



Natural Language Processing CA#5

Student Name:
Pouya Haji Mohammadi Gohari

SID:810102113

Date of deadline
Monday 10th June, 2024

Dept. of Computer Engineering

University of Tehran

Contents

1	Dataset	3
2	Part 1: Training BPE tokenizer and Pre-processing dataset	8
3	Part 2: Training LSTM ENCODER-DECODER	11
4	Part 3: Training Transformer ENCODER-DECODER	14
5	Evaluation Metric and Examine Test dataset	15
	References	20

1 Dataset

In this section, we are going to investigate through our translation dataset which can achieve from this link. To investigate data, we should first load it with utilizing following code:

```
1 drive.mount('/content/drive')
2 !wget -P /content/drive/My\ Drive/ https://object.pouta.csc.fi/OPUS-MIZAN/v1/moses
   ↪ /en-fa.txt.zip
3 !unzip "/content/drive/My Drive/en-fa.txt.zip"
```

First of all, we will allow 'colab' to access our google drive. Afterward, the dataset will be downloaded by wget command. With unzipping it, the data will be prepared to use.

Based on the two file created after unzipping (Persian and English text file) we will answer some questions as follow:

Q1: The number of lines for each file can be obtained with following command:

```
1 !wc -l MIZAN.en-fa.en
2 !wc -l MIZAN.en-fa.fa
```

These lines are shown in table 1:

File Name	Number of Lines
'MIZAN.en-fa.fa' file	1,021,597
'MIZNA.en-fa.en' file	1,021,597

Table 1: Number of Lines for Each File

In order to see what is three first lines of each, we utilize following command:

```
1 !head MIZAN.en-fa.en -n 3
2 !head MIZAN.en-fa.fa -n 3
```

These three lines have been illustrated in figure 1:

The story which follows was first written out in Paris during the Peace Conference from notes jotted daily on the march, strengthened by some reports sent to my chiefs in Cairo. Afterwards, in the autumn of 1919, this first draft and some of the notes were lost.

داستانی که از نظر شما می‌گذرد، ابتدا ضمن کنفرانس صلح پاریس از روی یادداشت‌هایی که به طور روزانه در حال خدمت در صف برداشته شده بودند و از روی گزارشاتی که برای رؤسای من در قاهره ارسال گردیده بودند نوشته شد. بعداً در پائیز سال ۱۹۱۹، این نوشته اولیه و بعضی از یادداشت‌ها، مفقود شدند.

Figure 1: First Three Lines from Each File

Q2: In order to tokenize with white spaces, we have implemented code below:

```
1 def number_of_tokens_whitespaces(example:list[str]) -> int:
2     """
3     This function will simply tokenize each sentence with respect to white spaces
4     Args:
5         example: It should be list of string as representation of a sentence
6     Output:
7         This function will return a integer since we want to
8     """
```

```

9     example = example.replace('\u200c', '')
10    example = example.replace('\u202C', '')
11    token_counts = len(example.split())
12    return token_counts

```

Since we are dealing with Persian texts in our dataset, we consider to delete some characters like '۰' and '۰' which stands for half spaces and right to left respectively. We will count the tokens in each line using length of the line where splitted function.

Finally with utilizing 'matplotlib' we can plot histogram for number of tokens for each line. The code below will help us to do so:

```

1 def plot_histogram(tokens_list:list[int], language:str) -> None:
2     """
3     Creating a plot for each language.
4     Args:
5         tokens_list: list of token's number for each line
6         language: Specify the language string
7     Output:
8         Creating a plot
9     """
10    plt.figure(figsize=(10,8))
11    plt.hist(tokens_list, bins=30, alpha=0.75)
12    plt.xlabel('Each line')
13    plt.ylabel('Frequency of tokens')
14    plt.title(f'Token Count Histogram for {language} Lines')
15    plt.grid(True)
16    plt.show()

```

From opening a file with specify path, we can calculate number of tokens of each line and finally plot histogram for each. The figure 2 and figure 3 are histogram for each file.

