

HAND GESTURE RECOGNITION USING COMPUTER VISION FOR SIGN LANGUAGE

A PROJECT REPORT
SUBMITTED IN PARTIAL FULFILMENT OF THE REQUIREMENTS
FOR THE AWARD OF THE DEGREE
OF

**BACHELOR OF TECHNOLOGY IN
COMPUTER ENGINEERING**

Submitted by:
Kinzang Tobgay (2K21/CO/241)
Sumpanna Acharya (2K21/CO/484)
Ugyen Wangchuk (2K21/CO/497)

Under the supervision of
Dr. R K Yadav
Associate Professor



DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING
DELHI TECHNOLOGICAL UNIVERSITY
(Formerly Delhi College of Engineering)
Bawana Road, New Delhi-110042

MAY, 2025

DELHI TECHNOLOGICAL UNIVERSITY
(Formerly Delhi College of Engineering)
Bawana Road, New Delhi-110042

CANDIDATE'S DECLARATION

We, **Sumpanna Acharya (2K21/CO/484), Ugyen Wangchuk (2K21/CO/497) & Kinzang Tobgay (2K21/CO/241)**, students of B. Tech Computer Engineering, hereby declare that the project Dissertation titled **“HAND GESTURE RECOGNITION USING COMPUTER VISION FOR SIGN LANGUAGE”** which is submitted by us to the Department of Computer Science & Engineering, Delhi Technological University, Delhi in partial fulfilment of the requirement for the award of the degree of Bachelor of Technology, is original and not copied from any source without proper citation. This work has not previously formed the basis for the award of any Degree, Diploma Associateship, Fellowship or other similar title or recognition.

Place: Delhi, India

Kinzang Tobgay

Date: 22th May, 2025

Sumpanna Acharya

Ugyen Wangchuk

DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING
DELHI TECHNOLOGICAL UNIVERSITY
(Formerly Delhi College of Engineering)
Bawana Road, New Delhi-110042

CERTIFICATE

I hereby certify that the Project Dissertation titled “**HAND GESTURE RECOGNITION USING COMPUTER VISION FOR SIGN LANGUAGE**” which is submitted by **Sumpanna Acharya (2K21/CO/484), Ugyen Wangchuk (2K21/CO/497) & Kinzang Tobgay (2K21/CO/241)** in partial fulfilment of the requirement for the award of Degree of Bachelor of Technology, as a record of the project work carried out by the students under my supervision. To my knowledge this work is not been submitted in part, or full for any Degree or Diploma to this University, or elsewhere.

Place: Delhi

Dr. R K Yadav

Date: 22th May, 2024

(SUPERVISOR)

ACKNOWLEDGEMENT

We would like to take this prospect to acknowledge the support and guidance of all those without whom the project would not have been possible. First and foremost, we would like to sincerely thank our mentor, Dr. R K Yadav, for his guidance, support, criticism, and encouragement which led to the completion of the project. The interactions and communications with him gave us a deep insight into this topic and advanced our ideas. His enthusiastic attitude helped us to sail through and stay on track.

We also express our gratefulness to all other faculty members of the department for their perceptive guidance, constant encouragement, and sincere support for this project work. We also place on record our sense of gratitude to one and all, for directly or indirectly having lent their helping hand in this venture.

Finally, we show our appreciation the effort of our parents and friends for supporting us through this project. Their trust in our competences helped us in providing our best.

ABSTRACT

The world's population now stands at around 8 billion and out of these many, there are millions of people who are unfortunately deprived of the ability to communicate in the same way as others even though there is a gap in communication between these people and the general abled people there are many types of sign languages which help in decreasing this gap. This paper describes a system that is aimed to build a bridge between the two groups of people it is a fact that this kind of technology has helped reduce the communication gap by converting different gestures to text or speech. The main focus of this project is to produce a real-time sign language translator that can be used in language translation. Our system will be able to recognize pre-stored gestures according to asl (American sign language) [21]. This would contain a user-friendly environment for the users by providing speech or text output for a sign gesture input and vice versa.

Keywords – communication, sign language, gestures, translator & technology.

TABLE OF CONTENTS

CANDIDATE’S DECLARATION	II
CERTIFICATE	III
AKNOWLEDGEMENT	IV
ABSTRACT	V
LIST OF FIGURES & TABLES	VIII
SYMBOLS AND ABBREVIATIONS	IX
CHAPTER 1	
INTRODUCTION	10
1.1 Background	
1.1.1 Introduction to Sign Language Communication	10
1.1.2 Role of Computer Vision in Gesture Recognition	11
1.1.3 Evolution of Human-Computer Interaction (HCI)	12
1.2 Literature Background	
1.2.1 Existing Sign Language Recognition Systems	13
1.2.2 Techniques Used in Hand Gesture Recognition	14
1.2.3 Datasets and Evaluation Metrics	15
1.3 Identification of Research Gaps	
1.3.1 Real-Time Recognition Limitations	16
1.3.2 Lack of Multilingual Sign Language Models	16
1.3.3 Environmental and Contextual Robustness	17
1.3.4 Public Awareness	17
CHAPTER 2	
THEORETICAL BACKGROUND	18
2.1 Literature Reviews	
2.1.1 Gesture-Based Approaches	18
2.1.2 Vision based techniques	19
2.1.3 Sensor-based approaches	19
2.1.4 Multimodal Approaches	20
2.1.5 Sign Language to Text/voice Translation	20

2.2	Prerequisites	
2.2.1	SSD MobileNetV2	21
2.2.2	TensorFlow	21
2.2.3	Open CV	21
2.2.4	Convolution Neural Network	22
2.2.5	PyTorch	24
2.2.6	Media Pipe	24
2.2.7	Transformers	25
2.2.8	Open Pose	25

CHAPTER 3

METHODOLOGY		26
3.1	Proposed Approach	26
3.2	System Specifications	26
3.3	System Architecture	27
3.4	Data Collection Method	
3.4.1	Collection of Primary Data	28
3.4.2	Data Sources	28
3.4.3	Data Collection Process	28
3.4.4	Ethical Considerations	29
3.5	Data Analysis Method	
3.5.1	Data Preprocessing	29
3.5.2	Feature Extraction	30
3.5.3	Model Architecture	31
3.5.4	Real-Time Gesture Prediction	32
3.5.5	Output	32
3.5.6	Parameters to be Evaluated	

CHAPTER 4

FINDINGS AND ANALYSIS		33
4.1	Evaluation Metrics	33
4.2	Performance Metrics	33
4.3	Comparative Analysis with Previous Research	34
4.4	Discussion	35

CHAPTER 5

CONCLUSION AND FUTURE WORK	36
REFERENCES	37

LIST OF FIGURES & TABLES

Figure 1. 1: Types of signs in ASL	10
Figure 1. 2: Palm key-points using Computer Vision	11
Figure 1. 3: Human-Computer Interaction (HCI)	12
Figure 1.4: Sign Language Recognition System	13
Figure 1. 5: Hand Gesture Recognition System	14
Figure 1. 6: Essential Classification Algorithms	15
Figure 1. 7: Real-Time Recognition Algorithm.....	16
Figure 2. 1: Gesture-Based Model	18
Figure 2. 2: Vision-Based Approach.	19
Figure 2.3: Sensor-Based Approach	19
Figure 2.4: Multimodal-Based Approach	20
Figure 2.5: Sign to Voice Translation.....	20
Figure 2.6: Convolution Neural Network Model.....	22
Figure 2.7: Convolution Table	22
Figure 2.8: Max-pooling	23
Figure 2.9: Average-pooling	23
Figure 2.10: Flattening	23
Figure 2.11: MediaPipe Hand Model	24
Figure 3.1: Model Architecture.....	27
Figure 3.2: Output	31
Figure 4.1: Evaluation Metrics	33
Figure 4.2: Confusion Matrix	34
.....	
Table 4.1: Performance Metrix	33
Table 4.2: Comparative Analysis.....	34

SYMBOLS AND ABBREVIATIONS

Symbol/Abbreviation	Meaning
AI	Artificial Intelligence
ML	Machine Learning
DL	Deep Learning
CV	Computer Vision
CNN	Convolutional Neural Network
ASL	American Sign Language
RGB	Red, Green, Blue (Colour channels in images)
FPS	Frames Per Second
TP	True Positive
FP	False Positive
TN	True Negative
FN	False Negative
Epoch	One full pass of the training dataset
Batch Size	Number of training samples used in one iteration
Dataset	Collection of data used for training and testing
Accuracy	$(TP + TN) / \text{Total Predictions}$
Precision	$TP / (TP + FP)$
Recall	$TP / (TP + FN)$
F1 Score	$2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$
GPU	Graphics Processing Unit
API	Application Programming Interface
YOLO	You Only Look Once (real-time object detection algorithm)
OpenCV	Open-Source Computer Vision Library
TTS	Text-to-Speech
STT	Speech-to-Text

CHAPTER 1

INTRODUCTION

1.1 Background

1.1.1 Introduction to Sign Language Communication

Sign language is an essential form of communication used by deaf and hard-of-hearing individuals to interact with others. It employs visual gestures, facial expressions, and body movements to convey meaning, making it fundamentally different from spoken languages, which rely on sound. One of the most well-known sign languages is American Sign Language (ASL) [15], used primarily in the United States and Canada, but many countries have their own unique sign languages, such as British Sign Language (BSL), French Sign Language (LSF) [16], and more.

The need for sign language is more than just communication; it is culture, identity and community for the deaf. Nevertheless, sign language is frequently overlooked in the technological world, making it less available than other forms of language. Often, deaf people encounter challenges trying to communicate with hearing people who do not use sign language, making them feel more isolated. Such isolating situations highlight the need for enabling technologies. For instance, creating devices that interpret gestures of sign language into texts or voice could help deaf and hearing people interact more without requiring them both to know sign language. With the reliance of millions on sign language, there is urgency for the development of recognition systems which provide real-time translation of sign language. Such systems would not only enhance the availability of communication, but also the interaction people have with machines and technology.

