

# 9. Transformer

## LING-581-Natural Language Processing 1

---

Instructor: Hakyung Sung

October 21, 2025

\*Acknowledgment: These course slides are based on materials from CS224N @ Stanford University; Dr. Kilho Shin @ Kyocera

# Table of contents

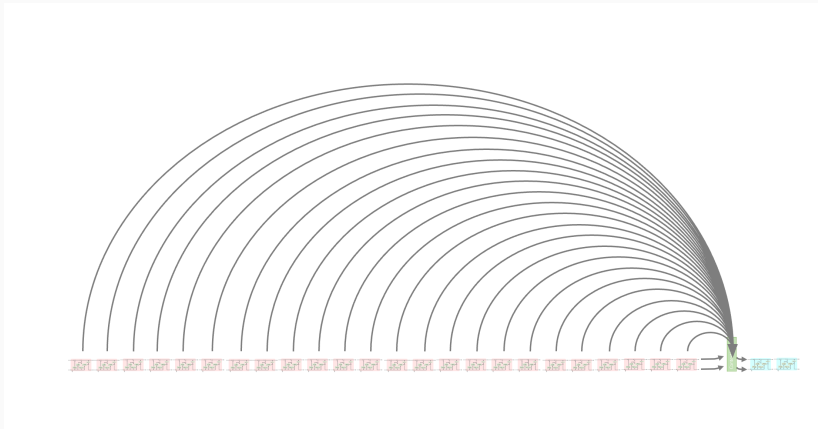
1. Review: Attention
2. The transformer model
3. Understanding the Transformer with example
4. Wrap up

## **Review: Attention**

---

## Seq2Seq: Problem

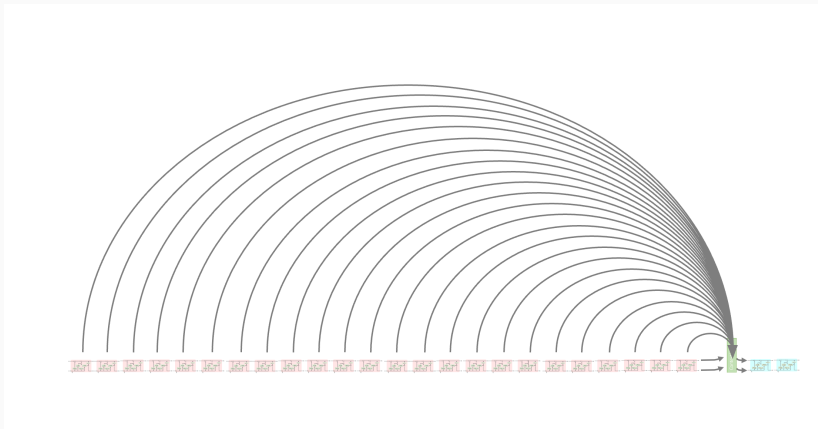
It becomes quite difficult to pack all the information of the input sequence into a fixed-length context vector.



This is called **Bottleneck** problem.

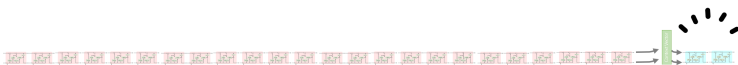
## Seq2Seq: Problem

The **attention** mechanism was introduced to address this problem.



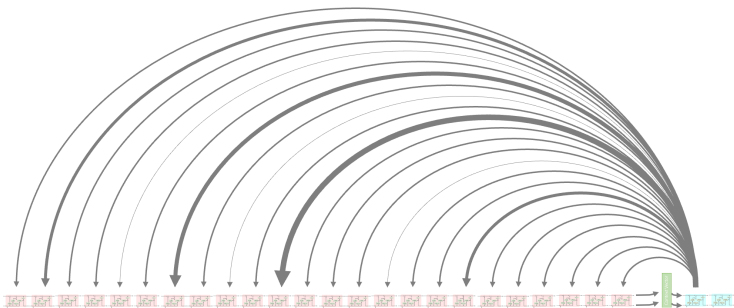
# Seq2Seq: Problem

When the decoder generates each word in the output sequence, the attention mechanism



## Seq2Seq: Problem

is an algorithm that makes it “attend” to which parts of the input sequence are important.



# Attention

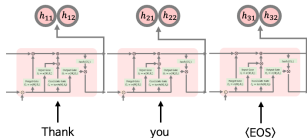
Suppose the input sequence comes in like this:





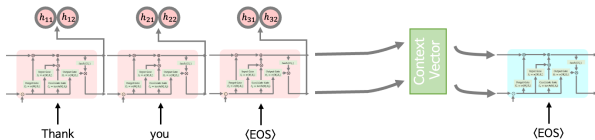
# Attention

We store the hidden state for each input word separately.



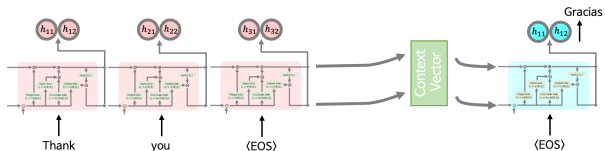
# Attention

We build a context vector and feed it to the decoder,



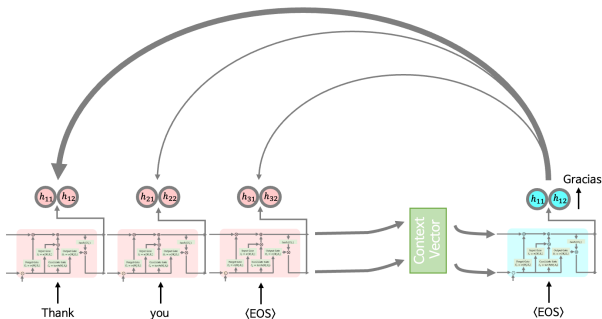
# Attention

and obtain the decoder's hidden state and output as follows.



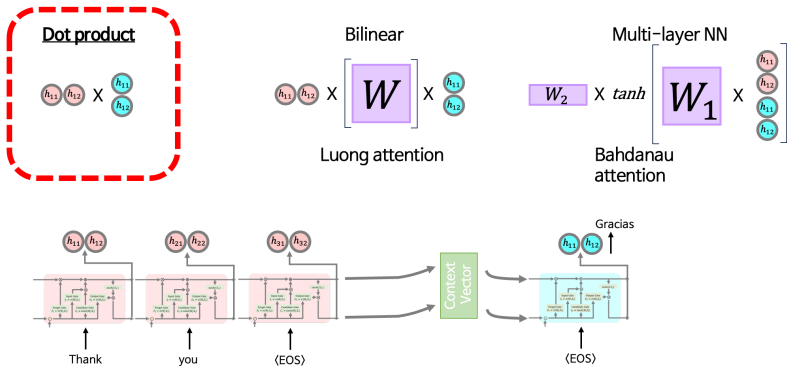
# Attention

The **similarity** between two vectors works as a measure that determines the relationship between two states.



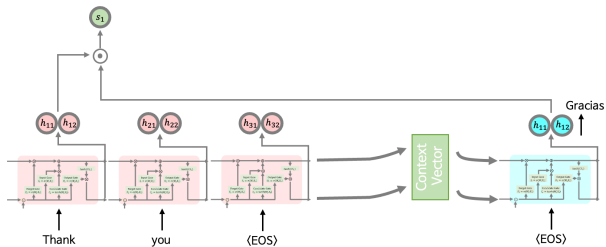
# Attention

We focused on the simplest method: the dot product.



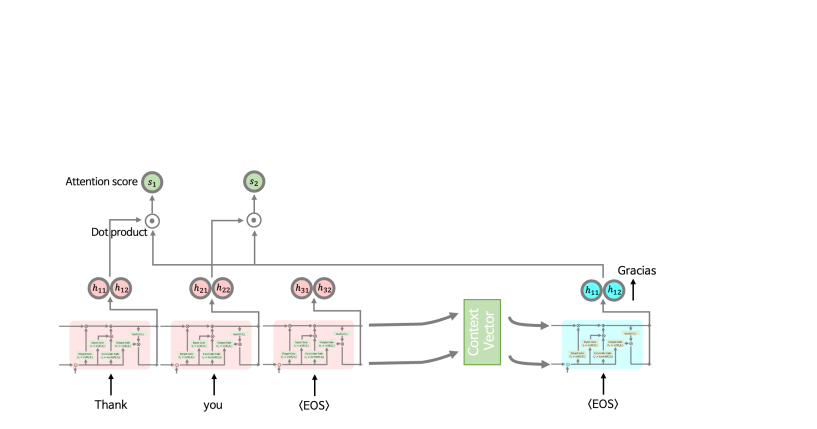
# Attention

For example:



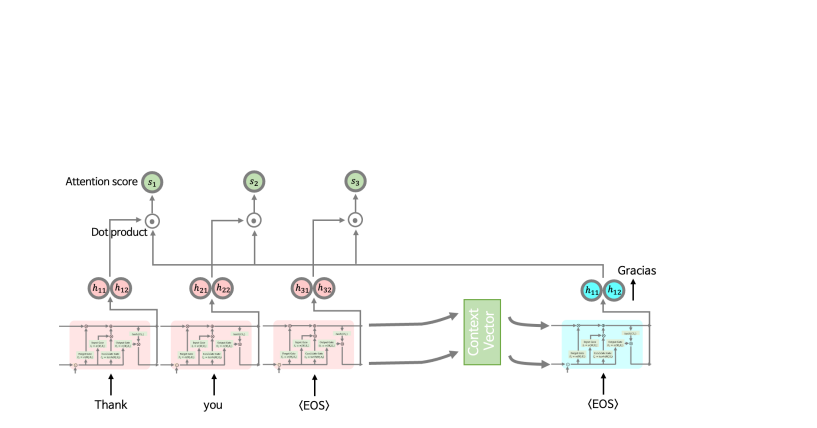
# Attention

For example:



# Attention

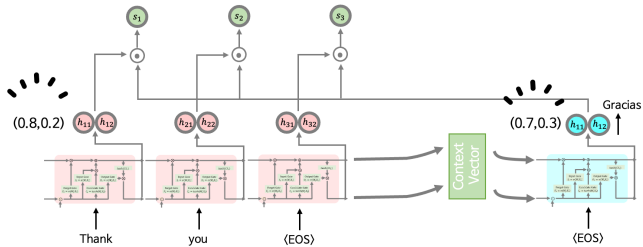
For example:





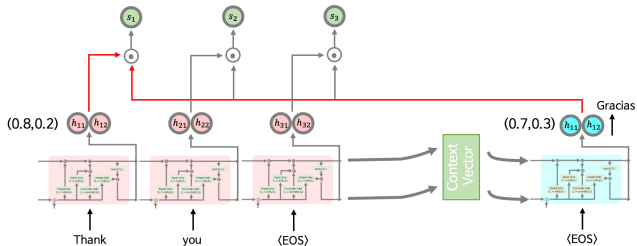
# Attention

Suppose the encoder hidden state is  $(0.8, 0.2)$  and the decoder hidden state is  $(0.7, 0.3)$ .



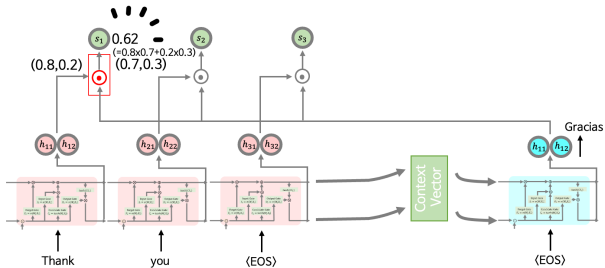
# Attention

Suppose the encoder hidden state is  $(0.8, 0.2)$  and the decoder hidden state is  $(0.7, 0.3)$ .



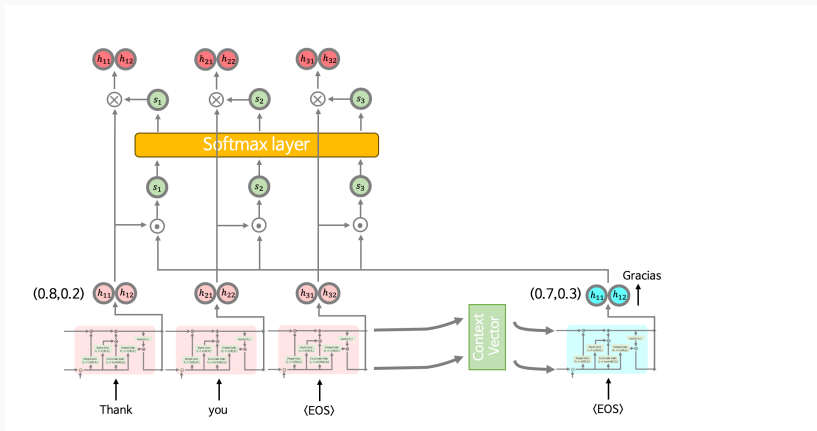
# Attention

Then the attention score  $s_1$  becomes 0.62 through the following dot-product computation.



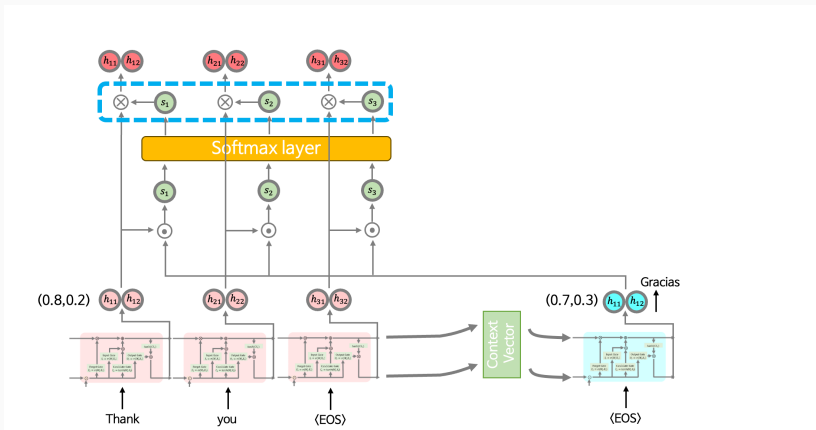
# Attention

Next, we compute the **softmax** of each attention score



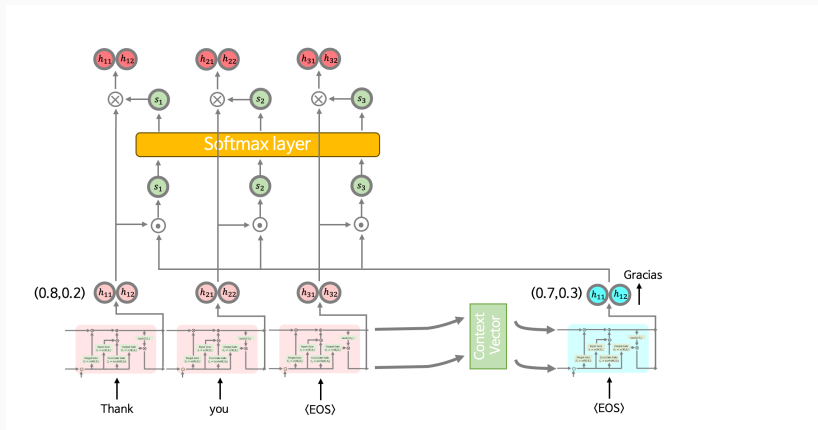
# Attention

to convert the attention scores into a **probability distribution** and **normalize** them.



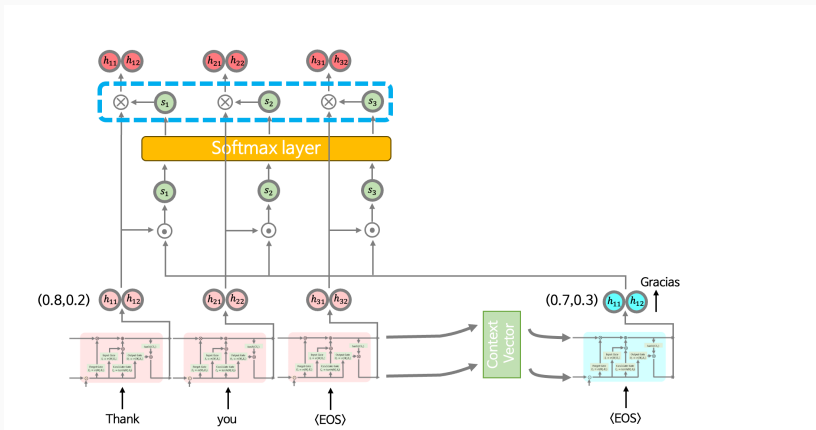
# Attention

Next, we multiply each attention score by its corresponding input hidden state.



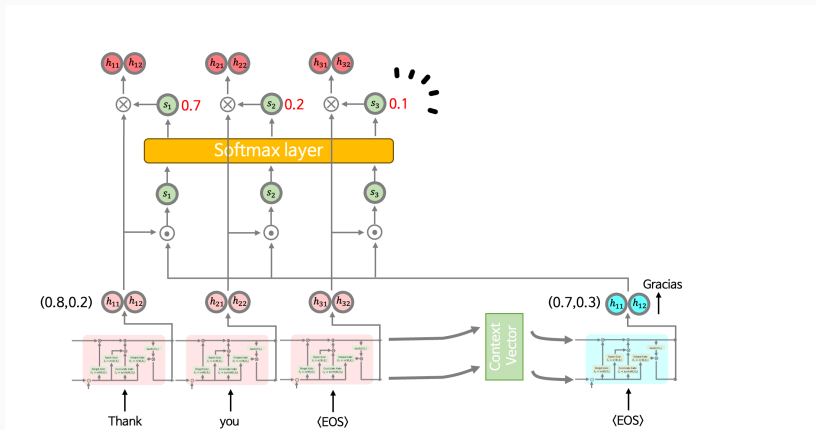
# Attention

The multiplication has the effect of **amplifying** the input hidden states according to their attention weights.



# Attention

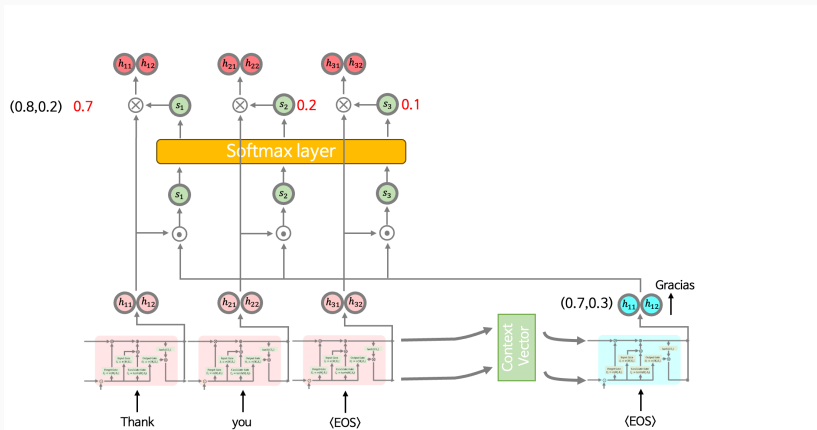
For example, if the attention scores after the softmax layer are as follows—





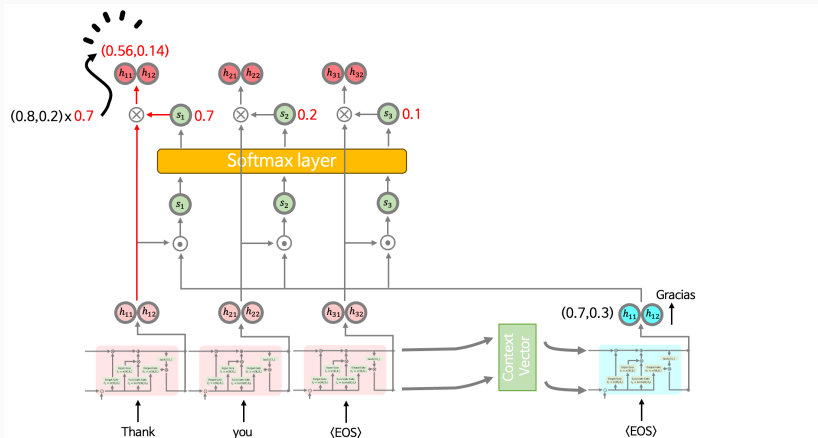
# Attention

For example, if the attention scores after the softmax layer are as follows—



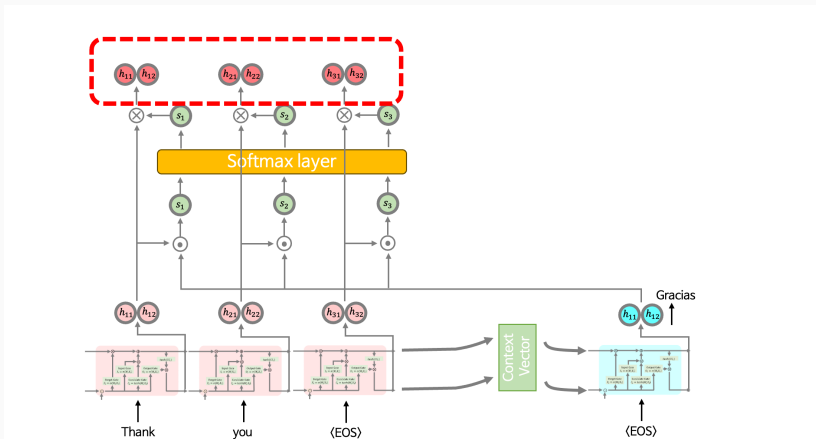
# Attention

Then the input state  $(0.8, 0.2)$  multiplied by  $0.7$  becomes  $(0.56, 0.14)$ .



