

**DESIGN OF AUGMENTATIVE AND ALTERNATIVE
COMMUNICATION IN INDIAN CONTEXT**

MINOR PROJECT-1 REPORT

Submitted by

RANJANI .S

ROHIT KUMAR

SAI SUNNYHITH .B

Under the Guidance of

Dr. KOUSHICK VENKATESH

in partial fulfillment for the award of the degree

of

BACHELOR OF TECHNOLOGY

in

ELECTRONICS & COMMUNICATION ENGINEERING



MAY 2024



Vel Tech

Rangarajan Dr. Sagunthala
R&D Institute of Science and Technology

(Deemed to be University Estd. u/s 3 of UGC Act, 1956)

BONAFIDE CERTIFICATE

Certified that this Minor project-1 report entitled "**Design of Augmentative and Alternative Communication in Indian Context**" is the bonafide work of "**Rohit Kumar (21UEEA0110), Ranjani .S (21UEEA0108) and B. Sai Sunnyhith (21UEEL0010)**" who carried out the project work under my supervision.

SUPERVISOR

Dr. KOUSHICK VENKATESH

Assistant Professor

Department of ECE

HEAD OF THE DEPARTMENT

Dr.A. SELWIN MICH PRIYADHARSON

Professor

Department of ECE

Submitted for Minor project-1 work viva-voce examination held on:-----

INTERNAL EXAMINER

EXTERNAL EXAMINER

ACKNOWLEDGEMENT

We express our deepest gratitude to our Respected Founder President and Chancellor **Col. Prof. Dr. R. Rangarajan**, Foundress President **Dr. R. Sagunthala Rangarajan**, Chairperson and Managing Trustee and Vice President.

We are very thankful to our beloved Vice Chancellor **Prof. Dr. S. Salivahanan** for providing us with an environment to complete the work successfully.

We are obligated to our beloved Registrar **Dr. E. Kannan** for providing immense support in all our endeavours. We are thankful to our esteemed Dean Academics **Dr. A. T. Ravichandran** for providing a wonderful environment to complete our work successfully.

We are extremely thankful and pay my gratitude to our Dean SoEC **Dr. R. S. Valarmathi** for her valuable guidance and support on completion of this project.

It is a great pleasure for us to acknowledge the assistance and contributions of our Head of the Department **Dr. A. Selwin Mich Priyadharson**, Professor for his useful suggestions, which helped us in completing the work in time and we thank him for being instrumental in the completion of third year with his encouragement and unwavering support during the entire course. We are extremely thankful and pay our gratitude to our Minor project -1 coordinator **Dr. Kanimozhi T**, for her valuable guidance and support on completing this project report in a successful manner.

We are grateful to our supervisor **Dr. Koushick Venkatesh**, Assistant Professor ECE for providing me the logistic support and his/her valuable suggestion to carry out our project work successfully.

We thank our department faculty, supporting staffs and our family and friends for encouraging and supporting us throughout the project.

RANJANI .S

ROHIT KUMAR

SAI SUNNYHITH .B

TABLE OF CONTENTS

ABSTRACT	vi
LIST OF FIGURES	vii
1 INTRODUCTION	1
1.1 HUMAN-COMPUTER INTERACTION	1
1.2 DEEP LEARNING FOR GESTURE RECOGNITION	2
1.2.1 Algorithm used for Gesture Recognition:	2
2 LITERATURE SURVEY	4
2.1 OVERVIEW	4
2.1.1 Hand gesture recognition	4
2.1.2 Pattern Forming and recognition	6
2.2 METHODOLOGY	7
2.2.1 Dataset	7
2.2.2 Pre-Processing	7
RELATED PAPERS	9
3 EXISTING DEVICES IN PUBLIC USE	12
3.1 INTRODUCTION	12
3.2 PRODUCTS AND DEVICES	12
3.2.1 Low-Tech AAC:	12
3.2.2 Mid-Tech AAC:	13
3.2.3 High-Tech AAC:	13
3.2.4 Beyond these categories, other noteworthy AAC tools include:	13
3.3 SHORTCOMINGS AND AREAS FOR IMPROVEMENT	13
4 PROPOSED MODEL	15
4.1 OVERVIEW	15
4.1.1 Pre-Processing	15
4.1.2 Hand Detection	17

4.2	SIMULATION	18
4.2.1	Hand Cropping and Normalization	18
4.2.2	Our 3D CNN Architecture	20
5	RESULT AND ANALYSIS	23
5.1	ACCURACY OF MODEL	23
5.2	FUTURE RESEARCH AREAS	24
6	CONCLUSION	28
	REFERENCES	28

ABSTRACT

Speech-impaired individuals rely on sign language to express their thoughts, emotions, and convey information in daily interactions. Despite comprising a significant percentage of the population, many individuals do not possess the necessary knowledge of sign languages to effectively communicate with them. Sign language recognition systems aim to bridge this communication gap by detecting and interpreting the gestures and motions of the human body, particularly focusing on hand movements, to understand the intended message.

In critical situations such as heart attacks, accidents, or emergencies, effective communication becomes paramount, yet speech-challenged individuals often face significant barriers due to the lack of accessible communication methods. Sign language recognition systems hold immense potential in such scenarios, as they can enable swift and accurate communication between speech-impaired individuals and others, potentially saving lives in urgent situations.

In countries like India, where a considerable portion of the population—approximately 6.3%—is sensory impaired, accessibility to affordable solutions is crucial. However, existing market solutions are often prohibitively expensive for middle-class individuals, rendering them inaccessible to those who need them the most. Moreover, many of these solutions lack effectiveness and mobility, further exacerbating the communication challenges faced by speech-impaired individuals.

Sign language recognition systems represent a promising avenue for addressing these challenges. By leveraging advanced technologies such as computer vision and machine learning, these systems can accurately interpret sign language gestures and translate them into spoken or written language, enabling seamless communication between speech-impaired individuals and emergency responders or bystanders. Our solution aims to offer affordable, effective, and mobile-friendly sign language recognition capabilities. By making these tools more accessible, we strive to empower speech-impaired individuals to communicate confidently and effectively in various contexts, ultimately enhancing their quality of life and ensuring their safety and well-being in critical situations. This capability has the potential to significantly improve response times and outcomes in emergency situations, ultimately saving lives and reducing the impact of disabilities on individuals' lives.

LIST OF FIGURES

1.1	Communication using sign languages	1
1.2	Sign language interface	3
2.1	Expressing ‘Thank you’ in a) Indian, and b) American sign languages	5
2.2	Patterns formed	6
2.3	Conventional Machine Learning	7
2.4	General process of sign language recognition system.	8
2.5	Data Processing Algorithm	9
4.1	The proposed HGR systems architecture.	16
4.2	Overall structure of CNN used for proposed HGR.	17
4.3	Hand detection via masking.	18
4.4	Data Processing Algorithm	19
4.5	Different Hand sign gestures collected	20
4.6	CNN Detection Model	21
4.7	Working Model	22
5.1	Hand Gesture Recognition Parameter	24
5.2	Hand Gesture Representation Model	25
5.3	Breakdown of studies from 2014 to 2020 focusing on hand gesture representations.	27

CHAPTER 1

INTRODUCTION

1.1 HUMAN-COMPUTER INTERACTION

Human-Computer Interaction (HCI) represents a dynamic interdisciplinary field dedicated to crafting computational technologies that facilitate interaction between humans and computers. Within HCI, hand gesture recognition emerges as a prominent sub-field, leveraging computer vision and artificial intelligence to enable nonverbal communication between individuals and computing systems. This technology identifies and interprets significant movements of human hands, thereby facilitating seamless interaction.

Despite the broad array of applications for hand gesture recognition, achieving accurate recognition

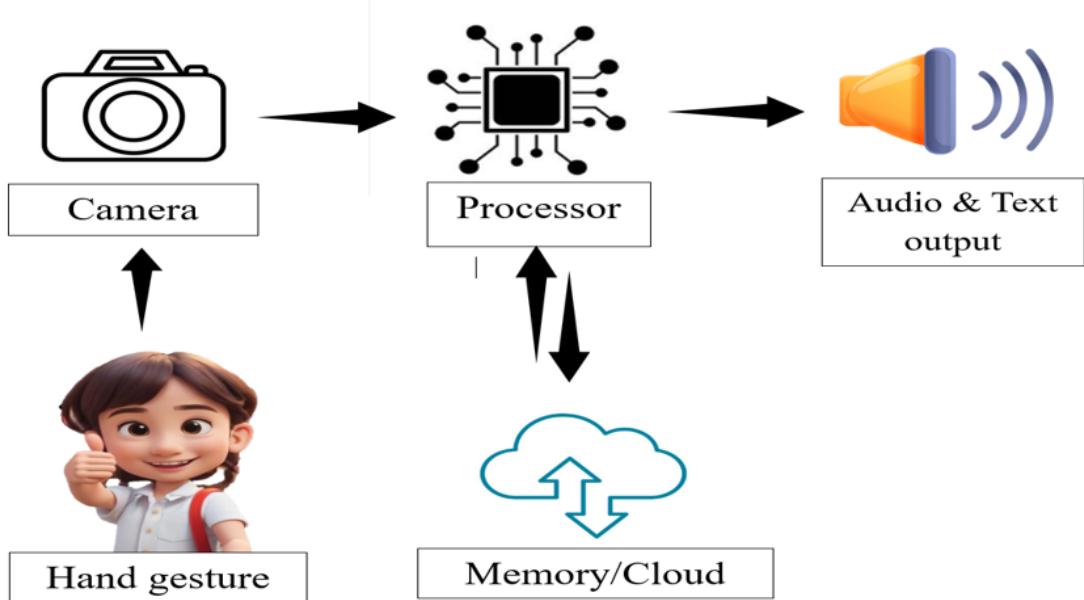


Figure 1.1: Communication using sign languages

poses a persistent challenge. Sign language recognition stands out as a primary application within this domain, enabling individuals, not solely limited to the deaf community, to convey thoughts and ideas

through nonverbal communication. Furthermore, the industrial sector has shown significant interest in hand gesture recognition, particularly in facilitating interactions between humans and robots in manufacturing settings and enhancing the capabilities of self-driving cars.

