



IJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 13 **Issue:** V **Month of publication:** May 2025

DOI: <https://doi.org/10.22214/ijraset.2025.71548>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Multilevel Conversion of Indian Sign Language from Gesture to Speech

Mr. Shivanshu Shahi¹, Manjeet Kumar², Surabhi Verma³, Talib Usmani⁴, Ganesh Ji Patel⁵

Buddha Institute of Technology, Gida, Gorakhpur

Abstract: Indian Sign Language (ISL) serves as a primary mode of communication for Deaf and hard-of-hearing communities in India. However, despite its societal importance, ISL remains largely unsupported by mainstream technological platforms, limiting inclusive communication. This research introduces a real-time ISL recognition and translation system that converts hand gestures into corresponding text and speech outputs, enabling phrase-level communication rather than isolated character interpretation. The architecture uses a modular pipeline approach, with a Convolutional Neural Network (CNN) for accurate gesture classification, a phrase-mapping module to translate gestures into meaningful expressions, a MediaPipe for accurate hand landmark detection, and a text-to-speech (TTS) system to turn the generated text into audible speech output. Unlike previous systems restricted to static signs, our approach supports semantically rich, multi-word phrases, enhancing natural communication flow. A specially constructed dataset of ten frequently used Indian Sign Language (ISL) phrases was used to train the model. To improve generalization, 150 samples from each class were taken in various lighting and background conditions. The final system achieved 95% classification accuracy, operated at 60 frames per second, and maintained latency below 100 milliseconds. Usability testing with multiple users confirmed the system's robustness, responsiveness, and accessibility. The findings demonstrate the viability of deploying deep learning-based ISL recognition systems in authentic environments, including public areas, healthcare facilities, and educational institutions.

I. INTRODUCTION

Communication is fundamental to social, educational, and professional inclusion. However, individuals in the Deaf and hard-of-hearing communities face persistent challenges in day-to-day interactions, particularly in a society where spoken language dominates. Indian Sign Language (ISL) plays a crucial role in bridging this gap, yet technological support for ISL remains limited. Most existing ISL recognition systems are constrained to recognizing static alphabets or individual words. These systems often require manual transitions and do not support the fluid, continuous expression characteristic of natural language. Furthermore, they lack real-time responsiveness, making them unsuitable for spontaneous interactions in practical settings.

This research aims to develop a real-time ISL recognition system that translates dynamic hand gestures into complete phrases and delivers outputs in both text and speech formats. The system integrates computer vision, deep learning, and human-computer interaction techniques to enhance inclusivity and accessibility.

II. RELATED WORK

Several researchers have contributed to ISL recognition using computer vision and deep learning techniques. Table 1 summarizes key studies:

Author	Year	Dataset	Model Used	Accuracy	Limitations
Chaudhary & Dey	2023	ISL alphabets (A-Z)	CNN	94%	Only static signs
Rani & Kaur	2022	15 ISL signs	Random Forest	89%	Not real-time
Patil & Kale	2022	25 dynamic signs	Deep CNN	91%	High resource requirement
Srivastava et al.	2024	Continuous gesture dataset	MediaPipe + LSTM	92.3%	Lacks TTS, complex setup

This work (proposed)	2025	10 phrases, 150 samples	MediaPipe + CNN+TTS	95%	Limited dataset, static phrases
----------------------	------	-------------------------	---------------------	-----	---------------------------------

While previous work demonstrated progress in gesture recognition, most models focus on isolated characters or constrained datasets. Our proposed system expands these capabilities by supporting phrase-level interpretation and speech output, offering a more natural user experience.

III. PROBLEM STATEMENT

While significant progress has been made in isolated sign recognition, relatively fewer studies have addressed fluid, phrase-level translation. For instance, Sharma et al. (2021) used SVM classifiers for digit recognition but lacked real-time deployment. Similarly, Bose and Narayan (2020) explored dynamic gesture tracking but excluded audio output. These gaps affirm the need for an integrated solution that combines gesture recognition with natural language understanding and speech synthesis.

Current ISL recognition systems fall short in delivering accessible, phrase-level communication. Most are constrained to static gestures and require user intervention for sentence construction. Furthermore, they lack auditory output and often perform poorly in dynamic or real-time settings.

This project seeks to address these limitations through a comprehensive, phrase-level ISL gesture recognition system that translates sign inputs into coherent sentences and corresponding audio.