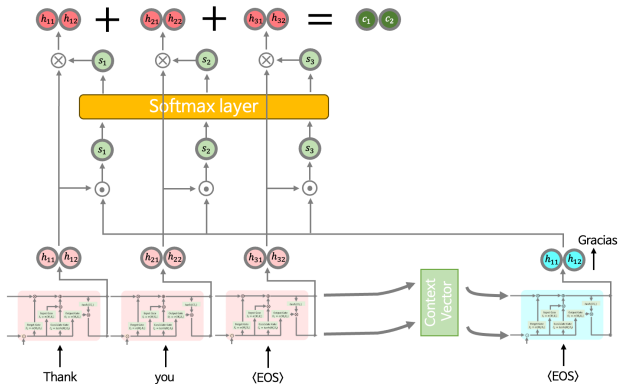
# Attention

Now, take these **attention-weighted** input hidden states,



# Attention

sum them up, and you obtain a **new context vector**.



## Attention is a *general* deep learning technique

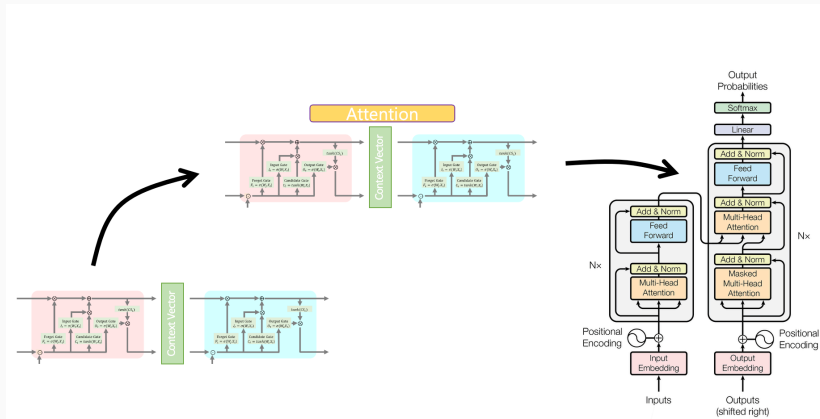
- Attention has become the powerful, flexible, general way pointer and memory manipulation in deep learning models. (A new idea from 2010).

# The transformer model

---

# Recall

Eventually led to the development of the **Transformer** model.



**Multi-head** attentions

<https://www.youtube.com/watch?v=eM1x5fFNoYc> Any hints from the video clip?



## RNNs: Long-dependency problem

- RNNs process input **sequentially**, one token at a time, passing information through hidden states.

## RNNs: Long-dependency problem

- RNNs process input **sequentially**, one token at a time, passing information through hidden states.
- We discussed the long-dependency problem.

## RNNs: Long-dependency problem

- RNNs process input **sequentially**, one token at a time, passing information through hidden states.
- We discussed the long-dependency problem.
- Somewhat solved through attention algorithms?

## RNNs: Lack of parallelizability

- RNNs process input step by step — each hidden state depends on the previous one.

## RNNs: Lack of parallelizability

- RNNs process input step by step — each hidden state depends on the previous one.
- This means both the forward and backward passes require  $O(\text{sequence length})$  sequential operations.

## RNNs: Lack of parallelizability

- RNNs process input step by step — each hidden state depends on the previous one.
- This means both the forward and backward passes require  $O(\text{sequence length})$  sequential operations.
- GPUs are great at performing many independent computations in parallel, but RNNs don't allow this because future states can't be computed until past states are done.

## RNNs: Lack of parallelizability

- RNNs process input step by step — each hidden state depends on the previous one.
- This means both the forward and backward passes require  $O(\text{sequence length})$  sequential operations.
- GPUs are great at performing many independent computations in parallel, but RNNs don't allow this because future states can't be computed until past states are done.
- As a result, training RNNs on very large datasets becomes **slow and inefficient**.

## RNNs+Attention!

- Seq2Seq (RNN) models: Attention connects the decoder to the encoder — each decoder step selectively focuses on encoder hidden states.  
(Cross-attention only; still sequential)



# RNNs+Attention!

- Seq2Seq (RNN) models: Attention connects the decoder to the encoder — each decoder step selectively focuses on encoder hidden states.  
(Cross-attention only; still sequential)
- In Transformers, attention also occurs **within a single sentence** — all words attend to all words in the previous layer.  
(Self-attention + Cross-attention)

# RNNs+Attention!

- Seq2Seq (RNN) models: Attention connects the decoder to the encoder — each decoder step selectively focuses on encoder hidden states.  
(Cross-attention only; still sequential)
- In Transformers, attention also occurs **within a single sentence** — all words attend to all words in the previous layer.  
(Self-attention + Cross-attention)
- This means every word can interact with every other word directly — no need to wait for sequential computation.

## RNNs+Attention!

- Seq2Seq (RNN) models: Attention connects the decoder to the encoder — each decoder step selectively focuses on encoder hidden states.  
(Cross-attention only; still sequential)
- In Transformers, attention also occurs **within a single sentence** — all words attend to all words in the previous layer.  
(Self-attention + Cross-attention)
- This means every word can interact with every other word directly — no need to wait for sequential computation.
- As a result, Transformers overcome both long-distance dependency and lack of parallelizability.

# RNNs+Attention!

- Seq2Seq (RNN) models: Attention connects the decoder to the encoder — each decoder step selectively focuses on encoder hidden states.  
(Cross-attention only; still sequential)
- In Transformers, attention also occurs **within a single sentence** — all words attend to all words in the previous layer.  
(Self-attention + Cross-attention)
- This means every word can interact with every other word directly — no need to wait for sequential computation.
- As a result, Transformers overcome both long-distance dependency and lack of parallelizability.
- *Notes.* This was NOT an entirely new ways of looking NLP problems (e.g., probabilistic language models → neural network), but made a huge progress in the field.

# Parallelization in Transformers

- **Self-Attention**: All tokens attend to all tokens in the same layer simultaneously.

# Parallelization in Transformers

- **Self-Attention:** All tokens attend to all tokens in the same layer simultaneously.
- Unlike RNNs, Transformer layers do not depend on previous hidden states; each token's representation is updated in parallel.

# Parallelization in Transformers

- **Self-Attention**: All tokens attend to all tokens in the same layer simultaneously.
- Unlike RNNs, Transformer layers do not depend on previous hidden states; each token's representation is updated in parallel.
- This allows GPUs to perform all attention computations at once using matrix multiplication.

# Parallelization in Transformers

- **Self-Attention:** All tokens attend to all tokens in the same layer simultaneously.
- Unlike RNNs, Transformer layers do not depend on previous hidden states; each token's representation is updated in parallel.
- This allows GPUs to perform all attention computations at once using matrix multiplication.
- Training and inference are therefore much faster, especially for long sequences and large datasets.



## From Dot Product to Soft Lookup

- Recall how we calculated the **dot product** between the decoder state (query) and each encoder hidden state (key) in the Seq2Seq attention.

## From Dot Product to Soft Lookup

- Recall how we calculated the **dot product** between the decoder state (query) and each encoder hidden state (key) in the Seq2Seq attention.
- That dot product measured how similar the query was to each key — giving us attention weights after softmax.

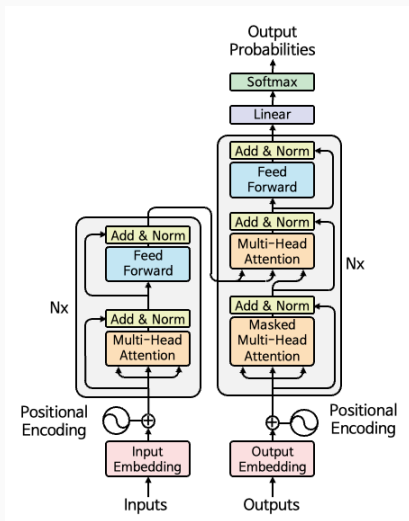
## From Dot Product to Soft Lookup

- Recall how we calculated the **dot product** between the decoder state (query) and each encoder hidden state (key) in the Seq2Seq attention.
- That dot product measured how similar the query was to each key — giving us attention weights after softmax.
- In Transformers, the same idea is extended: every token's representation acts as a **query**, looking up information from a set of keys and values (**soft, averaging lookup** in a key-value store)  $\Rightarrow$  The dot product attention becomes a **fuzzy retrieval mechanism** that allows each token to access information from all others in parallel.

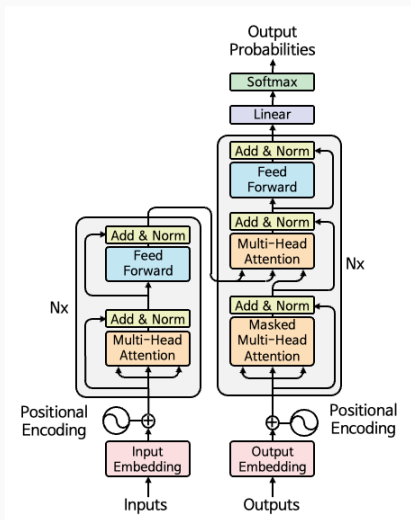
# **Understanding the Transformer with example**

---

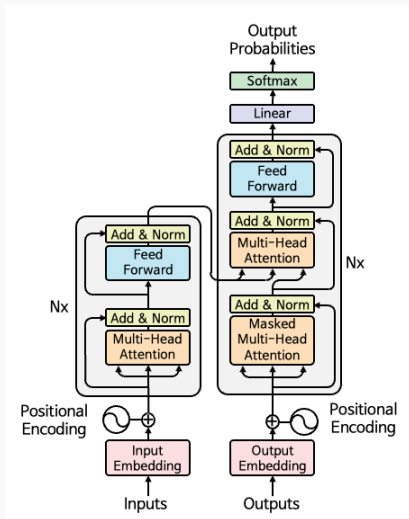
The structure of the Transformer (proposed by Vaswani et al., 2017) looks like this:



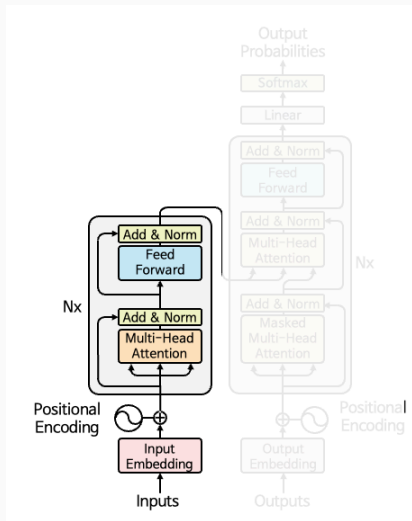
Before diving into the detailed calculations, let's first take a look at the overall structure.



The Transformer model can be broadly divided into two main parts.

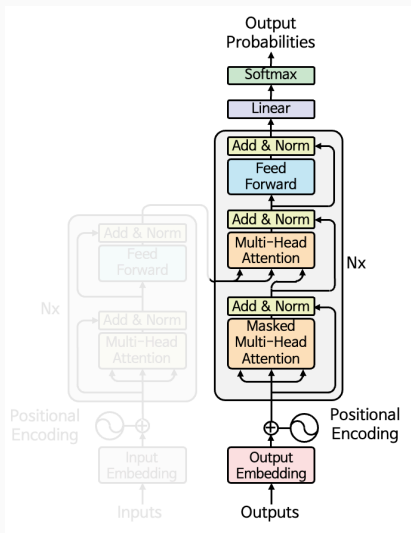


Here we see the **encoder** part,

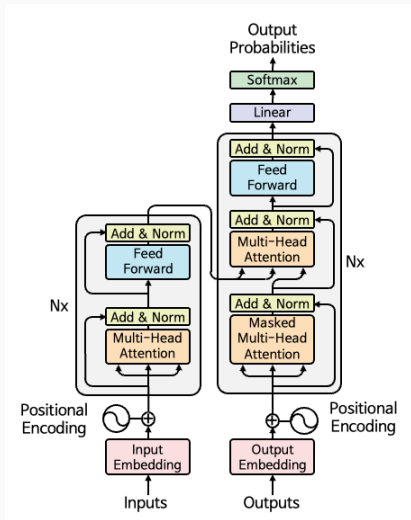




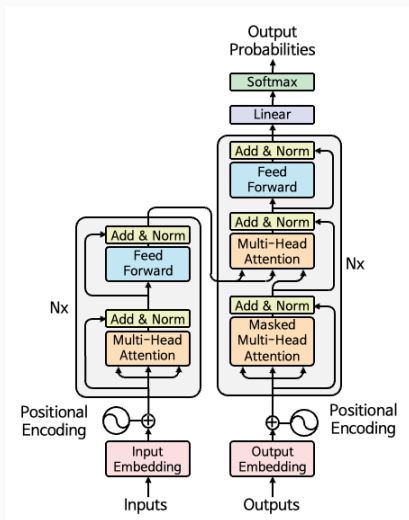
and the **decoder** part.



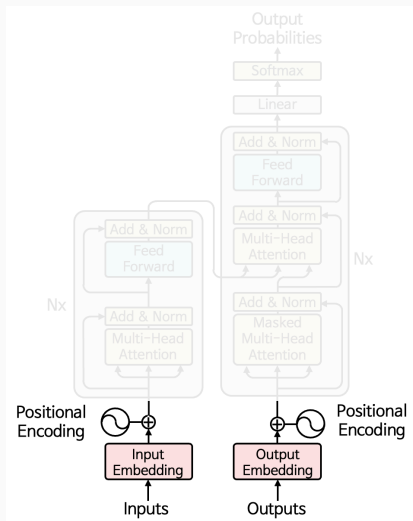
If we look closely at the entire Transformer model,



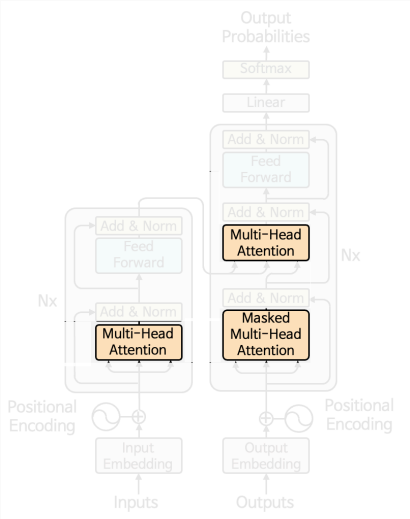
we can see that the same kinds of blocks are **repeatedly stacked**.



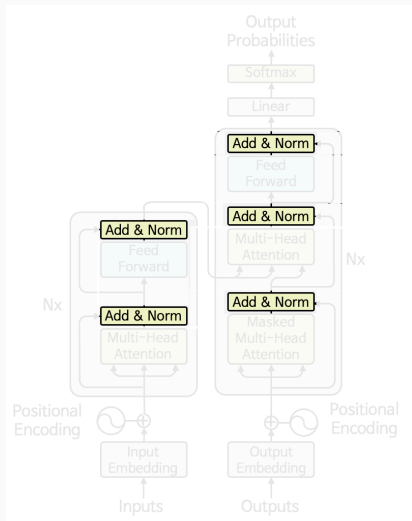
# Embedding layer



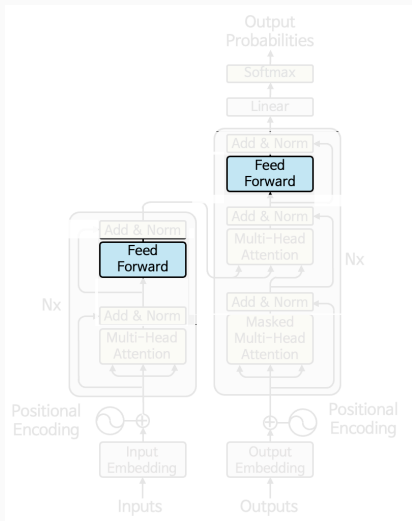
# Multi-head attention



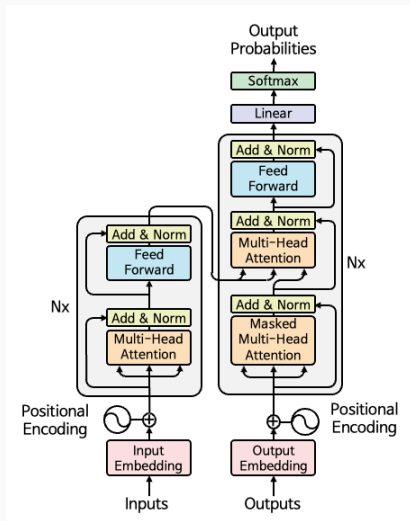
## Add & Norm layer



and also the Feed-Forward layer.

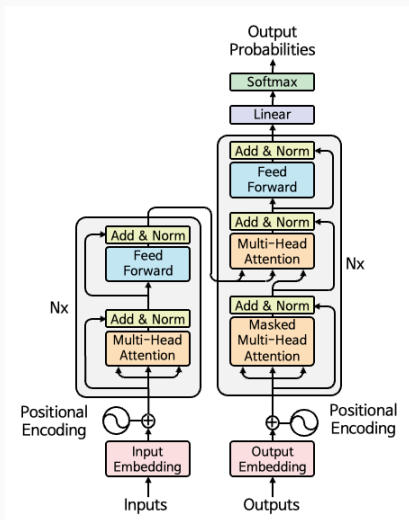


So, even though the Transformer may look complicated at first glance,

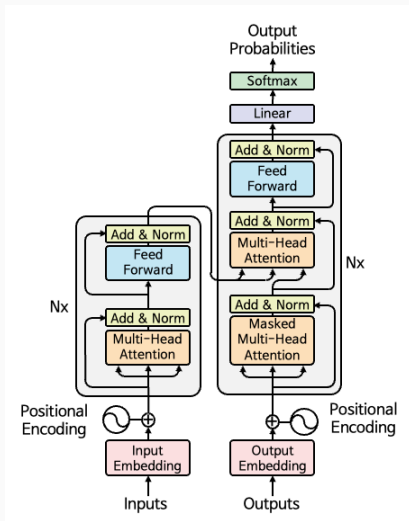




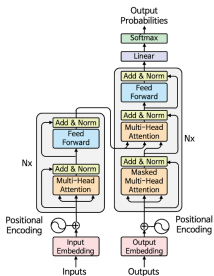
it's actually made up of a few components that are repeatedly stacked.



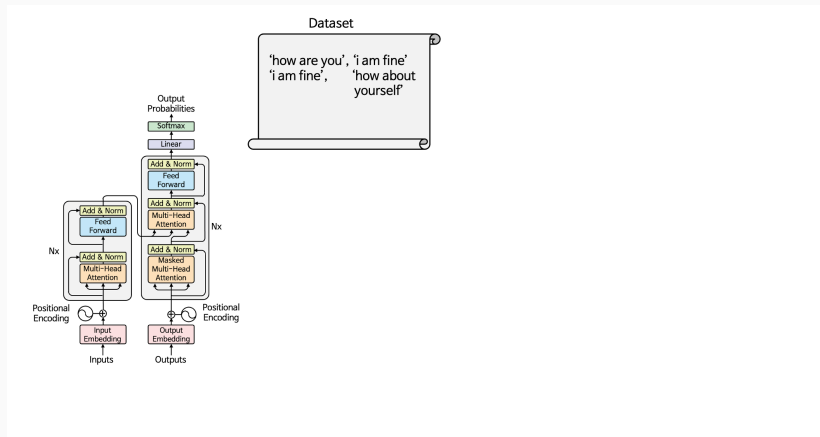
By examining each part step by step, we can fully understand how the model works.



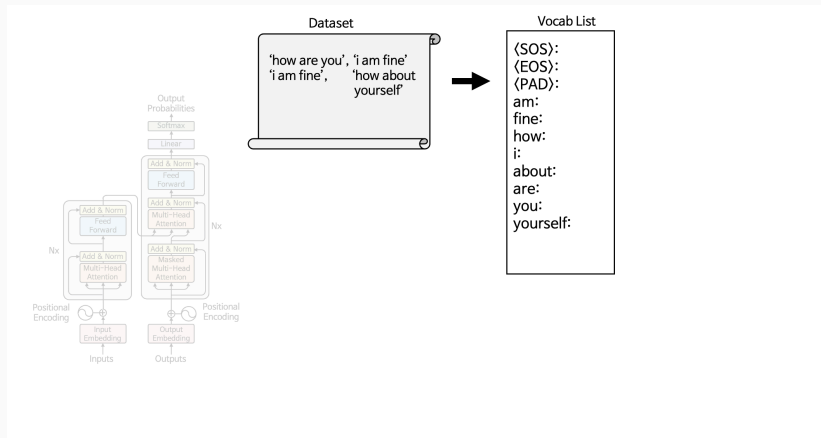
From now on, let's use a **simple example** to explore how the Transformer learns.



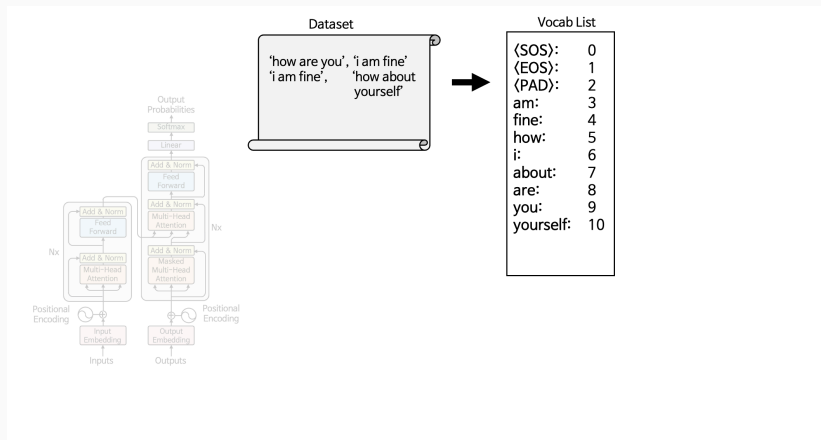
To train a Transformer model, the very first step is to **create a dataset**.



