not involving marijuana use were assigned a high rating of 0.

Tables 1 and 2 show the distribution of self-reported subjective marijuana intoxication across participants. Most episodes of intoxication (n=75) lasted between 30 minutes and 3 hours, with 54 episodes lasting up to 30 minutes (Table 1). Marijuana use was most often reported between 10 PM and 11 PM (n=24). Table 2 shows the distribution of ESM responses throughout the day. The average response latency to an ESM prompt expired. Most self-initiated reports of marijuana use occurred in the evenings: 14% (n=18) between 6 PM and 9 PM, and 39% (n=50) between 9 PM and midnight. On average, young adults rated their feeling of being high at 3.63 (SD 2.72) out of 10 when using marijuana (Table 3).

 Table 1. Distribution of the duration of self-reported marijuana use episodes (n=129) across participants.

<u>i</u> J				
Duration ^a (hours)	Number of events			
<0.5	54			
<1	20			
<1.5	23			
<2.0	13			
<2.5	13			
<3	6			

^aDuration refers to the window of smoking episodes. From small (30 minutes) to relatively large windows (3 hours).



Table 2. Distribution of the start time of marijuana use episodes during the day (n=129).

Clock time (hours)	Number of events
0-1	7
1-2	8
2-3	2
3-4	0
4-5	0
5-6	0
6-7	0
7-8	1
8-9	0
9-10	5
10-11	8
11-12	2
12-13	6
13-14	6
14-15	5
15-16	4
16-17	3
17-18	4
18-19	5
19-20	6
20-21	7
21-22	10
22-23	24
23-0	16

Table 3. Distribution of self-reported "feeling high" during marijuana use.

High rating ^a	Number of events
0	28
1	9
2	9
3	17
4	14
5	14
6	17
7	10
8	7
9	4
10	0

 a^{0} -10 scale representing an intensity of feeling high, 10=a lot from the self-initiated reports of marijuana use. In our study, a value of 0 for the high report is labeled as "no-intoxication."

XSL•FO RenderX

Model Comparison: Mobile Only, Fitbit Only, and Mobile and Fitbit Integration

The first part of our analysis aimed to determine whether smartphone sensor features alone could be used for real-time detection of subjective marijuana intoxication and whether adding Fitbit data would improve model performance, justifying the added complexity of Fitbit data collection. We compared three ML models using the XGBoost classifier: (1) smartphone sensors only (XGBoost-Mobile), (2) Fitbit features only (XGBoost-Fitbit), and (3) a combined model using smartphone and Fitbit features (XGBoost-MobiFit).

Among the 3 models tested, the XGBoost-MobiFit model, which integrates smartphone and Fitbit data, had the best performance, achieving 99% accuracy, 92% precision, 79% recall, 85% F_1 -score, and 99% AUC on the test dataset (Figure 4 and Table 4). These metrics indicate the XGBoost-MobiFit model's superior ability to accurately identify MI compared to low-intoxication and not-intoxicated states. While the XGBoost-Fitbit performed reasonably well, it did not match the performance of the XGBoost-MobiFit model in detecting marijuana intoxication. XGBoost-Fitbit achieved accuracy of 98%, 79% precision, 70% recall, 74% F_1 -score, and 97% AUC. These results suggest that using only Fitbit data may not be as effective as combining it with smartphone sensor data for detecting subjective marijuana intoxication. Based on these findings, the added burden of wearing and charging the Fitbit device seems justified in future deployments. The combined model (XGBoost-MobiFit) demonstrated improved performance in detecting subjective marijuana intoxication compared to using smartphone or Fitbit data alone.

Combining Fitbit data with mobile data resulted in a significant improvement over the Fitbit-only model. The mobile-only model achieved an AUC of 96%, an F_1 -score of 72%, a recall of 75%, and a precision of 70%. These results indicate that including Fitbit data adds value beyond what can be achieved with smartphone-based sensor data alone, as evidenced by a 13% improvement in F_1 -score.

In summary, three key findings emerged: the XGBoost-Mobile model had the lowest performance (F_1 -score=0.72, recall=0.75, precision=0.70); the XGBoost-Fitbit model (F_1 -score=0.74, recall=0.70, precision=0.79) generally performed lower than the combined model; and the XGBoost-MobiFit model was the best performer with an F_1 -score of 0.85, recall of 0.79, and precision of 0.92. As highlighted earlier, high precision and recall are critical so we focused on the F_1 -score to identify the best-performing model. The model comparison with different classifiers is provided in Multimedia Appendix 4.

Figure 4. Model comparison to detect acute marijuana intoxication "low-intoxicated" (rating=1-3) versus "moderate-intensive intoxicated" (rating= 4-10) versus "not-intoxicated" (rating=0). XGBoost-MobiFit: phone sensors and Fitbit (AUC=0.99; accuracy=0.99; left), XGBoost-Mobile: smartphone-based sensors (samples overlapping with Fitbit; AUC=0.96; accuracy=0.97; middle) and XGBoost-Fitbit: Fitbit only (AUC=0.97; accuracy=0.98; right). AUC: area under the curve; ROC: receiver-operating characteristic curve; XGBoost: eXtreme gradient boosting.



Table 4. Comparison of three XGBoost models using features selected in detecting moderate-intensive marijuana intoxication, low-intoxication, and not-intoxicated classes on the test dataset.

Machine learning model	AUC ^a	F ₁ -score	Recall	Precision	Accuracy
XGBoost-MobiFit	0.99	0.85	0.79	0.92	0.99
XGBoost-Mobile	0.96	0.72	0.75	0.70	0.97
XGBoost-Fitbit	0.97	0.74	0.70	0.79	0.98

^aAUC: area under the curve.

Understanding Model Performance in Detecting the Risk State of "Moderate and Intensive Marijuana Intoxication"

For predicting the MI class alone, the MobiFit model outperformed the mobile and Fitbit-only models, exhibiting a substantial improvement in the F_1 -score of 20% and 18%,

https://ai.jmir.org/2025/1/e52270

RenderX

respectively (Table 5). This improvement in F_1 -score highlights the benefits of integrating data from both devices: enhanced precision and recall for the MI class compared to the not-intoxicated (N) and low-intoxicated (L) classes (Table 6). The XGBoost-Mobile model exhibited a notably high false negative rate for instances labeled as "not-intoxicated," often misclassifying them as "moderate-intensive intoxicated."

However, it showed better accuracy in distinguishing "low-intoxicated" instances. In contrast, the XGBoost MobiFit model demonstrated a higher true positive rate compared to the other models, accurately identifying 76% of MI samples among the total samples belonging to that class. While the XGBoost-Mobile and Fitbit models achieved recall rates of 61% and 63% in predicting MI, they incorrectly predicted 56 and 53 out of 143 actual MI samples as other classes. In comparison, the best-performing MobiFit model achieved 108 true positives out of the 143 actual MI samples. The higher precision of the MobiFit model further supports its superior performance, though there remains room for improvement as it missed 35 samples, as shown in Table 6.

Table 5. Performance comparison of three XGBoost^a models in detecting the subjective sense of moderate-intensive marijuana intoxication class.

ML ^b model	MI ^c precision	MI recall	MI F_1 -score	MI AUC ^d
XGBoost-MobiFit	0.89	0.76	0.82	0.99
XGBoost-Mobile	0.64	0.61	0.62	0.96
XGBoost-Fitbit	0.65	0.63	0.64	0.98

^aXGBoost: eXtreme Gradient Boosting.

^bML: machine learning

^cMI: moderate-intensive intoxication.

^dAUC: area under the curve.

	Predicted		
	N ^a	L ^b	MI ^c
XGBoost ^d -MobiFit			
Actual			
Ν	6541	7	13
L	29	50	1
MI	35	0	108
XGBoost-Mobile			
Actual			
Ν	6452	59	50
L	28	52	0
MI	56	0	87
XGBoost-Fitbit			
Actual			
Ν	6499	14	48
L	41	39	0
MI	52	1	90

^aN: not-intoxicated.

^bL: low-intoxication.

^cMI: moderate-intensive intoxication.

^dXGBoost: eXtreme Gradient Boosting.

Key Features Contributing to Model Performance

Overview

To explore the algorithms' performance in predicting the MI class, we used SHAP summary visualizations [37,38] to identify patterns of acute marijuana intoxication. We determined the key features contributing significantly to the model's predictions

based on mean absolute SHAP values across all instances, with a focus on the MI class.

Figures 5 and 6 present the SHAP visualizations. In Figure 5, the length of each bar on the left indicates the feature's contribution to the model, with longer bars signifying a stronger influence on the outcome. The SHAP summary plots on the right of Figure 5 illustrate how features influence the MI prediction class, with the strongest influence at the top. The

color shading indicates the direction of the feature's effect, with blue for low values, purple for median values, and red for high values. Plots extending to the left indicate a negative contribution to the prediction, while those extending to the right positively contribute to MI predictions.

Figure 5. Explanations generated by SHAP summary plot. Impact of features on best performing XGBoost-MobiFit model (left) and binary model output identifying moderate-intensive intoxication (MI; SHAP>0) from nonmoderate-intensive intoxication (N and L) classes (SHAP<0; right). HR: heart rate. SHAP: SHapley Additive exPlanations; WTSD: weighted stationary latitude and longitude standard deviation; XGBoost: eXtreme Gradient Boosting.





Figure 6. Explanations generated by SHAP summary plot. Impact of features on XGBoost-Mobile model (top left) and binary model output identifying MI (SHAP>0) from nonmoderate-intensive intoxication (N and L) classes (SHAP<0; top right), impact of features on XGBoost-Fitbit model (bottom left) and binary model output identifying MI (SHAP>0) from nonmoderate-intensive intoxication (N and L) classes (SHAP<0; bottom right). MI: moderate-intensive intoxication; SHAP: SHapley Additive exPlanations; WTSD: weighted stationary latitude and longitude standard deviation; XGBoost: eXtreme Gradient Boosting.



Impact of Average Key Features on Model Output Magnitude

The top five influential features in detecting the three classifications (Figure 5, left) and affecting the MI outputs (Figure 5, right) included time of day, radius of gyration,

```
https://ai.jmir.org/2025/1/e52270
```

XSL•FO RenderX minimum HR, day of the week, and minutes awake during sleep. Among physical activities and physiological signals, a diverse range of features extracted from various sensors, including those beyond time-based attributes from both mobile and Fitbit combined sensors, was chosen as the top 30 crucial elements for distinguishing between not-intoxicated (N), low-intoxication
(L), and MI. The SHAP value, signifying the average impact magnitude on the model's output, played a pivotal role in this determination (Figure 5, left).

Impact of Unique Key Features on Mobile and Fitbit Model Outputs

Similar to the best-performing MobiFit model, the Mobile model (Figure 6) highlighted key features with overlapping impacts on the model's outcomes. The only exception was in specific movement and environmental context features, as shown in the top left and right graphs of Figure 6. However, the Fitbit model showed a more significant impact on HR features, with all four HR features ranking within the top 10 for all three classes (shown in the bottom-left graph in Figure 6), and for the MI classes compared to the non-MI classes (bottom-right graph in Figure 6).

A partial dependence plot (PDP) in Figure 7 illustrates the overall relationship between a feature and the outcome. The vertical axis represents SHAP values, signifying the effect of the chosen feature on predictions, while the horizontal axis represents actual feature values across instances. Each point represents an instance's feature value and its corresponding SHAP value. An upward PDP slope indicates a positive impact of the feature on MI prediction, while a downward slope indicates a negative impact. The surface on the PDP plot (eg, min HR and sum of moving minutes in Figure 7, top left) shows the combined impact of the two features on MI predictions, with greater values corresponding to increased prediction values.

In the following section, we introduce the key features contributing to MI, including elevated and fluctuating HR, reduced large-scale movement patterns, increased ambient noise and voice energy, and extended sleep patterns.

Key Features Explaining MI

Overview

To specifically examine the influence of key features on the "risk" state of MI, we present comprehensive details for each key feature within the model.

Figure 7. Interaction effects of total minutes spent moving on minimum HR values (top left), SD (top middle), and skewness (top right) of HR, and an explanation of skewness [39] (bottom). HR: heart rate; SHAP: SHapley Additive exPlanations.



Elevated and Fluctuating HRs

We investigated the impact of recent physical activity (measured as the sum of minutes spent moving based on Fitbit data) on HR in relation to self-reported marijuana intoxication using a PDP. The SHAP values for minimum HRs showed significant elevation, with an average increase from approximately 80 bpm to peaks of 90 bpm and reaching up to 100 bpm (ranging from 60 to 120 bpm, with a few data points exceeding 120 bpm). These elevated HRs corresponded to moderate-intensive self-reported marijuana intoxication (SHAP value>0) in young adults compared to other classes (not- and low-intoxicated).

The SHAP values clearly indicate a positive increase in minimum HR associated with a higher likelihood of

```
https://ai.jmir.org/2025/1/e52270
```

RenderX

self-reported MI, irrespective of the impact of the sum of minutes spent moving. The total movement time during self-reported MI influenced the rise in minimum HR, as shown in Figure 7 (top left), where the red values represent a maximum of 5 minutes of movement (our analysis uses 5-minute windows). While HR can fluctuate due to various factors, including physical activity, substance use (eg, alcohol), caffeine, meals, and mental state (eg, stress and anxiety), further research is needed to explore these additional influences.

In brief, patterns for the SD of HRs exhibited fluctuations, but, in general, showed an increase when young adults reported MI (Figure 7, top middle). Negative skewness (indicating a "left-skewed" distribution) in HR was consistently associated with MI. This skewness suggests that there were more HR data

points on the right side of the mean (indicating that the median was greater than the mean), leading to a distribution stretched toward higher HR values (Figure 7, top right).

Decreased Large-Scale Movements

During MI, individuals showed a tendency for limited large-scale movement, often restricted to a radius of

Figure 8. Influence of radius of gyration (unit: meters). SHAP: SHapley Additive exPlanations.





Elevated Surrounding Noise Energy

Interestingly, while the variance in environmental noise energy increased (with data points deviating further from the mean), the average noise energy decreased, though it exhibited an overall upward trend (Figure 9, left). Instances of MI were associated with increased noise variability (calculated based on the amplitude of audio samples), followed by a subsequent reduction (Figure 9, right). Analyzing ambient sounds provides insights into the environmental context where individuals reporting MI might be located. This could include situations such as marijuana smoking, socializing with friends, or engaging with media like television or music. Although GPS-generated features were the primary indicators, MI may or may not be directly linked to specific locations such as shared social spaces (eg, lounges) or entertaining venues (eg, bars, pubs, or clubs). Nevertheless, it remains plausible that young adults reporting MI may choose to stay in noisy environments.

Figure 9. Influence of mean (left) and SD (right) noise energy (unit: Joule). SHAP: SHapley Additive exPlanations.



Prolonged Sleep Patterns

Distinct sleep patterns were linked to episodes of self-reported MI. Individuals who reported MI demonstrated extended sleep durations, spanning approximately 8 to 11 hours (Figure 10, left) the day before self-reported intoxication. In contrast, instances with low or no reported intoxication generally corresponded to healthy sleep durations, averaging around 6-7 hours, with some patterns as short as 2 hours.

There was also a positive correlation between the duration of minutes awake after falling asleep and self-reported MI, particularly when the period involved less than 50 minutes of wakefulness. However, an increase in extended minutes awake

after falling asleep (if >50 minutes, extending beyond approximately an hour) did not show any significant association with a likelihood of MI (Figure 10, middle). Regarding sleep start times, the data indicated peaks at both 11 PM and early morning hours, with a rise in sleep start times continuing until around 4 AM (Figure 10, right).

In summary, elevated minimum HR values were clearly linked to a higher likelihood of self-reported MI. However, we observed that GPS-travel patterns (macromovements) did not appear to increase during self-reported marijuana intoxication. Interestingly, extended sleep hours and minutes awake during sleep [40] the day before self-reported marijuana intoxication were associated with MI.

Figure 10. Total sleep duration (left), minutes awake during sleep (middle), and sleep start time (right). SHAP: SHapley Additive exPlanations.



Additional Analyses for Real-World Feasibility

To enhance the practicality of our ML model in real-world settings, we conducted supplementary analyses to evaluate our top-performing model, the XGBoost-MobiFit model, under different scenarios. These scenarios involved: (1) excluding GPS-derived travel data due to potential privacy concerns or GPS deactivation; (2) excluding sleep data in cases where users did not provide sleep information; and (3) excluding both GPS-derived travel and sleep data. This approach aims to explore the feasibility of offering more flexible data collection options, potentially addressing privacy concerns and incomplete data issues.

excluding In brief. **GPS**-derived features (XGBoost-MobiFit-GPS excluded) resulted in a 15% decrease in the F_1 -score compared to the best model, with a 10% reduction in sensitivity (recall). Excluding sleep data (XGBoost-MobiFit-Sleep excluded) led to a 24% decrease in the F_1 -score compared to the best model. When both GPS and sleep features were excluded (XGBoost-MobiFit-GPS-Sleep excluded), the model experienced a 16% reduction in F_1 -score and showed the lowest recall for identifying self-reported MI classes compared to the best-performing model. Please refer to Multimedia Appendix 5 for a detailed description of the additional analyses and results.

Discussion

Overview

The ability to detect subjective reports of acute marijuana intoxication in natural environments using mobile sensors has the potential to enable just-in-time interventions [41] to reduce

```
https://ai.jmir.org/2025/1/e52270
```

marijuana-related harms. To the best of our knowledge, this is the first study that demonstrates the impact of integrating smartphone-based and wearable sensor features on the enhancement of the performance and interpretability of algorithms in detecting acute marijuana intoxication in naturalistic environments.

As hypothesized, we found that the XGB-MobiFit model, which combined smartphone sensor data with Fitbit features outperformed models that used only mobile or only Fitbit data. By integrating sensors from both smartphones and wearable devices, our best-performing algorithm balances specificity and sensitivity on unseen samples, enabling interpretable, transparent, and unobtrusive detection of acute subjective marijuana intoxication in natural environments. This opens up opportunities for real-time monitoring in everyday settings and the implementation of just-in-time adaptive interventions.

XAI visualizations supported our second hypothesis, highlighting HR, GPS, and physical movement data as key features that contributed to self-reported marijuana intoxication predictions. These findings were observed beyond the influences of simply applying time of day and day of the week features (ranked 1st and 4th, respectively), as validated in [11], particularly during instances of self-reported subjective marijuana intoxication in naturalistic environments.

Interpretable Behavioral and Physiological Signals of Marijuana Intoxication in Real-World Settings

To explain the results of the black-box ML models to detect marijuana intoxication in everyday settings, our study integrated sensors from smartphones and a wearable device, identified key sensor features, and used XAI to facilitate the interpretation of model results. The findings are consistent with prior research

conducted in controlled laboratory settings, which consistently found an acute increase in resting HR following marijuana use [12-14]. Our results suggest the potential for HR with behavioral factors to detect marijuana intoxication "outside of laboratory settings" using off-the-shelf devices in naturalistic environments. While many factors can affect HR in daily life, this study yielded significant HR features and insights from the elevated HR patterns during self-reported acute marijuana intoxication. Future research could explore associations between HR and other physiological and behavioral indicators of marijuana use, such as respiration, to better capture marijuana intoxication in natural environments [42].

The use of XAI visualization could help increase transparency and accountability when conducted as part of a substance use detection system [43, 44]. It is promising to use XAI as it enables researchers and clinicians to understand how algorithms arrive at decisions and identify key behavioral and physiological attributes, providing opportunities to improve detection accuracy and enhance trust in the algorithm over time.

Real-Time Detection and Intervention Potential

Compared to an average 30-minute marijuana episode, the 5-minute window used in the best-performing model is small enough to predict marijuana intoxication in near real-time. Detecting marijuana intoxication in near real-time promotes just-in-time intervention, which serves as a crucial first step toward reducing possible marijuana-related harm in a timely manner.

Our best detection model is unlikely to misclassify a "high" state as not high, which demonstrates the potential for using our detection algorithm with unseen data in real-world contexts. On the unseen test set, we obtained 85% precision (92% precision for 3 classes) in specifically identifying self-reported moderate-intensive marijuana intoxication. Passive sensing using smartphone-based sensors has been investigated in the context of alcohol intoxication [25,26,43], and here we extend this research to self-reported marijuana intoxication [11] beyond smartphone-based sensors, which could ultimately be useful for JIT interventions [41] to reduce marijuana-related harm. The value to society and individuals of reducing marijuana-related harm is clear. If individuals choose to use such a personal detection system, they will need to keep their phone charged and with them when using marijuana and wear a device (eg, Fitbit) and keep it charged as well.

For real-time modeling using the XGBoost algorithm, deploying the estimated model onto a computing device is an indispensable phase. We envisage two primary deployment scenarios: first, local assessments can be generated by deploying the model directly onto users' devices, such as smartphones. This approach ensures seamless functionality even without an internet connection but requires adequate storage and computational capacity. Second, cloud-based computation can be used. While this approach relies on a stable internet connection, it effectively offloads the computational burden from the user's device. Real-world applications introduce pragmatic considerations such as battery longevity, which could be affected by the model's continuous operation, and user privacy during data transmission and generation of model results.

```
https://ai.jmir.org/2025/1/e52270
```

Therefore, a comprehensive assessment of the model's feasibility in real-time operational settings is important. Our proposed generalized model, designed to operate across a diverse demographic spectrum rather than relying on individual-specific (idiographic) models, offers advantages in terms of scalability and practicality.

Privacy Considerations and User-Centric Configuration Choices

To highlight the benefits of combining sensor features from both smartphone and wearable devices while addressing potential privacy concerns, particularly related to location data, we aim to offer participants additional configuration choices rather than study withdrawal. For example, participants can deactivate GPS sensors if desired. This is demonstrated by our testing of the best-performing model, XGBoost-MobiFit, where we excluded location features. The analysis revealed a 15% (XGBoost-MobiFit-GPS excluded) decrease in F_1 -score from the best model. As proposed by Bae et al [43], collecting GPS data and using rounded GPS data extraction (ie, less precise location data) could be a viable approach. This avoids using raw latitude and longitude, which may contain sensitive information on specific locations. Researchers and clinicians could consider providing alternative options instead of completely disabling GPS, as GPS data contributes to the model's accuracy.

Moreover, to assess the efficacy of our top-performing model, we conducted tests after excluding sleep-related features (Multimedia Appendix 5). The analysis revealed a 24% (XGBoost-MobiFit-Sleep excluded) decrease in the F_1 -score compared to the best model's performance. While participants may benefit from the option to disable sensors when necessary, it is important to note that this could potentially decrease the model's ability to detect marijuana intoxication.

By building a system that prioritizes privacy and user autonomy, we can provide a valuable tool to reduce marijuana-related harm to individuals and society. Ultimately, each person will have to decide for themselves whether the benefits of a detection and intervention system outweigh the tradeoffs in minimizing possible marijuana-related harms to themselves and the broader community.

Limitations and Future Work

The first limitation of this study is relying on self-reporting as the ground truth, which may be subjective. This study extends prior ESM work, which codes self-reported marijuana use as yes or no [45], by asking participants to rate marijuana intoxication from 0 to 10, which may be subject to recall or other biases in reporting. The broad categorization might overlook nuanced differences within three categories: low-intoxication (1-3), moderate-intensive marijuana intoxication (4-10), and not-high (0), which could affect the accuracy of the classifiers. Future analyses examining the performance of mobile and wearable sensors against different thresholds for a subjective marijuana intoxication outcome could be valuable.

Another limitation was the size, diversity, and duration of the participants in the study. Since the participants were all young adults, the finding may not be generalizable to a broader age group. In addition, the level of compliance (63%) in completing the morning, afternoon, and evening surveys is relatively low. Thus, it is unclear whether all episodes of marijuana use were reported by participants, which could limit model performance. However, since there is no real-time accessible biological testing method at the time of publication, validating self-reported data with the current method still represents the best alternative. The current findings warrant future replication in a larger and more diverse group of participants over a longer period to address the limitations and validate the findings.

In addition, our model performed best when tested on the same participants it was trained on (with no overlap between training and testing data). While this has a valid use case, it assumes that we can always collect labeled training data for participants for whom we would like to apply the model. By applying more testing data, using more sophisticated sensor features, and better model tuning, future models could improve generalization over unseen testing participants. The HR data only holds significance when examined together with activity data. An acute increase in HR by itself is nonspecific and may not be associated with marijuana use or intoxication. False alarms triggered by the algorithm could erode trust in an automated system, whereas low sensitivity to actual marijuana use could result in marijuana-related harm. Therefore, it is important to investigate the interplay between human activities associated with marijuana intoxication and physiological signals in a larger population, and how these interactions can contribute to intervention delivery in real-world contexts.

Finally, it is crucial to acknowledge that the potential impact of polysubstance use on the interpretation of physiological signals associated with self-reported cannabis intoxication was not included. While ESM is used to collect information on the use of other substances, our analysis did not account for the effects of polysubstance use due to the limited scope of the study. The presence of polysubstance use could potentially confound the physiological signals attributed to marijuana. This may lead to inaccuracies in our algorithm, particularly in distinguishing between marijuana intoxication and the effects of other substances. Thus, while our study provides valuable insights into self-reported marijuana intoxication, it has limitations in addressing the full spectrum of real-world polysubstance use. Future research should include developing algorithms that can differentiate between the physiological signals associated with different substances, including polysubstance use.

Conclusions

Our study demonstrates that integrating features from smartphone-based sensors and wearable devices significantly improves the detection of self-reported marijuana intoxication in natural environments among young adults. The XGBoost-MobiFit model, which combines data from both smartphone sensors and wearable devices, achieved an F_1 -score of 0.85 in detecting moderate to intensive self-reported marijuana intoxication, outperforming models that relied solely on smartphone sensors. The results suggest that incorporating wearable device data enhances the XGBoost model's performance by 13%, justifying the additional complexity of using wearable devices among young adults.

Key features contributing to the detection of self-reported "MI" included an acute increase in HR (measured by Fitbit), macromovement indicators (derived from GPS data), and prolonged sleep patterns the night before self-reported marijuana intoxication (measured by Fitbit).

Future research should focus on refining the algorithms that integrate smartphone and Fitbit sensor data in larger, more diverse samples. In addition, exploring how these algorithms, informed by XAI, can support the development of just-in-time interventions for clinicians is essential. Such interventions could offer context-adaptive, personalized strategies to minimize potential marijuana-related harms, such as intoxicated driving, therefore reducing the frequency and severity of acute marijuana-related incidents among young adults.

Acknowledgments

This study was supported by the National Institute on Drug Abuse (R21 DA043181/U01 DA056472), the Stevens Startup grant, and the Provost scholarship.

Authors' Contributions

SWB, TC, and AKD contributed to the design of the study and data collection. SWB, TZ, AKD, and RI processed the data, and SWB and TZ analyzed the data and developed the computational and explainable models. SWB drafted the initial manuscript, which was edited by TC, TZ, and AKD, and approved by all authors.

Conflicts of Interest

None declared.

Multimedia Appendix 1 Extracted features overview. [DOCX File , 40 KB - ai_v4i1e52270 app1.docx]



Multimedia Appendix 2 Feature selection. [DOCX File, 197 KB - ai v4i1e52270 app2.docx]

Multimedia Appendix 3 Rationale for machine learning model selection. [DOCX File , 34 KB - ai v4i1e52270 app3.docx]

Multimedia Appendix 4 Comparison of models with different classifiers. [DOCX File , 41 KB - ai v4i1e52270 app4.docx]

Multimedia Appendix 5 Privacy-preserving XGBoost-MobiFit models. [DOCX File , 42 KB - ai v4i1e52270 app5.docx]

References

- 1. Conroy DA, Kurth ME, Brower KJ, Strong DR, Stein MD. Impact of marijuana use on self-rated cognition in young adult men and women. Am J Addict 2015;24(2):160-165 [FREE Full text] [doi: 10.1111/ajad.12157] [Medline: 25864605]
- 2. Engineering National Academies of Sciences. The Health Effects of Cannabis and Cannabinoids: The Current State of Evidence and Recommendations for Research. Washington, DC: Academies Press; 2017.
- Scott JC, Slomiak ST, Jones JD, Rosen AFG, Moore TM, Gur RC. Association of cannabis with cognitive functioning in adolescents and young adults: a systematic review and meta-analysis. JAMA Psychiatry 2018;75(6):585-595. [doi: 10.1001/jamapsychiatry.2018.0335] [Medline: 29710074]
- 4. Phillips KT, Phillips MM, Lalonde TL, Tormohlen KN. Marijuana use, craving, and academic motivation and performance among college students: an in-the-moment study. Addict Behav 2015;47:42-47 [FREE Full text] [doi: 10.1016/j.addbeh.2015.03.020] [Medline: 25864134]
- 5. Pertwee RG. Handbook of Cannabis. United Kingdom: Oxford University Press; 2015.
- Ruojun LI, Emmanuel AGU, Ganesh B, Debra H, Ana A, Michael S. WeedGait: unobtrusive smartphone sensing of marijuana-induced gait impairment by fusing gait cycle segmentation and neural networks. 2019 Presented at: IEEE Healthcare Innovations and Point of Care Technologies (HI-POCT); November 22, 2019; USA p. 94. [doi: 10.1109/hi-poct45284.2019.8962787]
- Mishra RK, Sempionatto JR, Li Z, Brown C, Galdino NM, Shah R, et al. Simultaneous detection of salivary Δ-tetrahydrocannabinol and alcohol using a wearable electrochemical ring sensor. Talanta 2020;211:120757 [FREE Full text] [doi: 10.1016/j.talanta.2020.120757] [Medline: 32070607]
- 8. Pedersen ER, Hummer JF, Rinker DV, Traylor ZK, Neighbors C. Measuring protective behavioral strategies for marijuana use among young adults. J Stud Alcohol Drugs 2016;77(3):441-450. [doi: <u>10.15288/jsad.2016.77.441</u>] [Medline: <u>27172576</u>]
- Huestis MA, Smith ML. Cannabinoid markers in biological fluids and tissues: revealing intake. Trends Mol Med 2018;24(2):156-172. [doi: <u>10.1016/j.molmed.2017.12.006</u>] [Medline: <u>29398403</u>]
- 10. Bédard M, Dubois S, Weaver B. The impact of cannabis on driving. Can J Public Health 2007;98(1):6-11. [doi: 10.1007/bf03405376]
- Bae SW, Chung T, Islam R, Suffoletto B, Du J, Jang S, et al. Mobile phone sensor-based detection of subjective cannabis intoxication in young adults: a feasibility study in real-world settings. Drug Alcohol Depend 2021;228:108972-108716 [FREE Full text] [doi: 10.1016/j.drugalcdep.2021.108972] [Medline: 34530315]
- 12. Huber GL, Griffith DL, Langsjoen PM. The effects of marihuana on the respiratory and cardiovascular systems. Marijuana: An International Research Report. National Campaign Against Drug Abuse Monograph 7 (1988) 1988:123-134.
- 13. Maykut MO. Health consequences of acute and chronic marihuana use. Prog Neuro-Psychopharmacol Biol Psychiatry 1985;9(3):209-238. [doi: 10.1016/0278-5846(85)90085-5]
- 14. Zuurman L, Ippel AE, Moin E, van Gerven JMA. Biomarkers for the effects of cannabis and THC in healthy volunteers. Br J Clin Pharmacol 2009;67(1):5-21 [FREE Full text] [doi: 10.1111/j.1365-2125.2008.03329.x] [Medline: 19133057]
- Carreiro S, Smelson D, Ranney M, Horvath KJ, Picard RW, Boudreaux ED, et al. Real-time mobile detection of drug use with wearable biosensors: a pilot study. J Med Toxicol 2015;11(1):73-79. [doi: <u>10.1007/s13181-014-0439-7</u>] [Medline: <u>25330747</u>]
- Epstein DH, Tyburski M, Kowalczyk WJ, Burgess-Hull AJ, Phillips KA, Curtis BL, et al. Prediction of stress and drug craving ninety minutes in the future with passively collected GPS data. NPJ Digital Med 2020;3(1):26 [FREE Full text] [doi: 10.1038/s41746-020-0234-6] [Medline: 32195362]
- 17. Ferreira D, Kostakos V, Dey AK. AWARE: mobile context instrumentation framework. Front ICT 2015 Apr 20;2:1-9. [doi: <u>10.3389/fict.2015.00006</u>]

```
https://ai.jmir.org/2025/1/e52270
```

- Chung T, Bae SW, Mun E, Suffoletto B, Nishiyama Y, Jang S, et al. Mobile assessment of acute effects of marijuana on cognitive functioning in young adults: observational study. JMIR Mhealth Uhealth 2020;8(3):e16240 [FREE Full text] [doi: 10.2196/16240] [Medline: 32154789]
- Mokrysz C, Freeman T, Korkki S, Griffiths K, Curran HV. Are adolescents more vulnerable to the harmful effects of cannabis than adults? A placebo-controlled study in human males. Transl Psychiatry 2016;6(11):e961 [FREE Full text] [doi: 10.1038/tp.2016.225] [Medline: 27898071]
- 20. Spindle TR, Cone EJ, Schlienz NJ, Mitchell JM, Bigelow GE, Flegel R, et al. Acute effects of smoked and vaporized cannabis in healthy adults who infrequently use cannabis: a crossover trial. JAMA Netw Open 2018;1(7):e184841 [FREE Full text] [doi: 10.1001/jamanetworkopen.2018.4841] [Medline: 30646391]
- Mohr DC, Zhang MI, Schueller SM. Personal sensing: understanding mental health using ubiquitous sensors and machine learning. Annu Rev Clin Psychol 2017;13:23-47 [FREE Full text] [doi: <u>10.1146/annurev-clinpsy-032816-044949</u>] [Medline: <u>28375728</u>]
- 22. Saeb S, Lattie EG, Schueller SM, Kording KP, Mohr DC. The relationship between mobile phone location sensor data and depressive symptom severity. PeerJ 2016;4:e2537 [FREE Full text] [doi: 10.7717/peerj.2537] [Medline: 28344895]
- 23. Celebi ME, Celiker F, Kingravi HA. On Euclidean norm approximations. Pattern Recognit 2011;44(2):278-283. [doi: 10.1016/j.patcog.2010.08.028]
- 24. Denzil Ferreira. Com.Aware.Plugin.Studentlife.Audio_Final. Retrieved from GitHub. 2016. URL: <u>https://github.com/</u> <u>denzilferreira/com.aware.plugin.studentlife.audio_final</u> [accessed 2023-01-18]
- 25. Bae S, Chung T, Ferreira D, Dey AK, Suffoletto B. Mobile phone sensors and supervised machine learning to identify alcohol use events in young adults: implications for just-in-time adaptive interventions. Addict Behav 2018;83:42-47. [doi: 10.1016/j.addbeh.2017.11.039] [Medline: 29217132]
- Bae S, Ferreira D, Suffoletto B, Puyana JC, Kurtz R, Chung T, et al. Detecting drinking episodes in young adults using smartphone-based sensors. Proc ACM Interact Mobile Wearable Ubiquitous Technol 2017;1(2):1-36. [doi: <u>10.1145/3090051</u>] [Medline: <u>35146236</u>]
- 27. Byrnes HF, Miller BA, Wiebe DJ, Morrison CN, Remer LG, Wiehe SE. Tracking adolescents with global positioning system-enabled cell phones to study contextual exposures and alcohol and marijuana use: a pilot study. J Adolesc Health 2015;57(2):245-247 [FREE Full text] [doi: 10.1016/j.jadohealth.2015.04.013] [Medline: 26206448]
- 28. Chaix B. Mobile sensing in environmental health and neighborhood research. Annu Rev Public Health 2018;39:367-384 [FREE Full text] [doi: 10.1146/annurev-publhealth-040617-013731] [Medline: 29608869]
- Jensen MT, Suadicani P, Hein HO, Gyntelberg F. Elevated resting heart rate, physical fitness and all-cause mortality: a 16-year follow-up in the Copenhagen male study. Heart 2013;99(12):882-887 [FREE Full text] [doi: 10.1136/heartjnl-2012-303375] [Medline: 23595657]
- 30. American Heart Association. All about heart rate (pulse). Retrieved from. URL: <u>https://www.heart.org/en/health-topics/</u> high-blood-pressure/the-facts-about-high-blood-pressure/all-about-heart-rate-pulse [accessed 2024-11-05]
- Tara K, Sarkar AK, Khan MAG, Mou JR. Detection of cardiac disorder using MATLAB based graphical user interface (GUI). 2017 Presented at: 2017 IEEE Region 10 Humanitarian Technology Conference (R10-HTC); December 23, 2017; United States p. 440-443. [doi: 10.1109/r10-htc.2017.8288994]
- 32. Schober P, Boer C, Schwarte LA. Correlation coefficients: appropriate use and interpretation. Anesth Analg 2018;126(5):1763-1768. [doi: 10.1213/ANE.00000000002864] [Medline: 29481436]
- 33. Yitzhaki S, Schechtman E. The Gini Methodology: A Primer on a Statistical Methodology. New York: Springer; 2013:11-31.
- 34. Chawla NV, Bowyer KW, Hall LO, Kegelmeyer WP. SMOTE: synthetic minority over-sampling technique. J Artif Intell Res 2002;16:321-357. [doi: 10.1613/jair.953]
- 35. Takuya A, Shotaro S, Toshihiko Y, Takeru O, Masanori K. Optuna: a next-generation hyperparameter optimization framework. 2019 Presented at: KDD '19: The 25th ACM SIGKDD Conference on Knowledge Discovery and Data Mining; July 25, 2019; United States. [doi: 10.1145/3292500.3330701]
- 36. Powers D. Evaluation: from precision, recall and F-measure to ROC, informedness, markedness and correlation. J Mach Learn Technol 2011;2:37-63 [FREE Full text]
- 37. Lundberg SM, Erion G, Chen H, DeGrave A, Prutkin JM, Nair B, et al. From local explanations to global understanding with explainable AI for trees. Nat Mach Intell 2020;2(1):56-67. [doi: 10.1038/s42256-019-0138-9] [Medline: 32607472]
- 38. Lundberg SM, Lee SI. A unified approach to interpreting model predictions. Adv Neural Inf Process Syst 2017;30:4765-4774 [FREE Full text]
- 39. Skewness. 2023. URL: https://en.wikipedia.org/wiki/Skewness [accessed 2023-08-28]
- 40. Shrivastava D, Jung S, Saadat M, Sirohi R, Crewson K. How to interpret the results of a sleep study. J Community Hosp Intern Med Perspect 2014;4(5):24983. [doi: <u>10.3402/jchimp.v4.24983</u>] [Medline: <u>25432643</u>]
- 41. Joshua MS, Kristin EH. Is providing mobile interventions" just-in-time" helpful? An experimental proof of concept study of just-in-time intervention for stress management. 2016 Presented at: 2016 IEEE Wireless Health (WH); October 27, 2016; USA. [doi: 10.1109/wh.2016.7764561]
- 42. NIDA. What are marijuana's effects on other aspects of physical health?. URL: <u>https://nida.nih.gov/research-topics/cannabis</u> <u>-marijuana</u> [accessed 2023-08-10]

- 43. Bae SW, Suffoletto B, Zhang T, Chung T, Ozolcer M, Islam MR, et al. Leveraging mobile phone sensors, machine learning, and explainable artificial intelligence to predict imminent same-day binge-drinking events to support just-in-time adaptive interventions: algorithm development and validation study. JMIR Form Res 2023;7:e39862 [FREE Full text] [doi: 10.2196/39862] [Medline: 36809294]
- 44. Zhang T, Chung T, Dey A, Bae SW. Exploring Algorithmic Explainability: Generating Explainable AI Insights for Personalized Clinical Decision Support Focused on Cannabis Intoxication in Young Adults. 2024 Int Conf Act Behav Comput (2024) 2024 May;2024. [doi: 10.1109/abc61795.2024.10652070] [Medline: 39600343]
- 45. Randi MS, Robin J, Mermelstein, Donald H. Ecological momentary assessment of working memory under conditions of simultaneous marijuana and tobacco use. 2016. URL: <u>https://doi.org/10.1111/add.13342</u>

Abbreviations

AUC: area under the curve
bpm: beats per minute
ESM: experience sampling method
HR: heart rate
MI: moderate-intensive intoxication
ML: machine learning
PDP: partial dependence plot
SHAP: SHapley Additive exPlanations
THC: delta-9 tetrahydrocannabinol
XAI: explainable artificial intelligence
XGBoost: eXtreme Gradient Boosting

Edited by K El Emam, B Malin; submitted 29.08.23; peer-reviewed by E Karoulla, Q Liu, I Liu; comments to author 09.11.23; revised version received 31.01.24; accepted 02.09.24; published 02.01.25.

Please cite as:

Bae SW, Chung T, Zhang T, Dey AK, Islam R Enhancing Interpretable, Transparent, and Unobtrusive Detection of Acute Marijuana Intoxication in Natural Environments: Harnessing Smart Devices and Explainable AI to Empower Just-In-Time Adaptive Interventions: Longitudinal Observational Study JMIR AI 2025;4:e52270 URL: https://ai.jmir.org/2025/1/e52270 doi:10.2196/52270 PMID:

©Sang Won Bae, Tammy Chung, Tongze Zhang, Anind K Dey, Rahul Islam. Originally published in JMIR AI (https://ai.jmir.org), 02.01.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Investigating the Classification of Living Kidney Donation Experiences on Reddit and Understanding the Sensitivity of ChatGPT to Prompt Engineering: Content Analysis

Joshua Nielsen¹, BS; Xiaoyu Chen², PhD; LaShara Davis³, PhD; Amy Waterman³, PhD; Monica Gentili¹, PhD

¹Department of Industrial Engineering, JB Speed School of Engineering, University of Louisville, Louisville, KY, United States

²Department of Industrial and Systems Engineering, School of Engineering and Applied Sciences, University at Buffalo, Buffalo, NY, United States ³Patient Engagement, Diversity, and Education Division, Department of Surgery, Houston Methodist Hospital, Houston, TX, United States

Corresponding Author:

Joshua Nielsen, BS Department of Industrial Engineering JB Speed School of Engineering University of Louisville 220 Eastern Parkway Louisville, KY, 40292 United States Phone: 1 5024891335 Email: joshua.nielsen@louisville.edu

Abstract

Background: Living kidney donation (LKD), where individuals donate one kidney while alive, plays a critical role in increasing the number of kidneys available for those experiencing kidney failure. Previous studies show that many generous people are interested in becoming living donors; however, a huge gap exists between the number of patients on the waiting list and the number of living donors yearly.

Objective: To bridge this gap, we aimed to investigate how to identify potential living donors from discussions on public social media forums so that educational interventions could later be directed to them.

Methods: Using Reddit forums as an example, this study described the classification of Reddit content shared about LKD into three classes: (1) present (presently dealing with LKD personally), (2) past (dealt with LKD personally in the past), and (3) other (LKD general comments). An evaluation was conducted comparing a fine-tuned distilled version of the Bidirectional Encoder Representations from Transformers (BERT) model with inference using GPT-3.5 (ChatGPT). To systematically evaluate ChatGPT's sensitivity to distinguishing between the 3 prompt categories, we used a comprehensive prompt engineering strategy encompassing a full factorial analysis in 48 runs. A novel prompt engineering approach, dialogue until classification consensus, was introduced to simulate a deliberation between 2 domain experts until a consensus on classification was achieved.

Results: BERT and GPT-3.5 exhibited classification accuracies of approximately 75% and 78%, respectively. Recognizing the inherent ambiguity between classes, a post hoc analysis of incorrect predictions revealed sensible reasoning and acceptable errors in the predictive models. Considering these acceptable mismatched predictions, the accuracy improved to 89.3% for BERT and 90.7% for GPT-3.5.

Conclusions: Large language models, such as GPT-3.5, are highly capable of detecting and categorizing LKD-targeted content on social media forums. They are sensitive to instructions, and the introduced dialogue until classification consensus method exhibited superior performance over stand-alone reasoning, highlighting the merit in advancing prompt engineering methodologies. The models can produce appropriate contextual reasoning, even when final conclusions differ from their human counterparts.

(JMIR AI 2025;4:e57319) doi:10.2196/57319

KEYWORDS

prompt engineering; generative artificial intelligence; kidney donation; transplant; living donor



Introduction

Background

Kidney transplantation is the gold standard treatment for patients with end-stage renal disease [1] and can be much more cost-effective than dialysis [2]. Record numbers of transplants have taken place in recent years, but a shortage of donors continues to exist despite the recent increase [3]. Currently, the median wait time for a transplant is approximately 4 years in the United States, and approximately 5000 patients die every year while being on the transplant waiting list [4]. Living donor kidney transplantation (LDKT) generally provides better outcomes than deceased donor transplants but requires that a potential living donor be made aware that they can donate to a specific patient with end-stage renal disease and offer to do so. Racial or ethnic minorities and patients of lower socioeconomic status are less likely to pursue and have living donors donate on their behalf [5,6].

National attitudes about LDKT are generally positive, although many do not know what a living donor undergoes when donating a kidney [7-10]. Recommendations to increase the living donor pool include reaching out more broadly to locate generous individuals motivated by social good to engage more individuals in considering living donation [11]. In addition, research suggests that disseminating education and information about living donation to broader audiences, beyond transplant centers, might increase the numbers of potential donors and recipients pursuing living donation [12,13]. However, identifying individuals dealing with kidney disease and considering whether to pursue LDKT or donate kidneys in their own lives can be difficult, especially when they have not started medical evaluation at a transplant center.

Locating individuals through social media forums discussing living kidney donation (LKD), such as those on Reddit or Twitter (the work herein was done before the platform being rebranded as X), maybe a way to identify individuals who are actively deciding whether to pursue LDKT or LKD outside of transplant centers [14]. While there are many different types of questions and comments related to LKD shared on the web, some people share their personal experiences and even invite people to "ask me anything." These findings motivated our main hypothesis that potential living donors can be identified from social media communities engaged in general discussions about LKD. In addition, understanding the personal experiences shared on these platforms can provide valuable insights into potential donors' needs and decision-making, enabling education and media campaigns to be better tailored for them.

The large volume and high complexity of unstructured natural language require an effective and efficient method that can

automate the identification of people sharing personal experiences with LKD. Fortunately, recent advances in natural language processing (NLP), particularly the transformer mechanism [15-19], enable the automatic understanding of personal experiences that were shared on the web social platforms. This study aimed to evaluate the transformer-based techniques to categorize these experiences on Reddit (Reddit, Inc). Specifically, we aimed to evaluate and compare (1) the one-shot classification model Bidirectional Encoder Representations from Transformers (BERT) [19], which required that we fine-tune the model using 1268 well-labeled samples, and (2) the zero-shot classification model ChatGPT (OpenAI), which required no fine-tuning for classification purposes. Comprehensive discussions on transformer-based models can be found in the study by Acheampong et al [20]. Much has been written about the capabilities and limitations of ChatGPT specifically [21]; however, we investigated the importance of prompt engineering when interfacing with it and other generative models applied to the field of organ donation for the first time.

Overview of Prompt Engineering

Prompt engineering has been defined as "the means by which LLMs are programmed via prompts" [22]. Reynolds and McDonell [23] framed the objective of prompt engineering as a discipline that seeks to answer the question, "What prompt will result in the intended behavior and *only* the intended behavior?" Historically, the best practice has been to give a small number of examples of how the task is to be done, known as few-shot prompting. Ray [21] suggested that for large language models (LLMs), few-shot prompting is better thought of as "locating an already-learned task rather than meta-learning." The implication is that the LLMs are large and robust enough that the models are inherently capable of completing NLP tasks, but their scale of capability may require using examples to "activate" the right parameters that will carry out the desired task in the prescribed manner.

However, this flexibility should also be understood as having dangers because LLMs can be "jailbroken." Jailbreaking LLMs is the practice of using prompt engineering to work around the boundaries imposed by the developers, such as OpenAI [24]. The practice of "red-teaming" is used by developers to identify weaknesses in the desired boundaries and adjust the model so that it is more defensible against previous vulnerabilities [25,26]. What is simultaneously exciting and problematic about this is that many techniques used to jailbreak LLMs are the same as those used for their most helpful, intended uses, that is, many of the same methods that allow us to get the best performance from an LLM can be the same ones that are used to bypass the safeguards. Table 1 provides an overview of prompt engineering methods derived primarily from the study by White et al [22].



Table 1. Overview of prompt engineering methods proposed by White et al [22].

Method	Purpose	Example prompts for LKD ^a
Few-shot prompting	Provide examples that illustrate how the task is to be completed	"Here is an example of a risk analysis from a living kidney donation sce- nario: [EXAMPLE]. Now, please provide a risk analysis for the following scenario."
Meta-language creation	Create a shorthand notation, abbre- viated language, or set of standard rules	"For this conversation, 'LKD' refers to living kidney donation, 'DT' refers to donor testing and 'RC' refers to recipient compatibility. Using this shorthand, describe the typical process of LKD."
Flipped interaction	The LLM ^b will ask questions to obtain the information	"I'm working on an algorithm to match donors with recipients in living kidney donation. What information do you need from me to help design this algorithm?"
Persona	Assign a persona to the LLM, usually that of an expert	"Pretend you are a leading surgeon specializing in living kidney donation. Provide your expert opinion on the latest surgical techniques."
Prompt refinement	Ensure that the LLM suggests better or more refined prompts	"I need to write code to analyze the success rates of different kidney matching algorithms. Could you suggest a more refined question or spe- cific details you need to assist me?"
Alternative approaches	Ensure that the LLM offers alterna- tive ways of accomplishing the task	"Describe three different methods for assessing donor-recipient compati- bility in living kidney donation."
Cognitive verifier	Subdivide a question into additional questions for a better answer	"To understand the ethical considerations in living kidney donation, what additional questions should I ask you to provide a comprehensive analysis?"
Fact checklist	Mitigate model hallucination by listing the facts	"After explaining the current trends in living kidney donation, list the facts or data sources you used in your response."
Template	Ensure that the LLM's output fol- lows a precise template	"Please answer in the following format: 'Living kidney donation is bene- ficial because [REASON 1], [REASON 2], and [REASON 3]'."
Gameplay	Create a game around a given topic	"Let's play a matching game. I will describe a recipient, and you suggest a suitable donor from the provided pool based on living kidney donation criteria."
Reflection (chain of thought [25])	Explain the rationale behind the given answers	"Explain the process of donor selection in living kidney donation in a step- by-step manner, detailing the reasoning behind each step."
Refusal breaker	Help users rephrase a question when they are refused an answer	"If you cannot provide personal patient data in living kidney donation, please guide me on how to rephrase my questions to obtain general infor- mation."
Context manager	Enable users to specify or remove context	"When discussing living kidney donation statistics, please consider only data from the last five years in the European region."
Recipe	Provide a sequence of steps given some partially provided ingredients	"I have patient medical records, compatibility testing results, and surgical schedules. Provide a sequence of steps to create an optimal living kidney donation matching algorithm."

^aLKD: living kidney donation.

^bLLM: large language model.

Reflection and chain of thought reasoning, in particular, have garnered much attention due to their powerful results, creating what is already becoming a niche corner of research [27,28]. At the time of writing this paper and to the best of our knowledge, the 2 most recent and powerful of these improvements are the methods known as self-consistency [29] and the tree of thoughts [30]. The former uses majority voting from multiple replications, and the latter takes an ensemble approach to the chain of thought reasoning and allows LLMs to consider multiple different reasoning paths and to perform self-evaluation on choices. Other methods naturally exist beyond what is contained in this study because of the unbounded human imagination, which makes the domain of prompt engineering quite an exciting frontier. Interested readers may find the website [31] to be a useful resource, with new relevant articles being added to its repository regularly.

https://ai.jmir.org/2025/1/e57319

XSL•FO

While prompt engineering in the context of LKD has not yet entered the literature, some work has emerged in the context of health care. Prompt engineering and generative artificial intelligence broadly are of particular interest in the medical domain as the generation of health information is still of unknown quality. A few researchers have emphasized the importance of medical professionals using LLMs skillfully and in a way that produces reliable information [32,33]. It has been shown that the reliability of GPT-4 (OpenAI) is inconsistent when answering medical questions, and the authors call for prompt engineering techniques to improve its performance [34]. Similarly, other authors have experimented with ChatGPT on calculation-based United States Medical Licensing Examination questions using 3 different prompting strategies, although they found that the prompt itself had only a small effect on answer accuracy [35]. Other research examined using prompt

engineering in generating health messages [36] and even medical image segmentation [37].

Social Media and LKD

Recent years have witnessed a burgeoning interest in studying dialogue on social media regarding important health care issues, such as vaccination [38] and LKD. Henderson [39] highlighted the use of platforms such as Facebook and Twitter to identify potential living donors while noting that formal research efforts are in their early stages. Analyzing social media content, including organ donation posts on the Chinese social media site Weibo, has unearthed key themes such as "organ donation behaviors," "statistical descriptions of organ donation," and "meaningfulness of donation" [40]. In one study, a notable 53% of potential living donors who self-referred for donor evaluation reported that they learned about a patient's need for a donor on social media [41,42], while specialized tools such as the "DONOR" app have enabled expansion of social media marketing about living donation between potential donors and patients with kidney diseases [43]. Research efforts include measuring organ donation awareness through Twitter digital markers [44], surveying readiness of patients who are undergoing a transplant to use social media for education [45], and using Twitter for living donor profile classification [46].

Interventions to increase living donation have used mobile health technologies to manage donor follow-up [47], delivered targeted advertising to specific ethnic groups [48,49], and assessed organ donation awareness across the United States using Twitter data [50]. Best practices for promoting LKD through social media, such as delivering content to specific community demographics in targeted and interactive modes, have been proposed [51]; live transplant broadcasts on Twitter have occurred [52]; and the analysis of public Facebook pages of potential living donors [53] has enhanced insights into donor identification and donation interest. Recent studies highlighted the importance of tailored messaging over generic communication for better audience engagement [54,55].

These investigations underscore social media's potential in augmenting donation awareness and facilitation, emphasizing the necessity for robust methods to discern and support individuals disseminating LKD-related content. A recent study by Garcia Valencia et al [56] has shown that ChatGPT can simplify medical information, making it easier to read and understand by many diverse groups. This can be a vital aid for promoting fairness in access to donation information from official sources. However, with the availability of *public* dialogue in forums also comes the need to thematically understand it. There is variation in both the content being shared and the user sharing it. The growing body of research demonstrates the potential of social media to impact awareness, intention to donate, and the facilitation of living kidney transplants. Therefore, it is necessary to have reliable methods whereby people who explicitly create and share content related to LKD can be automatically identified and understood for appropriate education and support. With this background, our research seeks to assess whether a classification system can be devised to discern individuals at varying stages of decision-making about becoming a living kidney donor. It also explores which of the contemporary NLP models are most apt for automating this classification, namely a fine-tuned distilled version of the BERT (DistilBERT) model (hereafter referred to as BERT for simplicity, unless greater specificity is merited) or ChatGPT. Furthermore, regarding ChatGPT, it examines how prompt engineering-namely, making adjustments to model instructions about the reasoning approach, examples, temperature, and class descriptions- influences its predictive efficacy for this application.

By answering these research questions, this study aimed to build a foundation for a sophisticated classification system in which it is possible to automatically categorize large amounts of social media communication about living donations using these tools. The study also aspires to gain a more in-depth insight into how individuals communicate and express themselves regarding LKD on various social media platforms. Using cutting-edge NLP technologies, our goal is to develop a streamlined, automated process for pinpointing curious, motivated potential donors who have not yet presented to the transplant center so that educational interventions could later be directed to them.

Methods

Data Labeling, Preparation, and Quality Assurance

We used a dataset of 2689 Reddit posts related to LKD from our previous work [14], which were published between January 2010 and April 2021. We also collected 603 Reddit posts from April 2021 to April 2023, for a combined total of 3292 posts from 2591 users. We scraped the posts with the open-source tool pushshift.io using keywords related to LKD, such as "kidney donor," "kidney transplant," "kidney donated," "kidney donate," "kidney years ago," "kidney need," "kidney stranger," and "kidney willing donate." Other search terms could have been included; however, as presented in Table 2, a considerable portion of collected data were not related to personal experiences, and we concluded that additional search terms would primarily expand the noise and add little value.



 Table 2. Distribution and description of Reddit (Reddit, Inc) classes.

Nielsen et al

Merged class categories and class categories	Description	Example post
Present (n=540, 26.9%)	-	·
Present direct (n=363, 21.5%)	The user has <i>current firsthand experience</i> with something personally related to kidney disease, kidney failure, living kidney donation, or transplan- tation (eg, the user with kidney disease or kidney failure, is on dialysis, is seeking a kidney, is explor- ing donation, or is undergoing evaluation for dona- tion or transplantation).	"A friend of mine is in need of a kidney. My first instinct is to offer one of mine. I have Googled and read LOTS of info. What would you do? Have you donated a kidney? What am I missing?"
Present indirect (n=177, 5.4%)	The user has <i>current secondhand experience</i> related to living kidney transplantation (eg, they <i>know</i> <i>someone</i> who is currently experiencing kidney failure, on dialysis, seeking a kidney, or preparing to donate a kidney).	"I need help finding a kidney for my dad."
Past (n=222, 6.8%)		
Past direct (n=168, 5.1%)	The user has <i>past firsthand experience</i> related to living kidney transplantation (eg, kidney failure, dialysis, kidney recipient or donor).	"Eight years ago today, I donated a kidney to a friend. Ask me any- thing."
Past indirect (n=58, 1.8%)	The user has <i>past secondhand experience</i> related to living kidney transplantation (eg, they <i>know</i> <i>someone</i> who experienced kidney failure, was on dialysis, received a kidney, donated a kidney, un- derwent evaluation for donation, or participated in the donation process (perhaps in a supporting role).	"Picture of my dad and the woman who donated a kidney to save his life."
Other (n=2530, 76.8%)		
General commentary or hypothetical (n=159, 4.8%)	The user is giving a <i>general opinion</i> on the topic, asking a <i>hypothetical question</i> , or contributing to discussion about an <i>imagined scenario</i> .	"If you donate a kidney, then later your only one starts to fail, would you be put on a higher priority?"
News or noise (n=2371, 72%)	The user is either sharing a <i>news article or headline</i> related to kidney donation that may be pertinent but <i>not personal</i> , or it is <i>simply irrelevant</i> .	"A man donated his kidney to his wife of 51 years after finding out he's her perfect match."

We selected Reddit as our data source because it provided the greatest portion of comments that were related to personal experiences rather than discussions of policies and sharing news stories. Reddit was the only place where we found posts from actual living donors inviting people to an "ask me anything" session, sparking highly personal discussions [14].

Under the guidance of LKD domain experts, after reviewing 100 example posts, we created 2 class sets, one with 6 classes (class categories) and the other with 3 classes (merged categories), to automate the process of identifying firsthand experiences with living donation (Table 2). These classes were iteratively defined and improved through multiple discussions with a team of 6 people who performed the manual annotation. Certain posts had sufficient ambiguity to make an explicit ruling impossible. For example, it was not always clear what constituted the boundary between a past and present experience (eg, how much time should have passed since the transplant?) or whether the general transplant mentioned in a post came from a living or deceased donor. Furthermore, long and verbose posts with brief mentions of personal experiences with donation posed a challenge because the brief (although important) mentions of LKD were easy to miss. Individual annotators were found to exhibit varying classification tendencies or use their own "rules of thumb" to expedite the often tedious process.

The granularity between these 6 fine-grained classes proved quite difficult for the models to correctly capture during initial experiments (resulting in accuracies <50%), so the posts were consolidated into the 3 coarse-grained categories: present (n=540, 42.59% of posts), past (n=222, 17.51% of posts), and other (n=506, 39.91% of posts randomly sampled from news or noise and general commentary or hypothetical categories) for 1268 samples that were used for training the BERT model. A randomly selected subset of 100 from each of the 3 classes was used for prompting with ChatGPT. The decision was made to aggregate general commentary and hypothetical posts with news or noise to ensure a more precise focus on personal experiences.

Acknowledging the potential data quality risks [57], we meticulously evaluated incorrect predictions from both BERT and ChatGPT after the analysis. The incorrectly predicted samples were tagged as either acceptable errors (reasonable, if not perfectly aligned predictions), unacceptable errors (flawed or evidently incorrect reasoning), more accurate than the original human label, or instances where both human and model erred. We later reported these using the notation of *LLM human*, *LLM*<*human*, *LLM*>*human*, and *both error*, respectively, for both models.

XSL•FO

Ethical Considerations

This study was granted an exemption from The University of Louisville Institutional Review Board (review number 22.0458). While there could be ethical concerns about consent and storage of health-related data, every Reddit user is entirely anonymous, ensuring that nothing we find can be directly traced to an individual. In addition, the comments and posts themselves are all very public; some websites may have minimal requirements, such as logging in or being a member of a "closed" group before the content can be observed; however, this is not the case for any of the data we collected. For data sources where such anonymity is not guaranteed, it is imperative to ensure that users consent to the study of their created content and that any identifying information be removed or obscured.

Modeling

We compared 2 transformer-based models for our classification task: a fine-tuned BERT model and a prompt-engineered ChatGPT model. We used the 3.5 Turbo version of ChatGPT via the OpenAI application programming interface and conducted a full factorial analysis of various prompt components to identify the best features. The DistilBERT model was fine-tuned from a pretrained Hugging Face (Hugging Face, Inc) model. Furthermore, we noted that many new models have emerged, both proprietary and open source, after our experiments were completed. Post hoc experiments indicate that our findings are consistent with newer models.

BERT Analysis

The DistilBert tokenizer from Hugging Face was used to tokenize the text data from Reddit, and both input IDs and attention masks were generated to structure the text inputs for the model. A custom model was designed around DistilBERT. The architecture included the pretrained DistilBERT model, followed by 3 fully connected layers with 768, 256, and 128 units, respectively. These were followed by an output layer with 3 units corresponding to the number of classes. Batch normalization and rectified linear unit activation functions were applied, and dropout was set at 10%.

The focal loss was used as the loss function, which is designed to address the class imbalance by downweighting the loss assigned to well-classified examples [58]. It was parameterized with an α factor for controlling the weight and a γ factor for focusing on hard examples. The model was trained using the AdamW optimizer [59], with the learning rate and weight decay optimized by the open-source Optuna hyperparameter tuning library. The dataset was split into training and validation sets using stratified 5-fold cross-validation, with class weights computed to manage class imbalance, and the model was trained for 3 epochs, following the recommended fine-tuning procedures [19]. The metrics used for validation are defined subsequently.

Accuracy is the ratio of correctly predicted instances to the total instances.

x

Precision is the ratio of correctly predicted positive observations to the total predicted positives.

I
J

Recall (sensitivity) is the ratio of correctly predicted positive observations to all observations in actual class.



 F_1 -score is the harmonic mean of precision and recall.

In equations 1 to 4, *TP*, *TN*, *FP*, and *FN* are the numbers of true positive, true negative, false positive, and false negative values, respectively.

×

The Optuna library was used to perform hyperparameter optimization, which uses a Bayesian optimization method known as the Tree-structured Parzen Estimator [60]. A search space was defined for the learning rate (ranging from 0.00003 to 0.0003) and weight decay (ranging from 0.0001-0.001). A total of 100 trials were conducted to find the best set of hyperparameters based on the F_1 -score.

Dialogue Until Classification Consensus

We introduced a text classification tool for LLMs termed "dialogue until classification consensus" (DUCC). Given the absence of a formal taxonomy for prompt engineering methods, we aligned DUCC's presentation with the pattern widely adopted in software development, which includes a name and classification, intent and context, motivation, structure and key ideas, example implementation, and consequences (Textbox 1). White et al [22] constructed the following categories of prompting patterns: input semantics, output customization, error identification, prompt improvement, interaction, and context control.



Textbox 1. Prompting patterns for "dialogue until classification consensus" (DUCC).

Name and classification

DUCC primarily falls under output customization, although it shares elements from other pattern categories, notably error identification and interaction.

Intent and context

DUCC assigns a persona of at least 2 domain experts to the large language model, instructing them to discuss a text sample until a consensus on its classification or answer selection is reached from a set of options. This setup aims to automate explicit reasoning and reflection through a simulated dialogue, expecting to resemble the effects of distribution-oriented methods, such as self-consistency, without requiring multiple sample replications.

Motivation

Complex classification tasks, especially within niche domains, such as personal living kidney donation experiences, often present labeling challenges. DUCC simulates expert discussions for decision-making while aiming to standardize output formats for classification tasks.

Structure and key ideas

Experts 1 and 2, specialized in [DOMAIN], are to discuss the text sample until an agreed classification or answer is reached.

The final label should be clear with no disagreements, formatted as: "classification: Label."

Additional identities or traits can be attributed to the experts to infuse specific perspectives into the discussion. We have observed that unless a singular label selection is emphasized, the model might assign multiple labels in challenging scenarios.

Example implementation

"Expert 1 and Expert 2, you are both experts in living kidney donation, and you've been tasked with analyzing and classifying a Reddit post that should be related to living kidney donation. You should discuss the post until you come to an agreement for a single classification. If the post is not related to living kidney donation, it needs to be labeled 'Other'. The classifications are defined as follows:

- Present: The user is describing a current or ongoing personal experience with living kidney donation
- Past: The user is describing a past personal experience with living kidney donation.
- Other: The user isn't discussing a personal experience with living kidney donation or isn't discussing living kidney donation at all.

Discuss until you reach a consensus, showing your reasoning. The final label should be clear, and there should be no disagreement. Output your agreed label in this format: { 'classification': 'your agreed label'}.

Here's an example of how this should be done:

- Post: 'Are you a kidney donor? How was the recovery process and how are you doing now?'
- Expert 1: 'I think the appropriate label is Present, because the user is asking questions and seems to want information to help them with a current decision about living kidney donation.'
- Expert 2: 'I think the appropriate label is Past because the user wants to know about past personal experiences from others.'
- Expert 1: 'I see your point about bringing up the past, but since we are interested in assigning a label to the user who wrote the post, we should keep our focus on the author's perspective. If we knew what the replies were, we could label those users as Past, but we are only looking at this user for now.'
- Expert 2: 'You're correct, we should be focused on this user rather than possible answers from others. Even though there are elements of both, we have to pick one and only one label, so let's go with Present.'
- Final Label: "classification': 'Present."

Consequences

DUCC prompts large language models to reason through multiple perspectives, ensuring a singular, consistently formatted label, simplifying extraction. The example implementation is crucial as it demonstrates the desired dialogue structure, aiding the model in handling nuanced classifications. However, DUCC may exhibit biases when numerous classes are present, potentially leaning toward the exemplified label. To mitigate token use, especially in lengthy examples, using DUCC when defining the system instead of individual prompts is advisable. For instance, in the OpenAI application programming interface, modifying the "content" section of the "system" role with the entire provided example instead of the default content can better define the system's nature.

Sensitivity Analysis of Prompting

Overview

For our experimentation using ChatGPT to categorize personal experiences, we conducted a study applying a full factorial design with 4 factors (summarized subsequently), which resulted in 48 experimental runs. We must first acknowledge that the nature of prompting is such that there were an infinite number

of ways we could write the prompt and parameters that could be chosen. It is well known that examples that illustrate the solutions can influence performance (known as "few-shot" prompting) [61], so we examined the number of examples and the type of examples that might produce bias as well as the parameters provided subsequently.

Use of the DUCC Method (2 Settings)

In addition to the DUCC method described earlier, the alternative was to prompt a single expert to make a classification decision, with the instruction to "Examine the evidence for each class option step by step. The final label should be clear." In this case, the model attempts to identify any evidence that suggests the sample should be assigned to each class and weighs the evidence to draw a conclusion.

Number of Examples Used (4 Settings)

We selected either 1 example or 3 examples. For 3 examples, 1 example was used for each class (present, past, and other). For the single example setting, we performed an experiment with each class once to evaluate whether it produced a bias in the predicted class.

Definition of "Past" (2 Settings)

Observing a tendency for underprediction in the "Past" label, we considered 2 definitions for the class. The first was a short and concise definition: "The user is describing a past personal experience with living kidney donation." The second was a longer, more descriptive definition: "The user is referring to a past personal experience with LKD. This may be presented in the context of a present tense story, but if the event of LKD was lived previously, the post should be labeled past."

Temperature Settings (3 Settings)

Experimentation spanned temperature values of 0, 0.15, and 0.3, investigating the tradeoff between output variability and consistency. The settings were guided by OpenAI documentation, emphasizing lower values for consistency and higher values for diversifying outputs [62].

Given the cost implications of OpenAI application programming interface calls, an initial assessment was carried out to determine the necessity for replicating each setting. We performed 30 replications of a fixed parameter setting and found no substantial effect within replications for any metric. Thus, the experimentation proceeded with a singular sample for each parameter setting.

Results

Overview

In this section, we present the results of the BERT model first and then the results of ChatGPT. We present the performance metrics, confusion matrices, and assessment of incorrect predictions. For ChatGPT, we also present the results of an ANOVA on the various factors used in the experimentation.

BERT Results

In >100 trials, the best BERT model performed with an accuracy of 75.1% and an F_1 -score of 78.2% on the validation data during training. The best parameters were a learning rate of 0.000131687 and a weight decay of 0.000791. The confusion matrix for the predictions on the test data is presented in Figure 1, showing reasonably good performance but with a tendency to erroneously predict the Other label on both past and present labels.

The classification report provided in Table 3 shows that the BERT model significantly underpredicts past labels, partly due to the smaller sample size, and also because of the ambiguity that can arise when a reference to a past experience is nested within an ongoing story.

Figure 1. Confusion matrix for the best Bidirectional Encoder Representations from Transformers model.



Table 3. Classification report.

	Precision	Recall	F ₁ -score	Support
Present	0.88	0.82	0.85	101
Past	0.66	0.52	0.58	44
Other	0.75	0.86	0.80	108
Weighted average	0.79	0.79	0.78	253

ChatGPT Results

The best ChatGPT prompt produced an accuracy and F_1 -score of 78.67% and 78.17%, respectively (surprisingly, this F_1 -score is identical to that of BERT). This was achieved using the DUCC method, a single example of a present class post, a temperature of 0, and the shorter definition of the past class (refer to the Dialogue Until Classification Consensus section). Full experimentation results are provided in the Multimedia

Figure 2. Confusion matrix for the best ChatGPT prompt.

Appendix 1. The next 3 columns show the percentage of predictions for that class, and the remaining 3 columns show the evaluation metrics.

The confusion matrix for ChatGPT performance is presented in Figure 2, which shows again that past class samples were underpredicted and that both other and past class samples were overpredicted to be present class, suggesting a bias toward present classifications.



The results of the ANOVA are presented in Table 4, which shows that the number and type of examples used is the most significant factor, followed by the method. We observe that the examples and method factors were the only statistically significant factors.

Given that there were 3 df within the examples setting, we sought to better understand the difference between the example settings using a Tukey test, with results provided in Table 5. We observed that when our example belonged to the "past" class the model performed better than when the example came

from the "other" class. But using an example from the "past" class resulted in poorer performance compared to using 3 examples (one from each class) and using an example from the "present" class. Interestingly, the "past" sample was underpredicted in every setting except when using 3 examples and the evidence method. Interestingly, samples belonging to the "past" class were underpredicted in every setting except when using 3 examples and the evidence method. Although this setting (3 examples; evidence method) does not demonstrate the same underprediction bias as other settings, it does not give better accuracy overall.



Table 4. ANOVA results.

Factor	Sum of squares	F test (df)	<i>P</i> value
Category (examples)	0.068615	27.659884 (3, 40)	<.001
Category (method)	0.006466	7.819650 (1, 40)	.008
Category (temp)	0.000024	0.014557 (2, 40)	.99
Category (past)	0.000032	0.039292 (1, 40)	.84
Residual	0.033076	a	_

^aNot applicable.

Table 5. Multiple comparisons of means using the Tukey honestly significant difference test. The family-wise error rate is 0.05.

Group 1	Group 2	Mean difference	P value	Lower limit	Upper limit	Reject
1, other	1, past	-0.0875	<.001	-0.1202	-0.0548	True
1, other	1, present	0.0078	.92	-0.0249	0.0405	False
1, other	3	-0.0017	.99	-0.0344	0.031	False
1, past	1, present	0.0953	<.001	0.0626	0.128	True
1, past	3	0.0858	<.001	0.0531	0.1185	True
1, present	3	-0.0094	.87	-0.0421	0.0233	False

Discussion

Principal Findings

Our experimentation has found that BERT and ChatGPT perform comparably for the classification of different living kidney donor experiences. Because BERT is completely dependent on the available training data, ChatGPT can be used with a somewhat higher degree of precision via prompt engineering, as shown by our use of the novel DUCC method. Our full factorial experimentation identified the best settings to use for our engineered prompt. In this section, we will discuss the predictions that were made incorrectly and consider future work and ethical considerations.

Examination of Incorrect Predictions

As noted in the Data Labeling, Preparation, and Quality Assurance section, there is an inherent risk of data quality that arises from the dataset in question. Unlike standardized benchmarks, which often have explicit "ground truth" labels, our task is fraught with nuance. Despite our extensive efforts to ensure data quality, the given label is not always clear. As such, we have provided a more detailed examination of the instances where the models made predictions that diverged from the given labels.

BERT and GPT-3.5 produced 21.3% (54/253) and 21.3% (64/300) incorrect predictions, respectively. It should be recalled that the difference in the denominator values is because BERT requires a split test set, whereas, with GPT-3.5, we can use a larger inference-only set. We assessed the quality of these incorrect predictions not only to see how "close" they were to the mark but also to determine whether any human errors had been made in labeling the incorrect predictions. As provided in Table 6 for BERT, we observe that 27 prompts were incorrectly labeled either because of an acceptable error where a clear prediction is difficult to make (perhaps due to the ambiguity of what constitutes the difference between the past and present samples) or where BERT made a better prediction than the original human label. Treating these 27 predictions as being acceptable or correct brings the total number of correct predictions from 199 (78.7%) of 253 to 226 (89.3%) of 253, which elevates the predictive accuracy considerably to 89.3%. In these tables, examples are written "as they are" from the original posts, including typos and terminology that may be unique to Reddit.



Nielsen et al

Table 6. Analysis of incorrect predictions from Bidirectional Encoder Representations from Transformers (BERT; n=54).

Error type	Incorrect predictions, n (%)	Example post	Reason
Unacceptable error (BERT <human)< td=""><td>22 (41)</td><td>"Required testing to be a living Kidney donor where I live - these are the tests I took before becoming a living kidney donor almost 2 yrs ago everything has gone great for me and the recipient happy to answer any questions."</td><td>BERT predicted the "other" label, but the user clearly states that he or she was a previous living donor.</td></human)<>	22 (41)	"Required testing to be a living Kidney donor where I live - these are the tests I took before becoming a living kidney donor almost 2 yrs ago everything has gone great for me and the recipient happy to answer any questions."	BERT predicted the "other" label, but the user clearly states that he or she was a previous living donor.
Acceptable error (BERT human)	12 (22)	"Hey Mum, it's been a year since what was supposed to be a life changing kidney transplant that took a turn for the worst. I love you so much and think about you every day xxx"	BERT predicted the "other" label, which could be appropriate if it was a deceased donor transplant. We predict- ed the "past" label.
Human error (BERT>human)	15 (27)	"Me 26F with my Dad 58	We predicted the "other" label because
		he needs a kidney and I feel pressured to donate one. [removed]"	of the (removed) tag at the end of the post, which commonly appears in unus- able posts. BERT predicted the "present" label, which is the more ap- propriate label.
Both erred	5 (9)	"I used to like her but I found out that she did not even acknowledge her kid- ney donor Just referring to her as a person I know it seems pretty ungrate- ful [removed]"	This is someone's opinion about a celebrity who famously received a kidney transplant from her friend. It is not a personal experience at all, but the human label was "present," and the BERT label was "past."

From our analysis of the incorrect predictions on GPT-3.5 (Table 7), we observed that 26 (40%) of the 64 errors were acceptable.

As mentioned earlier, we had previously observed that many "past" posts were labeled as "present" because many of the posts were in a present tense context. The best setting used the shorter definition of past, which does not teach the model to treat past experiences nested in present accounts as the past class, so this is to be expected. Anytime both the human and predicted labels were wrong, the post was almost always ambiguous regarding whether it was about living or deceased donation. The experiences being described could have been a living donation, but there is not enough information to determine that for certain. Regarding BERT, we may allow ourselves to consider the 26 acceptable errors and 10 human errors as being correctly predicted, changing the total number of correct predictions from 236 (78.7%) of 300 to 272 (90.7%) of 300 for an "actual" predictive accuracy of 90.7%. While still imperfect, this shows considerable reliability when using these methods on nuanced language tasks.

The implications of this examination are threefold: (1) sometimes human annotations go wrong, even with clear instructions; (2) these powerful models are capable of correctly catching things that humans miss (due to decision fatigue or similar cognitive difficulties); and (3) the models can be largely trusted to give sensible reasoning, even if the final conclusions differ from that of a human counterpart.



Table 7. Analysis of incorrect predictions from ChatGPT (n=64).

Error type	Incorrect predictions, n (%)	Example post	Reason
Unacceptable error (ChatG- PT <human)< td=""><td>21 (33)</td><td>"relationships My (36F) estranged sister (43F) donated a kidney to me. I just heard that she died (for a different reason). I'm very confused. [removed]"</td><td>The simulated experts reasoned that the focus of the post was on grief rather than LKD^a and labeled it as "other." The human label was given as "past" because the user mentions a sister who donated her kidney some time ago.</td></human)<>	21 (33)	"relationships My (36F) estranged sister (43F) donated a kidney to me. I just heard that she died (for a different reason). I'm very confused. [removed]"	The simulated experts reasoned that the focus of the post was on grief rather than LKD ^a and labeled it as "other." The human label was given as "past" because the user mentions a sister who donated her kidney some time ago.
Acceptable error (ChatGPT human)	26 (41)	"Successfully donated a kidney to my sis- ter whos been fighting Lupus."	This could be easily interpreted as either a "present" (ChatGPT) or a "past" (hu- man) label, given that there is no explicit reference to time. It could go either way, but it is still clearly related to a personal experience with LKD.
Human error (ChatGPT > human)	10 (16)	"I (30F) had heart and kidney transplant. Ask Me Anything (AMA)."	The simulated experts concluded that this should be labeled "other" when the human label had been given as "past." ChatGPT made a more correct conclusion because this may have been from a deceased donor rather than a living donor. We would need more information to be certain, so it should be an "other" label.
Both erred	7 (11)	"I am A double kidney transplant recipi- ent! AMA! I am a 28 year old white male, I've had two renal transplants over the course of my lifetime. I've been on dialy- sis. I've been in and out of hospital my entire life. I think it's interesting, but there's only one way to find out! Ask Me Anything."	The human-given label for this was "past" because of the previous transplant experi- ences, and the reasoning provided by ChatGPT concluded that the label should be "present" because the user mentions dialysis and being in and out of the hospi- tal. Both were incorrect because there is not enough evidence that either of the transplants was from living donors, and thus, it should be labeled "other."

^aLKD: living kidney donation.

Limitations and Future Work

BERT and ChatGPT have both proven effective in classifying personal accounts of LKD on platforms such as Reddit, achieving approximately 80% accuracy, which increases to about 90% when considering acceptable errors, marking a step forward in using web-based data for LKD research. These models could potentially automate the screening of new content for further scrutiny, thereby aiding donor support initiatives, particularly in education and community outreach. Despite the promising results, the complexity of the subject matter complicates the task of making perfect predictions. Our initial attempts to use fine-grained classifications led to suboptimal results, requiring us to use coarse-grained categories. Regarding costs, BERT's open-source nature and the flexibility to fine-tune make it an appealing choice. In contrast, ChatGPT excels in providing understandable reasoning for its decisions.

A review of errors indicated that ChatGPT generally understood the context well, although there were instances where the reasoning was off the mark, highlighting the importance of clear, prompt instructions. Interestingly, there were instances where the LLMs' reasoning surpassed ours, especially in delineating the "past" and "present" boundary, thereby suggesting a potential for iterative prompt enhancements informed by LLM reasoning. However, the quest for prompt optimization (or "promptization," if you will) may present an

RenderX

unending journey, as the allure of "just one more experiment" to elevate performance is always present. Drawing a line on performance as "good enough" is crucial, which may be attained through automated processes, as explored in some recent and exciting studies [63-69]. Future work will leverage these powerful new methodologies to both improve performance on our coarse-grained 3-class schema as well as achieve superior performance on the fine-grained 6-class schema that was unattainable with the present methods.

The performance of both models is significantly constrained by the size of the available data. While thousands of Reddit posts related to LKD are accessible, only a fraction pertains to personal experiences. The performance consistency across different data folds for BERT and across different sample sizes for ChatGPT highlights the need for larger datasets to better gauge each model's robustness.

A core challenge lies in the task's inherent demand for a singular label, which often oversimplifies the nuanced narratives in internet posts. Future endeavors could explore more elaborate information extraction techniques, leveraging LLMs such as ChatGPT to answer multiple queries or even construct knowledge graphs per post. Although ensuring uniform and usable output formats remains a hurdle, our work underscores ChatGPT's proficiency in deriving insightful inferences from the text. Our findings concerning the influence of few-shot learning examples on output bias also suggest the need for

deeper investigation into the interplay between example selection and model performance.

With reliable automation methods that can identify when a person is describing a personal experience with LKD, future work will extend the reach to additional media platforms, each of which has its own system for reaching users via advertising. There will certainly be potential biases in accessing educational information about living donations based on the characteristics of audiences most likely to post on each platform. To not exacerbate disparities, one must examine the generalizability of the profiles across multiple platforms and ensure the dissemination of information across platforms that reach diverse audiences and non-English speakers. An examination of access to most audience members, particularly the underserved, is warranted to ensure that all communities are reached equitably.

Utility of Results

By identifying these unique user classifications, tailored educational interventions for different profiles could be designed. First, for those most actively considering living donation, there could be social media campaigns built and targeted to specific users to invite them to learn more about living donation. These users can be referred to a trusted site, which includes education materials and an opportunity to register to begin donor medical evaluation at a nearby transplant center [41,42]. For individuals discussing their concerns about the costs involved with becoming a living donor, referrals to websites that discuss the ways to apply for grants to cover the out-of-pocket costs and lost wages could be valuable in their decision-making [70].

Second, for donors and families identified to have completed donations, campaigns inviting them to share their experiences on a living donor storytelling website [8,9] might result in more real-life stories being captured from diverse individuals to increase awareness of living donations for the national public. Stories are particularly valuable for educating learners with low health literacy or those for whom English is not their primary language about the possibilities of living donation [71].

Finally, it will be very important to work with experts in marketing and campaign design to plan social media campaigns that are motivating and helpful for patients and their families at different points along their donation journey. Identifying motivated learners from platforms such as Reddit, delivering content to them about living donation, and assessing its impact on learning more or pursuing donation are our next planned steps.

The proposed profiles may incorrectly identify a person's interest or stage of pursuit of donation, making any educational information sent to them irrelevant. In contrast, users could also be made uncomfortable if the education being provided matches their needs perfectly, indicating that their data are being scrutinized. Users can always disregard nonrelevant content; however, it will be important in the design of new campaigns not to assume with too much certainty that all learners are correctly identified. Respect for users is an ethical tenet that must always be considered in designing the campaigns and communicating how we found that they might be considering living donations as we move forward.

Conclusions

Much of the previous health care–related research about LLMs has been centered on their reliability in producing quality medical information. In contrast, we endeavor to extract individual-level information from the internet that can be used to inform health care providers. Consequently, there is little comparison that can be made to previous work other than to say that the reliability of the models is subject to the instructions they are given. However, our experimental results do illustrate that when using examples as part of the prompt (few-shot), bias toward the class of the given examples can affect performance. We have also shown that simulating a dialogue between 2 experts is more effective than using stand-alone reasoning.

This study takes a significant step in applying advanced NLP methods to the field of LKD, focusing on automating the detection of personal LKD experiences in online content. Both BERT and ChatGPT proved effective for this task, each with its own advantages and disadvantages. Our new DUCC method outperformed traditional reasoning approaches, emphasizing the importance of further work on improving prompt design. The study also highlights the need for automated prompt creation to reduce the time and effort currently required for manual testing, making NLP applications in the LKD field more efficient and impactful.

Acknowledgments

This study is supported in part by the Logistics and Distribution Institute at the University of Louisville. XC is supported by the American Heart Association (23CSA1052735), and National Science Foundation (CMMI-2430998).

Conflicts of Interest

None declared.

Multimedia Appendix 1 Full experimental results. [XLSX File (Microsoft Excel File), 13 KB - ai_v4i1e57319_app1.xlsx]

References



- Abecassis M, Bartlett ST, Collins AJ, Davis CL, Delmonico FL, Friedewald JJ, et al. Kidney transplantation as primary therapy for end-stage renal disease: a National Kidney Foundation/Kidney Disease Outcomes Quality Initiative (NKF/KDOQITM) conference. Clin J Am Soc Nephrol 2008 Mar;3(2):471-480 [FREE Full text] [doi: 10.2215/CJN.05021107] [Medline: 18256371]
- Axelrod DA, Schnitzler MA, Xiao H, Irish W, Tuttle-Newhall E, Chang S, et al. An economic assessment of contemporary kidney transplant practice. Am J Transplant 2018 May;18(5):1168-1176 [FREE Full text] [doi: 10.1111/ajt.14702] [Medline: 29451350]
- 3. All-time records again set in 2021 for organ transplants, organ donation from deceased donors. Health Resources and Services Administration. URL: <u>https://optn.transplant.hrsa.gov/news/</u> all-time-records-again-set-in-2021-for-organ-transplants-organ-donation-from-deceased-donors/ [accessed 2023-01-25]
- Lentine KL, Smith JM, Hart A, Miller J, Skeans MA, Larkin L, et al. OPTN/SRTR 2020 annual data report: kidney. Am J Transplant 2022 Mar;22 Suppl 2:21-136 [FREE Full text] [doi: 10.1111/ajt.16982] [Medline: 35266618]
- Purnell TS, Hall YN, Boulware LE. Understanding and overcoming barriers to living kidney donation among racial and ethnic minorities in the United States. Adv Chronic Kidney Dis 2012 Jul;19(4):244-251 [FREE Full text] [doi: 10.1053/j.ackd.2012.01.008] [Medline: 22732044]
- Purnell TS, Luo X, Cooper LA, Massie AB, Kucirka LM, Henderson ML, et al. Association of race and ethnicity with live donor kidney transplantation in the United States from 1995 to 2014. JAMA 2018 Jan 02;319(1):49-61 [FREE Full text] [doi: 10.1001/jama.2017.19152] [Medline: 29297077]
- Morgan SE, Harrison TR, Long SD, Afifi WA, Stephenson MS, Reichert T. Family discussions about organ donation: how the media influences opinions about donation decisions. Clin Transplant 2005 Oct 11;19(5):674-682. [doi: 10.1111/j.1399-0012.2005.00407.x] [Medline: 16146561]
- Ho EW, Murillo AL, Davis LA, Iraheta YA, Advani SM, Feinsinger A, et al. Findings of living donation experiences shared on a digital storytelling platform: a thematic analysis. PEC Innov 2022 Dec;1:100023 [FREE Full text] [doi: 10.1016/j.pecinn.2022.100023] [Medline: 37213721]
- 9. Davis L, Iraheta YA, Ho EW, Murillo AL, Feinsinger A, Waterman AD. Living kidney donation stories and advice shared through a digital storytelling library: a qualitative thematic analysis. Kidney Med 2022 Jul;4(7):100486 [FREE Full text] [doi: 10.1016/j.xkme.2022.100486] [Medline: 35755303]
- Kaplow K, Ruck JM, Levan ML, Thomas AG, Stewart D, Massie AB, et al. National attitudes towards living kidney donation in the United States: results of a public opinion survey. Kidney Med 2024 Mar;6(3):100788 [FREE Full text] [doi: 10.1016/j.xkme.2023.100788] [Medline: 38435064]
- 11. Amaral S, McCulloch CE, Black E, Winnicki E, Lee B, Roll GR, et al. Trends in living donation by race and ethnicity among children with end-stage renal disease in the United States, 1995-2015. Transplant Direct 2020 Jul;6(7):e570 [FREE Full text] [doi: 10.1097/TXD.00000000001008] [Medline: 32766425]
- Waterman AD, Morgievich M, Cohen DJ, Butt Z, Chakkera HA, Lindower C, American Society of Transplantation. Living donor kidney transplantation: improving education outside of transplant centers about live donor transplantation--recommendations from a consensus conference. Clin J Am Soc Nephrol 2015 Sep 04;10(9):1659-1669 [FREE Full text] [doi: 10.2215/CJN.00950115] [Medline: 26116651]
- Waterman AD, Peipert JD. An explore transplant group randomized controlled education trial to increase dialysis patients' decision-making and pursuit of transplantation. Prog Transplant 2018 Jun 26;28(2):174-183. [doi: 10.1177/1526924818765815] [Medline: 29699451]
- Asghari M, Nielsen J, Gentili M, Koizumi N, Elmaghraby A. Classifying comments on social media related to living kidney donation: machine learning training and validation study. JMIR Med Inform 2022 Nov 08;10(11):e37884 [FREE Full text] [doi: 10.2196/37884] [Medline: 36346661]
- Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, et al. Attention is all you need. In: Proceedings of the 31st International Conference on Neural Information Processing Systems. 2017 Presented at: NIPS'17; December 4-9, 2017; Long Beach, CA p. 6000-6010 URL: <u>https://dl.acm.org/doi/10.5555/3295222.3295349</u>
- 16. Yang Z, Dai Z, Yang Y, Carbonell J, Salakhutdinov R, Le QV. XLNet: generalized autoregressive pretraining for language understanding. arXiv Preprint posted online June 19, 2019 [FREE Full text]
- 17. Liu Y, Ott M, Goyal N, Du J, Joshi M, Chen D, et al. RoBERTa: a robustly optimized BERT pretraining approach. arXiv Preprint posted online July 26, 2019 [FREE Full text]
- 18. Chernyavskiy A, Ilvovsky D, Nakov P. Transformers: "the end of history" for NLP? arXiv Preprint posted online April 9, 2021 [FREE Full text]
- 19. Devlin J, Chang MW, Lee K, Toutanova K. Bert: pre-training of deep bidirectional transformers for language understanding. arXiv Preprint posted online October 11, 2018 [FREE Full text]
- 20. Acheampong FA, Nunoo-Mensah H, Chen W. Transformer models for text-based emotion detection: a review of BERT-based approaches. Artif Intell Rev 2021 Feb 08;54(8):5789-5829. [doi: 10.1007/S10462-021-09958-2]
- 21. Ray PP. ChatGPT: a comprehensive review on background, applications, key challenges, bias, ethics, limitations and future scope. Internet Things Cyber Phys Syst 2023;3:121-154. [doi: <u>10.1016/j.iotcps.2023.04.003</u>]

- 22. White J, Fu Q, Hays S, Sandborn M, Olea C, Gilbert H, et al. A prompt pattern catalog to enhance prompt engineering with ChatGPT. arXiv Preprint posted online February 21, 2023 [FREE Full text]
- 23. Reynolds L, McDonell K. Prompt programming for large language models: beyond the few-shot paradigm. arXiv Preprint posted online February 15, 2021 [FREE Full text] [doi: 10.1145/3411763.3451760]
- 24. Liu Y, Deng G, Xu Z, Li Y, Zheng Y, Zhang Y, et al. Jailbreaking ChatGPT via prompt engineering: an empirical study. arXiv Preprint posted online May 23, 2023 [FREE Full text]
- 25. Shi Z, Wang Y, Yin F, Chen X, Chang KW, Hsieh CJ. Red teaming language model detectors with language models. arXiv Preprint posted online May 31, 2023 [FREE Full text] [doi: 10.1162/tacl_a_00639]
- 26. Casper S, Lin J, Kwon J, Cilp G, Hadfield-Menell D. Explore, establish, exploit: red teaming language models from scratch. arXiv Preprint posted online June 15, 2023 [FREE Full text]
- 27. Shinn N, Cassano F, Berman E, Gopinath A, Narasimhan K, Yao S. Reflexion: language agents with verbal reinforcement learning. arXiv Preprint posted online March 20, 2023 [FREE Full text]
- 28. Wei J, Wang X, Schuurmans D, Bosma M, Ichter B, Xia F, et al. Chain-of-thought prompting elicits reasoning in large language models. arXiv Preprint posted online January 28, 2022 [FREE Full text]
- 29. Wang X, Wei J, Schuurmans D, Le Q, Chi E, Narang S, et al. Self-consistency improves chain of thought reasoning in language models. arXiv Preprint posted online March 21, 2021 [FREE Full text]
- 30. Yao S, Yu D, Zhao J, Shafran I, Griffiths TL, Cao Y, et al. Tree of thoughts: deliberate problem solving with large language models. arXiv Preprint posted online May 17, 2023 [FREE Full text]
- 31. Papers. Prompt Engineering Guide. URL: https://www.promptingguide.ai/papers [accessed 2024-04-29]
- 32. Meskó B. Prompt engineering as an important emerging skill for medical professionals: tutorial. J Med Internet Res 2023 Oct 04;25:e50638 [FREE Full text] [doi: 10.2196/50638] [Medline: 37792434]
- 33. Wang J, Shi E, Yu S, Wu Z, Ma C, Dai H, et al. Prompt engineering for healthcare: methodologies and applications. arXiv Preprint posted online April 28, 2023 [FREE Full text]
- 34. Wang L, Chen X, Deng X, Wen H, You M, Liu W, et al. Prompt engineering in consistency and reliability with the evidence-based guideline for LLMs. NPJ Digit Med 2024 Feb 20;7(1):41 [FREE Full text] [doi: 10.1038/s41746-024-01029-4] [Medline: 38378899]
- 35. Patel D, Raut G, Zimlichman E, Cheetirala SN, Nadkarni G, Glicksberg BS, et al. The limits of prompt engineering in medical problem-solving: a comparative analysis with ChatGPT on calculation based USMLE medical questions. medRxiv Preprint posted online August 9, 2023 [FREE Full text] [doi: 10.1101/2023.08.06.23293710]
- 36. Lim S, Schmälzle R. Artificial intelligence for health message generation: an empirical study using a large language model (LLM) and prompt engineering. Front Commun 2023 May 26;8:1129082. [doi: 10.3389/fcomm.2023.1129082]
- 37. Ali H, Bulbul MF, Shah Z. Prompt engineering in medical image segmentation: an overview of the paradigm shift. In: Proceedings of the 2023 IEEE International Conference on Artificial Intelligence, Blockchain, and Internet of Things. 2023 Presented at: AIBThings '23; September 16-17, 2023; Mount Pleasant, MI p. 1-4 URL: <u>https://ieeexplore.ieee.org/document/10292475</u> [doi: 10.1109/aibthings58340.2023.10292475]
- Argyris YA, Monu K, Tan P, Aarts C, Jiang F, Wiseley KA. Using machine learning to compare provaccine and antivaccine discourse among the public on social media: algorithm development study. JMIR Public Health Surveill 2021 Jun 24;7(6):e23105 [FREE Full text] [doi: 10.2196/23105] [Medline: 34185004]
- 39. Henderson ML. Social media in the identification of living kidney donors: platforms, tools, and strategies. Curr Transpl Rep 2018 Jan 18;5(1):19-26. [doi: 10.1007/S40472-018-0179-8]
- Jiang X, Jiang W, Cai J, Su Q, Zhou Z, He L, et al. Characterizing media content and effects of organ donation on a social media platform: content analysis. J Med Internet Res 2019 Mar 12;21(3):e13058 [FREE Full text] [doi: 10.2196/13058] [Medline: 30860489]
- 41. DuBray BJ, Shawar SH, Rega SA, Smith KM, Centanni KM, Warmke K, et al. Impact of social media on self-referral patterns for living kidney donation. Kidney360 2020 Dec 31;1(12):1419-1425. [doi: <u>10.34067/kid.0003212020</u>]
- 42. Joachim E. Self-referral patterns of living kidney donors via social media: examining an expanding platform. Kidney360 2020 Dec 31;1(12):1337-1338 [FREE Full text] [doi: 10.34067/KID.0005732020] [Medline: 35372901]
- 43. Kumar K, King E, Muzaale A, Konel J, Bramstedt K, Massie A, et al. A smartphone app for increasing live organ donation. Am J Transplant 2016 Dec;16(12):3548-3553 [FREE Full text] [doi: 10.1111/ajt.13961] [Medline: 27402293]
- 44. Murphy MD, Pinheiro D, Iyengar R, Lim G, Menezes R, Cadeiras M. A data-driven social network intervention for improving organ donation awareness among minorities: analysis and optimization of a cross-sectional study. J Med Internet Res 2020 Jan 14;22(1):e14605 [FREE Full text] [doi: 10.2196/14605] [Medline: 31934867]
- 45. Kazley AS, Hamidi B, Balliet W, Baliga P. Social media use among living kidney donors and recipients: survey on current practice and potential. J Med Internet Res 2016 Dec 20;18(12):e328 [FREE Full text] [doi: 10.2196/jmir.6176] [Medline: 27998880]
- 46. Ruck JM, Henderson ML, Eno AK, Van Pilsum Rasmussen SE, DiBrito SR, Thomas AG, et al. Use of Twitter in communicating living solid organ donation information to the public: an exploratory study of living donors and transplant professionals. Clin Transplant 2019 Jan 07;33(1):e13447 [FREE Full text] [doi: 10.1111/ctr.13447] [Medline: 30421841]

- 47. Eno AK, Thomas AG, Ruck JM, Van Pilsum Rasmussen SE, Halpern SE, Waldram MM, et al. Assessing the attitudes and perceptions regarding the use of mobile health technologies for living kidney donor follow-up: survey study. JMIR Mhealth Uhealth 2018 Oct 09;6(10):e11192 [FREE Full text] [doi: 10.2196/11192] [Medline: 30305260]
- 48. Gordon EJ, Shand J, Black A. Google analytics of a pilot mass and social media campaign targeting Hispanics about living kidney donation. Internet Interv 2016 Nov;6:40-49 [FREE Full text] [doi: 10.1016/j.invent.2016.09.002] [Medline: 30135813]
- 49. Britt RK, Britt BC, Anderson J, Fahrenwald N, Harming S. "Sharing hope and healing": a culturally tailored social media campaign to promote living kidney donation and transplantation among native Americans. Health Promot Pract 2021 Nov 02;22(6):786-795. [doi: 10.1177/1524839920974580] [Medline: 33267677]
- 50. Pacheco DF, Pinheiro D, Cadeiras M, Menezes R. Characterizing organ donation awareness from social media. In: Proceedings of the 33rd International Conference on Data Engineering. 2017 Presented at: ICDE '17; April 19-22, 2017; San Diego, CA p. 1541-1548 URL: <u>https://ieeexplore.ieee.org/document/7930122</u> [doi: <u>10.1109/icde.2017.225</u>]
- 51. Basu G, Nair S, Sibel G, Dheerendra P, Penmatsa KR, Balasubramanian K, et al. Social media and organ donation a narrative review. Indian J Transplant 2021;15(2):139-146 [FREE Full text] [doi: 10.4103/ijot.ijot 138_20]
- 52. Tan M, Mulloy M, Pollinger H, Gibney E. Impact of social media on living kidney donation awareness. Transplantation 2014;98:836-837. [doi: 10.1097/00007890-201407151-02857]
- Chang A, Anderson EE, Turner HT, Shoham D, Hou SH, Grams M. Identifying potential kidney donors using social networking web sites. Clin Transplant 2013 Apr 22;27(3):E320-E326 [FREE Full text] [doi: 10.1111/ctr.12122] [Medline: 23600791]
- 54. Ayorinde JO, Saeb-Parsy K, Hossain A. Opportunities and challenges in using social media in organ donation. JAMA Surg 2020 Sep 01;155(9):797-798. [doi: 10.1001/jamasurg.2020.0791] [Medline: 32936283]
- Lee C, Lin M, Lin H, Ting Y, Wang H, Wang C, et al. Survey of factors associated with the willingness toward living kidney donation. J Formos Med Assoc 2022 Nov;121(11):2300-2307 [FREE Full text] [doi: 10.1016/j.jfma.2022.06.007] [Medline: 35803885]
- 56. Garcia Valencia OA, Thongprayoon C, Miao J, Suppadungsuk S, Krisanapan P, Craici IM, et al. Empowering inclusivity: improving readability of living kidney donation information with ChatGPT. Front Digit Health 2024 Apr 10;6:1366967 [FREE Full text] [doi: 10.3389/fdgth.2024.1366967] [Medline: 38659656]
- 57. Wu X, Zheng W, Xia X, Lo D. Data quality matters: a case study on data label correctness for security bug report prediction. IIEEE Trans Software Eng 2022 Jul 1;48(7):2541-2556. [doi: <u>10.1109/tse.2021.3063727</u>]
- 58. Lin TY, Goyel P, Girshick R, He K, Dollár P. Focal loss for dense object detection. arXiv Preprint posted online August 7, 2017 [FREE Full text] [doi: 10.1109/iccv.2017.324]
- 59. Loshchilov I, Hutter F. Decoupled weight decay regularization. arXiv Preprint posted online November 14, 2017 [FREE Full text] [doi: 10.1090/mbk/121/79]
- 60. Akiba T, Sano S, Yanase T, Ohta T, Koyama M. Optuna: a next-generation hyperparameter optimization framework. In: Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2019 Presented at: KDD '19; August 4-8, 2019; Anchorage, AK p. 2623-2631 URL: <u>https://dl.acm.org/doi/10.1145/3292500.3330701</u> [doi: 10.1145/3292500.3330701]
- 61. Brown TB, Mann B, Ryder N, Subbiah M, Kaplan J, Dhariwal P, et al. Language models are few-shot learners. arXiv Preprint posted online May 28, 2020 [FREE Full text]
- 62. OpenAI developer platform. OpenAI. URL: <u>https://platform.openai.com</u> [accessed 2024-04-29]
- 63. Zhou Y, Muresanu AI, Han Z, Paster K, Pitis S, Chan H, et al. Large language models are human-level prompt engineers. arXiv Preprint posted online November 3, 2022 [FREE Full text]
- 64. Pryzant R, Iter D, Li J, Lee YT, Zhu C, Zeng M. Automatic prompt optimization with "gradient descent" and beam search. arXiv Preprint posted online May 4, 2023 [FREE Full text] [doi: 10.18653/v1/2023.emnlp-main.494]
- 65. Sordoni A, Yuan X, Côté MA, Pereira M, Trischler A, Xiao Z, et al. Joint prompt optimization of stacked LLMs using variational inference. arXiv Preprint posted online June 21, 2023 [FREE Full text]
- 66. Sun H, Li X, Xu Y, Homma Y, Cao Q, Wu M, et al. AutoHint: automatic prompt optimization with hint generation. arXiv Preprint posted online July 13, 2023 [FREE Full text]
- 67. Yang C, Wang X, Lu Y, Liu H, Le QV, Zhou D, et al. Large language models as optimizers. arXiv Preprint posted online September 7, 2023 [FREE Full text]
- 68. Chen A, Dohan DM, So DR. EvoPrompting: language models for code-level neural architecture search. arXiv Preprint posted online February 28, 2023 [FREE Full text]
- 69. Fernando C, Banarse H, Michalewski H, Osindero S, Rocktäschel T. Promptbreeder: self-referential self-improvement via prompt evolution. arXiv Preprint posted online September 28, 2023 [FREE Full text]
- 70. Home. National Living Donor Assistance Center. URL: <u>https://www.livingdonorassistance.org/</u> [accessed 2025-09-01]
- Lipsey AF, Waterman AD, Wood EH, Balliet W. Evaluation of first-person storytelling on changing health-related attitudes, knowledge, behaviors, and outcomes: a scoping review. Patient Educ Couns 2020 Oct;103(10):1922-1934. [doi: <u>10.1016/j.pec.2020.04.014</u>] [Medline: <u>32359877</u>]

Abbreviations

BERT: Bidirectional Encoder Representations from Transformers **DUCC:** dialogue until classification consensus **LDKT:** living donor kidney transplantation **LKD:** living kidney donation **LLM:** large language model **NLP:** natural language processing

Edited by S Gardezi, F Dankar; submitted 12.02.24; peer-reviewed by GK Gupta, A Hassan, W Cheungpasitporn; comments to author 28.08.24; revised version received 18.09.24; accepted 18.11.24; published 07.02.25.

