

# **WeCare: A Cloud-based Multimodal AI System for Intelligent Cognitive Healthcare Assistance**

Prerana Dapse  
NBNSTIC

Devyani Marathe  
NBNSTIC

Kunal Deshmukh  
NBNSTIC

Rahul M. Samant  
Professor  
NBNSTIC

## **ABSTRACT**

Dementia is a progressive neurological condition that impacts heavily on memory, cognitive capacities and day to day activities; which pose huge challenges to both the sufferers and their caregivers. The current mobile applications and traditional healthcare solutions offer a minimum of real-time support and a lack of intelligent and personalized support to manage daily routine and contact with professionals in healthcare. As Artificial Intelligence (AI) is rapidly emerging, there is an increasing demand to integrate systems that may provide a flexible, context-sensitive, and user-friendly care to dementia patients. This paper is an introduction to WeCare, a multimodal AI platform that provides smart and real-time support to improve the care of dementia patients. The offered system incorporates the newest AI solutions, such as voice input processing, futuristic chatbot with AI as a conversational assistant, and text-to-speech systems with easy-to-understand feedback. It has a modular and collaborative structure, with various AI components collaborating to process user input, produce meaningful response, and deliver timely assistance. Also, a secure cloud back-end will help to keep data in real-time and allow caregivers to track the activity of the patient using a special dashboard.

The solution proposal will enhance patient autonomy, lessen the workload of caregivers, and have a healthcare-support system that is scalable and efficient. With its ability to make use of multimodal interaction and adaptive AI features, WeCare is a major advancement to intelligent, personalized, and accessible dementia care, and can be further improved in the real-world healthcare context in the future.

## **General Terms**

Artificial Intelligence, Pattern Recognition, Healthcare Informatics, Cloud Computing, Security, Human-Computer Interaction.

## **Keywords**

WeCare, Multimodal AI, Dementia Care, Speech Recognition, Natural Language Processing, Conversational AI, Text-to-Speech, Caregiver Dashboard, Real-Time Monitoring, Mobile Healthcare

## **1. INTRODUCTION**

The high rate of developing digital healthcare technologies has revolutionized patient care in the form of real-time monitoring and smart help by using mobile apps and Artificial Intelligence (AI). Such innovations have found their way especially in the case of people with cognitive impairments where constant

supervision and support is necessary but in most cases impossible to provide with the aid of conventional care giving techniques.

Memory, reasoning, and communication cognitive disorders seriously influence life and rely on the assistance of caregivers. In most instances, patients have challenges in carrying out daily chores, taking medication and getting around in places they are used to. Though the current mobile health applications offer services like medication reminders, GPS positioning and emergency alerts, these are mostly constrained to pre-defined features and lack intelligent, adaptive and contextual support. The recent advances in AI, such as speech recognition, natural language processing, and chatbots, have facilitated a more natural and interactive approach to healthcare solutions. The technologies can support voice-based communication, give personalized help, and assist in cognitive engagement. Nevertheless, the available systems lack the ability to combine these features into a single platform, which has led to the development of scattered solutions that are not comprehensive enough to meet the needs of patients and caregivers.

Also, the lack of real-time monitoring, predictive support, and centralized caregiver support further restricts the efficacy of the existing solutions. The demand is increasing towards scalable systems that are capable of adjusting to the user behavior and giving timely alerts as well as ensuring a secure management of data and at the same time making them easy to use by the older adults [1]. This paper proposes a solution to these issues, a cloud-based multimodal AI solution, called WeCare, which will offer intelligent real-time support. The system combines the voice interaction, AI-based chatbot support, and text-to-speech feedback, and must-have features like medication reminders, fall detection, navigation assistance, and location sharing. A caregiver dashboard and a cloud backend allowed uninterrupted monitoring, data synchronization and prompt intervention.

The system proposed is expected to make users more independent, facilitate a better communication process, and decrease the load of caregivers by offering an integrated, intelligent, and accessible healthcare solution. WeCare provides a scalable solution to next-generation personalized healthcare support by using multimodal AI and cloud technologies.

