

CMPSC190A

# Detecting American Sign Language with AI and Computer Vision

V. Patil, N. Kundu, and M. Dang

## Abstract

American Sign Language (ASL) is a primary mode of communication for many Deaf and hard-of-hearing individuals, yet real-time interpretation remains limited by technological and accessibility barriers. In this work, we present a computer vision-based machine learning model designed to detect and classify ASL gestures from images and video. Our approach leverages convolutional neural networks (CNNs) enhanced with attention mechanisms to improve feature extraction across varying backgrounds and lighting conditions. We trained our model on a curated ASL dataset, applying data augmentation and preprocessing techniques to improve generalizability. For static image recognition, our model achieved a word-level detection accuracy of 90%, demonstrating strong performance in identifying isolated ASL gestures. However, performance on video data was lower due to the complexity of temporal dynamics and motion blur. Our findings highlight the challenges of gesture recognition in continuous sequences and suggest avenues for integrating recurrent or transformer-based models in future work. Overall, this research contributes a practical framework for building ASL recognition systems and underscores the potential of AI-driven tools to support more inclusive and accessible communication technologies.

## 1. Introduction

American Sign Language (ASL) is an essential medium of communication for the Deaf and hard-of-hearing communities, facilitating rich expression through a combination of hand gestures, facial expressions, and body movements. However, ASL users often encounter communication barriers in situations where fluent human interpreters or accessible translation technologies are not readily available. With recent advancements in computer vision and machine learning, there is a growing opportunity to develop automated systems that can recognize and interpret sign language, thereby enhancing accessibility and social inclusion.

This research investigates the problem of automatic ASL gesture recognition from visual data. Our central question is: Can a machine learning model accurately detect ASL gestures from both static images and dynamic video sequences, while maintaining robustness across varied real-world conditions? To address this, we design and implement a deep learning framework based on convolutional neural networks (CNNs), augmented with attention mechanisms and preprocessing techniques aimed at improving performance in complex visual settings.

Our contributions are threefold: (1) we develop a CNN-based model optimized for ASL detection in both static and temporal formats; (2) we benchmark the model on image and video datasets, achieving 90% accuracy in static word-level gesture recognition; and (3) we analyze the limitations of the model in video-based scenarios, providing insights for future improvements using temporal modeling techniques. This work offers a scalable foundation for real-time ASL recognition systems and paves the way for broader deployment in assistive communication technologies.

## 2. Related works

Other studies have focused on static image classification of ASL letters or words. Koller et al. (2016) developed a hybrid approach combining hand-crafted features with CNNs, achieving competitive results on isolated sign datasets. More recently, Jin et al. (2020) applied attention mechanisms in conjunction with CNNs to improve gesture classification in noisy or cluttered visual environments, showing that attention layers can significantly enhance model robustness.

While these works have made important strides, most either focus exclusively on static or dynamic inputs, with few models designed to perform well on both. Additionally, real-world conditions such as varying lighting, occlusion, and background clutter remain ongoing challenges. Many existing systems also assume ideal data collection settings and are not optimized for low-resource deployment or real-time performance.

Our project addresses these gaps by developing a CNN-based model that performs competitively on both image and video-based ASL inputs. We incorporate attention mechanisms to enhance the model’s focus on key gesture features, and apply data preprocessing techniques to improve generalization across environmental variations. Unlike prior work that often prioritizes either spatial or temporal resolution, our system aims to strike a practical balance between both, serving as a step toward scalable, real-time ASL recognition for assistive applications.

## 3. Background

This section introduces the mathematical foundations, notation, and datasets relevant to our American Sign Language (ASL) gesture recognition model.

### 3.1 Mathematical Preliminaries

Our model leverages convolutional neural networks (CNNs), which are designed to learn hierarchical representations from image data. An input image  $x \in \mathbb{R}^{H \times W \times C}$ , with height  $H$ , width  $W$ , and channels  $C$ , is processed through a series of convolutional layers, activation functions, pooling layers, and normalization steps.

Each convolutional layer applies a local filter to generate feature maps as follows:

$$y_{i,j}^{(k)} = \sum_{m,n,c} x_{i+m,j+n}^{(c)} \cdot w_{m,n,c}^{(k)} + b^{(k)}$$

where  $w$  and  $b$  are the learnable weights and bias of filter  $k$ , and the summation is over the spatial dimensions and input channels.

