

# Homework 4 Part 2

## Attention-based End-to-End Speech-to-Text Deep Neural Network

11-785: INTRODUCTION TO DEEP LEARNING (SPRING 2024)

OUT: **March 31, 2024, at 11:59 PM ET**

Preliminary Submission (and Canvas Quiz): **April 10, 2024, at 11:59 PM ET**

DUE: **April 26, 2024, at 11:59 PM ET**

### Start Here

- **Collaboration policy:**

- You are expected to comply with the University Policy on Academic Integrity and Plagiarism.
- You are allowed to talk with / work with other students on homework assignments
- You can share ideas but not code, you must submit your own code. All submitted code will be compared against all code submitted this semester and in previous semesters using MOSS.

- **Submission:**

- You need to go through the writeup and finish the Canvas quiz by the early submission deadline.
- You need to implement what is mentioned in the write-up (although you are free and encouraged to explore further) and submit your results to Kaggle . We will share a Google form or an Autolab link after the Kaggle competition ends for you to submit your code.

- **TL; DR:**

- In this homework you will work on a sequence-to-sequence conversion problem, in which you must train models to transcribe speech recordings into word sequences, spelled out alphabetically.
- As in previous HWP2s, this homework too is conducted as a Kaggle competition. You can access the kaggle for the homework at this link.
- You can download the training and test data directly from Kaggle using their API or this link. Your download will include the following:
  - \* There are two sets of training data, train-clean-100 and train-clean-50. They both include data pairs comprising of sequences of audio feature vectors and their transcriptions, in mfcc and transcripts folders respectively.
  - \* The validation data will be similarly found in the dev-clean folder.
  - \* The test data will be found in the test-clean folder. It will only contain an mfcc folder.
- You must use the train data to train the model. Data in the dev-clean folder are to be used only for validation and not for training.
- We recommend you using the train-clean-50 if you just want to try out things before using train-clean-100.
- You must transcribe the utterances in the test-clean folder and write them out into submission.csv following the format already specified in the file, and upload your results to Kaggle.
- You will be graded by your performance. Additional details will be posted on piazza.

## Homework objectives

If you complete this homework successfully, you would ideally have learned:

- To implement a Transformer architecture to solve a sequence-to-sequence problem.
  - How to setup an encoder with mutlihead-attention and self-attention.
  - How to setup a decoder with cross-attention, multihead-attention and self-attention.
  - How to use attention mechanism with masking.
  - How to use positional encoding.
- You would have learned how to address some of the implementation details
  - How to setup a teacher forcing for the decoder training.
  - How to use search techniques, such as greedy search, for inference in the transformer decoder to generate the output sequence.
  - How to pad-pack the variable length data
- To explore architectures and hyperparameters for the optimal solution
  - To tune parameters and hyperparameters that affect your solution.
  - To effectively explore the search space and strategically finding the best solution.
- As a side benefit you would also have learned how to construct a speech recognition system using the transformer encoder-decoder architecture.

## Important: A note on presentation style and pseudocode

The following writeup is intended to be reasonably detailed and explains several concepts through pseudocode. However there is a caveat; the pseudocode is intended to be *illustrative*, not *exemplary*. It conveys the concepts, but you cannot simply convert it directly to python code. For instance, most of the provided pseudocode is in the form of functions, whereas your python code would use classes (in fact, we've written the pseudocode as functions with the express reason that you *cannot* simply translate it to your code. Function-based DL code will be terribly inefficient. You *must* use classes). Please write reusable code wherever possible, and avoid redundant calculations to make your code more efficient.

Nonetheless, we hope that when you read the entire write-up, you will get a reasonable understanding of how to implement (at least the baseline architecture of) the homework.

# Contents

<b>1</b>	<b>Introduction</b>	<b>6</b>
1.1	Overview . . . . .	6
1.2	Problem specifics . . . . .	8
1.3	Getting Started Early . . . . .	10
1.4	Starter Notebook . . . . .	10
1.5	Resources . . . . .	11
<b>2</b>	<b>Notebook Walkthrough</b>	<b>12</b>
<b>3</b>	<b>Dataset</b>	<b>14</b>
3.1	Dataloader . . . . .	14
<b>4</b>	<b>Speech Transformer</b>	<b>15</b>
4.1	Encoder . . . . .	15
4.1.1	CNN-LSTM as Encoder and as Input Embedding . . . . .	16
4.1.2	Encoder Input Encoding . . . . .	17
4.1.3	Encoder Layers . . . . .	17
4.1.4	Transformer Encoder . . . . .	19
4.2	Decoder . . . . .	21
4.2.1	Decoder inputs . . . . .	21
4.2.2	Output Encoding . . . . .	21
4.2.3	Decoder Layers . . . . .	22
4.2.4	Transformer Decoder . . . . .	23
4.3	Combining it All . . . . .	24
<b>5</b>	<b>Training</b>	<b>25</b>
5.0.1	Loss Function . . . . .	25
5.0.2	Overall Training Procedure . . . . .	25
5.1	Inference . . . . .	26
5.1.1	Greedy Search . . . . .	26
<b>6</b>	<b>Common Errors and Fixes</b>	<b>28</b>
6.1	Gradient Clipping . . . . .	28
6.2	Debugging NaN Loss . . . . .	28
6.3	Out Of Memory (OOM) . . . . .	28
6.4	Dataloader and Dataset . . . . .	28
6.5	Concatenating the Context Vector with the Embeddings . . . . .	29
6.6	Training Duration and Monitoring . . . . .	29
<b>7</b>	<b>Conclusion</b>	<b>30</b>
	<b>Appendices</b>	<b>31</b>
<b>A</b>	<b>Character Based vs Word Based</b>	<b>31</b>
<b>B</b>	<b>Transcript Processing</b>	<b>31</b>
<b>C</b>	<b>Single Head Attention (Minibatch)</b>	<b>31</b>
<b>D</b>	<b>Multi-Head Attention</b>	<b>32</b>
<b>5</b>	<b>Levenshtein Distance</b>	<b>33</b>
<b>6</b>	<b>Optimizer and Learning Rate</b>	<b>33</b>

<b>7</b>	<b>Transformers and Pretraining</b>	<b>33</b>
<b>8</b>	<b>An important note on Attention</b>	<b>34</b>
8.1	An attention head . . . . .	34
8.2	Multi-head attention . . . . .	34
8.3	Masking . . . . .	35
8.4	Position Encoding . . . . .	35
8.4.1	Positional Encoding Layer in Transformers . . . . .	35

# Checklist

Here is a checklist page that you can use to keep track of your progress as you go through the write-up and implement the corresponding sections in your starter notebook. As you complete your unit tests for each component to verify your progress, you can check the corresponding boxes aligned with each section and go through this writeup and starter notebook simultaneously step by step.

1. Read Introduction and Notebook Walkthrough.
2. Read the Writeup, and finish the Canvas Quiz.
3. Dataset and Dataloader.  
[Dataset Description.](#)
4. Encoder  
[Encoder and Transformer Encoder Description.](#)  
[Implement CNN-LSTM Encoder \(Part I\).](#)  
[Implement Transformer Encoder \(Part II\).](#)
5. Decoder  
[Speech Transformer Decoder Description.](#)  
[Implement Transformer Decoder \(Decoder Layer\) \(Part I and II\).](#)  
[Implement Decoder \(combine Decoder modules\) \(Part I and II\).](#)
6. Training  
[Training Description.](#)

# 1 Introduction

**Key new concepts:** Sequence-to-sequence conversion, Encoder-decoder architectures, Attention, Transformer.

**Core Ideas:** Positional Encoding, Multihead-attention, Cross-attention, Self-attention, Masking.

**Mandatory implementation:** Encoder-decoder architecture, Attention. Your solution *must* implement an encoder-decoder architecture *with* attention. While it may be possible to solve the homework using other architectures, you will not get marks for implementing those.

**Restrictions:** You may not use any data besides that provided as part of this homework. You may not use the validation data for training. You cannot use an implementation from PyTorch or any other third-party package for implementing any type of Attention (and optionally beam search).

## 1.1 Overview

In this homework you will learn to build neural networks for sequence-to-sequence conversion or transduction.

”Sequence-to-sequence” conversion refers to problems where a sequence of inputs goes into the system, which emits a sequence of outputs in response. The need for such conversion arises in many problems, e.g.

- Speech recognition: A sequence of speech feature vectors is input. The output is a sequence of characters writing out what was spoken.
- Machine translation: A sequence of words in one language is input. The output is a sequence of words in a different language.
- Dialog systems: A user’s input goes in. The output is the system’s response.

For this homework, we will work on the speech recognition problem. Unlike HW1P2 and HW3P2, where the task was to predict *phonemes* in the speech, in this homework, we will learn to directly spell the spoken sentence out in the English alphabet.

A key characteristic of sequence-to-sequence problems is that there may be no obvious correspondence between the input and the output sequences. For instance, in a dialog system, a user input “my screen is blank” may elicit the output “check the power switch”. In a speech recognition system the input may be the (feature vector sequence for the) spoken recording for the phrase “know how”. The system output must be the character sequence “k”, “n”, “o”, “w”, “ ”, “h”, “o”, “w” (note the blank space “ ” between know and how – the output is supposed to be readable, and must include blank spaces as characters). There is no obvious audio corresponding to the silent characters “k” and “w” in *know*, or the blank space character between *know* and *how*, nevertheless these characters must be output.

As a consequence of the absence of correspondence between the input and output sequences, the usual “vertical” architectures that we have used so far in our prior homeworks (Figure 1a), where the input is sequentially processed by layers of neurons until the output is finally computed, can no longer be used.

Instead, we must use a two-component model, comprising two *separate* network blocks, one to process the input, and another to compute the output. Such architectures are called *encoder-decoder* architectures. The encoder processes the input sequence to compute feature representations (embeddings) of the input sequence. The decoder subsequently uses the input representations computed by the encoder to *sequentially* compute the output *de novo*, *i.e.* from scratch (Figure 1b).

