

Transformers II

Dr. Parlett-Pelleriti

Outline

- Overview of Transformer Architecture
- Attention
- Attention in Transformers
- Multi-Headed Attention
- Masking and Cross-Attention
- Regularization in Transformers
- Transformers You Might Know

Attention Is All You Need

Ashish Vaswani*

Google Brain

avaswani@google.com

Noam Shazeer*

Google Brain

noam@google.com

Niki Parmar*

Google Research

nikip@google.com

Jakob Uszkoreit*

Google Research

usz@google.com

Llion Jones*

Google Research

llion@google.com

Aidan N. Gomez* †

University of Toronto

aidan@cs.toronto.edu

Łukasz Kaiser*

Google Brain

lukaszkaizer@google.com

Illia Polosukhin* ‡

illia.polosukhin@gmail.com

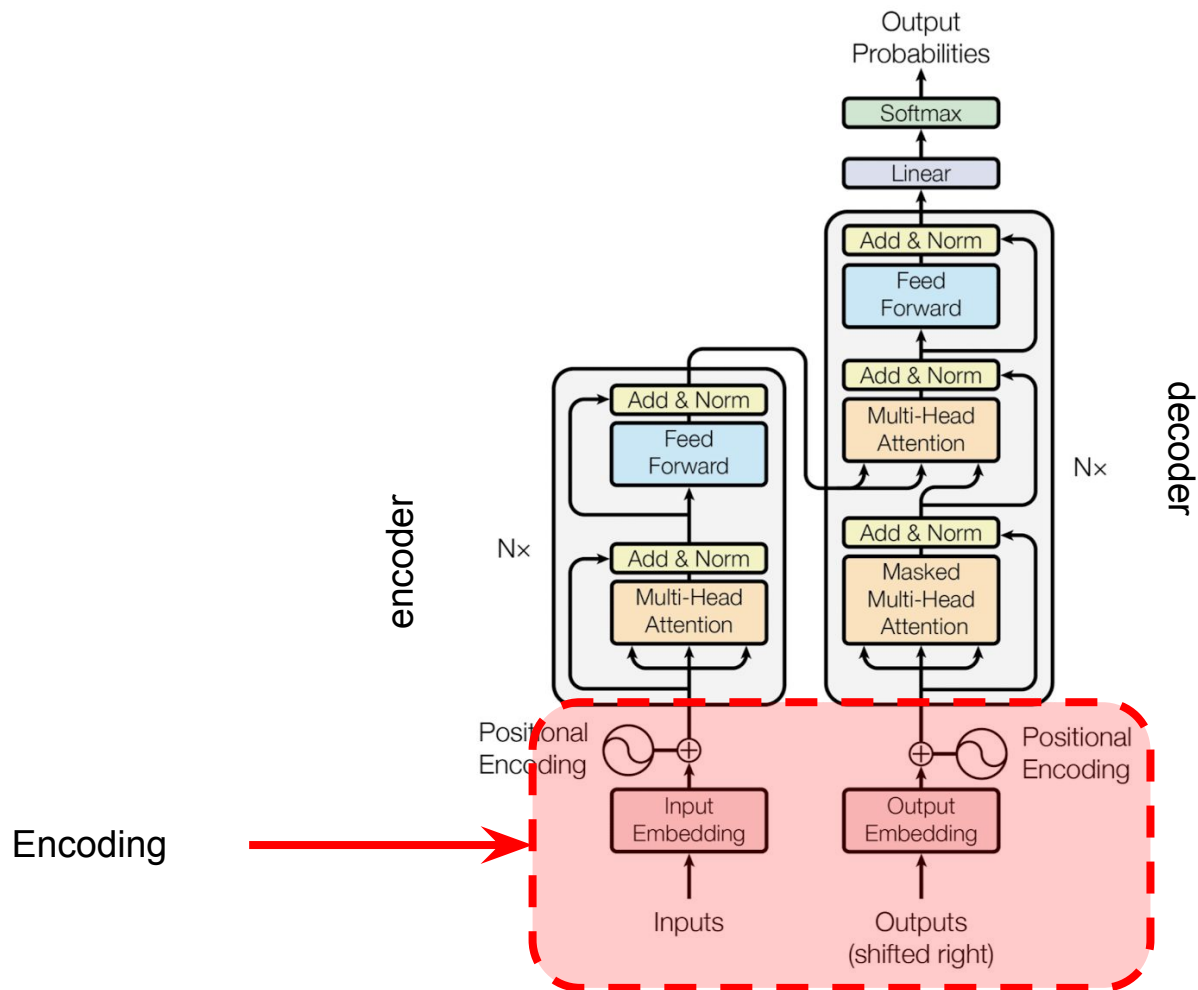


Figure 1: The Transformer - model architecture.

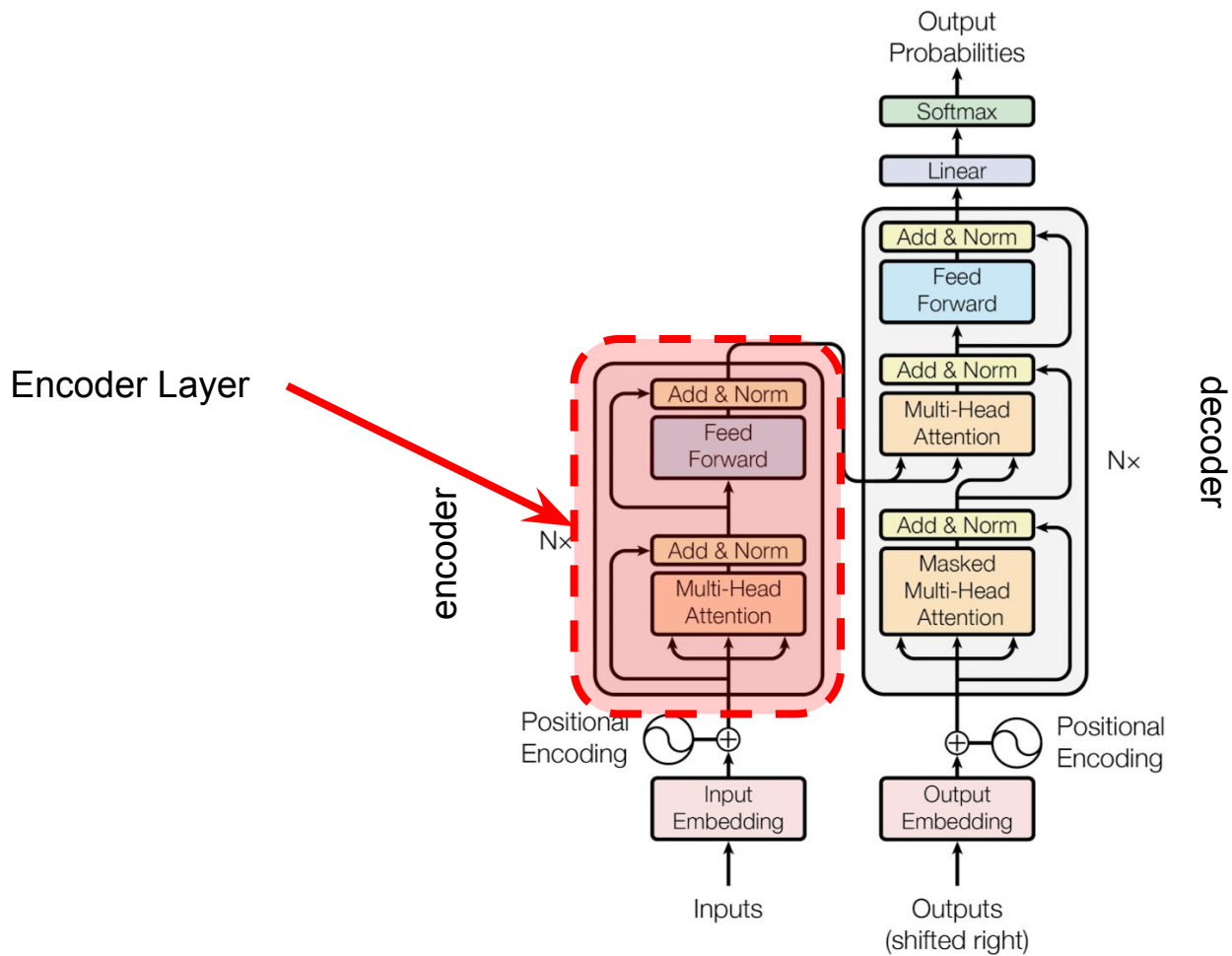


Figure 1: The Transformer - model architecture.

