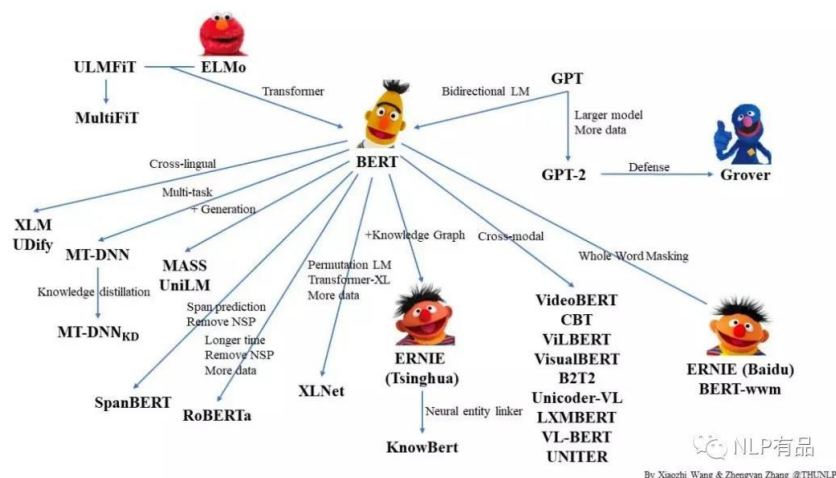


# Table of Contents

- 1 输入序列与嵌入表示
- ▼ 2 单词级和句子级的多任务学习
  - 2.1 遮蔽语言模型 (MLM)训练任务
  - 2.2 下一句预测任务
- ▼ 3 模型架构
  - 3.1 编码器结构
  - 3.2 多头注意力机制与缩放点积
  - 3.3 前馈神经网络
  - 3.4 pad掩码
  - 3.5 模型任务
- ▼ 4 PyTorch实现
  - 4.1 导包及参数设置
  - 4.2 数据处理
  - 4.3 嵌入表示
  - 4.4 编码器定义
  - 4.5 缩放点积、多头注意力机制与前馈神经网络
  - 4.6 模型定义
  - 4.7 模型训练及测试

## BERT原理及实现



## 1 输入序列与嵌入表示

Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing	[SEP]
Token Embeddings	$E_{[CLS]}$	$E_{my}$	$E_{dog}$	$E_{is}$	$E_{cute}$	$E_{[SEP]}$	$E_{he}$	$E_{likes}$	$E_{play}$	$E_{##ing}$	$E_{[SEP]}$
	+	+	+	+	+	+	+	+	+	+	+
Segment Embeddings	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$
	+	+	+	+	+	+	+	+	+	+	+
Position Embeddings	$E_0$	$E_1$	$E_2$	$E_3$	$E_4$	$E_5$	$E_6$	$E_7$	$E_8$	$E_9$	$E_{10}$

输入序列

$$input = ([CLS], s_1, s_2, \dots, s_m, [SEP], p_1, p_2, \dots, p_n, [SEP])$$

其中,  $s_i, p_j \in \mathbb{N}$  为输入符号表中的序号; 子序列  $(s_1, \dots, s_m)$  为句子对中前序句子; 子序列  $(p_1, \dots, p_n)$  为句子对中后续句子; 输入序列首标记  $[CLS]$  用作分类任务表示; 特殊标记  $[SEP]$  用作区分句子对各子句。

符号嵌入  $Embedding_{tok}(input) \in \mathbb{R}^{N \times d_{model}}$ , 其中,  $N$  为输入符号个数,  $d_{model}$  为符号嵌入维度。

句子分割嵌入  $Embedding_{seg}(input) \in \mathbb{R}^{N \times d_{model}}$ , 其中,  $N$  为输入符号个数,  $d_{model}$  为句子嵌入维度。对于单句输入, 仅使用句子A嵌入。

符号位置嵌入  $Embedding_{pos}(input) \in \mathbb{R}^{N \times d_{model}}$ , 其中,  $N$  为输入符号个数,  $d_{model}$  为符号位置嵌入维度。

## 2 单词级和句子级的多任务学习

### 2.1 遮蔽语言模型 (MLM) 训练任务

遮蔽语言模型可描述为给定单词上下文序列后, 当前单词出现的条件概率的乘积:

$$P(w_1^T) = \prod_{t=1}^T P(w_t | w_1^{t-1}, w_{t+1}^T)$$

其中,  $w_t$  是第  $t$  个单词,  $w_i^j = (w_i, w_{i+1}, \dots, w_{j-1}, w_j)$  是从第  $i$  个单词到第  $j$  个单词的子序列。

