

Машинное обучение

Лекция 10

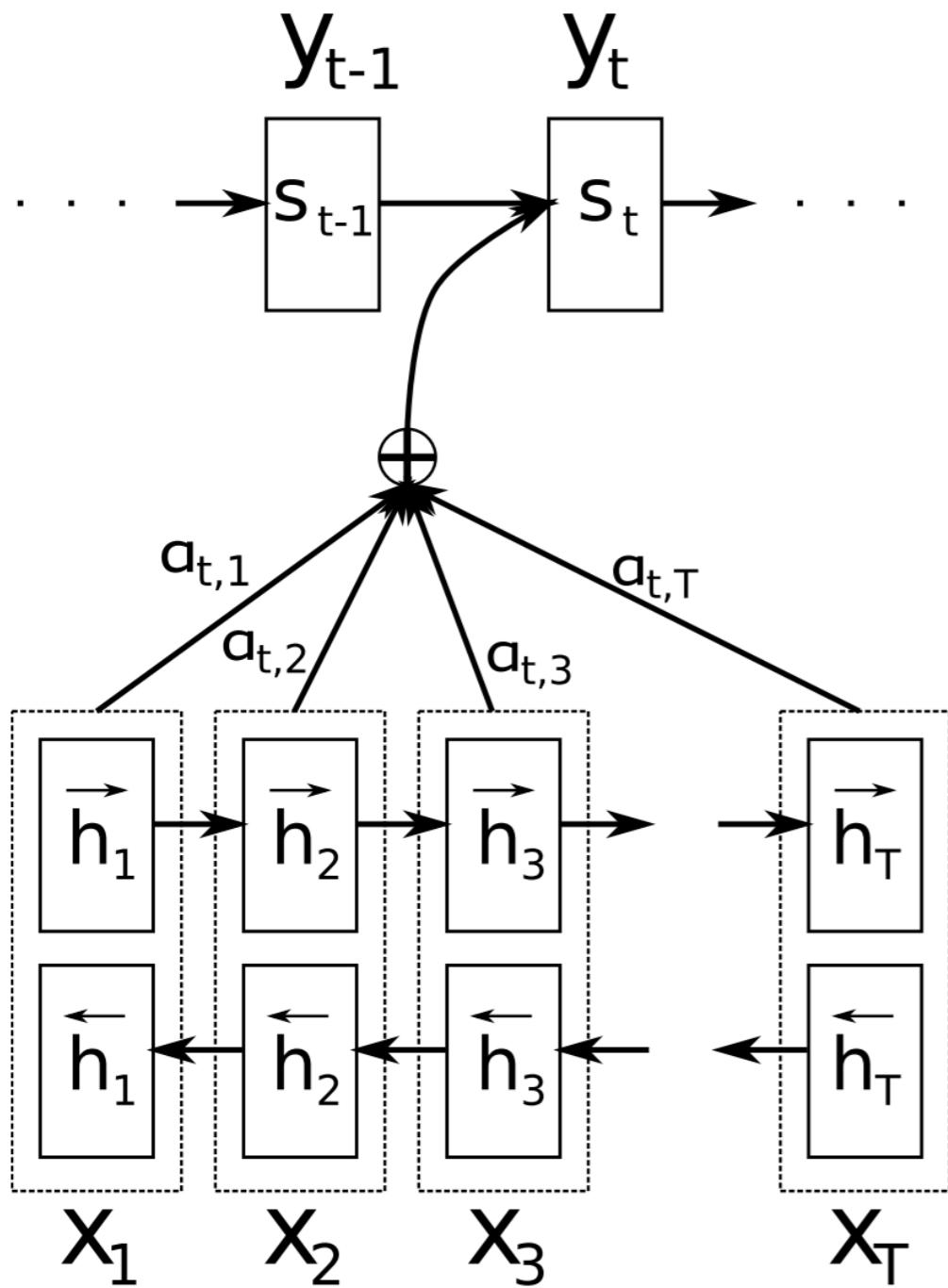
Seq2seq learning, Attention

Власов Кирилл Вячеславович



2019

Recap: Attention



Скрытый слой RNN

$$s_i = f(s_{i-1}, y_{i-1}, c_i).$$

Взвешенная сумма

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j.$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

$$e_{ij} = a(s_{i-1}, h_j)$$

А что если без RNN?

