# Index

# ABSTRACT

Hand gesture recognition plays a crucial role in enhancing human-computer interaction (HCI) by enabling intuitive and natural communication between users and machines. This project aims to develop a real-time hand gesture recognition system using state-of-the-art technologies, specifically MediaPipe and TensorFlow. The project begins with a preliminary analysis, identifying common limitations from 20 research papers in existing hand gesture recognition systems and proposing solutions leveraging MediaPipe. The methodology involves collecting datasets, preprocessing data, and training a neural network model using TensorFlow. The system's architecture and functionality are explained, highlighting the use of MediaPipe for hand tracking and TensorFlow for deep learning-based gesture recognition. The project's significance lies in its potential to address common challenges in hand gesture recognition, such as environmental sensitivity and real-time performance, ultimately improving the usability and effectiveness of HCI applications. Through this project, we aim to contribute to the advancement of hand gesture recognition technology, fostering seamless interaction between users and computers in various domains.

## Keywords

- Hand Gesture Recognition
- Neural Networks
- Deep Learning
- Machine Learning
- MediaPipe
- TensorFlow
- OpenCV
- Feature Extraction
- Classification
- Real-Time Processing
- Pose Estimation
- Hand Tracking
- Gesture Recognition

# INTRODUCTION

Hand gestures serve as a vital means of communication between humans and machines across diverse domains. In addition to robotics, where they enable precise control of robotic arms and manipulation of objects, hand gestures have significant applications beyond the industrial and domestic settings. Real-time recognition of hand gestures equips robots to respond swiftly to human commands, enhancing their effectiveness in various tasks. Moreover, hand gesture recognition holds immense potential in assisting individuals with communication challenges, particularly in educational settings.

Interactive applications utilizing hand gestures offer immersive learning experiences, benefiting children and individuals with learning disabilities. By engaging users in hands-on activities and providing visual cues, these applications enhance engagement and comprehension, making education more accessible and enjoyable for diverse learners. Furthermore, in assistive technology, hand gesture recognition can facilitate communication for individuals with disabilities, offering alternative methods of interaction with computers and devices.

This project introduces a real-time hand gesture recognition system using MediaPipe and TensorFlow, two cutting-edge technologies. MediaPipe offers accurate hand tracking even in challenging environments, while TensorFlow provides efficient deep learning models for gesture recognition.

Compared to traditional approaches such as CNN and RNN, which may struggle with latency issues and computational overhead in real-time processing, MediaPipe and TensorFlow offer superior performance and efficiency.

By synergistically combining MediaPipe and TensorFlow, this project aims to revolutionize human-computer interaction in robotics, assistive technology, and education. Leveraging real-time hand gesture recognition, it empowers users with improved control over robots and enhances learning experiences, redefining interaction across various fields.

Through the integration of these advanced technologies, the project seeks to address existing challenges and unlock new possibilities for intuitive communication between humans and machines.

# PRELIMINARY ANALYSIS

## Literature Review:

1. In (Haria et al., 2017), the authors proposed a hand gesture recognition system which utilizes a combination of techniques, including contours, convexity defects, and the Haar cascade classifier for hand gesture recognition. These methods enable gesture recognition without the need for markers or additional hardware, making the system cost-effective and user-friendly. While it covers various gestures for different functions, it may encounter difficulties with complex hand movements or lighting variations.

2. In (Mohamed et al., 2012), the authors proposed a hand gesture recognition system which uses background subtraction and logical heuristic equations to spot gestures in live video. Although it works well, particularly in tricky backgrounds, it sometimes wrongly identifies other objects as hand gestures. Even though it has improved finger tip detection by 52%, the system's performance varies depending on the complexity of the background and the contrast with the skin color. The authors suggested that further enhancements may be needed to address these limitations and enhance overall accuracy and robustness.

3. In (Yun et al., 2012), the authors addressed the limitations of traditional vision-based hand gesture recognition methods, particularly those based on Multi-Feature Fusion and Template Matching. These limitations include low recognition accuracy, system instability, and ambiguous recognition results. These issues hinder the effectiveness of human-computer interaction systems and other applications reliant on accurate hand gesture recognition.

4. In the study conducted by Sharma et al. (2020), the authors noted several limitations. These include the focus on static hand gestures rather than dynamic ones, potentially restricting its applicability in real-time scenarios. Additionally, the research primarily utilizes American Sign Language images, which may not fully represent the diversity of hand gestures across different cultures or contexts. The study also relies on predefined classifiers, which may not capture the complexity of all hand gestures accurately. Finally, while the proposed ORB approach shows promising results, there may still be room for improvement in

terms of overall accuracy and robustness, especially when dealing with noisy or ambiguous input data.

5. The study conducted by Elsayed et al. (2017) explores hand gesture recognition using Histogram of Oriented Gradients (HOG), demonstrating effectiveness for rapid deployment but encountering challenges with diverse hand angles and backgrounds. Variations in lighting conditions and hand size further impact performance, suggesting the need for enhanced robustness across different scenarios. The authors suggested that further research is warranted to address these limitations and optimize performance for real-world applications.

6. A method was proposed by Bhuiyan et al. (2017) for improving hand gesture recognition by reducing feature dimensions. Utilizing webcam-based image processing, they extracted hand features through discrete wavelet transformation and singular value decomposition. Genetic algorithm optimizes feature selection, enhancing recognition accuracy for American Sign Language numerals. Despite promising results, challenges persist in handling diverse backgrounds and hand angles. The authors suggested that further research is needed to enhance robustness and scalability for real-world applications.

7. In (Katkar, G., Ramteke, A., & Gaikwad, S., 2023), the authors introduced a technique to aid in comprehending Indian Sign Language (ISL) through a computer program named support vector machine (SVM). The goal of this technique was to enhance communication for deaf and mute individuals, even amidst noisy surroundings. However, the program encounters challenges in rapidly changing, real-time scenarios, hindering its ability to accurately interpret signs. This limitation may restrict its effectiveness. The authors suggested that further investigation is necessary to refine the program's capabilities, particularly in dynamic, real-world situations, ensuring accurate sign language interpretation across diverse conditions.

8. Morimoto et al. (2012) stated that the system is proficient at recognizing hand gestures using certain models and motion details. However, it struggles when the gesture data is unclear or too dynamic. Additionally, having limited data can impact its effectiveness. To enhance its performance, improvements are needed in translating motion into symbols and exploring the utilization of more advanced models.

9. In (Li, M., Wang, X., et al., 2012), the authors proposed a dynamic hand gesture recognition system employing Hidden Markov Models (HMMs) and features from Adaboost with Histogram of Oriented Gradients (HOG) for hand detection and tracking. While achieving better recognition than traditional methods, it faces limitations. The tracking algorithm's sensitivity to light affects performance, and detection is only reported once a gesture ends. The authors stated that future efforts will concentrate on enhancing the tracking algorithm's robustness and improving the system's ability to recognize gestures in real-time and under varying lighting conditions.

10. In their study, the authors (Fortuna, L., Nguyen Huu, P., & Ngoc Phung, T., 2021) developed a system that recognizes gestures for robots using Histogram of Oriented Gradients (HOG) and Support Vector Machine (SVM), achieving 80% accuracy. However, due to insufficient gesture examples, the system did not perform perfectly. Their future objectives include making the system faster, improving accuracy through better image data, and integrating neural networks with other systems, despite potential requirements for increased computational power and complexity.

11. In (Kumar, Sunil & Srivastava, Tanu & Singh, Raj., 2017), the authors proposed a hand gesture recognition system using PCA, skin color segmentation, and thresholding, achieving high accuracy in controlled conditions but lower accuracy in low brightness images. It faces limitations with variations in skin color and uncontrolled background, and it may not work well with dynamic gestures or non-skin-colored hands. The authors suggested that further improvements could address dynamic gestures and expand compatibility beyond static gestures and skin-colored hands.

