

---

## ▼ Lab 10 - Language Modeling with Transformers

In this lab, your task is to build a character-level language model with LSTM layer.

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

```
Mounted at /content/drive
```

```
1 cd "/content/drive/MyDrive/UCCD3074_Labs/UCCD3074_Lab10"
```

```
/content/drive/MyDrive/UCCD3074_Labs/UCCD3074_Lab10
```

```
1 import os
2 import math
3 import torch
4 import torch.nn as nn
5 from torch.nn import functional as F
6
7 import time
```

```
1 torch.manual_seed(1234)
2 device = 'cuda' if torch.cuda.is_available() else "cpu"
```

```
1 if not os.path.exists('input.txt'):
2     !wget 'https://raw.githubusercontent.com/karpathy/char-rnn/master/data/tinyshakespeare/input.txt'
```

---

## ▼ Preprocessing

Read the dataset into the string `raw_data`.

```
1 with open('./input.txt', 'r', encoding='utf-8') as f:
2     raw_data = f.read()
```

Create the vocabulary

```
1 vocab = sorted(list(set(raw_data)))
2 vocab_size=len(vocab)
3 print('vocab_size:', vocab_size)
```

```
    vocab_size: 65
```

```
1 stoi = {ch:i for i, ch in enumerate(vocab)}
2 itos = {i:ch for i, ch in enumerate(vocab)}
```

```
1 encode_text = lambda s : [stoi[c] for c in s]    # encode: take a string, output a list of integers
2 decode_text = lambda l : ''.join([itos[i] for i in l])
```

```
1 data = torch.tensor(encode_text(raw_data), dtype=torch.long)
```

The function `get_batch` randomly sample a block of text as input `x`. The label `y` is the block of text shifted by 1 position of `x`.

```
1 def get_batch(batch_size, block_size):
2     ix = torch.randint(len(data) - block_size, (batch_size,))
3     x = torch.stack([data[i:i+block_size] for i in ix])
4     y = torch.stack([data[i+1:i+block_size+1] for i in ix])
5     x, y = x.to(device), y.to(device)
6     return x, y
```

---

## ▼ Character-level language model with Transformer model

Now, let's build the character-level language model with LSTM. The network `LMNet` has the following layers:

Layer	Configuration	Shape
Input	-	(B, T)
Token Embedding	num_embedding = vocab_size, embedding_dim = d_model	(B, T, d_model)
Position Embedding	d_model = d_model	(B, T, d_model)
TransformerEncoder	d_model (default: 512), nhead (default = 6), dim_feedforward, batch_first (default: False)	(B, T, d_model)
fc	in_features = d_model, out_features = vocab_size	(B, T, vocab_size)

## ▼ Transformer in PyTorch

PyTorch offers the necessary [transformer layers](#) to built the transformer model, including:

1. [nn.Transformer](#) - a tranformer model. The module is constructed from `nn.TranformerEncoder` and `nn.TransformerDecoder` below.
2. [nn.TransformerEncoder](#) - a stack of  $N$  encoder layers
3. [nn.TransformerDecoder](#) - a stack of  $N$  decoder layers
4. [nn.TransformerEncoderLayer](#) - made up of multi-head (self) attention and feedforward network
5. [nn.TransformerDecoderLayer](#) - made up of multi-head (self) attention, multi-head (encoder) attention and feedforward network

## Transformer Language Model

To implement a language model, we need only the encoder portion of Transformer. It can be constructed with `nn.TransformerEncoder` and `nn.TransformerEncoderLayer`. The following diagram show the network that we shall use to build the language model.

The network is implemented in two steps.

1. Create an encoder layer with `nn.TransformerEncoderLayer`.

- `d_model` (int) – the number of expected features in the input (required).
- `nhead` (int) – the number of heads in the multiheadattention models (required).
- `dim_feedforward` (int) – the dimension of the feedforward network model (default=2048).
- `batch_first` (bool) – If True, then the input and output tensors are provided as (batch, seq, feature). Default: False (seq, batch, feature).