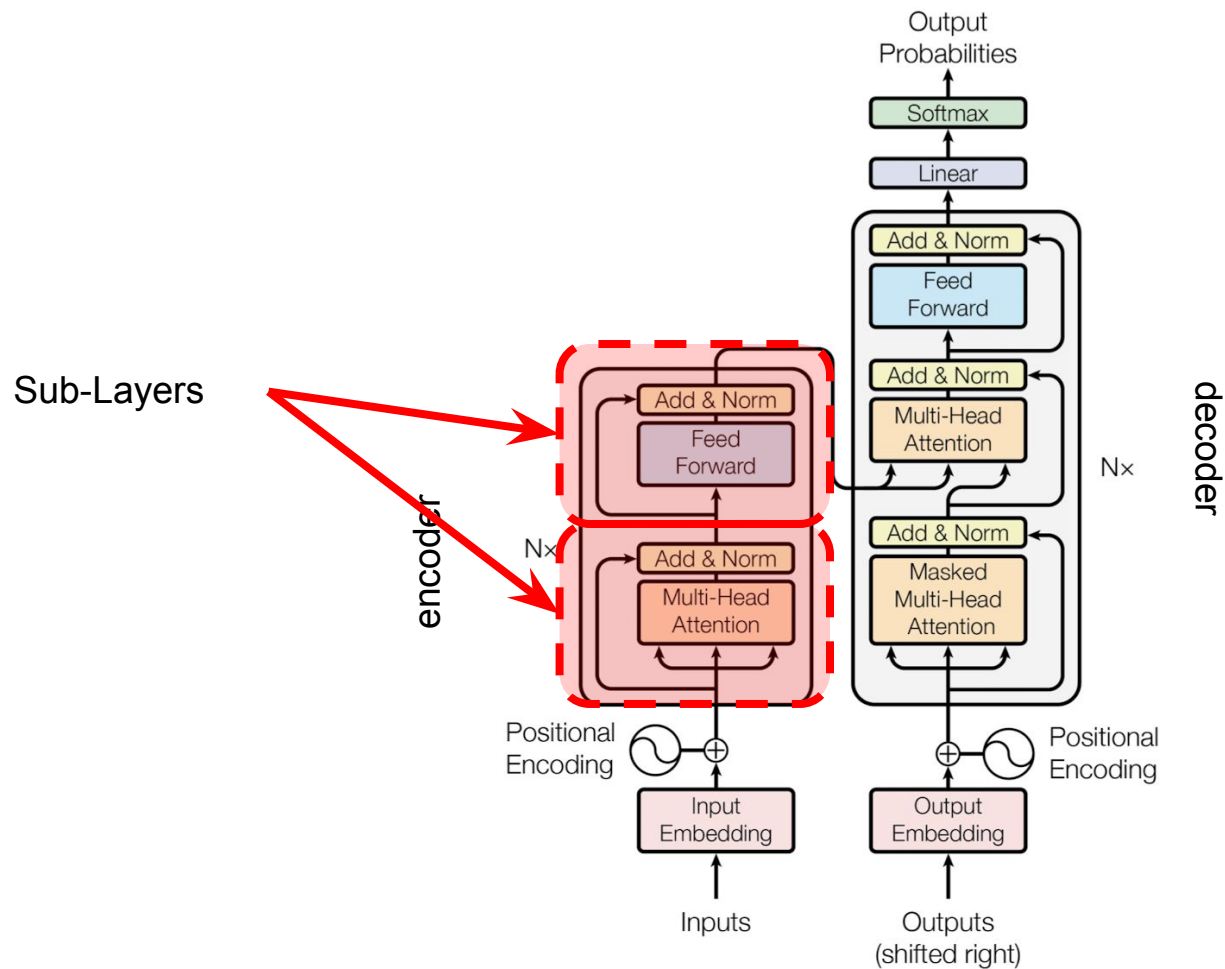


Figure 1: The Transformer - model architecture.

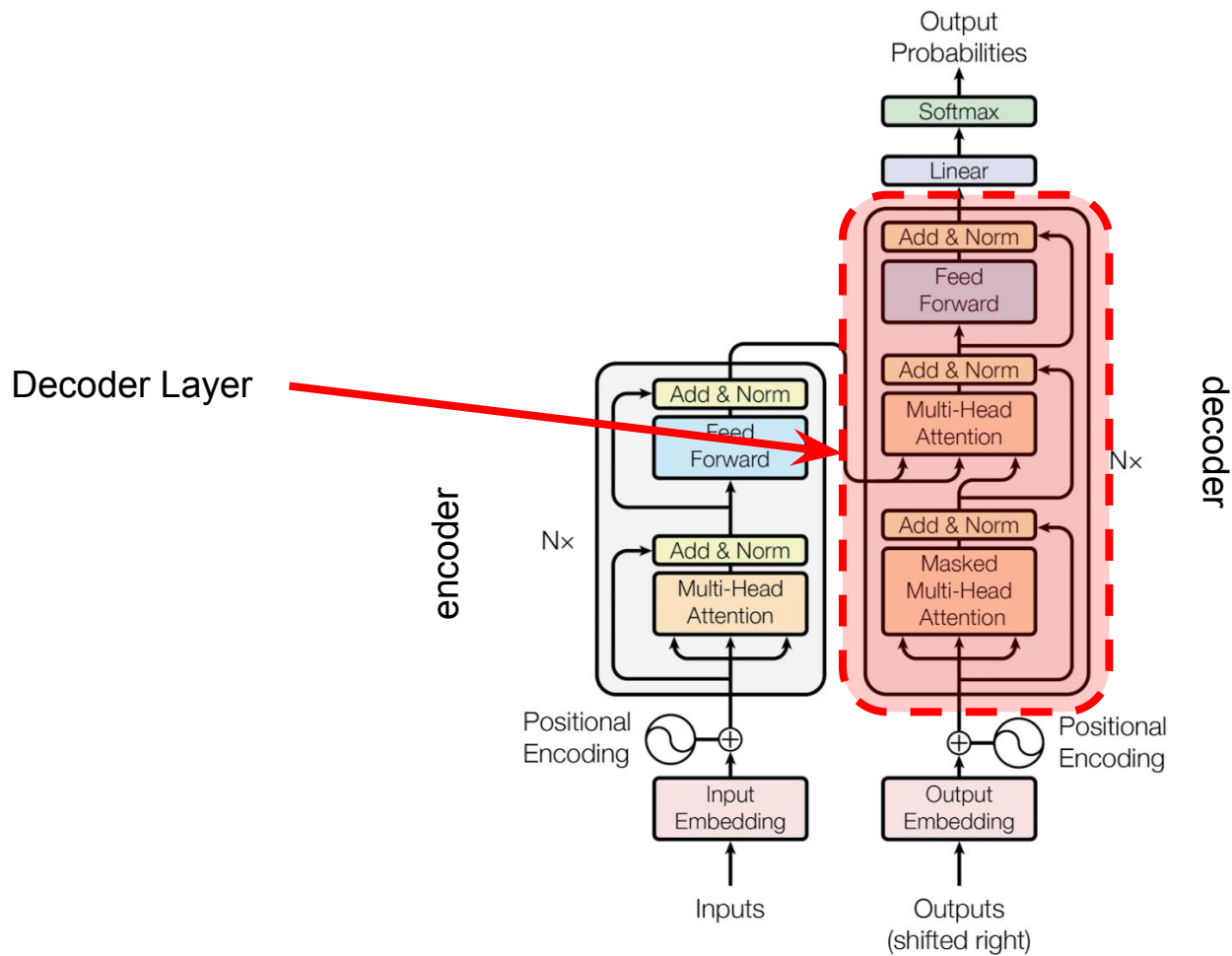


Figure 1: The Transformer - model architecture.

