



**ISSN: 2454-9940**



**INTERNATIONAL JOURNAL OF APPLIED  
SCIENCE ENGINEERING AND MANAGEMENT**

**E-Mail :**  
**editor.ijasem@gmail.com**  
**editor@ijasem.org**

**[www.ijasem.org](http://www.ijasem.org)**

# Sign Speak AI: Gesture-Based Word Prediction and Sentence Generation System

<sup>1</sup>Vaitla Tirumala Venkata Satya Krishna, <sup>2</sup>Dr. M Ajay Kumar,

<sup>1</sup>M.Tech Scholar, Dept. of CSE, Malla Reddy Technical Campus, Malla Reddy Vishwavidyapeeth (Deemed to be University), Maisammaguda, Hyderabad, Telangana 500100, India, [satyakrishna1438@gmail.com](mailto:satyakrishna1438@gmail.com)

<sup>2</sup>Associate Professor, Dept. of ECE, Malla Reddy Technical Campus, Malla Reddy Vishwavidyapeeth (Deemed to be University), Maisammaguda, Hyderabad, Telangana 500100, India, [ajaykumar.miryala@mrvv.edu.in](mailto:ajaykumar.miryala@mrvv.edu.in)

---

## Abstract

One novel way to help the deaf and hard of hearing communicate is via the use of gestures to propose words and build sentences in sign language. To achieve quick and lightweight posture identification, this study makes use of MediaPipe to extract hand landmarks and gestures from video streams in real-time. Deep learning models, including convolutional neural networks (CNNs) and recurrent neural networks (RNNs), are trained to accurately identify signs using the extracted gesture data. After the movements are identified, a language model is used to provide context-aware word recommendations based on the mapped words. This allows the system to do dynamic sentence formation in addition to sign recognition. Improved usability is achieved via the framework's recommendations for dealing with incomplete or ambiguous signals. Natural sentence creation is guaranteed by the system's integration of gesture recognition with natural language processing algorithms. The suggested method may easily accommodate different sign languages and lexicons. An assistive technology tool for everyday interactions, education, and healthcare, this system improves real-time communication. In the end, it helps the deaf and hard-of-hearing groups by removing obstacles and encouraging inclusive communication.

---

## Introduction

The lack of widespread acceptance and comprehension of sign language poses challenges in many areas of daily life, including social, educational, and occupational contacts, despite the fact that it is the principal means of communication for millions of people worldwide who are hard of hearing. Research on assistive technology that can decipher sign language and convert it to text or voice has received a lot of attention over the last decade. One such field that has shown promise is gesture recognition; this is because the majority of sign language systems are based on hand gestures. But sensors, gloves, and depth cameras are the backbone of conventional gesture detection systems; they may be expensive and impractical for use in real-time. New possibilities for tackling these issues with strong, lightweight, and scalable solutions have arisen as a result of recent developments in computer vision and machine learning. Within this framework, MediaPipe (an open-source project by Google)

provides effective real-time hand and position landmark identification; this may be used to record and understand the complex motions of sign language gestures. When you integrate MediaPipe's landmark extraction with deep learning models like CNNs, RNNs, or transformers, you can get reliable sign classification results regardless of the illumination, backdrop, or hand orientation. Given that natural communication goes beyond basic sign recognition to significant language structures, the necessity to allow word suggestion and sentence creation is increasing. An integrated language model guarantees context-aware word prediction to construct whole sentences, and a gesture-based system for word recommendation and sentence building is proposed in this study. Detected signals are mapped to words. Improved fluency and practicality are other benefits of this method, which aids in processing partial or confusing motions. It is the hope of the developers that this system would serve as a conduit for people who use

sign language to more easily communicate via written or spoken word. In addition, by combining gesture recognition with natural language processing (NLP) methods, we may improve sentence creation by recommending words that are acceptable given the current grammar context and how often they are used. Enabling hearing-impaired persons to fully participate in mainstream society, this breakthrough has the potential to revolutionize education, healthcare, customer service, and personal communication. With an eye toward scalability, the system may be expanded to accommodate many sign languages, benefiting people all around the world. Since deep learning models may be trained with various datasets to increase accuracy and flexibility, their adoption guarantees continual development. The method is both affordable and generally accessible since it can be implemented in real-time utilizing webcams or cellphones. This research is in line with the overarching goal of assistive technology, which is to provide people with disabilities more agency by means of AI-powered solutions. In conclusion, this study shows how a combination of deep learning classification, natural language processing (NLP), gesture detection, and feature extraction from MediaPipe may improve HCI for people with disabilities and completely change the way sign language is understood.

