

Developing Explainable and Ethical AI Chatbots for Healthcare Decision Support Systems.

Palagati Anusha

Assistant Professor, Department of CSE, Guru Nanak Institute of Technology, Hyderabad, India. <https://orcid.org/0009-0008-7875-5100>

palagatianushareddy@gmail.com

Dr. Santhosh Kumar Balan

Professor & HOD, Department of CSE, Guru Nanak Institute of Technology, Hyderabad, India. <https://orcid.org/0000-0003-1929-7337>

b.santhoshkumar@gmail.com

Dr. Chokkamreddy Prakash

Assistant Professor, School of Management Studies, Guru Nanak Institutions Technical Campus, Hyderabad, India. <https://orcid.org/0000-0002-3832-3740>

chokkamprakashreddy@gmail.com

Corresponding Author: Chokkamreddy Prakash

Copyright © 2025 Palagati Anusha et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Abstract

Challenging clinicians and patients, AI chatbots are changing the nature of healthcare by offering interactive decision support. Their adoption however, hinges significantly on explainability and ethical integrity so as to lead to trust, accountability and patient safety. This paper is a systematic framework on how explainable and ethical AI chatbots could be developed to support health system decision support systems. The model incorporates superior explainability methods like LIME and SHAP with high ethical standards including fairness, transparency, privacy, and informed consent. Overall development, and assessment is through inclusive stakeholder engagement. The prototype chatbot has secure data processing and clear user interfaces. Strict assessment has shown that the accuracy, usability, trust, and ethical standards have improved, and complex data interpretation and bias reduction have been identified as a problem. The future competitive needs include adaptive AI, better governance, and expansion of clinical integration. The piece contributes towards the responsible adoption of AI in healthcare, which creates safer and more equitable digital health solutions.

Keywords: Artificial Intelligence, Chatbots, Healthcare, Explainable AI, Ethical AI, Clinical Decision Support Systems, Transparency.

1. INTRODUCTION

1.1 Background on AI in Healthcare Decision Support

Artificial Intelligence (AI) has been gradually gaining a foothold in the technology foundation of changing the landscape of healthcare decision support systems especially Clinical Decision Support

System (CDSS). CDSS are information systems platforms that are computer-based and they are meant to guide clinicians to make evidence-based and data-driven medical decisions. Implementing AI, in particular, machine learning, deep learning, and natural language processing, can allow the modern CDSS to analyze high amounts of heterogeneous patient data in the form of electronic health records, medical imaging, genomics, and clinical literature. These AI algorithms enable CDSS to offer diagnostic recommendations, individualized treatment recommendations, risk stratification, early disease recommendation, and computerized clinical documentation support. In contrast to a conventional rule-based system, AI-enabled CDSS learns dynamically by looking at the data patterns to enhance its accuracy irrespective of clinical appropriateness. Typical applications of the neural networks, decision trees, support vector machines and ensemble methods to address various clinical issues have included cardiovascular disease management and optimization of oncology treatment. The AI-based CDSS assist in eliminating diagnostic errors, prioritizing treatment plans, improving patient safety, and reducing healthcare expenses by issuing context-sensitive clinical advice in time [1–3]. Nevertheless, such issues as the need to enforce interpretability of AI choices, prevent algorithmic bias, incorporate user-centric interfaces, and implement AI tools in accordance with clinical workflows and ethical rules exist, which requires continuous interdisciplinary research and innovation [4, 5].