具体的, 训练数据集中选择15%的遮蔽单词  $w_t$ , 并以特殊标记  $[MASK]$  进行替换。为减少  $[MASK]$  标记对微调的影响, 数据生成器将执行以下操作, 而不是始终用  $[MASK]$  替换所选单词:

- 80%的时间: 将单词替换为  $[MASK]$ ;
- 10%的时间: 用随机单词替换单词;
- 10%的时间: 保持单词不变。这样做的目的是使表示偏向实际观察到的单词。

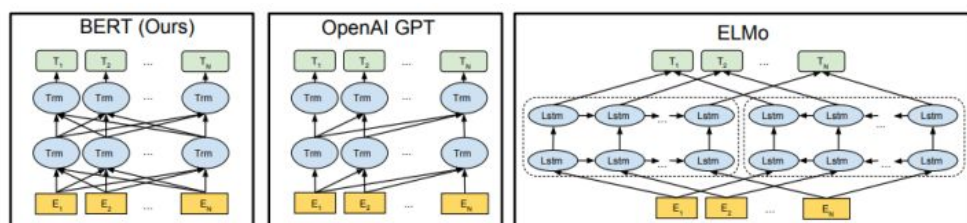
### 2.2 下一句预测任务

从语料库中生成二值化的下一句句预测任务。

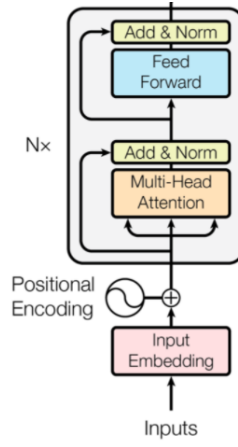
具体的, 当为每个预训练选择句子A和B时, B的50%的时间是跟随A的实际下一个句子, 而50%的时间是来自语料库的随机句子。

- input=[CLS] the man went to [MASK] store [SEP] he bought a gallon [MASK] milk [SEP] label = IsNext
- input=[CLS] the man [MASK] to the store [SEP] penguin [MASK] are flight ##less birds [SEP] label = NotNext

## 3 模型架构



### 3.1 编码器结构



编码器结构：

$$e_0 = \text{Embedding}_{tok}(\text{inputs}) + \text{Embedding}_{seg}(\text{inputs}) + \text{Embedding}_{pos}(\text{inputs})$$

$$e_l = \text{EncoderLayer}(e_{l-1}), l \in [1, n]$$

其中,  $e_0 \in \mathbb{R}^{N \times d_{model}}$  为编码器输入,  $\text{EncoderLayer}(\cdot)$  为编码器层,  $n$  为层数,  $e_l \in \mathbb{R}^{N \times d_{model}}$  为第  $l$  层编码器层输出。

编码器层  $\text{EncoderLayer}$ :

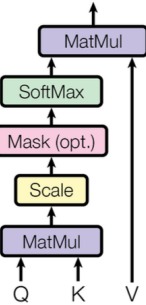
$$e_{mid} = \text{LayerNorm}(e_{in} + \text{MultiHeadAttention}(e_{in}))$$

$$e_{out} = \text{LayerNorm}(e_{mid} + \text{FFN}(e_{mid}))$$

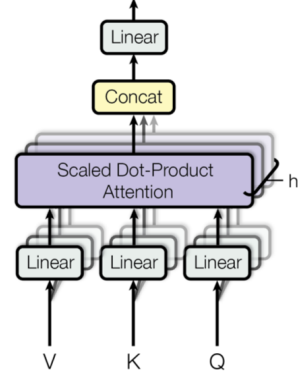
其中,  $e_{in} \in \mathbb{R}^{N \times d_{model}}$  为编码器层输入,  $e_{out} \in \mathbb{R}^{N \times d_{model}}$  为编码器层输出,  $\text{MultiHeadAttention}(\cdot)$  为多头注意力机制,  $\text{FFN}(\cdot)$  为前馈神经网络,  $\text{LayerNorm}(\cdot)$  为层归一化。

