

Homework Assignment No. 13:

HW No. 13: Transformers

submitted to:

Professor Joseph Picone
ECE 8527: Introduction to Pattern Recognition and Machine Learning
Temple University
College of Engineering
1947 North 12th Street
Philadelphia, Pennsylvania 19122

April 20th, 2023

prepared by:

Gavin Koma
Email: gavintkoma@temple.edu

A. INTRODUCTION

This homework is a bit different as we have not really spent time discussing transformer implementation, thus, I want to spend less of this time focusing on the code and instead explain the individual parts that make up this implementation for my own understanding. The code does work in its entirety, and I have attached it on the last few pages of this assignment. Since the code is explained within the tutorial itself, I will limit my own screenshots and include portions that I feel are imperative to my own understanding of a transformer implementation.

The tutorial itself was easy to implement and I did not have many problems albeit at the very start. To begin this tutorial, one needs to have the proper versions of modules installed. I regularly try to maintain my modules at their most contemporary version however, this code utilizes older modules that I had newer versions of. As a result, many of my dependencies would clash and I was met with a variety of errors. I am a bit grateful for this as it forced me to delve into a subject that I have not previously studied: virtual environments.

I first spend some time reading through the virtualenvironments documentation and created my own venv for this assignment that would work well given the code. To the right, one can see a few of the required modules needed and their respective versions. I previously wondered why and how engineers are able to run so many different applications without maintaining dependencies between modules, but this was a huge eureka moment for me.

Virtual environments are incredibly important because they allow developers to create an isolated environment that have their own individual dependencies and configurations. This allows us to manage a project's dependencies separately so distinct projects can have distinct dependencies with different package versions. Another key point regarding virtual environments is that they help to ensure reproducibility. In other words, using a VE will allow code to run in the same way on different machines while utilizing the same isolated environment for development.

As someone that doesn't have a computer science background, learning about VEs is essential for my own growth as someone that wants to work in the data science field.

B. MODEL ARCHITECTURE

The transformer algorithm is one type of deep learning model that has most notably been used in the field of natural language processing (NLP) and was first introduced by Vaswani in 2017. The algorithm itself is based on the concept of self-attention and thus allows the model to focus on different parts of the input

```
Last login: Thu Apr 20 09:55:00 on ttys000
[(venv) gavinkoma@Gavins-MBP-2 ~ % pip list
```

Package	Version
altair	4.2.2
attrs	23.1.0
blis	0.7.9
catalogue	2.0.8
certifi	2022.12.7
charset-normalizer	3.1.0
click	8.1.3
confection	0.0.4
cymem	2.0.7
entrypoints	0.4
filelock	3.12.0
GPUtil	1.4.0
idna	3.4
Jinja2	3.1.2
jsonschema	4.17.3
langcodes	3.3.0
MarkupSafe	2.1.2
mpmath	1.3.0
murmurhash	1.0.9
networkx	3.1
numpy	1.24.2
packaging	23.1
pandas	2.0.0
pathy	0.10.1
pip	23.1
prashed	3.0.8
pydantic	1.10.7
pyrsistent	0.19.3
python-dateutil	2.8.2
pytz	2023.3
requests	2.28.2
setuptools	67.6.1
six	1.16.0
smart-open	6.3.0
spacy	3.5.2
spacy-legacy	3.0.12
spacy-loggers	1.0.4
srsly	2.4.6
sympy	1.11.1
thinc	8.1.9
toolz	0.12.0
torch	2.0.0