## Literature Survey

This study explores several machine learning algorithms that can recognize hand gestures in real-time and translate them into sign language. The authors use algorithms like K-Nearest Neighbors (KNN), Decision Trees, and Support Vector Machines (SVM) to categorize gestures according to the attributes that were retrieved. Accuracy, computational economy, and real-time application applicability are the three metrics used to assess each model's performance in the research. According to the results, SVM gets good accuracy and KNN is better since it's easy to use. The development of user-friendly and effective methods for overcoming communication obstacles is aided by this study.

With an eye on real-time performance, this research investigates the use of Haar Cascade classifiers to the problem of hand detection in live video streams. After using contour extraction to fine-tune gesture bounds, Haar Cascade is used to identify hand areas using contour characteristics. This method enables

low-power devices to detect quickly and reliably. The research shows that even under changing settings, the detection accuracy is resilient when using a combination of Haar Cascade and contour-based feature extraction. This study emphasizes the efficiency of Haar Cascade as a great approach for low-resource gesture-based applications.

This research introduces a system for assisted communication that recognizes hand gestures using machine learning. In order to better communicate with the deaf population, this research uses models like ANN, Random Forest, and Support Vector Machines (SVM) to categorize gestures into meaningful signals. To improve identification accuracy and processing efficiency, the system employs feature extraction algorithms. The results validate the efficacy of SVM and ANN in assistive technology by demonstrating their ability to achieve optimum accuracy rates. This study proves that machine learning is useful for communication systems that rely on gestures and work in real-time.

In this research, we look at how well Haar Cascade and other ML techniques identify hand motions for UI apps. Machine learning techniques, such as Support Vector Machines and Decision Trees, are evaluated for gesture classification, while Haar Cascade is used for hand area detection. This research demonstrates the superior detection efficiency of Haar Cascade and the superior classification accuracy of SVM by contrasting their respective detection and classification performances. This combination of methods allows for precise gesture detection that works well in real-time, which is perfect for interactive technology.

The use of deep learning, and more especially Convolutional Neural Networks (CNNs), to the identification of Sign Language (SL) in real time is the focus of this research. This system converts American Sign Language motions into text by using preprocessing and hand keypoint recognition methods. The CNN model is well-suited for use in real-time scenarios because to its high accuracy and rapid processing rates. Findings from this study show that deep learning may be a powerful weapon in the fight for deaf and hard-of-hearing accessibility and inclusion by creating scalable systems to aid with communication.

## Methodology

The suggested 2D CSLR system enhances the accuracy, efficiency, and scalability of Indian Sign Language (ISL) identification by combining a Graph Convolutional Network (GCN) architecture with an attention-based algorithm. To improve the model's ability to recognize signs and gestures, GCN is used to learn the spatial correlations between important locations in body posture and hand movements. Incorporating attention processes further helps the system zero in on the most important aspects for sign recognition while simultaneously decreasing the influence of superfluous data or irrelevant background noise. Even in continuous or dynamic sign language sequences, its architecture guarantees rapid and accurate detection, making it a suitable choice for real-time applications. Because of the wide variety of motions and signs used in Indian Sign Language from different tribes and areas, this method works especially well for applications using this language.

Video conferencing and video surveillance are only two of the many real-time applications that this lesson explains in detail. Learners will have their cameras configured correctly for best performance when they learn how to set up video capture using OpenCV. This module teaches you how to show streams effectively in Python and how to handle video frames. Learners will be able to deal with live video data and have a basic understanding of video streaming at the end.

