# REAL-TIME SIGN LANGUAGE RECOGNITION USING DEEP LEARNING AND COMPUTER VISION

*Project Report*
*Submitted for the partial Fulfilment of the*
*Requirements for the  Award of the Degree of M.Sc.*
*Artificial Intelligence*

By

KARTHIKNAIR U

Under the Supervision of

ASST.PROF.JOHNSYMOL
JOY



Department of Computer Applications and Artificial Intelligence

SAINTGITS COLLEGE OF APPLIED SCIENCES,
PATHAMUTTOM,KOTTAYAM, KERALA

2022-2024

# SAINTGITS COLLEGE OF APPLIED SCIENCES, PATHAMUTTOM,KOTTAYAM,KERALA



## CERTIFICATE

This is certifying that the project report entitled **"REAL-TIME SIGN LANGUAGE RECOGNITION USING DEEP LEARNING AND COMPUTER VISION"** is a Bonafede report of the project work under taken by KARTHIKNAIR U(Reg:220011023344) fourth semester Artificial Intelligence student under my supervision and guidance, in partial fulfilment of the requirement for the award of the Degree of M.Sc. Artificial Intelligence of Mahatma Gandhi University, Kottayam during the period 2022- 2024.

Dr. Ambily Merlin Kuruvila                                    Asst. Prof. Johnsymol Joy

Head of the Department                                              Project Guide

Submitted for the viva-voice examination held on……………………………………………

Date:
External Examiner

2

# SAINTGITS COLLEGE OF APPLIED SCIENCES, PATHAMUTTOM,KOTTAYAM,KERALA

## DECLARATION

I KarthikNair U hereby declare that the project work entitled "REAL-TIME SIGN LANGUAGE RECOGNITION USING DEEP LEARNING AND COMPUTER VISION" is the original work done by me under the guidance of Asst. Prof. Johnsymol Joy for the partial fulfilment of the course M.Sc. Artificial Intelligence during the year 2022-2024.

I also declare that this report has been submitted by me fully for the award of degree before. Further, this is submitted on the partial fulfilment of the award of the degree of MSc Artificial intelligence of Mahatma Gandhi University, Kottayam.

# ACKNOWLEDGEMENT

First and foremost,I thankfully acknowledge my Principal Dr. Mathew Jacob for giving me an opportunity to present this project. The constant encouragement and timely support rendered by my Head of the Department, Dr. Ambily Merlin Kuruvila is deeply remembered. I express our heartfelt gratitude to my project guide, Asst. Prof. Johnsymol Joy PG Department of Computer Applications and AI, for her valuable guidance, support and encouragement during course of the project and preparation of the report.I have greatly benefited from her experience and knowledge.The help extended by all other staff members of the department is received with gratitude..I also remember with thanks to all my friends and well-wishes for their encouragement and support.

Above all, I would like to express our profound gratitude to God Almighty for his immense blessings upon me that led to the successful completion of the project.

KarthikNair U (Reg: 220011023344)

# CONTENTS

# 1. INTRODUCTION

## 1.1 ABSTRACT

Real-Time Sign Language Recognition Using Deep Learning and Computer Vision presents a comprehensive approach to developing a real-time sign language detection system using deep learning and computer vision techniques. The primary goal is to create a model capable of accurately recognizing and classifying various sign language gestures in real-time using a webcam. The process begins with the collection of a diverse dataset of sign language gestures, followed by training a convolutional neural network (CNN) using the MobileNetV2 architecture. Data augmentation techniques are employed to enhance the model's generalization capabilities. The trained model is then integrated into a real-time detection system that leverages the MediaPipe library for hand tracking and landmark detection. The system processes video input from a webcam, detects hand landmarks, and uses the trained model to predict the corresponding gesture. The effectiveness of the system is demonstrated through high accuracy in gesture recognition and smooth real-time performance. This project lays the foundation for more advanced sign language recognition applications, with potential future work including the expansion of the gesture vocabulary and optimization for deployment on mobile and embedded devices.

## 1.2 INTRODUCTION

Sign language serves as a primary means of communication for millions of individuals with hearing impairments worldwide. It enables these individuals to express themselves and interact with others effectively. However, a significant communication gap persists between sign language users and those who do not understand it, often leading to misunderstandings and social isolation for the hearing impaired. Bridging this gap is crucial for fostering inclusive and enhancing the quality of life for these individuals.

The project titled "Real-Time Sign Language Recognition Using Deep Learning and Computer Vision" aims to address this communication barrier by developing an advanced system capable of recognizing and translating sign language gestures into text in real-time. This system leverages the power of deep learning and computer vision, harnessing recent technological advancements to create a tool that can facilitate smoother interactions between sign language users and non-users.

The project is structured around three primary stages:

*1. Data Collection*: The initial stage involves gathering a diverse and comprehensive dataset of sign language gestures. This dataset forms the foundation of the project, as the accuracy and reliability of the recognition system heavily depend on the quality and diversity of the training data. The dataset includes various common gestures, ensuring that the model can recognize a wide range of signs.

*2. Model Training:* In the second stage, a convolutional neural network (CNN) is trained using the collected dataset. The MobileNetV2 architecture is chosen for its balance of efficiency and accuracy, making it well-suited for real-time applications. Data augmentation techniques, such as rotation, zoom, and horizontal flipping, are employed to enhance the model's ability to generalize across different conditions and users.

*3. Real-Time Detection:* The final stage integrates the trained model into a real-time detection system. This system uses a standard webcam to capture video input, which is processed to detect and recognize hand gestures. The MediaPipe library is utilized for hand tracking and

landmark detection, providing precise and reliable input to the model. The recognized gestures are then translated into text and displayed on the screen, enabling immediate understanding for non-sign language users.

The development of this system involves several key components and technologies:

1. *Computer Vision:* Techniques for hand detection and landmark identification are essential for accurately capturing the gestures in real-time. MediaPipe, an open-source framework developed by Google, provides robust tools for real-time hand tracking, ensuring that the system can operate smoothly and accurately under various conditions.

2. *Deep Learning*: The core of the gesture recognition system is the deep learning model, specifically a CNN based on the MobileNetV2 architecture. This model is trained to identify specific patterns and features within the hand gestures, enabling accurate classification of the signs.

3. *Data Augmentation:* To improve the model's robustness and ability to generalize, data augmentation techniques are applied. These techniques help create a more diverse training set by artificially increasing the variety of images, simulating different angles, lighting conditions, and backgrounds.

The effectiveness of the system is evaluated through rigorous testing, achieving high accuracy in gesture recognition and demonstrating smooth real-time performance. This project showcases the potential of deep learning and computer vision to create practical solutions for real-world problems, particularly in enhancing communication for the hearing impaired.

By bridging the communication gap, this project aims to make a significant positive impact on the lives of sign language users. Future work could involve expanding the gesture vocabulary to cover more signs, optimizing the system for deployment on mobile and embedded devices, and improving the robustness of the model to handle different environmental conditions and user variations. Ultimately, this project lays the groundwork for more advanced sign language recognition applications, contributing to a more inclusive and connected society.

## 1.3 OBJECTIVE AND SCOPE

The primary objective of this project, "Real-Time Sign Language Recognition Using Deep Learning and Computer Vision," is to develop an efficient and accurate system capable of recognizing and translating sign language gestures into text in real-time. This objective can be broken down into several specific goals:

*1. Data Collection:* To gather a comprehensive and diverse dataset of sign language gestures that includes multiple examples of each gesture under various conditions. This dataset forms the foundation for training a robust deep learning model.

*2.Model Development:* To design and train a convolutional neural network (CNN) using the MobileNetV2 architecture, which is known for its efficiency and accuracy, to accurately classify sign language gestures based on the collected dataset.

*3.Real-Time Processing*: To integrate the trained model with a real-time video feed from a webcam, enabling the system to detect and recognize hand gestures as they occur.

*4.User Interface:* To create a user-friendly interface that displays the recognized gestures in text form, making it easy for non-sign language users to understand the gestures.

*5.Performance Evaluation:* To evaluate the system's accuracy and efficiency in real-time scenarios and ensure it performs well under various conditions, such as different lighting and backgrounds.

*6.Future Expansion:* To lay the groundwork for future improvements, such as expanding the gesture vocabulary, optimizing for mobile and embedded devices, and enhancing the system's robustness to different users and environments.

The scope of this project encompasses several key aspects, each contributing to the overall goal of creating a reliable and practical sign language recognition system:

*1.Gesture Selection:* Focus on a specific set of commonly used sign language gestures that are essential for basic communication. The initial dataset includes gestures like "I Love You," "Victory," "Okay," "I Dislike It," "Goodbye," "Yes," "Live Long," and "Stop."

*2.Data Collection and Preprocessing:* Capture images of each gesture using a webcam, ensuring a variety of examples are collected to account for different angles, lighting conditions, and hand positions. Preprocess the images to standardize their size and normalize pixel values.

*3.Model Training:* Utilize the MobileNetV2 architecture to train a CNN on the collected dataset. Apply data augmentation techniques to enhance the model's ability to generalize and improve its robustness.

*4. Real-Time Detection:* Develop a system that integrates the trained model with a webcam feed, using the MediaPipe library for hand detection and landmark tracking. Ensure the system can process video input in real-time and accurately predict gestures.

*5.User Interface Development:* Create a graphical interface that displays the recognized gestures in text form, providing immediate feedback to users.

*6.Testing and Validation:* Conduct extensive testing to evaluate the system's accuracy and performance in real-time scenarios. Test under various conditions to ensure the system's reliability and robustness.

*7.Documentation and Reporting*: Document the entire development process, including data collection methods, model architecture, training procedures, and real-time implementation. Prepare a comprehensive project report detailing the methodology, results, and future work recommendations.

*8.Future Work:* Identify potential areas for improvement and expansion, such as increasing the number of recognizable gestures, optimizing the system for deployment on mobile devices, and improving the robustness to handle different users and environmental conditions.

By addressing these aspects, the project aims to create a functional and reliable sign language recognition system that can significantly enhance communication for sign language users and contribute to a more inclusive society.