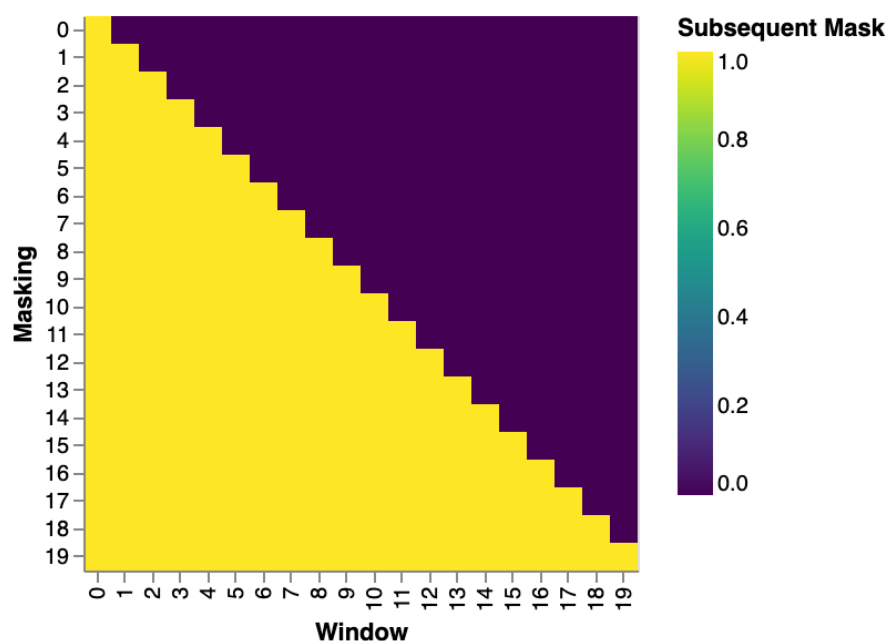
sequence when processing each output. The algorithm uses a series of self-attention layers (which are shown later in the code) and feedforward layers (also seen later) to process the given input data. When the model approaches a self-attention layer, the model functions to learn which parts of the input sequence to focus on and how much attention to give each part.

The first portion of this paper specifies that most competitive neural sequence transduction models rely on an encoder-decoder structure. Within said structure, the encoder will map an input sequence of symbol representations to another sequence of continuous representations. Given these values, the decoder then functions to generate an output sequence of symbols one variable at a time. One unique and important aspect in this architecture is that the model itself is auto-regressive. In other words, as the model progresses, it also consumes the previously generated output and utilizes it as an additional source of input when generating the subsequent output.

The encoder layer itself takes in the input sequence and applies self-attention to it and actively attempts to learn which parts of the sequence to focus on and how much attention to give to each part. The self-attention mechanism that we implement will allow the encoder to focus on dependencies between different parts of the input sequence. After this is completed, the encoder applies a feedforward neural network to the processed sequence and produces the previously mentioned encoded representation of the data.

The decoder layer then takes the encoded representation and applies another round of self-attention to it as well as attention to the input sequence. The decoder will use the attention mechanism to perform a similar action as the encoder: to focus on different parts of the encoded representation and input sequence to generate an output sequence. Just like the encoder, the decoder will then pass this to a feedforward neural network to process the data and output the next token in the sequence.

At this point of the tutorial, we output the same graph as shown in the post. This image shows an attention mask which shows the position that each tgt word is allowed to look at. It is used to control the attention mechanism found within the transformer algorithm and its purpose is to mask out certain elements within the input sequence so that the attention mechanism does not focus on them.