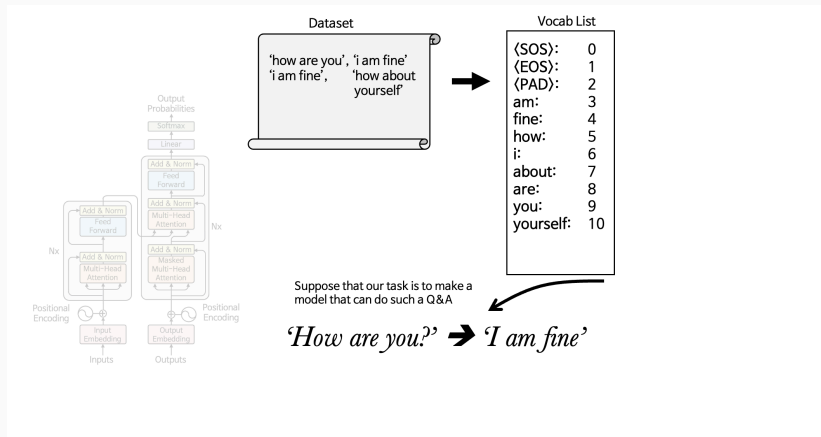
Next, we extract all the words from the dataset to build a vocabulary.



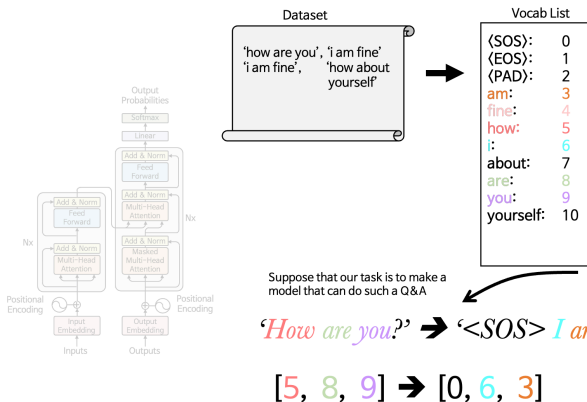
Then, we assign each word (or token) a unique number (index) so that the model can process it as input data.



We assign each word (token) a unique index to convert it into a form that the model can handle.

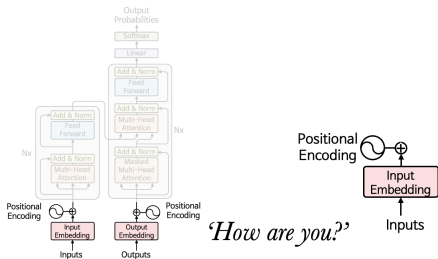


For simplicity, let's assume a very small number of tokens when assigning these indices.



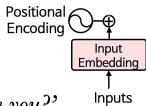


When an input sentence like this comes in, the first step is word embedding.

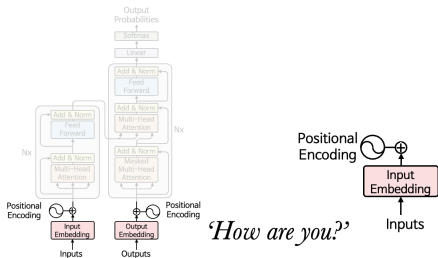


*How are you.?*

[5, 8, 9]



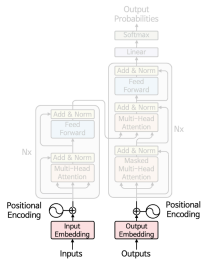
The input embedding block takes input tokens like [5, 8, 9] and



*‘How are you?’*

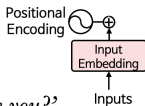
[5, 8, 9]

outputs the corresponding embedding vectors for each word.



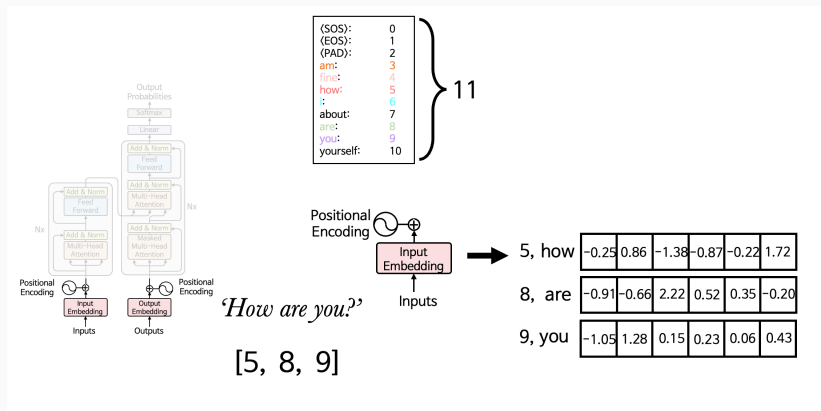
*How are you?*

[5, 8, 9]

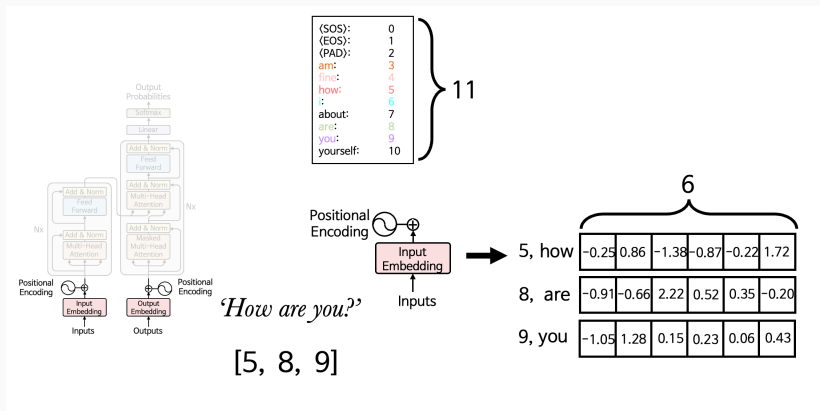


5, how	-0.25	0.86	-1.38	-0.87	-0.22	1.72
8, are	-0.91	-0.66	2.22	0.52	0.35	-0.20
9, you	-1.05	1.28	0.15	0.23	0.06	0.43

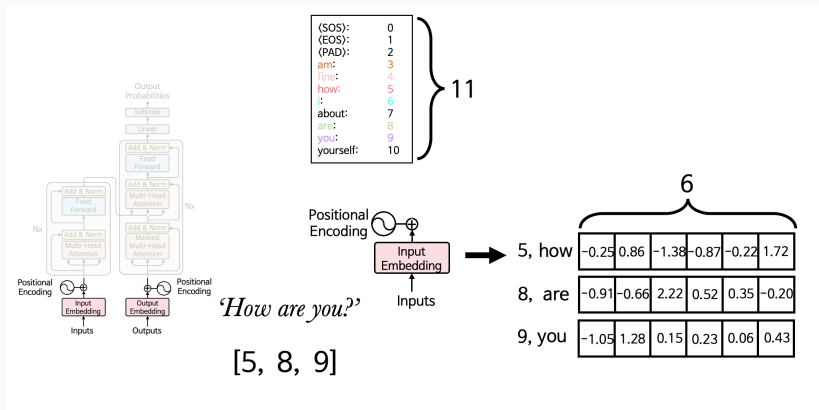
The embedding layer compresses, for example, 11 words into



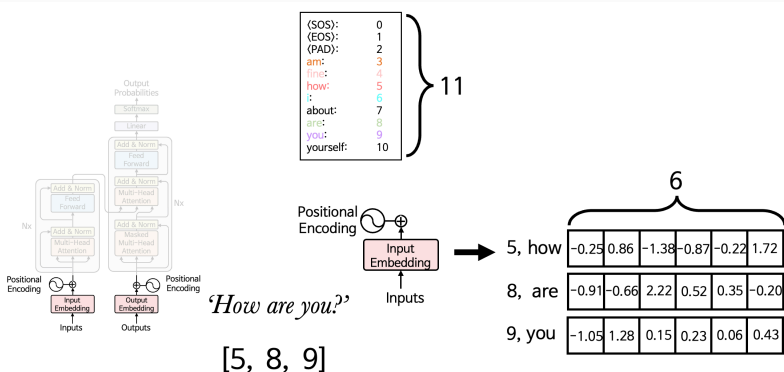
dense vectors of length 6.



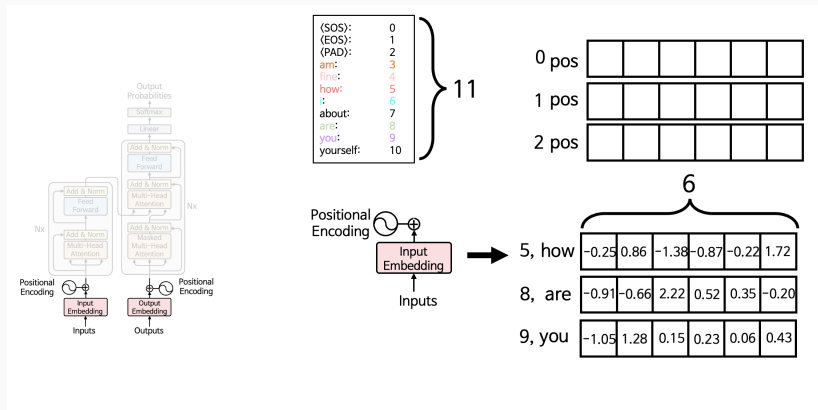
- In the original paper, 37,000 words were reduced to 512-dimensional embeddings. (So, it doesn't have to be 6 as in our example.)



- In the original paper, 37,000 words were reduced to 512-dimensional embeddings. (So, it doesn't have to be 6 as in our example.)
- It's important to understand that this dimensional embeddings allow the model to efficiently process and represent a large number of words internally.



Next is **positional** encoding. The Transformer uses a unique way to encode the position of words within a sentence.

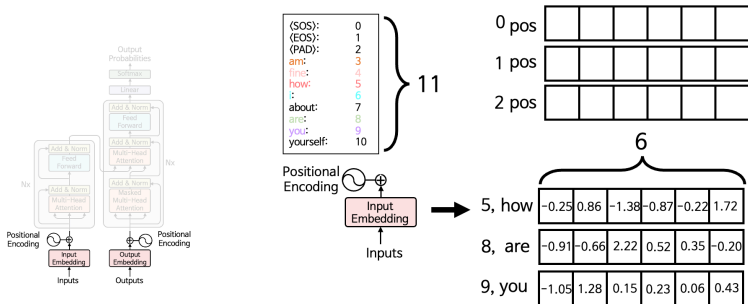




Why do we need this?

Why do we need this? Self-attention alone does not capture word order (e.g., RNNs); it treats inputs as a set, not a sequence.

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right) \quad PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$



We encode the position of each word based on the following formulas.

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) \quad PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

0 pos	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
1 pos	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
2 pos	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>

Let's take a closer look at these formulas.

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) \quad PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

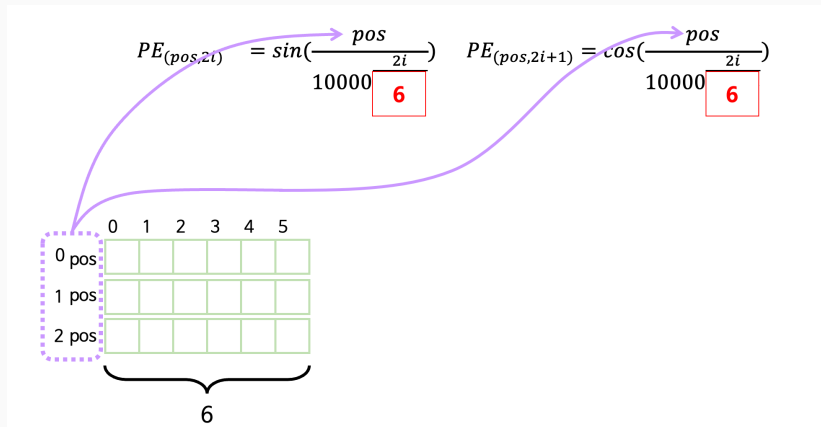
	0	1	2	3	4	5(d_model-1)
0 pos						
1 pos						
2 pos						

In this example, the value of  $d_{\text{model}}$  is 6 (the dimensionality of the model).

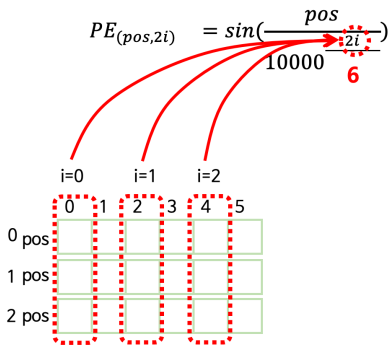
$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{6}}}\right) \quad PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{6}}}\right)$$

	0	1	2	3	4	5(d_model-1)
0 pos						
1 pos						
2 pos						

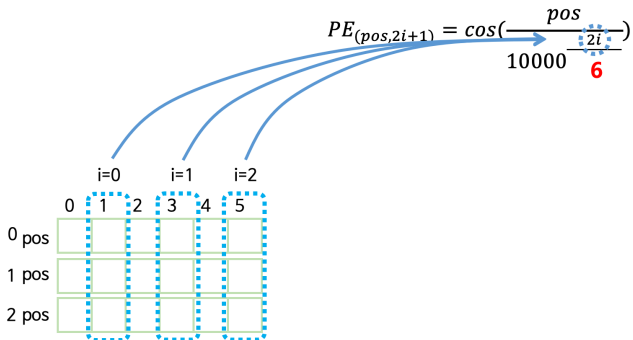
And the variable *pos* takes the values 0, 1, and 2 in order.



$i = 0, 1, 2$  denotes the dimension index used for even ( $2i$ ) and odd ( $2i+1$ ) components. For even dimensions, we apply this formula:



And for odd dimensions, we apply this formula:





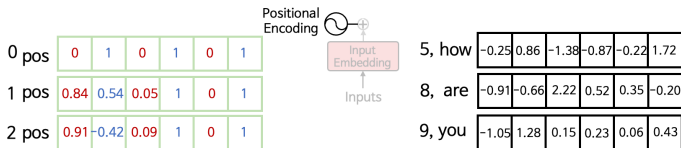
We can now compute the positional encoding values as follows:

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) \quad PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

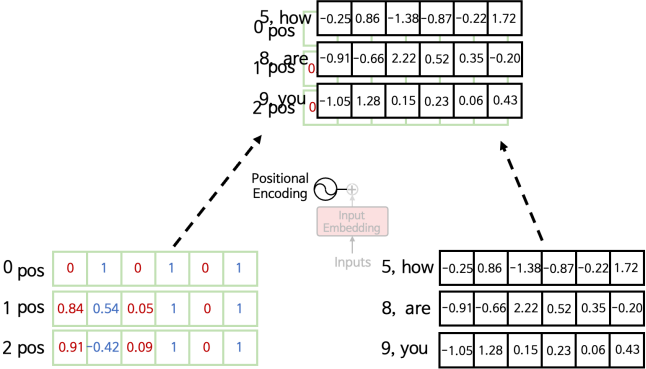
	0	1	2	3	4	5
0 pos	0	1	0	1	0	1
1 pos	0.84	0.54	0.05	1	0	1
2 pos	0.91	-0.42	0.09	1	0	1

└──────────────────┘  
6

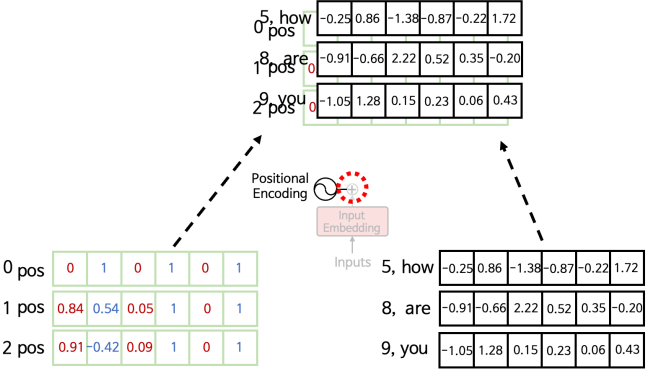
Once we have these positional embeddings, we simply add them to the input embeddings:



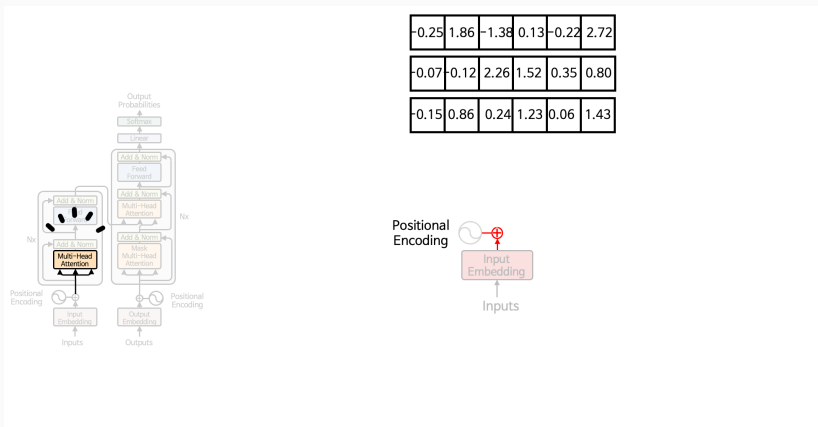
By adding them together, we obtain the combined input + positional embedding vectors.



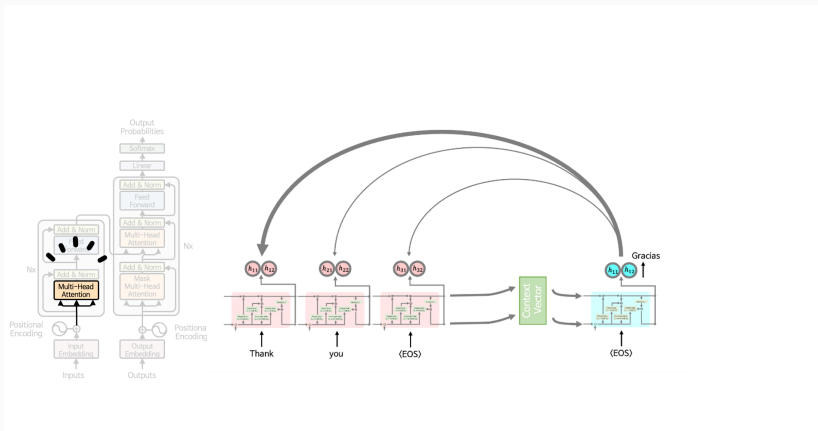
By adding them together, we obtain the combined input + positional embedding vectors.



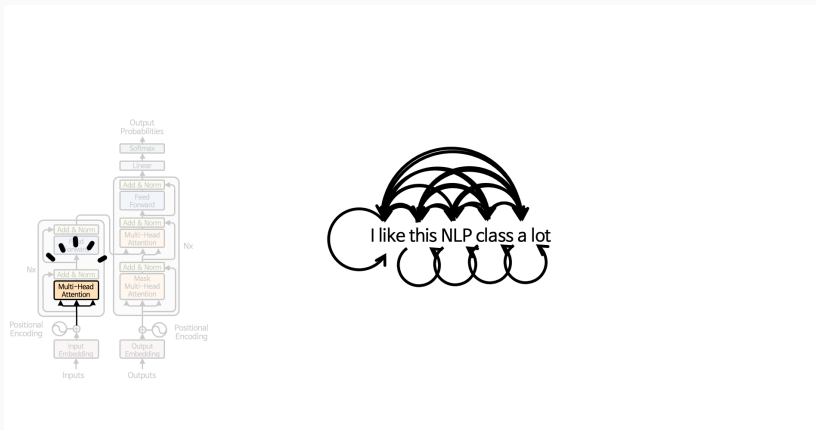
Now, it's time to feed this combined input-position matrix into the **multi-head attention**, which is the core component of the Transformer.



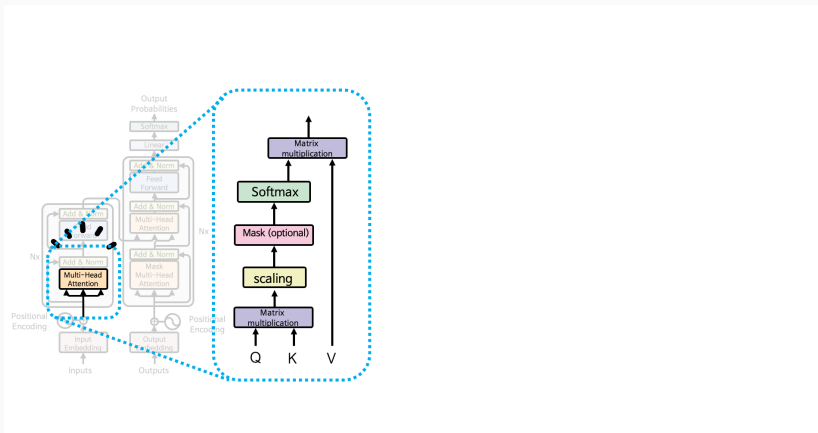
Again, the Transformer's multi-head attention is different from the attention mechanism used in traditional seq2seq models. While seq2seq attention focuses on aligning the input and output sequences,



the Transformer's attention captures the relationships between words within the same input sentence.

