

STA5007 Assignment 1

Li Xuanran
12312110

October 24, 2025

1 Q1

Tokenization is an important step in modern NLP pipelines. In this task,

1. (3 pts) Explain the Byte Pair Encoding (BPE) algorithm. Then illustrate with code to show how BPE works.
2. (2 pts) Write code to train a BPE model (vocabulary size = 10000) and apply the tokenizer to the dataset. You can use the Subword-NMT toolkit. Please refer to link for the training data. You need to write python code to extract the data from the `conversations` data in the jsonl file and follow the ChatML format.

Note you can refer to the following code MinBPE and Subword-NMT, as well as the original paper Byte pair encoding. You can also ask for help from GPTs, however, you are responsible for checking the correctness of the responses.

Solution.

1. BPE algorithm steps:
 - (a) Initialize the vocabulary from basic tokens (i.e. letter or other kind of chars).
 - (b) Calculate the frequencies of every nearby token-pairs in the corpus.
 - (c) Merge the adjacent tokens which has the highest frequency, and add it to the vocabulary.
 - (d) Follow the step repeatedly, until we reach the limit number of tokens.

The following main functions show how it works.

```
1 # Transfer every text into a list of token with '</w>'.
2 def extract_frequencies(texts):
3     tokens = Counter()
4     for text in texts:
5         text = ' '.join(text) + ' </w>'
6         tokens.update(text.split())
7     return tokens
8
9 # Calculate the frequency of adjacent tokens.
```

```

10 def frequency_of_pairs(frequencies):
11     pairs = Counter()
12     for token, freq in frequencies.items():
13         symbols = token.split()
14         for i in range(len(symbols) - 1):
15             pairs[symbols[i], symbols[i + 1]] += freq
16     return pairs
17
18 # Merge the pair and add it into vocab.
19 def merge_vocab(pair, vocab):
20     bigram = re.escape(' '.join(pair))
21     merged = ' '.join(pair)
22     new_vocab = Counter()
23     for token in vocab:
24         new_token = token.replace(bigram, merged)
25         new_vocab[new_token] = vocab[token]
26     return new_vocab
27
28 # Extract the original frequencies, then merge the most-frequency
    pair iteratively.
29 def encode_with_bpe(texts, num_merges):
30     vocab = extract_frequencies(texts)
31     for _ in range(num_merges):
32         pairs = frequency_of_pairs(vocab)
33         if not pairs:
34             break
35         most_frequent = pairs.most_common(1)[0][0]
36         vocab = merge_vocab(most_frequent, vocab)
37     return vocab

```

Listing 1: BPE Algorithm

2. Simply run `Q1_BPE.py`. The text following ChatML format is stored in `chatml_data.txt`, and the token-frequency pairs are stored in `bpe.codes`.

□

2 Q2

(3 pts) Complete the missing parts of `transformer_model.py`. Search for **TODO** in the code. There are 5 **TODOs** that correspond to key functions. Note that you can ask for help from GPTs, however, you are responsible for checking the correctness of the responses.

Solution.

Positional encoding: We use sin and cos function to set positional vectors:

$$\begin{aligned}
 \text{PE}(\text{pos}, 2i) &= \sin\left(\frac{\text{pos}}{10000^{\frac{2i}{d_m}}}\right), \\
 \text{PE}(\text{pos}, 2i + 1) &= \cos\left(\frac{\text{pos}}{10000^{\frac{2i}{d_m}}}\right),
 \end{aligned}$$

where pos is the position, i is the dimension, d_m is the hidden state dimension of the model. Then we add PE(\cdot) to embedding vectors and then do drop-out, which can let the model identify the order information.

```

1 class PositionalEncoding(nn.Module):
2     def __init__(self, embed_dim: int, max_len: int, dropout: float):
3         super().__init__()
4         self.dropout = nn.Dropout(dropout)
5
6         pe = torch.zeros(max_len, embed_dim)
7         pos = torch.arange(0, max_len, dtype=torch.float).unsqueeze(1)
8         div_term = torch.exp(torch.arange(0, embed_dim, 2).float() * (-
math.log(10000.0) / embed_dim))
9
10        pe[:, 0::2] = torch.sin(pos * div_term)
11        pe[:, 1::2] = torch.cos(pos * div_term)
12        pe = pe.unsqueeze(0)
13        self.register_buffer('pe', pe)
14
15    def forward(self, x):
16        # TODO: Add positional encoding self.pe to input x, then apply
dropout
17        x = x + self.pe[:, :x.size(1)]
18        return self.dropout(x)

