

10. LLMs in 2025

LING-581-Natural Language Processing 1

Instructor: Hakyung Sung

October 28, 2025

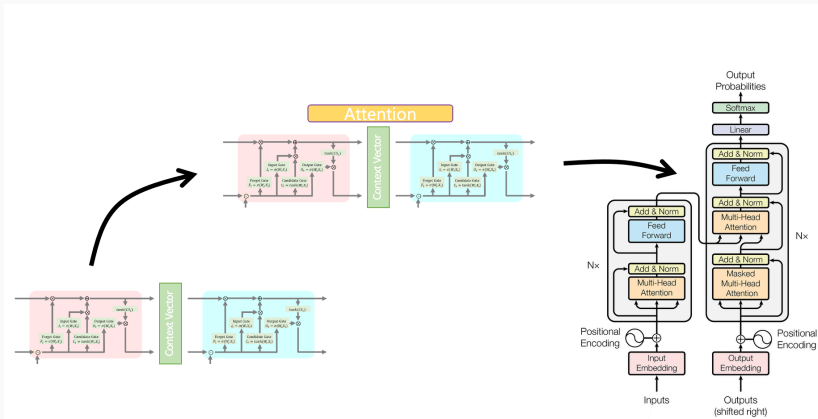
Table of contents

1. Review: Transformer
2. Different transformers
3. LLMs in 2025
4. Wrap up

Review: Transformer

Recall

Towards the development of the **Transformer** model.



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.
- 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.
- 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.

Transformer and self-attention

- In Transformers, attention occurs within a single sentence — all words attend to all words in the previous layer.
(Self-attention + Cross-attention)

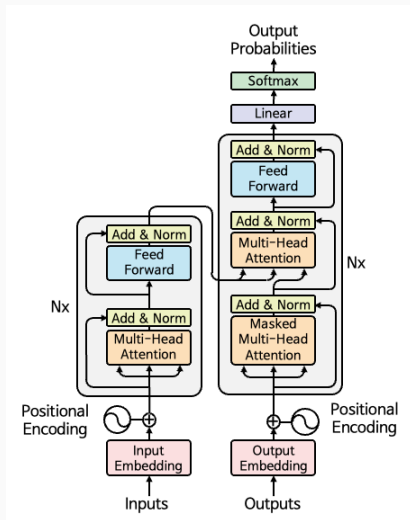
Transformer and self-attention

- In Transformers, attention occurs within a single sentence — all words attend to all words in the previous layer.
(Self-attention + Cross-attention)
- As a result, Transformers overcome both long-distance dependency and lack of parallelizability.

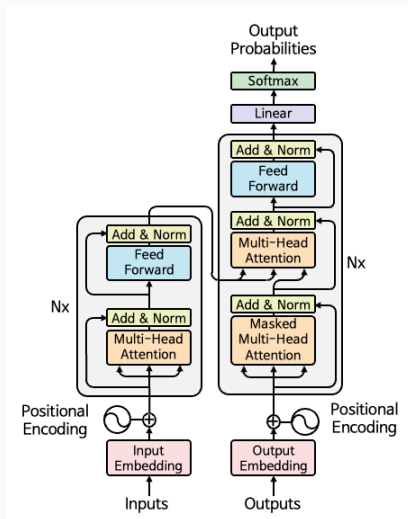
Transformer and self-attention

- In Transformers, attention occurs within a single sentence — all words attend to all words in the previous layer.
(Self-attention + Cross-attention)
- 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.

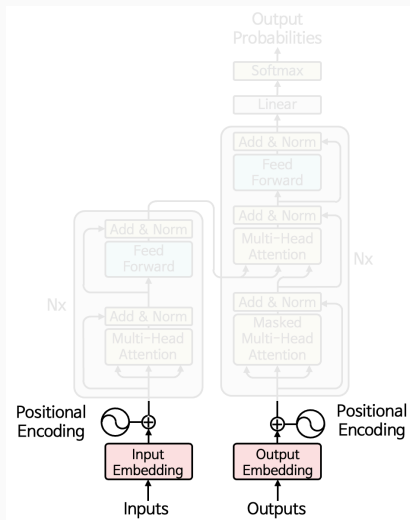
The structure of the Transformer (Vaswani et al., 2017):



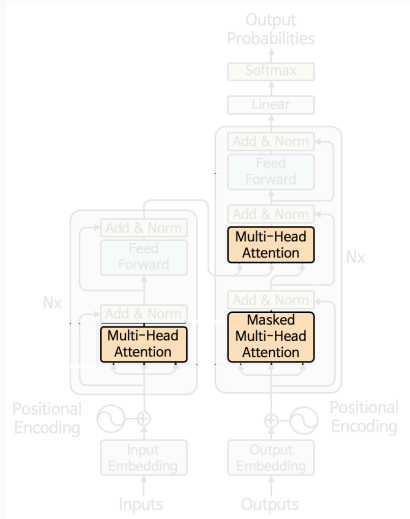
We can see that the same kinds of blocks are **repeatedly stacked**.



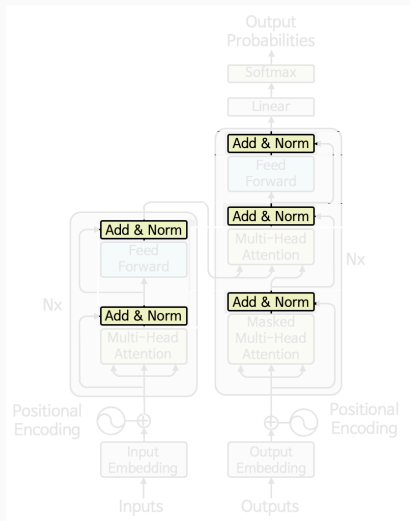
Embedding layer



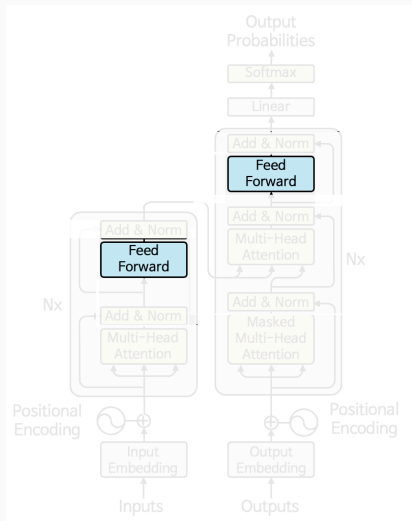
Multi-head attention



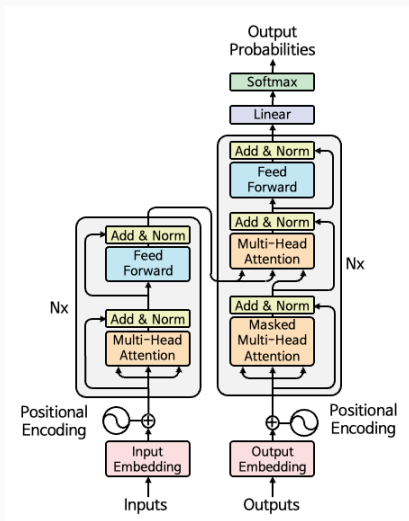
Add & Norm layer



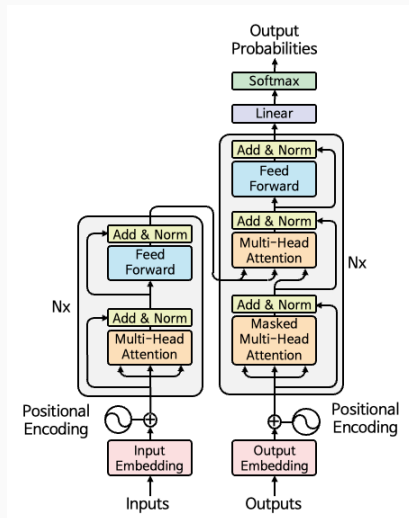
and also the Feed-Forward layer.



