



Transformer in NYC (created from photofunia)

## Solving Transformer by Hand: A Step-by-Step Math Example



Fareed Khan · Follow

Published in Level Up Coding · 13 min read · Dec 18, 2023

2.8K

38



I have already written a detailed blog on how transformers work using a very small sample of the dataset, which will be my best blog ever because it has elevated my profile and given me the motivation to write more. However, that blog is incomplete as it only covers 20% of the transformer architecture and contains numerous calculation errors, as pointed out by readers. After a considerable amount of time has passed since that blog, I will be revisiting the topic in this new blog.

My previous blog on transformer architecture (covers only 20%):



I plan to explain the transformer again in the same manner as I did in my previous blog (for both coders and non-coders), providing a complete guide with a step-by-step approach to understanding how they work.

### Table of Contents

- [Defining our Dataset](#)
- [Finding Vocab Size](#)
- [Encoding](#)
- [Calculating Embedding](#)
- [Calculating Positional Embedding](#)
- [Concatenating Positional and Word Embeddings](#)
- [Multi Head Attention](#)
- [Adding and Normalizing](#)
- [Feed Forward Network](#)
- [Adding and Normalizing Again](#)
- [Decoder Part](#)
- [Understanding Mask Multi Head Attention](#)
- [Calculating the Predicted Word](#)
- [Important Points](#)
- [Conclusion](#)

### Step 1 — Defining our Dataset

The dataset used for creating ChatGPT is 570 GB. On the other hand, for our purposes, we will be using a very small dataset to perform numerical calculations visually.

#### Dataset (corpus)

I drink and I know things.  
When you play the game of thrones, you win or you die.  
The true enemy won't wait out the storm, He brings the storm.

Our entire dataset contains only three sentences, all of which are dialogues taken from a TV show. Although our dataset is cleaned, in real-world scenarios like ChatGPT creation, cleaning a 570 GB dataset requires a significant amount of effort.

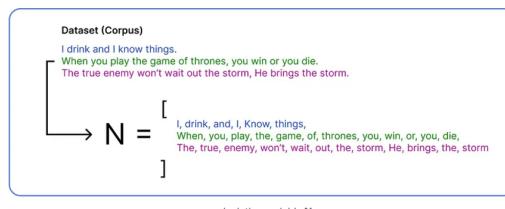
### Step 2— Finding Vocab Size

The vocabulary size determines the total number of unique words in our dataset. It can be calculated using the below formula, where N is the total number of words in our dataset.

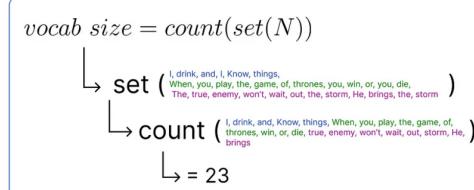
$$\text{vocab size} = \text{count}(\text{set}(N))$$

vocab\_size formula where N is total number of words

In order to find N, we need to break our dataset into individual words.



After obtaining N, we perform a set operation to remove duplicates, and then we can count the unique words to determine the vocabulary size.



Therefore, the vocabulary size is 23, as there are 23 unique words in our dataset.

### Step 3— Encoding

Now, we need to assign a unique number to each unique word.

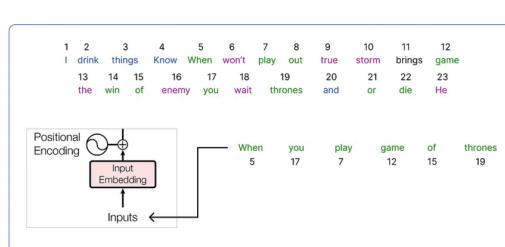


As we have considered a single token as a single word and assigned a number to it, ChatGPT has considered a portion of a word as a single token using this formula: 1 Token = 0.75 Word

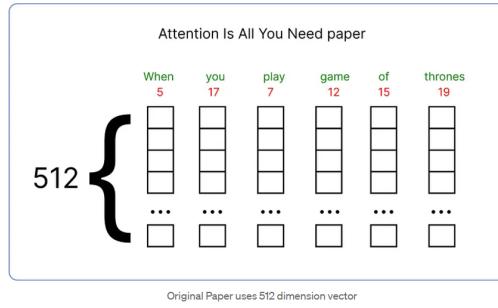
After encoding our entire dataset, it's time to select our input and start working with the transformer architecture.

### Step 4 — Calculating Embedding

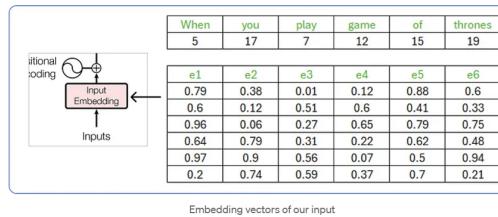
Let's select a sentence from our corpus that will be processed in our transformer architecture.



We have selected our input, and we need to find an embedding vector for it. The original paper uses a **512-dimensional embedding vector** for each input word.



