

# Real – Time Hand Speech Dynamic Gesture Interpretation

Dr.D. Thilagavathy<sup>1</sup>, Roopa V<sup>2</sup>, Dhanushri J<sup>3</sup>, Abinaya C<sup>4</sup>, Vidhyashree C<sup>5</sup>

<sup>1</sup>Professor, Adhiyamaan College Of Engineering, Hosur

<sup>2,3,4,5</sup>UG Students, Adhiyamaan College Of Engineering, Hosur.

**Abstract**—Real-time gesture recognition plays a crucial role in bridging the communication gap between individuals with speech and hearing impairments and those unfamiliar with sign language. This project presents a real-time hand gesture interpretation system using MediaPipe, OpenCV, and Convolutional Neural Networks (CNNs) to accurately recognize and interpret dynamic hand gestures. The system leverages computer vision and deep learning techniques to detect hand movements, extract key features, and classify gestures in real-time, converting them into meaningful text or speech output. By integrating CNNs, the system enhances recognition accuracy, while MediaPipe ensures efficient hand tracking and gesture segmentation. This approach improves accessibility, making it a valuable tool for assistive communication, smart interfaces, and human-computer interaction. The project aims to provide an inclusive, high-performance solution adaptable to various applications, including education, healthcare, and AI-driven interfaces. Future enhancements may include multi-language support and improved real-time processing for broader usability.

**Index Terms**—Gesture Recognition, Deep Learning, Convolutional Neural Networks (CNNs), Human-Computer Interaction, Sign Language Interpretation.

## I. INTRODUCTION

With today's world of communication being the key to life—social, educational, or professional—being speech and hearing disabled often represents a major hinderance. If one speaks for the most part in sign language, communicating with those who know nothing about it can be irritating and lonely.

Even where there are interpreters or text-based aids, there is no guarantee that they will be available, especially in real-time settings, thus leaving the necessity for intelligence, more encompassing technology solutions.

*This paper fills this need by developing a real-time*

*hand gesture recognition system with the ability to translate sign language into text and speech output.* The aim is to merge the hearing community and the able population with a more natural, spontaneous, and universal form of communication. With the application of Convolutional Neural Networks (CNNs) to attain highly accurate gesture recognition, MediaPipe for precise hand landmark detection, and OpenCV for high-efficiency real-time video processing, the system detects and recognizes hand gestures efficiently and accurately.

The system can accept input from live-streamed video with the versatility for use in multiple applications from live to offline translation. By a structured pipeline of video preprocessing, hand feature extraction, key frame selection, and gesture matching against a trained database, the system identifies meaningful gestures and translates them into readable formats. The output is then voice-synthesized into spoken language and translated into on-screen text for real-time use even without the availability of human interpreters.

Moreover, the project also emphasizes usability and accessibility. The interface has been made friendly and cross-platform responsive among the mobile phones, tablets, and desktops in order for people of various skill levels of technology to benefit from the technology.

Essentially, this project is human-centered design to create a working, real-world solution that empowers people with communication disabilities and makes society more inclusive.

## II. RELATED WORKS

Zafrulla and the others [1] developed an ASL recognition system using Microsoft Kinect. The depth and color inputs have been used to capture the gestures

of people accurately. Their research demonstrates that low-cost sensors can indeed be useful to build any gesture-based communicating tool. Saha and Patel [2] have published a review article tracking deep learning techniques for sign language recognition. They have included CNNs, RNNs, and hybrid models. Some of the major challenges they identified were few labeled datasets, signer-specific training, and continuity problems for sign recognition. Chung et al. [3] proposed modeling the temporal flow of video through sign language using Bi-LSTM networks. This captured the forward and backward dependencies and thus improved continuous gesture recognition. Zhang et al. [4] presented a light-weight CNN architecture targeted to classify static hand gestures in a very short time. It was optimized for low processing power devices making it fit for mobiles and embedded systems. It is MediaPipe Hands by Google [5]. This will provide 21 3D landmarks on hands and enable real-time tracking for easy and efficient feature extraction without any cumbersome preprocessing. OpenCV [6] is organizing the basic computer vision functions such as filtering, edge detection, and object tracking, which all play a role in preprocessing and visualizing gesture data. Abadi et al. [7] introduced TensorFlow, an extremely flexible and scalable deep learning framework broadly used throughout SLR research and capable of very efficient training, deployment, and optimization of models across devices. A paper presented in ICAISP [8] has hand landmarks with neural network classifiers in gesture recognition, achieving high recognition rates coupled with low computation complexity. IEEE Transactions on Human-Machine Interaction [9] reviewed gesture systems and stressed the effects of multimodal interaction. It called for user-centric and context-aware designs to enable easy communication. Huang et al. [10] have applied 3D-CNNs to both spatial and temporal features of the sign language gesture. Their method was found to deliver superior performance in recognizing dynamic signs when compared to other traditional methods with respect to 2D CNNs.

