

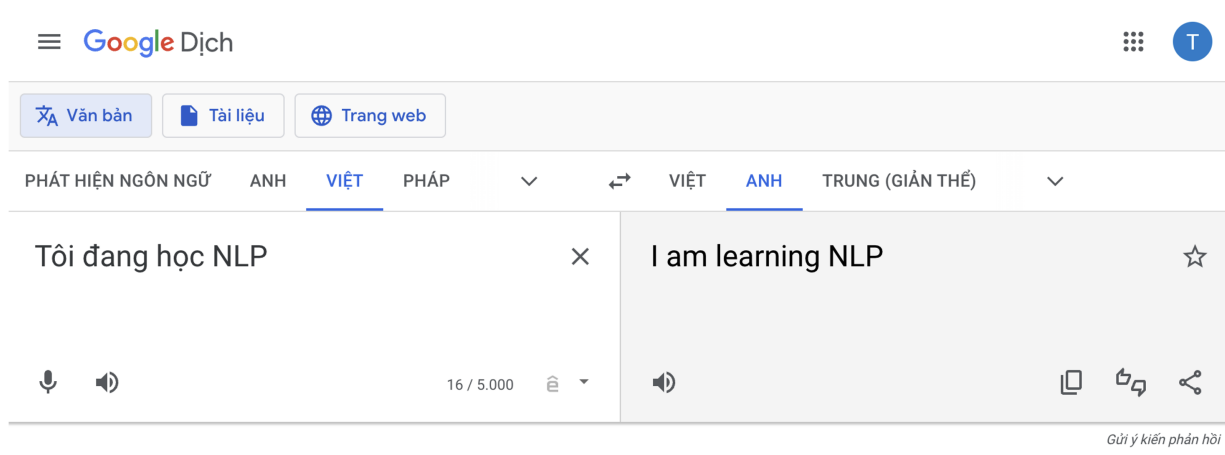
Project: Neural Machine Translation

Quoc-Thai Nguyen và Quang-Vinh Dinh

PR-Team: Minh-Châu Phạm, Hoàng-Nguyên Vũ và Đăng-Nhã Nguyễn

Ngày 19 tháng 2 năm 2024

Phần I. Giới thiệu



Hình 1: Hệ thống dịch máy Google.

Dịch Máy (Machine Translation) với mục đích tự động dịch văn bản hoặc lời nói từ ngôn ngữ tự nhiên này sang ngôn ngữ tự nhiên khác. Dịch máy sử dụng kết hợp nhiều ý tưởng và các kỹ thuật với nhau từ ngôn ngữ học, khoa học máy tính, xác suất thống kê và trí tuệ nhân tạo. Với mục tiêu của dịch máy là phát triển một hệ thống cho phép tạo ra bản dịch chính xác giữa các ngôn ngữ tự nhiên của con người.

Các hệ thống dịch máy điển hình hiện nay như: Google Translate, Bing Translator,... đã đạt được chất lượng bản dịch tốt và được tích hợp trong nhiều nền tảng ứng dụng khác nhau và có thể dịch tốt giữa hơn 100 các ngôn ngữ tự nhiên.

Như vậy, Input/Output của bài toán là:

- **Input:** Văn bản đầu vào của ngôn ngữ nguồn.
VD: Câu đầu vào tiếng việt: "Tôi đang học NLP"
- **Output:** Văn bản được dịch sang ngôn ngữ đích.
VD: Câu được dịch sang tiếng anh: " I am learning NLP"

Mô hình hoá bài toán dịch máy, với mục tiêu là huấn luyện và tối ưu các tham số mô hình θ với đầu vào văn bản từ ngôn ngữ nguồn $w^{(s)}$ và đầu ra văn bản từ ngôn ngữ đích tương ứng $w^{(t)}$:

$$\hat{w}^{(t)} = \underset{w^{(t)}}{\operatorname{argmax}} \theta(w^{(s)}, w^{(t)})$$

Vì vậy, để giải quyết tốt bài toán dịch máy, chúng ta cần quan tâm tối ưu:

- **Phần 1:** Thuật toán học tối ưu bộ tham số θ
- **Phần 2:** Thuật toán giải mã (decoding) để sinh ra bản dịch tốt nhất cho văn bản đầu vào

Hiện nay, có ba hướng tiếp cận chính cho bài toán này:

- **Hướng 1:** Dịch máy dựa vào luật (Rule-based Machine Translation - RBMT)
- **Hướng 2:** Dịch máy dựa vào thống kê (Statistical Machine Translation - SMT)
- **Hướng 3:** Dịch máy dựa vào mạng nơ-ron (Neural Machine Translation - NMT)