2. Create the encoder block with `nn.TransformerEncoder` by stacking multiple encoder layers

- `encoder_layer` – an instance of the `TransformerEncoderLayer()` class (required).
- `num_layers` – the number of sub-encoder-layers in the encoder (required).

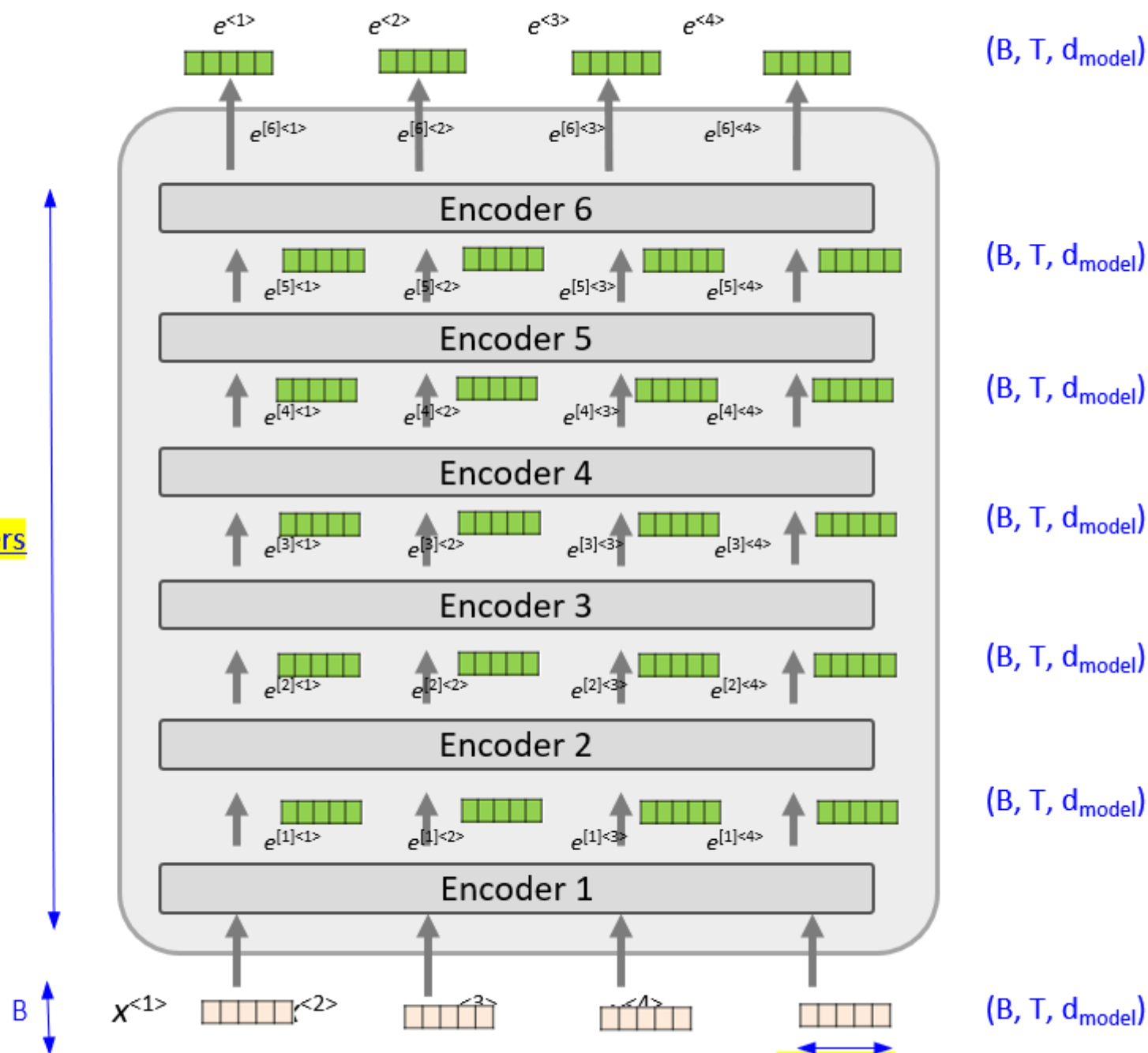
The following figure shows a Transformer's Encoder block. It consists of `num_layers` encoder layers. When `batch_size = True`, the input to the transformer model has a shape of  $(B, T, d_{model})$  where  $B$  is the batch size,  $T$  is the sequence length and  $d_{model}$  is the feature dimension of each token. The transformer model can handle different batch sizes  $B$  and sequence length  $T$  during inference. However, the `d_model` is a hyperparameter and is fixed during inference.