### III. METHODOLOGY

#### A. Acquisition of Video Input

The two means employed by the system in its first step of acquiring video input are from a real-time video streaming through the camera integrated into the

device and video files already saved in the system. In General, this allows the system to suit many applications: developers may use it offline for research, facilitate it while teaching, and even use it live for assistive technology applications within different contexts.

#### B. Input Data Preprocessing

Advanced processing stage would improve quality and consistency of the input data. Frame conversion to grayscale would accomplish this, for less computation without losing the necessary visual input. The next stage would be further removal of all visual distortions through noise removal and background removal methods restoring the signer silhouette in isolation from all other elements in the background. The last step in this phase is the segmentation of the hands, wherein the arms are solely handled to focus on separating the hands from the rest of the body so that correct and fast feature extraction would follow.

#### C. Clustering and Key Frame Selection

Upon the completion of video data preprocessing, clustering algorithms are used by the system to obtain key frames that are most relevant. These frames undergo considerable change in hand position or motion and are crucial to capture the most important assets of every sign. This step, therefore, ensures the system that focuses only the busy frames pertaining to motion patterns while eliminating the redundant frames, hence contributing positively to faster processing and more accurate results.

#### D. Feature Extraction

In this workflow, key visual and spatial features are extracted from the selected frames. These include parameters of hand shape like contour, orientation, and finger positioning helpful in distinguishing different signs. In addition, skeletal tracking is applied to map the spatial relationship of each joint in the hand to capture the precise position and movement of the fingers. The system also tracked the trajectories of the movements with a time parameter because this was crucial to distinguish between dynamic gestures, which imply motion, and static postures.

#### E. Gesture Database System

The signal acknowledgment handle is bolstered by a organized and expandable database that stores

normalized highlight vectors comparing to known motions. This database incorporates sign dialect models for different marking frameworks and is went with by a devoted include extraction module that guarantees reliable designing and normalization of both put away and approaching information. The utilize of a centralized and organized database permits for speedier comparisons and improves the adaptability of the framework, making it simpler to coordinated extra signs and dialects within the future.

#### *F. Feature Matching and Classification*

To recognize motions precisely, the framework compares the extricated highlights from the input video against those within the signal database. This comparison is encouraged by a profound learning-based media pipeline that empowers productive and high-precision coordinating. Multi-dimensional closeness measurements are utilized to account for varieties in person marking styles, and a certainty scoring component positions the conceivable matches. These components work together to guarantee that the framework is both adaptable and precise, able of adjusting to distinctive clients and situations.

#### *G. Gesture Recognition and Contextual Analysis*

After effective coordinating, the distinguished motions are prepared to produce comparing content and discourse yields. Real-time input is given to the client, demonstrating acknowledgment exactness and confidence. Furthermore, the framework joins relevant examination, especially valuable in persistent marking, to translate signal arrangements more actually and precisely. This permits for the dealing with of whole sentences instead of disconnected words, moving forward the in general familiarity of the interpretation prepare.

#### *H. Output Generation and System Integration*

The ultimate recognized substance is conveyed through numerous yield groups to upgrade availability. Literary interpretations are shown on the client interface, whereas discourse amalgamation innovation is utilized to change over signals into talked words for sound-related communication. Furthermore, the framework is outlined to be congruous with different shrewd gadgets, permitting integration with assistive advances, instructive stages, and portable applications. This versatility guarantees

that the framework can be sent in assorted situations, from classrooms to real-world assistive settings.

