

**VISVESVARAYA TECHNOLOGICAL UNIVERSITY**  
**Jnana Sangama, Belagavi - 590 018**



**PROJECT REPORT ON**  
**Real Time Sign Language Recognition**  
**System Using Machine Learning**

*This is submitted in partial fulfillment for the Award of Degree of*

**Bachelor of Engineering**  
in  
**Electronics and Communication Engineering**

**Submitted by**

Raghavendra M Hegde	1RN20EC066
Samarth Shinnur	1RN20EC078
Samrudh B R	1RN20EC079
Sanjay S B	1RN20EC081

*Under the Guidance of*  
**Dr. Leena Chandrashekhar**  
*Assistant Professor*



**DEPARTMENT OF ELECTRONICS AND COMMUNICATION ENGINEERING**

**RNS INSTITUTE OF TECHNOLOGY**  
(AICTE Approved, VTU Affiliated and NAAC 'A+' Grade' Accredited)  
(UG Programs - CSE, ECE, ISE, EIE and EEE have been Accrediated by  
NBA for the Academic Year 2022-25 )  
Channasandra, Dr.Vishnuvardhan Road,Bengaluru-560098  
2023 - 24

**VISVESVARAYA TECHNOLOGICAL UNIVERSITY**  
**Jnana Sangama, Belagavi - 590 018**



**PROJECT REPORT ON**  
**Real Time Sign Language Recognition**  
**System Using Machine Learning**

*This is submitted in partial fulfillment for the Award of Degree of*

**Bachelor of Engineering**  
in  
**Electronics and Communication Engineering**

**Submitted by**

Raghavendra M Hegde	1RN20EC066
Samarth Shinnur	1RN20EC078
Samrudh B R	1RN20EC079
Sanjay S B	1RN20EC081

*Under the Guidance of*  
**Dr. Leena Chandrashekhar**  
*Assistant Professor*



**DEPARTMENT OF ELECTRONICS AND COMMUNICATION ENGINEERING**

**RNS INSTITUTE OF TECHNOLOGY**  
(AICTE Approved, VTU Affiliated and NAAC 'A+' Grade' Accredited)  
(UG Programs - CSE, ECE, ISE, EIE and EEE have been Accrediated by  
NBA for the Academic Year 2022-25 )  
Channasandra, Dr.Vishnuvardhan Road,Bengaluru-560098  
2023 - 24

**RNS INSTITUTE OF TECHNOLOGY**  
(AICTE Approved, VTU Affiliated and NAAC ‘A+ Grade’ Accredited)  
(UG Programs - CSE, ECE, ISE, EIE and EEE have been Accrediated by  
NBA for the Academic Year 2022-25 )  
Channasandra, Dr.Vishnuvardhan Road,Bengaluru-560098

**DEPARTMENT OF ELECTRONICS AND COMMUNICATION ENGINEERING**



## **CERTIFICATE**

Certified that the Project work entitled “Real Time Sign Language Recognition System Using Machine Learning” is carried out by **Raghavendra M Hegde (1RN20EC066)**, **Samarth Shinnur (1RN20EC078)**, **Samrudh B R (1RN20EC079)** and **Sanjay S B (1RN20EC081)** in partial fulfillment for the award of degree of Bachelor of Engineering in **Electronics and Communication Engineering** of Visvesvaraya Technological University, Belagavi, during the year 2023-2024. It is certified that all corrections and suggestions indicated during internal assessment have been incorporated in the report. The project report has been approved as it satisfies the academic requirements in respect of the project work prescribed for the award of degree of **Bachelor of Engineering**.

**Dr. Leena Chandrashekhar**  
Assistant Professor

**Dr. Vipula Singh**  
Head of the Department

**Dr. Ramesh Babu H S**  
Principal

### **External Viva**

**Name of the examiners**

**Signature with date**

1 .....

2 .....

**RNS INSTITUTE OF TECHNOLOGY**  
(AICTE Approved, VTU Affiliated and NAAC ‘A+ Grade’ Accredited)  
(UG Programs - CSE, ECE, ISE, EIE and EEE have been Accredited by  
NBA for the Academic Year 2022-25 )  
Channasandra, Dr.Vishnuvardhan Road,Bengaluru-560098

DEPARTMENT OF ELECTRONICS AND COMMUNICATION ENGINEERING



## DECLARATION

We here by declare that the entire work in this project report titled, “**Real Time Sign Language Recognition System Using Machine Learning**” submitted to **Visvesvaraya Technological University Belagavi**, is carried out at the department of **Electronics and Communication Engineering, RNS Institute of Technology, Bengaluru** under the guidance of **Dr.Leena Chandrashekhar**, Assistant Professor. This report has not been submitted for the award of any Diploma or Degree of this or any other University.

Name	USN	Signature
1. Raghavendra M Hegde	1RN20EC066	.....
2. Samarth Shinnur	1RN20EC078	.....
3. Samrudh B R	1RN20EC079	.....
4. Sanjay S B	1RN20EC081	.....



## RNS INSTITUTE OF TECHNOLOGY

Autonomous Institution Affiliated to VTU, Recognized by GOK, Approved by AICTE,  
(NAAC 'A+' Grade) Accredited, NBA Accredited UG - CSE, ECE, ISE, EIE and EEE  
Dr. Vishnuvardhan Road, Channasandra, Rajarajeshwari Nagar, Bengaluru-560 098, Karnataka.  
Telephone: 080-28611880/1  
Website: [www.rnsit.ac.in](http://www.rnsit.ac.in)

Department of Electronics & Communication Engineering  
**IETE Sponsored Seventh National Conference**

on  
**Emerging Trends in Engineering, Science and Technology (NCETEST-7)**



## Certificate of Participation

This is to certify that Mr/Ms Raghavendra M Hegede \_\_\_\_\_  
of  
RNS Institute of Technology \_\_\_\_\_  
has participated and presented  
a paper entitled Real Time Sign Language Recognition System Using  
Machine Learning  
in the Seventh National Conference on  
"Emerging Trends in Engineering, Science & Technology" held on April 30 2024, Organized by the  
Department of Electronics & Communication Engineering, in collaboration with IETE.

  
\_\_\_\_\_  
**Dr. Nandini K S**  
Coordinator  
  
\_\_\_\_\_  
**Dr. Vipula Singh**  
HoD, ECE  
  
\_\_\_\_\_  
**Dr. Ramesh Babu HS**  
Principal  
  
\_\_\_\_\_  
**Dr. M K Venkatesha**  
Director  
  
\_\_\_\_\_  
**PRINCIPAL**  
**R.N.S.I.T., Bengaluru-98** **R.K. Chinnappa**



## RNS INSTITUTE OF TECHNOLOGY

Autonomous Institution Affiliated to VTU, Recognized by GOK, Approved by AICTE,  
 (NAAC 'A+' Grade) Accredited, NBA Accredited UG - CSE, ECE, ISE, EIE and EEE  
 Dr. Vishnuvardhan Road, Channasandra, Rajarajeshwari Nagar, Bengaluru-560 098, Karnataka.  
 Telephone: 080-28611880/1  
 Website: [www.rnsit.ac.in](http://www.rnsit.ac.in)



### Department of Electronics & Communication Engineering IETE Sponsored Seventh National Conference

on  
**Emerging Trends in Engineering, Science and Technology (NCETEST-7)**

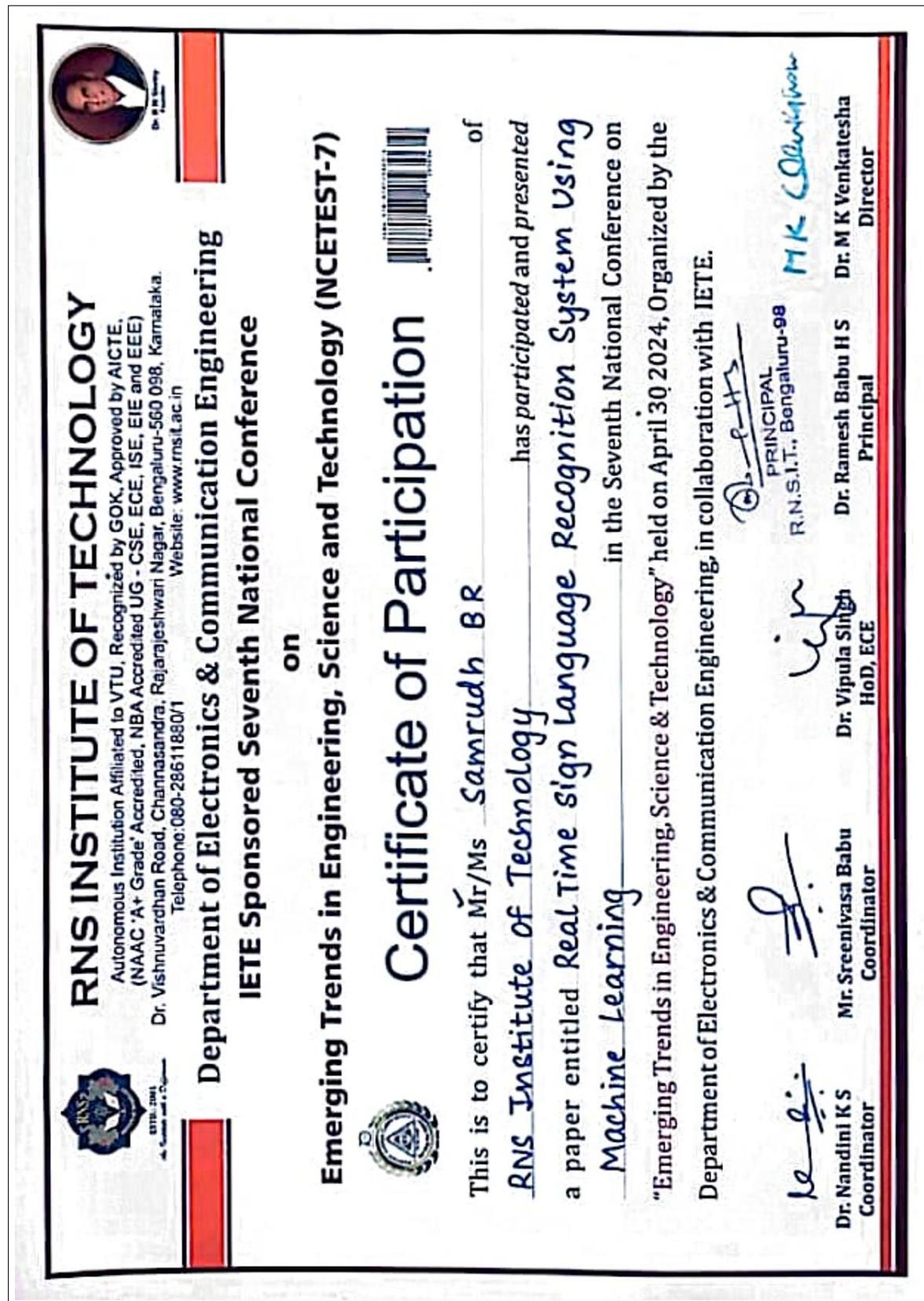


## Certificate of Participation

This is to certify that Mr/Ms Samarth S \_\_\_\_\_  
RNS Institute of Technology \_\_\_\_\_ has participated and presented  
 a paper entitled Real Time Sign Language Recognition System  
Using Machine Learning \_\_\_\_\_ in the Seventh National Conference on  
 "Emerging Trends in Engineering, Science & Technology" held on April 30 2024, Organized by the  
 Department of Electronics & Communication Engineering, in collaboration with IETE.

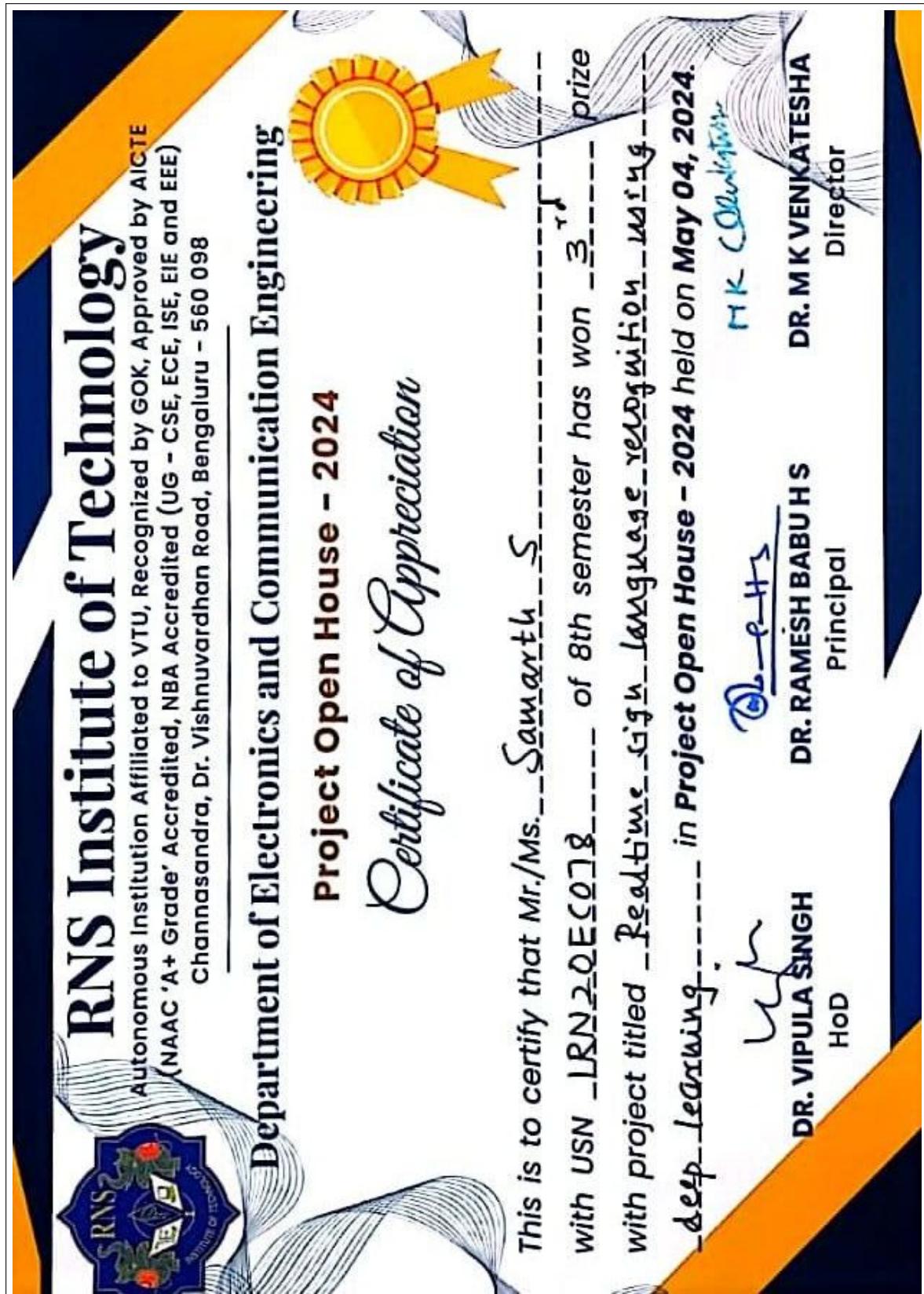
R.N.S.I.T., Bengaluru-98 M K Cheluvappa  
PRINCIPAL  
Dr. Ramesh Babu H S Dr. M K Venkatesha  
HoD, ECE Principal  
Mr. Sreenivasa Babu Director  
Coordinator

Dr. Nandini KS Dr. Vipula Singh  
Coordinator HoD, ECE











# RNS Institute of Technology

Autonomous Institution Affiliated to VTU, Recognized by GOK, Approved by AICTE  
 (NAAC 'A+' Grade' Accredited, NBA Accredited (UG - CSE, ECE, ISE, EIE and EEE)  
 Channasandra, Dr. Vishnuvardhan Road, Bengaluru - 560 098



Department of Electronics and Communication Engineering

## Project Open House - 2024

### Certificate of Appreciation

This is to certify that Mr./Ms. Sanjay S. B with USN JRN2OEC08L of 8th semester has won 3<sup>rd</sup> prize with project titled Real time sign language recognition using deep learning in Project Open House - 2024 held on May 04, 2024.

**DR. VIPULASINGH**  
 HOD

**DR. RAMESH BABU HS**  
 Principal

**DR. M K VENKATESHA**  
 Director

## Acknowledgement

The joy and satisfaction that accompany the successful completion of any task would be incomplete without thanking those who made it possible. We consider ourselves proud to be a part of RNS Institute of Technology, the institution which moulded us in all our endeavors.

We express our gratitude to our Chairman **Late Dr. R N Shetty** and Managing Director, **Mr. Satish Shetty**, for providing state of art facilities.

We would like to express our sincere thanks to **Dr. M K Venkatesha**, Director, **Dr. Ramesh Babu H S**, Principal, and **Dr. Vipula Singh**, Professor and HOD, Department of Electronics and Communcitaion Engineering, for their valuable guidance and encouragement throughout our program.

We wish to express our profound thanks to the coordinators for their invaluable insights and guidance during the project. Furthermore, our heartfelt gratitude goes to **Dr. Leena Chandrashekhar**, Assistant Professor in the Department of Electronics and Communication Engineering, for her consistent support, guidance, and encouragement, all of which were pivotal in the successful culmination of our project.

Finally, we take this opportunity to extend our earnest gratitude and respect to our parents, teaching and non-teaching staff of the department, the library staff and all our friends who have directly or indirectly supported us.

**Raghavendra M Hegde**

**Samarth Shinnur**

**Samrudh B R**

**Sanjay S B**

## **Abstract**

Sign language is a main mode of communication for vocally disabled. This language use set of representation which is finger sign, expression or mixture of both to express their information among others. Yet not everyone is able to interpret sign language which might result in miscommunication. This work presents a robust system that leverages machine learning and computer vision techniques to accurately detect and classify sign language gestures in real-time. This technology is made to accurately detect and identify sign language motions in real-time.

The system introduce a reliable hand sign tracking and recognition model that is accurate despite any complex environment, it can track and recognize RGB dynamic gestures. Firstly, videos are converted into frames. After images light intensity adjustment and noise removal, images are passed through CNN for Key-point extraction. From the acquired hand images, features are extracted covering entire palm. The optimized features are then passed to the Random Forest machine learning Algorithm for the classification of the Hand signs. The custom datasets are used for the validation of proposed systems. The experimental results over custom datasets demonstrate overall accuracies of 94.58% and 95.67% respectively. The experiments proves our system readability and suitability of our proposed model with the other state of the art techniques.

# Table of Contents

<b>Abstract</b>	ii
<b>Table of Contents</b>	iii
<b>List of Figures</b>	v
<b>List of Tables</b>	1
<b>1 Introduction</b>	2
1.1 Motivation . . . . .	3
1.2 Objectives . . . . .	3
1.3 Methodology . . . . .	3
1.4 Applications . . . . .	5
1.5 Benefits and Challenges . . . . .	5
<b>2 Literature Survey</b>	8
<b>3 Software Dependencies</b>	25
3.1 Programming with Python . . . . .	25
3.2 Sci-kit Learn . . . . .	26
3.3 Open CV . . . . .	27
3.4 Pandas . . . . .	28
3.5 Media Pipe . . . . .	29
<b>4 Proposed Classifier Implementation</b>	35
4.1 Data Set Collection and Preprocessing . . . . .	37
4.2 Mediaipe Feature Extractor . . . . .	38
4.3 Design of Random Forest Classifier . . . . .	40
<b>5 Results and Disscussion</b>	45
5.1 Feature Extraction Results . . . . .	46
5.2 Performance Metrics . . . . .	48
5.3 Comparative Analysis Across Lighting Conditions . . . . .	51
5.4 Speechify Text . . . . .	55
5.5 Comparison with State-of-Arts-Techniques . . . . .	57
5.6 Conclusion . . . . .	58
5.7 Future Scope . . . . .	59

<b>References</b>	<b>61</b>
<b>Appendices</b>	<b>65</b>
<b>A Code</b>	<b>65</b>
<b>B Conference Paper</b>	<b>75</b>

# List of Figures

1.1	Block Diagram	4
3.1	SSD Architecture	30
3.2	VGG16 Architecture	32
4.1	Flow Chart of Hand Sign Recognition System	35
4.2	Media Pipe Hand Landmark	39
5.1	Frequency of Occurance	48
5.2	Confusion Matrix for 200, 100 and 50 n-estimators	50
5.3	Daytime Lighting Evaluation	52
5.4	Dim Lighting Test Case	53
5.5	Dark Environment Evaluation	54
5.6	Speechify Text	56

# List of Tables

4.1	Parameters Grid Values . . . . .	44
4.2	GridSearchCV Result . . . . .	44
5.1	Feature Values for X-Coordinate . . . . .	46
5.2	Feature Values for Y-Coordinate . . . . .	47
5.3	Overall Results of Proposed technique . . . . .	50
5.4	Image Quality Parameters for Case 1 (Light) . . . . .	52
5.5	Image Quality Parameters for Case 2 (Dim) . . . . .	53
5.6	Image Quality Parameters for Case 3 (Dark) . . . . .	54
5.7	A Chronological Summary of Various Techniques in the Domain . . . . .	58

# Chapter 1

## Introduction

In our everyday lives, the interaction between humans and computers has become essential, thanks to the continuous advancements in technology. This interaction, known as Human-Computer Interaction (HCI), encompasses various methods through which we communicate with and control computers. From clicking a mouse to tapping on a touchscreen, HCI plays a significant role in how we navigate the digital world.

Among the diverse aspects of HCI, hand gesture recognition has emerged as a prominent area of interest for researchers. It offers a natural and intuitive means of interacting with computers, devoid of the complexities associated with traditional input devices. Unlike requiring training for devices like joysticks or remote controls, using hand gestures feels instinctive and user-friendly. Its applications span across numerous domains, from controlling home appliances to guiding robots and aiding in medical procedures[1]. However, despite its potential benefits, hand gesture recognition presents challenges, particularly in ensuring accurate recognition across varying environmental conditions. In critical fields such as healthcare, where precision is paramount, errors in gesture recognition can have serious consequences. Researchers are actively striving to develop robust systems that can reliably interpret gestures under diverse circumstances, but achieving this remains a formidable task[2-4].

The process of hand gesture recognition typically involves several stages, including capturing gestures, identifying hand positions, extracting relevant features, and classifying gestures. Utilizing technologies such as cameras and sensors, these stages enable the detection and analysis of both simple and complex gestures. In the realm of our research project, we are committed to advancing the field of hand gesture recognition by proposing innovative solutions to address existing challenges[5-7]. Our focus lies on enhancing the accuracy and reliability of gesture recognition systems, particularly in critical domains like healthcare and home automation. Additionally, our efforts aim to streamline communication between humans and computers, facilitating more seamless interactions. Through our project, we aspire not only to improve hand gesture recognition technology but also to foster a future where human-computer interaction is more intuitive and efficient. By bridging the gap between humans and machines, we aim to pave the way for a more connected and accessible digital world[8-9].

## 1.1 Motivation

The development of Sign Language Recognition Systems, fueled by compelling statistics on the prevalence of hearing impairments, aims to bridge the communication gap between deaf individuals and non-sign language users. With 2.42 million deaf and mute individuals in India and over 466 million globally, representing over 5% of the world's population, the need for such systems is urgent. Leveraging advanced image processing and machine learning, these systems interpret sign language gestures, enabling participation across education, healthcare, and workplaces, fostering inclusivity and understanding[1-8].