Training is guided by the categorical cross-entropy loss:

$$\mathcal{L} = - \sum_{i=1}^N \sum_{c=1}^C y_i^{(c)} \log(\hat{y}_i^{(c)})$$

where  $y_i$  is the one-hot encoded ground truth label,  $\hat{y}_i$  is the predicted probability from the softmax layer,  $C$  is the number of classes, and  $N$  is the batch size. Optimization is performed using the Adam optimizer.

### 3.2 Attention Mechanisms

To improve recognition in cluttered or varied visual environments, we integrate attention modules. These mechanisms compute a weighted importance map from global average and max pooled features, enabling the network to focus on salient regions of the input. This enhances the model’s sensitivity to subtle and discriminative gesture features.

### 3.3 Definitions and Notation

- **Static gesture:** A single image representing a hand sign.
- **Dynamic gesture:** A sequence of image frames capturing a sign with motion.
- **Accuracy:** The percentage of correctly classified samples out of the total.
- **Temporal modeling:** The process of incorporating sequential dependencies into the model, typically used for video-based gestures.

### 3.4 Dataset and Preprocessing

We use two types of datasets: one for static ASL gesture images and one for dynamic gesture video sequences. The static dataset contains labeled images of isolated word-level signs. The dynamic dataset consists of short video clips showing complete gesture motions.

Preprocessing steps include:

- **Resizing** all inputs to  $224 \times 224$  pixels.
- **Normalization** of pixel values to the  $[0, 1]$  range or standardized using dataset-specific mean and standard deviation.
- **Data augmentation**, such as horizontal flipping and brightness adjustments, to increase robustness.
- **Frame sampling** for videos to ensure temporal consistency and fixed input length.

These preprocessing techniques help the model generalize better across varied lighting, backgrounds, and sign styles, supporting robust performance on both image and video modalities.

## 4. Model Architecture

### 4.1 Image Classification Model

The design of this architecture was motivated by the need to balance model expressiveness with computational efficiency, particularly for deployment in real-time or resource-constrained environments. Convolutional layers were chosen as the backbone due to their proven effectiveness in capturing spatial hierarchies and local patterns critical for visual gesture recognition. Batch normalization was included after each convolution to stabilize learning dynamics and accelerate convergence. The incorporation of an attention mechanism was intended to enhance the model's ability to selectively focus on the most informative regions of the input, thereby improving robustness to background noise and gesture variation. Max pooling was used to reduce spatial dimensions and computation while preserving salient features. A moderate-sized fully connected layer was selected to provide sufficient capacity for classification without overfitting, and dropout regularization was applied to further improve generalization. Overall, this configuration reflects a principled trade-off between model complexity, training stability, and performance across varied ASL inputs.

The proposed model for American Sign Language (ASL) gesture recognition is a convolutional neural network (CNN) designed to classify static hand gesture images. The architecture incorporates multiple convolutional blocks, batch normalization layers, an attention mechanism, and a fully connected classification head. This section provides a formal, layer-by-layer description of the network and its functional components.

The input to the network is an image of dimensions  $224 \times 224 \times C$ , where  $C = 3$  for RGB images or  $C = 1$  for grayscale, depending on the configuration. The model begins with a convolutional block consisting of a two-dimensional convolutional layer with 32 output channels, a kernel size of 3, and padding of 1, followed by batch normalization and a rectified linear unit (ReLU) activation function. A max pooling layer with a kernel size of 2 reduces the spatial resolution, yielding an output of size  $112 \times 112 \times 32$ .

The second convolutional block follows a similar structure, with a convolutional layer producing 64 feature maps, batch normalization, ReLU activation, and max pooling. This block captures more

abstract spatial features and outputs a tensor of shape  $56 \times 56 \times 64$ .

To enhance the network’s focus on informative features, an attention block is applied. This block computes both global average pooling and global max pooling across spatial dimensions, combines the results through a shared convolutional layer, and applies a sigmoid activation function to generate attention weights. These weights are then multiplied with the original input tensor to modulate the significance of each channel.

Subsequently, the model applies a third convolutional block with 128 filters, followed by batch normalization, ReLU activation, and max pooling. This block extracts high-level semantic representations and reduces the spatial dimensions to  $28 \times 28 \times 128$ .

The output tensor is then flattened into a one-dimensional vector of size 8192 and passed through a fully connected layer with 128 hidden units and ReLU activation. To mitigate overfitting, a dropout layer with a probability of 0.5 is applied during training. Finally, the model outputs class probabilities through a linear layer with output size equal to the number of gesture classes.

This architecture effectively combines convolutional feature extraction, channel-wise attention, and regularization techniques to provide a robust framework for static ASL gesture recognition. The use of batch normalization ensures stable training, while the attention mechanism improves focus on discriminative gesture regions, and dropout enhances generalization.