## 1.4 PROBLEM STATEMENT

Despite the widespread use of sign language as a primary mode of communication for individuals with hearing impairments, there exists a significant communication barrier between sign language users and non-users. Non-users often struggle to understand sign language gestures, leading to misunderstandings and limited interaction opportunities for the hearing impaired.

The problem addressed by this project is the lack of efficient and accessible tools for real-time sign language recognition and translation. Existing solutions may be limited in accuracy, speed, or usability, making them less effective for practical use in everyday communication scenarios.

The project aims to develop a robust and efficient system capable of recognizing and translating sign language gestures into text in real-time, using a standard webcam. The system should be able to accurately identify a diverse set of common sign language gestures and display the recognized gestures in text form, enabling immediate understanding for non-sign language users.

Key challenges to be addressed include:

*1. Accuracy:* Ensuring that the system can accurately recognize a wide range of sign language gestures under various conditions, including different hand positions, lighting conditions, and backgrounds.

*2.Real-Time Performance*: Developing a system that can process video input from a webcam in real-time, with minimal delay between gesture recognition and display.

*3.User-Friendliness:* Creating an intuitive and user-friendly interface that allows non-sign language users to easily understand and interact with the recognized gestures.

*4.Robustness:* Enhancing the system's robustness to handle variations in hand shapes, sizes, and movements, as well as different users' signing styles and speeds.

*5.Scalability*: Designing the system to be scalable and adaptable, allowing for future expansion to support additional sign language gestures and optimization for deployment on mobile and embedded devices.

By addressing these challenges, the project aims to provide a practical and accessible solution for real-time sign language recognition, ultimately improving communication and interaction opportunities for individuals with hearing impairments.

# 2. SYSTEM CONFIGURATION

## 2.1 HARDWARE SPECIFICATION

- **Webcam**: A standard webcam with at least 720p resolution capable of capturing high-quality video input for real-time gesture detection.

  - CPU: Intel Core i5 or equivalent, or higher
  - GPU: NVIDIA GeForce GTX 1060 or equivalent, or higher (for GPU acceleration)
  - RAM: 8 GB or higher
  - Storage: SSD recommended for faster data access

- **Display**: A monitor or display screen for viewing the real-time output of the gesture recognition system.

## 2.2 SOFTWARE SPECIFICATION

- **Operating System**: The system should be compatible with commonly used operating systems, such as:

  - Windows 10 or higher
  - macOS 10.13 or higher
  - Ubuntu 18.04 LTS or higher

- **Software Dependencies**:

  - Python 3.x: The system relies on Python for implementing the deep learning model, real-time processing, and interface development.
  - TensorFlow: TensorFlow library is used for building and training the deep learning model.
  - OpenCV: OpenCV library is used for image processing, video capture, and displaying real-time output.
  - MediaPipe: MediaPipe library is utilized for hand detection and landmark tracking.
  - Matplotlib: Matplotlib library is used for visualization and plotting training/validation metrics.
  - Visual Studio Code (VS Code): Visual Studio Code is a lightweight yet powerful code editor with built-in support for Python development.

**2.3 ABOUT DEVELOPMENT TOOL**

The development of the real-time sign language recognition system relies on a combination of software tools and libraries that facilitate different aspects of the project, including model development, data preprocessing, real-time processing, and user interface design. Below are some essential development tools used in this project:

**1. Python Programming Language:**

Python serves as the backbone of the project due to its simplicity, readability, and extensive support for libraries and frameworks relevant to machine learning and computer vision. Its syntax is intuitive and easy to understand, making it ideal for rapid prototyping and development.

**2. TensorFlow:**

TensorFlow is a widely-used open-source machine learning framework that provides high-level APIs for building and training neural networks. In this project, TensorFlow is crucial for developing and training the deep learning model responsible for recognizing sign language gestures. Its flexibility and scalability make it suitable for handling complex model architectures and large datasets.

**3. OpenCV (Open Source Computer Vision Library):**

OpenCV is a powerful library for computer vision tasks, offering a wide range of functionalities for image and video processing. It plays a pivotal role in processing real-time video input from a webcam, performing tasks such as frame capture, image preprocessing, feature detection, and object tracking. OpenCV's efficient algorithms and extensive documentation make it indispensable for implementing real-time computer vision applications.

**4. MediaPipe:**

MediaPipe, developed by Google, provides pre-built solutions for various multimedia processing tasks, including hand tracking and pose estimation. In this project, MediaPipe is used for hand detection and landmark tracking, enabling precise localization of hand positions in real-time video streams. Its optimized algorithms and ease of integration make it an excellent choice for building real-time hand gesture recognition systems.

**5. Matplotlib:**

Matplotlib is a comprehensive plotting library for Python, offering a wide range of visualization options for analyzing data and displaying results. In this project, Matplotlib is utilized for visualizing training/validation metrics, such as accuracy and loss, during the model training process. Its customizable plots and interactive features facilitate insightful analysis of model performance and training dynamics.

**6. Visual Studio Code (VS Code):**

Visual Studio Code is a lightweight yet powerful code editor with built-in support for Python development. It offers features such as syntax highlighting, code completion, debugging, and version control integration, enhancing the productivity of developers. VS Code's extensibility through plugins and its intuitive user interface make it a preferred choice for developing machine learning and computer vision applications.

By leveraging these development tools, developers can efficiently implement, train, and deploy the real-time sign language recognition system, ensuring robust performance and user-friendly interaction. Each tool contributes unique capabilities and advantages, enabling the seamless integration of machine learning and computer vision techniques to address real-world challenges.

# 3. LITERATURE REVIEW

# 3.1 INTRODUCTION

Sign language serves as a crucial means of communication for individuals with hearing impairments, offering a rich and expressive mode of expression. However, the communication gap between sign language users and non-users remains a significant challenge in various social, educational, and professional settings. Recognizing the importance of bridging this gap, researchers and technologists have explored diverse methodologies and techniques to develop sign language recognition systems capable of translating sign gestures into text or speech in real-time.

This literature review aims to provide a comprehensive overview of the existing research landscape in the field of sign language recognition. By examining both traditional and cutting-edge approaches, this review seeks to identify key findings, challenges, and future directions in the domain. Through synthesizing insights from a wide range of studies, this review aims to inform the development of more accurate, efficient, and accessible sign language recognition systems.

The review begins by delving into traditional approaches to sign language recognition, which typically rely on handcrafted features and conventional machine learning algorithms. While these approaches have laid the groundwork for sign language recognition, they often face limitations in terms of robustness and scalability. Subsequently, the review explores the emergence of deep learning-based techniques, which have revolutionized the field by leveraging neural networks to automatically learn discriminative features from raw input data.Key findings from research in this area are examined, including the adoption of advanced convolutional neural network (CNN) architectures, innovative data augmentation strategies, real-time processing methodologies, and approaches for expanding the gesture vocabulary of recognition systems. Furthermore, the review sheds light on the challenges confronting sign language recognition, such as data scarcity, robustness issues, and real-time performance constraints.

By synthesizing insights from diverse studies and identifying areas for future exploration, this literature review aims to contribute to the ongoing advancement of sign language recognition technology. Ultimately, the development of more advanced and effective sign language recognition systems has the potential to enhance communication, foster inclusivity, and empower individuals with hearing impairments to participate fully in various aspects of society.

## 3.2 MODELS

**Deep Dive into the Models: Convolutional Neural Networks (CNNs), MobileNetV2 and MediaPipe Hands**

This project utilizes two models: a Convolutional Neural Network (CNN) for sign language detection, MobileNetV2 as a pre-trained model for feature extraction and MediaPipe Hands for hand landmark detection. Here's a detailed explanation with visualizations:

**1. Convolutional Neural Network (CNN) for Sign Language Detection**

*Concept:*CNNs are a type of deep learning architecture particularly well-suited for image recognition tasks. They achieve this through convolutional layers that process images using learnable filters. These filters progressively extract features from the image, such as edges, shapes, and textures, that are crucial for classification.

*Architecture:*The script suggests using a pre-trained MobileNetV2 model as the base architecture. Here's a simplified breakdown:

*MobileNetV2:* This pre-trained model consists of several convolutional layers followed by pooling layers for dimensionality reduction. It also includes depth wise separable convolutions for efficiency. These layers are not retrained in this project and serve as feature extractors.

*Custom Layers:* On top of the pre-trained MobileNetV2, the script adds custom layers to specialize in sign language classification:

*Global Average Pooling Layer:* Averages the spatial dimensions (width and height) of the feature maps from the pre-trained model, resulting in a single feature vector for each channel.

*Dense Layers:* Fully-connected layers with learnable weights and biases. These layers process the features from the pooling layer and learn higher-level representations specific to sign language gestures. The script mentions ReLU (Rectified Linear Unit) activation functions for these layers, which introduce non-linearity for improved learning.

*Output Layer:* The final layer has a number of neurons equal to the number of sign classes you want to recognize. It likely uses softmax activation to generate probability scores for each class, enabling multi-class classification (e.g., classifying between "I love you" and "Victory").

*Equations (Simplified):*

*Convolutional Layer:* Applies a filter (represented by learnable weights) to the input image, producing a feature map. Mathematically, it can be expressed as:

$$Output[i, j, c] = \Sigma (Input[m, n, d] * Filter[i - m, j - n, d, c]) + Bias[i, j, c]$$

*Pooling Layer:* Reduces the dimensionality of the feature maps. Common pooling functions include average pooling and max pooling. Average pooling takes the average of a specific region in the feature map, while max pooling selects the maximum value.

*Dense Layer:* Performs a linear transformation on the input vector using its weights and biases. Mathematically:

$$Output = Input * Weights + Biases$$

*Softmax Activation (Output Layer):* Converts the output logits from the dense layer into probability scores for each class. The sum of these probabilities equals 1.

$$Softmax(x\_i) = e^x\_i / \Sigma(e^x\_j) \text{ for all classes } j$$

## 2. MobileNetV2 Architecture:

MobileNetV2 is a pre-trained convolutional neural network known for its efficiency in terms of computational cost and memory usage. This efficiency makes it well-suited for deployment on mobile and embedded devices. Here's a simplified view of its key components:

*Inverted Residual Blocks:* These are the building blocks of MobileNetV2. Unlike traditional residual blocks, they use a depthwise separable convolution approach to achieve a significant reduction in parameters and computations.