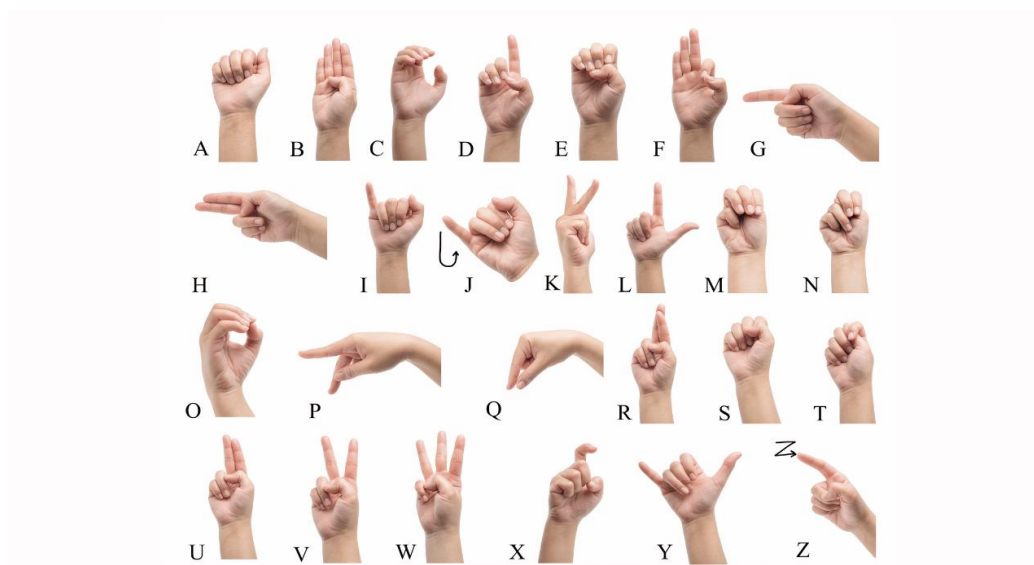


Figure 1. 1: Types of signs in ASL (Source: math4children.com)

1.1.2 Role of Computer Vision in Gesture Recognition

Computer vision (CV) is one of the emerging areas of artificial intelligence (AI) that focuses on allowing computers to make sense of the visual world [17]. CV technologies have the capabilities to identify objects, monitor activities, and even recognize actions by analyzing images or frames from videos. The Recognition of hand gestures is one of the most advance and promising applications of CV since it enables machines to communicate with users using gestures like those used in human-to-human interaction.

CV methods can be applied to dissect hand movements and relate them to specific signs in the case of sign language. Numerous methods have been used to create efficient systems for recognizing hand gestures, including Convolutional Neural Networks (CNNs) and a variety of feature extraction techniques [18]. To illustrate, the application of CNNs permits the identification of sophisticated patterns within images of hands, even when the images are captured in adverse circumstances such as with cluttered surroundings or changing lights. Furthermore, recognition of gestures in real time is possible nowadays thanks to improvements in computing power and algorithm development.

Computer vision has the benefit of a number of aspects over other gesture recognition systems, including wearable sensors. Sensor-based systems are dependent on physical devices that are not only inconvenient but also uncomfortable to wear, whereas CV-based systems only need cameras and hence are more convenient and easier to employ in various environments.

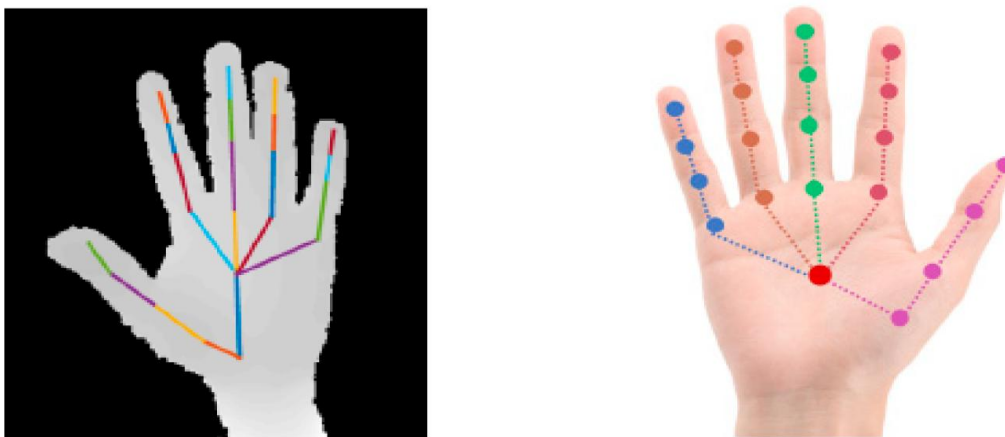


Figure 1. 2: Palm key-points using Computer Vision (Source: medium.com)

1.1.3 Evolution of Human-Computer Interaction (HCI)

The field of Human-Computer Interaction (HCI) has undergone significant change over the years. Early interfaces were mainly device-centric, such as keyboards, mice, and joysticks [19], which required learning specific input skills. Over time, the trend has been towards more user-oriented interfaces, such as touchscreens, that enabled more natural and intuitive interaction. The field is now moving towards more intuitive input devices, such as voice recognition, eye tracking, and gesture recognition.

Gesture recognition is a major innovation in the area of Human-Computer Interaction (HCI). This trend of interaction allows users to communicate or control devices or computer systems through body movements, e.g., hand gestures. The emergence of technologies like virtual reality (VR) and augmented reality (AR) [20] has fueled the development of gesture-based interfaces since these environments require new forms of interaction that go beyond the use of traditional input devices. Being a particular subset of gesture recognition, sign language recognition is an ideal candidate for this trend and is a vital area of research that endeavors to contribute to accessibility [21] in communicative and technological interfaces.

The emphasis of contemporary HCI study is natural and transparent interaction. Hand gesture recognition can make users interact with computers or other users without any need for physical touch or professional training, and therefore the systems are extremely accessible to the disabled.

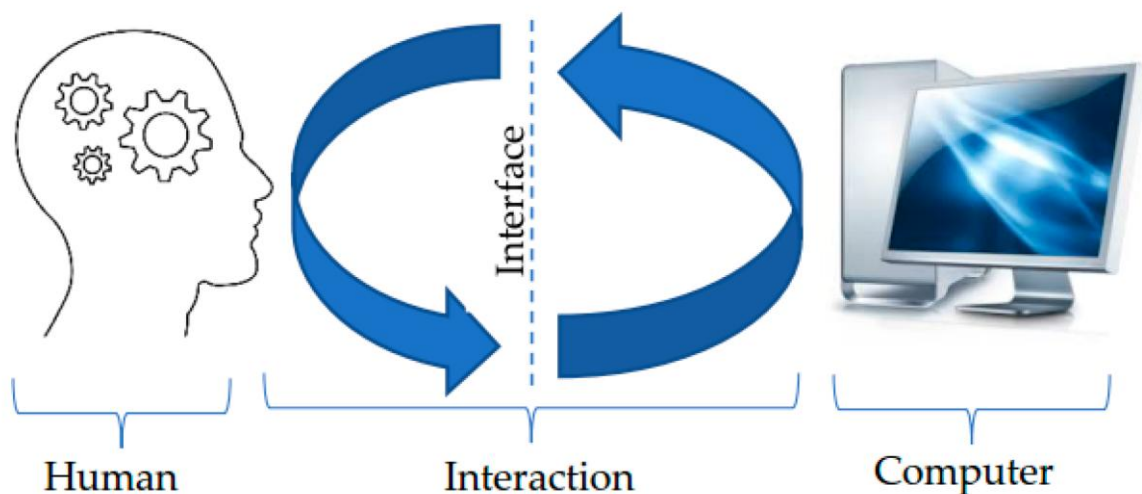


Figure 1. 3: Human-Computer Interaction (HCI) (Source: mdpi.com)

1.2 Literature Background

1.2.1 Existing Sign Language Recognition Systems

There have been various research studies and commercial systems towards the creation of sign language recognition technologies, especially for the American Sign Language (ASL). They typically employ computer vision and machine learning to track and decipher [22] hand movement. Some of the initial gesture recognition methods were considerably reliant on tracking hands and fingers with the help of external sensors or cameras, but with the improvement in computer vision algorithms, camera-based systems are now the norm.

One of the most well-documented examples of this is the "DeepASL" system, which applies deep learning techniques to the recognition of American Sign Language (ASL) hand gestures. The system uses a convolutional neural network (CNN) to process images, recognizing hand shapes, locations, and movements and thereby translating them into the corresponding words or sentences. Another instance is the "Vision-based Hand Gesture Recognition System," which uses a hybrid system that blends hand gesture recognition and object recognition to provide higher accuracy in dynamic environments.

Commercial systems, which include businesses such as SignAll, provide real-time sign language interpretation, thus facilitating communication between hearing impaired individuals [21] and people who are not familiar with sign language. Such systems apply advanced cameras in combination with machine learning algorithms to identify and interpret gestures.

Although there has been improvement in such systems, there are still challenges. Most systems currently in place still face difficulties with issues of environmental conditions, such as changes in lighting, ambient noise, or variations in users' hand shapes and sizes. Furthermore, real-time performance and accuracy enhancement are still areas that need attention.

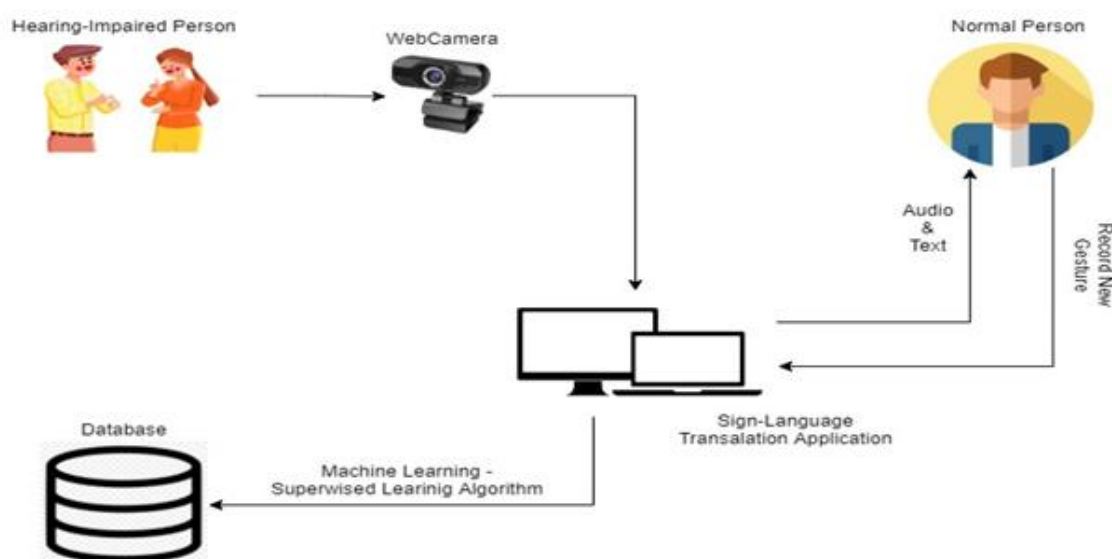


Figure 1. 4: Sign Language Recognition System (Source: mdpi.com)

1.2.2 Techniques Used in Hand Gesture Recognition

Several techniques are commonly used for hand gesture recognition in the context of sign language translation. Traditional image processing methods involve detecting the hand's shape, color, and motion within video frames. For example, skin color detection can help isolate the hand from the background. Once the hand is detected, techniques such as contour analysis, edge detection, and optical flow analysis can help identify the hand's position and movement.