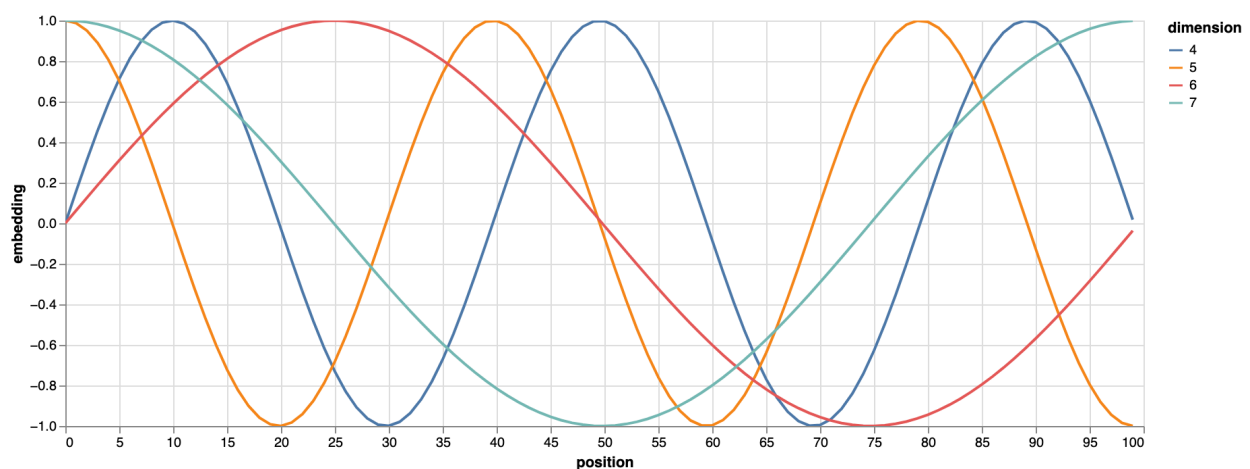
The attention function of our transformer is used to compute a weighted sum of the values. These weights are computed based on the similarity between the query and the keys themselves. Our given input consists of queries, keys, and their corresponding dimension. Our function computes the dot products of the query with all keys divided by the root of keys dimension and followed by a softmax function. The softmax function works to normalize the scores and ensures that they sum to a total of 1. These scores are then used as weights which will compute the previously mentioned weighted sum. This function is able to be visualized as a mechanism that allows the model to focus on different parts of the input sequence when processing each output. This function is used in both the encoder and decoder layers to best provide an output based on the input sequence and the encoded representation.

Our transformer specifically utilizes multi-head attention. This mechanism ensures that our model can attend to multiple positions within the input sequence at the same time. The query, keys, and value vectors are transformed to create multiple “heads” for attention. The queries come from the previous decoder layer and the memory keys and corresponding value vectors come from the output of the encoder. Doing so, allows every position in the decoder to attend over all positions in the input sequence. The encoder’s attention layers have a similar element. However, in these self-attention layers, the keys, values, and queries all come from the same place: the output of the previous layer in the encoder.

Our transformer then moves onto a position-wise feedforward neural network which can be found fully connected in each encoder and decoder. The feedforward neural network contains two linear transformations with a ReLU activation in between them. I will not go in-depth regarding this portion as we have worked with these functions in previous homework assignment.

Our model does not contain recurrence or convolution and as a byproduct, we must include some information regarding the relative position of tokens in the sequence. We do so by including “positional encodings” in the input embeddings at the bottoms of the encoder and decoder stacks. Doing so, allows us to provide a method for the model to encode the order of the input sequence.

The process of positional encoding involves adding a set of fixed-length vectors to the input embeddings of each word in the sequence. These vectors are calculated using a sinusoidal method and the dimension of the corresponding embedding space. These sinusoidal functions ensure that the positional encoding vectors are unique for each position in the sequence and that they also have a smooth and continuous pattern that permits the model to interpolate between given positions. At this point in the tutorial, we are given a chance to output the following graph that shows varying sinusoidal functions of an embedding given their position:



The full model of our transformer is defined below along with a test of inference to try to generate a prediction of our model:

```

def make_model(
    src_vocab, tgt_vocab, N=6, d_model=512, d_ff=2048, h=8, dropout=0.1
):
    """Helper: Construct a model from hyperparameters."""
    c = copy.deepcopy
    attn = MultiHeadedAttention(h, d_model)
    ff = PositionwiseFeedForward(d_model, d_ff, dropout)
    position = PositionalEncoding(d_model, dropout)
    model = EncoderDecoder(
        Encoder(EncoderLayer(d_model, c(attn), c(ff), dropout), N),
        Decoder(DecoderLayer(d_model, c(attn), c(attn), c(ff), dropout), N),
        nn.Sequential(Embeddings(d_model, src_vocab), c(position)),
        nn.Sequential(Embeddings(d_model, tgt_vocab), c(position)),
        Generator(d_model, tgt_vocab),
    )

    # This was important from their code.
    # Initialize parameters with Glorot / fan_avg.
    for p in model.parameters():
        if p.dim() > 1:
            nn.init.xavier_uniform_(p)
    return model