Beyond traditional applications, hand gesture recognition continues to evolve, giving rise to innovative applications such as sports-specific sign language recognition, precise monitoring of human activity, stance, and posture, as well as tracking and analyzing physical exercises. These advancements expand the scope of hand gesture recognition, offering solutions that cater to diverse needs across various domains.

1.2 DEEP LEARNING FOR GESTURE RECOGNITION

Deep learning in gesture recognition refers to using advanced computer algorithms to teach computers how to understand and interpret human gestures. These gestures could be hand movements, body motions, or facial expressions. Deep learning involves training a computer model to recognize patterns in large sets of data. In the case of gesture recognition, the model learns from examples of different gestures, such as waving hello, making a fist, or pointing. Once trained, the deep learning model can analyze new gestures and accurately identify what they mean. This technology is used in various applications, such as sign language recognition, controlling devices with hand gestures (like in gaming consoles), or even in industrial settings to communicate with robots. In simpler terms, deep learning in gesture recognition allows computers to understand and respond to the gestures humans make, opening up new ways for us to interact with technology.

1.2.1 Algorithm used for Gesture Recognition:

Convolutional Neural Networks (CNNs) are a type of deep learning model commonly used for processing visual data, such as images and videos. They excel at tasks like image classification and object detection. In recent years, researchers have extended their use to video classification due to their effectiveness in analyzing both images and their contents. Implementing deep learning models, including CNNs, for real-time applications, such as sign language recognition, poses challenges despite their fast and accurate performance. A typical implementation of sign language recognition involves gathering data (images or videos), preprocessing, feature extraction, and classification. Static gestures, represented by images, are classified directly, while dynamic gestures require sequences of images or videos to capture spatiotemporal features. Preprocessing involves selecting frames from videos, applying filters if needed, and resizing frames to fit the model's input. Features are then extracted using convolutional and pooling layers, and classification is performed using fully connected layers and a SoftMax layer to provide labels for each frame.

Sign language plays a crucial role in emergency communication, yet ordinary individuals may struggle to recognize signs for urgent matters like pain or calling for help. Sign language gestures can be categorized as static or dynamic, with the latter requiring sequences of images or videos to capture temporal and spatial features. CNNs and other deep learning algorithms have been instrumental in advancing computer vision applications, including sign language recognition, games, virtual reality, and human-computer interaction. Various approaches utilizing CNNs have been successful in classifying hand gestures, with some systems employing neural network classifiers for continuous sign language identification.[1] Dynamic sign language recognition primarily relies on deep learning techniques, such as CNNs, to extract discriminative hand movement characteristics. Static hand ges-

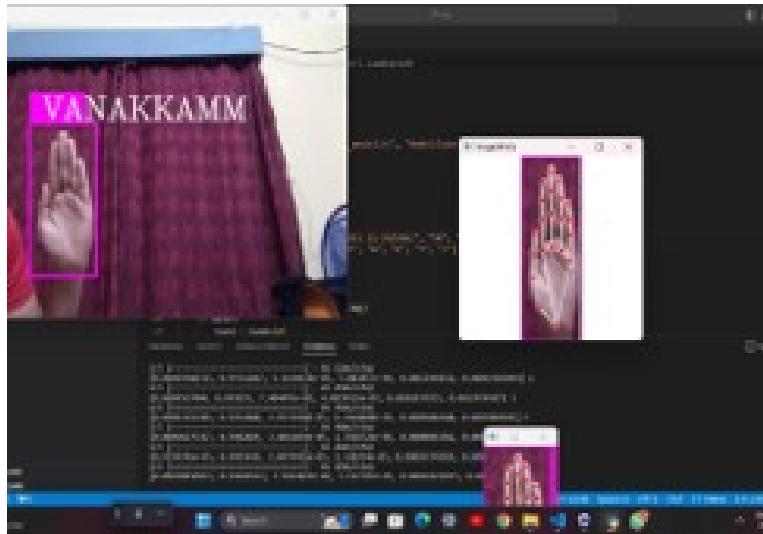


Figure 1.2: Sign language interface

tures are independent of time and order, while dynamic gestures require preserving previous states to learn long-term dependencies. Video datasets are commonly used to capture the dynamics of hand gestures effectively. Overall, CNNs and deep learning techniques have significantly contributed to the advancement of hand gesture recognition, particularly in the context of sign language communication.

CHAPTER 2

LITERATURE SURVEY

2.1 OVERVIEW

Human-Computer interaction is an interdisciplinary research area focusing on designing computational technologies to make the interaction between humans and computers possible. Hand gesture recognition is its sub-field in which computer vision and artificial intelligence have aided to provide nonverbal communication between humans and computers by identifying significant movements of the human hands. Though there are a variety of applications of hand gesture recognition, accurate recognition remains a challenging task.[2] Sign language recognition is a typical application of hand gesture recognition. It is often considered that only deaf people rely on sign languages for conveying their thoughts.

2.1.1 Hand gesture recognition

Like spoken languages, sign languages also develop naturally as a result of groups of people interacting with one another; regions and cultures also play an essential role in their development.[3] People who do not know the same language cannot communicate because most sign languages are not mutually intelligible. Common sign languages are American Sign Language (ASL), British Sign Language (BSL), and French Sign Language (FSL). Each Sign Language has its own syntax and ways to use it. For instance, ASL uses only one hand while languages like ISL and BSL use both hands in the gestures. Because of differences in syntax and expression, results obtained in recognizing one language cannot be used for another.[4] To clearly explain, take an example of a sign representing ‘Thank You.’ In ASL, it is expressed by starting with the fingers of the dominant hand near the lips and moving the hand forward and a bit down in the direction of the person we are thanking, shown in Figure 2.1(b). While in ISL, it is represented in the same way but with both hands, as shown in Figure 2.1(a).



Figure 2.1: Expressing ‘Thank you’ in a) Indian, and b) American sign languages

CNNs (Convolutional Neural Networks) are deep neural networks that process visual data. In recent years, researchers have extensively used them to solve image classification and detection problems. CNNs are also used for video classification due to their ability to classify and detect images and their contents. Deep learning models have been demonstrated to perform recognition and classification tasks quickly and accurately, but their use in real-time application settings is limited. A typical implementation of sign language recognition includes data collection (images or videos), pre-processing, feature extraction, and classification.[5] Static gestures require images, whereas dynamic gestures require a sequence of images or videos to extract spatio-temporal features. Pre-processing entails selecting multiple frames from the videos, applying different filters as needed, and resizing the frames based on the model’s input. Following preparation,(fig 2.2) the input is routed through a series of convolutional and pooling layers to extract features. To assign a label to each frame, the extracted feature output is routed to the fully connected layers and then to the SoftMax layer. To extract temporal characteristics and categorize a sequence of frames, RNN with long short-term memory (LSTM) is frequently used.

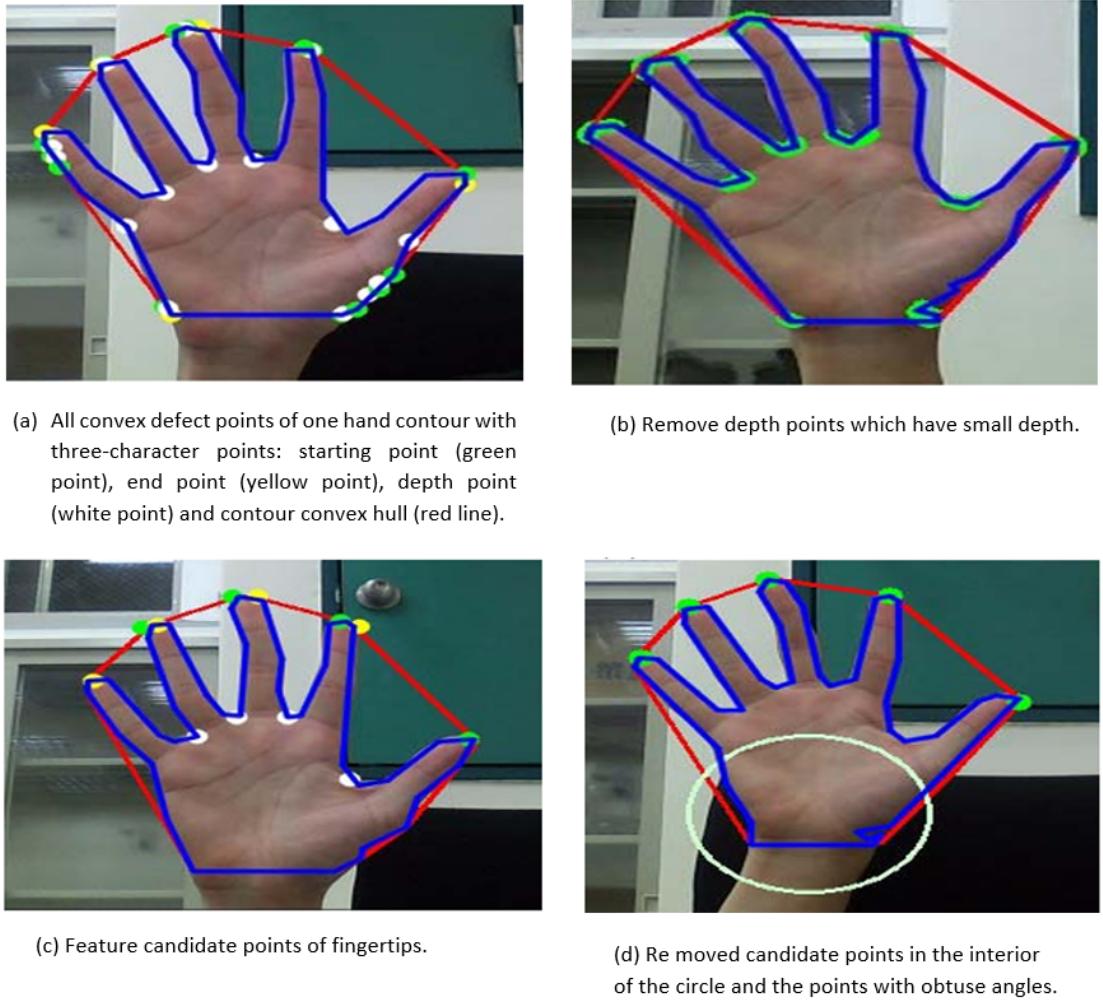


Figure 2.2: Patterns formed

2.1.2 Pattern Forming and recognition

Dynamic gesture recognition is an important aspect of HCI. This paper presents a real-time dynamic hand gesture recognition system. All hand gestures are established using dynamic video and OpenCV. The three convex defect character points of the hand contour are defined to calculate the angles between the fingers, while the fingertip positions are calculated to recognize hand gestures. (fig 2.3) Ten tested users generate 330 cases that recognize hand gestures. The eleven hand gestures representing the numbers one through nine have been recognized. The accurate recognition rate achieved more than 95.1% and each picture spends processing time is about 55 ms. The real-time dynamic hand gesture recognition system can be used in the application of HCI, such as wisdom home of our daily life, games, robots, and so on, in the future.