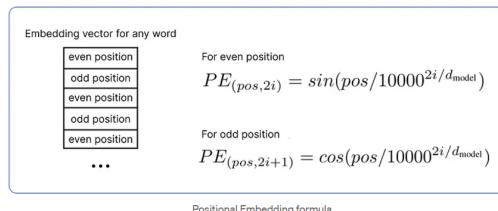
Since, for our case, we need to work with a smaller dimension of embedding vector to visualize how the calculation is taking place. So, we will be using a dimension of  $6$  for the embedding vector.



These values of the embedding vector are between 0 and 1 and are filled randomly in the beginning. They will later be updated as our transformer starts understanding the meanings among the words.

#### Step 5 — Calculating Positional Embedding

Now we need to find positional embeddings for our input. There are two formulas for positional embedding depending on the position of the  $i$ th value of that embedding vector for each word.



As you do know, our input sentence is “when you play the game of thrones” and the starting word is “when” with a starting index (POS) value is  $0$ , having a dimension ( $d$ ) of  $6$ . For  $i$  from  $0$  to  $5$ , we calculate the positional embedding for our first word of the input sentence.

When				
	5			
i	e1	Position	Formula	p1
0	0.79	Even	$\sin(0/10000^{(2*0/6)})$	0
1	0.6	Odd	$\cos(0/10000^{(2*1/6)})$	1
2	0.96	Even	$\sin(0/10000^{(2*2/6)})$	0
3	0.64	Odd	$\cos(0/10000^{(2*3/6)})$	1
4	0.97	Even	$\sin(0/10000^{(2*4/6)})$	0
5	0.2	Odd	$\cos(0/10000^{(2*5/6)})$	1

$d$  (dim)    6  
POS        0

Positional Embedding for word: When

Similarly, we can calculate positional embedding for all the words in our input sentence.

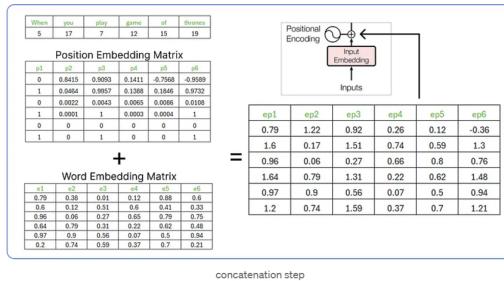


3	1	0.0001	1	0.0003	0.0004	1
4	0	0	0	0	0	0
5	1	0	1	0	0	1
d(dim)	6	6	6	6	6	6
POS	0	1	2	3	4	5

Calculating Positional Embeddings of our input (The calculated values are rounded)

### Step 6 — Concatenating Positional and Word Embeddings

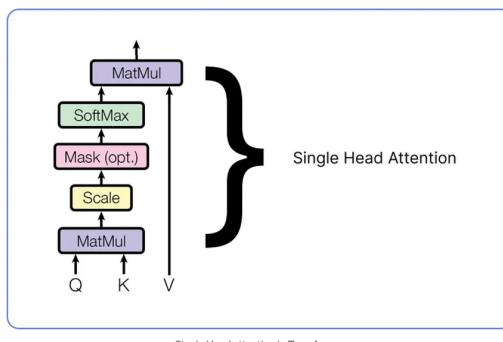
After calculating positional embedding, we need to add word embeddings and positional embeddings.



This resultant matrix from combining both matrices (Word embedding matrix and positional embedding matrix) will be considered as an input to the encoder part.

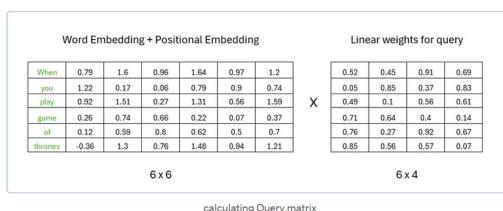
### Step 7 — Multi Head Attention

A multi-head attention is comprised of many single-head attentions. It is up to us how many single heads we need to combine. For example, LLaMA LLM from Meta has used 32 single heads in the encoder architecture. Below is the illustrated diagram of how a single-head attention looks like.

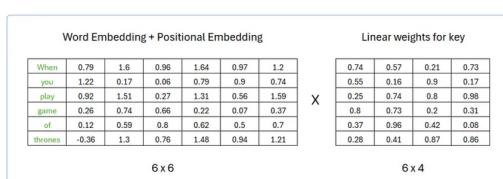


There are three inputs: query, key, and value. Each of these matrices is obtained by multiplying a different set of weights matrix from the Transpose of same matrix that we computed earlier by adding the word embedding and positional embedding matrix.

Let's say, for computing the query matrix, the set of weights matrix must have the number of rows the same as the number of columns of the transpose matrix, while the columns of the weights matrix can be any; for example, we suppose 4 columns in our weights matrix. The values in the weights matrix are between 0 and 1 randomly, which will later be updated when our transformer starts learning the meaning of these words.



Similarly, we can compute the key and value matrices using the same procedure, but the values in the weights matrix must be different for both.