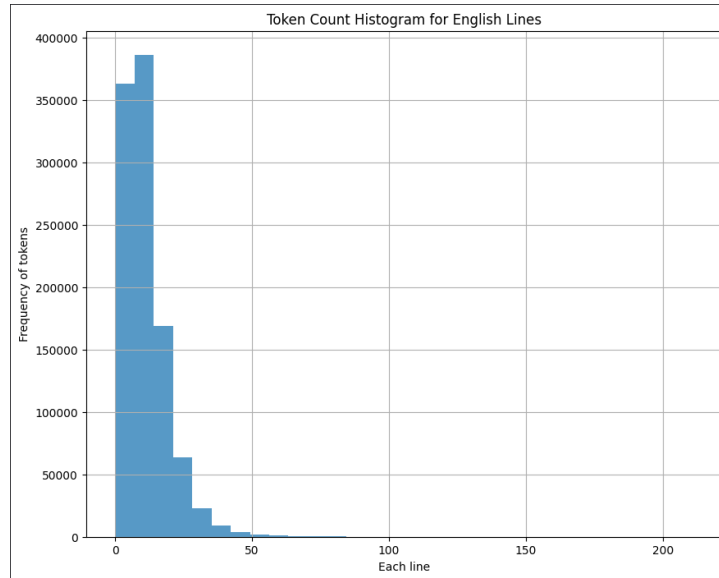


Figure 2: Token Count for English Lines

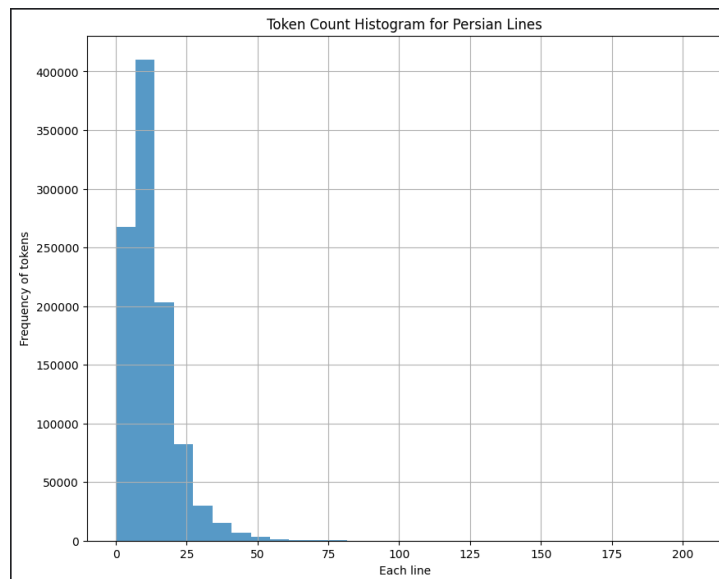


Figure 3: Token Count for Persian Lines

Q3:It would be very good idea to prune the outliers based on the target language and have less noisy number of tokens per line. So based on target language we can prune the lines having more than 50 tokens or less than 10. Based on the 'ignoring_outlier' function, we can create a dictionary for both target and source languages to save index alongside their number of tokens for each line.

```
1 def ignoring_outliers(src_dict:dict, trg_dict:dict, low_threshold:int,
2   ↪ high_threshold:int) -> tuple:
3     """
```

```

3 This function will remove outliers with respect to the low and high threshold
  ↳ based on the trg_dict
4 Args:
5     src_dict: a dictionary of source contains {index:number_of_tokens}
6     trg_dict: a dictionary of target just like src_dict
7     low_threshold: if number of tokens are less than this threshold it would
  ↳ be pruned.
8     high_threshold: if number of token are higher than this threshold it
  ↳ would be also condiser to be pruned.
9 Output:
0     This fucntion will return two dictionary in tuple data type.
1 """
2 revised_src_dict, revised_trg_dict = {}, {}
3 indices = []
4 for index, n_tokens in trg_dict.items():
5     if n_tokens > high_threshold or n_tokens < low_threshold:
6         continue
7     else:
8         revised_trg_dict[index] = n_tokens
9         indices.append(index)
10 for index in indices:
11     revised_src_dict[index] = src_dict[index]
12
13 return (revised_src_dict, revised_trg_dict)

```

Now from length of each dictionary, we can obtain filtered lines. Filtered lines can be chosen with respect to the corresponding index as follow:

```

1 with open('/content/MIZAN.en-fa.en', 'r', encoding='utf-8') as eng_file:
2     eng_lines = [line for line in eng_file.readlines()]
3 with open('/content/MIZAN.en-fa.fa', 'r', encoding='utf-8') as per_file:
4     per_lines = [line for line in per_file.readlines()]
5 filtered_eng_lines = [eng_lines[index] for index in revised_src_dict.keys()]
6 filtered_per_lines = [per_lines[index] for index in revised_trg_dict.keys()]

```

The filtered lines has been reported in table 2:

Language	Filtered Lines
Persian file	1,021,597
English file	1,021,597

Table 2: Number of Filtered Lines for Each File

Q4:In order to shuffle filtered dataset, we will zip filtered lines in order to pair them(since each English sentence translated into one Persian sentence where they are indeed parallel).Moreover, we will set random seed to 42 since for each run we obtain same shuffled data.The code below will clarify what we did:

```

1 import random
2 random.seed(42)

```

```

3 random.shuffle(paired_lines)
4 n_train_samples = 500000
5 n_validation_samples= 5000
6 n_test_samples = 10000
7 train_set = paired_lines[:n_train_samples]
8 val_set = paired_lines[n_train_samples :n_validation_samples+n_train_samples]
9 test_set = paired_lines[n_train_samples + n_validation_samples : n_train_samples +
    ↪ n_validation_samples + n_test_samples]

```

From implemented code above, we will assign 500,000 lines of first paired lines to the train set and next 5,000 sentences would be consider as validation set and finally the next 10,000 would be test set.

