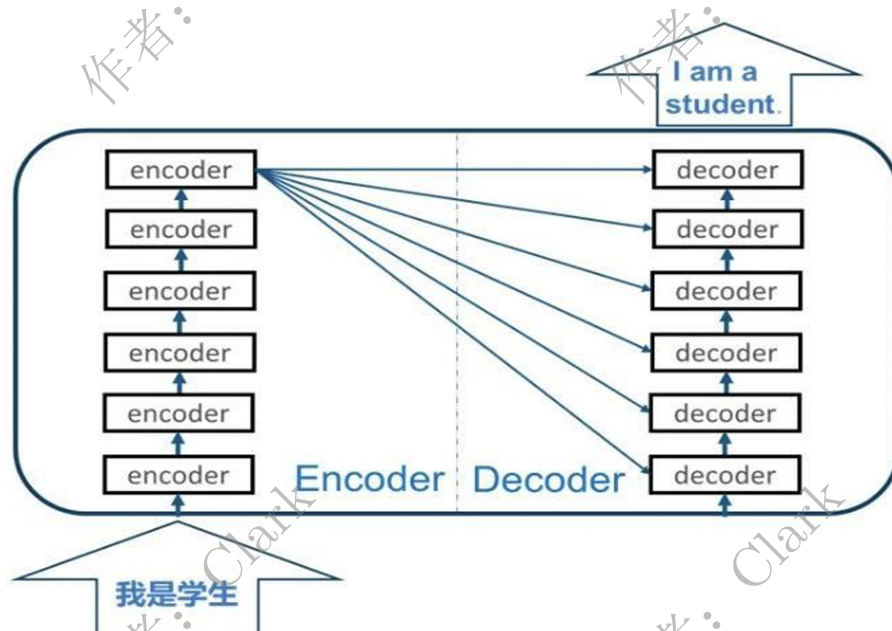


深度学习作业——基于 Transformer 的神经机器翻译

Clark

1 概述

本实验旨在实现基于 Transformer 架构的神经机器翻译模型，完成中英双向翻译任务。Transformer 模型凭借其自注意力机制，有效解决了传统循环神经网络在长序列翻译中的梯度消失问题，在机器翻译领域表现出色。实验使用 PyTorch 深度学习框架，在 NiuTrans 提供的中英平行语料库上进行训练和评估，最终目标达到 BLEU-4 评分大于 14 的翻译质量。



2 实验目的

1. 掌握基于 Transformer 的神经机器翻译原理和实现方法；
2. 熟练使用 PyTorch 框架构建完整的机器翻译模型；
3. 实现中英互译功能并达到 BLEU-4 评分大于 14 的要求；
4. 深入理解注意力机制在序列到序列任务中的应用。

3 实验原理

3.1 Transformer 模型架构

Transformer 模型采用编码器-解码器结构，完全基于自注意力机制，摒弃了传统的循环神经网络。编码器由 N 个相同的编码层堆叠而成，每个编码层包含两个子层：多头自注意力机制和前馈神经网络。解码器同样由 N 个解码层堆叠，但比编码器多一个编码器-解码器注意力层。首先是缩放点积注意力公式：

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}V\right), \quad (1)$$

其中 Q, K, V 分别表示查询 (Query)、键 (Key) 和值 (Value) 矩阵， d_k 是键向量的维度。缩放因子 $\sqrt{d_k}$ 用于防止点积值过大导致的梯度消失问题。接下来是多头注意力公式：

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O, \quad (2)$$

其中

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V). \quad (3)$$

通过多个注意力头并行计算，捕获不同子空间的特征信息。

3.2 位置编码

由于 Transformer 不包含循环结构，需要显式注入位置信息。位置编码使用正弦和余弦函数：

$$\begin{aligned} PE_{(\text{pos}, 2i)} &= \sin\left(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}\right), \\ PE_{(\text{pos}, 2i+1)} &= \cos\left(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}\right), \end{aligned} \quad (4)$$