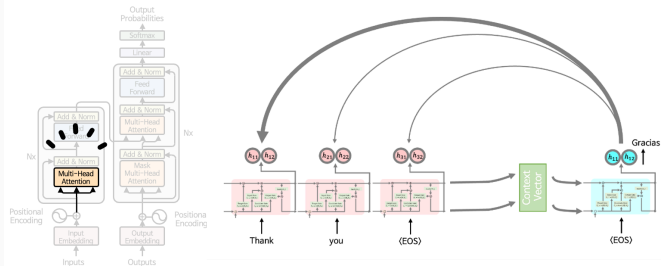
So, it's actually made up of a few components that are repeatedly stacked.



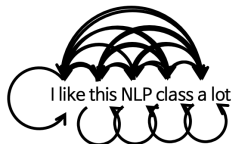
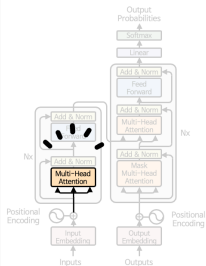
Which of these do you think is the most important feature of the Transformer model?



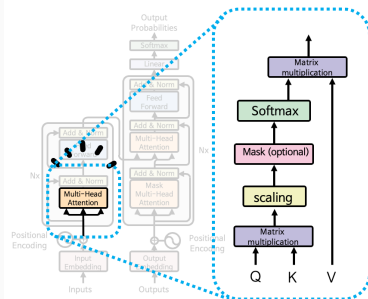
The Transformer's multi-head attention is (1) different from the attention mechanism used in traditional seq2seq models.



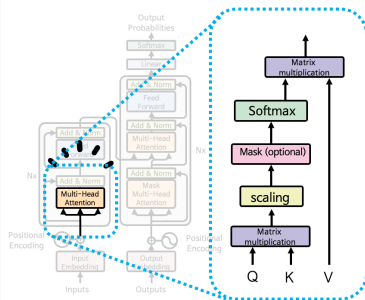
The Transformer's attention (2) captures the relationships between words within the same input sentence.



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.

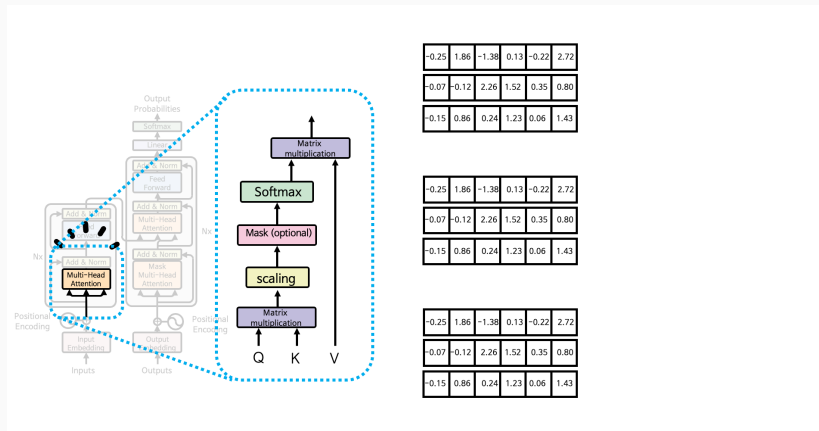


-0.25	1.86	-1.38	0.13	-0.22	2.72
-0.07	-0.12	2.26	1.52	0.35	0.80
-0.15	0.86	0.24	1.23	0.06	1.43

-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

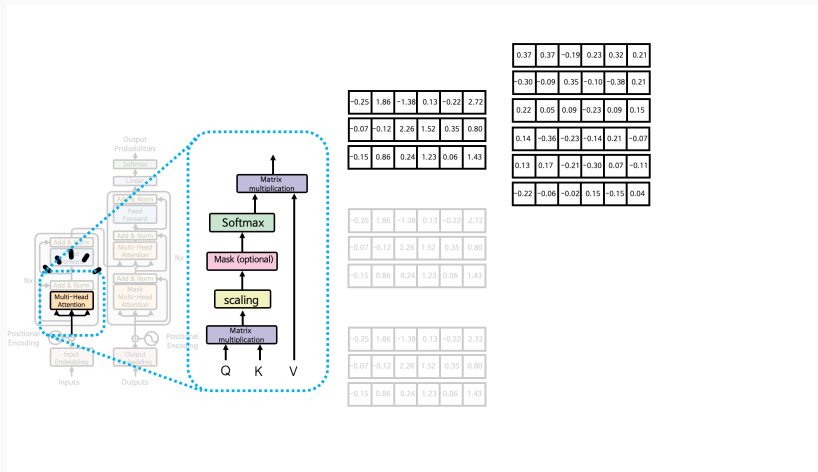
-0.25	1.86	-1.38	0.13	-0.22	2.72
-0.07	-0.12	2.26	1.52	0.35	0.80
-0.15	0.86	0.24	1.23	0.06	1.43

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

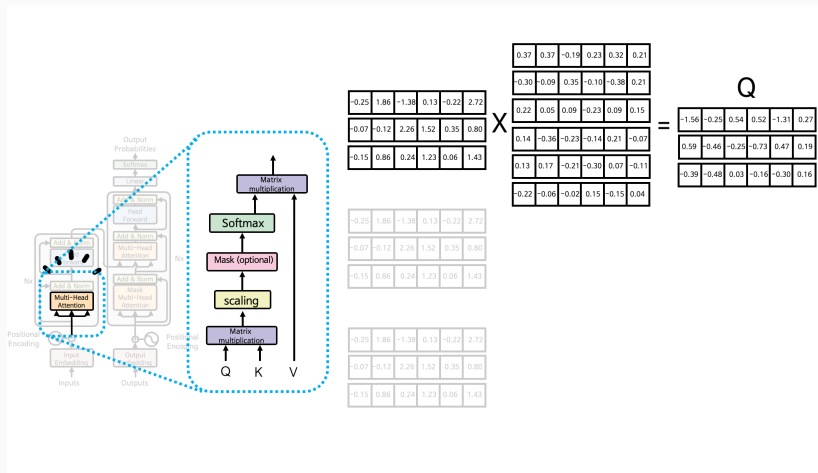


This is done to create **Query (Q)**, **Key (K)**, and **Value (V)** matrices, each representing a different projection of the same input for the attention mechanism.

