

# **RECOGNITION OF SIGN LANGUAGE GESTURES**

*A Project report submitted in partial fulfillment of the requirements for*

*the award of the degree of*

**BACHELOR OF TECHNOLOGY**

**IN**

**COMPUTER SCIENCE ENGINEERING (DATA SCIENCE)**

*Submitted by*

B Srimayee	320126551006
D Vishnu Vardhan	320126551013
CH Eswar Nadh	320126551010
G Koteswara Rao	320126551020

**Under the guidance of**

**Mrs. P Ratna Kumari**

**(Assistant Professor)**



**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

**(AI&ML, DS)**

**ANIL NEERUKONDA INSTITUTE OF TECHNOLOGY AND SCIENCES**

**(UGC AUTONOMOUS)**

*(Affiliated to AU, Approved by AICTE and Accredited by NBA)*

**Sangivalasa, Bheemili Mandal, Visakhapatnam**

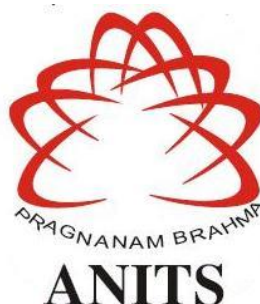
**dist. (A.P) 2020-2024**

# **Anil Neerukonda Institute of Technology & Sciences (Autonomous)**

(Permanent Affiliation by Andhra University & Approved by AICTE

Accredited by NBA (ECE, EEE, CSE, IT, Mech. Civil & Chemical) & NAAC)

Sangivalasa-531 162, Bheemunipatnam Mandal, Visakhapatnam District



## **BONAFIDE CERTIFICATE**

This is to certify that the project report entitled “**RECOGNITION OF SIGN LANGUAGE GESTURES**” submitted by Boyina Srimayee (320126551006), Dommeti Vishnu Vardhan (320126551013), Chillimuntha Eswar Nadh (320126551010), Gude Koteswara Rao (320126551020) in partial fulfillment of the requirements for the award of the degree of **Bachelor of Technology in Computer Science Engineering (Data Science)** of Anil Neerukonda Institute of technology and sciences, Visakhapatnam is a record of bonafide work carried out under my guidance and supervision.

Project Guide

**Mrs. P RATNA KUMARI**

Assitant Professor

Department of CSE(AI&ML, DS)

ANITS

Head of the Department

**Dr. K SELVANI DEEPTHI**

Professor

Department of CSE(AI&ML, DS)

ANITS

## DECLARATION

We **B.Srimayee, D.Vishnu Vardhan, CH.Eswar Nadh, G.Koteswara Rao**, of final semester B.Tech., in the department of CSE(AI&ML, DS) from ANITS, Visakhapatnam, hereby declare that the project work entitled **RECOGNITION OF SIGN LANGUAGE GESTURES** is carried out by us and submitted in partial fulfillment of the requirements for the award of Bachelor of Technology in Computer Science Engineering(Data Science), under Anil Neerukonda Institute of Technology & Sciences during the academic year 2020-2024 and has not been submitted to any other university for the award of any kind of degree.

B Srimayee	320126551006
D Vishnu Vardhan	320126551013
CH Eswar Nadh	320126551010
G Koteswara Rao	320126551020

## ACKNOWLEDGEMENT

We would like to express our deep gratitude to our project guide. Mrs. P Ratna Kumari, Assistant Professor Designation, Department of CSE (AI&ML, DS), ANITS, for guidance with unsurpassed knowledge and immense encouragement. We are grateful to Dr. K. Selvani Deepthi, Head of the Department, CSE (AI&ML, DS), for providing us with the required facilities for the completion of the project work. We are very much thankful to the Principal and Management, ANITS, Sangivalasa, for their encouragement and cooperation to carry out this work. We express our thanks to Project Coordinator Mrs. G. V. Gayathri, Asst. Professor, for continuous support and encouragement. We thank all teaching faculty of Department of CSE (AI&ML, DS), whose suggestions during reviews helped us in accomplishment of our project. We would like to thank Mrs. G. V. Gayathri of the Department of CSE (AI&ML, DS), ANITS for providing great assistance in accomplishment of our project. We also thank our technical staff for their support. We would like to thank our parents, friends, and classmates for their encouragement throughout our project period. At last, but not the least, we thank everyone for supporting us directly or indirectly in completing this project successfully.

B Srimayee	320126551006
D Vishnu Vardhan	320126551013
CH Eswar Nadh	320126551010
G Koteswara Rao	320126551020

# ABSTRACT

This paper presents a novel approach aimed at enhancing communication accessibility for individuals using sign language through the development of a Recognition of Sign Language Gestures system. The system leverages the WLASL dataset, a comprehensive collection of videos depicting American Sign Language (ASL) gestures, to train a deep learning model that can translate these gestures into text and speech. By integrating advanced technologies such as computer vision and deep learning, our approach opens new avenues for effective communication and inclusion for individuals with hearing impairments.

Our methodology incorporates OpenCV for efficient video processing, allowing us to handle large volumes of ASL gesture videos with ease. OpenCV's capabilities enable frame-by-frame analysis of gesture movements, ensuring high-quality input data for subsequent stages of processing. For feature extraction, we utilize ResNet50, a deep convolutional neural network renowned for its ability to capture complex patterns and fine-grained details within images and video frames.

The system's core analysis is performed using a CNN-GRU model, a hybrid architecture combining convolutional neural networks (CNNs) and gated recurrent units (GRUs). The CNN component excels in identifying spatial patterns within the video data, while the GRU component captures temporal dependencies and sequence information across multiple frames. This combination enables our system to understand and interpret the intricate dynamics of sign language gestures, resulting in high-quality translations from gestures to text and speech.

In terms of performance, our Recognition of Sign Language Gestures system achieves an impressive accuracy of 79% in converting ASL gestures into corresponding text and audio outputs. This level of accuracy demonstrates the efficacy of our approach in accurately translating sign language gestures, making communication more accessible for sign language users.

**Keywords:** Sign language recognition (SLR), Computer vision, Deep learning, Communication accessibility.

# CONTENTS

<b>ABSTRACT</b>	<b>i</b>
<b>LIST OF FIGURES</b>	<b>iv</b>
<b>LIST OF TABLES</b>	<b>vi</b>
<b>LIST OF ABBREVIATIONS</b>	<b>vii</b>
<b>CHAPTER 1 INTRODUCTION</b>	<b>01-05</b>
1.1 Introduction	01
1.2 Motivation of the work	04
1.3 Problem Statement	05
<b>CHAPTER 2 LITERATURE SURVEY</b>	<b>06-15</b>
2.1 Introduction	06
2.2 Existing Systems	07
2.2.1 Sign Language Recognition Using Machine Learning Algorithm	07
2.2.2 Sign Language Recognition using Mediapipe	08
2.2.3 Exploiting domain transformation and deep learning for hand gesture recognition using a low-cost dataglove	08
2.2.4 Gesture-Based Real-time Indian Sign Language Interpreter	10
2.2.5 Word-level Deep Sign Language Recognition from Video	10
2.2.6 Video Sign Language Recognition using Pose Extraction and Deep Learning Models	11
2.2.7 Sign Language Recognition using LSTM and Media Pipe	12
2.2.8 Sign language recognition using deep learning	12
2.2.9 Sign Language Recognition System	13
2.3 Proposed System	14
2.3.1 System Architecture	15

<b>CHAPTER 3 PROPOSED METHODOLOGY</b>	<b>16-32</b>
3.1 Introduction	16
3.2 Module Division	17
3.2.1 Dataset Collection	17
3.2.2 Approaches Prior to the Final Solution	20
3.2.3 Data Preprocessing	23
3.2.4 Feature Extraction	26
3.2.5 Model Building	27
3.2.6 Training the Model	29
3.2.7 Testing the Model	31
<b>CHAPTER 4 EXPERIMENTAL ANALYSIS AND RESULTS</b>	<b>33-69</b>
4.1 System Configuration	33
4.1.1 Software Requirements	33
4.1.1.1 Introduction to Python	33
4.1.1.2 Introduction to Kaggle	34
4.1.1.3 Introduction to Google Colab	35
4.1.1.4 Introduction to Jupyter Notebook and Jupyter Lab	36
4.1.1.5 Python Libraries	38
4.1.2 Hardware Requirements	47
4.2 Sample code	48
4.3 Accuracy Graphs	67
<b>CHAPTER 5 CONCLUSION AND FUTURE WORK</b>	<b>70-71</b>
5.1 Conclusion	70
5.2 Future Work	71
<b>BIBLIOGRAPHY</b>	<b>72-74</b>

## LIST OF FIGURES

<b>Fig No.</b>	<b>Topic Name</b>	<b>Page No.</b>
2.3.1.1	System Architecture for the Proposed Model	15
3.2.1.1	Examples highlighting the variety within our dataset, encompassing diverse backgrounds, lighting conditions, and signers with varying appearances.	18
3.2.1.2	Comparing our WLASL dataset with other ASL datasets. The "Mean" column shows how many videos, on average, are available for each sign.	19
3.2.3.1	Saving each video in its respective word folder	24
3.2.3.2	Breaking each video into frames	25
3.2.5.1	Model CNN+GRU	28
3.2.6.1	Accuracy	30
3.2.6.2	Confusion Matrix Graph	30
3.2.6.3	Accuracy vs Epoch Graph	31
4.2.1	Shape of loaded dataset	59
4.2.2	Words of WLASL100	59
4.2.3	Feature Extraction	59
4.2.4	Extracted Features	60
4.2.5	Model Summary	61
4.2.6	Test Accuracy	61
4.2.7	Confusion Matrix	62
4.2.8	Model Accuracy VS Epochs	62
4.2.9	Testing	63
4.2.10	User Interface	63
4.2.11	Action 1 – Preprocessing	63
4.2.12	Action 1 - Predicted Output	64
4.2.13	Action 2 - Capturing Webcam Input	64



4.2.14	Action 2 - Predicted Output	65
4.2.15	Action 3 – Capturing Webcam Input	65
4.2.16	Action 3 – Preprocessing	66
4.2.17	Action 3 – Predicted Output	66
4.3.1	Random Forest Model	67
4.3.2	Mediapipe + Baseline 2D Shallow CNN	67
4.3.3	Mediapipe + BiLSTM	68
4.3.4	Mediapipe + CNN	68
4.3.5	Mediapipe + SVM	69

## LIST OF TABLES

Table No.	Topic Name	Page No.
3.2.2	Comparison of our Model Accuracies	23

## **LIST OF ABBREVIATIONS**

SLR	Sign Language Recognition
CNN	Convolutional neural networks
ASL	American Sign Language
RNN	Recurrent neural networks
WLASL	World Level American Sign Language
TGCN	Temporal Graph Convolutional Networks
ISL	Indian Sign Language
GRU	Gated Recurrent Unit
BiLSTM	Bidirectional Long Short-Term Memory

# **CHAPTER 1 - INTRODUCTION**

## **1.1 INTRODUCTION**

The communication challenges encountered by the global mute community due to the widespread unfamiliarity with sign language among the general population are addressed by this project. A robust Sign Language Recognition (SLR) system is proposed to overcome this barrier. Leveraging advanced deep learning techniques and the WLASL dataset [7], sign language gestures are accurately identified and interpreted from video inputs. By translating these gestures into textual and spoken forms, accessibility and inclusivity are aimed to be enhanced, facilitating seamless communication between the mute community and society.

### **Video Processing**

Video processing encompasses a broad spectrum of techniques and algorithms aimed at manipulating and analyzing digital video data to extract meaningful insights or enhance visual quality. It is pivotal in various domains, including entertainment, healthcare, security, and education. At its core, video processing involves the manipulation of individual frames within a video sequence.

Preprocessing steps are often the initial stage of video processing, involving operations such as resizing, color space conversion, and noise reduction. These steps ensure consistency and improve the quality of subsequent analyses. Additionally, techniques like frame interpolation can be employed to generate additional frames, enhancing the smoothness of motion in videos.

One of the key aspects of video processing is video compression, where algorithms such as MPEG and H.264 are utilized to reduce the size of video files while preserving perceptual quality. These algorithms exploit temporal and spatial redundancies within the video sequence to achieve high compression ratios without significant loss of visual fidelity.

Object tracking is another important application of video processing, enabling the monitoring of specific objects or regions of interest within a video. This capability finds applications in surveillance, augmented reality, and human-computer interaction. Object

tracking algorithms utilize techniques such as feature extraction, motion estimation, and Kalman filtering to accurately predict the trajectory of moving objects.

Scene recognition algorithms analyze the content of video frames to categorize scenes based on their visual characteristics. Deep learning approaches, such as convolutional neural networks (CNNs), have shown remarkable performance in scene recognition tasks, enabling applications like video summarization, content-based retrieval, and automatic tagging.

In summary, video processing is a multifaceted field with applications ranging from entertainment to security. Its ability to extract valuable insights and facilitate efficient manipulation and dissemination of video content makes it indispensable in today's digital age.

## **Sign Language**

Sign language is a visual means of communication that uses hand gestures, facial expressions, and body movements to convey meaning. It is primarily used by individuals who are deaf or hard of hearing as their primary mode of communication. Sign languages are full-fledged languages with their own grammar and syntax, distinct from spoken languages.

There are many different sign languages around the world, each with its own unique vocabulary and grammar. For example, American Sign Language (ASL) is used primarily in the United States and Canada, while British Sign Language (BSL) is used in the United Kingdom. Other examples include Auslan (Australian Sign Language), LSF (French Sign Language), and JSL (Japanese Sign Language).

Sign languages are not universal; they vary from country to country and even region to region within the same country. However, there may be similarities and shared signs among sign languages, especially those that have historical or cultural connections.

Sign languages are not simply gestural representations of spoken languages; they are fully developed languages with their own linguistic structures. They involve the use of manual signs, facial expressions, and body movements to convey meaning. In addition to handshapes and movements, aspects such as palm orientation, facial expressions, and body posture also play important roles in sign language communication.

Sign languages are recognized as official languages in many countries, and efforts are ongoing to promote their recognition and use in various domains, including education, government, and media. Sign language interpreters play a crucial role in facilitating communication between deaf or hard-of-hearing individuals and those who do not use sign language.

### **Processing of Sign Language form Videos**

Sign language processing involves the analysis and interpretation of sign language gestures to facilitate communication for individuals who are deaf or hard of hearing. It encompasses various techniques and algorithms aimed at recognizing, translating, and generating sign language gestures.

Preprocessing steps are often employed to enhance the quality and consistency of sign language video inputs. These steps may include background subtraction, noise reduction, and hand segmentation to isolate and focus on the signer's gestures.

Sign language recognition algorithms analyze video frames to detect and interpret sign language gestures. These algorithms typically utilize computer vision techniques, such as feature extraction and pattern recognition, to identify key hand shapes, movements, and facial expressions associated with different signs.

Translation algorithms convert recognized sign language gestures into textual or spoken language, making the content accessible to individuals who do not understand sign language. These algorithms may employ natural language processing techniques to generate grammatically correct translations based on the recognized signs.

On the other hand, sign language generation algorithms produce sign language gestures from textual or spoken inputs. These algorithms often involve mapping linguistic elements to corresponding sign language gestures and generating fluent and natural-looking signing animations.

Sign language processing plays a crucial role in various applications, including education, accessibility, and communication. By enabling effective communication between individuals who are deaf or hard of hearing and those who do not understand sign language, sign language processing promotes inclusivity and equal access to information and services.

## **Gesture Recognition**

Gesture recognition is a field within sign language processing that focuses on identifying and understanding human gestures, particularly hand movements and body gestures, to interpret intent or communicate with electronic devices. It involves analyzing and interpreting the spatial and temporal characteristics of gestures to extract meaningful information.

Preprocessing steps in gesture recognition may involve similar techniques as those used in sign language processing, such as background subtraction, noise reduction, and hand segmentation. These steps help isolate the gestures from the surrounding environment and enhance the quality of input data.

Gesture recognition algorithms typically utilize computer vision techniques, such as feature extraction, pattern recognition, and machine learning, to analyze video frames or sensor data and identify specific gestures. These algorithms may also incorporate deep learning approaches, such as convolutional neural networks (CNNs) or recurrent neural networks (RNNs), to improve accuracy and robustness.

Once gestures are recognized, they can be translated into actionable commands or interpreted in context to enable interaction with devices or systems. For example, gesture recognition can be used in virtual reality systems to track hand movements and gestures, allowing users to interact with virtual objects or environments naturally.

Gesture recognition has applications in various domains, including human-computer interaction, gaming, healthcare, and automotive interfaces. For instance, it can be used in healthcare settings for gesture-based control of medical equipment or rehabilitation exercises, or in automotive interfaces for hands-free control of infotainment systems or driver assistance features.

## **1.2 MOTIVATION OF THE WORK**

The motivation for developing a Sign Language Recognition (SLR) system lies in addressing the longstanding communication barriers faced by the global mute community due to widespread unfamiliarity with sign language. Through the creation of a robust SLR system, individuals who are deaf or hard of hearing can be empowered to communicate effectively with the broader society, fostering inclusivity and reducing social isolation. Leveraging advanced deep learning techniques and comprehensive

datasets, such as WLASL, the communication gap is sought to be bridged, providing equitable access to education, employment, and social opportunities. This project represents a significant step toward creating a more inclusive and accessible world, where everyone, regardless of their communication abilities, can fully participate and thrive. Ultimately, the motivation stems from a commitment to equality, diversity, and the fundamental right of every individual to communicate and be understood.

### **1.3 PROBLEM STATEMENT**

The global mute community faces a significant communication barrier due to the widespread unfamiliarity with sign language among the general population. Mute individuals rely on sign language as their primary means of communication, which is not easily comprehensible to those without knowledge of it. This creates obstacles for mute individuals in various aspects of daily life, from professional settings to social interactions. This project addresses the need for a robust solution to bridge this communication gap by developing a highly efficient sign language recognition system. The key problem we aim to solve is accurate and efficient recognition of sign language gestures from inputs, with the goal of converting these gestures into both textual and spoken forms for seamless communication. This will facilitate more inclusive and equitable experiences for mute individuals across a range of environments.



## **CHAPTER 2 - LITERATURE SURVEY**

### **2.1 INTRODUCTION**

Sign language serves as a vital mode of communication for individuals with hearing impairments, offering a means to express thoughts, feelings, and ideas through hand gestures, facial expressions, and body movements. However, the inherent complexity and diversity of sign languages pose significant challenges for those who are not proficient in this visual language. To address this communication barrier and foster inclusivity, researchers and developers worldwide have been actively exploring innovative technologies to facilitate sign language recognition (SLR) and interpretation.

In recent years, several pioneering projects have emerged, each leveraging machine learning algorithms and advanced computer vision techniques to decode sign language gestures from video inputs. One such endeavor, led by Prof. Radha S. Shirbhate and a team from JSPM's BSIOTR – Wagholi, Pune, focuses on interpreting Indian Sign Language (ISL) alphabets. Their SLR system employs a methodology that encompasses data collection, image preprocessing, feature extraction, and classification using machine learning algorithms like Support Vector Machines (SVM) and Random Forest. Despite reported accuracies of 46.45% to 53.23%, their work highlights the potential for bridging the communication gap between speech-impaired individuals and those unfamiliar with sign language.

Similarly, the Sign Language Recognition system developed by Ketan Gomase and colleagues from Atharva College of Engineering, Mumbai, emphasizes real-time interpretation of American Sign Language (ASL) using Mediapipe. By employing hand and finger tracking solutions, their system achieves impressive accuracy ranging from 86% to 91%, thus offering promising prospects for enhancing communication accessibility for the deaf and hard-of-hearing community.

Building upon these endeavors, researchers Akshay Divkar and Rushikesh Bailkar, along with Dr. Chhaya S. Pawar, present a real-time Indian Sign Language interpreter that utilizes convolutional and recurrent neural networks (CNNs and RNNs). By leveraging spatial and temporal feature extraction, their system achieves a test accuracy of 54.2%, marking significant progress in the realm of ISL interpretation.

Moreover, the work by Gomase et al. from Atharva College of Engineering explores deep learning methods for word-level sign recognition, employing both appearance-based and pose-based approaches. By integrating neural network frameworks and innovative models like Temporal Graph Convolutional Networks (TGCN), their SLR system achieves remarkable accuracies for ASL letter recognition and word-level interpretation, offering a comprehensive solution for enhancing communication accessibility.