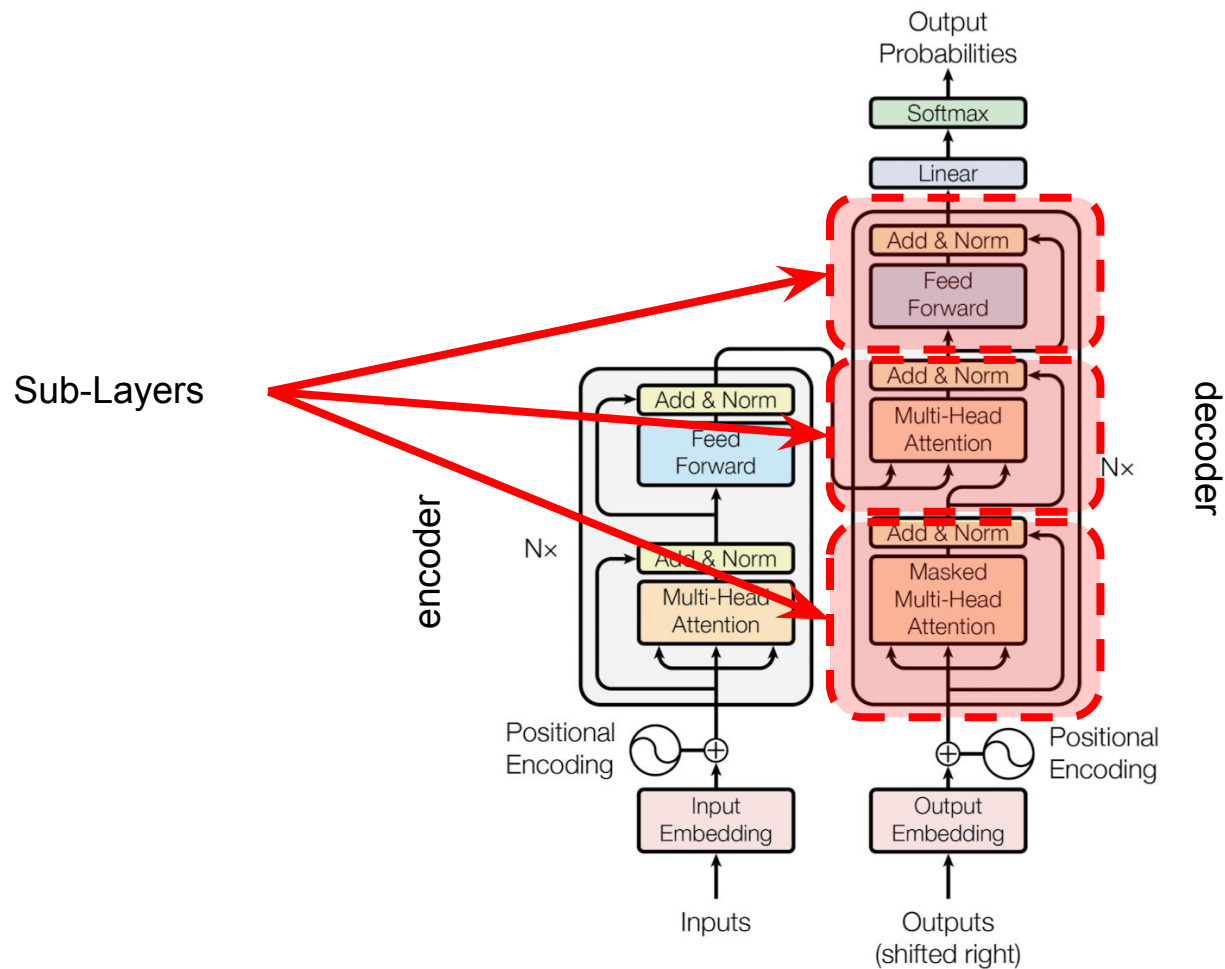


Figure 1: The Transformer - model architecture.

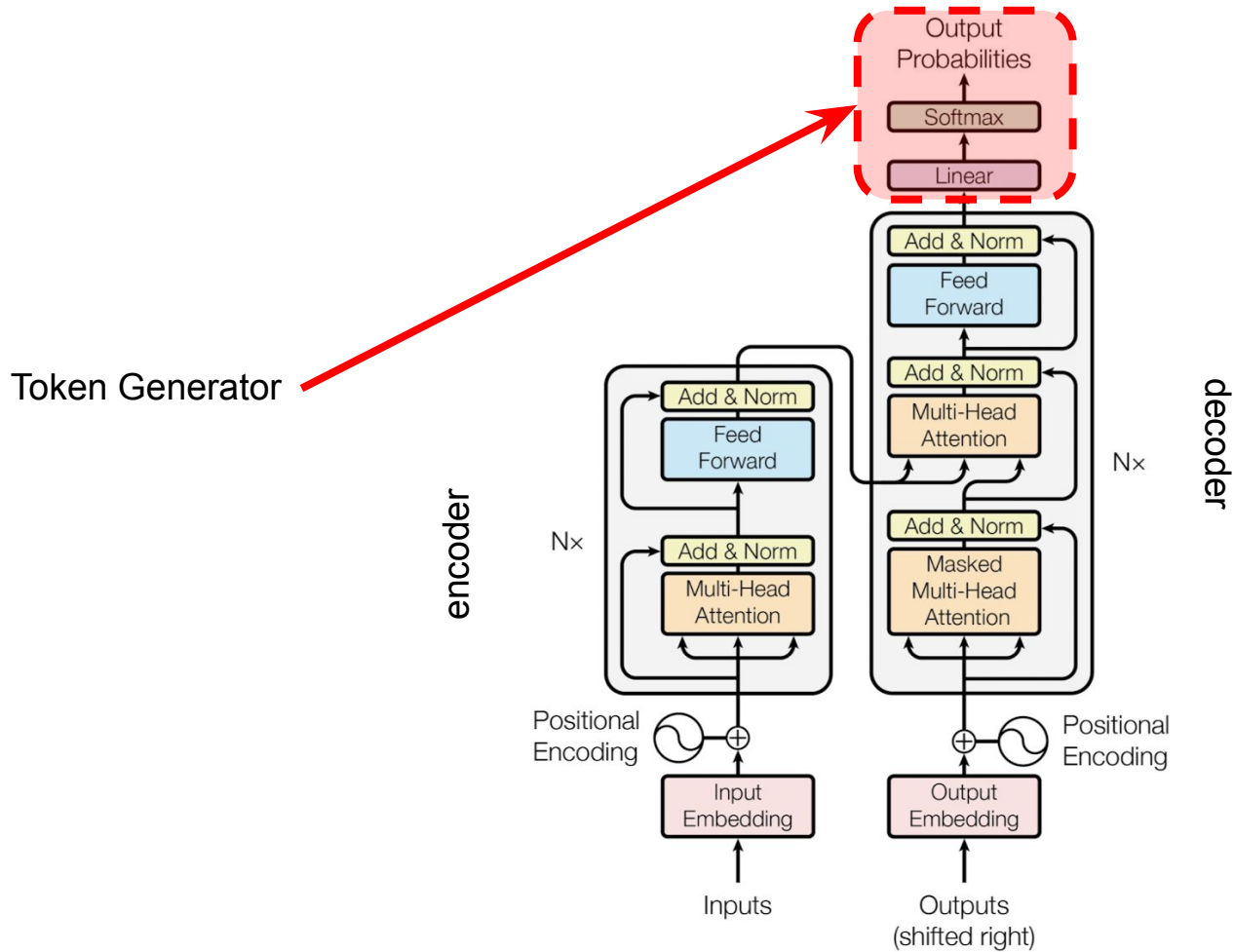


Figure 1: The Transformer - model architecture.

Vicki
@vboykis

...

They don't tell you this in the paper (well they do but you have to read it like 15 times)