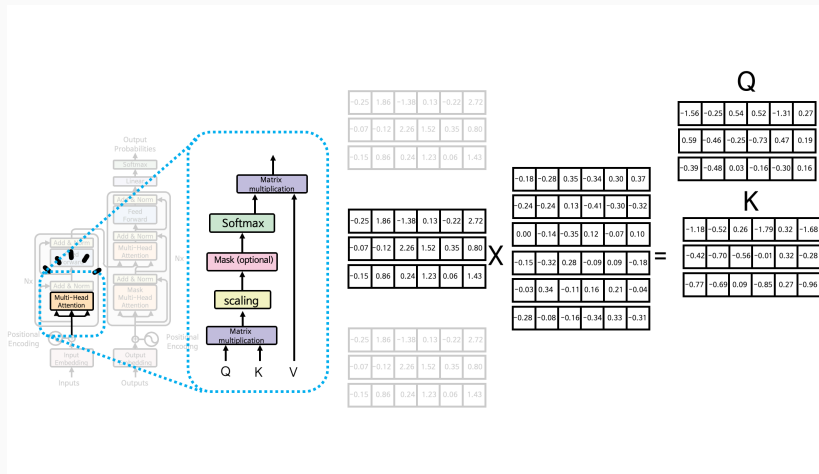
To obtain the Q matrix, we create the following 6×6 weight matrix (randomly initialized).



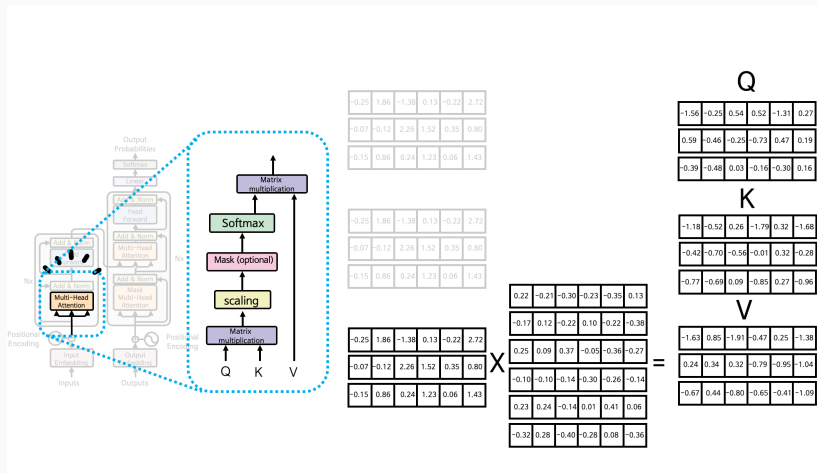
And then we perform matrix multiplication to obtain the Q (Query) matrix.



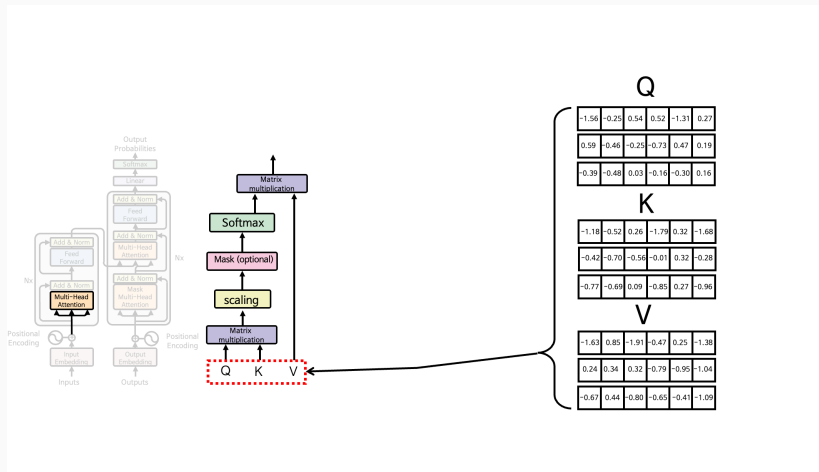
To compute the K (Key) matrix, we create another 6×6 weight matrix (randomly initialized) and multiply it with the input.



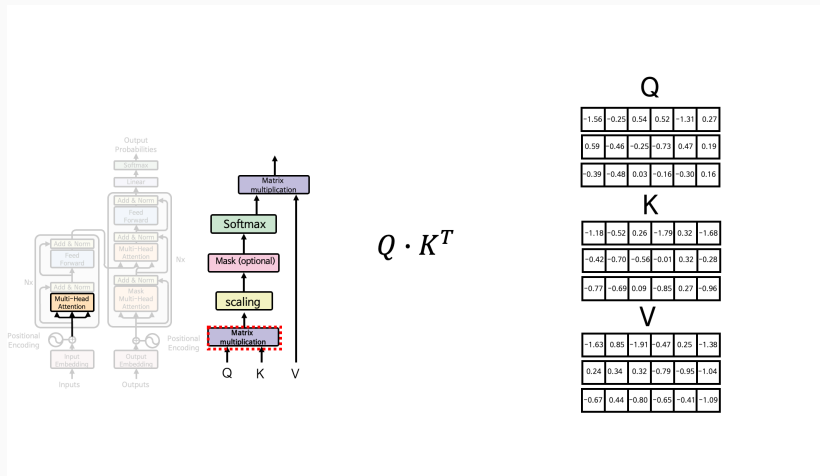
Using the same process, we perform another matrix multiplication to obtain the V (Value) matrix.



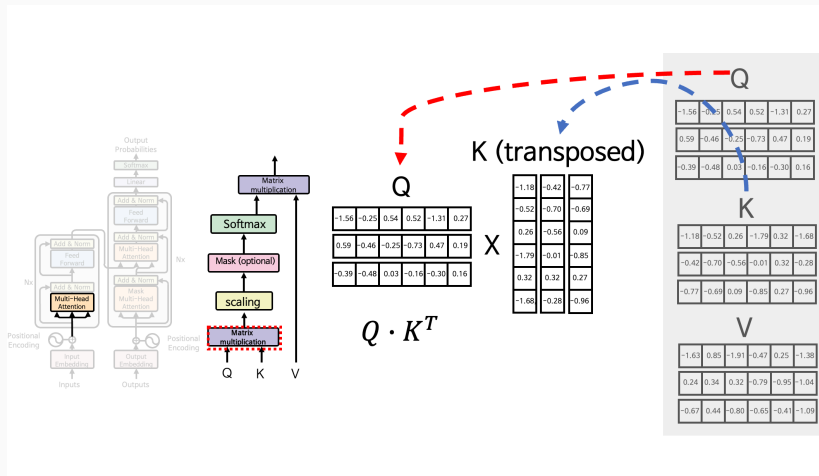
Now, we have the three inputs of the multi-head attention layer: Q (Query), K (Key), and V (Value).