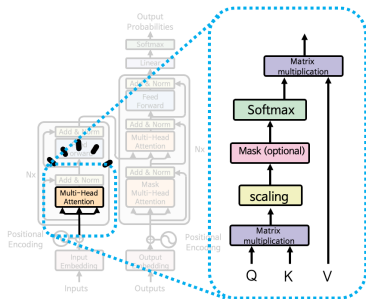


The structure of the multi-head attention mechanism used for self-attention looks like this:





# We make three copies of the input + positional encoding matrix.

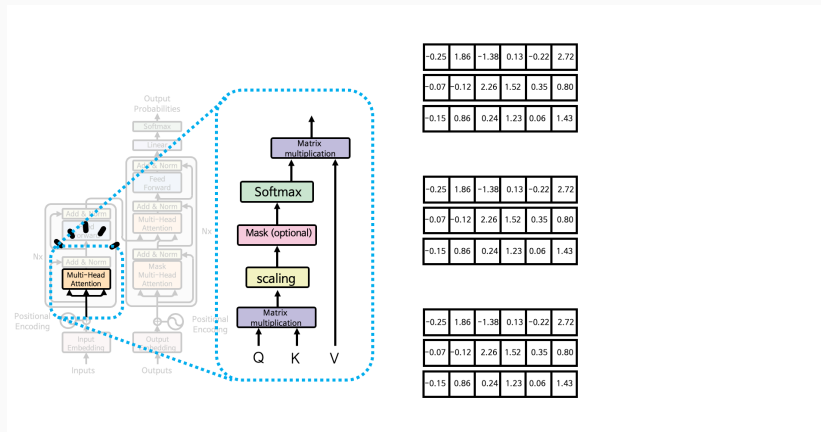


-0.25	1.86	-1.38	0.13	-0.22	2.72
-0.07	-0.12	2.26	1.52	0.35	0.80
-0.15	0.86	0.24	1.23	0.06	1.43

-0.25	1.86	-1.38	0.13	-0.22	2.72
-0.07	-0.12	2.26	1.52	0.35	0.80
-0.15	0.86	0.24	1.23	0.06	1.43

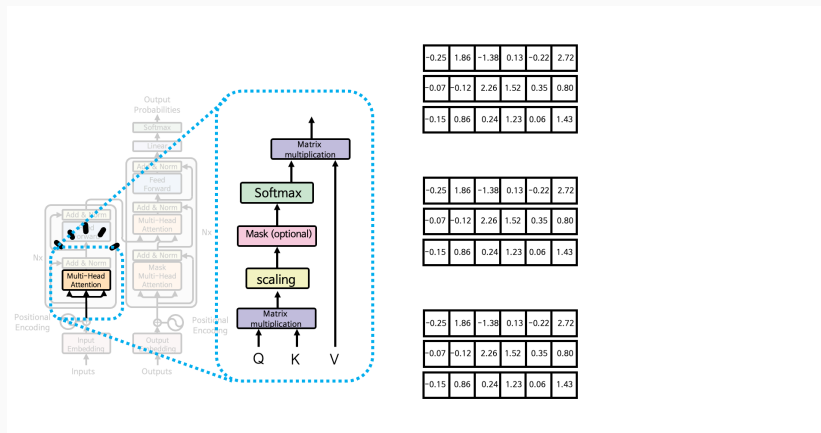
-0.25	1.86	-1.38	0.13	-0.22	2.72
-0.07	-0.12	2.26	1.52	0.35	0.80
-0.15	0.86	0.24	1.23	0.06	1.43

We make three copies of the input + positional encoding matrix.



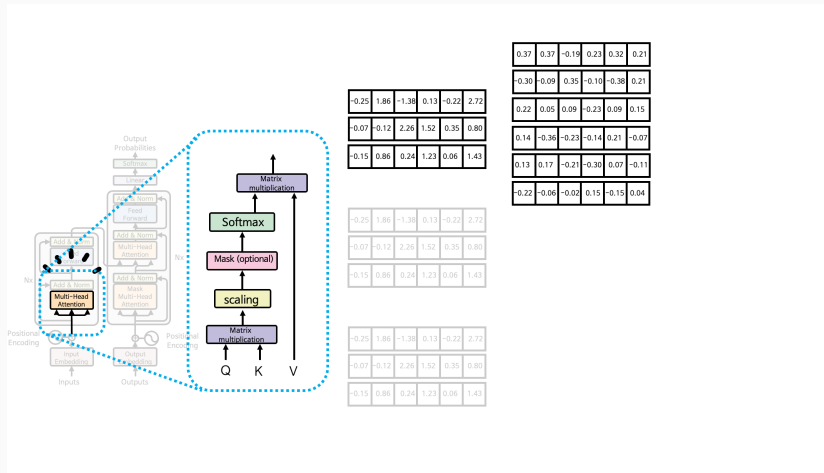
Why?

We make three copies of the input + positional encoding matrix.

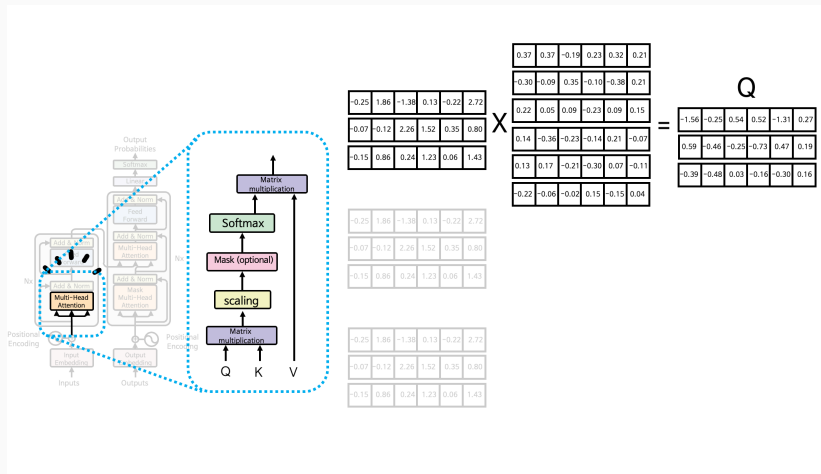


Why? This is done to create **Query (Q)**, **Key (K)**, and **Value (V)** matrices — each representing a different projection of the same input for the attention mechanism.

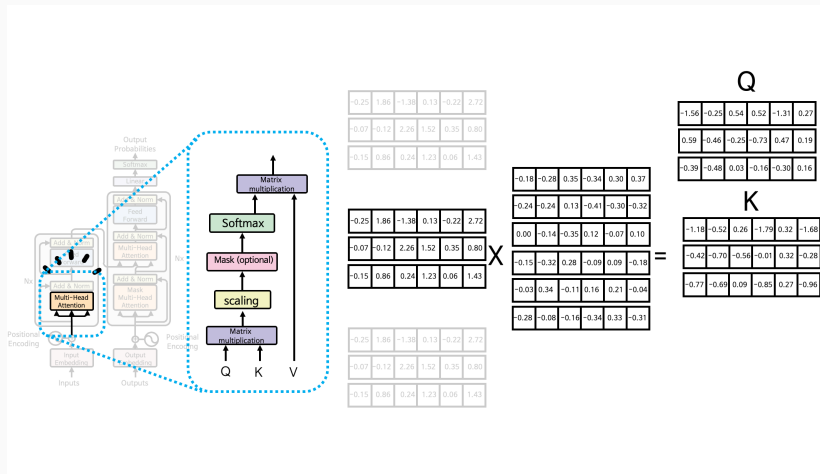
To obtain the Q matrix, we create the following 6×6 weight matrix (randomly initialized).



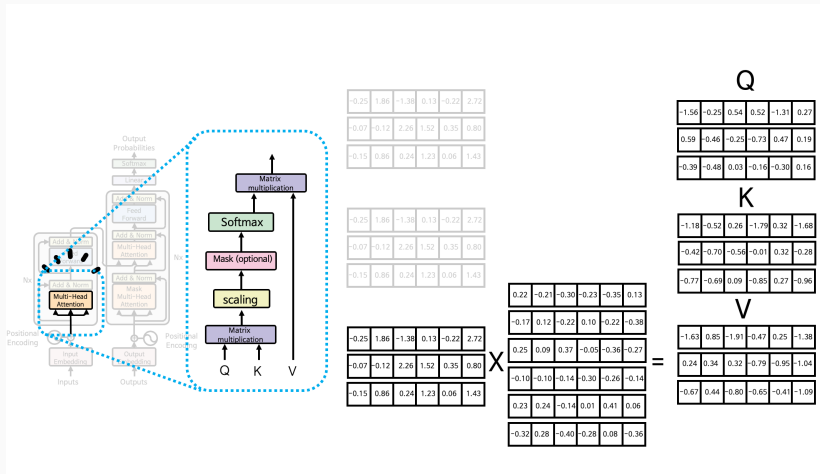
And then we perform matrix multiplication to obtain the Q (Query) matrix.



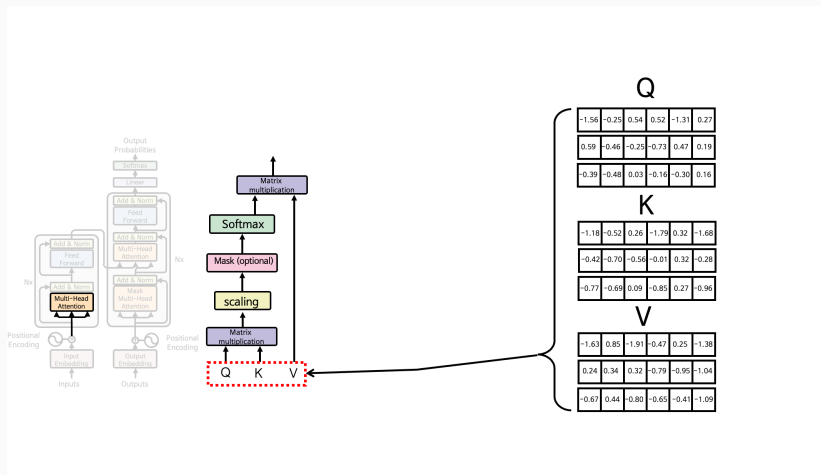
To compute the K (Key) matrix, we create another 6x6 weight matrix (randomly initialized) and multiply it with the input.



Using the same process, we perform another matrix multiplication to obtain the V (Value) matrix.

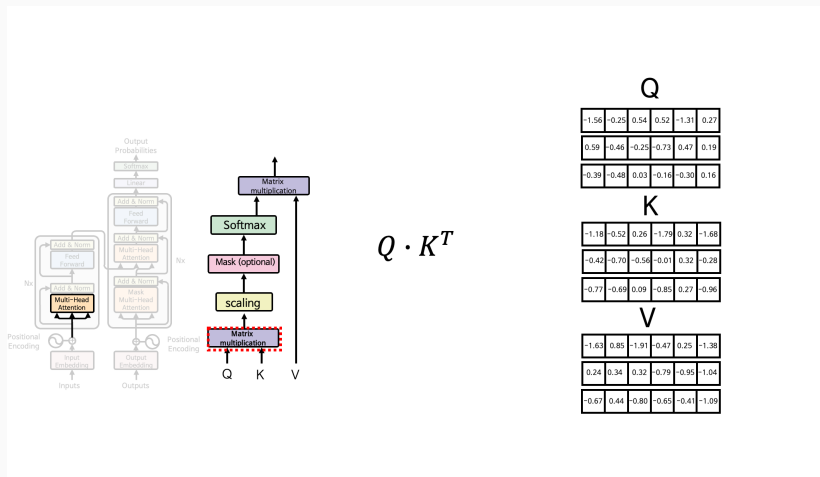


Now, we have the three inputs of the multi-head attention layer — Q (Query), K (Key), and V (Value).





Next, we perform the matrix multiplication between Q and K. The formula for this operation is as follows:



**Q**

-1.56	-0.25	0.54	0.52	-1.31	0.27
0.59	-0.46	-0.25	-0.73	0.47	0.19
-0.39	-0.48	0.03	-0.16	-0.30	0.16

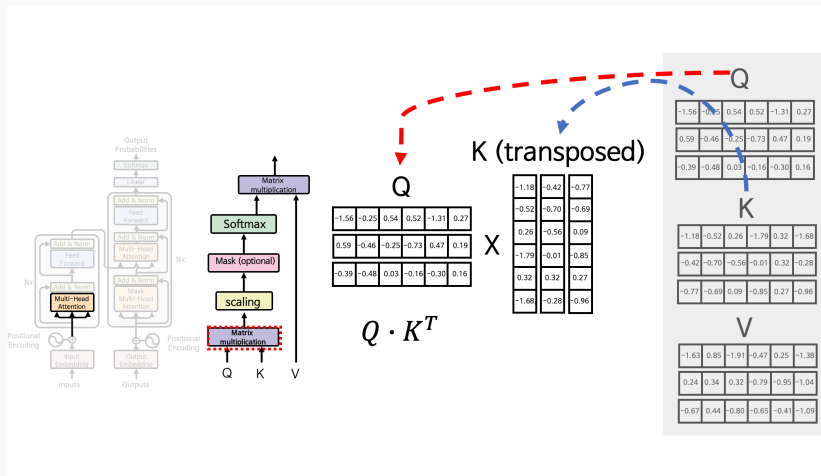
**K**

-1.18	-0.52	0.26	-1.79	0.32	-1.66
-0.42	-0.70	-0.56	-0.01	0.32	-0.28
-0.77	-0.69	0.09	-0.85	0.27	-0.96

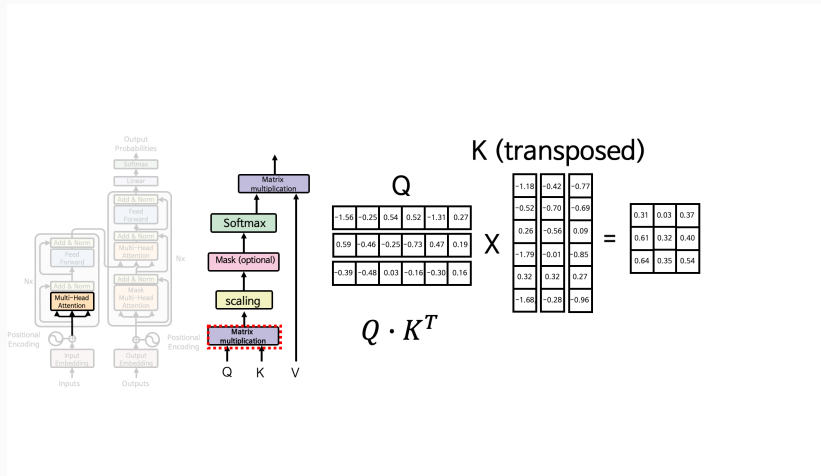
**V**

-1.63	0.85	-1.91	-0.47	0.25	-1.38
0.24	0.34	0.32	-0.79	-0.95	-1.04
-0.67	0.44	-0.80	-0.65	-0.41	-1.09

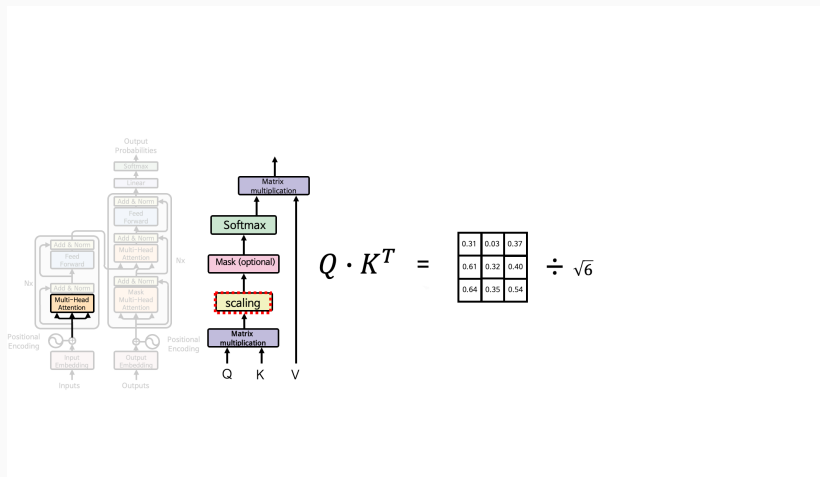
When we plug in the matrix values and calculate,



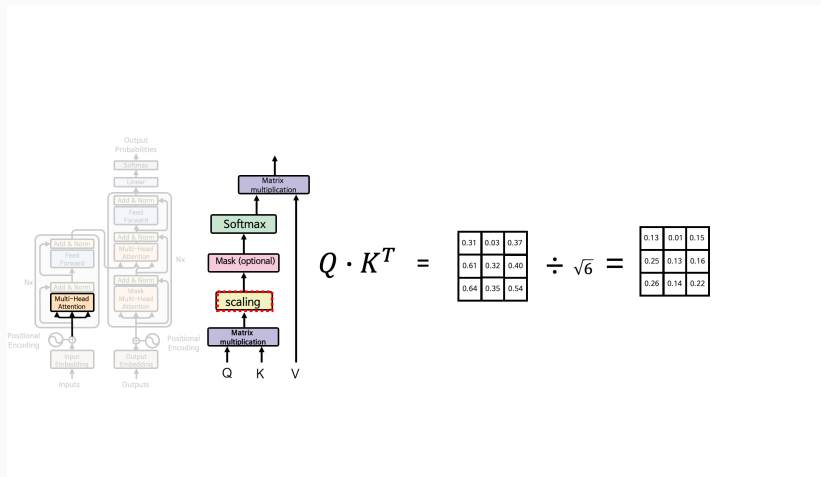
we can obtain the result as follows:



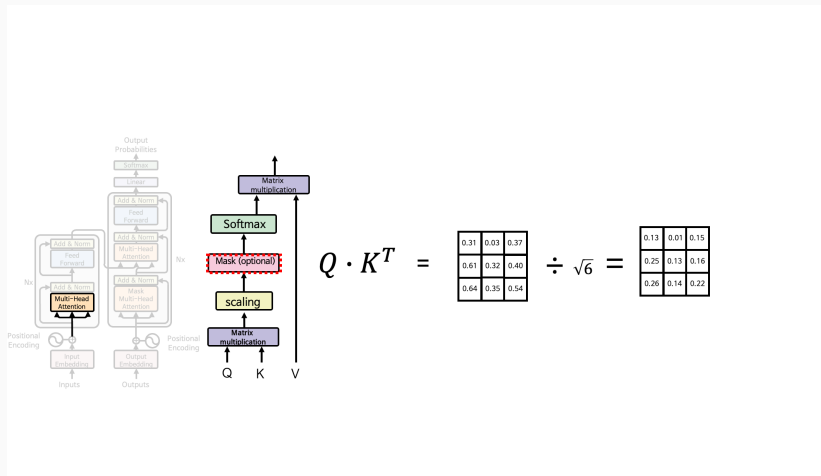
Next, we apply *scaling*, which divides the matrix by  $\sqrt{6}$ , since in this example the value of  $d_{\text{model}}$  is 6.



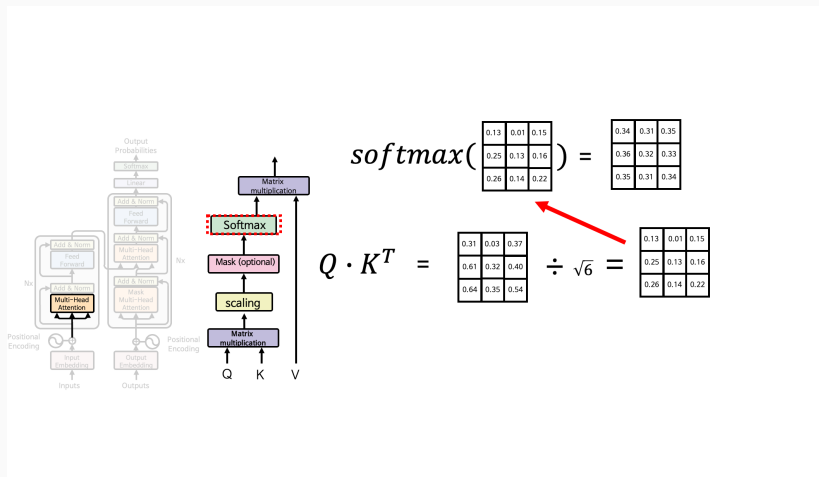
Next, we apply *scaling*, which divides the matrix by  $\sqrt{6}$ , since in this example the value of  $d_{\text{model}}$  is 6.



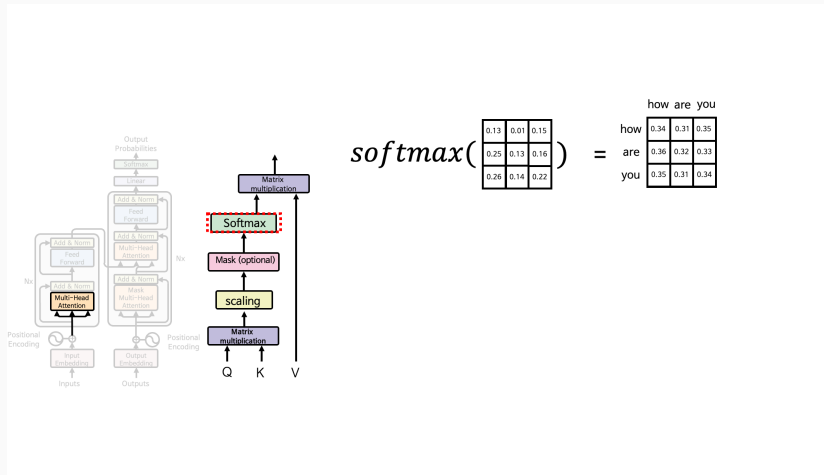
Since the mask layer is not used in the encoder, we will skip it here.