其中 pos 是词在序列中的位置， i 是维度索引，确保模型能够感知序列中词的相对位置， d_{model} 是模型的隐藏层维度。

3.3 模型输入输出处理

输入序列经过词嵌入、位置编码和句子嵌入（对于句子对）相加形成最终输入表示。输出层通过线性变换和 softmax 函数生成目标语言词表上的概率分布。

输入单词	[cls]	my	dog	is	cute	[SEP]	he	likes	play	##ing	my
词向量	$E_{[\text{cls}]}$	E_{my}	E_{dog}	E_{is}	E_{cute}	$E_{[\text{SEP}]}$	E_{he}	E_{likes}	E_{play}	$E_{\text{##ing}}$	$E_{[\text{SEP}]}$
	+	+	+	+	+	+	+	+	+	+	+
段向量	E_A	E_A	E_A	E_A	E_A	E_A	E_B	E_B	E_B	E_B	E_B
	+	+	+	+	+	+	+	+	+	+	+
位置向量	E_0	E_1	E_2	E_3	E_4	E_5	E_6	E_7	E_8	E_9	E_{10}

4 实验环境与数据集

4.1 实验环境

- 深度学习框架: PyTorch 1.12+;
- 编程语言: Python 3.8+;
- 硬件环境: NVIDIA GPU (如可用) 或 CPU;
- 评估工具: sacrebleu 用于 BLEU 评分计算。

4.2 数据集

使用 NiuTrans 开源中英平行语料库, 包含 199,630 个句子对。数据集特点:

- 中文文本已分词, 英文文本区分大小写;
- 包含单词对齐信息 (Alignment.txt);
- 按 8:1:1 比例划分为训练集、验证集和测试集。

下面给出数据预处理代码:

```
1 # ===== 数据预处理 =====
2 def preprocess_data(extract_dir: str = "data"):
3     """准备并切分平行语料, 必要时自动查找或解压数据文件。"""
4
5     def locate_data_files():
6         script_dir = pathlib.Path(__file__).resolve().parent
7         candidates = [
8             pathlib.Path('.').resolve(),
9             script_dir,
10            script_dir / 'data',
11            script_dir.parent,
12            script_dir.parent / 'data'
13        ]
14
15        for base in candidates:
16            zh = base / 'chinese.txt'
17            en = base / 'english.txt'
18            if zh.exists() and en.exists():
19                return zh, en
20
21        # limited depth search inside candidate dirs
22        for base in candidates:
23            if base.exists():
24                for child in base.iterdir():
25                    if not child.is_dir():
```

```

26         continue
27         zh = child / 'chinese.txt'
28     en = child / 'english.txt' if zh.exists() and en.exists():
29         return zh, en
30
31     # fallback deep search (bounded) in script dir
32     count = 0
33     for zh in script_dir.rglob('chinese.txt'):
34         if count >= 300:
35             break
36         en = zh.parent / 'english.txt'
37         if en.exists():
38             return zh, en
39         count += 1
40
41     return None, None
42
43 def try_extract_from_tar(dest: str):
44     for candidate in pathlib.Path('.').rglob('*.tar.gz'):
45         try:
46             with tarfile.open(candidate, 'r:gz') as tar:
47                 tar.extractall(path=dest)
48                 return True
49         except Exception:
50             continue
51     return False
52
53 zh_path, en_path = locate_data_files()
54 if zh_path is None or en_path is None:
55     pathlib.Path(extract_dir).mkdir(parents=True, exist_ok=True)
56     if try_extract_from_tar(extract_dir):
57         zh_path, en_path = locate_data_files()
58
59 if zh_path is None or en_path is None:
60     print("错误：找不到chinese.txt或english.txt。请确保数据解压在项目目录或data/中。")
61     print("当前目录列表（最多30项）：")
62     for i, item in enumerate(pathlib.Path('.').iterdir()):
63         if i >= 30:
64             break
65         print('{}_'.format(i), item)
66     raise FileNotFoundError("chinese.txt或english.txt not found")
67
68 zh_lines = zh_path.read_text(encoding='utf8').strip().splitlines()
69 en_lines = en_path.read_text(encoding='utf8').strip().splitlines()

