

# Attention Is All You Need

By: Subash Sah

# Why Transformers?

- RNNs/LSTMs are slow and sequential
- Hard to capture long-range dependencies
- Transformers allow full parallel processing
- Much better performance in translation and NLP

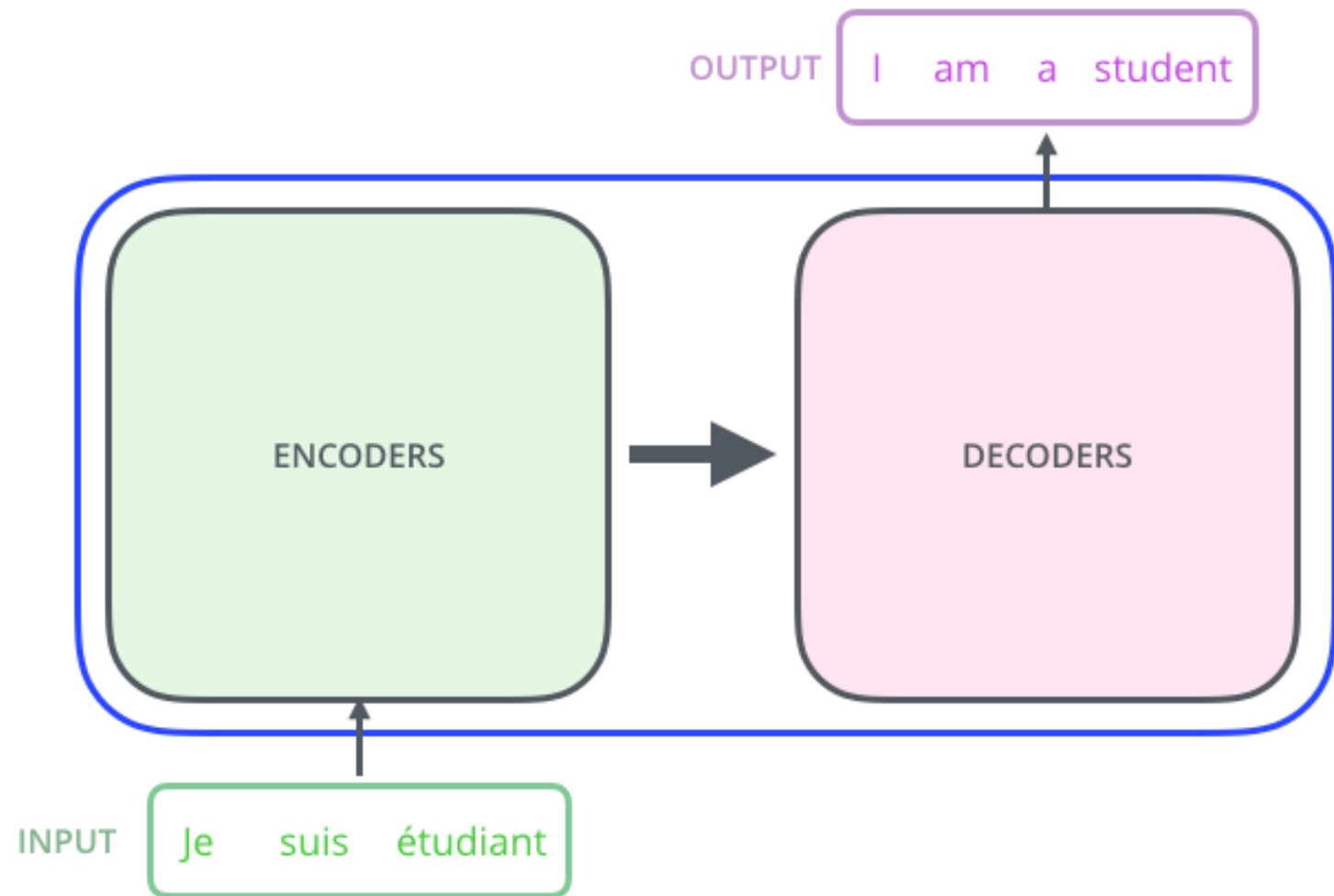
# What Makes Transformers Different

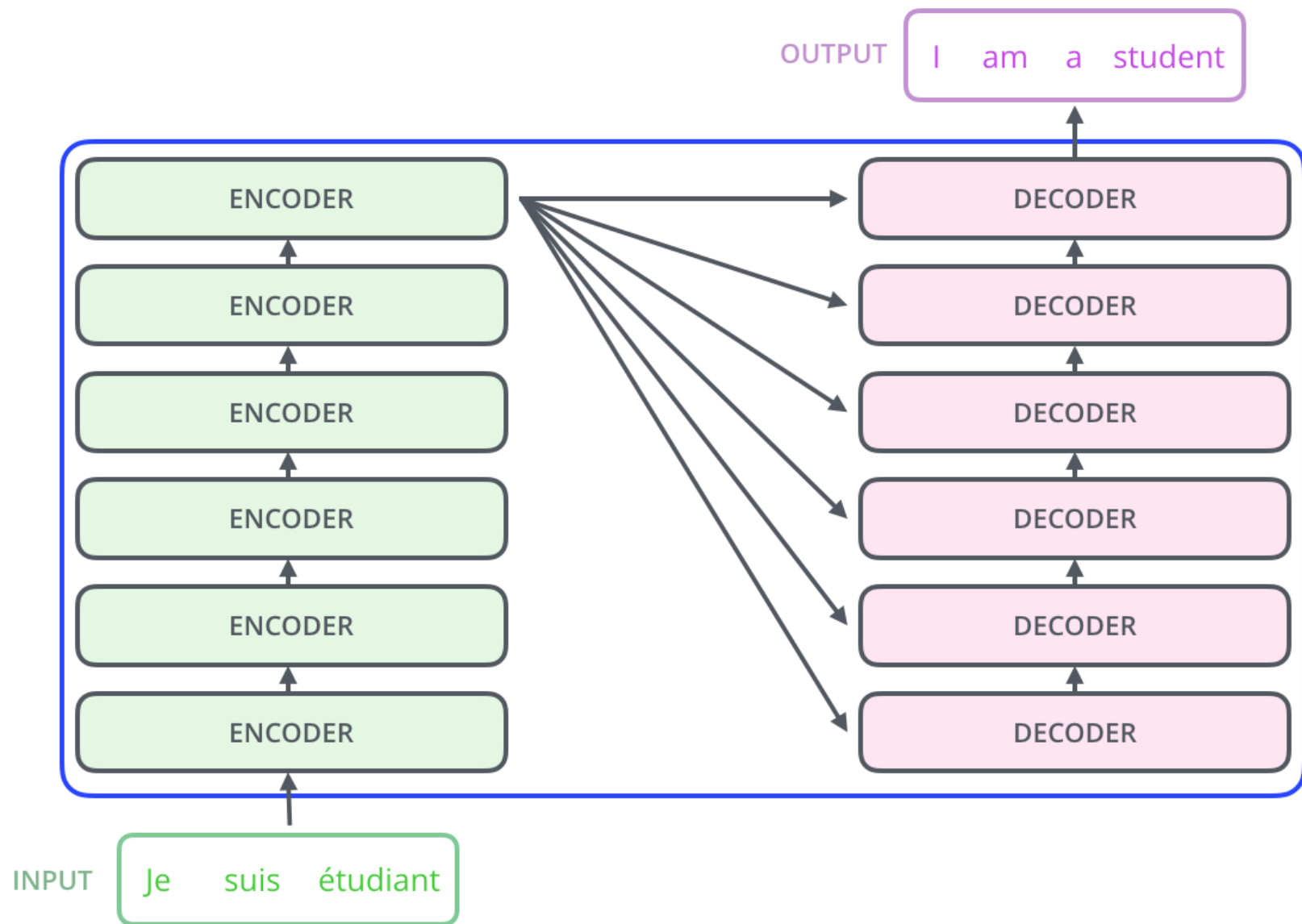
- No recurrence
- Use self-attention instead
- Handle whole sequences at once
- Scale well with data and model size
- Foundation for modern models like BERT and GPT

