# Attention

Tarushii Goel

April 2020

## 1 Overview

Recurrent Neural Networks (RNNs) provide outstanding results on a variety of Natural Language Processing tasks; however, they continue to struggle in capturing long-term dependencies in complex sentences. Several advancements, such as LSTMs & GRUs were designed to tackle the problem of learning long-term dependencies, but perhaps one of the most influential is Attention. This lecture will begin by discussing the origin of Attention as a improvement to Seq2Seq for Neural Machine Translation (NMT), discuss and compare a variety of attention mechanisms, and conclude by exploring self-attention multi-head attention, which serve as the groundwork for recent advancements in Transformers and BERT.
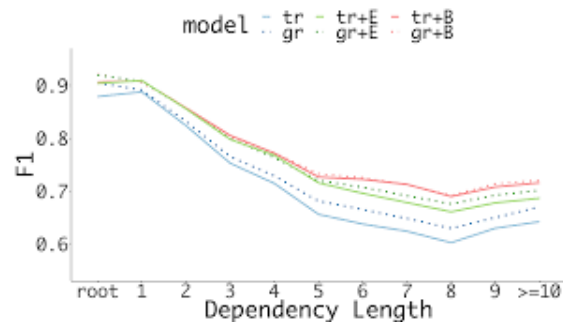


Figure 1: The Long-Term Dependency Problem (Graph of BERT variants, which are state-of-the-art)

## 2 Recap

To refresh you memory, I'll briefly review RNNs and Encoder-Decoders models. RNNs are recurrent; in each iteration they are fed the next input in a sequence and the previous hidden state, which carries accumulated information from other parts of the sentence, and produce a new hidden state.
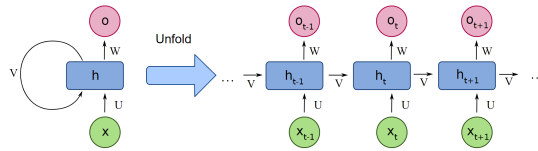
Figure 2: A Diagram of an RNN

Seq2Seq are a unique RNN architecture used for NMT. It is composed to two RNNs; an encoder, which 'encodes' the sentences into a 'context' vector (the last encoder hidden state, a latent representation of the entire sentence), and a decoder, which takes the 'context' vector and decodes it into a translated sentence.
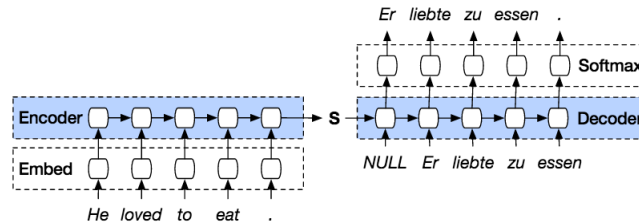


Figure 3: A Diagram of an RNN

# 3   What is Attention?

When we think about the English word "Attention", we know that it means directing your focus at something and taking greater notice. The Attention mechanism in Deep Learning is based off this concept of directing your focus, and it pays greater attention to certain factors when processing the data.

In broad terms, Attention is one component of a network's architecture, and is in charge of managing and quantifying the interdependence:

1. Between the input and output elements (General Attention)

2. Within the input elements (Self-Attention)

Say we have the sentence "How was your day", which we would like to translate to the French version - "Comment se passe ta journée". What the Attention component of the network will do for each word in the output sentence is map the important and relevant words from the input sentence and assign higher weights to these words, enhancing the accuracy of the output prediction.

Now, the question remains, how do we calculate attention weights?
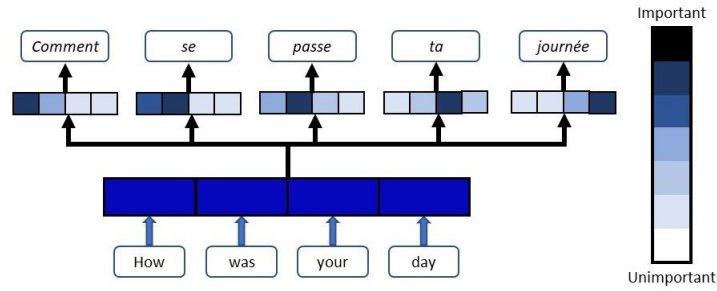
Figure 4: Example of Attention

