

# Transformer

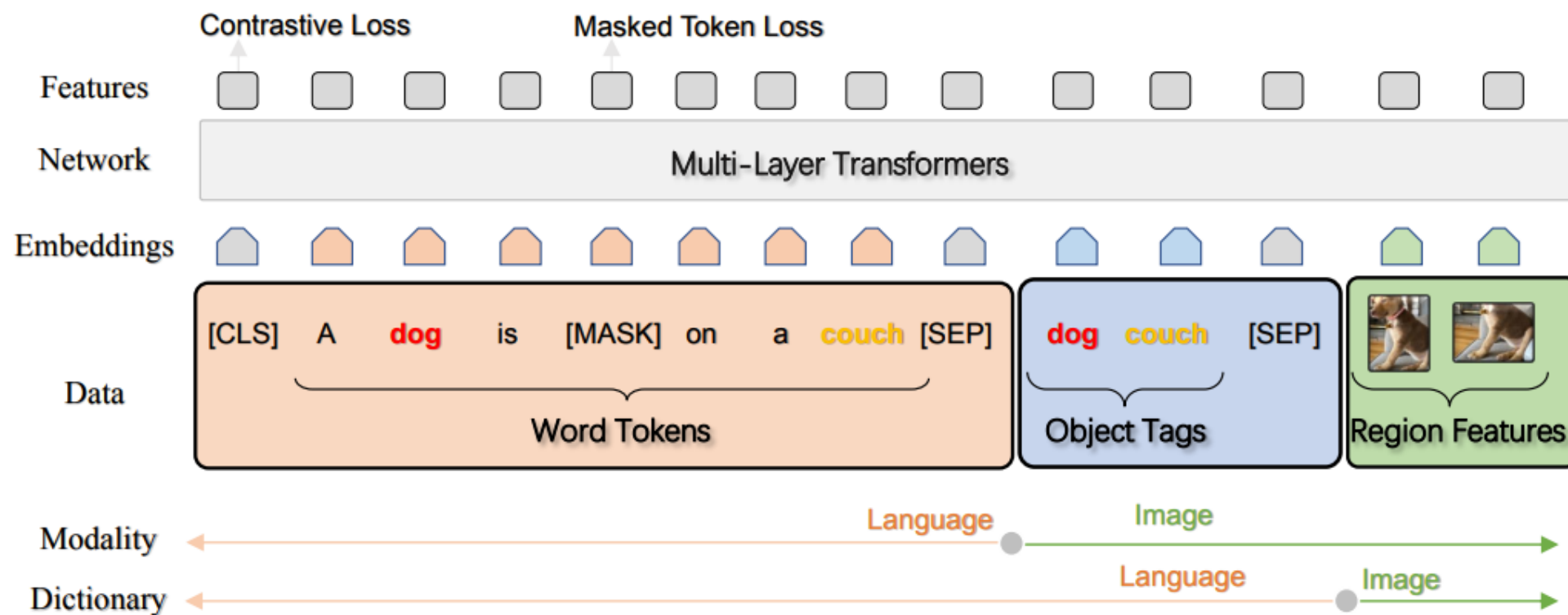
Machine Learning Study  
JinHo Kim

# Contents

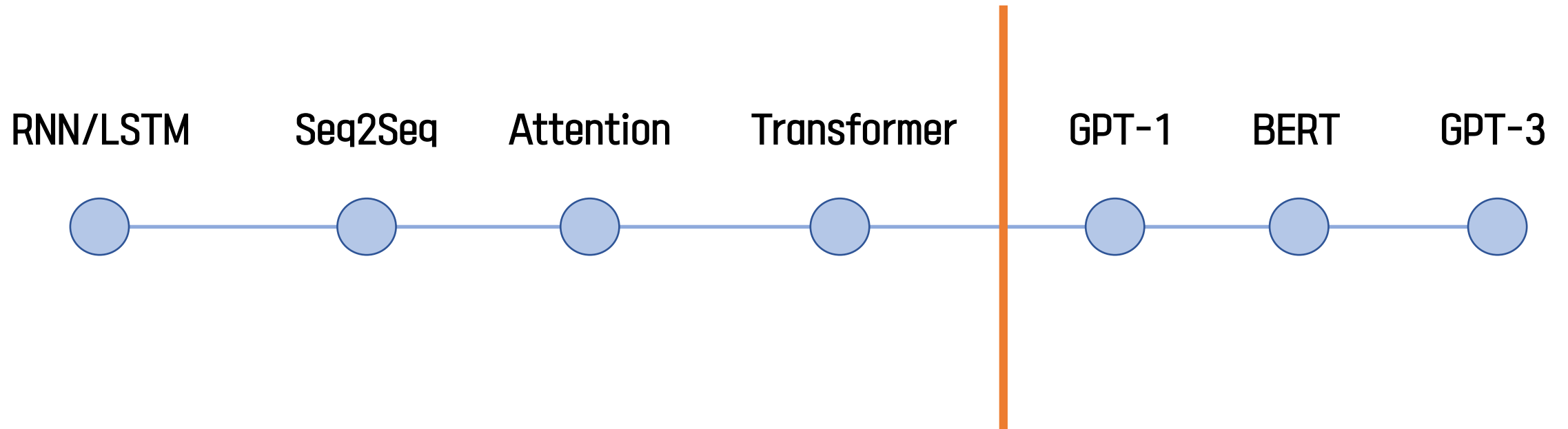
1. Introduction
2. Paper review
3. Code Practicae

# 1. Introduction

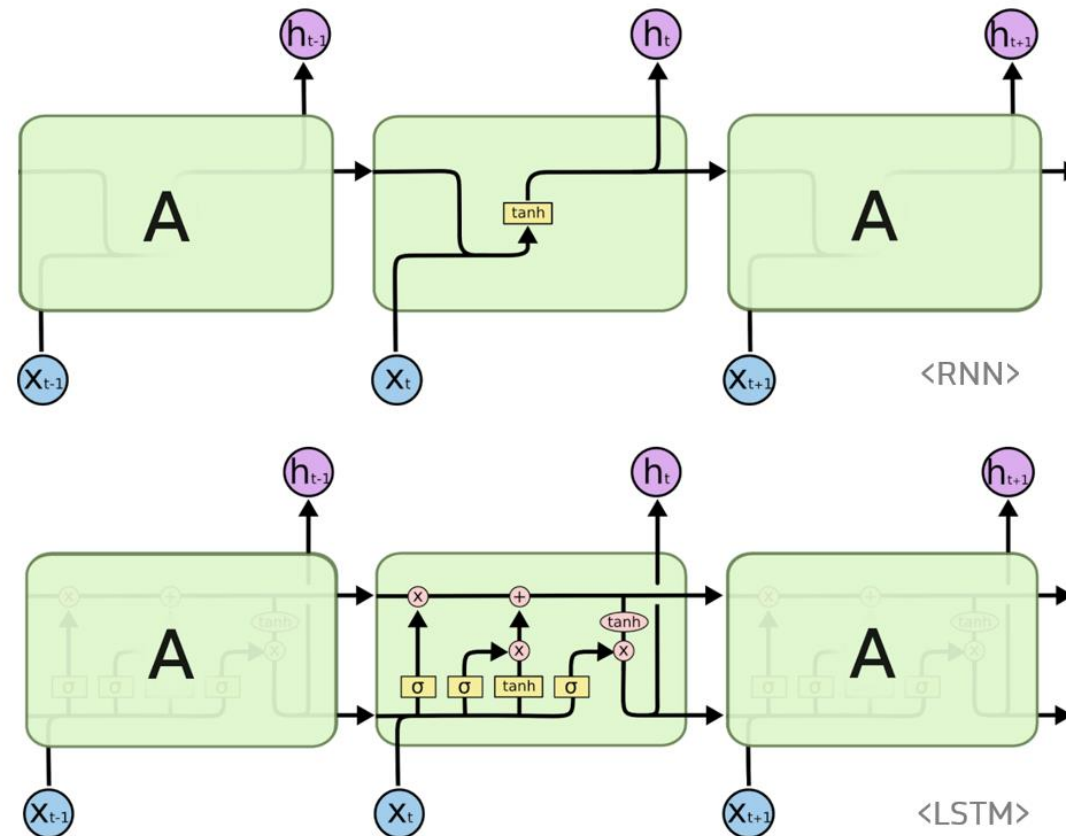
# OSCAR...?



