# AI Powered Virtual Assistant for Alerting Home / Hospital-Disabled / Paralytic Patients

Kanishka Raj[*1], N Harini[#2],

Shashank P[#3], Kamaleswari Pandurangan[#4],

*Department of Information Science and Engineering*

*CMR Institute of Technology Bengaluru, Karnataka 560037*
*kamale.csc@gmail.com, kamaleswari.p@cmrit.ac.in*

*Abstract*—**Healthcare assistance for paralyzed and physically challenged individuals requires reliable and continuous monitor- ing, as such patients sometimes lack in expressing the distress they are going through and even a minor inconvenience that their inner self might be feeling. Traditional patient monitoring systems depend largely on manual observation, which may lead to delayed responses in critical situations. This project presents an AI-powered virtual assistant designed to improve communication between patients and caretakers using computer vision, deep learning, and speech processing techniques which not only re- duces the need to rely on conventional systems but also automate the process of monitoring the patients in need. The system captures real-time facial expressions, hand gestures, and voice inputs from patients through a webcam and microphone. Facial emotions are identified using a Convolutional Neural Network (CNN), hand gestures are recognized using MediaPipe-based landmark detection combined with a deep learning classifier, and speech input is converted into text using speech-to-text techniques. The processed information is displayed on a web-based dashboard provided for the caretaker, developed using the Flask framework. Initial testing indicates that the system is capable of functioning in real-time under controlled conditions.**

*Keywords*—*Artificial Intelligence, Assistive Technology, Emotion Detection, Hand Gesture Recognition, Speech-to-Text, Computer Vision, DeepLearning*

## I. Introduction (*Heading 1*)

In recent years, the use of Artificial Intelligence in health-care has increased significantly, especially in applications that support patients with physical disabilities. Individuals suffering from paralysis or severe motor impairments often face difficulties in communicating their needs to caretakers. There are many day to day tasks like conveying about an emergency, expressing distress, random waves of emotions which becomes difficult for people with disability and patients with related problem to express and to show naturally.

Traditional assistive systems used in patient monitoring mainly rely on manual observation or very basic alert mech-anisms. Systems like that requires constant attention from caretakers and may fail to detect subtle emotional or behavioral changes in patients. In many cases, the patient may not be able to clearly express pain, anxiety, or emergency situations, leading to delayed responses. This highlights the need for an intelligent system that can continuously monitor patients and reduce the dependency on the systems which are basic and requires manual proceedings.

Advancements in computer vision, deep learning, and speech processing provide an opportunity to address these challenges. Technologies such as facial expression recognition, hand gesture recognition, and speech-to-text conversion can be combined to create a supportive communication environ- ment for physically challenged patients. By observing facial expressions, hand movements, and voice inputs, an AI-based system can interpret patient distress and convey meaningful information to caretakers in real time.

In this project, an AI-powered virtual assistant is developed to support paralyzed patients by enabling multiple modes of interaction. The system uses a webcam and microphone to capture real-time inputs from the patient. Facial expressions are analyzed to identify emotional states, hand gestures are recognized to detect predefined commands, and speech input is converted into text for communication. These inputs are processed using machine learning and deep learning models and displayed on a caretaker dashboard for monitoring and response.

This project aims to demonstrate how the integration of computer vision, speech processing, and web-based technolo-gies can improve patient–caretaker interaction. This system does not only provide a structured mechanism to monitor the patients with disabilities but also reduces the dependency on the systems that are basic and requires manual support of the users.Ease of Use

## II. RELATED WORK

The advancements and development in the field of artificial intelligence has resulted in development of assistive health-care systems improving the monitoring of the patients and enhanced communication between caretaker and the patient. This section reviews existing research related to facial emo-tion recognition, hand gesture-based interaction, speech-to- text systems, and integrated assistive platforms relevant to the proposed system.

### A. Emotion Recognition Using Computer Vision

Facial expression recognition has been widely studied as a procedure for understanding human emotional states through visual hints. Convolutional Neural Networks (CNNs) have shown an immense performance in extracting discriminative facial features such as eye movement, mouth shape structure,

and facial muscle variations. Goodfellow et al. [1] provided a benchmark dataset and evaluation framework for face expression recognition, highlighting challenges such as lighting variation and partial blockage. In healthcare related studies, expression recognition has been explored as a supportive tool for monitoring patient mental health and stress levels. These studies shows that CNN based approaches are successful for real-time emotion detection but require robust face detection and preprocessing for efficient and true results.