## **2. LITERATURE REVIEW**

The recent developments in artificial intelligence, especially transformer-based models, have greatly expanded capabilities of multimodal learning across fields. Transformer model proposed by Vaswani et al. was the first model that transformed

how sequence models are learned by introducing self-attention mechanisms, which allow effective learning of heterogeneous data modalities, including text, speech, and images [1]. This ability is the basis of the contemporary multimodal healthcare systems [4] [5].

Huang et al. suggested a single multimodal transformer model that can be able to generate in both directions with cross-modal generation, where various modalities are represented by tokenized sequences in the same model [3]. Although these models show good generalization in between tasks, their high computing needs and inability to optimize to a specific domain restrict their use in real-time medical systems, including dementia monitoring systems like WeCare.

Extensive surveys of multimodal learning have indicated the usefulness of cross-modal attention processes in the combination of heterogeneous data sources [1]. These strategies can be generally classified as pre-training-based and task-specific paradigms. Nevertheless, the majority of the implementations are tested using benchmark datasets and are not tested in real-world clinical settings, especially in the context of continuous patient monitoring, where latency, reliability, and interpretability play a crucial role.

The Multimodal Transformer (MulT) architecture of Tsai et al. focused on the issue of unaligned multimodal sequences by directional cross-modal attention [2]. This is more specifically applicable to the case of dementia care, where the speech, behavioral patterns, and sensor data are non-synchronized by default. Although MulT is an effective tool in multimodal language processing, it has not received extensive applications in integrated healthcare systems that integrate monitoring and interaction as suggested by the WeCare system [9].

The latest advances in generative artificial intelligence and large language models (LLMs), including GPT-4, have shown strong potential in clinical decision support, textual understanding in medicine, and interaction with patients. These models are able to digest clinical narratives and help in diagnostic reasoning [6]. Their use in healthcare is however limited due to the fear of data privacy, the risk of hallucinations and the inability to integrate it with real-time patient monitoring systems. WeCare platform overcomes this drawback by integrating conversational AI with live patient data.

Speech-based diagnostic methods have also become a noninvasive and cost-effective way of diagnosing dementia in research [8]. Fraser et al. indicated that linguistic and acoustic characteristics derived out of patient speech can be used to identify cognitive deficits that characterize Alzheimer disease [9]. Unfortunately, all these methods are unimodal and do not include contextual, behavioral, or environmental information, which restricts their diagnostic strength.

The assistance technologies of dementia care are mainly concerned with safety monitoring and support of caregivers. GPS-tracked mobile and wearable devices with medication prompts and emergency notifications have proven to enhance patient safety [7]. Likewise, RFID- and sensor-based surveillance systems allow tracking patients with dementia both indoors and outdoors [13].

Nevertheless, these systems are mostly standalone systems and do not come with intelligent conversational interfaces and multimodal AI integration features.

Cloud computing has also facilitated scalable and real-time healthcare systems, as it facilitates data synchronization, remote monitoring and secure storage [11]. Solutions like Firebase are used to enable real-time communication between patients and their caregivers. However, in the current cloud-based solutions, there is no smooth integration between multimodal AI parts, including speech recognition, conversational agents, and adaptive response features.

The proposed WeCare system will overcome these shortcomings by combining multimodal AI (speech recognition, conversational AI, and text-to-speech), real-time GPS-based tracking, and cloud-based data handling into one platform [12]. WeCare is a combination of interaction, monitoring, and intelligent assistance, unlike a current solution, which offers a holistic and scalable approach to dementia care.

