BME646 and ECE60146: Homework 9

Spring 2023

Arghadip Das

das169@purdue.edu

1. Introduction

Over the past few years, convolutional neural networks (CNNs) have dominated the field of computer vision, achieving state-of-the-art performance on a wide range of visual recognition tasks. However, these networks have several limitations, such as a fixed receptive field and a lack of attention mechanisms for modeling long-range dependencies. To address these issues, a new type of neural network architecture called the Vision Transformer (ViT) has been proposed. The ViT is based on the Transformer architecture originally developed for natural language processing (NLP) and uses self-attention mechanisms to model long-range dependencies between image patches. The ViT has quickly gained attention in the computer vision community due to its impressive performance on various image classification tasks, surpassing the previous state-of-the-art results achieved by CNNs. In addition, the ViT has shown promising results in other computer vision tasks such as object detection, segmentation, and generation.

In this homework, we have several goals. Firstly, we aim to gain a deeper understanding of the multi-headed self-attention mechanism and the transformer architecture. Secondly, we hope to comprehend how the transformer architecture, which was originally developed for language translation, can be readily adapted to process images in the Vision Transformer (ViT). Finally, we are required to implement our own ViT for image classification. By achieving these goals, we will not only learn more about the ViT architecture but also gain valuable insights into the underlying principles of modern neural networks for image processing.

2. Methodology

In this report, I detail the steps taken to complete the homework assignment. To prepare for this homework, we first reviewed Professor Kak's slides on self-attention and gained a better understanding of the self-attention mechanism and its implementation through matrix multiplication. Additionally, we learned how multiple self-attention heads can work in parallel in multi-headed attention to capture inter-word dependencies. We also studied the encoder-decoder structure of a transformer for sequence-to-sequence translation and experimented with seq2seq_with_transformerFg.py and seq2seq_with_transformerPreLN.py. Finally, we went through the ViT paper to understand the fundamental concepts behind the Vision Transformer (ViT) for image classification. In particular, we closely examined Figure 1, which provides a clear illustration of how an image can be transformed into a sequence of embeddings and processed using a transformer encoder. We paid special attention to how the class token is prepended as a learnable parameter to the input patch embedding sequence, as well as how the same token is taken from the final output sequence to produce the predicted label.

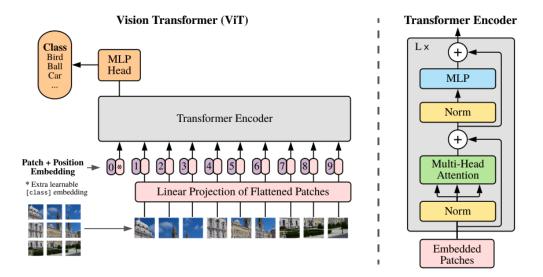


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable "classification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

To implement our ViT, we utilized the Google Colab platform and followed the instructions provided in the homework manual. We reused the dataset, training, and evaluation scripts from HW4, but replaced the HW4 CNN with our ViT implementation. After training the ViT, we generated a confusion matrix on our test set to evaluate its performance and compared it to CNN-based networks. In addition, we also attempted to implement a multi-headed self-attention mechanism using torch.einsum, which we were able to accomplish within just 10 lines of code.

3. Programming Tasks

3.1. Building Our Own ViT

Before starting, I reviewed the provided VitHelper.py file. This file includes all the necessary classes for us to build a transformer from scratch quickly. We can find these classes in DLStudio's Transformers.TransformerPreLN module. These classes are now standalone so that we can use them directly in our ViT implementation. Specifically, we will use the MasterEncoder class as the "Transformer Encoder" block, as shown in Figure 1.