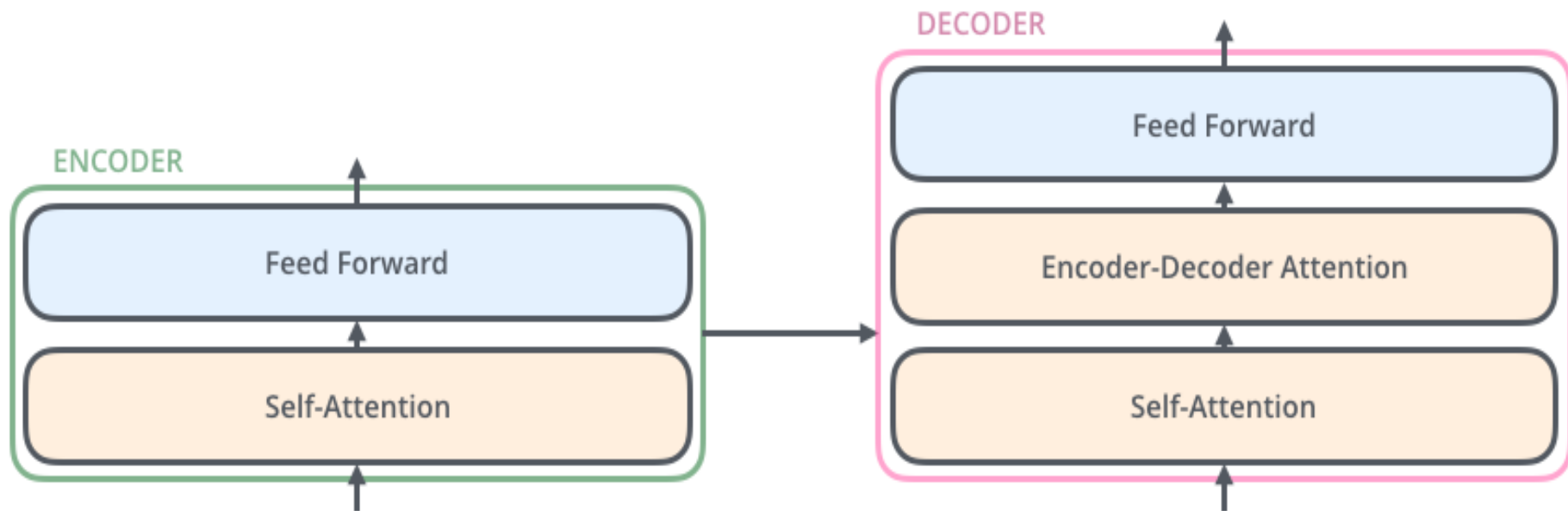
# Architecture of Transformer

# A High-Level Look



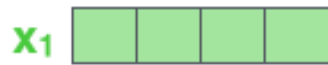




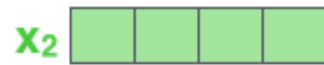




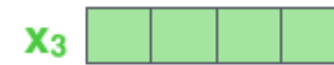
# Bringing Tensors/Vectors into the Picture