Furthermore, the master's project at San José State University spearheaded by a proficient research team delves into video sign language recognition using pose extraction and deep learning models. Through meticulous experimentation and optimization techniques, their approach achieves significant accuracy improvements on various datasets, underscoring the potential of leveraging advanced technologies to bridge communication gaps for individuals with hearing impairments.

In line with these groundbreaking efforts, our proposed system aims to revolutionize sign language interpretation by leveraging advanced deep learning techniques such as CNNs and RNNs. Our model prioritizes comprehensive word recognition, offering multimodal output in the form of text and speech to facilitate user-friendly communication for the mute community. By harnessing the power of machine learning and computer vision, we endeavor to bridge the communication gap, fostering greater understanding and inclusivity for all individuals, regardless of their proficiency in sign language.

## **2.2 EXISTING SYSTEMS**

### **2.2.1 Sign language Recognition Using Machine Learning Algorithm [1]**

The prototype or model, developed by Prof. Radha S. Shirbhate, Mr. Vedant D. Shinde, Ms. Sanam A. Metkari, Ms. Pooja U. Borkar, and Ms. Mayuri A. Khandge from JSPM's BSIOTR – Wagholi, Pune, is a Sign Language Recognition (SLR) system designed to interpret sign language gestures from video inputs, specifically focusing on Indian Sign Language (ISL) alphabets. This system aims to bridge the communication gap between speech-impaired individuals who use sign language and those who do not understand it. The methodology involves collecting a dataset, preprocessing images, segmenting skin parts, extracting relevant features using techniques like HU's moments, and employing

machine learning algorithms such as Support Vector Machines (SVM) and Random Forest for classification. The dataset used includes images of ISL gestures, with a portion being reserved for training and testing. Various algorithms and techniques are explored to achieve accurate recognition, with accuracies of 46.45% for Random Forest and 53.23% for hierarchical classification using SVM reported. Challenges addressed include the variance in sign language with locality and the lack of standard datasets for ISL. The system holds potential for further development to include recognition of words and sentences in ISL, as well as integration of additional feature extraction algorithms and classifiers to improve accuracy.

### **2.2.2 Sign Language Recognition using Mediapipe [2]**

The Sign Language Recognition (SLR) system, developed by Ketan Gomase, Akshata Dhanawade, Prasad Gurav, and Sandesh Lokare from the Department of Electronics & Telecommunication Engineering at Atharva College of Engineering, Mumbai, Maharashtra, India, focuses on interpreting sign language gestures using Mediapipe. The system aims to address the communication barriers faced by speech-impaired individuals by translating sign language into text or speech. The methodology involves data acquisition, preprocessing, feature extraction, segmentation, and evaluation of results. Mediapipe Hands, a hand and finger tracking solution using machine learning, is utilized to understand 3D local hand marks from images. The objective is to create a neural network capable of recognizing American Sign Language (ASL) alphabet characters from handwritten signatures. The system aims to improve communication accessibility for the deaf and hard-of-hearing community. Various hardware and software requirements are outlined, including the use of Python, OpenCV, TensorFlow, and Keras. The system achieves real-time hand detection and recognition of ASL letters with an average accuracy ranging from 86% to 91%. Future research directions include enhancing Human-Computer Interoperability (HCI) using more advanced algorithms.

### **2.2.3 Exploiting domain transformation and deep learning for hand gesture recognition using a low-cost dataglove [3]**

Hand gesture recognition has emerged as a crucial aspect of human-computer interaction, facilitating intuitive communication between individuals and devices. The study conducted by Md.AhasanAtick Faisal, Farhan FuadAbir, Mosabber UddinAhmed, and MdAtiqur RahmanAhad explores the efficacy of a low-cost dataglove for hand

gesture classification leveraging deep learning methodologies. This research endeavor is propelled by recent advancements in hardware accessibility and the proliferation of deep learning algorithms. The study methodology encompasses the development of a dataglove equipped with flex sensors, an inertial measurement unit (IMU), and a microcontroller for onboard processing and wireless connectivity. Data collection involved gathering samples from 25 subjects performing both static and dynamic American Sign Language (ASL) gestures, amounting to 24 static and 16 dynamic gestures. Notably, a novel Spatial Projection Image-based technique was proposed for dynamic hand gesture recognition, alongside the exploration of a parallel-path neural network architecture for enhanced multimodal data processing. The hardware configuration, including sensor specifications and processing capabilities, was meticulously detailed to ensure transparency and reproducibility. Ethical considerations were also addressed, with informed consent obtained from all participants and approval from the relevant institutional review board. The collected dataset was made publicly available to facilitate further research and innovation in the field. Through rigorous data preprocessing techniques such as gravity compensation, axis rotation, and normalization, the researchers aimed to refine the input data for subsequent analysis. The proposed Spatial Projection Image generation method aimed to mitigate challenges associated with dynamic sign language recognition, particularly temporal dependency and variability in hand size. This technique involved converting 3D acceleration data into 2D spatial projections, enabling pattern recognition irrespective of temporal variations. Furthermore, the study introduced a comprehensive architecture for gesture recognition, comprising distinct modules for static and dynamic gestures. For static gestures, parallel 1D ConvNet blocks were employed, while a combination of ConvNet blocks and MobileNetV2 architecture was utilized for dynamic gestures. Evaluation metrics included macro-averaged precision, recall, F1-score, and accuracy, with leave-one-out-cross-validation (LOOCV) employed for robust model assessment. Logistic regression yielded a final accuracy of 81.19% for static gestures and 94.05% for dynamic gestures. Random forest achieved an accuracy of 81.25% for static gestures and 94.05% for dynamic gestures. The 1D CNN model attained a final accuracy of 81.02% for static gestures and 89.27% for dynamic gestures. Overall, this research contributes to the advancement of generalized hand gesture recognition techniques, with implications for improved human-computer interaction and accessibility.

#### **2.2.4 Gesture Based Real-time Indian Sign Language Interpreter [4]**

The article, authored by Akshay Divkar, Rushikesh Bailkar, and Dr. Chhaya S. Pawar, presents the development of a real-time Indian Sign Language (ISL) interpreter based on hand gestures. It aims to address communication barriers for hearing and speech-impaired individuals. The system utilizes a vision-based approach, leveraging Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN) models for spatial and temporal feature extraction, respectively. The CNN model extracts spatial features from video frames, while the RNN model captures temporal dependencies for gesture recognition. The authors curated a dataset comprising 456 videos covering 38 ISL gestures, facilitating training and testing of the system. Through image processing techniques, hand gestures are segmented and processed to extract relevant features. The system achieved a test accuracy of 54.2%, demonstrating its potential effectiveness in real-time ISL interpretation. Future work includes enhancing accuracy for continuous sign language gestures and exploring hybrid models for improved performance.

#### **2.2.5 Word-level Deep Sign Language Recognition from Video [5]**

A new large-scale Word-Level American Sign Language (WLASL) video dataset has been introduced by Dongxu Li, Cristian Rodriguez Opazo, Xin Yu, and Hongdong Li, affiliated with The Australian National University, Australian Centre for Robotic Vision (ACRV). This dataset, comprising over 2000 words performed by more than 100 signers, has been made publicly available to facilitate research in word-level sign recognition. In their investigation, various deep learning methods are experimented with using the WLASL dataset. Holistic visual appearance-based approaches are compared with 2D human pose-based approaches to evaluate their performance in large-scale scenarios. Additionally, a novel pose-based temporal graph convolution network (Pose-TGCN) is proposed, designed to model spatial and temporal dependencies in human pose trajectories simultaneously. Through these experiments, the effectiveness and challenges of the WLASL dataset [7] are demonstrated. Word-level recognition accuracies of up to 62.63% at top-10 accuracy are achieved, showcasing the validity and potential of the dataset for advancing sign language recognition research. Furthermore, the dataset and baseline deep models are made available to the research community for further exploration and development. Overall, the paper contributes to the field of sign language recognition by providing a comprehensive evaluation of deep learning

methods on a large-scale dataset, paving the way for future advancements in word-level sign recognition algorithms.

### **2.2.6 Video Sign Language Recognition using Pose Extraction and Deep Learning Models [6]**

The master's project titled "Video Sign Language Recognition using Pose Extraction and Deep Learning Models" at San José State University led by a team of proficient researchers aimed to enhance communication accessibility for the deaf and hard of hearing community. Employing innovative methodologies and extensive experimentation, the project explored various deep learning techniques, including CNNs, RNNs, Transformers, and GCNs, to extract features and model temporal dependencies in sign language videos. Two primary approaches were pursued: appearance-based recognition and pose-based recognition. Baseline models yielded a modest accuracy of 43.02%, which significantly improved to 46.80% through meticulous optimization techniques such as data augmentation and hyperparameter tuning. Further experiments with different pose extraction methods, including OpenPose, MediaPipe, and SSTCN, revealed varying performance levels. MediaPipe emerged as the top performer on the WLASL100 dataset, achieving an accuracy of 55.31%, followed closely by OpenPose with 46.80%, while SSTCN achieved 42.53%. Ensemble techniques combining predictions from MediaPipe and OpenPose led to a significant accuracy boost, reaching 55.96%. Additional experiments on the AUTSL dataset validated the scalability and generalization of the proposed approach. On AUTSL, the Transformer model achieved a baseline accuracy of 76.57%, which increased to 78.6% with data augmentation and OpenPose, and further improved to 81.08% with MediaPipe. The ensemble of MediaPipe and OpenPose predictions yielded the highest accuracy of 82.90%. Comparative analyses with existing literature highlighted the project's achievements, with accuracies consistently surpassing previous works on both WLASL100 and AUTSL datasets. This comprehensive approach underscores the potential of leveraging pose extraction and deep learning models to bridge communication gaps and enhance accessibility for individuals with hearing impairments.

### **2.2.7 Sign Language Recognition using LSTM and Media Pipe [13]**

The Sign Language Recognition (SLR) system, developed by G. Mallikarjuna Rao, P. A. Sujasri, Cheguri Sowmya, Soma Anitha, Dharavath Mamatha, and Ramavath Alivela from Gokaraju Rangaraju Institute Of Engineering and Technology, Hyderabad, India, is a sophisticated solution aimed at interpreting sign language gestures with precision. Leveraging cutting-edge technologies like Mediapipe Hands for hand and finger tracking, the system effectively captures 3D local hand marks from images. Employing a neural network framework for ASL alphabet recognition, the researchers explore various deep learning methods for word-level sign recognition, including appearance-based and pose-based approaches. In the realm of appearance-based methods, the team retrained VGG backbone and GRU models, alongside fine-tuned I3D, demonstrating superior performance over the VGG-GRU baseline. Introducing a novel Temporal Graph Convolutional Network (TGCN) for pose-based methods, the system adeptly captures both temporal and spatial dependencies in pose trajectories. With an average accuracy ranging from 86% to 91%, the SLR system achieves real-time hand detection and ASL letter recognition. Furthermore, experimental findings reveal comparable classification performance on a dataset of 2,000 words, achieving up to 62.63% accuracy.

### **2.2.8 Sign language recognition using deep learning [14]**

Taras Shevchenko National University of Kyiv, Faculty of Mechanics and Mathematics, Department of Algebra and Computer Mathematics presents a course work thesis on Sign Language Recognition using Deep Learning by 3rd year students Roman Polishchenko and Valentyn Kofanov under the supervision of Professor Andriy Oliynyk and scientific consultant Olga Sinelnikova. The research introduces a visual-based system for Sign Language Recognition employing convolutional neural networks (CNN) for spatial features analysis and recurrent neural networks (RNN) with Long Short-Term Memory (LSTM) cells for temporal features analysis, achieving up to 90% accuracy for a small dataset subset (10 signs) and 83% for the entire dataset (64 signs). The study also includes the development of a cross-platform application using the Kivy Python framework for real-time gesture recognition from both saved videos and live camera feed. Analysis of the problem highlights the significance of automated sign language recognition and translation, exploring existing SLR approaches and neural

network architectures. The work utilizes Python 3 with OpenCV, Keras, NumPy, and Matplotlib libraries, alongside the Kivy framework, and leverages Google Colab for computational resources. Experimentation involves preprocessing the Argentinean Sign Language dataset, implementing a CNN-LSTM model architecture, and evaluating results. Despite computational and dataset limitations, the study underscores the efficacy of combining CNNs and RNNs for sign language recognition, with future plans to create a dataset for Ukrainian Sign Language and develop real-time translation applications.

### **2.2.9 Sign Language Recognition System [15]**

The "Indian Journal of Software Engineering and Project Management (IJSEPM)" presents a study on a Sign Language Recognition System, authored by Pooja M.R, Meghana M, Praful Koppalkar, Bopanna M J, Harshith Bhaskar, and Anusha Hullali. The paper introduces a novel technique aimed at addressing communication challenges faced by individuals with disabilities such as deafness and muteness. By leveraging Histogram Oriented Gradient (HOG) in conjunction with Artificial Neural Network (ANN), the system enables virtual communication without the need for sensors. Utilizing a webcam, users input gestures, which are then processed to recognize and interpret sign language, bridging the communication gap between impaired and normal individuals. The proposed approach facilitates two-way communication by converting sign language into text and voice, enhancing accessibility for the hearing and speech impaired.

The methodology section outlines the implementation, with a focus on frontend and backend development. ReactJS and Bootstrap are utilized for the frontend, offering a reusable component-based UI, while NodeJS and Express handle the backend, with data management handled by MongoDB. The system employs Python for algorithm development, leveraging machine learning libraries such as NumPy and OpenCV for training and processing image data.

Results and discussion showcase the system's performance through confusion matrices, demonstrating its ability to recognize sign language gestures accurately. The conclusion highlights the creation of an American Sign Language translator based on a web application employing a CNN classifier. While acknowledging limitations in dataset diversity, the study suggests that with additional data under natural conditions, models



could achieve higher accuracy, providing robust communication support for all sign language gestures.

The paper concludes with references to related works and author profiles, detailing the affiliations and research interests of the contributors involved in the development of the Sign Language Recognition System.

## **2.3 PROPOSED SYSTEM**

Our proposed system is designed with a focus on facilitating effective communication for the mute community, ensuring that they have a user-friendly tool for bridging the communication gap. By leveraging advanced deep learning techniques, such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), the system is capable of recognizing a wide range of sign language gestures from video inputs with high accuracy. The output from the model is multimodal, providing not just text representation but also converting sign language gestures into speech, thereby allowing for greater versatility in communication.

Unlike traditional sign language recognition systems that might focus solely on alphabet interpretation, our proposed system goes a step further by offering the capability to interpret complete words. This extended functionality significantly enhances the convenience for users, allowing for smoother and more efficient communication. By prioritizing comprehensive word recognition, our system aims to foster a more inclusive environment, making it easier for those who rely on sign language to interact with those who might not be familiar with it.

In addition to the deep learning backbone, our system incorporates a user-friendly interface, ensuring that the technology is accessible to a wider audience. This approach helps bridge the communication gap between sign language users and non-users, contributing to a more inclusive society. Through the integration of state-of-the-art techniques and a focus on user needs, the proposed system aims to promote greater understanding and accessibility, ultimately transforming the way people communicate across diverse communities.

### 2.3.1 SYSTEM ARCHITECTURE

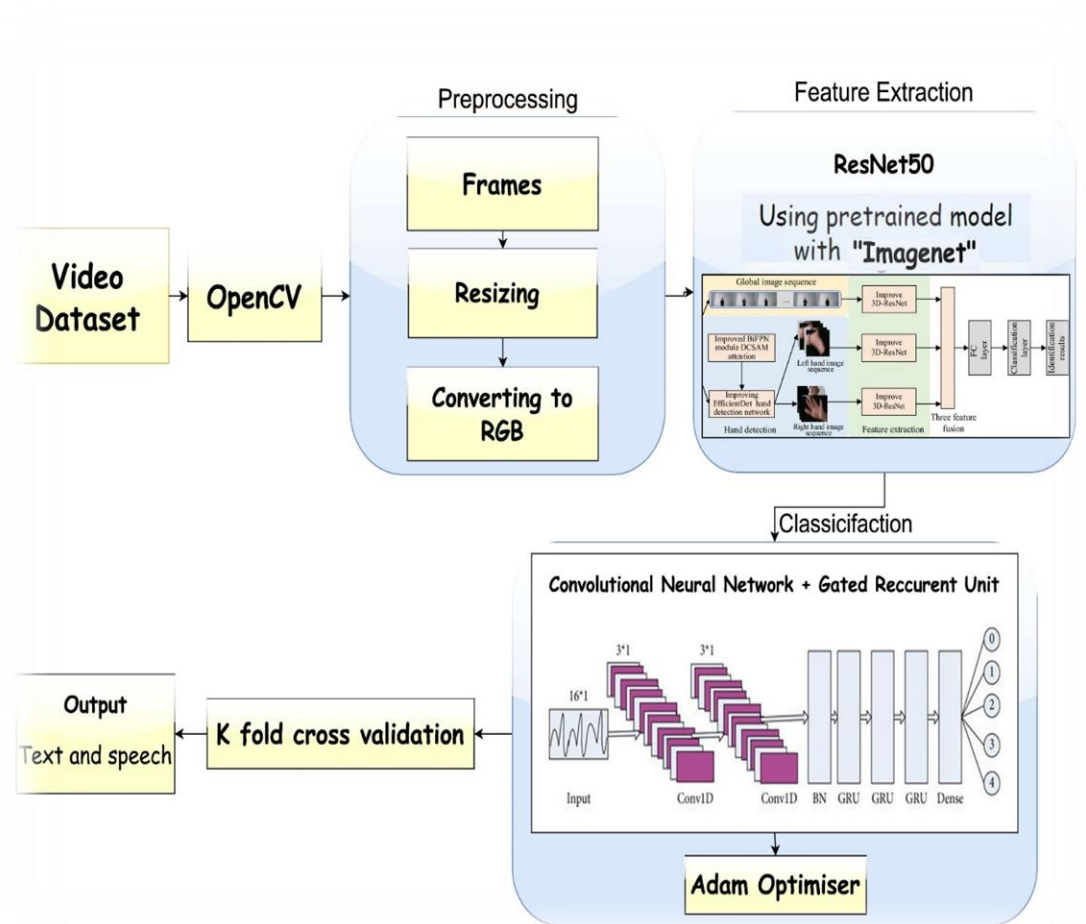


Fig 2.3.1.1: System Architecture for the Proposed Model

## CHAPTER 3 - PROPOSED METHODOLOGY

### 3.1 INTRODUCTION

Effective communication remains a formidable challenge for the global mute community, primarily due to the widespread unfamiliarity with sign language among the general population. This study aims to address this communication barrier by developing a robust Sign Language Recognition (SLR) system capable of accurately identifying and interpreting sign language gestures from video inputs. The proposed SLR system leverages state-of-the-art deep learning techniques to bridge the communication gap between the mute community and the broader society.

The architecture of the SLR system is intricately linked to the characteristics and requirements of the WLASL dataset [7]. Preprocessing of video frames from the WLASL dataset entails resizing and conversion to RGB format using OpenCV, ensuring compatibility with subsequent processing stages. The extracted frames serve as inputs for the feature extraction phase, where the ResNet50 model, pre-trained on the Imagenet dataset [12], is employed to extract high-level features from the sign language gestures. The utilization of a pre-trained model is essential for capturing complex visual patterns inherent in sign language expressions present in the WLASL dataset.

Furthermore, the Improved BiFPN module and attention mechanism enhance feature extraction by capturing subtle nuances and spatial relationships within the sign language gestures. The hand detection network, specifically the Improved EfficientDet model, plays a pivotal role in detecting and extracting features from left and right-hand image sequences within the WLASL dataset [7]. This step is crucial as hand gestures constitute a fundamental aspect of sign language communication, and accurate detection is paramount for effective recognition.

The proposed SLR system not only aims to accurately recognize sign language gestures but also to translate them into both textual and spoken forms, thereby enhancing accessibility for a wider audience. By employing advanced deep learning techniques and leveraging the rich annotations provided by the WLASL dataset [7], the proposed SLR system endeavors to facilitate seamless communication and foster inclusivity for the global mute community.

In summary, the SLR system represents a significant step forward in addressing the communication challenges faced by the mute community. Through the integration of advanced technologies and the utilization of a comprehensive dataset, the system holds promise in enhancing communication accessibility and promoting inclusivity on a global scale.

In the following sections, we will describe each component of the proposed SLR system in detail, including the preprocessing steps, feature extraction techniques, hand detection network, and classification layer. We will also provide experimental results and analysis to demonstrate the system's effectiveness and efficiency.