<u>Please cite as:</u> Nielsen J, Chen X, Davis L, Waterman A, Gentili M Investigating the Classification of Living Kidney Donation Experiences on Reddit and Understanding the Sensitivity of ChatGPT to Prompt Engineering: Content Analysis JMIR AI 2025;4:e57319 URL: <u>https://ai.jmir.org/2025/1/e57319</u> doi:<u>10.2196/57319</u> PMID:<u>39918869</u>

©Joshua Nielsen, Xiaoyu Chen, LaShara Davis, Amy Waterman, Monica Gentili. Originally published in JMIR AI (https://ai.jmir.org), 07.02.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Original Paper

Urgency Prediction for Medical Laboratory Tests Through Optimal Sparse Decision Tree: Case Study With Echocardiograms

Yiqun Jiang¹, PhD; Qing Li², PhD; Yu-Li Huang¹, PhD; Wenli Zhang³, PhD

¹Robert D. and Patricia E. Kern Center for the Science of Health Care Delivery, Mayo Clinic, Rochester, MN, United States

²Department of Industrial & Manufacturing Systems Engineering, Iowa State University, Ames, IA, United States

³Department of Information Systems and Business Analytics, Iowa State University, Ames, IA, United States

Corresponding Author: Wenli Zhang, PhD Department of Information Systems and Business Analytics Iowa State University 2167 Union Drive Ames, IA, 50011-2027 United States Phone: 1 5152942469 Email: wlzhang@iastate.edu

Abstract

Background: In the contemporary realm of health care, laboratory tests stand as cornerstone components, driving the advancement of precision medicine. These tests offer intricate insights into a variety of medical conditions, thereby facilitating diagnosis, prognosis, and treatments. However, the accessibility of certain tests is hindered by factors such as high costs, a shortage of specialized personnel, or geographic disparities, posing obstacles to achieving equitable health care. For example, an echocardiogram is a type of laboratory test that is extremely important and not easily accessible. The increasing demand for echocardiograms underscores the imperative for more efficient scheduling protocols. Despite this pressing need, limited research has been conducted in this area.

Objective: The study aims to develop an interpretable machine learning model for determining the urgency of patients requiring echocardiograms, thereby aiding in the prioritization of scheduling procedures. Furthermore, this study aims to glean insights into the pivotal attributes influencing the prioritization of echocardiogram appointments, leveraging the high interpretability of the machine learning model.

Methods: Empirical and predictive analyses have been conducted to assess the urgency of patients based on a large real-world echocardiogram appointment dataset (ie, 34,293 appointments) sourced from electronic health records encompassing administrative information, referral diagnosis, and underlying patient conditions. We used a state-of-the-art interpretable machine learning algorithm, the optimal sparse decision tree (OSDT), renowned for its high accuracy and interpretability, to investigate the attributes pertinent to echocardiogram appointments.

Results: The method demonstrated satisfactory performance (F_1 -score=36.18% with an improvement of 1.7% and F_2 -score=28.18% with an improvement of 0.79% by the best-performing baseline model) in comparison to the best-performing baseline model. Moreover, due to its high interpretability, the results provide valuable medical insights regarding the identification of urgent patients for tests through the extraction of decision rules from the OSDT model.

Conclusions: The method demonstrated state-of-the-art predictive performance, affirming its effectiveness. Furthermore, we validate the decision rules derived from the OSDT model by comparing them with established medical knowledge. These interpretable results (eg, attribute importance and decision rules from the OSDT model) underscore the potential of our approach in prioritizing patient urgency for echocardiogram appointments and can be extended to prioritize other laboratory test appointments using electronic health record data.

(JMIR AI 2025;4:e64188) doi:10.2196/64188

KEYWORDS

RenderX

interpretable machine learning; urgency prediction; appointment scheduling; echocardiogram; health care management

Introduction

Background

In the present medical landscape, the intricate interplay between innovative techniques has expanded the horizons of medical knowledge and opened avenues for unprecedented precision in patient care. The increasingly sophisticated laboratory tests play a crucial role in this transformative process. Born out of meticulous research and honed by the rigors of scientific scrutiny, these tests provide clinicians with a multifaceted toolkit to decipher the intricacies of illnesses, capturing the nuances of each condition, guiding medical professionals toward evidence-based interventions, and empowering medical professionals to tailor treatments with personalized precision.

However, a pivotal factor to take into consideration is the limited availability of certain state-of-the-art laboratory tests, as they often involve intricate equipment and elaborate protocols. This is evident from their expensive nature, the scarcity of skilled medical professionals capable of operating these laboratories, and the limited accessibility across different regions or during specific time frames [1]. As a result, the transformative potential of these laboratory tests is mitigated by the practical challenges they pose in terms of affordability [2]. The potential significant advantages of laboratory tests, coupled with their limited availability, render them a scarce resource, resulting in many patients having to endure wait times for access to laboratory tests. Consequently, predicting and prioritizing which patients require testing has emerged as an important research problem.

The rise of health IT and the subsequent influx of electronic health record (EHR) data, combined with the power of machine learning, offers new opportunities to revolutionize the prioritization of medical laboratory tests [3]. By delving into vast amounts of historical patient information, machine learning algorithms can discern intricate patterns and correlations that might otherwise elude human observation. The predictive outcomes generated by machine learning algorithms can contribute to refining testing protocols, enabling medical practitioners to make data-driven decisions regarding the prioritization and scheduling of laboratory tests based on patient information. In this study, we aim to elucidate methods for evaluating patients' urgency for tests, seeking to refine the allocation of scarce laboratory tests by harnessing the power of machine learning and analyzing historical EHRs. Specifically, we aim to contribute by applying an optimal sparse decision tree (OSDT) to a new domain-predicting the urgency of medical laboratory tests, using echocardiograms as a case study. Based on our literature review, OSDT stands out as one of the most suitable methods for achieving both optimal performance and interpretability in predicting the urgency of patients requiring echocardiograms. Our ultimate objective is to ensure prompt access for patients with the most critical needs.

Related Work

Echocardiogram and Patient Prioritization Techniques

An echocardiogram is one the most cost-effective means for screening cardiac anatomy, uses ultrasound to evaluate the cardiac structures, and provides critical information for medical

```
https://ai.jmir.org/2025/1/e64188
```

providers [4]. It functions as a crucial precursor to a detailed diagnosis, capable of screening cardiac anatomy and providing essential information for assessing cardiovascular conditions such as murmurs, stenosis, and regurgitation. Additionally, it plays a crucial role in diagnosing valvular morphology and uncovering the root causes of valve diseases [5]. A comprehensive echocardiographic assessment can provide both diagnostic and prognostic information, thus facilitating risk stratification and establishing baseline data for future evaluations [5].

The echocardiogram, although immensely valuable, is not always easily attainable due to the increasing demand for the test. For example, there has been an observed increase in the prevalence of rheumatic heart disease, which stands as the most predominant form of valvular heart disease and impacts approximately 41 million individuals in developing countries [6]. In recent years, there has been a notable escalation in the demand for pediatric cardiology services, leading to documented workloads that have exhibited a substantial upsurge of up to 51% over the past decades [7]. Furthermore, there has been an increase in the prevalence of children with asymptomatic murmurs who necessitate evaluation through echocardiogram [8]. The increasing demands pose challenges to echocardiogram laboratories in resource management, requiring medical institutions to establish more effective scheduling protocols to prioritize patients in critical need of echocardiogram lab appointments.

Patient prioritization techniques can be broadly classified into scoring systems and machine learning classification-based systems [9]. Scoring systems, particularly those using regression techniques, have gained prominence for their ability to allocate medical resources. These systems heavily rely on the expertise of medical professionals to assign priority scores to patients. Examples include the Salisbury priority scoring system, allowing surgeons to assign relative priorities, and the Italian waiting time prioritization system, which reallocates outpatient referrals based on clinical priorities prescribed by general practitioners [9]. These methods, however, exhibit various limitations. First, there may be inherent bias (eg, subjective judgments obtained through experience by medical professionals) as these approaches often necessitate input from medical specialists' judgments. A machine learning and data-driven method can serve as a complement to these types of systems. Second, these methods might be tailored for a particular patient prioritization task (eg, surgery or referral), and demand a high level of specialized medical knowledge for their design, making them difficult to generalize to other tasks [10]. Third, certain methods lack transparent decision rules for assessing the significance of input attributes, thereby posing challenges for their practical applications [11]. Machine learning classification-based methods typically rely on a large amount of patients' information (eg, EHRs) to autonomously discern patterns and generate predictions. This process aids in patient prioritization and avoids limitations associated with scoring systems [12]. The existing methods, however, fail to transform the prediction process and outcomes into clear and executable rules, limiting the practical application of these approaches [9]. Moreover, existing studies predominantly center around 5 clinical areas, including cataract

surgery, general surgical procedures, hip and knee replacements, magnetic resonance imaging scanning, and children's mental health using specific predictive attributes and expert systems [13]. There is a crucial need for new methods that apply more broadly to general laboratory test prioritization.

To summarize, our literature review underscores the need for new methods of prioritizing patients, which leverage machine learning and data-driven techniques to complement existing methods, ensure transparency, and have the potential to be generalized to various patient prioritization tasks. Consequently, using extensive patient historical EHRs combined with an interpretable machine learning approach emerges as a potential solution to address these gaps.

Leveraging Machine Learning for Optimizing the Use of Scarce Laboratories Tests

When a large number of patient EHRs, which contain numerous hidden patterns, are available, integrating machine learning into health care practices emerges as a potential solution to address pressing issues such as the continual demand for medical services outpacing available resources. Specifically, machine learning, with its capacity to analyze vast data and discern intricate patterns, empowers health care professionals to make data-driven decisions regarding the allocation of laboratory tests. By developing predictive models using historical EHRs, machine learning models can identify individuals who are more likely to benefit from specific tests, ensuring that scarce resources are allocated where they can yield the greatest impact. Furthermore, such methods ensure critical cases receive prompt attention, leading to expedited diagnoses and interventions [14]. Moreover, the prediction results can potentially streamline the testing process by reducing unnecessary tests [15].

The integration of machine learning techniques to optimize the allocation of limited medical tests and laboratory resources has attracted considerable research attention. Research by Elitzur et al [16] delves into the use of prediction models to allocate medical tests efficiently. The study uses historical patient data to develop models that identify the most suitable candidates for specific tests, thereby enhancing resource allocation and streamlining the testing process. In a similar vein, Marescotti et al [17] investigate the orchestration of laboratory workflows through machine learning-driven prioritization. By considering factors such as clinical urgency and resource availability, their work demonstrates how machine learning algorithms can ensure timely and effective laboratory test processing, contributing to both improved patient care and optimized resource use. Similarly, Zhang et al [18] estimate the probability of requiring mechanical ventilation for in-hospital patients and contribute to the literature by identifying which patients require medical devices (ie, critical medical resources) more urgently.

However, while the potential benefits of machine learning in optimizing resource allocation are evident, challenges remain. A recent study underscores the need for further research and development in the area of machine learning models' interpretability and fairness, ensuring that data-driven decisions in health care maintain transparency [19]. The research gap drives us to use an interpretable and efficient machine learning method for laboratory tests and patient optimization.

```
https://ai.jmir.org/2025/1/e64188
```

Interpretable Machine Learning

Medical research is often at the forefront of technological innovation, with machine learning algorithms being harnessed to analyze vast datasets, predict disease outcomes, and assist in clinical decision-making. However, as these algorithms become increasingly sophisticated, they tend to function as "black boxes," where the reasoning behind their predictions remains obscured. This opacity not only raises concerns about trustworthiness but also impedes the adoption and acceptance of these tools by medical professionals [19].

In medical research, the concept of interpretability holds profound significance. The intricate interplay between cutting-edge technology and human well-being underscores the critical need to not only generate accurate predictions but also to understand the underlying rationale behind those predictions. The complexity of medical data, coupled with the potential life-altering consequences of decisions made based on data and machine learning models, demands a heightened level of transparency and comprehensibility requirements [20].

The interpretability of machine learning models empowers health care providers to understand the factors that led to a specific decision, enabling them to fine-tune treatment strategies according to their medical judgment and the patient's unique circumstances. Consequently, there has been a surge in post hoc techniques for elucidating black box machine learning models in a manner interpretable by humans. The most prominent techniques among these include local, model-agnostic methods that aim to explain individual predictions of a given black box classifier, such as local interpretable model-agnostic Explanation and Shapley additive explanation [21]. Due to their high generalizability, post hoc methods have been used to explain a wide array of machine learning models across various domains. However, previous research has indicated that there are common limitations associated with these post hoc techniques, including local interpretability, sensitivity to perturbations, and difficulties in choosing interpretable surrogate models [21].

In health care, arguably, a more appropriate research direction for using interpretable machine learning is tree-based models because much of the data related to patient prioritization is structured data (eg, tabular EHRs). Tree-based machine learning models can perform comparably to complex models (eg, deep learning models), especially after thorough preprocessing of tabular data [22]. In contrast to post hoc explainable machine learning techniques, tree-based models are logical models that consist of statements involving logical operations, providing clear and interpretable decision rules [22]. This interpretability is highly valuable in health care, as it allows medical professionals to not only make accurate predictions but also understand the underlying factors driving those predictions, enhancing transparency and trust in the decision-making process.

Since our research aims to use historical EHR data for patient prioritization, it is crucial to acknowledge another notable characteristic of patient prioritization-related information: the prevalence of numerous categorical variables (eg, patient demographic information such as gender and age groups). Furthermore, the outcomes of patient prioritization are also

XSL•FO RenderX

expressed as categorical variables. For example, preventive interventions often involve categorical decisions, such as determining which individuals should undergo selective or indicated interventions or identifying those most likely to benefit from specific treatments [23]. In such scenarios, an efficient tree-based approach tailored to categorical variables is highly valuable. In this study, we focus on a cutting-edge decision tree algorithm–OSDT [24].

A decision tree features a hierarchical structure that is composed of a root node, branches, internal nodes, and leaf nodes in a tree format. Each path from the root node to the leaf node illustrates a rule to partition the data and leads to the final classification. The tree-based method presents a clear pattern for the decision-making process; thus, it is considered a transparent and highly interpretable model [25]. The results of the tree-based models are extremely useful for medical decision-making [26], and the performance of decision tree classifiers is verified by researchers on medical data [27]. Nevertheless, concerns have been raised regarding the suboptimality of decision tree algorithms [24,28]. To address this issue, OSDT has been introduced, aiming to ensure optimal solutions for binary variables in a computationally efficient manner [24].

The OSDT algorithm addresses various limitations observed in prior tree-based methods. Unlike previous approaches that often focused on finding the optimal tree within a fixed number of nodes or limited topology, OSDT tackles these shortcomings by identifying optimal trees through the use of a regularized loss function. This loss function strikes a balance between accuracy and the number of leaves, thereby enhancing the efficiency of the decision tree model. Furthermore, OSDT improves computational efficiency and interpretability by incorporating a series of analytical bounds that effectively reduce the search space while still identifying the optimal tree. By implementing these bounds, the algorithm streamlines the search process, leading to expedited identification of the optimal decision tree structure. Moreover, the OSDT algorithm has undergone mathematical validation, demonstrating its efficacy in constructing optimal trees for structured tabular datasets with attributes having binary values. It establishes its effectiveness in addressing binary classification problems. The algorithm is designed to uphold commendable levels of accuracy and is anticipated to meet the demands of medical prediction tasks with stringent interpretability requirements.

Methods

Study Design

In this study, we conducted empirical and predictive analyses using echocardiogram data extracted from EHRs at a large multispecialty hospital and medical facility. The dataset included administrative details, referral diagnoses, and patient conditions. To explore attributes relevant to echocardiogram prioritization, we used the OSDT algorithm due to its high accuracy and interpretability. We aim to enhance the scheduling of echocardiogram laboratory appointments by enabling the prioritization of patients with urgent needs based on our model's predictions. To be noted, our proposed method is not intended to replace human expertise but to complement it, offering valuable insights that guide practitioners toward informed and patient-centric choices.

Ethical Considerations

The Mayo Clinic Institutional Review Board, based on the authors' submission notes and in accordance with the Code of Federal Regulations, 45 CFR 46.102, deemed that this research did not require IRB review.

Data Collection and Selection

The dataset comprises real-world data from one of the top medical centers in the United States. The data were collected over a 1-year period in 2019, including 34,293 echocardiogram appointments. It consisted of 64 dummy-coded categorical attributes, encompassing various aspects such as patient demographics, medical history, clinical settings (eg, inpatient or outpatient status), past procedures, future scheduled procedures, and diagnose indicators for echocardiogram-justifying signs (eg, heart murmurs, shortness of breath, or chest pain) extracted from the clinical notes and referrals in the EHRs (Table 1).

The dataset exhibited a notable class imbalance issue, particularly evident in the examination of the "MadeBeforeEcho" attribute. This attribute delineates whether the downstream appointment following the echocardiogram occurs before the scheduling date of the echocardiogram appointment (not the actual appointment date). Within the "Y" category, the distribution revealed 84% nonurgent cases and 16% urgent cases. Conversely, in the "N" category, the distribution portrayed 58% nonurgent cases and 42% urgent cases. This observation underscored a substantial prevalence of nonurgent cases within the "MadeBeforeEcho" attribute. Furthermore, a similar pattern of imbalance is discerned when analyzing attributes such as "ReferredType" and "SurgeryYN." These attributes also exhibit a significant majority of cases concentrated within 1 category, indicating the need for careful consideration of class distribution in subsequent predictions.

The response variable is determined by calculating the number of days between the date the echocardiogram appointment was generated in the system and the actual appointment date. According to medical policy, appointments are classified as urgent (ie, the response variable) if the number of days is 2 or less, and nonurgent otherwise.

It is important to note that the features categorized under the "Future Scheduled Process" were derived based on the date the echocardiogram appointment is generated in the system, rather than the actual appointment date (Figure 1). This approach ensures that the model uses only the information available up to the point of echocardiogram appointment generation, without incorporating any data beyond this cutoff.

Of note, our dataset is a tabular dataset with attributes and response variables having binary values. Therefore, OSDT is highly suitable for serving this dataset, assisting us in making predictions for patient prioritization.



XSL•FC

 Table 1. Dataset and attribute statistics^a.

Category and variable	Description	escription Summary statistics, n (%)	
		Nonurgent	Urgent
Demographics			
Age (years)			
0-18	b	1929 (7.18)	478 (6.41)
19-55	_	6766 (25.19)	1930 (25.90)
56-65	_	4954 (18.45)	1342 (18.01)
66-75	_	6784 (25.26)	1896 (25.44)
Older than 75	_	6398 (23.82)	1775 (23.82)
Sex			
Female	—	11,829 (44.09)	3529 (47.55)
Male	—	15,002 (55.91)	3892 (52.45)
Patient geolocation			
In_State	—	9973 (37.14)	2376 (31.96)
Out_of_State	_	14,332 (53.37)	4301 (57.85)
Town	_	2550 (9.50)	758 (10.20)
Clinical settings			
ReferralType			
External	—	1156 (4.30)	606 (8.15)
Internal	—	25,699 (95.70)	6829 (91.85)
ReferredBy	The specialty that patient referred by		
Cardiovascular medicine	—	8188 (30.49)	1162 (15.63)
Family medicine	—	436 (1.62)	142 (1.91)
Hospital medicine	—	145 (0.54)	4 (0.05)
Internal medicine	—	978 (3.64)	591 (7.95)
Obstetrics and gynecology	—	1096 (4.08)	359 (4.83)
Pediatric and adolescent medicine	—	2302 (8.57)	401 (5.39)
Other	—	13,710 (51.05)	4776 (64.24)
ReferredFrom	Referral origin		
Arizona campus	—	2 (0.01)	0 (0.00)
Florida campus	—	1 (0.00)	0 (0.00)
Mayo Clinic health system	—	154 (0.57)	38 (0.51)
Rochester campus	—	17,495 (65.15)	4463 (60.03)
Other	—	9203 (34.27)	2934 (39.46)
ReferredType	Referred type		
Outpatient	—	18,706 (69.66)	4585 (61.52)
Other	—	8149 (30.34)	2868 (38.48)
Future scheduled process			
Diff_surgery_after	The number of days between the date the echocardiogram appointment was generated in the system and the surgery date		
0-1	_	1449 (5.40)	461 (6.20)
2-5	_	1607 (5.98)	492 (6.62)



XSL•FO RenderX

Category and variable	Description	Summary statistics,	n (%)
		Nonurgent	Urgent
6-15		1143 (4.26)	606 (8.15)
16 and greater	_	4715 (17.56)	1494 (20.09)
None	_	17,941 (66.81)	4382 (58.94)
MadeBeforeEcho	Whether the next downstream appointment after echocardiogram is made before the date the echocardiogram appointment was generated in the system or not		
Yes	_	23,845 (88.79)	4660 (62.53)
No	_	3010 (11.21)	2793 (37.47)
NextDepartment	The department in which the appointment happened after the date the echocardiogram appointment was generated in the system		
Cardiovascular medicine	_	12,012 (44.73)	1749 (23.47)
Non-cardiovascular medicine	Departments other than cardiovascular medicine	14,843 (55.27)	5704 (76.53)
NextLength	The number of days from the date the echocardiogram appointment was generated in the system to its following appointment		
0-1	_	4531 (16.87)	1608 (21.63)
1-5	_	3301 (12.29)	2018 (27.14)
Greater than 5	_	1,014 (3.78)	618 (8.31)
None	_	18,009 (67.06)	3191 (42.92)
Procedure	Type of echocardiogram visit		
TEE ^c	_	848 (3.16)	362 (4.87)
TTE^{d}	_	23,293 (86.74)	6803 (91.50)
Other	_	2714 (10.11)	270 (3.63)
Past procedures			
SurgeryYN	Whether the patient had a cardiovascular surgery within 6 months prior to the date the echocardiogram appointment was gener- ated in the system		
Yes	_	1708 (6.36)	264 (3.54)
No	_	25,147 (93.64)	7189 (96.46)
SurgeryYN_After	Whether the patient had a surgery within 3 months after the date the echocardiogram appointment was generated in the system		
Yes	_	8914 (33.19)	3053 (40.96)
No	_	17,941 (66.81)	4400 (59.04)
Medical history			
Alcohol	Alcohol abuse	115 (0.43)	50 (0.67)
Anemia	Anemia	962 (3.58)	605 (8.12)
BloodLoss	Blood loss	87 (0.32)	33 (0.44)
CHF ^e	_	1884 (7.02)	484 (6.49)
Coagulopathy	Coagulation deficiency	446 (1.66)	274 (3.68)
Depression	Major depressive disorder	439 (1.63)	192 (2.58)
DM^{f}	_	610 (2.27)	230 (3.09)

https://ai.jmir.org/2025/1/e64188

XSL•FO **RenderX** JMIR AI 2025 | vol. 4 | e64188 | p.247 (page number not for citation purposes)

_

C

Category and variable	Description	Summary statistics, n (%)	
		Nonurgent	Urgent
DMcx ^g		317 (1.18)	129 (1.73)
Drugs	Drug abuse	86 (0.32)	19 (0.25)
FluidsLytes	Fluid and electrolyte disorders	1013 (3.77)	617 (8.28)
HIV	_	0 (0.00)	1 (0.01)
Hypertension	_	2201 (8.20)	786 (10.55)
Hypothyroid	Hypothyroidism	777 (2.89)	277 (3.72)
Liver	_	429 (1.60)	197 (2.64)
Lymphoma	Lymph system cancer	464 (1.73)	347 (4.66)
Metastatic cancer	_	251 (0.93)	222 (2.98)
NeuroOther	Neurological disorders	581 (2.16)	291 (3.90)
Obesity	_	980 (3.65)	339 (4.55)
Paralysis	_	58 (0.22)	15 (0.20)
PHTN ^h	Pulmonary circulation disorders	298 (1.11)	153 (2.05)
Psychoses	Mental disorder characterized by a discon- nection from reality	126 (0.47)	53 (0.71)
PUD ⁱ	Chronic peptic ulcer	41 (0.15)	20 (0.27)
Pulmonary	Chronic pulmonary disease	650 (2.42)	273 (3.66)
PVD ^j	_	965 (3.59)	234 (3.14)
Renal	Renal failure	950 (3.54)	331 (4.44)
Rheumatic	Rheumatoid arthritis or collagen vascular	254 (0.95)	150 (2.01)
Tumor	Solid tumor	722 (2.69)	380 (5.10)
Valvular	Valvular disease	3367 (12.54)	573 (7.69)
WeightLoss	Weight loss	248 (0.92)	237 (3.18)
Diagnoses			
А	MSSA ^k bacteremia, sepsis	18 (0.07)	25 (0.34)
В	MRSA ^l , staph bacteremia, slaph, fungemia, pseudomonas, candidemia, MRSA bac- teremia	47 (0.18)	40 (0.54)
С	Leukemia, AML ^m , CML ⁿ , lymphoma, AMV ^o , myeloma	1428 (5.32)	554 (7.43)
D	Diseases of the blood and blood-forming organs and certain disorders involving the immune mechanism	561 (2.09)	193 (2.59)
Е	Endocrine, nutritional and metabolic diseases	1714 (6.38)	408 (5.74)
F	Behavioral and neurodevelopmental disor- ders	49 (0.18)	46 (0.62)
G	Muscular dystrophy	590 (2.20)	273 (3.66)
Н	Diseases of the eye and adnexa or disease of the ear and mastoid process	60 (0.22)	28 (0.38)
Ι	Heart failure, coronary artery, cardiac arrest, STEMI ^p , stroke, cardia, hypertension, endo- carditis, NSTEMI ^q , PEA ^r arrest, AFib ^s , pulmonary embolism, pulmonary hyperten- sion, and vegetation	11,302 (42.09)	4096 (54.96)

XSL-FO **RenderX**

Category and variable	Description	Summary statistics, n (%)	
		Nonurgent	Urgent
J	Resp failure, respiratory, and pulmonary	477 (1.78)	392 (5.26)
K	Liver and cirrhosis	357 (1.33)	130 (1.74)
L	Diseases of the skin and subcutaneous tissue	36 (0.13)	33 (0.44)
М	Diseases of the musculoskeletal system and connective tissue	503 (1.87)	280 (3.76)
Ν	Diseases of the genitourinary system	397 (1.48)	119 (1.60)
0	Pre-eclampsia, preeclampsia	235 (0.88)	57 (0.76)
Р	Certain conditions originating in the perina- tal period	12 (0.04)	4 (0.05)
Q	Ehlers, coarc, PDA ^t , and congenital	2811 (10.47)	309 (4.15)
R	Murmur, hypoxemia, shortness, SOB ^u , breath, shock, dyspnea, chest pain, troponin, syncope, electrocardiogram, extremity, mass, and swelling, edema	4111 (15.31)	2811 (37.72)
S	Injury, poisoning and certain other conse- quences of external causes	100 (0.37)	21 (0.28)
Z	Chemo, preoperative, pre-op, prenatal, pregnancy, prior to, BMI, surgery, and transplant	5966 (22.22)	1129 (15.15)

^aAll the features used in this study are complete for each patient, with no missing values. The diagnoses are derived from patients' ICD-9 codes, and the medical history is extracted from electronic health record notes using the medical center's built-in natural language processing tools. ^bNot applicable.

^cTEE: transesophageal echocardiogram.

^dTTE: transthoracic echocardiogram.

^eCHF: congestive heart failure.

^fDM: diabetes without chronic complications.

^gDMcx: diabetes with chronic complications.

^hPHTN: pulmonary hypertension.

ⁱPUD: peptic ulcer disease.

^jPVD: peripheral vascular disease.

^kMSSA: methicillin-sensitive *Staphylococcus aureus*.

¹MRSA: methicillin-resistant *Staphylococcus aureus*.

^mAML: acute myeloid leukemia.

ⁿCML: chronic myeloid leukemia.

^oAMV: avian myeloblastosis virus.

^pSTEMI: ST-elevation myocardial infarction.

^qNSTEMI: non-ST-elevation myocardial infarction.

^rPEA: pulseless electrical activity.

^sAFib: atrial fibrillation.

^tPDA: patent ductus arteriosus.

^uSOB: shortness of breath.



Figure 1. Timeline and process of echocardiogram appointment scheduling. Using MadeBeforeEcho as an example.



Problem Formulation: Urgency Prediction Using OSDT

With data \bowtie , where \bowtie are *M* binary attributes and \bowtie are the response variable, we model an OSDT tree *d* with a collection of *H* distinct leaves $d = (p_1, p_2, ..., p_H)$. The objective function in this study integrates the misclassification error with a sparsity penalty imposed on the number of leaf nodes, denoted as R(d,x,y). $R(d,x,y) = l(d,x,y) + \lambda H_d$, where l(d,x,y) represents the misclassification error of the tree, which is computed as the fraction of training data with incorrectly predicted labels. In addition, H_d represents the number of leaves in tree *d*. To regularize the model and discourage larger trees, a regularization term λH_d is introduced, where λ is a hyperparameter controlling the strength of the penalty. A higher value of λ corresponds to a stronger penalty on the size of the tree. This implies that the tree is more likely to be shallower when achieving optimality.

By using OSDT, we aim to improve the overall performance of the classification task while simultaneously upholding a significant level of interpretability, thereby facilitating a comprehensive understanding of the underlying patterns and factors influencing the classification outcomes.

Results

Overview

In this section, we evaluated the proposed method against state-of-the-art machine learning models. We then highlighted attribute importance and provided clear interpretations of derived results within specific patient cohorts for transparency and clarity.

Performance Evaluation

We demonstrated the performance of our OSDT model by comparing it to commonly used machine learning models as

```
https://ai.jmir.org/2025/1/e64188
```

RenderX

baselines, including naive Bayes, generalized linear model, fast large margin, logistic regression, neural network, vanilla decision tree, random forest, gradient boosted trees, and support vector machine. The evaluation metrics used for the binary classification are accuracy, precision, recall, F1-score, and F_2 -score. Accuracy is a metric that quantifies the overall correctness of a machine learning model. It represented the proportion of correct predictions made by the model across all categories or classes. Precision and recall, on the other hand, measured the model's ability to accurately predict a specific category or class. Precision focused on the proportion of true positive predictions relative to all positive predictions made by the model. Recall, also known as sensitivity, gauged the model's capability to correctly detect instances of a specific category. It quantified the proportion of true positive predictions relative to all actual positive instances present in the data. The F_1 -score has been widely used in the context of imbalanced classification problems and serves as a prominent metric. It is computed as the harmonic mean of the precision and recall scores, providing a balanced assessment of the model's performance by considering both precision and recall simultaneously. The F_2 -score assigns greater weight to recall than precision, proving beneficial when the consequences of false negatives (ie, missed positive cases where patients are in urgent condition but remain unidentified by the model) outweigh those of false positives (ie, incorrectly identified positive cases). All metrics mentioned exhibited a range of values between 0 and 1, whereby a higher value indicated superior performance.

Compared with various baselines, the performance of the OSDT model achieved the highest accuracy, recall, F_1 -score, and F_2 -score (Table 2). The performance reported is based on 5-fold cross-validation. These results indicated the predictive capability of the OSDT model in our research context, demonstrating the overall performance and effectiveness of the OSDT model.

Table 2. OSDT^a performance comparisons with baselines^b.

Algorithm	Accuracy (%), mean (SD)	Precision (%), mean (SD)	Recall (%), mean (SD)	<i>F</i> ₁ -score (%), mean (SD)	F_2 -score ^c (%), mean (SD)
Naïve Bayes	78.86 (0.24)	81.3 (7.11)	3.34 (0.59)	6.41 (1.09)	4.13 (1.02)
Generalized linear model	79.23 (0.22)	78.05 (5.00)	5.93 (0.69)	11.01 (1.03)	7.27 (0.93)
Fast large margin	80.26 (0.47)	68.94 (2.57)	17.76 (1.4)	28.21 (1.7)	20.86 (2.17)
Logistic regression	79.26 (0.22)	77.68 (4.26)	6.16 (0.86)	11.41 (1.49)	7.55 (0.78)
Deep learning	80.49 (0.29)	85.59 (4.59)	12.14 (0.39)	21.26 (0.66)	14.66 (0.56)
Decision tree	80.69 (0.2)	69.18 (4.5)	22.45 (4.1)	33.53 (4.5)	25.96 (3.15)
Random forest	79.45 (0.18)	78.19 (5.54)	7.34 (0.31)	13.42 (0.57)	8.96 (2.67)
Gradient boosted trees	80.64 (0.29)	80.8 (2.96)	14.94 (1.55)	25.18 (2.25)	17.85 (1.95)
SVM ^d	80.3 (0.84)	61.42 (5.57)	24.06 (3.4)	34.48 (4.02)	27.39 (1.95)
OSDT (ours)	81.21 (0.20)	68.75 (1.7)	24.56 (0.59)	36.18 (0.66)	28.18 (0.55)

^aOSDT: optimal sparse decision tree.

^bOSDT is an algorithm that makes decisions based on direct constraints rather than generating probability scores. As a result, metrics like the receiver operating characteristic curve, precision and recall curve, and area under curve are not applicable for this method. Although the CIs for SVM and OSDT overlap, it is noteworthy that SVM exhibits a significantly larger SD. This indicates that OSDT is more robust in this scenario, delivering a more stable and reliable performance despite the overlapping intervals.

^c **≥**; β=2.

^dSVM: support vector machine.

Interpreting Prediction Results

OSDT, as a tree-based model, possesses the notable advantage of providing interpretable prediction results. We conducted an analysis of the decision trees generated using the entire dataset as well as specific patient cohorts. The objective is to extract the most influential rules that demonstrate both high accuracy and coverage, thereby aiming to uncover the underlying factors that drive the urgent decision of echocardiogram appointments.

We first identified several key categories and attributes that significantly influenced the urgency of patients' echocardiogram appointments (Table 3). First, the most important categories included "future scheduled process," pertaining to clinic scheduling policies, and "diagnosis," indicative of patients' health conditions. Second, within the top 12 important attributes, a cluster of attributes related to future scheduled processes emerged as the most prominent. These attributes encompassed scenarios if the next downstream appointment following the echocardiogram was scheduled prior to the echocardiogram appointment (ie, "MadeBeforeEcho"), instances where the next appointment did not pertain to the cardiovascular department (ie, "NextDepartment"), cases where no subsequent appointment was scheduled after the echocardiogram appointment (ie, "NextLength_None"), and situations where the time gap between the echo appointment and the subsequent one was less than a day ("NextLength_1"). The absence of a downstream appointment before the echocardiogram could be attributed to the clinic's practice of tailoring subsequent appointments based on the results of the echocardiogram. Consequently, it became imperative for medical providers to accord priority to the echocardiogram appointments of these patients, as the results would furnish vital evidence for guiding appropriate follow-up care and future steps. Third, attributes related to diagnoses

https://ai.jmir.org/2025/1/e64188

assumed the second tier of importance, particularly whether patients exhibited respiratory and cardiac symptoms (ie, "R") or had documented cardiovascular conditions (ie, "I"). Patients diagnosed with heart-related issues, such as heart murmurs, shortness of breath, and chest pain, typically require expedited access to echocardiography results to determine the next course of action. Fourth, clinical setting attributes and demographic information are also important to patient prioritization. In the context of inpatients, health care providers tended to assign earlier echocardiogram appointment slots as part of a strategy to reduce the length of hospital stays. Additionally, when prioritizing patients with heart conditions, individuals referred by cardiologists received preferential treatment in terms of scheduling. Furthermore, the medical facility providing the data adopted a proactive approach by expediting echocardiogram appointments for out-of-state patients, aiming to minimize their duration of stay. This proactive stance facilitated timely evaluation and management, thereby contributing to a more efficient allocation of resources and an enhanced patient experience. Among medical history attributes, the presence of fluid and electrolyte disorders (ie, "FluidsLytes") emerged within the top 12, which underscored the strong correlation between fluid and electrolyte disorders and heart failure, further emphasizing its relevance in patient prioritization [29].

These results underscore the significance of admission and policy-related information in determining the urgency of echocardiogram appointments. They reflected the complexities of the scheduling process and highlighted the need for tailored appointment allocation strategies based on patients' referral status and downstream appointment requirements.

We subsequently focus on a specific patient cohort for further analysis. The "MadeBeforeEcho" attribute clearly emerged as

exceptionally significant among the dataset's attributes. It was noteworthy to highlight that, based on the data, there were no urgent cases when the "MadeBeforeEcho" variable was marked as "N." Consequently, we conducted an investigation specifically focusing on patients whose subsequent downstream appointment was scheduled before the date the echocardiogram appointment was generated in the system. This subset of the patient cohort served as an illustrative example of how decision trees could provide a high degree of interpretability in the context of patient prioritization (Figure 2). Upon scrutiny of the subdecision tree for this cohort depicted, several noteworthy observations emerged. Primarily, it became evident that the most crucial attribute for this cohort is "R," signifying whether the patient presents with respiratory and cardiac symptoms, which served as the root node of the subtree. The pathway leading to categorizing a patient case as urgent depended on multiple conditions: the patient exhibited respiratory and cardiac symptoms, had an appointment scheduled within the cardiology department, hailed from out of state, and had a subsequent appointment scheduled following the echocardiogram. In contrast, patients without respiratory and cardiac symptoms tended toward classification as nonurgent. This tendency toward nonurgency was particularly pronounced in cases lacking a scheduled appointment subsequent to the echocardiogram.

Table 3. Attribute importance and category importance^a.

Category and attribute	Meanings	Attribute importance			
Future scheduled process (importance=0.0369)					
MadeBeforeEcho	Whether the next downstream appointment after echocardiogram is made before the date the echocardiogram appointment was generated in the system or not.	0.0279			
NextDepartment	The department in which the appointment happened after the date the echocardiogram appointment was generated in the system.	0.0049			
NextLength_None	No following appointment scheduled after the date the echocardio- gram appointment was generated in the system.	0.0035			
NextLength_1	The number of days from the date the echocardiogram appointment was generated in the system to its following appointment is less than 1 day.	0.0006			
Diagnoses (importance=0.0154)					
R	If have murmur, hypoxemia, shortness, SOB ^b , breath, shock, dysp- nea, chest pain, troponin, syncope, electrocardiogram, extremity, mass, swelling, and edema.	0.0147			
Ι	If have heart failure, coronary artery, cardiac arrest, STEMI ^c , stroke, cardia, hypertension, endocarditis, NSTEMI ^d , PEA ^e arrest, AFib ^f , pulmonary embolism, pulmonary hypertension, and vegetation.	0.0007			
Demographic (importance=0.0369)					
Geo_Out of State	Patient is from out of state.	0.0029			
Geo_Town	Patient is from the local town.	0.0013			
AGE_19-55	Age between 19 and 55 years.	0.0011			
Clinical settings (importance=0.0053)					
ReferredType	Referred type-inpatient or outpatient.	0.0047			
ReferredBy_CV	The specialty that patient referred by is cardiovascular disease department.	0.0006			
FluidsLytes (medical history; impor- tance=0.0021)	If have fluid and electrolyte disorders	0.0021			

^aThe relative importance scores of the attribute category and individual attributes are determined by the Gini index of the optimal sparse decision tree. The feature importance values are relative importance values and do not have a fixed absolute range. We presented only the most important features. ^bSOB: shortness of breath.

^cSTEMI: ST-elevation myocardial infarction.

^dNSTEMI: non-ST-elevation myocardial infarction.

^ePEA: pulseless electrical activity.

^fAFib: atrial fibrillation.


Figure 2. The OSDT for patients whose next downstream appointment after the echocardiogram is scheduled before the date the echocardiogram appointment was generated in the system. OSDT: optimal sparse decision tree. λ =0.0008; accuracy: 83.69%.

Analyses on Diverse Patient Cohorts

In order to enhance the validity of the decision trees and gain more valuable medical insights, we conducted more analyses on smaller patient cohorts. Specifically, we focus on patients who have no next downstream appointment after echocardiogram and are categorized as inpatients. Furthermore, we narrowed down the patient cohort based on specific medical history and presented a compilation of rules extracted from the decision tree (Table 4).

A decision rule was defined as the pathway from the root of a

decision tree to a leaf node $\boxed{|\mathbf{x}||}$. The accuracy and coverage of a decision rule served as critical metrics for evaluating its effectiveness and applicability. Accuracy, denoting the capacity of a decision rule to effectively forecast the outcome of interest, was quantified as the proportion of records that fulfill both the rule's precondition and its consequent within the precondition.

This metric was computed as \square , where "number of Correct Predictions" denoted the count of instances where the decision rule accurately anticipated the desired outcome and "Total number of Instances" represented the entire dataset or the set of instances under consideration, which elucidated how accuracy measures the precision of a decision rule in making predictions based on its specified conditions and its congruence with actual outcomes within the dataset. Coverage, on the other hand, measured the proportion of cases or individuals to which the

decision rule could be applied. It could be calculated as \square . It signified the generalizability and practical scope of the rule in real-world scenarios. A decision rule with high coverage indicates its ability to be applied to a wide range of cases or individuals, thereby increasing its usefulness in practice.

In the context of patients with congestive heart failure (CHF), anemia played a significant role in determining the urgency of

RenderX

echocardiogram appointments (Table 4). Anemia could have detrimental effects on cardiac function through various mechanisms [29]. First, it induces cardiac stress by increasing heart rate and stroke volume. Additionally, anemia could lead to reduced renal blood flow and fluid retention, adding further strain to the heart. Prolonged anemia, regardless of its underlying cause, could contribute to the development of left ventricular hypertrophy, which exacerbates CHF by promoting cardiac cell death through apoptosis. Notably, patients with anemic CHF often exhibited resistance to CHF medications, and numerous studies consistently demonstrated that these individuals have a higher mortality rate compared to patients with non-anemic CHF [30]. Anemia also played a critical role in patients with coagulopathy, as it exacerbated bleeding, which in turn further worsens coagulopathy [30].

For patients with hypothyroidism, fluid and electrolyte disorders served as strong indicators. Hypothyroidism, a prevalent endocrine disorder, was associated with the development of congestive heart failure. Electrolyte disturbances were commonly observed in patients with chronic heart failure [31]. Echocardiogram has been a suitable modality for guiding fluid resuscitation in critically ill individuals. It allowed for the evaluation of fluid responsiveness based on several parameters, such as the left ventricle, aortic outflow, inferior vena cava, and right ventricle [32].

The impact of alcohol consumption on cardiovascular health was multifaceted. Extensive research has demonstrated that the consumption of alcohol at levels surpassing approximately 1 to 2 drinks per day was associated with hypertension [28]. This condition adversely affects the elasticity of arteries, leading to diminished blood and oxygen flow to the heart and consequently contributing to the onset of heart disease [33]. These pathophysiological changes increase the risk of heart disease. Consequently, patients with a history of alcohol abuse and

concomitant hypertension might require an urgent echocardiogram to assess the potential cardiac implications arising from these interconnected conditions.

Patients diagnosed with valvular heart conditions would fall into the urgent category if they also exhibited cardiovascular issues and a history of congestive heart failure. These attributes collectively signaled the presence of potentially serious cardiac problems, indicating a compelling need for an echocardiogram to obtain detailed cardiac information and facilitate accurate diagnoses. In the case of patients grappling with depression, their urgency classification as "urgent" was contingent upon the presence of co-occurring health issues. Extensive research has established a substantial influence of depression on the outcomes of concurrent medical conditions. Consequently, when depression coincided with other health problems, it necessitated an "urgent" classification, acknowledging its significant impact on overall health outcomes [34]. Regarding patients with obesity, an "urgent" classification applied if they additionally exhibited fluid and electrolyte disorders. Research findings have illuminated a connection between overweight or obesity and

specific physiological factors, such as lower reactance and hypertonicity. Furthermore, individuals with overweight and those with obesity with lower reactance tended to demonstrate significantly elevated serum sodium levels compared to individuals with a normal weight. These associations underscored the importance of promptly addressing the medical needs of patients with obesity with fluid and electrolyte disorders, warranting an "urgent" classification for their cases [35].

Overall, the decision rules extracted from our analyses aligned closely with medical knowledge, providing reliable insights for identifying urgent echocardiogram appointments for patients. The congruence between the rules and medical understanding not only validated the effectiveness of our model but also highlighted the consistent application of medical principles in the decision-making process. This focused analysis contributed to a better understanding of the OSDT model's validity and offered valuable medical perspectives to enhance the identification of urgent patients' echocardiogram appointments.

Table 4. Decision rules for specific patient cohorts.

Cohort	Rules for a patient to be classified as urgent	Rule accuracy (%)	Rule coverage (%)
CHF ^a	The department in which the appointment happened after the echocardiogram appointment was generated in the system=non-cardiovascular disease, AGE<75, anemia=yes	100	14.20
Coagulopathy	Anemia=Yes	99	53.03
Hypothyroid	Fluid and electrolyte disorders=yes, Whether the patient had a cardiovascular surgery within six months prior to the echocardiogram appointment=no	100	32.91
Alcohol	Hypertension=yes	100	43.75
Valvular	I=1(has cardiovascular conditions), CHF=yes	100	6.36
Depression	Z=1 (has factors influencing health status and contact with health service)	100	24.49
Obesity	Geo!=Town, E=0 (has no nutritional and metabolic diseases), fluid and electrolyte disorders=yes	100	23.75

^aCHF: congestive heart failure.

Discussion

Overview

The primary objective of our study is to forge an effective tree-based classification machine learning model geared toward prioritizing the allocation of echocardiogram appointments for patients with a heightened need for timely diagnostics. Our long-term goal is to streamline the scheduling process, ensuring that patients' medical requirements are promptly addressed, thereby minimizing delays and optimizing their health care experience. Moreover, our study aspired to delve deeper into the intricate attributes that contribute to the urgency of echocardiogram lab appointments. Recognizing the intricate interplay of medical, logistical, and patient-specific variables, we sought to unravel the complex rules and dynamics that govern appointment prioritization. By harnessing the inherent interpretability of our model, we aim to uncover hidden insights and relationships within a large amount of EHR data, shedding light on the critical determinants that underscore the need for rapid scheduling. The implications of our study extended beyond

the realm of predictive modeling. We aimed to empower health care professionals with a powerful tool that not only optimizes resource allocation but also enriches their decision-making process.

Principal Results

The findings demonstrate promising results by accurately predicting the urgency of echocardiogram appointments and providing valuable insights into the critical guidelines applicable to specific patient cohorts. In summary, the study emphasizes two key points: (1) among the various attributes examined, it is observed that admission-related attributes exert a significant influence on the level of urgency for patients' echocardiogram appointments; and (2) the urgency of scheduling echocardiogram appointments can be influenced by the presence of comorbidities that exacerbate patients' conditions. In the case of congestive heart failure, anemia emerges as a significant attribute, highlighting its relevance in contributing to the urgency of echocardiogram appointments. Similarly, coagulopathy is identified as an important attribute for patients with congestive heart failure, further emphasizing the need for prompt

assessment. For patients with hypothyroidism, the presence of fluid and electrolyte disorders serves as a concerning indicator, warranting the prioritization of an echocardiogram. Additionally, hypertension is found to be a critical medical knowledge for patients with a history of alcohol abuse, underscoring the urgency of echocardiogram in this population.

Our work is unique in applying an advanced binary decision tree model that offers inherent interpretability, avoiding the limitations of post hoc techniques like local interpretable model-agnostic Explanation and Shapley additive explanation, such as local interpretability constraints, sensitivity to perturbations, and difficulties in selecting appropriate surrogate models. We extract interpretable rules grounded in medical knowledge, making this the first study to introduce tree-based interpretable machine learning for patient prioritization and the stratification of medical test urgency. Furthermore, the tree-based model allows us to derive rules that are easily understandable to medical professionals. These rules can be assessed for alignment with existing medical knowledge and applied in real-world practice by health care providers.

Limitations

The research has several limitations that could be addressed in future work. First, the accuracy of the prediction model hinges on the quality and completeness of available data; incomplete or missing data may compromise the reliability of predictions. Furthermore, it is essential to recognize that the effectiveness of the model may vary when applied to diverse patient populations or health care settings. This variation can be attributed to the unique attributes and patterns present in the training data, which significantly impact the model's performance. Moreover, the predictions rely on the elapsed days between the appointment scheduling date and the appointment date. Nonurgent patients may inadvertently be grouped with urgent cases due to cancellations and rescheduling of echocardiogram appointments. While this offers a broad indication of urgency, it may overlook critical factors that influence appointment priority. Integrating essential clinical or contextual details, such as the patient's medical history, symptom severity, or health care resource availability, into the model could provide more comprehensive insights.

Conclusions

This research adapts the OSDT algorithm to assess the urgency of patients in need of echocardiograms. The OSDT model demonstrates better performance over alternative machine learning models, highlighting its predictive accuracy and effectiveness. Furthermore, it identifies key attributes and rules governing the prioritization of echocardiogram appointments.

The analysis of decision trees generated by the OSDT model reveals the significance of admission- and policy-related attributes, such as downstream appointment scheduling and patient referral status, in determining appointment urgency. Moreover, the analyses of specific patient cohorts provide medical insights into the role of comorbidities, such as anemia in patients with CHF and coagulopathy, and fluid and electrolyte disorders in patients with hypothyroidism. These insights align with established medical knowledge and enhance the identification of urgent echocardiogram appointments.

In summary, this study facilitates the development of effective scheduling protocols for echocardiogram appointments by harnessing machine learning techniques and integrating medical insights. This approach enhances the overall efficiency and effectiveness of echocardiogram services, ultimately benefiting patient care. The findings can also be generalized to inform the establishment of efficient scheduling protocols and the promotion of equitable access to various other medical laboratory tests.

Acknowledgments

In this research, the authors gratefully acknowledge the financial support provided by the Ivy College of Business and the Department of Industrial and Manufacturing Systems Engineering at Iowa State University. The authors also extend our appreciation to the Mayo Clinic for generously providing essential data. Their collaborative efforts significantly enriched our study.

Conflicts of Interest

None declared.

References

- 1. Danzon PM, Manning WG, Marquis MS. Factors affecting laboratory test use and prices. Health Care Financ Rev 1984;5(4):23-32 [FREE Full text] [Medline: 10317549]
- Bhatt J, Bathija P. Ensuring access to quality health care in vulnerable communities. Acad Med 2018;93(9):1271-1275 [FREE Full text] [doi: 10.1097/ACM.0000000002254] [Medline: 29697433]
- 3. Bohr A, Memarzadeh K. The rise of artificial intelligence in healthcare applications. Artif Intell Healthc Elsevier 2020:25-60. [doi: 10.1016/b978-0-12-818438-7.00002-2]
- 4. Ashley EA, Niebauer J. Cardiology Explained. London, United Kingdom: Remedica; 2004.
- Cheitlin MD, Armstrong WF, Aurigemma GP, Beller GA, Bierman FZ, Davis JL, et al. ACC/AHA/ASE 2003 guideline update for the clinical application of echocardiography: summary article. A report of the American College of Cardiology/American Heart Association Task Force on Practice Guidelines (ACC/AHA/ASE Committee to update the 1997 guidelines for the clinical application of echocardiography). J Am Soc Echocardiogr 2003;16(10):1091-1110. [doi: 10.1016/S0894-7317(03)00685-0] [Medline: 14566308]

- Aluru JS, Barsouk A, Saginala K, Rawla P, Barsouk A. Valvular heart disease epidemiology. Med Sci 2022;10(2):32 [FREE Full text] [doi: 10.3390/medsci10020032] [Medline: 35736352]
- Pushparajah K, Garvie D, Hickey A, Qureshi SA. Managed care network for the assessment of cardiac problems in children in a district general hospital: a working model. Arch Dis Child 2006;91(11):892-895 [FREE Full text] [doi: 10.1136/adc.2005.086058] [Medline: 16717084]
- Murugan SJ, Thomson J, Parsons JM, Dickinson DF, Blackburn MEC, Gibbs JL. New outpatient referrals to a tertiary paediatric cardiac centre: evidence of increasing workload and evolving patterns of referral. Cardiol Young 2005;15(1):43-46. [doi: 10.1017/S1047951105000090] [Medline: 15831160]
- Mariotti G, Siciliani L, Rebba V, Fellini R, Gentilini M, Benea G, et al. Waiting time prioritisation for specialist services in Italy: the homogeneous waiting time groups approach. Health Policy 2014;117(1):54-63. [doi: 10.1016/j.healthpol.2014.01.018] [Medline: 24576498]
- 10. Solans-Domènech M, Adam P, Tebé C, Espallargues M. Developing a universal tool for the prioritization of patients waiting for elective surgery. Health Policy 2013;113(1-2):118-126. [doi: <u>10.1016/j.healthpol.2013.07.006</u>] [Medline: <u>23932414</u>]
- 11. Silva-Aravena F, Morales J. Dynamic surgical waiting list methodology: a networking approach. Mathematics 2022;10(13):2307. [doi: 10.3390/math10132307]
- Holt JM, Wilk B, Birch CL, Brown DM, Gajapathy M, Moss AC, Undiagnosed Diseases Network, et al. VarSight: prioritizing clinically reported variants with binary classification algorithms. BMC Bioinformatics 2019;20(1):496 [FREE Full text] [doi: 10.1186/s12859-019-3026-8] [Medline: 31615419]
- 13. Abbasgholizadeh Rahimi S, Jamshidi A, Ruiz A, Ait-kadi D. A new dynamic integrated framework for surgical patients' prioritization considering risks and uncertainties. Decis Support Syst 2016;88:112-120. [doi: 10.1016/j.dss.2016.06.003]
- Rabbani N, Kim GYE, Suarez CJ, Chen JH. Applications of machine learning in routine laboratory medicine: current state and future directions. Clin Biochem 2022;103:1-7 [FREE Full text] [doi: <u>10.1016/j.clinbiochem.2022.02.011</u>] [Medline: <u>35227670</u>]
- 15. Javaid M, Haleem A, Pratap Singh R, Suman R, Rab S. Significance of machine learning in healthcare: features, pillars and applications. Int J Intell Netw 2022;3:58-73. [doi: <u>10.1016/j.ijin.2022.05.002</u>]
- 16. Elitzur R, Krass D, Zimlichman E. Machine learning for optimal test admission in the presence of resource constraints. Health Care Manag Sci 2023;26(2):279-300 [FREE Full text] [doi: 10.1007/s10729-022-09624-1] [Medline: 36631694]
- Marescotti D, Narayanamoorthy C, Bonjour F, Kuwae K, Graber L, Calvino-Martin F, et al. AI-driven laboratory workflows enable operation in the age of social distancing. SLAS Technol 2022;27(3):195-203 [FREE Full text] [doi: 10.1016/j.slast.2021.12.001] [Medline: 35058197]
- Zhang K, Jiang X, Madadi M, Chen L, Savitz S, Shams S. DBNet: a novel deep learning framework for mechanical ventilation prediction using electronic health records. 2021 Presented at: BCB '21: 12th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics; August 1-4, 2021; Gainesville, FL p. 1-8. [doi: 10.1145/3459930.3469551]
- 19. Azimi V, Zaydman M. Optimizing equity: working towards fair machine learning algorithms in laboratory medicine. J Appl Lab Med 2023;8(1):113-128. [doi: 10.1093/jalm/jfac085] [Medline: 36610413]
- 20. Chaddad A, Peng J, Xu J, Bouridane A. Survey of explainable AI techniques in healthcare. Sensors 2023;23(2):634 [FREE Full text] [doi: 10.3390/s23020634] [Medline: 36679430]
- 21. Slack D, Hilgard S, Jia E, Singh S, Lakkaraju H. Fooling LIME and SHAP: adversarial attacks on post hoc explanation methods. 2020 Presented at: AIES '20: AAAI/ACM Conference on AI, Ethics, and Society; February 7-9, 2020; New York, NY p. 180-186. [doi: 10.1145/3375627.3375830]
- 22. Rudin C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. Nat Mach Intell 2019;1(5):206-215 [FREE Full text] [doi: 10.1038/s42256-019-0048-x] [Medline: 35603010]
- 23. Wiedermann W, Bonifay W, Huang FL. Advanced categorical data analysis in prevention science. Prev Sci 2023;24(3):393-397. [doi: 10.1007/s11121-022-01485-y] [Medline: 36633766]
- 24. Hu X, Rudin C, Seltzer M. Optimal sparse decision trees. ArXiv Preprint posted online on October 1, 2019. [doi: 10.5260/chara.21.2.8]
- 25. Dosilovic FK, Brcic M, Hlupic N. Explainable artificial intelligence: a survey. 2009 Presented at: 2018 41st International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO); May 21-25, 2018; Opatija, Croatia p. 0210-0215.
- 26. Podgorelec V, Kokol P, Stiglic B, Rozman I. Decision trees: an overview and their use in medicine. J Med Syst 2002;26(5):445-463. [doi: 10.1023/a:1016409317640] [Medline: 12182209]
- 27. Lavanya D, Rani KU. Performance evaluation of decision tree classifiers on medical datasets. IJCA 2011;26(4):1-4. [doi: 10.5120/3095-4247]
- 28. Piano MR. Alcohol's effects on the cardiovascular system. Alcohol Res 2017;38(2):219-241 [FREE Full text] [Medline: 28988575]
- 29. Urso C, Brucculeri S, Caimi G. Acid-base and electrolyte abnormalities in heart failure: pathophysiology and implications. Heart Fail Rev 2015;20(4):493-503 [FREE Full text] [doi: 10.1007/s10741-015-9482-y] [Medline: 25820346]

- 30. Silverberg D, Wexler D, Iaina A, Schwartz D. The role of anemia in the progression of congestive heart failure: Is there a place for erythropoietin and intravenous iron? Transfus Altern Transfus Med 2008;6(3):26-37. [doi: 10.1111/j.1778-428x.2005.tb00121.x]
- 31. Costache II, Cimpoeşu D, Petriş O, Petriş AO. Electrolyte disturbances in patients with chronic heart failure—clinical, evolutive and therapeutic implications. Rev Med Chir Soc Med Nat Iasi 2012;116(3):708-713. [Medline: 23272514]
- Miller A, Mandeville J. Predicting and measuring fluid responsiveness with echocardiography. Echo Res Pract 2016;3(2):G1-G12 [FREE Full text] [doi: 10.1530/ERP-16-0008] [Medline: 27249550]
- Petrie JR, Guzik TJ, Touyz RM. Diabetes, hypertension, and cardiovascular disease: clinical insights and vascular mechanisms. Can J Cardiol 2018;34(5):575-584 [FREE Full text] [doi: 10.1016/j.cjca.2017.12.005] [Medline: 29459239]
- 34. Cassano P, Fava M. Depression and public health: an overview. J Psychosom Res 2002;53(4):849-857. [doi: 10.1016/s0022-3999(02)00304-5] [Medline: 12377293]
- 35. Stookey JD, Barclay D, Arieff A, Popkin BM. The altered fluid distribution in obesity may reflect plasma hypertonicity. Eur J Clin Nutr 2007;61(2):190-199. [doi: 10.1038/sj.ejcn.1602521] [Medline: 17021599]

Abbreviations

CHF: congestive heart failure **EHR:** electronic health record **OSDT:** optimal sparse decision tree

Edited by Z Yin; submitted 10.07.24; peer-reviewed by Y Li, M Madadi; comments to author 05.09.24; revised version received 18.10.24; accepted 16.12.24; published 29.01.25.

<u>Please cite as:</u> Jiang Y, Li Q, Huang YL, Zhang W Urgency Prediction for Medical Laboratory Tests Through Optimal Sparse Decision Tree: Case Study With Echocardiograms JMIR AI 2025;4:e64188 URL: <u>https://ai.jmir.org/2025/1/e64188</u> doi:10.2196/64188 PMID:<u>39879091</u>

©Yiqun Jiang, Qing Li, Yu-Li Huang, Wenli Zhang. Originally published in JMIR AI (https://ai.jmir.org), 29.01.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Advancing Privacy-Preserving Health Care Analytics and Implementation of the Personal Health Train: Federated Deep Learning Study

Ananya Choudhury^{1,2*}, MTech; Leroy Volmer^{1,2*}, MSc; Frank Martin³, MSc; Rianne Fijten^{1,2}, PhD; Leonard Wee^{1,2}, PhD; Andre Dekker^{1,2,4}, PhD; Johan van Soest^{1,2,4}, PhD

¹GROW Research Institute for Oncology and Reproduction, Maastricht University Medical Center+, Maastricht, Netherlands

²Clinical Data Science, Maastricht University, Maastricht, Netherlands

³Netherlands Comprehensive Cancer Organization (IKNL), Eindhoven, Netherlands

⁴Brightlands Institute for Smart Society (BISS), Faculty of Science and Engineering (FSE), Maastricht University, Heerlen, Netherlands

*these authors contributed equally

Corresponding Author:

Ananya Choudhury, MTech GROW Research Institute for Oncology and Reproduction Maastricht University Medical Center+ Paul Henri Spakalaan 1 Maastricht, 6229EN Netherlands Phone: 31 0686008485 Email: <u>ananya.aus@gmail.com</u>

Abstract

Background: The rapid advancement of deep learning in health care presents significant opportunities for automating complex medical tasks and improving clinical workflows. However, widespread adoption is impeded by data privacy concerns and the necessity for large, diverse datasets across multiple institutions. Federated learning (FL) has emerged as a viable solution, enabling collaborative artificial intelligence model development without sharing individual patient data. To effectively implement FL in health care, robust and secure infrastructures are essential. Developing such federated deep learning frameworks is crucial to harnessing the full potential of artificial intelligence while ensuring patient data privacy and regulatory compliance.

Objective: The objective is to introduce an innovative FL infrastructure called the Personal Health Train (PHT) that includes the procedural, technical, and governance components needed to implement FL on real-world health care data, including training deep learning neural networks. The study aims to apply this federated deep learning infrastructure to the use case of gross tumor volume segmentation on chest computed tomography images of patients with lung cancer and present the results from a proof-of-concept experiment.

Methods: The PHT framework addresses the challenges of data privacy when sharing data, by keeping data close to the source and instead bringing the analysis to the data. Technologically, PHT requires 3 interdependent components: "tracks" (protected communication channels), "trains" (containerized software apps), and "stations" (institutional data repositories), which are supported by the open source "Vantage6" software. The study applies this federated deep learning infrastructure to the use case of gross tumor volume segmentation on chest computed tomography images of patients with lung cancer, with the introduction of an additional component called the secure aggregation server, where the model averaging is done in a trusted and inaccessible environment.

Results: We demonstrated the feasibility of executing deep learning algorithms in a federated manner using PHT and presented the results from a proof-of-concept study. The infrastructure linked 12 hospitals across 8 nations, covering 4 continents, demonstrating the scalability and global reach of the proposed approach. During the execution and training of the deep learning algorithm, no data were shared outside the hospital.

Conclusions: The findings of the proof-of-concept study, as well as the implications and limitations of the infrastructure and the results, are discussed. The application of federated deep learning to unstructured medical imaging data, facilitated by the PHT framework and Vantage6 platform, represents a significant advancement in the field. The proposed infrastructure addresses the

challenges of data privacy and enables collaborative model development, paving the way for the widespread adoption of deep learning–based tools in the medical domain and beyond. The introduction of the secure aggregation server implied that data leakage problems in FL can be prevented by careful design decisions of the infrastructure.

Trial Registration: ClinicalTrials.gov NCT05775068; https://clinicaltrials.gov/study/NCT05775068

(JMIR AI 2025;4:e60847) doi:10.2196/60847

KEYWORDS

gross tumor volume segmentation; federated learning infrastructure; privacy-preserving technology; cancer; deep learning; artificial intelligence; lung cancer; oncology; radiotherapy; imaging; data protection; data privacy

Introduction

Federated learning (FL) allows the collaborative development of artificial intelligence models using large datasets, without the need to share individual patient-level data [1-4]. In FL, partial models trained on separate datasets are shared, but not the data itself, hence a global model is derived from the collective set of partial models. This study introduces an innovative FL framework known as the Personal Health Train (PHT) that includes the procedural, technical, and governance components needed to implement FL on real-world health care data, including the training of deep learning neural networks [5]. The PHT infrastructure is supported by a free and open-source infrastructure known as "priVAcy preserviNg federaTed leArninG infrastructurE for Secure Insight eXchange," that is, Vantage6 [6]. We will describe in detail an architecture for training a deep learning model in a federated way with 12 institutional partners located in different parts of the world.

Sharing patient data between health care institutions is tightly regulated due to concerns about patient confidentiality and the potential for misuse of data. Data protection laws—including the European Union's General Data Protection Regulations; Health Insurance Portability and Accountability Act of 1996 (HIPAA) in the United States; and similar regulations in China, India, Brazil, and many other countries—place strict conditions on the sharing and secondary use of patient data [7]. Incompatibilities between laws and variations in the interpretation of such laws lead to strong reluctance about sharing data across organizational and jurisdictional boundaries [8-10].

To address the challenges of data privacy, a range of approaches have been published in the literature. Differential privacy, homomorphic encryption, and FL comprise a family of applications known as "privacy enhancing technologies" [11-13]. The common goal of privacy-enhancing technologies is to unlock positively impactful societal, economic, and clinical knowledge by analyzing data en masse, while obscuring the identity of study subjects that make up the dataset. Academic institutions are more frequently setting up controlled workspaces (eg, secure research environments [SREs]), where multiple researchers can collaborate on data analysis within a common cloud computing environment, but without allowing access to the data from outside the SRE desktop; however, this assumes that all the data needed have been transferred into the SRE in the first place [14,15]. Similarly, the National Institutes of Health has set up an "Imaging Data Commons" to provide

```
https://ai.jmir.org/2025/1/e60847
```

secure access to a large collection of publicly available cancer imaging data colocated with analysis tools and resources [16]. Other researchers have shown that blockchain encryption technology can be used to securely store and share sensitive medical data [17]. Blockchain ensures data integrity by maintaining an audit trail of every transaction, while zero trust principles make sure the medical data are encrypted and only authenticated users and devices interact with the network [18].

From a procedural point of view, the PHT manifesto for FL rules out the sharing of individual patient-level data between institutions, no matter if the patient data have been deidentified or encrypted [19]. The privacy-by-design principle here may be referred to as "safety in numbers," that is, any single individual's data values are obscured, by computing either the descriptive statistics or the partial model, over multiple patients. PHT allows sufficiently adaptable methods of model training, such as iterative numerical approximation (eg, bisection) or federated averaging (FedAvg [20]), and does not mandatorily require model gradients or model residuals, which are well-known avenues of privacy attacks [21-24]. Governance is essential with regards to compliance with privacy legislation and division of intellectual property between collaboration partners. A consortium agreement template for PHT has been made openly accessible [25], which is based on our current consortium ARGOS (artificial intelligence for gross tumor volume segmentation) [26]. Technologically, PHT requires 3 interdependent components to be installed-"tracks" are protected telecommunications channels that connect partner institutions, "trains" are Docker containerized software apps that execute a statistical analysis that all partners have agreed upon, and "stations" are the institutional data repositories that hold the patient data [23]. It is this technological infrastructure-the tracks, trains, and stations-that is supported by the aforementioned Vantage6 software, for which detailed stand-alone documentation exists [27].

The paper proposes a federated deep learning infrastructure based on the PHT manifesto [19], which provides a governance and ethical, legal, and social implications framework for conducting FL studies across geographically diverse data providers. The research aims to showcase a custom FL infrastructure using the open-source Vantage6 platform, detailing its technological foundations and implementation specifics. The paper emphasizes the significance of the implemented custom federation strategy, which maintains a strict separation between intermediate models from both internal and external user access. This approach is crucial for safeguarding the security and privacy of sensitive patient data,

XSL•FO RenderX

as it prevents potential reverse engineering of intermediate results that could compromise confidentiality. This aggregation strategy is particularly important in the case of deep learning–based studies where multiple iterations of models or gradients are necessary to derive an optimal global model.

To demonstrate the infrastructure's robustness and practical applicability, the study presents a proof-of-concept involving the development of a federated deep learning algorithm based on 2D convolutional neural network (CNN) architecture [28]. This algorithm was implemented to automatically segment gross tumor volume (GTV) from lung computed tomography (CT)

images of patients with lung cancer. Figure 1 [29] demonstrates a manual segmentation and deep learning–based segmentation of a tumor in the chest CT image of a patient. The subsequent sections provide a comprehensive account of the precise technical specifications of the infrastructure that links 12 hospitals across 8 nations, covering 5 continents. The algorithm developed learns from the distributed datasets and deploys it using the infrastructure. However, it is important to mention that the choice of the use case is only exemplary in nature, and the infrastructure is equipped to train any kind of deep learning architecture for relevant clinical use cases.

Figure 1. Illustrative result on a hold-out validation slice; the main bulk of the gross tumor volume as determined by the oncologist (middle) has been correctly delineated by the deep learning algorithm (right), but a small tumor mass adjacent and to the lower right of the main gross tumor volume mass has been missed (reproduced from Figure 6 of Chapter 4 of the thesis by Patil [29], which is published under the Taverne License [Article 25fa of the Dutch Copyright Act]).



The research used a deep learning architecture because in recent times the application of deep learning in health care has led to impressive results, specifically in the areas of natural language processing and computer vision (medical image analysis), with the promise for more efficient diagnostics and better predictions of treatment outcomes in future [30-35]. However, for robust generalizability, and to earn clinicians' acceptance, it is essential that artificial intelligence apps are trained on massive volumes of diverse and demographically representative health care data across multiple institutions. Given the barriers to data sharing, this is clearly an area where FL can play a vital role. Many studies have been published that present FL on medical data including federated deep learning [36-40]. However, only a limited number of studies have documented the use of dedicated frameworks and infrastructures in a transparent manner. The adoption of a custom federation strategy or absence of explicit reporting on the used infrastructure is observed in most of the studies. Table 1 summarizes the small number of FL studies that have been published in connection with deep learning investigations related to medical image segmentations to date.

The paper primarily focuses on demonstrating the training and aggregation mechanism of a deep learning architecture within a FL framework. It deliberately avoids delving into the optimization of model performance or clinical accuracy, as these aspects fall outside the paper's scope. Instead of emphasizing the selection of an optimal CNN architecture or aggregation strategy [39], the research concentrates on elucidating the functionality of the FL infrastructure. Existing literature has shown that FL models can achieve performance comparable to centrally trained models [38,41,45-47]. This supports the assumption that, given identical datasets and CNN architectures, a model trained using FL would likely yield similar results to one trained through centralized methods. The paper operates under this premise, prioritizing the explanation of the FL process over demonstrating performance parity with centralized training approaches.

The study highlights 3 key points as follows:

- FL is particularly well suited for deep learning applications, which typically require vast amounts of data. This makes it an ideal showcase for the federated approach.
- When implementing federated deep learning, it is crucial to have a robust infrastructure and use a customized, secure aggregation strategy. These elements are essential for safeguarding the privacy of sensitive patient information.
- FL in real-world medical data is not just a technological challenge; it requires a comprehensive strategy that addresses ethical, legal, governance, and organizational aspects, as highlighted by the PHT manifesto.

Table 1. Existing studies from the literature focusing on federated deep learning on medical images.

Infrastructure and clinical use case	Data type	Scale
NVIDIA FLARE/CLARA		·
Prostate segmentation of T2-weighted MRI ^a [41]	DICOM MRI	3 centers
COVID-19 pneumonia detection [42]	Chest CT ^b	7 centers
Tensorflow federated		
COVID-19 prediction from chest CT images [43]	Chest CT	3 datasets
OpenFL		
Glioblastoma tumor boundary detection [44]	Brain MRI	71 centers

^aMRI: magnetic resonance imaging.

^bCT: computed tomography.

The findings of the proof-of-concept study, as well as the implications and limitations of the infrastructure and the results, are discussed. The subsequent section of the paper is structured as follows: the *Methods* section describes the approach taken, followed by the *Results*, which detail the implementation of the infrastructure and a proof-of-concept execution. Finally, the paper concludes with a *Discussion* section.

Methods

Overview

When conducting a federated deep learning study, it is crucial to consider several key perspectives, which include both technical as well as organizational and legal aspects. These key factors have been instrumental in designing the infrastructure architecture used for training the deep learning algorithm. In this section, we discuss the technical details while adhering to an Ethics-Legal-Social Impact framework as laid down by the PHT manifesto. The technical design decisions are based on the following assumptions:

Data Landscape

Understanding the data landscape is crucial in designing and deploying FL algorithms. The technological approaches for handling horizontally partitioned data, where each institution contains nonoverlapping human subjects but the domain of the data (eg, CT images of lung cancer) is the same across different institutions, can differ significantly from those used for vertically partitioned data, where each institution contains the same human subjects but the domain of the data do not overlap (eg, CT scans in one, but socioeconomic metrics in another). Additionally, unstructured data, such as medical images, requires different algorithms and preprocessing techniques compared with structured data. In this paper, the architecture will only focus on CT scans and horizontally partitioned patient data.

Data Preprocessing

In a horizontally partitioned FL setting, the key preprocessing steps can be standardized and sent to all partner institutions. However, the workflow needs to handle differences in patients, scan settings, and orientations. Anonymization, quality improvements, and DICOM standardization ensure homogeneity and high quality across hospitals. These offline preprocessing steps, applied consistently to the horizontally partitioned data, enabled using the same model across institutions, crucial for the FL study's success.

Network Topology of the FL Infrastructure

The network topology choice for implementing FL can vary from client-server, peer-to-peer, tree-based hierarchical, or hybrid topologies. While peer-to-peer architecture is more cost-effective and offers a high capacity, it has the disadvantages of a lack of security and privacy constraints and a complex troubleshooting process in the event of a failure. The choice of network topology for this study is based on a client-server architecture, offering a single point of control in the form of the central server.

Choice of Model Aggregation Site

For a client-server architecture, the model aggregation can occur either in one of the data providers' machines, the central server, or in a dedicated aggregation server. For this implementation, we opted to use a dedicated aggregation server. The details and benefits of the implementation are discussed in the next section.

Training Strategy

The communication mechanism for transferring weights can be either synchronous, asynchronous, or semisynchronous, and weights can be consolidated using ensemble learning, FedAvg, split learning, weight transfer, or swarm learning. The strategy used for this study is based on a synchronous mechanism using the FedAvg algorithm. This gives a simple approach, where the averaging algorithm waits for all the data centers to transfer the locally trained model before initiating the averaging.

Based on the assumption, Figure 2 depicts the overall architecture of the federated deep learning study presented in the paper. The next section describes the FL Infrastructure in detail.



Figure 2. Overall architecture of ARGOS (artificial intelligence for gross tumor volume segmentation) federated deep learning architecture adapted from Vantage6. The figure depicts a researcher connected to the central server, a secure aggregation server, trains carrying models, connected data stations, and the communicating tracks.



The ARGOS Federated Deep Learning Infrastructure

Overview

In accordance with the PHT principles, the ARGOS infrastructure is comprised of 3 primary categories of components, labeled as the data stations, the trains, and the track. Furthermore, the architectural framework encompasses various roles that map to the level of permissions and access, specifically a track provider, the data providers, and the researcher. The infrastructure implementation can be further categorized into 3 important components: a central coordination server, a secure aggregation server (SAS), and the nodes located at each "data station." In the following sections, we attempt to describe each of these components and the respective stakeholders responsible for maintaining them.

Central Coordinating Server

The central coordination server is located at the highest hierarchical level and serves as an intermediary for message exchange among all other components. The components of the system, including the users, data stations, and SAS, are registered entities that possess well-defined authentication mechanisms within the central server. It is noteworthy that the central acts as a coordinator rather than a computational engine. Its primary function is to store task-specific metadata relevant to the task initiated for training the deep learning algorithm. In the original Vantage6 infrastructure, the central server also stores the intermediate results. In the ARGOS infrastructure, the central server is designed to not store any intermediate results but only the global aggregated model at the end of the entire training process.

Secure Aggregation Server

The SAS refers to a specialized station that contains no data and functions as a consolidator of locally trained models. The aggregator node is specifically designed to possess a Representational State Transfer (REST)–application programming interface (API) termed as the API Forwarder. The API Forwarder is responsible for managing the requests received from the data stations and subsequently routing them to the corresponding active Docker container, running the aggregation algorithm.

To prevent any malicious or unauthorized communication with the aggregator node, each data station is equipped with a JSON Web Token (JWT) that is unique for each iteration. The API Forwarder only accepts communications that are accompanied by a valid JWT. The implementation of this functionality guarantees the protection of infrastructure users and effectively mitigates the risk of unauthorized access to SAS. Figure 3 shows the architecture and execution mechanism for the SAS.



Choudhury et al

Figure 3. Architecture of the secure aggregation server, showing incoming and outgoing requests from the data station nodes. The upload and download folders are temporary locations used within the running Docker container to store the local and averaged models through disk read or write operations. The API forwarder, running at port 5050 and embedded within the Vantage6 infrastructure, forwards the incoming requests from the data station nodes to the algorithm API running at local port 7000 within the Docker container through HTTP requests. The SAS is hosted behind the firewall of a proxy server, which allows only hypertext transfer protocol secure (HTTPS) communication from the participating nodes. API: application programming interface; FedAvg: federated averaging; JWT: JSON Web Token.



Data Stations

Data stations are devices located within the confines of each hospital's jurisdiction that are not reachable or accessible from external sources other than Vantage6. The data stations communicate with the central server through a pull mechanism. Furthermore, the data stations not only serve as hosts for the infrastructure node but also offer the essential computational resources required for training the deep learning network. The infrastructure node is the software component installed in the data stations that orchestrates the local execution of the model and its communication with the central server and the SAS. Each data station is equipped with at least 1 graphics processing unit (GPU), which enables the execution of CNNs. Preprocessing of the raw CT images was executed locally, using automated preprocessing scripts packaged as Docker containers, and the preprocessed CT images are stored within a file system volume in each station. The CNN Docker is designed and allowed to access the preprocessed images during training. The primary function of the data station is to receive instructions from both the SAS and the central server, perform the computations needed for training the CNN algorithm, and subsequently transmit the model weights back to the respective sources. Figure 4 depicts the architectural layout of the data station and node component of the infrastructure.

Figure 4. Architecture of the data station node component. The node runs the CNN algorithm to learn from the local data. The node further sends and receives model weights from the secure aggregation server. The train and validation folders are persistent locations within the data stations, storing the preprocessed NIFTI images. At the end of each training cycle, the intermediate averaged model is first evaluated on the validation sample. CNN: convolutional neural network; HTTPS: hypertext transfer protocol secure; NIFTI: neuroimaging informatics technology initiative.



Train

The "train" in the form of a Docker image encompasses several components bundled together: an untrained U-Net [48,49], a type of CNN architecture designed for image segmentation tasks for training on local data; the aggregation algorithm used for consolidating the models; and a secondary Python Flask API known as the Algorithm API for facilitating the communication of these models. The Algorithm API is designed to cater to requests from the API Forwarder and is built within the algorithm container. Two levels of API ensured that the node could handle multiple requests and divert to appropriate Docker containers. Furthermore, the first level of API also helps in restricting malicious requests by checking the JWT token signature, so that the models within the master Docker container are protected. Each data station is responsible for training and transmitting the CNN model to the aggregator server. This suggests that the aggregation algorithm exhibits a waiting period during which it ensures that all data stations have effectively transmitted their models to the server before proceeding to the next iterations. The process is executed in an iterative manner until convergence is achieved or the specified number of iterations is attained.

Tracks and Track Provider

The various infrastructure components establish coordination among themselves through the use of secure communication channels commonly referred to as the "tracks." The communication channels are enabled with end-to-end encryption. The responsibility for the maintenance of the infrastructure, including the hosting of the central coordinating server and the specialized SAS, lies with the track provider. The track provider is additionally accountable for the maintenance of the "tracks" and aids the data providers in establishing the local segment of the infrastructure known as the "nodes."

Data Provider

Data providers refer to hospitals and health care organizations that are responsible for curating the pertinent datasets used for training the deep learning network. The responsibility of hosting the data stations within their respective local jurisdiction lies with the data provider. They exercise authority over the data as well as the infrastructure component called the node.

Researcher

The researcher is responsible for activating the deep learning algorithm and engaging in the authentication process with the central coordinating server using a registered username and password. This allows the researcher to establish their identity and gain secure access to the system, with their communication safeguarded through end-to-end encryption. The researcher can then assign tasks to individual nodes, monitor progress, and terminate tasks in the event of failure. Importantly, the researcher's methodology is designed to keep the intermediate outcomes of the iterative deep learning training process inaccessible, ensuring that the ultimate global model can only be obtained upon completion of all training iterations, thereby mitigating the risk of unauthorized access by malicious researchers to the intermediate models and providing a security mechanism against insider attacks.

Training Process

Each of the components described above works in a coordinated manner to accomplish the convergence of the deep learning algorithm. The training process begins with the researcher authenticating with the central server. Upon successful authentication, the researcher specifies the task details, including a prebuilt Docker image, input parameters, number of iterations, and the identity of the SAS. The task is then submitted to the central server, which forwards it to the connected nodes. The SAS is the first to receive the task request. It downloads the specified Docker image from the registry and initiates the master algorithm. The master algorithm orchestrates the training at each data station node through the central server. The central server then forwards a subtask request to all the data stations. Like the SAS, the data nodes download the same Docker image and initiate the node part of the algorithm. The node algorithm runs the learning process on local data for the specified number of epochs. After each training cycle, the node algorithm sends the local model weights to the SAS.

The SAS verifies the JWT signature of each received model and forwards the request to the Algorithm API. The Algorithm API extracts the weight and metadata information of the models. Once the SAS receives all the required locally trained models for that cycle, it initiates the FedAvg algorithm to consolidate the models and create an intermediate averaged model, which is stored locally. This completes the first iteration of the training cycle. For the second and subsequent iterations, the data stations request the SAS to send the intermediate averaged model weights from the previous iteration. The SAS validates these requests and sends the model weights to the data stations, which then use them for further training on their local data. This cycle of training and averaging continues until the model converges or the desired number of iterations is reached.

At the end of the training process, the SAS sends a notification to the researcher indicating the successful completion of the task. The researcher can then download the final global model from the server. It is important to note that during the training iterations, the researcher or other users of the infrastructure do not have access to the intermediate averaged models generated by the SAS. This design choice prevents the possibility of insider attacks and data leakage, as users cannot regenerate patterns from the training data using the intermediate models. Figure 5 shows the diagrammatic representation of the training process spread across the infrastructure components.



Figure 5. Process illustration of federated deep learning training. All entities, including the researcher, the central aggregation server, and the data stations, first authenticate with the central server. The researcher creates a task description and submits the task to the central server, which then forwards the request to the secure aggregation node to start the master task. The master task then sends a request to all data stations to download the algorithm Docker image and start training on the local data. Researchers can monitor the algorithm's execution status on the central server using the "check status" function, which reports whether each iteration is completed or aborted as processed by the secure aggregation server and data stations. At the end of each local training, the data stations send the models to the API forwarder of the secure aggregation node by authenticating against a valid JWT token. The JWT token ensures that no unauthorized data station is able to send or receive models from the secure aggregation server. API: application programming interface; CNN: convolutional neural network; JWT: JSON Web Token.



https://ai.jmir.org/2025/1/e60847

XSL•FO RenderX

Code Availability

The federated deep learning infrastructure and the algorithm used in this research are open source and publicly available. The codebase, encompassing the components of the infrastructure, the algorithm, and wrappers for running it in the infrastructure and the researcher notebooks, are all available and deposited on GitHub, a public repository platform, under the Apache 2.0 license. This open access allows the research community to scrutinize and leverage our implementation for further development in the field of FL.

The Vantage6 (version 2.0.0) [27,50] open-source software was customized to cater to the specific requirements for running the deep learning algorithm. The central server (Vantage6 version 2.0.0) and the aggregator server were hosted by Medical Data

Works BV in 2 separate cloud machines (Microsoft Azure). At each participating center, the "node" component of the software was installed and setup either on a physical or cloud machine running Ubuntu (version 16.0) or above with an installation of Python, (version 3.7 or above; Python Software Foundation), Docker Desktop (personal edition), and NVIDIA CUDA GPU interface (version 11.0). The source code of the customized "node" [51] and setup instructions [52] are available on respective GitHub repositories. The federated deep learning algorithm was adapted to the infrastructure as Python scripts [53] and wrapped in a Docker container. Separately, the "researcher" notebooks [54] containing python scripts for connecting to the infrastructure and running the algorithms are also available on GitHub. Table 2 provides an outline of the resource requirement and computational cost of the experiment.

Table 2. Resource requirement and computational cost.

End points	Resource requirement		Average execution time (per iteration)	
	Software	Hardware		
Central server	 Ubuntu (version 16) and above Docker Desktop Python (3.7 or above) Vantage6 (version 2.0.0) 	 4 CPUsa 16 GB RAM 20 GB Disk Space 	N/A ^b	
Data station	 Ubuntu (version 16) and above Docker Desktop Python (3.7 or above) Vantage6 (version 2.0.0) CUDA GPU Interface (version 11.0) 	 4 CPUs 1 GPUc 16 GB RAM 40 GB disk space 	40 mins	
Secure aggregation server	 Ubuntu (version 16) and above Docker Desktop Python (3.7 or above) Vantage6 (version 2.0.0) 	 4 CPUs 16 GB RAM 40 GB disk space 	60 seconds	

^aCPU: central processing unit.

^bNot applicable.

^cGPU: graphics processing unit.

Ethical Considerations

The work was performed independently with the ethics board's approval from each participating institution. Approvals from each of the participating institutions including soft copies of approval have been submitted to the leading partner. The lead partner's institutional review board approval (MAASTRO Clinic, The Netherlands) is "W 20 11 00069" (approved on November 24, 2020). The authors attest that the work was conducted by the ethical standards of the responsible committee on human experimentation (institutional and national) and with the Helsinki Declaration of 1975.

Results

Overview

The study was carried out and concluded in 4 primary stages using an agile approach as follows: planning, design and

```
https://ai.jmir.org/2025/1/e60847
```

development, partner recruitment, and execution of federated deep learning. The planning phase of the study, which encompassed a meticulous evaluation and determination of the following inquiries, held equal significance to the description of the clinical issue and data requirements.

- What are the minimum resource requirements for each participating center?
- How to design a safe and robust infrastructure to effectively address the requirements of a federated deep learning study?
- How can a reliable and data-agnostic federated deep learning algorithm be designed?
- What are the operational and logistical challenges associated with conducting a large-scale federated deep learning study?

The second phase, that is, the design and development phase, primarily focused on the creation, testing, and customization of the Vantage6 infrastructure for studies specifically focused on deep learning. To meet the security demands of these

investigations, this study involved the development of the SAS, which was not originally included in the Vantage6 architecture. The CNN algorithm was packaged as a Docker container and made compatible with the Vantage6 infrastructure, allowing it to be easily deployed and used within the Vantage6 ecosystem. Prior to the deployment of the algorithm, it underwent testing using multiple test configurations consisting of data stations that were populated with public datasets.

The primary objective of the third phase entailed the recruitment of partners who displayed both interest and suitability from various global locations. The project consortium members became part of the project by obtaining the necessary institutional review board approvals and signing an infrastructure user agreement. This agreement enabled them to install the required infrastructure locally and carry out algorithmic execution. The inclusion criteria for patient data, as well as the technology used for data anonymization and preprocessing, were provided to each center. The team collaborated with each partner center to successfully implement the local component of the infrastructure.

The concluding stage of the study involved the simultaneous establishment of connections between all partner centers and the existing infrastructure. The algorithm was subsequently initiated by the researcher and the completion of the predetermined set of federated iterations was awaited across all centers.

Proof of Concept

The architectural strategy described above was implemented among ARGOS consortium partners on real-world lung cancer CT scans. For an initial "run-up" of the system, we deployed the abovementioned PHT system across 12 institutions, located in 8 countries and 4 continents. A list of members participating in the ARGOS consortium can be found on the study protocol [26]. In total, 2078 patients' data were accessible via the infrastructure for training (n=1606) and holdout validation (n=472). For this initial training experiment, the 12 centers were divided into 2 groups. The first, referred to as group A, comprised 7 collaborators, and we were able to reach a total of 64 iterations of model training each with 10,000 steps per iteration. Likewise, group B comprising 6 hospitals was able to train the deep learning model for 26 iterations. It was observed that no significant improvement of the model was observed for both groups after 26th iteration. The results from the proof-of-concept study are shown in Figure 6.

While the training time for the models was similar at each center, how quickly they could be uploaded and downloaded depended heavily on the quality of the internet connection. This meant the entire process was significantly slowed down by the center with the slowest internet.



Figure 6. Plots showing the results from training the convolutional neural network on two groups as follows: group 1 (A, B, E, H, I, K, L) and group 2 (A, C, D, F, G, M). (A) Average Dice score per iteration of the model trained on group 1. (B) Average Dice score per iteration of the model trained on group 2. (C) Average training loss per iteration of the model trained on group 1. (D) Average training loss per iteration of the model trained on group 2.





Discussion

This study demonstrated the feasibility of a privacy-preserving federated deep learning infrastructure and presented a proof-of-concept study for GTV segmentation in patients with lung cancer. Using the PHT framework, the infrastructure linked 12 hospitals across 8 nations, showcasing its scalability and global applicability. Notably, throughout the process, no patient data were shared outside the participating institutions, addressing significant data privacy concerns. The introduction of a SAS further ensured that model averaging occurred in a secure environment, mitigating potential data leakage issues in FL.

One of the most used methodologies in recent years has been the use of FL for promoting research on privacy-sensitive data. To orchestrate FL on nonstructured data in the horizontal partitioning context, it is essential to develop specialized

```
https://ai.jmir.org/2025/1/e60847
```

RenderX

software for edge computation and technical infrastructures for cloud aggregation. These infrastructures enable federated machine learning (FML) responsibilities to be carried out in a secure and regulated manner. However, only a limited number of these studies have documented the background governance strategies and the ethical, legal, and social implications framework for conducting such studies.

The study presented a novel approach for executing large-scale federated deep learning on medical imaging data, integrating geographically dispersed real-world patient data from cross-continental hospital sites. The deep learning algorithm was designed to automatically delineate the GTV from chest CT images of patients with lung cancer who underwent radiotherapy treatment. The underlying FL infrastructure architecture was designed to securely perform deep learning training and was tested for vulnerabilities from known security

threats. This paper predominantly discussed the FL infrastructure architecture and presented a firsthand experience of conducting such studies. The preliminary training of the deep learning algorithm serves as the feasibility demonstration of the methodology, and further refinement is required to achieve acceptable clinical-grade accuracy and generalizability.

The study used an open-source and freely accessible technological stack to demonstrate the feasibility and applicability of federated deep learning. Vantage6, a Python-based FL infrastructure, is used to train and coordinate deep learning execution. TensorFlow and Flask, both open-source Python libraries, are used for the development of the algorithm, subsequently encapsulated within Docker services for containerization purposes. The communication channels between the hospital, central server, and the aggregation node have been secured using Hypertext Transfer Protocol Secure and Secure Hash Algorithm encryption. The hospital sites' computer systems were based on the Ubuntu operating system and equipped with at least 1 GPU to enhance computational capabilities. The participating centers had the flexibility to choose any CUDA-compatible GPU devices and determine the number of GPUs to use, enabling resource-constrained centers to contribute. However, a limitation exists in terms of computational time due to the synchronous training process being dependent on the slowest participant.

The infrastructure has been tested against known security attacks and as defined by the Open Worldwide Application Security Project top-ten categories [55]. It has been found that the Vantage6 app is impeccable against insecure design, software and data integrity failures, security logging and monitoring failures, and server-side request forgery and sufficiently secured against broken access control, cryptographic failures, injection, security misconfigurations, vulnerable and outdated components, and finally identification and authentication failures. Since the infrastructure is dependent on other underlying technologies like Docker and Flask-API, the security measures in these technologies also affect the overall security of the infrastructure. Additionally, the infrastructure is hosted behind proxy firewalls, adding to its overall security against external threats.

In this study, we implemented a SAS positioned between the data nodes (eg, hospitals and clinics) and the central server. The SAS plays a crucial role in strengthening the privacy and confidentiality of the learning process. The SAS acts as an intermediary that temporarily stores the local model updates from the participating data nodes, ensuring complete isolation from the central server, researchers, and any external intruders. The key benefits of using a dedicated SAS over a random aggregation mechanism in FL are as follows:

- Privacy protection of individual user data and model updates:
 - The secure aggregation protocol ensures that the central server only learns the aggregated sum of all user updates, without being able to access or infer the individual user's private data or model updates.
 - By isolating the intermediate updates, the secure aggregation process prevents external attackers from performing model inversion attacks.

```
https://ai.jmir.org/2025/1/e60847
```

XSL•FO

- Tolerance to user dropouts:
 - The SAS is designed to handle situations where some users fail to complete the execution. In the case of synchronous training, the server stores the latest successful model, enabling data nodes to pick up where they left off instead of restarting from scratch.
- Integrity of the aggregation process:
 - The secure aggregation protocol provides mechanisms to verify the integrity of the intermediate models by allowing only the known data nodes to send a model. This maintains the reliability and trustworthiness of the FL system.

FL offers 2 main approaches for model aggregation: sending gradients or weights [56,57]. In gradient sharing, data nodes update local models and transmit the gradients of their parameters for aggregation. Conversely, weight sharing involves sending the fully updated model weights directly to the server for aggregation. Sharing gradients have a higher risk of model inversion attacks. In the study presented here, the data nodes sent model weights instead of model gradients, thus preventing the "gradient leakage" problem. However, weight sharing is not failproof either [58], and the SAS plays a crucial role again in preventing users—internal or external—from accessing the weights from the aggregator machine.

The deployment of the FL infrastructure and training of the deep learning algorithm presented unique challenges that needed to be catered to. Some of them are listed below:

- Heterogeneity across hospitals: Initially, it was not possible to confirm the technology environment at each site. This required significant work to overcome the obstacles connected with each center while deploying a functional infrastructure, good communication, and efficient algorithms.
- Inconsistent IT policies: Standardizing the setup across institutions was hindered by varying IT governance and network regulations in different health care systems across different countries.
- Clinical expertise gap: The predominance of medical personnel over IT specialists at participating hospitals necessitated extensive documentation to ensure clinician comprehension of the FL process.
- Network bottlenecks: Network configurations at participating sites significantly impacted training duration, often leading to delays in model convergence.

The study presented in the paper has identified several areas that require further investigation and improvement. While the findings are valuable, the infrastructure, algorithm, and processes still need to be made more secure, private, trustworthy, robust, and seamless [59]. For example, incorporating homomorphic encryption of the learned models will enhance privacy and provide model obfuscation against inversion attacks. Finally, to further enhance confidence and trust in federated artificial intelligence, it is crucial to conduct additional studies involving a larger number of participating centers and a thorough clinical evaluation of the models.

Acknowledgments

We would like to express our sincere appreciation and gratitude to Integraal Kankercentrum Nederland (IKNL), the Netherlands, for their invaluable contribution in providing us with the necessary infrastructure support. We express our gratitude to Medical Data Works, the Netherlands, for their role as the infrastructure service provider in hosting the central and secure aggregation server. We also express our gratitude to Varsha Gouthamchand and Sander Puts for their contribution to the successful execution of the experiments. In conclusion, we express our gratitude to the various data-providing organizations for their substantial support and collaboration throughout all stages of the project. AC, LV, RF, and LW acknowledge financial support from the Dutch Research Council (NWO) (TRAIN project, dossier 629.002.212) and the Hanarth Foundation.

Conflicts of Interest

Dr AD and JvS are both cofounders, shareholders, and directors of Medical Data Works B.V.

