

# Procesamiento de Lenguaje Natural

Clase 13 – Transformers

Ph.D. Rubén Manrique

[rf.manrique@uniandes.edu.co](mailto:rf.manrique@uniandes.edu.co)

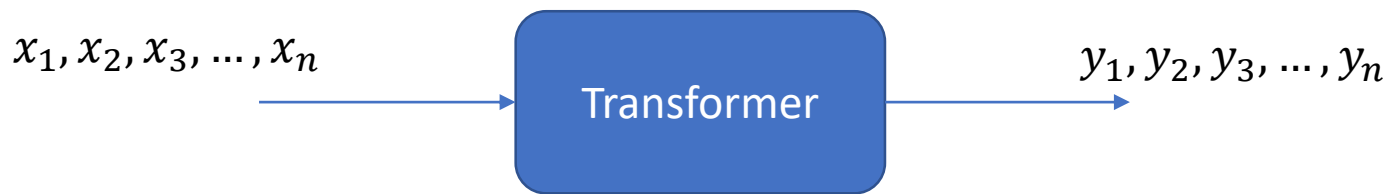
Maestría en Ingeniería de Sistemas  
y Computación



- Chapter 9/11: Transformers

# Transformers

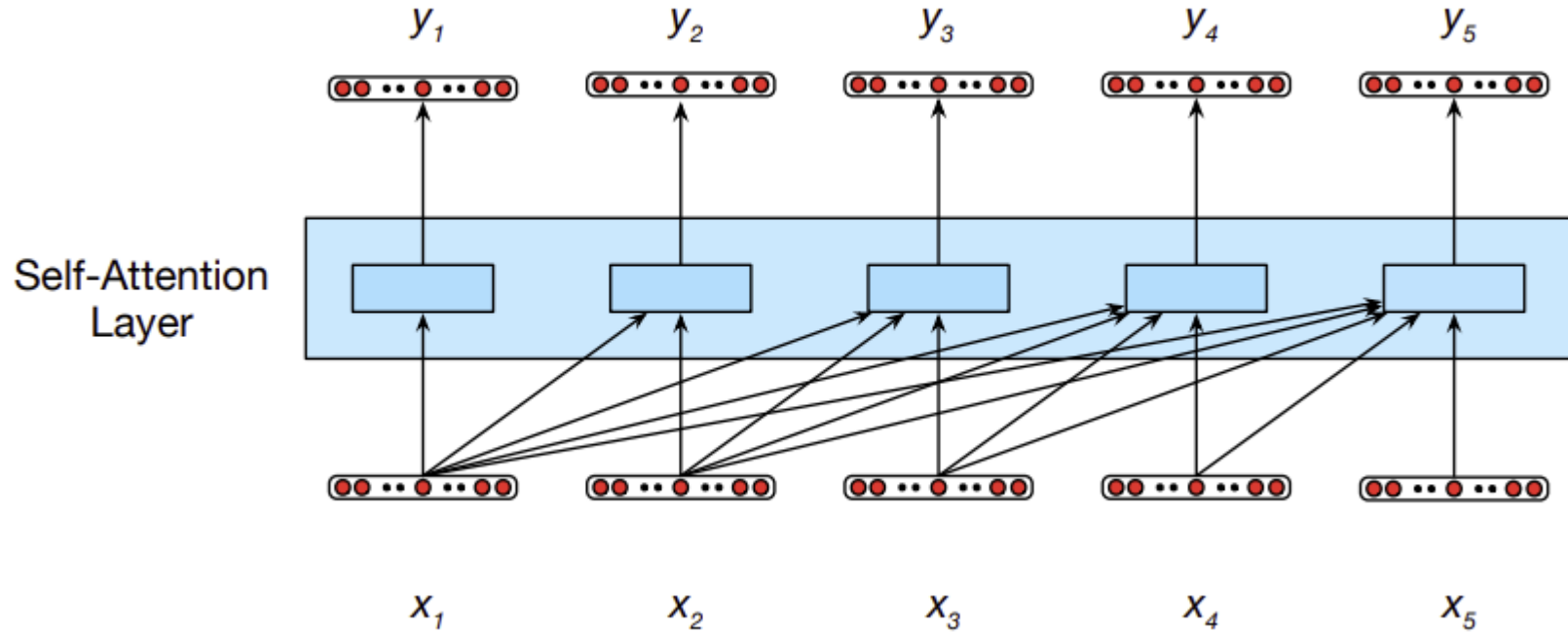
- An approach to sequence processing that eliminates recurrent connections and returns to architectures reminiscent of the fully connected networks.



Simple feedforward networks  
+  
**Self-attention** layers

# Self-attention (I)

- Self-attention allows a network to directly extract and use information from arbitrarily large contexts.
- There is not need to pass it through intermediate recurrent connections as in RNNs.



- The model has access to all of the inputs up to and including the one under consideration.
- No access to information about inputs beyond the current one.
- The computation of  $y_3$  is based on a set of **comparisons** between the input  $x_3$  and its preceding elements  $x_1$  and  $x_2$ , and to  $x_3$  itself.

# Self-attention (II): Ability to compare

- The attention-based approach is the ability to compare an item of interest to a collection of other items in way that reveals their relevance in the current context.
- The simplest form of comparison between elements in a self-attention layer is a dot product.

$$\textit{score}(x_i, x_j) = x_i \cdot x_j$$

- The larger the value the more similar the vectors that are being compared.

# Self-attention (II): Compute the output

- Suppose you want to compute  $y_3$ , so you compare the different inputs:

- $x_3 \cdot x_1$
  - $x_3 \cdot x_2$
  - $x_3 \cdot x_3$

*a set of comparisons to relevant items in some context*

- How to evaluate the proportion of relevance of each input?
  - Create a vector of weights  $\alpha$

$$\alpha_{ij} = \text{softmax}(\text{score}(x_i, x_j)) \quad \forall j \leq i$$

*Normalization to  
provide a probability  
distribution*

$$y_i = \sum_{j \leq i} \alpha_{ij} x_j$$

*weighted sum using  
this distribution*

# We need extra parameters....

- Transformers include additional parameters in the form of a set of weight matrices that operate over the input embeddings.
- Input embeddings roles:
  - **Query**: As the current focus of attention when being compared to all of the other preceding inputs.
  - **Key**: In its role as a preceding input being compared to the current focus of attention.
  - **Value**: as a value used to compute the output for the current focus of attention.
- Transformers introduce three sets of weights:
  - $W^Q, W^K, W^V$
  - These weights will be used to compute linear transformations of each input  $x$ .
  - The resulting values being used in their respective roles in subsequent calculations.



