

Deep Learning - Project

Gesture Based Game Control

- Objective
- Dataset: Jester
- Preprocessing Steps for Gesture Recognition Model
- Model Training
- Overfitting Analysis
- Results
 - a) 3DCNN-4 layers (elu, SGD)
 - b) 3DCNN-2 ConvLSTM-2 layers (elu, Adam)
 - c) 3DCNN-3 layers + ConvLSTM 3 layers (elu, rmsprop)
 - d) 3DCNN-4 layers + ConvLSTM 1 layer (relu, Adam)
 - e) 3DCNN-4 layers + ConvLSTM 2 layers (elu, Adadelata)
- Final Notes
- Deployment and Working

0. Objective:

The primary objective of the Gesture Based Game Control project is to develop an effective gesture recognition system using the Jester dataset. This system aims to enable intuitive human-computer interaction by accurately recognizing hand gestures and translating them into game control commands. By leveraging deep learning techniques, the project seeks to enhance the gaming experience, allowing users to interact with games seamlessly through gestures.

1.Dataset : Jester

<https://www.qualcomm.com/developer/software/jester-dataset>

Overview

The Jester dataset is designed for training machine learning models to recognize human hand gestures, particularly in the context of human-computer interaction. It enables the development of responsive and accurate gesture recognition systems capable of distinguishing between subtle differences in gestures.

Content:

The dataset consists of **148,092 labeled video clips** of individuals performing a variety of hand gestures in front of a camera or webcam.

Classes:

There are **27 distinct gesture classes**

Dataset Split:

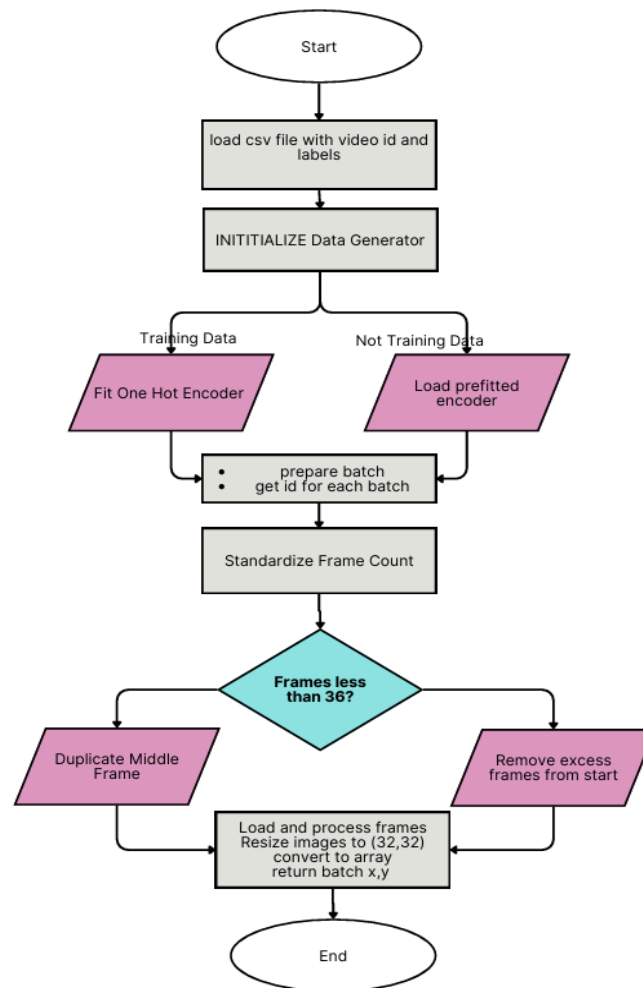
The dataset is divided into three parts to facilitate model training and evaluation:

- **Training Set:** 118,562 clips
- **Validation Set:** 14,787 clips
- **Test Set (without labels):** 14,743 clips

Quality and Format:

- **Image Quality:** Each video clip is represented as a series of JPG images, with a height of 100 pixels and a variable width depending on the content.
- **Frame Rate:** Videos were extracted at **12 frames per second**, and the number of JPGs per clip varies based on the length of the original video.
- **File Structure:** The dataset is provided in a TGZ archive format, split into parts with a maximum size of 1 GB. The total download size is approximately **22.8 GB**. Each directory within the archive corresponds to a single video, and filenames are sequentially numbered starting from **1.jpg**. The dataset was created with contributions from over **1,300 unique crowd actors**, ensuring a diverse representation of gestures.

2. Preprocessing Steps for Gesture Recognition Model



1. Load and Filter Annotations:

- Import necessary libraries.
- Define a list of gestures (`gesture_list`) that you want to keep for your training and validation datasets.
- Load the training annotations from the CSV file (`jester-v1-train.csv`) using Pandas.
- Filter the DataFrame to keep only the rows where the labels are in the `gesture_list`.
- Save the filtered DataFrame to a new CSV file (`new_jester_train.csv`).
- Repeat the same process for validation data by loading `jester-v1-validation.csv` and saving it as `new_jester_val.csv`.

2. Create a Data Generator for Batch Processing:

- Define a `DataGenerator` class that inherits from `tf.keras.utils.Sequence`. This class should:
 - Initialize the DataFrame, batch size, image dimensions, number of frames, and channel count.
 - Use one-hot encoding for the labels and store the encoder for later use.
 - Implement methods to:
 - Get the length of the dataset.
 - Generate batches of data, ensuring frames are standardized to a specific count per sample.

Initialization (`__init__` method)

- Initializes the generator with key parameters such as `batch_size`, `image_dim`, `frames_count`, and others.
- Loads the annotations from the specified CSV file and applies one-hot encoding to the labels using `OneHotEncoder`.
- Prepares the data for training or validation by loading the appropriate encoder based on whether the generator is for training or validation.

Length of Dataset (`__len__` method):

- Calculates and returns the total number of batches based on the size of the DataFrame and the batch size.

Batch Generation (`__getitem__` method):

- Retrieves the indexes for the current batch, collects the corresponding IDs, and calls the data generation method to create the batch data and labels.

Epoch Handling (`on_epoch_end` method):

- Shuffles the data indexes after each epoch to ensure that the model sees the data in a different order during each training pass.

Data Generation (`__data_generation` method):

- Initializes an empty array for the batch of images and labels.
- Iterates through the list of IDs, loading and preprocessing the corresponding frames (images).
- Resizes and converts each image to an array, ensuring that the images are in the correct format for the model.
- One-hot encodes the labels for the batch and returns both the batch of frames and the encoded labels.

Standardizing Frame Count (`standardize_frame_count` method):

Ensures that each sample contains a fixed number of frames (`frames_count`).

If a sample has fewer frames, duplicates frames from the middle to fill the gap. If there are too many frames, it removes excess frames.

Overview:

First, we loaded the CSV file containing image IDs and labels.

The DataGenerator is initialized with various parameters.