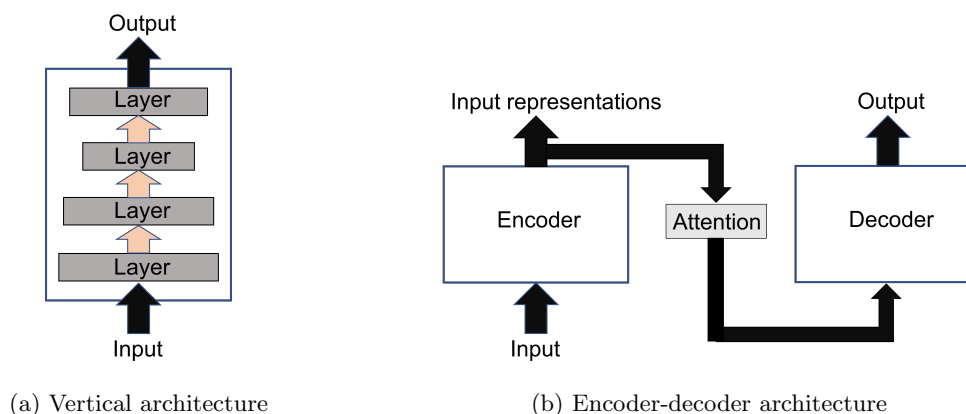


Figure 1: (a) Standard “vertical” architectures we have seen in previous homeworks. The input is sequentially passed through several layers until the output is computed from the final layer. *This kind of architecture is not useful for sequence-to-sequence problems with no input-output correspondence, such as the problem we deal with in this homework.* (b) General encoder-decoder architecture, comprising two separate modules, an encoder to process the input and a decoder to generate the output. The two are linked through an attention module.

In this homework, we will build a **transformer architecture**, which is an example of encoder-decoder architectures. The transformer encoder takes an input sequence and produces feature representations, while the decoder uses these representations to generate the target sequence. Unlike the LSTM and RNN models we saw in HW3P2, which process sequences step-by-step, the transformer can process the entire sequence at once. This allows for much faster training by taking advantage of parallel computing. **But how does the transformer achieve parallel computation?** It does so through the attention mechanism, a technique that allows the model to focus on specific parts of the input when producing the output. This mechanism dynamically weighs the importance of different input elements to capture their relationships. This features in two places:

- i. **Self-attention** – the attention that the encoder/decoder pay to themselves: In computing their respective outputs, both the encoder (which derives feature representations from the input) and the decoder (which predicts outputs) must also pay attention to their own internal context, which in the case of the encoder would be adjacent inputs, and for the decoder would be previous outputs. This is *self-attention* – the attention a module pays to its own internal variables.
- ii. **Cross-attention** – the attention the decoder pays to the encoder: Although the decoder’s output may not have a direct correspondence with the input, each individual output symbol (character) nonetheless relates more to some parts of the input than others, *e.g.* the blank character “ ” in our above speech recognition example corresponds primarily to the audio at the boundary between the words “know” and “how” in the recording. When outputting the “ ”, the decoder must somehow know to “focus” on, or “pay attention” to this region of the input. In practice, this is implemented as an “attention” that the decoder pays to feature representations of the inputs derived by the encoder. Since this involves attention paid by one module, the decoder, to representations derived by another module (the encoder) this is referred to as *cross-attention*.

The attention mechanism allows the model to selectively focus on relevant parts of the input and ignore irrelevant information. This is particularly useful in tasks like machine translation, where the relationship between input and output elements is not always linear, and different parts of the input may be relevant for different parts of the output.

**But how does the transformer know the order of the sequence?** The transformer model

understands the order of the sequence through a technique called **positional encoding**. Since the transformer processes the entire sequence at once and doesn't have any inherent sense of order, positional encoding is added to give the model information about the position of each element in the sequence. Positional encoding is typically implemented by adding a vector to each input embedding, where the vector is generated based on the position of the element in the sequence. These vectors follow a specific pattern that helps the model distinguish between different positions. For more details refer to section 8.4

In this homework, you will learn how to build an encoder to effectively extract features from a speech signal, how to construct a Transformer-based decoder that can generate the transcription of the audio, and how to implement an attention mechanism between the Transformer decoder and the encoder.

Although the specifics of the homework will be targeted towards speech recognition, please note that the general framework (with appropriate modification of the encoder and decoder architectures) should also apply to other types of sequence-to-sequence problems, or even to image or video captioning (you would only need to modify the encoder to derive features from images or video respectively).

## 1.2 Problem specifics

In this problem of speech recognition you will learn how to convert a sequence of feature vectors of speech to a sequence of characters that spells out its transcription. To achieve this, you will implement a Speech Transformer model which is effective in capturing complex patterns in audio data using its self-attention, multi-head attention and cross-attention mechanisms. You will be implementing the two major parts of a transformer architecture which are the encoder and decoder

- **Transformer Encoder.** In this part you will implement a module that is capable of extracting high level feature vectors from your inputs (MFCCs). This involves putting together various pieces of the encoder architecture which are convolutional modules, positional encoding, layer normalization, feed forward neural network, self-attention and multi-head attention.
- **LSTM Encoder:** In addressing the substantial data and computational demands associated with training Transformer models, our approach integrates Long Short-Term Memory (LSTM) networks and Transformer encoders in a novel hybrid architecture. Initially, LSTMs are utilized to encode sequential data into compact embeddings, capitalizing on their proficiency in capturing temporal dependencies with lower resource consumption. These LSTM-generated embeddings are subsequently fed into Transformer encoders, harnessing the Transformer's superior ability to model complex relationships in data. This strategic combination aims to marry the LSTM's resource efficiency with the Transformer's advanced processing prowess, offering a more accessible and scalable solution for handling sophisticated sequence modeling tasks in scenarios where resources are limited.
- **Transformer Decoder:** operates by taking the encoder's output and any previously generated outputs as inputs, processing them through multiple layers of self-attention and encoder-decoder attention mechanisms. This structure allows the decoder to focus on different parts of the speech input at different times, effectively capturing the context and dependencies within the speech. Similar to the encoder it is a combination of the positional encoding, layer normalization, self-attention, multi-head attention, cross-attention, feed forward neural networks and linear classifiers. You will be required to use a pre-trained neural network as encoder and also as input embedding.
- **Attention.** Both the encoder and decoder are expected to employ self-attention to compute their respective outputs. You will be required to implement masked multi-head attention and cross-attention between the encoder and decoder
- **Speech Transformer:** You will be expected to combine the encoder and decoder to build a speech-transformer capable of performing speech recognition task.



- **Training:** You will be training the transformer you built above using the training data provided. Since this model has no recurrence, you will be using a cross-entropy loss function rather than CTC loss as in HW3P2. You will use teacher forcing during training which involves always using the true previous output as input for the next step, ensuring direct learning.
- **Inference:** In previous assignments your forward method will give the full sequence outputs of your inputs. However here, during inference, the decoder behaves as an auto-regressive network meaning the predicted sequence so far serves as input to predict the next token and this behavior is different from the forward during training. You are tasked in implementing a greedy search algorithm for the decoding the decoder outputs.

In this homework, we will be implementing the speech transformer architecture based on Dong et al. [2018], Vaswani et al. [2017] (which we highly recommend that you read), and describe the key components of it later in Section 4. We will provide you with the code for some components, and you will implement the others yourself.

- **Components provided:**

- i. Dataset class and data loader.
- ii. Three masking functions to create masks for encoder and decoder training.
- iii. Scaled Dot-Product Module for computing attention scores between query, key, and value vectors.
- iv. Multi-Head Attention Module to create multiple sets of query, key, and value weight matrices based on the number of heads.
- v. Position Encoding Module to add positional information to the input embeddings.
- vi. Feed Forward Module, which consists of fully connected layers to increase the model's capacity.
- vii. Pre-trained CNN-LSTM encoder for you to get started and train your decoder. You might also use your HW3P2 as an encoder if you got good results in the assignment. Note this is only for making the decoder learn, performance will be improved when we leverage incremental training by adding transformer-encoder layers..
- viii. Speech Transformer Module that combines the encoder and decoder modules to create a full speech transformer architecture.
- ix. We will also provide the training, validation and testing code.

- **Your task is to implement these components:**

- i. Encoder Layer Module to use in the encoder module. Here you will use the Multi-Head Attention and Feed Forward Modules that we provided, along with a LayerNorm from PyTorch.
- ii. Encoder Module that stacks multiple encoder layers. Here you will use positional encoding and two masking functions that were provided to you.
- iii. Decoder Layer Module to use in the decoder module. Here you will use the Multi-Head Attention and Feed Forward Modules that we provided, along with a LayerNorm from PyTorch.
- iv. Decoder Module that stacks multiple decoder layers. Here you will use positional encoding and three masking functions that were provided to you.

Transformer models are difficult to train naively, and can result in suboptimal performance or even fail to converge, if improperly designed. It's crucial to carefully select and tune various hyperparameters to achieve optimal performance. However we assure you that there are many

configurations where they will be possible to train well. To find one of these you must explore the following:

- **Number of attention heads:** In the multi-head attention mechanism it is a critical hyperparameter, with a typical range being 2 to 10 heads. Increasing the number of heads allows the model to capture different aspects of the input data in parallel, potentially improving its ability to learn complex patterns.
- **Number of layers:** The Number of layers in both the encoder and decoder parts of the transformer also significantly impacts its performance. A range of 2 to 8 layers is common, with more layers providing greater representational capacity at the cost of increased computational complexity and potential over-fitting.
- **Model dimensions:** Including the convolution dimension (128-256), the feedforward dimension (512-2048), and the encoder-decoder attention dimension (256-1024), are crucial for determining the size and capacity of the model. Larger dimensions can enable the model to capture more information but also increase the risk of over-fitting and the computational burden.
- **Optimizers:** The choice of optimizer is another important factor, with options such as Adam, AdamW, and SGD (Stochastic Gradient Descent) being popular. Each optimizer has its own characteristics and may perform differently depending on the specific task and data.
- **Learning rate schedulers:** like the cosine annealing learning rate and ReduceLR, help in adjusting the learning rate during training to improve convergence and prevent overshooting. The cosine annealing scheduler reduces the learning rate following a cosine curve, while ReduceLR decreases the learning rate when a metric has stopped improving.

