

Translation MODEL: NN and DL Project

By Ved Anand (12021002018015)

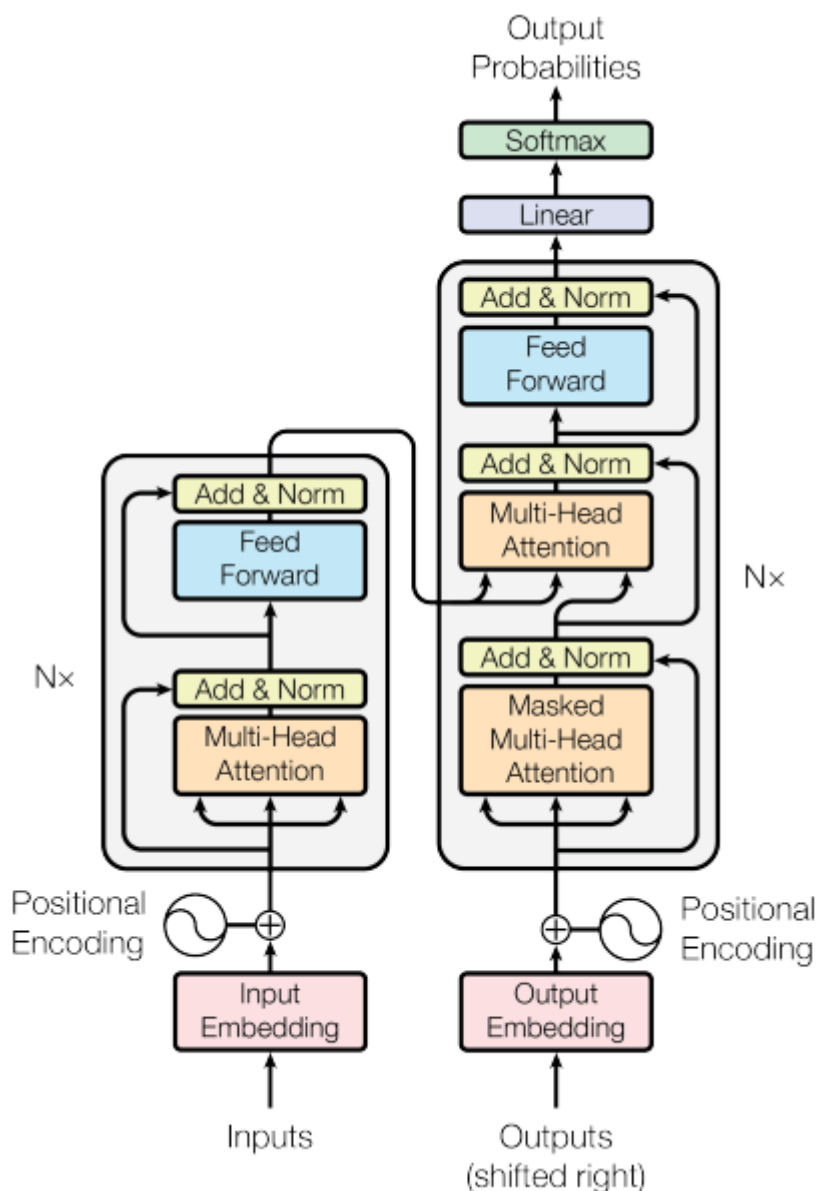
- Dataset Used: Multi-30K Dataset

- We are going to train our model on the Multi 30k data set from English to German.
- This data set already comes in a tokenized format for German, English, French.
- So, this is a very nice data set.
- This can be taken as a toy data set to implement and see how good our model works.
- Dataset Link: <https://github.com/multi30k/dataset>

- Model Used: Standard Transformer Model

- For the model, we'll use the Standard Transformer Model architecture.
- For a Good tutorial series on implementing seek to seek architectures from scratch in bluepython, one can refer to: [Ben Trevett Implementation](#)
- So, if one wants to understand the more minute details behind the Transformer architecture, they can refer to above link.