batch\_first = true

num\_layers



`d_model`

The following figure shows the block diagram for the encoder layer. Each layer contains a total of `n_head` heads. Each head consists of two modules:

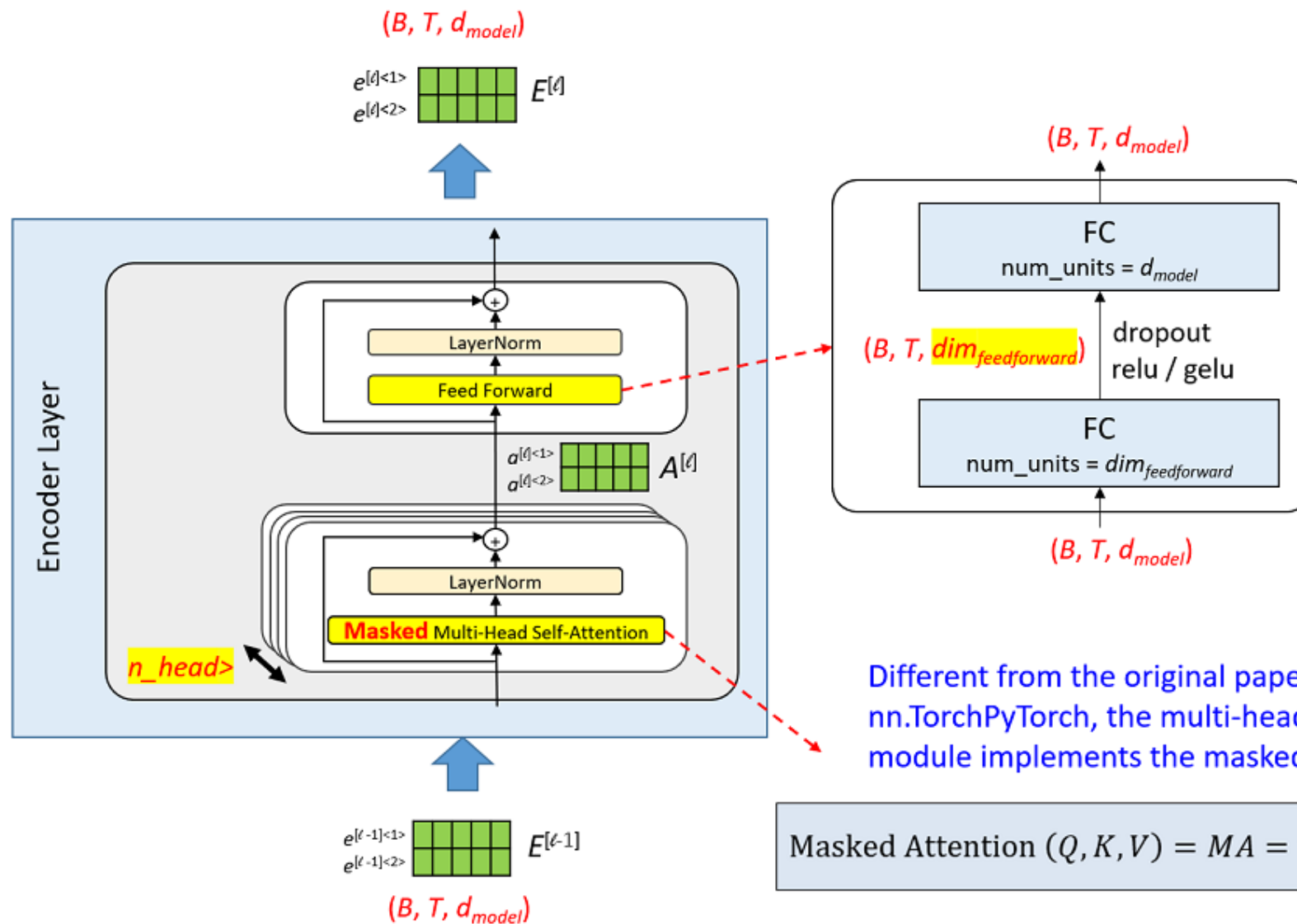
1. Self-attention module

Different from the origin transformer, PyTorch implementation replaces the *multi-head self-attention* layer with the **masked** version as shown below. This allows it the encoder be used for generative language modeling.

2. Feed forward module

This modules consists of two linear layer. The output feature length of the 1st linear layer is set to `dim_feedforward` while the 2nd linear layer is set to `d_model`.

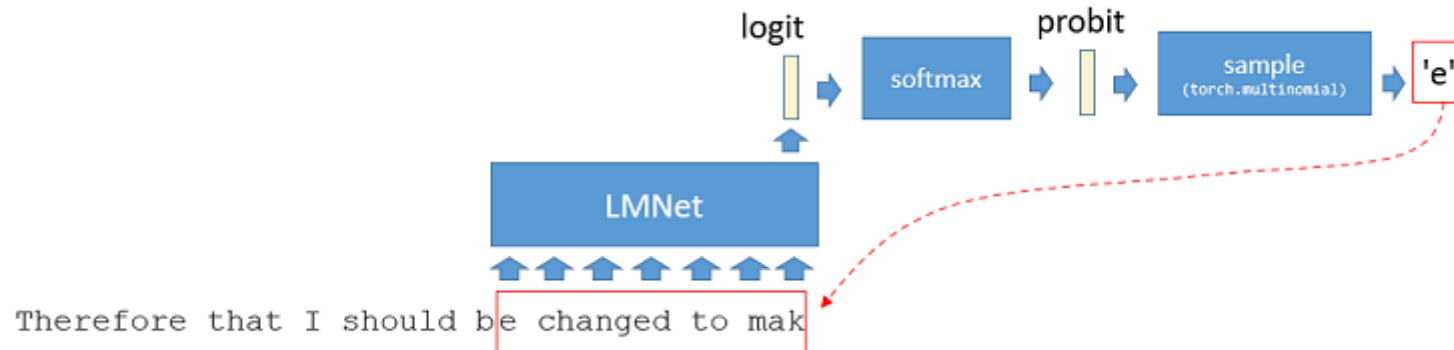




## ▼ Generating novel text

To generate novel text, we implement the method `generate`. Since the network is trained on a sequence of length  $T = \text{block\_size}$ , when generating the text, we feed the most recent `block_size` characters into the network to generate the next character. Here are the steps:

1. Crop the most recent `block_size` characters in the generated text
2. The cropped text is fed to the generative model to generate the next character. The network output the `logit` value.
3. Convert the logit of the network to `probit` value by performing `softmax` operation.
4. Sample a character from the `probit` by using [torch.multinorm](#)
5. Append the sampled character to the end of the generated text.
6. Repeat steps 1-5 for `text_len` times



