



IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 13 Issue: VII Month of publication: July 2025 DOI: https://doi.org/10.22214/ijraset.2025.73121

www.ijraset.com

Call: 🕥 08813907089 🔰 E-mail ID: ijraset@gmail.com



Text Conversion of Sign Language by Hand Gestures using Deep Learning

Gattadi Vikranth¹, G. Narasimham²

¹Post Graduate Student, M. Tech (Data Science), ²Associate Professor, Department of Information Technology, Jawaharlal Nehru Technological University Hyderabad, UCESTH

Abstract: Communication can be a major hurdle for individuals with hearing and speech impairments, particularly when interacting with people who do not understand sign language. Existing methods such as using human interpreters or written messages often fall short due to cost, inconvenience, or lack of real-time interaction. To address this challenge, we developed an intelligent system that recognizes American Sign Language (ASL) fingerspelling gestures and translates them into both text and speech in real time. Using computer vision and a Convolutional Neural Network (CNN), the system processes hand gestures captured via webcam and identifies corresponding alphabet characters. Hand landmark detection is carried out using the cvzone HandTrackingModule, and recognized letters are displayed within a graphical user interface created using Tkinter. The interface also offers suggested word predictions with the help of the Enchant dictionary and provides audio feedback through the pyttsx3 library. Additionally, gesture-based commands like space, clear, and backspace make the system more interactive and user-friendly. This solution aims to support seamless, accessible communication, especially in educational and assistive settings. Keywords: American Sign Language (ASL), Fingerspelling Recognition, Convolutional Neural Network (CNN), Computer Vision, Real-Time Communication, Hand Gesture Detection, Accessibility, Assistive Technology, Tkinter GUI, Text-to-Speech, cvzone Hand Tracking Module, Enchant Dictionary.

I. INTRODUCTION

Advancements in artificial intelligence and computer vision have significantly expanded the development of assistive technologies that can bridge communication barriers. One promising application of these innovations is in aiding individuals with hearing or speech impairments—particularly in situations where sign language is not widely understood or recognized.

This project introduces a real-time system designed to translate American Sign Language (ASL) fingerspelling gestures into both written text and spoken words. The goal is to support more inclusive and seamless communication by converting hand gestures into understandable formats for a broader audience. The system uses a webcam to capture hand gestures, which are then analysed using a Convolutional Neural Network (CNN) trained specifically for ASL alphabet recognition. For precise tracking and interpretation of hand movements, the solution utilizes the HandTrackingModule from the cvzone library. This ensures accurate identification of each fingerspelling gesture. A user-friendly interface developed with Tkinter provides a live video feed, displays the detected letters, constructs complete words, and offers spelling suggestions with the help of the Enchant dictionary. To further enhance the experience, a text-to-speech feature powered by the pyttsx3 library delivers real-time audio feedback. Additionally, users can perform special gestures to trigger actions such as inserting spaces, deleting characters, or confirming input.

By integrating gesture recognition, intelligent text prediction, speech output, and gesture-based commands, this application offers a comprehensive communication tool. It is particularly useful for individuals with speech or hearing difficulties and also serves as a valuable resource for learners of sign language. The system emphasizes fingerspelling-based input and focuses on building full words by identifying one gesture at a time.



Fig.1 ASL fingerspelling chart illustrating hand gestures for each alphabet letter A to Z.



International Journal for Research in Applied Science & Engineering Technology (IJRASET)

ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538 Volume 13 Issue VII July 2025- Available at www.ijraset.com

II. LITERATURE SURVEY

The rapid advancements in artificial intelligence—especially in deep learning and computer vision—have significantly propelled the development of gesture and sign language recognition systems. These technological strides have enabled the creation of highly accurate and fast systems capable of understanding human gestures, offering meaningful support for individuals with hearing or speech challenges. This review presents pivotal research contributions in hand gesture recognition and American Sign Language (ASL) interpretation, which serve as the basis for the current study.

A notable early contribution in this domain came from Molchanov et al. (2016), who introduced a convolutional neural network (CNN)-based approach for real-time dynamic hand gesture recognition. By extracting spatiotemporal features from video sequences, their method achieved superior performance compared to traditional image-based techniques. This research demonstrated the practical potential of deep learning in real-time gesture applications and set the stage for further innovations [1].

Further progress was made by Simon et al. (2017), who improved hand key point detection using a technique called Multiview bootstrapping. Their system was trained on 3D data and could accurately estimate 2D hand landmarks from single RGB images. This advancement was especially important for fine-grained gesture tasks, such as fingerspelling in ASL, where accurate finger position detection is crucial [2].