Furthermore, these systems reduce reliance on costly human interpreters, offering a scalable and cost-effective solution for interpretation services. By ensuring wider accessibility, they contribute to a more inclusive society, minimizing communication barriers and enhancing participation for deaf individuals. Ultimately, the goal is to enhance communication, accessibility, and inclusion, empowering deaf individuals to lead fulfilling lives and engage fully in societal activities[9-14].

## 1.2 Objectives

- To create a comprehensive dataset of sign language gestures by converting video sequences into frames.
- To extract 21 key points from the hand images to represent the sign language gestures effectively.
- To train the machine learning model using the extracted features and labeled data with the goal of optimizing the model's accuracy.
- To leverage the trained model to develop a text-to-speech application that can interpret written text and generate corresponding spoken output.

## 1.3 Methodology

Developing a Sign Language Recognition System involves several key steps as shown in Figure 1.1 .Initially, a comprehensive dataset of sign language gestures is collected, spanning various dialects. These gestures are then processed through image enhancement techniques to refine quality and isolate the hand region. Following this, features such as hand shape, movement, and position are extracted using methods like Mediapipe. Machine learning models like Convolutional Neural Networks (CNNs) or

Random Forest classifiers are then chosen and trained on the extracted features and labeled data. Subsequently, the model's performance is rigorously tested and evaluated using metrics like accuracy, precision, recall, and F1-score. Once validated, the model can be deployed for real-world applications, aiding in bridging communication barriers for the hearing impaired.

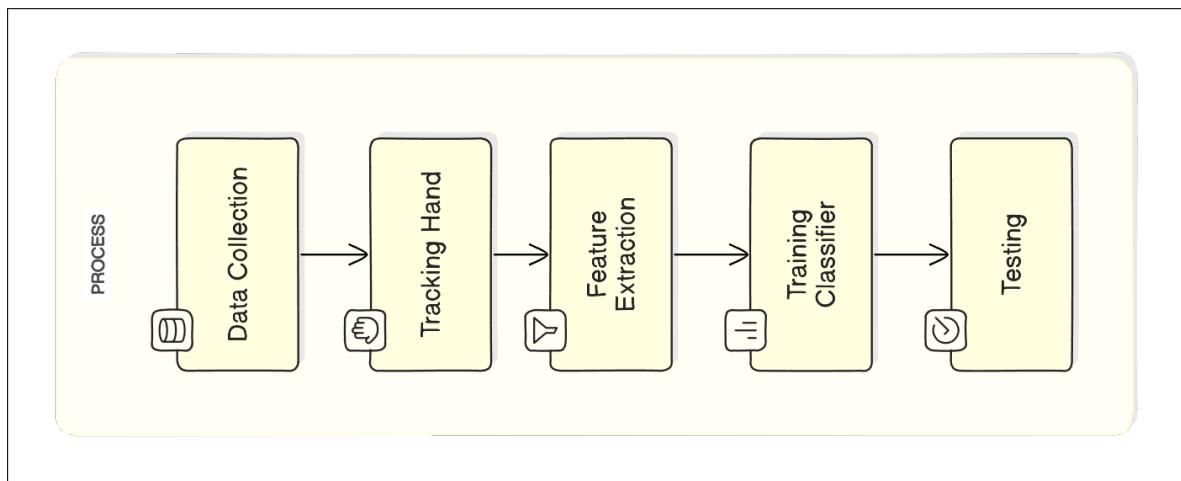


Figure 1.1: Block Diagram

- Data collection: The first step in developing a Sign Language Recognition System is to collect a large dataset of sign language gestures by converting video sequence into frames using open CV. This dataset should include a variety of sign language gestures from different sign language dialects.
- Image Processing: The next step is to preprocess the images of the sign language gestures to improve their quality, remove background noise, and isolate the hand region.
- Feature Extraction: After preprocessing, features need to be extracted from the images to represent the sign language gestures. Extract 21 holistic key points in 3D coordinates using Mediapipe. These features could include hand shape, hand movement, and hand position.
- Model Selection: The next step is to select an appropriate machine learning model that can classify sign language gestures accurately. Some of the popular machine learning models used for sign language recognition include Convolutional Neural Networks (CNNs), and Random Forest classifier.
- Model Training: The selected machine learning model is trained on the extracted features and labeled data. This process involves feeding the model with the data and adjusting its parameters to optimize its accuracy.

- Testing and Evaluation: After training the model, it is tested on a separate dataset to evaluate its performance. The performance metrics include accuracy, precision, recall, and F1-score.
- Deployment: Once the model has been tested and evaluated, it can be deployed in the real-world setting.

## 1.4 Applications

Sign Language Recognition Systems hold profound potential across various domains, notably in enabling real-time sign language interpretation. By facilitating seamless communication between sign language users and non-sign language users, these systems contribute significantly to the inclusion of deaf or hard of hearing individuals in society, ultimately enhancing their quality of life. In educational settings, such systems play a pivotal role by offering real-time interpretation in classrooms, thereby ensuring equitable access to education for deaf or hard of hearing students. Additionally, they empower teachers and instructors to better comprehend and interact with their deaf or hard of hearing students, fostering a more inclusive learning environment.

Moreover, Sign Language Recognition Systems have substantial implications in healthcare, where they enable real-time interpretation during medical appointments, enhancing accessibility for deaf or hard of hearing patients. This not only facilitates effective communication between medical professionals and patients but also promotes better health outcomes. Beyond healthcare and education, these systems contribute to enhancing accessibility across technology products, such as smartphones and computers, by integrating sign language-enabled interfaces. Furthermore, they serve as invaluable tools for training and evaluating sign language interpreters, ensuring the provision of accurate interpretation services, and offering constructive feedback to sign language learners to improve their proficiency. Through these multifaceted applications, Sign Language Recognition Systems play a pivotal role in fostering inclusivity and accessibility for the deaf and hard of hearing community.

## 1.5 Benefits and Challenges

### 1.5.1 Benefits

Sign language serves as a vital means for facilitating swift and efficient communication among deaf and hard of hearing individuals, fostering inclusivity and understanding.

By bridging communication gaps, sign language enables seamless interaction, empowering individuals to express thoughts, emotions, and desires confidently. Additionally, mastering sign language provides a powerful avenue for self-expression, promoting independence and boosting self-esteem, particularly for children facing hearing difficulties.

Moreover, Sign Language Recognition Systems offer a cost-effective solution to increase accessibility for deaf or hard of hearing individuals, facilitating effective communication with non-sign language users in real-time. By enabling prompt exchange of information, these systems prove invaluable, especially in time-sensitive situations where immediate communication is essential. Furthermore, the scalability of Sign Language Recognition Systems ensures adaptability to accommodate varying numbers of users, making them suitable for organizations or institutions catering to large populations of deaf or hard of hearing individuals, thereby ensuring widespread accessibility.

### 1.5.2 Challenges

Sign language, inherently reliant on hand gestures, poses challenges for individuals with limited hand mobility, potentially hindering their ability to effectively communicate through this medium. Additionally, while Random Forests offer enhanced accuracy in sign language recognition by employing a larger number of trees, this comes at the expense of increased computational resources and longer training times. Moreover, the scalability of Random Forests can be a concern when dealing with large datasets, leading to computational overhead and slower model performance, particularly in real-time applications where low latency is crucial.

Furthermore, interpreting feature importance in Random Forest models can be challenging, especially with correlated features, potentially leading to misconceptions about feature relevance. Additionally, the limited repertoire of recognized sign language gestures by Sign Language Recognition Systems may not fulfill all communication needs. Moreover, the system's accuracy can be impacted by environmental factors such as lighting conditions, background noise, and camera quality, posing further challenges in ensuring reliable recognition in various settings.

The remainder of the report is organized as follows: Chapter 2 delves into the Literature Survey summarizing the key findings of referenced papers, their methodologies and their contributions to shaping our project. Chapter 3 shifts focus to System Analysis, offering technical details such as software dependencies and a high-level design

overview. Chapter 4 details the Implementation process describing the step-by-step methodology employed in detail. Lastly Chapter 5 focuses on Simulation and Evaluation presenting the outputs and results obtained assessing the efficiency and accuracy of our model and suggesting avenues for future work.

# Chapter 2

## Literature Survey

This chapter explores various techniques that have been proposed to address the challenges in accurate hand gesture recognition and sign language translation, which are critical for bridging the communication barrier faced by individuals relying on sign language. The methods covered in this chapter include leveraging machine learning models, improving object detection frameworks, and employing signal processing approaches. These techniques aim to overcome issues such as complex backgrounds, lighting variations, occlusions, and variations in hand postures and signing styles that hinder effective communication through sign language. By examining these state-of-the-art approaches, this chapter provides a comprehensive overview of the current efforts to develop robust and accurate solutions for sign language recognition and translation.

The author in [1] describes the development and implementation of a gesture-based language translator to aid communication for individuals who are deaf and mute. The authors begin by discussing the challenges faced by people who are deaf and mute in communicating with others. They describe how current technologies are limited in their ability to address the complex and nuanced nature of sign language, and highlight the need for a more effective and efficient communication tool. A reliable hand gesture tracking and recognition model that can accurately track and recognize RGB dynamic gestures, with overall accuracies of 90.73% and 89.33% over Egogesture and Jester datasets, respectively. The proposed model is specifically designed for hand gesture recognition in the medical field and outperforms other state-of-the-art methods.

The methodology in [1] includes converting videos into frames, adjusting image light intensity, removing noise, extracting features from the full hand, using neural gas and locomotion thermal mapping to create the feature vector, passing the feature vector through a fuzzy optimiser, and using the Deep Belief Network (DBW) for gesture classification. The hand gesture recognition system in [1] achieved high overall accuracies by Conversion of videos into frames, Image light intensity adjustment and noise removal, Hand gesture extraction using CNN, Feature extraction from the full hand, Neural gas extraction for feature vector determination, Locomotion thermal mapping for capturing hand movements, Fuzzy optimization for feature vector optimization, Classification of gestures using Deep Belief Network (DBN), Evaluation of

performance metrics for system robustness. Overall in conclusion the authors evaluate the performance of the system using a custom dataset of hand gestures and a range of languages. They report high accuracy rates for both gesture recognition and language translation. The authors also discuss the potential applications of the system, including its use in hospitals and public spaces. They highlight the importance of making the system available and accessible and suggest future directions for research and development[1].

The author in [2] illustrates the challenges in hand gesture recognition and the proposed solution of using a multi-branch attention-based graph and general deep-learning model to recognize hand gestures, achieving high accuracy with the proposed model. The methodology involves proposing a multi-branch attention-based graph and general deep-learning model, using two graph-based neural network channels and one general neural network channel, constructing the final feature vector by concatenating spatial-temporal, temporal-spatial, and general features, including position embedding and mask operation for both spatial and temporal attention modules, developing the architecture based on the dynamic graph-based attention-based mechanism, considering two branches as the graph-based deep neural network branch, employing position embedding and masking operations for each attention at spatial and temporal domains, and recognizing hand gestures based on 3D hand skeleton data points.

Data augmentation techniques such as shifting, scaling, time interpolation, and adding noise, Implementation of the architecture in the PyTorch platform ,Random selection of eight frames for each video as input, Subtraction of every input frame sequence by the first frame palm position, Use of Adam optimizer with a learning rate of 0.001 for training the model, Setting the batch size to 32 and dropout rate to 0.1 and 0.2, Employment of an attention-based Multi-Branch Attention Based Graph and General Deep Learning approach for recognizing hand gestures. Overall this paper provides a summary of the challenges in hand gesture recognition, the proposed solution using a Multi-Branch Attention Based Graph and General Deep Learning Model, and the achievement of high accuracy with the proposed model when evaluated with three benchmark datasets. The proposed system has potential applications in various fields, including assistive technology for individuals with speech disabilities and as a tool for human-computer interaction. The system can be used to provide a more natural and intuitive way of communicating with computers and other electronic devices. The proposed system can also be extended to recognize a larger set of hand gestures and can be integrated with other speech recognition systems to provide a more comprehensive speech and gesture recognition system[2].

Hand detection plays a crucial role in various computer vision applications, ranging from gesture recognition to human-computer interaction. The study in [3] focuses on the significance of hand detection, particularly emphasizing the advancements made by Yolov7 and Yolov7x models in this domain. The primary objective of this research is to assess the performance of these models and compare their effectiveness in hand detection tasks. The methodology involves a comprehensive review of existing research papers related to hand detection techniques. Additionally, it encompasses an explanation of the training procedure adopted for Yolov7 and Yolov7x models, along with a detailed analysis of the outcomes. The evaluation process includes utilizing sample image sets to train and test the models, thereby facilitating a comparative assessment of their performance.

During training, Yolov7x with 200 epochs exhibited the most stable and highest performance, achieving impressive precision (84.7%), recall (79.9%), and mean Average Precision (mAP) (86.1%). Similarly, during testing, Yolov7x demonstrated superior performance compared to Yolov7, with precision, recall, and mAP values of 84.4%, 80%, and 86.3%, respectively. The discussion section provides a detailed analysis of the performance metrics obtained by Yolov7 and Yolov7x models across different epochs. Notably, the superior performance of Yolov7x with 200 epochs is highlighted, emphasizing specific accuracy results for hand detection tasks. Furthermore, insights are offered regarding the factors contributing to the observed performance differences between the two models[3].

Gesture recognition, especially hand gesture recognition, is a burgeoning field with applications ranging from human-computer interaction to prosthetic control systems. The study in [4] focuses on identifying specific hand gestures from Electromyography (EMG) signals acquired using sensor-based bands. EMG signals, representing muscle contractions, are preprocessed to remove noise artifacts before extracting eight time-domain features for classification. Machine learning techniques, specifically Support Vector Machines (SVM) and K-Nearest Neighbors (KNN) are employed for classification. The study aims to achieve high accuracy in recognizing hand gestures, offering efficient device control capabilities. The methodology begins with preprocessing raw EMG signals to eliminate noise using rectification, filtering, and detrending techniques. Rectification ensures signal polarity consistency, while filtering with a Butterworth bandpass filter removes unwanted frequency components. Time-domain features extracted from preprocessed signals form feature matrices for each hand gesture. SVM and KNN classifiers trained on Mio dataset and tested on Ninapro facilitate the workflow comprising preprocessing, feature extraction and classification phases implemented in MATLAB 2019a.

Despite achieving promising accuracy levels, the study has several limitations. Firstly, the reliance on surface EMG signals may pose challenges in accurately capturing subtle gestures or distinguishing between similar movements. Additionally the effectiveness of the classification algorithms may vary depending on factors such as dataset size, gesture complexity, and individual variability in muscle contractions. Furthermore the study focuses on hand gesture recognition potentially overlooking other aspects of gesture-based interaction such as finger movements or body gestures. The study reports an accuracy of 83% for SVM and 93% for KNN in recognizing hand gestures from EMG signals. Confusion matrices illustrate the classification performance, highlighting both correct predictions and misclassifications. These results demonstrate the efficacy of the proposed methodology in accurately identifying hand gestures for efficient device control and human-computer interaction[4].

The study in [5] presents a communication system designed to facilitate interaction between sign language users and non-users by recognizing hand gestures and translating them into text or speech in real-time. Leveraging machine learning algorithms including K-Nearest Neighbors (KNN), Support Vector Machine (SVM) and Convolutional Neural Network (CNN) the system aims to bridge the communication gap. The dataset consists of approximately 29,000 images representing American Sign Language (ASL) alphabets. KNN is employed for classification, with parameter optimization for noise reduction. SVM is utilized for classification and regression analysis, with parameter tuning for optimal model fitting. CNN extracts features from hand gesture images to achieve high accuracy. The methodology involves collecting a dataset of hand gesture images captured via webcam, preprocessing the images, and splitting them into training and testing sets. Three machine learning algorithms, namely KNN, SVM, and CNN, are employed for hand gesture recognition. KNN utilizes a feature space built on nearest training data points, with parameter optimization for noise reduction. SVM is utilized for classification and regression analysis, with parameter tuning for optimal model fitting. CNN is employed to extract features from hand gesture images. Model performance is evaluated based on accuracy and confusion matrices. Real-time translation of recognized hand gestures into text or speech is implemented using appropriate libraries.

Several limitations are encountered during the development and implementation of the system. These include the challenge of creating a large and diverse dataset of hand gesture images, which requires significant time and resources. Additionally, overfitting and underfitting of machine learning models pose challenges, requiring careful selection of model architectures and hyperparameters. The computational complexity of CNN models may limit real-time performance on resource-constrained devices.

Furthermore, the system's accuracy may be affected by variations in hand gestures, lighting conditions, and background clutter, necessitating robustness improvements. The system demonstrates promising results in recognizing hand gestures and translating them into text or speech. CNN achieves the highest accuracy of 98.49%, outperforming KNN (93.83%) and SVM (88.89%). Confusion matrices provide insights into the classification performance of each algorithm. Real-time translation capabilities enhance the system's usability for sign language communication. Despite challenges such as dataset creation and model overfitting, the system holds potential for improving communication between sign language users and non-users, fostering inclusivity and accessibility in social interactions. Further research and development are warranted to address limitations and enhance the system's robustness and scalability[5].

The paper [6] outlines the importance of gesture recognition in human-computer interaction, highlighting the limitations of traditional approaches and the pressing need for enhanced accuracy and robustness through advanced machine learning or deep learning algorithms. Its primary aim is to develop an algorithm for gesture detection and recognition that effectively balances precision and speed. The study introduces an innovative approach based on an improved YOLOv5 framework specifically tailored for complex background scenarios. It incorporates an efficient layer aggregation network module to refine YOLOv5 for better gesture feature extraction and to optimize network complexity thereby enhancing recognition speed. Furthermore the introduction of the CBAM attention mechanism module bolsters network robustness by discerning and prioritizing relevant features while filtering out irrelevant ones.

The method achieves recognition accuracies, with mAP@0.5: 0.95 scores of 75.6% and 66.8% for the EgoHands and TinyHGR datasets, respectively, while maintaining a commendable recognition speed of 64 FPS for  $640 \times 640$  input images. Notably, it outperforms benchmarks such as YOLOv5l, YOLOv7, and other comparative algorithms in both accuracy and robustness, demonstrating superior performance across diverse complex backgrounds[6].

Communication plays a crucial role in human life, but individuals with speech impairments rely on sign language, presenting challenges in communication with the majority. Automatic sign language recognition systems can bridge this gap by converting sign language into spoken or written language and vice versa. However, existing systems are limited, costly, and cumbersome. Researchers worldwide are working on developing such recognizers, prompting the need for robust and efficient systems. The paper [7] aims to train TensorFlow for sign language recognition, design a suitable algorithm for hand gesture recognition, and implement real-time sign language

detection. The proposed system facilitates communication between speech-impaired individuals who use hand signs and gestures and normal individuals, thus promoting inclusivity and accessibility in social interactions.

The methodology in [7] involves training TensorFlow for sign language recognition, designing a suitable algorithm for hand gesture recognition, and implementing real-time sign language detection. Existing literature on sign language recognition systems is reviewed to understand the current state-of-the-art approaches. TensorFlow is utilized for training neural networks for sign language recognition. An appropriate algorithm, such as Convolutional Neural Network (CNN), is designed for hand gesture recognition. Real-time sign language detection is implemented using image processing techniques and deep learning architectures.

Several limitations are encountered during the development and implementation of the proposed system. These include the need for a diverse and comprehensive dataset for training TensorFlow, challenges in accurately recognizing complex hand gestures, and computational complexity in real-time detection. Additionally, the system's performance may be affected by variations in lighting conditions, background clutter, and individual hand gestures, requiring robustness improvements. Existing literature presents various approaches to sign language recognition, including techniques based on Histogram of Oriented Gradients (HOG), Support Vector Machine (SVM), Convolutional Neural Network (CNN), and image processing. These approaches demonstrate promising results in recognizing sign language gestures and converting them into spoken or written language. However, challenges such as computational complexity, dataset diversity, and real-time performance need to be addressed. The proposed system aims to overcome these challenges and enhance communication between speech-impaired individuals and normal individuals, thus promoting inclusivity and accessibility in social interactions. Further research and development are warranted to improve the system's robustness and scalability[7].

The author in [8] introduces a real-time system designed to recognize poses and hand gestures in Indian Sign Language (ISL) using grid-based features. The primary motivation behind this system is to bridge the communication gap between individuals who are deaf or hard of hearing and the broader society. Existing solutions often suffer from low precision or lack real-time functionality, necessitating the development of a more accurate and efficient system. ISL, a complex sign language comprising 28 hand poses including digits and letters, is the focus of this research. The system employs techniques such as face detection, object stabilization, and skin color segmentation to process sign language captured from a smartphone camera. By utilizing grid-based

feature extraction and classification algorithms like k-Nearest Neighbors (k-NN) and Hidden Markov Models (HMMs), the system aims to achieve both high precision and real-time performance.

The methodology in [8] involves several key steps to achieve accurate and real-time recognition of ISL gestures. Initially, the captured sign language frames undergo pre-processing, including face removal, stabilization, and skin color segmentation to isolate the hand region. Hand detection and tracking are then performed using segmentation techniques. For hand pose recognition, grid-based feature extraction is employed, and the resulting feature vectors are classified using the k-NN algorithm. Gesture recognition involves observing intermediate hand positions and movements, encoding them into sequences, and utilizing HMMs for classification. Object stabilization, achieved through facial reference tracking, ensures accurate hand motion detection despite camera movement. The system architecture utilizes an Android smartphone for capturing sign language and a remote server for processing frames and performing recognition tasks. A dataset comprising thousands of images and gesture videos was used for training and testing the system.

Despite its effectiveness, the system has some limitations. One major limitation is the reliance on a steady camera position for accurate object stabilization, which may not always be feasible in practical scenarios. Additionally, the system's performance may vary depending on lighting conditions and background clutter, potentially affecting the accuracy of hand detection and tracking. Furthermore, the system's current implementation focuses on single-handed gestures, limiting its applicability to two-handed gestures without further extension. The system demonstrates promising results in hand pose and gesture recognition in ISL. With an accuracy of 96.4% in classifying hand poses and an average accuracy of 98.23% in recognizing gestures, the system outperforms existing approaches in terms of precision and real-time performance. The recognition time for hand poses is approximately 0.25 seconds, while for gestures, it is 0.0048 seconds, highlighting the system's efficiency. These results suggest that the proposed system can effectively bridge the communication gap between individuals who understand ISL and those who do not, providing a practical solution for real-time sign language translation[8].

As the demand for human-computer interaction grows, hand gesture recognition has become increasingly essential, finding applications in various fields like robot control and intelligent furniture. Hand gestures convey rich information, making them valuable for interaction with machines. The paper [9] presents a system for hand gesture recognition, focusing on segmentation, tracking, and recognition of hand gestures.

The system comprises three main components: hand gesture segmentation, tracking, and recognition. Hand gesture segmentation involves isolating hand gestures from video frames, a crucial step preceding recognition. Various techniques, including skin color-based segmentation, edge detection, motion analysis, and statistical templates, are employed. The paper [9] utilizes a fusion algorithm to achieve robust segmentation even in complex backgrounds. Hand gesture tracking ensures real-time monitoring and location of hand gestures within the video stream. Popular tracking algorithms such as mean-shift, Kalman filtering, and optical flow are considered. The paper opts for the cam shift algorithm, an improved mean-shift version, due to its balance of real-time performance and accuracy. For hand gesture recognition, the system focuses on recognizing digits 1-10. A dataset comprising 16,000 hand gesture images is collected, with 4,000 images allocated for testing. The LeNet-5 neural network architecture is employed for classification within specific regions of interest. The network is trained to classify hand gestures into the ten digit categories.

