

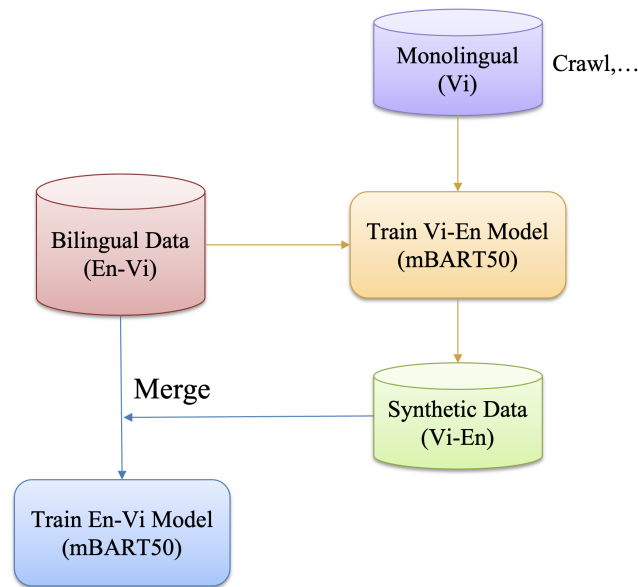
# Project: Low-Resource Machine Translation

Quoc-Thai Nguyen và Quang-Vinh Dinh

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

Ngày 25 tháng 2 năm 2024

## Phần I. Giới thiệu



Hình 1: Kỹ thuật dịch ngược áp dụng cho mô hình dịch Anh - Việt.

**Dịch Máy (Machine Translation)** với mục đích tự động dịch văn bản từ ngôn ngữ tự nhiên này sang ngôn ngữ tự nhiên khác. Một trong những thách thức lớn hiện nay của các hệ thống dịch là có ít tài nguyên để huấn luyện mô hình. Vì vậy, trong phần này, chúng ta sẽ tập trung vào cải thiện các mô hình dịch với ít tài nguyên (Low-resource machine translation), số lượng các cặp dữ liệu song ngữ hạn chế, ví dụ mô hình dịch huấn luyện với chỉ 20000 cặp câu tiếng anh và tiếng việt. Vì vậy, để cải tiến mô hình dịch trong trường hợp có ít tài nguyên, chúng ta tập trung vào 2 hướng tiếp cận như sau:

1. Hướng 1: Sử dụng mô hình pre-trained mBART50
2. Hướng 2: Sử dụng kỹ thuật dịch ngược (Back-translation) để tăng cường dữ liệu

## Phần II. Pre-trained mBART50

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)

Trong phần này, chúng ta sẽ sử dụng kỹ thuật fine-tuning để huấn luyện mô hình dịch theo chiều EN-VI.

### 1. Build Dataset

```
1 # install libs
2 !pip install -q transformers sentencepiece datasets accelerate evaluate sacrebleu
3
4 # import libs
5 import os
6 import numpy as np
7 import torch
8 from torch.utils.data import Dataset
9 from datasets import load_dataset
10 import evaluate
11 from transformers import (
12     MBart50TokenizerFast,
13     AutoModelForSeq2SeqLM,
14     DataCollatorForSeq2Seq,
15     Seq2SeqTrainingArguments,
16     Seq2SeqTrainer
17 )
18
19 # build dataset
20 class NMTDataset(Dataset):
21     def __init__(self, cfg, data_type="train"):
22         super().__init__()
23         self.cfg = cfg
24
25         self.src_texts, self.tgt_texts = self.read_data(data_type)
26
27         self.src_input_ids = self.texts_to_sequences(self.src_texts)
28         self.labels = self.texts_to_sequences(self.tgt_texts)
29
30     def read_data(self, data_type):
31         data = load_dataset(
32             "mt_eng_vietnamese",
33             "iwslt2015-en-vi",
34             split=data_type
35         )
36         src_texts = [sample["translation"][self.cfg.src_lang] for sample in data]
37         tgt_texts = [sample["translation"][self.cfg.tgt_lang] for sample in data]
38         return src_texts, tgt_texts
39
40     def texts_to_sequences(self, texts):
41         data_inputs = self.cfg.tokenizer(
42             texts,
43             padding='max_length',
44             truncation=True,
45             max_length=self.cfg.max_len,
46             return_tensors='pt'
47         )
```

```

48         return data_inputs.input_ids
49
50     def __getitem__(self, idx):
51         return {
52             "input_ids": self.src_input_ids[idx],
53             "labels": self.labels[idx]
54         }
55
56     def __len__(self):
57         return np.shape(self.src_input_ids)[0]

