



Monetize Your Data Analysis

Share data-science insights with 1M+ readers on TDS and collect payouts through the Author Payment Program.

S

ARTIFICIAL INTELLIGENCE

LATEST

Transformers Explained Vi (Part 2): How it works, step by step

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

A Gentle Guide to the Transformer under the hood, all step by step operation.

WRITE FOR TDS

Ketan Doshi

Sign in

Jan 2, 2021 12 min read

Submit an Article

INTUITIVE TRANSFORMERS SERIES NLP

Showcase Your Expertise

Share data-science insights with 1M+ readers and get compensation with the TDS Author Payment Program.

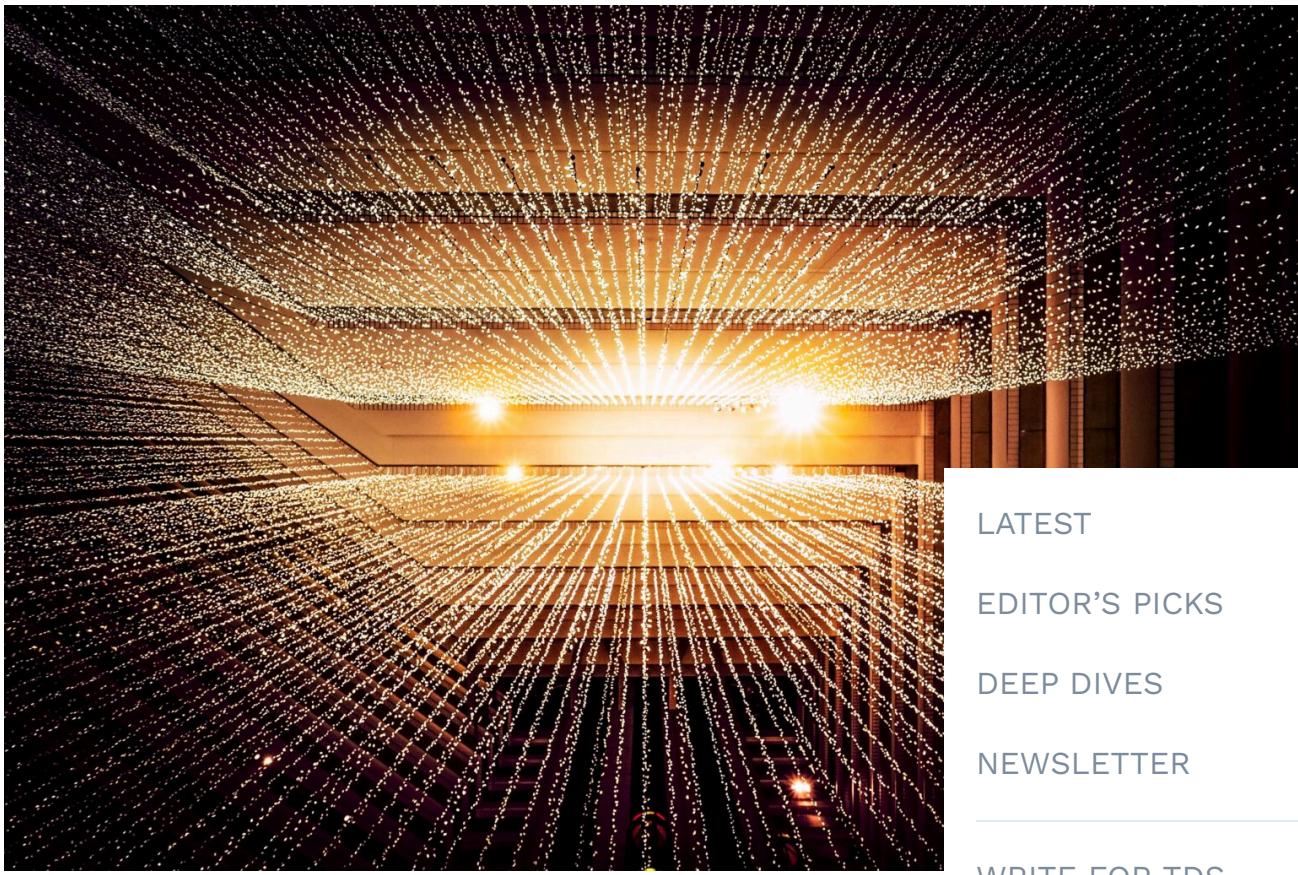


Photo by [Joshua Sortino](#) on [Unsplash](#)

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

[Sign in](#)

[Submit an Article](#)

This is the second article in my series on Transf [first article](#), we learned about the functionality of how they are used, their high-level architecture, advantages.

In this article, we can now look under the hood how they work in detail. We'll see how data flow system with their actual matrix representations understand the computations performed at each

Here's a quick summary of the previous and following the series. My goal throughout will be to understand something works but why it works that way.

1. Overview of functionality **** (How Transformers are used, and why they are better than RNNs. Components of the architecture, and behavior during Training and Inference)
2. **How it works – this article** (Internal operation end-to-end. How data flows and what computations are performed, including matrix representations)
3. Multi-head Attention (Inner workings of the Attention module throughout the Transformer) LATEST
4. Why Attention Boosts Performance (Not just does but why it works so well. How does Attention relationships between words in a sentence) EDITOR'S PICKS
DEEP DIVES
NEWSLETTER

And if you're interested in NLP applications in general, here are some other articles you might like.

WRITE FOR TDS

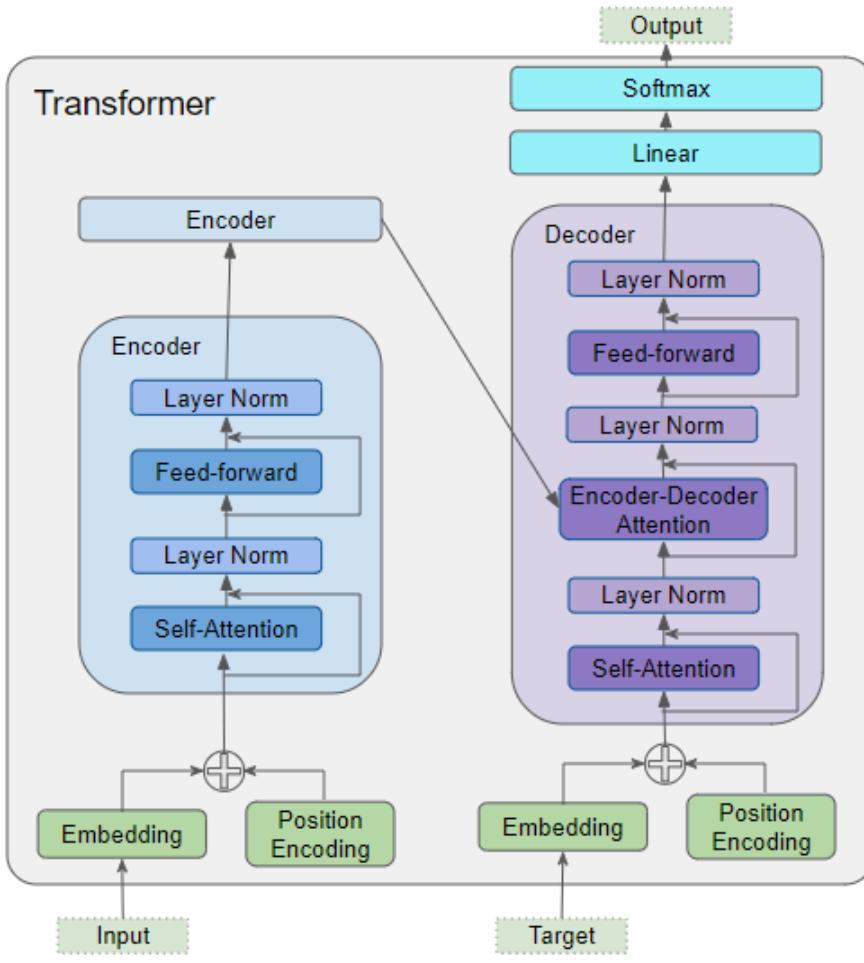
Sign in

1. Beam Search (Algorithm commonly used by NLP models to enhance predictions)
2. Bleu Score (Bleu Score and Word Error Rate metrics for NLP models)

Submit an Article

Architecture Overview

As we saw in Part 1, the main components of the Transformer are:



(Image by Author)

LATEST
EDITOR'S PICKS
DEEP DIVES
NEWSLETTER

WRITE FOR TDS

[Sign in](#)

[Submit an Article](#)

Data inputs for both the Encoder and Decoder, \

- Embedding layer
- Position Encoding layer

The Encoder stack contains a number of Encoder layers, each of which contains:

- Multi-Head Attention layer
- Feed-forward layer

The Decoder stack contains a number of Decoder layers, each of which contains:

- Two Multi-Head Attention layers
- Feed-forward layer

Output (top right) – generates the final output, and contains:

- Linear layer
- Softmax layer.

