# Models for FER

- Gaurav Sharma

---

All the below models are trained and tested by myself and the results are trustable.

Datasets Under Consideration -

1. FER-13 Cleaned
   Angry - 966 samples
   Sadness - 1326 samples
   Fear - 859 samples
   Neutral - 1446 samples
   Happy - 2477 samples

   This dataset has a total of **7074 images**. All the images are grayscale and have dimension (48x48x1). This dataset contains images taken in **un-controlled settings** and has a lot of variations.

2. FER-13 Aligned
   Angry - 803 samples
   Sadness - 1069 samples
   Fear - 653 samples
   Neutral - 1230 samples
   Happy - 1996 samples

   I applied face alignment to every image of the above dataset and removed all images which fail during this process. The dataset is reduced to a total of **5751 images**.
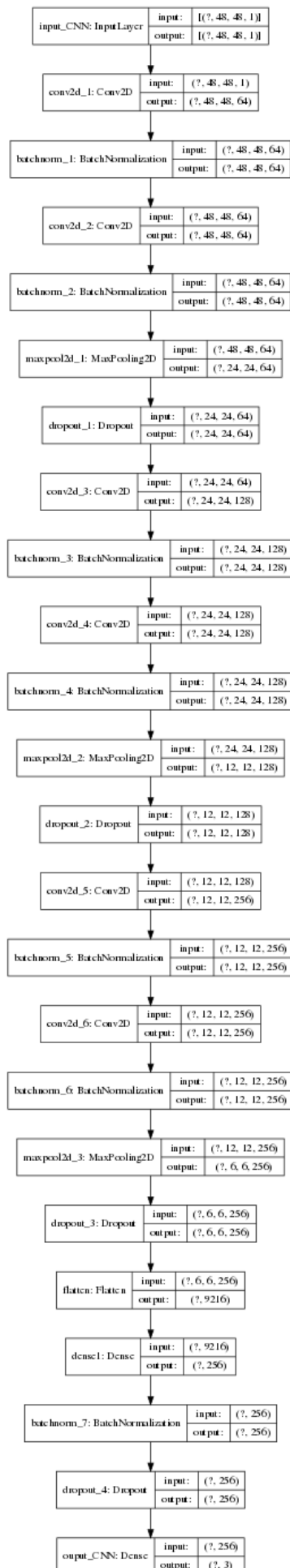
3. CK
   Fear - 75 samples
   Surprise - 249 samples
   Happy - 207 samples
   Sadness - 84 samples
   Angry - 135 samples

   This dataset has a total of **750 images**. All the images are grayscale and have dimension (48x48x1). This dataset is constructed in **controlled settings** and has very clear frontal images.

4. FER-13 Aligned + CK
   I simply combined the **2nd** and **3rd** datasets, it not only results in more data but also 1 more emotion. But the class **surprise** is very low as compared to others (we noticed this data imbalance during testing our model).

# 1. Convolutional Neural Network (CNN)



This is a naive CNN, nothing fancy. It has three blocks stacked sequentially, where each block consists of **conv -> batchnorm -> conv -> batchnorm -> maxpool ->dropout** and at the end, there is a single dense layer followed by output layer.

## Model evaluation on the following datasets:-

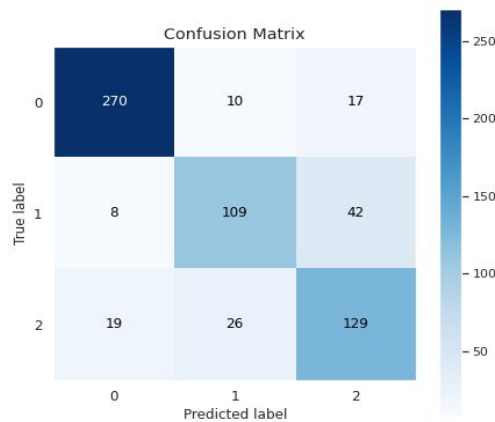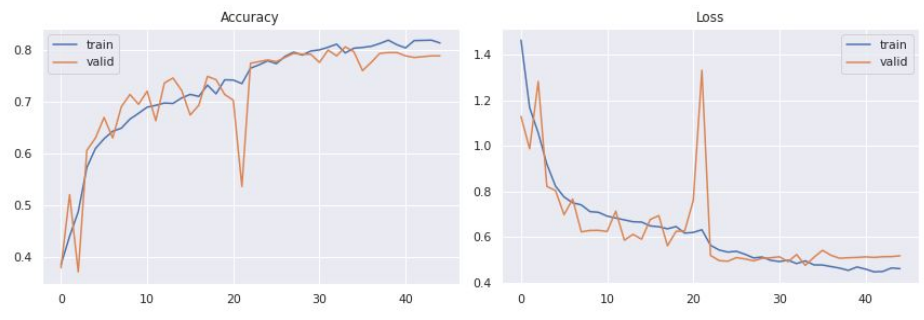**Dataset:** FER-13 Cleaned

**Emotions:** Happy, Sadness, Neutral
**Input:** Single image
**Output:** Emotion in that image

**Performance**: train-test split of 0.88
Best Train Accuracy ~ 82%
Best Test Accuracy ~ 80.09%

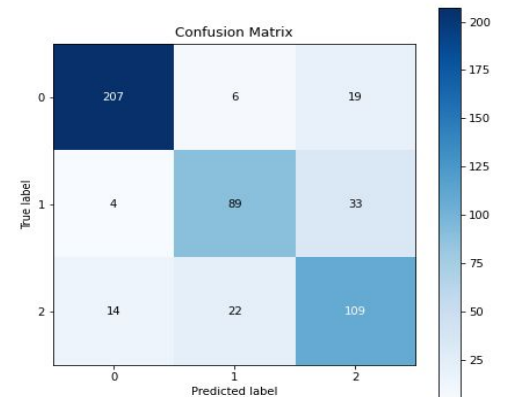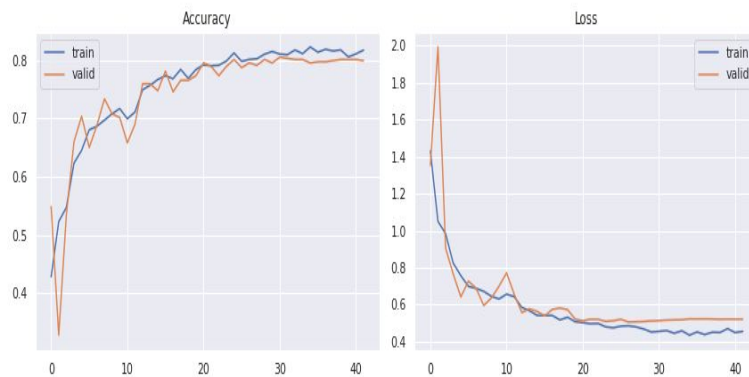Dataset: FER-13 Aligned