## 3.2 MODULE DIVISION

### 3.2.1 DATASET COLLECTION

The WLASL dataset [7] comprises video recordings depicting individuals executing American Sign Language (ASL) signs corresponding to English words. Originally introduced by Dongxu et al., this dataset features diverse recordings captured under various lighting conditions, contexts, and angles, ensuring comprehensive views of hand, elbow, and body motions during each gesture. Each video is meticulously annotated with consistent metadata, including gloss, signer's bounding box, and the start and end frames of the sign. Additionally, a unique identifier is assigned to differentiate between different signers. All metadata attributes are detailed below:

- **gloss:** The sign language word or gesture being depicted in the video instance.
- **bbox:** The bounding box coordinates (x\_min, y\_min, x\_max, y\_max) specifying the region containing the sign gesture within the video frame.
- **fps:** Frames per second, indicating the frame rate of the video.
- **frame\_start:** The starting frame number of the sign gesture within the video.
- **frame\_end:** The ending frame number of the sign gesture within the video.
- **variation\_id:** A unique identifier for different variations or renditions of the sign gesture.
- **instance\_id:** A unique identifier for each instance of the sign gesture within the dataset.

- **signer\_id**: The identifier for the signer or individual performing the sign gesture.
- **source**: The source or origin of the video, such as the dataset or website where it was obtained.
- **split**: Indicates the dataset split (e.g., train, validation, test) to which the video instance belongs.
- **url**: The URL or web link to access the video file.
- **video\_id**: A unique identifier for the video containing the sign gesture.

The authors of [21] partitioned the dataset into four distinct subsets, denoted as WLASL100, WLASL300, WLASL1000, and WLASL2000. These subsets were created by selecting the top-K (where  $k=100, 300, 1000$ , and  $2000$ ) glosses based on predefined criteria outlined in the paper. For the purposes of this project, our focus will be on utilizing the WLASL100 subset to showcase the outcomes of our approach and to establish a foundation for future endeavors involving other datasets.



**Fig 3.2.1.1: Examples highlighting the variety within our dataset, encompassing diverse backgrounds, lighting conditions, and signers with varying appearances.**

The WLASL dataset [7] provides a unique feature for users to access videos quickly: a set of URL links that point to downloadable video content. Although the majority of these links take users to videos on YouTube, there are also several instances where the videos are hosted on other platforms, such as alternative streaming services or cloud-based storage. The creators of the dataset have gone the extra mile by supplying source code that not only helps to load the dataset but also allows for direct video downloads from these links. This added functionality significantly improves the ease of use for researchers and developers, simplifying the entire process of accessing and working with the dataset's extensive video collection. This approach makes the WLASL dataset more versatile, allowing users to focus on their analysis or development tasks without getting bogged down by the complexities of downloading and managing video files from different sources.

Datasets	#Gloss	#Videos	Mean	#Signers	Year
Purdue RVL-SLLL [70]	39	546	14	14	2006
RWTH-BOSTON-50 [79]	50	483	9.7	3	2005
Boston ASLLVD [6]	2,742	9,794	3.6	6	2008
WLASL100	100	2,038	20.4	97	2019
WLASL300	300	5,117	17.1	109	2019
WLASL1000	1,000	13,168	13.2	116	2019
WLASL2000	2,000	21,083	10.5	119	2019

**Fig 3.2.1.2: Comparing our WLASL dataset with other ASL datasets. The “Mean” column shows how many videos, on average, are available for each sign.**

Due to several videos being deleted or having restricted access, our team was able to collect only a fraction of the intended WLASL100 dataset. Originally, the full dataset contained 2,038 videos spread across 100 different sign language glosses. Unfortunately, many of these videos were no longer accessible, leaving us with just 1,323 recoverable videos.

Despite this setback, we divided the available videos into training, validation, and testing sets. For training, we had 1,036 videos, while the validation and testing sets contained 166 and 121 videos, respectively. This distribution was crucial for training our sign language recognition model.

To ensure the model's accuracy and robustness, we searched for the missing videos elsewhere. By gathering replacement videos from other sources, we managed to supplement our dataset and ensure that our SLR system had enough data to learn from, allowing for a more comprehensive and accurate training process. This additional effort helped us make the most out of the remaining dataset and improved the overall quality of our sign language recognition system.

### **3.2.2 Approaches Prior to the Final Solution**

To find the best model for our sign language recognition (SLR) system, we explored multiple approaches, focusing on both preprocessing and deep learning techniques. Our dataset provided flexibility, allowing us to experiment with different methods to determine the most effective solution. Let's delve deeper into the different approaches we took, highlighting the key steps and outcomes.

#### **Preprocessing Approaches:**

We utilized two main preprocessing methods: one incorporating MediaPipe and another without it. MediaPipe, an open-source framework by Google, offers customizable machine learning solutions for multimedia processing tasks, including object detection, tracking, and hand and pose estimation. By leveraging MediaPipe's capabilities, we aimed to enrich our feature set, thus improving our model's accuracy.

#### **Deep Learning Models:**

Our exploration into deep learning involved training 13 different models to evaluate their accuracy. Each model had unique advantages and disadvantages, contributing to our understanding of how to build an effective SLR system.

1. **Random Forest with MediaPipe:** We started with a Random Forest classifier using MediaPipe, achieving a testing accuracy of 28%. This approach demonstrated the

impact of feature enrichment facilitated by MediaPipe, offering a baseline for our subsequent experiments.

2. **Shallow CNN with MediaPipe:** To improve upon the initial accuracy, we turned to deep learning methodologies, training a 2D Shallow CNN with MediaPipe features. This yielded a training accuracy of 82% and a testing accuracy of 42%, indicating that MediaPipe-enhanced feature extraction could enhance our model's performance.
3. **CNNs with Recurrent Layers:** We integrated recurrent neural networks (RNNs) to capture temporal dependencies in our data. A combination of CNNs and gated recurrent units (GRUs) resulted in a testing accuracy of 50%, showcasing the benefits of sequential modeling.
4. **BiLSTM with Attention and MediaPipe:** Next, we experimented with a bidirectional Long Short-Term Memory (BiLSTM) architecture with attention mechanisms, again utilizing MediaPipe features. This approach yielded a testing accuracy of 58%, highlighting the effectiveness of combining advanced deep-learning techniques with MediaPipe's enriched features.
5. **10-Fold Cross-Validation with BiLSTM:** To enhance model robustness, we used 10-fold cross-validation with BiLSTM models. This resulted in a testing accuracy of 70%, demonstrating the reliability of our approach across different data splits.
6. **Standalone LSTM with MediaPipe:** This approach involved using standalone Long Short-Term Memory (LSTM) architectures enhanced by MediaPipe features. With a testing accuracy of 76%, this method underscored the efficacy of our chosen techniques.



## Final Approach with Transfer Learning

Our ultimate approach focused on harnessing the benefits of transfer learning by using a pre-trained ResNet50 model [11], albeit without its top classification layer. This strategy allowed us to leverage the extensive features learned from ImageNet [12], a comprehensive image dataset, thereby providing a solid foundation for our model's learning process. By reusing these features, we could build on the model's inherent understanding of complex visual patterns.

To improve our model's ability to recognize and process sign language, we integrated the ResNet50 features with a combination of Convolutional Neural Networks (CNNs) and Gated Recurrent Units (GRUs). This hybrid approach enabled us to benefit from the strengths of both CNNs and GRUs, capturing both spatial and temporal information in the sign language sequences. The Adam optimizer was selected to fine-tune the learning process, ensuring efficient and effective convergence during training.

The results were promising. This final configuration achieved a testing accuracy of 80%, demonstrating the significant impact of transfer learning on our sign language recognition system. The high accuracy reflects the capability of our model to understand and categorize complex sign language gestures accurately.

Over the course of this project, we faced various challenges that required us to adapt and refine our approach. The final solution embodies our journey of experimentation, illustrating how combining advanced preprocessing techniques like MediaPipe with sophisticated deep-learning architectures can yield a robust and reliable sign language recognition system. This outcome represents a significant step forward in creating technology that bridges communication gaps and supports inclusive interaction for users of sign language.

<i><b>Model</b></i>	<i><b>Test Accuracy</b></i>
<i>Mediapipe + Random Forest</i>	28%
<i>Mediapipe + Baseline 2D Shallow CNN model</i>	42%
<i>CNN+GRU</i>	50%
<i>Mediapipe + RNN (Bidirectional LSTM with attention)</i>	58%
<i>Mediapipe + BiLstm + 10 cross fold</i>	70%
<i>Mediapipe + LSTM</i>	76%
<i><b>Resnet50 + CNN + GRU</b></i>	<b>80%</b>

**Table 3.2.2 Comparison of our Model Accuracies**

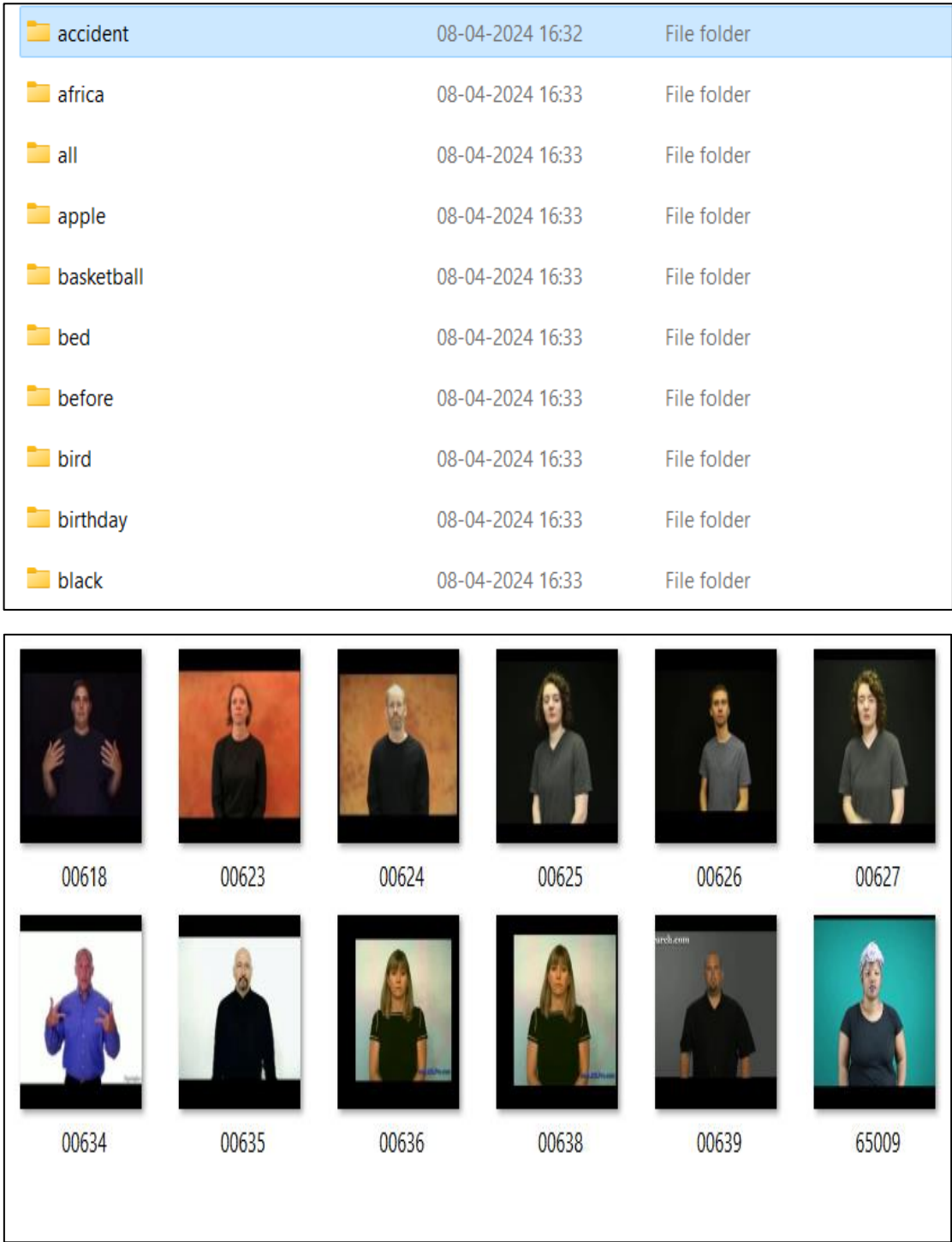
### **3.2.3 DATA PREPROCESSING**

#### **1. Video Processing**

In this phase, we focus on loading and processing videos using OpenCV, a widely recognized computer vision library. To access each video file, we employ OpenCV's `VideoCapture` function, allowing us to retrieve content from a wide range of video formats. Each video is then sorted into its corresponding category based on the specific word it represents. These categorized videos are stored in structured directories, ensuring that the intended meaning of each word is maintained. This systematic organization allows for easy retrieval and consistent analysis of video data.

The capabilities of OpenCV play a central role in this process. Its robust functions enable precise handling of video frames, including resizing, cropping, and color conversion. This versatility ensures that video processing is accurate and efficient. OpenCV's extensive tools and resources allow for seamless integration into various computer vision workflows, making it a preferred choice for tasks that involve video manipulation and analysis. This enhanced efficiency can be leveraged across a broad spectrum of applications, from surveillance systems that require real-time video analysis

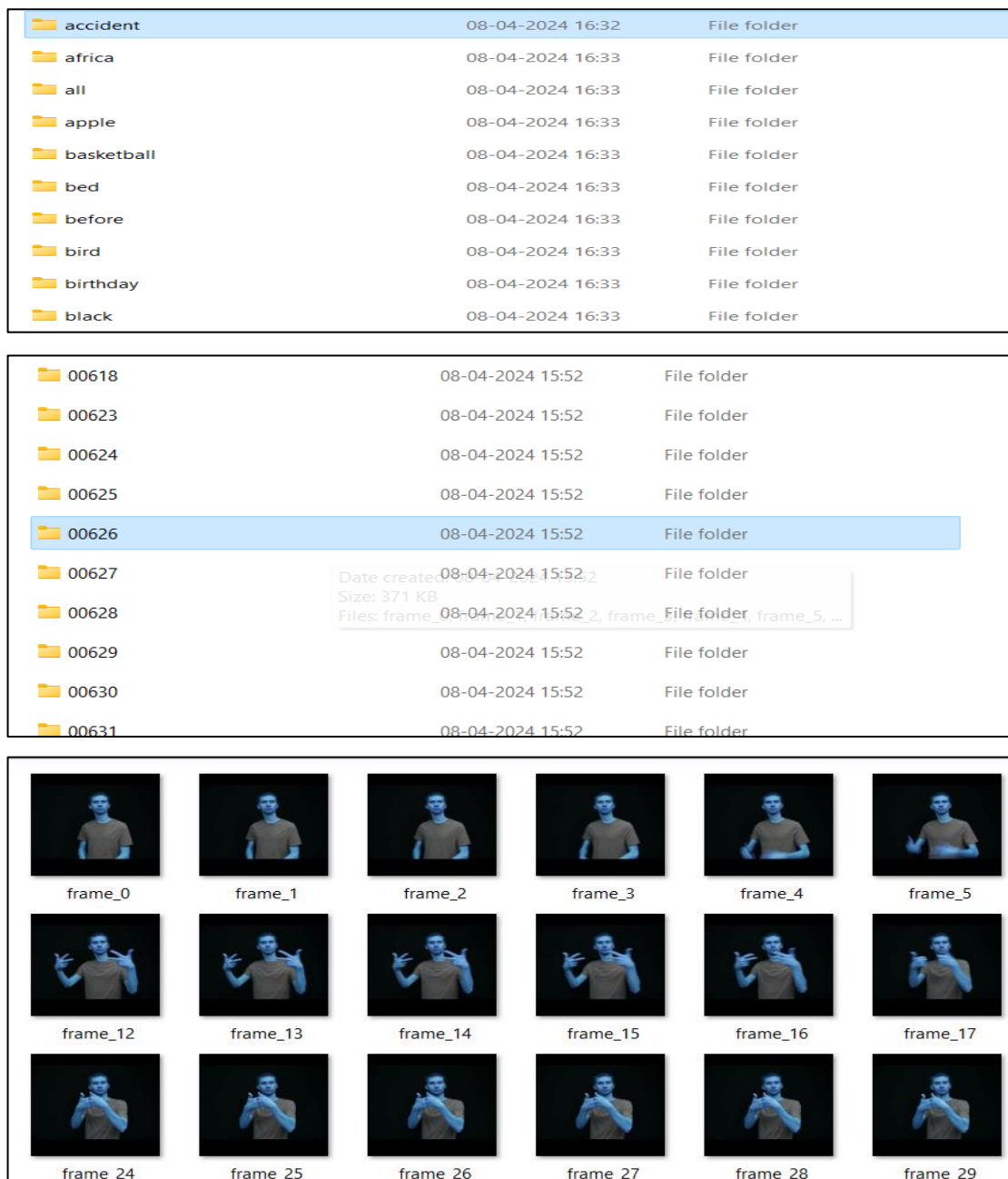
to augmented reality projects where precise video processing is crucial. Overall, OpenCV's features significantly contribute to the reliability and scalability of our video processing workflow, forming the foundation for further exploration and development in the field of computer vision.



**Fig 3.2.3.1 Saving each video in its respective word folder**

## 2. Extracting Frames

In this step, we're getting the videos ready to work with. These videos show sign language gestures. We want to turn each video into a bunch of pictures. First, we load the videos. Then, we go through each video one by one. For each video, we make a special folder to keep all the pictures from that video. Finally, we go through the video frame by frame, taking out each picture and saving it in the folder we made. This makes it easier for the computer to understand and work with the sign language gestures later.



**Fig 3.2.3.2 Breaking each video into frames**

### 3. Resizing

The frames taken from the videos are changed to a specific size. We call this size the "target size." Usually, it's set to (224, 224) unless we say otherwise. Changing the size makes sure that all frames have the same dimensions. This is crucial for training a machine-learning model because it needs consistency to learn properly. So, by resizing the frames, we guarantee that the model gets uniform inputs, which helps it understand and learn from the data more effectively.

### 4. Converting to RGB

During this step, the resized frames undergo a color space conversion. By default, OpenCV uses the BGR color space, but many machine-learning models prefer the RGB color space for input data. To align with this expectation, the frames are converted from BGR to RGB using the `cvtColor` function. This ensures compatibility between the frame data and the machine-learning model's requirements, allowing for seamless integration and accurate processing of the visual information contained within the frames.

#### 3.2.4 FEATURE EXTRACTION

Once we have extracted and preprocessed the frames, we use a pre-trained ResNet50 model to extract features from the frames. The ResNet50 model [11] is a deep convolutional neural network that has been pre-trained on the ImageNet dataset [12], which contains over 1.2 million images and 1000 classes. In our implementation, we exclude the top classification layer and freeze the weights of the pre-trained layers. This means that the weights of the pre-trained layers will not be updated during training, and only the newly added layers will be trained.

We define a function to extract features from the video frames using the ResNet50 model. The function preprocesses the input frames for the ResNet50 model [11] and extracts features using the model. The extracted features are then pooled across frames using average pooling.

We then define a function to extract features from all videos in the dataset. The function loops through each video in the dataset, extracts features using the `extract_features` function, and stores the extracted features in a list. The list is then converted to a numpy array and returned.

Finally, we extract features from the dataset and print the shape of the extracted features to the console. The extracted features can be used as input to a machine-learning model for sign language recognition.

### **3.2.5 MODEL BUILDING**

To build the classification model, we use the Keras Sequential API to define a model architecture that consists of Conv1D, GRU, GlobalAveragePooling1D, and Dense layers.

The first layer in our model is a Conv1D layer with 64 filters, a kernel size of 3, and a ReLU activation function. This layer is responsible for extracting spatial features from the input feature vectors. We set the input shape to be the same as the shape of the extracted feature vectors.

Next, we add a GRU layer with 64 units and set the `return_sequences` parameter to `True`. This allows us to retain the sequential output of the GRU layer, which is useful for the next layer. The GRU layer is added to capture the temporal dependencies in the feature vectors.

After capturing the temporal dependencies, we add a GlobalAveragePooling1D layer to compute the average of the sequential output from the GRU layer. This helps to reduce the dimensionality of the output and provides a fixed-length representation of the input feature vectors.

We then add a fully connected dense layer with 128 units and a ReLU activation function. This layer helps to learn more complex representations from the input data.

In the final stage of constructing our model, we add the output layer. This layer is a dense (fully connected) layer consisting of a number of units equal to `'num_classes'`, representing the distinct word categories in our sign language dataset. To generate the probability distribution for each category, we use the softmax activation function. This activation function ensures that the output values are between 0 and 1, making it easier to determine which word category is most likely based on the model's prediction.