# We need extra parameters....

- Linear transformations:

- $q_i = W^Q x_i, \quad W^Q \in \mathbb{R}^{d_q \times d_m}, \quad d_m=1024, d_q=64$

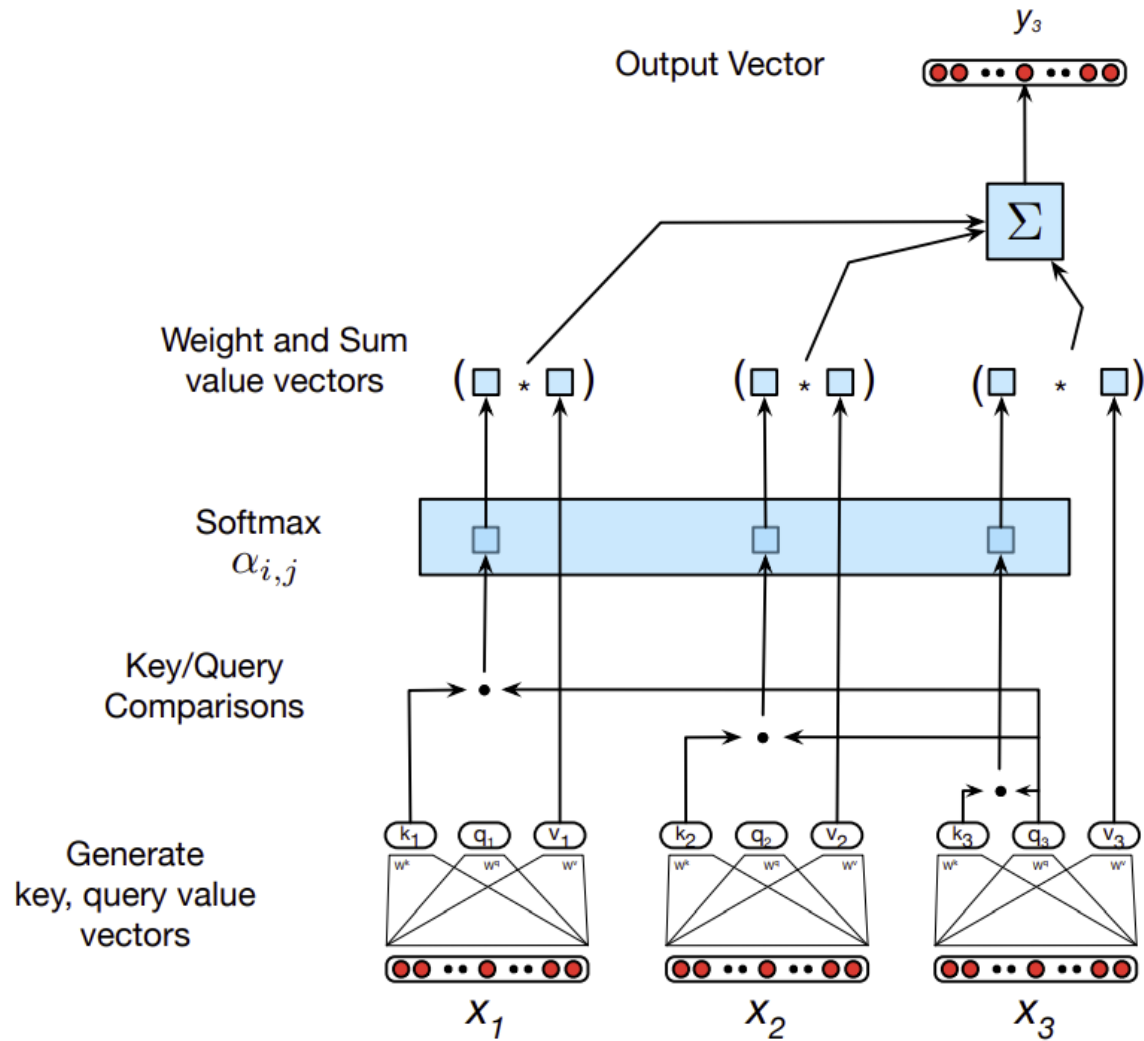
- $k_i = W^K x_i, \quad W^K \in \mathbb{R}^{d_k \times d_m}, \quad d_m=1024, d_k=64$

- $v_i = W^V x_i, \quad W^V \in \mathbb{R}^{d_v \times d_m}, \quad d_m=1024, d_v=64$

$$\textit{score}(x_i, x_j) = q_i \cdot k_j$$

$$\alpha_{ij} = \textit{softmax}(\textit{score}(x_i, x_j)) \quad \forall j \leq i$$

$$y_i = \sum_{j \leq i} \alpha_{ij} v_j$$



*Calculation of the value of the third element of a sequence using causal self-attention.*

# Score normalization

- Practical consideration about  $\alpha_{ij}$ 
  - Dot product can produce large values.
  - Exponentiating such large values can lead to numerical issues.
    - Effective loss of gradients during training.
- A typical approach is to divide the dot product by the square root of the dimensionality of the query and key vectors.

$$\text{score}(x_i, x_j) = \frac{q_i \cdot k_j}{\sqrt{d_k}}$$

# Single matrix packing...

- Produce matrices containing all the key, query and value vectors:

$$Q = W^Q X$$

$$K = W^K X$$

$$V = W^V X$$

- Can we then calculate all in a single step?

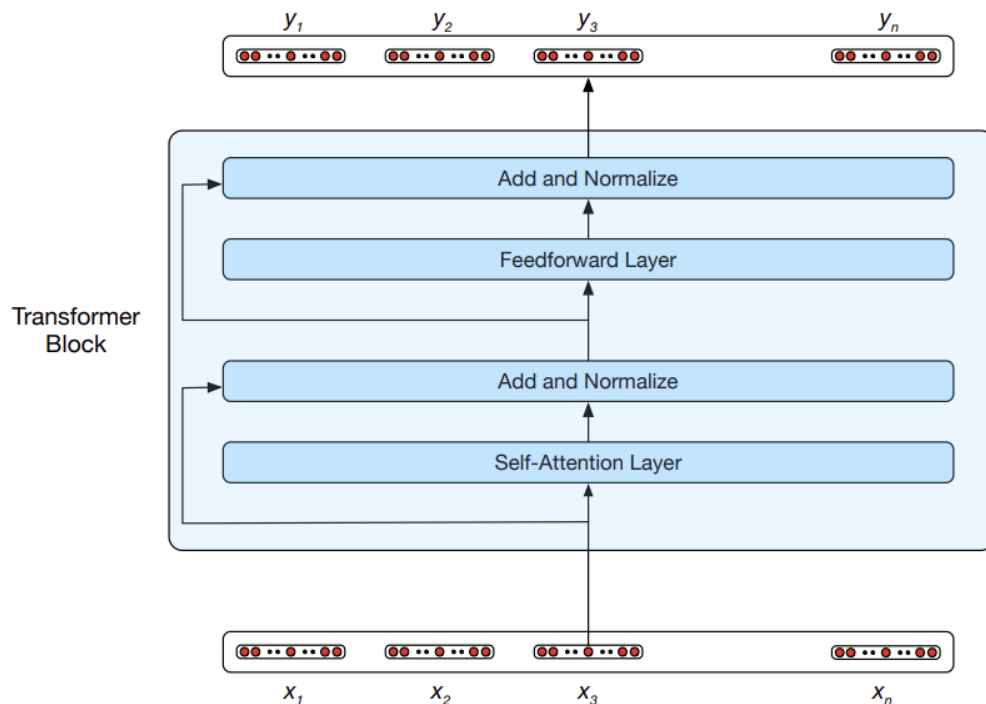
$$SelfAttention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