```

## 2. Tokenizer, Model, Metric

```

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    max_len = 75
13    add_special_tokens = True
14
15    # Model
16    model_name = "facebook/mbart-large-50-many-to-many-mmt"
17
18    # Training
19    device = torch.device('cuda') if torch.cuda.is_available() else torch.device('cpu')
20    learning_rate = 5e-5
21    train_batch_size = 32
22    eval_batch_size = 32
23    num_train_epochs = 2
24    save_total_limit = 1
25    ckpt_dir = f'./mbart50-{src_lang}-{tgt_lang}'
26    eval_steps = 1000
27
28    # Inference
29    beam_size = 5
30
31    cfg = NMTConfig()
32    cfg.tokenizer = MBart50TokenizerFast.from_pretrained(cfg.model_name)
33    model = AutoModelForSeq2SeqLM.from_pretrained(cfg.model_name)
34
35    # metric
36    metric = evaluate.load("sacrebleu")
37
38    def postprocess_text(preds, labels):
39        preds = [pred.strip() for pred in preds]
40        labels = [[label.strip()] for label in labels]
41
42        return preds, labels
43
44    def compute_metrics(eval_preds):
45        preds, labels = eval_preds
46        if isinstance(preds, tuple):

```

```

47     preds = preds[0]
48
49     preds= np.where(preds != -100, preds, cfg.tokenizer.pad_token_id)
50     decoded_preds = cfg.tokenizer.batch_decode(preds, skip_special_tokens=True,
51         clean_up_tokenization_spaces=True)
52
53     labels= np.where(labels != -100, labels, cfg.tokenizer.pad_token_id)
54     decoded_labels = cfg.tokenizer.batch_decode(labels, skip_special_tokens=True,
55         clean_up_tokenization_spaces=True)
56
57     decoded_preds, decoded_labels = postprocess_text(decoded_preds, decoded_labels)
58
59     result = metric.compute(predictions=decoded_preds, references=decoded_labels)
60     result = {"bleu": result["score"]}
61
62     prediction_lens = [np.count_nonzero(pred != cfg.tokenizer.pad_token_id) for pred
63         in preds]
64     result["gen_len"] = np.mean(prediction_lens)
65     result = {k: round(v, 4) for k, v in result.items()}
66
67     return result

```

### 3. Trainer

```

1 train_dataset = NMTDataset(cfg, data_type="train")
2 valid_dataset = NMTDataset(cfg, data_type="validation")
3 test_dataset = NMTDataset(cfg, data_type="test")
4
5 training_args = Seq2SeqTrainingArguments(
6     predict_with_generate=True,
7     evaluation_strategy="steps",
8     save_strategy='steps',
9     save_steps=cfg.eval_steps,
10    eval_steps=cfg.eval_steps,
11    output_dir=cfg.ckpt_dir,
12    per_device_train_batch_size=cfg.train_batch_size,
13    per_device_eval_batch_size=cfg.eval_batch_size,
14    learning_rate=cfg.learning_rate,
15    save_total_limit=cfg.save_total_limit,
16    num_train_epochs=cfg.num_train_epochs,
17    load_best_model_at_end=True,
18 )
19
20 data_collator = DataCollatorForSeq2Seq(
21     cfg.tokenizer,
22     model=model
23 )
24
25 trainer = Seq2SeqTrainer(
26     model,
27     training_args,
28     train_dataset=train_dataset,
29     eval_dataset=valid_dataset,
30     data_collator=data_collator,
31     tokenizer=cfg.tokenizer,
32     compute_metrics=compute_metrics
33 )
34
35 trainer.train()

