
Knowledge Discovery & Data Mining

— Transformers —

Instructor: Yong Zhuang

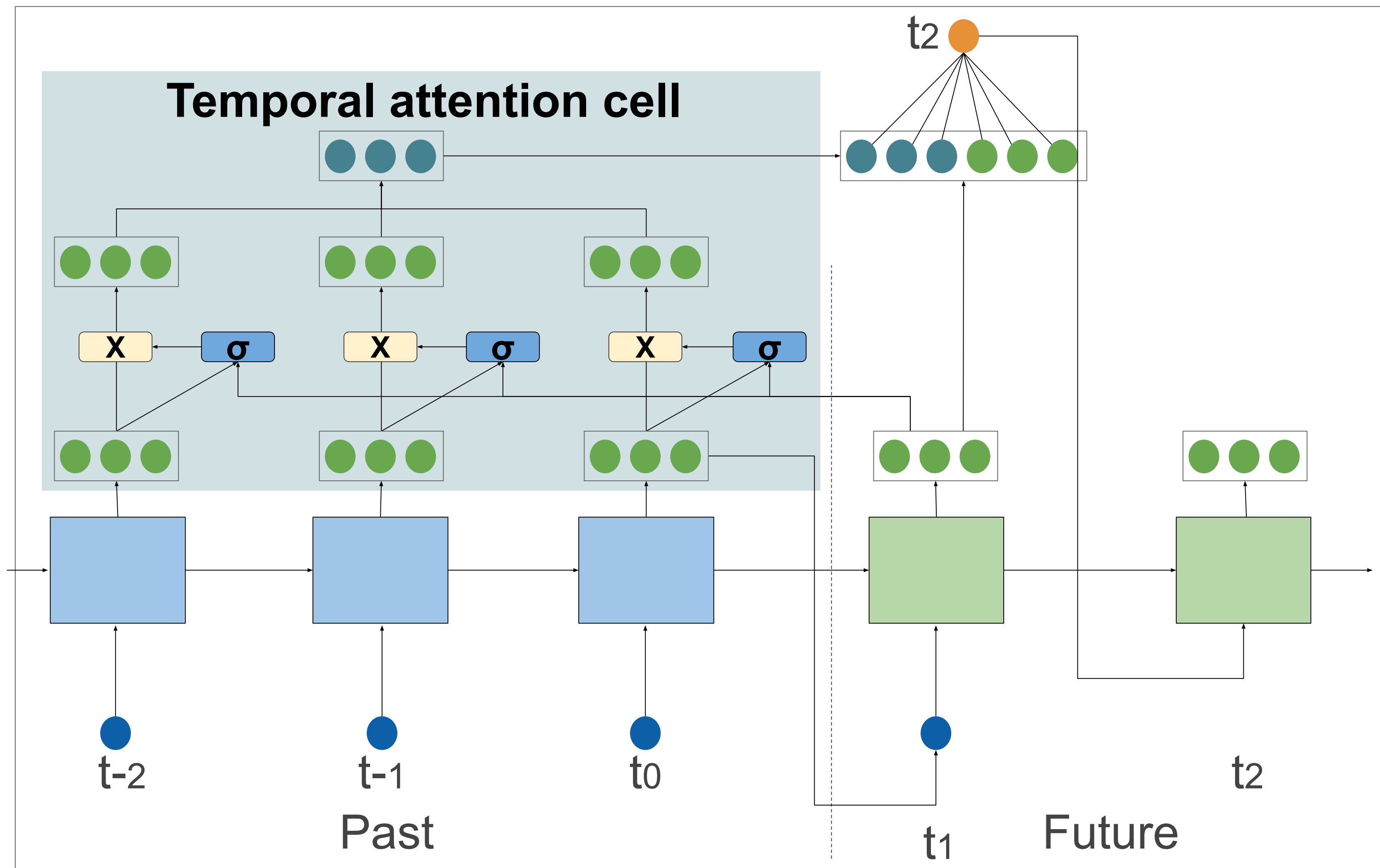
yong.zhuang@gvsu.edu

Based on the original version at <https://jalammar.github.io/illustrated-transformer/>

Outline

- Transformers
 - Self-Attention, Softmax,
 - Multi-Headed
 - Self-Attention,
 - Positional Encoding,
 - Residual Connection,
 - Layer Normalization,
 - Encoder-Decoder

Temporal attention injection



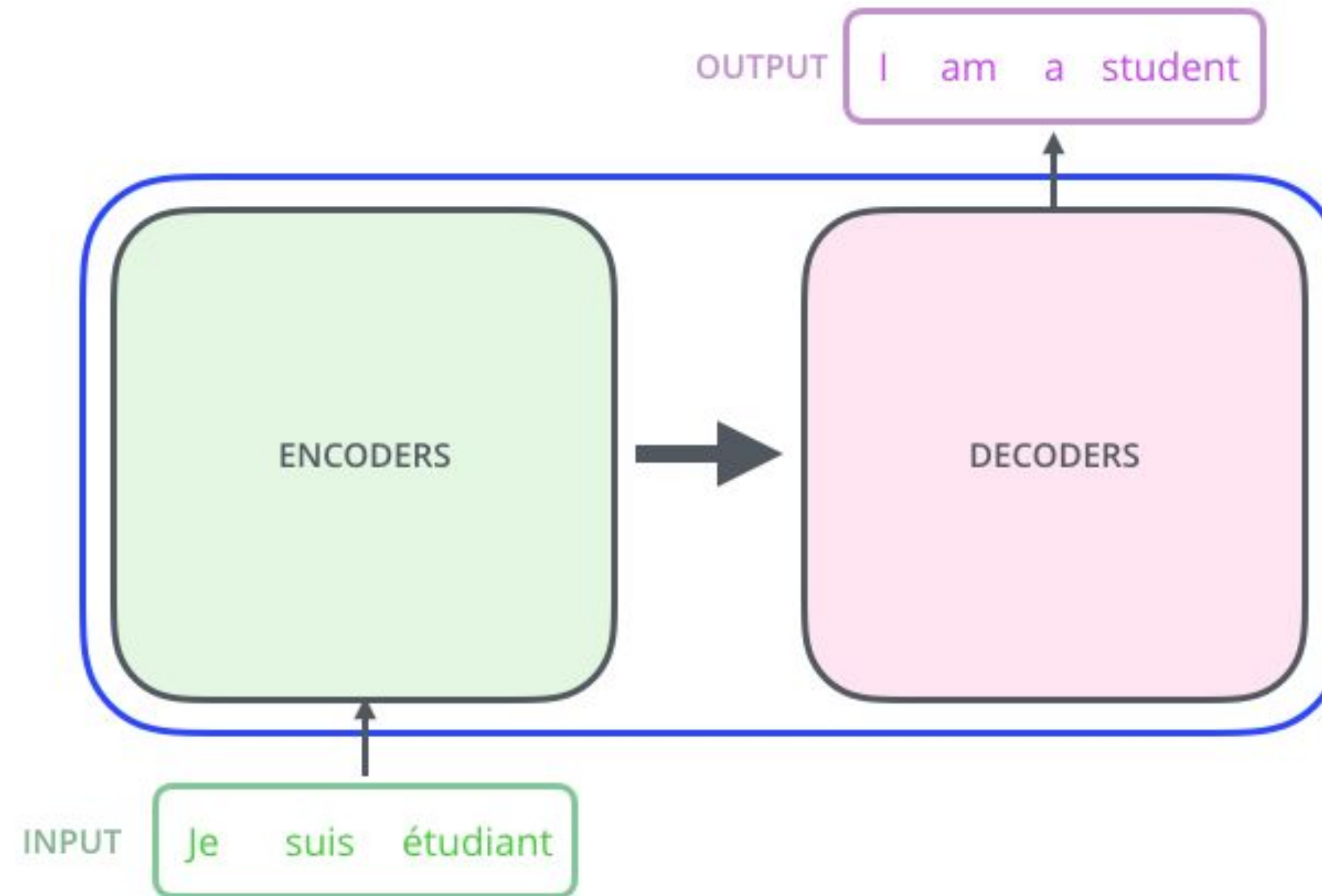
In light of the performance enhancements brought by attention mechanisms, is the RNN component still necessary in our models, or can attention mechanisms effectively replace the traditional role of RNNs?