With the architecture defined, we proceed to compile the model. For this, we select the Adam optimizer, known for its efficiency and adaptability in deep learning applications. As our loss function, we opt for sparse categorical cross-entropy, which is well-suited

for multi-class classification problems like ours. In addition, we set accuracy as the key metric to monitor the model's performance, providing a clear indication of how well the model is learning from the training data.

Once the model is compiled, we print a summary to inspect the overall architecture. This summary includes crucial information such as the total number of trainable parameters, the configuration of each layer, and the shape of the output at each step. This step serves as a checkpoint to ensure that the model is structured as intended and is ready for training and evaluation. By reviewing the model summary, we can confirm that the architecture aligns with our design goals and is optimized for recognizing sign language gestures with high accuracy.

Model: "sequential"

Layer (type)	Output Shape	Param #
conv1d (Conv1D)	(None, 48, 64)	393,280
gru (GRU)	(None, 48, 64)	24,960
global_average_pooling1d (GlobalAveragePooling1D)	(None, 64)	0
dense (Dense)	(None, 256)	16,640
dense_1 (Dense)	(None, 100)	25,700

Total params: 460,580 (1.76 MB)

Trainable params: 460,580 (1.76 MB)

Non-trainable params: 0 (0.00 B)

Fig 3.2.5.1: Model CNN+GRU

### 3.2.6 TRAINING THE MODEL

Training our sign language recognition model involves a step-by-step process, much like constructing a building where each component serves a unique role in establishing a stable and reliable structure. The foundation of our model begins with the Conv1D layer, which is adept at identifying the basic shapes and contours of sign language gestures. This layer plays a crucial role in interpreting the fundamental features of the signs.

Next, we introduce the GRU (Gated Recurrent Unit) layer, which adds a deeper level of understanding by capturing the sequence of movements in each sign. This is similar to mapping out the framework of a building, providing context and continuity to the visual data. The GlobalAveragePooling1D layer then consolidates all this information, reducing the dimensionality while preserving the key patterns. This makes it easier for the model to process the data and draw accurate conclusions.

Following these foundational steps, we proceed to the Dense layers, where the model refines its understanding of the intricate details in sign language gestures. The Dense layers act like the finishing touches on a building, allowing for more complex decision-making. Finally, the output layer is where the model makes its final prediction, indicating which sign language word or gesture it has recognized.

During the training phase, we divide our dataset into smaller groups for practice and validation, using a technique called Stratified K-Fold Cross-Validation. This method ensures that our model receives a well-balanced training experience by exposing it to various subsets of data. In each training cycle, the model learns from examples of sign language gestures, with additional guidance provided through class weights that reflect the distribution of different signs in the dataset. After each training session, we test the model's ability to recognize unseen signs, measuring its progress and adjusting where necessary.

We repeat this process several times to ensure our model gains enough experience to perform reliably. This iterative approach allows us to gather comprehensive results, enabling us to track the model's accuracy over time. By continuously evaluating and refining our model, we aim to improve its recognition capabilities.



After completing the training and validation process, we assess the model's overall performance. In the final round of validation, our model correctly identified approximately 79.80% of the signs. This rate, although impressive, indicates that there is still room for improvement. Across all rounds of testing, the model's overall accuracy stood at around 79.73%. These metrics provide valuable insights into the model's effectiveness in recognizing sign language and help guide future improvements to increase accuracy and reliability.

```
# Calculate average loss and accuracy
avg_loss = total_loss / k_folds
avg_accuracy = total_accuracy / k_folds

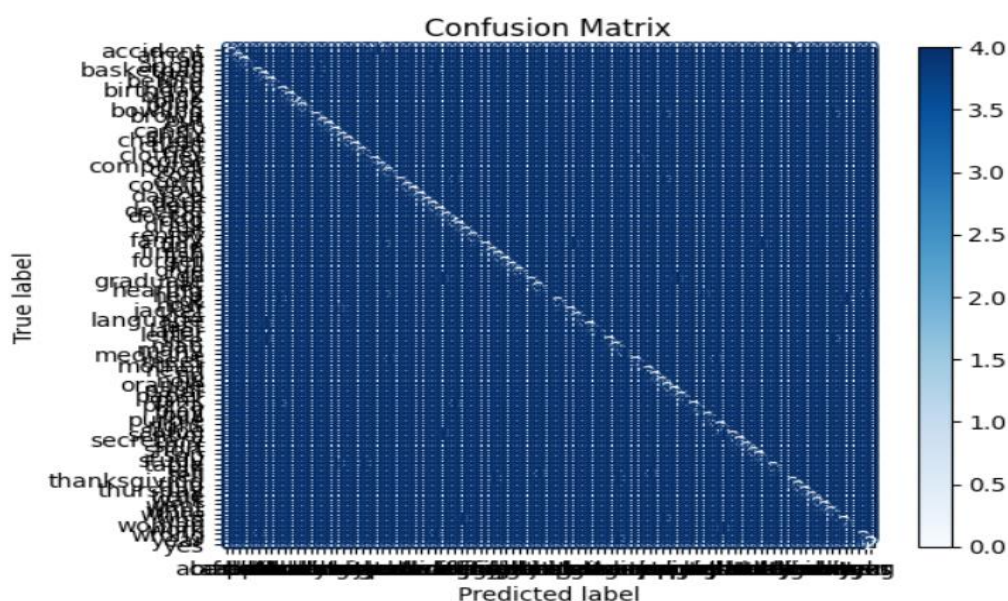
print("Accuracy:", avg_accuracy)

[14]
... Accuracy: 0.7973558485507966

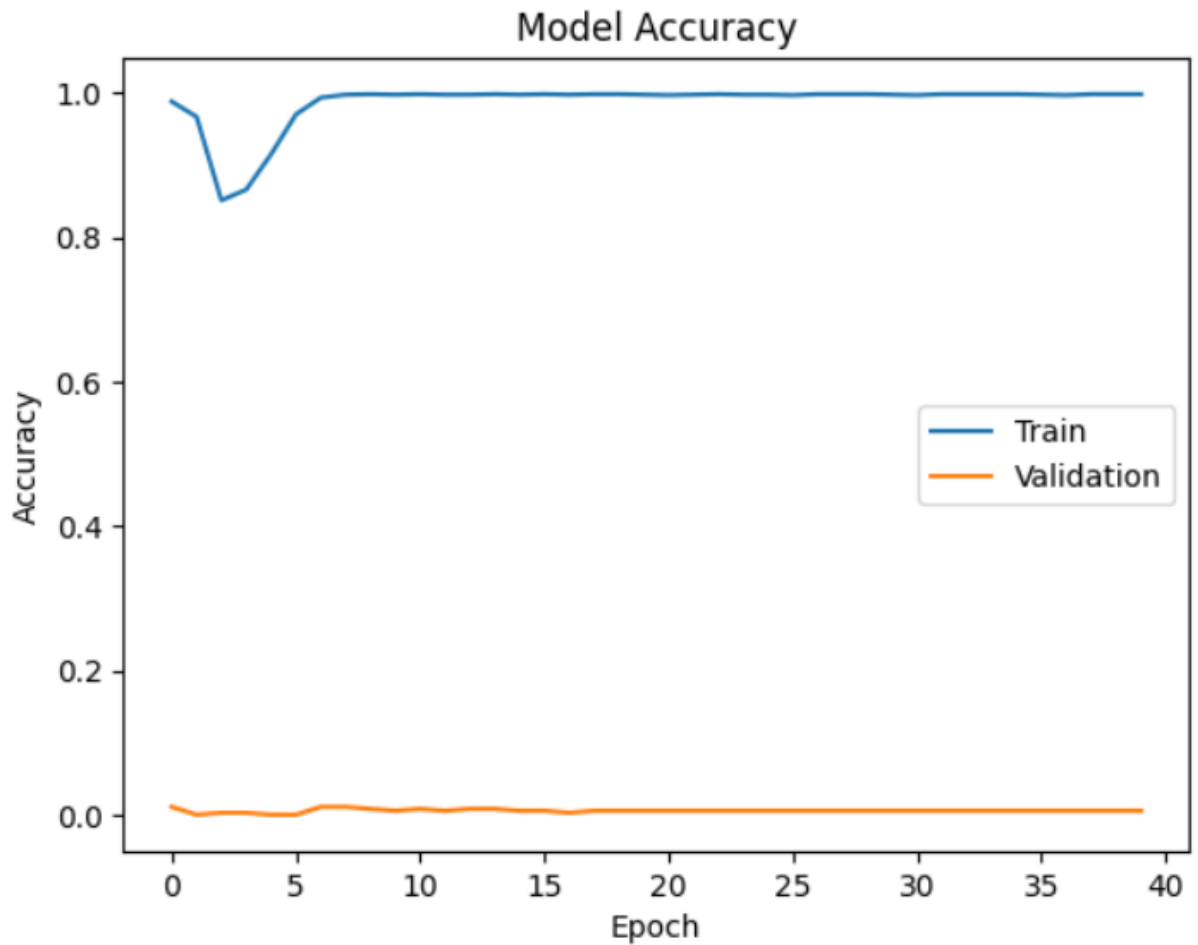
print("\nAverage Loss:", avg_loss)

[15]
...
Average Loss: 9.216709804534911
```

**Fig 3.2.6.1: Accuracy**



**Fig 3.2.6.2 Confusion Matrix Graph**



**Fig 3.2.6..3 Accuracy vs Epoch Graph**

### 3.2.7 TESTING THE MODEL

The process of testing our sign language model was like assembling a puzzle, with each step being crucial to the overall outcome. We began by preparing the video data, ensuring that every frame was properly sized and in the correct format. This preparation was essential to create a consistent starting point for our analysis.

Once the video data was prepped, we utilized ResNet50, a sophisticated deep learning architecture, to analyze each frame. ResNet50 helped us extract critical features and patterns, allowing us to better understand the visual elements within the video sequences. It was like having an expert take a detailed look at every frame, picking out the key details that would help us recognize specific signs.

With the data processed and essential information gathered, we moved on to feeding these insights into our sign language recognition model. This step was akin to presenting a series of clues to our model, challenging it to decipher the correct sign from the visual cues. The model's job was to make sense of the information and guess which sign was being represented in the video.

After the model made its predictions, we evaluated its accuracy by comparing its guesses to the correct signs. This entire process gave us a clear picture of how well our model could recognize and understand sign language. It was as if we were conducting a test to determine if the model was learning effectively and improving its ability to distinguish between different signs.

Through this approach, we were able to assess the model's performance, identify areas for improvement, and gain insights into the intricacies of sign language recognition. The journey from video preprocessing to model evaluation helped us fine-tune our system, ensuring it could accurately interpret sign language and contribute to more accessible communication.

# **CHAPTER 4 - EXPERIMENTAL ANALYSIS AND RESULTS**

## **4.1 SYSTEM CONFIGURATION**

### **4.1.1 Software Requirements**

#### 1. Software:

Python version 3 or above

#### 2. Operating System:

Windows10 or higher, Mac, Linux

#### 3. Tools:

Kaggle, Google Colab, Jupyter Notebook, JupyterLab (Any of these)

#### 4. Python Libraries:

Numpy, OpenCV, Keras, ResNet50, TensorFlow, SkLearn, Matplotlib

#### **4.1.1.1 Introduction to Python**

Python, a high-level and interpreted programming language, stands out for its simplicity and adaptability. It was created by Guido van Rossum and has been in use since 1991. Over the years, it has become one of the most widely used and beloved programming languages worldwide. This widespread adoption is due to its user-friendly design philosophy, which places a strong emphasis on code readability and a straightforward syntax.

One of the most distinctive aspects of Python is its use of whitespace for code indentation, which replaces the conventional curly braces or specific keywords used by many other languages. This unique feature fosters cleaner, more readable code, appealing to both novice and seasoned programmers. Additionally, Python's versatility allows it to support various programming paradigms, including procedural, object-oriented, and functional programming.

Python's comprehensive standard library is another reason for its popularity. This library offers a vast collection of modules and packages that address a multitude of programming needs, such as file input/output, network communication, and web

development. This broad array of built-in functionalities allows developers to tackle various projects without having to rely heavily on external libraries.

Python's adaptability is further demonstrated by the diversity of its applications. It is a common choice for web development, with frameworks like Django and Flask offering robust solutions for building websites and web applications. In the realm of scientific computing, Python excels thanks to libraries such as NumPy and SciPy, which provide tools for complex mathematical calculations and data analysis. Furthermore, Python plays a significant role in data science, with popular tools like Pandas for data manipulation and Matplotlib for data visualization.

The language is also a key player in the fields of artificial intelligence and machine learning, with TensorFlow and PyTorch serving as leading libraries for developing complex AI models. Python's flexibility extends to automation tasks, where modules like BeautifulSoup and Selenium allow for web scraping and automated browser interactions.

Python's success is greatly attributed to its vibrant community support and wealth of documentation. This extensive network of developers and resources ensures that anyone learning or working with Python can find guidance, tutorials, and libraries to aid their projects. With such a strong ecosystem backing it, Python continues to be a top choice for a wide range of programming tasks, from simple automation scripts to large-scale software development projects. Its straightforwardness, flexibility, and strong support structure make it a versatile and accessible language for virtually any application.

#### **4.1.1.2 Introduction to Kaggle**

Kaggle is a well-known online platform that has become a central hub for data science and machine learning enthusiasts. Established in 2010 and later acquired by Google in 2017, Kaggle has built a reputation for hosting a diverse range of competitions that challenge participants to solve real-world problems by creating innovative predictive models. These competitions attract a global community of data scientists, engineers, and statisticians, who compete for prizes, recognition, and opportunities to showcase their skills.

In addition to competitions, Kaggle provides a vast library of datasets covering a multitude of domains, including finance, healthcare, social sciences, technology, and

more. These datasets serve as invaluable resources for learning, practicing, and advancing one's expertise in data analysis and machine learning. Users can explore a wide array of data, allowing them to experiment with different techniques, build models, and test their hypotheses in a collaborative environment.

A unique feature of Kaggle is its integrated development environment (IDE), known as Kaggle Notebooks. This browser-based tool allows users to write, execute, and share code without the hassle of setting up a local development environment. Kaggle Notebooks support multiple programming languages, most notably Python and R, and come pre-equipped with popular data science libraries like Pandas, NumPy, TensorFlow, and scikit-learn. This eliminates the need for extensive software installations and provides a seamless coding experience.

Kaggle Notebooks also foster a sense of community by enabling users to collaborate on projects, share insights, and learn from each other's work. The platform's collaborative nature is further enhanced by its public forums, where users can discuss challenges, share code snippets, and seek advice from peers. This collaborative spirit has made Kaggle a popular destination for both novice data scientists who are eager to learn and seasoned professionals who wish to showcase their expertise and connect with like-minded individuals.

Moreover, Kaggle offers a wealth of educational content, including tutorials, articles, and courses designed to help users improve their data science skills. This educational aspect, combined with the vibrant community and extensive dataset repository, makes Kaggle an invaluable resource for anyone interested in data science, machine learning, and artificial intelligence. Whether you're a beginner looking to kickstart your journey or an experienced data scientist aiming to compete in high-stakes competitions, Kaggle has something to offer everyone.

#### **4.1.1.3 Introduction to Google Colab**

Google Colab, short for Google Colaboratory, is a free, cloud-based platform provided by Google that allows users to write and execute Python code in a browser environment. Launched by Google Research in 2017, Google Colab provides a Jupyter Notebook-like interface and integrates seamlessly with Google Drive.

Here are some key features and highlights of Google Colab:

**Jupyter Notebook Compatibility:** Google Colab supports Jupyter Notebooks, allowing users to create, edit, and run code cells in a sequential manner. This makes it easy to write and experiment with Python code, visualize data, and document the analysis process.

**Free GPU and TPU Support:** One of the most significant advantages of Google Colab is the availability of free access to Graphics Processing Units (GPUs) and Tensor Processing Units (TPUs). This makes it an attractive platform for training machine learning models, as GPUs and TPUs significantly accelerate computation, especially for tasks involving deep learning.

**Integration with Google Drive:** Google Colab seamlessly integrates with Google Drive, allowing users to save and access their notebooks directly from Google Drive. This simplifies collaboration and ensures that work is automatically saved and synced across devices.

**Pre-installed Libraries:** Google Colab comes with many popular Python libraries pre-installed, including NumPy, Pandas, Matplotlib, TensorFlow, and scikit-learn. This eliminates the need for users to install these libraries manually and enables them to start coding right away.

**Shareability and Collaboration:** Users can easily share their Google Colab notebooks with others by generating shareable links or collaborating in real-time, similar to Google Docs. This fosters collaboration and knowledge sharing within the data science and machine learning community.

Overall, Google Colab provides a convenient and powerful environment for data analysis, machine learning experimentation, and collaboration—all accessible through a web browser without the need for local installation or powerful hardware.

#### **4.1.1.4 Introduction to Jupyter Notebook and Jupyter Lab**

Jupyter Notebook and JupyterLab are interactive computing environments that allow users to create and share documents containing live code, equations, visualizations, and narrative text. Developed as part of the Jupyter Project, these tools have gained widespread popularity among data scientists, researchers, educators, and professionals in various fields.

## **Jupyter Notebook**

1. Jupyter Notebook is a web-based interactive computing environment that enables users to create and share documents, called notebooks, containing live code, equations, visualizations, and narrative text.
2. It supports multiple programming languages, including Python, R, Julia, and others, although Python is the most commonly used language.
3. Notebooks are organized into cells, which can contain code, markdown text, or raw text. Users can execute code cells interactively and see the results inline, making it ideal for data exploration and analysis.
4. Jupyter Notebook has a rich ecosystem of extensions and plugins, allowing users to customize their environment with additional features and functionalities.
5. It has become a popular tool for data science, machine learning, scientific computing, and education due to its ease of use, flexibility, and reproducibility.

## **JupyterLab**

1. JupyterLab is the next-generation web-based interactive development environment (IDE) built on top of Jupyter Notebook.
2. It offers a more versatile and powerful interface with support for multiple activities in a single window, including notebooks, text editors, terminals, file browsers, and more.
3. JupyterLab provides a flexible layout system, allowing users to arrange and customize their workspace according to their preferences.
4. It supports interactive widgets, code consoles, and advanced debugging capabilities, enhancing the productivity and workflow of users.
5. JupyterLab retains all the features of Jupyter Notebook and extends them with additional functionalities, making it a comprehensive and integrated environment for interactive computing.

In summary, Jupyter Notebook and JupyterLab are highly versatile platforms that support interactive computing, data analysis, and collaborative work. They offer an intuitive and flexible environment where users can create, share, and explore computational documents with ease. Whether you're a data scientist delving into complex analyses or an educator creating engaging content, these platforms provide a powerful toolkit to meet your needs.



#### **4.1.1.5 Python Libraries**

Python libraries are pre-written code packages that serve as building blocks for various tasks, ranging from data analysis and scientific computing to web development and machine learning. These libraries simplify complex operations by providing ready-made functions, modules, and classes that can be used in Python programs. With thousands of libraries available, Python's ecosystem offers a wide array of tools designed to meet the needs of developers, data scientists, researchers, and engineers.

One of the most popular types of Python libraries is those dedicated to data analysis and scientific computing. Libraries like NumPy, SciPy, and Pandas provide comprehensive support for numerical computations, data manipulation, and statistical analysis. These libraries enable users to process large datasets, perform complex mathematical operations, and analyze data with ease and efficiency. They form the backbone of many data-driven projects and are essential in scientific research and engineering applications.

#### **NumPy**

NumPy, an abbreviation for Numerical Python, is an essential library for numerical computation in Python. It specializes in handling large-scale, multi-dimensional arrays and matrices, allowing developers to perform complex mathematical operations with ease and efficiency. NumPy is a critical component within the Python scientific computing landscape, providing a robust foundation for an array of operations, from simple arithmetic to intricate linear algebra computations.

The flexibility of NumPy's data structures and the speed of its operations have made it indispensable for a wide range of applications, including data analysis, machine learning, artificial intelligence, and engineering. Its array operations are optimized for performance, enabling researchers and developers to process and manipulate data quickly and efficiently.

Moreover, NumPy offers a vast collection of mathematical functions, which can be applied to its arrays to perform transformations, statistical analyses, and other computations. This extensive functionality has cemented NumPy's status as a foundational library, paving the way for more advanced frameworks such as SciPy, Pandas, and scikit-learn, which build upon its capabilities to address specific needs in scientific research, data manipulation, and machine learning. Whether you're conducting

simple data manipulations or engaging in complex machine learning projects, NumPy's versatility and performance make it an invaluable tool in the Python ecosystem.