There's a check to determine if it's training data:

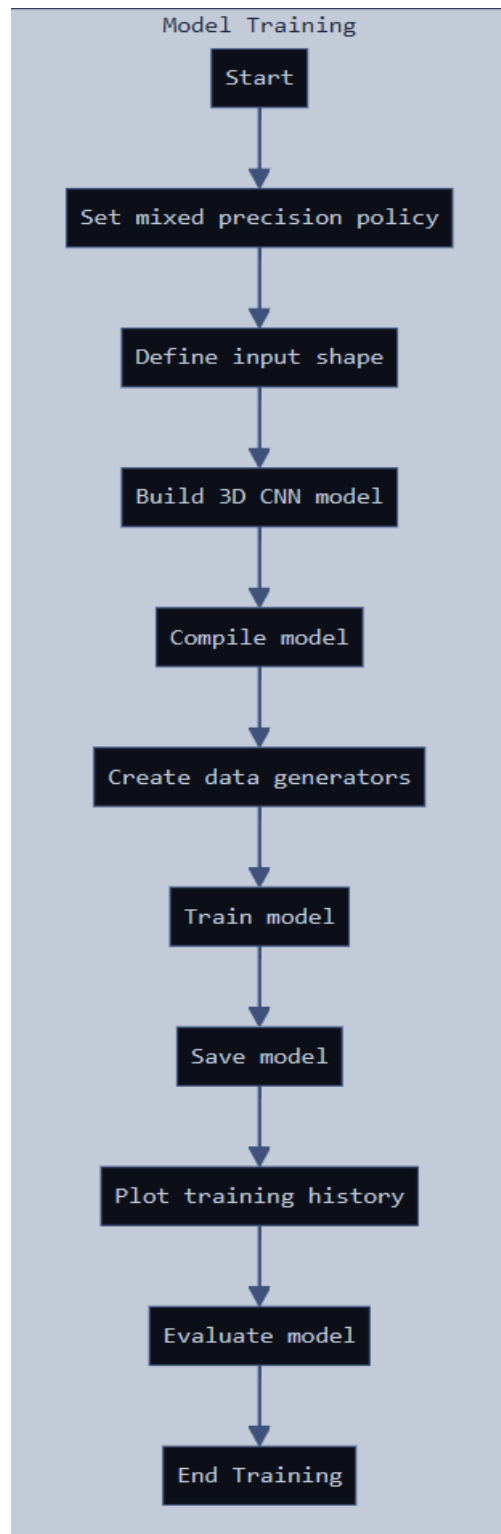
- If yes, a new OneHotEncoder is fitted on the labels and saved.
- If no, a pre-fitted encoder is loaded.

For each batch:

- File paths for each ID in the batch are retrieved.
- The frame count is standardized:
 - If there are fewer frames than required, the middle frame is duplicated.
 - If there are more frames than required, excess frames are removed from the start.
- Each frame is then loaded, resized, and converted to an array.

Finally, the labels are one-hot encoded, and the batch (X, y) is returned.

3. Model Training



Use of Mixed precision: It uses both 16-bit and 32-bit floating-point numbers in computations, improving speed and reducing memory usage.

Layers like convolutions, weights use 16-bit precision, while critical layers like batch normalization, loss values use 32-bit for stability. Although there's some reduction in precision, model accuracy remains stable.

This technique is effective with **NVIDIA GPUs** equipped with **Tensor Cores**. On CPUs it has very less benefit.

Refer:

<https://docs.nvidia.com/deeplearning/performance/mixed-precision-training/index.html>

Input shape : (36,32,32,3)

(36, 32, 32, 3) refers to:

- 36 frames (like a video sequence).
- Each frame is 32x32 pixels.
- 3 represents the RGB color channels.

(Tried using grayscale, i.e 1 channel, performance was low)

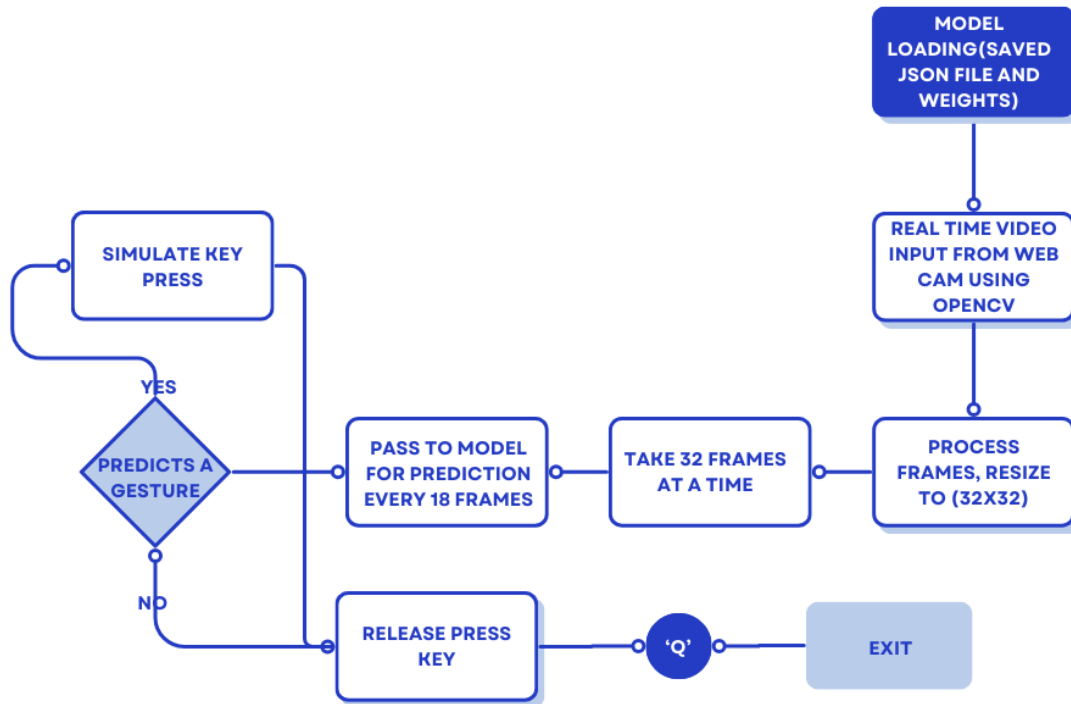
(Tried passing 16 frames together instead of 36, i.e , performance was low, nan error)

Create DataGenerators:

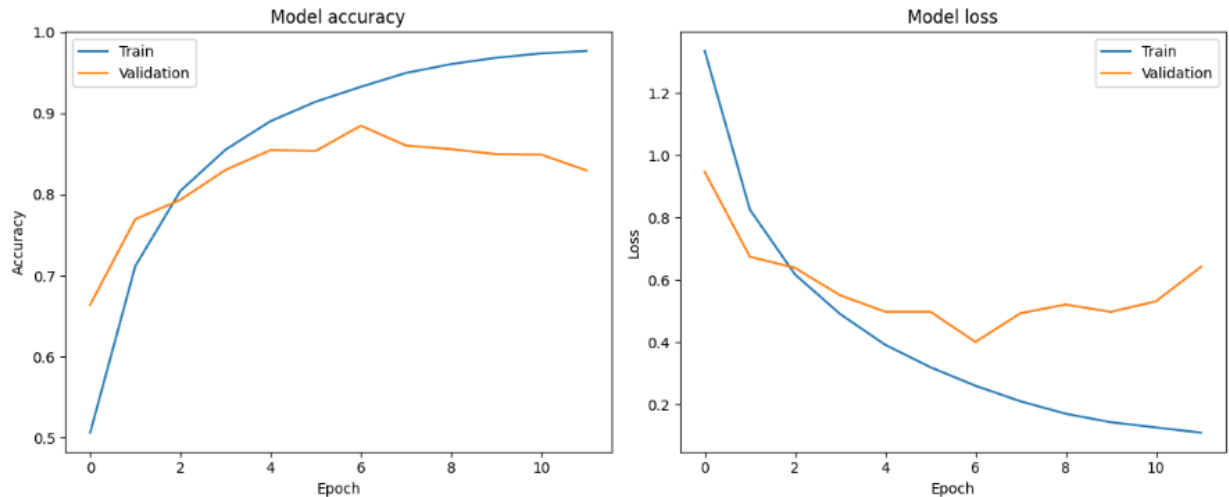
preprocessing(see above diagram)

Refer: <https://docs.nvidia.com/deeplearning/performance/mixed-precision-training/index.html>

3. real-time gesture recognition



4. Overfitting Analysis



We see that the training and validation curves in the graphs show signs of overfitting. While training accuracy increases steadily and reaches close to 1.0, validation accuracy plateaus and then declines slightly, indicating that the model is fitting too closely to the training data and struggling to generalize.