Transformer

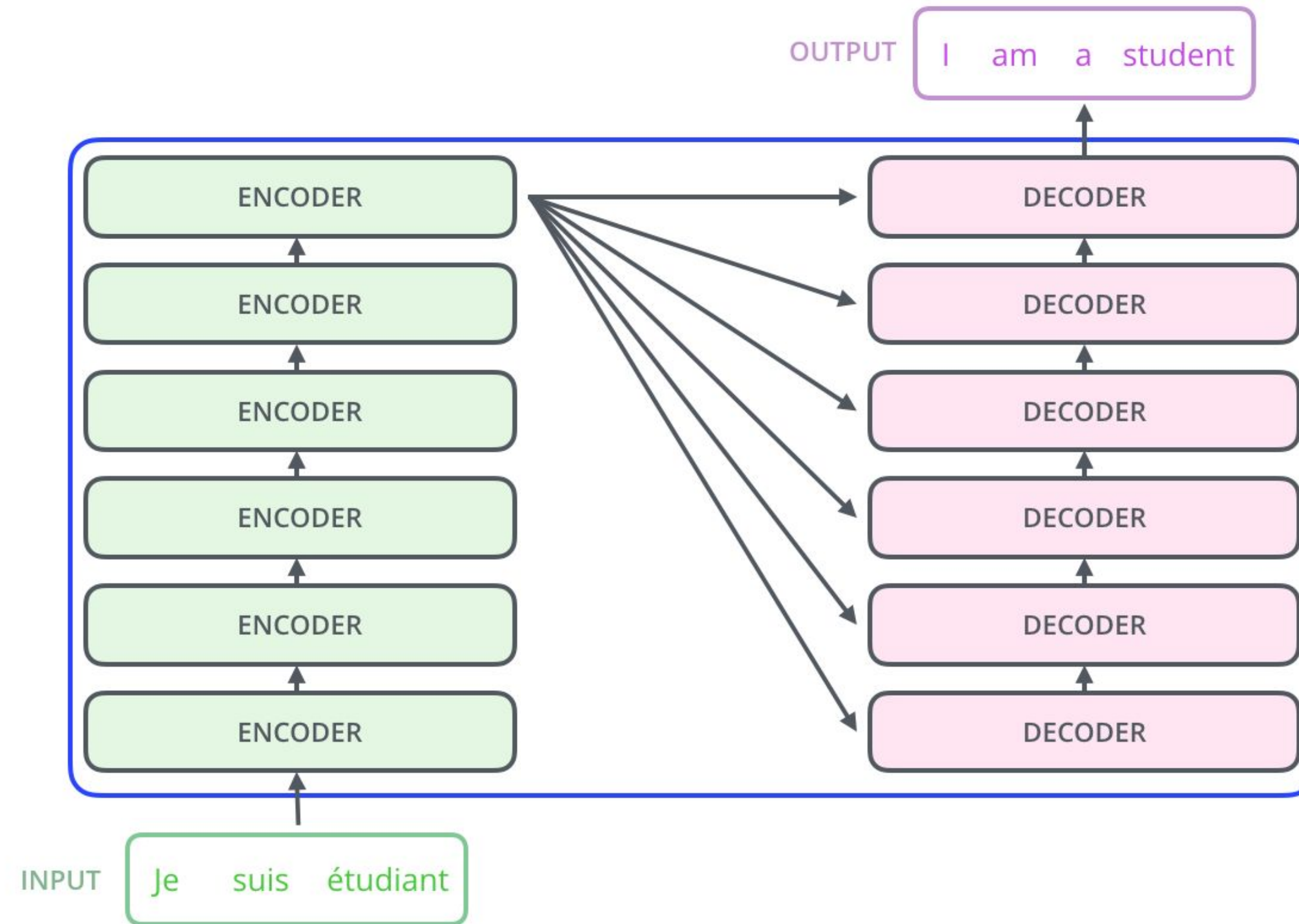


Attention is all you need

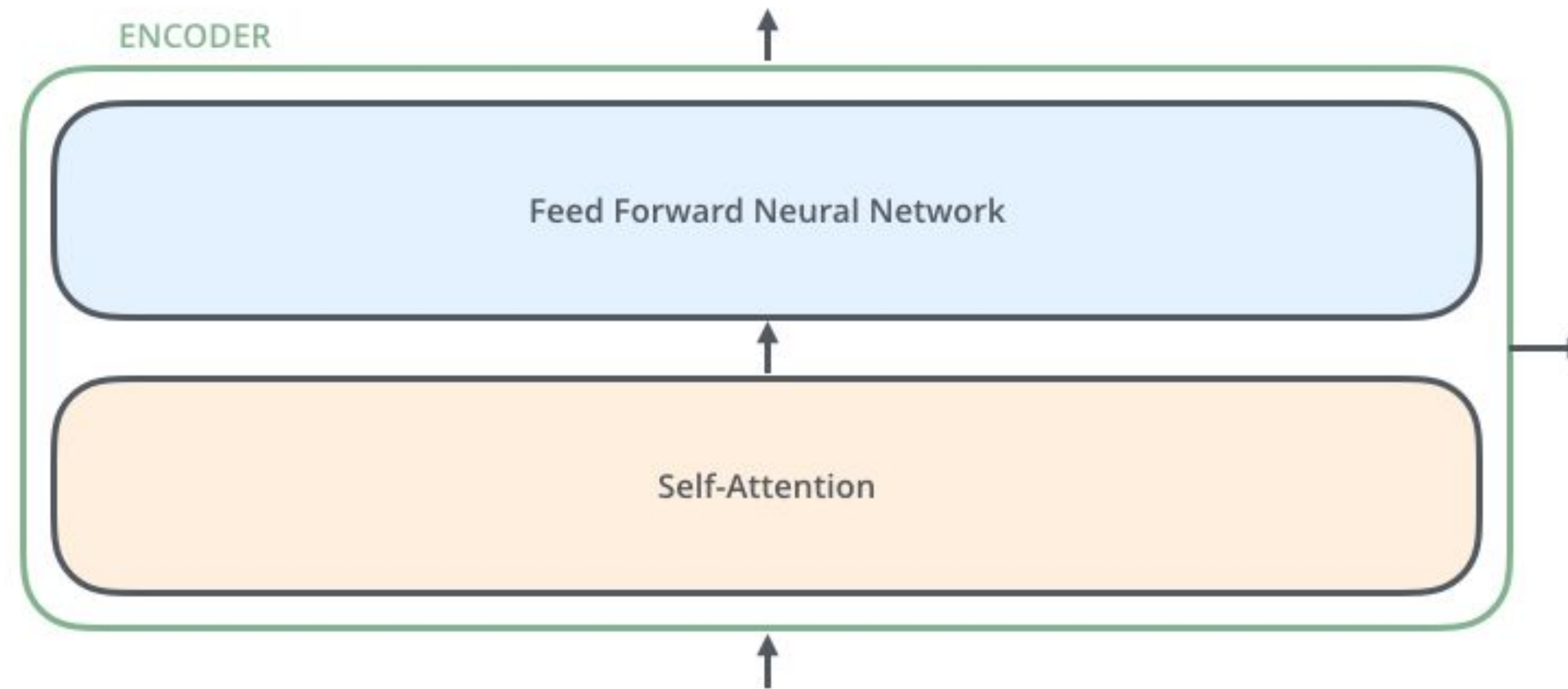
Encoder-Decoder Architecture



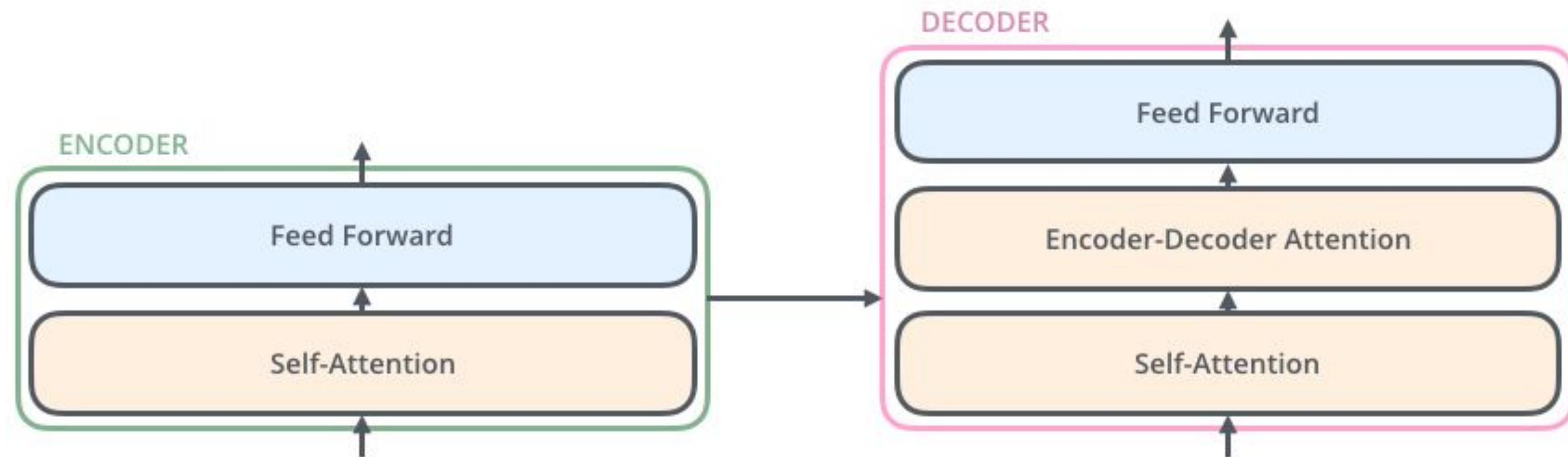
Encoder-Decoder Architecture



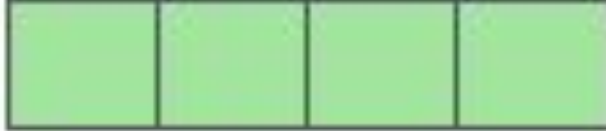
Encoder



Decoder



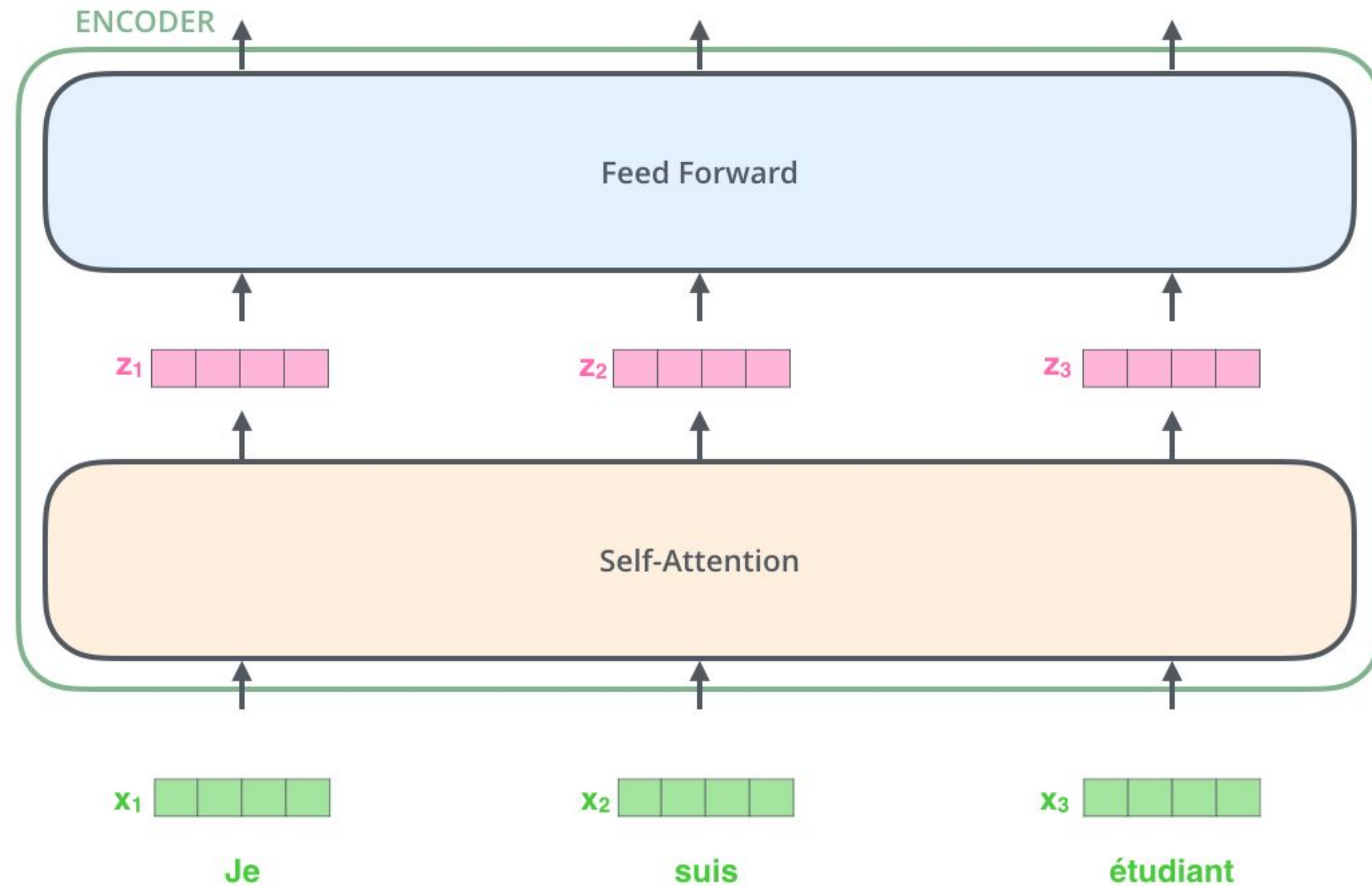
Word Embedding

x_1 
Je

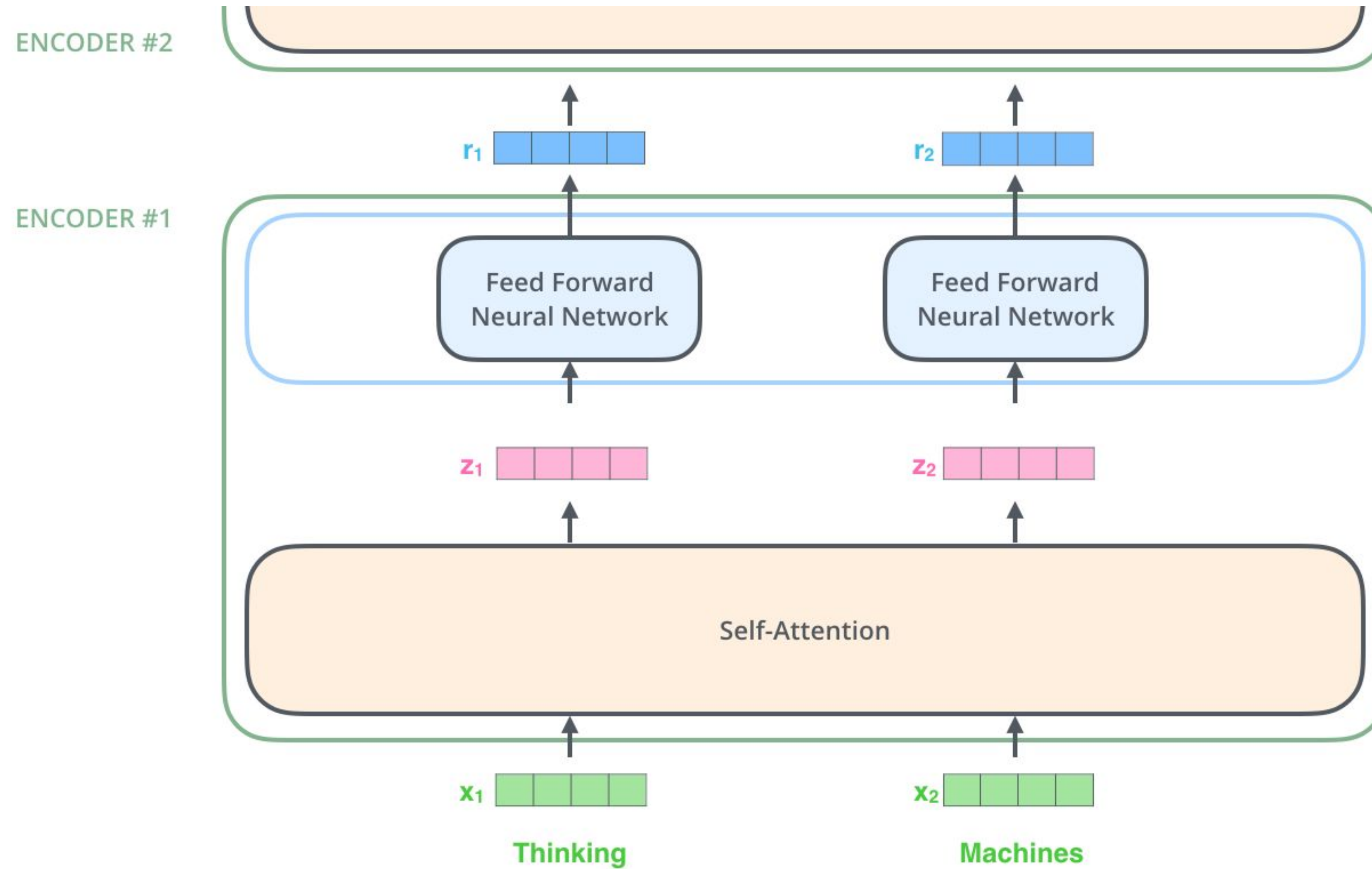
x_2 
suis

x_3 
étudiant

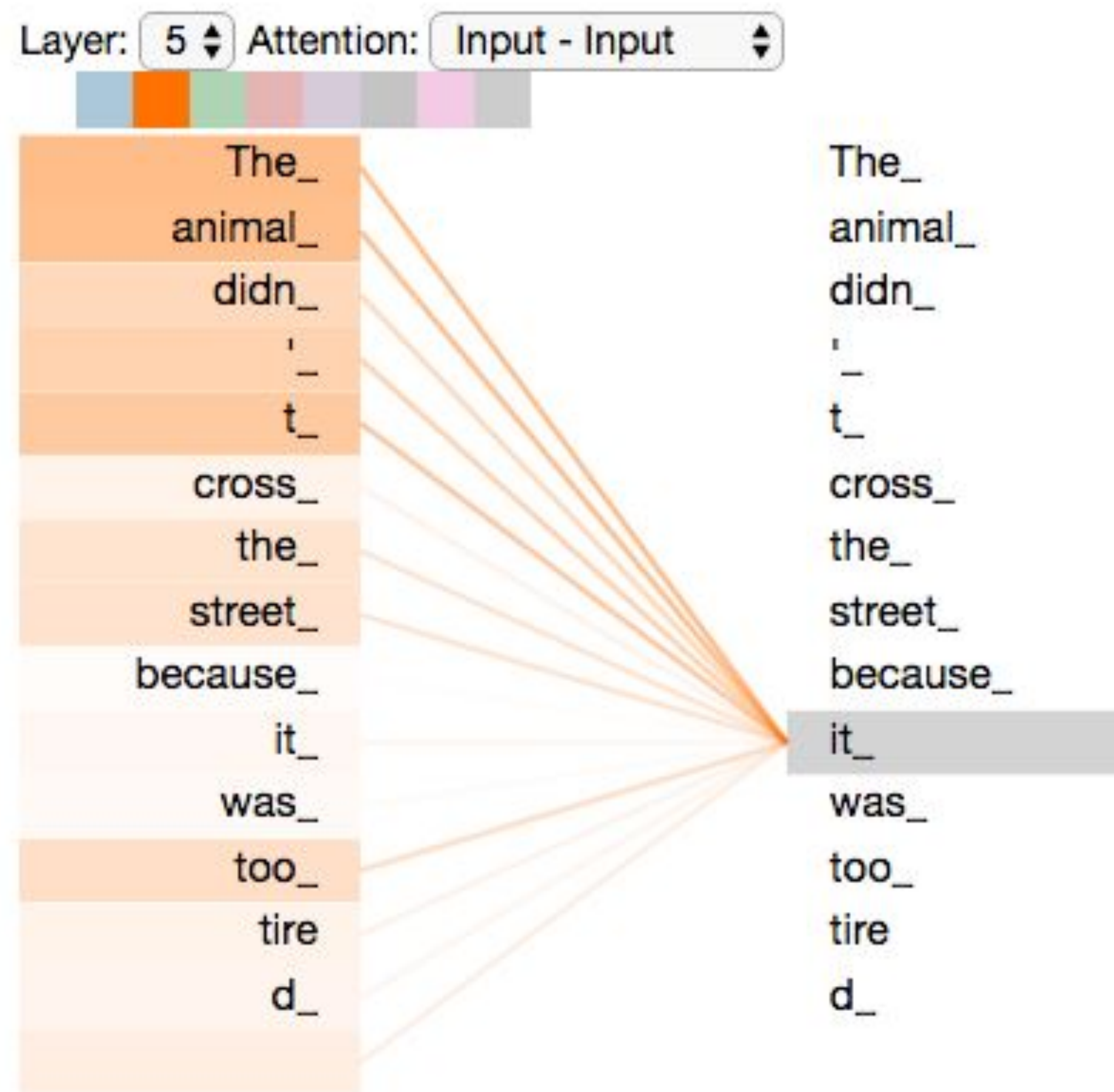
Feeding Word Embeddings into the Encoder



Encoding

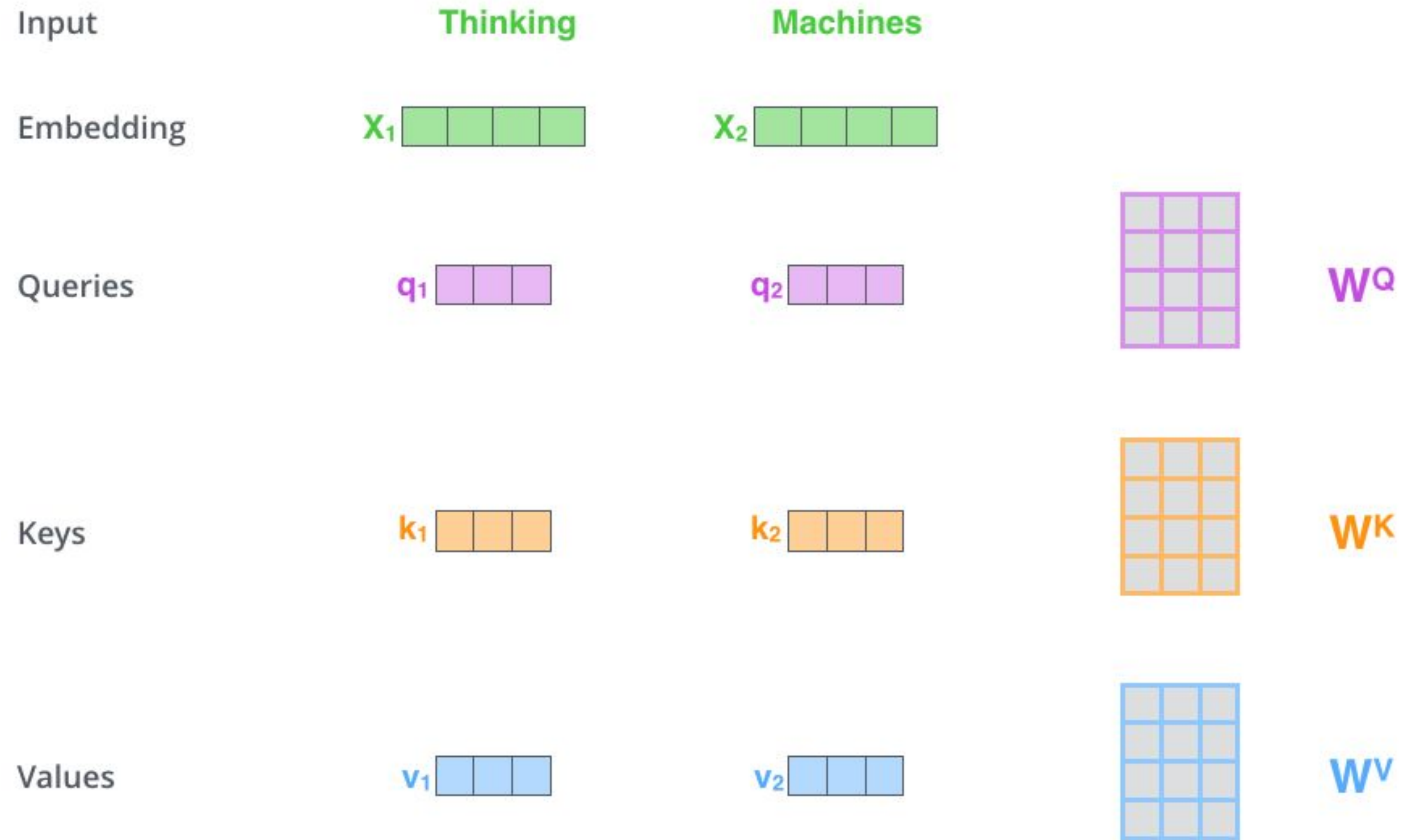


Self-Attention

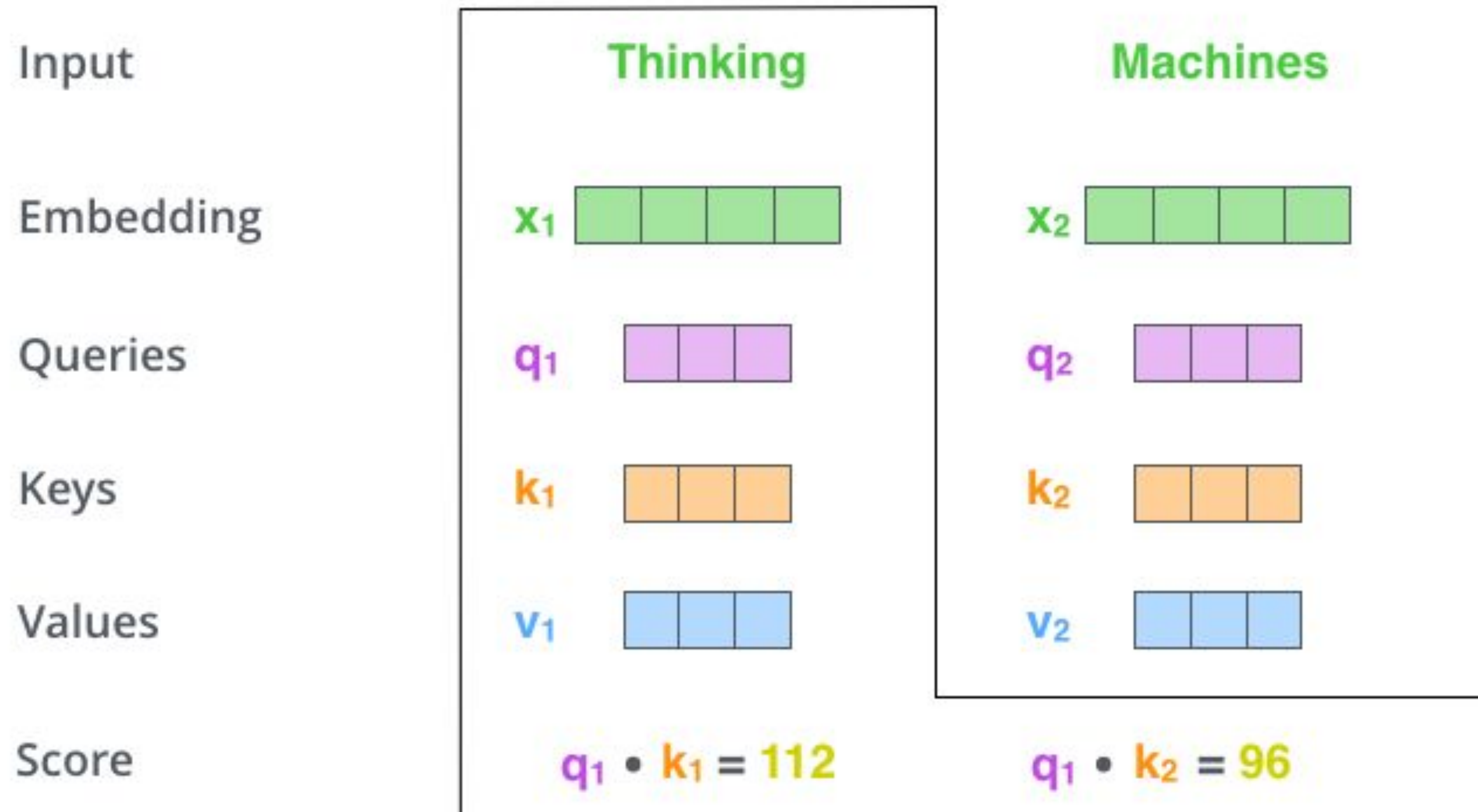


