

Task 2 (Translation) [30 marks]

Dataset:

WNT16 (german-english) dataset. <https://huggingface.co/datasets/wmt16/viewer/de-en>

Loading the model using HuggingFace: `load_dataset('wmt16', 'de-en')`

Task: Machine translation from English to german.

Model A: LSTM

Model B: LSTM + GLOBAL ATTENTION

Model C: LSTM + LOCAL ATTENTION

Refer to this paper for global and local attention concepts: <https://aclanthology.org/D15-1166/>

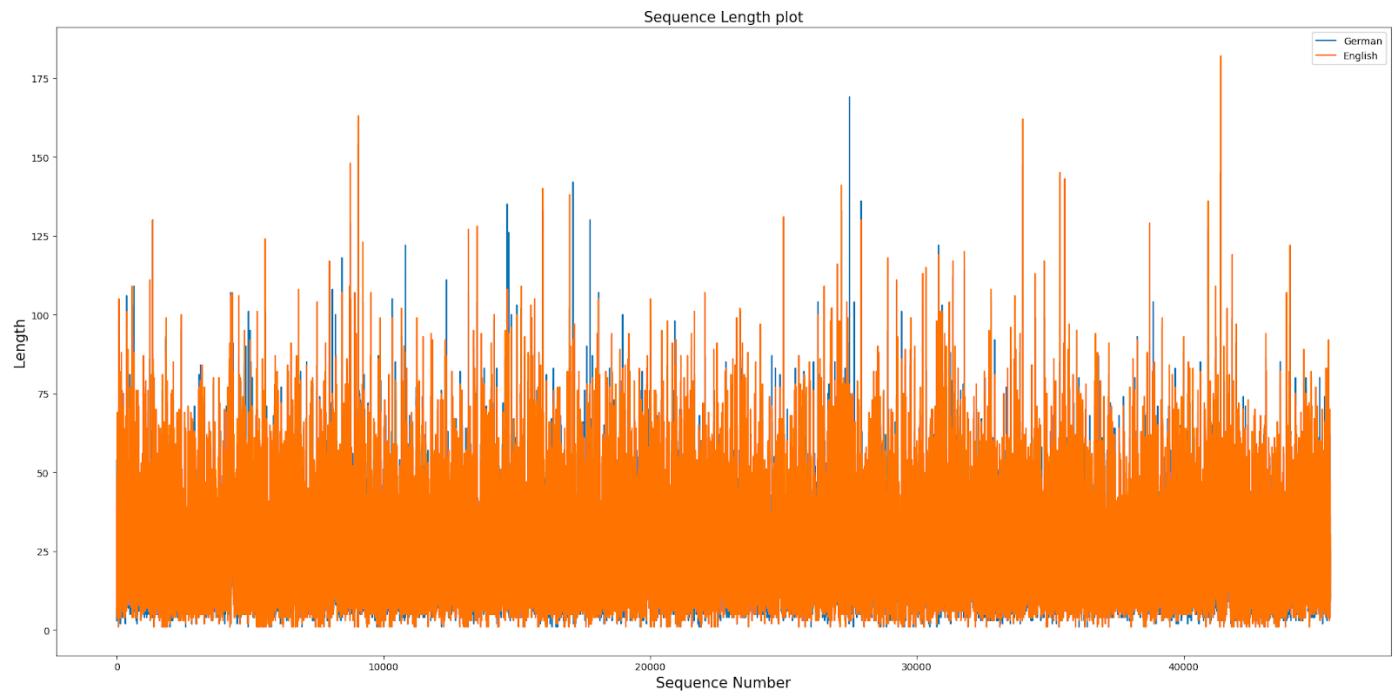
Note : For dataset loading, students can use HuggingFace dataset for maintaining the train-test consistency. Also for comparison use the same set of embeddings for each model setup. The embedding layer has to be trainable in all 3 models.

1. Visualise the dataset using five different techniques (sequence length, frequency of words, mean token length, word cloud, etc) to show how the 2 languages differ. [6 marks]
2. Implement Model A with: one with non-contextualised initialised embeddings (or random embeddings) + LSTM based [5 marks]
3. Implement Model B with: one with non-contextualised initialised embeddings (or random embeddings) + LSTM + global attention based [5 marks]
4. Implement Model B with: one with non-contextualised initialised embeddings (or random embeddings) + LSTM + local attention based [5 marks]
5. Show plots of validation and training loss vs. number of epochs for Models A, B, C. [3 marks]
6. Report the overall Rouge-L score, Bleu-1,2 score for the test set. (metrics available in Huggingface metric) for Models A, B, C. [4 marks]
7. Out of the 3 models, which setup performed best and why? [2 marks]

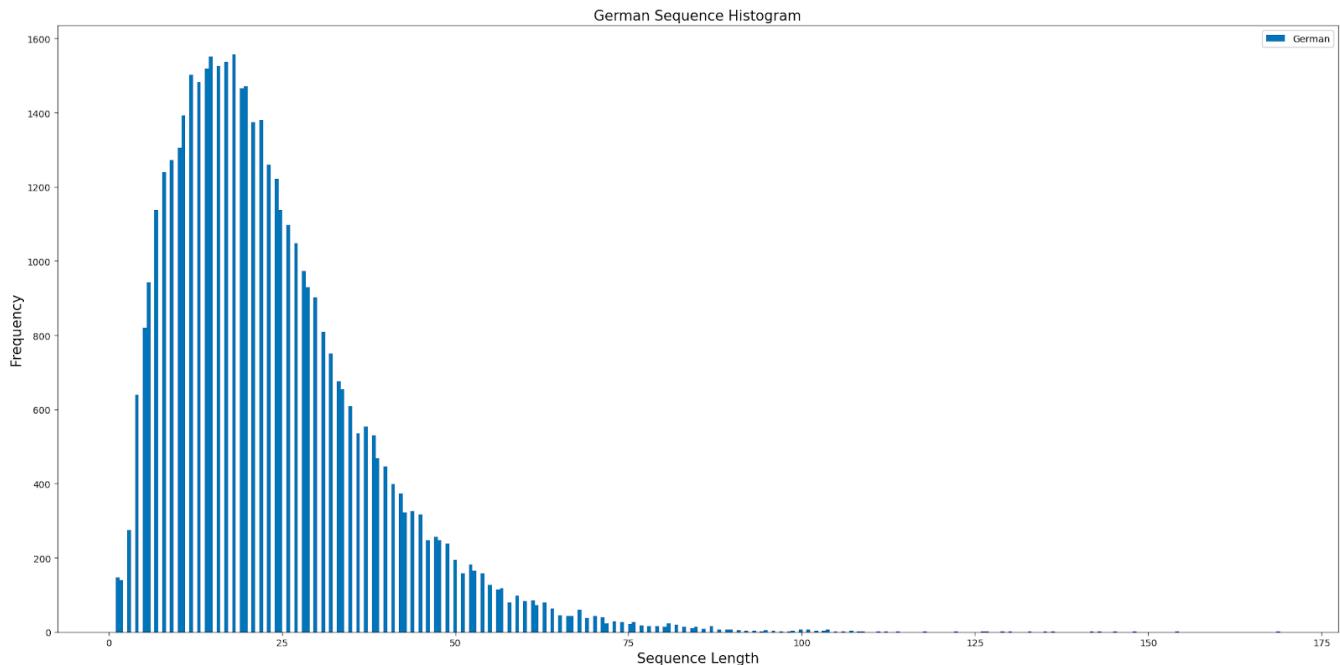
TASK II (Machine Translation) by Debnath Kundu