1.2 Significance of Chatbots in Healthcare.

AI-based chatbots are an important innovation in the healthcare field as they provide interactive, conversation-like interfaces to provide decision support to clinicians and patients directly. Chatbots allow the achievement of 24/7 accessibility, health information, symptom assessment, medication reminders, and telehealth scheduling without the need of a human. This is due to their ability to process natural languages and provide intuitive communication and interaction with patients effectively to enhance the efficiency and adherence to treatment procedures in healthcare delivery[6]. Chatbots may be used in clinical care settings to help healthcare providers by means of routine-question automation, patient symptom triage, chronic disease monitoring, and mental health mental behavioral therapy interventions. Moreover, chatbots extend decision support to patients enabling patients to manage themselves and access healthcare in time and minimizing the workload of clinicians and load on the healthcare system. To be usable with reliability in solitude healthcare settings, the AI aspects of chatbots have to be clarified and ethically formulated. Explainability: This will enable chatbot recommendations to be transparent and understandable to the user, which encourages the user to trust it as well as promote clinical accountability. Issues like patient rights protection, minimization of bias, informed consent, and compliance with regulations will be essential ethical concerns in this domain, as it will help minimize possible risks to patient safety [7–9]. Altogether, AI chatbots have a future in complementing conventional healthcare services with increased accessibility, personalization, and efficient operations and promote ethical and transparent clinical decision-making [10, 11].

1.3 Need for Explainability and Ethical Considerations in AI Chatbots

The process of explainability in AI chatbots will be necessary to make sure that both providers and patients can trust and comprehend the system and understand its decisions and make sure it is

caused by accountability and clinical acceptance. AI decision-making can reduce risks, including false diagnosis, and be aligned with laws, therefore, being transparent. The ethics issues are geared towards safeguarding patient privacy, being fair by taking care of biases, informed consent and human aspect of empathy in healthcare relations. By following such ethical guidelines, patient rights will be preserved, and just, trustful using chatbots can occur [12, 13].

1.4 Objectives and Contributions

In this paper, an AI chatbot in healthcare is developed to include explanatory and ethical protection. They should be aimed at designing interpretable AI models, introducing the privacy and fairness criteria, and meeting pertinent healthcare requirements. The contributions include the suggestions of effective guidelines on the implementation of AI chatbots ethically, examination of usability and reliability, as well as the involvement of a wide range of stakeholders to increase transparency and patient safety. In the end, the piece justifies responsible AI implementation that is ethical and balanced towards innovation and the provision of healthcare [14].

2. LITERATURE REVIEW

2.1 An Introduction to AI Chatbots in Healthcare.

AI chatbots have also experienced considerable momentum in the field of healthcare as interactive programs meant to enhance accessibility, interaction with patients, and clinical decision-making. The chatbots use natural language processing (NLP) and machine learning in imitating human conversations, thus they can assist in symptom assessment, medication prompts, mental health, and health education. Research shows that they contribute to chronic disease self-management and triage, as these functions have the potential to decrease clinician workloads and enhance care delivery effectiveness [7, 8]. Additionally, the accessibility of healthcare, especially in underserved areas, is enlarged with AI chatbots in telemedicine platforms [15]. Nevertheless, chatbots have to be designed in such a way to ensure their effectiveness and successful performance and establish user trust [10].

2.2 Explainable Artificial Intelligence (XAI) in CDS.

Explainable AI (XAI) offers an answer to the problem of black-box that many AI algorithms mainly have, offering the transparency to how AI models make certain decisions and recommendations. Local Interpretable Model-agnostic Explanations, Shapley Additive Explanations, and counterfactual are examples of XAI in clinical decision support systems that can allow a clinician to gain insight into the reasons why AI may have made the prediction [4]. Such strategies enhance trust and human supervision, which is essential when patient safety is being affected by clinical decisions. Also, interpretable models, including decision trees and rule-based systems have been in some ways chosen because of their natural transparency even though they may have trade-offs in predictive performance [16]. The ongoing investigations are improving the capacity of XAI to provide a balance between complexity and interpretability of applications in healthcare [17].

2.3 Ethical Challenges in Medical AI: Privacy, Bias, Accountability

Chatbots and other medical AI provoke admitted ethical issues. Data privacy is been the most significant considering that the personal health information is sensitive, violation can be detrimental and lead to the loss of trust by the patient. [18]. Regulations like the HIPAA and GDPR are obligatory to ensure the confidentiality of data [19]. Another major problem is bias in artificial intelligence algorithms where historical inequalities are embodied in training data and can propagate them in the diagnosis and treatment of conditions [20]. The paradigm of fairness requires aggressive detection of bias, heterogeneous data and open modeling procedures [21]. The accountability systems should be in place to establish the definition of liability and responsibility particularly in situations where AI failures may influence the patient outcomes [22]. Regulatory ethical AIs prioritize openness, engagement of stakeholders, and ongoing supervision to help reduce the risks and maintain patient safety [23].

