



# A Multimodal AI-Based Medical Chatbot for Voice and Vision-Enabled Healthcare Assistance

<sup>1</sup>Ganji Mahithi, <sup>2</sup>Goranti Ramya, <sup>3</sup>Dr.Shivani Yadao

<sup>1,2</sup>PG Scholars, <sup>3</sup>Asso.Professor

<sup>1</sup>Department of Computer Science Engineering,

<sup>1,2,3</sup>Stanley College of Engineering and Technology for Women, Hyderabad, India.

**Abstract:** Artificial Intelligence (AI) is increasingly being used to improve access to basic health care support. This work presents an AI-powered medical assistant that integrates voice interaction, image analysis, and conversational intelligence to provide preliminary health guidance. Users can describe their symptoms through speech, which is converted into text, and upload images of visible conditions such as skin issues or injuries. The system analyzes both textual and visual inputs using an advanced language model to generate informative responses, including possible causes and general health suggestions. To enhance accessibility, the responses are also delivered through text-to-speech output, enabling a more natural interaction. Experimental results indicate that the system performs effectively in understanding user input and providing useful feedback. Although it does not replace professional medical diagnosis, the proposed system can assist users in gaining better awareness of their health conditions and encourage timely medical consultation.

**Index Terms -** Artificial Intelligence, Health care Assistance, Speech Recognition, Image Processing, Multi-modal Interaction, Medical Diagnosis Support, Voice-Based Interface, Vision-Based Analysis

## I. INTRODUCTION

Health care is very important for keeping people and communities healthy. But many people find it hard to get the medical help they need when they need it. Some of the main problems include waiting a long time to see a doctor, expensive visits, and not having enough doctors, especially in areas that are far away or hard to reach. Because of this, people sometimes wait too long before getting help or look for information online from places that aren't trusted, which can make them misunderstand their health problems. With AI developing so quickly, there are new chances to solve these problems using smart and easy-to-use digital tools. AI can handle a lot of information and find useful patterns, which makes it great for uses in health care like predicting diseases, helping with diagnosis, and providing virtual health support. Especially, chatbots and big language models are being noticed because they can talk to people in a way that feels natural and easy to use. Traditional health care systems that use AI mostly depend on text, which can make it hard for them to fully understand a person's health situation. On the other hand, multi-modal AI systems use different types of input, like speech, text, and images, which allows for a better understanding of a person's symptoms. For example, images of skin issues can give important details that might not be clear just from words. In this we're developing a medical assistant powered by multi-modal AI that combines voice interaction, image analysis, and conversation understanding to offer basic health care support. Users can talk about their symptoms, and the system converts their speech into text using speech recognition. They can also upload pictures of medical problems, like rashes or wounds, for further evaluation. By looking at both text and images, the system can give helpful responses that include possible reasons for the symptoms and general advice on staying healthy. The main goal of this system is to make health care information more accessible and provide a user-friendly tool for initial guidance. Using voice helps people who might find it hard to type. The system also helps by explaining symptoms in simple terms and encourages users to consult a doctor when needed.

## II. LITERATURE SURVEY

The integration of Artificial Intelligence (AI) into health care has led to the development of advanced systems that support medical diagnosis, patient care, and health monitoring. Recent studies have increasingly focused on leveraging multi-modal AI, large language models, and conversational agents to enhance the efficiency and accessibility of health care services. These technologies aim to bridge gaps in traditional health care systems by providing intelligent and scalable solutions. A notable contribution in this area is the multi-modal medical chatbot proposed by [1] Agarwal et al. (2025), which incorporates large language models with multiple input formats, including text and images. This approach enables users to describe their symptoms more effectively, leading to improved interpretation and more relevant preliminary guidance. The study demonstrates that combining different input modalities enhances system performance and broadens accessibility. In a related direction, [2] Qiu et al. (2023) investigated the role

