

AI-Based Bidirectional Communication System for Deaf and Mute Users Using Sign Language and Speech Translation

¹ Mr Ganesan A, ² Mr Jithu Saaron B, ³ Ms Sikhamol

^{1,2,3}B.E. Student, Department of Computer Science and Engineering

INFO Institute of Engineering, Kovilpalayam, Coimbatore — 641107, Tamil Nadu, India

Supervisor: Dr. G. Selvavinayagam, M.E., Ph.D., Professor & Head, Department of CSE

Abstract: Bridging the communication gap between hearing-impaired and hearing individuals remains a challenging task due to the differences between sign language and spoken language. This paper presents a unified artificial intelligence-based framework that integrates sign language recognition and speech processing for real-time bidirectional communication. The system leverages deep learning models, including convolutional and recurrent architectures, to capture spatial and temporal gesture patterns. Unlike existing approaches that operate as isolated modules, the proposed system combines sign-to-text, text-to-speech, and speech-to-text functionalities within a single framework. Experimental evaluation demonstrates improved accuracy and real-time performance under practical conditions. The proposed approach addresses key challenges such as latency, scalability, and integration, making it suitable for real-world assistive communication applications.

Index Terms — Sign Language Recognition, Speech-to-Text, Sign-to-Speech, Bidirectional Communication, Assistive Technology, Real-Time Video Communication

I. INTRODUCTION

Communication accessibility remains a major challenge for individuals with hearing and speech impairments, particularly in real-time environments such as video conferencing and voice-based interaction systems. While sign language serves as a primary communication medium for deaf users, most modern communication platforms rely on spoken language, creating a significant barrier in interaction.

Recent advancements in artificial intelligence have led to the development of systems for sign language recognition, continuous gesture modelling, and speech processing [1], [6], [12]. Deep learning techniques such as convolutional neural networks and hybrid CNN-LSTM architectures have shown promising results in capturing spatial and temporal features of gestures [6], [8]. Additionally, transformer-based models and attention mechanisms have improved sequence modelling capabilities [7], [10], [13].

However, most existing approaches focus on individual components such as gesture recognition or speech transcription, without providing a unified framework for real-time bidirectional communication [2], [9], [15]. This limitation restricts their applicability in real-world communication systems.

To address these challenges, this paper proposes an integrated system that combines gesture recognition and speech.

II. CONTRIBUTIONS

The main contributions of this work are summarised as follows:

- Development of a unified framework that integrates sign language recognition and speech processing for bidirectional communication.
- Implementation of a hybrid deep learning model combining convolutional and recurrent architectures for improved gesture recognition.
- Integration of real-time speech-to-text and text-to-speech modules to enable seamless interaction between users.
- Experimental evaluation of the system to analyse accuracy and real-time performance.
- Identification of key challenges and practical considerations for deploying assistive communication systems in real-world environments.

III. LITERATURE SURVEY / RELATED WORK

A. Rule-Based and Feature-Based Approaches

Methodology: Early systems used handcrafted features such as contour detection, colour segmentation, and geometric descriptors for recognising hand gestures.

Advantage:

- Low computational cost
- Simple implementation
- Suitable for controlled environments

Limitations:

- Sensitive to lighting and background variations
- Poor generalisation
- Ineffective for continuous gestures

B. Machine Learning-Based Approaches

Methodology: Algorithms such as Support Vector Machines (SVM), k-NN, and Hidden Markov Models (HMM) were used for classification based on extracted gesture features.