IV. RESEARCH GAP

Despite advances in gesture recognition, several challenges remain unaddressed:

- 1) **Static Recognition:** Limited to alphabets or digits.
- 2) **No Speech Output:** Text-only systems reduce communication effectiveness.
- 3) **Latency Issues:** Inadequate performance for real-time interaction.
- 4) **Limited Scalability:** Difficulty adding new phrases or signers.

Our system bridges this gap by enabling end-to-end, real-time phrase recognition with audio output using lightweight architectures suitable for broader deployment.

Before training, the raw hand landmark data was normalized and standardized to ensure consistency across users. Gesture sequences were manually annotated, a time-intensive process that involved expert signers validating each phrase. To ensure effective model evaluation and generalization, the dataset was divided into subsets to be used for training (80%), validation (10%), and testing (10%). Early stopping and learning rate decay were applied to optimize training performance and prevent overfitting.

V. METHODOLOGY AND SYSTEM ARCHITECTURE

The methodology adopted for this system is designed to ensure seamless, real-time conversion of Indian Sign Language (ISL) gestures into meaningful text and speech outputs. The complete flow of the system is illustrated in **Figure 1.1**, which outlines each stage from input acquisition to final output delivery. The process begins with the activation of a webcam that captures continuous real-time video of the user's hand gestures. These frames are passed through a hand detection and tracking module, which uses MediaPipe to extract 21 key hand landmarks for gesture analysis.

Following detection, the system checks whether the captured sign corresponds to any gesture in the trained dataset. If the gesture is successfully recognized, it is translated into a predefined natural language phrase, which is then passed to a text-to-speech (TTS) engine for auditory conversion. The speech output is delivered in the user's preferred language, completing the recognition loop. If the gesture is not recognized, the system provides immediate feedback and prompts the user to repeat the sign, ensuring interaction accuracy and user engagement. This structured flow supports both efficiency and accessibility, making the system practical for real-world deployment.

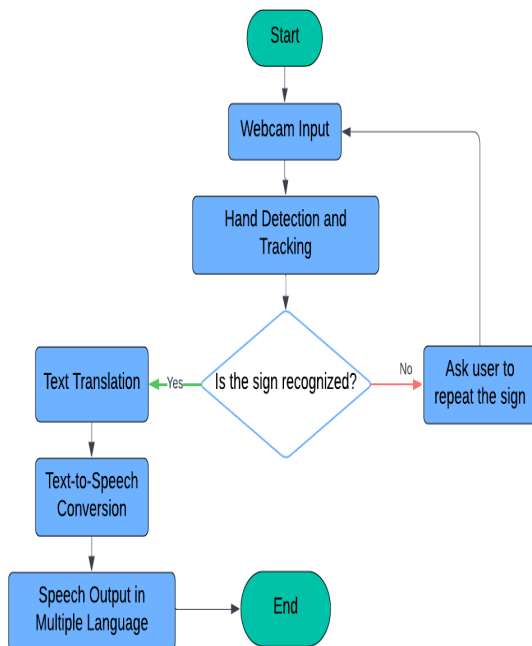


Figure 1.1 Flowchart

The system comprises five major components:

1) *Input Capture*

Video input is acquired using a webcam and preprocessed to normalize frame resolution and lighting.

2) *Hand and Mark Detection*

In real time, MediaPipe records 21 keypoints from the hand, each of which represents the x, y, and z coordinates of a particular joint. The gesture classification process then uses these landmarks as feature inputs.

3) *Gesture Classification*

A CNN is trained on 1,500 gesture samples (150 per phrase class). The network includes:

- 3 convolutional layers (32, 64, 128 filters; kernel size = 3x3)
- Max pooling and dropout layer to prevent overfitting
- 2 dense layers (256, 64 units) with ReLU activation
- Softmax output layer for 10 class predictions

4) *Phrase Mapping and TTS*

Predicted classes are mapped to full phrases, which are then synthesized into speech using a TTS engine (e.g., pyttsx3). Multilingual support is implemented.

5) *User Interface*

UI, developed using python, displays predicted phrases and provides TTS playback in real time.