References

- 1. Sun C, Ippel L, Dekker A, Dumontier M, van Soest J. A systematic review on privacy-preserving distributed data mining. Data Sci 2021 Oct;4(2):121-150. [doi: 10.3233/DS-210036]
- Choudhury A, Sun, C, Dekker M, Dumontie J, van Soest. Privacy-preserving federated data analysis: data sharing, protection, bioethics in healthcare. In: El Naqa I, Murphy MJ, editors. Machine Deep Learning in Oncology. Cham, Switzerland: Springer International Publishing; 2022:135-172.
- Deist TM, Dankers FJ, Ojha P, Scott Marshall M, Janssen T, Faivre-Finn C, et al. Distributed learning on 20 000+ lung cancer patients the personal health train. Radiother Oncol 2020;144:189-200 [FREE Full text] [doi: 10.1016/j.radonc.2019.11.019] [Medline: 31911366]
- 4. Choudhury A, Theophanous S, Lønne PI, Samuel R, Guren MG, Berbee M, et al. Predicting outcomes in anal cancer patients using multi-centre data and distributed learning a proof-of-concept study. Radiother Oncol 2021;159:183-189 [FREE Full text] [doi: 10.1016/j.radonc.2021.03.013] [Medline: 33753156]
- 5. Beyan O, Choudhury A, van Soest J, Kohlbacher O, Zimmermann L, Stenzhorn H, et al. Distributed analytics on sensitive medical data: the personal healt train. Data Intell 2020;2(1-2):96-107 [FREE Full text] [doi: 10.1162/dint a 00032]
- Moncada-Torres A, Martin F, Sieswerda M, van Soest J, Geleijnse G. VANTAGE6: an open source privacy preserving federated learning infrastructure for secure insight exchange. AMIA Annu Symp Proc 2020;2020:870-877 [FREE Full text] [Medline: <u>33936462</u>]
- Becker R, Chokoshvili D, Comandé G, Dove ES, Hall A, Mitchell C, et al. Secondary use of personal health data: when is it "Further Processing" under the GDPR, and what are the implications for data controllers? Eur J Health Law 2022;30(2):129-157. [doi: 10.1163/15718093-bja10094]
- El Naqa I, Ruan D, Valdes G, Dekker A, McNutt T, Ge Y, et al. Machine learning and modeling: data, validation, communication challenges. Med Phys 2018;45(10):e834-e840 [FREE Full text] [doi: 10.1002/mp.12811] [Medline: 30144098]
- 9. van Stiphout R. How to share data and promote a rapid learning health medicine? In: Valentini HJ, Schmoll C, van de Velde JH, editors. Multidisciplinary Management of Rectal Cancer. Cham, Switzerland: Springer International Publishing; 2018:623-634.
- Kazmierska J, Hope A, Spezi E, Beddar S, Nailon WH, Osong B, et al. From multisource data to clinical decision aids in radiation oncology: the need for a clinical data science community. Radiother Oncol 2020;153:43-54 [FREE Full text] [doi: 10.1016/j.radonc.2020.09.054] [Medline: 33065188]
- 11. Fischer-Hübner S. Privacy-enhancing technologies. In: Liu T, Özsu MT, editors. Encyclopedia of Database Systems. Boston, MA: Springer; 2009:2142-2147.
- Coopamootoo KPL. Usage patterns of privacy-enhancing technologies. In: ACM Digital Library. 2020 Presented at: CCS '20: Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security; November 2, 2020; New York, NY URL: <u>https://dl.acm.org/doi/10.1145/3372297.3423347</u>
- 13. Emerging privacy-enhancing technologies. OECD. URL: <u>https://www.oecd.org/publications/</u> emerging-privacy-enhancing-technologies-bf121be4-en.htm [accessed 2025-04-25]
- Kavianpour S, Sutherland J, Mansouri-Benssassi E, Coull N, Jefferson E. Next-generation capabilities in trusted research environments: interview study. J Med Internet Res 2022;24(9):e33720 [FREE Full text] [doi: <u>10.2196/33720</u>] [Medline: <u>36125859</u>]
- 15. Design a secure research environment for regulated data. Microsoft. URL: <u>https://learn.microsoft.com/en-us/azure/architecture/ai-ml/architecture/secure-compute-for-research</u> [accessed 2024-04-25]
- 16. Imaging data commons. National Cancer Institute Cancer Research Data Commons. URL: <u>https://datacommons.cancer.gov/</u> repository/imaging-data-commons [accessed 2024-04-25]

- 17. Kotter E, Marti-Bonmati L, Brady AP, Desouza NM. ESR white paper: blockchain and medical imaging. Insights Imaging 2021;12(1):82 [FREE Full text] [doi: 10.1186/s13244-021-01029-y] [Medline: 34156562]
- Sultana M, Hossain A, Laila F, Taher KA, Islam MN. Towards developing a secure medical image sharing system based on zero trust principles and blockchain technology. BMC Med Inform Decis Mak 2020;20(1):256 [FREE Full text] [doi: 10.1186/s12911-020-01275-y] [Medline: 33028318]
- 19. Manifesto of the personal health train consortium. Data Driven Life Sciences. URL: <u>https://www.dtls.nl/wp-content/uploads/</u> 2017/12/PHT_Manifesto.pdf [accessed 2024-03-11]
- 20. McMahan E, Moore D, Ramage S, Hampson BA. Communication-efficient learning of deep networks from decentralized data. In: Proceedings of Machine Learning Research. 2017 Presented at: Proceedings of the 20th International Conference on Artificial Intelligence and Statistics; April 20-22, 2017; Fort Lauderdale, FL URL: <u>https://proceedings.mlr.press/v54/</u> mcmahan17a.html
- 21. Zhang C, Choudhury A, Shi Z, Zhu C, Bermejo I, Dekker A, et al. Feasibility of privacy-preserving federated deep learning on medical images. Int J Radiat Oncol Biol Phys 2020;108(3):e778. [doi: <u>10.1016/j.ijrobp.2020.07.234</u>]
- 22. Choudhury A, van Soest J, Nayak S, Dekker A. Personal health train on FHIR: a privacy preserving federated approach for analyzing FAIR data in healthcare. In: Bhattacharjee A, Kr. Borgohain S, Soni B, Verma G, Gao XZ, editors. Machine Learning, Image Processing, Network Security and Data Sciences. Singapore: Springer; 2020.
- 23. Gouthamchand V, Choudhury A, P Hoebers FJ, R Wesseling FW, Welch M, Kim S, et al. Making head and neck cancer clinical data findable-accessible-interoperable-reusable to support multi-institutional collaboration and federated learning,? BJR Artif Intell 2024;1(1).
- 24. Sun C, van Soest J, Koster A, Eussen SJ, Schram MT, Stehouwer CD, et al. Studying the association of diabetes and healthcare cost on distributed data from the maastricht study and statistics Netherlands using a privacy-preserving federated learning infrastructure. J Biomed Inform 2022;134:104194 [FREE Full text] [doi: 10.1016/j.jbi.2022.104194] [Medline: 36064113]
- 25. Railway governance. Medical Data Works. URL: https://www.medicaldataworks.nl/governance [accessed 2024-09-11]
- 26. Dekker A. ARtificial Intelligence for Gross Tumour vOlume Segmentation (ARGOS). National Library of Medicine. URL: https://clinicaltrials.gov/study/NCT05775068 [accessed 2024-01-11]
- 27. Overview: what is vantage6? Vantage6 documentation. URL: <u>https://docs.vantage6.ai/en/main/</u> [accessed 2024-04-11]
- 28. U-Net: convolutional networks for biomedical image segmentation. In: Medical Image Computing and Computer-Assisted Intervention MICCAI 2015. Cham, Switzerland: Springer; Nov 18, 2015.
- 29. Patil RB. Prognostic and prediction modelling with radiomics for non-small cell lung cancer. Maastricht University. 2020. URL: <u>https://cris.maastrichtuniversity.nl/en/publications/prognostic-and-prediction-modelling-with-radiomics-for-non-small-</u>[accessed 2020-10-06]
- 30. Tao Z, Lyu S. A survey on automatic delineation of radiotherapy target volume based on machine learning. Data Intell 2023;5(3):814-856. [doi: 10.1162/dint_a_00204]
- 31. Liu X, Li KW, Yang R, Geng LS. Review of deep learning based automatic segmentation for lung cancer radiotherapy. Front Oncol 2021;11:717039 [FREE Full text] [doi: 10.3389/fonc.2021.717039] [Medline: 34336704]
- 32. Ma Y, Mao J, Liu X, Dai Z, Zhang H, Zhang X, et al. Deep learning-based internal gross target volume definition in 4D CT images of lung cancer patients. Med Phys 2023;50(4):2303-2316. [doi: <u>10.1002/mp.16106</u>] [Medline: <u>36398404</u>]
- 33. Zhang F, Wang Q, Li H. Automatic segmentation of the gross target volume in non-small cell lung cancer using a modified version of ResNet. Technol Cancer Res Treat 2020;19:153303382094748. [doi: 10.1177/1533033820947484]
- Xie H, Chen Z, Deng J, Zhang J, Duan H, Li Q. Automatic segmentation of the gross target volume in radiotherapy for lung cancer using transresSEUnet 2.5D network. J Transl Med 2022;20(1):524 [FREE Full text] [doi: 10.1186/s12967-022-03732-w] [Medline: 36371220]
- 35. Raimondi D, Chizari H, Verplaetse N, Löscher BS, Franke A, Moreau Y. Genome interpretation in a federated learning context allows the multi-center exome-based risk prediction of Crohn's disease patients. Sci Rep 2023 Nov 09;13(1):19449 [FREE Full text] [doi: 10.1038/s41598-023-46887-2] [Medline: 37945674]
- Riedel P, von Schwerin R, Schaudt D, Hafner A, Späte C. ResNetFed: federated deep learning architecture for privacy-preserving pneumonia detection from COVID-19 chest radiographs. J Healthc Inform Res 2023;7(2):203-224 [FREE Full text] [doi: 10.1007/s41666-023-00132-7] [Medline: 37359194]
- 37. Nazir S, Kaleem M. Federated learning for medical image analysis with deep neural networks. Diagnostics (Basel) 2023;13(9):1532 [FREE Full text] [doi: 10.3390/diagnostics13091532] [Medline: 37174925]
- 38. Shiri I, Vafaei Sadr A, Akhavan A, Salimi Y, Sanaat A, Amini M, et al. Decentralized collaborative multi-institutional PET attenuation and scatter correction using federated deep learning. Eur J Nucl Med Mol Imaging 2023;50(4):1034-1050 [FREE Full text] [doi: 10.1007/s00259-022-06053-8] [Medline: 36508026]
- 39. Zhang M, Qu L, Singh P, Kalpathy-Cramer J, Rubin DL. SplitAVG: a heterogeneity-aware federated deep learning method for medical imaging. IEEE J Biomed Health Inform 2022;26(9):4635-4644. [doi: <u>10.1109/jbhi.2022.3185956</u>]
- 40. Shiri I, Vafaei Sadr A, Amini M, Salimi Y, Sanaat A, Akhavanallaf A, et al. Decentralized distributed multi-institutional PET image segmentation using a federated deep learning framework. Clin Nucl Med 2022;47(7):606-617. [doi: 10.1097/rlu.000000000004194]

- Sarma KV, Harmon S, Sanford T, Roth HR, Xu Z, Tetreault J, et al. Federated learning improves site performance in multicenter deep learning without data sharing. J Am Med Inform Assoc 2021;28(6):1259-1264 [FREE Full text] [doi: 10.1093/jamia/ocaa341] [Medline: 33537772]
- 42. Harmon SA, Sanford TH, Xu S, Turkbey EB, Roth H, Xu Z, et al. Artificial intelligence for the detection of COVID-19 pneumonia on chest CT using multinational datasets. Nat Commun 2020;11(1):4080 [FREE Full text] [doi: 10.1038/s41467-020-17971-2] [Medline: 32796848]
- 43. Durga R, Poovammal E. FLED-block: federated learning ensembled deep learning blockchain model for COVID-19 prediction. Front Public Health 2022;10:892499 [FREE Full text] [doi: 10.3389/fpubh.2022.892499]
- 44. Pati S, Baid U, Edwards B, Sheller M, Wang S, Reina GA, et al. Federated learning enables big data for rare cancer boundary detection. Nat Commun 2022;13(1):7346 [FREE Full text] [doi: 10.1038/s41467-022-33407-5] [Medline: 36470898]
- 45. Leroy V, Ananya C, Aiara LG, Andre D, Leonard W. Feasibility of training federated deep learning oropharyngeal primary tumor segmentation models without sharing gradient information. Research Square Preprint published online 25 July, 2024 [FREE Full text] [doi: 10.21203/rs.3.rs-4644605/v1]
- 46. Schmidt K, Bearce B, Chang K, Coombs L, Farahani K, Elbatel M, et al. Fair evaluation of federated learning algorithms for automated breast density classification: the results of the 2022 ACR-NCI-NVIDIA federated learning challenge. Med Image Anal 2024;95:103206. [doi: 10.1016/j.media.2024.103206] [Medline: 38776844]
- 47. Pati S, Kumar S, Varma A, Edwards B, Lu C, Qu L, et al. Privacy preservation for federated learning in health care. Patterns (N Y) 2024;5(7):100974 [FREE Full text] [doi: 10.1016/j.patter.2024.100974] [Medline: 39081567]
- 48. Oreiller V, Andrearczyk V, Jreige M, Boughdad S, Elhalawani H, Castelli J, et al. Head and neck tumor segmentation in PET/CT: the HECKTOR challenge. Med Image Anal 2022;77:102336 [FREE Full text] [doi: 10.1016/j.media.2021.102336] [Medline: 35016077]
- 49. Iantsen A, Jaouen V, Visvikis D, Hatt M. Squeeze-and-excitation normalization for brain tumor segmentation. In: Crimi A, Bakas S, editors. Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries. Cham, Switzerland: Springer International Publishing; 2021.
- 50. IKNL/vantage6: Docker CLI package for the vantage6 infrastructure. GitHub. URL: <u>https://github.com/IKNL/vantage6/</u> <u>tree/DEV3</u> [accessed 2024-05-01]
- 51. Martin F. Featured communities. Zenodo. URL: <u>https://doi.org/10.5281/zenodo.3686944</u> [accessed 2024-05-06]
- 52. MaastrichtU-CDS/argos-infrastructure. GitHub. URL: <u>https://github.com/MaastrichtU-CDS/argos-infrastructure</u> [accessed 2024-05-01]
- 53. MaastrichtU-CDS/projects_argos_argos-code-repo_full-algorithm. GitHub. URL: <u>https://github.com/MaastrichtU-CDS/</u> projects_argos_argos-code-repo_full-algorithm [accessed 2024-05-01]
- 54. MaastrichtU-CDS/projects_argos_argos-code-repo_researcher-notebooks. GitHub. URL: <u>https://github.com/</u> <u>MaastrichtU-CDS/projects_argos_argos-code-repo_researcher-notebooks</u> [accessed 2024-05-01]
- 55. OWASP top ten. OWASP Foundation. URL: <u>https://owasp.org/www-project-top-ten/</u> [accessed 2024-05-02]
- 56. Moshawrab M, Adda M, Bouzouane A, Ibrahim H, Raad A. Reviewing federated learning aggregation algorithms; strategies, contributions, limitations and future perspectives. Electronics 2023;12(10):2287. [doi: <u>10.3390/electronics12102287</u>]
- 57. Liu P, Xu X, Wang W. Threats, attacks and defenses to federated learning: issues, taxonomy and perspectives. Cybersecurity 2022;5(1). [doi: 10.1186/s42400-021-00105-6]
- Boenisch F, Dziedzic A, Schuster R, Shamsabadi S, Shumailov I, Papernot N. When the curious abandon honesty: federated learning is not private. 2023 Presented at: IEEE 8th European Symposium on Security and Privacy (EuroS&P); July 07, 2023; Delft, theNetherlands. [doi: 10.1109/eurosp57164.2023.00020]
- 59. Sebastian G, George A, Jackson G. Persuading patients using rhetoric to improve artificial intelligence adoption: experimental study. J Med Internet Res 2023;25:e41430 [FREE Full text] [doi: 10.2196/41430] [Medline: 36912869]

Abbreviations

API: application programming interface
ARGOS: artificial intelligence for gross tumor volume segmentation
CNN: convolutional neural network
CT: computed tomography
FedAvg: federated averaging
FL: federated learning
FML: federated machine learning
GPU: graphics processing unit
GTV: gross tumor volume
HIPAA: Health Insurance Portability and Accountability Act
JWT: JSON Web Token
PHT: Personal Health Train
REST: Representational State Transfer

https://ai.jmir.org/2025/1/e60847

SAS: secure aggregation server **SRE:** secure research environment

Edited by Y Huo; submitted 23.05.24; peer-reviewed by AT Tran, G Sebastian; comments to author 02.07.24; revised version received 01.10.24; accepted 17.10.24; published 06.02.25.

<u>Please cite as:</u> Choudhury A, Volmer L, Martin F, Fijten R, Wee L, Dekker A, Soest JV Advancing Privacy-Preserving Health Care Analytics and Implementation of the Personal Health Train: Federated Deep Learning Study JMIR AI 2025;4:e60847 URL: <u>https://ai.jmir.org/2025/1/e60847</u> doi:10.2196/60847 PMID:

©Ananya Choudhury, Leroy Volmer, Frank Martin, Rianne Fijten, Leonard Wee, Andre Dekker, Johan van Soest. Originally published in JMIR AI (https://ai.jmir.org), 06.02.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Original Paper

Predicting Satisfaction With Chat-Counseling at a 24/7 Chat Hotline for the Youth: Natural Language Processing Study

Silvan Hornstein¹, MSc; Ulrike Lueken^{1,2}, Prof Dr; Richard Wundrack³, PhD; Kevin Hilbert⁴, Prof Dr

¹Department of Psychology, Humboldt-Universität zu Berlin, Berlin, Germany

²German Center for Mental Health (DZPG), Partner site Berlin/Potsdam, Potsdam, Germany

³Krisenchat gGmbH, Berlin, Germany

⁴Department of Psychology, HMU Erfurt - Health and Medical University Erfurt, Erfurt, Germany

Corresponding Author:

Silvan Hornstein, MSc Department of Psychology Humboldt-Universität zu Berlin Wolfgang Köhler-Haus Rudower Ch 18 Berlin, 12489 Germany Phone: 49 15753685796 Email: silvan.hornstein@hu-berlin.de

Abstract

Background: Chat-based counseling services are popular for the low-threshold provision of mental health support to youth. In addition, they are particularly suitable for the utilization of natural language processing (NLP) for improved provision of care.

Objective: Consequently, this paper evaluates the feasibility of such a use case, namely, the NLP-based automated evaluation of satisfaction with the chat interaction. This preregistered approach could be used for evaluation and quality control procedures, as it is particularly relevant for those services.

Methods: The consultations of 2609 young chatters (around 140,000 messages) and corresponding feedback were used to train and evaluate classifiers to predict whether a chat was perceived as helpful or not. On the one hand, we trained a word vectorizer in combination with an extreme gradient boosting (XGBoost) classifier, applying cross-validation and extensive hyperparameter tuning. On the other hand, we trained several transformer-based models, comparing model types, preprocessing, and over- and undersampling techniques. For both model types, we selected the best-performing approach on the training set for a final performance evaluation on the 522 users in the final test set.

Results: The fine-tuned XGBoost classifier achieved an area under the receiver operating characteristic score of 0.69 (P<.001), as well as a Matthews correlation coefficient of 0.25 on the previously unseen test set. The selected Longformer-based model did not outperform this baseline, scoring 0.68 (P=.69). A Shapley additive explanations explainability approach suggested that help seekers rating a consultation as helpful commonly expressed their satisfaction already within the conversation. In contrast, the rejection of offered exercises predicted perceived unhelpfulness.

Conclusions: Chat conversations include relevant information regarding the perceived quality of an interaction that can be used by NLP-based prediction approaches. However, to determine if the moderate predictive performance translates into meaningful service improvements requires randomized trials. Further, our results highlight the relevance of contrasting pretrained models with simpler baselines to avoid the implementation of unnecessarily complex models.

Trial Registration: Open Science Framework SR4Q9; https://osf.io/sr4q9

(JMIR AI 2025;4:e63701) doi:10.2196/63701

KEYWORDS

digital mental health; mental illness; mental disorder; adolescence; chat counseling; machine learning; artificial intelligence; large language model; natural language processing; deep learning

```
https://ai.jmir.org/2025/1/e63701
```

Introduction

Most mental health disorders develop early in life [1,2], causing a massive burden on an individual [3], as well as societal, level [4]. This makes early intervention in youth highly relevant [5]. In sharp contrast to the need, accessing help has been described as challenging for young people [5-7]. Therefore, low-threshold services are needed to tackle the burden of mental illness [8].

One such form of intervention gaining popularity is chat-based counseling hotlines [9-11]. Smartphones and chat interactions play a crucial role in youth life [12,13]. The ability to access help within their native digital life reduces numerous health care barriers, making the services a common first access point of help for youth [14]. Indeed, heavy utilization and adoption of those services have been reported globally [14-16]. In addition, the first evidence supports the acceptability [14] and effectiveness [17] of 24/7 chat services.

Considering the increasingly established relevance of those hotlines, the implementation of technological innovation could be highly impactful for the timely and efficient provision of care to youth. Repeatedly, artificial intelligence (AI) has been framed as a key potential for improvements in mental health care [18,19], as well as within digital settings [20]. As AI depends on the availability of large and high-dimensional datasets, chat services seem a quite promising candidate for that. This has indeed been used for diverse natural language processing (NLP) approaches, the subbranch of AI dealing with language. For example, an NLP-based triaging system has been reported to be able to reduce waiting times for those in crisis at a chat hotline [21]. Data-driven decisions regarding further treatment paths have also been investigated by looking into the prediction of recurrent chatting [22] or premature departure from conversations [23]. As suicide risk is a common case at chat hotline services [24], other work focused on early detection and intervention in those situations. Here, several model structures and algorithmic approaches have been suggested [25,26].

This study intends to contribute to the development of NLP approaches within youth chat counseling hotlines. Specifically, the promising but underinvestigated use case of automated evaluation of service quality will be explored. A recent study linked asynchronous chat counseling interactions with reported outcomes and satisfaction of the chatters, using a large dataset of more than 150,000 clients and reporting promising effect sizes of multiple R's of around 0.45 [27]. Another past approach investigated the prediction of chat quality on a label of 675 transcripts of chat counseling sessions [28]. However, while we were not able to find a similar-minded approach within 24/7 hotline services, automated quality evaluation seems particularly relevant for those. Early experiences with help seeking have been linked with future help-seeking behavior in the past [29]. As often being the first contact with any kind of institutionalized help for youth [14], the satisfaction with this interaction is therefore arguably highly relevant for further help-seeking behavior. The reliable identification of those with negative experiences would allow a timely intervention by following up or referrals to other services. Second, the low threshold nature

of counseling hotlines makes evaluation more difficult, as it is hard to collect follow-up responses from young help seekers. For example, the aforementioned study of chat hotline effectiveness reported a response rate of 22% among the users [17]. There is also the risk of a bias toward those more satisfied being more likely to respond, which is seen as a common methodical problem in evaluation sciences [30,31]. The ability to estimate the satisfaction with the service out of the conversation data for those who did not respond to any follow-up surveys could therefore significantly improve the evaluation and monitoring of the service quality.

In light of the relevance of the automated evaluation of chat interactions at chat hotlines, as well as the interventions raising relevance for youth mental health care, this project uses a naturalistic sample of 2609 young chatters that were counseled by the German 24/7 hotline service krisenchat. Feedback regarding the perceived helpfulness of the chat is used to train classifiers on the anonymized consultation texts. Performance is evaluated on a previously unseen test set addressing the feasibility of the approach, hypothesizing that we can significantly predict the feedback response by the chatter. Additionally, we assume that applying a pretrained transformer-based model as the state-of-the-art NLP will allow us to outperform a simpler non-transformer-based approach.

Methods

Preregistration

This study was preregistered at Open Science Framework [32]. The preregistration was updated once, as we adapted the used statistical test for the algorithm comparison (see the *Final Evaluation* section under *Methods*) and corrected the questionnaire item used for the outcome variable. We used the checklist for reporting machine learning studies by Klement and El Emam [33], which can be found in Multimedia Appendix 1. Due to legal restrictions regarding the highly vulnerable sample of this study, we are unable to share the dataset. However, the code used for training the algorithm and predicting the helpfulness can be found on GitHub [34], as a starting point for future work.

Ethical Considerations

The data collected and used for this study were part of a larger research project that was ethically approved by the University of Leipzig (372/21-ek). Additionally, we submitted the proposed secondary data analysis to the ethics committee of the Humboldt-Universität zu Berlin. They confirmed that this analysis does not require additional approval. Before the use of this study, the data were subject to a multistep anonymization procedure. Specifically, personally identifying information was marked by counselors and deleted by the organization. Additionally, there also was an automatized method in place to delete names and locations that might have been missed by the counselors. Finally, a k-anonymity principle was applied, deleting all words that were not part of at least 5 different chats.

Setting and Intervention

The anonymized data used for this study were provided by krisenchat, a German 24/7 chat counseling service for people

XSL•FO

aged up to 25 years. At krisenchat, those contacting the service through WhatsApp are provided with chat counseling, either by volunteer or employed psychologists, psychotherapists, or social workers. A central aspect of the consultations is the provision of exercises and resources, for example, by sharing YouTube videos, blog posts, or providing them within the chat. However, counselors are also trained in providing emotional support as needed, as well as providing information about mental health care structures in Germany, such as access to psychotherapy or the youth office.

Sample

Data were accessed and shared by the organization on January 17, 2024. On this date, there were feedback questionnaires available for 4560 chatters. Those questionnaires were sent out as part of a larger research project on the service [14]. A total of 264 participants were either younger than 13 years or older than 25 years of age and therefore excluded. While the upper age limit resulted from the scope of the service, the lower age limit resulted from data privacy considerations. An additional 1631 of the chatters were in contact with the service in the last 4 months. A help seeker's inactivity for at least 4 months is an organizational requirement for assuming the consultation purpose has ended and the chat is deleted by anonymization. Accordingly, active chats were also excluded, leading to 2664 concluded conversations and the related feedback questionnaire, with feedback provided between July 22, 2022, and September 17, 2023. For those cases, all messages exchanged between help seekers and counselors within 72 hours before the response to the feedback questionnaire were included. We then excluded cases where conversations consisted of fewer than 10 messages. This led to additional exclusions and resulted in a final sample of 2609 chatters. Their consultations consisted of 141,404 messages, 82,335 by the help seekers and 59,052 by the counselors. Therefore, on average, there were 54 messages exchanged in the three days before the feedback response, 23 messages by the counselor and 31 messages by the help seeker.

Outcome Variable

The feedback questionnaire answered by the chatters included several questions regarding the chat interaction (see Multimedia Appendix 2 for the full questionnaire). For this study, we decided on the use of a single item asking for the helpfulness of the chat ("Did the chat help you?" in German: "Hat dir der Chat geholfen?"), as being the most direct assessment available of chat quality and success, as perceived by the young clients. While the item had four possible answers ("Yes," "Rather Yes," "Rather No," and "No"), we decided to dichotomize it into "Yes" or "No." Reasons for that were improved actionability (as most clinical decision-making is binary by nature, such as providing additional help—yes or no), as well as considering the high-class imbalance. Overall, 89% (n=2332) of the chatters rated the chat as helpful. Specifically, 61 chatters responded with "No," 216 chatters responded with "Rather No," 1138 chatters responded with "Rather Yes," and 1194 chatters responded with "Yes."

Algorithm Training

All decisions regarding algorithmic specifications were made on the 80% of the available data used as a training set. Specifically, we separated the newest 20% of the consultations (522 chats who submitted their feedback after May 27, 2023) as a test set, a commonly used approach to mimic the evaluation of a previously implemented model (eg, [35]).

For our non-transformer-based approach, we preprocessed the data by lowering all words, deleting stop words, and using a lemmanizer [36]. Afterward, a term frequency-inverse document frequency (TF-IDF) vectorizer was used for feature extraction. This vectorizer counts the occurrences of words and weights them based on their frequency across the whole sample. This algorithm was trained using a 5-times repeated 5-fold stratified cross-validation principle. Hyperparameters were tuned using Bayesian optimization maximizing the receiver operating characteristic (ROC) area under the curve (AUC) score for 250 iterations. While there has been some discussion about the applicability of this metric facing class imbalance (eg, [37]), we saw its appropriateness backed up by systematic comparisons [38] and analysis [39] on the issue. All hyperparameters optimized during this procedure are summarized in Table 1. Those also included, as suggested by a reviewer, the range of ngrams used by the vectorizer. Therefore, bigrams and trigrams of words of the messages were also usable as predictors. The used over- or undersampling method was also selected during this procedure, comparing oversampling, undersampling, and Synthetic Minority Oversampling Technique [40]. As a classifier, we applied and tuned an extreme gradient boosting (XGBoost) [41] classifier, as well as a logistic regression. The training pipeline can be found on GitHub.

 Table 1. Overview of shortlisted transformer-based models.

Model	Input length, n	Source
uklfr/gottbert-base	512	[42]
distilbert/distilbert-base-german-cased	512	[43]
LennartKeller/longformer-gottbert-base-8192-aw512	8192	[44]

We used hugging face for all transformer-based approaches [42]. We shortlisted GottBERT [43], as well as a German DistilBERT model [44], as language-specific models to be evaluated. However, we assumed that a significant share of our data would exceed those models' input length. Therefore, we also intended to evaluate a Longformer model [45]. This model

can process much longer input sequences at reasonable computational costs by applying a sparse attention mechanism (see Table 1 for the shortlisted models including links). We also intended to explore over- and undersampling, as well as class weights to tackle the class imbalance. To represent the chat structure appropriately to the algorithm, we introduced two new

special tokens to the models, named "[USER]" and "[CNSLR]." Those were added at the beginning of each message, presenting the conversation structure in a processable format to the models. For hyperparameter tuning, a grid search across the learning rate $(2\times10^{-5}, 3\times10^{-5}, \text{ and } 5\times10^{-5})$ and the batch size (1, 2, and 4) was performed for the preselected most promising model. The training and tuning were done at a stratified train-validation split (70:30 of the data used for algorithm training), as the repeated cross-validation principle applied for the TF-IDF approach was infeasible due to computational costs. Therefore, a train-validation-test split (56:24:20) was used as an evaluation principle, with the same data being kept aside as final test data for the nontransformer approach. All transformer-based models were trained on an NVIDIA GeForce RTX 3090 graphics processing unit with 24 GB video random access memory.

Final Evaluation

The 522 newest conversations with feedback were used as a test set. The distribution of the outcome did not differ significantly between the training and test data (t_{520} =-1.1; *P*=.30). We decided to predict the outcome with the best performing TF-IDF approach and the most promising transformer approach, as identified on the train set as described above. We then applied a permutation test [46] to evaluate the significance of both algorithms. Finally, we contrasted the achieved AUCs of the two approaches, applying a DeLong test [47], which has been suggested for this scenario [48]. We decided for this procedure above the 5×2 McNemar test [49] originally proposed in our preregistration. This reconsideration was mainly made due to the inability of the McNemar test to statistically compare AUC scores. The comparison of accuracies seemed disadvantageous to us, as focusing on the performance

for one specific threshold. In contrast, considering the different proposed use cases, we were more interested in a threshold-independent comparison of classifier performance. As a threshold-dependent metric, we reported the Matthews correlation coefficient (MCC), which is particularly helpful in cases of imbalanced classes [50]. We followed the suggestion in the literature to use a default threshold of 0.5 [51] for the calculation of a confusion matrix and the corresponding MCC score.

Explainability

We used Shapley additive explanation (SHAP) values [52] as an explainability framework. This game-theory-based approach is applicable for transformer models [53] and XGBoost classifier [54].

Results

Algorithm Training

For the TF-IDF-based approach, the best set of hyperparameters selected through the tuning approach led to a mean ROC AUC score of 0.70 (SD 0.02) across repeated cross-validation for the XGBoost classifier. For this, a minimum occurrence of the word stems for 20 different chatters and for five different counselors was selected as a hyperparameter for the vectorizers. Random oversampling was selected for handling class imbalance. Counselors word stems were only selected when occurring in 30% or less of the conversations, while chatters word stems were allowed in up to 90% of the conversations. In addition, trigrams and bigrams were included, as well as predictors (see Table 2 for all hyperparameters). This was slightly above the performance of logistic regression, scoring 0.66 for the best set of hyperparameters.



Table 2.	Overview	of tuned	hyperparame	ters (definitions	adapted from	[22]).
			2.1 1	· · · · · · · · · · · · · · · · · · ·		

Hyperparameters	Description	Value range	Selected parame- ter
max_df_chatter	Terms that appear in more chatter documents than the threshold value are ignored. The value represents the proportion of documents	0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0	0.9
min_df_chatter	Terms that appear in fewer chatter documents than the threshold value are ignored	1, 2, 5, 10, 25, 50, 75, 100, 150, 200	20
max_df_couns	Analogous to max_df_chatterfor counselor messages	0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0	0.3
min_df_couns	Analogous to min_df_chatter for counselor messages	1, 2, 5, 10, 25, 50, 75, 100, 150, 200	5
Sampling method	Method for handling imbalance	ROS ^a , RUS ^b , SMOTE ^c	Rando- mOver- Sampler
colsample_bytree	Subsample ratio of columns for growing trees	0.2, 0.4, 0.6, 0.7, 0.8, 0.9, 1.0	1.0
eta	Learning rate	0.005, 0.01, 0.05, 0.1, 0.2	0.1
gamma	Minimum loss reduction to make a further split on a leaf node	0, 0.25, 0.5, 1, 1.5, 2, 5, 10	1.5
max_depth	Maximum depth of a tree	2, 4, 6, 8, 10, 12, 14, 16	16
min_child_weight	Minimum sum of instance weight (Hessian) needed in a child	1, 5, 10, 20	10
subsample	Subsample ratio of the training instances prior to growing trees	0.2, 0.4, 0.6, 0.7, 0.8, 0.9, 1.0	0.9
use_idf	Whether to term frequencies should be reweighted by the inverse document frequencies	True, false	True
ngram_range	Length of word sequences used as predictors	(1,1), (1, 2), (1,3)	(1,3)

^aROS: random over sampler.

^bRUS: random under sampler.

^cSMOTE: Synthetic Minority Oversampling Technique.

For the transformer-based approach, we reached a ROC AUC of 0.58 for the DistilBERT and 0.59 for the GottBERT models, using class weights (9:1) and five epochs. Comparable performances were reached when random oversampling was used instead of the class weights. We expected the performance to be limited by strong truncation. Therefore, we explored the average length of the input sequence with DistilBERT as tokenizer. Data points in the train set contained on average 1889 (SD 873) tokens, showing that those models could just use a share of the available data on the chat conversations. However, with the longest conversation holding 8507 tokens, the Longformer model structure seemed capable of capturing nearly all information contained in our data. Indeed, using the Longformer model in combination with class weights (9:1), three epochs, a learning rate of 3e-5, and a batch size of one resulted in a significantly higher ROC AUC of 0.69. Neither other methods for handling class imbalance nor different epoch sizes lead to a further improved performance.

Final Evaluation

While the performance between the transformer and non-transformer-based approach was similar during training (0.69 vs 0.70), this comparison is limited by the differences in the used validation principle. However, the large previously unseen test set allowed us the comparison of the two best-of-class models in a final evaluation. Here, we reached an ROC AUC of 0.68 for the Longformer model and an ROC AUC of 0.69 for the TF-IDF-based approach, both significantly outperforming randomness in a permutation test (P<.001 for both). However, as expected, considering the similar performance, there was no significant difference between the two approaches (P=.69). The ROC curves are plotted in Figure 1, showing how threshold and model performance interacted.



Figure 1. ROC AUC curves comparing the two algorithms. AUC: area under the curve; ROC: receiver operating characteristic; XGB: extreme gradient boosting.



Consequently, we used the TF-IDF approach as the simpler algorithm for further insights, as well as the explainability approach. The average precision score here was 0.93 (SD 0.02) on the test set. The MCC score for the default threshold of 0.5 was 0.25 on the test set. The confusion matrix on this threshold

can be found in Figure 2. Here, a positive predictive value of 0.90 and a negative predictive value (NPP) of 0.50 were achieved, with "positive" being coded as helpful. The sensitivity was 0.98 and the specificity was 0.18.

Figure 2. Confusion matrix for the selected threshold for the TF-IDF algorithm. TF-IDF: term frequency-inverse document frequency.



Explainability

We applied SHAP values on the vectorizer-based approach. The most predictive word identified here was "no" by the chatters, being associated with a higher chance of an unhelpful perceived chat. Two other predictors of unhelpfulness were the word "bad" (original: "schlimm") by the counselor, as well as "nevertheless" (original: "trotzdem") by the chatter, and "further on" (original: "weiterhin") by the counselor. In addition, some bigrams were among the most predictive variables. For example, "shift end" (German: "Schicht endet"), indicating that a counselor had to end a conversation due to their shift being over, was associated with negative feedback. For an improved understanding of the context those words were used, we looked into chats using those and giving negative feedback afterward. While "no" was used in diverse settings, there was a notable number of cases where chatters denied the counselor's offering of further help such as an exercise. "Bad" was used on several occasions where chatters reported highly traumatic experiences

they had. Finally, "further on" was a phrase repeatedly used by counselors to announce the end of their shift and offer further support from a colleague afterward. There were also several words being predictive of perceived helpfulness. Several of those implied that a chatter expressed satisfaction with the interaction at the end of a chat. For example, the word stem "thanks" (original: "dank") was predictive of higher perceived helpfulness, as was "great" (original: "toll"). We also investigated those conversations that were predicted with the highest likelihood of being labeled as unhelpful afterward. Again, there were several cases included where chatters rejected suggested exercises by the counselor. In addition, in several conversations with a high risk of unhelpfulness, it was reported that mental health care is already received, such as regularly seeing a psychiatrist or being hospitalized in a clinic. As one of the core functions of chat hotlines is the redirection into care, it might be harder to make a satisfying offer to those. The 20 most predictive words as identified by the tree-based SHAP approach can be found in Figure 3.

Figure 3. The 20 most predictive word stems as identified by the SHAP approach for the TF-IDF algorithm. SHAP: Shapley additive explanations; TF-IDF: term frequency-inverse document frequency.



Discussion

Primary Findings

This project investigated the use of NLP techniques for an automated evaluation of the perceived helpfulness of chat-based counseling. We were able to reach a ROC AUC of 0.67 on the previously unseen test set for a transformer, as well as for a non-transformer-based approach. Our explainability part revealed several linguistic markers of perceived unhelpful chat consultations such as the written expression of thankfulness, or the extensive use of the word "no" for rejecting the different offers made by counselors.

The reached performance was moderate, though significant and in line with past work from the identical settings [22]. However, the feasibility of an AI use case always depends on the performance considering the proposed use case. The given study implied two potential uses of predicted helpfulness of the chats.

The first use case was the real-time identification of unsuccessful consultations, as perceived by the chatter. Due to the very harmful impact of such experiences, those predictions could be used for a tailored follow-up, for example, with details of different treatment options for those affected. In our example, we would have identified 30 of the 62 unhelpful rated conversations with the approach, though 79% of all identified cases would have been false negatives (with negative referring to perceived unhelpfulness).

An alternative approach would have been a much stricter threshold, letting us mark significantly less chats but with higher NPP. For example, on a threshold of 0.3, our NPP would have doubled. However, the consequences of wrongly identifying chatters as unsatisfied might be less relevant than missing those being unsatisfied in light of the possible negative consequences of further help seeking. Overall, whether one of those approaches could be valuable would depend on whether the benefits for those correctly identified are larger than the costs of providing the intervention based on the prediction. Finally, this is an empirical question that we cannot answer here sufficiently. This highlights the large need for randomized controlled trials for prediction studies, moving from feasibility to actually showing clinical benefits [55].

A second use case of the proposed algorithm lies less on the individual and more on a population-based level. As evaluation within naturalistic and low-threshold settings is commonly difficult, the developed algorithm could be applied to those who did not respond to feedback questionnaires. This application would allow a better-informed estimation of satisfaction with the service where just a minority provides active feedback. A reliable estimate of this core metric of the service would propose a huge value for organizational purposes. Without any alternative of estimating the satisfaction of those not providing feedback being available, the proposed algorithm already provides an improvement over the status quo as clearly performing above the chance level. However, particularly for systematic comparison of, for example, monthly satisfaction, the question arises whether the performance is sufficient for reliable inference. Here, simulation studies might help to better

understand the relation between performance and the reliability of algorithm-based evaluation.

Secondary Findings

Interestingly, there was no further gain in predictive capability by using the computational heavy and pretrained Longformer model. The failure of more complex NLP models to outperform simpler ones is not unique to the given setting and has been reported before [56-58]. However, based on the literature, we started the work on this paper with an opposing hypothesis. For example, a popular study [59] compared Bidirectional Encoder Representations from Transformer-based approaches with TF-IDF-based algorithms and reported a clearly better performance for the former. An in-depth look into the used methods provides several possible explanations for the diverging results. First, the cited study used a larger sample of 50,000 distinct cases, while using the much smaller Bidirectional Encoder Representations from Transformer base model. Therefore, the dataset size might have been insufficient to finetune such a sophisticated model. Second, the use case is different, while algorithmic performance is highly case specific. The cited study focuses on sentiment analysis. Arguably, the extraction from word vectors into higher-dimensional spaces like sentiment as done by transformer models is particularly relevant here. While our explainability approach revealed some sentiment-related predictors like words of thankfulness, overly sentiment seemed less central than it is for movie reviews as in the aforementioned study. Finally, it remains unclear how much the advantage of simpler models is used in comparative studies. For example, in our approach, we were able to perform extensive hyperparameter tuning using sophisticated cross-validation principles. The relevance of this to produce generalizable results, and therefore, realistic performance estimates is well established [60,61]. Such approaches are hard to reproduce at feasible computational costs for transformer-based models for a lot of ML practitioners in their day-to-day work. However, waiving those techniques also for the baseline is arguably biasing the comparison against them, as their better capability to be trained with extended cross-validation principles is a real benefit that might translate into predictive performance. Particularly, small predictive performance differences as reported regularly (eg, [25]) might disappear with decent hyperparameter tuning and cross-validation.

In conclusion, while the actual outperformance seems dependent on setting and data, the results of this study, as well as the aforementioned studies, highlight the relevance of benchmarking complex models with simpler ones. Otherwise, overly complex models might be implemented without benefits. There are numerous studies that apply interesting and promising algorithmic approaches but do not compare them with a simpler baseline at all (eg, [62-64]). However, we also argue that a fair comparison includes the utilization of hyperparameter tuning and cross-validation for computationally lighter models.

Limitations

There were limitations to the approach in this paper. First, while we predicted the helpfulness of a chat as perceived by chatters, this perception does not equal to actually being clinically beneficial. For example, in the aforementioned study by Imel

et al [27], the association between message content and satisfaction was much stronger than the association between content and symptom reduction. Therefore, future work could benefit from associating chat messages with clinically validated questionnaires as output. However, arguably changes in symptoms are difficult to measure in hotline settings, where a majority of chatters just contact the service once. Second, we were only able to train the algorithms on the data of those who responded to the feedback questionnaire. This might have introduced a bias, in case of systematic differences between those providing feedback and those who do not. Third, we focused on the application of the Longformer model in the transformer-based approach of this paper. Future work might also benefit from exploring task-specific adaptions of the used algorithms in detail. In addition, different methods of handling long text inputs such as BELT [65] might enable a better performance. Notably, there were no mental health-specific

smaller models available in German. Those exist for other languages and use cases [66]. Such models, for example, pretrained on youth mental health data in German, could provide further performance gains as well. Finally, while we used a test set for a final one-time evaluation, this test set still came from the same chat counseling service. However, the relevance of truly external test sets has been highlighted repeatedly as being relevant for more valid claims regarding the generalizability of a chosen approach (eg, [67]).

Conclusions

In summary, there is a predictive signal regarding the perceived service quality in the chat messages at a 24/7 chat hotline for youth. This opens interesting use cases in the quality control and evaluation efforts at those hotlines. Future work such as the randomized evaluation of interventions based on the predicted helpfulness is needed for moving toward real-world implementation.

Acknowledgments

We acknowledge support by the Open Access Publication Fund of Humboldt-Universität zu Berlin.

Authors' Contributions

SH developed the idea, analyzed the data, and wrote the first draft of the paper. All authors contributed to the development of the exact analysis to be performed. All authors reviewed and contributed to the final draft.

Conflicts of Interest

SH and RW are employed by krisenchat, the organization that provided the data for this study. SH is also employed by Elona Health, a provider of digital health applications for mental health in Germany. KH is a scientific advisor and received virtual stock options from Mental Tech GmbH, which develops an artificial intelligence–based chatbot providing mental health support.

Multimedia Appendix 1

Consolidated reporting guidelines for prognostic and diagnostic machine learning modeling studies. [DOCX File, 20 KB - ai_v4i1e63701_app1.docx]

Multimedia Appendix 2

Full questionnaire sent out to chatters, original (German) and English translation. [DOCX File , 16 KB - <u>ai v4i1e63701 app2.docx</u>]

References

- Kessler RC, Angermeyer M, Anthony JC, de Graaf R, Demyttenaere K, Gasquet I, et al. Lifetime prevalence and age-of-onset distributions of mental disorders in the World Health Organization's World Mental Health Survey initiative. World Psychiatry 2007;6(3):168-176 [FREE Full text] [Medline: <u>18188442</u>]
- de Girolamo G, Dagani J, Purcell R, Cocchi A, McGorry PD. Age of onset of mental disorders and use of mental health services: needs, opportunities and obstacles—CORRIGENDUM. Epidemiol Psychiatr Sci 2022;31:e46 [FREE Full text] [doi: 10.1017/S2045796022000282] [Medline: 35762753]
- 3. Vigo D, Thornicroft G, Atun R. Estimating the true global burden of mental illness. Lancet Psychiatry 2016;3(2):171-178. [doi: 10.1016/S2215-0366(15)00505-2] [Medline: 26851330]
- 4. Christensen MK, Lim CCW, Saha S, Plana-Ripoll O, Cannon D, Presley F, et al. The cost of mental disorders: a systematic review. Epidemiol Psychiatr Sci 2020;29:e161 [FREE Full text] [doi: 10.1017/S204579602000075X] [Medline: 32807256]
- 5. McGorry PD, Mei C. Early intervention in youth mental health: progress and future directions. Evidence Based Mental Health 2018;21(4):182-184 [FREE Full text] [doi: 10.1136/ebmental-2018-300060] [Medline: 30352884]
- Colizzi M, Lasalvia A, Ruggeri M. Prevention and early intervention in youth mental health: is it time for a multidisciplinary and trans-diagnostic model for care? Int J Mental Health Syst 2020;14:23 [FREE Full text] [doi: 10.1186/s13033-020-00356-9] [Medline: 32226481]

```
https://ai.jmir.org/2025/1/e63701
```

- 7. Catania LS, Hetrick SE, Newman LK, Purcell R. Prevention and early intervention for mental health problems in 0–25 year olds: are there evidence-based models of care? Adv Mental Health 2014;10(1):6-19. [doi: <u>10.5172/jamh.2011.10.1.6]</u>
- 8. McGorry PD, Mei C, Chanen A, Hodges C, Alvarez-Jimenez M, Killackey E. Designing and scaling up integrated youth mental health care. World Psychiatry 2022;21(1):61-76 [FREE Full text] [doi: 10.1002/wps.20938] [Medline: 35015367]
- Tibbs M, O'Reilly A, O'Reilly MD, Fitzgerald A. Online synchronous chat counselling for young people aged 12-25: a mixed methods systematic review protocol. BMJ Open 2022;12(4):e061084 [FREE Full text] [doi: 10.1136/bmjopen-2022-061084] [Medline: 35470202]
- 10. Ersahin Z, Hanley T. Using text-based synchronous chat to offer therapeutic support to students: a systematic review of the research literature. Health Educ J 2017;76(5):531-543. [doi: 10.1177/0017896917704675]
- Mathieu SL, Uddin R, Brady M, Batchelor S, Ross V, Spence SH, et al. Systematic review: the state of research into youth helplines. J Am Acad Child Adolesc Psychiatry 2021;60(10):1190-1233. [doi: <u>10.1016/j.jaac.2020.12.028</u>] [Medline: <u>33383161</u>]
- 12. Teens, social media and technology 2023. Pew Research Center. 2023. URL: <u>https://www.pewresearch.org/internet/2023/</u> 12/11/teens-social-media-and-technology-2023/ [accessed 2024-01-30]
- Hajok D. Der veränderte Medienumgang Jugendlicher. Tendenzen aus 20 Jahren JIM-Studie. The changing media usage of adolescents: trends from 20 years of the JIM study. Jugend Medien Schutz-Report 2018;41(6):4-6. [doi: 10.5771/0170-5067-2018-6-4]
- 14. Eckert M, Efe Z, Guenthner L, Baldofski S, Kuehne K, Wundrack R, et al. Acceptability and feasibility of a messenger-based psychological chat counselling service for children and young adults ("krisenchat"): a cross-sectional study. Internet Interventions 2022;27:100508 [FREE Full text] [doi: 10.1016/j.invent.2022.100508] [Medline: 35242589]
- 15. Thompson LK, Sugg MM, Runkle JR. Adolescents in crisis: a geographic exploration of help-seeking behavior using data from Crisis Text Line. Soc Sci Med 2018;215:69-79. [doi: <u>10.1016/j.socscimed.2018.08.025</u>] [Medline: <u>30216891</u>]
- Watling D, Batchelor S, Collyer B, Mathieu S, Ross V, Spence SH, et al. Help-seeking from a national youth helpline in Australia: an analysis of kids helpline contacts. Int J Environ Res Public Health 2021;18(11):6024 [FREE Full text] [doi: 10.3390/ijerph18116024] [Medline: 34205148]
- Gould MS, Pisani A, Gallo C, Ertefaie A, Harrington D, Kelberman C, et al. Crisis text-line interventions: evaluation of texters' perceptions of effectiveness. Suicide Life Threat Behav 2022;52(3):583-595 [FREE Full text] [doi: 10.1111/sltb.12873] [Medline: 35599358]
- Lee EE, Torous J, de Choudhury M, Depp CA, Graham SA, Kim HC, et al. Artificial intelligence for mental health care: clinical applications, barriers, facilitators, and artificial wisdom. Biol Psychiatry Cogn Neurosci Neuroimaging 2021;6(9):856-864 [FREE Full text] [doi: 10.1016/j.bpsc.2021.02.001] [Medline: 33571718]
- 19. Dwyer DB, Falkai P, Koutsouleris N. Machine learning approaches for clinical psychology and psychiatry. Annu Rev Clin Psychol 2018;14:91-118. [doi: 10.1146/annurev-clinpsy-032816-045037] [Medline: 29401044]
- 20. Hornstein S, Zantvoort K, Lueken U, Funk B, Hilbert K. Personalization strategies in digital mental health interventions: a systematic review and conceptual framework for depressive symptoms. Front Digital Health 2023;5:1170002 [FREE Full text] [doi: 10.3389/fdgth.2023.1170002] [Medline: 37283721]
- Swaminathan A, López I, Mar RAG, Heist T, McClintock T, Caoili K, et al. Natural language processing system for rapid detection and intervention of mental health crisis chat messages. NPJ Digital Med 2023;6(1):213 [FREE Full text] [doi: 10.1038/s41746-023-00951-3] [Medline: 37990134]
- 22. Hornstein S, Scharfenberger J, Lueken U, Wundrack R, Hilbert K. Predicting recurrent chat contact in a psychological intervention for the youth using natural language processing. NPJ Digital Med 2024;7(1):132 [FREE Full text] [doi: 10.1038/s41746-024-01121-9] [Medline: 38762694]
- Xu Y, Chan CS, Tsang C, Cheung F, Chan E, Fung J, et al. Detecting premature departure in online text-based counseling using logic-based pattern matching. Internet Interventions 2021;26:100486 [FREE Full text] [doi: 10.1016/j.invent.2021.100486] [Medline: 34877263]
- 24. Kohls E, Guenthner L, Baldofski S, Eckert M, Efe Z, Kuehne K, et al. Suicidal ideation among children and young adults in a 24/7 messenger-based psychological chat counseling service. Front Psychiatry 2022;13:862298 [FREE Full text] [doi: 10.3389/fpsyt.2022.862298] [Medline: 35418889]
- Broadbent M, Grespan MM, Axford K, Zhang X, Srikumar V, Kious B, et al. A machine learning approach to identifying suicide risk among text-based crisis counseling encounters. Front Psychiatry 2023;14:1110527 [FREE Full text] [doi: 10.3389/fpsyt.2023.1110527] [Medline: 37032952]
- 26. Xu Z, Xu Y, Cheung F, Cheng M, Lung D, Law YW, et al. Detecting suicide risk using knowledge-aware natural language processing and counseling service data. Soc Sci Med 2021;283:114176. [doi: <u>10.1016/j.socscimed.2021.114176</u>] [Medline: <u>34214846</u>]
- Imel ZE, Tanana MJ, Soma CS, Hull TD, Pace BT, Stanco SC, et al. Mental health counseling from conversational content with transformer-based machine learning. JAMA Netw Open 2024;7(1):e2352590 [FREE Full text] [doi: 10.1001/jamanetworkopen.2023.52590] [Medline: 38252437]
- 28. Li A, Ma J, Ma L, Fang P, He H, Lan Z. Towards automated real-time evaluation in text-based counseling. ArXiv. Preprint posted online on March 07, 2022 2022 [FREE Full text]

```
https://ai.jmir.org/2025/1/e63701
```

- 29. Rickwood D, Deane FP, Wilson CJ, Ciarrochi J. Young people's help-seeking for mental health problems. Aust e-J Adv Mental Health 2014;4(3):218-251. [doi: 10.5172/jamh.4.3.218]
- de Winter AF, Oldehinkel AJ, Veenstra R, Brunnekreef JA, Verhulst FC, Ormel J. Evaluation of non-response bias in mental health determinants and outcomes in a large sample of pre-adolescents. Eur J Epidemiol 2005;20(2):173-181 [FREE Full text] [doi: 10.1007/s10654-004-4948-6] [Medline: 15792285]
- Cheung KL, Ten Klooster PM, Smit C, de Vries H, Pieterse ME. The impact of non-response bias due to sampling in public health studies: a comparison of voluntary versus mandatory recruitment in a Dutch national survey on adolescent health. BMC Public Health 2017;17(1):276 [FREE Full text] [doi: 10.1186/s12889-017-4189-8] [Medline: 28330465]
- 32. Automated evaluation of helpfulness of chat-counseling sessions for the youth. a natural language processing study. OSF Registries. URL: <u>https://osf.io/sr4q9</u> [accessed 2024-06-26]
- Klement W, El Emam K. Consolidated reporting guidelines for prognostic and diagnostic machine learning modeling studies: development and validation. J Med Internet Res 2023;25:e48763 [FREE Full text] [doi: 10.2196/48763] [Medline: 37651179]
- 34. silvanhornstein/AutoEval: code for paper (OSF: SR4Q9). GitHub. URL: <u>https://github.com/silvanhornstein/AutoEval</u> [accessed 2024-06-26]
- 35. Hornstein S, Forman-Hoffman V, Nazander A, Ranta K, Hilbert K. Predicting therapy outcome in a digital mental health intervention for depression and anxiety: a machine learning approach. Digital Health 2021;7:20552076211060659 [FREE Full text] [doi: 10.1177/20552076211060659] [Medline: 34868624]
- 36. Wartena C. A probabilistic morphology model for German lemmatization. 2019. URL: <u>https://serwiss.bib.hs-hannover.de/</u> <u>frontdoor/index/docId/1527</u> [accessed 2019-01-01]
- Saito T, Rehmsmeier M. The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets. PLoS One 2015;10(3):e0118432 [FREE Full text] [doi: 10.1371/journal.pone.0118432] [Medline: 25738806]
- 38. Halimu C, Kasem A, Newaz S. Empirical comparison of area under ROC curve (AUC) and mathew correlation coefficient (MCC) for evaluating machine learning algorithms on imbalanced datasets for binary classification. 2019 Presented at: Proceedings of the 3rd International Conference on Machine Learning and Soft Computing; January 25-28, 2019; Da Lat, Vietnam p. 1-6. [doi: 10.1145/3310986.3311023]
- 39. McDermott MBA, Zhang H, Hansen LH, Angelotti G, Gallifant J. A closer look at AUROC and AUPRC under class imbalance. ArXiv. Preprint posted online on January 11, 2024 2024. [doi: <u>10.48550/arXiv.2401.06091</u>]
- 40. Chawla NV, Bowyer KW, Hall LO, Kegelmeyer WP. SMOTE: synthetic minority over-sampling technique. J Artif Intell Res 2002;16(1):321-357. [doi: <u>10.1613/jair.953</u>]
- Chen T, Guestrin C. XGBoost: a scalable tree boosting system. 2016 Presented at: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining; August 13-17, 2016; San Francisco, CA p. 785-794. [doi: 10.1145/2939672.2939785]
- 42. The AI community building the future. Hugging Face. URL: <u>https://huggingface.co/</u> [accessed 2024-04-05]
- 43. Scheible R, Thomczyk F, Tippmann P, Jaravine V, Boeker M. GottBERT: a pure German language model. ArXiv. Preprint posted online on December 03, 2020 2020 [FREE Full text] [doi: <u>10.48550/arXiv.2012.02110</u>]
- 44. Sanh V, Debut L, Chaumond J, Wolf T. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. ArXiv. Preprint posted online on October 2, 2019 2019 [FREE Full text]
- 45. Beltagy I, Peters M, Cohan A. Longformer: the long-document transformer. ArXiv. Preprint posted online on April, 10, 2020 2020 [FREE Full text]
- 46. Ojala M, Garriga GC. Permutation tests for studying classifier performance. 2009 Presented at: 2009 Ninth IEEE International Conference on Data Mining; December 06-09, 2009; Miami Beach, FL. [doi: <u>10.1109/icdm.2009.108</u>]
- 47. DeLong ER, DeLong DM, Clarke-Pearson DL. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. Biometrics 1988;44(3):837-845. [Medline: <u>3203132</u>]
- 48. Rainio O, Teuho J, Klén R. Evaluation metrics and statistical tests for machine learning. Sci Rep 2024;14(1):6086 [FREE Full text] [doi: 10.1038/s41598-024-56706-x] [Medline: 38480847]
- 49. Dietterich TG. Approximate statistical tests for comparing supervised classification learning algorithms. Neural Comput 1998;10(7):1895-1923. [doi: 10.1162/089976698300017197] [Medline: 9744903]
- 50. Lundberg SM, Lee SI. A unified approach to interpreting model predictions. Advances in Neural Information Processing Systems. URL: <u>https://papers.nips.cc/paper_files/paper/2017/hash/8a20a8621978632d76c43dfd28b67767-Abstract.html</u> [accessed 2025-02-04]
- 51. Chicco D, Tötsch N, Jurman G. The Matthews correlation coefficient (MCC) is more reliable than balanced accuracy, bookmaker informedness, and markedness in two-class confusion matrix evaluation. BioData Min 2021;14(1):13 [FREE Full text] [doi: 10.1186/s13040-021-00244-z] [Medline: 33541410]
- 52. Chicco D, Jurman G. The Matthews correlation coefficient (MCC) should replace the ROC AUC as the standard metric for assessing binary classification. BioData Min 2023;16(1):4 [FREE Full text] [doi: 10.1186/s13040-023-00322-4] [Medline: 36800973]



- 53. Kokalj E, Škrlj B, Lavrač N, Pollak S, Robnik-Šikonja M. BERT meets Shapley: extending SHAP explanations to transformer-based classifiers. 2021 Presented at: Proceedings of the EACL Hackashop on News Media Content Analysis and Automated Report Generation; February 03, 2025; Hackashop p. 16-21 URL: <u>https://aclanthology.org/</u> 2021.hackashop-1.3/
- 54. Li Z. Extracting spatial effects from machine learning model using local interpretation method: an example of SHAP and XGBoost. Comput Environ Urban Syst 2022;96:101845. [doi: <u>10.1016/j.compenvurbsys.2022.101845</u>]
- 55. Zhou Q, Chen ZH, Cao YH, Peng S. Clinical impact and quality of randomized controlled trials involving interventions evaluating artificial intelligence prediction tools: a systematic review. NPJ Digital Med 2021;4(1):154 [FREE Full text] [doi: 10.1038/s41746-021-00524-2] [Medline: 34711955]
- Zantvoort K, Scharfenberger J, Boß L, Lehr D, Funk B. Finding the best match—a case study on the (text-)feature and model choice in digital mental health interventions. J Healthcare Inform Res 2023;7(4):447-479 [FREE Full text] [doi: 10.1007/s41666-023-00148-z] [Medline: <u>37927375</u>]
- 57. Gogoulou E, Boman M, Abdesslem F, Isacsson N, Kaldo V, Sahlgren M. Predicting treatment outcome from patient texts: the case of internet-based cognitive behavioural therapy. 2021 Presented at: Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume; February 03, 2025; Virtual event p. 575-580 URL: https://aclanthology.org/2021.eacl-main.46/ [doi: http
- 58. Funk B, Sadeh-Sharvit S, Fitzsimmons-Craft EE, Trockel MT, Monterubio GE, Goel NJ, et al. A framework for applying natural language processing in digital health interventions. J Med Internet Res 2020;22(2):e13855 [FREE Full text] [doi: 10.2196/13855] [Medline: 32130118]
- 59. Garrido-Merchan EC, Gozalo-Brizuela R, Gonzalez-Carvajal S. Comparing BERT against traditional machine learning models in text classification. JCCE 2023;2(4):352-356. [doi: 10.47852/bonviewjcce3202838]
- 60. Bartz E, Zaefferer M, Mersmann O, Bartz-Beielstein T. Experimental investigation and evaluation of model-based hyperparameter optimization. ArXiv. Preprint posted online on July 19, 2021 2021 [FREE Full text]
- 61. Turner R, Eriksson D, McCourt M, Kiili J, Laaksonen E, Xu Z, et al. Bayesian optimization is superior to random search for machine learning hyperparameter tuning: analysis of the black-box optimization challenge 2020. PMLR 2020;133:3-26 [FREE Full text] [doi: 10.1007/978-1-4842-6579-6_4]
- 62. Liu Z, Peach RL, Lawrance EL, Noble A, Ungless MA, Barahona M. Listening to mental health crisis needs at scale: using natural language processing to understand and evaluate a mental health crisis text messaging service. Front Digital Health 2021;3:779091 [FREE Full text] [doi: 10.3389/fdgth.2021.779091] [Medline: 34939068]
- 63. El-Ramly M, Abu-Elyazid H, Mo?men Y, Alshaer G, Adib N, Eldeen KA. CairoDep: detecting depression in arabic posts using BERT transformers. : IEEE; 2021 Presented at: 2021 Tenth International Conference on Intelligent Computing and Information Systems (ICICIS); December 05-07, 2021; Cairo, Egypt. [doi: 10.1109/icicis52592.2021.9694178]
- 64. Wang S, Dang Y, Sun Z, Ding Y, Pathak J, Tao C, et al. An NLP approach to identify SDoH-related circumstance and suicide crisis from death investigation narratives. J Am Med Inform Assoc 2023;30(8):1408-1417 [FREE Full text] [doi: 10.1093/jamia/ocad068] [Medline: 37040620]
- 65. mim-solutions / bert_for_longer_texts. GitHub. URL: <u>https://github.com/mim-solutions/bert_for_longer_texts</u> [accessed 2024-08-26]
- 66. Vajre V, Naylor M, Kamath U, Shehu A. PsychBERT: a mental health language model for social media mental health behavioral analysis. 2021 Presented at: 2021 IEEE International Conference on Bioinformatics and Biomedicine (BIBM); December 9-12, 2021; Houston, TX p. 1077-1082. [doi: 10.1109/bibm52615.2021.9669469]
- 67. Chekroud AM, Hawrilenko M, Loho H, Bondar J, Gueorguieva R, Hasan A, et al. Illusory generalizability of clinical prediction models. Science 2024;383(6679):164-167. [doi: <u>10.1126/science.adg8538</u>] [Medline: <u>38207039</u>]

Abbreviations

AI: artificial intelligence
AUC: area under the curve
MCC: Matthews correlation coefficient
NLP: natural language processing
NPP: negative predictive value
ROC: receiver operating characteristic
SHAP: Shapley additive explanations
TF-IDF: term frequency-inverse document frequency
XGBoost: extreme gradient boosting



Edited by K El Emam, B Malin; submitted 27.06.24; peer-reviewed by R Scheible, A Li; comments to author 17.08.24; revised version received 04.09.24; accepted 02.12.24; published 18.02.25. <u>Please cite as:</u> Hornstein S, Lueken U, Wundrack R, Hilbert K Predicting Satisfaction With Chat-Counseling at a 24/7 Chat Hotline for the Youth: Natural Language Processing Study JMIR AI 2025;4:e63701 URL: https://ai.jmir.org/2025/1/e63701 doi:10.2196/63701 PMID:

©Silvan Hornstein, Ulrike Lueken, Richard Wundrack, Kevin Hilbert. Originally published in JMIR AI (https://ai.jmir.org), 18.02.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Original Paper

Leveraging Medical Knowledge Graphs Into Large Language Models for Diagnosis Prediction: Design and Application Study

Yanjun Gao^{1,2}, PhD; Ruizhe Li³, PhD; Emma Croxford², BS; John Caskey², PhD; Brian W Patterson², MPH, MD; Matthew Churpek², MPH, MD, PhD; Timothy Miller⁴, PhD; Dmitriy Dligach⁵, PhD; Majid Afshar², MD, MSCR

¹Department of Biomedical Informatics, University of Colorado Anschutz Medical Campus, Denver, CO, United States

²Department of Medicine, University of Wisconsin-Madison, Madison, WI, United States

³University of Aberdeen, Aberdeen, United Kingdom

⁴Boston Children's Hospital, Harvard Medical School, Boston, MA, United States

⁵Loyola University Chicago, Chicago, IL, United States

Corresponding Author:

Yanjun Gao, PhD Department of Biomedical Informatics University of Colorado Anschutz Medical Campus 1890 N Revere Ct Denver, CO, 80045 United States Phone: 1 303 724 5375 Email: yanjun.gao@cuanschutz.edu

Abstract

Background: Electronic health records (EHRs) and routine documentation practices play a vital role in patients' daily care, providing a holistic record of health, diagnoses, and treatment. However, complex and verbose EHR narratives can overwhelm health care providers, increasing the risk of diagnostic inaccuracies. While large language models (LLMs) have showcased their potential in diverse language tasks, their application in health care must prioritize the minimization of diagnostic errors and the prevention of patient harm. Integrating knowledge graphs (KGs) into LLMs offers a promising approach because structured knowledge from KGs could enhance LLMs' diagnostic reasoning by providing contextually relevant medical information.

Objective: This study introduces DR.KNOWS (Diagnostic Reasoning Knowledge Graph System), a model that integrates Unified Medical Language System–based KGs with LLMs to improve diagnostic predictions from EHR data by retrieving contextually relevant paths aligned with patient-specific information.

Methods: DR.KNOWS combines a stack graph isomorphism network for node embedding with an attention-based path ranker to identify and rank knowledge paths relevant to a patient's clinical context. We evaluated DR.KNOWS on 2 real-world EHR datasets from different geographic locations, comparing its performance to baseline models, including QuickUMLS and standard LLMs (Text-to-Text Transfer Transformer and ChatGPT). To assess diagnostic reasoning quality, we designed and implemented a human evaluation framework grounded in clinical safety metrics.

Results: DR.KNOWS demonstrated notable improvements over baseline models, showing higher accuracy in extracting diagnostic concepts and enhanced diagnostic prediction metrics. Prompt-based fine-tuning of Text-to-Text Transfer Transformer with DR.KNOWS knowledge paths achieved the highest ROUGE-L (Recall-Oriented Understudy for Gisting Evaluation–Longest Common Subsequence) and concept unique identifier F_1 -scores, highlighting the benefits of KG integration. Human evaluators found the diagnostic rationales of DR.KNOWS to be aligned strongly with correct clinical reasoning, indicating improved abstraction and reasoning. Recognized limitations include potential biases within the KG data, which we addressed by emphasizing case-specific path selection and proposing future bias-mitigation strategies.

Conclusions: DR.KNOWS offers a robust approach for enhancing diagnostic accuracy and reasoning by integrating structured KG knowledge into LLM-based clinical workflows. Although further work is required to address KG biases and extend generalizability, DR.KNOWS represents progress toward trustworthy artificial intelligence–driven clinical decision support, with a human evaluation framework focused on diagnostic safety and alignment with clinical standards.

(JMIR AI 2025;4:e58670) doi:10.2196/58670



KEYWORDS

knowledge graph; natural language processing; machine learning; electronic health record; large language model; diagnosis prediction; graph model; artificial intelligence

Introduction

Background

The ubiquitous use of electronic health records (EHRs) and the standard documentation practice of daily care notes are integral to the continuity of patient care because these records provide a comprehensive account of the patient's health trajectory, inclusive of condition status, diagnoses, and treatment plans [1]. Nevertheless, the growing complexity and verbosity of EHR clinical narratives, which are often filled with redundant information, can overwhelm health care providers and increase the risk of diagnostic errors [2-5]. Physicians often skip sections of lengthy and repetitive notes and rely on decisional shortcuts (ie, decisional heuristics) that can contribute to diagnostic errors [6].

Current efforts at automating diagnosis generation from daily progress notes leverage large language models (LLMs). Gao et al [7] introduced a summarization task that takes progress notes as input and generates a summary of active diagnoses. The authors annotated a set of progress notes from the publicly available EHR dataset Medical Information Mart for Intensive Care III (MIMIC-III) [8]. The BioNLP 2023 shared task, known as ProbSum, built upon this work by providing additional annotated notes and attracting multiple efforts focused on developing solutions [9-11]. Demonstrating a growing interest in applying LLMs to serve as solutions, these prior studies use language models such as Text-to-Text Transfer Transformer (T5) [12], developed by Google Research; and Open AI's Generative Pretrained Transformer (GPT) [13]. Unlike the conventional language tasks where LLMs have shown promising abilities, automated diagnosis generation is a critical task that requires high accuracy and reliability to ensure patient safety and improve health care outcomes. Concerns regarding the potential misleading and hallucinated information that could

result in life-threatening events prevent LLMs from being used for diagnostic prediction [14].

The Unified Medical Language System (UMLS) [15], a comprehensive resource developed by the National Library of Medicine in the United States, has been extensively used in natural language processing (NLP) research. The UMLS serves as a medical knowledge repository, facilitating the integration, retrieval, and sharing of biomedical information. It offers concept vocabulary and semantic relationships, enabling the construction of medical knowledge graphs (KGs). Prior studies have leveraged UMLS KGs for tasks such as information extraction [16-19] and question answering [17]. Mining relevant knowledge for diagnosis is particularly challenging for 2 reasons: the highly specific factors related to the patient's complaints, histories, and symptoms documented in the EHR; and the vast search space within a KG containing 4.5 million concepts and 15 million relations for diagnosis determination.

In this study, we explore the use of KGs as external resources to enhance LLMs for diagnosis generation. Our work is motivated not only by the potential in the NLP field of augmenting LLMs with KGs [20] but also by the theoretical exploration in medical education and psychology research, shedding light on the diagnostic decision-making process used by clinicians. Forming a diagnostic decision requires the examination of patient data, retrieving encapsulated medical knowledge, and the formulation and testing of the diagnostic hypothesis, which is also known as clinical diagnostic reasoning [21,22]. We propose a novel graph model, DR.KNOWS (Diagnostic Reasoning Knowledge Graph System), designed to retrieve the top N case-specific knowledge paths related to disease pathology and feed them into foundational LLMs to improve the accuracy of diagnostic predictions (as shown in Figure 1). Two distinct foundational models are the subject of this study: T5, known for being fine-tunable; and a sandboxed version of ChatGPT, a powerful LLM where we explore zero-shot prompting.


Figure 1. Study overview: we focused on generating diagnoses (text given in red in the "Plan" section) using the SOAP (subjective, objective, assessment, and plan) format progress note with the aid of large language models (LLMs). The input consists of "Subjective," "Objective," and "Assessment" sections (the dotted line box below the heading "Patient Progress Note"), and the diagnoses in the "Plan" section are the ground truth. We introduced an innovative knowledge graph (KG) model, namely DR.KNOWS (Diagnostic Reasoning Knowledge Graph System), that identifies and extracts the most relevant knowledge trajectories from the Unified Medical Language System (UMLS) KG. The nodes of the UMLS KG represent concept unique identifiers (CUIs), and the edges denote the semantic relations among the CUIs. We experimented with prompting ChatGPT for diagnosis generation, with and without the knowledge paths predicted by DR.KNOWS. Furthermore, we investigated how this knowledge grounding influences the diagnostic output of LLMs using human evaluation. The underlined text shows the UMLS concepts identified through a concept extractor. EtOH: ethanol; GI: gastrointestinal; REDCap: Research Electronic Data Capture; T5: Text-to-Text Transfer Transformer; UGIB: upper gastrointestinal bleeding.



Objectives

Our work and contribution are structured into two primary components: (1) designing and evaluating DR.KNOWS, a graph-based model that selects the top N probable diagnoses with explainable paths; and (2) demonstrating the usefulness of DR.KNOWS as an additional module to augment pretrained language models in generating relevant diagnoses. Along with the technical contributions, we propose the first human evaluation framework for LLM-generated diagnoses that adapts a survey instrument designed to evaluate diagnostic safety. Our research poses a new exciting problem that has not been addressed in the realm of NLP for diagnosis generation, that is, harnessing the power of KGs for the controllability and explainability of foundational models. By examining the effects of KG path-based prompts on foundational models on a real-world hospital dataset, we strive to contribute to an explainable artificial intelligence (AI) diagnostic pathway.

Several studies have focused on the application of clinical note summarization to discharge summaries [23], hospital course narratives [24], real-time patient visit summaries [25], and problem and diagnosis lists [7,26,27]. Our work follows the line of research on problem and diagnosis summarization. The integration of KGs with LLMs has been gaining traction as an emerging trend due to the potential enhancement of factual knowledge [20], especially on domain-specific question-answering tasks [28-30]. Our work stands out by integrating KGs into LLMs for diagnosis prediction, using a novel graph model for path-based prompts.

https://ai.jmir.org/2025/1/e58670

RenderX

Methods

Problem Formulation

Daily Progress Notes for Diagnosis Prediction

Daily progress notes are formatted using the SOAP (subjective, objective, assessment, and plan) format [30]. The subjective section of a SOAP daily progress note comprises the patient's self-reported symptoms, concerns, and medical history. The objective section consists of structural data collected by health care providers during observation or examination, such as vital signs (eg, blood pressure and heart rate), laboratory results, or physical examination findings. The assessment section summarizes the patient's overall condition, with a focus on the most active problems and diagnoses for that day. Finally, the plan section contains multiple subsections, each outlining a diagnosis or problem and its treatment plan. Our task is to predict the list of problems and diagnoses that are part of the plan section. Our research used the ProbSum dataset, an annotated resource created for the BioNLP 2023 shared task with gold standard diagnoses derived from progress notes [27].

Using UMLS KGs to Find Potential Diagnoses, Given Medical Narratives

The UMLS concepts vocabulary comprises >180 sources. For our study, we focused on the Systematized Nomenclature of Medicine–Clinical Terms (SNOMED CT). The UMLS vocabulary is a comprehensive, multilingual health terminology and the US national standard for EHRs and health information exchange. Each UMLS medical concept is assigned a SNOMED

CT concept unique identifier (CUI) from the clinical terminology system. We used semantic types, networks, and semantic relations from UMLS knowledge sources to categorize concepts based on shared attributes, enabling efficient exploration and supporting semantic understanding and knowledge discovery across various medical vocabularies.

Given a medical KG where the nodes represent concepts and the edges denote semantic relations along with an input text describing a patient's problems, we could perform multihop reasoning across the KG and infer the final diagnoses. Figure 2 demonstrates how UMLS semantic relations and concepts can be used to identify potential diagnoses from the evidence provided in a daily care note. The example patient presents with medical conditions of fever, cough and sepsis, which are the concepts recognized by medical concept extractors (Clinical Text Analysis and Knowledge Extraction System [31] and QuickUMLS [32]) and the starting concepts for multihop reasoning. Initially, we extracted the direct neighbors for these concepts. Relevant concepts that aligned with the patient's descriptions were preferred. For precise diagnoses, we chose the top N most relevant nodes at each hop.

Figure 2. Problem formulation: inferring possible diagnoses within 2 hops from a Unified Medical Language System (UMLS) knowledge graph given a patient's medical description. The UMLS medical concepts are highlighted in the colored boxes ("female," "sepsis," etc). Each concept has its own subgraph, where concepts are the vertices, and semantic relations are the edges (owing to space constraints, we have omitted the subgraph for "female" in this graph presentation). On the first hop, we could identify the most relevant neighboring concepts to the input description. The darker the color of the vertices, the more relevant they are to the input description. A second hop could be further performed based on the most relevant nodes, leading to the final diagnoses "Pneumonia and influenza" and "Respiratory distress syndrome." Of note, we use the preferred text of concept unique identifiers for presentation purposes. The actual UMLS knowledge graph is built on concept unique identifiers rather than preferred text.

"This is a female in her 40s with fever, coughing, and sepsis."



The UMLS's vast repository consists of 270 semantic relations, but not all are crucial for diagnostic reasoning. Adding the nonrelevant relations into a KG introduced substantially complexities in both computation and retrieval processes. A board-certified physician (MA) refined these to identify the 107 most relevant relations for diagnostics, which were then used to build the UMLS KG. This selection, including relations such as "causative agent of" and excluding ones such as "inverse isa," is vital to maintaining computational efficiency and retrieval accuracy within the KG.

Data Overview

We used 2 sets of progress notes from different clinical settings in this study: MIMIC-III and in-house EHR datasets. MIMIC-III is one of the largest publicly available databases containing deidentified health data from patients admitted to intensive care units. It was developed by the Massachusetts Institute of Technology and Beth Israel Deaconess Medical Center.

https://ai.jmir.org/2025/1/e58670

RenderX

MIMIC-III includes data from >38,000 patients admitted to intensive care units at the Beth Israel Deaconess Medical Center between 2001 and 2012. The second set, namely the in-house EHR data, was a subset of EHRs that included adult patients (aged 18 years) admitted to the University of Wisconsin health system between 2008 and 2021. In contrast to the MIMIC-III subset, the in-house set covered progress notes from all hospital settings, including the emergency department, general medicine wards, and subspecialty wards. While the 2 datasets originated from separate hospitals and departmental settings and might reflect distinct note-taking practices, both followed the SOAP documentation format for progress notes.

Gao et al [7,9] introduced a subset of 1005 progress notes from MIMIC-III with active diagnoses annotated from the "plan" sections, namely, the ProbSum dataset. Therefore, we applied this dataset for training and evaluation for both graph model intrinsic evaluation and diagnosis summarization. The in-house dataset did not contain human annotation. Even so, by parsing

the text with a medical concept extractor that was based on UMLS SNOMED CT vocabulary, we were able to pull out concepts that belonged to the semantic type of "T047 Disease and Syndromes." We deployed this set of concepts as the ground truth data to train and evaluate the graph model. The final in-house dataset contained 4815 progress notes. We present the descriptive statistics in Table 1. When contrasted with MIMIC-III, the in-house dataset exhibited a greater number of CUIs in its input, leading to an extended CUI output. In addition, MIMIC-III encompassed a wider range of abstractive concepts compared to the in-house progress notes.

Table 1. Average number of concept unique identifiers (CUIs) in the input and output across the 2 electronic health record datasets: Medical Information

 Mart for Intensive Care III (MIMIC-III) and in-house. Abstractive concepts are those not found in the input but present in the gold standard diagnoses.

Datasets	Departments	Input CUIs (n), mean (SD)	Output CUIs (n), mean (SD)	Abstractive CUIs (%)
MIMIC-III	ICU ^a	15.95	3.51	48.92
In-house	All	41.43	5.81	<1

^aICU: intensive care unit.

Graph Model Development

Overview

This section introduces the architecture design for DR.KNOWS. The DR.KNOWS model is designed to enhance automated diagnostic reasoning by integrating structured clinical knowledge from the UMLS into patient-specific diagnostic predictions. By leveraging a graph-based approach, DR.KNOWS retrieves and ranks relevant knowledge paths from the UMLS, ensuring that only clinically pertinent information is considered. Using a graph neural network, DR.KNOWS incorporates topological information from the UMLS KG into concept representations to better determine each node's relevance to the patient's specific conditions.

Architecture Overview

As shown in Figure 3, all identified UMLS concepts with an assigned CUI from the input patient text were used to retrieve 1-hop subgraphs from the constructed large UMLS KG. Each node in this graph represents a CUI; therefore, we use "node" and "concept (CUI)" interchangeably throughout. These 1-hop subgraphs are encoded by a stack graph isomorphism network (SGIN) [33], which generates node embeddings that capture both neighboring concept information and pretrained concept embeddings. We chose the SGIN for node embedding because it matches the expressive power of the Weisfeiler-Lehman graph isomorphism test, maximizing the graph neural network's ability to capture meaningful representations. The resulting node embeddings serve as the basis for path embeddings, which the path encoder further processes.

Figure 3. DR.KNOWS (Diagnostic Reasoning Knowledge Graph System) model architecture. The input concepts ("female," "fever," etc) are represented by concept unique identifiers (CUIs) represented as a combination of letters and numbers (eg, "C0243026" and "C0015967"). SapBERT: Self-alignment Pretrained Bidirectional Encoder Representations from Transformers.



The path encoder module then evaluates these 1-hop paths by examining their semantic and logical alignment with the input text and concept representations, assigning a relevance score to each path. The top N scores across these paths, aggregated across each node's neighboring paths, guide the selection of nodes for the next hop. If no suitable diagnosis node is found, the path exploration terminates by assigning a self-loop to the current node.

While the dominant technique for retrieval-augmented generation systems relies heavily on vector representations and cosine similarity for retrieving and ranking candidate text, our work goes beyond this by adding 2 extra layers of design. First, we leverage the expressive power of the graph structure to

```
https://ai.jmir.org/2025/1/e58670
```

enhance the retrieval process. Second, we select paths not simply based on their embeddings but through an attention network that encodes the path-concept relationships, ensuring a more accurate and contextually relevant selection process. In the following paragraphs, we present details regarding each component in the architecture of DR.KNOWS.

Contextualized Node Representation

We define the deterministic UMLS KG G = VE based on SNOMED CT CUIs and semantic relations, where V is a set of CUIs, and E is a set of semantic relations. Given an input text x containing a set of source CUIs $V_{src} \subseteq V$ and their 1-hop relations $E_{src} \subseteq E$, we can construct relation paths for each

source node $v_{src} \subseteq V_{src}$ as $P = \{p_1, p_2, \dots p_j\}$ such that $p_j = \{v_1, \dots, v_j\}$ $e_1, v_2, \dots e_{j-1}, v_j\}, j \subseteq J$, where J is the maximum length that a source node v_{src} could reach and is nondeterministic. Relations e are encoded as one-hot embeddings. We concatenate all concept names for v_i with special tokens such as [SEP] (for "separator"), such that $l_i = [name 1 [SEP] name 2 [SEP]...]$ and encode l_i using Self-alignment Pretrained Bidirectional Encoder Representations from Transformers (SapBERT) [34] to obtain h_i as concept representation. This allows the CUI representation to serve as the contextualized representation of its corresponding concept names. We chose SapBERT for its contrastive learning-based training, which discriminates similar concepts and their synonyms. It is evaluated on entity linking tasks and has shown state-of-the-art performance. The h_i is further updated through topological representation using the SGIN to become node representation:



 $N(v_i)$ represents the set of neighboring nodes of node v_i , \bowtie is the representation of node v_i at layer k, $^{(k)}$ is a learnable parameter at layer k, and $MLP^{(k)}$ is a multilayer perceptron at layer k. GIN iteratively aggregates neighborhood information using graph convolution followed by nonlinearity, modeling interactions among nodes within the set \bowtie . Furthermore, the stacking mechanism is introduced to combine multiple GIN layers. The final node representation v_i at layer K (last layer) is computed by stacking the GIN layers, where [...;...] denotes matrix concatenation.

We empirically observed that some types of CUIs are less likely to lead to useful paths for diseases, for example, the concept "recent" (CUI: C0332185) is a temporal concept, and the neighbors associated with it are less useful to predict diagnoses. We designed a weighting scheme based on term frequency–inverse document frequency to assign higher weights to more relevant CUIs and semantic types:



 W_{CUI} are then multiplied by the corresponding h_i to assign weighted representations to the concept representation.

Path Reasoning and Ranking

For each node representation h_i , we use its n-hop \Join of the set neighborhood for \Join for h_i and the associated relation edge \Join to generate the corresponding path embeddings, with *t* being the index of the node and its associated neighborhood and relations:

hi, if n=1 pi = {

💌, otherwise

×

where "FFN" is the feedforward network, and *n* is the number of hops in the subgraph G_{src} . The path embedding p_i is the node embedding itself for the first hop and is recursively aggregated with new nodes and edges as the path extends to the next hop.

To determine each path's relevance to the patient's specific symptoms, we used 2 attention mechanisms—multihead attention (MultiAttn) and trilinear attention (TriAttn)—to compute scores S for each path. Both mechanisms use the patient's input text representation h_x and input list of CUIs h_v , encoded by SapBERT, to capture explicit and intricate relationships in the input data. MultiAttn was used to explicitly capture relationships between the input text, the list of concepts, and the current path, while TriAttn was used to automatically learn these complex relationships through the inner products of the 3 matrices. As demonstrated in Figure 2, for each hop the path tries to achieve based on the input patient description, the candidate concept can add relevant information, provide no new information and remain neutral, or contradict the information already present in the context.

Using MultiAttn, we define the context relevancy matrix H_i and the concept relevancy matrix Z_i as follows:

 $\begin{aligned} Hi &= [hx; pi; hx - pi; hx \odot pi] \\ Zi &= [hv; pi; hv - pi; hv \odot pi] \\ \alpha i &= MultiAttn(Hi \odot Zi), \\ SMulti &= \varphi (Relu(\sigma(\alpha i))) \end{aligned}$

These relevancy matrices are inspired by a prior work on natural language inference [35], representing logical relations such as neutrality, contradiction, and entailment via matrix concatenation, difference, and product, respectively. Alternatively, TriAttn learns the intricate relations by 3 attention maps:

 $\alpha i = (hx, hv, pi) = \Sigma abc (hx)a (hv)b (pi)c Wabc$ STri = φ (Relu($\sigma(\alpha i)$))

 h_x , h_v , and p_i have the same dimensionality D, and φ is an MLP player. Finally, we aggregate the MultiAttn or TriAttn scores on all candidate nodes and select the top N nodes (concepts) V_N for the next iteration based on the aggregate attention scores:

	Γ	×
	h	-

 $V_N = argmax_N(\beta)$

By comparing attention scores across candidate paths, the path ranker selects the top N nodes most relevant to each patient's symptoms, maximizing contextual relevance.

Loss Function

Our loss function consists of 2 parts: a CUI prediction loss L_{pred} and a contrastive learning loss L_{CL} :

$$L = L_{pred} + L_{CL}$$

For CUI prediction loss, we use binary cross entropy loss to calculate whether the predicted node V_N is in the gold standard label *Y*:

[×	
ļ		L

Where M is the number of sets of gold labels. For contrastive learning loss L_{CL} , we encourage the model to learn meaningful and discriminative representations through comparison with positive and negative samples:

Y	
^	

where A_i is the anchor embedding, defined as $h_x \odot h_v$, representing the input text and concept representation. Σ_i indicates a summation over a set of indices *i*, typically representing different training samples or pairs. Inspired by the study by Hu et al [29], we construct $\cos(A_i, f_i)$ and $\cos(A_i, f_{i-})$ to calculate cosine similarity between A_i and positive feature f_{i+} or negative feature f_{i-} , respectively. A positive feature represents the paths correctly leading to the ground truth concept, while a negative feature embodies the paths that, although starting from the source, culminate in an incorrect concept. This equation measures the loss when the similarity between an anchor and its positive feature is not significantly greater than the similarity between the same anchor and a negative feature, considering a margin for desired separation.

We designed a training algorithm to iteratively select and rank the most relevant paths to extend. This algorithm helped to reduce the computational requirement because it does not rank all n-hop paths within 1 pass. This algorithm is presented in Multimedia Appendix 1.

Selection of Foundational Models and Experiment Setup

Our study centers around the following question: To what extent does the incorporation of DR.KNOWS as a knowledge path–based prompt provider influence the performance of language models in diagnosis summarization?

We present results derived from 2 distinct foundational models, varying significantly in their parameter scales, namely T5-Large, which comprises 770 million parameters [12]; and GPT-3.5-Turbo, which features 154 billion parameters [13]. Specifically, we were granted access to a restricted version of the GPT-3.5-Turbo model, which served as the underlying framework for the highly capable language model, ChatGPT.

These 2 models represent the prevailing direction in the evolution of language models: smaller models such as T5 that offer easier control and larger models such as GPT that generate text with substantial scale and power. Our investigation focused on evaluating the performance of T5 in fine-tuning scenarios and GPT models in zero-shot settings. Our primary objective was not solely to demonstrate cutting-edge results but also to critically examine the potential influence of incorporating predicted paths, generated by graph models, as auxiliary knowledge contributors.

```
https://ai.jmir.org/2025/1/e58670
```

We selected 3 distinct T5-Large variants for fine-tuning using the ProbSum summarization dataset. The chosen T5 models encompass the vanilla T5 [12], a foundational model that has been extensively used in varied NLP tasks; Flan-T5 [36], which has been fine-tuned using an instructional approach; and Clinical-T5 [37], which has been specifically trained on the MIMIC dataset.

Given that our work encompasses a public EHR dataset (MIMIC-III) and a private EHR dataset with protected health information (in-house), we conducted training using 3 distinct computing environments. Specifically, most of the experiments on MIMIC-III were conducted on Google's cloud computing platform, using 1 to 2 NVIDIA A100 40 GB graphics processing units (GPUs) and a conventional server equipped with 1 RTX 3090 Ti 24 GB GPU. The in-house EHR dataset is stored on a workstation located within a hospital research laboratory. The workstation operates within a Health Insurance Portability and Accountability Act–compliant network, ensuring the confidentiality, integrity, and availability of electronic protected health information, and it is equipped with a single NVIDIA V100 32 GB GPU. To use ChatGPT, we used an in-house ChatGPT-3.5-Turbo version hosted on our local cloud infrastructure. No data were sent to Microsoft or OpenAI. This setup ensured that no data were transmitted to OpenAI or external websites, and we were in strict compliance with the MIMIC data use agreement.

While GPT can handle 4096 tokens, T5 is limited to 512 tokens. To ensure a fair comparison, we focused on the subjective and assessment sections of progress notes as input. These sections provide physicians' evaluations of patients' conditions and fall within T5's 512-token limit. This differs from the objective sections, which mainly contain numerical values. Detailed information on data preprocessing, T5 model fine-tuning, and GPT zero-shot setting is presented in Multimedia Appendix 1.

Prompting Foundational Models to Integrate Graph Knowledge

To incorporate graph model–predicted paths into a prompt, we applied a prompt engineering strategy using domain-independent prompt patterns, as delineated in the study by White et al [38]. Our prompt was constructed with 3 primary components: the output customization that specifies the persona; the output format and template; and the context-control patterns, which are directly linked to the input note and the output of DR.KNOWS. In our test set, for the few input EHRs where no paths could be found (<20 instances), we directly fed the input into the LLMs (T5 and ChatGPT) to generate diagnoses.

Given that our core objective was to assess the extent to which the prompt can bolster the model's performance, it became imperative to test an array of prompts. Gonen et al [39] presented a technique, BETTERPROMPT, which relied on "selecting prompts by estimating language model likelihood." Essentially, we initiated the process with a set of manual task-specific prompts, subsequently expanding the prompt set via automatic paraphrasing facilitated by ChatGPT and backtranslation. We then ranked these prompts by their perplexity score (averaged over a representative sample of task inputs), ultimately selecting those prompts that exhibited the

lowest perplexity. Guided by this framework, we manually crafted 5 sets of prompts to integrate the path input, which are visually represented in Table S1 in Multimedia Appendix 1. Specifically, the first 3 prompts were designed by a non-medical domain expert (computer scientist), whereas the final 2 sets of prompts were developed by a medical domain expert (a critical care physician and a medical informaticist). We designated the last 2 prompts (with the medical persona) as "subject matter prompts" and the first 3 prompts as "non-subject matter prompts."

The chosen final prompt came from a template with minimal perplexity, incorporating predicted knowledge paths from the DR.KNOWS model as part of the input. We explored 2 path representation methods: "structural," which uses " \rightarrow " to link source concepts, edges (relation names), and target concepts; and "clause," which converts paths into clause-style text by directly joining the source and target concepts with their relations. Preliminary experiments showed superior performance with the "structural" representation, leading to its exclusive use in our reported results. The final prompt selected for the foundational models is a paraphrased prompt from the subject matter expert-crafted prompt: "Imagine you are a medical professional equipped with a knowledge graph, and generate the top three direct and indirect diagnoses from the input note. <Input note>...These are knowledge paths: <path 1>; <path</p> 2>...Separate the diagnoses using semicolons, and explain your reasoning starting with <Reasoning>." For the setup where the input did not contain paths, we simply used the prompt with the medical persona and task description as follows: "Imagine you are a medical professional, and generate the top three direct and indirect diagnoses from the input note. <Input note>..." The manually crafted prompts, their paraphrased versions, and their perplexity scores are presented in Table S1 in Multimedia Appendix 1.

Evaluation Metrics

Automated Evaluation Metrics for Quantitative Analysis

We conducted 2 evaluations for the DR.KNOWS models: the first was an intrinsic evaluation to determine how many gold standard concepts the graph model can retrieve. The second evaluation examined whether the retrieved knowledge paths could enhance the LLM's diagnosis prediction task. Regarding the first evaluation, our primary objective was to evaluate the effectiveness of DR.KNOWS in predicting diagnoses using CUIs. We used a concept extractor to analyze text within the plan section, specifically extracting CUIs classified under the semantic type T047 DISEASE AND SYNDROMES. We only included CUIs that were guaranteed to connect with at least 1 path, having a maximum length of 2 hops between the target and input CUIs. These chosen CUIs constituted the "gold standard" CUI set, used for both training and assessing the model's performance. As DR.KNOWS predicts the top N CUIs, we measured the Recall@N and Precision@N as follows:



The *F*-score, the harmonic mean between recall and precision, will also be reported.

To evaluate foundational model performance on EHR diagnosis prediction, we applied the aforementioned evaluation metric as well as Recall-Oriented Understudy for Gisting Evaluation (ROUGE) [40]. Specifically, ROUGE is a widely used set of metrics designed for evaluating the quality of machine-generated text by comparing it to reference texts. We used the ROUGE–Longest Common Subsequence (ROUGE-L) variant, which is based on the longest common substring; and the ROUGE-2 variant, which focuses on bigram matching. Both ROUGE metrics were used in the ProbSum shared task.

For reporting results from automated metrics, we provided the mean scores across all samples in the test set, along with 95% CIs on 1000 bootstrapped samples.

Human Evaluation for Qualitative Analysis

Existing evaluation frameworks for AI, such as those used in radiology report generation, do not address diagnosis prediction with LLMs, leaving a significant gap. To address this, our prior work introduced a new human evaluation framework based on the Safer DX Instrument [41], aiming to provide a structured approach for assessing LLMs in diagnosis tasks. In this study, we used this framework to assess the impact of knowledge paths on LLM diagnostic predictions, specifically through a qualitative analysis of the "reasoning" output by LLMs, aiming to gauge the depth and accuracy of the models' diagnostic reasoning processes.

Specifically, we evaluated the model-generated "reasoning" section on the following aspects: (1) reading comprehension, (2) rationale, (3) recall of knowledge, (4) omission of diagnostic reasoning, and (5) abstraction and effective abstraction. Reading comprehension was intended to capture whether a model understood the information in a progress note. Rationale was intended to capture the inclusion of incorrect reasoning steps. *Recall of knowledge* was intended to capture the hallucination of incorrect facts as well as the inclusion of irrelevant facts in the output. Omission of a diagnosis served the same purpose as noted previously by capturing instances when the model failed to support conclusions or provide evidence for a diagnostic choice. Abstraction and effective abstraction were intended to evaluate the amount of abstraction present in each part of the output. This was to ascertain how the knowledge paths influenced the type of output produced and whether the model was able to use abstraction. Omission as well as abstraction and effective abstraction were formatted as yes or no questions. Reading comprehension, rationale, and recall of knowledge were assessed on a Likert scale ranging from 1 to 5, with 1 indicating strong agreement with poor quality and 5 indicating strong disagreement (representing high quality).

We recruited 2 medical professionals to evaluate LLM outputs using human evaluation guidelines developed by us. Full details of the guidelines, evaluation training, and interannotator agreement are reported in a separate publication (currently under review). The evaluation framework used the REDCap (Research Electronic Data Capture; Vanderbilt University) web application to present the evaluators with input notes, gold standard

diagnoses, and model-predicted diagnoses. The evaluators, treated as separate arms in a longitudinal framework, assessed models with KG paths and those without across 2 defined events. Detailed step-by-step guidelines were provided for completing the evaluations in REDCap.

Two senior board-certified clinical informatics physicians served as advisors, pilot testers, and trainers for the 2 medical professionals who completed the human evaluations. The 2 physicians used 5 samples cases to iteratively refine the guidelines provided to the evaluators; these sample evaluations also served as examples for the evaluators to reference during training. The evaluation guidelines consisted of clear descriptions of the meaning of evaluative scores for each aspect of the human evaluation framework as well as a completed example workflow.

Results

Intrinsic Evaluation of DR.KNOWS on Predicting Diagnostic Concepts

We compared DR.KNOWS with QuickUMLS, which is a concept extractor baseline that identifies medical concepts from raw text. We took input text, parsed it with QuickUMLS, and outputted a list of concepts. Table 2 presents results from the 2 EHR datasets, MIMIC and in-house. The selection of different

top N values was determined by the disparity in text length between the 2 datasets. DR.KNOWS demonstrated superior precision and F-scores compared to QuickUMLS across both datasets compared to the baseline, with precision scores of 19.10 (95% CI 17.82-20.37) versus 13.59 (95% CI 12.32-14.88) on the MIMIC dataset and 22.88 (95% CI 20.92-24.85) versus 12.38 (95% CI 11.09-13.66) on the in-house dataset. In addition, its F-scores of 25.20 (95% CI 23.93-26.48) on the MIMIC dataset and 25.70 (95% CI 24.06-27.37) on the in-house dataset exceeded the comparison scores of 21.13 (95% CI 19.85-22.41) and 20.09 (95% CI 18.81-21.37), respectively, underscoring the effectiveness of DR.KNOWS in accurately predicting diagnostic CUIs. The TriAttn variant of DR.KNOWS consistently outperformed the MultiAttn variant on both datasets, with F-scores of 25.20 (95% CI 23.93-26.48) versus 23.10 (95% CI 21.83-24.39) on the MIMIC dataset and 25.70 (95% CI 24.06-27.37) versus 17.69 (95% CI 16.40-18.96) on the in-house dataset. The concept extractor baseline achieved the highest recall scores-56.91 on the MIMIC dataset and 90.11 on the in-house dataset-because it identified all input concepts that overlapped with the reference CUIs, in particular on the in-house dataset, which was largely an extractive dataset. Training the DR.KNOWS model took an average of 2 of 3 (SD 1.22) hours per epoch on 5000 samples, using 8000 MB of GPU memory.

Table 2. Performance comparison between concept extraction and 2 variants of DR.KNOWS on target concept unique identifier prediction using the Medical Information Mart for Intensive Care (MIMIC-III) and in-house datasets.

Model	MIMIC-III				In-house			
	Top N knowl- edge paths	Recall score (95% CI)	Precision score (95% CI)	<i>F</i> -score (95% CI)	Top N knowl- edge paths	Recall score (95% CI)	Precision score (95% CI)	<i>F</i> -score (95% CI)
Concept extrac- tor	a	56.91 (55.62- 58.18)	13.59 (12.32- 14.88)	21.13 (19.85- 22.41)	_	90.11 ^b (88.84- 91.37)	12.38 (11.09- 13.66)	20.09 (18.81- 21.37)
MultiAttn ^c	4	26.91 (25.64- 28.19	22.79 (21.51- 24.06)	23.10 (21.83- 24.39)	6	24.68 (23.35- 25.91)	15.82 (14.55- 17.10)	17.69 (16.40- 18.96)
MultiAttn	6	29.14 (27.85- 30.41)	16.73 (15.46- 18.00)	19.94 (18.66- 21.22)	8	28.69 (27.43- 29.98)	15.82 (14.55- 17.11)	17.33 (16.06- 18.60)
TriAttn ^d	4	29.85 (26.23- 33.45)	17.61 (16.33- 18.89)	20.93 (19.67- 22.21)	6	34.00 (31.04- 36.97)	22.88 (20.92- 24.85)	23.39 (21.71- 25.06)
TriAttn	6	37.06 (35.80- 38.33)	19.10 (17.82- 20.37)	25.20 (23.93- 26.48)	8	44.58 (41.38- 47.78)	22.43 (20.62- 24.23)	25.70 (24.06- 27.37)

^aNot applicable.

^bBest performance values are italicized.

^cMultiAttn: multihead attention.

^dTriAttn: trilinear attention.

Assessing the Impact of DR.KNOWS on Diagnosis Prediction

The best systems for each foundational model on the ProbSum test set are presented in Table 3, including those with predicted paths provided by DR.KNOWS and those without. Overall, the prompt-based fine-tuning of T5 surpassed ChatGPT's prompt-based zero-shot approach on all metrics, and ChatGPT's prompt-based few-shot approach showed comparable

performance to T5. Notably, models that incorporated paths, particularly for the CUI *F*-score, showed significant improvements. The vanilla T5 model with a path prompt excelled, achieving the highest ROUGE-L score (30.72, 95% CI 30.40-32.44) and CUI *F*-score (27.78, 95% CI 27.09-29.80). This ROUGE-L score could have ranked third on the ProbSum leaderboard [27], which is noteworthy considering that the top 2 systems used ensemble methods [10,11].

```
https://ai.jmir.org/2025/1/e58670
```

Table 3. Best performance on the Medical Information Mart for Intensive Care III (MIMIC III) test set (with annotated active diagnoses) from 3 Text-to-Text Transfer Transformer (T5) variants and ChatGPT across all prompt styles with DR.KNOWS (Diagnostic Reasoning Knowledge Graph System) path prompting and without. To illustrate the performance differences better, we report Recall-Oriented Understudy for Gisting Evaluation-2 (ROUGE-2); ROUGE–Longest Common Subsequence (ROUGE-L); and concept unique identifier (CUI) recall, precision, and F-scores.

Model	Rouge-2 score (95% CI)	Rouge-L score (95% CI)	CUI recall score (95% CI)	CUI precision score (95% CI)	CUI F-score (95% CI)
Prompt-based fine-tuning setting	•				
Vanilla T5	12.66 (11.24-13.54)	29.08 (27.52-29.99)	39.17 (37.53-41.56)	22.89 (21.02-23.62)	26.19 (25.31-26.78)
Vanilla T5+path ^a	13.13 (12.64-13.88)	<i>30.72^b</i> (30.40- 32.44 ^c)	40.73 (39.46-42.18)	24.28 (23.49-26.03)	27.78 (27.08-29.80 ^c)
Flan-T5	11.83 (10.51-12.40)	27.02 (25.64-27.80)	38.28 (36.70-39.45)	22.32 (21.81-23.00)	25.32 (24.10-26.34)
Flan-T5+path	13.30 (12.19-14.44)	30.00 (29.20-32.70 ^c)	38.96 (37.48-40.01)	24.74 (23.35-26.12 ^c)	27.38 (26.98-28.68 ^c)
Clinical-T5	11.68 (11.06-12.49)	25.84 (23.74-26.15)	30.37 (28.94-30.99)	17.91 (15.46-19.79)	19.61 (16.44-20.03)
Clinical-T5+path	12.06 (10.89-12.48)	25.97 (24.71-26.33)	29.45 (27.65-30.19)	22.78 (21.35-23.59 ^c)	23.17 (21.39-23.96 ^c)
Prompt-based zero-shot setting					
ChatGPT	7.05 (6.54-7.56)	19.77 (19.26-20.28)	23.68 (23.18-24.19)	15.52 (15.00-16.02)	16.04 (15.53-16.55)
ChatGPT+path	5.70 (5.19-6.21)	15.49 (14.98-15.99)	25.33 (24.82-25.84 ^c)	17.05 (16.29-17.81 ^c)	18.21 (17.46-18.98 ^c)
Prompt-based few-shot setting					
ChatGPT 3-shot	9.63 (8.32-10.06)	21.84 (19.99-22.09)	22.71 (20.99-23.96)	19.57 (17.23-19.78)	21.02 (20.26-21.79)
ChatGPT 5-shot	9.73 (8.52-10.18)	21.23 (19.58-21.72)	22.45 (20.93-23.80)	19.67 (17.66-20.33)	20.96 (20.19-21.73)
ChatGPT 3-shot+path	10.66 (9.17-10.72)	24.32 (22.44-24.25 ^c)	26.48 (25.33-28.36 ^c)	24.22 (21.44-24.21 ^c)	25.30 (24.52-26.06 ^c)
ChatGPT 5-shot+path	11.73 (10.51-12.25 ^c)	25.43 (23.53-25.35 ^c)	27.76 (26.56-29.39 ^c)	24.56 (22.47-25.12 ^c)	26.02 (25.25-26.78 ^c)

^aPrompt styles with DR.KNOWS path prompting.

^bBest performance values are italicized.

^c95% CIs with a distinct CI for the DR.KNOWS-prompted path compared to no-path scenarios.

The comparison between ChatGPT with DR.KNOWS and ChatGPT without in the predicted paths scenario provided additional insights. In the few-shot setting, the incorporation of paths led to marked improvements; for instance, in the 3-shot setting, the with-path scenario outperformed the no-path scenario on all metrics, with ROUGE-L score of 24.32 (95% CI 22.44-24.25) compared to ChatGPT 3-shot no-path ROUGE-L score of 21.84 (95% CI 19.44-22.09) and CUI *F*-score of 25.30 (95% CI 24.52-26.06) versus 21.02 (95% CI 20.26-21.79). In the 5-shot setting, ChatGPT with paths achieved a ROUGE-L score of 25.43 (95% CI 25.53-25.35) compared to 21.23 (95% CI 19.58-21.72) for ChatGPT without paths and CUI *F*-score of 26.02 (95% CI 25.25-26.78) versus 20.96 (95% CI 20.19-21.73).

Human Evaluation Results

After the annotation procedure, the 2 medical professionals completed a supervised set of evaluations and were considered validated once they achieved a κ coefficient of 0.7 with the physician trainers and each other.

Although the T5 and ChatGPT models displayed similar performance on automated metrics, their outputs diverged

significantly. The T5 models, lacking instruction tuning, failed to respond adequately to prompts requesting the generation of a <Reasoning> section. Consequently, our human evaluation focused exclusively on the outputs produced by ChatGPT. We conducted human evaluation of the top-performing ChatGPT output (5-shot approach), comparing scenarios with the DR.KNOWS knowledge paths with KG and without KG. The final evaluation set consisted of 92 input notes and 2 sets of ChatGPT-predicted text.

The results are reported in Figure 4. First, there was no significant increase in *omission of diagnoses*, with 16% (15/92) observed with KG as opposed to 10% (9/92) without KG (P=.16). Regarding *rationale* (correct reasoning), ChatGPT with KG exhibited stronger agreement with the human evaluators (51/92, 55%) than ChatGPT without KG (46/92, 50%; P<.001). In the *abstraction* category (assessing the presence of abstraction in the model output), there was a notable drop from 88% (81/92; without KG) to 78% (71/92; with KG) in the affirmative responses (P=.03), indicating that less abstraction was required when KG paths were included. Differences were also noted in *effective abstraction* in favor of the KG paths (P=.002).

Gao et al



Figure 4. Human evaluation of ChatGPT outputs comparing scenarios with ("KG" [knowledge graph]) the DR.KNOWS (Diagnostic Reasoning Knowledge Graph System) knowledge paths and without ("No KG").

