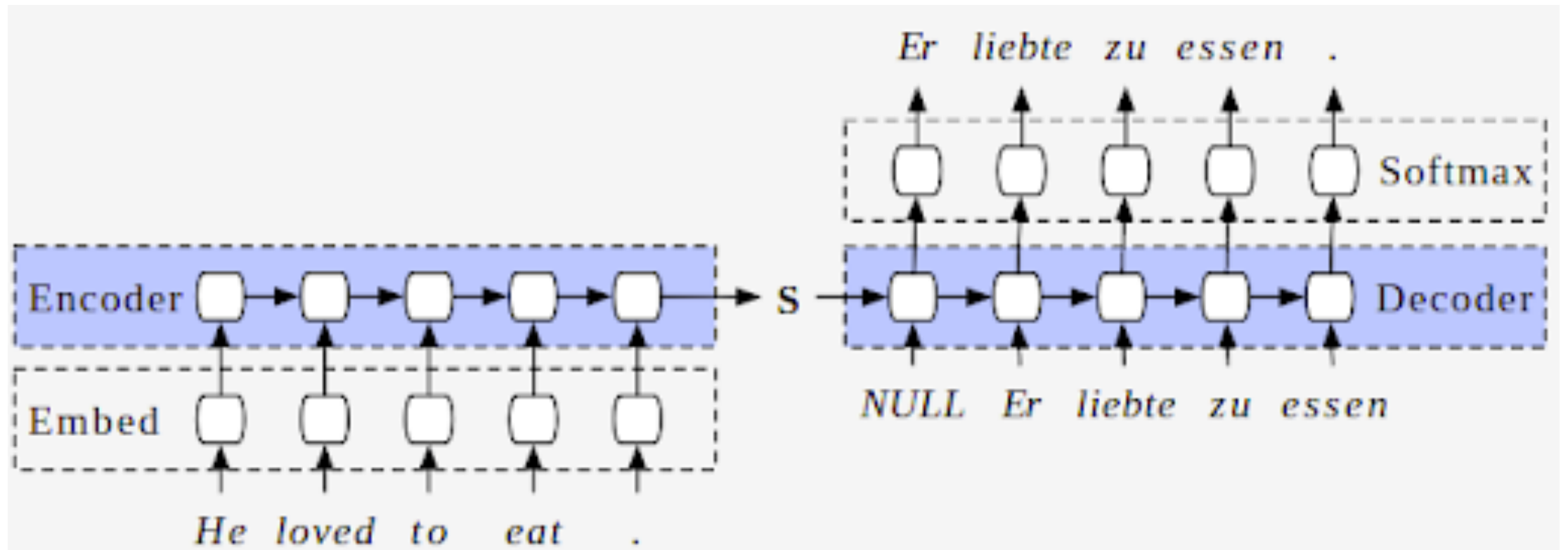
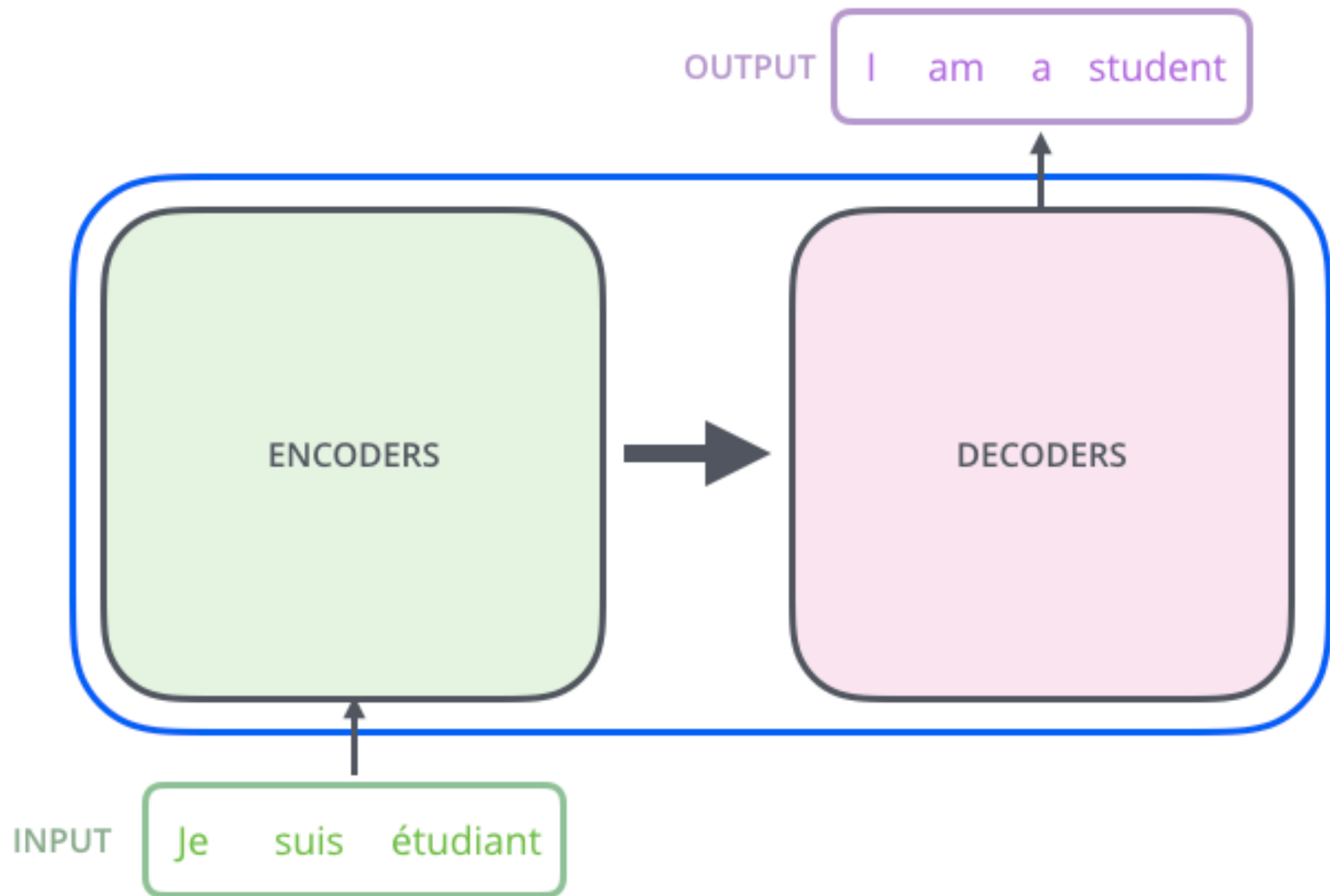


