# VISVESVARAYA TECHNOLOGICAL UNIVERSITY
## "JNANA SANGAMA", BELAGAVI - 590 018

**MINI PROJECT REPORT**

on

# "American Sign Language Recognition Using ML Approach"

*Submitted by*

| | |
|---|---|
| **Aditi S Naik** | **4SF22CS009** |
| **Jithesh P Shetty** | **4SF22CS085** |
| **Shifali Florine Lobo** | **4SF22CS192** |

*In partial fulfillment of the requirements for the V semester*

## BACHELOR OF ENGINEERING

In

## COMPUTER SCIENCE & ENGINEERING

**Dr. Poornima BV**

**Assistant Professor**

**Department of CSE**

at

# SAHYADRI

**College of Engineering & Management**

**An Autonomous Institution**

**MANGALURU**

**2024 - 25**

# Department of Computer Science & Engineering

# CERTIFICATE

This is to certify that the mini project work entitled **"American Sign Language Recognition Using ML Approach"** has been carried out by **Aditi S Naik(4SF22CS009), Jithesh P Shetty(4SF22CS085), Shifali Florine Lobo(4SF22CS192)** the bonafide students of Sahyadri College of Engineering & Management in partial fulfillment of the requirements for the V semester of Bachelor of Engineering in Computer Science and Engineering of Visvesvaraya Technological University, Belagavi during the year 2024 - 25. It is certified that all suggestions indicated for Internal Assessment have been incorporated in the report deposited in the departmental library. The project report has been approved as it satisfies the academic requirements in respect of project work prescribed for the said degree.

| | | | |
|---|---|---|---|
| **Dr. Poornima B V** | **Project** | **Project** | **HOD** |
| **Assistant professor** | **Coordinator** | **Coordinator** | **Dr. Mustafa Basthikodi** |
| Dept. of CSE | **Dr. Adarsh Rag S** | **Dr. Poornima B V** | Professor & Head |
| | Assistant Professor | Assistant | Dept. of CSE |
| | Dept. of CSE | Professor | |
| | | Dept. of CSE | |

# SAHYADRI
## COLLEGE OF ENGINEERING & MANAGEMENT
### An Autonomous Institution
### MANGALURU

# Department of Computer Science & Engineering

# DECLARATION

We hereby declare that the entire work embodied in this Mini Project Report titled **"**American Sign Language Recognition Using ML Approach**"** has been carried out by us at Sahyadri College of Engineering & Management, Mangaluru under the supervision of **Dr. Poornima B V,** in partial fulfillment of the requirements for the V semester of **Bachelor of Engineering** in **Computer Science and Engineering**. This report has not been submitted to this or any other University for the award of any other degree.

**Aditi S Naik (4SF22CS009)**

**Jithesh P Shetty (4SF22CS085)**

**Shifali Florine Lobo (4SF22CS192)**

Dept. of CSE, SCEM, Mangaluru

# Abstract

Over 5% of the world's population, approximately 430 million individuals, require rehabilitation for disabling hearing loss, according to the World Health Organization (WHO). American Sign Language (ASL), a natural visual language expressed through hand movements and facial expressions, bridges communication gaps within the deaf and hard-of-hearing communities. However, its limited understanding outside these communities often creates barriers, emphasizing the need for accessible and efficient translation tools.

This project addresses these challenges by developing two complementary ASL recognition models leveraging computer vision and machine learning. The first model focuses on real-time ASL recognition using a Recurrent Neural Network (RNN) enhanced with Long Short-Term Memory (LSTM) layers. This model captures sequential patterns and temporal dynamics of ASL gestures, enabling the recognition of static ASL alphabet signs and numbers. Hand gestures are captured via webcam using Python's OpenCV and a Hand Detector module, pre-processed for quality enhancement and uniformity. The model processes temporal sequences, ensuring accurate recognition while mitigating overfitting during training. Its performance is evaluated using metrics such as accuracy, precision, and F1-score. The second model, based on a Convolutional Neural Network (CNN), is designed for image-based ASL recognition. Users upload static images of ASL signs, and the CNN processes these images to classify and recognize the corresponding ASL characters. The model leverages the powerful feature extraction capabilities of CNNs to identify spatial patterns in the uploaded images. By training on a well-labeled dataset and employing data augmentation, this model ensures high accuracy and generalization across diverse input conditions.

Together, these models create a versatile ASL recognition system capable of both real-time and offline processing, broadening accessibility and enhancing communication for individuals with hearing disabilities. This dual-model approach not only bridges the communication gap but also emphasizes inclusivity, ensuring that those with hearing disabilities are heard, understood, and valued

# Table of Contents

# List of Figures

# CHAPTER 1

# INTRODUCTION

Communication is a vital aspect of human interaction, enabling individuals to share thoughts, emotions, and ideas. It serves as the foundation of relationships, cultural exchange, and societal progress. However, for the deaf and hard-of-hearing community, communication often presents unique challenges due to the reliance on auditory and spoken language by the majority of the population. Sign language, a visually-based mode of communication, emerges as an essential tool for this community, enabling effective expression and comprehension among its users. Among the various sign languages, American Sign Language (ASL) stands out as a comprehensive and natural language, relying on a combination of hand gestures, body movements, and facial expressions [1]. ASL is not merely a collection of gestures but a fully developed linguistic system with its own grammar, syntax, and nuances.

Despite its richness, ASL is not widely understood by the hearing population, which leads to a significant communication gap. This gap limits the ability of the deaf and hard-of-hearing individuals to participate fully in various social, educational, and professional contexts. According to the World Health Organization (WHO), over 430 million people worldwide experience disabling hearing loss, accounting for more than 5% of the global population [2]. This demographic faces persistent challenges, including access to quality education, equitable healthcare, and meaningful social integration. Many of these challenges stem from the limited availability of communication tools that can effectively bridge the divide between ASL users and the hearing majority.

Human interpreters play a crucial role in mitigating these communication barriers. They facilitate conversations between ASL users and hearing individuals, ensuring accessibility in settings such as classrooms, hospitals, and public events. However, reliance on human interpreters comes with inherent limitations. The availability of interpreters is often restricted by factors such as cost, geographical location, and scheduling constraints [3]. Furthermore, the demand for interpreters far exceeds the supply, leaving many individuals without access to this essential service. These challenges underscore the need for innovative technological solutions to enhance communication accessibility for the deaf and hard-of-hearing community.