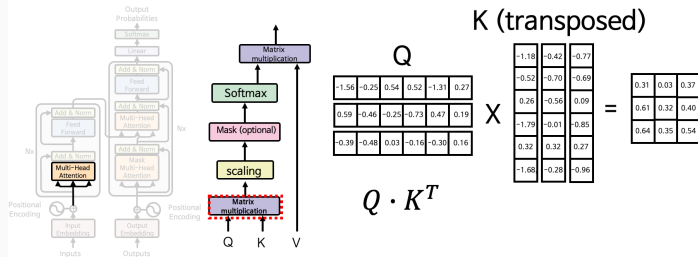
Next, we perform the matrix multiplication between Q and K. The formula for this operation is as follows:



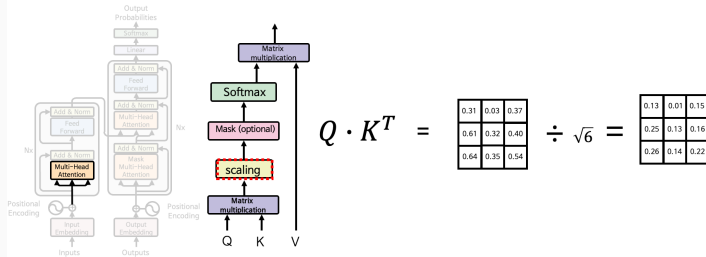
When we plug in the matrix values and calculate,



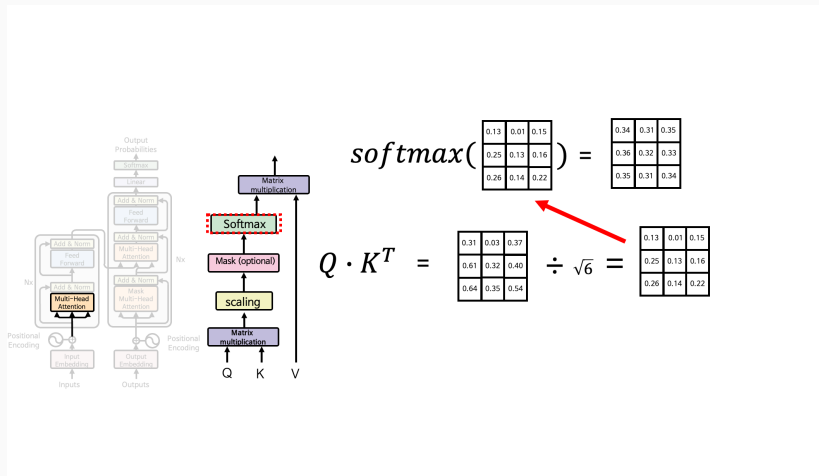
we can obtain the result as follows:



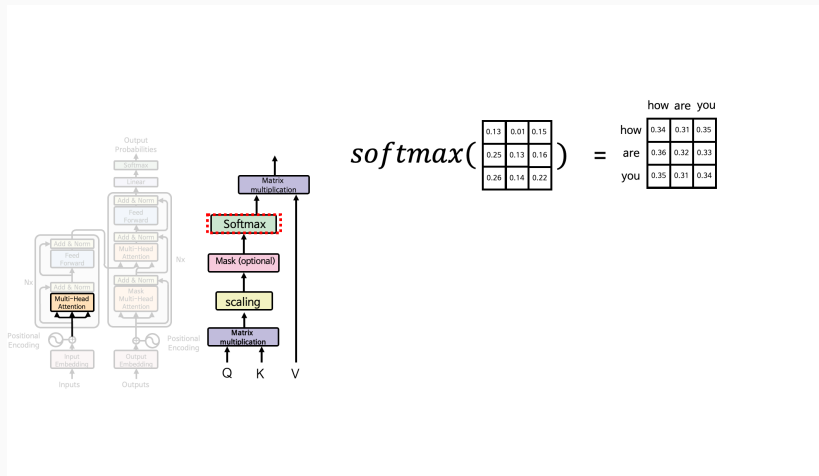
Next, we apply *scaling*, which divides the matrix by $\sqrt{6}$, since in this example the value of d_{model} is 6.



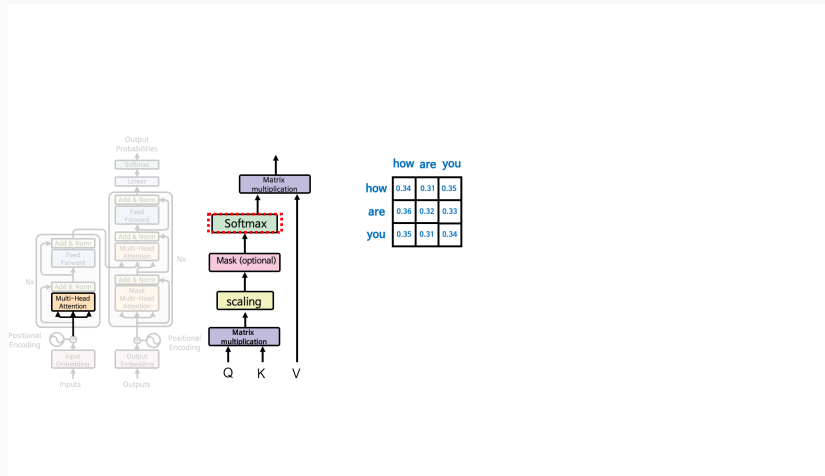
The next softmax layer converts the matrix values into probabilities.

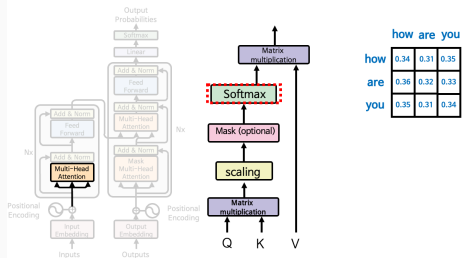


This 3x3 matrix represents the self-attention weights, which shows how each word in the input is related to every other.

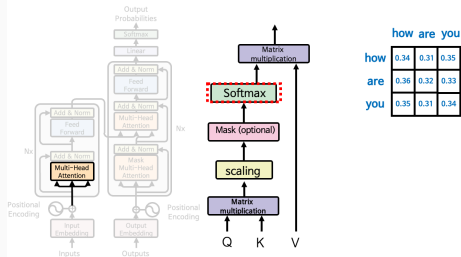


Word pairs with higher relevance receive higher attention values, while those with lower relevance receive smaller values. To be specific...



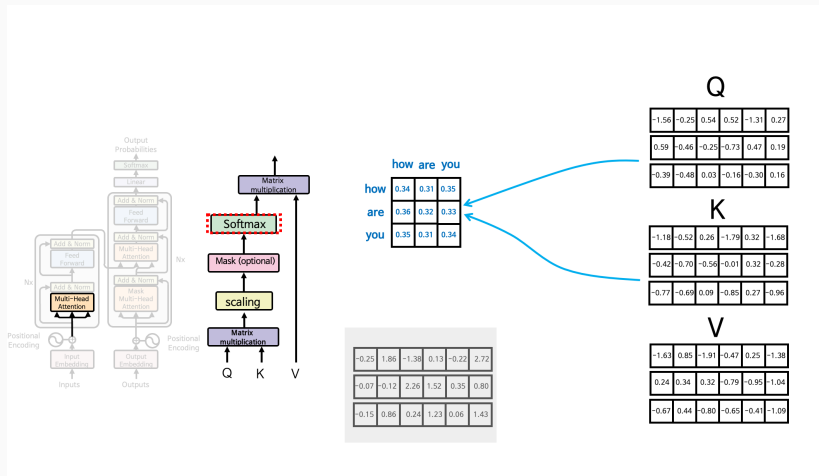


- Each *row*: how much a query word attends to other words (where it sends attention)