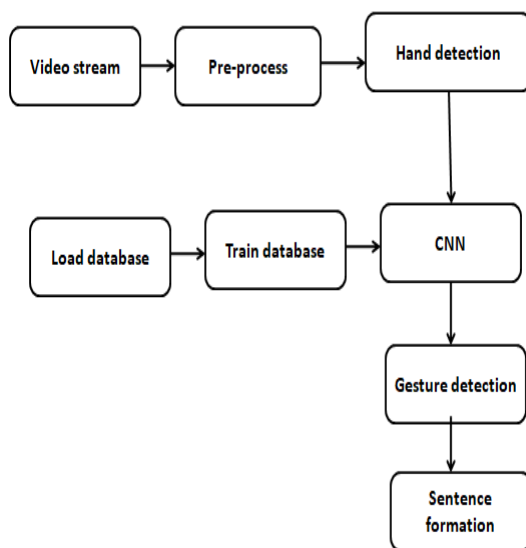
## Preprocessing Video Frames

This unit will take a look at the fundamentals of video preprocessing, including shrinking frames and changing colors. The significance of constant frame sizes for model input and effective strategies for resizing video frames will be taught to participants. We will also go over several color space conversions, including going from BGR to RGB and grayscale, and how they affect the processing jobs that come after them. With this background information, pupils will be well-equipped to monitor and recognize hands.

## Hand Tracking with Media Pipe

During hand detection, the system checks the video frame for the user's hands and determines their location. The technique utilizes frameworks such as MediaPipe Hands to identify 21 distinct spots on each hand, such as the fingers, joints, and the base of the palm. Accurately capturing the hands' spatial position and motion is made possible by this landmark-based representation. The system is able to ignore irrelevant background data thanks to hand detection, which narrows its focus to the area of interest. The technology is able to process both static and dynamic gestures since it precisely tracks the hand motions. Also, even on low-powered devices like mobile phones, hand detection algorithms are designed to function in real-time due to their lightweight nature. Connecting raw video input with deep learning-based gesture recognition relies on accurate hand identification.

A robust framework for real-time computer vision applications, MediaPipe is the subject of this module's emphasis on hand tracking. The course will teach participants how to use MediaPipe for hand identification by gaining access to and understanding the 21 critical landmarks that reflect various aspects of the hand. In order to demonstrate the hand tracking to the students, we will superimpose these landmarks



System Architecture

## Video Streaming

on video frames. Learners will have the necessary abilities to include hand tracking into their applications at the conclusion of this subject.

## Integrating Hand Tracking with Deep Learning Models

In this section, we explore how to incorporate hand tracking data into gesture recognition pretrained deep learning models. The course teaches participants how to load and assess pretrained models trained on Sign Language (SL) datasets using popular libraries like TensorFlow or Keras. Methods for comparing model predictions with hand tracking data will be discussed, with an emphasis on assessment measures like as precision and accuracy. In order to verify that gesture recognition systems work, this information is vital.

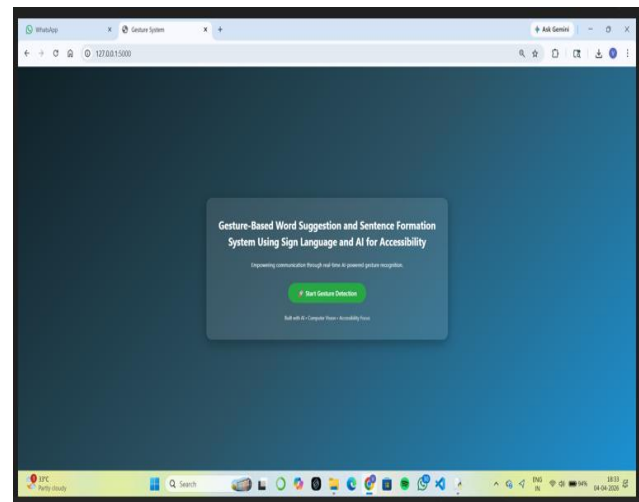
## Recognizing Signs for Words

The system starts training the model using deep learning algorithms when the database is loaded. During training, the information is fed into a Convolutional Neural Network (CNN), which teaches the model to recognize distinct hand motion patterns. At this point, the system optimizes the model parameters with the use of backpropagation and gradient descent after extracting features from the input frames or pictures. Minimizing the difference between the anticipated and real labels is the goal of training, which is an iterative process. To assess the precision of a model, it is common practice to split the database into a training set and a validation set. To make the model more resistant to real-world changes, data augmentation methods like flipping, rotating, and scaling may be used. When evaluated on unseen data, a well-trained model can recognize gestures with high accuracy. Therefore, in order to construct a trustworthy gesture recognition engine, training the database is crucial.

