

# Comparative Analysis of Machine Learning Models for Hand Gesture Recognition

Prof. Manisha Mathur

*Assistant Professor, Computer Engineering*  
*Vivekanand Education Society's Institute of Technology*  
Mumbai, India  
manisha.mathur@ves.ac.in

Jiya Gangwani

*Student, Computer Engineering*  
*Vivekanand Education Society's Institute of Technology*  
Mumbai, India  
2021.jiya.gangwani@ves.ac.in

Chirag Santwani

*Student, Computer Engineering*  
*Vivekanand Education Society's Institute of Technology*  
Mumbai, India  
d2021.chirag.santwani@ves.ac.in

Nikhil Dhanwani

*Student, Computer Engineering*  
*Vivekanand Education Society's Institute of Technology*  
Mumbai, India  
d2021.nikhil.dhanwani@ves.ac.in

Soham Panjabi

*Student, Computer Engineering*  
*Vivekanand Education Society's Institute of Technology*  
Mumbai, India  
2021.soham.panjabi@ves.ac.in

**Abstract**—Silent Cue is a hand gesture recognition system designed to facilitate non-verbal communication for deaf and non-verbal individuals. The system leverages OpenCV and MediaPipe for real-time hand tracking and gesture recognition, enabling users to communicate effectively using sign language. By employing an adaptive machine learning model, the system can recognize a wide range of gestures, facilitating hands-free interaction with various devices such as volume control, scrolling, and mouse cursor movement. This paper delves into the system's design, the methodology behind its development, and its potential applications, offering valuable insights into how technology can empower individuals with hearing impairments and improve communication.

**Index Terms**—Hand Gesture Recognition, Non-Verbal Communication, Deaf Communication, Machine Learning, Sign Language, Adaptive Systems, OpenCV, MediaPipe, Silent Cue, Assistive Technology

## I. INTRODUCTION

Communication is a fundamental aspect of human interaction, enabling individuals to share thoughts, ideas, and emotions. For individuals with hearing impairments or non-verbal communication needs, traditional means of communication, such as spoken language, pose significant challenges. Sign language serves as a crucial tool for deaf individuals, enabling effective communication within their community. However, a lack of widespread understanding and adoption of sign language can lead to communication barriers with individuals who do not use it.

The development of hand gesture recognition systems offers a promising solution to these challenges, as they facilitate non-verbal communication and interaction with digital devices.

Silent Cue, a hand gesture recognition system, was designed specifically for deaf and non-verbal individuals to help bridge this gap. Utilizing real-time hand tracking and gesture recognition technologies, Silent Cue enables users to interact with devices and communicate effectively without the need for speech or traditional sign language interpreters.

This paper presents an in-depth overview of the Silent Cue system, exploring its architecture, design considerations, and real-world applications. We examine the combination of OpenCV and MediaPipe for real-time hand tracking, and the machine learning models that enable gesture recognition. By incorporating these technologies, the system aims to offer a seamless and intuitive solution for non-verbal communication, enhancing accessibility and empowering individuals with hearing and speech impairments. The paper also discusses the potential applications of the system in various domains, from everyday device control to support in educational and professional settings.

## II. BACKGROUND

Communication barriers faced by deaf and non-verbal individuals have long been a significant challenge in society. According to the World Health Organization (WHO), over 5

Hand gesture recognition systems, as a technology for non-verbal communication, have seen considerable advancements in recent years. These systems can convert hand movements or gestures into commands or even textual representations, enabling individuals to communicate without the need for vocal speech. Gesture recognition typically relies on computer

vision techniques, machine learning models, and deep learning algorithms to detect and interpret gestures accurately.

In particular, systems based on OpenCV and MediaPipe have gained prominence for their efficiency and ability to provide real-time performance. OpenCV, an open-source computer vision library, provides tools for image processing and object tracking, while MediaPipe, developed by Google, offers state-of-the-art solutions for real-time hand and body tracking. These tools have become integral to the development of hand gesture recognition systems due to their high accuracy and low latency in real-time applications.

Previous work has been done in the domain of sign language recognition and gesture-based interaction, but these systems often require complex setups or extensive hardware. Moreover, many systems are focused on recognizing only a limited set of gestures, often specific to particular use cases. The Silent Cue system aims to bridge this gap by providing a flexible, adaptive, and user-friendly solution that can recognize a wide range of hand gestures, enabling non-verbal communication for daily use and beyond.