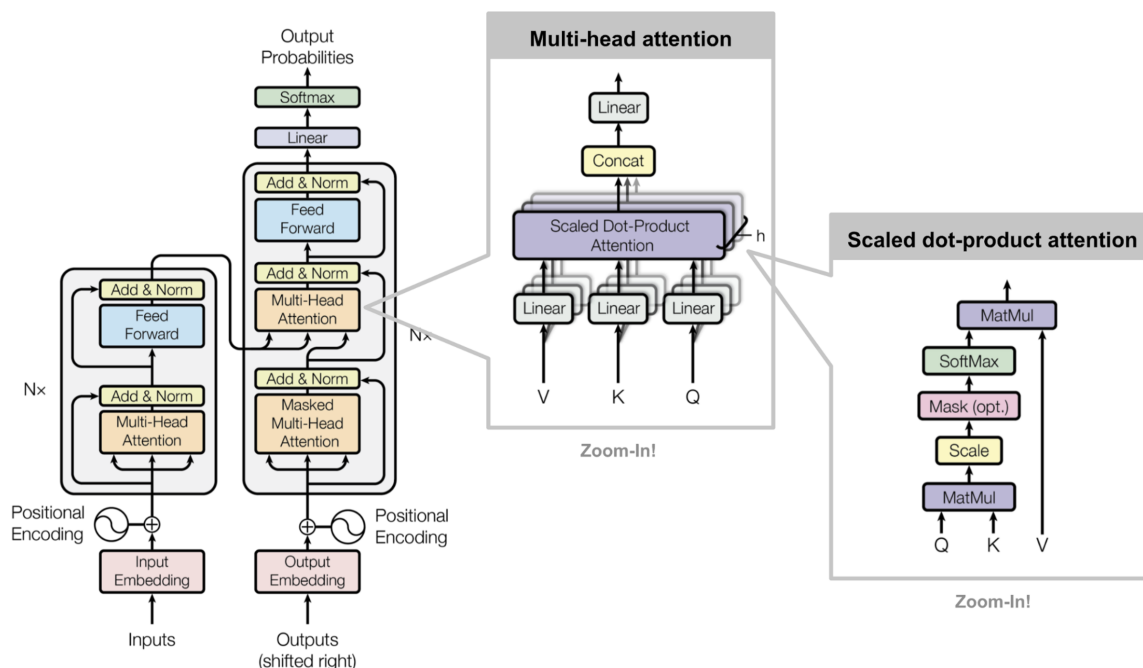
Trong các hướng tiếp cận này NMT đang ngày càng phát triển và cho chất lượng bản dịch tốt vượt trội. Vì vậy, phạm vi project tập trung vào các phương pháp dựa trên mạng nơ-ron gồm 2 nội dung chính:

- **Phương pháp 1:** Xây dựng mô hình dịch máy sử dụng kiến trúc Transformer
- **Phương pháp 2:** Xây dựng mô hình dịch máy sử dụng Pre-trained Language Model như BERT và GPT

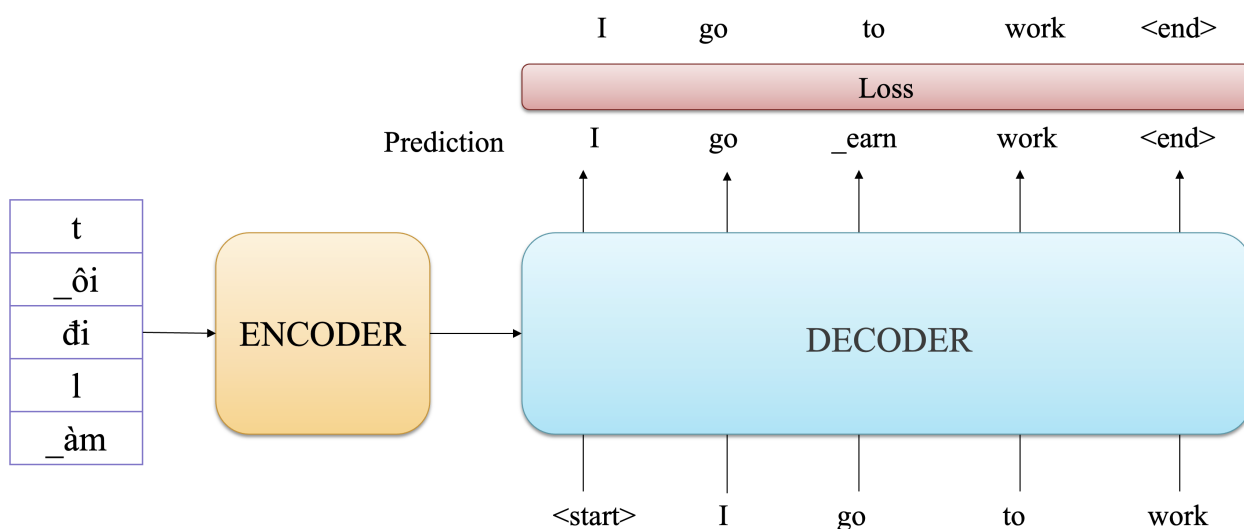
Phần II. Machine Translation using Transformer Model

Để xây dựng, đánh giá hiệu suất các mô hình chúng ta sử dụng bộ dữ liệu dịch **IWSLT'15 English -Vietnamese** với số lượng mẫu cho training: 133,317 cặp câu song ngữ, tập validation: 1,553 cặp câu song ngữ và tập test: 1,269 cặp câu song ngữ. Chiều dịch sẽ từ tiếng anh sang tiếng việt.

Metric để đánh giá chúng ta sử dụng: ScacreBleu (BLEU)



Hình 2: Mô hình Transformer.



Hình 3: Dịch máy sử dụng mô hình Transformer.