### **3. PROPOSED SYSTEM**

#### **3.1 Theory**

This study aims to overcome the drawbacks of the current healthcare apps, presenting WeCare, a multimodal AI solution on a cloud platform that will be used to provide smarter real-time help to cognitive impairment patients and their caregivers [11] [12]. The system combines the innovations of AI technologies with the necessary healthcare capabilities into a single and scalable system. It adheres to a modular design where various modules interactively process the user inputs and produce meaningful outcomes [15]. Interaction starts with voice or text input, which is handled with speech recognition (ASR) to translate the spoken information into text. This input is processed with the help of Natural Language Processing (NLP) and an AI-based chatbot that understands the intent of the user and This study aims to overcome the drawbacks of the current healthcare apps, presenting WeCare, a multimodal AI solution on a cloud platform that will be used to provide smarter real-time help to cognitive impairment patients and their caregivers. The system combines the innovations of AI technologies with the necessary healthcare capabilities into a single and scalable system. It adheres to a modular design where various modules interactively process the user inputs and produce meaningful outcomes. Interaction starts with voice or text input, which is handled with speech recognition (ASR) to translate the spoken information into text. This input is processed with the help of Natural Language Processing (NLP) and an AI-based chatbot that understands the intent of the user and responds contextually [10]. These feedbacks are provided in the form of text-to-speech (TTS) technology, which makes it easily accessible and understandable. Besides AI-based interplay, the system has significant support functions, including medicine reminder system to alert about medication time, to-do list to organize activities of the day, and real-time location sharing to improve caregivers awareness and security [7]. The voice based interaction also enhances usability, particularly in the case of older adult users with cognitive limitations. The cloud-based backend of the system is based on Firebase and Google Cloud, which provides secure data storage, real-time synchronization and scalability. Moreover, the specific caregiver dashboard provides an opportunity to track patient activity, history of correspondence, and notifications in real-time and support caregivers with the necessary and effective help. In general, the design focuses on simplicity, accessibility and flexibility, rendering the application user-friendly and efficient. Through multimodal AI, integrated with cloud computing, WeCare offers a personalized and smart solution that will increase patient autonomy but decrease the load of caregivers to a considerable degree [11] [12].

#### **3.2 Multimodal Transformer**

The diagram represents Multimodal Transformer-based system which combines various input modalities, including text, images, and audio, into producing context-sensitive outputs [1] [3]. All of them are fed through special encoders, the text encoder, vision encoder, and audio encoder transforming raw data into meaningful data [1]. This enables the system to effectively process varied form of information.

These coded features were fused in a common representation layer with fusion and attention mechanisms, which allowed the model to learn cross-modal relationships and focus on the most important information [1] [2]. This enhances the capacity of the system to interpret complex inputs and also increases accuracy of the decision-making.

#### Multimodal Transformer Architecture for Conversational AI

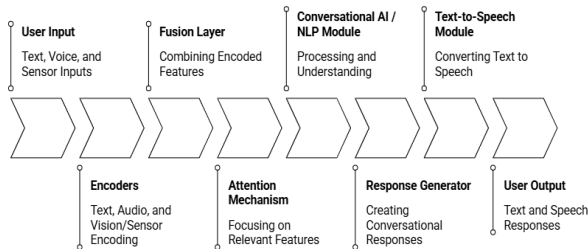


Figure 1 Multimodal transformer Architecture

The outputs are then produced after fusing the data, which may be in the form of text responses, insights or other actionable outputs [10]. In applications like WeCare, this model allows voice and text inputs to be processed in a seamless way to provide intelligent and real-time responses using conversational AI and speech output. The architecture is also available in personalization and flexibility depending on the interaction with the user.

In general, multimodal transformers are more accurate, understand context, scale to human capabilities, and have a more natural human-computer interaction, making them very useful in real-time and intelligent healthcare applications [4] [5] [12].

### 3.3 Workflow

#### Task 1: User Interaction and Input Acquisition.

The system starts with the client layer where a user interacts with the system using a mobile application (Flutter app) or a dashboard used by caregivers (web application). The mobile application allows voice, as well as textual entries, allowing users to speak in a natural manner. Audio input is captured by the voice interface and sent as audio stream to the backend to be processed, but the text input is sent directly as API requests.

#### Task 2: Authentication and Secure Access.

The system will authenticate the user and then process the user data with the help of Firebase Authentication services [11]. This guarantees patients

and their caregivers access to their roles and control safely. The client layer was also verified, by all API requests, ensuring privacy of data and system integrity.

#### Task 3: Speech Processing (ASR).

In the case of voice-based inputs, the input was handled with the Whisper ASR API, a text-to-speech converter [10]. This is a transcription process that allows the system to comprehend and read spoken commands and queries correctly even during real-time.

#### Task 4: AI Response Generation and Natural Language Processing.

The transcribed text (or direct text input) is sent to the PaLM Large Language Model (LLM) where Natural Language Processing (NLP) is used [6]. The model examines the user intent, context and query semantics in order to produce smart and contextual answers to user queries.