12. In (Arora et al., 2014), the authors proposed Neural Networks for recognizing hand gestures, particularly aiding non-verbal communication like sign language. Although Neural Networks perform effectively with small datasets, they might encounter issues with larger datasets.

13. Haroon et al. (2022) introduced a technique to enhance hand gesture recognition in various lighting conditions using the luminosity method with Artificial Neural Network. However, the approach relies on specific features, potentially limiting its ability to accurately detect complex gestures.

Furthermore, the study was conducted in controlled environments, raising concerns about its applicability in real-world settings. The authors suggested that future research should focus on addressing these limitations and exploring a wider range of scenarios to improve overall accuracy.

14. In (Nogales et al., 2023), the authors compared two approaches for finding features in hand gesture recognition. One approach involved human annotation using Artificial Neural Networks (ANN) and Support Vector Machines (SVM), which yielded good results but was slow. The other approach utilized Convolutional Neural Networks (CNN) and Bidirectional Long Short-Term Memory (BiLSTM) networks automatically, resulting in significantly better and faster performance. However, these automated methods may require more computational power and extensive datasets, posing challenges for scalability and deployment on a large scale.

15. In another study, Chatterjee et al. (2022) demonstrated the use of a Generative Adversarial Network (GAN) to generate synthetic hand gesture images for training Convolutional Neural Networks (CNNs). While the generated images exhibit similarity to real ones, they may not fully capture the variability and complexity of real-world scenarios. The dataset lacks diversity in hand shapes, sizes, and environmental conditions, potentially limiting the model's robustness in practical applications with varied factors. The authors suggested that future iterations should focus on enhancing data variability and addressing these limitations to improve the model's performance in diverse settings.

16. In (Tsinganos et al., 2022) , the authors introduced a Temporal Convolutional Network (TCN) model for recognizing hand gestures based on sEMG signals. While the model performed well in offline evaluations using complete sEMG sequences, its accuracy dropped significantly in real-time evaluations when provided with shorter segments of data. Comparisons with similar TCN models in the literature suggested that the model's performance could be improved with a tailored data-feeding procedure during training specifically designed for real-time applications. Additionally, employing a training procedure based on sliding windows improved the real-time performance, indicating the need for more effective training methods for real-time gesture recognition.

17. In (M. Z., Hossain et al., 2019), the authors introduced a method for static hand gesture recognition using Convolutional Neural Networks (CNNs) with data augmentation. While CNNs have shown high accuracy in recognizing hand gestures, limitations exist such as the reliance on a less complex background and assumption of gestures made with one hand. Additionally, the effectiveness of the model may decrease when applied to gestures in more complex backgrounds or involving both hands. Furthermore, the study does not explore the recognition of dynamic gestures or address potential issues related to scalability or real-time processing.

18. Obaid et al. (2020) proposed a method for hand gesture recognition in video sequences using deep convolutional and recurrent neural networks. However, the method relies heavily on specific data preprocessing techniques and augmentation methods, potentially limiting its applicability to real-world scenarios. Additionally, the study primarily focuses on the program's accuracy without considering factors such as its performance with noisy data or processing speed. Future research should explore alternative data preparation methods and evaluate the program under diverse conditions to enhance its effectiveness.

19. Toro-Ossaba et al. (2022) proposed a method using LSTM-RNN for recognizing hand gestures with EMG signals, which shows promise but has some limitations. It doesn't perform as well in real-time testing because it's sensitive to small changes in how the EMG armband is positioned and to signal noise. To improve its performance, we need to enhance how we collect and process the EMG data and find the right balance between data quality and model setup. In the future, we plan to upgrade the EMG armband, explore new data processing techniques, and adapt the model to work with a microcontroller for controlling hand prostheses with different hand gestures.

20. In (Koch et al., 2019), the authors investigated the use of Recurrent Neural Networks (RNNs) with accelerometer data for hand gesture recognition, which is important for affordable and mobile systems. Accelerometers capture forearm movements and can be easily integrated into mobile devices. The proposed neural network outperforms existing methods, especially for amputees, but may struggle with complex gestures or noisy environments.

# Common Limitations: Insights from 20 Research Papers

Here are some common limitations identified across the 20 research papers on hand gesture recognition:

## 1. Environmental Sensitivity:

Many systems are sensitive to variations in lighting conditions, background complexity, and noise, which can affect their accuracy and robustness in real-world scenarios.

## 2. Limited Accuracy in Challenging Scenarios:

Certain systems exhibit lower accuracy or reliability in challenging conditions such as low brightness, noisy backgrounds, or when dealing with non-skin-colored hands.

## 3. Hardware Dependency:

Some approaches rely on specific hardware components like webcams, EMG armbands, or accelerometers, which may limit their accessibility or deployment in practical settings.

## 4. Real-Time Performance:

Many systems encounter difficulties in real-time gesture recognition, such as processing speed, latency, or the ability to handle rapidly changing scenarios.

## 5. Data Limitations:

Some systems face challenges due to limited or unrepresentative datasets, impacting their ability to generalize across diverse scenarios or recognize a wide range of gestures accurately.

These limitations underscore the ongoing need for research and development efforts to address challenges in hand gesture recognition systems, aiming to enhance their accuracy, robustness, accessibility, and applicability across various real-world scenarios and user populations.

# Problem Statement:

The main objective of this study is to:-

1. Develop a system that recognizes hand gestures using MediaPipe and TensorFlow.
2. Make the system track hand movements in real-time from a camera.
3. Ensure the system can understand hand gestures even in different lighting and background noise.
4. Show the recognized hand gestures on the screen accurately.
5. The system should be able to identify different hand gestures.
6. The system should be able to handle fast changes in hand movements and backgrounds.

# Data Collection:

The dataset for this project is sourced from various channels to ensure its diversity and representativeness of real-world scenarios. To augment this dataset, numerous publicly available datasets from Kaggle, specifically designed for hand gesture recognition tasks, are utilized. These datasets contain labeled images and videos of hand gestures captured under controlled conditions, providing a foundational resource for model training and evaluation. Additionally, synthetic data generation techniques are employed to further enrich the dataset. By leveraging computer graphics software or simulation environments, artificial images or videos of hand gestures are generated, diversifying the dataset and enhancing model robustness.

Furthermore, the project takes advantage of pre-trained models provided by the MediaPipe framework for tasks such as hand detection and keypoint extraction, as well as a pre-trained TensorFlow model dedicated to gesture recognition. This strategic decision significantly streamlines the development process by eliminating the need to undertake the time-consuming task of collecting and annotating a dataset from scratch. By leveraging these readily available models, the project expedites the implementation of hand gesture recognition functionalities without the necessity of extensive data collection efforts.

# Actors in the system:

**Users:** Interacts with the hand gesture recognition system by performing gestures, such as controlling virtual environments or translating sign language.

