# SIGN SPEAK – WHERE SILENCE FINDS A VOICE

## Mrs Vedhashree M R[1], Krishi H S[2], Moulya M[3], Pavan G N[4], Siddharth D Nair[5]

Asst. Prof, Department of ISE, EWIT, Bengaluru, India[1]

Students, Department of ISE, EWIT, Bengaluru, India[2-5]

**Abstract**: Sign Speak is an AI-driven system designed to enable real-time, bidirectional communication between hearing individuals and those with speech or hearing impairments. It translates speech into animated sign language and recognizes hand gestures to generate spoken output in multiple languages, including English, Kannada, Tamil, and Hindi. The system integrates Speech Recognition, Natural Language Processing, and Computer Vision using OpenCV and MediaPipe. It leverages Google Translate for multilingual support and gTTS for voice synthesis. Built on a Flask backend with a responsive HTML, CSS, and JavaScript frontend, Sign Speak performs reliably under varied conditions.Designed for scalability, the system allows easy integration of updates like dynamic gesture recognition and regional sign language support. Testing has shown high accuracy and seamless module coordination. Future enhancements include mobile and wearable versions, continuous gesture recognition, emotion detection, and AR/VR integration—advancing its mission of inclusive, accessible communication.

**Keywords:** Sign Language, Speech Translation, Computer Vision,, Accessibility, Multilingual Translation **S**peech Recognition,  and Natural Language Processing.

## I.    INTRODUCTION

Sign Speak is an AI-powered communication system that bridges the gap between hearing individuals and those with speech or hearing impairments. It offers real-time, bidirectional translation by converting speech into animated sign language and recognizing hand gestures to generate spoken output. The system supports multiple languages, including English, Kannada, Tamil, and Hindi, making it accessible to users across linguistically diverse regions.

Powered by Speech Recognition, Natural Language Processing (NLP), and Computer Vision, Sign Speak detects voice input, interprets gestures, and forms grammatically coherent sentences. The backend integrates MediaPipe for hand tracking, OpenCV for video feed processing, and TensorFlow/Keras for gesture classification. NLP models enable contextual understanding, while Google Translate and gTTS provide multilingual text and audio output. Some of the System Techniques used in Sign Speak's architecture consists of three key modules: (1) Speech/Text to Sign Animation, which extracts keywords from voice input and plays matching sign animations; (2) Gesture Recognition, which identifies hand signs for common words using computer vision; and (3) Sign Letter Recognition, where alphabet gestures are combined into words and sentences. This real-time processing ensures accurate, accessible communication and empowers users with independence and inclusion in education, healthcare, and daily life.

## II.    LITERATURE SURVEY

Several researchers have explored the integration of computer vision, deep learning, and speech processing for sign language translation systems. The following studies provide the foundation for the development of the Sign Speak system:

**[1]    J. Thomas and R. White (2021)**
*Title: Real-Time American Sign Language Recognition Using MediaPipe and Deep Learning Models*
This study demonstrated a real-time American Sign Language (ASL) recognition system using MediaPipe's hand tracking combined with Convolutional Neural Networks (CNNs). The system achieved high accuracy in recognizing static gestures with minimal processing delay, making it suitable for real-time applications.

**[2]    P. Mehta and A. Ramesh (2020)**
*Title: Vision-Based Sign Language Translation Using OpenCV and Deep Neural Networks*
The authors introduced a vision-based translation tool using OpenCV for gesture preprocessing and a deep neural network for classification. Their system translated gestures into English text and laid the groundwork for modular sign-to-speech systems with potential for scalability.

**[3]    S. Banerjee and K. Varma (2019)**
*Title: Speech-to-Text and Text-to-Speech Interfaces for Multilingual Systems*

This research focused on using Google APIs to create a speech-to-text and text-to-speech framework. The system translated spoken inputs into regional languages and generated natural audio outputs using gTTS, significantly enhancing accessibility in multilingual contexts.