*Depthwise Convolution:* This layer applies a filter to each input channel independently, extracting features without increasing the number of channels. Mathematically, for an input feature map X of size (height, width, number of input channels) and a depthwise filter F of size (height, width, 1), the output Y can be expressed as:

$$Y[i, j, c] = \Sigma (X[i + k, j + l, c] * F[k, l])$$

*Pointwise Convolution (Linear Projection):* This layer uses a 1x1 convolution to project the output of the depthwise convolution to a new channel dimension. It effectively reduces the number of channels while preserving the spatial information.

*Bottleneck Layers:* These are optional layers placed within some inverted residual blocks. They further reduce the number of channels before and after the depthwise convolution, leading to even lower computational cost.

*Residual Connections:* Similar to traditional residual blocks, MobileNetV2 utilizes residual connections to skip some layers and add the original input to the output, facilitating the flow of gradients and improving model performance.

*Equations (Simplified):*

The core mathematical operations within MobileNetV2 involve depthwise separable convolutions and pointwise convolutions. The equations provided earlier for the depthwise convolution demonstrate the element-wise multiplication between the input feature map and the filter, resulting in a new feature map.

*Role in Sign Language Detection:* In this project, MobileNetV2 acts as a feature extractor. It takes pre-processed image frames (likely resized and normalized) as input and passes them through its convolutional layers. These layers progressively extract features from the image, such as edges, shapes, and textures relevant for sign language gestures. The final output of MobileNetV2 represents a compressed feature vector that captures the essential information from the input image.

**Benefits of MobileNetV2:**

*Efficiency:* MobileNetV2 requires fewer parameters and computations compared to traditional CNN architectures, making it suitable for real-time applications on mobile devices with limited resources.

*Transfer Learning:* By leveraging pre-trained weights from MobileNetV2, the script focuses training on the custom layers added on top, which are specifically designed for sign language classification. This reduces training time and improves performance.

By employing MobileNetV2 for feature extraction, the project achieves a balance between accuracy and efficiency in real-time sign language detection.

## 3. MediaPipe Hands

*Concept:*MediaPipe Hands is a pre-built solution from Google's MediaPipe framework. It's not a deep learning model, but a machine learning solution that utilizes non-deep learning techniques for real-time hand pose estimation. It detects hands in a frame and identifies specific landmark points on each finger (e.g., tip, base joint).

*Functionality:*MediaPipe Hands takes an RGB image frame as input.It employs image processing techniques, likely including background subtraction, skin color segmentation, and hand shape analysis, to locate hand regions in the frame.Once a hand is detected, it employs a hand keypoint detection model to identify specific landmark points on each finger (typically 21 points).

# 4. METHODOLOGY

# 4.1 INTRODUCTION

The methodology section outlines the process and techniques used to develop the sign language detection model. This involves three main phases: data collection, model development, and real-time detection. Each phase is crucial for the overall performance and functionality of the system.

## 1. Data Collection

The data collection phase involves capturing images of various hand gestures that represent different sign language symbols. The `data_collection.py` script is used to facilitate this process. Here's a detailed breakdown of the data collection methodology:

*Gesture Labels and Directory Structure:* A dictionary maps gesture labels to corresponding directory names where the images will be stored. The gestures included are 'I Love You', 'Victory', 'Okay', 'I Dislike It', 'Goodbye', 'Yes', 'Live Long', and 'Stop'.

*Capture Images:* For each gesture, a specified number of images are captured using a webcam. The script initializes the webcam and creates a directory for each gesture if it doesn't already exist.

*Image Capture Process:* A delay is provided to allow the user to get ready, after which images are captured and stored in the corresponding directory. The user can interrupt the capture process by pressing the 'q' key.

*Normalization and Storage:* The captured images are normalized and resized to 224x224 pixels to ensure uniformity and compatibility with the neural network input requirements.

## 2. Model Development

The model development phase involves building and training a convolutional neural network (CNN) to recognize hand gestures. The `sign_language_detection.py` script is used for this purpose. Here's a detailed breakdown of the model development methodology:

*Load and Preprocess Data:* The script loads images from the dataset directory, resizes them to 224x224 pixels, and normalizes the pixel values. The labels are one-hot encoded for use in the neural network.

*Data Augmentation:* To improve the robustness and generalization of the model, data augmentation techniques such as rotation, width and height shifts, shearing, zooming, and horizontal flipping are applied.

*Model Architecture:* The base model used is MobileNetV2, pre-trained on the ImageNet dataset and modified to exclude the top layer. Custom layers are added on top of the base model:

*GlobalAveragePooling2D:* Reduces each feature map to a single value.

*Dense Layer:* Fully connected layer with 1024 units and ReLU activation.

*Output Layer:* Dense layer with softmax activation to classify the gestures into one of the predefined categories.

*Compile the Model:* The model is compiled using the Adam optimizer with a learning rate of 0.001, categorical cross-entropy loss, and accuracy as the evaluation metric.

*Train the Model*: The model is trained using the augmented training data for 10 epochs, with a batch size of 32. The training process includes validation to monitor the model's performance on unseen data.

*Evaluate the Model:* After training, the model is evaluated on the test set to determine its accuracy and loss. Training and validation accuracy and loss are plotted to visualize the model's performance over epochs.

## 3. Real-time Detection

The real-time detection phase involves using the trained model to recognize hand gestures from live webcam feed. The `realtime_detection.py` script is used for this purpose. Here's a detailed breakdown of the real-time detection methodology:

*Load Trained Model:* The pre-trained model is loaded for use in real-time detection.

*Initialize MediaPipe Hands:* MediaPipe Hands, a framework for hand tracking, is initialized to detect hand landmarks in the video feed.

*Preprocess Frame:* Each frame from the webcam feed is preprocessed by resizing it to 224x224 pixels and normalizing the pixel values.

*Hand Landmark Detection:* MediaPipe Hands processes each frame to detect hand landmarks. If hands are detected, landmarks are drawn on the frame, and custom finger lines are added to enhance visualization.

*Gesture Recognition:* The preprocessed frame is fed into the trained model to predict the gesture. The predicted gesture label is displayed on the frame.

*Display and User Interaction*: The processed frame with detected landmarks and predicted gesture is displayed in a window. The user can exit the real-time detection by pressing the 'q' key.

This methodology ensures a systematic approach to building a robust sign language detection system, from data collection to real-time application.


## 4.2  DATASET

The dataset used for training the sign language detection model consists of images representing various hand gestures. This section provides an overview of the dataset, including the gestures captured, the data collection process, and the structure and characteristics of the dataset.

The dataset includes eight distinct gestures, each corresponding to a specific sign language symbol. The gestures and their respective labels are as follows:

1. I Love You (`i_love_you`)

2. Victory (`victory`)

3. Okay (`okay`)

4. I Dislike It (`i_dislike_it`)

5. Goodbye (`goodbye`)

6. Yes(`yes`)

7. Live Long (`live_long`)

8. Stop (`stop`)

Each gesture is mapped to a directory where the corresponding images are stored.

Dataset Summary

1. Number of Gestures: 8
2. Number of Images per Gesture: 100
3. Total Number of Images: 800
4. Image Dimensions: 224x224 pixels
5. Image Format: JPEG

The dataset is organized into directories, each corresponding to a specific gesture. Here is an example of the directory structure:
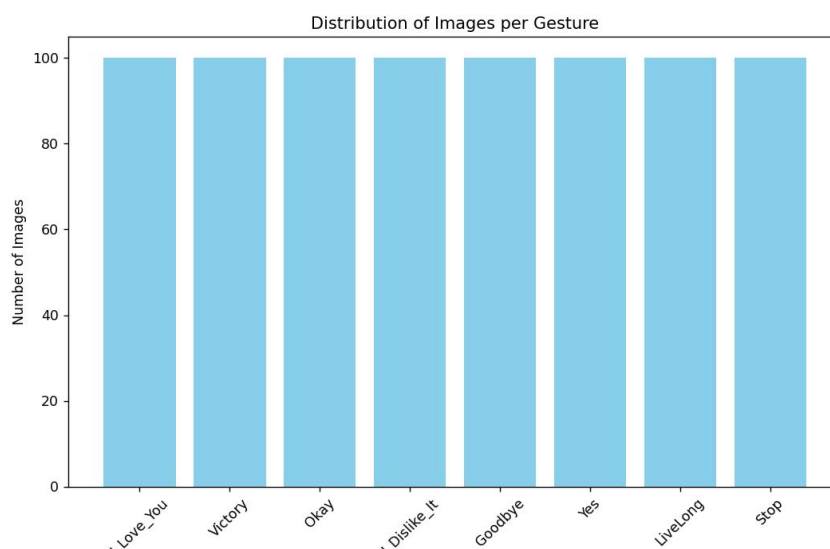