of large AI models in health informatics. Their work highlights how such models can process complex medical data sets and assist in tasks such as clinical decision-making, disease prediction, and biomedical data analysis. The findings suggest that large-scale AI systems have strong potential to support health care professionals by providing data-driven insights. One of the major challenges in AI-based health care systems is ensuring the accuracy and reliability of generated information. To address this concern, [3] Lecu et al. (2025) introduced a method that integrates knowledge graphs with retrieval-based techniques. This combination helps reduce incorrect or misleading outputs by grounding AI responses in verified medical knowledge, thereby improving trustworthiness. The application of generative AI in health care has also been widely explored. [5] Sai et al. (2024) provided a comprehensive overview of generative models and their role in medical applications, including diagnosis support, patient data management, and automated health care assistance. Their study emphasizes the growing importance of generative techniques in handling complex health care tasks. Another significant advancement is the incorporation of multilingual capabilities in AI-driven healthcare systems. [6] Bazzi Mohamed Salim et al. (2025) developed a retrieval-augmented framework designed to deliver accurate medication-related guidance in multiple languages. This approach enhances communication between patients and AI systems, particularly in linguistically diverse populations. AI-powered chatbots for disease prediction have also been extensively studied. [7] Chakraborty et al. (2022) designed a conversational system that analyzes user-reported symptoms using natural language processing and machine learning techniques to predict infectious diseases. The system demonstrates the effectiveness of interactive AI tools in providing early-stage health insights. In chatbot-based solutions, mobile health care applications have emerged as practical tools for improving access to medical services. [8] Kalnoor et al. (2025) introduced an AI-enabled mobile application that offers personalized health care support, including symptom-based doctor recommendations and health monitoring features. Such applications highlight the role of AI in delivering convenient and user-centric health care solutions. Addressing language diversity remains an important aspect of health care accessibility. [9] Tejasri and Lawanya Shri (2026) examined multilingual and code-switched conversational systems, emphasizing the importance of natural language processing techniques that enable effective communication across different languages and dialects. These systems play a key role in making health care support more inclusive. Large language models have further demonstrated their utility in biomedical research. [10] Li et al. (2026) explored the use of AI models for identifying causal relationships in biomedical literature. Their findings indicate that such models can assist researchers in extracting valuable insights from large-scale medical data. The comparative performance of traditional machine learning methods and advanced language models has also been studied. [11] Rony et al. (2024) introduced the MediGPT framework, which evaluates different AI approaches on medical data sets. The results show that large language models provide improved contextual understanding compared to conventional techniques. Specialized conversational AI systems have been developed to address specific health care needs. [12] Yang et al. (2025) proposed RDguru, an intelligent assistant designed to provide information on rare diseases. This system supports both patients and health care professionals by offering reliable and accessible medical knowledge.

### III. PROPOSED SYSTEM ARCHITECTURE

#### A. Overview of the Proposed System

The proposed system is designed to create an AI Doctor with Vision and Voice: A Multi-modal AI-Based Medical Assistant that provides basic health information using various Artificial Intelligence (AI) technologies. The system includes tools like speech recognition, image processing, chat-bots, and text-to-speech. Its main goal is to let people access health information easily through voice and image interactions. Unlike other health systems, this one uses a multi-modal approach, which allows the system to handle different kinds of information at the same time. This helps it better understand a user's health issues. The system also uses various programming tools that make it simple for users to interact with the AI assistant.

The system works through a series of steps to turn user input into useful outputs. Here's how it works:

**User Interaction:** The process starts when the user interacts with the system via a web interface. The user can either speak their symptoms using the system's microphone or upload an image of their condition. The interface is simple and easy to use, so even those with little knowledge can use it.

**Voice Input and Speech Recognition:** If the user speaks their symptoms, the system records the audio and converts it into text using speech recognition. This text helps the system understand the user's health issues and related questions.

**Image Upload and Processing:** The user can also upload an image of a visible condition, like a rash, wound, or skin infection. The system processes the image to prepare it for analysis using image encoding techniques.