The next softmax layer converts the matrix values into probabilities.

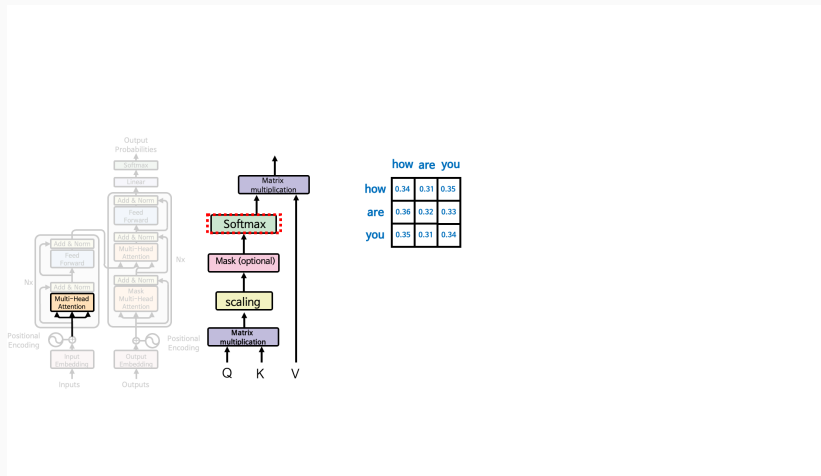


This 3×3 matrix represents the self-attention weights.

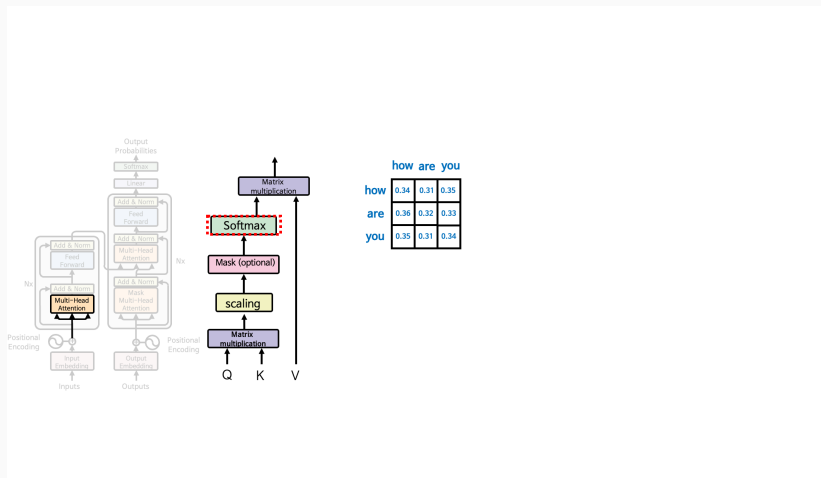




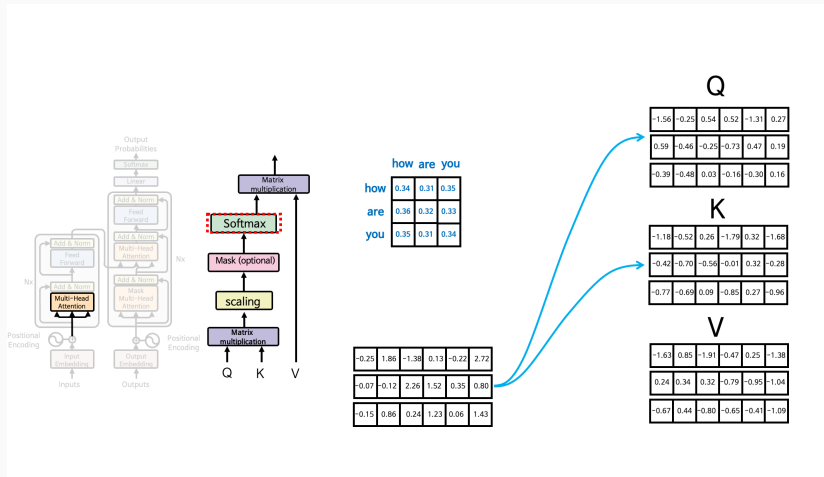
This matrix numerically shows how each word in the input is related to every other word.



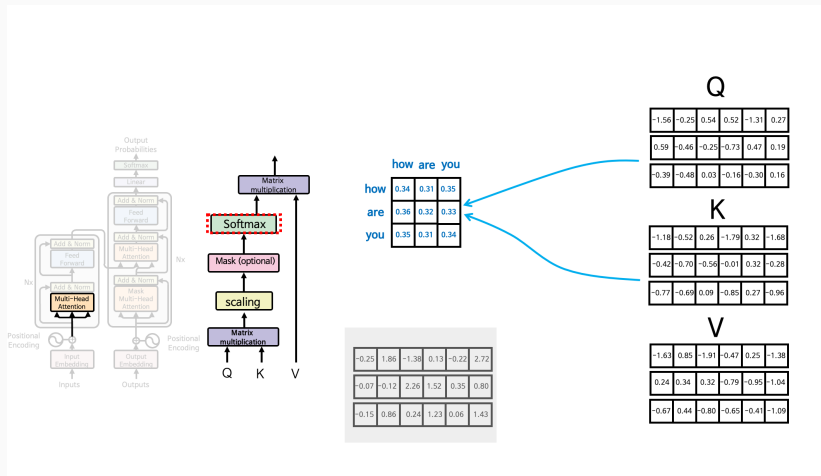
Thus, word pairs with higher relevance receive higher attention values, while those with lower relevance receive smaller values — as the model learns these relationships.



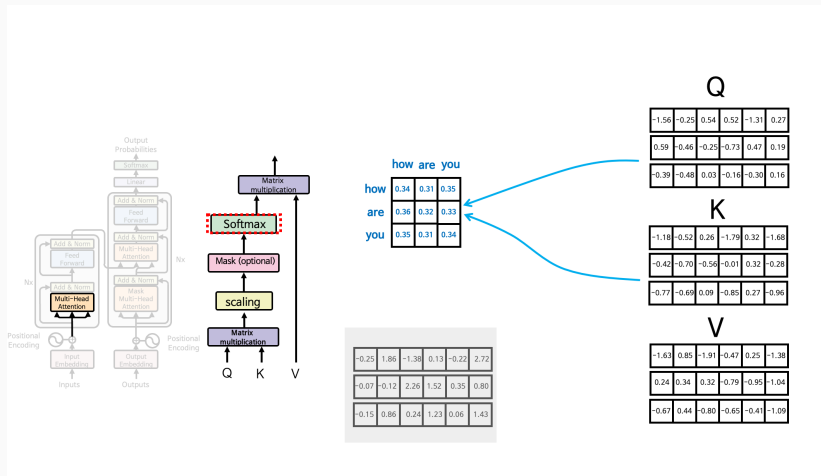
While *traditional* attention models focused on relationships between the input sequence and output sequence, self-attention takes the same input matrix and feeds it into two separate networks to produce the Q and K matrices.



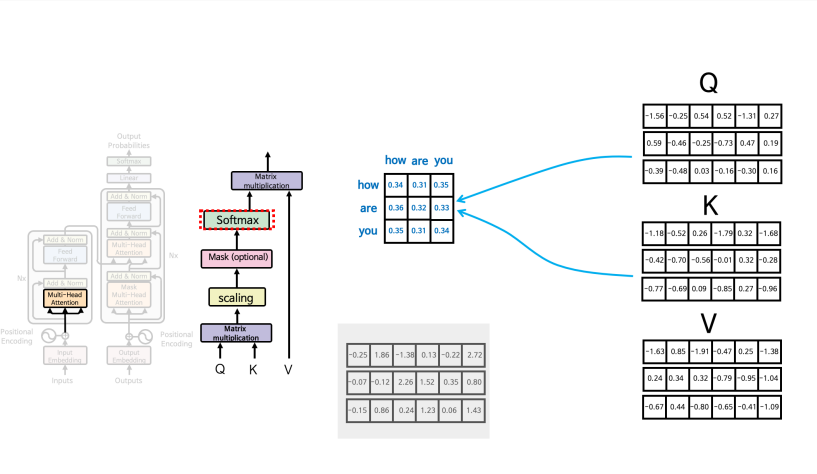
Amazingly, by simply multiplying these two matrices, the model can represent the correlations between words within a sentence.



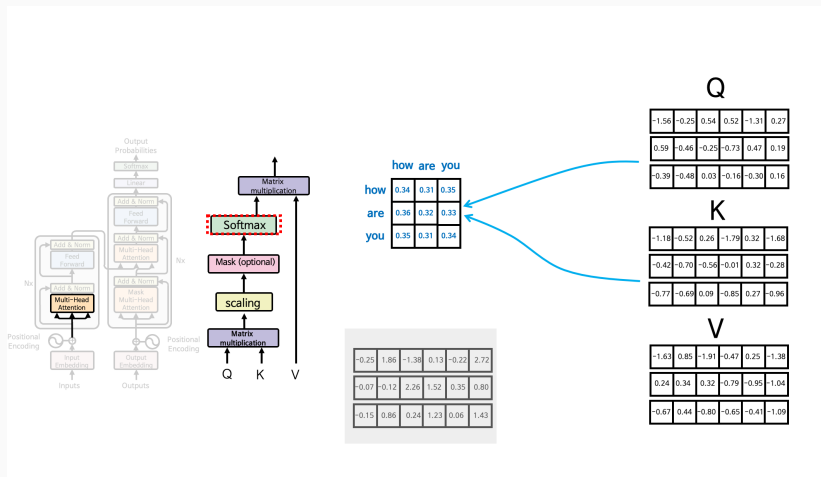
Because the model processes all words in parallel rather than sequentially, it achieves much faster computation.



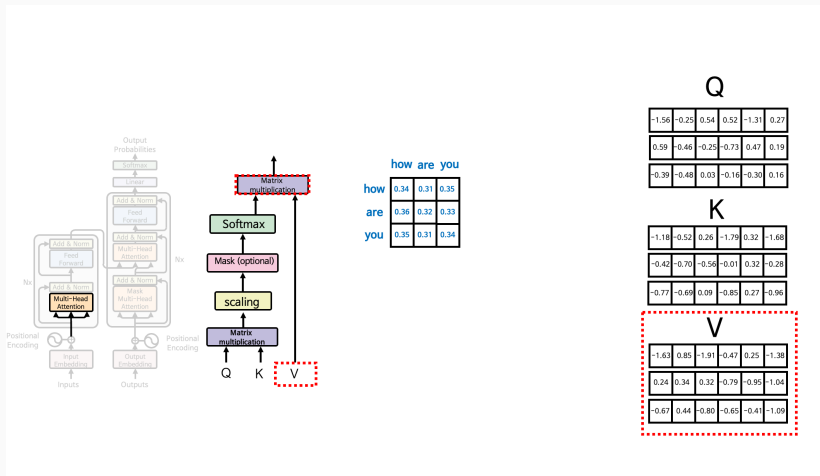
Even for long sentences, the model can calculate the attention between all pairs of words without bias or loss of information.



This self-attention mechanism is the core structure that made today's large language models (LLMs) possible.

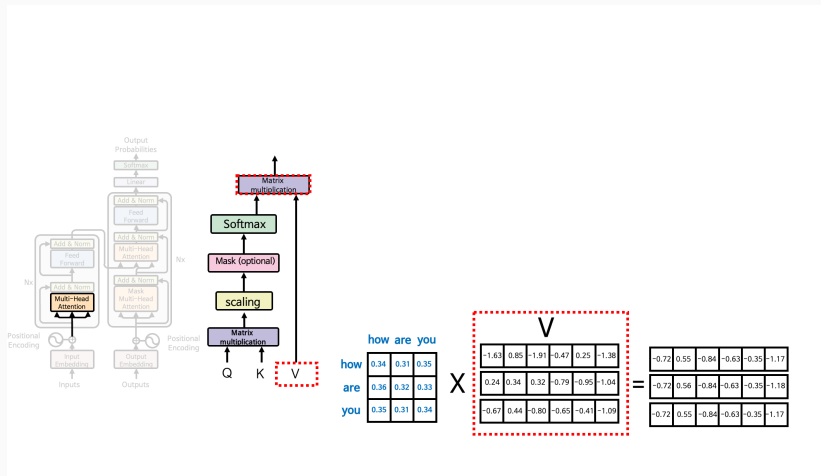


So, what about the final matrix multiplication?

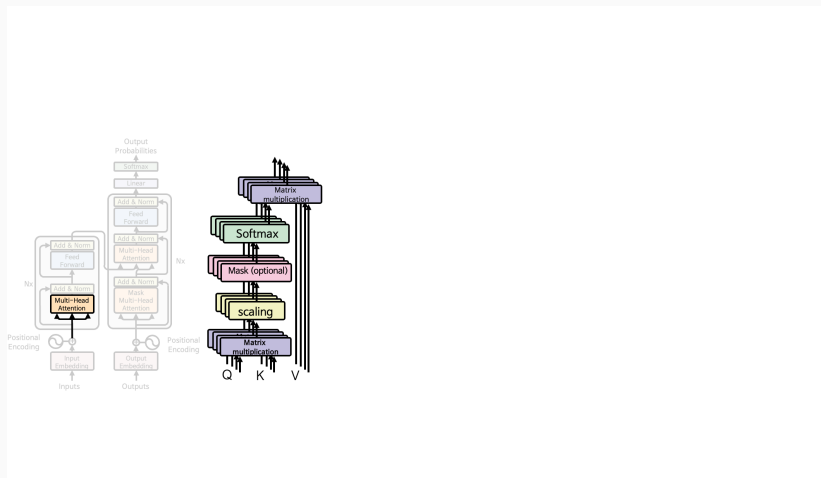




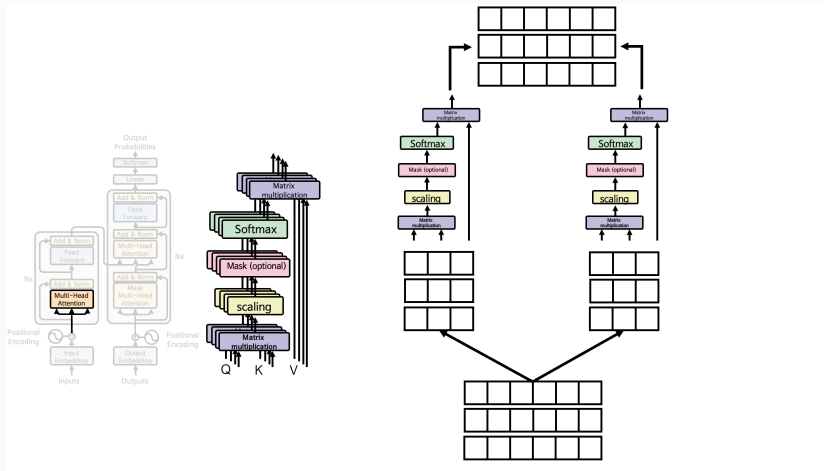
We multiply this by the V matrix to obtain the final output — a self-attended embedding that combines (1) input, (2) positional, and (3) attention information.



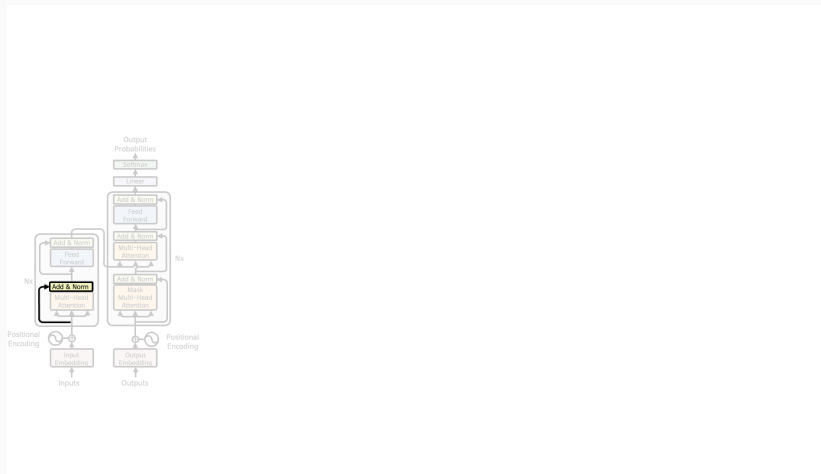
So far, we've assumed a single-head attention for simplicity, but in reality, the Transformer computes multi-head attention — the original paper uses 8 heads.



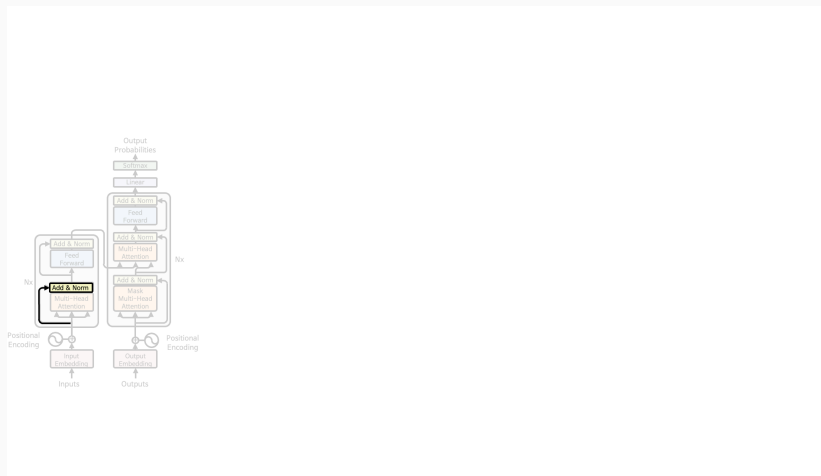
When using *multiple* heads, the model splits the dimensions of Q, K, and V according to the number of heads, performs self-attention separately for each head, then **concatenates** the resulting matrices and passes them through a fully connected layer to produce the final output of multi-head attention.



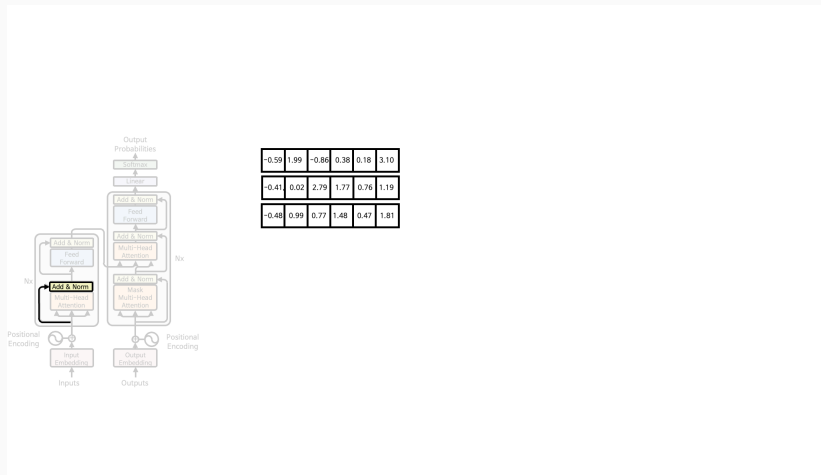
Next, we move on to the **addition and normalization** layer.



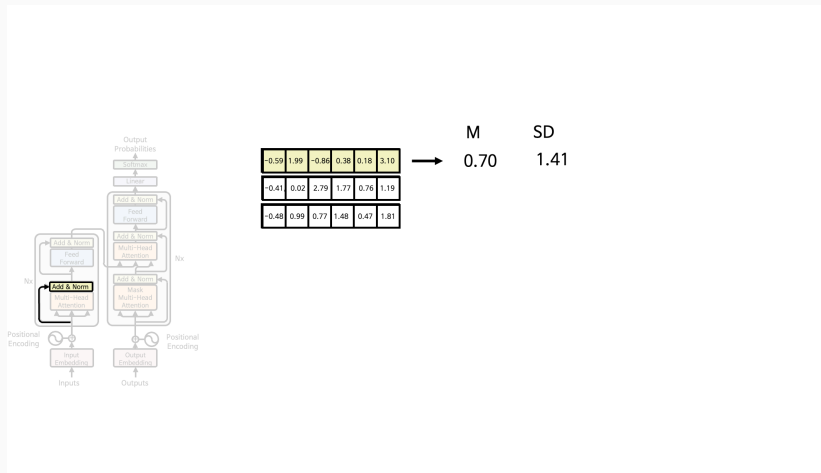
The **addition** step adds the multi-head output matrix to the initial input + positional embedding.

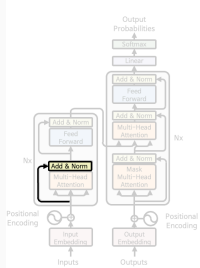


Let's assume that the resulting matrix looks like this:



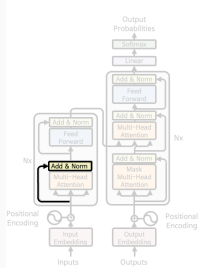
Next, for the normalization step, we first calculate the mean and standard deviation for each column.