```

```

def inference_test():
    test_model = make_model(11, 11, 2)
    test_model.eval()
    src = torch.LongTensor([[1, 2, 3, 4, 5, 6, 7, 8, 9, 10]])
    src_mask = torch.ones(1, 1, 10)
    memory = test_model.encode(src, src_mask)
    ys = torch.zeros(1, 1).type_as(src)
    for i in range(9):
        out = test_model.decode(
            memory, src_mask, ys, subsequent_mask(ys.size(1)).type_as(src.data)
        )
        prob = test_model.generator(out[:, -1])
        _, next_word = torch.max(prob, dim=1)
        next_word = next_word.data[0]
        ys = torch.cat(
            [ys, torch.empty(1, 1).type_as(src.data).fill_(next_word)], dim=1
        )
    print("Example Untrained Model Prediction:", ys)

def run_tests():
    for _ in range(10):
        inference_test()

```

C. MODEL TRAINING

This portion of the assignment was by far the most confusing for me to understand. I will try my best to explain the methods introduced in this article and to explain their implementation. The model we implement in this tutorial contain the following essential steps to properly train:

1. Batches and Masking Creation
2. Training Loop
3. Providing Training Data and Batching
4. Hardware for Training
5. Optimizer
6. Regularization

Batches and masking are considered fairly common tools that are used for the implementation of transformer models. Batches are used to improve the efficiency of our training process. This occurs by processing multiple input sequences simultaneously. During this portion, several input sequences are grouped together into one single batch and are processed at the same time instead of processing each batch individually. Doing so, allows us to reduce the time required to not only train the model but also allows us to utilize parallel computing architectures. Batches are typically processed in a parallel fashion in which the input sequences are stacked together into a singular matrix. Unfortunately, input sequences within a batch may have different lengths and it is common for researchers to include padding when working with the shorter sequences to ensure that they are the same length as the longest sequence.

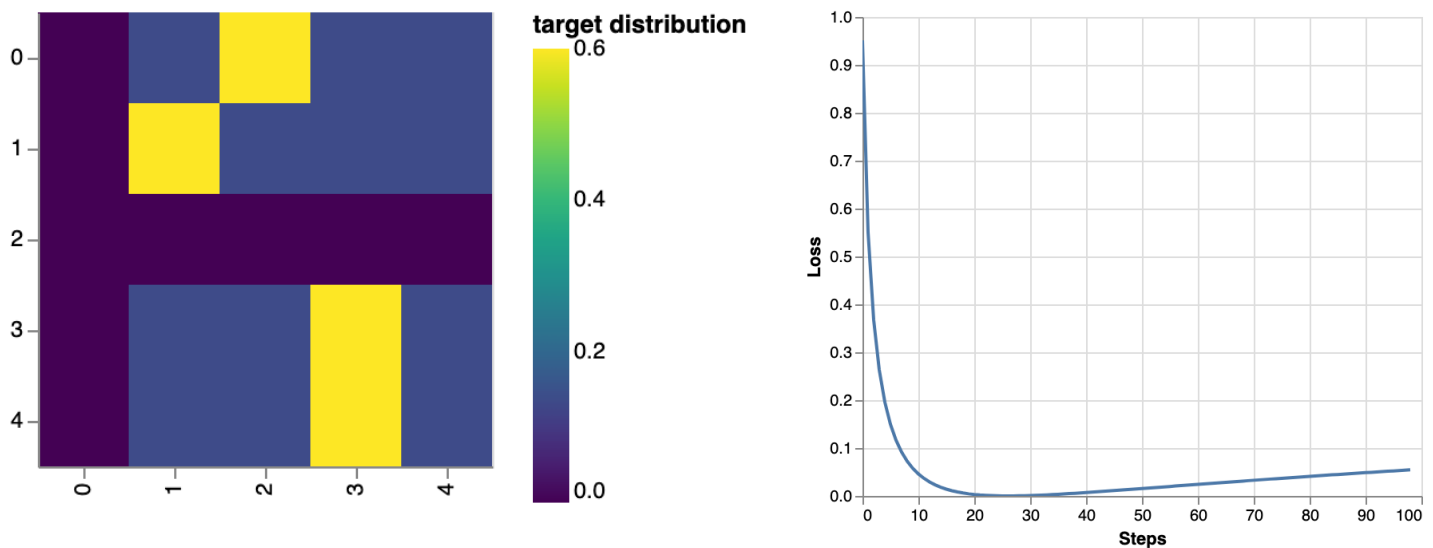
Masking is commonly used to selectively control which elements of the input sequence are attended to during the processing phase. There are two common types of masks that are used: padding masks and sequence masks. I will only focus on sequence masks as this is the type of mask that is utilized in our model. Sequence masks are binary matrices that are used to mask out future elements found in the input sequence. During the generation of an output sequence, it is important that the model only attends to the previous elements in the input sequence. This matrix is the same shape as the input or output sequence and each element within it is either a 1 or a 0. A value of 1 means that the corresponding element in the sequence *should* be attended to while a value of 0 allows the model to know that the element should be masked out.