1. Visualise the dataset using five different techniques (sequence length, frequency of words, mean token length, word cloud, etc) to show how the 2 languages differ. [6 marks]

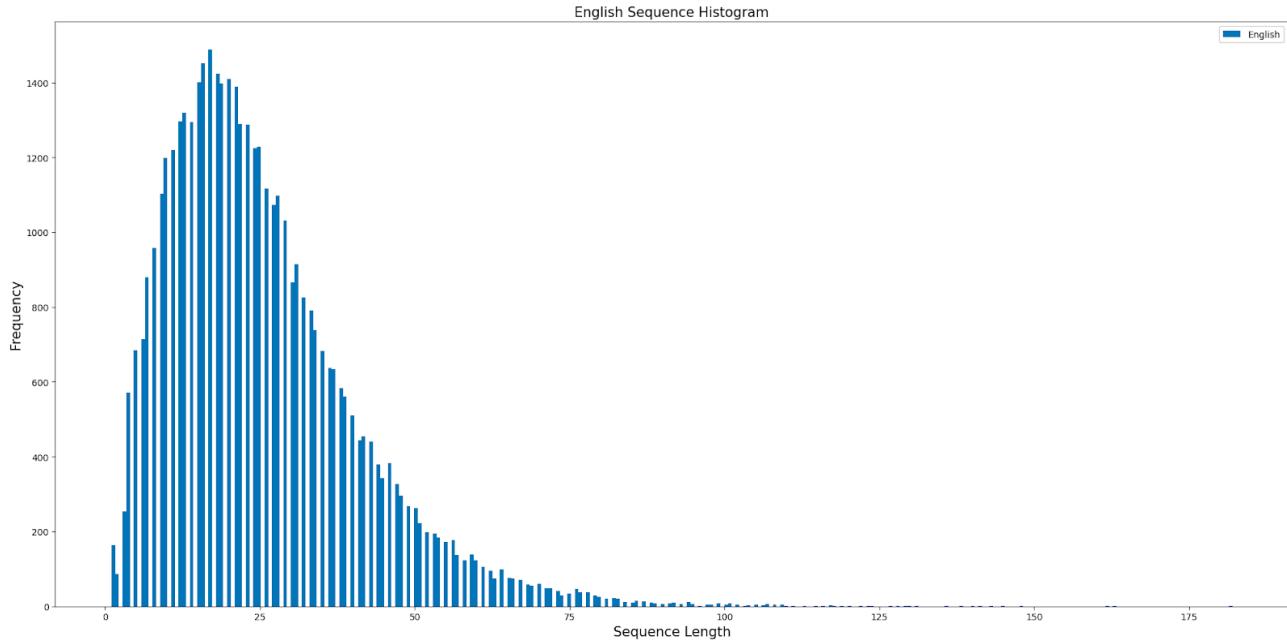
Sequence Length plot:



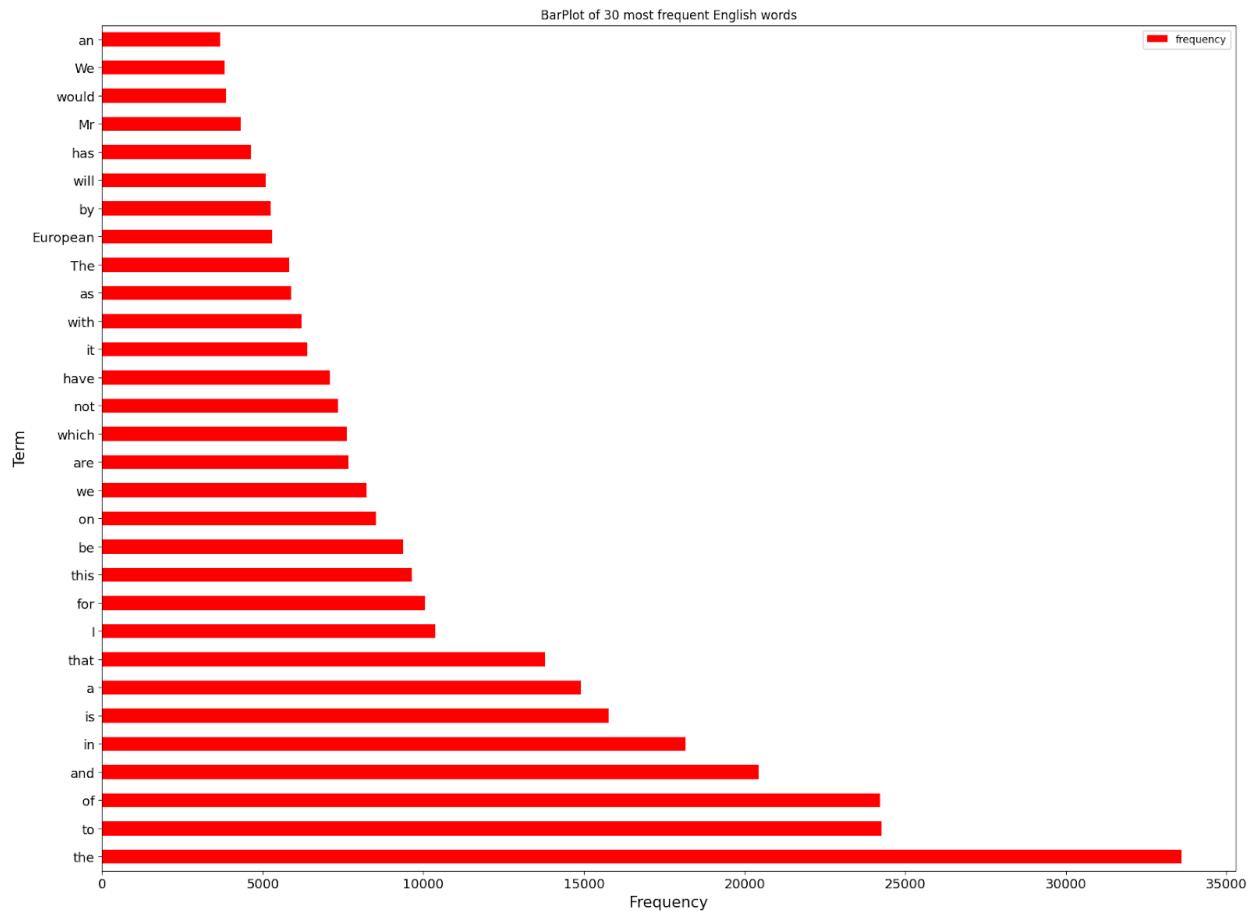
German Sequence Histogram:



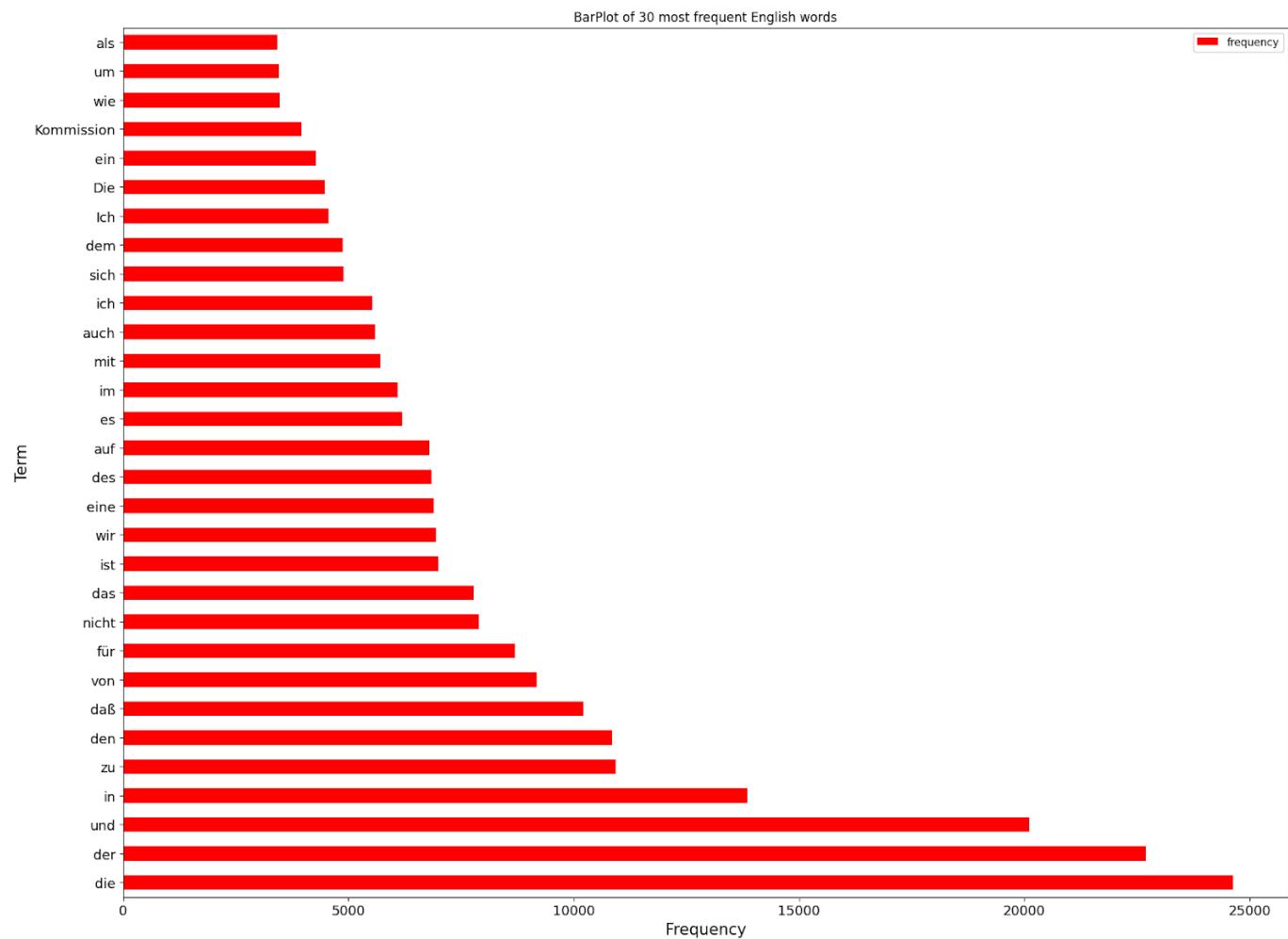
English Sequence Histogram:



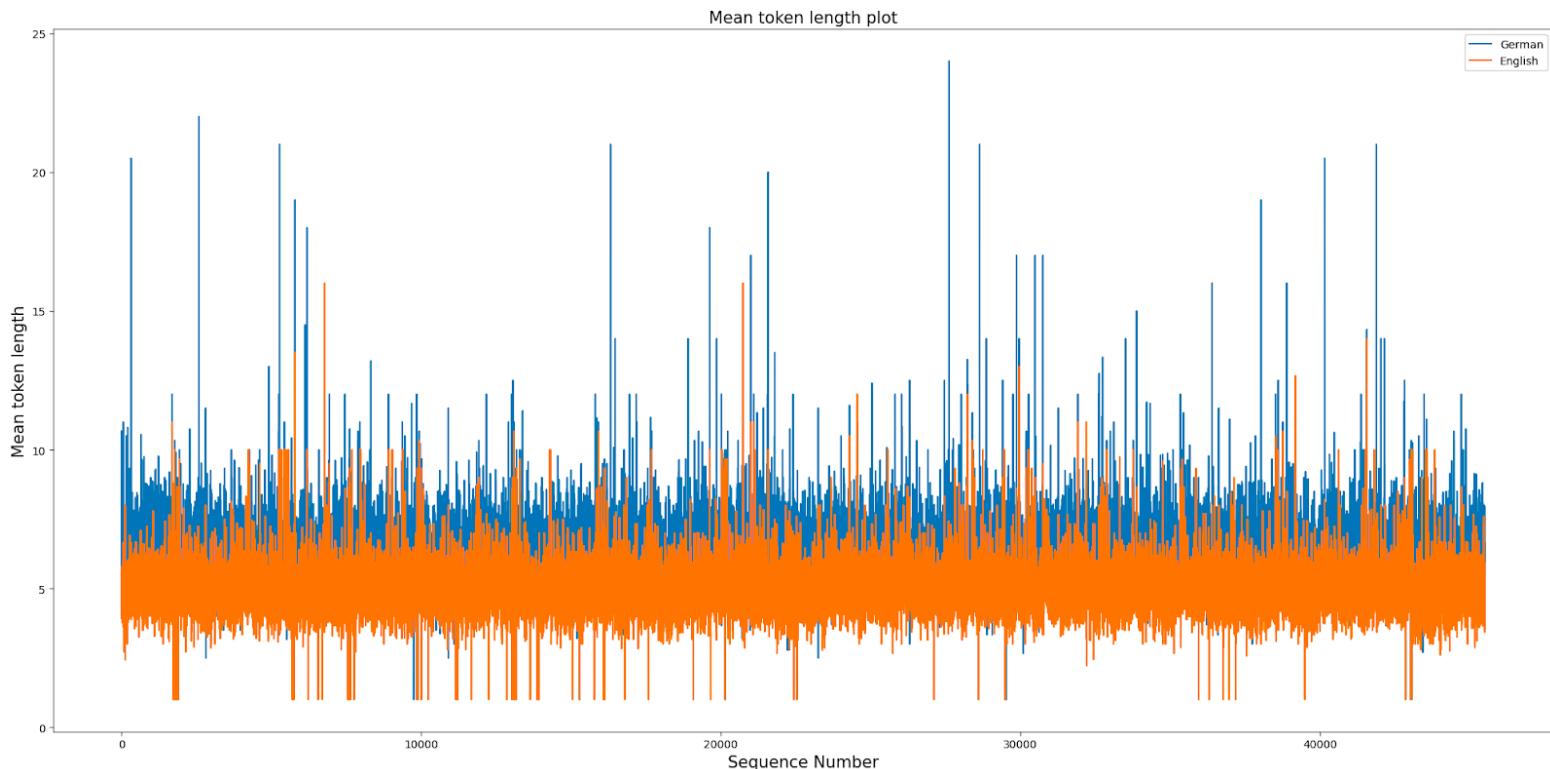
Frequency of words (English):