```

Listing 2: Positional encoding

In a FFL, for every token \mathbf{x}_i , we have

$$\text{FFN}(\mathbf{x}_i) = \text{Dropout}(\text{ReLU}(\mathbf{x}_i \mathbf{W}_1 + \mathbf{b}_1)) \mathbf{W}_2 + \mathbf{b}_2.$$

```

1 class FeedForward(nn.Module):
2     def __init__(self, embed_dim: int, hidden_dim: int, dropout: float):
3         super().__init__()
4         self.fc1 = nn.Linear(embed_dim, hidden_dim)
5         self.fc2 = nn.Linear(hidden_dim, embed_dim)
6         self.relu = nn.ReLU()
7         self.drop = nn.Dropout(dropout)
8
9     def forward(self, x):
10        # TODO: Implement two linear layers with ReLU and Dropout
11        return self.fc2(self.drop(self.relu(self.fc1(x))))

```

Listing 3: Feed forward

In multi-head self-attention, we first calculate the attention score α :

$$\alpha = \frac{QK^T}{\sqrt{d_k}},$$

then mask out all the coordination that mask = 0, then

$$\text{Attention}(Q, K, V) = \text{Softmax}(\alpha)V.$$

```

1 class MultiHeadAttention(nn.Module):
2     def __init__(self, embed_dim: int, num_heads: int, dropout: float):
3         super().__init__()
4         assert embed_dim % num_heads == 0
5         self.d_k = embed_dim // num_heads
6         self.num_heads = num_heads
7         self.q_linear = nn.Linear(embed_dim, embed_dim)
8         self.k_linear = nn.Linear(embed_dim, embed_dim)
9         self.v_linear = nn.Linear(embed_dim, embed_dim)
10        self.out_linear = nn.Linear(embed_dim, embed_dim)
11        self.drop = nn.Dropout(dropout)
12
13    def forward(self, q, k, v, mask=None):
14        B = q.size(0)
15        Q = self.q_linear(q).view(B, -1, self.num_heads, self.d_k).
16        transpose(1, 2)
17        K = self.k_linear(k).view(B, -1, self.num_heads, self.d_k).
18        transpose(1, 2)
19        V = self.v_linear(v).view(B, -1, self.num_heads, self.d_k).
20        transpose(1, 2)
21
22        # TODO: Compute attention scores
23        scores = Q @ K.transpose(-2, -1) / math.sqrt(self.d_k)
24        if mask is not None:
25            scores = scores.masked_fill(mask == 0, -1e9)
26        attn = torch.softmax(scores, dim=-1)
27        attn = self.drop(attn)
28        out = attn @ V
29        out = out.transpose(1, 2).contiguous().view(B, -1, self.num_heads
30        * self.d_k)
31        return self.out_linear(out)

```

Listing 4: Multi-head attention

Do LayerNorm to the input, then going through the sub-layer (Attention or FFN), and then add the residual:

$$\mathbf{x}^l = \mathbf{x}^{l-1} + \text{block}(\mathbf{x}^{l-1}).$$

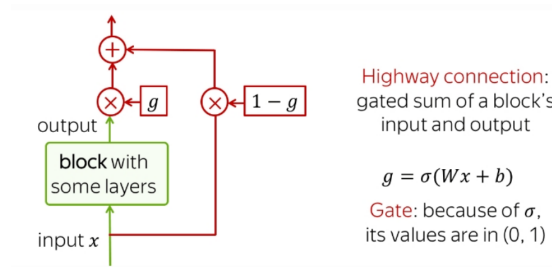


Figure 1: Residual connection

```

1 class Residual(nn.Module):

```

```

2     def __init__(self, dropout: float):
3         super().__init__()
4         self.norm = LayerNorm()
5         self.drop = nn.Dropout(dropout)
6
7     def forward(self, x, sublayer):
8         # TODO: Apply LayerNorm on sublayer(x), then add residual
9         # connection with input x
10        return x + self.drop(sublayer(self.norm(x)))

```

Listing 5: Residual connection

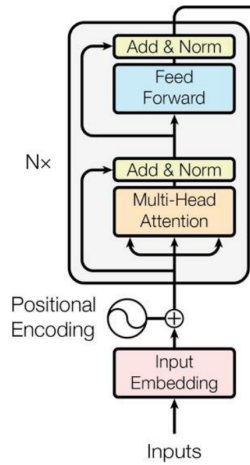


Figure 2: Encoder

According to Figure 2, every Encoder layer has two sub-layers: Layer 1 is self-attention with residual connection, Layer 2 is FFN with residual connection.

```

1 class EncoderLayer(nn.Module):
2     def __init__(self, embed_dim: int, num_heads: int, hidden_dim: int,
3         dropout: float):
4         super().__init__()
5         self.self_attn = MultiHeadAttention(embed_dim, num_heads, dropout)
6         self.ffn = FeedForward(embed_dim, hidden_dim, dropout)
7         self.res_layers = nn.ModuleList([Residual(dropout) for _ in range
8         (2)])
9
10    def forward(self, x, mask):
11        # TODO: Apply self-attention + residual, then feed-forward +
12        # residual
13        x = self.res_layers[0](x, lambda x: self.self_attn(x, x, x, mask))
14        x = self.res_layers[1](x, self.ffn)
15        return x