Q5: We will save these into 6 different file as follow:

- train.fa: Including Persian lines of train set.
- train.en: Including English lines of train set.
- valid.fa: Including Persian lines of validation set.
- valid.en: Including English lines of validation set.
- test.fa: Including Persian lines of test set.
- test.en: Including English lines of test set.

Following code will satisfy our purpose(We will save them in raw_data folder):

```

1 train_path_en = '/content/raw_data/train.en'
2 train_path_fa = '/content/raw_data/train.fa'
3 val_path_en = '/content/raw_data/valid.en'
4 val_path_fa = '/content/raw_data/valid.fa'
5 test_path_en = '/content/raw_data/test.en'
6 test_path_fa = '/content/raw_data/test.fa'
7 paths = [train_path_en, train_path_fa, val_path_en, val_path_fa, test_path_en,
    ↪ test_path_fa]
8 train_en, train_fa = zip(*train_set)
9 val_en, val_pa = zip(*val_set)
10 test_en, test_pa = zip(*test_set)
11 data = [train_en, train_fa, val_en, val_pa, test_en, test_pa]
12 for index, path in enumerate(paths):
13     with open(path, 'w', encoding='utf-8') as myfile:
14         for line in data[index]:
15             myfile.write(line)

```

2 Part 1: Training BPE tokenizer and Pre-processing dataset

In this section we will train a BPE¹ tokenizer by utilizing the sentence piece library. One will be trained on Persian train set while another will be trained on English train set. We will set the vocab size to 10K as description said. The following code will show you the way:

```
1 train_en_file = '/content/raw_data/train.en'
2 train_fa_file = '/content/raw_data/train.fa'
3 model_prefix_en = '/content/tokenizer/english_tokenizer'
4 model_prefix_fa = '/content/tokenizer/persian_tokenizer'
5 vocab_size = 10000
6 model_type = 'bpe'
7 character_coverage = 1.0
8 spm.SentencePieceTrainer.train(
9     input=train_en_file,
10    model_prefix=model_prefix_en,
11    vocab_size=vocab_size,
12    model_type=model_type,
13    character_coverage=character_coverage
14 )
15 spm.SentencePieceTrainer.train(
16    input=train_fa_file,
17    model_prefix=model_prefix_fa,
18    vocab_size=vocab_size,
19    model_type=model_type,
20    character_coverage=character_coverage
21 )
```

First we will define the input file path for trainer of sentence piece and output file path as input and model prefix arguments. Then vocab size will be set where in our case is 10k. We want to train BPE tokenizer therefore we pass 'bpe' as model type. Character coverage is 1, passing higher value for this parameter (ranging 0 to 1) means that more unique characters from the training data will be included in the vocabulary which can be important for languages with many unique characters which we consider to set it to 1.

After training our BPE tokenizer for each train set (English and Persian), we consider to make a directory called 'tokenized' and encode each file for raw_data directory. Firstly, we will load two trained models of BPE tokenizer with following code:

```
1 eng_tokenizer = spm.SentencePieceProcessor(model_file='/content/tokenizer/
    ↳ english_tokenizer.model')
2 per_tokenizer = spm.SentencePieceProcessor(model_file='/content/tokenizer/
    ↳ persian_tokenizer.model')
```

Then utilizing the following code we can tokenize each file and store it in specified path.

```
1 def tokenize_file(input_file:str, output_file:str, tokenizer:spm.
    ↳ SentencePieceProcessor):
2     """
```

¹Byte-Pair Encoding


```

3 This function will tokenize each file based on your tokenizer.
4 Args:
5     input_file: Specify the path of your raw data
6     output_file: Specify where you want to write the tokenized data
7     tokenizer: pass your trained tokenizer
8 Output:
9     A file will be created in specify location where each line is
10    ↪ tokenized format of raw data
11 """
12 with open(input_file, 'r', encoding='utf-8') as raw_file:
13     lines = raw_file.readlines()
14     cleaned_lines = [line.replace('\u200c', ' ') for line in lines]
15     super_cleaned_lines = [line.replace('\u202C', ' ') for line in cleaned_lines]
16 with open(output_file, 'w', encoding='utf-8') as tokenized_file:
17     for line in tqdm(super_cleaned_lines):
18         tokenized_line = tokenizer.encode(line.strip(), out_type=str)
19         tokenized_file.write(' '.join(tokenized_line) + '\n')