Initially, I utilized a linear layer for converting the patch into an embedding, as illustrated in Figure 1. But then, I incorporated a Conv2D based embedding which turned out to be significantly faster than the linear layer-based conversion. This is because we can apply Conv2D directly on the complete image without dividing it into patches first. However, it requires appropriate tensor reshaping to work effectively. We are working with images of size 64x64 and a patch size of 16x16, resulting in a total of 16 patches. To prepare the input sequence for ViT, we *prepend a*

class token to our patch sequence. To account for this, we set the maximum sequence length of the transformer to 17. For the class token, it must be set as learnable parameter, which is implemented using nn.Parameter. In contrast to the sinusoid-based position embedding commonly used in language processing, ViT makes use of learnable position embeddings. Therefore, besides initializing the class token, I have also defined the position embeddings as learnable parameters. We obtain the final class prediction by extracting the class token from the output sequence generated by the Transformer Encoder block. The class token is then passed through a Multilayer Perceptron (MLP) layer to obtain the logits for the 5 classes. Here in this report, we have included the libraries and helper functions provided in ViTHelper.py to ensure that the report can be fully understood without referring to any external files. Following this, the source code for ADVit is presented and it is sufficiently commented to make it easy to comprehend.

Source Code:

Importing libraries and getting the device:

Helper classes and functions:

```
super().__init__()
self.max_seq_length = max_seq_length
self.embedding_size = embedding_size
self.qkv_size = self.embedding_size // num atten heads
self.num_atten_heads = num_atten_heads
    max_seq_length, embedding_size, num_atten_heads) # (A)
self.norm1 = nn.LayerNorm(self.embedding_size) # (C)
self.W1 = nn.Linear(self.max_seq_length * self.embedding size,
self.norm2 = nn.LayerNorm(self.embedding size) # (E)
output_self_atten = self.self_attention_layer(
    normed_input_self_atten).to(device) # (F)
input for FFN = output self atten + input for self atten
    self.W1(normed_input_FFN.view(sentence_tensor.shape[0], -1))) # (K)
basic_encoder_out = basic_encoder_out.view(
    sentence_tensor.shape[0], self.max_seq_length, self.embedding_size)
basic encoder out = basic encoder out + input for FFN
  init (self, max seg length, embedding size, num atten heads):
```

```
super().__init__()
        self.max_seq_length = max_seq_length
        self.num_atten_heads = num atten heads
         self.qkv_size = self.embedding_size // num_atten_heads
range(num atten heads)]) # (A)
         concat_out_from_atten_heads = torch.zeros(sentence_tensor.shape[0], self.max_seq_length,
        self.attention_heads_arr[i](sentence_tensor_portion) # (D)
return concat_out_from_atten_heads
        __init__(self, max_seq_length, qkv_size):
super().__init__()
        self.max_seq_length = max_seq_length
                               max_seq_length * self.qkv size) # (B)
        self.WK = nn.Linear(max_seq_length * self.qkv_size,
        max_seq_length * self.qkv_size)
self.WV = nn.Linear(max_seq_length * self.qkv_size,
                               max seq length * self.qkv size) # (D)
         Q = self.WQ(sentence_portion.reshape(
             sentence_portion.shape[0], -1).float()).to(device) # (G)
         V = self.WV(sentence_portion.reshape(
         Q = Q.view(sentence_portion.shape[0],
                     self.max_seq_length, self.qkv_size) # (J)
                     self.max_seq_length, self.qkv_size) # (K)
        V = V.view(sentence portion.shape[0],
                     self.max\_seq\_length, self.qkv\_size) \quad \# \ (\texttt{L})
        A = K.transpose(2, 1) \# (M)
QK dot prod = Q @ A # (N)
        rowwise_softmax_normalizations = self.softmax(QK_dot_prod) # (0)
         Z = rowwise softmax normalizations @ V
```

Vision Transformer (ADVit) Implementation:

```
how_many_basic_encoders, mlp_dim, num_classes=5, patch_to_embed_method='Conv'):
   super().__init__()
   self.image_size = image_size
   self.embedding size = embedding size
   self.num patches = (image size // patch size) ** 2
    self.max_seq_length = self.num_patches + 1
     self.patch embeddings = nn.Linear(self.patch dim, embedding size)
```

```
self.patch embeddings = nn.Conv2d(3, embedding size, kernel size=patch size,
stride=patch size)
    self.encoder = MasterEncoder(self.max seq length, embedding size, how many basic encoders,
                                 nn.Linear(self.mlp dim, 64),
    elif self.patch_to_embed_method == 'Conv':
   x = self.patch_embeddings(x)
   x = x.view(x.size(0), -1, x.size(1))
    class token = self.class token.expand(x.size(0), -1, -1)
    x = self.encoder(x)
```

ADVit takes in an image and outputs a prediction of its class. The input image is assumed to be square with a specified image size and three channels for RGB. The image is divided into non-overlapping patches of a specified patch size, and each patch is embedded into a vector of a specified embedding size using <u>either a linear layer or a convolutional layer</u>. The module also learns learnable positional embeddings of the same size as the patch embeddings, and it includes an encoder made up of a specified number of basic encoders that each have a specified number of attention heads. Additionally, a learnable class token is added to the input embeddings, and the output of the class token is fed through an MLP to produce a prediction of the input image's class.