The training loop included in the program is fairly self-explanatory and will not be further explained here. The code itself is straight forward and similar methods have been implemented to train our models in the past in a loop-like model. At this point I did need to reassess my method of running this algorithm. The researchers utilized 8 NVIDIA P100 GPUs. Unfortunately, I do not have 8 GPUs casually laying around my PhD-student-only-household. I first attempted to run this code on my lab's remote GPU but training was taking well over an hour. I got a bit lucky here and have a close friend that pays for Google Collab's Premium Service which supplies great computer resources for heavy computational loads.

The final portion of training was to utilize the Adam optimizer. Adam, short for Adaptive Moment Estimation, is designed to update the weights of a neural network in an efficient and *adaptive* manner. It utilizes two main strategies to complete this task: adaptive learning rates and momentum-based updates. The adaptive learning rate strategy adjusts the learning rate for each weight in the network based on the magnitude of its gradients. Doing so helps to prevent the network from becoming stuck in steep regions when the calculated gradients are large. The momentum-based update strategy helps to accelerate the optimization process by using a moving average of the gradients to update the computed weights. This helps the optimization process to overcome small, local minima and to continue moving towards the global minimum. Adam combines these two strategies by using the first and second moments of the gradients to compute the adaptive learning rates and then to compute the moving average of the gradients to update the weights.

We conclude this training portion by performing label smoothing regularization using KL divergence loss. Label smoothing is a regularization technique that prevents our model from overfitting to the training data by artificially smoothing our target labels. Generally, label smoothing involves replacing the hard target labels with a smoothed distribution across all classes. The idea behind this concept is to introduce a small amount of noise into the training labels in an attempt to force the model to learn more robust features. In our model, we utilize the KL divergence loss to penalize the difference between the predicted probability distribution and the smoothed target distribution. This loss is then used as the objective function for training by the model. By minimizing this loss, we allow our model to complete our previously stated goal: to learn more robust features.

At this point in the tutorial, our code outputs the following example of label smoothing and loss graph:



D. SYNTHETIC EXAMPLE

The next portion of this assignment has us focus on a simple copy-task given a random set of input symbols from some small vocabulary. The goal of this portion is to generate these same symbols back to us. The code implemented generates synthetic data of some size. We introduce a simple loss compute and train functions and we predict our training using greedy decoding. Greedy decoding is a simple decoding strategy that is commonly used in sequence-to-sequence models and the strategy involves selecting the output token with the highest probability at each decoding step and appending it to the output sequence.

The code itself worked with no problems on my GPU but the article does not spend much time discussing this portion. Therefore, I also will not spend much time regarding this synthetic example but will instead focus more on the real-world example that is performed next in much greater detail. The code that was ran on my laptop can be found at the end of this submission in .pdf format with all output results maintained and included.