```

As we said in previous section we will replace two characters representing half spaces and rtl with space. Finally utilizing code above, we can tokenize each raw file into tokenized directory.

```

1 tokenize_file(train_path_en, train_path_en_tokenized, eng_tokenizer)
2 tokenize_file(val_path_en, val_path_en_tokenized, eng_tokenizer)
3 tokenize_file(test_path_en, test_path_en_tokenized, eng_tokenizer)
4 tokenize_file(train_path_fa, train_path_fa_tokenized, per_tokenizer)
5 tokenize_file(val_path_fa, val_path_fa_tokenized, per_tokenizer)
6 tokenize_file(test_path_fa, test_path_fa_tokenized, per_tokenizer)

```

Now we can proceed to use fair seq library in order to pre-process the tokenized data. Let's first use the following command and explain each arguments later.

```

1 !fairseq-preprocess --source-lang en --target-lang fa \
2                     --trainpref ./tokenized/train \
3                     --validpref ./tokenized/valid \
4                     --testpref  ./tokenized/test \
5                     --destdir   ./data_bin/ \
6                     --nwordssrc 10000\
7                     --nwordstgt 10000

```

First line will specify the source language that we are translating from with **--source-lang** where it is en(English) and the target language that we are translating to with **--target-lang** where it is fa(Persian). After that we will specify the train, validation and test path for **--trainpref**, **--validpref**, **--testpref** respectively. **--destdir** will get a path where this command will create some files for us which are useful for training phase. At the end **--nwordssrc**, **--nwordstgt** will be set to 10k as description said so. The usage of these two commands are as follow:

- **nwordssrc**: Setting this parameter to 10k, means the source language vocabulary will be limited to 10,000 most frequent tokens or words. Moreover, this parameter would mainly used to specify how many words will be retained for source language.
- **nwordtrg**: Like **nwordssrc**, this will be used to specify how many words will be retained for target language.

The usage and definition of these two parameters are according to documentation of fairseq. Now lets answer what files these preprocess command line will create for us. This function will create three main files as follow:

1. Dictionary Files: These files contain the vocabulary for source and target languages where each line of these files represents as token and their corresponding frequency.
2. Binary Files: For either source or target with respect to train, validation and test set will create multiple binary files. These files consist of pre-processed data in binary format. The binary format is efficient for input and output operation during training.
3. Index Files: These index files provide a mapping that correspond to their binary files where helps to access the data in efficient way.

Moreover, binary files include tokenized sentences from their corresponding files (for example train.fa and train.en). These tokens will be represented as their tokens ids based on the vocabulary file. (dict.en.text and dict.fa.text) Index files includes offsets that points to start of the sentences corresponding to their binary files. The overall files created by this function has been illustrated in figure 4:

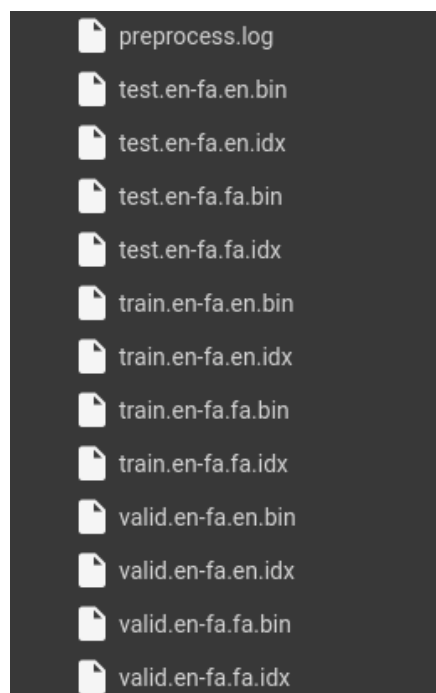


Figure 4: Total Files Created by fairseq-preprocess

3 Part 2: Training LSTM ENCODER-DECODER

After pre-process our dataset in previous section, it is time to proceed into training a encoder-decoder model based on LSTM. Let us introduce you to our implemented code then explain why we choose some parameters:

```
1 !fairseq-train \  
2   "./data_bin/" \  
3   --arch lstm --share-decoder-input-output-embed \  
4   --encoder-layers 6 --decoder-layers 6 \  
5   --optimizer adam --adam-betas '(0.9, 0.98)' --clip-norm 0.0 \  
6   --lr 2e-3 --lr-scheduler inverse_sqrt --warmup-updates 4000 --warmup-init-lr 5  
7   ↪ e-08 \  
8   --dropout 0.3 --weight-decay 0.0001 \  
9   --criterion label_smoothed_cross_entropy --label-smoothing 0.2 \  
10  --max-tokens 4096 \  
11  --max-sentences 128 \  
12  --update-freq 4 \  
13  --num-workers 2 \  
14  --eval-bleu \  
15  --eval-bleu-args '{"beam": 5, "max_len_a": 1.2, "max_len_b": 10}' \  
16  --eval-bleu-detok moses \  
17  --eval-bleu-print-samples \  
18  --best-checkpoint-metric bleu --maximize-best-checkpoint-metric \  
19  --fp16 --memory-efficient-fp16 \  
20  --max-epoch 5 \  
21  --save-dir ./data_bin/checkpoints/ \  
    --tensorboard-logdir ./data_bin/logs
```