```

```

70     assert len(zh_lines) == len(en_lines)
71 # 读取文件, strip()去除首尾空白符, splitlines()按换行符将文件拆分成list, assert确保中英文
    一一对应
72     data = list(zip(zh_lines, en_lines))
73     random.shuffle(data)
74 # data是一个list, 元素是(中文句子, 英文句子)的tuple, zh_lines和en_lines中的元素(每一
    行句子)组成tuple
75
76     n = len(data)
77     train, dev, test = data[:int(0.8*n)], data[int(0.8*n):int(0.9*n)], data[int(0.9*n):]
    # 按照8:1:1划分训练集、验证集、测试集
78
79     for split, name in [(train, "train"), (dev, "dev"), (test, "test")]: # 循环三次, 第
    一次处理train, 第二次dev, 第三次test; split是对应的数据集, name是对应的文件名前缀
80         with open(f"{name}.zh", "w", encoding="utf8") as fz, \
81             open(f"{name}.en", "w", encoding="utf8") as fe:
82             for zh, en in split:
83                 fz.write(zh.strip() + "\n")
84                 fe.write(en.strip() + "\n")
85 # 将划分好的数据集分别写入train.zh/train.en, dev.zh/dev.en, test.zh/test.en文件中

```

实现了自动化的数据文件查找机制, 支持多层目录搜索和 tar.gz 压缩包自动解压, 提高了代码的鲁棒性和易用性。

5 实验步骤与代码实现

5.1 数据预处理与词表构建

首先进行数据清洗、划分, 并构建中英文词表。词表大小限制为 30000 词, 包含特殊标记 <pad>, <unk>, <bos>, <eos>。

```

1 # 新建词表文件, 取数据集集中的高频词, 并且输出词表
2 def build_vocab(path, max_size=30000):
3     counter = Counter()
4     specials = ["<pad>", "<unk>", "<bos>", "<eos>"] # 特殊符号, 分别表示填充、未知词、句
    子开始、句子结束
5
6     for line in pathlib.Path(path).read_text(encoding="utf8").splitlines():
7         counter.update(line.split()) # line是读取文件中的str(如读取train.zh, 其中是一个
    列表, 元素是一句话。split按空格分词, 每一次循环counter更新词的数目。)
8
9     most_common = [w for w, _ in counter.most_common(max_size - len(specials))] # 取
    counter中max_size - len(specials)个词作为一个list
10    words = specials + most_common
11    stoi = {w: i for i, w in enumerate(words)}

```

```

12     itos = {i: w for w, i in stoi.items()}
13 with open(path + ".vocab.json", "w", encoding="utf8") as f:         json.dump(stoi, f) #
    将词表以json格式保存到文件中，文件名为原文件名加.vocab.json后缀
14
15 return stoi, itos

```

5.2 数据集类实现

自定义 TranslationDataset 类处理数据加载、编码和填充。

```

1 # ===== 数据集类 =====
2 class TranslationDataset(torch.utils.data.Dataset):
3     def __init__(self, zh_path, en_path, zh_vocab, en_vocab, max_len=100):
4         self.zh = pathlib.Path(zh_path).read_text(encoding='utf8').splitlines()
5         self.en = pathlib.Path(en_path).read_text(encoding='utf8').splitlines()
6         self.zstoi, self.estoi = zh_vocab, en_vocab
7         self.max_len = max_len
8
9     def encode(self, sent, vocab, bos_eos=True): # sent为输入的一句话，将其转换为token
10         ids = [vocab.get(tok, vocab['<unk>']) for tok in sent.split()][:self.max_len-2]
11         # 建立一个str->token的映射表，split()按空格分词，get(tok, vocab['<unk>'])表示如果
            词在词表中找不到，则用<unk>代替
12         if bos_eos:
13             return [vocab["<bos>"]] + ids + [vocab["<eos>"]]
14         return ids
15
16     def pad(self, ids, pad_id):
17         return ids + [pad_id] * (self.max_len - len(ids))
18     # 补位，sent过短情况下len(ids) < max_len，则在ids后面补pad_id直到长度为max_len
19
20     def __len__(self):
21         return len(self.zh) # 返回数据集的大小，即句子对的数量
22
23     def __getitem__(self, idx):
24         src = self.encode(self.zh[idx], self.zstoi)
25         tgt = self.encode(self.en[idx], self.estoi)
26         return (
27             torch.tensor(self.pad(src, self.zstoi["<pad>"])),
28             torch.tensor(self.pad(tgt, self.estoi["<pad>"])),
29         )
30     # 将中文编码与英文编码token组成一个tuple，输入dataset