Recent advancements in deep learning have significantly improved [23] the accuracy of hand gesture recognition. Convolutional Neural Networks (CNNs) are one of the most widely used models for image classification tasks, including hand gesture recognition. CNNs excel in learning spatial hierarchies of features from images, which is crucial for recognizing hand [24] shapes and motions accurately. Recurrent Neural Networks (RNNs), especially Long Short-Term Memory (LSTM) networks [21], are also used to analyse temporal data, such as the movement of hands over time.

Another powerful approach is the use of 3D hand models, which can capture the full geometry of the hand, allowing for more accurate gesture recognition even in dynamic scenarios. Recent systems use a combination of 2D and 3D convolutional layers to detect gestures with higher precision, even under different viewing angles.

Despite these advancements, challenges remain in handling occlusions (when parts of the hand are hidden from view), real-time recognition, and adapting to different individuals' hand shapes or gestures.

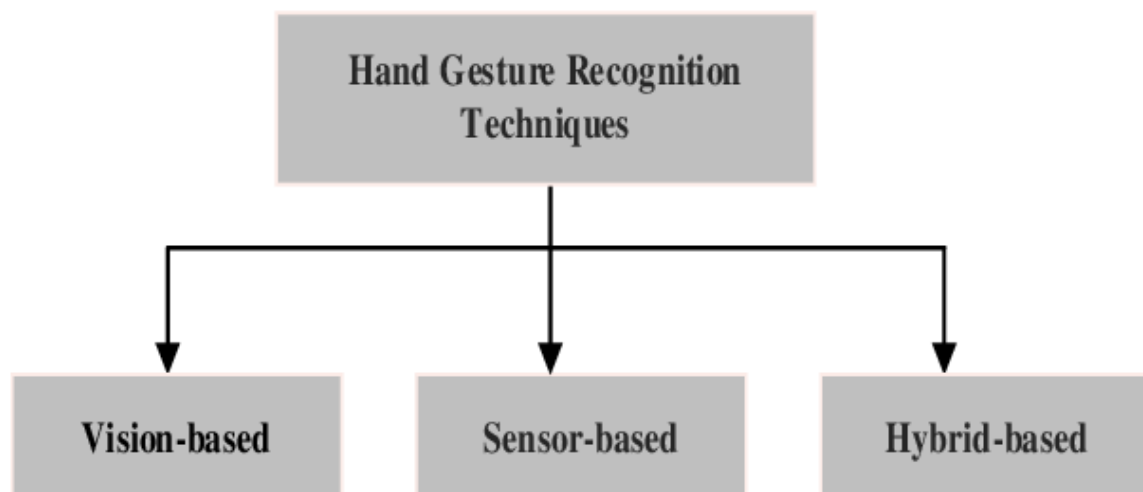


Figure 1. 5: Hand Gesture Recognition Technique (Source: researchgate.net)

1.2.3 Datasets and Evaluation Metrics

For research and testing to create gesture recognition systems, scientists use enormous collections of sign language gestures. Some of the notable datasets are the American Sign Language (ASL) dataset, which is characterized by a diverse range of ASL signs, and the RWTH-PHOENIX-Weather dataset, which consists of continuous sign language sentences in a conversation setting. These datasets offer a rich foundation for training machine learning models [25] and for testing their performance metrics.

Evaluation of gesture recognition systems is generally performed with measures such as accuracy, precision, recall, and F1-score [26]. Accuracy determines how often the system correctly identifies a gesture, while precision and recall are measures of how effectively the system constrains false positives and false negatives, respectively. Latency is also a factor of concern in real-time scenarios, since systems need to process and identify gestures in real-time in order to be effective.

Researchers also compare systems based on how well they can withstand environmental conditions, for instance, varying lighting, varied backgrounds, and background noise. Those models that generalize well across multiple contexts are most valuable, particularly for systems that will be used in real-world contexts.

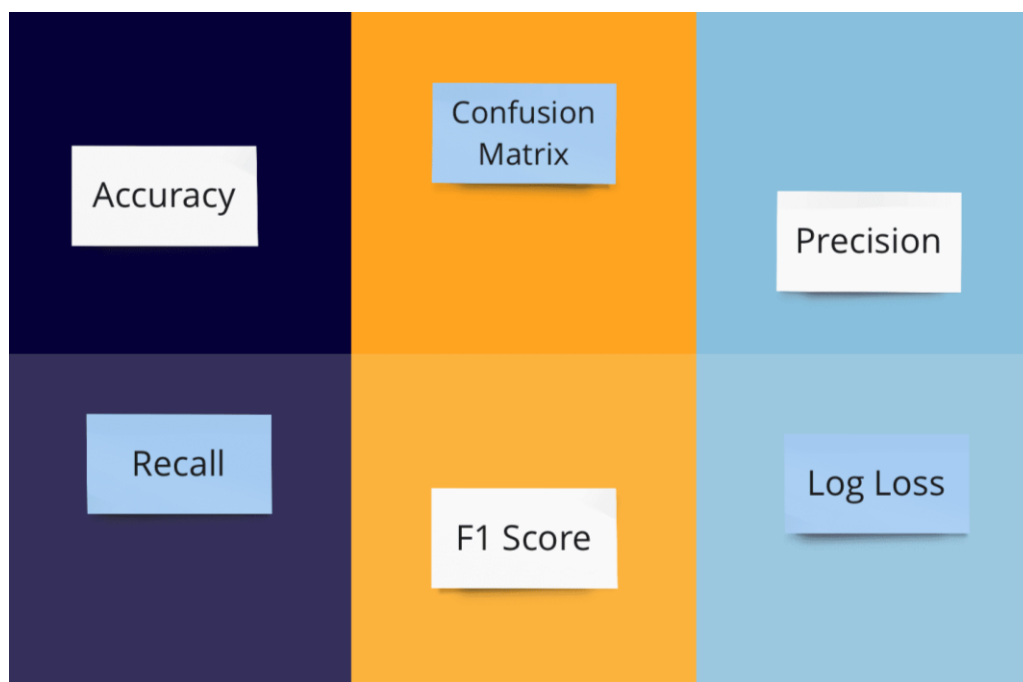


Figure 1. 6: Essential Classification Algorithms (Source: dataaspirant.com)

1.3 Identification of Research Gaps

1.3.1 Real-Time Recognition Limitations

Despite significant progress in hand gesture recognition, one of the major challenges in real-world applications is achieving real-time recognition with high accuracy. Real-time performance is crucial for systems intended to translate sign language into spoken or written [21] text, as delays can hinder communication. Current systems may struggle to process gestures quickly enough, leading to lag or inaccuracies, particularly in more complex or rapid sign language gestures.

Real-time recognition also requires powerful computational resources [27], which can be a limiting factor, especially for mobile devices or edge computing applications. While deep learning models like CNNs offer high accuracy, they also demand substantial computational power, which may not be available on all devices. One potential solution lies in optimizing models to reduce computational requirements while maintaining performance. Techniques such as model quantization, pruning, or using lightweight architectures like MobileNets may offer a balance between accuracy and efficiency.

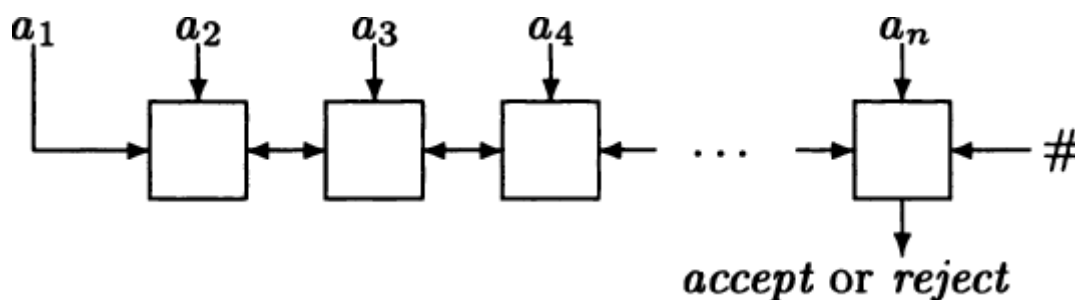


Figure 1. 7: Real-Time Recognition Algorithm (Source: researchgate.net)

1.3.2 Lack of Multilingual Sign Language Model

Another fundamental research gap is creating multilingual sign language recognition models. Current systems are mainly designed for a single language, i.e., ASL, but sign languages differ considerably from country to country and region to region. For example, BSL (British Sign Language) [15] is quite distinct from ASL, and even in the same nation, there are different dialects and types of the sign language.

Designing systems that can handle multiple sign languages or dialects is a big challenge, requiring large and varied datasets that capture the nuances present in each language. Additionally, sign language recognition systems will need to account, in addition to the gestures themselves, regional variations in signs and the cultural nuances of non-manual signals [28] (e.g., the position of the body and facial expressions). Researchers are actively pursuing methods for building multilingual models; however, this is a big area in the field left to be filled.

1.3.3 Environmental and Contextual Robustness

One of the primary limitations of current hand gesture recognition technologies is that they are environment-sensitive. Although such systems can be made highly accurate within research laboratories, they are significantly impaired in actual usage scenarios due to variations in lighting, environmental interference, and orientation of the users. For instance, a system will struggle to recognize a gesture when a hand is covered by another object or when there is poor lighting.

Furthermore, the environmental factors, such as the attire on the user or the pose of the body, can interfere with the system from being able to differentiate hand movement and recognize gestures accurately. More research needs to be conducted to develop robust systems that can generalize across different environments and still provide recognition accuracy. Implementation of techniques like effective feature extraction, sophisticated background subtraction, and deep learning models with the ability to generalize well across diverse environments is essential in achieving this requirement.

1.3.4 Public Awareness

Sign language recognition by hand sign, as explained in reference [21], is also dependent to a large degree upon public awareness, which is required for both its general acceptance and effectiveness. Despite much progress being achieved in gesture recognition and sign language translation, their usefulness can be limited if the general public is not aware of these technologies or has no knowledge of their applications. Public awareness in terms of sign language recognition systems, their numerous applications, and the far-reaching social contribution that can be made using them is essential in achieving the full potential of these technologies.

CHAPTER 2

THEORETICAL BACKGROUND OF THE STUDY FROM CONTEXT OF OPERATION RESEARCH

2.1 Literature reviews

This topic has gained popularity among the researchers due to the advent of new technologies such as computer vision, deep learning and Natural language processing. However, this project topic still has lots of unsolved challenges for providing a seamless and smooth communication between hearing impaired community and normal community. The following content highlights key studies and processes primary focusing on gesture-based, image-based approaches.