By experimenting with these hyperparameters and understanding their impact on the transformer model, it's possible to find a configuration that leads to effective training and strong performance on the target task.

Believe it or not, although this may seem like a lot, you will manage to get all this done in the course of this homework. In the following sections, we will outline how to complete the work.

### 1.3 Getting Started Early

Uniquely, HW4P2 weaves together your work from HW3P2 and HW4P1. As a result, it can, and most likely will take a significant amount of your time compared to previous homeworks, so please start early. With more moving parts, you are more likely to encounter errors, which will take significant time to identify and debug, the models also typically take longer to converge and produce good results that is why we leverage pretraining of of an CNN-LSTM encoder. This assignment can be completed within 5 hours.

### 1.4 Starter Notebook

The starter notebook is divided into two sections; Part I and Part II. In part 1 you will use an CNN-LSTM (we provide the weights in for you) or your HW3P2 Encoder as encoder to the speech transformer (You will need to load your own weights). Part 2, you will add encoder transformer layers. We approximate completing part 1 within 5 hours and part 2 within 6 hours.

The rest of this writeup is structured as follows – Section 3 describes the dataset and dataloader. In Section 2 we provide a brief overview of the starter notebook. In Section 4 we outline the baseline model, including the individual components (encoder, decoder and attention), how to train it, and how to perform inference. Section ?? in the Appendices outline debugging and specific implementation details.

## 1.5 Resources

An important list of resources that will help you implement the transformer architecture in your homework:

- i. The first transformer lecture in this course by Akshat Gupta. **Don't miss this amazing lecture.**
- ii. The Illustrated Transformer by Jay Alammar.
- iii. Batoool Haider did three episodes on transformers, providing in-depth explanations and practical insights into the architecture and its applications.
- iv. PyTorch's Transformer Documentation.
- v. TensorFlow's Transformer Model Tutorial.
- vi. An amazing lecture about attention and awareness from neuroscience perspective by Nancy Kanwisher.

By leveraging these resources, you will gain a comprehensive understanding of transformer architectures and how to implement them effectively in your homework.

## 2 Notebook Walkthrough

This section serves to provide a brief overview of the different sections in the starter notebook.

- i. **Configurations.** This section lays out the major hyperparameters inside the notebook. Ranges associated with some of the parameters are given to you to narrow down your experimental search space. Please get comfortable with the parameters listed and their significance to different aspects of the model. Key comments are added to aid you in your ablations.

Several notes: (1) the Transformer Encoder parameters "enc-num-layers" and "enc-num-heads" are for Part II of the homework, where you will add transformer blocks to your existing encoder. (2) the "d-model" and "d-ff" correspond to the hidden representation of the model and feed forward layers, respectively.

- ii. **Character-based LibriSpeech (HW4P2).** This section holds the Dataset classes as well as Dataset and DataLoader instantiation. Please read all comments within the dataset class, as there are some small, key changes from HW3P2. The major change is the need for both shifted and golden transcripts, which helps the model learn the right prediction at each timestep. Refer to Section 3 to learn more about the need for shifted and golden transcripts.
- iii. **Introduction.** This section provides a description of the Transformer architecture used in this homework as well as some rationale as to why attention mechanisms perform so well on nuanced data like speech recordings. Additionally, this section defines some necessary utilities and modules that are extremely helpful for the implementation of the model. Please read through the code and comments to understand all the utilities and modules given.

NOTE: You don't have to change any of the given code; however, you must know what each module is doing to correctly implement the model architecture and answer the Canvas quiz.

- iv. **Part I.** In this section, you will use a pretrained encoder (either the reference simple CNN-LSTM encoder given or your own HW3P2 encoder) and implement the Transformer Decoder. You will then connect these two modules together to create your baseline network.

In this section you are tasked with the following:

- A. implement the **DecoderLayer** class.
- B. implement sections of the **Decoder** class.
- C. implement the **SpeechTransformer** model by merging the **Encoder**, **Decoder** classes.
- D. load a pretrained encoder into the encoder module and begin training.

The pretrained encoder is a very important aspect of this homework. Starting with a simple RNN-based encoder allows the neural network to first learn the basic structures and patterns in the data. When we later introduce more complex mechanisms (like a Transformer Encoder in Part II), the model already has a foundational understanding of the data. This foundation can lead to a smoother and more effective learning process when the complexity of the model increases. In other words, we are using a simple, pretrained RNN encoder as the basis for our transformer decoder and incrementally adding complexity!

- v. **Part II.** In this section, you will build off of your Part I model by adding transformer blocks to your existing RNN encoder, thereby making it a true Transformer Encoder. This incremental update to the model complexity with multi-head self-attention blocks will greatly improve model performance by learning more nuanced, selective patterns in the dataset.

In this section you are tasked with the following:

- A. implement the **EncoderLayer** class.
- B. implement the **Encoder** class.
- C. implement the **FullTransformer** model by merging the **Encoder** class and **Decoder** class. Note that the Decoder here is Transformer Decoder partially trained from Part I.
- D. load the Encoder and Transformer Decoder from Part I into the model.
- E. Freeze the weights of all the components of the model except for the added transformer block(s) and continue the training process.
- F. After about 3 epochs, when the encoder layers must have adjusted to the pretrained parts of your model. Unfreeze weights and train the entire network..

The transformer encoder plays a crucial role by processing the input sequence into a high-dimensional space, capturing complex relationships between elements. It uses attention to understand the context around each feature, enabling the model to grasp subtle nuances and dependencies, crucial for our speech recognition task and improving performance.

- vi. **Part III.** This is the final, experimental portion of the homework, where you can experiment with other input embedding RNN encoders (i.e HW3P2 encoders) and incrementally adding more transformer blocks. This is more open-ended, and we encourage you to be creative!

### 3 Dataset

You will be working on a similar dataset as in HW1P2 and HW3P2. The training set will comprise input-output pairs, where each input consists of a sequence of 27-dimensional mel-spectral vectors derived from a spoken recording, and the corresponding output consists of the spelled out text transcription. The transcriptions are in terms of 31 characters (the 26 characters in the English alphabet, blank space, “,” [comma] and three special tokens: “<sos>”, “<eos>” and “<pad>”, representing the start and end of a sentence and padding respectively).

The audio in the different training instances are of different lengths, and the output sequences too are (obviously) of different lengths. We also provide you a Python array **VOCABULARY** (refer to the starter notebook) which contains a mapping from the letters to numeric indices. *You must use this list to convert letters from transcriptions to their indices* for the purposes of the competition.

#### 3.1 Dataloader

The **Dataset** class and **Dataloader** for this homework are pre-designed and will be provided. The dataloader prepares the input for the speech transformer’s decoder by shifting the transcriptions. Each decoder input sequence starts with a <SOS> token to signal the beginning and is used to process the subsequent characters. Conversely, the target sequence for comparison is shifted in the opposite direction to include the <EOS> token at the end, indicating the end of a sentence. This setup represents a full teacher forcing strategy, where the actual target sequences are always provided to the decoder during training. The decoder constructs the autoregressive matrix using an attention mechanism and a triangular mask to prevent future characters from influencing the prediction of the current character. The dataloaders will supply you with the shifted target sequence, e.g., “<SOS> H E L L O,” and the corresponding golden target, e.g., “H E L L O <EOS>.”.

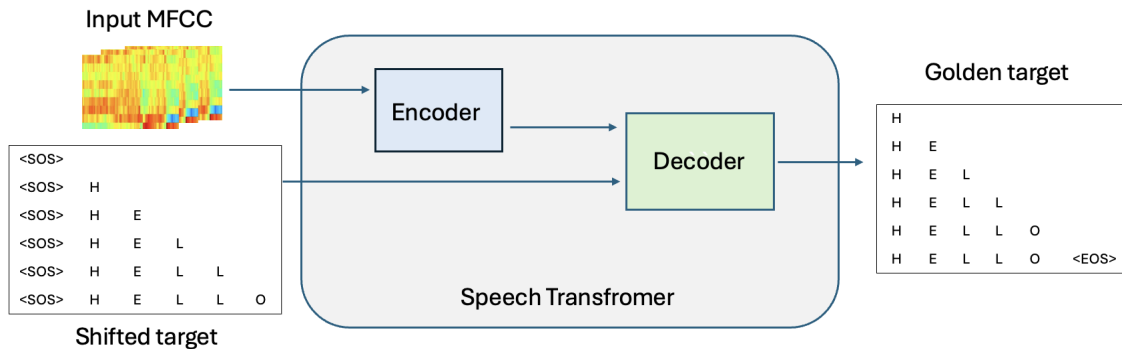


Figure 2: The input Mel-Frequency Cepstral Coefficients (MFCC) are fed into the encoder. The shifted target sequence, beginning with a start-of-sentence token (<SOS>) and excluding the end-of-sentence token (<EOS>), is used by the decoder to predict the subsequent token in the sequence. The golden target sequence is the true next token the model aims to predict at each timestep, including the <EOS> token but excluding the initial <SOS>.

## 4 Speech Transformer

In this homework, you will build a Transformer model Vaswani et al. [2017], Dong et al. [2018] for speech recognition. Unlike the original application of the Transformer for machine translation, we will adapt the architecture to transcribe spoken language into text. The Transformer model will be divided into three main components:

By the end of this homework, you will have implemented a complete Transformer model tailored for the task of speech recognition, gaining a deeper understanding of this powerful architecture and its applications in natural language processing. Specifically, we will be implementing the Speech-Transformer in Figure 3

We acknowledge that a transformer architecture consists of multiple components, which might be confusing. However, the subsequent sections will guide you through this transformer, step by step. We recommend familiarizing yourself with the architecture before beginning to code, as this understanding is crucial for working with the starter notebook without encountering problems.

