# Computer Vision Assignment 2

Authors (Group 2):

- Aditya Shourya (i6353515)

- Ashkan Karimi Saberi (i6365010)

  code at : https://github.com/adishourya/vit_fer2013

## Structure of the Report

- The focus of the report would not be on attaining the best test score but would be on experimenting with different architectures.
- And looking through how the forward pass looks in each case.

### Convolutional Neural Network

- we will first develop a simple convolutional neural network
- And reason the effectiveness of the netwrok on a simple task as FER 2013[https://www.kaggle.com/datasets/msambare/fer2013]

### A simple hand made vision transformer (shallow)

- we will try to replicate the paper "AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE" -- adosovitskiy et.al [https://arxiv.org/pdf/2010.11929v2] and show through one full forward pass.
- We will try to reason the rate of learning by giving it same number of epochs as our CNN

### Pretrained Vision Transformer

- we will use vit_b_16 [https://pytorch.org/vision/main/models/generated/torchvision.models.vit_b_16.html#vit-b-16] and unfreeze the last few layers to perform the training.
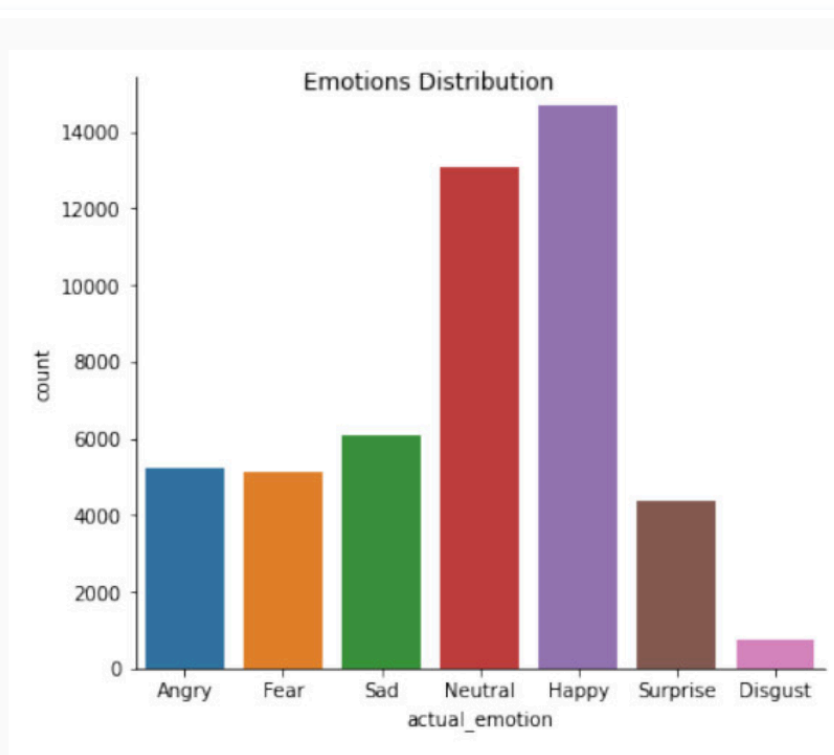
- And compare the rate of learning with shallow transformer by giving it the same number of epochs.

## Dataset

The dataset comprises **48x48** pixel grayscale images of faces. These faces have been automatically aligned to ensure that each face is roughly centered and occupies a similar amount of space in every image.

The objective is to classify each facial expression into one of seven emotional categories: 0 for Angry, 1 for Disgust, 2 for Fear, 3 for Happy, 4 for Sad, 5 for Surprise, and 6 for Neutral. The training set includes 28,709 examples, while the public test set contains 3,589 examples.

-- [https://www.kaggle.com/datasets/msambare/fer2013]



Class Diftribtion in FER2013, and the challenges asscociated with the dataset are :

The FER2013 dataset has several inherent issues that make it challenging for deep learning architectures to achieve optimal results. Key problems include imbalanced data, intra-class variation, and occlusion. Specifically, the database exhibits significant imbalance in the training data, with classes having vastly different numbers of samples. For example, the 'happy' emotion has over 13,000 samples, whereas 'disgust' has only about 600 samples, as shown in the figure above.