2.1.1 Gesture-Based Approaches

Zhao et al. (2015) [5] presented a paper which made use of hidden Markov Model (HMM) which was one of the oldest techniques used in this field. This made use of gesture-based model. The advantage was that it achieved a good accuracy value the real-time performance was poor and it was affected by the noise from the environment. Later Wang and lee (2017) [6] proposed another model called Hybrid model which was effectively a combination of CNN and HMM. The result was a improved real-time interpretation and accuracy since background noise was reduced.

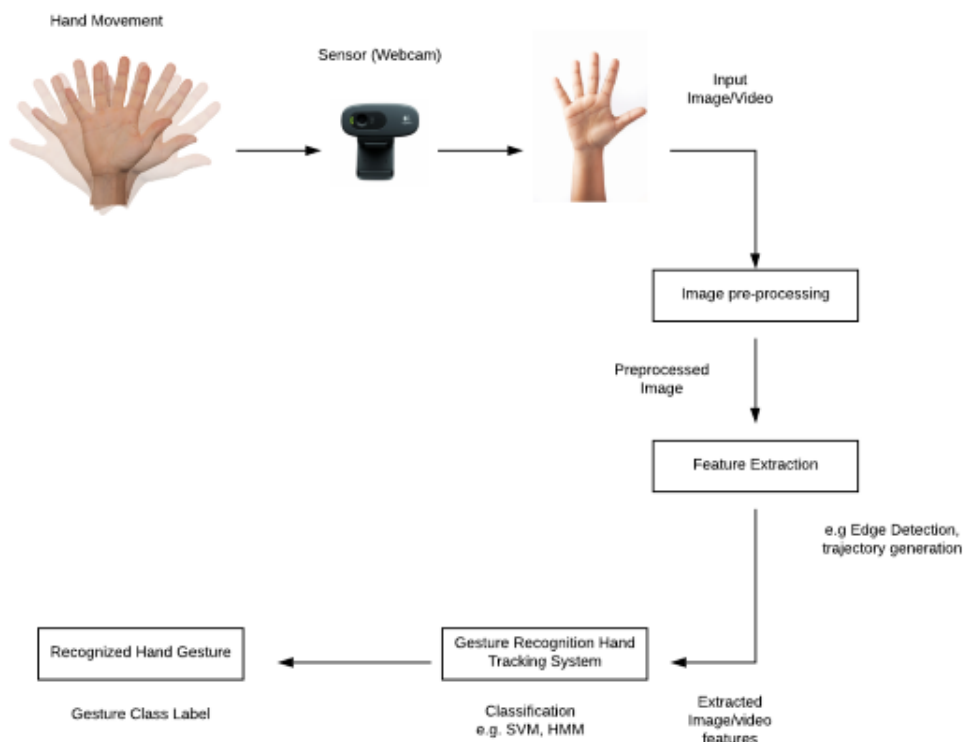


Figure 2. 1: Gesture-Based Model (Source: medium.com)

2.1.2 Vision based techniques

This technique involves use of camera-based systems for interpreting signs. Pujari et al. (2018) [7] achieved an astounding accuracy of 94%. Since every individual uses different styles, their work emphasized use of large annotated datasets. Another pioneer was Sharma et al. (2020) [8] who used CNN to identify individual signs from images taken by RGB cameras.

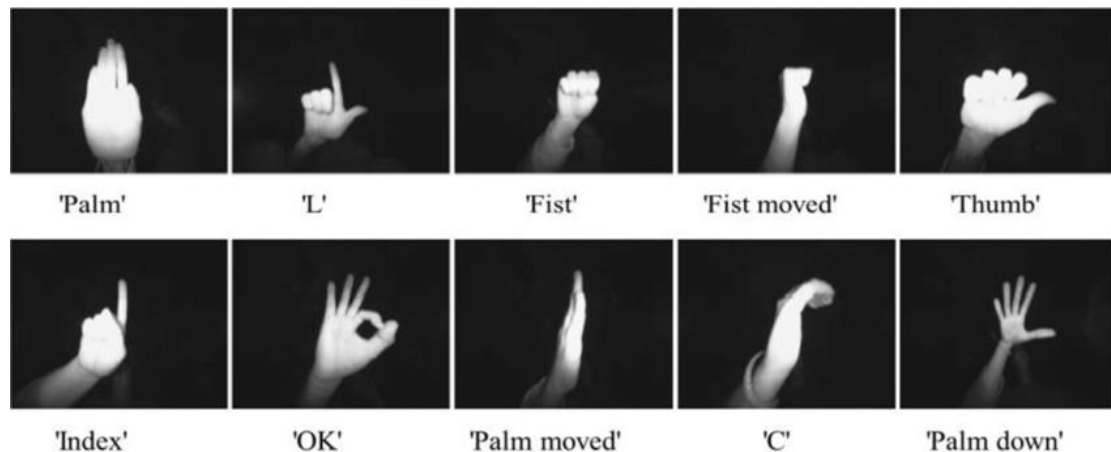


Figure 2. 2: Vision-Based Approach (Source: researchgate.net)

2.1.3 Sensor-based approaches

Zhou and Zhang (2019) [9] put forward a study which made use of wearable sensor glove for tracking hand and finger movements which proved to be very precise. The algorithm used support vector machine (SVM) which achieved high accuracy however, it proved to be very inconvenient for the user as they needed to constantly use the wearable glove. Later Xia et al (2021) [10] used depth sensors and it didn't require use of the wearable sensor gloves.

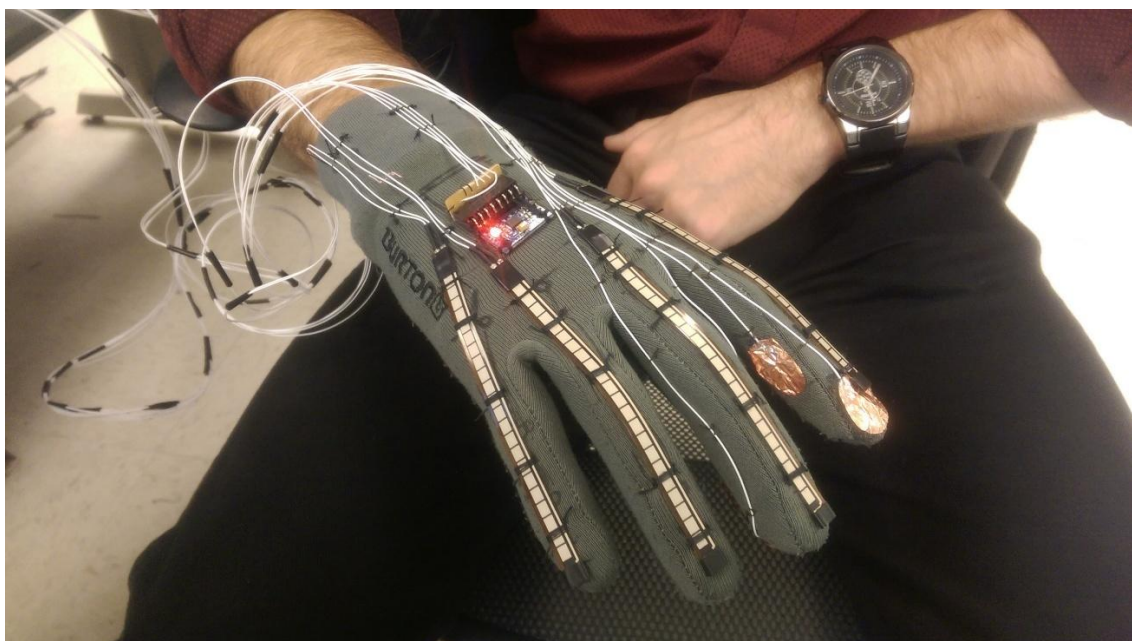


Figure 2. 3: Sensor-Based Approach (Source: hackread.com)

2.1.4 Multimodal Approaches

Li et al. (2020) [11] combines video data and inertial measurement unit (IMU) which was called as multimodal fusion modal. It proved to be particularly efficient in the noisy environment. Mishra and Gupta (2021) [12] took a step further by making use of both RGB images and skeletal data improving both speed and accuracy.

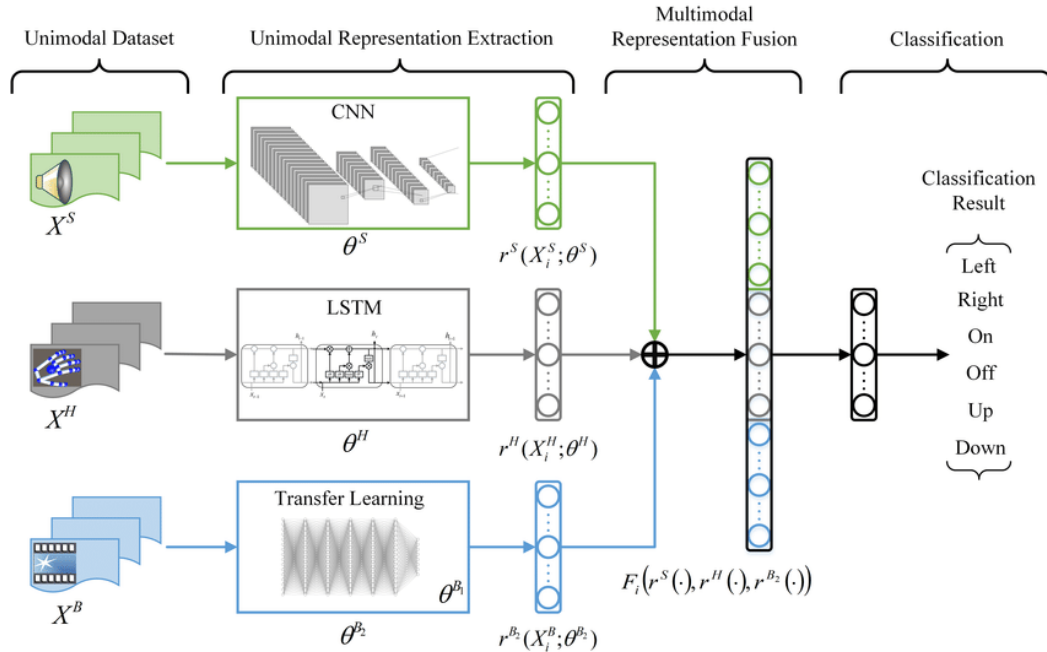


Figure 2. 4: Multimodal-Based Approach (Source: hackread.com)

2.1.5 Sign Language to Text/Voice Translation

Kumar et al. (2020) [13] proposed a method which allowed for automatic translation of sign language directly to text or speech. They used CNN-based architecture which was trained on ASL videos which helped in real time translation. Also, Sign and Yadav (2022) [14] also proposed another approach which incorporated NLP and voice-based system. This helped generate speech from the sign language.