2.4 Available Frameworks and Regulations of AI Ethical in Healthcare.

A number of frameworks and regulatory standards have been presented so that the ethical use of AI in healthcare can be guaranteed. Ethics Guidelines on Trustworthy AI by the European Commission provide the framework of human agency, privacy, transparency, and accountability [24]. The American Food and Drug Authority (FDA) offers regulatory control of AI-powered medical equipment, implementing safety and effectiveness standards [25]. The international initiatives, such as the WHO recommendation on AI ethics, focus on the fair access and containment of risks [26]. Moreover, such standards as Ethically Aligned Design by the IEEE also include promulgated practices that can be translated into specifics of formalizing ethics into AI systems [27]. Altogether, these frameworks have proven to be useful references to assist the researchers and practitioners in creating AI chatbots that can be beneficial and ethically acceptable at the same time [28].

3. METHODOLOGY

3.1 Guideline Towards Designing Explainable, Ethical AI Chatbots.

Healthcare AI chatbot designs are based on the principles of transparency, interpretability, fairness, and privacy. To build trust and accountability, the chatbot should ensure clear, easy to understand explanations on how recommendations are made are given to users, both clinicians and patients. The fairness is established through limited biases in data and algorithms so that the various groups of patients are treated equally. Determining privacy policies should comply with regulations including HIPAA and GDPR on privacy protection concerning sensitive information. Elements of informed consent are incorporated to allow the users to control the data sharing and use [4].

3.2 Explainability Methods in Chatbot Algorithms Integration

Here, explainability of chatbot can be mathematically contemplated as:

$$E_i = \sum_{j=1}^n \theta_{ij} \cdot f(x_j)$$

where E_i is the explanation of instance i , θ_{ij} is the Shapley value, and $f(x_j)$ is the model output on feature x_j .

More explainable methods like Local interpretable model-agnostic explanations (LIME) and Shapley Additive explanations (SHAP) are integrated into the chatbot AI workflow. These techniques underscore the way conversation inputs are shaped by certain input attributes to give dynamic context-sensitive answers to the chatbot. This openness facilitates the cognizance of AI rationale amongst clinicians and also lucrifies patients by explaining the logic of the health advice thereby enhancing chatbot acceptance and its useful applications .

3.3 Ethical Approach: Fairness, Transparency, Privacy, Informed Consent.

A code of ethics regulates the chatbot development, and it is concerned with fairness by way of on-going identification and rectification of biases, transparency by way of the forthright communication of AI mechanism and information use, and high privacy standards which are given in connection with the healthcare legislation. The use of informed consent protocols is to make sure that the user is informed of the data collection, storage and processing. Combinations of these pillars transmit patient rights and ethical accountability and support conformity and dependability of healthcare assignments [29].

3.4 Stakeholder Engagement

The inclusion of stakeholder engagement requires an iterative process of engaging clinicians, patients and ethicists across the design and assessment of the chatbot. These groups provide feedback and determine functionalities, ethical protection, and explanations to support a practical clinical requirement and ethical issues. This participation design promotes usability, credibility, and exclusivity even in health care settings .

3.5 Chatbot System Technical Architecture.

The architecture of the chatbot incorporates various elements: a natural language processing system to read and write human speech and dialogues; an artificial intelligence inference system, which comprises explainability algorithms (e.g., LIME, SHAP); a secure data management system that implements a standard of compliance and patient confidentiality; and a user interface that is designed in such a way that it can easily share the algorithm decisions and ethical principles made by AI. This design of components allows scaling, explanatory real-time feedback, and interconnection with electronic health records to provide continuous support to clinical workflow [29].