Here are some key features and highlights of NumPy:

1. **Multi-dimensional Arrays:** NumPy's primary data structure is the `ndarray` (n-dimensional array), which allows efficient storage and manipulation of homogeneous data. These arrays can have any number of dimensions and support various data types.
2. **Mathematical Functions:** NumPy provides a comprehensive set of mathematical functions for performing operations on arrays, such as trigonometric functions, statistical functions, linear algebra operations, Fourier transforms, and more. These functions are optimized for performance and can operate on entire arrays without the need for explicit looping constructs.
3. **Broadcasting:** NumPy's broadcasting mechanism enables arithmetic operations between arrays of different shapes and sizes. It automatically aligns arrays' dimensions and applies the operation element-wise, facilitating concise and efficient code writing.
4. **Vectorized Operations:** NumPy encourages vectorized operations, where mathematical operations are applied to entire arrays instead of individual elements. This approach leverages the underlying C implementation of NumPy, leading to faster execution times compared to equivalent Python loops.
5. **Integration with Other Libraries:** NumPy seamlessly integrates with other Python libraries and frameworks, such as SciPy (scientific computing), Pandas (data manipulation), Matplotlib (data visualization), and scikit-learn (machine learning), forming a powerful ecosystem for data analysis and computation.
6. **Efficiency and Performance:** NumPy's array operations are implemented in highly optimized C and Fortran code, making them significantly faster than equivalent Python implementations. This efficiency is crucial for handling large datasets and performing complex computations efficiently.

Overall, NumPy's rich functionality, efficient array operations, and seamless integration with other libraries make it an essential tool for numerical computing and data analysis in Python. Its versatility and performance have contributed to its widespread adoption and popularity among the scientific and engineering communities.

## OpenCV

OpenCV, short for Open Source Computer Vision Library, is an open-source computer vision and machine learning software library originally developed by Intel in 1999. It is designed to provide a comprehensive infrastructure for computer vision applications and has since become one of the most widely used libraries in the field.

Here are some key features and highlights of OpenCV:

1. **Rich Functionality:** OpenCV offers a vast array of functionalities for image and video processing, including image manipulation, feature detection, object recognition, face detection, optical character recognition (OCR), and more. It provides a rich set of functions and algorithms to perform various tasks in computer vision and machine learning.
2. **Cross-Platform:** OpenCV is cross-platform and supports multiple operating systems, including Windows, Linux, macOS, Android, and iOS. This allows developers to write code once and deploy it across different platforms seamlessly.
3. **Wide Language Support:** OpenCV is primarily written in C++ and provides interfaces for programming in C++, Python, Java, and MATLAB/Octave. This enables developers to use their preferred programming language while leveraging the capabilities of OpenCV.
4. **Performance Optimization:** OpenCV is optimized for performance and efficiency, with many of its core functions implemented in highly optimized C/C++ code. It also supports hardware acceleration through technologies like OpenCL and CUDA, enabling faster execution on GPUs and parallel processing units.
5. **Community and Support:** OpenCV has a large and active community of developers, researchers, and users who contribute to its development, provide support, and share resources and knowledge. This vibrant community ensures that OpenCV remains up-to-date, well-documented, and supported across various platforms and environments.
6. **Integration with Other Libraries:** OpenCV integrates seamlessly with other popular libraries and frameworks in the Python ecosystem, such as NumPy, SciPy, and Matplotlib. This interoperability allows developers to combine the

capabilities of OpenCV with other tools for data analysis, visualization, and machine learning.

Overall, OpenCV is a powerful and versatile library for computer vision applications, offering a wide range of functionalities, cross-platform support, performance optimization, and strong community support. It is widely used in various fields, including robotics, augmented reality, surveillance, medical imaging, and automotive safety, among others

## **Keras**

Keras is an open-source deep learning library written in Python that provides a high-level interface for building and training neural networks. It was developed with a focus on user-friendliness, modularity, and extensibility, allowing both beginners and experienced researchers to quickly prototype and experiment with deep learning models.

Here's an introduction to Keras and some of its modeling techniques, along with a specific focus on the ResNet50 architecture:

Keras Features:

1. **Simplicity:** Keras offers a simple and intuitive API for building neural networks, allowing users to create complex models with minimal code.
2. **Modularity:** Keras models are built using modular building blocks called layers, which can be easily combined to create complex architectures. This modularity facilitates experimentation and model customization.
3. **Flexibility:** Keras supports both sequential and functional model building paradigms. Sequential models are linear stacks of layers, while functional models allow for more complex architectures with multiple input and output connections.
4. **Extensibility:** Keras is highly extensible, allowing users to define custom layers, loss functions, and metrics. It also provides integration with external libraries and frameworks, such as TensorFlow and Theano, for advanced customization.

Modeling Techniques in Keras:

1. **Sequential Model:** The sequential model is the simplest type of model in Keras, consisting of a linear stack of layers. It is suitable for building feedforward neural networks where the data flows sequentially through each layer.

2. **Functional API:** The functional API allows for more flexible model architectures, including multi-input and multi-output networks, shared layers, and residual connections. It is suitable for building complex architectures such as multi-modal networks and neural network ensembles.
3. **Transfer Learning:** Keras supports transfer learning, a technique where pre-trained models are used as a starting point for training new models on tasks with limited data. This allows users to leverage the knowledge learned by pre-trained models on large datasets for tasks with smaller datasets.
4. **Custom Layers and Models:** Keras allows users to define custom layers and models, enabling the implementation of novel architectures and algorithms. Custom layers can be created by subclassing the Layer class, while custom models can be built by subclassing the Model class.

## **ResNet50**

ResNet50 is a convolutional neural network architecture that belongs to the ResNet (Residual Network) family of models. It was proposed by Microsoft Research in their paper "Deep Residual Learning for Image Recognition" in 2015. ResNet50 is renowned for its depth, featuring 50 layers, and its effectiveness in tackling the vanishing gradient problem commonly encountered in very deep neural networks.

Here's an introduction to ResNet50 and its processing with the ImageNet dataset [12]:

ResNet50 Architecture:

1. **Deep Architecture:** ResNet50 consists of 50 layers, making it significantly deeper than earlier convolutional neural network architectures. This depth enables ResNet50 to capture complex features and patterns from input images effectively.
2. **Residual Connections:** One of the key innovations of ResNet50 is the introduction of residual connections, also known as skip connections. These connections enable the network to learn residual mappings, which helps alleviate the vanishing gradient problem during training. Residual connections pass the input directly to deeper layers, allowing the network to focus on learning the residual information.
3. **Building Blocks:** ResNet50 is composed of multiple building blocks called residual blocks. Each residual block contains several convolutional layers followed by batch normalization and ReLU activation, along with shortcut connections. The shortcut

connections bypass one or more convolutional layers and add the input to the output of the block.

Processing with ImageNet:

1. **ImageNet Dataset:** The ImageNet dataset is a large-scale dataset containing millions of labeled images across thousands of categories. It is widely used for training and evaluating deep learning models for image classification and object recognition tasks.
2. **Pre-Trained Model:** ResNet50 is often used as a pre-trained model for image classification tasks with the ImageNet dataset. Pre-trained models are neural networks that have been trained on a large dataset (e.g., ImageNet) and then fine-tuned or used as feature extractors for other tasks.
3. **Transfer Learning:** Transfer learning involves using pre-trained models, such as ResNet50 trained on ImageNet, as a starting point for training models on specific tasks or datasets with limited amounts of labeled data. By leveraging the knowledge learned from the large-scale ImageNet dataset, transfer learning can significantly improve the performance of models on new tasks with smaller datasets.
4. **Fine-Tuning:** Fine-tuning involves retraining the pre-trained ResNet50 model on a new dataset (e.g., a subset of ImageNet or a different dataset altogether) to adapt it to the specific characteristics of the new task. Fine-tuning typically involves freezing the early layers of the network (which capture generic features) and only training the later layers (which capture task-specific features) to prevent overfitting and accelerate training.

In summary, ResNet50 is a powerful convolutional neural network architecture known for its depth and effectiveness in image recognition tasks. When processed with the ImageNet dataset [12] , ResNet50 can be used as a pre-trained model for transfer learning and fine-tuning on various image-related tasks, making it a valuable asset in the field of deep learning and computer vision.

## **TensorFlow**

TensorFlow is an open-source machine learning framework developed by Google Brain and widely used for building and training deep learning models. It provides a comprehensive ecosystem of tools, libraries, and resources for developing and

deploying machine learning applications across a variety of domains, including image and speech recognition, natural language processing, and recommendation systems.

Here's an introduction to TensorFlow and some of its key features:

1. **Simplicity:** Keras provides a high-level, user-friendly API for building and training neural networks. It abstracts away much of the complexity of TensorFlow, allowing you to focus on the design and experimentation of your models rather than low-level implementation details.
2. **Modularity:** Keras follows a modular design, allowing you to easily assemble neural network architectures from pre-built layers. You can stack layers sequentially or create more complex network architectures using the functional API, which enables building models with multiple input or output layers, shared layers, and branching architectures.
3. **Extensibility:** While Keras offers a wide range of built-in layers and functionalities, it also allows for customization and extension. You can create custom layers, loss functions, metrics, and callbacks, enabling you to tailor your models to specific tasks and requirements.
4. **Integration with TensorFlow:** Keras is tightly integrated with TensorFlow, making it easy to leverage TensorFlow's powerful backend for computations. You can seamlessly combine Keras with TensorFlow's low-level functionalities, such as custom operations and automatic differentiation, for advanced model building and optimization.
5. **Visualization and Debugging:** TensorFlow's TensorBoard can be used to visualize and monitor the training process of Keras models. You can track metrics, visualize model graphs, and analyze performance, facilitating model debugging and optimization.
6. **Performance:** TensorFlow's optimized backend ensures efficient execution of Keras models, whether running on CPUs, GPUs, or TPUs. This allows for fast training and inference, even with large-scale deep learning models and datasets.
7. **Transfer Learning and Pre-trained Models:** Keras provides easy-to-use APIs for transfer learning, allowing you to leverage pre-trained models and weights from TensorFlow Hub or other sources. You can fine-tune pre-trained models on your own datasets or use them as feature extractors for downstream tasks.

Overall, the combination of Keras and TensorFlow offers a powerful and user-friendly platform for deep learning projects. Whether you're a beginner experimenting with neural networks or an experienced researcher developing complex architectures, Keras with TensorFlow provides the flexibility, simplicity, and performance needed to build and deploy state-of-the-art deep learning models.

## **SkLearn**

Scikit-learn, commonly referred to as sklearn, is an open-source machine learning library for Python. It provides simple and efficient tools for data mining and data analysis, built on top of other Python libraries such as NumPy, SciPy, and matplotlib.

Here's an introduction to scikit-learn and some of its key features:

### **Key Features:**

1. **Simple and Consistent API:** Scikit-learn offers a uniform and easy-to-use API, making it accessible to both beginners and experienced machine learning practitioners. The library follows a consistent interface for various algorithms, making it easy to switch between different models and techniques.
2. **Comprehensive Set of Algorithms:** Scikit-learn includes a wide range of machine learning algorithms for classification, regression, clustering, dimensionality reduction, and model selection. These algorithms cover both supervised and unsupervised learning techniques, providing solutions for various tasks and problem domains.
3. **Model Evaluation and Selection:** Scikit-learn provides tools for model evaluation, including metrics for measuring classification accuracy, regression performance, and clustering quality. It also offers techniques for hyperparameter tuning, cross-validation, and model selection to optimize model performance and generalization.
4. **Preprocessing and Feature Engineering:** Scikit-learn includes utilities for data preprocessing and feature engineering, such as data normalization, feature scaling, categorical encoding, and missing value imputation. These tools help prepare raw data for training machine learning models and improve their performance.
5. **Integration with Other Libraries:** Scikit-learn seamlessly integrates with other Python libraries and frameworks, such as NumPy, SciPy, pandas, and matplotlib. This integration enables users to leverage the rich functionality of these libraries in conjunction with scikit-learn for data manipulation, visualization, and analysis.



6. **Scalability and Efficiency:** Scikit-learn is designed for scalability and efficiency, with optimized implementations of algorithms for large-scale datasets. It supports sparse matrices and provides parallelized implementations for efficient computation on multi-core processors.
7. **Community and Documentation:** Scikit-learn has a large and active community of developers, researchers, and users who contribute to its development and maintenance. The library is well-documented, with comprehensive documentation, tutorials, and examples to help users get started and master its features.

Overall, scikit-learn is a versatile and user-friendly machine-learning library that provides a powerful set of tools for building, evaluating, and deploying machine-learning models. Whether you're a beginner exploring machine learning concepts or an experienced practitioner developing production-ready solutions, scikit-learn offers the tools and resources needed to tackle a wide range of machine learning tasks.

## **Matplotlib**

Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python. It was created by John D. Hunter in 2003 and has since become one of the most widely used plotting libraries in the Python ecosystem.

Here's an introduction to Matplotlib and some of its key features:

### **Key Features:**

1. **Versatile Plotting:** Matplotlib offers a wide variety of plotting functions and styles for creating a diverse range of plots, including line plots, scatter plots, bar plots, histograms, heatmaps, and more. It provides flexibility in customization and allows users to create highly customizable and publication-quality visualizations.
2. **Support for Multiple Output Formats:** Matplotlib supports various output formats, including PNG, PDF, SVG, and interactive formats like HTML and notebook widgets. This versatility allows users to generate plots for different purposes and seamlessly integrate them into documents, presentations, and web applications.
3. **Integration with Jupyter Notebooks:** Matplotlib integrates seamlessly with Jupyter notebooks, allowing users to create, display, and interact with plots directly within the notebook environment. This makes it a popular choice for data exploration, analysis, and visualization in the Jupyter ecosystem.

4. **Customization and Styling:** Matplotlib provides extensive support for customization and styling of plots, including control over colors, linestyles, markers, fonts, labels, axes, and annotations. Users can customize every aspect of their plots to match their specific requirements and preferences.
5. **Multi-platform Support:** Matplotlib is cross-platform and works on various operating systems, including Windows, macOS, and Linux. It provides consistent behavior and appearance across different platforms, ensuring that plots look consistent regardless of the environment.
6. **Interactivity:** Matplotlib offers limited interactivity features for exploring and interacting with plots, such as zooming, panning, and mouse hover tooltips. While it's primarily designed for static visualizations, Matplotlib can be combined with other libraries like `mplcursors` or `mpl-interactions` to add additional interactivity.
7. **Extensibility and Integration:** Matplotlib is highly extensible and integrates seamlessly with other Python libraries and frameworks, such as NumPy, Pandas, SciPy, and Seaborn. It can be easily combined with these libraries to create complex and sophisticated visualizations for data analysis and scientific computing.

Overall, Matplotlib is a powerful and flexible plotting library that provides a comprehensive set of tools for creating high-quality visualizations in Python. Whether you're a beginner exploring data or an experienced data scientist creating publication-ready plots, Matplotlib offers the tools and flexibility needed to visualize data effectively.

#### **4.1.2 Hardware Requirements**

To ensure smooth operation and optimal performance of our application, certain hardware requirements must be met. These specifications guarantee that the system can handle data processing and maintain efficient execution of tasks, especially those involving computer vision, machine learning, and data-intensive operations. The following hardware components are recommended:

1. **Camera:**

A high-quality camera is crucial for capturing clear images and videos. We recommend a camera with at least a 3-megapixel (MP) sensor. This resolution is adequate for detailed image analysis and ensures that the data used for training and recognition is of sufficient quality.

## 2. RAM:

Random Access Memory (RAM) plays a significant role in multitasking and handling large data sets. A minimum of 8GB of RAM is required to run the application smoothly, especially when processing video data or running machine learning algorithms.

## 3. GPU:

A dedicated Graphics Processing Unit (GPU) with at least 4GB of memory is recommended. The GPU accelerates tasks like image processing and deep learning, enabling faster computations and smoother operation.

## 4. Processor:

The central processing unit (CPU) should be capable of handling multiple processes efficiently. A processor with a minimum clock speed of 850 MHz is the baseline requirement, but for optimal performance, particularly with multi-threaded tasks, a multi-core processor with higher clock speeds is recommended. Processors with four or more cores can significantly improve the speed and responsiveness of the application.

## 5. Storage:

Adequate storage is crucial for storing application data, models, and other assets. A hard disk drive (HDD) with a capacity of at least 20GB is required, but solid-state drives (SSDs) are preferred for their faster read/write speeds and improved reliability.

Meeting these hardware requirements ensures that the system runs efficiently and can handle the demands of the application without performance bottlenecks or interruptions. This foundation supports not only the current application but also future scalability and expansion.

## 4.2 SAMPLE CODE

### Importing and Preprocessing dataset:

```
import os
import cv2

# Function to load and preprocess the videos
```

```

def load_and_preprocess_dataset(dataset_folder, target_size=(224, 224),
max_frames=50):

    X = [] # List to store preprocessed video frames

    y = [] # List to store labels (word categories)

    for word_folder in os.listdir(dataset_folder): # Iterate through each word folder

        word_path = os.path.join(dataset_folder, word_folder)

        if os.path.isdir(word_path):

            # Iterate through each video file in the word folder

            for video_file in os.listdir(word_path):

                video_path = os.path.join(word_path, video_file)

                # Create a folder to store frames if it doesn't exist

                frames_folder = os.path.join("/kaggle/working/frames", word_folder,
video_file.split('.')[0])

                os.makedirs(frames_folder, exist_ok=True)

                cap = cv2.VideoCapture(video_path) # Load video using OpenCV

                frames = []

                frame_count = 0

                while cap.isOpened():

                    ret, frame = cap.read()

                    if not ret:

                        break

                    # Resize frame to target size and preprocess as needed

                    frame = cv2.resize(frame, target_size)

                    frame = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB) # Convert to RGB

                    frames.append(frame)

                    frame_filename = os.path.join(frames_folder, f"frame_{frame_count}.jpg")

                    cv2.imwrite(frame_filename, frame)

                    frame_count += 1

                if frame_count == max_frames:

                    break

```

```

        cap.release()

    if len(frames) < max_frames:
        frames.extend([frames[-1]] * (max_frames - len(frames)))
    else:
        frames = frames[:max_frames]

    # Stack frames into a single video tensor
    video_tensor = np.stack(frames)

    X.append(video_tensor)

    y.append(word_folder) # Assuming folder names are the labels

X = np.array(X)
y = np.array(y)
return X, y

dataset_folder = "/kaggle/input/signproject-ds/dataset/dataset"
X, y = load_and_preprocess_dataset(dataset_folder)

print("Dataset loaded.")

print("Shape of input data (videos):", X.shape)
print("Shape of output data (labels):", y.shape)

np.save("/kaggle/working/X.npy", X)
np.save("/kaggle/working/y.npy", y)

Feature Extraction:

import tensorflow as tf

from tensorflow.keras.applications import ResNet50
from tensorflow.keras.applications.resnet50 import preprocess_input

# Load the pre-trained ResNet50 model (excluding the top classification layer)
base_model = ResNet50(weights='imagenet', include_top=False, input_shape=(224,
224, 3))

base_model.trainable = False # Freeze the weights of the pre-trained layers

# Function to extract features from video frames using ResNet50
def extract_features(video_frames):

```

```

# Preprocess input frames for ResNet50 model
frames_preprocessed = preprocess_input(video_frames)
features = base_model.predict(frames_preprocessed) # Extract features
# Pool features across frames (e.g., average pooling)
pooled_features = tf.reduce_mean(features, axis=(1, 2))
return pooled_features

# Function to extract features from all videos in the dataset
def extract_features_from_dataset(X):
    features_list = []
    total_videos = len(X)
    for i, video_frames in enumerate(X):
        video_features = extract_features(video_frames)
        features_list.append(video_features)
        progress = (i + 1) / total_videos * 100
        print(f'Extracting features: {progress:.2f}% complete', end='\r', flush=True)
    return np.array(features_list)

X_features = extract_features_from_dataset(X)
print("\nFeatures extracted.")
print("Shape of extracted features:", X_features.shape)

```

### **Model Building:**

```

from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Conv1D, GRU, Dense, GlobalAveragePooling1D
# Define the model architecture
model = Sequential([Conv1D(filters=64, kernel_size=3, activation='relu',
input_shape=(X_features.shape[1], X_features.shape[2])),
    GRU(64, return_sequences=True),
    GlobalAveragePooling1D(),
    Dense(128, activation='relu'),
    Dense(num_classes, activation='softmax') ])