**Emotions:** Happy, Sadness, Neutral
**Input:** Single image
**Output:** Emotion in that image

**Performance:** train-test split of 0.88
Best Train Accuracy ~ 82%
Best Test Accuracy ~ 80%



Dataset: CK

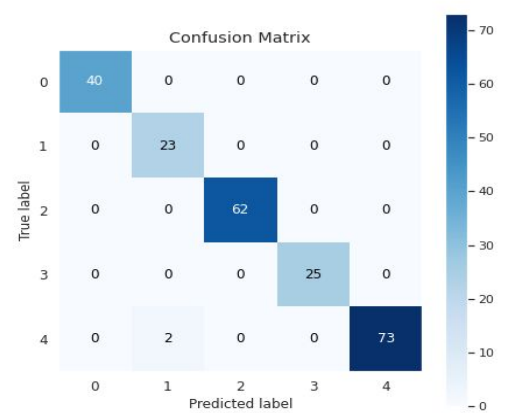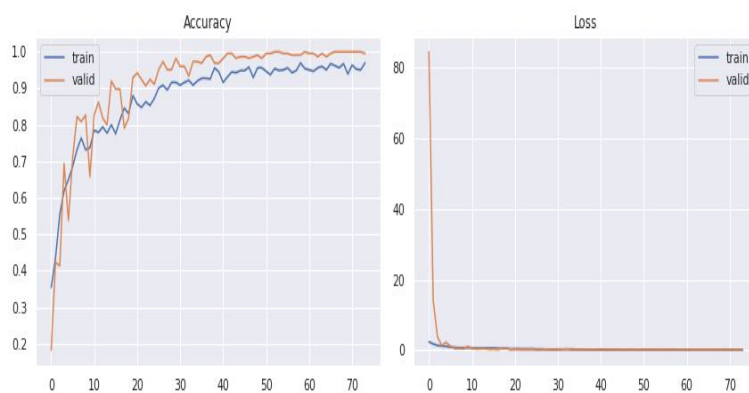**Emotions:** Happy, Sadness, Surprise, Fear, Angry
**Input:** Single image
**Output:** Emotion in that image

**Performance:** train-test split of 0.75
Best Train Accuracy ~ 100%
Best Test Accuracy ~ 96%



Dataset: FER-13 Aligned + CK
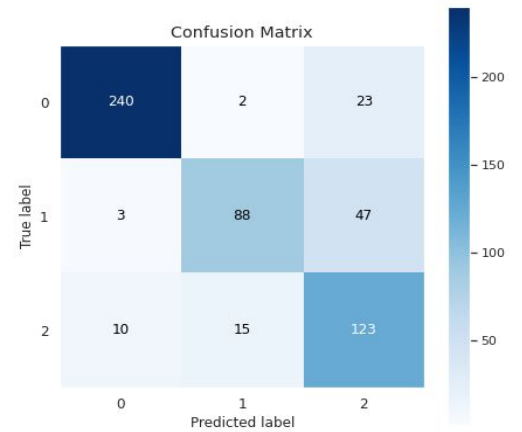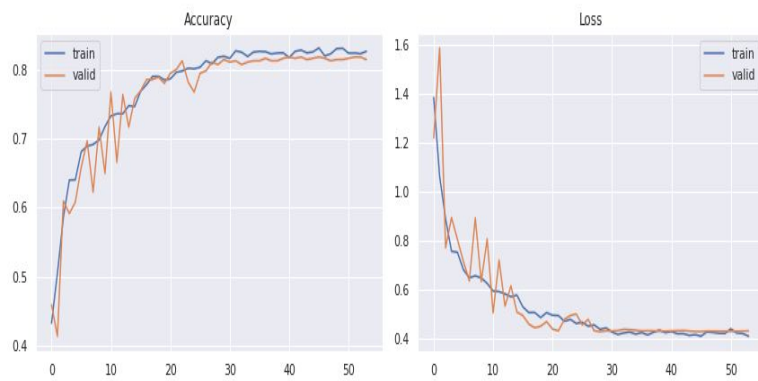**Input:** Single image
**Output:** Emotion in that image
**Emotions:** Happy, Sadness, Neutral

Performance: train-test split of 0.86
        Best Train Accuracy ~82%
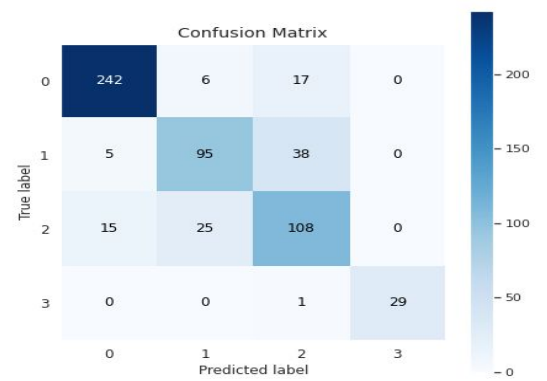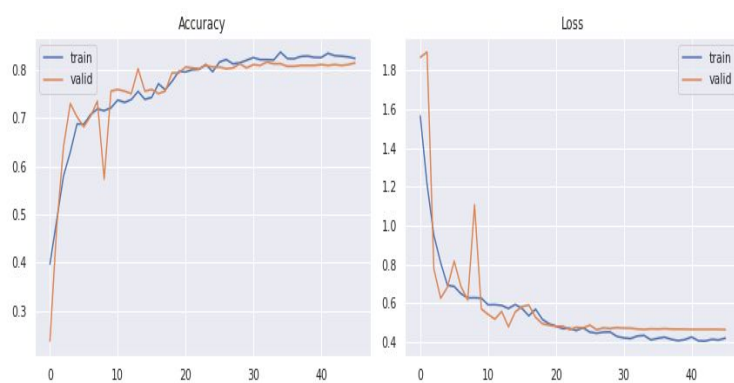        Best Test Accuracy ~ 81%



Emotions: Happy, Sadness, Neutral, Surprise
Performance: train-test split of 0.86
        Best Train Accuracy ~ 82%
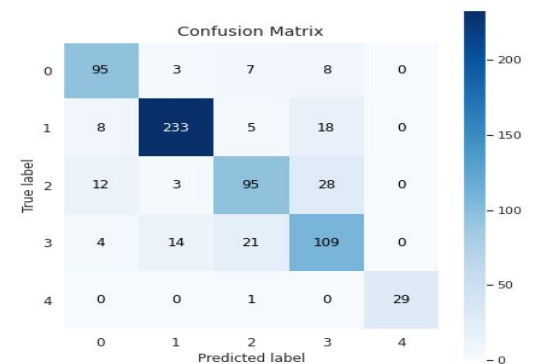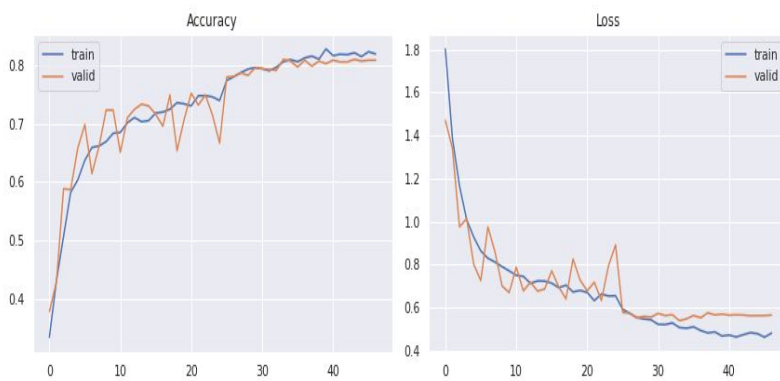        Best Test Accuracy ~ 81%