In the __init__ function, the relevant parameters are set, and the various components of the module are defined. The forward function takes in an input image and performs the necessary operations to convert the image into a sequence of embeddings that can be fed into the encoder. If linear patch embeddings are used, the input image is divided into patches using the unfold function, and each patch is embedded using the linear layer. If convolutional patch embeddings are used, the input image is passed through a convolutional layer, resulting in a tensor of size (batch_size, embedding_size, num_patches, num_patches), and the tensor is reshaped to (batch_size, num_patches, embedding_size). In both cases, the class token is added to the sequence of embeddings, and the positional embeddings are added to the embeddings. The resulting sequence is fed through the encoder, and the output of the class token is passed through an MLP to produce a prediction of the input image's class.

We get the following details of the ADVit. In the following code, we can see the <u>hyperparameters</u> <u>used in ADVit.</u> The mentioned shapes of the tensors and variables in comments in the earlier ADVit code are based on these hyperparameters.

The number of layers in the model: 122

The number of learnable parameters in the model: 41577829

AttentionHead-22	[-1, 17, 16]	
Linear-23	[-1, 272]	74,256
Linear-24 Linear-25	[-1, 272] [-1, 272]	
Softmax-26	[-1, 17, 17]	
AttentionHead-27	[-1, 17, 16]	
Linear-28	[-1, 272]	
Linear-29	[-1, 272]	74 , 256
Linear-30	[-1, 272]	
Softmax-31	[-1, 17, 17]	0
AttentionHead-32 Linear-33	[-1, 17, 16] [-1, 272]	
Linear-34	[-1, 272]	
Linear-35	[-1, 272]	
Softmax-36	[-1, 17, 17]	
AttentionHead-37	[-1, 17, 16]	
Linear-38	[-1, 272]	
Linear-39 Linear-40	[-1, 272] [-1, 272]	
Softmax-41	[-1, 2/2]	
AttentionHead-42	[-1, 17, 16]	0
SelfAttention-43	[-1, 17, 128]	
LayerNorm-44	[-1, 17, 128]	256
Linear-45	[-1, 4352]	
Linear-46	[-1, 2176]	9,472,128
BasicEncoder-47 LayerNorm-48	[-1, 17, 128] [-1, 17, 128]	
Linear-49	[-1, 17, 120]	
Linear-50	[-1, 272]	
Linear-51	[-1, 272]	74,256
Softmax-52	[-1, 17, 17]	
AttentionHead-53	[-1, 17, 16]	0
Linear-54	[-1, 272]	74,256
Linear-55 Linear-56	[-1, 272] [-1, 272]	
Softmax-57	[-1, 2/2] [-1, 17, 17]	
AttentionHead-58		
Linear-59	[-1, 17, 16] [-1, 272]	74,256
Linear-60	[-1, 272]	74,256
Linear-61	[-1, 272]	
Softmax-62	[-1, 17, 17]	U
AttentionHead-63 Linear-64	[-1, 17, 16] [-1, 272]	74 , 256
Linear-65	[-1, 272]	74,256
Linear-66	[-1, 272]	
Softmax-67	[-1, 17, 17]	
AttentionHead-68	[-1, 17, 16]	0
Linear-69 Linear-70	[-1, 272]	74,256
Linear-71	[-1, 272] [-1, 272]	74,256 74,256
Softmax-72	[-1, 17, 17]	
AttentionHead-73	[-1, 17, 16]	
Linear-74	[-1, 272]	74 , 256
Linear-75	[-1, 272]	
Linear-76 Softmax-77	[-1, 272]	
AttentionHead-78	[-1, 17, 17] [-1, 17, 16]	
Linear-79	[-1, 272]	
Linear-80	[-1, 272]	
Linear-81	[-1, 272]	74,256
Softmax-82	[-1, 17, 17]	
AttentionHead-83	[-1, 17, 16]	74 256
Linear-84 Linear-85	[-1, 272] [-1, 272]	74,256 74,256
Linear-85	[-1, 2/2] [-1, 272]	
Softmax-87	[-1, 17, 17]	
AttentionHead-88	[-1, 17, 16]	
SelfAttention-89	[-1, 17, 128]	
LayerNorm-90	[-1, 17, 128]	256
Linear-91	[-1, 4352]	
Linear-92	[-1, 2176]	9,472,128

```
BasicEncoder-93 [-1, 17, 128] 0
MasterEncoder-94 [-1, 17, 128] 0
Linear-95 [-1, 96] 12,384
Linear-96 [-1, 64] 6,208
Linear-97 [-1, 5] 325

Total params: 41,575,525
Trainable params: 0
Input size (MB): 0.05
Forward/backward pass size (MB): 0.43
Params size (MB): 158.60
Estimated Total Size (MB): 159.08
```