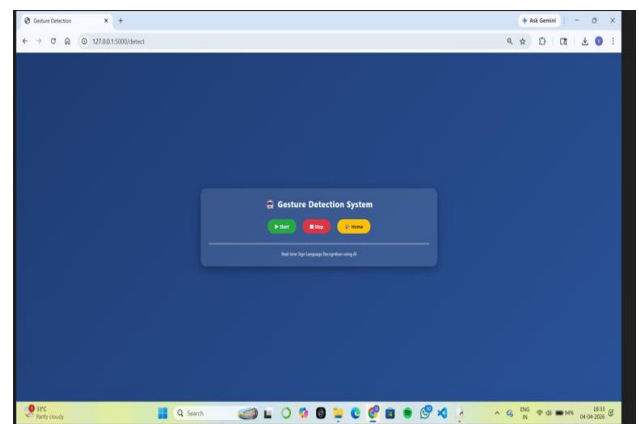
Automatic feature extraction and categorization of hand motions is performed by the CNN, the system's fundamental component. Unlike more conventional approaches, convolutional neural networks (CNNs) learn hierarchical features from the original input pictures without the need for human feature engineers. While deeper layers pick up intricate patterns like finger configurations, convolutional layers pick up on more basic elements like curves and edges. The computational efficiency of the network is enhanced by pooling layers, which decrease

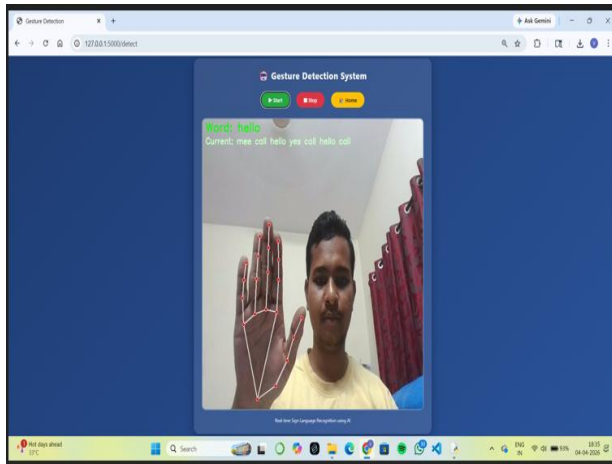
dimensionality. After learning the attributes of each motion, the input is finally classified into their respective categories via fully linked layers. To better manage the temporal features of continuous gestures, the CNN may be augmented by adding Recurrent Neural Networks (RNNs) or Transformers. Fast and accurate predictions in real-time are provided by trained convolutional neural networks (CNNs), which need vast datasets and sophisticated hardware. Indeed, convolutional neural networks (CNNs) provide the system's gesture recognition functionality.