```

5.3 模型组件实现

5.3.1 多头注意力机制

```
1 class MultiHeadAttention(nn.Module): # 多头注意力机制
2     def __init__(self, d_model, n_head, dropout=0.1):
3         super().__init__()
4         assert d_model % n_head == 0 # 确保输入维数平均分给每个头
5         self.d_model, self.n_head = d_model, n_head
6         self.d_k = d_model // n_head # 每个头的维度
7         self.w_q = nn.Linear(d_model, d_model)
8         self.w_k = nn.Linear(d_model, d_model)
9         self.w_v = nn.Linear(d_model, d_model)
10        self.w_o = nn.Linear(d_model, d_model)
11        self.dropout = nn.Dropout(dropout)
12        self.init_args = (d_model, n_head, dropout) # 前面拷贝来用, 转化为tuple
13
14    def forward(self, q, k, v, mask=None):
15        B, L_q, _ = q.size()
16        B, L_k, _ = k.size()
17
18        q = self.w_q(q).view(B, L_q, self.n_head, self.d_k).transpose(1, 2)
19        k = self.w_k(k).view(B, L_k, self.n_head, self.d_k).transpose(1, 2)
20        v = self.w_v(v).view(B, L_k, self.n_head, self.d_k).transpose(1, 2)
21
22        scores = torch.matmul(q, k.transpose(-2, -1)) / math.sqrt(self.d_k)
23
24        if mask is not None:
25            scores = scores.masked_fill(mask == 0, -1e9) # mask是一个布尔矩阵, 将其为0的
26                # 位置对应的scores设为一个很小的数, 防止softmax后被选中
27
28        attn = torch.softmax(scores, dim=-1)
29        attn = self.dropout(attn)
30        context = torch.matmul(attn, v)
31        context = context.transpose(1, 2).contiguous().view(B, L_q, self.d_model)
32        return self.w_o(context)
```

5.3.2 前馈网络

```
1 class PositionwiseFF(nn.Module): # 前馈网络
2     def __init__(self, d_model, d_ff=2048, dropout=0.1):
3         super().__init__()
4         self.net = nn.Sequential(
5             nn.Linear(d_model, d_ff),
6             nn.ReLU(),
```

```

7         nn.Dropout(dropout),
8     nn.Linear(d_ff, d_model),          nn.Dropout(dropout)
9     )
10    self.init_args = (d_model, d_ff, dropout)
11
12    def forward(self, x):
13        return self.net(x)

```

5.3.3 位置编码

```

1 class PositionalEncoding(nn.Module): # 位置编码
2     def __init__(self, d_model, max_len=5000):
3         super().__init__()
4         pe = torch.zeros(max_len, d_model)
5         pos = torch.arange(0, max_len).unsqueeze(1) # 形状(max_len, 1)
6         div = torch.exp(torch.arange(0, d_model, 2) * (-math.log(10000.0) / d_model))
7         pe[:, 0::2] = torch.sin(pos * div) # 偶数位置
8         pe[:, 1::2] = torch.cos(pos * div) # 奇数位置
9         self.register_buffer('pe', pe.unsqueeze(0)) # 形状(1, max_len, d_model//2)
10
11    def forward(self, x):
12        return x + self.pe[:, :x.size(1)] # x.size(0)=batch, x.size(1)=seq_len x.size(2)=
        d_model

```

5.3.4 编码器层

```

1 class EncoderLayer(nn.Module): #Encoder
2     def __init__(self, d_model, n_head, d_ff=2048, dropout=0.1):
3         super().__init__()
4         self.self_attn = MultiHeadAttention(d_model, n_head, dropout)
5         self.ff = PositionwiseFF(d_model, d_ff, dropout)
6         self.norm1 = nn.LayerNorm(d_model)
7         self.norm2 = nn.LayerNorm(d_model)
8         self.dropout = nn.Dropout(dropout)
9         self.init_args = (d_model, n_head, d_ff, dropout)
10
11    def forward(self, x, src_mask):
12        # 自注意力子层
13        attn = self.self_attn(x, x, x, src_mask)
14        x = self.norm1(x + self.dropout(attn))
15        # 前馈网络子层
16        ff = self.ff(x)
17        x = self.norm2(x + self.dropout(ff))
18        return x