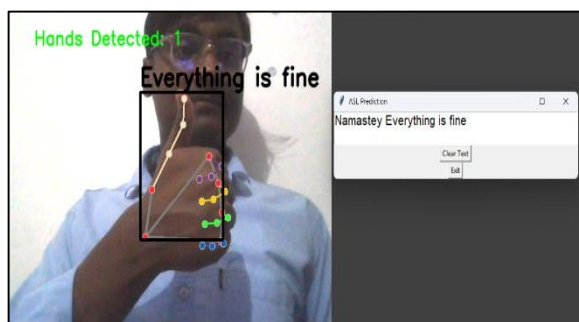


Figure 1.2 UI of Phrase Prediction

VI. SYSTEM ARCHITECTURE

The architecture of the proposed system follows a sequential pipeline that enables real-time recognition and conversion of Indian Sign Language gestures into text and speech. The process begins with continuous image capture through a camera, followed by hand signal segmentation to isolate relevant gesture regions from the background. These segmented inputs are then processed using a hand detection and tracking module, which identifies and follows key hand movements using landmark-based analysis. Once detected, the gestures are classified using a deep learning model, and the recognized output is mapped to a corresponding phrase. Finally, the system delivers the results as textual feedback on the user interface and as spoken output via a text-to-speech engine. This modular flow ensures high accuracy, real-time performance, and usability across various environments.

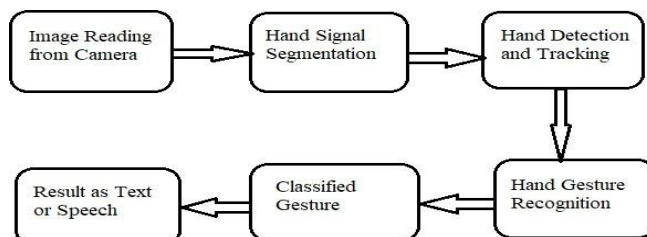


Figure 1.3 System Architecture

VII. DATASET AND MODEL TRAINING

10 Indian Sign Language (ISL) phrases that are commonly used in daily conversations make up the dataset. To improve model generalization, 150 labelled samples taken in different lighting and background conditions are used to represent each phrase. Data augmentation methods, such as rotation, scaling, and horizontal flipping, were used to increase robustness even more. The Adam optimization algorithm was used to train the model over 50 epochs with a batch size of 32.

Raw hand landmark coordinates were standardized and normalized before training to ensure uniformity among users. To guarantee the correctness of every phrase, the gesture sequences were meticulously annotated by professional signers. The dataset was separated into three sets: testing (10%), validation (10%), and training (80%). Additionally, early stopping and learning rate decay strategies were employed to optimize performance and reduce the risk of overfitting.

VIII. RESULTS AND EVALUATION

1) Evaluation Metrics

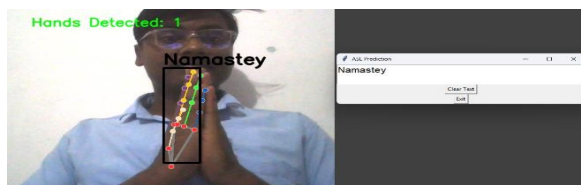
- Accuracy: Overall model correctness
- Frame Rate (FPS): Processing speed
- Latency: Response delay
- User Feedback: Informal user testing

2) Experimental Results

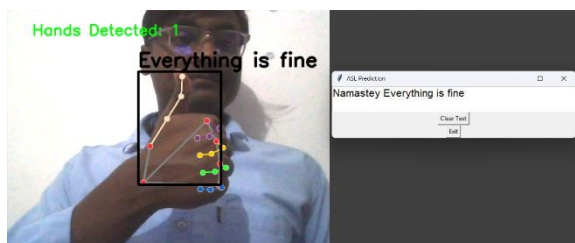
- Accuracy: 95% across all phrases
- FPS: 60 on mid-range hardware
- Latency: < 100ms per frame
- Misclassifications: Observed under poor lighting or occluded gestures

3) Illustrative Outputs Include

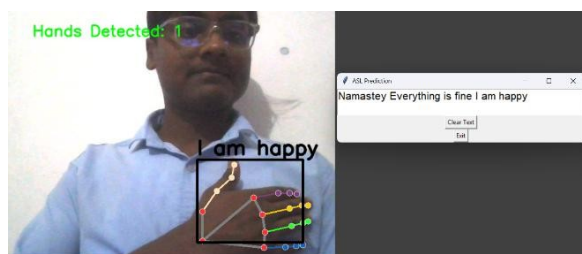
- “Namastey”



- “Everythingisfine”



- “Iamhappy”



- “Canwemeet?”