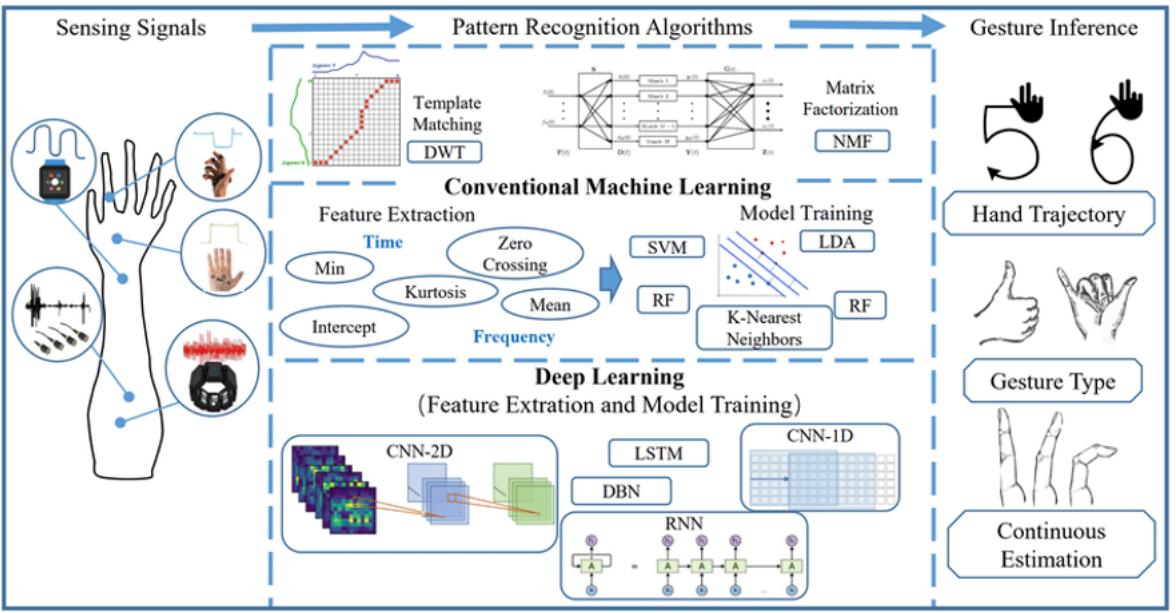


Figure 2.3: Conventional Machine Learning

2.2 METHODOLOGY

2.2.1 Dataset

Sign Language recognition has not been well-researched for ISL due to the non-availability of a standard publicly available dataset. For other languages, especially for ASL, there are many standard datasets. The three such word-level ASL datasets are Purdue RVL-SLLL ASL Database, Boston ASLLVD , and RWTH-BOSTON-50 . LSA64 is an Argentine word-level dataset, PSL Kinect 30 is a Polish word-level dataset, DEVISIGN is a Chinese, GSL is Greek, DGS Kinect is German, and LSE-sign is Spanish Sign Language dataset. In this study, a video-based ISL dataset is used that contains 412 videos. Researchers working on vision-based sign language recognition and hand gesture recognition will get benefitted from this dataset.[6] The datasets primary goal is to advance in the field of sign recognition since it has several applications in society, such as providing a platform for the deaf to communicate essential messages to authorities. Furthermore, the dataset may be used as a basic benchmark database for a collection of emergency ISL hand gestures. The dataset included eight hand gestures representing ISL words such as accident, call, doctor, help, hot, lose, pain, and thief ; often used to transmit information or request help in emergency scenarios . Out of a total of 412 videos, each sign is represented in 50 different videos on average.[7]

2.2.2 Pre-Processing

To identify the sign words, the current study uses a brief subsequence instead of the entire video. In order to effectively instruct the model on the movements' dynamics, we extracted 20 frames from every video as part of our strategy. We started by reducing the number of extracted frames (1-10)

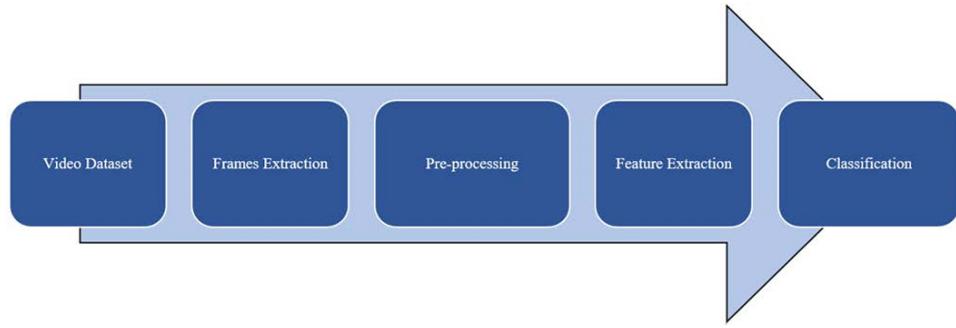


Figure 2.4: General process of sign language recognition system.

in order to shorten the training period. We found that 5 frames, taken at equal intervals from each video sequence, was enough to retract the gesture dynamics without sacrificing prediction accuracy. Additionally, earlier studies show that a mere 17 frames may be sufficient to recognize a human action. To symbolize the eight distinct classes, images were ranked from 0 to 7. The dataset was labeled in order to facilitate object detection prior to model application. [8]

Data labels and other information, like class ID and class to which it belongs, were stored in a text format using the YOLO format. The 500 by 600 pixel extracted frames have been resized to 150 by 150 pixels. In order to ensure uniform data distribution in every input pixel and hasten the convergence of the model training process, data normalization was also employed. The overall dataset (100%) is divided into three subsets: training (60%), validation (20%), and testing (20%).

Compared to gesture recognition, sign language recognition presents a more intricate and varied task. The translation of sign language through recognition systems serves as a crucial tool for fostering communication between the hearing impaired and those with functional hearing. Additionally, such recognition systems contribute to enhancing human-machine interaction intelligence within sign language education contexts. (fig 2.4) Traditionally, sign language instruction relied heavily on manual techniques and video materials. However, manual methods are constrained by limited availability of demonstrations due to manpower constraints. Moreover, video-based teaching lacks efficient feedback mechanisms for correcting improper sign language execution. Hence, there's an urgent demand for sign language teaching software that employs deep learning techniques to assess the accuracy of signs produced by sign language learners. [9]

In comparison to the relatively straightforward recognition of general gestures, sign language recognition poses a more complex and diverse challenge. Sign language serves as a primary mode of communication for many individuals with hearing impairments, and accurate translation of sign language into spoken or written language is crucial for facilitating effective communication between individuals who are deaf or hard of hearing and those with functional hearing abilities. Moreover, the application of sign language recognition extends beyond mere communication facilitation; it plays a pivotal role in advancing the field of human-machine interaction, particularly in educational settings focused on

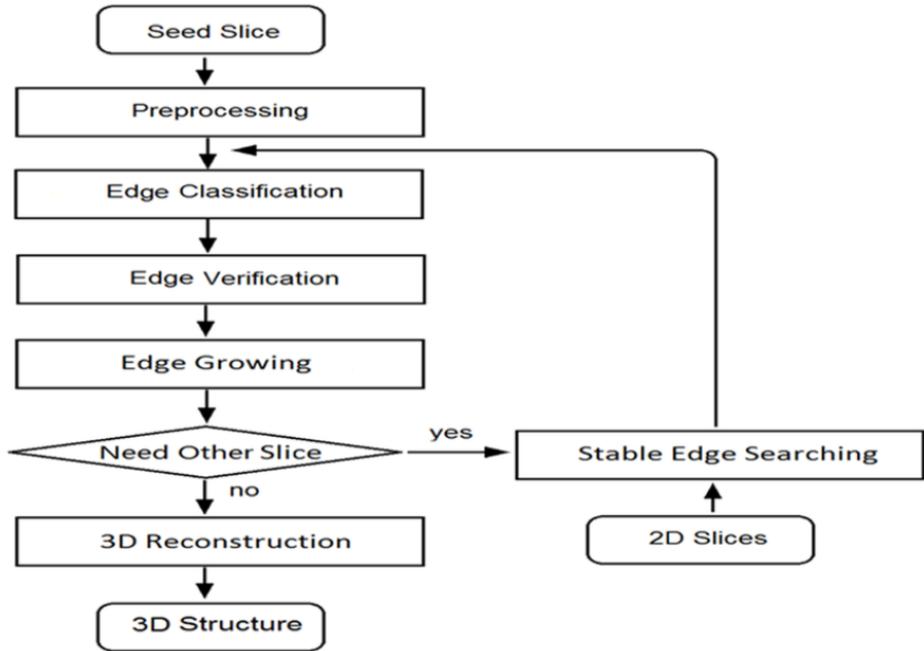


Figure 2.5: Data Processing Algorithm

sign language learning and instruction.