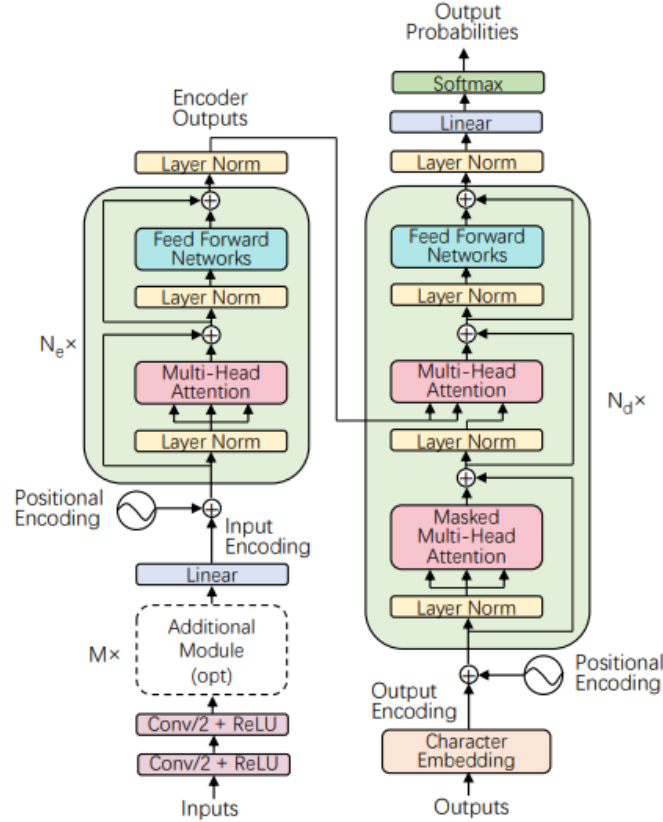


Figure 3: Baseline Model: Speech Transformer

The above architecture can be divided into two main parts. The Encoder and the Decoder.

### 4.1 Encoder

The encoder in a speech recognition Transformer model processes the input speech features, such as Mel-Frequency Cepstral Coefficients (MFCCs), and converts them into a sequence of embeddings that capture the contextual information of the speech. These embeddings capture meaningful information from the speech, which the decoder then uses to predict characters in

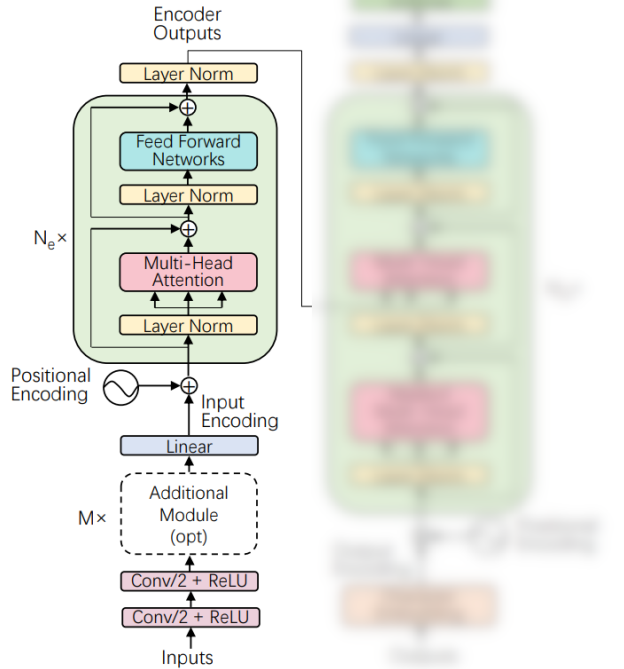


Figure 4: Encoder part of Transformer

the English alphabet. The encoder consists of convolutional layers, positional encoding, and multiple encoder layers each comprising self-attention mechanism and feed-forward networks. More specifically, this section walks through the portion of the architecture as shown in Figure 4

#### 4.1.1 CNN-LSTM as Encoder and as Input Embedding

In our speech processing architecture, Mel Frequency Cepstral Coefficients (MFCCs) are pivotal, transforming audio inputs into a highly informative representation for neural network analysis. To enrich this representation and enable nuanced interpretation of speech, we integrate a pretrained Long Short-Term Memory (LSTM) network. Pretraining ensures that our model can effectively leverage historical knowledge of sequential data, capturing the intricate temporal relationships inherent in speech.

We refine the input to the LSTM with a 1D convolutional layer, acting as an embedding mechanism that processes the MFCC features. This layer is specifically chosen for its proficiency in extracting meaningful patterns from the temporal sequence of MFCCs, effectively condensing and enhancing the feature set before it reaches the LSTM. The 1D convolution layer thus serves as a critical bridge, translating raw audio features into a format that the LSTM can efficiently process.

Freezing certain layers within the CNN-LSTM is a strategic move aimed at preserving the pre-trained model’s intrinsic understanding of speech dynamics. This approach ensures that the core competencies of the CNN-LSTM, honed on vast datasets, remain intact and focused on deciphering complex speech patterns. Meanwhile, the trainable parts of the network adapt to the nuances of the specific task at hand.

This design not only facilitates a deep understanding of speech but also provides a robust framework for experimentation and learning. By integrating a 1D convolutional layer with a pretrained and partially frozen CNN-LSTM, we create a model that is both sophisticated in its analytical capabilities and grounded in the practical realities of speech processing challenges.



### 4.1.2 Encoder Input Encoding

- i. **CNN-LSTM:** Instead of using convolutional modules followed by a linear layer as in the baseline model, we use a pretrained CNN-LSTM which has already been trained in extracting relevant features from MFCC data.
- ii. **Positional Encoding:** Positional encoding in speech transformers adds unique markers to input sequences, enabling the model to recognize word order and temporal relationships within speech data. This technique allows transformers, which inherently lack sequence awareness, to process speech with an understanding of time and sequence, crucial for capturing the nuances of spoken language. Refer to Section 8.4 for more details. We have also implemented this module for you in the starter notebook. If you have any doubts or confusion, please feel free to attend an OH or post a question on piazza for further clarification.

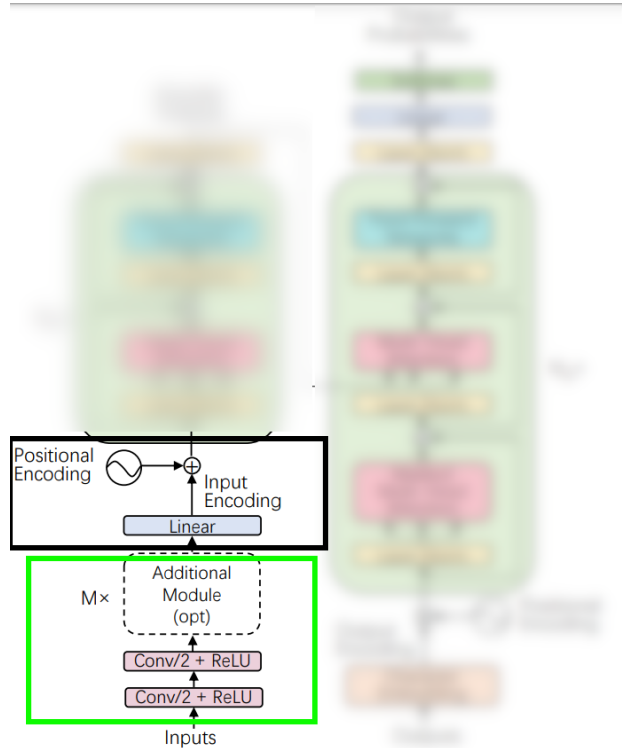


Figure 5: Encoder Input Encoding

### 4.1.3 Encoder Layers

Congratulations, now you are ready to extract features and add positional encoding for a transformer encoder. Your inputs are ready and we can now dive into the encoder layers which are the main components of the transformer encoder. As shown in Figure 6, this part of the architecture consist of 3 main modules which are:

- i. **Layer Normalization:** This helps stabilize the learning process. By normalizing the input to various modules, the models performance and ability to generalize can be improved.
- ii. **Multi-Head Attention:** This module facilitates the parallel processing of speech signals, enabling the model to focus on different aspects of the input sequence simultaneously. This design allows for a richer representation of the input by capturing various linguistic and acoustic features from multiple perspectives, thus improving the encoder's ability to understand and encode complex speech patterns and relationships within the data. Please refer to section D for more information.

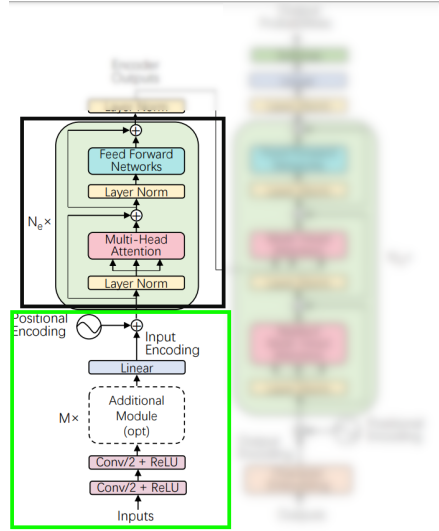


Figure 6: Encoder Layers

- iii. **Feed-Forward Network (FFNs):** serve as a critical component following the multi-head attention mechanism. Each FFN consists of two linear transformations with a nonlinear activation function in between. This structure allows the model to introduce additional complexity and nonlinearity, enabling the encoding of more abstract representations of the speech data. The FFNs operate on each position separately and identically, ensuring that the model can enhance its understanding of the speech input by applying the same transformation across all positions, thereby enriching the feature space before passing it onto the next layer.