Word Embedding + Positional Embedding						
When	0.79	1.6	0.96	1.64	0.97	1.2
you	1.22	0.17	0.06	0.79	0.9	0.74
play	0.92	1.51	0.27	1.31	0.56	1.59
game	0.26	0.74	0.66	0.22	0.07	0.37
of	0.12	0.59	0.8	0.62	0.5	0.7
throne	-0.36	1.3	0.76	1.48	0.94	1.21

X

6 x 6	Linear weights for value	6 x 4

Calculating Key and Value Matrices

So, after multiplying matrices, the resultant **query**, **key**, and **values** are obtained:

Query						
3.88	3.8	4.08	3.42			
2.55	1.86	2.77	1.78			
3.39	3.6	3.49	2.72			
1.02	1.18	1.24	1.3			
1.9	1.56	1.88	1.53			
3.04	2.9	2.73	2.22			

6 x 4

key						
3.71	4.04	4.15	3.41			
2.18	2.51	1.64	1.93			
3.28	3.11	3.65	3.01			
1.07	1.13	1.64	1.35			
1.49	1.97	2.14	1.81			
2.51	3.04	3.45	2.28			

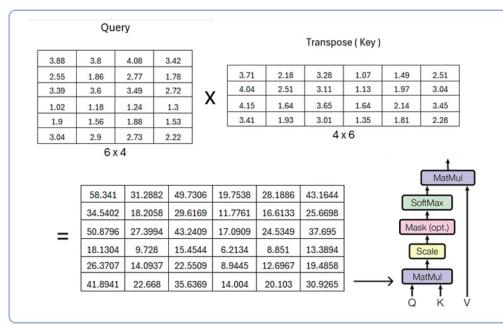
6 x 4

value						
3.88	3.8	4.08	3.42			
2.55	1.86	2.77	1.78			
3.39	3.6	3.49	2.72			
1.02	1.18	1.24	1.3			
1.9	1.56	1.88	1.53			
3.04	2.9	2.73	2.22			

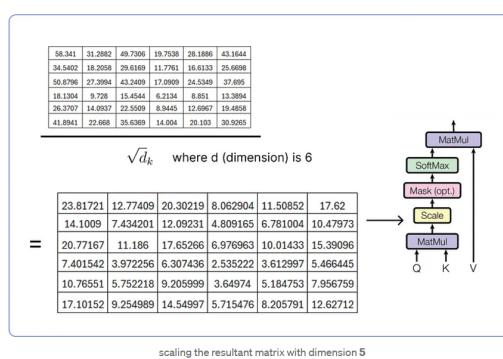
6 x 4

Query, Key, Value matrices

Now that we have all three matrices, let's start calculating single-head attention step by step.

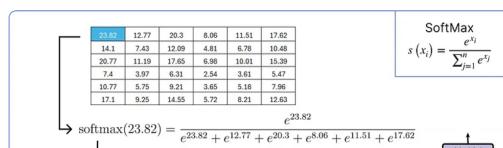


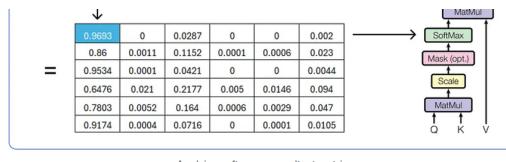
For scaling the resultant matrix, we have to reuse the dimension of our embedding vector, which is  $\sqrt{6}$ .



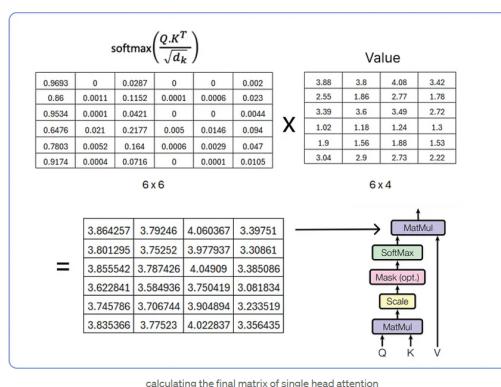
The next step of masking is optional, and we won't be calculating it. Masking is like telling the model to focus only on what's happened before a certain point and not peek into the future while figuring out the importance of different words in a sentence. It helps the model understand things in a step-by-step manner, without cheating by looking ahead.

So now we will be applying the softmax operation on our scaled resultant matrix.

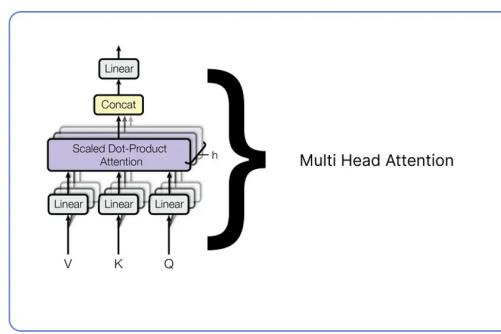




Doing the final multiplication step to obtain the resultant matrix from single-head attention.