**Problem:**  $QK^T$  results in a score for each query value to every key value, including those that follow the query.

**Solution:** the elements in the upper-triangular portion of the comparisons matrix are zeroed out.

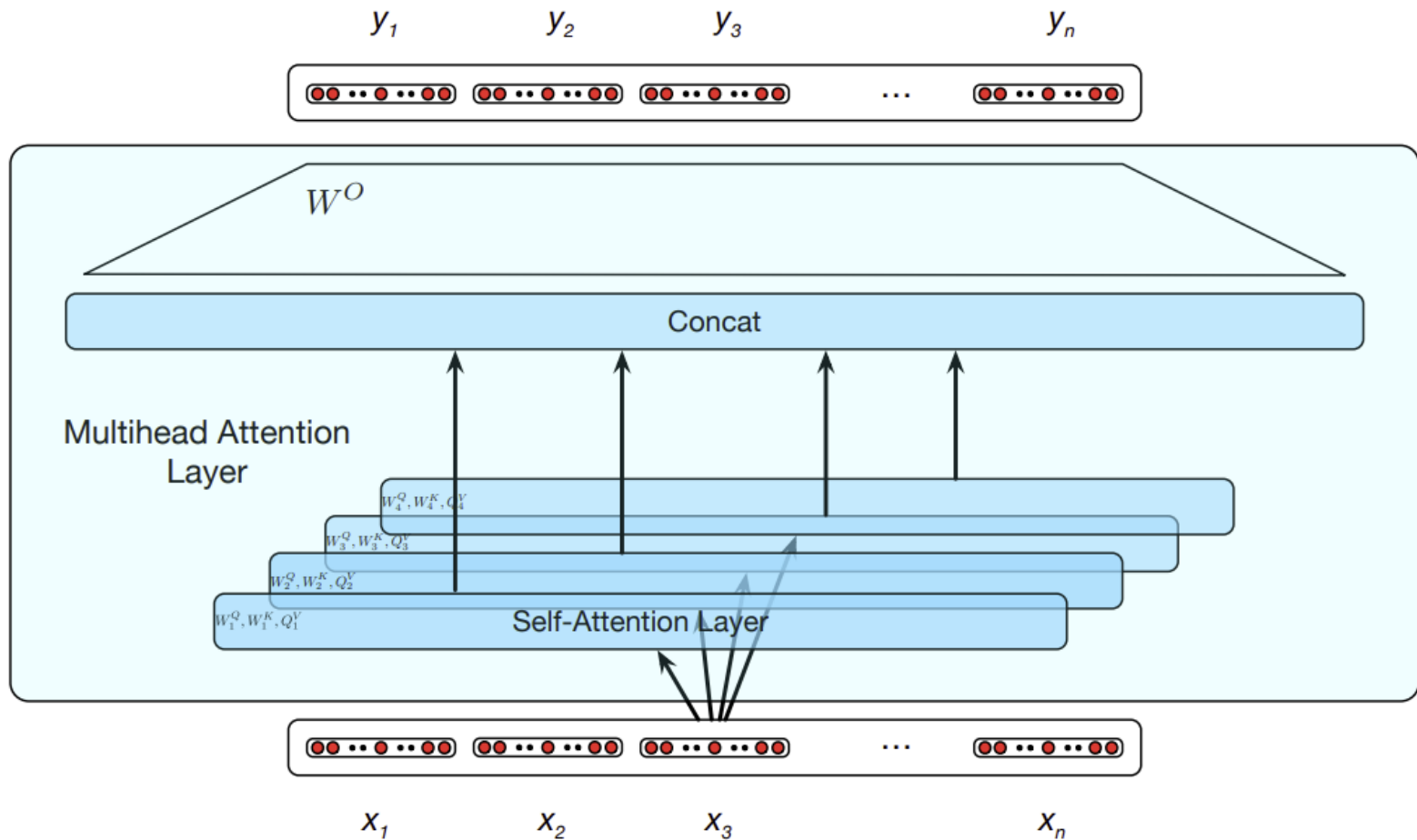
# Transformer block

- We just see the self attention layer.
- A transformer block also contains:
  - Normalization layers: Batch normalization accelerates training, in some cases by halving the epochs or better, and provides some regularization, reducing generalization error.
  - Feedforward layers.



# Multihead Attention Layers

- The different words in a sentence can relate to each other in many different ways simultaneously.
- Transformers address this issue with multihead self-attention layers.
- These are sets of self-attention layers, called heads, that reside in parallel layers at the same depth in a model, each with its own set of parameters.



$$MultiHeadAttn(Q, K, V) = W^O(head^1 \oplus head^2 \dots \oplus head^h)$$

Ojo Concatenación!