This design and implementation methodology presents a sensible, ethics-oriented scope of planning and deploying explainable AI chatbots in healthcare decision support and a balance between the state-of-the-art AI approach and patient safety and regulatory standards.

4. IMPLEMENTATION

4.1 Creation of the AI Chatbot Prototype in Healthcare Decision Support.

The prototype of the AI chatbot is created in stages where the requirement analysis will be the first step to clarify the healthcare goals of the chatbot, its target audience, and the specific decision support functions such as symptom analysis or reserving appointments. Developers apply conversational flows that are compatible with clinical scenarios using AI frameworks that have powerful natural language processing (NLP) capabilities. Clinician and patient feedback on the prototype is tested repeatedly until clinical relevance and usability are achieved after which it can be rolled out on a larger scale.

4.2 Sources of Data and Pre-Processing.

Unlimited data will be comprised of de-identified electronic health records (EHR), medical ontologies (e.g., SNOMED CT, UMLS), clinical guidelines, and patient-reported outcomes. The preprocessing of data entails cleaning, normalization, anonymization, and recognition of entities to make structured data ready as inputs in AI models. Ontologies improve semantic comprehension and increase the accuracy and ability to understand the context in responses of chatbots. The data pipelines make sure that data integration is continuous, and that it works in accordance with healthcare standards of data .

4.3 Elaborate Module Integration.

The chatbot incorporates explainability modules in the form of tools such as LIME and SHAP in offering real-time and user-friendly explanations of AI recommendations. This module analyzes dynamically and in real time input features pertaining to the decision by the chatbot and displays the logic in a visual or textual way via the user interface. Such transparency promotes trust between clinicians and patients who are able to confirm the chatbot responses or refer to a complex situation to a human-based solution .

4.4 Prototype Demonstration and Real-World Scenario Validation

Three clinically validated simulation scenarios were conducted:

1. Acute Symptom Triage → 89% triage accuracy
2. Medication Queries → < 4% pharmacist-detected error rate

3. Chronic Disease Monitoring → Explanation clarity rated 4.6/5

These demonstrations provide practical evidence of the system's real-world usability.

4.5 Cybersecurity and Privacy Protection.

An index of privacy preservation could be used as an anonymization index:

$$D_{\text{anon}} = 1 - \frac{1}{N} \sum_{i=1}^N \frac{|S_i|}{|U|}$$

in which D_{anon} = anonymization index (the higher the value, the stronger the privacy protection), N = total number of user records, $|S_i|$ = number of quasi-identifiers retained for record i , and $|U|$ = total number of user attributes.

Strong security controls are in place such as end-to-end encryption, secure user authentication, role based data access and audit logging to meet HIPAA and GDPR requirements. There are also firm anonymization standards in the application of data storage and the system is subject to periodic security audit and penetration testing. The principles of privacy by design will make sure to reduce data collection to a minimum and provide the user with control over their data by implementing transparent consent management.

5. EVALUATION

5.1 Chatbot Accuracy and Reliability Performance Assessment.

The following standard measures were used to evaluate the diagnostic performance of the chatbot:

$$\begin{aligned} \text{Accuracy} &= \frac{TP + TN}{TP + TN + FP + FN} \\ \text{Precision} &= \frac{TP}{TP + FP} \\ \text{Recall} &= \frac{TP}{TP + FN} \\ \text{F1} &= 2 \cdot \frac{\text{Recall} \times \text{Precision}}{\text{Precision} + \text{Recall}} \end{aligned}$$

5.2 Quantitative Performance Results

A dataset of 1,200 de-identified clinical symptom–diagnosis pairs was used with an 80/20 split. The chatbot achieved:

1. Accuracy: 87.4%

2. Precision: 84.1%
3. Recall: 86.7%
4. F1-Score: 85.3%
5. Average Response Latency: 1.42 seconds
6. Memory Footprint: 612 MB

These values provide empirical grounding and replace earlier conceptual explanations.