Intra-class variation refers to the differences within the same class. Reducing intra-class variation while increasing inter-class variation is crucial for effective classification. Variations, uncontrolled lighting conditions, and occlusions are common issues that face recognition systems encounter in real-world applications . These challenges often lead to a drop in accuracy compared to performance in controlled experimental settings. Occlusion occurs when an object blocks part of a person's face, such as a hand, hair, cap, or sunglasses. Although occlusion complicates face recognition, it can also provide useful information, as people often use their hands while communicating through gestures.

-- https://www.oaepublish.com/articles/ir.2021.16

## CNN

```
1    # define a small convolutional network
2    # see beautiful mnist in tinygrad .
3    import torch.nn as nn
4    import torch.nn.functional as F
5
6    # shape after operations n,n →(with padding p and stride s) (n + 2p - f +
     1 )/s + 1
7
8    class Net(nn.Module):
9        def __init__(self):
10           super().__init__()
11           self.conv1 = nn.Conv2d(1, 6, 5) # input channel 1 , 6 filter banks
     each of kernels size (5,5)
12           self.pool = nn.MaxPool2d(2, 2) # this is not a learnable operaation
     just performs downsampling
13           self.conv2 = nn.Conv2d(6, 16, 5) # 16 kernels
14           self.fc1 = nn.Linear(16 * 5 * 5, 120)
15           self.fc2 = nn.Linear(120, 84)
```

```
16            self.fc3 = nn.Linear(84, 7) # we have 7 classes
17
18        def forward(self, x):
19            x = self.pool(F.relu(self.conv1(x)))
20            x = self.pool(F.relu(self.conv2(x)))
21            x = torch.flatten(x, 1) # flatten all dimensions except batch
22            x = F.relu(self.fc1(x))
23            x = F.relu(self.fc2(x))
24            logits = self.fc3(x)
25            return logits
26
27
28    net = Net()
```

## forward pass in the convolutional layer
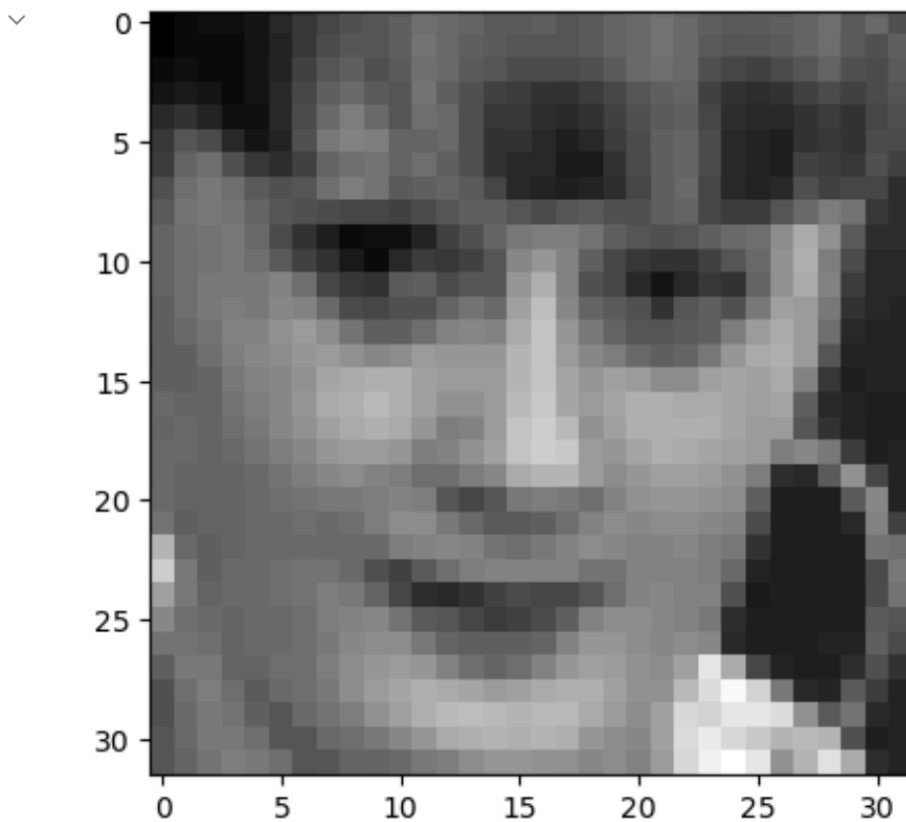
```
1    # lets visualize the first layer on a sample image
2    x = self.pool(F.relu(self.conv1(x)))
3
4    # 1,6 input channel = 1 (greyscale)
5    # output dimension = 6 (filterbanks)
6    # kernel size = (5,5) kernels
7    net.conv1 , net.conv1.weight.shape
8
9    # > out: (Conv2d(1, 6, kernel_size=(5, 5), stride=(1, 1)), torch.Size([6, 1,
     5, 5]))
```

```python
In 52    1   import einops
         2   sample_img = img[0]
         3   plt.imshow(
         4       einops.rearrange(sample_img,"1 h w → h w"),
         5       cmap="grey")
```
Executed at 2024.05.24 22:29:07 in 247ms