Figure 1: Dataset directory



Figure 2: Python generated histogram for the dataset

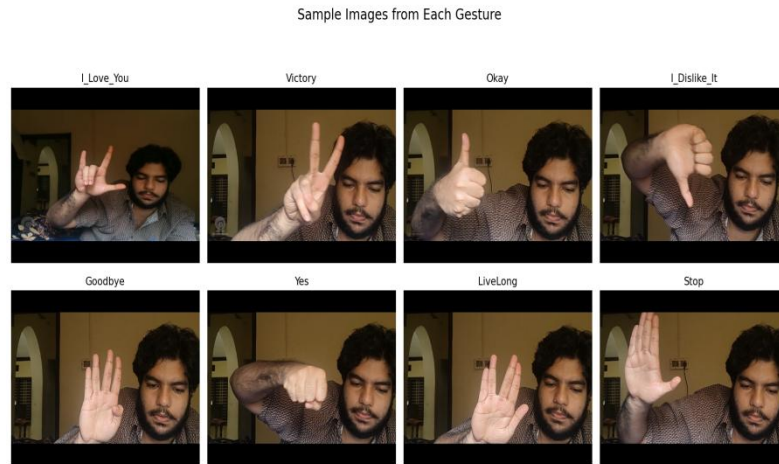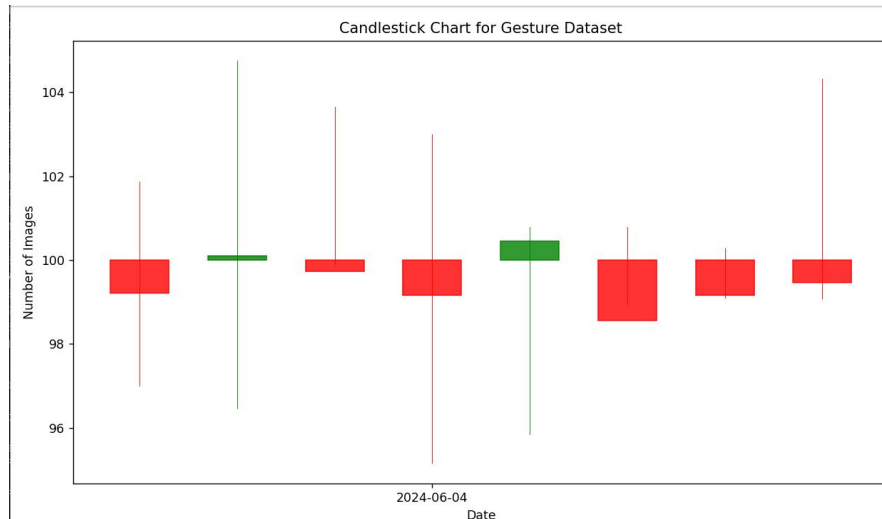Figure 3: Sample image from each Gestures



Figure 4: Candlestick Graph

## 4.3 Convolutional Neural Networks(CNNs)

Convolutional Neural Network is a deep learning technique which is developed from the inspiration of visual cortex which are the fundamental blocks of human vision. It is observed from the research that, the human brain performs a large-scale convolutions to process the visual signals received by eyes, based on this observation CNNs are constructed and observed to be outperforming all the prominent classification techniques Two major

operations performed in CNN are convolution (wT∗ X) and pooling(max()) and these blocks are wired in a highly complex fashion to mimic the human brain .The neural network is constructed in layers, where the increase in the number of layers increases the network complexity and is observed to improve the system accuracy .The CNN architecture consists of three operational blocks which are connected as a complex architecture.

The functional blocks of Convolutional Neural Network:

1.Input Layer

2.Convolutional Layer

3.Max Pooling layer

4.Activation Layer (ReLU)
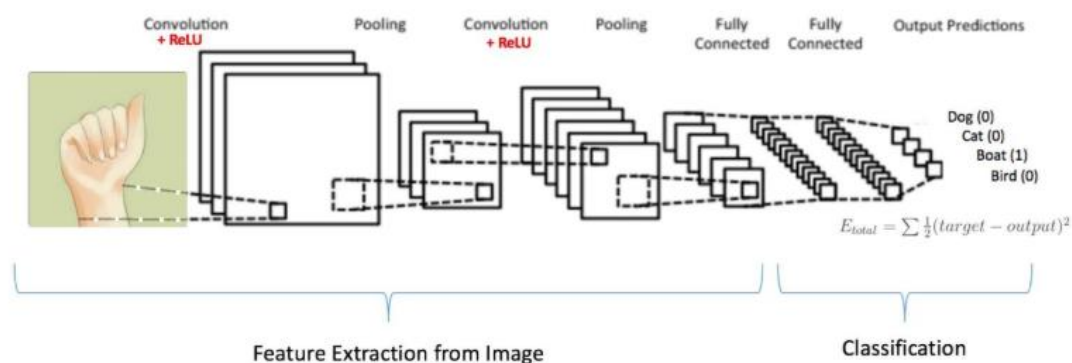
5.Fully-Connected layer

6.Output Layer



Figure 5:Working of CNN

**1. Input Layer:**

Imagine a video frame of a hand signing a letter. This frame is converted into a digital image represented by a grid of pixels. Each pixel has a value corresponding to its color intensity. This 2D array of pixel values forms the input layer of the CNN.

*Equation (Reshaping Input for some cases):*In some cases, depending on the framework or library used, the input layer might require reshaping the image data into a specific format. Here's an example equation for reshaping a grayscale image (width x height) into a format suitable for a specific CNN architecture (channels x width x height):

*Reshaped_Input = np.reshape(Original_Image, (1, width, height))*

Where:

31

- Original_Image is the 2D grayscale image data.
- Reshaped_Input is the reshaped data with one channel added at the beginning (becomes a 3D tensor).

**2. Convolutional Layers:**

The core of the CNN architecture. These layers use filters (small matrices) to learn patterns directly from the image data. The filter slides across the input image, performing element-wise multiplication between the filter weights and the corresponding pixel values. This generates a feature map that captures specific features like edges, corners, or shapes relevant for sign language recognition.

*Equation for Convolution:*Here, X represents the input image, W represents the filter weights, and b represents the bias term:

$$Output = ReLU(W * X + b)$$

The ReLU (Rectified Linear Unit) activation function introduces non-linearity, allowing the network to learn complex features.
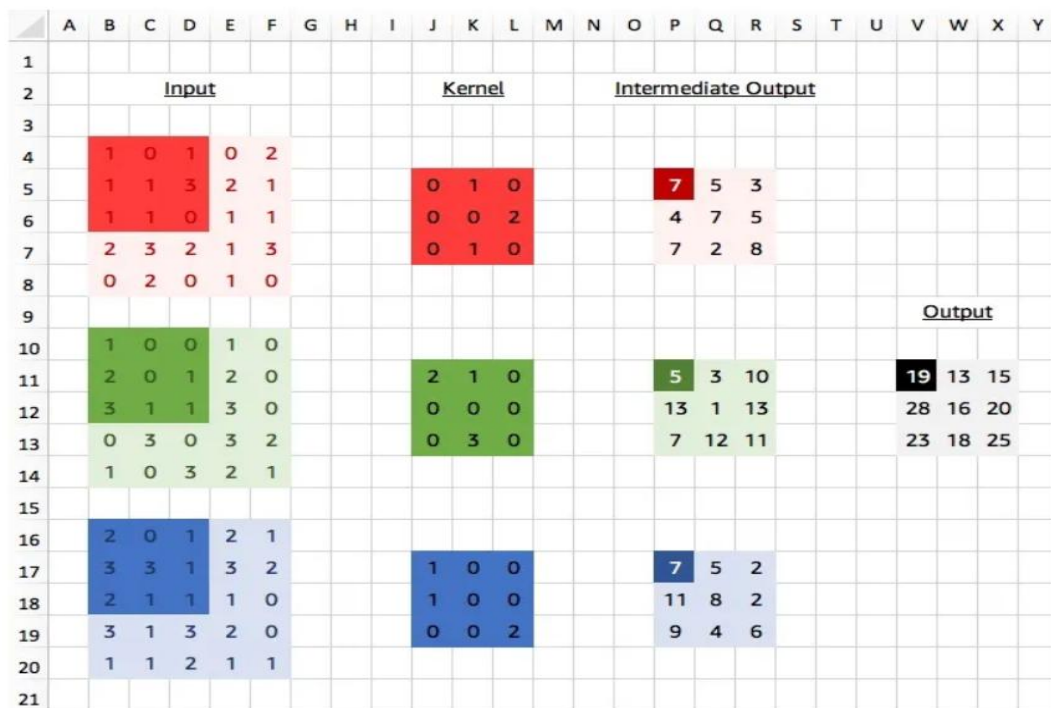


Figure 6: A 2D Convolution with a 3x3 kernel applied to an 3 channel RGB input of size 5x5 to give output of 3x3.

**3. Pooling Layers:**

These layers down sample the feature maps, reducing their dimensionality and computational cost. Techniques like max pooling select the maximum value from a sub-region of the feature map, summarizing the presence of a feature. This helps the network focus on important aspects and reduces sensitivity to small variations in the hand posture.

*Pooling Operations:*There are several common pooling operations used in CNNs:

*Max Pooling:* Takes the maximum value from a rectangular region of the input feature map. This captures the most dominant feature within that region.

*Average Pooling:* Calculates the average value from a rectangular region of the input feature map. This provides a summary of the average intensity within that region.

*Equation (Max Pooling Example):*

$$Output(x, y) = max(Input(x + i, y + j)) \quad \textit{for all i, j within the pooling window}$$

Where:

- (x, y) is the output location
- Input(x + i, y + j) is the corresponding input value within the pooling window
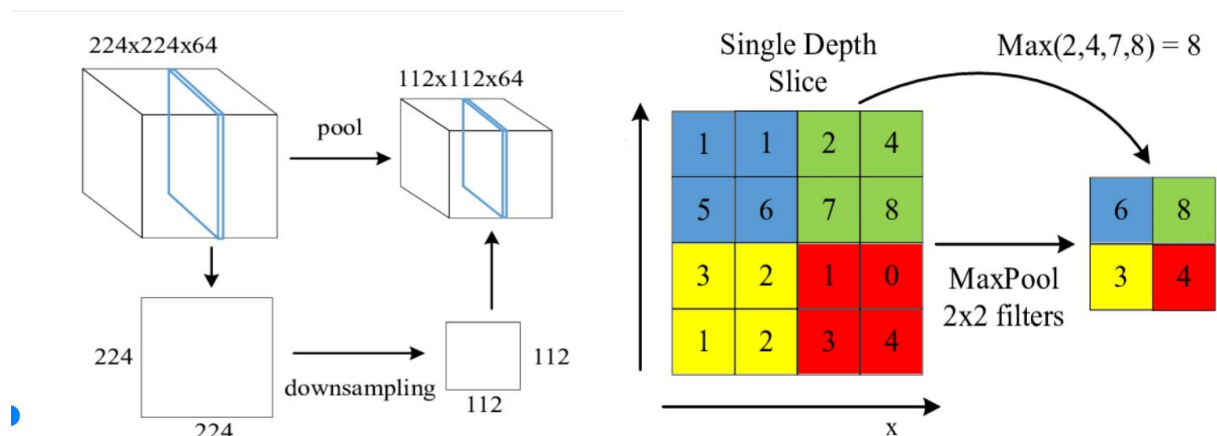


Figure 7:Max Pooling operation on feature map (2×2 window)

## 4.Activation Layer (ReLU)

The activation layer, specifically the ReLU (Rectified Linear Unit) function, plays a vital role in Convolutional Neural Networks (CNNs) designed for sign language detection.The ReLU function introduces non-linearity into the network. In simpler terms, it allows the network to

33

learn more complex patterns in the data compared to linear activation functions. This is crucial for sign language detection, where hand gestures can have intricate variations.