The Gaussian mixture model of skin color is utilized for hand gesture segmentation, yielding satisfactory results as shown in the experiments. The system achieves an impressive average accuracy of 98.3% in hand gesture recognition, with the exception of digits 7 and 9 due to their complexity. The limitation of the system lies in its inability to capture 3D information of hand gestures, which could enhance recognition accuracy further. The proposed system demonstrates robustness and accuracy in hand gesture recognition, offering potential applications in various interactive systems. By combining segmentation, tracking, and recognition techniques, the system effectively interprets hand gestures, facilitating seamless human-computer interaction. Further advancements in capturing 3D information could enhance the system's capabilities in recognizing complex hand gestures[9].

The author in [10] illustrates an imaging system designed for hand sign recognition and monitoring. This innovation marks a significant advancement in the field of gesture recognition. The study reports an impressive classification accuracy of 95.94% achieved using Convolutional Neural Networks (CNN). This underscores the effectiveness of CNN algorithms in hand sign recognition tasks. Methodologically the paper outlines the steps involved in developing and evaluating the imaging system. This includes data collection from eight subjects validation using a water tank system implementation and optimization of the CNN classification algorithm. The study compares the performance of the classification algorithm with other classifiers. This comparative analysis provides insights into the strengths and weaknesses of different approaches. Additionally the authors provide a step-by-step guide for implementing Support Vector Machines (SVMs) in practice.

The authors emphasize the significance of careful model selection and evaluation, particularly for SVMs. This discussion underscores the importance of avoiding overfitting to ensure the reliability of classification results. The paper elucidates the fundamentals of Electrical Impedance Tomography (EIT) and the process of database collection for EIT measurements. This foundational knowledge contextualizes the development of the EIT imaging system. The study details the optimization and evaluation of the CNN classification algorithm. This includes strategies to mitigate overfitting and comparisons of classification accuracies across different algorithms[10].

The study in [11] aims to recognize American Sign Language (ASL) characters using hand images captured from a webcam. It leverages the MediaPipe Hands algorithm to estimate hand joint coordinates from RGB hand images. Based on these joint coordinates, two types of features are generated for classification: distances between joint points and angles between vectors and 3D axes. The employment of two machine learning classifiers, Support Vector Machine (SVM) and Light Gradient Boosting Machine (GBM), to classify the ASL characters using three datasets: the ASL Alphabet dataset, the Massey dataset, and the Finger Spelling dataset. The use of two different classifiers SVM and GBM allows the researchers to compare the performance of these machine learning models on the ASL recognition task. SVMs find the optimal hyperplane to separate data points into classes, while GBMs are ensemble methods that combine many weak models into a powerful predictive model. The proposed system achieves high accuracy rates, with 99.39% for the Massey dataset, 87.60% for the ASL Alphabet dataset, and 98.45% for the Finger Spelling A dataset.

The problem addressed in this study [11] is the recognition of American Sign Language (ASL) characters using hand images captured from a webcam. While sensor-based approaches historically yield higher accuracy, they can be costly. Vision-based methods, although more accessible, may trade accuracy for affordability. Thus, this study aims to develop a cost-effective and accurate system utilizing vision-based methods. Leveraging the MediaPipe Hands algorithm, hand joint coordinates are estimated from RGB hand images. Features are then extracted from these coordinates, considering distances between joints and angles relative to 3D axes. The objective is to effectively represent the intricate shapes of hand gestures and enhance recognition accuracy. By employing Support Vector Machine (SVM) and Light Gradient Boosting Machine (GBM) classifiers, the study evaluates the performance across three datasets: the ASL Alphabet dataset, the Massey dataset, and the Finger Spelling A dataset. The ultimate goal is to provide a system that requires no special sensors or devices, is computationally efficient, and outperforms previous studies in accuracy[11].

Sign language serves as a crucial means of communication between deaf-mute individuals and the hearing population through hand gestures. Visual-based gesture recognition systems play a vital role in bridging this communication gap. While various techniques have been proposed in this field, many existing methods still face compatibility issues and limitations. Hand segmentation, as the initial step in processing input images for recognition, encounters challenges in distinguishing the hand from the face region, especially with skin detection methods. Motivated by these challenges the paper [12] reviews and discusses the progress of feature extraction methods in sign language recognition over the past decade. By focusing on feature extraction, the aim is to identify the most effective and compatible method for sign language recognition systems, facilitating future research progress.

The paper [12] reviews literature from previously published international papers discussing sign language recognition. The primary focus is on studying feature extraction methods to improve hand gesture recognition accuracy. Hand segmentation methods discussed include the canny edge detector for locating hand regions based on detected edges, as well as skin detection algorithms that leverage color models or machine learning to identify skin-colored regions corresponding to the hand and arm. Additionally the use of gloves for hand tracking is explored, ranging from coloured glove-based methods to sensory glove-based approaches. Principal Component Analysis (PCA), Scale Invariant Feature Transform (SIFT) algorithm, Serial Particle Filter algorithms for recursively estimating hand poses over time using a particle filter framework, and the Microsoft Kinect depth sensor for providing 3D hand position and joint data. Each method's principles, advantages, and reported accuracies are summarized to provide insights into their effectiveness and potential applications.

The feature extraction methods in sign language recognition highlights significant progress and diverse approaches in the field. While challenges such as hand segmentation and compatibility persist, researchers have explored various techniques to improve accuracy and efficiency. The use of gloves, particularly sensory gloves, has simplified hand tracking algorithms, while PCA offers dimensionality reduction for effective data representation. Additionally, robust algorithms like the SIFT algorithm and Serial Particle Filter show promise in tracking hand movements and reducing noise in recognition systems. Microsoft Kinect, with its 3D sensor camera, provides efficient hand gesture recognition capabilities, overcoming challenges such as lighting conditions. Overall, the review underscores the importance of feature extraction methods in advancing sign language recognition systems and identifies avenues for future research, including real-time systems and further exploration of markerless passive sensors[12].

Hand gesture recognition plays a vital role in Natural User Interfaces (NUI), enabling intuitive interaction with technology. The first step in hand gesture recognition is extracting hand features accurately and efficiently. The paper [13] proposes a real-time method for hand feature recognition using three cameras. The hand region is extracted using background subtraction, and features such as arm angle and finger positions are calculated based on variations in the vertical contour image. Wrist detection is achieved by calculating the distance between a baseline and the hand contour. The proposed method is evaluated on a dataset comprising approximately 1800 images captured from three videos, demonstrating its effectiveness and efficiency. The method utilizes three cameras to capture hand gestures in real time. Background subtraction is employed to extract the hand region, followed by the calculation of arm angle and finger positions using variations in the vertical contour image. Wrist detection is achieved by determining the maximum distance between a baseline and the hand contour. Experiments conducted on a dataset consisting of 1800 images from three videos validate the performance of the proposed method.

The real-time method in [13] for hand feature recognition demonstrates promising results in hand gesture recognition tasks. By leveraging background subtraction and contour analysis techniques, the method efficiently extracts hand regions and computes key features such as arm angle, finger positions, and wrist detection. Despite challenges such as motion blur and occlusion, the method performs well and offers a low-resource solution for implementing finger detection with good results. While improvements can be made to address issues such as misplaced wrist points and occlusion, the proposed method offers a practical and efficient approach for scenarios where speed is prioritized over precision in finger detection. Overall, the method presents a valuable contribution to the field of hand gesture recognition, offering a cost-effective and efficient solution for real-time applications[13].

With the continuous advancements in smartphone technology, opportunities arise to revolutionize healthcare practices by facilitating remote monitoring, clinical activities tracking, and preemptive measures for patients under medical supervision. Smartphones, being ubiquitous in modern life, serve as the foundation for monitoring health conditions within Human Activity Recognition (HAR) systems. The paper [14] introduces a Smartphone-Based Human Activities Recognition System employing the Random Forest Algorithm (RFA). The RFA classifier, comprising multiple decision trees, enhances projection accuracy by averaging across subclasses of the dataset. Evaluation parameters such as F1 Score, accuracy, precision, and sensitivity are utilized to assess the system's performance, with the proposed RFA achieving an accuracy of 98.34% compared to existing classifiers.

The method in paper [14] revolves around utilizing smartphones as a platform for Human Activities Recognition (HAR) through the implementation of the Random Forest Algorithm (RFA). The RFA, a classification algorithm consisting of multiple decision trees, operates by averaging across subclasses of the dataset to enhance projection accuracy. Four evaluation parameters - F1 Score, accuracy, precision, and sensitivity - are employed to assess the performance of the system. Through extensive experimentation and comparison with existing classifiers, the proposed RFA's accuracy is determined to be 98.34%, outperforming other methods. This paper presents a novel approach to Human Activities Recognition (HAR) using smartphones, leveraging the Random Forest Algorithm (RFA). By harnessing the capabilities of smartphones, the proposed system offers remote monitoring and tracking of clinical activities, providing preemptive measures for individuals under medical supervision[14].

Effective communication serves as the cornerstone of human relationships, whether personal or professional, facilitating societal integration and cooperation. Verbal communication relies on a shared language understood by both parties, yet around 26% of the disabled population in India, who rely on sign language for communication, face significant barriers in interacting with the general public. To address this communication gap, there is a compelling need to develop innovative solutions. The proposed idea in [15] aims to bridge this gap by creating a pair of sensor gloves capable of detecting Indian Sign Language (ISL) gestures and translating them into audible speech.

The solution given in [15] involves the development of sensor gloves equipped with various components, including Arduino Nano microcontrollers, flex sensors, touch sensors, Inertial Measurement Units (IMU), RF, and Bluetooth modules. These sensors enable the quantification of hand movements, capturing the state of both hands in numerical data sequences. Additionally, the sensor data is fed into a Machine Learning classification algorithm, enhancing the accuracy of gesture recognition. Furthermore, an accompanying smartphone application is designed to convert the recognized gestures into audible speech output, thus facilitating communication between the speech-impaired individuals and the general public. The proposed solution addresses the critical need to enhance communication accessibility for individuals using Indian Sign Language (ISL). By leveraging sensor technology and Machine Learning algorithms, the sensor gloves can accurately detect ISL gestures, enabling real-time translation into audible speech[15].

The author in [16] presents a method for Traffic Sign Detection and Recognition (TSDR), a crucial component in various intelligent transportation systems. The proposed method operates in three primary steps to effectively detect and recognize

traffic signs in real-time scenarios. The first step involves image segmentation based on thresholding of components in the HSI color space. Subsequently, the second step focuses on detecting traffic signs by processing the extracted regions of interest (ROIs). Finally, the third step undertakes the recognition of the detected traffic signs. The methodology comprises three main steps. Initially, the image is segmented to extract ROIs based on color information, ensuring only significant ROIs are considered based on predefined size and aspect ratio constraints. In the second step, shapes within the ROIs, including circular, rectangular and triangular shapes are detected using invariant geometric moments. This approach eliminates the need for machine learning algorithms for shape classification. In the recognition stage, a novel descriptor is formed by extending Histogram of Oriented Gradients (HOG) features to the HSI color space and combining them with Local Self-Similarity (LSS) features. This descriptor is then utilized in conjunction with classifiers like Random Forest and Support Vector Machine (SVM) to recognize the traffic signs. This paper proposes a comprehensive three-stage system for real-time Traffic Sign Detection and Recognition (TSDR). The first stage involves segmenting images into ROIs based on color information and size constraints. The second stage focuses on detecting shapes within the ROIs using invariant geometric moments, thereby eliminating the need for machine learning algorithms[16].

The paper [17] focuses on Automatic Sign Language Recognition (ASLR) and emphasizes the critical role of robust hand tracking and detection for accurate recognition. It introduces a new dataset comprising depth data from continuously signed American Sign Language (ASL) sentences. The analysis presented in the paper highlights the limitations of the Microsoft Kinect Skeleton Tracker (MKST), particularly in scenarios where hands are close to the body, close to each other, or when the arms cross. These limitations underscore the need for an alternative method for hand detection. The paper proposes a method called Domain-Driven Random Forest Regression (DDRFR) for predicting real-world 3D hand locations using features extracted from depth images. This method aims to overcome the shortcomings of the MKST. DDRFR employs a random forest regression approach that is tailored to the specific domain of ASL recognition. By training on the new dataset and utilizing domain-specific features, DDRFR achieves significant improvement over the MKST in hand detection accuracy.

The methodology also explores the concept of user-adaptive training, which involves collecting a small amount of sample data from an unknown signer to retrain the prediction trees. The motivation behind this approach is to enhance the system's robustness by quickly adapting to different signing styles. However, the paper notes

that user-adaptive training may not be effective when the signer's style significantly differs from those in the training data. Despite attempts to adapt, the predictions may still be inaccurate due to insufficient training data. The paper presents a detailed analysis of the limitations of the MKST in ASL recognition, particularly in scenarios involving intricate hand movements. It introduces DDRFR as a novel method for hand detection, which outperforms the MKST in terms of accuracy, especially when trained with a larger and more diverse dataset[17].

The author in [18] describes the critical role of hand gestures in communication and the burgeoning advancements in technology to interpret and utilize them effectively. As the need for automated systems to understand and respond to hand gestures escalates, the study focuses on analyzing various machine learning algorithms using datasets like MNIST and sign MNIST. These datasets are pivotal for gesture recognition, offering a wide array of images encompassing different gestures. The analysis spans across different methodologies, including feature selection, feature extraction techniques, and the evaluation of classification models, presenting a comprehensive approach to tackle the challenges posed by gesture recognition in human-computer interaction.

The study in [18] adopts a systematic methodology to evaluate the efficacy of various machine learning algorithms in classifying hand gestures. Leveraging datasets like MNIST (Modified National Institute of Standards and Technology database) and sign MNIST the research meticulously selects features and extracts discriminative attributes from the images to enhance classification accuracy. The analysis encompasses a range of machine learning algorithms including SVC, KNN, Logistic Regression, Naïve Bayes Classifier, SGDC (Stochastic Gradient Descent Classifier), CNN, Random Forest and XGBoost. The methodology employs rigorous evaluation criteria, such as confusion matrices, classification reports and accuracy measures, to compare the performance of different algorithms and derive meaningful insights into their suitability for gesture recognition tasks. Interpreting hand gestures in real-time is crucial for applications like human-computer interaction and accessibility aids, but classifying hand gestures accurately faces challenges. These include selecting features, extracting rotation and scale-invariant features and choosing appropriate classification algorithms despite advances in machine learning and computer vision. Through a meticulous analysis of various machine learning algorithms using datasets like MNIST and sign MNIST, the study highlights the algorithms in accurately classifying hand gestures for diverse applications. The findings underscore the efficacy of CNN in achieving superior classification accuracy, while also showcasing the strengths and limitations of other algorithms like SVC, KNN, logistic regression, and Naïve Bayes[18].

Effective communication, whether verbal or gestural, is essential for human interaction. However, for individuals who are deaf or mute, conveying messages through sign language presents a unique challenge, particularly when interacting with those who are unfamiliar with sign language. While interpreters can bridge this communication gap, they are often expensive and may not be readily available throughout a deaf person's life. Thus, advancements in hand gesture recognition technology hold promise for facilitating communication within the deaf and mute community, by enabling direct interpretation of sign language without the need for human intermediaries. The paper [19] aims to explore the utilization of machine learning techniques, specifically a lightweight VGG16 feature extractor and Random Forest as an ensemble classifier, to develop a robust hand gesture recognition system. By leveraging these technologies, the system seeks to overcome the communication barriers faced by individuals who rely on sign language for communication.

The methodology adopted in [19] involves the development of a hand gesture recognition system using machine learning techniques. Initially, various potential options for feature extractors, including MobileNetV2, VGG16, ResNet50V2, InceptionV3, InceptionResNetV2, DenseNet169, Xception, and NASNetMobile, were evaluated through pre-work studies. Hyperparameter settings were configured for each model to assess their performance. Based on the results, the VGG16 model was selected as the feature extractor due to its high accuracy and efficiency. A lightweight VGG16 architecture was used, leveraging transfer learning by fine-tuning weights pre-trained on ImageNet. For classification, the Random Forest ensemble method was chosen due to its ability to handle high-dimensional input features and provide robust performance.

The system's performance was evaluated on three datasets - ASL dataset, ASL Digits dataset, and NUS (National University of Singapore) Hand Posture dataset - containing hand gesture images captured under varying conditions. This comprehensive testing highlights the need for cost-effective and accurate hand gesture recognition systems that can interpret sign language from input images, overcoming limitations of traditional methods like human interpreters or expensive devices. The system achieved promising results across three datasets: ASL dataset, ASL Digits dataset, and NUS Hand Posture dataset. Specifically, the ASL dataset, comprising 69,600 training images and 8,700 test images, achieved an accuracy of 99%. The ASL Digits dataset, containing 87,000 images, attained an accuracy of 98%. Finally, the NUS Hand Posture dataset, with 8,000 training images and 1,000 test images, achieved a perfect accuracy score of 100%. These results demonstrate the effectiveness of the proposed approach in accurately recognizing hand gestures[19].

The evolution from button-type interfaces to touchscreen-type interfaces has broadened the spectrum of control methods, with hand gestures emerging as an intuitive means to interact with computers or devices. However, the inherent variability in hand postures and rotations poses significant challenges for feature selection, hindering the effective representation of all hand gestures. Previous approaches, like color filtering, struggle to adapt to varying environmental conditions and distinguish hand gestures from complex backgrounds. While technologies like IR sensors or Kinect offer solutions, this paper focuses on vision-based bare hand detection in cluttered backgrounds, given the greater ease of integration into devices. The paper [20] proposes a real-time hand gesture recognition system using a camera, leveraging random forest for hand detection and Linear Discriminant Analysis (LDA) for gesture pattern implementation, and defining three deterministic hand postures to classify gestures effectively.

The system employs a two-step process: hand detection and validation/recognition. Hand detection utilizes random forest, trained on a dataset of hand images, to detect hand regions. Here, the dataset consists of three classes, each with 1000 image data, collected under various environmental conditions. The images are resized to 20x20 resolution, from which 79,800 normalized pixel difference (NPD) features are extracted per image. These NPD features capture local intensity differences, enabling robust detection regardless of changes in brightness or size. The decision trees in the random forest are trained to classify samples into two groups based on the most discriminative feature at each node, simplifying multi-class detection into binary classification. After hand regions are detected, validation and recognition are performed using Linear Discriminant Analysis (LDA) to classify hand gestures into three defined postures. The primary challenge is accurately detecting hand gestures in cluttered backgrounds with varying environmental conditions, where traditional methods like color filtering or template matching struggle to adapt. The experimental results reveal the performance of the proposed hand gesture detection and recognition system across various image resolutions. As the resolution increases from 128x96 to 640x480, there is a notable trend in the detection time, with values escalating from 0.00869 to 0.15513. Despite this increase, the detection accuracy remains relatively consistent, ranging from 0.87 to 0.9. Similarly, the recognition time exhibits a slight rise with higher resolutions, reaching 0.000400 at 640x480 resolution, while the recognition accuracy remains consistently high, ranging from 0.90 to 0.93. These results underscore the robustness and effectiveness of the proposed system across different image resolutions, with satisfactory detection and recognition performance[20].

## Summary

The primary problem addressed is the communication barrier faced by individuals who are deaf or mute and rely on sign language, as traditional methods like human interpreters or expensive devices are often inaccessible or unaffordable. The accurate detection and recognition of hand gestures in complex environments pose challenges due to factors such as cluttered backgrounds, lighting variations, hand occlusions, and variations in hand postures and rotations. The inherent variability in signing styles and the nuances of sign language make it difficult to capture and interpret the intricate movements and expressions involved in sign language communication. Existing technologies struggle to adapt to the complex and nuanced nature of sign language, limiting their ability to provide effective communication tools. Furthermore, robust and accurate hand tracking and detection are crucial, as limitations in existing methods, particularly when hands are close to the body, close to each other, or when arms cross, hinder the overall performance of gesture recognition systems.

The methods employed in these papers span various techniques including leveraging machine learning models, improving object detection frameworks and signal processing approaches. One prominent method involves using lightweight models like VGG16 as feature extractors combined with ensemble classifiers such as Random Forest for accurate hand gesture recognition. Additionally multi-branch attention-based graph neural networks and general deep learning models have been proposed to recognize hand gestures from 3D skeleton data leveraging spatial-temporal and temporal-spatial features. To address the challenges of complex backgrounds and varying conditions researchers have focused on improving object detection models. This includes assessing the performance of advanced models like Yolov7 and Yolov7x for hand detection tasks as well as introducing improved frameworks like an enhanced YOLOv5 with efficient layer aggregation network modules and attention mechanisms for better feature extraction and network robustness. Signal processing techniques have also been explored particularly in the context of electromyography (EMG) signals. This involves preprocessing EMG signals to remove noise artifacts and extracting time-domain features which are then fed into machine learning classifiers like Support Vector Machines (SVM) and K-Nearest Neighbors (KNN) for hand gesture recognition. Furthermore machine learning algorithms such as K-Nearest Neighbors (KNN), Support Vector Machines (SVM), and Convolutional Neural Networks (CNN) have been employed on datasets of sign language alphabets for recognition and translation purposes aiming to bridge the communication gap between sign language users and non-users.

# Chapter 3

## Software Dependencies

### 3.1 Programming with Python

Python is a popular, high-level, and versatile programming language created by Guido van Rossum and first released on February 20, 1991. The name “Python” is derived from the British comedy troupe Monty Python’s Flying Circus. Python is a general-purpose language widely used for developing websites, software, task automation, data analysis and data visualization.

Python offers several features that make it an attractive choice for developers. It has a simple and clean syntax, making it easy to learn and code. The emphasis on readability and use of indentation make the code more understandable and maintainable. Python is freely available and open-source, encouraging community contribution and widespread adoption. It comes with a vast standard library that provides a wide range of modules and functionalities out-of-the-box. Additionally, Python is a cross-platform language, meaning code written on one platform can run on other platforms with minimal or no modifications. Python supports both object-oriented and procedural programming paradigms, providing flexibility to developers. Furthermore, Python can be extended with modules and libraries written in other languages like C and C++ allowing for performance optimization and integration with existing code-bases[21].

In our project Python version 3.11 was employed. This was chosen because it offers several improvements and features that are beneficial for machine learning tasks. First it introduces performance enhancements such as faster startup times improved memory usage and more efficient handling of data structures leading to faster execution of machine learning models and data processing pipelines. Secondly Python 3.11 includes better error handling and reporting mechanisms which allows for more granular handling of exceptions beneficial for robust machine learning workflows. By leveraging the latest version of Python (3.11) and its improvements our project aims to develop an efficient and reliable sign language recognition system using machine learning techniques.

## 3.2 Sci-kit Learn

Scikit-learn, often abbreviated as sklearn, is an open-source machine learning library for Python. It's built on top of other popular Python libraries like NumPy, SciPy and matplotlib, providing a cohesive ecosystem for machine learning tasks. Its primary focus is on ease of use, efficiency and accessibility making it an ideal choice for both beginners and seasoned machine learning practitioners.

