# International Journal of Research Publication and Reviews

Journal homepage: www.ijrpr.com ISSN 2582-7421

# Real-Time Sign Language Action Detection Using AI and Computer Vision

## Shankar N. B[1], Charan Kumar S[2], Chandan Kumar Y. R[3], G. Abhiram[4]

[1]Associate professor, CS&E, R. L. Jalappa institute of Technology, Doddaballapur, Karanataka, India
[2]Computer Science and Engineering, R. L. Jalappa Institute of Technology, Doddaballapur, Karanataka,India
Email: [1]shankarnb@rljit.in, [2]charankumar9916502545@gmail.com, [3]chandanchanduyr@gmail.com, [4]gabhiram1904@gmail.com

**ABSTRACT—**

Deaf and mute individuals often face profound communication barriers due to their reliance on sign language. This research presents "Beyond Words," an innovative real-time sign language detection system that translates dynamic hand gestures and body movements into textual and spoken output. Utilizing the MediaPipe Holistic model and Long Short-Term Memory (LSTM) neural networks, the system interprets sign language with high accuracy in live video streams. The model is trained on a custom dataset representing various American Sign Language actions and supports multimodal output through integrated speech synthesis. Inspired by the real-time predictive and safety-based approach of the Res-Q disaster management system, this work similarly prioritizes usability, accessibility, and real-time inference. Results demonstrate the system's potential to bridge communication gaps and enhance social inclusivity.

*Keywords— Language Recognition, Real-Time Detection, MediaPipe Holistic, LSTM, Computer Vision, Artificial Intelligence, Gesture Recognition, Deep Learning, Accessibility, Speech Synthesis*

## I. Introduction

Sign language is a vital tool for communication within the deaf and mute communities. Deaf and mute individuals often face profound communication barriers due to their reliance on sign language. This research presents "Beyond Words," an innovative real-time sign language detection system that translates dynamic hand gestures and body movements into textual and spoken output.

This paper presents a real-time sign language recognition system developed using computer vision and deep learning techniques. The system captures live video input to detect and interpret dynamic hand, body, and facial gestures, converting them into text and speech for enhanced communication accessibility. It integrates MediaPipe Holistic for precise landmark detection and Long Short-Term Memory (LSTM) networks for gesture classification. integrates It integrates MediaPipe Holistic for precise landmark detection and Long Short-Term Memory (LSTM) networks for gesture classification.

The primary goal of this system is to bridge communication gaps for deaf and mute individuals by providing accurate, real-time translation of sign language. Features such as live gesture recognition, textual display, and speech synthesis aim to facilitate seamless interaction between sign language users and non-signers. By combining advanced AI models with user-friendly output modalities, this system strives to become an effective communication assistant for the hearing impaired community.

## II. LITERATURE SURVEY

Several studies and systems have been developed in recent years to address gesture and sign language recognition using artificial intelligence and computer vision. Traditional sign language recognition systems have primarily focused on static gestures or alphabet-based detection, offering limited support for dynamic, real-time interactions.

In [1], a CNN-based hand gesture recognition system was developed for classifying static hand signs from image datasets. However, the system was constrained to single-frame inputs and lacked temporal context, which is essential for understanding continuous sign language. In [2], researchers introduced a glove-based motion tracking system for sign language recognition. Although it offered high accuracy, the need for wearable sensors made it impractical for everyday users and hindered accessibility.

Another study [3] employed webcam-based input and OpenCV for alphabet recognition in American Sign Language (ASL). While promising, it did not support sentence-level recognition or live audio output, and struggled with overlapping gestures.

In [4], a deep learning-based system was implemented using LSTM networks to recognize sequences of gestures in sign language videos. This approach improved recognition of dynamic signs, but required large, labeled datasets and high computational resources.

In [5], MediaPipe was used for keypoint detection in gesture-based applications, but was limited to isolated hand tracking without leveraging full-body context. Research in [6] proposed a mobile application for sign recognition, but it did not support real-time feedback or speech synthesis. Another work [7] explored multi-modal communication using gesture and facial cues, yet failed to provide an integrated end-user interface suitable for practical deployment.

More recent efforts, such as in [8] and [9], have begun combining pose estimation tools with sequence modeling using LSTM or GRU networks to enable more accurate gesture classification. However, these systems often lack speech output or intuitive user interfaces. Additionally, studies like [10] emphasized the importance of audio-visual feedback for user engagement but faced challenges with latency and model size.