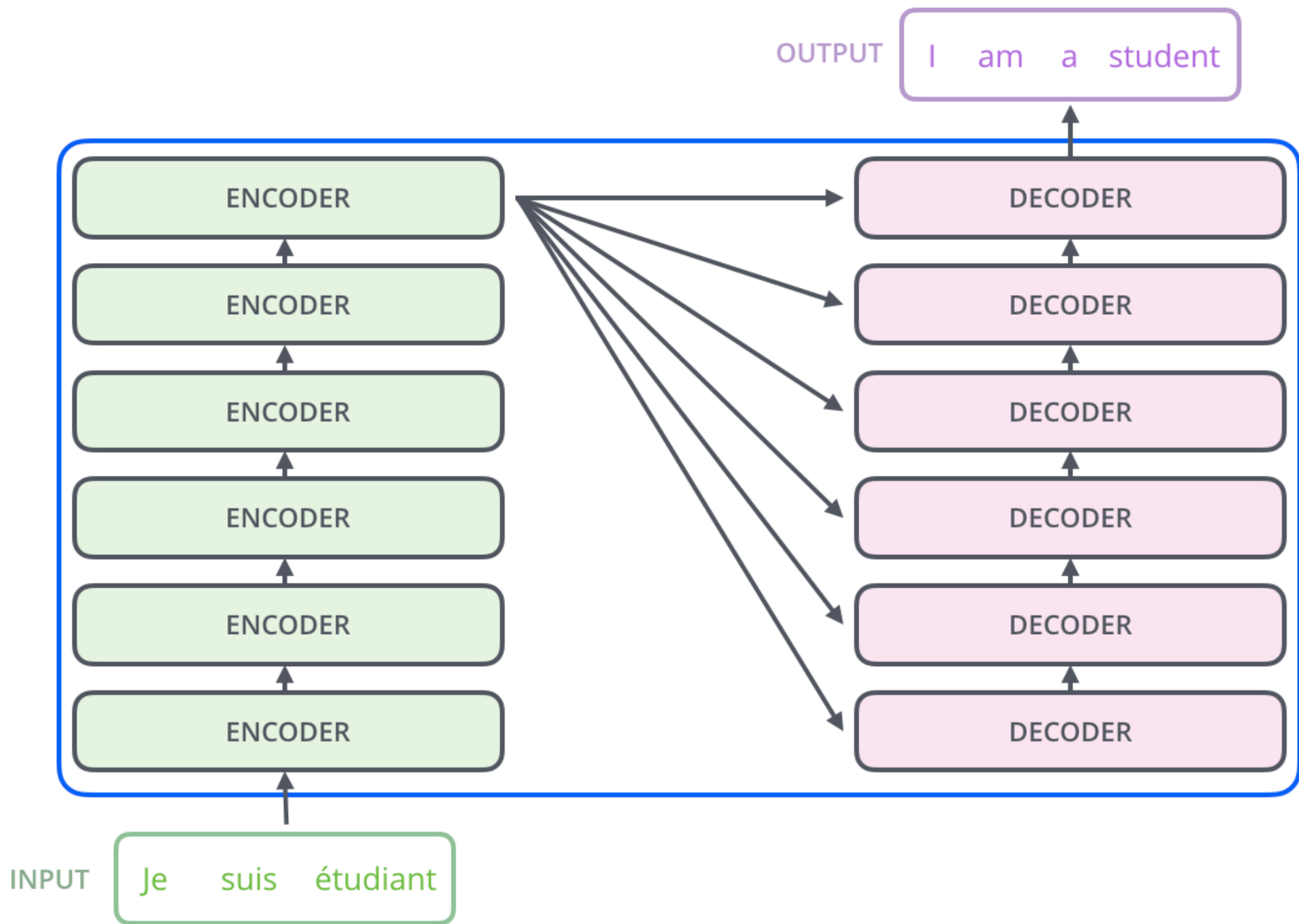


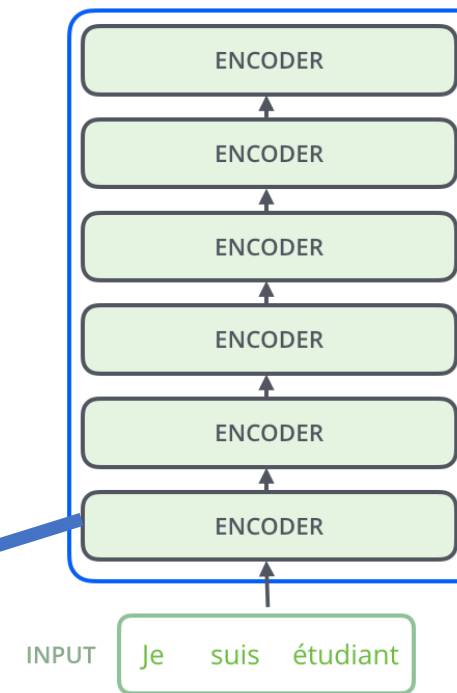
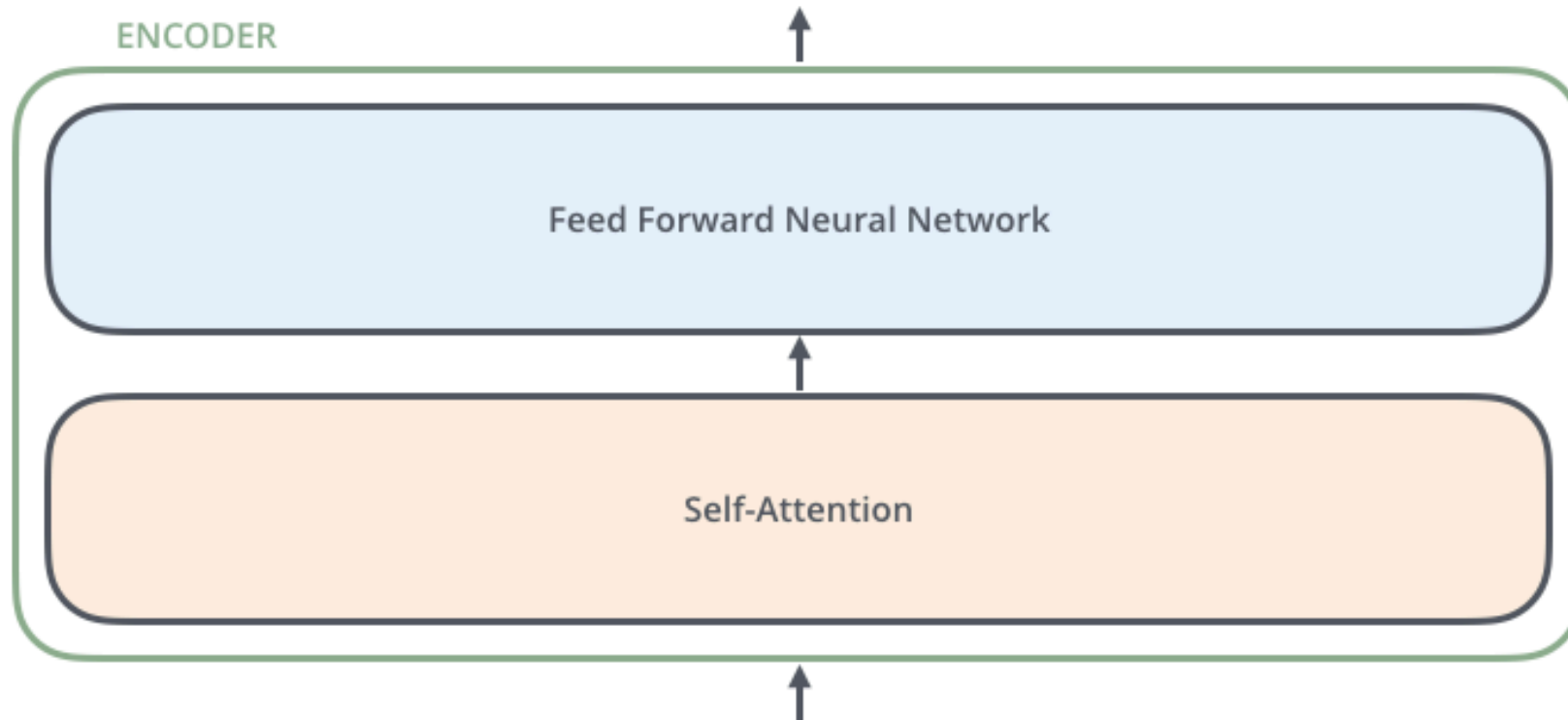
Transformer