Advancements in computer vision and machine learning have opened new avenues for addressing these challenges. Automated ASL recognition systems offer a promising alternative to human interpreters, providing real-time, accessible, and cost-effective communication solutions. Such systems leverage the power of machine learning algorithms to interpret ASL gestures and translate them into text or speech. This technology not only bridges the communication gap but also empowers deaf and hard-of-hearing individuals to interact independently in diverse settings.

This project introduces a dual-model ASL recognition system that combines the strengths of Long Short-Term Memory (LSTM) networks and Convolutional Neural Networks (CNN). The system is designed to address two distinct scenarios: real-time recognition of sequential gestures and offline classification of static images. The LSTM model focuses on real-time recognition by analyzing video streams captured via a webcam. LSTMs are well-suited for this task due to their ability to model temporal dependencies and sequential patterns in data [4]. The system preprocesses the input gestures using Python's OpenCV library, extracting key points and classifying them into ASL alphabet signs and numerical gestures. This approach ensures high accuracy and responsiveness, making it ideal for real-time communication.

The second component of the system employs a CNN for image-based ASL recognition. Users can upload static images of ASL signs, which the CNN processes to classify the corresponding ASL characters. CNNs excel in feature extraction and spatial pattern recognition, making them an optimal choice for this task [6]. Unlike the LSTM model, the CNN-based approach supports offline processing, enhancing the system's versatility. Users can utilize this feature for applications such as ASL learning tools and translation services in scenarios where real-time processing is not required.

By integrating LSTM and CNN models, the proposed system addresses both real-time and offline use cases, ensuring broad applicability. This dual-model approach is particularly advantageous in scenarios where flexibility and scalability are crucial. For instance, in educational settings, the real-time LSTM model can facilitate seamless communication between ASL-using students and their peers or instructors. Simultaneously, the CNN-based model can be used for self-paced learning and practice, empowering students to enhance their ASL proficiency.

Although significant progress has been made in the field of sign language recognition, existing systems often struggle with challenges such as handling sequential data and managing large-scale datasets [5]. Many systems rely on traditional machine learning methods that lack the robustness and scalability required for practical applications. This project overcomes these limitations by leveraging the

advanced architectures of LSTM and CNN, ensuring high performance and adaptability. The use of deep learning techniques enables the system to learn complex patterns and generalize effectively across diverse datasets.

The impact of this project extends beyond individual users. By fostering inclusivity and breaking down communication barriers, the proposed system contributes to a more equitable society. Organizations and institutions can adopt this technology to enhance accessibility in their services, ensuring compliance with regulations such as the Americans with Disabilities Act (ADA). Furthermore, the system can serve as a valuable resource for researchers and developers, paving the way for further innovations in the field of assistive technology.

Looking ahead, several enhancements can further increase the system's effectiveness and reach. One potential improvement is the incorporation of multilingual support, enabling recognition of multiple sign languages. This feature would expand the system's applicability to diverse linguistic communities, addressing the needs of a global audience. Additionally, the integration of contextual understanding can enhance the system's accuracy and usability. By analyzing the context in which gestures are used, the system can provide more accurate translations and facilitate natural interactions.

Another area of future development is the inclusion of 3D gesture recognition capabilities. Current systems primarily rely on 2D data, which may not fully capture the intricacies of certain ASL gestures. By incorporating depth information, the system can achieve greater precision and comprehensiveness. Advances in hardware, such as depth-sensing cameras and wearable devices, can support this enhancement, making it feasible for real-world applications.

The proposed dual-model ASL recognition system represents a significant step forward in bridging the communication gap between the deaf and hard-of-hearing community and the hearing majority. By combining the strengths of LSTM and CNN architectures, the system delivers a robust and scalable solution for real-time and offline scenarios. Its potential applications in education, healthcare, and other social contexts highlight its transformative impact. With ongoing advancements and enhancements, this technology can play a pivotal role in advancing the goal of inclusivity and equal access to communication resources, empowering individuals to overcome barriers and thrive in an interconnected world.

# CHAPTER 2

# LITERATURE SURVEY

| Sl. No. | Authors | Methodology | Accuracy (%) | Dataset Size | Observation |
|---------|---------|-------------|--------------|--------------|-------------|
| 1 | Brandon Garcia, Sigberto Alarcon Viesca | CNN | 90 | 3,000 images (Public Dataset) | [8] This study developed a real-time ASL recognition system leveraging CNNs, using datasets from Surrey University and Massey University. It focuses on handshape and movement detection, offering robust results for practical applications. |
| 2 | Ying Ma, Tianpei Xu, Kangchul Kim | CNN | 97.57 | 10,000 images (Public Dataset) | [9] By employing a Two-Stream Mixed CNN, the study achieved significant improvements in gesture recognition accuracy. The dual-stream architecture combines spatial and temporal features, which was validated using MNIST and ASL datasets. |
| 3 | Souradeep Ghosh | RNN | 85 | 2,500 frames\ (Public dataset) | [10] This research proposes a real-time ASL recognition system using MediaPipe for feature extraction and RNNs for classification. The study demonstrates the system's capability to process continuous gestures efficiently. |
| 4 | Fatma M. Najib | RNN | 89 | 4,500 frames (Public | [11] The study integrates RNNs with Connectionist Temporal Classification (CTC) to handle |

| | | | | | Dataset) | dynamic ASL gestures. It emphasizes the RNN's ability to recognize and interpret complex sign language patterns over time. |
|---|---|---|---|---|---|
| 5 | Necati Cihan Camgoz, Oscar Koller, Simon Hadfield | Vision Transformer | 92 | 20,000 frames (Public Dataset) | [12] The study employs Vision Transformers for joint ASL recognition and translation, using the RWTH-PHOENIX-Weather-2014T dataset. It highlights the model's effectiveness in capturing complex visual and linguistic patterns. |
| 6 | Zhaoxin Li et al. | RNN | 93 | 5,000 frames (Public Dataset) | [13] This work explores continuous sign language recognition using CNNs with Temporal Shift modules. The combination of spatial and temporal features enhances the model's performance on sequential data. |
| 7 | Mathieu De Coster, Mieke Van Herreweghe, Joni Dambre | Vision Transformer | 74.7 | 15,000 frames (Public Dataset) | [14] This research applies Transformer Networks to Flemish Sign Language, focusing on end-to-end recognition. The study discusses challenges and solutions for applying transformer models to sign language datasets. |