Out 52    <matplotlib.image.AxesImage at 0x15796c460>



```python
 1   import torch.nn.functional as F
 2   with torch.no_grad():
 3       # first layer convolution
 4       out1 =F.conv2d(sample_img,net.conv1.weight,
 5                   bias=None, stride=1, padding=0)
 6       print(out1.shape) # input channel = 1 filter banks = 6
 7
 8
 9   plt.imshow(
10       einops.rearrange(out1,"out_c h w → h (out_c w)"),
```
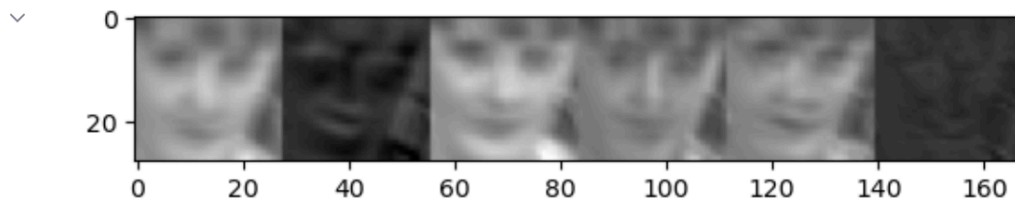
```
11          cmap="grey"
12      )
13      # visualizing the output of convolution from the first learnt layer.
14      # notice how some kernels are vastly different
```

```
torch.Size([6, 28, 28])
```

Out 63          `<matplotlib.image.AxesImage at 0x157e77550>`



- we will not show the pooling operation here (check notebook) as it is not a learnable parameter

```
1   with torch.no_grad():
2       # first layer convolution
3       out2 =F.conv2d(out1_p,net.conv2.weight,
4                   bias=None, stride=1, padding=0)
5       print(out2.shape) # input channel = 1 filter banks = 6
6
7   plt.figure(figsize=(15,8))
8   plt.imshow(
9       einops.rearrange(out2,"out_c h w → h (out_c w)"),
10      cmap="grey"
11  )
```

```
torch.Size([16, 23, 23])
```