One of scikit-learn's standout features is its extensive collection of machine learning algorithms. It covers a wide range of tasks, including classification, regression, clustering, dimensionality reduction and model selection. For classification tasks, scikit-learn offers algorithms like Support Vector Machines (SVM), Decision Trees, Random Forests and k-Nearest Neighbors (kNN). Similarly, for regression tasks, it provides Linear Regression, Ridge Regression, Lasso Regression and more. This diverse selection of algorithms caters to various needs and preferences, allowing users to choose the most suitable approach for their specific problem. We have used scikit-learn version 1.2.0 leveraging its robust tools and algorithms to build and evaluate our machine learning models[22].

Scikit-learn prioritizes simplicity and consistency in its API design. The library follows a uniform interface across all its algorithms, making it intuitive for users to switch between different models and experiment with various techniques seamlessly. This consistency extends to the configuration and hyperparameter tuning process where scikit-learn provides standardized methods for setting parameters and conducting cross-validation. As a result users can focus more on experimenting with models and less on dealing with intricate implementation details. Scikit-learn also provides utilities for data preprocessing, feature extraction, and model evaluation, further streamlining the machine learning workflow[23]. These utilities include functions for standardizing or scaling data, handling missing values, encoding categorical variables, and splitting data into training and testing sets. Additionally, scikit-learn offers built-in tools for assessing model performance through various metrics such as accuracy, precision, recall and F1-score

**NumPy:** As the fundamental package for numerical computing in Python, NumPy provides essential tools for creating and manipulating multidimensional arrays and matrices, along with a vast collection of mathematical functions to operate on these arrays efficiently. Its array-oriented computing capabilities enable concise and expressive code for tasks such as mathematical operations, linear algebra and statistical analysis making it an indispensable component of the scientific Python ecosystem.

Matplotlib: This is a library in Python that facilitates the creation of high-quality static, interactive, and animated visualizations for data analysis and presentation. With a versatile and user-friendly interface, Matplotlib allows users to generate a wide range of plots, including line plots, scatter plots, bar plots, histograms, heatmaps, and more, with extensive customization options to tailor the appearance and style of plots to specific requirements. Its seamless integration with other scientific computing libraries, coupled with its extensive documentation and active community support, makes Matplotlib an indispensable tool for data visualization tasks in Python.

sklearn.model\_selection: This module provides tools for model selection and evaluation, including functions for splitting data into training and testing sets, performing cross-validation, tuning hyperparameters using grid search or randomized search, and evaluating model performance using various metrics. sklearn.metrics module offers a wide range of metrics for evaluating the performance of machine learning models. It includes metrics for classification, regression, clustering, and ranking tasks, such as accuracy, precision, recall, F1-score, mean squared error, and ROC curves.

sklearn.ensemble: This submodule provides a collection of ensemble learning techniques in scikit-learn, including Random Forests, Gradient Boosting, AdaBoost, Extra Trees, and Voting classifiers/regressors. Ensemble methods combine predictions from multiple base estimators to improve model performance and robustness. With these algorithms, users can leverage the diversity of individual models to achieve superior predictive accuracy across various machine learning tasks, such as classification, regression, and anomaly detection. sklearn.ensemble plays a pivotal role in enhancing model generalization and handling complex data patterns effectively.

### 3.3 Open CV

OpenCV (Open Source Computer Vision Library) is an open source computer vision and machine learning software library. OpenCV was built to provide a common infrastructure for computer vision applications and to accelerate the use of machine perception in the commercial products. OpenCV is a huge open-source library for computer vision, machine learning, and image processing. OpenCV supports a wide variety of programming languages like Python, C++, Java, etc. It can process images and videos to identify objects, faces, or even the handwriting of a human[24].

OpenCV was built for maximum efficiency and performance of computing-intensive vision tasks. Therefore, it has a strong focus on real-time applications of AI vision.

The software is written in optimized C and is able to take advantage of multicore processors (multi-threading). OpenCV provides a module called ml that has many machine learning algorithms bundled into it. Some of the algorithms include Bayes Classifier, K-Nearest Neighbors, Support Vector Machines, Decision Trees, Neural Networks. Parameters of Open CV:

- Src - It denotes the source of the image.
- Dst- It denotes the destination image of the same size.
- Depth - It denotes the output image depth.
- Ksize - It blurs the kernel size.
- Anchor - It denotes the anchor points.
- Normalize - It is the flag, specifying whether the kernel should be normalized or not.

OpenCV offers support for the image formats Windows bitmap (bmp), portable image formats (pbm, pgm, ppm) and Sun raster (sr, ras). OpenCV includes GPU module that contains all GPU accelerated stuff. Supported by NVIDIA the work on the module, started in 2010 prior to the first release in Spring of 2011. OpenCV is open source and released under the Apache 2 License. It is free for commercial use. We can use any OS—macOS, Windows, and Linux-based OS—with this book. recommend to have at least 4 GB RAM in your system. Opencv-python version 4.7.0.68 is used in our project providing Python bindings for OpenCV functionalities, allowing developers to work with OpenCV within Python environments and it enables easy access to camera devices and real-time image processing making it a popular choice for developing camera applications across various platforms and domains[25].

### 3.4 Pandas

Pandas a widely-used Python library plays a pivotal role in hand sign detection projects by facilitating efficient data management and analysis. In these projects datasets often encompass various attributes and annotations for each sign, such as image features, labels and metadata. Its DataFrame structure allows seamless handling, loading, and parsing of datasets from various file formats like CSV files or databases. With Pandas, researchers can easily inspect dataset characteristics such as class distributions and identify missing or erroneous data, facilitating thorough exploration and manipulation[26].

Once the data is loaded, Pandas aids in data preprocessing tasks critical for preparing the dataset for model training. It offers powerful tools for cleaning data, handling missing values, and transforming features. For instance, Pandas can be utilized to normalize image features, encode categorical labels, or extract relevant features from raw data. Furthermore, Pandas facilitates data exploration and analysis, enabling researchers to gain insights into the dataset's properties and distributions. Visualization techniques like histograms, scatter plots, or box plots can be employed to visualize the distribution of sign features and identify any underlying patterns or correlations between variables.

Pandas version 2.0.3 is utilized for efficient data management and analysis including handling pickle files. Pandas seamlessly integrates with other Python libraries commonly employed in hand sign detection projects, such as scikit-learn for model training and evaluation, and Matplotlib for visualizing results. Data preprocessed and organized using Pandas can be seamlessly fed into machine learning models for training and testing, enabling the development of accurate and robust hand sign detection systems.

### 3.5 Media Pipe

Mediapipe, a versatile open-source framework developed by Google, serves as an invaluable tool in the realm of sign language detection. Its primary aim is to streamline the creation of machine learning pipelines, offering a rich repository of pre-built components and models specifically tailored for various perception tasks, including but not limited to hand tracking and pose estimation. In addition to hand tracking, Mediapipe encompasses components dedicated to pose estimation. By understanding the positioning and orientation of body parts, Mediapipe enhances the overall accuracy of sign language recognition systems, ensuring more robust and reliable performance. OpenPose is more suitable for detecting images, while MediaPipe Pose is better suited for detecting videos. MANO (Model for Articulated Hands with Objectives) is a widely used hand model in the vision and graphics community, offering detailed 3D hand shapes and poses. While MediaPipe's hand tracking model version 0.9.0.1 is a machine learning model that can detect and track the 21 3D landmarks of a hand from a video stream in real-time [27].

As shown in the Figure 3.1 the model is based on a combination of a single-shot detector (SSD) network and a pose estimation model that predicts the 3D coordinates of the hand landmarks. The hand landmark model bundle contains a palm detection

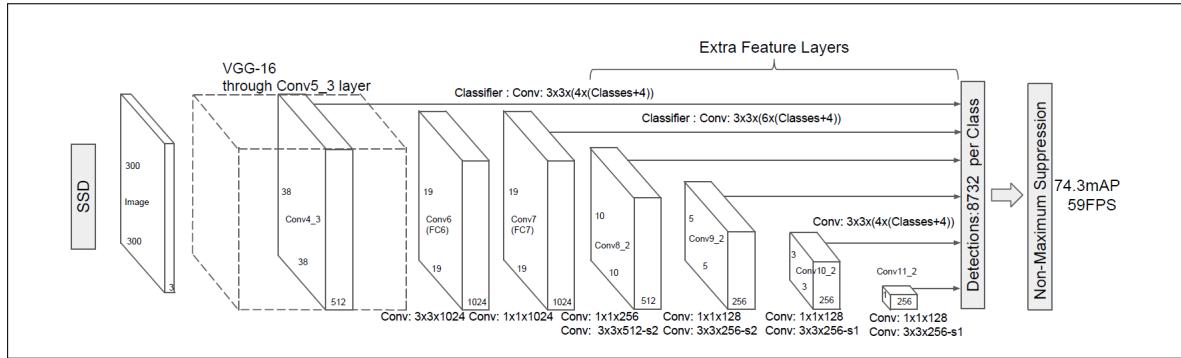


Figure 3.1: SSD Architecture

model and a hand landmarks detection model. The palm detection model detects the presence of hands and returns an oriented hand bounding box, while the hand landmark model operates on the cropped image region defined by the palm detector and returns 3D landmarks of the hand. The MediaPipe Model uses a combination of deep learning architectures. It uses a Single Shot Detector (SSD) network for palm detection, Hand landmark model based on a fully convolutional neural network (FCN). The SSD network is a popular object detection architecture that can detect objects in an image with a single forward pass of the network. The hand landmark model is based on a fully convolutional neural network (FCN) that predicts the 3D coordinates of the hand landmarks. The FCN is trained on a large dataset of annotated hand images and is optimized for speed and accuracy.

At the core of the SSD architecture lies the VGG-16 network, a well-established and extensively validated convolutional neural network (CNN) that serves as the backbone for feature extraction. The VGG-16 network is employed up to the Conv5\_3 layer, leveraging its proven ability to capture rich visual features from input images through a series of convolutional and pooling layers. These feature maps serve as the foundation for subsequent object detection and localization tasks. Building upon the VGG-16 backbone, the SSD architecture introduces a series of additional feature layers that facilitate multi-scale object detection. These extra feature layers operate in parallel, allowing the model to effectively detect objects of varying sizes and scales within the input image. The first set of extra feature layers comprises the Conv6 and Conv7 layers, which are responsible for generating feature maps at different resolutions. Specifically, the Conv6 layer produces a 19x19 feature map, while the Conv7 layer generates a 10x10 feature map. These feature maps capture both fine-grained and coarse-grained visual information, enabling the detection of small and large objects, respectively.

Subsequent to these feature extraction layers, the SSD architecture employs a series of convolutional predictors, which are tasked with classifying the detected objects and regressing their bounding box coordinates. These predictors operate on the extracted feature maps at various scales, allowing the model to make predictions at multiple resolutions simultaneously. The Figure 3.1 provides a detailed illustration of the convolutional predictors employed in the SSD architecture, such as the Conv8\_2 layer, a 4x4 convolutional filter applied to the Conv7 feature map, and the Conv9\_2 layer, a 3x3 convolutional filter applied to the Conv6 feature map. These convolutional predictors are meticulously designed to classify objects and predict their bounding boxes based on the extracted visual features.

Furthermore, the SSD architecture [28] incorporates additional convolutional predictors, such as Conv10\_2 and Conv11\_2, which operate on feature maps of lower resolutions (e.g., 5x5 and 3x3). These predictors enhance the model's ability to detect small objects by leveraging the higher-resolution feature maps, thereby improving overall detection accuracy. In addition to the convolutional predictors, the SSD architecture employs a series of classifiers, denoted as “Classifier: Conv: 3x3x(4x(Clases+4))” and “Classifier: Conv: 3x3x(6x(Clases+4))”. These classifiers are responsible for predicting the class probabilities and bounding box offsets for the detected objects, based on the features extracted by the convolutional predictors. A key aspect of the SSD architecture is its ability to generate a large number of bounding box predictions, known as default bounding boxes or anchor boxes, with varying aspect ratios and scales. Specifically, the SSD model creates 8,732 bounding box predictions per object, which are then filtered and refined through a non-maximum suppression (NMS) block. This block removes duplicate predictions by analyzing the confidence scores of each bounding box and retaining the top 200 predictions with the highest confidence scores.

Further to evaluate the accuracy of the predicted bounding boxes, the SSD model employs the Intersection-over-Union (IoU) metric. This metric measures the overlap between the predicted bounding box and the ground truth bounding box, which is provided as part of the input data. The bounding box with the maximum IoU score, indicating the highest degree of overlap with the ground truth, is selected as the final prediction for the object's orientation and location. The 3.1 further highlights the computational efficiency of the SSD architecture, with the Detection SSD model achieving a remarkable 74.3 mean Average Precision (mAP) at 59 Frames Per Second (FPS) on a specific benchmark. This impressive performance is attributed to the careful design and optimization of the various layers and components within the SSD architecture, enabling real-time object detection capabilities.

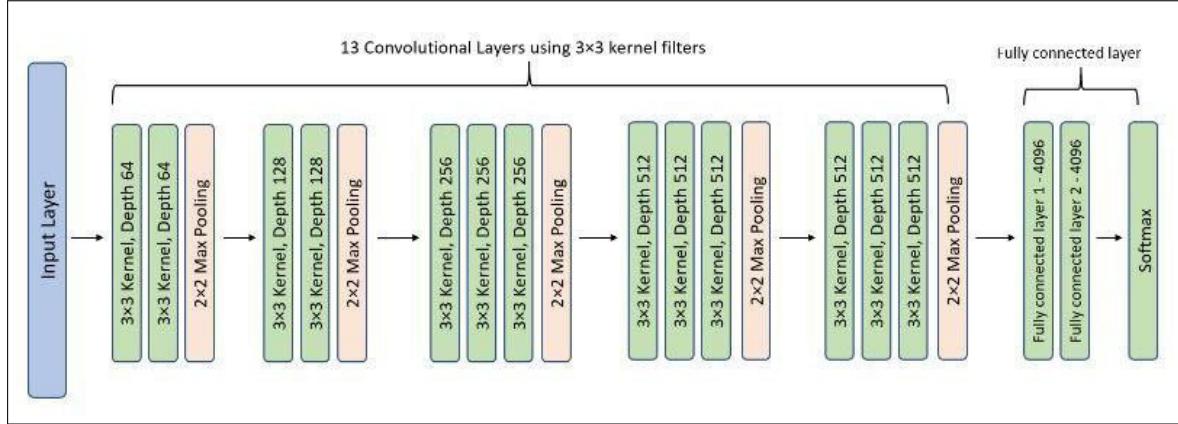


Figure 3.2: VGG16 Architecture

Overall the Single Shot Detector (SSD) architecture presents a comprehensive and effective approach to object detection by leveraging the strengths of convolutional neural networks and introducing specialized layers for multi-scale feature extraction and object prediction. The seamless integration of the VGG-16 backbone, extra feature layers, convolutional predictors, classifiers, and non-maximum suppression block enables the SSD model to achieve exceptional accuracy and real-time performance, making it a highly influential and widely adopted architecture in the field of object detection across diverse industrial and research domains.

The 3.2 depicts VGG16 architecture [29], proposed by researchers at the University of Oxford, has become a seminal deep learning model for image classification tasks. This architecture employs a principled approach to convolutional neural network design, leveraging a series of well-structured layers that progressively extract and refine visual features from input images. At its inception, the VGG16 model ingests input images of a fixed size, 224 by 224 pixels, in the RGB color space. This initial layer ensures a consistent input format, which is a prerequisite for the subsequent convolutional operations. The fixed input size allows the model to learn spatial relationships and patterns within a standardized image representation.

The heart of the VGG16 architecture lies in its 13 convolutional layers, which are responsible for extracting visual features from the input images. These layers are organized into several blocks, with each block comprising two or more convolutional layers followed by a max-pooling layer. The convolutional layers within each block employ 3x3 kernel filters, a strategically chosen size that enables the network to capture intricate patterns while maintaining computational efficiency. The first two convolutional layers in the VGG16 model contain 64 filters each, tasked with detecting low-level

visual features such as edges, corners, and simple shapes. As the network progresses through subsequent blocks, the number of filters gradually increases, reaching 128, 256, and eventually 512. This progressive increase in filter count allows the model to capture increasingly complex and abstract visual features, enabling it to learn higher-level representations of the input data.

Interspersed between the convolutional layer blocks are max-pooling layers, which serve a crucial role in reducing the spatial dimensions of the feature maps. By down-sampling the feature maps, these layers introduce translation invariance and robustness to small distortions or shifts in the input image. This property is particularly valuable in image classification tasks, where the model needs to generalize effectively across diverse input samples. Upon completing the convolutional and max-pooling layers, the VGG16 architecture transitions to a series of three fully connected layers. These layers are responsible for integrating the high-level visual features extracted by the previous layers and performing the final classification task. The first two fully connected layers have 4096 units each, while the final layer contains 1000 units, corresponding to the 1000 class labels in the ImageNet dataset, which was used for pretraining the VGG16 model.

The output of the final fully connected layer is then passed through a softmax activation function, which converts the raw output values into a probability distribution over the target classes. This softmax layer ensures that the sum of the output probabilities equals 1, allowing for the model to make a confident classification decision based on the highest probability class. Throughout the VGG16 architecture, the Rectified Linear Unit (ReLU) activation function is employed after each convolutional and fully connected layer. ReLU introduces non-linearity to the network, enabling it to learn complex, non-linear mappings between the input data and the target classes. This non-linearity is crucial for the model to capture the intricate relationships present in real-world image data effectively. The VGG16 architecture has demonstrated remarkable performance on various image classification benchmarks, thanks to its well-structured design and ability to learn hierarchical visual representations. However, it is worth noting that the model's efficacy comes at the cost of a substantial number of parameters (approximately 138 million), which can pose challenges in terms of computational resources and memory requirements, particularly in resource-constrained environments.

Overall the VGG16 architecture presents a principled approach to convolutional

neural network design, leveraging a series of convolutional, pooling, and fully connected layers to progressively extract and refine visual features from input images. Its strategic use of 3x3 kernel filters, progressive filter count increase, and max-pooling layers contribute to its effectiveness in capturing complex visual patterns and achieving state-of-the-art performance on image classification tasks.

# Chapter 4

## Proposed Classifier Implementation

This chapter presents the implementation of the Sign Language Classifier encompassing the project's workflow and design methodology outlining the step-by-step techniques employed which provides a thorough understanding of the system's development process.

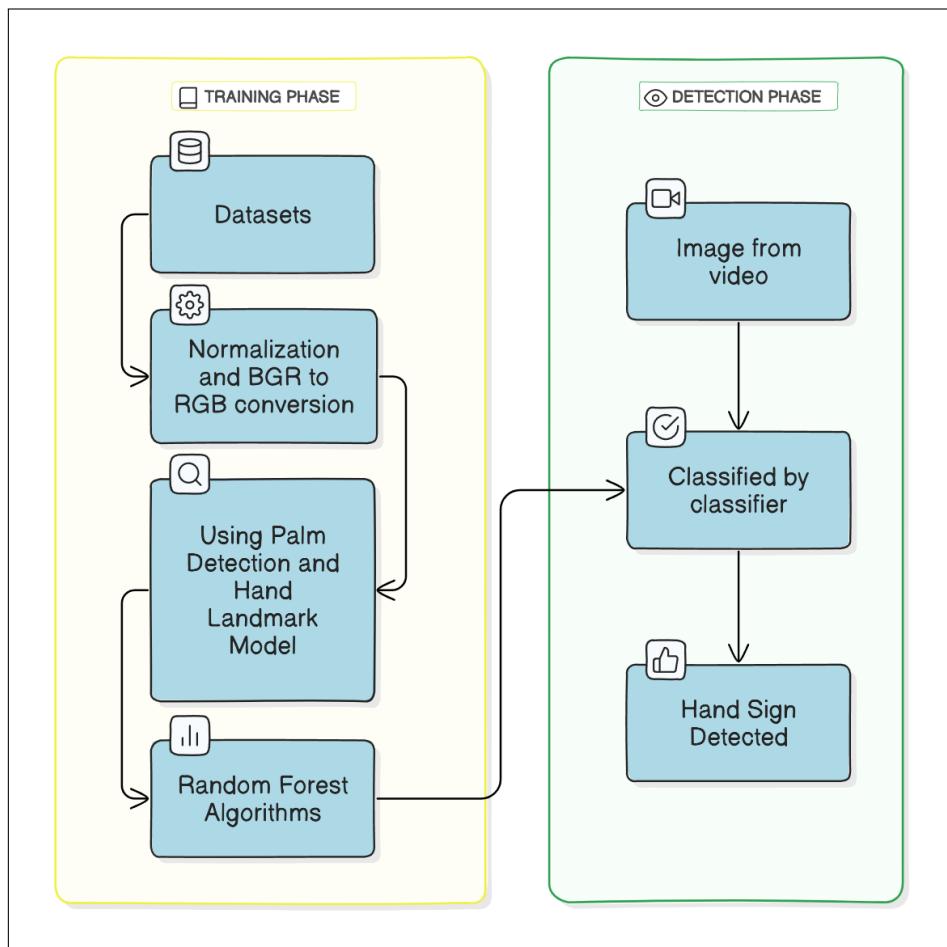


Figure 4.1: Flow Chart of Hand Sign Recognition System

The Figure 4.1 illustrates a comprehensive workflow for a computer vision task involving hand gesture recognition using the MediaPipe library and machine learning techniques. The workflow begins with importing the necessary dependencies, such as MediaPipe (version 0.9.0.1), scikit-learn (version 1.2.0), OpenCV (version 4.7.0.68), and pandas (version 2.0.3). These libraries provide functionalities for hand tracking,

machine learning, image processing, and data manipulation, respectively. The next step is capturing a custom dataset, which is a crucial component for training the hand gesture recognition model. The dataset consists of 26 classes, one for each alphabet of the American Sign Language (ASL). For each class, 100 images are captured, resulting in a total of 2,600 images. The images have a resolution of 640x340 pixels and are in JPG format. Data augmentation techniques, such as rotation, scaling, flipping, or adding noise to images, are performed in the workflow, and final dataset comprises 200 images. OpenCV is used for image capturing, and custom key-strokes are employed to initiate, continue, and terminate the capturing process for each class.

Pre-processing steps are applied to the extracted features, including converting the image color space from BGR to RGB and normalizing the feature values. Feature extraction is a critical step in the workflow, where relevant information is extracted from the input images to represent the hand gestures effectively. The MediaPipe Hand solutions are employed, which include transfer learning models, hand models, and various backend models like SSD, YOLO, and CNN. These models facilitate hand detection, landmark localization, and drawing utilities. The static images are processed to extract 21 key points of the hand, and a total of 42 features are derived from each hand, representing the x and y coordinates of the landmarks. Post-processing operations involve data splitting, where the dataset is divided into training and testing subsets, and data serialization, where the features and corresponding labels are stored in a suitable format for further processing.

The next stage is training a classifier using the extracted features and labels. The workflow suggests using a Random Forest Classifier with an Ensemble approach. Specific hyperparameters are provided, including the train\_test\_split ratio (20:80), the impurity measures (gini or entropy), the number of decision trees (100), the maximum samples for splitting nodes (max\_samples\_split=2 and max\_samples\_leaf=1), bootstrapping (TRUE), and the random state (42). Model evaluation is an essential step to assess the performance of the trained classifier. Metrics such as accuracy, precision, recall, and F1-score are computed, and a detailed classification report is generated. After training and evaluation, the workflow proceeds to testing, where the trained model is loaded and initialized using the MediaPipe hands module. During testing, the model predicts the hand gestures based on the input images, and the predictions are visualized, allowing for the display of unknown gestures and their corresponding probabilities.