To mitigate overfitting, I have employed several techniques in the model:

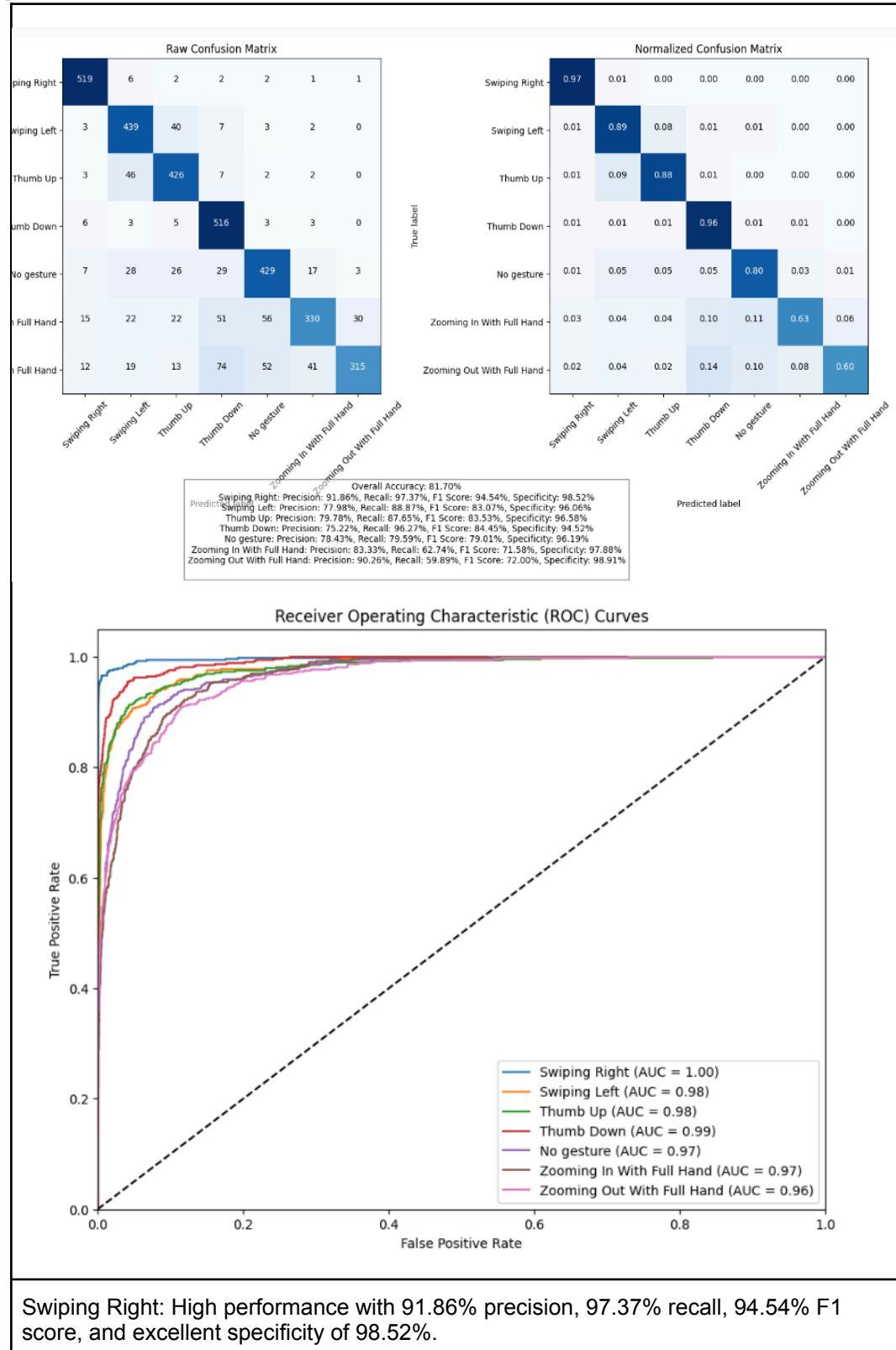
- **Dropout Layers:** I have experimented with dropout layers at different positions and by different amounts. Helps reduce reliance on specific neurons by randomly deactivating them during training.
- **Early Stopping:** Early stopping monitors validation loss and halts training once it no longer improves, ensuring the model doesn't continue to overfit after reaching optimal performance.
- **Regularization:** I have used L2 regularization in the dense layer, so that it penalizes overly complex weights, reducing the tendency for the model to fit noise in the training data.

Despite these strategies, there is still room for improvement in validation performance.

further tuning might help optimize generalization.

5. Results:

a)3DCNN-4 layers(elu, SGD)



Swiping Left: Lower precision (77.98%) and recall (88.87%) compared to Swiping Right, suggesting the model struggles more with left swipes.

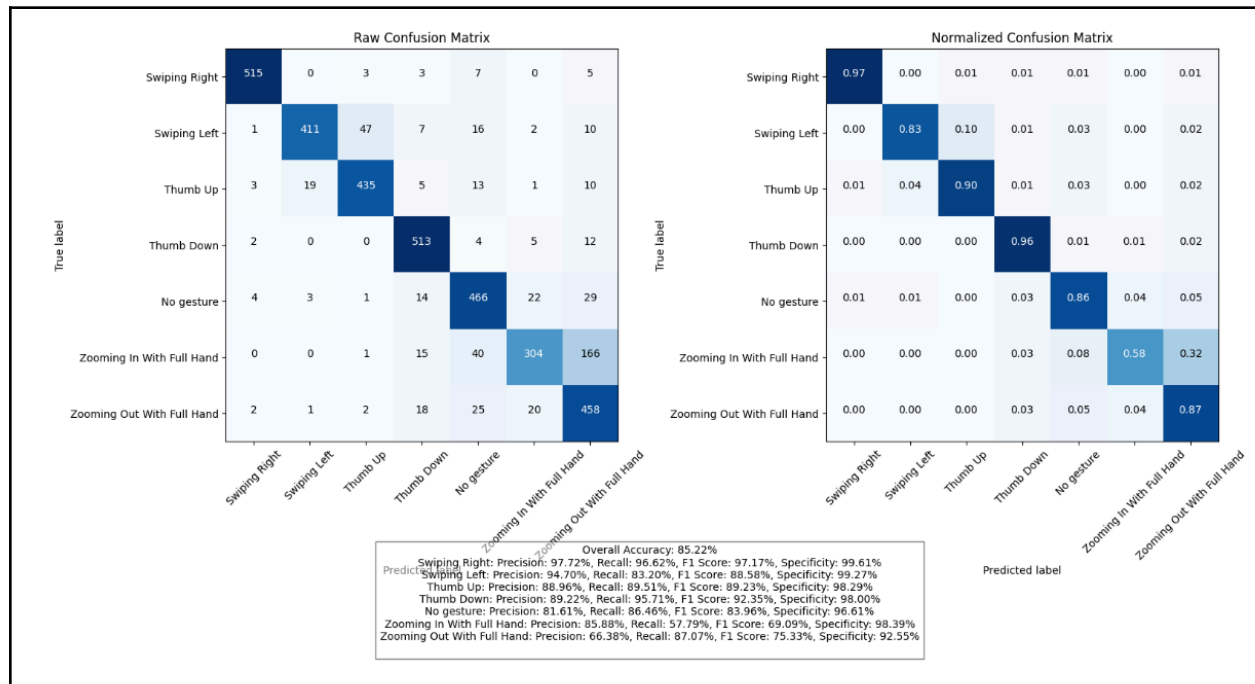
Thumb Up: Shows the largest misclassification with a lower precision (79.78%) and recall (82.65%).

