

16824 Homework 3

Jiajun Wan
Andrew ID: jiajunw2

Task 1: Understanding VQA

- 1.1 Why do you think that it's computationally costly to use all answers as separate classes instead of keeping only the most frequent ones?

We have to use a large number of output neurons, that add computations on final sigmoid/softmax calculation, loss computation, and backward propagation. Keeping only the most frequent ones will largely reduce the amount of calculations.

- 1.2 Complete the `__init__` method of the dataset class for VQA in `vqa_dataset.py`. Specifically, initialize the VQA API and anything you need from that.

```
1 self._vqa = VQA(annotation_file=annotation_json_file_path,  
2                  question_file=question_json_file_path)
```

- 1.3 Implement the `__len__` method of the VQADataset class. Should the size of the dataset be equal to the number of images, questions or the answers?

```
1 def __len__(self):  
2     # total number of questions  
3     return len(self._vqa.qqa)
```

The size of the dataset be equal to the number of questions.

- 1.4 Complete the `__getitem__` method

```
1 q_anno = self._vqa.qa[idx]  # load annotation  
2 q_str = self._vqa.qqa[idx]['question']  # question in str format
```

Task 2: Building a pipeline for VQA

- 2.1 What should be the output dimension of the trainable linear layer?
Complete the corresponding part in the `__init__` method of BaselineNet in `models.py`.

```
1 self.classifier = nn.Linear(
2     self.text_encoder.config.hidden_size + 512,
3     n_answers,
4 )
```

It should be `n_answers` (5127)

- 2.2 Implement `compute_vis_feats` that featurizes the image (it can be implemented as an one-liner!).

```
1 def compute_vis_feats(self, image):
2     """Convert image tensors to feature tensors."""
3     return torch.nn.functional.adaptive_avg_pool2d(self.vis_encoder(image), 1).
4        squeeze() # feed to vis_encoder and then mean pool on spatial dims
```

- 2.3 Implement the forward pass of BaselineNet. Make sure to use `compute_vis_feats` and `compute_text_feats`.

```
1 def forward(self, image, question):
2     """Forward pass, image (B, 3, 224, 224), qs list of str."""
3     text_feats = self.compute_text_feats(question)
4     vis_feats = self.compute_vis_feats(image)
5     concat_feats = torch.hstack((text_feats, vis_feats))
6     logits = self.classifier(concat_feats)
7     return logits
```

- 2.4 What is the loss for multi-label classification? (Hint: in HW1 you also tackled a multi-class classification problem)

Binary Cross-Entropy Loss

- 2.5 Implement the loss call and the optimization code in the `train_test_loop` method.

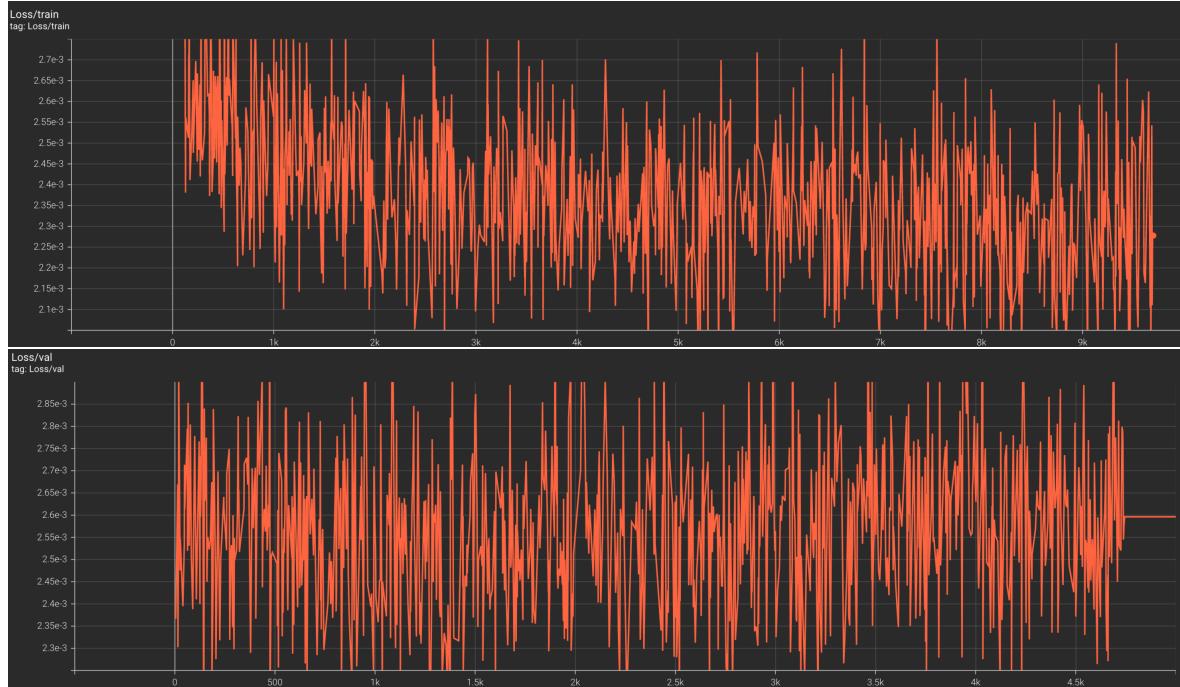
```
1 loss = F.binary_cross_entropy_with_logits(scores, answers)
2
3 # Update
4 if mode == 'train':
5     # optimize loss
6     loss.backward()
7     self.optimizer.step()
8     self.optimizer.zero_grad()
```