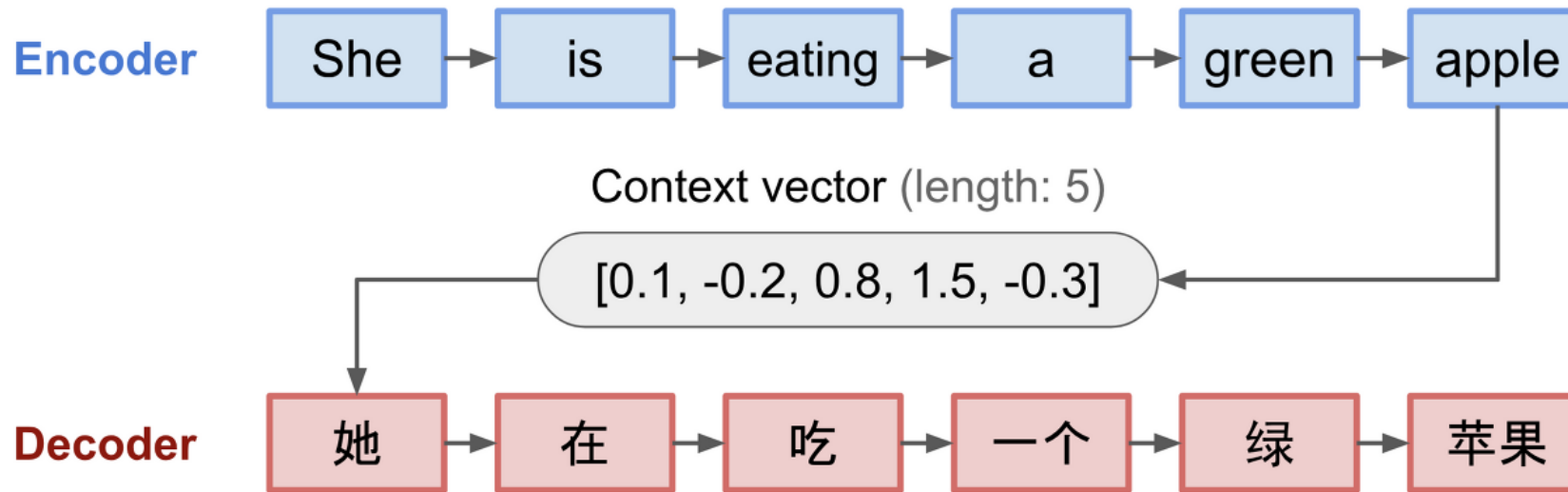
# NLP Model History



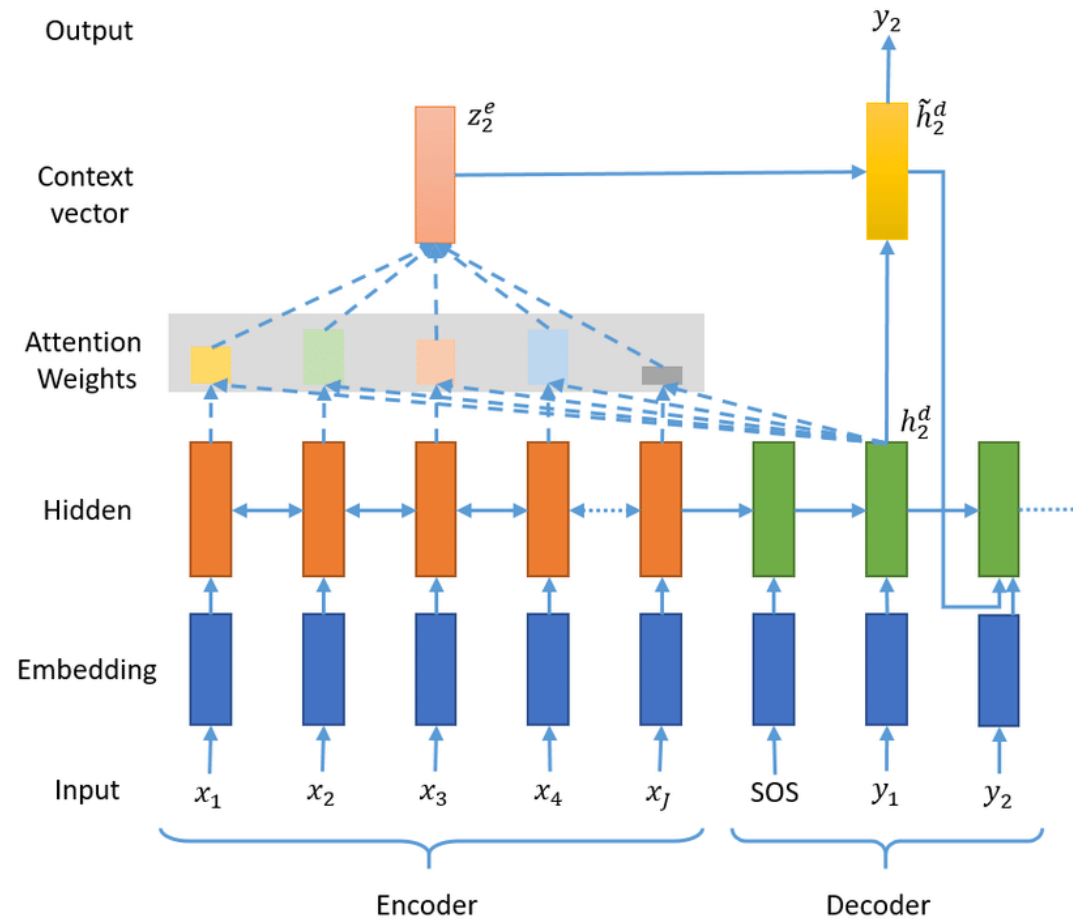
# RNN / LSTM



# Seq2Seq



# Seq2Seq with Attention





## 2. Paper Review

**‘Attention is All You Need’**

# ‘Attention is All You Need’

the Transformer,  
base solely on **attention mechanism**,  
dispense with **recurrence** and **convolutions** entirely

→ establishes a new model **state-of-the-art**

# ‘Attention is All You Need’

Parallelizable ?  Sequential ?

# 'Attention is All You Need'

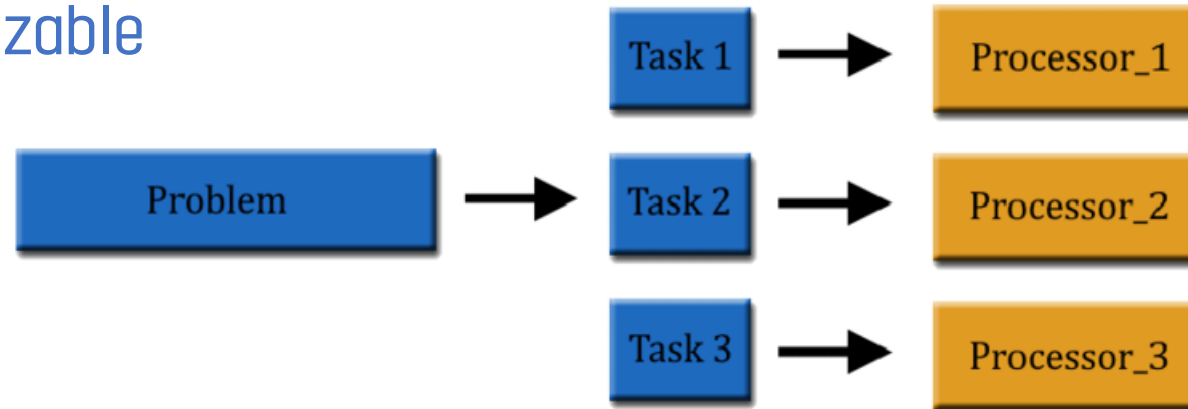
Sequential (Serial)