Now that we understand the various parts of the encoder layers let us see how to put it together. In the starter notebook we have created a module called the EncoderLayer which you are to complete.

- i. **Module Initialization:** First we need to initialize all the modules needed in the encoder layer. This involve creating 2 Layer Norm, 1 Multi-Head attention Module and one Feed-Forward Neural Network Module. All the modules have been implemented for you check the module definitions and see how to initialize them. The parameters to initialize them should be passed to the init method of your EncoderLayer module.

## ii. Forward Method

- A. **Layer Norm 1:** As mentioned above normalizing our inputs to a module is very important. In the forward method of the encoder layer, we receive the input encoding which is the combination of the extracted features and positional encoding as described in the previous sections. This input encoding passes through the first layer normalization module.
- B. **Multi-Head Self-Attention Module:** Next we need to pass the output of our first layer norm to the multi-head self-attention module. In this section, make a minor modification by using a mask with the multi-head attention module. The reason for this is, our inputs are padded and to be able to attend well, we need to mask those padding we added to the input. Hence we also receive the pad mask and the self attention mask as parameters to the forward method of the EncoderLayer. The padding mask is set to one for all tokens except for the pad token, which is set to zero. The self-attention mask is also a padding mask but for the attention weights. However, in this case, it is the opposite: the positions of non-pad tokens are set to zero, and the pad token is set to one. This is because, during the attention computation, we will fill the positions with ones (representing pad tokens)

with infinity to ensure that the attention scores for padding tokens are zero. As a result, padding tokens do not contribute to the encoder representation. After the layer norm, we pad the outputs, and use this as parameters to our multi-head attention module forward method which then produces the attention scores.

- C. **Residual Connection 1:** As indicated in the architecture, the output of the multi-head attention is combined with the input encoding. This connection helps in alleviating the vanishing gradient problem, enabling deeper models to be trained more effectively.
- D. **Layer Norm 2:** Similar to the first layer norm, the input and the multi-head attention output combined is passed through the second layer normalization.
- E. **Feed Forward Networks:** The results from layer norm 2 is passed through the feed forward network. Its benefits are as described above.
- F. **Residual Connection 2:** Finally, the output of the feed forward network is combined with the residual connection 1 result which was the combination of the inputs and the attention output.

This concludes the encoder layer. Take a moment to reflect on the process through which the input encoding passes through. Does it make sense? If not please visit OH or ask a question of Piazza to get clarification.

#### 4.1.4 Transformer Encoder

Now that we understand all the various building blocks for the transformer encoder, let us put them together to form our Speech-Transformer Encoder.

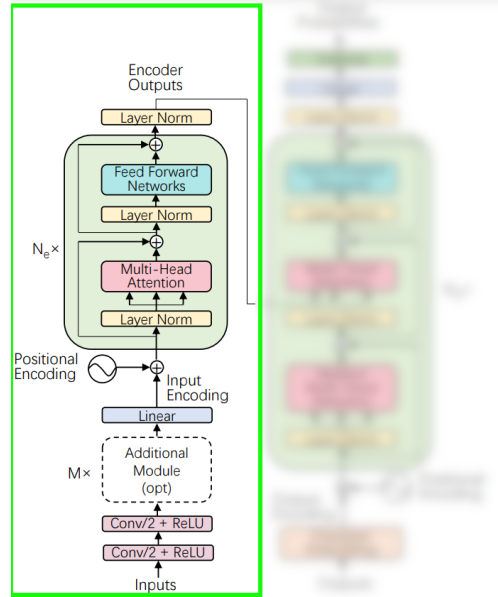


Figure 7: Speech-Transformer Encoder

The encoder part which we will now put together is as shown on Figure 7. To complete this module, we divide the process into two main steps:

##### i. Module Initialization:

- A. **Convolutional Module:** Check the convolutional module implementation provided to understand how to it needs to be initialized. Think about some of the parameters you can play with, for example the number of convolutional layers.

- B. **Positional Encoding:** Check the positional encoding module provided to understand how to initialize it. Think about some of the parameters you can play with.
- C. **Encoder Layers:** Lucky you! We have made the initialization for you. However, make sure you understand what is done because you will do the initialization in the decoder part. Can we increase/decrease the number of encoder layers ? What are the effects and why do you think this impacts the performance ?
- D. **Layer Normalization:** After the encoder layer we have one more layer normalization to prepare the encoder outputs which are used in the decoder section.

## ii. Forward Method

- A. **Convolutional Module:** We need to pass the input to our encoder to the convolutional module we initialized earlier. How is this done in Pytorch given the convolutional module is a neural network module.
- B. **Positional Encoding:** We create a positional encoding of the input using the Positional encoding module we initialized earlier.
- C. **PAD Mask:** We have provided a function to help create a padding mask. Check the utilities section of the notebook and also ask yourself why are we creating this.
- D. **Attention Mask:** We have provided a function to create the attention mask, please have a look at it and ask yourself why we create an attention mask. Remember we are dealing with Multi-Head self-attention.
- E. **Encoder Layers:** We initialized our encoder layers earlier. If we have more than one encoder layer, the output of the first layer serves as input to the second layer.
- F. **Layer normalization:** The output of the final encoder layer is passed through a layer norm layer to prepare it for our decoder.

Congratulations! You have successfully built a transformer encoder. Take a moment to reflect through the steps, what happen in each module and how the modules are connected. Feel free to ask any questions for more clarification if there is any part you don't understand.

## 4.2 Decoder

The decoder in a speech recognition Transformer model takes as input the embeddings generated by the encoder and decodes them into sequences of characters representing the transcribed text. It includes positional encoding, multiple decoder layers each comprising self-attention and cross-attention mechanisms, and feed-forward networks. During training, the decoder predicts the next token in the target sequence based on the previously generated tokens. This is achieved through self-attention and cross-attention mechanisms, where the decoder attends not only to the encoder outputs but also to its own previously generated outputs.

In this section we will be exploring the different components in the decoder part of the architecture as shown in Figure 8

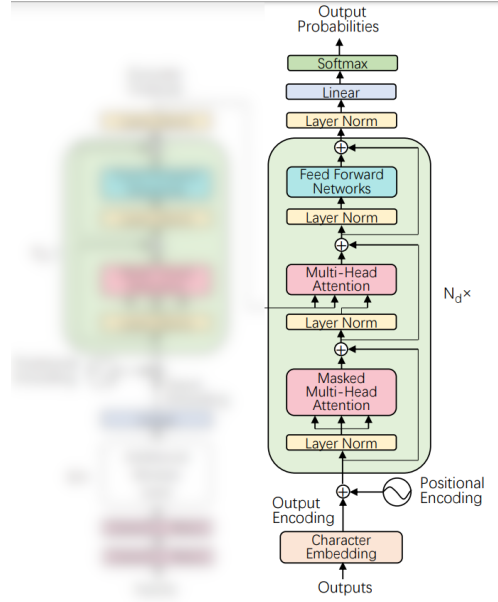


Figure 8: Decoder Part of Speech Transformer

Similar to section ?? on the encoder, we will go through the components that make up the decoder one after the other. Let's get started

### 4.2.1 Decoder inputs

The decoder learns to decode the inputs to some target output. During training, the decoder accepts inputs from the encoder as well as a shifted target (indicated as outputs in the decoder side of the architecture as in Figure 8). During inference, outputs are the generated tokens from the decoder itself. For more on shifted targets, please refer to Section 3.

### 4.2.2 Output Encoding

Similar to the encoder encoding, we need to encode one of the inputs to the decoder. We do this by first encoding the targets using character embedding. Since this representation is already standard and does not need to be learned, we use an already existing neural network. The character embedding is combined with the positional encoding to produce the output encoding which is one of the inputs to the decoder layers.

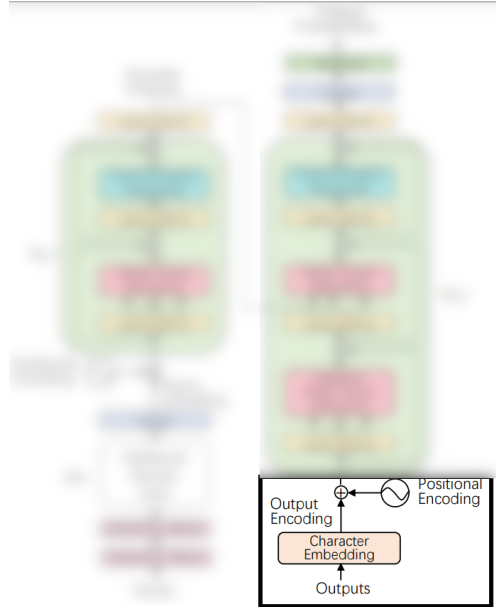


Figure 9: Decoder Output Embedding

#### 4.2.3 Decoder Layers

The Transformer Decoder is designed to take into account the output of the Transformer Encoder and generate the target sequence. It also contains self-attention layers, but with a slight modification to prevent positions from attending to subsequent positions. This masking ensures that the predictions for position  $i$  can only depend on the known outputs at positions less than  $i$ . We discuss this in more details in 8.3

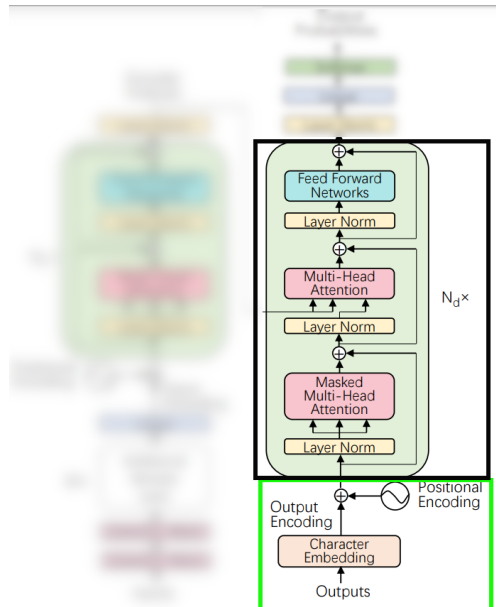


Figure 10: Decoder Layer

There are four main modules in a single decoder layer;