**Camera:** Captures real-time video frames of the hand gestures performed by the users.

**Computer:** Processes the video frames and executes the hand gesture recognition algorithm.

**Admin:** Manages the system's deployment, infrastructure, and user access, ensuring smooth operation and performance.

# Hardware & Software used:

**Hardware:**

- Camera
- Installed RAM: 16.0 GB (14.9 GB usable).
- System Type: 64-bit operating system, x64-based processor.
- Additional Information:
    - Color Format: RGB
    - Color Space: SDR (Standard Dynamic Range)
    - Refresh Rate: 1.000 Hz

**Software:**

- Anaconda 2023.09
- PyCharm 2023.3.2
- Python-3.11
- OpenCV (Open Source Computer Vision Library) - Version 4.5
- MediaPipe - Version 0.8.5
- TensorFlow - Version 2.5.0
- NumPy - Version 1.19.3

## Project Schedule

| Project Task | Start Date | End Date | Duration |
|---|---|---|---|
| Literature review and research problem formulation. | 09th December 2023 | 19th December 2023 | 11 |
| Defining research problem and stating the objective. | 20th December 2023 | 23rd December 2023 | 4 |
| Understanding Neural Network and its architecture. | 24rd December 2023 | 25th December 2023 | 2 |
| Understanding MediaPipe and its architecture. | 26th December 2023 | 30th December 2023 | 5 |
| Dataset Collection | 30th December 2023 | 31st December 2023 | 2 |
| Drafting UMLs | 01st January 2024 | 25th January 2024 | 26 |
| Model Building and Testing | 26th January 2024 | 26th February 2024 | 31 |
| Building UI | 27th February 2024 | 08th March 2024 | 11 |
| Statistical testing | 09th March 2024 | 20th March 2024 | 12 |
| Final Documentation | 09th December 2023 | 09th April 2024 | 103+20 |
| | | | 123 days |

# Grantt Chart:

**PHASE - I**

| | | |
|---|---|---|
| Literature review and research problem formulation. | 9th Dec — 19th Dec | |
| Defining research problem and stating the objective. | 20th Dec — 23rd Dec | |
| Understanding Neural Network and its architecture. | 24th Dec — 25th Dec | |
| Understanding MediaPipe and its architecture. | 26th Dec — 30th Dec | |
| Dataset Collection | 30th Dec — 31st Dec | |
| Drafting UMLs | 1st Jan — 25th Jan | |

**2023** November | December | January **2024**

**PHASE-II**

| | | |
|---|---|---|
| Model Building and Testing | 26th Jan — 26th Feb | |
| Building UI | 27th Feb — 8th March | |
| Statistical testing | 9th March — 20th March | |
| Final Documentation | 9th Dec — 9th April | |

**2023** Nov | Dec | Jan | Feb | March | April **2024**

## Pert Chart:

| Project Task (Activity) | Expected Time (te) | | | |
|---|---|---|---|---|
| | | Actual Time (tm) | Shortest Time (to) | Longest Time (tp) |
| Literature review and research problem formulation. | 10 | | | 11 |
| Defining research problem and stating the objective. | 8 | | 4 | |
| Understanding Neural Network and its architecture. | 5 | | 2 | |
| Understanding MediaPipe and its architecture. | 5 | 5 | | |
| Dataset Collection | 2 | 2 | | |
| Drafting UMLs | 20 | | | 26 |
| Model Building and Testing | 30 | | | 31 |
| Building UI | 10 | | | 11 |
| Statistical testing | 20 | 12 | | |
| Final Documentation | 5 | | | 20 |
| Total Time | 115 days | | 123 days | |

# Frameworks used for Hand Gesture Recognition:

Gesture recognition is an active research field in Human-Computer Interaction technology with numerous applications in virtual environment control and robot control. In this Hand Gesture Recognition project, a real-time Hand Gesture Recognizer is developed using the MediaPipe framework and TensorFlow within the OpenCV and Python environment.

OpenCV is a real-time computer vision and image processing framework originally built in C/C++, but in this project, Python is utilized through the OpenCV-python package.

## Understanding OpenCV, A Powerful Tool for Real-Time Computer Vision:

OpenCV, short for Open Source Computer Vision Library, is an essential toolkit for anyone working with computer vision and machine learning. The OpenCV Library has a suite of over 2500 optimized algorithms. The deep learning aspect of OpenCV includes neural networks, with support for frameworks such as TensorFlow.

## Flowchart of OpenCV:-

**OpenCV's architecture:-**

OpenCV's architecture is designed to handle a wide range of tasks in computer vision and machine learning.



**OpenCV's Architecture**

At its core is a component called CXCore, which contains the main functions and algorithms. This setup reduces redundancy and makes operations more efficient.

In addition to the core, there are other components like CV, which deals with image processing and vision algorithms, and MLL, which includes tools for statistical classifiers and clustering. Another important component is HighGUI, which focuses on graphical user interface functions and image/video input/output.

This modular design allows different machine learning features to be easily integrated, making OpenCV flexible and versatile for various applications.

**Understanding TensorFlow, A Deep Learning Framework:**

TensorFlow is an open-source library for machine learning and deep learning developed by the Google brains team. It can be used across a range of tasks but has a particular focus on deep neural networks. TensorFlow supports various types of neural networks, including Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), in addition to traditional neural networks. This versatility allows developers and researchers to leverage TensorFlow for a wide range of machine learning and deep learning tasks.

**Neural Networks** are also known as artificial neural networks. It is a subset of machine learning and the heart of deep learning algorithms. The concept of Neural networks is inspired by the human brain. It mimics the way that biological neurons send signals to one another. Neural networks are composed of node layers, containing an input layer, one or more hidden layers, and an output layer.

**The input layer** is the first layer of the neural network. It consists of input nodes, where each node represents a feature or input variable. The number of nodes in the input layer is determined by the dimensionality of the input data.

**Hidden layers** are intermediate layers between the input and output layers. Each hidden layer contains one or more layers of nodes, known as hidden nodes or neurons. The number of hidden layers and the number of nodes in each hidden layer are design choices and can vary depending on the complexity of the problem and the data.

**The output layer** is the final layer of the neural network. It consists of output nodes, where each node represents a predicted output or target variable. The number of nodes in the output layer depends on the type of task the neural network is performing.

**Convolution neural network:-**



**Feature Learning:** CNNs are like smart filters that learn to recognize patterns in images. These patterns can be things like edges, textures, or shapes.

**Layers:** A CNN is made up of layers. The first layer looks at tiny parts of the image, like pixels. Then, deeper layers combine these parts to recognize more complex features.

**Convolution:** The key idea is convolution, where small filters slide across the image, looking for patterns. Each filter finds specific features, like edges or corners.

**Pooling:** After convolution, pooling layers reduce the size of the image. They keep the important information while discarding less relevant details.

**Fully Connected Layers:** Finally, the features learned are passed to fully connected layers. These layers combine all the information to make a final decision, like classifying the image as a cat or dog.

**Understanding MediaPipe Framework:**

MediaPipe stands as a robust tool engineered by Google to facilitate developers in harnessing the power of machine learning. As an open-source framework, it invites individuals from diverse backgrounds to leverage its capabilities and even tailor them to their specific needs. Despite its extensive array of functionalities, MediaPipe remains impressively lightweight, ensuring smooth operation across various devices and platforms.





One of the standout features of MediaPipe is its provision of pre-built solutions for common tasks in the realm of machine learning. These encompass a wide range of functionalities, from facial recognition and pose estimation to hand gesture recognition and object detection in both images and videos. By offering these ready-to-use tools, MediaPipe streamlines the development process, empowering

developers to create innovative applications without having to build everything from scratch. Ultimately, MediaPipe strives to democratize access to advanced machine learning technologies, fostering a culture of creativity and ingenuity across different domains.

**Hand Perception Pipeline Overview:**



**Palm Detector Model:**

MediaPipe framework includes a pre-trained model called BlazePalm to detect the initial position of the palm.

BlazePalm uses a method called non-maximum suppression to find the most likely location of the palm, using square bounding boxes to simplify the detection process.

It also uses an encoder-decoder feature extraction technique to better understand the context of the scene, even when dealing with small objects.

During training, BlazePalm minimizes focal loss to effectively handle a large number of different hand sizes.