In 2019, Google Research introduced Media Pipe Hands, a framework that brought efficient and accurate on-device hand tracking to the forefront. By integrating palm detection with a 21-point hand landmark model, this tool enabled real-time gesture recognition with low computational demands, making it an ideal solution for applications like sign language interpretation [3].

Kaur and Singh (2020) developed a deep learning-based system aimed at recognizing ASL alphabets using webcam input. Their CNN model translated static hand gestures into corresponding text in real-time, focusing on both performance and user-friendly design. Their work further demonstrated the power of combining deep learning and computer vision in developing accessible tools for the hearing-impaired community [4].

III. OBJECTIVE

The primary objective of this project is to design an integrated real-time system that recognizes American Sign Language (ASL) hand gestures and translates them into both text and speech, aiming to enhance communication for individuals with hearing and speech impairments. The system utilizes a webcam for gesture input and employs a Convolutional Neural Network (CNN) along with hand landmark tracking via the cvzone HandTrackingModule to accurately identify ASL alphabets. To ensure an interactive and user-friendly experience, a custom graphical interface built with Tkinter displays live video, recognized characters, constructed words, and suggested terms using the enchant dictionary. Furthermore, the application incorporates text-to-speech functionality for voice output and supports intuitive gesture-based commands such as space, backspace, and clear, making it a practical and accessible communication aid for both daily use and educational purposes.

IV. SYSTEM ANALYSIS

A. Existing System

Existing sign language recognition systems typically use CNN-based models combined with tools like Media Pipe or cvzone for detecting hand landmarks from webcam input, focusing mainly on recognizing static ASL alphabets and converting them into text. While some systems offer basic GUI interfaces and integrate text-to-speech for audio output, they often suffer from key limitations. These include poor performance in low lighting or cluttered backgrounds, dependence on high-quality webcam input, and the need for precise hand positioning. Most models lack support for dynamic gestures and sentence-level understanding, leading to fragmented communication. Additionally, high CPU usage on low-end devices and a restricted gesture vocabulary—limited to trained signs—reduce the scalability and real-world effectiveness of such systems.

B. Proposed System

The proposed system offers a real-time solution for translating American Sign Language (ASL) hand gestures into both text and speech, aiming to improve communication for individuals with hearing or speech impairments. It captures live hand movements using a webcam and accurately detects hand landmarks through the cvzone HandTrackingModule. These gestures are processed using a Convolutional Neural Network (CNN) to identify ASL alphabets. The system features a user-friendly Tkinter-based interface that displays the video feed, recognized characters, sentence construction, and word suggestions using the enchant dictionary. It also includes gesture-based controls such as space, backspace, and next for smooth interaction. To further enhance accessibility, a text-to-speech component (pyttsx3) provides voice feedback. This integrated approach ensures high accuracy, and supports natural communication without requiring any additional hardware, making it suitable for both educational applications.



User Layer User Layer User processing Layer Processing Layer Hand Rol data Rol data Rol data Rol data Processing Rol data Rol

V. SYSTEM ARCHITECTURE AND METHODOLOGY

A. System Architecture



The architecture of the Sign Language Recognition System is organized into three main layers: User Layer, Processing Layer, and Prediction Layer. In the User Layer, a webcam captures real-time video of hand gestures performed by the user. The OpenCV module processes this video stream by extracting individual frames. In the Processing Layer, each frame undergoes hand detection through a Hand Extraction module, where the hand region is isolated and pre-processed via resizing, normalization, and background removal to ensure consistent input quality. The refined frame is then passed to the Prediction Layer, where a trained Convolutional Neural Network (CNN) analyses it and predicts the corresponding alphabet or character. This prediction is returned to the interface and displayed visually, while gesture-based commands like space and backspace are also managed. Additionally, the system features a text-to-speech (TTS) component that converts recognized text into spoken words, enabling users to receive both visual and audio feedback, thus supporting efficient and accessible communication.