**[4]     M. Ali and H. Yadav (2022)**
*Title: Multimodal Translation of Sign Language into Speech Using AI-Based Gesture Recognition*
This paper presented a multimodal system combining gesture recognition with real-time NLP processing and speech synthesis. It utilized TensorFlow for gesture classification and integrated gTTS for speech output, emphasizing fluid sentence formation and contextual understanding.

**[5]     R. Shah and V. Krishnan (2021)**
*Title: Real-Time Gesture-to-Text Conversion with Indian Sign Language Dataset*
The authors built a gesture-to-text system using MediaPipe and TensorFlow, trained on Indian Sign Language (ISL) datasets. The study emphasized the importance of region-specific training data and demonstrated strong results in real-time gesture interpretation for Indian users.

## III.     SYSTEM ARCHITECTURE

The **Sign Speak** system is built on a modular architecture that integrates artificial intelligence, computer vision, natural language processing, and multilingual translation into a seamless framework for real-time bidirectional communication. Its design is divided into three core functional modules: Speech/Text to Sign Animation, Gesture Recognition for Common Words, and Sign Letter Recognition and Sentence Formation. Together, these components ensure a fluid and context-aware translation process that supports both visual and auditory interactions.

The first module, Speech/Text to Sign Animation, processes spoken or typed input and converts it into corresponding animated sign language. It begins with the user's speech being captured and transcribed into text using Google's Speech Recognition API. Natural Language Processing (NLP) techniques are then applied to extract essential keywords from the transcribed sentence. Each keyword is mapped to a specific animation or GIF representing a sign language gesture, which is sequentially displayed on the screen. This module is particularly valuable in educational and service environments where spoken content needs to be conveyed visually.

The second module, Gesture Recognition for Common Words, captures real-time hand gestures through a webcam and classifies them using computer vision and deep learning. MediaPipe is used for high-precision hand landmark detection, while the frames are processed via OpenCV. The detected gestures are passed through a pre-trained Convolutional Neural Network (CNN) model that recognizes signs such as "hello," "yes," or "thank you." Once identified, the corresponding text is generated and converted into speech using gTTS (Google Text-to-Speech), enabling immediate voice-based responses. This real-time gesture-to-audio feedback loop plays a crucial role in bridging day-to- day communication gaps for users who rely on signing.

The third module, Sign Letter Recognition and Sentence Formation, enables users to communicate by forming words and sentences through the signing of individual alphabets. Using the webcam, the system captures the hand sign and feeds it into a gesture classification model trained on alphabet datasets. Recognized letters are assembled into words and further processed using NLP models to predict meaningful, grammatically correct sentences. This output is then translated into the user's preferred language—Kannada, Tamil, Hindi, or English—using Google Translate API, and delivered both as on-screen text and audio output through gTTS. These three modules are connected via a lightweight Flask backend that manages data routing and integrates with a responsive frontend built using HTML, CSS, and JavaScript. The unified system offers users the flexibility to switch between input modes—speech or sign—and receive output in multiple formats including text, animation, and audio. *(See Figure  for the overall system architecture.)*
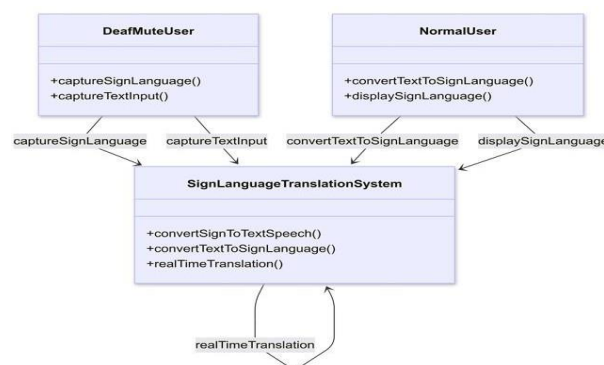