Serial Computing



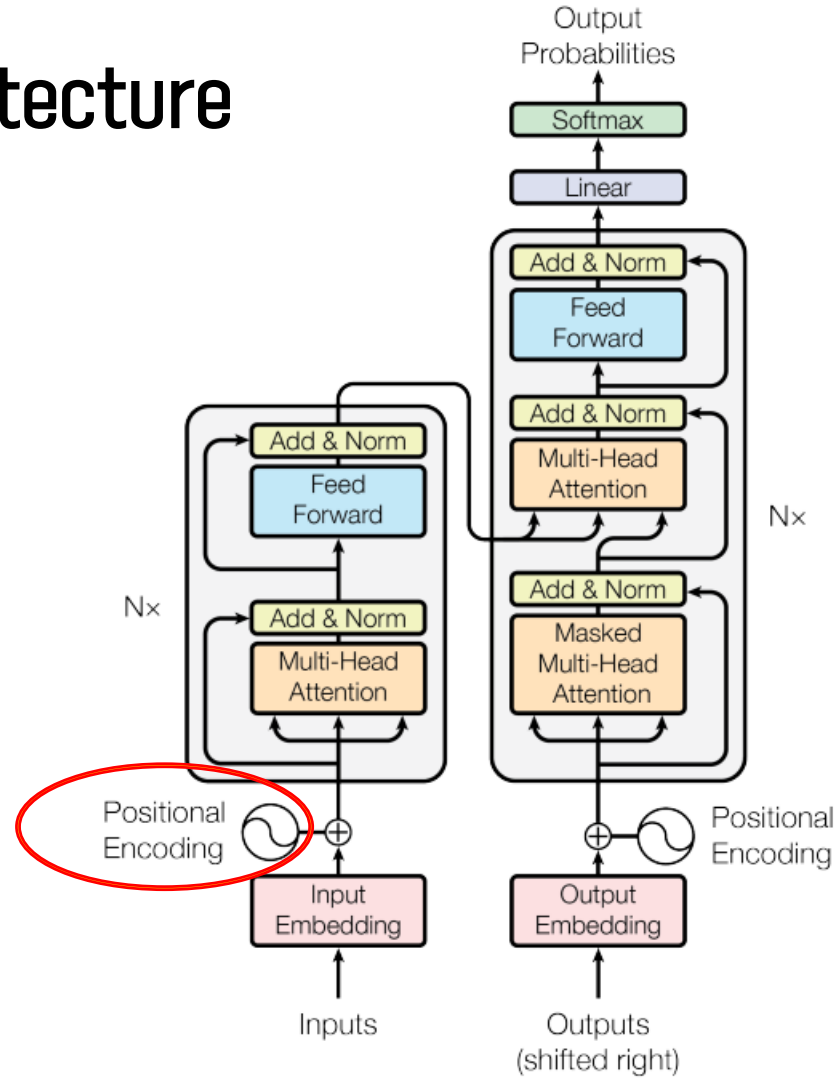
Parallel Computing

Parallelizable



# Transformer architecture

Encoder

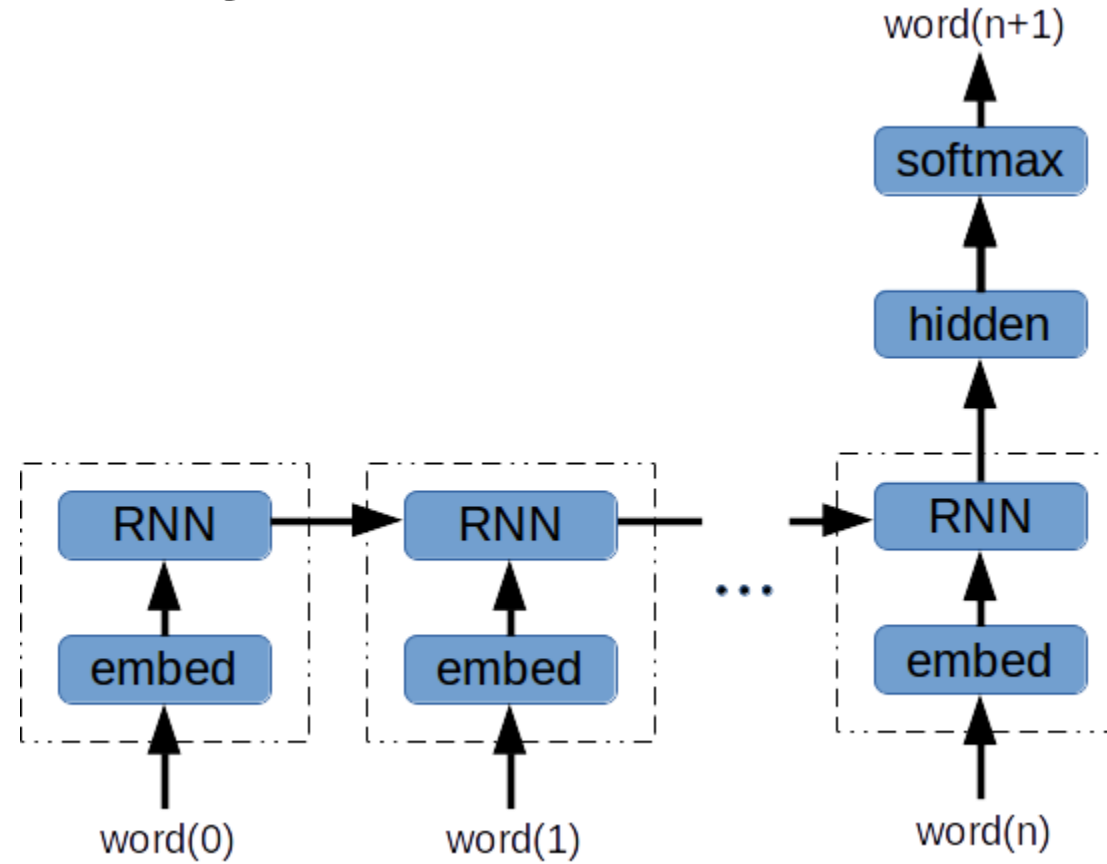


Decoder

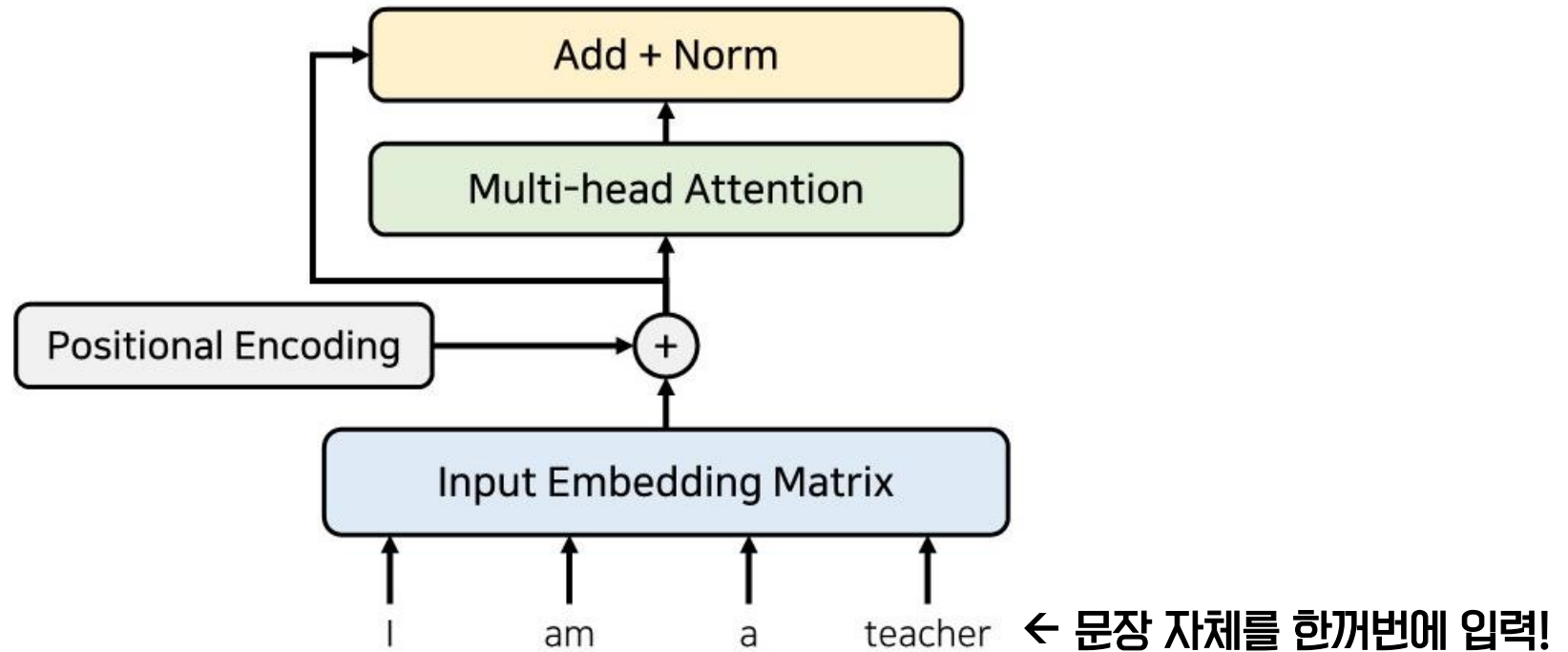
Figure 1: The Transformer - model architecture.