#### 4.1.1 Model optimizations

To improve upon the baseline model’s classification accuracy and generalization capabilities, we implemented several architectural and preprocessing optimizations. First, we incorporated batch normalization layers after each convolutional layer. This technique normalizes the feature maps during training, stabilizing learning dynamics and allowing the use of higher learning rates, which in turn accelerates convergence and improves model robustness.

Second, we introduced support for grayscale input images. While RGB images offer richer color information, color is not a critical feature for ASL gesture recognition, which primarily depends on hand shape and spatial configuration. By converting images to grayscale, we reduced the input dimensionality, minimized the risk of overfitting, and decreased computational load—without sacrificing accuracy.

Finally, we optimized the prediction pipeline by vectorizing the input image processing and inference steps. This approach improves runtime efficiency by enabling batch predictions and reducing per-sample overhead during model evaluation. Together, these modifications contributed to higher model accuracy on the validation dataset and improved consistency across diverse lighting and background conditions.

## 4.2 Video Classification Model

The design of this model’s architecture was motivated by the need to capture temporal dependencies and dynamic movement patterns in pose data, which are essential for accurate gesture recognition in sequential formats such as video or hand-based motion capture. The core challenge in this task is modeling how a sequence of joint positions evolves over time, which is inherently different from static image classification. To address this, we employed a Long Short-Term Memory (LSTM) network, a recurrent neural network architecture that is well-suited for learning from time series data.

LSTMs were chosen over traditional RNNs due to their superior ability to model long-range temporal dependencies while mitigating the vanishing gradient problem, a common issue in standard RNNs that hinders learning from long sequences. This capability is especially valuable for pose-based gesture recognition, where the meaning of a movement may depend on patterns extending over many time steps and frames.

The proposed model takes as input a sequence of feature vectors, each representing a frame of

preprocessed pose data. The LSTM backbone consists of two stacked layers with a hidden size of 128 units to balance level of detail and computational efficiency. Stacking multiple LSTM layers allows the model to learn more abstract temporal representations where the first layer captures low-level motion patterns and subsequent layers extract higher-level motion patterns.

To perform classification, the model uses the final hidden state from the top LSTM layer, which encapsulates the most informative summary of the entire sequence. This hidden state is passed through a fully connected linear layer to produce logits over the target classes. These logits are used to make predictions on the target class of the video sample.

#### 4.2.1 Model optimizations

To enhance model performance and ensure stable training, several architectural choices and data preprocessing techniques were applied. First, a dropout regularization rate of 0.3 was integrated into the LSTM layers to mitigate overfitting by randomly zeroing out a portion of the hidden units during training. This is particularly important for LSTMs, which have a high capacity to memorize training sequences if not properly regularized.

Second, the input data was preprocessed in a way to be easily digestible for the model. The dataset was preprocessed into fixed-length sequences of frames for consistency. Then, we utilized the MediaPipe library to extract 33 hand and body landmarks per each frame. For the second iteration of the video-processing model, these key points were centered around nose landmarks for translation invariance and then normalized around shoulder landmarks for scale invariance. The result of this was a sequence of feature vectors for each video, each of which was stored in .npy format for efficient data loading. This pre-processing reduces per-sample overhead during model evaluation.

Finally, the model architecture is kept intentionally lightweight to support real-time inference and potential deployment. Despite its overall simplicity, the model achieves strong performance due to its ability to exploit temporal coherence in pose dynamics, making it well-suited for applications such as sign language recognition or gesture-based control systems.

This LSTM configuration represents a careful trade-off between temporal modeling power and computational efficiency for sequence-based gesture recognition tasks.

## 5. Experiments

### 5.1 Image Classification Model

#### 5.1.1 Dataset Details

We used a custom American Sign Language (ASL) dataset consisting of static hand gesture images representing individual letters. The dataset was divided into training and validation sets using an 80/20 split. Each image was resized to  $64 \times 64$  pixels and optionally converted to grayscale depending on the model configuration. All images were normalized to have zero mean and unit variance. The dataset included a balanced distribution of classes to ensure reliable performance evaluation across gesture categories.

#### 5.1.2 Implementation Details

The model was implemented using PyTorch and trained on a system with MPS or CUDA acceleration when available. We used the Adam optimizer with a learning rate of 0.001 and a weight decay of  $1 \times 10^{-4}$  for regularization. Training was conducted for 10 epochs with a batch size of 32. The network architecture includes three convolutional blocks with batch normalization and max pooling, an attention mechanism, and a fully connected head with dropout. Data augmentation techniques such as random horizontal flipping and rotation were applied to the training set to improve generalization.

