

American Sign Language (ASL) Letter Recognition and Real-Time Translation

Michael Yeung, Richard Huang, Jaya Singh
Professor Dr. Hamed Pirsiavash - ECS174 – Computer Vision
University of California, Davis
December 7, 2025

Abstract

001 We developed a real-time American Sign Language (ASL)
002 letter recognition system that combines computer vision
003 preprocessing with a lightweight convolutional neural net-
004 work trained on the SignMNIST dataset. Our model
005 achieves 94.70% test accuracy on the clean dataset, with
006 12 letters reaching 100% per-class accuracy. However, cer-
007 tain letters, particularly T (62.10%), M (83.50%), and R
008 (85.42%)—remain challenging due to subtle finger config-
009 urations that are difficult to distinguish in 28×28 grayscale
010 images. We deployed the system with an 8-step prepro-
011 cessing pipeline to handle real-world webcam input, tem-
012 poral smoothing using 7-frame majority voting, and acces-
013 sibility features including text-to-speech. The system runs
014 in real time on a CPU-only laptop and achieves approx-
015 imately 80–90% accuracy in real-world conditions. This
016 report documents the methodology, results, limitations, and
017 directions for future work.

018 1. Introduction

019 American Sign Language is used by a large community of
020 deaf and hard-of-hearing individuals in the United States.
021 ASL relies on hand shapes, positions, and movements rather
022 than acoustic speech. Many hearing people never learn
023 ASL, which creates a communication barrier in everyday
024 settings such as classrooms, clinics, and workplaces. Com-
025 puter vision and machine learning provide an opportunity
026 to partially bridge this gap by recognizing hand signs and
027 translating them into text or speech. This project focuses on
028 ASL fingerspelling, where letters are spelled one at a time
029 using hand poses. The goal is to build an end-to-end sys-
030 tem that recognizes 24 static ASL letters (A–Y; J and Z are
031 excluded because they require motion) from a live webcam
032 feed. The system captures video frames, preprocesses them
033 to isolate the hand, feeds a 28×28 grayscale image into a
034 convolutional neural network (CNN), and outputs real-time
035 letter predictions. Temporal smoothing is used to stabilize

predictions, and recognized letters are appended to a buffer
to form words. An optional text-to-speech (TTS) feature
pronounces letters aloud. The work is motivated by acces-
sibility and by the desire to apply concepts from computer
vision and deep learning to a realistic task. The project
demonstrates that a relatively small CNN can reach 94.70%
test accuracy on the SignMNIST dataset, but also shows that
real-world accuracy is lower because of domain differences
between curated data and noisy webcam input. The report
emphasizes both the technical design of the system and the
practical challenges that arise when deploying models out-
side of ideal conditions.

2. Dataset & Model

2.1. Sign MNIST Dataset

The system is trained on the SignMNIST dataset, which
contains grayscale images of hands performing ASL let-
ters. Each image is 28×28 pixels and centered on the hand.
The dataset covers 24 letters (A–Y), excluding J and Z be-
cause those letters are defined by hand motion instead of a
static pose. There are tens of thousands of images in to-
tal, with each class represented by roughly 1,700–2,500 ex-
amples across training, validation, and test splits. SignM-
NIST is a “clean” dataset: background clutter is minimal,
lighting is relatively uniform, and the hand is clearly visi-
ble and centered. These properties make it ideal for training
and benchmarking a classifier, but they also mean that the
dataset does not fully represent real-world webcam condi-
tions. This gap between training data and deployment data
is a central theme in the project.

2.2. Simple CNN Architecture

To classify the 28×28 grayscale images, we designed a com-
pact CNN called SimpleCNN. The model is lightweight
enough to run in real time on a CPU, but has enough ca-
pacity to achieve high accuracy on SignMNIST. The archi-
tecture is:

- Input: 28×28 grayscale image normalized to image.jpg

072	• Block 1: 1→32 convolution (3×3), batch normalization,	on the central area where the user is expected to place	121
073	ReLU, 2×2 max pooling → 14×14 feature map	their hand.	122
074	• Block 2: 32→64 convolution (3×3), batch normalization,	2. Grayscale conversion: The cropped frame is converted	123
075	ReLU, 2×2 max pooling → 7×7 feature map	to grayscale, reducing the input from three color chan-	124
076	• Block 3: 64→128 convolution (3×3), batch normaliza-	nels to one intensity channel and focusing on shape and	125
077	tion, ReLU, 2×2 max pooling → 3×3 feature map	contrast.	126
078	• Classifier: flatten (1,152 features) → fully connected	3. Gaussian blur: A Gaussian filter (e.g., 5×5 kernel with	127
079	layer with 256 units + ReLU + dropout(0.4) → fully con-	= 2) is applied to smooth out noise and small texture	128
080	nected layer with 26 outputs → softmax at evaluation	details that could interfere with thresholding.	129
081	time	4. Histogram equalization: The grayscale image is equal-	130
082	The model has about 400,000 trainable parameters and oc-	ized to normalize brightness and improve contrast,	131
083	cupies less than 2 MB on disk. This small footprint allows	which helps in both thresholding and visual consistency	132
084	the network to be loaded quickly and to run inference at	under different lighting conditions.	133
085	webcam frame rates on a typical laptop CPU.	5. Otsu thresholding: Otsu’s method is used to compute	134
086	2.3. Model Training	an automatic threshold that separates the foreground	135
087	The model was implemented in PyTorch and trained in a	(hand) from the background. The result is a binary im-	136
088	Python 3.10 virtual environment on a MacBook Pro with an	age.	137
089	M1 CPU. The training setup used:	6. Contour detection: Contours are found in the binary	138
090	• Optimizer: Adam	image, and the largest contour is assumed to correspond	139
091	• Learning rate: 1e-3	to the user’s hand.	140
092	• Loss: cross-entropy	7. Bounding box and resize: A bounding box around the	141
093	• Batch size: 128	largest contour is extracted. This region is then resized	142
094	• Epochs: 10	to 28×28 pixels to match the input size expected by the	143
095	• Device: CPU	CNN.	144
096	At the end of each epoch, the model was evaluated on a val-	8. Normalization: Pixel values in the 28×28 crop are nor-	145
097	idation split. Whenever validation accuracy improved, the	malized to the range and reshaped into the format re-	146
098	current weights were saved as the “best” checkpoint. After	quired by the model.image.jpg	147
099	training, the best checkpoint was loaded and evaluated once		
100	on the held-out test set. Training converged quickly. By the	This pipeline is a crucial part of the system. It is what makes	148
101	second epoch, validation accuracy had already reached very	it possible to reuse a model trained on SignMNIST (which	149
102	high values and remained stable, indicating that the model	assumes clean, centered, segmented hands) on noisy web-	150
103	could fit the dataset without difficulty.	cam data.	151
104	3. Pre-processing for Webcam Input		
105	3.1. Motivation	4. Quantitative Result	152
106	The SignMNIST images are small, centered, and clean, but		
107	real webcam frames are not. A raw webcam image includes	4.1. Overall Test Accuracy	153
108	the user’s face, background objects, varying lighting, and		
109	other noise. If these frames are directly downsampled to	Using the best validation checkpoint, the model achieved	154
110	28×28, the hand might only occupy a small part of the im-	94.70% accuracy on the SignMNIST test split of 7,172	155
111	age, and the model will fail. To make the webcam images	images. This result confirms that SimpleCNN is strong	156
112	resemble SignMNIST as closely as possible, the system ap-	enough to solve the static letter classification task on this	157
113	plies a sequence of preprocessing steps to each frame. The	dataset while remaining computationally light. However,	158
114	goal is to isolate the hand, normalize contrast, and reduce	this 94.70% figure corresponds to the curated benchmark	159
115	background clutter before passing the image to the CNN.	setting. When the model is used with webcam input and the	160
116	3.2. Eight-Step Pipeline:	full preprocessing + smoothing pipeline, the effective accu-	161
117	Each frame from the webcam is processed as follows:	racency observed in practice is lower, typically in the 80–90%	162
118	1. Center-crop: The original 640×480 frame is cropped	range depending on lighting, background, and signer con-	163
119	to a 480×480 square around the center of the image.	sistency.	164
120	This removes some peripheral background and focuses		