By leveraging machine learning algorithms such as Random Forests for classification and utilizing real-time hand tracking, Silent Cue enables a hands-free interaction experience. This approach offers the potential to create a more inclusive communication ecosystem, enhancing the lives of deaf and non-verbal individuals by making it easier to communicate with others, interact with technology, and reduce reliance on external assistance.

#### A. Hand Gesture Recognition Techniques

Hand gesture recognition systems can be broadly categorized into vision-based and sensor-based approaches. Vision-based systems [1] rely on computer vision techniques such as image processing, feature extraction, and machine learning to detect and interpret hand gestures in real-time. These systems generally use cameras or depth sensors to capture the movements and gestures of the user. Vision-based techniques, particularly those using deep learning models, are becoming increasingly popular due to their accuracy and ability to recognize a wide range of gestures without the need for specialized hardware. On the other hand, sensor-based systems [2] use accelerometers, gyroscopes, and other wearable sensors to capture hand movements. While these systems may offer higher precision in some cases, they often require specialized hardware and are less flexible than vision-based systems.

#### B. Real-Time Hand Tracking Using OpenCV and MediaPipe

OpenCV and MediaPipe are two of the most widely used tools for real-time hand tracking and gesture recognition. OpenCV [3], a popular open-source computer vision library, provides an extensive suite of image processing techniques that can detect, track, and recognize hand gestures in real time. MediaPipe [4], developed by Google, offers state-of-the-art hand tracking capabilities with high accuracy and minimal latency. By using a series of hand landmarks, MediaPipe can track the motion of individual fingers and the palm of the

hand, allowing for robust gesture detection even in dynamic environments. These tools are often used in conjunction to create seamless hand gesture recognition systems, where OpenCV processes the captured image data, and MediaPipe extracts the hand landmarks that are used for further gesture classification.

#### C. Machine Learning for Gesture Classification

Machine learning plays a crucial role in gesture classification, transforming raw hand tracking data into meaningful actions. In the context of hand gesture recognition, supervised learning algorithms such as Random Forests [5] are often employed for classification tasks. These models are trained on labeled datasets, where each gesture is associated with specific features such as hand position, movement trajectory, and the relative angles between fingers. Once trained, the model can predict the gesture a user is performing in real-time based on new input data. The advantage of using Random Forests is their ability to handle large datasets and provide interpretable results, which are essential for developing user-friendly systems like Silent Cue. Other machine learning algorithms, such as Support Vector Machines (SVM) and Neural Networks, are also explored for improving classification accuracy and handling more complex gesture sets.

### III. METHODOLOGY

This section outlines the methodology used to develop and evaluate the hand gesture recognition system for facilitating non-verbal communication for deaf and non-verbal individuals. The approach aims to assess the accuracy and efficiency of various hand tracking and gesture recognition models, as well as evaluate the system's performance in real-world applications. Key aspects of the methodology include model selection, data collection, the pre-processing of hand gesture data, and the evaluation metrics used to measure the system's success. Additionally, the fine-tuning of machine learning models to improve gesture classification accuracy is explored. The goal is to offer insights into how different models perform under real-world conditions and to contribute to the improvement of hand gesture recognition technologies.

#### A. Model Selection

The performance of hand gesture recognition systems greatly depends on the model architecture and the ability to process input data accurately. Several models were considered for this project, each with unique characteristics suited to different aspects of hand gesture recognition:

- **Model 1: MediaPipe Hand Tracking**

Developed by Google, MediaPipe [1] provides real-time hand tracking and gesture recognition capabilities. MediaPipe tracks 21 hand landmarks and is capable of detecting gestures in dynamic environments. This model was chosen for its simplicity, real-time performance, and robust hand landmark detection. It can be used as the primary tool for detecting and tracking hand positions and movements.

- **Model 2: OpenCV Hand Tracking with Haar Cascades**

OpenCV [2] is a widely used library for computer vision tasks. The Haar Cascade Classifier is an object detection algorithm used for real-time hand tracking. It was selected for this project due to its lightweight nature and fast processing, making it suitable for lower-powered devices. OpenCV helps in pre-processing the captured image frames to identify hand contours and positions.

- **Model 3: Random Forest Classifier for Gesture Recognition**

A machine learning model, Random Forest [3] is an ensemble method that combines multiple decision trees to improve classification accuracy. This model is particularly useful for classifying hand gestures based on features like hand position, velocity, and trajectory. It was trained using a custom dataset of hand gestures to recognize various sign language symbols and non-verbal gestures.