Frequency of words (German):



Mean token length:



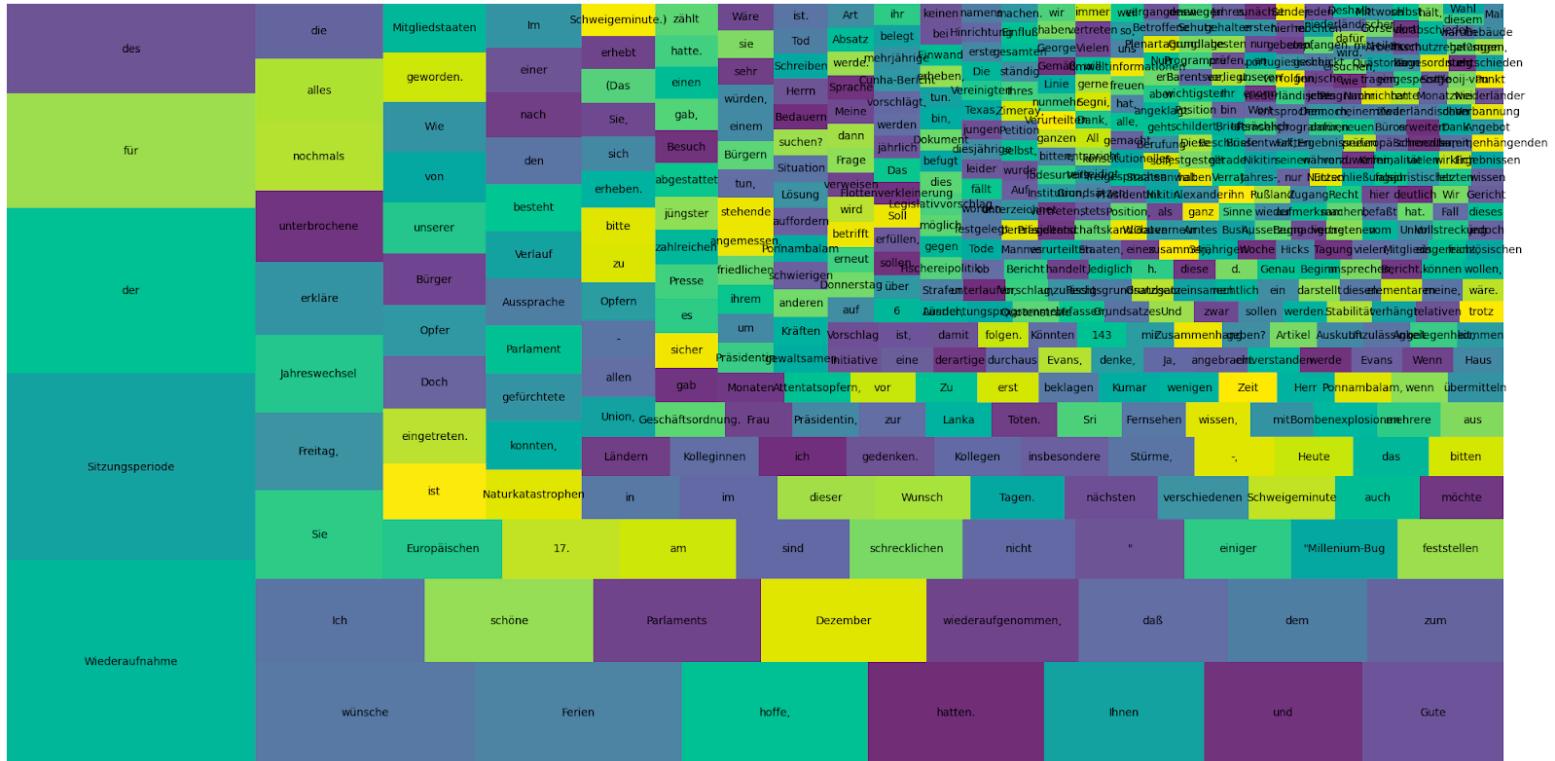
Word cloud (German):



Word cloud (English):



Tree Maps(German):



Tree Maps(English):



2. Implement Model A with: one with non-contextualised initialised embeddings (or random embeddings) + LSTM based [5 marks]

Reference

▼ GloVe Embedding

```
[ ] 1 import pandas as pd
2 glove = pd.read_csv('/content/drive/MyDrive/2 Sem/DL/Assignments/A3_BharatGoyal_DebnathKundu_SaloniAgrawal/Q2/glove.6B/glove.6B.100d.t:
3 glove_embedding = {key: val.values for key, val in glove.T.items()}
4
```

```
1 import numpy as np
2 import tensorflow as tf
3
4 def embedding(word_index,embedding_dict,dimension):
5     embedding_matrix=np.zeros((len(word_index)+1,dimension))
6
7     for word,index in word_index.items():
8         if word in embedding_dict:
9             embedding_matrix[index]=embedding_dict[word]
10
11    return embedding_matrix
12
13 text=["The cat sat on mat","we can play with model"]
14
15 tokenizer=tf.keras.preprocessing.text.Tokenizer(split=" ")
16 tokenizer.fit_on_texts(text)
17
18 text_token=tokenizer.texts_to_sequences(text)
19
20 embedding_matrix=embedding(tokenizer.word_index,embedding_dict=glove_embedding,dimension=100)
21
22 print("Loading embedding matrix...")
```