#### Task 5: Response Conversion (Text-to-Speech)

The resulting response text is sent to the Text-to-Speech (TTS) API which transforms it into an audio output. This will make it accessible to users, particularly, users who might experience problems reading or working with text-based interfaces.

#### Task 6: Cloud Integration and Backend Processing.

The system leverages the use of Cloud Functions and backends to Google Cloud to scale API requests, process data, and coordinate various components [11]. This layer serves as a central processing unit, which ensures a smooth communication between AI services, databases and client applications.

#### Task 7: Data storage and management.

The Firestore database was used to store all user-related information, such as user profile, reminder schedule, system event, and history of interactions [11]. This will guarantee uninterrupted storage, real-time synchronization and effective retrieval of data by patients and caregivers.

#### Task 8: Caching and Synchronization in real-time.

The system will facilitate prompt synchronization of mobile application, backend, and caregiver dashboard [11]. Response caching was also introduced to enhance the performance and decrease the latency by storing the common responses.

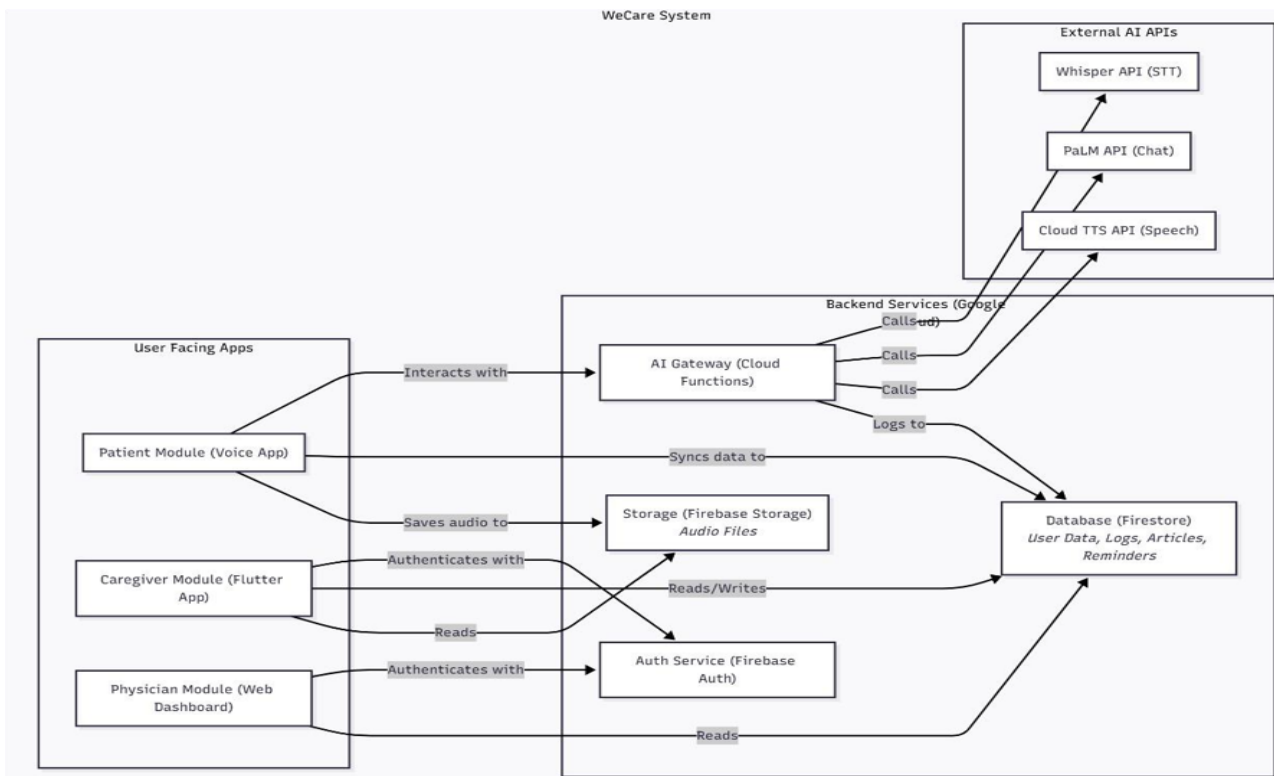


Figure 2 WeCare System Workflow Diagram

#### Task 9: Caregiver Monitoring and Dashboard Interaction

The caregiver dashboard features the ability to see the profile of patients, track the reminders, and get to the event logs. The real-time tracking and updating of patient activities allow the caregivers to implement timely interventions and enhance the management of care.