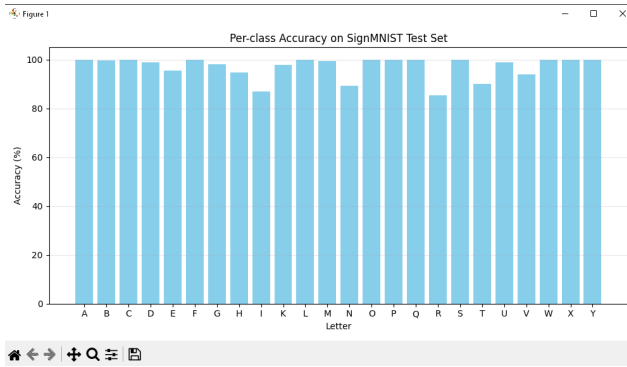


Figure 1. Each letter class' accuracy with the SignMNIST dataset.

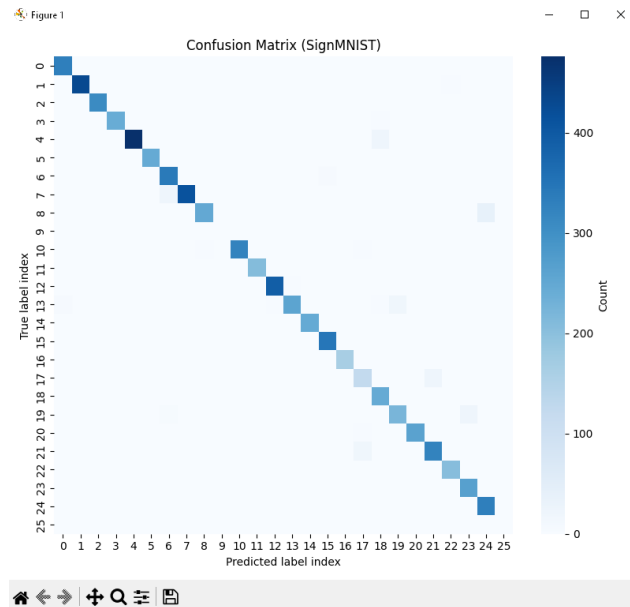


Figure 3. Confusion Matrix heatmap - row : true label index — column : predicted label index, color shows how often each pair occurs

4.2. Per-Class Accuracy

Per-class accuracy reveals which letters are easy or hard for the model. The test results from the .venv310 run are:

Letter	Accuracy	Samples	Letter	Accuracy	Samples
A	100.00%	331	N	88.32%	291
B	100.00%	432	O	97.56%	246
C	99.03%	310	P	100.00%	347
D	100.00%	245	Q	100.00%	164
E	99.80%	498	R	85.42%	144
F	100.00%	247	S	91.06%	246
G	99.43%	348	T	62.10%	248
H	95.41%	436	U	93.98%	266
I	99.31%	288	V	100.00%	346
K	93.66%	331	W	91.26%	206
L	100.00%	209	X	94.01%	267
M	83.50%	394	Y	88.25%	332

Figure 2. Statistics of all 24 letter's accuracy and sample size

These numbers show that:

- Twelve letters are essentially perfect (100% accuracy): A, B, D, F, L, P, Q, V plus several near-100% letters like C, E, G, and I.
- A group of letters (H, K, O, U, W, X, Y, N, R, S) are recognized well but not perfectly, generally in the mid-80s to upper-90s.
- Two letters, T and M, stand out as significantly more difficult, with T at 62.10% and M at 83.50

4.3. Further Analysis

The confusion matrix makes it clear that the errors are not random. Instead, the model tends to confuse letters that are visually similar in static images:

- T is often confused with L and R. All of these involve configurations where fingers cross or overlap, and the subtle differences between them are hard to capture at the 28×28 resolution after contour cropping.
- M is confused with N, R, and Y. These letters all involve multiple fingers being extended or partially folded, making them structurally similar.

These patterns suggest that the limitation is not just model capacity but information content: static 28×28 grayscale snapshots cannot always capture the details needed to distinguish certain finger arrangements. In real ASL, movement and temporal cues also help differentiate some of these letters.