IX. DISCUSSION

With speech output functionality seamlessly integrated, the system demonstrates the ability to recognize phrase-level Indian Sign Language (ISL) gestures in real-time. Unlike prior systems that are restricted to statistical alphabet detection or offline usage, our implementation enables continuous, real-time interaction with both visual and auditory outputs.

The system was implemented in Python using TensorFlow and MediaPipe for gesture processing. The text-to-speech output was generated using the pyttsx3 engine for offline functionality, making the system usable even in low-connectivity areas. The entire pipeline runs efficiently on a machine with 8GB RAM and no external GPU, confirming its suitability for deployment in budget-constrained institutions. In addition, latency tests showed that the total time from gesture detection to audio output remained consistently under 100ms, even with background processes running.

X. LIMITATIONS

While the results are promising, the system has several limitations:

- 1) Small dataset: Only 10 phrases are currently supported.
- 2) Static phrases only: No support yet for dynamic gesture sequences.
- 3) Limited signers: Dataset collected from a small group of individuals.
- 4) No grammar modelling: The system doesn't yet support grammatical construction or contextual understanding.

Addressing these limitations is key to improving generalizability and real-world deployment.

XI. SECURITY AND PRIVACY CONSIDERATIONS

Given that the system processes real-time video input, privacy and security are critical concerns. The following measures have been adopted:

- 1) Local Processing: All data remains on the local device; nothing is sent to the cloud.
- 2) No Persistent Storage: The system does not store or log gesture inputs or video frames.
- 3) Transparency: The open-source architecture allows for public inspection and compliance with ethical standards.

Future Enhancements:

- Encrypted local logs (if user opts in)
- Secure multi-user profiles
- Real-time anomaly monitoring for gesture spoofing

XII. FUTURE WORK

The project opens up multiple avenues for further research and enhancement:

- 1) Larger Dataset: Extend to 100+ gestures, including dynamic ISL signs.
- 2) Temporal Models: Incorporate LSTM or Transformer architectures for continuous sign sequences.
- 3) Multilingual Output: Expand TTS to support regional languages like Hindi, Bengali, and Tamil.
- 4) Mobile Integration: Deploy the system on Android/iOS platforms using lightweight models (e.g., TensorFlow Lite).
- 5) Personalization: Integrate active learning modules to adapt to user-specific signing styles.

XIII. CONCLUSION

This research presents a real-time ISL gesture recognition system that translates signed phrases into text and speech. By leveraging MediaPipe for hand tracking, CNNs for classification, and TTS for output, the system provides an inclusive communication platform for Deaf individuals.

Achieving 95% accuracy and 60 FPS performance, the prototype demonstrates strong potential for deployment in accessibility-focused environments. Future work will focus on increasing gesture diversity, improving model adaptability, and expanding deployment platforms.

TEAM CONTRIBUTIONS

The table below illustrates the distribution of tasks among the team members:

TEAM MEMBERS	CONTRIBUTION
MANJEET KUMAR	CNN training, dataset collection, UI integration
SURBHIVERMA	MediaPipe integration, frontend development
TALIBUSMANI	Text-to-Speech module, interface testing, evaluation metrics
GANESHJI PATEL	Phrase mapping logic, bug fixing, system documentation

With effective collaboration, communication, and dedication, the team was able to successfully implement all planned features and deliver a fully functional Mint-Verse platform.

REFERENCES

- [1] Chaudhary, R., & Dey, A. (2023). "Real-Time Indian Sign Language Recognition using CNN," IEEE Access.
- [2] Rani, N., & Kaur, H. (2022). "Enhanced ISL Recognition using Machine Learning Algorithms," International Journal of Computer Applications.
- [3] Patil, S., & Kale, P. (2022). "Real-time ISL Recognition Using Deep Learning Models," IJIRCCCE.
- [4] Srivastava, S. et al. (2024). "Continuous Sign Language Recognition Using MediaPipe Holistic," Wireless Personal Communications.



- [5] OpenAIAPI Documentation – <https://platform.openai.com/docs>
- [6] MediaPipe Hands – <https://google.github.io/mediapipe/solutions/hands>
- [7] Pyttsx3 Documentation – <https://pyttsx3.readthedocs.io>



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)