The ReLU function operates element-wise on the input values. Mathematically, it can be represented as:

$$ReLU(x) = max(0, x)$$

Where:

x is the input value to the activation layer (often the output from the previous convolutional or pooling layer).

ReLU(x) is the output value after applying the ReLU function.

*Effect of ReLU*: ReLU outputs the input value itself if it's positive (x > 0),ReLU outputs zero if the input value is non-positive (x <= 0).
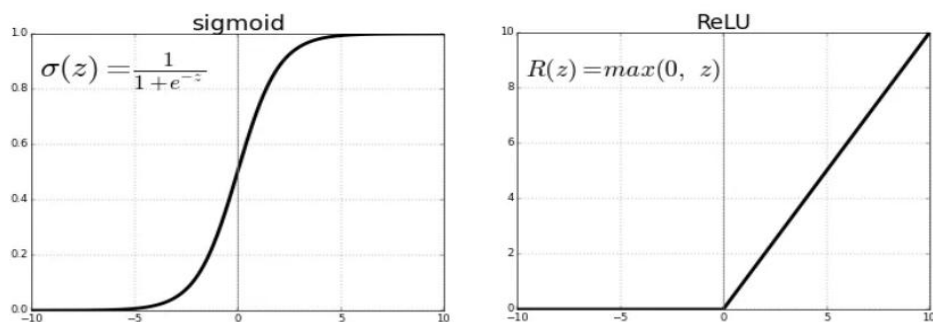


Figure 8: ReLU v/s Logistic Sigmoid

## 5. Fully Connected Layers:

After several convolutional and pooling layers, the network uses fully connected layers to classify the extracted features. These layers are similar to traditional artificial neural networks, where each neuron in a layer is connected to all neurons in the previous layer. Here, the neurons learn to combine the features from the previous layers to identify the specific sign language being presented.The FC layer takes the output from the previous layers (typically a flattened feature map) and performs classification. It connects all the neurons from the previous layer to each neuron in the output layer. This dense connectivity allows the network to learn complex relationships between the extracted features and the corresponding sign classes.

34

*Flattening:* The output from the previous convolutional or pooling layer (usually a 3D tensor) is reshaped into a 1D vector. This process essentially combines all the feature maps into a single long list of values.

*Matrix Multiplication:* This vector is then multiplied by a weight matrix and a bias vector. The weight matrix contains learnable weights that capture the relationships between the features and the output classes. The bias vector adjusts the activation levels of the neurons in the FC layer.

*Activation Function:* Finally, an activation function (often softmax) is applied to the output of the matrix multiplication. The softmax function converts the output values into probabilities, where each element represents the probability of the input image belonging to a specific sign class.

*Equation (Simplified):*

$$Output = Softmax( W * Flattened\_Input + bias )$$

- Output: A vector containing probability scores for each sign class.
- Softmax: The activation function used for multi-class classification.
- W: The weight matrix of the FC layer.
- Flattened_Input: The 1D vector obtained by flattening the previous layer's output.
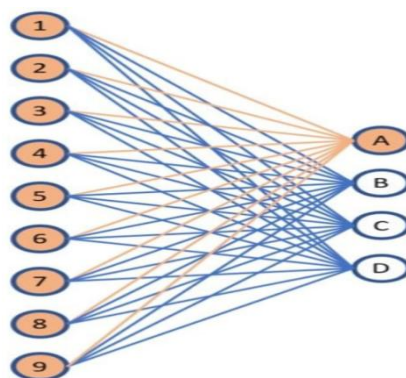- bias: The bias vector of the FC layer.



Figure 9:Illustration of a fully connected layer.

**6. Output Layer:**

The final layer has one neuron for each signs. The network assigns probabilities to each sign based on the processed features. The sign with the highest probability is recognized as the output.

## 4.4 MobileNetV2

MobileNetV2 is a pre-trained convolutional neural network (CNN) architecture known for its efficiency in terms of computational cost and memory usage. This makes it particularly well-suited for deployment on mobile devices and embedded systems with limited resources.
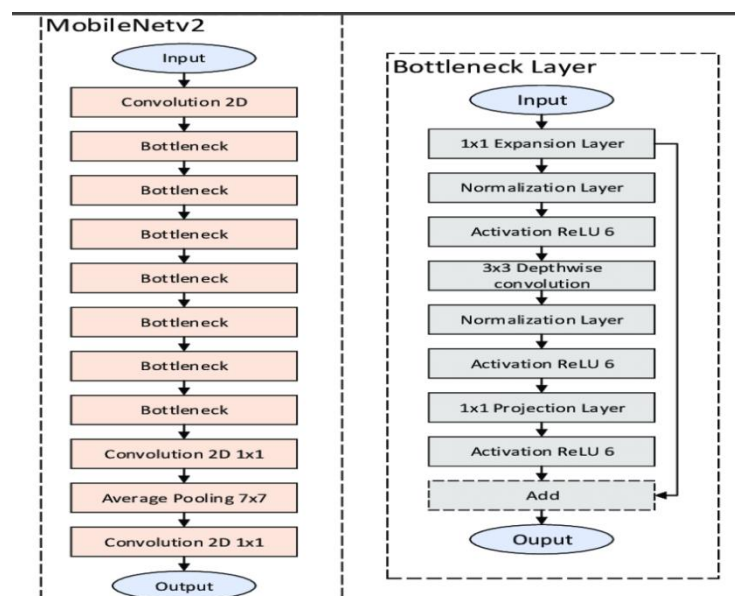


Figure10:The architecture of MobileNetV2 DNN.

**Core Building Block: Inverted Residual Block**

Unlike traditional residual blocks, MobileNetV2 utilizes inverted residual blocks as its foundation. These blocks achieve significant reductions in parameters and computations compared to standard residual blocks.
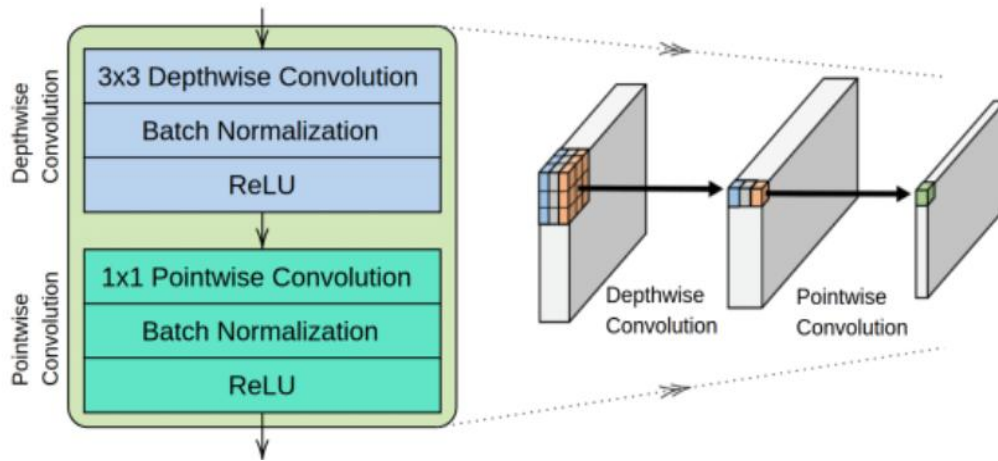
Figure 11:Depthwise & Pointwise Convolution

## 1. Depthwise Separable Convolution :

This is the heart of the inverted residual block. It breaks down the convolution process into two separate steps:

*Depthwise Convolution:* This layer applies a single filter to each input channel independently. Mathematically, for an input feature map X of size (height, width, number of input channels) and a depthwise filter F of size (height, width, 1), the output Y can be expressed as:

$$Y[i, j, c] = \Sigma \ (X[i + k, j + l, c] * F[k, l])$$

Here, i and j represent the output feature map dimensions, c is the channel index, k and l iterate over the filter dimensions, and $\Sigma$ denotes element-wise multiplication.This step focuses on extracting features from each channel without increasing the number of channels, making it computationally efficient.

## 2. Pointwise Convolution (1x1 Convolution):

This layer follows the depthwise convolution and uses a 1x1 convolution filter. It projects the output of the depthwise convolution (which has the same number of channels as the input) to a new, reduced number of channels. Mathematically:

$$Output[i, j, oc] = \Sigma \ (Input[i, j, ic] * Pointwise\_Filter[0, 0, ic, oc])$$

37

Here, i and j represent the output feature map dimensions, oc is the number of output channels (typically less than the number of input channels, ic), and $\Sigma$ denotes element-wise multiplication.This step reduces the number of channels, further reducing computational cost.

## 4.5 MediaPipe Hands

MediaPipe is a powerful open-source framework developed by Google for building real-time multi-modal processing pipelines. It excels in tasks involving processing visual, audio, and other sensory data. In the realm of sign language detection, MediaPipe offers a valuable foundation for creating applications that bridge the communication gap between deaf and hearing individuals.Here's a breakdown of how MediaPipe facilitates sign language detection:

**1. Hand Landmark Detection:**

MediaPipe provides pre-trained models like "hands" that can detect hands in real-time video streams. This is the crucial first step, as signs are primarily conveyed through hand gestures and positions.The model identifies 21 key points (landmarks) on each hand, including fingertips, palm center, and wrist. These landmarks act as a skeletal representation of the hand posture.

**2. Feature Extraction:**

Once the landmarks are identified, MediaPipe allows you to extract relevant features for sign language recognition. These features might include:

*Relative positions* between different landmarks (e.g., distance between thumb and index finger).