- **Model 4: Support Vector Machine (SVM) for Gesture Classification**

Support Vector Machines (SVM) [4] are powerful classifiers that find an optimal hyperplane to separate classes of data. SVM was used to classify hand gestures by mapping hand gesture features (such as finger angles and hand position) into a higher-dimensional space. The SVM model was trained on labeled gesture data to recognize specific sign language gestures.

- **Model 5: Convolutional Neural Networks (CNN) for Gesture Recognition**

CNNs [5] are deep learning models that excel at identifying patterns and features in image data. A CNN was trained to recognize hand gestures directly from images captured by the camera. The CNN model is designed to process raw pixel data and automatically learn high-level features for gesture classification, providing a more flexible and robust gesture recognition solution.

- **Model 6: TensorFlow for Custom Gesture Classification**

TensorFlow [6] is an open-source deep learning framework that was used to develop custom models for recognizing a set of predefined gestures. TensorFlow was utilized to fine-tune a pre-trained neural network for gesture classification, enabling the system to identify gestures with high accuracy and minimal latency. TensorFlow's scalability and adaptability made it a key component of the project for real-time gesture recognition.

Table I provides a summary of the models examined in this study. After providing an overview of each model, the process of preparing the data is delved into, describing the dataset used for assessment and any processing measures taken to ensure consistency and precision in the evaluation of the model's performance.

## B. Dataset

**Silent Cue Gesture Dataset:** The dataset used in the Silent Cue project consists of hand gesture frames captured

TABLE I  
OVERVIEW OF HAND GESTURE RECOGNITION MODELS

Model Name	Description	Architecture	Size(Min)
MediaPipe Hand Tracking	Real-time hand tracking with 21 hand landmarks, used for gesture recognition	Transformer-based hand tracking model	20MB
OpenCV Hand Tracking	Hand tracking using Haar Cascade Classifier, optimized for real-time applications	Haar Cascade Classifier (Traditional)	10MB
Random Forest Classifier	Ensemble method for gesture classification, based on hand movement features	Machine Learning (Random Forest)	5MB
Support Vector Machine (SVM)	Classifies gestures based on geometric hand features	Support Vector Machine (SVM)	100MB
TensorFlow Gesture Recognition	Fine-tuned deep learning models for custom hand gesture classification	Deep Learning (CNN, Custom Architecture)	150MB

for American Sign Language (ASL) alphabets and various custom gestures. Each gesture in the dataset corresponds to a specific ASL alphabet character and other common non-verbal communication gestures, including greetings like "Hi," "Bye," and others tailored for specific applications. The dataset was created by capturing real-time frames using OpenCV and MediaPipe for accurate hand tracking, ensuring a comprehensive and diverse set of samples. These frames are annotated with the respective gesture labels, facilitating the training of machine learning models for gesture recognition. This dataset is crucial for enhancing non-verbal communication for deaf and non-verbal individuals, enabling hands-free interaction with devices through gesture recognition.

### Data Splits:

- Train: 2,600 images (26 alphabets x 100 images each)
- Validation: 300 images (for cross-validation)
- Test: 300 images (testing with publicly available ASL datasets)

### Sample Instance:

- "id": "123456",  
"gesture":  
– "Hand gesture: 'A' in ASL - 'Open hand with thumb extended Left' "  
"label":  
– "0",

## C. Evaluation Strategy

In this project, the effectiveness of the hand gesture recognition system is evaluated using several standard performance metrics, with the goal of assessing the accuracy and reliability of the gesture classification. The metrics used to evaluate the

model's performance include accuracy, precision, recall, F1-score, confusion matrix, and cross-validation. These metrics help assess the model's ability to classify hand gestures accurately and handle misclassifications.

- 1) **Accuracy** This metric calculates the proportion of correctly classified gestures out of the total gestures predicted. It is defined as:

$$\text{Accuracy} = \frac{\text{Correctly Predicted}}{\text{Total Predicted}}$$

Higher accuracy indicates better overall performance in identifying gestures from the dataset.

- 2) **Precision:** Precision measures the proportion of true positive predictions (correctly identified gestures) out of all the gestures predicted as a specific class. It is defined as:

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

Precision is important when minimizing false positives is a priority, ensuring that the model only classifies gestures as a certain class when it is confident in the prediction.

- 3) **Recall:** Recall calculates the proportion of true positive predictions (correctly identified gestures) out of all the actual instances of that gesture in the dataset. It is defined as:

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

Recall is important when it is critical to identify all instances of a particular gesture, even at the cost of some false positives.