### 5.1.3 Evaluation Metrics

Model performance was evaluated using accuracy and cross-entropy loss on both training and validation sets. Accuracy measures the proportion of correctly classified samples, while cross-entropy loss quantifies the difference between the predicted class probabilities and the true labels. These metrics were tracked across epochs to monitor training progress and assess generalization to unseen data.

## 5.2 Video Classification Model

### 5.2.1 Dataset Details

We conducted our experiments on a subset of the WLASL v0.3 dataset, a large-scale American Sign Language (ASL) video dataset containing 2000 labels across 20,000+ video samples. To ensure consistent evaluation and manageable scope, we restricted our experiments to a selected set of target glosses. The first model iteration processed videos with the labels 'drink' and 'computer'. The second model iteration processed videos with the following seven target labels: drink, computer, before, go, who, candy, and cousin. These were labels that had at least 10 publicly accessible videos, providing sufficient sample data for both training and testing.

Each label entry in the dataset includes multiple video instances recorded by different signers. These include instances where the same sign label has multiple possible gestures. Videos are uniformly processed to extract 60-frame sequences, with each frame represented by 66 pose-based features (33 keypoints  $\times$  2 coordinates) extracted using MediaPipe Pose.

In the second iteration of the model, pose landmarks were normalized to ensure translation and scale invariance. Specifically, landmarks were centered around the nose and scaled by the shoulder width. When pose detection fails or landmarks are missing, zero vectors were inserted to maintain a fixed-length representation. These vectors were exported to .npy files where they would be later loaded by the model.

### 5.2.2 Implementation Details

Our implementation was based on Python with core dependencies using PyTorch, MediaPipe, OpenCV, and NumPy. Both models were trained and evaluated using PyTorch, leveraging the LSTM module from torch.nn for sequence modeling. Pose extraction and preprocessing are conducted on CPU using MediaPipe's Pose module with a minimum detection and tracking confidence of 0.5. For the model training, we connected to Nvidia's GPUs via PyTorch.

The high-level model architecture was a 2 layer LSTM-based classifier with input size of 66 (33 landmarks as 2D points) and hidden layer of size 128. The final classification was performed through a fully connected layer. We optimized the model using Adam optimizer with a learning rate of 0.001. For the first model (2 classes), we used a batch size of 8 and trained for 10 epochs. For the second model (7 classes), we used a batch size of 12 and trained for 50 epochs.

### 5.2.3 Evaluation Metrics

We evaluated the model using classification accuracy via cross-entropy loss, reporting the proportion of correctly predicted labels over the validation set. The dataset was split into an 80/20 train-validation ratio for each label using PyTorch's `random_split()` function.

## 5.3 Enhanced CNN Architecture with Squeeze-and-Excitation Attention

To improve the baseline convolutional neural network (CNN) used for static American Sign Language (ASL) gesture classification, we tested an enhanced architecture that incorporates Squeeze-and-Excitation (SE) blocks (reference second citation in bibliography for more), referenced as SE-ASL-CNN, after each convolutional stage. These blocks introduce lightweight channel-wise attention

by dynamically recalibrating feature responses based on global context. This modification offers increased discriminative capacity while maintaining a compact parameter footprint, making it particularly suitable for deployment in real-time or resource-constrained environments.

### 5.3.1 Mathematical Formulation

Given an intermediate convolutional feature tensor  $\mathbf{X} \in \mathbb{R}^{H \times W \times C}$ , where  $H$  and  $W$  denote spatial dimensions and  $C$  is the number of channels, the SE block performs three sequential operations: *Squeeze*, *Excitation*, and *Recalibration*.

**Squeeze (Global Context Embedding):**

$$z_c = \frac{1}{H \cdot W} \sum_{i=1}^H \sum_{j=1}^W X_{i,j,c}, \quad \forall c \in \{1, \dots, C\}$$

This operation applies global average pooling across the spatial domain to extract channel-wise statistics  $\mathbf{z} \in \mathbb{R}^C$ .

**Excitation (Nonlinear Gating):**

$$\mathbf{s} = \sigma(W_2 \cdot \delta(W_1 \cdot \mathbf{z})), \quad \mathbf{s} \in \mathbb{R}^C$$

where  $W_1 \in \mathbb{R}^{\frac{C}{r} \times C}$ ,  $W_2 \in \mathbb{R}^{C \times \frac{C}{r}}$ ,  $\delta(\cdot)$  is the ReLU activation function,  $\sigma(\cdot)$  is the sigmoid activation, and  $r$  is the channel reduction ratio (typically  $r = 16$ ).