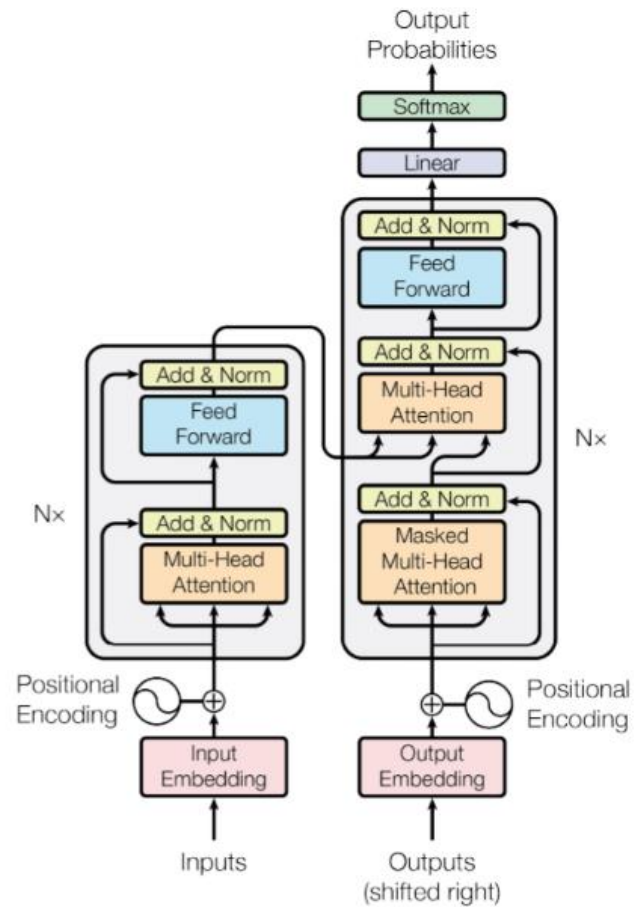
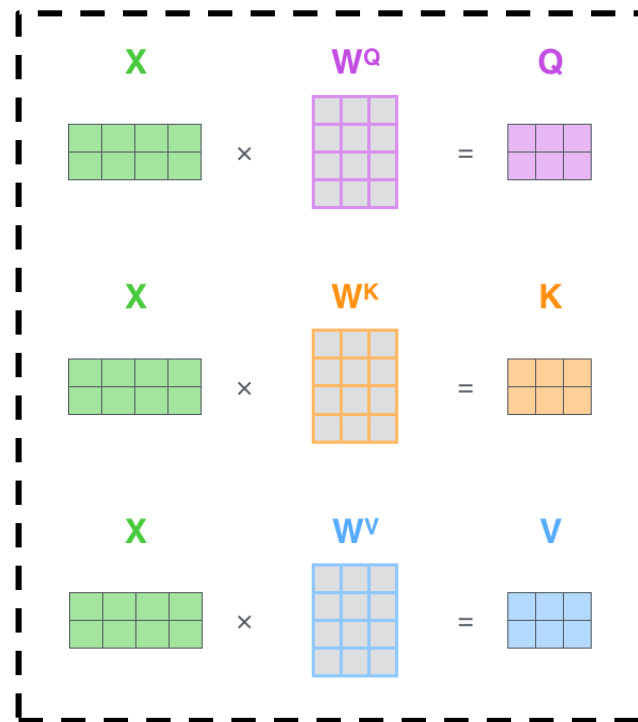
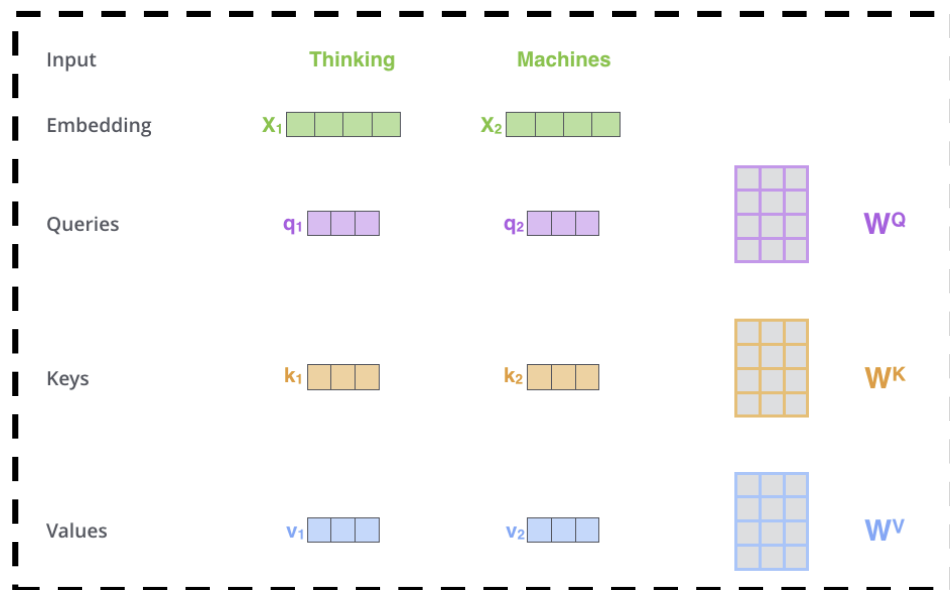


Figure 1: The Transformer - model architecture.



# Self attention review (I)



# Self attention review (II)

Input

Embedding

Queries

Keys

Values

Score

Divide by 8 (  $\sqrt{d_k}$  )