As we are encoding the word "it" in encoder #5 (the top encoder in the stack), part of the attention mechanism was focusing on "The Animal", and baked a part of its representation into the encoding of "it".

Self-Attention in Detail



Self-Attention in Detail



Self-Attention in Detail

Input

Embedding

Queries

Keys

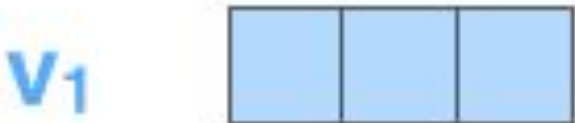
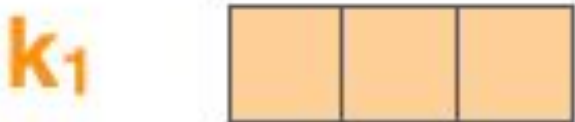
Values

Score

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

Softmax

Thinking

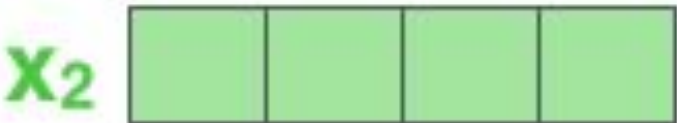


$q_1 \cdot k_1 = 112$

14

0.88

Machines



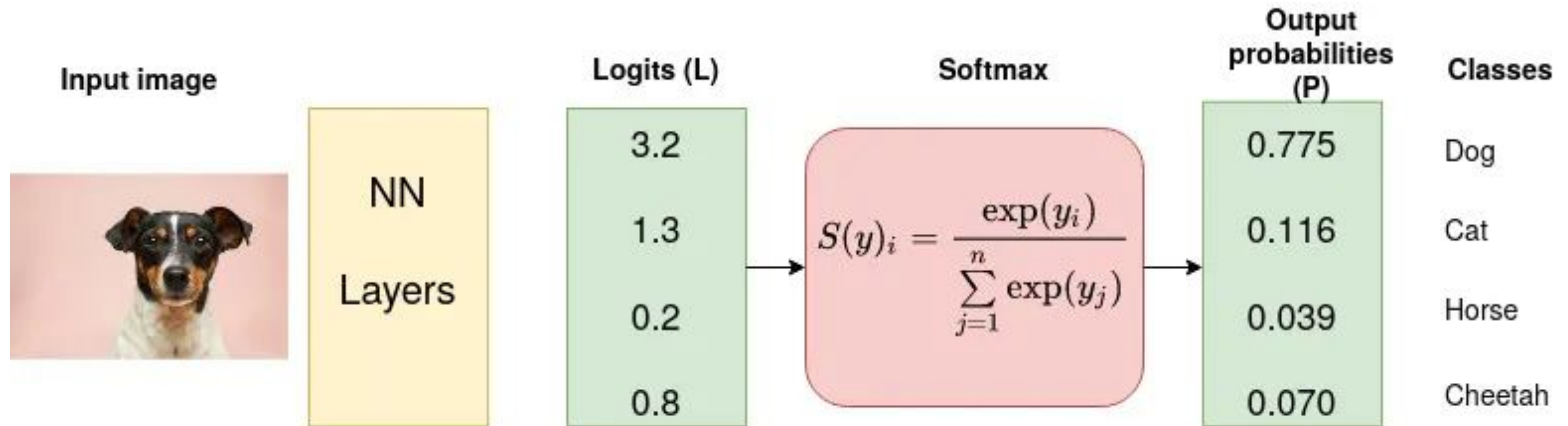
$q_1 \cdot k_2 = 96$

12

0.12

Divide the scores by 8 (the square root of the dimension of the key vectors used in the paper – 64). This leads to having more stable gradients.

