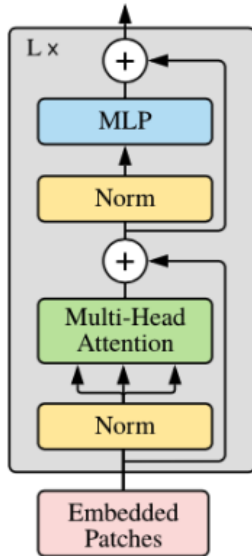## ∨  P1: Implement Encoder Block in Transformer (40%)

**Transformer Encoder**



Transformer is the most popular architecture in many fields (computer vision, speech recognition, natural language processing, etc). The image above shows a typical transformer encoder.

In this problem, you will implement the same transformer encoder as shown in the image. Particularly you will implement multi-head attention and use provided functions to finish the missing parts in the encoder block.

Some things to note:

- Use LayerNorm for the Norm as shown in the figure above.
- FeedFoward is MLP in the figure.
- The gray rectangle illustrates the structure of EncoderBlock, and you only need to implement that.

```
 1 import torch
 2 import torch.nn as nn
 3 from torch.nn import functional as F
 4 import numpy as np
 5
 6 # hyperparameters
 7 batch_size = 1
 8 block_size = 32 # Maximum context length
 9 device = 'cuda' if torch.cuda.is_available() else 'cpu'
10 n_embd = 2 # emgedding dimension
11 n_head = 2 # number of heads in multi-head attention
12 # ------------
13 torch.manual_seed(568)
```

        <torch._C.Generator at 0x7c3bc43f9410>

```python
1  class Head(nn.Module):
2      """ one head of self-attention """
3
4      def __init__(self, head_size):
5          super().__init__()
6          self.key = nn.Linear(n_embd, head_size)
7          self.query = nn.Linear(n_embd, head_size)
8          self.value = nn.Linear(n_embd, head_size)
9
10     def forward(self, x):
11         # TODO: implement the single-head attention
12         selfAttention = torch.bmm(self.query(x), self.key(x).transpose(-2, -1)) / (x.shape[2] ** 0.5)
13         selfAttention = F.softmax(selfAttention, dim=-1)
14         out = torch.bmm(selfAttention, self.value(x))
15         return out
16
17 class MultiHeadAttention(nn.Module):
18     """ multiple heads of self-attention in parallel """
19
20     def __init__(self, num_heads, head_size):
21         super().__init__()
22         self.heads = nn.ModuleList([Head(head_size) for _ in range(num_heads)])
23
24     def forward(self, x):
25         # TODO: implement the multi-head attention
26         out = torch.cat([head(x) for head in self.heads], dim=-1)
27         return out
28
29 class FeedFoward(nn.Module):
30     """ a simple linear layer followed by a non-linearity """
31
32     def __init__(self, n_embd):
33         super().__init__()
34         self.net = nn.Sequential(
35             nn.Linear(n_embd, 4 * n_embd),
36             nn.ReLU(),
37             nn.Linear(4 * n_embd, n_embd),
38         )
39
40     def forward(self, x):
41         return self.net(x)
42
43 class EncoderBlock(nn.Module):
44     def __init__(self, n_embd, n_head):
45         # n_embd: embedding dimension, n_head: the number of heads
46         super().__init__()
47         self.sa = MultiHeadAttention(n_head, n_embd // n_head)# TODO: initialize multi-head self attention
48         self.mlp = FeedFoward(n_embd) # TODO: create feed forward network
49         self.norm1 = nn.LayerNorm(n_head) # TODO: create layer norm
50         self.norm2 = nn.LayerNorm(n_embd) # TODO: create layer norm
51
52
53     # initialize the weights with an identity matrix to make sure we get the same weight everytime
54     # it is used just for reproducing the results
55     def _initialize_weights(self):
56         for m in self.modules():
57             if isinstance(m, nn.Linear):
58                 nn.init.eye_(m.weight)
59                 if m.bias is not None:
60                     nn.init.zeros_(m.bias)
61             elif isinstance(m, nn.LayerNorm):
62                 nn.init.ones_(m.weight)
63                 nn.init.zeros_(m.bias)
64
65     def forward(self, x):
66         # TODO: reference the provided figure of transformer encoder to implement the forward function
67         out = x + self.sa(self.norm1(x))
68         out = out + self.mlp(self.norm2(out))
69         return out
```