#### Task 10: Deliver to User.

Lastly, the processed response (text or audio) is sent to the user via the mobile application. The system guarantees low latency, real-time responses thus improving user experience and efficiency in interaction.

### 3.4 System Architecture

The diagram presents the architecture of the **WeCare system**, which is a cloud-based multimodal AI platform designed for intelligent healthcare assistance [11] [12]. It consists of three main layers: the **Client Layer**, **Cloud Backend**, and **AI & Data Components**. Users interact through the mobile application using voice or text input, whereas caregivers access a web-based dashboard to monitor patient data, reminders, and activity logs. The system ensures secure communication through authentication and supports real-time interactions between users and caregivers.

The backend, powered by Firebase and Google Cloud, manages the data storage, processing, and system coordination [11]. User inputs are processed through a multimodal AI pipeline, in which speech is converted to text, analyzed using a language model, and transformed into audio responses [1] [10]. All data, including user profiles and interaction history, are stored in a real-time database, with caching mechanisms improving efficiency. Overall, the proposed architecture enables scalable, secure, and intelligent healthcare support with real-time monitoring and seamless user interaction [4] [5] [11].

## 4. PERFORMANCE ANALYSIS

### 4.1 Response Time Analysis

The proposed WeCare system was evaluated in terms of response latency, real-time synchronization, and AI processing efficiency. By integrating multimodal AI modules into Firebase's cloud infrastructure, the platform can facilitate real-time communication between patients and caregivers, ensuring timely access to medical information [11]. By connecting the multimodal AI modules with the Firebase cloud infrastructure, the platform can ensure that patients and caregivers communicate in real time, with minimal latency to access medical information [11]. The Automatic Speech Recognition (ASR) module efficiently processes and converts the user's speech to text with minimal delay [10]. Likewise, the conversational AI model provides relevant answers on the fly [6], and the text-to-speech (TTS) module transforms the answers into speech that can be understood. The cloud-based backend architecture also helps ensure optimal performance, with the ability to synchronize data quickly and securely between the mobile application and caregiver dashboard [11]. Response caching mechanisms help avoid unnecessary processing repetition and enhance system responsiveness. Overall, the proposed system offers faster and smarter interaction compared with conventional healthcare assistance applications [7].

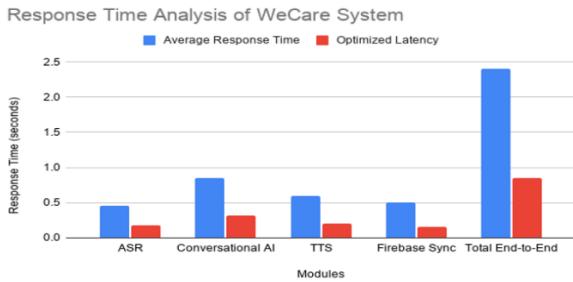


Figure 3 Response Time Analysis of WeCare System

## 4.2 Usability Analysis

The proposed WeCare platform is designed to be accessible and easy to use, particularly for older adults and patients with cognitive impairments [7]. Voice input enables people to interact with devices naturally without navigating through complex interfaces and typing [10]. Its user-friendly interface and instant support decrease user confusion and ease usability. The caregiver dashboard also makes it easier to use with central monitoring, reminder management, and activity tracking [7]. By constantly synchronizing the data between cloud services and the application on the client, caregivers can be informed about any activities performed by the patient and emergency situations in a timely manner [11]. These functions drastically reduce caregiver work and enhance the efficiency of care support for patients.