1. Build Dataset

```

1 # import libs
2 !pip install -q datasets sacrebleu
3
4 # download dataset
5 from datasets import load_dataset
6
7 data = load_dataset(
8     "mt_eng_vietnamese",
9     "iwslt2015-en-vi"
10 )
11
12 # tokenization
13 from torchtext.data.utils import get_tokenizer
14 from torchtext.vocab import build_vocab_from_iterator
15
16 SRC_LANGUAGE = 'en'
17 TGT_LANGUAGE = 'vi'
18
19 token_transform = {}
20 vocab_transform = {}
21
22 token_transform[SRC_LANGUAGE] = get_tokenizer('basic_english')
23 token_transform[TGT_LANGUAGE] = get_tokenizer('basic_english')
24
25 UNK_IDX, PAD_IDX, BOS_IDX, EOS_IDX = 0, 1, 2, 3
26 special_symbols = ['<unk>', '<pad>', '<bos>', '<eos>']
27
28 def yield_tokens(data_iter, lang):
29     for data_sample in data_iter['translation']:
30         yield token_transform[lang](data_sample[lang])
31
32
33 for lang in [SRC_LANGUAGE, TGT_LANGUAGE]:
34     train_iter = data['train']
35
36     # Create torchtext's Vocab object
37     vocab_transform[lang] = build_vocab_from_iterator(
38         yield_tokens(train_iter, lang),
39         min_freq=1,
40         specials=special_symbols,
41         special_first=True
42     )
43
44     vocab_transform[lang].set_default_index(UNK_IDX)
45
46 # dataloader
47 import torch
48 from torch.nn.utils.rnn import pad_sequence
49
50 # helper function to club together sequential operations
51 def sequential_transforms(*transforms):
52     def func(txt_input):
53         for transform in transforms:
54             txt_input = transform(txt_input)
55         return txt_input
56     return func
57
58 # function to add BOS/EOS and create tensor for input sequence indices

```

```

59 def tensor_transform(token_ids):
60     return torch.cat((torch.tensor([BOS_IDX]),
61                          torch.tensor(token_ids),
62                          torch.tensor([EOS_IDX])))
63
64 # 'src' and 'tgt' language text transforms to convert raw strings into tensors
65 indices
66 text_transform = {}
67 for lang in [SRC_LANGUAGE, TGT_LANGUAGE]:
68     text_transform[lang] = sequential_transforms(
69         token_transform[lang], # Tokenization
70         vocab_transform[lang], # Numericalization
71         tensor_transform # Add BOS/EOS and create tensor
72     )
73 # function to collate data samples into batch tensors
74 def collate_fn(batch):
75     src_batch, tgt_batch = [], []
76     for sample in batch:
77         src_sample, tgt_sample = sample[SRC_LANGUAGE], sample[TGT_LANGUAGE]
78         src_batch.append(text_transform[SRC_LANGUAGE](src_sample).to(dtype=torch.int64))
79         tgt_batch.append(text_transform[TGT_LANGUAGE](tgt_sample).to(dtype=torch.int64))
80
81     src_batch = pad_sequence(src_batch, padding_value=PAD_IDX, batch_first=True)
82     tgt_batch = pad_sequence(tgt_batch, padding_value=PAD_IDX, batch_first=True)
83     return src_batch, tgt_batch
84
85 from torch.utils.data import DataLoader
86
87 BATCH_SIZE = 8
88
89 train_dataloader = DataLoader(
90     data['train']['translation'],
91     batch_size=BATCH_SIZE,
92     collate_fn=collate_fn
93 )
94
95 valid_dataloader = DataLoader(
96     data['validation']['translation'],
97     batch_size=BATCH_SIZE,
98     collate_fn=collate_fn
99 )
100
101 test_dataloader = DataLoader(
102     data['test']['translation'],
103     batch_size=BATCH_SIZE,
104     collate_fn=collate_fn
105 )

```

2. Modeling

```

1 from torch import Tensor
2 import torch
3 import torch.nn as nn
4 from torch.nn import Transformer
5 import math
6 DEVICE = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
7

```

```

8 # helper Module that adds positional encoding to the token embedding to introduce a
  notion of word order.
9 class PositionalEncoding(nn.Module):
10     def __init__(self,
11                   emb_size: int,
12                   dropout: float,
13                   maxlen: int = 5000):
14         super(PositionalEncoding, self).__init__()
15         den = torch.exp(- torch.arange(0, emb_size, 2)* math.log(10000) / emb_size)
16         pos = torch.arange(0, maxlen).reshape(maxlen, 1)
17         pos_embedding = torch.zeros((maxlen, emb_size))
18         pos_embedding[:, 0::2] = torch.sin(pos * den)
19         pos_embedding[:, 1::2] = torch.cos(pos * den)
20         pos_embedding = pos_embedding.unsqueeze(-2)
21
22         self.dropout = nn.Dropout(dropout)
23         self.register_buffer('pos_embedding', pos_embedding)
24
25     def forward(self, token_embedding: Tensor):
26         return self.dropout(token_embedding + self.pos_embedding[:token_embedding.size
27 (0), :])
28 # helper Module to convert tensor of input indices into corresponding tensor of token
  embeddings
29 class TokenEmbedding(nn.Module):
30     def __init__(self, vocab_size: int, emb_size):
31         super(TokenEmbedding, self).__init__()
32         self.embedding = nn.Embedding(vocab_size, emb_size)
33         self.emb_size = emb_size
34
35     def forward(self, tokens: Tensor):
36         return self.embedding(tokens.long()) * math.sqrt(self.emb_size)
37
38 # Seq2Seq Network
39 class Seq2SeqTransformer(nn.Module):
40     def __init__(self,
41                   num_encoder_layers: int,
42                   num_decoder_layers: int,
43                   emb_size: int,
44                   nhead: int,
45                   src_vocab_size: int,
46                   tgt_vocab_size: int,
47                   dim_feedforward: int = 512,
48                   dropout: float = 0.1):
49         super(Seq2SeqTransformer, self).__init__()
50         self.transformer = Transformer(d_model=emb_size,
51                                       nhead=nhead,
52                                       num_encoder_layers=num_encoder_layers,
53                                       num_decoder_layers=num_decoder_layers,
54                                       dim_feedforward=dim_feedforward,
55                                       dropout=dropout,
56                                       batch_first=True)
57         self.generator = nn.Linear(emb_size, tgt_vocab_size)
58         self.src_tok_emb = TokenEmbedding(src_vocab_size, emb_size)
59         self.tgt_tok_emb = TokenEmbedding(tgt_vocab_size, emb_size)
60         self.positional_encoding = PositionalEncoding(
61             emb_size, dropout=dropout)
62
63     def forward(self,
64               src: Tensor,