# RNN - Word Embedding

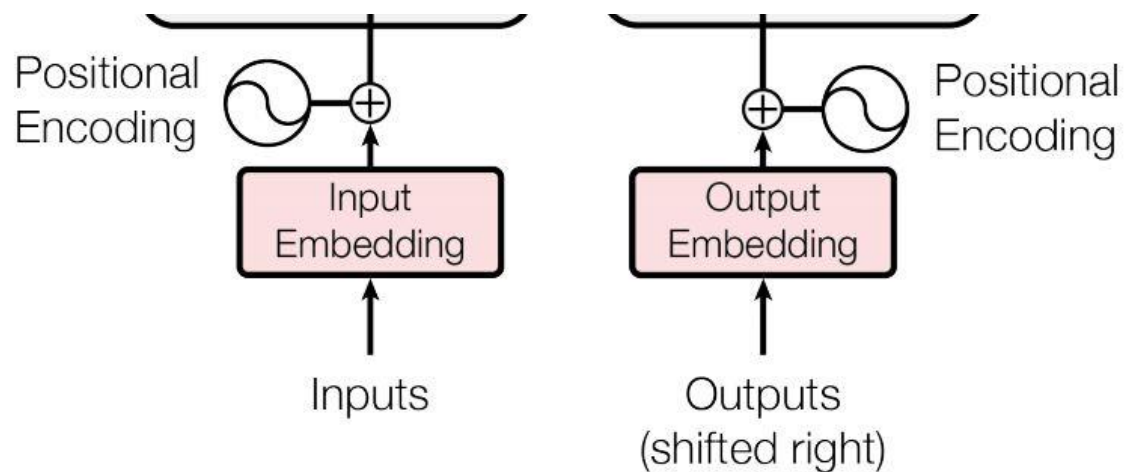
word order!



# Postional Encoding



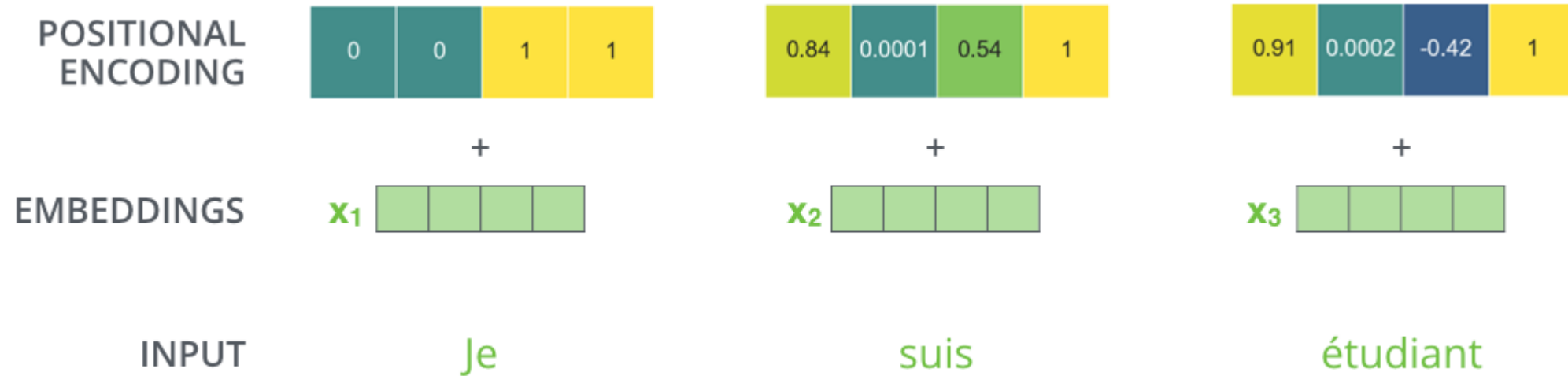
# Positional Encoding



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



# Positional Encoding



# Transformer architecture

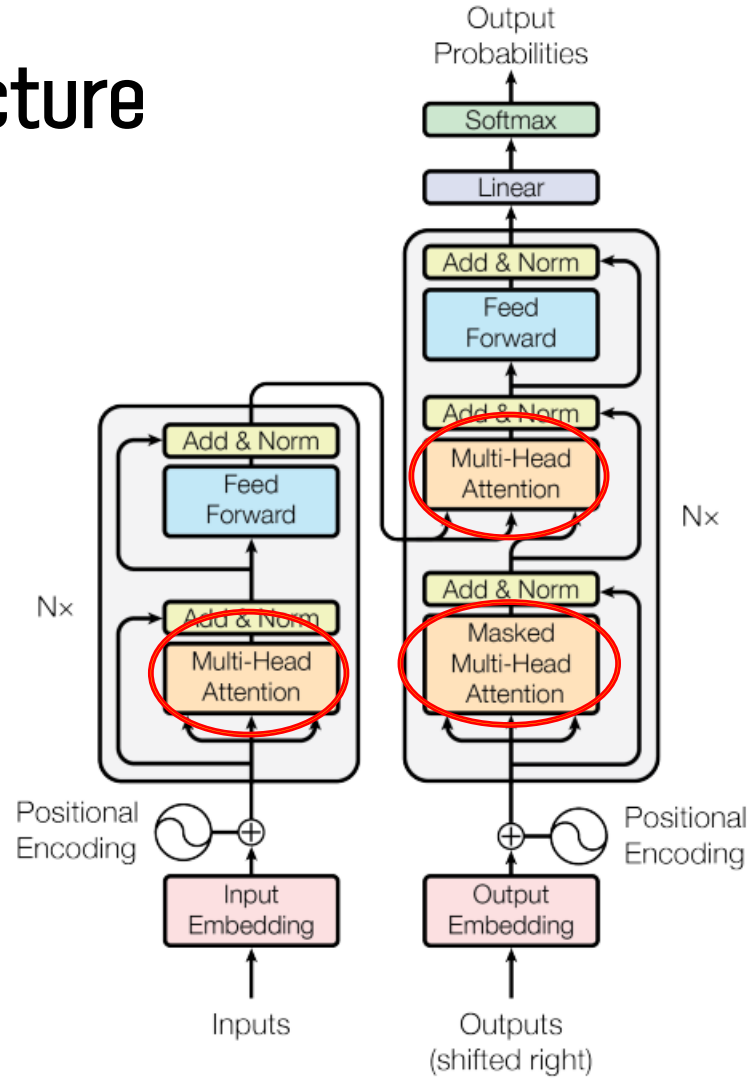
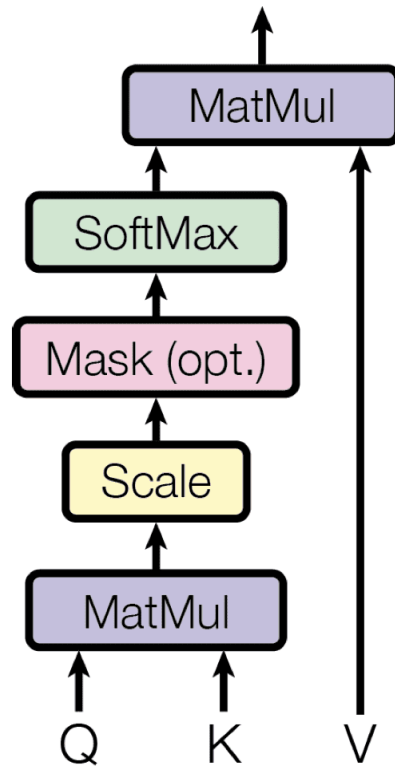


Figure 1: The Transformer - model architecture.