Emotions: Happy, Sadness, Neutral, Surprise, Angry
Performance: train-test split of 0.86
        Best Train Accuracy ~ 82%
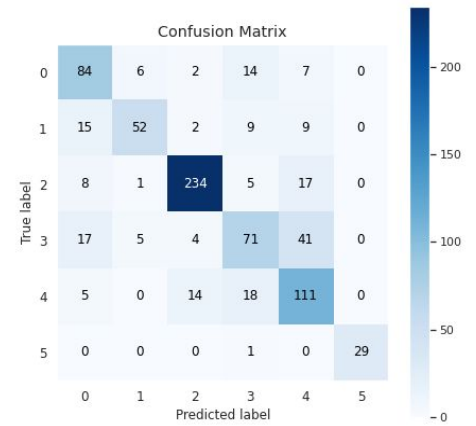        Best Test Accuracy ~ 80%
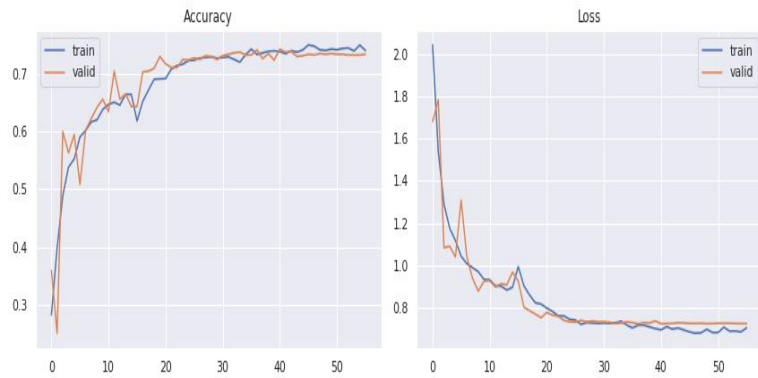
**Emotions:** Happy, Sadness, Neutral, Surprise, Angry, Fear
**Performance:** train-test split of 0.86
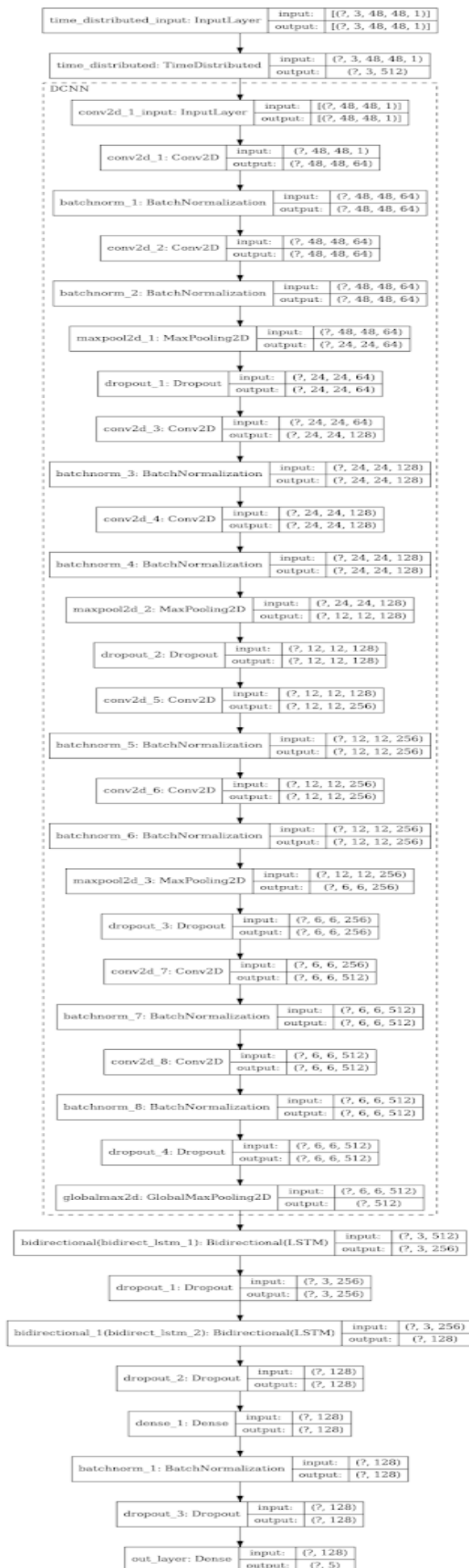Best Train Accuracy ~ 76%
Best Test Accuracy ~ 75%



Pros:
- A very simple model (nothing fancy)
- Fast to train and test
- Easy to debug

Cons:
- Model is simple as compared to complex task in hand
- Sometimes focuses on unwanted things as well like nose

# 2. TimeDistributed Convolutional Bidirectional LSTM



This is a **memory-based model**. In this model, all frames of a video are first passed to CNN, after that the output is passed to bidirectional LSTM and then to dense layers. Because this model requires sequential (video-based data) so I only applied this on CK. The CK dataset we have in our hand is the cropped one, in which for each video we have the images of the last three frames. So I stacked those 3 frames in time dimension resulting in a **3 frame video** (reverse engineering).