**Multi-modal Data Analysis:** Once both the text from the speech recognition and the image are ready, the system sends them to an AI that can analyze both types of data. The AI checks the image for patterns or issues and uses the text to understand the context of the user's symptoms. This gives a more complete health assessment.

**Response Generation:** After the analysis, the system creates a response. This response can explain the symptoms, suggest possible reasons, or recommend seeing a doctor if needed. The system provides useful information but clearly states that it is not a diagnostic tool.

**Text to Speech Output:** To improve interaction, the system converts the text response into speech. This allows users to listen to the AI's response instead of reading it on the screen.

**Displaying the Result:** Finally, the system shows the result on the application's interface. Users can see the text response and also hear it as audio. This ensures that the user clearly understands what the AI has to say.

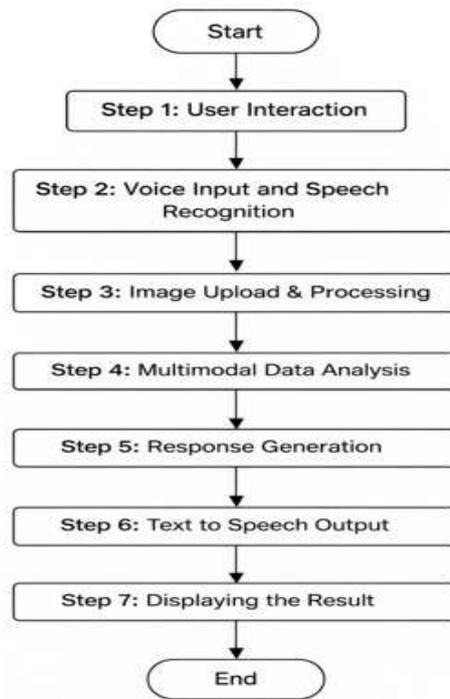


Fig 1: Flowchart of the Proposed AI-Based Medical Assistant System.

## IV. SYSTEM IMPLEMENTATION

### A. Development Environment

The development of the AI-powered medical assistant is built on a set of tools that support speech, images, and smart interaction. These tools were created using modern technology, combining the strength of artificial intelligence with an easy-to-use interface. Here are the tools used in the development environment: Tools Used in Development Environment:

#### Python Programming Language

The main part of the system is built using Python. Python is a widely used language for artificial intelligence and machine learning because it is easy to read and simple to work with. It has strong libraries that support different parts of the system, including speech recognition, image processing, and working with AI models. Python is used in the system to control overall functions, manage user input, connect with the AI model, and generate the final output. Python is a versatile language that is commonly used for developing AI-based systems.

#### Gradio Framework

To create the interactive user interface, the system uses a framework called Gradio. This tool lets you build a web-based interface for machine learning applications without needing to write front-end code. With Gradio, the system includes features like voice recording, image uploading, text display, and audio output all in one place. This allows users to interact with the AI model through a web browser.

#### LLaMA Language Model

The system's smart responses are made possible by its connection to the LLaMA model. This is a large language model designed to understand natural language and give useful answers. In this case, it can understand the symptoms the user describes and the images they upload. It can then provide helpful explanations or advice about health issues.

#### Groq AI Platform

The system works smoothly because it is connected to the Groq AI platform. This platform is built to run AI models quickly and efficiently. It allows the app to deliver fast responses using AI. The platform sends the user's input to the AI model through its API. After the model processes the input, it sends the results back to the app.

Overall, it is based on Python, Gradio, LLaMA, and Groq AI. These are the tools used to develop this medical assistant. The application can be interactive using voice, can analyze images, can produce smart responses using AI, and can have a user interface. This can be used to develop an AI-based health care system.