5. Further Discussion

5.1. Real-Time System Behavior

Below are the analysis of components to aid in the processing of data into the model, and their uses:

1. **Temporal Smoothing:** Without temporal smoothing, the letter predictions flicker between classes whenever the hand moves slightly or the preprocessing step produces small variations in the crop. To address this, the system keeps a window of the last several predictions (for example, the last 7 frames) and selects the most

common letter in that window as the current output.

2. **Letter Buffer and TTS:** The temporal-smoothed predictions are fed into a letter buffer. When the prediction changes and the model is confident (e.g., probability above 0.5), the new letter is appended to the buffer. Users can see the sequence of letters accumulate into a word or phrase.
3. **Domain Gap in Practice:** In the real-time demo, the model behaves best when:
 - The background is simple and not too cluttered
 - Lighting is even and not too dim or harsh
 - The hand is well centered and reasonably close to the camera

In more challenging conditions (busy backgrounds, low light, very small or off-center hand), recognition quality degrades. The observed accuracy under good conditions is roughly 80–90%, in line with expectations based on the domain difference between SignMNIST and live webcam data.

5.2. Limitations

The system has several important limitations:

- It does not handle J and Z, which require motion; it only recognizes static letters A–Y.
- It is sensitive to the assumption that the hand is the largest white region after thresholding, which may fail in cluttered scenes.
- It has difficulty with letters that are inherently ambiguous in static, low-resolution images (T and M in particular).
- It has not been evaluated extensively across different users, skin tones, hand sizes, ages, and camera setups.

5.3. Future Directions

Future work can address these issues by:

- Replacing contour-based hand detection with a robust hand tracking solution such as MediaPipe Hands, which would provide accurate hand keypoints and reduce failures in cluttered scenes.
- Applying more aggressive data augmentation during training (rotations, translations, brightness variations, elastic distortions) to improve robustness to real-world conditions.
- Exploring deeper architectures (e.g., ResNet18, MobileNetV3) or transfer learning from ImageNet to push test accuracy beyond 95% while maintaining real-time performance.
- Adding temporal models (LSTM, temporal CNN, or 3D CNN) that process sequences of frames to recognize motion-based letters J and Z and to incorporate dynamic cues into classification.
- Conducting user studies with members of the ASL community to evaluate usability, fairness, and real-world performance.

6. Conclusion

This project implemented and evaluated an end-to-end ASL letter recognition system that runs in real time on a CPU using a webcam. A custom CNN trained on the SignMNIST dataset achieved 94.70% test accuracy and near-perfect recognition for many letters. At the same time, per-class analysis highlighted a small set of inherently difficult letters, especially T and M, where static low-resolution images do not provide enough information for reliable classification. By combining the trained model with an 8-step preprocessing pipeline, temporal smoothing, a letter buffer, and text-to-speech, the system demonstrates how computer vision and deep learning can be integrated into a practical assistive technology prototype. It also illustrates the gap between controlled benchmark performance and real-world deployment and motivates continued work on robustness, motion modeling, and inclusive evaluation.

6.1. Contributions

Michael: Main programmer to the final submitted code for this project, accumulated data + graphs, aided in reporting and presentation. Richard: Formatted final report in LaTeX, aided in presentation and programming versions of the project to compare with other attempts by team. Jaya Singh: Contributed heavily to the report + presentation, created programs of which were used to contrast the final submitted code's accuracy.

7. References

- [1] Daniel Kuzmanovski and Kaggle. Sign Language MNIST (SignMNIST) dataset, 2018. Available: <https://www.kaggle.com/datasets/datamunge/sign-language-mnist>.
- [2] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In International Conference on Learning Representations (ICLR), 2015.
- [3] Nobuyuki Otsu. A threshold selection method from gray-level histograms. IEEE Transactions on Systems, Man, and Cybernetics, 9(1):62–66, 1979.
- [4] Alexey Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, et al. PyTorch: An imperative style, high-performance deep learning library. In Advances in Neural Information Processing Systems (NeurIPS), 2019.
- [5] Gary Bradski. The OpenCV library. Dr. Dobb's Journal of Software Tools, 2000.
- [6] Fan Zhang, Valentin Bazarevsky, Andrey Vakunov, Alexander Tkachenka, George Sung, Chuo-Ling Chang, and Matthias Grundmann. MediaPipe Hands: On-device real-time hand tracking. arXiv:2006.10214, 2020.