E. REAL WORLD EXAMPLE

To conclude this assignment, we utilize our model on a real-world example. In this case, we are undertaking the German-English translation task. We first load the data, which is self-explanatory as we have performed this task in nearly every assignment since the beginning of the semester. Once our data is loaded, we see that we have the following vocabulary sizes for German and for English respectively: 59981 and 36745.

As previously discussed, it is incredibly important that we include batching to maximize the speed of our algorithm. We do want to evenly divide our batches and we want to do so with minimal padding. The code provided does this for us. It is important to have minimal padding in a transformer as it can lead to an increased efficiency and a faster training time. By minimizing these two factors, we are able to decrease the amount of computation required to process the sequences and it allows the model to focus more on the actual data and less on the padding. This also will hopefully lead to better accuracy and an overall faster convergence.

Training the system took a notable amount of time, but included below is a small snippet of the output to show the general trend of loss from our model:

```

| ID | GPU | MEM |
-----
| 0 | 79% | 26% |
[GPU0] Epoch 0 Validation ====
(tensor(3.8904, device='cuda:0'), <__main__.TrainState object at 0x7fd3af1372b0>)
[GPU0] Epoch 1 Training ====
Epoch Step:      1 | Accumulation Step:  1 | Loss:    3.83 | Tokens / Sec: 1955.9 | Learning Rate: 2.4e-04
Epoch Step:     41 | Accumulation Step:  5 | Loss:    4.01 | Tokens / Sec: 1680.1 | Learning Rate: 2.6e-04
Epoch Step:     81 | Accumulation Step:  9 | Loss:    3.87 | Tokens / Sec: 1632.0 | Learning Rate: 2.7e-04
Epoch Step:    121 | Accumulation Step: 13 | Loss:    3.60 | Tokens / Sec: 1658.1 | Learning Rate: 2.8e-04
Epoch Step:    161 | Accumulation Step: 17 | Loss:    3.63 | Tokens / Sec: 1674.1 | Learning Rate: 2.9e-04
Epoch Step:    201 | Accumulation Step: 21 | Loss:    3.57 | Tokens / Sec: 1663.8 | Learning Rate: 3.0e-04
Epoch Step:    241 | Accumulation Step: 25 | Loss:    3.61 | Tokens / Sec: 1659.0 | Learning Rate: 3.1e-04
Epoch Step:    281 | Accumulation Step: 29 | Loss:    3.54 | Tokens / Sec: 1676.2 | Learning Rate: 3.2e-04
Epoch Step:    321 | Accumulation Step: 33 | Loss:    3.57 | Tokens / Sec: 1670.8 | Learning Rate: 3.3e-04
Epoch Step:    361 | Accumulation Step: 37 | Loss:    3.54 | Tokens / Sec: 1673.2 | Learning Rate: 3.4e-04
Epoch Step:    401 | Accumulation Step: 41 | Loss:    3.39 | Tokens / Sec: 1678.1 | Learning Rate: 3.5e-04
Epoch Step:    441 | Accumulation Step: 45 | Loss:    3.26 | Tokens / Sec: 1654.8 | Learning Rate: 3.6e-04
Epoch Step:    481 | Accumulation Step: 49 | Loss:    3.37 | Tokens / Sec: 1635.7 | Learning Rate: 3.7e-04
Epoch Step:    521 | Accumulation Step: 53 | Loss:    3.28 | Tokens / Sec: 1626.0 | Learning Rate: 3.8e-04
Epoch Step:    561 | Accumulation Step: 57 | Loss:    3.37 | Tokens / Sec: 1647.0 | Learning Rate: 4.0e-04
Epoch Step:    601 | Accumulation Step: 61 | Loss:    3.22 | Tokens / Sec: 1673.1 | Learning Rate: 4.1e-04
Epoch Step:    641 | Accumulation Step: 65 | Loss:    3.04 | Tokens / Sec: 1681.3 | Learning Rate: 4.2e-04
Epoch Step:    681 | Accumulation Step: 69 | Loss:    2.93 | Tokens / Sec: 1651.9 | Learning Rate: 4.3e-04
Epoch Step:    721 | Accumulation Step: 73 | Loss:    2.93 | Tokens / Sec: 1688.6 | Learning Rate: 4.4e-04
Epoch Step:    761 | Accumulation Step: 77 | Loss:    2.89 | Tokens / Sec: 1671.6 | Learning Rate: 4.5e-04
Epoch Step:    801 | Accumulation Step: 81 | Loss:    2.99 | Tokens / Sec: 1664.0 | Learning Rate: 4.6e-04
Epoch Step:    841 | Accumulation Step: 85 | Loss:    3.11 | Tokens / Sec: 1658.2 | Learning Rate: 4.7e-04
Epoch Step:    881 | Accumulation Step: 89 | Loss:    2.73 | Tokens / Sec: 1670.7 | Learning Rate: 4.8e-04