```

```

65         trg: Tensor,
66         src_mask: Tensor,
67         tgt_mask: Tensor,
68         src_padding_mask: Tensor,
69         tgt_padding_mask: Tensor,
70         memory_key_padding_mask: Tensor):
71     src_emb = self.positional_encoding(self.src_tok_emb(src))
72     tgt_emb = self.positional_encoding(self.tgt_tok_emb(trg))
73     outs = self.transformer(src_emb, tgt_emb, src_mask, tgt_mask, None,
74                             src_padding_mask, tgt_padding_mask,
memory_key_padding_mask)
75     return self.generator(outs)
76
77     def encode(self, src: Tensor, src_mask: Tensor):
78         return self.transformer.encoder(self.positional_encoding(
79             self.src_tok_emb(src)), src_mask)
80
81     def decode(self, tgt: Tensor, memory: Tensor, tgt_mask: Tensor):
82         return self.transformer.decoder(self.positional_encoding(
83             self.tgt_tok_emb(tgt)), memory,
84             tgt_mask)
85
86     def generate_square_subsequent_mask(sz):
87         mask = (torch.triu(torch.ones((sz, sz), device=DEVICE)) == 1).transpose(0, 1)
88         mask = mask.float().masked_fill(mask == 0, float('-inf')).masked_fill(mask == 1,
float(0.0))
89         return mask
90
91
92     def create_mask(src, tgt):
93         src_seq_len = src.shape[1]
94         tgt_seq_len = tgt.shape[1]
95
96         tgt_mask = generate_square_subsequent_mask(tgt_seq_len)
97         src_mask = torch.zeros((src_seq_len, src_seq_len), device=DEVICE).type(torch.bool)
98
99         src_padding_mask = (src == PAD_IDX)
100        tgt_padding_mask = (tgt == PAD_IDX)
101        return src_mask, tgt_mask, src_padding_mask, tgt_padding_mask

```

3. Trainer

```

1  import time
2
3  def train_epoch(model, optimizer, criterion, train_dataloader, device):
4      model.train()
5      losses = []
6
7      for src_ids, tgt_ids in train_dataloader:
8          src_ids = src_ids.to(device)
9          tgt_ids = tgt_ids.to(device)
10
11         tgt_input = tgt_ids[:, :-1]
12         tgt_output = tgt_ids[:, 1:]
13
14         src_mask, tgt_mask, src_padding_mask, tgt_padding_mask = create_mask(src_ids,
tgt_input)
15         try:
16             output = model(
17                 src_ids, tgt_input, src_mask, tgt_mask, src_padding_mask,

```

```

18         tgt_padding_mask, src_padding_mask
19     )
20     except:
21         print(src_ids.shape, tgt_input.shape)
22
23     optimizer.zero_grad()
24
25     loss = criterion(
26         output.reshape(-1, output.shape[-1]),
27         tgt_output.reshape(-1))
28     loss.backward()
29
30     optimizer.step()
31     losses.append(loss.item())
32
33     return sum(losses) / len(losses)
34
35 def evaluate(model, data_loader, criterion, device):
36     model.eval()
37     losses = []
38     with torch.no_grad():
39         for src_ids, tgt_ids in data_loader:
40             src_ids = src_ids.to(device)
41             tgt_ids = tgt_ids.to(device)
42
43             tgt_input = tgt_ids[:, :-1]
44             tgt_output = tgt_ids[:, 1:]
45
46             src_mask, tgt_mask, src_padding_mask, tgt_padding_mask = create_mask(
47                 src_ids, tgt_input)
48             output = model(
49                 src_ids, tgt_input, src_mask, tgt_mask, src_padding_mask,
50                 tgt_padding_mask, src_padding_mask)
51             loss = criterion(
52                 output.reshape(-1, output.shape[-1]),
53                 tgt_output.reshape(-1))
54             losses.append(loss.item())
55
56     return sum(losses) / len(losses)
57
58 def train(model, train_dataloader, valid_dataloader, optimizer, criterion, device,
59           epochs):
60     for epoch in range(1, epochs+1):
61         start_time = time.time()
62         train_loss = train_epoch(model, optimizer, criterion, train_dataloader, device)
63
64         valid_loss = evaluate(model, valid_dataloader, criterion, device)
65         end_time = time.time()
66         print(f"Epoch: {epoch}, Train loss: {train_loss:.3f}, Val loss: {valid_loss:.3f}, "
67               f"Epoch time = {(end_time - start_time):.3f}s")

```

4. Training

```

1 SRC_VOCAB_SIZE = len(vocab_transform[SRC_LANGUAGE])
2 TGT_VOCAB_SIZE = len(vocab_transform[TGT_LANGUAGE])
3 EMB_SIZE = 512
4 NHEAD = 8
5 FFN_HID_DIM = 512
6 BATCH_SIZE = 128

```



```

7 NUM_ENCODER_LAYERS = 3
8 NUM_DECODER_LAYERS = 3
9
10 transformer = Seq2SeqTransformer(NUM_ENCODER_LAYERS, NUM_DECODER_LAYERS, EMB_SIZE,
11                                NHEAD, SRC_VOCAB_SIZE, TGT_VOCAB_SIZE, FFN_HID_DIM)
12 transformer = transformer.to(DEVICE)
13
14 criterion = torch.nn.CrossEntropyLoss(ignore_index=PAD_IDX)
15
16 optimizer = torch.optim.Adam(transformer.parameters(), lr=0.0001, betas=(0.9, 0.98),
17                               eps=1e-9)
18
19 epochs = 5
20 train(transformer, train_dataloader, valid_dataloader, optimizer, criterion, DEVICE,
21        epochs)