Dataset: Ck
  Emotions: Happy, Surprise, Angry, Sadness, Fear
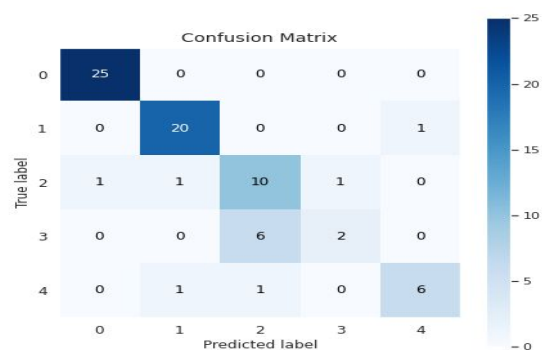  Input: Sequence of images(frames of video)
  Output: Single emotion in that entire image sequence
  Performance: train-test split of 0.7
              Best Train Accuracy ~ 87%
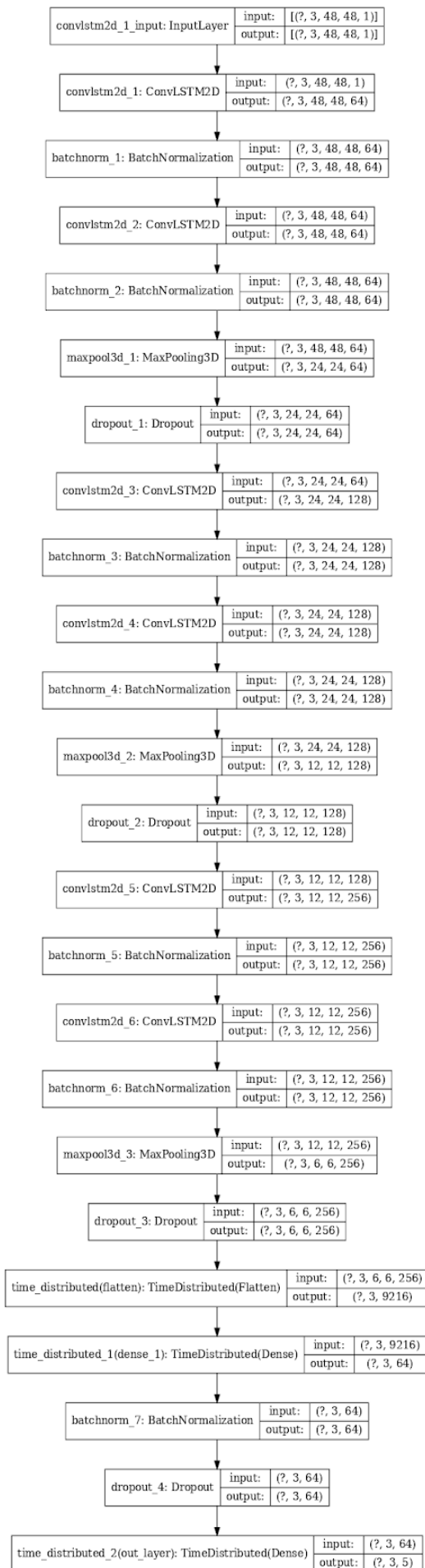              Best Test Accuracy ~ 85%





Pros:
  ● Learn the context from past frames.
Cons:
  ● Takes more time to train on large data

# 3. ConvLSTM2D



This is a **memory-based model**. In this model, all frames of a video are first passed to convLSTM2D and then to time distributed dense layers.
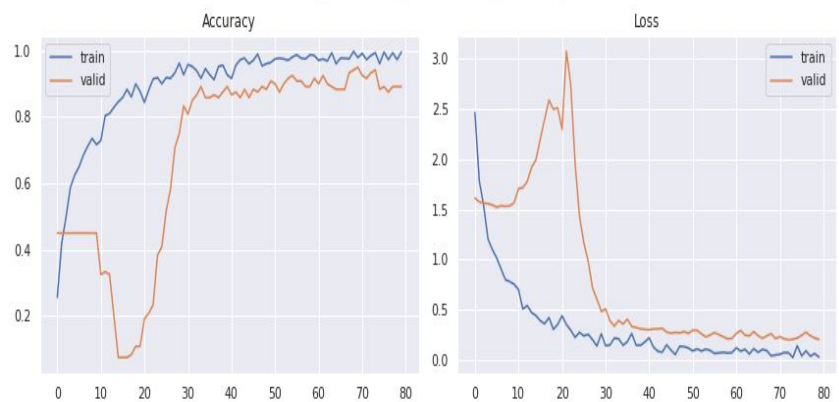
Dataset Used: CK
    **Emotions:** Happy, Surprise, Angry, Sadness, Fear
    **Input:** Sequence of images(frames of video)
    **Output:** One emotion for each frame in that sequence

**Performance:** As the dataset is very small and the results are not very great for now so I just showed the end results because until it is tested on good quantity data these results may miss-lead.



Pros:
- Learn the context from past frames.

Cons:
- Takes more time to train on large data

# 4. CNN Wrapped with NSL

This model is the same as the previous CNN but there is an extra regularization added to the model.

Dataset: FER-13 Cleaned
  Emotions: Happy, Sadness, Neutral
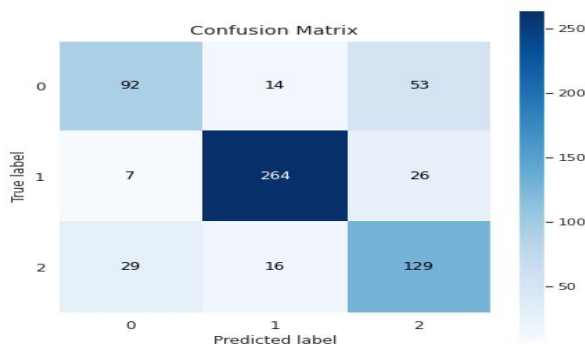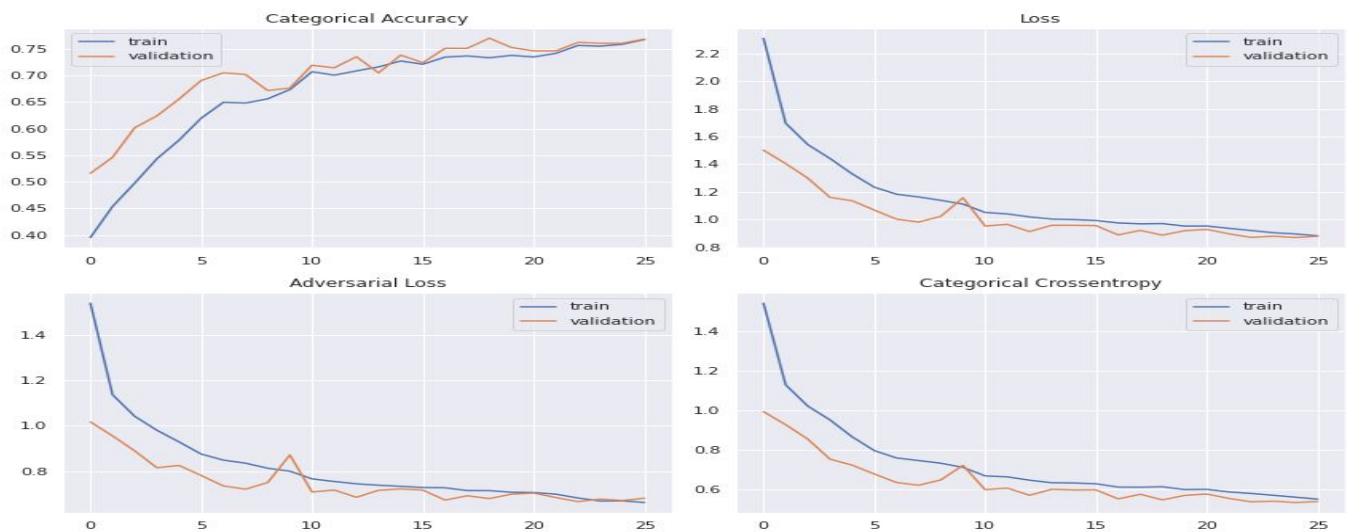  Input: Single image
  Output: Emotion in that image
  Performance: train-test split of 0.7
    Best Train Accuracy ~ 76%
    Best Test Accuracy ~ 76%