# Scaled Dot-Product Attention



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

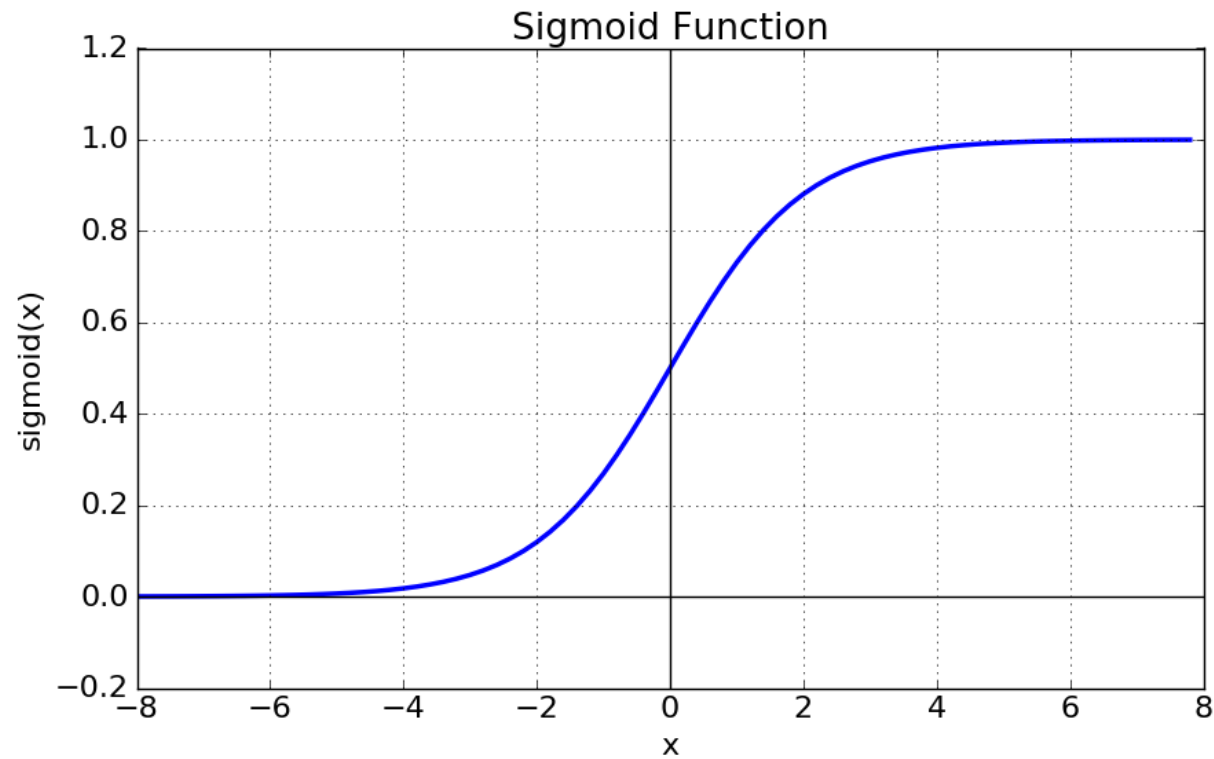
# Dot-Product Attention

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

$\frac{1}{\sqrt{d_k}}$   **'Scale'**

extremely small gradients

# Sigmoid function

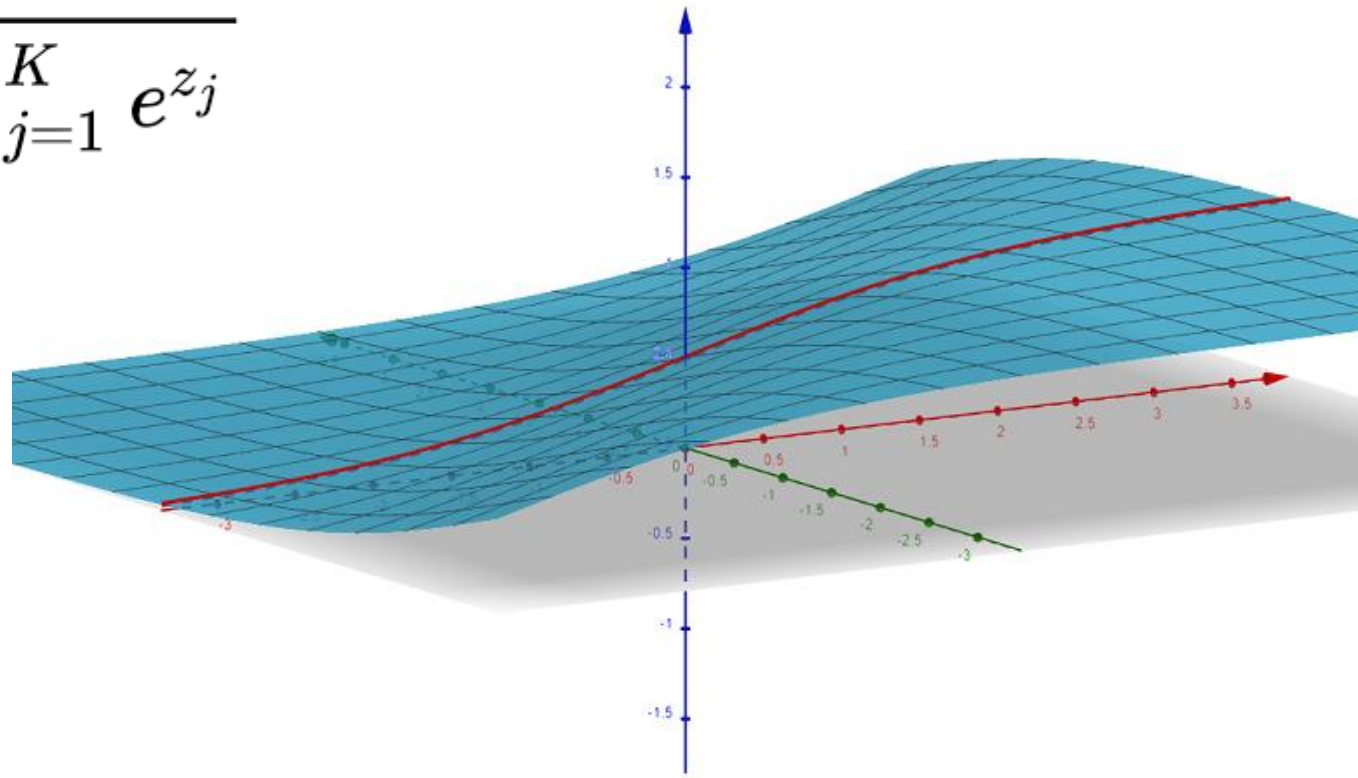


$$S(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1}$$

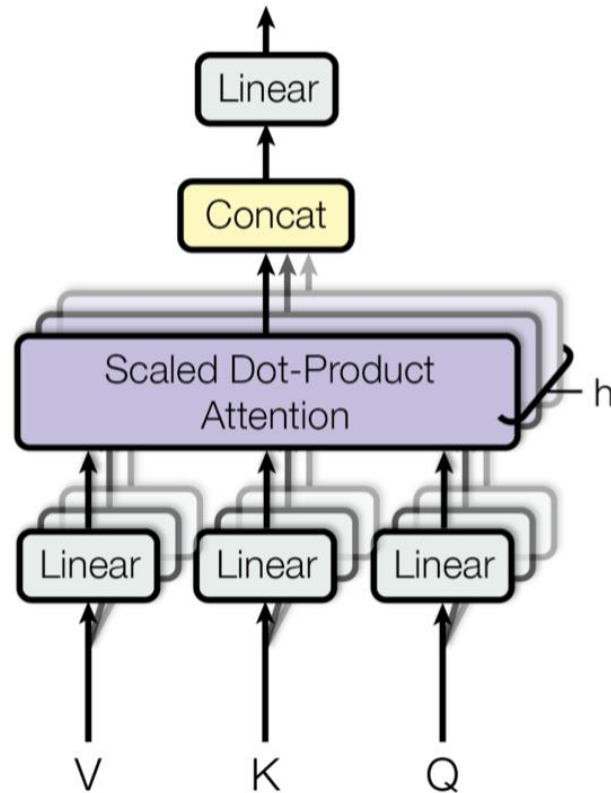
# Softmax function

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

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$



# Multi-Head Attention



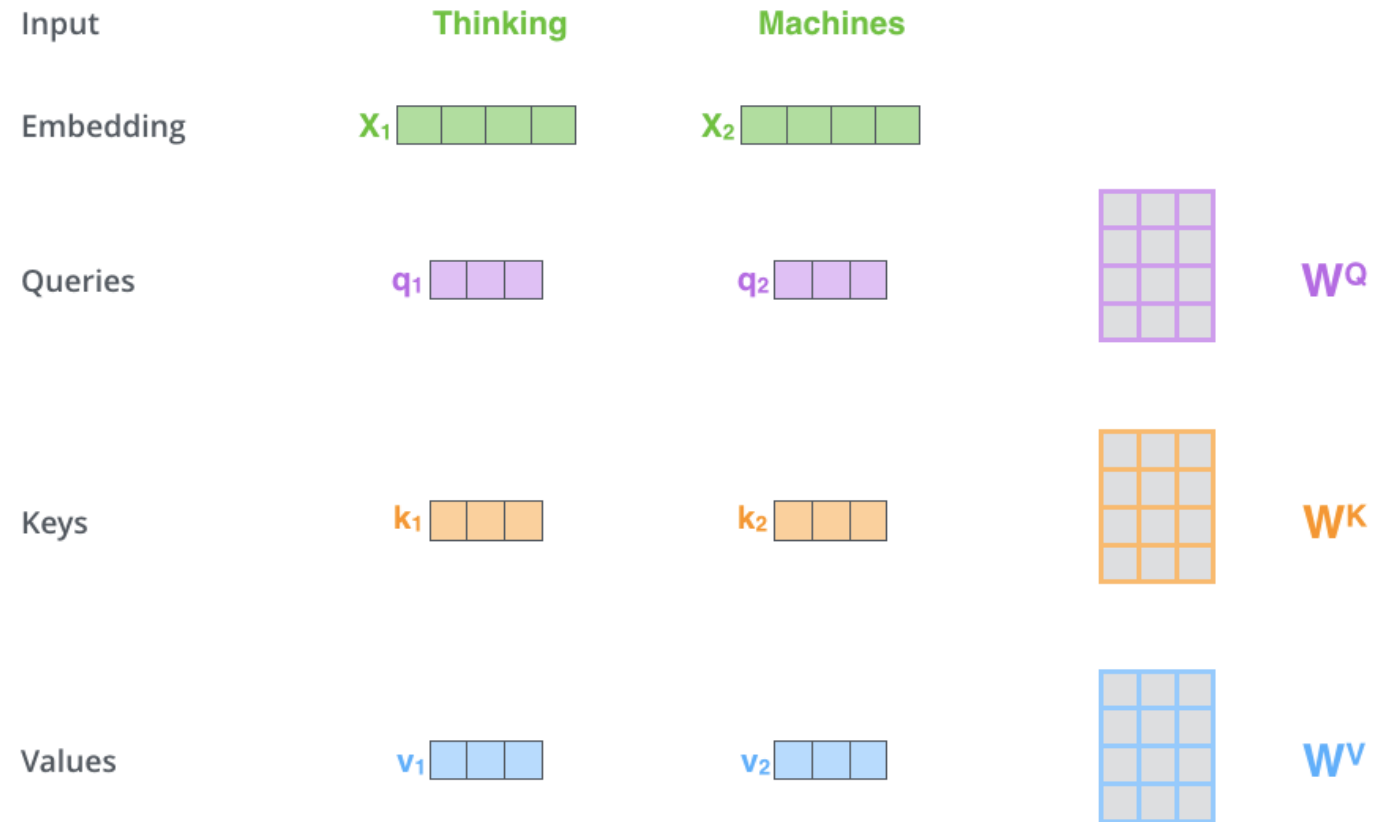
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where  $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