In addition to facilitating communication, sign language recognition systems can be integrated into public places such as transportation ticket booths, banks, and healthcare facilities (fig 2.5). When sign language is translated into spoken language and used in video conferencing platforms, barriers to real-time communication can be broken down in a variety of settings, including meetings, lectures, and social interactions.[10]

Deep learning models can be trained to recognize regional variations and dialects within sign languages, promoting inclusivity in sign language education. Additionally, integrating facial expression analysis into sign language recognition systems can provide a more nuanced understanding of conveyed messages, enriching communication experiences. Deep learning-enabled sign language teaching software can offer multimodal feedback through various channels, enhancing the

RECENT DISCOVERY IN THIS FIELD

[1] Wearable Data Glove for Indian Sign Language Recognition Using Deep Residual Networks

Authors: *Pranjal Kumar, Ashish K. Dey, Sudipta Roy (2023)*

In this recent paper, a wearable data glove-based deep learning approach for ISL recognition is proposed. It investigates how this technology might be used to make effective and cozy AAC devices. In this paper, we propose an Indian Sign Language (ISL) alphabet recognition wearable data glove embedded with deep residual networks (ResNets). The proposed system utilizes sensor data from the glove to capture hand postures and finger movements, which are then fed into the ResNet architecture for classification. This approach offers a promising solution for developing user-friendly and accurate AAC devices for deaf and speech-impaired individuals in India.

[2] A Review of Hand Gesture and Sign Language Recognition Techniques

Authors: *Abhijit Ghosh, Moumita Dutta, and Nilanjan Sinha (2023)*

This review paper addresses several methods for recognizing hand gestures and sign language, including those that are pertinent to ISL. It offers a more comprehensive framework for comprehending the advancement of AAC devices in this field. A key area of research for bridging the communication gap between the deaf and hearing communities is sign language recognition (SLR). For Indian Sign Language (ISL) to be used effectively in AAC devices, reliable and efficient SLR systems must be developed.

[3] Investigating the Efficacy of Leap Motion Controller for Recognizing Dynamic Gestures in Indian Sign Language

Authors: *Deepak Kumar Vishwakarma, Mukesh Kumar, Ashutosh Kumar (2022)* **Description:**

Building upon their previous work, this paper focuses on the effectiveness of the Leap Motion Controller in recognizing dynamic ISL gestures. This research is crucial for developing AAC devices that can capture the full complexity of sign language. **Excerpt:** "This paper presents an investigation on the efficacy of Leap Motion Controller (LMC) for recognizing dynamic gestures in Indian Sign Language (ISL). Recognizing dynamic gestures is essential for seamless communication using sign language, and this research paves the way for developing more comprehensive AAC devices for the Indian population."

[4] **Development of a Low-Cost Real-time Hand Gesture Recognition System for Assistive Devices**

Authors: *Sunil Kumar, Gaurav Dogra, and Suman Kumar (2021)*

The purpose of this paper is to discuss the need for affordable AAC devices. It explores the development of a low-cost real-time hand gesture recognition system with potential applications for users in India. Cost is a crucial factor for the widespread adoption of AAC devices, particularly in developing countries like India. This paper proposes a low-cost real-time hand gesture recognition system that can be integrated into AAC devices, making them more accessible to a wider range of users.

[5] **A Deep Learning Framework for Sign Language Recognition using Indian Sign Language Dataset**

Authors: *Priyanka Goyal, Kanupriya Goswami, Garima Malik, and Ekta Gupta (2020)*

This paper delves into deep learning techniques for sign language recognition using an ISL dataset. It explores the potential of these methods for improving the accuracy and efficiency of AAC devices. Deep learning has shown promising results in various computer vision tasks, including sign language recognition. This paper proposes a deep learning framework for recognizing Indian Sign Language (ISL) signs using a publicly available ISL dataset. This research contributes to developing more accurate and efficient AAC devices for deaf and speech-impaired individuals in India.

CHAPTER 3

EXISTING DEVICES IN PUBLIC USE

3.1 INTRODUCTION

Communication is the cornerstone of human interaction, allowing us to express needs, share ideas, and build relationships. However, for individuals with speech and language impairments, traditional communication methods may be challenging or impossible. Augmentative and Alternative Communication (AAC) steps in to bridge this gap, providing a range of tools and strategies to empower these individuals to participate actively in their world.[9][11]

This chapter explores the diverse landscape of AAC devices currently used around the world. We will delve into the functionalities of various products, analyze their strengths, and identify potential shortcomings that researchers and developers are working to overcome.

3.2 PRODUCTS AND DEVICES

The world of AAC encompasses a spectrum of tools, catering to different needs and abilities. Here's a closer look at some prevalent categories:

3.2.1 Low-Tech AAC:

- **Picture Boards:** These are simple boards featuring symbols or pictures representing words, phrases, or actions. Users point to the desired symbol to communicate. Their portability and ease of use make them ideal for various situations, especially for individuals with limited motor skills or cognitive abilities.
- **Object Symbols:** Actual objects can serve as communication tools. For example, a cup might symbolize thirst, and a key might represent the need to go outside. This approach is particularly helpful for individuals who may not grasp abstract pictures.
- **Communication Books:** These bound collections of pages contain picture symbols or written words. Users flip through the pages and point to express themselves. Communication books can

be customized for specific needs and situations, such as requesting food at school or expressing emotions at home.

3.2.2 Mid-Tech AAC:

- **Speech-Generating Devices (SGDs):** These electronic devices are equipped with buttons or touchscreens. Users select symbols or words on the screen to build messages that are then vocalized through synthesized speech. SGDs often offer features like customized voices, vocabulary sets, and message prediction to enhance communication efficiency.

3.2.3 High-Tech AAC:

- **Tablets and Computers with AAC Apps:** These devices leverage the power of tablets and computers with dedicated AAC applications. Apps can provide a variety of functionalities, including synthesized speech, text prediction, and access to vast symbol sets. The flexibility and customization options offered by AAC apps surpass traditional SGDs, allowing for personalized communication strategies.
- **Eye Gaze Technology:** This technology empowers individuals with limited mobility to control a computer cursor or communication software using their eye movements. Users can select symbols or words on the screen to communicate, offering a hands-free communication approach.

3.2.4 Beyond these categories, other noteworthy AAC tools include:

- **Body Language and Gestures:** These nonverbal cues can be a powerful form of AAC, especially for individuals who are nonverbal or in the early stages of communication development.
- **Sign Language:** Formal sign languages are complete communication systems widely used by deaf communities. They can be a primary or secondary mode of AAC for individuals with hearing impairments.

3.3 SHORTCOMINGS AND AREAS FOR IMPROVEMENT

While AAC devices have revolutionized communication accessibility, there's still room for improvement. Here are some challenges that researchers and developers are actively addressing:

- **Limited Vocabulary and Customization:** Some AAC devices, particularly low-tech options, may have limitations in vocabulary size. This can restrict the ability to express complex ideas or emotions. Advancements in technology aim to provide more comprehensive vocabulary sets and customizable features to cater to diverse needs.
- **Cost and Accessibility:** High-tech AAC devices can be expensive, potentially limiting access for individuals with financial constraints. Initiatives are underway to develop more affordable AAC solutions and increase insurance coverage for these devices.

- **Social Stigma:** There can be a social stigma associated with using AAC devices. Raising awareness and promoting the positive impact of AAC can help break down these barriers.
- **Natural Language Processing Integration:** Current AAC systems may not fully capture the nuances of natural language, hindering the flow and expressiveness of communication. Integrating advanced natural language processing (NLP) techniques can enable more natural and intuitive communication experiences.
- **Limited Environmental Integration:** While some AAC devices can connect with other assistive technologies, there's still a gap in seamlessly integrating AAC into various environments. Future advancements aim to create a more interconnected ecosystem where AAC tools can interact with smart homes, educational tools, and other technologies for a more inclusive environment.

The field of AAC is constantly evolving, with researchers and developers striving to overcome existing limitations and create more robust, user-friendly communication solutions. As AAC technology continues to develop, we can expect to see a future where individuals with communication needs have access to a wider range of tools that empower them to participate meaningfully in all aspects of life.

CHAPTER 4

PROPOSED MODEL

4.1 OVERVIEW

This section provides a detailed explanation of the suggested system methodology. The five phases of the proposed HGR make it an effective system for hand gesture detection. Prior to applying morphological operators like dilation and erosion, a low-resolution grayscale image is first pre-processed to adjust its size. Next, noise is reduced to remove extra information. Finally, the image light intensity is adjusted by sharpening and enhancing it using filters. Finally, the image is converted into binary. Second, the image's directional hand gesture region is identified. Thirdly, convexity removal is followed by an examination of extreme point localization over the detected hand via a convex hull, leading to landmark extraction. Fourthly, on hand gestures, geometric features and distance are detected over extreme points. Finally, CNN is used to classify various hand gestures using landmark-based features.[12]

4.1.1 Pre-Processing

Images of hand gestures are pre-processed to effectively extract hands from noisy, low-resolution greyscale images. Two datasets, MINIST and ASL, are taken into consideration for the proposed HGR system. ASL uses 200x200 resolution, whereas the MINIST dataset uses a 28x28 resolution grayscale image size. Greyscale images are first created from the ASL dataset. Unwanted information can be seen in the images. Gaussian noise in the MINIST dataset is eliminated. The random value that is added to the initial pixel value is known as Gaussian noise.