I was expecting the performance to be better than the base CNN model but the accuracy achieved is slightly low. Maybe we can tweak the parameter more and see the results. I tried 9 different configurations for the two NSL parameters and got almost the same result for all configurations. You can find my results here. As this NSL thing is recently introduced so there are no resources out there except official TensorFlow blogs.
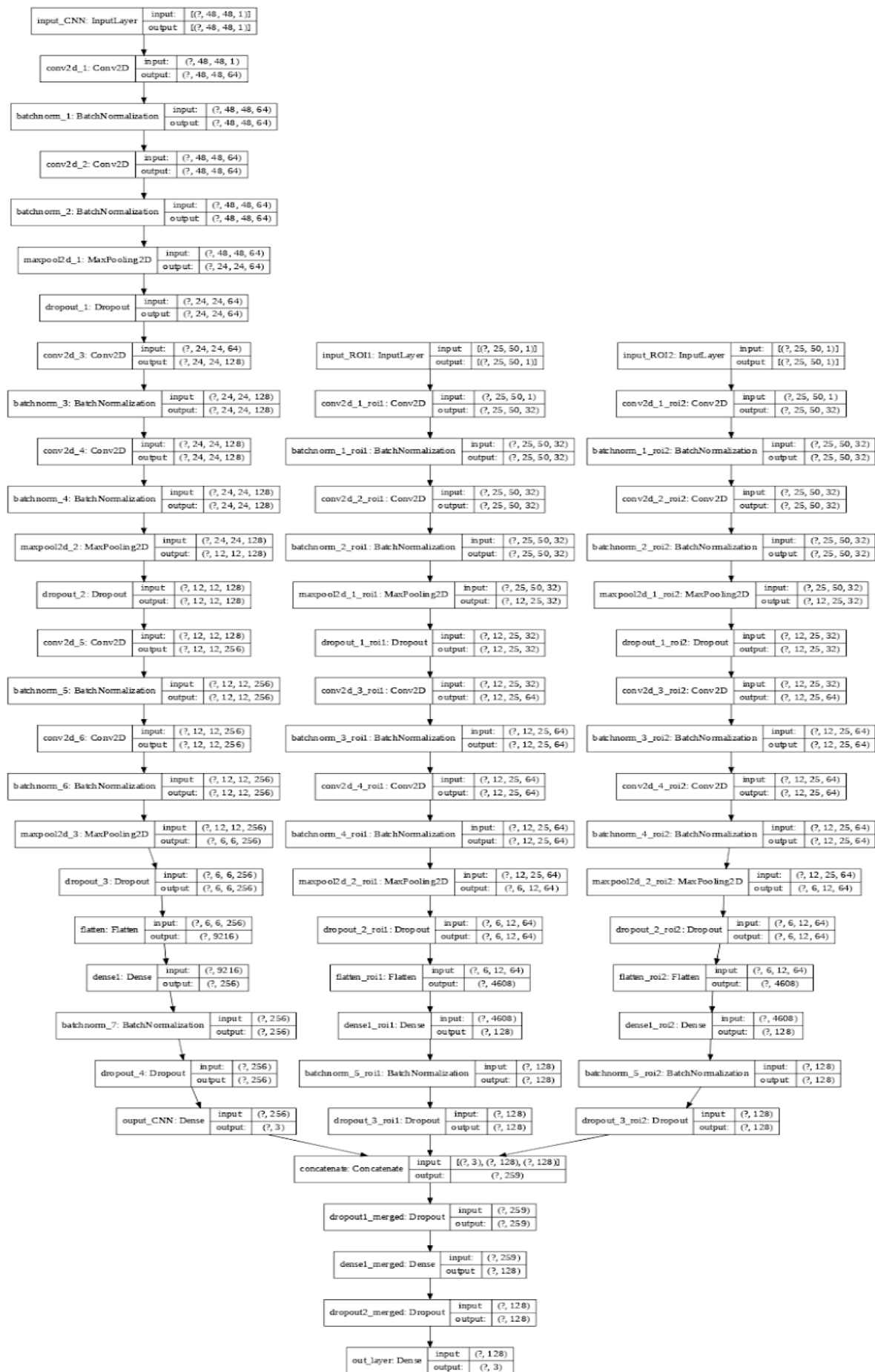
Pros:
- Add an extra regularization term
- Good for adversarially perturbed images

Cons:
- As it is introduced very recently so fewer resources are out there.
- Generally overfits
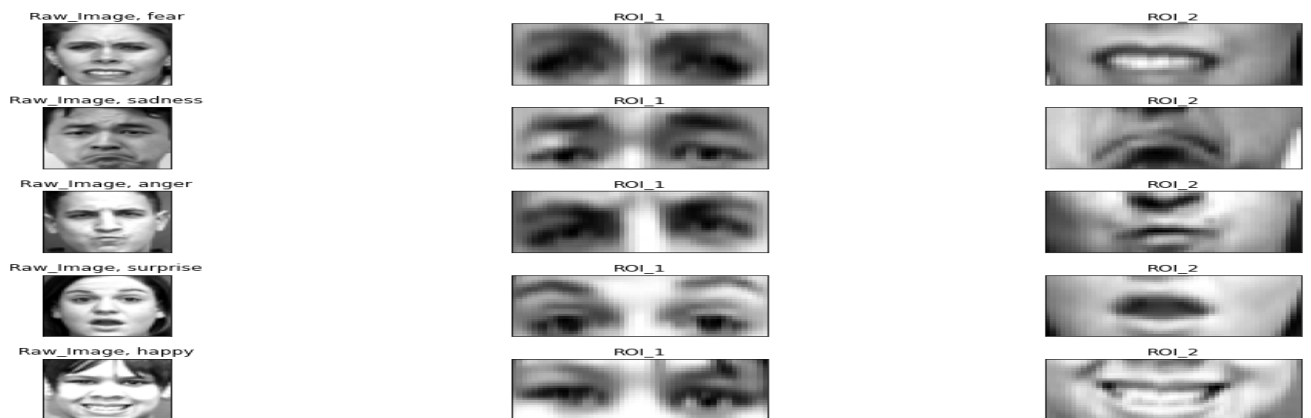- Takes more time to train

# 5. Raw Image + ROI_1 + ROI_2

This is a **multi-input model** with 3 image inputs,
1. Raw image, input_shape = 48x48x1
2. Region of interest 1 (ROI_1), this contains the region of eyebrows and eyes, input_shape = 25x50x1
3. Region of interest 2 (ROI_2), this contains region of lips, input_shape = 25x50x1

In the below image the first column is the raw image whereas the second and third columns are the ROI_1 and ROI_2 of that image respectively.