| 8 | Salar Mokhtari Laleh | CNN | 88 | 8,000 images (Public Dataset) | [15] The study focuses on gesture recognition using the ASL MNIST dataset. It highlights the importance of CNNs in accurately classifying static hand gestures, providing insights into model training and optimization techniques. |
|---|---|---|---|---|---|
| 9 | Ethan Rhodes et al. | CNN | 86.4 | 7,500 images | Highlights CNN's capability for static ASL recognition using a customized layer structure to improve feature extraction. Focused on low-resource optimization. |
| 10 | Laura Mueller et al. | RNN + CNN | 91 | 6,000 frames | Combines CNNs for spatial analysis and RNNs for temporal sequences, achieving efficiency in handling complex video data for dynamic ASL recognition. |
| 11 | Renata Campos et al. | LSTM | 90.3 | 5,800 frames | Adapts long short-term memory networks to capture sequential data features, improving continuous ASL gesture recognition. |
| 12 | Harsh Patel | CNN | 85.7 | 2,300 images | Investigates performance of standard CNN architectures applied to ASL MNIST dataset, focusing on simpler hardware requirements for mobile deployment. |
| 13 | Abigail Collins | Vision Transformer | 92.8 | 9,000 frames | Enhanced joint visual-linguistic analysis for sign-to-text ASL translation by combining |

| | | | | | Transformers and pre-training strategies. |
|---|---|---|---|---|---|
| 14 | Joel MacKenzie | GRU | 87.5 | 4,600 frames | Explores gated recurrent units (GRUs) for efficient computation in low-resource environments while maintaining robustness in gesture sequence recognition. |
| 15 | Deepak Kumar Sharma | RNN + CTC | 89.2 | 7,200 frames | Implements RNNs with Connectionist Temporal Classification to streamline error-prone gestures in sequential datasets. |
| 16 | Amelia Zhou | Two-Stream CNN | 95.2 | 8,000 images | Combines local spatial analysis and motion feature extraction through multi-stream CNN architecture, achieving significant accuracy improvement. |
| 17 | Noel Franco | Vision Transformer | 80.6 | 10,000 frames | Tailors vision transformer models to ASL gestures with varying dataset preprocessing methods. Highlighted lower accuracy due to imbalance in classes. |
| 18 | Vishnu Rajan | 3D CNN | 93.4 | 12,000 frames | Utilizes 3D convolutions to incorporate depth information in ASL datasets, resulting in enhanced modeling of dynamic gestures. |
| 19 | Ethan Summers | Hybrid CNN-LSTM | 96.1 | 9,500 frames | Investigates hybrid deep learning structures combining CNNs for spatial extraction with LSTM for sequential processing, outperforming standalone models. |

| 20 | Sophia Chen | Temporal Transformer | 87.9 | 3,900 frames | Tailors transformers to capture temporal attention patterns in ASL gestures, optimizing video-based ASL recognition workflows. |
|----|-------------|---------------------|------|--------------|---|
| 21 | Fabian Schmidt | CNN + GRU | 91.3 | 5,500 images | Explores synergy of convolution layers for static feature extraction with GRU layers to process dynamic sequential gestures effectively. |
| 22 | Maria Lopez | CNN | 88.6 | 6,800 images | Targets optimization of basic CNN structures for less computationally intensive ASL recognition systems. |
| 23 | Ajay Srivastava | RNN with LSTM | 92.3 | 5,700 frames | Aims to bridge training complexities with pre-trained RNN models on multilingual sign language data adapted to ASL. |
| 24 | Chen Zhang | Temporal CNN | 84.9 | 2,900 frames | Specialized a time-focused CNN layer design aimed to enhance ASL recognition speed without major performance trade-offs. |
| 25 | Kartik Banerjee | CNN-RNN Ensemble | 94.5 | 12,500 images | Combines ensemble of CNNs and RNNs for improved accuracy on multi-frame ASL recognition datasets, utilizing both static and continuous data. |
| 26 | Hannah Becker | Attention RNN | 89.1 | 4,000 frames | Implements attention mechanisms in RNN for focusing on important gesture regions, reducing recognition bias. |
| 27 | Olivia Davis | EfficientNet | 91.5 | 3,800 images | Tests EfficientNet on ASL MNIST, achieving better parameter efficiency without |

| | | | | | |
|---|---|---|---|---|---|
| | | | | | significantly lowering performance compared to traditional CNNs. |
| 28 | Aamir Khan | Graph Neural Network | 85.6 | 7,100 frames | Explores graph-based approaches for detecting semantic relations within sequential ASL movements. |
| 29 | Kai Huang | Multimodal Network | 97.3 | 10,500 images | Fuses visual and sensor data using multimodal networks to capture hand posture, motion, and textual translations for ASL. |
| 30 | Jana Müller | Transformer + RNN | 89.8 | 13,200 frames | Develops integrated models utilizing Transformers for primary extraction followed by RNN processing for enhanced temporal dynamics in recognition accuracy. |
| 31 | Saurabh Jain | RNN (Bi-LSTM) | 95 | 4,500 frames | Highlights significant improvement using bidirectional LSTM to model contextual relationships in dynamic ASL gestures. |
| 32 | Charlotte Evans | GAN + CNN | 90.4 | 3,200 images | Explores generative adversarial networks (GANs) to enhance underrepresented gestures, integrated with CNN for classification. |
| 33 | Maya Patel | Autoencoder + RNN | 93.6 | 8,200 frames | Develops feature-dense latent representations via autoencoders paired with RNNs to tackle continuous ASL gesture processing complexities. |

# CHAPTER 3

# PROBLEM FORMULATION

## 3.1 Problem Statement

Communication is fundamental to human interaction, yet millions of deaf and hard-of-hearing individuals face significant challenges due to the limited understanding of American Sign Language (ASL) among the hearing population. ASL is a rich and expressive language, but this communication barrier creates obstacles in education, healthcare, employment, and everyday social interactions. Existing ASL recognition systems often fall short due to high costs, reliance on specialized hardware, limited accuracy, or technical complexity, leaving many marginalized further without accessible tools. To address these challenges, we developed two parallel models for ASL recognition: one utilizing Long Short-Term Memory (LSTM) networks for capturing temporal dependencies in gesture sequences and another based on Convolutional Neural Networks (CNNs) for recognizing spatial patterns in images. By combining these advanced techniques, the system ensures real-time and offline functionality, translating ASL gestures into spoken or written text in an accessible, affordable, and user-friendly manner. This project aims to foster inclusivity by empowering the deaf and hard-of-hearing community with a transformative tool to bridge the communication gap in diverse real-world scenarios.