B. Methodology



Fig.3 Block diagram: ALS Gesture Recognition System



International Journal for Research in Applied Science & Engineering Technology (IJRASET) ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538 Volume 13 Issue VII July 2025- Available at www.ijraset.com

- 1. Gesture Capture
 - The process begins with real-time gesture input captured through a webcam. This live video stream provides continuous frames that represent the user's hand movements for American Sign Language (ASL).
- 2. Preprocessing
 - Each captured frame is standardized by resizing it to a consistent dimension, ensuring compatibility with the neural network. Techniques like histogram equalization are applied to enhance contrast, helping the system to better distinguish hand features even in challenging lighting environments.
- 3. Data Storage
 - The pre-processed images, particularly those in RGB format, are organized and stored in a dataset. These images serve as training and testing samples for developing the gesture recognition model.
- 4. Model Training
 - The dataset is split into training and validation sets. A Convolutional Neural Network (CNN) is trained using the training data to learn spatial patterns and features associated with each ASL alphabet. This stage is crucial for enabling accurate gesture classification.
- 5. Model Evaluation and Optimization
 - The trained model is evaluated using test data to measure its accuracy. If results are below expectations, the system undergoes fine-tuning, where adjustments such as learning rate, number of layers, and filter sizes are made to improve performance.
- 6. Gesture Prediction and Output Generation
 - Once optimized, the CNN model is used for real-time prediction of ASL gestures from incoming video frames. The recognized character is shown on a user-friendly graphical interface, and a Text-to-Speech (TTS) engine optionally converts the text into spoken output, facilitating smooth, interactive, and inclusive communication.
- A. Results

VI. RESULTS AND ANALYSIS



Fig.4 Predicted Alphabet A



Fig.5 system ability to append next Character to form a sentence



International Journal for Research in Applied Science & Engineering Technology (IJRASET) ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538 Volume 13 Issue VII July 2025- Available at www.ijraset.com

Text Conversion of	Sign Language
	Character : Sentence : A Suggestions : A As An Ar
¥.	Clear Speak

Fig.6 Captures ambiguous signs in the system



Fig.7 Demonstrate Backspace Functionality



Fig.8 Highlights the system Sentence Building Functionality



International Journal for Research in Applied Science & Engineering Technology (IJRASET) ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538 Volume 13 Issue VII July 2025- Available at www.ijraset.com

Text Conversion of Sign Language



Fig.9 System Forming Complete Sentence

B. Analysis

S. No	Transformation Task	Observed Performance	Existing System Accuracy	Proposed System Accuracy
Ι	Alphabet Recognition	High Accuracy with minimal	91.2%	97.6%
		classification		
II	Real Time Video	Smooth frame capture	88.7%	95.1%
	Capture			
III	Hand Detection Zone	Consistent Detection	90.4%	96.9%

Table 1.Accuracy Score of System

VII. CONCLUSION AND FUTURE SCOPE

This project presents an effective ASL recognition system that uses deep learning to convert hand gestures into both text and speech. Real-time gesture detection through a webcam and classification via a CNN model ensures accurate alphabet recognition. The system includes a user-friendly Tkinter-based GUI that displays predicted characters, forms sentences, and provides voice output. Word suggestions further enhance sentence construction and communication efficiency. The integration of visual and audio feedback makes it an accessible tool for individuals with hearing or speech impairments. Overall, the system offers a reliable and interactive platform for assistive communication. The system can be extended to support dynamic gestures and full sign language sentences using advanced models like RNNs or Transformers. Future versions may also include facial expression analysis to add emotional context to communication. Multilingual output and customizable gesture sets could enhance its versatility across users. Improving the GUI with gesture previews and responsive controls can boost user interaction. Deployment as a mobile or web application with GPU acceleration will improve performance and accessibility. Incorporating evaluation tools like real-time feedback and confusion matrices can further increase accuracy and reliability.

VIII. ACKNOWLEDGEMENT

I sincerely thank Sri G. Narasimham sir for his continual support in helping me carry out this project. Also extend my heartfelt gratitude to the authors of the research papers reviewed in this study.

REFERENCES

- [1] Molchanov, P., Gupta, S., Kim, K., & Kautz, J. (2016). Real-time recognition of hand gestures using 3D convolutional neural networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (pp. 1–7). https://doi.org/10.1109/CVPRW.2015.7301347
- [2] Simon, T., Joo, H., Matthews, I., & Sheikh, Y. (2017). Estimating hand key points from single RGB images using Multiview bootstrapping. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 1145–1153. https://doi.org/10.1109/CVPR.2017.125
- [3] Lugaresi, C., Tang, J., Nash, H., McClanahan, C., Ceze, L., & Taylor, J. (2019). Media Pipe: A versatile framework for building real-time perception pipelines. arXiv preprint arXiv:1906.08172. https://arXiv.org/abs/1906.08172
- [4] Kaur, A., & Singh, M. (2020). Deep learning approach for real-time recognition of American Sign Language alphabets. International Journal of Engineering and Advanced Technology (IJEAT), 9(3), 136–140. https://doi.org/10.35940/ijeat.C4918.029320











45.98



IMPACT FACTOR: 7.129







INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089 🕓 (24*7 Support on Whatsapp)