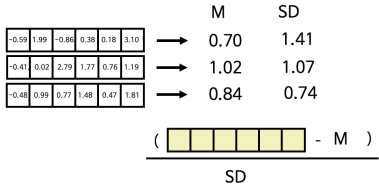
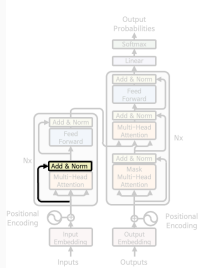


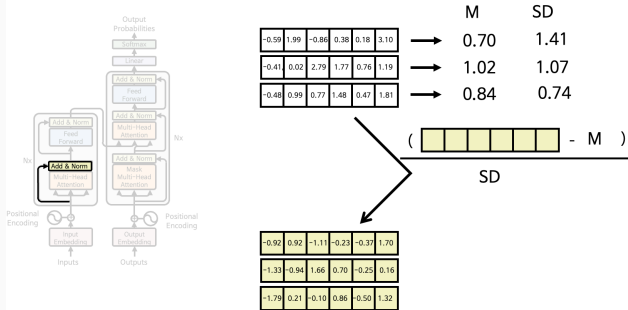
							M	SD
	-0.59	1.99	-0.86	0.38	0.18	3.10	→ 0.70	1.41
	-0.41	0.02	2.79	1.77	0.76	1.19	→ 1.02	1.07
	-0.48	0.99	0.77	1.48	0.47	1.81		



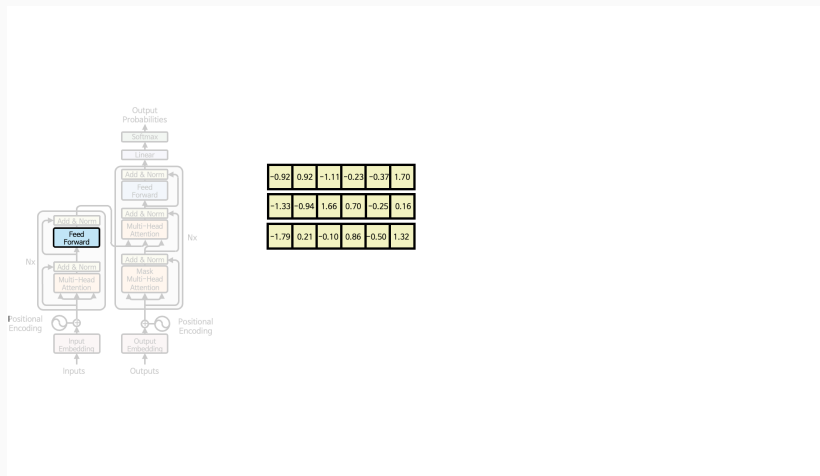


						M	SD	
	-0.59	1.99	-0.86	0.38	0.18	3.10	→ 0.70	1.41
	-0.41	0.02	2.79	1.77	0.76	1.19	→ 1.02	1.07
	-0.48	0.99	0.77	1.48	0.47	1.81	→ 0.84	0.74



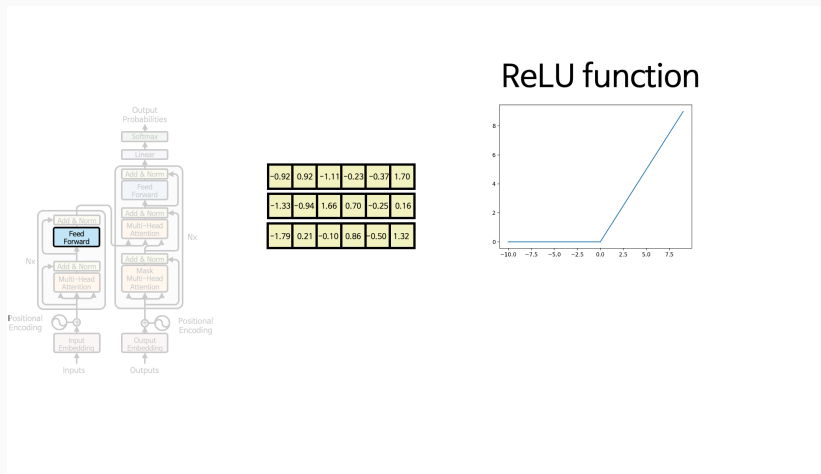


Next is the **Feed-Forward** layer.

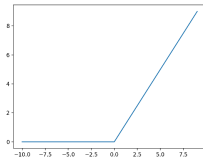


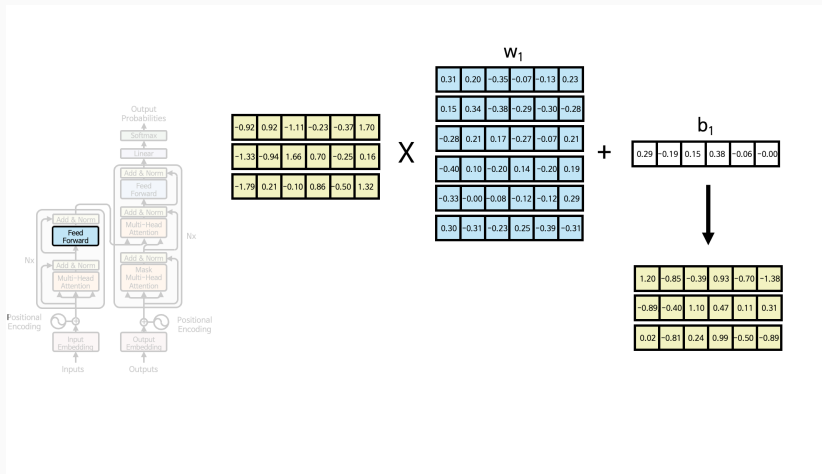
The Feed-Forward layer is a simple neural network consisting of **two layers** and using **ReLU()** as the activation function.

$$\text{FFN}(x) = \text{ReLU}(xW_1 + b_1)W_2 + b_2$$

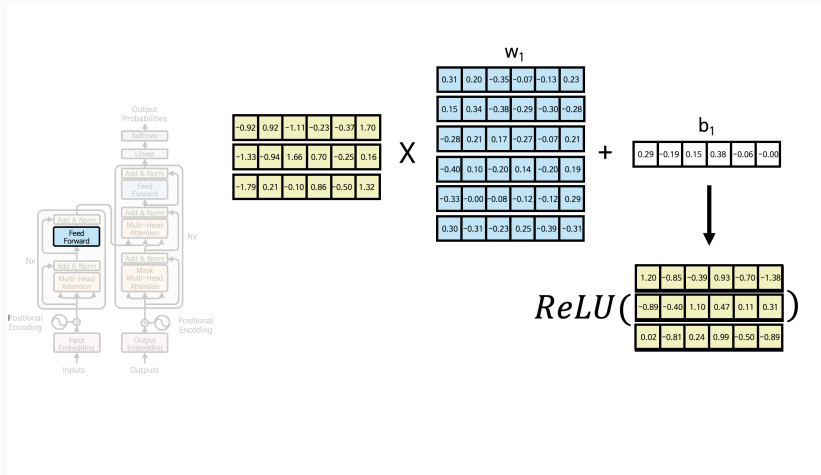


ReLU function





# Applying the ReLU activation function...



# Applying the ReLU activation function...

The diagram illustrates the internal structure of a Transformer decoder block. It shows the flow from **Inputs** through **Input Embedding** and **Positional Encoding** into a stack of  $N \times$  layers. Each layer contains a **Multi-Head Attention** sub-layer followed by an **Add & Norm** operation. This is followed by a **Feed Forward** sub-layer, another **Add & Norm** operation, and a **Multi-Head Attention** sub-layer that also receives **Mask Multi-Head Attention** as input. A final **Add & Norm** operation leads to the **Outputs**, which are processed by **Output Embedding** and **Positional Encoding**. The output then passes through a **Linear** layer and a **Softmax** layer to produce **Output Probabilities**.

The main calculation shown is the application of the ReLU activation function to the product of a weight matrix  $w_1$  and a bias vector  $b_1$ . The weight matrix  $w_1$  is a 6x6 matrix:

0.31	0.20	-0.35	-0.07	-0.13	0.23
0.15	0.34	-0.38	-0.29	-0.30	-0.28
-0.28	0.21	0.17	-0.27	-0.07	0.21
-0.40	0.10	-0.20	0.14	-0.20	0.19
-0.33	-0.00	-0.08	-0.12	-0.12	0.29
0.30	-0.31	-0.23	0.25	-0.39	-0.31

The bias vector  $b_1$  is a 1x6 vector:

0.29	-0.19	0.15	0.38	-0.06	-0.00
------	-------	------	------	-------	-------

The resulting vector before the ReLU function is:

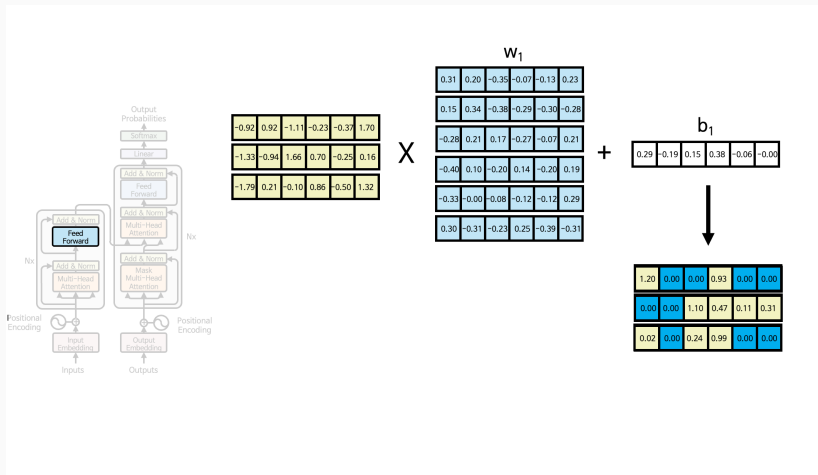
1.20	-0.85	-0.39	0.93	-0.70	-1.38
-0.89	-0.40	1.10	0.47	0.11	0.31
0.02	-0.81	0.24	0.99	-0.50	-0.89

The ReLU function is applied to this vector, as shown in the graph below:

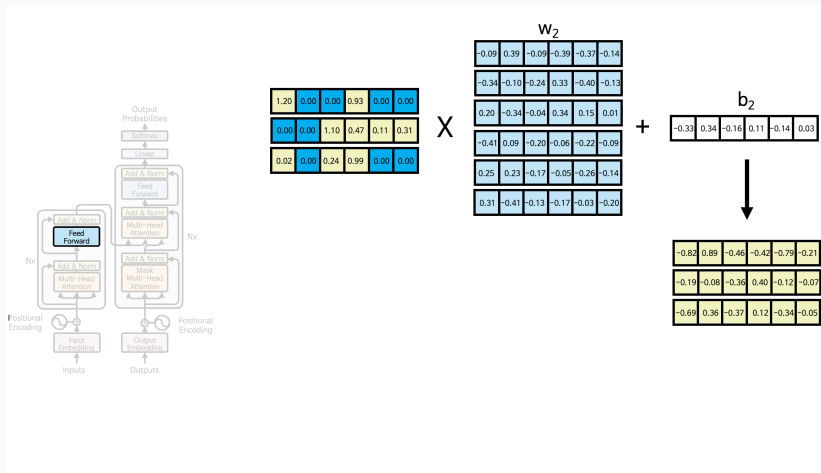
The graph shows the ReLU function:  $y = \max(0, x)$ . The x-axis ranges from -10.0 to 7.5, and the y-axis ranges from 0 to 8. The function is zero for  $x \leq 0$  and follows the line  $y = x$  for  $x > 0$ .



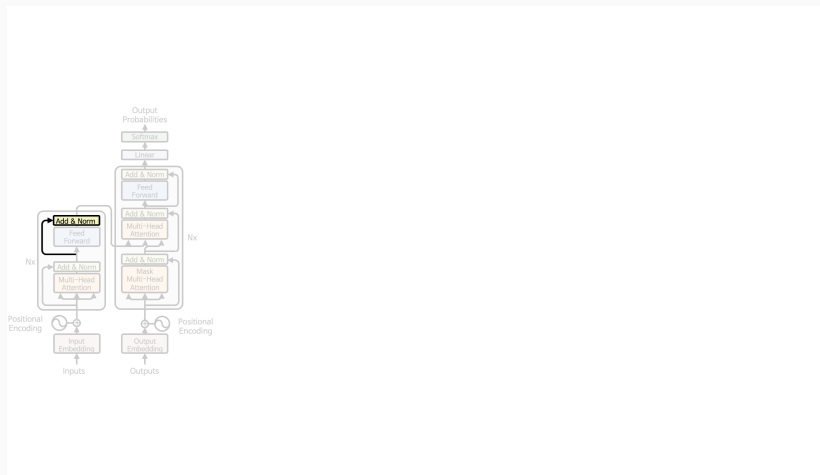
Negative values become 0.



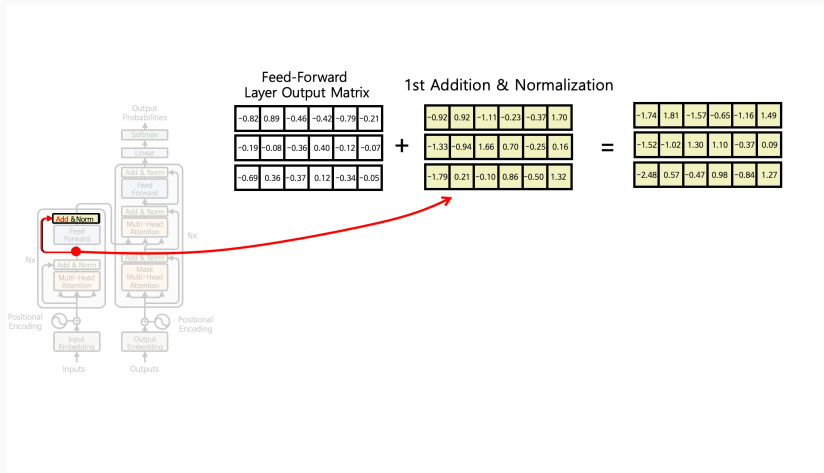
Then, by multiplying with the second layer's weights and biases, we obtain the final output of the Feed-Forward layer. The purpose of this layer is to increase **non-linearity**, allowing the network to better process and distinguish complex patterns (recall multi-layer perceptrons!)



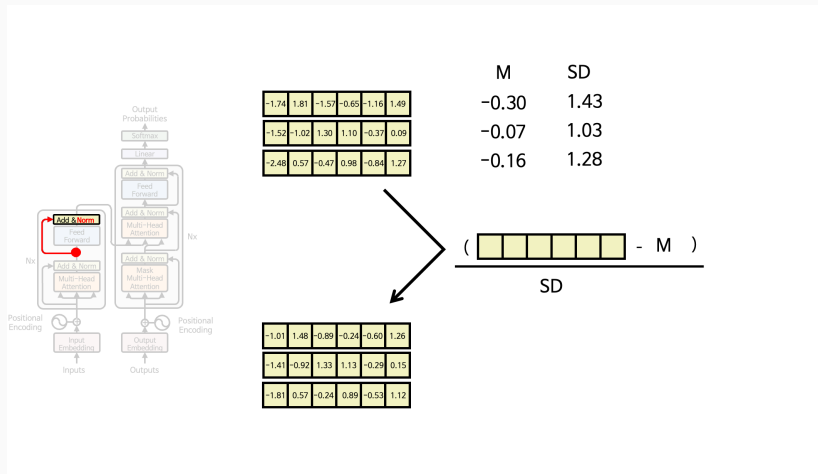
Next, we repeat the same **Add & Normalize** process.



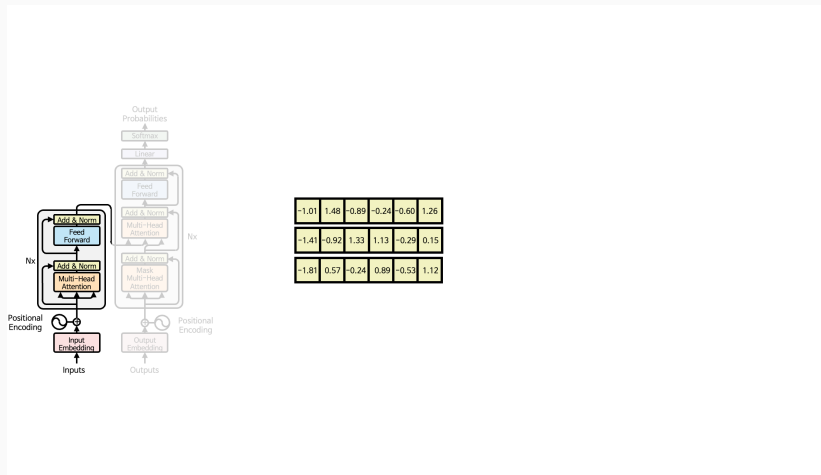
As before, we first calculate the sum of the two matrices.



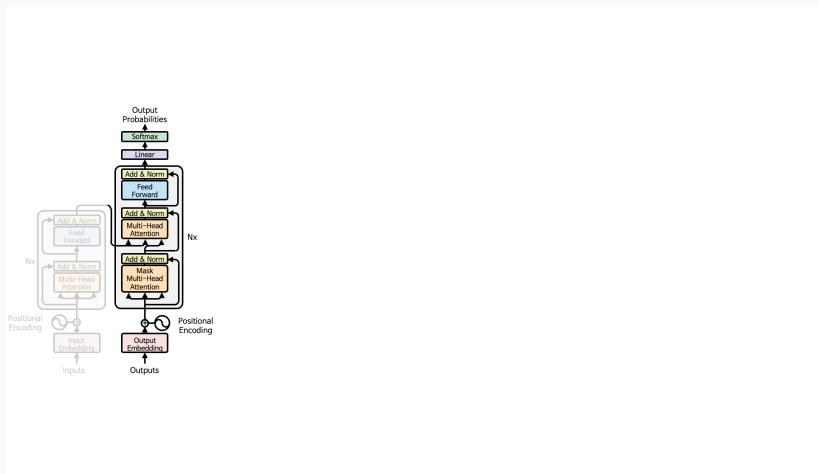
Then, we normalize the matrix again:



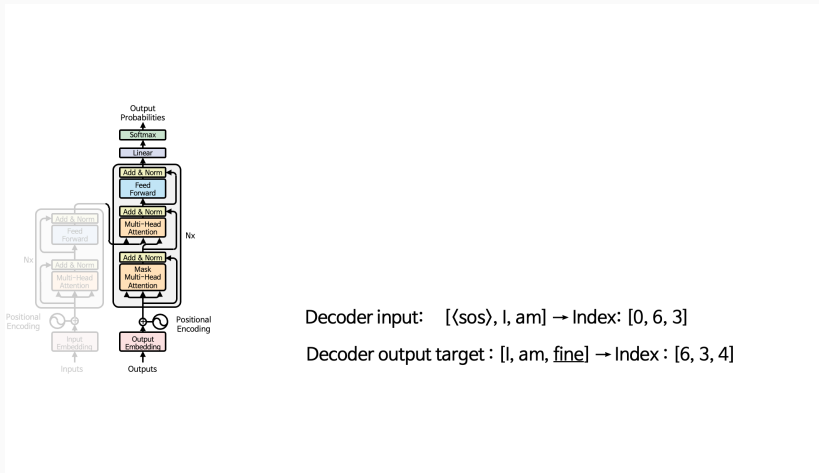
This resulting matrix is the **final output of the encoder**.



Now, let's move on to the **Decoder**.



Just like the encoder, the decoder first performs **word encoding** and **positional encoding** of the output (target) sequence.



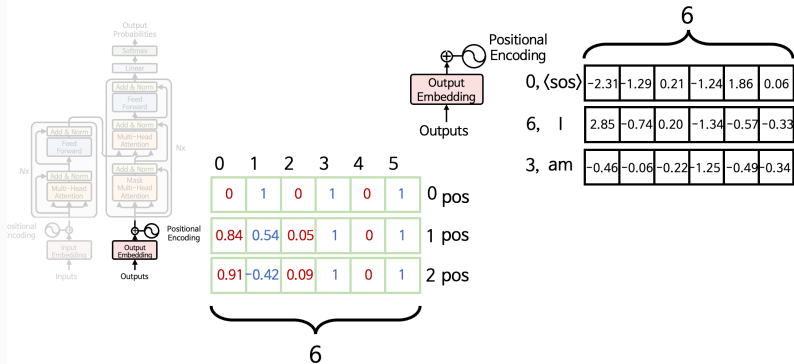
Decoder input: [⟨sos⟩, I, am] → Index: [0, 6, 3]

Decoder output target : [I, am, fine] → Index : [6, 3, 4]

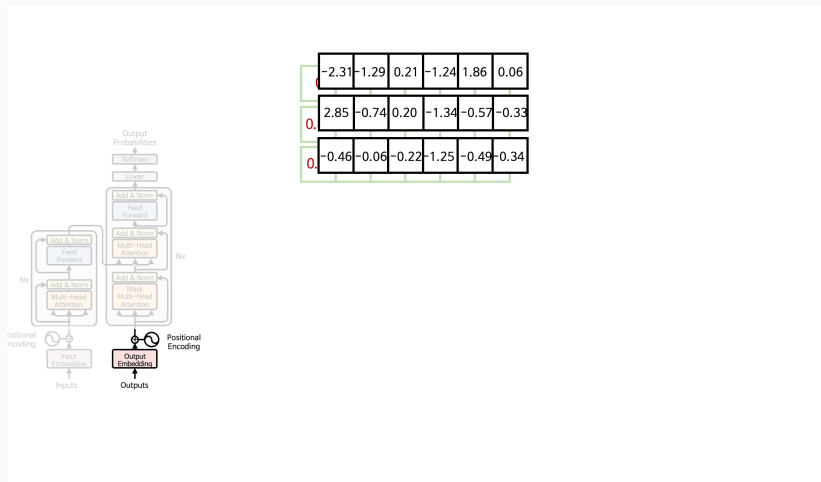


The process of word encoding and positional encoding is identical to that in the encoder. In fact, the positional encoding values can be reused from the encoder.