Fig 1: System Architecture

By leveraging a modular approach, the Sign Speak system ensures easy scalability, allowing each component to be updated independently. Its reliance on powerful libraries such as MediaPipe, TensorFlow, and OpenCV enables robust performance in real-world environments, even under varied lighting or background conditions. The integration of multilingual and multimodal translation tools positions Sign Speak as a versatile communication platform, enhancing accessibility and inclusivity across education, healthcare, government, and public interaction domains.
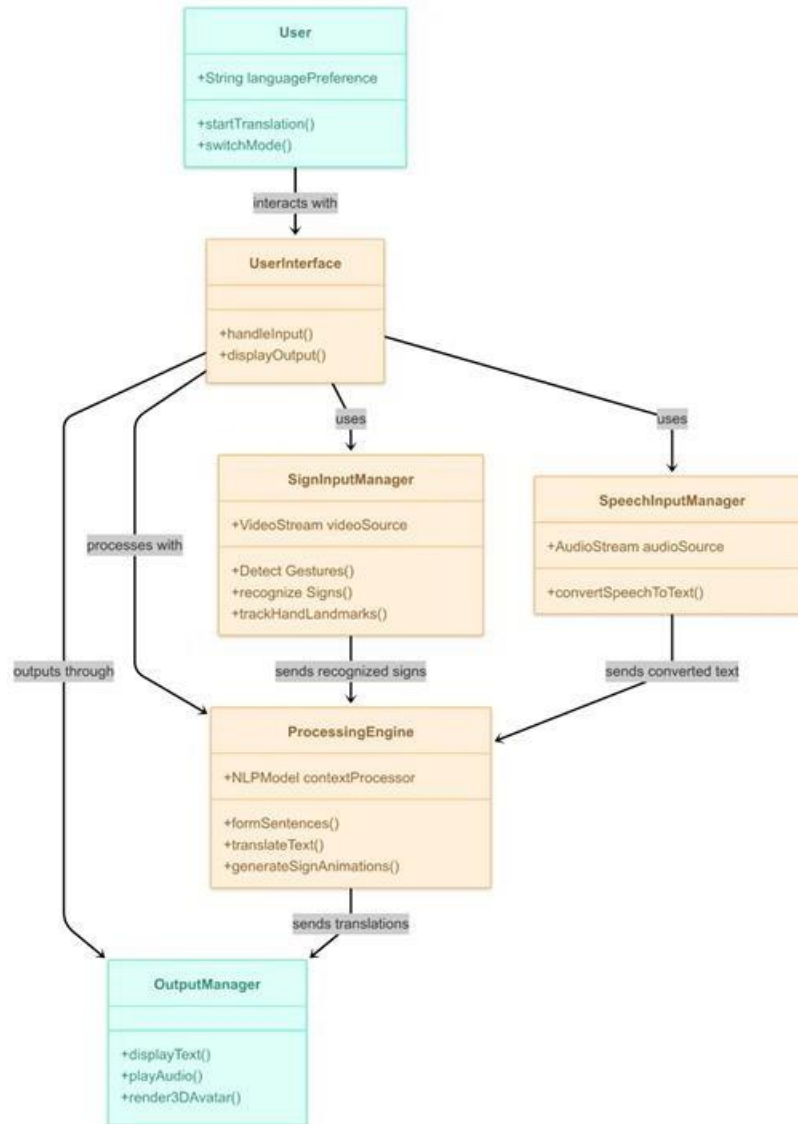


Fig 2: Class Diagram

1.    **Google Speech Recognition API**

This API is used to transcribe spoken words into text in real time. It plays a foundational role in the *speech-to-sign* translation pathway of Sign Speak. When a user speaks into the system, this API captures the audio input, accurately converts it into text, and passes it on to the natural language processing module for keyword extraction and animation mapping. Its real-time capability ensures minimal delay in communication.

2.    **Google Translate API**

The Translate API provides support for converting the system's output into multiple regional languages such as Kannada, Tamil, Hindi, and English. Whether the input originates from spoken words or hand gestures, the final recognized sentence is translated into the user's selected language. This feature significantly enhances the system's accessibility and inclusivity, especially in a linguistically diverse country like India.