The proposed Hand Gesture Recognition (HGR) system entails pre-processing hand gesture images to effectively extract hands from noisy and low-resolution grayscale images. This process involves the utilization of two primary datasets: the MINIST dataset and the American Sign Language (ASL) dataset.[13]

The ASL dataset comprises images with a resolution of 200x200 pixels, while the MINIST dataset contains images with a smaller resolution of 28x28 pixels. Initially, grayscale images are created from

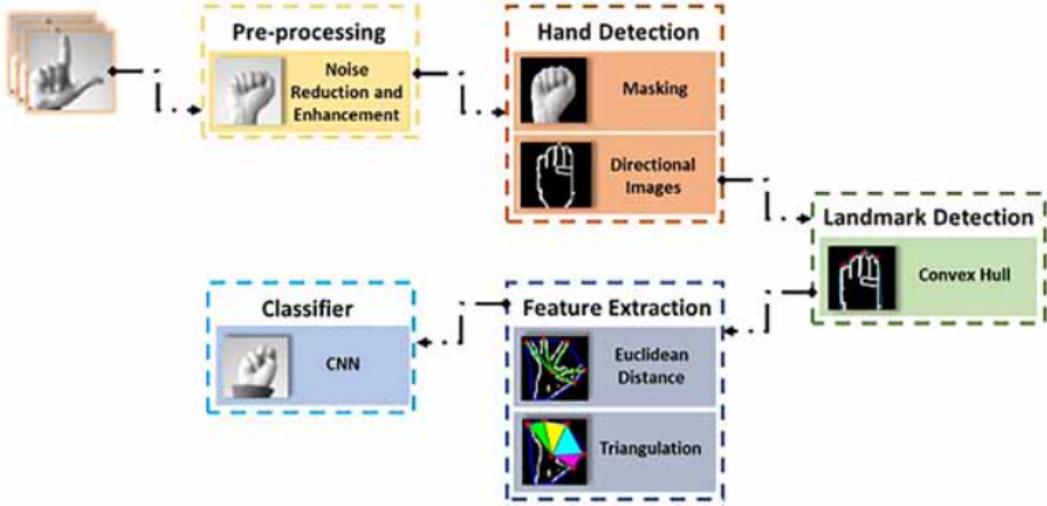


Figure 4.1: The proposed HGR systems architecture.

the ASL dataset, wherein unwanted information may be present due to factors such as background clutter or artifacts.

To enhance the quality of images from the MINIST dataset, Gaussian noise, a common type of noise characterized by its random distribution, is addressed. Gaussian noise is effectively eliminated through noise reduction techniques, ensuring cleaner images for subsequent processing stages. This noise reduction process involves modifying pixel values by adding random values, which are typically generated from a Gaussian distribution, to the initial pixel values.

By pre-processing images from both datasets to mitigate noise and unwanted information, the HGR system can facilitate more accurate and reliable hand gesture recognition. These pre-processing steps are essential for optimizing the quality and suitability of input data, thereby improving the overall performance of the HGR system in recognizing and interpreting hand gestures from diverse sources.[14]

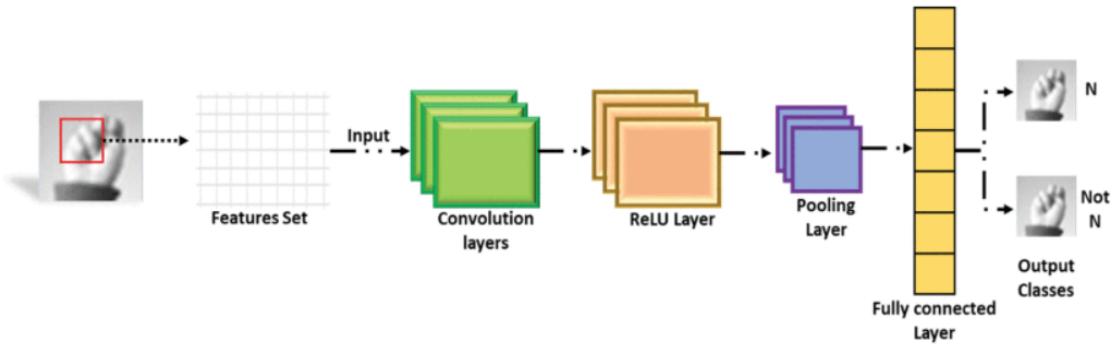


Figure 4.2: Overall structure of CNN used for proposed HGR.

4.1.2 Hand Detection

The hand region is extracted for further model processing following the pre-processing technique. A two-way method is used to detect the hand. The image is first transformed into binary using the threshold method of Otsu. (fig 4.2) The binary image is then subjected to morphological processes like erosion and dilation. To obtain the intended result in the case of dilation, the hand region's area—which includes the sphere structuring element—is examined.

S, the organizing element, is moved by h and image. S and I have to share at least one thing in common. The dilation process shrank the background, which had a 0-pixel value, and increased the pixel value in the hand region by 1.

The hand region's boundary pixel is lost and the background is enlarged when the structuring element S is translated by h. Following the processes of erosion and dilation, a connected component is used to determine the precise hand boundary. The hand is then covered by a mask over the binary image.[15]



Figure 4.3: Hand detection via masking.

4.2 SIMULATION

The input videos are converted into RGB frame sequences of varying lengths, and linear sampling is used for temporal dimension normalization. Only 16 frames are linearly selected from each video sequence.

For hand gesture recognition, the sequence order should be preserved because it encodes highly discriminative features (fig 4.3), so linear sampling is the preferred technique. Bag-of-visual-words techniques, for example, are very efficient when the sequence order is of low importance for discrimination, as in video event and human action recognition.[16]

4.2.1 Hand Cropping and Normalization

This method makes use of openpose, an open-source real-time human pose estimation framework based on deep learning that detects each individual’s 2D key points in an image. This framework enhances the machine’s understanding of human activity in an image or video sequence. It accepts an RGB image as input and returns a list of (x,y) coordinates for each key point on the human body. From the whole list of returned key points, only the wrist and elbow joints are used for cropping the hand region.

Typically, when using transfer learning, some of the architecture layers are iteratively fine-tuned on the target domain data to adjust their parameters. The other layers, on the other hand, are frozen in order to keep their parameter values unchanged. To determine the optimal case, we investigated how changing the number of trainable layers affects the performance of the C3D architecture. This optimization was done in the signer-independent mode. All the samples that were recorded by the first 32 signers (80% of the samples), were used for training the architecture. The remaining 1600 samples, that were recorded by the other eight signers (20% of the samples), were used for evaluation.[17]

We linearly sampled 16 frames from each sequence, each of which contained the entire gesture space. Following the replacement of the last two FC layers and the classification layer, the C3D architecture underwent end-to-end training. The model was fitted using mini-batch gradient descent with a learning rate of 104, weight decay of 106, and momentum of 0.9 over 100 iterations. The batch size was 16 samples. We repeated the experiment, changing the number of trainable and frozen layers

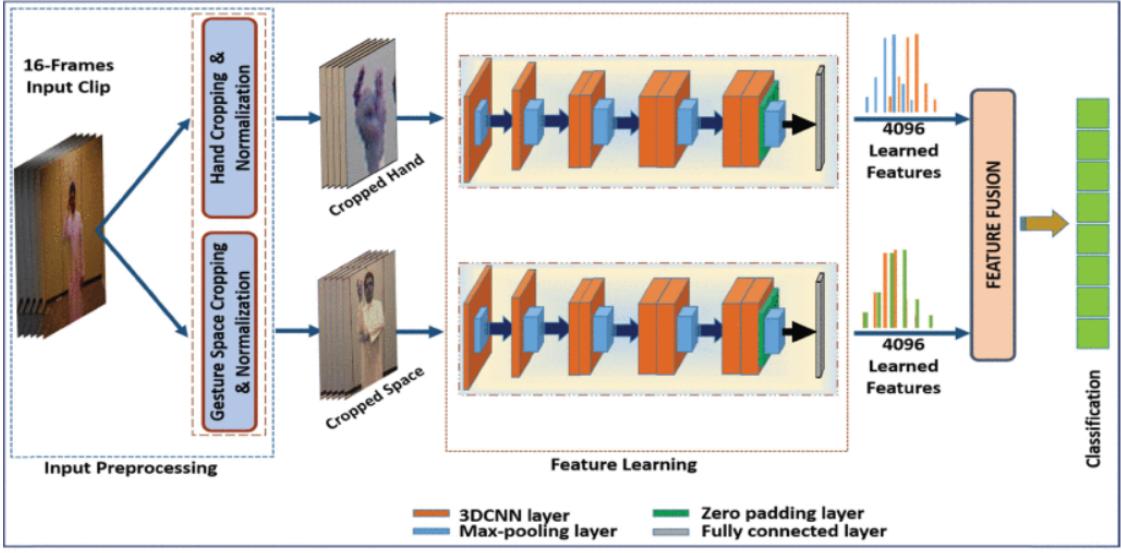


Figure 4.4: Data Processing Algorithm

each time, to determine the best level for knowledge transfer. We began by training only the final 3DCNN layer with the FC layer and the classification layer, leaving the remaining layers frozen.

Then, in each repetition, we increased the number of trainable layers by activating the layer nearest to the previously activated ones. It depicts the results of the experiment in terms of evaluation loss and recognition accuracy. It demonstrates that increasing the number of trainable layers improves model performance, as long as the first layer is frozen. That is, the best performance (80.94%) was achieved by fine-tuning all the layers except the first one. (fig 4.4) This result supports the intuition that the first layer learns common preliminary motifs in both the source and target domains. As a result, the parameters of this layer were optimized well on the source data and there was no need to distort them by a small and maybe noisy data of the target domain.[18]

As we mentioned before, CNNs are capable of learning features from raw image automatically, however GMM-HMM is not. That means we need to extract hand-crafted features from sign language video and then use those features to train GMM-HMM. On observation of sign motion, we find that both the change of hand-shape and the trajectory of body movement are two of the most important features to describe a sign motion. So we extract trajectory and hand-shape features to train GMM-HMM for recognition. We use the idea proposed by to crop the hand shape from the background. In addition, Kinect was used as an input device, and an effective algorithm for hand segmentation and tracking was proposed. The algorithm considers both color and depth information, with no requirement for a uniformly colored or stable background. We use a 32×32 hand-shape image to calculate 36-D HOG. The 32×32 image is regarded as one cell. Each cell has 16×16 pixels and the gradient is divided into 9 orientation bins. A block consists of four adjacent cells. Each image includes one block.[19]Simultaneously, we use 3D coordinate positions of key skeleton joints as trajectory features for simplicity. Skeleton joints include the right and left hands, right and left wrists, right and



Figure 4.5: Different Hand sign gestures collected

left elbows, right and left shoulders, shoulder centers, and heads. Trajectory features are obtained as a $3 \times 10 = 30$ dimension vector. After combining these two types of features, we get a 66-dimensional vector that we use to train GMM-HMM. For the recognition stage, we extract a 66-D vector from each frame of video and use the trained GMM-HMM to generate semantics.