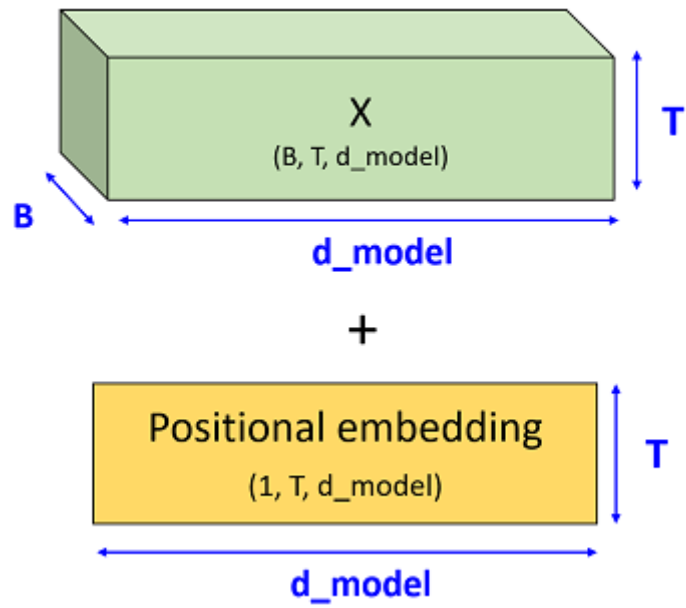
## ▼ Positional Embedding

The following `PositionalEncoding` module injects some information about the relative or absolute position of the tokens in the sequence. The positional encodings have the same dimension as the embeddings so that the two can be summed. To do so, we can use sine and cosine functions of different frequencies to embed the distance.

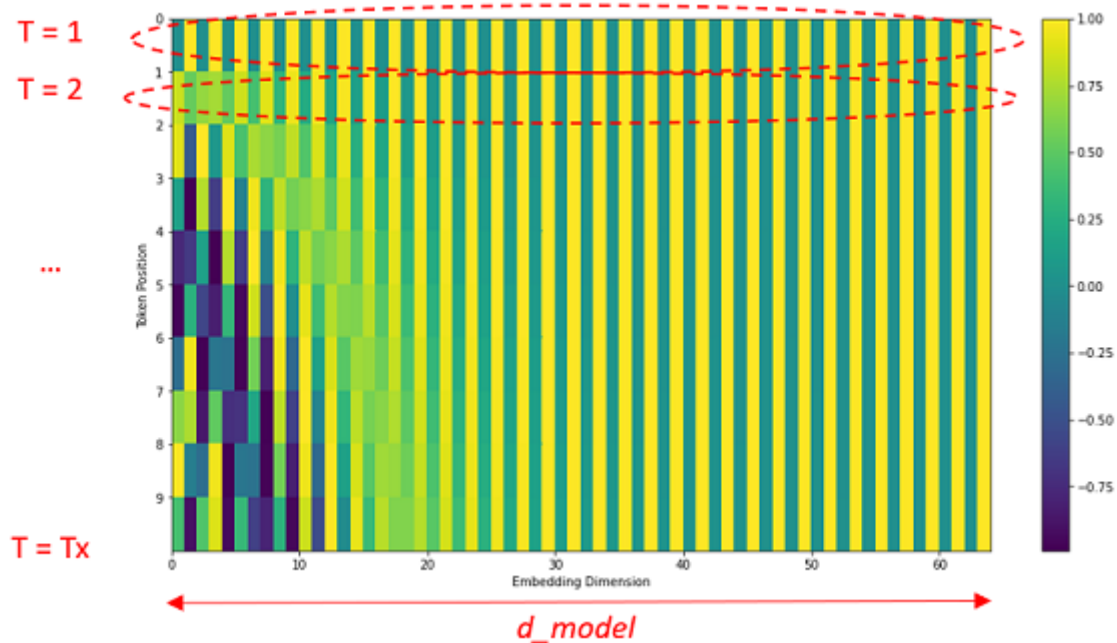
$$PE_{(t,2i)} = \sin(t/10000^{2i/d_{model}})$$

$$PE_{(t,2i+1)} = \cos(t/10000^{2i/d_{model}})$$

where  $t$  is the (time) index of the token in the input sequence and  $i$  is the index of the embedding feature representing the token



The following diagram shows how the position embedding looks like for different time steps.