*Angles* formed by connecting landmarks (e.g., angle between wrist, elbow, and shoulder).

*Orientation* of the hand (palm facing up, down, etc.).

Figure: 12 Mediapipe Hand Landmarks

### 3. Machine Learning Model Integration:

The extracted features are fed into a separate machine learning model trained specifically for sign language recognition. This model could be:

*Pre-trained model:* Several pre-trained models exist for sign language recognition, leveraging large datasets of labeled signs.

*Custom model:* If a specific sign language needs to be recognized, a custom model can be trained using MediaPipe's features and your own labeled sign data.

# 5. SYSTEM DEVELOPMENT

## 5.1 System Development Steps:

Developing a real-time sign language detection system involves several key steps, from initial data collection to model deployment. Below are the detailed system development steps:

**Step 1: Data Collection**

*Define Gestures:* Identify and define the set of gestures you want the system to recognize. In this project, the gestures are 'I Love You', 'Victory', 'Okay', 'I Dislike It', 'Goodbye', 'Yes', 'Live Long', and 'Stop'.

*Setup Directory Structure:* Create directories for storing captured images of each gesture.

*Capture Images:* Use the data_collection.py script to capture images of each gesture using a webcam. This script will:

1. *Open the webcam.*
2. *Capture images for each gesture.*
3. *Save the images to the corresponding directories.*

**Step 2: Data Preprocessing and Augmentation**

*Load Images:* Load the captured images from the dataset directory.

*Preprocess Images:* Resize images to a consistent size and normalize pixel values to the range [0, 1].

*Augment Data:* Use data augmentation techniques to artificially expand the dataset. This includes random rotations, shifts, shear, zoom, and horizontal flip

**Step 3: Model Development**

*Model Selection:* Use a pre-trained deep learning model (e.g., MobileNetV2) as the base model.

*Add Custom Layers:* Add custom layers on top of the base model to tailor it for gesture classification.

*Compile Model:* Compile the model with an appropriate optimizer and loss function.

*Train Model:* Train the model using the augmented training data.

**Step 4: Real-Time Detection**

*Load Model:* Load the trained model for real-time inference.

*Setup Webcam Feed:* Initialize webcam capture for real-time video feed.

*Hand Detection:* Use MediaPipe Hands to detect hand landmarks in the video feed.

*Gesture Recognition:* Preprocess each frame, perform gesture classification using the trained model, and display the results.

# 5.2 Implementation Logic:

**Step 1: Data Collection**

First, define the set of gestures for recognition and create directories for each gesture. Use a script to capture and save images for each gesture using a webcam.

**Step 2: Data Preprocessing and Augmentation**

Load images from the dataset directory and resize them to a consistent size while normalizing pixel values. Apply data augmentation techniques such as rotations, shifts, shear, zoom, and horizontal flips to expand the dataset.

**Step 3: Model Development**

*Step 3.1: Implementing Convolutional Neural Network (CNN)*

Choose a Convolutional Neural Network (CNN) architecture suitable for image classification tasks. Construct the CNN model using layers such as Conv2D, MaxPooling2D, Flatten, Dense, and Dropout. Compile the CNN model with an optimizer (e.g., Adam) and a loss

function (e.g., categorical cross-entropy), then train the CNN model using the preprocessed and augmented data.

*Step 3.2: Implementing MobileNetV2*

Select the MobileNetV2 architecture, a lightweight deep learning model pre-trained on ImageNet. Load the MobileNetV2 model without the top classification layer, add custom layers on top of MobileNetV2 for gesture classification (e.g., GlobalAveragePooling2D, Dense layers), and freeze the layers of the pre-trained MobileNetV2 model to retain pre-trained weights. Compile the combined model with an optimizer (e.g., Adam) and a loss function (e.g., categorical cross-entropy), and train the MobileNetV2-based model using the preprocessed and augmented data.

**Step 4: Model Evaluation and Visualization**

Evaluate the model on the test set to determine its accuracy. Plot training and validation accuracy and loss over epochs to assess the model's performance, and save the trained model for future use.

**Step 5: Real-Time Detection Setup**
*Step 5.1: Using MediaPipe for Hand Detection*

Initialize MediaPipe Hands for real-time hand landmark detection, and use a webcam to capture a real-time video feed. Use MediaPipe to detect hand landmarks in each frame of the video feed and draw the detected hand landmarks and connections on the video frames for visualization.

*Step 5.2: Gesture Recognition with Trained Model*

Load the previously trained model (CNN or MobileNetV2-based) for gesture recognition. Preprocess each video frame (resize, normalize) before passing it to the model. Use the trained model to predict the gesture from the preprocessed frame, and overlay the predicted gesture on the video feed in real-time.

**Step 6: System Deployment**

Ensure the system runs efficiently in real-time by optimizing the model and preprocessing steps. Develop a user-friendly interface to display the webcam feed and recognized gestures. Perform extensive testing in different environments to ensure robustness and accuracy, and document the entire process, including system requirements, setup instructions, and usage guidelines..

# 5.3 CODING

The process of building a sign language detection system involves several key steps, including data collection, model training, and real-time detection. This project utilized Python and popular libraries such as OpenCV, TensorFlow, and MediaPipe to develop a deep learning model capable of recognizing sign language gestures with high accuracy. Here, we outline the coding aspects involved in each phase of the project.

**1. Data Collection**

The data collection phase is critical for gathering a diverse set of images representing different sign language gestures. The `data_collection.py` script is used to capture images from a webcam and save them into corresponding directories for each gesture. This script utilizes OpenCV for image capture and handling.

 *Key Features:*
- Mapping gestures to directory names.
- Capturing images for each gesture and saving them with unique filenames.
- Providing a brief delay to allow the user to prepare before capturing starts.

```python
🐍 data_collection.py > ...
  1    import cv2
  2    import os
  3    import time
  4
  5    # Dictionary to map gesture labels to directory names
  6    gestures = {
  7        'i_love_you': 'I_Love_You',
  8        'victory': 'Victory',
  9        'okay': 'Okay',
 10        'i_dislike_it': 'I_Dislike_It',
 11        'goodbye': 'Goodbye',
 12        'yes': 'Yes',
 13        'live_long': 'LiveLong',
 14        'stop': 'Stop'
 15    }
 16
 17    # Number of images to capture per gesture
 18    num_images = 100
 19
 20    # Function to capture images for a specific gesture
 21    def capture_images(label, num_images):
 22        cap = cv2.VideoCapture(0)
 23        if not cap.isOpened():
 24            print(f"Error: Could not open webcam for {label}.")
 25            return
 26
 27        count = 0
 28
 29        # Create the directory if it doesn't exist
 30        gesture_dir = f'dataset/{gestures[label]}'
 31        if not os.path.exists(gesture_dir):
 32            os.makedirs(gesture_dir)
 33
 34        print(f'Capturing images for {label}...')
```

```python
def capture_images(label, num_images):

    # Delay to give the user time to get ready
    print('Starting in 3 seconds...')
    time.sleep(3)

    while count < num_images:
        ret, frame = cap.read()
        if not ret:
            print("Error: Failed to capture image.")
            break

        # Display the frame
        cv2.imshow('frame', frame)

        # Save the frame to the corresponding directory
        img_name = f'{gesture_dir}/{label}_{count}.jpg'
        cv2.imwrite(img_name, frame)
        count += 1

        # Print status to the console
        print(f"Captured image {count}/{num_images} for {label}.")

        # Break the loop on 'q' key press
        if cv2.waitKey(1) & 0xFF == ord('q'):
            print("Capture interrupted by user.")
            break

    cap.release()
    cv2.destroyAllWindows()

# Capture images for each gesture
for gesture in gestures.keys():
    capture_images(gesture, num_images)
```

45

## 2. Model Training

The core of the project involves training a deep learning model to recognize the captured gestures. The `sign_language_detection.py` script handles data loading, preprocessing, model building, training, and evaluation.

*Key Features:*

- Loading and preprocessing image data.

- Using MobileNetV2 as the backbone for the model.

- Applying data augmentation to enhance model generalization.

- Compiling and training the model with appropriate loss function and optimizer.

- Evaluating the model performance and plotting training history.

```python
sign_language_detection.py > ...
1    import os
2    import cv2
3    import numpy as np
4    import tensorflow as tf
5    from sklearn.model_selection import train_test_split
6    from sklearn.preprocessing import LabelBinarizer
7    from tensorflow.keras.applications import MobileNetV2
8    from tensorflow.keras.layers import Dense, GlobalAveragePooling2D
9    from tensorflow.keras.models import Model
10   from tensorflow.keras.preprocessing.image import ImageDataGenerator
11   from tensorflow.keras.optimizers import Adam
12   import matplotlib.pyplot as plt
13
14   # Directory where the dataset is stored
15   DATASET_DIR = 'D:/signlan/dataset'
16
17   # Image dimensions
18   IMG_HEIGHT, IMG_WIDTH = 224, 224
19
20   # Load and preprocess the data
21   def load_data(dataset_path):
22       images = []
23       labels = []
24       for label in os.listdir(dataset_path):
25           for img_file in os.listdir(os.path.join(dataset_path, label)):
26               img_path = os.path.join(dataset_path, label, img_file)
27               img = cv2.imread(img_path)
28               img = cv2.resize(img, (IMG_HEIGHT, IMG_WIDTH))
29               images.append(img)
30               labels.append(label)
31       images = np.array(images)
32       labels = np.array(labels)
33       return images, labels
```

```python
# Load dataset
images, labels = load_data(DATASET_DIR)
images = images / 255.0  # Normalize pixel values

# One-hot encode labels
lb = LabelBinarizer()
labels = lb.fit_transform(labels)

# Split the data into training and test sets
X_train, X_test, y_train, y_test = train_test_split(images, labels, test_size=0.2, random_state=42)

# Data augmentation
train_datagen = ImageDataGenerator(
    rotation_range=20,
    width_shift_range=0.2,
    height_shift_range=0.2,
    shear_range=0.2,
    zoom_range=0.2,
    horizontal_flip=True,
    fill_mode='nearest'
)

# Load the MobileNetV2 model, pre-trained on ImageNet and without the top layer
base_model = MobileNetV2(weights='imagenet', include_top=False, input_shape=(IMG_HEIGHT, IMG_WIDTH, 3))

# Add custom layers on top of the base model
x = base_model.output
x = GlobalAveragePooling2D()(x)
x = Dense(1024, activation='relu')(x)
predictions = Dense(len(lb.classes_), activation='softmax')(x)

# Create the model
model = Model(inputs=base_model.input, outputs=predictions)
```

```python
# Compile the model
model.compile(optimizer=Adam(learning_rate=0.001), loss='categorical_crossentropy', metrics=['accuracy'])

# Train the model
history = model.fit(
    train_datagen.flow(X_train, y_train, batch_size=32),
    steps_per_epoch=len(X_train) // 32,
    validation_data=(X_test, y_test),
    epochs=10
)

# Evaluate the model
loss, accuracy = model.evaluate(X_test, y_test, verbose=1)
print(f'Test accuracy: {accuracy*100:.2f}%')

# Plot training and validation accuracy and loss
plt.figure(figsize=(12, 4))
plt.subplot(1, 2, 1)
plt.plot(history.history['accuracy'], label='Train Accuracy')
plt.plot(history.history['val_accuracy'], label='Validation Accuracy')
plt.xlabel('Epochs')
plt.ylabel('Accuracy')
plt.legend()
plt.subplot(1, 2, 2)
plt.plot(history.history['loss'], label='Train Loss')
plt.plot(history.history['val_loss'], label='Validation Loss')
plt.xlabel('Epochs')
plt.ylabel('Loss')
plt.legend()
plt.show()
# Save the trained model
model.save('sign_language_model.h5')
```