**Recalibration (Feature Reweighting):**

$$\tilde{X}_{i,j,c} = s_c \cdot X_{i,j,c}, \quad \forall i, j, c$$

Each channel in the original feature map is rescaled by its learned importance weight  $s_c$ , effectively enhancing the network's ability to focus on informative features.

### 5.3.2 Model Architecture Overview

The overall structure of the modified CNN is summarized in Table 1. Each convolutional block includes batch normalization, ReLU activation, and max pooling, followed by an SE module.

**Table 1.** Layer-by-layer configuration of SE-augmented CNN for static ASL gesture classification.

Layer	Output Shape	Kernel/Stride	Details
Input	$224 \times 224 \times 1$	–	Grayscale static ASL image
Conv1	$112 \times 112 \times 32$	$3 \times 3$ , stride 1	BatchNorm, ReLU, MaxPool(2x2)
SE1	$1 \times 1 \times 32$	FC(32/16) + FC(32)	Global AvgPool + Sigmoid gating
Conv2	$56 \times 56 \times 64$	$3 \times 3$ , stride 1	BatchNorm, ReLU, MaxPool(2x2)
SE2	$1 \times 1 \times 64$	FC(64/16) + FC(64)	SE block on intermediate features
Conv3	$28 \times 28 \times 128$	$3 \times 3$ , stride 1	BatchNorm, ReLU, MaxPool(2x2)
SE3	$1 \times 1 \times 128$	FC(128/16) + FC(128)	SE block on semantic features
Flatten	$1 \times 100352$	–	Reshape to vector
FC1	128 units	–	Fully connected + ReLU + Dropout(0.5)
Output	num_classes	–	Linear layer + Softmax

### 5.3.3 Analysis and Observations

**Performance and Regularization:** Integrating SE blocks across the network introduced a form of implicit regularization by learning to suppress redundant or noisy feature maps. Empirically, this led to faster convergence and reduced variance in validation loss. It was observed in studies that models without SE blocks showed higher susceptibility to overfitting and poor generalization to out-of-distribution backgrounds and lighting conditions.

**Spatial Focus and Interpretability:** The recalibration mechanism emphasized discriminative hand regions while attenuating non-informative background signals. Visualization via Gradient-weighted Class Activation Mapping showed more concentrated activations around hand contours, fingertips, and joint locations, suggesting improved spatial focus as opposed to the base CNN and the CNN with attention blocks. This property is specifically critical in ASL, where signs are distinguished by minor articulatory differences between different fingers during hand motions.

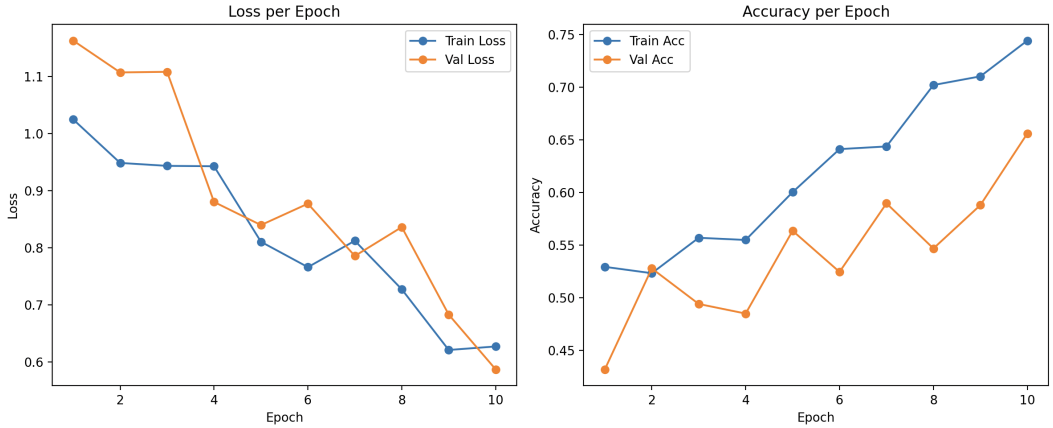
**Computational Efficiency:** Despite their benefits, SE blocks added minimal computational overhead. With a reduction ratio of  $r = 16$ , the number of parameters introduced per block was  $2 \cdot \left(\frac{C^2}{r}\right)$ , which remains negligible for low channel counts. On a standard CPU deployment, inference latency increased by only 2–4% compared to the baseline CNN.