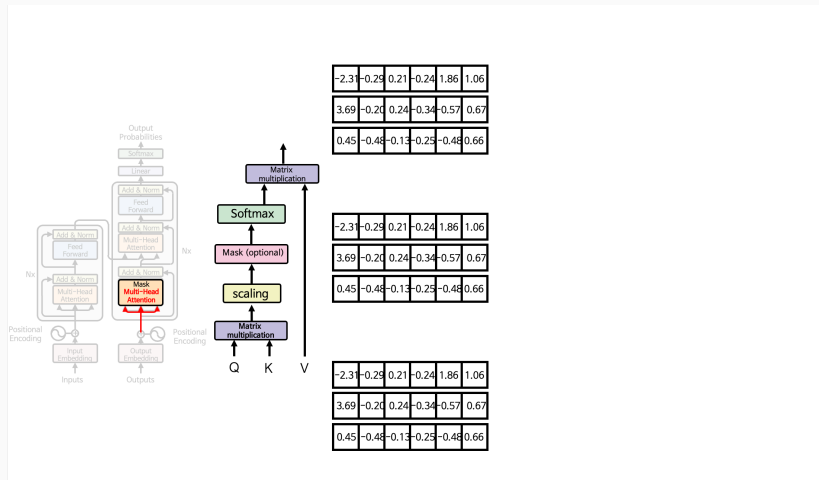
$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) \quad PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$



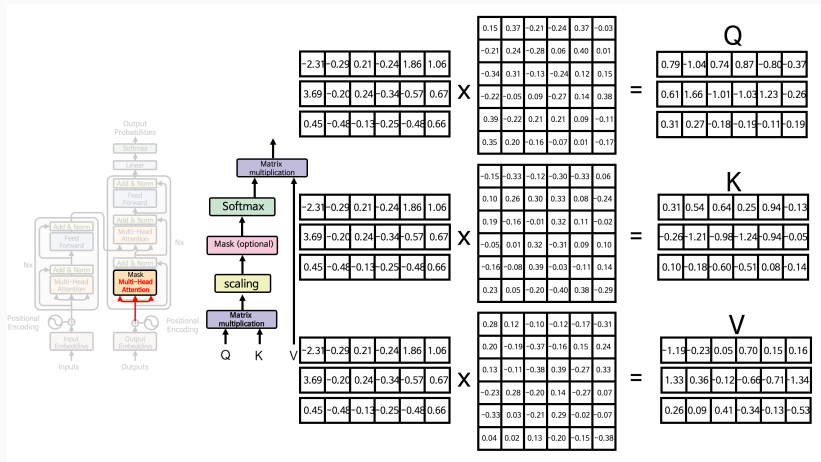
Here's how it looks:



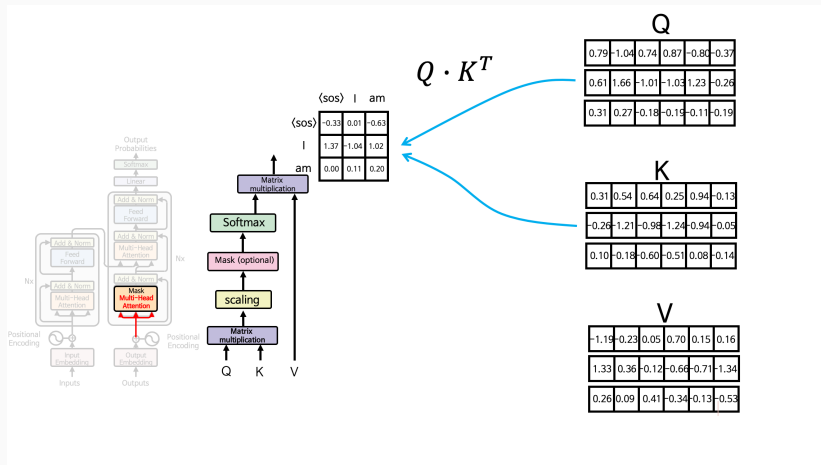
The **Masked Multi-Head Attention** operation in the decoder is almost the same as in the encoder.



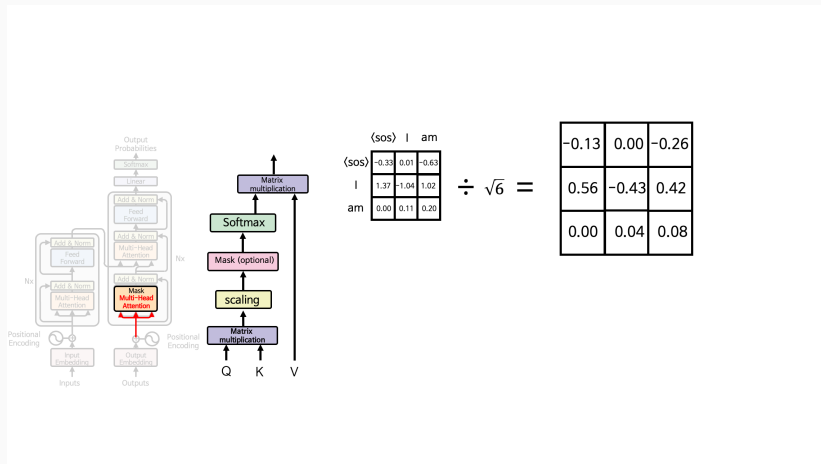
The **Masked Multi-Head Attention** operation in the decoder is almost the same as in the encoder.



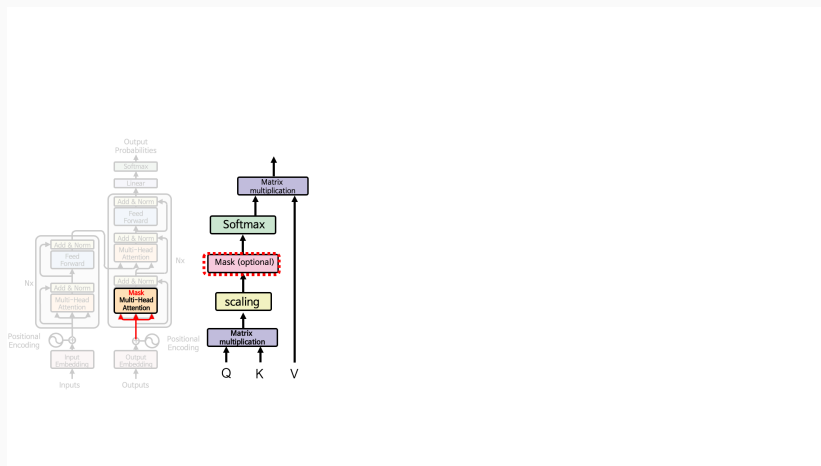
We multiply the Q and K matrices to create the attention matrix.



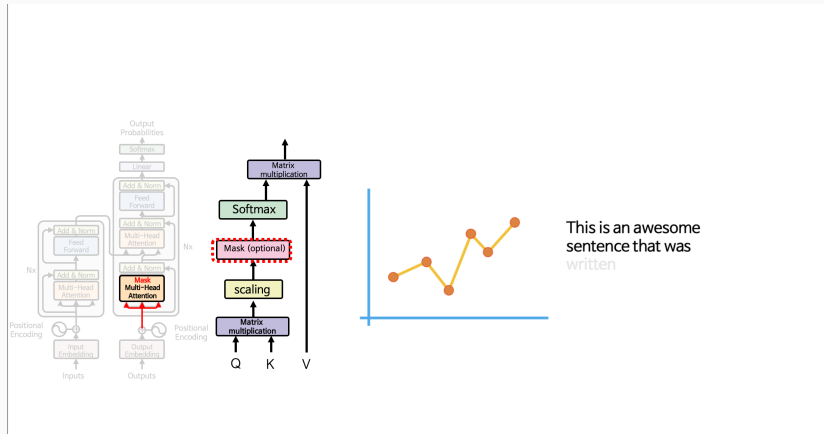
We then scale the matrix by dividing it by  $\sqrt{6}$ , just as before, so that the range of values changes accordingly.



Now, let's explore the key concept of the decoder's multi-head attention — **masking**.

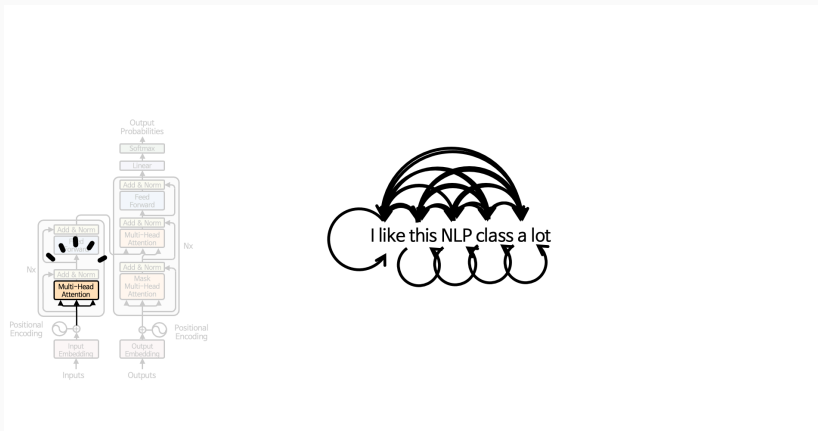


The goal of the transformer decoder is to generate the output word sequence, one token at a time (recall: language modeling).

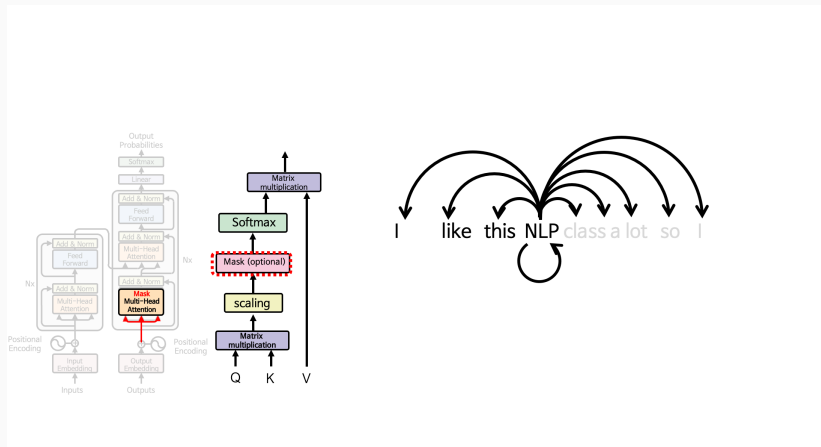




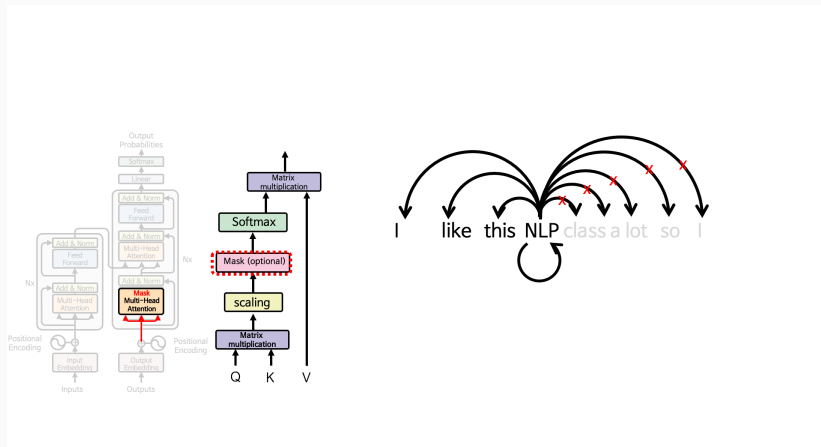
While the encoder needs to consider all tokens in the input sentence to understand the full meaning,



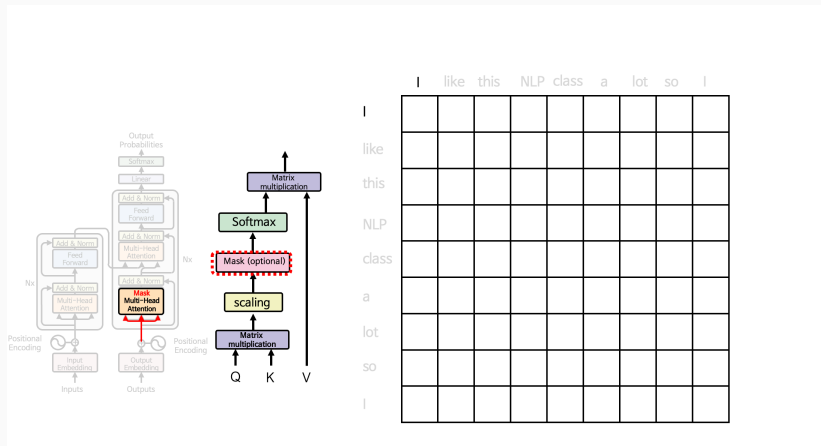
the decoder generates output **one word at a time**.



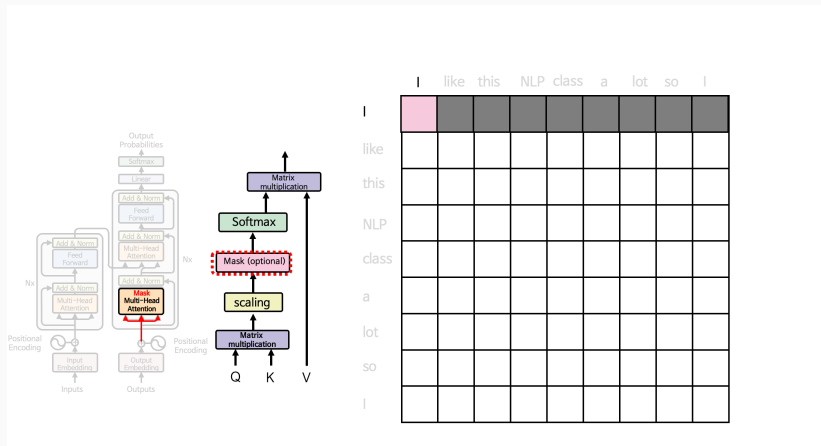
Therefore, it's natural that the decoder should NOT attend to words that haven't been generated yet.

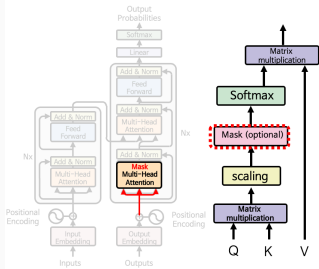


To reflect this characteristic, the decoder applies a masking mechanism during training.



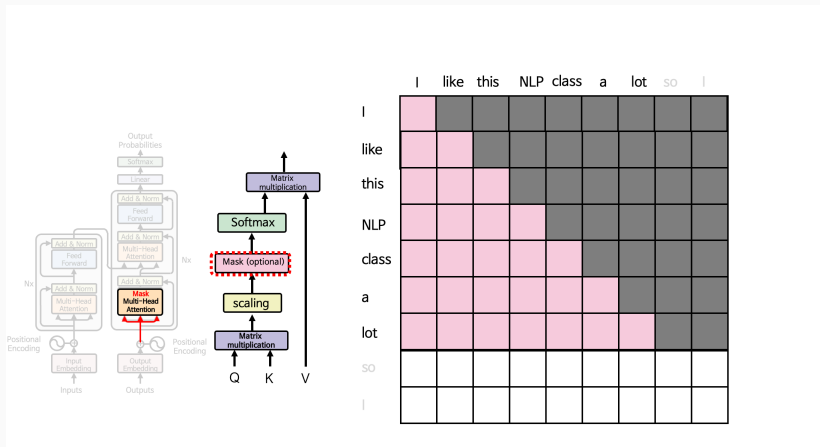
The key idea is to hide future tokens so they do not affect the current prediction.

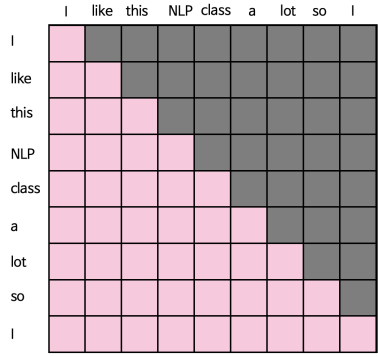
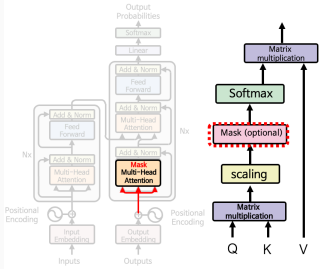




I like this NLP class a lot so I

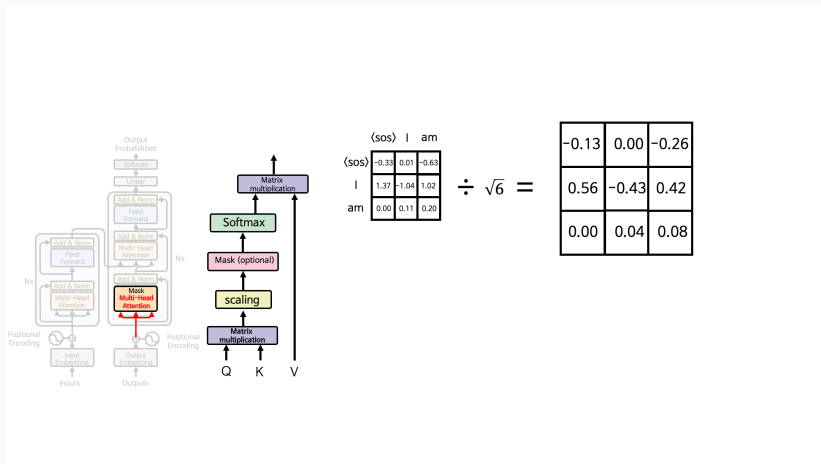
I								
like								
this								
NLP								
class								
a								
lot								
so								
I								



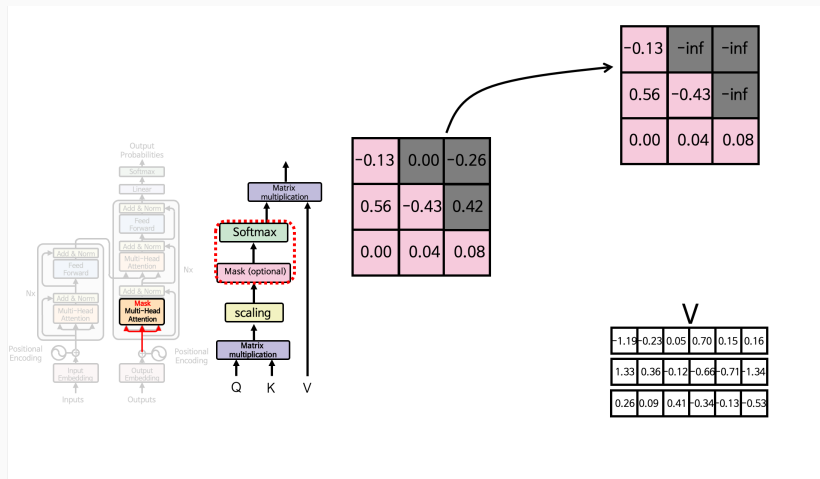




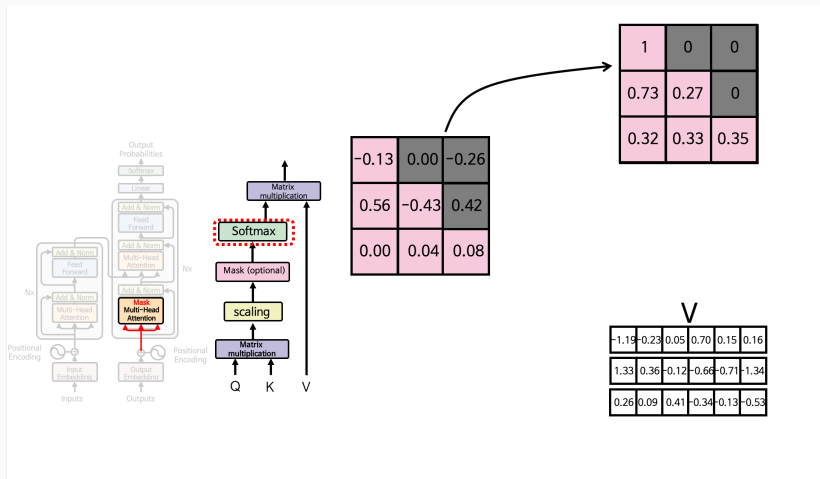
When this masking algorithm is applied to the attention matrix, we get:



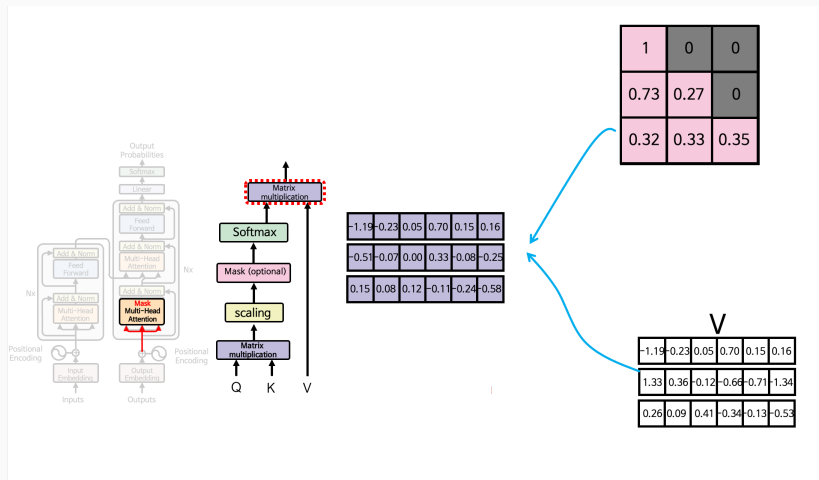
We add  $-\text{inf}$  to the masked positions because, after passing through the softmax layer,  $-\text{inf}$  becomes 0, effectively eliminating attention to those positions.