Softmax



Input image source: Photo by [Victor Grabarczyk](#) on [Unsplash](#) . Diagram by author.

Self-Attention

Input

Embedding

Queries

Keys

Values

Score

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

Softmax

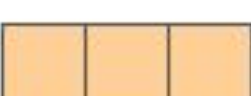
Softmax
X
Value

Sum

Thinking

x_1 

q_1 

k_1 

v_1 

$q_1 \cdot k_1 = 112$

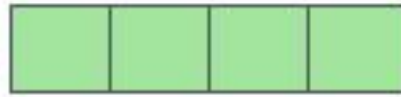
14

0.88

v_1 

z_1 

Machines

x_2 

q_2 

k_2 

v_2 

$q_1 \cdot k_2 = 96$

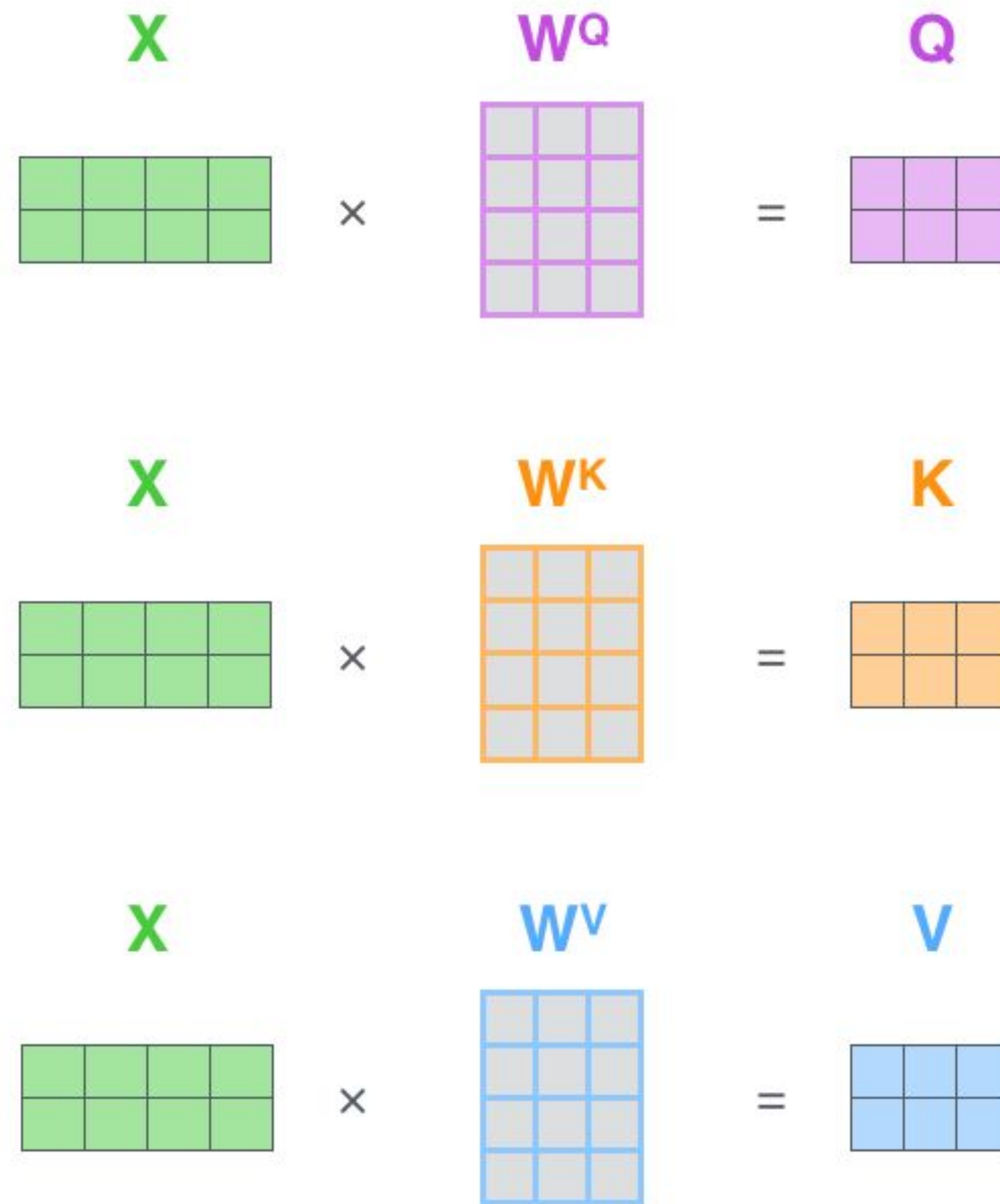
12

0.12

v_2 

z_2 

Matrix Calculation of Self-Attention



Matrix Calculation of Self-Attention

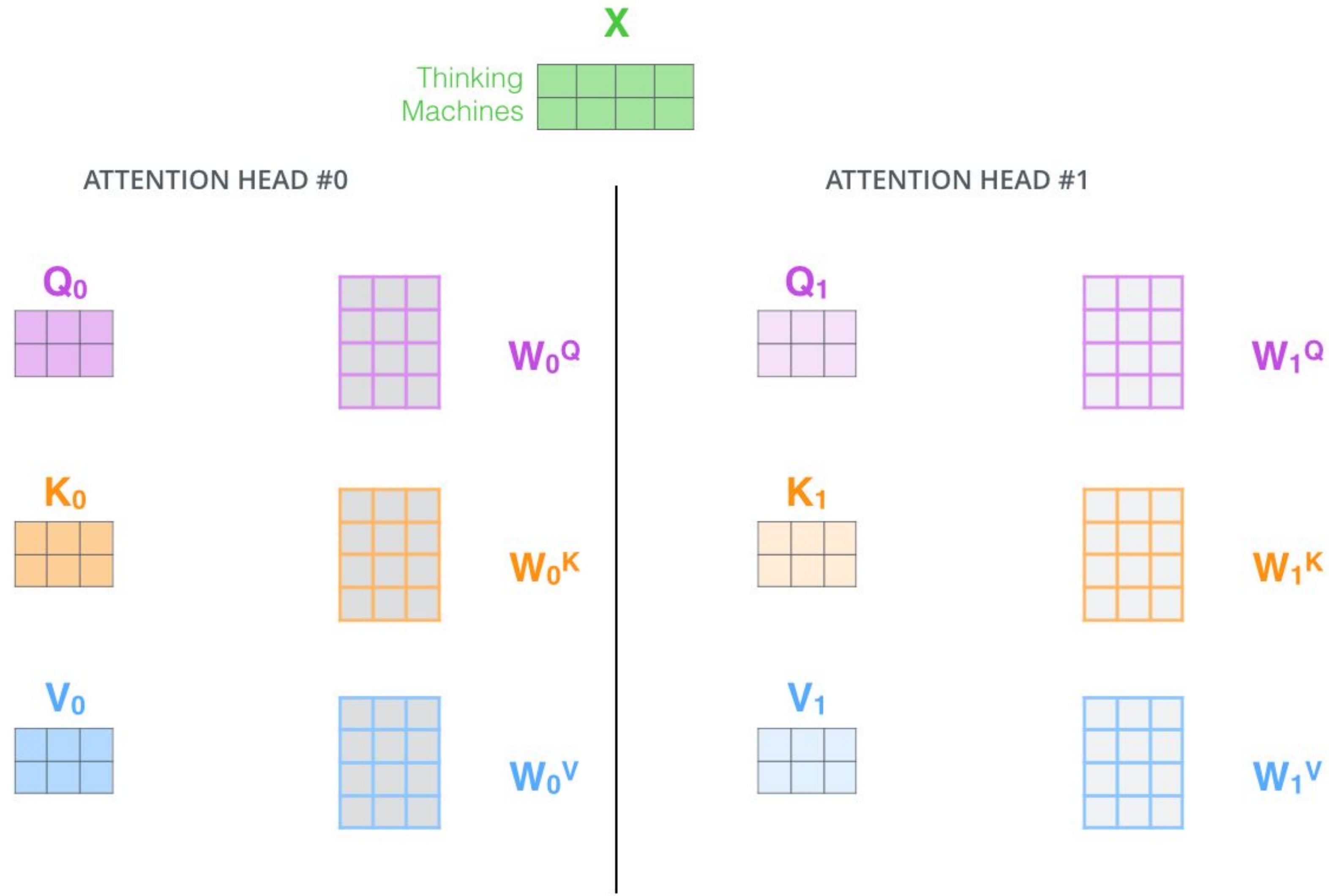
$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{array}{|c|c|} \hline \square & \square \\ \hline \square & \square \\ \hline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