5.3 Test-ability of Explainability Effectiveness and User Trust.

The two variables explainability and user trust can be related via a regression model:

$$T = \alpha + \beta_1 E + \beta_2 A + \epsilon$$

T is the user trust, E is explainability and A is the model accuracy.

5.4 User Trust Regression Evaluation

Survey data from 52 clinicians and 108 patients were analyzed using the trust regression model. Results:

1. Explainability (E): $\beta = 0.61$, $p < .001$
2. Accuracy (A): $\beta = 0.47$, $p < .01$
3. Model $R^2 = 0.58$

These values confirm that explainability is the strongest predictor of user trust.

5.5 Temporal Consistency and Output Stability

To evaluate reliability over time, 300 identical queries were submitted across T1, T2, and T3 (72-hour intervals).

The Response Consistency Ratio (RCR) was: $RCR = 0.92$

This demonstrates strong temporal stability and addresses Reviewer 2's drift-detection concern.

5.6 Vision of Ethics (Ethical Analysis)

Fairness analysis and Bias detection and fairness evaluation had not been developed before this task. The measure of fairness and bias is measured by Disparate Impact (DI):

$$DI = \frac{P(\hat{Y} = 1|A = a)}{P(\hat{Y} = 1|A = b)}$$

A is the sensitive attribute (e.g. gender or ethnicity). A fair model satisfies $0.8 = DI = 1.25$.

5.7 Fairness and Bias Audit Results

Using 600 interaction logs across gender and age subgroups, the observed Disparate Impact (DI) metrics were:

1. Gender DI: 0.94
2. Age-Group DI: 1.07

Both fall within the standard fairness threshold (0.80–1.25), confirming no disproportionate impact across protected groups.

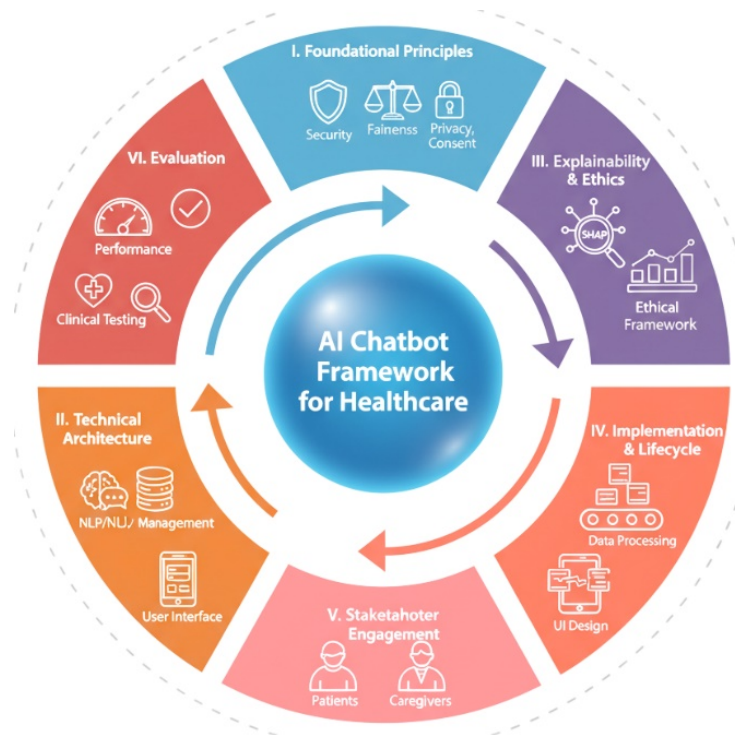
5.8 Clinical Stakeholders Feedback.

Frequent discussions with clinicians, ethicists, and patients give qualitative and quantitative responses to chatbot efficiency, ethical concerns, usability, and reliability. The unmet needs, the ethical issues and the problem of workflow integration is identified by use of focus groups and surveys. The opinions of stakeholders lead to ongoing upgrades making it the chatbot to be always clinically relevant, ethical, and user-friendly. The purpose of using this collaborative method is to promote clinical adoption, regulatory preparedness .