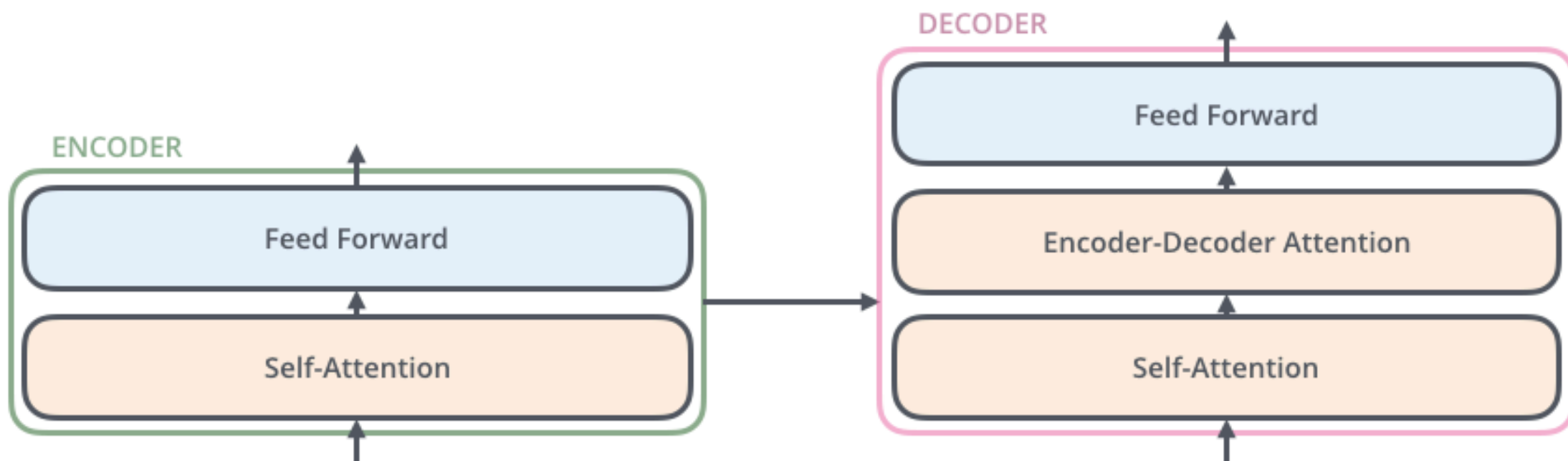
Seq2Seq Model





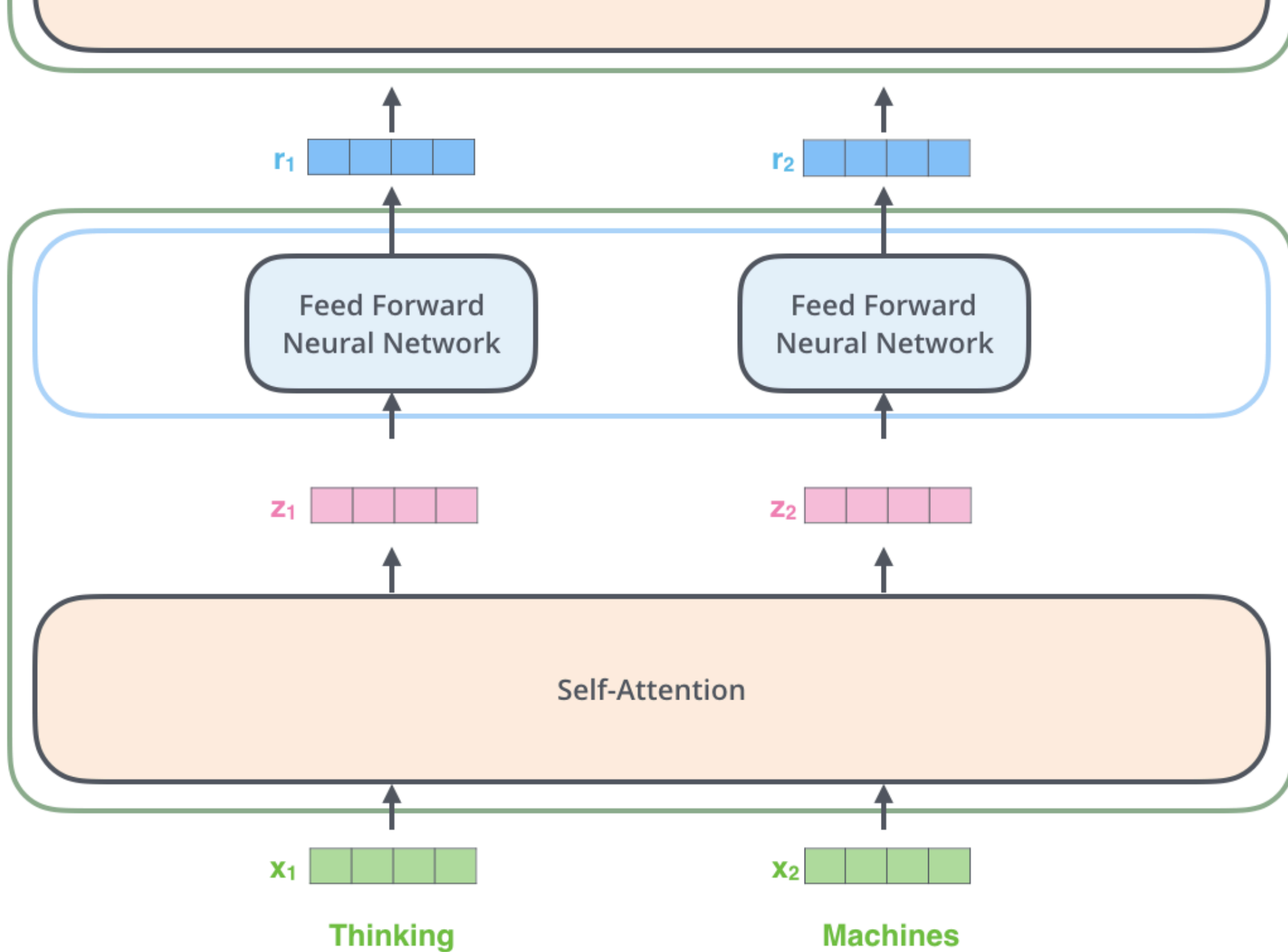




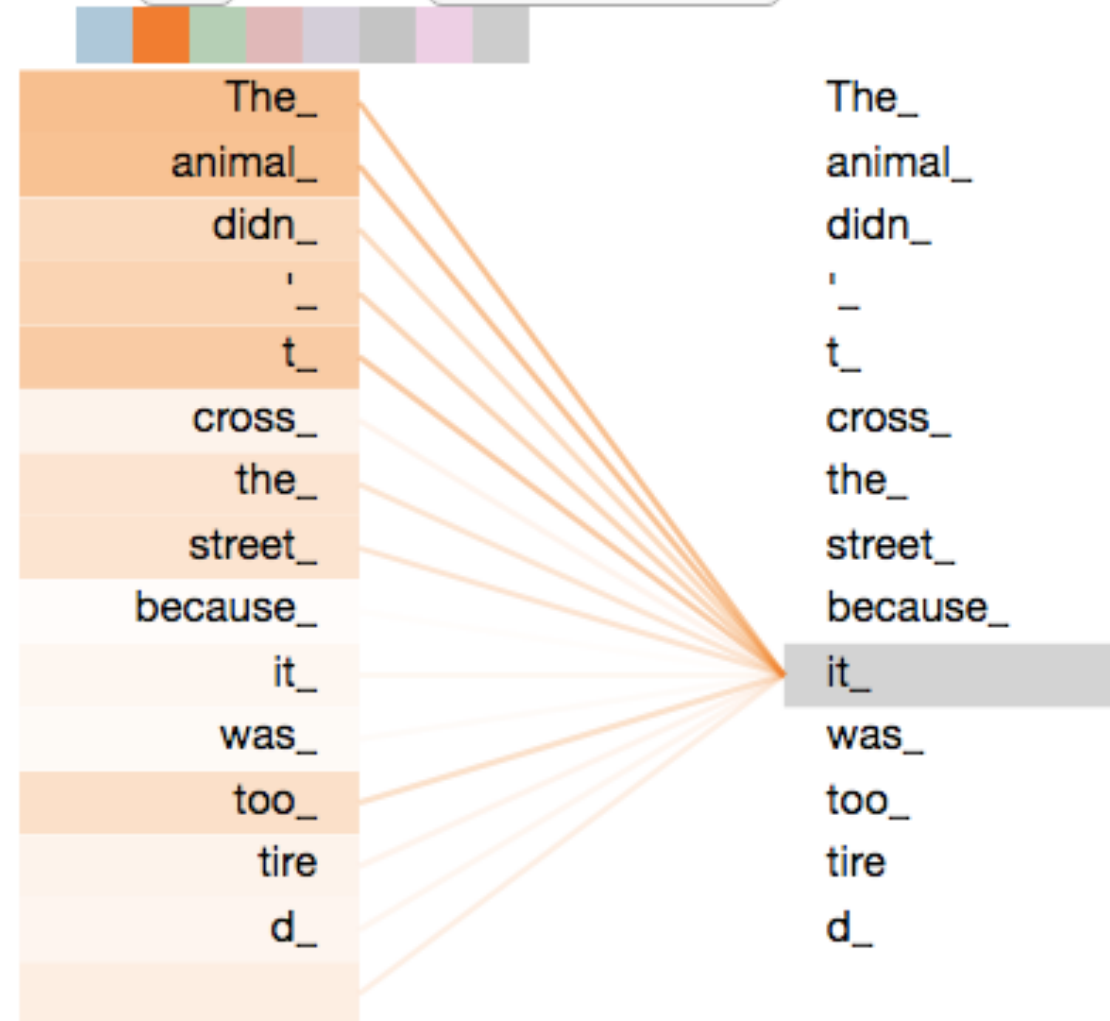


ENCODER #2

ENCODER #1



Layer: 5 Attention: Input - Input

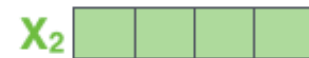
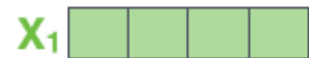