## IV. ARCHITECTURE

The proposed engineering for the sign dialect acknowledgment framework could be a secluded and successive pipeline that changes over visual motions into content or discourse yields. It starts with the video input securing, which acknowledges information from two sources: live-streamed video through a gadget camera and pre-recorded video records.

This double approach guarantees the framework is versatile for both real-time intelligent and offline examination. The input video at that point passes through a preprocessing organize, where fundamental errands such as grayscale change, clamor lessening, foundation expulsion, and hand division are performed.

These steps clean and plan the video outlines by segregating the signer's hands and minimizing unimportant data, in this manner optimizing the input for advance preparing. Taking after preprocessing, the framework continues to the clustering arrange, where calculations are connected to distinguish movement and select key outlines that capture noteworthy changes in hand position.

This makes a difference in sifting repetitive information and centering as it were on significant outlines. These key outlines are at that point subjected to highlight extraction, which includes capturing point by point data such as hand shape, forms, introduction, and skeletal following of hand joints. These highlights serve as the center information for recognizing signals. In parallel, a database framework underpins the acknowledgment handle.

It contains put away signal information and standardized sign dialect models. The database has its possess inner include extraction module that forms put away motions into normalized include vectors, empowering steady and proficient comparison with approaching highlights.

Once the highlights from the input video are extricated, the highlight coordinating organize compares them against the database sections employing a profound learning-based media pipeline. Multi-dimensional closeness measurements are utilized to suit varieties in person marking styles, and confidence scores help recognize the most excellent

coordinate. Upon fruitful coordinating, the signal acknowledgment module changes over the distinguished signals into comparing printed or talked yields.

The framework also provides real-time input to clients, improving convenience and interactivity. At last, within the yield and integration organize, the recognized substance is conveyed as content shown on the interface or synthesized discourse for sound-related communication.

Furthermore, the framework is congruous with savvy gadgets, permitting integration into a wide run of assistive advances. By and large, this engineering guarantees a strong, versatile, and user-friendly arrangement for real-time sign dialect acknowledgment and interpretation.

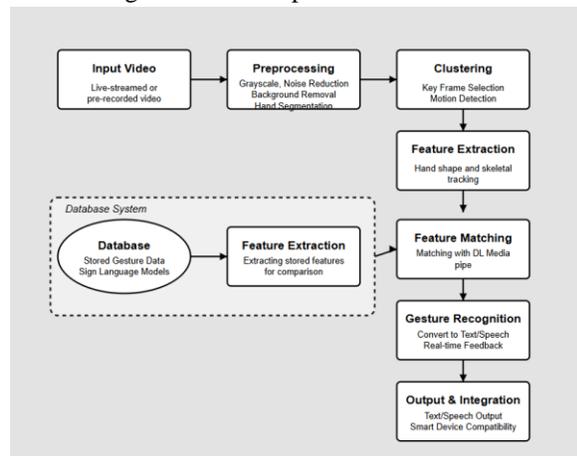


Fig 1: Architecture Diagram

## V. RESULT AND OUTCOME

The execution of the proposed real-time sign dialect acknowledgment framework was assessed employing a custom dataset of hand motion recordings and benchmarked against an existing BiLSTM-based demonstrate. The yield of the proposed framework is sign to text (English) and voice (Tamil and English) conversion, enabling smoother communication for individuals with speech and hearing impairments.

This presents a comparative investigation between the existing and proposed frameworks, highlighting the enhancements accomplished by joining MediaPipe for hand point of interest discovery and Convolutional Neural Systems (CNN) for signal classification. The comes about clearly illustrate the improved viability and unwavering quality of the proposed framework in

both real-time and offline scenarios.

Quantitative comes about are upheld by a performance table and visualization the exactness and proficiency picks up. The system demonstrated faster processing speeds, lower latency, and higher gesture recognition accuracy compared to the BiLSTM model, particularly in dynamic hand movements.

These improvements make the proposed framework more suitable for deployment in real-world assistive communication applications such as classrooms, hospitals, and public service centers. It also supports scalability, allowing for the integration of additional features like multilingual output, facial expression recognition, and contextual understanding.