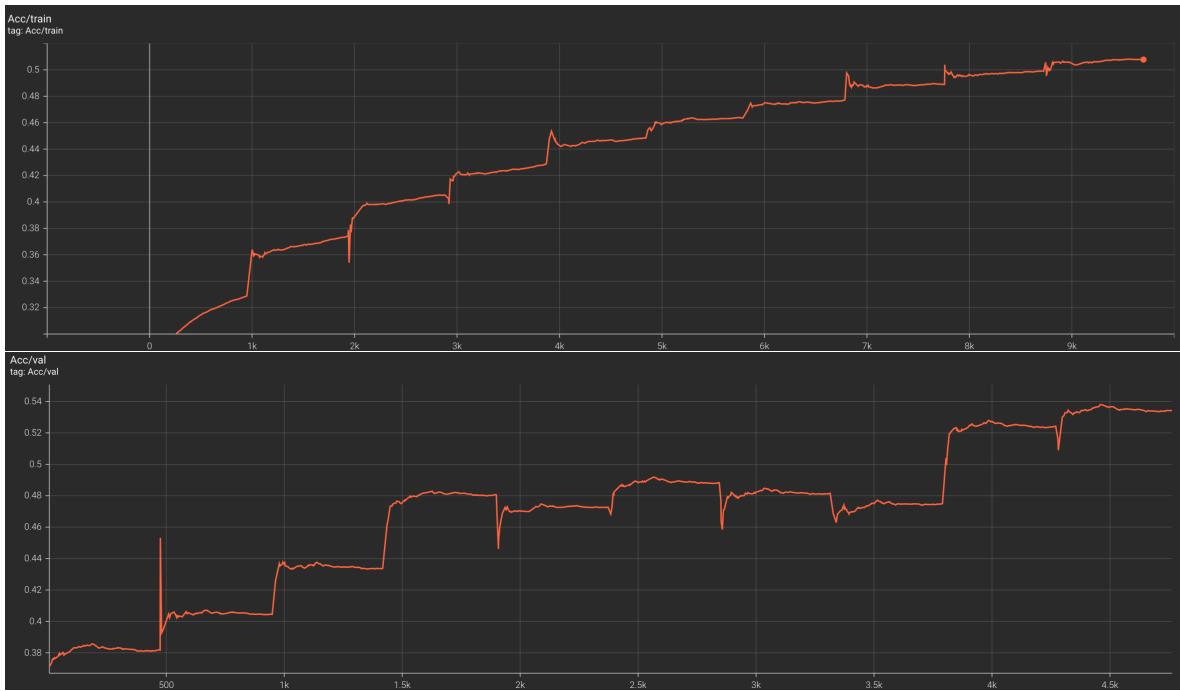
2.6 Complete the `train_test_loop` method to monitor the performance in Tensorboard. Plot the loss and accuracy and include these in your report. Additionally show multiple image-question pairs (at least 3) with the respective answers (predicted and ground-truth).

```

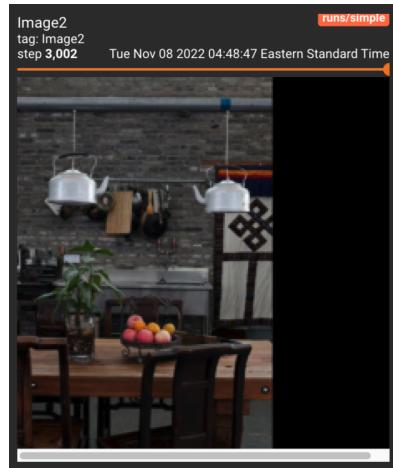
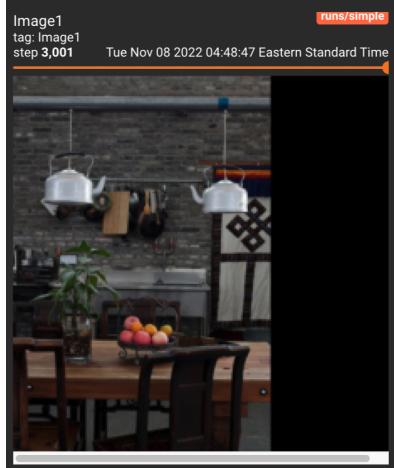
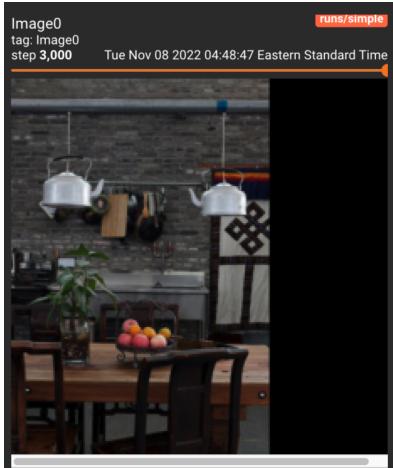
1 # add code to show the question
2 self.writer.add_text(
3     'Question%d' % i, data['question'][i],
4     epoch * _n_show + i
5 )
6 # the gt answer
7 self.writer.add_text(
8     'GT Answer%d' % i, self._id2answer[(data['answers'][i] == 1).nonzero(
9         as_tuple=True)[0][0].item()],
10    epoch * _n_show + i
11 )
12 # and the predicted answer
13 self.writer.add_text(
14     'Predicted Answer%d' % i, self._id2answer[scores.argmax(1)[i].item()],
15     epoch * _n_show + i
16 )
17 # add code to plot the current accuracy
18 self.writer.add_scalar(
19     'Acc/' + mode, n_correct / n_samples,
20     epoch * len(self.data_loaders[mode]) + step
21 )