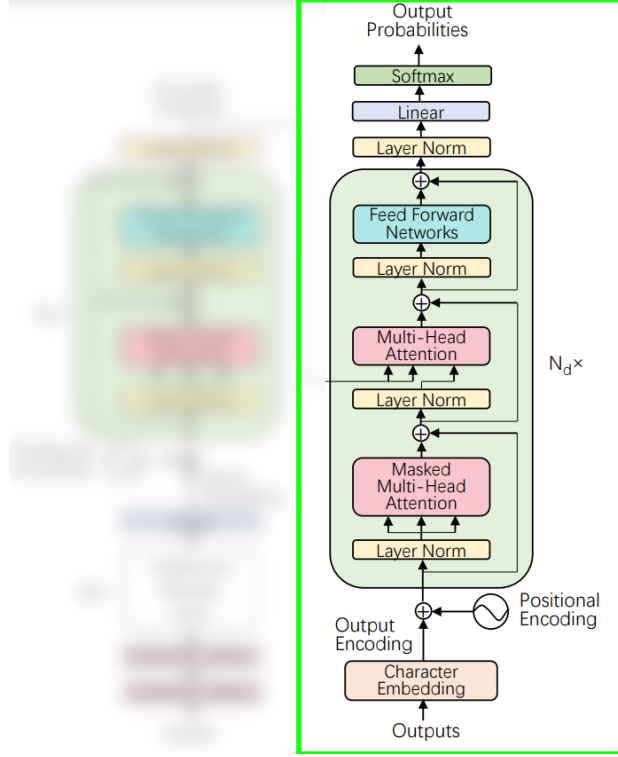


Figure 11: Transformer-Decoder

i. **Layer Normalization:** Same as explained in the Encoder

ii. **Masked Multi-Head Self-Attention:**

In the decoder part, the masked multi-head self-attention module uses a lookahead mask, also known as a causal mask, in addition to the padding mask, unlike the encoder which was using attention mask. This lookahead mask is crucial for ensuring that the decoder can only attend to earlier positions in the output sequence, preventing it from seeing future tokens during training and inference. This is essential for maintaining the autoregressive property of the decoder, where each token is predicted based on the preceding tokens.

iii. **Feed-Forward Network (FFNs):** Same as explained in the encoder section.

iv. **Multi-Head Self-Attention (Cross-Attention):** Cross-attention plays a crucial role in this process by allowing each step of the decoder to attend to different parts of the encoder's output. This is achieved by calculating attention scores that determine how much focus should be placed on each part of the encoder output when generating the current element of the decoder output. For more information please check Section 8

#### 4.2.4 Transformer Decoder

Now that you have all your modules ready. To build the decoder part, we follow the same approach as with the encoder. Initializing all modules required and connecting them as shown in the figure 11

If you have problems with the connection please check how it is done in the encoder section. However, note that the outputs of the decoder layer passes through another layer normalization followed by a linear layer and finally a Softmax to output the probabilities. How do we convert these probabilities to actual tokens ? please check the section 5.1 on inference.

### 4.3 Combining it All

Now we have all parts of our transformer implemented and working. We need to combine the encoder and decoder. We have provided a module called Speech Transformer which does this combination. It is as easy as initializing the encoder and decoder module, writing a forward function which receives inputs and target output, passes the inputs to the encoder layer and get the encoder outputs, use these encoder outputs in the decoder.



## 5 Training

Training an encoder-decoder transformer model involves several key steps. The transformer model is based entirely on self-attention mechanisms and does not use recurrent neural networks. During the forward pass, the encoder processes the input sequence to produce a set of representations. These representations are then passed to the decoder, which takes both the encoder's output and the target sequence (shifted right) as inputs to generate the output sequence. The decoder uses both self-attention to understand the relationships within the target sequence and cross-attention to focus on relevant parts of the input sequence.

### 5.0.1 Loss Function

The loss function is typically the cross-entropy loss between the predicted output sequence and the ground truth sequence. The loss is computed for each position in the sequence and then averaged.

### 5.0.2 Overall Training Procedure

The overall training procedure can be summarized as follows:

---

```
# X is the input sequence, T_shift is the ground-truth output sequence shifted right, T_gold
#   is the ground-truth output sequence.
# model is the transformer model (encoder-decoder)
function gradient = ComputeGradient(X, T_shifted, T_golden, model)
    encoded = model.encoder(X)
    decoded = model.decoder(encoded, T_shifted)
    L = CrossEntropyLoss(decoded, T_golden)
    gradient = BackPropagate(L)
end
```

---

Here, `model.encoder` and `model.decoder` represent the forward passes through the encoder and decoder, respectively. `CrossEntropyLoss` computes the loss, and `BackPropagate` performs the gradient backpropagation for updating the model parameters.

## 5.1 Inference

Inference in the Speech Transformer model is a sequential process, particularly during the generation of output sequences. Initially, the encoder encodes the entire input audio sequence into a set of hidden representations using self-attention.

During the decoding phase, the model generates the transcription sequence (sequence of characters) one element at a time in an autoregressive manner. This means that the generation of each element in the output sequence is conditioned on the previously generated elements. The decoder uses the encoded representations and the partially generated output sequence to predict the next element. This process continues iteratively until a predefined end-of-sequence ( $\langle \text{EOS} \rangle$ ) token is produced or a maximum output length is reached.

The autoregressive nature of the inference process ensures that the model takes into account the context provided by the preceding elements in the output sequence, which is crucial for maintaining coherence and relevance in the generated sequence.

So the entire process can be encapsulated by the following pseudocode:

---

**Listing 1** Inference

---

```
# X is the input audio feature vector sequence
# T is the <SOS> token

function O = Inference(X)
    E = Encoder(X)
    T = [<SOS>]
    for i in range(maximum_sequence_length):
        NT = Decoder(E, T)
        T += NT
        if NT == <EOS>:
            stop
    end
```

---

### 5.1.1 Greedy Search

The naive approach is greedy search, also the easiest decoding method for inference. All you would have to do is to draw the character with the highest probability at each time step from the output probability distribution for the entire batch of predictions you get from your model. Pseudocode would look like this:

---

**Listing 2** Greedy Decoding

---

```
# X is the input audio feature vector sequence
# T is the <SOS> token

function O = Inference(X)
    E = Encoder(X)
    T = [<SOS>]
    for i in range(maximum_sequence_length):
        NT = Decoder(E, T)
        NT = argmax(NT)
        T += NT
        if NT == <EOS>:
            stop
    end
```

---

While this usually generates reasonable results, due to the autoregressive nature of the model, the greedy method does not however guarantee that the final output sequence is the most likely one. Hence, in order to get the most likely final output, exploring other techniques like random search and beam search, are introduced in Appendix ?? and Appendix ?? respectively.

## 6 Common Errors and Fixes

During the implementation of the Transformer model for speech recognition, students may encounter several common issues. Here are some guidelines to help address these problems:

### 6.1 Gradient Clipping

Gradient clipping is a technique used to prevent the exploding gradient problem in neural networks. It is recommended to implement gradient clipping with a norm of 1 or 2. This will add a ceiling to your gradient, ensuring that it does not exceed the maximum norm value. To clarify, gradient clipping is optional but highly recommended for stabilizing training.

### 6.2 Debugging NaN Loss

Encountering NaN (Not a Number) loss during training can be a frustrating issue. Here are some steps to debug and fix this problem:

- (a) **Check Learning Rate:** A high learning rate can cause the model to diverge, leading to NaN values. Try reducing the learning rate and observe if the issue persists.
- (b) **Inspect Weight Initialization:** Improper weight initialization can lead to NaN loss. Ensure that the weights are initialized using a method appropriate for the activation functions in your network (e.g., Xavier initialization for sigmoid activations).
- (c) **Monitor Gradients:** Check if the gradients are exploding, which can lead to NaN values. You can print the gradients or use a tool like TensorBoard to visualize them. If you notice large gradients, consider implementing gradient clipping to prevent them from exceeding a certain threshold.
- (d) **Use a Different Optimizer:** Sometimes, switching to a different optimizer can help stabilize training. For example, if you're using SGD, try using Adam or RMSprop and see if the issue is resolved.
- (e) **Add Small Constants:** In operations that involve division or logarithms, adding a small constant (e.g.,  $1e-8$ ) to the denominator or argument can prevent division by zero or taking the logarithm of zero, which can lead to NaN values.
- (f) **Check for Inf Values:** Ensure that there are no infinite values in your data or intermediate computations. Infinite values can propagate through the network and result in NaN loss.
- (g) **Reduce Model Complexity:** A very complex model might be more prone to numerical instability. Try simplifying your model architecture and see if the issue persists.

By systematically addressing each of these potential causes, you can identify and fix the source of NaN loss in your training process.

### 6.3 Out Of Memory (OOM)

Out of Memory errors occur when your model requires more memory than is available on your GPU. To resolve this, you can reduce the batch size, use gradient accumulation, or simplify your model architecture.

### 6.4 Dataloader and Dataset

For training, use a combination of 100 and 360 hours of data. Ensure that the dataloader is correctly configured to provide the expected values, including the correct input feature dimensions and batch size.

## 6.5 Concatenating the Context Vector with the Embeddings

To concatenate the context vector with the embeddings, you can use a function like `torch.cat()` in PyTorch. Ensure that the dimensions of the context vector and embeddings align correctly for concatenation.

## 6.6 Training Duration and Monitoring

The training time for one epoch of a Transformer model can vary based on factors such as the size of the model, the size of the dataset, and the hardware used for training. On Google Colab, using a standard GPU, you can expect an approximate training time of 1-2 hours per epoch for a moderate-sized Transformer model on a medium-sized speech dataset. However, this is just an estimation, and actual training times may vary.

