# Part 1: Introduction

In this lab, we are asked to implement a Vision Transformer (ViT) to classify images in a 5 class dataset. We reuse the data loading code that was created in Lab 4. Furthermore, we reuse the image dataset that was created in Lab 4. After implementing the ViT, we train for 10 epochs and then provide a confusion matrix of the resulting model on the test dataset. We compare the results from the ViT with that of the CNN based model from Lab 4.

# Part 2: Transformer Implementation

**Transformer Implementation - Embedding Generation:**

In order to implement the ViT on a dataset of 64 x 64 x 3 (H x W x C) images, we need to represent each "image patch" as an embedding. An image patch is a square region within the image. In this lab, 16 image patches are generated by applying a 3 channel to 100 channel convolution using a stride of 16 and kernel size of 16. The 100 channels are important because that will represent the size of the embedding for a particular patch. This will use a 3 x 100 x 16 x 16 kernel. Since the kernel is 16 x 16 with a stride of 16, there are 16 pixels in each resulting channel of the convolution. Then, a stack of pixels along the channel dimension is taken, resulting in a 1 x 100 embedding vector for each patch. In total, this process creates 16 distinct 1 x 100 embedding vectors. Put into a tensor together, that is a 16 x 100 tensor. This process is achieved via the following code:

```python
class PatchEmbed(nn.Module): #Note, from Meta DINO ViT code ()
    """ Image to Patch Embedding
    """
    def __init__(self, img_size=64, patch_size=16, in_chans=3, embed_dim=100):
        super().__init__()
        num_patches = (img_size // patch_size) * (img_size // patch_size)
        self.img_size = img_size
        self.patch_size = patch_size
        self.num_patches = num_patches

        self.proj = nn.Conv2d(in_chans, embed_dim, kernel_size=patch_size, stride=patch_size)

    def forward(self, x):
        B, C, H, W = x.shape
        y = self.proj(x)
        #print('conv output before flattening and transpose', y.shape)
        x = self.proj(x).flatten(2).transpose(1, 2)
        return x
```

**Transformer Implementation - Adding Positional and Class Embeddings**

The ViT paper specifies to add a learnable class embedding as a row on top of the 16 x 100 tensor. This class embedding is created via the nn.Parameter command using an input as a

random vector of size 1 x 100. Furthermore, a learnable position embedding is added to the now 17 x 100 tensor using nn.Parameter (but an input as a random tensor of the size 17 x 100).

Code for this is below:

```
3 class MasterEncoder(nn.Module):
4     def __init__(self, max_seq_length, embedding_size, how_many_basic_encoders, num_atten_heads):
5         super().__init__()
6         self.class_embedding = nn.Parameter(torch.rand(size = (1, 1, embedding_size))) #create class embedding.
7         #this class_embedding will be the first row of the embedding matrix, where axis 0 is each patch embedding
8         #Note that the first 1 in the size parameter above corresponds to the batch size.
9         self.pos_embedding = nn.Parameter(torch.rand(size = (1, max_seq_length, embedding_size)))
10        self.patch_generator = PatchEmbed()
11        self.max_seq_length = max_seq_length
12        self.basic_encoder_arr = nn.ModuleList([BasicEncoder(
13            max_seq_length, embedding_size, num_atten_heads) for _ in range(how_many_basic_encoders)])  # (A)
14        self.mlp_head = nn.Linear(embedding_size, 5)
15
```

**Transformer Implementation - Query, Key, Value logic with multiple attention heads**

Now, we enrich these embeddings using multiple head Query, Key, and Value logic specified in the original transformer paper. The learnable Query, Key, and Value matrices are embedding_size x embedding_size dimension (WQ, WK, QV). Specifically, Q = embedding · WQ, K = embedding · WK, and V = embedding · WV. The resulting Q, K, and V tensors are then used to product a probability mass function interpretation of embeddings by applying a softmax to Q · K.T and dividing the result by the length of the square root of the embedding. That result is then multiplied by V.

After the attention blocks are processed and concatenated (because of the multiple heads), the output is subject to a Layer Normalization and then followed by a fully connected network that preserves the shape of the input and another Layer Normalization. A series of four consecutive attention blocks, layer normalizations, fully connected layers, and layer normalizations comprise the entire encoder network.

Note, most of the code for this was taken from ViTHelper.py which was provided in the lab documents.