=

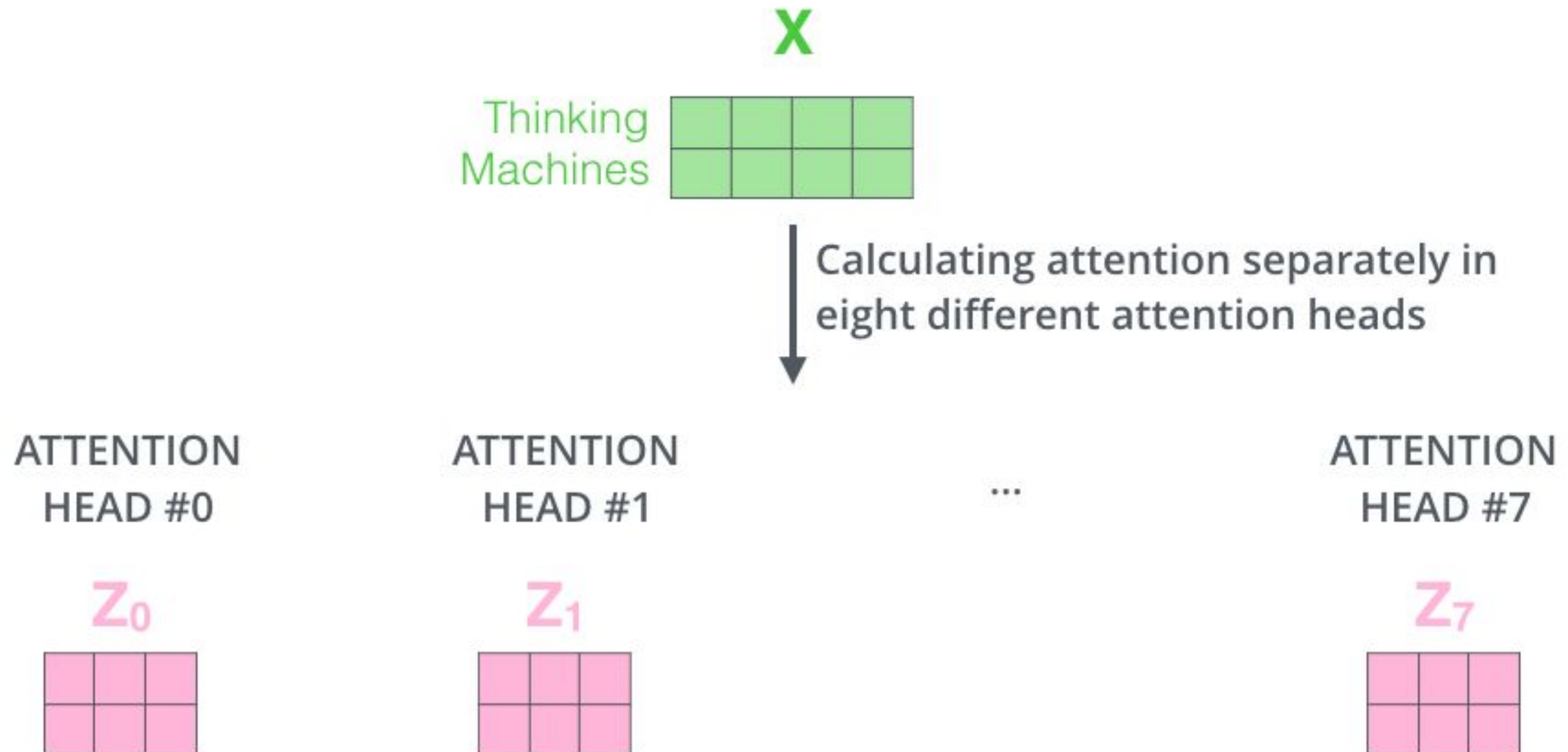
Z

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

Multi-Headed Self-Attention



Multi-Headed Self-Attention



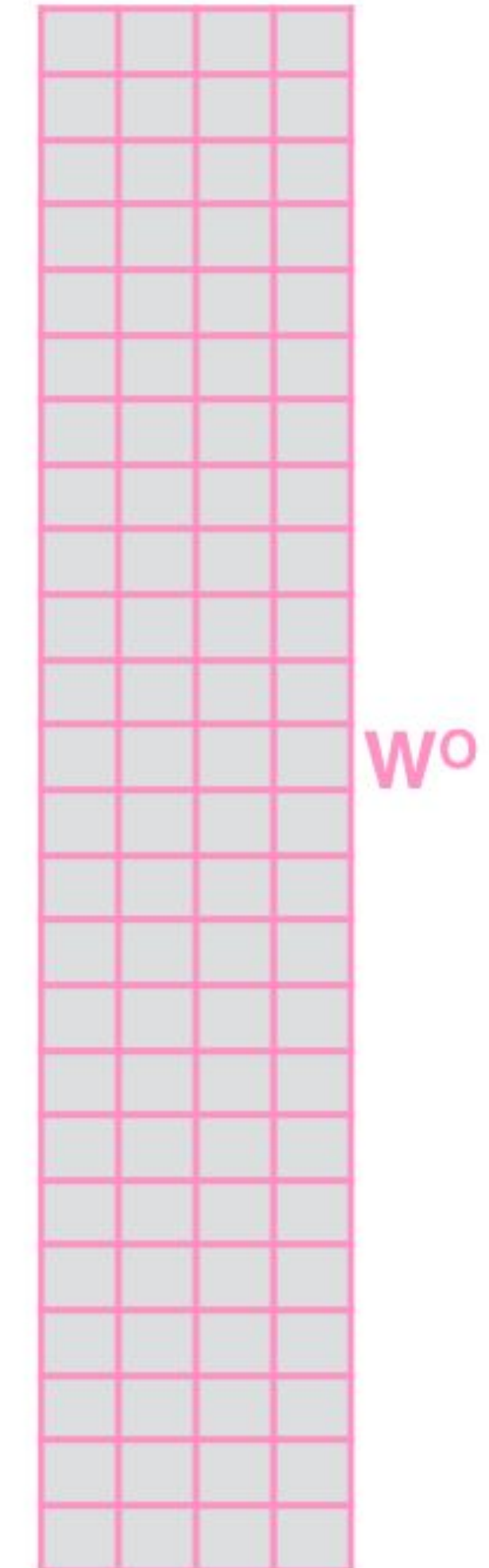
Multi-Headed Self-Attention

1) Concatenate all the attention heads

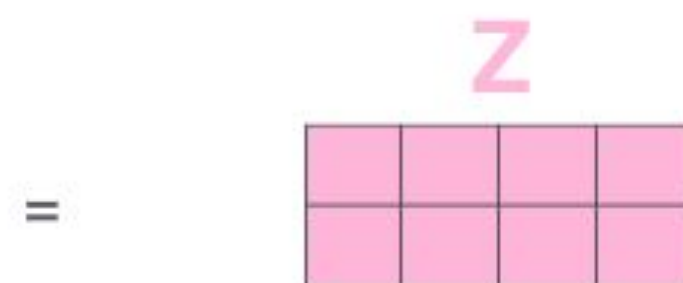


2) Multiply with a weight matrix W^O that was trained jointly with the model

\times



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

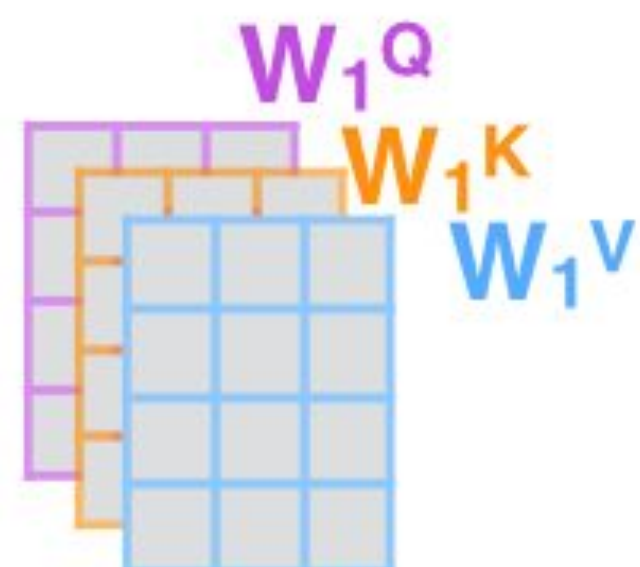
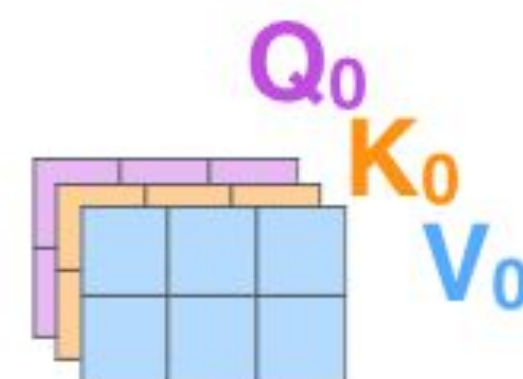
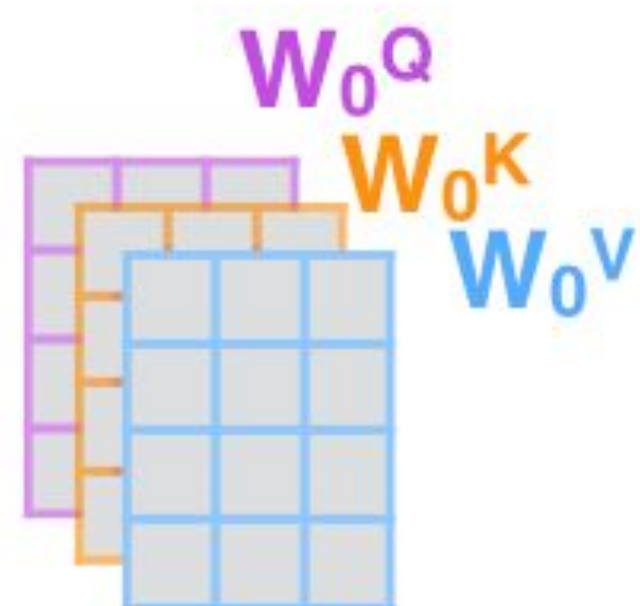
4) Calculate attention using the resulting $Q/K/V$ matrices

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

Thinking
Machines



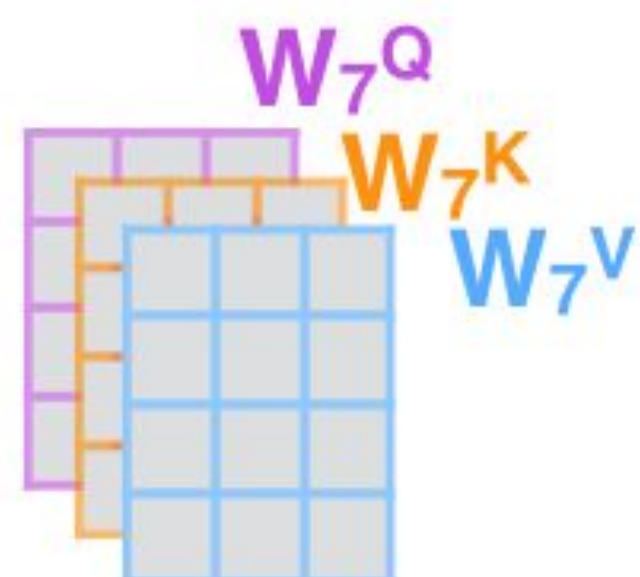
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



...

...

...