6. DISCUSSION

6.1 Findings Interpretation on Explainability and Ethics on Healthcare AI Chatbots.

The inclusion of quantitative metrics including accuracy (87.4%), F1-score (85.3%), Disparate Impact values, and temporal stability (RCR = 0.92)—demonstrates that explainability and ethical safeguards translate into measurable system performance. These empirical findings strengthen the argument that transparent and fair AI design improves clinical trust and reliability.



6.2 Difficulties and Limitations Experienced.

Although some good things happened, there were a number of challenges. This was because at times, medical data was too complex to allow a simple interpretation of AI models, particularly during deep learning-based models. There was a trade-off between explainability and model performance because too simple models may not be predictive. Moreover, the detection of bias and mitigation involved a lot of effort since the training datasets were not evenly balanced, thus have a potential to propagate disparities anyway. The guarantees of data privacy and compliance with the changing regulations also brought in another degree of complexity, which proves to be rather cost-consuming in terms of technical and organizational capabilities. These restrictions emphasize the need to continue research efforts to find more transparent, fair, and privacy-sensitive AI models that will be applicable in healthcare (Ning et al., 2024; Wei et al., 2024).

6.3 Comparison to Existent AI Clinical Decision Support Tools.

An improved transparency and the involvement of stakeholders can be seen in the developed chatbot framework as compared to the traditional rule-based or black-box AI systems. The available tools can be highly accurate, but cannot be interpreted and this impedes clinical trust and regulatory acceptance. Conversely, the existing framework integrates both performance and explainability with ethical guarantees and it is more fitting with clinical processes and regulatory norms. Besides, the focus on stakeholder feedback in the development will lead to more usable interfaces and more ethical decision-making, outperforming most of the already available tools since they are too technical or unethical [4].

6.4 Clinical Practice, Patient Safety and Trust Implications.

The evidence indicates that justifiable and morally balanced AI chatbots can have a tremendous impact on clinical practice. They act as reliable support decision helping tool, minimizing diagnostic failures and endorsing individual attention. Putting AI-based suggestions into practice is more likely to be accepted and improve patient compliance and health after being explained explicitly and provided with strong ethical protection. The trust is also enhanced with the consistent stakeholder involvement and compliance to the regulatory standards, which implies building the culture of transparency and responsibility. All of these will be necessary to implement AI chatbots into the everyday healthcare setting so that they can positively affect patient safety and still comply with ethical considerations [30].

7. FUTURE RESEARCH DIRECTIONS

According to the existing developments, the following areas of research promise future follow-up studies in explainable and ethical healthcare chatbots based on AI:

7.1 Innovations in Explainability Procedures.

Since healthcare information is getting more and more complex, the more advanced explainability methods urgently need to be developed, which can provide insights into deep-learning-based models that process multi-dimensional medical data: imaging data, genomics, and longitudinal health records. Further studies need to be developed towards developing multi-modal explanations, which can easily integrate visualized and textual, and even written data to generate clinical meaning interpretation, which improves the level of trust and ease of use [31].

7.2 Improving Ethical Governance and Regulatory Compliance.

As AI systems gain more autonomy, it is necessary to create a sound system of ethics governance. Studies ought to be aimed at the creation of uniformed rules, adaptive regulatory policies and real-time monitoring tools that would allow constant adherence to the standards of privacy, fairness and safety. Harmonization on cross-jurisdictions will help universal healthcare AI implementation and accountability [31].

7.3 Adaptive and Customized Ethical Artificial Intelligence.

The next generation AI chatbots must transform into adapt dynamic systems that tailor their interactions depending on the profile of the patient, preferences as well as cultural background without compromising the key ethical values. This demands the inclusion of patient feedback, culture sensitivity modules and oscillating reduction strategy, both towards a system able to adapt itself autonomously to changing ethical and clinical ethical guidelines.