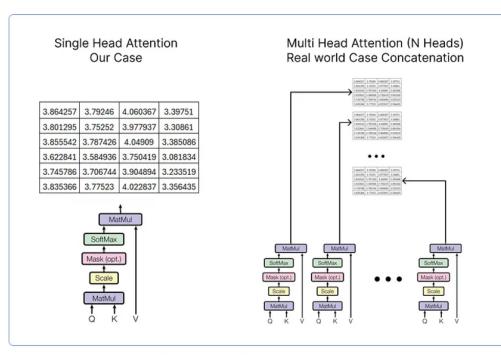


We have calculated single-head attention, while multi-head attention comprises many single-head attentions, as I stated earlier. Below is a visual of how it looks like:

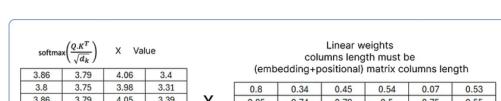


Each single-head attention has three inputs: **query**, **key**, and **value**, and each three have a different set of weights. Once all single-head attentions output their resultant matrices, they will all be concatenated, and the final concatenated matrix is once again transformed linearly by multiplying it with a set of weights matrix initialized with random values, which will later get updated when the transformer starts training.

Since, in our case, we are considering a single-head attention, but this is how it looks if we are working with multi-head attention.



In either case, whether it's single-head or multi-head attention, the resultant matrix needs to be once again transformed linearly by multiplying a set of weights matrix.



Q.WE	V.WE	Q.P.E	V.P.E
3.75	3.71	3.9	3.23
3.84	3.78	4.02	3.36

4 x 6

10.84	9.45	7.33	7.8	6.09	7.66
10.65	9.28	7.22	7.67	5.99	7.51
10.83	9.43	7.33	7.79	6.08	7.65
10.08	8.77	6.85	7.25	5.67	7.09
10.48	9.12	7.11	7.54	5.89	7.38
10.77	9.38	7.29	7.75	6.05	7.6

normalizing single head attention matrix

Make sure the linear set of weights matrix number of columns must be equal to the matrix that we computed earlier (**word embedding + positional embedding**) matrix number of columns, because the next step, we will be adding the resultant normalized matrix with (**word embedding + positional embedding**) matrix.

Output of Multi Head attention					
10.84	9.45	7.33	7.8	6.09	7.66
10.65	9.28	7.22	7.67	5.99	7.51
10.83	9.43	7.33	7.79	6.08	7.65
10.08	8.77	6.85	7.25	5.67	7.09
10.48	9.12	7.11	7.54	5.89	7.38
10.77	9.38	7.29	7.75	6.05	7.6

6 x 6

Output matrix of multi head attention

As we have computed the resultant matrix for multi-head attention, next, we will be working on adding and normalizing step.

### Step 8 — Adding and Normalizing

Once we obtain the resultant matrix from multi-head attention, we have to add it to our original matrix. Let's do it first.

Word Embedding + Positional Embedding					
When	0.79	1.6	0.96	1.64	0.67
play	0.23	0.17	0.96	0.79	0.6
game	0.92	1.51	0.27	1.31	0.56
at	0.26	0.74	0.66	0.22	0.07
browser	0.12	0.59	0.8	0.62	0.5
	-0.36	1.3	0.76	1.48	0.84

6 x 6

+

=

Output of Multi Head attention					
10.84	9.45	7.33	7.8	6.09	7.66
10.65	9.28	7.22	7.67	5.99	7.51
10.83	9.43	7.33	7.79	6.08	7.65
10.08	8.77	6.85	7.25	5.67	7.09
10.48	9.12	7.11	7.54	5.89	7.38
10.77	9.38	7.29	7.75	6.05	7.6

6 x 6

Adding matrices to perform add and norm step

To normalize the above matrix, we need to compute the mean and standard deviation row-wise for each row.

$mean = \frac{\sum_{i=1}^N X_i}{N}$	$standard\ dev = \sqrt{\frac{\sum_{i=1}^N (X_i - \mu)^2}{N}}$				
Row Wise Implementation					
11.63	11.05	8.29	9.44	7.06	8.86
11.87	9.45	7.28	8.46	6.89	8.25
11.75	10.94	7.6	9.1	6.64	9.24
10.34	9.51	7.51	7.47	5.74	7.46
10.6	9.71	7.91	8.16	6.39	8.08
10.41	10.68	8.05	9.23	6.99	8.81

Mean Standard Deviation

calculating meand and std.

we subtract each value of the matrix by the corresponding row mean and divide it by the corresponding standard deviation.

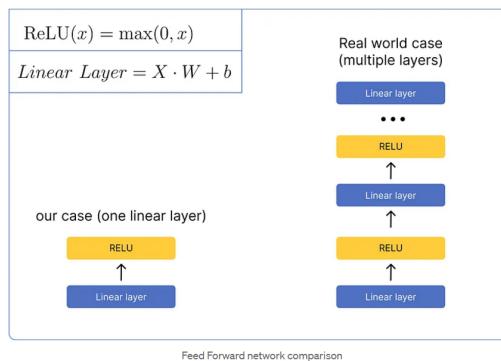
$\frac{value - mean}{std + error} = \frac{11.63 - 9.26}{1.57 + 0.0001}$	$Mean\ Std$
$1.51$	9.26 1.57
$2.02$	8.56 1.64
$1.54$	9.04 1.76
$1.64$	7.96 1.51
$1.65$	8.37 1.35
$1.16$	8.93 1.28

normalizing the resultant matrix

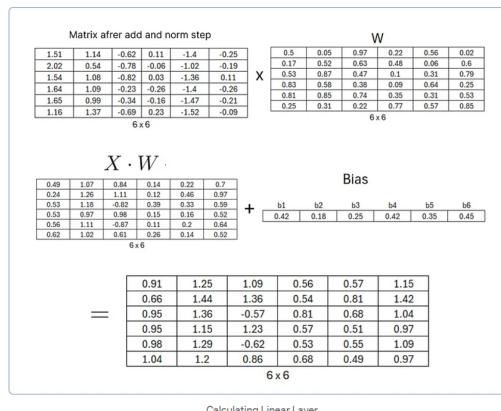
Adding a small value of error prevents the denominator from being zero and avoids making the entire term infinity.