Softmax

Softmax

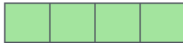
X  
Value

Sum

Thinking

Machines

$x_1$  

$x_2$  


$q_1$  

$q_2$  

$k_1$  

$k_2$  

$v_1$  

$v_2$  

$q_1 \cdot k_1 = 112$

$q_1 \cdot k_2 = 96$

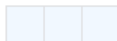
14

12

0.88

0.12

$v_1$  

$v_2$  

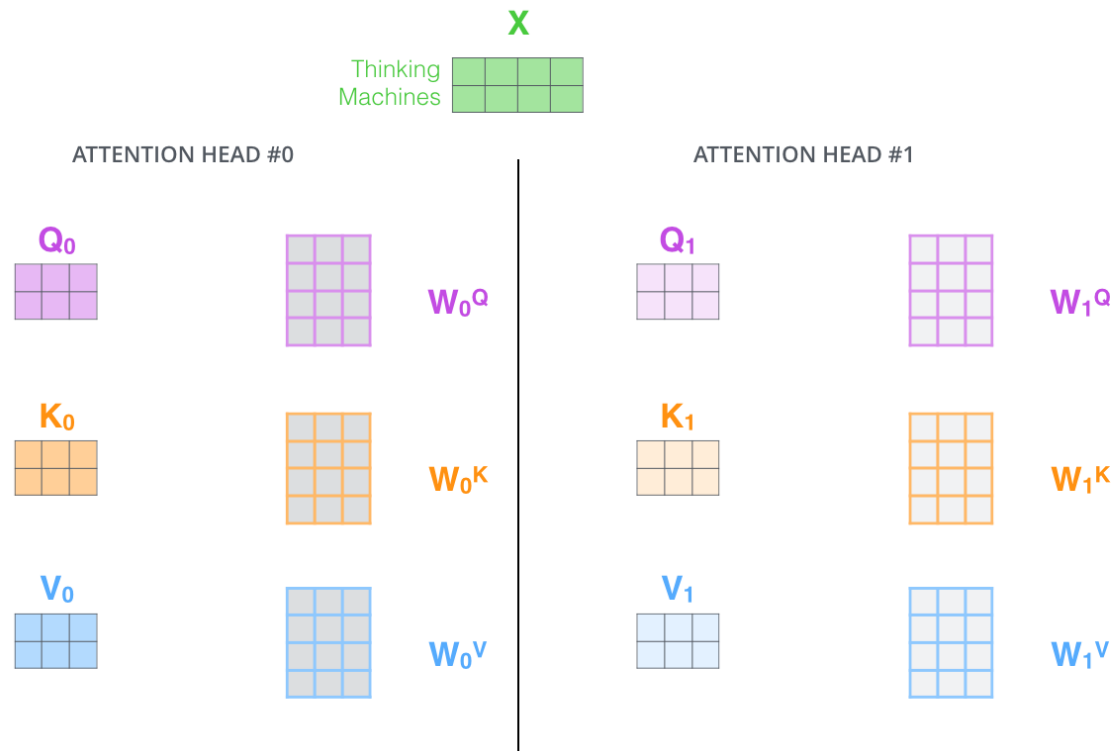
$z_1$  

$z_2$  

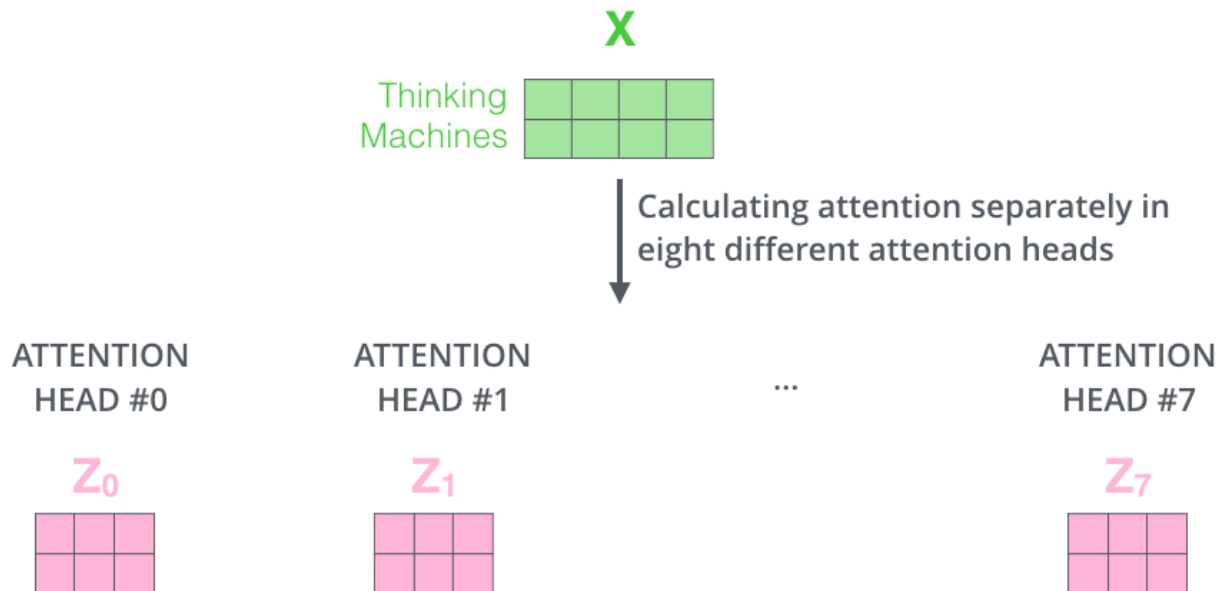
$$\text{softmax} \left( \frac{\begin{matrix} Q \\ \text{3x2 grid} \end{matrix} \times \begin{matrix} K^T \\ \text{2x3 grid} \end{matrix}}{\sqrt{d_k}} \right) \begin{matrix} V \\ \text{3x1 grid} \end{matrix}$$

$$= \begin{matrix} Z \\ \text{3x2 grid} \end{matrix}$$

# Self attention review (III): Multihead



# Self attention review (IV): Multihead



The diagram shows a 2x16 grid of cells. Above the grid, labels  $Z_0$  through  $Z_7$  are positioned over pairs of columns. Each pair of columns contains two adjacent cells, one light pink and one dark pink, representing a 2-bit vector element.

$\begin{array}{|c|c|c|c|} \hline & & & \\ \hline \end{array}$

X

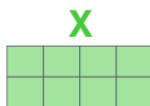
Wo

# Self-attention: Summary

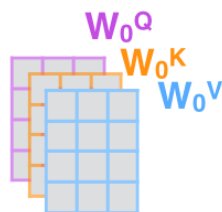
1) This is our  
input sentence\*

Thinking  
Machines

2) We embed  
each word\*



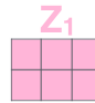
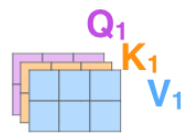
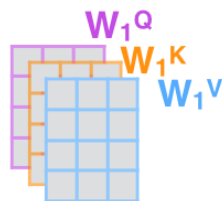
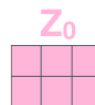
3) Split into 8 heads.  
We multiply  $X$  or  
 $R$  with weight matrices



4) Calculate attention  
using the resulting  
 $Q/K/V$  matrices



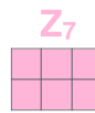
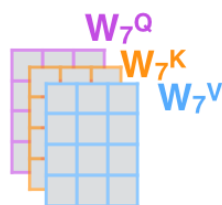
5) Concatenate the resulting  $Z$  matrices,  
then multiply with weight matrix  $W^O$   
to produce the output of the layer



...

...

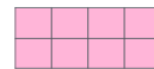
...



$W^O$

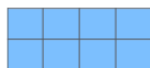


$Z$

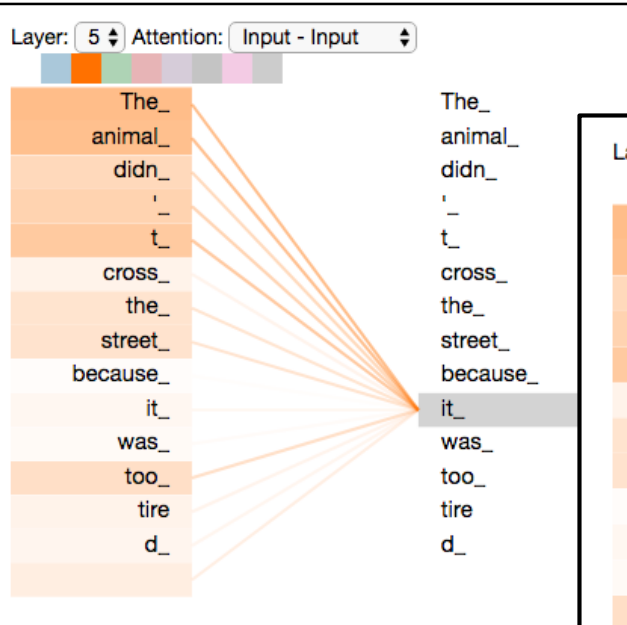


\* In all encoders other than #0,  
we don't need embedding.  
We start directly with the output  
of the encoder right below this one