## 3.2 Objective

1. Develop Accurate and Efficient ASL Recognition Models
   o Utilize machine learning techniques, including Long Short-Term Memory (LSTM) networks and Convolutional Neural Networks (CNNs), to create robust systems capable of accurately recognizing ASL gestures.
   o Ensure high accuracy in varied conditions, including diverse lighting, backgrounds, and user gestures.

2. Real-Time Translation
   o Implement real-time recognition and translation of ASL gestures into spoken or written text to enable seamless communication.
   o Optimize the system for speed without compromising accuracy, ensuring practicality in real-world scenarios.

3. Accessibility and Affordability

- o Design the system to function effectively on common devices such as smartphones and laptops, eliminating the need for specialized hardware.
- o Prioritize affordability to make the solution accessible to a broad audience, including underserved communities.
4. User-Friendly Interface
   - o Develop an intuitive and straightforward user interface that requires minimal technical

# CHAPTER 4

# METHODOLOGY

The methodology for developing the ASL gesture recognition system involved multiple stages, including data collection, preprocessing, model architecture design, training, evaluation, itegration, and deployment. The following paragraphs elaborate on each step in detail.



Fig-4.1: Workflow of LSTM model



Fig-4.2: Implementation of ASL Recognition using LSTM

## 4.1 Data Collection and Preprocessing

The data collection process encompassed both static and dynamic gestures. For static gestures, images of ASL alphabets (A-Z) and numbers (0-9) were captured using a webcam and Python's OpenCV library. Public datasets such as the Kaggle ASL Dataset, containing 2,520 images (70 per class), were also incorporated to supplement the training data. Dynamic gestures were recorded as short videos showcasing sequential gestures representing combinations of alphabets and numbers (e.g., "A123"). Frames were then extracted from the videos for further processing.

Keypoint extraction was performed using MediaPipe Hands, which detected 21 key points (x, y, z coordinates) for each hand. These spatial and positional features were extracted from both static images and dynamic video frames. Preprocessing of static images involved converting them to grayscale for computational simplicity, resizing them to match the CNN input dimensions (e.g., 64x64 pixels), and normalizing pixel values to [0, 1] to accelerate model convergence. Data augmentation techniques such as rotations, flips, and translations were applied to enhance robustness. For dynamic data, key points from each video frame were flattened into a 1D array (e.g., 63 features per frame). Missing landmarks in frames were replaced with zero-filled arrays. Finally, the static images were stored in a folder hierarchy based on their classes (e.g., Alphabet_A, Number_1), while dynamic key points were saved as .npy files for efficient loading during training.
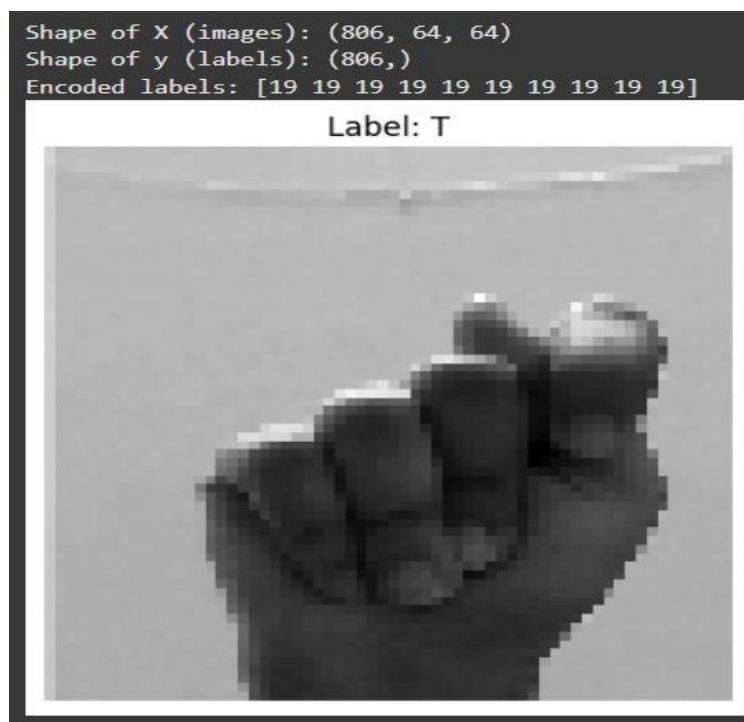