(Based on: Goodfellow et al., 2013; Li and Deng, 2020)

B. Face Detection Techniques

Accurate face detection is a critical preprocessing step for expression recognition systems. Viola and Jones [2] proposed a quick object detection framework using Haar-like features and a cascade classifier, which has been pro actively used for its speed and computation power. Many real-time applications continue to embrace this approach, especially in webcam-based systems, where efficiency is given priority over complex deep learning detectors. Even though modern deep learning-based detectors exist, cascade-based methods remain most widely used and practical for lightweight healthcare monitoring systems.

(Based on: Viola and Jones, 2001)

C. Hand Gesture Recognition Using Landmark-Based Meth- ods

Hand gesture recognition provides an efficient method for the patients or people who cannot communicate and express themselves. Recent research has drifted from image-based gesture recognition to landmark-based approaches. Zhang et al. [3] introduced MediaPipe Hands, a real-time hand tracking framework which has the ability of detecting 21 hand landmarks using a single RGB(Red,Green,Blue) camera. Foregoing studies have used these landmarks as input to machine learning and deep learning classifiers to acknowledge predefined gestures. Landmark-based methods provides better stability, reduced complexity in computations, and suitability for real-time applications.

(Based on: Zhang et al., 2020; Biswas et al., 2023)

D. Speech-to-Text Systems in Assistive Applications

Speech-to-text technology has been widely used in assistance and healthcare systems to provide communication and documentation. Zolnoori et al. [4] evaluated automatic speech recognition systems in patient and nurse communication and reported that speech-to-text tools can notably reduce manual effort, although performance might show variations under conditions such as noisy environment. Other studies have exhibited the use cases of speech-to-text interfaces in assistance providing devices for people with limited movement. These results support the addition of speech-based interaction as a complementary communication mode in assistive system and provide a support system that is robust enough.

(Based on: Zolnoori et al., 2024; Rodriguez et al., 2025)

E. Integrated AI-Based Assistive Healthcare Systems

Several researchers have proposed the systems that is combination of computer vision, deep learning models and automated architectures for providing support to healthcare systems and reduction to manual efforts. These systems aim to improve the process of monitoring the patients and caregiver response by combining multiple AI outputs into a single interface. However, many existing systems focus on a single interaction methodology or require setting up of complex hardware solutions. There remains a need for lightweight and multimodal systems that combine expression recognition, gesture-based commands, and speech-to-text conversion in a single platform suitable for real-time healthcare assistance.

(Based on: Akhtar et al., 2024; Dhobale et al., 2025)

## III. PROPOSED SYSTEM

The proposed system is an AI-powered assistive virtual assistant designed to help paralyzed and physically challenged patients communicate effectively with caretakers. The sys- tem integrates computer vision, deep learning, and speech processing techniques to monitor patient emotions, recognize hand gestures, and convert voice input into text. The primary objective of the system is to provide a practical, real-time communication mechanism that reduces manual supervision and supports timely caretaker response.

A. System Architecture

The system architecture is designed using an approach that is based on modules to support real-time processing and the simplicity of integration . The major components of the system are explained as follows:

• Input Module: Captures real-time video and audio input from the patient using a webcam and microphone. Video input is used for analyzing the gestures and emotions of the users, while the input from the audio is used for conversion speech-to-text .

• Face Detection Module: The Haar Cascade algorithm is applied to detect the facial portion from the frames of the video. This step helps in the elimination background noise and improves the accuracy of detection of emotions.

• **Emotion Recognition Module:** Uses a Convolutional Neural Network (CNN) for the analysis of facial expres- sions and for the classification of the patient's emotional state into predefined categories.

• **Hand Gesture Recognition Module:** Extracts 21 hand landmarks using MediaPipe and processes them using a deep learning classifier to recognize hand gestures that are already defined within the system.

• **Speech Processing Module:** this module is used for converting patient voice input into text using speech-to- text techniques, enabling verbal communication without manually providing the inputs.