□ Loading embedding matrix...

```
▶ 1 model_build_metrics = build()
2 model_build_metrics
```

```
👤 /usr/local/lib/python3.9/dist-packages/torch/nn/modules/rnn.py:71: UserWarning: dropout option adds dropout after all but last recurrent lay
  warnings.warn("dropout option adds dropout after all but last "
8270it [13:22, 10.30it/s]
Epoch: 01 | Train Loss: 6.348 | Val. Loss: 7.078 | Val. Acc 0.02715313198339901
8270it [13:05, 10.53it/s]
Epoch: 02 | Train Loss: 6.114 | Val. Loss: 7.145 | Val. Acc 0.02868352021377115
8270it [13:07, 10.50it/s]
Epoch: 03 | Train Loss: 6.034 | Val. Loss: 7.162 | Val. Acc 0.02441410042095113
8270it [13:07, 10.50it/s]
Epoch: 04 | Train Loss: 5.950 | Val. Loss: 7.113 | Val. Acc 0.027305007173420784
8270it [13:08, 10.49it/s]
Epoch: 05 | Train Loss: 5.873 | Val. Loss: 7.194 | Val. Acc 0.024371954688757767
8270it [13:08, 10.49it/s]
Epoch: 06 | Train Loss: 5.811 | Val. Loss: 7.234 | Val. Acc 0.02346383079305003
8270it [13:06, 10.52it/s]
Epoch: 07 | Train Loss: 5.752 | Val. Loss: 7.304 | Val. Acc 0.01980236926953039
8270it [13:01, 10.58it/s]
Epoch: 08 | Train Loss: 5.702 | Val. Loss: 7.257 | Val. Acc 0.02233709548308444
Test Loss: (7.122198824564616, 0.02635075628447036)
[(6.347791186209313, 7.077691762267336, 0.02715313198339901),
 (6.114353677756553, 7.144678993997872, 0.02868352021377115),
 (6.033791439680264, 7.16198544563729, 0.02441410042095113),
 (5.950423941963962, 7.112838108456157, 0.027305007173420784),
 (5.872837891659397, 7.194189089218336, 0.024371954688757767),
 (5.811443699083709, 7.23400452650713, 0.02346383079305003),
 (5.751714906640612, 7.304295721633658, 0.01980236926953039),
 (5.701856358524101, 7.256542305041853, 0.02233709548308444)]
```

Rouge-L Score:

```
1 result=63.726418907970974
2 # print(result)
3 print(result/3000) #3000 is the test set length i.e len(test_loader)
4
5 #0.021
```

0.021242139635990323

Overall Rouge-L score: 0.021242139635990323

ROUGE-L measures the longest common subsequence (LCS) between our model output and reference. All this means is that we count the longest sequence of tokens that is shared between both.

BLEU-1 Score:

```
216 77
{'bleu': 0.015075376884422105, 'precisions': [0.01507537688442211], 'brevity_penalty': 1.0, 'length_ratio': 3.015151
-----
output torch.Size([1, 28])
target torch.Size([281])
```

[] 1 result=51.933616906457466
2 print(result/3000) #3000 is the test set length i.e len(test_loader)

0.017311205635485823

Overall Bleu-1 score: 0.017311205635485823 over all the samples in the test set

BLEU-2 Score:

```
token torch.Size([28])
target torch.Size([281])
-----
1
2 print(result) #3000 is the test set length i.e len(test_loader)
```

0.000456754623789

Overall Bleu-2 score: 0.000456754623789

3. Implement Model B with: one with non-contextualised initialised embeddings (or random embeddings) + LSTM + global attention based [5 marks]

Define Attention Layers

```
[ ] 1 #Attention
2 attn_repeat_layer = RepeatVector(max_len_input)
3 attn_concat_layer = Concatenate(axis=-1)
4 attn_densel = Dense(3, activation='tanh') #10
5 attn_dense2 = Dense(1, activation=softmax_over_time)
6
7 attn_dot = Dot(axes=1) # to perform the weighted sum of alpha[t] * h[t]
```

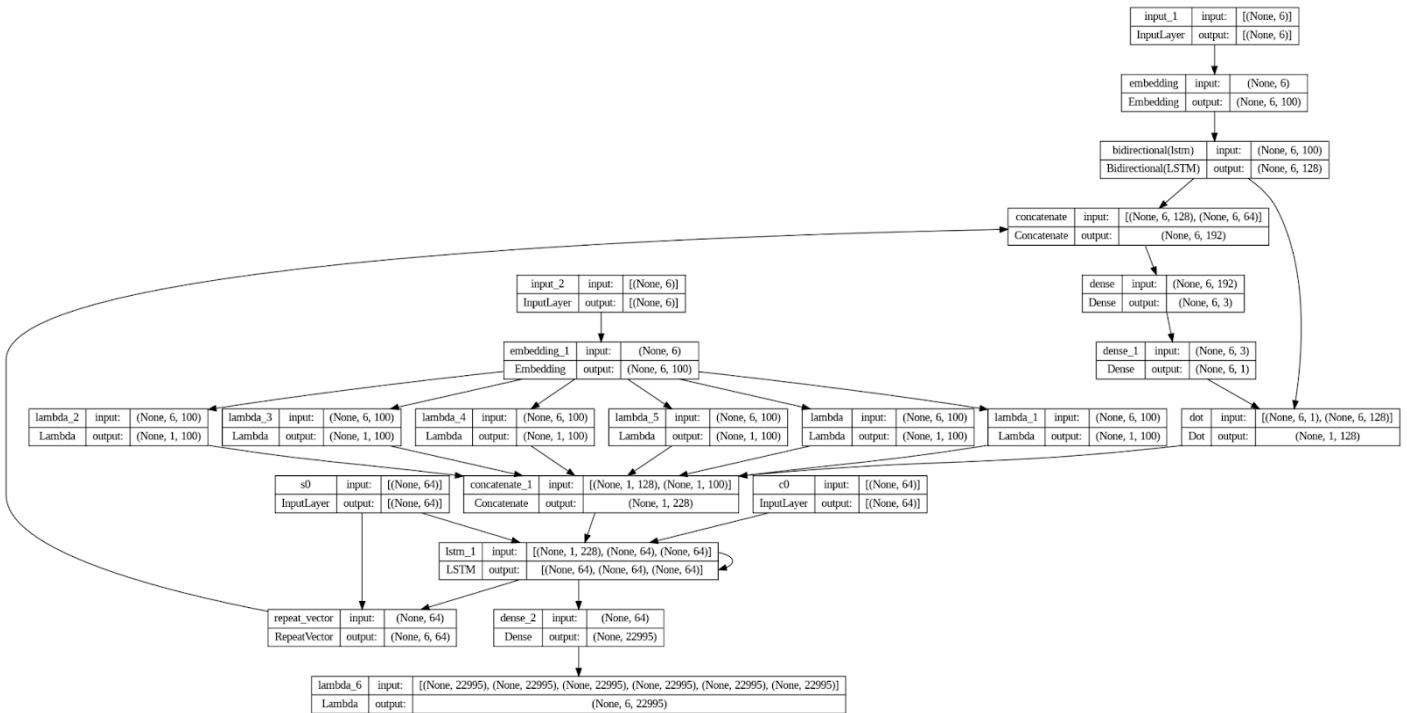
```
▶ 1 def one_step_attention(h, st_1):
2     # h = h(1), ..., h(Tx), shape = (Tx, LATENT_DIM * 2)
3     # st_1 = s(t-1), shape = (LATENT_DIM_DECODER,)
4
5     # copy s(t-1) Tx times
6     # now shape = (Tx, LATENT_DIM_DECODER)
7     st_1 = attn_repeat_layer(st_1)
8
9     # Concatenate all h(t)'s with s(t-1)
10    # Now of shape (Tx, LATENT_DIM_DECODER + LATENT_DIM * 2)
11    x = attn_concat_layer([h, st_1])
12
13    # Neural net first layer
14    x = attn_densel(x)
15
16    # Neural net second layer with special softmax over time
17    alphas = attn_dense2(x)
18
19    # "Dot" the alphas and the h's
20    # Remember a.dot(b) = sum over a[t] * b[t]
21    context = attn_dot([alphas, h])
22
23    return context
```