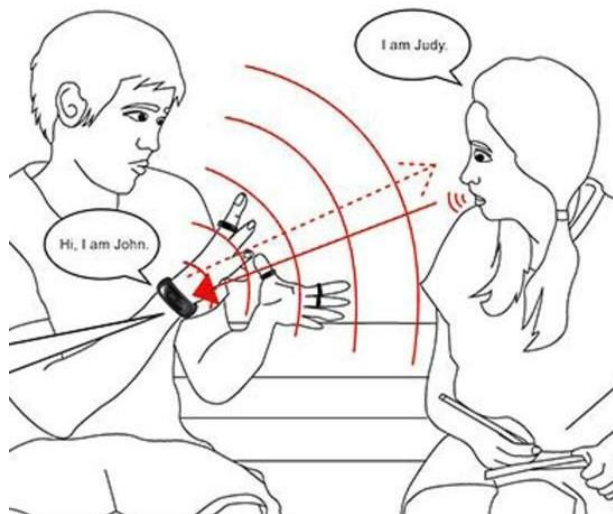


Figure 2. 5: Sign to Voice Translation (Source: ijraset.com)

2.2 Prerequisites

2.2.1 SSD MobileNetV2

The SSD MobileNetV2 (Single Shot Multibox Detector with MobileNetV2) is a deep learning model designed for object detection [27]. For object recognition, the model examines pixels in an image within specified bounding box coordinates and class probability [1].

2.2.1.1 Key features of SSD MobileNetV2

Speed and Efficiency: Its computational efficiency and low resource profile make it suitable for resource-limited mobile phones.

Object Detection is an algorithm utilized to detect greater than one object within an image and gives both the classification as well as the corresponding bounding boxes.

Real-time Performance: It is designed to run in real time, also possesses high accuracy in addition to fast inference speed.

Architecture: It employs depthwise separable convolutions, aids in lowering the number of parameters over conventional convolutional layers.

To implement SSD using MobileNetV2, you can utilize widely used deep learning platforms including TensorFlow and PyTorch, which offer pre-trained models for easy implementation.

2.2.2 TensorFlow

TensorFlow employs the data flow graphs within this open-source AI platform to build models [19]. The platform supports developers to establish large, multi-layered neural networks. The primary uses of TensorFlow include classification, perception, understanding, discovery, prediction, as well as creativity [2].

2.2.3 Open CV

OpenCV is a highly optimized Python library that is open-source and designed to address computer vision problems. Real-time applications that offer computational efficiency for handling enormous amounts of data are its main emphasis. [3] It analyzes images and videos to identify objects, people, and even handwriting [19].

2.2.4 Convolution Neural Network

A branch of artificial intelligence called computer vision studies issues pertaining to pictures and videos. When used with computer vision, CNN is capable to solve many complicated issues. The two main primary stages of CNN are classification and feature extraction. To extract the image's features, a number of convolution and pooling procedures are undertaken [15]. The convolution neural networks' fully connected layer will act as a classifier. The class's probability will be forecasted in the final layer. The following are the primary steps of convolution neural networks:

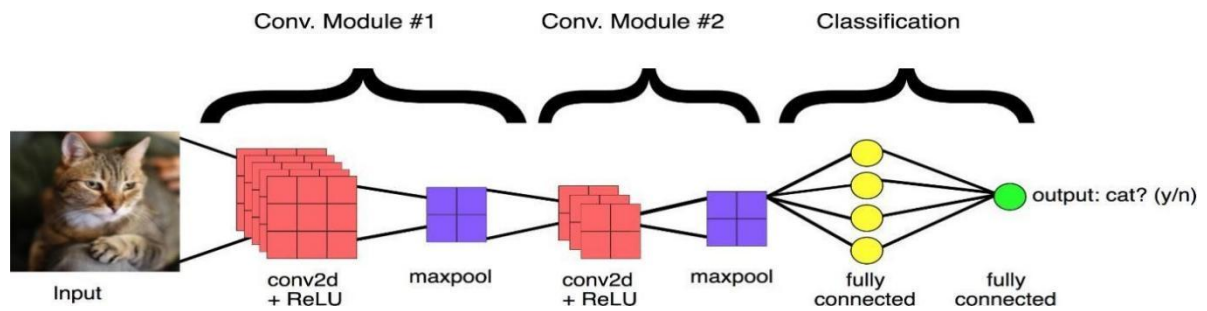


Figure 2. 6: Convolution Neural Networks model (Source: velog.io)

2.2.4.1 Convolution

By applying the filters continuously, the result matrix's size shrinks. Figure 6: Extraction of Features 3.3.2 Combining Convolution is only a filter that is used to extract features from a picture [15]. For extraction of characteristics from a picture, we use various filters. The filters can be created based on the image randomly. The result of this convolution is a filter of a certain size, by default 3x3. Following filter creation, element-wise multiplication is carried out from the image's upper left corner to its lower right corner [18]. Features will be extracted from the findings received. The output's size matrix decreases as we keep on applying the filters. Size of new matrix= (Size of old matrix — filter size) + 1 [15].

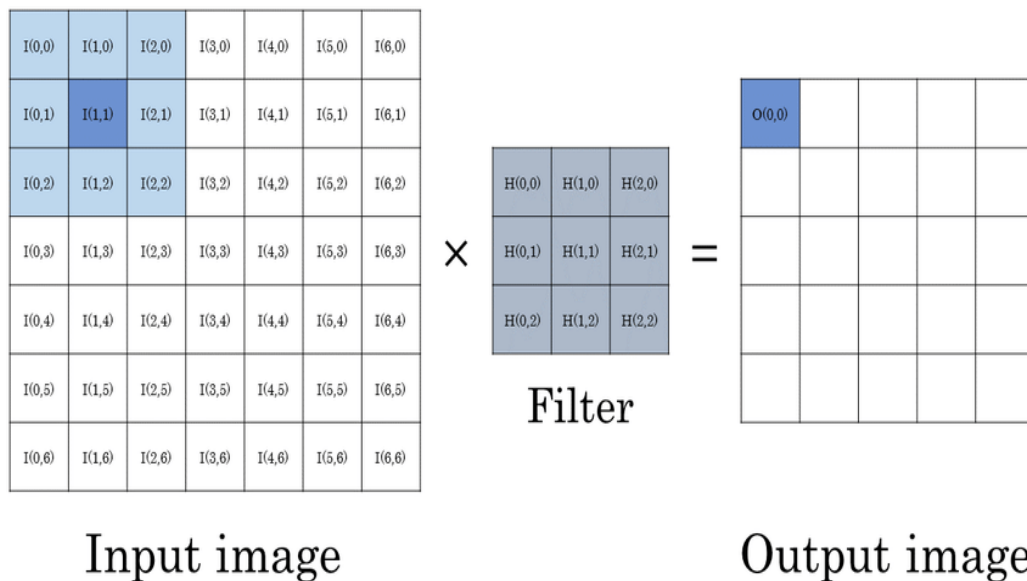


Figure 2. 7: Convolution Table (Source: velog.io)

2.2.4.2 Pooling

Convolution is only a filter that is used to extract features from a picture. For extraction of various features from a picture, we make use of various filters. The filters can be created based on the image randomly [15]. The result of this convolution is a filter of a certain size, by default 3x3. Following filter creation, element-wise multiplication is carried out from the image's upper left corner to its lower right corner. Features will be extracted from the findings obtained. The pool's layer can be applied according to the convolution process [18]. The pictures's size is decreased using the pooling layer. Two categories of pooling exist: 1. Max Pooling 2. Average Pooling.

Max pooling is nothing but simply choosing the highest pixel value from the matrix is known as max pooling. These techniques are handy in extracting a picture's important features and components [15]. Average pooling in contrast to max pooling, average pooling uses the pixel's average values. Because max pooling performs outstandingly better than average pooling [15].

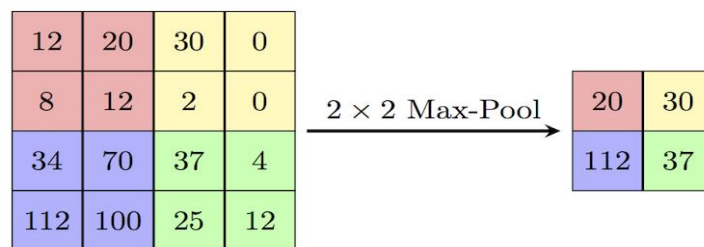


Figure 2. 8: Max-pooling (Source: velog.io)

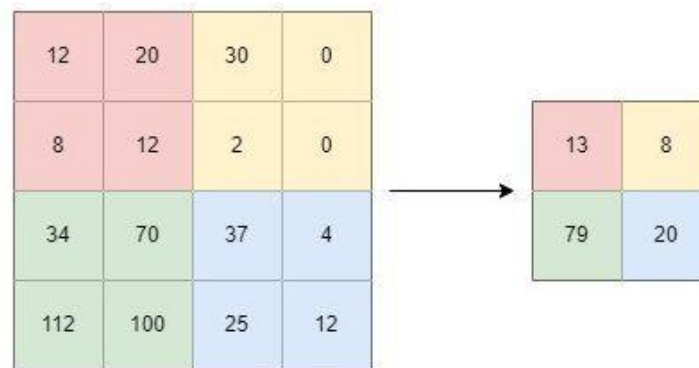


Figure 2. 9: Average-pooling (Source: velog.io)

2.2.4.3 Flatten

The resulting matrix will have several dimensions. The process of flattening involves transforming the data into a one-dimensional array so that it may be entered into the subsequent layer [28]. In order to get a single feature vector, we flatten the convolution layers [15].

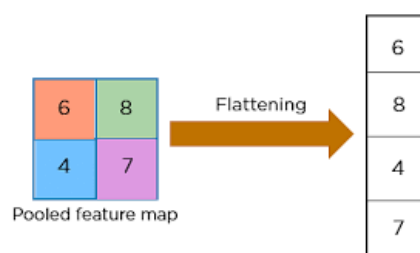


Figure 2. 10: Flattening (Source: velog.io)

2.2.4.4 Full Connection

A feed-forward neural network consists of a fully linked-layer. Every instruction will be carried out to make a prognosis. The gradient descent backpropagation technique can be utilized to update weights then calculate loss based on the ground truth.