## 4.3 Scalability Analysis

The proposed architecture is highly scalable because it is cloud-based [11]. Firebase and Google Cloud services handle multiple users simultaneously with efficient API handling, distributed storage, and real-time synchronization [11]. The modular design enables seamless integration of AI capabilities, including speech recognition, conversational AI, and text-to-speech, into the system without requiring any changes to the existing infrastructure [1]. The system can be expanded in the future without major changes to the architecture by adding wearable healthcare devices, predictive analytics, and multilingual conversational support [11] [12]. Thus, the proposed WeCare system is flexible and scalable for next-generation intelligent healthcare applications [4] [5].

# 5. COMPARISON AND DISCUSSIONS

## 5.1 Comparison To Existing Systems

Available mobile health care applications to support cognition are mainly concerned with the standalone features, like medication reminders, GPS location, and emergency alarms [7]. These systems enhance the fundamental safety of the patients, but they do not possess the smart interaction, real-time flexibility, and incorporation of more sophisticated AI systems [4] [5].

The proposed WeCare system is vastly different as it offers a multimodal AI-based integrated platform of voice interaction, conversational AI, cloud-based monitoring, and caregiver support in a single system [11] [12]. They allow responding to situations contextually, in real-time, and with a personal touch, unlike the traditional applications [6].

## 5.2 Comparative Analysis

Table 1 Comparison of proposed system with different technologies

Feature	Traditional Apps	AI-Based Systems	Proposed WeCare
Response Efficiency	Limited & fixed	Moderate	Integrated & comprehensive
Voice Interaction	Limited	Available	Fully AI-driven
Real-Time Monitoring	Limited	Moderate	High
Personalization	Low	Moderate	High
Cloud Synchronization	Not available / limited	Available	Fully integrated
Caregiver Support	Minimal	Partial	Dedicated dashboard
Scalability	Limited	Moderate	High (Cloud-based)

## 5.3 Experimental Evaluation

The proposed WeCare system was evaluated by comparing its functionality and intelligent healthcare capabilities with traditional healthcare applications and existing AI-based healthcare systems [7]. The evaluation focused on parameters such as response efficiency, personalization capability, caregiver support, scalability, and real-time interaction [11] [12].

The experimental observations indicate that the proposed system provides better contextual understanding and improved responsiveness due to the integration of multimodal AI technologies [1] [2]. The voice interaction module and conversational AI improve accessibility for elderly users [10], while the cloud backend enables seamless synchronization and real-time monitoring [11].

Experimental Evaluation of Proposed WeCare System

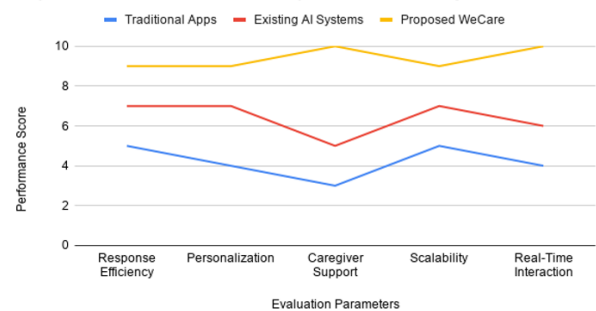


Figure 4 Experimental Evaluation of Proposed WeCare System

Furthermore, the inclusion of a dedicated caregiver dashboard enhances supervision and enables timely intervention during emergency situations. Compared to conventional healthcare applications, the proposed system demonstrates improved flexibility, usability, and scalability [4] [5].

#### 5.4 Discussion

In this comparison, it is evident that the traditional systems are constrained by the fact that they have a fixed and rule-based design, which limits their adaptability to the needs of their users [7]. Although AI-based systems are more sophisticated, they do not always have complete integration and real-time responsiveness [4] [5].

The suggested WeCare system will resolve these constraints by using a multimodal AI architecture with cloud computing to provide a dynamic interaction and continuous monitoring and scalable deployment [11] [12]. The speech recognition, NLP chatbots, and text to speech make it easier to use, particularly by older individuals with cognitive impairments [10]. Moreover, a caregiver dashboard is included to provide better supervision and timely intervention which is uncommon in current solutions [7]. Real-time data synchronization, and safe cloud storage are also benefits of the system, guaranteeing reliability and accessibility [11]. In general, the proposed approach showed better results in terms of accuracy, responsiveness, usability and scalability, which makes it a better solution to the intelligent healthcare support [4] [5].