To understand what each component does, let's work through the Transformer while we are training a translation problem. We'll use one sample of our which consists of an input sequence ('You are w English) and a target sequence ('De nada' in Spa

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

Embedding and Position Encoding

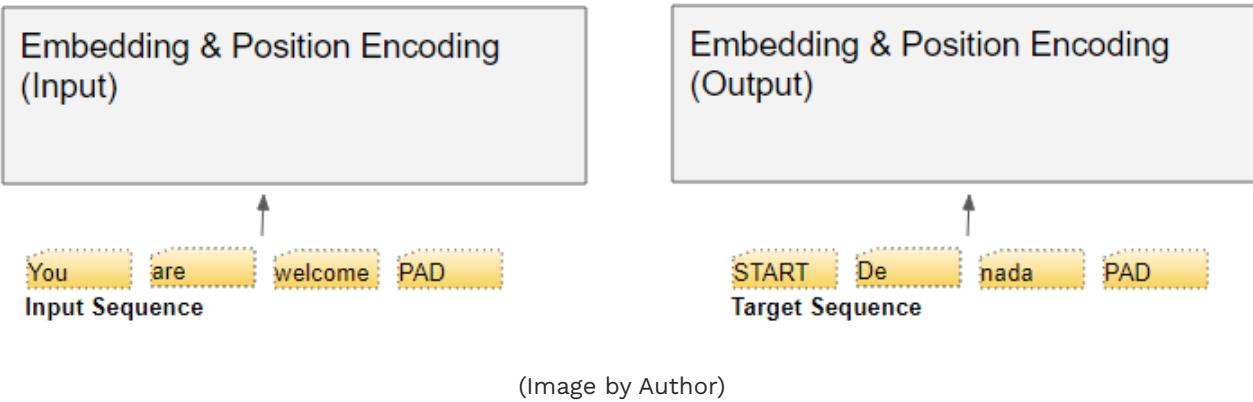
Like any NLP model, the Transformer needs two things for each word – the meaning of the word and its position in the sequence.

- The Embedding layer encodes the meaning of the word.
- The Position Encoding layer represents the position of the word.

The Transformer combines these two encodings to produce the final output.

Embedding

The Transformer has two Embedding layers. The input sequence is fed to the first Embedding layer, known as the Word Embedding.



(Image by Author)

The target sequence is fed to the second Embedding layer. This layer performs two main operations: shifting the targets right by one position and inserting a special `START` token in the first position. Note that, during Inference, if there is no target sequence and we feed the output sequence back to the second layer in a loop, as we learned in [Part 1](#). This is called the Output Embedding.

LATEST

EDITOR'S PICKS

DEEP DIVES

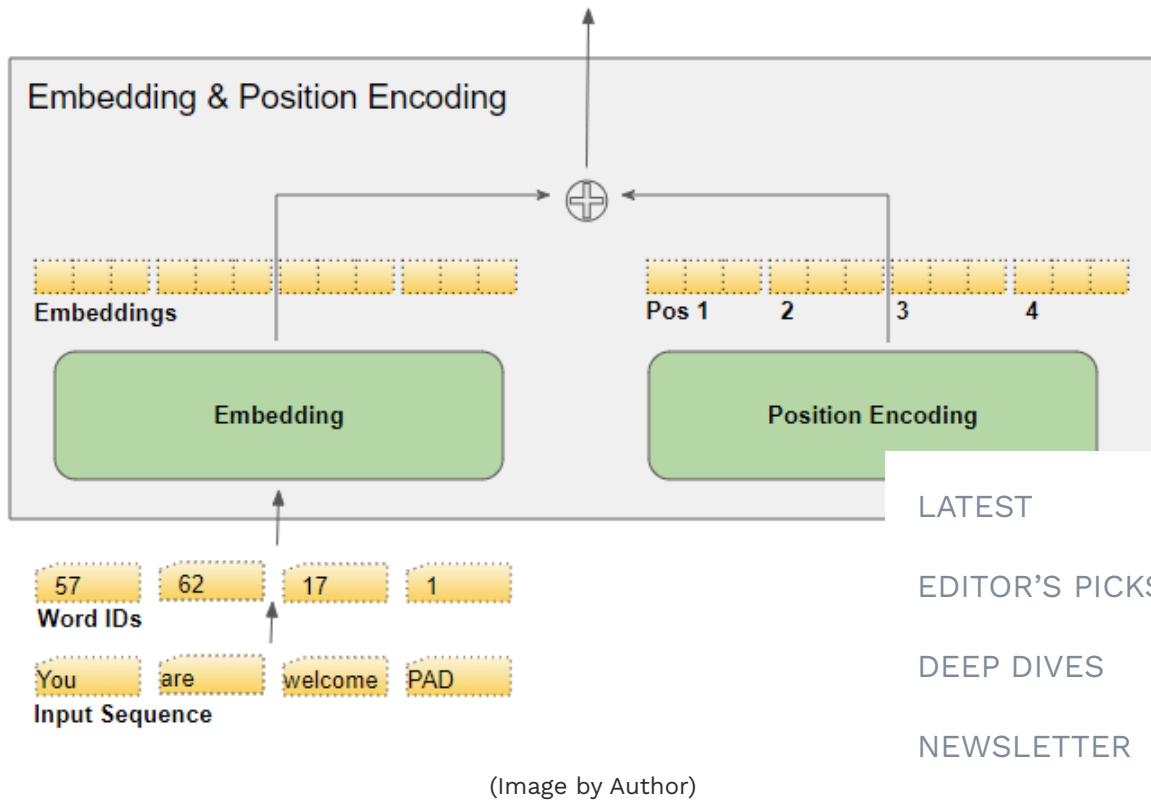
NEWSLETTER

WRITE FOR TDS

[Sign in](#)

[Submit an Article](#)

The text sequence is mapped to numeric word IDs based on the vocabulary. The embedding layer then maps each word ID to an embedding vector, which is a richer representation of the meaning of that word.



(Image by Author)

WRITE FOR TDS

Sign in

Submit an Article

Position Encoding

Since an RNN implements a loop where each word is processed sequentially, it implicitly knows the position of each word.

However, Transformers don't use RNNs and all words in a sequence are input in parallel. This is its major advantage over the RNN architecture, but it means that the position information is lost, and has to be added back in separately.

Just like the two Embedding layers, there are two Position Encoding layers. The Position Encoding is computed independently of the input sequence. These are designed to depend only on the max length of the sequence.

- the first item is a constant code that indicates the position

- the second item is a constant code that indicates the second position,
- and so on.

These constants are computed using the formula below, where

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

- pos is the position of the word in the sequence
- d_{model} is the length of the encoding vector (embedding vector) and
- i is the index value into this vector.

LATEST

EDITOR'S PICKS

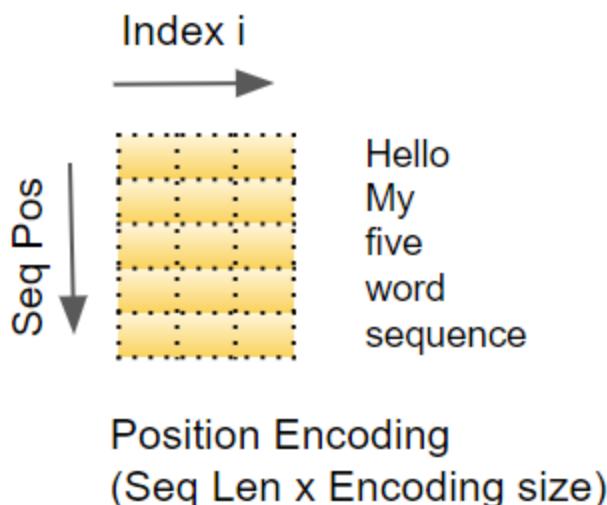
DEEP DIVES

NEWSLETTER

WRITE FOR TDS

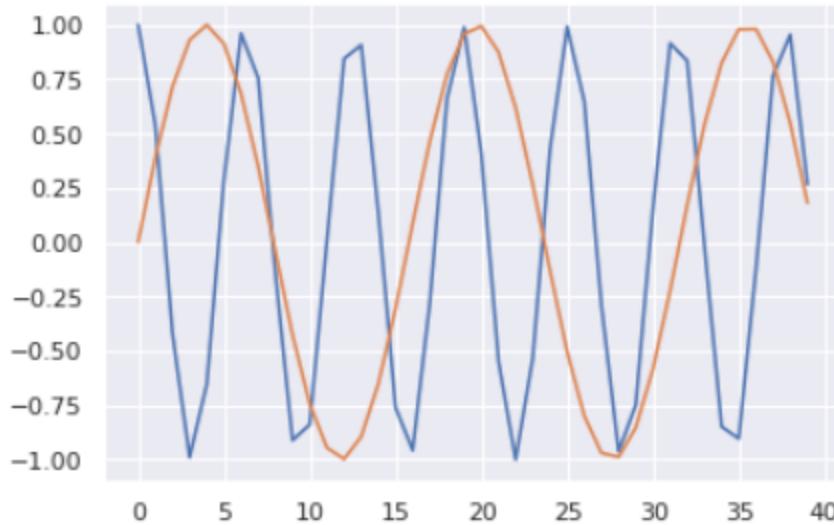
Sign in

Submit an Article



(Image by Author)

In other words, it interleaves a sine curve and a cosine curve for all even indexes and cosine values for odd indexes. As an example, if we encode a sequence, we can see below the encoding values for a few (word, encoding_index) combinations.