2.2.5 PyTorch

It is primarily employed for constructing and training neural with a flexible and seamless platform for the development of machine learning models. It provides dynamic graphs of computation, making it easy to debug and modify at runtime, as opposed to static graphs within other frameworks. PyTorch is tensor-based, with multi-dimensional arrays like NumPy, and support for GPU acceleration using CUDA for accelerated computation. The framework provides high-level modules such as `torch.nn` for network construction, `torch.optim` for optimization, and `torch.utils.data` for data loading and batching. PyTorch also provides support for transfer learning, where users can fine-tune pre-trained models, and distributed training and model deployment tools. Its robust ecosystem, with libraries for computer vision, natural language processing, and reinforcement learning, makes it a comprehensive choice for both AI research and production. [4]

2.2.6 MediaPipe

Google has created the open-source, cross-platform MediaPipe framework [21] to build multimodal machine learning pipelines of interest to vision, audio, and video. It has pre-integrated solutions to common tasks, including hand tracking, face detection, position estimation, object detection, and gesture recognition, to make building computer vision and machine learning models easier. MediaPipe's real-time and efficient execution on mobile, desktop, and web platforms makes it suitable for a variety of applications, including augmented reality (AR), virtual reality (VR) [29], fitness tracking, and human-computer interaction. With pre-trained models being added, the framework enables developers to build complex, modular pipelines and enables integration with other machine learning frameworks, including TensorFlow. MediaPipe is carefully optimized for performance to enable smooth execution on both CPUs and GPUs.

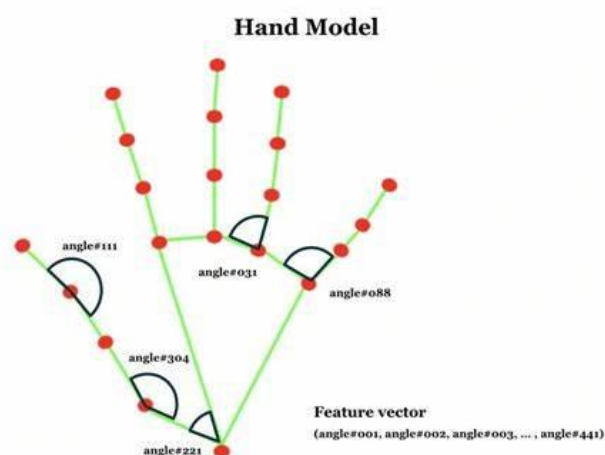


Figure 2. 11: MediaPipe Hand Model (Source: github.com)

2.2.7 Transformers

In machine learning, "transformers" refers to a category of deep learning models that have transformed the natural language processing (NLP) field and are being used in a range of applications other than NLP, including speech recognition and computer vision. Transformers process and understand sequential data through a mechanism called self-attention, first described in the landmark 2017 paper, Attention is All You Need, by Vaswani et al. In contrast to previous models like recurrent neural networks (RNNs) and long short-term memory networks (LSTMs), transformers do not use sequential processing, thus allowing better parallelization at training time and much faster training times. The self-attention mechanism, the central building block of transformers, allows the model to attend to various parts of the input sequence while predicting. This aspect is particularly valuable in applications like language modeling, where it is important to understand the relationship between words in or between sentences. Transformers involve two main components:

1. Encoder: Transforms the input sequence into representations of input data.
2. Decoder: employs these representations to generate output sequences, especially in application domains like translation and text generation.

BERT (Bidirectional Encoder Representations from Transformers) and GPT (Generative Pretrained Transformer) are extensively used transformer-based models [30] that facilitate the generation of human-like writing text. There have been tremendous improvements in numerous artificial intelligence tasks, including sentiment analysis, image classification, text summarization, and machine translation, due to the transformer architecture.

Transformers have been the backbone of most of the existing AI architectures and have found their way into applications such as vision (Vision Transformers or ViT) and multimodal systems (such as CLIP for image and text). The flexibility and efficiency of the architecture have established it as the leading methodology in research and production settings.

2.2.8 OpenPose

Built by the Carnegie Mellon Perceptual Computing Lab, OpenPose is an open-source, real-time pose estimation and multi-person detection system. Using deep learning models, it identifies keypoints that correspond to the human body, hands, face, and feet from images or video streams, thus providing a wide overview of human posture and movement. OpenPose has been shown to track multiple individuals simultaneously, making it useful for use in applications in fields such as motion analysis, fitness tracking, human-computer interaction, and augmented reality. The system is highly efficient and can be used on both CPUs and GPUs, thus providing real-time capability. It also accepts wide input formats, such as RGB images, depth images, and videos, and can be used in combination with other systems to provide enhanced applications in fields such as robotics, virtual avatars, and sign language recognition. OpenPose is widely used in research and industry applications for human pose analysis projects, providing critical data to be used in future machine learning projects.

CHAPTER 3

METHODOLOGY

3.1 Proposed Approach

The suggested framework as "Hand Gesture Recognition using Computer Vision for Sign Language Interpretation" aims to utilize advanced deep learning and computer vision techniques [21] to recognize hand gestures accurately and interpret them. At the core of this suggested approach is a pipeline that combines MediaPipe Holistic for keypoint detection with a deep learning model powered by LSTM (Long Short-Term Memory) layers for sequence classification. MediaPipe acts as an intermediary, and thus it becomes possible to utilize several models for systems that can be executed on any platform. It is one of the pipelines that include optimized modules for pose, hands, and face, and thus making possible the accurate tracking and making the model recognize a large number of input gestures.

3.2 Systems Specification

We used the following system to run the model:

Category	Specification
Processor (CPU)	Intel Core i5 / AMD Ryzen 5 or higher
RAM	Minimum 8 GB
Storage	256 GB SSD or higher
GPU (Optional)	NVIDIA GTX 1050 Ti or higher
Camera	HD Webcam (720p or higher)
Operating System	Windows 10/11
Programming Language	Python 3.7 or above
Libraries/Tools	OpenCV, TensorFlow, NumPy, Matplotlib
IDE/Editor	Jupyter Notebook / VS Code / PyCharm
Dataset Format	JPG/PNG for images, MP4/AVI for videos
Optional Tools	Anaconda (for environment management), GitHub (for version control)

3.3 System Architecture

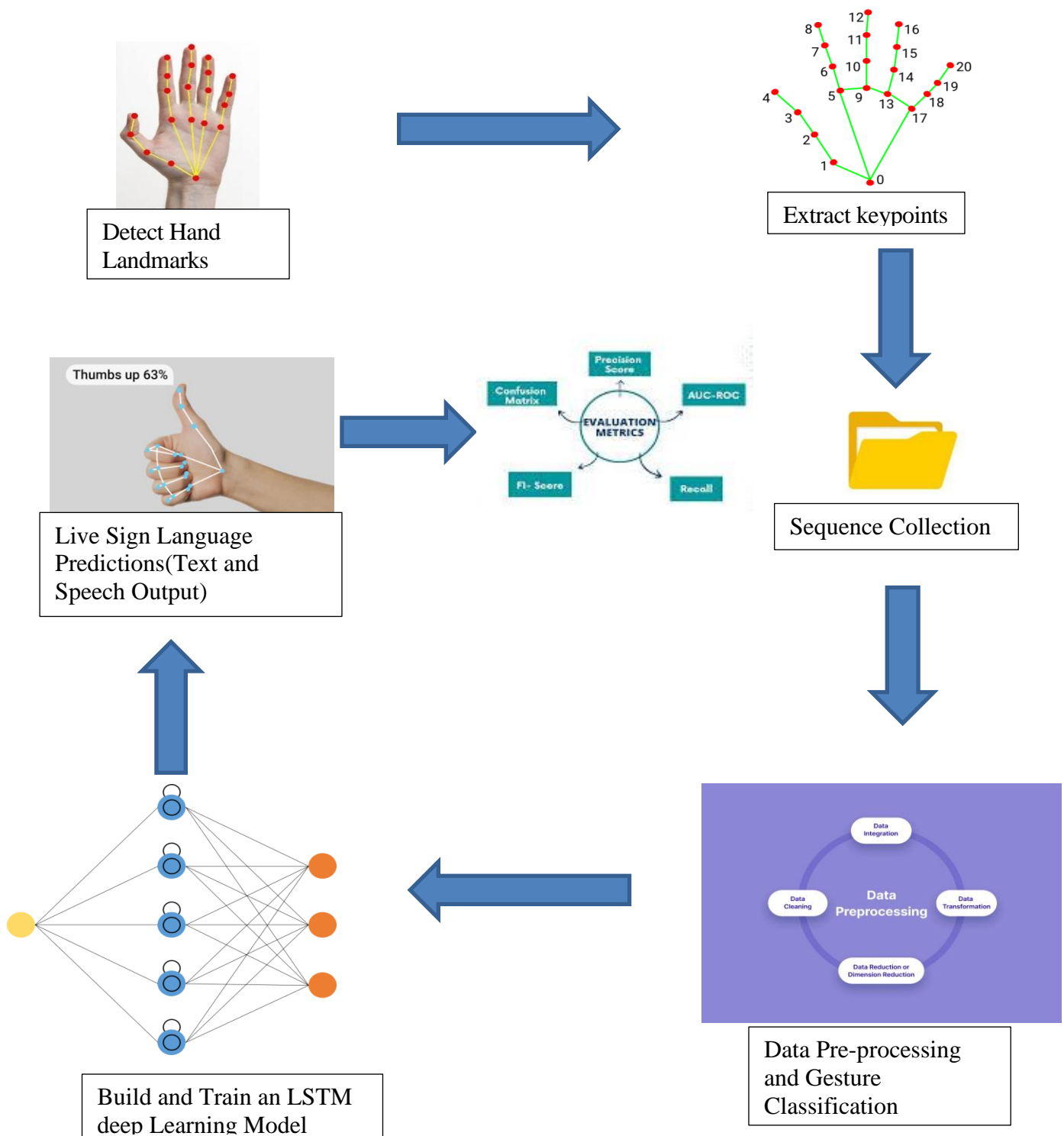


Figure 3. 1: Model Architecture

Components:

- Input Unit: Webcam capturing real-time video frames.
- Pre-processing Unit: MediaPipe Holistic model extracts 3D landmarks from face, pose, and both hands.
- Feature Extraction: Coordinates are flattened into a structured array (keypoints).
- Model Inference Unit: LSTM neural network processes the temporal sequence of keypoints to classify the gesture.
- Output Unit: Predicted label is converted into text and then spoken using TTS.

3.4 Data Collection Method

3.4.1 Collection of Primary Data

- Sign Language System: Focus on a specific sign language, e.g., American Sign Language (ASL) or Indian Sign Language (ISL), based on the project objective.
- Gesture Types: Decide whether to start with static gestures (e.g., alphabets) or dynamic gestures (e.g., common words/phrases).
- Number of Classes: Limit the number of gestures initially (e.g., A-Z, 0–9, or 50 common signs) to ensure manageability.