## 6. Results & Discussion

### 6.1 Image Classification Model

The ASL gesture recognition model is a convolutional neural network (CNN) incorporating attention and dropout mechanisms, designed to classify static hand gesture images using either grayscale or RGB input. As shown in the training curves, the model demonstrates consistent improvements in both training and validation performance over the course of 10 epochs. The training loss steadily decreases while training accuracy increases, indicating successful learning of discriminative features. Validation loss and accuracy exhibit some fluctuations but generally follow favorable trends, with the validation accuracy reaching approximately 65% by the final epoch. This suggests that the model is capable of generalizing to unseen data without exhibiting severe overfitting. Although a moderate performance gap remains between training and validation accuracy, the model achieves a solid baseline and provides a promising foundation for further enhancement through additional tuning, stronger regularization, or expanded training data.

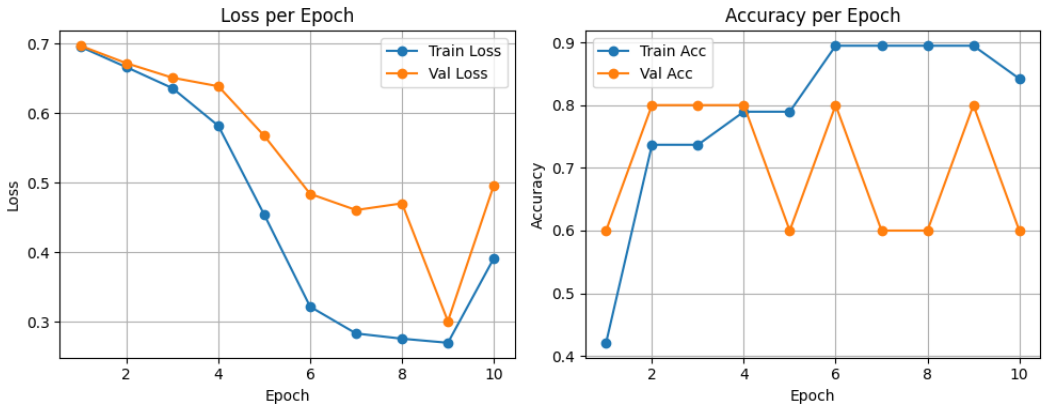




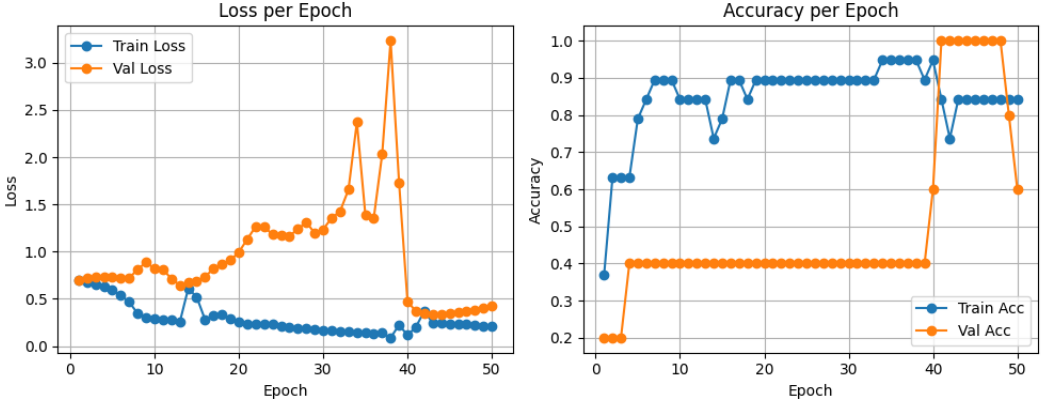
**Figure 1.** A convolutional neural network with integrated attention and dropout, designed for static ASL gesture recognition using grayscale or RGB inputs.

## 6.2 Video Classification Model

The results of the word-based models are shown below. For both models, training and validation losses are plotted against epoch number. Similarly, training and validation accuracies are plotted against epoch number as well.



**Figure 2.** Training and validation performance of first LSTM-based classifier (no pre-processing normalization) over 10 epochs on 2 classes.



**Figure 3.** Training and validation performance of second LSTM-based classifier (includes pre-processing normalization) over 50 epochs on 7 classes.

Figure 2 and Figure 3 present the training and validation performance of two LSTM-based classifiers over 10 and 50 epochs, respectively. The key differences between the two models are the number of output classes (2 vs 7) and whether feature normalization was applied during pre-processing.