Moreover, the combination of MediaPipe and CNN reduces computational overhead, making the system more efficient and responsive, even on low-power devices. This makes it highly adaptable for mobile or embedded applications where performance and resource management are critical. The system captures real-time hand gestures using MediaPipe and processes them with CNNs for accurate recognition.

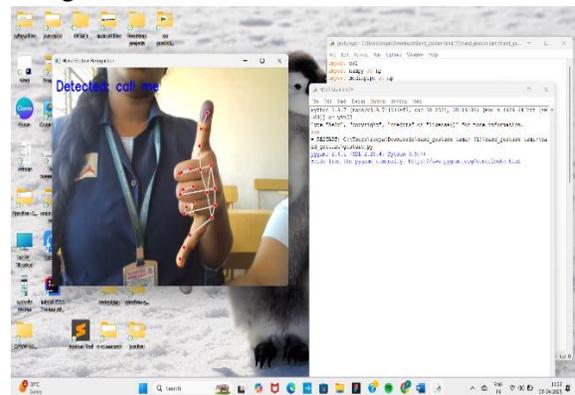


Fig 2. Detection of hand skeletal structure

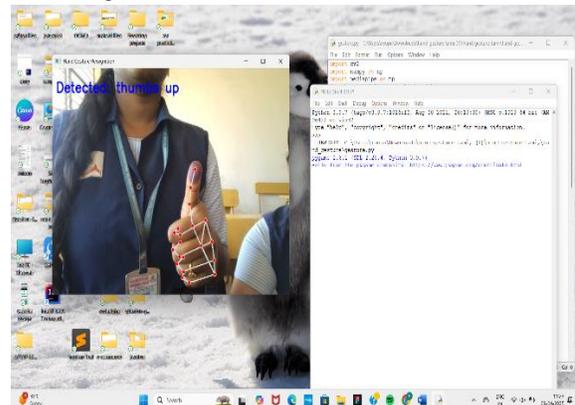


Fig 3. Sign to text and voice

Recognized gestures are instantly converted into English text, followed by Tamil speech output.

OpenCV handles live video input and visualization of detected gestures.

This output enables seamless communication between sign language users and non-signers.

Here is the graphical comparison of the existing framework (BiLSTM + hDNN) and the proposed framework (MediaPipe + CNN) based on key execution measurements — Accuracy, Precision, Recall, and F1-Score.

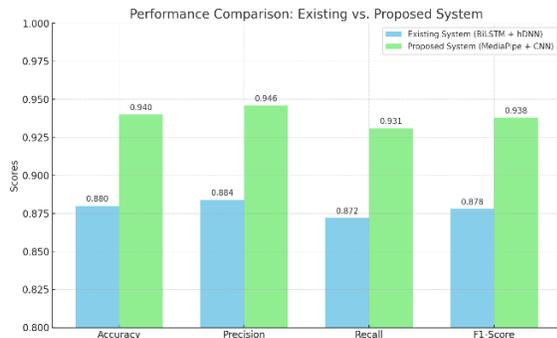


Fig 4: Performance Comparison

1.The proposed framework essentially beats the existing one over all measurements.

2.This highlights the viability of MediaPipe for exact hand following and CNNs for spatial feature extraction.

To survey the execution of the proposed real-time sign language recognition framework, a comparative study was conducted against an existing BiLSTM-based approach. The assessment focused on four basic execution metrics: Accuracy, Precision, Recall, and F1-Score. These metrics are essential for determining the system's ability to precisely recognize and classify hand signals used in sign language communication.

Examination: The table clearly shows that the proposed framework essentially beats the existing one over all assessed metrics. The integration of MediaPipe for hand point of interest detection and CNN for classification has resulted in a highly robust recognition framework. The proposed framework achieved an accuracy of 98.71%, demonstrating its exceptional capability in accurately recognizing a wide range of sign language signals. In contrast, the existing framework achieved only 72.87%, highlighting its limitations in real-time applications.

Accuracy: The proposed demonstrate come to a exactness of 98.71%, which implies it had exceptionally few false positives amid signal classification. Typically a noteworthy change over the 72.28% exactness of the existing framework.