Initialize the encoder block.

```
1 encoder_block = EncoderBlock(n_embd, n_head).to(device)
2 encoder_block._initialize_weights()
```

Now, download data_p1.pt from Canvas and upload to your working directory.

```
1 data_dict_loaded = torch.load('data_p1.pt')
2 ans = data_dict_loaded['ans']
3 data = data_dict_loaded['data']
```

Forward pass the data to the encoder block and collect the reuslts into a list.

```
1 results = []
2 for datum in data:
3     results.append(encoder_block(datum.to(device)).detach().cpu())
```

```
1 ans
```

```
                [-1.7702,  2.3148],
                [-0.9737,  1.1272],
                [ 1.6078,  0.1909],
                [ 1.6731, -0.3200],
                [ 3.1005,  0.5838],
                [-1.2745,  1.0775],
                [ 2.7161, -0.6739],
                [ 2.8139,  0.6781],
                [ 0.6114, -0.4760],
                [-1.7369,  0.6209],
                [-1.5547,  1.4092],
                [ 0.6270, -0.4385],
                [ 3.5181, -0.6285],
                [-2.1928,  1.6743],
                [ 3.6704, -0.4864],
                [-1.0429,  1.2313],
                [-1.6441,  0.2750],
                [ 1.7162,  0.3108],
                [-1.1280,  1.5796],
                [ 0.2795,  1.3356],
                [-0.0457,  1.4466],
                [ 3.4281,  0.3123],
                [ 2.9692,  0.7618],
                [ 3.3695, -1.2879],
                [ 3.2384,  0.6489],
                [ 3.3479, -0.6007],
                [ 0.5302,  2.4495]]]),
        tensor([[[-1.0259,  0.6152],
                [-2.2288, -0.3617],
                [-0.9368,  0.4891],
                [ 1.9080, -0.0684],
                [-1.5225,  0.8566],
                [-0.9295,  1.3268],
                [-1.3441,  0.9111],
                [-1.5945, -0.4248],
                [ 1.1780, -3.0424],
                [-0.5267,  0.8443],
                [-0.7441,  0.6056],
                [ 1.6322, -0.4910],
                [-1.2116,  0.3364],
                [ 1.9748, -0.4159],
                [ 1.2068, -0.3376],
                [-1.7814,  0.6242],
                [-2.9207,  0.3983],
                [ 2.1858, -0.2490],
                [ 2.4302,  0.8357],
```

```
                [ 1.0421,   0.0202]],
                [-1.2421,   0.2006],
                [-1.1448,  -0.9206],
                [ 2.0967,   0.0365],
                [ 2.2502,   1.0912],
                [ 1.8053,  -0.3889],
                [-0.1921,   1.3671],
                [ 2.1004,   0.7266]]]),
```