Input

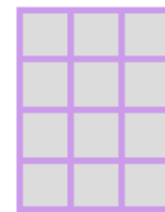
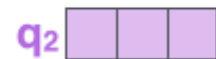
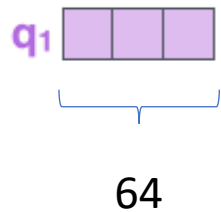
Thinking

Machines

Embedding

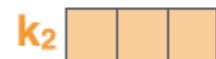
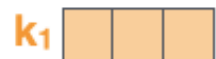


Queries



W^Q

Keys

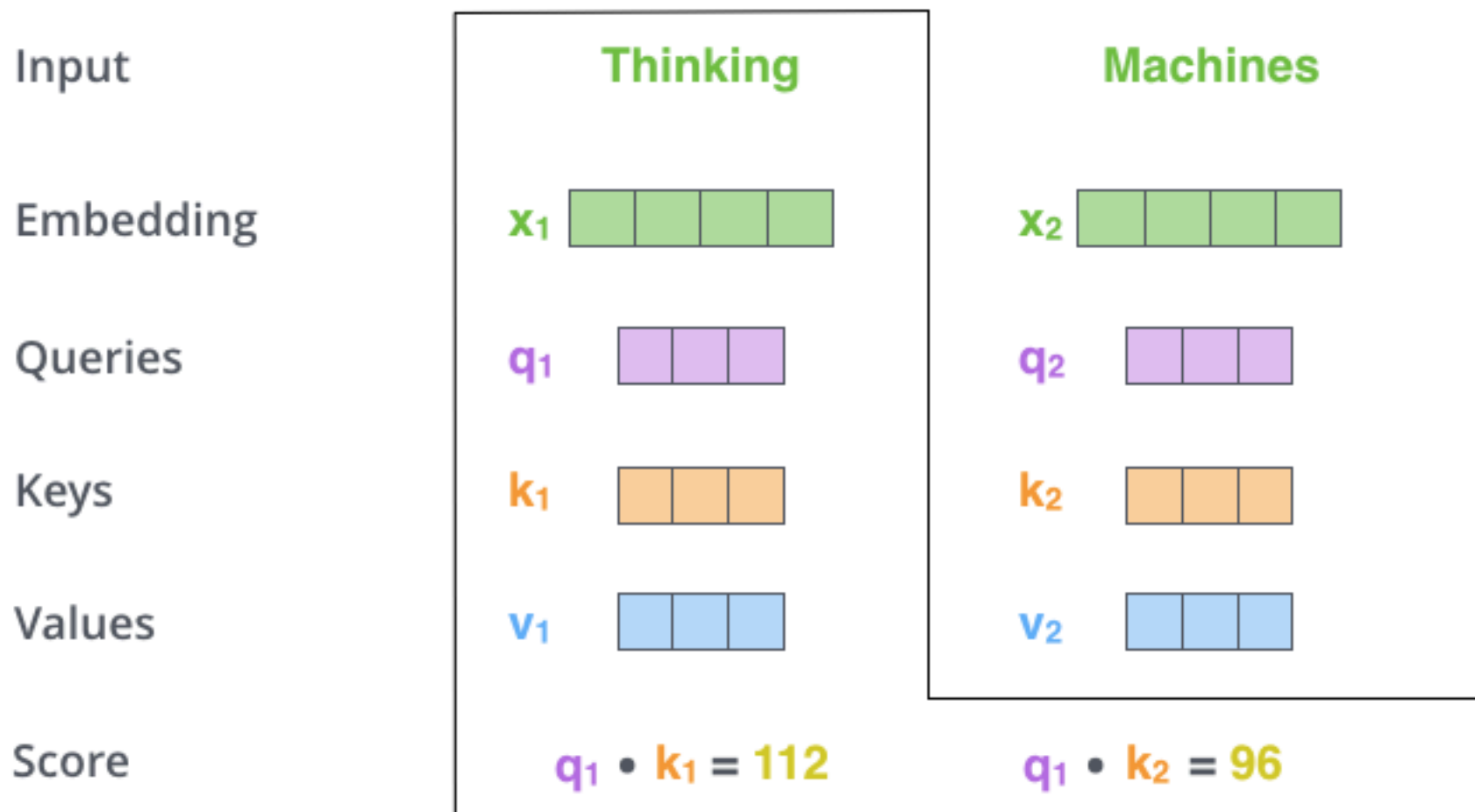


W^K

Values



W^V

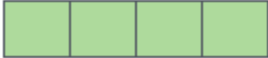


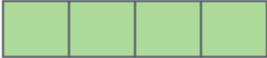
Input

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

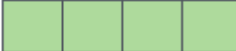
0.12

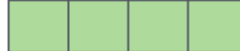
Input

Thinking

Machines

Embedding

\mathbf{x}_1 

\mathbf{x}_2 

Queries

\mathbf{q}_1 

\mathbf{q}_2 

Keys

\mathbf{k}_1 

\mathbf{k}_2 

Values

\mathbf{v}_1 

\mathbf{v}_2 

Score

$\mathbf{q}_1 \cdot \mathbf{k}_1 = 112$

$\mathbf{q}_1 \cdot \mathbf{k}_2 = 96$

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

14

12

Softmax

0.88

0.12

Softmax

X

Value

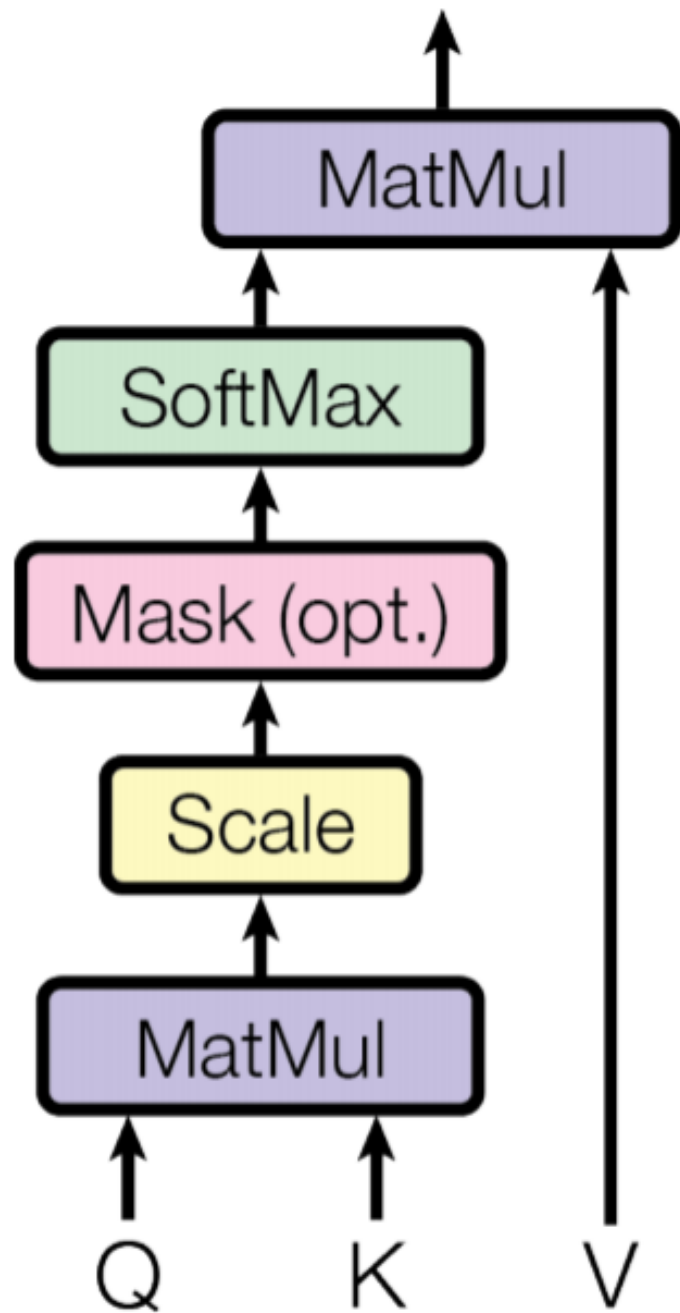
\mathbf{v}_1 

\mathbf{v}_2 

Sum

\mathbf{z}_1 

\mathbf{z}_2 



Scaled Dot-Product Attention

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

X **W^Q** **Q**