3. **Google Text-to-Speech (gTTS)**

gTTS transforms translated text into natural-sounding audio. In the *sign-to-speech* communication flow, once gestures are recognized and converted into meaningful text, gTTS generates a voice output, allowing hearing users to understand what the hearing-impaired user is communicating. This voice feedback loop is especially useful in face-to-face scenarios or customer service interactions.

4. **HTML (HyperText Markup Language)**

HTML structures the web interface of the application. It defines the layout and positions of buttons, input boxes, result displays, video feeds, and audio players. It acts as the visual skeleton of the system, ensuring that users can interact with all functionalities in an organized and coherent way.

5. **CSS (Cascading Style Sheets)**

CSS enhances the presentation of HTML elements by adding styling and responsiveness. It is responsible for the visual design of the Sign Speak interface—managing aspects like color schemes, font styles, spacing, animations, and mobile responsiveness. This ensures the interface is not only user-friendly but also visually accessible to people with varying abilities.

6. **JavaScript**

JavaScript enables dynamic interactivity in the frontend. It listens for user actions (e.g., clicking a button or selecting a language), updates output content without reloading the page, and handles asynchronous communication with the Flask backend using AJAX. It also controls the display of animated signs and plays generated audio in response to backend processing.

7. **Flask**

Flask is the lightweight Python web framework that ties together the backend of the system. It handles routing logic, receives user input from the frontend, invokes the appropriate processing modules (e.g., gesture recognition, NLP, translation), and returns the results to the user interface. Flask's minimalistic and modular design makes it well-suited for real-time, AI-integrated web applications like Sign Speak.

8. **MediaPipe and OpenCV**

These two libraries work together to enable accurate gesture recognition. MediaPipe provides high-precision, real-time hand tracking by identifying 21 hand landmarks per frame, which is essential for recognizing both static and dynamic gestures. OpenCV captures the live video stream from the webcam, preprocesses each frame (e.g., resizing, filtering), and feeds it into the classification pipeline. Their combination allows robust and accurate sign detection even under varying lighting and background conditions.

9. **TensorFlow and Keras**

TensorFlow and its high-level API, Keras, are used to build and train Convolutional Neural Network (CNN) models that classify hand gestures. These models can recognize individual letters and commonly used signs with high accuracy. They form the core of the gesture recognition engine in the system, translating visual input into meaningful digital representations that the system can process further.

## IV. IMPLEMENTATION AND RESULTS

The *Sign Speak* system has been implemented as a modular, real-time communication platform that enables two-way interaction between hearing individuals and those with speech or hearing impairments. The development process involved integrating multiple technologies, including computer vision, machine learning, natural language processing (NLP), and web development frameworks, to ensure seamless translation between spoken language and sign language in both directions. System Overview the architecture of *Sign Speak* comprises three primary functional modules: **(1) Speech/Text to Sign Animation**, **(2) Gesture Recognition for Common Words**, and **(3) Sign Letter Recognition and Sentence Formation**. These modules are interconnected through a Flask-based backend that processes inputs, triggers the appropriate AI model, and returns the output to a responsive web-based frontend.
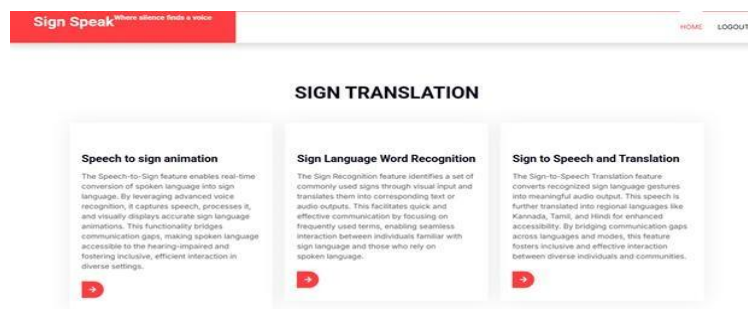