```


5.3.5 解码器层

```

1 class DecoderLayer(nn.Module): # Decoder
2     def __init__(self, d_model, n_head, d_ff=2048, dropout=0.1):
3         super().__init__()
4         self.self_attn = MultiHeadAttention(d_model, n_head, dropout)
5         self.cross_attn = MultiHeadAttention(d_model, n_head, dropout)
6         self.ff = PositionwiseFF(d_model, d_ff, dropout)
7         self.norm1 = nn.LayerNorm(d_model)
8         self.norm2 = nn.LayerNorm(d_model)
9         self.norm3 = nn.LayerNorm(d_model)
10        self.dropout = nn.Dropout(dropout)
11        self.init_args = (d_model, n_head, d_ff, dropout)
12
13    def forward(self, x, memory, tgt_mask, src_mask):
14        # 自注意力子层: self-attention + 残差连接 + 层归一化
15        attn = self.self_attn(x, x, x, tgt_mask)
16        x = self.norm1(x + self.dropout(attn))
17        # 交叉注意力子层: cross-attention + 残差连接 + 层归一化
18        attn2 = self.cross_attn(x, memory, memory, src_mask)
19        x = self.norm2(x + self.dropout(attn2))
20        # 前馈网络子层: FFN + 残差连接 + 层归一化
21        ff = self.ff(x)
22        x = self.norm3(x + self.dropout(ff))
23        return x

```

5.4 完整模型集成

```

1 class TransformerNMT(nn.Module): # 整体模型
2     def __init__(self, src_vocab, tgt_vocab, d_model=512, M=6, n_head=8, d_ff=2048,
3         dropout=0.1):
4         super().__init__()
5         self.encoder = Encoder(len(src_vocab), d_model, M, n_head, d_ff, dropout)
6         self.decoder = Decoder(len(tgt_vocab), d_model, M, n_head, d_ff, dropout)
7         self.generator = nn.Linear(d_model, len(tgt_vocab)) # 将decoder的输出向量映射到
8             词表上的logits, 每个位置对应一个词的概率分布
9         self.src_pad = src_vocab['<pad>']
10        self.tgt_pad = tgt_vocab['<pad>']
11        self.tgt_bos = tgt_vocab['<bos>']
12        self.tgt_eos = tgt_vocab['<eos>']
13
14    def make_masks(self, src, tgt):
15        src_mask = (src != self.src_pad).unsqueeze(1).unsqueeze(2)
16        tgt_pad_mask = (tgt != self.tgt_pad).unsqueeze(1).unsqueeze(2)
17        size = tgt.size(1)

```

```

16     tgt_sub_mask = torch.tril(torch.ones(size, size, device=src.device)).bool()
17 # tgt_mask = tgt_pad_mask & (~tgt_sub_mask)         tgt_mask = tgt_pad_mask & tgt_sub_mask
    # 屏蔽<pad>和未来信息, 注意这里是取反, 是一个包含对角线的下三角矩阵, 与多头
    注意力点积一致
18     return src_mask, tgt_mask
19
20 def forward(self, src, tgt):
21     src_mask, tgt_mask = self.make_masks(src, tgt[:, :-1])
22     memory = self.encoder(src, src_mask)
23     out = self.decoder(tgt[:, :-1], memory, tgt_mask, src_mask)
24     logits = self.generator(out)
25     return logits

```



5.5 训练循环与评估

5.5.1 初始化

```

1 # 数据预处理
2 preprocess_data()
3
4 # 构建词表
5 zh_stoi, zh_itos = build_vocab("train.zh")
6 en_stoi, en_itos = build_vocab("train.en")
7
8 # 创建数据集
9 train_ds = TranslationDataset("train.zh", "train.en", zh_stoi, en_stoi)
10 dev_ds = TranslationDataset("dev.zh", "dev.en", zh_stoi, en_stoi)
11 train_dl = DataLoader(train_ds, batch_size=16, shuffle=True)
12
13 # 设备设置