(Image by Author)

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

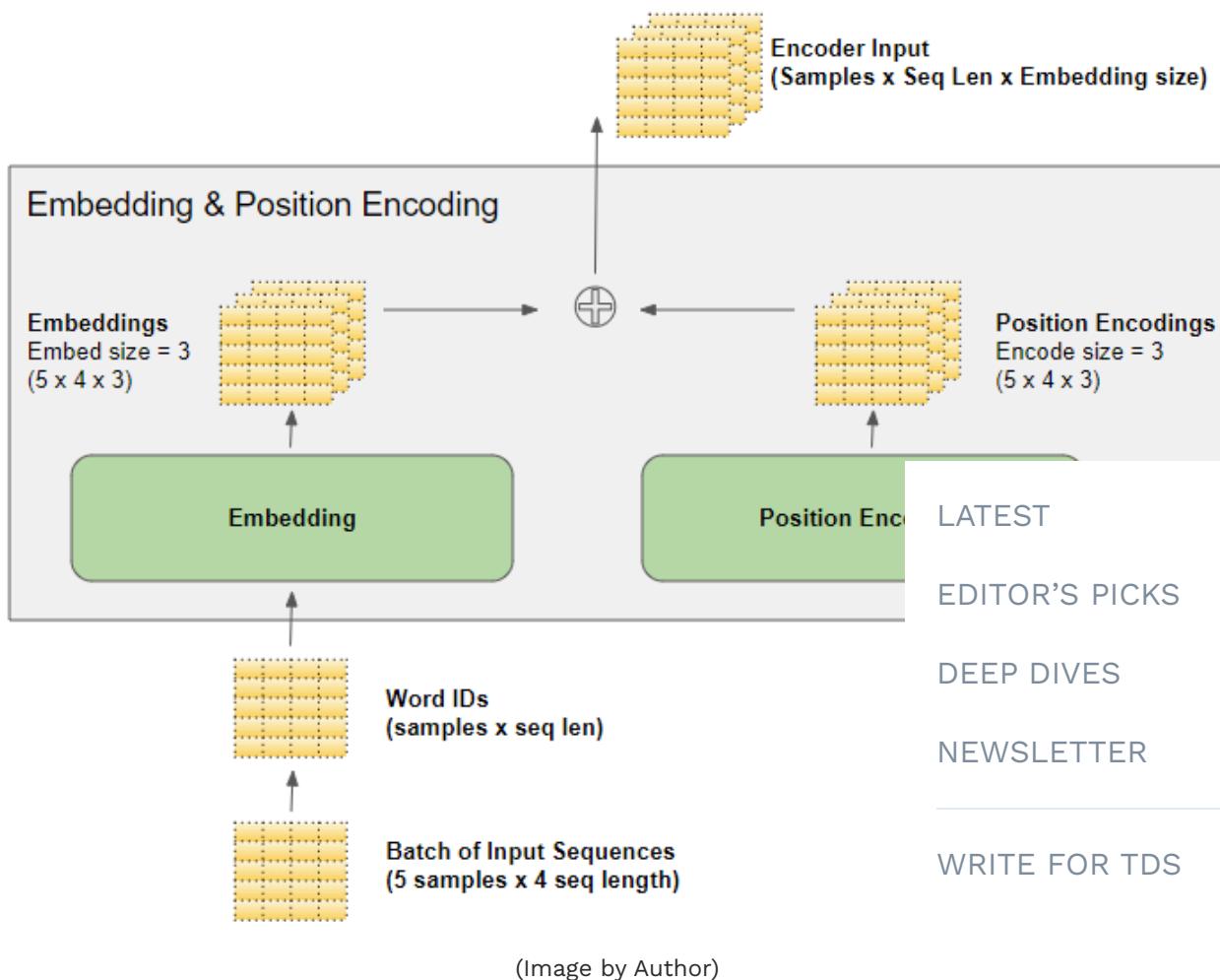
WRITE FOR TDS

Sign in

Matrix Dimensions

Submit an Article

As we know, deep learning models process a batch of samples at a time. The Embedding and Position Encoding operate on matrices representing a batch of sequences. The Embedding takes a (batch_size, sequence_length) shaped matrix of word IDs. It encodes each word ID into a vector whose length is the embedding size, resulting in a (batch_size, sequence_length, embedding_size) shaped output matrix. Position Encoding uses an encoding size that is equal to the embedding size. So it produces a similarly shaped matrix that can be added to the embedding matrix.



(Image by Author)

LATEST
EDITOR'S PICKS
DEEP DIVES
NEWSLETTER

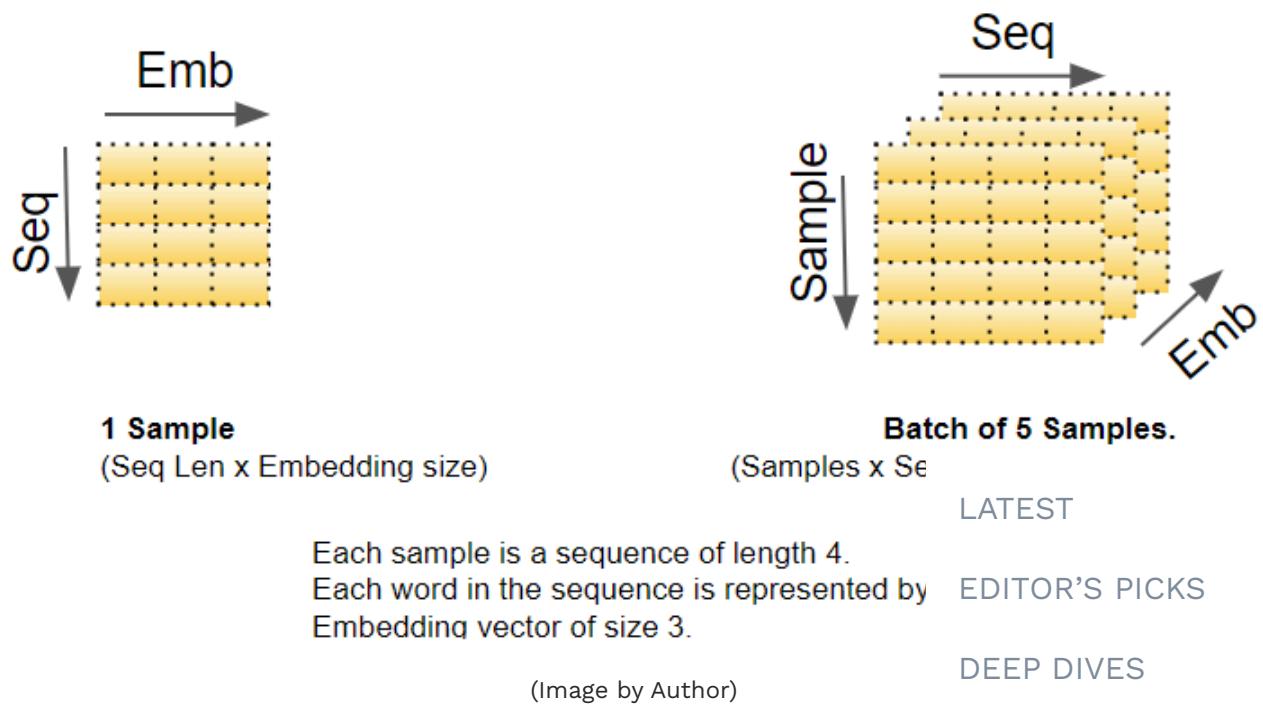
WRITE FOR TDS

Sign in

Submit an Article

The (samples, sequence length, embedding size) is maintained by the Embedding and Position Encoding layers through the Transformer, as the data flows through the Encoder and Decoder Stacks until it is reshaped by the final layers.