```





Final Val Acc: 0.5341694647442228



GT Answer0

GT Answer0/text_summary
tag: GT Answer0/text_summary

step 3,000

wood

GT Answer1

GT Answer1/text_summary
tag: GT Answer1/text_summary

step 3,001

yes

GT Answer2

GT Answer2/text_summary
tag: GT Answer2/text_summary

step 3,002

kettle

Predicted Answer0

Predicted Answer0/text_summary
tag: Predicted Answer0/text_summary

step 3,000

Other

Predicted Answer1

Predicted Answer1/text_summary
tag: Predicted Answer1/text_summary

step 3,001

yes

Predicted Answer2

Predicted Answer2/text_summary
tag: Predicted Answer2/text_summary

step 3,002

5
Other

Question0

Question0/text_summary
tag: Question0/text_summary

step 3,000

What is the table made of?

Question1

Question1/text_summary
tag: Question1/text_summary

step 3,001

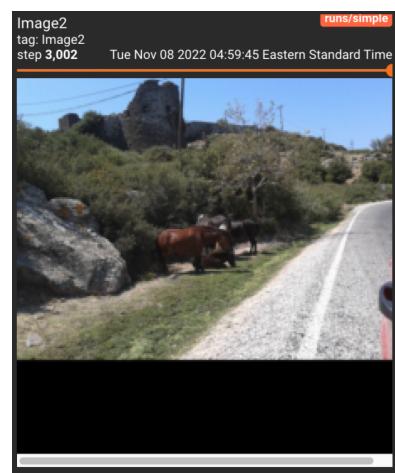
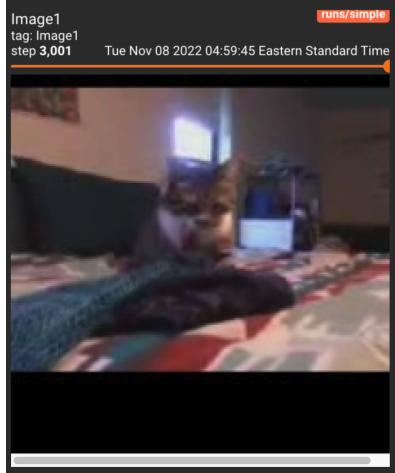
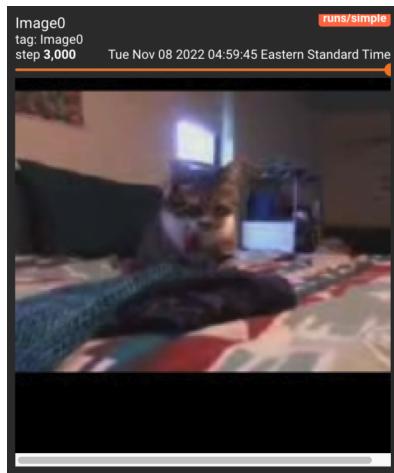
Is the food napping on the table?