- 4) **F1-Score:** The F1-score is the harmonic mean of precision and recall, providing a balanced measure that considers both false positives and false negatives. It is defined as:

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

The F1-score is especially useful when there is an uneven class distribution or when you need a balance between precision and recall.

- 5) **Confusion Matrix:** The confusion matrix is used to visualize the classification performance of the model. It provides a table that shows the number of correct and incorrect predictions for each gesture class. This matrix helps in analyzing the types of errors made by the model and which gestures are frequently misclassified.
- 6) **Cross-Validation:** Cross-validation is performed during the training phase to ensure that the model generalizes well across different subsets of the dataset. This technique splits the dataset into multiple training and validation sets, ensuring that the model's performance is consistent across different data samples and helping prevent overfitting.

These evaluation metrics are applied to assess the model's performance in recognizing and classifying hand gestures. The

evaluation is conducted on both the custom dataset created for the project (which includes hand gestures for ASL alphabets and other gestures) and any publicly available datasets used for testing. The combination of these metrics provides a comprehensive understanding of the model's accuracy, precision, recall, and overall classification performance.

**Sample Instance:** Accuracy scores calculated for a sample gesture recognition result using the following metrics:

- Accuracy:
  - 1) Accuracy (A): 0.85
- Precision (P):
  - 1) Precision (P): 0.87
- Recall (R):
  - 1) Recall (R): 0.82
- F1-Score (F):
  - 1) F1-score (F): 0.84

Although accuracy and other evaluation metrics offer important insights into how effectively the gesture recognition model performs, it's essential to understand their limitations. These metrics, such as accuracy, precision, recall, and F1-score, focus primarily on the quantitative aspect of performance, such as the number of correct predictions or the overlap between predicted and actual gestures. However, they may fail to fully capture the complexity of human gesture recognition, such as the nuanced differences between similar hand movements or subtle variations in gesture performance.

One of the major drawbacks is that these metrics do not account for the semantic similarity between gestures. For example, two different hand movements may convey the same meaning but may have slight variations in hand shape, position, or speed, which may not be reflected accurately in the performance scores. Additionally, factors like user experience, gesture recognition speed, and real-time adaptability—critical in real-world applications—are not always fully captured by these metrics.

Thus, in order to gain a more comprehensive understanding of the model's performance, these quantitative metrics should be combined with qualitative evaluation, such as user feedback, real-time testing, and expert assessment of gesture accuracy and contextual relevance. This combined approach will provide a fuller picture of how well the gesture recognition system is functioning, both from a technical and practical standpoint.

#### D. Fine-Tuning

In machine learning, fine-tuning is the process of adapting a pre-trained model to a specific dataset or task by retraining it with task-specific data [?]. In the context of gesture recognition, fine-tuning involves adjusting a model that has been pre-trained on a large, generic dataset to better understand and classify specific hand gestures, such as those used in American Sign Language (ASL). Fine-tuning enables the model to learn specialized patterns and features relevant to the gestures in your dataset, improving performance on the target task.

**Motivation for Fine-Tuning:** Pre-trained models may not always perform optimally for specialized tasks, such as recognizing specific hand gestures. Generic models are often trained on broad datasets and may not account for the nuances in hand shapes, movement patterns, and orientation present in sign language gestures. Fine-tuning helps the model adapt to these nuances, improving its ability to accurately recognize gestures in the target dataset, such as your custom ASL dataset.

**Use of Fine-Tuning:** Fine-tuning is particularly beneficial in applications like gesture recognition, where the model needs to adjust to the specific variations of hand shapes, positions, and motions used in sign language. By fine-tuning the pre-trained model using your custom dataset of ASL gestures, the model can learn more precise representations of each gesture and handle the variability in user performances.

**How Fine-Tuning Improves the Model:** Fine-tuning enhances the model's accuracy and efficiency in recognizing hand gestures by training it on specific, labeled data that reflect the unique characteristics of the gestures. The model's parameters are adjusted during the fine-tuning process to better represent these specialized features, leading to improved performance in classifying and recognizing each gesture. This fine-tuning process ensures that the model can deliver more accurate, contextually appropriate results in a real-world setting, making it more reliable for tasks such as hand gesture-based communication in ASL.