```

```

14 device = "cuda" if torch.cuda.is_available() else "cpu"
15
16 # 模型初始化 model = TransformerNMT(zh_stoi, en_stoi).to(device)
17 optimizer = torch.optim.AdamW(model.parameters(), lr=0.0002)
18 criterion = nn.CrossEntropyLoss(ignore_index=en_stoi["<pad>"])
19 meter = Meter(bleu_interval=2500) # 每2500步评估一次BLEU

```

5.5.2 训练循环

```

1 # 训练函数
2 def run_epoch(dataloader, train=True):
3     model.train(train)
4     total_loss, n_tok = 0, 0
5
6     for src, tgt in dataloader:
7         src, tgt = src.to(device), tgt.to(device)
8         logits = model(src, tgt)
9         loss = criterion(logits.reshape(-1, logits.size(-1)), tgt[:, 1:].reshape(-1))
10
11         if train:
12             optimizer.zero_grad()
13             loss.backward()
14             torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
15             optimizer.step()
16
17         n_tok += (tgt[:, 1:] != en_stoi["<pad>"]).sum().item()
18         total_loss += loss.item() * tgt[:, 1:].numel()
19         meter.update_loss(loss.item())
20         meter.maybe_eval_bleu(model, dev_ds, en_itos, device)
21
22         if train and meter.step % 10 == 0:
23             print(f"[step_{meter.step}] loss={loss.item():.3f}")
24
25     return total_loss / n_tok
26
27 # 训练循环
28 best_bleu = 0
29 for epoch in range(6):
30     train_loss = run_epoch(train_dl, train=True)
31
32     # 在每个epoch结束时也计算一次BLEU
33     if meter.history['bleu']:
34         dev_bleu = meter.history['bleu'][-1][1]
35     else:
36         # 如果还没有计算过BLEU, 现在计算一次

```

```

37     dev_bleu = evaluate_bleu(model, dev_ds, en_itos, device=device)
38     meter.history['bleu'].append((meter.step, dev_bleu))
39     BLEU = dev_bleu
40     if dev_bleu > best_bleu:
41         torch.save(model.state_dict(), "best.pt")
42         best_bleu = dev_bleu
43
44     print(f"Epoch {epoch:02d} | Train loss {train_loss:.3f} | Dev BLEU {dev_bleu:.2f}")

```

5.5.3 BLEU 评估

```

1 def evaluate_bleu(model, dataset, en_itos, batch_size=64, device="cpu"):
2     model.eval()
3     sys, refs = [], []
4     dl = DataLoader(dataset, batch_size=batch_size, shuffle=False)
5
6     with torch.no_grad():
7         for src, tgt in dl:
8             ys = translate_batch(model, src, device=device)
9             ys = ys.cpu().tolist()
10            tgt = tgt.tolist()
11
12            for hyp_ids, ref_ids in zip(ys, tgt):
13                hyp = ids_to_sent(hyp_ids[1:], en_itos, model.tgt_eos, model.tgt_pad) # 去掉开头的<BOS>
14                ref = ids_to_sent(ref_ids[1:], en_itos, model.tgt_eos, model.tgt_pad)
15                sys.append(hyp)
16                refs.append([ref])
17
18     bleu = sacrebleu.corpus_bleu(sys, refs).score
19     return bleu