The diagram illustrates the matrix multiplication $X \times W^Q = Q$. Matrix X is a 2x4 grid of green squares. Matrix W^Q is a 4x4 grid of light purple squares. Matrix Q is a 2x4 grid of light purple squares. The multiplication is represented by a large 'x' between X and W^Q , and an equals sign between W^Q and Q .

Diagram illustrating matrix multiplication: X (2x4 green grid) multiplied by W^K (4x4 orange grid) equals K (2x3 orange grid).

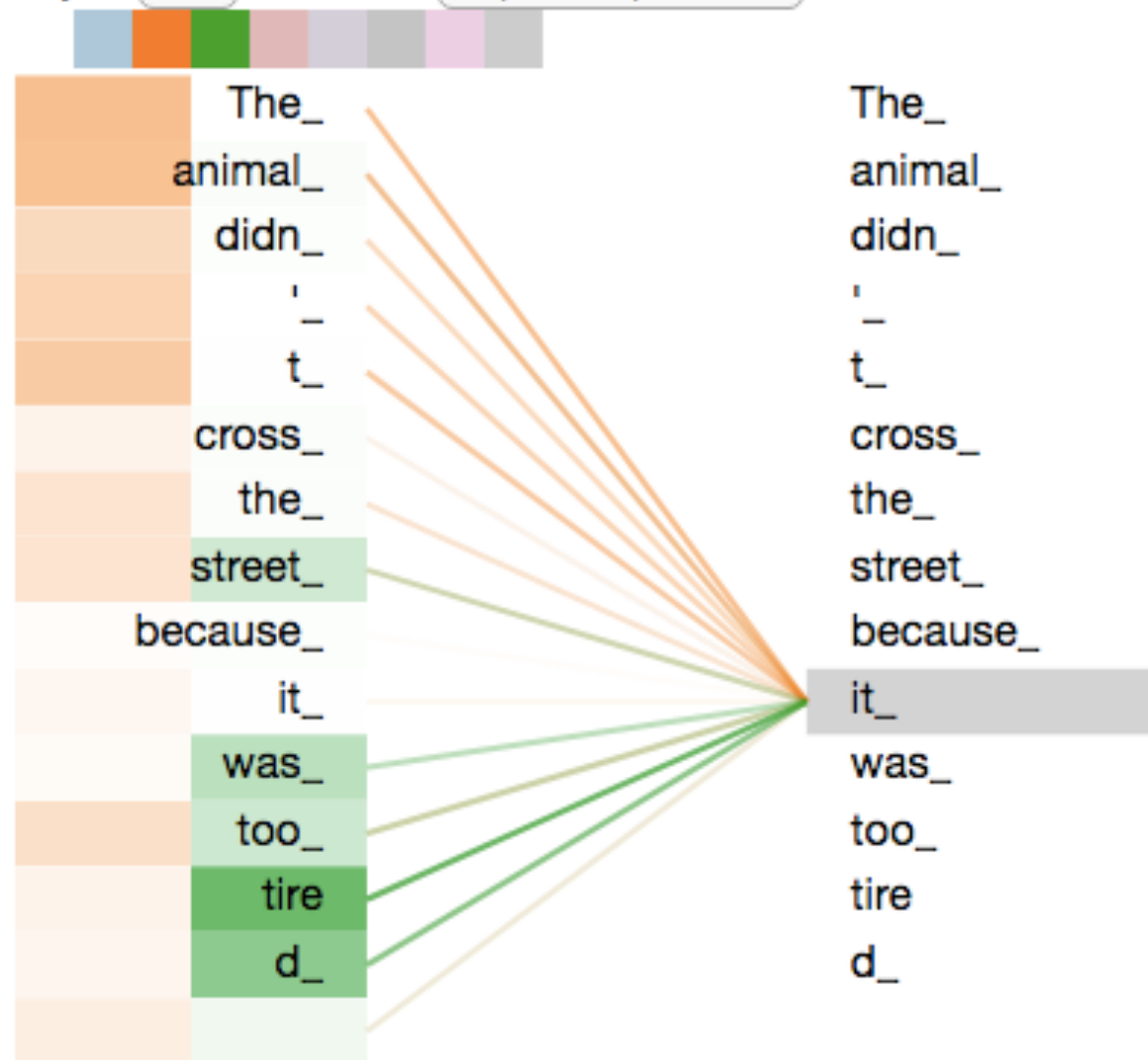
X **W^v** **V**

The diagram illustrates the multiplication of two matrices. On the left is matrix **X**, a 2x4 grid of green squares. In the middle is matrix **W^v**, a 4x3 grid of blue squares. To the right of **W^v** is an equals sign, followed by matrix **V**, a 2x3 grid of blue squares. A large 'x' symbol is placed between **X** and **W^v**, and another large 'x' symbol is placed between **W^v** and **V**.

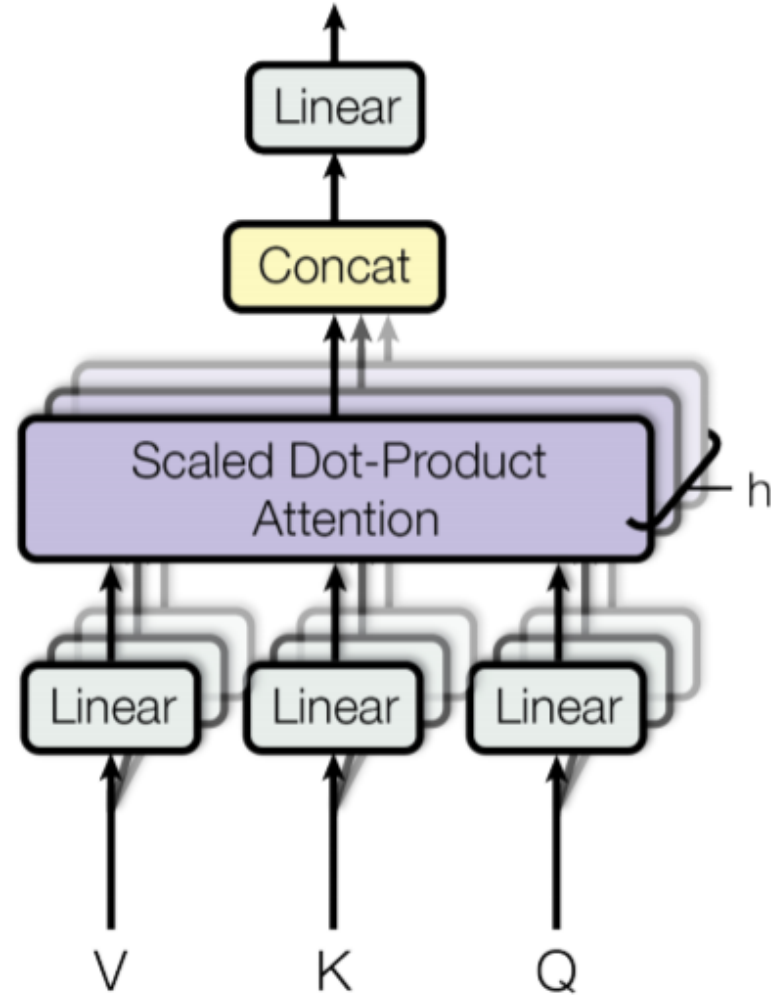
$$\text{softmax} \left(\frac{
 \begin{array}{c} \text{Q} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \times \begin{array}{c} \text{K}^T \\ \begin{array}{|c|c|} \hline \square & \square \\ \hline \square & \square \\ \hline \square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