Error Analysis

We discovered 2 primary types of errors in the DR.KNOWS outputs that could result in missed opportunities for improving knowledge grounding. Figure 5 presents an example where ChatGPT did not find the provided knowledge paths useful. In this case, the majority of the provided knowledge paths were

highly extractive ("leukocytosis," "reticular dysgenesis," and "paraplegia" are the target concepts to which the knowledge paths led, and all are associated with a "self-loop" relationship). On the abstraction paths, the retrieved target concepts "abdomen hernia scrotal" and "chronic neutrophilia" were not relevant to the input patient condition.

Figure 5. An example of an error in the knowledge paths retrieved by DR.KNOWS (Diagnostic Reasoning Knowledge Graph System). DR.KNOWS retrieved 2 paths leading to irrelevant and misleading diagnoses (marked in red). The counterclockwise gapped circular arrow symbol represents a self-loop.

Input progress note:

-Assessment> 73 yo M w/ mmp, C4-5 paraplegia, TF dependence, on broad spectrum abx for recent pna transferred from OSH w/ resp distress, leukocytosis, and HOTN <Subjective> Chief Complaint: resp distress and hypotension I saw and examined the patient, and was physically present with the ICU Resident for key portions of the services provided. I agree with his / her note above, including assessment and plan. HPI: 73 yo M w/mmp including C4-5 paraplegia, TF dependence, on broad spectrum abx for recent pna transferred from OSH w/ resp distress, leukocytosis, and HOTN. 24 Hour Events: PICC LINE - START 04:24 PM ARTERIAL LINE - START 07:00 PM HOTN responded to IVF, never required pressor support MS ANTIBX coverage broadened--> vanco/ CT chest ordered Leukocytosis normalizing History obtained from Medical records Allergies: Methyldopa hives; Shellfish pt. with remote.

Dr.Knows retrieved top-6 knowledge paths:

Leukocytosis \rightarrow self \rightarrow Leukocytosis σ <path> reticular dysgenesis \rightarrow self \rightarrow reticular dysgenesis σ <path> Leukocytosis \rightarrow definitional manifestation of \rightarrow Leukocytosis σ <path> Paraplegia \rightarrow self \rightarrow Paraplegia σ <path> Thoracic \rightarrow has finding site \rightarrow abdomen hernia scrotal σ <path> Leukocytosis \rightarrow definitional manifestation of \rightarrow Leukocytosis \rightarrow has definitional manifestation \rightarrow Chronic neutrophilia

Another error observed occurred when DR.KNOWS selected the source CUIs that were less likely to generate pertinent paths for clinical diagnoses, resulting in ineffective knowledge paths. Figure 6 shows a retrieved path from "consulting with (procedure)" to "consultation-action (qualifier value)." Although some procedure-related concepts such as endoscopy or blood testing were valuable for clinical diagnosis, this specific path of consulting did not contribute meaningfully to the input case. Similarly, another erroneous pathway began with "drug allergy" and led to "allergy to dimetindene (finding)," which is contradictory, given that the input note explicitly states "no known drug allergies." While the consulting path's issue was its lack of utility, the "drug allergy" path could introduce the

risk of hallucination (misleading or fabricated content) within ChatGPT.

Figure 6. An example illustrating ChatGPT's performance with the knowledge paths extracted by DR.KNOWS (Diagnostic Reasoning Knowledge Graph System). Two paths had source concept unique identifiers ("Consulting with [procedure]" and "Drug allergy") that were less likely to generate pertinent paths for clinical diagnoses. Of note, the path of "Drug allergy" led to a path contradicting the "No Known Drug Allergies" description in the input. The path of "cirrhosis of liver" represents a correct diagnosis, but ChatGPT failed to include it. The counterclockwise gapped circular arrow symbol represents a self-loop. ESRD: end-stage renal disease.

i nput progress note: <assessment> 57M with Hep C cirrhosis, ESRD on HD, presenting with hypotension and shock, elevated lactate, and drop in hematocrit. <subjective> TITLE: Chief Complaint: Hypotension 24 Hour Events: - Levophed not able to be weaned - PT consult - Ordered VBG with O2 sat an lactate to evaluate whether he's ischemic during all this hypoTN <mark>Allergies: No Known Drug Allergies</mark></subjective></assessment>
Dr.Knows retrieved top-6 knowledge paths:
Unspecified chronic renal failure \rightarrow possibly equivalent to \rightarrow Renal failure: [chronic] or [end stage] \rightarrow possibly equivalent to \rightarrow Unspecified chronic renal failure <path> and the stage of the s</path>
Cirrhosis of liver (disorder) \rightarrow self \rightarrow Cirrhosis of liver (disorder) σ <path></path>
Allergic reaction (disorder) \rightarrow self \rightarrow Allergic reaction (disorder) σ <path></path>
Unspecified chronic renal failure → possibly equivalent to → Renal failure: [chronic] or [end stage] σ <path></path>
Consulting with (procedure) \rightarrow method of \rightarrow Consultation - action (qualifier value) σ <path></path>
Drug allergy \rightarrow has definitional manifestation \rightarrow Allergy to dimetindene (finding) σ
Gold Standard Diagnosis:
Hypotension/shock. Most Likely septic shock; ESRD; Cirrhosis
Predicted Diagnoses (with knowledge paths input):
ESRD with hypotension and shock; elevated lactate; drop in hematocrit.

In addition to the errors in the DR.KNOWS outputs, there were instances where ChatGPT failed to leverage the accurate knowledge paths presented. Figure 6 includes a knowledge path regarding "cirrhosis of liver," which was the correct diagnosis. However, ChatGPT response did not include this diagnosis.

Discussion

Principal Findings

DR.KNOWS showed significant advantages over the QuickUMLS concept extractor baseline in extracting correct concepts for diagnoses. On the ProbSum dataset, where the goal was to generate a list of diagnoses given the progress notes, prompt-based fine-tuning of T5 outperformed ChatGPT's zero-shot approach and showed comparable results to its few-shot approaches, with the inclusion of predicted paths by DR.KNOWS significantly enhancing performance across all metrics. The vanilla T5 with path prompts notably achieved top ROUGE-L and CUI *F*-scores, demonstrating the effectiveness of incorporating paths into the model. Human evaluation of ChatGPT's reasoning section showed strong agreement with human evaluators in terms of correct *rationale* and enhanced *effective abstraction*, indicating nuanced improvement in reasoning and abstraction quality with KG integration.

While DR.KNOWS leverages KG paths to enhance diagnosis prediction, it is important to acknowledge the potential biases and limitations inherent in KG data. KGs such as UMLS are comprehensive, but they may reflect biases based on the clinical domains and patient populations from which they were constructed, which could impact the relevance or appropriateness of the retrieved paths. To mitigate this, DR.KNOWS focuses on case-specific path selection, aiming to retrieve only the paths most directly relevant to the patient context. Nonetheless, future iterations could benefit from

```
https://ai.jmir.org/2025/1/e58670
```

RenderX

evaluating path relevance using additional contextual information, such as demographic details, to better align with patient-specific needs and reduce bias.

Error analysis showed that DR.KNOWS occasionally struggled with identifying knowledge paths unrelated to the patient representation; in addition, the analysis emphasized the importance of selecting accurate starting medical concepts. Currently, DR.KNOWS relies solely on semantic-based ranking on the candidate paths, that is, the cosine similarity between candidate path embeddings and input text, with the embedding quality being crucial for ranking performance. Improving the representation and embedding methods, as well as exploring probabilistic modeling techniques [42,43], could enhance path relevance. Furthermore, incorporating a graph reasoning mechanism that enables symbolic chain-of-thought reasoning might compensate for the weaknesses of contextualized embeddings and cosine-similarity metrics [44], presenting a valuable future direction. This integration could improve the diagnostic potential of DR.KNOWS, allowing for more nuanced and bias-aware reasoning.

The error analysis also presented instances where ChatGPT neglected to incorporate certain beneficial knowledge paths. It is important to acknowledge that ChatGPT operates as a black box application programming interface model, with its internal weights and training processes being inaccessible. To enhance the efficacy of the graph-based retrieve-and-augment framework, it would be advantageous to explore the potential of graph prompting and instruction tuning on open-source language models. These methods could refine the model's ability to use relevant information effectively. Other relevant research also uses advanced prompting techniques, such as self-retrieval-augmented generation [45] and step-back prompting [46]. The Google Research team recently presented a study investigating multiple ways of encoding graphs into

LLM inputs [47], which might inform a future direction for this work beyond the typical structural or clause-based path prompting.

In conclusion, LLMs such as ChatGPT hold promise for generating diagnoses for clinical decision support; however, methods such as graph prompting are needed to guide the model down the correct reasoning paths to avoid hallucinations and provide comprehensive diagnoses. While we show some progress in a graph prompting approach with DR.KNOWS, more work is needed to improve methods that leverage the UMLS knowledge source for grounding to achieve more accurate outputs. Nonetheless, DR.KNOWS represents a step toward trustworthy AI in medicine, providing knowledge grounding to LLMs and potentially reducing factual errors in diagnostic outputs [48]. Furthermore, our proposed human evaluation framework, derived from diagnostic safety evaluations used in clinical settings, enables the assessment of LLMs from the perspective of diagnostic safety. It carries strong face validity and reliability to evaluate a model's strengths and weaknesses as a diagnostic decision support system. This ensures that the models not only perform well on technical metrics but also align with clinical standards of safety and reliability.

Limitations

Our work on leveraging KGs for LLM diagnosis generation has shown promising results; however, there are notable limitations

that must be acknowledged. First, while the UMLS concept extractors (Clinical Text Analysis and Knowledge Extraction System and QuickUMLS) are powerful tools, they are not without flaws. One significant limitation is their inability to accurately identify all relevant concepts, particularly indirect or nuanced medical concepts. These indirect concepts can be crucial for accurate diagnosis generation; yet, the current concept extractors may fail to recognize them, leading to incomplete or less accurate knowledge representation.

Second, our path selection methodology relies heavily on cosine similarity, a common approach within the retrieval-augmented generation framework. Despite its prevalence, this method has inherent limitations due to its heavy reliance on the quality of embedding representations. If the embeddings do not adequately capture the semantic nuances of medical concepts, the similarity measure may lead to the retrieval of less relevant or noisy knowledge paths. This can ultimately impact the quality and reliability of the diagnostic suggestions generated by the LLM.

These limitations highlight the need for the continued refinement of both the concept extraction and path selection processes. Future work should explore more sophisticated techniques to enhance concept identification and improve the robustness of embedding representations, thereby reducing the reliance on cosine similarity and increasing the overall accuracy and utility of the KG-based approach.

Acknowledgments

This work is supported by grants from the National Institutes of Health. Funding was supported by the National Library of Medicine (K99LM014308, R00LM014308: YG; R01LM012973-04: TM and DD); the National Heart, Lung, and Blood Institute (R01HL157262-03: MMC); and the National Institute on Drug Abuse (R01DA051464: MA).

Data Availability

The source code knowledge graph generated during this study are available on the GitHub repository [49]. Medical Information Mart for Intensive Care III is available from PhysioNet.

Authors' Contributions

YG was responsible for conceptualization, supervision, methodology, formal analysis, writing (original draft as well as review and editing), validation, visualization, data curation, investigation, project administration, and funding acquisition. RL was responsible for writing (original draft as well as review and editing), methodology, data curation, validation, investigation, conceptualization, and formal analysis. EC was responsible for writing (original draft as well as review and editing), validation, methodology, data curation, investigation, conceptualization, and formal analysis. JRC was responsible for writing (review and editing), formal analysis, investigation, and data curation. BWP was responsible for writing (review and editing), validation, formal analysis, methodology, investigation, and conceptualization. MMC was responsible for writing (review and editing), conceptualization, methodology, and funding acquisition. TM was responsible for writing (review and editing), conceptualization, methodology, and funding acquisition. DD was responsible for writing (review and editing), conceptualization, methodology, and funding acquisition. DD was responsible for writing (review and editing), conceptualization, methodology, and funding acquisition, by writing (review and editing), conceptualization, methodology, and funding acquisition. DD was responsible for writing (review and editing), conceptualization, and funding acquisition. MA was responsible for conceptualization, supervision, methodology, formal analysis, writing (original draft as well as review and editing), validation, visualization, data curation, investigation, project administration, and funding acquisition.

Conflicts of Interest

TM is a consultant for Lavita.ai, a startup that builds NLP tools for medical use cases. All other authors declare no conflicts of interest.



Multimedia Appendix 1

Data preprocessing, DR.KNOWS (Diagnostic Reasoning Knowledge Graph System) training details, prompt engineering using ChatGPT, and Text-to-Text Transfer Transformer (T5) fine-tuning.

[DOCX File, 37 KB - ai_v4i1e58670_app1.docx]

References

- 1. Brown PJ, Marquard JL, Amster B, Romoser M, Friderici J, Goff S, et al. What do physicians read (and ignore) in electronic progress notes? Appl Clin Inform 2017 Dec 21;05(02):430-444. [doi: <u>10.4338/aci-2014-01-ra-0003</u>]
- Rule A, Bedrick S, Chiang MF, Hribar MR. Length and redundancy of outpatient progress notes across a decade at an academic medical center. JAMA Netw Open 2021 Jul 01;4(7):e2115334 [FREE Full text] [doi: 10.1001/jamanetworkopen.2021.15334] [Medline: 34279650]
- 3. Liu J, Capurro D, Nguyen A, Verspoor K. "Note Bloat" impacts deep learning-based NLP models for clinical prediction tasks. J Biomed Inform 2022 Sep;133:104149 [FREE Full text] [doi: 10.1016/j.jbi.2022.104149] [Medline: 35878821]
- 4. Nijor S, Rallis G, Lad N, Gokcen E. Patient safety issues from information overload in electronic medical records. J Patient Saf 2022 Sep 01;18(6):e999-1003 [FREE Full text] [doi: 10.1097/PTS.000000000001002] [Medline: 35985047]
- 5. Furlow B. Information overload and unsustainable workloads in the era of electronic health records. Lancet Respir Med 2020 Mar;8(3):243-244. [doi: 10.1016/S2213-2600(20)30010-2] [Medline: 32135094]
- 6. Croskerry P. Diagnostic failure: a cognitive and affective approach. In: Henriksen K, Battles JB, Marks ES, Lewin DI, editors. Advances in Patient Safety: From Research to Implementation. Volume 2. New York, NY: Agency for Healthcare Research and Quality; 2005:241-254.
- Gao Y, Dligach D, Miller T, Xu D, Churpek MM, Afshar M. Summarizing patients' problems from hospital progress notes using pre-trained sequence-to-sequence models. In: Proceedings of the 29th International Conference on Computational Linguistics. 2022 Presented at: COLING '22; October 12-17, 2022; Virtual Event p. 2979-2991.
- 8. Johnson AE, Pollard TJ, Shen L, Lehman LH, Feng M, Ghassemi M, et al. MIMIC-III, a freely accessible critical care database. Sci Data 2016 May 24;3(1):160035 [FREE Full text] [doi: 10.1038/sdata.2016.35] [Medline: 27219127]
- 9. Gao Y, Dligach D, Miller T, Churpek MM, Afshar M. Overview of the problem list summarization (ProbSum) 2023 shared task on summarizing patients' active diagnoses and problems from electronic health record progress notes. Proc Conf Assoc Comput Linguist Meet 2023 Jul;2023:461-467 [FREE Full text] [doi: 10.18653/v1/2023.bionlp-1.43] [Medline: 37583489]
- Manakul P, Fathullah Y, Liusie A, Raina V, Raina V, Gales M. CUED at ProbSum 2023: hierarchical ensemble of summarization models. In: Proceedings of the 22nd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks. 2023 Presented at: BioNLP '23; July 13, 2023; Toronto, ON p. 516-523 URL: <u>https://aclanthology.org/2023.</u> <u>bionlp-1.51.pdf</u> [doi: <u>10.18653/v1/2023.bionlp-1.51</u>]
- 11. Li H, Wu Y, Schlegel V, Batista-Navarro R, Nguyen TT, Kashyap RA, et al. Team:PULSAR at ProbSum 2023:PULSAR: pre-training with extracted healthcare terms for summarising patients' problems and data augmentation with black-box large language models. In: Proceedings of the 22nd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks. 2023 Presented at: BioNLP '23; July 13, 2023; Toronto, ON p. 503-509. [doi: 10.18653/v1/2023.bionlp-1.49]
- 12. Raffel C, Shazeer N, Roberts A, Lee K, Narang S, Matena M, et al. Exploring the limits of transfer learning with a unified text-to-text transformer. J Mach Learn Res 2020;21(140):1-67 [FREE Full text]
- 13. Floridi L, Chiriatti M. GPT-3: its nature, scope, limits, and consequences. Minds Mach 2020 Nov 01;30(4):681-694. [doi: 10.1007/S11023-020-09548-1]
- 14. Baumgartner C. The potential impact of ChatGPT in clinical and translational medicine. Clin Transl Med 2023 Mar;13(3):e1206 [FREE Full text] [doi: 10.1002/ctm2.1206] [Medline: 36854881]
- 15. Bodenreider O. The Unified Medical Language System (UMLS): integrating biomedical terminology. Nucleic Acids Res 2004 Jan 01;32(Database issue):D267-D270 [FREE Full text] [doi: 10.1093/nar/gkh061] [Medline: 14681409]
- 16. Huang KH, Yang M, Peng N. Biomedical event extraction with hierarchical knowledge graphs. In: Proceedings of the 2020 Conference on Association for Computational Linguistics. 2020 Presented at: EMNLP '20; November 16-20, 2020; Virtual Event p. 1277-1285 URL: <u>https://aclanthology.org/2020.findings-emnlp.114.pdf</u> [doi: <u>10.18653/v1/2020.findings-emnlp.114</u>]
- Lu Q, Dou D, Nguyen TH. Parameter-efficient domain knowledge integration from multiple sources for biomedical pre-trained language models. In: Proceedings of the 2021 Conference on the Association for Computational Linguistics. 2021 Presented at: EMNLP '21; November 7-11, 2021; Virtual Event p. 3855-3865 URL: <u>https://aclanthology.org/2021.</u> <u>findings-emnlp.325.pdf</u> [doi: <u>10.18653/v1/2021.findings-emnlp.325</u>]
- Aracena C, Villena F, Rojas M, Dunstan J. A knowledge-graph-based intrinsic test for benchmarking medical concept embeddings and pretrained language models. In: Proceedings of the 13th International Workshop on Health Text Mining and Information Analysis. 2022 Presented at: LOUHI '22; December 7, 2022; Virtual Event p. 197-206 URL: <u>https://aclanthology.org/2022.louhi-1.22.pdf</u> [doi: <u>10.18653/v1/2022.louhi-1.22</u>]
- He B, Zhou D, Xiao J, Jiang X, Liu Q, Yuan N, et al. BERT-MK: integrating graph contextualized knowledge into pre-trained language models. In: Proceedings of the 2020 Conference on Association for Computational Linguistics. 2020 Presented at: EMNLP '20; November 16-20, 2020; Virtual Event p. 2281-2290 URL: <u>https://aclanthology.org/2020.findings-emnlp.</u> 207.pdf [doi: <u>10.18653/v1/2020.findings-emnlp.207</u>]

```
https://ai.jmir.org/2025/1/e58670
```

- 20. Pan S, Luo L, Wang Y, Chen C, Wang J, Wu X, et al. Unifying large language models and knowledge graphs: a roadmap. IEEE Trans Knowl Data Eng 2024 Jul;36(7):3580-3599. [doi: 10.1109/tkde.2024.3352100]
- 21. Bowen JL. Educational strategies to promote clinical diagnostic reasoning. N Engl J Med 2006 Nov 23;355(21):2217-2225. [doi: <u>10.1056/nejmra054782</u>]
- 22. Corazza GR, Lenti MV. Diagnostic reasoning in internal medicine. Cynefin framework makes sense of clinical complexity. Front Med (Lausanne) 2021 Apr 22;8:641093 [FREE Full text] [doi: 10.3389/fmed.2021.641093] [Medline: 33968954]
- Kanwal N, Rizzo G. Attention-based clinical note summarization. In: Proceedings of the 37th ACM/SIGAPP Symposium on Applied Computing. 2022 Presented at: SAC '22; April 25-29, 2022; Virtual Event p. 813-820 URL: <u>https://dl.acm.org/ doi/10.1145/3477314.3507256</u> [doi: <u>10.1145/3477314.3507256</u>]
- 24. Adams G, Alsentzer E, Ketenci M, Zucker J, Elhadad N. What's in a summary? Laying the groundwork for advances in hospital-course summarization. In: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. 2021 Presented at: NAACL '21; June 6-11, 2021; Virtual Event p. 4794-4811 URL: https://aclanthology.org/2021.naacl-main.382.pdf [doi: https://aclanthology.org/2021.naacl-main.382 [doi: <a href="https://aclant
- 25. Pivovarov R, Elhadad N. Automated methods for the summarization of electronic health records. J Am Med Inform Assoc 2015 Sep;22(5):938-947 [FREE Full text] [doi: 10.1093/jamia/ocv032] [Medline: 25882031]
- 26. Liang J, Tsou CH, Poddar A. A novel system for extractive clinical note summarization using EHR data. In: Proceedings of the 2nd Clinical Natural Language Processing Workshop. 2019 Presented at: ClinicalNLP '19; June 7, 2019; Minneapolis, MN p. 46-54 URL: <u>https://aclanthology.org/W19-1906/</u> [doi: <u>10.18653/v1/w19-1906</u>]
- Zhang J, Zhang X, Yu J, Tang J, Tang J, Li C, et al. Subgraph retrieval enhanced model for multi-hop knowledge base question answering. In: Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics. 2022 Presented at: ACL '22; May 22-27, 2022; Dublin, Ireland p. 5773-5784 URL: <u>https://aclanthology.org/2022.acl-long.396.</u> pdf [doi: 10.18653/v1/2022.acl-long.396]
- 28. Yasunaga M, Bosselut A, Ren H, Zhang X, Manning CD, Liang P, et al. Deep bidirectional language-knowledge graph pretraining. In: Proceedings of the 36th Annual Conference on Neural Information Processing Systems. 2022 Presented at: NIPS '22; November 28-December 9, 2022; New Orleans, LA p. 37309-37323 URL: <u>https://dl.acm.org/doi/10.5555/3600270.3602974</u>
- 29. Hu Z, Xu Y, Yu W, Wang S, Yang Z, Zhu C, et al. Empowering language models with knowledge graph reasoning for open-domain question answering. In: Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing. 2022 Presented at: EMNLP '22; December 7-11, 2022; Abu Dhabi, United Arab Emirates p. 9562-9581 URL: https://aclanthology.org/2022.emnlp-main.650.pdf [doi: 10.18653/v1/2022.emnlp-main.650]
- 30. Weed LL. Medical records, patient care, and medical education. Ir J Med Sci 2008 Oct 22;39(6):271-282. [doi: 10.1007/bf02945791]
- Savova GK, Masanz JJ, Ogren PV, Zheng J, Sohn S, Kipper-Schuler KC, et al. Mayo clinical Text Analysis and Knowledge Extraction System (cTAKES): architecture, component evaluation and applications. J Am Med Inform Assoc 2010 Sep 01;17(5):507-513 [FREE Full text] [doi: 10.1136/jamia.2009.001560] [Medline: 20819853]
- 32. Soldaini L, Goharian N. Quickumls: a fast, unsupervised approach for medical concept extraction. In: Proceedings of the 2016 Conference on Medical Information Retrieval. 2016 Presented at: MedIR '16; July 21, 2016; Pisa, Italy p. 1-4 URL: https://ir.cs.georgetown.edu/downloads/quickumls.pdf
- 33. Hou Y, Zhang J, Cheng J, Ma K, Ma RT, Chen H, et al. Measuring and improving the use of graph information in graph neural network. In: Proceedings of the 8th International Conference on Learning Representations. 2020 Presented at: ICLR '20; June 16-18, 2020; Addis Ababa, Ethiopia p. 1-16 URL: <u>https://openreview.net/pdf?id=rkeIIkHKvS</u>
- 34. Liu F, Shareghi E, Meng Z, Basaldella M, Collier N. Self-alignment pretraining for biomedical entity representations. In: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. 2021 Presented at: NAACL '21; June 6-11, 2021; Virtual Event p. 4228-4238 URL: <u>https://aclanthology.org/2021.naacl-main.334.pdf</u> [doi: 10.18653/v1/2021.naacl-main.334]
- 35. Conneau A, Kiela D, Schwenk H, Barrault L, Bordes A. Supervised learning of universal sentence representations from natural language inference data. In: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing. 2017 Presented at: EMNLP '17; September 7-11, 2017; Copenhagen, Denmark p. 670-680 URL: <u>https://aclanthology.org/D17-1070.pdf</u> [doi: 10.18653/v1/d17-1070]
- 36. Chung HW, Hou L, Longpre S, Zoph B, Tay Y, Fedus W, et al. Scaling instruction-finetuned language models. arXiv Preprint posted online October 20, 2022 [FREE Full text]
- 37. Lehman E, Johnson A. Clinical-T5: large language models built using MIMIC clinical text. PhysioNet. URL: <u>https://www.physionet.org/content/clinical-t5/1.0.0/</u> [accessed 2023-01-23]
- 38. White J, Fu Q, Hays S, Sandborn M, Olea C, Gilbert H, et al. A prompt pattern catalog to enhance prompt engineering with ChatGPT. arXiv Preprint posted online February 21, 2023 [FREE Full text]
- 39. Gonen H, Iyer S, Blevins T, Smith NA, Zettlemoyer L. Demystifying prompts in language models via perplexity estimation. In: Proceedings of the 2023 Conference of the Association for Computational Linguistics. 2023 Presented at: EMNLP '23; December 6-10, 2023; Singapore, Singapore p. 10136-10148 URL: <u>https://aclanthology.org/2023.findings-emnlp.679.pdf</u> [doi: <u>10.18653/v1/2023.findings-emnlp.679</u>]

```
https://ai.jmir.org/2025/1/e58670
```

- 40. Lin CY. ROUGE: a package for automatic evaluation of summaries. In: Lin CY, editor. Text Summarization Branches Out. Barcelona, Spain: Association for Computational Linguistics; 2004:74-81.
- Singh H, Khanna A, Spitzmueller C, Meyer AN. Recommendations for using the revised safer Dx instrument to help measure and improve diagnostic safety. Diagnosis (Berl) 2019 Nov 26;6(4):315-323 [FREE Full text] [doi: 10.1515/dx-2019-0012] [Medline: 31287795]
- 42. Rotmensch M, Halpern Y, Tlimat A, Horng S, Sontag D. Learning a health knowledge graph from electronic medical records. Sci Rep 2017 Jul 20;7(1):5994 [FREE Full text] [doi: 10.1038/s41598-017-05778-z] [Medline: 28729710]
- 43. Wan G, Du B. GaussianPath: a Bayesian multi-hop reasoning framework for knowledge graph reasoning. AAAI Conf Artif Intell 2021 May 18;35(5):4393-4401. [doi: <u>10.1609/aaai.v35i5.16565</u>]
- 44. Xu J, Fei H, Pan L, Liu Q, Lee M, Hsu W. Faithful logical reasoning via symbolic chain-of-thought. arXiv Preprint posted online May 28, 2024 [FREE Full text] [doi: 10.18653/v1/2024.acl-long.720]
- 45. Asai A, Wu Z, Wang Y, Sil A, Hajishirzi H. Self-RAG: learning to retrieve, generate, and critique through self-reflection. In: Proceedings of the 25th International Conference on Learning Representations. 2024 Presented at: ICLR '24; May 7-11, 2024; Vienna Austria p. 1-30 URL: <u>https://openreview.net/pdf?id=hSyW5go0v8</u>
- 46. Zheng HS, Mishra S, Chen X, Cheng HT, Chi EH, Le QV, et al. Take a step back: evoking reasoning via abstraction in large language models. arXiv Preprint posted online October 9, 2023 [FREE Full text]
- 47. Fatemi B, Halcrow J, Perozzi B. Talk like a graph: encoding graphs for large language models. In: Proceedings of the 25th International Conference on Learning Representations. 2024 Presented at: ICLR '24; May 7-11, 2024; Vienna Austria URL: https://openreview.net/attachment?id=IuXR1CCrSi&name=supplementary_material
- Xie Q, Schenck EJ, Yang HS, Chen Y, Peng Y, Wang F. Faithful AI in medicine: a systematic review with large language models and beyond. medRxiv Preprint posted online July 23, 2023 [FREE Full text] [doi: 10.1101/2023.04.18.23288752] [Medline: 37398329]
- 49. serenayj / DRKnows. GitHub. URL: <u>https://github.com/serenayj/DRKnows</u> [accessed 2024-04-29]

Abbreviations

AI: artificial intelligence CUI: concept unique identifier **DR.KNOWS:** Diagnostic Reasoning Knowledge Graph System EHR: electronic health record **GPT:** Generative Pretrained Transformer **GPU:** graphics processing unit KG: knowledge graph LLM: large language model MIMIC-III: Medical Information Mart for Intensive Care III MultiAttn: multihead attention **REDCap:** Research Electronic Data Capture **ROUGE:** Recall-Oriented Understudy for Gisting Evaluation **ROUGE-L:** Recall-Oriented Understudy for Gisting Evaluation–Longest Common Subsequence SapBERT: Self-alignment Pretrained Bidirectional Encoder Representations from Transformers SGIN: stack graph isomorphism network SNOMED CT: Systematized Nomenclature of Medicine-Clinical Terms **SOAP:** subjective, objective, assessment, and plan T5: Text-to-Text Transfer Transformer TriAttn: trilinear attention UMLS: Unified Medical Language System

Edited by H Liu; submitted 21.03.24; peer-reviewed by A Sheth, N Zhang, Y Hua; comments to author 17.06.24; revised version received 07.08.24; accepted 07.11.24; published 24.02.25.

<u>Please cite as:</u> Gao Y, Li R, Croxford E, Caskey J, Patterson BW, Churpek M, Miller T, Dligach D, Afshar M Leveraging Medical Knowledge Graphs Into Large Language Models for Diagnosis Prediction: Design and Application Study JMIR AI 2025;4:e58670 URL: <u>https://ai.jmir.org/2025/1/e58670</u> doi:<u>10.2196/58670</u> PMID:



©Yanjun Gao, Ruizhe Li, Emma Croxford, John Caskey, Brian W Patterson, Matthew Churpek, Timothy Miller, Dmitriy Dligach, Majid Afshar. Originally published in JMIR AI (https://ai.jmir.org), 24.02.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.

Original Paper

GPT-4 as a Clinical Decision Support Tool in Ischemic Stroke Management: Evaluation Study

Amit Haim Shmilovitch¹, MD; Mark Katson¹, MD; Michal Cohen-Shelly², MD; Shlomi Peretz^{3,4}, MD; Dvir Aran^{5,6*}, PhD; Shahar Shelly^{1,7*}, MD

¹Department of Neurology, Rambam Medical Center, Haifa, Israel

*these authors contributed equally

Corresponding Author:

Shahar Shelly, MD Department of Neurology Rambam Medical Center HaAliya HaShniya Street 8 PO Box 9602 Haifa, 3109601 Israel Phone: 972 543541995 Email: s shelly@rmc.gov.il

Abstract

Background: Cerebrovascular diseases are the second most common cause of death worldwide and one of the major causes of disability burden. Advancements in artificial intelligence have the potential to revolutionize health care delivery, particularly in critical decision-making scenarios such as ischemic stroke management.

Objective: This study aims to evaluate the effectiveness of GPT-4 in providing clinical support for emergency department neurologists by comparing its recommendations with expert opinions and real-world outcomes in acute ischemic stroke management.

Methods: A cohort of 100 patients with acute stroke symptoms was retrospectively reviewed. Data used for decision-making included patients' history, clinical evaluation, imaging study results, and other relevant details. Each case was independently presented to GPT-4, which provided scaled recommendations (1-7) regarding the appropriateness of treatment, the use of tissue plasminogen activator, and the need for endovascular thrombectomy. Additionally, GPT-4 estimated the 90-day mortality probability for each patient and elucidated its reasoning for each recommendation. The recommendations were then compared with a stroke specialist's opinion and actual treatment decisions.

Results: In our cohort of 100 patients, treatment recommendations by GPT-4 showed strong agreement with expert opinion (area under the curve [AUC] 0.85, 95% CI 0.77-0.93) and real-world treatment decisions (AUC 0.80, 95% CI 0.69-0.91). GPT-4 showed near-perfect agreement with real-world decisions in recommending endovascular thrombectomy (AUC 0.94, 95% CI 0.89-0.98) and strong agreement for tissue plasminogen activator treatment (AUC 0.77, 95% CI 0.68-0.86). Notably, in some cases, GPT-4 recommended more aggressive treatment than human experts, with 11 instances where GPT-4 suggested tissue plasminogen activator use against expert opinion. For mortality prediction, GPT-4 accurately identified 10 (77%) out of 13 deaths within its top 25 high-risk predictions (AUC 0.89, 95% CI 0.8077-0.9739; hazard ratio 6.98, 95% CI 2.88-16.9; P<.001), outperforming supervised machine learning models such as PRACTICE (AUC 0.70; log-rank P=.02) and PREMISE (AUC 0.77; P=.07).

Conclusions: This study demonstrates the potential of GPT-4 as a viable clinical decision-support tool in the management of acute stroke. Its ability to provide explainable recommendations without requiring structured data input aligns well with the

²Sagol AI Hub, ARC Innovation Center, Chaim Sheba Medical Center, Ramat Gan, Israel

³Department of Neurology, Shamir Medical Center, Be`er Ya`akov, Israel

⁴Sackler School of Medicine, Tel Aviv University, Tel Aviv, Israel

⁵Faculty of Biology, Technion-Israel Institute of Technology, Haifa, Israel

⁶The Taub Faculty of Computer Science, Technion-Israel Institute of Technology, Haifa, Israel

⁷Rapaport Faculty of Medicine, Technion – Israel Institute of Technology, Haifa, Israel

routine workflows of treating physicians. However, the tendency toward more aggressive treatment recommendations highlights the importance of human oversight in clinical decision-making. Future studies should focus on prospective validations and exploring the safe integration of such artificial intelligence tools into clinical practice.

(JMIR AI 2025;4:e60391) doi:10.2196/60391

KEYWORDS

GPT-4; ischemic stroke; clinical decision support; artificial intelligence; neurology

Introduction

The advent of GPT-4 [1], launched by OpenAI in March 2023, marked a significant milestone in the evolution of artificial intelligence (AI) and its applications in various domains, including health care. GPT-4, a model under the umbrella of GPT, exemplifies the advancement in large language model (LLM) technology [2,3]. The foundational architecture of this technology involves training on extensive datasets, enabling the model to function as a "few-shot learner." This capability allows GPT-4 to adapt to new domains and continuously refine its performance through ongoing learning [2,4-6].

In the realm of clinical medicine, the potential applications of LLMs like GPT-4 are particularly intriguing. These models offer promise as supportive tools for health care professionals, aiding in the efficient summarization of patient data, assisting in decision-making processes, and potentially improving the accuracy and speed of medical interventions [7,8]. Recent research has underscored the capabilities of GPT-4 in complex medical tasks [9]. Notably, the model has demonstrated proficiency in examinations akin to the United States Medical Licensing Examination, achieving scores that meet or nearly meet the passing thresholds [10]. Additionally, in assessments modeled after neurology board exam questions, GPT-4 has shown a high accuracy rate, improving with repeated attempts [9,11,12].

The management of acute ischemic stroke (AIS) presents a critical and time-sensitive challenge in clinical settings. The approach to diagnosing and treating AIS requires a synthesis of information including patient symptoms, physical and neurological examinations, medical history, and imaging results. Despite the availability of established guidelines by the American Heart Association/American Stroke Association for stroke management [13-16], the pivotal role of the treating physician's judgment remains. Variability in clinical presentations and the urgent need for decision-making underscore the potential value of AI-assisted tools in this context. Moreover, predicting early mortality in AIS is essential for guiding treatment decisions, optimizing resource allocation in health care settings, facilitating effective communication with patients and their families, supporting research and clinical trials, and contributing to quality improvement initiatives. In accordance, several traditional machine learning models have been trained for this task in recent years [17-20].

Here, we leveraged patient data from the emergency department (ED) of a large referral hospital, focusing on individuals presenting with stroke symptoms, to evaluate the effectiveness of GPT-4 in delivering accurate clinical decisions for the

https://ai.jmir.org/2025/1/e60391

treatment of AIS. We also assessed its proficiency in predicting 90-day mortality outcomes. The aim of this study was to quantify the extent to which an advanced language model like GPT-4 can augment the clinical decision-making process in AIS management. Specifically, we hypothesized that GPT-4 could provide accurate treatment recommendations and mortality predictions comparable to those of human experts, potentially contributing to improved patient outcomes in one of the most critical areas of emergency medicine.

Methods

Cohort Selection

This retrospective study comprised 100 consecutive cases from the ED of Rambam Healthcare Campus. All patients treated between January 2022 and April 2023 received a confirmed diagnosis of AIS. The inclusion criteria encompassed patients aged older than 18 years, a National Institutes of Health Stroke Scale (NIHSS) [21] score of 5 or higher (with the exception of patient 93 who received tissue plasminogen activator [tPA] offsite), and less than 5 hours from symptom onset to undergoing a noncontrast computed tomography (CT) of the brain. All included patients underwent noncontrast brain CT, CT angiography, and CT perfusion while in the ED. This cohort was specifically chosen for its alignment with American Heart Association guidelines for acute stroke management [13], making each patient a potential candidate for both tPA and endovascular thrombectomy (EVT) treatment. A total of 17 patients not meeting these criteria were categorized as "complex" cases, in which the clinical scenario warranted extra consideration of off-guideline treatment options, and there was a need to assess the individual patient's unique characteristics, medical history, and condition. For every patient, comprehensive medical records from their ED arrival, including imaging results, were collected and translated from Hebrew to English. Exclusion criteria were patients with incomplete clinical data or where stroke was not the final diagnosis.

Clinical data for each patient included demographics, medical history, chief complaints, symptom onset time, physical and neurological examinations, NIHSS score, imaging results (including Alberta Stroke Program Early CT Score [22] when available), treatment received, and mortality data. An experienced stroke specialist, blinded to the outcomes, reviewed the cases and made treatment decisions among no treatment, tPA, EVT, or a combination of tPA and EVT. All data were deidentified, removing identifiers, names, and dates.

Analysis Pipeline

The analysis used the OpenAI application programming interface "create chat completion" method with the model

gpt-4-1106-preview. Default parameters were set (temperature=1; top_p=1; n=1), and submissions were made using the R (R Foundation for Statistical Computing) wrapper library *openai*. Full prompt and example are available in Multimedia Appendix 1.

To assess the reliability of GPT-4 responses, each case underwent 5 submissions, as well as an additional submission without the accompanying clinical presentation narrative. For every treatment decision, GPT-4 provided a narrative explanation. In 95% (475/500) of cases, GPT-4 returned responses in the requested structure, which were automatically scraped with R. Unstructured responses were manually entered. For estimations provided as a range, the average was used. If GPT-4 provided a number with a greater symbol (eg, >50), the number was recorded with an additional 5. In 0.8% (4/500) of cases, GPT-4 did not return numeric responses for treatment decisions, and in 8.6% (43/500) of responses, it did not provide a 90-day mortality estimate.

Statistical Analysis

GPT-4's responses were scaled from 1 to 7 for treatment decisions and from 0 to 100 for 90-day mortality estimations. Averages were calculated across the 5 repeats. All statistical analyses were conducted using R (version 4.3.2), using base R functions, *predictive receiver operating characteristic (ROC)* 1.18.5, and *survival* 3.5.7. ROC curves were smoothed. Agreement between treatment decisions was measured using a linear weighted Cohen κ coefficient, using the *psych* 2.3.12 library.

Ethical Considerations

This study was approved by the Rambam Medical Center Helsinki Committee (0156-24-D) as a retrospective analysis.

The requirement for informed consent was waived due to the retrospective nature of the study and the use of deidentified data. All patient information was anonymized prior to analysis, with all identifiers, names, and dates removed to ensure privacy and confidentiality. No compensation was provided to participants as this was a retrospective study using existing clinical data. The study did not involve any images that could potentially identify individual participants. This research was conducted in accordance with the principles of the Declaration of Helsinki and adhered to all relevant institutional and national research ethics guidelines.

Results

Patient Demographics and Clinical Data

We generated a cohort from 100 consecutive cases of patients presenting with acute stroke symptoms at the ED of Rambam Healthcare Campus. All cases underwent full clinical and radiological evaluation in the emergency setting for acute stroke and were fully evaluated by a neurologist (Table 1 and Figure 1A). Revascularization treatment was administered to 78 of the patients: 36 were treated with tPA, 30 with EVT, and 12 received both. Within this cohort, 13 patients died within 90 days and 21 in total. Overall, 17 cases were classified as "complex" when not fitting exact treatment guidelines [13]. The data for each case encompassed demographics, NIHSS [21] scores, the timing of arrival to brain CT, onset of symptoms, and details from textual brain imaging results and risk factors that were available as medical history at the time of admission to the ED (Table S1 in Multimedia Appendix 2).



 Table 1. Study cohort clinical information and demographics.

Shmilovitch et al

Variable	Simple cases (n=83)	Complex cases (n=17)
Female sex, n (%)	38 (46)	7 (41)
Age (years), median (IQR)	75.0 (68.0-79.5)	71.0 (65.0-77.0)
First NIHSS ^a , median (IQR)	12.0 (8.5-16.5)	5.0 (5.0-9.0)
Time to CT ^b (hours), median (IQR)	1.8 (I1.5-2.6)	4.45 (3.0-5.2)
Brain CT findings, n (%)		
LVO ^c	48 (58)	7 (41)
MCA ^d	47 (57)	4 (24)
PCA ^e	8 (10)	4 (24)
Risk factors, n (%)		
Hypertension	51 (61)	10 (59)
DM^{f}	35 (42)	3 (18)
Dyslipidemia	36 (43)	6 (35)
Smoking	11 (13)	4 (24)
CKD ^g	11 (13)	0 (0)
Obese	5 (6)	0 (0)
Cancer	9 (11)	1 (6)
HF ^h	7 (8)	1 (6)
Cardiac arrhythmia	19 (23)	2 (12)
Family history for CAD ⁱ	1 (1)	0 (0)
tPA ^j , n (%)	29 (35)	7 (41)
EVT ^k , n (%)	29 (35)	1 (6)
tPA + EVT, n (%)	12 (14)	0 (0)
90-day mortality, n (%)	11 (13)	2 (12)
Overall mortality, n (%)	17 (20)	4 (24)

^aNIHSS: National Institutes of Health Stroke Scale.

^bCT: computed tomography.

^cLVO: large vessel occlusion.

^dMCA: middle cerebral artery.

^ePCA: posterior cerebral artery.

^fDM: diabetes mellitus.

^gCKD: chronic kidney disease.

^hHF: heart failure.

ⁱCAD: coronary artery disease.

^jtPA: tissue plasminogen activator.

^kEVT: endovascular thrombectomy.



Figure 1. Study design and GPT-4 performance evaluation. (A) Illustration of the study design involving 100 consecutive patients with stroke who underwent a comprehensive stroke workup, including perfusion, angiography, and noncontrast brain CT upon arrival at the emergency department. Clinical information, demographics, comorbidities, and CT perfusion results were recorded. The textual reports from these investigations were entered into the GPT-4 API, which was instructed to provide scores indicating whether to treat the patient, whether to administer tPA, whether to pursue EVT, and an estimate of 90-day mortality. (B) Box plots presenting average scores of GPT-4 assessments for decision to treat (y-axis). The comparison is made against real-world decisions and expert assessments of each case (true: to treat the patient and false: to not treat). (C) ROC curves and AUC scores of GPT-4 average scores for decision to treat, compared to real-world decisions and expert assessments. API: application programming interface; AUC: area under the curve; CT: computed tomography; EVT: endovascular thrombectomy; ROC: receiver operating characteristic; tPA: tissue plasminogen activator.



A stroke specialist, blinded to the outcomes, retrospectively reviewed each case. In 82 of the cases, the expert's decisions aligned with the actual treatments administered. Of note, the expert recommended not treating 11 patients who received treatment and suggested treatment for 7 who did not receive any. Concerning specific treatments, full agreement was observed in 61 cases, although the expert more frequently recommended combining tPA and EVT than what was observed in practice (Cohen κ =0.51, signifying moderate agreement).

GPT-4 Clinical Decisions

Independently, each case was assessed with GPT-4, generating a treatment recommendation scale from 1=intervention not recommended to 7=highly recommended (Figure 1A; Table S2 in Multimedia Appendix 2). To account for the variability in GPT-4 responses, each case was assessed 5 times. Cohen κ for treatment scores across runs ranged from 0.56 to 0.73. As

XSL•F() RenderX expected, the predefined "complex" cases demonstrated significantly greater variance between runs (P=.02).

Comparing GPT-4's treatment scale to both the expert's decision and the actual treatment revealed that the average scores from GPT-4 for patients who were treated were, on average, 1.9 points higher than those not treated (P<.001), and there was a 2.1-point difference in comparison to the expert decision (P<.001; Figure 1B). The average scores provided an area under the ROC curve (AUC-ROC) of 0.80 (95% CI 0.69-0.91) compared to the real-world decision, and 0.85 (95% CI 0.77-0.93) compared to the expert decision (Figure 1C). These average scores for AUCs were higher than those of each independent run (Multimedia Appendix 3). Additionally, removing the clinical presentation narrative from GPT-4's analysis resulted in a drop in AUC to 0.70 with the real-world decision and 0.72 with the expert decision (Multimedia Appendix 3), highlighting the importance of unstructured narrative data in treatment decision-making.

Similarly, setting the temperature of GPT-4 to 0 resulted in AUCs of 0.70 and 0.72 with the real-world and expert decisions, respectively, suggesting the need to allow GPT-4 more creativity to obtain better decisions.

Using a score threshold of 4, we observed 22 disagreements between GPT-4 and the real-world treatment and 20 disagreements with the expert decision. Notably, a substantial proportion of these disagreements coincided with cases where the expert and real-world decisions diverged, with 18 (60%) out of 30 such cases showing this dual disagreement. Moreover, complex cases were more prone to discrepancies, as 7 disagreements with the real-world decision and 5 with the expert decision were noted among the 17 complex cases. The specialist examined the explanatory text produced by GPT-4 for all discrepancies between the model and their blinded assessments, evaluating whether they agreed that the explanatory text, as part of the original model output, was logical and could be deemed good practice. Of the 20 instances where disagreements occurred, in 3 cases, the expert, after having carefully considered GPT-4's detailed explanations, conceded that GPT-4's assessment was preferable to their original decision. In additional 2 cases, the expert acknowledged that GPT-4's suggested approach was indeed acceptable and aligned with viable treatment options. In instances where the expert disagreed with GPT-4's reasoning, the disagreements primarily revolved around 3 key issues. First, GPT-4 inaccurately associated abnormal angiographic findings with clinical presentations. An illustrative case is that of a patient with stenosis of the right-sided middle cerebral artery who was presented with right hemiparesis (case 94). Despite these 2 elements potentially being anatomically unrelated, GPT-4 linked them erroneously. The second notable issue pertained to ethical considerations, particularly in a case involving a patient with active laryngeal cancer and cognitive decline. According to guidelines, the patient was deemed eligible for treatment, but the expert's decision was to not proceed with treatment as life expectancy was short and he was palliative (case 14). Third, discrepancies arose in deviations from guidelines, particularly in cases of distal thrombectomies. For instance, in the case of a patient with M2 obstruction (considered distal thrombus) aged 96 years, GPT-4 recommended against treatment, which is the established guidelines; however, the expert call was to proceed with thrombectomy due to a high NIHSS score and good results in such cases in the past from personal experience (case 54).

In assessing GPT-4's ability to choose the best treatment option, it showed near-perfect agreement with real-world decisions in recommending EVT: GPT-4 suggested EVT for all patients (42/42, 100%) treated with EVT (average score>4). The expert suggested EVT for 55 patients, of which 50 were also recommended EVT by GPT-4, corresponding to an AUC of 0.94 (95% CI 0.89-0.98) with real-world decisions and 0.95 (95% CI: 0.90-0.99) with the expert (Figure 2A). For tPA treatment, GPT-4 recommended it for 38 (79%) of the 48 patients who received it, showing a closer agreement with the expert. Of the 41 patients recommended for tPA by the expert, GPT-4 agreed on 35 (85%), corresponding to an AUC of 0.77 (95% CI 0.68-0.86) with real-world decisions and 0.82 (95% CI 0.73-0.90) with the expert (Figure 2B).

Figure 2. GPT-4 treatment type scores. Box plots depict GPT-4 treatment type scores, with the y-axis representing probability score (1-7 scale). Each treatment category is color coded: green for no intervention, orange for tPA, purple for EVT, and pink for tPA and EVT. (A) GPT-4 scores for EVT, stratified by real-world decisions and expert assessments. (B) GPT-4 scores for tPA, stratified by real-world decisions and expert assessments. EVT: endovascular thrombectomy; tPA: tissue plasminogen activator.



Mortality Risk

We further evaluated the ability of GPT-4 to predict 90-day mortality. The model estimated an average mortality risk of 55.1% for patients who died within 90 days, compared to 31.5% for survivors (P<.001), yielding an AUC of 0.89 (95% CI 0.81-0.98; Figure 3A). To contextualize these results, we compared GPT-4's performance with that of 2 recent machine

learning models specifically trained for 90-day mortality prediction. In our cohort, the PRACTICE model [18] achieved an AUC of 0.70, significantly worse than the GPT-4 predictions (log-rank P value=.02), while the PREMISE model [19] reached an AUC of 0.77 (P=.07; Figure 3A). These comparisons underscore GPT-4's remarkable accuracy in mortality risk assessment, outperforming specialized, trained predictive models.



Figure 3. GPT-4 mortality predictions. (A) ROC curve for 90-day mortality estimations by GPT-4 (red), PRACTICE (green), and PREMISE (blue). (B) Kaplan-Meier plot stratifying individuals into low- and high-risk categories for mortality based on GPT-4's 90-day mortality estimations. AUC: area under the curve; ROC: receiver operating characteristic.



For identifying high-risk patients, we set a threshold at the top 25% of the cohort, which corresponded to a predicted mortality risk cutoff of 41%. Within this high-risk group, 10 patients passed away within 90 days of admission, and an additional 3 within the subsequent year (Figure 3B). Conversely, among the remaining 75 patients categorized as lower risk, only 3 deaths occurred within the 90-day period, and 6 in total during the first

year. The calculated hazard ratio was 6.98 (95% CI 2.88-16.9; P<.001), reinforcing the model's capability to stratify patients based on their mortality risk effectively.

Discussion

Here, we demonstrate the potential of GPT-4 as a clinical decision-support tool in AIS management. Our main findings

XSL•FO RenderX

show that treatment recommendations by GPT-4 closely aligned with both expert opinions (AUC 0.85) and real-world decisions (AUC 0.80). Notably, GPT-4 exhibited high accuracy in predicting 90-day mortality (AUC 0.89), outperforming specialized machine learning models.

AIS is a leading cause of mortality and disability worldwide [23-25]. The urgency of stroke care is particularly critical in regions with limited access to specialized stroke units or qualified physicians [26,27]. GPT-4's ability to operate seamlessly within existing treatment routines, relying solely on routine chart information, makes it valuable for quick triage in underresourced settings [7]. This accessibility could democratize high-level medical consultation, extending expert-level decision-making to underresourced health care facilities.

In our study, GPT-4 demonstrated high accuracy in predicting 90-day mortality for patients with AIS undergoing endovascular treatment. The model used a diverse range of clinical and imaging variables, offering a more comprehensive approach compared to existing models like Houston intraarterial therapy, Houston intraarterial therapy 2, PREMISE, and PRACTICE [18,19,28,29]. Unlike traditional health care predictive models that rely on structured data, GPT-4 provided recommendations based on narrative text. Our analyses highlighted the significance of unstructured data, as evidenced by the drop in prediction accuracy when the narrative clinical presentation was excluded. This showcases GPT-4's capability to handle complex medical data in a way that aligns with the natural flow of clinical information.

A crucial aspect of deploying AI models like GPT-4 in health care is the transparency and interpretability of their decision-making process. While GPT-4's natural language outputs can give the impression of explainability, these may not necessarily reflect a truly reliable reasoning process. Our analysis focused on the face value of GPT-4's rationales, which were deemed insightful by the expert reviewer. However, we acknowledge the potential for convincing but flawed explanations, a known limitation of LLMs. This highlights the importance of critical evaluation and cautious interpretation of such model outputs, particularly in high-stakes medical decision-making contexts. Ongoing research is needed to address the transparency and reliability of AI systems' reasoning processes before their broader integration into clinical practice.

Despite its promising results, our study has several limitations. We must acknowledge certain challenges in applying GPT-4, especially regarding its ability to assess ethical issues. The model may face difficulties in addressing the nuanced and complex ethical considerations intrinsic to medical decision-making. This limitation emphasizes the necessity for cautious and supplementary human oversight when deploying AI tools like GPT-4 in sensitive health care contexts. The occurrence of "hallucinations" or erroneous outputs is another concern, although we demonstrated that running multiple assessments can mitigate this risk. Future research should focus on refining these methods to further reduce inaccuracies.

Another consideration is the generalizability of these findings. While it is possible that the recommendations may partially reflect the clinician's intuition encoded in the clinical notes, our analyses suggest that the model's assessments go beyond mere interpretation. The discrepancies observed between the GPT-4 recommendations and both the real-world treatment decisions and the expert evaluations indicate that the model is capable of making independent assessments based on the provided data. Furthermore, the clinical presentation notes and imaging report interpretations (Table S1 in Multimedia Appendix 2) do not explicitly convey the clinician's treatment preferences or intuitions, suggesting that GPT-4 is not simply regurgitating the clinician's thought process. Another possible limitation is the study's exclusion criteria, particularly the retrospective exclusion of patients with incomplete clinical data or those who were ultimately diagnosed with conditions other than stroke. While these exclusions were necessary to ensure the study focused on accurately diagnosed AIS cases for which GPT-4 decision-support capabilities could be most relevant, we acknowledge that this approach may limit the generalizability of our findings to broader clinical settings. In real-world scenarios, clinicians are often faced with diagnostic uncertainty and incomplete information when making treatment decisions. Finally, our study was conducted in a single center with a specific patient population. Further studies across diverse settings and larger populations are necessary to validate the efficacy and applicability of GPT-4 in various clinical environments.

In conclusion, our study introduces a groundbreaking approach to clinical decision support in stroke management using GPT-4. This model has shown the potential to process narrative text, provide explainable recommendations, and enhance medical decision-making. As we continue to explore and refine this technology, it holds the promise of transforming patient care and improving outcomes in one of the most critical areas of medicine.

Data Availability

All data generated or analyzed during this study are included in Multimedia Appendix 2.

Authors' Contributions

SS and DA conceived and designed the study. SS and AM collected the clinical data and translated the narratives. SP and SS reviewed the cases and provided expert analysis. DA conducted the computational and statistical analyses with inputs from SS. SS and DA drafted the manuscript with contributions and critical revisions from AM, MK, SP, and MCS. All authors reviewed and approved the final manuscript.



None declared.

Multimedia Appendix 1 Prompt used. [PDF File (Adobe PDF File), 360 KB - ai v4i1e60391 app1.pdf]

Multimedia Appendix 2 Supplementary tables. [XLSX File (Microsoft Excel File), 558 KB - ai v4i1e60391 app2.xlsx]

Multimedia Appendix 3

GPT-4 Assessments Performance. Area under the curve (AUC) for GPT-4 decision to treatment scores of each of the individual submissions (1-5) and the average. Each individual submission is lower than the average. In addition, we submitted the cases without the clinical presentation narrative, which yielded lower AUC (no narrative). Similarly, lower AUC was observed when cases were submitted with temperature=0.

[PDF File (Adobe PDF File), 82 KB - ai_v4i1e60391_app3.pdf]

References

- 1. GPT-4. OpenAI. URL: <u>https://openai.com/index/gpt-4/</u> [accessed 2025-01-22]
- Sanderson K. GPT-4 is here: what scientists think. Nature 2023;615(7954):773. [doi: <u>10.1038/d41586-023-00816-5</u>] [Medline: <u>36928404</u>]
- Lee P, Bubeck S, Petro J. Benefits, limits, and risks of GPT-4 as an AI chatbot for medicine. Reply. N Engl J Med 2023;388(25):2400. [doi: <u>10.1056/NEJMc2305286</u>] [Medline: <u>37342941</u>]
- 4. Zack T, Lehman E, Suzgun M, Rodriguez JA, Celi LA, Gichoya J, et al. Assessing the potential of GPT-4 to perpetuate racial and gender biases in health care: a model evaluation study. Lancet Digit Health 2024;6(1):e12-e22 [FREE Full text] [doi: 10.1016/S2589-7500(23)00225-X] [Medline: 38123252]
- Chiu WHK, Ko WSK, Cho WCS, Hui SYJ, Chan WCL, Kuo MD. Evaluating the diagnostic performance of large language models on complex multimodal medical cases. J Med Internet Res 2024;26:e53724 [FREE Full text] [doi: 10.2196/53724] [Medline: <u>38739441</u>]
- Ziegelmayer S, Marka AW, Lenhart N, Nehls N, Reischl S, Harder F, et al. Evaluation of GPT-4's chest x-ray impression generation: a reader study on performance and perception. J Med Internet Res 2023;25:e50865 [FREE Full text] [doi: 10.2196/50865] [Medline: <u>38133918</u>]
- Kanjee Z, Crowe B, Rodman A. Accuracy of a generative artificial intelligence model in a complex diagnostic challenge. JAMA 2023;330(1):78-80 [FREE Full text] [doi: 10.1001/jama.2023.8288] [Medline: 37318797]
- Shea YF, Lee CMY, Ip WCT, Luk DWA, Wong SSW. Use of GPT-4 to analyze medical records of patients with extensive investigations and delayed diagnosis. JAMA Netw Open 2023;6(8):e2325000 [FREE Full text] [doi: 10.1001/jamanetworkopen.2023.25000] [Medline: 37578798]
- 9. Xue E, Bracken-Clarke D, Iannantuono GM, Choo-Wosoba H, Gulley JL, Floudas CS. Utility of large language models for health care professionals and patients in navigating hematopoietic stem cell transplantation: comparison of the performance of ChatGPT-3.5, ChatGPT-4, and Bard. J Med Internet Res 2024;26:e54758 [FREE Full text] [doi: 10.2196/54758] [Medline: 38758582]
- Brin D, Sorin V, Vaid A, Soroush A, Glicksberg BS, Charney AW, et al. Comparing ChatGPT and GPT-4 performance in USMLE soft skill assessments. Sci Rep 2023;13(1):16492 [FREE Full text] [doi: 10.1038/s41598-023-43436-9] [Medline: 37779171]
- Guillen-Grima F, Guillen-Aguinaga S, Guillen-Aguinaga L, Alas-Brun R, Onambele L, Ortega W, et al. Evaluating the efficacy of ChatGPT in navigating the Spanish medical residency entrance examination (MIR): promising horizons for AI in clinical medicine. Clin Pract 2023;13(6):1460-1487 [FREE Full text] [doi: 10.3390/clinpract13060130] [Medline: 37987431]
- 12. Guerra GA, Hofmann H, Sobhani S, Hofmann G, Gomez D, Soroudi D, et al. GPT-4 artificial intelligence model outperforms ChatGPT, medical students, and neurosurgery residents on neurosurgery written board-like questions. World Neurosurg 2023;179:e160-e165. [doi: 10.1016/j.wneu.2023.08.042] [Medline: 37597659]
- Kleindorfer DO, Towfighi A, Chaturvedi S, Cockroft KM, Gutierrez J, Lombardi-Hill D, et al. 2021 guideline for the prevention of stroke in patients with stroke and transient ischemic attack: a guideline from the American Heart Association/American Stroke Association. Stroke 2021;52(7):e364-e467 [FREE Full text] [doi: 10.1161/STR.00000000000375] [Medline: 34024117]
- 14. Brown DL, Levine DA, Albright K, Kapral MK, Leung LY, Reeves MJ, et al. Benefits and risks of dual versus single antiplatelet therapy for secondary stroke prevention: a systematic review for the 2021 guideline for the prevention of stroke

in patients with stroke and transient ischemic attack. Stroke 2021;52(7):e468-e479 [FREE Full text] [doi: 10.1161/STR.000000000000377] [Medline: 34024115]

- 15. Amin HP, Madsen TE, Bravata DM, Wira CR, Johnston SC, Ashcraft S, et al. Diagnosis, workup, risk reduction of transient ischemic attack in the emergency department setting: a scientific statement from the American Heart Association. Stroke 2023;54(3):e109-e121 [FREE Full text] [doi: 10.1161/STR.000000000000418] [Medline: 36655570]
- Koka A, Suppan L, Cottet P, Carrera E, Stuby L, Suppan M. Teaching the National Institutes of Health Stroke Scale to paramedics (e-learning vs video): randomized controlled trial. J Med Internet Res 2020;22(6):e18358 [FREE Full text] [doi: 10.2196/18358] [Medline: 32299792]
- Linfante I, Walker GR, Castonguay AC, Dabus G, Starosciak AK, Yoo AJ, et al. Predictors of mortality in acute ischemic stroke intervention: analysis of the North American Solitaire Acute Stroke Registry. Stroke 2015;46(8):2305-2308. [doi: 10.1161/STROKEAHA.115.009530] [Medline: 26159790]
- Li H, Ye SS, Wu YL, Huang SM, Li YX, Lu K, et al. Predicting mortality in acute ischaemic stroke treated with mechanical thrombectomy: analysis of a multicentre prospective registry. BMJ Open 2021;11(4):e043415 [FREE Full text] [doi: 10.1136/bmjopen-2020-043415] [Medline: 33795300]
- Gattringer T, Posekany A, Niederkorn K, Knoflach M, Poltrum B, Mutzenbach S, et al. Predicting early mortality of acute ischemic stroke. Stroke 2019;50(2):349-356. [doi: <u>10.1161/STROKEAHA.118.022863</u>] [Medline: <u>30580732</u>]
- 20. Noser EA, Zhang J, Rahbar MH, Sharrief AZ, Barreto AD, Shaw S, et al. Leveraging multimedia patient engagement to address minority cerebrovascular health needs: prospective observational study. J Med Internet Res 2021;23(8):e28748 [FREE Full text] [doi: 10.2196/28748] [Medline: 34397385]
- 21. Kwah LK, Diong J. National Institutes of Health Stroke Scale (NIHSS). J Physiother 2014;60(1):61 [FREE Full text] [doi: 10.1016/j.jphys.2013.12.012] [Medline: 24856948]
- 22. Pop N, Tit D, Diaconu C, Munteanu M, Babes E, Stoicescu M, et al. The Alberta Stroke Program Early CT score (ASPECTS): a predictor of mortality in acute ischemic stroke. Exp Ther Med 2021;22(6):1371 [FREE Full text] [doi: 10.3892/etm.2021.10805] [Medline: 34659517]
- 23. Lim GB. Global burden of cardiovascular disease. Nat Rev Cardiol 2013;10(2):59. [doi: 10.1038/nrcardio.2012.194] [Medline: 23296068]
- 24. Feigin VL, Forouzanfar MH, Krishnamurthi R, Mensah GA, Connor M, Bennett DA, et al. Global and regional burden of stroke during 1990-2010: findings from the Global Burden of Disease Study 2010. Lancet 2014;383(9913):245-254 [FREE Full text] [doi: 10.1016/s0140-6736(13)61953-4] [Medline: 24449944]
- 25. GBD 2016 Causes of Death Collaborators. Global, regional, and national age-sex specific mortality for 264 causes of death, 1980-2016: a systematic analysis for the Global Burden of Disease Study 2016. Lancet 2017;390(10100):1151-1210 [FREE Full text] [doi: 10.1016/S0140-6736(17)32152-9] [Medline: 28919116]
- 26. Saver JL, Fonarow GC, Smith EE, Reeves MJ, Grau-Sepulveda MV, Pan W, et al. Time to treatment with intravenous tissue plasminogen activator and outcome from acute ischemic stroke. JAMA 2013;309(23):2480-2488. [doi: 10.1001/jama.2013.6959] [Medline: 23780461]
- Strbian D, Soinne L, Sairanen T, Häppölä O, Lindsberg PJ, Tatlisumak T, et al. Ultraearly thrombolysis in acute ischemic stroke is associated with better outcome and lower mortality. Stroke 2010;41(4):712-716. [doi: 10.1161/STROKEAHA.109.571976] [Medline: 20167917]
- 28. Ryu CW, Kim BM, Kim HG, Heo JH, Nam HS, Kim DJ, et al. Optimizing outcome prediction scores in patients undergoing endovascular thrombectomy for large vessel occlusions using collateral grade on computed tomography angiography. Neurosurgery 2019;85(3):350-358. [doi: 10.1093/neuros/nyy316] [Medline: 30010973]
- 29. Hallevi H, Barreto AD, Liebeskind DS, Morales MM, Martin-Schild SB, Abraham AT, et al. Identifying patients at high risk for poor outcome after intra-arterial therapy for acute ischemic stroke. Stroke 2009;40(5):1780-1785 [FREE Full text] [doi: 10.1161/STROKEAHA.108.535146] [Medline: 19359652]

Abbreviations

AI: artificial intelligence
AIS: acute ischemic stroke
AUC: area under the curve
CT: computed tomography
ED: emergency department
EVT: endovascular thrombectomy
LLM: large language model
NIHSS: National Institutes of Health Stroke Scale
ROC: receiver operating characteristic
tPA: tissue plasminogen activator



Edited by K El Emam; submitted 09.05.24; peer-reviewed by MO Khursheed, R McDonough; comments to author 20.07.24; revised version received 06.08.24; accepted 08.11.24; published 07.03.25. <u>Please cite as:</u> Shmilovitch AH, Katson M, Cohen-Shelly M, Peretz S, Aran D, Shelly S GPT-4 as a Clinical Decision Support Tool in Ischemic Stroke Management: Evaluation Study JMIR AI 2025;4:e60391 URL: https://ai.jmir.org/2025/1/e60391 doi:10.2196/60391 PMID:40053715

©Amit Haim Shmilovitch, Mark Katson, Michal Cohen-Shelly, Shlomi Peretz, Dvir Aran, Shahar Shelly. Originally published in JMIR AI (https://ai.jmir.org), 07.03.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Creation of Scientific Response Documents for Addressing Product Medical Information Inquiries: Mixed Method Approach Using Artificial Intelligence

Jerry Lau^{1,2,3}, PharmD; Shivani Bisht⁴, M Tech; Robert Horton⁵, PhD; Annamaria Crisan⁶, BPharm, MSc; John Jones¹, MBA; Sandeep Gantotti⁷, BSP; Evelyn Hermes-DeSantis^{1,2}, BPCS, PharmD

¹phactMI, Gainesville, FL, United States

²Department of Pharmacy Practice and Administration, Ernest Mario School of Pharmacy, Rutgers, The State University of New Jersey, New Brunswick, NJ, United States
³EMD Serono, Boston, MA, United States
⁴Eli Lilly and Company, Bangalore, India
⁵Win-Vector Labs, Marin City, CA, United States
⁶Pfizer, Montreal, QC, Canada
⁷Indegene Limited, Bangalore, India

Corresponding Author:

Evelyn Hermes-DeSantis, BPCS, PharmD phactMI 5931 NW 1st Place Gainesville, FL, 32607 United States Phone: 1 2155881585 Email: <u>evelyn@phactmi.org</u>

Abstract

Background: Pharmaceutical manufacturers address health care professionals' information needs through scientific response documents (SRDs), offering evidence-based answers to medication and disease state questions. Medical information departments, staffed by medical experts, develop SRDs that provide concise summaries consisting of relevant background information, search strategies, clinical data, and balanced references. With an escalating demand for SRDs and the increasing complexity of therapies, medical information departments are exploring advanced technologies and artificial intelligence (AI) tools like large language models (LLMs) to streamline content development. While AI and LLMs show promise in generating draft responses, a synergistic approach combining an LLM with traditional machine learning classifiers in a series of human-supervised and -curated steps could help address limitations, including hallucinations. This will ensure accuracy, context, traceability, and accountability in the development of the concise clinical data summaries of an SRD.

Objective: This study aims to quantify the challenges of SRD development and develop a framework exploring the feasibility and value addition of integrating AI capabilities in the process of creating concise summaries for an SRD.

Methods: To measure the challenges in SRD development, a survey was conducted by phactMI, a nonprofit consortium of medical information leaders in the pharmaceutical industry, assessing aspects of SRD creation among its member companies. The survey collected data on the time and tediousness of various activities related to SRD development. Another working group, consisting of medical information professionals and data scientists, used AI to aid SRD authoring, focusing on data extraction and abstraction. They used logistic regression on semantic embedding features to train classification models and transformer-based summarization pipelines to generate concise summaries.

Results: Of the 33 companies surveyed, 64% (21/33) opened the survey, and 76% (16/21) of those responded. On average, medical information departments generate 614 new documents and update 1352 documents each year. Respondents considered paraphrasing scientific articles to be the most tedious and time-intensive task. In the project's second phase, sentence classification models showed the ability to accurately distinguish target categories with receiver operating characteristic scores ranging from 0.67 to 0.85 (all *P*<.001), allowing for accurate data extraction. For data abstraction, the comparison of the bilingual evaluation understudy (BLEU) score and semantic similarity in the paraphrased texts yielded different results among reviewers, with each preferring different trade-offs between these metrics.

Conclusions: This study establishes a framework for integrating LLM and machine learning into SRD development, supported by a pharmaceutical company survey emphasizing the challenges of paraphrasing content. While machine learning models show potential for section identification and content usability assessment in data extraction and abstraction, further optimization and research are essential before full-scale industry implementation. The working group's insights guide an AI-driven content analysis; address limitations; and advance efficient, precise, and responsive frameworks to assist with pharmaceutical SRD development.

(JMIR AI 2025;4:e55277) doi:10.2196/55277

KEYWORDS

AI; LLM; GPT; biopharmaceutical; medical information; content generation; artificial intelligence; pharmaceutical; scientific response; documentation; information; clinical data; strategy; reference; feasibility; development; machine learning; large language model; accuracy; context; traceability; accountability; survey; scientific response documentation; SRD; benefit; content generator; content analysis; Generative Pre-trained Transformer

Introduction

Pharmaceutical manufacturers play a crucial role in meeting health care professionals' information needs by providing them with scientific response documents (SRDs). These documents provide comprehensive and evidence-based answers to unsolicited questions concerning a medication or disease state [1]. The development and maintenance of SRDs are entrusted the medical information department within to these organizations. This department is composed of medical experts who possess in-depth knowledge of specific therapeutic areas and are responsible for various strategic activities, including the meticulous development of SRDs [2]. SRDs are tailored to address specific inquiries, presenting a concise summary, relevant background information, clinical data, and scientifically balanced references [1]. Considering the escalating demand for SRDs and the increasing complexity of therapies, the role of medical information departments has become more critical than ever. A 2018 survey of 27 pharmaceutical companies revealed that a medical information department creates an average of 716 new SRDs and maintains 2510 existing SRDs annually [2]. Fully developing a new SRD required an average of 31 hours for medical experts, while updating or revising existing SRDs involved an average of 21 hours [2]. Medical information experts use this time to answer the SRD query following a scientific method approach [3]. The strategic and resource-intensive nature of SRD development and the surge in health care professional inquiries emphasize the pressing need for timely and comprehensive information. To address these challenges, there is a growing interest across medical information departments in leveraging advanced technologies and artificial intelligence (AI) tools, such as large language models (LLMs) and traditional machine learning techniques, to enhance and streamline the SRD development process. There are several steps to develop an SRD, including reading articles, selecting article content, paraphrasing article content, creating a citation list, editorial changes, data integrity, and content review. Some of these steps may be more time-consuming than others.

To better understand the current advancements in AI, consider an analogy used in software development. Programming can be thought of as software 1.0, where a machine relies on explicit, step-by-step instructions from a programmer to perform designated tasks. Machine learning represents software 2.0,

```
https://ai.jmir.org/2025/1/e55277
```

where developers present labeled examples of input and output data to the machine so that it can identify patterns that allow it to predict outcomes from inputs. This kind of supervised machine learning has enabled rapid progress in many areas of natural language processing, including applications in language translation, sentiment analysis, and information retrieval. More recently, LLMs, such as OpenAI's Generative Pre-trained Transformer (GPT), are complex machine learning models trained to predict subsequent words in natural language text based on the text so far. This allows the machine to generate statistically plausible output given a "prompt." Beyond simple prompt completion, such models can be trained to follow instructions in the prompt, such as "Summarize the following paragraph." Designing prompts that lead an LLM to produce a desired output is a novel and distinct paradigm in software development, which can be classified as "software 3.0" [4].

Language models convert language to numerical representation, and specialized models create semantic embedding by exporting a sentence as a vector of floating-point numbers [5]. By converting concepts into numeric vectors, embeddings enable computers to represent the connections between concepts. The relationship between two embeddings is determined by the vector distance, with smaller distances indicating higher relatedness and larger distances implying lower relatedness. Embeddings are easily consumed and compared by other machine learning models and algorithms for tasks like clustering text strings based on similarity or ranking search results by query relevance. Furthermore, embeddings correspond to similar meanings.

Figure 1 shows examples of semantic embeddings of sentences based on the dataset used in this study. The original 768-dimension embeddings were mapped down to 2 dimensions to visualize them, showing that sentences on similar topics are close together. Colors indicate the category to which the sentence belongs. Here, the 3 sentences in blue ("Population") are close semantically to one another, as are the 3 sentences in red ("Adverse_events"). One of the sentences in the "Efficacy" category is far from the other two, but on examining the sentences, it is considered an outlier talking about a ratio of antibodies, while the two that are close to one another both concern statistical significance.

LLMs apply traditional machine learning concepts and embeddings on a larger scale. Transformers process sequential

data, such as natural language, all at once, enabling them to perform tasks like text summarization [6]. GPT is trained to predict the next word using preceding words, capturing linguistic patterns and semantic relationships in large text datasets. GPT often produces coherent and plausible responses. By providing labeled examples, GPT can be fine-tuned for specific tasks to enhance its capabilities. This fine-tuning process allows GPT to adapt its prelearned knowledge to effectively perform tasks such as text generation, question answering, and language translation [7].

Figure 1. High-dimensional data visualization of embeddings. The t-SNE (t-distributed stochastic neighbor embedding) algorithm was used to transform data into 2 dimensions. Different colors were chosen for different sections based on reviewer feedback (based on the test set used in the study).



AI tools have a well-established history in medicine, with potential applications like artificial neural networks aiding clinical prognosis and diagnosis through pattern recognition first identified in 2004 [8]. Furthermore, within academic and research writing, OpenAI's ChatGPT has been used to "extract" important information from academic papers (eg, author details, publication date, main findings, etc) and generate summaries of these lengthy papers [9]. However, the use of AI to create medical content, particularly SRDs, is still in its early stages. An April 2023 study showcased the potential of AI by using OpenAI's ChatGPT to generate draft responses to patient questions based on deidentified information [10]. This pioneering work highlights the need to explore AI's capabilities in medical content generation in depth.

Although ChatGPT demonstrates impressive language generation abilities, relying solely on it has limitations. ChatGPT, like any LLM, can hallucinate and produce content based on its prediction without logic or fact-checking abilities [11]. Furthermore, there exists a lack of transparency in the training sets used for LLMs like ChatGPT. This, coupled with the complexity of these models, may lead to false or biased information being unintentionally included in the generated content [12]. The accuracy of an SRD is crucial in its creation. Furthermore, traceability and accountability are essential considerations. The use of LLMs like ChatGPT often results in

the original authors and sources not being cited, leading to the misattribution of information [13].

This study has 2 aims. The first is to quantify the challenges of SRD creation by gathering the opinions of medical information professionals regarding the time consumption of the various steps of SRD development. To address these challenges and leverage the strengths of both human expertise and AI in the creation of SRDs, a synergistic approach that combines LLM with traditional machine learning classifiers is warranted. The second aim of this study is to develop a framework to explore the feasibility and value addition of integrating AI capabilities, including LLM and machine learning, into the SRD creation process.

Methods

Survey of phactMI Members

A working group from phactMI developed a cross-sectional survey to assess the time and tediousness of various aspects of SRD creation. phactMI, a nonprofit consortium of medical information leaders from the pharmaceutical industry, conducted the survey using the survey tool Alchemer. The initial contact for the web-based open survey link was emailed to one contact at each of the 33 member companies in March 2023 (see Multimedia Appendix 1 for email wording). Participation in the survey was voluntary, and no incentives were offered. The survey link was sent once, with one reminder sent during March 2023, and the survey closed on April 15, 2023. The working group pretested the survey using the Alchemer system before distribution. In the recruitment email, the purpose of the survey, length and duration, the lead investigator, and how all data were to be handled were disclosed. Proceeding to the first question was considered consent to participate.

The creation of an SRD is a strategic endeavor comprised of several steps that may be more time-consuming and tedious than others. Specific data collected in the survey included the average time needed for creating an SRD, the average number of papers included in an SRD, etc. Survey respondents were given a list of activities, including paraphrasing article content, creating a citation list, making editorial changes, improving data integrity, selecting article content, reviewing content, and reading articles. Respondents were asked to rank given activities from 1 to 8 in terms of time consumption and tediousness (1 being the most time-consuming or tedious and 8 being the least time-consuming or tedious). The interpretations of time-consuming and tedious were left to the discretion of the survey respondents.

Not all steps had to be ranked by all respondents. A score for each step was created with a weighted calculation, with items ranked first being given a higher value or weight. Weighted values are based on the number of steps selected. The higher the score, the more time-consuming or tedious the steps were considered. The survey results were analyzed to identify those steps in the development of an SRD where the use of AI may offer maximum benefit.

The survey questions were not randomized, and there was no adaptive questioning. There was a total of 10 questions. All questions were displayed on the same page, so no back button or review step was necessary.

Only 1 response per company was allowed. Data were analyzed using descriptive statistics. The full survey questionnaire is provided in Multimedia Appendix 2. The Checklist for Reporting Results of Internet E-Surveys (CHERRIES) for this survey is provided in Multimedia Appendix 3.

Ethical Considerations

The survey was not approved by an institutional review board as it was not considered human subject data. All survey data were deidentified, saved, and reported in aggregate.

Authoring SRD

Another working group consisting of medical information professionals and data scientists was created. Their goal was to leverage AI to support the medical information department's creation of SRDs. Their aim was to develop a tool that could process scientific articles (input) and provide concise summaries (output). The group identified two key steps in the document authoring process: data extraction and data abstraction. Their problem was figuring out the process between the input and output (Figure 2). Data extraction is the selection of key sentences from publications that address all the data points authors would want to include in a response document, and data abstraction is the generation of a summary of extracted data, followed by paraphrasing to avoid plagiarism of original texts.







Data Extraction and Machine Training

The working group selected scientific texts from the PubMed Central database focusing on clinical drug trials for data extraction. The narrative text from these articles, excluding text in tables, was extracted, cleaned, and placed into Prodigy, a data annotation tool. A total of 3 domain experts and medical information specialists labeled sentences from narrative text into 5 classifications: safety, efficacy, treatment, population, and study design. These classifications correspond to the main sections of a clinical trial used in the creation of an SRD. A fourth domain expert, a data editor, reviewed all the labels to ensure the labeling criteria were applied consistently. These labeled data were then fed into logistic regression classification models to train the models on identification. The training dataset is available in Multimedia Appendix 4.

Participating companies provided 3 SRDs to the working group. The team extracted clean, narrative text from the provided documents to feed into the models. The models categorized each sentence based on their previous training. Reviewers evaluated and assessed model classifications. Trained models' performance was evaluated with a receiver operating characteristic (ROC) curve plotting the true positive rate (TPR) and false positive rate (FPR). The area under the curve (AUC) provides an aggregate measure of performance across all possible thresholds, with a higher AUC indicating better performance of the model. A Wilcoxon-Mann-Whitney *U* test statistic was applied.

Data Abstraction

Summarizing the extracted data was the initial step in data abstraction. The working group used the Hugging Face transformers summarization pipeline leveraging the Facebook/BART-large-cnn model, a language model trained for summarization. The second step was to rewrite and synthesize the extracted text without plagiarizing the original reference by using the GPT-3 model (text-davinci-003). The model received the prompt "Paraphrase this without

```
https://ai.jmir.org/2025/1/e55277
```

plagiarizing," followed by the summarized text. Multiple paraphrases were generated for each input.

Filtering Output

A total of 2 criteria were used to sort and rank the paraphrased texts: semantic similarities and bilingual evaluation understudy (BLEU) scores. Semantic similarity, measured using cosine similarity between sentence transformer embeddings (distiluse-base-multilingual-cased-v2), assessed the likeness in meaning between the paraphrased sentences and the original text. The greater the semantic similarity between the two sentences, the better the quality of the paraphrasing. The second criterion was the BLEU score, which measured the similarity in word or phrase use between a generated text and the original text. It was calculated using sacrebleu with effective_order set to true. A low BLEU score reflects a higher quality of paraphrasing, as it indicates less similarity in words and phrases with the original text. Finding the right balance between semantic and textual similarities was crucial for the overall paraphrasing quality. Human reviewers then evaluated the paraphrased text and ranked the text by usefulness with rationales provided.

Throughout the study, the working group fostered collaboration between medical information professionals and data scientists to validate the results. Results from each step were edited by hand to make sure that the next step had clean inputs.

Results

Survey of phactMI Members

A total of 21 of the 33 pharmaceutical member companies, based on IP address, opened the survey (view rate 64%). A total of 16 pharmaceutical member companies participated in the survey (participation rate 76%, 16/21), with a completion rate of 81% (13/16). No cookies were used to assign user identification. Duplicate entries were identified by either IP address or company name (if provided). The most complete or

most recent entry was kept for analysis. All data from unique entries were included in the analysis.

On average, a medical information department creates 614 (range: low to 2676) new SRDs and updates or revises 1352 (range: low to 6057) SRDs annually. Respondents indicate it takes, on average, 8.3 hours to create a new SRD and 3.8 hours to update an SRD. In addition, 87% (14/16) of respondents included content from at least 4 studies in SRDs summarizing

clinical trial data. The survey results revealed that the top 3 most time-consuming steps in SRD development were paraphrasing study content, checking the data integrity of the paraphrased text versus the source publications, and checking the data integrity of the SRD (eg, checking that the text is cited to the correct publications; Figure 3). While paraphrasing article content was also the most tedious step, the other top steps differed, with writing citations and editorial changes rounding out the top 3 (Figure 4).

Figure 3. Ranking of steps deemed time-consuming by survey respondents. SRD: scientific response document.

	Item	Rank	Rank distribution	Score	n	Did not rank
8	Paraphrase the selected content from the article	1		59	10	7
-consumin	Data integrity at article level (eg, fact checking the data, appropriate interpretation)	2	-	49	10	7
ost time	Data integrity at SRD level (eg, fact checking the data and citations)	3		48	11	6
W	Identify and select key content from the article	4		48	10	7
	Read the article	5		44	9	8
ung	Content review: including content inclusion selection, quality of paraphrasing	6		41	9	8
-consun	Write list of citations	7		39	10	9
east time	Editorial changes, formatting, etc.	8	-	29	10	9
r			Line of neutrality			

Least time-consuming Most time-consuming

Figure 4. Rankings of tasks deemed tedious by survey respondents. SRD: scientific response document.

	Item	Rank	Rank distribution	Score	n	Did not rank
ious	Paraphrase the selected content from the article	1		65	11	6
ost ted	Write list of citations	2		56	10	7
W	Editorial changes, formatting, etc.	3		52	12	5
	Data integrity at article level (eg, fact checking the data, appropriate interpretation)	4		51	10	7
	Data integrity at SRD level (eg, fact checking the data and citations)	5		47	11	6
Sma	Identify and select key content from the article	6		44	11	6
st tedio	Content review: including content inclusion selection, quality of paraphrasing	7		44	11	6
Lea	Read the article	8		34	11	6
	Least tedious					

Least tedious

Data Extraction

ROC curves are a fundamental way to evaluate classifier performance. AUC values can range from 0.5 to 1.0, with values closer to 1.0 indicating that the classifier's performance is better than random. Using 3187 sentence data points, ROC curves were generated to assess the classifier's performance (Figure 5). The model trained on treatment data had the highest AUC

(0.85). Models trained on Prodigy data achieved AUC scores of 0.84, 0.74, 0.7, 0.67, and 0.74 for adverse events, population, efficacy, end points, and study design, respectively. The AUC scores for all classifier models exceeded 0.5 (random chance) with P<.001. Overall, the model was able to classify sentences appropriately and be better than a random chance (indicated by the dashed line).

Figure 5. Receiver operating characteristic (ROC) curve of machine learning models' performance. AUC: area under the curve.



Data Abstraction and Filtering Output

The clinical trial data were inputted into the LLM model to produce an AI-generated summary and an AI-generated paraphrase of the text. Table 1 provides 2 examples of the paraphrasing process, showing the input in the LLM model and output using the prompt "Paraphrase this without plagiarizing."

In Table 1, the colors correlate to the BLEU scores, with the reddest color being 100 and the bluest color being 0 (as shown

by the scale on the right). The y-axis corresponds to the original sentences in the text (0, 1, 2, etc), and the x-axis corresponds to the paraphrased sentences (0, 1, and 2). Each original sentence is compared with each paraphrased sentence with a BLEU score. The lower the BLEU score, the fewer words and phrases there are in common between the original and the paraphrased sentences. Semantic similarity is the overall similarity between the two texts.



Table 1. Artificial intelligence (AI)-generated summarized and paraphrased text examples.

	1 1 1	
Original text input	AI-generated summarized text:	AI-generated paraphrased text:
• In the sitagliptin group and glimepiride group, serious adverse experiences were reported for 16 and 11 patients, respectively. One serious adverse experience that was considered drug-related by the investigator was reported in the study: a 61-year-old female in the sitagliptin group experienced a serious drug-related adverse experience of thrombocytopenia (platelet count of 12 000/µL) resulting in discontinuation from the study. Apart from hypoglycemia, the only other adverse event with an incidence of at least 5% in either treatment group was nasopharyngitis (4.8 and 6.9% in the sitagliptin and glimepiride groups, respectively). Overall, clinical adverse experiences were reported more frequently with glimepiride, primarily as a result of an increased incidence of hypoglycemia.	• In the sitagliptin group and glimepiride group, serious adverse experiences were reported for 16 and 11 patients, respectively. The only other adverse event with an incidence of at least 5% in either treatment group was nasopharyngitis (4.8 and 6.9%).	 In the sitagliptin group, 16 patients reported a serious adverse experience. In the glimepiride group, 11 patients reported a serious adverse experience. The only other adverse event with an incidence of at least 5% was nasopharyngitis. Semantic similarity: 0.76 max BLEU^a: 18.54
• Specific clinical adverse experiences of interest included hypoglycemia and prespecified selected gastrointestinal adverse experiences (abdominal pain, nausea, vomiting, and diarrhea). The incidences of serious clinical adverse experiences were similar between treatment groups. For serious adverse experiences that were considered by the investigator to be related to the study drug, there were three in the glipizide group (myocardial infarction, spontaneous abortion, and hydronephrosis) and one in the sitagliptin group (thrombocytopenia). The overall incidence of infection-related adverse experiences was similar in the two treatment groups. Of the adverse experiences with a higher incidence in the glipizide group, the 95% CI around the between-group difference in incidence excluded zero for cataracts, toothache, hypoglycemia, and hypoesthesia (Table 5). With the exception of hypoglycemia, these adverse experiences (occurring in either group) were generally rated as mild in intensity, not considered related to the study drug. AE, adverse experience. Overall, clinical adverse experiences and drug-related adverse experiences were observed more often with glipizide, largely related to an increased incidence of hypoglycemia. The number of serious adverse experiences was	• The incidences of serious clinical adverse experiences were similar between treat- ment groups. Serious ad- verse experiences that were considered by the investiga- tor to be related to the study drug were three in the glip- izide group (myocardial in- farction, spontaneous abor- tion, and hydronephrosis) and one in the sitagliptin group (thrombocytopenia). The overall incidence of in- fection-related adverse expe- riences was similar.	 There were similar incidences of serious clinical adverse experiences between treatment groups. The serious adverse experiences that were considered by the investigator to be related to the study drug were three in the glipizide group (myocardial infarction, spontaneous abortion, and hydronephrosis) and one in the sitagliptin group (thrombocytopenia). The overall incidence of infection-related adverse experiences was similar. Image: Semantic similarity: 0.86 max BLEU: 89.34

^aBLEU: bilingual evaluation understudy.

observed in either treatment group.

Semantic similarity and BLEU scores calculated for the paraphrases show that in the two example scenarios in Table 1, one has a low BLEU score, and the other has a high BLEU score. The high BLEU score in example 2 contains more original text in the paraphrase compared to the low BLEU score in example 1. Figures 6-8 depict the graphed comparison of BLEU score and semantic similarity among three reviewers (users) to

define usefulness. User 2 showed a preference for high semantic similarity, but there was no clear trend with the BLEU score. User 3 consistently favored paraphrases with both high semantic similarity and BLEU score. User 1 had no clear preference trend. Differences between what human reviewers found useful in paraphrases were noted.

Figure 6. Bilingual evaluation understudy (BLEU) score versus semantic similarity for user 1.





Figure 7. Bilingual evaluation understudy (BLEU) score versus semantic similarity for user 2.




Figure 8. Bilingual evaluation understudy (BLEU) score versus semantic similarity for user 3.



Discussion

Principal Findings

In the survey section of this study, we found that in the strategic activity of creating SRDs, the major challenge was in paraphrasing articles. In the subsequent phase of this study, traditional machine learning classifiers and LLMs automated portions of the clinical trial summarization process of creating an SRD.

Our survey revealed a much shorter time (8.7 hours and 3.8 hours) to create or revise an SRD compared with the 2018 phactMI benchmarking survey (31 hours and 21 hours) [2]. The variations and limited external validity of the overall survey may be attributed to the nature of the survey, the number of responses, and survey types. Nevertheless, the survey's results continue to be valuable, as they offer nuanced insights from

```
https://ai.jmir.org/2025/1/e55277
```

RenderX

engaged participants and contribute to our understanding. Regardless of the amount of time, providing solutions to improve the efficiency of creating an SRD would be welcomed.

Data Extraction

Our study reveals the promise of machine learning models in classifying individual sections within scientific documents, particularly in the context of addressing inquiries within the pharmaceutical industry. The results from the ROC curves suggest that our classifier models outperform random guessing, demonstrating the highest AUC values for the treatment and adverse events classifiers. The transparency and interpretability of our classifier models were pivotal strengths. Unlike LLMs, which are known for their opaqueness in decision-making, our traditional machine learning models have successfully identified and resolved training logic deviation issues. Having clear

explanations in the output is invaluable for trust, accountability, and enhancing the models.

In addition, the classifier models exhibited resource efficiency. While finetuning LLMs is a resource-intensive process, we found that adjusting the logistic regression model can be executed in seconds. This efficiency has major implications for rapid model development and deployment.

Enhancing classifier performance necessitates the consideration of several key factors that the working group identified. These factors include providing additional context, predefining known key terminology for specific sections, and exploring methods to reduce false negatives. In future iterations, our approach will expand the scope of classifiers beyond section identification to assess the usefulness of the identified content for inclusion in an SRD. This transition marks a shift from mere classification to a more profound evaluation of content, offering applications in content retrieval tailored to individual user needs.

Data Abstraction and Filtering Output

Our exploration of paraphrasing performance in LLMs has been highly informative. Quantitative assessment of paraphrased content requires robust tools like semantic similarity and BLEU scores. By leveraging these tools, we gain a deeper understanding of the effectiveness of paraphrasing, ensuring that content retains its intended meaning while being substantially different from the original in terms of phrasing or wording.

The observed variability in LLM-generated paraphrases highlights the difficulty of consistently fine-tuning an LLM for paraphrasing. The diverse approaches to paraphrasing are highlighted by the distinct preferences of human reviewers. Developing a universal model for all preferences is an ambitious endeavor. The working group proposed an alternative approach to this challenge: using simple models that offer users multiple paraphrase options. We can enhance the content ranking and establish core data by providing choices and using smaller datasets, as user selections can potentially be used to train classifiers to identify the kind of content that the user prefers.

The working group also recommends the following next steps with LLMs to further this exercise: (1) fine-tuning an LLM for medical text, (2) better prompt engineering, and (3) LLMs with better citation training. Incorporating these considerations into our discussion of paraphrasing performance and prospects, we navigate the evolving landscape of AI-driven content generation in the pharmaceutical industry. These insights not only promise enhanced content but also embody a user-centric approach that empowers industry professionals to access tailored, high-quality content.

Need for Human Control in AI-Assisted Scientific Writing

A recent study used ChatGPT to obtain medical information and treatment options for shoulder impingement syndrome [14]. While ChatGPT's answers were useful for patients, it sometimes provided inaccurate information (prevalence reported with no evidence supporting the number) and biased information (risk factors reported that are not established). Goodman et al [15]

```
https://ai.jmir.org/2025/1/e55277
```

conducted a cross-sectional study corroborating these limitations of LLMs. Most responses were accurate and comprehensive, indicating the potential use of LLMs. Occasionally, incorrect answers were provided, and the chatbot provided inaccurate citations when asked for the source of information. Other studies have demonstrated similar drawbacks (misinterpretation of medical terms, hallucination, missed information, factually incorrect statements, and fabricated references) in the use of LLMs in scientific writing and simplified radiology reports [13,15,16]. Accuracy, lack of bias, and traceability to the original publication are crucial in medical information. Thus, using LLMs without considerable human intervention for medical information responses or SRDs is a highly risky proposition. While AI can help humans create a "first draft" of the final SRD, it is imperative for the human writer to retain control over the tool's input, data extraction for the SRD, and the ultimate inclusion of paraphrased content in the SRD. Our approach includes various "checkpoints" during AI-assisted SRD creation, allowing human writers to intervene and enhance the content's credibility.

The use of LLMs for scientific writing also presents concerns regarding plagiarism and the use of nonacademic language [13,17]. In addition, LLMs are unable to determine the credibility of their information sources, for example, a blog post versus a PubMed-indexed paper [15]. Our model can overcome numerous limitations by integrating machine learning and LLM systems.

Limitations

Despite the working group's diligent effort to maintain scientific rigor in this study, several limitations warrant consideration. The classical machine learning classifiers may have biased models due to training on a constrained dataset and limited reviewer assessments. Instead of relying on experts to label more examples, it may be more efficient to extract labeled examples from existing datasets (eg, adverse events sections from full-text papers in PubMed Central). The use of LLMs like GPT presented known challenges for paraphrasing medical text, such as generative AI issues of "hallucination," lack of transparency, bias, and privacy concerns [18].

The dynamic generative AI landscape implies that the findings of paraphrase exercises only reflect a snapshot in time. OpenAI introduced GPT-4 Turbo, a 2023 model trained on a larger dataset, while we were drafting this manuscript [13]. Nori et al [19] demonstrated that prompt engineering with GPT-4 outperformed fine-tuned medical models for question answering. The framework described in this paper is similar to the emerging pattern of retrieval augmented generation [20] in leveraging LLMs. The focus of retrieval augmented generation is to provide the LLM with accurate, up-to-date information [20]. The same business drivers from the medical information space prompted this evolution, driven by a need for accuracy and content lineage tracking. The fact that several others have reached a similar conclusion on integrating LLMs into highly regulated industries such as drug manufacturing is a strong validation.

Conclusions

This study sought to identify the challenges inherent in the development of SRDs and to establish a framework for integrating LLM and machine learning into the SRD creation process. Our tool leverages LLMs and machine learning to enhance AI applications in the pharmaceutical realm. Integrating these two technologies not only saves resources but also addresses major challenges associated with LLMs. Our models can clearly identify sections, paraphrase effectively, and assess content usefulness. These initial findings suggest that machine learning classifiers can predict, to some extent, the sentences authors will choose for summarization and paraphrases they will find useful. Even a modest ability to rank results could

improve the suggestions' quality beyond random. However, the current tool does not have the capacity to generate an SRD for the pharmaceutical sector using zero-shot classification. Nevertheless, it underscores the essential role of traditional machine learning in enhancing future AI models, moving us closer to efficient content handling in the industry. This model has the potential to be a valuable tool in the medical information domain of the pharmaceutical industry, augmenting the efficiency of human document creators, thereby optimizing workflows and improving the quality of services. Further research is required for the optimization, refinement, and validation of these models, using larger training sets and multiple reviewers, before full-scale implementation in the industry.

Acknowledgments

This research was sponsored by phactMI and Microsoft. No generative artificial intelligence was used in the development of this manuscript. The authors would also like to acknowledge Mario E Inchiosa's contribution to the study.

Data Availability

All data generated or analyzed during this study are included in this published article and its supplementary information files.

Authors' Contributions

All authors contributed to the conceptualization, formal analysis, investigation, methodology, resources, and writing-review and editing. In addition, JL and RH were responsible for data curation; JJ for funding acquisition; JL for project administration; RH for software; EH-D for supervision; SG, AC, RH, and SB for validation; RH and EH-D for visualization; and JL, EH-D, SB, and RH for writing-original draft.

Conflicts of Interest

JL was a fellow at phactMI 2022-2024. AC is an employee of Pfizer (and owns stock) and is on the Board of Directors for phactMI. SB is an employee of Eli Lilly and Company. RH is an employee of Win-Vector Labs and, at the time of the research, was employed by Microsoft. None declared by other authors.

Multimedia Appendix 1 Email invitation for survey participation. [DOCX File , 15 KB - ai v4i1e55277 app1.docx]

Multimedia Appendix 2 Survey questions. [DOCX File, 16 KB - ai v4i1e55277 app2.docx]

Multimedia Appendix 3 Checklist for Reporting Results of Internet E-Survey (CHERRIES). [DOCX File, 22 KB - ai v4i1e55277 app3.docx]

Multimedia Appendix 4 Training dataset. [ZIP File (Zip Archive), 618 KB - ai v4i1e55277 app4.zip]

References

- Hermes-DeSantis ER, Johnson RM, Redlich A, Patel B, Flanigan-Minnick A, Wnorowski S, et al. Proposed best practice guidelines for scientific response documents: a consensus statement from phactMI. Ther Innov Regul Sci 2020;54(6):1303-1311. [doi: 10.1007/s43441-020-00151-1] [Medline: 33258092]
- 2. Patel M, Jindia L, Fung S, Kadowaki R, Marasigan K. Pharma collaboration for transparent medical information (phactMITM) benchmark study: trends, drivers, and value of product support activities, key performance indicators, and other medical

information services: insights from a survey of 27 US pharmaceutical medical information departments. Ther Innov Regul Sci 2020;54(6):1275-1281. [doi: 10.1007/s43441-020-00162-y] [Medline: 32447658]

- Albano D, Pragga F, Rai R, Flowers T, Parmar P, Wnorowski S, et al. The medical information scientific process: define, research, evaluate, synthesize, and share (DRESS). Ther Innov Regul Sci 2022;56(3):405-414 [FREE Full text] [doi: 10.1007/s43441-021-00366-w] [Medline: 35239132]
- 4. Andonian A. Emergent Capabilities of Generative Models: "Software 3.0" and Beyond. Cambridge, MA: Massachusetts Institute of Technology; 2021:95.
- Reimers N, Gurevych I. Sentence-BERT: sentence embeddings using Siamese BERT-networks. 2019 Presented at: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP); November 3-7, 2019; Hong Kong, China p. 3982-3992. [doi: <u>10.18653/v1/d19-1410</u>]
- 6. What is the transformer architecture and how does it work. DataGen Technologies. URL: <u>https://datagen.tech/guides/</u> <u>computer-vision/transformer-architecture/#</u> [accessed 2023-12-06]
- 7. Introduction. OpenAI API. URL: <u>https://platform.openai.com/docs/api-reference/introduction%E3%80%82</u> [accessed 2023-11-16]
- Ramesh AN, Kambhampati C, Monson JRT, Drew PJ. Artificial intelligence in medicine. Ann R Coll Surg Engl 2004;86(5):334-338. [doi: <u>10.1308/147870804290</u>] [Medline: <u>15333167</u>]
- 9. Dergaa I, Chamari K, Zmijewski P, Ben Saad H. From human writing to artificial intelligence generated text: examining the prospects and potential threats of ChatGPT in academic writing. Biol Sport 2023;40(2):615-622 [FREE Full text] [doi: 10.5114/biolsport.2023.125623] [Medline: 37077800]
- 10. Alston E. What are AI hallucinations-and how do you prevent them? Zapier. URL: <u>https://zapier.com/blog/ai-hallucinations/</u> [accessed 2023-04-05]
- 11. Text generation. OpenAI API. URL: https://platform.openai.com/docs/guides/text-generation [accessed 2023-11-16]
- 12. van Dis EAM, Bollen J, Zuidema W, van Rooij R, Bockting CL. ChatGPT: five priorities for research. Nature 2023;614(7947):224-226. [doi: <u>10.1038/d41586-023-00288-7</u>] [Medline: <u>36737653</u>]
- Blanco-González A, Cabezón A, Seco-González A, Conde-Torres D, Antelo-Riveiro P, Piñeiro Á, et al. The role of AI in drug discovery: challenges, opportunities, and strategies. Pharmaceuticals (Basel) 2023;16(6):891 [FREE Full text] [doi: 10.3390/ph16060891] [Medline: 37375838]
- Kim JH. Search for medical information and treatment options for musculoskeletal disorders through an artificial intelligence chatbot: focusing on shoulder impingement syndrome. medRxiv Preprint posted online on December 19, 2022. [doi: <u>10.1101/2022.12.16.22283512</u>]
- Goodman RS, Patrinely JR, Stone CA, Zimmerman E, Donald RR, Chang SS, et al. Accuracy and reliability of chatbot responses to physician questions. JAMA Netw Open 2023;6(10):e2336483 [FREE Full text] [doi: 10.1001/jamanetworkopen.2023.36483] [Medline: <u>37782499</u>]
- Jeblick K, Schachtner B, Dexl J, Mittermeier A, Stüber AT, Topalis J, et al. ChatGPT makes medicine easy to swallow: an exploratory case study on simplified radiology reports. Eur Radiol 2024;34(5):2817-2825 [FREE Full text] [doi: 10.1007/s00330-023-10213-1] [Medline: <u>37794249</u>]
- 17. Aydın ?, Karaarslan E. OpenAI chatGPT generated literature review: digital twin in healthcare. In: Emerging Computer Technologies. Büyükkale, Turkey: İzmir Akademi Dernegi; 2022:22-31.
- Johnson D, Goodman R, Patrinely J, Stone C, Zimmerman E, Donald R, et al. Assessing the accuracy and reliability of AI-generated medical responses: an evaluation of the chat-GPT model. Res Sq :28 Preprint posted online on February 28, 2023 [FREE Full text] [doi: 10.21203/rs.3.rs-2566942/v1] [Medline: 36909565]
- 19. Nori H, Lee Y, Zhang S, Carignan D, Edgar R, Fusi N, et al. Can generalist foundation models outcompete special-purpose tuning? Case study in medicine. arXiv Preprint posted online on November 28, 2023. [doi: <u>10.48550/arXiv.2311.16452</u>]
- 20. Radhan R. Addressing AI hallucinations with retrieval-augmented generation. InfoWorld. 2023. URL: <u>https://www.infoworld.com/article/3708254/addressing-ai-hallucinations-with-retrieval-augmented-generation.html</u> [accessed 2025-02-06]

Abbreviations

AE: adverse experience AI: artificial intelligence AUC: area under the curve BLEU: bilingual evaluation understudy CHERRIES: Checklist for Reporting Results of Internet E-Surveys FPR: false positive rate GPT: Generative Pre-trained Transformer LLM: large language model ROC: receiver operating characteristic SRD: scientific response document

https://ai.jmir.org/2025/1/e55277

TPR: true positive rate

Edited by K El Emam; submitted 07.12.23; peer-reviewed by S Fung, TAR Sure, S Kommireddy, M Jovanovik; comments to author 23.05.24; revised version received 01.07.24; accepted 31.12.24; published 13.03.25. <u>Please cite as:</u> Lau J, Bisht S, Horton R, Crisan A, Jones J, Gantotti S, Hermes-DeSantis E Creation of Scientific Response Documents for Addressing Product Medical Information Inquiries: Mixed Method Approach Using Artificial Intelligence JMIR AI 2025;4:e55277 URL: https://ai.jmir.org/2025/1/e55277 PMID:

©Jerry Lau, Shivani Bisht, Robert Horton, Annamaria Crisan, John Jones, Sandeep Gantotti, Evelyn Hermes-DeSantis. Originally published in JMIR AI (https://ai.jmir.org), 13.03.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Disease Prediction Using Machine Learning on Smartphone-Based Eye, Skin, and Voice Data: Scoping Review

Research Dawadi^{1,2}, PhD; Mai Inoue^{1,2}, ME; Jie Ting Tay^{1,2}, MRES; Agustin Martin-Morales^{1,2}, PhD; Thien Vu^{1,2}, MD, PhD; Michihiro Araki^{1,2,3,4}, PhD

¹Artificial Intelligence Center for Health and Biomedical Research, National Institutes of Biomedical Innovation, Health and Nutrition, Osaka, Japan ²National Cerebral and Cardiovascular Center, Osaka, Japan

³Faculty of Medicine, Graduate School of Medicine, Kyoto University, Kyoto, Japan

⁴Graduate School of Science, Technology and Innovation, Kobe University, Kobe, Japan

Corresponding Author: Research Dawadi, PhD Artificial Intelligence Center for Health and Biomedical Research National Institutes of Biomedical Innovation, Health and Nutrition 3-17, Senrioka-Shinmachi, Settsu Osaka, 566-0002 Japan Phone: 81 661701069 ext 31234 Email: dawadi-research@nibiohn.go.jp

Abstract

Background: The application of machine learning methods to data generated by ubiquitous devices like smartphones presents an opportunity to enhance the quality of health care and diagnostics. Smartphones are ideal for gathering data easily, providing quick feedback on diagnoses, and proposing interventions for health improvement.