```

5. Evaluation

```

1 def greedy_decode(model, src, src_mask, max_len, start_symbol):
2     src = src.to(DEVICE)
3     src_mask = src_mask.to(DEVICE)
4
5     memory = model.encode(src, src_mask)
6     ys = torch.ones(1, 1).fill_(start_symbol).type(torch.long).to(DEVICE)
7     for i in range(max_len-1):
8         memory = memory.to(DEVICE)
9         tgt_mask = (generate_square_subsequent_mask(ys.size(1))
10                  .type(torch.bool)).to(DEVICE)
11         out = model.decode(ys, memory, tgt_mask)
12         out = out.transpose(0, 1)
13         prob = model.generator(out[:, -1])
14         _, next_word = torch.max(prob, dim=1)
15         next_word = next_word[-1].item()
16
17         ys = torch.cat([ys,
18                       torch.ones(1, 1).type_as(src.data).fill_(next_word)], dim=1)
19         if next_word == EOS_IDX:
20             break
21     return ys
22
23
24 # actual function to translate input sentence into target language
25 def translate(model: torch.nn.Module, src_sentence: str):
26     model.eval()
27     src = text_transform[SRC_LANGUAGE](src_sentence).view(1, -1)
28     num_tokens = src.shape[1]
29     src_mask = (torch.zeros(num_tokens, num_tokens)).type(torch.bool)
30     tgt_tokens = greedy_decode(
31         model, src, src_mask, max_len=num_tokens + 5, start_symbol=BOS_IDX).flatten()
32     return " ".join(
33         vocab_transform[TGT_LANGUAGE].lookup_tokens(list(tgt_tokens.cpu().numpy()))).
34         replace("<bos>", "").replace("<eos>", "")
35
36 translate(transformer, "i go to school") # => toi den truong
37
38 # evaluate on test set
39 from tqdm import tqdm
40 import sacrebleu
41
42 pred_sentences, tgt_sentences = [], []
43 for sample in tqdm(data['test']['translation']):

```

```
42     src_sentence = sample[SRC_LANGUAGE]
43     tgt_sentence = sample[TGT_LANGUAGE]
44
45     pred_sentence = translate(transformer, src_sentence)
46     pred_sentences.append(pred_sentence)
47
48     tgt_sentences.append(tgt_sentence)
49
50 bleu_score = sacrebleu.corpus_bleu(pred_sentences, [tgt_sentences], force=True)
51 bleu_score => 46.8/16.8/6.4/2.4
```

Phần III. Machine Translation using Pre-trained LMs

Các pre-trained language model dựa vào kiến trúc transformer được huấn luyện trên tập dữ liệu lớn. Vì vậy các mô hình ngôn ngữ đạt hiệu quả tốt cho các tác vụ NLU (Natural Language Understanding). Trong đó điển hình là BERT cốt lõi là các khối Transformer-Encoder và GPT (Cụ thể trong project này sẽ sử dụng GPT2) cốt lõi là các khối Transformer-Decoder. Dựa trên kiến trúc đó, chúng ta hoàn toàn có thể sử dụng trong số các mô hình như BERT và GPT2 để khởi tạo trọng số cho mô hình Transformer.

Ở trong phần này, chúng ta sẽ sử dụng các pre-trained LMs model là BERT và GPT2 để khởi tạo cho Transformer thông qua thư viện "transformers" của huggingface. Gồm có 2 cách chính:

- **Cách 1:** BERT cho Transformer-Encoder và BERT cho Transformer-Decoder
- **Cách 2:** BERT cho Transformer-Encoder và GPT2 cho Transformer-Decoder

1. Build Dataset

```

1 # install libs
2 !pip install -q datasets sacrebleu accelerate>=0.20.1
3
4 # import libs
5 import os
6 import numpy as np
7 import sacrebleu
8 import torch
9 from torch.utils.data import Dataset
10 from datasets import load_dataset, load_metric
11 from transformers import *
12
13 # load dataset
14 class NMTDataset(Dataset):
15     def __init__(self, cfg, data_type="train"):
16         super().__init__()
17         self.cfg = cfg
18
19         self.src_texts, self.tgt_texts = self.read_data(data_type)
20
21         self.src_input_ids, self.src_attention_mask = self.texts_to_sequences(self.
src_texts)
22         self.tgt_input_ids, self.tgt_attention_mask, self.labels = self.
texts_to_sequences(
23             self.tgt_texts,
24             is_src=False
25         )
26
27     def read_data(self, data_type):
28         data = load_dataset(
29             "mt_eng_vietnamese",
30             "iwslt2015-en-vi",
31             split=data_type
32         )
33         src_texts = [sample["translation"][self.cfg.src_lang] for sample in data]
34         tgt_texts = [sample["translation"][self.cfg.tgt_lang] for sample in data]
35         return src_texts, tgt_texts