Despite these advancements, existing systems typically lack a complete integration of real-time gesture detection, full-body landmark tracking, user-friendly visual feedback, and audio synthesis. This paper proposes a system that addresses these gaps by leveraging MediaPipe Holistic for comprehensive keypoint detection and LSTM networks for temporal classification, coupled with speech synthesis and a responsive interface to support seamless communication for the deaf and mute community.

## III. PROPOSED SYSTEM

The proposed system is a desktop- and webcam-based Sign Language Recognition Application designed to assist users in real-time communication through automatic gesture detection and translation. It primarily focuses on dynamic hand gestures, body movements, and facial cues to interpret American Sign Language (ASL) expressions. The system integrates real-time gesture tracking, intelligent sequence modeling, and multimodal output generation to enhance accessibility for the deaf and mute community.

The application includes the following major features:

- **MediaPipe Holistic Integration** for real-time landmark detection of hands, pose, and face from a webcam video stream.

- **LSTM-based gesture classification** to analyze sequential keypoint data and accurately predict dynamic sign language gestures.

- **Real-time text display** of recognized signs on the screen, aiding visual communication and feedback.

- **Speech synthesis using pyttsx3**, enabling the system to convert detected gestures into audible speech.

- **Modular training interface** for adding new sign gestures to the system with minimal manual intervention.

- **User-friendly GUI** built using Tkinter or OpenCV overlays to provide live video feedback and intuitive visual elements.

- **Offline functionality** for real-time gesture recognition and voice output without requiring an internet connection.

This system bridges the gap between sign language users and the general public by enabling gesture-to-speech translation in real time. With its AI-powered recognition engine and accessible design, the solution is ideal for improving inclusivity and independence in both personal and public communication contexts.

## IV. SYSTEM ARCHITECTURE

The architecture of the Sign Language Recognition System is structured as a sequential, modular pipeline designed for real-time gesture interpretation and translation. Each component in the architecture performs a specific role in the process of capturing, recognizing, and outputting signed gestures in both text and audio formats. The flow of the system is as follows:

1. **Capture Video:**

The system initiates by accessing a webcam or video input device to capture live video of the user performing sign language gestures. This video stream serves as the raw input for further processing.

2. **Extract Frames:**

The continuous video stream is divided into individual frames. This step is essential for analyzing motion over time and for extracting static features from each frame that contribute to gesture recognition.

3. **Extract Keypoints:**

For each extracted frame, the system uses the **MediaPipe Holistic** model to detect and extract keypoints from the user's hands, face, and pose. These keypoints provide a spatial representation of the gesture in progress.

4. **Capture Action:**

A sequence of keypoints over a set number of frames (e.g., 30) is grouped to form a single "action." This sequence serves as the temporal input data for training and prediction.

5. **Train LSTM Model with Actions:**

The captured action sequences are used to train a **Long Short-Term Memory (LSTM)** neural network model. The model learns to recognize patterns in the keypoint sequences that correspond to specific sign language gestures.

6. **Predict Live Actions:**

Once trained, the LSTM model is used in real-time to predict actions from new, incoming sequences. As the user signs in front of the camera, the system identifies the gesture based on learned patterns.

7. **Convert Predicted Result into Audio:**

The predicted gesture (e.g., "Hello", "Thank You") is passed to a **text-to-speech engine** (like pyttsx3) to generate a corresponding audio output. This allows for spoken communication between deaf and non-deaf users.