```

#### 4. Training

Sau khi huấn luyện chúng ta thu được kết quả đánh giá như sau:

Step	Training Loss	Validation Loss	Bleu	Gen Len
1000	0.501900	0.599626	31.826500	32.365600
2000	0.493200	0.581667	32.680700	32.416900
3000	0.484500	0.572604	32.876400	32.588700
4000	0.484600	0.561723	33.089600	32.848700
5000	0.473300	0.555238	33.772600	32.662700
6000	0.471100	0.551521	33.693500	32.728900
7000	0.464800	0.545575	33.738100	32.953500
8000	0.453600	0.534770	34.213600	32.956700
9000	0.366800	0.546183	34.096400	32.688700
10000	0.366000	0.546823	33.752300	32.802200
11000	0.361400	0.544369	33.953200	32.760400
12000	0.359000	0.540036	34.301600	32.951100
13000	0.359900	0.535568	34.730100	32.772300
14000	0.355300	0.535669	34.714100	32.805400
15000	0.355800	0.532964	34.919300	32.687200
16000	0.352400	0.530680	34.871400	32.792000

Hình 2: Huấn luyện mô hình mBART50 trên bộ dữ liệu Anh - Việt.

#### 5. Evaluation

```

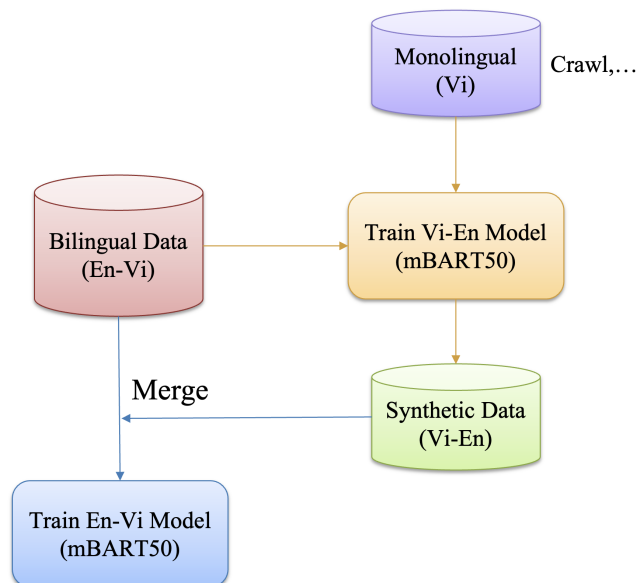
1 def inference(
2     text,
3     tokenizer,
4     model,
5     device="cpu",
6     max_length=75,
7     beam_size=5
8 ):
9     inputs = tokenizer(
10         text,
11         padding="max_length",
12         truncation=True,
13         max_length=max_length,
14         return_tensors="pt"
15     )
16     input_ids = inputs.input_ids.to(device)
17     attention_mask = inputs.attention_mask.to(device)

```

```
18     model.to(device)
19
20     outputs = model.generate(
21         input_ids,
22         attention_mask=attention_mask,
23         max_length=max_length,
24         early_stopping=True,
25         num_beams=beam_size,
26         length_penalty=2.0
27     )
28
29     output_str = tokenizer.batch_decode(outputs, skip_special_tokens=True)
30
31     return output_str
32
33 sentence = 'i go to school'
34 inference(sentence, cfg.tokenizer, model)
```

# Phần III. Back-Translation

Trong phần này chúng ta sẽ cải tiến mô hình dịch với kỹ thuật dịch ngược (Back-Translation)



Hình 3: Kỹ thuật dịch ngược áp dụng cho mô hình dịch.