head : h \* scaled dot-product attention

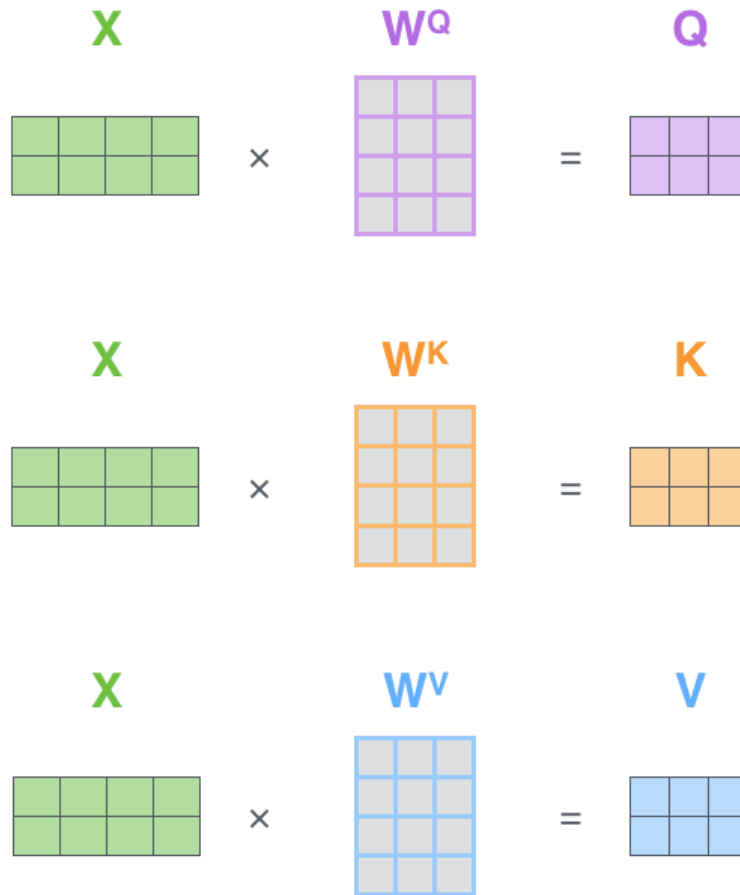
model parameter  $\rightarrow h = 8$  (base)

# Scaled Dot-Product Attention (1)





## Scaled Dot-Product Attention (2)



## Scaled Dot-Product Attention (3)

$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{array}{|c|c|} \hline \square & \square \\ \hline \square & \square \\ \hline \square & \square \\ \hline \end{array} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$

=

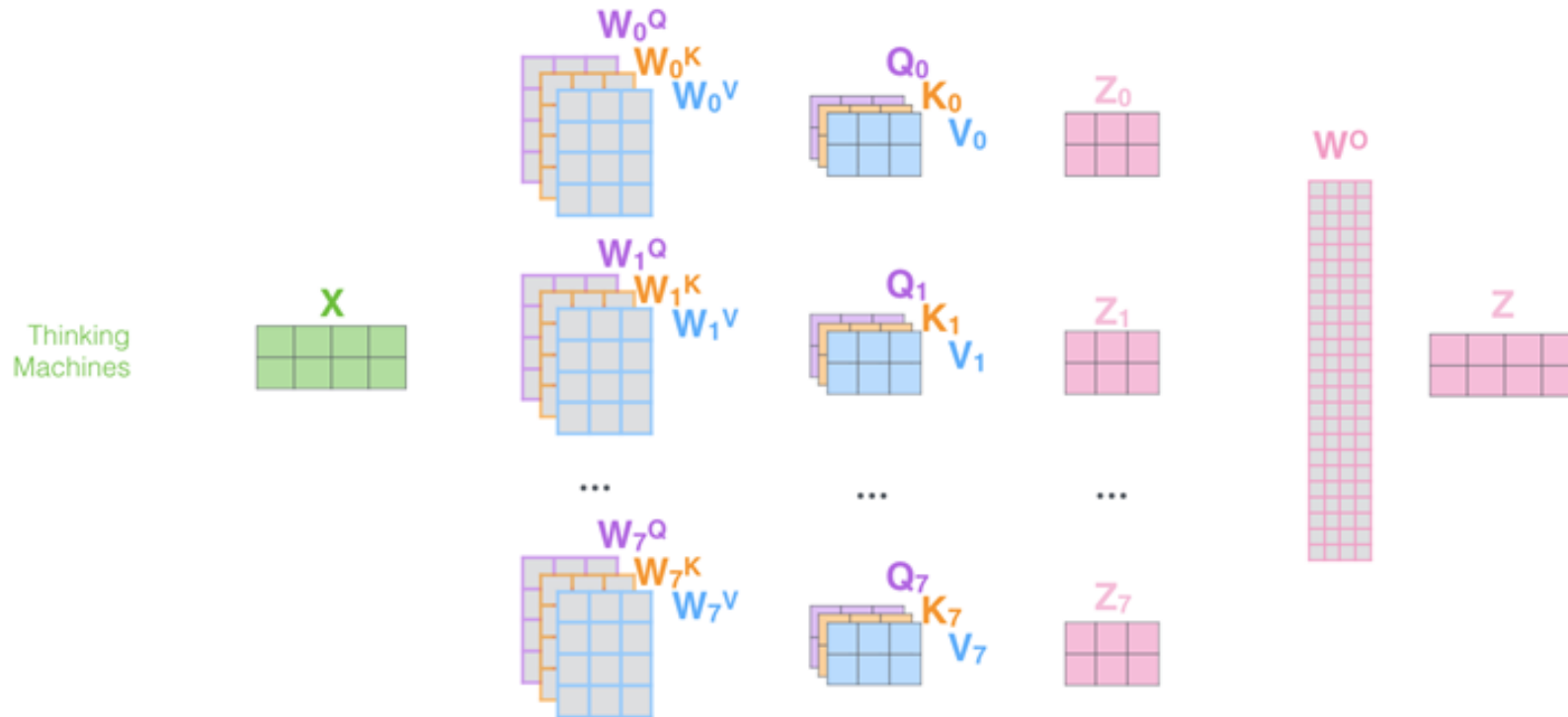
$\begin{matrix} \text{Z} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$

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

# Multi-Head Attention

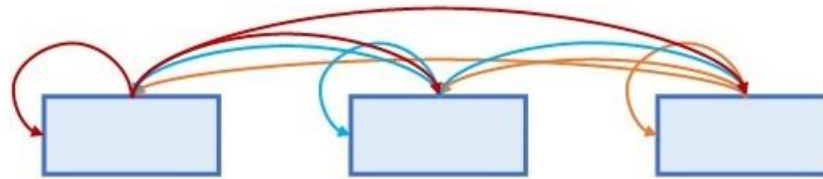
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

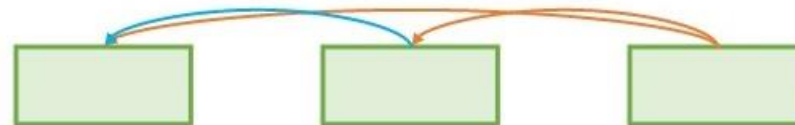


## multi-head attention in 3 different ways :

Encoder Self-Attention:



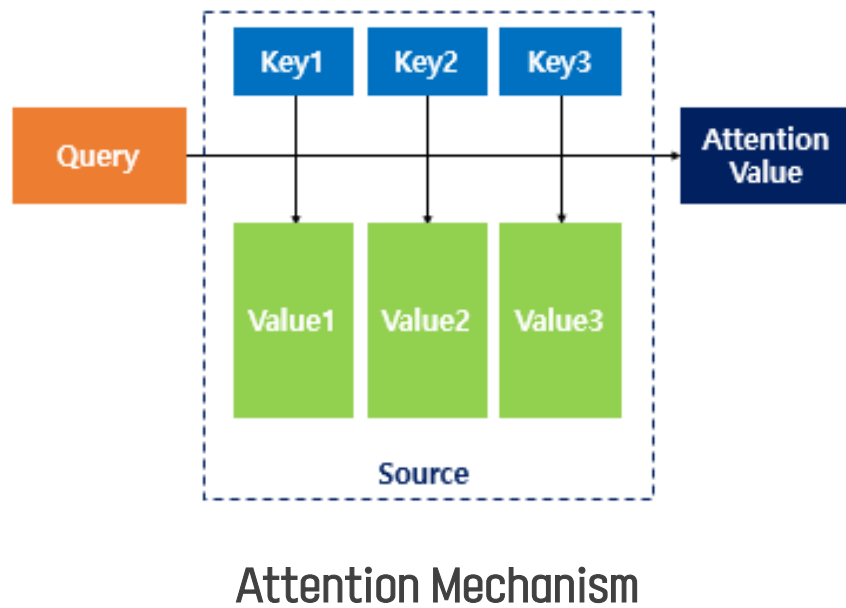
Masked Decoder Self-Attention:



Encoder-Decoder Attention:

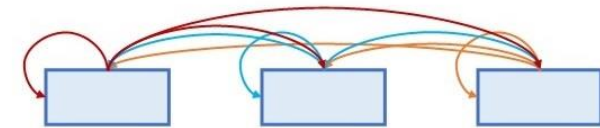


# Self-Attention



In Self-Attention,  
Query = Key = Vector

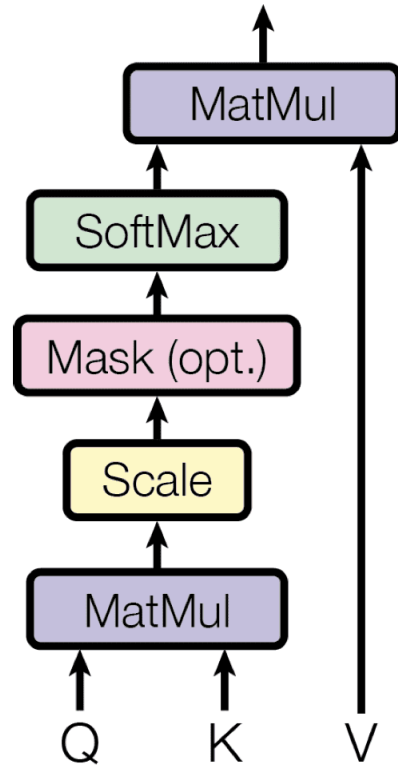
Encoder Self-Attention:



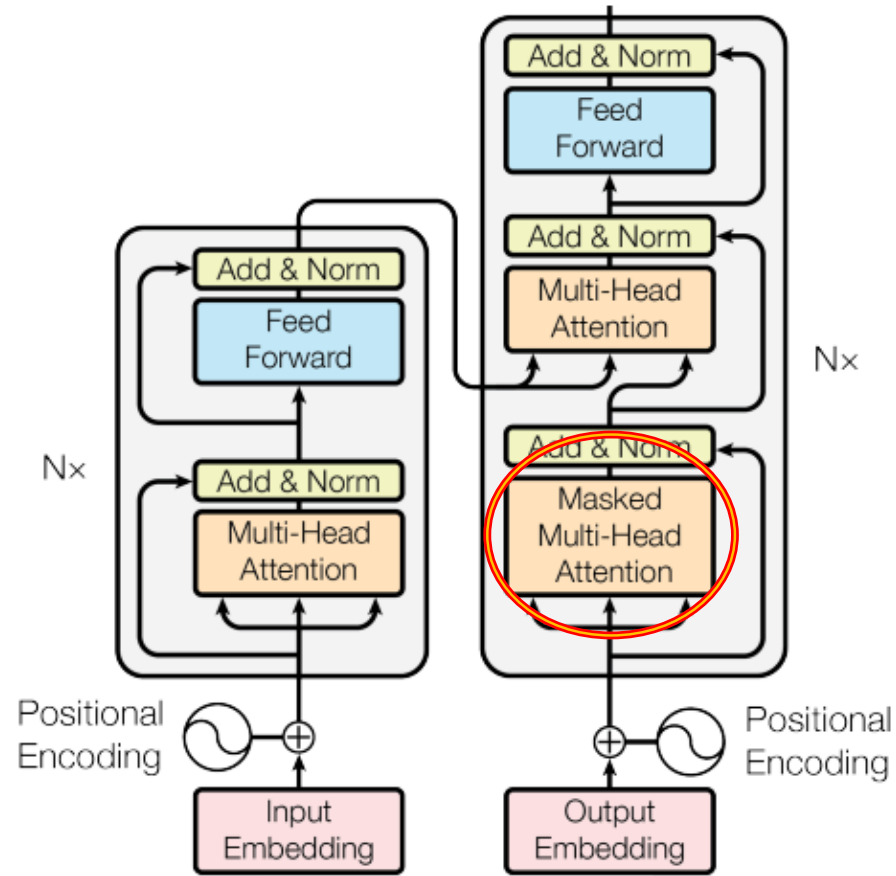
Masked Decoder Self-Attention:



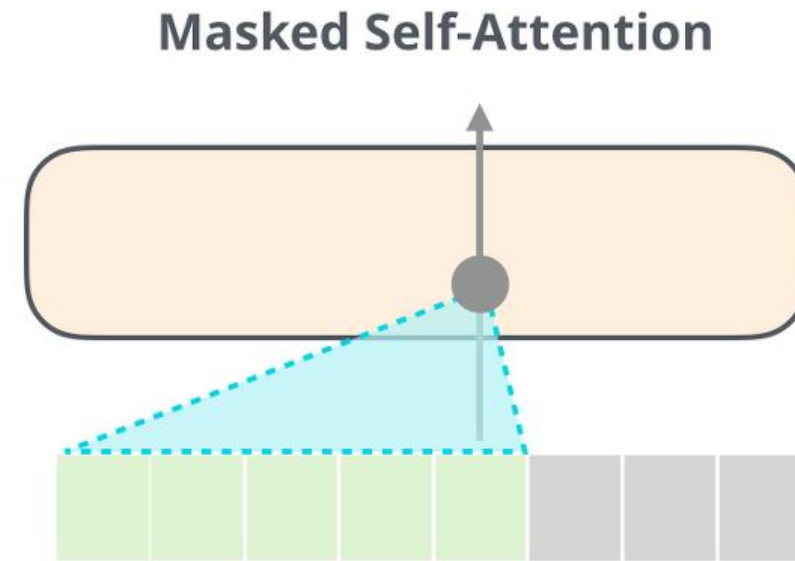
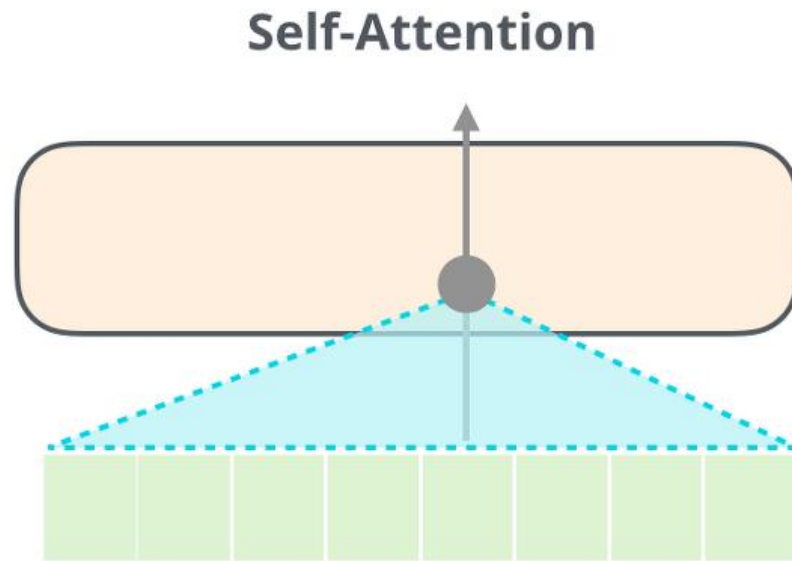
# 'Masked' Self-Attention