Overall the workflow encompasses importing dependencies, capturing a custom

dataset, feature extraction using MediaPipe Hand solutions, pre-processing and post-processing steps, training a Random Forest Classifier, model evaluation, testing with the trained model and results.

## 4.1 Data Set Collection and Preprocessing

The dataset collection process for the hand sign recognition project is meticulously designed to systematically gather images representing various hand signs forming the foundational dataset for subsequent model training and evaluation. The process unfolds in a series of steps, each contributing to the systematic acquisition and organization of the dataset.

To commence, we initialized and configured a directory structure to systematically store the collected images. This structured approach ensured orderly storage and facilitated easy retrieval during subsequent stages of the project. Upon setup the system interfaced with the user through the utilization of keystrokes. The user prompted by on-screen instructions could initiate the data collection process for each hand sign class by pressing ‘S’. Alternatively they had the option to navigate to the next hand sign class by pressing ‘N’ or exit the program altogether by pressing ‘Q’.

Once data collection for a particular hand sign class was initiated, the system seamlessly transitioned into capturing images in real-time using the webcam feed. During this phase, the user’s interaction with the system allowed us to collect a pre-determined number of images per class, which we set to 100 images. As images were captured, they were promptly stored in JPG format within designated directories corresponding to their respective hand sign classes. This systematic naming convention ensured proper organization and simplified the retrieval process during subsequent stages of model development and evaluation. To ensure diversity in our dataset, we collected images of 3 different individuals, capturing hand signs under varying lighting conditions. Additionally, we accounted for the diverse skin tone colors of Indian people ensuring that our dataset represented a broad range of variations.

Then these images undergo several preprocessing steps to prepare the data for model training. The images are converted from the BGR (Blue-Green-Red) to the RGB (Red-Green-Blue). This conversion is necessary because neural networks typically expect input data in the RGB format discussed later. After the color space conversion data augmentation techniques are applied to enhance the diversity and robustness of the dataset. Data augmentation is a crucial step in improving the model’s performance

and generalization ability. This process involves applying various transformations to the existing images, such as rotation, flipping and scaling. The augmented images along with the original ones form the final preprocessed dataset of 200 images per class.

## 4.2 Mediaipe Feature Extractor

The feature extraction phase is a critical component of the hand sign recognition project, responsible for extracting meaningful and representative features from input images or video frames. This phase leverages the power of MediaPipe's hand tracking model, which employs a sophisticated multi-step process integrating deep learning algorithms to detect and localize hand landmarks accurately. Before delving into the core feature extraction process the input frame undergoes a preprocessing step as discussed in previous section 4.1 to optimize performance and align with the mediapipe's model requirements. The step include converting the color space(BGR to RGB). Pre-processing ensures that the input data is in a format that the model can efficiently process, improving overall accuracy and computational efficiency.

The process begins with hand landmark extraction using MediaPipe's hands.process method [27] . This method processes the input image and extracts hand landmarks. The first major step in the feature extraction process is hand detection. MediaPipe's hand tracking model employs a Single Shot Detector (SSD) architecture for this purpose as discussed earlier in section 3.5. The SSD architecture is a powerful object detection algorithm that analyzes the input image at multiple scales to identify potential hand regions based on visual cues such as skin tone, shape and texture. The input image is passed through a base Convolutional Neural Network (CNN), which extracts feature maps at multiple scales. At each location in these feature maps, the SSD network predicts a set of bounding boxes with different aspect ratios and scales along with their associated confidence scores for the presence of a hand within each bounding box. Non-maximum suppression is applied to filter out predicted bounding boxes with low confidence scores or significant overlap with higher-scoring boxes, retaining only the most likely bounding boxes containing hands.

Once the hand regions are identified, the module proceeds to localize key points or landmarks on the hand. These landmarks totaling 21 keypoints as shown in Figure 4.2. For each input image, the module extracts a total of 42 coordinate values – 21 for the x-coordinates and 21 for the y-coordinates encompass critical hand features such as fingertips, knuckles, palm center and wrist. The accurate localization of these landmarks is crucial for building a robust and reliable hand sign recognition system. The

landmark localization process is achieved through a combination of convolutional neural network (CNN) layers and regression techniques tailored for landmark localization tasks. Convolutional layers apply a set of learnable filters to the input image, computing dot products between the filter weights and the input values, and producing feature maps as output. Pooling layers downsample the feature maps by aggregating information from neighboring regions, reducing the spatial dimensions and introducing invariance to small translations or distortions. Fully connected layers take the final feature maps, flatten them, and perform complex non-linear transformations to output the predicted landmark coordinates.

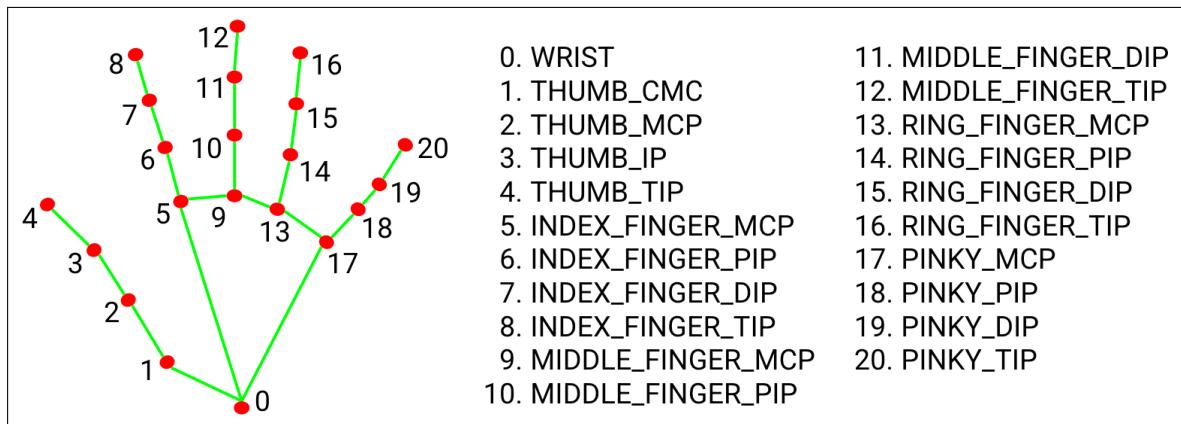


Figure 4.2: Media Pipe Hand Landmark

The resulting `hand_landmarks` object contains the normalized coordinates ( $x$ ,  $y$ ) of each landmark. During data preprocessing, the code iterates through the detected hands and landmarks, extracting the normalized  $x$  and  $y$  coordinates for each landmark. These coordinates are stored in separate arrays ( $x_{\cdot}$  and  $y_{\cdot}$ ) and then normalized by subtracting the minimum value from each coordinate to ensure all values remain between 0 and 1. This normalization is crucial because the dataset includes images of varying hand sizes and orientations. The normalized coordinates are stored in the `data_aux` array representing the features for a single hand instance.

In data serialization, the extracted features (`data_aux`) and corresponding labels (`dir_`) are appended to separate lists (`data` and `labels`). Once all images have been processed, the `data` and `labels` lists are serialized using the `pickle` module and stored in a file (`test.pickle`). The `data_aux` list stores the normalized coordinates of the hand landmarks for a single hand instance, with each pair of consecutive elements representing the  $x$  and  $y$  coordinates of a landmark. For example, if `data_aux` contains  $[x_1, y_1, x_2, y_2, \dots, x_{21}, y_{21}]$ , the first two elements ( $x_1, y_1$ ) correspond to the coordinates of the first landmark, the next two elements ( $x_2, y_2$ ) correspond to the second

landmark, and so on, for a total of 21 landmarks. The data.pickle file contains a dictionary with two keys: ‘data’ and ‘labels’. The ‘data’ key maps to a list where each element is a data\_aux list representing the features for a single hand instance, and the ‘labels’ key maps to a list of corresponding labels (class names) for each hand instance. These pickle files serve as the input data for further processing and analysis, including training a machine learning model for hand sign recognition.

## 4.3 Design of Random Forest Classifier

A Random Forest classifier is an advanced ensemble learning technique primarily used for classification and regression tasks. Introduced by Leo Breiman in 2001 it builds upon the concept of decision trees to enhance prediction accuracy and robustness. The Random Forest algorithm works by first creating multiple decision trees from randomly selected subsets of the training data. For each decision tree it randomly selects a subset of features from the total features when growing the tree. Each decision tree grows to its maximum depth without any pruning. To make a prediction for a new data instance the Random Forest passes the instance to each of the decision trees. The key advantages of Random Forest Classifiers include high accuracy due to combining multiple decision trees robustness to noise and outliers, ability to handle high dimensional data and relative simplicity compared to other ensemble methods. Additionally, Random Forests provide feature importance estimates which can help identify the most relevant features in a dataset[30].

### 4.3.1 Model Parameters

1. Splitting Criterion: Decision trees split nodes based on impurity measures. ‘Gini’ is one of the criteria used to measure node impurity, and it’s commonly employed in Random Forest models. For this model, the splitting criterion has been set to ‘gini’.
2. Maximum Depth: This parameter controls the maximum depth of the decision trees. For this model, the max\_depth parameter has been set to None, which allows the trees to expand until all leaves are pure or contain the minimum number of samples specified by min\_samples\_split.
3. Maximum Features: This parameter determines the number of features to consider when looking for the best split at each node. For this model, the max\_features parameter has been set to ‘sqrt’, which means that the square root of the total number of features is considered.

4. Minimum Leaf Nodes: This parameter specifies the minimum number of leaf nodes that must be present in each decision tree. For this model this parameter has been set to 1 meaning that each decision tree in the Random Forest must have at least one leaf node.
5. Max Leaf Nodes: This parameter specifies the maximum number of leaf nodes that can be created in each decision tree. Setting this parameter can help control the size and complexity of the individual trees in the Random Forest ensemble. For this model this parameter has been set to 2, which means that each decision tree in the Random Forest can have a maximum of 2 leaf nodes.
6. Number of Estimators: This parameter defines the number of decision trees in the Random Forest ensemble. For this model, the n\_estimators parameter has been set to 100 which means that the Random Forest ensemble consists of 100 decision trees.

For other parameters their default values are used, like bootstrapping ('True') and class weights ('balanced'), promote model diversity and reduces biases towards dominant classes. These settings enhance robustness and generalization by incorporating randomness in tree construction and adjusting class weights based on frequency.

#### 4.3.2 The Decision Making Process

The Random Forest classifier operates using an ensemble learning technique called bagging (Bootstrap + Aggregating). The process begins with bootstrap sampling where multiple decision trees are trained on different bootstrap samples of the training data[31]. In bootstrap sampling data points are randomly selected with replacement resulting in some points appearing multiple times while others may not appear at all in a sample. This creates diversity within the ensemble as each tree is trained on a slightly different subset of the original data.

Feature selection or feature bagging is another crucial step where a random subset of features from the input dataset is selected for each decision tree. This reduces correlation between trees and helps prevent overfitting. The number of features considered at each split point is controlled by the max\_features parameter, which can be set to 'sqrt', 'log2' or a specific integer value.

In the tree construction phase each decision tree is built using a bootstrap sample. At each node the algorithm searches for the best split among the subset of features determined by max\_features. The quality of the split is evaluated using criteria such

as Gini impurity or entropy. Gini impurity measures the probability of incorrectly classifying a randomly chosen element if labeled according to the class distribution in the set with lower values indicating more homogeneous sets. The formula for Gini impurity  $I_g$  for a set of data with K classes is given by below equation 4.1.

$$I_g = 1 - \sum_{i=1}^K P_i^2 \quad (4.1)$$

$P_i$  is the probability of randomly selecting an element of class i from the set. The sum is taken over all K classes. This formula calculates the Gini impurity by summing the squared probabilities of each class and subtracting the result from 1. A lower Gini impurity indicates a more homogeneous set of samples with respect to the target variable. Entropy is a measure of impurity or disorder in a set of data points. It quantifies the uncertainty associated with a given set of data points. The formula for Entropy  $I_h$  for a set of data with K classes is given by below equation 4.2.

$$I_h = - \sum_{i=1}^K P_i \log_2(P_i) \quad (4.2)$$

$P_i$  is the probability of randomly selecting an element of class i from the set. The sum is taken over all K classes. Minimizing the entropy of the resulting subsets is equivalently to maximizing the information gain. Gini impurity is generally faster to compute than entropy because it does not involve logarithmic calculations. This makes the tree-building process quicker. Hence gini is preferred over entropy. During the prediction phase each tree in the forest independently predicts the class of the input sample. For classification tasks the result (most frequent class) of the predictions from all the trees is taken as the final prediction. This process is called aggregation. In other words the class that receives the most votes among all the trees is selected as the predicted class.

### 4.3.3 Evaluation Of Model

The evaluation pipeline in the project relies on two fundamental components: the StratifiedKFold function and the GridSearchCV function, both pivotal for robust model assessment and hyperparameter tuning.

StratifiedKFold is a cross-validation technique widely used in machine learning to ensure fair and unbiased model evaluation, especially with imbalanced datasets. It divides the dataset into ‘k’ folds while maintaining the original class distribution across each fold. This ensures that each fold has a proportional representation of the different classes, minimizing biased evaluations[32]. The value for k=5 the dataset is

---

split into 5 folds, each having a class distribution similar to the original dataset.

The process begins with shuffling the dataset to eliminate any potential ordering bias that might exist. Once shuffled, the data is partitioned into ‘k’ approximately equal-sized folds. Unlike standard k-fold cross-validation, where data is partitioned randomly into folds without considering class distribution, StratifiedKFold ensures that each fold contains a balanced representation of classes. During each iteration of the cross-validation process, one of the folds is held out as the validation set, while the remaining folds are used for training the model. This process is repeated ‘k’ times, with each fold being used once as a validation set. By rotating the roles of training and validation sets across folds, StratifiedKFold trains and evaluates the model multiple times on different data subsets. Its key advantage is providing reliable performance estimates across various class distributions, crucial for assessing the model’s generalization to unseen data, especially with imbalanced classes. Maintaining a representative mix of classes in each fold helps detect issues like overfitting or underfitting.

GridSearchCV is a technique used for hyperparameter tuning in machine learning, essential for optimizing model performance by systematically exploring a predefined grid of hyperparameters. The process involves exhaustively searching through all possible combinations of hyperparameters and evaluating each combination using cross-validation to identify the set that produces the best results[33]. Firstly a grid of hyperparameters and their corresponding values given in Table 4.1 are defined. For instance, in a decision tree classifier, parameters like maximum depth, minimum samples split, and criterion are included in the grid. GridSearchCV then performs cross-validation by splitting the training data into ‘k’ folds. It trains the model on ‘k-1’ folds and evaluates its performance on the remaining fold, repeating this process ‘k’ times. During each iteration, GridSearchCV tests every combination of hyperparameters from the defined grid. The model is trained and evaluated using each combination, and the performance metric specified (such as accuracy, precision, or F1-score) is computed. This comprehensive evaluation allows GridSearchCV to identify the hyperparameter combination that yields the best performance across all folds.

After evaluating all combinations GridSearchCV selects the set of hyperparameters that maximizes the chosen performance metric in terms of accuracy. Some of the combination are shown in Table 4.2. This optimal set of hyperparameters is then used to train the final model on the entire training dataset. By fine-tuning hyperparameters based on cross-validated performance in terms accuracy. The total number of unique hyperparameter combinations is calculated by multiplying the number of options for

Table 4.1: Parameters Grid Values

Hyperparameter	Values
n_estimators	50, 100, 200
max_depth	None, 10, 20
min_samples_split	2, 5, 10
min_samples_leaf	1, 2, 4
max_features	sqrt, log2
criterion	gini, entropy, log_loss

each hyperparameter which results in 486 combinations. Therefore when fitting 5 folds for each of these 486 combinations the model is trained and evaluated a total of 2430 times ensuring a thorough and robust evaluation across different hyperparameter settings and data splits. GridSearchCV enhances model effectiveness and generalization, ensuring that the model is well-suited to the characteristics of the dataset and the learning task at hand.

Table 4.2: GridSearchCV Result

Parameters	Combination 1	Combination 2	Combination 3	Combination 4
Criterion	Entropy	Gini	Entropy	Gini
Max Depth	10	20	None	None
Max Features	sqrt	log2	sqrt	log2
Min Leaf	2	4	1	1
Max Leaf	5	10	100	2
N-Estimators	100	200	50	200

Combination 4 in Table 4.2represents the best combination, showcasing a strategic approach to constructing a decision tree ensemble. Notably, it sets the maximum depth to ‘None’ enabling trees to dynamically adjust their complexity, guarding against overfitting without arbitrary depth limitations. Additionally it establishes a minimum leaf size of 1 fostering refined decision boundaries while maintaining model interpretability. This combination strikes a balance between complexity and accuracy, making it the optimal choice for model performance and generalization.

Combination 4 strategically selects the ‘log2’ criterion for feature selection, balancing randomness and model performance. This criterion allows the algorithm to explore an optimal number of features without overwhelming computational resources or risking overfitting. Employing the ‘gini’ criterion prioritizes splits that minimize impurity, enhancing the model’s generalization capabilities. Furthermore, Combination 4 utilizes 200 estimators, leveraging the power of ensemble learning to mitigate biases and errors through averaging. This approach bolsters robustness and generalization without imposing excessive computational burdens.

## Chapter 5

# Results and Discussion

The main objective in this project is to develop a system capable of real-time hand tracking and interpretation. The ability to accurately identify and interpret hand signs in real-time has significant applications in various fields, including human-computer interaction, assistive technology, and communication aids for the differently-abled. Through the implementation of the random forest algorithm, we aimed to achieve robust and efficient hand sign recognition enabling seamless interaction between users and digital interfaces.

The project successfully achieved its objectives by creating a comprehensive dataset of sign language gestures and developing an accurate hand sign recognition system. A diverse dataset was meticulously constructed by converting video sequences into frames, initially containing 100 images for each distinct hand sign class. To enhance the dataset's diversity and robustness, data augmentation techniques were applied, resulting in a total of 5,200 images. An effective method was implemented to extract 21 key points from the hand images using MediaPipe's hand module, which employed machine learning techniques to accurately detect and extract key points representing hand landmarks within each image. These extracted key points served as effective features for representing the sign language gestures. For the machine learning model the Random Forest Classifier from Scikit-learn was selected, offering robustness and effectiveness for the classification task.

In the following sections, proposed technique present the results obtained from our implementation of the random forest algorithm for hand sign recognition, focusing on key performance metrics such as confusion matrix analysis, accuracy, precision, recall, and the F1 score. These results provide valuable insights into the efficacy and robustness of our approach, shedding light on its potential applications and areas for improvement.

## 5.1 Feature Extraction Results

The method was effectively implemented to extract 21 key points from hand images using MediaPipe's hand module which utilizes machine learning techniques to accurately identify and extract key points representing hand landmarks in each image. These extracted key points provide effective features for representing sign language gestures. The dataset consists of 26 classes, each representing a different sign language gesture, with 200 images per class, resulting in a total of 5,200 images. For each image, the MediaPipe hand module extracts 21 key points, with each key point having both x and y coordinates, resulting in 42 features per image. Consequently, for the entire dataset, the total number of features extracted is 2,18,400 (5,200 images multiplied by 42 features per image). This detailed representation of hand landmarks, comprising 21 x-coordinates and 21 y-coordinates for each image, enables precise and reliable recognition of sign language gestures.

Table 5.1: Feature Values for X-Coordinate

label	B	C	E	F	Z
x_0	0.083042	0	0.033975	0.090524	0.181193
x_1	0.716633	0.268186	0.292384	0.416548	0.347525
x_2	0.161132	0.072728	0.081073	0.137211	0.196509
x_3	0.656237	0.218068	0.26425	0.385668	0.302531
x_4	0.199316	0.108759	0.115817	0.180108	0.171798
x_5	0.316523	0.178883	0.194274	0.332283	0.242698
x_6	0.149894	0.139728	0.112814	0.209398	0.121938
x_7	0.309226	0.167105	0.138214	0.296423	0.212806
x_8	0.096974	0.177369	0.073794	0.195053	0.073016
x_9	0.372197	0.144292	0.128221	0.251326	0.210256
x_10	0.164193	0.053097	0.093877	0.138803	0.130446
x_11	0.373891	0.079631	0.100589	0.222505	0.136305
x_12	0.161548	0.092838	0.110522	0.169932	0.069641
x_13	0.224524	0.02337	0.019127	0.166133	0.065847
x_14	0.153317	0.12987	0.10385	0.182675	0.031256
x_15	0.129418	0.004915	0.056136	0.192937	0.032251
x_16	0.143041	0.164155	0.09283	0.18421	0
x_17	0.045214	0.004401	0.102174	0.230162	0
x_18	0.10854	0.041254	0.060646	0.105112	0.101954
x_19	0.358713	0.084202	0.092411	0.204753	0.167403
x_20	0.106382	0.088708	0.063453	0.112551	0.054431

Table 5.1 contains the x-coordinate values for the 21 key points (x\_0 through x\_20) extracted from hand images across the different classes labeled as B, C, E, F and Z. The values in each row is between 0 and 1 because they likely represent

normalized x-coordinate values for that particular gesture. These values are determined by considering a rectangle bounding box around the hand in the image. The origin is taken at the top left corner of the image of the hand. In the “C” category the value of x\_0 is 0 indicating that the wrist point is aligned with the origin or in other words the x-coordinate of the wrist point is 0. Similarly, in the “B” category the value of x\_4 is 0.199, indicating that the forefinger tip is 0.199 units away from the origin along the x-axis. This approach allows for a precise description of the spatial coordinates of different points on the hand relative to the origin at the top left corner of the hand image.