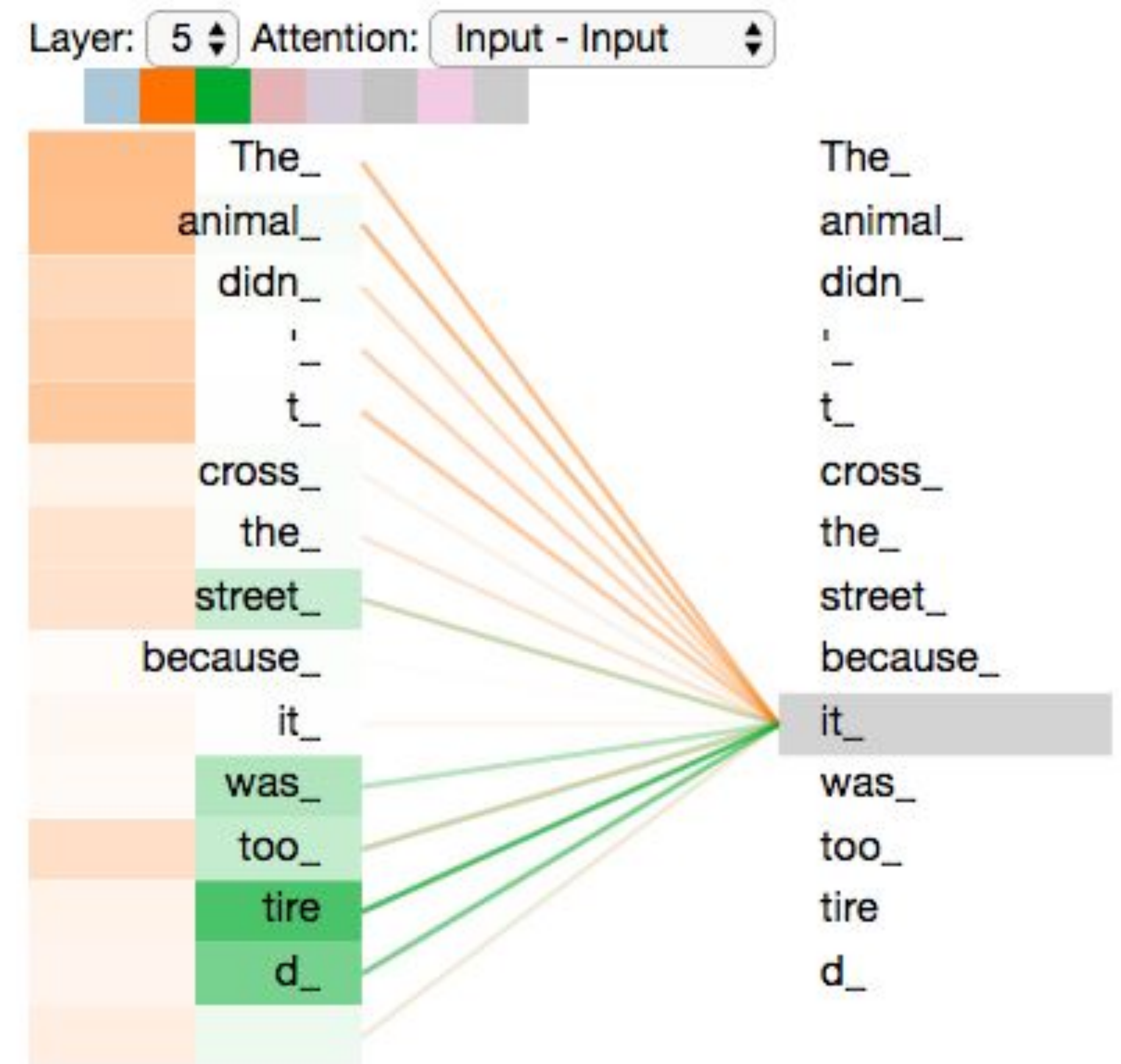


W^O

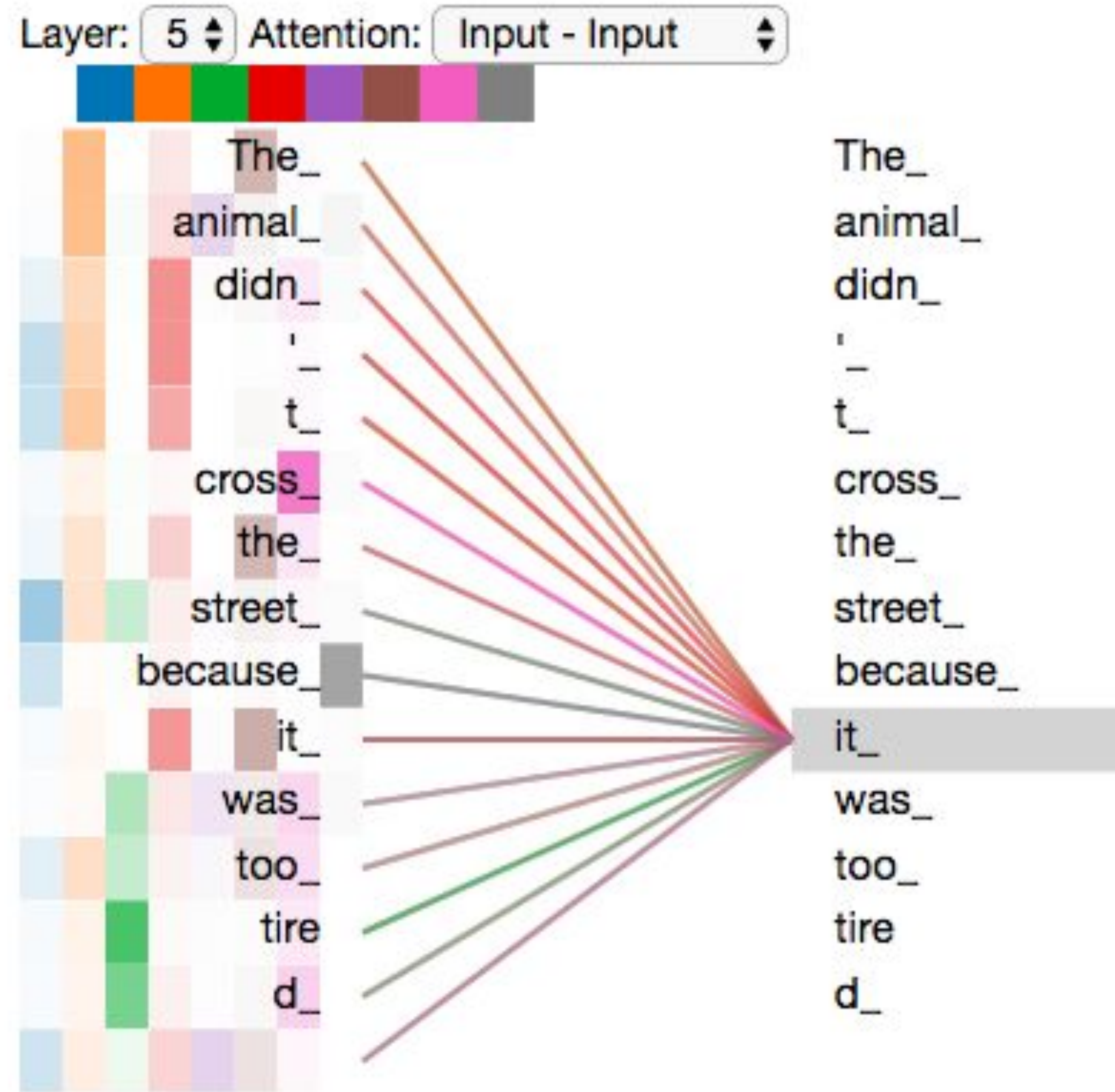


Multi-Headed Self-Attention

As we encode the word "it", one attention head is focusing most on "the animal", while another is focusing on "tired" -- in a sense, the model's representation of the word "it" bakes in some of the representation of both "animal" and "tired".



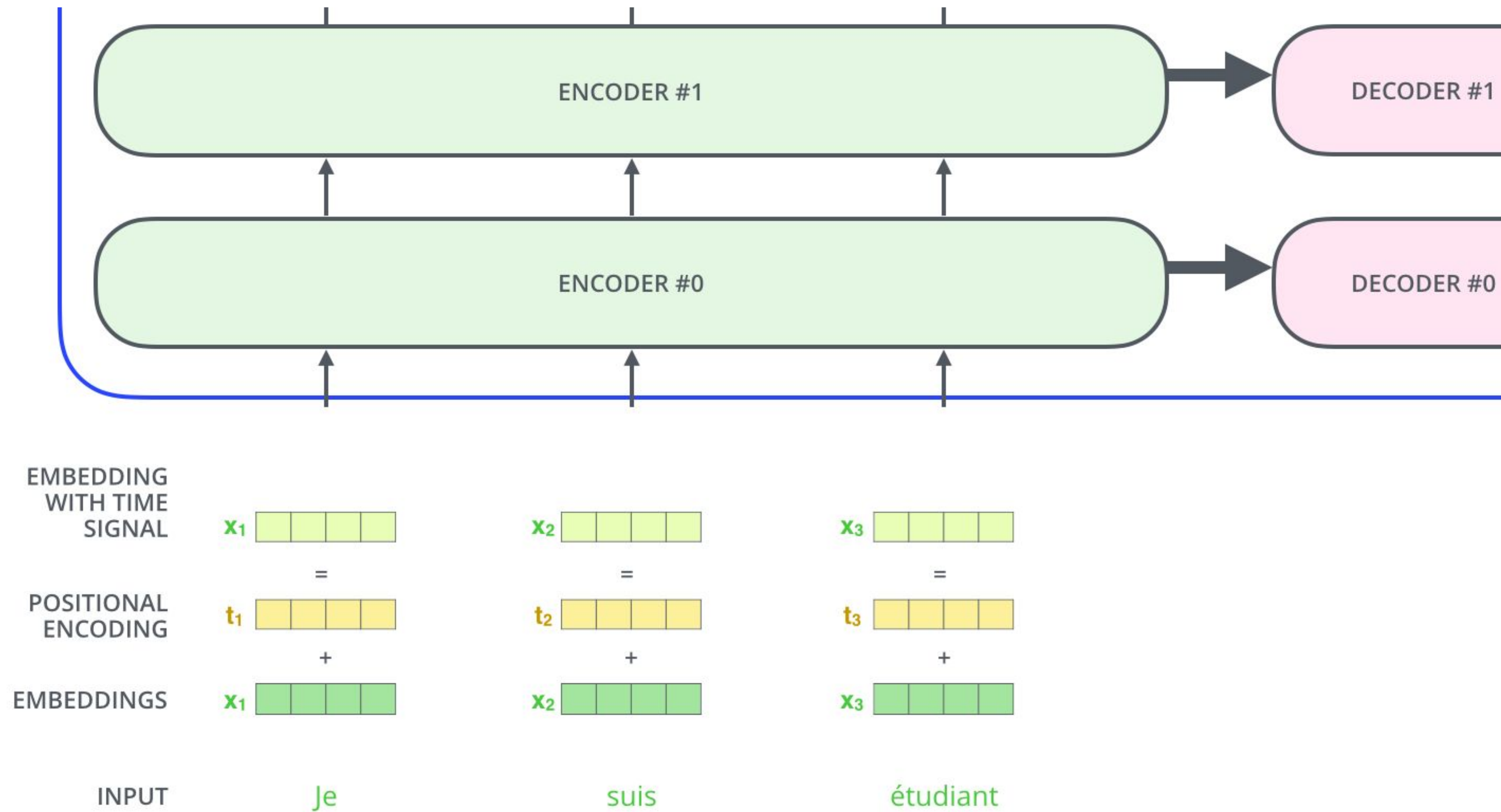
Multi-Headed Self-Attention



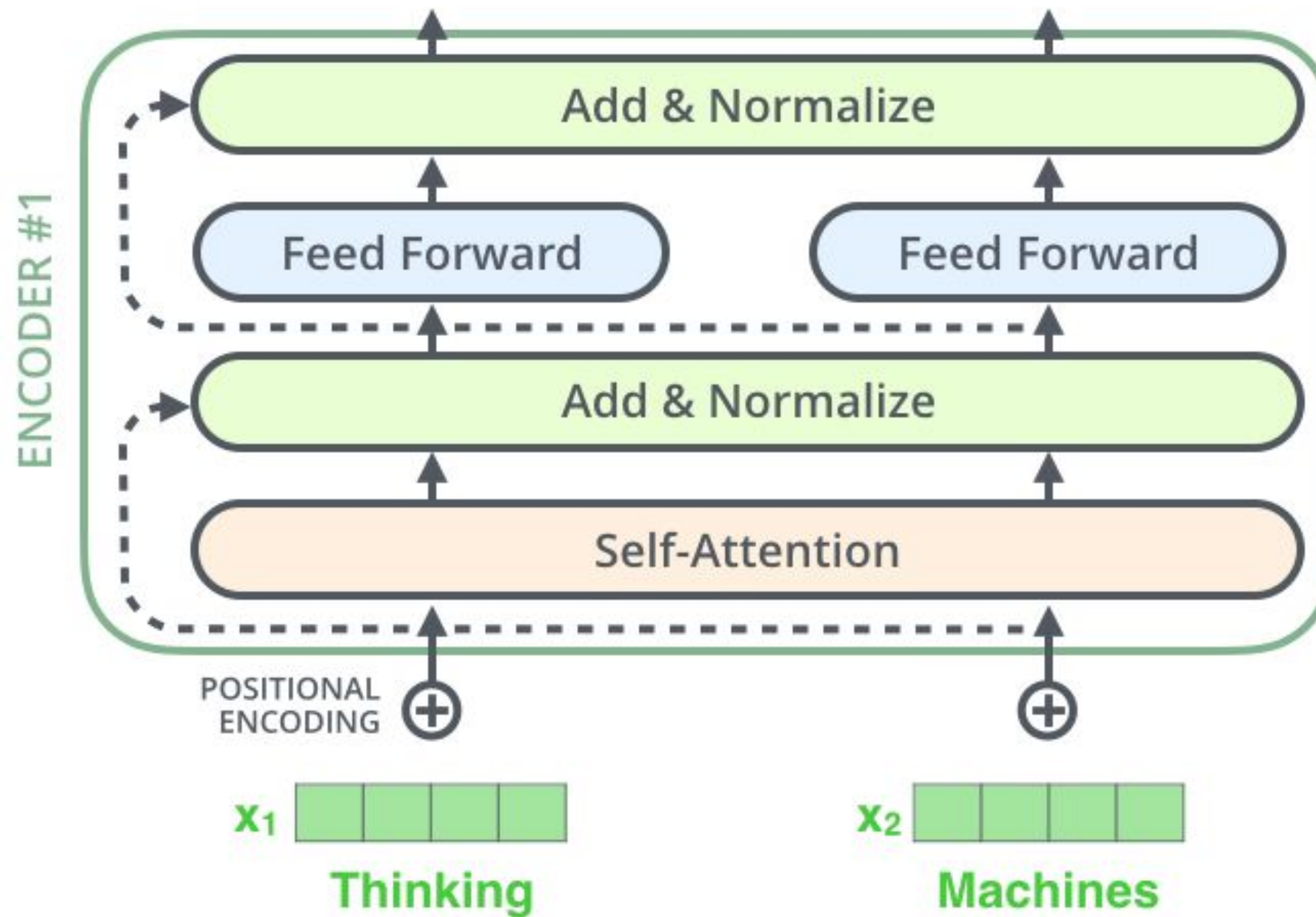


In considering the importance of the positional information of elements in sequence data, how can the Transformer model effectively utilize this positional information in its processing?