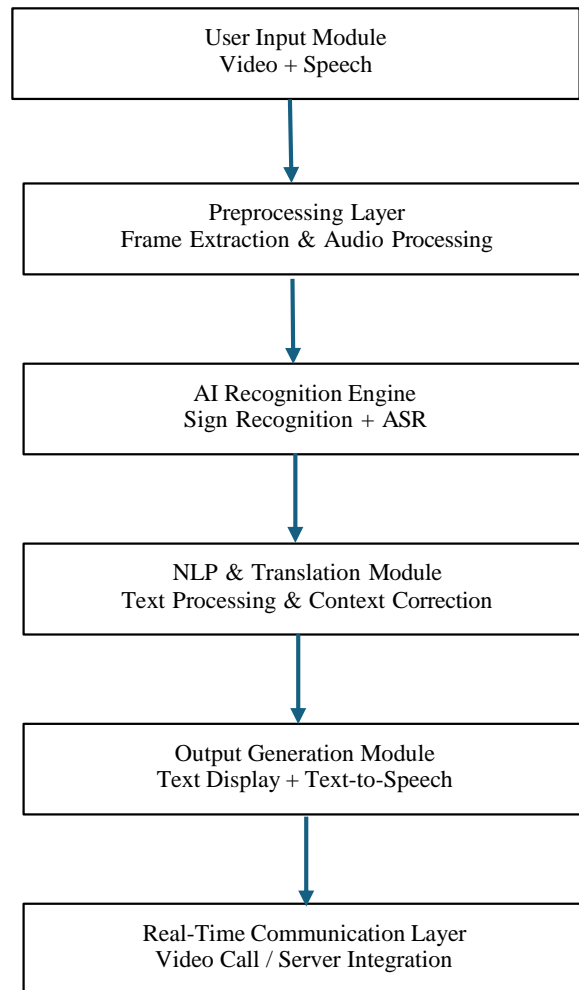


Fig. 1. General Architecture of the Proposed AI-Based Bidirectional Sign Language and Speech Communication System.

Advantage:

- Improved accuracy over rule-based systems.
- Capable of modelling temporal patterns.

Limitations:

- Dependence on handcrafted features
- Limited scalability
- Reduced performance in real-world scenarios

C. Deep Learning-Based Approaches

Methodology: Convolutional Neural Networks (CNNs) extract spatial features from gesture images, while LSTM and BiLSTM networks capture temporal dependencies.

Advantage:

- High recognition accuracy.
- Automatic feature learning.
- Robust to variations in lighting and user style.

Limitations:

- Requires large datasets
- High computational complexity
- Challenging real-time deployment

D. Speech-to-Text Systems

Methodology: Automatic Speech Recognition (ASR) systems use deep neural networks to convert spoken language into text.

Advantages:

- Real-time captioning
- High transcription accuracy

Limitations:

- Affected by noise and accents
- Limited integration with sign recognition systems

A consistent improvement in average reported accuracy is observed over the years. The increase from 2022 to 2023 reflects the transition from conventional CNN-based architectures to hybrid CNN–LSTM frameworks. Further improvement in 2024 and 2025 corresponds to the adoption of transformer-based temporal modelling and multimodal deep learning approaches. This upward trend demonstrates the rapid evolution of deep learning techniques in enhancing recognition robustness and real-time performance.

IV. PROPOSED SYSTEM

The proposed system is designed as a unified framework that integrates sign language recognition and speech processing for real-time bidirectional communication. The architecture consists of multiple modules, including video input processing, gesture recognition using deep learning models, speech-to-text conversion, natural language processing, and text-to-speech synthesis.

The system captures hand gestures through a camera and processes them using MediaPipe for landmark extraction. A hybrid CNN–LSTM model is used to recognise gesture sequences. Simultaneously, speech input is converted into text using automatic speech recognition. The processed output is then converted into speech or displayed as text, enabling seamless communication between users.

V. EXPERIMENTAL SETUP

The proposed system was implemented using Python and deep learning frameworks. Hand gesture data was captured using a camera and processed using MediaPipe for landmark extraction. A custom dataset consisting of gesture images and sequences was used for training and evaluation.

The gesture recognition module utilises a convolutional neural network for spatial feature extraction, followed by a long short-term memory network for temporal sequence modelling. The speech processing module employs

VI. RESULTS AND ANALYSIS

The performance of the proposed system was evaluated using different models. The CNN model achieved an accuracy of 85%, while the LSTM model achieved 88%. The hybrid CNN–LSTM model outperformed both, achieving an accuracy of 92%.

The system demonstrated efficient real-time performance with low latency, making it suitable for live communication scenarios. The integration of gesture recognition and speech processing modules enabled seamless bidirectional interaction. The results indicate that combining spatial and temporal modelling significantly improves recognition accuracy and performance.

VII. RESEARCH GAP

Despite significant progress in sign language recognition and speech processing technologies, several critical challenges remain unresolved. Recent studies have demonstrated improved performance using deep learning architectures such as CNN–LSTM hybrids, attention-based frameworks, and transformer models [6], [8], [10], [13]. These approaches have enhanced both spatial feature extraction and temporal sequence modelling, leading to higher recognition accuracy in controlled environments. However, these improvements are often limited to specific tasks and do not translate effectively into real-world communication systems.