• **Application Server:** This module is implemented using the Flask framework to manage routing, user authen- tication, role-based access control, and communication between AI modules and the user interface.

• **Database Module:** This module forms the base of the system where all the information regarding the system is stored.Uses SQLite to store patient details, caretaker information, detected emotions, gestures, voice messages, and task schedules with timestamps.

International Journal of Advanced Multidisciplinary Research and Educational Development
Volume 2, Issue 1 | January – February 2026 | www.ijamred.com

ISSN: **3107-6513**

• **Caretaker Dashboard:** Displays real-time alerts to caretakers, patient messages, emotional status, and task up- dates, allowing caretakers to respond immediately and effectively.

Figure 1 is the visual representation of layered architecture of the proposed system.Patient inputs are processed with the help of parallel emotion recognition and gesture classification modules,which are put together with speech processing and managed by a centralized backend.

*B. Technology Stack*

The system is developed using the technologies as given below:

• **Frontend:** HTML, CSS, and JavaScript for building responsive user dashboard which provides all the necessity information for patients and caretakers

• **Backend:** Python with Flask framework for routing, authentication, and integration of the enitre system

• **Computer Vision:** OpenCV for video processing in real time and Haar Cascade for face detection

• **Emotion Detection:** Convolutional Neural Network (CNN) enforced using deep learning libraries

• **Gesture Recognition:** MediaPipe Hand Landmark Detection with a deep learning classifier

• **Speech Processing:** Speech-to-text approach for voice-based communication

• **Database:** SQLite for storing patient data, messages, emotions, gestures, etc everything regarding the patient

*C. Algorithmic Approach*

The algorithmic workflow of the system consists of the following steps as mentioned below:

1) **Data Acquisition:** Patients real time gestures and emotions are captured with the help of a webcam and audio detector.

2) **Face Detection:** There is an application of Haar Cascade algorithm in order to process each video frame to capture the region of the face for emotion analysis.

3) **Emotion Classification:** The detected face is first processed and passed to a CNN model to classify the patient's emotional state.

4) **Hand Landmark Detection:** There is use of MediaPipe to extract 21 hand landmarks from video frames.

5) **Gesture Classification:** A deep learning model is used to analyze landmark coordinates to identify predefined hand gestures.

6) **Temporal Validation:** Gesture predictions are validated across multiple frames to ensure stability and accuracy.

7) **Speech-to-Text Conversion:** Input from audio is processed for the conversion of voice to text format.

8) **Data Storage and Display:** All detected outputs are stored in the database and displayed on the caretaker dashboard in real time.

## IV. METHODOLOGY

The methodology provided below covers the step by step flow of the entire system. It includes all the aspects of the system which includes collection, processing and integration of the modules. It also covers the design considerations which is needed to be kept in mind for the system to work effectively.

A. System Workflow

The overall workflow of the system ensures structured procedure to be followed to provide real time assistance to the patient and ensure communication between the caretaker and the patient. The major steps involved are as follows:

1) Input Capture: The system captures real-time video and audio input from the patient using a webcam and microphone which forms the input base.

2) Face Detection: Haar Cascade algorithm captures each of the video frames for the detection of facial emotions followed by its analysis.

3) Emotion Recognition: The detected face is prepro- cessed and then provided to a Convolutional Neural Network (CNN) to classify the patient's emotional state.

4) Hand Landmark Detection: MediaPipe is used to extract hand landmarks from video frames for gesture analysis.

5) Gesture Classification: A deep learning model clas- sifies landmarks that is extracted into predefined hand gestures.

6) Temporal Validation: Gesture predictions are validated across multiple consecutive frames to ensure stability.

7) Speech-to-Text Conversion: Audio input is then con- verted into texts using the speech to text libraries.

8) Data Storage and Display: Emotions , gestures , and voice that is detected is them stored in suitable format in the database and is displayed in the user dashboard.