Kỹ thuật dịch ngược bao gồm các bước sau:

- Bước 1: Xây dựng bộ dữ liệu đa ngữ (Monolingual data) dành cho tiếng việt.
- Bước 2: Huấn luyện mô hình mBART50 với chiều dịch là Việt - Anh.
- Bước 3: Tạo bộ dữ liệu tổng hợp sử dụng mô hình mBART50 Vi-En trên bộ dữ liệu đơn ngữ
- Bước 4: Kết hợp bộ dữ liệu tổng hợp với bộ dữ liệu gốc để huấn luyện mô hình En-Vi.

## 1. Vietnamese Monolingual Data Preparing

Bộ dữ liệu đơn ngữ dành cho tiếng việt có thể tạo ra bằng cách thu thập từ các trang tin tức,... Ở đây để đơn giản, chúng ta sử dụng tập test của bộ dữ liệu **PhoMT** lấy phần text của ngôn ngữ tiếng việt (Có thể tải về [tại đây](#)).

## 2. Train Vi-En Model

Chúng ta sử dụng bộ dữ liệu gốc huấn luyện mô hình mBART50 theo chiều Vi-En. Các bước huấn luyện tương tự với phần 2.

### 2.1. Build Dataset

```

1 # build dataset
2 class NMTDataset(Dataset):
3     def __init__(self, cfg, data_type="train"):
4         super().__init__()
5         self.cfg = cfg
6
7         self.src_texts, self.tgt_texts = self.read_data(data_type)
  
```

```

8
9     self.src_input_ids = self.texts_to_sequences(self.src_texts)
10    self.labels = self.texts_to_sequences(self.tgt_texts)
11
12    def read_data(self, data_type):
13        data = load_dataset(
14            "mt_eng_vietnamese",
15            "iwslt2015-en-vi",
16            split=data_type
17        )
18        src_texts = [sample["translation"][self.cfg.src_lang] for sample in data]
19        tgt_texts = [sample["translation"][self.cfg.tgt_lang] for sample in data]
20        return src_texts, tgt_texts
21
22    def texts_to_sequences(self, texts):
23        data_inputs = self.cfg.tokenizer(
24            texts,
25            padding='max_length',
26            truncation=True,
27            max_length=self.cfg.max_len,
28            return_tensors='pt'
29        )
30        return data_inputs.input_ids
31
32    def __getitem__(self, idx):
33        return {
34            "input_ids": self.src_input_ids[idx],
35            "labels": self.labels[idx]
36        }
37
38    def __len__(self):
39        return np.shape(self.src_input_ids)[0]

```

## 2.2. Tokenizer, Model, Metric

```

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 = 'vi'
11    tgt_lang = 'en'
12    max_len = 75
13    add_special_tokens = True
14
15    # Model
16    model_name = "facebook/mbart-large-50-many-to-many-mmt"
17
18    # Training
19    device = torch.device('cuda') if torch.cuda.is_available() else torch.device('cpu')
20    learning_rate = 5e-5
21    train_batch_size = 32
22    eval_batch_size = 32
23    num_train_epochs = 2
24    save_total_limit = 1

```



```

25     ckpt_dir = f'./mbart50-{src_lang}-{tgt_lang}'
26     eval_steps = 1000
27
28     # Inference
29     beam_size = 5
30
31     cfg = NMTConfig()
32     cfg.tokenizer = MBart50TokenizerFast.from_pretrained(cfg.model_name)
33     model = AutoModelForSeq2SeqLM.from_pretrained(cfg.model_name)
34
35     # metric
36     metric = evaluate.load("sacrebleu")
37
38     def postprocess_text(preds, labels):
39         preds = [pred.strip() for pred in preds]
40         labels = [[label.strip()] for label in labels]
41
42         return preds, labels
43
44     def compute_metrics(eval_preds):
45         preds, labels = eval_preds
46         if isinstance(preds, tuple):
47             preds = preds[0]
48
49         preds = np.where(preds != -100, preds, cfg.tokenizer.pad_token_id)
50         decoded_preds = cfg.tokenizer.batch_decode(preds, skip_special_tokens=True,
51             clean_up_tokenization_spaces=True)
52
53         labels = np.where(labels != -100, labels, cfg.tokenizer.pad_token_id)
54         decoded_labels = cfg.tokenizer.batch_decode(labels, skip_special_tokens=True,
55             clean_up_tokenization_spaces=True)
56
57         decoded_preds, decoded_labels = postprocess_text(decoded_preds, decoded_labels)
58
59         result = metric.compute(predictions=decoded_preds, references=decoded_labels)
60         result = {"bleu": result["score"]}
61
62         prediction_lens = [np.count_nonzero(pred != cfg.tokenizer.pad_token_id) for pred
63             in preds]
64         result["gen_len"] = np.mean(prediction_lens)
65         result = {k: round(v, 4) for k, v in result.items()}
66
67         return result