This gives a sense of the 3D matrix dimensions involved in the Transformer. However, to simplify the visualization, we will drop the first dimension (for the sample representation for a single sample).



The Input Embedding sends its outputs into the
Similarly, the Output Embedding feeds into the

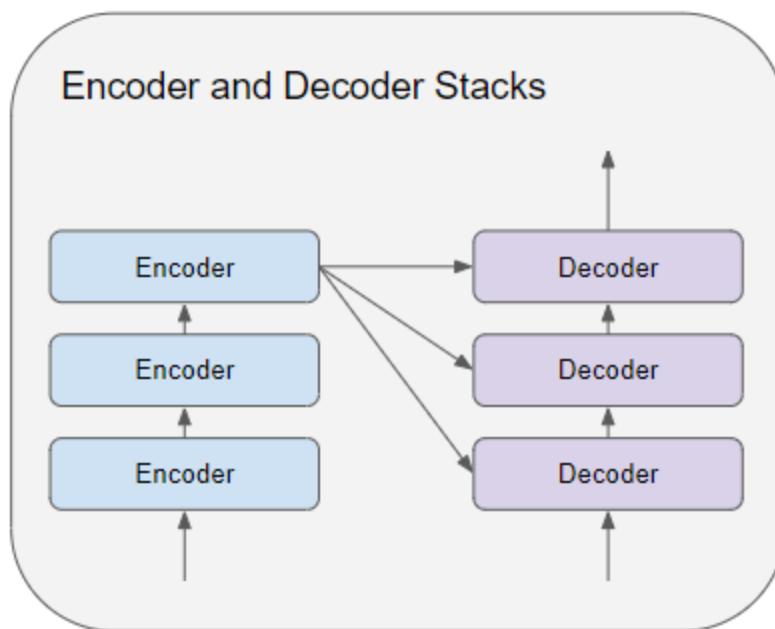
WRITE FOR TDS

Sign in

Submit an Article

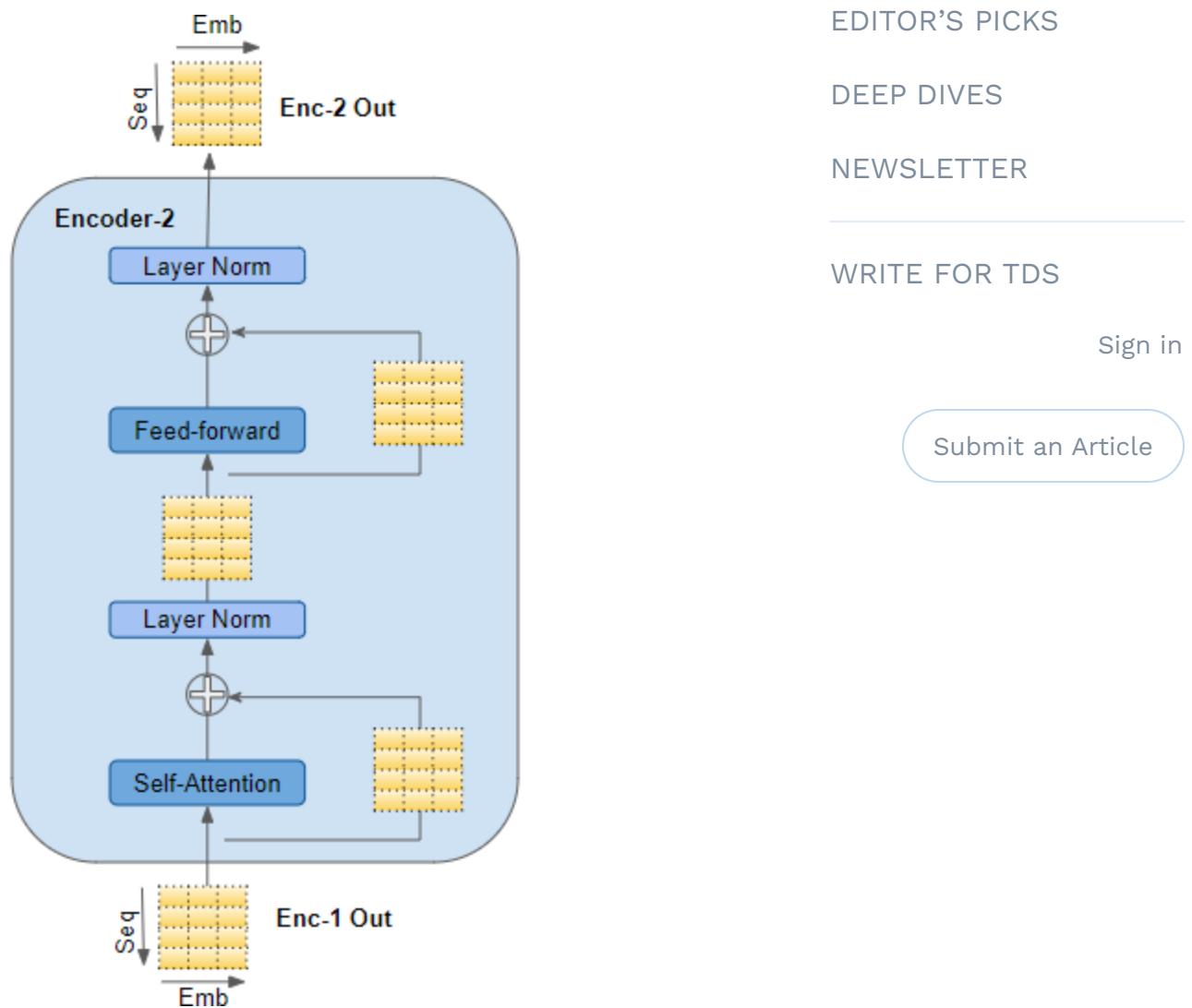
Encoder

The Encoder and Decoder Stacks consists of several
Encoders and Decoders respectively, connected



The first Encoder in the stack receives its input from the Embedding and Position Encoding. The other Encoders in the stack receive their input from the previous Encoder.

The Encoder passes its input into a Multi-head Self-attention layer. The Self-attention output is passed into a Feed-forward layer, which then sends its output upwards to the next Encoder.



Both the Self-attention and Feed-forward sub-layers include a residual skip-connection around them, followed by a Layer Norm layer.

Normalization.

The output of the last Encoder is fed into each Decoder in the Decoder Stack as explained below.

Decoder

The Decoder's structure is very similar to the Encoder's but with a couple of differences.

Like the Encoder, the first Decoder in the stack from the Output Embedding and Position Encoding. Decoders in the stack receive their input from the previous Decoder.

The Decoder passes its input into a Multi-head layer. This operates in a slightly different way than the Encoder. It is only allowed to attend to earlier positions in the sequence. This is done by masking future positions, which we will talk about shortly.

LATEST

EDITOR'S PICKS

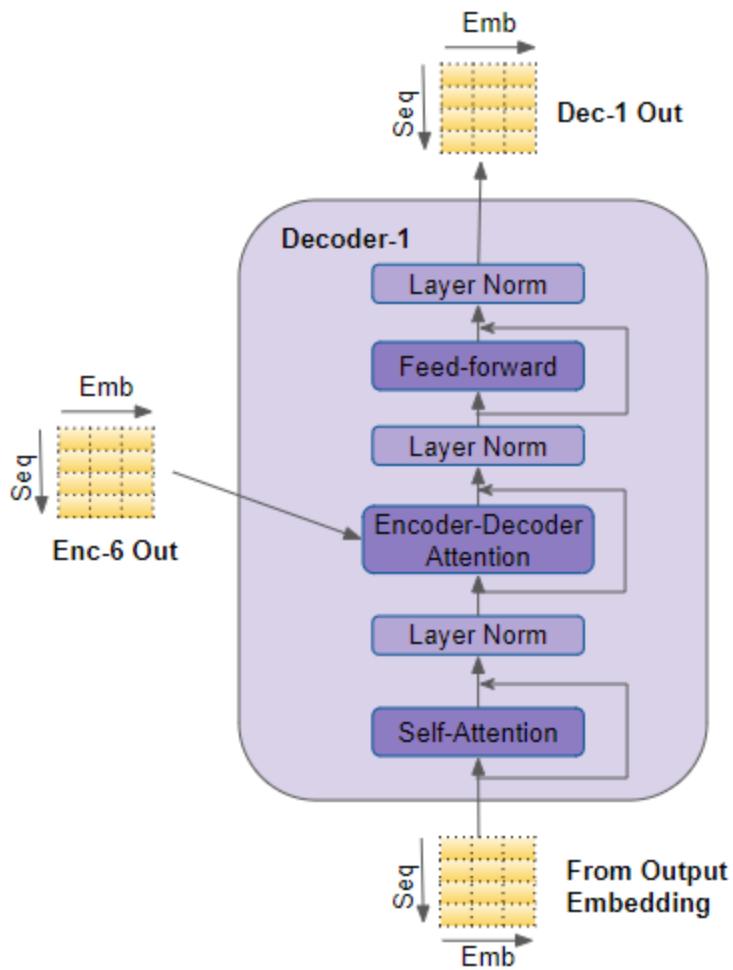
DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article



LATEST
EDITOR'S PICKS
DEEP DIVES
NEWSLETTER

WRITE FOR TDS

Sign in

[Submit an Article](#)

(Image by Author)

Unlike the Encoder, the Decoder has a second **N** attention layer, known as the Encoder-Decoder : The Encoder-Decoder attention layer works like except that it combines two sources of inputs – attention layer below it as well as the output of stack.