**Transformer Implementation - MLP Class Prediction**

After the entire encoder network is processed, the resulting class embedding is transformed from a 100 dimensional vector to a 5 dimensional vector representing the class probabilities. This is achieved using the nn.Linear() function. During training, the resulting 5 dimensional vector is used in the Cross Entropy loss computation in order to optimize the network parameters. For testing, the maximum argument of the resulting 5 dimensional vector is used as the class prediction from the network.
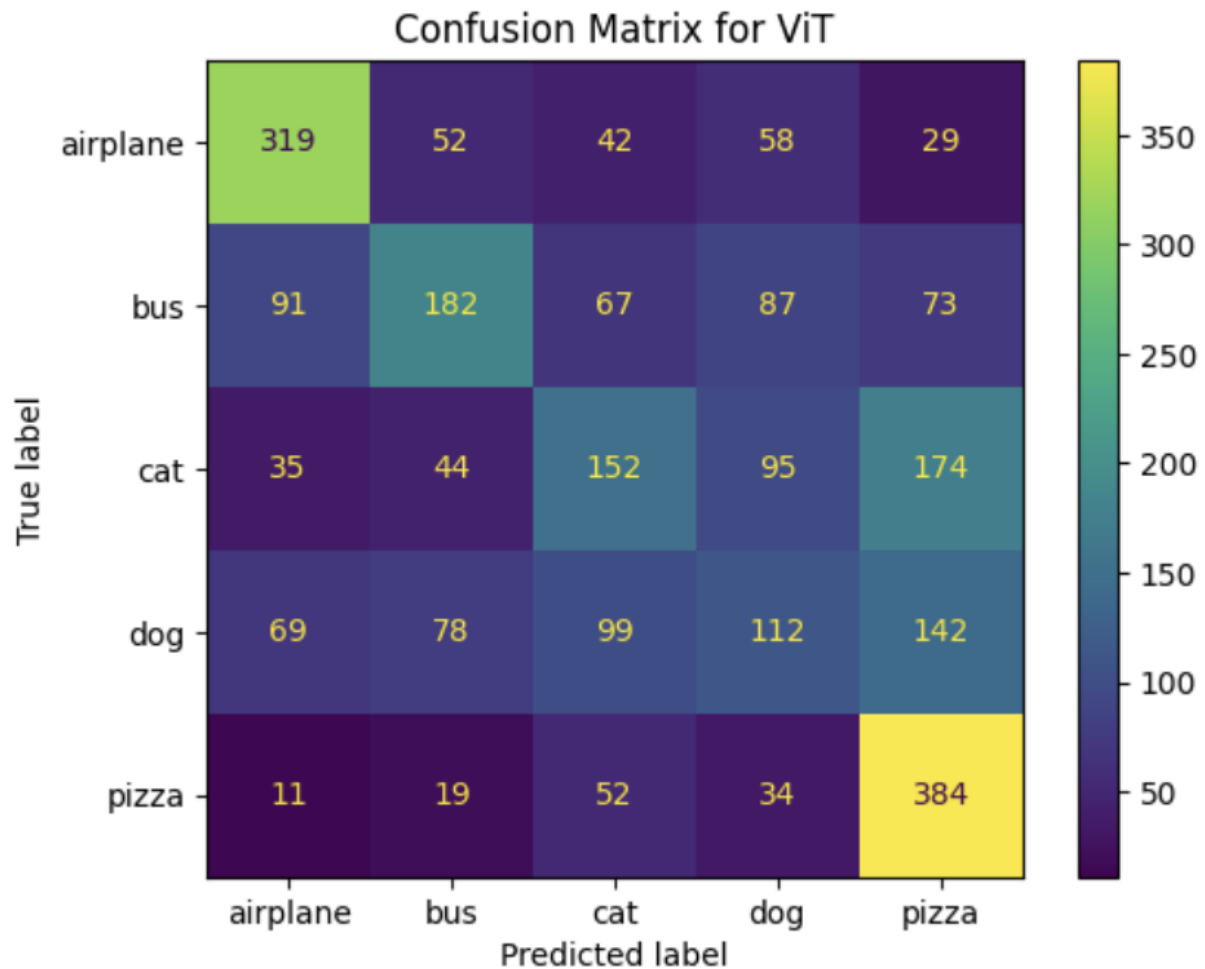
Vision Transformers

Note that unlike sequence to sequence predictions, the ViT only has an encoder network. The decoder is not needed because we are transforming the original data to the same shape - we are just doing a classification problem.