```

### 2.3. Trainer

```

1 #training
2 train_dataset = NMTDataset(cfg, data_type="train")
3 valid_dataset = NMTDataset(cfg, data_type="validation")
4 test_dataset = NMTDataset(cfg, data_type="test")
5
6 training_args = Seq2SeqTrainingArguments(
7     predict_with_generate=True,
8     evaluation_strategy="steps",
9     save_strategy='steps',
10    save_steps=cfg.eval_steps,
11    eval_steps=cfg.eval_steps,
12    output_dir=cfg.ckpt_dir,
13    per_device_train_batch_size=cfg.train_batch_size,
14    per_device_eval_batch_size=cfg.eval_batch_size,

```

```

15     learning_rate=cfg.learning_rate,
16     save_total_limit=cfg.save_total_limit,
17     num_train_epochs=cfg.num_train_epochs,
18     load_best_model_at_end=True,
19 )
20 data_collator = DataCollatorForSeq2Seq(
21     cfg.tokenizer,
22     model=model
23 )
24 trainer = Seq2SeqTrainer(
25     model,
26     training_args,
27     train_dataset=train_dataset,
28     eval_dataset=valid_dataset,
29     data_collator=data_collator,
30     tokenizer=cfg.tokenizer,
31     compute_metrics=compute_metrics
32 )
33 trainer.train()

```

### 3. Synthetic Data Generation

Chúng ta sử dụng mô hình đã huấn luyện ở mục 2 để dịch các bản dịch ở dữ liệu đơn ngữ.

```

1 # load vi data from phomt dataset
2 with open('./test.vi', 'r') as f:
3     vi_phomt = f.readlines()
4 # inference a sentence or batch
5 def inference(
6     text,
7     tokenizer,
8     model,
9     device="cpu",
10    max_length=75,
11    beam_size=5
12 ):
13     inputs = tokenizer(
14         text,
15         padding="max_length",
16         truncation=True,
17         max_length=max_length,
18         return_tensors="pt"
19     )
20     input_ids = inputs.input_ids.to(device)
21     attention_mask = inputs.attention_mask.to(device)
22     model.to(device)
23     outputs = model.generate(
24         input_ids,
25         attention_mask=attention_mask,
26         max_length=max_length,
27         early_stopping=True,
28         num_beams=beam_size,
29         length_penalty=2.0
30     )
31     output_str = tokenizer.batch_decode(outputs, skip_special_tokens=True)
32     return output_str[0]
33
34 en_phomt_preds = []
35 for sent in vi_phomt:
36     en_sent = inference(sent, cfg.tokenizer, model)
37     en_phomt_preds.append(en_sent)

```

## 4. Train En-VI model

Trong phần này chúng ta kết hợp bộ dữ liệu tổng hợp và bộ dữ liệu song ngữ ban đầu để huấn luyện mô hình.