International Journal of Advanced Multidisciplinary Research and Educational Development
Volume 2, Issue 1 | January – February 2026 | www.ijamred.com

ISSN: **3107-6513**

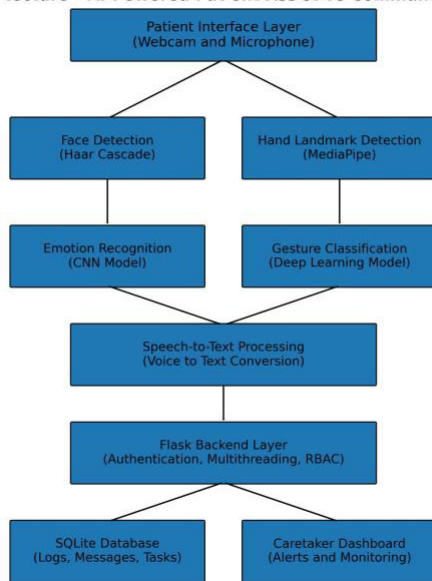Fig. 1: System Architecture of the AI powered virtual assistant



FIG. 1: SYSTEM ARCHITECTURE OF THE AI POWERED VIRTUAL ASSISTANT

B. Implementation Details

The system is implemented using Python as the primary programming language. OpenCV is used for real-time video processing and detection of the facial emotions. A CNN-based model is put to work for facial emotion recognition, while MediaPipe Hand Landmark Detection provides hand tracking with efficient accuracy. Gesture classification is performed using a deep learning model trained on hand landmark co-ordinates.

Speech-to-text functionality is integrated to support voice-based communication. The Flask web framework is used to manage the operations in the backend, which includes routing, authentication, and dashboard rendering. Multithreading is implemented to allow AI processing to run in parallel with web application services. An SQLite database is used to store patient data, detected events, and task-related information.

C. System Design Considerations

During the development of the system there were several design considerations that were considered for the efficient working of it and these were as follows:

• Real-Time Performance: The system is designed to capture and process video and audio inputs with minimal delay to support on time caretaker response.

• Modularity: Each functional component is developed independently as modules, allowing easier maintenance and future enhancements.

• Reliability: Temporal smoothing and confidence thresh- olds are used to reduce false gesture detection.

• Usability: To make the system flexible and easier to understand for caretakers, the user interface is kept simple and more informative.

• Scalability: The architecture is kept flexible in order to allow additional AI modules to be integrated in future without the need of major system redesign.

## V. RESULTS AND DISCUSSION

This section discusses the results obtained after performing preliminary tests of the proposed system.Since the system is still under development, these observations are made based on the controlled environment and test cases and does not portray the final performance the system might show in future after further development.

*A. Experimental Setup*

The system was tested using a standard laptop which has webcam and microphone to generate real time patient mon-itoring environment.The experiments were performed under controlled environment with stable indoor lighting to ensure stable capturing of the video.Test cases were simulated patient gestures, expressions of the face, and voice commands.The dashboard of the caretaker was evaluated through a web browser on the same local network to observe real time updates

B. Performance Evaluation

1) Emotion Recognition Performance: The emotion recog-nition module was tested by observing its ability to classify basic facial expressions such as Happy, Sad, Angry, and Neu-tral. During testing, the system was able to detect and update emotional states in real time when facial expressions were clearly visible. Performance was observed to be consistent under stable lighting conditions, while it showed variations under the conditions when lighting was not appropriate and was dull as it lead to difficulty in capturing the patient . These observations highlight the importance of a consistent environment for vision-based emotion analysis.

2) Hand Gesture Recognition Performance: Hand gesture recognition was tested using some of the predefined ges- tures which were associated with specific commands. The MediaPipe-based hand landmark detection proficiently tracked hand movements in real time. The use of temporal validation across multiple frames helped in reducing false detections caused by small or unintentional hand movements. Gestures performed slowly and clearly resulted in more reliable recog-nition compared to rapid or partially obstructive hand move-ments.

3) Speech-to-Text Performance: The speech-to-text module was tested using short commands and messages such as "emergency". The system was able to convert voice input into readable text and display it on the caretaker dashboard with very less delay. Performance was of satisfaction level in quiet environments, while background noise lead to casual transcription errors. Despite this limitation, the module effec-tively served as an additional communication medium.