Table 5.2: Feature Values for Y-Coordinate

label	B	C	E	F	Z
y_0	0.198569	0.011766	0	0.109356	0.13397
y_1	0.103743	0.135474	0.063585	0.119715	0.070382
y_2	0.089096	0	0.051351	0.05079	0.191981
y_3	0.097862	0.173328	0.064195	0.12327	0.096438
y_4	0	0.0145	0.100187	0	0.223389
y_5	0.058815	0.033682	0.029941	0.072134	0.080319
y_6	0.37279	0.100081	0.095787	0.210768	0.204874
y_7	0.051067	0.080339	0.02898	0.065428	0.037865
y_8	0.221216	0.027121	0.014158	0.123324	0.182628
y_9	0.053802	0.126571	0.037004	0.064573	0.057401
y_10	0.118617	0.013436	0.059629	0.068375	0.236478
y_11	0.051761	0.161868	0.041462	0.065632	0.08452
y_12	0.036774	0.013027	0.108455	0.019307	0.262466
y_13	0.006917	0.03462	0	0.041048	0.068415
y_14	0.408131	0.125249	0.108023	0.238259	0.241683
y_15	0	0.066457	0.003915	0.019915	0.034019
y_16	0.284088	0.064537	0.047312	0.17702	0.232111
y_17	0.001722	0.096101	0.012564	0.007778	0.050411
y_18	0.199139	0.032346	0.068125	0.135017	0.263138
y_19	0	0.016885	0	0.099518	0.077218
y_20	0.125643	0.012565	0.102703	0.095423	0.27875

Similarly Table 5.2 contains the corresponding y-coordinate values for the same 21 key points (y\_0 through y\_20) across the different gesture classes. These values are also between 0 and 1. The value of y\_18 for the “Z” class representing the pinky proximal interphalangeal (PIP) joint is 0.26. This indicates that the y-coordinate of the pinky PIP joint is 0.26 units from the origin along the y-axis. Similarly, for the “F” category, the value of y\_10 representing the middle finger proximal interphalangeal (PIP) joint is 0.068 indicating that the y-coordinate of the middle finger PIP joint is 0.068 units from the origin along the y-axis. Together these two tables provide the 42 features

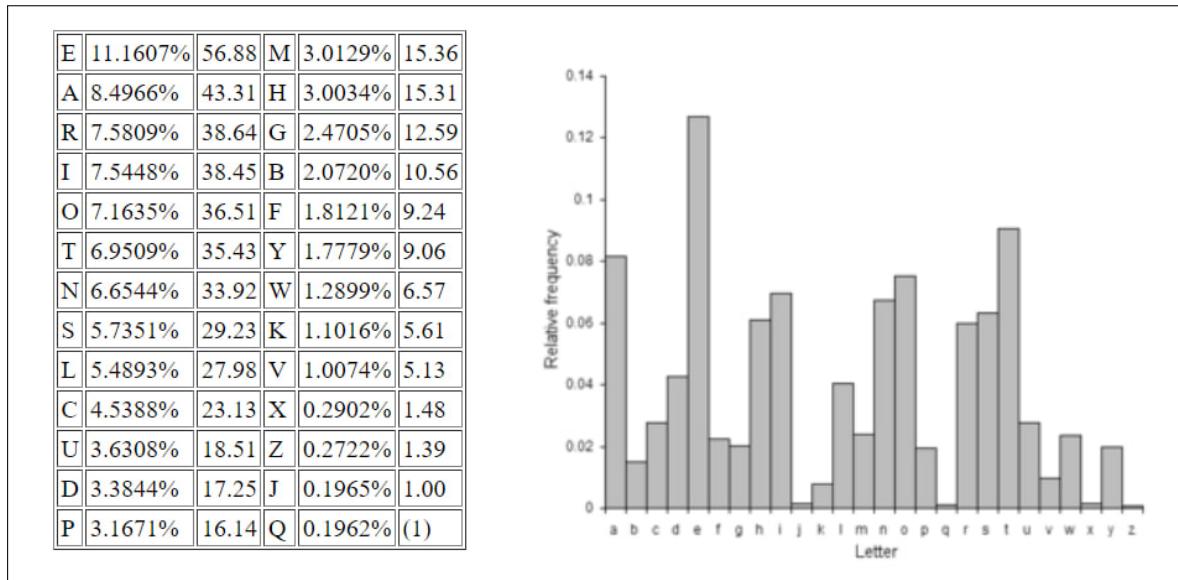


Figure 5.1: Frequency of Occurance

(21 x-coordinates and 21 y-coordinates) extracted for each image in the dataset which can be used as input to a machine learning model for sign language gesture recognition.

The Figure 5.1 shows relative frequency or percentage data for different letters of the alphabet. Each row displays a letter, a percentage value and rank or index value. The letters are ordered from highest percentage to lowest with E having the highest percentage of 11.1607% and Q having the lowest percentage of 1%.

The bar chart in Figure 5.1 visualizes the same data from that table representing with x-axis simply lists the letters in their alphabetical order and relative frequency or percentage of each letter on the y-axis[34]. From the chart we can clearly see that the letters E, A, R, I, O and T have the highest frequencies forming the tallest bars. The letters towards the end such as J, Z, X and Q have the lowest frequencies represented by the shortest bars.

## 5.2 Performance Metrics

Accuracy is the ratio of correctly predicted instances to the total instances in the dataset as shown with equation 6.1. It is a measure of the overall correctness of the model across all classes. Our model attained an accuracy of 95.67% with 200 estimators, demonstrating its high performance in hand sign recognition. Precision is the ratio of correctly predicted positive observations to the total predicted positive observations as shown with equation 6.2. It measures the proportion of correctly identified

positive cases among all cases that were predicted as positive. Our model attained an Precision of 95.83% with 200 estimators.

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}} \quad (5.1)$$

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}} \quad (5.2)$$

Recall also known as sensitivity or true positive rate, is the ratio of correctly predicted positive observations to all actual positives in the dataset as given in equation 6.3. With 200 estimators, our model achieved an impressive accuracy of 95.67%. The F1 score is the harmonic mean of precision and recall as shown with equation 6.4. It provides a balance between precision and recall and is often used as a single metric for evaluating classification models, especially when there is an uneven class distribution.proposed technique attained an score of 95.67% .

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}} \quad (5.3)$$

$$\text{F1 Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (5.4)$$

Table 5.3 presents the overall results of a proposed technique using different numbers of estimators 50, 100 and 200 in a machine learning model. As the number of estimators increases from 50 to 200, we observe an increasing trend across various evaluation metrics. The accuracy which measures the overall correctness of the model's predictions increases steadily. The precision, which quantifies the proportion of positive predictions that are truly positive also exhibits an upward trend. This implies that the model becomes more reliable in identifying positive instances correctly as the number of estimators grows. Furthermore the recall,representing the proportion of actual positive instances that are correctly identified improves with more estimators. Consequently the F1-score the harmonic mean of precision and recall also increases, indicating a balanced improvement in the model's ability to correctly classify positive and negative instances.

A confusion matrix is a table that is often used to describe the performance of a classification model on a set of test data for which the true values are known. It allows visualization of the performance of an algorithm and helps in understanding

Table 5.3: Overall Results of Proposed technique

No.of Estimators	Accuracy	Percision	Recall	F1 score	True samples	False samples	Total samples
50	0.9404	0.9421	0.9404	0.9404	3979	181	4160
100	0.9483	0.9497	0.9483	0.9483	3958	202	4160
200	0.9767	0.9583	0.9567	0.9567	3944	216	4160

how well the model is performing in terms of classifying different categories. In the presented confusion matrices, each corresponding to varying numbers of n-estimators (200, 100 and 50) the performance of the classification models is evaluated based on their ability to correctly classify samples. The x-axis and y-axis of the matrix represent the predicted labels and true labels respectively. Each row in the matrix corresponds to an actual label while each column represents a predicted label.

The diagonal elements (from the top-left to the bottom-right) of the matrix represent the correctly classified samples or true samples where the predicted label matches the true label. These values are typically highlighted or colored differently to make them visually distinct. The off-diagonal elements represent the misclassified samples or false samples. The values in these cells indicate the number of instances that were incorrectly classified as belonging to a different class.

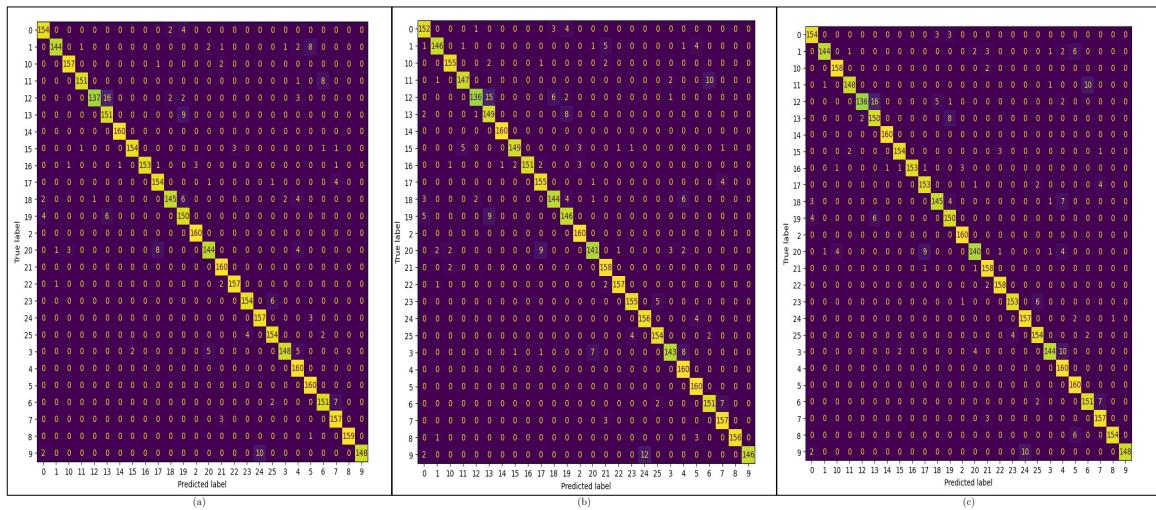


Figure 5.2: Confusion Matrix for 200, 100 and 50 n-estimators

In the confusion matrix shown in Figure 5.2(a) for 200 n-estimators for instance if we look at the row corresponding to the true label “10” (row 10) we can see that 157 samples were correctly classified as “10” which is “J” while 1 sample was misclassified as “15” which represent class “O” and 2 samples was misclassified as “21” which is “U”. Out of a total of 4160 samples 3944 were correctly classified as true samples

while 216 were misclassified as false samples. Similarly for the confusion matrix shown in Figure 5.2(b) corresponding to 100 n-estimators out of the same total sample count of 4160, 3958 were accurately identified as true samples, with 202 instances incorrectly categorized as false samples. Examining the confusion matrix in Figure 5.2(c) for 50 n-estimators out of 4160 samples, 3979 were correctly classified as true samples while 181 were labeled as false samples.

### 5.3 Comparative Analysis Across Lighting Conditions

This section includes three different test scenarios aimed at assessing the detection probability across different lighting conditions (Light, Dim, Dark). The different image quality descriptors are used for analysis of these mentioned scenarios. The used image quality descriptors are Luminance, Brightness, Contrast, Noise level.

Luminance, the first metric, is quantified in pixel intensity values within a range of 0 to 255 for 8-bit grayscale images. The luminance is calculated by converting the image to grayscale and then taking the mean value. Brightness, the second metric, shares the same unit as luminance, with a range of 0 to 255 for 8-bit color images. The brightness is calculated as the mean value of all pixels in the image. Unlike luminance, which focuses on grayscale images brightness evaluates the overall intensity across all pixels in the image, encompassing all color channels.

Contrast the third metric is measured in the standard deviation of pixel intensity values. The contrast is calculated by first computing a histogram of grayscale pixel values using, and then taking the standard deviation of the histogram. Theoretically, its range spans from 0 to approximately 73.74 for 8-bit images, as the standard deviation can extend to the value of a uniform distribution ranging from 0 to 255. Contrast reflects the extent to which pixel values deviate from the mean pixel value within the grayscale histogram. Higher contrast values indicate more pronounced deviations and greater image contrast.

Noise Level, the final metric, is denoted in pixel intensity values, with a range varying based on image content but typically falling within 0 to 255 for 8-bit images. The noise level is estimated by calculating the standard deviation of the difference between the grayscale image and a blurred version of the image using a 5x5 Gaussian kernel. This metric estimates the variability in pixel values attributed to noise. The results for luminance, brightness, contrast, and noise level are tabulated in Tables 5.4,

5.5, 5.6. These discussed assessment offers insight into the image quality, providing valuable information for image quality evaluation and processing adjustments.

## Case 1

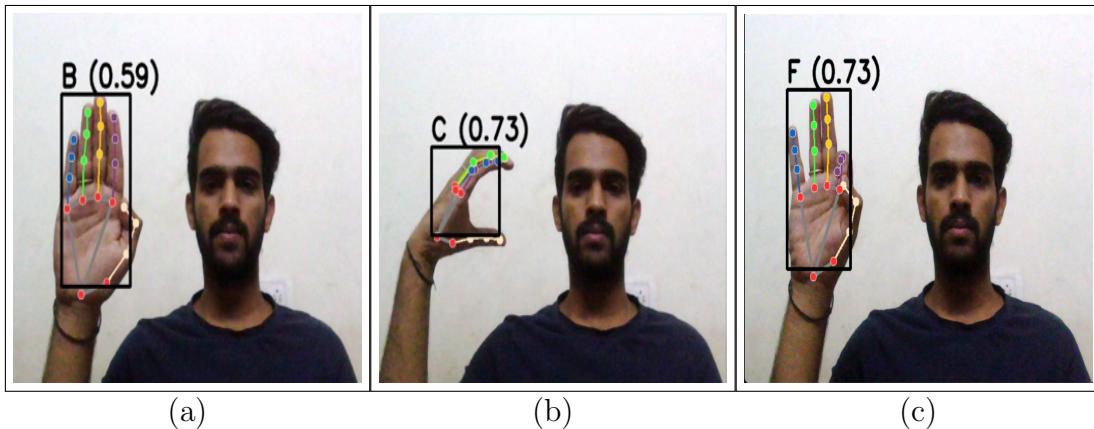


Figure 5.3: Daytime Lighting Evaluation

Table 5.4: Image Quality Parameters for Case 1 (Light)

Image	Luminance	Brightness	Contrast	Noise Level
Fig. 5.3(a)	161.42	161.54	402.21	4.46
Fig. 5.3(b)	167.54	167.69	464.55	4.01
Fig. 5.3(c)	159.29	159.38	2020.20	3.27

Referring to Figure 5.3 these images are captured in daytime conditions with sufficient luminance and brightness levels showing a person displaying three different hand signs labeled as B, C and F. The luminance values range from 159.29 to 167.54, indicating good overall brightness. The brightness values are also high, ranging from 159.38 to 167.69, confirming the well-lit nature of these images.

The contrast values show a wide variation, with one image Figure 5.3(c) having an exceptionally high contrast value of 2020.20, while the others have lower contrast values of 402.21 and 464.55. This variation in contrast could indicate differences in the lighting conditions or scene compositions within this case. The noise levels are slightly higher than Case 2, ranging from 3.27 to 4.46, which is acceptable because it can be compensated with high luminance and brightness. Overall, the Figures 5.3 in Case 1 appear to have very high luminance and brightness compared to other cases, but the wide variation in contrast values and slightly elevated noise levels suggest potential challenges in achieving consistent image quality across this case.

In the Figures 5.3 there is a probability value enclosed in parentheses which represents the model's confidence in classifying that particular hand sign. Figure 5.3(a) is daytime image the person is displaying a hand sign that the model has classified as "B" with a probability of 0.59 or 59%. Figure 5.3(b) the person is showing a hand sign that the model has classified as "C" with a probability of 0.73 or 73%.Figure 5.3(c) the person is displaying a hand sign that the model has classified as "F" with a probability of 0.73 or 73%.

## Case 2

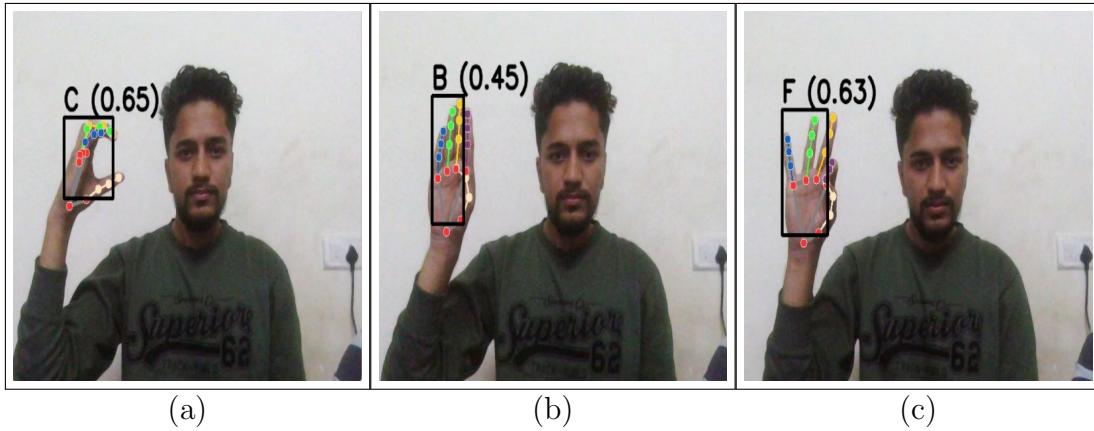


Figure 5.4: Dim Lighting Test Case

Table 5.5: Image Quality Parameters for Case 2 (Dim)

Image	Luminance	Brightness	Contrast	Noise Level
Fig. 5.4(a)	144.08	143.57	1136.97	3.57
Fig. 5.4(b)	146.05	145.70	1161.04	3.56
Fig. 5.4(c)	145.33	144.89	1201.20	3.58

Referring to Figure 5.4 these images appear to be captured in the evening or under dimmer lighting conditions. The luminance values range from 144.08 to 146.05, and the brightness values range from 143.57 to 145.70, which are slightly lower than Case 1 but still within an acceptable range. However, the contrast values are exceptionally high, ranging from 1136.97 to 1201.20. Such high contrast values suggest the presence of very distinct bright and dark regions within the images, which could potentially lead to loss of detail in the adequate brightness or darkness areas. The noise levels are relatively low, ranging from 3.56 to 3.58, indicating that these images have minimal noise or unwanted artifacts.

However an important factor to note is that this person's dataset was not included in the training data used to build the model. This means that the model has not been explicitly trained on this individual's hand gestures or appearances. Figure 5.4(a) the person is displaying a hand sign that the model has classified as "B" with a probability of 0.45 or 45%. Figure 5.4(b) the person is showing a hand sign that the model has classified as "C" with a probability of 0.65 or 65%. Figure 5.4(c) the person is displaying a hand sign that the model has classified as "F" with a probability of 0.63 or 63%.

Overall, the images in Case 2 appear to be dim-lighted with adequate luminance and brightness, but the high contrast values may warrant further investigation to ensure detail preservation in both bright and dark regions.

### Case 3

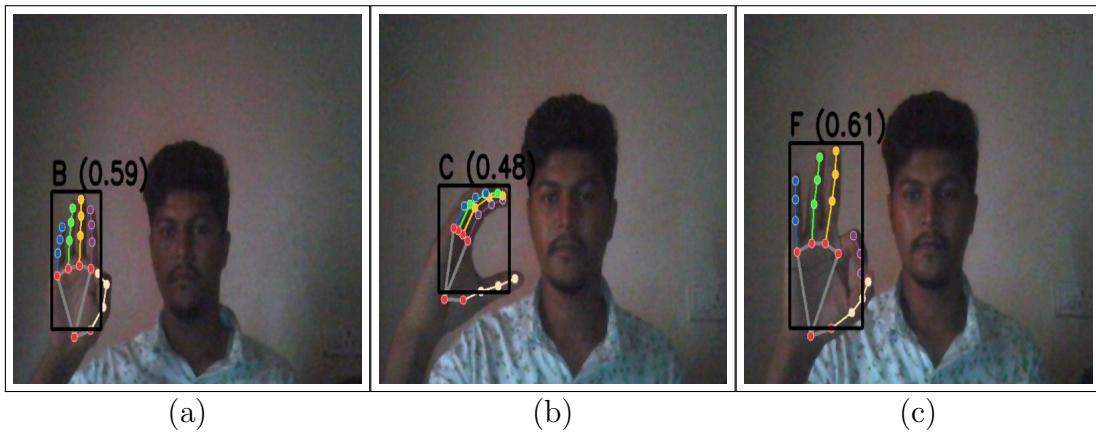


Figure 5.5: Dark Environment Evaluation

Table 5.6: Image Quality Parameters for Case 3 (Dark)

Image	Luminance	Brightness	Contrast	Noise Level
Fig. 5.5(a)	75.37	75.59	1171.11	3.00
Fig. 5.5(b)	79.34	79.40	1144.59	3.12
Fig. 5.5(c)	77.52	77.48	1175.62	3.36

The Figures 5.5 in Case 3 represent low-light or dark conditions. The luminance values range from 75.37 to 79.34, and the brightness values range from 75.59 to 79.40, which are significantly lower than the previous two cases, indicating darker overall conditions. Interestingly, the contrast values for this case are relatively high, ranging from 1144.59 to 1175.62, which could be attributed to the presence of both bright and dark regions within the low-light scenes. The noise levels are the lowest among

the three cases, ranging from 3.00 to 3.36, which might be due to the lower overall brightness levels in these images. Overall, the images in Case 3 exhibit low luminance and brightness levels, as expected for dark conditions. The high contrast values could pose challenges in preserving detail in both bright and dark areas, while the relatively low noise levels are a positive aspect.

In the test data represented by Figure 5.5 were captured in dark or low-light conditions. Similarly this person's dataset was not included in the training data used to build the model. Figure 5.5(a) the person is displaying a hand sign that the model has classified as "B" with a probability of 0.59 or 59%. Figure 5.5(b) the person is displaying a hand sign that the model has classified as "C" with a probability of 0.48 or 48%. Figure 5.5(c) is showing a hand sign that the model has classified as "F" with a probability of 0.61 or 61%. However the lower probability values for some hand signs especially in the case of the "C" sign in Figure 5.3(b) and 5.5(b) suggest that the model's performance may be slightly impacted due to the lightning condition.

Overall, the analysis revealed that daytime images (Case 1) had high luminance and brightness with variable contrast, dim lighting images (Case 2) had high contrast and low noise, while dark condition images (Case 3) exhibited low brightness but relatively high contrast and low noise. The model's classification confidence varied with lighting conditions, generally decreasing in dimmer settings.

## 5.4 Speechify Text

The primary purpose of our application is to enable individuals who are unable to speak to express themselves effectively using sign language. Through the webcam, users can perform hand signs, which are then interpreted and translated into audible speech output. This innovative approach not only facilitates communication with others but also enhances accessibility and inclusivity in various social and professional settings.

The OpenCV library is employed for capturing frames from the webcam, processing these frames, and displaying them with overlays such as text annotations. Additionally, OpenCV facilitates color space conversions and basic image manipulations. MediaPipe library is utilized for hand landmark detection. This library offers pre-trained machine learning models to accurately locate key points on the hand in real-time, enabling precise tracking of hand gestures and movements. Although not explicitly shown in the provided code snippet, a machine learning model is utilized

to classify detected hand gestures. Typically, this involves training a model on a dataset of hand gesture images annotated with their corresponding labels (e.g., letters or words). Popular algorithms for this task include Support Vector Machines (SVM), Convolutional Neural Networks (CNNs), or other classifiers. The pyttsx3 library is utilized for text-to-speech conversion. This library provides an interface to the platform-specific TTS engines installed on the system, allowing the script to generate spoken output from text strings. The script interacts with the user through keyboard commands. These commands toggle system readiness, trigger speech synthesis, add spaces, deselect words and delete characters from the constructed word.