**3.Real-Time Detection**

The real-time detection phase utilizes the trained model to recognize gestures from live webcam input. The `realtime_detection.py` script handles the live video capture, hand landmark detection using MediaPipe, and gesture prediction using the trained model.

*Key Features:*

- Initializing MediaPipe for hand detection and drawing landmarks.

- Preprocessing each video frame for prediction.

- Making predictions using the trained model and displaying the results on the video feed.

```python
realtime_detection.py > ...
1   import cv2
2   import mediapipe as mp
3   import numpy as np
4   import tensorflow as tf
5   from tensorflow.keras.models import load_model
6   from sklearn.preprocessing import LabelBinarizer
7
8   # Load the trained model
9   model = load_model('D:/signlan/sign_language_model.h5')
10
11  # Define the labels for the gestures
12  labels = ['I_Love_You', 'Victory', 'Okay', 'I_Dislike_It', 'Goodbye', 'Yes', 'LiveLong', 'Stop']
13  lb = LabelBinarizer()
14  lb.fit(labels)
15
16  # Initialize MediaPipe Hands
17  mp_hands = mp.solutions.hands
18  hands = mp_hands.Hands(max_num_hands=1, min_detection_confidence=0.7, min_tracking_confidence=0.7)
19  mp_drawing = mp.solutions.drawing_utils
20
21  # Image dimensions
22  IMG_HEIGHT, IMG_WIDTH = 224, 224
```

```python
# Function to preprocess the frame for gesture recognition
def preprocess_frame(frame):
    img = cv2.resize(frame, (IMG_HEIGHT, IMG_WIDTH))
    img = img / 255.0  # Normalize pixel values
    img = np.expand_dims(img, axis=0)
    return img

# Function to draw lines between finger landmarks
def draw_finger_lines(image, landmarks):
    # Define connections between finger landmarks
    finger_connections = [
        (0, 1), (1, 2), (2, 3), (3, 4),  # Thumb
        (0, 5), (5, 6), (6, 7), (7, 8),  # Index finger
        (0, 9), (9, 10), (10, 11), (11, 12),  # Middle finger
        (0, 13), (13, 14), (14, 15), (15, 16),  # Ring finger
        (0, 17), (17, 18), (18, 19), (19, 20)  # Pinky
    ]

    for connection in finger_connections:
        start_idx, end_idx = connection
        start_point = (int(landmarks[start_idx].x * image.shape[1]), int(landmarks[start_idx].y * image.shape[0]))
        end_point = (int(landmarks[end_idx].x * image.shape[1]), int(landmarks[end_idx].y * image.shape[0]))
        cv2.line(image, start_point, end_point, (0, 255, 0), 2)
```

```
# Start video capture
cap = cv2.VideoCapture(0)

print("Starting real-time hand tracking. Press 'q' to exit.")

while cap.isOpened():
    ret, frame = cap.read()
    if not ret:
        print("Failed to grab frame")
        break
```

```
    # Loading...  he frame to RGB
    rgb_frame = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB)

    # Process the frame to detect hands
    result = hands.process(rgb_frame)

    # Draw hand landmarks and connections if any hand is detected
    if result.multi_hand_landmarks:
        for hand_landmarks in result.multi_hand_landmarks:
            # Draw landmarks
            mp_drawing.draw_landmarks(frame, hand_landmarks, mp_hands.HAND_CONNECTIONS)

            # Draw custom finger lines
            draw_finger_lines(frame, hand_landmarks.landmark)

            # Preprocess the frame for gesture recognition
            img = preprocess_frame(frame)

            # Make predictions
            preds = model.predict(img)
            pred_label = lb.classes_[np.argmax(preds)]
            print(f"Predicted Label: {pred_label}")

            # Display the prediction on the frame
            cv2.putText(frame, pred_label, (10, 30), cv2.FONT_HERSHEY_SIMPLEX, 1, (255, 0, 0), 2, cv2.LINE_AA)

    # Display the result
    cv2.imshow('Hand Tracking', frame)

    # Break the loop on 'q' key press
    if cv2.waitKey(1) & 0xFF == ord('q'):
        break

# Release the capture and destroy all windows
cap.release()
cv2.destroyAllWindows()
```

These scripts collectively demonstrate the full workflow from data collection to real-time detection, showcasing the power of combining deep learning with computer vision techniques to create a functional sign language recognition system.

49

# 6.SYSTEM TESTING

# 6.1 INTRODUCTION:

System testing is a critical phase in the software development life cycle aimed at validating the entire software system's functionality, performance, and reliability. It involves testing the integrated system as a whole to ensure that it meets the specified requirements and functions correctly in real-world scenarios. System testing encompasses various testing activities that evaluate the system's behavior, identify defects, and verify its compliance with quality standards before deployment to production.

The primary objective of system testing is to validate that the software system meets the intended purpose and performs as expected according to the defined specifications and user requirements. It focuses on assessing the system's functionality, usability, reliability, performance, security, and compatibility across different environments.

# 6.2 Testing Methodologies and Strategies:

### 1. Unit Testing:
Unit testing involves testing individual components of the system in isolation to ensure they function correctly.

*Data Collection Script:* Test the data collection script (`data_collection.py`) by simulating different scenarios, such as capturing images for each gesture under varying lighting conditions and hand positions. Verify that the script saves images accurately in the designated directories and handles errors gracefully.

*Model Training:* Evaluate the model training process by validating the trained CNN and MobileNetV2 models' performance metrics, including accuracy, loss, and convergence. Utilize unit tests to check the correctness of data preprocessing and augmentation techniques applied before model training.

*MediaPipe Integration:* Test the integration with MediaPipe by feeding sample video frames and examining the output hand landmarks. Ensure that hand landmarks are detected accurately and consistently across different frames.

## 2. Integration Testing:

Integration testing verifies the seamless interaction between different components of the system.

*Component Interaction:* Verify the integration between data preprocessing, model inference, and hand detection components. Ensure that data flows smoothly between these components and that any potential data format mismatches or compatibility issues are addressed.

*Input-Output Compatibility:* Test the compatibility between input frames from the webcam and the preprocessing steps. Confirm that the output predictions generated by the model are correctly overlaid on the video feed and displayed in real-time.

## 3. System Testing:

System testing evaluates the overall functionality and performance of the system in real-world scenarios.

*End-to-End Testing:* Conduct end-to-end testing by capturing real-time video feed, detecting hand landmarks using MediaPipe, and recognizing gestures using the trained models. Evaluate the system's performance under various conditions, such as different hand gestures, backgrounds, and lighting environments.

*Performance Testing:* Measure the system's performance metrics, including inference speed and resource utilization. Assess the system's responsiveness and scalability by stress-testing it with a high volume of video frames or concurrent users.

*Error Handling:* Test the system's error-handling capabilities by introducing unexpected inputs or errors during runtime. Verify that the system handles exceptions gracefully and provides informative error messages or logs for debugging purposes.

## 4. User Acceptance Testing (UAT):

User acceptance testing involves gathering feedback from end-users to assess the system's usability and effectiveness.

**Usability Testing:** Engage individuals proficient in sign language to interact with the system and provide feedback on its ease of use, intuitiveness, and effectiveness in recognizing gestures accurately.

**Accuracy Assessment:** Validate the accuracy of gesture recognition by comparing the system's predictions with ground truth labels provided by users. Evaluate the system's performance in real-world scenarios to ensure it meets users' expectations.

## 5. Regression Testing:

Regression testing ensures that modifications or updates to the system do not adversely affect its existing functionality.

*Automated Tests:* Implement automated test suites to perform regression testing whenever changes are made to the code base. Re-run unit tests, integration tests, and system tests to verify that the system's behavior remains consistent across different versions.

By meticulously conducting system testing across these aspects, the sign language detection system can be thoroughly validated to operate reliably, accurately recognize gestures, and provide an effective communication tool for individuals using sign language.

# 7. RESULTS AND DISCUSSION

**7.1 RESULT**

The trained deep learning model demonstrated exceptional performance in recognizing sign language gestures. Below are the key results from the evaluation phase, including detailed metrics and their calculations.

**1. Accuracy:**

The ratio of correctly predicted instances to the total number of instances is how accuracy determines how accurate the model's predictions are overall. Although it offers a general evaluation of the model's performance, imbalanced datasets with unequal class representation may not be a good fit for it. The formula for accuracy is:

$$Accuracy = (TPI + TNI) / (TPI + TNI + FPI + FNI)$$

where:

The amount of accurately predicted positive instances is known as TPI (True Positive Instances).