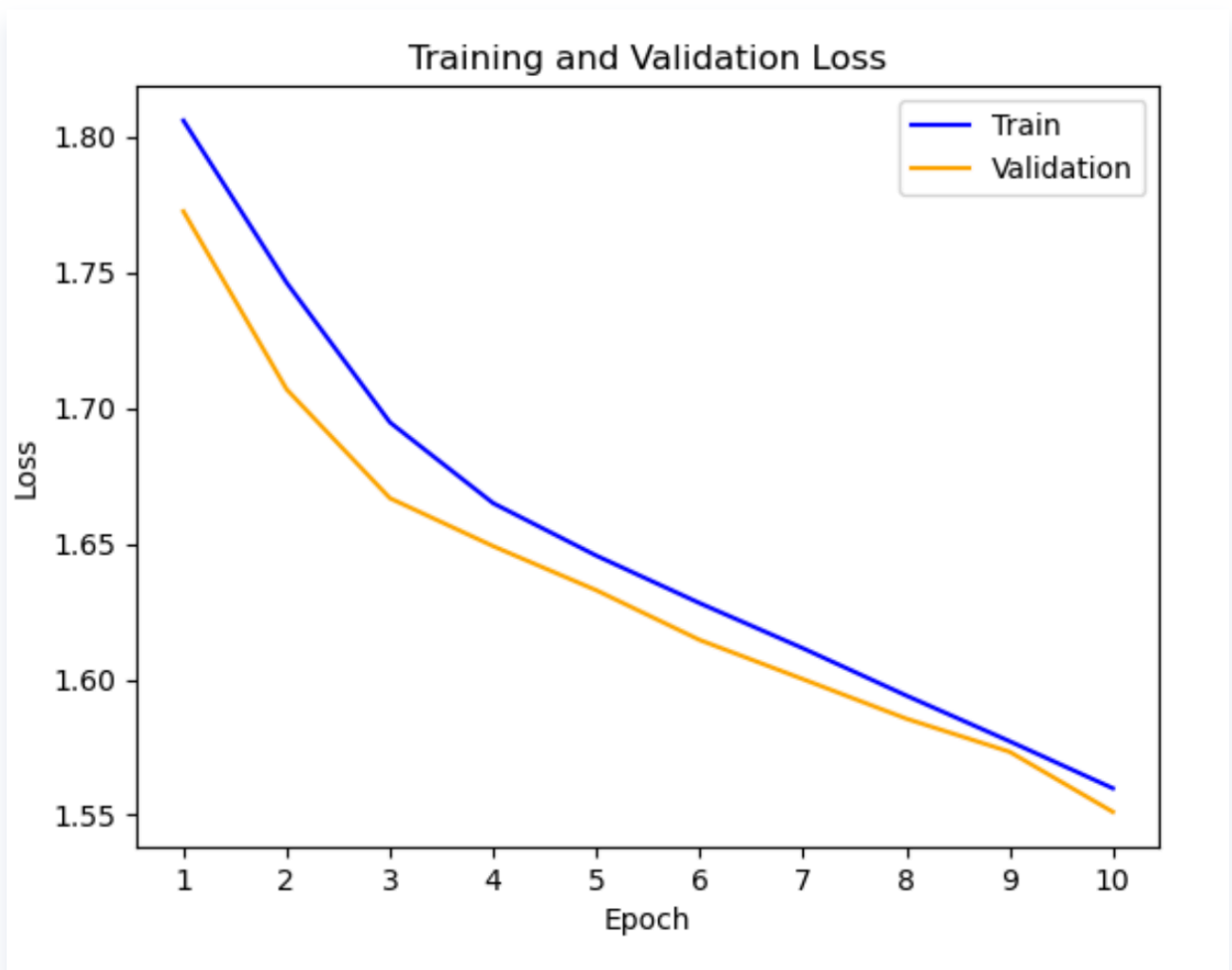`<matplotlib.image.AxesImage at 0x157fcbdc0>`



- and then we maxpool it and then flatten it to pass to a feed forward network to arrive at the logits

for classification

- we only train for 10 epochs as attaining the best test accuracy is not our main goal of this report