```

Listing 6: Encoder layer

□

3 Q3

ImageBind-LLM: Multi-modality Instruction Tuning is a paper about large language model. The Figure 3 shows a bind network which consists of some feedforward network.

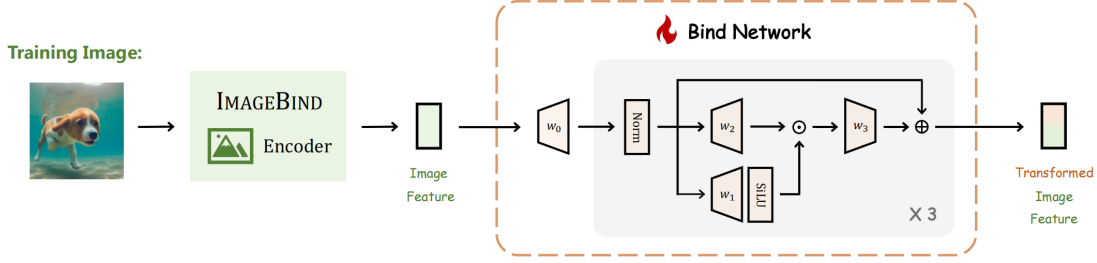


Fig. 3 Details of the Bind Network. Referring to the Feed-Forward Network (FFN) in LLaMA [11], we adopt cascaded blocks of RMSNorm [43], SiLU activation functions [44], and residual connections [45]. This aims to align the image feature from ImageBind [12] with LLaMA’s word embeddings.

Figure 3: Details of the Bind Network

(2 pts) Write the pytorch code for the bind network, the input is `image_feature` and output is `transformed_image_feature`. Define the `__init__()` and `forward()` function. Note that the bind network has three blocks ($\times 3$).

Solution. According to the paper, the main architecture would be:

Input: $F_I \in \mathbb{R}^{1 \times C_I}$: global image feature extracted by the ImageBind encoder.

Output: $T_I \in \mathbb{R}^{1 \times C}$: transformed image feature aligned to LLaMA’s feature dimension C

Structure: Linear projection ($F'_I = F_I W_0$), projecting ($C_I \rightarrow C$). Then pass through 3 repeated blocks, each composed of:

1. RMSNorm
2. Two linear projections (w_1, w_2, w_3)
3. SiLU activation
4. Residual connection

Computation for each block i :

$$F_I^{i+1} = F_I^i + (F_I^i W_2 \cdot \text{SiLU}(F_I^i W_1)) W_3$$

where $W_1, W_2 \in \mathbb{R}^{C \times C_h}$, $W_3 \in \mathbb{R}^{C_h \times C}$.

The `__init__()` and `forward()` function are as below.

```

1 class BindNetwork(nn.Module):
2     def __init__(self, input_dim, output_dim, hidden_dim=None, num_blocks
      =3):
3         super(BindNetwork, self).__init__()

```

```

4     self.input_dim = input_dim
5     self.output_dim = output_dim
6     self.hidden_dim = hidden_dim or output_dim * 4
7     self.num_blocks = num_blocks
8
9     self.proj = nn.Linear(input_dim, output_dim)
10    self.blocks = nn.ModuleList([
11        nn.ModuleDict({
12            "norm": nn.RMSNorm(output_dim),
13            "w1": nn.Linear(output_dim, self.hidden_dim),
14            "w2": nn.Linear(output_dim, self.hidden_dim),
15            "w3": nn.Linear(self.hidden_dim, output_dim)
16        })
17        for _ in range(num_blocks)
18    ])
19
20    def forward(self, image_feature):
21        x = self.proj(image_feature)
22        for blk in self.blocks:
23            h = blk["norm"](x)
24            h1 = blk["w1"](h)
25            h2 = blk["w2"](h)
26            h_mul = h2 * F.silu(h1)
27            h3 = blk["w3"](h_mul)
28            x = x + h3
29        return x

```

Listing 7: Bind Network

□