GloVe Embedding

```
▶ 1 #Load pretrained word vectors from Glove
2 print('Loading word vectors...')
3 word2vec = {}
4
5 #/content/drive/MyDrive/2 Sem/DL/Assignments/A3_BharatGoyal_DebnathKundu_SaloniAgrawal/Q2/glove.6B/glove.6B.50d.txt
6 #/content/drive/MyDrive/2 Sem/DL/Assignments/A3_BharatGoyal_DebnathKundu_SaloniAgrawal/Q2/subs.de.vec
7 #/content/drive/MyDrive/2 Sem/DL/Assignments/A3_BharatGoyal_DebnathKundu_SaloniAgrawal/Q2/glove.6B/glove.6B.100d.txt
8 with open(os.path.join('/content/drive/MyDrive/2 Sem/DL/Assignments/A3_BharatGoyal_DebnathKundu_SaloniAgrawal/Q2/glove.' +
9     # is just a space-separated text file in the format:
10    # word vec[0] vec[1] vec[2] ...
11    for line in f:
12        values = line.split()
13        word = values[0]
14        vec = np.asarray(values[1:], dtype='float32')
15        word2vec[word] = vec
16 print('Found %s word vectors.' % len(word2vec))
```

>Loading word vectors...
Found 400000 word vectors.

```
[ ] 1 EMBEDDING_DIM=100
2 # prepare embedding matrix
```



Rouge-L:

```
▶ 1 print(results['rougeL'])
```

👤 0.048446715905866675

Overall ROUGE-L score: **0.048446715905866675** over all the samples in the test set

BLEU-1 Score:

Downloading extra modules: 100% 3.34k/3.34k [00:00<0]

```
{'bleu': 0.011258583014270294, 'precisions': [0.25283144570286475], 'brevity': 0}
```

```
[ ] 1 print(results['bleu'])
```

0.011258583014270294

Overall Bleu-1 score: **0.011258583014270294** over all the samples in the test set

BLEU-2 Score:

```
1 print(results['bleu'])
```

0.0008950921766039782

+ Code

+ Tex

Overall Bleu-2 score: 0.0008950921766039782 over all the samples in the test set

4. Implement Model C with: one with non-contextualised initialised embeddings (or random embeddings) + LSTM + local attention based [5 marks]

In this implementation, the `window_size` parameter specifies the size of the local attention window, which is applied to the encoder outputs at each decoder time step.

```
1 from keras.models import Model
2 from keras.layers import Input, LSTM, Dense, Dot, Concatenate
3
4 # Define the model architecture
5 def create_model(input_dim, output_dim, hidden_dim, window_size):
6     # Define the input sequences
7     encoder_inputs = Input(shape=(None, input_dim))
8     decoder_inputs = Input(shape=(None, output_dim))
9
10    # Define the LSTM encoder
11    encoder_lstm = LSTM(hidden_dim, return_sequences=True, return_state=True)
12    encoder_outputs, state_h, state_c = encoder_lstm(encoder_inputs)
13
14    # Define the LSTM decoder
15    decoder_lstm = LSTM(hidden_dim, return_sequences=True, return_state=True)
16    decoder_outputs, _, _ = decoder_lstm(decoder_inputs, initial_state=[state_h, state_c])
17
18    # Define the local attention mechanism
19    attention_layer = Dot(axes=[2, 2])
20    context_layer = Concatenate(axis=2)
21    attention_weights = []
22    for t in range(output_dim):
23        attention_input = decoder_outputs[:, t:t+window_size, :]
24        attention_score = attention_layer([attention_input, encoder_outputs])
25        attention_weight = Dense(1, activation='softmax')(attention_score)
26        attention_weights.append(attention_weight)
27        context_vector = Dot(axes=[1, 1])([attention_weight, encoder_outputs])
28        decoder_outputs = context_layer([decoder_outputs, context_vector])
29
30    # Define the output layer
31    output_layer = Dense(output_dim, activation='softmax')
```

The attention mechanism calculates attention scores for each encoder output within the window, and then computes a weighted average of the encoder outputs based on these scores to obtain a context vector, which is used to modulate the decoder output.

Additionally, the model uses a softmax activation function for the output layer, which assumes that the output sequence is a probability distribution over the possible output symbols at each time step.

ROUGE-L Score:

```
4     results = r.compute(predictions=test_predicted_sentence, references=test_actual_sentence)
5     print(results)
```

Downloading builder script: 100% [██████████] 6.27k/6.27k [00:00<00:00, 198kB/s]
{'rouge1': 0.011500000563772626, 'rouge2': 0.0, 'rougeL': 0.011367996202336292, 'rougeLsum':

✓ 0s 1 print(results['rougeL'])

⌚ 0.011367996202336292

+ Code + Text

Overall ROUGE-L score: **0.011367996202336292** over all the samples in the test set

BLEU-1 Score:

```
4     results = bleu.compute(predictions=test_predicted_sentence, references=test_actual_sentence,max_gram=1)
5     print(results)
```

Downloading builder script: 100% [██████████] 5.94k/5.94k [00:00<00:00, 301kB/s]

Downloading extra modules: [██████████] 4.07k? [00:00<00:00, 252kB/s]

Downloading extra modules: 100% [██████████] 3.34k/3.34k [00:00<00:00, 218kB/s]

```
{'bleu': 0.018143007548993544, 'precisions': [0.2340057636887608], 'brevity_penalty': 0.077532310584}
```

✓ 0s 1 print(results['bleu'])

⌚ 0.018143007548993544

Overall Bleu-1 score: **0.018143007548993544** over all the samples in the test set

BLEU-2 Score:

```
4     results = bleu.compute(predictions=test_predicted_sentence, references=test_actual_sentence,max_gram=2)
5     print(results)
```