# 4 Types of Attention Calculation

## 4.1 General Attention

General Attention is used to improve Seq2Seq models. Traditionally, the only information carried over from the encoder to the decoder is the last encoder hidden state, but with Attention, a weighted sum of the encoder hidden states of every time-step is appended to the input. Given a query q and a set of key-
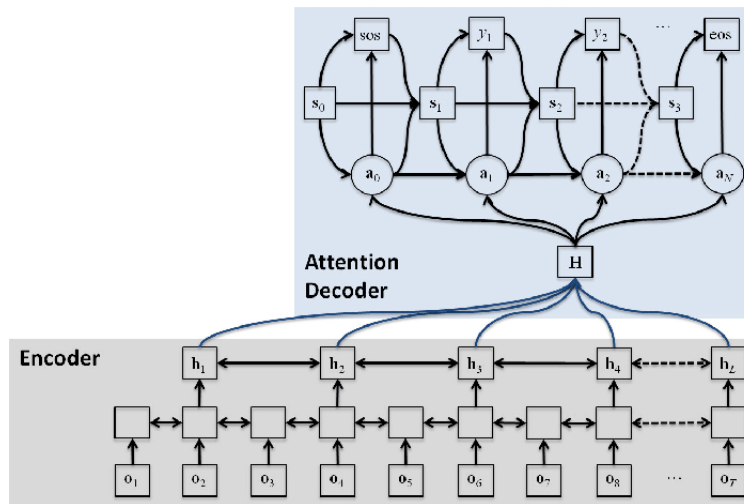


Figure 5: General Attention

value pairs (K, V), attention can be generalised to compute a weighted sum of the values dependent on the query and the corresponding keys. The query determines which values to focus on; we can say that the query 'attends' to the values. In the Seq2Seq, the query is the previous decoder hidden state $s_i - 1$ while the set of encoder hidden states $h_0$ to $h_n$ represented both the keys and the

values. The alignment model, is a row of weights applied to the decoder hidden state and the encoder hidden states. To calculate an 'alignment score' for any, calculate the dot-product $s_i^T * h_i$. Once the 'alignment scores are calculated for $h_{(}0 - n)$, apply a softmax to get the weights.

## 4.2 Self-Attention

With self-attention, each hidden state attends to the previous hidden states of the same RNN. Here $s_t$ is the query while the decoder hidden states $s_0$ to $s_t - 1$ represent both the keys and the values.
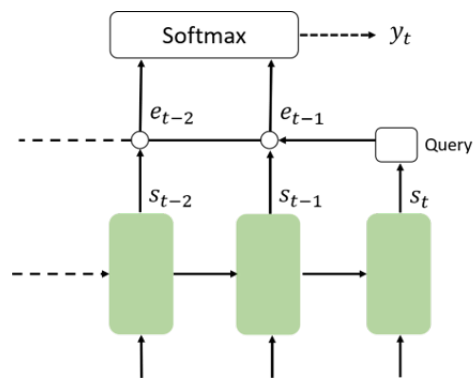


Figure 6: Self-Attention

## 4.3 Multi-Head Attention

When we have multiple queries q, we can stack them in a matrix Q. Each query can be thought of as a different 'head.' In the figure below, two attention heads are displayed (biege and green). The query word is "it". The first attention head is focusing more on the "animal" while the second in green is focusing on "tired". These multi-head mechanisms have been found to help in translation tasks.
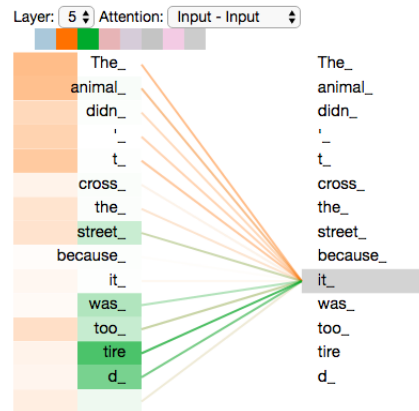
Figure 7: Multi-Head Attention

# 5 Code

I coded a standard Seq2Seq model with attention in PyTorch and put the code here: `https://github.com/2022tgoel/random_ml/blob/master/seq2seq.ipynb`. I highly recommend implementing the techniques I discussed in this lecture on your own and using this as a reference. Let me know if you have any questions!

# 6 References

I got a lot of the explanation and graphics for Attention from this blog post: https://blog.floydhub.com/attention-mechanism/.