## ▼ Implement the Position Embedding

```

1 class PositionEmbedding(nn.Module):
2
3     def __init__(self, d_model: int, batch_first: bool = False, max_len: int = 5000, dropout: float = 0.1):
4         super().__init__()
5         self.dropout = nn.Dropout(p=dropout)
6         self.batch_first = batch_first
7
8         # Shape of position: (max_len, 1)
9         position = torch.arange(max_len).unsqueeze(1)
10
11         # Shape of div_term: 1d tensor of shape (d_model/2, )
12         div_term1 = torch.exp(torch.arange(0, d_model, 2) * (-math.log(10000.0) / d_model))
13         div_term2 = div_term1[:-1] if d_model % 2 == 1 else div_term1
14
15         self.pe = torch.zeros(max_len, d_model)

```

```

15     # snake or pe: (max_len, 1, d_model)
16     if self.batch_first:
17         pe = torch.zeros(1, max_len, d_model)
18         pe[0, :, 0::2] = torch.sin(position * div_term1)
19         pe[0, :, 1::2] = torch.cos(position * div_term2)
20         self.register_buffer('pe', pe)
21     else:
22         pe = torch.zeros(max_len, 1, d_model)
23         pe[:, 0, 0::2] = torch.sin(position * div_term1)
24         pe[:, 0, 1::2] = torch.cos(position * div_term2)
25         self.register_buffer('pe', pe)
26
27     def forward(self, x):
28         """
29         Arguments:
30             x: Tensor, shape ``[seq_len, batch_size, embedding_dim]``
31         """
32         # Shape of self.pe: (T, 1, embedding_dim)
33         # shape of x      : (T, B, embedding_dim)
34         if self.batch_first:
35             x = x + self.pe[:, :x.size(1), :]
36         else:
37             x = x + self.pe[:x.size(0)]
38         return self.dropout(x)

```

Test the position embedding layer.

```

1 batch_size = 4
2 seq_len    = 8
3 d_model    = 512
4
5 x = torch.randn(batch_size, seq_len, d_model)
6 print('Shape of x:', x.shape)

```

Shape of x: torch.Size([4, 8, 512])

```

1 pos_embedding = PositionEmbedding(d_model, batch_first = True)
2 xp = pos_embedding(x)
3 print('Shape of xp:', xp.shape)

```

Shape of xp: torch.Size([4, 8, 512])

## ▼ Implement the Language Model

```

1 class LMNet(nn.Module):
2     def __init__(self, vocab_size, d_model=512, nhead=8, dim_feedforward=2048, num_layers=6, batch_first=True):
3         super().__init__()
4
5         # token embedding
6         self.token_embedding = nn.Embedding(vocab_size, d_model)
7
8         # positional embedding layer
9         self.pos_embedding = PositionEmbedding(d_model, batch_first)
10
11        # create the encoder block
12        encoder_layer = nn.TransformerEncoderLayer(d_model, nhead, dim_feedforward, batch_first=batch_first)
13        self.encoder_net = nn.TransformerEncoder(encoder_layer, num_layers)
14
15        # create the linear layer
16        self.fc = nn.Linear(d_model, vocab_size)
17
18    def forward(self, x, mask=None):          # (B,T)
19
20        # convert each token into its embedding vector
21        x = self.token_embedding(x)         # (B,T,d_model)
22
23        # add position embedding to x
24        x = self.pos_embedding(x)          # (B,T,d_model)
25
26        # perform sequence model inference
27        if mask is not None:

```

```

28         x = self.encoder_net(x, mask) # (B,T,d_model)
29     else:
30         x = self.encoder_net(x)
31
32     # generate logits
33     x = self.fc(x) # (B,T,vocab_size)
34
35     return x
36
37 def generate(self, text_len, block_size):
38
39     text = torch.zeros((1,1), dtype=torch.long).to(device) # text token (B, T) where B is fixed to 1, T = 1 initially
40
41     # repeat until the length of text = "text_len"
42     for _ in range(text_len): # (B, T)
43
44         # crop text to the last block-size tokens
45         text_cond = text[:, -block_size:] # (B, T)
46
47         # get the predictions
48         seq_len = text_cond.shape[-1]
49         mask = create_mask(seq_len)
50         yhat = self(text_cond, mask) # logits: (B, T, C)
51
52         # focus only on the last time step
53         yhat = yhat[:, -1, :] # becomes (B, C)
54
55         # apply soft max to get probabilities
56         probs = F.softmax(yhat, dim=-1) # (B, C)
57
58         # sample from distribution
59         next_token = torch.multinomial(probs, num_samples=1) # (B, 1)
60
61         # append sampled index to the running sequence
62         text = torch.cat((text, next_token), dim=1) # (B, T+1)
63
64         # print the sample