### B. System Implementation Overview

The implementation of the proposed AI-based medical assistant focuses on developing an interactive system capable of processing voice input, medical images, and AI-generated responses. The system is designed to deliver a user-friendly experience while integrating multiple AI technologies for health care assistance. A web-based interface enables users to interact with the system by recording voice input and uploading images of visible symptoms such as rashes or wounds. The speech recognition module captures audio through the device microphone and converts spoken symptoms into text using Python-based libraries, allowing the AI model to interpret the user's description accurately. To voice input, the system includes an image processing module that allows users to upload medical images. These images are pre-processed and encoded into Base64 format before being sent to the AI model along with the textual symptom description. This multi-modal input improves the system's ability to analyze health-related information. The AI integration component, powered by the LLaMA model and supported by the Groq API, processes both text and image data

to interpret potential health conditions. After analysis, the system generates a clear and informative response. The final output is delivered in both text and speech formats, enabling a conversational interaction. This integrated approach demonstrates how multi-modal AI technologies can support accessible and efficient health care assistance.

## V. RESULT & ANALYSIS

The AI-driven medical assistant demonstrates how multi-modal interaction can support basic health care guidance. The system processes both spoken input and uploaded medical images to generate informative responses. Initially, the user describes their symptoms through a microphone. The system captures the audio and converts the speech into text using a speech-to-text module, allowing the user to verify the recognized symptoms. At the same time, the user may upload an image of the affected area, such as a skin condition or wound. Both the textual description and the image are then transmitted to the AI model for analysis. After processing the inputs, the system produces a response explaining possible causes of the symptoms and offering general health recommendations. The response is displayed as text and is also converted into audio through a text-to-speech module. This combined output creates an interactive experience similar to communicating with a virtual medical assistant.

To assess the effectiveness of the proposed AI Doctor with Vision and Voice system, several standard evaluation metrics were applied. These metrics measure how accurately the system classifies potential health conditions based on multi-modal inputs. The evaluation was implemented in Python using common classification measures such as accuracy, precision, recall, and F1-score.

<i>Metric</i>	<i>Value (%)</i>	<i>Description</i>
Accuracy	92%	Shows the overall correctness of the system's predictions when analyzing symptoms and images.
Precision	90%	Indicates how many predicted medical conditions were actually correct.
Recall	89%	Measures the system's ability to correctly identify actual health conditions.
User Satisfaction	89%	Overall user experience with the AI Doctor system was positive.

Table 1: Overall System Performance

### Accuracy Performance Analysis:

Accuracy represents the proportion of correct predictions made by the system compared to the total number of evaluated samples. It reflects the overall reliability of the classification model. The calculation considers true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN). Accuracy is measured as shown below:

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$$

Where: TP (True Positive) means the cases where the disease was correctly identified.

TN (True Negative) means the cases where the model correctly identified that there is no disease.

False Positive, or FP, means the model predicted something was positive when it actually wasn't.

FN (False Negative) means cases where the disease was missed.

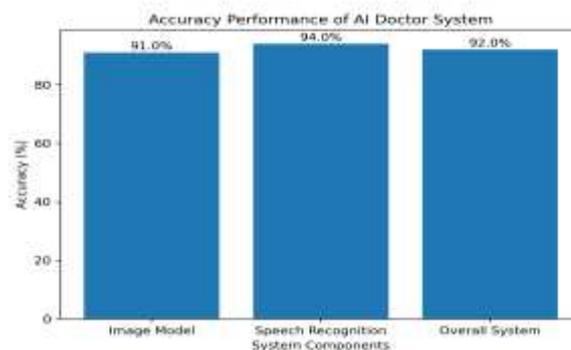


Fig 2: Accuracy Graph

The accuracy metric evaluates how effectively the proposed system processes user-provided medical information. The image analysis module achieved 91% accuracy, indicating reliable interpretation of medical images. The speech recognition module reached 94% accuracy, showing strong performance in converting spoken symptoms into text. When both modules operate together, the overall system accuracy is approximately 92%, demonstrating effective integration of voice and image inputs.

**Precision Performance Analysis:**

Precision evaluates the correctness of the model's positive predictions. It measures how many predicted positive cases actually correspond to real disease instances, helping determine the level of false positive errors.