### 3.2 多头注意力机制与缩放点积

Scaled Dot-Product Attention



Multi-Head Attention



输入向量序列  $e_{in} = (e_{in1}, e_{in2}, \dots, e_{inN}) \in \mathbb{R}^{N \times d_{model}}$ , 分别得到查询向量序列  $Q = e_{in}$ , 键向量序列  $K = e_{in}$ , 值向量序列  $V = e_{in}$ 。

多头注意力机制

$$\text{MultiHeadAttention}(e_{in}) = \text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

其中, 多头输出  $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$ , 可学习的参数矩阵

$$W_i^Q \in \mathbb{R}^{d_{model} \times d_k}, W_i^K \in \mathbb{R}^{d_{model} \times d_k}, W_i^V \in \mathbb{R}^{d_{model} \times d_v}, W^O \in \mathbb{R}^{hd_v \times d_{model}}$$

使用缩放点积作为打分函数的自注意力机制

$$\text{Attention}(QW_i^Q, KW_i^K, VW_i^V) = \text{softmax}\left(\frac{QW_i^Q(KW_i^K)^T}{\sqrt{d_k}}\right) VW_i^V$$

### 3.3 前馈神经网络

$$FFN(e_{mid}) = GELU(e_{mid}W_1 + b_1)W_2 + b_2$$

其中, 参数矩阵  $W_1 \in \mathbb{R}^{d_{model} \times d_{ff}}$ ,  $W_2 \in \mathbb{R}^{d_{ff} \times d_{model}}$ , 偏置  $b_1 \in \mathbb{R}^{d_{ff}}$ ,  $b_2 \in \mathbb{R}^{d_{model}}$

### 3.4 pad掩码

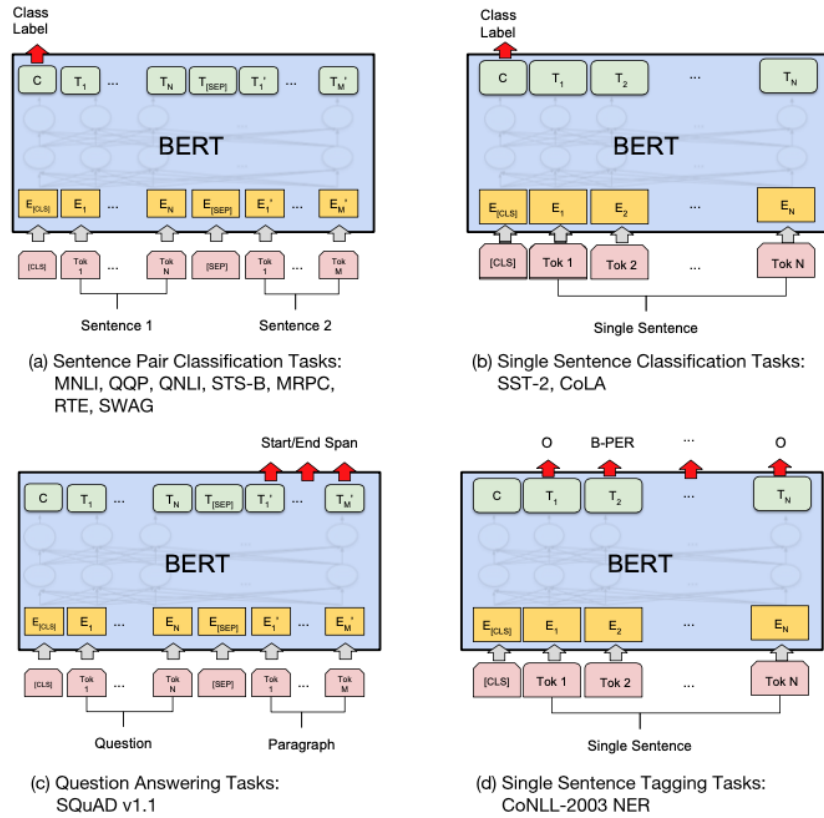
其中,

$$enc\_pad\_mask_j = (e_{j1}, e_{j2}, \dots, e_{jp}, \dots, e_{jN})$$

$$e_{jp} = \begin{cases} True, & i_p = 0 \\ False, & i_p \neq 0 \end{cases} \quad j = 1, 2, \dots, N$$