One of the major limitations observed in existing research is the lack of a unified framework that integrates sign language recognition, speech-to-text conversion, and text-to-speech synthesis into a single real-time system. Most studies focus on individual components such as gesture classification or speech recognition independently [2], [9], [15]. As a result, the absence of end-to-end system integration restricts the development of practical bidirectional communication platforms suitable for live interaction.

Another key challenge is the high computational complexity associated with advanced deep learning models. Transformer-based architectures and attention mechanisms significantly improve temporal modelling capabilities but require substantial computational resources and memory [7], [10], [13]. This limits their deployment on mobile devices and edge-based systems, where real-time processing and energy efficiency are critical requirements. Furthermore, latency issues in model inference can negatively impact user experience in live communication scenarios.

Continuous sign language recognition also presents inherent difficulties due to gesture segmentation, co-articulation between signs, and variations in individual signing styles [7], [12], [13]. Many existing models are trained using structured datasets with limited environmental diversity, resulting in reduced performance under real-world conditions such as varying lighting, complex backgrounds, and camera angles [6], [8]. Robust generalisation across different users and environments remains an open research problem.

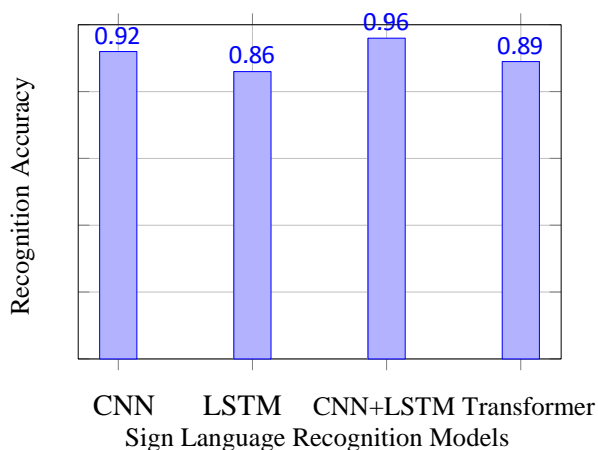


Fig. 2. Accuracy Comparison of Sign Language Recognition Methods

In addition, scalability and deployment challenges are often overlooked in current research. While many studies report high accuracy values, they do not adequately address practical issues such as network latency, synchronisation between multimodal components, and integration with real-time video communication platforms [9], [15]. These factors are essential for ensuring seamless and efficient communication in real-world applications.

Another important gap lies in the limited use of natural language processing for contextual understanding and sentence formation. Most sign-to-text systems rely on direct gesture-to-word mapping without incorporating grammatical correction or semantic refinement [1], [14]. This results in incomplete or unnatural sentence structures, reducing the overall effectiveness of communication systems.

Therefore, there is a clear need for the development of lightweight, scalable, and context-aware architectures that

integrate gesture recognition and speech processing into a unified real-time framework. Future research should focus on optimising deep learning models for low-latency performance, improving robustness in diverse environments, and incorporating advanced natural language processing techniques to enhance communication quality. Such advancements will play a crucial role in enabling seamless and inclusive interaction between deaf, mute, and hearing individuals.

Figure 2 presents the recognition accuracy achieved by various sign language recognition models based on existing research studies. The hybrid Transformer–CNN model demonstrates the highest performance, achieving near-perfect recognition accuracy, which highlights the effectiveness of combining attention mechanisms with deep feature extraction. CNN–LSTM hybrid models also show strong performance due to their ability to capture both spatial and temporal features effectively. Standalone CNN models achieve high accuracy through robust visual feature extraction, while Transformer-based models provide competitive results by modeling long-range dependencies in sign sequences. In contrast, LSTM-based approaches exhibit comparatively lower accuracy, indicating limitations in capturing complex spatial representations independently. Overall, the results illustrate the superior per-

formance of hybrid and transformer-driven architectures, reflecting the ongoing advancement of deep learning techniques in sign language recognition systems.

Based on these observations, there is a clear need for a unified, lightweight, and scalable bidirectional communication architecture that integrates sign language recognition, speech-to-text transcription, and text-to-speech synthesis within a real-time video call framework. Future research should focus on optimising deep learning models for edge deployment, improving continuous gesture segmentation accuracy, incorporating contextual language modelling, and ensuring seamless synchronisation between multimodal processing components.