### Step 9 — Feed Forward Network

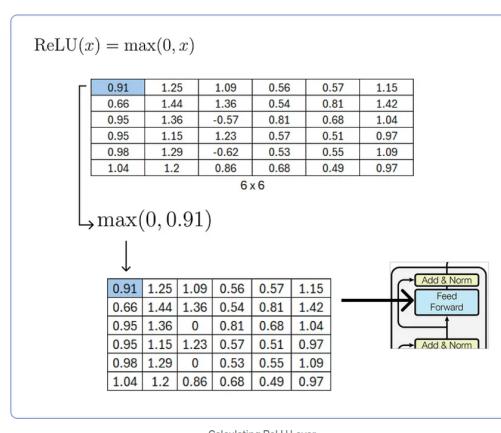
After normalizing the matrix, it will be processed through a feedforward network. We will be using a very basic network that contains only one linear layer and one ReLU activation function layer. This is how it looks like visually:



First, we need to calculate the linear layer by multiplying our last calculated matrix with a random set of weights matrix, which will be updated when the transformer starts learning, and adding the resultant matrix to a bias matrix that also contains random values.

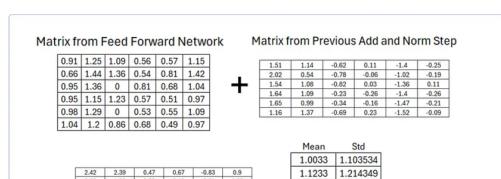


After calculating the linear layer, we need to pass it through the ReLU layer and use its formula.



### Step 10 — Adding and Normalizing Again

Once we obtain the resultant matrix from feed forward network, we have to add it to the matrix that is obtained from previous add and norm step, and then normalizing it using the row wise mean and standard deviation.



=	<table border="1"> <tr><td>2.49</td><td>2.44</td><td>-0.82</td><td>0.84</td><td>0.68</td><td>1.15</td></tr> <tr><td>2.59</td><td>2.28</td><td>1</td><td>0.31</td><td>-0.89</td><td>0.71</td></tr> <tr><td>2.63</td><td>2.28</td><td>-0.34</td><td>0.37</td><td>-0.92</td><td>0.69</td></tr> <tr><td>2.5</td><td>2.87</td><td>0.17</td><td>0.34</td><td>-1.03</td><td>0.68</td></tr> </table>	2.49	2.44	-0.82	0.84	0.68	1.15	2.59	2.28	1	0.31	-0.89	0.71	2.63	2.28	-0.34	0.37	-0.92	0.69	2.5	2.87	0.17	0.34	-1.03	0.68	→	<table border="1"> <tr><td>0.9933</td><td>1.289055</td></tr> <tr><td>0.8167</td><td>1.306016</td></tr> <tr><td>0.95</td><td>1.320773</td></tr> </table>	0.9933	1.289055	0.8167	1.306016	0.95	1.320773						
2.49	2.44	-0.82	0.84	0.68	1.15																																		
2.59	2.28	1	0.31	-0.89	0.71																																		
2.63	2.28	-0.34	0.37	-0.92	0.69																																		
2.5	2.87	0.17	0.34	-1.03	0.68																																		
0.9933	1.289055																																						
0.8167	1.306016																																						
0.95	1.320773																																						
=	<table border="1"> <tr><td>1.28</td><td>1.26</td><td>-0.48</td><td>-0.3</td><td>-1.66</td><td>-0.09</td></tr> <tr><td>1.28</td><td>0.71</td><td>-0.45</td><td>-0.53</td><td>-1.1</td><td>0.09</td></tr> <tr><td>1.22</td><td>1.18</td><td>-1.32</td><td>-0.05</td><td>-1.22</td><td>0.19</td></tr> <tr><td>1.24</td><td>0.97</td><td>0.01</td><td>-0.53</td><td>-1.46</td><td>-0.22</td></tr> <tr><td>1.39</td><td>1.12</td><td>-0.69</td><td>-0.34</td><td>-1.33</td><td>0.05</td></tr> <tr><td>0.95</td><td>1.23</td><td>-0.59</td><td>-0.03</td><td>-1.5</td><td>-0.05</td></tr> </table>	1.28	1.26	-0.48	-0.3	-1.66	-0.09	1.28	0.71	-0.45	-0.53	-1.1	0.09	1.22	1.18	-1.32	-0.05	-1.22	0.19	1.24	0.97	0.01	-0.53	-1.46	-0.22	1.39	1.12	-0.69	-0.34	-1.33	0.05	0.95	1.23	-0.59	-0.03	-1.5	-0.05	→	
1.28	1.26	-0.48	-0.3	-1.66	-0.09																																		
1.28	0.71	-0.45	-0.53	-1.1	0.09																																		
1.22	1.18	-1.32	-0.05	-1.22	0.19																																		
1.24	0.97	0.01	-0.53	-1.46	-0.22																																		
1.39	1.12	-0.69	-0.34	-1.33	0.05																																		
0.95	1.23	-0.59	-0.03	-1.5	-0.05																																		