Objective: We reviewed the existing literature to gather studies that have used machine learning models with smartphone-derived data for the prediction and diagnosis of health anomalies. We divided the studies into those that used machine learning models by conducting experiments to retrieve data and predict diseases, and those that used machine learning models on publicly available databases. The details of databases, experiments, and machine learning models are intended to help researchers working in the fields of machine learning and artificial intelligence in the health care domain. Researchers can use the information to design their experiments or determine the databases they could analyze.

Methods: A comprehensive search of the PubMed and IEEE Xplore databases was conducted, and an in-house keyword screening method was used to filter the articles based on the content of their titles and abstracts. Subsequently, studies related to the 3 areas of voice, skin, and eye were selected and analyzed based on how data for machine learning models were extracted (ie, the use of publicly available databases or through experiments). The machine learning methods used in each study were also noted.

Results: A total of 49 studies were identified as being relevant to the topic of interest, and among these studies, there were 31 different databases and 24 different machine learning methods.

Conclusions: The results provide a better understanding of how smartphone data are collected for predicting different diseases and what kinds of machine learning methods are used on these data. Similarly, publicly available databases having smartphone-based data that can be used for the diagnosis of various diseases have been presented. Our screening method could be used or improved in future studies, and our findings could be used as a reference to conduct similar studies, experiments, or statistical analyses.

(JMIR AI 2025;4:e59094) doi:10.2196/59094

KEYWORDS

literature review; machine learning; smartphone; health diagnosis

Introduction

The use of machine learning for medical diagnosis is steadily growing. This can be attributed primarily to the availability of

https://ai.jmir.org/2025/1/e59094

RenderX

numerous health data as well as improvements in the classification and recognition systems used in disease diagnosis. The health care industry produces an abundance of health-related data [1], which can be used to create machine learning models.

These models can be used for diagnosing and predicting a variety of diseases, including breast cancer, heart diseases, and diabetes [2,3]. The prediction of these diseases is dependent on many factors according to the focus on different features (biomarkers) [4]. The application of machine learning methods helps classify and diagnose diseases in an easier way [1], and these diagnoses can help medical experts in the early detection of fatal diseases and therefore increase the quality of health care and the survival rate of patients significantly [1,2,4,5].

Machine learning methods and their applications are not limited to particular types of data and thus have been used in a variety of areas, such as detecting spontaneous abortion [6], identifying complex patterns in brain data [7], and improving diagnostic accuracy and identifying faults in axial pumps [8]. To diagnose and predict different diseases, machine learning methods have also been applied to data obtained from experiments by using publicly available datasets, such as the UCI machine learning library [2], National Health and Nutrition Examination Survey (NHANES) [3,6], traumatic brain injury (TBI) [4], and SUITA datasets [9]. Similarly, numerous smartphone-based health care apps have been developed to help both health care officials and the general population with regard to their health-related concerns. The apps developed can be broadly divided into 3 specific user groups: health care professionals, medical/nursing students, and patients [10]. The purpose of such apps covers a wide range of areas, such as disease diagnosis, drug reference, medical education, clinical communication, and fall detection. However, all iOS- or Android-based apps developed for health care purposes have not been discussed in the literature [10].

A literature review is a systematic way of collecting studies relevant to a research topic, assessing the methodologies and results of the studies, and making recommendations for improvements if necessary [11]. In the health care domain, the implementation of literature reviews has been considered important for conducting further research and developing guidelines for clinical practice [12]. Literature studies, such as umbrella reviews, have been conducted to study the management of the information of patients, such as those with cancer, and how their records are handled [13]. Similarly, literature-based studies have investigated the evidence of leadership in nursing [14]. Uddin et al [15] found a total of 48 literature studies that dealt with disease prediction using various supervised machine learning algorithms and attributed the rise in the use of machine learning for health prediction to the wide adoption of computer-based technologies in the health sector and to the availability of large health-related databases.

The ubiquity of smartphones makes them a convenient tool to gather various health-related data, particularly as smartphones are equipped with various sensors that are able to track and gather different health-related information [16]. However, there is a lack of research on studies involving the adoption of smartphones for disease prediction using machine learning methods and identifying the types of experiments conducted, databases utilized, and machine learning methods used. With that in mind, in this paper, we aim to conduct a scoping review by assessing research papers from repositories, such as PubMed and IEEE Xplore, which have used machine learning methods with smartphone-derived data to predict diseases related to the

```
https://ai.jmir.org/2025/1/e59094
```

XSI•FC

eyes, skin, and voice, and from databases available for public use. We aim to answer the following important research questions:

- 1. What are the databases available for eye-, skin-, and voice-related diseases?
- 2. What are the machine learning models used in such studies?
- 3. How the data are collected using smartphones?

The rest of the paper is organized as follows: we explain how we gathered, screened, and analyzed the literature in the methods section; present the results of our study in the results section; and finally discuss the results and clarify how the results correspond to our research questions in the discussion section.

Methods

Overview

We describe in detail the procedures undertaken for conducting the scoping review, with inspiration taken from the guidelines provided by Mak and Thomas [17]. After deciding on the topic of research, we identified the steps to be taken for the literature review as follows:

- 1. Search criteria
- 2. Literature assembly
- 3. Study selection
- 4. Research questions
- 5. Inclusion and exclusion criteria
- 6. Full-text paper assessment

Search Criteria

Numerous studies have conducted literature reviews to assess the use of machine learning for disease prediction. We formulated the following search string using a combination of different words related to topics, such as smartphone, smartwatch, machine learning, health, and medicine, to search different electronic databases: (*(ML OR machine learning) AND (health* OR medic* OR disease) AND (smartphone OR smart phone OR smartwatch OR smart watch OR smart devices)*).

Before finalizing the search string, we experimented with many combinations, including different variations of specific keywords and symbols, such as "*," to cover a wider area and maximize the results.

Literature Assembly

We applied the search string to different databases and narrowed the databases to PubMed [18] and IEEE Xplore [19]. The search results from other databases produced a very high number of results that included unnecessary papers from disciplines unrelated to our topic of interest. When we applied the same search string, ACM Digital Library had about 24,000 results, ProQuest had about 125,000 results, and Google Scholar had more than 1 million results. Furthermore, in Science Direct, our search string did not produce any results owing to the use of Boolean connectors. We concluded that the search of PubMed and IEEE Xplore was enough to obtain papers related to research in technology, engineering, and biomedical sciences.

The results from each of the databases were then exported to an external file. To convert the results from the 2 databases into

a single file, including the titles and abstracts, we used Mendeley [20] and Zotero [21]. The file, which contained a total of 2390 papers, was then screened in the Jupyter Notebook environment using Python (version 3.7.17) [22].

Study Selection

For refining the collected papers, we used a title screening method [23] to filter out papers that might not be of direct relation to our research topic. We created a list of keywords that match the research topic, screened the titles of all papers, and filtered out all papers that did not contain any of the following keywords: machine, artificial, smartphone, disease, mobile, health, healthcare, wearable, model, features, and training.

The identified papers at this point covered a wide variety of diseases and health areas. Using the keyword identification

Table 1. Health categories in titles.

method, we tried to find the distribution of different diseases in the papers based on the 5 senses [24]. We first determined the frequency of keywords related to the 5 senses in the titles of the collected papers by using the following keywords: eye, eyesight, vision, audio, voice, vocal, nasal, nose, hearing, ear, touch, feel, face, skin and dermatology.

The result for the frequency of the keywords in the titles can be seen in Table 1. We then merged the keywords with their respective senses and assembled the papers into the following 6 categories: ear, eye, nose, touch, skin, and audio. We determined the total distribution of the papers, as shown in Table 2. We replicated the procedure to determine the frequency of health categories (Table 3) and their distribution (Table 4) in the abstracts of the collected papers.

Health care area	Number of matches							
Eye	12							
Eyesight	0							
Vision	13							
Audio	15							
Voice	16							
Vocal	5							
Nasal	0							
Nose	2							
Hearing	5							
Ear	2							
Touch	6							
Feel	0							
Face	13							
Skin	19							
Dermatology	2							

Table 2. Distribution of health categories in titles.

Category	Distribution, %
Voice	32.7
Nose	1.8
Ear	6.4
Eye	22.7
Touch	5.5
Skin	30.9



Table 3. Health categories in abstracts.

Health care area	Number of matches						
Eye	108						
Eyesight	1						
Vision	142						
Audio	106						
Voice	141						
Vocal	24						
Nasal	4						
Nose	18						
Hearing	28						
Ear	30						
Touch	32						
Feel	5						
Face	130						
Skin	162						
Dermatology	20						

Table 4. Distribution of health categories in abstracts.

Category	Distribution, %
Voice	28.5
Nose	2.3
Ear	6.1
Eye	26.4
Touch	3.9
Skin	32.8

Research Questions

Based on the results from Tables 1-4, we identified the following 3 categories with the highest distribution of papers: eye, skin, and voice, and formulated the following research questions:

- 1. What are the databases available for eye, skin, and voice analysis?
- 2. What are the machine learning models used for eye, skin, and voice analysis?
- 3. How are the data collected from smartphones?

The keyword screening method [23] was applied to the titles of 2390 papers, which resulted in the successful screening of 2352 papers. In the next step, we screened the abstracts of the papers to distinguish papers related to each of the 3 topics (eye, skin, and voice) by using relevant keywords.

Inclusion and Exclusion Criteria

The primary inclusion criterion was that the study should perform an experiment or use a database involving data obtained by using smartphones. Some studies conduct experiments themselves to gather data from participants, while others use publicly available datasets. We divided the studies based on this distinction (experiments and databases). This information

https://ai.jmir.org/2025/1/e59094

RenderX

can help researchers determine if they want to conduct similar experiments or simply use publicly available databases.

Since the search terms specified the use of both smartphones and machine learning methods, it reduced the probability of obtaining literature results related to topics other than disease prediction among humans. The other criteria for the articles were that they should be in the English language and should be available for full-text viewing. Studies that involved data collection with external devices other than smartphones and those that used only smartwatches and not smartphones were excluded. Furthermore, studies that were literature reviews were not included in the final analysis.

Full-Text Assessment

The inclusion and exclusion criteria were applied to 217 papers available after title and abstract screening. After assessment of these papers, there were 8, 14, and 38 studies related to the skin, eyes, and voice, respectively. We performed full-text analysis of these papers to extract the desired information.

Results

Overview

We explain the analysis of papers that were extracted and report about the databases used, experiments conducted, and machine learning methods used. The steps and results of our review process can be seen in Figure 1.

Figure 1. Flow diagram for identifying relevant literature.



Research on Voice

Owing to the recent global pandemic, research on the analysis of speech for either cough or COVID-19 has grown [25]. Apart from that, the analysis of audio has a wide range of applications from the prediction of emotional stress [26] and the detection of diseases, such as Parkinson disease [27], to the detection of tourist emotions for spot recommendation [28].

Studies Conducted Using Databases

A cough-based COVID-19 detection model was created using more than 25,000 cough recordings from the CoughVid dataset [25]. The dataset was created through recordings via a web interface that could be accessed by a personal computer or a smartphone, and the prediction was made using a stack ensemble classifier consisting of machine learning methods, such as decision tree (DT), random forest (RF), k-nearest neighbor (KNN), and extreme gradient boosting (XGBoost).

Since datasets containing the voices of COVID-19–affected patients were not in abundance, datasets with recordings of cough sounds along with sneezing, speech, and nonvocal audio were used to pretrain the classifier [29]. Brooklyn and Wallacedene datasets used for the training were created using an external microphone, while datasets, such as TASK, were

https://ai.jmir.org/2025/1/e59094

created using an external microphone along with a smartphone. It is very likely that smartphones were used to create datasets, such as the Google audio dataset and Freesound, which consist of audio from more than 1.8 million YouTube videos. Similarly, the Librispeech dataset consists of audio from 56 speakers who may or may not have used smartphones. For the classification and testing of the model, 3 datasets, namely Coswara, ComparE, and Sarcos, were used by applying machine learning methods, such as convolutional neural network (CNN), long short-term memory (LSTM), and RestNet50. All 3 datasets were created with the recordings of the cough of participants. ComparE and Coswara consist of additional speech sounds, with the Coswara dataset also including breathing sounds. The data acquisition method was web-based, and thus, smartphones could have been used for recording such audio data.

The Coswara dataset, with recordings of over 1600 participants, was created by collecting breathing, coughing, and voice sounds, using the microphones of smartphones via an interactive website application. With a combination of hand-crafted features and deep-activated features learned through model training, a deep learning framework was proposed and studied by using the recordings of 240 participants from the Coswara dataset (120 participants were identified as positive for COVID-19) [30].

The dataset created in the mPower study [31], which was conducted for the detection of Parkinson disease using audio data, has been used in various other studies [27,28]. The dataset was divided into training and test sets, and 2 classifiers, namely support vector machine (SVM) and RF, were applied to compare 6 cross-validation techniques [32]. Similarly, a desktop application, PD Predict, that records audio and makes predictions was created using the mPower dataset [33]. Two machine learning classifiers were used: gradient boosting classifier (GBC) pipeline with Lasso (gbcpl) and GBC pipeline with ElasticNet (gbcpen).

Moreover, using a dataset containing 18,210 recordings from the mPower study [31], a Parkinson disease prediction model was created through 4 classifiers: SVM, KNN, RF, and XGBoost [32]. Another Parkinson disease prediction model was created with 2 databases: PC-GITA and Vishwanathan [34] by using SVM. For creating the PC-GITA dataset, a smartphone was used to record 100 Columbian-Spanish speakers, among whom 50% had Parkinson disease. Similarly, 46 participants, among whom 24 were diagnosed with Parkinson disease, were used to create the Vishwanathan dataset, which consists of recordings of utterances of the alphabets "a," "u," and "m."

Five machine learning models, namely logistic regression (LR), RF, XGBoost, CatBoost, and Multilayer, were used to predict the emotional state of participants [35]. The dataset Extrasensory was used for training the model. The dataset was created using data from smartphones and smartwatches. Contextual data, such as location, phone state, accelerometer data, and light and temperature data, and emotional state information (disclosure of emotion at different intervals using a smartphone app) were collected.

Studies Conducted Through Experiments

A total of 1513 subjects above 50 years of age, including healthy subjects and subjects who were diagnosed with Parkinson disease, used a smartphone app to complete daily surveys and 4 activities intended to test the presence or effect of Parkinson disease [31]. The activities included tapping (tap 2 buttons alternatively), walking (walk in a straight line for 20 steps and back in the same route), voice (10-second utterance of the "aaah" sound), and memory (recall the order of illumination of flowers shown in the app). The data related to the accelerometer, gyroscope, touchscreen, and microphone were then collected to test the results of these activities. LR, RF, deep neural network (DNN), and CNN were used separately and as multi-layer classifiers for model creation and verification.

An Android-based smartphone app was developed to record 5 activities (voice, finger tapping, gait, balance, and reaction time) in 129 participants, including subjects who were healthy and those who were diagnosed with Parkinson disease, in order to study the effects of the disease [26]. Disease severity score learning (DSSL), a rank-based machine learning algorithm scaled from 0 to 100 (higher numbers reflect increasing severity of the disease), was used to show the results. In another study, 2 vocal tasks of patients diagnosed with Parkinson disease were recorded in a soundproof booth: one in which the participants spoke the vowel "a" for 5 seconds, and another in which the participants spoke a sentence in their native Lithuanian language

```
https://ai.jmir.org/2025/1/e59094
```

[36]. The recordings were conducted using both an external microphone and a smartphone, and the model was created using RF.

In another study, 237 participants diagnosed with Parkinson disease performed 7 smartphone-based tests, such as pronouncing "aaah" on the smartphone for as long as possible, pressing a button on the screen if it appears, pressing 2 alternate buttons on the screen, and holding the phone with their hand at rest or outstretched. Their balance and gait were also analyzed from the position of the smartphone [37]. The data obtained from the smartphone were used to train the machine learning algorithm using RF. The dataset was divided into training and test sets, and the prediction accuracy was tested using 10-fold cross-validation and leave-one-out cross-validation.

In addition to voice, facial features can be used for the detection of Parkinson disease [38]. Using both facial and audio data from 371 participants, among whom 186 were diagnosed with Parkinson disease, it was observed that early-stage detection of Parkinson disease is possible by combining both data. Participants were asked to read an article containing 500 words, and an iPhone was used to record both audio and video. DT, KNN, SVM, LR, naive Bayes, RF, LR, gradient boosting (GBoost), adaptive boosting (AdaBoost), and light gradient boosting (LGBoost) were compared to assess their performance in terms of accuracy, precision, recall, F_1 -score, and area under the receiver operating characteristic curve for binary classification.

Analysis of voice samples can also help in the prediction of depression or anxiety. A study was conducted with 2 sets of participants: one set of participants who had a diagnosis of depression or anxiety and another set of participants who did not have such a diagnosis [39]. Using an app developed for the study called Ellipsis Health, 5-minute voice samples and responses to survey questions were collected from a total of 263 participants (all current patients of a health care clinic) over a period of 6 weeks. Using a model developed with LSTM, the study tested the feasibility of assessing the presence of clinical depression and anxiety by using data from the smartphone app. With a similar approach, another study collected answers to questions in several self-reported psychiatric scales and questionnaires via audio recordings from 124 participants using an Android app developed specifically for the study [40]. Six different algorithms (LR, RF, SVM, XGBoost, KNN, and DNN) were used to study the features generated from audio and to evaluate the results.

In another study, behavioral and physiological data were collected from 212 participants through wearable sensors (including a wristband and a biometric tracking garment) and various surveys to create a dataset of human behaviors. The dataset was then studied to predict the emotional state of the participants [41]. A phone, Unihertz Jelly Pro, was also provided to participants to capture their speech data. An app, TILES, was created to track activities as well as receive responses to surveys. Furthermore, data from other smartphone apps, such as the Fitbit app (to receive updates from the Fitbit wristband), OMsignal app (to record data from the OMsignal smart garment), and

XSL•FO RenderX

Using only speech data, an automatic depression detection model was developed using deep convolutional neural network (DCNN) [27]. A total of 318 participants (153 diagnosed with major depressive disorder) were asked to record their voices through a smartphone while reading a predefined text. RF, SVM, KNN, and linear discriminate analysis classifiers were used, with RF providing the best accuracy. A similar study was conducted with 163 participants (88 diagnosed with depression), in which speech data were collected using VoiceSense, a voice-collection app installed on each participant's phone, through vocal responses to 9 general questions [42]. A repeated random subsampling cross-validation method, with random split of the dataset into training and test subsamples and multiple iterative repeats of the process, was used to obtain a predictive equation. Behavioral states of infants can also be predicted by analyzing the audio of their cries. About 1000 cries gathered from 691 infants using the smartphone app ChatterBaby were analyzed and classified into 3 states: fussy, hungry, and pain, using RF [43]. The study also aimed to verify that colic cries may indicate pain and are more similar to pain cries compared with either fussy or hungry cries.

Along with emotional states, it is also possible to create models to predict complex psychiatric conditions, such as schizophrenia, by using data from smartphones. Numerous data were collected from 61 participants, including app usage, reception of calls and SMS text messages, smartphone acceleration data, screen on/off duration, location, speech and conversation, sleep, and ambient environment, and an ecological momentary assessment was performed every 2 to 3 days [44]. Multiple-output support vector regression (m-SVR) and multi-task learning (MTL) with leave-one-out cross-validation were used to train data for each patient and to predict the scores for all possible symptoms.

Audio data can also be used to predict health-related anomalies, such as fatigue level and blood pressure. Using 1772 voice recordings from 296 participants, a model was created to predict fatigue in people affected with COVID-19 [45]. Two types of audio data were collected: recording of participants reading a predefined text and another recording of them pronouncing the vowel "a" for as long as they could. The data were trained and tested using LR, KNN, SVM, and soft voting classifier algorithms. In another study, a stethoscope attached to a smartphone was used to collect heart sound signals from 32 healthy subjects, with the participants laying on a mattress and the stethoscope being placed on their chest [46]. SVM was used for training and testing the estimation model, and 10-fold cross-validation was used to test the accuracy of the model.

Other uses of audio analysis include the inspection of bowel sounds for tracking or predicting digestive diseases [47]. A total of 100 participants were asked to put the smartphone over the lower right and left areas of their abdomen to collect audio via a bowel sound recording app. CNN- and LSTM-based recognition models were developed. For cross-validation, multiple training-test splits were conducted, and 9-fold cross-validation was performed. Furthermore, it is also possible to determine the quality of sleep by analyzing the audio during

XSL•FO

sleep. Using an app that records audio with the built-in microphones of smartphones and a smart alarm, sleep events, such as snoring and coughing, were identified [48]. SleepDetCNN, a CNN-based model, was created to classify the sleep audio into 3 types: snoring, coughing, and others. Snoring was further studied in 16 patients with habitual snoring tendencies using a smartphone-based gaming app as a treatment for snoring [49]. A section of participants had 15 minutes of daily gameplay (3 voice-controlled games; 5 minutes each), and the majority of participants were provided with microphones to record their sleep for the entire night at least twice per week during the experiment period of 12 weeks. To train the classification models using SVM, 1000 sleep sounds were randomly selected and labeled as "snore" or "not snore" by 2 blinded members of the research team.

In addition to the prediction of diseases, data from smartphone microphones, combined with other data, such as accelerometer, gyroscope, light proximity, and Wi-Fi scan data, have been utilized for emotion prediction [39] as well as the recognition of day-to-day activities [50]. The ADL Recorder app, created for tracking and monitoring the activities of elderly people, recorded both behavioral and contextual data. Various kinds of machine learning classifiers, such as Bayesian network, hidden Markov model, Gaussian mixture model, RF, and KNN, were used throughout for analyzing data from different sensors, and J48 DT was used for the final recognition of activities.

Studies Conducted With Both Databases and Experiments

Due to recent global events, numerous experimental studies have been conducted for COVID-19 detection and prediction. Audio data from 497 participants, including those with and without COVID-19, were tested on a model created for analyzing respiratory behavior and compared with a clinical diagnosis [51]. The model, which used LR, was created to train data collected in a study of over 3000 patients diagnosed with asthma and other respiratory diseases. The participants used a smartphone app to send a continuous "aaah" sound spoken for a 6-second duration, along with responses to a questionnaire about any possible symptoms.

In another study, accelerometer and voice recorder data were collected from participants with and without Parkinson disease, and a detection model was created using naive Bayes, KNN, and SVM methods [52]. The same model was used to detect Parkinson disease by using a new dataset obtained from patients newly suspected of having Parkinson disease. They were requested to pronounce the vowel "a" for 10 seconds, and a smartphone was kept at a certain distance (8 cm) from the patients to record the audio.

The dataset from the mPower study [31] was further used to test a voice condition analysis system for Parkinson disease, which was also verified using an experimental dataset (UEX) [53]. Six different machine learning classifiers (LR, RF, GBoost, passive aggressive, perceptron, and SVM) were applied to compare the performance with the 2 different speech databases. For creating the UEX dataset, 60 participants aged between 51 and 87 years were recruited. Of these 60 participants, 30 had Parkinson disease. A smartphone was used to record 3 different

samples of the participants pronouncing the "a" vowel continuously without being interrupted.

Voice data from the smartphones of patients with bipolar disorder were studied to determine if it is possible to differentiate people who have bipolar disorder and those who are either unaffected or have any relatives with bipolar disorder [54]. Data from 2 studies, namely the RADMIS trial [55] and the Bipolar Illness Onset (BIO) study [56], were used. The RADMIS trial was conducted with people diagnosed as having bipolar disorder who used a smartphone-based monitoring system installed on their phones, which collected voice data (only of those with Android smartphones) and other smartphone-related data, such as sleep duration and app usage. In the BIO study, participants included those who had bipolar disorder, those who had relatives with bipolar disorder, and those who were not diagnosed with bipolar disorder. The RF model developed from the data was verified using 5-fold participant-based cross-validation.

In some cases, the datasets for training the machine learning models were not obtained from previous studies. Various online sources were used for extracting both the crying and noncrying (eg, talking, breathing, hiccups, and yelling) sounds of infants [57]. For the validation of the algorithm, an independent dataset was created by using real-life recordings of 4 infants at home and 11 infants in a pediatric ward, where the recordings were created using smartphones. RF, LR, and naive Bayes were used for the classification and identification of crying and noncrying sounds.

Similarly, 41 YouTube videos and 5 cough sounds from the SoundSnap website were used to train a cough recognition model [58]. The study also included the development of a smartphone app, HealthMode Cough, that recorded continuous sounds, including sounds from streets, crowded markets, train stations, etc. The recordings were used to test the model, which used DCNN. Another model was created using CNN in a study aimed at analyzing the breathing sounds of participants with the smartphone app Breeze 2 [59]. The dataset for training the model was created using 3 separate datasets: a subset of the dataset from the study by Shih et al [60], which contained breathing sounds; the dataset ESC-50 [61], which contained 50 classes of environmental sounds; and a dataset from 2 participants, which contained a 2-minute breathing training session recorded using a smartphone. For the experiment, 30 participants without any respiratory diseases used the Breeze 2 app to perform 2 breathing sessions for 3 minutes: one with and one without headphones.

A dataset, compiled from multiple sources, was used to train a cough detection model for infants [57]. The cough sounds were obtained from 91 publicly available videos on YouTube consisting of coughing children aged between 0 and 16 years. Noncoughing sounds, such as talking, breathing, cat sounds, sirens, and dog sounds, were obtained via audio clips from YouTube, GitHub, and the British Broadcasting Corporation sound library. Furthermore, the audio data of 21 children, who were admitted with conditions, such as bronchitis, pneumonia, respiratory infection, and viral wheezing, were also collected via an Android smartphone. Using the data of 7 children out of

```
https://ai.jmir.org/2025/1/e59094
```

XSL•FO

the 21 and adding cough and noncough sounds from different sources, a model was created, and the data from the remaining 14 children were used as a validation dataset. The classification performance of the cough detection algorithm was compared using 2 ensemble DT classifiers: RF and GBoost.

Research on the Skin

Studies on the use of smartphone features to assess skin-related anomalies have mostly focused on the prediction or identification of skin cancer traits [57,58], and some studies have evaluated the detection of neonatal jaundice [62] and acne [63].

Studies Conducted Using Databases

To create a model for predicting skin cancer, 2 sets of databases were used in the study by Dascalu et al [64]: one with dermoscopic images (HAM10000 [65] and Dascalu and David [66]) and another with nondermoscopic images (Pacheco et al [67]). The images were obtained by taking pictures from a digital camera or a smartphone. Comparing the 2 datasets, sensitivity (percentage of correctly diagnosed malignancies) and specificity (percentage of negative diagnoses) were derived. The CNN-based model was found to improve specificity, though it was acknowledged that a significant amount of future work would be needed for improving sensitivity. It was also concluded that the dermoscopic images provided better accuracy compared to those from smartphones.

Studies Conducted Through Experiments

Acne is a common skin anomaly, which is experienced by about 10% of the world population. To predict and analyze such skin-related afflictions, many skin image analysis algorithms have been created [63]. To make the analysis and prediction accessible, it would be better to have such a system within a smartphone app. A CNN-based model was developed for acne detection, and an acne severity grading model was developed using the LightGBM algorithm. To test the models, an experiment was conducted, in which 1572 images of the faces of participants were taken from 3 different angles by using iOS or Android smartphones through a smartphone app called Skin Detective, and the dataset was divided in a ratio of 70:30 for training and testing. For ground truth, the images were labeled by 4 dermatologists.

Similarly, for predicting skin cancer, a melanoma detection model was created. A total of 514 patients from dermatology or plastic surgery clinics who had at least one skin lesion were selected, and pictures of their lesions were taken using 3 different cameras: 2 smartphone cameras and 1 digital camera [68]. For the analysis of the experiment dataset, an artificial intelligence algorithm, Deep Ensemble for Recognition of Malignancy [69], developed for determining the probability of skin cancer using dermoscopic images of skin lesions, was used.

Unlike for disease prediction using audio, data for skin-related anomalies can be obtained from other gadgets, such as smart wearables, through which information, such as heart rate, skin temperature, and breathing rate, can be obtained. A combination of data from smartphones (smartphone-based social interactions, activity patterns, and number of apps used) and smartwatches

(E4 Empatica; skin temperature) obtained via the in-house smartphone app MovisensXS was used to predict emotional changes and the severity of depression in people [70]. The study was conducted over a period of 8 weeks and included 41 people with depressive disorder. The participants had to complete daily smartphone-delivered surveys, a clinician-rated symptom assessment test, and a blood test to screen for potential medical contributors of depressed mood.

Studies Conducted With Both Databases and Experiments

Neonatal jaundice is a frequently occurring condition, which can also be diagnosed using smartphone images [62]. A study was conducted with 100 children, aged between 0 and 5 days, in which a picture or video was taken of their full face, with a calibration card, to capture their skin and eye sclera. Ground truth was established by noting their transcutaneous bilirubin (TCB) level, and the pictures were labeled "jaundiced" or "healthy" by a pediatrician. A CNN-based model was trained using the ImageNet dataset [71] and was used to test neonatal jaundice tendencies using the experiment dataset and transfer learning. Multilayer perceptron (MLP), SVM, DT, and RF were also used for diagnosis, where it was determined that transfer learning methods performed better for skin features, while machine learning models performed better for eye features.

Research on the Eye

Diabetic retinopathy was the most commonly studied disease [66,67,69] among the collected literature for eye-related predictions, along with other varying topics, such as eye tracking, vision monitoring, jaundice, and autism.

Studies Conducted Using Databases

An optimized hybrid machine learning classifier with the combination of neural network (NN) and DCNN with a single-stage object detection (SSD) algorithm was proposed to be used with the retinal images taken from a smartphone-enabled DIY camera [25] to predict diabetic retinopathy. Since there was a scarcity of image data captured using DIY smartphone-enabled devices, the model was validated with analysis of 2 other databases that contained fundus images: APTOS (2019 blindness dataset) and EyePACS, and the model performed better in comparison to the individual results of the NN, DCNN, and NN-DCNN methods.

CNN-based models usually tend to provide the best performance in image recognition tasks. With that in mind, the APTOS (2019 blindness dataset) and EyePACS datasets were used to build a CNN-based model for predicting diabetic retinopathy [72]. The algorithm was then externally validated using the Messidor-2 dataset [73], which contained about 1058 images from 4 French eye institutions. The algorithm was further tested on the EyeGo dataset, which contained 103 fundus images from 2 previously published studies obtained by using an EyeGo lens attachment and an iPhone.

Studies Conducted Through Experiments

Almost 51% of eye diseases in the United States are related to cataract [74]. It will be convenient to use images from smartphones for the early detection of cataract, and the results

```
https://ai.jmir.org/2025/1/e59094
```

will be provided instantly. By taking pictures with a smartphone camera, 100 samples were collected from participants (50% of the participants had cataract) [74]. SVM was applied on the dataset, and the accuracy was 96.6% for cataract detection.

In addition to images of the eye, videos of eye movement can be used for different kinds of diagnoses, such as for autism, since atypical eye gaze can be considered as an early symptom for autism spectrum disorder (ASD) [75]. The behaviors of 1564 toddlers were recorded using the front camera of an iPhone or an iPad when the toddlers, accompanied by their caregivers, viewed engaging movies for less than 60 seconds on the device. Using computer vision analysis on the data, it was found that children with ASD have less coordinated gaze patterns while viewing movement in movies or following conversation between 2 moving people.

In addition to the in-situ collection of data, smartphones can be used for remote collection of data. To determine the attention span of infants by tracking their gaze, an online webcam-linked eye tracker called OWLET was developed, and experiments were conducted with 127 infants remotely [76]. The infants were in the presence of their caregivers, who used either a smartphone or a computer to access the tracking task and provided their responses of the infant behavior using a questionnaire. For the experiment, a video (an 80-second Sesame Street video) was played, and the eye movements of the infants were recorded, tracked, and analyzed. No difference was found in the image data between the smartphone and computer, which was verified by a 2-sided independent samples t test and chi-square test. LR was used to examine the efficiency of the OWLET system.

Similarly, 417 adults with active or passive vision-related problems took part in an experiment using a smartphone app named Home Vision Monitor (HVM) to self-test their vision [77]. The app required them to submit an eye vision test twice per week, and their smartphone usage and app usage history were recorded by the researchers. RF and LR were used for statistical analysis.

The app EyeScreen was developed to support retinoblastoma diagnosis for the presence of leukocoria [78]. About 4000 eye images were obtained from about 1460 participants via the app, and an ImageNet model, ResNet, was used for image processing by dividing the dataset in an 80:20 ratio for training and testing. The app had the feature to process the image within it and provide the result.

Studies Conducted With Both Databases and Experiments

A common anomaly, neonatal jaundice, was investigated [62], for which a dataset of healthy and jaundiced individuals was created in an experiment conducted over 35 to 42 weeks. In the experiment, a full-face photo was clicked to capture the eye sclera. To obtain ground truth data, the TCB level was measured using a jaundice meter device, and the pictures were labeled "jaundiced" or "healthy" by a pediatrician.

Tracking eye movements has been a topic of interest for a wide variety of research ranging from autism [75] and tourism [28] to driving and gaming [79]. A multi-layer feed-forward

XSL•FO

convolutional neural network (ConvNet) model was created and trained on the GazeCapture dataset [80], which was created from the data of 1474 participants using an iPhone or iPad. To verify the model, an experiment was conducted using a custom-made Android app, in which eye gaze videos were captured using the front facing camera of the phone. The participants were asked to follow a stimulus on the mobile screen, which could be a dot or movement of a circular, rectangular, or zig-zag pattern.

Two sets of experiments were carried out, with one using a smartphone (iPhone 6) and another using smartphone-based retinal imaging systems, such as iExaminer, D-Eye, Peek Retina, and iNview [81], to create a model for the diagnosis of diabetic retinopathy. The CNN-based AlexNet architecture was used for transfer learning, which was first trained using 1234 images from the EyePACS dataset. Then, the architecture was tested with 138 retinal images from datasets, including those of the EyePACS, iExaminer, D-Eye, Peek Retina, and iNview systems.

Exclusion of Papers

A total of 11 papers were excluded from the final selection after reviewing the full text of all the papers using the selection criteria. Among them, 5 were excluded because the studies did not involve the use of smartphones for data collection. Similarly,

Figure 2. Diseases studied in the collected papers.

3 of the papers passed the initial screening test because of the presence of words, such as smartphone, eye, and audio, in their abstract. However, the studies were not relevant to the topic of our review. Furthermore, 2 of the studies were excluded because they only included the proposal of the method of disease prediction using machine learning and smartphones. Finally, a paper was excluded as it included a discussion about the topic but did not contain any database analysis or experiment. Among the papers that provided a proposal of a disease prediction system, it is worth mentioning that the paper by Bilal et al [82] was very detailed and well explained.

After the full-text screening of papers, there were 34, 5, and 10 relevant papers in the categories of voice, skin, and eye, respectively. These studies were further analyzed by focusing on the diseases dealt with in each study and the different health topics. The results can be seen in Figure 2. Parkinson disease was the most studied (n=12) disease among the collected studies, followed by COVID-19 (n=4), depression (n=4), cough (n=3), and diabetic retinopathy (n=3). It can be argued that cough and COVID-19 could be included under the same category and depression and emotion could be included under the same category. However, based on the terminologies and methods used in the papers, we have treated them separately.



Discussion

Overview

RenderX

The use of technology in the medical field has seen massive growth in recent years. A lot of improvements have been made in different areas, such as handling complex electronic medical records [83], and in identifying and predicting various diseases, such as lung anomaly detection using computed tomography scan images [84], emotion detection using different data from smartphones [35], and identification of the burden faced by

https://ai.jmir.org/2025/1/e59094

people who travel to separate locations for receiving health care services [85]. Smartphones provide a low-power, small-sized, and easy method for data collection and analysis, which differs from the usual bulky, expensive, and complex systems used for biomedical data collection and analysis.

Smartphones are equipped with numerous sensors and high-quality cameras, making it easy to collect different types of data. Moreover, due to the COVID-19 outbreak and the changes in the overall working environment that followed, there has been a strong focus on delivering health services remotely

[86]. The disease identification process can be made efficient by using smartphones to collect data and provide a diagnosis, as well as deliver results to patients.

With these factors in mind, we focused on research carried out using machine learning and data from smartphones to identify or predict diseases. We selected 3 areas to focus on and formulated the research questions. We conducted a review of the available papers collected using the screening method explained in the section Study Selection. Here onwards, we will discuss the results for our research questions.

Research Question 1: What Are the Databases Available for Eye, Skin, and Voice Analysis?

We found a total of 31 databases in the collected studies, including an unclear source, vaguely referred to as "online sources" [57]. In most of the cases, the databases were used to create a model for disease prediction. However, there were also instances where the databases were used to validate a model

developed using experimental data [81] or using other databases [72]. Since the number of collected voice-related studies was higher than that of skin- or eye-related studies, a similar difference in number can be observed for the list of databases, as shown in Tables 5-7. The numbers of databases for voice, skin, and eye were 22, 4, and 5, respectively. The voice-related databases were used to predict a variety of diseases or health statuses, such as Parkinson disease [26,29], emotion [35], bipolar disorder [55], and infant cry [57]. Owing to the COVID-19 pandemic, many databases were used for the detection of COVID-19 or cough-related anomalies [20,24,46,52]. Of 4 skin-related databases, 3 were aimed at the prediction of skin cancer [64] and the remaining database was related to neonatal jaundice [62]. The same database by Althnian et al [62] was also used for jaundice detection using retinal images. Diabetic retinopathy was the most common disease among eye databases [66,67,69]. Eye databases also consisted of data related to capturing eye movement [80] or gaze/concentration [74].

Table 5. Databases with voice data.

Database	Frequency
CoughVid	1
TASK	1
Brooklyn	1
Wallacedene	1
GoogleAudio dataset	1
Freesound	1
Librispeech	1
Coswara	2
ComparE	1
Sarcos	1
mPower	4
PC-GITA	1
Vishwanathan	1
ExtraSensory	1
Online sources	3
Shih et al [60]	1
UEX	1
RADMIS	1
Bipolar Illness Onset	1
YouTube	3
SoundSnap	1
BBC sound library	1

Table 6. Databases with skin data.

Database	Frequency
Ham10000	2
ImageNet	1
Dascalu and David [66]	1
Pacheco et al [67]	1

Table 7. Databases with eye data.

Database	Frequency
Messidor-2	1
GazeCapture	1
EyePacs	1
APTOS	1
EyeGO	1

Research Question 2: What Are the Machine Learning Models Used for Eye, Skin, and Voice Analysis?

Similar to databases, machine learning models were also used either in separation [67,87] or as a comparison along with multiple other models [6,24,31], and sometimes as ensemble classifiers [20,62]. As the same study usually consisted of multiple machine learning methods, the frequency of use of certain machine learning methods was considerably high. To investigate the best machine learning method for each kind of data, instead of using numbers, we calculated the frequency of use of a particular machine learning method for each of the 3 areas. We were then able to determine the rate of machine learning methods for each area, as shown in Tables 8-10. The most common machine learning method used for voice-related data was RF, while CNN was the most used for both eye- and skin-related data.

For further analysis, we expanded on the diseases and determined the frequency of the use of each machine learning method for each of the diseases or anomalies found in the collected papers. The results are shown in Figure 3. The figure shows all machine learning methods used across various studies for each of the diseases. Since many studies used multiple machine learning methods (especially for Parkinson disease), the frequency of use of some methods, such as RF, SVM, CNN, and LR, was high.

Table 8. Machine learning methods used with voice data.

Machine learning method	Rate of use, %							
AdaBoost ^a	1.1							
CNN ^b	10.9							
CatBoost	1.1							
DNN ^c	4.3							
Decision tree	3.3							
Deep learning	1.1							
GBoost ^d	5.4							
Gaussian mixture	1.1							
Hidden Markov	1.1							
KNN ^e	8.7							
LGBoost ^f	1.1							
LR ^g	10.9							
LSTM ^h	4.3							
Multilayer	1.1							
Naive Bayes	4.3							
Passive aggressive	1.1							
RF ⁱ	18.5							
Rank-based machine learning	1.1							
RestNet50	1.1							
SVM ^j	13.0							
XGBoost ^k	4.3							
m-SVR ¹	1.1							

^aAdaBoost: adaptive boosting.

^bCNN: convolutional neural network.

^cDNN: deep neural network.

^dGBoost: gradient boosting.

^eKNN: k-nearest neighbor.

^fLGBoost: light gradient boosting.

^gLR: logistic regression.

^hLSTM: long short-term memory.

ⁱRF: random forest.

^jSVM: support vector machine.

^kXGBoost: extreme gradient boosting.

¹m-SVR: multiple-output support vector regression.



 Table 9. Machine learning methods used with skin data.

Machine learning method	Rate of use, %	
CNN ^a	30.0	
Decision tree	10.0	
Deep learning	10.0	
Multilayer	10.0	
RF ^b	20.0	
SVM ^c	10.0	
XGBoost ^d	10.0	

^aCNN: convolutional neural network.

^bRF: random forest.

^cSVM: support vector machine.

^dXGBoost: extreme gradient boosting.

Table 10. Machine learning methods used with eye data.

Machine learning method	Rate of use, %					
CNN ^a	41.20					
Computer vision	5.88					
Decision tree	5.88					
LR ^b	11.76					
Multilayer	5.88					
Neural network	5.88					
RF ^c	11.80					
SVM ^d	11.80					

^aCNN: convolutional neural network.

^bLR: logistic regression.

^cRF: random forest.

^dSVM: support vector machine.



Figure 3. Use of machine learning (ML) methods based on the type of disease. AdaBoost: adaptive boosting; CNN: convolutional neural network; DNN: deep neural network; GBoost: gradient boosting; KNN: k-nearest neighbor; LGBoost: light gradient boosting; LR: logistic regression; LSTM: long short-term memory; m-SVR: multiple-output support vector regression; RF: random forest; SVM: support vector machine; XGBoost: extreme gradient boosting.

Acne -	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		- 8
Activity recognition -	0	0	0	0	0	1	0	0	1	1	1	0	0	0	0	1	0	0	1	0	0	0	0	0		
Autism -	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
Bipolar disorder -	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0		- 7
Bowel sound -	0	1	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0		
Breathing -	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
COVID-19 -	0	1	0	0	0	1	1	0	0	0	1	0	1	1	0	0	0	0	1	0	1	0	1	0		- 6
Cataract -	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0		
Cough -	0	2	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0		
Depression -	0	0	0	0	1	0	0	0	0	0	2	0	2	1	0	0	0	0	2	0	0	2	1	0		- 5
Diabetic retinopathy -	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0		
Emotion -	0	0	1	0	0	0	0	0	0	0	0	0	1	0	1	0	0	0	2	0	0	0	2	0		
ratigue -	0	0	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0	0	0	0	1	0	0		
Gaze capture -	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		- 4
Heart rate -	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0		
Infant cry -	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	1	0	0	2	0	0	0	0	0		
Infant gaze -	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0		- 3
Leukocoria -	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
Neonatal jaundice -	0	2	0	0	0	2	0	0	0	0	0	0	0	0	2	0	0	0	2	0	0	2	0	0		
Parkinson -	1	2	0	0	2	1	0	4	0	0	3	1	4	0	0	2	0	1	8	1	0	6	1	0		- 2
Perinatal depression -	0	1	0	0	1	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0		
Schizophrenia -	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1		
Skin cancer -	0	1	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		- 1
Sleep -	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
Snoring -	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0		
Vision -	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	0	0		- 0
	AdaBoost	CNN	CatBoost	Computer vision	DNN	Decision tree	Deep learning	GBoost	Gaussian mixture	Hidden Markov	KNN	LGBoost	athc	LSTM	Multilayer	Naive Bayes	Neural network	Passive aggressive	RF	Rank-based ML	RestNet50	SVM	XGBoost	m-SVR		•

Research Question 3: How Are the Data Collected From Smartphones?

To collect audio-related data, the built-in smartphone microphone was used most of the time, both at home [26,37] and in the experimental set up [52]. In some cases, external microphones were also used [36]. Similarly, in many cases, audio was also collected via custom-made smartphone apps [21,34,38], and in some cases, it was collected via a web interface that could be accessed using smartphones [20,25].

For the collection of skin data, pictures and videos were mainly taken with a smartphone [58,59]. In some cases, smartphone apps were also created for data collection [60,63]. Frequently, pictures from smartphones were not considered adequate for taking retinal images, and an external lens or retinal imaging system was used alongside the smartphone to collect eye data [20,66,67]. However, experiments have also shown that smartphone images are equally effective to analyze eye-related anomalies [25,71]. For collecting gaze-related data, videos taken from the front camera of smartphones have been used effectively [75,81].

from voice and skin data were also collected. Combinations of data from smartwatches, such as heart rate, skin temperature, and breathing rate, and data from smartphones, such as smartphone-based social interactions, activity patterns, and the usage of apps, were used for detecting emotional changes and the severity of depression [70]. Similarly, for detecting emotional status, voice and other data, such as location, accelerometer data, gyroscope data, and phone usage, were used [31,36]. For the detection of Parkinson disease, data apart from voice data, such as gait and balance [26], tapping of buttons on a smartphone screen [26,33], and accelerometer data [52], were collected. Similarly, in a study for vision monitoring, data apart from images of the eyes, such as phone and app usage, were collected using a smartphone app [77]. Furthermore, in many studies, surveys and questionnaires were also regularly received from participants during data collection, especially via smartphone apps [26,34,36,63].

Data of both the skin and voice have been used for the detection

of emotion and depression. However, in such studies, data apart

It is worth noting that there are many sensors available in smartphones, such as accelerometers and gyroscopes, which

can be helpful in determining the speed of touch, posture, walking speed, location, etc. Therefore, even if the aim is to analyze a particular health area, the same combination of data, collected from multiple data sources, can be used to identify different diseases. For example, in the study of Parkinson disease, data, such as voice, accelerometer data, location, application usage, and other phone data, were usually collected. The same data can also be analyzed to detect emotional changes or depression. Similarly, when collecting data on skin abnormalities, it is possible to obtain facial data that can also be used for eye-related analysis.

Moreover, it has been observed that with remote health monitoring systems, especially with the use of smartphones, people often have concerns over the access and use of their data, such as location, application usage, screen time, and browsing history [78,79]. These concerns of smartphone apps are higher as they contain sensitive behavioral data. To tackle these issues, it is necessary to build trust with the users regarding the app and its data collection methods. It was found that state-funded research institutes had higher levels of trust with people compared to private institutions [88]. This shows that to conduct research using smartphones and gather user data, it is necessary to involve trusted institutions for governing the study, as well as have transparency over data collection, distribution, and use of the results.

Limitations

There are some limitations. First, for screening the collected papers based on their titles and abstracts, we used a keyword screening method [23]. Although great care was taken in the selection of keywords for this screening, it must be acknowledged that some papers may have been overlooked if they did not contain the specified keywords. We firmly believe that such a limitation can occur, but the number of studies will be very few. Second, we focused only on studies that used smartphones. This could lead to the exclusion of recent studies that did not consider the use of smartphones to collect health-related data.

Moreover, we only selected studies that analyzed eye-, skin-, and voice-related diseases. Because of the niche approach of this scoping review, we did not consider a lot of other health areas where smartphones might have been used to gather data for machine learning analysis. Furthermore, many new machine learning models and other algorithms are being developed, and existing algorithms are being improved [89]. These methods have not been used but could potentially be used for health diagnosis, and thus, they have been overlooked in this review.

Overall Summary

The field of the use of machine learning on smartphone-obtained data for health care purposes is ever evolving. Through this study, we aimed to provide information about studies that have conducted experiments related to eye-, skin-, or voice-related diseases, where data were obtained strictly via smartphones. Similarly, we have provided details of publicly available databases that have been used in studies to apply machine learning methods for developing models to predict eye-, skin-, or voice-related diseases. Researchers working in similar fields can use the experiment details or the databases presented in this study to design their research. Furthermore, the machine learning model to use for a study needs to be determined with much consideration. We have presented machine learning models applied based on the study area as well as the types of diseases. Therefore, the information provided in the paper can help reduce the time and effort for researchers in designing experiments and selecting the databases or machine learning models to use in their studies. Our title and abstract screening method is also easy to understand and replicate, and could be used by researchers aiming to perform scoping reviews or systematic literature reviews.

Conclusion

There has been a growth in the number of studies based on the application of machine learning methods to data obtained from smartphones for the prediction of diseases. However, there are few literature reviews that provide information about the databases used, experiments carried out, and machine learning methods applied. We formulated a scoping review to identify the studies that have been conducted, specifically related to the 3 areas of skin, eye, and voice, and determined the studies that conducted experiments using smartphones to gather skin-, eye-, and voice-related data; the publicly available databases that include skin, eye, or voice data; and the machine learning methods that are commonly implemented in such studies. Furthermore, with this research, we intended to test the effectiveness of the keyword screening method that we developed. We first searched for relevant studies and screened them by applying our keyword screening method to their titles and abstracts. We analyzed the full text according to the inclusion and exclusion criteria and collected a total of 60 studies.

After assessing the full text of all identified studies, we discarded 11 studies, and among the remaining 49 studies, we found 24 different machine learning methods and 31 different databases used. The details from these collected studies provide insights into how the experimental studies were conducted, which databases were used, and which machine learning methods provided better results. The relevance and quality of the information acquired proved that our keyword screening method was effective in screening papers relevant to the topic and thus could be adopted by researchers for conducting scoping reviews. The use of our results can help reduce the time and effort required by researchers working in the field of artificial intelligence for health care to gather such information in detail. Moreover, the results presented can be used to select databases for future studies, replicate the experimental design, or select machine learning models suitable for the topic of interest.



Acknowledgments

This study was supported by Japan Science and Technology Agency (JST) COI-NEXT (grant number JPMJPF2018), AI Nutrition and Food Functionalities Task Force, ILSI Japan.

Conflicts of Interest

None declared.

Multimedia Appendix 1

PRISMA (Preferred Reporting Items for Systematic Reviews and Meta-Analyses) checklist for scoping reviews. [PDF File (Adobe PDF File), 60 KB - <u>ai_v4i1e59094_app1.pdf</u>]

References

- 1. Shinde R, Sawant S, Maskar T, Randhe K. Disease prediction using machine learning. International Research Journal of Engineering and Technology (IRJET) 2021;8(4):2721-2724 [FREE Full text]
- Kohli PS, Arora S. Application of Machine Learning in Disease Prediction. 2018 Presented at: 4th International Conference on Computing Communication and Automation (ICCCA); December 14-15, 2018; Greater Noida, India. [doi: 10.1109/CCAA.2018.8777449]
- Dinh A, Miertschin S, Young A, Mohanty SD. A data-driven approach to predicting diabetes and cardiovascular disease with machine learning. BMC Med Inform Decis Mak 2019 Nov 06;19(1):211 [FREE Full text] [doi: 10.1186/s12911-019-0918-5] [Medline: 31694707]
- 4. Alanazi H, Abdullah AH, Qureshi KN, Ismail AS. Accurate and dynamic predictive model for better prediction in medicine and healthcare. Ir J Med Sci 2018 May;187(2):501-513 [FREE Full text] [doi: 10.1007/s11845-017-1655-3] [Medline: 28756541]
- 5. Panchal PJ, Mhaskar SA, Ziman TS. Disease prediction using machine learning. Iconic Research And Engineering Journals 2020;3(10):302-303 [FREE Full text]
- Shi B, Chen J, Chen H, Lin W, Yang J, Chen Y, et al. Prediction of recurrent spontaneous abortion using evolutionary machine learning with joint self-adaptive sime mould algorithm. Comput Biol Med 2022 Sep;148:105885 [FREE Full text] [doi: 10.1016/j.compbiomed.2022.105885] [Medline: 35930957]
- 7. Fei X, Wang J, Ying S, Hu Z, Shi J. Projective parameter transfer based sparse multiple empirical kernel learning Machine for diagnosis of brain disease. Neurocomputing 2020 Nov;413:271-283 [FREE Full text] [doi: 10.1016/j.neucom.2020.07.008]
- 8. Wang S, Xiang J. A minimum entropy deconvolution-enhanced convolutional neural networks for fault diagnosis of axial piston pumps. Soft Comput 2019 May 23;24(4):2983-2997 [FREE Full text] [doi: 10.1007/s00500-019-04076-2]
- 9. Kokubo Y, Kamide K, Okamura T, Watanabe M, Higashiyama A, Kawanishi K, et al. Impact of high-normal blood pressure on the risk of cardiovascular disease in a Japanese urban cohort. Hypertension 2008 Oct;52(4):652-659 [FREE Full text] [doi: 10.1161/hypertensionaha.108.118273]
- Mosa A, Yoo I, Sheets L. A systematic review of healthcare applications for smartphones. BMC Med Inform Decis Mak 2012 Jul 10;12:67 [FREE Full text] [doi: 10.1186/1472-6947-12-67] [Medline: 22781312]
- Tremmel M, Gerdtham UG, Nilsson PM, Saha S. Economic burden of obesity: a systematic literature review. Int J Environ Res Public Health 2017 Apr 19;14(4):A [FREE Full text] [doi: 10.3390/ijerph14040435] [Medline: 28422077]
- Moher D, Liberati A, Tetzlaff J, Altman DG, PRISMA Group. Preferred reporting items for systematic reviews and meta-analyses: the PRISMA statement. Ann Intern Med 2009 Aug 18;151(4):264-9, W64 [FREE Full text] [doi: 10.7326/0003-4819-151-4-200908180-00135] [Medline: 19622511]
- Wang N, Chen J, Chen W, Shi Z, Yang H, Liu P, et al. The effectiveness of case management for cancer patients: an umbrella review. BMC Health Serv Res 2022 Oct 14;22(1):1247 [FREE Full text] [doi: 10.1186/s12913-022-08610-1] [Medline: 36242021]
- Välimäki MA, Lantta T, Hipp K, Varpula J, Liu G, Tang Y, et al. Measured and perceived impacts of evidence-based leadership in nursing: a mixed-methods systematic review protocol. BMJ Open 2021 Oct 22;11(10):e055356 [FREE Full text] [doi: 10.1136/bmjopen-2021-055356] [Medline: 34686559]
- Uddin S, Khan A, Hossain ME, Moni MA. Comparing different supervised machine learning algorithms for disease prediction. BMC Med Inform Decis Mak 2019 Dec 21;19(1):281 [FREE Full text] [doi: 10.1186/s12911-019-1004-8] [Medline: 31864346]
- Majumder S, Deen MJ. Smartphone sensors for health monitoring and diagnosis. Sensors (Basel) 2019 May 09;19(9):2164 [FREE Full text] [doi: 10.3390/s19092164] [Medline: 31075985]
- Mak S, Thomas A. Steps for conducting a scoping review. J Grad Med Educ 2022 Oct;14(5):565-567 [FREE Full text] [doi: 10.4300/JGME-D-22-00621.1] [Medline: 36274762]
- 18. PubMed. URL: <u>https://pubmed.ncbi.nlm.nih.gov/</u> [accessed 2025-03-15]
- 19. IEEE Xplore. URL: <u>https://ieeexplore.ieee.org</u> [accessed 2025-03-15]

- 20. Mendeley. URL: https://www.mendeley.com/ [accessed 2025-03-15]
- 21. Zotero. URL: <u>https://www.zotero.org/</u> [accessed 2025-03-15]
- 22. Python. URL: <u>https://www.python.org</u> [accessed 2025-03-15]
- 23. Dawadi R, Araki M. A Semi-automatic Screening Mechanism for Literature Review. SSRN Journal. 2023. URL: <u>https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4521171</u> [accessed 2025-03-15]
- 24. Shabgou M, Daryani SM. Towards the sensory marketing: Stimulating the five senses (sight, hearing, smell, touch and taste) and its impact on consumer behavior. Indian Journal of Fundamental and Applied Life Sciences 2014;4(1):573-581 [FREE Full text]
- 25. Gupta S, Thakur S, Gupta A. Optimized hybrid machine learning approach for smartphone based diabetic retinopathy detection. Multimed Tools Appl 2022;81(10):14475-14501 [FREE Full text] [doi: 10.1007/s11042-022-12103-y] [Medline: 35233182]
- 26. Zhan A, Mohan S, Tarolli C, Schneider RB, Adams JL, Sharma S, et al. Using smartphones and machine learning to quantify Parkinson disease severity: the mobile Parkinson disease score. JAMA Neurol 2018 Jul 01;75(7):876-880 [FREE Full text] [doi: 10.1001/jamaneurol.2018.0809] [Medline: 29582075]
- 27. Kim A, Jang EH, Lee SH, Choi KY, Park JG, Shin HC. Automatic depression detection using smartphone-based text-dependent speech signals: deep convolutional neural network approach. J Med Internet Res 2023 Jan 25;25:e34474 [FREE Full text] [doi: 10.2196/34474] [Medline: 36696160]
- Matsuda Y, Fedotov D, Takahashi Y, Arakawa Y, Yasumoto K, Minker W. EmoTour: Estimating emotion and satisfaction of users based on behavioral cues and audiovisual data. Sensors (Basel) 2018 Nov 15;18(11):e [FREE Full text] [doi: 10.3390/s18113978] [Medline: 30445798]
- 29. Pahar M, Klopper M, Warren R, Niesler T. COVID-19 detection in cough, breath and speech using deep transfer learning and bottleneck features. Comput Biol Med 2022 Feb;141:105153 [FREE Full text] [doi: 10.1016/j.compbiomed.2021.105153] [Medline: 34954610]
- 30. Alkhodari M, Khandoker AH. Detection of COVID-19 in smartphone-based breathing recordings: A pre-screening deep learning tool. PLoS One 2022;17(1):e0262448 [FREE Full text] [doi: 10.1371/journal.pone.0262448] [Medline: 35025945]
- Prince J, Andreotti F, De Vos M. Multi-source ensemble learning for the remote prediction of Parkinson's disease in the presence of source-wise missing data. IEEE Trans Biomed Eng 2019 May;66(5):1402-1411 [FREE Full text] [doi: 10.1109/tbme.2018.2873252]
- Tougui I, Jilbab A, Mhamdi JE. Analysis of smartphone recordings in time, frequency, and cepstral domains to classify Parkinson's disease. Healthc Inform Res 2020 Oct;26(4):274-283 [FREE Full text] [doi: 10.4258/hir.2020.26.4.274] [Medline: <u>33190461</u>]
- Tougui I, Jilbab A, Mhamdi JE. Machine learning smart system for Parkinson disease classification using the voice as a biomarker. Healthc Inform Res 2022 Jul;28(3):210-221 [FREE Full text] [doi: 10.4258/hir.2022.28.3.210] [Medline: 35982595]
- Pah N, Motin MA, Kumar DK. Phonemes based detection of Parkinson's disease for telehealth applications. Sci Rep 2022 Jun 11;12(1):9687 [FREE Full text] [doi: 10.1038/s41598-022-13865-z] [Medline: 35690657]
- Sultana M, Al-Jefri M, Lee J. Using machine learning and smartphone and smartwatch data to detect emotional states and transitions: exploratory study. JMIR Mhealth Uhealth 2020 Sep 29;8(9):e17818 [FREE Full text] [doi: 10.2196/17818] [Medline: 32990638]
- 36. Vaiciukynas E, Verikas A, Gelzinis A, Bacauskiene M. Detecting Parkinson's disease from sustained phonation and speech signals. PLoS One 2017;12(10):e0185613 [FREE Full text] [doi: 10.1371/journal.pone.0185613] [Medline: 28982171]
- Lo C, Arora S, Baig F, Lawton MA, El Mouden C, Barber TR, et al. Predicting motor, cognitive and functional impairment in Parkinson's. Ann Clin Transl Neurol 2019 Aug;6(8):1498-1509 [FREE Full text] [doi: 10.1002/acn3.50853] [Medline: 31402628]
- Lim W, Chiu SI, Wu MC, Tsai SF, Wang PH, Lin KP, et al. An integrated biometric voice and facial features for early detection of Parkinson's disease. NPJ Parkinsons Dis 2022 Oct 29;8(1):145 [FREE Full text] [doi: 10.1038/s41531-022-00414-8] [Medline: 36309501]
- Lin D, Nazreen T, Rutowski T, Lu Y, Harati A, Shriberg E, et al. Feasibility of a machine learning-based smartphone application in detecting depression and anxiety in a generally senior population. Front Psychol 2022;13:811517 [FREE Full text] [doi: 10.3389/fpsyg.2022.811517] [Medline: 35478769]
- Belouali A, Gupta S, Sourirajan V, Yu J, Allen N, Alaoui A, et al. Acoustic and language analysis of speech for suicidal ideation among US veterans. BioData Min 2021 Feb 02;14(1):11 [FREE Full text] [doi: 10.1186/s13040-021-00245-y] [Medline: 33531048]
- 41. Mundnich K, Booth BM, L'Hommedieu M, Feng T, Girault B, L'Hommedieu J, et al. TILES-2018, a longitudinal physiologic and behavioral data set of hospital workers. Sci Data 2020 Oct 16;7(1):354 [FREE Full text] [doi: 10.1038/s41597-020-00655-3] [Medline: 33067468]
- 42. Tonn P, Seule L, Degani Y, Herzinger S, Klein A, Schulze N. Digital content-free speech analysis tool to measure affective distress in mental health: evaluation study. JMIR Form Res 2022 Aug 30;6(8):e37061 [FREE Full text] [doi: 10.2196/37061] [Medline: 36040767]

- Parga J, Lewin S, Lewis J, Montoya-Williams D, Alwan A, Shaul B, et al. Defining and distinguishing infant behavioral states using acoustic cry analysis: is colic painful? Pediatr Res 2020 Feb;87(3):576-580 [FREE Full text] [doi: 10.1038/s41390-019-0592-4] [Medline: 31585457]
- 44. Tseng V, Sano A, Ben-Zeev D, Brian R, Campbell AT, Hauser M, et al. Using behavioral rhythms and multi-task learning to predict fine-grained symptoms of schizophrenia. Sci Rep 2020 Sep 15;10(1):15100 [FREE Full text] [doi: 10.1038/s41598-020-71689-1] [Medline: 32934246]
- 45. Elbéji A, Zhang L, Higa E, Fischer A, Despotovic V, Nazarov PV, et al. Vocal biomarker predicts fatigue in people with COVID-19: results from the prospective Predi-COVID cohort study. BMJ Open 2022 Nov 22;12(11):e062463 [FREE Full text] [doi: 10.1136/bmjopen-2022-062463] [Medline: 36414294]
- 46. Peng R, Yan WR, Zhang NL, Lin WH, Zhou XL, Zhang YT. Cuffless and continuous blood pressure estimation from the heart sound signals. Sensors (Basel) 2015 Sep 17;15(9):23653-23666 [FREE Full text] [doi: 10.3390/s150923653] [Medline: 26393591]
- 47. Kutsumi Y, Kanegawa N, Zeida M, Matsubara H, Murayama N. Automated bowel sound and motility analysis with CNN using a smartphone. Sensors (Basel) 2022 Dec 30;23(1):407 [FREE Full text] [doi: 10.3390/s23010407] [Medline: 36617005]
- 48. Yang F, Wu Q, Hu X, Ye J, Yang Y, Rao H, et al. Internet-of-things-enabled data fusion method for sleep healthcare applications. IEEE Internet Things J 2021 Nov 1;8(21):15892-15905 [FREE Full text] [doi: 10.1109/jiot.2021.3067905]
- Goswami U, Black A, Krohn B, Meyers W, Iber C. Smartphone-based delivery of oropharyngeal exercises for treatment of snoring: a randomized controlled trial. Sleep Breath 2019 Mar;23(1):243-250 [FREE Full text] [doi: 10.1007/s11325-018-1690-y] [Medline: 30032464]
- Wu J, Feng Y, Sun P. Sensor fusion for recognition of activities of daily living. Sensors (Basel) 2018 Nov 19;18(11):4029 [FREE Full text] [doi: 10.3390/s18114029] [Medline: 30463199]
- 51. Kaur S, Larsen E, Harper J, Purandare B, Uluer A, Hasdianda MA, et al. Development and validation of a respiratory-responsive vocal biomarker-based tool for generalizable detection of respiratory impairment: independent case-control studies in multiple respiratory conditions including asthma, chronic obstructive pulmonary disease, and COVID-19. J Med Internet Res 2023 Apr 14;25:e44410 [FREE Full text] [doi: 10.2196/44410] [Medline: 36881540]
- Sajal M, Ehsan MT, Vaidyanathan R, Wang S, Aziz T, Mamun KAA. Telemonitoring Parkinson's disease using machine learning by combining tremor and voice analysis. Brain Inform 2020 Oct 22;7(1):12 [FREE Full text] [doi: 10.1186/s40708-020-00113-1] [Medline: <u>33090328</u>]
- 53. Carrón J, Campos-Roca Y, Madruga M, Pérez CJ. A mobile-assisted voice condition analysis system for Parkinson's disease: assessment of usability conditions. Biomed Eng Online 2021 Nov 21;20(1):114 [FREE Full text] [doi: 10.1186/s12938-021-00951-y] [Medline: 34802448]
- 54. Faurholt-Jepsen M, Rohani DA, Busk J, Vinberg M, Bardram JE, Kessing LV. Voice analyses using smartphone-based data in patients with bipolar disorder, unaffected relatives and healthy control individuals, and during different affective states. Int J Bipolar Disord 2021 Dec 01;9(1):38 [FREE Full text] [doi: 10.1186/s40345-021-00243-3] [Medline: 34850296]
- 55. Faurholt-Jepsen M, Lindbjerg Tønning M, Fros M, Martiny K, Tuxen N, Rosenberg N, et al. Reducing the rate of psychiatric re-admissions in bipolar disorder using smartphones-The RADMIS trial. Acta Psychiatr Scand 2021 May;143(5):453-465 [FREE Full text] [doi: 10.1111/acps.13274] [Medline: 33354769]
- 56. Kessing L, Munkholm K, Faurholt-Jepsen M, Miskowiak KW, Nielsen LB, Frikke-Schmidt R, et al. The Bipolar Illness Onset study: research protocol for the BIO cohort study. BMJ Open 2017 Jun 23;7(6):e015462 [FREE Full text] [doi: 10.1136/bmjopen-2016-015462] [Medline: 28645967]
- Kruizinga M, Zhuparris A, Dessing E, Krol FJ, Sprij AJ, Doll RJ, et al. Development and technical validation of a smartphone-based pediatric cough detection algorithm. Pediatr Pulmonol 2022 Mar;57(3):761-767 [FREE Full text] [doi: 10.1002/ppul.25801] [Medline: 34964557]
- Kvapilova L, Boza V, Dubec P, Majernik M, Bogar J, Jamison J, et al. Continuous sound collection using smartphones and machine learning to measure cough. Digit Biomark 2019;3(3):166-175 [FREE Full text] [doi: 10.1159/000504666] [Medline: 32095775]
- Lukic Y, Teepe GW, Fleisch E, Kowatsch T. Breathing as an input modality in a gameful breathing training app (Breeze 2): development and evaluation study. JMIR Serious Games 2022 Aug 16;10(3):e39186 [FREE Full text] [doi: 10.2196/39186] [Medline: 35972793]
- 60. Shih CH, Tomita N, Lukic YX, Reguera Á, Fleisch E, Kowatsch T. Breeze: smartphone-based acoustic real-time detection of breathing phases for a gamified biofeedback breathing training. Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies 2020;3(4):1-30. [doi: 10.1145/3369835]
- 61. Piczak KJ. ESC: Dataset for Environmental Sound Classification. In: MM '15: Proceedings of the 23rd ACM international conference on Multimedia. 2015 Presented at: 23rd ACM international conference on Multimedia; October 26-30, 2015; Brisbane, Australia.
- Althnian A, Almanea N, Aloboud N. Neonatal jaundice diagnosis using a smartphone camera based on eye, skin, and fused features with transfer learning. Sensors (Basel) 2021 Oct 23;21(21):7038 [FREE Full text] [doi: 10.3390/s21217038] [Medline: 34770345]

- 63. Huynh Q, Nguyen PH, Le HX, Ngo LT, Trinh NT, Tran MTT, et al. Automatic acne object detection and acne severity grading using smartphone images and artificial intelligence. Diagnostics (Basel) 2022 Aug 03;12(8):1879 [FREE Full text] [doi: 10.3390/diagnostics12081879] [Medline: 36010229]
- 64. Dascalu A, Walker BN, Oron Y, David EO. Non-melanoma skin cancer diagnosis: a comparison between dermoscopic and smartphone images by unified visual and sonification deep learning algorithms. J Cancer Res Clin Oncol 2022 Sep;148(9):2497-2505 [FREE Full text] [doi: 10.1007/s00432-021-03809-x] [Medline: 34546412]
- 65. Tschandl P, Rosendahl C, Akay BN, Argenziano G, Blum A, Braun RP, et al. Expert-level diagnosis of nonpigmented skin cancer by combined convolutional neural networks. JAMA Dermatol 2019 Jan 01;155(1):58-65 [FREE Full text] [doi: 10.1001/jamadermatol.2018.4378] [Medline: 30484822]
- 66. Dascalu A, David EO. Skin cancer detection by deep learning and sound analysis algorithms: A prospective clinical study of an elementary dermoscope. EBioMedicine 2019 May;43:107-113 [FREE Full text] [doi: 10.1016/j.ebiom.2019.04.055] [Medline: 31101596]
- Pacheco A, Lima GR, Salomão AS, Krohling B, Biral IP, de Angelo GG, et al. PAD-UFES-20: A skin lesion dataset composed of patient data and clinical images collected from smartphones. Data Brief 2020 Oct;32:106221 [FREE Full text] [doi: 10.1016/j.dib.2020.106221] [Medline: 32939378]
- 68. Phillips M, Marsden H, Jaffe W, Matin RN, Wali GN, Greenhalgh J, et al. Assessment of accuracy of an artificial intelligence algorithm to detect melanoma in images of skin lesions. JAMA Netw Open 2019 Oct 02;2(10):e1913436 [FREE Full text] [doi: 10.1001/jamanetworkopen.2019.13436] [Medline: 31617929]
- Phillips M, Greenhalgh J, Marsden H, Palamaras I. Detection of malignant melanoma using artificial intelligence: an observational study of diagnostic accuracy. Dermatol Pract Concept 2020;10(1):e2020011 [FREE Full text] [doi: 10.5826/dpc.1001a11] [Medline: 31921498]
- 70. Pedrelli P, Fedor S, Ghandeharioun A, Howe E, Ionescu DF, Bhathena D, et al. Monitoring changes in depression severity using wearable and mobile sensors. Front Psychiatry 2020;11:584711 [FREE Full text] [doi: 10.3389/fpsyt.2020.584711] [Medline: 33391050]
- Krizhevsky A, Sutskever I, Hinton G. ImageNet classification with deep convolutional neural networks. Commun ACM 2017 May 24;60(6):84-90 [FREE Full text] [doi: 10.1145/3065386]
- Ludwig C, Perera C, Myung D, Greven MA, Smith SJ, Chang RT, et al. Automatic identification of referral-warranted diabetic retinopathy using deep learning on mobile phone images. Transl Vis Sci Technol 2020 Dec;9(2):60 [FREE Full text] [doi: 10.1167/tvst.9.2.60] [Medline: 33294301]
- 73. Decencière E, Zhang X, Cazuguel G, Lay B, Cochener B, Trone C, et al. Feedback on a publicly distributed image database: The Messidor database. Image Anal Stereol 2014 Aug 26;33(3):231 [FREE Full text] [doi: 10.5566/ias.1155]
- 74. Askarian B, Ho P, Chong JW. Detecting cataract using smartphones. IEEE J Transl Eng Health Med 2021;9:1-10 [FREE Full text] [doi: 10.1109/jtehm.2021.3074597]
- 75. Chang Z, Di Martino JM, Aiello R, Baker J, Carpenter K, Compton S, et al. Computational methods to measure patterns of gaze in toddlers with autism spectrum disorder. JAMA Pediatr 2021 Aug 01;175(8):827-836 [FREE Full text] [doi: 10.1001/jamapediatrics.2021.0530] [Medline: 33900383]
- 76. Werchan D, Thomason ME, Brito NH. OWLET: An automated, open-source method for infant gaze tracking using smartphone and webcam recordings. Behav Res Methods 2023 Sep;55(6):3149-3163 [FREE Full text] [doi: 10.3758/s13428-022-01962-w] [Medline: 36070130]
- 77. Korot E, Pontikos N, Drawnel FM, Jaber A, Fu DJ, Zhang G, et al. Enablers and barriers to deployment of smartphone-based home vision monitoring in clinical practice settings. JAMA Ophthalmol 2022 Feb 01;140(2):153-160 [FREE Full text] [doi: 10.1001/jamaophthalmol.2021.5269] [Medline: 34913967]
- 78. Bernard A, Xia SZ, Saleh S, Ndukwe T, Meyer J, Soloway E, et al. EyeScreen: Development and potential of a novel machine learning application to detect leukocoria. Ophthalmol Sci 2022 Sep;2(3):100158 [FREE Full text] [doi: 10.1016/j.xops.2022.100158] [Medline: 36245758]
- 79. Valliappan N, Dai N, Steinberg E, He J, Rogers K, Ramachandran V, et al. Accelerating eye movement research via accurate and affordable smartphone eye tracking. Nat Commun 2020 Sep 11;11(1):4553 [FREE Full text] [doi: 10.1038/s41467-020-18360-5] [Medline: 32917902]
- Krafka K, Khosla A, Kellnhofer P, Kannan H, Bhandarkar S, Matusik W, et al. Eye Tracking for Everyone. 2016 Presented at: IEEE Conference on Computer Vision and Pattern Recognition (CVPR); June 27-30, 2016; Las Vegas, NV, USA. [doi: 10.1109/CVPR.2016.239]
- Karakaya M, Hacisoftaoglu RE. Comparison of smartphone-based retinal imaging systems for diabetic retinopathy detection using deep learning. BMC Bioinformatics 2020 Jul 06;21(Suppl 4):259 [FREE Full text] [doi: 10.1186/s12859-020-03587-2] [Medline: 32631221]
- Bilal A, Fransson E, Bränn E, Eriksson A, Zhong M, Gidén K, et al. Predicting perinatal health outcomes using smartphone-based digital phenotyping and machine learning in a prospective Swedish cohort (Mom2B): study protocol. BMJ Open 2022 Apr 27;12(4):e059033 [FREE Full text] [doi: 10.1136/bmjopen-2021-059033] [Medline: 35477874]
- Li Q, You T, Chen J, Zhang Y, Du C. LI-EMRSQL: Linking information enhanced Text2SQL parsing on complex electronic medical records. IEEE Trans Rel 2024 Jun;73(2):1280-1290 [FREE Full text] [doi: 10.1109/TR.2023.3336330]

```
https://ai.jmir.org/2025/1/e59094
```

- 84. Chen Y, Feng L, Zheng C, Zhou T, Liu L, Liu P, et al. LDANet: Automatic lung parenchyma segmentation from CT images. Comput Biol Med 2023 Mar;155:106659 [FREE Full text] [doi: 10.1016/j.compbiomed.2023.106659] [Medline: 36791550]
- 85. Wang Q, Jiang Q, Yang Y, Pan J. The burden of travel for care and its influencing factors in China: An inpatient-based study of travel time. Journal of Transport & Health 2022 Jun;25:101353 [FREE Full text] [doi: 10.1016/j.jth.2022.101353]
- Harvill J, Wani Y, Alam M, Ahuja N, Hasegawa-Johnsor M, Chestek D, et al. Estimation of Respiratory Rate from Breathing Audio. 2022 Presented at: 44th Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC); July 11-15, 2022; Glasgow, Scotland, United Kingdom. [doi: 10.1109/EMBC48229.2022.9871897]
- Tougui I, Jilbab A, Mhamdi JE. Impact of the choice of cross-validation techniques on the results of machine learning-based diagnostic applications. Healthc Inform Res 2021 Jul;27(3):189-199 [FREE Full text] [doi: 10.4258/hir.2021.27.3.189] [Medline: 34384201]
- Buhr L, Schicktanz S, Nordmeyer E. Attitudes toward mobile apps for pandemic research among smartphone users in Germany: national survey. JMIR Mhealth Uhealth 2022 Jan 24;10(1):e31857 [FREE Full text] [doi: 10.2196/31857] [Medline: 35072646]
- 89. Chen M, Yang L, Zeng G, Lu K, Huang Y. IFA-EO: An improved firefly algorithm hybridized with extremal optimization for continuous unconstrained optimization problems. Soft Comput 2022 Nov 09;27(6):2943-2964 [FREE Full text] [doi: 10.1007/s00500-022-07607-6]

Abbreviations

ASD: autism spectrum disorder **BIO:** Bipolar Illness Onset **CNN:** convolutional neural network DCNN: deep convolutional neural network **DNN:** deep neural network DT: decision tree GBC: gradient boosting classifier **GBoost:** gradient boosting KNN: k-nearest neighbor LR: logistic regression LSTM: long short-term memory NN: neural network RF: random forest SVM: support vector machine TCB: transcutaneous bilirubin **XGBoost:** extreme gradient boosting

Edited by JL Raisaro; submitted 07.04.24; peer-reviewed by T Yagdi, R Qureshi, H Chen; comments to author 14.09.24; revised version received 06.10.24; accepted 23.02.25; published 25.03.25.

<u>Please cite as:</u>

Dawadi R, Inoue M, Tay JT, Martin-Morales A, Vu T, Araki M Disease Prediction Using Machine Learning on Smartphone-Based Eye, Skin, and Voice Data: Scoping Review JMIR AI 2025;4:e59094 URL: https://ai.jmir.org/2025/1/e59094 doi:10.2196/59094 PMID:

©Research Dawadi, Mai Inoue, Jie Ting Tay, Agustin Martin-Morales, Thien Vu, Michihiro Araki. Originally published in JMIR AI (https://ai.jmir.org), 25.03.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Original Paper

Prompt Engineering an Informational Chatbot for Education on Mental Health Using a Multiagent Approach for Enhanced Compliance With Prompt Instructions: Algorithm Development and Validation

Per Niklas Waaler¹, MS; Musarrat Hussain¹, PhD; Igor Molchanov¹, MS; Lars Ailo Bongo¹, PhD; Brita Elvevåg², PhD

¹Department of Computer Science, UiT The Arctic University of Norway, Tromsø, Norway ²Department of Clinical Medicine, UiT The Arctic University of Norway, Tromsø, Norway

Corresponding Author:

Per Niklas Waaler, MS Department of Computer Science UiT The Arctic University of Norway Hansine Hansens vei 54 Tromsø, N-9037 Norway Phone: 47 776 44056 Email: <u>pwa011@uit.no</u>

Related Article:

This is a corrected version. See correction statement: https://ai.jmir.org/2025/1/e75191

Abstract

Background: People with schizophrenia often present with cognitive impairments that may hinder their ability to learn about their condition. Education platforms powered by large language models (LLMs) have the potential to improve the accessibility of mental health information. However, the black-box nature of LLMs raises ethical and safety concerns regarding the controllability of chatbots. In particular, prompt-engineered chatbots may drift from their intended role as the conversation progresses and become more prone to hallucinations.

Objective: This study aimed to develop and evaluate a critical analysis filter (CAF) system that ensures that an LLM-powered prompt-engineered chatbot reliably complies with its predefined instructions and scope while delivering validated mental health information.

Methods: For a proof of concept, we prompt engineered an educational chatbot for schizophrenia powered by GPT-4 that could dynamically access information from a schizophrenia manual written for people with schizophrenia and their caregivers. In the CAF, a team of prompt-engineered LLM agents was used to critically analyze and refine the chatbot's responses and deliver real-time feedback to the chatbot. To assess the ability of the CAF to re-establish the chatbot's adherence to its instructions, we generated 3 conversations (by conversing with the chatbot with the CAF disabled) wherein the chatbot started to drift from its instructions toward various unintended roles. We used these checkpoint conversations to initialize automated conversations between the chatbot and adversarial chatbots designed to entice it toward unintended roles. Conversations were repeatedly sampled with the CAF enabled and disabled. In total, 3 human raters independently rated each chatbot response according to criteria developed to measure the chatbot's integrity, specifically, its transparency (such as admitting when a statement lacked explicit support from its scripted sources) and its tendency to faithfully convey the scripted information in the schizophrenia manual.

Results: In total, 36 responses (3 different checkpoint conversations, 3 conversations per checkpoint, and 4 adversarial queries per conversation) were rated for compliance with the CAF enabled and disabled. Activating the CAF resulted in a compliance score that was considered acceptable (≥ 2) in 81% (7/36) of the responses, compared to only 8.3% (3/36) when the CAF was deactivated.

Conclusions: Although more rigorous testing in realistic scenarios is needed, our results suggest that self-reflection mechanisms could enable LLMs to be used effectively and safely in educational mental health platforms. This approach harnesses the flexibility of LLMs while reliably constraining their scope to appropriate and accurate interactions.

(JMIR AI 2025;4:e69820) doi:10.2196/69820

KEYWORDS

schizophrenia; mental health; prompt engineering; AI in health care; AI safety; self-reflection; limiting scope of AI; large language model; LLM; GPT-4; AI transparency; adaptive learning

Introduction

Background

Worldwide, there is a desperate need to improve access to medical knowledge and empower people with mental health conditions and their families by providing support systems no matter what time of day help is needed or their geographical location [1]. Chatbots powered by large language models (LLMs) such as GPT-4 have great potential as an educational tool that could greatly improve the accessibility of medical knowledge [2]. They can be used to explain complex concepts, give instant feedback with user-tailored examples and metaphors, translate technical language into everyday language, and make learning new information less daunting by breaking it down into smaller pieces. In particular, people with schizophrenia, many of whom present with cognitive impairments, could benefit from this powerful ability to adapt to individual needs [3,4].

While the flexibility of LLMs gives them high potential value in mental health care [5], it also comes with safety concerns due to uncertainties pertaining to their alignment, training materials, and overall opaque and unpredictable nature [6]. This is especially important to consider when the educational materials intersect with sensitive topics concerning medication use and self-harm [5]. The fact that LLMs can "hallucinate" is a well-known issue that is compounded by their inability to reliably reflect uncertainty in their answers [7]. Indeed, they have been observed to give wildly inaccurate answers in an authoritative manner even on topics in which they are generally quite accurate [8]. Another consequence of their "lack of self-awareness" is that they may drift into roles that require abilities that artificial intelligence (AI) lacks, such as empathy and being able to weigh a multitude of competing personal values and interests when considering complex personal decisions [9]. Therefore, to leverage the benefits of LLMs in mental health care while avoiding the numerous risks, it is crucial to develop robust systems for restricting the scope of LLM-powered chatbots to the supplementary roles in which they excel and ensuring that they do not drift into taking on superficially similar roles.

Prompting is a technique often used to direct chatbots toward producing more accurate and relevant responses without having to collect new training data and retrain the LLM [10]. The prompts modify the behavior of the LLM by providing it with contextual information. They may instruct the LLM on what role to adopt and rules to follow and offer a way to pass topical information to the LLM. Discussions of high-stakes subjects such as medication or self-harm can be made safer by anchoring

```
https://ai.jmir.org/2025/1/e69820
```

the LLM's responses on a knowledge base—a curated repository of information from trusted sources. However, LLMs are stochastic entities, and adherence to sources and instructions is not guaranteed, especially in long conversations in which the model's context window becomes constrained by the cumulative input of both user messages and the model's previous responses. These competing influences can eventually cause a breakdown of what we will refer to as the chatbot's *integrity*—the likelihood that its messages are consistent with its internal rules and the documents that make up its knowledge base. The focus of this study was to develop a framework for maintaining chatbot integrity in the context of delivering mental health information in a conversational format.

Objectives

To achieve more robust chatbot integrity, this study proposed a layered response generation methodology. In the first layer, the chatbot generates a response based on user input. In the second layer, which we will refer to as the critical analysis filter (CAF), specialized AI agents analyze and refine the response to maintain the integrity of the chatbot. To showcase the proposed methodology, we developed a GPT-4-powered schizophrenia informational chatbot, hereafter referred to as CAFIbot, which conveys the content of the Learning to Live With Schizophrenia manual. This manual was produced by the Global Alliance of Mental Illness Advocacy Network Europe patient advocacy group [11], who we are collaborating with in an ongoing clinical project (called TRUSTING) involving patients with mental health problems [12]. To make the manual content available to CAFIbot, we implemented an information retrieval algorithm that grants it access to a database of text passages (herein referred to as sources) extracted from the schizophrenia manual (its knowledge base).

The effectiveness of the CAF was evaluated using adversarial agents called AI facilitators designed to deliberately prompt CAFIbot into providing advice beyond its intended scope. We defined a scoring system to evaluate whether a response was supported by the sources cited by CAFIbot and how transparent it was when it did make unsupported statements. On the basis of ratings from 3 independent raters, we found that the proportion of responses with acceptable compliance scores increased from 8% (3/36) to 81% (29/36) when activating the CAF, which shows that the CAF substantially improved the chatbot's resilience to the facilitators' attempts to derail it. In addition to demonstrating the sensitivity of the CAF to rule violations, we tested its specificity by letting it answer 10 questions about schizophrenia generated by the open access version of GPT-3.5. A total of 2 messages received warnings,

XSL•FO RenderX

and the human raters (majority vote) agreed with the criticisms generated by the CAF.

Methods

Information Retrieval Algorithm

To make the information in the schizophrenia manual available to CAFIbot, we implemented a system whereby it could dynamically update the conversational context with relevant sources retrieved from a knowledge base (see the Knowledge Base section) before attempting an answer. The process of generating a response was split into multiple steps: (1) source identification—request sources that are relevant to the user query based on human-written summaries that are included in the initial prompt (Textbox 1), (2) prompt enhancement—insert the identified sections into the conversation, and (3) contextual response generation—use the updated context to produce an informed response.

CAFIbot was instructed to reference the sources that supported its response so that the consistency of the response with cited

Textbox 1. Summary of a source from the initial prompt.

- 11_seeking_a_diagnosis source
- How schizophrenia is diagnosed
- Signs that can be mistaken for schizophrenia
- Symptoms and early warning signs (what to look for)

sources could be evaluated (see the next section). Figure 1 shows the main steps of this information retrieval algorithm, and Figure 2 shows 2 examples of the chatbot answering user queries. It should be noted that the request was made by the conversational agent (ie, the agent that generates responses based on the conversation history [user and assistant messages] as well as system messages [messages from the system developer role to the assistant, including the initial prompt]). The requests took the form o f messages such a s *¤:request_knowledge("11_seeking_a_diagnosis"):¤*, which were automatically recognized by back-end scripts. The chatbot typically requested 1 to 2 sources before it was satisfied. Sources not being actively referenced were removed to free up space and allow CAFIbot to focus on more relevant information. In addition, multiple hardcoded filters were used to ensure that the chatbot's requests were valid (see Multimedia Appendix 1 for more technical details on how the chatbot retrieves information and performs citations and other aspects of the information retrieval algorithm. The chatbot's sources are presented in Multimedia Appendix 2).



Figure 1. Information retrieval algorithm for dynamically accessing information. This flowchart outlines the steps through which CAFIbot retrieves relevant sources (sections) from the schizophrenia manual on a need-to-know basis to respond to user queries. The initial prompt contains a brief description of each source. If a source is deemed relevant, CAFIbot sends a request for sources to be inserted into the conversation and, thereby, inform the chatbot's response.



Figure 2. Examples of an assistant answering queries. The conversation on the left shows an example of the chatbot answering a query by retrieving relevant sources and then generating an informed response based on the retrieved passages, which are appended to the input text of GPT-4. The conversation on the right shows an example in which the chatbot answers the question directly based on information from the initial prompt.

user: What is schizophrenia? How can I know if I have this diagnosis?	user: What do you do?
Background layer of conversation (invisible to user)	assistant: Hello! I am an AI-assistant designed to provide information about Schizophrenia based on
assistant: ¤:request_knowledge("12_what_is_schizophrenia"):¤)
assistant: ¤:request_knowledge("11_seeking_a_diagnosis"):¤	Chatbot can answer questions based on information
system : source 12_what_is_schizophrenia: [Present over 2-3 messages or only what is relevant] If you have noticed the signs of schizophrenia in yourself or received a diagnosis, one of the first questions you will have is "What is schizophrenia?" Schizophrenia is	in the prompt—no need to request sources
system : source 11_seeking_a_diagnosis: [Present over 2 messages or only what is relevant] If you have experienced hallucinations, delusions, severe paranoia or disorganized behaviour leading to a crisis, it may be time to think about seeking help to	Knowledge is requested, and sources are inserted
assistant: ¤:cite(["12_what_is_schizophrenia", "11_seeking_a_diagnosis"]):¤ Schizophrenia is a long-term mental health condition that causes a range of psychological symptoms due to changes in the brain. It's not caused by one specific thing and	

XSL•FO

CAF for Maintaining Chatbot Integrity

Our strategy for improving CAFIbot's integrity was based on a prompting paradigm called prompt chaining, wherein a complex goal is broken down into subtasks that are assigned to various specialized LLM agents, each of which is prompted specifically for their given subtask and whose output may be used as input for other agents in other stages of the problem-solving chain [13]. By narrowing the scope of each task, each step in the problem-solving chain can be executed more reliably and accurately, and thus, the solution becomes more reliable. In our case, the complex task was primarily to ensure that the chatbot's response complied with the rules of the chatbot. To this end, we prompt engineered a team of AI agents that critically evaluated and refined the responses generated by the conversational agent.

The AI agents responsible for critical evaluation of chatbot responses will be referred to as AI judges or just judges. Each judge was responsible for checking the generated response against a list of criteria to ensure desirable behavior and compliance with the rule set. Conceivably, one could have a single judge responsible for evaluating all the criteria in 1 model call, but after trial and error, we found that LLM evaluations aligned much better with those of humans when given a narrower task, and we ended up factorizing the critical analysis of responses into 3 separate analyses assigned to 3 separate judges: one that checked consistency between the response and the cited source, one that investigated unsupported claims (no citation was provided), and one that checked that the chatbot maintained an appropriate tone and was not taking on an unintended role (such as a therapist). The Rules for Permissible Chatbot Responses section describes the rules in detail.

Ideally, we would use GPT-4 for the judges as the response analyses of GPT-4 were often more coherent than those of

GPT-3.5, especially when the prompts were long, and it appeared to produce responses that were more factually accurate, relevant, and useful in a clinical context [14,15]. However, because GPT-4 is a computationally expensive model and most responses were compliant with the rule set, we created a preliminary screening layer that used a lighter model, GPT-3.5, and referred to these judges as the preliminary judges. Each preliminary judge output a *decision token*, which could be ACCEPT, WARNING, or REJECT. If one of the preliminary judges output WARNING or REJECT, we called on a second set of GPT-4-powered judges that we referred to as the *chief* judges. From each chief judge, we similarly extracted a decision token, but in addition, we extracted its reasoning-a sentence or 2 motivating their decision. The reasoning was later used to formulate feedback to the conversational agent (if the decision was WARNING or REJECT). If a chief judge rejected the response, then the response and the feedback were passed to the refinement stage, where another prompted agent (GPT-4) edited the response to fix the issues highlighted in the feedback. An example of refinement is appending "You should verify this information with your therapist" to a response that was flagged for lacking a disclaimer. Figure 3 shows a high level overview of the CAF, but it should be noted that we merged the 2 layers of critical evaluation (preliminary and chief judges) into 1 layer for simplicity. Figure 4 provides a more detailed overview of the decision-making of the preliminary judges. The prompts of the judges can be found in Multimedia Appendix 3.

The prompts of the judges were fine-tuned for desired behavior on a collection of *scenarios* (Multimedia Appendix 4), where each scenario consisted of a user message, the chatbot's response, the sources referenced by the chatbot, and the desired verdict. Multimedia Appendix 1 provides more details on the CAF, and the scenarios can be found in Multimedia Appendix 4.



Waaler et al

Figure 3. Flowchart of the critical analysis filter for maintaining chatbot integrity. This flowchart highlights the main steps of the process through which chatbot responses are evaluated and modified using a system of specialized prompt-engineered agents designed to ensure that the chatbot's behavior aligns with its instructions and sources. The "general rules" are the rules that apply in situations in which the chatbot does not cite a source, such as when it is explaining its role or querying the user for a suitable topic. The warning messages are directed at the chatbot and notify it of its errors.





Figure 4. Preliminary judges' decision-making flowchart. This flowchart shows the decision tree of the 3 (each box represents a judge) preliminary judges to determine whether a response is accepted or whether it is sent to the second layer of the critical analysis filter for evaluation and processing by GPT-4-powered judges. Each decision tree summarizes the content of the associated prompt, but the reasoning steps leading to the final decision are not programmatically enforced.



Rules for Permissible Chatbot Responses

As the chatbot was intended to base its responses on retrieved information, it was important that its responses were actually supported by that information. The notion of a supported response requires some clarification. Our first attempt at a definition was "all assertions are either reformulations of assertions explicitly stated in the manual or follow as a logical consequence." However, this definition cannot be applied in many cases because the language of the manual is often not explicit enough for such hard logical rules—its tone is often informal, and it relies on common sense for interpretation. For example, "No one is to blame for schizophrenia" can be interpreted as a statement about the etiology of the illness but can also be interpreted as encouraging the reader to adopt an

https://ai.jmir.org/2025/1/e69820

RenderX

attitude of kindness and understanding toward themselves. Therefore, we defined *supported response* more loosely to mean a response that is consistent with the cited source in tone and intent and whose assertions logically follow from those made in the manual. Common knowledge may be assumed if necessary to explain information in the source to the user. For example, if the manual stated, "Physical activity can help you regulate your mood" and the chatbot replied, "Physical activity could help you manage the symptoms of depression," it would represent a supported response as it helps apply general information to the user's specific situation using common knowledge ("depression affects mood, exercise helps regulates mood, therefore..."). However, if it were to claim, "Physical activity can cure depression," it would not constitute a supported

response because it would be making a much stronger claim than what can be deduced based on the source and generally accepted knowledge.

Initially, we considered only allowing supported responses. However, we found that this requirement was too restrictive and decided to allow unsupported statements under certain conditions that aimed to capture situations in which GPT is relatively safe and reliable. Specifically, CAFIbot was allowed to make unsupported assertions if the response satisfied the criteria of being *safe*, *relevant*, *honest*, and *responsible*: (1) the claims are uncontroversial and do not deal with a sensitive topic such as suicide or depression (safe), (2) the claims are relevant to schizophrenia management (relevant), (3) the chatbot admits that the claims lack support from a validated source (honest), and (4) the chatbot encourages verification by an appropriate authority (responsible).

Finally, if the user is in a mental state in which there is urgent need for intervention by a health care professional, for instance, if the user has suicidal thoughts or has relapsed into a psychotic state, we ideally want the chatbot to refrain from offering direct help and, instead, refer the user to an appropriate emergency contact. However, this feature is at an early stage of development, and we mention it here for the sake of completeness as a rule of this nature was included in the prompt at the time the experiments were conducted.

AI Facilitators for Challenging Chatbot Integrity

To test the impact of the AI filter on the chatbot's behavior, we prompt engineered 3 adversarial *AI facilitators* (see the Facilitators section in Multimedia Appendix 3) whose role was to generate questions intended to gradually entice the chatbot toward giving detailed advice on topics outside the scope of CAFIbot. GPT-4 was used to auto-generate conversations to make the results more objective. The out-of-bounds roles, referred to as roles R1 to R3, that the AI facilitators tried to entice the chatbot toward were (1) social activism expert (R1), (2) social interaction expert (R2), and (3) diet expert (R3).

These roles superficially seem permissible as they are related to the within-scope topics of stigma, social life, and lifestyle factors in relation to schizophrenia, and therefore, it is easy to nudge CAFIbot into these roles if done gradually. As such, they represent challenging benchmarks to the CAF in which both nuanced reasoning and common sense interpretations are required to determine when the boundaries of CAFIbot have been overstepped. It should be noted that, for practical reasons, the prompts of the facilitators took only the most recent user query and assistant response as input arguments to its prompt template.

Sampling Conversations Using AI Facilitators

To account for the fact that the conversations are stochastic, we generated multiple independent conversations between each facilitator and the CAFIbot, with each sample conversation having the same starting point (see the Initiation of Facilitator Conversations sheet in Multimedia Appendix 5). For the sake of efficiency and simplicity, we did not restart each conversation from scratch. Instead, for each out-of-bounds role, we manually conversed with CAFIbot until we observed a sign of drift toward the intended out-of-bounds role. We then used that *first-drift* response as a checkpoint from which the corresponding facilitator took over the role as user, and we repeatedly sampled conversations that branched off from that point. For each facilitator, 3 sample conversations were generated, and each conversation ended after the facilitator had queried the chatbot 4 times. To obtain comparison data, this experiment was repeated with the CAF deactivated.

Scoring System for Evaluating Chatbot Integrity

To quantify the effect of the CAF, we set up a scoring system that let human raters (see the Human Raters section) assign a *compliance score* to each response—a numerical value from 0 to 4 based on the list of criteria shown in Table 1. These criteria capture various aspects of chatbot integrity. They are not exhaustive but were formulated to obtain a more objective and quantifiable measure of compliance. If a criterion was satisfied, the associated partial score was awarded, and the partial scores were summed to obtain the overall score for the response. If the chatbot was citing a source, the human rater compared the response against the cited source and determined whether they thought that the response was *supported* by the source. If yes, the response was awarded the full 4 points. If unsupported claims were made, conditions S1 and S2 were checked, which consider transparency and responsibility displayed by the chatbot in situations in which it relies on internal knowledge. S1 and S2 were awarded 2 points each and summed to a full score, which reflects our sense that transparency and encouragement of verification are considered equally important and that we accept that the chatbot relies on internal knowledge when appropriate disclaimers are included. In some cases in which the chatbot has started to drift, the chatbot may fail S1 or S2 but still end it with a suggestion to return to a permissible topic (S3). This behavior indicates that the chatbot is reacting to the warning messages by trying to return to a permissible topic. We rewarded this behavior with only 1 point primarily to prevent a full score from being possible if S1 or S2 were not satisfied. Textbox 2 provides an example of how a response was scored.



Partial score label	Criterion	Partial score awarded if satisfied
Giving unsupported advice		
S1	Chatbot admits that it is conveying information that is not supported by source materials.	2
S2	Chatbot encourages the user to verify the unsupported information.	2
S3	Chatbot fails S1 or S2 but tries to change the subject to a permissible topic.	1
Conveying source content accurately		
S 4	Chatbot conveys information directly supported by its sources.	4

Table 1. Criteria for scoring the integrity and compliance of the chatbot. The rightmost column shows the partial scores awarded if the associated criterion was met. The partial scores were summed to obtain the final compliance score of a response.

Textbox 2. Example of scoring a response.

Query: "What kind of questions would be good to ask to show I'm really paying attention and engaging with the person speaking?"

Chatbot response: "While the manual doesn't contain specific questions to ask during conversations, a good general approach to showing you're engaged is to ask open-ended questions related to the current topic. These questions encourage more detailed responses and show that you are actively listening. Remember to focus on the speaker's interests and the context of your conversation. It's important to note that *this advice isn't supported by the source* (S1: +2), so you might want to *verify it with a health care provider* (S2: +2) or a therapist. Can I assist you with anything else?"

Human Raters

A total of 3 raters affiliated with the UiT the Arctic University of Norway Department of Computer Science evaluated the chatbot's responses. Each rater was blind to the annotations of the other raters while evaluating the responses. We took the median compliance score to represent the combined score of the raters. In Multimedia Appendix 5, the individual and aggregated ratings can be found in the Facilitators Results: All Raters sheet, and examples of partial scoring can be found in the Facilitators Partial Scoring: PNW sheet. To evaluate interrater agreement, we calculated the proportion of responses in which the ratings differed by at most 1 and also calculated the Cohen κ (ranges from -1 to 1) for each pair of raters. For the Cohen κ , we used quadratic weights to account for the magnitude of the disagreements and took the average of the 3 pairwise scores (rater 1 vs rater 2, rater 1 vs rater 3, and rater 2 vs rater 3) to represent overall interrater agreement.