The Self-attention output is passed into a Feed-forward layer, which then sends its output upwards to the next

Each of these sub-layers, Self-attention, Encoder-Decoder attention, and Feed-forward, have a residual skip connection around them, followed by a Layer-Normalization

Attention

In [Part 1](#), we talked about why Attention is so important while processing sequences. In the Transformer, Attention is used in three places:

- Self-attention in the Encoder – the input sequence pays attention to itself
- Self-attention in the Decoder – the target sequence pays attention to itself
- Encoder-Decoder-attention in the Decoder – the target sequence pays attention to the input sequence

The Attention layer takes its input in the form of three parameters, known as the Query, Key, and Value

- In the Encoder's Self-attention, the Encoder uses all three parameters, Query, Key, and Value

LATEST

EDITOR'S PICKS

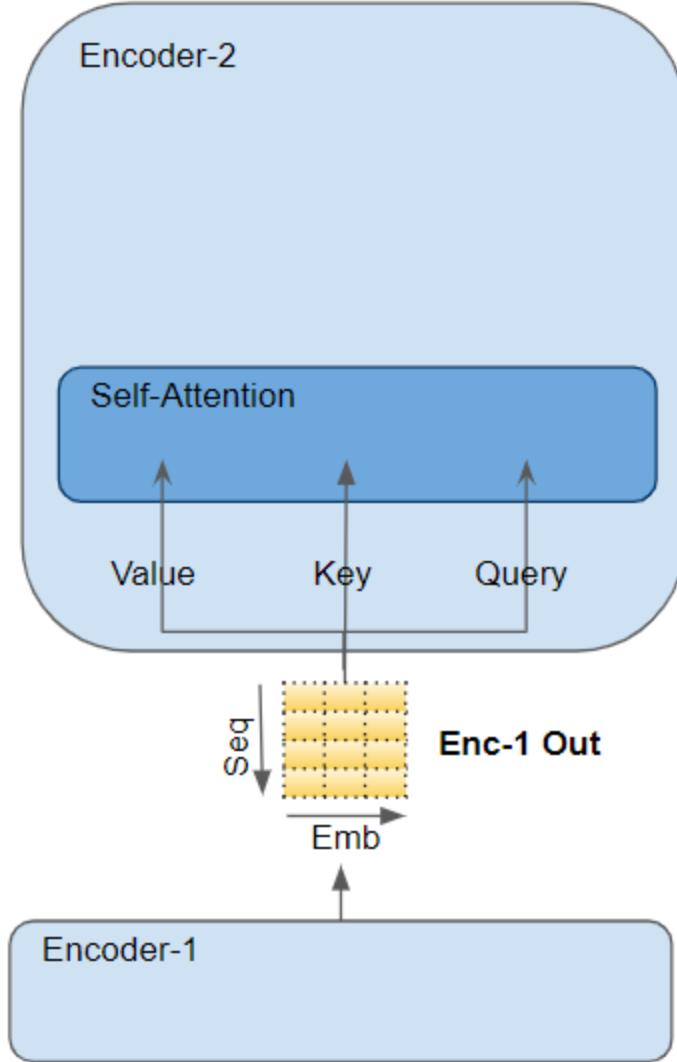
DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article



LATEST
EDITOR'S PICKS
DEEP DIVES
NEWSLETTER

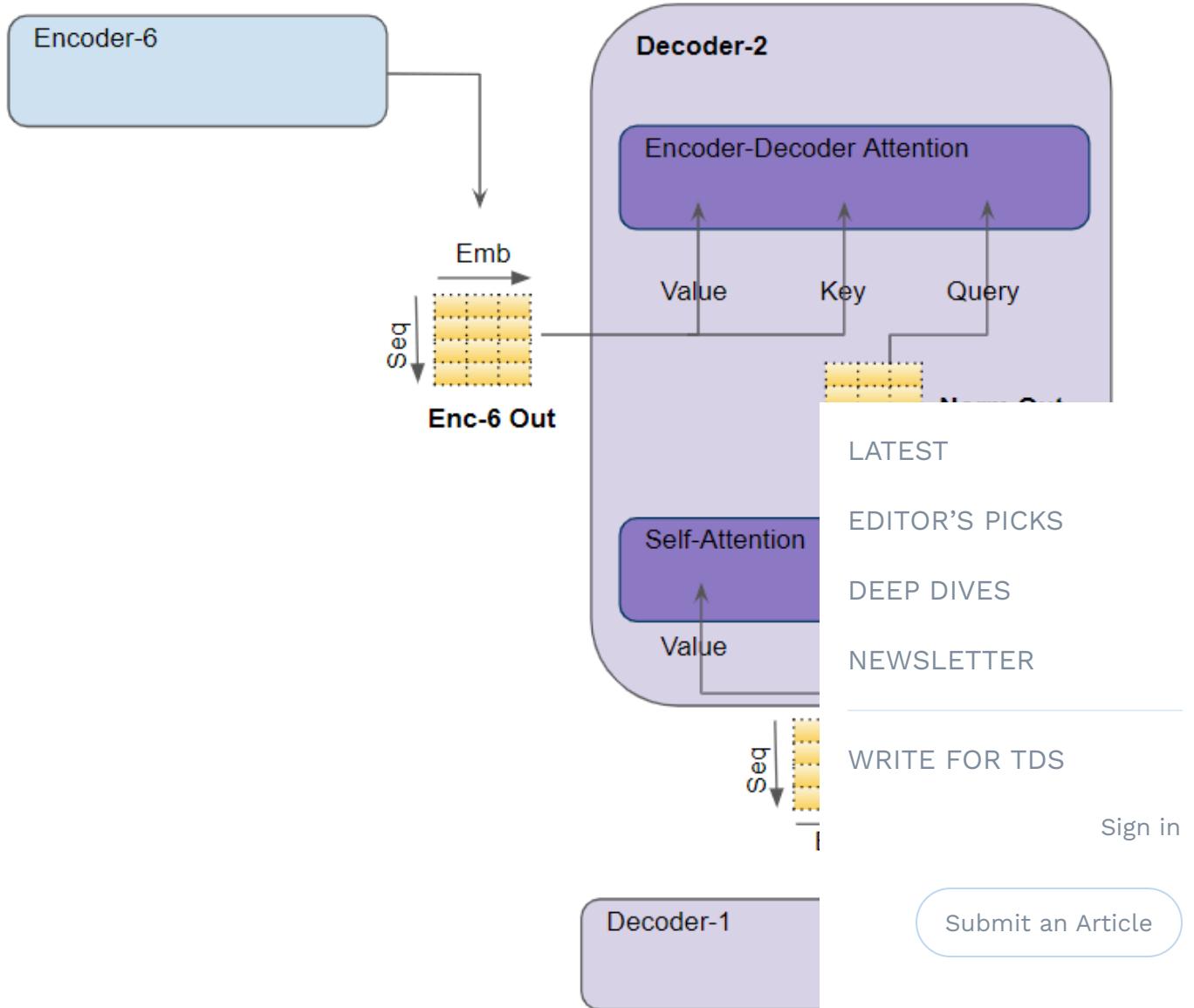
WRITE FOR TDS

[Sign in](#)

[Submit an Article](#)

(Image by Author)