```

```

# Compile the model and Print the model summary
model.compile(optimizer='adam',
              loss='sparse_categorical_crossentropy',
              metrics=['accuracy'])

model.summary()

Training the Model:

from sklearn.model_selection import StratifiedKFold
from sklearn.utils.class_weight import compute_class_weight
from sklearn.preprocessing import LabelEncoder

# Initialize variables to store total loss and accuracy
total_loss = 0
total_accuracy = 0

# Define the number of folds for cross-validation
k_folds = 10

skf = StratifiedKFold(n_splits=k_folds, shuffle=True, random_state=42)

# Initialize LabelEncoder and Encode the string labels to integer representations
label_encoder = LabelEncoder()
y_encoded = label_encoder.fit_transform(y)

# Get the unique class labels and Compute class weights
unique_classes = np.unique(y_encoded)

class_weights = compute_class_weight('balanced', classes=unique_classes,
y=y_encoded)

# Convert class weights to a dictionary
class_weights_dict = dict(zip(unique_classes, class_weights))

# Perform k-fold cross-validation
fold = 1

for train_index, test_index in skf.split(X_features, y_encoded):
    print(f"\nTraining Fold {fold}/{k_folds}")

    X_train, X_test = X_features[train_index], X_features[test_index]

```

```

y_train, y_test = y_encoded[train_index], y_encoded[test_index]

# Train the model with class weights

history = model.fit(X_train, y_train, batch_size=42, epochs=19,
validation_split=0.2, class_weight=class_weights_dict)

# Evaluate the model on the test set

loss, accuracy = model.evaluate(X_test, y_test)

print(f'Fold {fold} Test Loss:', loss)

print(f'Fold {fold} Test Accuracy:', accuracy)

total_loss += loss

total_accuracy += accuracy

fold += 1

avg_loss = total_loss / k_folds # Calculate average loss

avg_accuracy = total_accuracy / k_folds # Calculate average accuracy

print("Accuracy:", avg_accuracy)

```

### **Plotting Confusion Matrix and Accuracy Graph:**

```

import matplotlib.pyplot as plt

from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay

def plot_conf_matrix(y_true, y_pred, classes): # Plot confusion matrix

    cm = confusion_matrix(y_true, y_pred)

    disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=classes)

    disp.plot(cmap=plt.cm.Blues)

    plt.title('Confusion Matrix')

    plt.show()

def plot_accuracy(history): # Plot accuracy graph

    plt.plot(history.history['accuracy'])

    plt.plot(history.history['val_accuracy'])

    plt.title('Model Accuracy')

    plt.ylabel('Accuracy')

    plt.xlabel('Epoch')

```



```

plt.legend(['Train', 'Validation'], loc='lower left')

plt.show()

# Obtain predictions and Plot confusion matrix , accuracy graph

y_pred = model.predict(X_test)

plot_conf_matrix(y_test, np.argmax(y_pred, axis=1), classes=label_encoder.classes_)

plot_accuracy(history)

```

### **Testing the Model:**

```

import os

import cv2

import tensorflow as tf

from tensorflow.keras.models import load_model

from tensorflow.keras.applications.resnet50 import preprocess_input

# Define the function to preprocess a single video

def preprocess_video(video_path, target_size=(224, 224), max_frames=50):

    frames = []

    cap = cv2.VideoCapture(video_path) # Load video using OpenCV

    frame_count = 0

    while cap.isOpened():

        ret, frame = cap.read()

        if not ret:

            break

        frame = cv2.resize(frame, target_size) # Resize frame to target size

        frame = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB) # Convert to RGB

        frames.append(frame)

        frame_count += 1

        if frame_count == max_frames:

            break

    cap.release()

    if len(frames) < max_frames: # Pad or truncate frames to ensure fixed length

```

```

        frames.extend([frames[-1]] * (max_frames - len(frames)))
    else:
        frames = frames[:max_frames]

    video_tensor = np.stack(frames)  # Stack frames into a single video tensor
    return video_tensor

# Define the function to extract features from video frames using ResNet50
def extract_features(video_frames, base_model):
    # Preprocess input frames for ResNet50 model
    frames_preprocessed = preprocess_input(video_frames)

    features = base_model.predict(frames_preprocessed)  # Extract features

    # Pool features across frames (e.g., average pooling)
    pooled_features = tf.reduce_mean(features, axis=(1, 2))  # Average pooling over
    spatial dimensions

    return pooled_features

# Load the saved model
saved_model_path = '/kaggle/working/sign_language_rec_model.h5'
model = load_model(saved_model_path)

# Load the pre-trained ResNet50 model (excluding the top classification layer)
base_model = tf.keras.applications.ResNet50(weights='imagenet', include_top=False,
input_shape=(224, 224, 3))

# Define the path to your test video and Preprocess the test video
test_video_path = '/kaggle/input/signproject-ds/dataset/dataset/black/06473.mp4'
test_video = preprocess_video(test_video_path)

test_features = extract_features(test_video, base_model) # Extract features from the
test video

# Reshape the features to match the input shape expected by the model
test_features = np.expand_dims(test_features, axis=0)

# Perform inference using the model and Get the predicted class index
predictions = model.predict(test_features)

predicted_class_index = np.argmax(predictions)

```

```
# Define a dictionary mapping class indices to words

class_index_to_word = { 0: 'accident', 1: 'africa', 2: 'all', 3: 'apple', 4: 'basketball', 5:
'bed', 6: 'before', 7: 'bird', 8: 'birthday', 9: 'black', 10: 'blue', 11: 'book', 12: 'bowling',
13: 'brown', 14: 'but', 15: 'can', 16: 'candy', 17: 'chair', 18: 'change', 19: 'cheat', 20:
'city', 21: 'clothes', 22: 'color', 23: 'computer', 24: 'cook', 25: 'cool', 26: 'corn', 27:
'cousin', 28: 'cow', 29: 'dance', 30: 'dark', 31: 'deaf', 32: 'decide', 33: 'doctor', 34: 'dog',
35: 'drink', 36: 'eat', 37: 'enjoy', 38: 'family', 39: 'fine', 40: 'finish', 41: 'fish', 42: 'forget',
43: 'full', 44: 'give', 45: 'go', 46: 'graduate', 47: 'hat', 48: 'hearing', 49: 'help', 50: 'hot',
51: 'how', 52: 'jacket', 53: 'kiss', 54: 'language', 55: 'last', 56: 'later', 57: 'letter', 58:
'like', 59: 'man', 60: 'many', 61: 'medicine', 62: 'meet', 63: 'mother', 64: 'need', 65: 'no',
66: 'now', 67: 'orange', 68: 'paint', 69: 'paper', 70: 'pink', 71: 'pizza', 72: 'play', 73: 'pull',
74: 'purple', 75: 'right', 76: 'same', 77: 'school', 78: 'secretary', 79: 'shirt', 80: 'short', 81:
'son', 82: 'study', 83: 'table', 84: 'tall', 85: 'tell', 86: 'thanksgiving', 87: 'thin', 88:
'thursday', 89: 'time', 90: 'walk', 91: 'want', 92: 'what', 93: 'white', 94: 'who', 95:
'woman', 96: 'work', 97: 'wrong', 98: 'year', 99: 'yes'}
```

```
# Get the predicted word
```

```
predicted_word = class_index_to_word.get(predicted_class_index, 'Unknown')
```

```
print("Predicted word:", predicted_word)
```

### **User interface:**

```
import cv2
```

```
import gradio as gr
```

```
import tensorflow as tf
```

```
from tensorflow.keras.models import load_model
```

```
from tensorflow.keras.applications.resnet50 import preprocess_input
```

```
from gtts import gTTS
```

```
from io import BytesIO
```

```
# Define the function to preprocess a single video
```

```
def preprocess_video(video_path, target_size=(224, 224), max_frames=50):
```

```
    frames = []
```

```
    cap = cv2.VideoCapture(video_path)    # Load video using OpenCV
```

```
    frame_count = 0
```

```
    while cap.isOpened():
```

```
        ret, frame = cap.read()
```

```
        if not ret:
```

```

        break

    # Resize frame to target size and preprocess as needed
    frame = cv2.resize(frame, target_size)

    frame = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB) # Convert to RGB
    frames.append(frame)

    frame_count += 1

    if frame_count == max_frames:
        break

cap.release()

# Pad or truncate frames to ensure fixed length
if len(frames) < max_frames:
    frames.extend([frames[-1]] * (max_frames - len(frames)))
else:
    frames = frames[:max_frames]

video_tensor = np.stack(frames) # Stack frames into a single video tensor
return video_tensor

# Define the function to extract features from video frames using ResNet50
def extract_features(video_frames, base_model):
    # Preprocess input frames for ResNet50 model
    frames_preprocessed = preprocess_input(video_frames)

    # Extract features using ResNet50
    features = base_model.predict(frames_preprocessed)

    # Pool features across frames (e.g., average pooling)
    pooled_features = tf.reduce_mean(features, axis=(1, 2)) # Average pooling over
    spatial dimensions

    return pooled_features

saved_model_path = '/content/sign_language_rec_model(1).h5'
model = load_model(saved_model_path)

# Load the pre-trained ResNet50 model (excluding the top classification layer)

```

```

base_model = tf.keras.applications.ResNet50(weights='imagenet', include_top=False,
input_shape=(224, 224, 3))

# Define the function to predict sign language from webcam frames
def predict_sign_language(frame):
    video_tensor = preprocess_video(frame)    # Preprocess the frame
    video_features = extract_features(video_tensor, base_model)
    video_features = np.expand_dims(video_features, axis=0)    # Reshape the features
    predictions = model.predict(video_features)    # Perform prediction
    predicted_class_index = np.argmax(predictions)    # Get the predicted class index
    # Get the predicted word
    predicted_word = class_index_to_word.get(predicted_class_index, 'Unknown')
    # Convert the predicted word to speech
    tts = gTTS(predicted_word)
    byte_io = BytesIO()
    tts.write_to_fp(byte_io)
    byte_io.seek(0)
    speech_output = byte_io.getvalue()
    return predicted_word, speech_output

# Create a Gradio interface
webcam_input = gr.Video(label="Please make sign language for 5 seconds")
predicted_word_output = gr.Label(label="Predicted Word")
speech_output = gr.Audio(label="Speech")

gr.Interface(predict_sign_language, webcam_input, [predicted_word_output,
speech_output], title="Sign Language Recognition").launch()

```

## Outputs

```
Dataset loaded.  
Shape of input data (videos): (2038, 50, 224, 224, 3)  
Shape of output data (labels): (2038,)
```

**Fig 4.2.1 Shape of loaded dataset**

```
['accident' 'africa' 'all' 'apple' 'basketball' 'bed' 'before' 'bird'  
'birthday' 'black' 'blue' 'book' 'bowling' 'brown' 'but' 'can' 'candy'  
'chair' 'change' 'cheat' 'city' 'clothes' 'color' 'computer' 'cook'  
'cool' 'corn' 'cousin' 'cow' 'dance' 'dark' 'deaf' 'decide' 'doctor'  
'dog' 'drink' 'eat' 'enjoy' 'family' 'fine' 'finish' 'fish' 'forget'  
'full' 'give' 'go' 'graduate' 'hat' 'hearing' 'help' 'hot' 'how' 'jacket'  
'kiss' 'language' 'last' 'later' 'letter' 'like' 'man' 'many' 'medicine'  
'meet' 'mother' 'need' 'no' 'now' 'orange' 'paint' 'paper' 'pink' 'pizza'  
'play' 'pull' 'purple' 'right' 'same' 'school' 'secretary' 'shirt'  
'short' 'son' 'study' 'table' 'tall' 'tell' 'thanksgiving' 'thin'  
'thursday' 'time' 'walk' 'want' 'what' 'white' 'who' 'woman' 'work'  
'wrong' 'year' 'yes']
```

**Fig 4.2.2 Words of WLASL100**

```
2/2 _____ 0s 38ms/step  
2/2 _____ 0s 38ms/step  
2/2 _____ 0s 38ms/step  
2/2 _____ 0s 37ms/step  
2/2 _____ 0s 37ms/step  
2/2 _____ 0s 38ms/step  
2/2 _____ 0s 38ms/step  
2/2 _____ 0s 37ms/step  
2/2 _____ 0s 38ms/step  
2/2 _____ 0s 38ms/step  
2/2 _____ 0s 38ms/step  
2/2 _____ 0s 38ms/step  
2/2 _____ 0s 37ms/step  
2/2 _____ 0s 38ms/step  
2/2 _____ 0s 38ms/step  
Extracting features: 100.00% complete  
Features extracted.  
Shape of extracted features: (2038, 50, 2048)
```

**Fig 4.2.3 Feature Extraction**

```
print(x_features)
```

```
[[[0.00000000e+00 2.63927102e-01 4.53065217e-01 ... 1.20960128e+00
    1.93759009e-01 3.39217991e-01]
 [0.00000000e+00 2.61086226e-01 4.48085308e-01 ... 1.20344710e+00
    1.92992210e-01 3.39572996e-01]
 [0.00000000e+00 2.57667184e-01 4.32862461e-01 ... 1.17658782e+00
    1.93931788e-01 3.45444709e-01]
 ...
 [5.02856195e-01 3.96880746e-01 1.03041328e-01 ... 1.43230200e+00
    4.37320977e-01 4.95182276e-01]
 [5.25891423e-01 4.45999026e-01 5.98202236e-02 ... 1.28450525e+00
    4.15642262e-01 4.84339744e-01]
 [5.28617144e-01 4.30488735e-01 1.03728607e-01 ... 1.47595525e+00
    4.82345372e-01 4.61182684e-01]]]

[[[6.96059406e-01 6.28638566e-01 3.46326157e-02 ... 2.18346858e+00
    1.27872065e-01 7.34374747e-02]
 [7.39025354e-01 4.92727757e-01 2.81873010e-02 ... 2.22387218e+00
    1.49335846e-01 7.76762962e-02]
 [8.07803869e-01 3.97068381e-01 1.17912963e-02 ... 2.26917958e+00
    1.52607262e-01 1.00927673e-01]
 ...
 [1.07739222e+00 0.00000000e+00 0.00000000e+00 ... 2.07771873e+00
    1.12211376e-01 1.01911470e-01]
 [1.07380462e+00 9.47292708e-03 0.00000000e+00 ... 2.11444116e+00
    1.03420958e-01 1.33759350e-01]
 ...
 [2.58258581e-01 7.37199366e-01 0.00000000e+00 ... 5.00693917e-01
    8.35448410e-03 3.33888292e-01]
 [1.87069952e-01 8.67471874e-01 0.00000000e+00 ... 4.49102014e-01
    1.12365140e-02 3.39235723e-01]]]
```

**Fig 4.2.4 Extracted Features**

Model: "sequential"

Layer (type)	Output Shape	Param #
conv1d (Conv1D)	(None, 48, 64)	393,280
gru (GRU)	(None, 48, 64)	24,960
global_average_pooling1d (GlobalAveragePooling1D)	(None, 64)	0
dense (Dense)	(None, 128)	8,320
dense_1 (Dense)	(None, 100)	12,900

Total params: 439,460 (1.68 MB)

Trainable params: 439,460 (1.68 MB)

Non-trainable params: 0 (0.00 B)

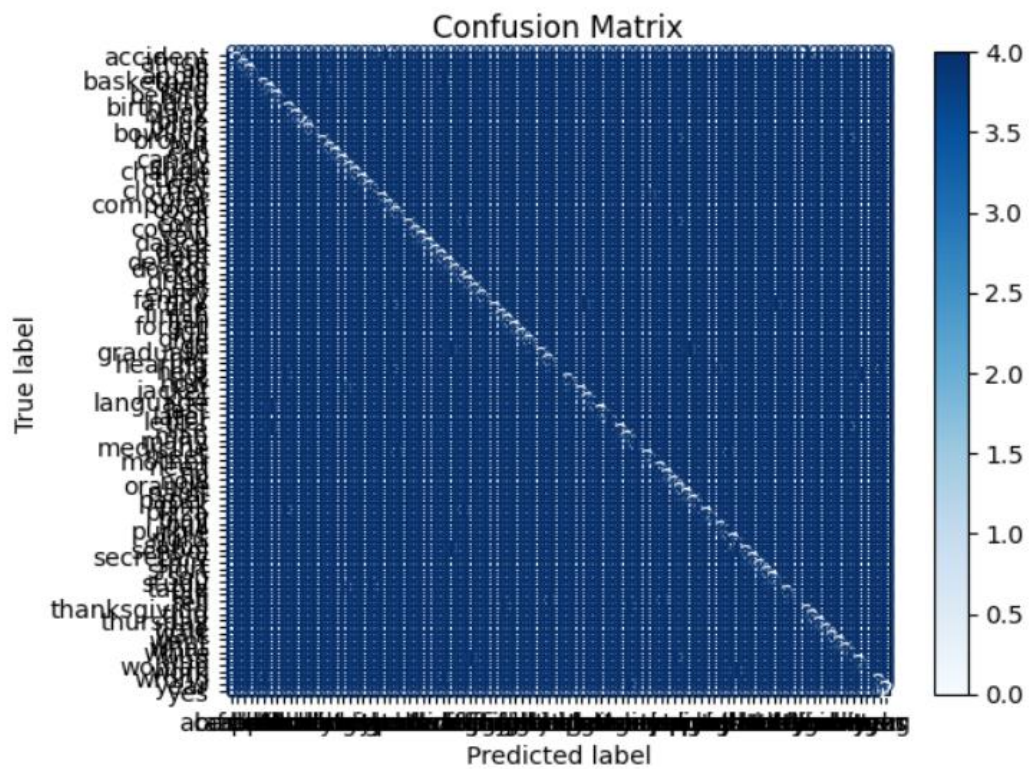
**Fig 4.2.5 Model Summary**

```
print("Accuracy:", avg_accuracy)
```