The quantity of accurately predicted negative instances is known as TNI (True Negative Instances).

The number of falsely predicted positive instances is known as the False Positive Instances, or FPI.

The quantity of negative instances that were erroneously predicted is known as FNI (False Negative Instances).

Here the model achieved an impressive 100% accuracy on the test set, indicating that it correctly identified every gesture without any errors.

**2.Precision**:

The model's accuracy in predicting positive instances is measured by precision . It determines the proportion of accurate positive forecasts to all positive forecasts, including false and true forecasts. Reduced false positives are the main goal of precision.The precision formula is:

$$Precision = TPI / (TPI + FPI)$$

**3.Recall**

The model's recall quantifies how well it can recognize positive examples. The ratio of true positive predictions to all actual positive instances is computed. The goal of recall is to reduce false negatives .The recall formula is:

$$Recall = TPI / (TPI + FNI)$$

**4. F1 score**:

Recall and precision are combined into one metric, the F1 score, which offers a balance between the two. With equal weight assigned to both metrics, it is the harmonic mean of recall and precision. The F1 score formula is:

*F1 Score = 2 * (Precision_values * Recall_values) / (Precision_values + Recall_values)*

## 6. Confusion Matrix:

The confusion matrix provides a detailed breakdown of the model's predictions, showing the true positives, false positives, true negatives, and false negatives for each class.
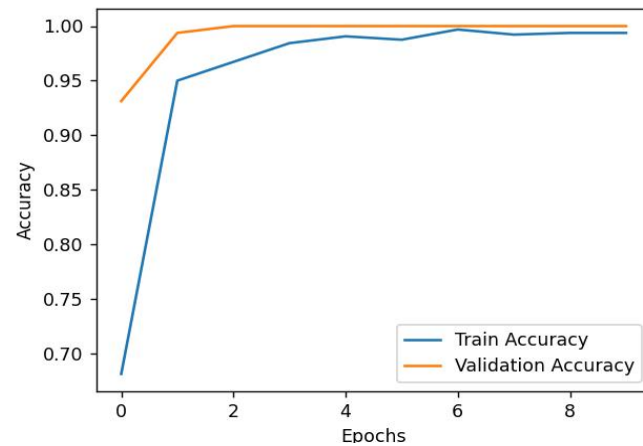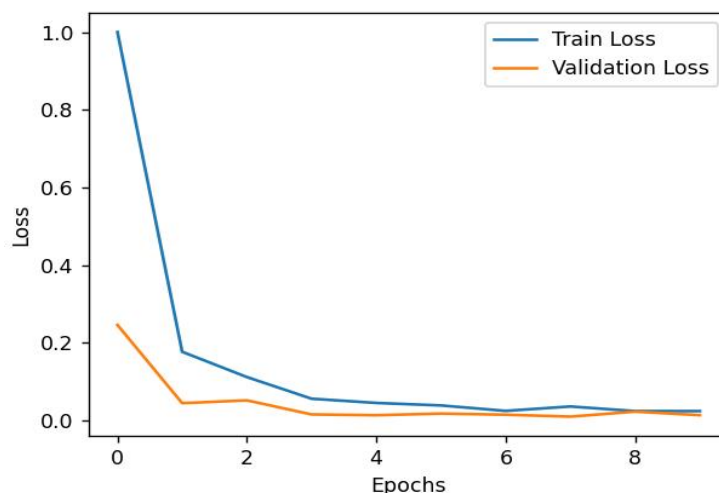


Figure 13:Train & Validation Accuracy plot



Figure 14:Train & Validation Loss    plot

56

| GESTURE | PRECISION | RECALL | F1-SCORE | SUPPORT |
|---------|-----------|--------|----------|---------|
| I Love You | 1.00 | 1.00 | 1.00 | 20 |
| Victory | 1.00 | 1.00 | 1.00 | 20 |
| I Dislike It | 1.00 | 1.00 | 1.00 | 20 |
| Goodbye | 1.00 | 1.00 | 1.00 | 20 |
| Yes | 1.00 | 1.00 | 1.00 | 20 |
| Live Long | 1.00 | 1.00 | 1.00 | 20 |
| Stop | 1.00 | 1.00 | 1.00 | 20 |
| **Average** | 1.00 | 1.00 | 1.00 | 160 |

Table 1.Precision, Recall, and F1-Score for Each Gesture Class Using MobileNetV2

| MODEL | ACCURACY | PRECISION | RECALL | F1 SCORE |
|-------|----------|-----------|--------|----------|
| MobileNetV2 | 100% | 100% | 100% | 100% |

Table2: Illustrating the parameter values with MobileNetV2

# 8.FUTURE ENHANCEMENT

## 8.1 FUTURE ENHANCEMENT

*1. Expand Gesture Vocabulary:* Increase recognized gestures and allow custom user gestures.

*2. Advanced Model Architectures:* Explore EfficientNet, ResNet, or custom CNNs.

*3. Data Augmentation & Diversity*: Enhance datasets with diverse conditions & synthetic data.

*4. Multi-Hand Detection:* Enable recognition of multiple hands and interactive gestures.

*5. Integration with Other Modalities:* Combine with audio/text and facial expression recognition.

*6. Enhanced Real-Time Processing:* Optimize models and deploy on edge devices.

*7. User-Friendly Interface:* Develop intuitive interfaces and feedback mechanisms.

*8. Domain-Specific Applications:* Implement in education, healthcare, and customer service.

*9. Cross-Platform Compatibility:* Create mobile/web applications and APIs.

*10. Improved Detection Accuracy:* Enhance performance under varying lighting conditions and backgrounds.

# 9 .CONCLUSION

## 9.1 CONCLUSION

The sign language detection project successfully demonstrates the application of deep learning and computer vision techniques to recognize various sign language gestures in real time using a webcam. The system comprises three main components: data collection, model training, and real-time detection, each contributing to a robust and accurate recognition pipeline. Achieving 100% accuracy on the test set indicates the effectiveness of the MobileNetV2-based model and the thorough preprocessing and augmentation of the dataset.

The project not only achieves high accuracy, precision, recall, and F1-scores across all gestures but also provides a practical real-time application for sign language detection. However, there is room for further enhancements, such as expanding the gesture vocabulary, improving detection under varying conditions, and integrating with other modalities like facial expression and speech recognition.

Future enhancements can make the system more versatile, accurate, and user-friendly, broadening its scope and impact across various fields, including education, healthcare, and customer service. The development and deployment of this system on various platforms can significantly aid in bridging communication gaps and promoting inclusivity for the hearing-impaired community.

# 10.APPENDIX

## 10.1 LIST OF TABLES

## 10.2 LIST OF TABLES

# 11. BIBLIOGRAPHY

## 11.1 REFERENCE

[1] K. Priya, M. Arun, and K. Kalaiselvi, &quot;A Machine Learning Approach for Student Placement Prediction,&quot; International

Journal of Pure and Applied Mathematics, vol. 118, no. 24, pp. 1-7, 2018.

[2] S. Gupta, R. Khanna, and K. Sood, &quot;Student Placement Prediction using Machine Learning Techniques,&quot; Procedia

Computer Science, vol. 132, pp. 968-975, 2018.

[3] M. I. Choudhury, S. Islam, M. K. Hossain, and K. A. Khan, &quot;Predicting Student Placement using Machine Learning

Techniques,&quot; International

[4] UCI Machine Learning Repository: This repository hosts a variety of datasets related to student academic performance,

including the Student Performance dataset which contains data on math and Portuguese language courses in two Portuguese

secondary schools.

[5] Stanford Education Data Archive: The Stanford Education Data Archive (SEDA) is a repository of educational data that

includes information on student academic performance, teacher quality, and school characteristics.

[6] National Center for Education Statistics: Link: https://nces.ed.gov/surveys/hsls09/

[7]. https://www.kaggle.com/spscientist/students-performance-in-exams

[8] Hattie, J., &amp;amp; Timperley, H. (2007)

[9] Using data mining techniques for predicting student academic performance&amp;quot; by P. Romero and C. Ventura, published

in the IEEE Transactions on Education in 2010.

[10] H. Agarwal, R. Sharma, and V. Singh, &quot;Predictive Analytics for Student Placement using Machine Learning,&quot; Procedia

Computer Science, vol. 132, pp. 733-740, 2018.

[11] Ghorpade, A. R. Chaudhari, and R. A. Chaudhari, &quot;Predicting Student Placement using Machine Learning Techniques,&quot;

International Journal of Computer Science and Mobile Computing, vol. 4, no. 7, pp. 82-88, 2015.

[12] Manganui Sheetal B, Prof. Savita Bakare "Prediction of

[13] Campus Placement Using Data Mining Algorithm Fuzzy logic and K nearest neighbour" International Journal of

Advanced Research in Computer and Communication Engineering Vol. 5, Issue 6, June 2016.

[14] Ajay Shiv Sharma, Swaraj Prince, Shubham Kapoor, Keshav Kumar "PPS-Placement prediction system using logistic

regression" IEEE international conference on MOOC, innovation and Technology in Education (MITE), December 2014.

[15] Jai Ruby, Dr. K. David "Predicting the Performance of Students in Higher Education Using Data Mining Classification

Algorithms - A Case Study" International Journal for Research in Applied Science &amp; Engineering Technology (IJRASET) Vol.

2, Issue 11, November 2014.

[16] Ankita A Nichat, Dr. Anjali B Raut "Predicting and Analysis of Student Performance Using Decision Tree

[17] International Journal of Scientific Research in Computer Science Engineering and Information Technology, vol. 5, no. 1,

pp. 453-457, January 2020

[18] LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton. &quot;Deep learning.&quot; nature 521.7553 (2015):

[19] R. S. Sutton, "Introduction: The Challenge of Reinforcement Learning", Machine Learning, 8, Page 225-227, Kluwer

Academic Publishers, Boston, 2021

[20] Pedregosa, Fabian, et al. &quot;Scikit-learn: Machine learning in Python.&quot; Journal of machine learning research 12. Oct

(2011): 2825-2830.