А что если без RNN?

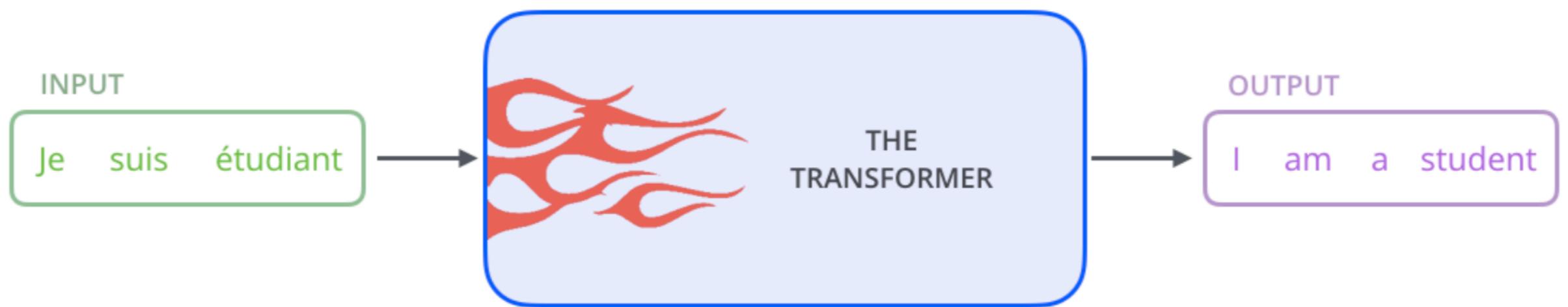
Но ведь нам нужно обрабатывать последовательности...