Je

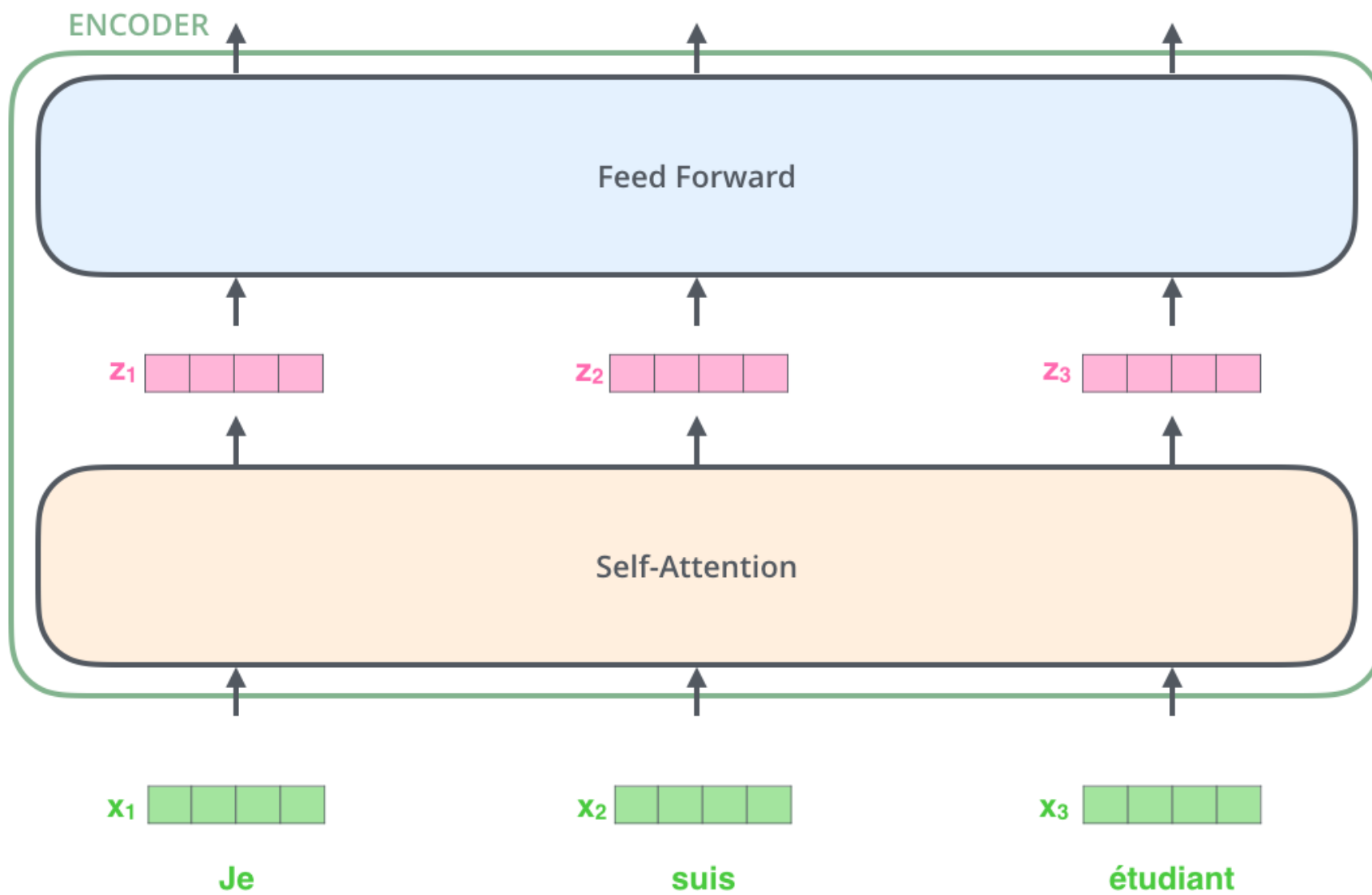


suis



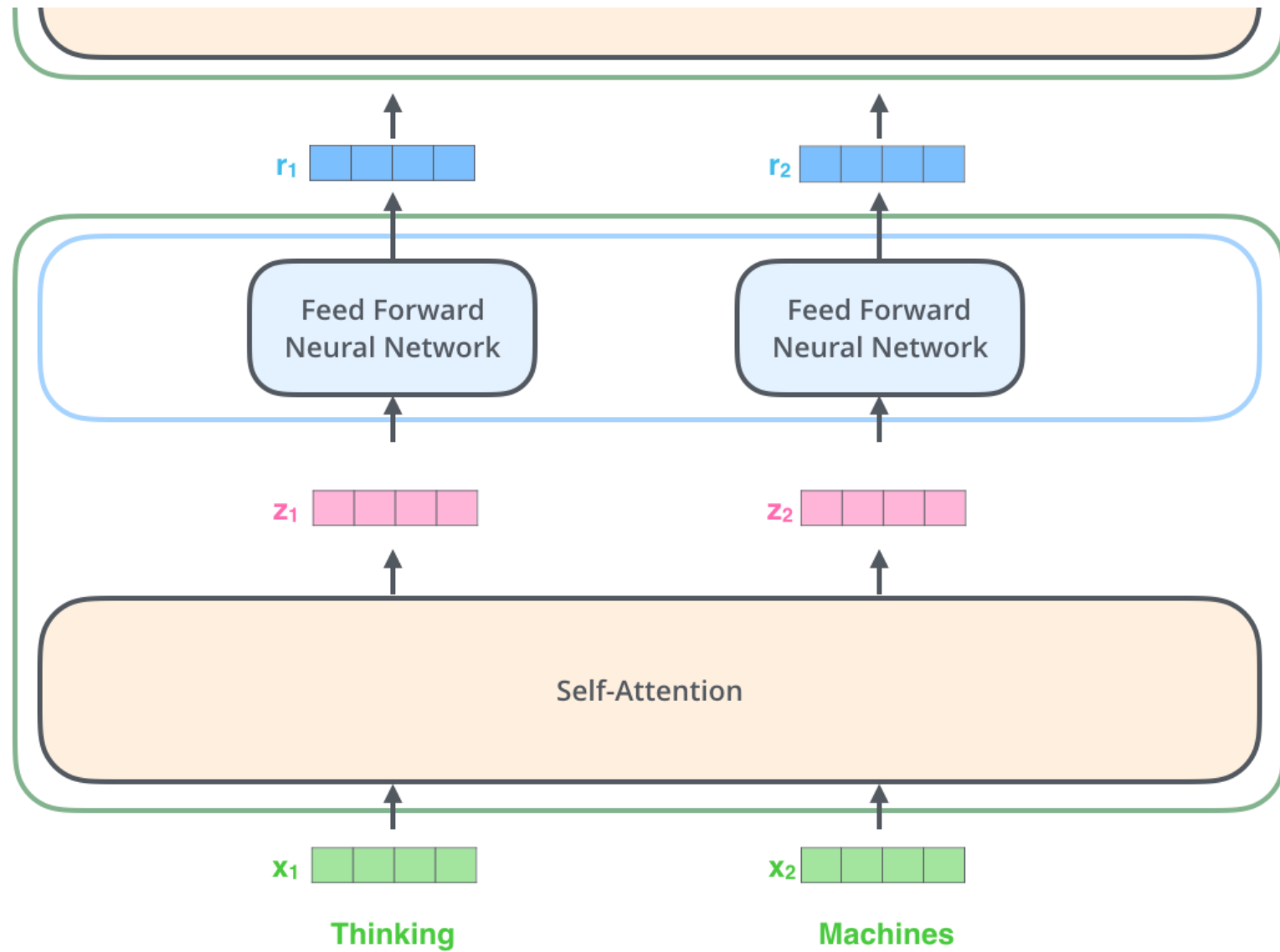
étudiant

Vector 512 dimensions



ENCODER #2

ENCODER #1



# Self-Attention at High Level

- Self-attention is a mechanism where each word in a sentence looks at every other word (including itself) to decide which ones are important for understanding its meaning.

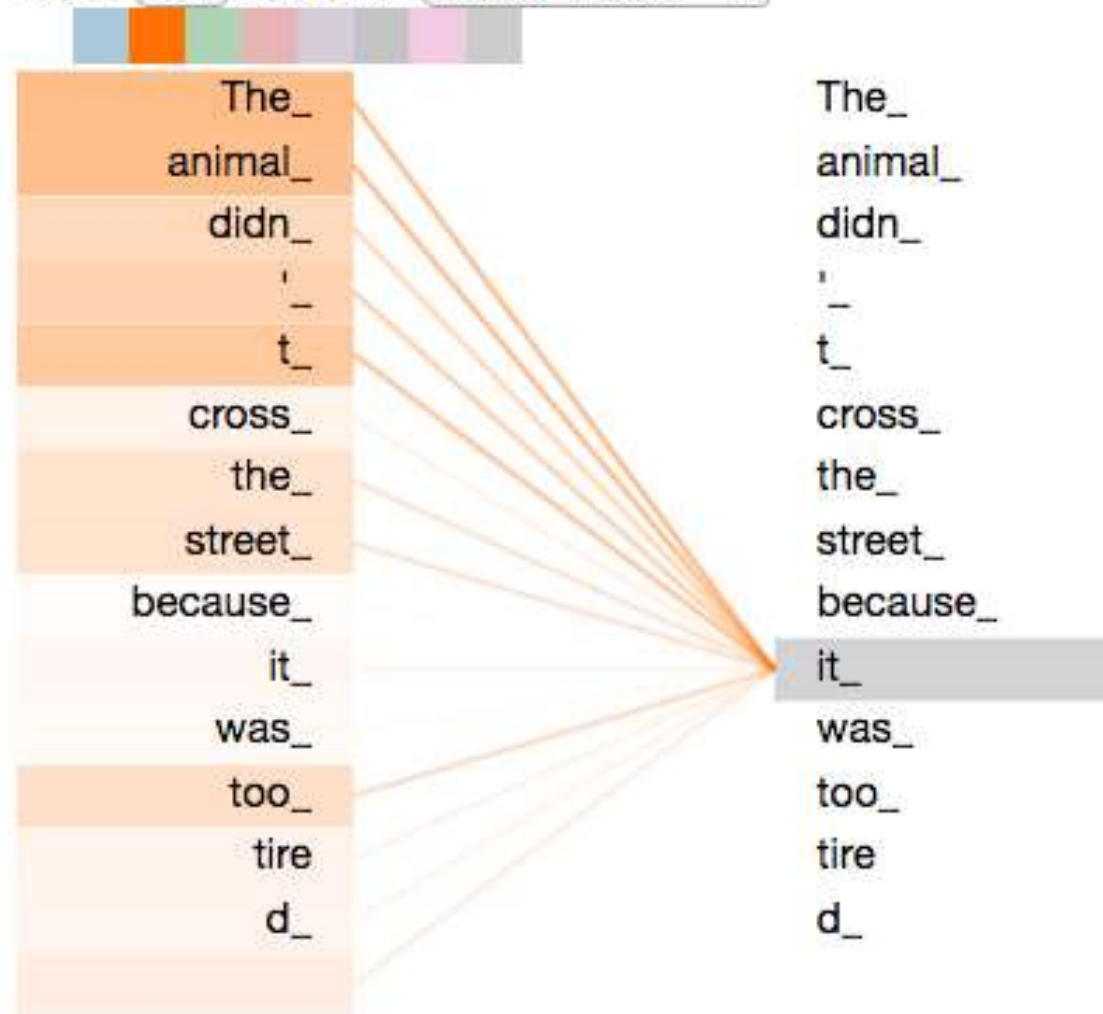
## Why we need it?

- When reading the sentence:  
    “The animal didn’t cross the street because it was too tired.”  
    To understand “it”, we must look at “animal”, not “street”.
- Self-attention lets each word pull in information from other relevant words.

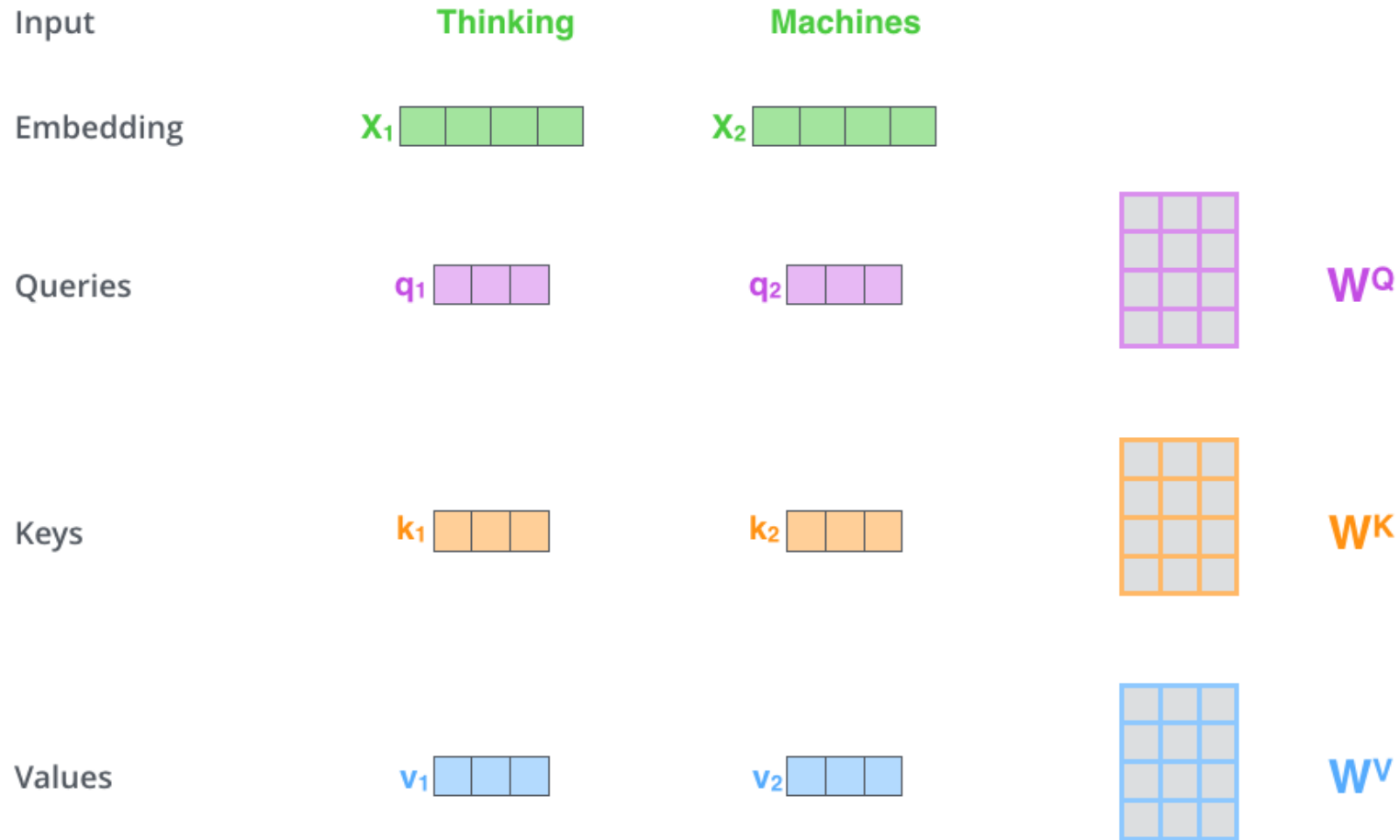
For each word, self-attention asks:

- “Which words in this sentence should I pay attention to?” and
- “How much should I pay attention to each one?”

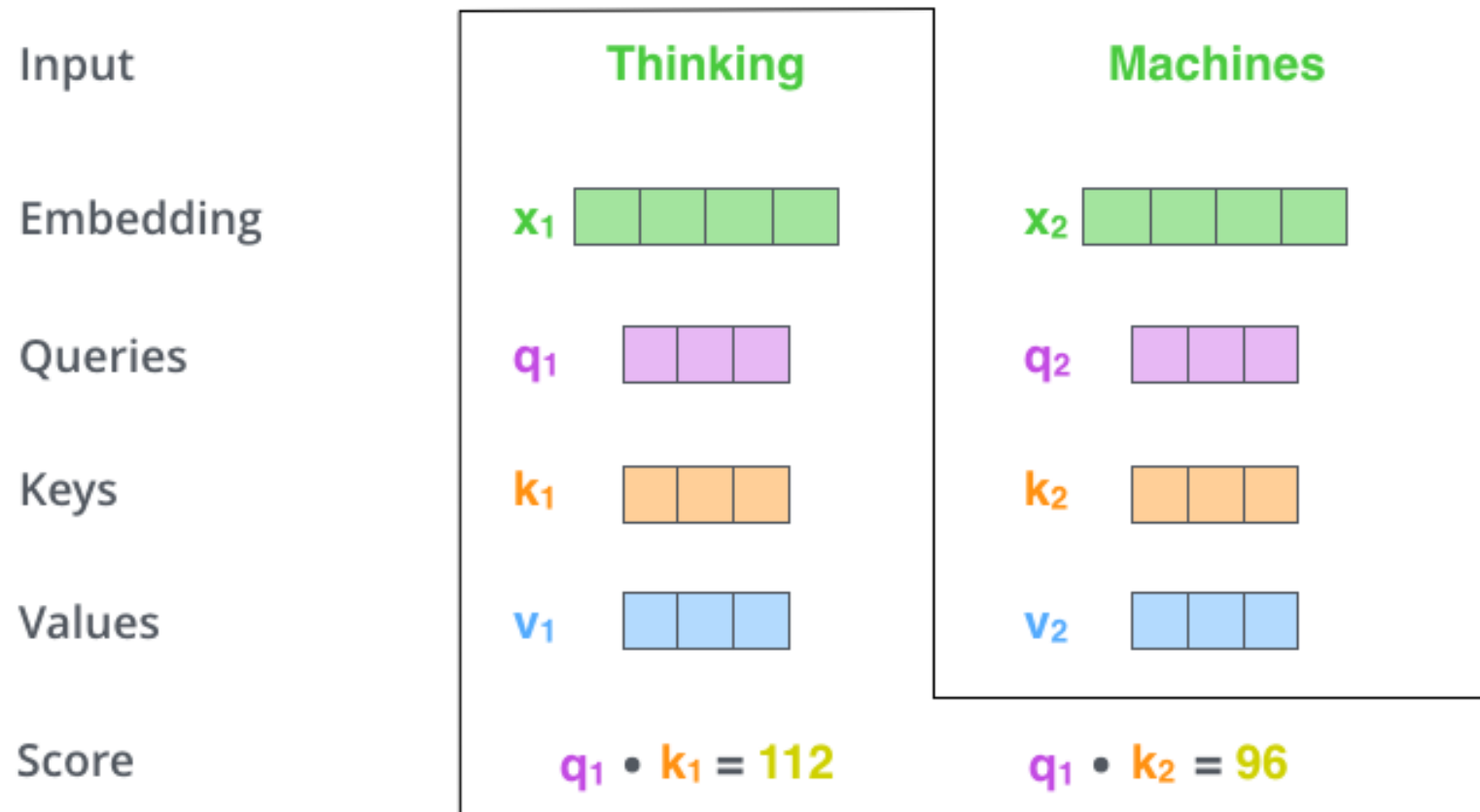
Layer: 5 Attention: Input - Input



# Self-Attention in Detail



The first step in calculating self-attention is to create three vectors from each of the encoder's input vectors.



The second step in calculating self-attention is to calculate a score.



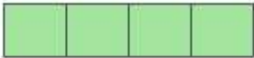
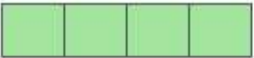
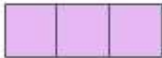
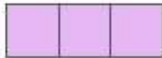


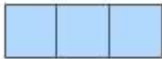
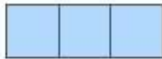
# Interpretation of Score value

- High score  $\rightarrow$  The second word is very relevant to the first
- Low score  $\rightarrow$  The word is not important
- Near zero  $\rightarrow$  The words are unrelated
- Negative (sometimes)  $\rightarrow$  They should be ignored

So while computing attention for the query of “it”, we expect:

- $q_{it} \cdot k_{animal} \rightarrow$  highest score
- $q_{it} \cdot k_{street} \rightarrow$  lower
- $q_{it} \cdot k_{it} \rightarrow$  moderate
- $q_{it} \cdot k_{otherwords} \rightarrow$  lower

In one sentence: The attention score ( $q \cdot k$ ) is a learned measure of how much one word should pay attention to another.

Input	Thinking	Machines
Embedding	$x_1$ 	$x_2$ 
Queries	$q_1$ 	$q_2$ 
Keys	$k_1$ 	$k_2$ 
Values	$v_1$ 	$v_2$ 
Score	$q_1 \cdot k_1 = 112$	$q_1 \cdot k_2 = 96$
Divide by 8 ( $\sqrt{d_k}$ )	14	12
Softmax	0.88	0.12

Step 3 is to divide the score by 8( square root of dimension of key vectors) and step 4 is to use softmax operation.

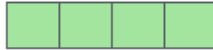
- Why divide by 8?
  - Without dividing, the dot product becomes too large and can cause exploding gradient problems( training becomes unstable)
- Why softmax?
  - Softmax takes raw scores and turns them into a probability distribution.
  - That means:
    - all values become positive
    - they add up to 1
    - bigger scores become bigger weights
    - smaller scores become smaller weights
  - It answers the question: “How much attention should I give to each word?”

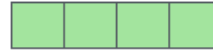
Input

Thinking

Machines

Embedding

$x_1$  

$x_2$  

Queries

$q_1$  

$q_2$  

Keys

$k_1$  

$k_2$  

Values

$v_1$  

$v_2$  

Score

$q_1 \cdot k_1 = 112$

$q_1 \cdot k_2 = 96$

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

14

12

Softmax

0.88

0.12

Softmax

X

Value

$v_1$  

$v_2$  

Sum

$z_1$  

$z_2$  

Fifth step is to multiply each value vector by softmax score. And 6<sup>th</sup> step is to sum up the weighted value vectors.