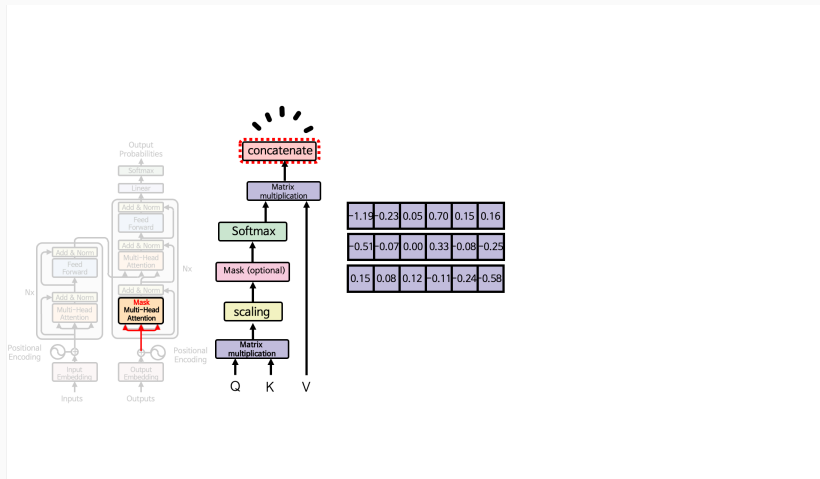
Feeding this matrix into the softmax layer gives us:



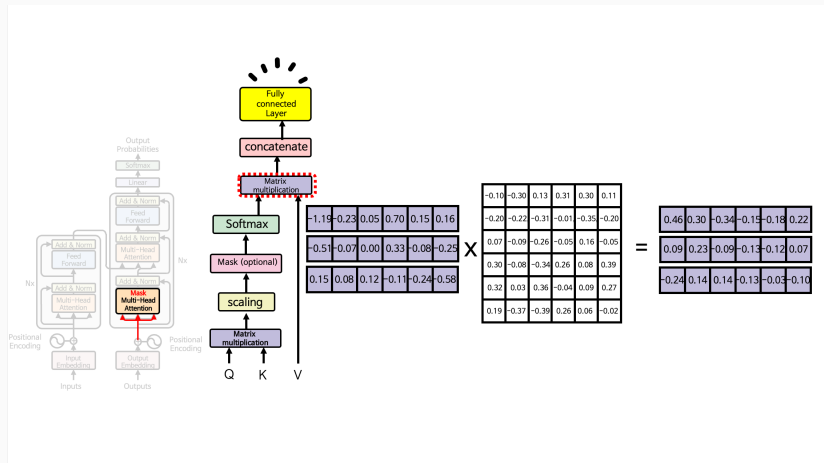
Then we multiply the two matrices as follows:



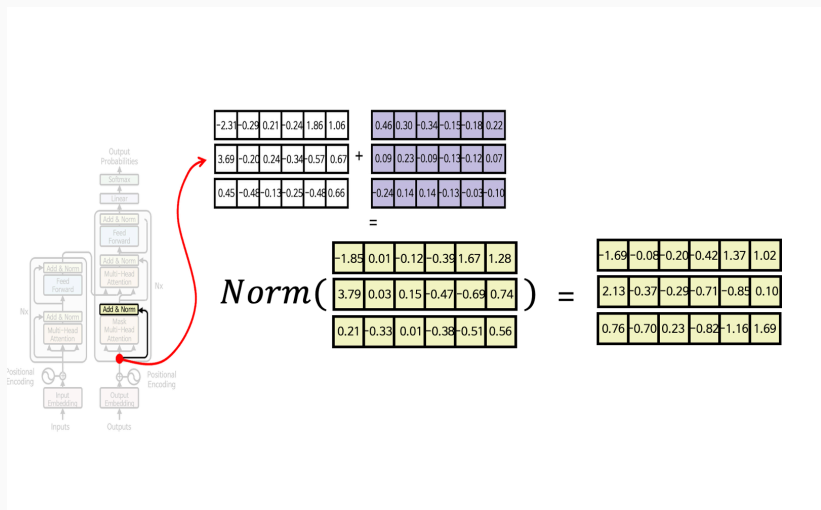
Next, we concatenate the resulting matrices (if needed):



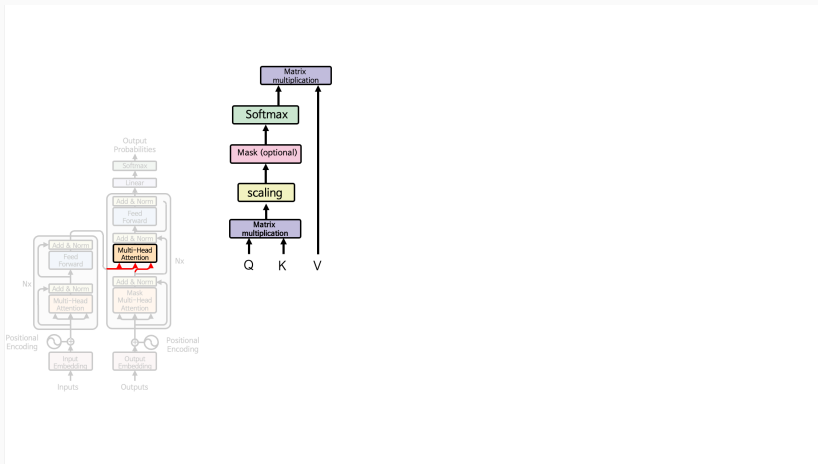
We then multiply again to produce the final matrix of the masked multi-head attention.



We then apply the same Add & Normalize process again.

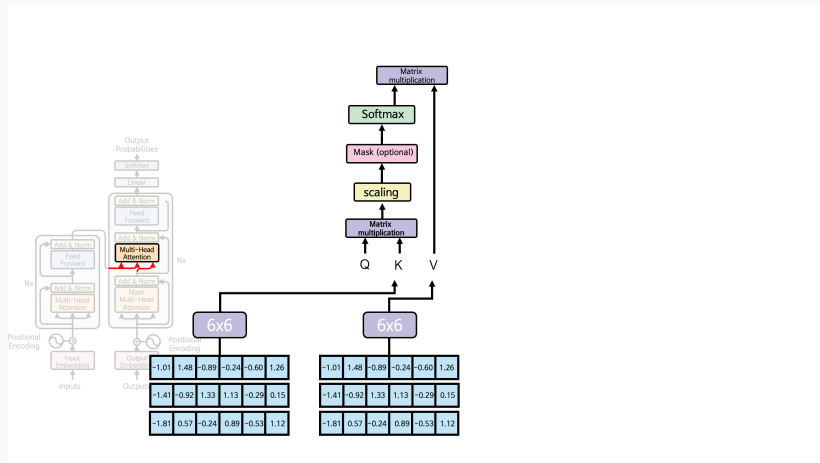


The decoder's second multi-head attention operates the same way as the encoder's, except for the inputs.

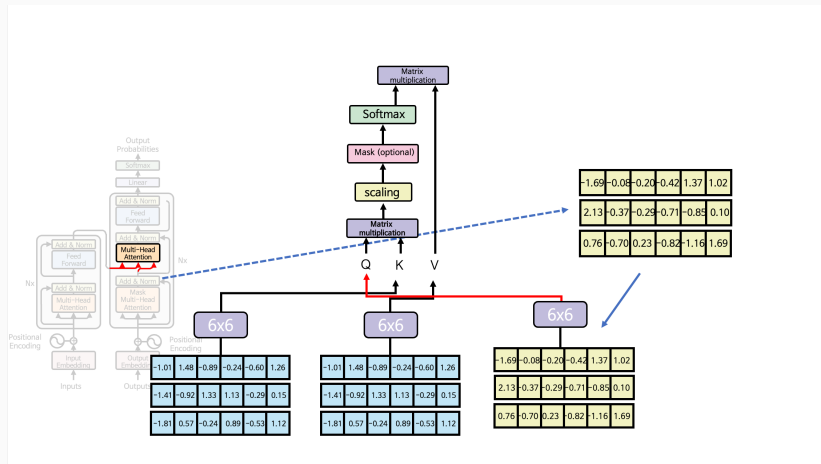




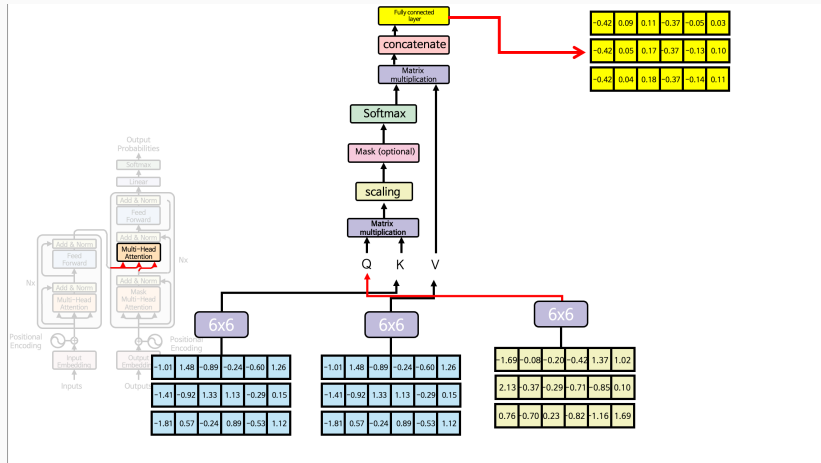
Here, the values of  $K$  and  $V$  are derived from the **encoder's final output**, multiplied by a  $6 \times 6$  matrix.



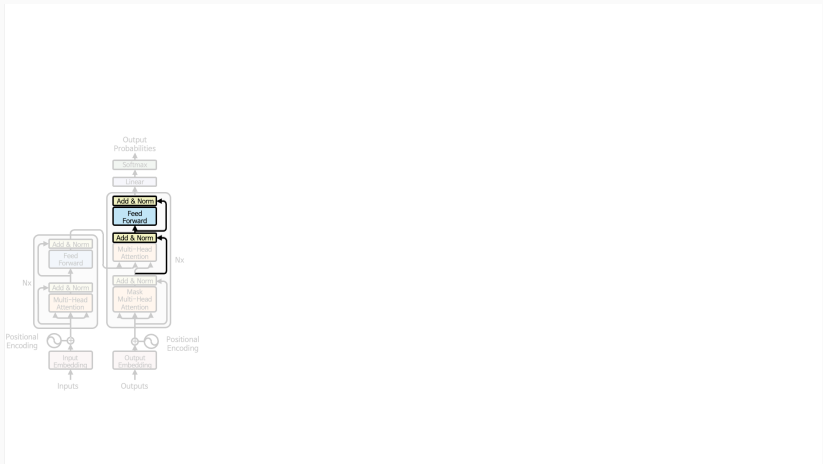
The value of  $Q$  comes from the output of the decoder's previous Add & Norm layer. **In other words, the decoder determines which parts of the encoder's output (K, V) to attend to, based on the context it has generated so far (Q)-similar to attention in RNNs.**



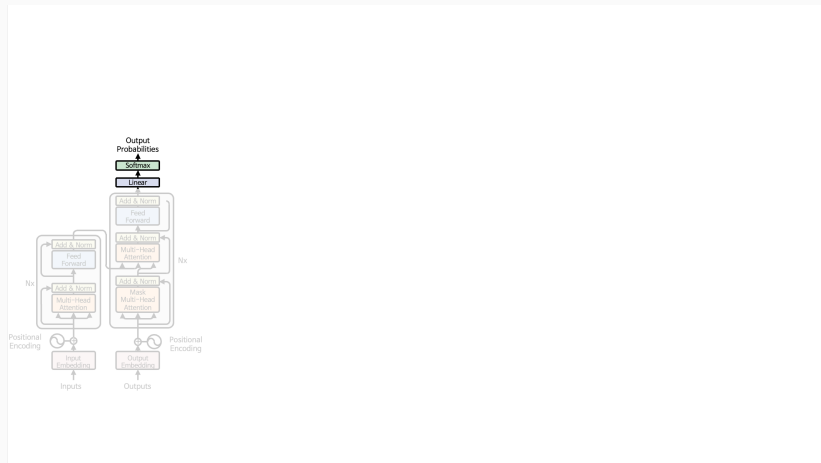
Assume the final output of this multi-head attention looks as follows:



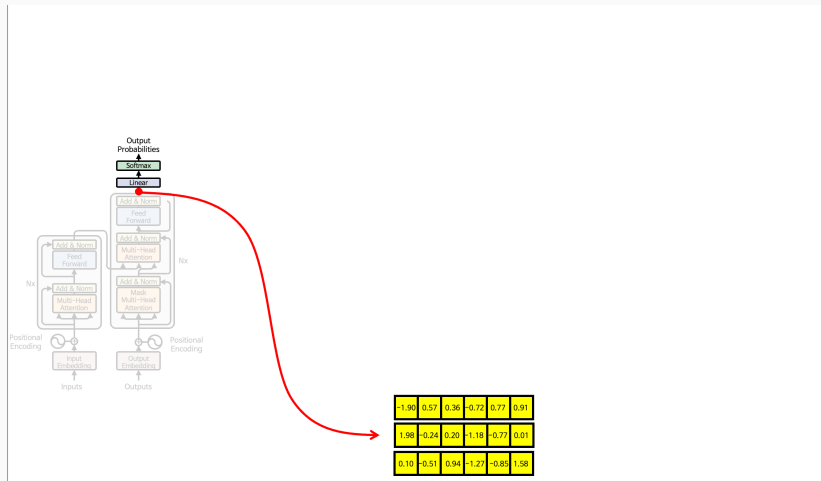
The same computational process is repeated, so details are omitted here.



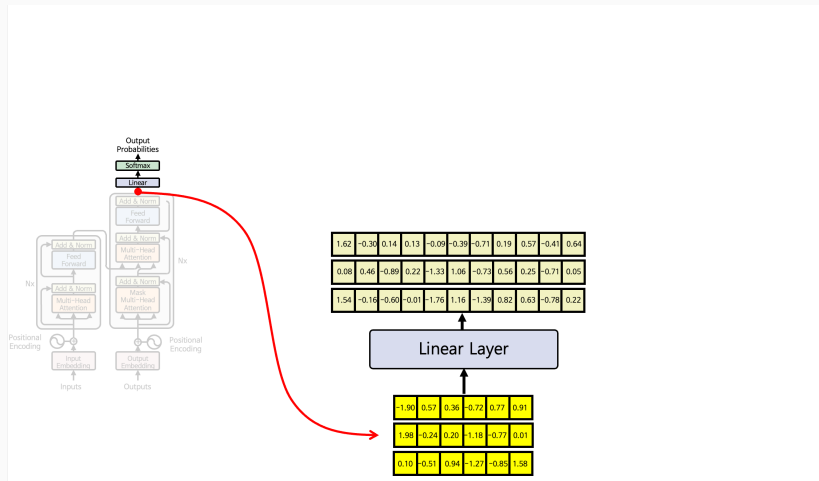
Finally, the last linear and softmax layers produce the final output probabilities.



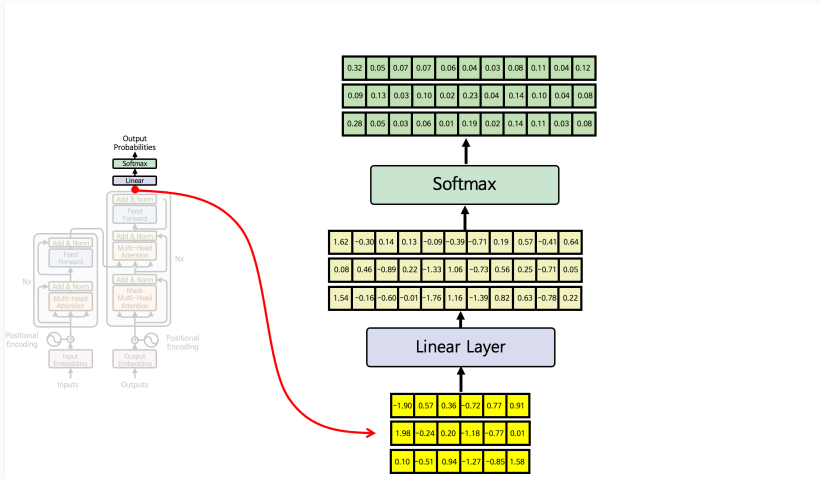
If the decoder output matrix at this stage looks like this:



The linear layer projects it back to the full vocabulary size (=11).

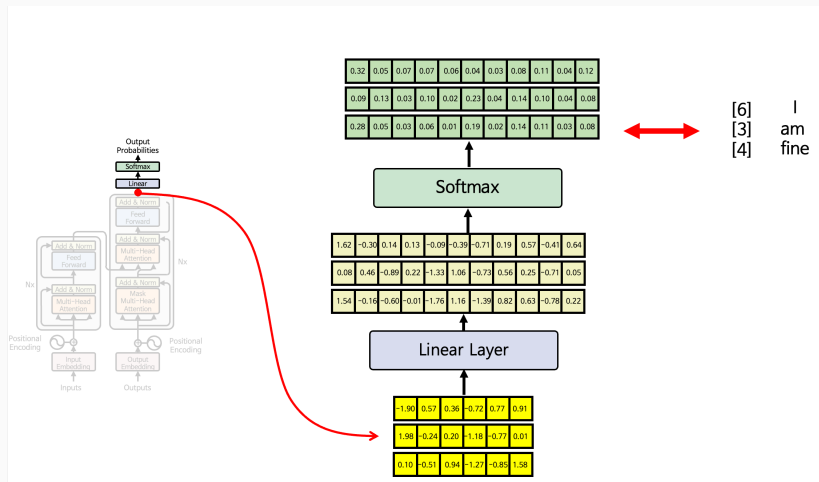


The softmax function then produces the final output probabilities.

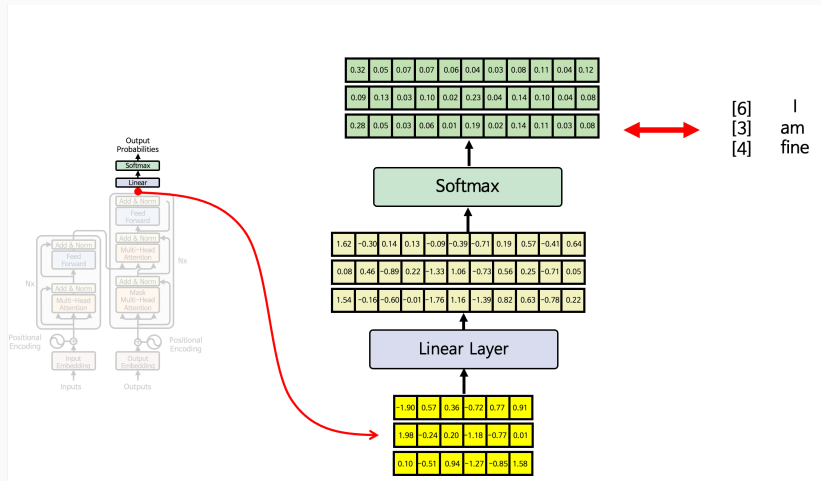


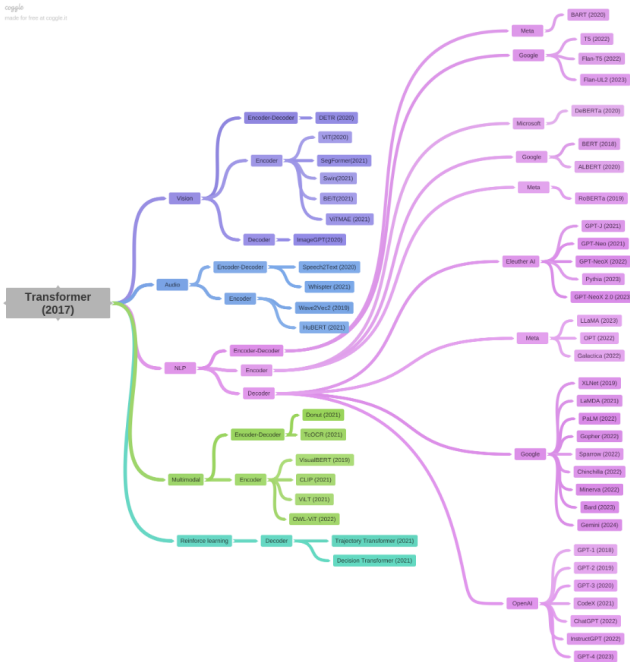


Finally, the predicted outputs are compared with the ground-truth labels (e.g., 6, 3, 4).



Using a **loss function** (e.g., cross-entropy) and **backpropagation**, the model updates all weight parameters across every layer — this is the learning process of the Transformer.





## **Wrap up**

---

- We explored the architecture of the Transformer model.

## Wrap-up

- We explored the architecture of the Transformer model.
- Do we need to train this model entirely from scratch?

## Wrap-up

- We explored the architecture of the Transformer model.
- Do we need to train this model entirely from scratch?
- **No.** We can take advantage of powerful *pre-trained language models* that have already learned general language patterns from massive datasets.

## Wrap-up

- We explored the architecture of the Transformer model.
- Do we need to train this model entirely from scratch?
- **No.** We can take advantage of powerful *pre-trained language models* that have already learned general language patterns from massive datasets.
- We will experiment with these models using *Hugging Face* (open-source library) in the lab session.



# Attention in transformers

Would be a nice recap:

<https://www.youtube.com/watch?v=eM1x5fFNoYc> Any hints from the video clip?