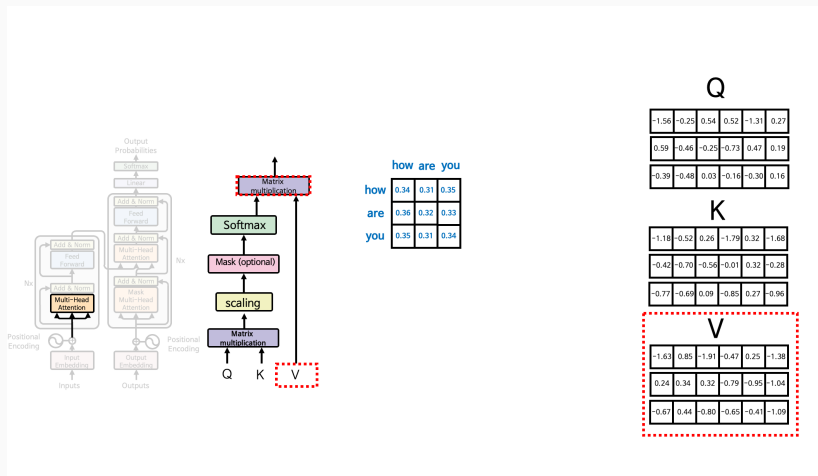


- Each *row*: how much a query word attends to other words (where it sends attention)
- Each *column*: how much a key words is attended to by other words (where it receives attention)

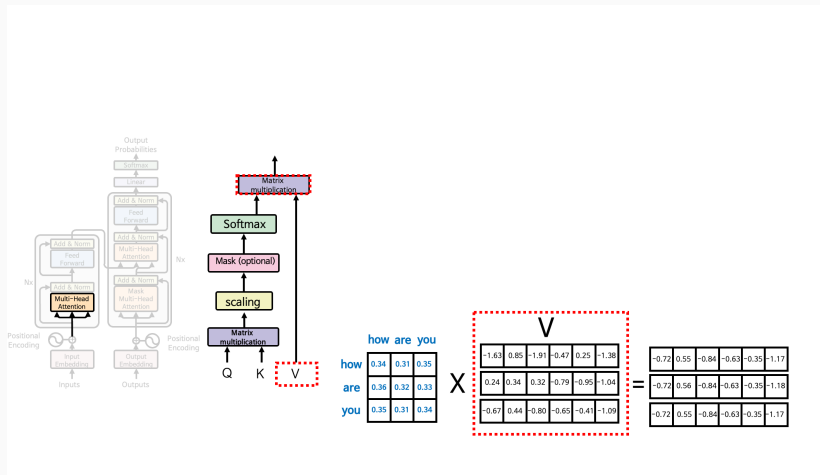
This **self-attention mechanism** is the core structure.



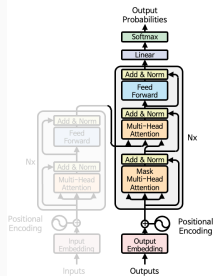
So, what about the final matrix multiplication?



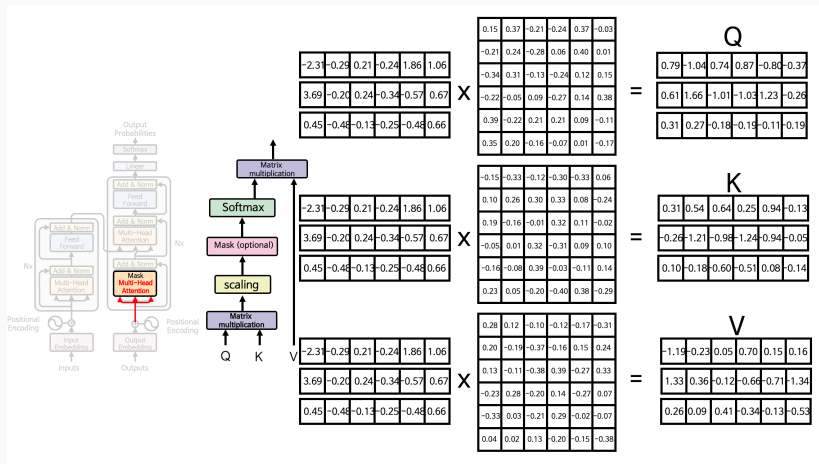
We multiply this by the V matrix to obtain the final output — a self-attended embedding that combines (1) input, (2) positional, and (3) attention information.



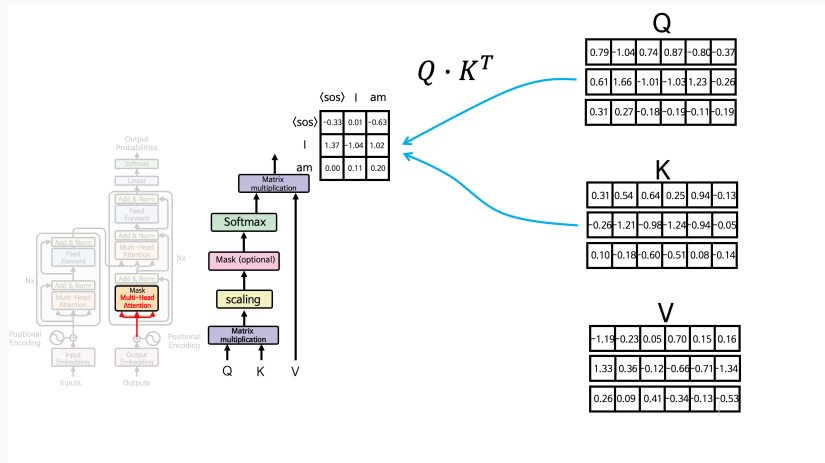
Focus on to the **Decoder**.



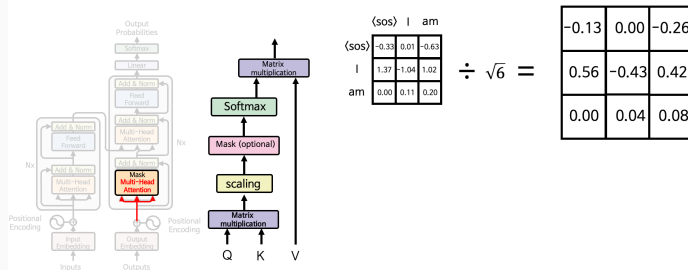
The **Masked Multi-Head Attention** operation in the decoder is almost the same as in the encoder.



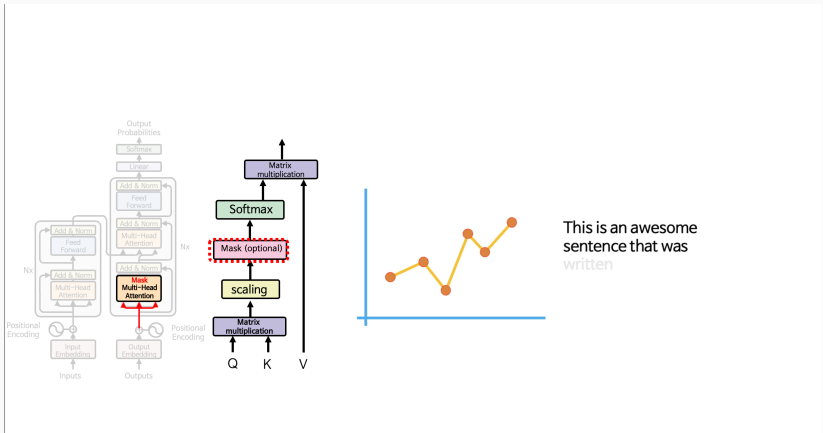
We multiply the Q and K matrices to create the attention matrix.