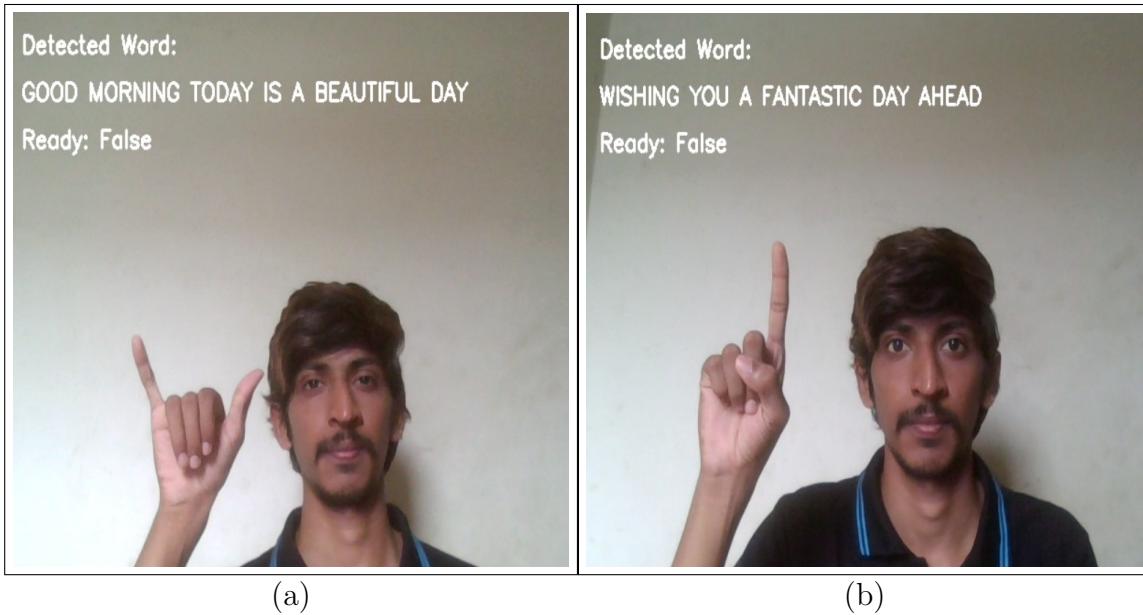


Figure 5.6: Speechify Text

The application demonstrated in Figure 5.6(a) and 5.6(b) are an innovative tool designed to bridge the communication gap for individuals who use sign language or have difficulties with verbal expression. By leveraging advanced computer vision and gesture recognition technologies, it accurately interprets hand signs and converts them into textual and auditory output. With the implementation of key strokes like backspace key likely allows the user to delete the last entered character, the delete key enables the user to remove the entire detected text, the “Ready” state displayed in the images suggests the presence of a key or command that toggles the system’s readiness to capture and interpret hand gestures. When set to “True” the application actively monitors and translates the user’s hand signs into text. While not explicitly shown, it is reasonable to assume the inclusion of a key or command that triggers the text-to-speech functionality. This would allow the detected text to be vocalized, providing an audible output for improved accessibility.

In the Figure 5.6(a) the detected text based on the hand gestures reads “good morning today is a beautiful day”. However, the “Ready” state is set to False indicating that the system is not actively capturing or interpreting hand signs at that moment. Figure 5.6(b) shows a different sentence detected as “wishing you a fantanstic day ahead again with the “Ready” state set to False. Both images feature the same person demonstrating hand signs or gestures against a plain background.

Overall, this system integrates computer vision for hand gesture detection, machine learning for gesture classification, and text-to-speech technology for spoken output. It offers a practical framework for real-time translation of hand signs into audible speech, facilitating communication for individuals with hearing impairments.

## 5.5 Comparison with State-of-Arts-Techniques

The most recent entry [1] in Table 5.7 is the proposed work in 2024, which utilized a Random Forest technique on a custom Hand Gesture Dataset, achieving an accuracy of 95.73%. In 2023, Hira Ansir et al. [2] employed Deep Belief Nets and Convolutional Neural Networks (CNNs) on a custom Hand Gesture Dataset, attaining an accuracy of 90.73%. Another work in 2023 by Dewi et al. [3] used Ye3b7 and Yebbx models on the American Sign Language (ASL) dataset, yielding a mean average precision of 86.3% and recall of 79.9%.

In 2021, Ms. Greeshma Pala et al. [4] applied Support Vector Machines (SVM), K-Nearest Neighbors (KNN), and Convolutional Neural Networks (CNN) on the American Sign Language (ASL) Dataset, achieving an accuracy of 98.49% for CNN, 83% for SVM, and 93% for KNN. Chen et al. [5] proposed YOLOv5 in 2023, achieving accuracies of 75.6% and 66.8% on EgoHands and TinyHGR datasets, respectively. Mohit Raut et al. [8] utilized K-Nearest Neighbors (KNN) and Hidden Markov Models (HMMs) on the Indian Sign Language (ISL) Dataset in 2020, obtaining an accuracy of 96.4% for HMM and 98.23% for KNN. Shin et al. [11] employed Support Vector Machines (SVM) and Gradient Boosting Machine (GBM) on ASL dataset, Massey dataset, and Finger Spelling A dataset in 2021, achieving an average accuracy of 96% for SVM and 93% for GBM.

In 2022, A.Pothuri et al. [15] and Lee et al. [19] proposed Random Forest techniques on the Indian Sign Language (ISL) dataset and American Sign Language (ASL) dataset and NUS Hand dataset, respectively, achieving accuracies of 96% for ISL and

Table 5.7: A Chronological Summary of Various Techniques in the Domain

S.No	Year	Author	Detection Technique	Dataset	Results
1	2024	<b>proposed work</b>	Random Forest	Custom Hand Gesture Dataset	Accuracy 95.73% achieved
2	2023	Hira Ansar et al. [1]	Deep Belief Network and CNN	Custom Hand Gesture Dataset	Accuracy 90.73% achieved
3	2023	Dewi et al. [3]	Yolo7 and Yolo7x	American Sign Language (ASL) dataset	Mean Average Precision 86.3% and recall 79.9%
4	2021	Ms Greeshma Pala et al. [5]	Support Vector Machines (SVM), K-Nearest Neighbors (KNN) and Convolutional Neural Network (CNN)	American Sign Language (ASL) Dataset	Accuracy of 98.49%, for CNN, 83% for SVM and 93% for KNN
5	2023	Chen et al. [6]	YOLOv5	EgoHands and TinyHGR datasets	Accuracies 75.6% and 66.8% is achieved.
6	2020	Mohit Patil et al. [8]	K-Nearest Neighbors (KNN) and Hidden Markov Models (HMMs)	Indian Sign Language (ISL) Dataset	Accuracy of 96.4% for HMM and 98.23% for KNN.
7	2021	Shin et al. [11]	Support Vector Machines (SVM) and Gradient Boosting Machine (GBM)	ASL dataset ,Massey dataset and Finger Spelling A dataset	Average accuracy of 96% for SVM and 93% for GBM.
8	2022	A.Potluri et al. [15]	Random Forest	Indian Sign Language(ISL) dataset	Accuracy of 96% is achieved.
9	2022	Lee et al. [19]	Random Forest	American Sign Language (ASL) dataset and NUS Hand dataset	Accuracy of 99% for ASL and 98% for NUS Hand dataset.
10	2015	Sangjun et al. [20]	Random Forest	Custom Hand Gesture Dataset	Accuracy of 90% for 160x120 resolution and 93% for 640x480 resolution.

99% for ASL and 98% for NUS Hand dataset. Sangum et al. [20] used Random Forest on a Custom Hand Gesture Dataset in 2015, attaining an accuracy of 90% for 160x120 resolution and 94% for 640x480 resolution.

## 5.6 Conclusion

The field of sign language recognition has seen significant advancements in recent years, with the application of image processing and machine learning techniques. Based on the results obtained from the implementation of various n-estimators in the

random forest algorithm for hand sign recognition, it is evident that increasing the number of estimators generally improves the accuracy of classification. The highest accuracy of 95.67% was achieved with 200 n-estimators, showcasing the robustness of the model in distinguishing between different hand signs. However, it's important to note that while higher n-estimators may lead to improved accuracy, there might be a trade-off with computational resources and efficiency. Furthermore, the analysis of confusion matrices revealed consistent trends across different numbers of n-estimators, with the majority of samples correctly classified as true positives. This indicates the effectiveness of the random forest algorithm in accurately recognizing hand signs.

In conclusion, the project successfully demonstrated the feasibility and effectiveness of using the random forest algorithm for hand sign recognition, laying the foundation for further research and development in this field. The insights gained from this project contribute to the advancement of gesture recognition technology, with potential applications in human-computer interaction, assistive technology, and accessibility solutions.

## 5.7 Future Scope

Integration of hand sign recognition capabilities into wearable devices such as smart glasses or wristbands could enable hands-free interaction and communication for individuals with disabilities. Additionally incorporating hand sign recognition into Internet of Things (IoT) devices could enable gesture-based control of smart home appliances, entertainment systems and other connected devices. While there is great potential for hand sign recognition systems in assistive technology applications, including communication aids for individuals with disabilities the accessibility and affordability of such technologies may limit their widespread adoption.

Machine learning-powered educational applications for sign language learning and training hold promise creating immersive virtual reality (VR) or augmented reality (AR) environments with realistic hand sign recognition capabilities may still be challenging. Developing adaptive learning platforms that cater to individual learning styles and abilities also requires sophisticated machine learning algorithms and user interface design. Although real-time hand sign recognition systems have potential applications in healthcare and rehabilitation integrating these systems into clinical practice may face regulatory hurdles and require rigorous validation and testing. Additionally ensuring the accuracy, reliability and safety of machine learning algorithms in critical healthcare settings remains a significant challenge.

While these future scopes represent exciting opportunities for advancing real-time hand sign recognition systems using machine learning, addressing the associated challenges will require collaborative efforts across multiple disciplines, including computer science, engineering, healthcare, and social sciences. Continued research innovation and investment in these areas are essential to realizing the full potential of hand sign recognition technology in improving accessibility, communication and quality of life for individuals worldwide.

# References

- [1] Alonazi, Mohammed, Hira Ansar, Naif Al Mudawi, Saud S. Alotaibi, Nouf Abdullah Almujally, Abdulwahab Alazeb, Ahmad Jalal, Jaekwang Kim and Moohong Min. "Smart Healthcare Hand Gesture Recognition Using CNN-Based Detector and Deep Belief Network." *IEEE Access* 11 (2023): 84922-84933.
- [2] Miah, Abu Saleh Musa, Md. Al Mehedi Hasan and Jungpil Shin. "Dynamic Hand Gesture Recognition Using Multi-Branch Attention Based Graph and General Deep Learning Model." *IEEE Access* 11 (2023): 4703- 4716.
- [3] Dewi, Christine, Abbott Po Shun Chen and Henoch Juli Christanto. "Deep Learning for Highly Accurate Hand Recognition Based on Yolov7 Model." *Big Data Cogn. Comput.* 7 (2023): 53..
- [4] Athira Devaraj and Aswathy K Nair. "Hand Gesture Signal Classification using Machine Learning." International Conference on Communication and Signal Processing, July 28 - 30, 2020, India:978-1-7281-4988-2 ©2020 IEEE.
- [5] Ms. Greeshma Pala,Ms. Jagruti Bhagwan Jethwani,Mr. Satish Shivaji Kumbhar,Ms. Shruti Dilip Patil. "Machine Learning-based Hand Sign Recognition." Proceedings of the International Conference on Artificial Intelligence and Smart Systems (ICAIS-2021) IEEE Xplore Part Number: CFP21OAB-ART; ISBN: 978-1-7281-9537-7.
- [6] Chen, Renxiang and Xia Tian. "Gesture Detection and Recognition Based on Object Detection in Complex Background." *Applied Sciences* (2023): n. pag.
- [7] Kaushik Kumar Jha, Balram Kumar Jha, Suyash Ajit Byale, Vivek Sitaram Dadewa, Prof. Pushpan- jali, S.Sajjanshetti."Sign Language Detection Using CNN." *International Journal for Research in Applied Science Engineering Technology (IJRASET)*:Volume 10 Issue XI Nov 2022.
- [8] Mohit Patil, Pranay Pathole, Hrishikesh Patil, Ashutosh Raut, Prof. S S Jadhav. "Indian Sign Language Recognition." *International Journal of Scientific Research Engineering Trends* Volume 6, Issue 4, July-Aug-2020, ISSN (Online): 2395-566X.

- [9] Pranjali Manmode, Rupali Saha, Manisha N. Amnerkar."Real-Time Hand Gesture Recognition." International Journal of Scientific Research in Computer Science, Engineering and Information Technology (IJSRC- SEIT)Volume 7, Issue 3 Page Number: 618-624:May-June-2021
- [10] Ben Atitallah, Bilel, Zheng Hu, Dhouha Bouchaala, Mohammed Hussain, Amir Ismail, Nabil Derbel and Olfa Kanoun. "Hand Sign Recognition System Based on EIT Imaging and Robust CNN Classification." IEEE Sensors Journal 22 (2022): 1729-1737.
- [11] Shin, Jungpil, Akitaka Matsuoka, Md. Al Mehedi Hasan and Azmain Yakin Srizon. "American Sign Language Alphabet Recognition by Extracting Feature from Hand Pose Estimation." Sensors (Basel, Switzerland) 21 (2021): n. pag.
- [12] Suharjito, F. Wiryan, G. P. Kusuma and A. Zahra, "Feature Extraction Methods in Sign Language Recognition System: A Literature Review," 2018 Indonesian Association for Pattern Recognition International Conference (INAPR), Jakarta, Indonesia, 2018, pp. 11-15, doi: 10.1109/INAPR.2018.8626857.
- [13] Cruz Bautista, Antonio Guadalupe, José-Joel González-Barbosa, Juan Bautista Hurtado-Ramos, Francisco-Javier Ornelas-Rodriguez and Erick-Alejandro González-Barbosa. "Hand features extractor using hand contour – a case study." Automatika 61 (2019): 108 - 99.
- [14] V. Radhika, C. R. Prasad and A. Chakradhar, "Smartphone-Based Human Activities Recognition System using Random Forest Algorithm," 2022 International Conference for Advancement in Technology (ICONAT), Goa, India, 2022, pp. 1-4, doi: 10.1109/ICONAT53423.2022.9726006.
- [15] A. S, A. Potluri, S. M. George, G. R and A. S, "Indian Sign Language Recognition Using Random Forest Classifier," 2021 IEEE International Conference on Electronics, Computing and Communication Technologies (CONECCT), Bangalore, India, 2021, pp. 1-6, doi: 10.1109/CONECCT52877.2021.9622672.
- [16] Ellahyani, Ayoub, Mohamed El Ansari and Ilyas El Jaafari. "Traffic sign detection and recognition based on random forests." Appl. Soft Comput. 46 (2016): 805-815.
- [17] Zafrulla, Zahoor et al. "Hand detection in American Sign Language depth data using domain-driven random forest regression." 2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG) 1 (2015): 1-7.

- [18] Bhushan, S.; Alshehri, M.; Keshta, I.; Chakraverti, A.K.; Rajpurohit, J.; Abugabah, A. "An Experimental Analysis of Various Machine Learning Algorithms for Hand Gesture Recognition." *Electronics* 2022, 11, 968. <https://doi.org/10.3390/electronics11060968>.
- [19] Ewe, E.L.R.; Lee, C.P.; Kwek, L.C.; Lim, K.M. "Hand Gesture Recognition via Lightweight VGG16 and Ensemble Classifier." *Appl. Sci.* 2022, 12, 7643. <https://doi.org/10.3390/app12157643>.
- [20] Sangjun, O. Mallipeddi, Rammohan Lee, Minho. (2015). " Real Time Hand Gesture Recognition Using Random Forest and Linear Discriminant Analysis." [10.1145/2814940.2814997](https://doi.org/10.1145/2814940.2814997).
- [21] Python Coding Standards and Best Practices for Code Quality. Available : <https://www.zenesys.com/python-coding-standards-best-practices>.
- [22] Scikit-Learn Essentials for Machine Learning in Python. Available: <https://www.datacamp.com/tutorial/machine-learning-python>.
- [23] Practical Machine Learning with Scikit-Learn. Available : <https://www.simplilearn.com/tutorials/python-tutorial/scikit-learn>.
- [24] Getting Started with OpenCV. Available : [https://docs.opencv.org/4.x/db/d05/tutorial\\_config\\_reference.html](https://docs.opencv.org/4.x/db/d05/tutorial_config_reference.html).
- [25] OpenCV Demystified: An Introductory Guide. Available : <https://github.com/opencv/opencv/wiki/0E-34.-Named-Parameters>
- [26] Pandas Ecosystem and Integration. Available : <https://pandas.pydata.org/docs/>
- [27] MediaPipe:<https://mediapipe.readthedocs.io/en/latest/solutions/hands.html>
- [28] Exploring VGGNet: Unraveling the Depths of Very Deep Convolutional Networks:<https://viso.ai/deep-learning/vgg-very-deep-convolutional-networks/>.
- [29] Rapid Response Object Detection with SSD Network. Available : <https://towardsdatascience.com/review-ssd-single-shot-detector-object-detection-851a94607d11>.
- [30] Random Forest <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>
-

- [31] Understanding Random Forest with Examples:<https://www.analyticsvidhya.com/blog/2021/06/understanding-random-forest/>
- [32] StratifiedKFold:[https://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.StratifiedKFold.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.StratifiedKFold.html)
- [33] Grid Search Cv:[https://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.GridSearchCV.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html)
- [34] Frequency Of Alphabets <https://www3.nd.edu/~busiforc/handouts/cryptography/letterfrequencies.html>

# Appendix A

## Code

### Capture Image

```
# Capture Images
# pip install mediapipe==0.9.0.1 scikit-learn==1.2.0
# opencv-python==4.7.0.68 pandas==2.0.3
!pip install -r requirements.txt

import os
import cv2
DATA_DIR = './data'
number_of_classes = 26 # Assuming you have 26 classes
dataset_size = 100
if not os.path.exists(DATA_DIR):
    os.makedirs(DATA_DIR)

cap = cv2.VideoCapture(0)

try:
    if not cap.isOpened():
        print("Error: Could not open camera.")
        exit()

    for j in range(number_of_classes):
        if not os.path.exists(os.path.join(DATA_DIR,
                str(j))):
            os.makedirs(os.path.join(DATA_DIR, str(j)))

        print('Collecting data for class {}'.format(j))

        while True:
            ret, frame = cap.read()
```

```

        cv2.putText(frame, 'Ready? Press "S" to start
                     capturing, "N" to skip, or "q" to exit!',
                     (50, 50), cv2.FONT_HERSHEY_SIMPLEX, 1.0,
                     (0, 255, 0), 3, cv2.LINE_AA)
        cv2.imshow('frame', frame)
        key = cv2.waitKey(25)

        if key == ord('s'):
            print("Capturing data for class
                  {}".format(j))
            break
        elif key == ord('n'):
            print("Skipping class {}".format(j))
            break
        elif key == ord('q'):
            print("Exiting the program.")
            cap.release()
            cv2.destroyAllWindows()

        if key == ord('n'):
            continue # Skip to the next class

        counter = 0
        while counter < dataset_size:
            ret, frame = cap.read()
            cv2.imshow('frame', frame)
            cv2.waitKey(25)
            cv2.imwrite(os.path.join(DATA_DIR, str(j),
                                     '{}.jpg'.format(counter)), frame)
            counter += 1

    except Exception as e:
        print("An error occurred:", e)
    finally:
        cap.release()
        cv2.destroyAllWindows()

```

Listing A.1: Code for capturing hand images from webcam

## Feature Extraction

```
import os
import pickle
import mediapipe as mp
import cv2
import matplotlib.pyplot as plt

mp_hands = mp.solutions.hands
mp_drawing = mp.solutions.drawing_utils
mp_drawing_styles = mp.solutions.drawing_styles

hands = mp_hands.Hands(static_image_mode=True,
    min_detection_confidence=0.3)

DATA_DIR = './data'

data = []
labels = []

for dir_ in os.listdir(DATA_DIR):
    for img_path in os.listdir(os.path.join(DATA_DIR,
        dir_)):
        data_aux = []

        x_ = []
        y_ = []

        img = cv2.imread(os.path.join(DATA_DIR, dir_,
            img_path))
        img_rgb = cv2.cvtColor(img, cv2.COLOR_BGR2RGB)

        results = hands.process(img_rgb)
        if results.multi_hand_landmarks:
            for hand_landmarks in
                results.multi_hand_landmarks:
                for i in
                    range(len(hand_landmarks.landmark)):
                    x = hand_landmarks.landmark[i].x
```

```

y = hand_landmarks.landmark[i].y

x_.append(x)
y_.append(y)

for i in
range(len(hand_landmarks.landmark)):
    x = hand_landmarks.landmark[i].x
    y = hand_landmarks.landmark[i].y
    data_aux.append(x - min(x_))
    data_aux.append(y - min(y_))

data.append(data_aux)
labels.append(dir_)

f = open('data.pickle', 'wb')
pickle.dump({'data': data, 'labels': labels}, f)
f.close()

```

Listing A.2: Code for extracting features

## Training Classifier and Model Summary

```

import os
import pickle
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score,
    precision_score, recall_score, f1_score,
    classification_report
import numpy as np
from sklearn.metrics import confusion_matrix,
    ConfusionMatrixDisplay

def train_model_and_evaluate(model, data, labels,
    test_size=0.8, random_state=42):

```

```
# Split the data into training and testing sets
x_train, x_test, y_train, y_test =
    train_test_split(data, labels, test_size=test_size,
                      random_state=random_state, shuffle=True,
                      stratify=labels)

# Train the model
model.fit(x_train, y_train)

# Make predictions on the testing data
y_predict = model.predict(x_test)

# Calculate evaluation metrics
accuracy = accuracy_score(y_test, y_predict)
precision = precision_score(y_test, y_predict,
                             average='weighted')
recall = recall_score(y_test, y_predict,
                       average='weighted')
f1 = f1_score(y_test, y_predict, average='weighted')

# Generate classification report
report = classification_report(y_test, y_predict)

cm = confusion_matrix(y_test,y_predict)
cmd = ConfusionMatrixDisplay(confusion_matrix=cm,
                             display_labels=model.classes_)
cmd.plot()

return accuracy, precision, recall, f1, report

data_dict = pickle.load(open('./data.pickle', 'rb'))
data = np.asarray(data_dict['data'])
labels = np.asarray(data_dict['labels'])

# Define models to train
models = [
```

```

("Random_Forest_Classifier",
    RandomForestClassifier(bootstrap=True,
    n_estimators=200,
    random_state=42, min_samples_split=2,
    min_samples_leaf=1, max_features='sqrt',
    max_depth=10)])

# Train and evaluate each model
for model_name, model in models:
    print(f"Training and evaluating {model_name}...")
    accuracy, precision, recall, f1, report =
        train_model_and_evaluate(model, data, labels)

    print("Evaluation Metrics:")
    print(f"Accuracy: {accuracy:.4f}")
    print(f"Precision: {precision:.4f}")
    print(f"Recall: {recall:.4f}")
    print(f"F1-score: {f1:.4f}")
    print("\nClassification Report:")
    print(report)
    print("-----")

# Save the trained model to a file
os.makedirs('./all_models/', exist_ok=True)
with open(f'./all_models/{model_name}_model.p', 'wb') as f:
    pickle.dump({'model': model}, f)

```

Listing A.3: Code for Training the Classifier and Model Summary

## Testing

```

import pickle, cv2, mediapipe as mp, numpy as np

model = model_dict['model']
model_dict =
    pickle.load(open('./all_models/RF_Classifier_model.p',
    'rb'))