# Part 3: Results

**Results on ViT**

After training the network, the following confusion matrix is obtained:
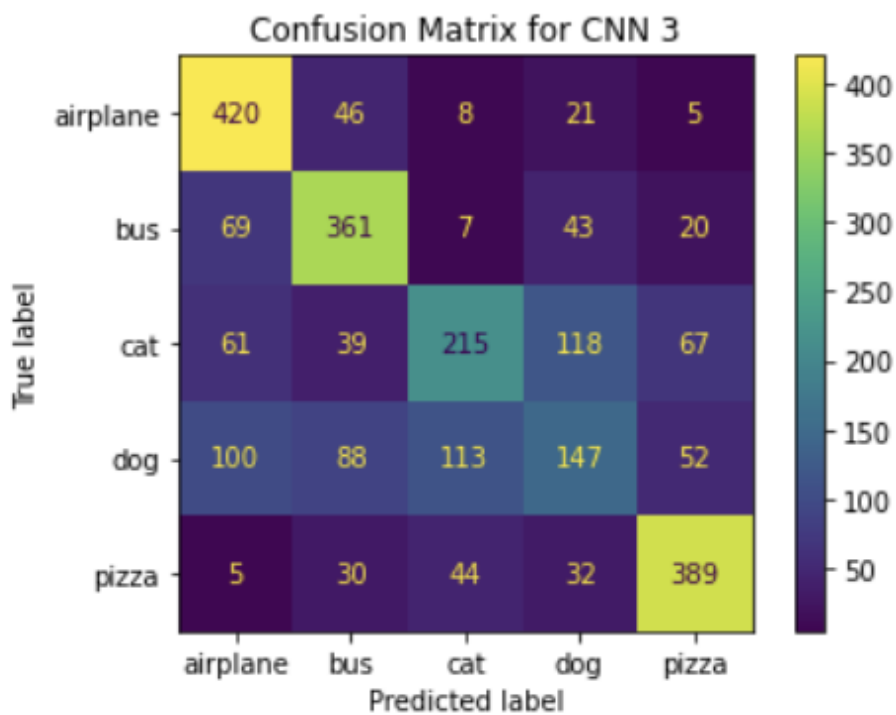


This is slightly worse performance than the CNN based network we did in Lab 4 (45% overall accuracy vs 61% overall accuracy). The reason is likely due to the large increase in parameter size while keeping the training data the same size, which is resulting in overfitting the dataset and reducing generalizability during test time.

**Results on CNN from Lab 4:**

Below is the Confusion Matrix for the best performing CNN in Lab 4:

Accuracy of the network on the val images: 61 %

Confusion Matrix for CNN 3

|            | airplane | bus | cat | dog | pizza |
|------------|----------|-----|-----|-----|-------|
| airplane   | 420      | 46  | 8   | 21  | 5     |
| bus        | 69       | 361 | 7   | 43  | 20    |
| cat        | 61       | 39  | 215 | 118 | 67    |
| dog        | 100      | 88  | 113 | 147 | 52    |
| pizza      | 5        | 30  | 44  | 32  | 389   |