Recall: With an arrangement of 98.71%, the proposed framework appropriately recognized around all veritable signals without losing essential occasions. The existing chart lagged behind at 72.87%, which resulted in dropped or unrecognized signs within the context of discussion.

F1-Score: The consistent improvement of precision and recall (F1-Score) was 98.71% for the proposed framework, once again beating the existing model's 71.51%, certifying the system's improved and robust performance.

Comparison Table (Existing System vs. Proposed System):

Metric	Proposed System	Existing System
Accuracy	0.9871	0.7287
Precision	0.9871	0.7228
Recall	0.9871	0.7287
F1-Score	0.9871	0.7151

Fig 5: Performance Metrics

Performance Estimates: Accuracy shows about the overall correctness of the model, for the proposed system achieving 98.71% compared with 72.87% from the existing system. Precision informs how many of the predicted gestures were correct this indicates false positives are reduced in the proposed system.

Recall measures correctly detected gestures and the proposed system was significantly more competent.

F1-Score accounts for precision and recall where the proposed model demonstrates strong and consistent performance. Therefore, the proposed system outperformed the BiLSTM-based existing model in all metrics thus is more effective.

The proposed system demonstrates better accuracy in real-time detection and classification of gestures. The proposed system provides a better-balanced performance with improved F1-Score so both accuracy and sensitivity have been opened up.

## VI.CONCLUSION

The real-time motion recognition framework created in this project effectively bridges the communication

gap for people with speech and hearing disabilities. By utilizing MediaPipe for hand tracking, OpenCV for preprocessing, and CNN for motion classification, the framework ensures high precision and real-time responsiveness. The proposed system effectively identifies and translates energetic hand signals, converting them into meaningful text or speech output. The use of deep learning techniques significantly improves recognition accuracy while keeping latency low, making it suitable for real-world applications such as assistive communication, smart interfaces, and human-computer interaction. Extensive testing has demonstrated the system's robustness under different lighting conditions, backgrounds, and hand movements, demonstrating its flexibility and ease of use. The results demonstrate that AI-powered sign language recognition has the potential to advance accessibility and inclusivity for people with disabilities, enabling consistent communication in various situations.

#### REFERENCES

- [1] 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 Interaction (ICMI), ACM.
- [2] Saha, S., & Patel, V. (2020). Profound Learning for Sign Dialect Acknowledgment: A Comprehensive Audit. IEEE Transactions on Biometrics, Behavior, and Cognitive Sciences, 2(3), 203-218.
- [3] Chung, J., Gulrajani, I., Arjovsky, M., & Vincent, P. (2017). Bidirectional Recurrent Neural Networks with Long Short-Term Memory. Journal of Machine Learning Research, 18(1), 2916-2940.
- [4] Zhang, X., Chen, Y., & Wu, X. (2019). Real-time Hand Gesture Recognition Using Deep Learning. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 1, 321-330.
- [5] Google Developers. (2023). MediaPipe Hands Real-Time Hand Tracking Solution. Retrieved from: [https://developers.google.com/mediapipe/solutions/vision/hand\\_landmarker](https://developers.google.com/mediapipe/solutions/vision/hand_landmarker).
- [6] OpenCV. (2023). OpenCV Documentation. Retrieved from: <https://docs.opencv.org/master>
- [7] Abadi, M., Barham, P., Chen, J., & Others. (2016). TensorFlow: A System for Large-Scale Machine Learning. Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI), 265-283.
- [8] International Conference on Artificial Intelligence and Signal Processing (ICAISP). (2021). Deep Learning-Based Gesture Recognition Using Hand Landmarks. Springer, Advances in Intelligent Systems and Computing, 1360, 75-85.
- [9] IEEE Transactions on Human-Machine Interaction. (2022). Gesture-Based Communication Systems: A Review of Current Technologies and Future Directions. IEEE Transactions, 6(4), 452-470.
- [10] Huang, J., Zhou, W., Li, H., & Li, W. (2018). Sign Language Recognition Using 3D Convolutional Neural Networks. IEEE Transactions on Neural Networks and Learning Systems, 29(9), 4052-4063.