Question2

Question2/text_summary
tag: Question2/text_summary

step 3,002

What has been upcycled to make lights?



GT Answer0

GT Answer0/text_summary
tag: GT Answer0/text_summary

step 3,000

camera

GT Answer1

GT Answer1/text_summary
tag: GT Answer1/text_summary

step 3,001

no

GT Answer2

GT Answer2/text_summary
tag: GT Answer2/text_summary

step 3,002

3

Predicted Answer0

Predicted Answer0/text_summary
tag: Predicted Answer0/text_summary

step 3,000

yes

Predicted Answer1

Predicted Answer1/text_summary
tag: Predicted Answer1/text_summary

step 3,001

yes

Predicted Answer2

Predicted Answer2/text_summary
tag: Predicted Answer2/text_summary

step 3,002

2

6

Question0

Question0/text_summary
tag: Question0/text_summary

step 3,000

What is causing the picture to be blurry?

Question1

Question1/text_summary
tag: Question1/text_summary

step 3,001

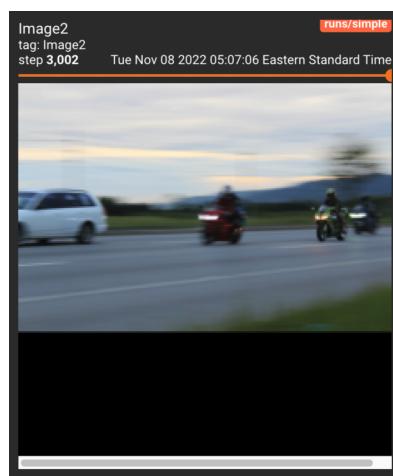
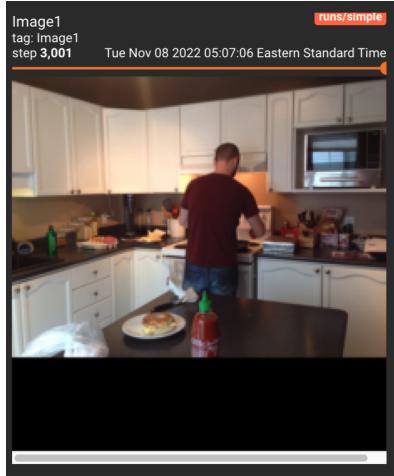
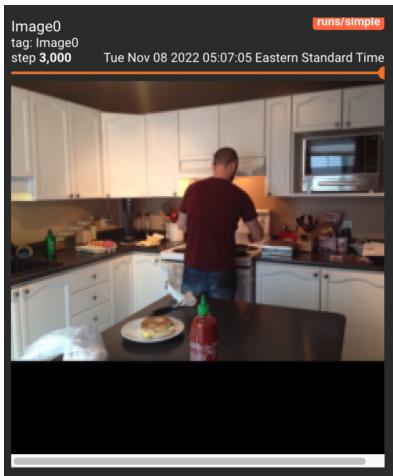
Is it night time?

Question2

Question2/text_summary
tag: Question2/text_summary

step 3,002

How many horses are shown?



GT Answer0

GT Answer0/text_summary
tag: GT Answer0/text_summary

step 3,000

ketchup

GT Answer1 GT Answer0

GT Answer1/text_summary
tag: GT Answer1/text_summary

step 3,001

yes

GT Answer2

GT Answer2/text_summary
tag: GT Answer2/text_summary

step 3,002

yes

Predicted Answer0

Predicted Answer0/text_summary
tag: Predicted Answer0/text_summary

step 3,000

kitchen

Predicted Answer1

Predicted Answer1/text_summary
tag: Predicted Answer1/text_summary

step 3,001

kitchen

Predicted Answer2

Predicted Answer2/text_summary
tag: Predicted Answer2/text_summary

step 3,002 7

2

Question0

Question0/text_summary
tag: Question0/text_summary

step 3,000

What is in the bottle on the kitchen Isle?