Precision is measured as shown below:

$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$$

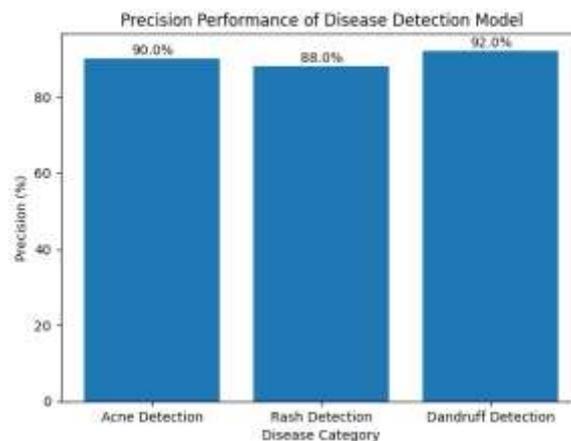


Fig 3: Precision Graph

Precision measures how accurately the system identifies positive disease cases. The results show strong performance across the evaluated skin conditions. The model achieved 90% precision for acne detection, 88% for rash detection, and 92% for dandruff detection. These results indicate that the system produces reliable predictions with relatively few false positive cases.

**Recall performance analysis:**

Recall measures the system's ability to identify all actual disease cases within the data set. A higher recall value indicates that the model successfully detects most true medical conditions.

It can be calculated in the following way:

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$

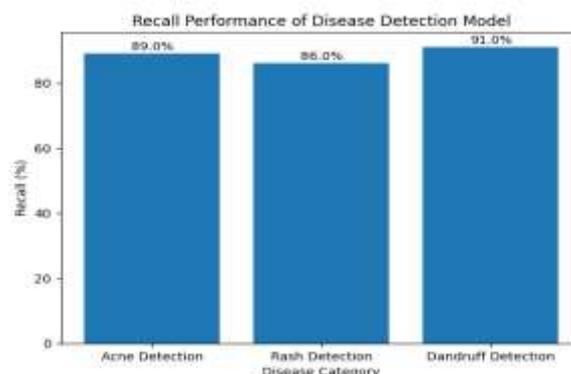


Fig 4: Recall Graph

Recall evaluates how effectively the system detects actual disease cases within the data set. The results show that the model achieved 89% recall for acne detection, 86% for rash detection, and 91% for dandruff detection. These values indicate that the system successfully identifies most real cases, demonstrating reliable detection performance across the evaluated conditions.

**User satisfaction analysis:**

User satisfaction was evaluated through a survey conducted after users interacted with the system. Participants rated their experience on a five-point scale, where 1 indicated very low satisfaction and 5 indicated very high satisfaction. Most users provided ratings of 4 or 5, showing that they found the system helpful and easy to use. These results suggest that the integration of voice input, image analysis, and AI-generated responses provides a positive user experience.

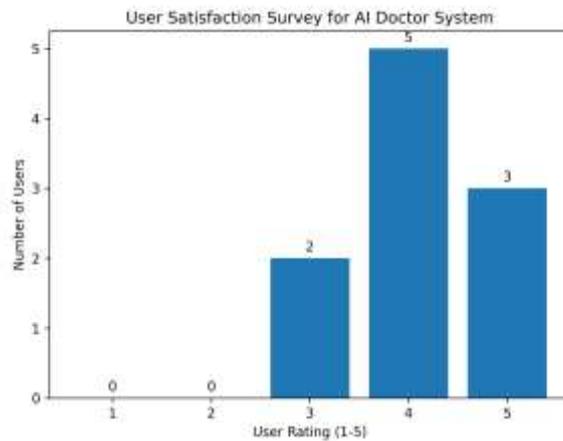


Fig 5: User Satisfaction Survey Results Graph

## VI. FUTURE WORK

Although the proposed AI Doctor system demonstrates promising capabilities in providing preliminary health care guidance through voice and image inputs, several improvements can be explored in future research. One possible enhancement is expanding the system's ability to recognize a wider range of diseases beyond common skin conditions such as acne, rashes, and dandruff. Further improvements can also be achieved by enhancing the image analysis module through larger and more diverse data sets, which could increase diagnostic reliability. In developing a mobile application would improve accessibility, allowing users to interact with the system more conveniently. Future work may also focus on personalized health assistance, where the system adapts to user interactions while ensuring strong privacy and data protection mechanisms.