3.2.Image Classification with ADVit

To accomplish this task, we will reuse the training and evaluation scripts we developed in HW4. Instead of using the HW4 CNN, we will replace it with our ViT model. We will also use the COCO-based dataset that we created for HW4, which contains 64 × 64 images from five classes. To ensure that our report is self-contained and can be used to reproduce our results, we have included the entire process of creating the dataset, defining the dataloader class, and the training and testing routines from HW4. We will present each step in detail below.

Creating Our Own Image Classification Dataset

In this task we need to create our own image classification dataset by taking images from the MS-COCO dataset. For that, I took the help of python version of the COCO API. We are using 2014 Train images. I downloaded all the images from the following link and uploaded it to the Google Drive. I also downloaded the annotation files: 2014 Train/Val annotations. For image classification task, instances_train2014.json file is used. The following script is used to read the images from the MS-COCO dataset, resize it to 64x64 images using PIL module and save the images to another directory. It is made sure that there are no duplicate images.

Source Code:

```
dataDir = '/content/drive/MyDrive/Arghadip/DL/Datasets'  # Directory of annotation and
entire MS-COCO dataset
dataType='train2014'  # We are working with
"train2014" version
annFile='{}/annotations/instances_{}.json'.format(dataDir,dataType)  # Name of the annotation file
coco=COCO((annFile)  # initialize COCO api for
instance annotations

classes = ['airplane', 'bus', 'cat', 'dog', 'pizza']  # We are interested in these
5 classes

for cat in classes:  # Loop over all the classes
catIds = coco.getCatIds(catNms=[cat])  # Get class ids
imgIds = coco.getImgIds(catIds=catIds)  # Get image ids corresponding
to the class

# Shuffle images
indices = list(range(len(imgIds)))
random shuffle(indices)
```

```
# Take first 1500 images for training
train_dir = os.path.join(dataDir, 'coco', 'train', cat)  # Directory to write training
images
  os.mkdir(train_dir)  # Make the directory
for i in range(1500):
    img = coco.loadImgs(imgIds[indices[i]])[0]  # Get image details
    I = Image.open(os.path.join(dataDir, dataType, img['file_name'])) # Load as PIL image from
the MS-COCO directory
    I = I.convert('RGB')  # Convert to RGB
    im_resized = I.resize((64,64), Image.BOX)  # Resize to 64 x 64
    im_resized.save(os.path.join(train_dir, cat + '_' + str(i) + '.jpg')) # Save with class name
and sequence number

# Take next 500 images for testing
test_dir = os.path.join(dataDir, 'coco', 'test', cat)  # Directory to write test
images
  os.mkdir(test_dir)  # Make the directory
for i in range(1500,2000):
    img = coco.loadImgs(imgIds[indices[i]])[0]  # Get image details
I = Image.open(os.path.join(dataDir, dataType, img['file_name'])) # Load as PIL image from
the MS-COCO directory
    I = I.convert('RGB')  # Convert to RGB
    im_resized = I.resize((64,64), Image.BOX)  # Resize to 64 x 64
    im_resized.save(os.path.join(test_dir, cat + '_' + str(i-1500) + '.jpg')) # Save with class
name and sequence number
```

Output Directory Structure:

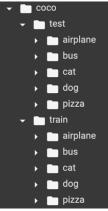


Figure 2. Dataset directory structure

Once the images are saved in proper directory, the next step is to implement our own dataset class to provide the necessary support to the torch.utils.data.DataLoader class. The source code is given below. The dataset class My_COCO_Dataset contains the relevant information such as root directory, split (train or test), number of images per class etc. It also performs proper transform to covert PIL images to CNN supported tensor input. I haven't used any data augmentation techniques, but the scope is already present in the given source code. One method named _get_filenames_and_labels() is defined to get all the filenames and corresponding labels, so that later it can be used during query (i.e. __getitem__()). The __len__() method is also overwritten to return the total number of images in the dataset. During training feeding images from different classes in a single batch leads to better training. Therefore, shuffling support is also provided (shuffle=True). The labels are integers (0 to 4) for 5 classes.

Source Code:

```
super().__init__() # Part of the definition is obtained from parent class
      self.num images per class = 1500
      self.num_images_per class = 500
    for i in range (len (self.classes)):
      for j in range(self.num images per class):
self.list_of_files_and_labels.append([os.path.join(self.path, self.classes[i],
self.classes[i] + '_' + str(j) + '.jpg'), i])
    image = Image.open(self.list of files and labels[index][0])
```

Training routine:

We use Adam optimizer with *learning rate* = 1e-3, β_1 = 0.9 and β_2 = 0.99 for training. It also saves the model parameters in every epoch as a dictionary that is later used for testing using validation dataset. The model is trained for 50 epochs.

Test routine and script to plot confusion matrix:

The test routine is inspired from the routine present in DLStudio. The seaborn package is useful for the proper display of the confusion matrix.

Main code:

```
num atten heads=8,
                embedding size=128, \
                how_many_basic_encoders=2, \
                mlp dim=96, \
batch size = 64
                            # Batch size
num_workers = 2  # Number of parallel threads for data loading
train_dataset = My_COCO_Dataset(split='train', shuffle=True)  # Training dataset
num workers = 2
num workers=num workers) # Loader
test dataloader = DataLoader(test dataset, batch size=batch size, shuffle=False,
num workers=num workers) # Loader
learning rate = 1e-3
                            # Learning rate
num_epochs = 50
                           # Number of training epochs
print_frequency=print_frequency) # Loss
fig, axes = plt.subplots(ncols=1, figsize = (8,4))
axes.plot(loss, label=label)
axes.set_title('Training loss vs. batches')
axes.set_xlabel('Batches in hundreds, batch size = ' + str(batch_size))
axes.set_ylabel('Training loss')
```

Training Loss over Training Iterations:

The following plots show how training progresses along the epochs.

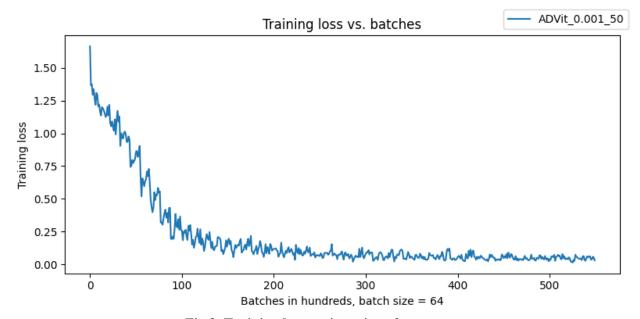


Fig.3. Training loss vs. iterations for ADVit

Confusion matrix and overall accuracy:

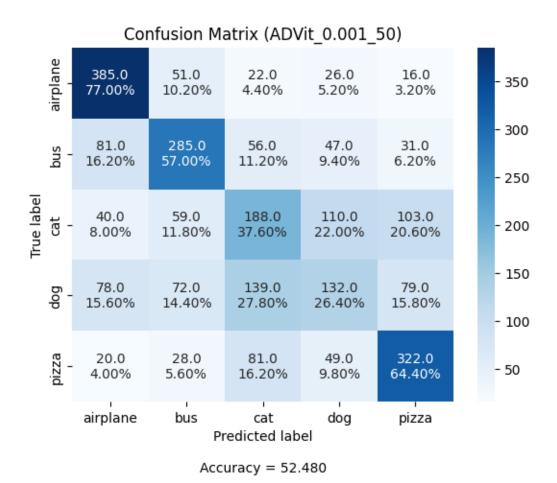


Fig.6. Confusion matrix for ADVit over HW4-COCO testset

Overall accuracy: 52.48%

Comparison with HW4 CNN based network:

For better comparison, I present the confusion matrices of 3 different CNN tasks from HW4. These networks were trained with the same hyperparameters (Adam optimizer with learning rate 1e-3, $\beta_1 = 0.9$ and $\beta_2 = 0.99$) for 50 epochs with batch size of 8.