4.2.2 Our 3D CNN Architecture

A variety of CNN architectures can be created using the 3D convolution described above. The following describes a 3D CNN architecture that we created for sign language recognition. We use the Microsoft Kinect as an input device, which provides a color video stream, a depth video stream, and the ability to track users' body movements at the same time. Color information consists of three RGB channels. We obtain five different types of input data in total, including depth and body skeleton. In our CNN architecture , we consider nine 64×48 frames centered on the current frame for each type of visual source as input to the 3D CNN (fig 4.6). This yields five feature maps, denoted by color-R, color-G, color-B, depth, and body-skeleton. Each feature map consists of nine stacked frames from the corresponding channel arranged in a cube. Multiple feature maps as input typically result in better performance than a single gray-scale intensity input.[16]

Our architecture is made up of eight layers, including an input layer. Following the input layer, there are four layers: convolution (C1), sub-sampling (S1), and convolution (C2), which is then followed by sub-sampling (S2). This is followed by a third convolution layer (C3) that does not include subsampling. This is followed by two fully connected layers that house the output layer. It is

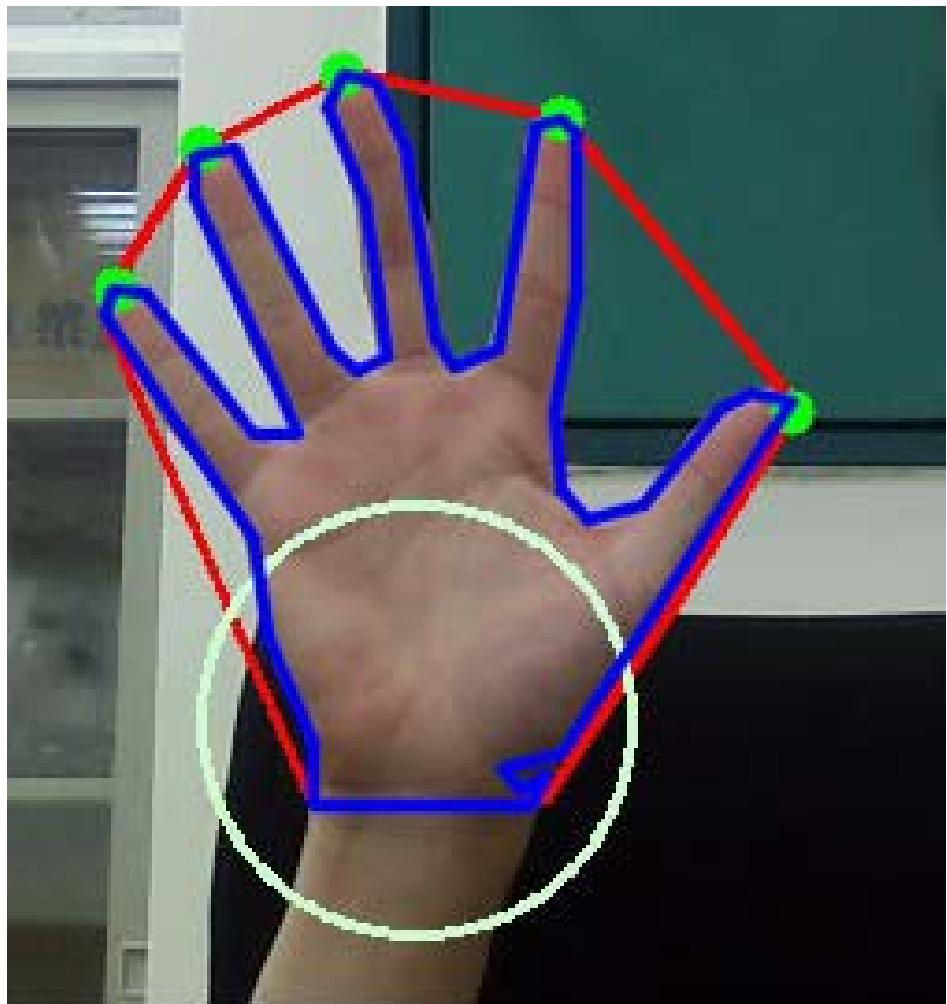


Figure 4.6: CNN Detection Model

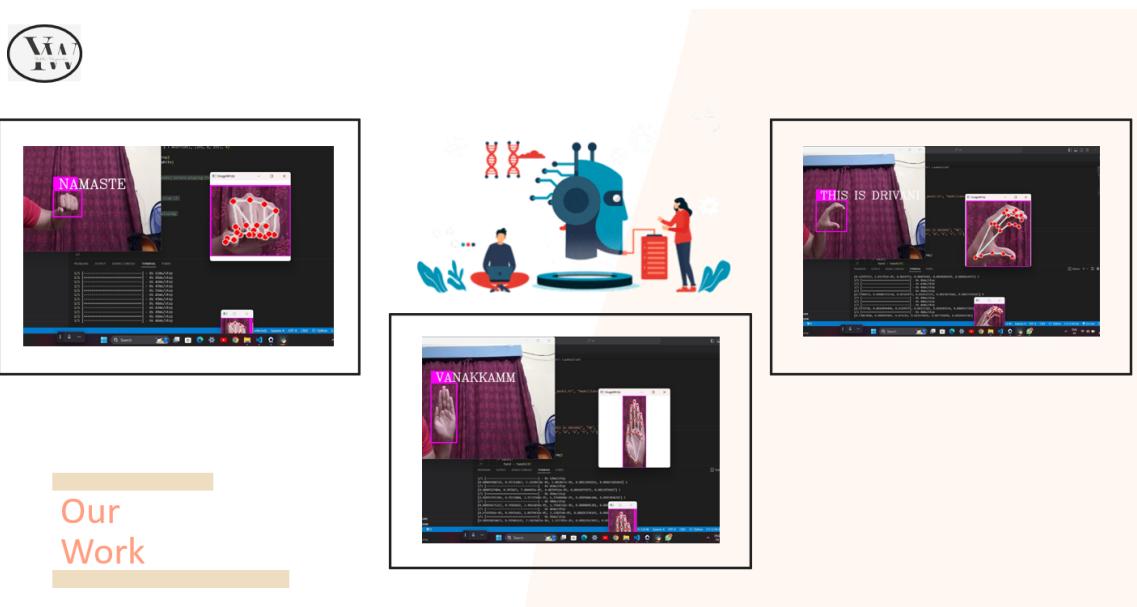


Figure 4.7: Working Model

critical in architecture to design kernel sizes in various layers. First, we apply 50 different 3D kernels size of $7 \times 7 \times 5$ (7×7 in spatial dimension and 5 in temporal dimension) on all five channels to obtain the first convolutional layer C1, which consists of 50 features maps of size $58 \times 42 \times 5$ (58×42 in spatial dimension and 5 in temporal dimension). In the S1 subsampling layer, we apply 2×2 subsampling to each feature map in the C1 layer, resulting in the same number of feature maps with a reduced spatial resolution. Our model is more resistant to small spatial distortions thanks to subsampling. As a result, the subsequent layers recognize patterns at larger spatial scales and lower resolution. Thus, a CNN with multiple subsampling layers can process large inputs with relatively few free weights. The second convolution layer C2 includes 50 feature maps measuring $23 \times 15 \times 3$. This layer is created using $7 \times 7 \times 3$ 3D kernels and factor 2 sub-sampling. The third convolution layer C3 has 10 $8 \times 4 \times 1$ feature maps created by applying 3D kernels of size $5 \times 5 \times 3$ to layer S2.[17][18]

CHAPTER 5

RESULT AND ANALYSIS

The accuracy, state-of-the-art, and future directions of hand gesture recognition systems are covered in this chapter. First, the recognition accuracy of different gesture types in different environments is described, emphasizing the preference for isolated gestures over continuous gestures and finger spelling; signer-dependent recognition is generally more accurate than signer-independent recognition, especially in controlled environments where accuracies are generally high.

5.1 ACCURACY OF MODEL

Most works focused on recognizing isolated gestures at 67% compared to dynamic recognition for continuous gestures at only 21%. Only 12% of the works used finger spelling words and alphabets.