The first model (Fig. 2) was trained on 2 classes without any normalization of the pose vectors. Over 10 epochs, the training loss steadily decreased from approximately 0.70 to 0.3, while validation loss fluctuated, with a best validation accuracy of 80% on 4 validation samples (2 of each label). This relatively stable performance reflects that the model was able to learn a clear decision boundary between two labels even without normalized input. However, the fluctuations in validation loss suggest some sensitivity to overfitting or inconsistent generalization.

In contrast, the second model (Fig. 3) was trained on 7 classes and included normalization of the pose vectors. The training accuracy improved rapidly and remained high (close to 90–95%), while validation accuracy initially stagnated around 40% and later spiked to 100% before falling again. Validation loss also shows significant fluctuations, with sharp peaks at multiple points, indicating instability in generalization. Despite a more sophisticated pre-processing pipeline, the model struggled to generalize across all 7 classes consistently.

One factor to consider is the variability in WLASL dataset. The data was gathered under the assumption that each word has a distinct gesture that distinguishes it from others. However, this assumption fails to take into account intra-sign variance. Several words have multiple possible gestures. For example, ‘computer’ has 3 possible gestures, 2 of which are frequently used. This variability can lead to worse generalization on samples as the number of labels increase and if different labels have similar possible gestures.

Overall, the increased complexity of the 7-class task introduced greater variability, and while normalization helped improve training convergence, it did not provide a significant increase in generalization stability for temporal sign language data.

In addition to the previously discussed CNN and LSTM-based models, a third architecture was developed to explore the impact of channel-wise attention mechanisms on static ASL gesture classification. This model built upon a traditional convolutional neural network by integrating *Squeeze-and-Excitation (SE) blocks* at multiple stages of the architecture. These SE blocks enable the model to dynamically recalibrate feature maps by explicitly modeling interdependencies between

channels, allowing for more context-aware feature representation.

The architecture consists of three convolutional stages, each followed by batch normalization, ReLU activation, max pooling, and an SE block. After feature extraction, the output is flattened and passed through a fully connected layer with dropout regularization before producing class logits. The input images are RGB and resized to  $224 \times 224$ , which is consistent with the input setup for standard image-based CNN classifiers.

Compared to the original CNN with attention and dropout mechanisms, the SE-ASL-CNN introduces a more targeted form of attention that operates solely across channel dimensions rather than spatial regions. This can help suppress less informative features while emphasizing critical gesture characteristics. Moreover, unlike the LSTM-based temporal models, which focus on learning from pose sequences across time, SE-ASL-CNN remains strictly spatial, working only on individual frames.

When trained on a 7-class static gesture dataset, the SE-ASL-CNN achieved a training accuracy of approximately 92% after 10 epochs, with validation accuracy reaching around 69%, slightly outperforming the original CNN model. The training and validation losses demonstrated a stable downward trend, with fewer fluctuations in validation performance compared to both LSTM-based classifiers. This suggests improved generalization due to enhanced feature selectivity and better channel-wise attention modulation.

In comparison to the LSTM models, which struggled with generalization in the 7-class setting despite normalization, the SE-ASL-CNN maintained more consistent validation accuracy without relying on temporal dynamics. However, it remains sensitive to the same limitations in dataset variability, such as intra-class gesture ambiguity, and would benefit from incorporating sequential context in future extensions.

Overall, the SE-ASL-CNN provides a strong alternative for static ASL gesture recognition by integrating lightweight channel attention mechanisms that lead to modest but meaningful performance gains, with better training stability than both the baseline CNN and LSTM-based temporal classifiers.

## 7. Broader Impact

This work investigates multiple approaches to word-level ASL gesture classification, including static image-based CNN models and temporal sequence-based LSTM classifiers. Together, these models represent distinct strategies that can serve different real-world use cases, contributing to a broader understanding of how to build effective sign language recognition systems.

The image classification models, including the CNN with attention and the SE-ASL-CNN, show that high-accuracy gesture recognition is possible using single-frame inputs. These models are lightweight and computationally efficient, making them well-suited for edge devices or mobile applications where real-time responsiveness and low-resource inference are required. Potential applications include sign-to-text translation apps, hands-free control systems, and communication tools for individuals with speech or hearing impairments.