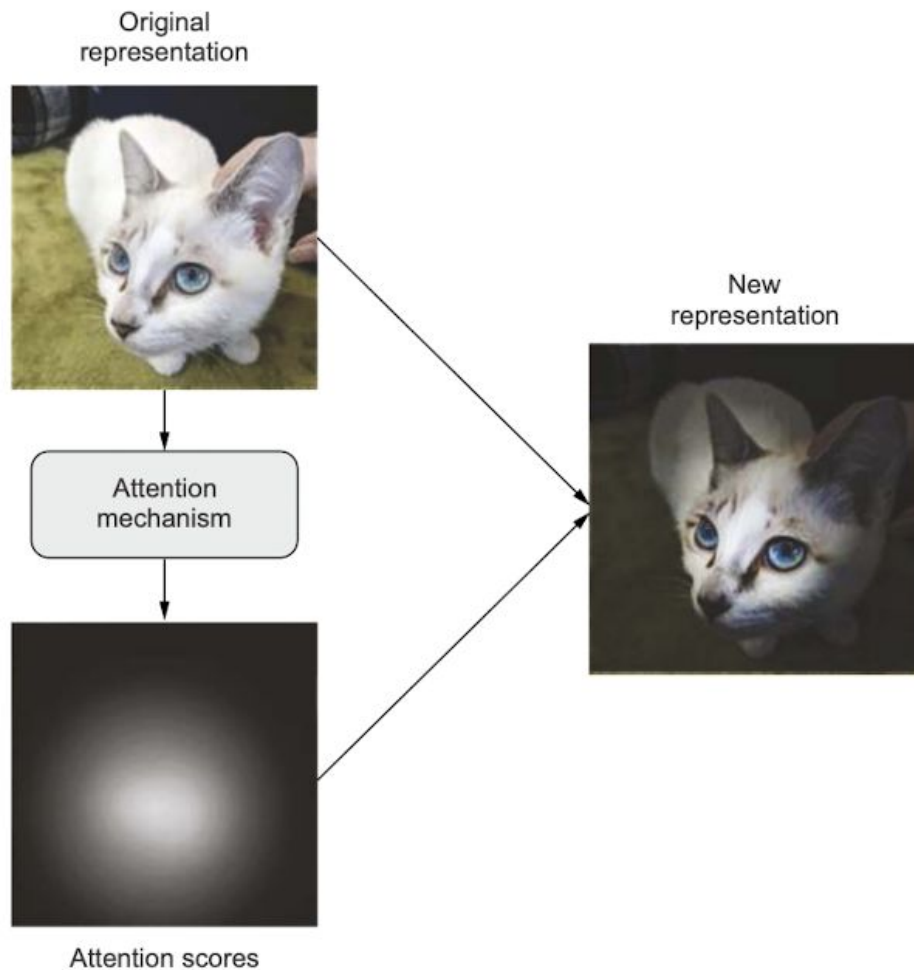


Multiplying
a lot of vectors
a lot of times
with scaled softmax



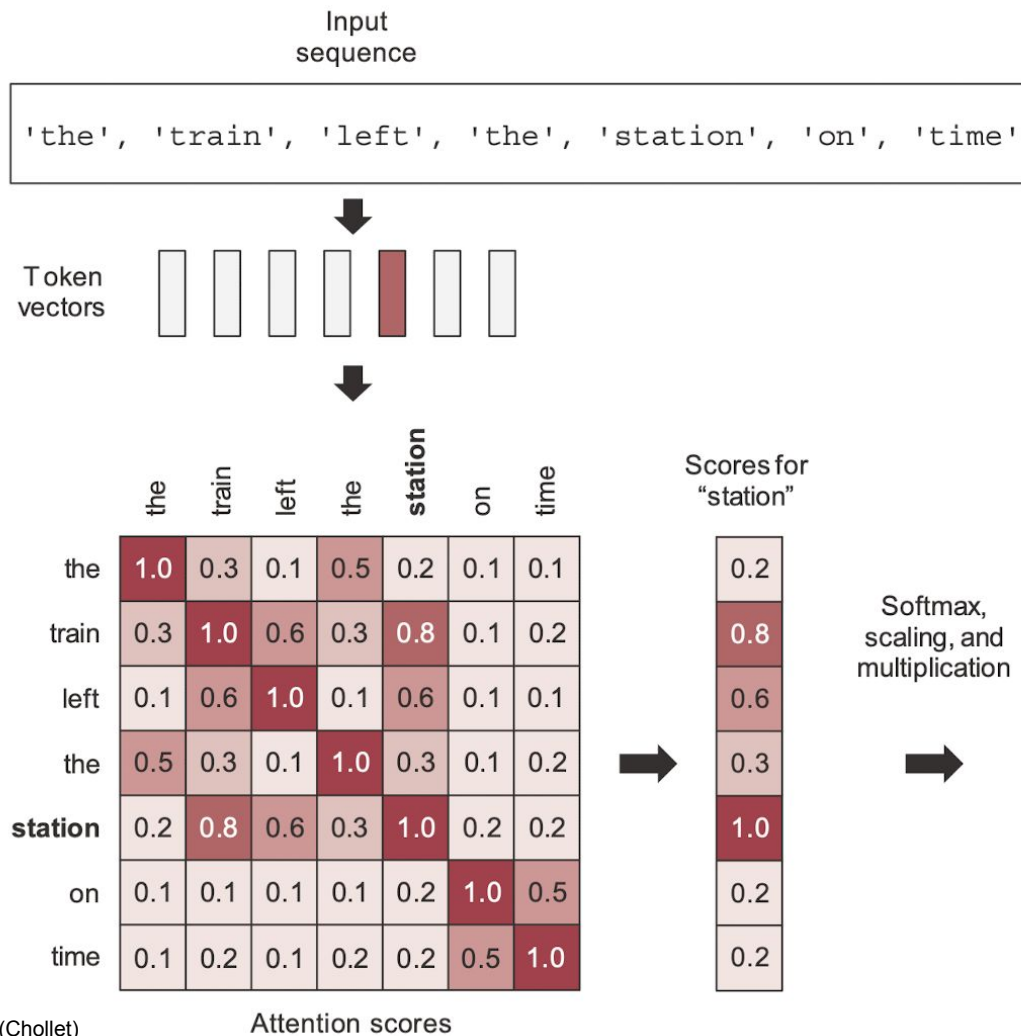
Attention

Attention



Attention

Attention

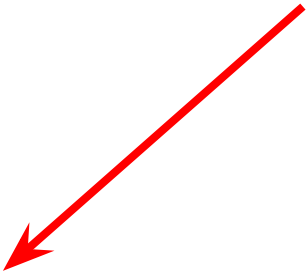


Attention

$$x'_i = \sum_{j=1}^n w_{ji} \cdot x_j$$

Attention

Attention weights

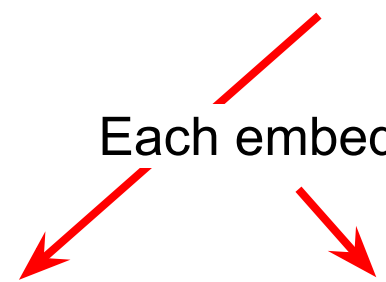
$$x'_i = \sum_{j=1}^n w_{ji} \cdot x_j$$