- Decoding Colab Notebook File: [TranslationModel.ipynb](#)



➤ **Importing the Dataset (Multi-30K):**

- So, our first step would be to get the data set.
- Get the data set from [here](#).
- And then we install a couple of python modules.

➤ **Installing Required Python Modules:**

- One is WandB for login into [wandb.ai](#),
- And Sacremoses for tokenizing text,
- And fair-seq for training and evaluating our model.

➤ **To use W and B, Creating and Logging in Account**

- If you want to use wnb, you have to create your account and then give the command specified in Notebook.
- Authorize your login Credentials by putting your API key from WandB Profile.
- API key for this project: [b910992fda2c68afc731d31d0623442d3bbe5df0](#)
- And it directly shows up to your particular project, all the parameters, all the values that it keeps track of during the training process.

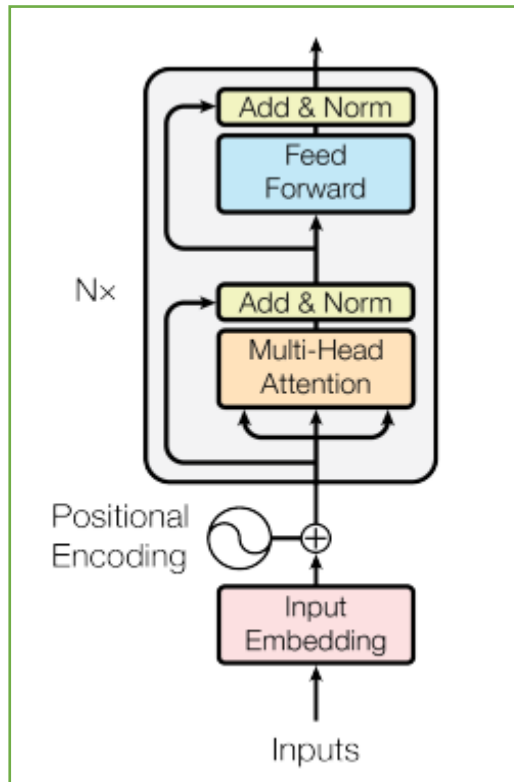
➤ **Pre-process and Binarize to build Vocabularies:**

- After that we have logged into our wnb account, the first step is the preprocessing step.
- This step takes all the tokenized data in our data set, which is a multi 30k data set, which is there in the task1 under tokenized folder.
- It takes all the text and then try to learn vocabulary from it and then binarize the data so that it becomes faster during the training to batch them and use them.
- The parameters that we are passing for the fair-seq preprocess is the source language that we want to translate and the target language that we want to translate it to.
- The train prefix, valid prefix and test prefix are the entire path of the train data set, validation data set and the test data set except the last extensions for the language.
- For the test preference, we are using the “2018 flickr Data set”.
- The destination directory for all the vocabulary learnt by the fair-seq preprocess and the binary data set are stored in data bin.
- Threshold SRC and TGT what they specify is: The minimum number of times (minimum threshold) a token repeats so that it is added into the vocabulary.