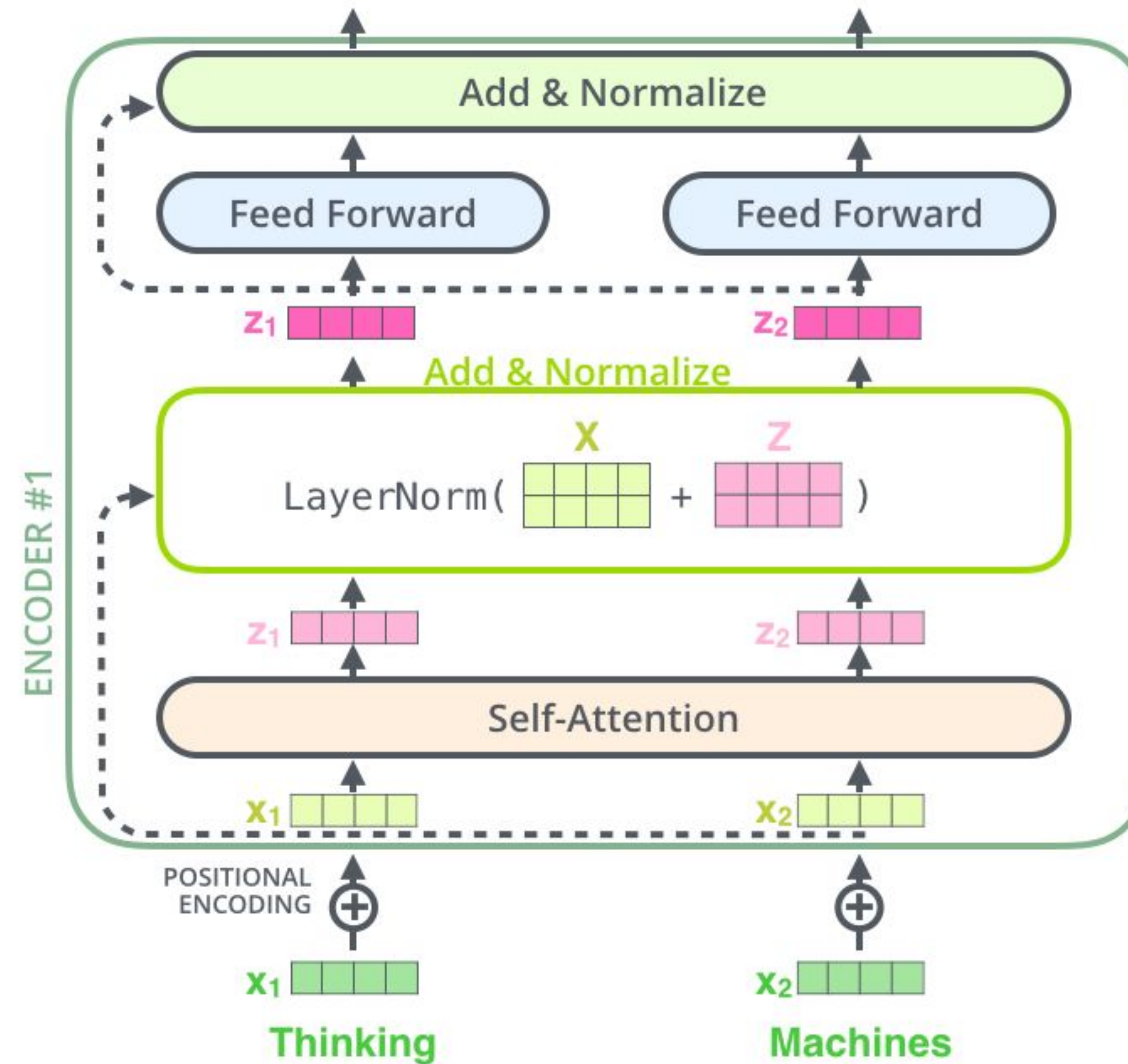
Positional Encoding



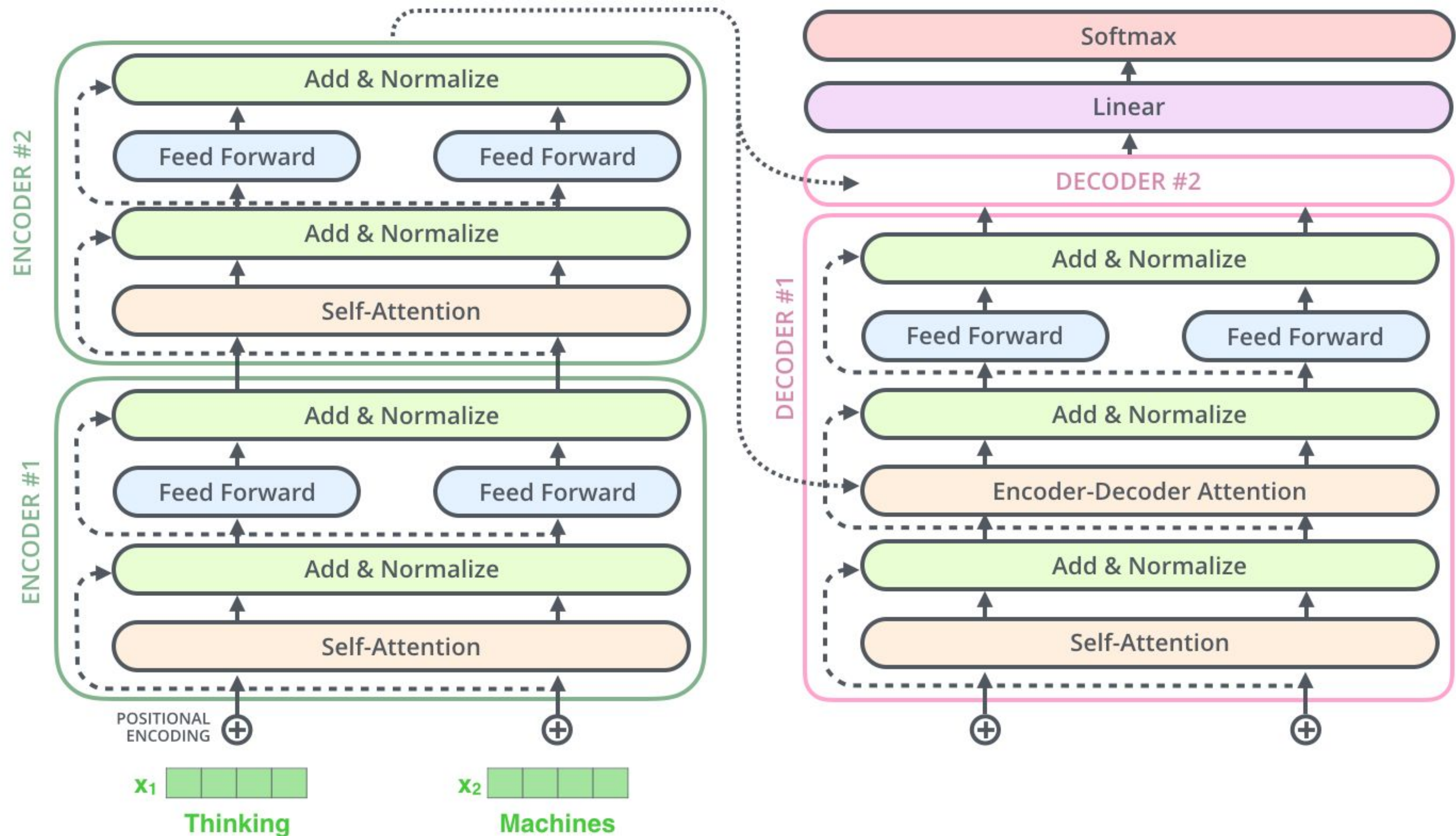
Add Residual Connections



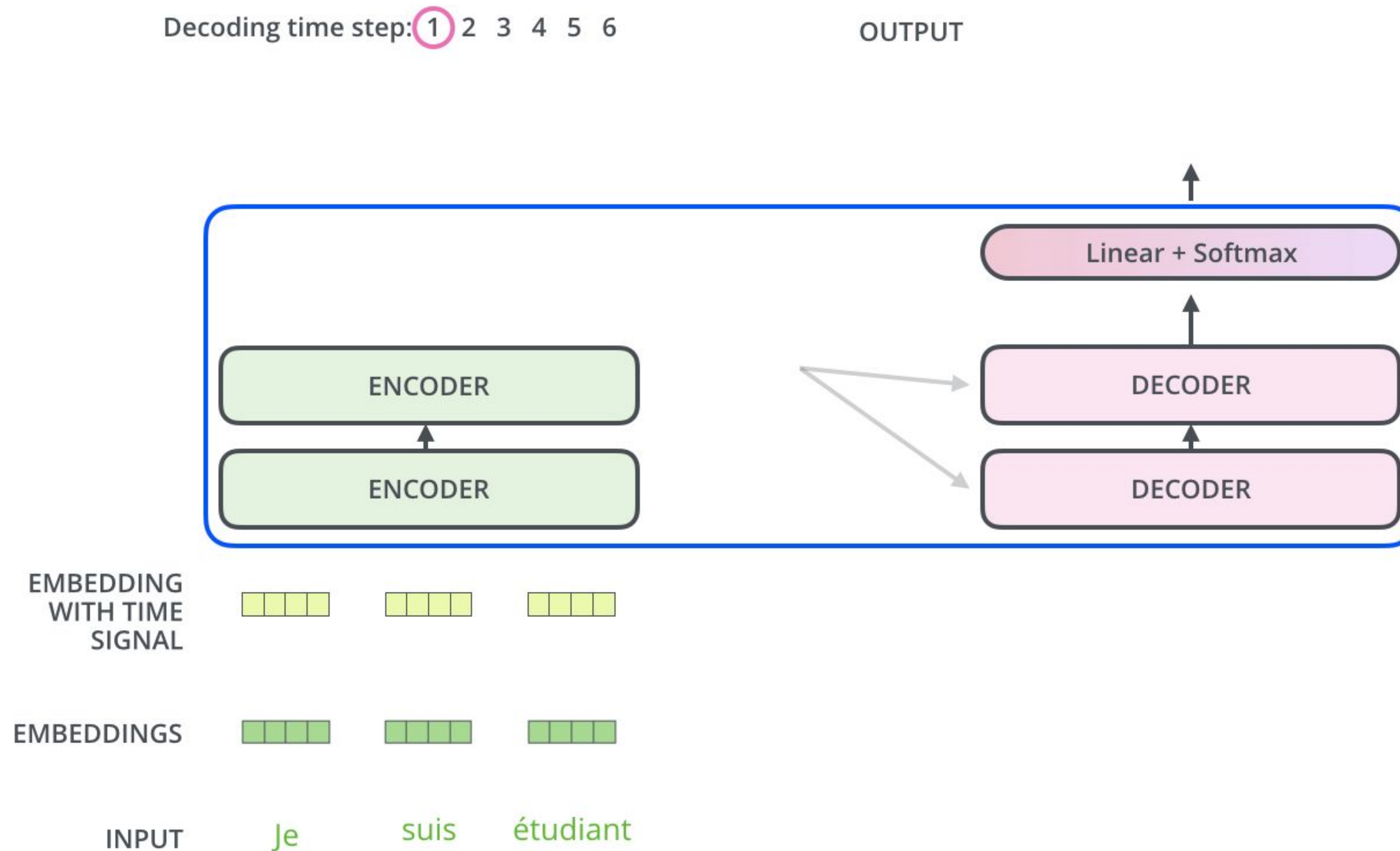
Layer Normalization



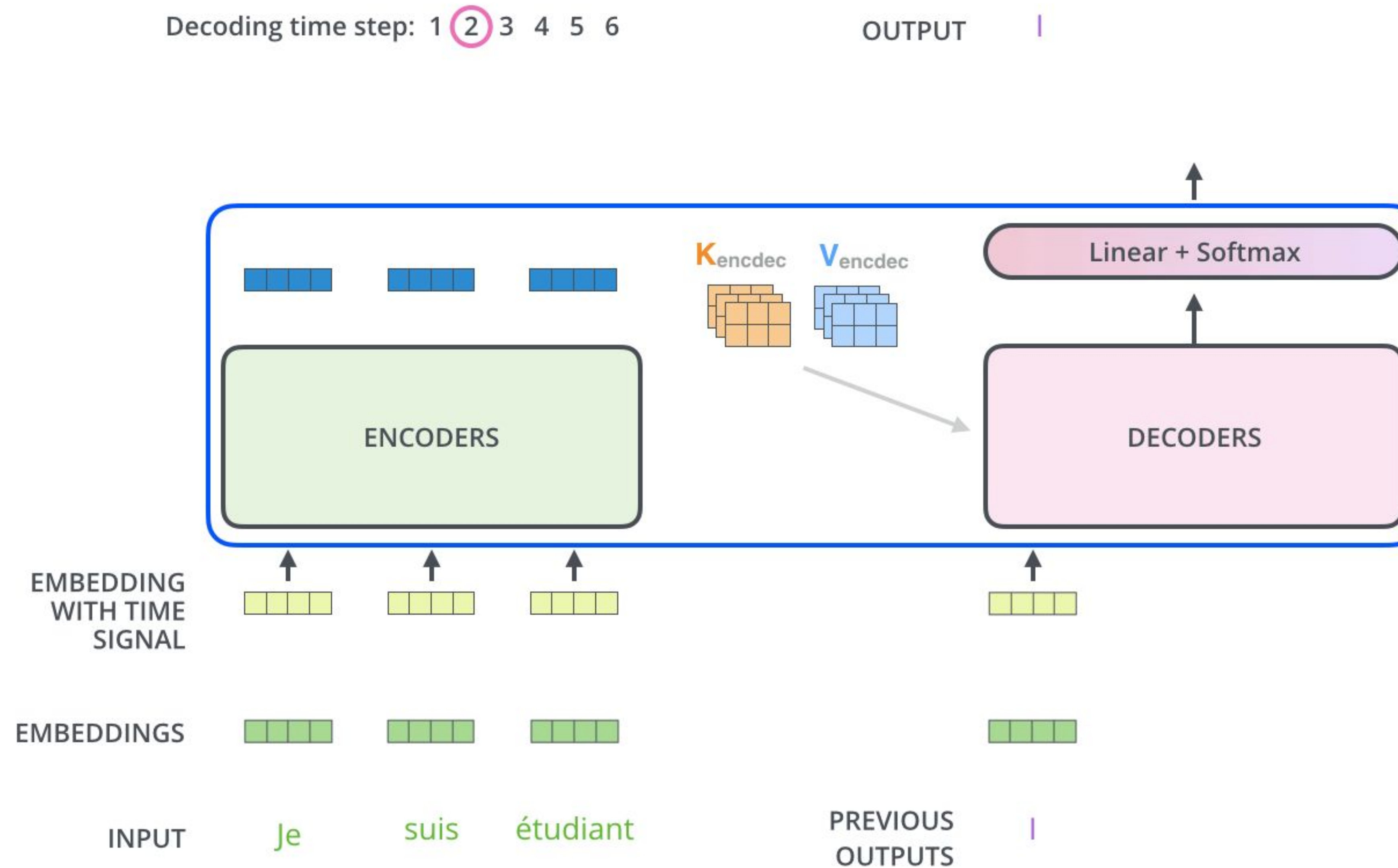
Encoder-Decoder



Encoder-Decoder



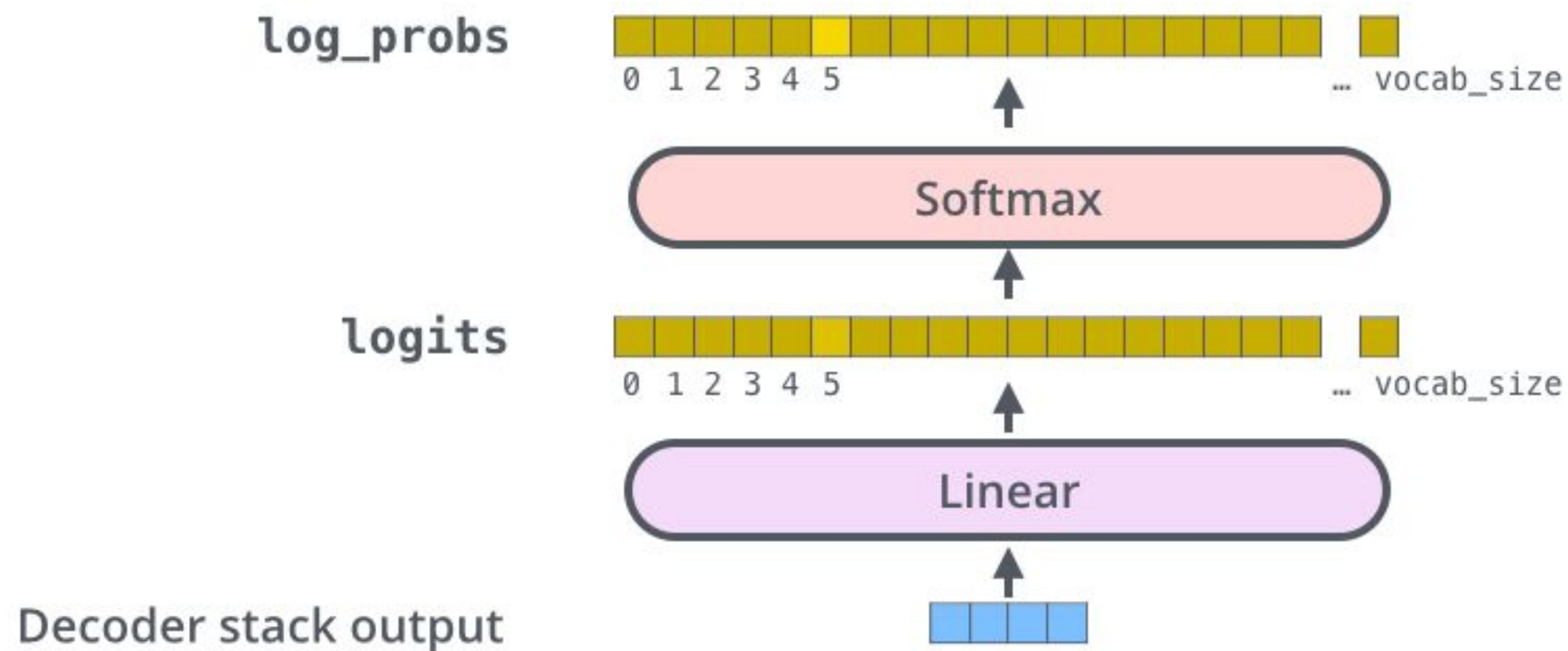
Encoder-Decoder



The Final Linear and Softmax Layer

Which word in our vocabulary
is associated with this index?

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



Example

Output Vocabulary

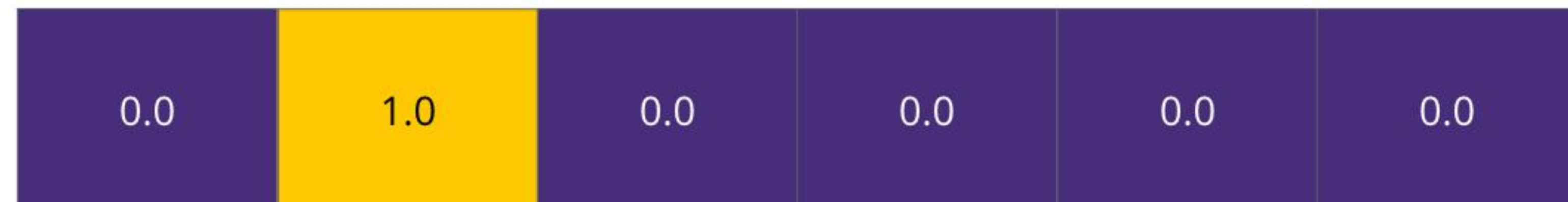
WORD	a	am	I	thanks	student	<eos>
INDEX	0	1	2	3	4	5

Example

Output Vocabulary

WORD	a	am	I	thanks	student	<eos>
INDEX	0	1	2	3	4	5

One-hot encoding of the word "am"



Example

Untrained Model Output