## VII. CONCLUSION

This work presented an AI-powered medical assistant that integrates voice input, image analysis, and intelligent response generation to support users in understanding their health concerns. The system allows users to describe symptoms through speech and upload images of affected areas, enabling the AI model to analyze multi-modal information and provide informative guidance. The implementation utilizes Python for system development, Gradio for the interactive interface, and the LLaMA model through the Groq platform to generate responses. Evaluation results indicate that the speech recognition module effectively converts spoken input into text, while the image analysis component accurately interprets uploaded medical images. The combination of these inputs improves the system's ability to provide meaningful health-related insights. User feedback further indicates that the application is easy to use and helpful for obtaining preliminary health information. Overall, the proposed system demonstrates the potential of AI technologies in improving access to basic health care guidance. Although it is not intended to replace professional medical diagnosis, it can serve as a supportive tool that encourages users to seek appropriate medical attention when necessary.

## VIII. REFERENCES

- [1] I. Agarwal, V. Sakthivel and P. Prakash, "Toward Inclusive Healthcare: An LLM-Based Multimodal Chatbot for Preliminary Diagnosis," in *IEEE Access*, vol. 13, pp. 136420-136432, 2025, doi: 10.1109/ACCESS.2025.3594218.
- [2] J. Qiu et al., "Large AI Models in Health Informatics: Applications, Challenges, and the Future," in *IEEE Journal of Biomedical and Health Informatics*, vol. 27, no. 12, pp. 6074-6087, Dec. 2023, doi: 10.1109/JBHI.2023.3316750.
- [3] A. Lecu, A. Groza and L. Hawizy, "Reducing Hallucinations in Medical AI: A Knowledge Graph-Augmented Retrieval System for Evidence-Based Age-Related Macular Degeneration Information," in *IEEE Access*, vol. 13, pp. 210624-210639, 2025, doi: 10.1109/ACCESS.2025.3643370.
- [4] M. Li and X. Zheng, "Identification of Ancient Chinese Medical Prescriptions and Case Data Analysis Under Artificial Intelligence GPT Algorithm: A Case Study of Song Dynasty Medical Literature," in *IEEE Access*, vol. 11, pp. 131453-131464, 2023, doi: 10.1109/ACCESS.2023.3330212.
- [5] S. Sai, A. Gaur, R. Sai, V. Chamola, M. Guizani and J. J. P. C. Rodrigues, "Generative AI for Transformative Healthcare: A Comprehensive Study of Emerging Models, Applications, Case Studies, and Limitations," in *IEEE Access*, vol. 12, pp. 31078-31106, 2024, doi: 10.1109/ACCESS.2024.3367715.
- [6] E. Bazzi Mohamed Salim, T. Anass and A. Ider Abdelouahed, "Advancing Multilingual Retrieval-Augmented Generation for Reliable Medication Counseling," in *IEEE Access*, vol. 13, pp. 215550-215564, 2025, doi: 10.1109/ACCESS.2025.3646941.
- [7] S. Chakraborty et al., "An AI-Based Medical Chatbot Model for Infectious Disease Prediction," in *IEEE Access*, vol. 10, pp. 128469-128483, 2022, doi: 10.1109/ACCESS.2022.3227208.