### 4.1. Build Dataset

```

1 # build dataset
2 class NMTDataset(Dataset):
3     def __init__(self, cfg, src_augs=[], tgt_augs=[], data_type="train"):
4         super().__init__()
5         self.cfg = cfg
6
7         self.src_texts, self.tgt_texts = self.read_data(data_type)
8         if data_type == 'train':
9             self.src_texts.extend(src_augs)
10            self.tgt_texts.extend(tgt_augs)
11
12            self.src_input_ids = self.texts_to_sequences(self.src_texts)
13            self.labels = self.texts_to_sequences(self.tgt_texts)
14
15        def read_data(self, data_type):
16            data = load_dataset(
17                "mt_eng_vietnamese",
18                "iwslt2015-en-vi",
19                split=data_type
20            )
21            src_texts = [sample["translation"][self.cfg.src_lang] for sample in data]
22            tgt_texts = [sample["translation"][self.cfg.tgt_lang] for sample in data]
23            return src_texts, tgt_texts
24
25        def texts_to_sequences(self, texts):
26            data_inputs = self.cfg.tokenizer(
27                texts,
28                padding='max_length',
29                truncation=True,
30                max_length=self.cfg.max_len,
31                return_tensors='pt'
32            )
33            return data_inputs.input_ids
34
35        def __getitem__(self, idx):
36            return {
37                "input_ids": self.src_input_ids[idx],
38                "labels": self.labels[idx]
39            }
40
41        def __len__(self):
42            return np.shape(self.src_input_ids)[0]

```

### 4.2. Tokenizer, Model, Metric

```

1 # config
2 class BaseConfig:
3     """ base Encoder Decoder config """
4
5     def __init__(self, **kwargs):
6         for k, v in kwargs.items():
7             setattr(self, k, v)

```

```

8
9 class NMTConfig(BaseConfig):
10     # Data
11     src_lang = 'en'
12     tgt_lang = 'vi'
13     max_len = 75
14     add_special_tokens = True
15
16     # Model
17     model_name = "facebook/mbart-large-50-many-to-many-mmt"
18
19     # Training
20     device = torch.device('cuda') if torch.cuda.is_available() else torch.device('cpu')
21     learning_rate = 5e-5
22     train_batch_size = 32
23     eval_batch_size = 32
24     num_train_epochs = 2
25     save_total_limit = 1
26     ckpt_dir = f'./mbart50-{src_lang}-{tgt_lang}-backtranslation-2'
27     eval_steps = 1000
28
29     # Inference
30     beam_size = 5
31
32 cfg = NMTConfig()
33
34 cfg.tokenizer = MBart50TokenizerFast.from_pretrained(cfg.model_name)
35 model = AutoModelForSeq2SeqLM.from_pretrained(cfg.model_name)
36
37 train_dataset = NMTDataset(cfg, en_phomt_preds, vi_phomt, data_type="train")
38 valid_dataset = NMTDataset(cfg, data_type="validation")
39 test_dataset = NMTDataset(cfg, data_type="test")
40
41 # metric
42 metric = evaluate.load("sacrebleu")
43
44 def postprocess_text(preds, labels):
45     preds = [pred.strip() for pred in preds]
46     labels = [[label.strip()] for label in labels]
47     return preds, labels
48
49 def compute_metrics(eval_preds):
50     preds, labels = eval_preds
51     if isinstance(preds, tuple):
52         preds = preds[0]
53
54     preds = np.where(preds != -100, preds, cfg.tokenizer.pad_token_id)
55     decoded_preds = cfg.tokenizer.batch_decode(preds, skip_special_tokens=True,
56         clean_up_tokenization_spaces=True)
57
58     labels = np.where(labels != -100, labels, cfg.tokenizer.pad_token_id)
59     decoded_labels = cfg.tokenizer.batch_decode(labels, skip_special_tokens=True,
60         clean_up_tokenization_spaces=True)
61
62     decoded_preds, decoded_labels = postprocess_text(decoded_preds, decoded_labels)
63
64     result = metric.compute(predictions=decoded_preds, references=decoded_labels)
65     result = {"bleu": result["score"]}