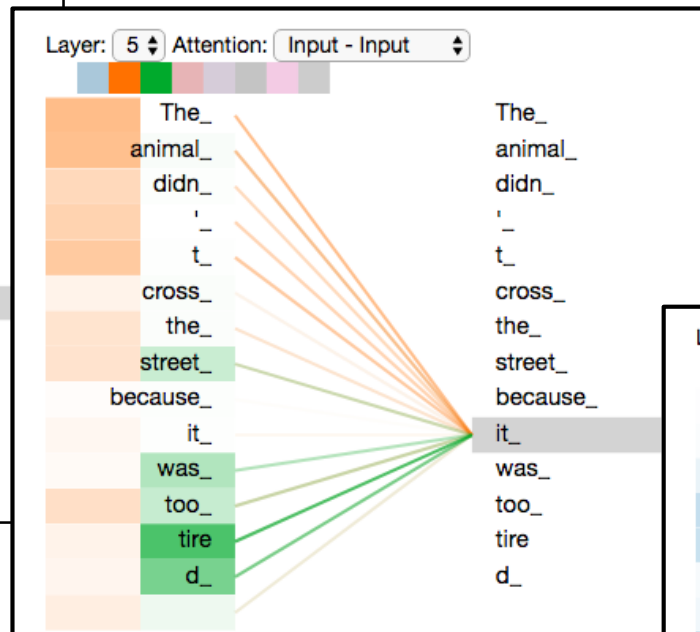
$R$



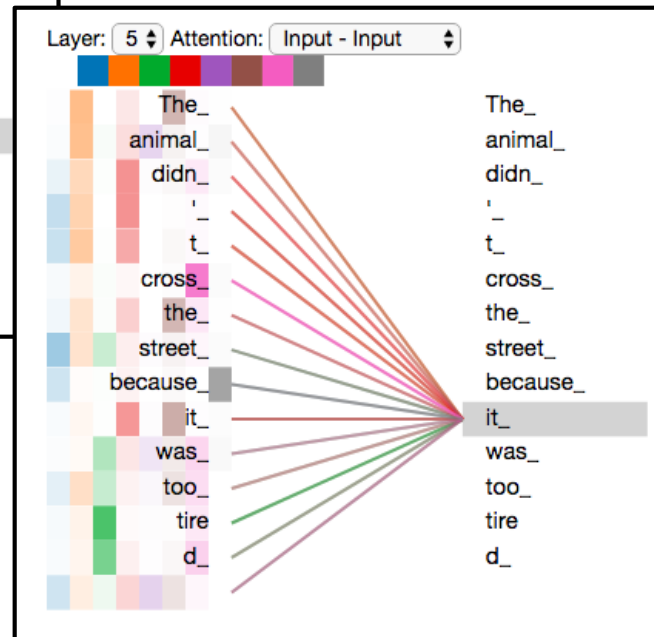
# One head



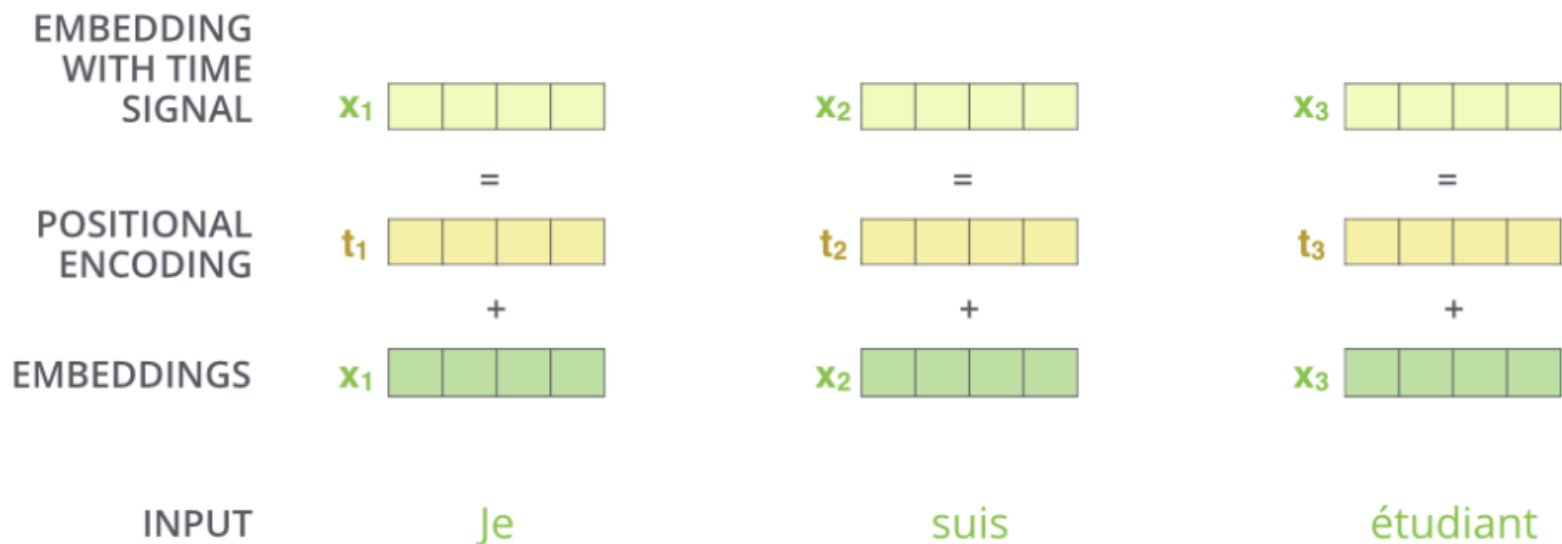
# Two heads



# All heads



# Representing The Order of The Sequence Using Positional Encoding



But the Transformer architecture ditched the recurrence mechanism in favor of multi-head self-attention mechanism. Avoiding the RNNs' method of recurrence will result in massive speed-up in the training time. As each word in a sentence simultaneously flows through the Transformer's encoder/decoder stack, The model itself doesn't have any sense of position/order for each word.



# Using Positional Encoding

- The first idea that might come to mind is to assign a number to each time-step within the  $[0, 1]$  range in which 0 means the first word and 1 is the last time-step.
  - **Problem:** you can't figure out how many words are in the sentence.
- Another idea is to assign a number to each time-step linearly. That is, the first word is given "1", the second word is given "2", and so on.
  - **Problem:** values could get quite large, but also our model can face sentences longer than the ones in training.

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

$$\omega_k = \frac{1}{10000^{2k/d}}$$

Geometric progression.

$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

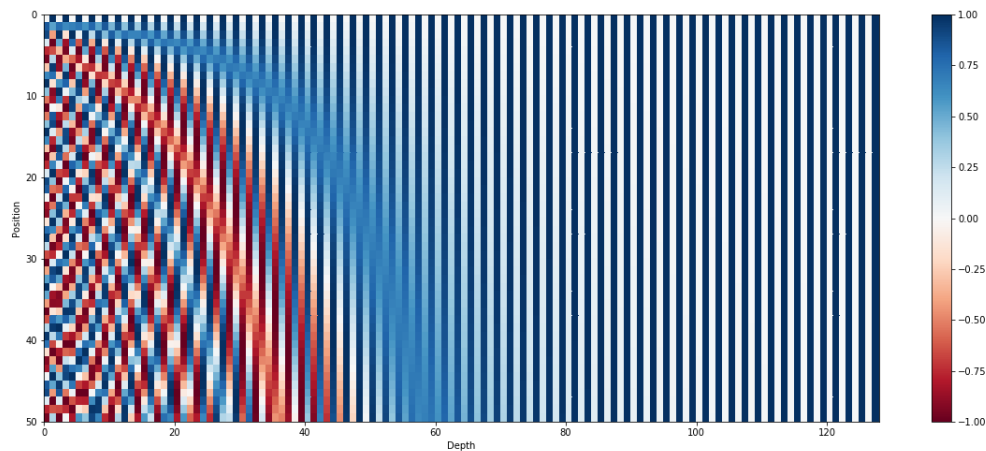
## Using Positional Encoding

- Let  $t$  be the desired position in an input sentence,  $p_t \in \mathbb{R}^d$  the positional encoding and  $d$  the encoding dimension. Then  $f: \mathbb{N} \rightarrow \mathbb{R}^d$  will be the function that produces the output vector  $p_t$ .