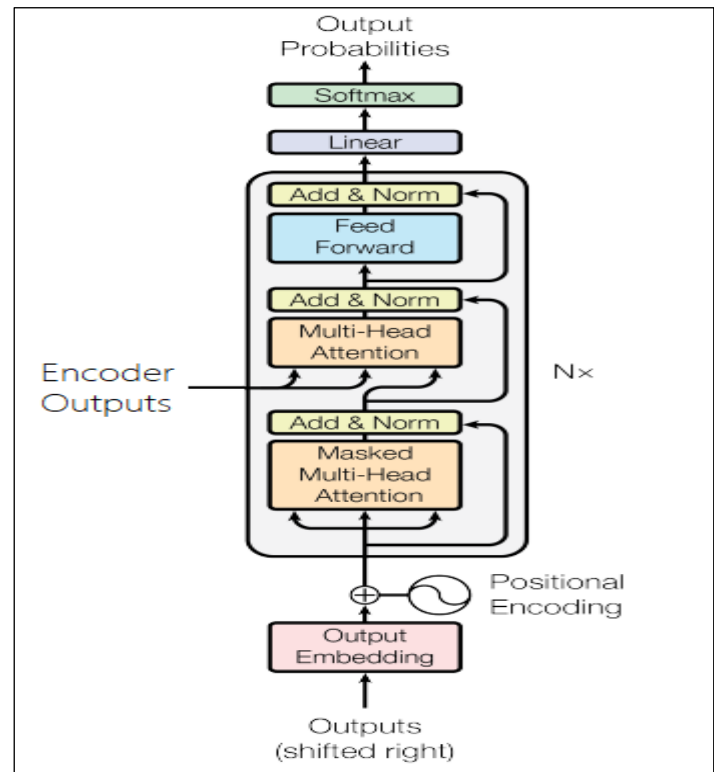
➤ **Training the model:**

- In the train step, we have specified a lot of parameters. Let's discuss them one by one.
- So, most of the hyperparameters are chosen or directly taken from [Ben Trevett tutorial](#).
- We will try explain what they do with the help of the [fairseq documentation](#).
- Dropout, attention-dropout and activation-dropout are three types of dropouts in the entire transformer architecture.
- Attention dropout is done after the data passes through the attention layer and the other one in the above diagram in the feed forward network.
- There are two feed forwards, and there are three attentions that you can see in above diagram and then the activation dropout comes.
- In the below picture architecture, the feed forward part block is made of two layers.
- The size of the hidden layer in the middle the second layer is the hyperparameter that we are setting here to 512.

- Number of encoder layers is the repetition of each encoder layer, i.e., all these encoder layers are stacked on top of each other and all the decoder layers are stacked on top of each other.
- That is what this Nx in the picture signifies.
- So, we have used three encoder layers, and similarly, three decoder layers and it is important that “The number of encoder layers = The Number of decoder layers”.



Encoder



Decoder

- And the rest parameters/commands are referred from [fairseq documentation](#):
- **Evaluating the Model:**

Additional command-line arguments

--dropout	dropout probability
--encoder-embed-dim	encoder embedding dimension
--encoder-embed-path	path to pre-trained encoder embedding
--encoder-layers	encoder layers [(dim, kernel_size), ...]
--decoder-embed-dim	decoder embedding dimension
--decoder-embed-path	path to pre-trained decoder embedding
--decoder-layers	decoder layers [(dim, kernel_size), ...]
--decoder-out-embed-dim	decoder output embedding dimension
--decoder-attention	decoder attention [True, ...]
--share-input-output-embed	share input and output embeddings (requires --decoder-out-embed-dim and --decoder-embed-dim to be equal) Default: False

- 'fairseq-generate' is a command within the fairseq Library, which can be used to generate translations.
- After running the translations of our model over the test-data set, at the end we would get a Bleu Score.
- We are running the fairseq-generate on the checkpoint_best.
- Checkpoint_best is the checkpoint where the validation loss is the lowest.
- On the bleu score we are getting 25.99 which is actually a very decent Bleu Score if compared to the other state of the art models.
- Being a NOOB, our MODEL is not fully optimized for all our hyperparameters to fit the data perfectly.
- Given that the highest bleu score is 38.4 in 2020, we have still managed to reach 25.99 with very less amount of code. So, this even is a great plus point for US.
- Meanwhile we can improve the bleu score further using a few more techniques like:
 - ❖ Beam size is a technique which is used during the decoding of the transformer model.
 - ❖ Increasing the Beam size from 5 to 10, **BLEU SCORE increased from 25.99 to 26.17.**
 - ❖ Another great technique that we have that is **Averaging Checkpoints.**
 - ❖ What it does is, it takes all the Last five checkpoints and take an average of it and save it to a different file and in most of the time the average checkpoint acts like an ensemble of all the last five models, and hence provides better results usually.
 - ❖ After performing it, we got a **BLEU SCORE RISE from 26.17 to 26.74** for beam of 5.
 - ❖ **Our Final BLEU SCORE is = 26.74 and Highest in World is 49.3.**

➤ Exporting the Model:

- In order to export the model, we need to have a dictionary of both the languages and then just our model checkpoint.
- And while exporting we'll exporting the Averaged Checkpoint Model of beam 5 as for that we got the highest Bleu Score.
- Then, there is another command which is called as fairseq-interactive in which you can type sentences as inputs as you go and it will translate and keep running.