```

```
65         print(itos[next_token.item()], end='')
66         time.sleep(0.1)
67
```

Test the model without any masking

```
1 d_model=512
2 nhead=8
3 dim_feedforward=2048
4 num_layers=6
5 vocab_size = len(vocab)
6 batch_first = True
7 batch_size=4
8 block_size=8
```

```
1 model = LMNet(vocab_size, d_model, nhead, dim_feedforward, num_layers, batch_first).to(device)
```

```
1 x, y = get_batch(batch_size, block_size)
2 x, y = x.to(device), y.to(device)
3
4 yhat = model(x)
5
6 print('Shape of x: ', x.shape)
7 print('Shape of yhat: ', yhat.shape)
```

```
Shape of x:      torch.Size([4, 8])
Shape of yhat:   torch.Size([4, 8, 65])
```

Test the model with masking

```
1 def create_mask(Tx):
2     mask = torch.triu(torch.ones(Tx, Tx) * float('-inf'), diagonal=1).to(device)
3     mask = mask.to(device)
```



```
4     return mask
5
6 mask = create_mask(block_size)
7 print('Shape of mask:', mask.shape)
```

```
Shape of mask: torch.Size([8, 8])
```

```
1 yhat = model(x, mask)
2 print('Shape of yhat:', yhat.shape)
```

```
Shape of yhat: torch.Size([4, 8, 65])
```

---

## ▼ Train the model

```
1 max_iters      = 5000
2 batch_size     = 128
3 block_size     = 256
4 lr             = 5e-4
5 max_iters      = 8000
6 show_interval  = 200
7
8 d_model=256
9 nhead=4
10 dim_feedforward=512
11 num_layers=6
12 vocab_size = len(vocab)
```

Create the model

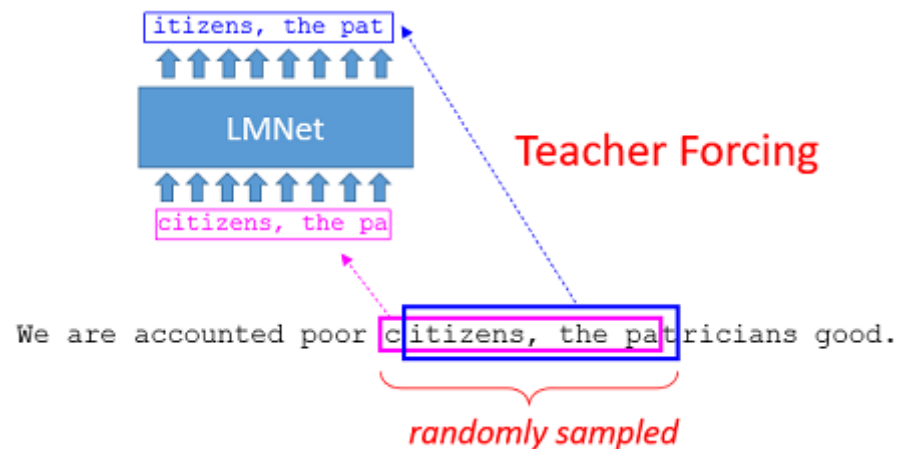
```
1 model = LMNet(vocab_size, d_model, nhead, dim_feedforward, num_layers).to(device)
```

## Create the optimizer

```
1 optimizer = torch.optim.AdamW(model.parameters()), lr=lr)
```

## Train the model

To train the model, we use **teacher forcing** where the predicted output is simply the 1-shifted sequence of the input sequence. We shall train the network with sentence sequence of length `block_size`. Since the network is trained on sequences of length `block_size`, during inference, the generative model should use input sequence of similar length to get good results.



During training, the network is based on the many-to-many architecture. However, during inference (Figure at generating novel text), the network is based on many-to-one architecture.