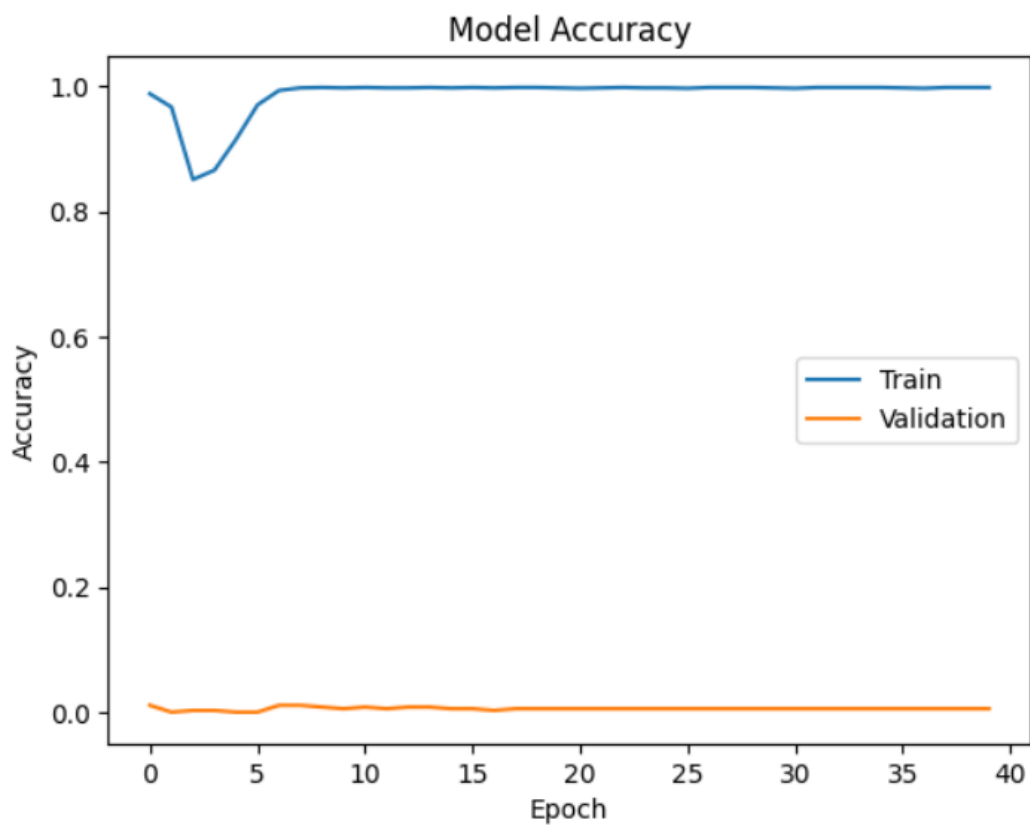
```
Accuracy: 0.7973558485507966
```

**Fig 4.2.6 Test Accuracy**

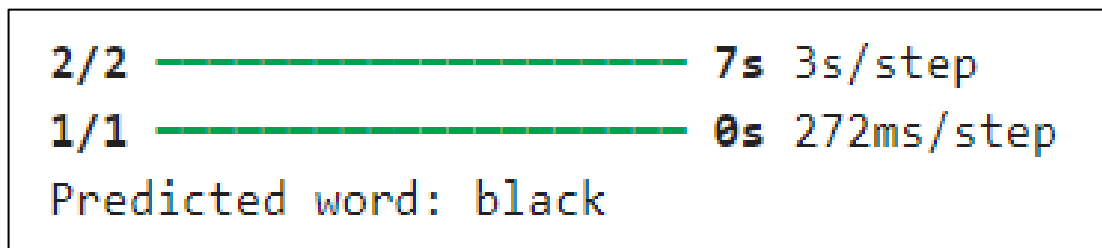




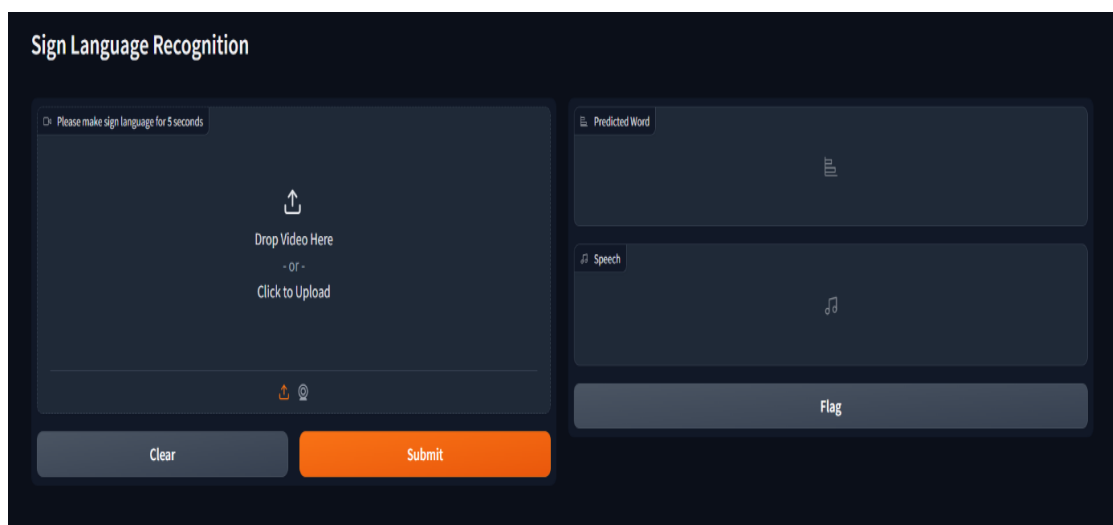
**Fig 4.2.7 Confusion Matrix**



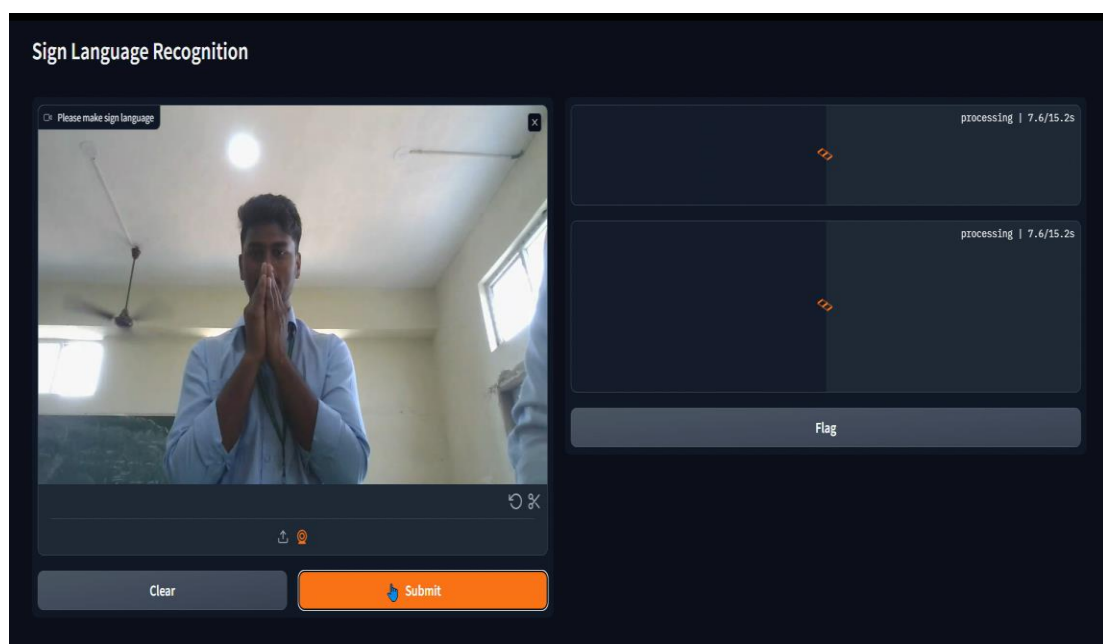
**Fig 4.2.8 Model Accuracy VS Epochs**



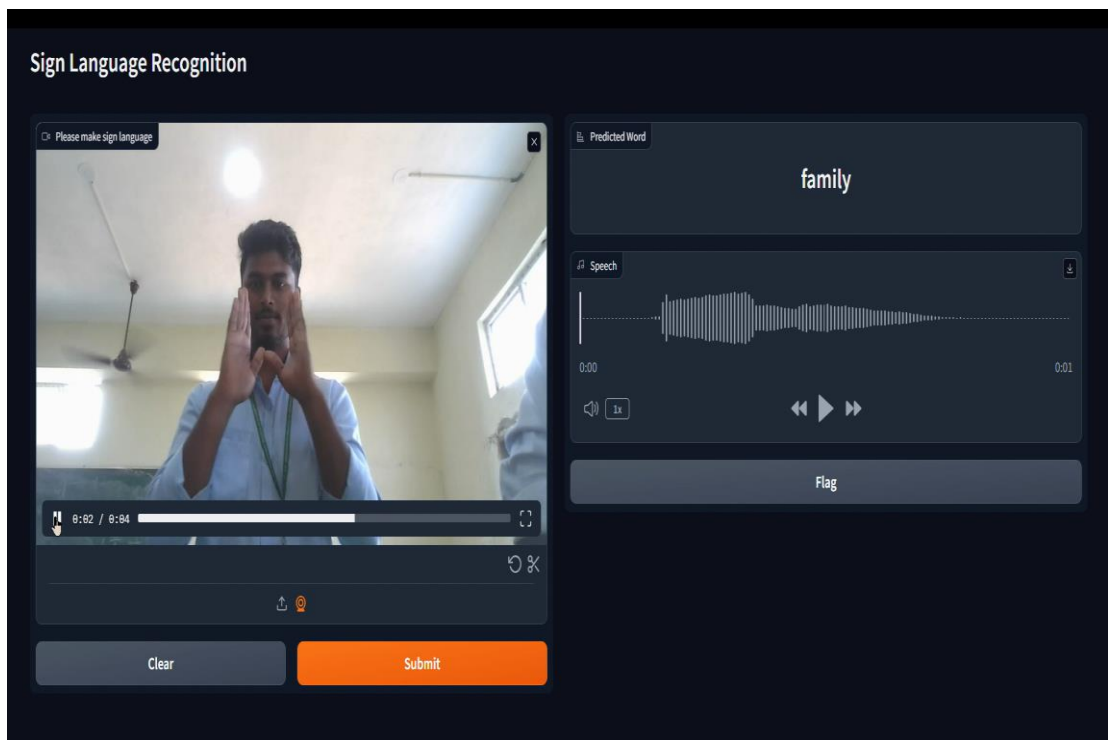
**Fig 4.2.9 Testing**



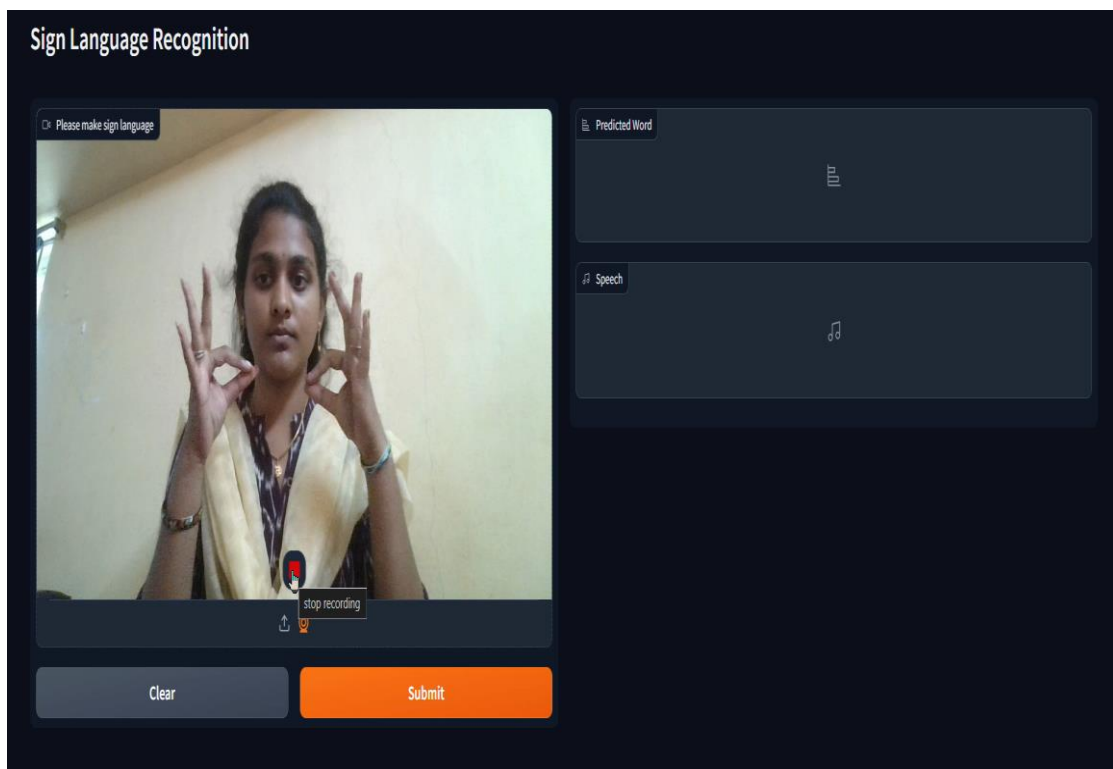
**Fig 4.2.10 User Interface**



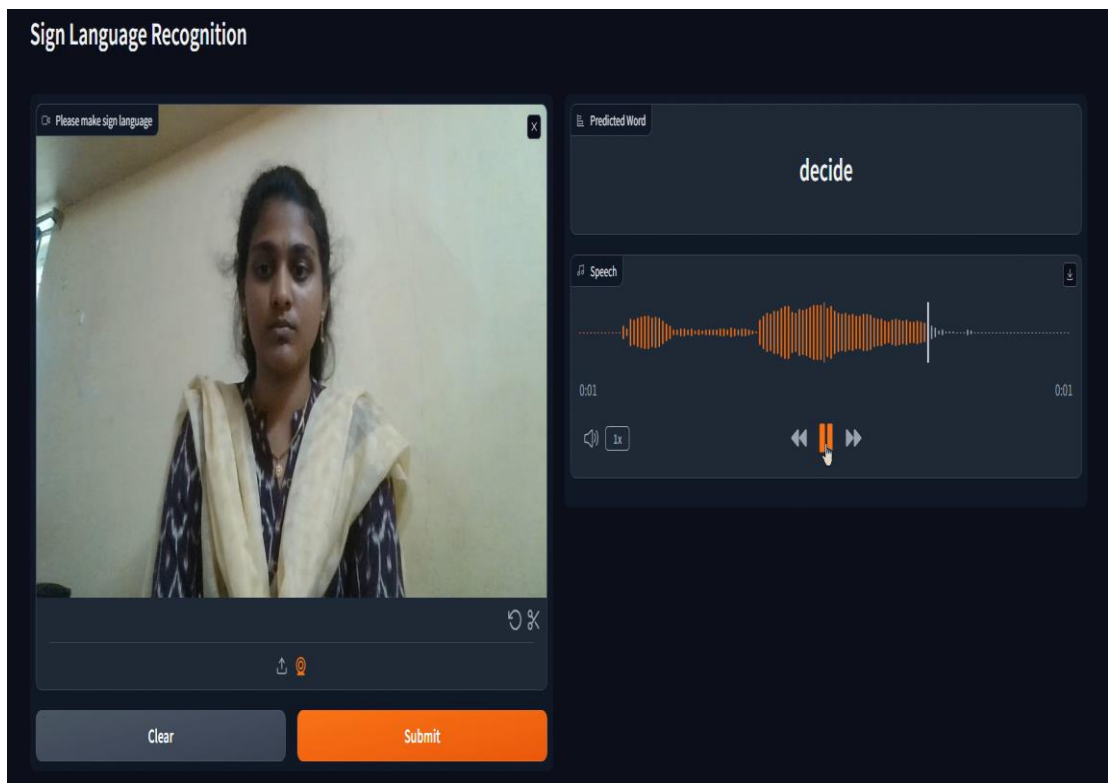
**Fig 4.2.11 Action 1 - Preprocessing**



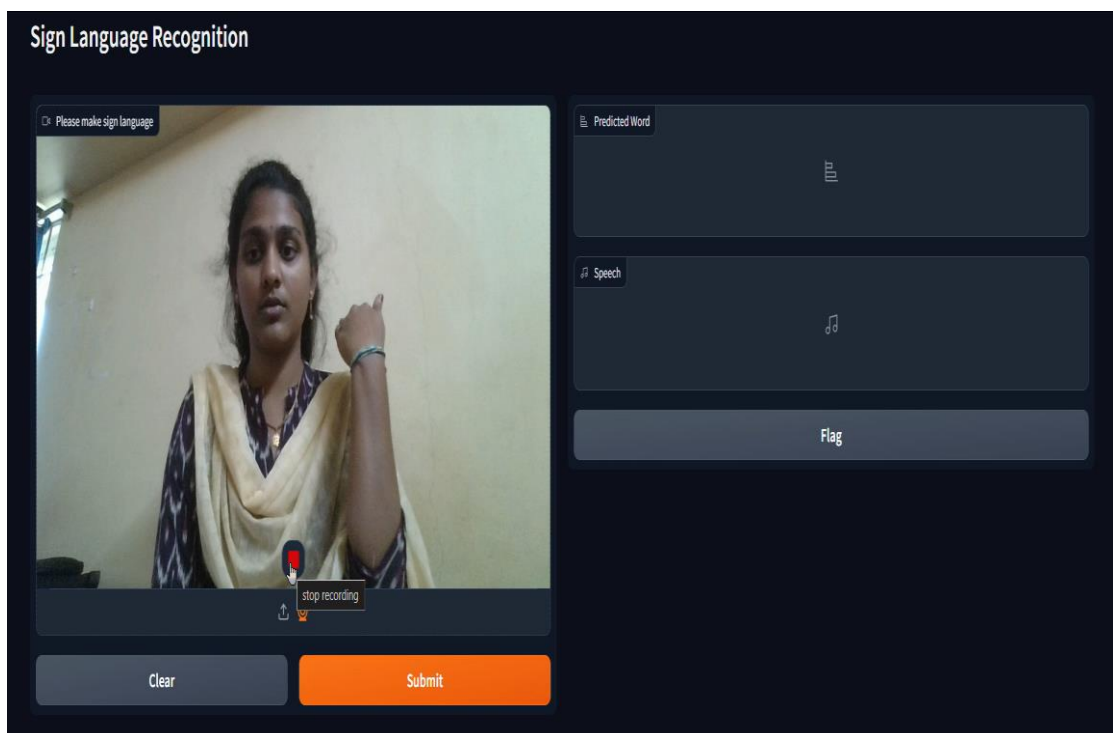
**Fig 4.2.12 Action 1 - Predicted Output**



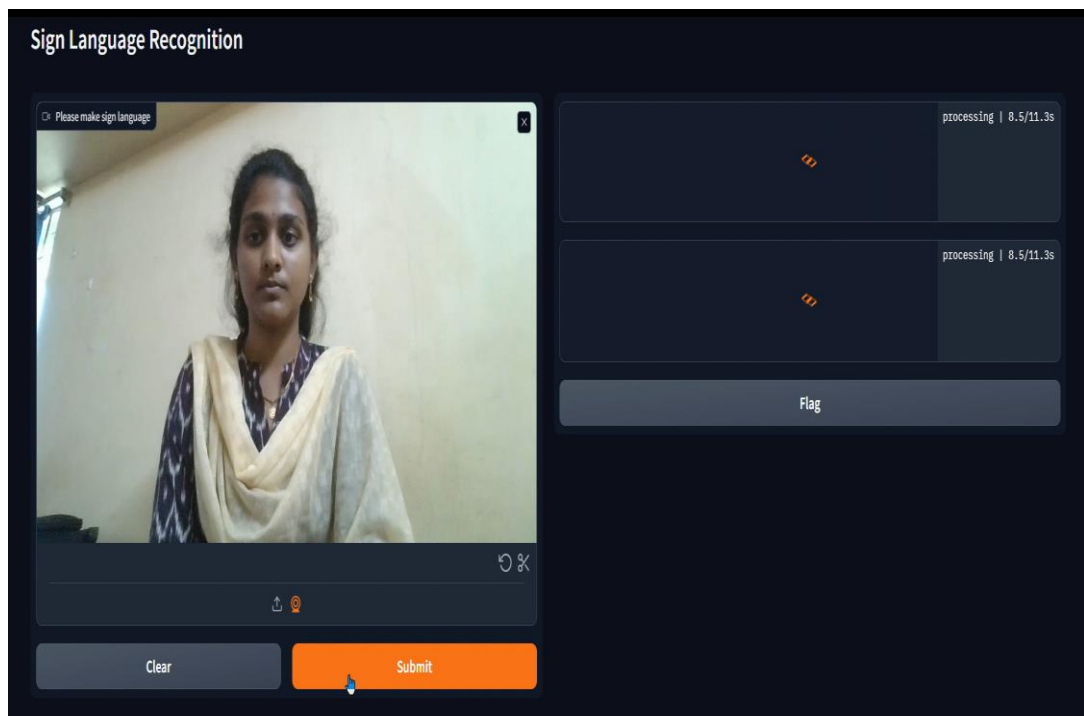
**Fig 4.2.13 Action 2 - Capturing Webcam Input**



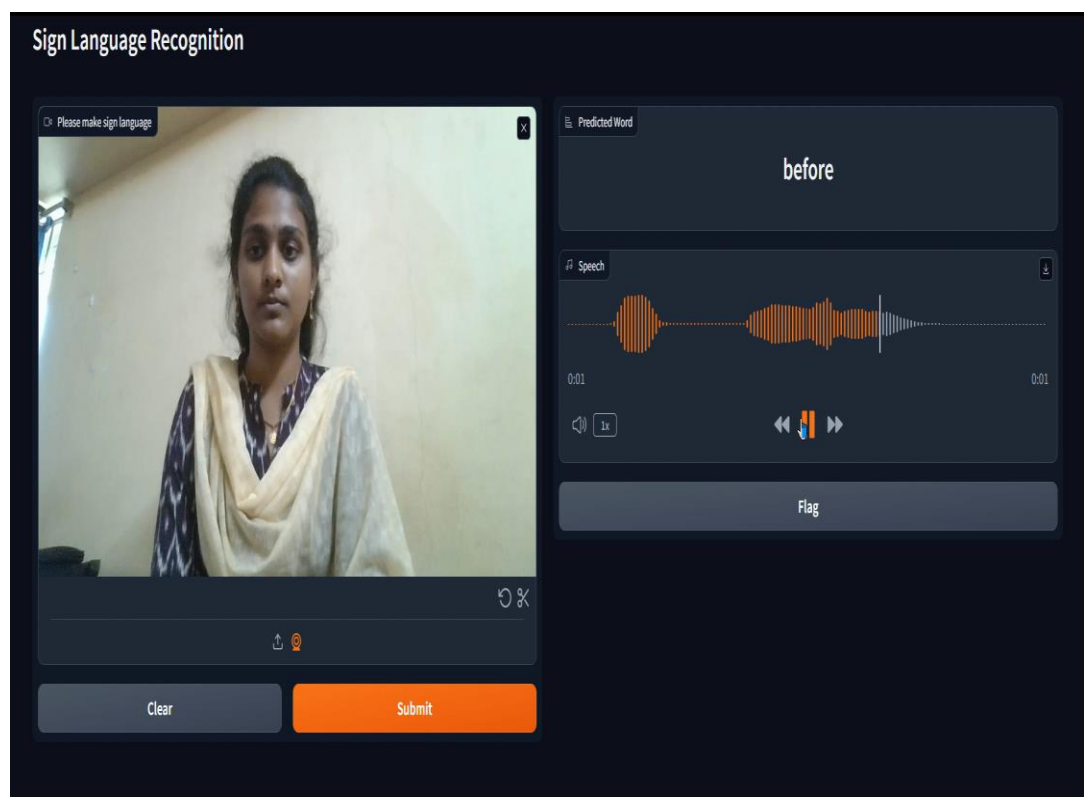
**Fig 4.2.14 Action 2 - Predicted Output**