```

```

36
37 def texts_to_sequences(self, texts, is_src=True):
38     if is_src:
39         src_inputs = self.cfg.src_tokenizer(
40             texts,
41             padding='max_length',
42             truncation=True,
43             max_length=self.cfg.src_max_len,
44             return_tensors='pt'
45         )
46         return (
47             src_inputs.input_ids,
48             src_inputs.attention_mask
49         )
50
51     else:
52         if self.cfg.add_special_tokens:
53             texts = [
54                 ' '.join([
55                     self.cfg.tgt_tokenizer.bos_token,
56                     text,
57                     self.cfg.tgt_tokenizer.eos_token
58                 ])
59                 for text in texts
60             ]
61             tgt_inputs = self.cfg.tgt_tokenizer(
62                 texts,
63                 padding='max_length',
64                 truncation=True,
65                 max_length=self.cfg.tgt_max_len,
66                 return_tensors='pt'
67             )
68
69             labels = tgt_inputs.input_ids.numpy().tolist()
70             labels = [
71                 [
72                     -100 if token_id == self.cfg.tgt_tokenizer.pad_token_id else
token_id
73                     for token_id in label
74                 ]
75                 for label in labels
76             ]
77
78             labels = torch.LongTensor(labels)
79
80             return (
81                 tgt_inputs.input_ids,
82                 tgt_inputs.attention_mask,
83                 labels
84             )
85
86 def __getitem__(self, idx):
87     return {
88         "input_ids": self.src_input_ids[idx],
89         "attention_mask": self.src_attention_mask[idx],
90         "decoder_input_ids": self.tgt_input_ids[idx],
91         "decoder_attention_mask": self.tgt_attention_mask[idx],
92         "labels": self.labels[idx]
93     }
94

```

```

95     def __len__(self):
96         return np.shape(self.src_input_ids)[0]

```

2. Tokenizer

```

1 def postprocess_text(preds, labels):
2     preds = [pred.strip() for pred in preds]
3     labels = [[label.strip()] for label in labels]
4     return preds, labels
5
6 def load_tokenizer(model_name_or_path):
7     if 'bert' in model_name_or_path.split('-'):
8         return BertTokenizerFast.from_pretrained(model_name_or_path)
9     elif 'gpt2' in model_name_or_path.split('-'):
10        return GPT2TokenizerFast.from_pretrained(model_name_or_path)
11    else:
12        return AutoTokenizer.from_pretrained(model_name_or_path)

```

3. Trainer

```

1 class Manager():
2     def __init__(self, cfg, is_train=True):
3         self.cfg = cfg
4
5         print("Loading Tokenizer...")
6         self.get_tokenizer()
7
8         print("Loading Model...")
9         self.get_model()
10
11        print("Loading Metric...")
12        self.bleu_metric = load_metric("sacrebleu")
13
14        print("Check Save Model Path")
15        if not os.path.exists(self.cfg.ckpt_dir):
16            os.mkdir(self.cfg.ckpt_dir)
17
18        if is_train:
19            # Load dataloaders
20            print("Loading Dataset...")
21            self.train_dataset = NMTDataset(self.cfg, data_type="train")
22            self.valid_dataset = NMTDataset(self.cfg, data_type="validation")
23
24        print("Setting finished.")
25
26    def get_tokenizer(self):
27        if self.cfg.load_model_from_path:
28            self.cfg.src_tokenizer = load_tokenizer(self.cfg.ckpt_dir)
29            self.cfg.tgt_tokenizer = load_tokenizer(self.cfg.ckpt_dir)
30        else:
31            self.cfg.src_tokenizer = load_tokenizer(self.cfg.src_model_name)
32            self.cfg.tgt_tokenizer = load_tokenizer(self.cfg.tgt_model_name)
33            if "bert" in self.cfg.tgt_model_name.split('-'):
34                self.cfg.add_special_tokens = False
35                self.cfg.bos_token_id = self.cfg.tgt_tokenizer.cls_token_id
36                self.cfg.eos_token_id = self.cfg.tgt_tokenizer.sep_token_id
37                self.cfg.pad_token_id = self.cfg.tgt_tokenizer.pad_token_id
38            else:
39                self.cfg.add_special_tokens = True
40                self.cfg.tgt_tokenizer.add_special_tokens(

```

```

41         {
42             "bos_token": "[BOS]",
43             "eos_token": "[EOS]",
44             "pad_token": "[PAD]"
45         }
46     )
47     self.cfg.bos_token_id = self.cfg.tgt_tokenizer.bos_token_id
48     self.cfg.eos_token_id = self.cfg.tgt_tokenizer.eos_token_id
49     self.cfg.pad_token_id = self.cfg.tgt_tokenizer.pad_token_id
50     self.cfg.src_tokenizer.save_pretrained(
51         os.path.join(self.cfg.ckpt_dir, f"{self.cfg.src_lang}_tokenizer_{cfg.
src_model_name}")
52     )
53
54     self.cfg.tgt_tokenizer.save_pretrained(
55         os.path.join(self.cfg.ckpt_dir, f"{self.cfg.tgt_lang}_tokenizer_{cfg.
tgt_model_name}")
56     )
57
58     def get_model(self):
59         if self.cfg.load_model_from_path:
60             save_model_path = os.path.join(self.cfg.ckpt_dir, self.cfg.ckpt_name)
61             self.model = EncoderDecoderModel.from_pretrained(save_model_path)
62         else:
63             self.model = EncoderDecoderModel.from_encoder_decoder_pretrained(
64                 self.cfg.src_model_name,
65                 self.cfg.tgt_model_name
66             )
67             self.model.decoder.resize_token_embeddings(len(self.cfg.tgt_tokenizer))
68             self.model.config.decoder_start_token_id = self.cfg.bos_token_id
69             self.model.config.eos_token_id = self.cfg.eos_token_id
70             self.model.config.pad_token_id = self.cfg.pad_token_id
71             self.model.config.vocab_size = len(self.cfg.tgt_tokenizer)
72             self.model.config.max_length = self.cfg.max_length_decoder
73             self.model.config.min_length = self.cfg.min_length_decoder
74             self.model.config.no_repeat_ngram_size = 3
75             self.model.config.early_stopping = True
76             self.model.config.length_penalty = 2.0
77             self.model.config.num_beams = self.cfg.beam_size
78
79     def train(self):
80         print("Training...")
81         if self.cfg.use_eval_steps:
82             training_args = Seq2SeqTrainingArguments(
83                 predict_with_generate=True,
84                 evaluation_strategy="steps",
85                 save_strategy='steps',
86                 save_steps=self.cfg.eval_steps,
87                 eval_steps=self.cfg.eval_steps,
88                 output_dir=self.cfg.ckpt_dir,
89                 per_device_train_batch_size=self.cfg.train_batch_size,
90                 per_device_eval_batch_size=self.cfg.eval_batch_size,
91                 learning_rate=self.cfg.learning_rate,
92                 weight_decay=0.005,
93                 num_train_epochs=self.cfg.num_train_epochs
94             )
95         else:
96             training_args = Seq2SeqTrainingArguments(
97                 predict_with_generate=True,
98                 evaluation_strategy="epoch",