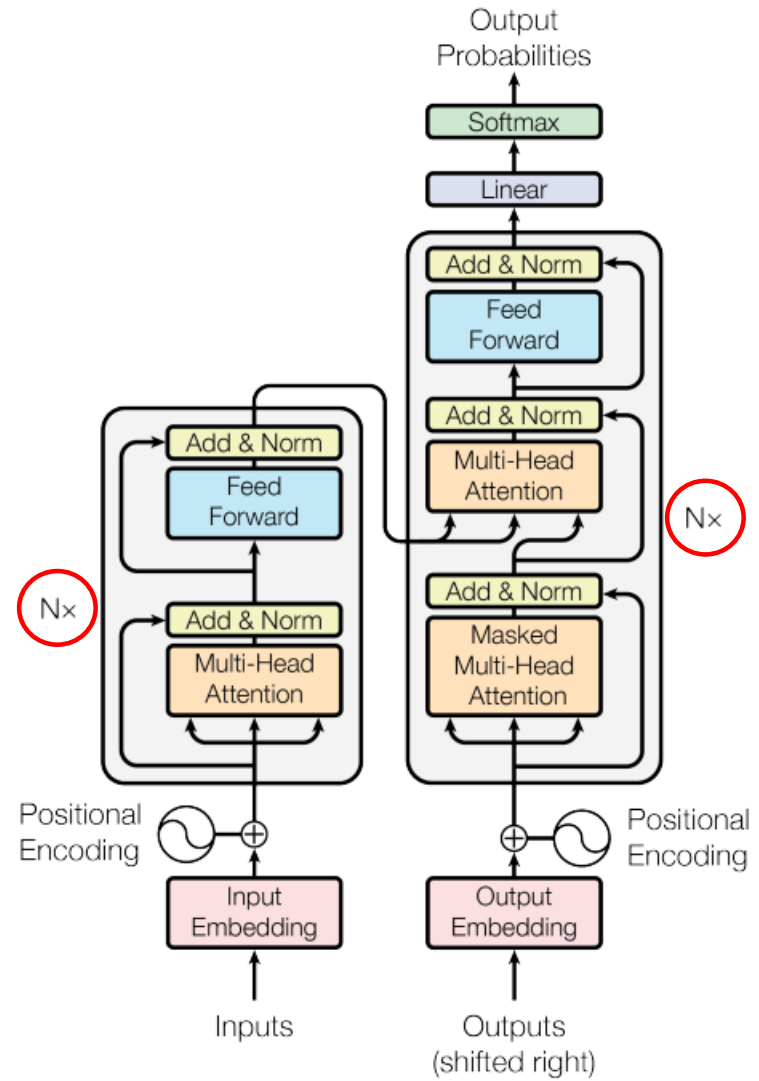
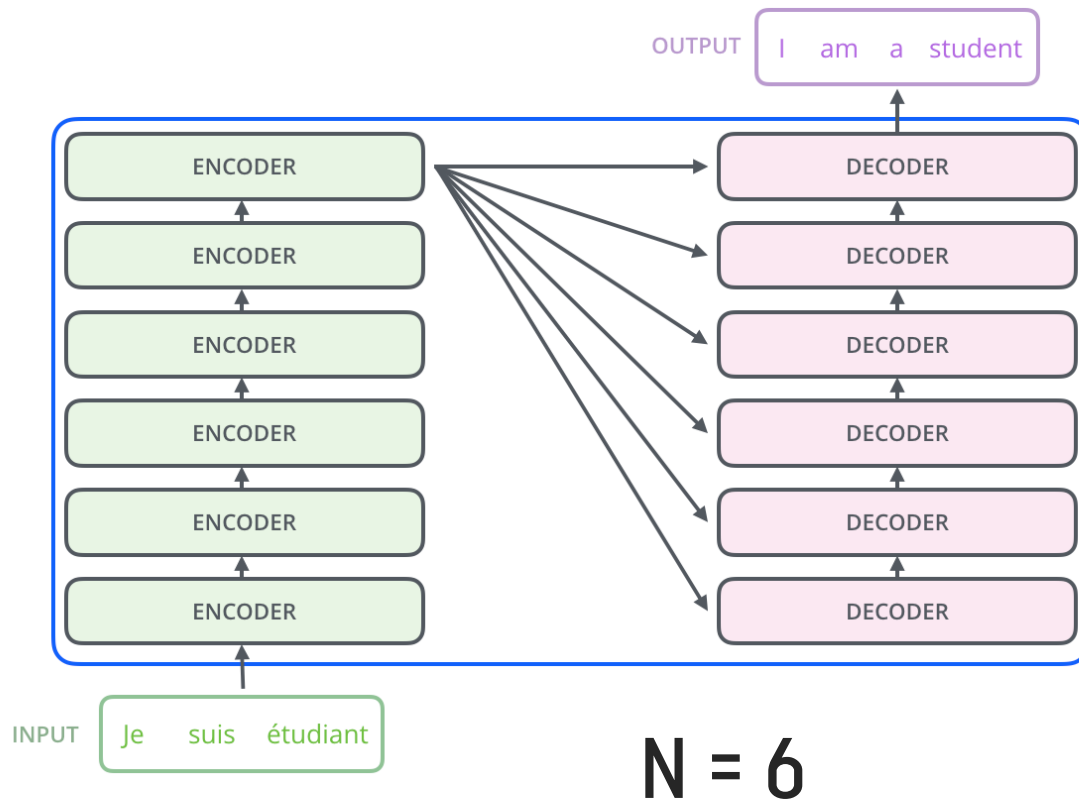
Scaled Dot-Product Attention



# 'Masked' Self-Attention



# Encoder Decoder Attention



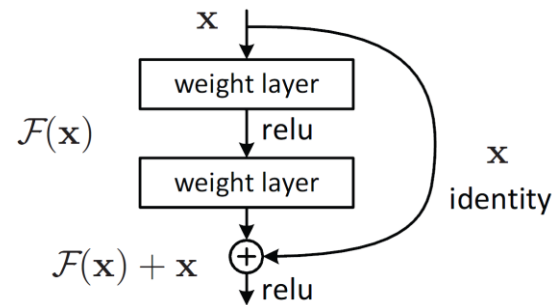


# Feed Forward network

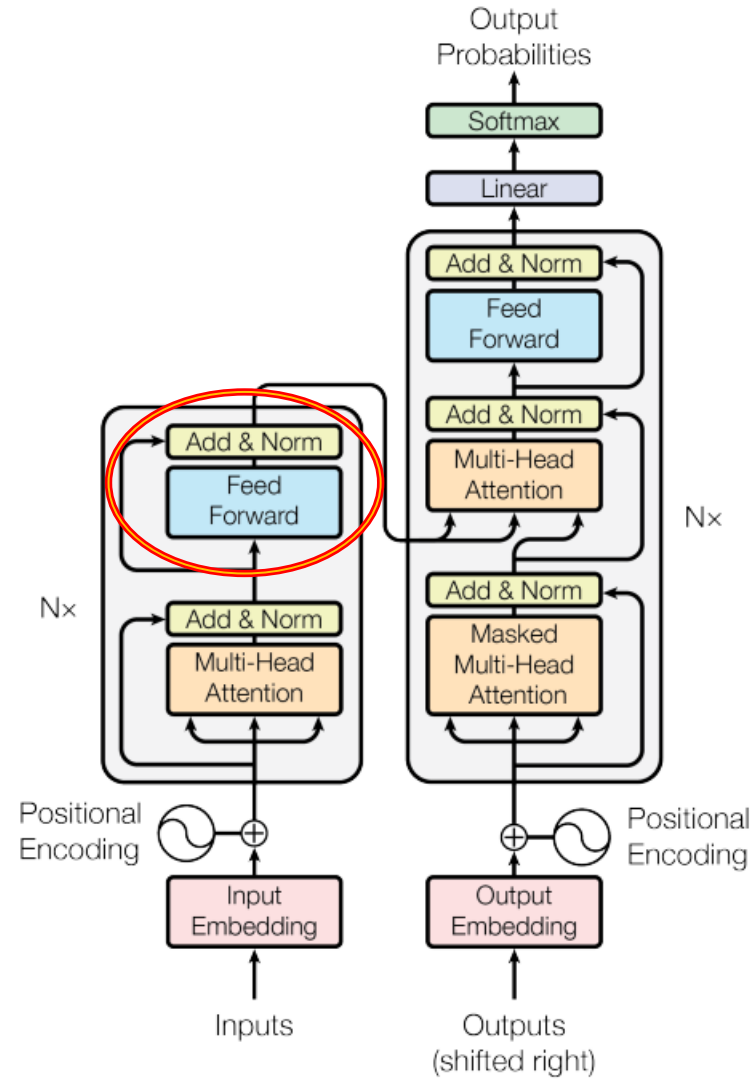
$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

ReLU Function  $\rightarrow$  ELU, GeLU, etc.,

## Residual connection



In ResNet



# 3. Code Practice

**QnA?**

**Thank You!**