Confusion Matrix for HW4 Task 1:

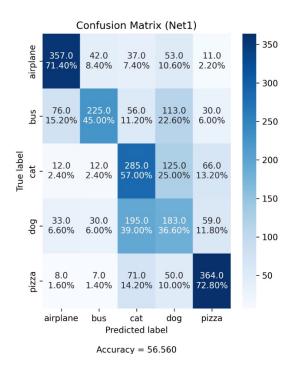


Fig.4. Confusion matrix for HW4 Net1 (Overall test accuracy = 56.56%)

Confusion Matrix for HW4 Task 2:

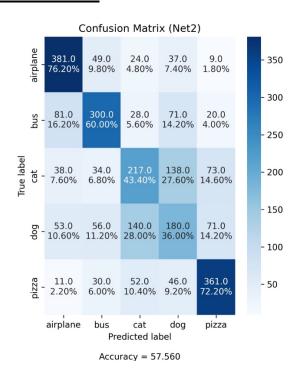


Fig.5. Confusion matrix for HW4 Net2 (Overall test accuracy = 57.56%)

Confusion Matrix for HW4 Task 3:

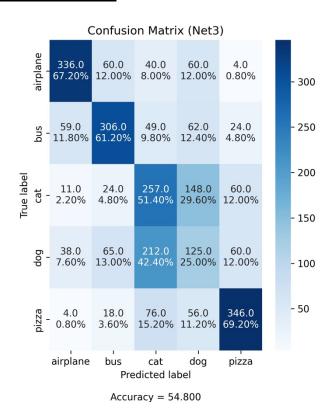


Fig.6. Confusion matrix for HW4 Net3 (Overall test accuracy = 54.80%)

The transformer model (ADVit) performs worse than the CNN-based models, with only 52.48% accuracy on the test set. In comparison, the CNNs had the following accuracy in decreasing order: Net2 (57.56%) > Net1 (56.56%) > Net3 (54.80%). *This is because ADVit is overfitting the small training set*, which consists of only 7500 images of size 64x64. The total number of training points is 7500x(3x64x64) = 92160000, while the number of parameters in the transformer model (41577829) is comparable with the size of the training set. Although I have tried to reduce the number of parameters by adjusting the hyperparameters, the model still overfits. To address this issue, regularization, dropout, and early termination methods can be used, but they are not within the scope of this homework. On the other hand, the low accuracies of the CNN models in HW4 were due to underfitting caused by their small model orders.