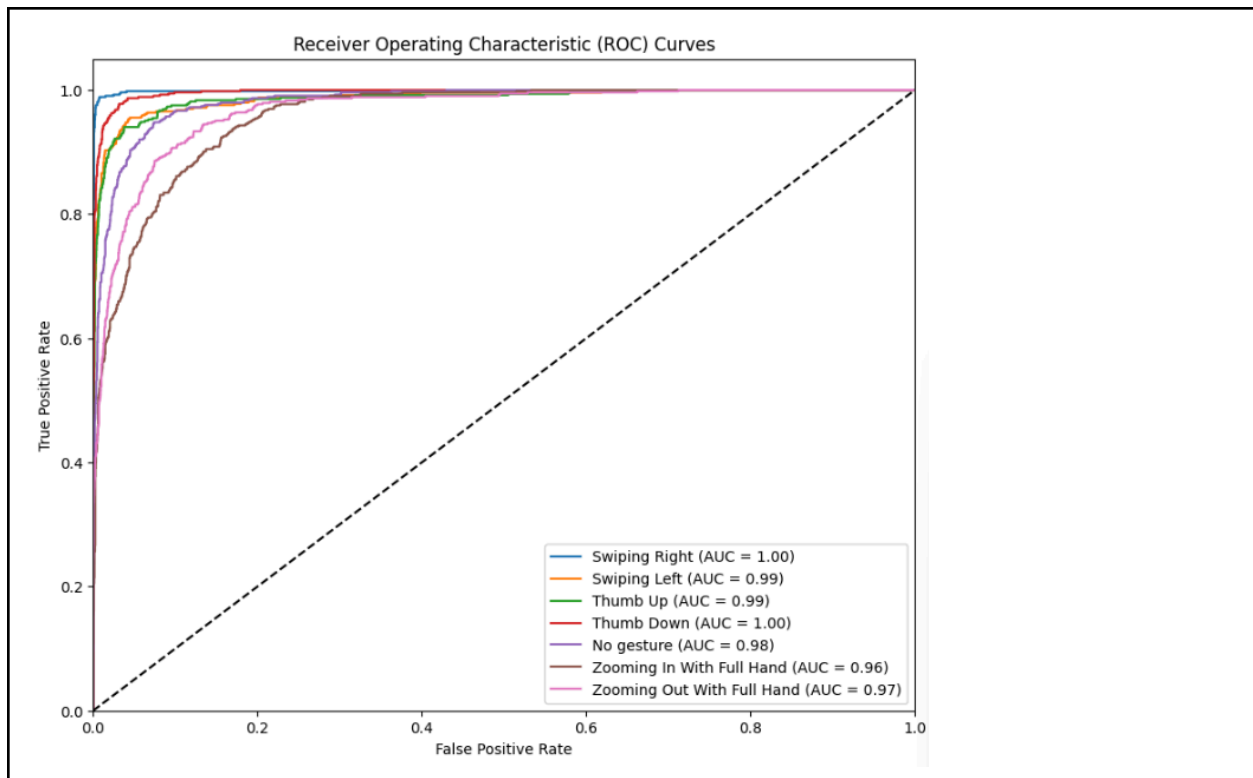
Zooming Gestures: These show the lowest performance, especially for Zooming Out, which has a 59.89% recall and 72.00% F1 score. There's significant confusion between Zooming In and Out, which might suggest the model struggles to differentiate them clearly.

The model performs well for gestures like **Swiping Right**, **Swiping Left**, and **Thumb Down**.

It struggles with differentiating between **Zooming In** and **Zooming Out**, and there is considerable confusion between **No gesture** and other gestures.

b)3DCNN-2 ConvLSTM-2 layers(elu, Adam)