- Why multiply?

Because the Value vectors represent the actual content of each word.

- High score  $\rightarrow$  keep more of that word's content
- Low score  $\rightarrow$  keep less of it
- Zero score  $\rightarrow$  ignore that word completely

- And then finally add all the value vectors to give Z which is the output of self attention.

- The Z vector contains:
  - The word itself
  - Plus relevant information from other words
  - Weighted by how important they are

# Matrix Calculation of Self-Attention

$$\begin{matrix} \text{X} \\ \begin{array}{|c|c|c|c|} \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{W}^{\text{Q}} \\ \begin{array}{|c|c|c|c|} \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \end{array} \end{matrix} = \begin{matrix} \text{Q} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$

$$\begin{matrix} \text{X} \\ \begin{array}{|c|c|c|c|} \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{W}^{\text{K}} \\ \begin{array}{|c|c|c|c|} \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \end{array} \end{matrix} = \begin{matrix} \text{K} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$

$$\begin{matrix} \text{X} \\ \begin{array}{|c|c|c|c|} \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{W}^{\text{V}} \\ \begin{array}{|c|c|c|c|} \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \end{array} \end{matrix} = \begin{matrix} \text{V} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$

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

The self-attention calculation in matrix form.

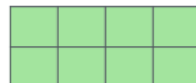
# Multi-Head Attention

- Why multi-head attention?
  - One head cannot capture all the patterns at once.
  - Multiple heads allow the transformer to understand language from many angles at the same time.
- In the example sentence:
  - “The animal didn’t cross the street because it was tired.”, different kinds of relationships matter.
  - Head 1: “it”  $\leftrightarrow$  “animal” (pronoun resolution)
  - Head 2: grammatical structure
  - Head 3: adjective  $\rightarrow$  noun and so on.



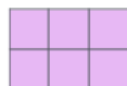
X

Thinking  
Machines



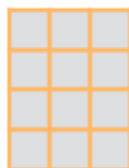
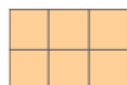
ATTENTION HEAD #0

$Q_0$



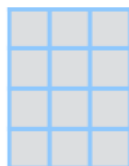
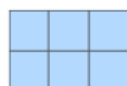
$W_0^Q$

$K_0$



$W_0^K$

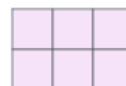
$V_0$



$W_0^V$

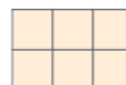
ATTENTION HEAD #1

$Q_1$



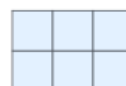
$W_1^Q$

$K_1$



$W_1^K$

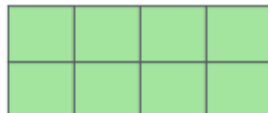
$V_1$



$W_1^V$

X

Thinking  
Machines

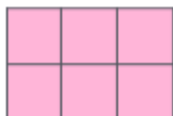


Calculating attention separately in  
eight different attention heads



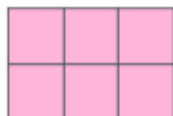
ATTENTION  
HEAD #0

$Z_0$



ATTENTION  
HEAD #1

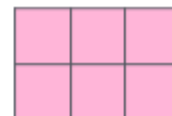
$Z_1$



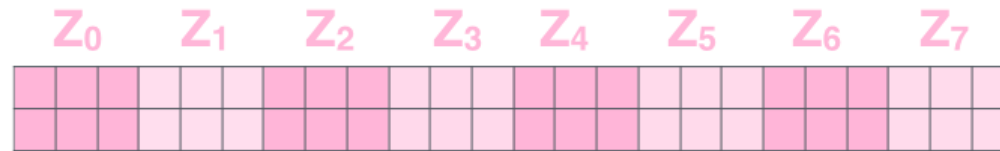
...

ATTENTION  
HEAD #7

$Z_7$



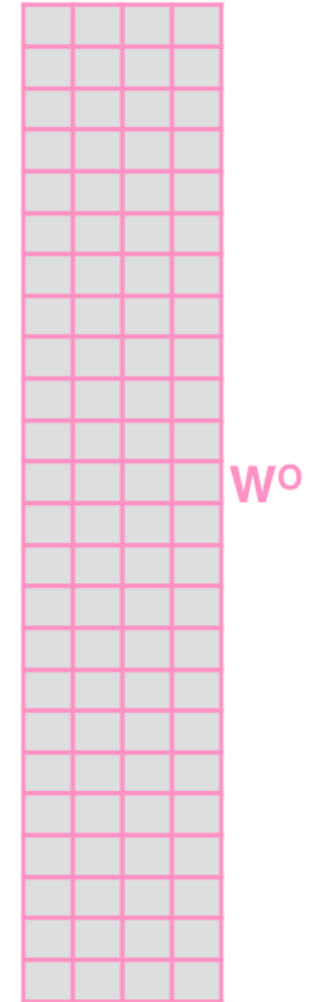
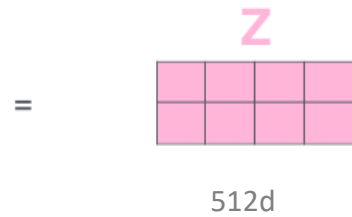
1) Concatenate all the attention heads



2) Multiply with a weight matrix  $W^O$  that was trained jointly with the model

X

3) The result would be the  $Z$  matrix that captures information from all the attention heads. We can send this forward to the FFNN



1) This is our input sentence\*

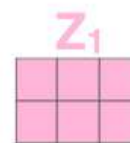
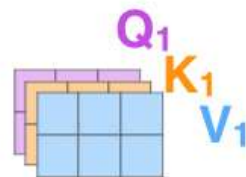
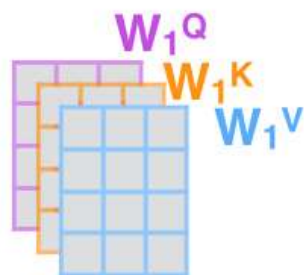
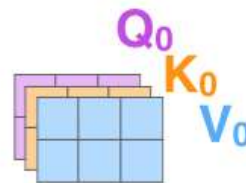
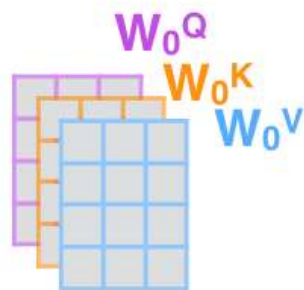
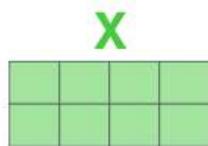
2) We embed each word\*

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

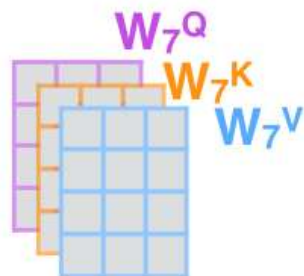
Thinking  
Machines



...

...

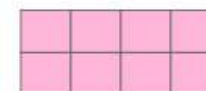
...