**Hand Landmark Detection:**



| | | | |
|---|---|---|---|
| **0.** | WRIST | **11.** | MIDDLE_FINGER_DIP |
| **1.** | THUMB_CMC | **12.** | MIDDLE_FINGER_TIP |
| **2.** | THUMB_MCP | **13.** | RING_FINGER_MCP |
| **3.** | THUMB_IP | **14.** | RING_FINGER_PIP |
| **4.** | THUMB_TIP | **15.** | RING_FINGER_DIP |
| **5.** | INDEX_FINGER_MCP | **16.** | RING_FINGER_TIP |
| **6.** | INDEX_FINGER_PIP | **17.** | PINKY_MCP |
| **7.** | INDEX_FINGER_DIP | **18.** | PINKY_PIP |
| **8.** | INDEX_FINGER_TIP | **19.** | PINKY_DIP |
| **9.** | MIDDLE_FINGER_MCP | **20.** | PINKY_TIP |
| **10.** | MIDDLE_FINGER_PIP | | |

Once the palm is detected, the Hand Landmark model within MediaPipe accurately identifies 21 key points on the hand, including the 3D coordinates of the hand knuckles. This model employs regression techniques to predict the coordinates of these landmarks directly within the detected hand regions. MediaPipe streamlines the process of hand detection and key point extraction, which are essential for gesture analysis. With MediaPipe, users can effortlessly identify hands within images or video streams and obtain detailed information about their positions and movements. Specifically, MediaPipe provides a comprehensive set of 21 key points for each detected hand, significantly enhancing the accuracy and efficiency of real-time hand gesture recognition systems.

# Methodology

## Collection of Datasets:

This project utilizes two different approaches for hand gesture recognition.

The first approach involves using pretrained convolutional neural networks (CNNs) model on MediaPipe framework that have been trained by me on datasets available from Kaggle. These datasets consist of labeled images and videos showcasing different hand gestures captured in controlled environments. These datasets serve as a solid foundation for training and evaluating the models. To enhance the dataset further, synthetic data generation techniques are applied. This involves using computer graphics software or simulation environments to create artificial images or videos of hand gestures. These synthetic data additions help increase the diversity of the dataset, making the models more robust and capable of handling various real-world scenarios.

In the second approach, a pretrained model from MediaPipe is utilized for hand gesture recognition. This pretrained model, available within the MediaPipe Gesture Recognizer task, is specifically designed to recognize hand gestures in real-time. The model leverages machine learning techniques to detect hand landmarks and classify gestures based on hand geometry. By utilizing this pretrained model, the project benefits from a robust and efficient solution for hand gesture recognition without the need for extensive custom training.

### First Approach (Training CNN and using Mediapipe on it):-

**Dataset Acquisition:** Utilized publicly available datasets from Kaggle and synthetic datasets tailored for hand gesture recognition tasks. Acquired labeled images and videos capturing various hand gestures under controlled conditions. Ensured dataset diversity in terms of hand poses, angles, and lighting conditions.

**Data Preprocessing:** Standardized image sizes and pixel values to facilitate model training. Augmented the dataset using techniques such as rotation, translation, scaling, and flipping to enhance diversity. Splited the dataset into training, validation, and testing sets for robust model evaluation.

**Model Selection and Pretraining:** Selected a suitable pretrained convolutional neural network (CNN) architecture for hand gesture recognition tasks. Loaded

pretrained CNN model weights and froze layers to retain learned features during training. Customized the architecture to align with the specific requirements of hand gesture recognition.

**Model Training:** Trained the CNN model on the preprocessed dataset for hand gesture recognition. Fine-tuned the model by selectively unfreezing layers to adapt to the target task. Monitored training progress using established metrics and adjusted hyperparameters accordingly.

**Evaluation and Performance Analysis:** Conducted thorough evaluation of the trained model on validation and testing datasets. Assessed metrics such as accuracy, precision, recall, and F1-score to gauge model performance. Analyzed model predictions and identified areas for improvement based on evaluation results.

**Integration of MediaPipe:** Integrated the trained model into Pycharm and utilized MediaPipe for optimized deployment on target platforms.

**Result and Observations:** The trained model exhibited satisfactory accuracy, showcasing its proficiency in recognizing hand gestures, however, challenges arose with variations in hand poses and angles, impacting the model's prediction accuracy. Particularly, the model struggled to generalize across diverse poses, highlighting the importance of training data diversity. Additionally, limitations were observed in the model's ability to track and recognize dynamic hand movements seamlessly, posing potential challenges in real-time gesture recognition applications.

## Second Approach: Utilizing Pretrained Model from MediaPipe within the MediaPipe Framework:-

**Model Integration:** Incorporated a pretrained model from the MediaPipe framework, specifically designed for hand gesture recognition tasks.

**MediaPipe Framework:** Leveraged the MediaPipe framework's capabilities to implement real-time hand gesture recognition efficiently.

**Ease of Implementation:** Benefitted from the ready-to-use solution provided by MediaPipe, reducing the need for extensive custom training or model development.

**Functionality:** The pretrained model facilitated the detection of hand landmarks and classification of gestures based on hand geometry, ensuring robust recognition performance.

**Deployment:** Seamlessly integrated the pretrained model within the MediaPipe framework, enabling straightforward deployment and utilization for various applications requiring hand gesture recognition.

**Result and Observations:** The integration of the pretrained model from the MediaPipe framework yielded impressive results. The hand gesture recognition system exhibited smooth and accurate performance, surpassing expectations in real-time applications. Leveraging the MediaPipe framework's capabilities, the model demonstrated efficient detection of hand landmarks and classification of gestures, ensuring robust recognition performance.

Moreover, the implementation of the pretrained model significantly reduced the complexity of development, benefiting from MediaPipe's ready-to-use solution. This streamlined approach minimized the need for extensive custom training or model development, thereby accelerating the deployment process.

Overall, the deployed system showcased remarkable functionality, delivering fast and precise hand gesture recognition. It outperformed the first approach in terms of smoothness and speed, highlighting the effectiveness of utilizing pretrained models within the MediaPipe framework for real-time applications.

**A Deep Dive into MediaPipe:-**

MediaPipe provides a suite of libraries and tools to quickly apply artificial intelligence (AI) and machine learning (ML) techniques in applications. Users can plug these solutions into their applications immediately, customize them to their needs, and use them across multiple development platforms. User can further customize the their codes using MediaPipe to meet their application needs.

These libraries and resources provide the core functionality for each MediaPipe Solution:

**MediaPipe Tasks:** Cross-platform APIs and libraries for deploying solutions.

**MediaPipe Models:** Pre-trained, ready-to-run models for use with each solution. These tools let users customize and evaluate solutions.

**MediaPipe Model Maker:** Customize models for solutions with your data.

**MediaPipe Studio:** Visualize, evaluate, and benchmark solutions in your browser.

## Working of project :-

**Import Necessary Packages:** The project begins by importing essential Python packages, including OpenCV, NumPy, MediaPipe, and TensorFlow. These libraries provide functionalities for image processing, machine learning, and hand gesture recognition.

**Initialize Models:** Two main models are initialized: one for hand detection using the MediaPipe framework and another for hand gesture recognition using a pre-trained TensorFlow model. The hand detection model identifies the landmarks of the detected hand, while the gesture recognition model predicts the gesture based on these landmarks.

**Read Frames from Webcam:** The project captures frames in real-time from the webcam using the OpenCV library. These frames serve as input for hand detection and gesture recognition.

**Detect Hand Keypoints:** Each frame captured from the webcam is processed to detect the keypoints (landmarks) of the hand using the MediaPipe framework. The detected keypoints represent the spatial positions of various parts of the hand, such as fingertips and joints.

**Recognize Hand Gestures:** Once the hand keypoints are obtained, they are fed into the pre-trained TensorFlow model for gesture recognition. The model predicts the gesture based on the spatial configuration of the hand landmarks. The predicted gesture is then displayed on the frame using OpenCV.