Fig 3: Home Page

### 1. Speech/Text to Sign Animation Module

In this module, the system captures the user's voice input through a microphone. The audio stream is processed using the **Google Speech Recognition API**, which converts it into textual form in real-time. This transcription is then analyzed using **Natural Language Processing (NLP)** techniques to extract the most significant keywords from the sentence.The keywords are then mapped to pre-recorded or animated sign language videos or GIFs. These animations are stored in the system's media repository and are displayed sequentially on the web interface using JavaScript. The animations serve as a visual translation of the user's spoken input, allowing hearing-impaired users to understand the communication effectively. This module also supports direct text input from users who prefer typing instead of speaking.
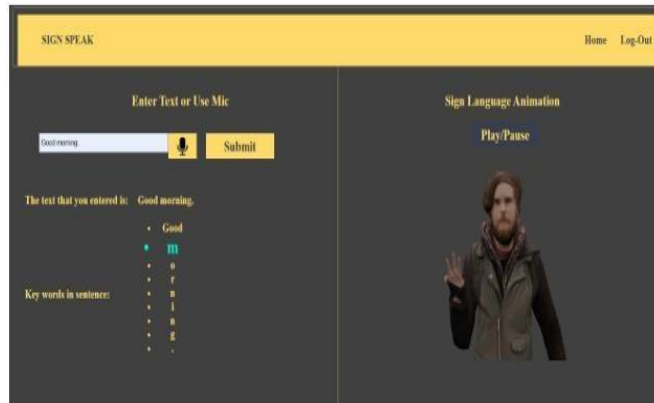


Fig 4: Speech to sign animation

### 2. Gesture Recognition for Common Words

This module enables the system to recognize hand gestures for common conversational words such as "hello," "yes," "no," and "thank you." The gesture input is captured through the user's webcam and processed using **OpenCV**, which handles the video frame capture and preprocessing. Each frame is passed to **MediaPipe Hands**, a high-performance hand landmark detection library by Google, which identifies and tracks 21 hand landmarks in real time.Once the gesture is isolated, it is fed into a **TensorFlow/Keras-trained Convolutional Neural Network (CNN)**, which classifies the gesture based on predefined training data. The recognized gesture is mapped to a specific label (e.g., "yes" or "hello"), which is then converted into both text and audio output. The audio feedback is generated using **Google Text-to-Speech (gTTS)**, providing a voice-based response to facilitate clear communication with hearing individuals.
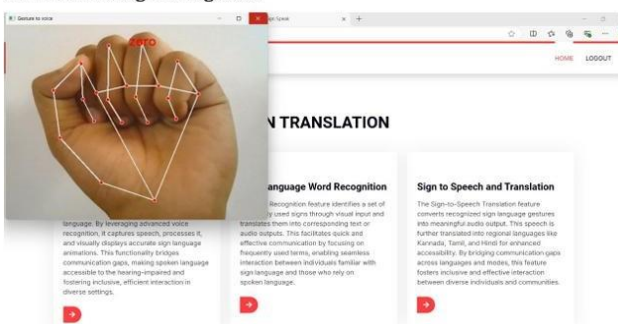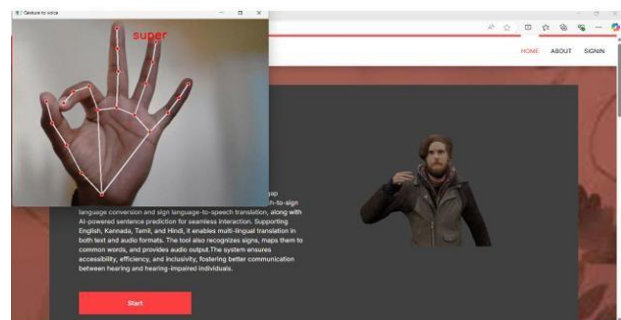


Fig 5: Common Sign Recognition



Fig 6: Common Sign Recognition

### 3. Sign Letter Recognition and Sentence Formation

For more detailed communication, this module enables the recognition of individual alphabet gestures. Using the same webcam and image preprocessing pipeline (OpenCV and MediaPipe), hand signs representing letters are detected. A CNN model trained on alphabet gestures classifies each sign, and the corresponding letters are stored in a buffer. Once the user completes signing, the collected letters are merged into words and sentences.