#### IV. RESULTS

##### A. Gesture Recognition Results

The performance of different gesture recognition models is presented in Table ?? . The evaluation metrics used for these models include Accuracy, Precision, Recall, and F1 Score. These metrics were computed on a custom ASL gesture dataset consisting of 100 images per gesture, covering the 26 letters of the ASL alphabet and several common phrases. The accuracy metric was calculated based on the percentage of correct predictions made by each model during testing.

**Model A** achieved the highest overall accuracy (92

**Model B** showed good precision (80

**Model C** demonstrated a balanced performance across all metrics, achieving an accuracy of 89

**Model D** also performed well with an accuracy of 91

Overall, these results suggest that the models are capable of recognizing ASL gestures with high accuracy. Model A outperforms others in terms of overall performance, but other models such as Model C and Model D show promising results as well.

##### B. Fine-Tuning Results

After fine-tuning the models on the custom ASL gesture dataset, improvements were observed across all evaluation metrics, as summarized below:

Fine-tuning the models led to improvements in the overall recognition accuracy and model performance. Specifically:

**Model A** showed the highest improvement, with accuracy increasing from 92

**Model B** also saw improvements in its precision and recall, with accuracy increasing to 87

**Model C** performed better after fine-tuning, with a slight increase in recall and accuracy. The model's F1 score improved significantly from 85

**Model D** showed notable improvements, with accuracy increasing from 91

These results show that fine-tuning can significantly improve the performance of gesture recognition models, particularly in adapting them to specialized datasets like the custom ASL gesture dataset.

#### V. CONCLUSION

In conclusion, the gesture recognition models evaluated in this study showed promising results in recognizing ASL gestures. The models demonstrated strong performance, with **Model A** achieving the highest overall accuracy and performance. Fine-tuning these models further enhanced their ability to recognize gestures accurately, with improvements in accuracy, precision, recall, and F1 score across all models.

These findings underscore the importance of model customization, especially through fine-tuning, in enhancing performance for specific tasks such as ASL gesture recognition. Future work could explore further optimizations and investigate additional models to continue improving recognition accuracy and robustness.

#### REFERENCES

- [1] "OpenCV: Open Source Computer Vision." Available: <https://opencv.org/>
- [2] "MediaPipe: Cross-platform library for building pipelines to process video, audio, and other multimedia types." Available: <https://google.github.io/mediapipe/>
- [3] "Random Forest Classifier – Scikit-learn." Available: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>
- [4] "MediaPipe: Cross-platform library for building pipelines to process video, audio, and other multimedia types." Available: <https://github.com/google/mediapipe>
- [5] "Convolutional 2D Layer – TensorFlow." Available: [https://www.tensorflow.org/api\\_docs/python/tf/keras/layers/Conv2D](https://www.tensorflow.org/api_docs/python/tf/keras/layers/Conv2D)
- [6] M. Lewis et al., 'BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension', arXiv [cs.CL]. 2019. Available: <https://arxiv.org/pdf/1910.13461.pdf>
- [7] M. Wu, A. Waheed, C. Zhang, M. Abdul-Mageed, and A. F. Aji, 'LaMini-LM: A Diverse Herd of Distilled Models from Large-Scale Instructions', arXiv [cs.CL]. 2024. Available: <https://arxiv.org/abs/2304.14402>
- [8] B. Gliwa, I. Mochol, M. Biesek, and A. Wawer, 'SAMSum Corpus: A Human-annotated Dialogue Dataset for Abstractive Summarization', in Proceedings of the 2nd Workshop on New Frontiers in Summarization, 2019. Available: <https://arxiv.org/pdf/1911.12237.pdf>
- [9] "MediaPipe: Cross-platform library for building pipelines to process video, audio, and other multimedia types." Available: <https://github.com/google/mediapipe>
- [10] A. Vaswani et al., 'Attention is All You Need,' in Advances in Neural Information Processing Systems, 2017. Available: <https://arxiv.org/abs/1703.03208>
- [11] K. Simonyan and A. Zisserman, 'Very Deep Convolutional Networks for Large-Scale Image Recognition', in Advances in Neural Information Processing Systems, 2014. Available: <https://arxiv.org/abs/1612.05350>
- [12] "Random Forest Classifier – Scikit-learn." Available: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>
- [13] "TensorFlow: Open source machine learning framework." Available: <https://github.com/tensorflow/tensorflow>

- [14] X. Chen, L. Zhu, and X. Xie, 'Learning Deep Representations of Fine-grained Visual Descriptions,' in International Conference on Computer Vision, 2015. Available: <https://arxiv.org/abs/1511.06722>