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

Layer: 5 Attention: Input - Input



Multi-Head Attention



1) This is our input sentence*

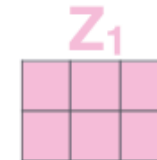
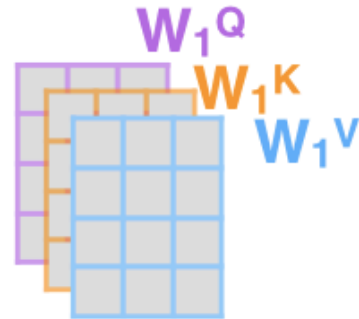
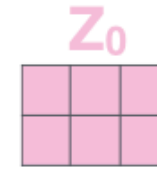
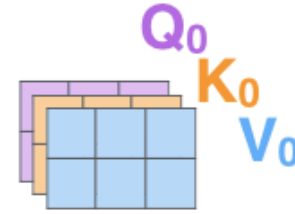
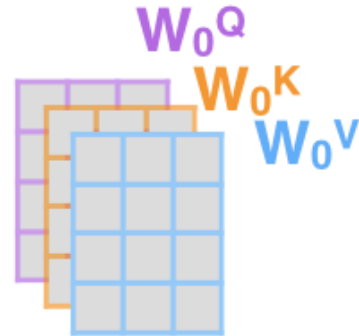
2) We embed each word*

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

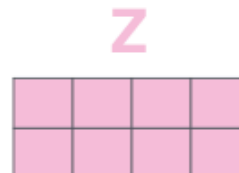
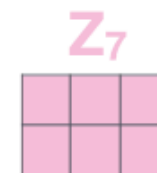
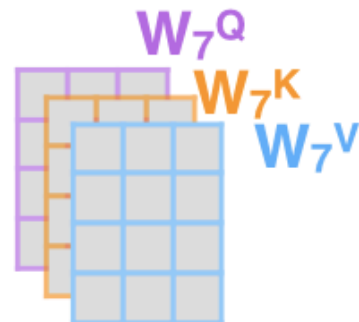
Thinking
Machines



...

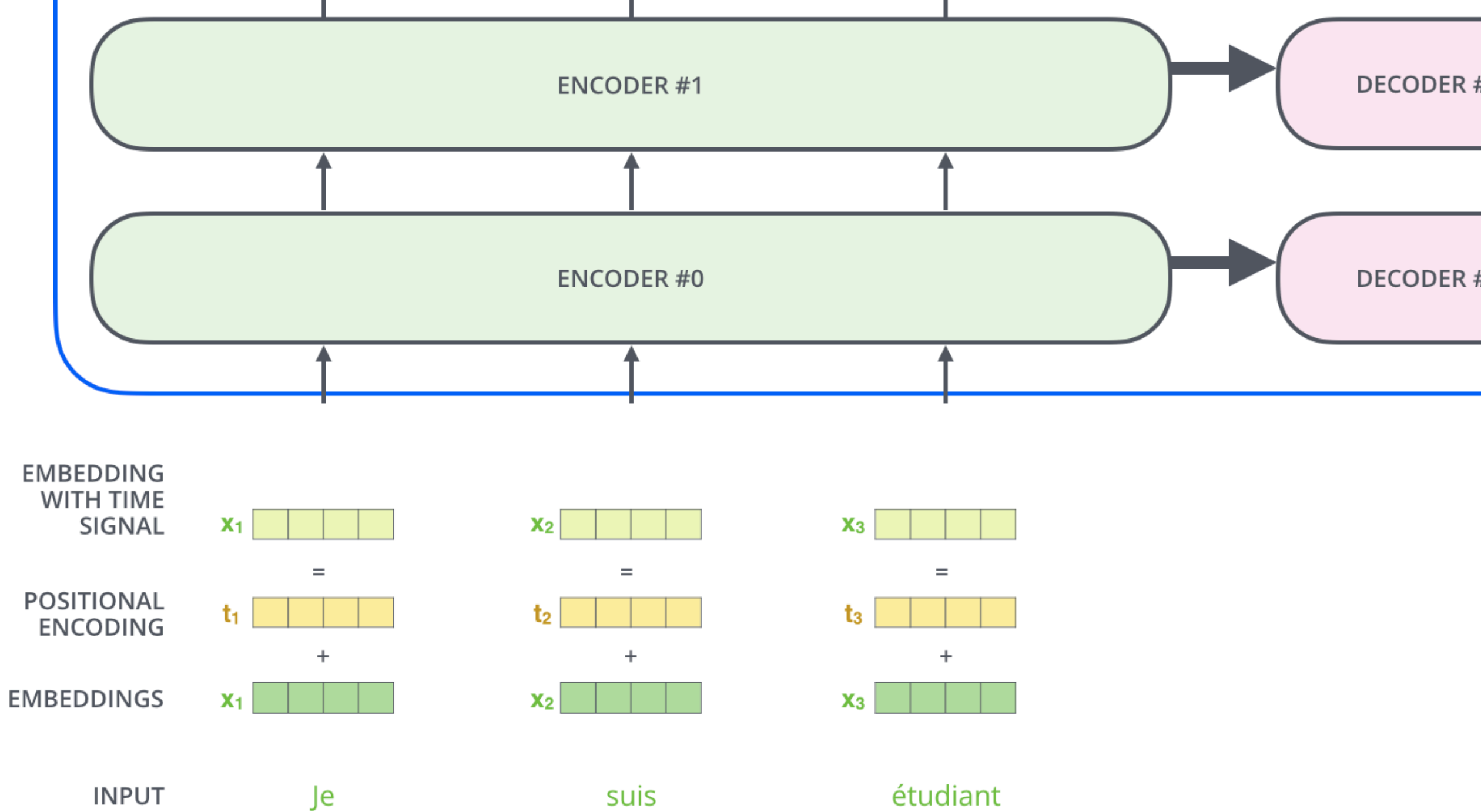
...

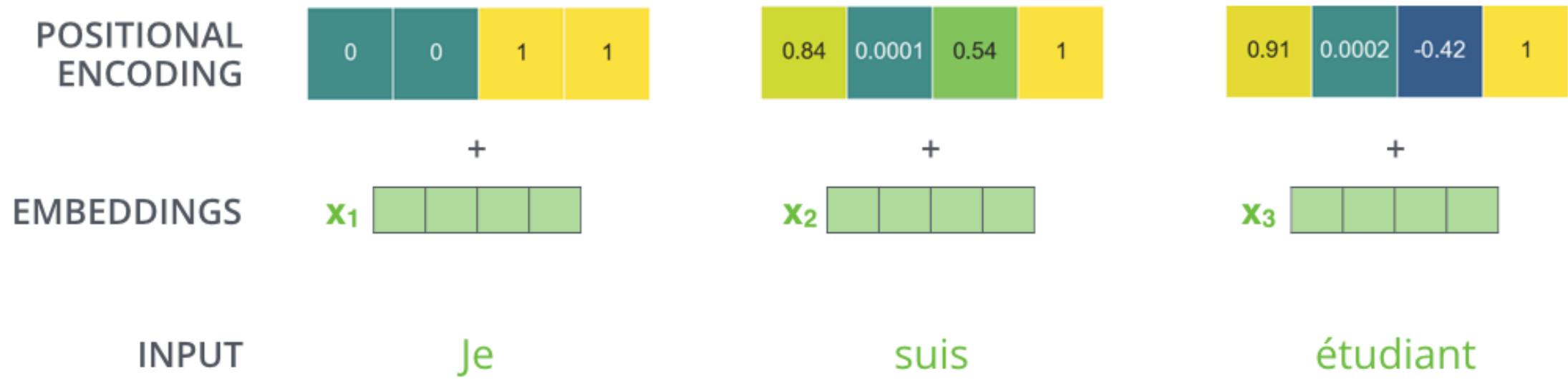
...