Training and Validation Loss

```
Executed at 2024.05.24 21:24:47 in 2s 129ms

100%|████████| 113/113 [00:02<00:00, 53.20it/s]

Test Loss: 1.556882
```

```
 1    Test Accuracy of Classes (Generalization scores of our Convolutional Model)
 2
 3    Angry : 10%    (54/491)
 4    Disgust : 0%   (0/55)
 5    Fear  : 7%    (41/528)
 6    Happy : 80%    (710/879)
 7    Sad : 32%    (191/594)
 8    Surprise  : 47%    (197/416)
 9    Neutral : 33%    (208/626)
10
11    Test Accuracy of Dataset:    39%    (1401/3589)
12
```

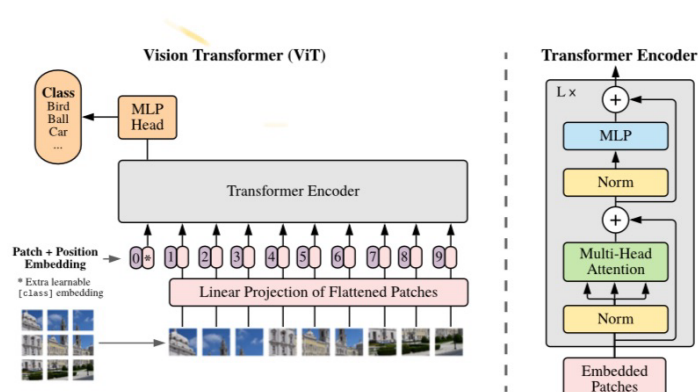# Shallow Transformer Replicating paper



Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable "classification token" to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).
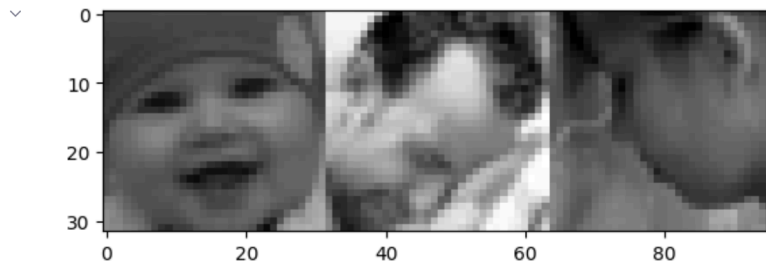
3 METHOD