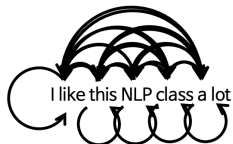
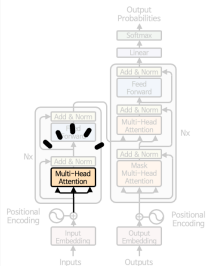
We then scale the matrix by dividing it by $\sqrt{6}$, just as before, so that the range of values changes accordingly.



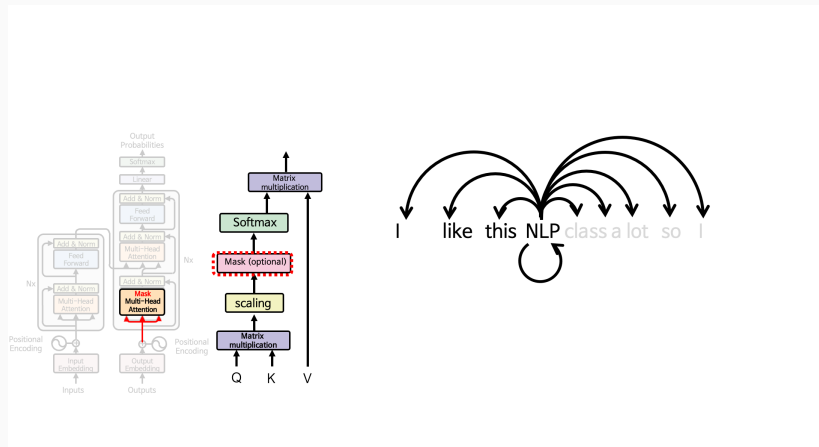
The goal of the transformer *decoder* is to generate the output word sequence, one token at a time (recall: language modeling).



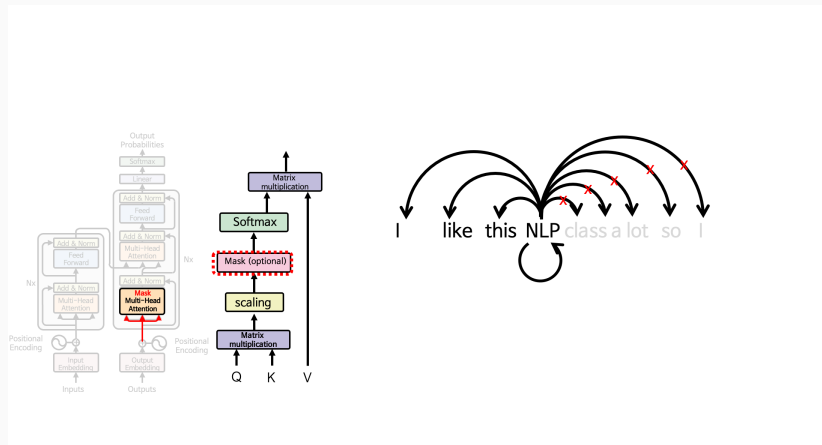
While the encoder needs to consider all tokens in the input sentence to understand the full meaning,



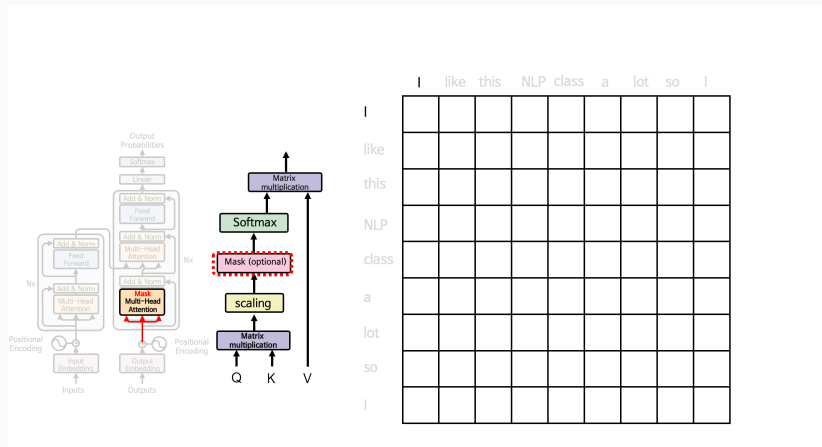
the decoder generates output **one word at a time**.



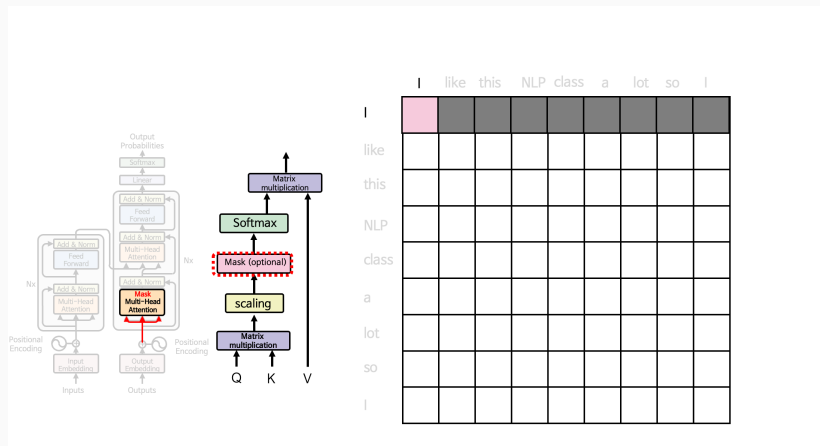
Therefore, it's natural that the decoder should NOT attend to words that haven't been generated yet.

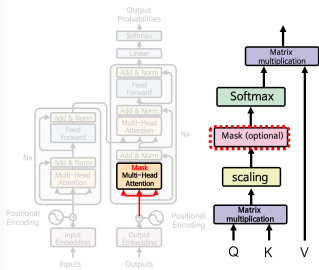


To reflect this characteristic, the decoder applies a masking mechanism during training.

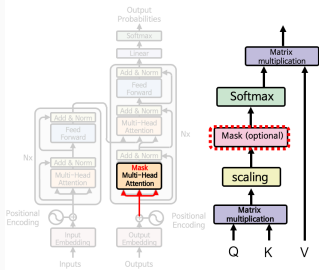


The key idea is to hide future tokens so they do not affect the current prediction.