The screenshot shows a code editor with the following command in the terminal:

```
fairseq-interactive \
--path trained_model/model.pt \
--source-lang de --target-lang en \
--tokenizer mooses \
--task translation --cpu \
--beam 5 \
trained_model/
```

The output shows various TensorFlow warnings and debug messages, followed by the translation results for the input sentence "eine frau mit lockigen haaren und einer blauen jacke geht auf einem bürgersteig in der stadt an einem schaufenster vorbei .".

```
... 2023-11-24 15:36:53.654352: E tensorflow/compiler/xla/stream_executor/cuda/cuda_dnn.cc:9342] Unable to register cuDNN factory: Attempting to register factory for plugin cuDNN when one
2023-11-24 15:36:53.654432: E tensorflow/compiler/xla/stream_executor/cuda/cuda_fft.cc:609] Unable to register cuFFT factory: Attempting to register factory for plugin cuFFT when one h
2023-11-24 15:36:53.654477: E tensorflow/compiler/xla/stream_executor/cuda/cuda_blas.cc:1518] Unable to register cuBLAS factory: Attempting to register factory for plugin cuBLAS when c
2023-11-24 15:36:53.665368: I tensorflow/core/platform/cpu_feature_guard.cc:182] This TensorFlow binary is optimized to use available CPU instructions in performance-critical operation
To enable the following instructions: AVX2 FMA, in other operations, rebuild TensorFlow with the appropriate compiler flags.
2023-11-24 15:36:55.087094: W tensorflow/compiler/tf2tensorrt/utils/py_utils.cc:38] TF-TRT Warning: Could not find TensorRT
DEBUG:hydra.core.utils:Setting JobRuntime:name=UNKNOWN_NAME
DEBUG:hydra.core.utils:Setting JobRuntime:name=utils
INFO:fairseq_cli.interactive:{'name': None, 'common': {'name': None, 'no_progress_bar': False, 'log_interval': 100, 'log_format': None, 'log_file': None, 'aim_repo': None, 'aim_run_h
INFO:fairseq.tasks.translation:[de] dictionary: 7864 types
INFO:fairseq.tasks.translation:[en] dictionary: 5928 types
INFO:fairseq_cli.interactive:loading model(s) from trained_model/model.pt
INFO:fairseq_cli.interactive:NOTE: hypothesis and token scores are output in base 2
INFO:fairseq_cli.interactive:Type the input sentence and press return:
eine frau mit lockigen haaren und einer blauen jacke geht auf einem bürgersteig in der stadt an einem schaufenster vorbei .
INFO:fairseq.tasks.fairseq_task:can_reuse_epoch_itr = True
INFO:fairseq.tasks.fairseq_task:reuse_data_loader = True
INFO:fairseq.tasks.fairseq_task:rebuild_batches = False
INFO:fairseq.tasks.fairseq_task:creating new batches for epoch 1
S-0  eine frau mit lockigen haaren und einer blauen jacke geht auf einem bürgersteig in der stadt an einem schaufenster vorbei .
W-0  0.176 seconds
H-0  -0.40713638067245483 a woman with curly curly hair and a blue jacket walks past a window on a city sidewalk .
D-0  -0.40713638067245483 a woman with curly curly hair and a blue jacket walks past a window on a city sidewalk.
P-0  -0.0750 -0.0205 -0.0964 -0.0246 -0.5516 -0.4122 -0.1106 -0.0484 -0.0083 -0.0019 -1.2137 -0.3577 -0.0364 -1.4694 -0.9580 -1.1975 -0.1993 -0.0275 -1.3336 -0.0002
```

The status bar at the bottom indicates the execution time is 3m 43s and the current line is 1.