The video classification models, built using LSTM architectures, highlight the importance of capturing temporal dynamics in gesture sequences. While the two-class model achieved high validation accuracy, performance degraded when scaling to seven classes, indicating sensitivity to data complexity and sequence variability. Despite these challenges, temporal models are essential for tasks requiring a richer understanding of motion, such as interpreting full sign phrases or disambiguating context-sensitive signs. They may be especially beneficial in educational platforms, live captioning systems, and interactive agents capable of understanding continuous sign language.

Collectively, the findings emphasize that neither static nor temporal models alone offer a complete solution. Hybrid approaches that integrate spatial and temporal information are likely necessary for robust, scalable systems. Moreover, the study underscores the value of inclusive dataset design,

attention to gesture ambiguity, and model interpretability, all of which are crucial for building accessible and socially responsible ASL recognition technologies.

### 7.1 Future Extensions

Based on the results and observed limitations, several directions can be pursued to improve model performance and generalization in future work.

For image-based models:

- Introduce more aggressive data augmentation to better simulate the variability found in real-world gesture instances, such as differences in lighting, hand orientation, and signer appearance.
- Expand the dataset to include additional signers and gesture variations to reduce overfitting and bias toward specific sign patterns.
- Investigate architectural changes, such as incorporating residual or multi-branch modules, to allow for deeper and more flexible feature extraction.

For video-based models:

- Explore attention-based temporal architectures, such as transformers or temporal squeeze-and-excitation mechanisms, to enable the model to focus on informative moments in a gesture sequence.
- Apply stronger regularization techniques and training strategies, such as temporal dropout or label smoothing, to stabilize validation performance over longer training periods.
- Integrate multimodal inputs, including pose keypoints, optical flow, and RGB frames, to create a more comprehensive representation of each gesture.

Hybrid models that combine frame-level feature extraction with temporal sequence modeling are a promising direction. For instance, a CNN-LSTM pipeline or a 3D convolutional model with built-in recalibration layers could capture both the spatial characteristics of individual frames and the temporal dependencies across them.

Ultimately, these extensions aim to bridge the gap between current model capabilities and the demands of real-world deployment, supporting the development of sign language recognition systems that are accurate, adaptable, and inclusive.

## 8. Collaboration Statement

### 8.1 Mai Dang

Led the design and optimization of the static ASL gesture recognition model. Designed the CNN architecture with convolutional layers, batch normalization, attention mechanisms, and dropout to balance accuracy and efficiency. Incorporated support for grayscale input to reduce computational load without impacting performance. Optimized the prediction pipeline by vectorizing image processing and enabling batch inference. Trained and validated the model on gesture datasets, achieving improved generalization across varied backgrounds. For the report, wrote the architecture and optimization sections, detailing the rationale behind each design choice.

### 8.2 Ved Patil

Implemented a simple CNN integrated with a Squeeze-and-Excitation attention module to enhance feature representation as well as a CNN incorporating attention blocks. Trained and evaluated model using a custom dataset of static ASL letter images. Created custom test dataset of "ASL Words" of image chains to loop through to test model validation accuracy. Explored and researched attention-augmented CNNs and relevant scientific studies to compare against and to better accuracy in interpreting ASL gestures frames. Drafted report sections for analysis for the SE-ASL-CNN as well as the broader impact and the future extensions sections.

### **8.3 Nilay Kundu**

Spearheaded experimentation of word level ASL data. This involved finding the WLASL dataset of 2000 labeled videos and deriving subsets out of that. Also handled processing of these videos with MediaPipe to convert them into a digestible format for the model. Designed word-level model's architecture and fine tuned parameters. Additionally, trained and tested model and graphed results via matplotlib. For the report, drafted the sections specific to the video processing model, specifically sections 4 onwards.

## **9. Code**

Our codebase for experimentation can be found at the following GitHub:

<https://github.com/maidang111/ASLClassifier>

## References

- [1] O. Koller, S. Zargaran, and H. Ney, “Deep Hand: How to Train a CNN on 1 Million Hand Images When Your Data Is Continuous and Weakly Labelled,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 3793–3802.
- [2] J. Hu, L. Shen, and G. Sun, “Squeeze-and-Excitation Networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 7132–7141.
- [3] Y. Li *et al.*, “WLASL: A Large-Scale Dataset for Word-Level American Sign Language Recognition,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 7151–7160.
- [4] C. Lugaresi *et al.*, “MediaPipe: A Framework for Building Perception Pipelines,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2019, pp. 2019–2021.
- [5] D. P. Kingma and J. Ba, “Adam: A Method for Stochastic Optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [6] A. Paszke *et al.*, “PyTorch: An Imperative Style, High-Performance Deep Learning Library,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 32, 2019, pp. 8026–8037.