As description said we have chose some parameters as follow:

1. Number of Epochs: The total epoch we consider model to be trained is 5 as description said.
2. Criterion: label_smoothed_cross_entropy with label_smoothing of 0.2.
3. Optimizer: The adam was set for optimizer with parameter beta set to 0.9 and 0.98.
4. Layers: Decoder and Encoder layers must be set to 6 as said.

As said we have chosen same hyper parameters for each model we are willing to train. Thus after training each model with same parameters, we decided to choose learning rate of $2e - 3$ with initial learning rate of $5e - 8$ and 4000 for warm up steps. (This parameters was set in way that both model will be effective in their own way. Moreover, we have tried so many parameters for each model and obtain the best parameters not just for model based on LSTM but also best for model that is based on Transformers)

We will use blue score for validation set to see the performance with beam width of 5. Each epoch will saved under the 'data_bin/checkpoints/' and logs would be saved under the 'data_bin/checkpoints/' in order to use to future purposes with tensor board. The best checkpoint will be saved under the same directory of checkpoints.

Some other hyper paramterers and their explanation:

1. max tokens: This parameter sets a limit for maximum number of tokens in each batch. This parameter will help the GPU memory usage.

2. max sentences: According to the fairseq documentation the batch size and max sentences are equal parameters since they represent how much example(sentences) must be in a batch. Furthermore, this parameter will set a maximum limit for how many sentences should be in a batch. This will help us when we dealing a dataset where the variation of sequence length is high.
3. update freq: This parameter is as same as gradient accumulation. For example in our scenario it has been set to 4 meaning that loss of 4 batches will be accumulated then weights will be updated.
4. number of workers: This parameter sets number of worker threads for data loading. More workers will speed up data loading by parallelizing the process. When dealing with large dataset this can help us to load data more efficient.

Moreover, max tokens and max sentences are two maximum limit of number of tokens and number of sentences can be seen with in a batch. If a batch reach max sentences first the batch will be processed even with lower max tokens and vice versa.

In order to answer the provided question, we have tried so many values for these two parameters and even with higher values for them, the loss' scale would be decreased so much leading to change the order for loss of training. Thus we should have a trade off between the speed of training and order of how loss decreasing. Even though the higher values for these parameters are possible and could be benefit us to dealing with large dataset but they must choose wisely where they to do not effect the training phase. In summary, in order to control the dataset's volume we can indeed use these two parameters where can help the GPU memory usage with respect to limiting the number of tokens and number of sentences within a batch.

Since we could not downmload the loss from log for training and validation set, we will use following function to do so: (These csv files will be uploaded alongside with code and report)

```
1 import os
2 import pandas as pd
3 from tensorboard.backend.event_processing import event_accumulator
4
5 log_dir = './data_bin/logs/train'
6
7 ea = event_accumulator.EventAccumulator(log_dir)
8 ea.Reload()
9 print("Available scalar keys:")
10 print(ea.Tags()['scalars'])
11
12
13 scalars = ea.Scalars('loss')
14 df = pd.DataFrame([(s.step, s.value) for s in scalars], columns=['step', 'loss'])
15 df.to_csv('lstm_training_loss.csv', index=False)
16
17 log_dir = './data_bin/logs/valid'
18
19 ea = event_accumulator.EventAccumulator(log_dir)
20 ea.Reload()
21 print("Available scalar keys:")
22 print(ea.Tags()['scalars'])
23
```

```

24
25 scalars = ea.Scalars('loss')
26 df = pd.DataFrame([(s.step, s.value) for s in scalars], columns=['step', 'loss'])
27 df.to_csv('lstm_validation_loss.csv', index=False)
28
29 %load_ext tensorboard
30 %tensorboard --logdir ./data_bin/logs --port 12345

```

Note that we have downloaded these files from 'colab'. The loss for training alongside with validation loss has been illustrated in figure 5:

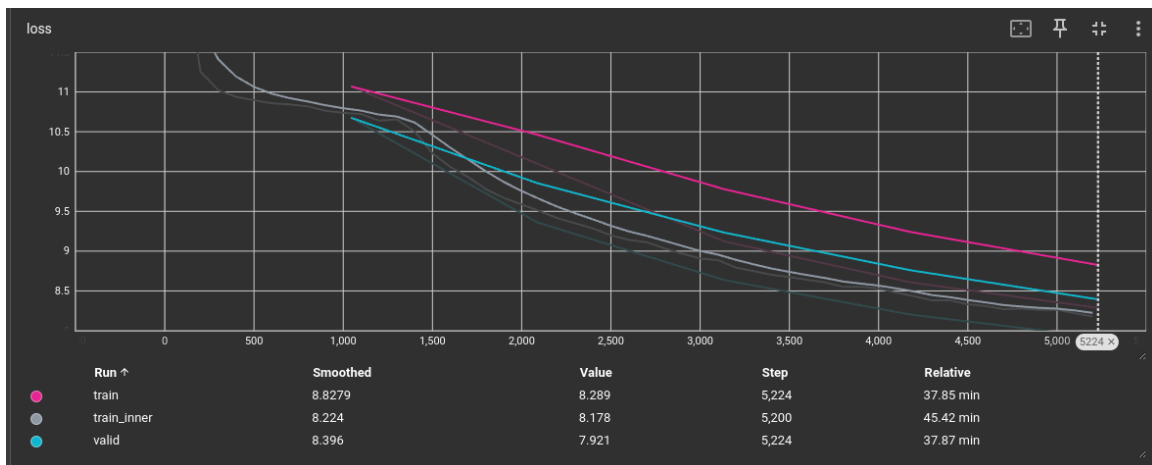


Figure 5: Training and Validation Loss

According to the loss for train set and validation set, we could have more epochs to train it more but as description 5 epochs would be enough.

4 Part 3: Training Transformer ENCODER-DECODER

We will avoid the explanation of following code for training encoder-decoder model based on the transformer since whole parameters are same as previous section. The implemented code is as follow:

```
1 !fairseq-train \  
2   "./data_bin/" \  
3   --arch transformer --share-decoder-input-output-embed \  
4   --encoder-layers 6 --decoder-layers 6 \  
5   --optimizer adam --adam-betas '(0.9, 0.98)' --clip-norm 0.0 \  
6   --lr 2e-3 --lr-scheduler inverse_sqrt --warmup-updates 4000 --warmup-init-lr 5  
7   ↪ e-08 \  
8   --dropout 0.3 --weight-decay 0.0001 \  
9   --criterion label_smoothed_cross_entropy --label-smoothing 0.2 \  
10  --max-tokens 4096 \  
11  --max-sentences 128 \  
12  --update-freq 4\  
13  --num-workers 2\  
14  --eval-bleu \  
15  --eval-bleu-args '{"beam": 5, "max_len_a": 1.2, "max_len_b": 10}' \  
16  --eval-bleu-detok moses \  
17  --eval-bleu-print-samples \  
18  --best-checkpoint-metric bleu --maximize-best-checkpoint-metric \  
19  --fp16 --memory-efficient-fp16 \  
20  --max-epoch 5 \  
21  --save-dir ./data_bin/checkpoints2/ \  
    --tensorboard-logdir ./data_bin/logs2
```

We will consider checkpoint2 and logs2 for this model. The loss for training and validation are depicted in figure 6:

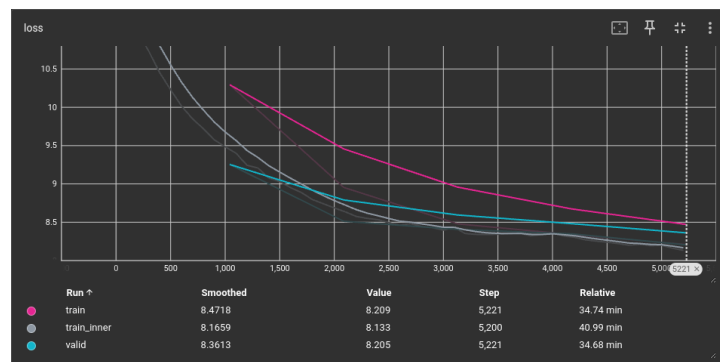


Figure 6: Transformer loss for validation and training set

The loss for validation and training has been uploaded alongside with code and report. As depicted in figure 6, but we did not yield good performance with BLEU score where will be discussed in the later sections. (Evaluation section)

5 Evaluation Metric and Examine Test dataset

In final section, we will dive into examine the performance on two score called BLEU and COMBAT scores. Where these metric's definition are as follow:

1. BLEU score: This is a metric for evaluating quality of a machine translating a language to another. More precisely, BLEU calculates the n-gram precision, between the candidate and the reference translation. These n-grams are commonly unigrams, bigrams, trigrams, 4-grams. This score equation is as follow:

$$BLEU = BP * \exp \sum_{i=1}^n w_n \log p_n$$

where p_n is the modified precision for n-grams of length n, w_n is the weight for n-grams of length n and N is the maximum length of n-grams which considered. The BP is brevity prenalty(These are according to the slide of NLP course)

2. COMBAT: This metric is based on neural network model to evaluate the quality of machine translation. Moreover, this metric will use a pre-trained multilingual encoder to predict the quality of a translation by comparing it with reference translation and the source. COMET takes triple input consist of the source sentence, the reference translation and the candidate translation. These models has been trained on the human judgement quality estimation datasets. In the end for each triplet input these scores can be averaged to produce an overall evaluation score for current model we have been training.

