# Neural Networks Project – Gesture Recognition

## Team

Gurpreet Singh (Group Facilitator) **android.gurpreet@gmail.com**

Karan Prinja (Group Member) **karan.prinja@rakuten.com**

## Problem Statement

Imagine you are working as a data scientist at a home electronics company which manufactures state of the art smart televisions. You want to develop a cool feature in the smart-TV that can recognize five different gestures performed by the user which will help users control the TV without using a remote.

Each video is a sequence of 30 frames (or images).

The gestures are continuously monitored by the webcam mounted on the TV. Each gesture corresponds to a specific command:

- Thumbs up:  Increase the volume
- Thumbs down: Decrease the volume
- Left swipe: 'Jump' backwards 10 seconds
- Right swipe: 'Jump' forward 10 seconds
- Stop: Pause the movie

## Dataset

The training data consists of a few hundred videos categorized into one of the five classes. Each video (typically 2-3 seconds long) is divided into a sequence of 30 frames(images). These videos have been recorded by various people performing one of the five gestures in front of a webcam - similar to what the smart TV will use.

## The Models

Exploring the efficacy of various deep learning models for real-time gesture recognition in smart TV applications, we delve into the characteristics and performance of eight distinct models, culminating in the selection of the optimal model that balances accuracy, generalization, and computational efficiency.

| Model | Type | Batch Size | Image Size | Frames | Dense Neurons | Total Params | Trainable Params | Non-trainable Params | Categorical Accuracy | Validation Categorical Accuracy |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Conv3D | 24 | 100 x 100 | 30 | 64 | 592,773 | 592,037 | 736 | 96% | 80% |

The model's strength lies in its high training accuracy, suggesting effective learning, but the disparity in validation accuracy points towards a need for better generalization to unseen data.

| Model | Type | Batch Size | Image Size | Frames | Dense Neurons | Total Params | Trainable Params | Non-trainable Params | Categorical Accuracy | Validation Categorical Accuracy |
|---|---|---|---|---|---|---|---|---|---|---|
| 2 | Conv3D | 20 | 120 x 120 | 20 | 64 | 1,113,925 | 1,112,933 | 992 | 99% | 76% |

The high training accuracy indicates effective learning capabilities, yet the considerable drop in validation accuracy suggests issues with overfitting

| Model | Type | Batch Size | Image Size | Frames | Dense Neurons | Total Params | Trainable Params | Non-trainable Params | Categorical Accuracy | Validation Categorical Accuracy |
|---|---|---|---|---|---|---|---|---|---|---|
| 3 | Conv3D | 20 | 100 x 100 | 30 | 64 | 900,933 | 899,941 | 992 | 98% | 90% |

The model exhibits strong learning with its high training accuracy, and the closer alignment between training and validation accuracies indicates better generalization capabilities compared to the previous models. This balance suggests a more effective model in handling unseen data.

| Model | Type | Batch Size | Image Size | Frames | Dense Neurons | Total Params | Trainable Params | Non-trainable Params | Categorical Accuracy | Validation Categorical Accuracy |
|---|---|---|---|---|---|---|---|---|---|---|
| 4 | Conv3D | 20 | 150 x 150 | 30 | 128 | 1,638,213 | 1,637,221 | 992 | 99% | 74% |

Despite its excellent training accuracy, the significant gap between training and validation accuracies suggests a strong tendency towards overfitting.

| Model | Type | Batch Size | Image Size | Frames | Dense Neurons | Total Params | Trainable Params | Non-trainable Params | Categorical Accuracy | Validation Categorical Accuracy |
|---|---|---|---|---|---|---|---|---|---|---|
| 5 | Conv3D | 22 | 110 x 110 | 30 | 136 | 9,40,029 | 9,39,005 | 1024 | 99% | 75% |

The model's high training accuracy is indicative of effective learning from the training data. However, the large discrepancy between training and validation accuracies suggests potential overfitting

| Model | Type | Batch Size | Image Size | Frames | Dense Neurons | Total Params | Trainable Params | Non-trainable Params | Categorical Accuracy | Validation Categorical Accuracy |
|---|---|---|---|---|---|---|---|---|---|---|
| 6 | CNN-LSTM | 20 | 120 x 120 | 30 | 64 | 7488069 | 7487621 | 448 | 51% | 49% |

The model's balanced training and validation accuracies indicate consistency in performance on both training and unseen data. However, the relatively low accuracy levels suggest that the model might benefit from further tuning, possibly including adjustments in architecture, hyperparameters, or training data augmentation to enhance learning efficiency.

| Model | Type | Batch Size | Image Size | Frames | Dense Neurons | Total Params | Trainable Params | Non-trainable Params | Categorical Accuracy | Validation Categorical Accuracy |
|---|---|---|---|---|---|---|---|---|---|---|
| 7 | Transfer Learning (ResNet50) | 20 | 120 x 120 | 30 | 64 | 24,129,861 | 538,053 | 23,591,808 | 62% | 50% |

The use of transfer learning with ResNet50 provides a complex model structure, evident from the high number of non-trainable parameters. While the training accuracy is moderate, the drop in validation accuracy indicates challenges in generalizing to new data.

| Model | Type | Batch Size | Image Size | Frames | Dense Neurons | Total Params | Trainable Params | Non-trainable Params | Categorical Accuracy | Validation Categorical Accuracy |
|---|---|---|---|---|---|---|---|---|---|---|
| 8 | CNN-LSTM | 20 | 100 x 100 | 30 | 64 | 8,87,141 | 8,86,149 | 992 | 99% | 81% |

The model demonstrates excellent learning capabilities as reflected in the high training accuracy. However, the noticeable gap between training and validation accuracies suggests a potential overfitting issue.

For the specific application of gesture recognition for a smart TV using a webcam, **Model 3** (Conv3D) would be the most suitable choice. Here is why:

**Balance Between Accuracy and Generalization:** Model 3 strikes an excellent balance with 98% training accuracy and 90% validation accuracy. This indicates not only a high degree of learning from the training dataset but also a strong ability to generalize well to new, unseen data, which is crucial for real-time gesture recognition.

**Appropriate Complexity for Real-Time Processing:** While Model 3 is complex enough to capture the nuances in gesture recognition (which can be quite subtle and varied), it doesn't have an excessively high number of parameters. This balance is essential for ensuring that the model can process data in real-time, a key requirement for a smooth user experience in a smart TV context.

**Suitability for Video Data:** Given that gesture recognition involves analyzing sequences of images (video data), the Conv3D architecture of Model 3 is well-suited for this task. It can effectively capture the spatial and temporal dynamics of gestures, which is essential for accurate recognition.

**Potential for Real-World Robustness:** The higher validation accuracy suggests that Model 3 will be more robust in a real-world environment, dealing effectively with the variations and unpredictability inherent in how different users might perform gestures.

In conclusion, for gesture recognition in a smart TV setup, Model 3's combination of high accuracy, good generalization, appropriate complexity, and suitability for video data makes it the optimal choice.