We also looked at the recognition accuracy based on the data environment as shown in Table. For restricted environment, the average recognition accuracy for signer dependent is 88%, and for signer independent is 77%. However, for the uncontrolled environment (with only three articles reporting the results), the average recognition accuracy for signer dependent is 98%, and for signer independent is 90%. Though it does not provide conclusive evidence, it can be observed that the research on hand gestures in an uncontrolled environment shows promising results. [15][16]

Table shows the recognition accuracy for hand gesture recognition. For isolated gestures, the

Input method	Recognition accuracy		
	Min	Max	Average
Single camera	SD = 70 SI= 39	SD = 98 SI = 97	SD = 88 SI = 79
Active techniques	SD = 69 SI = 56	SD = 97 SI = 97	SD = 90 SI = 77

SD=Signer Dependent, SI=Signer Independent

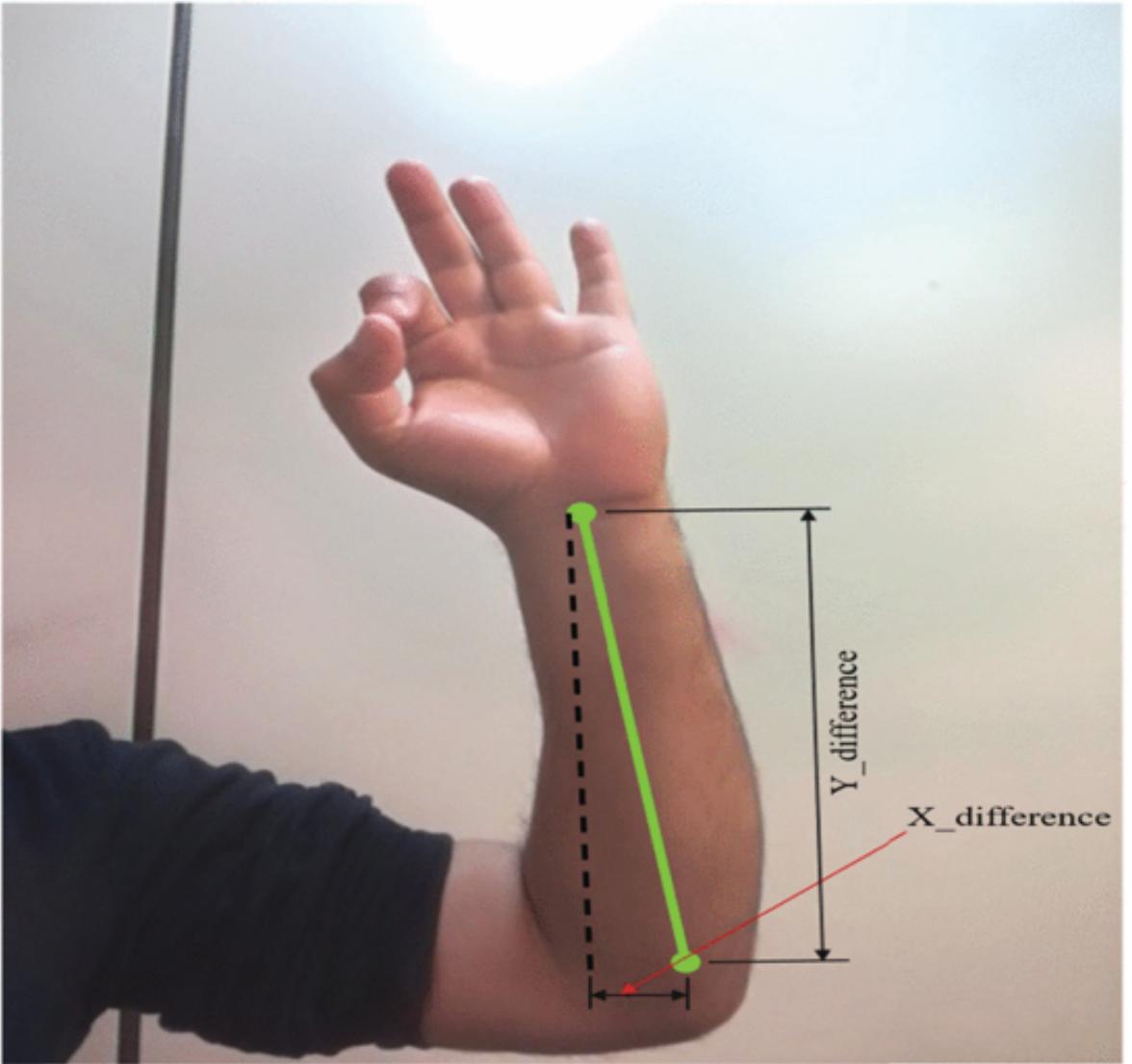


Figure 5.1: Hand Gesture Recognition Parameter

average recognition accuracy for signer dependent is 92%, and for signer independent is 77%. For continuous gestures, the average recognition accuracy for signer dependent is 84%, and for signer independent is 82% (fig 5.2) . On the other hand, for the finger spelling, the average recognition accuracy for signer dependent is 81%, and for signer independent is 71%. [20]

5.2 FUTURE RESEARCH AREAS

Currently, smartphone has become the most popular electronic device for most people because it provides us with various functions and enriches our lives. In addition to standard audio features, smartphones have the ability to transmit and receive ultrasonic signals, enabling them to function as active sonar sensing devices. Based on this sensing system, numerous intriguing applications for hand gesture recognition have recently been developed and implemented in a variety of contexts. Even though these systems perform well in terms of recognition, there are still a lot of problems that need

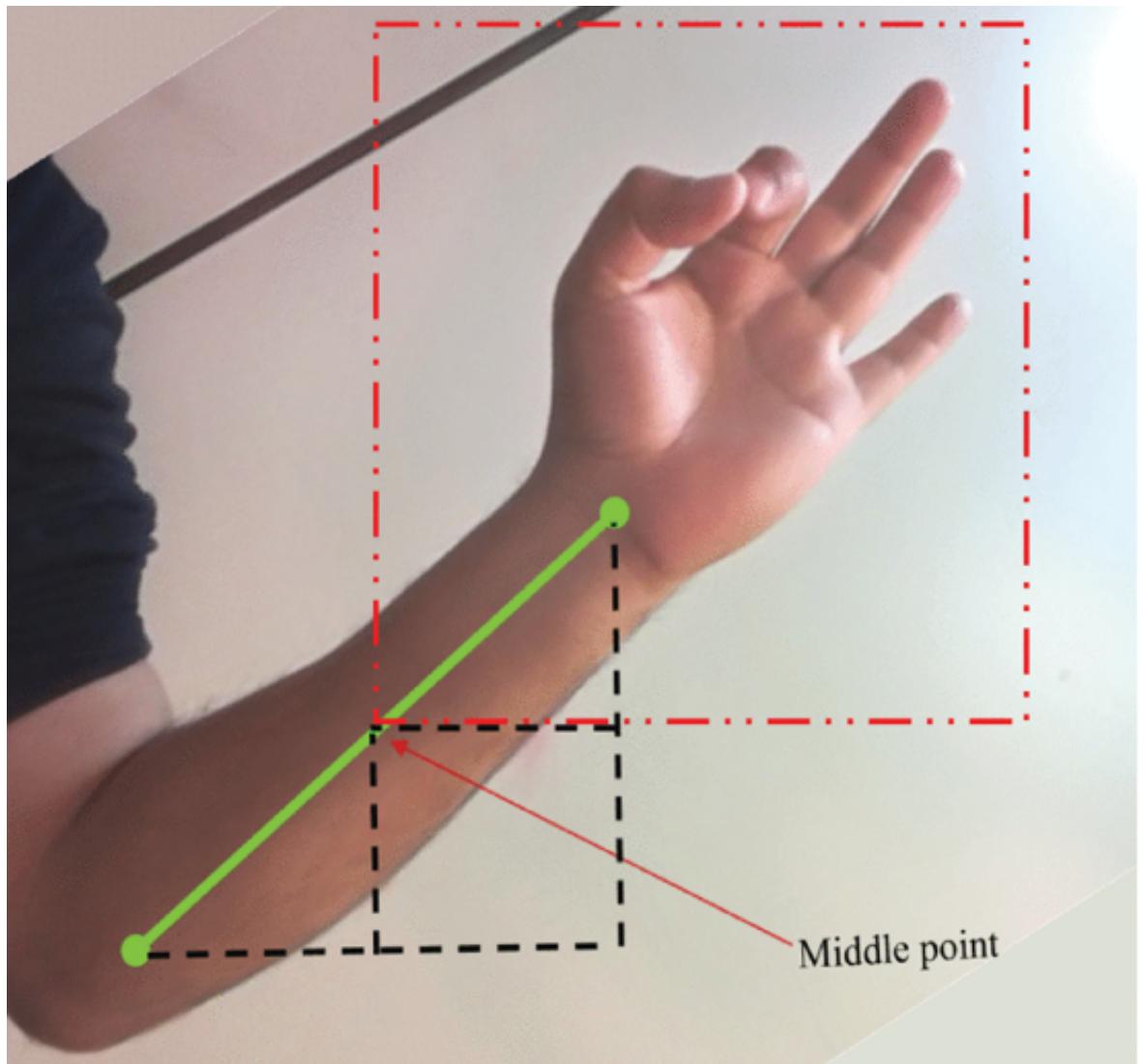


Figure 5.2: Hand Gesture Representation Model

to be fixed in order to increase identification accuracy. The primary problems and potential fixes for smartphone-based active sonar hand gesture sensing applications are discussed in this section.

A. Multi-Data Fusion

The majority of applications in this paper use the smartphone's built-in speakers and microphones to detect hand gestures. The quantity of sensors restricts the range of gestures possible as well as the size of the measurement data. Therefore, these systems recognize a few types of dynamic gestures and track hand in 1D and 2D space. UltraGesture proves that adding extra sensors can achieve fine-grained gesture recognition and identify more hand postures. To recognize more hand postures and track hand in a 3D space, using more speakers and microphones would be a fine solution . Besides, the fusion of many features, including Doppler shift, phase, ToF, and CIR would be another simple and effective method. Other sensors embedded in the smartphone (e.g., accelerometer , gyroscope) can also provide more useful information to improve the quality of HCI and recognition accuracy.[11]

B. Robustness

Current research verifies system performance in a range of noise-leveled environments. However, it seems difficult to compare these levels of noise because they are much different. Besides, human movements near the participant severely affect recognition accuracy because the movements generate complicated multipath effect . In addition, the distance between the smartphone and the user is also an important factor that affects the recognition performance . Moreover, the hand size, the hand that participants used (left or right), the hand with or without occlusion can also impact the accuracy of the application . Therefore, how to eliminate the impacts of nearby environmental changes and interference is a challenging problem.[14]

C. Standard Dataset

Currently, many systems evaluate the performance by using different types of actions. For dynamic gesture recognition, the number of gestures is very different . It varies from 5 to 24. For hand tracking, the shapes drawn by hand contain some simple shapes (e.g., diamond, triangle, circle) or letters and words . The variation of the shapes generates significant difficulty when evaluating the system performance because of the difference of the experimental conditions. To effectively compare the performance of various algorithms, the standard dataset is required.