Attention

$$x'_i = \sum_{j=1}^n w_{ji} \cdot x_j$$

Attention weights

Each embedding



Attention

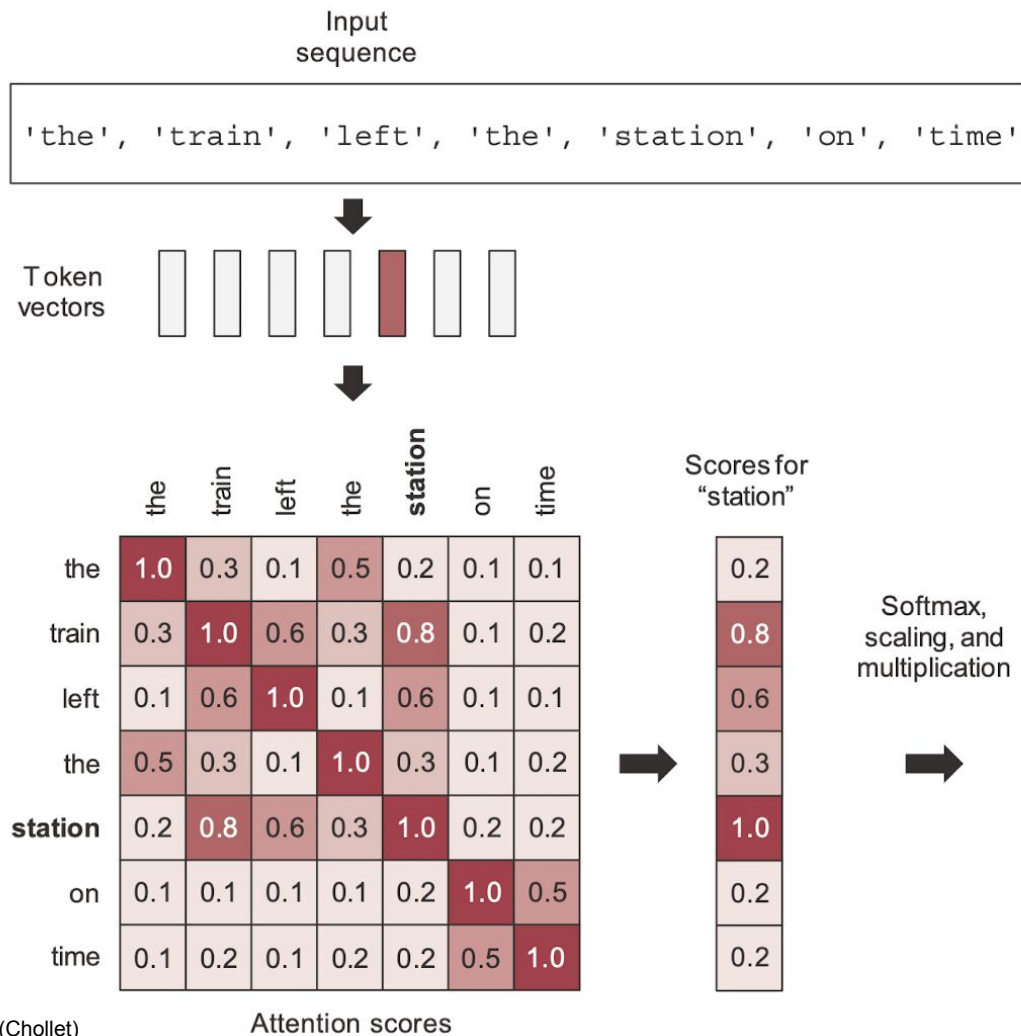
Attention weights

Each embedding

$$x'_i = \sum_{j=1}^n w_{ji} \cdot x_j$$

Context aware
vector

Attention



Attention in Transformers

Queries, Keys, and Values

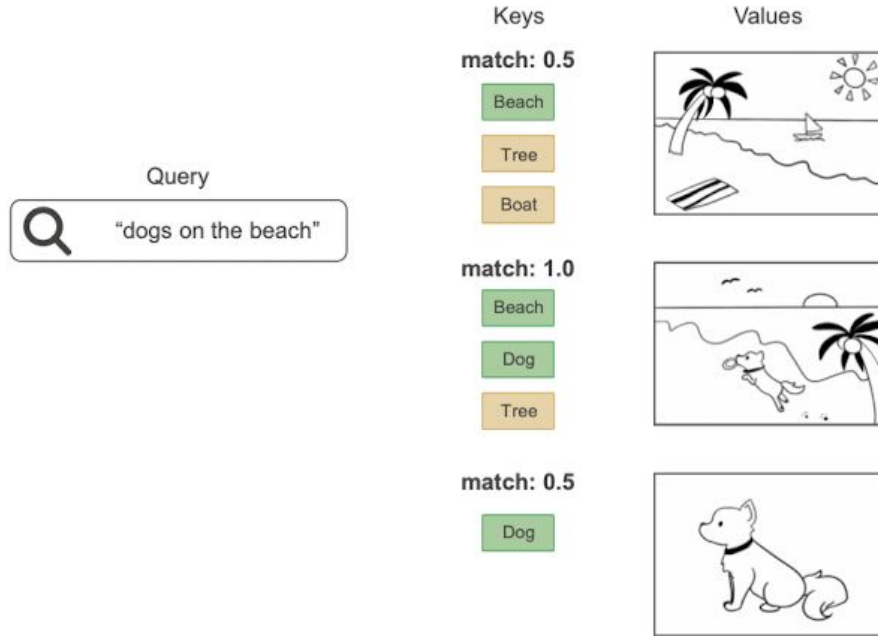


Figure 11.7 Retrieving images from a database: the “query” is compared to a set of “keys,” and the match scores are used to rank “values” (images).

Queries, Keys, and Values

$$Q = E * W_q$$

$$K = E * W_k$$

$$V = E * W_v$$

Learnable
Parameters

Queries, Keys, and Values

$$Q = E * W_q$$

$$K = E * W_k$$

$$V = E * W_v$$

Queries, Keys, and Values

Learnable
Parameters

$$Q = E * W_q$$

$$K = E * W_k$$

$$V = E * W_v$$

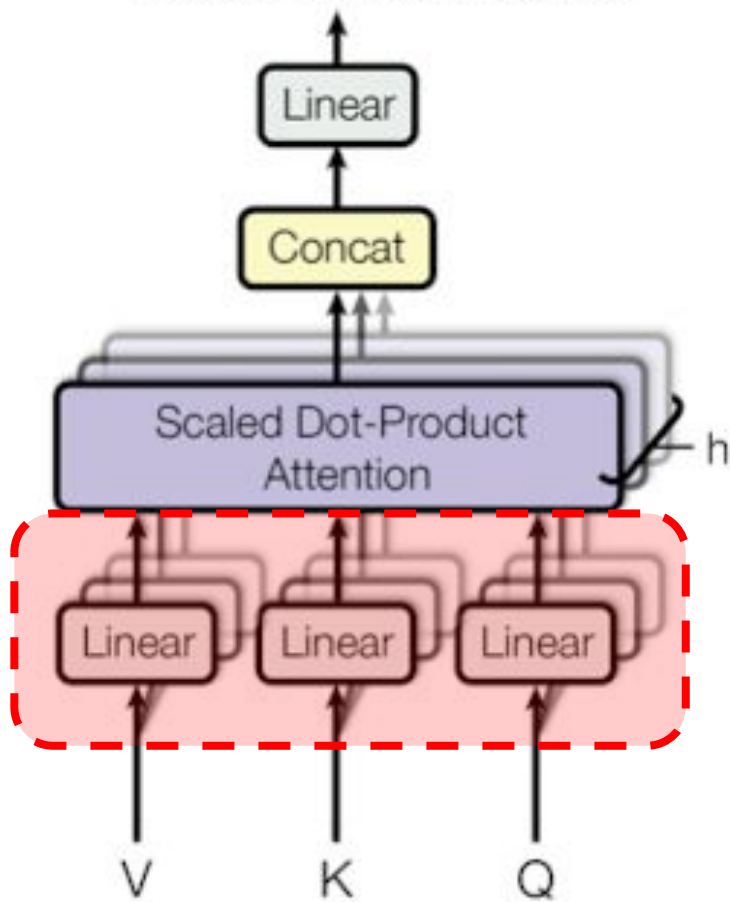
W_q

W_k

W_v

Must be the same
size

Multi-Head Attention



Queries, Keys, and Values

$$\textit{Attention}(Q, K, V) = \textit{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Vicki
@vboykis

...

They don't tell you this in the paper (well they do but you have to read it like 15 times)



Multiplying
a lot of vectors
a lot of times
with scaled softmax



Attention

Queries, Keys, and Values

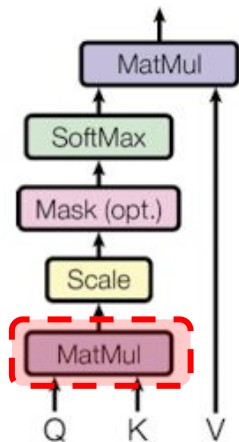
Attention Weights

$$\textit{Attention}(Q, K, V) = \textit{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Queries, Keys, and Values

$$\textit{Attention}(Q, K, V) = \textit{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

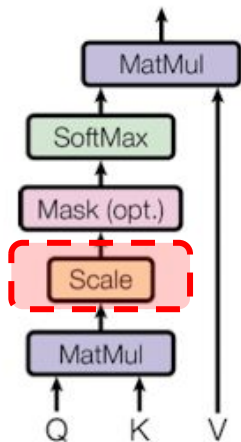
Scaled Dot-Product Attention



Queries, Keys, and Values

$$\textit{Attention}(Q, K, V) = \textit{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

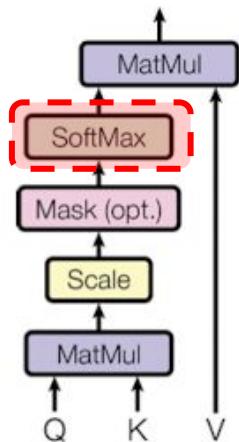
Scaled Dot-Product Attention



Queries, Keys, and Values

$$\textit{Attention}(Q, K, V) = \textit{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

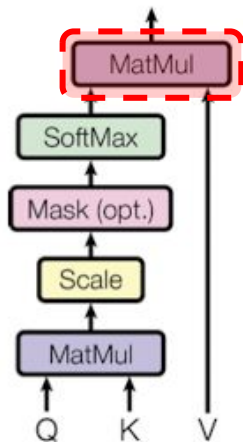
Scaled Dot-Product Attention



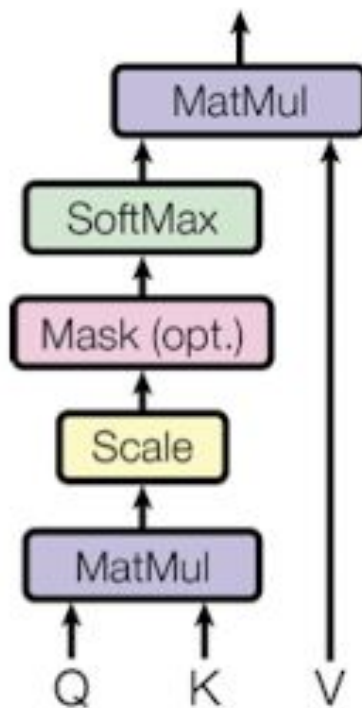
Queries, Keys, and Values

$$\textit{Attention}(Q, K, V) = \textit{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention

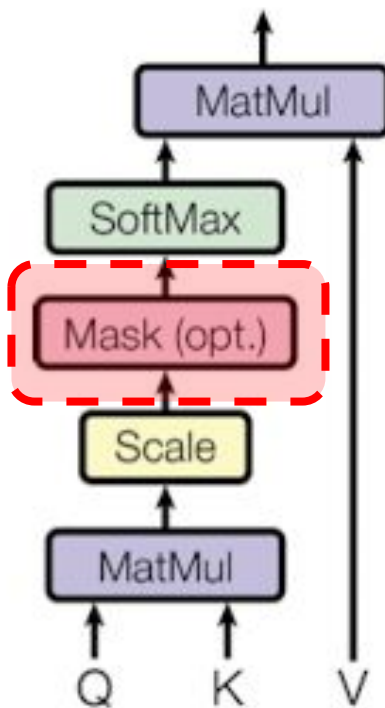


Scaled Dot-Product Attention



Masking

Scaled Dot-Product Attention



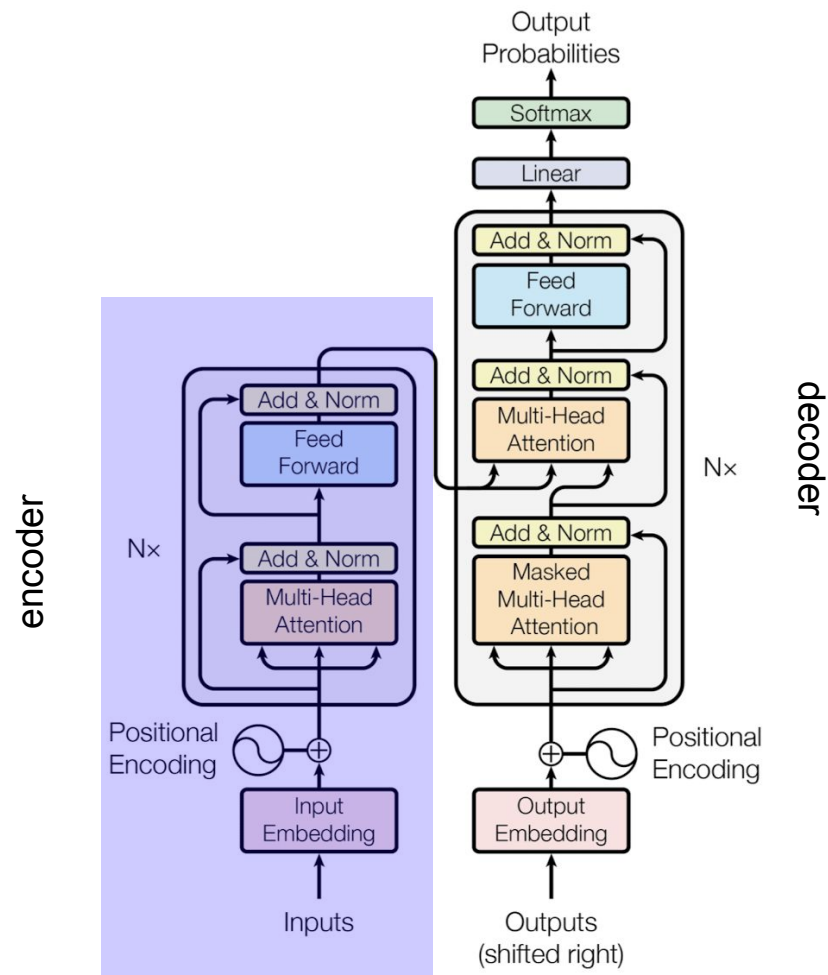


Figure 1: The Transformer - model architecture.

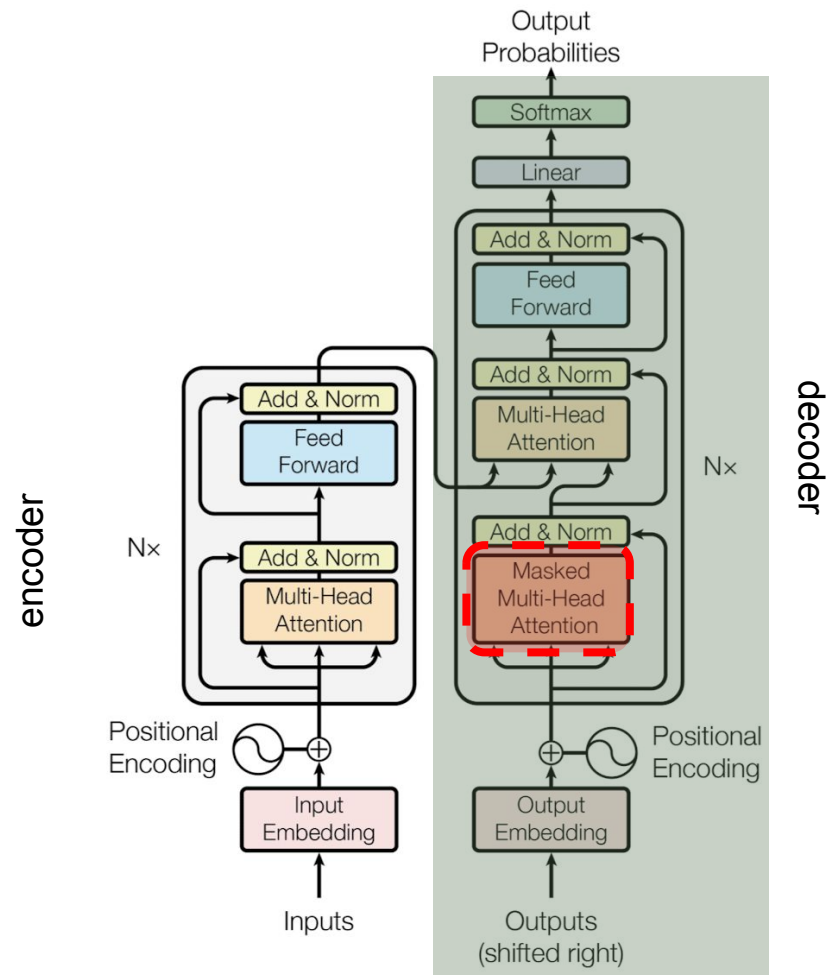


Figure 1: The Transformer - model architecture.

$$\begin{bmatrix} 1 & 0.2 & 0.6 \\ 0.1 & 1 & 0.8 \\ 0.2 & 0.1 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & -inf & -inf \\ 0.1 & 1 & -inf \\ 0.2 & 0.1 & 1 \end{bmatrix}$$