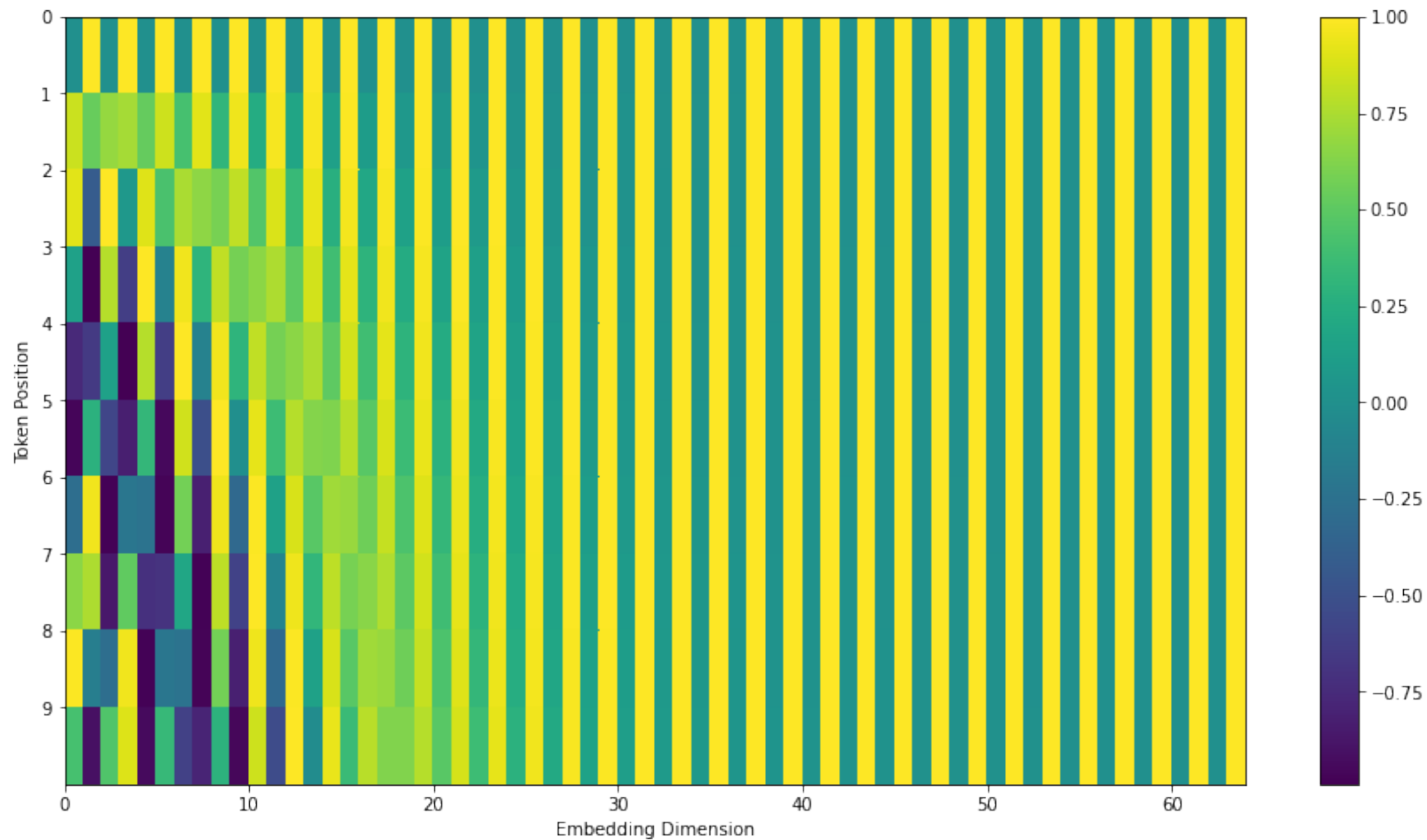


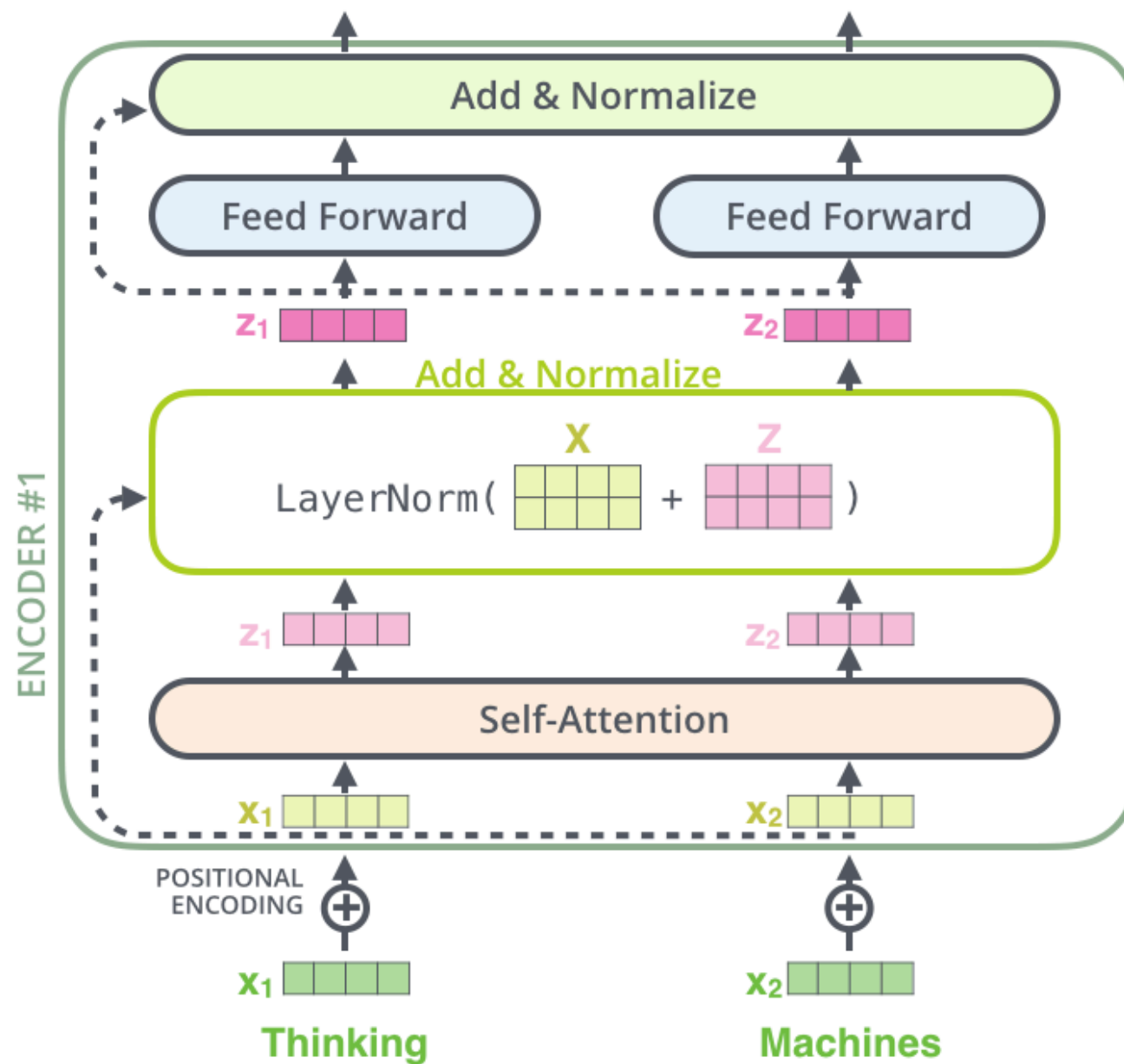
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