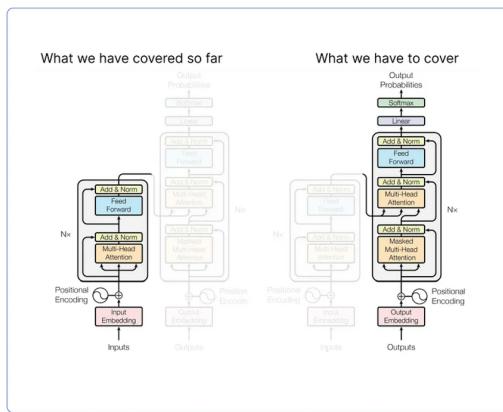
Add and Norm after Feed Forward Network

The output matrix of this add and norm step will serve as the query and key matrix in one of the multi-head attention mechanisms present in the decoder part, which you can easily understand by tracing outward from the add and norm to the decoder section.

### Step 11 — Decoder Part

The good news is that up until now, we have calculated **Encoder part**, all the steps that we have performed, from encoding our dataset to passing our matrix through the feedforward network, are unique. It means we haven't calculated them before. But from now on, all the upcoming steps that is the remaining architecture of the transformer (**Decoder part**) are going to involve similar kinds of matrix multiplications.

Take a look at our transformer architecture. What we have covered so far and what we have to cover yet:

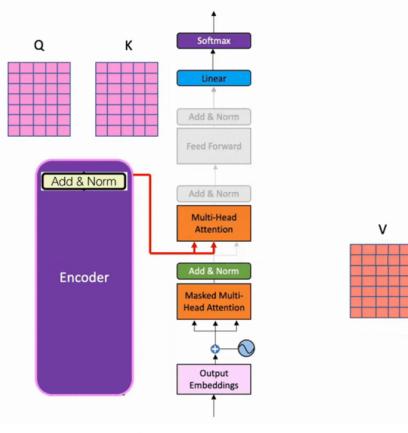


Upcoming steps illustration

We won't be calculating the entire decoder because most of its portion contains similar calculations to what we have already done in the encoder. Calculating the decoder in detail would only make the blog lengthy due to repetitive steps. Instead, we only need to focus on the calculations of the input and output of the decoder.

When training, there are two inputs to the decoder. One is from the encoder, where the output matrix of the last add and norm layer serves as the **query** and **key** for the second multi-head attention layer in the decoder part. Below is the visualization of it (from [batool haider](#)):

Top highlight



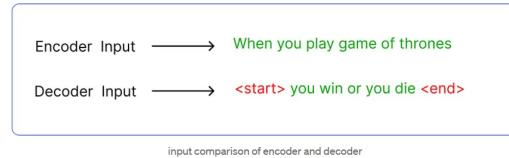
Visualization is from [Batool Haider](#)

While the value matrix comes from the decoder after the first **add and norm** step.

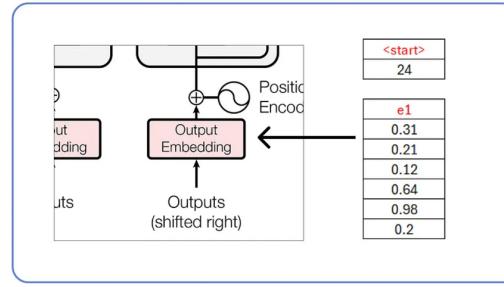
The second input to the decoder is the predicted text. If you remember, our input to the encoder is `when you play game of thrones` so the input to the

decoder is the predicted text, which in our case is `you win or you die`.

But the predicted input text needs to follow a standard wrapping of tokens that make the transformer aware of where to start and where to end.

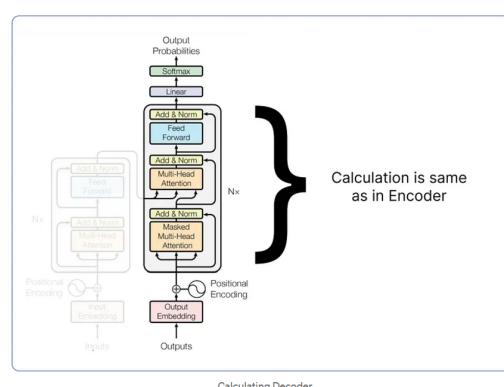


Where `<start>` and `<end>` are two new tokens being introduced. Moreover, the decoder takes one token at a time. It means that `<start>` will be served as an input, and `you` must be the predicted text for it.



As we already know, these embeddings are filled with random values, which will later be updated during the training process.