Additionally, it is important to save checkpoints regularly during training. This practice helps to avoid losing progress in case of interruptions and allows for resuming training from the last saved state. Implementing a checkpointing mechanism is especially important when training on platforms like Google Colab, where sessions can be terminated after a certain period of inactivity.

Debugging common issues involves carefully examining the loss curves, checking the gradients for any abnormalities, and ensuring that the data pipeline is functioning correctly. If the model is not converging, consider adjusting the learning rate, experimenting with different optimization algorithms, or revising the model architecture.

## 7 Conclusion

This homework is a true litmus test for your understanding of concepts in Deep Learning. It is not easy to implement, and the competition is going to be tough. But we guarantee you, that if you manage to complete this homework by yourself, you can confidently claim that you are now a **Deep Learning Specialist!** Good luck and enjoy the challenge!

# Appendices

## A Character Based vs Word Based

We are giving you raw text in this homework, so you have the option to build a character-based or word-based model. Word-based models won't have incorrect spelling and are very quick in training because the sample size decreases drastically. The problem is, it cannot predict rare words (i.e words in testing set that are not present in your training set), thus resulting in a drop in performance. You would have to come up with strategies to handle OOV (Out of Vocabulary) words.

The paper describes a character-based model. Character-based models are known to be able to predict some really rare words but at the same time they are slow to train because the model needs to predict character by character. In this homework, we **strongly recommend** you to implement a character based model, since most TAs are familiar with the character-based model, so you can receive more help for debugging.

## B Transcript Processing

HW4P2 transcripts contain letters. You should make use of the `LETTER_LIST` array (consisting of all characters/letters in the vocabulary) provided in the starter notebook and convert the transcripts into their corresponding numerical indices to train your model.

Each transcript/utterance is a separate sample that is a variable length. We want to predict all characters, so we need a [start] and [end] character added to our vocabulary, which are `<eos>` and `<eos>`, respectively.

Hints/Food for thought:

- Think about special characters that are included in the vocabulary.
- You need to build a mapping from char to index, and vice versa. Reverse mapping is required to convert output, i.e., index, back to char.
- How do you deal with the length of the final output when evaluating on validation/test sets?

## C Single Head Attention (Minibatch)

The following pseudocode describes context computation, as outlined in Section ??, for a minibatch. We use loops in the pseudocode – actual python code would be smarter and use tensor operations.

---

**Listing 3** Attention (minibatch)

---

```
# E(:, :, :) = batchsize x max_sequence_length x dim array of embeddings for a minibatch
# mask(:, :, :) = batchsize x max_sequence_length array of binary masks (1 for valid, 0 for
    invalid)
# q(:, :, :) = batchsize x Dq matrix of query vectors
# C(:, :, :) = batchsize x Dv matrix of computed context vectors
# attention_weights = batchsize x max_sequence_length array of attention weights
function C, attention_weights = Attend(E, D, mask=None)
    # WK, WV, WQ must also be specified in actual code
    query_length = length(D(1, :))
    for b = 0:batchsize-1 # Goes over the entire minibatch
        instance_embeddings = E(b, :, :)
        K = instance_embedding * WK      # WK is a dim x Dq matrix
        V = instance_embedding * WV      # WV is a dim x Dv matrix
        Q = D(b, :, :) * WQ
```

---

```

raw_weights = (1/sqrt(query_length)) * Q * transpose(K)
if mask != None:
    # Mask attention weights manually here if not using ignore_index, or if you
    # require a causal mask in the case of a transformer
attention_weight = softmax(raw_weights) # Renormalize to sum to 1
attention_weights(b,:) = attention_weight
C(b,:) = attention_weight * V
end
return C, attention_weights
end

```

---

## D Multi-Head Attention

The attention computed in (Section ??) computes a *single* context vector  $C_k$  to compute output  $O_k$ .  $C_k$  derives and focuses on a single specific aspect of the input that is most relevant to  $O_k$ . It is reasonable to assume that there are *multiple* separate aspects of the input that must all be considered to compute  $O_k$ .

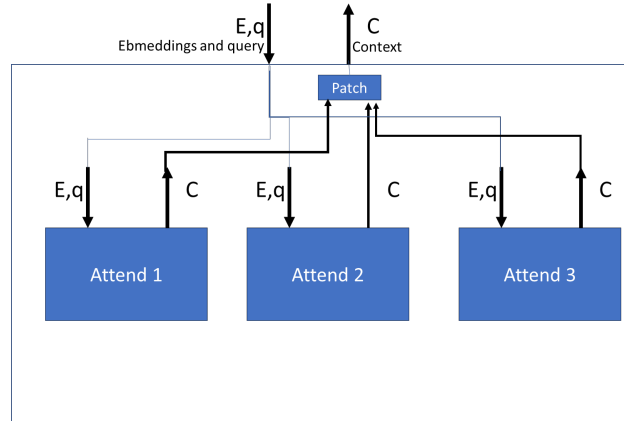


Figure 12

*Multi-head* attention works on this principle to derive multiple context vectors  $C_k^1, C_k^2, \dots, C_k^L$ . Each context vector  $C_k^l$  is computed using a separate attention module of the form explained above (Section ??), with its own learnable parameters  $\mathbf{W}_K^l, \mathbf{W}_V^l$ , and  $\mathbf{W}_Q^l$ . The complete module that computes an individual context vector  $C_k^l$  is called an *attention head*, thus the overall attention framework that computes multiple context vectors is referred to as *multi-head* attention.

The final context vector  $C_k$  that is used to compute  $O_k$  is obtained by concatenating the individual context vectors:  $C_k = [C_k^1, C_k^2, \dots, C_k^L]$  (where  $L$  is the number of attention heads).

Following is the pseudocode for multi-head attention:

---

### Listing 4 Multihead Attention

---

```

# E = (sequence_length, dim) array of speech embeddings for \textit{single} input
# D = (sequence_length, dim) array of character embeddings for \textit{single} input
# context = context vector computed (1,Dv) row vector
# no_of_heads = number of attention heads
function Attention = MultiHead(E, D, mask=None)
    key = E * WK # WK is a (dim, Dk) matrix
    value = E * WV # WV is a (dim, Dv) matrix
    query = D * WQ # WQ is a (dim, Dq) matrix

```



```

# Reshape Q,K,V to compute attention in parallel for each head
query = reshape(query, (1, no_of_heads, Dq//no_of_heads))
key = reshape(key, (sequence_length, no_of_heads, Dk//no_of_heads))
value = reshape(value, (sequence_length, no_of_heads, Dv//no_of_heads))
# Transpose to (no_of_heads, 1, Dkv//no_of_heads)
query = transpose(query)
# Transpose to (no_of_heads, sequence_length, Dkv//no_of_heads)
key = transpose(key)
value = transpose(value)
mask = ? #Similar to single head attention, but now have one for each head
context, attention_weights = Attend(query,key,value,mask) #Assuming Attend works on
    minibatches, Attend here just computes the dotproducts instead of projecting the
    matrices again
context = reshape(context, (1, Dv))
context = context * W_0 #W_0 is a (Dv,Dv) matrix
return context, attention_weights
end

```

---

## 5 Levenshtein Distance

Kaggle will evaluate your outputs with the Levenshtein distance, aka edit distance : number of character modifications needed to change the sequence to the gold sequence. To use that on the validation set, you can use the python-Levenshtein package.

Implementing it is fairly similar to how you did it in HW3P2. You would first need to create the predicted transcript based on your inference method for each item in the batch. Calculate and sum the Levenshtein distance between the predicted transcript and the ground truth for each pair, during your inference loop. Finally, average the sum to obtain the levenshtein distance for the entire batch. Make sure you do not include `<eos>` and `<sos>` tokens in your transcripts when calculating the distance and test predictions.

## 6 Optimizer and Learning Rate

From our empirical studies, we recommend you use Adam or AdamW and start with a learning rate of  $1e^{-3}$ , and experiment with the scheduling methods. On the other hand, the LAS paper uses ASGD as the optimizer with a learning rate of 0.2. The learning schedule involved a geometric decay of 0.98 per 3M utterances (i.e., 1/20-th of an epoch). Feel free to try that if it works for you.

## 7 Transformers and Pretraining

In the context of transformer models and their pre-training, a common strategy involves selectively freezing parts of the model while training others to optimize learning for specific tasks. For instance, when working with transformer-based models, you might focus initially on pre-training the encoder component to capture rich contextual representations from the input data. This pre-training could involve tasks such as masked language modeling, where certain words or tokens are hidden from the model, and it must predict them based on the context provided by the surrounding tokens.

Subsequently, to specialize the model for a particular downstream task (e.g., text classification, question-answering), you might then freeze the encoder to preserve the learned contextual embeddings and focus on training the decoder or a specific task-oriented head with a relevant objective. This approach allows the model to leverage the generalized understanding acquired during pre-training while adapting its output layers to perform well on the target task. By freezing parts of the model, computational resources are conserved, and training time is reduced since only a portion of the model's parameters are updated. This technique of staged training—starting with pre-training on a large corpus to learn

a broad representation, followed by fine-tuning on task-specific data—has proven effective across a variety of NLP applications.

## 8 An important note on Attention

Since the *key* feature of this homework is the *attention* mechanism, it is useful to first explain it briefly. You will need this information for your homework. You have already encountered attention in Lecture 18. We assume that you have an idea of attention, and the concepts behind it. If not, we recommend that you review lectures 18 and 19 and recitation 10. Assuming this prior knowledge, we quickly recap the key concepts of attention that apply to this homework problem before we continue.

### 8.1 An attention head

The encoder in an encoder-decoder model takes in an input sequence of vectors  $X_0, \dots, X_{N-1}$  and produces a sequence of *embeddings*,  $E_0, \dots, E_{M-1}$ . The length of the embedding sequence  $M$  need not be the same as that of the input sequence,  $N$ , but each embedding vector  $E_i$  effectively represents a section of the input. Assuming all the embeddings to be *row* vectors (in the standard python format), the set of embedding vectors can be vertically stacked into a matrix  $\mathbf{E}$ .