$$N = \frac{HW}{P^2} = \frac{32 \times 32}{16 \times 16} = 4$$

$$P \in \mathbb{R}^{16, 16, C=1}$$

- Images to 16 * 16 patches



```
In 105   1 v # now what we want is to turn them into patches
         2   # patches = einops.rearrange(xb,"b 1 (p1 h) (p2 w) → b (p1 p2) (h w)",p1=2,p2=2)
         3   patches = einops.rearrange(xb,"b 1 (p1 h) (p2 w) → h (b p1 p2 w)",p1=2,p2=2) # increase p1,p2 to
             get 16 x 16 patches
         4   # here p1 and p2 must be 14 x 14
         5   plt.figure(figsize=(150,80))
         6   plt.imshow(patches, cmap="gray")
         7   plt.show()
```
Executed at 2024.05.19 11:53:46 in 90ms



```python
1    # check appendix for unit test of this class
2    class SingleHead(nn.Module):
3        """
4        Implements a single head of attention (unmasked)
5        """
6
7        def __init__(self,n_embed=32,head_size=8):
8            super().__init__()
9            # single head
10           self.head_size = torch.tensor(head_size)
11           self.n_embed = torch.tensor(n_embed)
12           self.Q = nn.Parameter( torch.randn(self.n_embed,head_size) *
     (1/torch.tensor(2.82)))
13           self.K = nn.Parameter( torch.randn(self.n_embed,head_size) *
     (1/torch.tensor(2.82)))
14           self.V = nn.Parameter( torch.randn(self.n_embed,head_size) *
     (1/torch.tensor(2.82)))
15
16       def forward(self,x):
17           query = x @ self.Q
18           key =   x @ self.K
```

```python
19            value=  x @ self.V
20
21            # hand implementation
22            # scale ⟹ sqrt head size
23            scale = 1 / torch.sqrt(self.head_size)
24
25            # we will not use any masking here as its an image
26            # and no dropout consideration in this implementation
27            comm = query @ key.transpose(-2,-1)
28            comm = comm* scale
29            soft_comm = torch.softmax(comm, dim=2)
30            att = soft_comm @ value
31
32            return att
```

```python
1  class Multihead(nn.Module):
2      def __init__(self,n_embed,n_heads):
3          super().__init__()
4
5          self.n_embed = n_embed
6          self.n_heads = n_heads
7          self.head_size = self.n_embed // self.n_heads
8
9          self.multiheads = nn.ModuleList(
10              [SingleHead(self.n_embed,self.head_size)
11               for _ in range(self.n_heads)]
12          )
13
14      def forward(self,x):
15          return torch.cat([head(x) for head in self.multiheads],dim=2)
```

```python
1  # only multihead →  skip connection → layernorm
2  # Batch norm : couples examples in and normalizes it .. (also has a
   regularization effect) but we need to keep a running mean to track new mean
   and sigma
3  # layernorm : normalizes the features of each example (does not couple
   examples across the batch) more popular in transformers
4
5  class TranformerBlock(nn.Module):
```

```python
    def __init__(self, n_embed,n_head):
        super().__init__()
        self.multi_head = Multihead(n_embed,n_head)
        # i am not going to implement my own layer norm it wont be
efficient and will be janky at best
        self.norm = nn.LayerNorm(n_embed) # we want to normalize feeatures
(each patch gets normalized)

    def forward(self,x):
        # pass through multihead
        attention = self.multi_head(x)
        # skip connection and non linarity
        attention = torch.relu( x + attention)
        # layer norm
        attention = self.norm(attention)
        return attention # B , n_patch , n_embed
        ...
```

```python
# most of the comments are pasted verbatim from the paper
class SmallVIT(nn.Module):

    def __init__(self):
        super().__init__()
        # patches
        # embedding
        self.vocab_size = torch.tensor(256) # 0 to 255 pixels
        # each patch will only get one n_embed representation
        self.n_embed = 32 # we will project each patch 16*16 to a 32
dimensional representation
        # so the lookup table would be of the shape
        # unlike in nlp where we embed token to a vector like below we
would project matrix of patch size to a vector
        # self.C = nn.Embedding(self.vocab_size, self.n_embed)

        self.C = nn.Parameter(torch.randn(self.vocab_size, self.n_embed) *
1/torch.sqrt(self.vocab_size) )

        # positional embedding
        # the paper says Position embeddings are added to the patch
embeddings to retain positional information. We use standard learnable 1D
position embeddings, since we have not observed significant performance
gains from using more advanced 2D-aware position embeddings
```

```python
        self.pe = nn.Parameter(torch.randn(1,4,self.n_embed)) # each pach
and representation should get positional embedding

        # we use the standard approach of adding an extra learnable
"classification token" to the sequence
        self.classification_token = nn.Parameter(torch.randn(1, 1,
self.n_embed))
        # we will keep the step above optional .. i dont understand why we
should use it yet.

        # transformer block
        self.n_heads = 4 # we will use 4 heads for now
        self.transformer_block = TranformerBlock(self.n_embed,self.n_heads)

        # MLP Head for final logit calculation
        # n_patch * n_embed → fer["emotion"].nuinque() : 7
        self.mlp_head = nn.Parameter(torch.randn(4*32, 7) *
torch.sqrt(torch.tensor(4*32)))




        ...

    def forward(self,X):
        batch_size = X.shape[0]
        patches = X.view(-1,4,256) # B , p_num , 16*16
        # B , p_num , 256 @ 256 , 32
        emb = patches @ self.C # B , p_num , n_embed
        emb = emb + self.pe # kind of acts like a bias towards each
patches.

        # 2 transformers
        tf = self.transformer_block(emb)
        tf = self.transformer_block(tf)

        # flatten it : across patches
        tf = tf.view(batch_size,-1)

        # logits
        logits = tf @ self.mlp_head