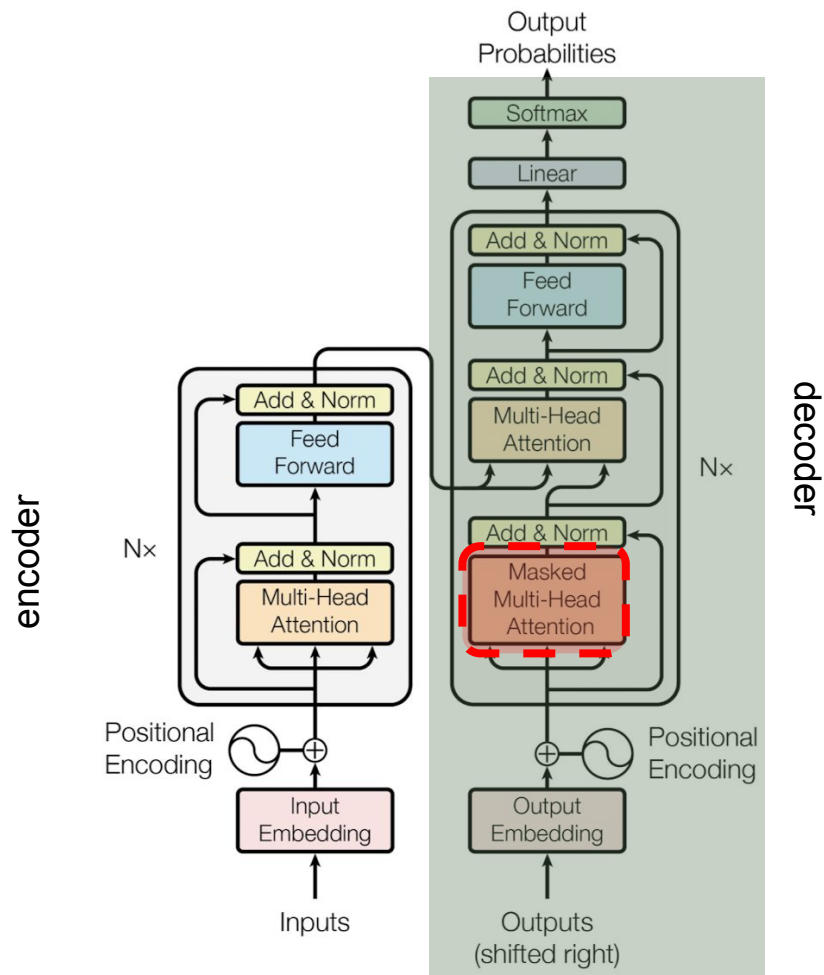


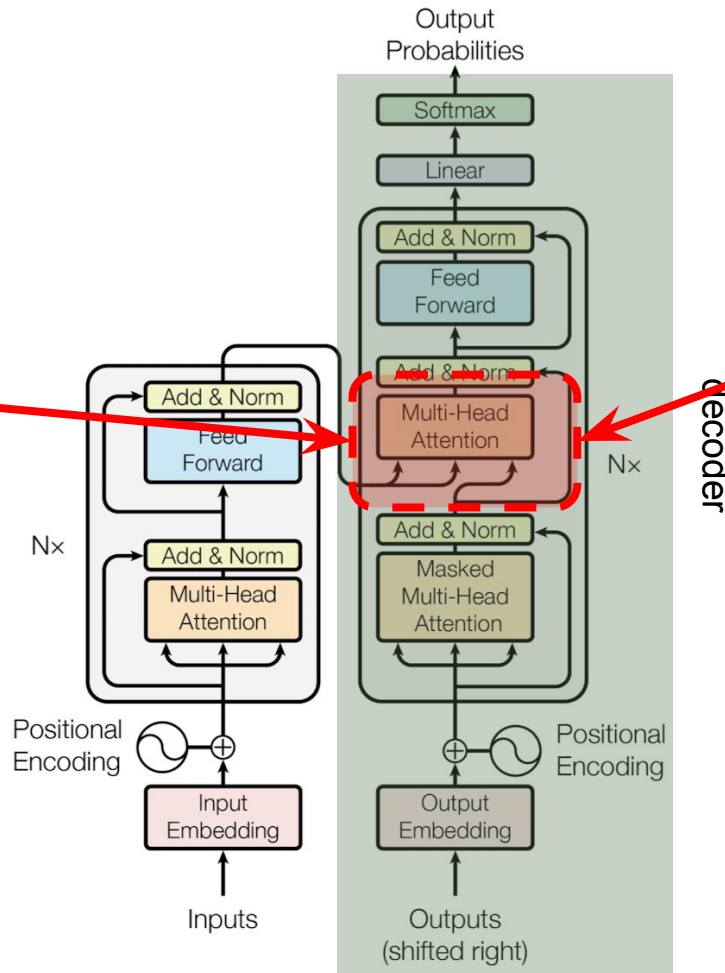
Figure 1: The Transformer - model architecture.

Cross-Attention

K, V coming from encoder

Q coming from decoder

encoder



decoder

Figure 1: The Transformer - model architecture.

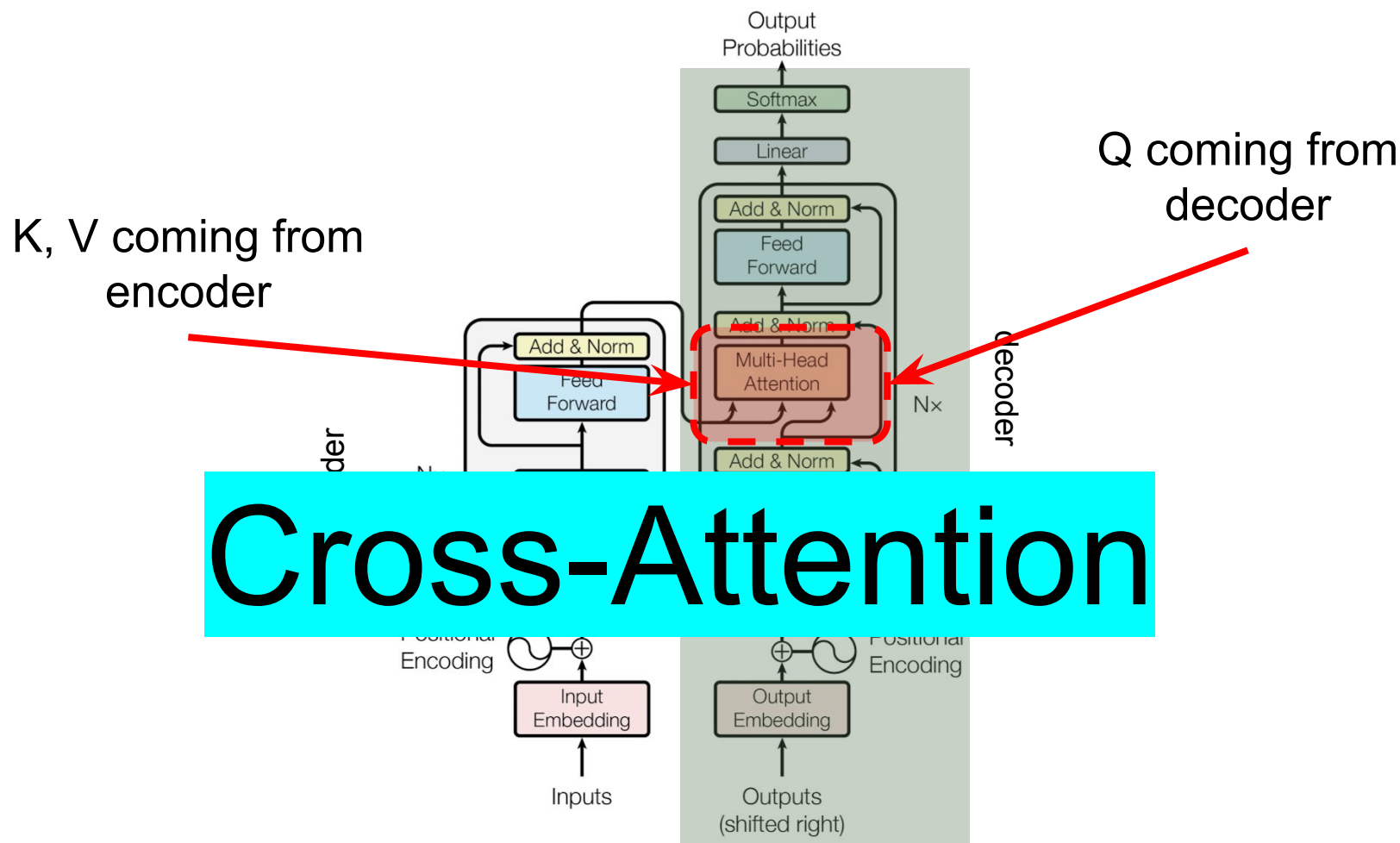
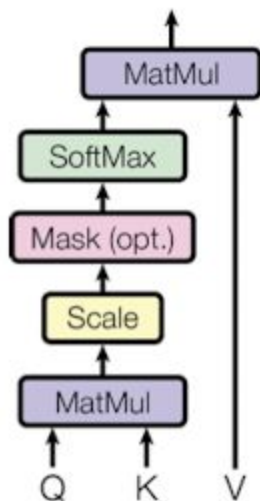


Figure 1: The Transformer - model architecture.

Scaled Dot-Product Attention



Multi-Head Attention

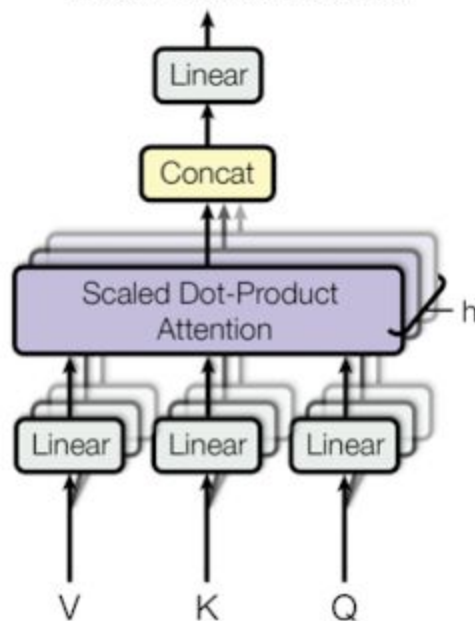


Figure 2: (left) Scaled Dot-Product Attention. (right) Multi-Head Attention consists of several attention layers running in parallel.