**Display Output:** The final output frame contains visualizations of the detected hand landmarks and the recognized hand gesture. This output is continuously updated as new frames are captured from the camera.

**User Interaction:** The user can interact with the system by performing hand gestures in front of the webcam. The system processes these gestures in real-time and provides corresponding feedback on the screen.

## Working of MediaPipe in diagram:-



The proposed learning-based framework combines the MediaPipe pre-trained models for hand/object detection and tracking with a multi-class classifier for object recognition based on the hand keypoints.

### Flow Diagram of MediaPipe working:-

# System Analysis:

## Tools and Technologies used:-

Python (Pycharm) - Frontend and Backend (both).

Tensorflow - Used for building and training machine learning models, particularly neural networks, for tasks such as gesture recognition.

MediaPipe - Utilized for real-time hand gesture recognition by providing pre-trained models and tools for efficient implementation.

OpenCv - Employed for computer vision tasks, including reading frames from a webcam, image processing, and displaying the results, in the context of hand gesture recognition.

Draw.io - For making UML diagrams.

Excel - For statistical analysis.

## Event Table:-

| Event | Source | Trigger | Activity | Response | Destination |
|---|---|---|---|---|---|
| Initialize camera | User | Program start | Create VideoCapture object for webcam | Success | User |
| Read Frame | Camera | Continuous Loop | Read frame from webcam | Success | User |
| Flip Frame | Program | Read Frame | Vertically flip the Frame | Success | User |
| Display Frame | Program | Read Frame | Show the frame with captured video using OpenCV imshow | Success | Display |
| Check for Quit Command | User | Key press event (every iteration) | Check if 'q' key is pressed | If 'q' key pressed, | Program |

| | | | | break loop and release resources | |
|---|---|---|---|---|---|
| Process Hand Landmarks | Program | Read Frame | Convert frame from BGR to RGB format hand landmark detection | Success | User |
| Detect Hand Keypoints | MediaPipe | Process Hand Landmarks | Process RGB frame to detect hand landmarks | Success, hand landmark detected | User |
| Draw Hand Landmarks | Program | Detect Hand Keypoints | Draw detected hand landmarks on the frame | Success, landmarks drawn on the frame | Display |
| Predict Hand Gestures | TensorFlow Model | Detect Hand Keypoints | Pass hand landmarks to TensorFlow model for gesture prediction | Success, predicted gesture class | Program |
| Get Gesture Prediction | Program | Predict Hand Gestures | Retrieve predicted gesture class index | Success, gesture class index | Program |
| Display Gesture | Program | Get Gesture Prediction | Display predicted gesture label on the frame | Success, gesture label displayed on the frame | Display |
| Release Camera | Program | Quit Command detected | Release webcam resource | Success, webcam resource released | System |
| Close OpenCV Windows | Program | Quit Command detected | Close all active OpenCV windows | Success, OpenCV windows closed | System |

**Block Diagram:-**

## Uml Diagrams:-

- **System Design**



- **Data Flow Diagram**

- **Use Case Diagram**



- **Sequence Diagram**

# Code and Implementation

## Code Snippet:-

### Step 1 – Import necessary packages:

```
# import necessary packages for hand gesture recognition project using Python OpenCV
import cv2
import numpy as np
import mediapipe as mp
import tensorflow as tf
from tensorflow.keras.models import load_model
```

### Step 2 – Initialize models:

### Initialize MediaPipe:

```
# initialize mediapipe
mpHands = mp.solutions.hands
hands = mpHands.Hands(max_num_hands=1, min_detection_confidence=0.7)
mpDraw = mp.solutions.drawing_utils
```

### Note:-

The Mp.solution.hands module performs the hand recognition algorithm.

An object is created and stored in mpHands to utilize this module.

Using the mpHands.Hands method, the model is configured.

The first argument, max_num_hands, determines the maximum number of hands that the model will detect in a single frame.

Although MediaPipe can detect multiple hands in a single frame, only one hand will be detected at a time in this project.

The Mp.solutions.drawing_utils library is employed to draw the detected key points automatically, eliminating the need for manual drawing.

**Initialize Tensorflow:**

```
# Load the gesture recognizer model
model = load_model('mp_hand_gesture')
# Load class names
f = open('gesture.names', 'r')
classNames = f.read().split('\n')
f.close()
print(classNames)
```

**Note:-**

The load_model function is utilized to load the pre-trained TensorFlow model.

The Gesture.names file includes the names of the gesture classes.

Initially, the file is opened using Python's built-in open function, and then its contents are read using the read() function.

**Step 3 – Read frames from a webcam:**

```
# Initialize the webcam for Hand Gesture Recognition Python project
cap = cv2.VideoCapture(0)
while True:
  # Read each frame from the webcam
  _, frame = cap.read()
x , y, c = frame.shape
  # Flip the frame vertically
  frame = cv2.flip(frame, 1)
  # Show the final output
  cv2.imshow("Output", frame)
  if cv2.waitKey(1) == ord('q'):
```

```
                    break

# release the webcam and destroy all active windows

cap.release()

cv2.destroyAllWindows()
```

**Note:-**

A VideoCapture object is instantiated by passing the argument '0', representing the camera ID of the system. If multiple webcams are connected, the camera ID can be adjusted accordingly. The cap.read() function is then used to read each frame from the webcam. Subsequently, the cv2.flip() function is applied to flip the frame, and cv2.imshow() displays the frame on a new OpenCV window. Finally, the cv2.waitKey() function is used to keep the window open until the key 'q' is pressed.

**Step 4 – Detect hand keypoints:**

```
framergb = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB)

  # Get hand landmark prediction

  result = hands.process(framergb)


className = ''

  # post process the result

  if result.multi_hand_landmarks:

        landmarks = []

        for handslms in result.multi_hand_landmarks:

        for lm in handslms.landmark:

                # print(id, lm)

                lmx = int(lm.x * x)

                lmy = int(lm.y * y)


                landmarks.append([lmx, lmy])

        # Drawing landmarks on frames
```

mpDraw.draw_landmarks(frame, handslms,

mpHands.HAND_CONNECTIONS)

**Note:-**

MediaPipe operates with RGB images, while OpenCV reads images in BGR format. Therefore, the cv2.cvtColor() function is used to convert the frame to RGB format. The process() function accepts an RGB frame and returns a result class. Subsequently, the script checks if any hand is detected using the result.multi_hand_landmarks method. Upon detection, the script iterates through each detection and stores the coordinates in a list called landmarks. Additionally, the coordinates are scaled by multiplying them with the image height (y) and width (x) since the model returns normalized values between 0 and 1. Finally, the mpDraw.draw_landmarks() function is used to draw all the landmarks on the frame.

**Step 5 – Recognize hand gestures:**

# Predict gesture in Hand Gesture Recognition project

prediction = model.predict([landmarks])

print(prediction)

classID = np.argmax(prediction)

className = classNames[classID]

 # show the prediction on the frame

cv2.putText(frame, className, (10, 50), cv2.FONT_HERSHEY_SIMPLEX,

1, (0,0,255), 2, cv2.LINE_AA)

**Note:-**

The model.predict() function receives a list of landmarks and returns an array containing 10 prediction classes for each landmark. The output is represented as a multidimensional array, where each inner array corresponds to the prediction probabilities for a particular landmark. After obtaining the prediction probabilities, np.argmax() is applied to find the index of the maximum probability, which corresponds to the predicted class. The class name corresponding to this index is then retrieved from the classNames list. Finally, the cv2.putText function is utilized to display the detected gesture on the frame.

# Code Snapshot:

1.



2.

3.



4.