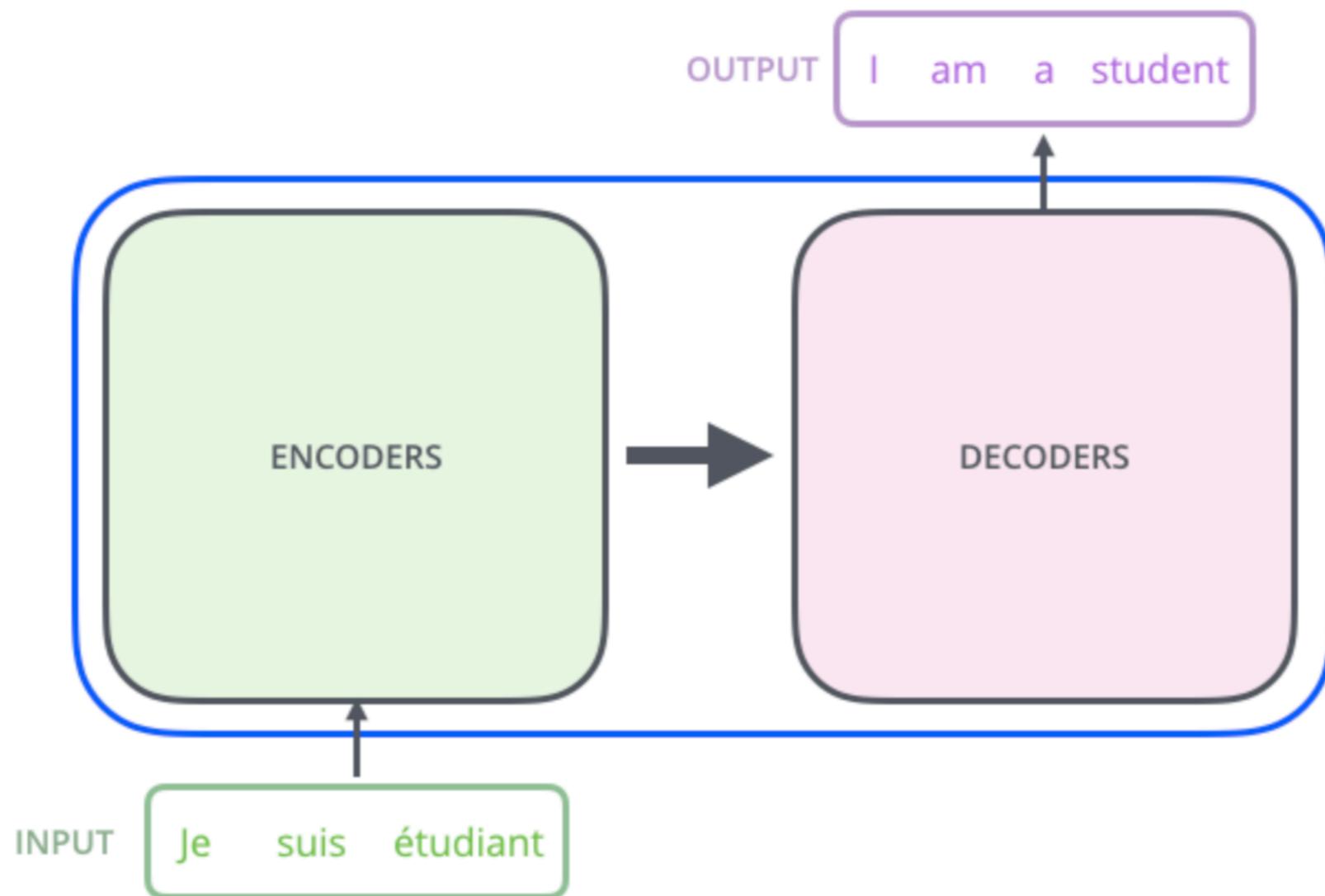
Self-Attention

The Illustrated Transformer