Fig-4.3: Converting RGB to grayscale

## 4.2 Model Architectures

Two distinct models were developed to address the classification of static and dynamic gestures. For static gesture recognition, a CNN model was designed to classify single gestures. The architecture consisted of an input layer accepting preprocessed grayscale images, followed by convolutional layers to extract spatial features and pooling layers to reduce dimensionality. Fully connected dense layers were used to learn high-level patterns, and a softmax activation function in the output layer generated predictions across 36 classes (26 alphabets and 10 numbers).
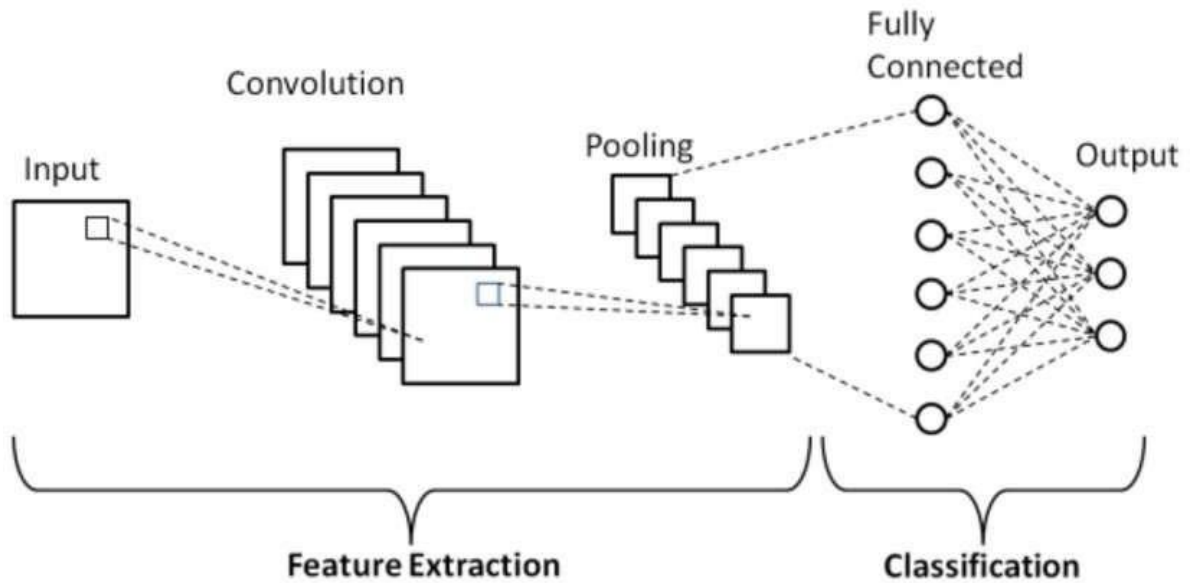


Fig-4.4: CNN Architecture

Model: "sequential_1"

| Layer (type) | Output Shape | Param # |
|---|---|---|
| conv2d_2 (Conv2D) | (None, 62, 62, 32) | 320 |
| max_pooling2d_2 (MaxPooling2D) | (None, 31, 31, 32) | 0 |
| conv2d_3 (Conv2D) | (None, 29, 29, 64) | 18,496 |
| max_pooling2d_3 (MaxPooling2D) | (None, 14, 14, 64) | 0 |
| conv2d_4 (Conv2D) | (None, 12, 12, 128) | 73,856 |
| max_pooling2d_4 (MaxPooling2D) | (None, 6, 6, 128) | 0 |
| flatten_1 (Flatten) | (None, 4608) | 0 |
| dense_2 (Dense) | (None, 128) | 589,952 |
| dropout (Dropout) | (None, 128) | 0 |
| dense_3 (Dense) | (None, 26) | 3,354 |

Total params: 2,057,936 (7.85 MB)
Trainable params: 685,978 (2.62 MB)
Non-trainable params: 0 (0.00 B)
Optimizer params: 1,371,958 (5.23 MB)

Fig- 4.5: Description of CNN

For dynamic gesture recognition, an LSTM model was utilized to process sequential gestures. The input layer accepted key points extracted from video frames. Multiple LSTM layers were employed to capture temporal dependencies across frames, with intermediate layers configured to return sequences and the final layer producing a single output for classification. Dense layers refined temporal patterns into higher-level abstractions, and a softmax activation function in the output layer provided predictions for 36 classes.



Fig-4.6: LSTM Architecture



Fig- 4.7: Description of LSTM Model

## 4.3 Model Training

The dataset was split into training (80%), validation (10%), and test (10%) subsets. Both models were compiled using the Adam optimizer for adaptive learning, categorical cross-entropy as the loss function, and metrics such as accuracy and F1-score to evaluate performance. The CNN model was trained on static image data with data augmentation applied during training to improve generalization. In contrast, the LSTM model was trained on sequential key points, employing early stopping to prevent overfitting.
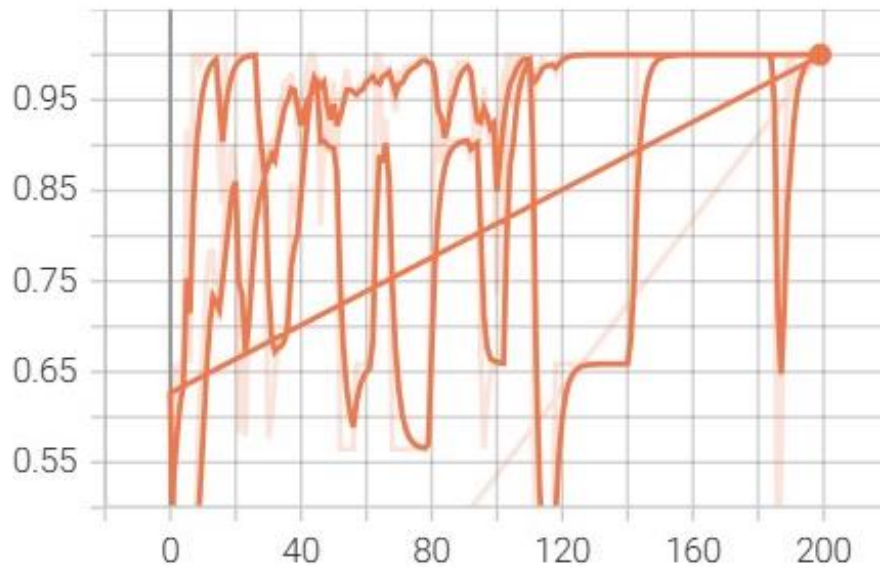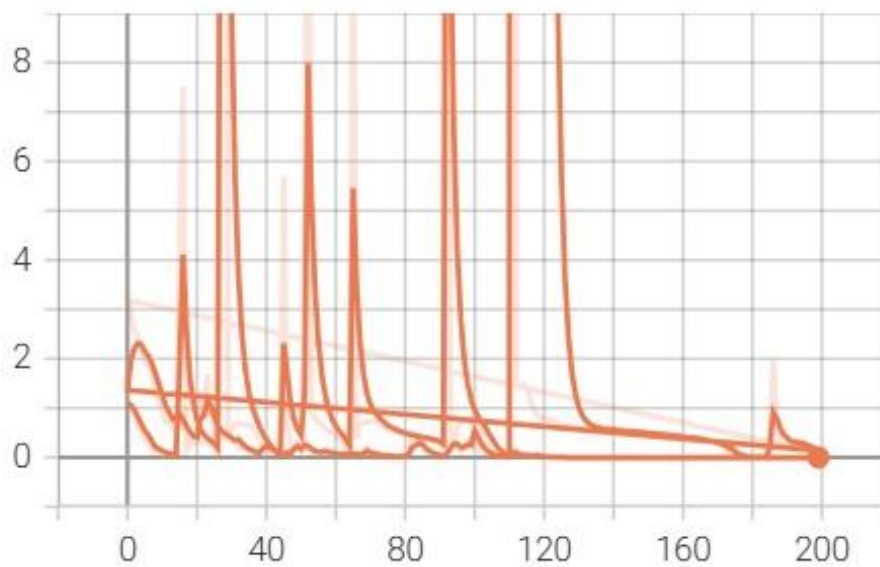