Dataset Used: CK
      Emotions: Happy, Surprise, Anger, Sadness, Fear
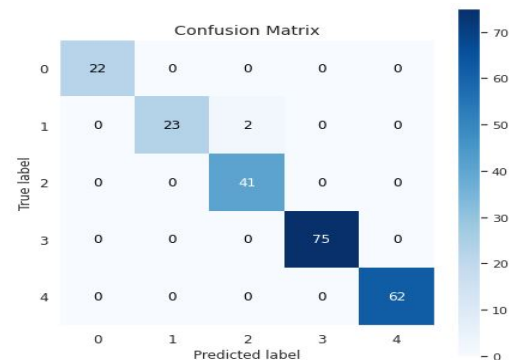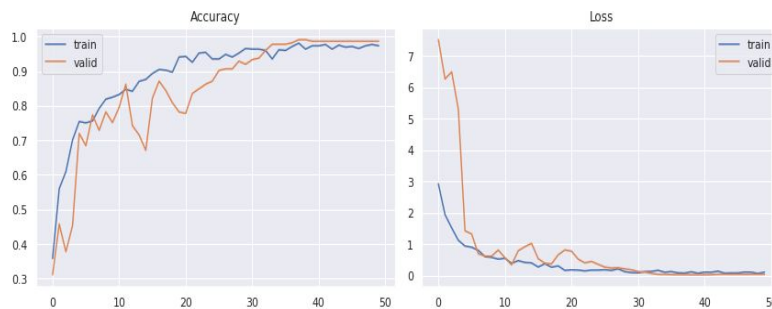      Input: Single image
      Output: Emotion in that image
      Performance: train-test split of 0.7
            Best Train Accuracy ~ 97%
            Best Test Accuracy ~ 99%



Dataset Used: FER-13 Aligned + CK
      Input: Single image
      Output: Emotion in that image
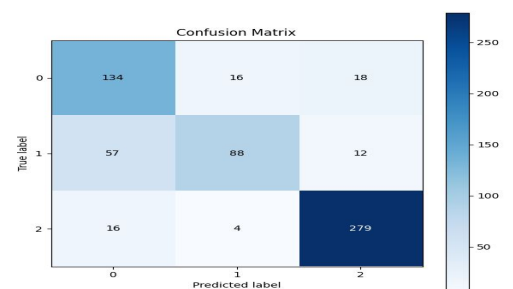
      Emotions: Happy, Sadness, Neutral
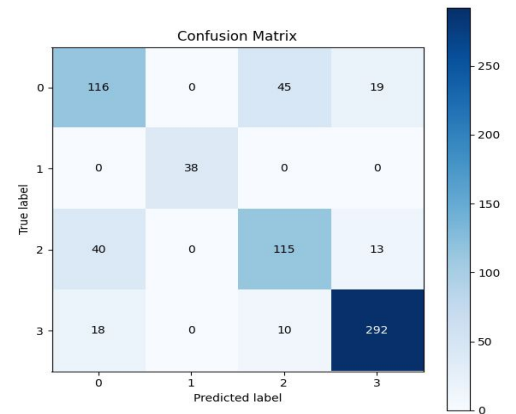      Performance: train-test split of 0.86
               Best Train Accuracy ~86%
               Best Test Accuracy ~ 80%
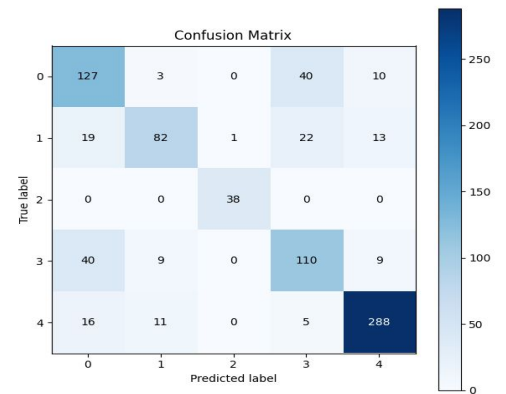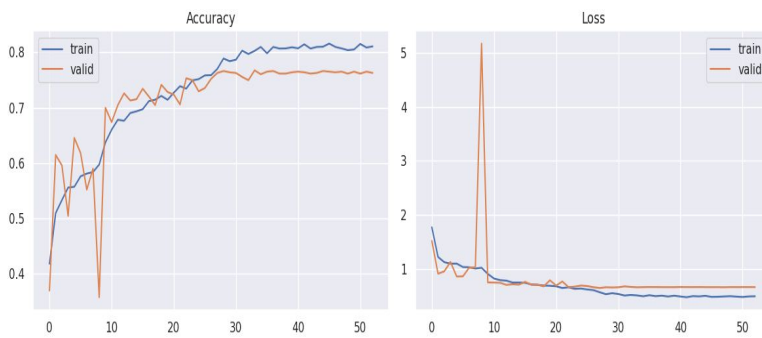      Emotions: Happy, Sadness, Neutral, Surprise

**Performance**: train-test split of 0.85
        Best Train Accuracy ~ 85%
        Best Test Accuracy ~ 80%



**Emotions**: Happy, Sadness, Neutral, Surprise, Angry
**Performance**: train-test split of 0.85
        Best Train Accuracy ~ 81%
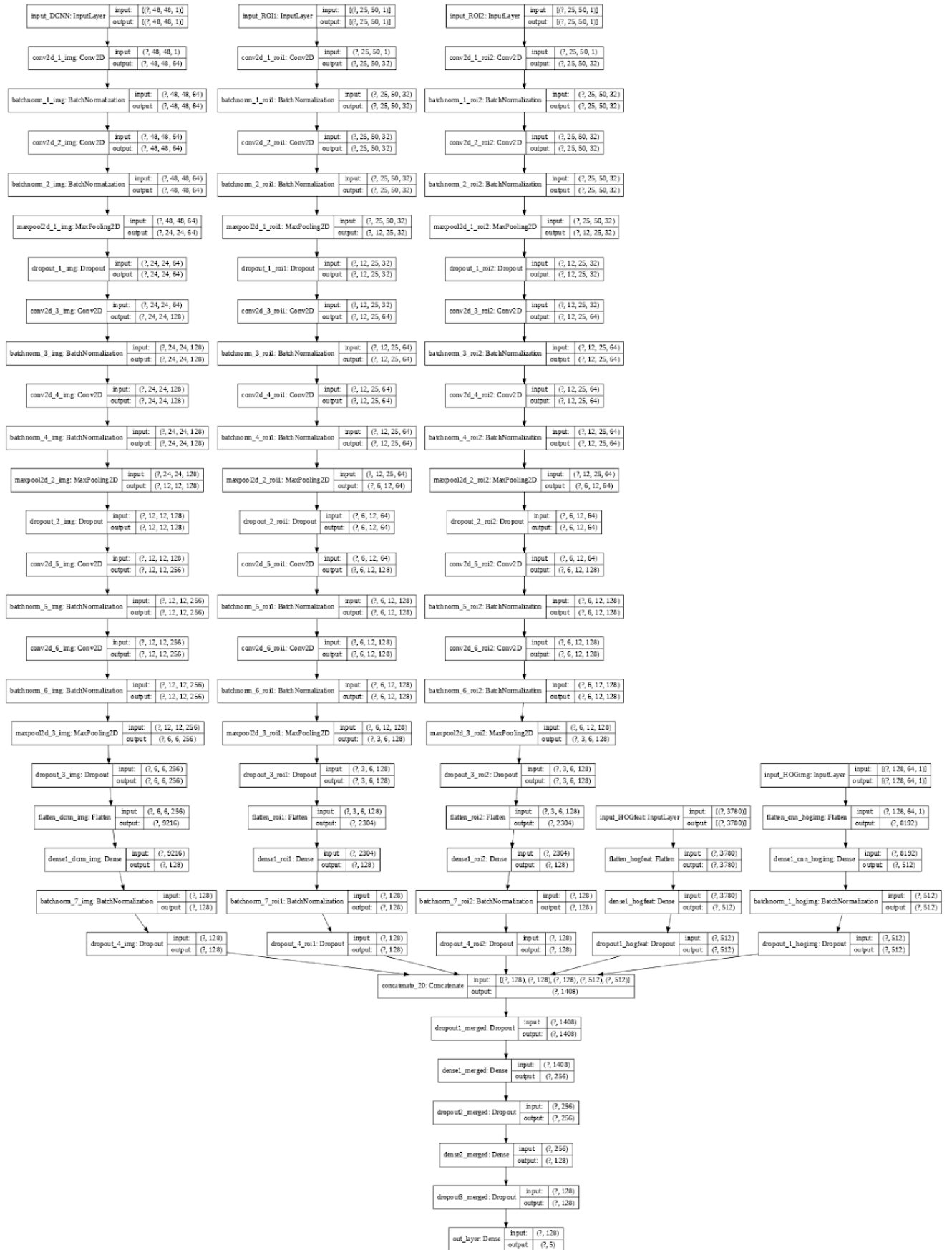        Best Test Accuracy ~ 76%