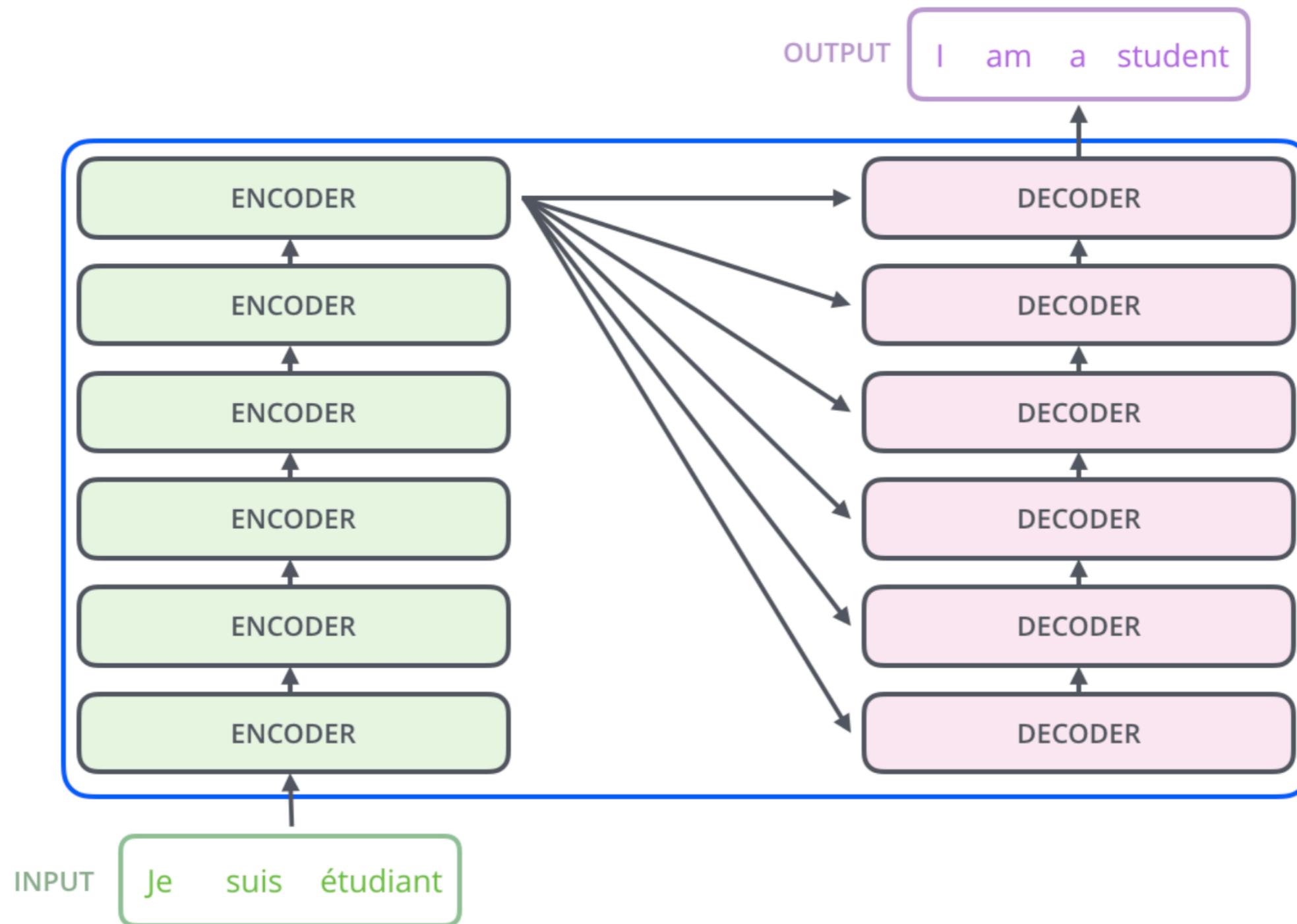
Self-Attention



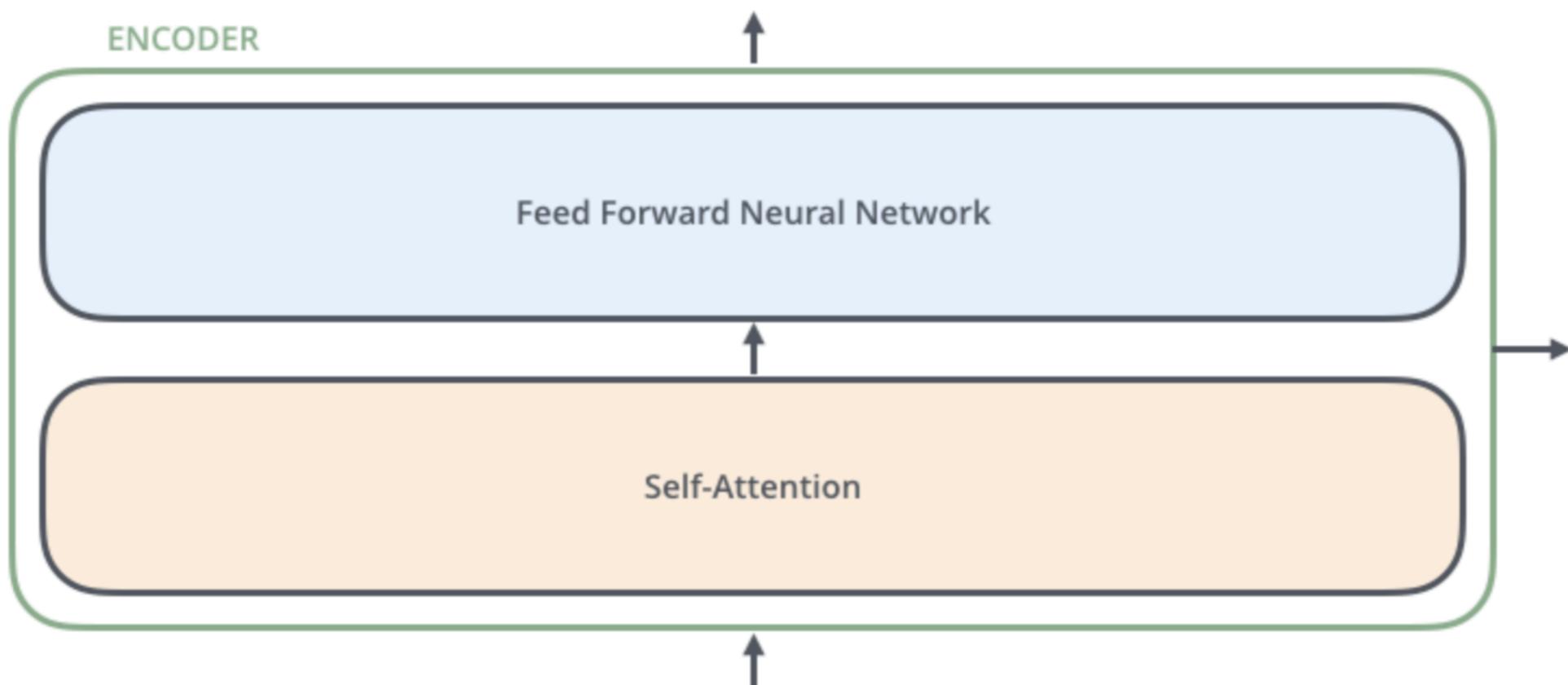
Self-Attention