$W^O$



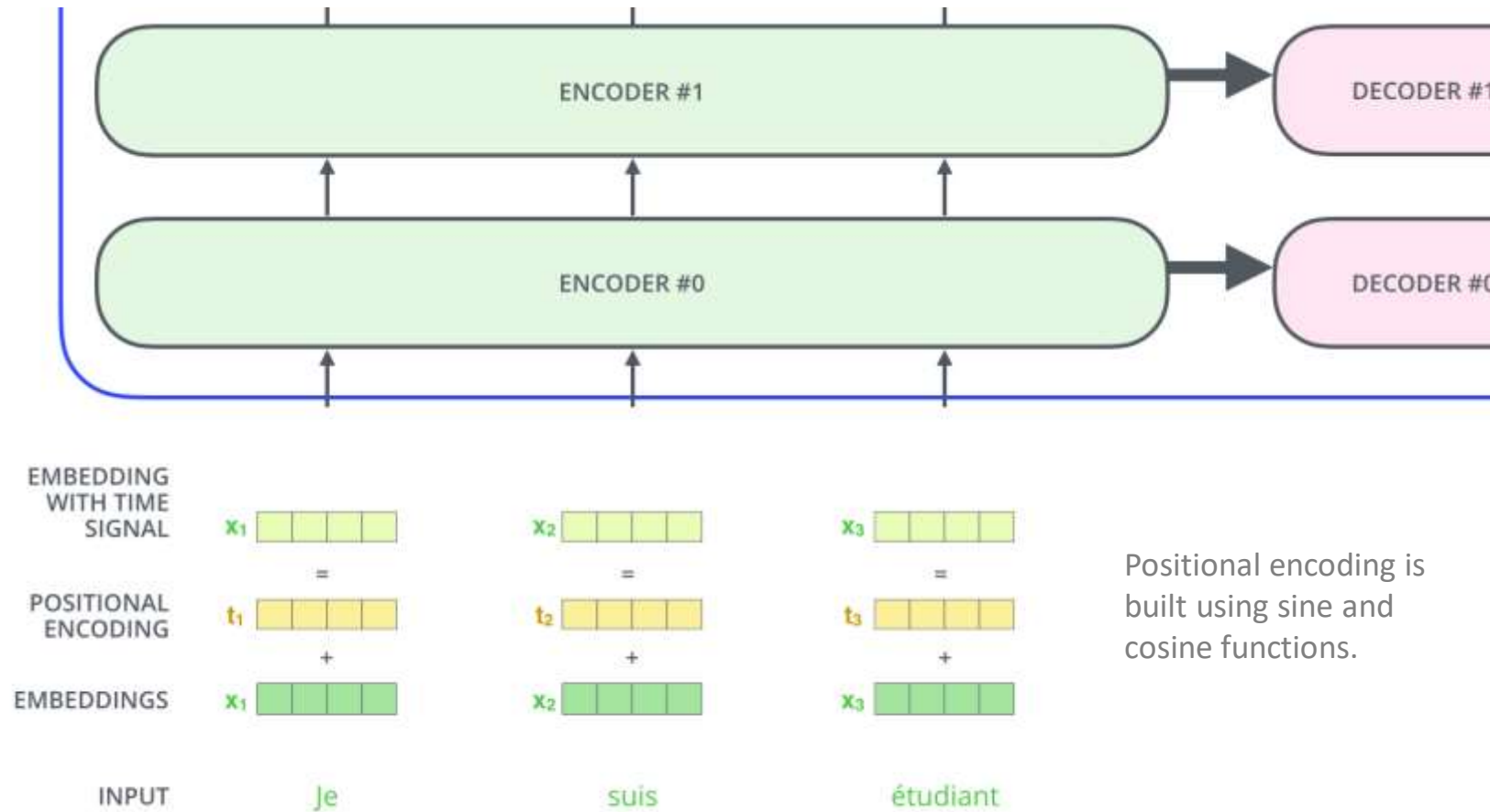
$Z$



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

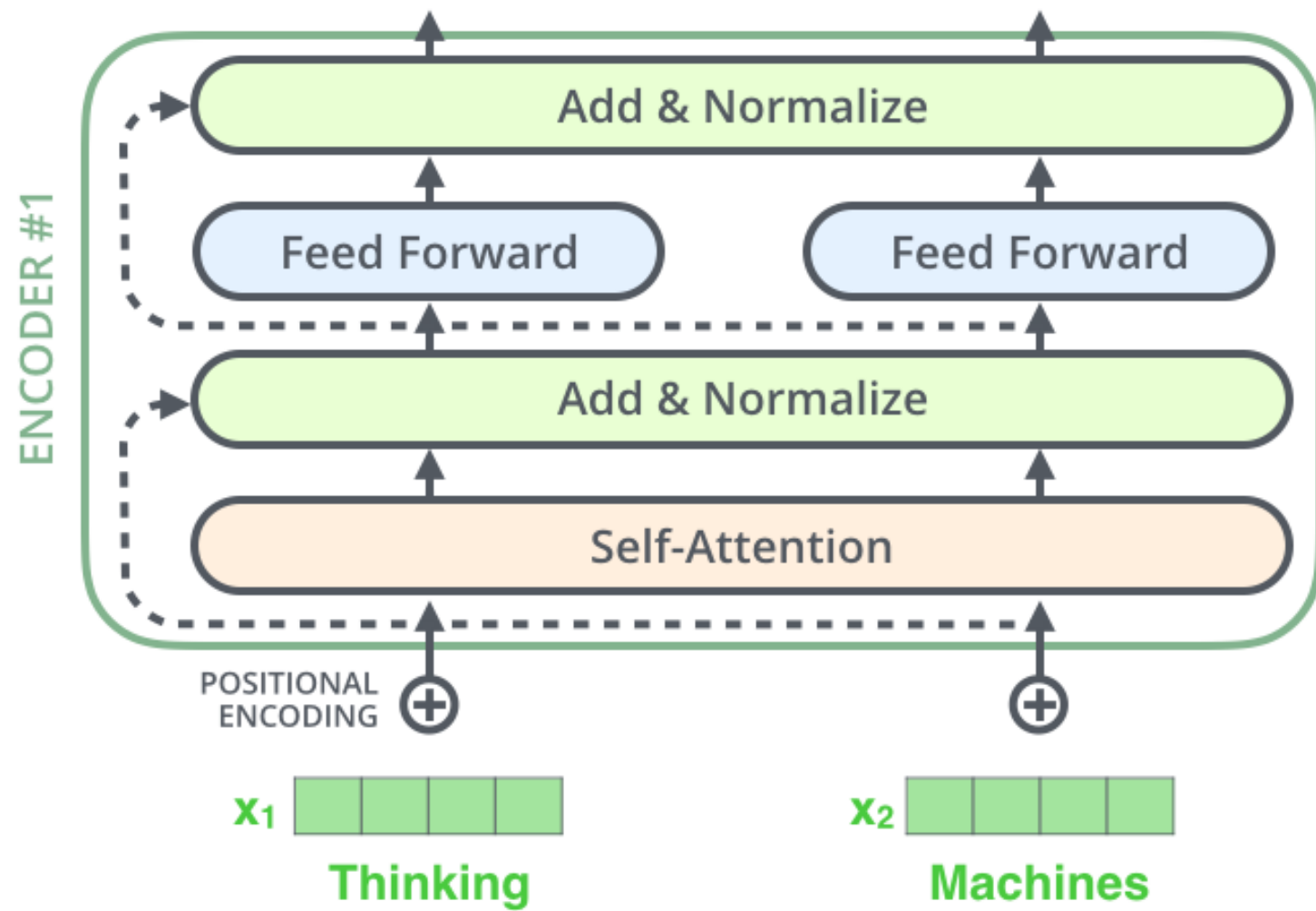


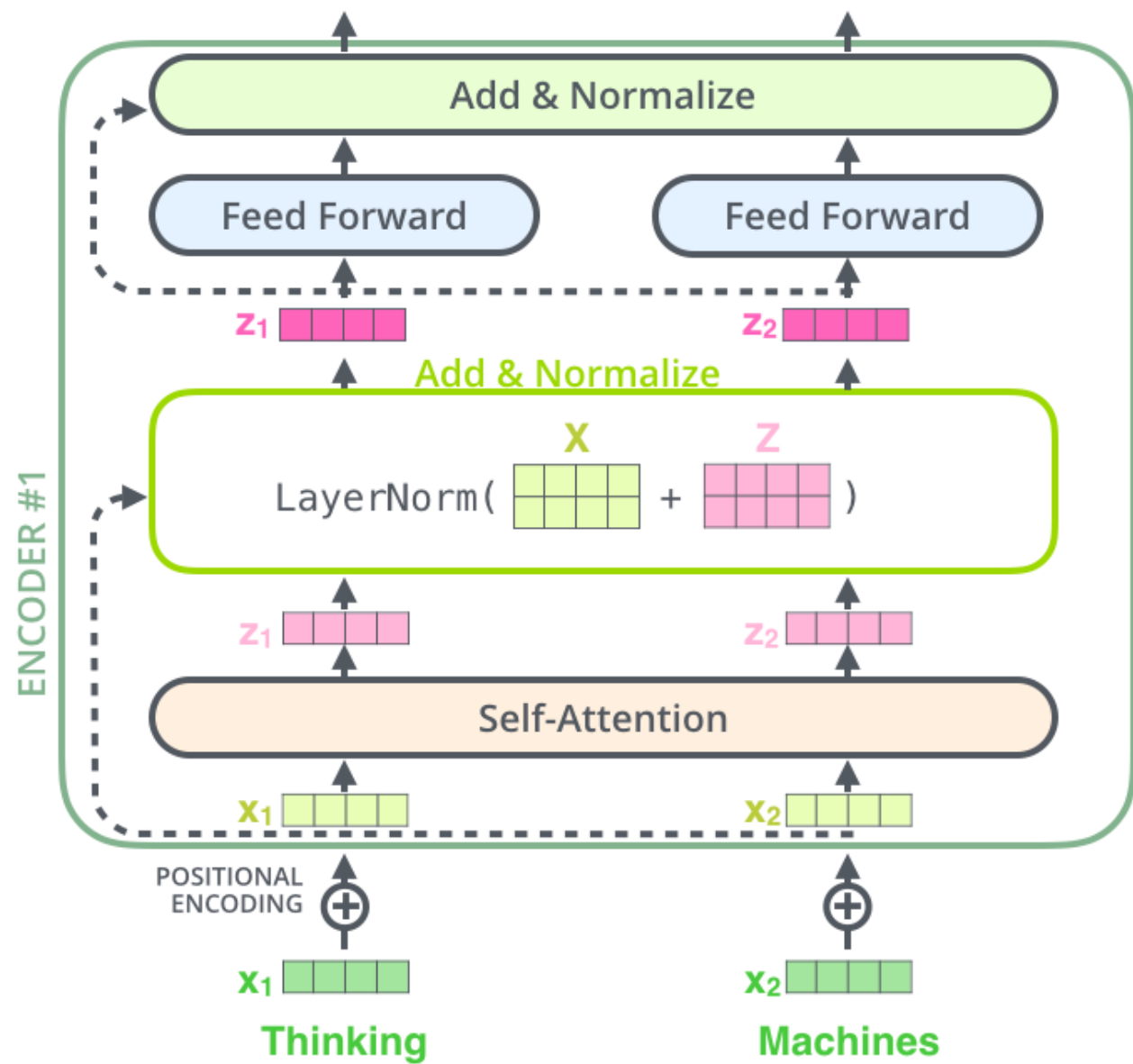
# Representing The Order of The Sequence Using Positional Encoding