True label / Predicted label

# Part 4: Source Code
**Source code - Master Encoder Network (mostly from ViTHelper.py)**

## Vision Transformers

```python
3  class MasterEncoder(nn.Module):
4      def __init__(self, max_seq_length, embedding_size, how_many_basic_encoders, num_atten_heads):
5          super().__init__()
6          self.class_embedding = nn.Parameter(torch.rand(size = (1, 1, embedding_size))) #create class embedding.
7          #this class_embedding will be the first row of the embedding matrix, where axis 0 is each patch embedding
8          #Note that the first 1 in the size parameter above corresponds to the batch size.
9          self.pos_embedding = nn.Parameter(torch.rand(size = (1, max_seq_length, embedding_size)))
10         self.patch_generator = PatchEmbed()
11         self.max_seq_length = max_seq_length
12         self.basic_encoder_arr = nn.ModuleList([BasicEncoder(
13             max_seq_length, embedding_size, num_atten_heads) for _ in range(how_many_basic_encoders)])  # (A)
14         self.mlp_head = nn.Linear(embedding_size, 5)
15
16     def forward(self, img_patch):
17         #out_tensor = sentence_tensor
18         img_embedding = self.patch_generator(img_patch) #pass in img_patch, which has convo2d applied to it and
19         #results in the img embedding.
20
21         img_embedding = torch.cat(tensors = (self.class_embedding, img_embedding), dim = 1)
22         img_embedding = img_embedding + self.pos_embedding
23         for i in range(len(self.basic_encoder_arr)):  # (B)
24             img_embedding = self.basic_encoder_arr[i](img_embedding)
25
26
27         img_embedding = self.mlp_head(img_embedding[0,0,:]) #use the class embedding vector to do the prediction.
28         #here, we take the 100-dimension embedding to 5d, so we can get log probabilities and then do NLLLoss
29         # all using CE loss in the training loop.
30         return img_embedding
31
32
```

```python
33  class BasicEncoder(nn.Module):
34      def __init__(self, max_seq_length, embedding_size, num_atten_heads):
35          super().__init__()
36          self.max_seq_length = max_seq_length
37          self.embedding_size = embedding_size
38          self.qkv_size = self.embedding_size // num_atten_heads
39          self.num_atten_heads = num_atten_heads
40          self.self_attention_layer = SelfAttention(
41              max_seq_length, embedding_size, num_atten_heads)  # (A)
42          self.norm1 = nn.LayerNorm(self.embedding_size)  # (C)
43          self.W1 = nn.Linear(self.max_seq_length * self.embedding_size,
44                              self.max_seq_length * 2 * self.embedding_size)
45          self.W2 = nn.Linear(self.max_seq_length * 2 * self.embedding_size,
46                              self.max_seq_length * self.embedding_size)
47          self.norm2 = nn.LayerNorm(self.embedding_size)  # (E)
48
49      def forward(self, sentence_tensor):
50          input_for_self_atten = sentence_tensor.float()
51          normed_input_self_atten = self.norm1(input_for_self_atten)
52          output_self_atten = self.self_attention_layer(
53              normed_input_self_atten).to(device)  # (F)
54          input_for_FFN = output_self_atten + input_for_self_atten
55          normed_input_FFN = self.norm2(input_for_FFN)  # (I)
56          basic_encoder_out = nn.ReLU()(
57              self.W1(normed_input_FFN.view(sentence_tensor.shape[0], -1)))  # (K)
58          basic_encoder_out = self.W2(basic_encoder_out)  # (L)
59          basic_encoder_out = basic_encoder_out.view(
60              sentence_tensor.shape[0], self.max_seq_length, self.embedding_size)
61          basic_encoder_out = basic_encoder_out + input_for_FFN
62          return basic_encoder_out
63
```

Vision Transformers

## Attention Logic - from ViTHelper.py

```python
68  class SelfAttention(nn.Module):
69      def __init__(self, max_seq_length, embedding_size, num_atten_heads):
70          super().__init__()
71          self.max_seq_length = max_seq_length
72          self.embedding_size = embedding_size
73          self.num_atten_heads = num_atten_heads
74          self.qkv_size = self.embedding_size // num_atten_heads
75          self.attention_heads_arr = nn.ModuleList([AttentionHead(self.max_seq_length,
76                                              self.qkv_size) for _ in range(num_atten_heads)])  # (A)
77
78      def forward(self, sentence_tensor):  # (B)
79          concat_out_from_atten_heads = torch.zeros(sentence_tensor.shape[0], self.max_seq_length,
80                                      self.num_atten_heads * self.qkv_size).float()
81          for i in range(self.num_atten_heads):  # (C)
82              sentence_tensor_portion = sentence_tensor[:,
83                                          :, i * self.qkv_size: (i+1) * self.qkv_size]
84              concat_out_from_atten_heads[:, :, i * self.qkv_size: (i+1) * self.qkv_size] =          \
85                  self.attention_heads_arr[i](sentence_tensor_portion)  # (D)
86          return concat_out_from_atten_heads
87
88
89
90
91  class AttentionHead(nn.Module):
92      def __init__(self, max_seq_length, qkv_size):
93          super().__init__()
94          self.qkv_size = qkv_size
95          self.max_seq_length = max_seq_length
96          self.WQ = nn.Linear(max_seq_length * self.qkv_size,
97                              max_seq_length * self.qkv_size)  # (B)
98          self.WK = nn.Linear(max_seq_length * self.qkv_size,
99                              max_seq_length * self.qkv_size)  # (C)
100         self.WV = nn.Linear(max_seq_length * self.qkv_size,
101                             max_seq_length * self.qkv_size)  # (D)
102         self.softmax = nn.Softmax(dim=1)  # (E)
103
104     def forward(self, sentence_portion):  # (F)
105         Q = self.WQ(sentence_portion.reshape(
106             sentence_portion.shape[0], -1).float()).to(device)  # (G)
107         K = self.WK(sentence_portion.reshape(
108             sentence_portion.shape[0], -1).float()).to(device)  # (H)
109         V = self.WV(sentence_portion.reshape(
110             sentence_portion.shape[0], -1).float()).to(device)  # (I)
111         Q = Q.view(sentence_portion.shape[0],
112                 self.max_seq_length, self.qkv_size)  # (J)
113         K = K.view(sentence_portion.shape[0],
114                 self.max_seq_length, self.qkv_size)  # (K)
115         V = V.view(sentence_portion.shape[0],
116                 self.max_seq_length, self.qkv_size)  # (L)
117         A = K.transpose(2, 1)  # (M)
118         QK_dot_prod = Q @ A  # (N)
119         rowwise_softmax_normalizations = self.softmax(QK_dot_prod)  # (O)
120         Z = rowwise_softmax_normalizations @ V
121         coeff = 1.0/torch.sqrt(torch.tensor([self.qkv_size]).float()).to(device)  # (S)
122         Z = coeff * Z  # (T)
123         return Z
```