Question1

Question1/text_summary
tag: Question1/text_summary

step 3,001

Is the sauce spicy?

Question2

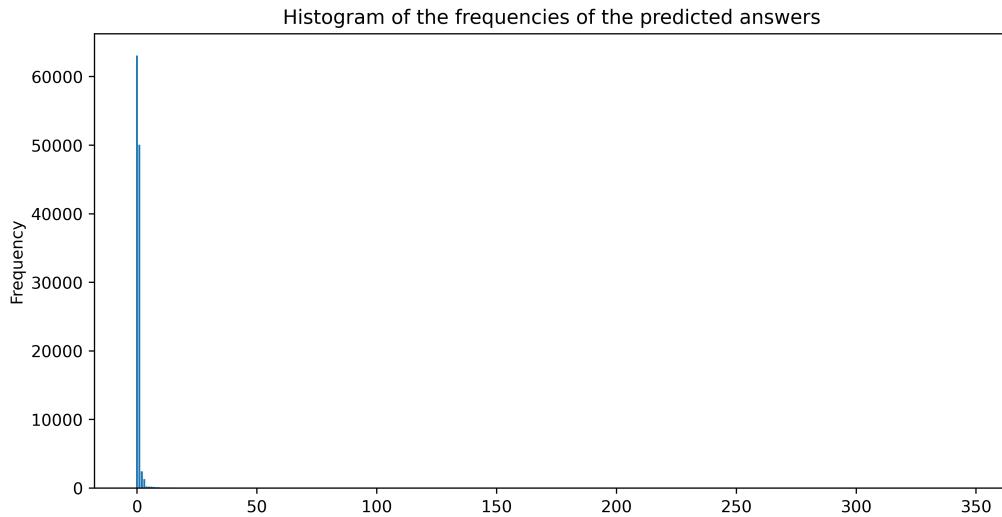
Question2/text_summary
tag: Question2/text_summary

step 3,002

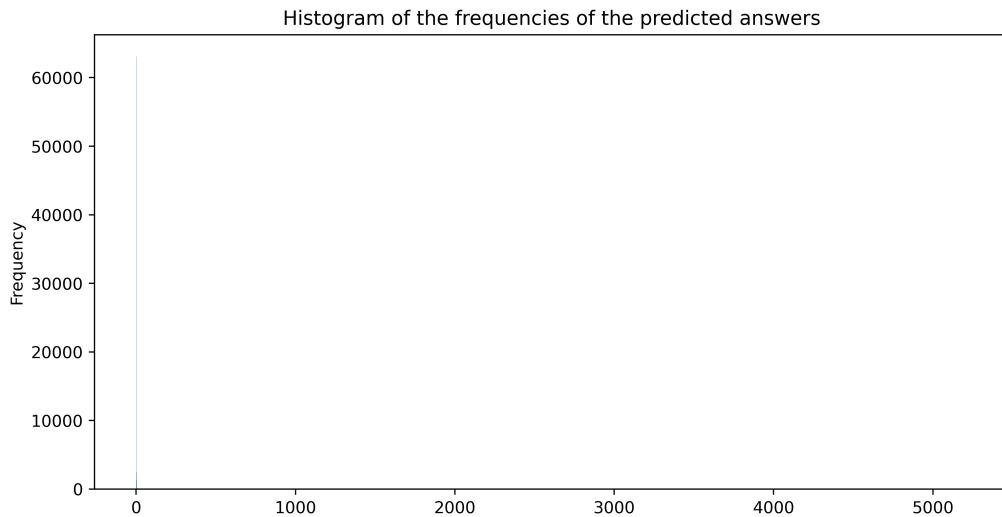
How many headlights are visible?

2.7 Make a histogram of the frequencies of the predicted answers (all 5217) without labels.

Show only the label with predictions



Show all the 5217 labels



2.8 What are the 10 most frequently predicted classes? Why do you think this happens? You can optionally visualize the frequency of ground-truth answers in the training set.

10 most frequent classes are: ['Other', 'yes', '2', 'no', 'baseball', 'broccoli', 'red', 'sheep', '3', 'surfing']
Because network is not expressive enough to capture the causal factors the data, it will try to capture the output distribution. That is, because the output distribution, the distribution of the labels in the

training dataset, is largely skewed, having a majority of "other", "yes", and so on as we can see in the histogram, the network will try to predict these top frequent classes as much as possible to get better accuracy.