```
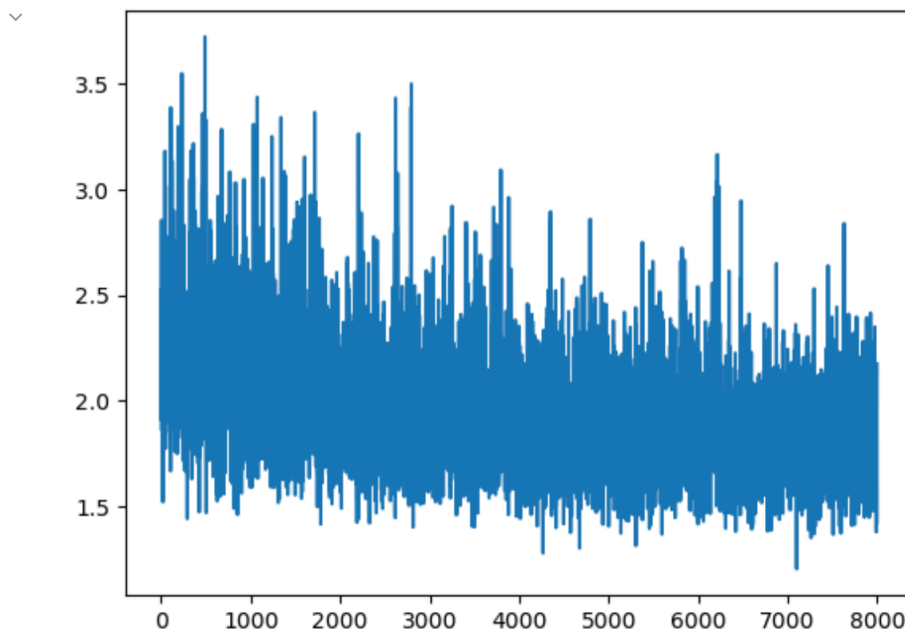
```
56
57
58          # broadcasting steps in the above command
59          # B,p_n , p*p , n_embed
60          #1,4,32
61
62
63          return logits
64
65
```

In 101  1  `plt.plot(track_loss[-8000:])`
Executed at 2024.05.19 21:54:47 in 129ms

Out 101     [<matplotlib.lines.Line2D at 0x1436d3250>]