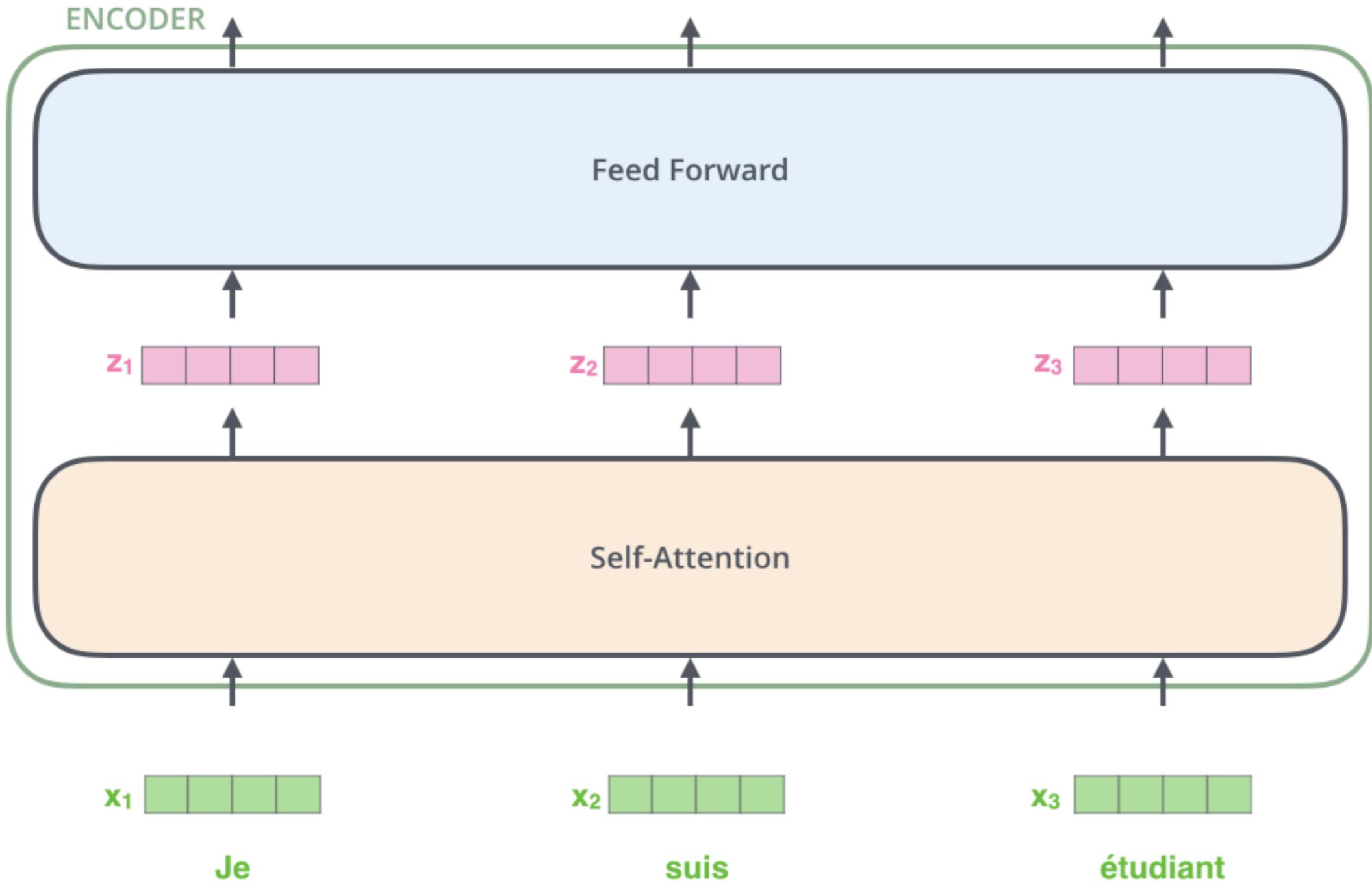
Self-Attention



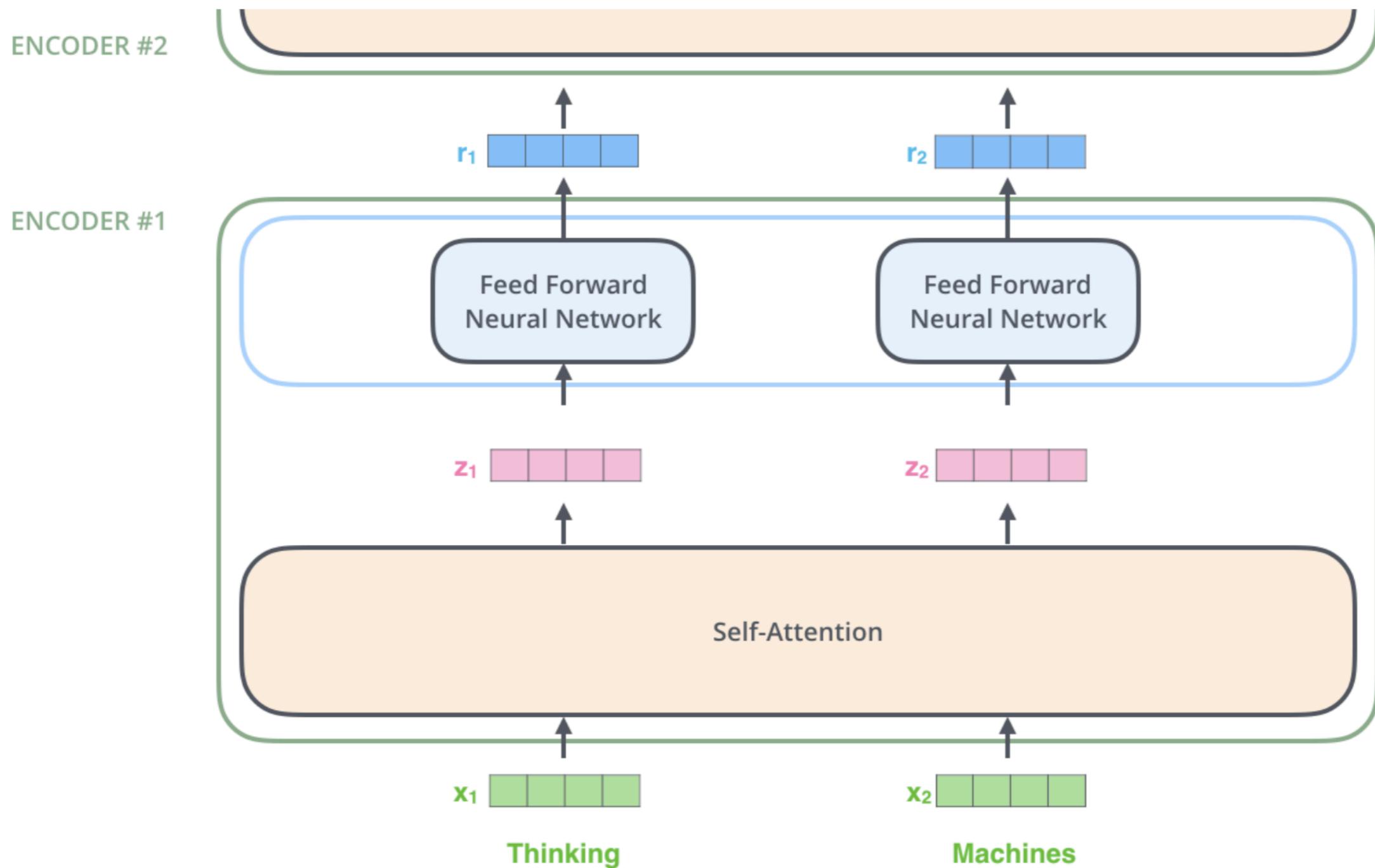