Compute rest of the blocks in the same way that we computed earlier in the encoder part.



Before diving into any further details, we need to understand what masked multi-head attention is, using a simple mathematical example.

### Step 12 — Understanding Mask Multi Head Attention

In a Transformer, the masked multi-head attention is like a spotlight that a model uses to focus on different parts of a sentence. It's special because it doesn't let the model cheat by looking at words that come later in the sentence. This helps the model understand and generate sentences step by step, which is important in tasks like talking or translating words into another language.

Suppose we have the following input matrix, where each row represents a position in the sequence, and each column represents a feature:

$$\text{Input Matrix} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$$

input matrix for masked multi head attentions

Now, let's understand the masked multi-head attention components having two heads:

- Linear Projections (Query, Key, Value):** Assume the linear projections for each head: Head 1:  $Wq1, Wk1, Wv1$  and Head 2:  $Wq2, Wk2, Wv2$
- Calculate Attention Scores:** For each head, calculate attention scores using the dot product of Query and Key, and apply the mask to prevent attending to future positions.
- Apply Softmax:** Apply the softmax function to obtain attention weights.
- Weighted Summation (Value):** Multiply the attention weights by the Value to get the weighted sum for each head.
- Concatenate and Linear Transformation:** Concatenate the outputs from both heads and apply a linear transformation.

**Let's do a simplified calculation:**

Assuming two conditions

- $Wq1 = Wk1 = Wv1 = Wq2 = Wk2 = Wv2 = I$ , the identity matrix.
- $Q=K=V=\text{Input Matrix}$

Head 1:	Head 2:
$Q_1 = K_1 = V_1 = \text{Input Matrix}$	$Q_2 = K_2 = V_2 = \text{Input Matrix}$
$A_1 = Q_1 \cdot K_1^T$	$A_2 = Q_2 \cdot K_2^T$
$A_1 = \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix}$	$A_2 = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$
$A_1 = \begin{bmatrix} 1 & 0 & 0 \\ 4 & 5 & 0 \\ 7 & 8 & 9 \end{bmatrix}$ (Masked)	$A_2 = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 5 & 6 \\ 0 & 0 & 9 \end{bmatrix}$ (Masked)
$W_1 = \text{softmax}(A_1)$	$W_2 = \text{softmax}(A_2)$
$O_1 = W_1 \cdot V_1$	$O_2 = W_2 \cdot V_2$

**Concatenate and Linear Transformation:**

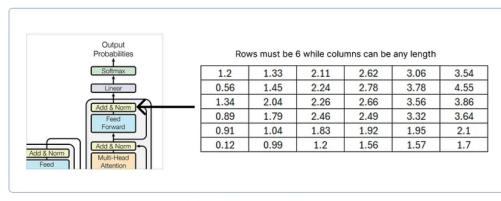
Concatenate( $[O_1, O_2]$ )  
(Apply Learnable Linear Transformation)

Mask Multi Head Attention (Two Heads)

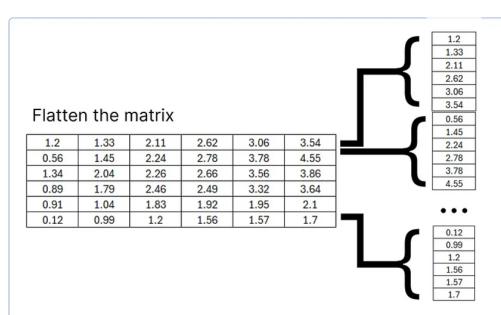
The concatenation step combines the outputs from the two attention heads into a single set of information. Imagine you have two friends who each give you advice on a problem. Concatenating their advice means putting both pieces of advice together so that you have a more complete view of what they suggest. In the context of the transformer model, this step helps capture different aspects of the input data from multiple perspectives, contributing to a richer representation that the model can use for further processing.

### Step 13 — Calculating the Predicted Word

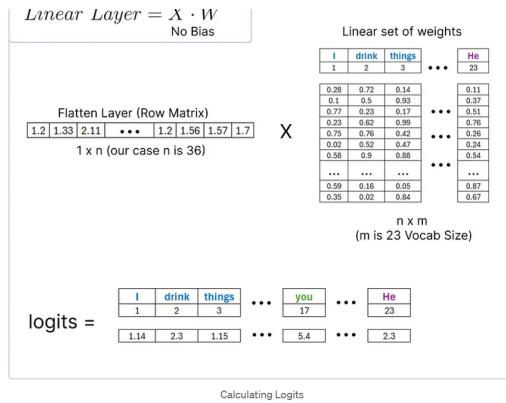
The output matrix of the last add and norm block of the decoder must contain the same number of rows as the input matrix, while the number of columns can be any. Here, we work with 6.