```

```
mp_hands = mp.solutions.hands
mp_drawing = mp.solutions.drawing_utils
mp_drawing_styles = mp.solutions.drawing_styles

hands = mp_hands.Hands(static_image_mode=True,
    min_detection_confidence=0.3)
# labels from 0 to 12
labels_dict = {0:'A', 1:'B', 2:'C', 3:'D', 4:'E', 5:'F',
    6:'G', 7:'H', 8:'I', 9:'J', 10:'K', 11:'L', 12:'M',
    13:'N', 14:'O', 15:'P', 16:'Q', 17:'R', 18:'S', 19:'T',
    20:'U', 21:'V', 22:'W', 23:'X', 24:'Y', 25:'Z'}

import cv2

# Initialize the webcam
cap = cv2.VideoCapture(0)

while True:
    try:
        # Data and variables initialization
        data_aux = []
        x_ = []
        y_ = []

        # Read a frame from the webcam
        ret, frame = cap.read()

        # Get the dimensions of the frame
        H, W, _ = frame.shape

        # Convert the frame to RGB color space
        frame_rgb = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB)

        # Process the frame to detect hand landmarks
        results = hands.process(frame_rgb)
```

```
# If hand landmarks are detected
if results.multi_hand_landmarks:
    for hand_landmarks in
        results.multi_hand_landmarks:
            mp_drawing.draw_landmarks(
                frame, # image to draw
                hand_landmarks, # model output
                mp_hands.HAND_CONNECTIONS, # hand
                connections
                mp_drawing_styles.get_default_hand_
                landmarks_style(),
                mp_drawing_styles.get_default_hand_
                connections_style())

    for i in
        range(len(hand_landmarks.landmark)):
        x = hand_landmarks.landmark[i].x
        y = hand_landmarks.landmark[i].y

        x_.append(x)
        y_.append(y)

    for i in
        range(len(hand_landmarks.landmark)):
        x = hand_landmarks.landmark[i].x
        y = hand_landmarks.landmark[i].y
        data_aux.append(x - min(x_))
        data_aux.append(y - min(y_))

    x1 = int(min(x_) * W) - 10
    y1 = int(min(y_) * H) - 10

    x2 = int(max(x_) * W) - 10
    y2 = int(max(y_) * H) - 10

    # Make predictions using the model
    prediction =
        model.predict([np.asarray(data_aux)])
```

```
prediction_prob =
    model.predict_proba([np.asarray(data_aux)])  
  
# If predictions are available
if prediction.shape[0] > 0:
    predicted_label_index = int(prediction[0])
    predicted_character =
        labels_dict.get(predicted_label_index,
                        "Unknown")
    prediction_probability =
        np.max(prediction_prob) # Get the
        maximum probability  
  
else:
    predicted_character = "Unknown"
    prediction_probability = 0.0 # Set
        probability to zero for unknown
        prediction  
  
# Draw bounding box and label on the frame
cv2.rectangle(frame, (x1, y1), (x2, y2), (0,
    0, 0), 4)
if prediction_probability > 0.4:
    cv2.putText(frame,
        f'{predicted_character}'
        f'({prediction_probability:.2f})', (x1,
        y1 - 10), cv2.FONT_HERSHEY_SIMPLEX,
        1.3, (0, 0, 0), 3, cv2.LINE_AA)
else:
    cv2.putText(frame, f'{ "Unknown"}'
        f'({prediction_probability:.2f})', (x1,
        y1 - 10), cv2.FONT_HERSHEY_SIMPLEX,
        1.3, (0, 0, 0), 3, cv2.LINE_AA)  
  
# Display the frame
cv2.imshow('frame', frame)
```

```
# Check for user input to exit the loop
c = cv2.waitKey(1)
if c == ord('q'):
    break

except Exception as e:
    # Print the exception
    print("An error occurred:", str(e))

# Release the webcam and close all windows
cap.release()
cv2.destroyAllWindows()
```

Listing A.4: Code for Testing

# Appendix B

## Conference Paper

IETE Sponsored 7<sup>th</sup> National Conference on Emerging Trends in Engineering, Science and Technology (NCTEST-7)

### Real Time Sign Language Recognition System Using Machine Learning

1<sup>st</sup> Raghavendra M Hegde

*Electronics and Communication  
R N Shetty Institute of Technology*

Bengaluru, India

1rn20ec066.raghavendramanjunathegde@rnsit.ac.in | 1rn20ec078.samarthshinnur@rnsit.ac.in | 1rn20ec079.samrudhbr@rnsit.ac.in

2<sup>nd</sup> Samarth S

*Electronics and Communication  
R N Shetty Institute of Technology*

Bengaluru, India

1rn20ec078.samarthshinnur@rnsit.ac.in | 1rn20ec079.samrudhbr@rnsit.ac.in

3<sup>rd</sup> Samrudh BR

*Electronics and Communication  
R N Shetty Institute of Technology*

Bengaluru, India

4<sup>th</sup> Sanjay SB

*Electronics and Communication  
R N Shetty Institute of Technology*

Bengaluru, India

1rn20ec081.sanjaysb@rnsit.ac.in

5<sup>th</sup> Dr. Leena chandrashekhar

*Electronics and Communication  
R N Shetty Institute of Technology*

Bengaluru, India

leenac@rnsit.ac.in

**Abstract**—Abstract—Sign language serves as a primary means of communication for individuals who are vocally disabled, utilizing a combination of finger signs, facial expressions, and gestures to convey information. However, the interpretation of sign language poses challenges for many, leading to potential miscommunication. To address this issue, a system leveraging advanced technologies such as deep learning, machine learning, Random Forest classifier, computer vision, Mediapipe, and Python is developed. This technology aims to accurately detect and identify sign language motions in real-time. The system introduces a robust hand sign tracking and recognition model capable of accurately tracking and recognizing RGB dynamic gestures, even in complex environments. The process involves converting videos into frames, adjusting image light intensity, removing noise, extracting features using CNN for key-point extraction, and employing the Random Forest machine learning algorithm for classification. Custom datasets are utilized to validate the proposed system, with experimental results demonstrating overall accuracies of 94.58 percent and 95.67 percent respectively. These experiments affirm the readability and effectiveness of the proposed model compared to other state-of-the-art techniques.

**Index Terms**—sign language, dynamic gesture, machine learning.

#### I. INTRODUCTION

In contemporary society, the interaction between humans and computers has become indispensable due to continuous technological advancements. This interaction, referred to as Human-Computer Interaction (HCI)[2], encompasses various methods facilitating communication with and control of computers. From utilizing a mouse to tapping on a touchscreen, HCI plays a vital role in navigating the digital realm.

Hand gesture recognition has emerged as a focal point of interest among researchers within the broader scope of HCI. It provides a natural and intuitive means of interacting with computers, free from the complexities associated with traditional input devices. Unlike devices such as joysticks or remote controls, hand gestures feel instinctive and user-friendly. Their applications span diverse domains, including

home appliance control, robotic guidance, and medical assistance. However, despite its potential benefits, hand gesture recognition[5] presents challenges, particularly in ensuring accurate recognition under varying environmental conditions. In fields like healthcare, where precision is critical, errors in gesture recognition can have severe consequences. Researchers are actively working on developing robust systems capable of reliably interpreting gestures[7] in diverse contexts, although achieving this goal remains challenging.

The process of hand gesture recognition involves several stages, including gesture capture, hand position identification, feature extraction, and gesture classification. Leveraging technologies like cameras and sensors, these stages enable the detection and analysis of both simple and complex gestures. In our research project, we aim to advance the field of hand gesture recognition by proposing innovative solutions to address existing challenges. Our primary focus is on enhancing the accuracy and reliability of gesture recognition systems[1], particularly in critical domains such as healthcare and home automation. Additionally, our efforts aim to streamline communication between humans and computers[3], facilitating smoother interactions. Through our project, we seek to not only enhance hand gesture recognition technology but also contribute to a future where human-computer interaction is more intuitive and efficient. By bridging the gap between humans and machines, we strive to create a more interconnected and accessible digital world.

#### II. METHODOLOGY

##### A. Data Set Collection

The dataset collection process for hand sign recognition is systematically designed, starting with the initialization of a structured directory to store images. Users interact with the system through keystrokes, initiating data collection for each hand sign class and navigating between classes. Real-time webcam feed captures images, with users prompted to collect

IETE Sponsored 7<sup>th</sup> National Conference on Emerging Trends in Engineering, Science and Technology (NCTEST-7)

a predetermined number per class, typically set to 100 images. Images are promptly stored in JPG format within designated directories based on hand sign classes, ensuring organized data storage. Robust error-handling mechanisms are integrated to manage unforeseen exceptions, enhancing system reliability for acquiring a diverse dataset essential for training robust hand sign recognition models.

#### B. Feature Extraction

The feature extraction stage in hand sign recognition initiates with data augmentation, which enhances dataset diversity and reduces overfitting by applying transformations like rotation and scaling. Subsequently, the Mediapipe hand tracking model is employed to extract landmarks, utilizing deep learning algorithms for detection and localization. Pre-processing steps standardize input formats, followed by hand detection using SSD architecture to identify potential hand regions across various image scales. Landmarks, such as fingertips and palm center, are then localized using CNN layers and regression techniques, with coordinates normalized for consistency across different hand sizes and orientations. This enriched dataset, resulting from augmentation and landmark extraction, serves as the input for training a robust hand sign recognition model, essential for accurate gesture interpretation and classification in real-world scenarios

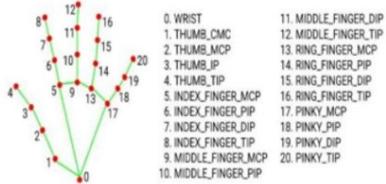


Fig. 1. Media Pipe Hand Landmark

#### C. Training Classifier

In the training phase, the machine learning pipeline begins with data splitting, where the dataset is divided into separate training and testing sets using the train-test-split function from Scikit-learn. Approximately 80 percent of the data is allocated for training, while the remaining 20 percent is reserved for testing the model's performance. A Random Forest Classifier model is then trained, chosen for its versatility and effectiveness in classification tasks, as it employs an ensemble of decision trees. These decision trees are generated through a bootstrap process, which involves randomly sampling the training data to create multiple subsets. Each decision tree recursively partitions the feature space based on feature thresholds, contributing to the ensemble's predictive power. By aggregating predictions from multiple trees through majority voting, the model's ability to generalize is enhanced, reducing the risk of overfitting and ensuring reliability and stability.

#### D. Testing

The hand sign recognition system evaluates its accuracy by processing live input from the webcam using the Mediapipe library's hand module to detect landmarks. Once landmarks are detected, features are extracted and normalized, serving as input for the trained model to predict hand signs for each frame. Predictions, along with corresponding probabilities, are seamlessly integrated into the system's output for real-time visualization, enabling users to observe predicted hand signs instantly and assess system performance effectively.

### III. FLOW DIAGRAM

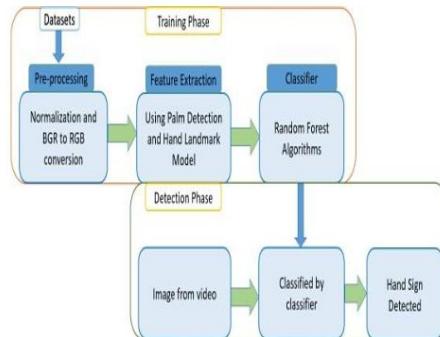


Fig. 2. Flow Chart of Hand Sign Recognition System

### IV. MODEL SUMMARY

The Random Forest model, a robust ensemble learning technique, comprises multiple decision trees trained on bootstrap samples of the dataset, injecting randomness and diversity into the ensemble. During training, decision trees split nodes based on impurity measures like Gini or entropy, with the ensemble aggregating predictions through majority voting or averaging for classification and regression tasks, respectively, thereby enhancing generalization and mitigating overfitting. Key parameters such as the number of estimators, maximum depth, and maximum features govern the model's behavior and performance.

Bootstrap sampling involves randomly selecting samples with replacement from the dataset, fostering diversity in the trees when activated. Class weights allow for balancing the influence of imbalanced class distributions during training. Decision trees split nodes based on impurity measures, commonly using 'Gini' in Random Forest models. The maximum depth parameter controls tree depth, while maximum features determine the number of features considered for the best split. Minimum samples per leaf sets the minimum samples required at a leaf node, and the number of estimators defines the

IETE Sponsored 7<sup>th</sup> National Conference on Emerging Trends in Engineering, Science and Technology (NCTEST-7) ensemble size, affecting both performance and computational cost.

Random Forest employs bagging to create multiple decision trees trained on bootstrap samples, diminishing overfitting and bolstering model stability and accuracy. Feature selection randomly chooses feature subsets for each tree, reducing correlation and overfitting. Decision trees are built recursively on bootstrap samples, with parameters like max depth and min samples split regulating tree growth. Gini impurity guides the algorithm in selecting optimal splits for homogeneous groups. During prediction, trees independently forecast, and the mode or average of predictions is adopted for classification or regression tasks, respectively. Out-of-bag evaluation estimates model performance using samples excluded from training, providing intrinsic validation for Random Forest models.

$$I_G = 1 - \sum_{i=1}^K p_i^2$$

Fig. 3. Gini Formula

The evaluation pipeline relies on StratifiedKFold for unbiased cross-validation, preserving class distributions and enhancing model assessment across diverse data subsets. GridSearchCV facilitates exhaustive hyperparameter tuning, iterating over parameter combinations to identify optimal configurations for the Random Forest Classifier. By splitting the dataset and evaluating performance for each parameter combination, GridSearchCV ensures robust model selection based on specified scoring metrics. The comparison of different parameter combinations reveals variations in criteria like maximum depth and minimum samples leaf, impacting model complexity and generalization. Despite similarities in certain hyperparameters, variations across combinations highlight the importance of selecting the optimal configuration for maximizing predictive performance.

## V. RESULT

- Our objective was to establish a reliable and effective hand sign recognition system through the utilization of the random forest algorithm, enhancing the interaction between users and digital interfaces seamlessly.
  - We aimed to facilitate the advancement of more precise and inclusive recognition systems by meticulously assembling and annotating a varied collection of hand sign images, addressing the requirements of diverse user demographics.
  - Our focus on evaluating the performance of the random forest algorithm for hand sign recognition encompassed thorough analysis of key metrics such as confusion matrix analysis, accuracy, precision, recall, and the F1 score. These findings offer valuable insights into the efficiency and resilience of our methodology, highlighting potential applications and areas for enhancement.



Fig. 4. Result

#### A. Performance Metrics

The accuracy metric represents the proportion of correctly predicted instances to the total instances within the dataset, serving as a comprehensive indicator of the model's correctness across all classes. Our model achieved an impressive accuracy rate of 95.67 percent, utilizing 200 estimators, showcasing its exceptional performance in hand sign recognition.

Additionally, a confusion matrix serves as a valuable tool for evaluating the performance of a classification model on a test dataset with known true values. It provides a visual representation of the model's performance, offering insights into its ability to classify different categories accurately.

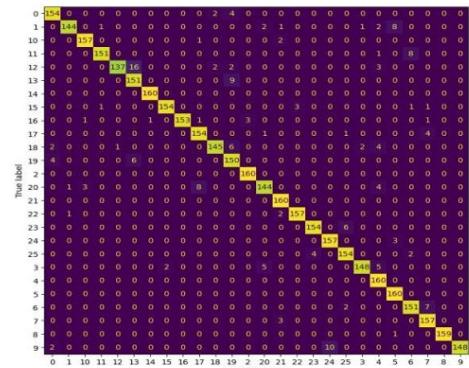


Fig. 5. Confusion Matrix for 200 n-estimators

In the confusion matrix referred in figure.5 for 200 n-estimators, out of a total of 4160 samples, 3944 were correctly classified as true samples, while 216 were misclassified as false samples. For the confusion matrix shown in figure.6 corresponding to 100 n-estimators, out of the same total sample count of 4160, 3958 were accurately identified as true samples, with 202 instances incorrectly categorized as false samples. Examining the

IETE Sponsored 7<sup>th</sup> National Conference on Emerging Trends in Engineering, Science and Technology (NCTEST-7)

confusion matrix in figure.7 for 50 n-estimators, out of 4160 samples, 3979 were correctly classified as true samples, while 181 were labeled as false samples.

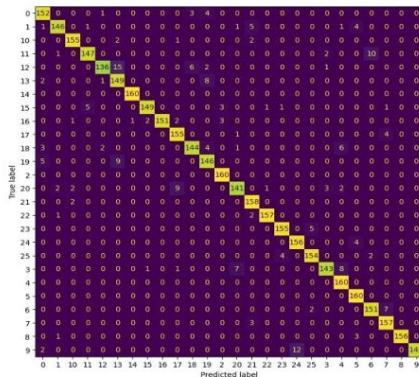


Fig. 6. Confusion Matrix for 100 n-estimators

Our model achieved impressive performance metrics, including a precision of 95.83 percent and a recall of 95.67 percent with 200 estimators, highlighting its ability to correctly identify positive cases and capture all actual positives in the dataset. Additionally, the F1 score, which balances precision and recall, was notably high at 95.67 percent. These results demonstrate the model's effectiveness in accurately classifying hand signs while considering both the proportion of correctly identified positive cases and the ability to capture all actual positives. Overall, the model's robust performance across multiple evaluation metrics underscores its reliability and suitability for hand sign recognition tasks.

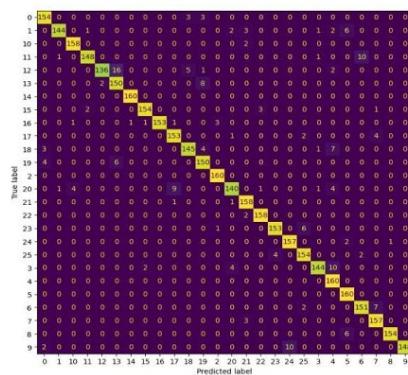


Fig. 7. Confusion Matrix for 50 n-estimators

No. of Estimators	Accuracy	Precision	Recall	F1-score	True Sample	False Sample
50	0.9404	0.9421	0.9404	0.9404	3979	181
100	0.9483	0.9497	0.9483	0.9483	3958	202
200	0.9567	0.9583	0.9567	0.9567	3944	216

Fig. 8. Results For Different Trees

Hence the performance metrics, such as accuracy, precision, and recall, achieved by various decision tree models referred in figure.8 . The comparison highlights the effectiveness of different tree-based algorithms in accurately classifying the data.

## VI. CONCLUSION AND FUTURE SCOPE

Recent advancements in sign language recognition, utilizing image processing and machine learning methodologies, have demonstrated considerable potential. Experimentation with different numbers of estimators in the random forest algorithm has revealed a trend wherein increasing the number of estimators generally improves classification accuracy. The highest accuracy of 95.67 percent was observed with 200 estimators, underscoring the model's capability to effectively distinguish between various hand signs. Nevertheless, it's imperative to consider the balance between higher estimator numbers and computational resources. Examination of confusion matrices across varying estimator counts consistently highlights the random forest algorithm's efficacy in accurately identifying hand signs.

In summary, this project effectively showcases the viability and efficiency of employing the random forest algorithm for hand sign recognition, laying a robust groundwork for further exploration and advancement in this domain. The insights garnered contribute significantly to the progression of gesture recognition technology, with potential applications spanning human-computer interaction, assistive technology, and accessibility solutions.

Looking ahead, integrating hand sign recognition into wearable devices such as smart glasses or IoT devices holds promise for enabling hands-free interaction and control of connected appliances. However, challenges related to accessibility and affordability may hinder widespread adoption within assistive technology. Educational applications present opportunities, yet creating immersive VR or AR environments with realistic hand sign recognition poses significant challenges. In the healthcare sector, regulatory hurdles and the need for thorough validation for clinical integration must be addressed.

IETE Sponsored 7<sup>th</sup> National Conference on Emerging Trends in Engineering, Science and Technology (NCTEST-7)

## REFERENCES

- [1] Alonazi, Mohammed, Hira Ansar, Naif Al Mudawi, Saud S. Alotaibi, Nour Abdullah Almualljy, Abdulwahab Alazeab, Ahmad Jalal, Jaekwang Kim and Moohong Min. "Smart Healthcare Hand Gesture Recognition Using CNN-Based Detector and Deep Belief Network." *IEEE Access* 11 (2023): 84922-84933.
- [2] Miah, Abu Saleh Musa, Md. Al Mehedi Hasan and Jungil Shin. "Dynamic Hand Gesture Recognition Using Multi-Branched Attention Based Graph and Generalized Ewi, Christine, Abbott Po Shun Chen and Henoch Juli Christanto. "Deep Learning for Highly Accurate Hand Recognition Based on Yolov7 Model." *Big Data Cogn. Comput.* 7 (2023): 53..
- [3] Athira Devraj and Aswathy K Nair. "Hand Gesture Signal Classification using Machine Learning." *International Conference on Communication and Signal Processing*, July 28 - 30, 2020, India. 978-1-7281-4988-2 ©2020 IEEE.
- [4] Ms. Geeshma Pala, Ms. Jagruti Bhagwan Jethwani, Mr. Satish Shivaji Kumbhar, Ms. Shruti Dilip Patil. "Machine Learning-based Hand Sign Recognition." *Proceedings of the International Conference on Artificial Intelligence and Smart Systems (ICAIS-2021)* IEEE Xplore Part Number: CFP21OAB-ART; ISBN: 978-1-7281-9537-7.
- [5] Chen, Renxiang and Xia Tian. "Gesture Detection and Recognition Based on Object Detection in Complex Background." *Applied Sciences* (2023): n. pag.
- [6] Kaushik Kumar Jha, Balram Kumar Jha, Suyash Ajit Byale, Vivek Sitaram Dadewa, Prof. Pushpanjali, S. Sajianshetti. "Sign Language Detection Using CNN." *International Journal for Research in Applied Science Engineering Technology (IJRASET)*; Volume 10 Issue XI Nov 2022.
- [7] Mohit Patil, Pramay Pathole, Hrishikesh Patil, Ashutosh Raut, Prof. S S Jadhav. "Indian Sign Language Recognition." *International Journal of Scientific Research Engineering Trends* Volume 6, Issue 4, July-Aug-2020, ISSN (Online): 2395-566X.
- [8] Pranali Manmode, Rupali Saha, Manisha N. Ammerkar. "Real-Time Hand Gesture Recognition." *International Journal of Scientific Research in Computer Science, Engineering and Information Technology (IJSRC- SEIT)* Volume 7, Issue 3 Page Number: 618-624: May-June-2021.
- [9] Ben Attallah, Bilel, Zheng Hu, Dhouha Bouchaala, Mohammed Hussain, Amir Ismail, Nabil Derbel and Olfa Kanoun. "Hand Sign Recognition System Based on EIT Imaging and Robust CNN Classification." *IEEE Sensors Journal* 22 (2022): 1729-1737.

**DrillBit**

The Report is Generated by DrillBit Plagiarism Detection Software

---

**Submission Information**

Author Name	Raghavendra Hegde, Samarth Shinnur, Samrudh B R, Sanjay S B
Title	Real Time Sign Language Recognition System Using Machine Learning
Paper/Submission ID	1869983
Submitted by	rnsce2001@gmail.com
Submission Date	2024-05-25 19:55:54
Total Pages, Total Words	75, 23353
Document type	Project Work

**Result Information**

**Similarity 15 %**

**Sources Type**

Sources Type	Percentage
Student Paper	0.15%
Journal/Publication	6.96%
Internet	7.89%

**Report Content**

Report Content	Percentage
Quotes	0.42%
Words < 14	5.97%

**Exclude Information**

Quotes	Not Excluded
References/Bibliography	Not Excluded
Source: Excluded < 14 Words	Not Excluded
Excluded Source	<b>0 %</b>
Excluded Phrases	Not Excluded

**Database Selection**

Language	English
Student Papers	Yes
Journals & publishers	Yes
Internet or Web	Yes
Institution Repository	Yes

A Unique QR Code use to View/Download/Share Pdf File