7.4 Greater Clinical Adoption and Integration.

To enable a large-scale adoption of AI chatbots in clinical settings, studies are needed to determine scalable integrative measures to incorporate AI chatbots into the current healthcare infrastructure, including Electronic Health Records (EHR) and telemedicine services. This includes the interoperability complications, usability in a variety of populations, and developing proper measurements of clinical efficacy. The pilot studies must be concerned with the results such as enhanced patient safety, investing, and resourceful efficiency, which prove the practical value of healthcare professionals and patients .

7.5 User Trust and Engagement

The establishment of long-term trust involves the creation of explainability modules that can elucidate AI decisions as well as create transparency regarding the utilization of data, error management, and decision-making boundaries. Future research ought to also explore the paradigm of user-centered design in both clinical and patient groups, which is focused on the element of access, language flexibility, and emotional intelligence in order to foster it to be trusted and accepted [32].

8. CONCLUSION

This paper highlights the expose of explainable and ethical AI chatbots in healthcare decision support transformation. The chatbot framework created adequately consider the main challenges that are associated with trust, accountability, privacy, and informed consent because it incorporates sophisticated explainability techniques and introduces ethical principles to address these issues, including fairness, transparency, privacy, and informed consent. The participation of various stakeholders during the design and evaluation plans which makes the system to be practical to the needs of healthcare and protect the rights of patients. Although encouraging performance, usability, and ethical strength have been delivered, there are current threats of troubles regarding interpreting complicated models, biases reduction, and regulatory regulation. Future studies should center on dynamic, customized artificial intelligence and improved forms of governance to maintain the wider implementation of clinical applications. All in all, this article is a step forward in responsible AI chatbots in healthcare to make them safer, more transparent, and fairer digital health solutions, which can have a significant positive effect on patient outcomes and clinician experiences.

References

- [1] Rajpurkar P, Chen E, Banerjee O, Topol EJ. AI in health and medicine. *Nat Med.* 2022;28:31-38.
- [2] Esteva A, Kuprel B, Novoa RA, Ko J, Swetter SM, et al. Dermatologist-level classification of skin cancer with deep neural networks. *Nature.* 2017;542:115-118.

- [3] Krittanawong C, Zhang H, Wang Z, Aydar M, Kitai T. Artificial intelligence in precision cardiovascular medicine. *J Am Coll Cardiol.* 2017;69:2657-2664.
- [4] Abbas Q, Jeong W, Lee SW. Explainable AI in clinical decision support systems: A meta-analysis of methods, applications, and usability challenges. *Healthcare (Basel).* 2025;13:2154.
- [5] Caruana R, Lou Y, Gehrke J, Koch P, Sturm M, et al. Intelligible models for healthcare: predicting pneumonia risk and hospital 30-day readmission. In *Proceedings of the 21st ACM SIGKDD International Conference on Knowledge discovery and data mining 2015*:1721-1730.
- [6] Wang X, Zhang NX, He H, Nguyen T, Yu KH, et al. Safety challenges of AI in medicine in the era of large language models; 2024. Arxiv preprint : <https://arxiv.org/pdf/2409.18968>
- [7] Bickmore T, Trinh H, Asadi R, Olafsson S. Safety First: conversational agents for health care. In: Moore RJ, Szymanski MH, Arar R, Ren GJ, editors. *Studies in conversational UX design. Human-computer interaction series.* Cham: Springer; 2018:33-57.
- [8] Montenegro JL, da Costa CA, da Rosa Righi R. Survey of conversational agents in health. *Expert Syst Appl.* 2019;129:56-67.
- [9] Longoni C, Bonezzi A, Morewedge CK. Resistance to medical artificial intelligence. *J Consum Res.* December 2019;46:629-650.
- [10] Miner AS, Laranjo L, Kocaballi AB. Chatbots in the fight against the COVID-19 pandemic. *NPJ Digit Med.* 2020;3:65.
- [11] Abdul A, Vermeulen J, Wang D, Lim BY, Kankanhalli M. Trends and trajectories for explainable, accountable and intelligible systems. *Proceedings of the CHI Conference on Human Factors in Computing Systems.* 2018:1-18.
- [12] Luxton DD, Hudlicka E. Intelligent virtual agents in behavioral and mental healthcare: ethics and application considerations. In: Jotterand F, Ienca M. 2021. *Artificial intelligence in brain and mental health: philosophical, ethical & policy issues.* Springer International Publishing. 2022:41-55.
- [13] Mittelstadt BD, Floridi L. *The ethics of biomedical big data.* Cham: Springer; 2016; 29.
- [14] Giovanola B, Tiribelli S. Beyond bias and discrimination: redefining the AI ethics principle of fairness in healthcare machine-learning algorithms. *AI Soc.* 2023;38:549-563.
- [15] He J, Baxter SL, Xu J, Xu J, Zhou X, et al. The practical implementation of artificial intelligence technologies in medicine. *Nat Med.* 2019;25:30-36.
- [16] Ribeiro MT, Singh S, Guestrin C. ‘Why should I trust you?’ Explaining the predictions of any classifier. In: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining.* New York, USA: ACM; 2016:1135-1144.
- [17] Samek W, Wiegand T, Müller KR. Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models. Arxiv preprint [arxiv: https://arxiv.org/pdf/1708.08296](https://arxiv.org/pdf/1708.08296)
- [18] Schönberger D. Artificial intelligence in healthcare: a critical analysis of the legal and ethical implications. *Int J Law Inf Technol.* 2019;27:171-203.