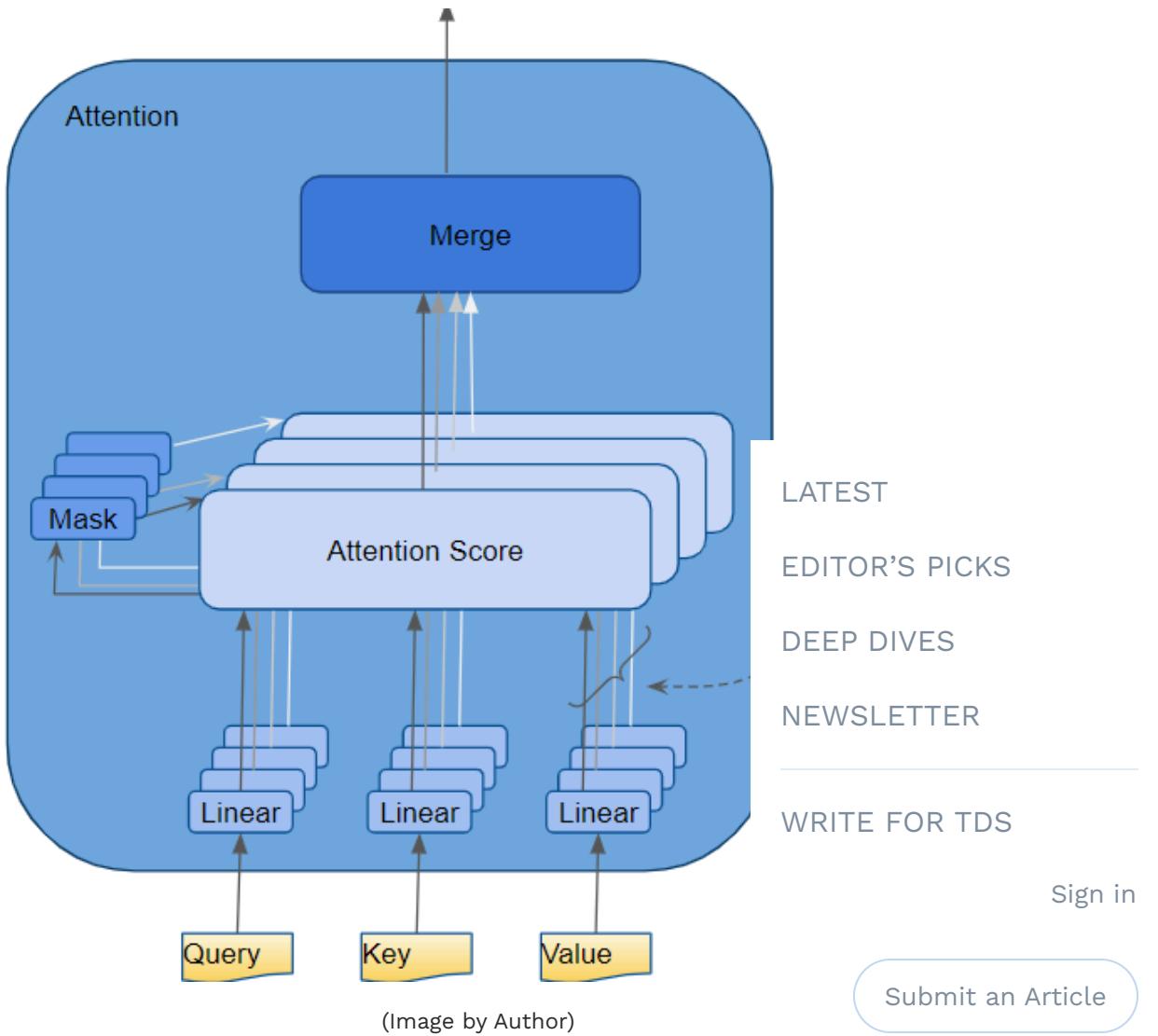
- In the Decoder's Self-attention, the Decoder is passed to all three parameters, Query, Key, and Value.
- In the Decoder's Encoder-Decoder attention, the final Encoder in the stack is passed to the three parameters. The output of the Self-attention (with Layer Norm) module below it is passed to the Query parameter.



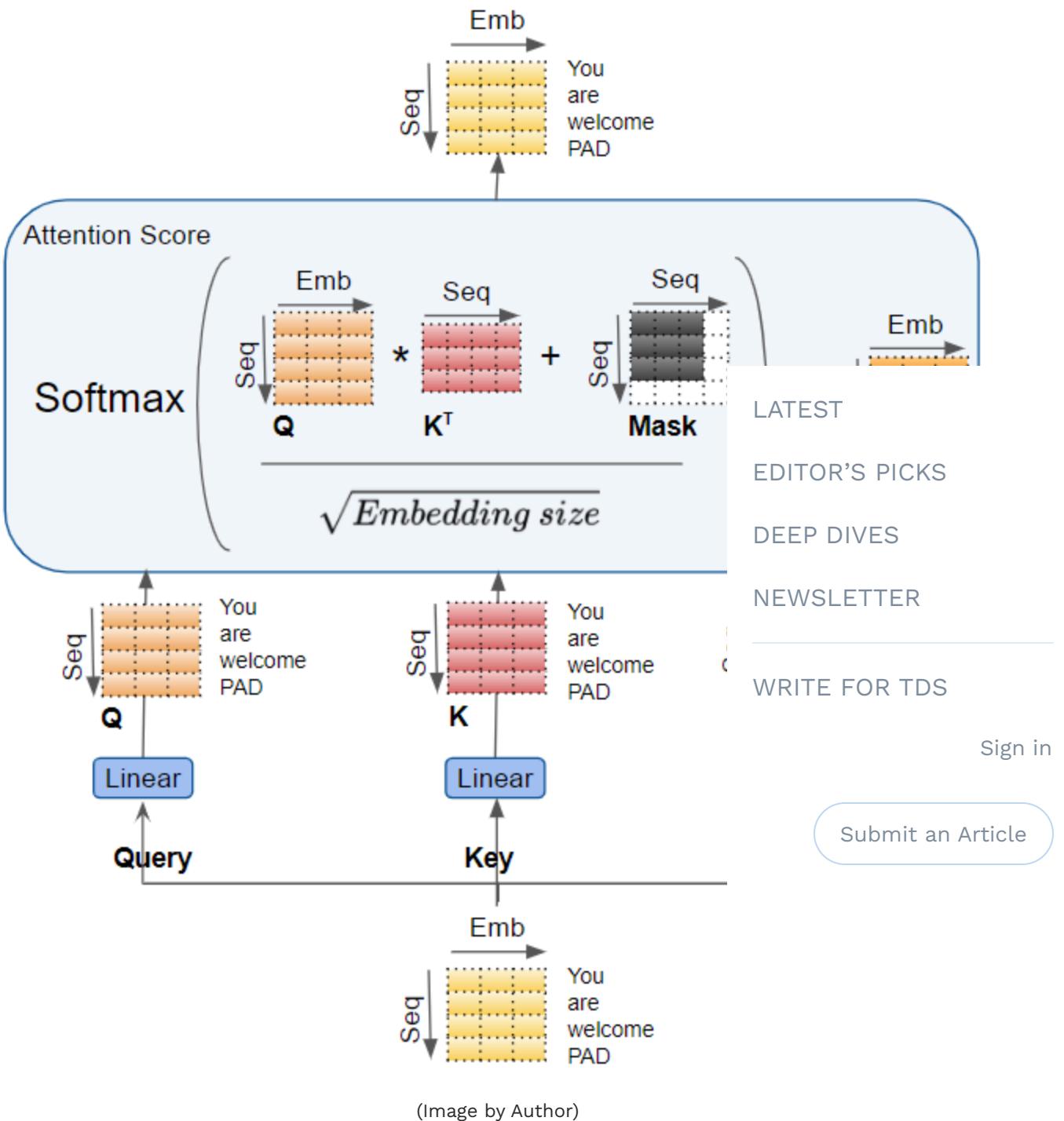
(Image by Author)

Mult-head Attention

The Transformer calls each Attention processor and repeats it several times in parallel. This is k head attention. It gives its Attention greater power discrimination, by combining several similar Attention calculations.



The Query, Key, and Value are each passed through Linear layers, each with their own weights, producing results called Q, K, and V respectively. These are then together using the Attention formula as shown to produce the Attention Score.



The important thing to realize here is that the Q carry an encoded representation of each word in the sequence. The Attention calculations then combine each word with every other word in the sequence, so that the Attention mechanism can assign a score for each word in the sequence.

When discussing the Decoder a little while back, we briefly mentioned masking. The Mask is also shown in the Attention diagrams above. Let's see how it works.