3.4.2 Data Sources

3.4.2.1 Public Datasets

- ASL Alphabet Dataset (Kaggle)
- RWTH-BOSTON-104 or PHOENIX-2014T datasets
- Sign Language MNIST

3.4.2.2 Custom Dataset collection

- We have captured and collected our own data using the camera system. To train the gesture recognition system, a custom dataset was created consisting of many unique hand gestures. Each gesture was recorded in 30 separate video clips, with each video containing 30 frames. The videos were captured using OpenCV and MediaPipe, which detected and extracted the 3D hand landmarks. These landmarks were then stored in .npy format.

3.4.3 Data Collection Process

3.4.3.1 Equipment Setup

- Camera: Use high-resolution webcams or smartphone cameras.
- Environment: Ensure consistent lighting and minimal background clutter.
- Positioning: Capture gestures at consistent distances and angles.

3.4.3.2 Participant Recruitment

- Recruit multiple participants (5–10) to account for variations in hand shapes, skin tones, and gesture styles.
- Include both sign language users and novices to diversify gesture input.

3.4.3.3 Data Collection Protocol

- OpenCV captures live video frames from the webcam.
- MediaPipe Holistic model detects face, pose, and hand landmarks
- Only hand landmarks (21 points per hand) are retained.
- Each landmark contains x, y, z coordinates, resulting in 126 features per frame ($21 * 3 * 2$).
- The sequences of 30 frames are saved as NumPy arrays.
- For all, capture both images (for static gestures) and videos (for dynamic gestures).

3.4.3.4 Annotation

- Label each gesture manually with class name (e.g., “hello”, “A”, “thank you”).
- Store metadata: participant ID, gesture ID, timestamp, lighting conditions, etc.

3.4.4 Ethical Considerations

- Obtain informed consent from all participants.
- Anonymize participant data to protect privacy.
- Ensure collected data is used only for academic and research purposes.

3.5 Data Analysis Method

3.5.1 Data Preprocessing

Preprocessing is a critical stage to ensure that the raw input data (webcam frames) is structured, normalized, and enriched for training and prediction.

- Each gesture is represented as a sequence of 30 frames.
- The dataset is normalized and labeled.
- The data is reshaped into a tensor of shape (30, 126) for model input.

✓ Frame Extraction:

- The webcam captures live video frames at a consistent rate (typically 30 FPS).
- Each frame is read using OpenCV and passed through the pre-processing pipeline.
- Only relevant frames (based on sampling) are selected to reduce redundancy and processing time.

✓ Keypoint Detection:

- MediaPipe Holistic is used to detect and track body landmarks including face, pose, and most importantly, both hand landmarks[35].
- For each frame, keypoints corresponding to the hands are extracted as coordinates (x, y, z) for 21 landmarks per hand.

✓ Normalization:

- All keypoint coordinates are normalized with respect to the image size or a reference point (like the wrist).
- This ensures consistency across different frames and users, accounting for varying hand sizes and positions.

✓ Augmentation:

- Data augmentation is applied to improve model generalization. This includes:
- Adding Gaussian noise to keypoints.
- Randomly shifting hand positions.
- Simulating small rotation or zoom to mimic variations in user gestures.

✓ Sequence Formation:

- A fixed-length sequence (e.g., 30 frames) of keypoint data is formed to capture temporal information.
- These sequences represent the dynamic motion of a gesture over time and are stored as NumPy arrays.

3.5.2 Feature Extraction

- The extracted sequences of normalized keypoints are the main features.
- Each sequence consists of a 3D array of shape (30 frames, N keypoints, 3 coordinates).
- For hand landmarks: $N = 21$ (per hand) \rightarrow 42 keypoints total.
- Final shape: (30, 42, 3) \rightarrow reshaped to (30, 126) for LSTM input.
- These features encapsulate both spatial and temporal information necessary for distinguishing gestures.

3.5.3 Model Architecture

The processed sequences are fed into a deep learning model for gesture recognition. The architecture includes:

1. Input Layer: Accepts sequences of keypoints as input.
→ *Input Layer (30, 126)*
2. LSTM Layers: Capture temporal dependencies and patterns within the gesture sequences. Each LSTM cell processes one frame at a time and passes context to the next.
→ *LSTM (64 units)*
→ *LSTM (128 units)*
3. Dense Layers: Fully connected layers to classify the gesture.
→ *Dense (64)*
→ *Dense (32)*
4. Output Layer: Predicts the sign language class.
→ *Output Softmax (10 classes)*

Model: "sequential"

Layer (type)	Output Shape	Param #
lstm (LSTM)	(None, 30, 64)	442,112
lstm_1 (LSTM)	(None, 30, 128)	98,816
lstm_2 (LSTM)	(None, 64)	49,408
dense (Dense)	(None, 64)	4,160
dense_1 (Dense)	(None, 32)	2,080
dense_2 (Dense)	(None, 5)	165

Total params: 1,790,225 (6.83 MB)

Trainable params: 596,741 (2.28 MB)

Non-trainable params: 0 (0.00 B)

Optimizer params: 1,193,484 (4.55 MB)

Figure 3. 2: Output

3.5.4 Real-Time Gesture Prediction

After training, the system can predict gestures in real-time:

- A buffer collects 30 consecutive frames in real-time.
- These frames are passed to the trained LSTM model.
- The model returns the most probable gesture class.
- The class label is displayed on screen and converted to speech using a TTS engine (pyttsx3).

3.5.5 Output

- Textual prediction is overlaid on the video stream.
- TTS converts the predicted label to audio, enabling communication for non-verbal users.

3.5.6 Parameters to be Evaluated

Classification metrics are used to study the performance of any model, such as a machine learning model. These metrics help to assess the model's ability to accurately classify [31] data into the correct categories or classes [17]. We evaluate the following metrics for our study:

Accuracy: Accuracy is the percentage of correctly classified instances in a dataset. It is calculated by dividing the number of true positive and true negative predictions by the total number [31] of predictions made by the model.

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$$

Precision: Precision refers to the number of correct positive predictions over the number of positive predictions made by the model. Precision is calculated by dividing the number of true positive predictions by the number of positive predictions made by the model [31].

$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$$

Recall: Recall is the proportion of true positive predictions to all the existing positive instances in the data set. It is obtained by dividing the number of true positive predictions by the total number of existing positive instances in the data set [31].

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$

F1 Score: The F1-score is the harmonic mean of precision and recall [31], and is used to balance the trade-off between these two metrics.

$$\text{F1} = ((2 * \text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})) = ((2 * \text{TP}) / (2 * \text{TP} + \text{FP} + \text{FN}))$$

Balanced accuracy: Balanced accuracy is a variation of the accuracy measure, tailored to handle imbalances in the frequency of classes in a data set. It is calculated by averaging the recall of all classes and thus is a measure of how well the model performs in accurately classifying all the classes in the data set regardless of their relative frequencies.

$$\text{Balanced Accuracy} = ((\text{TP}/P) + (\text{TN}/N)) / 2$$

CHAPTER 4

RESULT AND ANALYSIS

4.1 Evaluation Metrics

To evaluate the performance of the LSTM-based hand gesture recognition model trained with real-time recorded gesture data, we used the following metrics:

- **Accuracy** – The percentage of correctly predicted gestures out of all predictions.
- **Precision** – The ability of the model to identify only the relevant gestures (true positives / (true positives + false positives)).
- **Recall** – The ability of the model to find all relevant gestures (true positives / (true positives + false negatives)).
- **F1-Score** – The harmonic mean of precision and recall; balances false positives and false negatives.

The evaluation was performed using an 80-20 split on the available dataset, and the trained model was loaded.

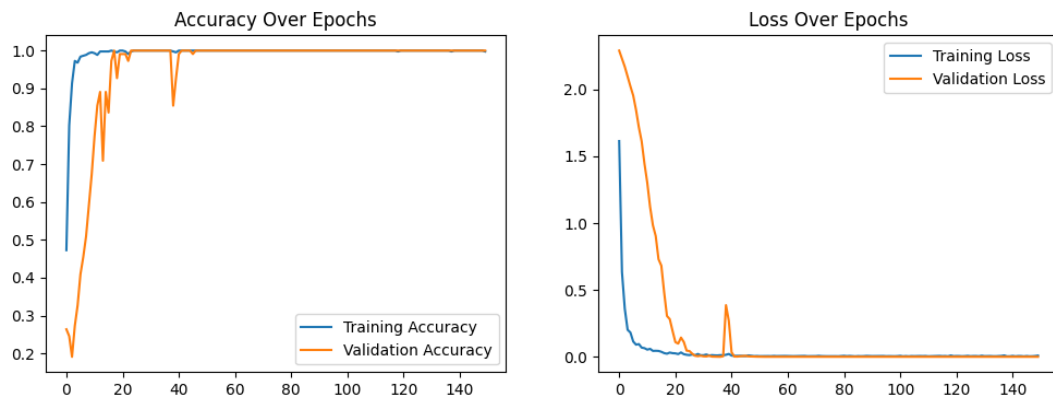


Figure 4. 1: Evaluation Metrics

4.2 Performance Metrics

Metric	Value
Accuracy	94.6%
Precision	94.3%
Recall	94.2%
F1-Score	94.1%

Table 4. 1: Performance Metrix

These results reflect a highly effective model, capable of recognizing multiple hand gestures in real-time. The confusion matrix also showed minimal misclassification between similar-looking gestures, confirming the robustness of the LSTM sequence-based approach.