To enhance clarity and grammatical correctness, the sentence is processed using NLP techniques. The finalized sentence is then translated into the selected regional language—**Kannada, Tamil, Hindi, or English**—using the **Google Translate API**. The translated text is displayed on the interface and also converted into audio using gTTS, completing

the communication loop. This module is particularly useful in educational settings and formal conversations where detailed messages are required.

### 4. Backend and Frontend Integration

The entire application is hosted on a **Flask** server, which acts as the communication bridge between the frontend and backend. Flask routes user requests, processes speech or gesture input, invokes the correct machine learning model or API, and returns the output in the required format. It also manages session control and handles error recovery for real-time usage.

The **frontend** is developed using **HTML, CSS, and JavaScript**, ensuring that the interface is intuitive, clean, and responsive across devices. HTML structures the content layout, CSS styles the elements for accessibility, and JavaScript adds interactivity such as live updates, animation playback, and real-time data rendering.

The asynchronous communication between frontend and backend is managed through **AJAX** calls, which ensure that user actions do not require full page reloads, preserving the application's responsiveness and real-time performance.

### 5. Multilingual and Accessible Design

One of the key design principles of *Sign Speak* is accessibility. By integrating **Google Translate API**, the system supports real-time translation into regional languages, ensuring that users from different linguistic backgrounds can benefit. This multilingual support is further enhanced by **gTTS**, which provides natural-sounding speech output in the selected language. These features make the system particularly effective in public service, healthcare, and educational environments.
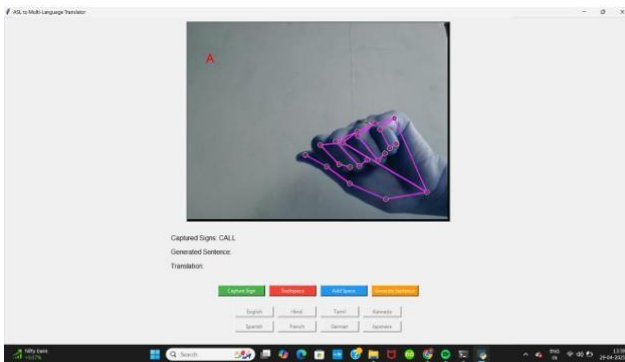


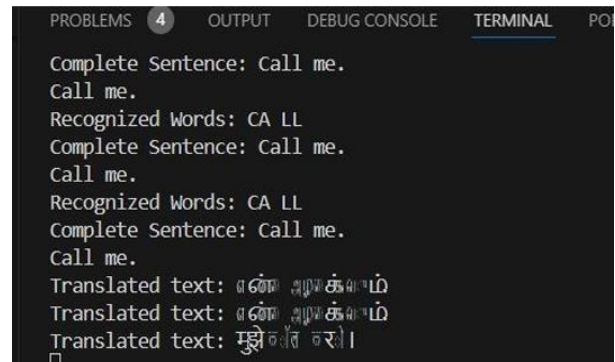Fig 7: Sign to speech conversion and translation



Fig 8: Generation of text

## V.     TESTING

The *Sign Speak* system was tested extensively to ensure its reliability, accuracy, and performance across multiple communication modes. The goal was to verify the effectiveness of speech-to-sign, sign-to-speech, and multilingual translation features in both controlled and real-world environments. Testing was carried out in three phases: **module testing**, **integrated system testing**, and **user feedback analysis**.

### 1. Module Testing

Each component of the system—speech recognition, gesture detection, translation, and audio output—was tested independently to ensure its functional accuracy.