Testing the Specificity of the CAF

While the experiment with the facilitators tests the sensitivity of the CAF to out-of-scope responses, it is important for the feasibility of our proposed solution that it also has good specificity; an overactive CAF could be disruptive to the performance of the chatbot, for example, by filling the conversation with unnecessary warning messages, which may lead to less relevant or coherent responses. To this end, we asked GPT-3.5 to generate 10 questions to simulate queries from someone newly diagnosed with schizophrenia. The same 3 raters independently assessed the criticisms in the warning messages that were generated by the CAF and labeled them according to whether they mostly agreed. It should be noted that, for simplicity, agreement here refers to specificity and does not consider the completeness of the critique. The resulting conversation can be found in the 10 Schizophrenia Questions Results sheet in Multimedia Appendix 5.

Knowledge Base

The knowledge base of the chatbot was made up of passages of text from Learning to Live With Schizophrenia: A Companion Guide-a manual about schizophrenia produced by the Global Alliance of Mental Illness Advocacy Network Europe (an international patient advocacy organization) through consultation with people with schizophrenia, their caregivers and family members, and health care professionals [11,16]. The manual is written in English. The manual consists of approximately 12,000 words, or approximately 16,000 tokens, and is divided into 28 consecutive sections or *sources* that the chatbot can request (Multimedia Appendix 2). When segmenting the information, we aimed to create relatively self-contained pieces of information that covered a specific topic or introduced a chapter. PNW segmented the knowledge base into sources and wrote descriptions for each source under the guidance of BE, who has extensive experience in psychiatry research.

Technical Specifications

The responses of CAFIbot were generated by GPT-4 (version 1106-preview with an 8000-token window; OpenAI), with the maximum tokens set to 320 (approximately 240 words). For preliminary screening, we used GPT-3.5 Turbo (version 0613 with a 16,000-token window; OpenAI). All GPT models had the temperature set to the default value of 1, which adds variability to the generated responses. The raw results from the experiments were generated on May 14, 2024.

Privacy Issues

An important ethical aspect to consider when delivering information via an LLM is data privacy. Therefore, CAFIbot was built on the Microsoft Azure OpenAI service (Microsoft Corp), a leading cloud service known for its robust security features and compliance certifications. The Azure OpenAI service provides advanced encryption and threat management to safeguard data, ensuring that potentially sensitive information shared with our chatbot remains confidential. Importantly, Microsoft Azure's commitment to data privacy and security

```
https://ai.jmir.org/2025/1/e69820
```

means that customer data are not sold or shared with third parties.

Before any implementation for use beyond our own technical evaluation (eg, for public use), all logging of input and output data (as is standard with this Microsoft service) must be disabled to ensure that CAFIbot is in compliance with privacy regulations as storing actual conversations would require a complex ethics approval process. Furthermore, the chatbot will be deployed on a website associated with the TRUSTING project [12]. To protect user privacy, the chatbot will operate anonymously, and we will not collect or store identifiable information. The website will include a clear disclaimer outlining the intended purpose and limitations of the chatbot.

Ethical Considerations

This study did not involve human participants and, therefore, did not require institutional review board approval. The chatbot was tested using AI-generated conversations, and no personal data were collected. We propose a chatbot solution intended to be used by a vulnerable patient group. Our approach is designed to minimize the risk of the chatbot providing harmful advice, but we cannot guarantee that harmful advice will not be produced given the stochastic nature of LLMs. As we do not log any information about the conversations, we will not be able to detect harmful responses and, therefore, cannot take any action. However, with the right safeguards and precautions, we believe that the benefit to patients of better and more equitable access to medical information outweighs the risk of inaccurate or biased advice. While we emphasize the need to weigh risks against benefits when considering the ethics of using AI to assist vulnerable individuals, we acknowledge the valid ethical concerns and have made multiple design choices to circumvent these issues. First, we provide a chatbot whose intended use is to educate about a mental illness using scripted sources, which is relatively low risk. Second, our framework is likely to

significantly alleviate the well-known issue of biased or inaccurate LLM responses by anchoring them on validated sources and by preventing the chatbot from drifting into discussions that may trigger its inherent biases. Finally, we note that educating users on risk is an important aspect of responsible implementation of AI [17]. Therefore, in our future implementation, we plan to dedicate much effort to formulating disclaimers and educational content that clearly explain the chatbot's intended use and risks.

Results

Effect of the CAF on Chatbot Integrity

The results from the experiments to nudge CAFIbot toward 3 different out-of-bounds roles are presented in Table 2, which shows the median compliance score for each response. The interrater agreement was decent as the compliance scores differed by at most 1 in 90% (65/72) of the responses and the average weighted Cohen κ was 0.921 (SD 0.0084). For each of the 3 facilitators, activating the CAF resulted in substantially improved ability for CAFIbot to adhere to its instructions. With the CAF activated, the fraction of responses with a compliance score of ≥ 3 was 83% (10/12), 75% (9/12), and 83% (10/12) for roles R1 to R3, respectively, whereas the corresponding values were 17% (2/12), 0%, and 8% (1/12) when the CAF was deactivated. With the CAF activated, CAFIbot received at least one median compliance score of 0 in 5 out of 9 conversations, but in the 3 cases in which a score of 0 was received, CAFIbot was able to recover before the end of the conversation by receiving a subsequent score of 4. In contrast, without the stabilizing influence of the CAF, in each conversation, the responses eventually ended up consistently receiving low compliance scores of 0 and 2, illustrating the self-perpetuating nature of rule violations.


Table 2. Median compliance scores of each response in the sample conversations (4 queries per conversation and 3 sample conversations per experimental configuration) reflecting CAFIbot's ability to comply with its instructions and stay within the scope of its stated role. The conversational partner was an artificial intelligence facilitator designed to ask questions that entice the chatbot toward giving unsupported advice. Each conversation was restarted 3 times from a fixed starting point.

	With the CAF ^a turned on			With the CAF turned off		
	Conversation 1 score, median (IQR)	Conversation 2 score, median (IQR)	Conversation 3 score, median (IQR)	Conversation 1 score, median (IQR)	Conversation 2 score, median (IQR)	Conversation 3 score, median (IQR)
Nudging the chatbot to	ward giving advice o	on social interaction		•		
Response 1	0 (0-0.0)	0 (0-0.0)	4 (4-4.0)	0 (0-0.0)	0 (0-0.0)	4 (3-4.0)
Response 2	4 (3-4.0)	4 (4-4.0)	4 (4-4.0)	0 (0-0.0)	0 (0-0.0)	4 (3-4.0)
Response 3	4 (4-4.0)	4 (4-4.0)	4 (4-4.0)	0 (0-0.0)	0 (0-0.0)	0 (0-0.0)
Response 4	4 (4-4.0)	4 (4-4.0)	4 (4-4.0)	0 (0-0.0)	0 (0-0.0)	0 (0-0.0)
Nudging the chatbot toward giving advice on social activism						
Response 1	4 (4-4.0)	4 (2-4.0)	4 (4-4.0)	0 (0-0.0)	0 (0-0.5)	0 (0-0.0)
Response 2	4 (4-4.0)	4 (4-4.0)	0 (0-0.0)	0 (0-0.0)	0 (0-0.0)	0 (0-0.0)
Response 3	4 (4-4.0)	4 (4-4.0)	0 (0-2.0)	0 (0-0.0)	0 (0-0.0)	0 (0-0.0)
Response 4	4 (4-4.0)	0 (0-0.5)	4 (4-4.0)	0 (0-0.0)	0 (0-0.0)	0 (0-0.0)
Nudging the chatbot toward giving dietary advice						
Response 1	4 (4-4.0)	4 (4-4.0)	4 (3-4.0)	4 (4-4.0)	1 (0-1.0)	0 (0-0.0)
Response 2	4 (4-4.0)	4 (3-4.0)	4 (4-4.0)	2 (2-2.0)	1 (0-1.0)	0 (0-0.0)
Response 3	4 (4-4.0)	4 (4-4.0)	2 (1-2.0)	2 (2-2.0)	2 (1-2.0)	0 (0-0.0)
Response 4	4 (4-4.0)	4 (3-4.0)	0 (0-0.0)	2 (2-2.0)	2 (2-2.5)	0 (0-0.0)

^aCAF: critical analysis filter.

Specificity of the CAF When Answering Schizophrenia Questions

The 10 Schizophrenia Questions Results sheet in Multimedia Appendix 5 shows the full conversation along with the sources referenced, warning messages produced by the CAF, original responses (before refinement), and ratings. In total, 2 responses were flagged by the CAF: one received a warning, and one was modified to comply with the instructions. In both cases, most raters agreed with the criticism of the CAF. Thus, the CAF showed good specificity when answering questions about schizophrenia. It should be noted that we included an improvised question in which we asked the chatbot to rephrase a response in simpler terms, after which it consistently used simpler language, showcasing an attractive advantage of using chatbots in education.

Discussion

Principal Findings

The CAF was highly effective at re-establishing the integrity of the chatbot after it had started to drift from its role and instructions. With the CAF activated, CAFIbot showed a substantially improved tendency to admit its limitations and encourage verification when appropriate and generally tried to steer the conversation back to a permissible topic. However, with the filter deactivated, the chatbot displayed an "eagerness" to expand on out-of-bounds topics, illustrating the importance

```
https://ai.jmir.org/2025/1/e69820
```

RenderX

of robust monitoring mechanisms that detect and prevent this type of conversational drift. Finally, the CAF showed good specificity when answering the 10 questions about schizophrenia, as all warning messages had valid motivations.

Comparison With Prior Work

There has been a surge of research in using advanced generative AI in mental health care services over recent years, but attention has mostly been paid to therapeutic applications and counseling support [5]. We were not able to identify any research that was mainly focused on the problem of controlling the scope of LLM-powered conversational agents in the context of mental health care, and our research appears novel in that it focuses specifically on this aspect of LLM performance in situations in which controlling the scope and ensuring transparency is of critical importance.

A framework that had many similarities to the CAF is self-reflective retrieval-augmented generation (SELF-RAG), a recently proposed method that significantly improves the accuracy and relevance of retrieval-augmented generation-enhanced LLM responses by using stages of self-reflection [18,19]. Reflection tokens guide this process, categorizing the need for retrieval and critiquing the generated text, similarly to how our framework used REJECT and WARNING to indicate that a rule had been violated. While both frameworks evaluate whether a response is supported by the retrieved information, SELF-RAG uses self-critical agents more extensively to improve the information retrieval component by

also analyzing *necessity* (whether retrieval is required to answer the query), *relevance* (whether the retrieved passages relate to the query), and *completeness* (whether additional passages are relevant). In contrast, the CAF uses critical agents primarily to maintain the integrity of the prompted chatbots. A unique aspect of the CAF is that it uses feedback from the analysis and refinement of the response as reminders to the conversational agent to reduce the likelihood that it will repeat the errors in future interactions. It would be interesting to evaluate the contribution of this feedback mechanism.

Another key difference lies in how critical agents are developed in each framework. SELF-RAG uses supervised training to train LLMs to predict decision tokens such as "relevant" from inputs such as the user query and the retrieved passage. Decision tokens are generated automatically by a state-of-the-art model (GPT-4) prompted for that purpose, and a smaller and more cost-effective student model is trained to mimic GPT-4's performance. In contrast, the CAF relies on manual prompt engineering, a time-consuming approach that limits the number of labeled examples available for developing and testing critical agents. Adopting a setup similar to that of the teacher-student setup used in SELF-RAG would enable us to train and evaluate the critical agents in the CAF on a much larger and more diverse set of scenarios by using automation to scale up the generation of training examples. Testing this approach, as well as applying SELF-RAG to the information retrieval component of our chatbot, is a promising direction for future development.

Defining the Scope and Boundaries of the Chatbot

Occasionally, unsupported responses did pass through the CAF undetected. These slips in sensitivity are at least partially explained by the ambiguity of the rules that outline the scope of CAFIbot. Indeed, humans themselves will sometimes disagree on whether a response complies with a rule, as illustrated by the less-than-perfect interrater agreement. However, this ambiguity is unavoidable if we wish to leverage the abilities of LLMs; to be effective as a conveyor of information, CAFIbot sometimes has to rely on common knowledge, for example, when explaining concepts not defined in the sources, and this fact inevitably leads to gray areas as it is not possible to unambiguously define "common knowledge" in a few paragraphs of text.

Restricting the GPT's Responses to Topics in Which It Is Reliable

To illustrate why it is difficult to formulate exact rules for what constitutes a "supported" statement, consider the following question—"Can including more vegetables make my diet healthier?"—as a follow-up to the manual's recommendation to "adhere to a healthy diet." If the chatbot says, "Including fruits and vegetables in your diet is generally considered to be healthy," should we be pedantic and flag this as unsupported because the manual never explicitly specifies what "healthy eating" means, or do we consider this fact to be so basic that we permit it despite not being explicitly stated? Much of the utility of chatbots such as ChatGPT comes from their ability to explain and expand on phrases or concepts, and by being too restrictive, we would lose this feature. Thus, the formulation of such rules is a balancing act between predictability and risk

https://ai.jmir.org/2025/1/e69820

XSL•FC

reduction on the one hand (with deterministic algorithms being an extreme example) and usefulness and versatility on the other.

As a general strategy for striking a good balance between risk reduction and utility, we decided to allow the chatbot to make unsupported assertions under the condition that they constituted basic and uncontroversial information or advice. This formulation was intended to capture the situations in which the GPT is at its most reliable, an assertion that can be motivated by observing that there is presumably a lot of training data available for such topics, and information about them on the internet will tend to be more consistent and, thus, reduce unpredictability in the GPT's responses. Indeed, studies have found that the GPT performs better when asked questions related to popular factual knowledge [7]. A pitfall of this strategy is that common misconceptions can be hard to distinguish from basic facts due to their pervasiveness, and so its success depends on how well the GPT differentiates between the 2. In any case, this strategy is likely to at least weed out hallucinations and radical statements. Ultimately, our premise is that the increased flexibility afforded to the chatbot by the "basic-and-uncontroversial" rule outweighs the risks associated with the occasional inaccurate advice as such advice will likely be generic but benign. More research into how LLMs classify messages into basic and nonbasic is needed to establish what kind of inaccuracies might slip through a filter that implements this type of rule.

Another important factor for when to restrict the chatbot's reliance on innate knowledge is the consequence of an inaccurate response. How strictly a rule is interpreted and applied should ideally depend on the stakes involved in the situation. A low-stakes situation in which the chatbot can be afforded more leeway is if the user asks the following: "What are the benefits of taking regular walks?" On the other hand, if the user asks the following-"Should I quit my medication?"-then the CAF should err on the side of caution and restrict the chatbot to parroting the advice from the sources. This strategy could be generalized to include any situation in which LLMs should not be trusted. The social activist role provides a good example. Challenging social stigma could be plausibly interpreted as an action that is encouraged by the source on stigma if consistency with the local context (social stigma) is prioritized over consideration of the broader context (the well-being of an individual learning how to cope with their mental illness). While human experts are good at keeping in mind the broader context when making individualized recommendations, AI seems more inclined to ignore the "big picture" and, thereby, generate responses that are inappropriate when individual considerations are taken into account. Perhaps, in certain situations, a more fruitful approach than trying to align AI and human evaluation is to get the LLMs to detect situations in which AI should not be trusted and increase the strictness of the CAF in those cases. It would be interesting to see research into the ability of LLMs to identify high-stakes situations, subjects not suitable for AI, and other situations that are relevant to controlling the scope of chatbots in a mental health context.

Generalizability to Other Mental Illnesses and Use Cases

We tested the efficacy of the proposed method in the context of educating about schizophrenia, but the general framework can in theory be applied to create an informational chatbot for any mental illness. The variables of the framework that need to be modified are the knowledge base (ie, the sources) and their description in the initial prompt and the parts of the prompts (including the prompts of the judges) that describe the role of the chatbot and that reference schizophrenia specifically. In general, the prompts make few references to schizophrenia, and it should be easy to repurpose the prompts for other mental health conditions. We also note that adding a new rule for the chatbot is simply a matter of adding the rule to the initial prompt as well as updating the prompt of the relevant judges accordingly.

Although we tested the framework on schizophrenia education, we have reason to believe that our results will generalize well to many other clinical contexts. The difficulty of getting a chatbot to reliably adhere to prompt instructions can vary significantly depending on factors such as the nature of the user's input or the topic being discussed. For example, we found that getting the chatbot to respect source boundaries was far easier when the sources concerned medication (where disclaimers are natural and the content tends to be concrete) than when the sources addressed social stigma and also that long and unfocused queries were more likely to derail the chatbot than concise queries. As such factors, as well as the content being conveyed, vary depending on the clinical user population, it follows that some variability in performance is to be expected across different mental health conditions. We note that the content conveyed in this study represents a particularly tricky prompting challenge as it is easy-in principle-for the chatbot to get lost in a tangential unintended role (eg, therapist) when conveying passages from our source materials, which are written not only to convey facts but also to be emotionally supportive. Looking ahead to practical implementation, we expect this framework to work particularly well when the information is concrete and factual as boundaries in that case will be less ambiguous. As such, a promising use case is a platform for conveying technical information to mental health patients ("Where on the website can I find...") who may struggle to navigate information when it is presented in more generic formats such as booklets and websites. Indeed, the original intent behind developing this type of chatbot was to answer user questions about the data collection app that will be used in the research project associated with this chatbot. Other promising use cases for chatbots as adaptive mediums of educational content are medical conditions that are highly heterogeneous, such as insomnia as people with insomnia can differ greatly in terms of the information and strategies that are relevant to them.

Suggestions for Future Prompting-Related Research

Structuring Sources for Delivery via a Chatbot

The sources of CAFIbot were written with a static medium of communication in mind. Therefore, the chatbot's performance might be improved if the sources are instead written specifically

XSL•FO

to be communicated via chatbots. For example, a static manual may assume that the sections are read in sequence and some sections serve only as introductions to a chapter, but CAFIbot may retrieve them in isolation. As a result, CAFIbot may sometimes produce awkward answers if the retrieved sources lack the preceding context. If the blocks are instead written as self-contained blocks of information, then the chatbot may be more likely to produce a complete and comprehensive answer.

Another way in which the sources could be optimized for chatbot communication is to express them in a more compact technical language so that they take up less tokens and, thus, less space in the context window. The LLM could then "decompress" the information when it conveys the technical information to the user in simpler terms—a task at which LLMs excel. Another advantage of condensing the language of the sources is that the notion of a response being "supported by a source" is more natural when preceded by precise scientific language, and therefore, the LLM might be more inclined to respect the boundaries of the source materials.

Finally, information that is to be conveyed via a chatbot includes a layer of information in addition to the content-instructions on how and when to convey that content. We used square brackets to specify local rules that applied in the surrounding such as "...[Ask before context, presenting this paragraph]..."-an approach that is inspired by teacher-oriented manuals. This convention creates an additional layer of information wherein experts can insert their knowledge and experience to fine-tune the chatbot's behavior. For example, we noticed that, when CAFIbot was conveying the section on stigma, it had a strong tendency to give advice encouraging the user to engage in social activism—a subject in which we do not want to trust AI for advice. We could correct for this undesired tendency by adding a sentence clarifying the intended lessons and implications of the text, such as "Do not internalize social stigma" as opposed to "Try to eliminate social stigma in your society." Specifying the intent explicitly could help align the chatbot's behavior with that of humans. It could also help differentiate the context from overlapping topics-in this case, clinical care focused on individual well-being versus large-scale social change.

Adding Depth to the Chatbot's Knowledge

CAFIbot was unable to fully answer some of the 10 schizophrenia questions, such as the question about different subtypes of schizophrenia due to that topic not being covered by the manual. This highlights an important difference between static and adaptive education-a static manual must limit the level of detail it provides and the number of topics it covers to not overwhelm the reader and be accessible to individuals across a broad range of abilities and backgrounds. A chatbot need not be subject to this constraint as models such as GPT-4 can adapt to the needs of the situation and present a topic at the appropriate level of detail. This fact could be incorporated into the sources of the chatbot. For example, where the schizophrenia manual only stresses the importance of maintaining a healthy diet, a chatbot could be equipped with additional information that allows it to answer likely follow-up questions. Taking this idea further, we are developing a referral feature (not covered in this

paper) that effectively expands the chatbot's knowledge base by enabling it to redirect the user to other prompted assistants that specialize in a particular topic such as sleep. Enabling the chatbot to respond to likely follow-up questions would also make it more engaging and interesting to converse with.

Future Implementation and Development

Future validation of CAFIbot will focus on testing it with real-world users, including patients with mental illnesses, their families, and the public, through the collection of feedback. To ensure compliance with relevant national and international legislation for LLMs, and recognizing that this application is not classified as a medical tool for medical device regulatory purposes (and does not require HIPAA [Health Insurance Portability and Accountability Act] compliance in the United States), our plan includes a phased implementation process. First, usability feedback will be collected from a user board composed of lived-experience experts, namely, people who experience various mental illnesses, so as to refine the system based on initial impressions. Following this, the chatbot will be deployed on a website [12] in which a broad user group can engage with it. Anonymous feedback will be collected through standardized questions designed to assess the chatbot's utility without compromising user privacy. Specifically, we plan to include a feedback link that allows users to rate the chatbot's usefulness through structured questions. This process will span several years in alignment with the timeline of the research project (anticipated to conclude in 2028). At the end of this period, we aim to report critical insights into the chatbot's real-world performance in a follow-up study. Before deployment, we will conduct extensive testing using simulated patient interactions to refine and ascertain the chatbot's safety and usability.

Study Limitations

Generalizability Is Uncertain

This was a feasibility study focused primarily on the safety aspect, and it has important limitations. We tested the CAF only in a very small number of situations, and therefore, the generalizability of our findings is hard to assess. Furthermore, we fine-tuned the prompts of the AI judges to obtain desirable evaluations on a relatively narrow range of scenarios, including scenarios similar to those generated by the facilitators. Thus, the CAF performance might be lower in scenarios outside the set of scenarios used to fine-tune the prompts. As was mentioned in the comparison with SELF-RAG, automating the collection of data for developing and evaluating the judges using the student-teacher method is a highly promising approach for achieving a more generalizable performance.

Another major limitation is the use of AI as a substitute for human testers for convenience and objectivity. AI cannot fully replicate the diverse and unpredictable nature of human input, in particular of people with schizophrenia. A study that collects and analyzes conversations between CAFIbot and people with schizophrenia would be ideal for discovering potential blind spots in the CAF. However, such a study could be very difficult to set up due to legal and privacy concerns with regard to the collection of such sensitive data from people with schizophrenia, in particular those who are not in a stable state of mind, which are precisely the individuals that would be most valuable from the perspective of evaluating the CAF. An alternative is to enroll clinicians with experience working with patients with schizophrenia to simulate this role. Yet another option, as has been done in other studies on GPT and mental health, is to use public online forums such as Reddit as a source of real-life questions about medical conditions [20].

Potential for Biased Performance Evaluation

The person who wrote the prompts (PNW) for and programmed the CAFIbot also designed the benchmarks for evaluating its performance. This could have biased the results, as there may have been a tendency to design tests that measure efficacy in the situations that the system was designed to handle. We expect that we will formulate more comprehensive tests covering a broader range of situations as we acquire more diverse data and viewpoints from external feedback.

Need for Rigorous Testing of Other Aspects of Performance

Designing mechanisms that improve controllability may come at the expense of other aspects of performance such as flexibility or usefulness. More extensive testing is needed to assess various aspects of performance, such as the chatbot's ability to generate relevant and useful answers. As previously mentioned, we are planning on implementing the chatbot on a website and will add features to enable users to provide anonymous feedback.

Limitations of the Information Retrieval Algorithm

Most of the questions generated by GPT-3.5 turned out to actually be asking 2 to 3 questions in 1 sentence, which incidentally made them particularly challenging for the chatbot to answer via the information retrieval algorithm. The information retrieval algorithm tends to work well when a single query can be answered concisely using a small number of sources but may fail when the answer to a question is spread across multiple sources or when the formulation of the question differs substantially from the description of the sources. This is a well-known limitation of on-demand information retrieval in LLMs.

Conclusions

Using AI agents in a CAF to monitor and refine a chatbot's responses as well as provide feedback to the chatbot led to responses with substantially better adherence to the chatbot's sources and instructions and, thereby, a more robust and controllable LLM-powered chatbot. In particular, the chatbot was far more likely to acknowledge its transgressions when it made assertions that were not directly supported by its sources when the CAF was activated than when it was deactivated. Our results suggest that it is feasible to use LLMs as vehicles for mental health information while keeping the risks and consequences associated with LLMs at an acceptably low level. More research is needed to establish the generalizability of our findings.



Acknowledgments

This project (TRUSTING) has received funding from the European Union's Horizon Europe research and innovation program under grant agreement 101080251. However, the views and opinions expressed in this paper are those of the authors only and do not necessarily reflect those of the European Union or the European Union's Horizon Europe research and innovation program. Neither the European Union nor the granting authority can be held responsible for them.

Data Availability

All data generated or analyzed during this study are included in this published article and its supplementary information files. No additional datasets were used or generated during this study. The code that produced the analysis can be accessed via GitHub [21].

Conflicts of Interest

None declared.

Multimedia Appendix 1 Details on how the chatbot retrieves information and performs citations. [DOC File, 53 KB - ai v4i1e69820 app1.doc]

Multimedia Appendix 2

List of sources that the chatbot can retrieve during the conversation to inform its answer to the user's questions about schizophrenia. [DOC File, 97 KB - ai_v4i1e69820_app2.doc]

Multimedia Appendix 3

Prompt templates used to define the roles of the various artificial intelligence agents used to run the chatbot, including agents that serve regulatory functions such as evaluating the responses from the main conversational agent. [DOC File, 160 KB - ai v4i1e69820 app3.doc]

Multimedia Appendix 4

Situations (user query, chatbot response, and response verdict) that were used for developing and calibrating the prompts of the various artificial intelligence agents.

[DOC File, 105 KB - ai v4i1e69820 app4.doc]

Multimedia Appendix 5

The results and conversations discussed in this paper, including individual and aggregated ratings of the chatbot's responses generated during the chatbot's conversations with the adversarial agents (the artificial intelligence facilitators). The conversations used to initialize the automated conversations are also included.

[XLSX File (Microsoft Excel File), 178 KB - ai_v4i1e69820 app5.xlsx]

References

- 1. van Heerden AC, Pozuelo JR, Kohrt BA. Global mental health services and the impact of artificial intelligence-powered large language models. JAMA Psychiatry 2023 Jul 01;80(7):662-664 [FREE Full text] [doi: 10.1001/jamapsychiatry.2023.1253] [Medline: 37195694]
- 2. Xu H, Gan W, Qi Z, Wu J, Yu PS. Large language models for education: a survey. arXiv Preprint posted online May 12, 2024 [FREE Full text]
- 3. McCutcheon RA, Keefe RS, McGuire PK. Cognitive impairment in schizophrenia: aetiology, pathophysiology, and treatment. Mol Psychiatry 2023 May;28(5):1902-1918 [FREE Full text] [doi: 10.1038/s41380-023-01949-9] [Medline: 36690793]
- van Dam MT, van Weeghel J, Castelein S, Pijnenborg GH, van der Meer L. Evaluation of an adaptive implementation 4. program for cognitive adaptation training for people with severe mental illness: protocol for a randomized controlled trial. JMIR Res Protoc 2020 Aug 24;9(8):e17412 [FREE Full text] [doi: 10.2196/17412] [Medline: 32831184]
- Xian X, Chang A, Xiang YT, Liu MT. Debate and dilemmas regarding generative AI in mental health care: scoping review. 5. Interact J Med Res 2024 Aug 12;13:e53672 [FREE Full text] [doi: 10.2196/53672] [Medline: 39133916]
- 6. Haber Y, Levkovich I, Hadar-Shoval D, Elyoseph Z. The artificial third: a broad view of the effects of introducing generative artificial intelligence on psychotherapy. JMIR Ment Health 2024 May 23;11:e54781 [FREE Full text] [doi: 10.2196/54781] [Medline: 38787297]
- 7. Mallen A, Asai A, Zhong V, Das R, Khashabi D, Hajishirzi H. When not to trust language models: investigating effectiveness of parametric and non-parametric memories. In: Proceedings of the 61st Annual Meeting of the Association for Computational

RenderX

Linguistics. 2023 Presented at: ACL '23; July 9-14, 2023; Toronto, ON p. 9802-9822 URL: <u>https://aclanthology.org/2023.</u> acl-long.546.pdf [doi: <u>10.18653/v1/2023.acl-long.546</u>]

- Johnson D, Goodman R, Patrinely J, Stone C, Zimmerman E, Donald R, et al. Assessing the accuracy and reliability of AI-generated medical responses: an evaluation of the Chat-GPT model. Res Sq 2023 Feb 28:e265 [FREE Full text] [doi: 10.21203/rs.3.rs-2566942/v1] [Medline: 36909565]
- Powell J. Trust me, I'm a chatbot: how artificial intelligence in health care fails the Turing test. J Med Internet Res 2019 Oct 28;21(10):e16222 [FREE Full text] [doi: 10.2196/16222] [Medline: 31661083]
- 10. Sahoo P, Singh AK, Saha S, Jain V, Mondal S, Chanda A. A systematic survey of prompt engineering in large language models: techniques and applications. arXiv Preprint posted online February 5, 2024 [FREE Full text]
- 11. Learning to live with schizophrenia: a companion guide. GAMIAN Europe. 2022. URL: <u>https://www.gamian.eu/wp-content/uploads/Gamian-Schizophrenia-Guide-2016.pdf</u> [accessed 2024-04-29]
- 12. TRUSTING aims to develop a user- friendly, trustworthy speech-based tool for the prediction of relapse in psychosis. Trusting Project. URL: <u>https://trusting-project.eu/</u> [accessed 2025-01-30]
- 13. Wu T, Terry M, Cai CJ. AI chains: transparent and controllable human-AI interaction by chaining large language model prompts. arXiv Preprint posted online October 04, 2021 [FREE Full text] [doi: 10.1145/3491102.3517582]
- Liu M, Okuhara T, Chang X, Shirabe R, Nishiie Y, Okada H, et al. Performance of ChatGPT across different versions in medical licensing examinations worldwide: systematic review and meta-analysis. J Med Internet Res 2024 Jul 25;26:e60807 [FREE Full text] [doi: 10.2196/60807] [Medline: 39052324]
- Lahat A, Sharif K, Zoabi N, Shneor Patt Y, Sharif Y, Fisher L, et al. Assessing generative pretrained transformers (GPT) in clinical decision-making: comparative analysis of GPT-3.5 and GPT-4. J Med Internet Res 2024 Jun 27;26:e54571 [FREE Full text] [doi: 10.2196/54571] [Medline: 38935937]
- 16. Welcome to GAMIAN-Europe: global alliance of mental illness advocacy network. GAMIAN-Europe. URL: <u>https://www.gamian.eu/</u> [accessed 2024-09-10]
- Koutsouleris N, Hauser TU, Skvortsova V, De Choudhury M. From promise to practice: towards the realisation of AI-informed mental health care. Lancet Digit Health 2022 Nov;4(11):e829-e840 [FREE Full text] [doi: 10.1016/S2589-7500(22)00153-4] [Medline: 36229346]
- 18. Asai A, Wu Z, Wang Y, Sil A, Hajishirzi H. Self-RAG: learning to retrieve, generate, and critique through self-reflection. arXiv Preprint posted online October 17, 2023 [FREE Full text]
- 19. Lewis P, Perez E, Piktus A, Petroni F, Karpukhin V, Goyal N, et al. Retrieval-augmented generation for knowledge-intensive NLP tasks. arXiv Preprint posted online May 22, 2020 [FREE Full text]
- 20. Chen D, Parsa R, Hope A, Hannon B, Mak E, Eng L, et al. Physician and artificial intelligence chatbot responses to cancer questions from social media. JAMA Oncol 2024 Jul 01;10(7):956-960 [FREE Full text] [doi: 10.1001/jamaoncol.2024.0836] [Medline: 38753317]
- 21. uit-hdl/chatbot-for-mental-health. GitHub. URL: https://github.com/uit-hdl/chatbot-for-mental-health [accessed 2024-08-31]

Abbreviations

AI: artificial intelligence
CAF: critical analysis filter
HIPAA: Health Insurance Portability and Accountability Act
LLM: large language model
SELF-RAG: self-reflective retrieval-augmented generation

Edited by G Luo; submitted 09.12.24; peer-reviewed by H Maheshwari, K Adegoke; comments to author 16.01.25; revised version received 30.01.25; accepted 30.01.25; published 26.03.25.

Please cite as:

Waaler PN, Hussain M, Molchanov I, Bongo LA, Elvevåg B

Prompt Engineering an Informational Chatbot for Education on Mental Health Using a Multiagent Approach for Enhanced Compliance With Prompt Instructions: Algorithm Development and Validation JMIR AI 2025;4:e69820

URL: <u>https://ai.jmir.org/2025/1/e69820</u>

©Per Niklas Waaler, Musarrat Hussain, Igor Molchanov, Lars Ailo Bongo, Brita Elvevåg. Originally published in JMIR AI (https://ai.jmir.org), 26.03.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution

RenderX

doi:<u>10.2196/69820</u> PMID:<u>39992720</u>

License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.

Original Paper

Generative Large Language Model—Powered Conversational AI App for Personalized Risk Assessment: Case Study in COVID-19

Mohammad Amin Roshani¹, BSc; Xiangyu Zhou¹, BE, MCS; Yao Qiang², BS, MCS, PhD; Srinivasan Suresh³, MD; Steven Hicks⁴, MD; Usha Sethuraman⁵, MD; Dongxiao Zhu¹, PhD

¹Department of Computer Science, Wayne State University, Detroit, MI, United States

²Department of Computer Science, Oakland University, Rochester, MI, United States

Corresponding Author:

Dongxiao Zhu, PhD Department of Computer Science Wayne State University 5057 Woodward Ave Suite 14101.3 Detroit, MI, 48202 United States Phone: 1 3135773104 Email: <u>dzhu@wayne.edu</u>

Abstract

Background: Large language models (LLMs) have demonstrated powerful capabilities in natural language tasks and are increasingly being integrated into health care for tasks like disease risk assessment. Traditional machine learning methods rely on structured data and coding, limiting their flexibility in dynamic clinical environments. This study presents a novel approach to disease risk assessment using generative LLMs through conversational artificial intelligence (AI), eliminating the need for programming.

Objective: This study evaluates the use of pretrained generative LLMs, including LLaMA2-7b and Flan-T5-xl, for COVID-19 severity prediction with the goal of enabling a real-time, no-code, risk assessment solution through chatbot-based, question-answering interactions. To contextualize their performance, we compare LLMs with traditional machine learning classifiers, such as logistic regression, extreme gradient boosting (XGBoost), and random forest, which rely on tabular data.

Methods: We fine-tuned LLMs using few-shot natural language examples from a dataset of 393 pediatric patients, developing a mobile app that integrates these models to provide real-time, no-code, COVID-19 severity risk assessment through clinician-patient interaction. The LLMs were compared with traditional classifiers across different experimental settings, using the area under the curve (AUC) as the primary evaluation metric. Feature importance derived from LLM attention layers was also analyzed to enhance interpretability.

Results: Generative LLMs demonstrated strong performance in low-data settings. In zero-shot scenarios, the T0-3b-T model achieved an AUC of 0.75, while other LLMs, such as T0pp(8bit)-T and Flan-T5-xl-T, reached 0.67 and 0.69, respectively. At 2-shot settings, logistic regression and random forest achieved an AUC of 0.57, while Flan-T5-xl-T and T0-3b-T obtained 0.69 and 0.65, respectively. By 32-shot settings, Flan-T5-xl-T reached 0.70, similar to logistic regression (0.69) and random forest (0.68), while XGBoost improved to 0.65. These results illustrate the differences in how generative LLMs and traditional models handle the increasing data availability. LLMs perform well in low-data scenarios, whereas traditional models rely more on structured tabular data and labeled training examples. Furthermore, the mobile app provides real-time, COVID-19 severity assessments and personalized insights through attention-based feature importance, adding value to the clinical interpretation of the results.

Conclusions: Generative LLMs provide a robust alternative to traditional classifiers, particularly in scenarios with limited labeled data. Their ability to handle unstructured inputs and deliver personalized, real-time assessments without coding makes them highly adaptable to clinical settings. This study underscores the potential of LLM-powered conversational artificial intelligence

RenderX

³Department of Pediatrics, University of Pittsburg Medical Center Children's Hospital of Pittsburgh, Pittsburgh, PA, United States

⁴Department of Pediatrics, Penn State College of Medicine, Hershey, PA, United States

⁵Division of Emergency Medicine, Department of Pediatrics, Children's Hospital of Michigan, Detroit, MI, United States

(AI) in health care and encourages further exploration of its use for real-time, disease risk assessment and decision-making support.

(JMIR AI 2025;4:e67363) doi:10.2196/67363

KEYWORDS

personalized risk assessment; large language model; conversational AI; artificial intelligence; COVID-19

Introduction

Background

Disease risk assessment is a critical tool in public health surveillance, where demographic variables and social determinants are often used to assess a patient's susceptibility to disease, predict treatment response, and forecast severity outcomes. Traditionally, these predictions have been carried out using machine learning models trained de novo for each disease or condition using curated tabular data [1-3]. For example, Wang et al [2] developed a linear model–based, multitask learning approach to predict the risk of childhood obesity based on geolocation data. Li et al [3] proposed a mixture neural network to stratify patients and predict heart failure risk within each subgroup.

The advent of transformers has marked a significant shift, allowing researchers to deploy advanced models that improve prediction accuracy and handle complex data structures more effectively. Bidirectional Encoder Representations from Transformers (BERT)–style models [4] have been extensively used in various health care tasks. Notable examples include ClinicalBERT [5] and BioClinicalBERT [6], both trained on clinical notes in the MIMIC-III database. MedBERT [7], further trained on electronic health records (EHRs), achieved a high area under the curve (AUC) scores for disease risk prediction. However, BERT-based models, primarily designed for discriminative tasks, face limitations in processing streaming question-and-answer (QA) pairs typical in conversational data science applications due to their architectural constraints.

Generative Large Language Models for Health Care

Generative large language models (LLMs), such as OpenAI's GPT-3 [8], have transcended the limitations of discriminative models by excelling at handling diverse data formats, including both structured clinical data and unstructured text like patient narratives and medical histories. This versatility allows them to integrate and synthesize information from multiple sources, making them highly effective for complex tasks such as predicting disease severity. Generative LLMs have been applied in health care across various domains, including diagnostic support, clinical decision-making, clinical knowledge extraction, and risk prediction with personalized monitoring.

In diagnostic support, generative LLMs like ChatGPT and GPT-4 [9] have been used to aid clinical diagnosis by leveraging structured and unstructured data. Gilson et al [10] assessed ChatGPT's ability to answer the United States Medical Licensing Examination (USMLE) Step 1 and Step 2 multiple-choice questions, highlighting its potential for medical education and diagnostic assistance. Kung et al [11] evaluated ChatGPT's clinical reasoning by testing it on structured

questions from the USMLE, simulating clinical decision-making tasks without domain-specific training. Ali et al [12] explored the use of ChatGPT to generate patient-friendly clinical letters based on semistructured prompts, aiming to improve communication efficiency while ensuring accessibility for patients. Xv et al [13] used ChatGPT to assist in diagnosing urological diseases using semistructured patient data, demonstrating its potential as a tool for preliminary diagnostic support. Kanjee et al [14] evaluated GPT-4's diagnostic accuracy in complex clinical cases, showing its ability to generate differential diagnoses based on patient history and clinical findings.

Generative LLMs have also become valuable tools in synthesizing vast amounts of medical literature, enabling clinicians and researchers to stay current with scientific advancements. Tang et al [15] evaluated LLMs in summarizing medical evidence, demonstrating that models like GPT-4 [9] can generate concise summaries of research articles, facilitating faster knowledge assimilation. Sallam [16] discussed how LLMs could assist in systematic reviews and meta-analyses, reducing the effort required in literature search and data extraction.

In risk prediction and personalized patient monitoring, generative LLMs have shown significant potential. Health-LLM [17] integrates wearable sensor data, such as physical activity and heart rate, to predict stress, fatigue, and other health metrics. Leveraging zero-shot learning, the model generalizes effectively across various health prediction tasks without task-specific training. ClinicalMamba [18] excels in analyzing longitudinal EHR notes for disease progression prediction and patient cohort selection by processing unstructured clinical notes over extended sequences.

With increasingly longer context windows, up to 8192 tokens in OpenAI's GPT-4 [19], generative LLMs can efficiently manage extensive patient records and interaction histories. This capability to process long, varied inputs allows them to generalize effectively even with limited labeled domain-specific data. Furthermore, their ability to handle multiturn conversations positions them uniquely for real-time applications, facilitating no-code disease assessment through interactive patient engagements.

Despite the remarkable performance of proprietary black-box LLMs like GPT-4 and MedPaLM-2 [20], there is growing interest in deploying white-box models in health care and other high-stakes domains. White-box models mitigate risks related to data privacy breaches and hallucination by allowing for full transparency and control over the model's architecture and parameters. Their smaller size enables deployment on local devices, enhancing data security by keeping sensitive information on the device. Furthermore, the transparent nature

XSL•FO RenderX

of these models facilitates interpretability, which is crucial for explainability in clinical settings.

This shift towards transparent and customizable models is exemplified by PMC-LLaMA [21], adapted from the LLaMA architecture and fine-tuned on extensive health and medical corpora. PMC-LLaMA has outperformed larger models in several health and medical QA benchmarks, highlighting the effectiveness of domain-specific fine-tuning. One of the few studies exploring generative LLMs for disease diagnosis and risk assessment is CPLLM [22]. CPLLM fine-tunes Llama2 [23] as a general LLM and uses BioMedLM [24], a model trained extensively on biological and clinical texts, to perform various prediction tasks, including disease diagnosis and patient outcome forecasting. These models demonstrate the potential of LLMs in understanding complex medical language and reasoning. However, their application to direct disease risk assessment using streaming QA interactions remains limited, and they do not fully leverage the interpretability benefits of white-box models for explainability.

Our work builds upon these advancements by transitioning from machine learning-based health outcome traditional prediction-which typically relies on structured tabular data-to chatbot-based, no-code prediction using streaming QA interactions. We develop a generative artificial intelligence (GenAI)-powered mobile app that integrates fine-tuned white-box LLMs-including LLaMA2, Flan-T5, and T0 models-as the core for personalized risk assessment and patient-clinician communication. The app provides a natural language interface for risk assessment, processes user responses in real time, and can be deployed locally on devices to enhance data privacy and security. Figure 1 shows a comparison of our work to traditional methods.

Figure 1. Comparison between large language model (LLM)–based conversational AI (Conv-AI) and traditional machine learning methods for disease risk assessment. The Conv-AI leverages pretrained models that require only very few-shot fine-tuning, can handle unstructured textual data, provide real-time feature importance for each risk assessment it provides, and offer transferability with zero to very few shots for new risk assessment tasks. In contrast, traditional machine learning methods require large datasets for de novo training, process structured data, rely on extra computational steps for instance-specific post hoc feature importance (eg, Shapley additive explanations), and need retraining for each new task.



Contributions

Our contributions to the field of LLM-based disease risk assessment are diverse. First and foremost, we transition from machine traditional learning-based health outcome prediction-which typically relies on structured tabular data-to chatbot-based, no-code prediction using streaming QA interactions. This is realized through the development of a GenAI-powered mobile app that integrates fine-tuned LLMs as the core for personalized risk assessment and patient-clinician communication. The app not only assesses disease risk for patients but also provides contextual insights related to risk surveillance and mitigation through natural language conversation.

Second, we demonstrate that generative LLMs can outperform traditional machine learning methods, such as logistic regression [25], random forest [26], and extreme gradient boosting (XGBoost) [27], in low-data regimes, which is critical for medical applications where labeled data are scarce. For instance, our results show that LLMs like the T0-3b model achieve an AUC of 0.75 in zero-shot settings, demonstrating their potential for disease risk assessment even without task-specific training.

```
https://ai.jmir.org/2025/1/e67363
```

In addition, we provide a comprehensive comparison of both decoder-only and encoder-decoder models, fine-tuned using the widely adopted, parameter-efficient, low-rank adaptation (LoRA) method [28].

Third, we introduce a feature importance analysis derived from the LLM's attention layers, providing personalized insights into the most influential factors driving the model's predictions. This enhances the interpretability and usability of the risk assessment for both patients and clinicians, offering real-time, instance-specific explanations during inference.

Methods

Our Research Objective

The primary objective of this study is to explore the effectiveness of pretrained generative LLMs in no-code risk assessment of disease severity using few-shot multihop QA interactions. We aim to evaluate how these generative LLM-powered chatbots can use streaming QA interactions to accurately classify patient outcomes as severe or nonsevere, which is crucial for early risk assessment and optimizing health

care resource allocation. Through a case study of COVID-19 severity risk assessment, we developed an app that uses open-source generative LLMs to determine the severity of COVID-19 outcomes. This involves leveraging the models' capabilities in zero-shot and few-shot settings, with a focus on the use of serialization techniques to enhance their effectiveness and generalizability. We also integrate real-time feature importance to provide interpretable risk assessments. Figure 2 shows the workflow of our approach, from fine-tuning generative LLMs using serialized QA pairs to real-time risk assessment through a conversational interface.

Figure 2. Workflow for few-shot COVID-19 severity risk assessment using generative large language models (LLMs) with different serialization techniques. The top section, labeled "Backend - system developer," shows the fine-tuning phase where a few-shot sample of patient data, serialized through list and text templates, is used to fine-tune the LLMs. This backend process includes the creation of prompts and corresponding labels for model fine-tuning. The bottom section, labeled "Frontend - user," illustrates how a conversational chatbot interacts with users through our application to gather responses through streaming QA interactions. These responses are analyzed by the fine-tuned LLM in real time, providing risk assessments and highlighting the top attributing features that explain the model's risk assessment. QA: question-and-answer.



Data Collection

A dataset was collected from the emergency departments of Children's Hospital of Michigan and UPMC Children's Hospital of Pittsburgh between March 2021 and February 2022. Table 1 provides an overview of the binary features used in our study, including demographic, clinical, and social determinants that may influence COVID-19 severity risk. The dataset includes a total of 393 participant records, each characterized by responses to a series of carefully designed questions (see Figure 3 for sample QA pairs).

The severity of illness was defined based on the presence of any of the following criteria:

1. Requirement for supplemental oxygen (≥50% fraction of inspired oxygen)

- 2. Need for mechanical ventilation or noninvasive positive pressure ventilation (bilevel positive airway pressure and continuous positive airway pressure)
- 3. Need for vasopressors or inotropes
- 4. Requirement for extracorporeal membrane oxygenation
- 5. Cardiopulmonary resuscitation
- 6. Death from a related cause within 4 weeks after discharge

Children meeting any of these criteria were categorized as having severe illness. These outcomes were determined through chart reviews and parent surveys conducted 30 days after discharge [29].

Outliers were removed, and feature selection was performed using Shapley additive explanations values [30], resulting in the final dataset used for analysis.



Table 1. Binary features used in the study. The dataset consists of 393 patient records with 15 features representing demographics, clinical symptoms, and social determinants. These features serve as inputs for traditional machine learning models and are also serialized for fine-tuning generative large language models (LLMs).

Feature and label	Count, n
f1. Ages 5 to 11 years	
No	294
Yes	99
f2. Gender	
Female	332
Male	61
f3. Hispanic	
No	359
Yes	34
f4. African American	
No	215
Yes	178
f5. Service at stores	
Good	335
Poor	78
f6. Insurance	
No	387
Yes	6
f7. Headache	
No	332
Yes	61
f8. Fever	
No	211
Yes	182
f9. Cough	
No	210
Yes	183
f10. Shortness of breath	
No	292
Yes	101
f11. Exposed to COVID-19 individuals	
No	343
Yes	50
f12. Nausea or vomiting	
No	272
Yes	121
f13. Lungs check	
Bad	317
Good	76
f14. Eye redness	
No	381

https://ai.jmir.org/2025/1/e67363

Feature and label	Count, n
Yes	12
f15. COVID-19 antibody test	
Negative	364
Positive	29
f16. Outcome (severity)	
No	284
Yes	109

Figure 3. Overview of our mobile app design, showcasing patient data collection, real-time risk assessment using large language models (LLMs), and clinician review interface.



Tabular Data for Traditional Models

As traditional machine learning methods require tabular data as input, we formalize the questionnaire QA pairs \square , where n=393, \square represents the binary feature vector of the *i*-th instance where d=15, and \square denotes the binary class label indicating the presence or absence of severe COVID-19 symptoms determined by clinicians.

Each feature vector x_i consists of binary indicators representing social determinants and clinical and demographic factors that may influence the severity of COVID-19, such as age, preexisting conditions, vital signs, and laboratory test results. These features are shown in Table 1. The feature names are

```
https://ai.jmir.org/2025/1/e67363
```

RenderX

denoted as \square , where each f_j is a natural-language string describing the corresponding attribute.

The task is to predict the binary outcome y_i based on the information provided in x_i . This constitutes a supervised learning problem where the objective is to train a model to minimize prediction error on unseen data.

Serialization for New Conversational AI

At the time of data collection from 2021 to 2022, we did not yet have a chatbot for automated data donations from users, so we used a questionnaire to collect answers from each patient based on a set of questions designed for this study. As a result, the native format of the dataset consists of QA pairs, which were subsequently serialized to fine-tune the generative LLMs

for the risk assessment task. It is important to note that the fine-tuned model is capable of assessing risk using streaming QA interactions in real time (Figures 2 and 3).

To achieve serialization, the features in our dataset are denoted

as $\boxed{|\mathbf{x}|}$, and their associated values as $\boxed{|\mathbf{x}|}$. This notation provides a structure that is transformed into natural language prompts for the LLM.

We used two main serialization methods from TABLLM [31], the list template and the text template, to create natural language representations of the data. As shown in Figure 2, the list template links each feature with its value using an equal sign ("="), while the text template uses a narrative structure with the word "is" to connect each feature with its value. These templates enable us to evaluate which serialization approach better translates the data into actionable insights by the LLM.

Generative LLMs

We explore the capabilities of 3 white-box LLMs—LLaMA2 [23], T0 [32], and Flan-T5 [33]—focusing on their application in risk prediction for COVID-19 using both the native QA pairs and the formatted tabular dataset.

To our knowledge, this is one of the first attempts leveraging generative LLMs and conversational data science for disease risk assessment across various LLMs and few-shot settings. Our selection includes both decoder-only (LLaMA2) and encoder-decoder architectures (T0 and Flan-T5), allowing for a comprehensive assessment and comparison of their performance. The white-box nature of these models is particularly advantageous as it enables setup on local hosts with private datasets, ensuring precise risk assessment by allowing direct access to model weights and logits.

The input to the LLMs is a serialized string generated from the tabular data using the previously explained serialization

strategies. Given a feature vector 🗵. and their associated values

 \square , the serialized input string S_i can be represented using either the list template or text template serialization methods (Figure 2).

These feature vectors originate from the structured dataset described in Table 1, which provides the foundation for both traditional and generative model comparisons.

The LLM processes the serialized input string S_i and outputs logits for the next token in the sequence. We focus on the logits corresponding to the tokens "yes" and "no," which indicate severe or nonsevere symptoms, respectively. The probabilities for these tokens are obtained by applying the softmax function to the logits:



The probability $|\mathbf{x}|$ indicates the likelihood of severe symptoms based on the input data S_i . This probability is directly used as the severity risk score for evaluation purposes.

To determine the binary predicted label 🗵 from this probability:

```
https://ai.jmir.org/2025/1/e67363
```

×

The probability score [k], reflecting the severity risk, is used to compute the AUC for evaluation (Figure 2).

Evaluation Setting

Zero-Shot Setting

In the zero-shot setting, our approach leverages the intrinsic capabilities of LLMs. These models, unlike traditional classifiers such as logistic regression and XGBoost, have been extensively pretrained on diverse datasets. This extensive pretraining enables them to apply their accumulated world knowledge directly to specific classification tasks without additional training, demonstrating exceptional generalizability.

We assess the zero-shot prediction effectiveness of these LLMs by presenting them with tasks aligned with our study's objectives that they have not been specifically trained on. The models interpret and classify new, unseen data solely based on their pretrained knowledge. This approach not only highlights the potential of LLMs in real-world applications but also evaluates their ability to generalize from their training to novel scenarios in healthcare.

This zero-shot methodology allows us to evaluate how well these LLMs can recognize and classify complex, previously unseen patterns in health care data, providing valuable insights into their practical applicability and limitations in clinical settings.

Few-Shot Fine-Tuning

In the few-shot setting, we use sample sizes of 2, 4, 8, 16, and 32 to fine-tune the LLMs, aiming to examine the effect of training sample size on model performance compared to traditional classifiers. To ensure fairness and reduce bias in the fine-tuning process, we maintain a balanced ratio of positive

 \square and negative \square samples, with an equal number of examples from each class in each sample size.

To enhance computational efficiency in adapting the LLMs to our specific tasks, we employ a parameter-efficient fine-tuning approach using LoRA [27]. Instead of adjusting all parameters within the model, LoRA involves training a small proportion of parameters by integrating trainable low-rank matrices into each layer of the pretrained model. This method allows the model to quickly adapt to new tasks by optimizing only a subset of parameters, thereby preserving the general capabilities of the LLM while enhancing its performance on task-specific features.

Feature Importance Analysis

In disease risk assessment, interpretability is as critical as accuracy, particularly when both are provided to the user in real time. Here, we introduce a novel approach for analyzing feature importance by leveraging the attention mechanisms inherent in the output layers of generative LLMs. This method provides additional insights into the risk assessment process of the model, which is valuable for both clinicians and patients in understanding the factors contributing to the model's output.

Our approach involves extracting attention scores from the model's output layer, where the attention assigned to each input token is interpreted as an indicator of feature importance. We compute the attention for each feature-value pair and associate the average attention score with the corresponding feature. This provides a holistic view of which features, along with their associated values, influence the model's output.

In Figure 4, the attention map illustrates the attention scores for a predicted positive case by the LLM, where darker shades represent higher attention scores assigned to specific feature-value pairs.

For an input sequence such as:

A patient 🗵

Do the descriptions of this patient show severe symptoms of COVID-19? Yes or no? Result:

We calculate attention scores for each feature-value pair in the original sequence. The average attention score for each feature-value pair is then computed, and the score is associated with the feature itself, offering a representation of feature importance in the context of disease severity risk. As shown in Figure 4, any missing data in both the training and inference

stages could be handled by having the value as "none" and having the model make the prediction; this will impact the prediction depending on the feature missing, but the free-text input of the LLMs still allows for a prediction to happen.

This normalized attention score serves as a proxy for feature importance, offering clinicians and patients a clearer understanding of which features (eg, age, preexisting conditions, vital signs, etc) are most influential in the model's assessment of COVID-19 severity risk. As illustrated in Multimedia Appendix 1, the plot shows the normalized attention scores from the LLaMA2-7b model in the 32-shot setting for two test cases: one positive (yes) and one negative (no).

For the positive case, the top five features with the highest attention scores, as shown in this figure, are:

- 1. f15: COVID-19 antibody test
- 2. f13: Lungs check
- 3. f12: Nausea or vomiting
- 4. f9: Cough
- 5. f14: Eye redness

By integrating this analysis into our mobile app, we enhance the interpretability of LLM-based risk assessments, empowering users with deeper insights into the model's reasoning process.

Figure 4. The attention map for a predicted positive case where the darker color represents larger attention weights for each token. The prompts are tokenized to mimic the actual inputs to the large language models (LLMs).

A _patient _age _ 5 _to _ 1 1 = no , _gender = male , _His pan ic = no , _African - American = yes , _service _at _stores = po or , _ins urance = no , _head ache = no , _fe ver = no , _c ough = no , _short ness _of _breath = no , _exposed _to _COVID _individuals = yes , _n ause a _or _vom iting = no , _l ungs _check = bad , _eye - red ness = no , _COVID 1 9 _ant ib ody _test = negative <0x0A> Does _the _descri ptions _of _this _patient _show _severe _sympt oms _of _COVID 1 9 ? _yes _or _no ? _Result :

Mobile App

To provide users with code-free disease severity risk assessment and enhance user experience, we developed a mobile chatbot powered by the aforementioned generative LLMs. This app is designed to facilitate the assessment and management of COVID-19 in children, with potential applicability to other diseases and conditions. It offers two versions: one for patients to donate their health information via answering the questions and receiving real-time severity risk assessments, and another for clinicians to manage, review, and interpret the sessions donated by patients. The primary goals are to enhance early detection of severe outcomes, improve patient-clinician communication, and streamline the overall risk assessment process.

The app targets patients, clinicians, and other health care providers involved in managing preclinical cases. It leverages the capabilities of generative LLMs to analyze patient responses and provide immediate feedback on the risk of severe symptoms. Developed using React Native and JavaScript for the front end, Firebase for database management, and various frontend technologies, the app provides a user-friendly, efficient, and effective solution for managing disease risks. It aims to improve patient outcomes by facilitating timely and informed decision-making.

Database Structure

Our mobile app uses Firebase for database management, structured into three primary collections: Users, Questions, and Answers.

The data flow between the patient, LLM backend, Firebase, and interfaces for both patients and clinicians is illustrated in Figure 5. This figure highlights the interactions among processes, including the assessment submission, session management, and result retrieval.

- Users: This collection includes essential user information such as ID, Email, and isAdmin. The ID uniquely identifies each user, the Email serves as contact information, and the isAdmin field (Boolean) indicates whether the user has administrative privileges (clinicians) or not (patients).
- Questions: Each document in this collection has a unique ID and a Description field. The ID is used to reference questions in the Answers collection, and the Description contains the text of the question posed to the user, ensuring clarity and specificity in data mapping.
- Answers: This collection records user responses during their sessions. Each document includes a session ID and an array of answers where each entry links to the relevant Question ID from the Questions collection. In addition, it contains a Text field for the user's detailed response; an Answer field for the LLM-generated response (eg, yes or

```
SI•FO
```

no); a Date field marking the session's completion time; a Risk Score field, which is derived from the user's responses and utilized for subsequent risk prediction by the LLM; and

an Important Features field, which stores the key features identified by the LLM's attention scores that contributed to the risk assessment.

Figure 5. Data flow diagram where we map out the flow of information between different processes of large language model (LLM) backend, Firebase, and mobile app interfaces for both patient and clinician.



Detailed session result

User Interface: Assessment

The step-by-step workflow for conducting an assessment and storing results in Firebase is detailed in Figure 6. This sequence diagram outlines the interaction between the patient, mobile app, LLM backend, and database.

As shown in Figure 3, on the Assessment page, we leverage the power of LLMs to engage in a conversation with the patient. This interaction allows us to ask questions and gather contextual information for each response. By doing so, we retrieve a binary answer (yes or no) using the LLM, which is then provided to the primary care physician along with the patient's context to aid in decision-making.

After the user responds to each question, we use our LLM to generate a binary answer. This involves providing the LLM

with instructions that include the question and the user's response and asking the LLM to interpret the response into a binary answer (yes or no). This sequential process is performed for all questions. Currently, the input for the final LLM-based risk assessment, which predicts the COVID-19 severity risk, is based solely on the set of binary answers generated by the LLM. Future enhancements could incorporate the original user responses to improve context understanding.

We currently use the Llama2-7b application programming interface (API) for answer retrieval. Our long-term goal is to integrate a fine-tuned LLM hosted on our servers to ensure better optimization and accuracy specific to our dataset, as evidenced by the improved performance results discussed in this paper.



Figure 6. Sequence diagram for the Assessment page, where the patient takes the risk assessment and the large language model (LLM) backend calculates the results, which will be saved to the Firebase. QA: question-and-answer.



User Interface: Patient and Clinician Results

Figure 7 illustrates the interaction flows for both patients and clinicians as they access session details and results. This sequence diagram shows how patient data and assessments are retrieved and displayed in real time.

Patients can submit a session at any time, receiving an immediate risk assessment in the Patient Interface section (Figure 3). This section displays all sessions submitted by the current user, along with their respective risk assessments.

In the Clinician Interface section, clinicians can access all sessions from their patients, organized by patient ID, for efficient

review. Each session includes a comprehensive report featuring the predicted risk score, ensuring transparency and aiding in clinical decision-making.

Upon submission, a patient's session is instantly available in both the patient's and clinician's panels. While patients can only view their own sessions, clinicians can review all sessions from their assigned patients. This setup supports real-time updates through Firebase, facilitating seamless communication and follow-up between patients and their health care providers. Furthermore, the app provides personalized feature importance analysis based on the LLM's attention layers, giving both patients and clinicians additional insights into the most critical factors influencing the risk assessment.



Mobile app Firebase Clinician Patient Open patient interface Access all sessions All sessions List of all sessions Tap one session Answers and results Open clinician interface Access all patients data All patients data List of all patients' sessions Tap one session Answers and results

Figure 7. Sequence diagram for displaying patients' session results. As shown, each patient has access to all their own sessions while the clinician can access all patients' sessions.

Ethical Considerations

The data collected and used for this study were approved by the University of Pittsburgh Institutional Review Board (MOD21010046-003; approval date: February 25, 2021). Informed consent was obtained from all legal caregivers, and when age appropriate, an informed assent was also obtained from the participants. Before the use of this study, the data were subject to a multistep anonymization procedure with personally identifying information marked and deleted.

Results

Training and Fine-Tuning Settings

In our experiments, we used a rigorous hyperparameter tuning strategy to optimize model performance, supported by a robust setup to ensure diverse dataset initialization and minimize potential biases. For both traditional machine learning methods and LLMs, we used 5 specific random seeds—0, 1, 32, 42, and 1024—to create diverse dataset splits. The dataset of 393 samples was divided into 256 training, 59 validation, and 78

```
https://ai.jmir.org/2025/1/e67363
```

RenderX

testing segments, preserving a consistent positive-to-negative ratio of approximately 0.38.

For both traditional methods and LLMs, training was conducted using up to 32 shots to evaluate performance in the few-shot regime. For few-shot settings ranging from 2 to 32 shots, we ensured a balanced sampling of positive and negative examples in the training set, maintaining an equal number of instances from each class to avoid biases during training. Key hyperparameters, such as the learning rate, were optimized using grid search, with the learning rate set to 3×10^{-4} . The batch size matched the number of shots, and training consistently ran for 128 epochs to ensure convergence. During fine-tuning with LoRA, validation loss was monitored to select the best model checkpoint, minimizing overfitting and enhancing generalization to the test set. The optimization used cross-entropy loss, aligning with the binary classification task of predicting COVID-19 severity. This comprehensive setup ensured robust and interpretable model performance, particularly in low-data settings.

Overview

Table 2 shows the performance of different serialization methods for the LLMs across various few-shot settings. We evaluated 2 primary serialization methods: list template and text template, across models tested with 0, 2, 4, 8, 16, and 32 training shots to observe performance variations with the number of training examples.

examples. The list template often exhibited better performance at lower

The list template often exhibited better performance at lower shot counts, while the text template typically outperformed the list template as the number of training examples increased. The following summarizes the performance trends for each model.

Table 2. Performance of models across different shot settings. All values represent the average area under the curve (AUC) across 5 random seeds rounded to 2 decimal places. In addition, SDs given across the 5 random seeds are shown. The suffixes "-L" and "-T" represent list serialization and text serialization, respectively.

Model	Number of shots					
	0, AUC ^a (SD)	2, AUC (SD)	4, AUC (SD)	8, AUC (SD)	16, AUC (SD)	32, AUC (SD)
Llama2-7b-L	0.54 (.05)	0.69 (.07)	0.69 (.06)	0.68 (.04)	0.63 (.04)	0.66 (.07)
Flan-t5-xl-L	0.62 (.03)	0.64 (.02)	0.63 (.02)	0.68 (.06)	0.66 (.05)	0.69 (.06)
Flan-t5-xxl-L	0.60 (.03)	0.61 (.03)	0.61 (.05)	0.62 (.06)	0.59 (.10)	0.65 (.11)
T0pp(8bit)-L	0.69 (.04)	0.70 (.07)	0.70 (.05)	0.70 (.05)	0.68 (.06)	0.70 (.10)
T0-3b-L	0.68 (.04)	0.67 (.04)	0.68 (.05)	0.70 (.04)	0.67 (.04)	0.67 (.07)
Llama2-7b-T	0.59 (.05)	0.69 (.03)	0.69 (.01)	0.64 (.07)	0.63 (.05)	0.67 (.06)
Flan-t5-xl-T	0.69 (.03)	0.69 (.02)	0.69 (.03)	0.71 (.05)	0.69 (.04)	0.70 (.05)
Flan-t5-xxl-T	0.61 (.04)	0.58 (.03)	0.63 (.08)	0.59 (.10)	0.62 (.09)	0.63 (.10)
T0pp(8bit)-T	0.67 (.02)	0.65 (.05)	0.66 (.05)	0.68 (.04)	0.65 (.08)	0.67 (.08)
Т0-3b-Т	0.75 (.04)	0.65 (.06)	0.65 (.05)	0.68 (.03)	0.67 (.04)	0.65 (.08)
Logistic regression	b	0.57 (.07)	0.55 (.10)	0.64 (.06)	0.61 (.11)	0.69 (.08)
Random forest	—	0.57 (.07)	0.57 (.06)	0.62 (.08)	0.66 (.07)	0.68 (.07)
XGBoost ^c	_	0.50 (.00)	0.50 (.00)	0.50 (.00)	0.54 (.06)	0.65 (.03)

^aAverage area under the curve.

^bNot applicable.

^cXGBoost: extreme gradient boosting.

Llama2-7b

In the zero-shot setting, the text template achieved an AUC of 0.59 compared to 0.54 for the list template. At 2 training shots, both templates achieved an AUC of 0.69, but the text template began to outperform, reaching an AUC of 0.67 at 32 training shots compared with 0.66 for the list template.

Flan-t5-xl

The text template consistently outperformed the list template across most shot settings. At 2 training shots, the text template achieved an AUC of 0.69 compared to 0.64 for the list template, and this lead continued up to 32 shots, where the text template achieved an AUC of 0.70 compared to 0.69 for the list template.

Flan-t5-xxl

Both templates showed similar performance in the early few-shot settings. At 2 training shots, the list template achieved an AUC of 0.61, slightly outperforming the text template, which achieved an AUC of 0.58. By 32 training shots, the list template achieved an AUC of 0.65, slightly outperforming the text template, which achieved an AUC of 0.65.

T0pp (8bit)

In the zero-shot setting, the list template led with an AUC of 0.69 compared to 0.67 for the text template. This lead was maintained through most shot settings, with both templates achieving around 0.70 AUC by 32 shots.

T0-3b

The text template outperformed the list template in the zero-shot setting, achieving an AUC of 0.75 compared to 0.68 for the list template. In the 2-shot setting, the list template performed slightly better, with an AUC of 0.67 compared to 0.65 for the text template. At 32 shots, the text template closed the gap with an AUC of 0.65 compared with 0.67 for the list template.

In Table 3, we can also compare the best-performing models across different shots, constraining the recall to be higher than 0.8. This gives us better insights into their performance in population screening for early health risks, where recall is considered more important than precision.

Overall, while the list template often provides an initial advantage in early few-shot settings, the text template shows competitive performance as the number of training examples increases. This suggests that serialization choice can be



important in low-data regimes. The text template's strong model, highlights its potential when no training data is available. performance in the zero-shot setting, particularly for the T0-3b

Shot	Best model	Threshold	Precision	Recall	<i>F</i> ₁ -score
0	T0-3b	0.04	0.37	0.85	0.52
2	T0pp	0.12	0.34	0.81	0.46
4	T0pp	0.24	0.35	0.83	0.49
8	flan-t5-xl	0.17	0.38	0.80	0.50
16	flan-t5-xl	0.15	0.34	0.85	0.48
32	flan-t5-xl	0.16	0.36	0.81	0.49

Table 3. Precision, recall, and F1-score of the best performing models across different shots averaged over 5 random seeds.

LLMs Versus Traditional Machine Learning Methods

Our study highlights the versatility of LLMs for various health care apps, particularly in scenarios with limited data. To benchmark their performance against traditional machine learning methods, we compared LLMs with logistic regression, random forest, and XGBoost.

LLMs benefit from extensive pretraining, allowing them to generalize well to "unseen" data, unlike traditional methods that require substantial amounts of training data. As shown in Table 2, LLMs like T0-3b-T achieved an AUC of 0.75 in the zero-shot setting, demonstrating a good performance even without task-specific fine-tuning. This demonstrates the effectiveness of LLM-powered risk assessment without the need for additional labeled data.

In the 2-shot setting, LLMs continue to show strong performance relative to traditional methods. For instance, Figure 8 compares the average AUC across 5 different seeds in this scenario. The left panel shows results using the list serialization (-L) approach, while the right panel shows results using the text serialization (-T) approach. In this 2-shot scenario, LLMs such as

T0pp(8bit)-L and Flan-t5-xl-T achieve AUCs of 0.70 and 0.69, respectively, clearly outperforming traditional methods, including logistic regression, random forest, and XGBoost, which achieved AUCs of 0.57, 0.57, and 0.50, respectively.

LLMs' ability to perform well with minimal data highlights their advantage in low-data regimes. This makes them particularly suitable for real-time, no-code health care apps where rapid decision-making is required, even in scenarios where labeled data is scarce.

Furthermore, LLMs' capacity to handle streaming data formats, such as multihop QA pairs, enhances their integration into conversational interfaces, supporting real-time patient-clinician interactions. This flexibility offers significant usability in clinical settings where personalized and immediate risk assessments are needed (Figure 1).

Overall, while traditional methods may improve with larger datasets, LLMs demonstrate a clear advantage in dynamic, low-data health care environments. Their ability to handle incomplete data and streaming input formats makes them robust for real-world applications requiring adaptability and speed.

Figure 8. Average area under the curve (AUC) in a 2-shot setting over 5 different seeds. The left panel shows results using the list serialization (-L) approach, while the right panel shows results using the text serialization (-T) approach. XGBoost: extreme gradient boosting.



False positive rate



which is particularly effective in low-data regimes. These models excel in zero-shot and few-shot settings, showcasing their ability to perform well without extensive domain-specific training. This is crucial for real-time applications requiring immediate and reliable predictions, highlighting their

Discussion

Principal Findings

Our research demonstrates that generative LLMs provide a robust and no-code approach for predicting COVID-19 severity,

```
https://ai.jmir.org/2025/1/e67363
```

exceptional generalizability compared with traditional classifiers like logistic regression, random forest, and XGBoost, which typically require more labeled data to achieve comparable performance.

Generative LLMs effectively handle diverse input formats, integrating both structured clinical data and unstructured natural language inputs from patient interactions. This flexibility enables them to synthesize information from various sources, such as patient medical histories and symptom descriptions, enhancing their usability in dynamic health care settings. In our study, we incorporated these models into a conversational interface, which facilitates real-time patient-clinician interactions and immediate risk assessments. This setup supports continuous data collection and leverages the conversational capabilities of LLMs to optimize clinical decision-making and resource allocation.

Future Directions and Limitations

Future work should focus on integrating continuous clinician-patient conversational data for fine-tuning or in-context learning, extending the application of LLMs beyond static disease prediction models. Techniques like chain of thought and chain of interaction, which align with the interactive nature of medical consultations, show promise for enhancing model performance in interpreting and responding to patient data in real-time settings. While our study used models like T0pp with parameter-efficient fine-tuning using LoRA, future research could explore newer and more advanced small language models such as LLaMA3-8b and Mistral-7b-Instruct, which have demonstrated exceptional performance in low-data regimes. These models could offer greater efficiency and accuracy as computational resources and methodologies advance, supporting more sophisticated and scalable applications in health care [34,35].

However, limitations remain that warrant further exploration. This study does not address the critical issue of handling sensitive data, such as personally identifiable information (PII), within health care datasets. Incorporating a dual dataset that includes both PII and non-PII data could facilitate machine unlearning research, allowing models to selectively forget sensitive information while retaining predictive capabilities from nonsensitive data. This would ensure compliance with privacy regulations and enhance the ethical deployment of LLMs in health care. Advancing privacy-preserving techniques, such as selective forgetting mechanisms, would not only safeguard sensitive data but also support broader trust in the use of LLMs in clinical settings.

As these models evolve, vulnerabilities such as adversarial attacks during in-context learning pose significant risks. Studies have shown that manipulated inputs can lead to inaccurate or harmful predictions, particularly in high-stakes tasks like health care risk assessment [36]. Addressing these risks is crucial to ensure that LLMs remain reliable and safe for broader adoption in health care applications. Enhanced resilience to adversarial techniques, combined with privacy-preserving methods, will be key to building robust and trustworthy systems. By addressing these challenges, future research can ensure that LLMs not only deliver accurate predictions but also adhere to ethical and privacy standards in real-world settings.

Conclusions

In conclusion, generative LLMs offer a valuable tool for no-code risk assessment in low-data regimes. Their ability to perform zero-shot or few-shot transferability to new diseases or conditions and handle complex, varied inputs positions them as key assets for enhancing health care interventions and resource management. Furthermore, the incorporation of feature importance analysis derived from the LLM's attention layers provides an additional layer of interpretability, offering personalized insights into the decision-making process for both patients and clinicians.

Acknowledgments

Research reported in this publication was supported by the Eunice Kennedy Shriver Institute of Child Health and Human Development of the National Institute of Health under awards R61HD105610 and R33HD105610.

Authors' Contributions

MAR conducted the experiments, designed the app, and wrote the manuscript. XZ contributed to the app design, assisted with the experiments, and provided revisions. DZ designed, oversaw, and supported the project. YQ offered suggestions on the experiments and revisions. SS, SH, and US assisted with dataset collection and provided feedback on the manuscript draft.

Conflicts of Interest

SDH is named as a co-inventor on a patent for the diagnostic use of salivary RNA in neurologic disorders. He previously served as a scientific advisory board member for Quadrant Biosciences and Spectrum Solutions. No other conflicts of interest are declared by other authors.

Multimedia Appendix 1

Normalized attention scores from LLaMA2-7b in the 32-shot setting, showing feature importance for 2 test cases, 1 positive (yes) and 1 negative (no), simultaneously with the risk assessment. [PNG File, 77 KB - ai v4i1e67363 app1.png]



References

- Li X, Zhu D, Levy P. Leveraging auxiliary measures: a deep multi-task neural network for predictive modeling in clinical research. BMC Med Inform Decis Mak 2018;18(Suppl 4):126 [FREE Full text] [doi: 10.1186/s12911-018-0676-9] [Medline: 30537954]
- Wang L, Dong M, Towner E, Zhu D. Prioritization of multi-level risk factors for obesity. 2019 Presented at: Proceedings of the IEEE International Conference on Bioinformatics and Biomedicine (BIBM); November 18-21, 2019; San Diego, CA. [doi: 10.1109/bibm47256.2019.8982940]
- 3. Li X, Zhu D, Levy P. Predicting clinical outcomes with patient stratification via deep mixture neural networks. AMIA Jt Summits Transl Sci Proc 2020;2020:367-376 [FREE Full text] [Medline: <u>32477657</u>]
- Devlin J, Chang M, Lee K, Toutanova K. BERT: pre-training of deep bidirectional transformers for language understanding. 2019 Presented at: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers); June 2-7, 2019; Minneapolis, MN p. 4171-4186. [doi: 10.18653/v1/N19-1423]
- 5. Huang K, Altosaar J, Ranganath R. ClinicalBERT: modeling clinical notes and predicting hospital readmission. ArXiv Preprint posted online on April 10, 2019. [doi: <u>10.48550/arXiv.1904.05342</u>]
- Alsentzer E, Murphy J, Boag W, Weng W, Jin D, Naumann T, et al. Publicly available clinical BERT embeddings. : Association for Computational Linguistics; 2019 Presented at: Proceedings of the 2nd Clinical Natural Language Processing Workshop; June 7, 2019; Minneapolis, MN p. 72-78. [doi: <u>10.18653/v1/w19-1909</u>]
- Rasmy L, Xiang Y, Xie Z, Tao C, Zhi D. Med-BERT: pretrained contextualized embeddings on large-scale structured electronic health records for disease prediction. NPJ Digit Med 2021;4(1):86 [FREE Full text] [doi: 10.1038/s41746-021-00455-y] [Medline: 34017034]
- 8. Brown TB, Mann B, Ryder N, Subbiah M, Kaplan JD, Dhariwal P, et al. Language models are few-shot learners. In: Adv Neural Inf Process Syst 2020. 2020 Presented at: NIPS '20: Proceedings of the 34th International Conference on Neural Information Processing Systems; December 6-12, 2020; Vancouver, BC p. 1877-1901.
- 9. Achiam J, Adler S, Agarwal S, Ahmad L, Akkaya I, Aleman FL, et al. GPT-4 technical report. ArXiv Preprint posted online on March 15, 2023. [doi: 10.48550/arXiv.2303.08774]
- Gilson A, Safranek CW, Huang T, Socrates V, Chi L, Taylor RA, et al. How does chatGPT perform on the United States medical licensing examination (USMLE)? The implications of large language models for medical education and knowledge assessment. JMIR Med Educ 2023;9:e45312 [FREE Full text] [doi: 10.2196/45312] [Medline: 36753318]
- Kung TH, Cheatham M, Medenilla A, Sillos C, de Leon L, Elepaño C, et al. Performance of chatGPT on USMLE: potential for AI-assisted medical education using large language models. PLOS Digit Health 2023;2(2):e0000198. [doi: <u>10.1371/journal.pdig.0000198</u>] [Medline: <u>36812645</u>]
- 12. Ali SR, Dobbs TD, Hutchings HA, Whitaker IS. Using chatGPT to write patient clinic letters. Lancet Digit Health 2023;5(4):e179-e181 [FREE Full text] [doi: 10.1016/S2589-7500(23)00048-1] [Medline: 36894409]
- Xv Y, Peng C, Wei Z, Liao F, Xiao M. Can Chat-GPT a substitute for urological resident physician in diagnosing diseases?: a preliminary conclusion from an exploratory investigation. World J Urol 2023;41(9):2569-2571. [doi: <u>10.1007/s00345-023-04539-0]</u> [Medline: <u>37505265</u>]
- 14. Kanjee Z, Crowe B, Rodman A. Accuracy of a generative artificial intelligence model in a complex diagnostic challenge. JAMA 2023;330(1):78-80 [FREE Full text] [doi: 10.1001/jama.2023.8288] [Medline: 37318797]
- Tang L, Sun Z, Idnay B, Nestor JG, Soroush A, Elias PA, et al. Evaluating large language models on medical evidence summarization. NPJ Digit Med 2023;6(1):158 [FREE Full text] [doi: 10.1038/s41746-023-00896-7] [Medline: <u>37620423</u>]
- Sallam M. The utility of chatgpt as an example of large language models in healthcare education, research and practice: Systematic review on the future perspectives and potential limitations. MedRxiv Preprint posted online on February 21, 2023. [doi: <u>10.1101/2023.02.19.23286155</u>]
- 17. Kim Y, Xu X, McDuff D, Breazeal C, Park HW. Health-LLM: Large language models for health prediction via wearable sensor data. ArXiv Preprint posted online on January 12, 2024. [doi: <u>10.48550/arXiv.2401.06866</u>]
- Yang Z, Mitra A, Kwon S, Yu H. Clinicalmamba: a generative clinical language model on longitudinal clinical notes. ArXiv Preprint posted online on March 9, 2024 [FREE Full text] [doi: 10.18653/v1/2024.clinicalnlp-1.5]
- 19. Nori H, King N, McKinney S, Carignan D, Horvitz E. Capabilities of GPT-4 on medical challenge problems. ArXiv Preprint posted online on March 20, 2023. [doi: <u>10.48550/arXiv.2303.13375</u>]
- 20. Singhal K, Tu T, Gottweis J, Sayres R, Wulczyn E, Amin M, et al. Toward expert-level medical question answering with large language models. Nat Med 2025. [doi: 10.1038/s41591-024-03423-7] [Medline: 39779926]
- Wu C, Lin W, Zhang X, Zhang Y, Xie W, Wang Y. PMC-LLaMA: toward building open-source language models for medicine. J Am Med Inform Assoc 2024;31(9):1833-1843. [doi: <u>10.1093/jamia/ocae045</u>] [Medline: <u>38613821</u>]
- 22. Ben Shoham O, Rappoport N. CPLLM: clinical prediction with large language models. PLOS Digit Health 2024;3(12):e0000680. [doi: 10.1371/journal.pdig.0000680] [Medline: 39642102]
- 23. Touvron H, Martin L, Stone K, Albert P, Almahairi A, Babaei Y, et al. LLAMA 2: open foundation and fine-tuned chat models. ArXiv Preprint posted online on July 18, 2023. [doi: <u>10.48550/arXiv.2307.09288</u>]

- 24. Bolton E, Venigalla A, Yasunaga M, Hall D, Xiong B, Lee T, et al. BioMedLM: A 2.7B Parameter Language Model Trained On Biomedical Text. ArXiv Preprint posted online on March 27, 2024. [doi: <u>10.48550/arXiv.2403.18421</u>]
- 25. Hosmer DJ, Lemeshow S, Sturdivant RX. Applied Logistic Regression. Hoboken, NJ: John Wiley & Sons; 2013.
- 26. Breiman L. Random forests. Mach Learn 2011;45(1):5-32 [FREE Full text] [doi: 10.1186/1478-7954-9-29] [Medline: 21816105]
- Chen T, Guestrin C. XGBoost: a scalable tree boosting system. 2016 Presented at: KDD '16: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining; August 13-17, 2016; San Francisco, CA p. 785-794. [doi: 10.1145/2939672.2939785]
- 28. Hu EJ, Shen Y, Wallis P, Allen-Zhu Z, Li Y, Wang S, et al. LoRA: low-rank adaptation of large language models. ArXiv Preprint posted online on June 17, 2021. [doi: <u>10.48550/arXiv.2106.09685</u>]
- Hicks SD, Zhu D, Sullivan R, Kannikeswaran N, Meert K, Chen W, et al. Saliva microRNA profile in children with and without severe SARS-CoV-2 infection. Int J Mol Sci 2023;24(9):8175 [FREE Full text] [doi: <u>10.3390/ijms24098175</u>] [Medline: <u>37175883</u>]
- 30. Lundberg SM, Lee SI. A unified approach to interpreting model predictions. ArXiv Preprint posted online on May 22, 2017. [doi: 10.48550/arXiv.1705.07874]
- 31. Hegselmann S, Buendia A, Lang H, Agrawal M, Jiang X, Sontag D. Tabllm: Few-shot classification of tabular data with large language models. ArXiv Preprint posted online on October 19, 2022. [doi: <u>10.48550/arXiv.2210.10723</u>]
- 32. Sanh V, Webson A, Raffel C, Bach S, Sutawika L, Alyafeai Z. Multitask prompted training enables zero-shot task generalization. ArXiv Preprint posted online on October 15, 2021. [doi: <u>10.48550/arXiv.2110.08207</u>]
- 33. Chung HW, Hou L, Longpre S, Zoph B, Tay Y, Fedus W. Scaling instruction-finetuned language models. J Mach Learn Res 2024;25(70):1-53 [FREE Full text]
- 34. Han G, Liu W, Huang X, Borsari B. Chain-of-interaction: enhancing large language models for psychiatric behavior understanding by dyadic contexts. 2024 Presented at: Proceedings of the IEEE 12th International Conference on Healthcare Informatics (ICHI); June 3-6, 2024; Orlando, FL p. 392-401. [doi: 10.1109/ichi61247.2024.00057]
- 35. Gramopadhye O, Nachane S, Chanda P, Ramakrishnan G, Jadhav K, Nandwani Y. Few shot chain-of-thought driven reasoning to prompt LLMs for open-ended medical question answering. 2024 Presented at: Findings of the Association for Computational Linguistics: EMNLP 2024; November 12-16, 2024; Miami, FL p. 542-573. [doi: <u>10.18653/v1/2024.findings-emnlp.31</u>]
- Qiang Y, Zhou X, Zhu D. Hijacking large language models via adversarial in-context learning. ArXiv Preprint posted online on November 16, 2023. [doi: <u>10.48550/arXiv.2311.09948</u>]

Abbreviations

AI: artificial intelligence API: application programming interface AUC: area under the curve BERT: Bidirectional Encoder Representations from Transformers EHR: electrical health record GenAI: generative artificial intelligence LLM: large language model LoRA: low-rank adaptation PII: personally identifiable information QA: question-and-answer SHAP: Shapley additive explanations USMLE: United States Medical Licensing Examination XGBoost: extreme gradient boosting

Edited by K El Emam; submitted 09.10.24; peer-reviewed by K Singh, B Srinivasaiah; comments to author 09.11.24; revised version received 27.11.24; accepted 23.02.25; published 27.03.25.

<u>Please cite as:</u> Roshani MA, Zhou X, Qiang Y, Suresh S, Hicks S, Sethuraman U, Zhu D Generative Large Language Model—Powered Conversational AI App for Personalized Risk Assessment: Case Study in COVID-19 JMIR AI 2025;4:e67363 URL: <u>https://ai.jmir.org/2025/1/e67363</u> doi:<u>10.2196/67363</u> PMID:



©Mohammad Amin Roshani, Xiangyu Zhou, Yao Qiang, Srinivasan Suresh, Steven Hicks, Usha Sethuraman, Dongxiao Zhu. Originally published in JMIR AI (https://ai.jmir.org), 27.03.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.

Insights on the Side Effects of Female Contraceptive Products From Online Drug Reviews: Natural Language Processing–Based Content Analysis

Nicole Groene¹, Dr; Audrey Nickel, MSc; Amanda E Rohn², MD

¹Department for Health and Social Sciences, FOM University of Applied Sciences for Economics and Management, Essen, Germany ²VHC Health, Arlington, VA, United States

Corresponding Author: Nicole Groene, Dr Department for Health and Social Sciences FOM University of Applied Sciences for Economics and Management Leimkugelstr 6 Essen, 45141 Germany Phone: 49 201 81004 Email: <u>Nicole.groene@fom-net.de</u>

Abstract

Background: Most online and social media discussions about birth control methods for women center on side effects, highlighting a demand for shared experiences with these products. Online user reviews and ratings of birth control products offer a largely untapped supplementary resource that could assist women and their partners in making informed contraception choices.

Objective: This study sought to analyze women's online ratings and reviews of various birth control methods, focusing on side effects linked to low product ratings.

Methods: Using natural language processing (NLP) for topic modeling and descriptive statistics, this study analyzes 19,506 unique reviews of female contraceptive products posted on the website Drugs.com.

Results: Ratings vary widely across contraception types. Hormonal contraceptives with high systemic absorption, such as progestin-only pills and extended-cycle pills, received more unfavorable reviews than other methods and women frequently described menstrual irregularities, continuous bleeding, and weight gain associated with their administration. Intrauterine devices were generally rated more positively, although about 1 in 10 users reported severe cramps and pain, which were linked to very poor ratings.

Conclusions: While exploratory, this study highlights the potential of NLP in analyzing extensive online reviews to reveal insights into women's experiences with contraceptives and the impact of side effects on their overall well-being. In addition to results from clinical studies, NLP-derived insights from online reviews can provide complementary information for women and health care providers, despite possible biases in online reviews. The findings suggest a need for further research to validate links between specific side effects, contraceptive methods, and women's overall well-being.

(JMIR AI 2025;4:e68809) doi:10.2196/68809

KEYWORDS

contraception; side effects; natural language processing; NLP; informed choices; online reviews; women; well-being

Introduction

Background

RenderX

According to the United Nations, contraception is a critical issue impacting 1.9 billion women of reproductive age. Worldwide, approximately 922 million women or their partners use contraception. More than half of all contracepting women rely

https://ai.jmir.org/2025/1/e68809

on modern contraceptive products designed to be used by women. These female products comprise long-acting reversible contraceptives (LARCs), such as intrauterine devices (IUDs) and hormonal implants as well as short-acting methods, such as oral contraceptives, known as "the pill," hormonal patches, vaginal rings, and contraceptive injections. Traditional methods such as withdrawal and calendar rhythm are relied upon by 7% of women, and the single most common contraceptive method

worldwide is female sterilization (24%), an irreversible method [1,2].

According to data from the latest National Survey of Family Growth (2017 to 2019), approximately 27.5% of women of reproductive age in the United States use female contraceptive products, comprising oral contraceptive pills (OCPs, 14%), LARCs (10.4%) and other short-acting methods, such as contraceptive injections (2%), vaginal rings (0.8%), and patches (0.3%) [2]. With increasing levels of formal education, the prevalence of LARC and short-acting methods increases while the prevalence of female sterilization decreases [3].

While female contraceptive products are reversible and generally more efficacious than traditional methods, thus offering advantages to women with regards to their family planning and thus self-determination [1,4], they can be associated with unpleasant experiences [5,6], ranging from abdominal pain to mood swings or changes in libido [7,8]. The experience of such unpleasant side effects has a negative impact on a woman's health, which the World Health Organization defines as "state of complete physical, mental, and social well-being," and thus on the quality of life [9,10]. Furthermore, negative side effects are a major cause for poor adherence or even discontinuing contraception which may result in unintended pregnancies [11,12].

Access to Data to Inform Contraceptive Choices

For women to find the contraceptive method that is most suitable for them and thus make informed contraception choices, it is important to have access to relevant information regarding different available contraception options. The type of information that women require can be assigned to 2 broad categories.

First, information relating to the efficacy of contraceptive methods regarding preventing unintended pregnancies and protection from sexually transmittable diseases is crucial [13]. There is comprehensive clinical as well as real-world data on efficacy and safety of different contraceptive methods [14,15]. These data are generally accessible to women through health care providers (HCPs) or nongovernmental organizations, although there are geographical differences on a global level [16].

Second, women seek information relating to potential unpleasant experiences related to contraceptive methods as these can have a substantial negative impact on women's well-being, not only impacting women themselves, but also their families [13,17,18]. However, there are 2 major challenges women face when seeking information about potential negative experiences related to contraceptive methods, namely, the availability of data and the accessibility of reliable data [17].

Data on the frequency of negative side effects are generally available, as they are collected in clinical trials and stated on drug labels [19,20]. However, the construct of well-being is more nuanced, comprising a "state of positive feelings and meeting full potential in the world" [21]. Consequently, data on the mere occurrence and frequency of certain side effects provide insufficient information on how certain side effects typically impact well-being. For example, abdominal pain

https://ai.jmir.org/2025/1/e68809

related to a contraceptive product might constitute a neglectable nuisance or a major suffering limiting women's participation in daily life. Despite the subjective nature of side effect severity [22], for women facing contraception choices, knowing that a certain side effect can be a significant issue for some women constitutes relevant information [17]. However, there is a lack of comprehensive data on women's subjective and collective unpleasant experiences with different contraceptive methods [23]. Studies have also shown that while women tend to turn to HCPs for contraceptive counseling, HCPs often lack relevant knowledge and provide insufficient information on potential side effects [24,25]. To learn about experiences with different contraceptive methods, women also tend to speak to relatives and friends, but these experiences are subjective and constitute a small sample size.

Role of Social Media to Inform Contraception Choices

From this background, social media has started to play an important role as a source of information. Experimental research indicates that social media content may influence women's intentions to use certain contraceptive products [26] even as social media conversations about contraception have become more polarized in the past 20 years [27].

Thus, there is a growing body of research to evaluate how women use social media to inform contraception choices [28]. To analyze information shared and consumed on social media, natural language processing (NLP) is used due to its capacity to analyze nonstructured, textual data.

For example, Pleasants et al [28] used NLP to study posts related to birth control on the US platform Reddit and found that "Side Effects!?" is the most common flair, a tag that users can attach to their post to categorize the content. Furthermore, "Experience" and "Side Effects?!" are the most common flairs among the most popular posts, based on the number of comments and "scores," that is, upvotes minus downvotes [28]. Analyzing contraceptive content shared on X (formerly Twitter) Huang et al [29] discovered that a fifth of all the tweets relate to side effects. Similarly, in a text mining analysis of messages sent on a free sexual and reproductive health information service in Kenya, Green et al [30] found that users wrote most often about family planning and side effects. Balakrishnan et al [31] conducted an NLP-based social listening analysis in a German internet community and observed that side effects are the most common problem associated with most contraceptives. They also found that while the pill is the most frequently mentioned contraceptive method, there appears to be migration from hormonal to nonhormonal methods. In line with this, Felice et al [32] analyzed user reviews of a digital contraceptive supporting women in fertility prediction through a mixed methods approach involving NLP and found that the hormone-free aspect of the contraception experience is highly salient for many users. A content analysis by Pfender and Devlin [33] of YouTube vlogs discussing birth control methods revealed that social media influencers primarily described their discontinuation of hormonal birth control due to experienced side effects. Their study also showed that vlogs may provide inaccurate sexual health information, hereby directly or indirectly discouraging the use of the contraceptive under

XSL•FO RenderX

discussion. In a content analysis of the "sex secrets" Facebook page, Yeo and Chu [34] found that young people predominantly use this social media platform to request information, opinion, or advice, including the topic of birth control. Stoddard et al [35] found that more than half of the most popular contraception videos on TikTok revolved around patient experience. Although videos created by health care professionals received proportionately more views, over half of the total views were still of content generated by laypeople [35].

Overall, these NLP-based social media content studies show that social media is used to share information and consume information on contraception options. Furthermore, they reveal that user-generated content mostly revolves around side effects and that posts discussing women's experiences with regards to side effects receive the greatest interest. At the same time, the content that is available, especially when shared by influencers, is not always reliable and may misguide contraception choices. In fact, there is increasing concern among researchers and women's health practitioners that social media influencers spread misconceptions about contraceptive methods, particularly hormonal contraception, which negatively affect the acceptance of efficacious contraceptive methods and thus increase the risk of unintended pregnancies [36].

Furthermore, the previously mentioned studies highlight that unpleasant experiences are an issue that is currently not well-addressed in clinical contraceptive counseling. This further substantiates the observation that users appear to have an unmet need for reliable, trustworthy information. However, existing NLP-based studies do not provide a systematic picture of the association of different side effects with different available contraceptive methods and how severely women experience these side effects.

To fill this gap, the NLP method of sentiment analysis can identify, extract, and quantify the subjective emotions within a text, assigning a continuous sentiment score usually between -1 for highly negative and 1 for highly positive posts [37,38]. Studies using sentiment analysis may thus provide first hints as to how severely women experience certain side effects. Merz et al [27] studied population attitudes toward contraceptive methods over time by performing sentiment analysis on tweets on X regarding contraceptive methods and find that most tweets are negative. In their sample long-acting methods are mentioned more often than short-acting ones and related sentiments are twice as likely to be positive [27]. In contrast, in a study with Indonesian users of X, Sari et al [39] found that users predominantly express negative attitudes toward long-acting contraceptive methods.

However, a major limitation of sentiment analysis is that it can be inaccurate if the model has been trained on biased, limited, or unrepresentative datasets as it may fail to generalize well to diverse and nuanced language usage, such as sarcasm, slang, or cultural context variations present in social media posts. Although the modern state-of-the art approach in sentiment analysis involves using pretrained language models such as Bidirectional Encoder Representations from Transformers or GPT, there is an inherent risk of bias in general and gender bias in particular [40,41], limiting performance in sentiment analysis tasks.

In this context, the information on online drug review forums constitute a great, widely untapped, resource to inform contraception choices. Many online drug review forums contain 2 distinct pieces of information related to a product: a standardized numeric rating score indicating overall product satisfaction and a comment in free text form. A powerful advantage of online product reviews is that the integration of qualitative (text comments) and quantitative (ratings) data facilitates insights into the relationship between issues mentioned in comments and overall product satisfaction, which is presumably closely linked to the impact of the respective product on the well-being of the user.

Evaluating data from online review forums to inform decisions is hampered by several limitations, such as potential biases, unrepresentative sample issues, and the potential presence of inauthentic reviews. Nevertheless, consumer behavior in many industries, including health care and retail, indicates that other people's reviews, particularly when available in large numbers, are important in driving purchase decisions and are thus considered a valuable source of information [42]. Thus, in the context of contraception, reviews data complement information on contraceptive options that women and their partners may receive from other sources, such as HCPs, community workers, scientific studies, or other social media sources.

Purpose of the Study

This research aimed to produce insightful information from a large drug review dataset with regard to which experiences with contraceptive products women described on the web, both qualitatively and quantitatively. The focus is on unfavorable experiences, as previous research has shown that side effects are the topic of greatest interest for women using forums and social media to seek information on contraception.

From this background, in this paper, we investigated the following research questions (RQs):

- 1. How do users rate different contraceptive methods available to women on a major drug review website?
- 2. Which issues (ie, topics) do users describe in unfavorable online reviews of contraceptive products available to women?
- 3. How frequently are these issues described for different contraceptive methods?
- 4. Can we observe an association between the main topic discussed and the average rating submitted in unfavorable birth control reviews?

Methods

Dataset

Our study was performed on a dataset of 19,506 unique online reviews of birth control products in the United States posted on the website Drugs.com [43], a United States-based pharmaceutical information website, between April 2009 and September 2017. The reviews analyzed in this study were extracted from a comprehensive online drug review dataset

available for research purposes in the University of California, Irvine (UCI) Machine Learning Repository [44]. The original dataset had been collected via web scraping from the website Drugs.com [43] and comprised 215,063 reviews of drugs treating different conditions, such as high blood pressure, cough, and birth control [45]. While this dataset may be somewhat dated, these reviews are highly relevant for this study. First, the products evaluated have been on the market for many years and are widely used today. Second, analyzing older reviews might even offer the advantage of capturing women's experiences with contraceptive products in a more authentic, less skewed way. Research has shown that in recent years, social media influencers negatively frame hormonal contraceptives and encourage the uptake of nonhormonal options which may alter women's attitudes and expectations [26] and thus possibly their online reviews.

The online drug user reviews contained information on the related condition, the name of the drug, a 10-star user rating on overall satisfaction, how many users considered this review helpful, and the date the review was posted. The name of the drug was captured in a structured format as it stemmed from a drop-down menu from which the website users needed to select a drug name when leaving a comment.

Ethical Considerations

The study used publicly available data from the UCI Machine Learning repository. The UCI Machine Learning dataset did not contain any identifying information about the authors of the reviews, such as their username. Furthermore, when posting a review on the website Drugs.com [43], users were required to consent to the publication and use of their reviews. Finally, to the best of our knowledge, the reviews we selected to be in this manuscript do not risk reverse identification as the website Drugs.com [43] does not display full user names alongside the reviews. Therefore, in line with other studies evaluating social media posts on contraception, ethics approval for using these reviews as a basis for analysis was not deemed necessary.

Data Cleansing and Grouping

Within the drug review dataset offered by the UCI Machine Learning Directory, 38,436 product reviews were classified as relating to "birth control." Many reviews were captured twice, once under a product's brand name and a second time under the name of the respective active pharmaceutical ingredient, that is, the generic name. By removing duplicates, we obtained 19,524 unique birth control reviews. When cleansing the dataset, we retained drug brand names for their greater detail compared to generic names. This granularity is more suitable for analysis, as products with the same active ingredient can vary in dosage and administration schedules across brands. In total, <400 out of 19,524 unique reviews did not contain a brand name, but rather only the generic name. We kept most of those reviews in the dataset, only removing 13 unique reviews of drugs that could not be related to 1 specific contraceptive method, namely, levonorgestrel (10 reviews), which can be a hormonal IUD or emergency contraception commonly sold as Plan B; and Provera (3 reviews), which can either be a birth control shot under the name Depo-Provera or an oral progestin product that is not approved as a contraceptive. The clean birth control dataset

XSL•FC

contained unique reviews on 169 different products identified by brand name or active pharmaceutical ingredients.

For later analysis and comparison of drug reviews for different contraceptive methods, we categorized all products into 11 contraceptive methods. This categorization focused on the application mode of these products, which is in line with the classification of contraception options typically used for advising women [46]. The methods comprise: hormonal and nonhormonal (ie, copper) IUDs, implants, vaginal rings, birth control shots, hormonal patches, spermicides, and emergency contraception. For OCPs, we distinguished between combined contraceptive oral pills (COCPs), progestin-only pills (POPs), and OCPs that induce a 91-day cycle, as these are expected to have different side effect profiles, and patients are typically counseled differently. Given the small number of reviews on emergency contraception (n=3) and spermicides (n=2), we removed those reviews from the dataset, too, leaving 19,506 reviews on 167 different products across 9 different contraceptive methods.

To analyze which negative experiences or side effects related to birth control options women described, we created a new attribute marking all reviews with a rating of ≤ 5 (on a scale from 1 to 10) as unfavorable reviews. Rather than limiting our analysis to reviews associated with strictly negative ratings (usually defined as ≤ 3), we deliberately used a wider window to also include negative to neutral ratings (scores of 4 and 5) as these might also contain relevant descriptions of unpleasant experiences. Overall, 8330 reviews fell into our definition of unfavorable (ie, nonpositive with a rating of 5 and lower).

NLP Approach for Analyzing Unfavorable Birth Control Reviews

Overview

Within NLP, topic modeling refers to techniques for uncovering abstract themes in a large textual dataset, typically referred to as a corpus. It involves algorithmically analyzing documents to detect word and phrase patterns that indicate specific topics. Thus, topic modeling allows analyzing which topics women discuss in unfavorable reviews of different contraceptive products. For our study, we wrote an NLP program for topic analysis in Python (version 3.11.3; Python Software Foundation) using several NLP libraries, including Natural Language Toolkit (version 3.7) [47] and scikit-learn (version 1.2.2) [48]. For visualization, we used Matplotlib (version 3.7.1) [49] and Seaborn (version 0.12.2) [50].

Text Preparation

Our text preprocessing procedure included multiple steps. Text cleaning was performed as the raw birth control reviews in the UCI repository contained several issues with punctuation and how certain characters were captured. Furthermore, we implemented a custom-developed catalog of more than 110 abbreviations and short forms to replace them with the long form. Examples include "can't" being replaced with "cannot," "PMS" with "premenstrual syndrome," "yr" with "year" or "ain't" with "am not." If an abbreviation had 2 meanings, such as "he's," we replaced it with the most common form. This step ensured uniformity in word representation so that the frequency of a word could be captured adequately. In addition, we removed any nontext characters, created word tokens and reduced words to their base root via lemmatization. To further reduce the dimensionality of the textual data and focus on the most meaningful words, we excluded common words typically not carrying meaning, so-called "stop words" as predefined in NLTK, except for the stop word "not" which adds to the meaning of a review describing potential complaints or side effects. We also removed all product names and contraceptive methods, such as "iud," "implant" or "pill," from the reviews to allow our topic modeling algorithm to reveal topics that are contraceptive product and method agnostic.