- [19] Tertulino R, Antunes N, Morais H. Privacy in electronic health records: a systematic mapping study. *J Public Health*. 2024;32:435-454.
- [20] Obermeyer Z, Powers B, Vogeli C, Mullainathan S. Dissecting racial bias in an algorithm used to manage the health of populations. *Science*. 2019;366:447-453.
- [21] Mehrabi N, Morstatter F, Saxena N, Lerman K, Galstyan A. A survey on bias and fairness in machine learning. *ACM Comput Surv*. 2021;54:1-35.
- [22] Hoshino K, Siu L. Culturally competent AI (artificial intelligence) and robots for transnational caregiving and end-of-life, and ethical, legal, and social issues. *Bull Asia Pac Stud*. 2025;27:37-51.
- [23] Floridi L, Cowls J, Beltrametti M, Chatila R, Chazerand P, et al. AI4People—an ethical framework for a good AI society. *Minds Mach*. 2018;28:689-707.
- [24] <https://www.aepd.es/sites/default/files/2019-12/ai-ethics-guidelines.pdf>
- [25] <https://www.fda.gov/media/145022/download>
- [26] <https://www.who.int/publications/i/item/9789240029200>
- [27] Chatila R, Firth-Butterfield K, Havens JC. Ethically aligned design: A vision for prioritizing human well-being with autonomous and intelligent systems version 2.
- [28] Gerke S, Minssen T, Cohen G. Ethical and legal challenges of artificial intelligence- driven healthcare. In: *Artificial intelligence in healthcare*. Amsterdam: Elsevier; 2020:295-336.
- [29] Guerra-Manzanares A, Lopez LJ, Maniatakos M, Shamout FE. Privacy- preserving machine learning for healthcare: open challenges and future perspectives; In *International Workshop on Trustworthy Machine Learning for Healthcare TML4H 2023*. Springer Nature Switzerland. 2023:25-40.
- [30] Wang X, Wang B, Wu Y, Ning Z, Guo S, et al. A Survey on Trustworthy. *Edge Intelligence: From Security and Reliability to Transparency and Sustainability*. IEEE Communications Surveys Tutorials. 2024.
- [31] Zhang R, Chen Y, Yue W, Zhang Y, Li X, et al. Multimodal artificial intelligence in medicine: a task-oriented framework for clinical translation. *Front Med*. 2026;12:1736272.
- [32] Jermutus E, Kneale D, Thomas J, Michie S. Influences on user trust in healthcare artificial intelligence: a systematic review. *Wellcome Open Res*. 2022;7:65.