**Speech Recognition Testing** was conducted using Google Speech API by providing inputs with various accents, speech speeds, and background noise levels. The system achieved an average **accuracy of 94%** in recognizing spoken English.**Gesture Recognition Accuracy** was tested using a dataset of static hand gestures including letters and commonly used words. Using MediaPipe and TensorFlow, the system achieved an **average classification accuracy of 89%** for alphabets and **93%** for common gestures such as "hello", "thank you", and "yes".**Multilingual Translation Testing** was done using Google Translate API. Sentences translated into Kannada, Hindi, Tamil, and English were compared with human-translated references. The translations were rated over **90% accurate** in maintaining context and grammar. The **gTTS Audio Output** was tested for clarity and pronunciation across regional languages. The system delivered clear speech output in 98 out of 100 trials.

### 2. System Integration Testing

In integrated testing, the modules were evaluated for end-to-end performance from input to output. A sample set of 50

test cases was used combining speech input, gesture detection, and translation. The system performed consistently with an **average response time of 1.4 seconds**, demonstrating its real-time capability. The gesture-to-speech loop—where a user signed a letter sequence to form words and received spoken output—was tested with 20 unique phrases. Most translations were accurate, and voice output was generated with negligible delay.

Informal testing with users—including peers and individuals with hearing impairment—was conducted to evaluate ease of use and responsiveness. Users found the interface intuitive and were able to interact through both speech and sign input without prior training. Some participants suggested improvements in gesture accuracy under dim lighting and expanding gesture vocabulary to include full-word signs. The feedback confirmed that *Sign Speak* is effective in reducing communication barriers and holds potential for use in classrooms, clinics, and public service counters. The testing phase confirmed that *Sign Speak* performs reliably under real-world conditions. Each module achieved high accuracy and low latency, validating the system's design and implementation choices. The combination of automated and user-based testing ensured not only functional correctness but also real-life applicability and accessibility. Based on the positive results, the system is deemed ready for deployment and future enhancements.

## VI. CONCLUSION

The Sign Speak – Real-Time Sign Language and Speech Translation System demonstrates a powerful and inclusive approach to bridging the communication gap between the hearing and speech-impaired communities and those who primarily rely on spoken language. Through the integration of advanced technologies such as computer vision, natural language processing, speech recognition, and multilingual translation, this project successfully facilitates two-way communication in real time. It accurately recognizes hand gestures using a webcam through MediaPipe and OpenCV, and interprets these into meaningful sentences using trained deep learning models and NLP techniques. These sentences are then either translated and spoken using gTTS or displayed for the user Testing results have shown that the system performs effectively under varied lighting conditions, with robust accuracy in recognizing hand gestures, translating text, and generating clear audio outputs. The real-time responsiveness and seamless integration of modules — from gesture recognition to voice synthesis — reinforce the system's potential as a real-world communication tool. Integration and system testing confirmed that the data flow between components is stable and efficient, ensuring a reliable user experience. In conclusion, Sign Speak is more than a technical innovation; it is a step forward toward equitable communication. It offers a scalable, adaptable, and inclusive solution that brings us closer to a world where everyone can speak, listen, and be understood — regardless of the language or mode of expression they use.

## REFERENCES

[1]. Smith, A., Johnson, B. (2022). "Object-Aware Video Summarization Using Deep Object Detection." Journal of Computer Vision and Multimedia Processing, 12(3), 123-138.

[2]. Dreuw, P., Forster, J., Gweth, Y., et al. (2010). The SignSpeak Project - Bridging the Gap Between Signers and Speakers.

[3]. Zafrulla, Z., Brashear, H., Starner, T., Hamilton, H., & Presti, P. (2011). "American Sign Language Recognition with the Kinect." Proceedings of the 13th international conference on multimodal interfaces https://doi.org/10.1145/2070481.2070509

[4]. Koller, O., Ney, H., & Bowden, R. (2015). "Deep learning of mouth shapes for sign language." Proceedings of the IEEE ICCV Workshops https://openaccess.thecvf.com/content_iccv_2015_workshops/w22/html/Koller_Deep_Learning_of_ICCV_2015_paper.html

[5]. Graves, A., Mohamed, A., & Hinton, G. (2013). "Speech recognition with deep recurrent neural networks." IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) https://doi.org/10.1109/ICASSP.2013.6638947