Attention Masks

While computing the Attention Score, the Attention module implements a masking step. Masking serves two purposes:

LATEST

EDITOR'S PICKS

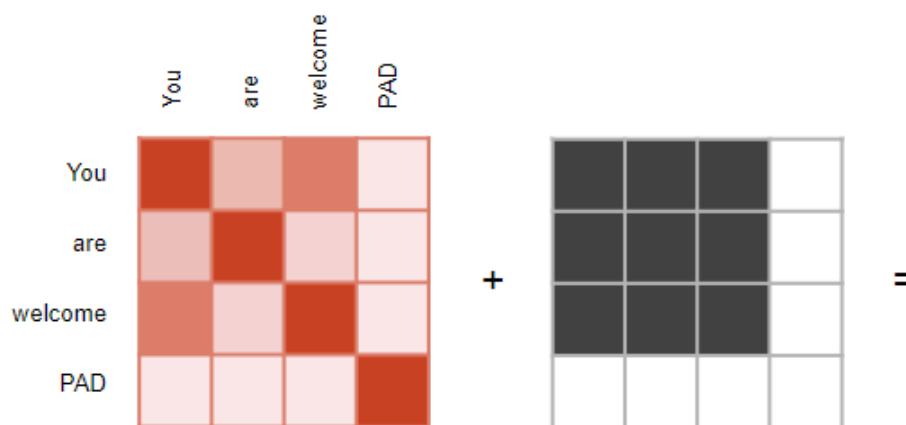
DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

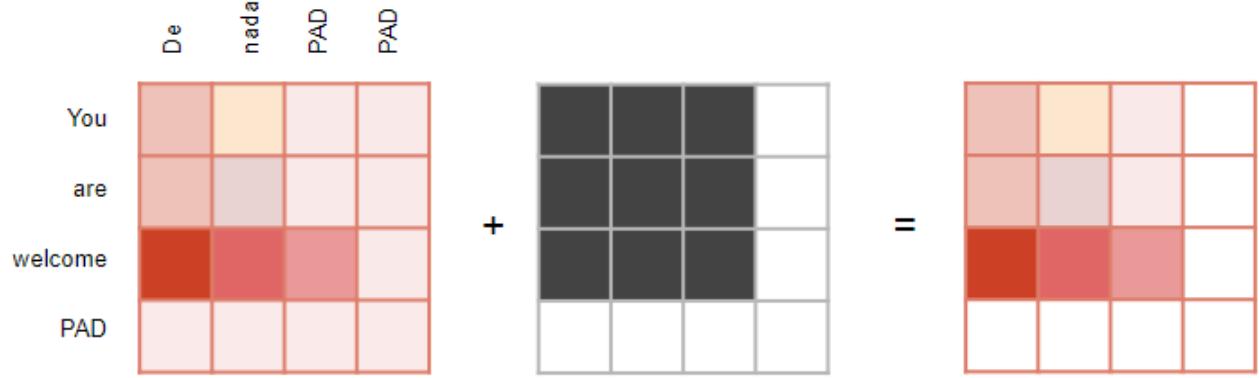
Submit an Article



Encoder Self-Attention Scores

(Image by Author)

Similarly for the Encoder-Decoder attention.



Encoder-Decoder Attention Scores

(Image by Author)

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

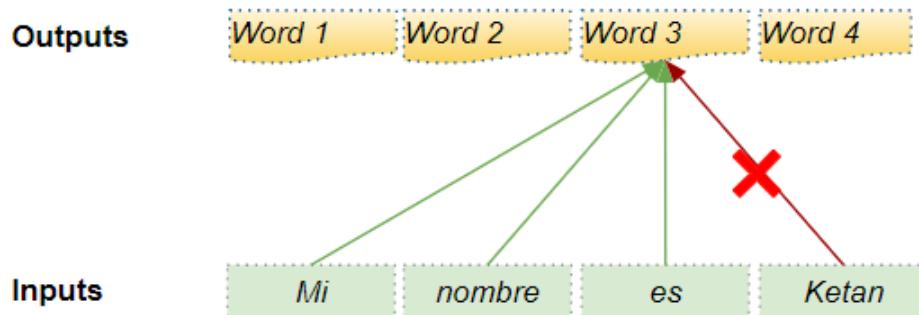
Sign in

[Submit an Article](#)

In the Decoder Self-attention: masking serves to prevent the decoder from ‘peeking’ ahead at the rest of the sequence when predicting the next word.

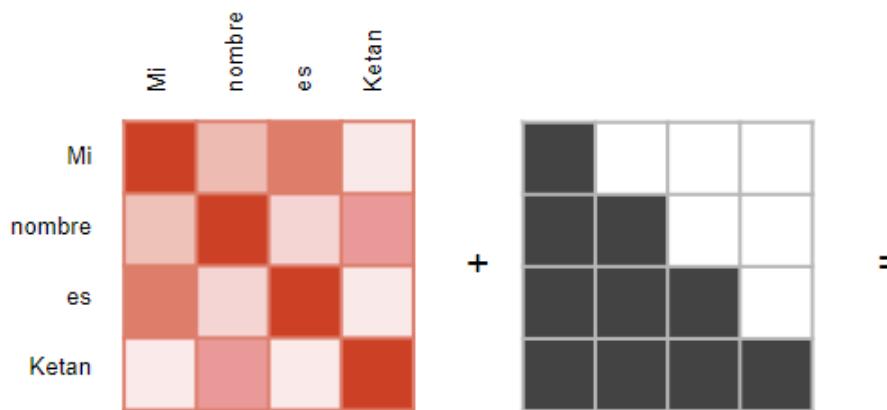
The Decoder processes words in the source sequence to predict the words in the destination sequence. During training, this is done via Teacher Forcing, where the target sequence is fed as Decoder inputs. Therefore, when predicting a word at a certain position, the Decoder attends to it the target words preceding that word as well as the words following that word. This allows the Decoder to predict using target words from future ‘time steps’.

For instance, when predicting ‘Word 3’, the Decoder attends only to the first 3 input words from the target sequence, i.e., the word ‘Ketan’.



(Image by Author)

Therefore, the Decoder masks out input words $t \in \text{LATEST}$
in the sequence.



Decoder Self-Attention Scores

(Image by Author)

When calculating the Attention Score (refer to the section showing the calculations) masking is applied to just before the Softmax. The masked out elements (e.g., squares) are set to negative infinity, so that Softmax values to zero.

Generate Output

The last Decoder in the stack passes its output component which converts it into the final output.

The Linear layer projects the Decoder vector into Word Scores, with a score value for each unique word in the target vocabulary, at each position in the sentence. For instance, if our final output sentence has 7 words and the target Spanish vocabulary has 10000 unique words, we generate 10000 score values for each of those 7 words. The score values indicate the likelihood of occurrence for each word in the vocabulary in that position of the sentence.

LATEST

The Softmax layer then turns those scores into (which add up to 1.0). In each position, we find the word with the highest probability, and then map corresponding word in the vocabulary. Those words form the output sequence of the Transformer.