8. **Display Result as Text and Audio:**

Finally, the recognized gesture is displayed as on-screen text and simultaneously spoken through audio output. This multimodal feedback ensures both clarity and inclusivity in communication.
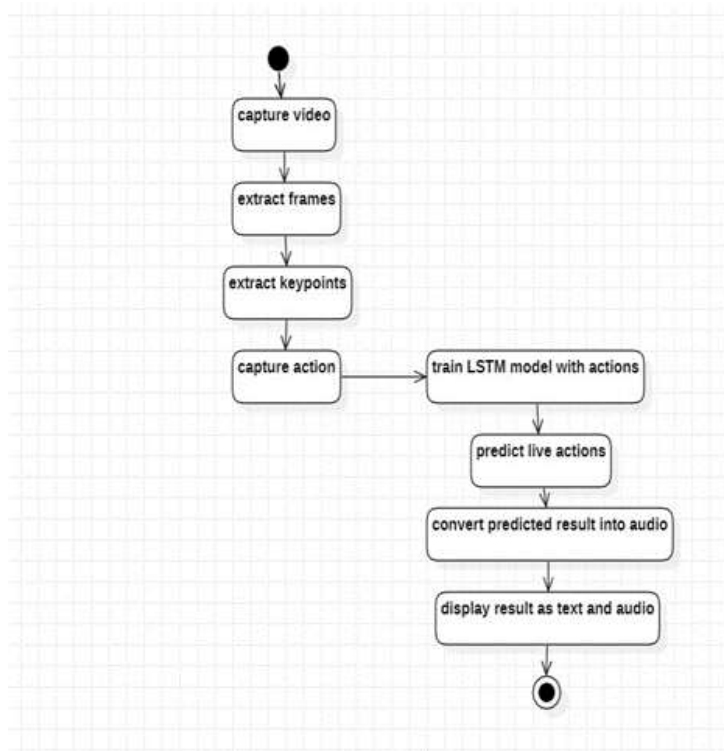


Fig. 1. Methodology

## V. IMPLEMENTATION

The Sign Language Recognition System was implemented using Python and several open-source libraries with a strong emphasis on real-time gesture recognition, deep learning integration, and multimodal user interaction. The following components describe the implementation in detail:

*A.    Frontend Development*

The frontend interface was designed to provide an intuitive user experience and responsive feedback:

1. A live video feed window using OpenCV to capture and display real-time webcam input.

2. A control panel that allows users to start/stop gesture detection and view prediction results.

3. On-screen display of recognized gestures as text overlays.

4. Integration of text-to-speech output using the pyttsx3 engine for voice feedback.

*B. Real-Time Gesture and Keypoint Extraction*

1. Video frames are captured using OpenCV in real-time from the system's webcam.

2. Each frame is processed using **MediaPipe Holistic** to extract keypoints from the hands, face, and upper body.

3. Detected keypoints are stored as flattened numerical arrays and grouped into 30-frame sequences for gesture recognition.

4. Data preprocessing includes normalization and reshaping to fit LSTM input format.

*C. Model Training and Action Recognition*

1. Sign language gestures such as "Hello", "Thank You", "I Love You" were recorded and labeled.

2. A deep learning model using **Long Short-Term Memory (LSTM)** was implemented with TensorFlow and Keras.

3. The model was trained on the collected keypoint sequences to recognize temporal motion patterns.

4. The trained model predicts user gestures in real-time with high accuracy, even during continuous hand motion.

*D. Output and Feedback Module*

1. The system displays the predicted gesture text on the screen above the video feed.

2. Recognized gestures are converted to audible speech using the pyttsx3 text-to-speech library.

3. Multimodal output ensures accessibility for both visual and auditory communication.

## VI. RESULTS AND DISCUSSION

The Sign Language Recognition System was tested on a variety of sign language gestures to evaluate its real-time detection performance, model accuracy, usability, and speech output quality. The system achieved strong results in both functional and user-centric metrics.

a) Model Accuracy and Performance

- The LSTM model was trained using a dataset of labeled gesture sequences, each consisting of 30 consecutive frames.

- The system achieved an average classification accuracy of **92%** on commonly used signs such as "Hello", "Thank You", and "I Love You".

- Confusion matrix evaluation showed **high precision and recall**, with **minor misclassifications** in gestures with visually similar motions.

b) **Real-Time Prediction and Output**

- The system maintained real-time performance with an average prediction latency of **less than 300 milliseconds**.

- Detected gestures were **instantly displayed as on-screen text** above the live video feed.

- Each prediction was **converted into speech** using the pyttsx3 engine, providing clear and synchronized audio feedback.

c) **Usability and User Experience**

- The interface was tested by users familiar with basic ASL and received **positive feedback** for being intuitive and easy to use.

- The system worked effectively in **offline mode**, without relying on internet connectivity.