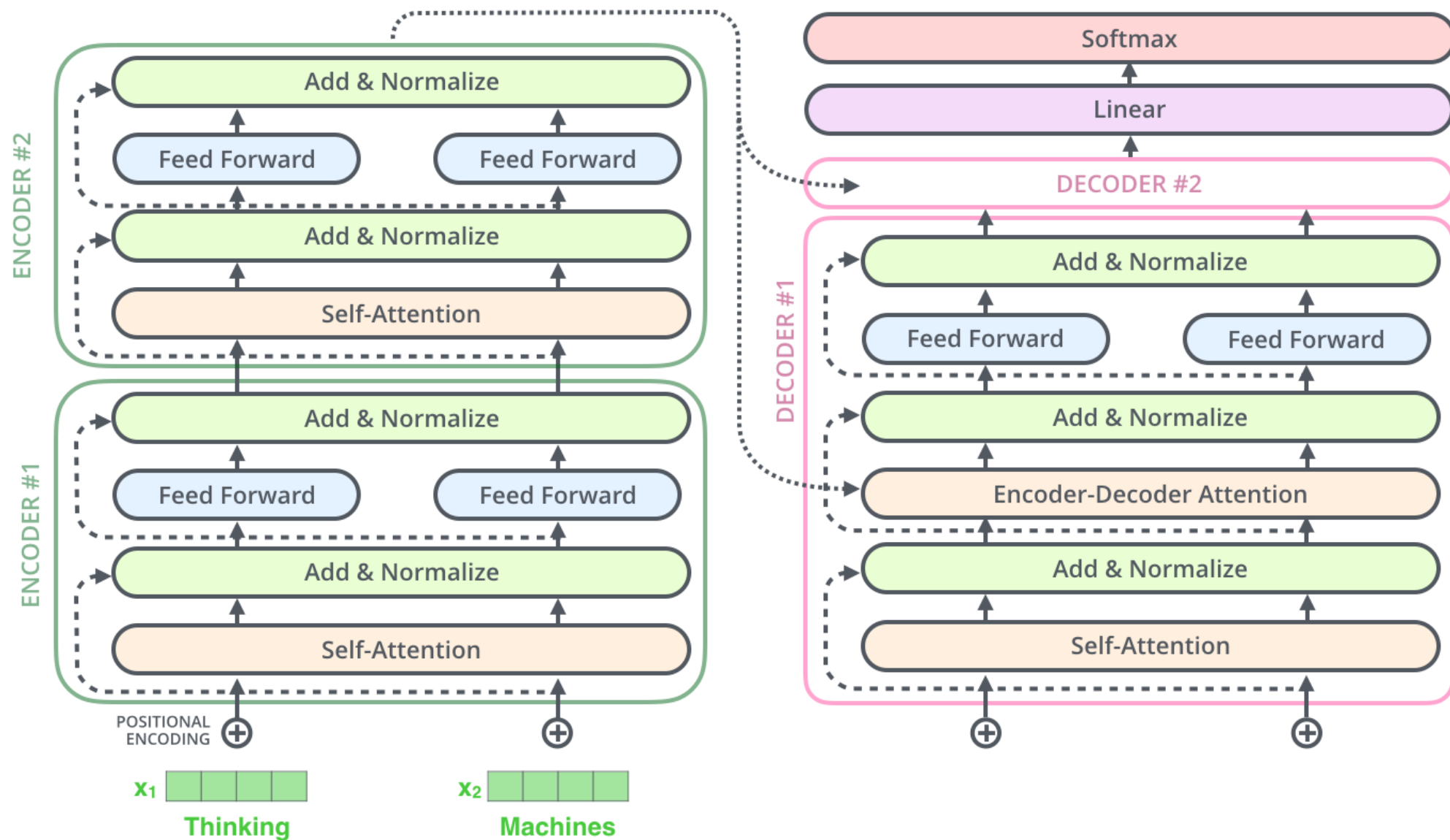
# The Residuals

- It is a skip connection where the original input is added back to the output of:
  - Self-attention
  - Feed-forward network
- Why do we need residuals?
  - To avoid vanishing gradient problems.