EDITOR'S PICKS

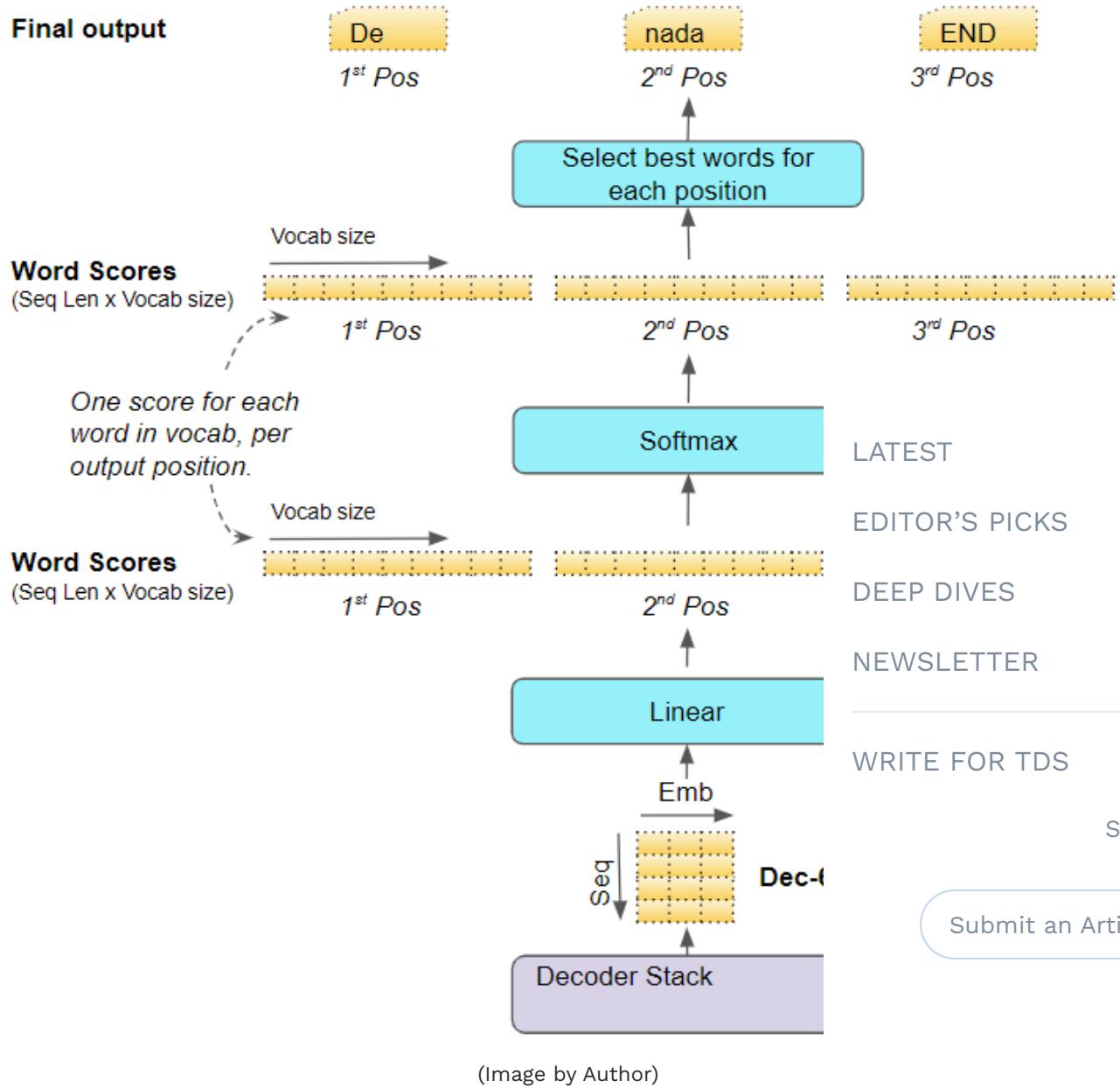
DEEP DIVES

NEWSLETTER

WRITE FOR TDS

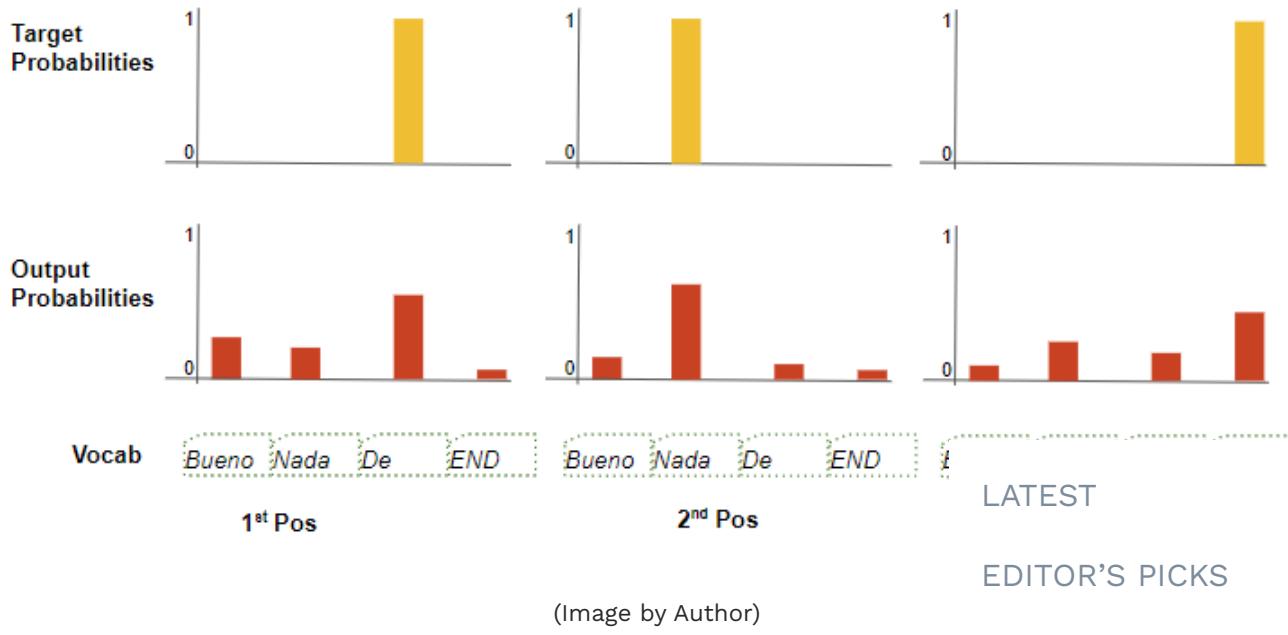
[Sign in](#)

[Submit an Article](#)



Training and Loss Function

During training, we use a loss function such as cross-entropy to compare the generated output probability distribution to the target sequence. The probability distribution gives the probability of each word occurring in that position.



Let's assume our target vocabulary contains just the words "Bueno", "Nada", "De", and "END". Our goal is to produce a probability distribution that matches the expected target sequence "De nada END".

[NEWSLETTER](#)

[WRITE FOR TDS](#)

[Sign in](#)

This means that the probability distribution for the first word-position should have a probability of 1 for "De" and 0 for all other words in the vocabulary being 0. Similarly, the second word-position should have a probability of 1 for "nada" and 0 for all other words. The third word-position should have a probability of 1 for "END" and 0 for all other words.

[Submit an Article](#)

As usual, the loss is used to compute gradients and backpropagate through the Transformer via backpropagation.

Conclusion

Hopefully, this gives you a feel for what goes on under the hood of the Transformer during Training. As we discussed in our [previous article](#), it runs in a loop during Inference but the underlying processing remains the same.

The Multi-head Attention module is what gives the Transformer its power. In the [next article](#), we will continue our journey and go one step deeper to really understand the details of how Attention is computed.

And finally, if you liked this article, you might also enjoy my other series on Audio Deep Learning, Geolocation Machine Learning, and Image Caption architectures.

LATEST

[**Audio Deep Learning Made Simple \(Part 1\): St Techniques**](#)

EDITOR'S PICKS

DEEP DIVES

[**Leveraging Geolocation Data for Machine Lea Techniques**](#)

NEWSLETTER

WRITE FOR TDS

[**Image Captions with Deep Learning: State-of Architectures**](#)

Sign in

[Submit an Article](#)

Let's keep learning!

• • •

WRITTEN BY

Ketan Doshi

[See all from Ketan Doshi](#)

Artificial Intelligence

Data Science

Deep Learn

[Machine Learning](#)[NLP](#)**Share This Article**

Towards Data Science is a community publication. Submit your insights to reach our global audience and earn through the TDS Author Payment Program!

[LATEST](#)[Write for TDS](#)[EDITOR'S PICKS](#)[DEEP DIVES](#)[NEWSLETTER](#)[WRITE FOR TDS](#)[Sign in](#)[Submit an Article](#)[ARTIFICIAL INTELLIGENCE](#)

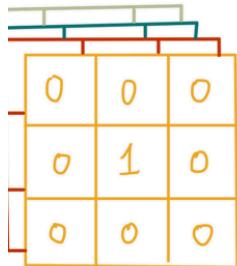
What Do Large Language Models “Understand”?

A deep dive on the understanding of LLMs

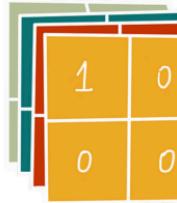
Tarik Dzekman

August 21, 2024

3



max
pooling

[ARTIFICIAL INTELLIGENCE](#)

Implementing Convolutional Neural Networks in TensorFlow

Step-by-step code guide to building a Convolutional Neural Network

Shreya Rao

August 20, 2024 6 min read



ARTIFICIAL INTELLIGENCE

How to Forecast Hierarchical Time Series

A beginner's guide to forecast reconciliation

Dr. Robert Kübler

August 20, 2024 13 min read



DATA SCIENCE

Hands-on Time Series Detection using Python

LATEST

Here's how to use detect signals with lines of...

EDITOR'S PICKS

Piero Paialunga

August 21, 2024 1

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article



MACHINE LEARNING

3 AI Use Cases (That Are Not a Chatbot)

Feature engineering, structuring unstructured data, and lead scoring

Shaw Talebi

August 21, 2024 7 min read

DATA SCIENCE

Solving a Complex Scheduling Problem with Quantum Annealing

Solving the resource scheduling problem with Wave's hybrid co-optimization model (CQM)

Luis Fernando PÉREZ

August 20, 2024 1



DATA SCIENCE

Back To Basics, Part Uno: Linear Regression and Cost Function

An illustrated guide on essential machine learning concepts

Shreya Rao

February 3, 2023 6 min read

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

[Sign in](#)



[Submit an Article](#)

Your home for data science and AI. The world's leading publication in analytics, data engineering, machine learning, and artificial intelligence.

© Insight Media Group, LLC 2025

[Subscribe to Our Newsletter](#)

[WRITE FOR TDS](#) · [ABOUT](#) · [ADVERTISE](#) · [PRIVACY POLICY](#)

[COOKIES SETTINGS](#)