## Dataset Class (which is later wrapped in the Dataloader class):

## Vision Transformers

```python
1  ### Create data_loader ###
2  root_train = 'train/'
3  root_val = 'val/'
4  catNms=['airplane','bus','cat', 'dog', 'pizza']
5
6  class MyDataset(torch.utils.data.Dataset):
7      def __init__(self, root, catNms):
8          super(MyDataset).__init__()
9          self.root = {} #dictionary for main directory which holds all the images of a category
10         self.filenames = {} #dictionary for filenames of a given category
11         for cat in catNms:
12             self.root[cat] = root + cat + '/'
13         for cat in catNms:
14         #create list of image files in each category that can be opened by __getitem__
15             self.filenames[cat] = os.listdir(self.root[cat])
16
17         self.rand_max = len(os.listdir(self.root[catNms[0]])) - 1 #number of files in directory
18
19         self.mapping = {0 : 'airplane',
20                         1: 'bus',
21                         2: 'cat',
22                         3: 'dog',
23                         4: 'pizza'} #makes it easy to convert between index and name of a category.
24
25         self.one_hot_encoding = {0: torch.tensor(np.array([1, 0, 0, 0, 0])),
26                         1: torch.tensor(np.array([0, 1, 0, 0, 0])),
27                         2: torch.tensor(np.array([0, 0, 1, 0, 0])),
28                         3: torch.tensor(np.array([0, 0, 0, 1, 0])),
29                         4: torch.tensor(np.array([0, 0, 0, 0, 1]))} #one hot encode each category.
30
31
32         self.to_Tensor_and_Norm = tvt.Compose([tvt.ToTensor(),tvt.Resize((64,64)) ,
33                             tvt.Normalize([0], [1]) ,  tvt.ColorJitter(0.75, 0.75) ,
34                             tvt.RandomHorizontalFlip( p = 0.75),
35                             tvt.RandomRotation(degrees = 45)]) #normalize and resize in case the resize op
36  #       wasn't done. Note that resizing here may not have any impact as the resizing was done previously.
37
38
39      def __len__(self):
40          count = 0
41          for cat in catNms:
42              temp_num = os.listdir(self.root[cat])
43              count = count + len(temp_num)
44          return count #. Will be 2500 if the root=val/ and 7500 if root=train/
45
```