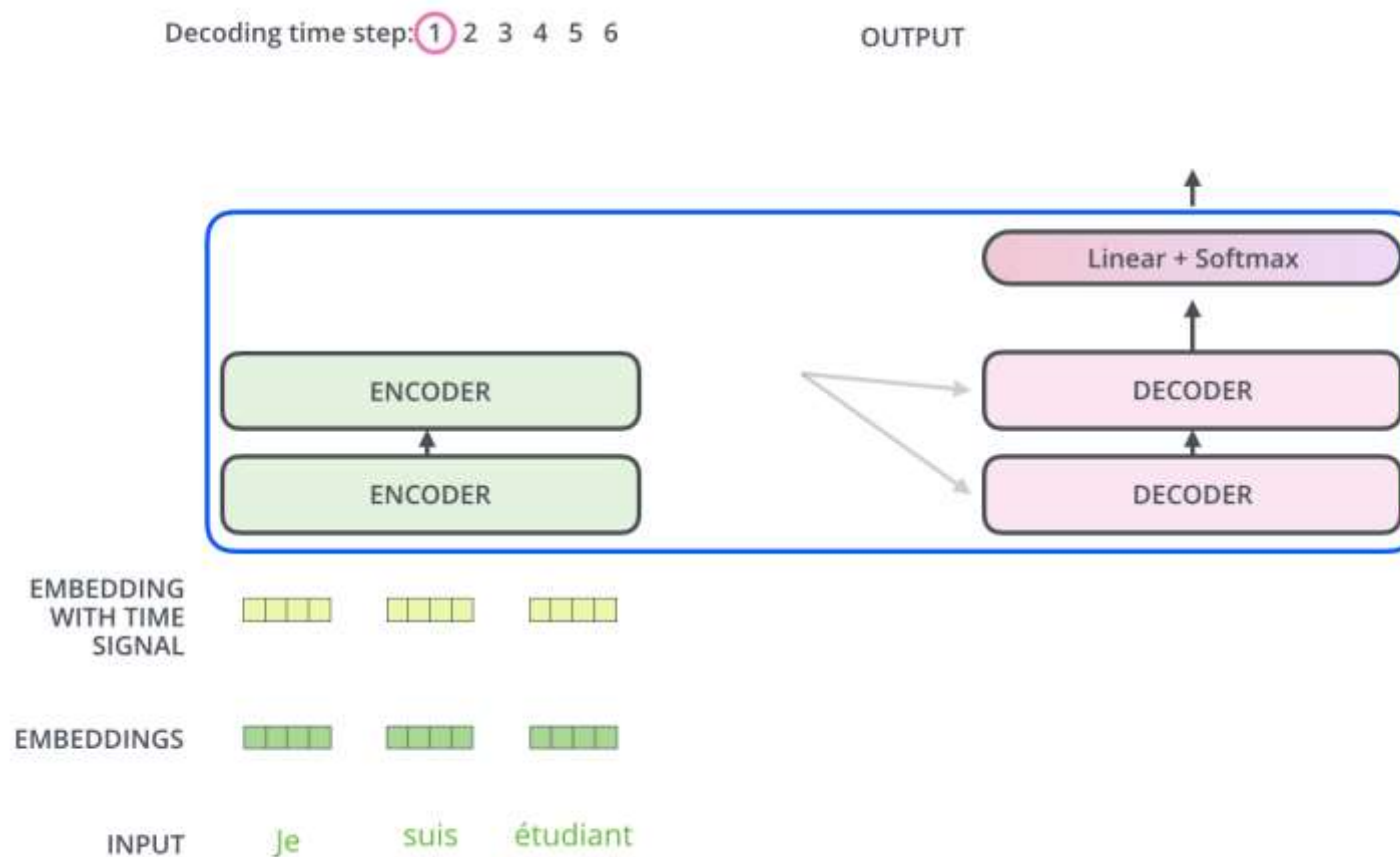






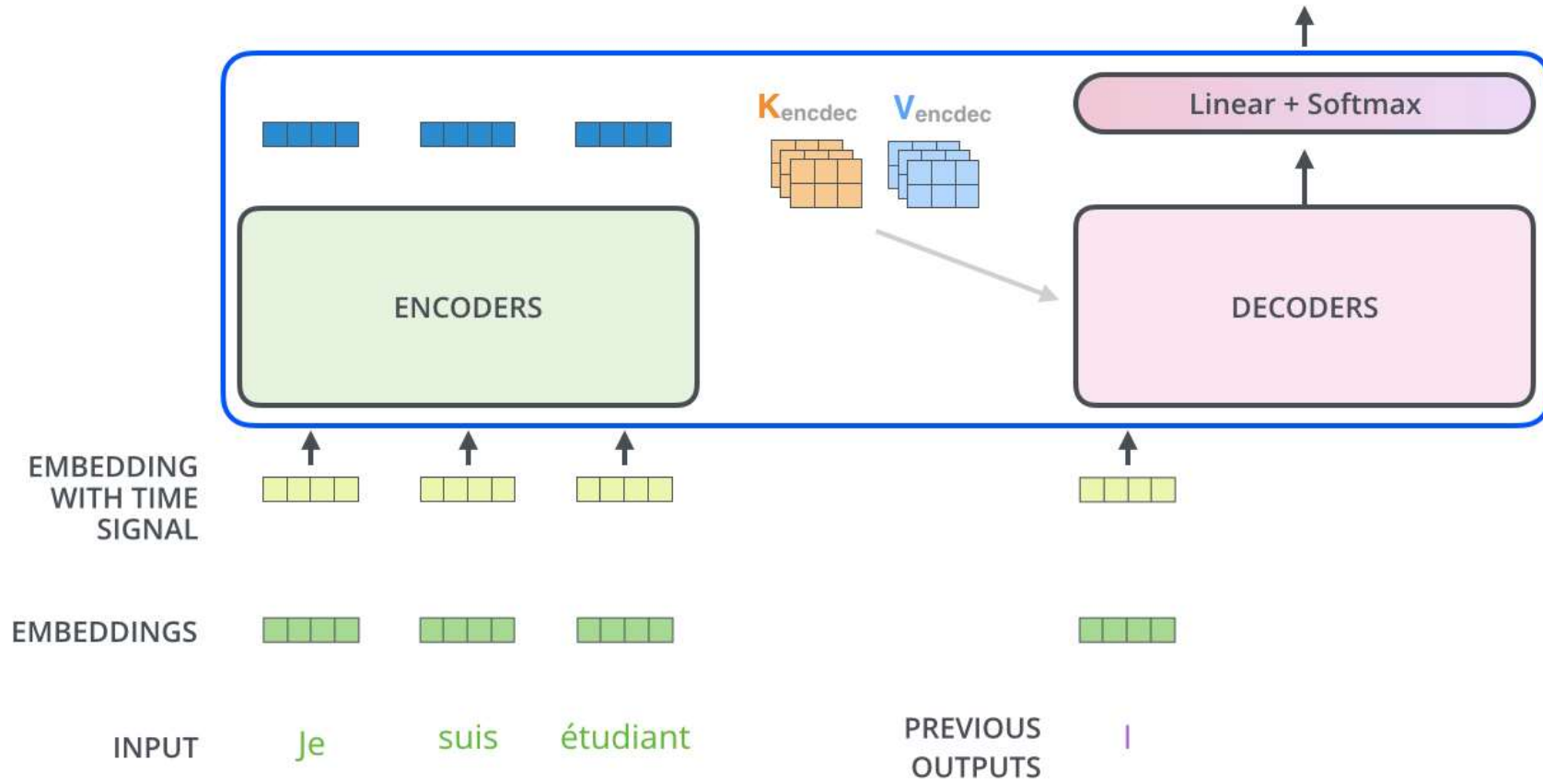


# The Decoder Side

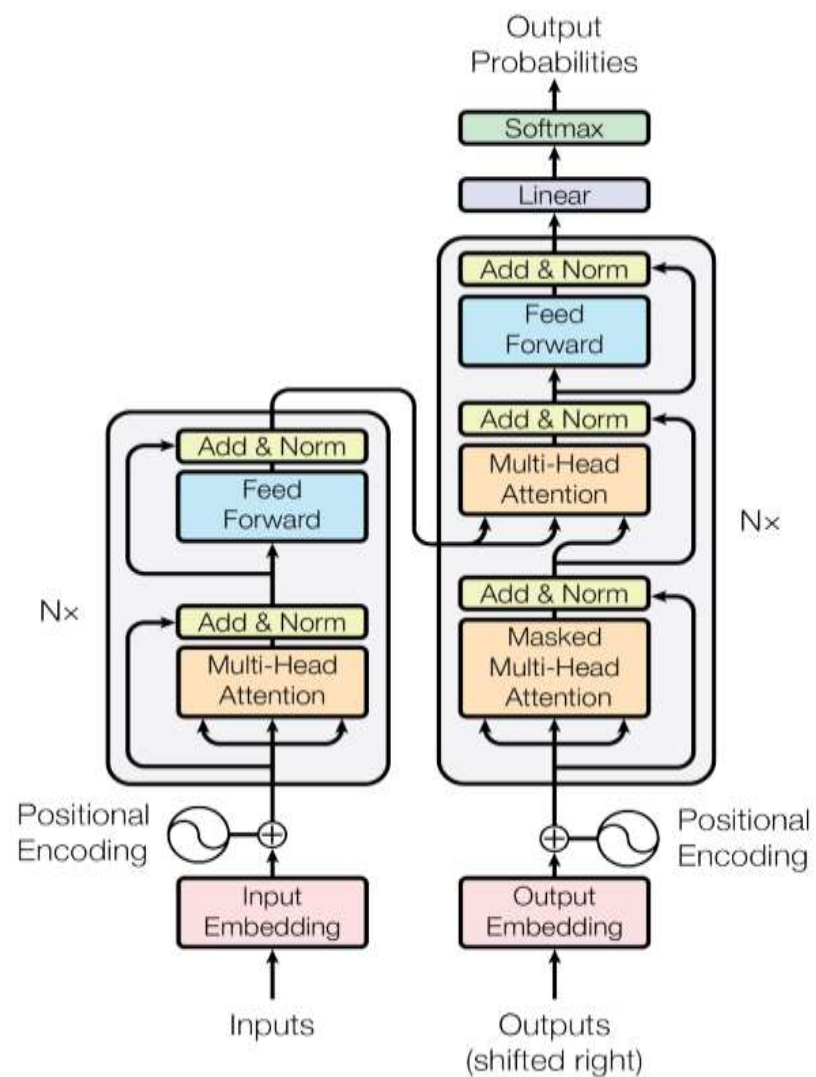


Decoding time step: 1 2 3 4 5 6

OUTPUT



# The architecture from the paper



# Reference

- <https://jalammar.github.io/illustrated-transformer/>

Thank You!