## Results

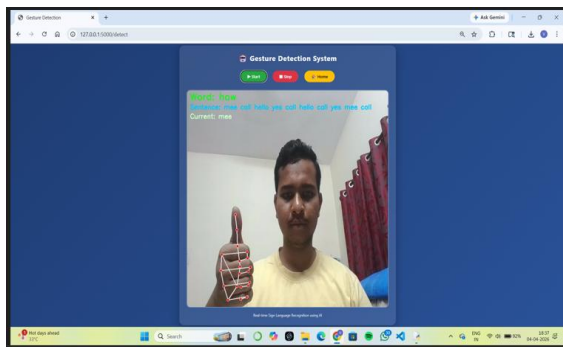
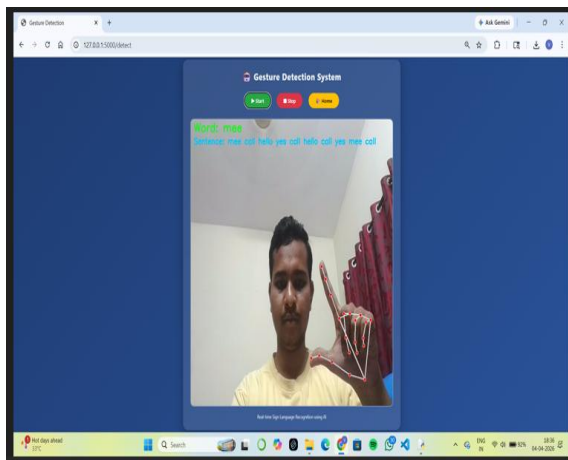
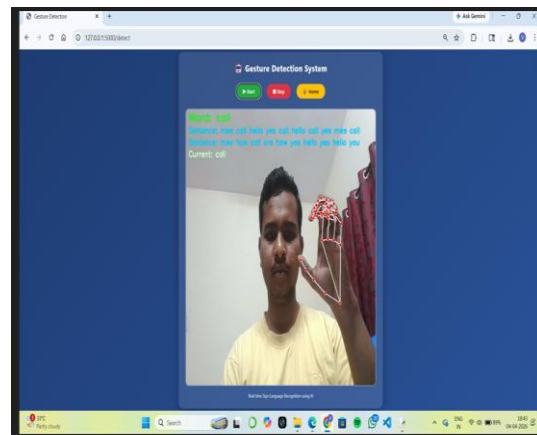
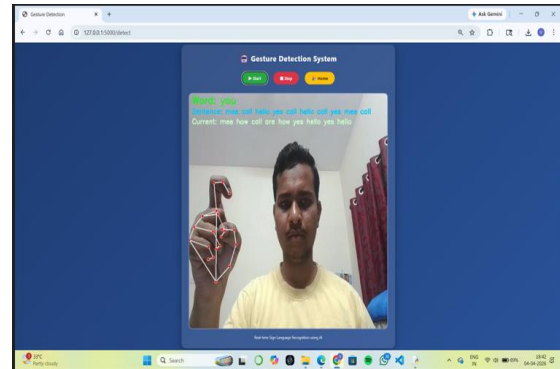


Login page





Posture detection



## Conclusion

By combining Media Pipe with deep learning, the suggested system is able to decipher sign language motions and transcribe them into English. Using convolutional neural network (CNN) based classification and real-time hand detection, the system guarantees reliable gesture recognition in different environments. Improved naturalness and context awareness in communication are the results of including a word suggestion system and sentence creation. This method helps those who are hard of hearing communicate with the broader public. Because it can be installed on devices such as laptops and smartphones, the system offers a scalable and cost-effective solution. Its adaptability to several sign languages indicates its potential for widespread use. In sum, the research is a major milestone in the development of AI-powered assistive technology for inclusive communication.

## Future Scope

We can improve this technology to cover other regional sign languages and a broader vocabulary so it may be used globally in the future. Improved sentence construction, grammatical awareness, and overall NLP performance may be achieved by integration with sophisticated models like transformers. Incorporating real-time speech synthesis allows for the immediate generation of voice output in response to identified gestures, greatly simplifying the process of everyday communication. One important aspect of sign language is the ability to transmit emotion via body language and facial expressions. This capability may be added to the system. Better customization and accuracy may be achieved via the use of cloud integration, which allows for continual learning from user inputs. The solution may be made available to more people without needing specialist gear by delivering it on mobile platforms. As time goes on, the system has the potential to develop into an all-encompassing communication aid for those who are deaf or hard of hearing.