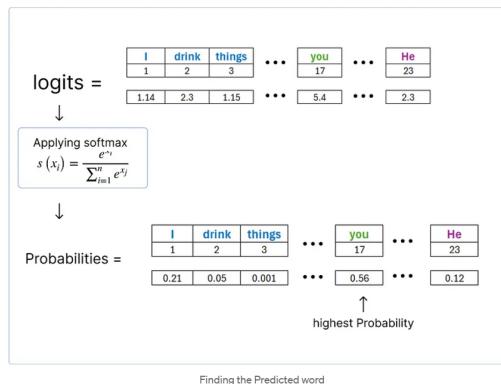
The last add and norm block resultant matrix of the decoder must be flattened in order to match it with a linear layer to find the predicted probability of each unique word in our dataset (corpus).



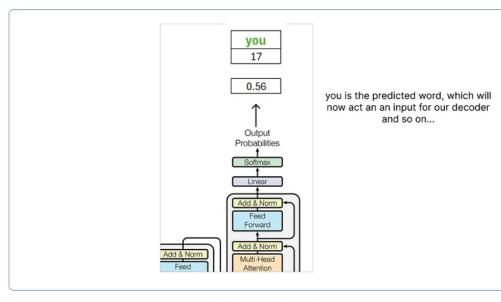
This flattened layer will be passed through a linear layer to compute the logits (scores) of each unique word in our dataset.



Once we obtain the logits, we can use the softmax function to normalize them and find the word that contains the highest probability.



So based on our calculations, the predicted word from the decoder is `you`.



This predicted word `you`, will be treated as the input word for the decoder, and this process continues until the `<end>` token is predicted.

### Important Points

1. The above example is very simple, as it does not involve epochs or any other important parameters that can only be visualized using a programming language like Python.
2. It has shown the process only until training, while evaluation or testing cannot be visually seen using this matrix approach.
3. Masked multi-head attentions can be used to prevent the transformer from looking at the future, helping to avoid overfitting your model.

### Conclusion

In this blog, I have shown you a very basic way of how transformers mathematically work using matrix approaches. We have applied positional encoding, softmax, feedforward network, and most importantly, multi-head attention.

In the future, I will be posting more blogs on transformers and LLM as my core focus is on NLP. More importantly, if you want to build your own million-parameter LLM from scratch using Python, I have written a blog on it which has received a lot of appreciation on Medium. You can read it here:



Have a great time reading!

Data Science Artificial Intelligence Machine Learning Python Deep Learning

2.8K 38

W+ ↑



### Written by Fareed Khan

31K Followers · Writer for Level Up Coding

AI Engineer, I write on AI <https://www.linkedin.com/in/fareed-khan-dev/>

Follow



### More from Fareed Khan and Level Up Coding



Fareed Khan

#### Modern GUI using Tkinter

There are two things to remember:

Sep 4, 2022 1K 8



Bret Cameron in Level Up Coding

#### Automate Your Release Notes with AI

How to save time every week using GitLab, OpenAI, and Node.js

1d ago 15

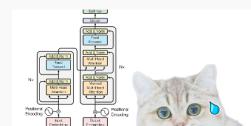


Jennifer Fu in Level Up Coding

#### Ace Your Microservices Interviews

Essential concepts and insights of microservices

1d ago 359 4



Fareed Khan

#### Understanding Transformers: A Step-by-Step Math Example—Part 1

I understand that the transformer architecture may seem scary, and you might...

Jun 5, 2023 2.8K 60

See all from Fareed Khan See all from Level Up Coding

### Recommended from Medium

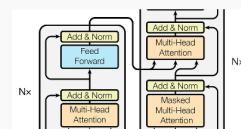


Abdur Rahman in Stackademic

#### Python is No More The King of Data Science

5 Reasons Why Python is Losing Its Crown

Oct 23 7K 29



Code Thulo

#### Top Large Language Model (LLM) Interview Question | Basic LLM...

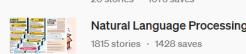
Large Language Models (LLMs) are deep learning models that are trained on huge...

Oct 17 5

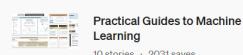
### Lists



Predictive Modeling w/  
Python  
20 stories · 1676 saves



Natural Language Processing  
1815 stories · 1428 saves



Practical Guides to Machine  
Learning  
10 stories · 2031 saves



ChatGPT  
21 stories · 879 saves

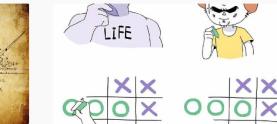


👤 Cristian Leo

### The Math Behind Transformers

Deep Dive into the Transformer Architecture, the key element of LLMs. Let's explore its...

👉 Jul 25 🎉 884 💬 6

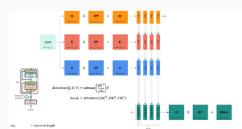


👤 Amit Yadav

### Computer Vision Project Ideas With Code

Not a Medium member? Read the full story by clicking here.

👉 Aug 7 🎉 114



👤 Vipra Singh

### LLM Architectures Explained: Coding a Transformer (Part 7)

Deep Dive into the architecture & building real-world applications leveraging NLP...

👉 Nov 10 🎉 285 💬 3



👤 Fareed Khan in Level Up Coding

### Building a Million-Parameter LLM from Scratch Using Python

A Step-by-Step Guide to Replicating LLaMA Architecture

👉 Dec 7, 2023 🎉 2.7K 💬 33

See more recommendations

Help Status About Careers Press Blog Privacy Terms Text to speech Teams