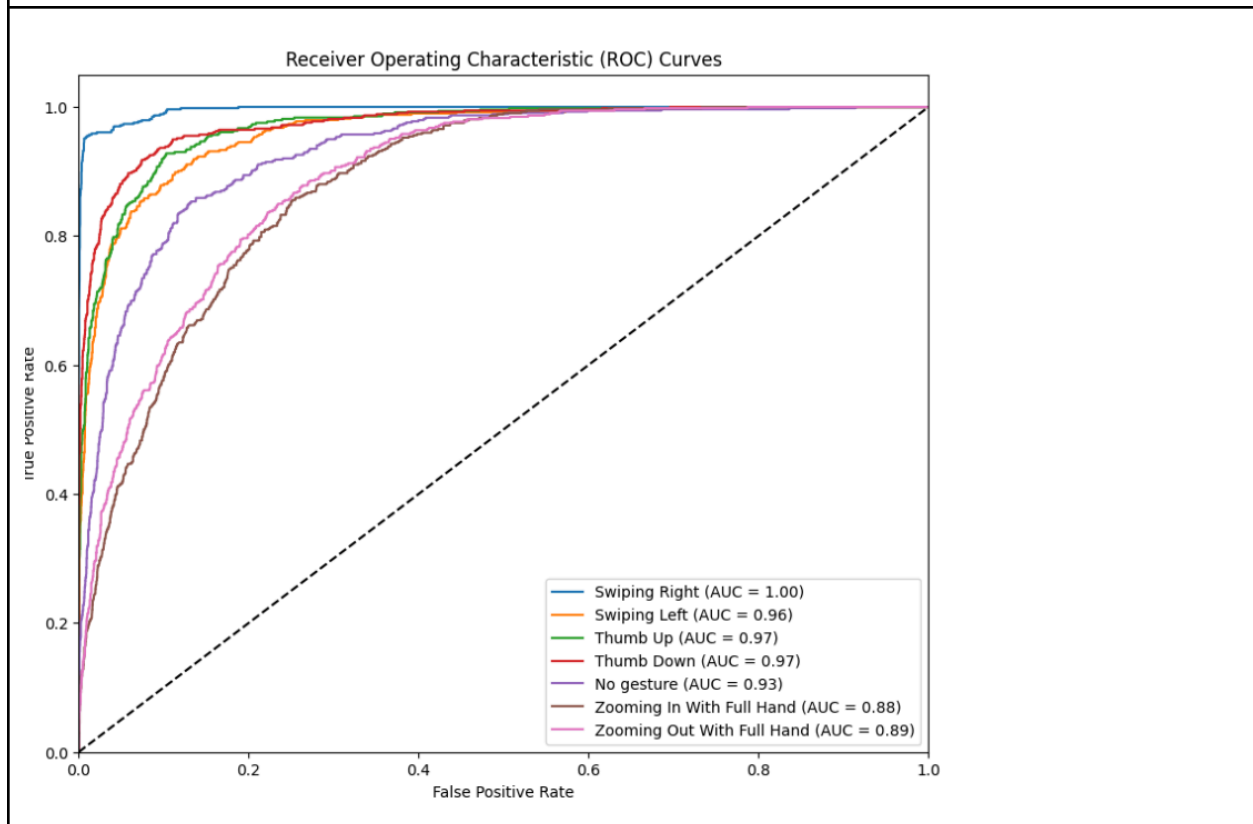
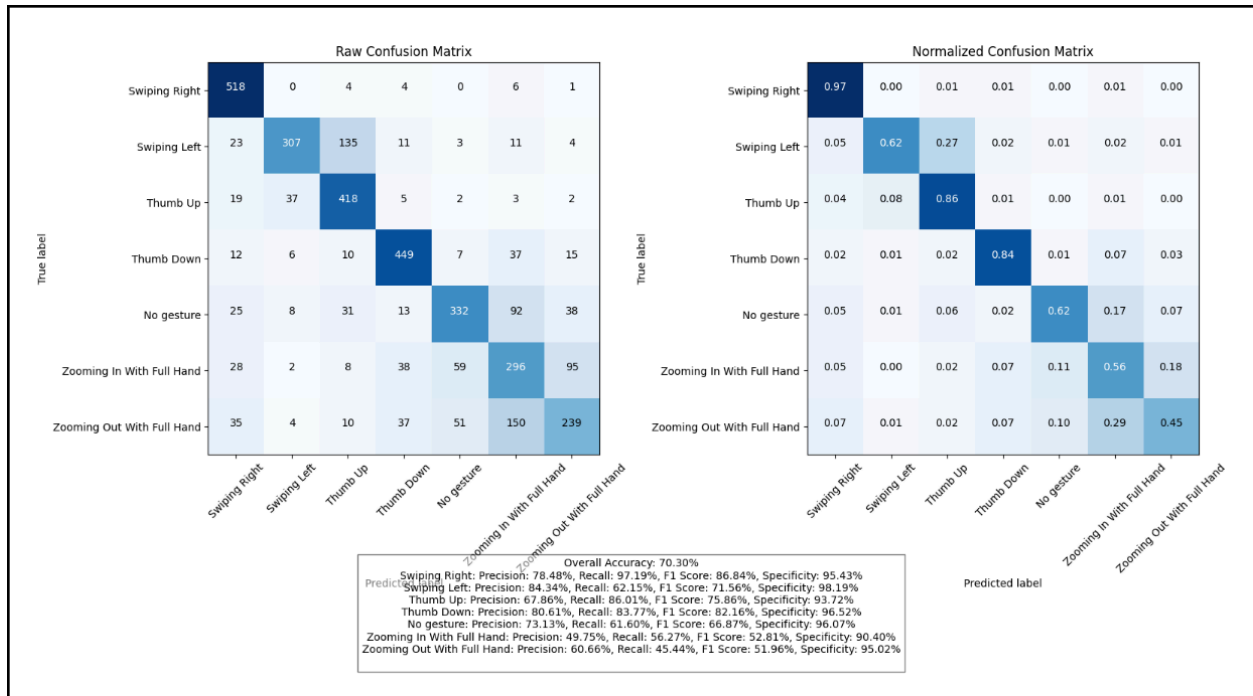
Better Class Separation: There is a clear improvement in performance across most classes compared to the previous model. The model shows better separation between similar gestures, especially for "Thumb Up," "Thumb Down," and "No Gesture."

Zooming Out With Full Hand still struggles a bit with lower precision, meaning the model often confuses it with other gestures, but its recall is quite good. Lesser accuracy for zoom in.

The changes in the model architecture, such as the increased number of filters, and the use of **ConvLSTM2D** layers, allow it to better capture temporal dependencies in the data. This likely contributes to the improvement in accuracy and performance for most gestures.

Also, the usage of the **Adam optimizer** with a small learning rate seems to be helping the model converge better during training.

c)3DCNN-3 layers + ConvLSTM 3 layers (elu, rmsprop)