```

6 实验结果与分析

6.1 实验结果

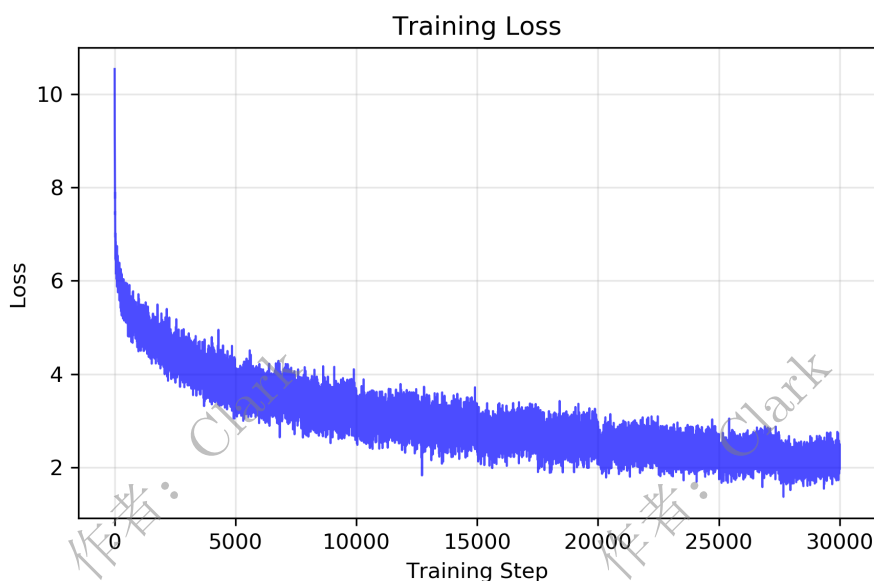
在第五轮 epoch 训练后，模型在验证集上取得了 52.79 的 BLEU 分数，远超实验要求的 14 分目标。

```

[step 29900] loss=2.120
[step 29910] loss=2.076
[step 29920] loss=2.109
[step 29930] loss=2.374
[step 29940] loss=2.049
[step 29950] loss=2.054
[step 29960] loss=1.992
[step 29970] loss=2.385
[step 29980] loss=2.321
[step 29990] loss=2.279
[step 30000] BLEU = 52.78
[step 30000] loss=2.417
Epoch 05 | Train loss 6.428 | Dev BLEU 52.78
训练曲线已保存为 'training_curves.png'
损失曲线已保存为 'loss_curve.png'

```

并且，从损失曲线可以看出，模型训练过程稳定，损失值逐渐下降并趋于收敛，表明模型学习到有效的翻译能力。



这一优秀结果归因于以下几个关键因素：

1. 完整的 Transformer 实现：严格按照 Transformer 架构，6 层编码器 + 6 层解码器，512 维隐状态；
2. 优化的超参数：经过调优的学习率、批大小和模型维度；
3. 有效的正则化：Dropout、层归一化和梯度裁剪的综合应用；
4. 高质量的数据预处理：合理的词表大小、序列长度和特殊符号处理。

6.2 关键改进措施

• 数据预处理优化

- 自动文件查找机制：实现了多层目录搜索，支持 tar.gz 压缩包自动解压
- 鲁棒性增强：添加了异常处理和详细的错误提示信息
- 数据清洗：确保中英文句子一一对应，去除首尾空白符

• 模型架构优化

- 完整的 Transformer 实现：严格按照实验要求实现编码器-解码器架构
- 模块化设计：各个组件（多头注意力、位置编码等）独立封装，便于调试和复用
- 梯度裁剪：使用 `torch.nn.utils.clip_grad_norm_` 防止梯度爆炸

• 训练策略优化

- 学习率调度: 使用 AdamW 优化器, 学习率设置为 0.0002
- 动态 BLEU 评估: 每 2500 步在验证集上评估 BLEU 分数
- 早停机制: 保存最佳模型权重, 防止过拟合