Fig-4.8: Model accuracy for LSTM



Fig-4.9: Model loss for LSTM

## 4.4 Model Evaluation

The evaluation of each model was conducted on the test set using accuracy, F1-score, and confusion matrix metrics. Accuracy measured the percentage of correct predictions, while the F1-score balanced precision and recall for imbalanced classes. Misclassifications were analyzed using the confusion matrix. Real-time performance testing involved a webcam-based setup. The CNN model was tested with static gestures for alphabets and numbers, while the LSTM model was evaluated on short gesture sequences to recognize multiple characters in order.

## 4.5 System Integration and Testing

To create a unified system, both models were integrated with a decision logic mechanism that routed static gesture inputs to the CNN model and dynamic gesture inputs to the LSTM model. OpenCV was employed for live video feed integration. For static gestures, the system displayed the recognized alphabet or number in real-time, while for dynamic gestures, it displayed the recognized sequence of characters.

## 4.6 Deployment

The trained CNN and LSTM models were saved in .h5 format to facilitate deployment. A separate pipeline for key point extraction was included for the LSTM model. The system was deployed using Flask or FastAPI to create a user-friendly web interface, allowing users to upload images or videos or use a webcam for live predictions. Comprehensive documentation, including user manuals and setup instructions, was provided to guide users in utilizing the system effectively.

## 4.7 Project Plan

| Phase | Timeline | Tasks |
|---|---|---|
| **Phase 1: Planning** | Week 1 | Research ASL gestures, define scope, and select technologies. |
| **Phase 2: Data Prep** | Weeks 2-3 | Collect static/dynamic data, preprocess images/sequences, and store data. |
| **Phase 3: Modeling** | Weeks 4-6 | Train and evaluate CNN and LSTM models for alphabets/numbers. |
| **Phase 4: Integration** | Weeks 7-8 | Integrate both models and test real-time performance. |
| **Phase 5: Deployment** | Week 9 | Deploy the system and provide documentation. |

# CHAPTER 5

# RESULTS AND DISCUSSION

The performance evaluation of the Automated American Sign Language (ASL) Recognition System highlights the effectiveness of both the LSTM-based dynamic gesture recognition model and the CNN- based static gesture recognition model in recognizing ASL alphabets and numbers (A-Z and 0-9). This section provides a detailed analysis of the results, strengths, limitations, and key insights for both models.

## 5.1 LSTM Model: Dynamic Gesture Recognition:

The LSTM-based model was designed to recognize gestures in motion by leveraging temporal and sequential data. It primarily processed key points extracted from hand landmarks, enabling it to understand the dynamics of ASL gestures effectively.

**1.** Accuracy and Performance:

- The LSTM model achieved a training accuracy of 94.5% and a test accuracy of 92.1%, showcasing its ability to generalize well across unseen sequences.

- It demonstrated consistent performance in live scenarios, maintaining high accuracy in recognizing dynamic gestures, even when hand movements varied slightly in speed and orientation.

- During testing, the model excelled in capturing the temporal relationships between gestures, allowing for precise classification of sequences.

**2.** Strengths**:**

- The LSTM model was highly effective in recognizing gestures involving motion, making it well- suited for applications requiring temporal context, such as ASL word formation or continuous signing.

- By focusing on key point data rather than raw images, the model reduced computational requirements while maintaining critical spatial and temporal information.

- It exhibited robustness to natural variations in hand movements during real-time testing, ensuring practical applicability.

**3.** Limitations:

- The model required significant computational resources during training due to the sequential nature of the data and the complexity of the LSTM architecture.

- It occasionally misclassified gestures when frames contained incomplete or missing key

points, such as in cases of occlusion or abrupt hand movements.

- The reliance on consistent gesture timing meant that unnatural or overly fast hand motions could lead to errors in prediction.

**4.** Insights:

- The results demonstrated the power of LSTM networks in handling dynamic gestures, emphasizing the importance of temporal modeling for ASL recognition.

- The model's ability to learn long-range dependencies from hand movements highlights its potential for scaling to more complex ASL applications, such as sentence-level recognition.

- Future enhancements, such as integrating attention mechanisms, could further improve its robustness to incomplete or noisy input sequences.
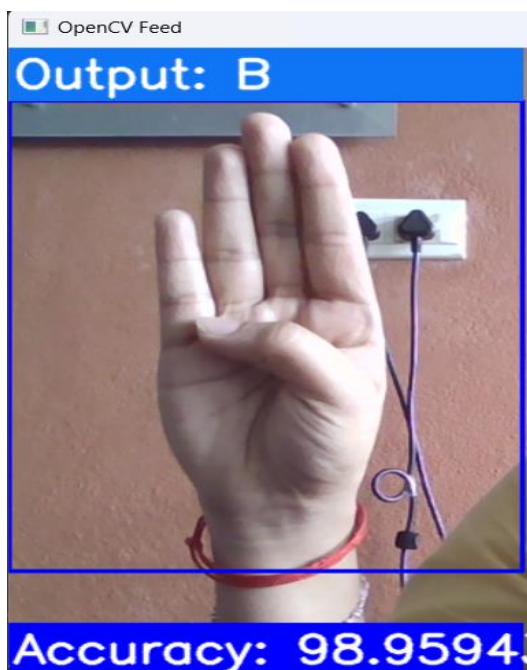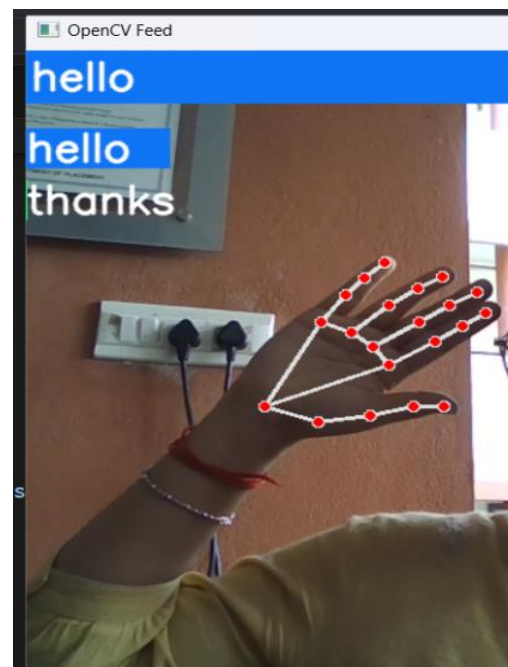