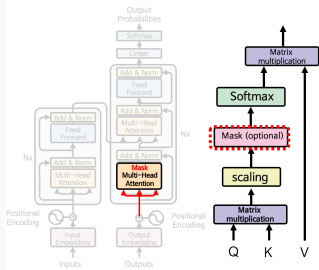




	I	like	this	NLP	class	a	lot	so	I
I									
like									
this									
NLP									
class									
a									
lot									
so									
I									

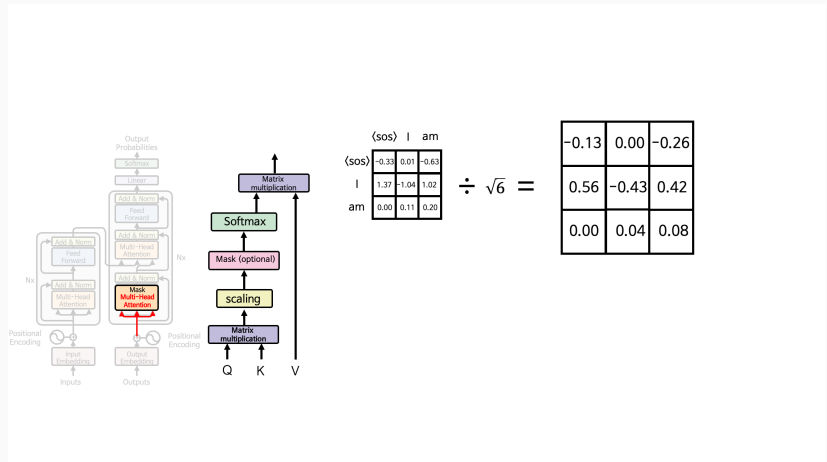


	I	like	this	NLP	class	a	lot	so	I
I									
like									
this									
NLP									
class									
a									
lot									
so									
I									

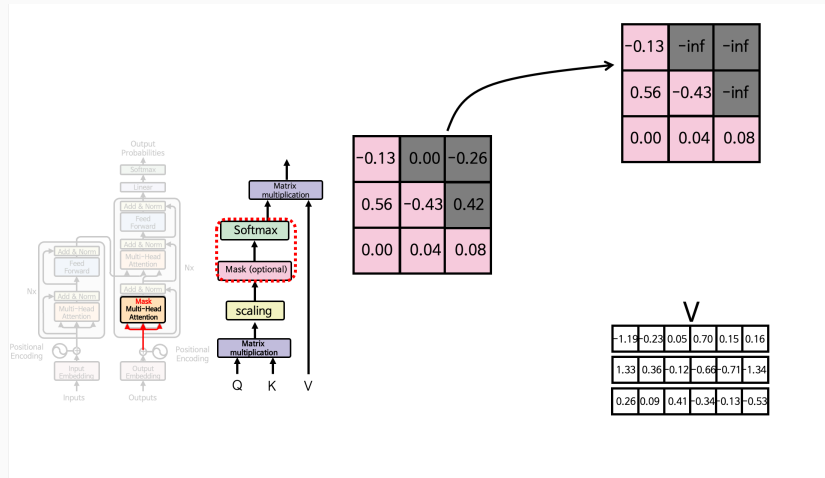


	I	like	this	NLP	class	a	lot	so	I
I									
like									
this									
NLP									
class									
a									
lot									
so									
I									

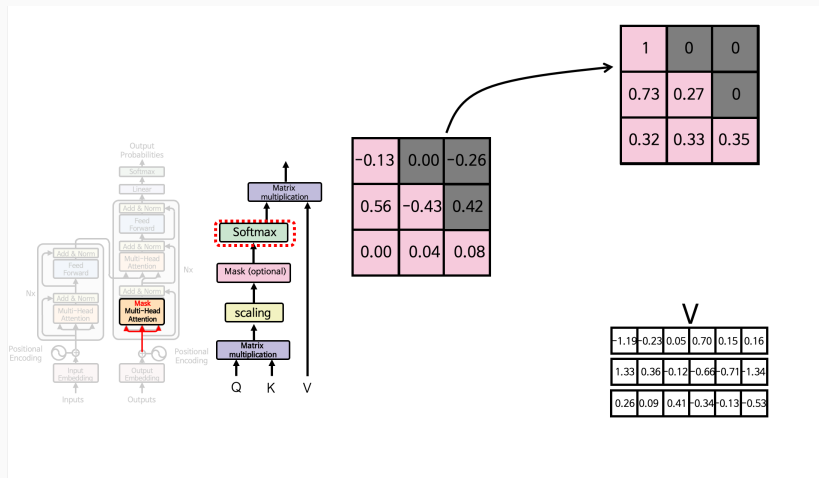
When this masking algorithm is applied to the attention matrix, we get:



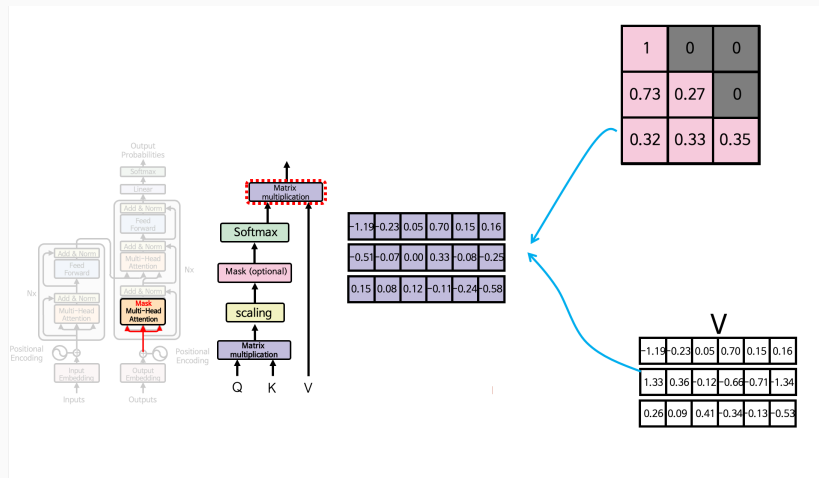
We add $-\text{inf}$ to the masked positions because, after passing through the softmax layer, $-\text{inf}$ becomes 0, effectively eliminating attention to those positions.



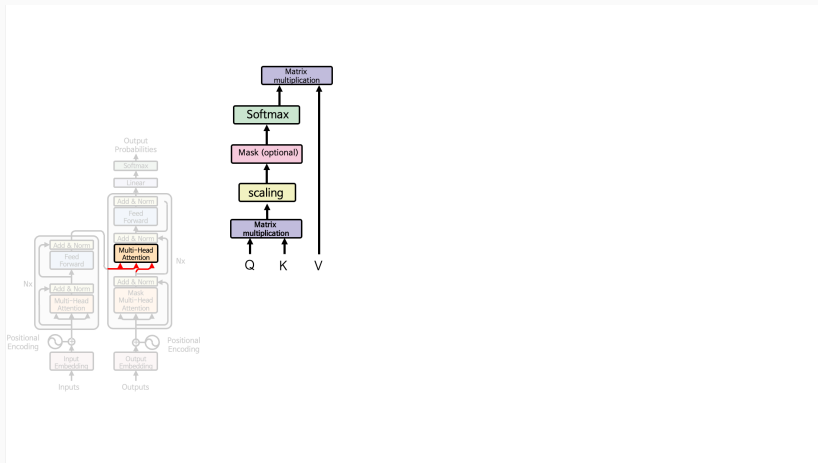
Feeding this matrix into the softmax layer gives us:



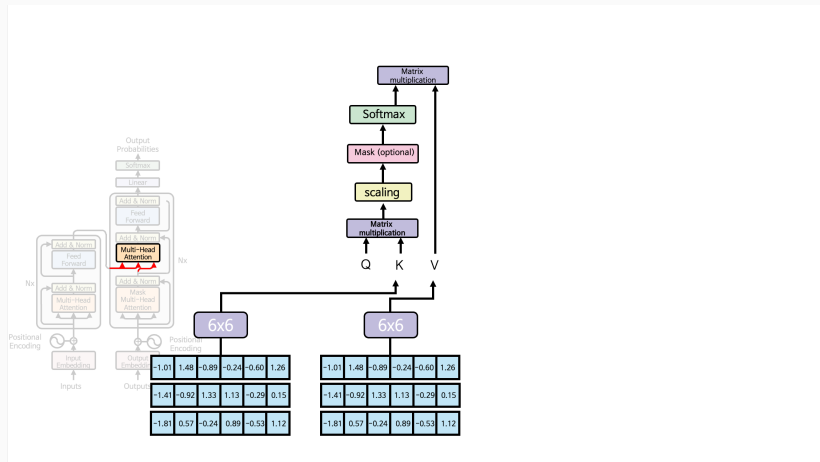
Then we multiply the two matrices as follows:



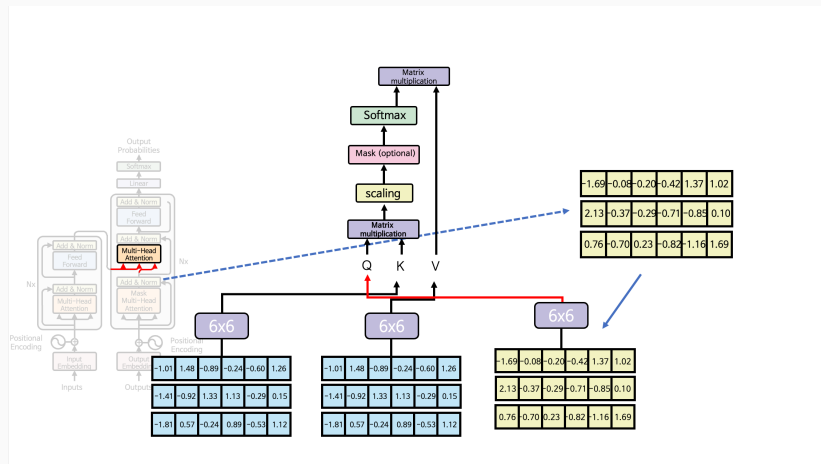
The decoder's **second** multi-head attention operates the same way as the encoder's, except for the inputs.



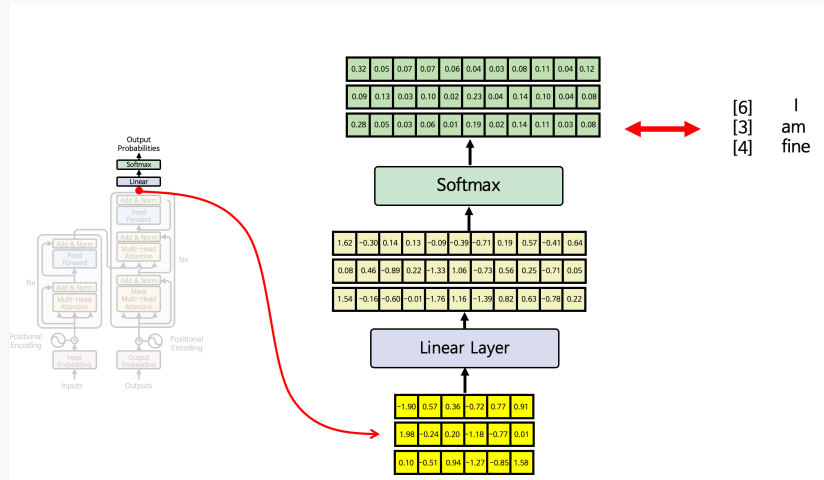
Here, the values of K and V are derived from the **encoder's final output**, multiplied by a 6×6 matrix.

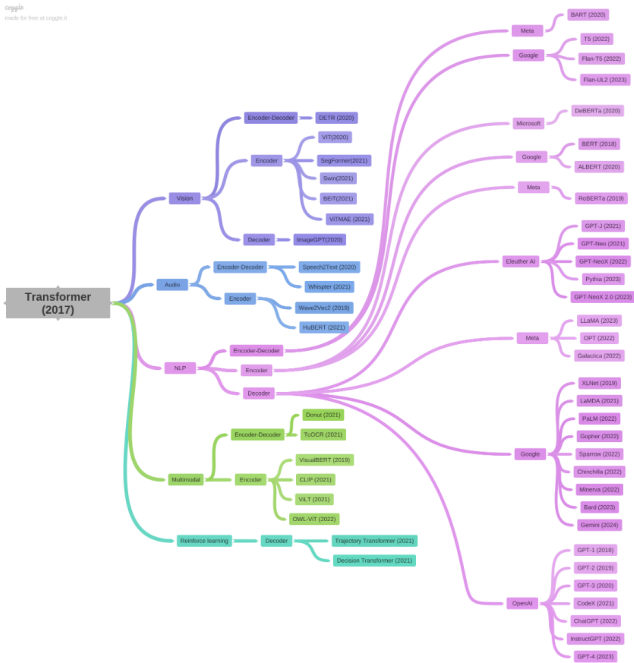


The value of Q comes from the output of the decoder's previous Add & Norm layer. In other words, the decoder determines which parts of the encoder's output (K, V) to attend to, based on the context it has generated so far (Q)-similar to attention in RNNs.



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.





Different transformers

- **Encoder-only models** use only the encoder stack to **understand text**.
 - Bidirectional attention (see both left and right context)
 - Tasks: text classification, NER, POS tagging
 - Examples: **BERT**, RoBERTa, ALBERT

- **Decoder-only models** use only the decoder stack to **generate text**.
 - Left-to-right (causal) attention — cannot see future tokens
 - Tasks: text generation, dialogue, code completion
 - Examples: **GPT**, LLaMA, Gemini, Claude

- **Encoder-Decoder models** combine both for **sequence-to-sequence tasks**.
 - Encoder encodes the input; decoder generates the output
 - Tasks: translation, summarization
 - Examples: **T5**, BART, mBART

LLMs in 2025

Sharing Thoughts

- Take a few minutes to share your thoughts or reflections on today's session.

Sharing Thoughts

- Take a few minutes to share your thoughts or reflections on today's session.
- Contribute one or two points to the **shared slide deck** (click [here](#), or a link in the course website) - this could be:

Sharing Thoughts

- Take a few minutes to share your thoughts or reflections on today's session.
- Contribute one or two points to the **shared slide deck** (click [here](#), or a link in the course website) - this could be:
 - Something new you learned about LLMs or Transformers

Sharing Thoughts

- Take a few minutes to share your thoughts or reflections on today's session.
- Contribute one or two points to the **shared slide deck** (click [here](#), or a link in the course website) - this could be:
 - Something new you learned about LLMs or Transformers
 - A question or idea you'd like to explore further

Sharing Thoughts

- Take a few minutes to share your thoughts or reflections on today's session.
- Contribute one or two points to the **shared slide deck** (click [here](#), or a link in the course website) - this could be:
 - Something new you learned about LLMs or Transformers
 - A question or idea you'd like to explore further
 - An observation about how these models relate to real-world NLP tasks

Wrap up

Wrap-up

- We explored recent advances in large language models (LLMs) as of 2025.

Wrap-up

- We explored recent advances in large language models (LLMs) as of 2025.
- In the lab session, we will further experiment with these models using *Ollama*, an open-source platform for running LLMs locally.