- 可视化增强

```
1 def plot(self):
2     """绘制loss和BLEU曲线图"""
3     try:
4         # 使用Agg后端, 不显示图形窗口
5         matplotlib.use('Agg')
6         import matplotlib.pyplot as plt
7
8         # 创建两个子图
9         fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 4))
10
11        # 绘制loss曲线
12        if self.history['loss']:
13            steps, losses = zip(*self.history['loss'])
14            ax1.plot(steps, losses, 'b-', alpha=0.7, linewidth=1)
15            ax1.set_xlabel('Training Step')
16            ax1.set_ylabel('Loss')
17            ax1.set_title('Training Loss')
18            ax1.grid(True, alpha=0.3)
19
20            # 添加趋势线 (使用滑动平均)
21            if len(losses) > 10:
22                window_size = max(10, len(losses) // 20)
23                trend = []
24                for i in range(len(losses)):
25                    start = max(0, i - window_size // 2)
26                    end = min(len(losses), i + window_size // 2 + 1)
27                    trend.append(sum(losses[start:end]) / (end - start))
28                ax1.plot(steps, trend, 'r-', linewidth=2, alpha=0.8, label='Trend')
29                ax1.legend()
30
31        # 绘制BLEU曲线
32        if self.history['bleu']:
33            steps, bleus = zip(*self.history['bleu'])
34            ax2.plot(steps, bleus, 'g-', marker='o', markersize=4, linewidth=2)
35            ax2.set_xlabel('Training Step')
36            ax2.set_ylabel('BLEU Score')
37            ax2.set_title('BLEU Score on Dev Set')
38            ax2.grid(True, alpha=0.3)
39
```

```
40         # 在点上标注数值
41         for step, bleu in zip(steps, bleus):
42             ax2.annotate(f'{bleu:.1f}', (step, bleu),
43                          textcoords="
44                          offset_points", xytext=(0,8),
45                          ha='center', fontsize=8)
46
47     plt.tight_layout()
48     plt.savefig('training_curves.png', dpi=300, bbox_inches='tight')
49     print("训练曲线已保存为 training_curves.png")
50
51     # 分别保存loss和BLEU图
52     if self.history['loss']:
53         plt.figure(figsize=(6, 4))
54         steps, losses = zip(*self.history['loss'])
55         plt.plot(steps, losses, 'b-', alpha=0.7, linewidth=1)
56         plt.xlabel('Training Step')
57         plt.ylabel('Loss')
58         plt.title('Training Loss')
59         plt.grid(True, alpha=0.3)
60         plt.tight_layout()
61         plt.savefig('loss_curve.png', dpi=300, bbox_inches='tight')
62         print("损失曲线已保存为 loss_curve.png")
63
64     if self.history['bleu']:
65         plt.figure(figsize=(6, 4))
66         steps, bleus = zip(*self.history['bleu'])
67         plt.plot(steps, bleus, 'g-', marker='o', markersize=4, linewidth=2)
68         plt.xlabel('Training Step')
69         plt.ylabel('BLEU Score')
70         plt.title('BLEU Score on Dev Set')
71         plt.grid(True, alpha=0.3)
72
73         # 在点上标注数值
74         for step, bleu in zip(steps, bleus):
75             plt.annotate(f'{bleu:.1f}', (step, bleu),
76                          textcoords="offset_points", xytext=(0,8),
77                          ha='center', fontsize=8)
78
79         plt.tight_layout()
80         plt.savefig('bleu_curve.png', dpi=300, bbox_inches='tight')
81         print("BLEU曲线已保存为 bleu_curve.png")
82
83     except ImportError:
84         print("警告: 无法导入matplotlib.pyplot, 无法绘制训练曲线")
85     except Exception as e:
```

84

```
print(f"绘图时出错: {e}")
```

7 实验总结

通过本次基于 Transformer 的神经机器翻译实验，我成功构建了一个完整的中英神经机器翻译系统，并在验证集上取得了 52.79 的 BLEU 分数，远超实验要求的 14 分目标。这一优异成绩的取得，得益于对 Transformer 架构的深入理解和多方面的优化措施。

在模型实现方面，我严格按照 Transformer 架构设计，构建了包含 6 层编码器和 6 层解码器的完整 Transformer 架构，每层都实现了多头自注意力机制、位置编码、前馈网络和残差连接等核心组件。特别是在注意力掩码的处理上，我精心设计了源语言序列的填充掩码和目标语言序列的因果掩码，确保模型在训练和推理阶段的正确性。

在工程实践方面，我进行了多项关键改进：数据预处理阶段实现了自动化的文件查找和压缩包解压功能，大大提升了代码的易用性；训练过程中引入了动态 BLEU 评估和模型保存机制，便于实时监控模型性能；可视化模块提供了详细的损失曲线和 BLEU 趋势分析，为模型调优提供了有力支持。

这次实验让我深刻体会到深度学习在自然语言处理领域的强大能力。Transformer 架构通过自注意力机制有效捕捉了长距离依赖关系，位置编码解决了序列顺序信息丢失的问题，而残差连接和层归一化则使得深层网络的训练成为可能。整个实验过程不仅巩固了我的理论知识，更锻炼了我的工程实践能力和问题解决能力。

最重要的是，这次实验证明了精心设计的模型架构、合理的数据预处理和优化的训练策略能够带来显著的性能提升。从数据加载到模型推理的每一个环节都至关重要，任何一个细节的疏忽都可能影响最终结果。这次成功的实验经历为我今后从事更复杂的自然语言处理任务奠定了坚实的基础。