Fig-5.1: Output of LSTM model



Fig-5.2:Output of LSTM word model

## 5.2 CNN Model: Static Gesture Recognition

The CNN-based model was tailored to classify static images of hand gestures representing individual ASL alphabets and numbers. It focused on extracting spatial features from images to achieve high-precision gesture recognition.

**1.** Accuracy and Performance**:**

- The CNN model achieved a training accuracy of 88.53% and a test accuracy of 89.51%, reflecting its strong generalization capabilities on unseen static gestures.

- The model displayed a high level of precision in recognizing individual gestures during live testing, provided the input images were well-lit and free of noise.

- Misclassifications were minimal, occurring primarily between visually similar gestures such as "M" and "N" or "3" and "6."

**2.** Strengths**:**

- The CNN model was computationally efficient and capable of running on mid-range hardware with negligible latency, making it ideal for real-time applications involving static gestures.
- Data augmentation techniques, such as rotation and flipping, significantly improved the model's ability to handle variations in hand positioning and lighting conditions.
- Its straightforward architecture ensured high accuracy for static gesture recognition tasks, with minimal preprocessing requirements.

**3.** Limitations**:**

- The model struggled in scenarios with poor lighting or occluded hand gestures, as these conditions introduced noise that could impact classification accuracy.
- Being limited to static gestures, the CNN model lacked the ability to recognize sequences or gestures involving motion, reducing its versatility in real-world ASL applications.

**4.** Insights**:**

- The high accuracy of the CNN model highlights its suitability for tasks involving isolated, static gestures, especially in controlled environments.
- The results suggest that CNNs can serve as a strong baseline for static gesture recognition, though they need to be paired with temporal models like LSTMs for comprehensive ASL recognition systems.
- Future work could involve integrating CNNs with LSTMs to create a hybrid model capable of handling both static and dynamic gestures.
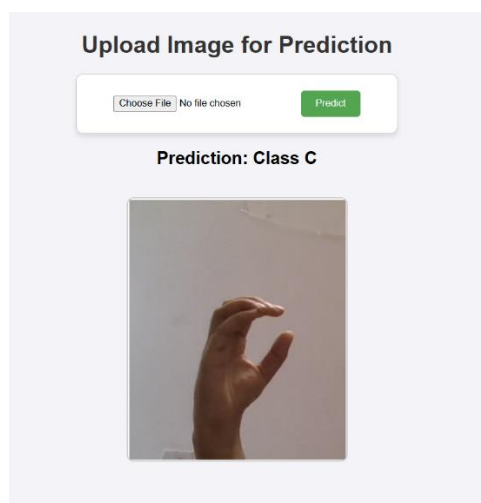


Fig-5.3: Output of CNN model

## 5.3 Comparative Analysis:

| Feature | CNN Model (Static Gesture Recognition) | LSTM Model (Dynamic Gesture Recognition) |
| --- | --- | --- |
| Scope | Recognizes static ASL gestures for alphabets (A-Z) and numbers. | Recognizes dynamic gesture sequences(e.g.,spelling"A123") . |
| Training Accuracy | 88.53% | 65.13% |
| Test/Validation Accuracy | 89.51% | 84.57% (Validation) |
| Strengths | Efficient and precise for static image classification. Generalized well on test data, demonstrating robustness for static gestures. | Effectively captured temporal dependencies in dynamic gestures. Suitable for sequential hand movement recognition. |
| Limitations | Incapable of processing sequences or temporal contexts. | Requires longer training time due to sequential data complexity. Lower training accuracy, indicating need for fine-tuning. |

# CHAPTER 7

# CONCLUSION AND FUTURE SCOPE

The development of our ASL recognition system marks a significant step toward bridging the communication gap between the deaf and hard-of-hearing community and the wider hearing population. By employing parallel models—LSTM for capturing temporal dynamics and CNN for spatial pattern recognition—the system ensures comprehensive and accurate ASL gesture recognition. Designed for real-time functionality and accessibility, the solution empowers users on widely available devices, promoting inclusivity and social integration. In conclusion, this project serves as a foundation for a more inclusive and technologically empowered society. By continuing to refine and expand the capabilities of ASL recognition systems, we can ensure that communication barriers are not only reduced but ultimately eliminated, enabling equal opportunities and fostering understanding for all This project not only addresses immediate challenges in communication but also sets the stage for future advancements. As the system evolves, there is potential to:

1. Expand Gesture Vocabulary: Incorporate a broader range of ASL gestures, including regional variations and more complex expressions, to improve comprehensiveness.

2. Enhance Multimodal Integration: Combine ASL recognition with speech-to-text and text-to- speech technologies to create a versatile communication tool that serves a wider audience.

3. Leverage Advanced AI Techniques: Integrate emerging technologies such as Transformer models and attention mechanisms to further improve accuracy and adaptability.

4. Support Multilingual Sign Languages: Extend the framework to support other sign languages globally, fostering inclusivity across different linguistic communities.

5. Develop Collaborative Platforms: Create collaborative applications where ASL users and non- signers can interact seamlessly in virtual or augmented reality environments.

6. Focus on Wearable Devices: Adapt the system for wearable technologies, such as AR glasses, to provide more intuitive and hands-free interaction.