$enc\_pad\_mask \in \mathbb{R}^{N \times N}$ ,  $i_p$  为输入序列  $inputs$  对应位置序号。

### 3.5 模型任务



## 4 PyTorch实现

### 4.1 导包及参数设置

```
In [2]: 1 import math
2 import re
3 from random import *
4 import numpy as np
5 import torch
6 import torch.nn as nn
7 import torch.optim as optim
8 from torch.autograd import Variable
9
10 maxlen = 30
11 batch_size = 6
12 max_pred = 5
13 n_layers = 6
14 n_heads = 12
15 d_model = 768
16 d_ff = 768 * 4
17 d_k = d_v = 64
18 n_segments = 2
```

executed in 1.52s, finished 11:39:41 2020-06-28

## ▼ 4.2 数据处理

```

In [3]: 1 text = (
2         'Hello, how are you? I am Romeo.\n'
3         'Hello, Romeo. My name is Juliet. Nice to meet you.\n'
4         'Nice to meet you too. How are you today?\n'
5         'Great. My baseball team won the competition.\n'
6         'Oh Congratulations, Juliet\n'
7         'Thanks you Romeo'
8     )
9
10 sentences = re.sub("[.,!?\\-]", ' ', text.lower()).split('\n')
11 word_list = list(set(" ".join(sentences).split()))
12 word_dict = {'[PAD]': 0, '[CLS]': 1, '[SEP]': 2, '[MASK]': 3}
13 for i, w in enumerate(word_list):
14     word_dict[w] = i + 4
15 number_dict = {i : w for i, w in enumerate(word_dict)}
16 vocab_size = len(word_dict)
17
18 token_list = list()
19 for sentence in sentences:
20     arr = [word_dict[s] for s in sentence.split()]
21     token_list.append(arr)
22
23 def make_batch():
24     batch = []
25     positive = negative = 0
26     while positive != batch_size / 2 or negative != batch_size / 2:
27         tokens_a_index, tokens_b_index = randrange(len(sentences)), randrange(len(s
28         tokens_a, tokens_b = token_list[tokens_a_index], token_list[tokens_b_index]
29         input_ids = [word_dict['[CLS]']] + tokens_a + [word_dict['[SEP]']] + tokens
30         segment_ids = [0] * (1 + len(tokens_a) + 1) + [1] * (len(tokens_b) + 1)
31
32         n_pred = min(max_pred, max(1, int(round(len(input_ids) * 0.15))))
33         cand_maked_pos = [i for i, token in enumerate(input_ids)
34                         if token != word_dict['[CLS]'] and token != word_dict['[S
35         shuffle(cand_maked_pos)
36         masked_tokens, masked_pos = [], []
37         for pos in cand_maked_pos[: n_pred]:
38             masked_pos.append(pos)
39             masked_tokens.append(input_ids[pos])
40             if random() < 0.8:
41                 input_ids[pos] = word_dict['[MASK]']
42             elif random() < 0.5:
43                 index = randint(0, vocab_size - 1)
44                 input_ids[pos] = word_dict[number_dict[index]]
45
46         n_pad = maxlen - len(input_ids)
47         input_ids.extend([0] * n_pad)
48         segment_ids.extend([0] * n_pad)
49
50         if max_pred > n_pred:
51             n_pad = max_pred - n_pred
52             masked_tokens.extend([0] * n_pad)
53             masked_pos.extend([0] * n_pad)
54
55         if tokens_a_index + 1 == tokens_b_index and positive < batch_size / 2:
56             batch.append([input_ids, segment_ids, masked_tokens, masked_pos, True])
57             positive += 1
58         elif tokens_a_index + 1 != tokens_b_index and negative < batch_size / 2:
59             batch.append([input_ids, segment_ids, masked_tokens, masked_pos, False])
60             negative += 1
61
62     return batch

```

executed in 21ms, finished 11:40:49 2020-06-28

### 4.3 嵌入表示

```
In [9]: 1 class Embedding(nn.Module):
2         def __init__(self):
3             super(Embedding, self).__init__()
4             self.tok_embed = nn.Embedding(vocab_size, d_model)
5             self.pos_embed = nn.Embedding(maxlen, d_model)
6             self.seg_embed = nn.Embedding(n_segments, d_model)
7             self.norm = nn.LayerNorm(d_model)
8
9         def forward(self, x, seg):
10            seq_len = x.size(1)
11            pos = torch.arange(seq_len, dtype=torch.long)
12            pos = pos.unsqueeze(0).expand_as(x)
13            embedding = self.tok_embed(x) + self.pos_embed(pos) + self.seg_embed(seg)
14            return self.norm(embedding)
```