Pros:
- Process important regions separately
- Performs well especially for CK dataset

Cons:
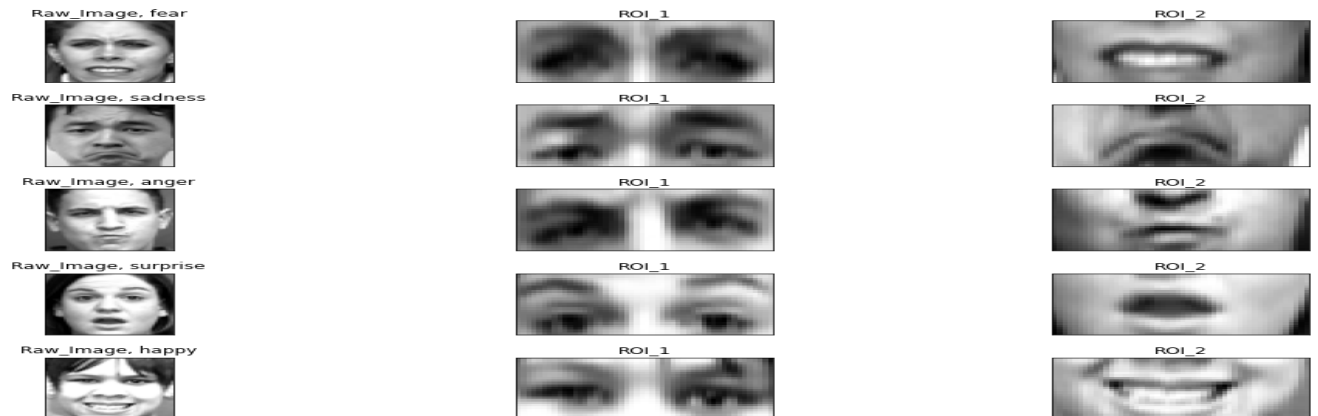- Doesn't work well with unaligned images

# 6. Raw Image + ROI_1 + ROI_2 + HOG_FEAT + HOG_IMG

This is a **multi-input model** with 5 inputs,
1. Raw image, input_shape = 48x48x1
2. Region of interest 1 (ROI_1), this contains the region of eyebrows and eyes, input_shape = 25x50x1
3. Region of interest 2 (ROI_2), this contains the region of lips, input_shape = 25x50x1
4. HOG features (HOG_FEAT), this contains hog feature descriptors, input_shape = 3780x1
5. HOG img (HOG_IMG), this contains hog feature image, input_shape = 128x64x1

In the below image the first column is the raw image whereas the second and third columns are the ROI_1 and ROI_2 of that image respectively.