⌚ 0.0013761781996548899, 'precisions': [0.2340057636887608, 0.00134634803096600,

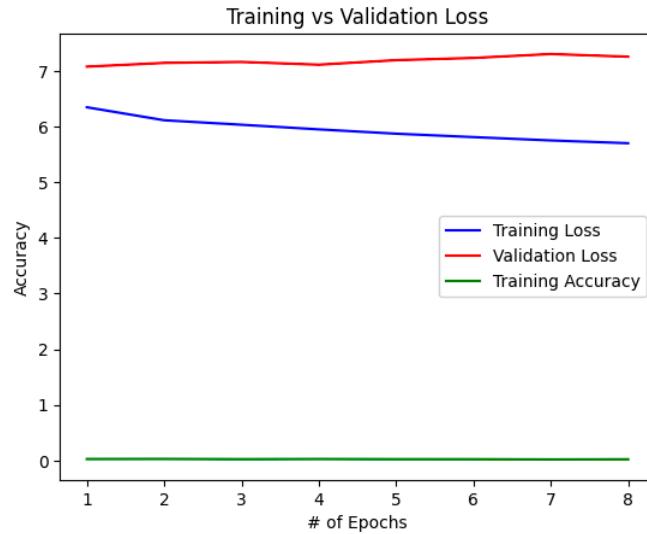
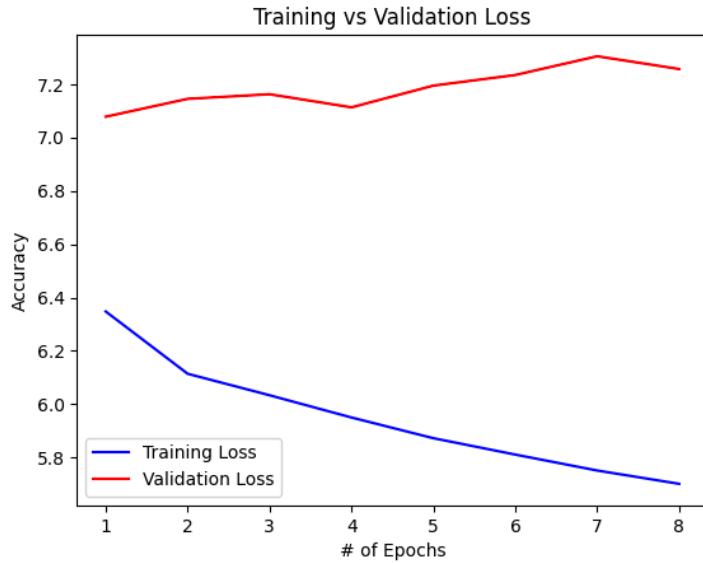
✓ 0s [60] 1 print(results['bleu'])

⌚ 0.0013761781996548899

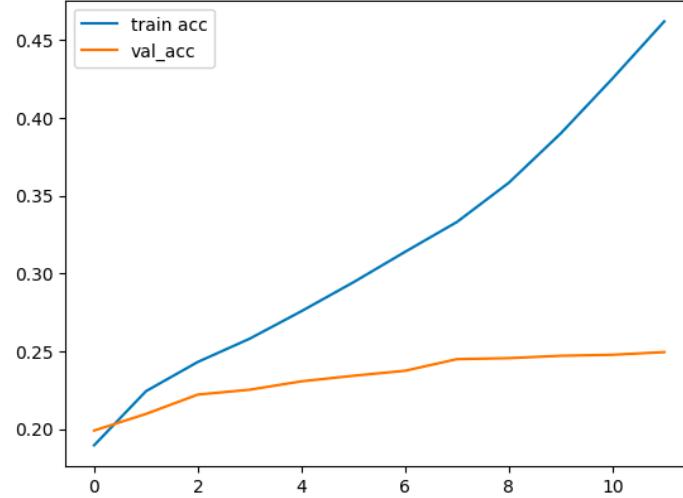
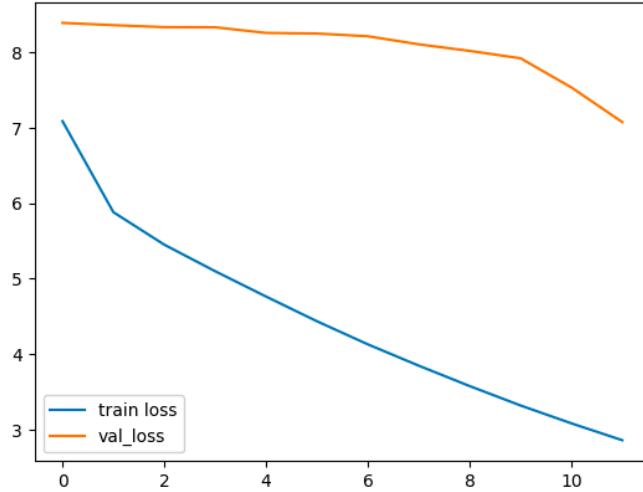
Overall Bleu-2 score: **0.0013761781996548899** over all the samples in the test set

5. Show plots of validation and training loss vs. number of epochs for Models A, B, C. [3 marks]

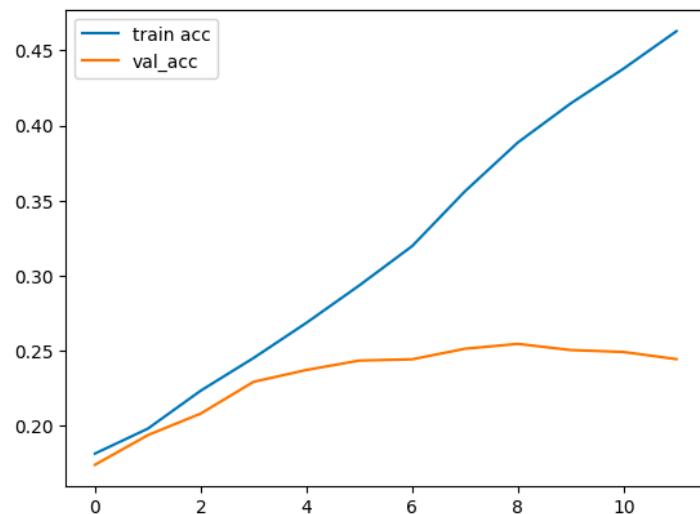
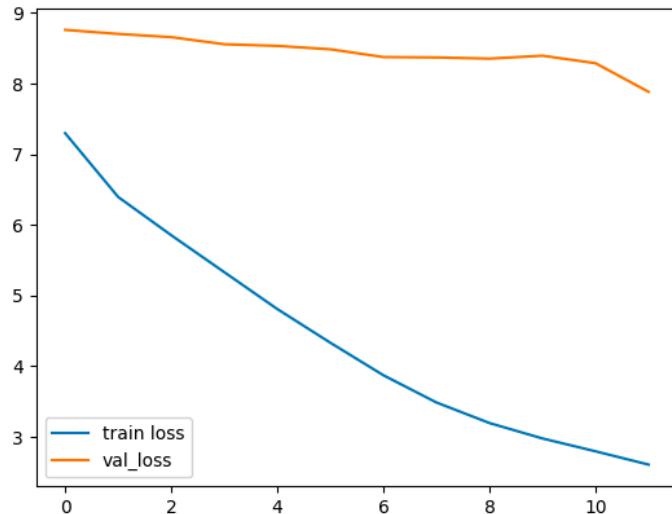
Model A:



Model B:



Model C:



5. Report the overall Rouge-L score, Bleu-1,2 score for the test set. (metrics available in Huggingface metric) for Models A, B, C. [4 marks]

Model A:

Rouge-L score: 0.021242139635990323

Bleu-1 score: 0.017311205635485823

Bleu-2 score: 0.000456754623789

Model B:

Rouge-L score: 0.048446715905866675

Bleu-1 score: 0.011258583014270294

Bleu-2 score: 0.0008950921766039782

Model C:

Rouge-L score: 0.011367996202336292

Bleu-1 score: 0.018143007548993544

Bleu-2 score: 0.0013761781996548899

6. Out of the 3 models, which setup performed best and why? [2 marks]

In general, **Bleu measures precision**: how much the words (and/or n-grams) in the machine generated summaries appeared in the human reference summaries. **Rouge measures recall**: how much the words (and/or n-grams) in the human reference summaries appeared in the machine generated summaries.