# Output:

1. After clicking on input button-



2. Camera starts-

3. Starts recognizing hand gestures.

## Discussion:-

The central focus of this project was to thoroughly investigate the functionality and effectiveness of the MediaPipe framework alongside its pretrained 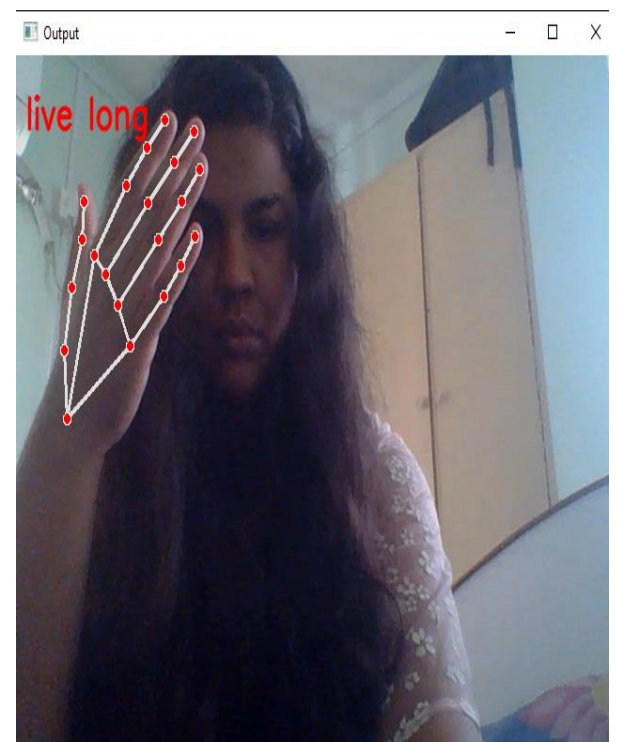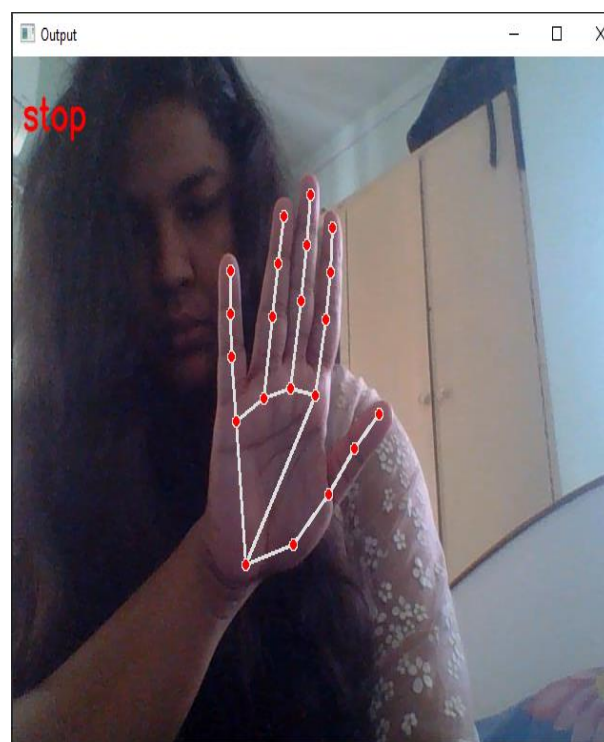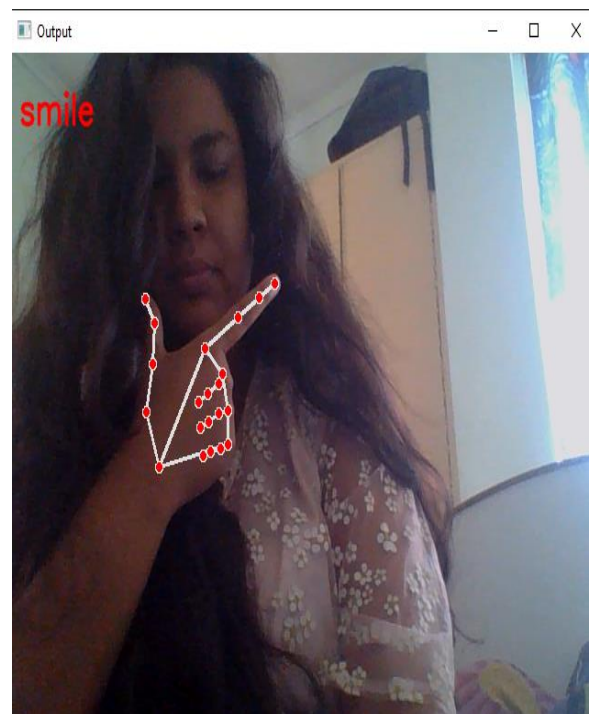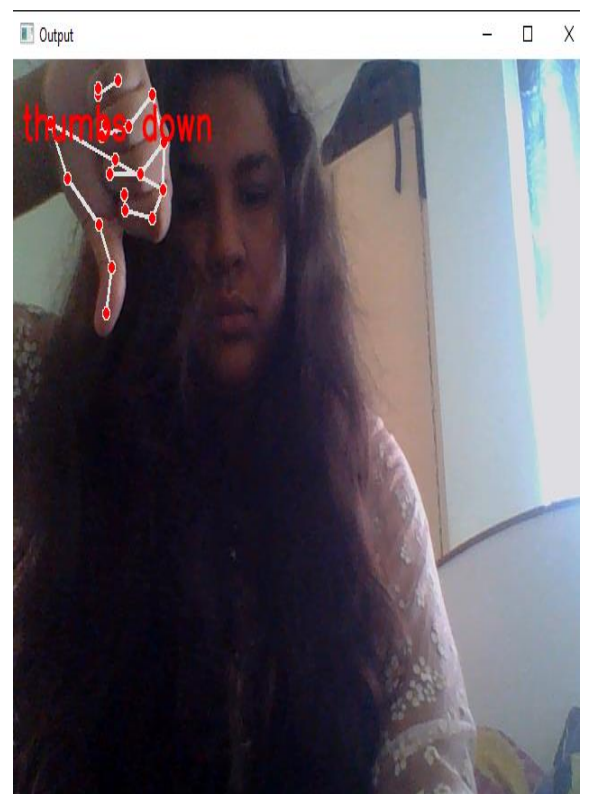model for hand gesture recognition. The overarching goal was to conduct a comprehensive comparison between the utilization of a custom-trained Convolutional Neural Network (CNN) and the adoption of a pretrained model from MediaPipe within the MediaPipe framework., with the aim of gaining deeper insights into the intricacies of MediaPipe and its functionalities.

In the pursuit of employing a custom CNN approach, a series of challenges emerged. The model encountered notable difficulties in smoothly recognizing hand gestures, especially when there were variations in hand positions. Despite concerted efforts to train the model effectively, it exhibited relatively accurate results; however, it struggled with timely recognition, particularly in response to changes in hand positions.

Conversely, the integration of the pretrained MediaPipe model yielded highly promising outcomes. This pre-built solution, meticulously designed for real-time hand gesture recognition, showcased exemplary performance characterized by its seamless and accurate functionality. The incorporation of MediaPipe into the project significantly simplified the process, eliminating the arduous task of extensive training while delivering consistently superior results.

In comparison to CNN methodologies, MediaPipe offers a myriad of advantages. Its pretrained models boast user-friendly interfaces and excel in real-time applications, particularly in scenarios such as webcam implementations. Furthermore, MediaPipe streamlines the development process by providing readily available, pre-built solutions, thereby significantly enhancing deployment efficiency and expediting the overall development lifecycle.

This project serves as a testament to the remarkable effectiveness of MediaPipe for hand gesture recognition tasks. Through its pretrained models and streamlined functionalities, MediaPipe offers a pragmatic and efficient solution that eclipses the traditional CNN methodologies in terms of ease of use, real-time performance, and deployment efficiency.