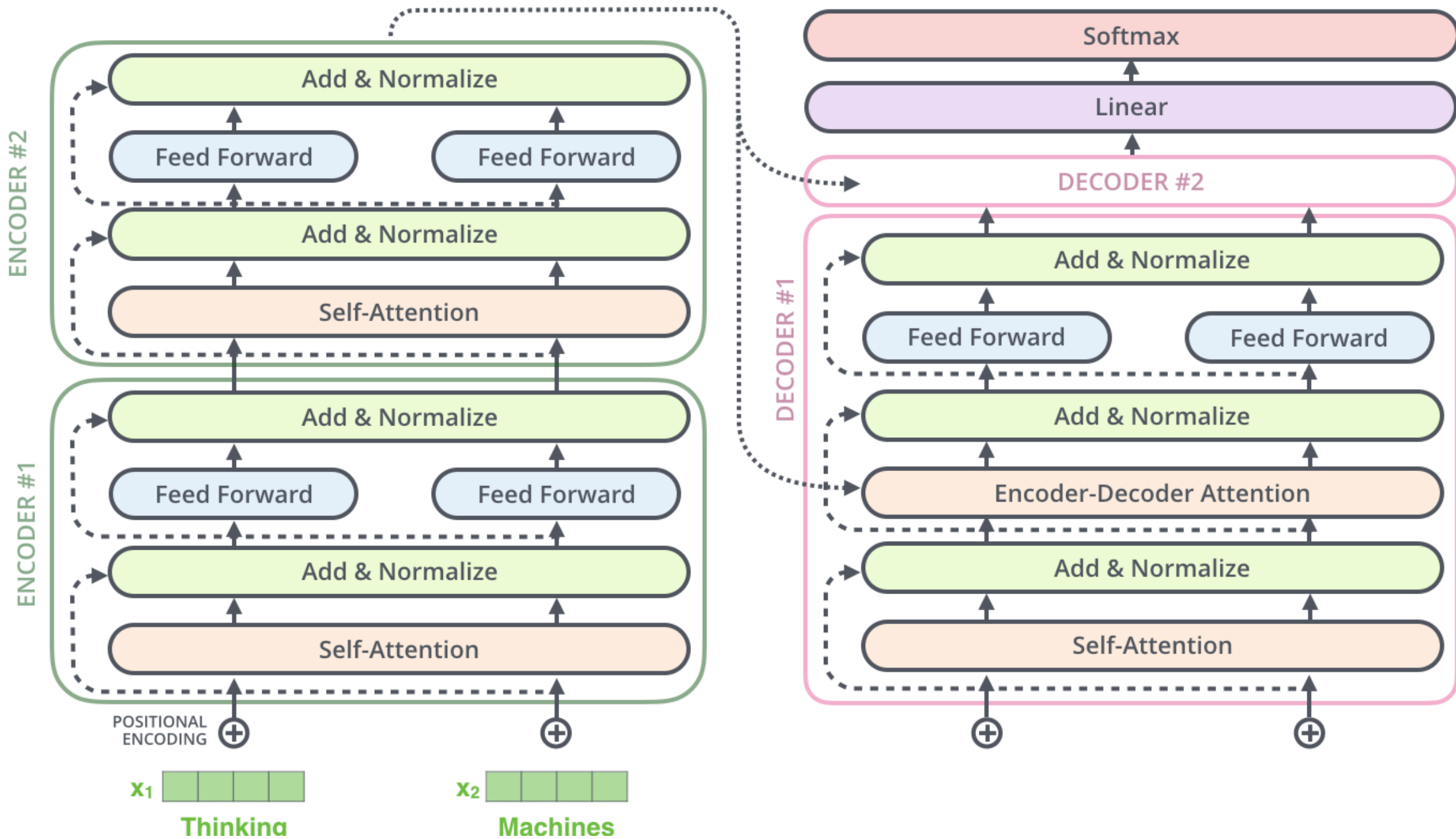






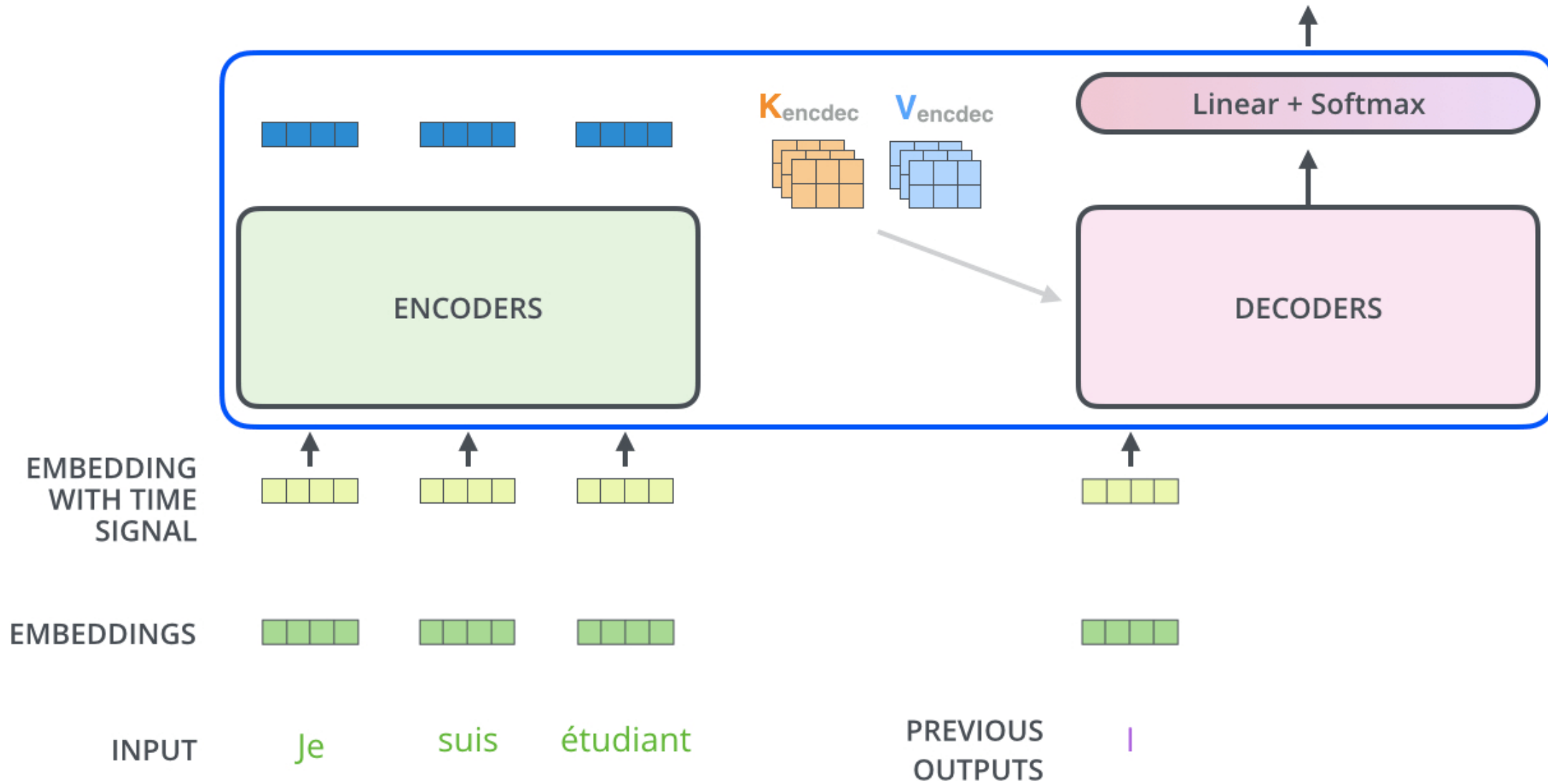






Decoding time step: 1 2 3 4 5 6

OUTPUT |



Which word in our vocabulary
is associated with this index?

Get the index of the cell
with the highest value
(**argmax**)

log_probs



am

5

Softmax

logits



Linear

Decoder stack output





GPT-2

DECODER

...

DECODER

DECODER



BERT

ENCODER

...

ENCODER

ENCODER



TRANSFORMER XL

RECURRENT DECODER

...

RECURRENT DECODER

RECURRENT DECODER

Reference

- “Attention Is All You Need” paper.
<https://arxiv.org/pdf/1706.03762.pdf>
- Visualization: <http://jalammar.github.io/illustrated-transformer/>