Regularization Methods

- Layer Norms
- Residual Connections

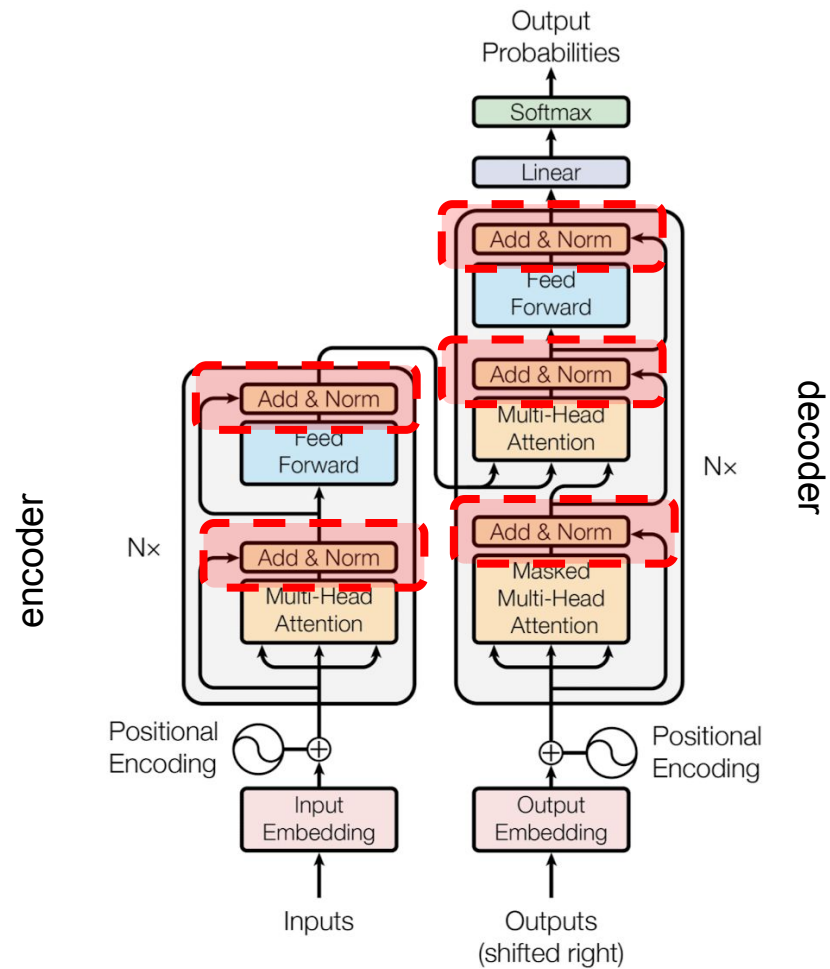


Figure 1: The Transformer - model architecture.

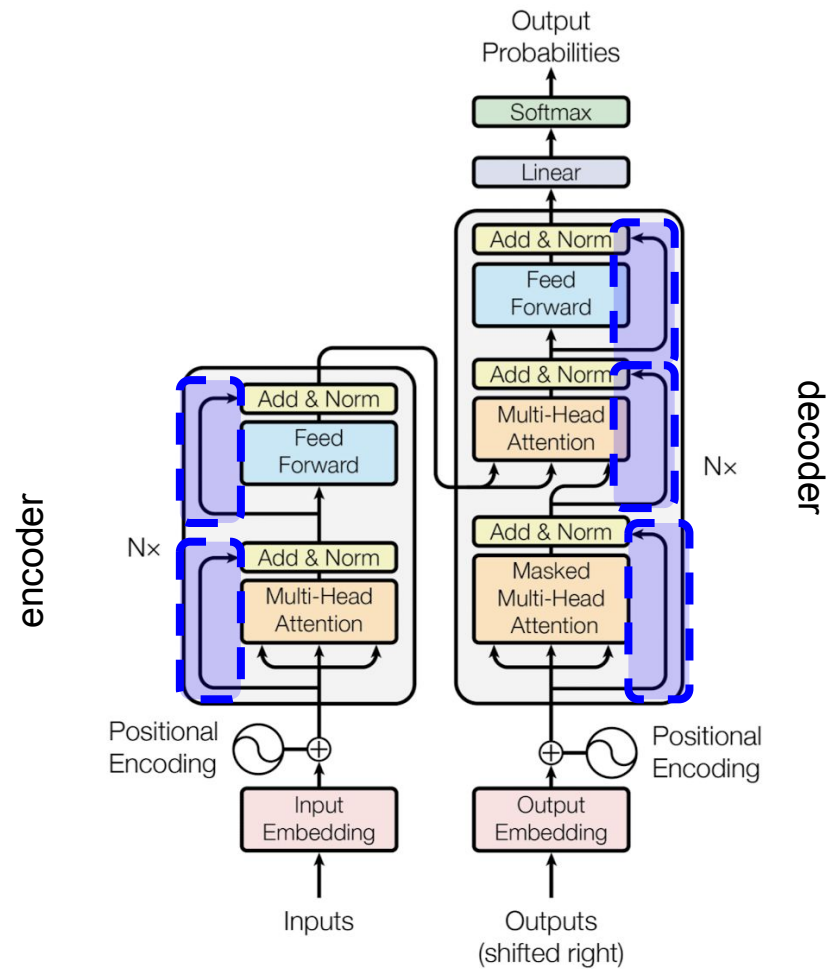


Figure 1: The Transformer - model architecture.

Overview

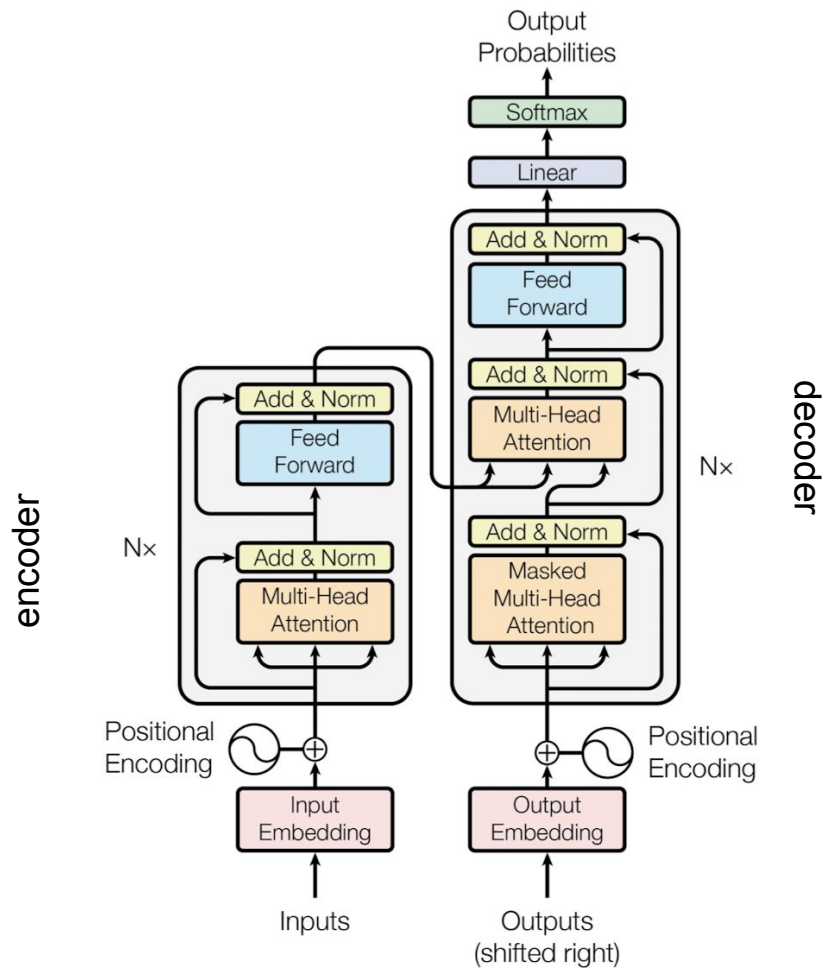


Figure 1: The Transformer - model architecture.

Transformer Models You Might Know

- **Encoder:** BERT
- **Decoder:** GPT
- **Encoder/Decoder:** Machine Translation