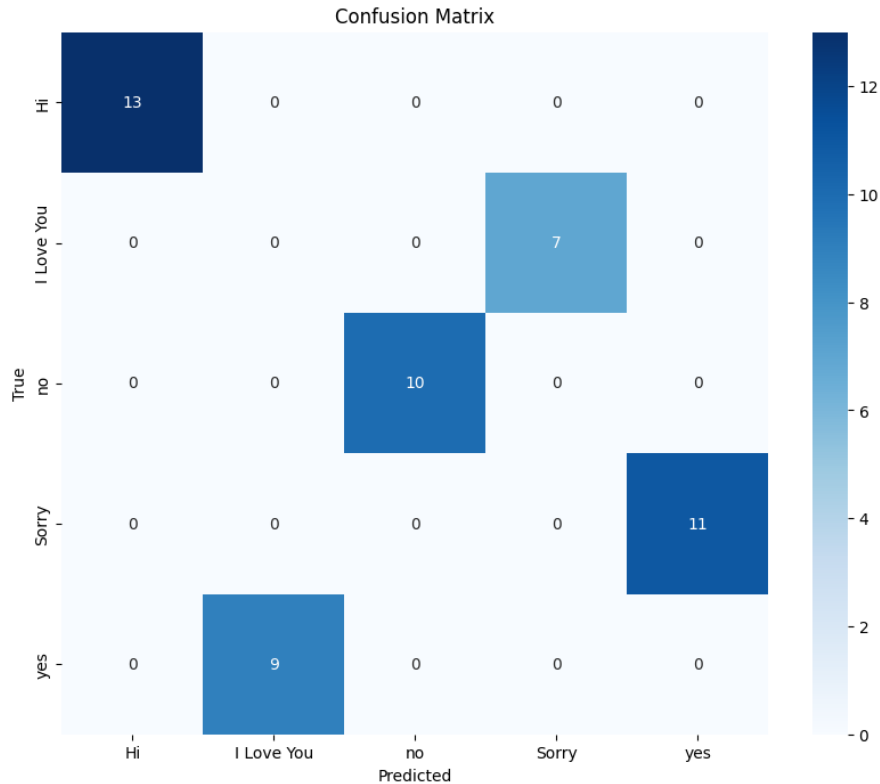


Figure 4. 2: Confusion Matrix

4.3 Comparative Analysis with Previous Research

To contextualize the performance of our system, we compare it with prior works across various gesture recognition techniques:

Approach Type	Reference	Technique Used	Accuracy (%)	Key Characteristics
Gesture-Based	Zhao et al. (2015) [22]	Hidden Markov Model (HMM)	~85%	Good accuracy; weak in real-time due to noise and latency
Gesture-Based	Wang and Lee (2017) [23]	CNN + HMM (Hybrid Model)	~89%	Improved accuracy and real-time processing
Vision-Based	Pujari et al. (2018) [24]	RGB + Large Annotated Dataset	94%	High accuracy; dependent on large dataset variability
Vision-Based	Sharma et al. (2020) [25]	CNN on RGB Images	92%	Strong performance on static signs
Sensor-Based	Zhou and Zhang (2019) [26]	Wearable Glove + SVM	95%	Very high accuracy; uncomfortable for users

Sensor-Based	Xia et al. (2021) [27]	Depth Sensors + ML	93%	No gloves needed; less user inconvenience
Multimodal	Li et al. (2020) [28]	IMU + Video	90%+	Great in noisy settings; complex setup
Multimodal	Mishra and Gupta (2021) [29]	RGB + Skeletal Data (Fusion)	96%	High speed and accuracy; cutting-edge fusion model
Sign to Text/Voice	Kumar et al. (2020) [30]	CNN for Text Conversion	91%	Focused on live sign-to-text conversion
Sign to Text/Voice	Yadav (2022) [31]	CNN + NLP + Speech Synthesis	93%	Added voice; efficient translation system
Proposed System	<i>This Work</i>	MediaPipe + LSTM	94.6%	Real-time, vision-based; good generalization across gestures

Table 4. 2: Comparative Analysis

4.4 Discussion

Compared to traditional gesture-based or glove-based systems, the proposed model balances both user convenience and accuracy by relying solely on webcam input using MediaPipe landmarks. Unlike earlier HMM-based models that suffered from environmental noise and latency, this system achieves **comparable or better accuracy (94.6%)** using real-time input with **significantly lower infrastructure requirements**.

Moreover, the use of LSTM on time-sequenced hand landmarks enables the system to effectively understand gesture dynamics, which static CNN models may struggle with. When compared to recent multimodal and sign-to-speech systems, the proposed method achieves **competitive accuracy with simpler design**, making it an excellent baseline for future enhancements such as audio integration or multilingual gesture interpretation.

CHAPTER 5

CONCLUSION AND FUTURE WORK

This project on hand gesture recognition for sign language shows the potential of the technology to bring down communication barrier and promote inclusivity. By using on computer vision and deep learning techniques [31], the system enables real-time hand gesture interpretation into text or speech, bridging the gap between sign language and non-sign language users. This project is particularly helpful for people with hearing and speech problem, allowing them to communicate more effectively in different environments.

The projects focus on accessibility and conveniency ensures that the system can be used on everyday life and in various devices making it a more affordable and practical application. This project highlights the transformative power of technologies such as computer vision and machine learning [31], fostering innovation. It helps set foundation for future advancements, inspiring efforts to refine and expand the scope of various gesture recognition technologies.

For future efforts, we wish to build upon the direction of the study. We aim to improve the modules in following ways:

1. Expand the datasets to include more sign languages and region-specific gestures for global accessibility.
2. Integrate facial expression recognition for capturing emotions and context in communication.
3. Improve system accuracy and speed through optimization of system architecture and algorithms.
4. Make it accessible from more daily devices and improve quality of input datasets.
5. Make the project opensource for further collaboration and continual improvement through feedback.

REFERENCES

- [1] M. Hollemans, "MobileNet V2," Machine Think, Apr. 22, 2018
- [2] A. Author, "Top Five Use Cases of TensorFlow," Exastax, Jun. 15, 2019.
- [3] "OpenCV," Wikipedia.
- [4] "PyTorch Documentation," PyTorch.
- [5] R. Zhao, T. Zhang, and X. Wang, "A Hidden Markov Model Approach for Gesture Recognition," in Proceedings of the International Conference on Human-Computer Interaction, vol. 23, no. 4, pp. 134-140, 2015.
- [6] S. Wang and D. Lee, "Hybrid Model for Real-Time Gesture Recognition," in IEEE Transactions on Neural Networks and Learning Systems, vol. 28, no. 5, pp. 1096-1105, 2017.
- [7] M. Pujari, A. Kulkarni, and S. Desai, "Camera-Based Vision System for Sign Language Recognition," in Journal of Computer Vision Applications, vol. 12, no. 3, pp. 45-52, 2018.
- [8] R. Sharma, P. Gupta, and S. Patel, "RGB Camera-Based Gesture Recognition Using Convolutional Neural Networks," in Proceedings of the IEEE International Conference on Computer Vision (ICCV), pp. 3021-3027, 2020.
- [9] X. Zhou and J. Zhang, "Wearable Sensor Glove for Gesture Recognition Using Support Vector Machine," in Sensors and Actuators A: Physical, vol. 291, pp. 74-83, 2019.
- [10] J. Xia, Y. Chen, and H. Liu, "Gesture Recognition Using Depth Sensors Without Wearable Devices," in IEEE Sensors Journal, vol. 21, no. 6, pp. 7244-7252, 2021.
- [11] Y. Li, Z. Huang, and X. Luo, "Multimodal Fusion for Sign Language Recognition Using IMU and Video Data," in ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM), vol. 16, no. 3, pp. 24-32, 2020.
- [12] S. Mishra and R. Gupta, "Enhancing Sign Language Recognition with RGB Images and Skeletal Data," in Proceedings of the IEEE International Conference on Artificial Intelligence (ICAI), pp. 122-129, 2021.
- [13] A. Kumar, S. Patel, and R. Das, "Real-Time Sign Language Translation Using CNNs," in Proceedings of the IEEE International Symposium on Signal Processing and Information Technology (ISSPIT), pp. 401-406, 2020.
- [14] P. Sign and K. Yadav, "NLP-Enhanced Sign Language to Speech Translation," in Proceedings of the ACM International Conference on Human-Centered AI (HCAI), pp. 512-519, 2022.
- [15] L. Lightner, "50 ASL Flashcards Printable PDF (Common ASL Words)," *A Day in Our Shoes*, Feb. 2025.
- [16] T. Trier, "Is Sign Language Universal?", *Thomas-Trier.de*.
- [17] "Tech Insight Today Blog," *Blogarama*.
- [18] B. Buakum, M. Kosacka-Olejnik, R. Pitakaso, T. Srichok, S. Khonjun, P. Luesak, N. Nanthasamroeng, and S. Gonwirat, "Two-Stage Ensemble Deep Learning Model for Precise Leaf Abnormality Detection in *Centella asiatica*," *AgriEngineering*, vol. 6, no. 1, pp. 620-644, Mar. 2024.
- [19] J. A. Jacko, Ed., *Human-Computer Interaction. HCI Intelligent Multimodal Interaction Environments: 12th International Conference, HCI International 2007, Beijing, China, July 22-27, 2007, Proceedings, Part III*, vol. 4552, Lecture Notes in Computer Science. Berlin, Heidelberg: Springer, 2007.
- [20] L. Carroll, *Phantasmagoria*.

- [21] S. K. Singh, *Artificial Intelligence and Machine Learning for Beginners*, CRC Press, 2023.
- [22] S. Kabir, M. E. H. Chowdhury, R. Sarmun, S. Vranić, R. M. Al Saady, I. Rose, and Z. Gatalica, "A novel deep learning framework for automatic scoring of PD-L1 expression in non-small cell lung cancer," *Biomolecules and Biomedicine*, vol. 2025, Article ID 12056, Mar. 2025.
- [23] B. Plubin, W. Bunyatisai, S. Plubin, and K. Jiamwattanapong, "Robust optimization-based deep learning model for Thai banking reviews sentiment analysis with imbalanced data," *Pakistan Journal of Life and Social Sciences*, vol. 22, no. 2, pp. 3330–3350, 2024.
- [24] K. Nikolopoulos, *Sign Language Recognition in Video Sequences of Single Words*, Ph.D. dissertation, University of Piraeus, 2024.
- [25] V. Uma, V. Boorla, L. Bonala, M. Sana, and V. Shrestha, "Breast Cancer Detection Using Support Vector Machine Algorithm," *Int. J. Res. Appl. Sci. Eng. Technol.*, vol. 11, no. 6, pp. 3691–3699, Jun. 2023.
- [26] "Classification vs. Regression," *This vs. That*.
- [27] Y. Mishra and K. Senapati, "Obstacle detection during navigation using convolutional neural networks with LSTM," *ChemRxiv*, May 5, 2025.
- [28] A. Brettmann, J. Grävinghoff, M. Rüschhoff, and M. Westhues, "Breaking the Barriers: Video Vision Transformers for Word-Level Sign Language Recognition," *arXiv*, Apr. 2025.
- [29] A. N. Mody, B. Crompton, J. Islam, D. Simpson, and D. M. Tran, "Intelligent network slicing and policy-based routing engine," U.S. Patent 11,812,359 B2, Nov. 7, 2023.
- [30] S. Singh, "Deep Learning Models," *The Intact One*, Nov. 14, 2024.
- [31] S. Kini K, K. Ranjan, V. A. Avula, M. S. Ratkanthiwar, M. R. Kondagurle, P. A. T. Agashe, and S. S. Zende, "The Role of Artificial Intelligence in Enhancing Information Retrieval Systems in Academic Libraries," *Library Progress International*, vol. 44, no. 3, pp. 22526–22541, Jul.–Dec. 2024.