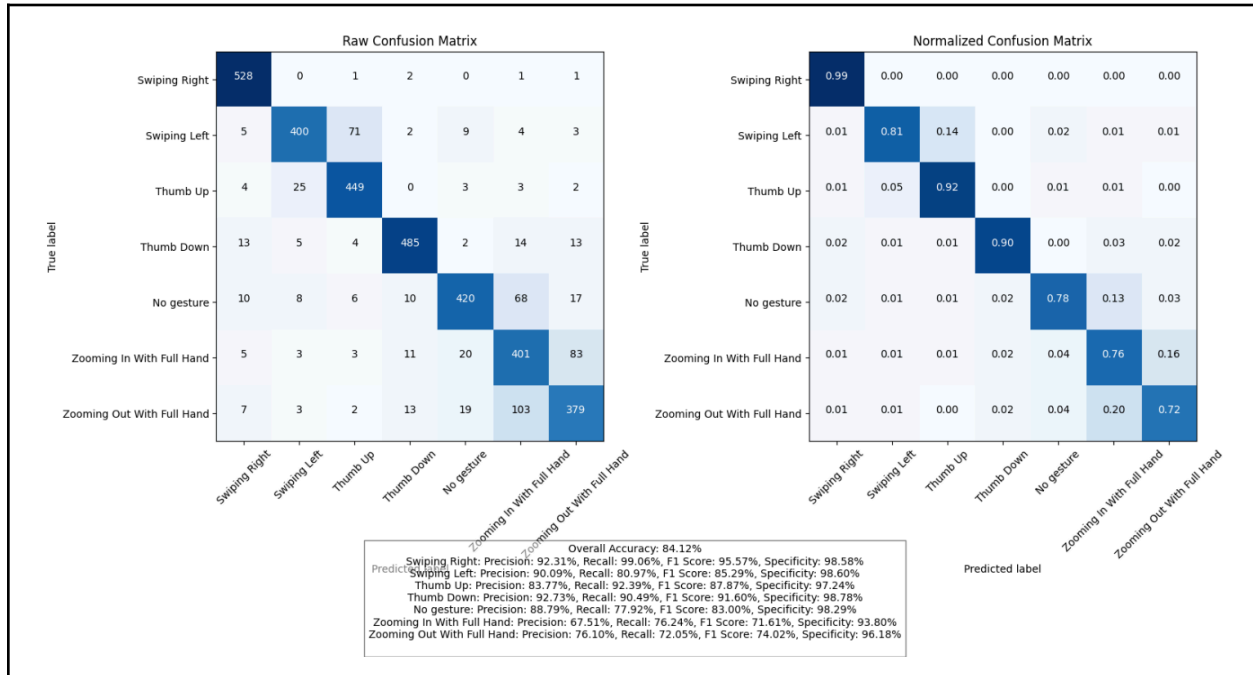


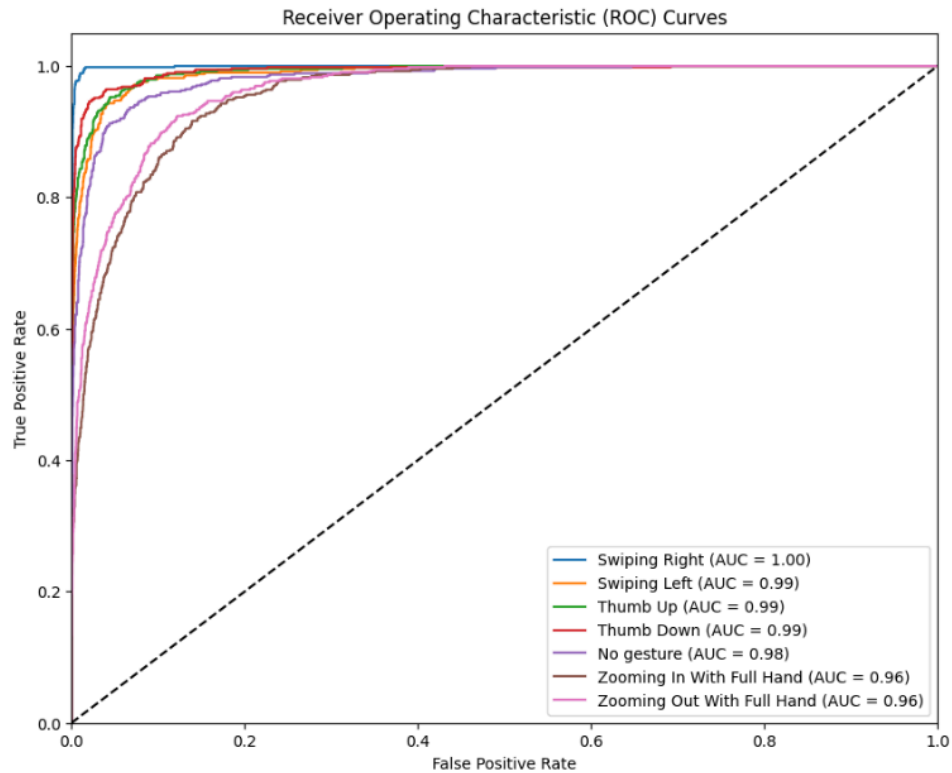
Performance: The model demonstrated not so great decent performance in recognizing various gestures. It did not show significant improvement compared to some of the other architectures.

Class Separation: There were still challenges in accurately distinguishing between certain gestures, particularly those with subtle differences.

No improvement in results

d)3DCNN-4 layers + ConvLSTM 1 layers (relu, Adam)

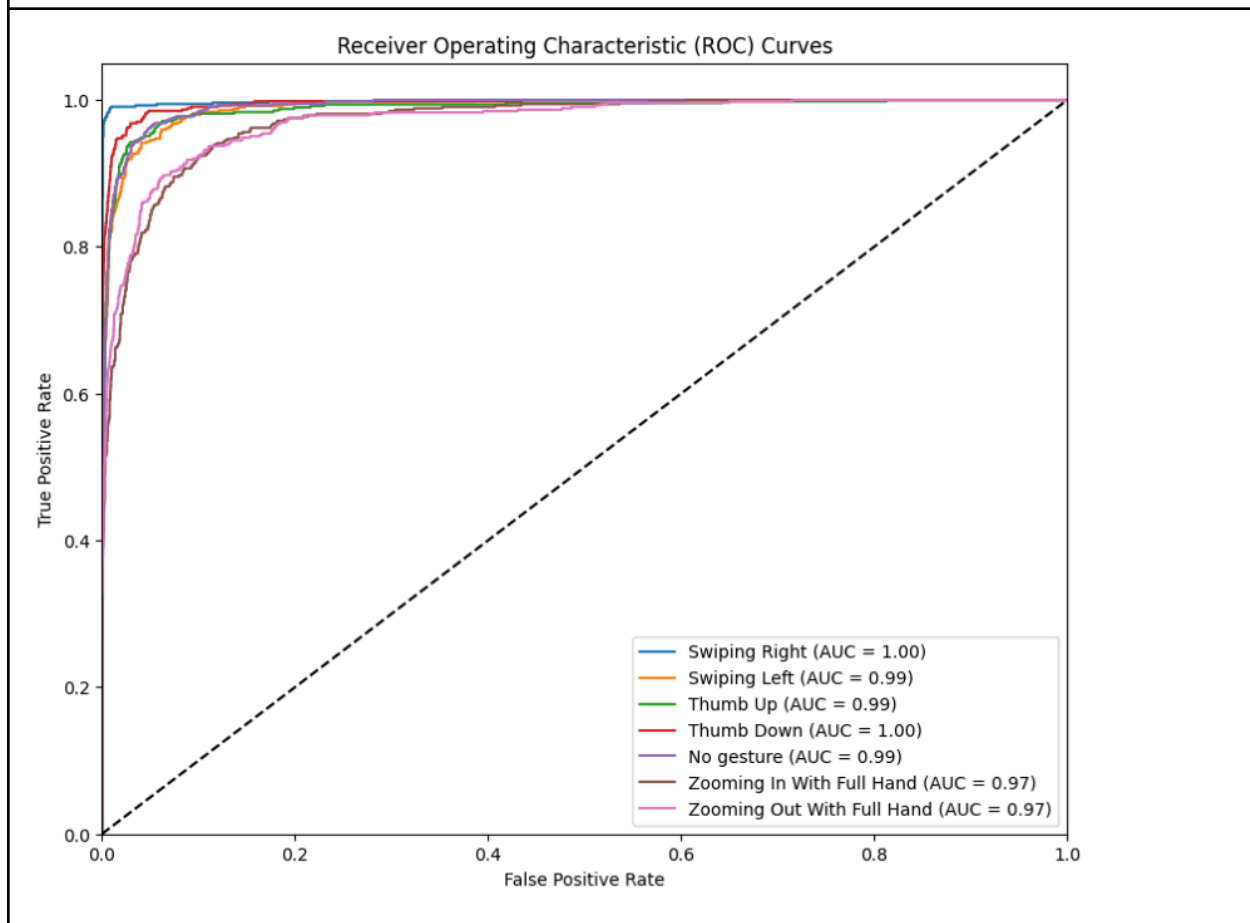
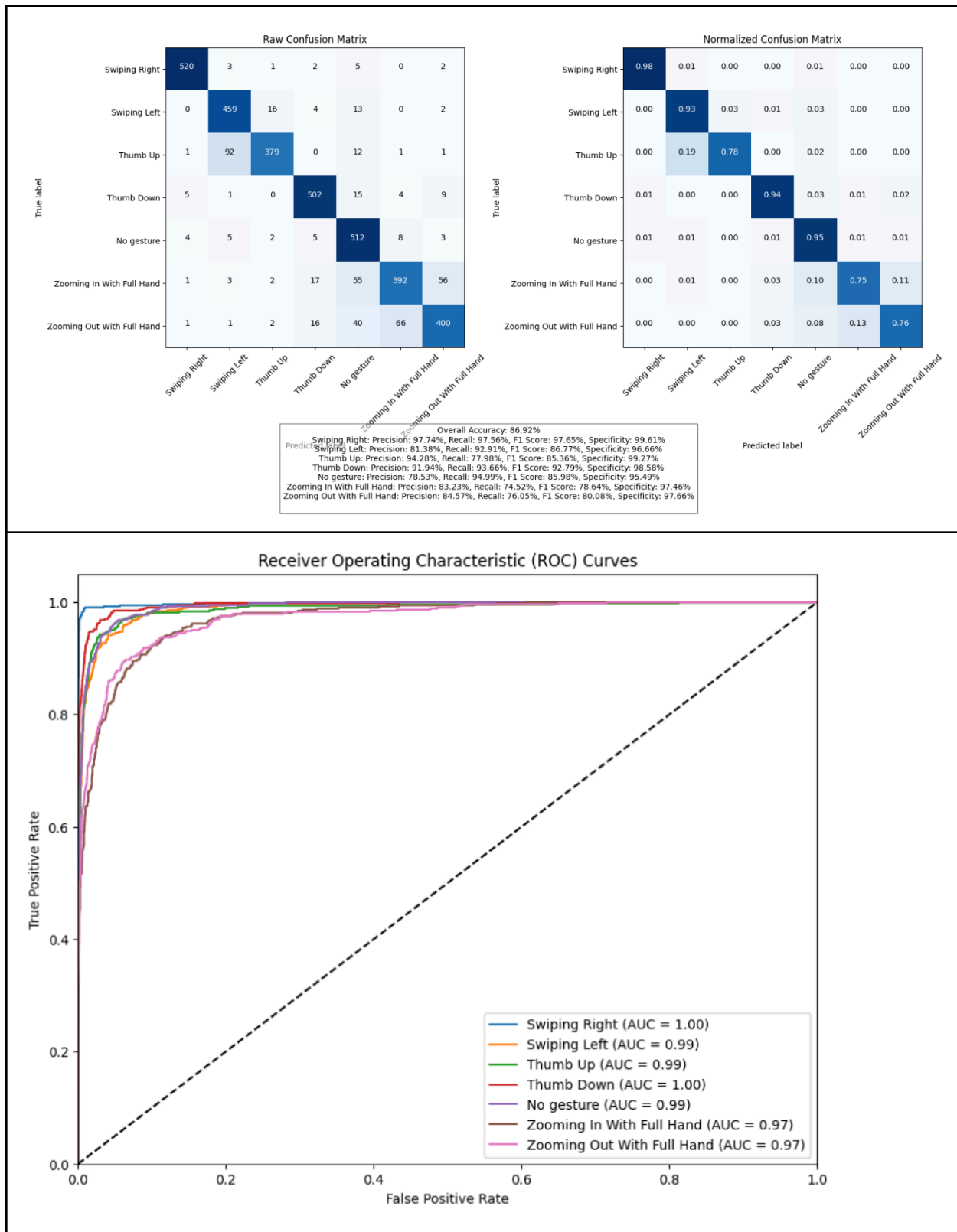