Dataset Used: Ck
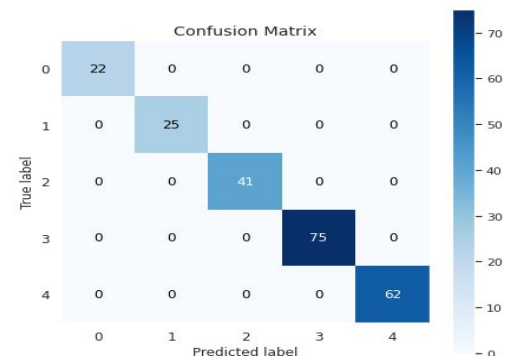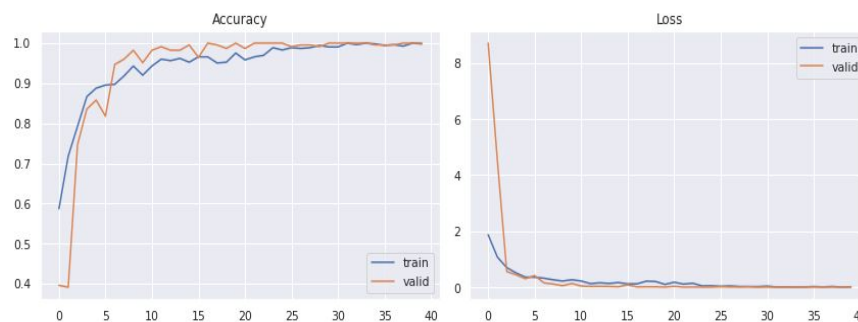Emotions: Happy, Surprise, Angry, Sadness, Fear
Input: Single image
Output: Emotion in that image
Performance: train-test split of 0.7
      Best Train Accuracy = 100%
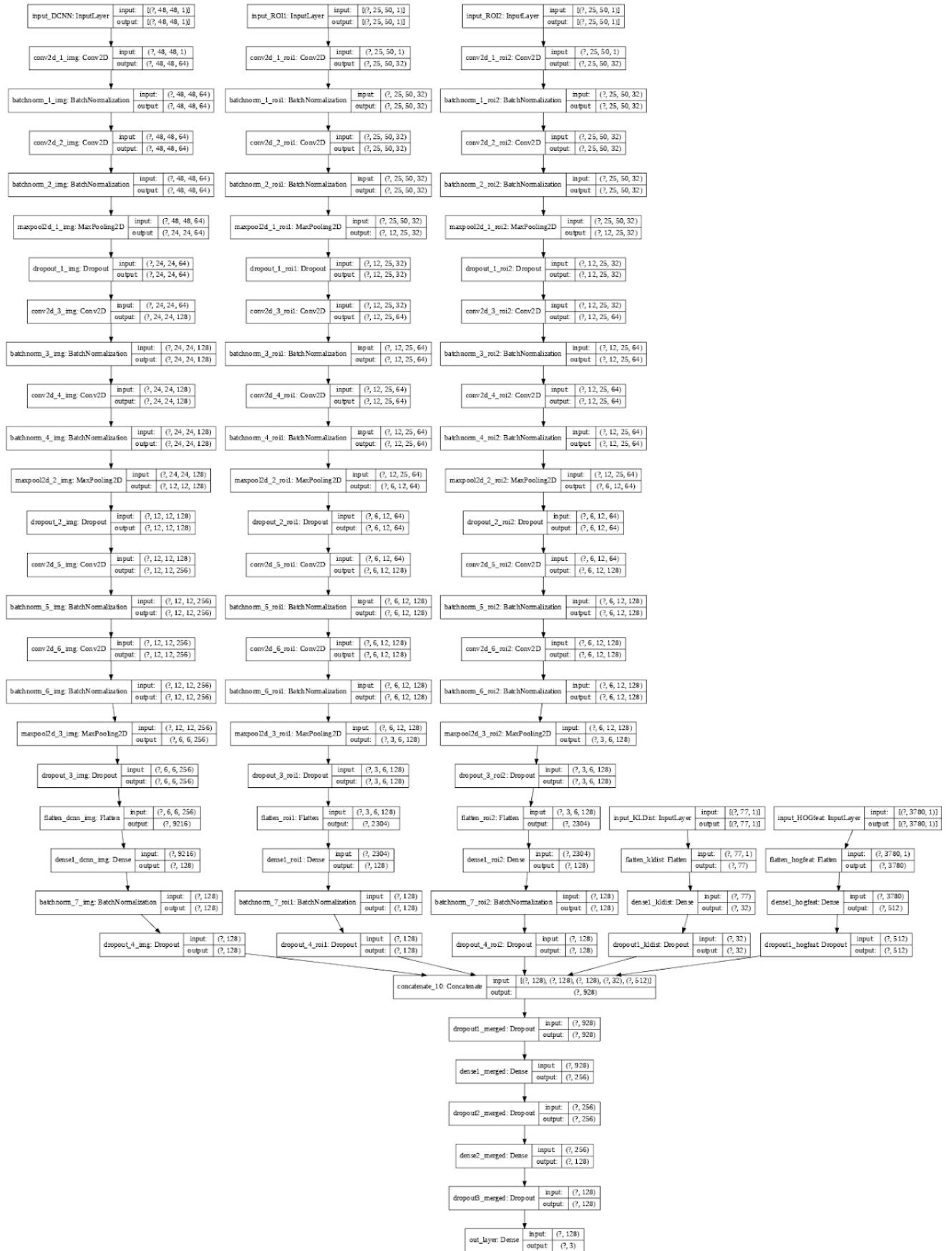      Best Test  Accuracy = 100%



Pros:
- Takes other important features of the image.
- Works really good with CK dataset

Cons:
- Performance is not that good on FER dataset
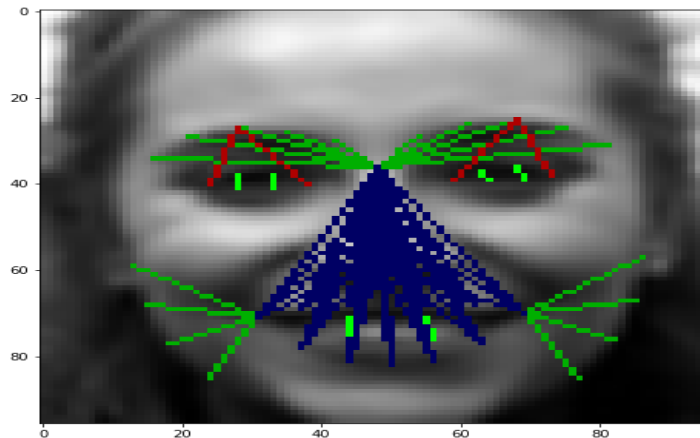- Takes more time to train and test

# 7. Raw Image + ROI_1 + ROI_2 + KL_DIST + HOG_FEAT

This is a **multi-input model** with 5 inputs,

1. Raw image, input_shape = 48x48x1
2. Region of interest 1 (ROI_1), this contains the region of eyebrows and eyes, input_shape = 25x50x1
3. Region of interest 2 (ROI_2), this contains region of lips, input_shape = 25x50x1
4. Key Landmarks distance (KL_DIST), this is the distance of important facial landmarks, input_shape = 77x1
5. HOG features (HOG_FEAT), this contains hog feature descriptors, input_shape = 3780x1

ROI_1 and ROI_2 are the same as before.



These are 77 important euclidean distances from a face. These contain distance between upper and lower lip, the distance between eyebrow and eyes, etc.

**Dataset:** FER-13 Cleaned
**Emotions:** Happy, Sadness, Neutral
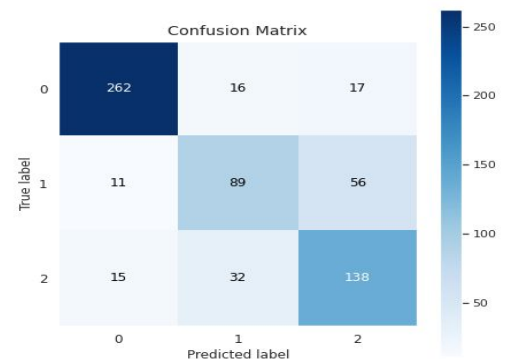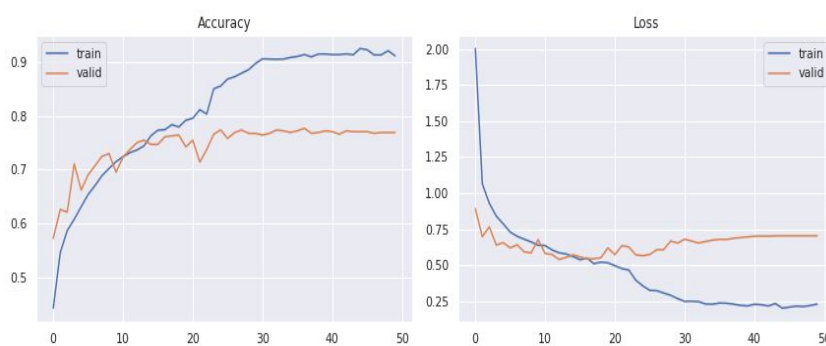**Input:** Single image
**Output:** Emotion in that image
**Performance:** train-test split of 0.86
        Best Train Accuracy ~ 92%
        Best Test Accuracy ~ 77%



Pros:
- Calculate distances between import landmarks
- Works well on CK dataset

Cons:
- These models work well on video-based datasets not image-based because of the change in landmarks shares context with previous and future frames.
- Overfits in case of FER-13

Apart from these I also tried many state-of-the-art architectures like **ResNet**, **DenseNet**, **GoogLeNet**, **XceptionNet**. You can find them [here](here).