In the typical implementation of attention, from the embeddings  $E_i$  the encoder computes two terms, a *key*  $K_i$  and a *value*  $V_i$ . Typically these are computed through a linear transform. Stacking all the key vectors (assumed to be row vectors) vertically into a matrix  $\mathbf{K}$ , and the value vectors into a matrix  $\mathbf{V}$ , these can be computed as  $\mathbf{V} = \mathbf{E}\mathbf{W}_V$  and  $\mathbf{K} = \mathbf{E}\mathbf{W}_K$ , where  $\mathbf{W}_V$  and  $\mathbf{W}_K$  are learnable parameters.

In the decoder, corresponding to each output  $O_k$ , we first compute a *query* vector  $\mathbf{q}_k$ , which is used to “query” the encoder outputs to determine which portion of the input to pay most attention to, to generate  $O_k$ . When the decoder is a recurrent network, the query is typically computed as a linear transform of the recurrent state  $\mathbf{s}_{k-1}$ , and optionally also  $O_{k-1}$ :  $\mathbf{q}_k = \mathbf{s}_{k-1}\mathbf{W}_Q$ , or  $\mathbf{q}_k = [\mathbf{s}_{k-1}, O_{k-1}]\mathbf{W}_Q$  where  $\mathbf{W}_Q$  is a learnable parameter,  $O_{k-1}$  is the (embedding or one-hot representation of the)  $(k-1)^{\text{th}}$  output, and  $[\mathbf{s}_{k-1}, O_{k-1}]$  represents the concatenation of  $\mathbf{s}_{k-1}$  and  $O_{k-1}$ .

In order to produce the  $k^{\text{th}}$  output symbol  $O_k$ , the decoder computes a *context*  $C_k$  from the encoder-derived value vectors:  $C_k = \mathbf{w}_k\mathbf{V}$ .  $\mathbf{w}_k$  is a vector of *attention weights* for the  $k^{\text{th}}$  output, computed as  $\mathbf{w}_k = \text{softmax}(\mathbf{q}_k\mathbf{K}^\top)$ , and has the property that (a) it has as many components as the number of vectors in  $\mathbf{E}$ , (b) all of its components are non-negative, and (c) they sum to 1.0. Ideally, the components of  $\mathbf{w}_k$  corresponding to Value vectors representing the region of the input most relevant to  $O_k$  will be high, and the rest will be low. The context  $C_k$  is input to the decoder (along with any other recurrent state values and previous outputs that your architecture considers) to generate  $O_k$ .

### 8.2 Multi-head attention

The attention computed using the above method (Section 8.1) computes a *single* context vector  $C_k$  to compute output  $O_k$ .  $C_k$  derives and focuses on a single specific aspect of the input that is most relevant to  $O_k$ . It is reasonable to assume that there are, in fact, *multiple* separate aspects of the input that must all be considered to compute  $O_k$ .

*Multi-head* attention works on this principle to derive multiple context vectors  $C_k^1, C_k^2, \dots, C_k^L$ . Each context vector  $C_k^l$  is computed using a separate attention module of the form explained above (Section 8.1), with its own learnable parameters  $\mathbf{W}_K^l, \mathbf{W}_V^l$ , and  $\mathbf{W}_Q^l$ . The complete module that computes an individual context vector  $C_k^l$  is called an *attention head*, thus the overall attention framework that computes multiple context vectors is referred to as *multi-head* attention.

The final context vector  $C_k$  that is used to compute  $O_k$  is obtained by concatenating the individual context vectors:  $C_k = [C_k^1, C_k^2, \dots, C_k^L]$  (where  $L$  is the number of attention heads).

In this homework you will also investigate the relative advantages of multi-head attention over simple attention and the benefits to be obtained through increasing the number of attention heads.

### 8.3 Masking

In transformers, the concept of masking plays a crucial role in controlling the flow of information, particularly for sequence processing. There are three primary types of masks utilized:

1. **Causal Mask:** The causal mask is utilized within the decoder to ensure that the prediction for a given token can only be influenced by the tokens that precede it. This maintains the autoregressive property of language generation. It is mathematically represented as:

$$M_{ij} = \begin{cases} 0 & \text{if } i \geq j, \\ -\infty & \text{otherwise.} \end{cases} \quad (1)$$

Here,  $i$  and  $j$  denote the positions in the sequence. The use of zero values permits attention to previous tokens, while  $-\infty$  blocks attention to future tokens.

2. **Attention Mask:** This mask directly modifies the attention scores, enabling the model to selectively focus on or ignore specific parts of the input. This mechanism is versatile, allowing for dynamic adjustments of the model's focus during training and inference.
3. **Padding Mask:** Since sequences are batched together for efficient computation, they are padded to a uniform length. The padding mask ensures that the model does not attend to these padding tokens, which carry no meaningful information. This is crucial for maintaining the integrity of the model's output. The padding mask is applied as follows:

$$P_{ij} = \begin{cases} 0 & \text{for actual sequence positions,} \\ 1 & \text{for padding positions.} \end{cases} \quad (2)$$

Together, these masks ensure that the transformer model's attention mechanism operates correctly, focusing only on relevant parts of the input and adhering to the desired sequence generation constraints.

### 8.4 Position Encoding

Positional encoding describes the location or position of an entity in a sequence so that each position is assigned a unique representation. There are many reasons why a single number, such as the index value, is not used to represent an item's position in transformer models. For long sequences, the indices can grow large in magnitude. If you normalize the index value to lie between 0 and 1, it can create problems for variable length sequences as they would be normalized differently. There are several requirements for the positional encodings such as; they should have some representation of time, they should be unique for each position and they should be bounded.

In the original transformer paper, they use a smart positional encoding scheme, where each position/index is mapped to a vector. Hence, the output of the positional encoding layer is a matrix, where each row of the matrix represents an encoded object of the sequence summed with its positional information.

#### 8.4.1 Positional Encoding Layer in Transformers

Suppose you have an input sequence of length  $L$  and require the position of the  $k$ th object within this sequence. The positional encoding is given by sine and cosine functions of varying frequencies:

$$P(k, 2i) = \sin\left(\frac{k}{n^{2i/d}}\right)$$
$$P(k, 2i + 1) = \cos\left(\frac{k}{n^{2i/d}}\right)$$

where

$k$ : position of an object in the input sequence,  $0 \leq k \leq l/2$

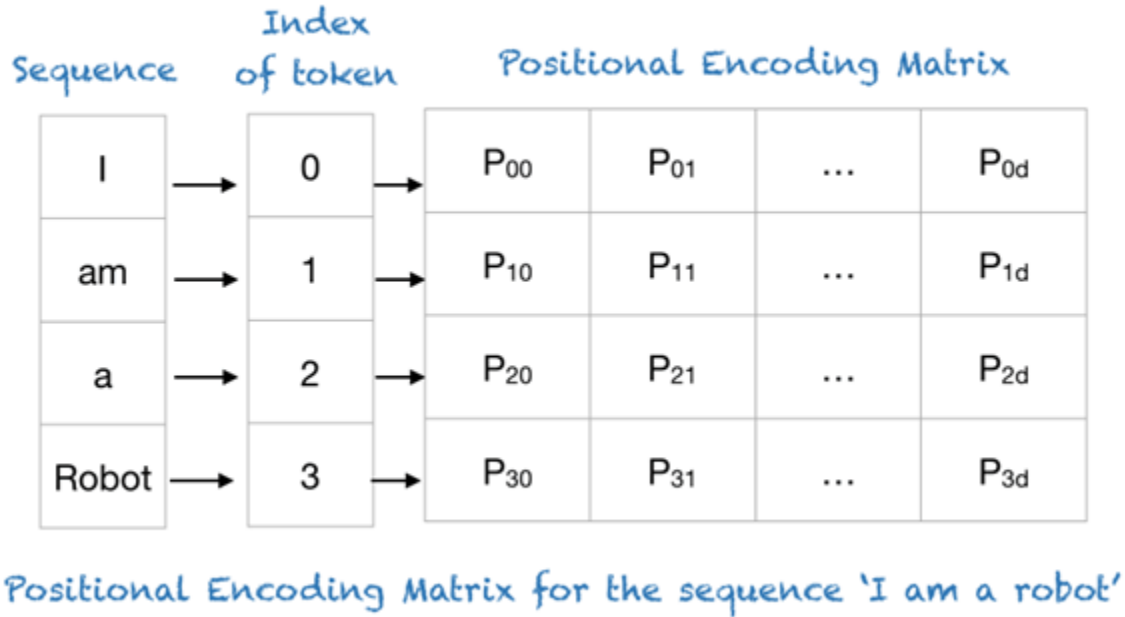


Figure 13: Adopted from machinelearningmastery

d: dimension of the output embedding

$P(k, j)$ : Position function for mapping a position  $k$  in the input sequence to index  $(k, j)$  of the positional matrix.

n: User-defined scalar, set to 10,000 by the authors of Attention Is All You Need.

i: Used for mapping to column indices, with a single value of maps to both sine and cosine functions.

In the above expression, you can see that even positions correspond to a sine function and odd positions correspond to cosine functions.

---

#### Listing 5 Positional Encoding

---

```
function P = PositionEncoding(seq_len, d, n)
    P = # Initialize P(seq_len, d) matrix to hold the positional encoding
    for t = 1:seq_len
        for i = 1:d/2
            denominator = # compute the denominator
            P[k, 2i] = # compute the positional embedding using sin
            P[k, 2i+1] = # compute the positional embedding using cos
        end
    end
    return P
end
```

---

## References

Linhao Dong, Shuang Xu, and Bo Xu. Speech-transformer: a no-recurrence sequence-to-sequence model for speech recognition. In 2018 IEEE international conference on acoustics, speech and signal processing (ICASSP), pages 5884–5888. IEEE, 2018.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.