**Fig 4.2.15 Action 3 - Capturing Webcam Input**

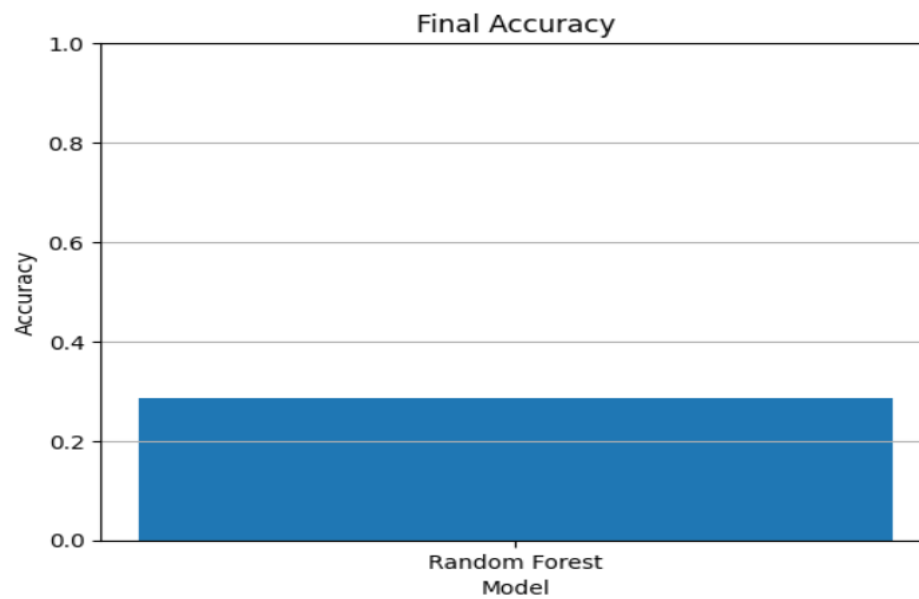


**Fig 4.2.16 Action 3 - Preprocessing**

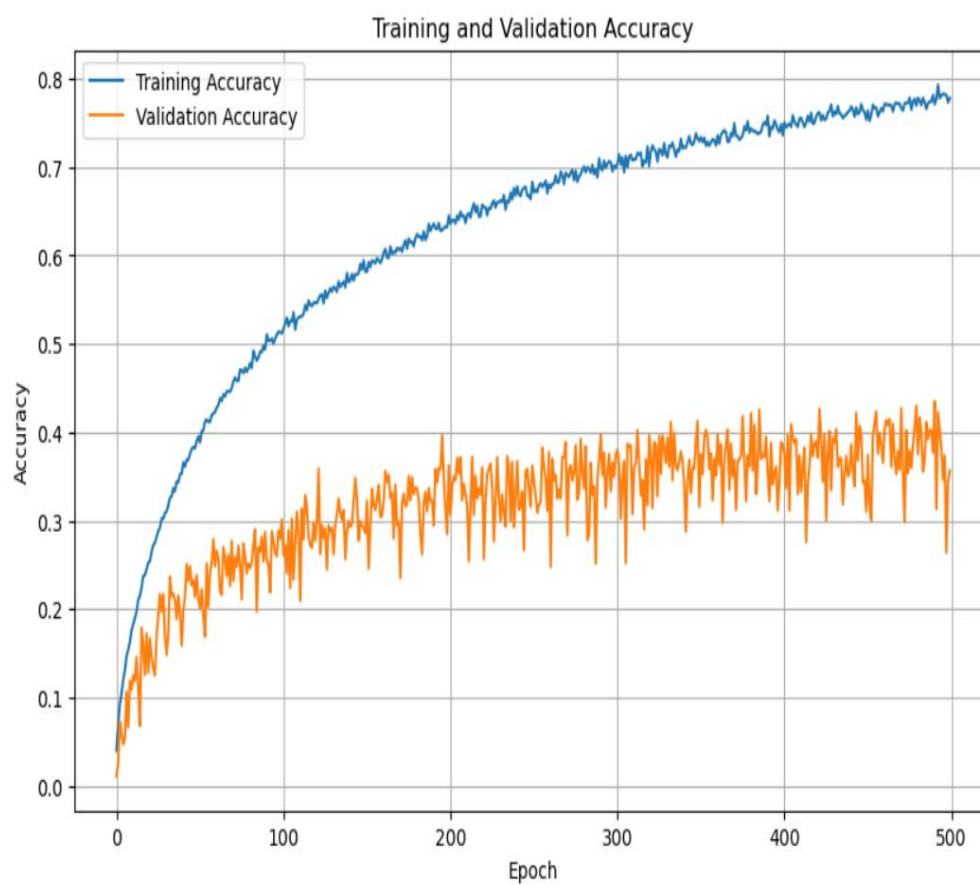


**Fig 4.2.17 Action 3 - Predicted Output**

### 4.3 Accuracy Graphs

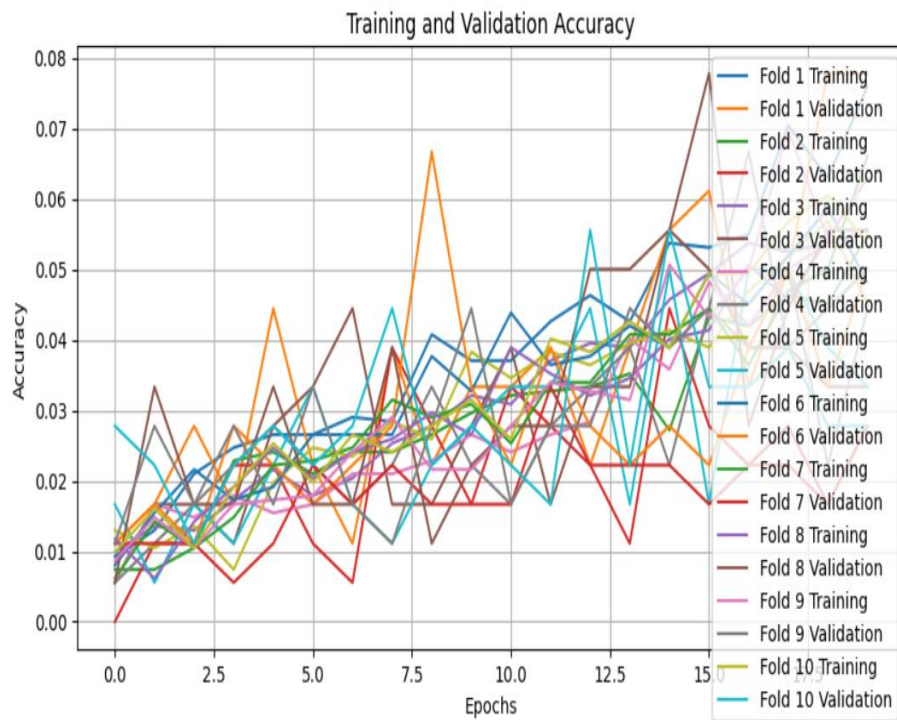


**Fig 4.3.1 Random Forest Model**

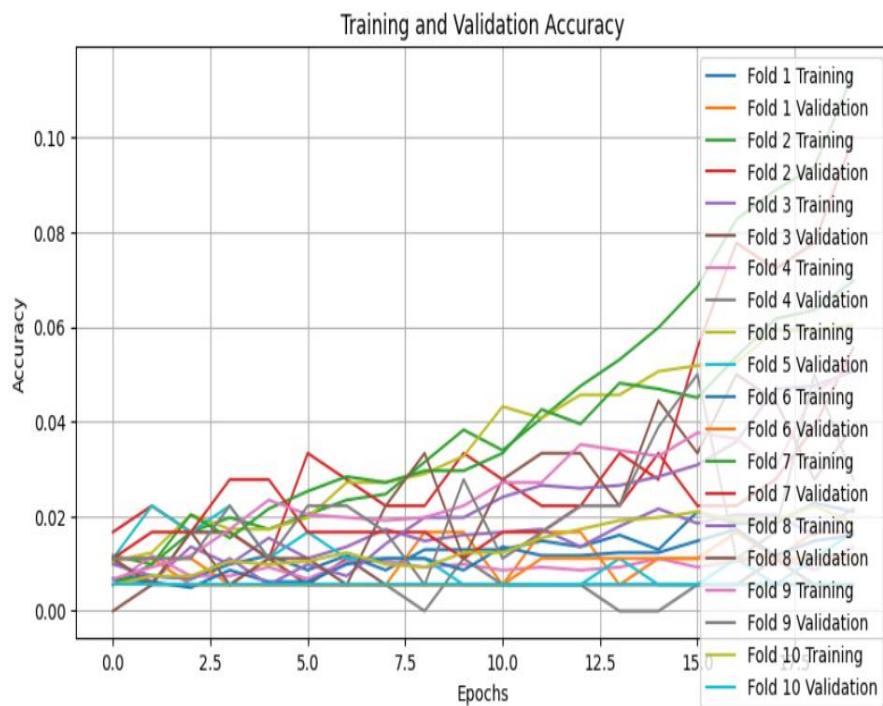


**Fig 4.3.2 Mediapipe + Baseline 2D Shallow CNN**

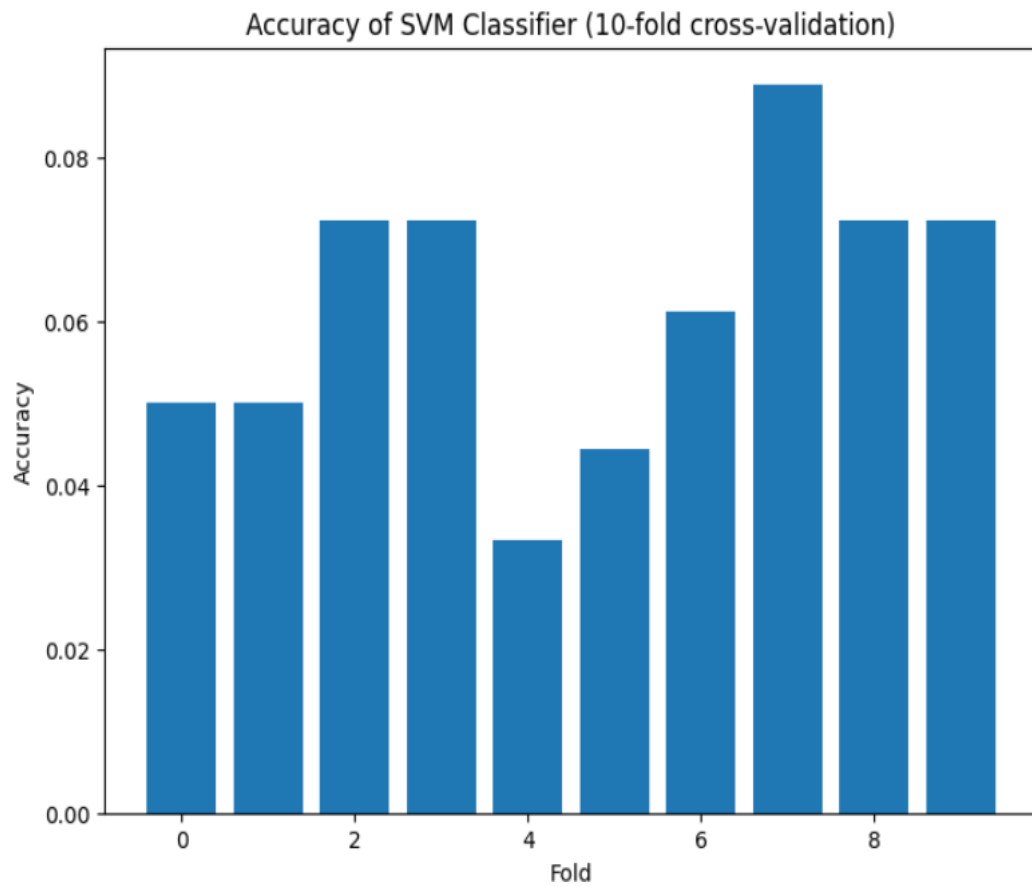




**Fig 4.3.3 Mediapipe + BiLSTM**



**Fig 4.3.4 Mediapipe + CNN**



**Fig 4.3.5 Mediapipe + SVM**



## **CHAPTER 5 - CONCLUSION AND FUTURE SCOPE**

### **5.1 CONCLUSION**

Our team has been diligently working on developing a robust sign language recognition (SLR) system, and we're thrilled to announce that it has achieved an accuracy rate of 79.7%. This significant milestone reflects our commitment to bridging the communication gap between those who use sign language and those who don't. Our system is designed to interpret sign language gestures and map them to the appropriate words, providing an essential tool for inclusive communication.

When examining the accuracy graph, you can visualize the steady progression our system made as we trained it. It's akin to watching a plant grow – it starts small, but with consistent care and attention, it flourishes over time. The accuracy graph serves as a testament to our system's learning capacity and resilience, indicating that our model is becoming increasingly proficient at recognizing sign language gestures. This upward trend is a promising sign that our SLR system has the potential for further growth and improvement.

The confusion matrix is another valuable tool that provides insights into our system's performance. Although it might sound complex, it's essentially a detailed chart that reveals where our system excels and where it could use some refinement. The confusion matrix is like a diagnostic map, allowing us to pinpoint the specific areas where our system might struggle or misinterpret gestures. By analyzing this data, we can identify patterns and focus our efforts on addressing these challenges, ultimately enhancing the system's overall accuracy and reliability.

Why is all this important? Our SLR system isn't just another tech project – it's a meaningful step toward creating a more inclusive world. By facilitating communication between sign language users and those who don't know sign language, we're helping to break down barriers that can lead to isolation and misunderstanding. Our system empowers people to connect and communicate more effectively, promoting a sense of unity and inclusivity. It's as if we're building a bridge that brings people closer together, using technology to ensure that everyone is understood and valued.

As we continue to improve our SLR system, we're driven by the goal of making communication accessible to everyone, regardless of their ability to hear or use sign

language. Through our efforts, we're demonstrating that technology can play a pivotal role in fostering a more inclusive and compassionate society, where everyone has the opportunity to be heard and understood.

## **5.2 FUTURE WORK**

In our study, we've explored sign language recognition (SLR) with a vocabulary of 100 words, giving us a solid start. However, there's potential to improve by expanding our system to include a broader range of signs and words, allowing us to meet the needs of a more diverse user base.

For future work, we could experiment with different types of technology to enhance SLR accuracy. This could include using advanced cameras and sensors, as well as leveraging sophisticated software for better sign recognition. A personalized approach could also help, where the system learns to adapt to an individual's signing style, making communication more effective and intuitive.

Additionally, collaborating with other researchers and accessing larger datasets could improve our system's learning process. By sharing resources and data, we can create a more robust and comprehensive SLR system. This would pave the way for broader applications, from education to customer service.

Another promising direction involves real-time translation features. By enabling our SLR system to translate sign language into text or speech in real time, we could help sign language users communicate more easily with those who don't use sign language. This would make interactions smoother and foster greater inclusivity.

These future directions provide a roadmap for enhancing our SLR system, with the ultimate goal of improving communication and creating a more inclusive environment for everyone.

## BIBLIOGRAPHY

- [1] Prof. Radha S. Shirbhate , Mr. Vedant D. Shinde , Ms. Sanam A. Metkari , Ms. Pooja U. Borkar , Ms. Mayuri A. Khandge. “Sign language Recognition Using Machine Learning Algorithm.” International Research Journal of Engineering and Technology (2020)
- [2] Ketan Gomase , Akshata Dhanawade , Prasad Gurav , Sandesh Lokare. “Sign Language Recognition using Mediapipe” International Research Journal of Engineering and Technology (IRJET) (2022)
- [3] Md.AhasanAtick Faisal , Farhan FuadAbir , Mosabber UddinAhmed , MdAtiqur RahmanAhad. “Exploiting domain transformation and deep learning for hand gesture recognition using a low-cost dataglove.” Scientific Reports (2022)
- [4] Akshay Divkar, Pranav Kulkarni, Pranav Patil, and Pranav Kulkarni. “ Gesture Based Real-time Indian Sign Language Interpreter.” International Journal of Scientific Research in Computer Science, Engineering and Information Technology (IJSRCSEIT) (2021).
- [5] Dongxu Li , Cristian Rodriguez Opazo, Xin Yu, Hongdong Li. “Word-level Deep Sign Language Recognition from Video: A New Large-scale Dataset and Methods Comparison.”(2020)
- [6] Shayla Luong “Video Sign Language Recognition using Pose Extraction and Deep Learning Models.” San Jose State University SJSU ScholarWorks (2023)
- [7] <https://dxli94.github.io/WLASL/>
- [8] <https://www.analyticsvidhya.com/blog/2022/02/k-fold-cross-validation-technique-and-its-essentials/#:~:text=K%2Dfold%20cross%2Dvalidation%20is,estimate%20the%20model's%20generalization%20performance.>
- [9] <https://www.analyticsvidhya.com/blog/2021/04/create-interface-for-your-machine-learning-models-using-gradio-python-library/>

- [10] Akriti Goyal, Deepanshu Dhar, Paras A Nair, Chirag Saini, Supreetha S M, Chetana Prakash “Sign Language Recognition System” International Journal of Engineering Research & Technology (IJERT) 2022
- [11] Benedicta Nana Esi Nyarko, Wu Bin, Jinzhi Zhou, George K. Agordzo . “Comparative Analysis of AlexNet, Resnet-50, and Inception-V3 Models on Masked Face Recognition.” 2022 IEEE World AI IoT Congress (AIIoT)
- [12] <https://en.wikipedia.org/wiki/ImageNet>
- [13] G. Mallikarjuna Rao , Cheguri Sowmya , Dharavath Mamatha , P. A.Sujasri , Soma Anitha , Ramavath Alivela “Sign Language Recognition using LSTM and Media Pipe” (2023)
- [14] Andriy Oliynyk , 3rd year students Roman Polishchenko Valentyn Kofanov “Sign language recognition using deep learning” Taras Shevchenko National University of Kyiv (2020)
- [15] Pooja M.R, Meghana M, Praful Koppalkar, Bopanna M J, Harshith Bhaskar, Anusha Hullali “Sign Language Recognition System” Indian Journal of Software Engineering and Project Management (IJSEPM) 2022
- [16] Yulius Obi , Kent Samuel Claudio , Vetri Marvel Budiman , Said Achmada, Aditya Kurniawan “Sign language recognition system for communicating to people with disabilities” 7th International Conference on Computer Science and Computational Intelligence 2022
- [17] Suharjito , Ricky Anderson , Fanny Wiryana , Meita Chandra Ariesta , Gede Putra Kusuma “Sign Language Recognition Application Systems for Deaf-Mute People: A Review Based on Input-Process-Output” 2nd International Conference on Computer Science and Computational Intelligence 2017, ICCSCI
- [18] R Rumana , Reddygari Sandhya Rani , Mrs. R. Prema “A Review Paper on Sign Language Recognition for The Deaf and Dumb” International Journal of Engineering Research & Technology (IJERT) 2021

[19] Rachana Patil , Vivek Patil , Abhishek Bahuguna , and Mr. Gaurav Datkhile “Indian Sign Language Recognition using Convolutional Neural Network” ITM Web of Conferences (2021)

[20] Prof. Mrs. Maheshwari Chitampalli, Dnyaneshwari Takalkar, Gaytri Pillai, Pradnya Gaykar, Sanya Khubchandani “REAL TIME SIGN LANGUAGE DETECTION” International Research Journal of Modernization in Engineering Technology and Science (April-2023)