```

```

65     prediction_lens = [np.count_nonzero(pred != cfg.tokenizer.pad_token_id) for pred
    in preds]
66     result["gen_len"] = np.mean(prediction_lens)
67     result = {k: round(v, 4) for k, v in result.items()}
68
69     return result

```

### 4.3. Trainer

```

1  #training
2  training_args = Seq2SeqTrainingArguments(
3      predict_with_generate=True,
4      evaluation_strategy="steps",
5      save_strategy='steps',
6      save_steps=cfg.eval_steps,
7      eval_steps=cfg.eval_steps,
8      output_dir=cfg.ckpt_dir,
9      per_device_train_batch_size=cfg.train_batch_size,
10     per_device_eval_batch_size=cfg.eval_batch_size,
11     learning_rate=cfg.learning_rate,
12     save_total_limit=cfg.save_total_limit,
13     num_train_epochs=cfg.num_train_epochs,
14     load_best_model_at_end=True,
15 )
16 data_collator = DataCollatorForSeq2Seq(
17     cfg.tokenizer,
18     model=model
19 )
20 trainer = Seq2SeqTrainer(
21     model,
22     training_args,
23     train_dataset=train_dataset,
24     eval_dataset=valid_dataset,
25     data_collator=data_collator,
26     tokenizer=cfg.tokenizer,
27     compute_metrics=compute_metrics
28 )
29 trainer.train()

```

## 5. Experiment

Kết quả đánh giá trên tập test của các thực nghiệm được mô tả như bảng sau:

Experiment	Model	ScoreBLEU
#1	Standard Transformer (Greedy Search)	24.66
#2	BERT-to-BERT (Greedy Search)	25.41
#3	BERT-to-GPT2 (Greedy Search)	23.56
#4	mBART50	34.87
#5	Back-Translation (Monolingual)	35.22

Hình 4: So sánh kết quả thực nghiệm.

## Phần 4. Câu hỏi trắc nghiệm

**Câu hỏi 1** Mô hình BART có kiến trúc là gì?

- a) Transformer-Encoder
- b) Transformer-Decoder
- c) Transformer-Encoder-Decoder
- d) BiLSTM

**Câu hỏi 2** Hàm mục tiêu nào sau đây không phải của BART?

- a) Token Masking
- b) Token Deletion
- c) Sentence Permutation
- d) Next Sentence Prediction

**Câu hỏi 3** Mô hình BART-Base có bao nhiêu lớp encoder?

- a) 6
- b) 12
- c) 24
- d) 48

**Câu hỏi 4** Mô hình BART-Large có bao nhiêu lớp encoder?

- a) 6
- b) 12
- c) 24
- d) 48

**Câu hỏi 5** mBART50 được huấn luyện trên bao nhiêu ngôn ngữ?

- a) 5
- b) 25
- c) 50
- d) 100

**Câu hỏi 6** mBART50 là mô hình tiền huấn luyện đa ngữ cho bài toán nào?

- a) Text Summarization
- b) Text Classification
- c) Question-Answering
- d) Machine Translation

**Câu hỏi 7** Tập dữ liệu đơn ngữ tiếng việt được lấy từ bộ dữ liệu nào sau đây?

- a) PhoMT
- b) IMDB-Review

- c) C4
- d) ROOT

**Câu hỏi 8** Số lượng sample trong bộ dữ liệu đơn ngữ tiếp việt là?

- a) 19150
- b) 19151
- c) 19152
- d) 19153

**Câu hỏi 9** Mô hình nào sau đây có kiến trúc tương tự mô hình BART?

- a) BERT
- b) GPT
- c) RoBERTa
- d) T5

**Câu hỏi 10** Mô hình T5 sử dụng tiền tố nào cho bài toán dịch?

- a) translate English to Vietnamese
- b) translate to Vietnamese from English
- c) Cả 2 đáp án đều đúng
- d) Cả 2 đáp án đều sai

**Câu hỏi 11** Độ đo đánh giá được sử dụng trong phần thực nghiệm là?

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

- Hết -