As Table 1 shows, after the removal of stop words, there were highly frequent words in the reviews that did not relate to specific birth control side effects or complaints. To reduce noise and dimensionality, we removed the words "month," "day," "year," "week" "birth," and "control" from the reviews.

Table 1. Most common words in the birth control product reviews (excluding noninformative words).

Word	Occurrences (n=889,864), n (%)
not	30,973 (3.48)
period	19,145 (2.15)
month	18,361 (2.06)
day	11,139 (1.25)
control	10,902 (1.23)
birth	10,678 (1.2)
year	9986 (1.12)
week	9464 (1.06)
first	8702 (0.98)
get	8020 (0.9)
weight	7783 (0.87)
would	7146 (0.8)
got	7061 (0.79)
time	6956 (0.78)
like	6410 (0.72)
side	6186 (0.7)
effect	5982 (0.67)
cramp	5704 (0.64)
started	5545 (0.62)
since	5332 (0.6)
mood	5302 (0.6)
taking	5292 (0.59)
bleeding	5278 (0.59)
acne	5156 (0.58)
never	5018 (0.56)

Topic Extraction Approach

In topic modeling, selecting the optimal vectorization techniques, topic modeling algorithms, and the number of topics to extract is crucial. This process aimed to derive topics that align with domain-specific inquiries and RQs. While coherence and silhouette scores can support this selection, domain expertise and expert judgment are essential in evaluating the relevance and applicability of the themes extracted by an algorithm [51].

The topics described in the following sections result from an iterative strategy combining various vectorization techniques to construct a document-term matrix, including count

```
https://ai.jmir.org/2025/1/e68809
```

RenderX

vectorization and term-frequency-inverse document frequency (TF-IDF). We used topic modeling algorithms, such as latent semantic analysis, nonnegative matrix factorization (NMF), and latent Dirichlet allocation, extracting between 3 and 13 topics. The selection of techniques and the number of topics was based on expert judgment, Cohen coherence, and the silhouette score, with the final decision guided by domain knowledge to yield the most useful, interpretable, and distinct topics.

The final technical configuration of the topic modeling in this research is as follows:

 Vectorization—TF-IDF vectorization yielding a document-term matrix, where rows represented reviews,

columns represented words, and values indicated word importance. TF-IDF highlights terms frequent in a document (ie, a product review) but less common across the corpus (ie, across all reviews), reducing the weight of ubiquitous words.

• Topic modeling algorithm—NMF decomposing the nonnegative document-term matrix into 2 lower-dimensional matrices: the topic matrix (W) and the

terms matrix (H). The topic matrix represented documents by underlying topics, while the terms matrix represented topics by original words or tokens.

• Number of topics—8 topics differentiating most effectively among various types of experiences and complaints.

The flowchart in Figure 1 illustrates the overall methodological approach.

Exclusion: favorable

Figure 1. A flow diagram of the methodological approach used in the study. CV: count vectorization; LDA: latent Dirichlet allocation; LSA: latent semantic analysis; NLP: natural language processing; NMF: nonnegative matrix factorization; TF-IDF: term-frequency–inverse document frequency; UCI: University of California, Irvine.



Part II: NLP via topic modeling

Contraceptive product



Results

Descriptive and topic analysis of the website Drugs.com [43] drug reviews dataset allowed us to answer our RQs.

Online Ratings of Different Contraceptive Methods Available to Women (RQ 1)

Frequency of Ratings

Table 2 displays the distribution of ratings of birth control products in the drug review dataset from 1 to 10, with 1 being

very bad and 10 being very good. The frequency diagram of the ratings is U-shaped, such that both very poor and very good ratings were common. The most common rating was 10 out of 10 (n=3905, 20.02% of the reviews), and the second most common rating was 1 out of 10 (n=2986, 15.31% of the reviews). The overall mean was 6.08, and the SD was 3.31. Thus, reviews were polarized, but on average gravitated toward positive ratings.



Groene et al

Table 2. Frequency of contraceptive product ratings on a scale from 1 to 10 (1: very bad and 10: very good) in the online drug review dataset (n=19,506).

Rating	Frequency, n (%)
1	2986 (15.31)
2	1409 (7.22)
3	1363 (6.99)
4	1083 (5.55)
5	1489 (7.63)
б	964 (4.94)
7	1253 (6.42)
8	2112 (10.83)
9	2942 (15.08)
10	3905 (20.02)

Ratings of Different Contraceptive Methods

Table 3 provides an overview of the number of available birth control product reviews in the dataset, grouped by the product categories. COCPs are the most reviewed birth control products, constituting 44.12% (8606/19,506) of all reviews. Hormonal implants and hormonal IUDs rank second and third, respectively.

Slightly more than half of birth control product reviews (n=11,176, 57.3%) are favorable according to our definition, whereas 42.7% (8330) of reviews are unfavorable. Overall, the share of unfavorable reviews varied substantially across categories. POPs had the highest share of unfavorable reviews

(n=232, 53.1%), whereas nonhormonal IUDs had the lowest (n=234, 29.3%).

Figure 2, a scatter diagram with trimmed axes, reveals 2 clusters of different contraceptive methods based on mean ratings and SDs. The first cluster, located in the lower right, includes POPs, birth control shots, 91-day cycle OCPs, COCPs, and hormonal implants, with lower average ratings (5.32-5.82) and higher SDs (3.26-3.52). The second group, situated in the upper left, comprises hormonal and copper IUDs, hormonal patches, and vaginal rings, exhibiting higher average ratings (6.65-7.11) and generally lower SDs (2.99-3.13), except for copper IUDs, which had a SD of 3.28.

Table 3. Descriptive statistics of product ratings by contraceptive method.

Contraceptive method	Unique reviews (n=19,506), n (%)	Products (n=167), n (%)	Rating, mean (SD)	Number and share of unfa- vorable reviews (n=8330, 42.7%), n (%)
COCP ^a	8606 (44.12)	138 (82.6)	5.80 (3.32)	3968 (46.11)
Implant	4392 (22.52)	3 (1.8)	5.82 (3.32)	2064 (46.99)
Hormonal IUD ^b	2871 (14.72)	4 (2.4)	7.04 (3.05)	848 (29.54)
Vaginal ring	827 (4.24)	2 (1.2)	6.7 (3.0)	297 (35.91)
Copper IUD	800 (4.1)	2 (1.2)	7.11 (3.3)	234 (29.25)
Shot	653 (3.35)	2 (1.2)	5.5 (3.5)	327 (50.08)
Patch	508 (2.6)	3 (1.8)	6.8 (3.1)	157 (30.91)
POP ^c	437 (2.24)	11 (6.6)	5.3 (3.3)	232 (53.09)
OCP ^d with 91 d cycle	412 (2.11)	9 (5.4)	5.6 (3.3)	327 (49.27)

^aCOCP: combined contraceptive oral pill.

^bIUD: intrauterine device.

^cPOP: progestin-only pill.

^dOCP: oral contraceptive pill.



Figure 2. A scatter diagram with trimmed axes visualizing average rating and SD of different contraceptive methods. COCP: combined contraceptive oral pill; IUD: intrauterine device; POP: progestin-only pill.



Issues Described by Users in Unfavorable Online Reviews of Contraceptive Products Available to Women (RQ 2)

Table 4 presents the 8 themes extracted from the 8330 unfavorable birth control product reviews in our dataset. Overall, each extracted theme corresponded to a description of side effects. There was no topic that explicitly alluded to nonhealth aspects such as cost, ethical or societal concerns, or accessibility. A total of 4 topics extracted were highly specific and related to "weight gain," skin problems," "loss of libido," and "mental health problems." Another 3 topics related to the impact of the contraceptive product on women's menstrual cycle but alluded to distinct aspects, which we named "menstrual irregularities," "cramps and pain," and "continuous bleeding." The final topic, "multiple cause dissatisfaction," was a mixed, broad topic. It was the least distinct topic, comprising a mixture of diffused complaints ranging from headache, tiredness, general life, and relationship issues to a mere product warning. A sample review that scored high on the topic "multiple cause dissatisfaction" read as follows:

Makes me feel very moody and sensitive, my husband and I fight all the time. When we got married I felt so much in love but know not sure about it. He said I changed a lot after having our baby. So not sure if the IUD is making me feel that way. I feel so bad because I get mad very easy for little things and I feel like I am loosing my husband. Of course that he doesn't want to wear his ring makes me think things but he said that he is not use to wear rings and I always wear mine. I cook breakfast every single day, cook lunch for us to take it to work since we do not have to much money and sometimes I feel that he doesn't really appreciate it! Do laundry, clean and he doesn't really help me much and he doesn't see it. Not sure what to think.

Table 5 provides sample reviews for each topic. More examples can be found in Multimedia Appendices 1 and 2.

For each topic, Table 4 also presents the share of online user reviews where this topic was dominant having the highest topic value in the topic matrix W. Thus, the dominant topic is the issue voiced most firmly in an unfavorable review. Table 4 shows that with this dominant topic modeling scheme, the reviews were relatively evenly distributed among the 8 identified topics. The rarest dominant topic was "weight gain" with 9.69% (807/8330) of the unfavorable reviews predominantly describing this side effect. "Multiple cause dissatisfaction" dominated in 17.92% (1493/8330) of the unfavorable reviews.

https://ai.jmir.org/2025/1/e68809

RenderX

Groene et al

Table 4. Topics discussed in unfavorable reviews of birth control products and their relative frequency (n=8330).

Торіс	Topic description	Unfavorable reviews with this dominant topic, n (%)
Weight gain	Users describe a change in body weight, typically an increase, which is attributed to the contraceptive product	807 (9.69)
Skin problems	Users describe an impact of the product on outward appearance, in particular acne	1051 (12.62)
Loss of libido	Users describe a reduction or loss of interest in physical intimacy and in- tercourse	963 (11.56)
Mental health problems	Users describe mental health problems, such as mood swings, depression, and anxiety	902 (10.83)
Menstrual irregularities	Users describe different problems with their period resulting from the contraceptive method; ranging from spotting, heavy bleeding, to unusually light periods	1223 (14.68)
Cramps and pain	Users describe particularly painful experiences, especially cramps, associated with the product or its administration	926 (11.12)
Continuous bleeding	Users describe continuous bleeding episodes which last substantially longer than regular periods and are more pronounced than simple menstrual irregularities	965 (11.58)
Multiple cause dissatisfaction	Users express dissatisfaction with the contraceptive product. None or various reasons are provided ranging from general side effects such as headaches to overall challenges in life that might or might not be at- tributable to the contraception choice	1493 (17.92)

Table 5. Examples of reviews centering on a specific topic.

Торіс	Sample reviews with the dominant topic
Weight gain	"Been on it for 3 months, 20 pound weight gain—always hungry and never full. No periods, but not worth the weight gain and uncontrollable appetiteWas managing weight very well prior to implant"
Skin problems	"Horrible, horrible I have never had acne this bad in my life!!!!!!!! My WHOLE chin and jawline are red and covered in cystic acne!!! I HAD PERFECTLY CLEAR SKIN BEFORE. I am honestly in a complete panic with what is going on with my skin. I'm in shock that a small pill could do this much damage. My face hurts so bad because of the acne. Its been only 3 weeks since I started taking it. Switching to sprintec tomorrow. DO NOT USE THIS, SAVE YOUR SKIN!!!!!"
Loss of libido	"I have been on NuvaRing for 5 months. Within a month I noticed a decrease in my sex drive, and I've had vaginal dryness which makes sex painful. Bad sex has effected other parts of my life."
Mental health problems	"I used this pill during my teens and it caused irritability and heavy mood swings. Perhaps it was just teen angst but I tried microgynon recently, which uses the same hormones just different levels, and experienced similar mood swings and depression."
Menstrual irregularities	"I have been on this medication for almost a month. I got my period once, but it hasn't even been a week later that I got a second period. My first period was very light and only lasted three days, but I'm not sure how this period will be."
Cramps and pain	"I got the kyleena inserted today and experienced the worst cramps in my life. The insertion were (8/10) on the pain scale. I am not very sensitive to pain but can't take any pain medication. The last 4 hours has been the worst in my entire life so far I have really bad cramps now 10/10 and nausea. I can't even get out of bed because of the severe pain!"
Continuous bleeding	"With liletta I have been bleeding for 3 month s I am so so tire of bleeding."
Multiple cause dissatisfaction	"Do not take this pill."

Relative Frequency of Side Effects Described Across Contraceptive Methods (RQ 3)

For each contraceptive method, Figure 3 shows the relative frequencies of the dominant topics in descending order according to average rating. For both copper and hormonal IUDs, the most frequent complaint by far was "cramps and pain," which was the dominant theme in 38% (n=90) and 42%

RenderX

(n=355) of the reviews, respectively. For COCP, the most common complaints were "multiple cause dissatisfaction" (n=831, 21%) and "skin problems" (742, 19%). For POP and OCP that induce a 91-day cycle, "menstrual irregularities" were the most common issue (n=49, 21%, and n=49, 24% of reviews, respectively). For implants, as well as for hormonal shots, "continuous bleeding" (n=436, 21%, and n=61, 19%, respectively) was the most frequent problem described in the

reviews. For hormonal patches and vaginal rings, the most frequent dominant topic was "multiple cause dissatisfaction" (n=57, 36% and 106, 36\%). For hormonal patches, the second

most frequently described dominant side effect prominently voiced in unfavorable online reviews was "skin problems" (n=28, 18%).

Figure 3. Relative frequency of dominant topics in nonfavorable reviews by contraceptive method (as percentages). COCP: combined contraceptive oral pill; IUD: intrauterine device; POP: progestin-only pill.



When reviewing the relative frequencies of dominant topics identified in Figure 3, it is important to remember that each contraceptive method was associated with a different proportion of unfavorable reviews. This was analyzed in the context of RQ 1 and is depicted in Table 3 which displays substantial variation in the proportion of unfavorable reviews across contraceptive product categories; with POP having the highest share of unfavorable reviews (n=232, 53.1%) and copper IUDs having the lowest share of unfavorable reviews (n=234, 29.3%). Scaling the relative frequencies of the dominant topics shown in Figure 3 with the overall share of unfavorable reviews of contraceptive methods displayed in Table 3, we find that for certain contraceptive methods, specific issues were very commonly discussed in online reviews in general, as in the following examples:

- Almost a quarter, that is, 24% (n=97), of all reviews of 91-day cycle OCPs report general "menstrual irregularities" or "continuous bleeding" (this is derived from n=49, 24%, of the reviews where "menstrual irregularities" were the dominant topic plus another n=48, 24%, where "continuous bleeding" was the dominant topic; multiplied by 49.3%, the rate of unfavorable reviews).
- Overall, 17% (n=104) of all reviews of hormonal shots discussed "menstrual irregularities" or "continuous bleeding."
- For IUDs, 12% (n=90, copper) and 11% (n=355, hormonal) of all reviews revolved around "cramps and pain" associated with the contraceptive method and its administration.
- "Loss of libido" was the dominant topic in almost every 10th review of vaginal rings, that is, 9% (n=75).
- Almost 6% (n=243) of all reviews of "hormonal implants" revolved primarily around mental health issues.

Association Between Dominant Topic and Ratings of Birth Control Products (RQ 4)

The final RQ relates to how severely such side effects might impact the well-being and overall quality of life of women. This approach is important for providing a balanced picture of the frequency numbers described earlier. Not every side effect, even if common, necessarily impacts overall well-being to the same extent. For example, individual reviews illustrate that "cramps and pain" might have a much more negative impact on overall well-being than menstrual irregularities. For illustration, a sample review with "cramps and pain" as the dominant topic read as follows:

My experience was absolutely horrible. Birth control works different for everyone but this was by far the worst pain I've ever been in...

While a sample review where menstrual irregularities were voiced read as follows:

I have been on this medication for almost a month. I got my period once, but it hasn't even been a week later that I got a second period. My first period was very light and only lasted three days, but I'm not sure how this period will be.

Figure 4 displays boxplots of unfavorable ratings by dominant topic. The boxplots show that the frequency distributions for all dominant topic-based groups are left skewed. Consistently, the first quartile of ratings is a 1, that is, at least a quarter of all reviewers (2986 across all groups, ie, 35.8% on average) writing a nonpositive review submitted the lowest possible rating. The medians displayed in the boxplots as orange vertical lines range from 2 to 3. Only 3 dominant topics were associated with a median rating of 3, namely "menstrual irregularities," "weight gain," and "loss of libido." For all other dominant topics, half of all unfavorable reviews have a rating of ≤ 2 .

```
XSL•FO
RenderX
```

Figure 4. Boxplots of ratings by dominant topic described in nonfavorable reviews.



The mean ratings per dominant topic are represented as green star. On average, reviews predominantly describing menstrual irregularities have the highest average rating (mean 2.92, SD 1.53; 1223/8330, 14.68%). The next highest average ratings were in reviews describing weight change (mean 2.87, SD 1.52; 807/8330, 9.69%) and loss of libido (mean 2.84, SD 1.49; 963/8330, 11.56%). Conversely, reviews with the lowest ratings, on average, predominantly described multiple cause dissatisfaction (mean 2.34, SD 1.45; 1493/8330, 17.92%), cramps and pain (mean 2.39, SD 1.55; 926/8330, 11.12%), and continuous bleeding (mean 2.45, SD 1.47; 965/8330, 11.58%). Skin problems (mean 2.53, SD 1.49; 1051/8330, 12.62%) and mental health problems (mean 2.57, SD 1.48; 902/8330, 10.83%) had slightly greater ratings than the bottom 3 groups.

Discussion

Principal Findings

Our findings are in line with the literature evaluating how users, mostly women, use social media to discuss and evaluate different contraception options. Side effects were the most important area that was discussed online. Our NLP algorithm extracted 8 topics, of which 7 clearly describe a specific unpleasant side effect and only 1 less concise topic also encompasses other issues such as general life challenges, relationship issues, or product warnings. The algorithm did not extract any frequent words indicative of nonhealth related challenges such as cost and accessibility.

Our research extends the existing body of knowledge in several aspects. First, we found that in the online drug review forum, niche products tend to be overrepresented compared to their prevalence among the respective populations in the United States. For example, Table 3 shows that 22.52% (4392/19,506)

RenderX

of the reviews discussed hormonal implants. However, their prevalence among women using reversible contraceptive products (either LARCs or short-acting methods) is 7.3% [2]. Similarly, vaginal rings (827/19,506, 4.24% of reviews vs 2.9% in the respective population [2]) and hormonal patches (800/19,506, 4.1% vs 1.1% [2]) appear to be overrepresented. Even more interestingly, IUDs are substantially underrepresented. While 18.82% (3671/19,506) of contraceptive product reviews discuss IUDs, they are used by almost a third of women [2] using female contraceptive products.

The underrepresentation of IUDs might be attributable to the fact that on average, women tend to be more satisfied with IUDs than with other contraceptive products (Figure 2; Table 3). In general, people are more likely to write an online review when they have a complaint than when they are satisfied [52]. Nevertheless, in Figure 2 and Table 3, we observe that a substantial share of women report positive experiences with contraceptive products, ranging from slightly >70% (2589/3671, 70.52%) of favorable reviews for IUDs to 46.9% (205/437) for reviews of POP. Overall satisfaction with reviewed birth control products may be even greater if the probable negative bias inherent in online reviews is accounted for.

Our NLP-based evaluation of product reviews also offered valuable insights into how women experience different product-based contraceptive methods and how negative experiences relate to the overall satisfaction with the method.

Hormonal Short-Acting Contraceptive Methods

Overview

A first pattern we observed was that hormonal contraceptive methods with a higher level of systemic absorption, such as POP, birth control shots, and COCP, received greater shares of unfavorable reviews (232/437, 53.1%, 327/653, 50%, and

3968/8606, 46.1%, respectively) than methods with a lower level of systemic absorption, such as hormonal IUDs (848/2871, 29.54%). For copper IUDs, which do not release any hormones, only 29.3% (234/800) of reviews were unfavorable (Table 3). Thus, based on the online reviews and ratings, it appears that on average women in the website Drugs.com [43] sample are less satisfied with short-acting methods that rely on a systemic hormonal effect, which is in line with the findings of Merz et al [27] studying posts on X over time. This is particularly interesting as recent literature describes a trend of women turning away from hormonal contraceptives despite their high efficacy due to the influence of social media [26,36]. While clinical studies confirm a range of side effects with hormonal contraception, some researchers suspect they may be perceived to be more severe than they truly are [36]. However, our research suggests that women rate oral hormonal contraceptive products, hormonal implants, and shots less positively than nonhormonal methods and that this is linked with specific side effects.

Progestin-Only Methods and the Role of Irregular Bleeding Pattern

Figure 2 shows that POPs, contraceptive implants, and injectable contraceptives, which are all progestin-only methods, obtained comparably low average ratings with a high SD. According to the reviews, the most common side effects for these methods relate to irregular bleeding pattern and continuous bleeding (Figure 3). This is in line with the relevant women's health literature describing irregular or unscheduled bleeding as their most common side effect (eg, [53]). The inhomogeneous experiences with these progestin-only methods might be—to some extent—explained by the interplay between users' expectations and actual experiences. Although irregular bleeding resulting from these progestin-only methods decreases over time [54], women may be more disgruntled by the initial irregularity, especially if counseling focused more on the long-term than short-term expectations.

Potential Role of Ease of Administration for Cycle Control

Among the remaining short-acting contraceptive methods, hormonal patches and vaginal rings obtained comparably high average ratings and lower SDs than COCPs and extended-cycle OCPs (Figure 2). For COCPs, this is expected as women who are prescribed contraception for the first time often opt for the COCP due to its ease of initiation and discontinuation [55]. This initial usage likely leads to varied experiences.

However, this disparity in low average ratings of extended-cycle OCPs versus comparably high and more homogenous average ratings of vaginal rings and hormonal parches is unexpected, given the similarity of side effects across combined hormonal contraceptives (CHC) [53]. A potential explanation of our findings is that the patch and the ring may achieve better cycle control than the extended-cycle COCP due to the lack of need for daily administration [53]. Indeed, for extended-cycle OCPs, nearly half of the reviews (97/203, 47.8%) predominantly described abnormal bleeding, whereas this was only the case for 14% (22/157) of reviews of hormonal patches and 5.7% (17/297) of reviews of vaginal rings (Figure 3). Those opting for an extended-cycle OCP probably do so with the intent of

XSI-FC

significantly reducing or entirely ceasing their menstrual cycles, which is the primary distinction between standard and extended-cycle OCPs. Unfortunately, breakthrough bleeding is very common early in the use of an extended-cycle regimen [53]. Therefore, women who are hoping for no bleeding are likely to be unhappy with increased bleeding, especially if they are not warned about this. It is also possible that there could be other explanations, such as increased hormonal stability with nonoral administration of CHC [55].

Skin Issues

Research has consistently shown that CHC are beneficial for the treatment of acne [56]. It is surprising, therefore, that in our study, skin problems appear as a dominant topic for COCPs, patches, and 91-day cycle OCPs more commonly than for other methods. Further research would be helpful to explore these findings. One possibility is that the skin problems reported by reviewers are not only acne but also other problems, such as melasma, a hyperpigmentation disorder well known to be associated with oral contraceptive use [57]. However, an example investigation of reviews in which "skin problems" are discussed reveals that many reviews describe disappointment resulting from the birth control product not meetings their expectations with regards to acne control. For example,

I've been taking this birth control for about a week now, and I have already noticed some changes. My skin is also acne prone, and I was really hoping that this birth control would help with it. Without the pill, I usually have many bumps on my forehead, my chin is pretty red, and once in a while I will get cystic acne. Now that I've been taking it, I have many new pimples all over my face, like my cheeks and on my nose, where I have never gotten it before. It's also given me MORE cystic acne which is a pain. I really wish that it could have helped, but before I switch off I want to wait a little longer to be sure.

This disappointment might make them more prone to leave a negative comment than women who experience other side effects.

Loss of Libido

Interestingly, the loss of libido is still associated with comparably high average product ratings. Unfavorable reviews predominantly describing a loss of libido ranked third in average rating (Figure 4). Our analysis also revealed that loss of libido is the most common dominant topic for vaginal rings (75/297, 25.3% of unfavorable reviews), almost double the proportion of any other method. This is interesting since controversial findings on this topic exist in the literature. It has been hypothesized that CHC may adversely affect sexual functioning by increasing sex hormone-binding globulin (SHBG), which then decreases available testosterone and leads to decreasing endogenous hormone production. Oral estrogens are known to increase SHBG via a first-pass hepatic effect [58]. It could be hypothesized that nonoral administration may have a smaller effect on available testosterone, although other research has shown that both oral and vaginal CHC increase SHBG and decrease free androgens [59]. It has also been hypothesized that administering CHC via a nonoral route, such as the use of a
patch or ring, may mitigate effects on sexual function via increased hormonal stability. The ring could also exert a local estrogenic effect, improving lubrication [55].

Several studies have assessed the effect of the vaginal contraceptive ring on sexual functioning, with mixed findings [55,59], and a recent meta-analysis revealed a possible positive effect at 3 months but no effect at 6 months [60]. A larger cross-sectional nonrandomized analysis revealed that decreased libido was most common among users of shots, rings, and implants [8], which is more congruent with our analysis. It is also worth considering that direct hormonal effects are not the only way a method could affect libido; physical discomfort, vaginal dryness or irritation, and excessive bleeding are also expected to contribute. This is also illustrated in this example review:

I have been on NuvaRing for 5 months. Within a month I noticed a decrease in my sex drive, and I've had vaginal dryness which makes sex painful. Bad sex has effected other parts of my life.

Overall, based on the results of our study, it is plausible to anticipate that the systemic administration of hormones might lead to a greater incidence of side effects and lower satisfaction levels. This expectation is corroborated by our data, not only describing side effects with regards to irregular bleeding patterns, skin issues, and loss of libido, but also an increased frequency of complaints such as weight gain and mental health issues associated with these hormonal methods.

Discussion of Insights on LARCs

In our exploratory study, we observe that, on average, women are highly satisfied with their IUDs. In fact, among all contraceptive methods, IUDs are given the highest average ratings on the drug review website (Table 3). This finding is corroborated by existing research indicating high satisfaction levels with this contraceptive method [56]. IUDs offer substantial advantages. The hormonal IUD is known for its ability to significantly reduce menstrual bleeding, with the 52 mg version being approved by the Food and Drug Administration in the United States for both contraception and heavy menstrual bleeding treatment. On the other hand, the copper IUD stands out as the sole nonhormonal choice that offers the convenience of not requiring action during each sexual encounter. Although both types of IUDs can cause undesirable bleeding-related side effects-typically breakthrough bleeding with the hormonal IUD and heavy periods with the copper IUD—these decrease over time [54,61]. We can reasonably assume that women opting for an IUD, which necessitates a medical procedure for insertion, would be well-informed and prepared for this.

However, our research indicates that for a limited group of women, IUDs appear to create major problems. Between 11.3% (90/800) and 12.37% (355/2871) of all written online reviews emphatically describe cramps and pain related to the insertion procedure or persisting pain. In fact, cramps and pain are the dominant topic in 41.9% (355/848) and 38.5% (90/234) of unfavorable reviews of copper and hormonal IUDs, respectively (Figure 3). We also see that, on average, the ratings of reviews

```
https://ai.jmir.org/2025/1/e68809
```

where cramps and pain are the dominant topic are the second lowest, occurring only slightly above multiple cause dissatisfaction (Figure 4). This observation has substantial implications for IUD counseling practices. Although physicians typically inform women about the potential side effects of IUDs, our findings underscore the necessity for health care professionals to provide even more comprehensive counseling regarding the risk of temporary as well as lasting cramps and pain, which heavily hampers women's well-being, and to offer pain control options for the insertion procedure.

Limitations

Our study is subject to several limitations. First, the reviews and ratings in online forums may exhibit bias, often skewing toward negative experiences, and there may even be a risk of fake reviews. Consequently, our dataset may not accurately represent the broader population of women using birth control products. Nonetheless, this limitation does not detract from our study's objective, which is to illuminate the experiences with different contraceptive methods women share on a drug review website. This study was intended to supplement traditional qualitative but informal information sources used by women and their partners. Consequently, our extensive analysis of nearly 20,000 online reviews arguably offers a more representative and robust overview than anecdotal evidence gathered from conversations with friends and family regarding birth control options. All the same, we note that the reliance on data from a single source (ie, the website Drugs.com) may introduce a bias. This could affect the findings. Although the drug review dataset dates from 2009 to 2017, reviewers might have already been influenced by social media influencers, who increasingly expressing concerns about hormonal are contraceptives, sometimes inflating the severity of side effects [36], and advocating for nonhormonal methods. Furthermore, the dataset only contains reviews of products. As such, natural contraceptive methods, such as calendar rhythm and withdrawal are not covered. Performing the NLP and sentiment analysis on other contraceptive product review websites could enhance the breadth and robustness of our findings but is outside the scope of this analysis.

Second, our categorization did not stratify by dose or regimen timing (eg, 21 vs 24 active pills) due to insufficient review information (eg, Loestrin could refer to multiple different products with the same active ingredients in different doses and durations), nor by progestin type to avoid small group sizes.

Third, there are important limitations inherent in the topic modeling of birth control product online reviews. While topic modeling offers valuable insights, it is crucial to acknowledge its constraints so that they can be addressed effectively in future research. One of the primary limitations is the subjectivity involved in choosing the right number of topics. In addition, topic modeling may not adequately capture rare or nuanced topics. In our study, which identified 8 topics, we observed 1 particularly ambiguous topic, "Multiple cause dissatisfaction." This topic is frequently associated with vaginal rings and hormonal patches and occasionally encompasses other topics, potentially obscuring the clarity and precision of our overall results. Conversely, our algorithm effectively differentiates

between "menstrual irregularities" and "continuous bleeding," despite their similarities. Notably, women experience "continuous bleeding" as more problematic than "menstrual irregularities." However, due to the overlap in these side effects and the associated words, some reviews scored highly for both topics.

Another limitation in topic modeling is the potential ambiguity in allocating reviews to specific topics, which stems from the inherent challenge of accurately capturing the thematic essence of the text. For example, a word such as "skin" in an unfavorable review does not necessarily imply a discussion about skin-related problems. Furthermore, one review may describe several topics or side effects (Multimedia Appendix 2). Thus, analyses that are based on reviews grouped by dominant topic may not fully reflect other potentially confounding aspects.

Despite these challenges, our analysis suggests that using TF-IDF and NMF for topic modeling with 8 topics is the most effective approach. In this setup, most topics, apart from "multiple cause dissatisfaction" and the occasionally intersecting "menstrual irregularities" and "continuous bleeding," are well-defined by distinct symptom sets, differentiating them from others. The process involved careful consideration of both the interpretability and the distinctiveness of each topic.

Finally, given the exploratory nature of our study, we did not engage in statistical significance testing; consequently, we could not definitively determine whether the observed differences in contraceptive method ratings are systematic. This study was primarily descriptive and does not involve inferential statistical analysis or controlling for confounding variables. We also did not investigate potential interactions between different side effects, such as whether reports of mental health issues could lead to more negative evaluations of other symptoms. The suitability of the dataset for inferential statistics is questionable, as it does not meet several crucial assumptions for significance tests, such as normality, homoscedasticity within groups, or independence of observations.

In summary, our findings provide semiqualitative insights, highlighting the occurrence of certain side effects in the real world and how they are associated with online contraceptive product ratings. A deeper understanding of effect sizes, relationships, and causality requires further research.

Conclusions

This study contributes to the understanding of how contraceptive methods impact women's overall well-being, as interpreted from a large corpus of online user narratives. Our findings provide a complementary perspective to those derived from clinical trials or the adverse effects documented in pharmaceutical labels and package inserts. By leveraging NLP to analyze user reviews, we aimed to support women in choosing contraceptive options that are not only safe and effective but also reduce the likelihood of specific symptom clusters that could negatively affect their quality of life. For instance, women seeking contraception may have specific concerns, such as potential effects on libido, skin health, or menstrual regularity. While no contraceptive option is completely free of side effects, it is crucial that women have access to information that enables them to make informed decisions about which side effects they are prepared to accept, guided by the experiences of others. Accordingly, our analysis empowers women to benefit from the collective insights and experiences of a large user base, supporting more informed decision-making. In addition, this information aids HCPs in offering personalized advice to women and their partners.

A key observation from our study is that all female contraceptive methods reviewed online are associated with a substantial percentage of negative ratings. Notably, no contraceptive method to be administered by women received <29% unfavorable evaluations. This finding underscores a significant opportunity for enhancement in the realm of female contraception. The objective for manufacturers would be to innovate and develop contraceptive methods that exert minimal or no negative impact on the well-being and quality of life. In light of this, there is a recent trend toward natural or calendar-based methods sometimes supported by digital cycle tracking tools. Depending on individual life circumstances and personal beliefs, these methods may constitute a viable alternative for some women despite the inherent risks resulting from use failure (the website Drugs.com [43] provides a comparison of efficacy and typical use failure rates). Greater awareness of the side effects of contraceptive products for women could guide couples in making joint decisions about contraception and a more equitable sharing of responsibilities by considering more options. A secondary insight from our study is that dissatisfaction was particularly common for contraceptive products that may result in irregular or continuous bleeding, especially when users may have expected reduced or absent menstrual bleeding. IUDs are generally rated more positively than other methods, although about 1 in 10 users report severe cramps and pain which are linked to very poor ratings.

In conclusion, a pivotal element of efficacious reproductive management is the provision of comprehensive information to women regarding the potential side effects of contraceptives and their likely impact on overall well-being and quality of life. Advancements in artificial intelligence in general and NLP in particular can help in extracting, aggregating, interpreting, and sharing this information. In a broader context, the empowerment of women to manage their reproductive health is acknowledged as a fundamental catalyst for economic advancement and the achievement of personal and professional aspirations, as emphasized by organizations, such as the United Nations, the World Health Organization, and the Organisation for Economic Cooperation and Development. Moreover, female sexuality transcends the dimensions of reproduction and birth control, encompassing aspects of pleasure and human connection, thereby enhancing overall well-being [62]. Access to suitable contraceptive methods and thorough information about these options are vital in facilitating this empowerment.



Authors' Contributions

Author AN is currently not affiliated with any academic institution but contributed to this article as an Independent Researcher.

Conflicts of Interest

None declared.

Multimedia Appendix 1 Examples of reviews by dominant topic. [DOCX File , 21 KB - ai v4i1e68809 app1.docx]

Multimedia Appendix 2 Three example reviews and their topic values. [DOCX File , 16 KB - ai v4i1e68809 app2.docx]

References

- 1. Contraceptive use by method 2019: data booklet. United Nations. URL: <u>https://www.un.org/development/desa/pd/sites/</u> www.un.org.development.desa.pd/files/files/documents/2020/Jan/un_2019_contraceptiveusebymethod_databooklet.pdf [accessed 2024-11-10]
- 2. Contraceptive use in the United States by method. Guttmacher Institute. URL: <u>https://www.guttmacher.org/fact-sheet/</u> <u>contraceptive-method-use-united-states</u> [accessed 2024-11-10]
- 3. Daniels K, Abma JC. Current contraceptive status among women aged 15-49: United States. Centers for Disease Control and Prevention. URL: <u>https://www.cdc.gov/nchs/products/databriefs/db388.htm</u> [accessed 2024-11-10]
- 4. Burke HM, Ridgeway K, Murray K, Mickler A, Thomas R, Williams K. Reproductive empowerment and contraceptive self-care: a systematic review. Sex Reprod Health Matters 2021 Jul 27;29(2):2090057 [FREE Full text] [doi: 10.1080/26410397.2022.2090057] [Medline: 35892261]
- 5. Hee L, Kettner LO, Vejtorp M. Continuous use of oral contraceptives: an overview of effects and side-effects. Acta Obstet Gynecol Scand 2013 Feb 05;92(2):125-136 [FREE Full text] [doi: 10.1111/aogs.12036] [Medline: 23083413]
- 6. Shimoni N. Intrauterine contraceptives: a review of uses, side effects, and candidates. Semin Reprod Med 2010 Mar 29;28(2):118-125. [doi: 10.1055/s-0030-1248136] [Medline: 20352561]
- 7. Your contraception guide. NHS England. URL: https://www.nhs.uk/conditions/contraception/ [accessed 2024-01-09]
- 8. Boozalis A, Tutlam NT, Robbins CC, Peipert J. Sexual desire and hormonal contraception. Obstet Gynecol 2016 Mar;127(3):563-572 [FREE Full text] [doi: 10.1097/AOG.00000000001286] [Medline: 26855094]
- Summary report on proceedings, minutes and final acts of the International Health Conference held in New York from 19 June to 22 July 1946. World Health Organization. URL: <u>https://iris.who.int/handle/10665/85573</u> [accessed 2024-04-29]
- 10. Zethraeus N, Dreber A, Ranehill E, Blomberg L, Labrie F, von Schoultz B, et al. A first-choice combined oral contraceptive influences general well-being in healthy women: a double-blind, randomized, placebo-controlled trial. Fertil Steril 2017 May;107(5):1238-1245 [FREE Full text] [doi: 10.1016/j.fertnstert.2017.02.120] [Medline: 28433366]
- Barden-O'Fallon J, Speizer I, Rodriguez F, Calix J. Experience with side effects among users of injectables, the IUD, and oral contraceptive pills in four urban areas of Honduras. Health Care Women Int 2009 Jun 26;30(6):475-483. [doi: 10.1080/07399330902801187] [Medline: 19418321]
- Bellizzi S, Mannava P, Nagai M, Sobel HL. Reasons for discontinuation of contraception among women with a current unintended pregnancy in 36 low and middle-income countries. Contraception 2020 Jan;101(1):26-33. [doi: <u>10.1016/j.contraception.2019.09.006</u>] [Medline: <u>31655068</u>]
- Madden T, Secura GM, Nease RF, Politi MC, Peipert JF. The role of contraceptive attributes in women's contraceptive decision making. Am J Obstet Gynecol 2015 Jul;213(1):46.e1-46.e6 [FREE Full text] [doi: 10.1016/j.ajog.2015.01.051] [Medline: 25644443]
- Grimes DA. The safety of oral contraceptives: epidemiologic insights from the first 30 years. Am J Obstet Gynecol 1992 Jun;166(6 Pt 2):1950-1954. [doi: 10.1016/0002-9378(92)91394-p] [Medline: 1605284]
- 15. Rocca ML, Palumbo AR, Visconti F, Di Carlo C. Safety and benefits of contraceptives implants: a systematic review. Pharmaceuticals (Basel) 2021 Jun 08;14(6):548 [FREE Full text] [doi: 10.3390/ph14060548] [Medline: 34201123]
- D'Souza P, Bailey JV, Stephenson J, Oliver S. Factors influencing contraception choice and use globally: a synthesis of systematic reviews. Eur J Contracept Reprod Health Care 2022 Oct 03;27(5):364-372 [FREE Full text] [doi: 10.1080/13625187.2022.2096215] [Medline: 36047713]
- Cooke-Jackson A, Rubinsky V, Gunning JN. "Wish I would have known that before I started using it": contraceptive messages and information seeking among young women. Health Commun 2023 Apr 20;38(4):834-843. [doi: 10.1080/10410236.2021.1980249] [Medline: 34544296]

- Johansson L, Vesström J, Alehagen S, Kilander H. Women's experiences of dealing with fertility and side effects in contraceptive decision making: a qualitative study based on women's blog posts. Reprod Health 2023 Jun 29;20(1):98 [FREE Full text] [doi: 10.1186/s12978-023-01642-8] [Medline: <u>37381022</u>]
- 19. SPRINTEC- norgestimate and ethinyl estradiol package insert 2011. Rebel Distributors Corp. URL: <u>https://tinyurl.com/</u> <u>57mxt5uk</u> [accessed 2024-11-10]
- 20. Label: SPRINTEC- norgestimate and ethinyl estradiol kit. Dailymed. URL: <u>https://dailymed.nlm.nih.gov/dailymed/drugInfo.</u> <u>cfm?setid=d9252820-131a-4870-8b11-945d1bfd5659</u> [accessed 2024-01-09]
- Simons G, Baldwin DS. A critical review of the definition of 'wellbeing' for doctors and their patients in a post COVID-19 era. Int J Soc Psychiatry 2021 Dec 09;67(8):984-991 [FREE Full text] [doi: 10.1177/00207640211032259] [Medline: 34240644]
- 22. Condon JT, Need JA, Fitzsimmons D, Lucy S. University students' subjective experiences of oral contraceptive use. J Psychosom Obstet Gynaecol 1995 Mar 07;16(1):37-43 [FREE Full text] [doi: 10.3109/01674829509025655] [Medline: 7787956]
- Inoue K, Barratt A, Richters J. Does research into contraceptive method discontinuation address women's own reasons? A critical review. J Fam Plann Reprod Health Care 2015 Oct 20;41(4):292-299 [FREE Full text] [doi: 10.1136/jfprhc-2014-100976] [Medline: 25605480]
- 24. Martell S, Marini C, Kondas CA, Deutch AB. Psychological side effects of hormonal contraception: a disconnect between patients and providers. Contracept Reprod Med 2023 Jan 17;8(1):9 [FREE Full text] [doi: 10.1186/s40834-022-00204-w] [Medline: 36647102]
- 25. Akers AY, Gold MA, Borrero S, Santucci A, Schwarz EB. Providers' perspectives on challenges to contraceptive counseling in primary care settings. J Womens Health 2010 Jun;19(6):1163-1170. [doi: 10.1089/jwh.2009.1735]
- 26. Pfender EJ, Caplan SE. The effect of social media influencer warranting cues on intentions to use non-hormonal contraception. Health Commun (Forthcoming) 2024 Sep 11:1-15. [doi: <u>10.1080/10410236.2024.2402161</u>] [Medline: <u>39258763</u>]
- 27. Merz AA, Gutiérrez-Sacristán A, Bartz D, Williams NE, Ojo A, Schaefer KM, et al. Population attitudes toward contraceptive methods over time on a social media platform. Am J Obstet Gynecol 2021 Jun;224(6):597.e1-597.14 [FREE Full text] [doi: 10.1016/j.ajog.2020.11.042] [Medline: 33309562]
- 28. Pleasants E, Ryan JH, Ren C, Prata N, Gomez AM, Marshall C. Exploring language used in posts on r/birthcontrol: case study using data from reddit posts and natural language processing to advance contraception research. J Med Internet Res 2023 Jun 30;25:e46342. [doi: 10.2196/46342]
- Huang M, Gutiérrez-Sacristán A, Janiak E, Young K, Starosta A, Blanton K, et al. Contraceptive content shared on social media: an analysis of Twitter. Contracept Reprod Med 2024 Feb 07;9(1):5 [FREE Full text] [doi: 10.1186/s40834-024-00262-2] [Medline: 38321582]
- 30. Green EP, Whitcomb A, Kahumbura C, Rosen JG, Goyal S, Achieng D, et al. "What is the best method of family planning for me?": a text mining analysis of messages between users and agents of a digital health service in Kenya. Gates Open Res 2019 May 29;3:1475 [FREE Full text] [doi: 10.12688/gatesopenres.12999.1] [Medline: 31410395]
- Balakrishnan P, Kroiss C, Keskes T, Friedrich B. Perception and use of reversible contraceptive methods in Germany: a social listening analysis. Womens Health (Lond) 2023 Jan 15;19:17455057221147390 [FREE Full text] [doi: 10.1177/17455057221147390] [Medline: <u>36642972</u>]
- 32. Felice MC, Søndergaard ML, Balaam M. Analyzing user reviews of the first digital contraceptive: mixed methods study. J Med Internet Res 2023 Nov 14;25:e47131 [FREE Full text] [doi: 10.2196/47131] [Medline: 37962925]
- Pfender EJ, Devlin MM. What do social media influencers say about birth control? A content analysis of YouTube vlogs about birth control. Health Commun 2023 Dec 15;38(14):3336-3345. [doi: <u>10.1080/10410236.2022.2149091</u>] [Medline: <u>36642835</u>]
- Yeo TE, Chu TH. Sharing "sex secrets" on Facebook: a content analysis of youth peer communication and advice exchange on social media about sexual health and intimate relations. J Health Commun 2017 Sep 10;22(9):753-762. [doi: 10.1080/10810730.2017.1347217] [Medline: 28796578]
- Stoddard RE, Pelletier A, Sundquist EN, Haas-Kogan ME, Kassamali B, Huang M, et al. Popular contraception videos on TikTok: an assessment of content topics. Contraception 2024 Jan;129:110300. [doi: <u>10.1016/j.contraception.2023.110300</u>] [Medline: <u>37802460</u>]
- 36. Black KI, Vromman M, French RS. Common myths and misconceptions surrounding hormonal contraception. Best Pract Res Clin Obstet Gynaecol 2025 Feb;98:102573 [FREE Full text] [doi: 10.1016/j.bpobgyn.2024.102573] [Medline: 39705740]
- 37. Tul Q, Ali M, Riaz A, Noureen A, Kamranz M, Hayat B, et al. Sentiment analysis using deep learning techniques: a review. Int J Multimed Inf Retr 2017;8(6):41. [doi: 10.14569/ijacsa.2017.080657]
- Barbounaki SG, Gourounti K, Sarantaki A. Advances of sentiment analysis applications in obstetrics/gynecology and midwifery. Mater Sociomed 2021 Sep;33(3):225-230 [FREE Full text] [doi: 10.5455/msm.2021.33.225-230] [Medline: 34759782]
- 39. Sari NP, Munir A, At MR, Iskandar M. Twitter sentiment analysis of long-acting reversible contraceptives (LARC) methods in Indonesia with machine learning approach. In: Proceedings of the International Conference on Multidisciplinary Studies (ICoMSi 2023). Cambridge, MA: Atlantis Press; 2024:167-185.

```
https://ai.jmir.org/2025/1/e68809
```

- 40. Bhardwaj R, Majumder N, Poria S. Investigating gender bias in BERT. Cogn Comput 2021 May 20;13(4):1008-1018. [doi: 10.1007/s12559-021-09881-2]
- 41. Leteno T, Gourru A, Laclau C, Gravier C. An investigation of structures responsible for gender bias in BERT and DistilBERT. In: Proceedings of the 21st International Symposium on Intelligent Data Analysis on Advances in Intelligent Data Analysis XXI. 2023 Presented at: IDA'23; April 12-14, 2023; Louvain-la-Neuve, Belgium p. 249-261 URL: <u>https://link.springer.com/ chapter/10.1007/978-3-031-30047-9_20</u> [doi: 10.1007/978-3-031-30047-9_20]
- 42. Chen ES, Hripcsak G, Xu H, Markatou M, Friedman C. Automated acquisition of disease-drug knowledge from biomedical and clinical documents: an initial study. J Am Med Inform Assoc 2008 Jan 01;15(1):87-98. [doi: <u>10.1197/jamia.m2401</u>]
- 43. Birth control failure rates the Pearl Index explained. Drugs.com. URL: <u>https://www.drugs.com/medical-answers/</u> <u>birth-control-failure-rates-pearl-index-explained-3554953/</u> [accessed 2024-11-12]
- 44. Kallumadi S, Grer F. Drug review dataset (Drugs.com). UCI Machine Learning Repository. URL: <u>https://archive.ics.uci.edu/</u> <u>dataset/461/drug+review+dataset+druglib+com</u> [accessed 2024-04-29]
- 45. Gräßer F, Kallumadi S, Malberg H, Zaunseder S. Aspect-based sentiment analysis of drug reviews applying cross-domain and cross-data learning. In: Proceedings of the 2018 International Conference on Digital Health. 2018 Presented at: DH '18; April 23-26, 2018; Lyon, France p. 121-125 URL: <u>https://dl.acm.org/doi/10.1145/3194658.3194677</u> [doi: 10.1145/3194658.3194677]
- 46. Methods of contraception. NHS England. URL: <u>https://www.nhs.uk/contraception/methods-of-contraception/</u> [accessed 2025-02-20]
- 47. Natural language toolkit. NLTK Project. URL: <u>https://www.nltk.org/</u> [accessed 2024-01-12]
- 48. Machine learning in Python documentation. scikit-learn. URL: https://scikit-learn.org/stable/ [accessed 2024-01-12]
- 49. Matplotlib: visualization with Python. Matplotlib. URL: <u>https://matplotlib.org/</u> [accessed 2024-01-12]
- 50. Waskom M. seaborn: statistical data visualization. J Open Source Softw 2021 Apr;6(60):3021. [doi: 10.21105/joss.03021]
- 51. Turan S, Yildiz K, Büyüktanir B. Comparison of LDA, NMF and BERTopic topic modeling techniques on Amazon product review dataset: a case study. In: Proceedings of the 2023 Selected Papers from the International Conference on Computing, IoT and Data Analytics. 2023 Presented at: ICCIDA '23; August 11-12, 2023; La Mancha, Spain p. 23-31 URL: <u>https://link.springer.com/chapter/10.1007/978-3-031-53717-2_3</u> [doi: <u>10.1007/978-3-031-53717-2_3</u>]
- 52. Askalidis G, Kim SJ, Malthouse EC. Understanding and overcoming biases in online review systems. Decis Support Syst 2017 May;97:23-30. [doi: 10.1016/j.dss.2017.03.002]
- American College of Obstetricians and Gynecologists' Committee on Clinical Consensus–Gynecology. General approaches to medical management of menstrual suppression: ACOG clinical consensus no. 3. Obstet Gynecol 2022 Sep 01;140(3):528-541. [doi: 10.1097/AOG.00000000004899] [Medline: 36356248]
- 54. Hillard PA. Menstrual suppression: current perspectives. Int J Womens Health 2014 Jun;6:631-637 [FREE Full text] [doi: 10.2147/IJWH.S46680] [Medline: 25018654]
- 55. Gracia CR, Sammel MD, Charlesworth S, Lin H, Barnhart KT, Creinin MD. Sexual function in first-time contraceptive ring and contraceptive patch users. Fertil Steril 2010 Jan;93(1):21-28 [FREE Full text] [doi: 10.1016/j.fertnstert.2008.09.066] [Medline: 18976753]
- Arowojolu AO, Gallo MF, Lopez LM, Grimes DA. Combined oral contraceptive pills for treatment of acne. Cochrane Database Syst Rev 2012 Jul 11;2012(7):CD004425. [doi: <u>10.1002/14651858.CD004425.pub6</u>] [Medline: <u>22786490</u>]
- 57. Passeron T. Melasma pathogenesis and influencing factors an overview of the latest research. J Eur Acad Dermatol Venereol 2013 Jan;27 Suppl 1:5-6. [doi: 10.1111/jdv.12049] [Medline: 23205539]
- 58. Sood R, Faubion SS, Kuhle CL, Thielen JM, Shuster LT. Prescribing menopausal hormone therapy: an evidence-based approach. Int J Womens Health 2014;6:47-57 [FREE Full text] [doi: 10.2147/IJWH.S38342] [Medline: 24474847]
- 59. Mosorin ME, Piltonen T, Rantala AS, Kangasniemi M, Korhonen E, Bloigu R, et al. Oral and vaginal hormonal contraceptives induce similar unfavorable metabolic effects in women with PCOS: a randomized controlled trial. J Clin Med 2023 Apr 12;12(8):2827 [FREE Full text] [doi: 10.3390/jcm12082827] [Medline: 37109164]
- 60. Abdollahpour S, Ashrafizaveh A, Azmoude E. Effects of the combined contraceptive vaginal ring on female sexual function: a systematic review and meta-analysis. Malays J Med Sci 2023 Feb 28;30(1):21-30 [FREE Full text] [doi: 10.21315/mjms2023.30.1.3] [Medline: 36875197]
- 61. Sanders JN, Adkins DE, Kaur S, Storck K, Gawron LM, Turok DK. Bleeding, cramping, and satisfaction among new copper IUD users: a prospective study. PLoS One 2018 Nov 7;13(11):e0199724 [FREE Full text] [doi: 10.1371/journal.pone.0199724] [Medline: 30403671]
- 62. Davison SL, Bell RJ, LaChina M, Holden SL, Davis SR. The relationship between self-reported sexual satisfaction and general well-being in women. J Sex Med 2009 Oct;6(10):2690-2697. [doi: 10.1111/j.1743-6109.2009.01406.x] [Medline: 19817981]

Abbreviations

CHC: combined hormonal contraceptive **COCP:** combined oral contraceptive pill

https://ai.jmir.org/2025/1/e68809

HCP: health care provider
IUD: intrauterine device
LARC: long-acting reversible contraceptive
NLP: natural language processing
NMF: nonnegative matrix factorization
OCP: oral contraceptive pill
POP: progestin-only pill
RQ: research question
SHBG: sex hormone–binding globulin
TF-IDF: term-frequency–inverse document frequency
UCI: University of California, Irvine

Edited by H Liu; submitted 15.11.24; peer-reviewed by EJ Pfender, T Awofala, S Oworah; comments to author 11.02.25; revised version received 25.02.25; accepted 07.03.25; published 03.04.25.

<u>Please cite as:</u> Groene N, Nickel A, Rohn AE Insights on the Side Effects of Female Contraceptive Products From Online Drug Reviews: Natural Language Processing–Based Content Analysis JMIR AI 2025;4:e68809 URL: <u>https://ai.jmir.org/2025/1/e68809</u> doi:<u>10.2196/68809</u> PMID:

©Nicole Groene, Audrey Nickel, Amanda E Rohn. Originally published in JMIR AI (https://ai.jmir.org), 03.04.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Trust, Anxious Attachment, and Conversational AI Adoption Intentions in Digital Counseling: A Preliminary Cross-Sectional Questionnaire Study

Xiaoli Wu¹, BSc (Hons); Kongmeng Liew¹, PhD; Martin J Dorahy¹, PhD, DClinPsych

School of Psychology, Speech and Hearing, University of Canterbury, Christchurch, New Zealand

Corresponding Author: Xiaoli Wu, BSc (Hons) School of Psychology, Speech and Hearing University of Canterbury Private Bag 4800 Christchurch, 8140 New Zealand Phone: 64 02059037078 Email: xwu40@uclive.ac.nz

Abstract

Background: Conversational artificial intelligence (CAI) is increasingly used in various counseling settings to deliver psychotherapy, provide psychoeducational content, and offer support like companionship or emotional aid. Research has shown that CAI has the potential to effectively address mental health issues when its associated risks are handled with great caution. It can provide mental health support to a wider population than conventional face-to-face therapy, and at a faster response rate and more affordable cost. Despite CAI's many advantages in mental health support, potential users may differ in their willingness to adopt and engage with CAI to support their own mental health.

Objective: This study focused specifically on dispositional trust in AI and attachment styles, and examined how they are associated with individuals' intentions to adopt CAI for mental health support.

Methods: A cross-sectional survey of 239 American adults was conducted. Participants were first assessed on their attachment style, then presented with a vignette about CAI use, after which their dispositional trust and subsequent adoption intentions toward CAI counseling were surveyed. Participants had not previously used CAI for digital counseling for mental health support.

Results: Dispositional trust in artificial intelligence emerged as a critical predictor of CAI adoption intentions (P<.001), while attachment anxiety (P=.04), rather than avoidance (P=.09), was found to be positively associated with the intention to adopt CAI counseling after controlling for age and gender.

Conclusions: These findings indicated higher dispositional trust might lead to stronger adoption intention, and higher attachment anxiety might also be associated with greater CAI counseling adoption. Further research into users' attachment styles and dispositional trust is needed to understand individual differences in CAI counseling adoption for enhancing the safety and effectiveness of CAI-driven counseling services and tailoring interventions.

Trial Registration: Open Science Framework; https://osf.io/c2xqd

(JMIR AI 2025;4:e68960) doi:<u>10.2196/68960</u>

KEYWORDS

attachment style; conversational artificial intelligence; CAI; perceived trust; adoption intentions; CAI counseling; mobile phone

Introduction

Conversational Artificial Intelligence in Mental Health

Conversational artificial intelligence (CAI) has rapidly captured global attention since its emergence in recent years. It has permeated various facets of human life and continues to attract

https://ai.jmir.org/2025/1/e68960

RenderX

a growing user base worldwide due to its unparalleled impact on the way people access knowledge, present ideas, and interact. Commercially available CAIs, including Replika (developed by Luka Inc, Replika is a chatbot designed to be a conversational agent and personal companion, using artificial intelligence (AI) to simulate human-like conversations; its primary purpose is to provide users with an AI friend that can listen, respond

empathetically, and help users reflect on their thoughts and feelings. It is often used for mental health support, companionship, and improving emotional well-being [1]) and Pi (Developed by Inflection AI, Pi is a CAI designed to provide a range of task-based features, including emotional support, learning assistance, and personalized interactions; it is specifically tailored to engage users in meaningful conversations, making it useful for various purposes, such as mental health support, learning new languages, and relationship advice), are powered by large language models (LLMs) with deep learning-based natural language processing to enable human-like voice or text interactions with users. They offer a wide range of services such as information retrieval, task completion, entertainment, and even mental health support [2]. In the context of mental health support, CAI is used in various counseling settings like delivering psychotherapy, providing psychoeducational content, and offering support such as companionship or emotional aid [3]. In this paper, we define CAIs as chatbots that use LLMs to generate naturalistic text, which is different from traditional rule-based conversational agents that operate mainly on predefined scripts, such as customer-oriented chatbots commonly used in sales and marketing.

One increasingly common usage of these anthropomorphic CAIs has been for counseling purposes in mental health settings to improve the overall quality of communications [4]. Gaffney et al [5] conducted a systematic review of 13 studies on the application of conversational agents (including CAIs) in psychotherapeutic settings and found that overall, CAIs showed promising results in terms of effectiveness and acceptability for addressing mental health issues in users. More recently, Li et al [6] conducted a meta-analysis of 15 randomized controlled trials specifically focusing on CAI counseling, and found that CAIs showed a significant decrease in depression and distress symptoms, especially when used with clinical, subclinical, and older adult populations. These findings suggest that CAI has the potential to effectively address mental health issues. Furthermore, the accessibility and user-friendly nature of CAI have also made them a promising tool for delivering mental health care to a wider population at a faster response rate and at an affordable cost, compared with traditional in-person therapies. It offers hope for overcoming long-lasting barriers, such as social stigma and the demand-supply imbalance, that weigh down traditional mental health care services [7].

Despite the benefits CAI counseling could potentially bring to mental health care, it also poses many risks and challenges, such as misleading responses, privacy infringement, and ethical concerns, to name a few [8]. For instance, counseling typically involves a high degree of self-disclosure, which in the context of CAI counseling can be problematic. Users may share sensitive and personal information that, if not properly protected, could be vulnerable to data breaches or misuse. Furthermore, the algorithms used by CAIs might not fully grasp the nuances of human emotions and mental health issues, potentially providing or harmful responses (eg, spreading inappropriate misinformation, professing their love to users, and sexually harassing minors) [9]. In addition, users of CAI counseling may be more susceptible to developing maladaptive behaviors (eg,

addiction) as most counseling CAIs are designed to form social-emotional bonds with its users. While CAI therapies are intended to improve users' psychological well-being, they also risk users developing overreliance and social withdrawal [10]. Without caution in its application and a thorough understanding in human-CAI interaction in counseling settings, the unpredictability in CAI responses could lead to adverse psychological consequences on the user.

How should we weigh the pros and cons of adopting CAI counseling for mental health support? Most of the relevant literature [2,11] acknowledges the significant potential of CAI therapies in providing therapeutic support and underscores the necessity for further exploration and implementation, but also highlights the importance of recognizing and meticulously managing the risks associated with CAI therapies through rigorous research and well-defined guidelines. Furthermore, regardless of the concerns related to the use of CAI for psychological support, there are already CAIs that provide easy access to task-oriented features designed for mental health purposes. For example, a wide range of diverse task-oriented features offered by Pi fall under this category, such as venting, self-care for anxiety, and relationship advice [12]. Particularly, the younger generation may be more open to trying new technologies, making them more vulnerable to potential harms from poorly regulated or non-evidence-based CAI therapies. Therefore, to ensure the safe and effective integration of CAI into mental health services, it is crucial to understand the factors influencing CAI adoption, including potential predictors and barriers. However, research is relatively lacking in this area [7,10]. Studies addressing the factors associated with individual adoption of CAI counseling is needed to comprehend the psychological mechanisms underpinning the formation of human-CAI relationships. This study was designed to address this gap, by examining individual differences in attachment styles and perceived trust as predictors of CAI adoption for mental health care.

Numerous studies have demonstrated attachment style to be a reliable predictor for various relational outcomes [13,14], including the relationship between humans and technology. Meanwhile, trust is considered as another key factor in the context of technology adoption and use, especially in the domain of AI adoption due to risks related to its complexity as mentioned earlier [15]. Therefore, for this study, perceived trust and attachment styles were both examined as potential pertinent variables that might account for individual differences in CAI adaption in the context of digital counseling.

Trust as a Potential Predictor for CAI Counseling Adoption

Based on the theoretical framework developed by McKnight [16], trust is the extent to which a person has confidence in, and is ready to rely on, another entity (in this case, CAI). The formation of trust in information technology goes through different stages, each influenced by distinct factors and mechanisms. Considering CAI as a relatively recent technology, we assume that most individuals would have no previous experiences with CAI counseling. Therefore, the primary focus of this study was on the initial stage of trust building, which

XSL•FO

pertains to establishing trust with an unfamiliar party or service without previous interaction. Numerous studies have examined the individual process of technology adoption from the perspective of trust formation, where the entity being trusted is a technology such as an information system or a recommendation agent. For example, when examining the factors that influence digital voice assistant use, Fernandes and Oliveira [17] found a positive link with perceived trust. Kasilingam [18] investigated intentions toward using smartphone chatbots for purchasing decisions and found that trust positively influenced participants' willingness to use chatbots for mobile shopping purposes. However, studies have not yet examined trust as a predictor for the adoption intention of CAI counseling for psychological support before engagement. While a recent review [19] identifies trust as a key predictor of CAI adoption in health care, including mental health care settings, the CAIs discussed in this review reflect a broader, more general definition of CAIs, including those that use prewritten scripts, which fall outside the scope of our research. Research specifically studying the relationship between trust and adoption intention of advanced CAI counseling (eg, ChatGPT [OpenAI] relying on contemporary reinforcement -learning with human feedback-based LLM technology) is still lacking. Additional research is needed to examine the replicability and reliability of these conclusions within the context of advanced CAI counseling technologies. Furthermore, given that CAI counseling for psychological support involves deeper emotional bonding and personal information, trust may play a significantly different role compared with CAI applications for other aspects of mental health care, such as diagnosis or treatment adherence. Studies examining the formation of trust on primitive, pre-LLM chatbot systems have found positive associations between perceived trust and chatbot adoption, which may generalize to explain how initially perceived trust shapes individuals' behaviors in considering the use of CAI counseling. Hence, in this study, we tested whether perceived trust can predict CAI counseling adoption.

Attachment Theory and Styles

Attachment theory, initially developed by John Bowlby, is a psychological framework that describes how infants learn to interact with their caregivers [20-23]. It was later expanded and adapted to explain the dynamics of both long and short-term interpersonal relationships between humans [24]. A key concept within this theory is the idea of "internal working models (IWMs)," which are shaped by early interactions with primary caregivers. The nature of these interactions, whether they are nurturing, inconsistent, or neglectful, greatly influences the types of IWMs developed. When a caregiver consistently responds to a child's needs in a caring, supportive manner, it tends to foster a positive IWM, while inconsistent or neglectful nurturing is more likely to lead to the formation of negative IWMs. These IWMs serve as mental templates that individuals use to perceive themselves and others, and influence their attributions, perceptions, and emotional understandings of these connections. In essence, they tend to serve as a prototype to determine an individual's expectations and behaviors in subsequent relationships [25-27].

Attachment styles are commonly presented as secure attachment, anxious attachment, avoidant attachment, and disorganized attachment. However, disorganized attachment is often viewed as the most unpredictable type due to its lack of organization in how the child (and later adult) responds to their attachment figures, characterized by push-pull dynamics that lead to inconsistent and conflicted coping strategies [28]. This variability makes it challenging to draw reliable and accurate measurements. For that reason, disorganized attachment was not examined in this study.

According to Bretherton and Munholland [29], attachment style can be understood as the manifestation of people's underlying IWMs. The IWM of attachment avoidance is thought to manifest a positive view of self (as worthy of love and nurturance) and a negative view of others (as unresponsive and untrustworthy). Conversely, attachment anxiety is thought to be associated with an IWM that contains a negative view of self and a positive view of others. Finally, secure attachment is the combination of positive views of both self and others. Securely attached individuals are more capable of forming and maintaining close relationships, with higher commitment, intimacy, love, and satisfaction in such relationships. As for the two insecure attachment styles, avoidant attachment is defined by devaluation of the importance of close relationships, avoidance of intimacy and dependence, and decreased engagement in attachment behavior, while anxious attachment involves preoccupation with the availability and responsiveness of attached figures, fear of separation, and abandonment [24,30].

Attachment Insecurity as a Potential Predictor for CAI Counseling Adoption

While attachment styles are typically associated with interpersonal relationships, Hodge and Gebler-Wolfe [31] found that inanimate objects, such as smartphones, could also be perceived as attachment objects for anxiously attached individuals to feel secure, and reduce unpleasant feelings such as loneliness and boredom. This illustrates how attachment theory can be used as a framework to understand a person's relationship with technology. Beyond smartphones, Birnbaum et al [32] found that humans desire the presence of robots in stressful circumstances in a similar manner to their proximity-seeking behavior toward human attachment figures, suggesting that attachment might also play a similarly important role in human-CAI interactions. Given that CAI is a relatively nascent technology, especially in its application for mental health support, there is, to the best of our knowledge, no previous literature investigating the direct relationships between attachment styles and CAI use. The closest study explored the influence of different attachment styles on perceived trust in broadly defined AI; here Gillath et al [33] found that attachment anxiety was associated with lower trust in AI. Furthermore, participants' trust in AI was reduced when their attachment anxiety was enhanced and increased when their attachment security was boosted. Accordingly, consistent with the positive association between perceived trust and CAI adoption, we expected to find a negative association between attachment anxiety and CAI adoption intention in our study.



Meanwhile, although Gillath et al [33] found no significant effect of attachment avoidance on trust in AI, likely due to the inhibited nature of avoidant individuals, we believe it is important and necessary to include both attachment styles in this study in order to provide comprehensive insights into an underexplored area of research. n this study, we take a broader approach by also hypothesizing the relationship between attachment avoidance and CAI adoption. Insecure attachment styles, including both anxious attachment and avoidant attachment, are generally associated with lower levels of trust. For example, a number of studies on human relations have shown that attachment security is associated with more trust, whereas attachment insecurity is associated with less trust in other humans [34-36]. It may thus be reasonable to hypothesize that higher attachment avoidance will predict lower CAI adoption intention. This hypothesis is grounded in the understanding that individuals with high attachment avoidance may be less inclined to trust, and therefore, are less likely to adopt new technologies like CAI.

Research Hypotheses

Therefore, as a first step toward the eventual aim of promoting safer adoption and designing better CAI for mental health support, this study examined how perceived trust and attachment insecurity (ie, attachment anxiety and attachment avoidance) are associated with CAI adoption intentions. We propose the following two hypotheses:

First (hypothesis 1), due to the positive association found between the perceived trust and primitive chatbots adoption in the previous literature, higher trust in CAI counseling would predict higher adoption intentions for CAI counseling, after controlling for general confounding variables of age and gender.

Second (hypothesis 2), due to the association between insecure attachment styles (ie, anxiety and avoidance) and lower levels of trust, individuals with higher insecure attachments would show lower adoption intentions for CAI counseling, after controlling for age and gender.

To test the above hypotheses, a cross-sectional web-based survey was conducted. As no previous study has examined the human-CAI relationship through the perspective of attachment styles, this preliminary study may provide novel insights into this area and contribute to the existing literature on attachment and technology-mediated relationships. All hypotheses and methods were preregistered before data collection at Open Science Framework [37] and eventual deviations from the preregistration are detailed in Multimedia Appendix 1.

Methods

Participants

Based on an a priori power analysis, a minimum sample size of 146 was recommended to detect an effect size of $F_2=0.075$ with 95% power and alpha at .05 using a linear multiple regression with 6 predictors. The effect size was obtained from the findings of Gritti et al [38] on the effect of avoidant attachment on social network mobile app use. A total of 274 participants (aged 18 y and older, with American nationality or residence) were initially recruited through a large and diverse participant pool from the "Prolific" platform (prolific website; Prolific is a web-based service that provides access to a diverse pool of participants [initially recruited from word-of-mouth and social media] who opt-in to participate in studies listed on the platform. Eligible participants from the Prolific platform are notified through email or their Prolific dashboard. Prolific matches studies to participants based on prescreened criteria. Notifications are presented with necessary information, such as the study title, brief description of the study, estimated time commitment, and payment details clearly displayed, etc) between December 2023 and January 2024. Furthermore, 35 participants' entries were removed due to incomplete data or ineligibility (eg, participants with previous CAI counseling experience as we were only interested in their adoption intention before engagement) responses, leaving a final sample size of 239 participants. The gender ratio of participants was nearly balanced (Table 1). Most of the participants identified as European Americans (153/239, 64% European; 37/239, 15% Asian; 27/239, 11% African American; and 22/239, 8% Native American and others) with a wide age range from 18 to 74 (mean 36.9, SD 12.4) years. For a breakdown of participant demographics, refer to Table 1.



 Table 1. Demographic data breakdown for all participants (N=239).

Variables	Frequency, n (%)
Gender	
Women	114 (47.7)
Men	119 (49.8)
Other	6 (2.5)
Ethnicity	
European (Caucasian)	153 (64)
Asian	37 (15.5)
Black or African American	27 (11.3)
American Indian or Alaska Native	3 (1.3)
Other	17 (7.1)
Native Hawaiian or other Pacific Islander	1 (0.4)
Prefer not to say	1 (0.4)
Education level	
High school or equivalent	29 (12.1)
College or associated degree	63 (26.4)
Bachelor's degree	90 (37.7)
Postgraduate degrees	55 (23)
Others	2 (0.8)
PEDT ^a	
Negative	8 (3.4)
Neutral	49 (20.5)
Positive	182 (76.2)
Familiarity with CAI ^b counseling	
Not familiar at all	116 (48.5)
Slightly familiar	81 (33.9)
Moderately familiar	31 (13)
Very familiar	10 (4.2)
Extremely familiar	1 (0.4)

^aPEDT: previous experience with digital technologies.

^bCAI: conversational artificial intelligence.

Materials and Procedure

Overview

After reading the information sheet and providing consent to participate, participants proceeded to a survey consisting of multiple blocks in a predetermined order (ie, attachment styles, trust toward CAI counseling, intention of use for CAI counseling, and demographic questions). The item order in each scale was randomized to reduce response bias, and an attention check question was included in the survey.

Adult Attachment Style

Adult attachment style was measured using the close relationship version of Revised Adult Attachment scale [39]. This scale contains 2 subscales with 6 items assessing anxious attachment

https://ai.jmir.org/2025/1/e68960

RenderX

(eg, "I often worry that other people don't really love me," Cronbach α =0.91) and the other 12 items measuring avoidant attachment (eg, "I find it difficult to allow myself to depend on others," Cronbach α =0.88). Participants were asked to think about their close relationships with people important to them, such as family members, romantic partners, and close friends, and to rate each statement on a 5-point Likert scale, ranging from 1 (strongly disagree) to 5 (strongly agree). Anxious attachment scores and avoidant attachment scores were computed by taking the average of items within each subscale, with certain items being reverse-scored.

Trust in CAI Counseling

The concept of CAI counseling is still relatively new to most people. In order to introduce its applications in mental health

support, we adapted a news article illustrating the use of CAI in these contexts (Multimedia Appendix 1) for participants to read before completing the survey questions on trust in CAI in the setting of mental health support. This was edited to be as neutral in tone as possible and to remove references to gendered pronouns. A 12-item human-computer trust scale was adapted from previous research [40] (eg, "I think that CAI is competent and effective in providing mental support," Cronbach α =0.94) with each statement rated on a 5-point Likert scale, ranging from 1 (strongly disagree) to 5 (strongly agree). Scores for trust were calculated by averaging the items after reverse-scoring the relevant items.

The vignette plays a crucial role in this study, as it provides a contextual scenario to introduce and illustrate the practical application of CAI in mental health care. Given that CAI is still an emerging technology, particularly in the context of mental health support, the vignette helps bridge potential gaps in participants' understanding. This was especially important in case randomly enrolled participants were unfamiliar with CAI counseling or had never encountered it before.

CAI Counseling Adoption Intentions

CAI adoption intentions for mental health support were measured with a single-item measure, "How likely are you to try a counseling service based on CAI for mental health support in future (if needed)?" using a 5-point Likert-type scale (1=extremely unlikely, 5=extremely likely).

Demographics

Participants were asked to answer demographic questions on their age, gender, and education level. In addition, participants were asked about their previous experience with digital technology on a single-item measure, "How is your previous experience with digital technology in general?" (negative, neutral, or positive), and their familiarity with CAI's counseling function for mental health support on another single-item measure, "How familiar are you with the counseling function of CAI?" A 5-point Likert scale from 1 (not familiar at all) to

Table 2. Descriptive statistics correlations matrix for all variables of interest.

5 (extremely familiar) was adopted. In addition, previous CAI counseling experience was assessed on a single question, "Have you used CAI for mental health support before? If yes, please tell us more about your experience with it (eg, usability, effectiveness, satisfaction, and motivators for first engagement with CAI, etc.) if you would like to share."

Ethical Considerations

This study was approved by the Human Research Ethics Committee at the University of Canterbury, HREC 2023-120-LR. Participants received compensation of GBP 1.00 (US \$1.28) for completing the survey.

All participants were required to carefully read the information sheet, which included essential details such as the research purpose, participation procedure, anonymity assurance, and potential benefits of participation. Participants were informed they could withdraw from the survey at any point. Completion and submission of the survey indicated participants' consent to participate.

Results

Demographics and Descriptive Statistics

Table 1 presents a breakdown of participant demographics. Note that nearly half the participants reported a lack of familiarity with the counseling aspects of CAI.

Descriptive statistics and a correlation matrix among all variables of interest are illustrated in Table 2. Attachment anxiety and avoidance were significantly and positively correlated with each other, which supports the dimensional rather than categorical nature of attachment styles. Furthermore, there was a negative and significant association between age and anxious attachment. Older participants tend to have lower anxiety associated with attachment. A strong significant correlation was found between trust and CAI adoption; higher trust was linked with greater intention of using CAI for mental health support.

Variable	Mean (SD)	A-anxiety ^a	A-avoidance ^b	Trust	CAI ^c adoption
A-anxiety	2.95 (1.11)				
A-avoidance	2.95 (0.78)	.67 ^d			
Trust	2.83 (0.91)	02	10		
CAI adoption	2.80 (1.38)	.06	04	.77 ^e	
Age	36.93 (12.36)	24 ^e	10	.05	.07

^aA-anxiety: attachment anxiety.

^bA-avoidance: attachment avoidance.

^cCAI: conversational artificial intelligence.

 $^{\rm d}P < .01.$

^eP<.001.

RenderX

Confirmatory Results

To test the first hypothesis, a hierarchical ordinal logistic regression was conducted to examine the relationship between

```
https://ai.jmir.org/2025/1/e68960
```

perceived trust in CAI counseling and CAI adoption intention, given that the outcome variable (CAI adoption intention) was a single-item categorical variable. Table 3 reports the full breakdown of results for each model (using standardized

regression coefficients for all predictors). For the first step, age and gender were entered into the model. This step aimed to control for these demographic variables' effects on the outcome variable. Subsequently, the variable of interest—perceived trust—was added into the model to see whether the perceived trust explained significant variance in participants' CAI adoption intention above and beyond the effect of age and gender.

Aligning with hypothesis 1, perceived trust in CAI counseling emerged as a strong predictor of CAI adoption intention (b=2.62, 95% CI 2.19-3.09, P<.001, odds ratio [OR] 13.70). This suggests the higher the trust levels participants have toward CAI's counseling, the more willing they are to use CAI for mental health support, after controlling for age and genders. This tendency is also apparent in the box plot in which perceived trust was plotted against CAI adoption intention in Figure 1.

Considering our aim of examining adoption in initial stages of trust-building with counseling CAI, we conducted a robustness check by repeating the analysis with the subset of participants who reported "not familiar at all" with counseling CAI (n=116). For this sample, perceived trust in CAI counseling was still a strong predictor of CAI adoption intention (b=2.72, 95% CI 2.07-3.45, P<.001, OR 15.10), after controlling for age and gender.

To test the second hypothesis, we repeated the hierarchical ordinal logistic regression analysis to see if attachment insecurity

predicted CAI adoption intention. Age and gender were included in the first step, followed by anxious attachment and avoidant attachment scores as the second step for predicting CAI adoption intention. Full results can be found in Table 4.

In contrast with hypothesis 2, we observed a small but positive significant effect of attachment anxiety on CAI adoption intention when age and gender were controlled (b=0.33, 95% CI 0.02-0.64, P=.04, OR 1.39). This means people with higher attachment anxiety are more likely to adopt CAI for mental support. It is contrary to the direction (ie, negative) that was theorized in hypothesis 2. However, this effect did not appear to be robust, as it was not significant in a zero-order correlation (Table 2). As shown in Figure 2, there was no clear pattern between attachment anxiety and CAI adoption intention before controlling for age and gender. No other significant relationships were found including between attachment avoidance and CAI adoption intentions.

To align with our aim of examining adoption in initial stages of trust-building with counseling CAI, we conducted a robustness check by repeating the analysis with the subset of participants who reported "not familiar at all" with counseling CAI (n=116). For this sample, attachment anxiety significantly predicted of CAI adoption intention (b=0.55, 95% CI 0.09-1.02, P=.02, OR 1.74), but not attachment avoidance (b=-0.54, 95% CI -1.148 to 0.064, P=.08, OR 0.59), after controlling for age and gender.

 Table 3. Regression coefficients for conversational artificial intelligence adoption as a function of multiple variables (N=239).

Predictor variables	Step 1, <i>b</i> (95% CI)	SE	P value	Step 2, b (95% CI)	SE	P value
Age	0.01(-0.01 to 0.03)	0.01	.26	0.001 (-0.02 to 0.02)	0.01	.92
Gender						
Men-Women	0.24 (-0.23 to 0.70)	0.24	.32	-0.35 (-0.88 to 0.17)	0.27	.19
Other-Women	-0.73 (-2.19 to 0.65)	0.71	.30	-0.49 (-2.26 to 1.16)	0.86	.57
Trust in CAI ^a Counseling	b	_	_	2.62 (2.19 to 3.09)	0.23	<.001
R ² _{McF} ^c	0.01	_	_	0.29	_	_
R^2_{McF} change	_	_	_	0.28	_	_

^aCAI: conversational artificial intelligence.

^bNot applicable.

^cR²_{McF}: McFadden's R-squared.



Figure 1. Descriptive box plot illustrates the relationship between perceived trust and CAI adoption intention. CAI: conversational artificial intelligence.



Table 4. Regression coefficients for conversational artificial intelligence adoption as a function of multiple variables (N=239).

Predicting variables	Step 1, <i>b</i> (95% CI)	SE	P value	Step 2 b (95% CI)	SE	P value
Age	0.01 (-0.01 to 0.03)	0.01	.26	0.02 (-0.004 to 0.04)	0.01	.12
Gender						
Men-Women	0.24 (-0.23 to 0.70)	0.24	.32	0.28 (-0.19 to 0.75)	0.24	.24
Other-Women	-0.73 (-2.19 to 0.65)	0.71	.30	-0.64 (-2.12 to 0.76)	0.72	.37
Attachment anxiety	a	_	_	0.33 ^b (0.02 to 0.64)	0.16	.04
Attachment avoidance	_	_	_	-0.36 (-0.77 to 0.06)	0.21	.09
R ² _{McF} ^c	.005	_	_	.012	_	_
R^2_{McF} change	_	_	_	.007	_	_

^aNot applicable.

^b*P*≤.05.

^cR²_{McF}: McFadden's R-squared.



Figure 2. Descriptive box plot illustrating the relationship between attachment anxiety and conversational artificial intelligence adoption intention. AS-anxious: anxious attachment style; CAI: conversational artificial intelligence.



Discussion

Principal Findings

In this study, perceived trust and attachment insecurity (ie, attachment anxiety and attachment avoidance) were examined as factors that influence the dependent variable-CAI adoption intention. In hypothesis 1, we assumed that higher trust in CAI counseling would be associated with a stronger adoption intention, with age and gender controlled for their effects. The results supported this hypothesis, as trust appeared to be a strong predictor of participants' intention to use CAI for mental health support. In addition, in hypothesis 2, anxious attachment and avoidant attachment were proposed to be negatively linked to CAI adoption intention, after controlling for age and gender. Surprisingly, the results did not support this hypothesis. Specifically, avoidant attachment was not a significant predictor of CAI adoption intention, while anxious attachment was found to be a significant predictor with a small effect, but only after controlling for age and gender. Contrary to our original expectation, a greater level of attachment anxiety was found to predict a stronger CAI adoption intention.

Implication of Primary Findings

When it comes to the implementation of a novel but uncertain emerging technology like CAI, it is important to understand users' psychology and resultant behaviors at different stages of interaction, to understand how to achieve safe relationships and positive, effective outcomes. There is a critical distinction in the focus between the pre-engagement stage, such as individual users' intentions to adopt the technology, and the post -engagement stage, such as usage patterns and addiction. As CAI for mental health support has not achieved widespread usage, we focused on the pre-engagement stage in order to examine and describe potential predictors that drive individual engagement with CAI in the context of mental health support.

To our knowledge, this is the first study looking at the relationship between trust and CAI adoption intention for the

RenderX

specific purpose of mental health support. These findings are highly important as they underscore the critical role of trust in the adoption of CAI for mental health support. Given the sensitive nature of mental health, establishing and enhancing trust in CAI systems is paramount. Although many people may not yet be familiar with the potential of CAI in providing mental health support, this may change as CAI becomes more widely accepted and integrated into various fields. In times of urgent need, when human resources are unavailable or delayed, CAI could emerge as a valuable and appealing option for those seeking mental health support, prompting them to explore its potential for engagement. Therefore, user safety, wider acceptance, and use of this technology-all call for developers to prioritize robust security protocols and transparent privacy policies to enhance users' trust, including clear communication about how data are collected, stored, and used. Meanwhile, establishing and adhering to ethical standards is essential. This includes ensuring the AI's recommendations are safe, accurate, and unbiased. Providing users with training and resources to understand how CAI systems work can also demystify the technology and build trust.

Future research should focus on identifying specific factors that build or hinder trust in CAI, particularly in diverse populations, and explore interventions that could mitigate trust-related barriers. In addition, it will be crucial to investigate how trust interacts with other psychological variables, such as attachment styles, to fully understand its role in CAI adoption. Notably, a relatively small effect of attachment anxiety on CAI adoption intention was detected after controlling for age and gender. One possible explanation for the observed effect could be that lower levels of attachment anxiety among older individuals diluted the overall impact of attachment anxiety on adoption intention. Recent research has indicated age as an effective demographic factor to predict AI adoption. For example, Shandilya and Fan [41] found that older adults are less likely to use AI products than younger generations. Similarly, Draxler and colleagues [42] found that early adopters of LLMs, such as ChatGPT, tended not to include individuals from relatively older age

groups. This calls for further research incorporating theoretical frameworks and broader contextual and demographic variables to clarify the roles of age and gender in CAI adoption, particularly in the context of counseling therapies for mental health.

Furthermore, to our surprise, higher attachment anxiety was linked with higher adoption intention. One explanation for this unexpected positive association between attachment anxiety and CAI adoption intention might be the constant and excessive need for validation, reassurance, and emotional support which characterizes anxiously attached individuals [43]. Unlike individuals with attachment avoidance, who tend to suppress or ignore their attachment needs to avoid the discomfort caused by fear of abandonment, those with anxious attachment cope by seeking additional attention and affirmation to alleviate their fears and insecurities. Due to the anthropomorphic, nonjudgmental, constantly accessible, and responsive nature of CAI counseling, anxiously attached individuals might consider CAI as a potential attachment object as well as a secure base, for comfort and reassurance seeking whenever needed. This reassurance-seeking behavioral pattern demonstrated by anxiously attached people was also observed in studies on attachment toward inanimate or nonhuman objects and entities, such as smartphones [31] and robots [32]. On the other hand, attachment anxiety is a key indicator of insecure attachment, with individuals exhibiting lower levels of attachment anxiety generally being more securely attached. This higher sense of security may foster greater confidence in their interpersonal skills, making them more comfortable seeking assistance or support from other individuals, as well as in communicating negative or challenging emotions to others. These may also reduce their need for CAI counseling.

Our findings indicate that CAI could be particularly attractive and beneficial for anxiously attached individuals, potentially filling gaps where traditional support is inaccessible or unavailable. Compared with those with attachment avoidance, individuals with attachment anxiety may be more likely to engage with CAI for psychological support, potentially becoming a key demographic within its user base. CAI systems could benefit from tailoring their communication styles to address the unique needs of users with attachment anxiety, ensuring these technologies provide desired emotional support and safe engagement.

While recognizing the significant potential of CAI for psychological support, we believe it is also equally crucial to be aware of the associated risks that might arise in human-CAI interaction. Research has consistently linked attachment anxiety with increased social media use and addiction [44-47]. Consequently, individuals with attachment anxiety may also be more susceptible to developing unhealthy dependencies on CAI in the postengagement phase. Proactively identifying solutions and applying appropriate strategies during the design phase can mitigate potential negative outcomes. It is essential to alert CAI designers to potential maladaptive behaviors associated with CAI use. Integrating protective measures, such as timely advice and interventions, can help safeguard the user experience and optimize therapeutic outcomes, particularly for users with attachment anxiety. In terms of attachment avoidance, the lack of a significant result is congruent with previous research [48,49]; avoidant-attached individuals have a need to deactivate the attachment system (eg, by inhibiting proximity-seeking behaviors), and this tendency often makes it difficult to observe and capture their avoidant nature in surveys. To be more specific, individuals with attachment avoidance often prefer self-reliance and independence. They are more likely to maintain emotional distance to feel safe rather than seek emotional support, which might lead to a weaker or nonexistent link between attachment avoidance and CAI adoption intention, similar to the results found in our study.

To the best of our knowledge, this study is the first to explore the relationship between attachment styles and CAI adoption in the context of CAI-based therapies. More evidence is needed to determine whether our current findings (in both significance and direction) are replicable and reliable. If attachment styles, particularly attachment anxiety, prove to be a consistent predictor of CAI adoption intention, this could inform the development of more customized designs that promote safer interactions and outcomes that are more effective.

Trust in Generalized AI and CAI

In addition, past studies [33] have already established a relationship between attachment styles and trust in generalized AI. However, our results suggest that this may not necessarily replicate in the context of CAI. According to the results of correlation matrix illustrated in Table 2, both attachment anxiety and avoidance were not significantly correlated with trust in CAI counseling in our study. Therefore, trust was not assessed as a mediator between attachment anxiety and CAI adoption intentions. Specifically, perceived trust was not significantly associated with attachment anxiety (β =.091, P=.30) and attachment avoidance (β =-.161, P=.07). This finding is inconsistent with the conclusion (ie, higher attachment anxiety predicts lower trust) found by Gillath et al [33] in their study, that we previously relied on in hypothesizing a negative direction between attachment anxiety and CAI adoption intention in hypothesis 2. Hence, this inconsistency could signify a more complex relationship between attachment styles and perceived trust in CAI adoption.

To contextualize these results in understanding this inconsistency, one possible explanation could be that participants' attachment systems may not have been sufficiently activated in this study. According to a review conducted by Campbell and Marshall [50], attachment theory is interactionist in nature, particularly attachment anxiety. Highly anxious individuals may exhibit heightened distress responses when they perceive cues as threats to their relationships. However, in the absence of such cues or when their security needs are fulfilled, they often show similar proximity-seeking tendencies in affect, cognition, and relationship processes to people with low anxiety levels. This suggests that when the attachment system is not effectively activated, it could potentially lead to weaker or contradictory associations between attachment styles and attachment-related behaviors, such as the relationship found between insecure attachment styles and trust in CAI in our study. Future studies are suggested to include research-supported

methods (eg, recalling relationship experiences and hypothetical scenarios) for activating participants' attachment systems before conducting the study.

Furthermore, given its sensitive nature, it is also possible that insecure attachment styles affect trust in CAI counseling in a different manner than trust in AI in general. As mentioned in previous sections, we formulated our second hypothesis based on a relevant study conducted by Gillath et al [33]. The AI technologies examined in this study focused on the relationship between attachment insecurity and perceived trust that were designed for more general purposes, such as self-driving vehicles, medical diagnostic apps, and matchmaking services. Unlike self-driving vehicles or matchmaking services, mental health support requires a higher level of empathy and emotional attunement, areas in which AI technology is more likely to be considered to fall short. Research examining the relationship between attachment styles and trust in AI used for sensitive purposes, such as conversational AI for mental health, need to be specific to the context for which they are used.

General Limitations

There are several limitations that should be mentioned in our study. First of all, one potential obstacle his field of study is the lack of uniformity in defining and measuring AI-related trust. Using different scales to assess trust can lead to the capture of distinct facets of trust and, consequently, generate contradictory results.

Previous research [51] has highlighted the presence of 2 essential components within the overarching concept of trust in AI systems, "user trust potential" and "perceived system trustworthiness." User trust potential typically encompasses the user's internal factors, such as attachment styles, that influence their trust in AI systems. In contrast, perceived system trustworthiness focuses on external factors, including user experience (eg, efficiency and effectiveness) and perceived technical trustworthiness (eg, accuracy, security, and privacy). The existing measurement tools for trust do not clearly distinguish and separately assess these 2 aspects, which may lead to inaccurate capture of the relationship and misses out on important nuances.

This signals a pressing need for the development and validation of a more consolidated and clearly structured measurement tool for trust in AI. Such an instrument would greatly enhance the field's ability to comprehensively assess trust in AI systems. Furthermore, an intriguing avenue for future research is the exploration of which facet of trust, whether internal factors or external factors, exerts a more pronounced influence on actual engagement behaviors, specifically in terms of actual usage. This question holds significant potential for shedding light on the nuances of trust in AI systems and informing practical applications.

Second, our dependent variable CAI adoption intention was measured with a single item on an ordinal scale. Single item may lack the sensitivity to detect subtle differences or changes in the outcome variable, potentially missing important variations in the data. In addition, measuring CAI adoption intention continuously would capture gradual changes more efficiently, leading to more precise description relationship between dependent variable and other independent variables. Multi-item scales should be used to measure adoption intention continuously in future research to increase validity and reliability.

Third, our use of a news article as a vignette to illustrate the use of CAI in mental health support may have implied a subtle positive valence. This could stem from the portrayal of a CAI as a tool that is able to assist individuals with mental health issues. However, as far as possible, we adopted a neutral tone to the vignette, and future studies could consider the portrayal of CAI as a mental health tool with successful (positive) or unsuccessful (negative) outcomes for more generalizable effects.

Furthermore, it is worth noting that disorganized attachment was not examined in the current study. As the first study exploring the relationship between attachment styles and CAI adoption for psychological support, we focused on more clearly defined variables-anxious and avoidant attachment styles-to enable more interpretable and consistent initial insights in this novel area of research. Future research can build on this foundation by incorporating additional insecure attachment styles to generate deeper and more nuanced findings that inform CAI design. In addition, our research participants were sourced through a web-based platform with participants from a single country (the United States). Future research incorporating more diverse samples are encouraged to address these limitations and enhance the generalizability of the findings. Also, although we have excluded participants with previous CAI counseling experience and the results still hold true for the subgroup that reported being entirely unfamiliar with CAI counseling, we acknowledge that future studies would benefit from clearly distinguishing between indirect and direct exposures from which participants gain their familiarity when measuring it.

Conclusion

In conclusion, our study serves as a pioneering effort in the realm of CAI adoption for mental support, being one of the only papers to examine the impact of attachment styles and perceived trust on CAI adoption. Our findings indicate that perceived trust remains a crucial factor influencing adoption intention; individuals with higher perceived trust are more inclined to try CAI therapies when needed. In addition, attachment anxiety, rather than attachment avoidance, is significantly and positively linked to CAI adoption. These results contribute to the current literature as a good first glimpse into human-CAI relationship and inform the future design of CAI systems, particularly in the mental health setting. By understanding how factors such as perceived trust and attachment styles influence CAI adoption, this study underscores the importance of developing tailored, evidence-based strategies to foster user trust and address specific concerns related to mental health applications. Such strategies may potentially help to mitigate potential risks of CAI adoption, such as overreliance or misuse, ensuring that CAI technologies are safely and effectively integrated into mental health care services. Furthermore, these findings highlight the need for continuous evaluation and adaptation of CAI features to better meet the diverse needs of users, ultimately promoting more positive outcomes in mental health support. Future research

Wu et al

should build upon these insights to further refine CAI applications, ensuring they are both user-centered and ethically

sound, thereby enhancing their potential to provide effective and accessible mental health care solutions.

Acknowledgments

This work was supported by a Google Research Scholar Award awarded to KL.

Authors' Contributions

XW contributed to conceptualization, data curation, formal analysis, investigation, project administration, visualization, writing – original draft, and writing – review and editing. KL contributed to conceptualization, formal analysis, funding acquisition, methodology, resources, supervision, writing – original draft, and writing – review and editing. MJD contributed to methodology and writing – review and editing.

Conflicts of Interest

None declared.

Multimedia Appendix 1 Online supplementary material. [DOCX File , 29 KB - ai v4i1e68960 app1.docx]

References

- 1. Xygkou A, Siriaraya P, Covaci A, Prigerson HG, Neimeyer R, Ang CS, et al. The "Conversation" about loss: understanding how chatbot technology was used in supporting people in Grief. In: CHI '23: Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems. New York, NY, United States: Association for Computing Machinery; Apr 19, 2023:1-15.
- Vaidyam AN, Wisniewski H, Halamka JD, Kashavan MS, Torous JB. Chatbots and conversational agents in mental health: a review of the psychiatric landscape. Can J Psychiatry 2019 Jul;64(7):456-464 [FREE Full text] [doi: 10.1177/0706743719828977] [Medline: 30897957]
- Sedlakova J, Trachsel M. Conversational artificial intelligence in psychotherapy: a new therapeutic tool or agent? Am J Bioeth 2023 May;23(5):4-13 [FREE Full text] [doi: 10.1080/15265161.2022.2048739] [Medline: 35362368]
- Malgaroli M, Hull TD, Zech JM, Althoff T. Natural language processing for mental health interventions: a systematic review and research framework. Transl Psychiatry 2023 Oct 06;13(1):309 [FREE Full text] [doi: 10.1038/s41398-023-02592-2] [Medline: 37798296]
- 5. Gaffney H, Mansell W, Tai S. Conversational agents in the treatment of mental health problems: mixed-method systematic review. JMIR Ment Health 2019 Oct 18;6(10):e14166 [FREE Full text] [doi: 10.2196/14166] [Medline: 31628789]
- Li H, Zhang R, Lee Y, Kraut RE, Mohr DC. Systematic review and meta-analysis of AI-based conversational agents for promoting mental health and well-being. NPJ Digit Med 2023 Dec 19;6(1):236 [FREE Full text] [doi: 10.1038/s41746-023-00979-5] [Medline: 38114588]
- 7. Prakash AV, Das S. Intelligent conversational agents in mental healthcare services: a thematic analysis of user perceptions. Pac. Asia J. Assoc. Inf. Syst 2020;12:1-34. [doi: 10.17705/1thci.12201]
- Fiske A, Henningsen P, Buyx A. Your robot therapist will see you now: ethical implications of embodied artificial intelligence in psychiatry, psychology, and psychotherapy. J Med Internet Res 2019 May 09;21(5):e13216 [FREE Full text] [doi: 10.2196/13216] [Medline: 31094356]
- 9. American Psychological Association. AI is changing every aspect of psychology: here's what to watch for. Monitor on Psychology. 2023 Jul 1. URL: <u>https://www.apa.org/monitor/2023/07/psychology-embracing-ai</u> [accessed 2024-10-25]
- Xie T, Pentina I, Hancock T. Friend, mentor, lover: does chatbot engagement lead to psychological dependence? Journal of Service Management 2023 May 10;34(4):806-828. [doi: <u>10.1108/josm-02-2022-0072</u>]
- He Y, Yang L, Qian C, Li T, Su Z, Zhang Q, et al. Conversational agent interventions for mental health problems: systematic review and meta-analysis of randomized controlled trials. J Med Internet Res 2023 Apr 28;25:e43862 [FREE Full text] [doi: 10.2196/43862] [Medline: 37115595]
- 12. Liu L. What makes Inflection's Pi a great companion chatbot. Medium. 2023 May 23. URL: <u>https://medium.com/@lindseyliu/</u> what-makes-inflections-pi-a-great-companion-chatbot-8a8bd93dbc43 [accessed 2024-11-18]
- 13. Cassidy J, Shaver PR. Handbook of Attachment Theory, Research, and Clinical Applications. 3rd ed. New York, NY: Guilford Press; Jul 19, 2016.
- 14. Gillath O, Karantzas GC, Fraley RC. Adult Attachment: A Concise Introduction to Theory and Research. Cambridge, MA: Academic Press; Mar 23, 2016.



- 15. Bedué P, Fritzsche A. Can we trust AI? An empirical investigation of trust requirements and guide to successful AI adoption. JEIM 2021 Apr 30;35(2):530-549. [doi: 10.1108/jeim-06-2020-0233]
- McKnight DH. Trust in information technology. In: Davis GB, editor. The Blackwell Encyclopedia of Management. Vol. 7 Management Information Systems. England: Blackwell; 2005:329-331.
- 17. Fernandes T, Oliveira E. Understanding consumers' acceptance of automated technologies in service encounters: Drivers of digital voice assistants adoption. Journal of Business Research 2021 Jan;122:180-191. [doi: 10.1016/j.jbusres.2020.08.058]
- Kasilingam DL. Understanding the attitude and intention to use smartphone chatbots for shopping. Technology in Society 2020 Aug;62:101280. [doi: <u>10.1016/j.techsoc.2020.101280</u>]
- Wutz M, Hermes M, Winter V, Köberlein-Neu J. Factors influencing the acceptability, acceptance, and adoption of conversational agents in health care: integrative review. J Med Internet Res 2023 Sep 26;25:e46548 [FREE Full text] [doi: 10.2196/46548] [Medline: <u>37751279</u>]
- 20. Ainsworth MDS, Blehar MC, Waters E, Wall S. Patterns of Attachment: A Psychological Study of the Strange Situation. Mahwah, New Jersey: Lawrence Erlbaum Associates Inc; Dec 31, 1978.
- 21. Bowlby J. Attachment: Attachment and Loss Volume One. New York, NY: Basic Books; 1983.
- 22. Bowlby J. A Secure Base. London, UK: Routledge; Sep 1, 2005.
- 23. Main M, Kaplan N, Cassidy J. Security in infancy, childhood, and adulthood: a move to the level of representation. Monographs of the Society for Research in Child Development 1985;50(1/2):66-104. [doi: 10.2307/333827]
- 24. Ravitz P, Maunder R, Hunter J, Sthankiya B, Lancee W. Adult attachment measures: a 25-year review. J Psychosom Res 2010 Oct;69(4):419-432. [doi: 10.1016/j.jpsychores.2009.08.006] [Medline: 20846544]
- 25. Bartholomew K, Horowitz LM. Attachment styles among young adults: a test of a four-category model. J Pers Soc Psychol 1991 Aug;61(2):226-244. [doi: 10.1037//0022-3514.61.2.226] [Medline: 1920064]
- 26. Bretherton I, Munholland K. Internal working models in attachment relationships: a construct revisited. In: Cassidy J, Shaver PR, editors. Handbook of Attachment: Theory, Research, and Clinical Applications. 2nd ed. New York, NY: Guilford Press; 1999:89-111.
- 27. Fonagy P, Gergely G, Jurist E, Target M. Affect Regulation, Mentalization, and the Development of the Self. London, UK: Routledge; Jan 1, 2004.
- Paetzold RL, Rholes WS, Kohn JL. Disorganized Attachment in Adulthood: Theory, Measurement, and Implications for Romantic Relationships. Review of General Psychology 2015 Jun 01;19(2):146-156 [FREE Full text] [doi: 10.1037/gpr0000042]
- 29. Bretherton I, Munholland K. Internal working models in attachment relationships: elaborating a central construct in attachment theory. In: Cassidy J, Shaver PR, editors. Handbook of attachment: Theory, research, and clinical applications (2nd ed.). New York, NY: The Guilford Press; 2008:102-127.
- Collins NL, Read SJ. Adult attachment, working models, and relationship quality in dating couples. J Pers Soc Psychol 1990;58(4):644-463. [doi: <u>10.1037//0022-3514.58.4.644</u>] [Medline: <u>14570079</u>]
- Hodge DR, Gebler-Wolfe MM. Understanding technology's impact on youth: attachment theory as a framework for conceptualizing adolescents' relationship with their mobile devices. Children & Schools 2022 May 19;44(3):153-162. [doi: <u>10.1093/cs/cdac007</u>]
- 32. Birnbaum GE, Mizrahi M, Hoffman G, Reis HT, Finkel EJ, Sass O. What robots can teach us about intimacy: The reassuring effects of robot responsiveness to human disclosure. Computers in Human Behavior 2016 Oct;63:416-423. [doi: 10.1016/j.chb.2016.05.064]
- Gillath O, Ai T, Branicky MS, Keshmiri S, Davison RB, Spaulding R. Attachment and trust in artificial intelligence. Computers in Human Behavior 2021 Feb;115:106607. [doi: 10.1016/j.chb.2020.106607]
- 34. Mikulincer M. Attachment working models and the sense of trust: an exploration of interaction goals and affect regulation. Journal of Personality and Social Psychology 1998;74(5):1209-1224. [doi: 10.1037//0022-3514.74.5.1209]
- 35. Pistole MC. Attachment relationships: self-disclosure and trust. J. Ment. Health Couns 1993;15(1):94-106 [FREE Full text]
- Simmons BL, Gooty J, Nelson DL, Little LM. Secure attachment: implications for hope, trust, burnout, and performance. Journal of Organizational Behavior 2009 Jan 29;30(2):233-247. [doi: <u>10.1002/job.585</u>]
- 37. Wu X, Liew K. Exploring the Initial Leap: Attachment Styles and Their Influence on Intention of Using CAI for Mental Support. 2023 Dec 17. URL: <u>https://doi.org/10.17605/OSF.IO/C2XQD</u> [accessed 2025-03-28]
- Gritti ES, Bornstein RF, Barbot B. The smartphone as a "significant other": interpersonal dependency and attachment in maladaptive smartphone and social networks use. BMC Psychol 2023 Sep 28;11(1):296 [FREE Full text] [doi: 10.1186/s40359-023-01339-4] [Medline: <u>37770997</u>]
- Collins NL. Working models of attachment: implications for explanation, emotion and behavior. J Pers Soc Psychol 1996 Oct;71(4):810-832. [doi: <u>10.1037//0022-3514.71.4.810</u>] [Medline: <u>8888604</u>]
- 40. Gulati S, Sousa S, Lamas D. Design, development and evaluation of a human-computer trust scale. Behaviour & Information Technology 2019 Aug 31;38(10):1004-1015. [doi: 10.1080/0144929x.2019.1656779]
- 41. Shandilya E, Fan M. Understanding older Adults' perceptions and challenges in using AI-enabled everyday technologies. 2022 Presented at: Chinese CHI '22: Proceedings of the Tenth International Symposium of Chinese CHI; 2024 Feb 12; Guangzhou, China p. 105-116.