In 95  1  `np.median(`
       2      `track_loss[:-10]`
       3  `)`
Executed at 2024.05.19 21:53:49 in 20ms
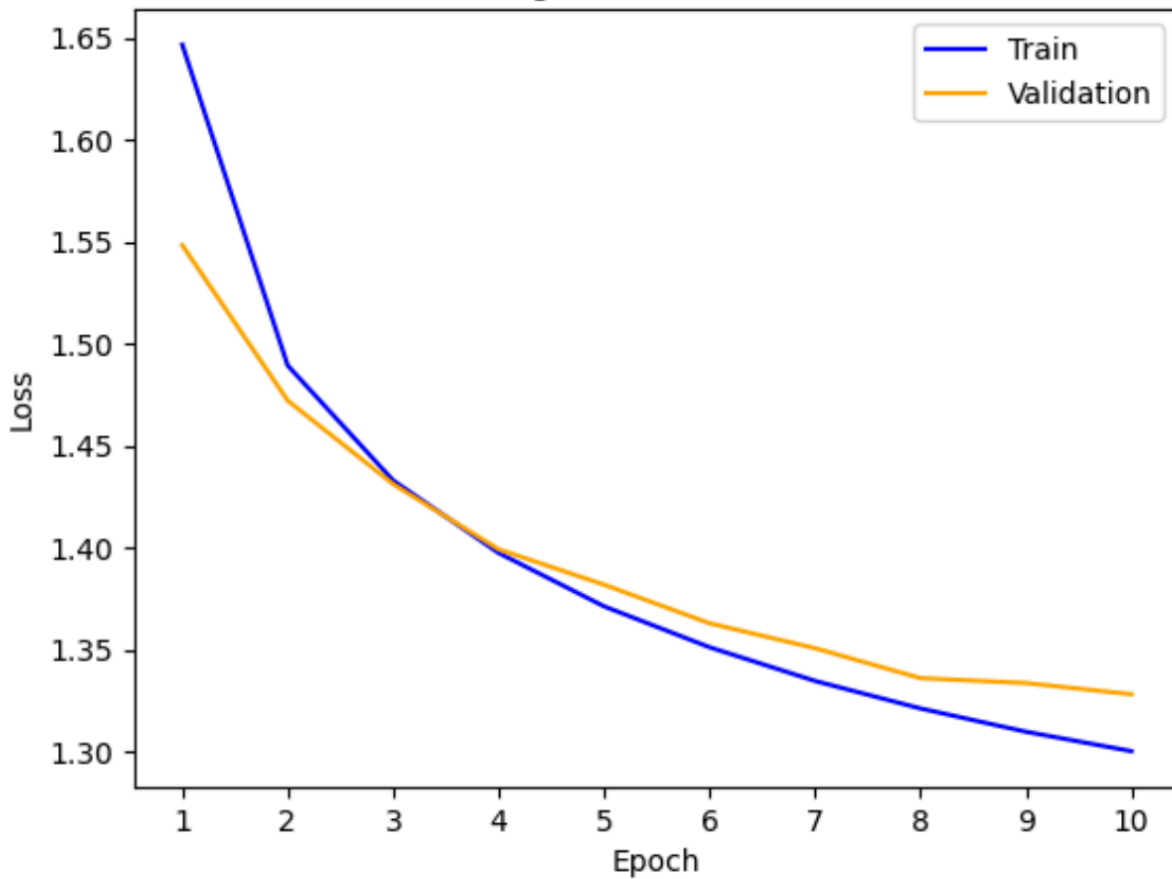
Out 95     1.9874691367149353

# Pretrained Vision Transformer (perform transfer learning)

```
52
53   # output test loss statistics
54   print('Test Loss: {:.6f}'.format(test_loss))
55
```

```
100%|██████████| 113/113 [36:49<00:00, 19.55s/it]

Test Loss: 1.314458
```

Training and Validation Loss

```
1    ## Results for each class
2    Disgust : 1%    (1/55)
3    Fear  : 29%    (155/528)
4    Happy : 78%    (688/879)
5    Sad : 35%    (211/594)
6    Surprise  : 62%    (260/416)
7    Neutral : 54%    (344/626)
8
9    Test Accuracy of Dataset:    50%   (1799/3589)
```

## Remarks

## Appendix

# unit test for self attention with pytorch's implementation

```python
def single_head(query, key,value):
    head_size = torch.tensor(query.shape[-1])
    # hand implementation
    # scale ⇒ sqrt head size
    scale = 1 / torch.sqrt(head_size)

    # we will not use any masking here as its an image
    # and no dropout consideration in this implementation
    comm = query @ key.transpose(-2,-1)
    comm = comm* scale
    soft_comm = torch.softmax(comm, dim=2)
    att = soft_comm @ value
    print(att.shape)
    return att


g=torch.Generator().manual_seed(123)
query, key, value = torch.randn(2, 3, 8 , generator = g), torch.randn(2, 3,
8, generator = g), torch.randn(2, 3, 8 , generator = g)

# our implementation
sh  = single_head(query,key,value)

# pytorch implementation
py_sa = nn.functional.scaled_dot_product_attention(query, key, value)

# > torch.allclose(py_sa , sh) prints True
```