# Intuition

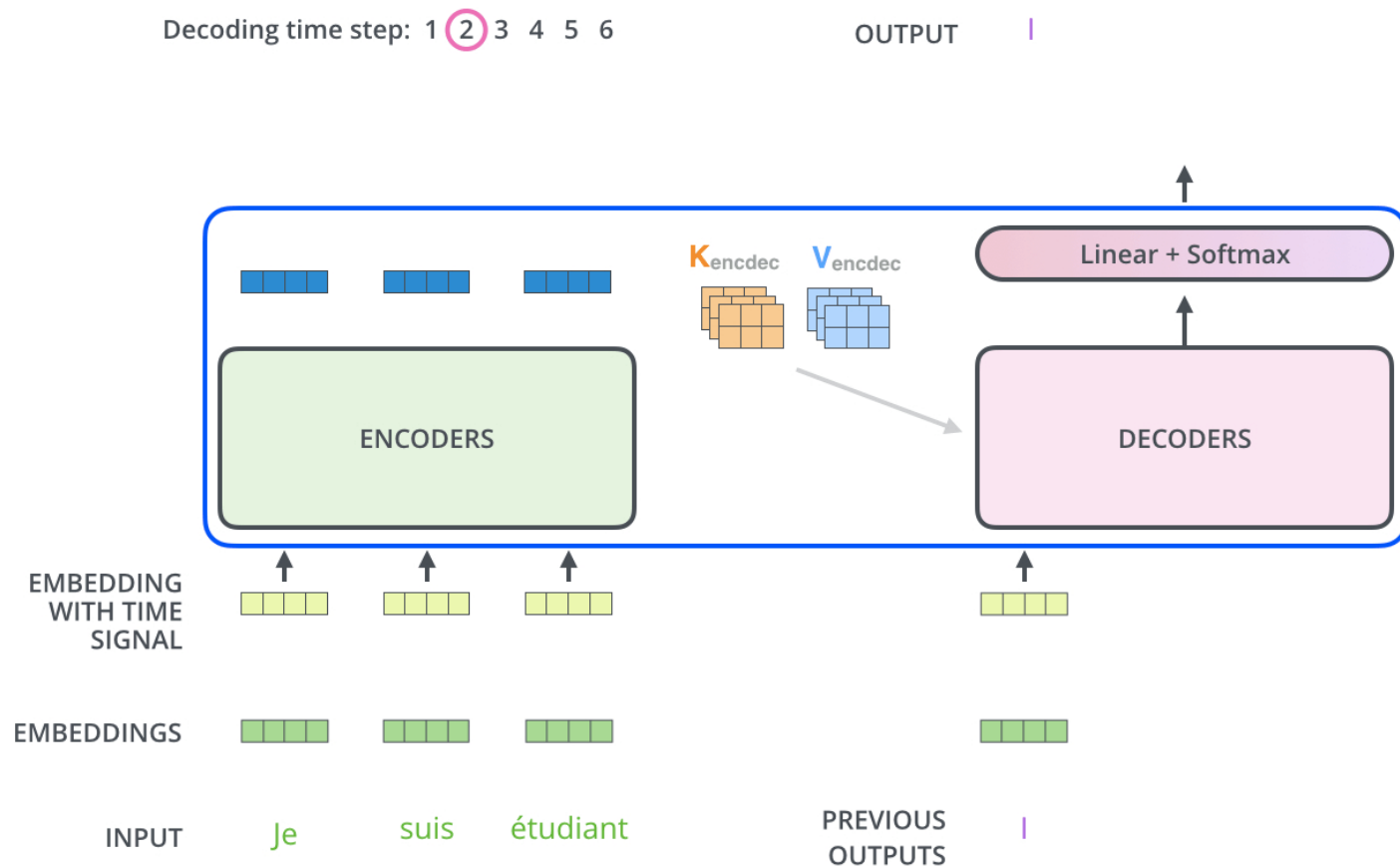
0 :	0	0	0	0	8 :	1	0	0	0
1 :	0	0	0	1	9 :	1	0	0	1
2 :	0	0	1	0	10 :	1	0	1	0
3 :	0	0	1	1	11 :	1	0	1	1
4 :	0	1	0	0	12 :	1	1	0	0
5 :	0	1	0	1	13 :	1	1	0	1
6 :	0	1	1	0	14 :	1	1	1	0
7 :	0	1	1	1	15 :	1	1	1	1

LSB bit is alternating on every number, the second-lowest bit is rotating on every two numbers, and so on.