```

F. RESULTS

The results of this model can be found below. We output three separate examples prioritizing three different methods: encoder self-attention, decoder self-attention, and decoder source-attention. Doing so, allows us to compare the three different attention mechanisms. Encoder self-attention is used in the encoder to calculate the context vector for each token in the input sequence. In this situation, the input sequence is passed to the model and each token attends to all other tokens in the same input sequence in order to calculate its own context vector. In this situation, the context vector only depends on the input sequence and not on the output sequence. Decoder self-attention is used in the decoder to calculate the context vector for each token in the output sequence. Here, the model generates the output sequence one token at a time and each token will attend to all the previous tokens that were generated by the decoder in order to calculate its context vector. In this situation, the context vector depends only on the previously generated tokens and not on the input sequence. Finally, there is the decoder source-attention. This method is used in the decoder to attend to the encoder's output sequence. In this situation, the decoder will attend to all the tokens in the encoder's output sequence to calculate the context vector for the current token being generated. Here, the decoder will use information from the input sequence to generate the output sequence.

Encoder Self-Attention:

Example 0 =====

```
Source Text (Input)      : <s> Zwei Männer mit <unk> sprechen miteinander auf einem Freiluftmarkt . </s>
Target Text (Ground Truth) : <s> Two men in <unk> 's have a discussion in an outdoor market . </s>
Model Output             : <s> Two men are talking to each other in an open air market . </s>
```

Decoder Self-Attention:

Example 0 =====

```
Source Text (Input)      : <s> Ein Mann mit Rastalocken <unk> sein Ohr , um ein <unk> hören zu können . </s>
Target Text (Ground Truth) : <s> A man with dreadlocks is plugging his ear to hear a phone call . </s>
Model Output             : <s> A man with a <unk> hat is seen his ear to listen to be a music stand . </s>
```

Decoder Source-Attention:

Example 0 =====

```
Source Text (Input)      : <s> Ein Mann isst in einem Restaurant zu Mittag . </s>
Target Text (Ground Truth) : <s> A man in a restaurant having lunch . </s>
Model Output             : <s> A man eats lunch to a restaurant . </s>
```

I have initially not included the heat maps for these outputs as they can be generally confusing to read/interpret. I have included them, instead, at the very end of this document where they are found in the output of the code implementation. From the above results, it seems like the encoder self-attention is providing the best output for German-English translation. It appears to be the easiest to understand of the three. There are simple mistakes and the other two methods are still legible and their intention is clear, the encoder self-attention is the most accurate.

G. CONCLUSION

This was, difficult, to say the least. I didn't realize how in-depth some of these algorithms could be or how complicated they can get. In many classes, the algorithms implemented don't even approach half of the intense work that would be required to design something like this from scratch. As always though, I find myself learning something from these assignments. The code itself would have been miserable to comment out line-by-line but I hope my overall explanations for each section were satisfactory.

I wrote this document primarily from the intent of learning whatever I could and teach myself more about concepts that I don't know in-depth. I feel like this assignment facilitated me learning a lot regarding encoder/decoders, the importance of model training methods, and as I stated in the beginning of this document, virtual environments.

Attached on the next few pages is the entire print out of my completed and functioning implementation of this code.