executed in 6ms, finished 11:43:43 2020-06-28

#### ▼ 4.4 编码器定义

```
In [4]: 1 class EncoderLayer(nn.Module):
2         def __init__(self):
3             super(EncoderLayer, self).__init__()
4             self.enc_self_attn = MultiHeadAttention()
5             self.pos_ffn = PoswiseFeedForwardNet()
6
7         def forward(self, enc_inputs, enc_self_attn_mask):
8             enc_outputs, attn = self.enc_self_attn(enc_inputs, enc_inputs, enc_inputs,
9             enc_outputs = self.pos_ffn(enc_outputs)
10            return enc_outputs, attn
```

executed in 5ms, finished 11:41:59 2020-06-28

#### ▼ 4.5 缩放点积、多头注意力机制与前馈神经网络

```

In [5]: 1 class ScaledDotProductAttention(nn.Module):
2         def __init__(self):
3             super(ScaledDotProductAttention, self).__init__()
4
5         def forward(self, Q, K, V, attn_mask):
6             scores = torch.matmul(Q, K.transpose(-1, -2)) / np.sqrt(d_k)
7             scores.masked_fill_(attn_mask, -1e9)
8             attn = nn.Softmax(dim=-1)(scores)
9             context = torch.matmul(attn, V)
10            return context, attn
11
12 class MultiHeadAttention(nn.Module):
13     def __init__(self):
14         super(MultiHeadAttention, self).__init__()
15         self.W_Q = nn.Linear(d_model, d_k * n_heads)
16         self.W_K = nn.Linear(d_model, d_k * n_heads)
17         self.W_V = nn.Linear(d_model, d_v * n_heads)
18
19     def forward(self, Q, K, V, attn_mask):
20         residual, batch_size = Q, Q.size(0)
21         q_s = self.W_Q(Q).view(batch_size, -1, n_heads, d_k).transpose(1, 2)
22         k_s = self.W_K(K).view(batch_size, -1, n_heads, d_k).transpose(1, 2)
23         v_s = self.W_V(V).view(batch_size, -1, n_heads, d_v).transpose(1, 2)
24
25         attn_mask = attn_mask.unsqueeze(1).repeat(1, n_heads, 1, 1)
26
27         context, attn = ScaledDotProductAttention()(q_s, k_s, v_s, attn_mask)
28         context = context.transpose(1, 2).contiguous().view(batch_size, -1, n_heads * d_v)
29         output = nn.Linear(n_heads * d_v, d_model)(context)
30         return nn.LayerNorm(d_model)(output + residual), attn
31
32 class PoswiseFeedForwardNet(nn.Module):
33     def __init__(self):
34         super(PoswiseFeedForwardNet, self).__init__()
35         self.fc1 = nn.Linear(d_model, d_ff)
36         self.fc2 = nn.Linear(d_ff, d_model)
37
38     def forward(self, x):
39         return self.fc2(gelu(self.fc1(x)))
40
41 def gelu(x):
42     return x * 0.5 * (1.0 + torch.erf(x / math.sqrt(2.0)))

```

executed in 20ms, finished 11:42:02 2020-06-28

## ▼ 4.6 模型定义



```

In [7]: 1 class BERT(nn.Module):
2         def __init__(self):
3             super(BERT, self).__init__()
4             self.embedding = Embedding()
5             self.layers = nn.ModuleList([EncoderLayer() for _ in range(n_layers)])
6             self.fc = nn.Linear(d_model, d_model)
7             self.activ1 = nn.Tanh()
8             self.linear = nn.Linear(d_model, d_model)
9             self.activ2 = gelu
10            self.norm = nn.LayerNorm(d_model)
11            self.classifier = nn.Linear(d_model, 2)
12
13            embed_weight = self.embedding.tok_embed.weight
14            n_vocab, n_dim = embed_weight.size()
15            self.decoder = nn.Linear(n_dim, n_vocab, bias=False)
16            self.decoder.weight = embed_weight
17            self.decoder_bias = nn.Parameter(torch.zeros(n_vocab))
18
19        def forward(self, input_ids, segmetn_ids, masked_pos):
20            output = self.embedding(input_ids, segmetn_ids)
21            enc_self_attn_mask = get_attn_pad_mask(input_ids, input_ids)
22            for layer in self.layers:
23                output, enc_self_attn = layer(output, enc_self_attn_mask)
24            h_pooled = self.activ1(self.fc(output[:, 0]))
25            logits_clsf = self.classifier(h_pooled)
26
27            masked_pos = masked_pos[:, :, None].expand(-1, -1, output.size(-1))
28            h_masked = torch.gather(output, 1, masked_pos)
29            h_masked = self.norm(self.activ2(self.linear(h_masked)))
30            logits_lm = self.decoder(h_masked) + self.decoder_bias
31
32            return logits_lm, logits_clsf
33
34        def get_attn_pad_mask(seq_q, seq_k):
35            batch_size, len_q = seq_q.size()
36            batch_size, len_k = seq_k.size()
37
38            pad_attn_mask = seq_k.data.eq(0).unsqueeze(1)
39            return pad_attn_mask.expand(batch_size, len_q, len_k)