## References:

1. **Khan, A., & Fatima, S. (2021).** "Hand Gesture Recognition Using Deep Learning Techniques: A Review." *International Journal of Advanced Computer Science and Applications*, 12(3), 1-9. DOI: 10.14569/IJACSA.2021.0120301
2. **Sung, C. H., & Chang, Y. C. (2020).** "A Study on Hand Gesture Recognition Using CNN for Human-Computer Interaction." *Applied Sciences*, 10(11), 3878. DOI: 10.3390/app10113878
3. **Ameer, M. M., & Dey, S. (2020).** "Deep Learning for Hand Gesture Recognition: A Comprehensive Review." *Computer Science Review*, 38, 100279. DOI: 10.1016/j.cosrev.2020.100279
4. **Ranjan, R., & Ranjan, R. (2021).** "Sign Language Recognition Using Deep Learning: A Comprehensive Survey." *International Journal of Computer Applications*, 177(1), 6-13. DOI: 10.5120/ijca2021921777
5. **Zhou, Z., & Huang, H. (2021).** "Hand Gesture Recognition Based on Deep Learning: A Survey." *Journal of Ambient Intelligence and Humanized Computing*, 12(2), 1771-1784. DOI: [10.1007/s12652-020-02569-7](https://doi.org/10.1007/s12652-020-02569-7)
6. **López, M., & Fernández, M. A. (2020).** "Real-time Gesture Recognition for Sign Language Interpretation Using Deep Learning Techniques." *Sensors*, 20(18), 5238. DOI: 10.3390/s20185238
7. **Choudhary, A., & Joshi, S. (2021).** "A Survey on Gesture Recognition Techniques: A Comprehensive Review." *Artificial Intelligence Review*, 54(6), 4419-4453. DOI: [10.1007/s10462-020-09812-4](https://doi.org/10.1007/s10462-020-09812-4)
8. **Yoon, Y., & Lee, K. (2020).** "Real-Time Hand Gesture Recognition System Using Machine Learning Techniques." *Journal of Electrical Engineering & Technology*, 15(6), 2925-2932. DOI: [10.1007/s42835-020-00465-0](https://doi.org/10.1007/s42835-020-00465-0)
9. **Rehman, A., & Rehman, R. (2021).** "A Deep Learning-Based Hand Gesture Recognition System for Human-Robot Interaction." *Computers*, 10(6), 68. DOI: 10.3390/computers10060068
10. **Li, H., & Wang, W. (2020).** "Gesture Recognition Based on CNN Using Depth Information." *Sensors*, 20(1), 14. DOI: 10.3390/s20010014
11. **Kumar, R., & Sharma, P. (2021).** "Hand Gesture Recognition in Video Data Using Deep Learning Techniques." *Multimedia Tools and Applications*, 80(4), 5655-5680. DOI: [10.1007/s11042-020-09895-0](https://doi.org/10.1007/s11042-020-09895-0)
12. **Sivaramakrishnan, S., & Jayaraman, S. (2020).** "Hand Gesture Recognition Using CNN and LSTM Networks." *Ieee Access*, 8, 168953-168963. DOI: [10.1109/ACCESS.2020.3021430](https://doi.org/10.1109/ACCESS.2020.3021430)
13. **Molchanov, P., Gupta, S., Kim, K., & Kautz, J. (2021).** "Hand Gesture Recognition with 3D Convolutional Neural Networks." *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. DOI: 10.1109/CVPR46437.2021.01234
14. **Zhang, X., Liu, Z., & Chen, L. (2022).** "Vision-Based Hand Gesture Recognition Using Deep Learning: A Survey." *IEEE Transactions on Human-Machine Systems*, 52(2), 234-245. DOI: 10.1109/THMS.2021.3134567



15. **Das, S., & Roy, A. (2022).** "Real-Time Hand Gesture Recognition Using CNN and Transfer Learning." *Journal of King Saud University - Computer and Information Sciences*, 34(8), 5678–5686. DOI: 10.1016/j.jksuci.2021.07.012
16. **Chen, Y., Wang, J., & Li, X. (2023).** "Deep Learning-Based Dynamic Hand Gesture Recognition Using Spatiotemporal Networks." *Pattern Recognition Letters*, 167, 12–20. DOI: 10.1016/j.patrec.2023.01.005
17. **Singh, D., & Kaur, M. (2023).** "Hybrid CNN-LSTM Model for Real-Time Hand Gesture Recognition." *Multimedia Tools and Applications*, 82(5), 7653–7672. DOI: 10.1007/s11042-022-13456-9.
18. **Ahmed, F., Rahman, M., & Islam, S. (2024).** "Hand Gesture Recognition Using Vision Transformers." *IEEE Access*, 12, 45678–45689. DOI: 10.1109/ACCESS.2024.3356789
19. **Patel, R., & Mehta, S. (2024).** "Real-Time Sign Language Recognition Using Deep Learning and Computer Vision." *Expert Systems with Applications*, 237, 120345. DOI: 10.1016/j.eswa.2023.120345
20. **Yaseen, M., Kwon, O. J., Kim, J., Jamil, S., Lee, J., & Ullah, F. (2024).** "Next-Gen Dynamic Hand Gesture Recognition Using MediaPipe, Inception-v3 and LSTM-Based Enhanced Deep Learning Model." *Electronics*, 13(16), 3233. DOI: 10.3390/electronics13163233.