Task 3: Transformer Network

3.1 Complete the CrossAttentionLayer

```
1 class CrossAttentionLayer(nn.Module):
2     """Self-/Cross-attention between two sequences."""
3
4     def __init__(self, d_model=256, dropout=0.1, n_heads=8):
5         """Initialize layers, d_model is the encoder dimension."""
6         super().__init__()
7
8         # Self-attention for seq1
9         self.sa1 = nn.MultiheadAttention(
10             d_model, n_heads, dropout=dropout, batch_first=True
11         ) # use batch_first=True everywhere!
12         self.dropout_1 = nn.Dropout(dropout)
13         self.norm_1 = nn.LayerNorm(d_model)
14
15         # Self-attention for seq2
16         self.sa2 = nn.MultiheadAttention(
17             d_model, n_heads, dropout=dropout, batch_first=True
18         )
19         self.dropout_2 = nn.Dropout(dropout)
20         self.norm_2 = nn.LayerNorm(d_model)
21
22         # Cross attention from seq1 to seq2
23         self.cross_12 = nn.MultiheadAttention(
24             d_model, n_heads, dropout=dropout, batch_first=True
25         )
26         self.dropout_12 = nn.Dropout(dropout)
27         self.norm_12 = nn.LayerNorm(d_model)
28
29         # FFN for seq1
30         self.ffn_12 = nn.Sequential(
31             nn.Linear(d_model, 1024),
32             nn.ReLU(),
33             nn.Dropout(dropout),
34             nn.Linear(1024, d_model),
35             nn.Dropout(dropout),
36         )
37         self.norm_122 = nn.LayerNorm(d_model)
38
39         # Cross attention from seq2 to seq1
40         self.cross_21 = nn.MultiheadAttention(
41             d_model, n_heads, dropout=dropout, batch_first=True
42         )
43         self.dropout_21 = nn.Dropout()
44         self.norm_21 = nn.LayerNorm(d_model)
45
46         # FFN for seq2
47         self.ffn_21 = nn.Sequential(
```

```

48         nn.Linear(d_model, 1024),
49         nn.ReLU(),
50         nn.Dropout(dropout),
51         nn.Linear(1024, d_model),
52         nn.Dropout(dropout),
53     )
54     self.norm_212 = nn.LayerNorm(d_model)
55
56     def forward(self, seq1, seq1_key_padding_mask, seq2,
57                 seq2_key_padding_mask,
58                 seq1_pos=None, seq2_pos=None):
59         """Forward pass, seq1 (B, S1, F), seq2 (B, S2, F)."""
60         # Self-attention for seq1
61         q1 = k1 = v1 = seq1
62         if seq1_pos is not None:
63             q1 = q1 + seq1_pos
64             k1 = k1 + seq1_pos
65         seq1b = self.sa1(
66             query=q1,
67             key=k1,
68             value=v1,
69             attn_mask=None,
70             key_padding_mask=seq1_key_padding_mask # (B, S1)
71         )[0]
72         seq1 = self.norm_1(seq1 + self.dropout_1(seq1b))
73
74         # Self-attention for seq2
75         q2 = k2 = v2 = seq2
76         if seq2_pos is not None:
77             q2 = q2 + seq2_pos
78             k2 = k2 + seq2_pos
79         seq2b = self.sa2(
80             query=q2,
81             key=k2,
82             value=v2,
83             attn_mask=None,
84             key_padding_mask=seq2_key_padding_mask # (B, S2)
85         )[0]
86         seq2 = self.norm_2(seq2 + self.dropout_2(seq2b))
87
88         # Create key, query, value for seq1, seq2
89         q1 = k1 = v1 = seq1
90         q2 = k2 = v2 = seq2
91         if seq1_pos is not None:
92             q1 = q1 + seq1_pos
93             k1 = k1 + seq1_pos
94         if seq2_pos is not None:
95             q2 = q2 + seq2_pos
96             k2 = k2 + seq2_pos
97
98         # Cross-attention from seq1 to seq2 and FFN