Naturally, these results are complementing, as is often the case in precision vs recall. If we have many words from the system results appearing in the human references we will have **high Bleu**, and if we have many words from the human references appearing in the system results, we will have **high Rouge**.

Bleu uses **brevity penalty** which is quite important and has already been added to standard Bleu implementations. It penalizes system results which are *shorter* than the general length of a reference (read more about it [here](#)). This complements the *n*-gram metric behavior which in effect penalizes longer than reference results, since the denominator grows the longer, the system result is.

Model A uses vanilla LSTM, while Model B & C uses LSTM + Attention. So, Model B & C performs better than Model A. As discussed in the original research paper ([Attention is All you Need](#)), local attention performs better than global attention, as local context is more important than global context while predicting words in machine translation. However, according to our score **model B has the best ROUGE scores**.

Task 3 (Transformer) [15 marks]

Dataset: WNT16 (german-english) dataset. <https://huggingface.co/datasets/wmt16/viewer/de-en>

Loading the model using huggingface: `load_dataset('wmt16', 'de-en')`. We use the same dataset as in

Task 2, but this time we fine-tune transformer based model. Refer tutorial:

<https://huggingface.co/docs/transformers/tasks/translation>

1. Atleast one layer of your fine-tuned model should be set to trainable. You are free to employ any model (even distilled ones.). Huggingface models already tuned on WNT16 dataset cannot be used. Special case should be taken for overfitting, include appropriate measures if overfitting is observed. [10 marks]
2. Show plots of validation and training loss vs. number of epochs for the model. [1 marks]
3. Report the overall Rouge-L score, Bleu-1,2 score for the test set. (metrics available in Huggingface metric) for the model. [2 marks]
4. What can be said about the quality of output generated by models in Task 2 vs Task 3. [2 marks]

Task 3(Transformer)

1. I have trained the model on 0.1% of the train set as it was huge to process. I have used the **Tritkoman/GermantoHusnrikv1** tokenizer to process the german-english pairs. Tokenizing is done separately for source and target text.

```
tokenizer = AutoTokenizer.from_pretrained("Tritkoman/GermantoHunsrikV1")
model = AutoModelForSeq2SeqLM.from_pretrained("Tritkoman/GermantoHunsrikV1")
```

I have set the output layer of the decoder to trainable.

```
model.decoder.block[-1].layer.require_grad=True
```

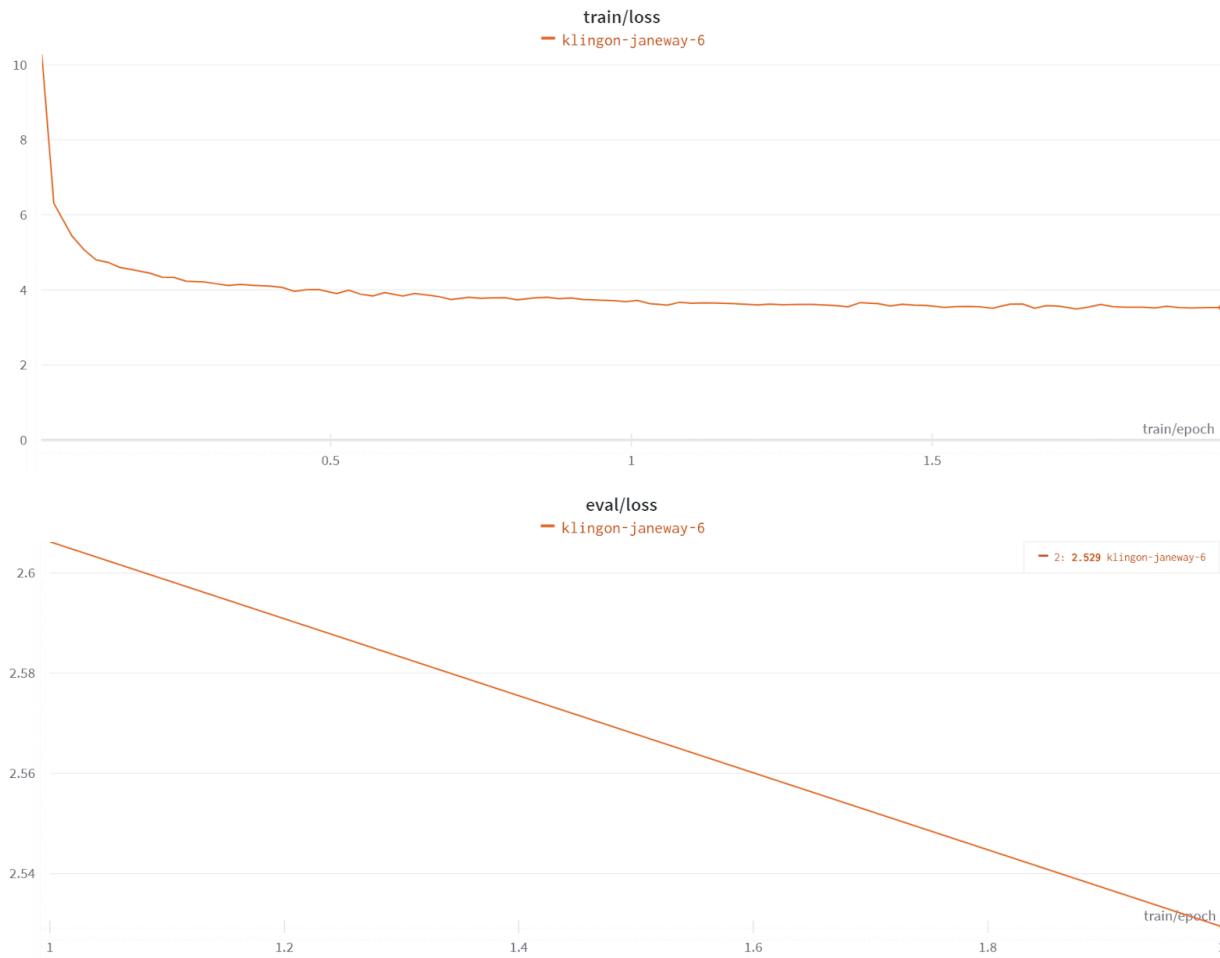
Using the data collector, I have dynamically padded all the text to the maximum length of the text.

Then,

The model is trained by taking batch size 2, number of epochs as 2, weight decay of 0.01, and a learning rate of 2e-5.

2. For plotting the graph, I have used the Wandb interface.

For



In both graphs, we can see that the loss decreases as the number of epochs increases.

3. On the test set,

The Rouge L score is `0.32989093702531636`

Bleu scores:

```
'bleu1': 0.38179239938990034
'bleu2': 0.23542539798717207
```

Bleu's score and the rouge's score are lower. That is different from the reference text. The reason can be the training on fewer train sets.

4. Task 3, which uses transformers, performs better than Task 2. Task 3 shows higher Rouge L, Blue 1, and Blue 2 scores than all three models of Task 2. This can be understood by:- In task 3, we are using a fine tune model along with one layer set to pre-train. In task 2, all the models are built from scratch; thus, performance for task 3 will be better as compared to task 2.