We see that this model has the better accuracy for zooming in and zooming out

Gesture Confusion: "No Gesture", "Zooming In With Full Hand" and "Zooming Out With Full Hand" show lower precision, meaning they are often misclassified with each other. This could be due to similarities in visual patterns

e)3DCNN-4 layers + ConvLSTM 2 layers (elu, Adadelta)



Better Class Separation: The previous model misclassified “No Gesture”, and **zooming in and out**, but this performs significantly better for “No Gesture”.

High Precision & Recall: The model performed exceptionally well on gestures like "Swiping Right" and "No Gesture," achieving high specificity and balanced precision-recall scores. This shows reliability in distinguishing these gestures accurately.

Challenges with Complex Gestures: Gestures like "Zooming In/Out" with a full hand had slightly lower F1 scores and recall, indicating potential difficulty in distinguishing fine-grained hand movements. But this model has improved accuracy for the same compared to all previous models.

Optimization: The switch to the Adadelata optimizer with a balanced learning rate contributed to better convergence and stability during training.

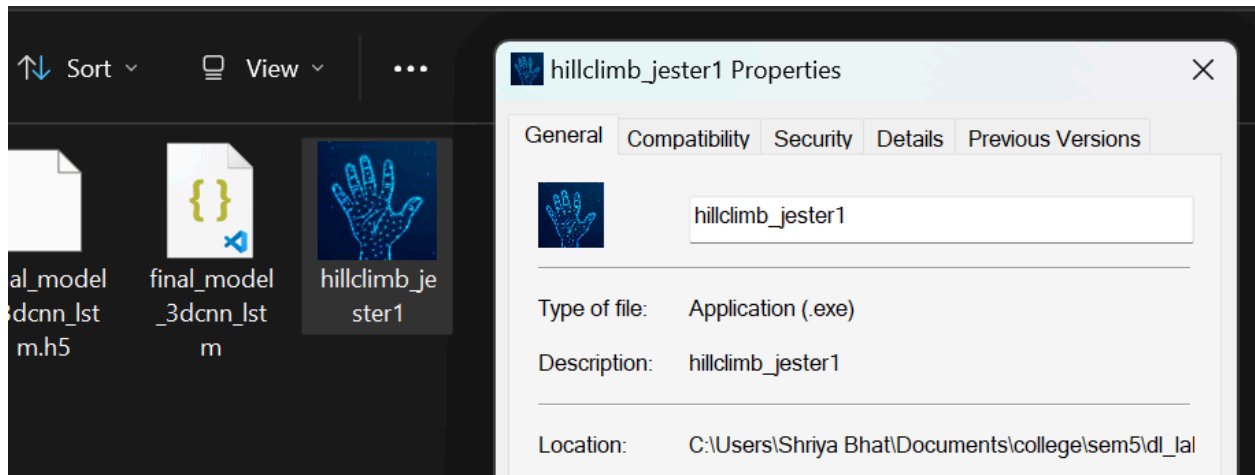
Model	Optimizer	Accuracy	Notable Observations
3DCNN-4 layers (elu, SGD)	SGD(0.001)	81.7 %	Struggles with "Zooming In/Out" distinction and "Thumb Up" misclassification
3DCNN-2 ConvLSTM-2 layers (elu)	Adam(0.0001)	85.21 %	Improved class separation, better performance in "Thumb Down" and "No Gesture"
3DCNN-3 layers + ConvLSTM 3 layers (elu)	RMSprop(0.0001)	83.95 %	No significant improvement
3DCNN-4 layers + ConvLSTM 1 layer (relu)	Adam(0.0001)	87.85 %	Improved results, least misclassifications(Zoom in, zoom out not good, misclassifies)
3DCNN-4 layers + ConvLSTM 2 layers (elu)	Adadelata(1.0)	90.55 %	Best accuracy Better Class Separation than previous models

Final Notes:

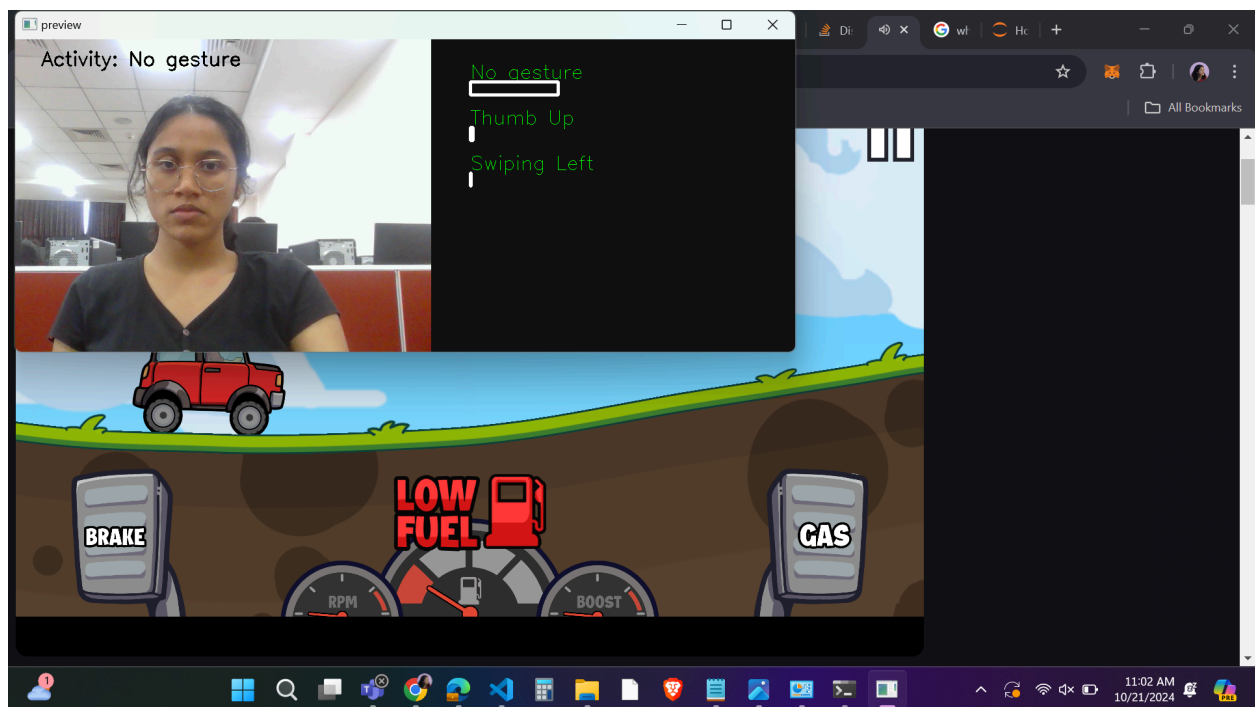
- **Model 1 (3DCNN-4 layers, SGD):** This model struggles with distinguishing between "Zooming In" and "Zooming Out," with low recall for "Zooming Out" gestures.
- **Model 2 (3DCNN-2 ConvLSTM-2 layers, Adam):** The incorporation of ConvLSTM layers helped capture temporal dependencies better, leading to higher accuracy and improved performance for many gesture classes.
- **Model 3 (3DCNN-3 layers + ConvLSTM 3 layers, RMSprop):** The model showed a decent performance, though slightly lower than some of the other models. The reduced dropout rate was used to maintain better generalization, and ConvLSTM layers helped capture temporal dependencies effectively.
- **Model 4 (3DCNN-4 layers + ConvLSTM 1 layer, Adam):** .This model performed much better, with a higher accuracy compared to the 3DCNN-3 layers + ConvLSTM 3 layers model. The higher dropout (0.5) helped with regularization, preventing overfitting, and the Adam optimizer contributed to better convergence
- **Model 5 (3DCNN-4 layers + ConvLSTM 2 layers, Adadelta):** The **Adadelta model (3DCNN-4 layers + ConvLSTM 2 layers)** demonstrates significant advancements in gesture recognition accuracy, achieving **90.55% accuracy**, the highest among all models evaluated. Key factor contributing to this success => **Improved Class Separation:** This model exhibited a notable ability to distinguish between the "No Gesture" class and other gestures, addressing the misclassification issues seen in previous models

5. Deployment and working:

- Run the Application (.exe) file in your local system without needing to install any dependencies.
- You can now do your gestures and simulate key press events.(In the background, predictions are printed in the terminal)

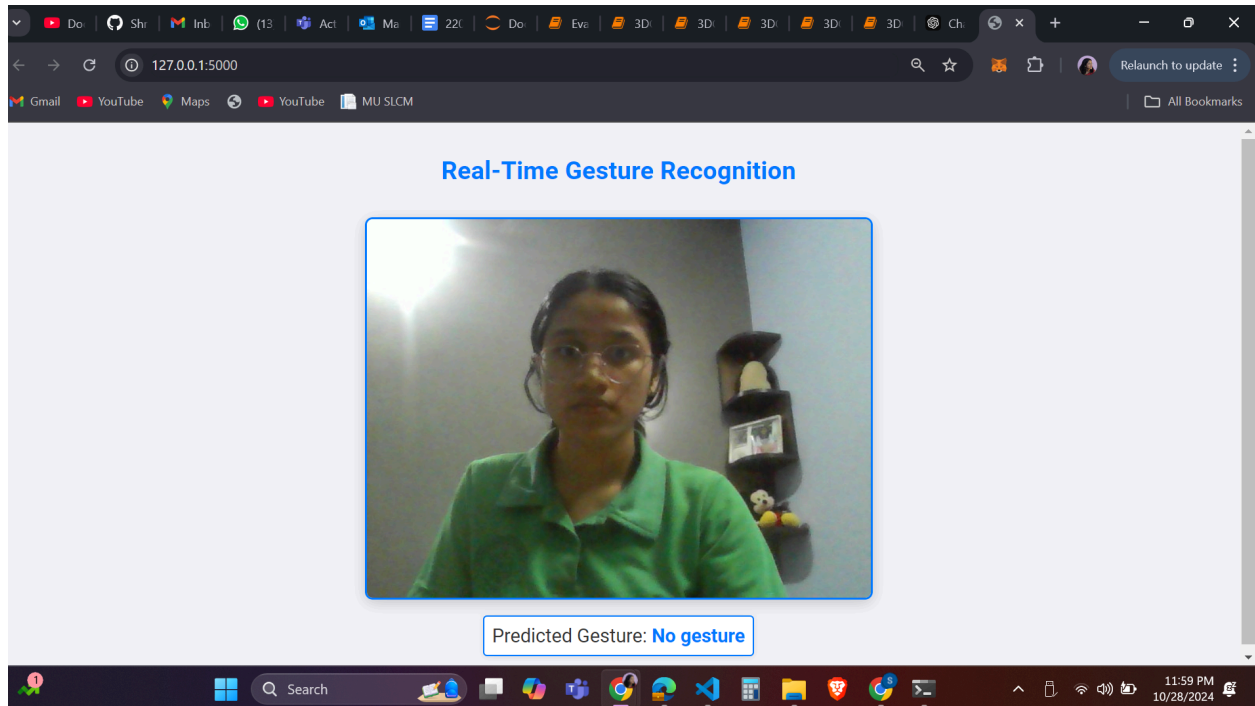


- The preview interface gives you the results' display in real time. It shows you the top 3 predictions. While this runs in the background, You can play your game.
- The gestures of 'swiping left' 'swiping right' 'zoom in' 'zoom out' are all mapped to left arrow key, right arrow key, and esc keys respectively.
- You can also watch a video, while easily forwarding 5 seconds and going back by 5 seconds.



Deployment using Flask:

(Frontend is just made to see the gesture predicted. It is supposed to run in the background, hence no frontend actually necessary.)



[Link to Github](#)

END