4) System Responsiveness: The overall system demon- strated a satisfactory real-time responsiveness due to the use of multithreading. Video processing modules ran concurrently with backend services, preventing interface lag. The caretaker dashboard reflected updates within a short time interval, en-abling near real-time monitoring of patient activity.

C. Comparative Analysis

Compared to traditional patient monitoring systems that rely heavily on observations that is done manually, the proposed system provides automation through multiple AI-based interaction modes. While advanced commercial systems may offer higher accuracy, the proposed solution provides a lightweight and practical alternative suitable for academic and prototype-level healthcare environments which is open sourced and is available to be used by anyone.Table I shows the comparison between the two based on some specified features.

D. Case Study Examples

1) Gesture-Based Emergency Request: In this emergency response situation scenario, the predefined hand signal was performed by the patient. The system successfully detected the signal and alerted the caregiver on the dashboard. The result proved that when verbal communication is not possible, gesture communication can indeed work effectively.

2) Emotion Monitoring Scenario: During the prolonged surveillance process, changes in facial expression patterns were observed and noted on the caretaker dashboard. This helped care takers realize the discomfort in emotion even when there was no specific expression or voice instruction indicated. Such implications point towards the usage of emotion recog- nition in facilitating this surveillance process.

E. Limitations and Challenges

Despite great initial observations, there were several limitations that the system faced:

• The performance of emotion and gesture recognition showed varied performance based on the positioning of the webcam.

• Speech-to-text accuracy decreases in environments where there is noise hindering the system.

• The system currently is less adaptable as it relies on predefined gestures and commands associated.

• Real world testing of the system extensively has not been done yet.

These challenges indicate areas for further refinement and optimization in future development stages of the proposed system.

REFERENCES

[1] R. Malviya and S. Rajput, "Empowering disabled people with help of AI," in *Advances and Knowledge into AI-Created Disability Supports*, Singapore: Springer Nature Singapore, pp. 43–60, 2025.

[2] V. Kumar, S. Barik, S. Aggarwal, D. Kumar, and V. Raj, "The use of artificial intelligence for people with disability: a bright and promising future ahead," *Disability and Rehabilitation: Assistance Technology*, vol. 19, no. 6, pp. 2415–2417, 2024.

[3] H. Madaan and S. Gupta, "AI improving the lives of physically disabled people," in *Proc. Int. Conf. on Soft Computation and Pattern Recognition*, Cham, Switzerland: Springer, pp. 103–112, Dec. 2020.

[4] R. Nacheva and M. Czaplewski, "Artificial intelligence in helping people with disabilities: opportunities and challenges," *HR and Technologies*, vol. 1, pp. 102–124, 2024.

[5] A. E. A. Ali, M. Mashhour, A. S. Salama, R. Shoitan, and H. Shaban, "Development of an intelligent personal assistant system based on IoT for helping people with disabilities," *Sustainability*, vol. 15, no. 6, p. 5166, 2023.

[6] R. Daher, "Integrating AI literacy into teacher education: a very critical perspective paper," *Discover Artificial Intelligence and its benefits*, vol. 5, no. 1, p. 217, 2025.

[7] M. Maashi, N. Aljohani, Y. A. Alsahafi, and M. Rizwanullah, "Assistive communication system using deep sparse autoencoder with feature learning to assist people with hearing disabilities," *Scientific Reports and Research*, 2025.

[8] S. Ferebee, "AI and accessibility: breaking hurdles for people with disabilities," *Premier Journal of Artificial Intelligence*, vol. 2, p. 100012, 2025.

[9] K. Selvi, N. Deepak, J. E. Infanto, M. Gouthamraj, G. Jaivishwak, and R. S. Kumar, "AI powered virtual assistant for enhanced accessibility in the visually impaired community of people," in *Proc. 4th Int. Conf. on Innovative Mechanisms for Industry Applications (ICIMIA)*, IEEE, pp. 933–937, 2025.

[10] B. R. Bhavana, K. Ankolekar, and B. H. Usha, "Artificial intelligence for accessibility: a comprehensive systematic review and impact framework for assistive technologies," 2024.

[11] S. K. Dwivedi, R. Amin, and A. Ghosh, "Artificial intelligence-based assistive technologies for healthcare applications," *IEEE Accessed*, vol. 9, pp. 123456–123468, 2021