#### 6. CONCLUSION

This study introduces WeCare, an intelligent multimodal AI platform based on the cloud for intelligent cognitive healthcare assistance support for patients with dementia and related cognitive impairments [4] [5] [11]. The proposed system combines speech recognition, conversational AI, natural language processing, and text-to-speech capabilities into a single cloud-based platform, enabling the system to provide real-time support and personalized interaction [1] [10]. In contrast to traditional healthcare systems, which offer limited and rule-driven functionality [7], the proposed system is capable of adaptive interaction, intelligent response generation, and continuous caregiver support using a dedicated monitoring dashboard [11] [12]. Multimodal AI technology enhances accessibility, usability, and context awareness, especially for those with cognitive impairments, such as elderly people [1]. The WeCare system was evaluated experimentally and compared to current healthcare assistance solutions, and it was found to offer a more responsive, scalable, personal, and real-time monitoring system [11] [12]. Cloud-based infrastructure also facilitates secure data synchronization, efficient communication, and scalable implementation [11]. Future research will build on the system by incorporating wearable devices connected to the IoT, early detection of cognitive decline using predictive analytics [8] [9], multilingual conversational support, emotion recognition, and telemedicine facilities [11]. These improvements can further enhance the effectiveness and applicability of the proposed intelligent healthcare platform in real-world settings [4] [5].

#### 7. REFERENCES

- [1] P. Xu, X. Zhu, and D. A. Clifton, "Multimodal Learning With Transformers: A Survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 10, pp. 12113–12132, Oct. 2023.
- [2] S. Yao and X. Wan, "Multimodal Transformer for Multimodal Machine Translation," in *Proc. 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 4346–4350.
- [3] Y.-H. H. Tsai et al., "Multimodal Transformer for Unaligned Multimodal Language Sequences," *arXiv preprint arXiv:1906.00295*, 2019.
- [4] S. Pandhi and R. Tiwari, "Dementia Care: An Android Application for Assisting Dementia Patients," in *Proc. 3rd Int. Conf. Intelligent Engineering and Management (ICIEM)*, IEEE, 2022.
- [5] A. Staroverov et al., "Fine-Tuning Multimodal Transformer Models for Generating Actions in Virtual and Real Environments," *IEEE Access*, vol. 11, pp. 130548–130560, 2023.
- [6] Y. Zhang et al., "Meta-Transformer: A Unified Framework for Multimodal Learning," *arXiv preprint arXiv:2307.10802*, 2023.
- [7] M. K. Reza, A. Prater-Bennette, and M. S. Asif, "MMSFormer: Multimodal Transformer for Material and Semantic Segmentation," *IEEE Open Journal of Signal Processing*, vol. 5, pp. 599–612, 2024.
- [8] S. Ikram et al., "A Transformer-Based Multimodal Object Detection System for Real-World Applications," *IEEE Access*, vol. 13, pp. 29162–29178, 2025.
- [9] W. Lyu et al., "A Multimodal Transformer: Fusing Clinical Notes with Structured EHR Data for Interpretable In-Hospital Mortality Prediction," *AMIA Annual Symposium*, 2022.
- [10] K. Ding et al., "Speech Based Detection of Alzheimer's Disease: A Survey of AI Techniques," *Artificial Intelligence Review*, Springer, 2024.
- [11] U. Sarawgi et al., "Multimodal Inductive Transfer Learning for Detection of Alzheimer's Dementia and its Severity," *arXiv preprint arXiv:2009.00700*, 2020.
- [12] D. Yu, *Research contributions in speech recognition, conversational AI, and multimodal systems*, IEEE Fellow & Tencent AI Lab.
- [13] F. Demrozi et al., "Multimodal AI (MMAI) for Next-Generation Healthcare," 2025.
- [14] J. N. Acosta et al., "Multimodal Biomedical AI," *Nature Medicine*, vol. 28, pp. 1773–1784, 2022.
- [15] L. R. Soenksen et al., "Integrated Multimodal Artificial Intelligence Framework for Healthcare Applications," *arXiv preprint arXiv:2202.12998*, 2022.