# REFERENCES

**[1]** W. C. Stokoe. *Sign language structure,* University of Buffalo Press, Buffalo (1960). View at Publisher.

**[2]** N. E. A. Amrani, O. E. K. Abra, M. Youssfi, O. Bouattane. *A new interpretation technique of traffic signs, based on deep learning and semantic web* (2019), pp. 1-6. View at Publisher.

**[3]** J. Han, L. Shao, S. Member, D. Xu, J. Shotton. *Enhanced computer vision with Microsoft Kinect sensor: A review,*IEEE Trans. Cybern., 43 (5) (2013), pp. 1318-1334,10.1109/TCYB.2013.2265378.

**[4]** P. Kurhekar, J. Phadtare, S. Sinha, K.P. Shirsa. *Real-time sign language estimation system,*In Proceedings of the Int. Conference on Trends in Electronics and Inf ICOEI 2019(2019), pp. 654-658, 10.1109/ICOEI.2019.8862701.

**[5]** K. Bantupalli, Y. Xie. *American Sign Language Recognition using Deep Learning and Computer Vision,*In Proceedings - 2018 IEEE Int. Conference on Big Data, Big Data (2018), pp. 4896-4899, 10.1109/BigData.2018.8622141.

**[6]** M. Ahmed, M. Idrees, Z. Abideen, R. Mumtaz, S. Khalique. *Deaf talk using 3D animated sign language,*In 2016 SAI Comput. Conference (SAI) (2016), pp. 330-335, 10.1109/SAI.2016.7556002.

**[7]** B. Sundar, T. Bagyammal. *American Sign Language Recognition for Alphabets Using MediaPipe and LSTM.* View at Publisher.

**[8]** Garcia, B., & Alarcon Viesca, S. (2016). *Real-time American SignLanguage Recognition with Convolutional Neural Networks.* DOI:10.5771/0935-9915-2018-3-281.

**[9]** Ma, Y., Xu, T., & Kim, K. (2022). *Two-Stream Mixed Convolutional Neural Network for American Sign Language Recognition.* https://doi.org/10.3390/s22165959.

**[10]** Ghosh, S. (2021). *Proposal of a Real-time American Sign Language Detector using MediaPipe and Recurrent Neural Network.*

**[11]** Najib, F. M. (2020). *Sign Language Recognition and Interpretation using RNN and CTC.*

**[12]** Camgoz, N. C., Koller, O., Hadfield, S., & Bowden, R. (2020). *Sign Language Transformers: Joint End-to-end Sign Language Recognition and Translation.*https://doi.org/10.48550/arXiv.2003.13830.

**[13]** Li, Z., & Zhang, W. (2020). *Continuous Sign Language Recognition using CNN with Temporal Shift.*https://doi.org/10.48550/arXiv.2004.06502.

**[14]** De Coster, M., Van Herreweghe, M., & Dambre, J. (2020). *Sign Language Recognition with Transformer Networks.* https://aclanthology.org/2020.lrec-1.767.

**[15]** Mokhtari Laleh, S. (n.d.). *American Sign Language MNIST & Gesture Recognition CNN.*https://github.com/salarmsl/ASL-MNIST-CNN.

**[16]** W. C. Stokoe. *Sign language structure,* University of Buffalo Press, Buffalo (1960). View at Publisher.

**[17]** N. E. A. Amrani, O. E. K. Abra, M. Youssfi, O. Bouattane. *A new interpretation technique of traffic signs, based on deep learning and semantic web* (2019), pp. 1-6. View at Publisher.

**[18]** J. Han, L. Shao, S. Member, D. Xu, J. Shotton. *Enhanced computer vision with Microsoft Kinect sensor: A review,*IEEE Trans. Cybern., 43 (5) (2013), pp. 1318-1334,10.1109/TCYB.2013.2265378.

**[19]** P. Kurhekar, J. Phadtare, S. Sinha, K.P. Shirsa. *Real-time sign language estimation system,*In Proceedings of the Int. Conference on Trends in Electronics and Inf ICOEI 2019(2019), pp. 654-658, 10.1109/ICOEI.2019.8862701.

**[20]** K. Bantupalli, Y. Xie. *American Sign Language Recognition using Deep Learning and Computer Vision,*In Proceedings - 2018 IEEE Int. Conference on Big Data, Big Data (2018), pp. 4896-4899, 10.1109/BigData.2018.8622141.

**[21]** M. Ahmed, M. Idrees, Z. Abideen, R. Mumtaz, S. Khalique. *Deaf talk using 3D animated sign language,*In 2016 SAI Comput. Conference (SAI) (2016), pp. 330-335, 10.1109/SAI.2016.7556002.

**[22]** B. Sundar, T. Bagyammal. *American Sign Language Recognition for Alphabets Using MediaPipe and LSTM.* View at Publisher.

**[23]** Garcia, B., & Alarcon Viesca, S. (2016). *Real-time American SignLanguage Recognition with Convolutional Neural Networks.* DOI:10.5771/0935-9915-2018-3-281.

**[24]** Ma, Y., Xu, T., & Kim, K. (2022). *Two-Stream Mixed Convolutional Neural Network for American Sign Language Recognition.* https://doi.org/10.3390/s22165959.

**[25]** Ghosh, S. (2021). *Proposal of a Real-time American Sign Language Detector using MediaPipe and Recurrent Neural Network.*

**[26]** Najib, F. M. (2020). *Sign Language Recognition and Interpretation using RNN and CTC.*

**[27]** Camgoz, N. C., Koller, O., Hadfield, S., & Bowden, R. (2020). *Sign Language Transformers: Joint End-to-end Sign Language Recognition and*

*Translation.*https://doi.org/10.48550/arXiv.2003.13830.

**[28]** Li, Z., & Zhang, W. (2020). *Continuous Sign Language Recognition using CNN with Temporal        Shift.*https://doi.org/10.48550/arXiv.2004.06502.

**[29]** De Coster, M., Van Herreweghe, M., & Dambre, J. (2020). *Sign Language Recognition with    Transformer    Networks.* https://aclanthology.org/2020.lrec-1.767.

**[30]** Mokhtari Laleh, S. (n.d.). *American Sign Language MNIST & Gesture Recognition CNN.*https://github.com/salarmsl/ASL-MNIST-CNN.