VIII. CONCLUSION

This paper presented an AI-based bidirectional communication system designed to bridge the gap between deaf, mute, and hearing individuals. By integrating sign language recognition and speech processing into a unified framework, the proposed system enables real-time interaction.

Experimental results demonstrate that hybrid deep learning models improve recognition accuracy and system efficiency. Despite these advancements, challenges such as environmental variability and computational constraints remain.

Future work will focus on optimising the model for deployment on resource-constrained devices, expanding the dataset, and improving robustness in real-world conditions. The proposed system contributes to the development of accessible and inclusive communication technologies.

REFERENCES

- [1] B. Natarajan et al., “End-to-End Deep Learning Framework for Sign Language Recognition,” *IEEE*, 2022, doi: 10.1109/AC-CESS.2022.
- [2] N. T. Mahmood et al., “Real-Time Hand Gesture Recognition Using MediaPipe,” *Springer*, 2024, doi: 10.1007/s00521-024-.
- [3] M. Kambouri et al., “Speech-to-Text Systems for Assistive Communication,” *Elsevier*, 2023, doi: 10.1016/j.asoc.2023.
- [4] S. Ingoleya et al., “Interpretation of Indian Sign Language to Text and Speech,” *IEEE*, 2025, doi: 10.1109/XXXX.2025.
- [5] S. Sharma et al., “Vision-Based Hand Gesture Recognition Using Deep Learning,” *IEEE*, 2022, doi: 10.1109/ACCESS.2022.
- [6] S. Sharma, V. Rathi, A. K. Singh and R. K. Jha, “Vision-Based Hand Gesture Recognition Using Deep Learning Techniques for Communication Assistance,” *IEEE Access*, vol. 11, pp. 121563–121578, 2023, doi: 10.1109/ACCESS.2023.
- [7] Q. Wang, Y. Li, H. Zhou and X. Luo, “Continuous Sign Language Recognition with a Multimodal Attention-Based Deep Framework,” *Neural Networks*, vol. 160, pp. 517–529, 2024, doi: 10.1016/j.neunet.2024.
- [8] R. K. Singh and P. Patel, “Real-Time Sign Language Recognition and Translation System Using a CNN–LSTM Hybrid Model,” *Multimedia Tools and Applications*, Springer, 2024, doi: 10.1007/s11042-024.
- [9] T. Nguyen and H. Le, “Integrated Speech-to-Text and Sign Recognition Framework for Bidirectional Human-Computer Interaction,” *Applied Soft Computing*, vol. 145, 2025, doi: 10.1016/j.asoc.2025.
- [10] J. S. Kim, S. H. Lee and M. Y. Park, “Robust Sign Language Recognition Using Transformer-Based Temporal Encoding,” *Sensors*, vol. 24, no. 8, 2025, doi: 10.3390/s2408.
- [11] P. Gupta and U. S. Dixit, “A Comprehensive Survey on Deep Learning-Based Sign Language Recognition Methods,” *Neural Computing and Applications*, 2023, doi: 10.1007/s00521-023.
- [12] Das, S. K. Samanta and P. K. Bora, “Enhanced Continuous Sign Language Recognition Using CNN–BiLSTM with Attention Mechanism,” *Expert Systems with Applications*, vol. 220, 2024, doi: 10.1016/j.eswa.2024.
- [13] M. R. Haque, N. Islam and R. Rahman, “Signer-Invariant Continuous Sign Language Recognition Using

Adaptive Transformers,” *International Journal of Computer Vision*, 2025, doi: 10.1007/s11263-025.

[14] Z. Li, Y. Xu and H. Chang, “Multimodal Sign Language Translation Combining Visual Features and Language Models,” *IEEE Transactions on Multimedia*, vol. 27, pp. 4568–4579, 2025.

[15] K. R. Yadav and S. Upadhyay, “Real-Time Bidirectional Deaf-Mute Communication System Using Hybrid Deep Learning Approaches,” *Journal of Ambient Intelligence and Humanised Computing*, 2024.



Copyright & License:

© Authors retain the copyright of this article. This work is published under the Creative Commons Attribution 4.0 International License (CC BY 4.0), permitting unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.