1 results

```
                [ 2.1111,   0.1221]],
                [-0.3857,   2.3982],
                [ 1.6878,   0.1076],
                [ 1.3279,  -1.0041],
                [-0.8110,   0.3060]]]),
    tensor([[[ 2.6773,   1.5188],
                [-0.2199,   0.9685],
                [-1.7405,   1.7295],
                [-1.3190,   0.6758],
                [ 2.2541,  -1.1568],
                [ 2.7457,   1.0191],
                [ 1.6217,   0.1682],
                [ 3.2371,  -0.6552],
                [-1.3366,   1.7308],
                [ 3.1672,   0.7356],
                [-2.2987,  -0.1141],
                [-1.1696,   1.5834],
                [-1.5938,   0.0408],
                [ 2.4058,  -0.0688],
                [ 2.3611,   0.3585],
                [-1.9571,  -0.3031],
                [-1.9340,   0.5315],
                [-1.7469,   0.6587],
                [ 1.7582,  -1.0192],
                [-1.7132,   0.6404],
                [-2.1351,   0.6079],
                [ 0.2586,   3.3612],
                [ 2.8285,   1.5993],
                [-1.0132,   0.0518],
                [-1.5383,   0.3088],
                [-1.5515,   1.9327],
                [ 1.5559,  -0.1682],
                [-2.2796,   0.9164],
                [ 1.5258,  -0.2152],
                [ 1.8087,  -0.9843],
                [-1.5651,   0.9921],
                [-1.7276,   0.3684]]]),
    tensor([[[ 1.2099,   0.1582],
                [ 3.2397,   0.0430],
                [-2.2824,   1.4001],
                [ 2.7031,  -0.1075],
                [ 2.8890,   0.7328],
                [-1.0576,   1.4579],
                [-2.2220,   0.6479],
                [ 2.2413,   1.1112],
                [-0.5619,   0.4465],
                [-1.3867,   1.8416],
                [-0.1727,   0.9043],
                [ 0.2314,   1.3074],
                [ 1.7920,   0.4774],
                [ 0.5770,  -1.2372],
                [-0.3655,   0.7035],
                [ 0.3585,  -1.1778],
                [-2.3048,   1.1799],
                [-0.2198,   0.9215],
                [ 2.1537,  -1.7752],
                [ 1.5642,   0.0756],
                [ 0.9455,  -1.2763],
                [-2.2637,   0.0240],
                [ 2.00,   0.0]
```

Make sure your model produce the same result for every test case.

```
1 ## Note, I changed the relative tolerance due to the fact that the values only ever go upto the 0.0001 place anyway.
2 ## going further gives not use.
3
4 for i in range(len(data)):
5     assert torch.allclose(results[i], ans[i], rtol=0.0001), f"Your implementation is wrong for the {i}-th example."
```

```
6 print('Pass')

    Pass
```

Save your model to a file named **enc_block.pth** and include **enc_block.pth** in your submission. We will test your encoder block with different test cases for you to get full credits.

```
1 PATH = './enc_block.pth'
2 torch.save(encoder_block.state_dict(), PATH)
```

In the decoder of transformer models, masking in self-attention is crucial for two primary reasons:

- Preventing Future Information Leakage: In tasks like language translation or text generation, the future tokens (words that come later in the sequence) should not influence the prediction of the current token. Masking ensures that while predicting a particular word, the model only has access to previously generated words and not to any future words. This is essential for the model to learn a proper language structure and generate coherent and contextually appropriate text.

- Maintaining Autoregressive Property: Transformers, especially in tasks like text generation, operate in an autoregressive manner, meaning they generate one part of the sequence at a time. By masking future tokens in the self-attention mechanism, the model respects this autoregressive nature, ensuring that each step of generation or translation is based solely on already known or generated information.

Please see https://arxiv.org/abs/1706.03762 for more details on Transformer.