# Statistical Testing And Analysis

## Hypothesis Testing:-

**H0 (Null Hypothesis):** There is no significant difference in the accuracy of hand gesture recognition using the proposed method (MediaPipe framework with TensorFlow in OpenCV) compared to other existing methods.

**H1 (Alternative Hypothesis):** The proposed method (MediaPipe framework with TensorFlow in OpenCV) shows a significant improvement in the accuracy of hand gesture recognition compared to other existing methods.

To test the hypothesis, I am using one-way ANOVA and paired t-test.

ANOVA is chosen because there are multiple methods (more than two) for hand gesture recognition. One-way ANOVA is chosed because there is only one independent variable (the method).

Paired t-test is selected because I have implemented two models on the MediaPipe framework (a custom-trained CNN and a pre-trained model from MediaPipe) in my project.

**One way Anova:**

Total Variation (SS Total) = Sum of Squares Between (SSB) + Sum of Squares Within (SSW)

Where;

SS Total (Total Sum of Squares): Measures the total variability in the data.

SSB (Sum of Squares Between): Measures the variability between the group means.

SSW (Sum of Squares Within): Measures the variability within each group.

The degrees of freedom (df) for each component are:

df Total = N - 1 (where N is the total number of observations)

df Between = k - 1 (where k is the number of groups or levels of the independent variable)

df Within = N - k (where N is the total number of observations and k is the number of groups)

The mean squares (MS) for each component are calculated by dividing the sum of squares by its respective degrees of freedom:

MS Between = SSB / df Between

MS Within = SSW / df Within

Finally, the F-statistic is calculated as the ratio of the mean squares:

F = MSB/MSW

Here,

F = coefficient of ANOVA

MSB = Mean sum of squares between the groups

MSW = Mean sum of squares within groups

The F-statistic follows an F-distribution with degrees of freedom df Between and df within.

This F-statistic is then used to determine if there is a statistically significant difference between the group means. If the calculated F-value exceeds the critical F-value from the F-distribution for a chosen significance level (usually 0.05), then the null hypothesis (that there are no significant differences between group means) is rejected.

After performing calculations, I found that the calculated F-value exceeds the critical F-value from the F-distribution for a chosen significance level (0.05). Therefore, the null hypothesis is rejected.

**Paired t-test:**

I have implemented two models on the MediaPipe framework (a custom-trained CNN and a pre-trained model from MediaPipe) in my project.

So here,

Null Hypothesis (H0): There is no significant difference in the performance (accuracy, for example) between the custom-trained CNN model and the pre-trained model from MediaPipe when applied to hand gesture recognition using the MediaPipe framework.

Alternative Hypothesis (H1): There is a significant difference in the performance (accuracy) between the custom-trained CNN model and the pre-trained model from MediaPipe when applied to hand gesture recognition using the MediaPipe framework.

Formula:-

$$t = \frac{\sum d}{\sqrt{\frac{n(\sum d^2) - (\sum d)^2}{n-1}}}.$$

Here, $\sum d$ is the sum of the differences between paired observations.

and, n is the number of samples of pairs or samples.

Here, the calculated T-value exceeds the critical value from the t-distribution for the chosen significance level (0.05), the null hypothesis (H0) is rejected.

So, the alternative hypothesis (H1) is accepted. Therefore, it means, there is a significant difference in the performance (accuracy) between the custom-trained CNN model and the pre-trained model from MediaPipe when applied to hand gesture recognition using the MediaPipe framework.

Since the null hypotheses are rejected in both the one-way ANOVA and paired t-test, it indicates that there is evidence to support the alternative hypothesis. **Therefore, we can conclude that the proposed method (MediaPipe framework with TensorFlow in OpenCV) shows a significant improvement in the accuracy of hand gesture recognition compared to other existing methods.**

# MediaPipe: A Superior Choice

Using the pre-trained module from MediaPipe offers several advantages in hand gesture recognition projects. Firstly, the pre-trained model comes with the expertise and knowledge of Google's machine learning team, leveraging their extensive research and development in the field. This ensures a high level of accuracy and performance, as the model has been trained on large datasets and fine-tuned through rigorous testing. Additionally, the pre-trained module significantly reduces the need for extensive training and tuning, saving time and resources for developers. This makes it particularly advantageous for projects with limited resources or tight deadlines.

Furthermore, the MediaPipe framework itself provides a comprehensive and versatile platform for developing real-time computer vision and machine learning applications. It offers a wide range of pre-built solutions, including hand recognition, pose estimation, face detection, and object tracking, among others. This modular approach allows developers to quickly integrate complex functionalities into their projects without the need to build everything from scratch. Additionally, MediaPipe is open-source and cross-platform, making it accessible to a broad community of developers and researchers worldwide.

When considering the choice between MediaPipe and traditional machine learning approaches such as CNNs, RNNs, or SVMs for hand gesture recognition, several factors come into play. Firstly, MediaPipe provides a more streamlined and user-friendly solution, abstracting away many of the complexities involved in training and deploying custom models. This simplifies the development process and allows developers to focus more on the application logic rather than low-level implementation details. Additionally, MediaPipe models are optimized for real-time performance, making them well-suited for applications requiring fast inference speeds, such as augmented reality, robotics, and human-computer interaction.

Opting for the MediaPipe framework and leveraging its pre-trained modules offers numerous benefits for hand gesture recognition projects. From high accuracy and performance to ease of development and real-time capabilities, MediaPipe provides a robust and efficient solution for developers looking to integrate gesture recognition into their applications. By choosing MediaPipe, developers can expedite their development process, reduce development overhead, and ultimately deliver more sophisticated and immersive user experiences.

# Limitations of MediaPipe:

**Environment Dependency:** MediaPipe installation can sometimes be challenging due to dependencies and environment issues, leading to difficulties for users.

**System Requirements:** MediaPipe may not be compatible with very old systems and may require relatively modern hardware for optimal performance.

**Limited Platform Support:** While MediaPipe supports multiple platforms, it may not be as versatile as some other frameworks, limiting its usability in certain environments.

**Learning Curve:** MediaPipe's learning curve can be steep for beginners, particularly those without prior experience in computer vision or machine learning.

**Limited Pre-trained Models:** While MediaPipe offers some pre-trained models for common tasks, the selection may be limited compared to other frameworks, restricting the range of applications without additional training data and model development.

**Resource Intensive:** MediaPipe's real-time processing capabilities may consume significant computational resources, leading to high CPU/GPU utilization and potential performance bottlenecks on certain devices.

**Development Overhead:** Developing custom solutions in MediaPipe may require a deep understanding of computer vision algorithms and frameworks, leading to longer development cycles and higher development costs.

**Community Support:** MediaPipe's community support may be less extensive compared to more established frameworks like OpenCV or TensorFlow, making it harder to find solutions to specific problems or access community-contributed resources.

# Future Scope

The future scope of MediaPipe is vast and holds immense potential for further advancements and innovations in the fields of computer vision, machine learning, and human-computer interaction. By staying at the forefront of technology and collaborating with the developer community, MediaPipe can continue to empower creators and drive the next wave of transformative applications and experiences.

Below are some key areas where MediaPipe can continue to evolve and make significant contributions:

**Expansion of Pre-Built Solutions:** MediaPipe can expand its repertoire of pre-built solutions to cover a wider range of computer vision and machine learning tasks. This includes additional modules for object detection, semantic segmentation, activity recognition, and more. By providing a comprehensive suite of ready-to-use models and pipelines, MediaPipe can empower developers to build complex applications with minimal effort.

**Integration with Emerging Technologies:** MediaPipe can explore integration with emerging technologies such as augmented reality (AR), virtual reality (VR), and mixed reality (MR). By providing tools and frameworks tailored for AR/VR development, MediaPipe can facilitate the creation of immersive and interactive experiences that leverage hand gesture recognition, facial tracking, and object detection.