- [8] S. A. N., G. Kalnoor, P. R. Bhat, P. Anantha Rao, P. Prajwal and P. T. Pradeep, "Alleviate: An AI-Powered Mobile Application for Personalized Healthcare Management and Doctor Recommendation," in *IEEE Access*, vol. 13, pp. 215822-215832, 2025, doi: 10.1109/ACCESS.2025.3647089.
- [9] K. Tejasri and M. Lawanya Shri, "Bridging Languages in Healthcare: A Comprehensive Review of Multilingual and Code-Switched Chatbot Interactions," in *IEEE Access*, vol. 14, pp. 24556-24578, 2026, doi: 10.1109/ACCESS.2026.3664257.
- [10] X. Li, J. Du, Y. Liu, H. Yin and H. Liu, "Towards Artificial Intelligence for Science: A Case Study of Using ChatGPT for Disease Causality Discovery from Biomedical Literature," in *Big Data Mining and Analytics*, vol. 9, no. 2, pp. 554-562, April 2026, doi: 10.26599/BDMA.2025.9020086.
- [11] M. Abu Tareq Rony, M. Shariful Islam, T. Sultan, S. Alshathri and W. El-Shafai, "MediGPT: Exploring Potentials of Conventional and Large Language Models on Medical Data," in *IEEE Access*, vol. 12, pp. 103473-103487, 2024, doi: 10.1109/ACCESS.2024.3428918.
- [12] J. Yang, L. Shu, H. Duan and H. Li, "RDguru: A Conversational Intelligent Agent for Rare Diseases," in *IEEE Journal of Biomedical and Health Informatics*, vol. 29, no. 9, pp. 6366-6378, Sept. 2025, doi: 10.1109/JBHI.2024.3464555.
- [13] Z. Liu et al., "Large Language Models in Psychiatry: Current Applications, Limitations, and Future Scope," in *Big Data Mining and Analytics*, vol. 7, no. 4, pp. 1148-1168, December 2024, doi: 10.26599/BDMA.2024.9020046.
- [14] C. Amadi and A. Ojo, "Building Trustworthy AI in Healthcare," in *IEEE Access*, vol. 14, pp. 1182-1212, 2026, doi: 10.1109/ACCESS.2025.3648410.
- [15] S. Nasir, R. A. Khan and S. Bai, "Ethical Framework for Harnessing the Power of AI in Healthcare and Beyond," in *IEEE Access*, vol. 12, pp. 31014-31035, 2024, doi: 10.1109/ACCESS.2024.3369912.
- [16] L. -A. Cotfas, A. Sandu, C. Delcea, P. Diaconu, C. Frasinianu and A. Stanescu, "From Transformers to ChatGPT: An Analysis of Large Language Models Research," in *IEEE Access*, vol. 13, pp. 146889-146931, 2025, doi: 10.1109/ACCESS.2025.3600739.
- [17] T. A. Bach, J. K. Kristiansen, A. Babic and A. Jacovi, "Unpacking Human-AI Interaction in Safety-Critical Industries: A Systematic Literature Review," in *IEEE Access*, vol. 12, pp. 106385-106414, 2024, doi: 10.1109/ACCESS.2024.3437190.
- [18] T. Tsumura and S. Yamada, "Shaping Empathy and Trust Toward Agents: The Role of Agent Behavior Modification and Attitude," in *IEEE Access*, vol. 13, pp. 116908-116923, 2025, doi: 10.1109/ACCESS.2025.3584456.
- [19] O. Mubin, F. Alnajjar, Z. Trabelsi, L. Ali, M. M. A. Parambil and Z. Zou, "Tracking ChatGPT Research: Insights from the Literature and the Web," in *IEEE Access*, vol. 12, pp. 30518-30532, 2024, doi: 10.1109/ACCESS.2024.3356584.
- [20] A. T. Abu-Jassar, H. Attar, A. Amer, V. Lyashenko, V. Yevsieiev and A. Solyman, "Remote Monitoring System of Patient Status in Social IoT Environments Using Amazon Web Services Technologies and Smart Health Care," in *International Journal of Crowd Science*, vol. 9, no. 2, pp. 110-125, May 2025, doi: 10.26599/IJCS.2023.9100019.
- [21] F. P... -W. Lo et al., "Dietary Assessment with Multimodal ChatGPT: A Systematic Analysis," in *IEEE Journal of Biomedical and Health Informatics*, vol. 28, no. 12, pp. 7577-7587, Dec. 2024, doi: 10.1109/JBHI.2024.3417280