```

```

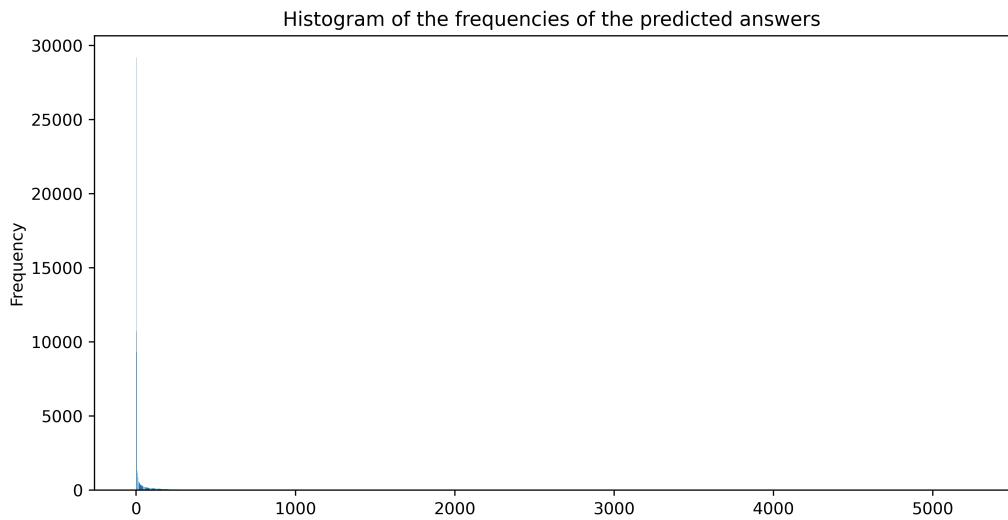
99     seq1b = self.cross_12(
100        query=q1,
101        key=k2,
102        value=v2,
103        attn_mask=None,
104        key_padding_mask=seq2_key_padding_mask # (B, S2)
105    )[0]
106    seq1 = self.norm_12(seq1 + self.dropout_12(seq1b))
107
108    # FFN for seq1
109    seq1 = self.norm_122(seq1 + self.ffn_12(seq1))
110
111    # Cross-attention from seq2 to seq1 and FFN
112    seq2b = self.cross_21(
113        query=q2,
114        key=k1,
115        value=v1,
116        attn_mask=None,
117        key_padding_mask=seq1_key_padding_mask # (B, S1)
118    )[0]
119    seq2 = self.norm_21(seq2 + self.dropout_21(seq2b))
120
121    # FFN for seq2
122    seq2 = self.norm_212(seq2 + self.ffn_21(seq2))
123
124    return seq1, seq2

```

3.2 Load the trained weights and reproduce our result.

Final Val Acc: 0.6762130489169794

3.3 How does the histogram of answers look now?



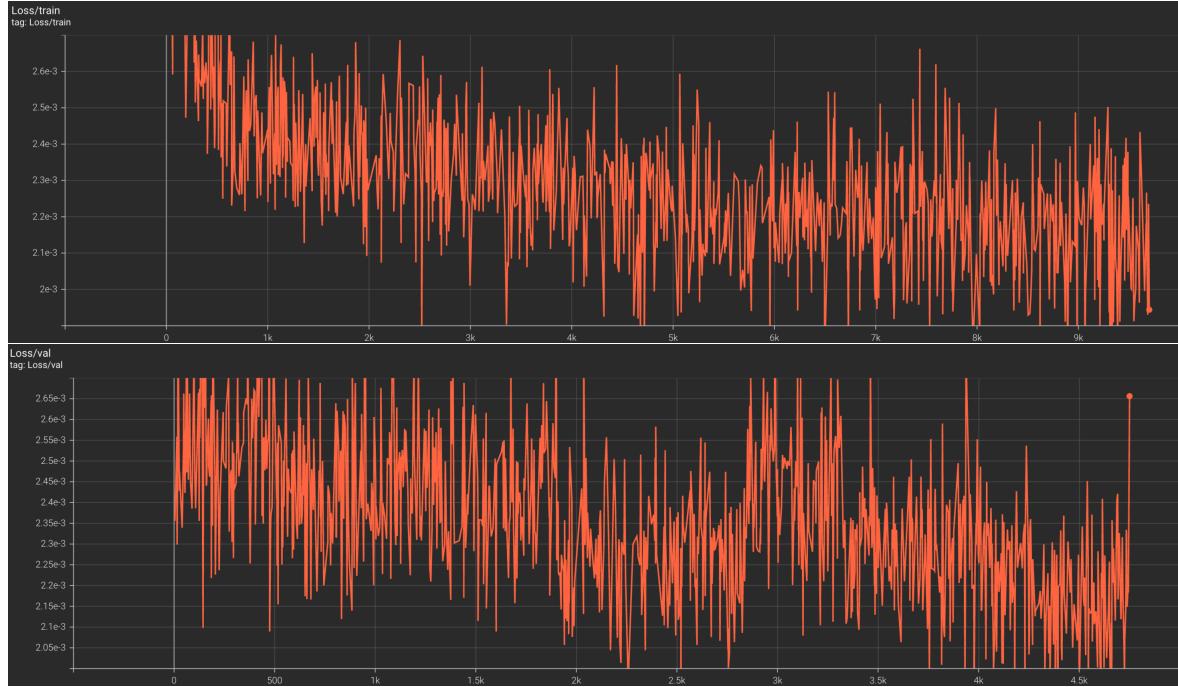
The histogram of answers look much better now with a peak that is much lower, more different sub-peaks, and more less frequent answers (not appearing just once). The top answer is "yes", not "other" in previous task.