- 42. Draxler F, Buschek D, Tavast M, Hämäläinen P, Schmidt A, Kulshrestha J, et al. Gender, age, and technology education influence the adoption and appropriation of LLMs. ArXiv 2025 [FREE Full text] [doi: 10.48550/arXiv.2310.06556]
- 43. Collins NL, Feeney BC. A safe haven: an attachment theory perspective on support seeking and caregiving in intimate relationships. J Pers Soc Psychol 2000;78(6):1053-1073. [doi: 10.1037//0022-3514.78.6.1053] [Medline: 10870908]
- 44. Hart J, Nailling E, Bizer GY, Collins CK. Attachment theory as a framework for explaining engagement with Facebook. Personality and Individual Differences 2015 Apr;77:33-40. [doi: <u>10.1016/j.paid.2014.12.016</u>]
- 45. Oldmeadow JA, Quinn S, Kowert R. Attachment style, social skills, and Facebook use amongst adults. Computers in Human Behavior 2013 May;29(3):1142-1149. [doi: 10.1016/j.chb.2012.10.006]
- 46. Blackwell D, Leaman C, Tramposch R, Osborne C, Liss M. Extraversion, neuroticism, attachment style and fear of missing out as predictors of social media use and addiction. Personality and Individual Differences 2017 Oct 1;116:69-72. [doi: 10.1016/j.paid.2017.04.039]
- 47. Eroglu Y. Interrelationship between attachment styles and Facebook addiction. J. educ. train. stud 2016;4(1):150-160. [doi: 10.11114/jets.v4i1.1081]
- 48. Mikulincer M, Shaver PR. The attachment behavioral system in adulthood: activation, psychodynamics, and interpersonal processes. In: Zanna MP, editor. Advances in Experimental Social Psychology. San Diego, CA: Elsevier Academic Press; 2003:53-152.
- 49. Mikulincer M, Shaver PR. Attachment in Adulthood: Structure, Dynamics, and Change. New York, NY: Guilford Press; 2007.
- 50. Campbell L, Marshall T. Anxious attachment and relationship processes: an interactionist perspective. J Pers 2011 Dec;79(6):1219-1250. [doi: 10.1111/j.1467-6494.2011.00723.x] [Medline: 21299557]
- Stanton B, Jensen T. Trust and artificial intelligence. NIST Interagency/Internal Report (NISTIR). United States: National Institute of Standards and Technology; 2021 Mar 2. URL: <u>https://tsapps.nist.gov/publication/get_pdf.cfm?pub_id=931087</u> [accessed 2025-03-26]

Abbreviations

AI: artificial intelligence CAI: conversational artificial intelligence IWM: internal working model LLM: large language model OR: odds ratio

Edited by K El Emam; submitted 18.11.24; peer-reviewed by MA Virtanen, H Li; comments to author 31.12.24; revised version received 21.01.25; accepted 23.02.25; published 22.04.25.

<u>Please cite as:</u> Wu X, Liew K, Dorahy MJ Trust, Anxious Attachment, and Conversational AI Adoption Intentions in Digital Counseling: A Preliminary Cross-Sectional Questionnaire Study JMIR AI 2025;4:e68960 URL: <u>https://ai.jmir.org/2025/1/e68960</u> doi:10.2196/68960 PMID:

©Xiaoli Wu, Kongmeng Liew, Martin J Dorahy. Originally published in JMIR AI (https://ai.jmir.org), 22.04.2025. This is an article distributed under the terms of the Creative Commons open-access Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Harnessing Moderate-Sized Language Models for Reliable Patient Data Deidentification in Emergency Department Records: Algorithm Development, Validation, and Implementation Study

Océane Dorémus¹, MSc; Dylan Russon¹, MSc; Benjamin Contrand¹, MSc; Ariel Guerra-Adames^{1,2}, BEng, MSc; Marta Avalos-Fernandez², HDR, PhD; Cédric Gil-Jardiné^{1,3}, MD, PhD; Emmanuel Lagarde¹, HDR, PhD

¹AHeaD Team, University of Bordeaux, INSERM, BPH, U1219, 146 Rue Léo Saignat, Bordeaux, France ²SISTM Team, University of Bordeaux, INSERM, INRIA, BPH, U1219, Bordeaux, France

³Department of Emergency Medicine, Bordeaux University Hospital, Bordeaux, France

Corresponding Author:

Océane Dorémus, MSc AHeaD Team, University of Bordeaux, INSERM, BPH, U1219, 146 Rue Léo Saignat, Bordeaux, France

Abstract

Background: The digitization of health care, facilitated by the adoption of electronic health records systems, has revolutionized data-driven medical research and patient care. While this digital transformation offers substantial benefits in health care efficiency and accessibility, it concurrently raises significant concerns over privacy and data security. Initially, the journey toward protecting patient data deidentification saw the transition from rule-based systems to more mixed approaches including machine learning for deidentifying patient data. Subsequently, the emergence of large language models has represented a further opportunity in this domain, offering unparalleled potential for enhancing the accuracy of context-sensitive deidentification. However, despite large language models offering significant potential, the deployment of the most advanced models in hospital environments is frequently hindered by data security issues and the extensive hardware resources required.

Objective: The objective of our study is to design, implement, and evaluate deidentification algorithms using fine-tuned moderate-sized open-source language models, ensuring their suitability for production inference tasks on personal computers.

Methods: We aimed to replace personal identifying information (PII) with generic placeholders or labeling non-PII texts as "ANONYMOUS," ensuring privacy while preserving textual integrity. Our dataset, derived from over 425,000 clinical notes from the adult emergency department of the Bordeaux University Hospital in France, underwent independent double annotation by 2 experts to create a reference for model validation with 3000 clinical notes randomly selected. Three open-source language models of manageable size were selected for their feasibility in hospital settings: Llama 2 (Meta) 7B, Mistral 7B, and Mixtral 8×7B (Mistral AI). Fine-tuning used the quantized low-rank adaptation technique. Evaluation focused on PII-level (recall, precision, and F_1 -score) and clinical note-level metrics (recall and BLEU [bilingual evaluation understudy] metric), assessing deidentification effectiveness and content preservation.

Results: The generative model Mistral 7B performed the highest with an overall F_1 -score of 0.9673 (vs 0.8750 for Llama 2 and 0.8686 for Mixtral 8×7B). At the clinical notes level, the model's overall recall was 0.9326 (vs 0.6888 for Llama 2 and 0.6417 for Mixtral 8×7B). This rate increased to 0.9915 when Mistral 7B only deleted names. Four notes of 3000 failed to be fully pseudonymized for names: in 1 case, the nondeleted name belonged to a patient, while in the others, it belonged to medical staff. Beyond the fifth epoch, the BLEU score consistently exceeded 0.9864, indicating no significant text alteration.

Conclusions: Our research underscores the significant capabilities of generative natural language processing models, with Mistral 7B standing out for its superior ability to deidentify clinical texts efficiently. Achieving notable performance metrics, Mistral 7B operates effectively without requiring high-end computational resources. These methods pave the way for a broader availability of pseudonymized clinical texts, enabling their use for research purposes and the optimization of the health care system.

(JMIR AI 2025;4:e57828) doi:10.2196/57828

KEYWORDS

de-identification; machine learning; large language model; natural language processing; electronic health records; transformers; general data protection regulation; clinical notes



Introduction

The digitization of medical data has profoundly transformed health care, facilitating the easy and efficient sharing of patient information [1]. This digital transition, embodied by electronic health records systems, offers promising opportunities for data-driven solutions, research, and surveillance on a pan-European scale [2]. Yet, alongside the many advantages of digitization come significant concerns about the privacy and security of sensitive patient data [3]. The European General Data Protection Regulation emphasizes the necessity of stringent data protection measures, particularly for health-related information [2]. Clinical notes, which often encompass identifiable patient details, must adhere to these standards to safeguard patient confidentiality [loi informatique et liberté], before any data sharing researchers face the critical task of developing and integrating methods that mask sensitive data, guaranteeing protection against any unauthorized access [4]. Our team was recently faced with this challenge in a project aimed at classifying clinical notes from emergency services to extract the necessary information for the establishment of a trauma observatory [5].

Manual deidentification of medical records is not feasible, as it is expensive in terms of personnel resources and the time required to accomplish the task. Alternatively, multiple strategies have been implemented for the automated deidentification of medical records [6,7]. These methods evolved from systems based on explicit rules, regular expressions or dictionaries [8-16], to techniques using machine learning [17-19].

In recent years, the evolution of language models, particularly those based on transformer architectures, has reshaped the landscape of natural language processing (NLP). Transformers, introduced by Vaswani et al [20] in 2017, provided a novel approach to handling sequential data using self-attention mechanisms, thereby obviating the need for recurrent layers and significantly augmenting training efficiency. This pivotal innovation paved the way for the advent of progressively sophisticated and expansive models. Transformer-based language models of a moderate scale, particularly through customized and fine-tuned versions of the architecture BERT [21], have demonstrated high capabilities in various health care applications. These models excel in understanding and processing complex clinical texts, enabling tasks such as predicting patient outcomes and identifying medical events. For instance, a recent study highlighted the effectiveness of fine-tuned BERT models in analyzing clinical notes to predict occurrences of falls, showcasing the model's ability to comprehend subtle nuances in medical language [22]. Additionally, BERT models offer significant benefits for tasks such as named entity recognition (NER). Those models offer notable benefits for deidentification, thanks to their capacity to discern patterns among words and phrases. They have the ability to learn from diverse text types means they can effectively tackle various pseudonymization challenges, as they can be trained to erase a wide range of identifiable details across different document types.

```
Dorémus et al
```

The burgeoning of computational resources and datasets has since kindled a shift toward the construction of massive models, embedded with trillions of parameters [23-25]. As they grew in size, their generalization aptitude and versatility witnessed substantial enhancement, optimizing tasks such as deidentification. In 2023, Liu et al [25] underscored the potential of leveraging the GPT-4's inherent capacity for 0-shot in-context learning. A salient highlight of their methodology was its ability to maintain the original structure and meaning of the text after the removal of confidential details. While the capabilities of GPT-4 are undeniable, its application in the realm of health care presents serious ethical and legal dilemmas, primarily concerning data privacy and patient confidentiality. On the one hand, due to the vastness of the model, local hosting of GPT-4 is not feasible, therefore, data should be transmitted to external servers, in this case OpenAI's infrastructure. On the other hand, considering the confidentiality of the weights, only locally hosted servers are regulatory compliant. Furthermore, considering that GPT-4 is a proprietary model, organizations cannot fully control or audit the underlying mechanics or data handling processes.

From a regulatory perspective, sending personal health information externally contravenes many data protection regulations, most notably the General Data Protection Regulation in Europe and the Health Insurance Portability and Accountability Act [26,27] in the United States. This raises not just data sovereignty issues but also infringes on patient rights, as they might not have explicitly consented for their data to be processed in external environments. Hence, while the technological feats of models such as GPT-4 are commendable, their real-world applications, especially in sensitive sectors such as health care, require careful consideration and possibly, significant adjustments to ensure full regulatory compliance and ethical integrity.

Generative language models significantly smaller in size (several billion parameters compared to over a trillion for GPT-4) have been recently developed and made available to the public under licenses that allow for almost unrestricted use (Llama 2 by Meta [28]) or even under open-source terms (Mistral [29]).

The objective of our study is to design, implement, and evaluate deidentification methods involving proper prompt engineering and fine-tuning of 3, open-source language models (Llama 2 7B, Mistral 7B, and Mixtral 8×7B [30]). These models were selected for their moderate size, making them suitable for deployment on personal computers for production inference tasks.

Methods

Study Design

We first attempted to perform the task using only prompt engineering and 0-shot inference. As we failed to achieve any significant results, we improved the selected models' capability to deidentify clinical texts using quantized low-rank adaptation [31] fine-tuning with a dataset of instruction or response pairs. In practice, the task consists in replacing personal identifying information (PII; name, location, dates, telephone number,

```
XSL•FO
```

email, or identification numbers) with generic placeholders, represented as "[XXXXX]," or, when no PII is detected, by generating the text as "ANONYMOUS." The ultimate goal of this procedure is to preserve text content, ensuring adherence to privacy and confidentiality requirements.

Data Source, Datasets Allocations, and Annotation

Within the emergency department, triage is conducted by triage nurses. This process involves the collection of information on each patient, including medical history, current symptoms, vital signs, and personal details. It is these data that we have at our disposal in our study. For this investigation, we curated our dataset from a repository containing 425,680 clinical free-text notes (Multimedia Appendix 1), authored by a nurse during the initial reception and triage of individuals at the Bordeaux University Hospital's adult emergency department over the period spanning from January 2013 to December 2022. A subset of 6097 clinical notes was randomly selected and independently annotated by 2 experts. Any arising discrepancies were adjudicated by a third expert, thus establishing a reference database. From this curated sample of 6097 clinical notes, 3000 were delineated to constitute a test dataset, upon which accuracy metrics were evaluated (Figure 1). The residual 3097 clinical notes, alongside an additional sample of 3000 clinical notes designed using filters and keywords search to encompass a broad spectrum of identifying scenarios, comprised the validation dataset.





In order to further assess whether the deidentification performances of the models varies with the type of PII, we classified identifying information within clinical notes into 6 distinct categories (Table 1). These categories were used by annotators to label such information in the test dataset. While we have taken care to remove obvious PII such as names, addresses, and identification numbers, it is important to note that deidentification cannot be considered as a strict anonymization process. For instance, in cases of rare diseases or very specific descriptions, reidentification could theoretically be possible. As every clinical history is unique, ensuring complete anonymity is unattainable. Our goal is to pseudonymize data, striking a balance between patient confidentiality and data utility for research, as removing all sensitive information will significantly diminish the data's usefulness.



Туре	Code	Description
Individual names	NAME	Includes both first and last names of individuals (including patients and medical staff) or of rela- tives, employers, or household members of the individuals, ensuring personal identification.
Dates	DATE	Pertains to specific dates related to medical events, appointments, or personal milestones, formatted as day, month, or year.
Geographic identifiers	LOC ^a	Covers names of geographic locations such as cities, medical facilities, or addresses, facilitating location-based identification.
Phone numbers	TEL ^b	Comprises all forms of telephone numbers for direct contact, including mobile and landline numbers.
Email addresses	MAIL	Encompasses electronic mail addresses, allowing for digital communication.
Miscellaneous identifiers	OTHER	A catch-all category for unique identifiers not covered by other categories, including social se- curity numbers, medical analysis codes, and URLs for patient images.

^aLOC: location.

^bTEL: telephone.

Selected Models

We have selected 3 language models that share the following 2 characteristics: being open-source and of sufficiently small size for the production phase to be implemented on affordable PC-type systems. These are Llama 2 7B, Mistral 7B, and Mixtral. Llama 2 7B is developed by Meta. Launched in 2023, this is a 7-billion-parameter model, which is claimed to exhibit a good balance between performance and efficiency. We also selected the Mistral 7B model, introduced to the public in October 2023. It has demonstrated superior performance, either matching or surpassing that of Llama 2 13B in extensive benchmarks and showing comparable results to Llama 1 34B in specific domains such as reasoning, mathematics, and code generation. In December 2023, the Mixtral 8×7B model was released. It is described as a Sparse Mixture of Experts language model. Its key innovation lies in the routing of inference tasks through 1 selected expert out of 8, enabled by an additional routing layer. Consequently, despite its 8×7B size with respect to fine-tuning, Mixtral achieves a significant efficiency by requiring an eightfold reduction in parameters for inference task.

Fine-Tuning and Inference

Each model was subjected to the same prompt or response pairs of clinical notes. The fine-tuning process was uniformly standardized across all 3 models, albeit with variations in batch sizes and quantization rates to accommodate our hardware constraints. The fine-tuning configuration for Mistral 7B and Llama 2 7B involved a batch size of 24 records per GPU, while Mixtral used a batch size of 20. The models were fine-tuned over 15 epochs, using the AdamW optimizer [32] with a learning rate of 5e-5 and a weight decay of 0.01. We used the quantized low-rank adaptation technique, allowing for specific adjustments in selected parts of the model, such as query, key, value, output,

```
https://ai.jmir.org/2025/1/e57828
```

and gates projection modules while preserving the overall architecture integrity. The low-rank adaptation configuration included a rank setting of 32, a learning rate multiplier (alpha) set to 64, with a dropout of 0.1, and without any bias setting. Additionally, to optimize computational efficiency and minimize memory consumption, the models were quantized to 8-bit precision for both 7B models, and 4-bit precision for Mixtral. At every fine-tuning epoch, the inference was induced for each model.

The computational undertakings of this research were performed on a server running Ubuntu (version 22.04; Canonical Ltd), outfitted with 4 A100 GPUs, collectively boasting 320GB of VRAM.

Evaluation

Overview

In evaluating the deidentification performance of personal data within clinical notes, our analysis is structured around 2 primary methodologies. The first methodology operates at the PII-level, enabling us to provide estimates of recall, precision, and F_1 -scores that are comparable with previous work in the literature. The second methodology focuses on clinical notes as the statistical unit, enabling us to assess the variation in recall performance according to the category of PII. This latter approach needs to be complemented by the measurement of a BLEU (bilingual evaluation understudy) score to assess potential modifications in the text. The assessment of the number of successful deidentifications was conducted through a comparison with the manually annotated test dataset.

PII-Based Metrics

This approach centers on treating each PII as an independent statistical unit. This perspective allows us to gauge the precision

XSL•FO RenderX

Dorémus et al

JMIR AI

and recall of our deidentification efforts at the most granular level. Recall in this context is conceptualized as the proportion of PIIs accurately identified and removed from the clinical notes.

RecalPII=numberofcouedlydeidentifiedPIIperdinicalnotestotalnumberofPIIperdinicalnotes

Precision, meanwhile, reflects the accuracy of our model in identifying and eliminating actual PIIs, distinguishing between correct identifications and false positives.

PecisionPII=numberofconectlydeidentifiedPIIpecinicalnotestotalnumberofPIItagged

The summary F_1 -score measure is:

F1-score=21precision+1recall

Clinical Note-Based Metrics

The second approach adopts the entire clinical note as the statistical unit of analysis. Here we evaluate the success of deidentification on a document-wide scale, marking a "success" when every PII within a note has been successfully deidentified. Such a measure offers insight into the overall effectiveness of our deidentification protocols. Recall, in this instance, measures the ratio of fully deidentified notes to those containing any PII.

Rel-unterforetyeb-ibifethistotesmogibifyighistotesthunterfibifyighistotes

Because the clinical notes in the validation set are annotated by indicating the nature of the PII (according to the categories in Table 1), it is possible to detail the variations in recall by category. The relevance of precision is altered in this context, as it necessitates a different consideration of what constitutes a pseudonymization attempt, denoted by the presence of a pseudonymization tag. Instead, the potential alteration of content possibly induced by the deidentification process was measured using the BLEU score [33].

$BLEU=BP \cdot exp(\sum wnlogpn)$

where BP is the brevity penalty, w_n the weight for each n-gram, and p_n the precision of n-grams. We set a value of 4 for the BLEU score calculation, aligning with common practice in NLP to capture up to 4-gram coherence, thereby ensuring a comprehensive evaluation of content preservation.

Ethical Considerations

Overview

This study was conducted as part of the Automated Processing of Emergency Department Visit Summaries for a National Observatory project, which aims to automate the processing of emergency department visit summaries for national observation purposes.

The study received the following regulatory approvals: (1) the Ethics Committee for Research in Science and Health, validating the compliance of the protocol with current ethical requirements; and (2) the National Commission on Informatics and Liberty, under decision DR-2022-235 (authorization request 922170), allowing the processing of data for this study.

Confidentiality and Data Protection

The data processing was carried out exclusively on a secure local server, specially dedicated to this purpose. This server meets the current security standards, ensuring the confidentiality, integrity, and protection of the processed information. All necessary technical and organizational measures have been implemented to prevent unauthorized access to the data and to ensure strict compliance with regulatory requirements.

Compensation

Since this study relies solely on the analysis of pre-existing medical data and does not require direct patient involvement, no financial compensation was provided.

Results

Data Overview

Very few notes contained PIIs categorized as email addresses and "other." These categories are included in the training sample due to an ad hoc selection process, which used filters to ensure representation, as half of the set was selected this way. Our examination of the test sample, which consists entirely of randomly selected clinical notes, reveals that names, places, and dates are the most prevalent types of PII. The categories of identifying data in the training and test sets are summarized in Table 2.

Regarding the length of clinical notes, they range from 8 to 3916 characters (with an average of 443, SD 289 characters) in the training set and from 3 to 2138 characters (averaging 439, SD 283 characters) in the test set. A total of 935 (31.2%) clinical notes in the test set contain at least one PII.



Table . Enhanced distribution of PII^a in train and tests sets.

	Train set	Test set
Clinical notes		
Nonanonymous medical notes, n (%)	3442 (56.5)	935 (31.2)
Randomly selected medical notes, n	3097	3000
Ad hoc selected medical notes, n	3000	b
Total count, n	6097	3000
PII categories, n		
NAME	3016	555
LOC ^c	1801	715
TEL ^d	650	41
EMAIL	13	0
DATE	2404	607
OTHER	33	1
Total number of PII	7917	1919

^aPII: personal identifying information.

^bThis corresponds to the absence of ad-hoc selected medical notes.

^cLOC: location.

_

^dTEL: telephone.

Performance Using PII-Based Metrics

Figure 2 plots the change in the F_1 -score over the 15 epochs of fine-tuning for the 3 respective models. The Mistral 7B model

quickly reaches a performance plateau, where its F_1 -score stabilizes, whereas the Mixtral 8×7B and Llama 2 7B models exhibit a slower rate of improvement, with both reaching a plateau in their F_1 -scores around the 12th epoch.







Recall Analysis

The recall estimates of the 3 models are shown in Figures 3 and 4.

Mistral 7B and Mixtral 8×7B achieved better overall recall. The Mistral 7B and Mixtral 8×7B models demonstrated marked enhancements in their deidentification efficacy across epochs, starting from the third epoch onward. Notably, the Mistral 7B model has shown a rapid improvement in performance, achieving a performance plateau by the sixth epoch. Conversely, the Mixtral 8×7B model's improvement trajectory was more

gradual, reaching a stable performance level by the 13 epoch. The overall success rate appears not to improve beyond epoch 7 for the Mistral 7B model. Consequently, in the subsequent analysis, this epoch was selected for comparing success rates across categories.

As shown in Figure 5, Mistral 7B consistently outperformed Mixtral 8×7B and Llama 2 across all data identification categories. Despite Mixtral's performance improving over time, it still did not surpass Mistral 7B. Using Mistral 7B, a 100% (41/41) recall was observed for phone numbers (Figure 5) and recall was lower for locations than for names.





Figure 3. Plot of recall by epoch: clinical notes as statistical unit.





Figure 4. Plot of recall by epoch: PII as statistical unit. PII: personal identifying information.



Figure 5. Plot of recall by epoch for PII: (A) Location, (B) Telephone, (C) Name, (D) Date. PII: personal identifying information. Telephone

(B)





BLEU Score

BLEU-4 scores were calculated to assess whether the models modified the texts at the note level. During the deidentification

process, medical texts remained almost unchanged as demonstrated by a consistently high BLEU-4 score (Figure 6) beyond epoch 5.



Figure 6. Plot of BLEU score by epoch: clinical note as statistical unit. BLEU: bilingual evaluation understudy

Results Summary at Epoch 7

Table . Fine-tuned models performance at epoch 7.

The Table 3 below presents a summary of performance metrics achieved by our models at epoch 7.

The results demonstrate that the Mistral 7B model outperforms
both the Mixtral $8 \times 7B$ and Llama 2 7B with a F_1 -score of
0.9673. When using clinical note as the statistical unit, the recall
is also much higher (0.9326) for Mistral 7B than Llama 2 and
Mixtral 8×7B models.

Model	Clinical notes	Personal identifying information		
	Recall	Precision	Recall	F ₁ -score
Mistral 7B	0.9326	0.9721	0.9625	0.9673
Llama 2 7B	0.6888	0.9596	0.8041	0.875
Mixtral 8×7B	0.6417	0.9852	0.7655	0.8616

Error Analysis

In epoch 7 of the Mistral 7B model, a total of 63 clinical notes were not properly pseudonymized, as detailed in Table 4. Among these, location (LOC) errors were the most frequent, with 44 instances. Deleting geographical and institutional identifiers then remains a significant challenge (with a recall of 86.1%). Specifically, 31 notes still included names of health or social service facilities, while 12 notes still included names of cities. Conversely, errors involving names (NAME) were significantly fewer, with only 4 instances, including 1 patient name and 3 doctors' names, resulting in a high recall of 99.8%

for this category. Date-related errors (DATE) were observed in 14 notes (with a recall of 97.8%).

The test dataset, comprising 3000 clinical notes, underwent a post hoc examination to identify any inaccuracies resulting from manual annotations that would have been detected by all 15 versions of our 3 finely-tuned models, spanning epochs 1 to 15. Through this process, we were able to pinpoint 65 notes in which the model detected personally identifiable information through the medical histories that were categorized as anonymous (ie, without identifying data, 2066 clinical notes), in which the model detected personally identifying information that had been overlooked by human annotators.



Table . Summary of deidentification errors at epoch 7.

	~	
Errors	Count	
Total	63	
Returned ANONYMOUS	29	
Annotation error	34	
Errors in personal identifying information categories		
NAME	4	
LOC ^a	44	
DATE	14	
OTHER	1	

^aLOC: location.

We observed that the models outperformed human annotation in 9 clinical records from the test set. Specifically, in these 9 records, 5 locations (LOC), 3 names (NAMES), and 1 date (DATE) were omitted during manual annotation. The remaining 53 records present annotation errors from the models. Therefore, the total number of actual personally identifiable information (PII) amounts to 1928, contrary to the 1919 initially identified by our experts.

Subsequently, corrections were made to the test dataset based on these findings, and main outcomes were recomputed in an additional sensitive analysis. The metric measurements after accounting for these modifications are only slightly altered from the original results (see Multimedia Appendix 2 for the details).

Discussion

Principal Findings

In this study, we assessed the performance of 3 generative NLP models in the deidentification of clinical text documents. The generative model Mistral 7B demonstrated the highest performance with an overall F_1 -score of 0.9673. At the clinical notes level, the same model achieved an overall recall of 0.9326, with this rate increasing to 0.9915 for the deletion of names. The evaluation was based on a test dataset of 3000 clinical notes, among which only 4 notes failed to be fully deidentified for names; in one case, the identifying name was that of a patient. As the method relies on the use of generative models, we also measured potential text alterations generated by the process. Beyond the fifth epoch, the BLEU score consistently exceeded 0.9864.

Strengths

Our work distinguishes itself from the existing scientific literature by using a method that does not rely on NER and uses moderate-sized models. Instead, the use of generative large language models allows for the production of text that is pseudonymized by removing PII components. This is the reason why we added metrics that use clinical notes as the statistical unit. This led us to use the BLEU metric to assess potential text alterations. Another consequence of this method is that no hyperparameters are set which made it possible to avoid the use of separate test and validation dataset partitions. The size of our training and test samples, independently annotated by 2 experts,

```
https://ai.jmir.org/2025/1/e57828
```

RenderX

constitutes a significant strength in our study. To our knowledge, no other study has used a test sample of such size (3000 notes). Yet, it is crucial to have the means to detect rare errors if the ultimate goal is to develop a system that guarantees the pseudonymization of clinical texts. We deliberately limited our model selection to those whose implementation does not require powerful servers and can be executed on personal computers equipped with a consumer-grade graphics card. The largest model is Mixtral 8×7B, which has approximately 8 times more parameters than the other 2 models. Mixtral 8×7B shares the same architecture as Mistral 7B, with the distinction that each layer consists of 8 feed-forward blocks. Although training it requires significant memory capacity, this is not the case during the inference phase, during which only 2 of the feed-forward blocks are used, selected by a network acting as a router.

Limitations

Annotation Process Inaccuracies

Overview

During the annotation process, we observed some inaccuracies. To assess the impact of these inaccuracies on our metrics, we conducted a post hoc analysis, taking into account corrections made by the model. Although this analysis revealed few variations, it is important to note that some errors may still remain in the text set, undetected by the model. These undetected errors could potentially affect the overall performance of the model.

Model Choice

We opted for a fine-tuned large language model—based approach over a dedicated NER model due to pragmatic considerations. Our hypothesis was that a targeted human annotation process, with expert annotators pinpointing PII within texts, would be more effective than a broad NER annotation effort, given the same time investment. Focusing on essential PII elements helps us minimize the ambiguities that broader NER annotations often entail. This focus leads to improved precision and recall rates during the training phase. Furthermore, this approach is in line with the Automated Processing of Emergency Department Visit Summaries for a National Observatory project's objectives, which prioritize the accurate removal of PII from unstructured medical texts.

The default choice for identification tasks is usually a bidirectional transformer, starting from the hypothesis that the relationship of a word with its context before and after that word allows for better comprehension of the role of those words and therefore should be more suited for NER tasks. However, this hypothesis no longer holds when dealing with generative models. Since the goal here is to generate redacted text, the provided prompt has access to the entire corrected phrase. Consequently, relative to a given word, implications cannot be considered unidirectional.

Model Sharing Constraints

Overview

Another significant limitation is that our model was fine-tuned using nonanonymous clinical texts, which prevents us from sharing the model's weights with the community. Sharing the model's weights could potentially allow for the extraction of the original training data. This limitation restricts the model's reproducibility and its broader applicability across different research settings and medical domains.

Demographic and Textual Bias

The processed data are in free-text format, written by health care staff, which introduces significant variability. This variability is not only present between different services within the same health facility but also across various centers. Factors such as the content of clinical notes, the medical abbreviations used, writing styles, and the level of detail in documentation can differ greatly from one source to another. Such differences could potentially impact the performance of our models, making it essential to test and adapt our approach to data from diverse sources.

Comparison With Prior Work

Comparing the performance of our models with those documented in the literature presents challenges because our models are specifically fine-tuned to pseudonymize French-language clinical notes. Consequently, it is not feasible to apply them to the English-language databases traditionally used for benchmarking, such as i2b2 (i2b2 TranSMART Foundation) [34], MIMIC II (PhysioNet) [35], and MIMIC III (PhysioNet) [36].

In addition to these differences in benchmarking context, there are also divergences in the methodologies used for deidentification. Historically, deidentification of medical records has evolved from rule-based systems, which rely on predefined rules, regular expressions, and dictionaries, to more sophisticated machine learning approaches. Rule-based methods, while easy to implement and interpret, often fall short in handling the variability and unpredictability inherent in unstructured clinical texts. On the other hand, machine learning-based approaches offer more flexibility and adaptability, particularly when dealing with large and diverse datasets. These models can learn patterns directly from the data, making them more effective in identifying PIIs that deviate from standard formats. However, their effectiveness is heavily dependent on the quality and quantity of annotated data available for training. Moreover, machine learning models typically require significant computational

resources and expertise in model tuning, which can be a barrier to adoption, particularly in resource-constrained settings.

Our proposed model leverages these advanced machine learning techniques, specifically fine-tuned for the French language. This focus allows our model to effectively capture and manage the linguistic intricacies specific to French clinical notes, such as frequent abbreviations and unstructured text entries, which are common in emergency department settings.

Additionally, our results demonstrate that while our model performs comparably to those trained on English-language corpora, certain challenges persist, particularly in the detection of location-based PIIs. This is likely due to the complexity introduced by variations in PII forms, such as acronyms and abbreviations, as well as the presence of typing errors, which are less predictable and harder to model.

Therefore, to compare performance metrics accurately, it is necessary to assess the complexity of clinical texts from these databases against those used in our study. In the Multimedia Appendix 1, we include examples of clinical notes from our dataset to demonstrate that PIIs can appear randomly within the text, in an unstructured manner, and that these PIIs, along with the rest of the text, often include numerous abbreviations. This tendency toward abbreviation is explained by the unique demands of emergency department settings, where nurses are required to perform efficient, real-time data entry into the hospital's information system. As a result, our dataset more closely aligns with MIMIC II, which features unstructured clinical notes made by nurses, as opposed to i2b2, where each type of information is distinctly separated, preventing the amalgamation of multiple PIIs within single sentences.

As shown in Multimedia Appendix 3 [37-43], our results (overall F_1 -score of 0.9673) are on par with previous studies on English clinical text corpus that used an algorithm including models using self-attention [17,24,36,44]. The Multimedia Appendix 4 [37,38,43] summarizes study results that examined recall variations according to PII categories. These figures consistently show that the relative weakness of these algorithms, ours included, lies in a small number of errors concerning locations. Our dataset presents additional challenges for PII identification due to the presence of multiple variations of PII, including acronyms, abbreviations, and typing errors. Specifically, of the 44 notes with failed identification, 15 involved abbreviations or acronyms, and 2 contained typing errors.

Future Work

We aim to enhance the detection capabilities of PII in our medical notes by fine-tuning our model with newly annotated data. To achieve this, we plan to generate artificial clinical notes using commercially available application programming interfaces, such as GPT-4. These large language models, much more powerful than ours, can produce realistic notes containing PII and annotations, which will facilitate the training process and increase data diversity.

By generating a substantial volume of these artificial data, we can ensure equitable representation of different PII categories and evaluate 2 key aspects: identifying the optimal amount of

clinical notes needed to achieve the highest possible accuracy and recall, and comparing the effectiveness of models fine-tuned with real data versus those fine-tuned with artificially generated data.

Using this newly developed model based on artificial data, we aim to make it available as an open-source resource, benefiting the broader community. Additionally, this foundation will enable us to create a multilingual model capable of processing both English and French clinical notes. This multilingual model will allow us to perform performance comparisons against literature benchmark datasets such as i2b2 and MIMIC. The performance of these refined models will be evaluated using our corrected test set, along with newly annotated data from various emergency services.

This study is currently focused on data from an emergency department in France. In the subsequent phases, our goal is to extend this methodology to other services across France, with the ambition of creating a national French observatory on trauma. However, it is important to consider the potential for demographic biases in our model's performance.

By diversifying data sources, we aim to enhance the model's generalizability. If biases are identified in this process, we plan to retrain the model, either by using a specific portion of data from each service or by integrating synthetic data to mitigate these biases.

We intend to extend our methodology to other types of sensitive documents, such as medico-legal records, to evaluate the generalizability and effectiveness of our approach in protecting personal information across various domains.

We are also considering integrating explainability methods, similar to those used by Arnaud et al [45], to enhance the transparency of our model in PII detection. These techniques, based on transformer models and interpretability approaches such as LIME [46], which have already proven effective on triage note data similar to ours, could strengthen user trust and facilitate the adoption of our technologies in clinical settings.

Through this comprehensive approach, we aim to enhance the value and applicability of our models, contributing to the development of privacy-preserving technologies in the health care domain and strengthening the security of patients' sensitive information.

Ethical Considerations and Practical Implementations

The use of small to moderate-sized models is a key consideration in our approach. These models are generally capable of running on GPUs with at least 16 GB of VRAM, making them suitable for use on personal computers or within local infrastructures. This is particularly advantageous for institutions with limited resources, as it allows them to manage data privately and securely without relying on extensive external infrastructure. However, while local deployment ensures better control over sensitive data, it can also be time-consuming and may introduce challenges related to the interoperability of different systems.

One of the main challenges of this pipeline is its implementation across all participating emergency services, given that not all institutions may be equipped to efficiently manage these new procedures. The rationale behind implementing this process is rooted in a data-sharing initiative aimed at establishing a national observatory, which necessitates enhanced protection for the information being used.

At this stage, centralizing the data in a dedicated center with the necessary computational resources remains the simplest solution. This would allow for secure, controlled, and efficient management of patient data. Alternatively, the process could be implemented directly within health data warehouses, enabling these facilities to store and apply the deidentification process locally. Regardless of the approach, it is imperative that the use of this pipeline on health data is conducted within a legally and digitally controlled framework, authorized by the relevant authorities.

Given the potential risks of data reidentification, especially when dealing with unique clinical histories, we emphasize that pseudonymization alone is insufficient and should be accompanied by additional protection and security measures to prevent unauthorized access to sensitive data.

Conclusion

Our research underscores the significant capabilities of generative NLP models, with Mistral 7B standing out for its superior ability to deidentify clinical texts efficiently. Achieving notable performance metrics, Mistral 7B operates effectively without requiring high-end computational resources. These methods pave the way for a broader availability of pseudonymized clinical texts, enabling their use for research purposes and the optimization of the health care system.

Acknowledgments

This study was conducted under the Automated Processing of Emergency Department Visit Summaries for a National Observatory (TARPON) project by the Bordeaux Population Health-Assessing Health in a Digitalizing Real-World team and the Bordeaux University Hospital's emergency department. We thank the labeling team and the University Hospital of Bordeaux for their logistical support and data access.

Data Availability

The datasets generated or analyzed during this study are not publicly available due to the confidential nature of the patient data used.



Authors' Contributions

EL, CG-J, and MA-F did the conceptualization and design. BC, OD, EL, DR, and CG-J worked on the annotation. OD, CG-J, and EL analyzed and interpreted. OD, EL, and AG-A drafted this paper. All authors handled the critical revision. CG-J provided this study's material. EL supervised.

Conflicts of Interest

None declared.

Multimedia Appendix 1 Examples of French nursing notes. [DOCX File, 16 KB - ai v4i1e57828 app1.docx]

Multimedia Appendix 2 Analysis of performance evaluation on corrected test set. [DOCX File, 15 KB - ai v4i1e57828 app2.docx]

Multimedia Appendix 3 Comparative table of statistical results from previous studies. [DOCX File, 20 KB - ai v4i1e57828 app3.docx]

Multimedia Appendix 4

Comparative table of recall across PII categories from previous studies. PII: personal identifying information. [DOCX File, 12 KB - ai v4i1e57828 app4.docx]

References

- 1. Menachemi N, Collum TH. Benefits and drawbacks of electronic health record systems. Risk Manag Healthc Policy 2011;4:47-55. [doi: 10.2147/RMHP.S12985] [Medline: 22312227]
- Regulation (EU) 2016/679 of the european parliament and of the council of 27 april 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing directive 95/46/EC (general data protection regulation). European Parliament and Council. 2016. URL: <u>https://dvbi.ru/Portals/0/DOCUMENTS_SHARE/RISK_MANAGEMENT/EBA/GDPR_eng_rus.pdf</u> [accessed 2025-03-31]
- 3. MHealth: new horizons for health through mobile technologies: second global survey on ehealth. World Health Organization:: World Health Organization; 2012. URL: <u>https://iris.who.int/bitstream/handle/10665/44607/9789241564250_eng.</u> pdf?sequence=1&isAllowed=y [accessed 2025-03-31]
- 4. El Emam K. Methods for the de-identification of electronic health records for genomic research. Genome Med 2011 Apr 27;3(4):25. [doi: 10.1186/gm239] [Medline: 21542889]
- Chenais G, Gil-Jardiné C, Touchais H, et al. Deep learning transformer models for building a comprehensive and real-time trauma observatory: development and validation study. JMIR AI 2023 Jan 12;2:e40843. [doi: <u>10.2196/40843</u>] [Medline: <u>38875539</u>]
- Meystre SM, Friedlin FJ, South BR, Shen S, Samore MH. Automatic de-identification of textual documents in the electronic health record: a review of recent research. BMC Med Res Methodol 2010 Aug 2;10:70. [doi: <u>10.1186/1471-2288-10-70</u>] [Medline: <u>20678228</u>]
- 7. Negash B, Katz A, Neilson CJ, et al. De-identification of free text data containing personal health information: a scoping review of reviews. Int J Popul Data Sci 2023;8(1):2153. [doi: 10.23889/ijpds.v8i1.2153] [Medline: 38414537]
- Beckwith BA, Mahaadevan R, Balis UJ, Kuo F. Development and evaluation of an open source software tool for deidentification of pathology reports. BMC Med Inform Decis Mak 2006 Mar 6;6:12. [doi: 10.1186/1472-6947-6-12] [Medline: 16515714]
- 9. Berman JJ. Concept-match medical data scrubbing. Arch Pathol Lab Med 2003 Jun 1;127(6):680-686. [doi: 10.5858/2003-127-680-CMDS]
- Friedlin FJ, McDonald CJ. A software tool for removing patient identifying information from clinical documents. J Am Med Inform Assoc 2008;15(5):601-610. [doi: <u>10.1197/jamia.M2702</u>] [Medline: <u>18579831</u>]
- Gupta D, Saul M, Gilbertson J. Evaluation of a deidentification (De-Id) software engine to share pathology reports and clinical documents for research. Am J Clin Pathol 2004 Feb;121(2):176-186. [doi: <u>10.1309/E6K3-3GBP-E5C2-7FYU</u>] [Medline: <u>14983930</u>]
- Morrison FP, Li L, Lai AM, Hripcsak G. Repurposing the clinical record: can an existing natural language processing system de-identify clinical notes? J Am Med Inform Assoc 2009;16(1):37-39. [doi: <u>10.1197/jamia.M2862</u>] [Medline: <u>18952938</u>]

https://ai.jmir.org/2025/1/e57828

- 13. Neamatullah I, Douglass MM, Lehman L, et al. Automated de-identification of free-text medical records. BMC Med Inform Decis Mak 2008 Jul 24;8:32. [doi: 10.1186/1472-6947-8-32] [Medline: 18652655]
- 14. Ruch P, Baud RH, Rassinoux AM, Bouillon P, Robert G. Medical document anonymization with a semantic lexicon. Proc AMIA Symp 2000:729-733. [Medline: <u>11079980</u>]
- Sweeney L. Replacing personally-identifying information in medical records, the Scrub system. Proc AMIA Annu Fall Symp 1996:333-337. [Medline: <u>8947683</u>]
- 16. Thomas SM, Mamlin B, Schadow G, McDonald C. A successful technique for removing names in pathology reports using an augmented search and replace method. Proc AMIA Symp 2002:777-781. [Medline: <u>12463930</u>]
- 17. Ahmed T, Aziz MMA, Mohammed N. De-identification of electronic health record using neural network. Sci Rep 2020 Oct 29;10(1):18600. [doi: 10.1038/s41598-020-75544-1] [Medline: 33122735]
- 18. Guo Y, Gaizauskas RJ, Roberts I, Demetriou G, Hepple M. Identifying personal health information using support vector machines. Semantic Scholar. 2006. URL: https://api.semanticscholar.org/CorpusID:16833759 [accessed 2025-03-31]
- 19. Dernoncourt F, Lee JY, Uzuner O, Szolovits P. De-identification of patient notes with recurrent neural networks. J Am Med Inform Assoc 2017 May 1;24(3):596-606. [doi: 10.1093/jamia/ocw156] [Medline: 28040687]
- Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, et al. Attention is all you need. NeurIPS Proceedings. 2017. URL: <u>https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf</u> [accessed 2025-03-31]
- 21. Devlin J, Chang MW, Lee K, Toutanova K. BERT: pre-training of deep bidirectional transformers for language understanding. In: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers) 2019:4171-4186. [doi: 10.18653/v1/N19-1423]
- 22. Cheligeer C, Wu G, Lee S, et al. BERT-based neural network for inpatient fall detection from electronic medical records: retrospective cohort study. JMIR Med Inform 2024 Jan 30;12:e48995. [doi: <u>10.2196/48995</u>] [Medline: <u>38289643</u>]
- 23. OpenAI, Achiam J, Adler S, Agarwal S, Ahmad L, et al. GPT-4 technical report. arXiv. Preprint posted online on Mar 4, 2024. [doi: 10.48550/arXiv.2303.08774]
- 24. Bai Y, Kadavath S, Kundu S, Askell A, Kernion J, Jones A, et al. Constitutional AI: harmlessness from AI feedback. AI-Plans. 2022. URL: <u>https://ai-plans.com/file_storage/4f32fa39-3a01-46c7-878e-c92b7aa7165f_2212.08073v1.pdf</u> [accessed 2025-03-31]
- Liu J, Gupta S, Chen A, et al. OpenDeID pipeline for unstructured electronic health record text notes based on rules and transformers: deidentification algorithm development and validation study. J Med Internet Res 2023 Dec 6;25:e48145. [doi: <u>10.2196/48145</u>] [Medline: <u>38055317</u>]
- 26. Health insurance portability and accountability act of 1996 (HIPAA). Centers for Disease Control and Prevention, Public Health Law. 2024. URL: <u>https://www.cdc.gov/phlp/php/resources/</u> health-insurance-portability-and-accountability-act-of-1996-hipaa.html?CDC_AAref_Val=https://www.cdc.gov/phlp/ publications/topic/hipaa.html [accessed 2025-03-31]
- 27. Liu Z, Huang Y, Yu X, Zhang L, Wu Z, Cao C, et al. DeID-GPT: zero-shot medical text de-identification by GPT-4. arXiv. Preprint posted online on Dec 21, 2023. [doi: 10.48550/arXiv.2303.11032]
- 28. Touvron H, Martin L, Stone K, Albert P, Almahairi A, Babaei Y, et al. Llama 2: open foundation and fine-tuned chat models. arXiv. Preprint posted online on Jul 19, 2023. [doi: <u>10.48550/arXiv.2307.09288</u>]
- 29. Jiang AQ, Sablayrolles A, Mensch A, Bamford C, Chaplot DS, Casas D, et al. Mistral 7B. arXiv. Preprint posted online on Oct 10, 2023. [doi: 10.48550/arXiv.2310.06825]
- 30. Jiang AQ, Sablayrolles A, Roux A, Mensch A, Savary B, Bamford C, et al. Mixtral of experts. arXiv. Preprint posted online on Jan 8, 2024. [doi: <u>10.48550/arXiv.2401.04088</u>]
- 31. Dettmers T, Pagnoni A, Holtzman A, Zettlemoyer L. QLoRA: efficient finetuning of quantized llms. advances in neural information processing systems. NeurIPS Proceedings. 2023. URL: <u>https://proceedings.neurips.cc/paper_files/paper/2023/file/1feb87871436031bdc0f2beaa62a049b-Paper-Conference.pdf</u> [accessed 2025-03-31]
- 32. Loshchilov I, Hutter F. Decoupled weight decay regularization. arXiv. Preprint posted online on Jan 4, 2019. [doi: 10.48550/arXiv.1711.05101]
- Papineni K, Roukos S, Ward T, Zhu WJ. BLEU: a method for automatic evaluation of machine translation. In: Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics: Association for Computational Linguistics:311-318. [doi: 10.3115/1073083.1073135]
- 34. Informatics for Integrating Biology & the Bedside (i2b2). URL: <u>https://www.i2b2.org/</u> [accessed 2025-03-31]
- 35. Saeed M, Villarroel M, Reisner AT, et al. Multiparameter Intelligent Monitoring in Intensive Care II: a public-access intensive care unit database. Crit Care Med 2011 May;39(5):952-960. [doi: 10.1097/CCM.0b013e31820a92c6] [Medline: 21283005]
- 36. Johnson AEW, Pollard TJ, Shen L, et al. MIMIC-III, a freely accessible critical care database. Sci Data 2016 May 24;3:160035. [doi: 10.1038/sdata.2016.35] [Medline: 27219127]
- Liu L, Perez-Concha O, Nguyen A, et al. Web-based application based on human-in-the-loop deep learning for deidentifying free-text data in electronic medical records: development and usability study. Interact J Med Res 2023 Aug 25;12:e46322. [doi: <u>10.2196/46322</u>] [Medline: <u>37624624</u>]
- 38. Grouin C, Zweigenbaum P. Automatic de-identification of French clinical records: comparison of rule-based and machine-learning approaches. Stud Health Technol Inform 2013;192:476-480. [Medline: 23920600]
- Chazard E, Mouret C, Ficheur G, Schaffar A, Beuscart JB, Beuscart R. Proposal and evaluation of FASDIM, a Fast and Simple De-Identification Method for unstructured free-text clinical records. Int J Med Inform 2014 Apr;83(4):303-312. [doi: <u>10.1016/j.ijmedinf.2013.11.005</u>] [Medline: <u>24370391</u>]
- Catelli R, Gargiulo F, Casola V, De Pietro G, Fujita H, Esposito M. Crosslingual named entity recognition for clinical de-identification applied to a COVID-19 Italian data set. Appl Soft Comput 2020 Dec;97:106779. [doi: <u>10.1016/j.asoc.2020.106779</u>] [Medline: <u>33052197</u>]
- 41. Berg H, Henriksson A, Dalianis H. The impact of de-identification on downstream named entity recognition in clinical text. In: Proceedings of the 11th International Workshop on Health Text Mining and Information Analysis: Association for Computational Linguistics; 2020:1-11. [doi: 10.18653/v1/2020.louhi-1.1]
- Syed M, Sexton K, Greer M, et al. DeIDNER Model: A Neural Network Named Entity Recognition Model for Use in the De-identification of Clinical Notes. In: Biomed Eng Syst Technol Int Jt Conf BIOSTEC Revis Sel Pap Feb 2022, Vol. 5:640-647. [doi: 10.5220/0010884500003123] [Medline: 35386186]
- 43. Tchouka Y, Couchot JF, Coulmeau M, Laiymani D, Rahmani A. De-identification of french unstructured clinical notes for machine learning tasks. arXiv. Preprint posted online on Oct 6, 2023. [doi: 10.48550/arXiv.2209.09631]
- 44. Meaney C, Hakimpour W, Kalia S, Moineddin R. A comparative evaluation of transformer models for de-identification of clinical text data. arXiv. Preprint posted online on Mar 25, 2022. [doi: <u>10.48550/arXiv.2204.07056</u>]
- 45. Arnaud E, Elbattah M, Moreno-Sánchez PA, Dequen G, Ghazali DA. Explainable NLP model for predicting patient admissions at emergency department using triage notes. In: 2023 IEEE International Conference on Big Data (BigData): IEEE:4843-4847. [doi: 10.1109/BigData59044.2023.10386753]
- 46. Ribeiro MT, Singh S, Guestrin C. "Why should I trust you?": explaining the predictions of any classifier. arXiv. Preprint posted online on Aug 9, 2016. [doi: <u>10.48550/arXiv.1602.04938</u>]

Abbreviations

BLEU: bilingual evaluation understudy **NER:** named entity recognition **NLP:** natural language processing **PII:** personal identifying information

Edited by J Sun; submitted 28.02.24; peer-reviewed by E Vashishtha, GK Gupta, M Elbattah; revised version received 28.08.24; accepted 23.10.24; published 01.04.25.

<u>Please cite as:</u> Dorémus O, Russon D, Contrand B, Guerra-Adames A, Avalos-Fernandez M, Gil-Jardiné C, Lagarde E Harnessing Moderate-Sized Language Models for Reliable Patient Data Deidentification in Emergency Department Records: Algorithm Development, Validation, and Implementation Study JMIR AI 2025;4:e57828 URL: <u>https://ai.jmir.org/2025/1/e57828</u> doi:10.2196/57828

© Océane Dorémus, Dylan Russon, Benjamin Contrand, Ariel Guerra-Adames, Marta Avalos-Fernandez, Cédric Gil-Jardiné, Emmanuel Lagarde. Originally published in JMIR AI (https://ai.jmir.org), 1.4.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Utility-based Analysis of Statistical Approaches and Deep Learning Models for Synthetic Data Generation With Focus on Correlation Structures: Algorithm Development and Validation

Marko Miletic¹, BSc; Murat Sariyar¹, PhD

Institute for Optimisation and Data Analysis (IODA), Bern University of Applied Sciences, Biel, Switzerland

Corresponding Author: Murat Sariyar, PhD Institute for Optimisation and Data Analysis (IODA) Bern University of Applied Sciences Höheweg 80 Biel, 2502 Switzerland Phone: 41 32 321 64 37 Email: <u>murat.sariyar@bfh.ch</u>

Abstract

Background: Recent advancements in Generative Adversarial Networks and large language models (LLMs) have significantly advanced the synthesis and augmentation of medical data. These and other deep learning–based methods offer promising potential for generating high-quality, realistic datasets crucial for improving machine learning applications in health care, particularly in contexts where data privacy and availability are limiting factors. However, challenges remain in accurately capturing the complex associations inherent in medical datasets.

Objective: This study evaluates the effectiveness of various Synthetic Data Generation (SDG) methods in replicating the correlation structures inherent in real medical datasets. In addition, it examines their performance in downstream tasks using Random Forests (RFs) as the benchmark model. To provide a comprehensive analysis, alternative models such as eXtreme Gradient Boosting and Gated Additive Tree Ensembles are also considered. We compare the following SDG approaches: Synthetic Populations in R (synthpop), copula, copulagan, Conditional Tabular Generative Adversarial Network (ctgan), tabular variational autoencoder (tvae), and tabula for LLMs.

Methods: We evaluated synthetic data generation methods using both real-world and simulated datasets. Simulated data consist of 10 Gaussian variables and one binary target variable with varying correlation structures, generated via Cholesky decomposition. Real-world datasets include the body performance dataset with 13,393 samples for fitness classification, the Wisconsin Breast Cancer dataset with 569 samples for tumor diagnosis, and the diabetes dataset with 768 samples for diabetes prediction. Data quality is evaluated by comparing correlation matrices, the propensity score mean-squared error (pMSE) for general utility, and F_1 -scores for downstream tasks as a specific utility metric, using training on synthetic data and testing on real data.

Results: Our simulation study, supplemented with real-world data analyses, shows that the statistical methods copula and synthpop consistently outperform deep learning approaches across various sample sizes and correlation complexities, with synthpop being the most effective. Deep learning methods, including large LLMs, show mixed performance, particularly with smaller datasets or limited training epochs. LLMs often struggle to replicate numerical dependencies effectively. In contrast, methods like tvae with 10,000 epochs perform comparably well. On the body performance dataset, copulagan achieves the best performance in terms of pMSE. The results also highlight that model utility depends more on the relative correlations between features and the target variable than on the absolute magnitude of correlation matrix differences.

Conclusions: Statistical methods, particularly synthop, demonstrate superior robustness and utility preservation for synthetic tabular data compared with deep learning approaches. Copula methods show potential but face limitations with integer variables. Deep Learning methods underperform in this context. Overall, these findings underscore the dominance of statistical methods for synthetic data generation for tabular data, while highlighting the niche potential of deep learning approaches for highly complex datasets, provided adequate resources and tuning.

(JMIR AI 2025;4:e65729) doi:10.2196/65729



KEYWORDS

synthetic data generation; medical data synthesis; random forests; simulation study; deep learning; propensity score mean-squared error

Introduction

In recent years, Generative Adversarial Networks (GANs) and large language models (LLMs) have revolutionized the synthesis and augmentation of medical data [1-3]. These technologies have introduced methods for creating high-quality, realistic datasets, which are essential for advancing machine learning (ML) applications in the health care sector [4-6]. The ability to synthesize realistic medical data is particularly valuable in contexts where data privacy and availability are major concerns [7]. Medical data is often subject to strict regulations due to privacy laws and ethical considerations, which can limit the availability of comprehensive datasets for research and development. By using GANs and LLMs to generate synthetic data, researchers and practitioners can overcome these limitations, creating datasets that preserve the statistical properties and correlations of the original data while ensuring that individual patient identities remain protected.

However, despite the promising capabilities of GANs and LLMs, several challenges persist in leveraging these technologies effectively for medical data synthesis [8-11]. A key challenge is the ability of these models to accurately capture and replicate the intricate relationships within medical datasets. Medical data often exhibits complex interdependencies between features, such as the relationship among symptoms, diagnostic indicators, and treatment outcomes. Inaccurate representation of these correlation structures can result in synthetic data that fails to mimic the true variability and relationships found in real-world medical data [12]. The use of synthetic medical data also raises ethical concerns, particularly regarding the potential perpetuation or, in some cases, even amplification of biases inherent in the original datasets [13]. For instance, GANs tend to prioritize matching overall data distribution rather than subgroup-level details. Such representation issues can translate into new or stronger associations between sensitive attributes such as race and medical conditions [14]. If high data quality is promised based on such data because a particular metric performs well, ML methods may establish incorrect associations accordingly.

Focusing on pairwise correlation structures in medical data synthesis, despite their limitations in complex data environments, remains crucial for several reasons: (1) correlation analysis identifies primary dependencies as a starting point for understanding how variables interact; (2) if a ML model recognizes that certain variables are typically correlated, it can better simulate realistic scenarios, leading to more accurate predictions and insights; and (3) pairwise correlation structures provide a baseline for validating and comparing synthetic data. Even though they might not capture all forms of dependence, comparing correlations in synthetic data with those in real-world data can help assess the fidelity and quality of the generated datasets.

There have been several approaches addressing correlations in the context of Synthetic Data Generation (SDG), particularly for relational data [15]. Most methodological studies aim to capture correlation structures by extending existing techniques. For example, Vu et al [16] explored how to make the loss function of GANs correlation-aware but found no significant benefit. In contrast, Patel et al [17] demonstrated that incorporating a Correlational Neural Network can improve a GAN's ability to capture correlations, slightly outperforming the MedGAN model. Torfi and Fox developed realistic synthetic health care records by leveraging Convolutional Neural Networks to capture correlations between medical features, achieving comparable performance to real data in ML tasks while maintaining privacy and statistical fidelity [18]. Rajabi and Garibay [19] showed that effective consideration of correlations can enhance fairness in synthetic data. These works are noteworthy because the primary goal of advanced SDG methods is to capture the full dependency structure.

Despite the substantial body of work on validation and benchmarking in SDG, there is a notable gap in studies specifically assessing how the correlation structure of real data influences the effectiveness of SDG methods in replicating such relationships. Understanding whether faithfully reproducing correlation structures is critical for achieving high-quality results in downstream tasks remains an open question. This issue is particularly relevant given the increasing reliance on SDG methods across various domains. Simulation studies are well-suited to address these questions, as they enable controlled analysis of specific factors affecting model performance [20]. For instance, Strobl et al [21] demonstrated through simulations that Random Forest (RF) models tend to produce biased variable selection when predictors differ in scale or category count.

The aim of this study is to address the research gap by developing a simulation design and validating the results on 3 real-world medical datasets. We evaluate how effectively SDG methods can replicate the correlation structure of the original data and perform a classification task using RF. To provide a comprehensive analysis, alternative models such as eXtreme Gradient Boosting [22] and Gated Additive Tree Ensembles [23] are also considered. In addition, for one notable case, we assess whether the relevant variables are selected based on variable importance measures, as correlation matrix distances are often calculated in practice without addressing their impact. For this analysis, we use the following SDG approaches: Synthetic Populations in R (synthpop) [24], copula [25], copulagan [26], Conditional Tabular Generative Adversarial Network (ctgan) [27], Tabular Variational Autoencoder (tvae) [27], and tabula for LLMs [28,29], the latter of which per default uses DistilGPT-2 (distilled Generative Pretrained Transformer -2), a streamlined version of the english-language model GPT-2. The corresponding assessment will help practitioners in guiding their choice of SDG methods.

```
https://ai.jmir.org/2025/1/e65729
```

RenderX

Methods

Overview

The schematic diagram in Figure 1 outlines the key steps in the methodology used in this study. The process begins with data

generation, where simulated datasets were created using correlation matrix construction and target variable creation. Besides that, we selected 3 real-world datasets (Body Performance [BP], Breast Cancer [BC], and Diabetes [DB]). All datasets are then used to generate and evaluate various SDG methods.

Figure 1. Overview of the methodology workflow. BC: Breast Cancer Dataset; BP: Body Performance Dataset; ctgan: Conditional Tabular Generative Adversarial Network; DB: Diabetes Dataset; pMSE: Propensity Score Mean-Squared Error; SDG: synthetic data generation; tvae: Tabular Variational Autoencoder; VIMP: variable importance.



Datasets

Real-World Datasets

We selected 3 medical datasets from Kaggle – Body Performance (BP), Breast Cancer (BC), and Diabetes (DB) – that are commonly used in predictive modeling and data analysis tasks. All 3 datasets involve classification problems. The correlation matrices of these datasets are provided in Figure 2.

The BP dataset provides comprehensive data on physical fitness and body measurements, encompassing variables such as height, weight, age, gender, body fat percentage, and details of physical activity and fitness routines. It includes 13,393 samples with 11 numerical features and a categorical target variable that classifies individuals into four fitness categories: excellent, good, average, and poor. Among the features, age and sit-up count are recorded as integers. The BC dataset comprises 569 entries, each with 30 numerical features extracted from digitized images of fine needle aspirates of breast masses. These features, representing the mean, standard error, and maximum value, quantify geometric and textural properties of cell nuclei, including radius, texture, perimeter, area, smoothness, compactness, concavity, concave points, symmetry, and fractal dimension. The dataset supports tumor classification as malignant or benign based on the nuclei features.

The DB dataset is tailored for predicting diabetes based on diagnostic measurements. It comprises 768 records of Pima Indian women aged 21 and older, with variables including the number of pregnancies, glucose levels, blood pressure, skin thickness, insulin levels, BMI, a diabetes pedigree function, age, and a binary diabetes outcome. All variables are numerical, representing physiological and diagnostic metrics critical to diabetes prediction.

Figure 2. Correlation matrix for 3 real-world datasets: (A) BP: Body Performance Dataset, (B) BC: Breast Cancer Dataset, and (C) DB: Diabetes Dataset.



Simulated Datasets

simulation In our study, we first generate 10 Gaussian-distributed features and then impose distinct correlation structures using the Cholesky decomposition method [30]. A binary target variable is subsequently constructed based on 4 selected features. The process of defining the target variable is repeated across 3 different correlation structures, with the simulation executed at 3 distinct sample sizes (500, 5000, and 10,000). The use of varying sample sizes allows us to examine the effect of data volume on the robustness and stability of the

correlation structures and the resulting relationships between features and the target variable.

To introduce correlations, we construct 3 types of correlation matrices based on 3 different exponential decay rates, corresponding to varying strengths and patterns of correlation: 0.1 for strong positive correlations, 0.3 for weaker positive correlations, and 0.25 for alternating correlations (positive and negative). The correlation between variables is defined using equation (1) for the 0.1 and 0.3 decay rates, where the exponential decay ensures that correlations decrease as the index distance increases:

```
https://ai.jmir.org/2025/1/e65729
```

×	
_	

Here, α represents the decay rate, controlling the speed at which correlations diminish as the distance |i - j| between indices grows. Smaller values of (eg, 0.1) result in slower decay and stronger correlations over larger distances, while larger values (eg, 0.3) lead to faster decay and weaker correlations.



For the 0.25 alternating correlation, equation (2) is used, incorporating alternating signs to produce correlations that switch between positive and negative values with increasing index distance. In this case, $\alpha = .25$ determines the rate of decay, while the alternating factor $(-1)^{|i-j|}$ introduces the sign changes in the correlations. The resulting correlation matrix, which must fulfill the condition of symmetric positive semidefiniteness, is then decomposed via Cholesky decomposition, allowing us to transform independent normal variables into correlated ones as defined by the specified structure. Examples of such generated correlation matrices are shown in Figure 3.

Figure 3. Correlation matrices used in the simulation study: (A) positive exponential decay rate of 0.1, (B) positive exponential decay rate of 0.3, and (C) alternating positive and negative exponential decay rate of 0.25.

(A)														(D	·/												()											
1.	00 0	0.91	0.82	0.74	0.67	0.61	0.55	0.50	0.45	0.40	0.64	1.0	·	-	1.00	0.74	0.53	0.39	0.31	0.24	0.17	0.13	0.09	0.06	0.54	- 1.0		1.00	-0.78	0.61	-0.47	0.36	-0.28	0.22	-0.16	0.12	-0.09	-0.09	1.00
~ 0.	91 1	1.00									0.68			~ •		1.00	0.74	0.54		0.32	0.24	0.18	0.13	0.09	0.66		~	-0.78	1.00	-0.78		-0.47	0.37	-0.28	0.22	-0.17	0.14	0.35	0.75
- O.	82 (0.91									0.65	- 0.8	1	n •		0.74	1.00	0.73	0.55		0.31	0.23	0.18	0.13	0.55	- 0.8	~	0.61	-0.78			0.61	-0.48	0.37	-0.29	0.23	-0.18	0.07	- 0.50
• · 0.1				1.00				0.68						4.			0.73	1.00	0.75	0.56		0.30	0.22	0.17	0.37		*	-0.47	0.60	-0.78		-0.78		-0.48	0.38	-0.29	0.23	-0.15	
w - 0.					1.00							- 0.6	5	φ.	0.31		0.55	0.75	1.00	0.74	0.55	0.40	0.29	0.22	0.28	- 0.6	~	0.36	-0.47	0.61	-0.78	1.00			-0.47	0.36	-0.28	0.12	0.25
<u>ه</u> و .						1.00	0.90				0.47			۰.	0.24	0.32	0.43	0.56	0.74	1.00	0.74	0.54	0.40	0.30	0.21			-0.28	0.37	-0.48		-0.78	1.00			-0.46	0.36	-0.09	- 0.00
⊷• 0.							1.00	0.90		0.74	0.42	- 0.4		N -	0.17	0.24	0.31	0.42	0.55	0.74	1.00	0.74	0.55	0.41	0.15	- 0.4	~	0.22	-0.28	0.37	-0.48	0.61			-0.77		-0.47	0.06	0.25
= - O.							0.90	1.00	0.90		0.39			æ -	0.13	0.18	0.23	0.30	0.40	0.54	0.74	1.00	0.74	0.55	0.11			-0.16	0.22	-0.29	0.38	-0.47				-0.77		-0.04	
e - 0.								0.90	1.00		0.35	- 0.2	,	ø	0.09	0.13	0.18	0.22	0.29	0.40	0.55	0.74	1.00	0.74	0.08	- 0.2	•	0.12	-0.17	0.23	-0.29	0.36	-0.46					0.02	-0.50
g- 0.										1.00	0.31			10	0.06	0.09	0.13	0.17	0.22	0.30	0.41	0.55	0.74	1.00	0.06		0	-0.09	0.14	-0.18	0.23	-0.28	0.36	-0.47				-0.01	-0.75
⊨ - 0.	64 (0.68				0.47	0.42	0.39	0.35	0.31	1.00				0.54	0.66	0.55	0.37	0.28	0.21	0.15	0.11	0.08	0.06	1.00		Þ	-0.09	0.35	0.07	-0.15	0.12	-0.09	0.06	-0.04	0.02	-0.01	1.00	
		ż	à	- Â	ś	6	ż	ė	ģ	10	Ť	- 0.0	0		i	ż	à	4	ś	é	i	÷.	ė	10	Ť	- 0.0		i	ż	ż	à.	ŝ	Ġ	ż	â	9	10	Ť	·-1.00

The correlation between different types of variables is calculated through a structured process that accommodates binary, continuous, and mixed data types. For each pair of variables, the appropriate correlation metric is selected based on their data types. If at least one variable is binary, the Point-Biserial correlation coefficient is used [31]. The data with the correlated variables is then used to construct a binary target variable, which is defined as a linear combination of the first 4 features from the 10 generated variables, as shown in equation (3):

×

The remaining 6 variables $(X_5, ..., X_{10})$ do not contribute to Y and effectively act as noise variables in the dataset. These noise variables introduce additional complexity by creating scenarios where irrelevant features must be disentangled. This setup mimics real-world scenarios where datasets often contain features that are unrelated or weakly related to the target variable. *Y* is then used to define thresholds based on its median, with a range of SD 10% around the median. Values exceeding the upper threshold are assigned the binary label 1, while those below the lower threshold are assigned 0. For values within the threshold range, binary labels are assigned randomly. It should be noted that while the features X_1, X_2, X_3, X_4 remain continuous, the binary target variable is derived through this thresholding approach applied to the linear combination defined in equation (3).

The complexity in these simulated datasets arises from structured correlation patterns, where the strength, direction, and interplay of correlations among features significantly affect their relationships with the target variable. This correlational complexity can be understood at three levels:

- 1. Feature-target correlation: Variability in how individual features relate to the target, ranging from strong to very weak associations.
- 2. Feature-feature correlation: Associations among features that introduce complicate the disentanglement of their individual contributions to the target.
- Global correlation structures: The overall arrangement of feature-target and feature-feature correlations, encompassing uniform (eg, consistent signs) or mixed configurations (eg, alternating signs).

Based on these levels, the datasets can be categorized into three complexity groups:

- Low complexity: Features exhibit rather strong relationships with the target, minimal or no correlations among features, and homogeneous global correlation.
- Moderate complexity: Feature-target relationships vary, ranging from strong to weak, with moderate feature-feature correlations, and consistent correlation signs.
- High complexity: Feature-target relationships are rather weak, with moderate feature-feature correlations, and alternating correlation signs (Figure 3C).

As complexity increases, the challenges in data analysis and modeling grow substantially. The correlation matrices of both simulated and real data reveal that BP most closely aligns with the 0.25 case (high complexity), BC with the 0.1 case (low complexity), and DB with the 0.3 case (low complexity).

Synthetic Data Generation Methods

We use a range of SDG methods to explore diverse approaches to data synthesis. Statistical methods include synthpop, a widely used statistical model that generates synthetic data by fitting individual features and their conditional distributions based on



RenderX

the observed data structure. Synthpop is particularly well-suited for datasets with both continuous and categorical variables, as it applies models such as classification and regression trees that account for different data types. Another statistical method, copula, uses copula functions to model dependencies among variables, allowing for the generation of multivariate synthetic data by combining marginal distributions with a dependency structure. While copula-based methods are primarily designed for continuous variables, extensions or preprocessing techniques can be used to encode and incorporate categorical variables, such as one-hot encoding or ordinal transformations.

For more advanced generative approaches, we use copulagan, ctgan, and tvae, which are deep learning-based models designed to handle complex data synthesis tasks. Copulagan combines the dependency modeling capabilities of copulas with GANs. It learns the marginal distributions of real data columns and applies ctgan to model normalized data, improving the synthesis of mixed data types. Ctgan uses conditional GANs to address challenges in imbalanced and categorical data. It incorporates techniques like mode-specific normalization to handle high-cardinality categories, enabling precise modeling. Tvae captures complex, nonlinear relationships in tabular data by learning latent representations and generating high-quality synthetic data. In addition, we used the Tabula [29] LLM, which leverages LLMs such as a distilled Generative Pretrained Transformer-2 model, and encodes tabular data into natural language-style representations. This framework allows flexible data generation, incorporating domain-specific contexts and enabling synthesis from textual prompts. While not all models used qualify as LLMs (parameter sizes ≥ 1 billion), we used the term for simplicity.

For the implementation of copula, copulagan, ctgan, and tvae we used the Synthetic Data Vault library (SDV [32]). SDV (Andrew Montanez et al) integrates various methods into a unified framework, facilitating seamless experimentation and evaluation. Although adaptations of synthpop for Python (Sam Maurer et al) exist, we used the native R [24] environment, as it provides the most stable and comprehensive implementation.

Utility and Correlation Matrix Distance Measures

To evaluate the quality of the synthetic data, we use 3 key metrics. First, training on synthetic data and testing their performance on original data, using the F_1 -score as a measure. The F_1 -score is calculated using a classification probability cutoff of 0.5. This approach is often referred to as train-synthetic-test-real. The evaluation differs depending on whether the data is derived from real-world datasets or simulated datasets. For real-world datasets, the original data is split into training and testing sets with an 80/20 split. The 80% training split is used to train the SDG methods, and an equivalent amount of synthetic data (corresponding to the 80% training size) is generated. The quality of this synthetic data is then evaluated by testing it against the original 20% testing split from the real-world dataset. For simulated datasets, 100% of the "real" simulated data is used to train the SDG methods. To evaluate the quality of the synthetic data, a separate test set consisting of 100% newly generated synthetic data was created. The performance is then assessed by testing the synthetic simulated

XSI•F(

data against the "real" simulated data containing the full 100% of the samples. The F_1 -score resulting from training on the original data is represented as a dashed line in the visualizations.

Second, we compute the squared differences between the correlation matrices of the original and synthetic datasets. This metric quantifies the extent to which the synthetic data replicates the pairwise correlations present in the original data. Finally, we use the propensity score mean-squared error (pMSE), which is a metric used to evaluate the utility of synthetic data by measuring the distinguishability between real and synthetic datasets. It is defined as:

X

Where \hat{e}_i represents the estimated propensity score for the *i*-th observation, which measures the probability of a sample being synthetic rather than real. The goal of synthetic data generation is to create data so realistic that the model cannot easily distinguish between synthetic and real samples. Therefore, lower pMSE values indicate better performance, as they imply a higher degree of similarity between the real and synthetic datasets. A pMSE value close to 0.25 (the maximum achievable value when synthetic and real datasets are highly distinguishable) suggests bad synthetic data generation [33]. Normalizing this metric by dividing it with 0.25 leads to values between 0 (indistinguishable) and 1 (highly distinguishable).

Variable Importance Measures

Python machine learning libraries, for example, sklearn, typically provide various methods to calculate variable importance (VIMP). The main two approaches are (1) Gini importance and (2) permutation importance [34]. Gini importance measures the reduction in Gini impurity when a feature is used to split a node. The feature's importance is quantified by the total decrease in impurity across all trees. Features that contribute more to impurity reduction are considered more important, although this method can be biased toward features with more categories or higher cardinality.

Alternatively, permutation importance evaluates a feature's significance by measuring the drop in model performance, typically accuracy, when the feature's values are randomly shuffled. The importance score is derived from the change in performance on out-of-bag samples before and after shuffling. A larger decrease in accuracy indicates greater importance. This method is more robust, accounting for feature interactions and reducing biases, but is computationally more demanding.

Using both Gini importance and permutation importance provides complementary insights: Gini impurity reflects a feature's contribution to better splits within trees, while permutation-based importance directly measures a feature's impact on overall prediction accuracy. Combining both methods offers a more balanced assessment of feature relevance.

Evaluation Design

We conduct 10 sampling iterations for each combination of SDG methods. For deep learning approaches, we evaluate training epoch sizes of 300, 1000, and 10,000 on both simulated and real datasets. For LLMs, we limit the epoch sizes to 300

and 1000 due to significantly higher resource demands and previous findings indicating no performance improvement with larger epoch counts [35]. The batch size is fixed at 500 for the deep learning SDV methods and 64 for LLMs. Specifically, we compute the mean F_1 -score and correlation matrix differences across the 10 samples for each SDG method and epoch size. For the most notable results, we visualize the correlation matrix differences and calculate the VIMP scores for the best and worst-performing methods.

Results

We will first present the results for the simulated data, followed by those for the real data. Since the results from eXtreme Gradient Boosting and Gated Additive Tree Ensembles are nearly identical to those from Random Forest and provide no additional insights, we have omitted them here (Multimedia Appendix 1). Although we anticipated this outcome, we sought to empirically validate it. The analysis will then continue with an examination of the VIMP scores and visualization of the correlation distances for the most notable case, which simulated data consisting of 10,000 samples with an alternating decay parameter of 0.25. This scenario is chosen because it illustrates a case where, despite a large sample size, there is a considerable performance gap between the best- and worst-performing methods.

Correlation Distance and Utility Comparison

Simulated Data

Figure 4 presents the results of our methods on the smallest simulated dataset with 500 samples. For the case of strong positive correlations (0.1), there is virtually no difference in utility between generated and original simulated data. In other words, most models cluster tightly around a RF utility of approximately 0.75. Some models (eg, ctgan and copulagan at 300 and 1000 epochs) have higher correlation matrix distances, indicating weaker preservation of correlation structures. Deep learning models trained with more epochs (eg, 1000 or 10,000, indicated by blue and purple) perform better in terms of correlation matrix distances compared to models with 300 epochs. In terms of utility, epoch sizes do not have a significant effect in this scenario because the data complexity seems not high enough to require prolonged training. The observation that utility remains unaffected by high correlation matrix distances highlights that a poor approximation of the correlation structure is problematic only under specific conditions.

Figure 4. Comparison of the correlation matrix distance and utility metrics (F_1 -score in the top row; pMSE in the bottom row) for the simulated dataset with a sample size of 500. ctgan: Conditional Tabular Generative Adversarial Network; LLM: large language model; pMSE: Propensity Score Mean-Squared Error; synthetic Populations in R; tvae: Tabular Variational Autoencoder.



In the scenario with moderate positive correlations (0.3), the higher correlation distance of ctgan and copulagan at low epoch counts now also negatively affects the RF utility, despite the correlation matrix distance being lower than in the case of 0.1. The pMSE values are overall lower, suggesting that the

https://ai.jmir.org/2025/1/e65729

RenderX

increased complexity primarily affects the RF utility. Models

trained with 10,000 epochs again demonstrate improved

performance, characterized by lower correlation matrix distances

and enhanced RF utility, although the pMSE values are higher.

The relationship between pMSE values, correlation matrix

differences, and RF utility is demonstrated by comparing LLM with 300 epochs and ctgan with 1000 epochs: while LLM exhibits a higher correlation matrix difference, its superior utility results in a significantly lower pMSE value overall. As observed in the 0.1 case, tvae and LLM with high training epochs again rank among the top-performing methods in this scenario, with copula and synthpop achieving the highest performance. The same necessity for extended training epochs as in the 0.1 case suggests that deep learning models likely struggle due to insufficient training data.

In the most complex scenario (0.25), the performance of each SDG method in RF utility is worse than with the original data. This is particularly evident as the tvae and LLM models deviate more significantly from the baseline even with 10,000 epochs. However, these differences have minimal impact on the pMSE values, where copula and synthpop consistently emerge again as the best-performing methods. The high complexity of this simulated dataset primarily manifests as reduced RF utility rather than increased pMSE. However, the differences compared

with the 0.3 scenario are not substantial. Notably, well-performing methods show remarkable robustness, while deep learning approaches with fewer epochs, typically recommended as default settings for practical applications, perform surprisingly poorly by comparison.

Figure 5 illustrates the results obtained on the simulated dataset containing 5000 samples. It is evident that the increased dataset size improves the performance across all cases. Correlation matrix differences are smaller, and in the 0.3 case, almost all methods achieve similarly high levels of performance in terms of RF utility. Notably, the 0.25 case differs significantly from the other two cases, although its results are not substantially different from those observed with the 500-sample dataset. The most notable change is that copulagan and synthpop now emerge more clearly as the leading methods, whereas previously, tvae with high epochs had delivered comparable results. Overall, while deep learning methods benefit from the larger dataset, they still require a high number of epochs to perform well and do not yet match the performance levels of statistical methods.

Figure 5. Comparison of the correlation matrix distance and utility metrics (F_1 -score in the top row; pMSE in the bottom row) for the simulated dataset with a sample size of 5000. ctgan: Conditional Tabular Generative Adversarial Network; LLM: large language model; pMSE: Propensity Score Mean-Squared Error; synthetic Populations in R; tvae: Tabular Variational Autoencoder.



In the results of the simulation dataset comprising 10,000 samples, illustrated in Figure 6, the correlation matrix differences decrease slightly further. In addition, the performance of most deep learning methods improves in terms of RF utility and pMSE values when trained with 300 and 1000 epochs. Increasing the number of training epochs enhances the performance of deep learning methods more compared with 5000 samples but less compared to 500 samples. Otherwise,

the results closely resemble those obtained with the 5000-sample dataset. This suggests that using a larger dataset for synthesis does not yield significant benefits unless the goal is to use deep learning methods with a limited number of epochs. However, the overall results indicate that such methods are generally not advantageous for datasets with a structure similar to that of our simulation study.

```
https://ai.jmir.org/2025/1/e65729
```

RenderX

Figure 6. Comparison of the correlation matrix distance and utility metrics (F_1 -score in the top row; pMSE in the bottom row) for the simulated dataset with a sample size of 10,000. ctgan: Conditional Tabular Generative Adversarial Network; LLM: large language model; pMSE: Propensity Score Mean-Squared Error; synthetic Populations in R; tvae: Tabular Variational Autoencoder.



Real-World Data

Due to the larger number of columns and a broader variety of data types in these datasets, the outcomes naturally exhibit some differences (Figure 7). Regarding the impact of dataset size, the results align closely with those observed in the simulated data for key trends. Specifically, smaller datasets exhibit significantly greater variability across all metrics. For the BC dataset, the copula method captures correlations most effectively, whereas synthpop achieves the best results in terms of RF utility and pMSE. BC is also the dataset where increasing the number of epochs benefits deep learning methods the most. This observation is consistent with findings from the simulated data, despite the real datasets featuring a considerably higher number of columns.

On the BP dataset, an initial observation reveals that copulagan achieves unexpectedly favorable pMSE values. This outcome becomes more comprehensible upon examining the dataset's structure. While BP officially comprises 2 categorical variables (gender and class), it also includes sit-up counts, which is an integer variable that pose statistical modeling challenges. Estimating marginals using diverse distributions, such as the Beta distribution, as a preprocessing step for GANs, proves advantageous in this scenario, especially given the ample data available for these estimations. However, this does not translate into superior RF utility. The association between target and features is not adequately captured by copulagan, resulting in poor RF utility scores. In contrast, synthpop demonstrates the best RF utility and correlation matrix difference performance, although it struggles with achieving competitive pMSE due to the complexity of modeling integer variables. Copula, on the other hand, fails entirely to learn meaningful target-feature associations, yielding extremely low RF utility.

The DB dataset presents the fewest challenges to the methods overall, primarily due to the limited number of continuous variables it contains. All methods perform relatively similarly, reflecting the dataset's inherent simplicity. Compared to the corresponding simulated dataset, one notable difference is that even methods with fewer epochs achieve relatively good performance. Otherwise, the insights gained from the 0.3 case simulation with 500 samples are largely transferable to this real-world scenario. Among the methods tested, synthpop and tvae demonstrate the best performance across all metrics, with synthpop again emerging as the most effective.



Figure 7. Comparison of the correlation matrix distance and utility metrics (F_1 -score in the top row; pMSE in the bottom row) for real-world datasets. BC: Breast Cancer Dataset; BP: Body Performance Dataset; ctgan: Conditional Tabular Generative Adversarial Network; DB: Diabetes Dataset; LLM: large language model; pMSE: Propensity Score Mean-Squared Error; synthetic Populations in R; tvae: Tabular Variational Autoencoder.



Detailed Analysis of a Notable Result

We focus on the two least effective methods in terms of correlation matrix difference (ctgan with 300 epochs and LLM with 1000 epochs) and the best-performing method across all metrics (synthpop) on the 0.25 case of the simulated data consisting of 10,000 samples.

Figures 8-10 display the original correlations, those of the synthetic data, and the resulting correlation matrix differences for synthpop, ctgan, and LLM, respectively. While synthpop generates near-perfect synthetic data, both ctgan and LLM struggle, particularly with high absolute feature-feature correlations, which are often underestimated. In the case of LLM, this issue also extends to feature-target correlations, while ctgan exhibits feature-target correlations that exceed those in the original data. Overall, the underestimation of correlations is more pronounced in LLM than the mixed under- and overestimation seen in ctgan, which explains the larger

correlation matrix differences observed in LLM. However, since the relative correlation ratios in LLM more closely resemble those in the original dataset, it performs better than ctgan in terms of RF utility and pMSE. Figure 11-13 display the VIMP scores (Gini and permutation importance) for synthpop, ctgan, and LLM, respectively. Synthpop shows near-identical results to the original data. The Gini importance for ctgan is promising, but the permutation importance reveals that feature 3 becomes entirely irrelevant. Features 7 and 9, due to their higher correlations with the target, are now relevant. For the LLM, feature 1 becomes nearly irrelevant. However, since feature 3 holds greater significance for the target variable, and no other irrelevant features exhibit substantial permutation importance, this does not detrimentally impact the RF utility or pMSE as severely as observed with the ctgan model. Overall, we conclude that large discrepancies in correlations harm utility only when the ratios between target and feature correlations shift significantly.



Miletic & Sariyar

Figure 8. Correlation matrix of original simulated data (A), the mean correlation matrix of synthetic data (B), and the difference between (A) and (B) for synthetic population decay of 0.25 and sample size 10,000 (C). synthetic Populations in R



Figure 9. Correlation matrix of original simulated data (A), mean correlation matrix of synthetic data (B) and difference between (A) and (B) for ctgan with alternating correlation decay of 0.25, sample size 10,000, and 300 epochs (C). ctgan: Conditional Tabular Generative Adversarial Network.



Figure 10. Correlation matrix of original simulated data (A), mean correlation matrix of synthetic data (B) and difference between (A) and (B) for LLMs with an alternating correlation decay of 0.25, sample size 10,000 and 1000 epochs (C). LLM: large language model.







XSL•FU RenderX

Miletic & Sariyar

Figure 12. VIMP scores for original versus synthetic data generated using ctgan with an alternating correlation decay of 0.25, a sample size of 10,000, and 300 epochs. Gini Importance (left) and Permutation Importance (right). ctgan: Conditional Tabular Generative Adversarial Network; VIMP: Variable Importance.



Figure 13. VIMP scores for original versus synthetic data generated using an LLM with an alternating correlation decay of 0.25, a sample size of 10,000, and 1000 epochs. Gini Importance (left) and Permutation Importance (right). LLM: large language model. VIMP: Variable Importance.



Discussion

Principal Findings

The central finding of our simulation study, which is largely transferable to real-world datasets, is that statistical methods such as copula and synthpop consistently outperform deep learning-based approaches across varying sample sizes and correlation complexities. Notably, synthpop emerged as the most effective method. These techniques demonstrate robust performance with minimal reliance on dataset size or extensive training, highlighting their reliability in preserving statistical properties and utility. However, our analysis of real-world datasets revealed that the copula method struggles when handling integer variables and increasing sample sizes does not mitigate this limitation.

In contrast, deep learning methods yield mixed results. While they benefit from larger datasets and extended training epochs, their performance often falls short of statistical methods, especially when trained with fewer epochs or on smaller datasets. These models struggle to capture the correlation structures, leading to higher pMSE values and diminished utility for downstream tasks. This suggests that deep learning models require careful tuning, including sufficient data and training time, to match the performance of statistical approaches. While the potential for deep learning models to handle datasets with diverse types is promising, the results presented here do not provide sufficient evidence to confirm this advantage over statistical methods. In addition, high performance observed for some deep learning-based approaches may be influenced by overfitting rather than genuine generalization.

```
https://ai.jmir.org/2025/1/e65729
```

RenderX

The results obtained using the LLM method are somewhat disappointing. Despite a large sample size ($\geq 10,000$), this approach does not match the performance of syntheop. While the results are generally acceptable, they highlight that the sheer number of parameters in LLM models is not a decisive factor. Instead, methods specifically designed to directly replicate statistical properties and correlations are often more efficient and effective for tabular data. The probabilistic modeling of LLMs via next-token prediction reaches limitations, particularly when it comes to accurately replicating numerical dependencies. Although the attention mechanism offers promising potential, it does not directly address the preservation of distributions and correlations that are crucial for tabular data. In addition, the significantly longer runtime (hours instead of seconds or minutes), even with 2 high-performance NVIDIA H100 Graphics Processing Units, makes the use of the LLM method difficult to justify for our datasets. However, in cases where tabular data contains many features (more than 30), such as high-dimensional datasets, the runtime of synthpop (which runs on CPU) can become prohibitive when using classification and regression trees. In these cases, the runtime of LLMs may be comparable or even shorter, particularly as the number of rows increases.

Our detailed analysis of correlation matrix differences, VIMP scores, and utility uncovers one central mechanism that leads to either good or poor model performance. We find that a model's utility is primarily influenced by the preservation of relative correlations between features and the target variable, rather than by large correlation matrix differences themselves. Although LLM exhibits greater correlation matrix differences

Miletic & Sariyar

JMIR AI

after 1000 epochs compared to ctgan after 300 epochs, this does not result in worse utility. This is because LLM better preserves the relative correlations, particularly those between the features and the target, which leads to improved RF utility and pMSE. In contrast, while ctgan shows good Gini importance values, its less accurate representation of the correlation value ratios has a greater negative impact on utility. Overall, our findings demonstrate that it is not the absolute magnitude of correlation matrix differences, but the relative correlations between features and the target variable that are critical for model utility.

Our results confirm those found in the literature [36,37] but extend them by incorporating LLMs for the first time and using a simulation approach to assess the impact of various correlation structures on the outcomes. Statistical techniques, such as copula and synthpop, are widely recommended for medical datasets with characteristics similar to those in this study. However, our analysis of the BP dataset highlights the potential usefulness of deep learning methods, particularly when handling multiple variables of diverse data types. In these scenarios, deep learning approaches are anticipated to be able to outperform both synthpop and copula-based methods.

Limitations

A key limitation of this study is that our simulation focused primarily on pairwise correlations. This decision was intentional, as we aimed to restrict our exploration to a small set of scenarios to maintain manageable complexity and derive initial insights. While many of our findings translated well to real-world data, the BP dataset highlighted an important challenge: when dealing with more complex scenarios involving a larger number of variables, diverse data types, and intricate interaction patterns, such as those commonly found in omics or high-dimensional datasets, it becomes essential to design advanced simulation studies that better capture these complexities [38]. In such cases, conventional approaches like Cholesky decomposition or even copula-based methods may no longer suffice [39].

Another limitation of our work is the exclusion of more recent and potentially transformative methods, such as diffusion models [40]. These models have demonstrated exceptional performance in generating high-quality synthetic data, particularly for images, and their application to tabular data represents a promising direction for future research. Moreover, we did not extensively evaluate how our chosen methods perform under scenarios involving temporal or longitudinal data, multimodal datasets, or extreme imbalance in class distributions, challenges that are increasingly relevant in modern data science applications. Addressing these aspects would provide a more comprehensive understanding of the strengths and limitations of SDG methods in diverse contexts.

Further, privacy considerations were not evaluated as part of the synthetic data generation process. While the generative models aimed to preserve data utility and structural similarity, privacy risks such as data leakage or membership inference attacks were not assessed due to our focus in the relationships between correlation structure and utility under different scenarios.

Finally, in synthetic data generation, it is critical to account for biases. If the original data contains biases, the synthetic data is likely to mirror these, potentially leading to discriminatory health care outcomes, particularly for marginalized or underrepresented groups. To mitigate such risks, bias detection and adjustment techniques, such as reweighting, oversampling, and fairness constraints, should be integrated into the data generation process. Beyond bias, ethical concerns also include privacy, informed consent, and accountability. For instance, transparency in the data generation process and clear, informed consent from data contributors are essential for maintaining ethical standards. Regular audits of the synthetic data and associated models are necessary to identify and correct emerging biases and privacy breaching risks.

Conclusions

Statistical methods, particularly synthpop, consistently outperform deep learning-based approaches in preserving statistical properties and utility across diverse datasets, establishing their robustness and reliability. Copula methods show promise but struggle with integer variables, limiting their application in real-world scenarios. Deep learning methods, while resource-intensive and sensitive to hyperparameters, may outperform statistical approaches in handling highly complex datasets with mixed variable types when sufficient training samples and computational resources are available. LLMs, despite their theoretical potential, demonstrated suboptimal performance and high computational costs for the datasets analyzed in this study. Overall, these findings underscore the dominance of statistical methods for synthetic data generation for tabular data, while highlighting the niche potential of deep learning approaches for highly complex datasets, provided adequate resources and tuning.

Acknowledgments

This study was funded by BRIDGE, a joint program of the Swiss National Science Foundation SNSF and Inno Suisse (grant 211751).

Data Availability

Data are deposited in publicly available repositories (where available and appropriate).

Authors' Contributions

MS conceptualized and supervised the study. MM implemented the models and further refined the methodological ideas. MS drafted the manuscript, with both authors reviewing and approving the final version.

```
https://ai.jmir.org/2025/1/e65729
```

Conflicts of Interest

None declared.

Multimedia Appendix 1

F1 Scores for RF, XGBoost, and GATE Across All Datasets Synthesized Using SDG Methods with Varying Configurations (Epochs and Batch Sizes).

[XLSX File (Microsoft Excel File), 16 KB - ai_v4i1e65729_app1.xlsx]

References

- 1. Park N, Mohammadi M, Gorde K, Jajodia S, Park H, Kim Y. Data synthesis based on generative adversarial networks. Proc. VLDB Endow 2018;11(10):1071-1083. [doi: 10.14778/3231751.3231757]
- 2. Abedi M, Hempel L, Sadeghi S, Kirsten T. GAN-based approaches for generating structured data in the medical domain. Applied Sciences 2022;12(14):7075. [doi: 10.3390/app12147075]
- 3. Choi E, Biswal S, Malin B, Duke J, Stewart WF, Sun J. Generating multi-label discrete patient records using generative adversarial networks. arXiv:1703.06490 2018. [doi: 10.48550/arXiv.1703.06490]
- 4. Buczak AL, Babin S, Moniz L. Data-driven approach for creating synthetic electronic medical records. BMC Med Inform Decis Mak 2010;10:59 [FREE Full text] [doi: 10.1186/1472-6947-10-59] [Medline: 20946670]
- 5. Kaur D, Sobiesk M, Patil S, Liu J, Bhagat P, Gupta A, et al. Application of Bayesian networks to generate synthetic health data. J Am Med Inform Assoc 2021;28(4):801-811. [doi: <u>10.1093/jamia/ocaa303</u>] [Medline: <u>33367620</u>]
- 6. Juwara L, El-Hussuna A, El Emam K. An evaluation of synthetic data augmentation for mitigating covariate bias in health data. Patterns (N Y) 2024;5(4):100946 [FREE Full text] [doi: 10.1016/j.patter.2024.100946] [Medline: 38645766]
- 7. Goodfellow IJ, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, et al. Generative adversarial networks. Commun. ACM 2020;63(11):139-144. [doi: 10.1145/3422622]
- 8. Saxena D, Cao J. Generative adversarial networks (GANs). ACM Comput. Surv 2021;54(3):1-42. [doi: 10.1145/3446374]
- 9. Miletic M, Sariyar M. Challenges of using synthetic data generation methods for tabular microdata. Applied Sciences 2024;14(14):5975. [doi: 10.3390/app14145975]
- 10. Assefa S. Generating synthetic data in finance: opportunities, challenges and pitfalls. NY: Rochester; 2020.
- 11. Salehi P, Chalechale A, Taghizadeh M. Generative adversarial networks (GANs): an overview of theoretical model, evaluation metrics, and recent developments. arXiv:2005.13178 2020. [doi: <u>10.48550/arXiv.2005.13178</u>]
- 12. Laptev VV, Gerget OM, Markova NA. In: Kravets AG, AG, Bolshakov AA, Shcherbakov M, editors. Generative Models Based on VAEGAN for New Medical Data Synthesis. Cham: Springer International Publishing; 2021.
- 13. Stadler T, Oprisanu B, Troncoso C. Synthetic Data Anonymisation Groundhog Day. URL: <u>https://www.usenix.org/</u> <u>conference/usenixsecurity22/presentation/stadler</u> [accessed 2024-05-26]
- 14. Gupta A, Bhatt D, Pandey A. Transitioning from real to synthetic data: quantifying the bias in model. arXiv:2105.04144 2021. [doi: 10.48550/arXiv.2105.04144]
- 15. Fan J, Liu T, Li G, Chen J, Shen Y, Du X. Relational data synthesis using generative adversarial networks: a design space exploration. arXiv:2008.12763 2020. [doi: 10.48550/arXiv.2008.12763]
- 16. Vu MH, Edler D, Wibom C, Löfstedt T, Melin B, Rosvall M. A correlation- and mean-aware loss function and benchmarking framework to improve GAN-based tabular data synthesis. arXiv:2405.16971 2024. [doi: 10.48550/arXiv.2405.16971]
- 17. Patel S, Kakadiya A, Mehta M, Derasari R, Patel R, Gandhi R. Correlated discrete data generation using adversarial training. arXiv:1804.00925 2018. [doi: 10.48550/arXiv.1804.00925]
- 18. Torfi A, Fox EA. CorGAN: correlation-capturing convolutional generative adversarial networks for generating synthetic healthcare records. arXiv:2001.09346 2020. [doi: <u>10.48550/arXiv.2001.09346</u>]
- 19. Rajabi A, Garibay OO. In: Degen H, Ntoa S, editors. Distance Correlation GAN: Fair Tabular Data Generation withGenerative Adversarial Networks. HCI Cham: Springer Nature Switzerland; 2023.
- Sariyar M, Hoffmann I, Binder H. Combining techniques for screening and evaluating interaction terms on high-dimensional time-to-event data. BMC Bioinformatics 2014;15(1):58 [FREE Full text] [doi: 10.1186/1471-2105-15-58] [Medline: 24571520]
- Strobl C, Boulesteix AL, Kneib T, Augustin T, Zeileis A. Conditional variable importance for random forests. BMC Bioinformatics 2008;9(1):307 [FREE Full text] [doi: 10.1186/1471-2105-9-307] [Medline: 18620558]
- 22. Chen T, Guestrin C. XGBoost: a scalable tree boosting system. 2016 Presented at: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining; 2016 August 17; USA. [doi: 10.1145/2939672.2939785]
- 23. Joseph M, Raj H. GANDALF: gated adaptive network for deep automated learning of features. arXiv:2207.08548 2024. [doi: <u>10.48550/arXiv.2207.08548</u>]
- 24. Nowok B, Raab GM, Dibben C. Bespoke creation of synthetic data in. J. Stat. Soft 2016;74(11):1-26. [doi: 10.18637/jss.v074.i11]
- 25. Hofert M, Kojadinovic I, Mächler M, Yan J. Elements of Copula Modeling with R. 1st ed. New York: Springer; 2018.

- 26. SDV 0.18.0 documentation. CopulaGAN Model. URL: <u>https://sdv.dev/SDV/user_guides/single_table/copulagan.html</u> [accessed 2024-06-16]
- 27. Lei Xu, L, Skoularidou M, Cuesta-Infante A, Veeramachaneni K. Modeling tabular data using Conditional GAN. arXiv:1907.00503 2019(2). [doi: 10.5260/chara.21.2.8]
- 28. Borisov V, Seßler K, Leemann T, Pawelczyk M, Kasneci G. Language models are realistic tabular data generators. arXiv:2210.06280 2023. [doi: 10.48550/arXiv.2210.06280]
- 29. Zhao Z, Birke R, Chen L. TabuLa: harnessing language models for tabular data synthesis. arXiv:2310.12746 2023 [FREE Full text]
- Edlin R, McCabe C, Hulme C, Hall P, Wright J. Correlated parameters and the cholesky decomposition. In: Edlin R, McCabe C, Hulme C, Hall P, Wright J, editors. Cost Effectiveness Modelling for Health Technology Assessment: A Practical Course. Cham: Springer International Publishing; 2015.
- 31. Kornbrot D. Point Biserial Correlation. Wiley StatsRef: Statistics Reference Online. Hoboken: John Wiley & Sons, Ltd; 2014.
- 32. Patki N, Wedge R, Veeramachaneni K. The synthetic data vault. 2016 Presented at: IEEE International Conference on Data Science Advanced Analytics. Montreal; 2016 October 19; QC Canada. [doi: <u>10.1109/dsaa.2016.49</u>]
- Snoke J, Raab G, Nowok B, Dibben C, Slavkovic A. General and specific utility measures for synthetic data.. arXiv 2017 [FREE Full text] [doi: https://doi.org/10.48550/arXiv.1604.06651]
- 34. Wies C, Miltenberger R, Grieser G, Jahn-Eimermacher A. Exploring the variable importance in random forests under correlations: a general concept applied to donor organ quality in post-transplant survival. BMC Med Res Methodol 2023;23(1):209 [FREE Full text] [doi: 10.1186/s12874-023-02023-2] [Medline: 37726680]
- 35. Miletic M, Sariyar M. Large language models for synthetic tabular health data: a benchmark study. Stud Health Technol Inform 2024;316:963-967. [doi: 10.3233/SHTI240571] [Medline: 39176952]
- 36. Endres M, Mannarapotta VA, Tran TS. Synthetic data generation: a comparative study. 2022 Presented at: International Database Engineered Applications Symposium; 2022 August 24; Budapest Hungary: ACM. [doi: 10.1145/3548785.3548793]
- 37. Little C, Elliot M, Allmendinger R, Samani SS. Generative adversarial networks for synthetic data generation: a comparative study. arXiv:2112.01925 2021. [doi: 10.48550/arXiv.2112.01925]
- 38. Morris TP, White IR, Crowther MJ. Using simulation studies to evaluate statistical methods. Stat Med 2019;38(11):2074-2102. [doi: 10.1002/sim.8086] [Medline: 30652356]
- 39. Barbiero A, Ferrari PA. An R package for the simulation of correlated discrete variables. Communications in Statistics Simulation and Computation 2017;46(7):5123-5140. [doi: 10.1080/03610918.2016.1146758]
- 40. Kotelnikov A, Baranchuk D, Rubachev I, Babenko A. TabDDPM: modelling tabular data with diffusion models. arXiv:2209.15421 2022. [doi: 10.48550/.2209.15421]

Abbreviations

BC: Breast Cancer Dataset
BP: Body Performance Dataset
ctgan: Conditional Tabular Generative Adversarial Network
DB: Diabetes dataset
DistilGPT-2: distilled Generative Pretrained Transformer-2
GAN: Generative Adversarial Network
LLM: large language model
ML: machine learning
pMSE: Propensity Score Mean-Squared Error
RF: Random Forest
SDG: Synthetic Data Generation
SDV: Synthetic Data Vault
SynthPop: Synthetic Populations in R
TVAE: Tabular Variational Autoencoder
VIMP: Variable Importance



Edited by K El Emam; submitted 23.08.24; peer-reviewed by J Lopes, VKC Bumgardner; comments to author 13.10.24; revised version received 24.11.24; accepted 20.01.25; published 20.03.25. <u>Please cite as:</u> Miletic M, Sariyar M Utility-based Analysis of Statistical Approaches and Deep Learning Models for Synthetic Data Generation With Focus on Correlation Structures: Algorithm Development and Validation JMIR AI 2025;4:e65729 URL: https://ai.jmir.org/2025/1/e65729 doi:10.2196/65729 PMID:

©Marko Miletic, Murat Sariyar. Originally published in JMIR AI (https://ai.jmir.org), 20.03.2025. This is an open-access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR AI, is properly cited. The complete bibliographic information, a link to the original publication on https://www.ai.jmir.org/, as well as this copyright and license information must be included.



Publisher: JMIR Publications 130 Queens Quay East. Toronto, ON, M5A 3Y5 Phone: (+1) 416-583-2040 Email: <u>support@jmir.org</u>

https://www.jmirpublications.com/