Correct and desired output



a

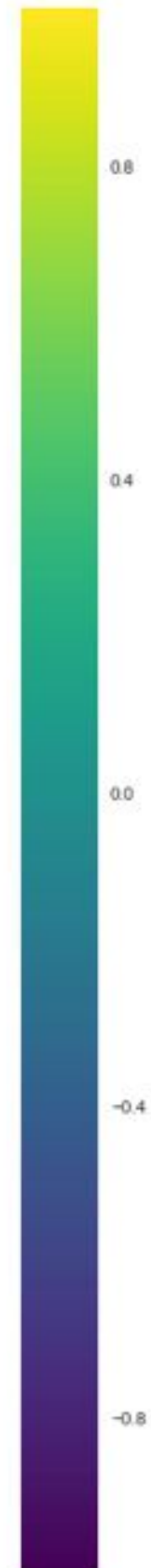
am

I

thanks

student

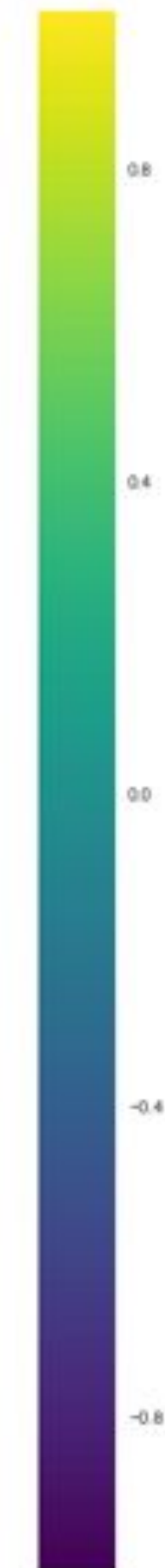
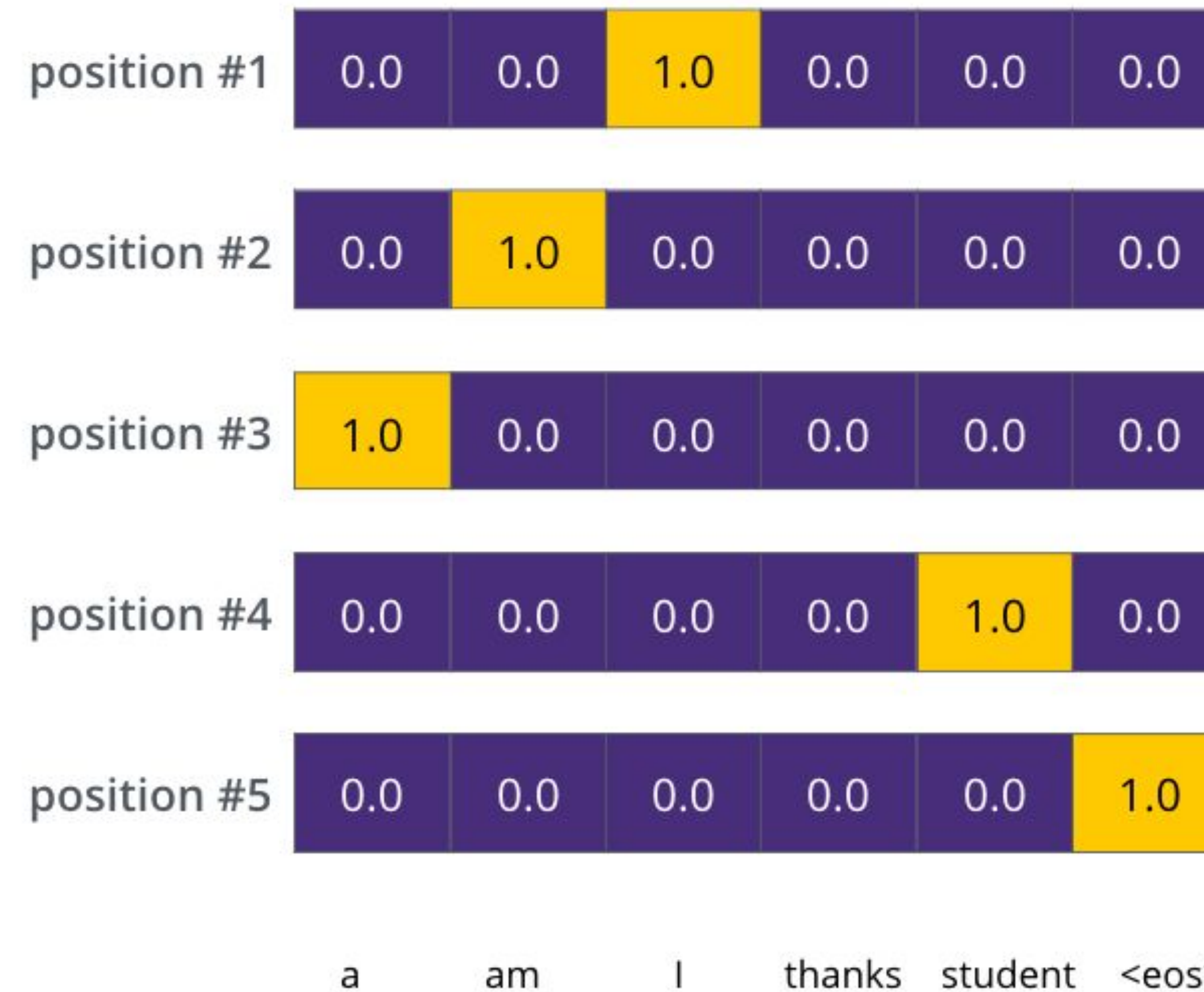
<eos>



Example

Target Model Outputs

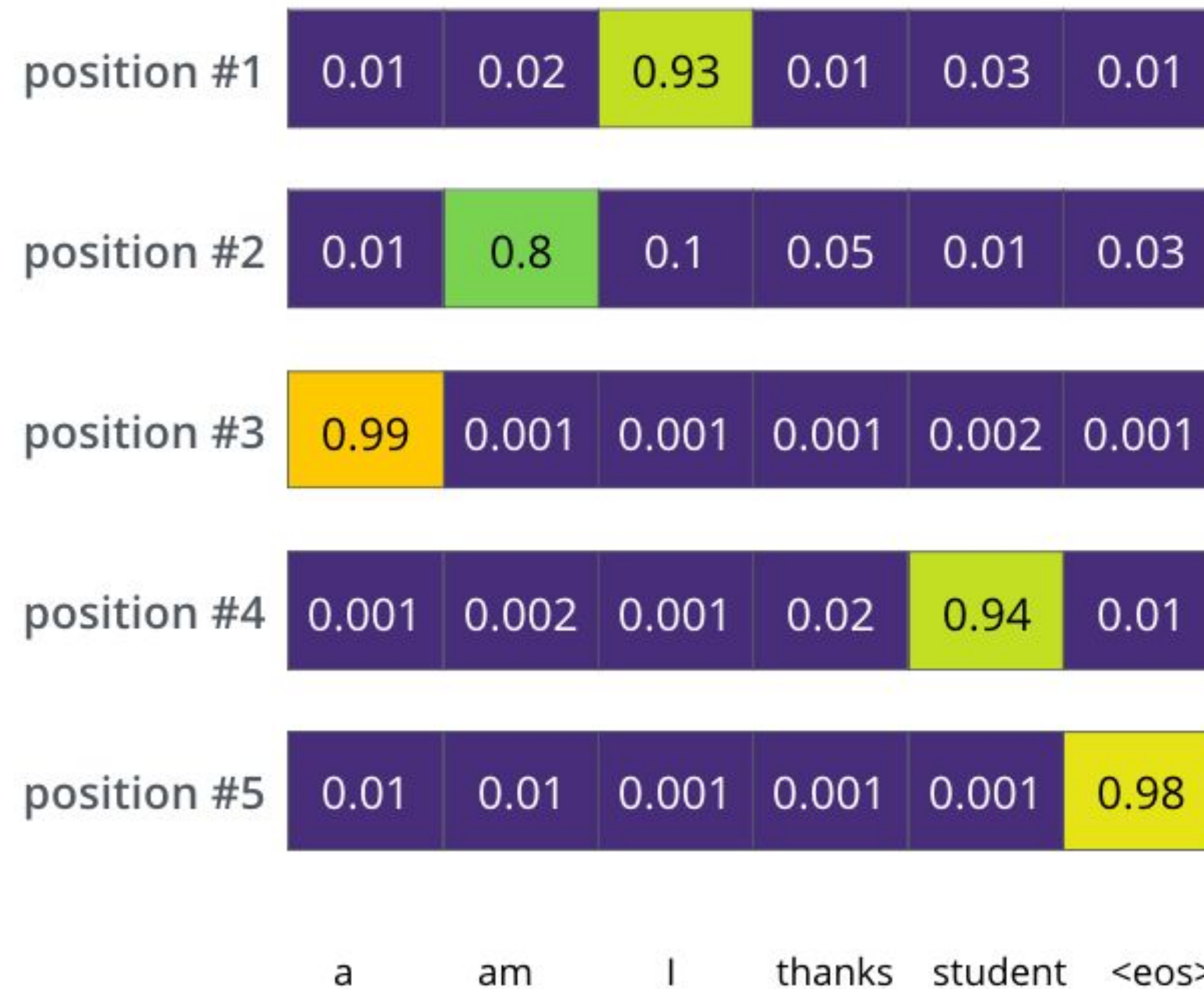
Output Vocabulary: a am I thanks student <eos>



Example

Trained Model Outputs

Output Vocabulary: a am I thanks student <eos>



References

- [Understanding Long Short-Term Memory Networks \(LSTMs\) - Chris Olah](#)
- [The Illustrated Transformer - Jay Alammar](#)
- [Attention is all you need](#)

Summary

- Transformers
 - Self-Attention, Softmax,
 - Multi-Headed
 - Self-Attention,
 - Positional Encoding,
 - Residual Connection,
 - Layer Normalization,
 - Encoder-Decoder