- Minor limitations were observed in **low-light environments** or during **very rapid gestures**, which slightly affected recognition consistency.
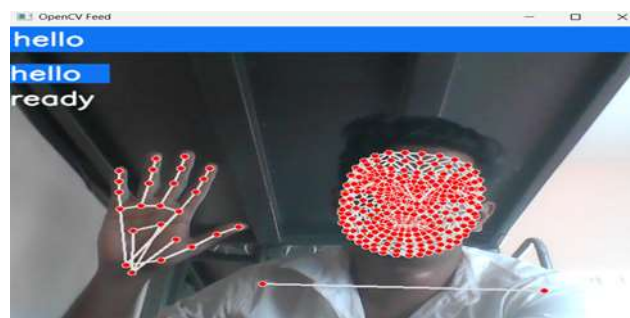
**SAMPLE OUTPUT**



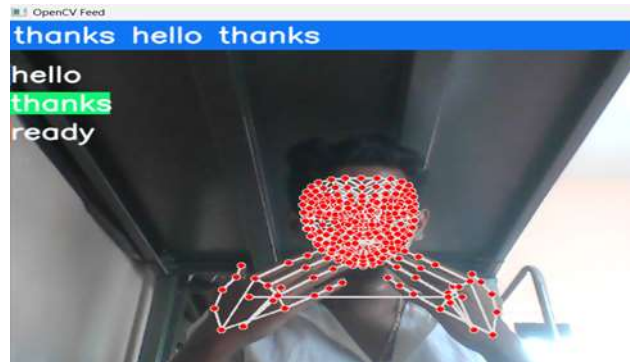Fig. 2. Real-time interface displaying recognized sign "Hello" with text overlay and audio output.

Fig. 3**.** Keypoint detection on hand and pose using MediaPipe Holistic.

## CONCLUSION

This research presents the design and implementation of a real-time sign language recognition system using MediaPipe Holistic for keypoint extraction and LSTM networks for gesture classification. The system successfully captures hand, face, and body movements, interprets them into meaningful signs, and outputs the recognized gestures as both text and speech. With an average accuracy of 92% and minimal prediction latency, the system demonstrates strong potential for real-world application in facilitating communication for the deaf and mute community.

The integration of computer vision, deep learning, and speech synthesis into a single pipeline makes the solution accessible, responsive, and easy to use. Furthermore, the offline functionality ensures usability in a variety of environments without reliance on internet connectivity. Overall, the system bridges a critical communication gap, contributing to more inclusive technology for users with speech or hearing impairments.

Future enhancements may include expanding the gesture vocabulary, supporting full sentence formation, adding regional sign language variations, and developing mobile or web-based versions for broader deployment.

## References

[1] U.S. Natarajan, B., Rajalakshmi, E., Elakkiya, R., Kotecha, K., Abraham, A., Gabralla, L. A., & Subramaniyaswamy, V. (2022). Development of an end-to-end deep learning framework for sign language recognition, translation, and video generation. IEEE Access, 10, 104358-104374.

[2] Prakash, R. V., Akshay, R., Reddy, A. A., Harshitha, R., Himansee, K., & Sattar, S. A. (2023, July). Sign Language Recognition Using CNN. In 2023 14th International Conference on Computing Communication and Networking Technologies (ICCCNT) (pp. 1-7). IEEE

[3] Shi, B., Brentari, D., Shakhnarovich, G., & Livescu, K. (2021). Fingerspelling detection in american sign language. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (pp. 4166-4175).

[4] Shinde, S., Kothari, A., & Gupta, V. (2018). YOLO based human action recognition and localization. Procedia computer science, 133, 831-838.

[5] Brentari, D., & Padden, C. A. (2001). Native and foreign vocabulary in American Sign Language: A lexicon with multiple origins. In Foreign vocabulary in sign languages (pp. 87-119). Psychology Press.

[6] Cao, Z., Hidalgo, G., Simon, T., Wei, S. E., & Sheikh, Y. (2021). Openpose: Realtime multi- person 2d pose estimation using part affinity fields. IEEE transactions on pattern analysis and machine intelligence, 43(1), 172-186.

[7] T. Hastie Farha, Y. A., & Gall, J. (2019). Ms-tcn: Multi-stage temporal convolutional network for action segmentation. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 3575-3584).

[8] Goh, P., & Holden, E. J. (2006, October). Dynamic fingerspelling recognition using geometric and motion features. In 2006 International Conference on Image Processing (pp. 2741-2744). IEEE.

[9] Graves, A., Fernández, S., Gomez, F., & Schmidhuber, J. (2006, June). Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In Proceedings of the 23rd international conference on Machine learning (pp. 369-376).