## Vision Transformers

```
39      def __len__(self):
40          count = 0
41          for cat in catNms:
42              temp_num = os.listdir(self.root[cat])
43              count = count + len(temp_num)
44          return count #. Will be 2500 if the root=val/ and 7500 if root=train/
45
46      def __getitem__(self, index):
47          file_index = index % self.rand_max + 1
48          class_index = index % 5
49
50          img_file = self.filenames[self.mapping[class_index]]
51
52          try:
53              item = Image.open(self.root[self.mapping[class_index]] + img_file[file_index])
54          except IndexError: #for debugging
55              print('these are the indices for the line above when shape is correct', class_index , file_index)
56
57          np_img = np.array(item)
58          shape = np_img.shape
59          while shape != (64, 64 ,3): #handle if the image from COCO is grayscale.
60              #print('found a grayscale image, fetching an RGB!')
61              another_rand = random.randint(0,self.rand_max)  #generate another rand num
62              #print('another_rand is', another_rand)
63              try:
64                  item = Image.open(self.root[self.mapping[class_index]] + img_file[another_rand])
65              except IndexError: #for debugging
66                  print('these are the indices for the line above when shape is incorrect', another_rand , class_index)
67              np_img = np.array(item)
68              shape = np_img.shape
69
70          img = self.to_Tensor_and_Norm(item)
71          class_label = self.one_hot_encoding[class_index].type(torch.FloatTensor) #convert to Float
72          return img, class_label
73
```

## Training Logic:

```
1  epochs = 10
2  criterion = nn.CrossEntropyLoss()
3  encoder = MasterEncoder(max_seq_length=17, embedding_size=100, how_many_basic_encoders=4, num_atten_heads=4)
4  optimizer = torch.optim.Adam(encoder.parameters(), lr = 1e-3, betas = (0.9, 0.99))
5  loss_running_list_encoder = []
6  running_loss = 0.0
7  encoder = encoder.to(device)
8  for i in range(epochs):
9      for n, data in enumerate(my_train_dataloader):
10          #Create encoder network:
11          #print(n)
12          optimizer.zero_grad() #Sets gradients of all model parameters to zero. We want to compute fresh gradients
13          #based on the new forward run.
14          img, GT = data
15          GT = torch.argmax(GT)
16
17          img = img.to(device)
18          GT = GT.to(device)
19
20          out = encoder(img)
21          loss = criterion(out, GT) #input, then target for arg order
22
23          loss.backward() #compute derivative of loss wrt each gradient.
24          optimizer.step() #takes a step on hyperplane based on derivatives
25          running_loss += loss.item()
26          if (n+1) % 500 == 0:
27              print("[epoch: %d, batch: %5d] loss: %3f" % (i + 1, n + 1, running_loss / 500))
28              loss_running_list_encoder.append(running_loss/500)
29              running_loss = 0.0
30
31
```

Vision Transformers

## Testing Logic and Confusion Matrix generation:

```python
2 correct = 0
3 total = 0
4 y_pred = []
5 y_label = []
6 mapping = { 0: 'airplane',
7             1: 'bus',
8             2: 'cat',
9             3: 'dog',
10            4: 'pizza'}
11
12
13 with torch.no_grad():
14     for n, data in enumerate(my_val_dataloader):
15         images, labels = data
16         images = images.to(device)
17         outputs = encoder(images)
18
19         predicted = torch.max(outputs.data)
20         predicted_class = torch.argmax(outputs.data)
21
22         total += labels.size(0) #add to total count of ground truth images so we can calculate total accuracy
23         #print("total images in val set", total)
24         for n, i in enumerate(labels):
25             temp = np.array(i) #arrays are one hot encoded, we need to convert it into a human readable label for
26             #display in the confusion matrix
27             label_arg = np.argmax(temp) #get the argument of the one hot encoding
28             y_label.append(mapping[label_arg]) #apply the argument to the mapping dictionary above. For example
29             # if the argument is 3, then, that corresponds to a label of dog in the mapping dictionary
30             t = int(predicted_class) #get integer representation of prediction from network (will
31             #be an int from 0 to 4.
32             y_pred.append(mapping[t]) #append the predicted output of this label to the prediction list, but,
33             #via the mapping dictionary definition so that the y_pred list is human readable.
34
35             if label_arg == t:
36                 correct = correct + 1 #add to total count of correct predictions so we can calculate total accuracy
37
38
39 print('Accuracy of the network on the val images: %d %%' % (
40     100 * correct / total))
41 from sklearn.metrics import confusion_matrix
42
43 y_true = y_label
44 y_pred = y_pred
45 confusion_matrix=confusion_matrix(y_true, y_pred, labels = [ "airplane", "bus", "cat", "dog", "pizza"])
46 disp = ConfusionMatrixDisplay(confusion_matrix, display_labels = [ "airplane", "bus", "cat", "dog", "pizza"])
47 disp.plot()
48 disp.ax_.set_title("Confusion Matrix for ViT")
49 plt.show()
```