```

```

99         save_strategy='epoch',
100         output_dir=self.cfg.ckpt_dir,
101         per_device_train_batch_size=self.cfg.train_batch_size,
102         per_device_eval_batch_size=self.cfg.eval_batch_size,
103         learning_rate=self.cfg.learning_rate,
104         weight_decay=0.005,
105         num_train_epochs=self.cfg.num_train_epochs
106     )
107
108     data_collator = DataCollatorForSeq2Seq(
109         self.cfg.tgt_tokenizer,
110         model=self.model
111     )
112
113     trainer = Seq2SeqTrainer(
114         self.model,
115         training_args,
116         train_dataset=self.train_dataset,
117         eval_dataset=self.valid_dataset,
118         data_collator=data_collator,
119         tokenizer=self.cfg.tgt_tokenizer,
120         compute_metrics=self.compute_metrics
121     )
122
123     trainer.train()
124
125     def compute_metrics(self, eval_preds):
126         preds, labels = eval_preds
127         if isinstance(preds, tuple):
128             preds = preds[0]
129         decoded_preds = self.cfg.tgt_tokenizer.batch_decode(preds, skip_special_tokens=True)
130
131         labels = np.where(labels != -100, labels, self.cfg.tgt_tokenizer.pad_token_id)
132         decoded_labels = self.cfg.tgt_tokenizer.batch_decode(labels,
133             skip_special_tokens=True)
134
135         decoded_preds, decoded_labels = postprocess_text(decoded_preds, decoded_labels)
136
137         result = self.bleu_metric.compute(
138             predictions=decoded_preds,
139             references=decoded_labels
140         )
141
142         result = {"bleu_score": result["score"]}
143
144         prediction_lens = [np.count_nonzero(pred != self.cfg.tgt_tokenizer.
145             pad_token_id) for pred in preds]
146         result["gen_len"] = np.mean(prediction_lens)
147         result = {k: round(v, 4) for k, v in result.items()}
148
149         return result

```

4. Training

```

1 class BaseConfig:
2     """ base Encoder Decoder config """
3
4     def __init__(self, **kwargs):

```

```

5         for k, v in kwargs.items():
6             setattr(self, k, v)
7
8     class NMTConfig(BaseConfig):
9         # Data
10        src_lang = 'en'
11        tgt_lang = 'vi'
12        src_max_len = 75
13        tgt_max_len = 75
14
15        # Model
16        src_model_name = "bert-base-multilingual-cased"
17        tgt_model_name = "bert-base-multilingual-cased"
18
19        # Training
20        load_model_from_path = False
21        device = torch.device('cuda') if torch.cuda.is_available() else torch.device('cpu')
22        learning_rate = 3e-5
23        train_batch_size = 16
24        eval_batch_size = 8
25        num_train_epochs = 5
26        ckpt_dir = src_model_name + '_to_' + tgt_model_name
27        use_eval_steps = False
28        eval_steps = 2000
29
30        # Inference
31        max_length_decoder = 75
32        min_length_decoder = 25
33        beam_size = 1
34
35    cfg = NMTConfig()
36
37    manager = Manager(cfg, is_train=True)
38
39    manager.train()

```

5. Prediction

```

1 # load model
2 def load_model(cfg, checkpoint_name):
3     # Load Tokenizer
4     src_tokenizer_save_path = f"{cfg.ckpt_dir}/{cfg.src_lang}_tokenizer_{cfg.src_model_name}"
5     src_tokenizer = BertTokenizerFast.from_pretrained(src_tokenizer_save_path)
6
7     tgt_tokenizer_save_path = f"{cfg.ckpt_dir}/{cfg.tgt_lang}_tokenizer_{cfg.tgt_model_name}"
8     tgt_tokenizer = GPT2TokenizerFast.from_pretrained(tgt_tokenizer_save_path)
9
10    # Load Model
11    model_save_path = f"{cfg.ckpt_dir}/{checkpoint_name}"
12    model = EncoderDecoderModel.from_pretrained(model_save_path)
13
14    # Inference Param
15    device = torch.device('cuda') if torch.cuda.is_available() else torch.device('cpu')
16
17    return src_tokenizer, tgt_tokenizer, model, device
18