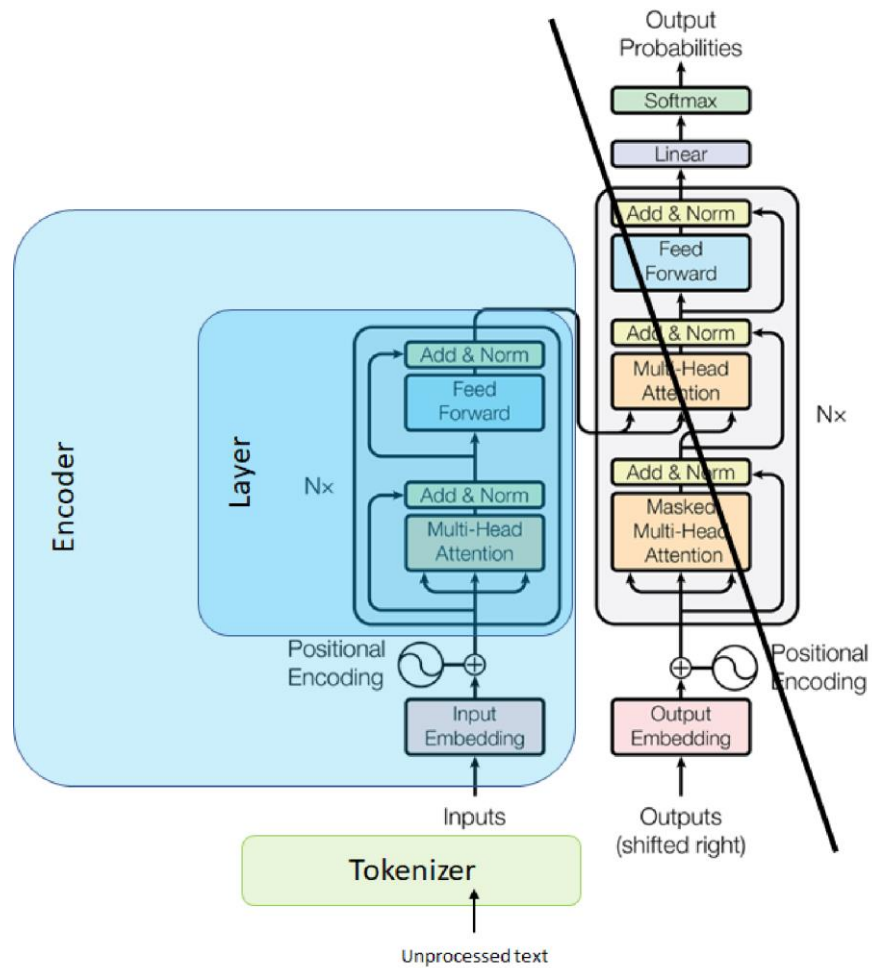


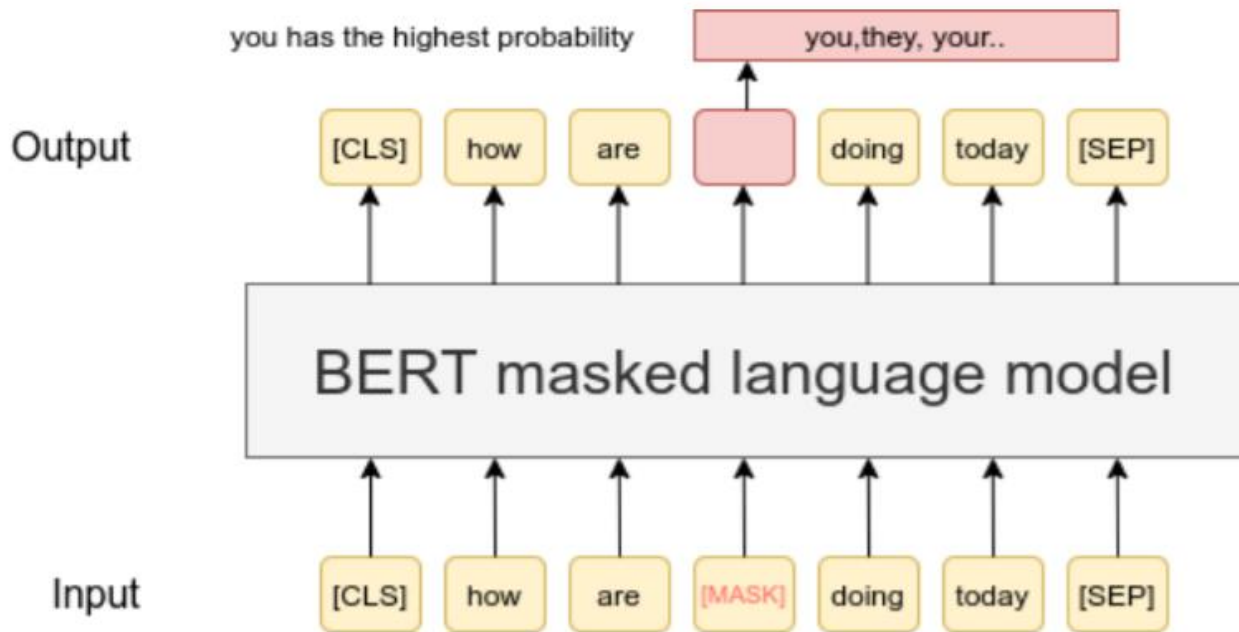
The 128-dimensional positional encoding for a sentence with the maximum length of 50.

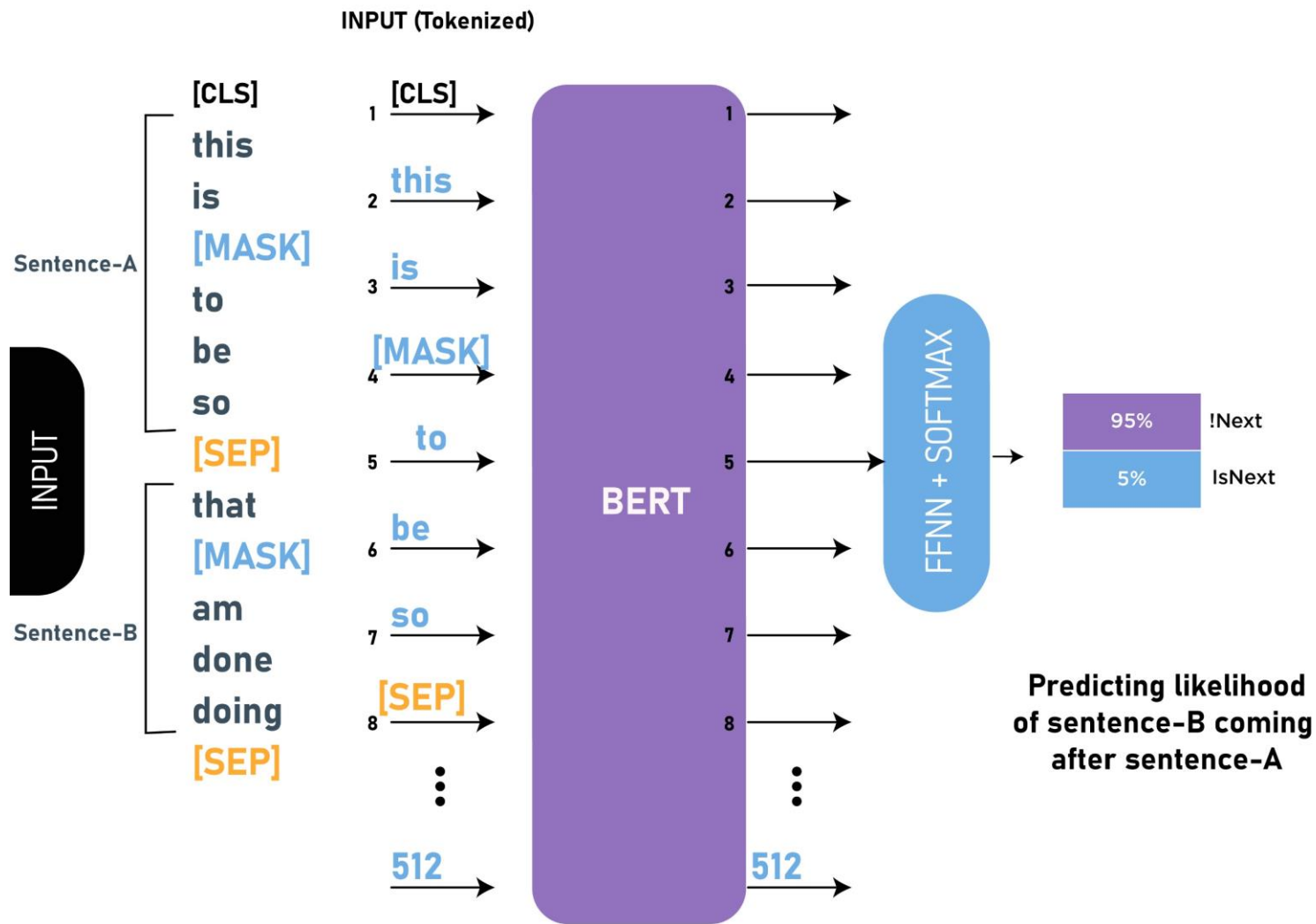
# And the decoder?



In the decoder, the self-attention layer is only allowed to attend to earlier positions in the output sequence.







# Original paper

- <https://ai.googleblog.com/2017/06/accelerating-deep-learning-research.html>
  - [Tensor2Tensor](#)
- <https://arxiv.org/pdf/1706.03762.pdf>



Gracias por la atención

¿Tiene alguna pregunta?