There is some strength in COMBAT metric since it can evaluate the model based on semantic understanding where the BLEU will not capture them. BLUE score can be indeed insensitive to word order where the COMBAT not also capture semantic similarities but also aligns more closely to human judgement(since they are trained based on human judgements!)

The best BLEU score has been reported in table 3:

Model	Best BLEU
LSTM encoder-decoder	4.76
Transformer encoder-decoder	1.34

Table 3: Best BLEU score for two trained models

These BLEU score for validation per epoch for these two model has been illustrated in figure 7 and figure 8 respectively for LSTM and Transformer:

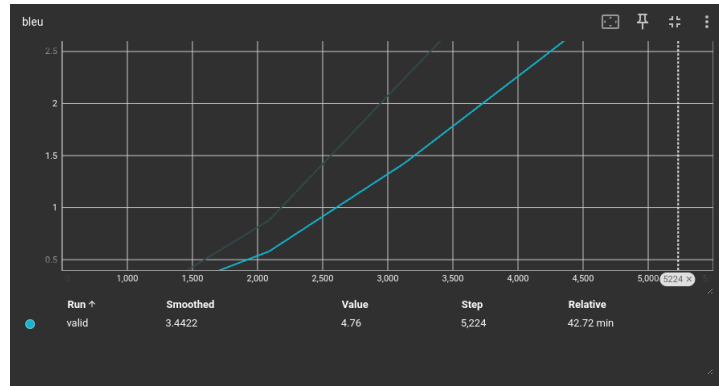


Figure 7: LSTM BLEU score per epoch of validation



Figure 8: Transformer BLEU score per epoch of validation

As you can see for lstm BLEU score is increasing where the BLEU score for transformer is not.(Maybe it is due that the same BLEU score function for these model is not good metric while we might get better score with some other metric such as COMBAT).

Attention: We have downloaded the loss for train and validation for both model with following code:

```
1 log_dir = './data_bin/logs2/train'
2
3 ea = event_accumulator.EventAccumulator(log_dir)
4 ea.Reload()
5 print("Available scalar keys:")
6 print(ea.Tags()['scalars'])
7
8 scalars = ea.Scalars('loss')
9 df = pd.DataFrame([(s.step, s.value) for s in scalars], columns=['step', 'loss'])
10 df.to_csv('transformer_training_loss.csv', index=False)
11
12 log_dir = './data_bin/logs2/valid'
13
14 ea = event_accumulator.EventAccumulator(log_dir)
15 ea.Reload()
```



```

16 print("Available scalar keys:")
17 print(ea.Tags()['scalars'])
18
19 scalars = ea.Scalars('loss')
20 df = pd.DataFrame([(s.step, s.value) for s in scalars], columns=['step', 'loss'])
21 df.to_csv('transformer_validation_loss.csv', index=False)

```

After the CSV files has been created, we will download it via 'colab' env. We will generate the the hypothesis for test set via following command:

```

1 !fairseq-generate data_bin \
2   --path data_bin/checkpoints/checkpoint_best.pt \
3   --batch-size 64 --beam 5 \
4   --results-path ./LSTM_result

```

This command simply use a batch size of 64 and beam width of 5 based on the best checkpoints the model has. The last line of generated text consist BLEU score where this score for each model is reported in table 4:

Model	BLEU score for test set
LSTM encoder-decoder	4.15
Transformer encoder-decoder	1.28

Table 4: Best BLEU score for two trained models

Again the idea of consider which model is better can not be obtained from this score since it is based on the n-grams precision. A better score where defined earlier in this section must be considered. This COMBAT score is indeed very powerful metric since, it aligns more closely with human judgements. Before jumping to the COMBAT score, the following implemented code will help us to get the BLUE score from the generated text by fairseq generate command.

```

1 def get_blue_score(file_path:str) -> float:
2     """
3     This method will get the the file path for the output of fairseq-generate
4     ↪ function and return the blue screen
5     output is blue score and its data type is float.
6     input is string contains folder path.
7     """
8     with open(file_path, 'r', encoding='utf-8') as read_file:
9         lines = read_file.readlines()
10        tokens = lines[-1].strip().split()
11        index = tokens.index('BLEU4')
12        return float(tokens[index+2].replace(',',''))

```

After this we must prepare these generated text via previous command for the combat models. These triplet inputs must be specified:

1. Source: In generated text, source will be marked as 'S-' for each line.
2. Target: The target sentence which is a reference sentence of translated source is marked as 'T-' at the first of each line.