3.3. Extra Credit: Multi-headed self-attention using torch. einsum

The implementation code with detailed comments is given below.

```
super().__init__()
self.max_seq_length = max_seq_length
        self.num_atten_heads = num_atten_heads
       self.qkv size = self.embedding size // num atten heads
        self.attention heads arr = nn.ModuleList([AttentionHead einsum(self.max seq length,
range(num atten heads)]) # Forming multiple attention heads
       concat out from atten heads = torch.zeros(sentence tensor.shape[0], self.max seq length,
               self.attention heads arr[i](sentence tensor portion) # Concatenate attention
   def __init__(self, max_seq_length, qkv_size):
    super().__init__()
    self.qkv_size = qkv_size
       self.max_seq_length = max_seq_length
       self.WQ = nn.Linear(max_seq_length * self.qkv size,
                           max seq length * self.qkv size)
       self.WK = nn.Linear(max_seq_length * self.qkv_size,
       max_seq_length * self.qkv_size)
self.WV = nn.Linear(max_seq_length * self.qkv_size,
                           max_seq_length * self.qkv size)
           sentence_portion.shape[0], -1).float()).to(device) # Query (1, 272)
       K = self.WK(sentence portion.reshape(
           sentence_portion.shape[0], -1).float()).to(device) # Key (1, 272)
       V = self.WV(sentence_portion.reshape(
           sentence portion.shape[0], -1).float()).to(device) # Value (1, 272)
       Q = Q.view(sentence_portion.shape[0],
       K = K.view(sentence_portion.shape[0],
       V = V.view(sentence_portion.shape[0],
        # This line creates a tensor A with the shape (1, 16, 16) by performing a batch matrix
multiplication
       # between the key tensor K (shape: (1, 17, 16)) and a tensor of ones with the same shape
as the query
        # tensor Q (shape: (1, 17, 16)).
         = torch.einsum("bnm,bnk->bmk", K, torch.ones_like(Q)) # (1, 16, 16)
```

```
This line performs a batch matrix multiplication between the query tensor Q (shape: (1,
        # and the attention weight tensor A (shape: (1, 16, 16)) to obtain a dot product tensor
of shape
       # (1, 17, 16). This tensor represents the pairwise dot product similarity scores between
each query
        # vector and key vector for each position in the input sentence.
       QK_dot_prod = torch.einsum("bnm,bmk->bnk", Q, A)
                                                                    # (1, 17, 16)
        # This line applies the softmax function along the second dimension of the dot product
tensor.
        # resulting in a row-wise normalized tensor with shape (1, 17, 16). Each row in the
        # represents the attention distribution for a given query vector, i.e. how much attention
to pay to
       rowwise softmax normalizations = self.softmax(QK_dot_prod) # (1, 17, 16)
        # This line performs a batch matrix multiplication between the attention weight tensor
       \# (shape: (1, 16, 16)) and the value tensor V (shape: (1, 17, 16)) to obtain the output
       # tensor Z with shape (1, 17, 16). This tensor represents the final output of the
attention head.
       Z = torch.einsum("bnm,bnk->bnk", rowwise_softmax_normalizations, V)
       # This line computes a scalar coefficient to normalize the output of the attention head
by dividing
       # by the square root of the dimension of the key, query and value vectors.
       coeff = 1.0/torch.sqrt(torch.tensor([self.qkv size]).float()).to(device) # (1)
       # This line normalizes the output tensor Z by scaling it with the coefficient.
       Z = coeff * Z
       return Z
```

The shape of each tensor is given in the comments assuming the following hyperparameter values.

- max seq length = 17
- embedding size = 128
- num atten heads = 8

The above code implements a multi-headed attention mechanism for a given input sentence portion. It does this by matrix-multiplying the embedding vector for each word in the sentence by the WQ, WK, and WV matrices to produce the query vector Q, key vector K, and value vector V for each word in the input sentence. The dot product of Q and K is then used to calculate the attention weights, which are applied to the value vectors to get the final attention output.

Suggested improvements:

To improve the performance of a transformer on a small training set, several techniques can be employed. Firstly, transfer learning can be used by pretraining the transformer on a larger dataset and then fine-tuning it on the small dataset. This can lead to better generalization and faster convergence on the small dataset. Another technique is to use data augmentation to artificially increase the size of the training set. This can be done by applying random transformations to the input data such as cropping, flipping, or adding noise. This can help the model learn more robust features and reduce overfitting. Regularization techniques such as dropout and weight decay can also be used to prevent overfitting on the small training set. Dropout randomly drops out some neurons during training, while weight decay adds a penalty to the loss function for large weights. Finally, ensembling multiple models can also help improve performance on a small dataset. By training multiple models with different initializations or architectures and combining their predictions, the overall performance can be boosted. However, this approach may require more computational resources.

4. Lessons Learned

In this Vision transformer homework, we learned how to implement a transformer model for image classification. We used the PyTorch library to build a transformer model that consists of an encoder, where the encoder extracts features from an image. We also learned about the importance of attention mechanisms in transformer models, which allow the model to focus on important parts of the input. Additionally, we experimented with different hyperparameters and techniques to improve the performance of our model. Overall, this homework provided a practical introduction to transformer models and their application in computer vision tasks. Through experimentation and fine-tuning, we gained insights into how different components and hyperparameters of the model affect its performance. These lessons can be applied to other transformer-based models and tasks and can help us develop more accurate and efficient deep learning models.