```

executed in 12ms, finished 11:43:33 2020-06-28

## ▼ 4.7 模型训练及测试

```

In [12]: 1 model = BERT()
2 criterion = nn.CrossEntropyLoss()
3 optimizer = optim.Adam(model.parameters(), lr=0.001)
4
5 batch = make_batch()
6 input_ids, segment_ids, masked_tokens, masked_pos, isNext = zip(*batch)
7 input_ids, segment_ids, masked_tokens, masked_pos, isNext = \
8     torch.LongTensor(input_ids), torch.LongTensor(segment_ids), torch.LongTensor(m
9     torch.LongTensor(masked_pos), torch.LongTensor(isNext)
10
11 for epoch in range(100):
12     optimizer.zero_grad()
13     logits_lm, logits_clsfc = model(input_ids, segment_ids, masked_pos)
14     loss_lm = criterion(logits_lm.transpose(1, 2), masked_tokens) # for masked LM
15     loss_lm = (loss_lm.float()).mean()
16     loss_clsfc = criterion(logits_clsfc, isNext) # for sentence classification
17     loss = loss_lm + loss_clsfc
18     if (epoch + 1) % 10 == 0:
19         print('Epoch:', '%04d' % (epoch + 1), 'cost =', '{:.6f}'.format(loss))
20     loss.backward()
21     optimizer.step()
22
23 # Predict mask tokens ans isNext
24 input_ids, segment_ids, masked_tokens, masked_pos, isNext = batch[0]
25 print(text)
26 print([number_dict[w] for w in input_ids if number_dict[w] != '[PAD]'])
27
28 logits_lm, logits_clsfc = model(torch.LongTensor([input_ids]), \
29                                 torch.LongTensor([segment_ids]), torch.LongTensor([m
30 logits_lm = logits_lm.data.max(2)[1][0].data.numpy()
31 print('masked tokens list : ', [pos for pos in masked_tokens if pos != 0])
32 print('predict masked tokens list : ', [pos for pos in logits_lm if pos != 0])
33
34 logits_clsfc = logits_clsfc.data.max(1)[1].data.numpy()[0]
35 print('isNext : ', True if isNext else False)
36 print('predict isNext : ', True if logits_clsfc else False)

```

executed in 1m 41.2s, finished 11:55:23 2020-06-28

```

Epoch: 0010 cost = 23.786695
Epoch: 0020 cost = 16.356310
Epoch: 0030 cost = 15.282166
Epoch: 0040 cost = 11.731420
Epoch: 0050 cost = 4.075703
Epoch: 0060 cost = 5.408222
Epoch: 0070 cost = 7.387890
Epoch: 0080 cost = 6.051405
Epoch: 0090 cost = 7.315189
Epoch: 0100 cost = 4.481243
Hello, how are you? I am Romeo.
Hello, Romeo My name is Juliet. Nice to meet you.
Nice meet you too. How are you today?
Great. My baseball team won the competition.
Oh Congratulations, Juliet
Thanks you Romeo
['[CLS]', 'great', 'my', 'baseball', 'team', 'won', 'the', 'competition', '[SEP]', 'g
reat', '[MASK]', 'baseball', '[MASK]', 'won', 'romeo', 'competition', '[SEP]']
masked tokens list : [23, 28, 27]
predict masked tokens list : [23, 23, 23]
isNext : False
predict isNext : True

```

In [ ]:

1