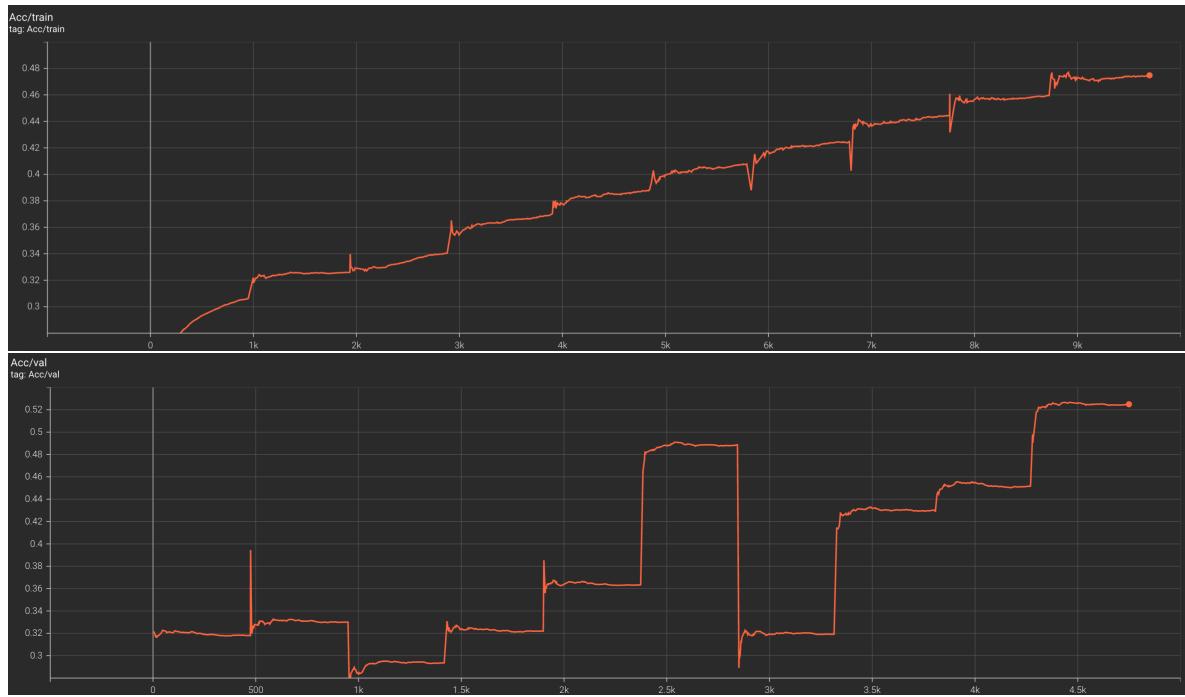
Task 4: Bells and Whistles

4.1 Is it a matter of capacity? In the BaselineNet replace the linear classifier with a deeper MLP. Report the accuracy plot of this model.

No it is not a matter of capacity. I replaced the linear classifier with two hidden layer MLP with hidden size of 2048.

```
1 self.classifier = nn.Sequential(  
2     nn.Linear(  
3         self.text_encoder.config.hidden_size + 512,  
4         2048,  
5     ),  
6     nn.ReLU(),  
7     nn.Linear(  
8         2048,  
9         2048,  
10    ),  
11    nn.ReLU(),  
12    nn.Linear(  
13        2048,  
14        n_answers,  
15    )  
16 )
```





The final val accuracy I got is 0.5249522680887484, which is even slightly worse than the baseline. The validation accuracy plot looks very unstable with jumping backwards among epochs. This is probably because the classifier is too large to train sufficiently. It shows that it is not a matter of capacity.