**Enhanced Performance and Efficiency:** Continuous optimization and improvement of MediaPipe's underlying algorithms and pipelines can lead to enhanced performance and efficiency. This includes optimizing inference speed, reducing resource consumption, and improving accuracy across different hardware platforms. By leveraging advancements in hardware acceleration and algorithmic optimizations, MediaPipe can deliver real-time performance for a wide range of applications.

**Support for Customization and Extension:** MediaPipe can enhance its support for customization and extension, allowing developers to tailor pre-built solutions to their specific needs. This includes providing APIs and tools for fine-tuning models, adding custom layers, and integrating external libraries. By fostering a vibrant ecosystem of third-party extensions and contributions, MediaPipe can cater to diverse use cases and domains.

**Integration with Cloud Services:** MediaPipe can explore integration with cloud services and platforms, enabling scalable deployment and management of computer vision pipelines. This includes support for distributed inference, model versioning, and automatic scaling based on workload demands. By leveraging cloud infrastructure, MediaPipe can address the needs of large-scale applications and facilitate collaboration among distributed teams.

**Advancements in Privacy and Security:** As privacy and security concerns become increasingly important, MediaPipe can prioritize advancements in data protection and privacy-preserving techniques. This includes implementing robust encryption, anonymization, and access control mechanisms to safeguard sensitive data. By adhering to industry best practices and compliance standards, MediaPipe can build trust among users and developers alike.

**Integration with Edge Devices and IoT:** MediaPipe can extend its reach to edge devices and Internet of Things (IoT) platforms, enabling gesture recognition and computer vision capabilities in resource-constrained environments. This includes optimizing models and pipelines for deployment on edge devices such as smartphones, wearables, and embedded systems. By providing lightweight and efficient solutions, MediaPipe can unlock new opportunities in smart homes, robotics, and wearable technology.

# References

[1] Haria, Aashni & Subramanian, Archanasri & Asokkumar, Nivedhitha & Poddar, Shristi & Nayak, Jyothi. (2017). Hand Gesture Recognition for Human Computer Interaction. Procedia Computer Science. 115. 367-374. 10.1016/j.procs.2017.09.092.

[2] Mohamed, Hazem Khaled & Sayed, Samir & Saad, Elsayed & Ali, Hossam. (2012). Hand gesture recognition using average background and logical heuristic equations. INTERNATIONAL JOURNAL OF COMPUTERS & TECHNOLOGY. 11. 2634-2640. 10.24297/ijct.v11i5.1149.

[3] Yun, Liu & Lifeng, Zhang & Shujun, Zhang. (2012). A Hand Gesture Recognition Method Based on Multi-Feature Fusion and Template Matching. Procedia Engineering. 29. 1678-1684. 10.1016/j.proeng.2012.01.194.

[4] Sharma, Ashish & Mittal, Anmol & Singh, Savitoj & Awatramani, Vasudev. (2020). Hand Gesture Recognition using Image Processing and Feature Extraction Techniques. Procedia Computer Science. 173. 181-190. 10.1016/j.procs.2020.06.022.

[5] Elsayed, Rania & Sayed, Mohammed & Abdalla, Mahmoud. (2017). Hand gesture recognition based on dimensionality reduction of histogram of oriented gradients. 119-122. 10.1109/JEC-ECC.2017.8305792.

[6] Bhuiyan, Rasel & Tushar, Abdul Kawsar & Ashiquzzaman, Akm & Shin, Jungpil & Islam, Dr. MD Rashedul. (2017). Reduction of gesture feature dimension for improving the hand gesture recognition performance of numerical sign language. 1-6. 10.1109/ICCITECHN.2017.8281833.

[7] Katkar, G., Ramteke, A., & Gaikwad, S. (2023). Hand Gesture Recognition Using Support Vector Machine Learning Algorithm. International Journal of Creative Research Thoughts (IJCRT), 2023, Manuscript ID: IJCRT2310227.

[8] Morimoto, Carlos & Yacoob, Yaser & Davis, Larry. (2012). Recognition of Hand Gestures Using Hidden Markov Models. Proceedings - International Conference on Pattern Recognition. 3. 10.1109/ICPR.2012.546990.

[9] Li, M., Wang, X., Xia, M., Cai, H., Gao, Y., & Cattani, C. (2012). Hidden-Markov-Models-Based Dynamic Hand Gesture Recognition. Mathematical Problems in Engineering, 2012, 986134. https://doi.org/10.1155/2012/986134

[10] Fortuna, L., Nguyen Huu, P., & Ngoc Phung, T. (2021). A hand gesture recognition algorithm using SVM and HOG model for control of robotic systems. Journal of Robotics, 2021, 3986497. doi:10.1155/2021/3986497.

[11] Kumar, Sunil & Srivastava, Tanu & Singh, Raj. (2017). Hand Gesture Recognition Using Principal Component Analysis. Asian Journal of Engineering and Applied Technology. 6. 2249-68. 10.51983/ajeat-2017.6.2.820.

[12] Arora, K., Suri, S., Arora, D., & Pandey, V. (2014). Hand Gesture Recognition Using Artificial Neural Network. International Journal of Computer Sciences and Engineering. Volume-2, Issue-4, E-ISSN: 2347-2693.

[13] Haroon, M.; Altaf, S.; Ahmad, S.; Zaindin, M.; Huda, S.; Iqbal, S. Hand Gesture Recognition with Symmetric Pattern under Diverse Illuminated Conditions Using Artificial Neural Network. Symmetry 2022, 14, 2045. https://doi.org/10.3390/sym14102045

[14] Nogales, R.E.; Benalcázar, M.E. Hand Gesture Recognition Using Automatic Feature Extraction and Deep Learning Algorithms with Memory. Big Data Cogn. Comput. 2023, 7, 102. https://doi.org/10.3390/bdcc7020102

[15] Chatterjee, K.; Raju, M.; Selvamuthukumaran, N.; Pramod, M.; Krishna Kumar, B.; Bandyopadhyay, A.; Mallik, S. HaCk: Hand Gesture Classification Using a Convolutional Neural Network and Generative Adversarial Network-Based Data Generation Model. Information 2022, 15, 85. https://doi.org/10.3390/info15020085

[16] Tsinganos, P.; Jansen, B.; Cornelis, J.; Skodras, A. Real-Time Analysis of Hand Gesture Recognition with Temporal Convolutional Networks. Sensors 2022, 22, 1694. https://doi.org/10.3390/s22051694

[17] M. Z., Hossain, Islam, M. S., Islam, R. U., & Andersson, K. (2019). Static Hand Gesture Recognition using Convolutional Neural Network with Data Augmentation. Expert Systems, 32(5), 563–577.

[18] Obaid, Falah & Babadi, Amin & Yoosofan, Ahmad. (2020). Hand Gesture Recognition in Video Sequences Using Deep Convolutional and Recurrent Neural Networks. Applied Computer Systems. 25. 57-61. 10.2478/acss-2020-0007.

[19] Toro-Ossaba, A.; Jaramillo-Tigreros, J.; Tejada, J.C.; Peña, A.; López-González, A.; Castanho, R.A. LSTM Recurrent Neural Network for Hand Gesture Recognition Using EMG Signals. Appl. Sci. 2022, 12, 9700. https://doi.org/10.3390/app12199700

[20] Koch P, Dreier M, Maass M, Bohme M, Phan H, Mertins A. A Recurrent Neural Network for Hand Gesture Recognition based on Accelerometer Data. Annu Int Conf IEEE Eng Med Biol Soc. 2019 Jul;2019:5088-5091. doi: 10.1109/EMBC.2019.8856844. PMID: 31947003.