**Question**: What should we modify if we want to implement masked attention head in transformer's decoder block? Please implement the masked attention by modifying your previous class Head.

```
1 class MaskedHead(nn.Module):
2     """ one head of self-attention """
3
4     def __init__(self, head_size):
5         super().__init__()
6         self.key = nn.Linear(n_embd, head_size)
7         self.query = nn.Linear(n_embd, head_size)
8         self.value = nn.Linear(n_embd, head_size)
9         self.register_buffer('mask', torch.tril(torch.ones(block_size, block_size)))
10
11     def forward(self, x):
12         # TODO: implement the single-head attention
13         _, maskShape, size = x.shape
14         selfAttention = torch.bmm(self.query(x), self.key(x).transpose(-2, -1)) / (size ** 0.5)
15         selfAttention = selfAttention.masked_fill(self.mask[:maskShape, :maskShape] == 0, float('-inf'))
16         selfAttention = F.softmax(selfAttention, dim=-1)
17         out = torch.bmm(selfAttention, self.value(x))
18         return out
```

```
1 encoder_block = EncoderBlock(n_embd, n_head).to(device)
2 encoder_block._initialize_weights()
3
4 data_dict_loaded = torch.load('data_p1.pt')
5 ans = data_dict_loaded['ans']
6 data = data_dict_loaded['data']
7
8 results = []
9 for datum in data:
10     results.append(encoder_block(datum.to(device)).detach().cpu())
```

```
1 results
```

```
           [-0.2199,  0.9685],
           [-1.7405,  1.7295],
           [-1.3190,  0.6758],
           [ 2.2541, -1.1568],
           [ 2.7457,  1.0191],
           [ 1.6217,  0.1682],
           [ 3.2371, -0.6552],
           [-1.3366,  1.7308],
           [ 3.1672,  0.7356],
           [-2.2987, -0.1141],
           [-1.1696,  1.5834],
           [-1.5938,  0.0408],
           [ 2.4058, -0.0688],
           [ 2.3611,  0.3585],
           [-1.9571, -0.3031],
           [-1.9340,  0.5315],
           [-1.7469,  0.6587],
           [ 1.7582, -1.0192],
           [-1.7132,  0.6404],
           [-2.1351,  0.6079],
           [ 0.2586,  3.3612],
           [ 2.8285,  1.5993],
           [-1.0132,  0.0518],
           [-1.5383,  0.3088],
           [-1.5515,  1.9327],
           [ 1.5559, -0.1682],
           [-2.2796,  0.9164],
           [ 1.5258, -0.2152],
           [ 1.8087, -0.9843],
           [-1.5651,  0.9921],
           [-1.7276,  0.3684]]]),
   tensor([[[ 1.2099,  0.1582],
           [ 3.2397,  0.0430],
           [-2.2824,  1.4001],
           [ 2.7031, -0.1075],
           [ 2.8890,  0.7328],
           [-1.0576,  1.4579],
           [-2.2220,  0.6479],
           [ 2.2413,  1.1112],
           [-0.5619,  0.4465],
           [-1.3867,  1.8416],
           [-0.1727,  0.9043],
           [ 0.2314,  1.3074],
           [ 1.7920,  0.4774],
           [ 0.5770, -1.2372],
           [-0.3655,  0.7035],
           [ 0.3585, -1.1778],
           [-2.3048,  1.1799],
           [-0.2198,  0.9215],
           [ 2.1537, -1.7752],
           [ 1.5642,  0.0756],
           [ 0.9455, -1.2763],
           [-2.2637,  0.0240],
```

1 ans

```
                [ 1.8053, -0.3889],
                [-0.1921,  1.3671],
                [ 2.1094,  0.7266]]]),
        tensor([[[ 0.8151, -0.9453],
                [ 2.5775, -1.7437],
                [-1.6758,  0.1372],
                [-1.5354,  0.0322],
                [-0.8102,  1.0794],
                [ 1.6140,  0.6365],
                [ 2.0891,  1.0063],
                [ 2.2872,  0.7041],
                [ 4.4790,  0.8228],
                [ 2.6819,  0.9185],
                [-0.6281,  1.5366],
                [ 0.0072,  1.1369],
                [-1.3940,  0.2573],
                [ 2.3398,  0.7757],
                [ 2.8084, -0.2453],
                [ 1.9982,  0.5587],
                [-1.8077,  1.1651],
                [ 3.3140,  1.9080],
                [-1.3583,  0.5102],
                [ 2.5500,  0.8867],
                [ 2.6635,  1.0557],
                [-1.8199,  1.7828],
                [ 3.2988,  1.5809],
                [-0.5473,  1.7100],
                [-1.6121,  1.0381],
                [-1.6970,  1.2286],
                [ 1.8611,  0.6345],
                [ 1.4477,  0.2314],
                [-0.3857,  2.3982],
                [ 1.6878,  0.1076],
                [ 1.3279, -1.0041],
                [-0.8110,  0.3060]]]),
        tensor([[[ 2.6773,  1.5188],
                [-0.2199,  0.9685],
```

1