```



```
19 from tqdm import tqdm
20 def inference(
21     text,
22     src_tokenizer,
23     tgt_tokenizer,
24     model,
25     device="cpu",
26     max_length=75,
27     beam_size=5
28 ):
29     inputs = src_tokenizer(
30         text,
31         padding="max_length",
32         truncation=True,
33         max_length=max_length,
34         return_tensors="pt"
35     )
36     input_ids = inputs.input_ids.to(device)
37     attention_mask = inputs.attention_mask.to(device)
38     model.to(device)
39
40     outputs = model.generate(
41         input_ids,
42         attention_mask=attention_mask,
43         max_length=max_length,
44         early_stopping=True,
45         num_beams=beam_size,
46         length_penalty=2.0
47     )
48
49     output_str = tgt_tokenizer.batch_decode(outputs, skip_special_tokens=True)
50
51     return output_str
52
53 def inference_batch(
54     texts,
55     src_tokenizer,
56     tgt_tokenizer,
57     model,
58     device="cpu",
59     max_length=75,
60     beam_size=5,
61     batch_size=32
62 ):
63
64     pred_texts = []
65
66     if len(texts) < batch_size:
67         batch_size = len(texts)
68
69     for x in tqdm(range(0, len(texts), batch_size)):
70         text = texts[x:x+batch_size]
71
72         inputs = src_tokenizer(
73             text,
74             padding="max_length",
75             truncation=True,
76             max_length=max_length,
77             return_tensors="pt"
78         )
```

```

79         input_ids = inputs.input_ids.to(device)
80         attention_mask = inputs.attention_mask.to(device)
81         model.to(device)
82
83         outputs = model.generate(
84             input_ids,
85             attention_mask=attention_mask,
86             max_length=max_length,
87             early_stopping=True,
88             num_beams=beam_size,
89             length_penalty=2.0
90         )
91
92         output_str = tgt_tokenizer.batch_decode(outputs, skip_special_tokens=True)
93         pred_texts.extend(output_str)
94         torch.cuda.empty_cache()
95
96     return pred_texts
97
98 class BaseConfig:
99     """ base Encoder Decoder config """
100
101     def __init__(self, **kwargs):
102         for k, v in kwargs.items():
103             setattr(self, k, v)
104
105 class NMTConfig(BaseConfig):
106     # Data
107     src_lang = 'en'
108     tgt_lang = 'vi'
109     src_max_len = 75
110     tgt_max_len = 75
111
112     # Model
113     src_model_name = "bert-base-multilingual-cased"
114     tgt_model_name = "bert-base-multilingual-cased"
115
116     # Training
117     load_model_from_path = False
118     device = torch.device('cuda') if torch.cuda.is_available() else torch.device('cpu')
119
120     learning_rate = 3e-5
121     train_batch_size = 16
122     eval_batch_size = 8
123     num_train_epochs = 15
124     ckpt_dir = src_model_name + '_to_' + tgt_model_name
125     use_eval_steps = False
126     eval_steps = 2000
127
128     # Inference
129     max_length_decoder = 75
130     min_length_decoder = 25
131     beam_size = 5
132
133     cfg = NMTConfig()
134
135     # load data
136     data = load_dataset("mt_eng_vietnamese", "iwslt2015-en-vi", split="test")
137     src_texts = [sample["translation"]["en"] for sample in data]

```

```
138 tgt_texts = [sample["translation"]["vi"] for sample in data]
139
140 src_tokenizer, tgt_tokenizer, model, device = load_model(cfg, checkpoint_name="
    checkpoint-41665")
141
142 pred_texts = inference_batch(src_texts, src_tokenizer, tgt_tokenizer, model, device,
    beam_size=1)
143
144 sacrebleu.corpus_bleu(pred_texts, [tgt_texts]) # => 25.41 53.8/31.8/19.8/12.3
```

Phần 4. Câu hỏi trắc nghiệm

Câu hỏi 1 Mục tiêu của bài toán dịch máy là gì?

- a) Xây dựng mô hình dịch từ ngôn ngữ nguồn sang ngôn ngữ đích
- b) Xây dựng mô hình phân loại văn bản
- c) Xây dựng mô hình tóm tắt văn bản
- d) Xây dựng mô hình phân tích cảm xúc

Câu hỏi 2 Số lượng sample tập test của bộ dữ liệu sử dụng trong thực nghiệm là?

- a) 1267
- b) 1268
- c) 1269
- d) 1270

Câu hỏi 3 Hệ thống nào sau đây không phải là hệ thống dịch?

- a) Bing Translator
- b) Google Translate
- c) ChatGPT
- d) Dall-E

Câu hỏi 4 Thư viện sentencepiece dùng để làm gì?

- a) Xây dựng Transformer-Encoder
- b) Xây dựng Transformer-Decoder
- c) Xây dựng Subword-Based Tokenization
- d) Huấn luyện mô hình Transformer

Câu hỏi 5 BERT sử dụng kiến trúc nào sau đây?

- a) Bi-LSTM
- b) RNN
- c) Transformer-Encoder
- d) Transformer-Decoder

Câu hỏi 6 BERT-Large sử dụng bao nhiêu khối encoder?

- a) 6
- b) 12
- c) 18
- d) 24

Câu hỏi 7 Trong số các mask tokens trong quá trình huấn luyện BERT. Tỷ lệ mask nào là tối ưu nhất cho BERT?

- a) 80% thay bởi token [MASK] : 20% thay bởi token bất kỳ : 0% giữ nguyên không đổi

- b) 80% thay bởi token [MASK] : 0% thay bởi token bất kỳ : 20% giữ nguyên không đổi
- c) 80% thay bởi token [MASK] : 10% thay bởi token bất kỳ : 10% giữ nguyên không đổi
- d) 80% thay bởi token [MASK] : 15% thay bởi token bất kỳ : 5% giữ nguyên không đổi

Câu hỏi 8 GPT sử dụng kiến trúc nào sau đây?

- a) Bi-LSTM
- b) RNN
- c) Transformer-Encoder
- d) Transformer-Decoder

Câu hỏi 9 GPT-Small sử dụng bao nhiêu khối decoder?

- a) 6
- b) 12
- c) 18
- d) 24

Câu hỏi 10 Độ đo thường sử dụng cho mô hình dịch máy là?

- a) F1
- b) BLEU
- c) Accuracy
- d) ROUGE

- Hết -