3. Hypothesis or machine-translated target: The translated sentence by the machine translated are marked 'H-' at the beginning of each line.

Now with the following code we can have dictionary that store each of them:

```
1 def create_dictionary(file_path:str) -> dict:
2     """
3     Simply get a path file and create a dictionary based on sources and targets
4     ↪ and hypothesis exist in that file.
5     Output: dictionary contain three key: hypothesis, sources and targets.
6     Input: path for a file.
7     """
8     with open(file_path, 'r', encoding='utf-8') as read_file:
9         lines = read_file.readlines()
10    output_dict = {
11        'hypothesis': [],
12        'sources': [],
13        'targets': []
14    }
15    for line in lines:
16        tokens = line.strip().split()
17        if tokens[0][0:2] == 'S-':
18            output_dict['sources'].append(' '.join(tokens[1:]))
19        elif tokens[0][0:2] == 'T-':
20            output_dict['targets'].append(' '.join(tokens[1:]))
21        elif tokens[0][0:2] == 'H-':
22            output_dict['hypothesis'].append(' '.join(tokens[2:]))
23        else:
24            continue
25    return output_dict
```

After saving each line of sources, targets and hypothesis sentences in a dictionary data type, we can decode them with pre-trained bpe model by sentence piece library where was implemented in section Part 1. As we translating English to Persian our hypothesis and target must be decoded by the Persian pre-trained model while for the source we must decode them with English pre-trained model. The following code will illustrate the way:

```
1 def decode_each(hypo_trg_src:dict) -> dict:
2     """
3     This function will decode each with our bpe.
4     """
5     decode_dict = {key:[] for key in hypo_trg_src.keys()}
6     for hypo in hypo_trg_src['hypothesis']:
7         tokenized_hypo = hypo.strip().split()
8         decode_dict['hypothesis'].append(per_tokenizer.decode(tokenized_hypo, out_type
9             ↪ =str))
10    for src in hypo_trg_src['sources']:
11        tokenized_src = src.strip().split()
12        decode_dict['sources'].append(eng_tokenizer.decode(tokenized_src, out_type=str
13            ↪ ))
```

```

12 for trg in hypo_trg_src['targets']:
13     tokenized_trg = trg.strip().split()
14     decode_dict['targets'].append(per_tokenizer.decode(tokenized_trg, out_type=
        ↪ str
        ↪ ))
15 return decode_dict

```

Now we have a dictionary where each key in it has decoded line! Now based on implemented code for the Unbabel-Comet library we find out we must give the input in the format of a list contains dictionaries. Thus we consider to do this:

```

1 data = []
2 for src, ref, mt in zip(decoded_lstm_output['sources'], decoded_lstm_output['
    ↪ targets'], decoded_lstm_output['hypothesis']):
3     data.append({
4         "src": src,
5         "ref": ref,
6         "mt": mt
7     })

```

Now our data has been prepared for combat models to predict the scores for us. Based on the frequently asked questions part in the COMET documentation link we find out for general purpose MT evaluation they recommend us to use 'Unbabel/wmt22-comet-da'. Now utilizing the following code we will load this model in order to predict the scores for us:

```

1 from comet import download_model, load_from_checkpoint
2 model_path = download_model("Unbabel/wmt22-comet-da")
3 model = load_from_checkpoint(model_path)
4 output = model.predict(data, batch_size=8, gpus=1)

```

Now let's proceed into obtain the combat score for each model where has been reported in table 5:

Model	combat score
LSTM encoder-decoder	0.52
Transformer encoder-decoder	0.43

Table 5: COMBAT score for two trained models

As you can see in the table5, the combat score is very good for these two models but the combat score for LSTM is higher than transformer which can indicate that the encoder-decoder based on RNN(LSTM, GRU) can some times outperform the Transformer models. In the class, when Dr. Faili was teaching us machine translation, I remembered that he said, when the parallel data volume is low, encoder-decoder model based on RNN can outperform the transformers. Note that we are telling this based on just 5 epochs and for more investigating which model has better performance we must try both model trained same amount of epoch. Then evaluate them again.

Based on the two last table 5 and table 4 both metrics ranging 0 to 1. (Note that the BLEU score is the percentage). Both of these metric agree that the LSTM encoder-decoder performs better than the transformer encoder-decoder on the test set.

The BLUE score as said before focused on ngram overlap and might not fully capture semantic and contextual accuracy.

The COMET score provides a better alignment with human judgement and confirms that the LSTM model will generate higher quality translation.