D. Low-Latency

For gesture recognition applications, latency usually is not a crucial metric because we pay more attention to recognition accuracy. However, for hand tracking applications, we may concentrate on the latency since many applications require a swift response to hand position. For example, low-latency is a crucial factor for real-time game control. However, a few applications consider latency feature when evaluating system performance, such as Strata and LLAP . Most others do not analyze the latency

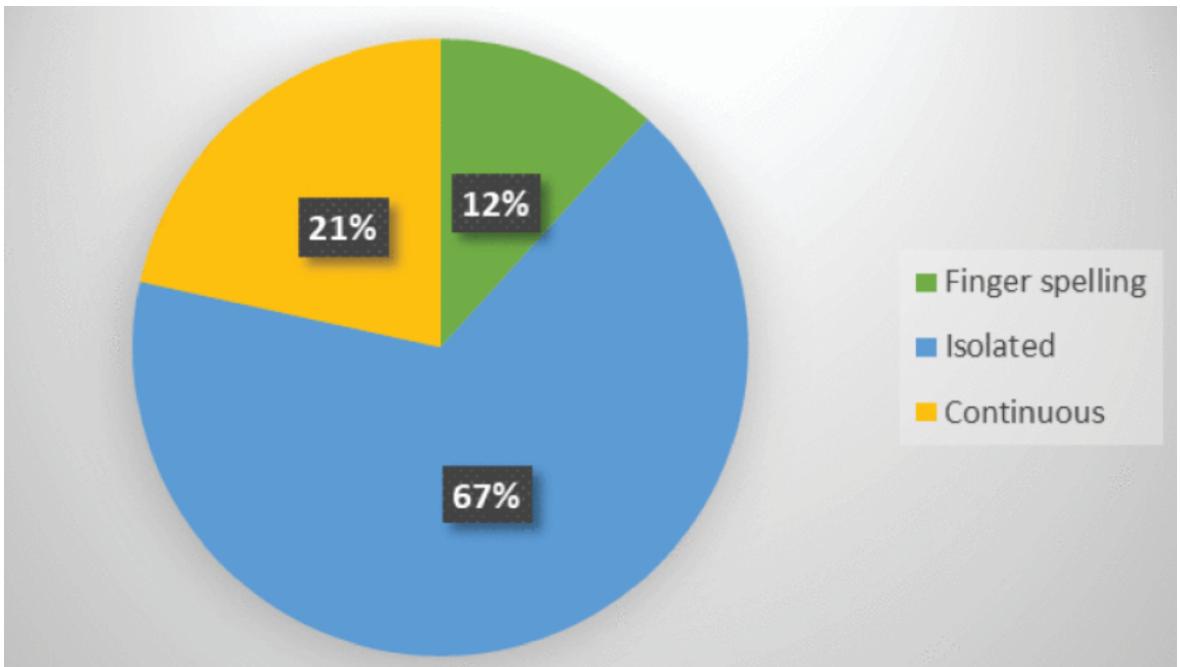


Figure 5.3: Breakdown of studies from 2014 to 2020 focusing on hand gesture representations.

factor when testing system performance. Besides that, the signal types and feature calculation are must be considered when developing a hand tracking algorithm.

E. Security Issues

We now have a novel approach to HCI with hand gesture recognition based on the ultrasonic signal of smartphones. Because the systems do not require any additional hardware and give us easy ways for users to interact, researchers are focusing more on creating applications using ultrasonic signals. However, there are security concerns associated with this method as well. Given that our systems use sound signals with frequencies between 16 and 24 kHz, and that the average person's hearing range is within 16 kHz, there is a possibility that an ultrasonic signal attack could be launched. Simultaneously, this technique could be used to steal personal data, deduce Android unlock patterns, and even perform inaudible sound attacks. As a result, in order to address these problems, some useful strategies must be researched in the future.[18]

CHAPTER 6

CONCLUSION

This chapter examines the evolution and challenges of vision-based hand gesture recognition systems over a seven-year period. It emphasizes the importance of feature selection, data acquisition, and training environments in determining the effectiveness of these systems, as well as the widespread use of databases derived from restricted environments. However, there is a growing recognition that more diverse datasets are required to improve the real-world applicability of these systems. Within this realm, the comparison between dynamic gesture recognition and hand tracking reveals distinct challenges and applications. Dynamic gesture recognition, akin to a classification problem, contrasts with the intricate task of hand tracking, which involves continuous localization in 2D or 3D space. The topics covered in the discussions are signal acquisition, processing, and performance evaluation, with an emphasis on the special traits and analytical insights of each sort of application.

The proposed systems employ optimized features and convolutional neural networks (CNN) to achieve encouraging improvements in accuracy enhancement. Techniques such as adaptive median filtering for noise reduction and convex hull for landmark extraction have significantly improved accuracy issues. However, there are still issues that necessitate further advancement in feature extraction methods, such as noise sensitivity and potential convex hull errors. The use of smartphone-based ultrasonic sensing for dynamic gesture tracking and recognition, which offers a natural and comfortable method of human-computer interaction, is a noteworthy development in the interim. This method leverages the microphones and speakers pre-installed on smartphones to facilitate dynamic gesture recognition and open up a plethora of Internet of Things applications.[17]

Additionally, a thorough analysis of acoustic signal-based hand gesture recognition systems is carried out, classifying smartphone-based audio sensing techniques into passive and active categories. A suggested framework for ultrasonic signal-based hand gesture recognition is also presented, along with a performance analysis of the system and processing methods. The fact that these systems can recognize sign language gestures better than cutting-edge techniques demonstrates how effective they are. Future research will concentrate on investigating temporal modeling strategies, maximizing input clip length, and testing for real-time recognition with the goal of improving the effectiveness and usability of hand gesture recognition systems in various real-world contexts.

REFERENCES

- [1] P. Kumar, A. K. Dey, S. Roy, "Wearable Data Glove for Indian Sign Language Recognition Using Deep Residual Networks", May 2022.
- [2] A. Ghosh, M. Dutta, N. Sinha, "A Review of Hand Gesture and Sign Language Recognition Techniques", Jun. 2013.
- [3] D. K. Vishwakarma, M. Kumar, A. Kumar, "Investigating the Efficacy of Leap Motion Controller for Recognizing Dynamic Gestures in Indian Sign Language", Mar. 2018.
- [4] S. Kumar, G. Dogra, S. Kumar, "Development of a Low-Cost Real-time Hand Gesture Recognition System for Assistive Devices", Jul. 2019.
- [5] P. Goyal, K. Goswami, G. Malik, E. Gupta, "A Deep Learning Framework for Sign Language Recognition using Indian Sign Language Dataset" ,May 2020.
- [6] J. Shi, Y. Li, Y. Liu, "Towards robust sign language recognition with dynamic attention and spatial constraints", Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Jun. 2021.
- [7] H. Martinez, A. Juan, M. Rudzsky, "Sign language recognition with ensemble deep learning models", Pattern Recognition Letters, vol. 138, pp. 222-228, May 2020.
- [8] M. Jadhav, P. Mahajan, "A review on sign language recognition using deep learning techniques", Artificial Intelligence Review, vol. 54, no. 2, pp. 1483-1521, Jun. 2020.
- [9] J. Wang, Y. Cui, Z. Wan, T. Long, "Continuous sign language recognition with a deep residual network", IEEE Transactions on Circuits and Systems for Video Technology, May 2020.
- [10] H. Cooper, R. Bowden, J. Kittler, "Sign language recognition from a single viewpoint", Computer Vision and Image Understanding, vol. 78, no. 1, pp. 133-157, Jul. 2000.
- [11] J. Chen, H. Yang, Y. Xu, "Sign language finger alphabet recognition using a wrist-mounted EMG sensor and a convolutional neural network", Biomedical Signal Processing and Control, vol. 52, 101902, Mar. 2023.
- [12] H. Wu, X. Chen, Y. Li, Y. Wang, H. Li, "Sign language translation with a hierarchical attention mechanism", IEEE Transactions on Circuits and Systems for Video Technology, Mar. 2022.

- [13] O. Koller, V. Krichevsky, R. Breakstone, "Sign language recognition using 3D convolutional neural networks", 2020 IEEE International Conference on Image Processing (ICIP), Oct. 2020.
- [14] Y. Jang, J. Kim, M. Kim, "Sign language recognition with deep convolutional neural networks", 2017 14th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS), Aug. 2017.
- [15] A. Plouffe, R. Bouchareb, "Deep learning for sign language gesture recognition", 2016 International Joint Conference on Neural Networks (IJCNN), Jul. 2016.
- [16] A. Mittal, M. Garg, A. Gupta, "Sign language recognition with partial hand occlusion handling using a transformer-based approach", Pattern Recognition Letters, vol. 166, pp. 348-355, Mar. 2022.
- [17] Z. Li, J. Ding, Y. Li, "Sign language recognition with privacy protection using federated learning", Neural Processing Letters, pp. 1-17, Jan. 2023.
- [18] J. Siha, Y. Wang, H. Li, "Sign language emotion recognition with graph convolutional neural networks", IEEE Transactions on Circuits and Systems for Video Technology, vol. 32, no. 3, pp. 838-851, Mar. 2021.
- [19] W. Chen, Z. Li, J. Wang, "Sign language recognition with attention-based capsule networks", Knowledge and Information Systems, vol. 65, no. 2, pp. 754-772, Feb. 2023.
- [20] Y. Unal, F. Temizel, "Sign language word segmentation and recognition using a multi-stream CNN architecture", Neurocomputing, vol. 507, pp. 130-141, Mar. 2022.