Self-Attention



Self-Attention

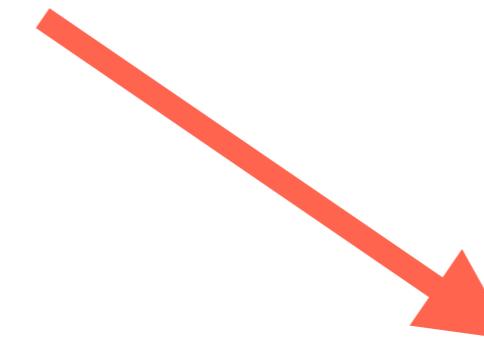


Self-Attention



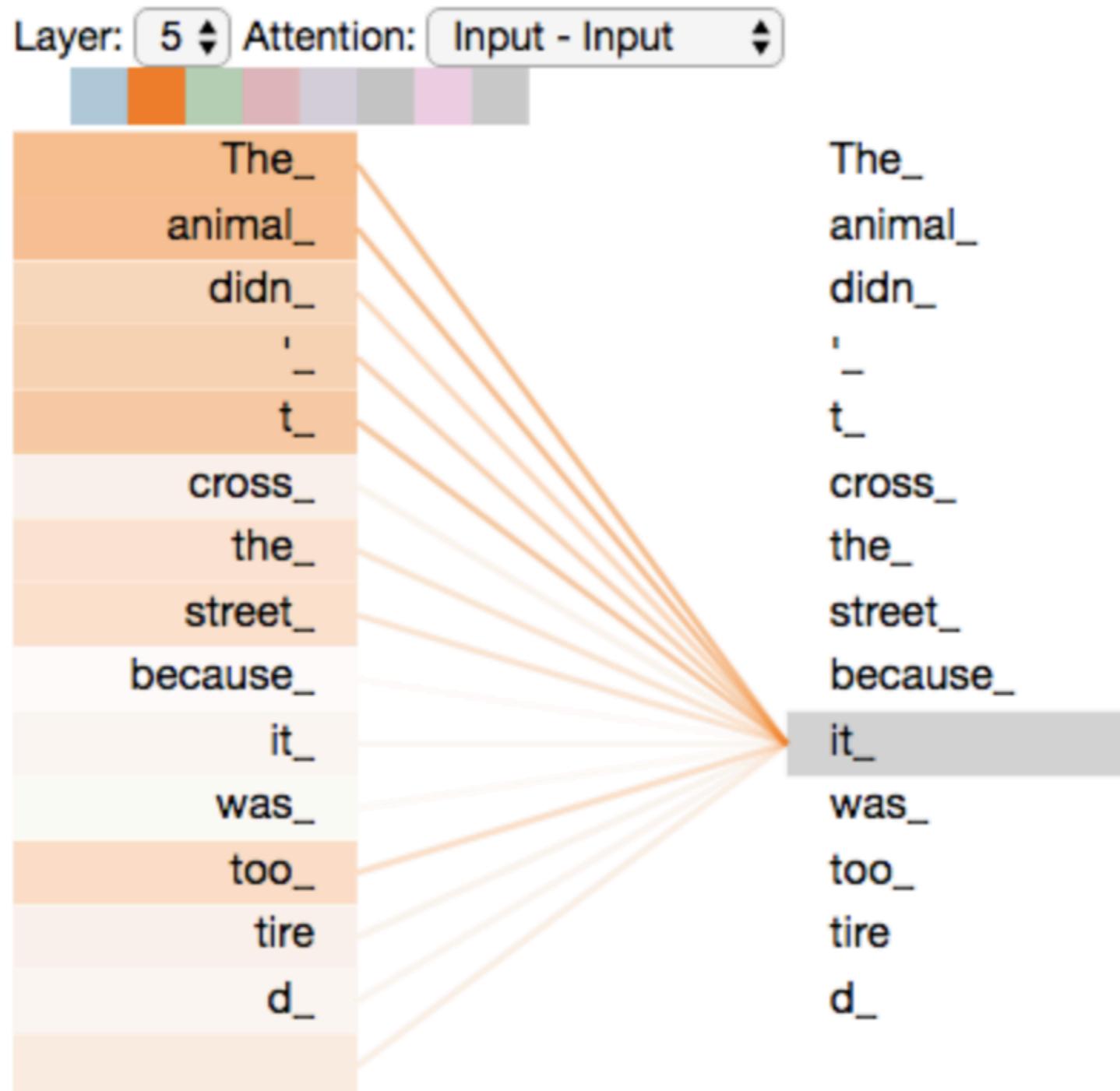
The animal didn't cross the street because it was too tired

?

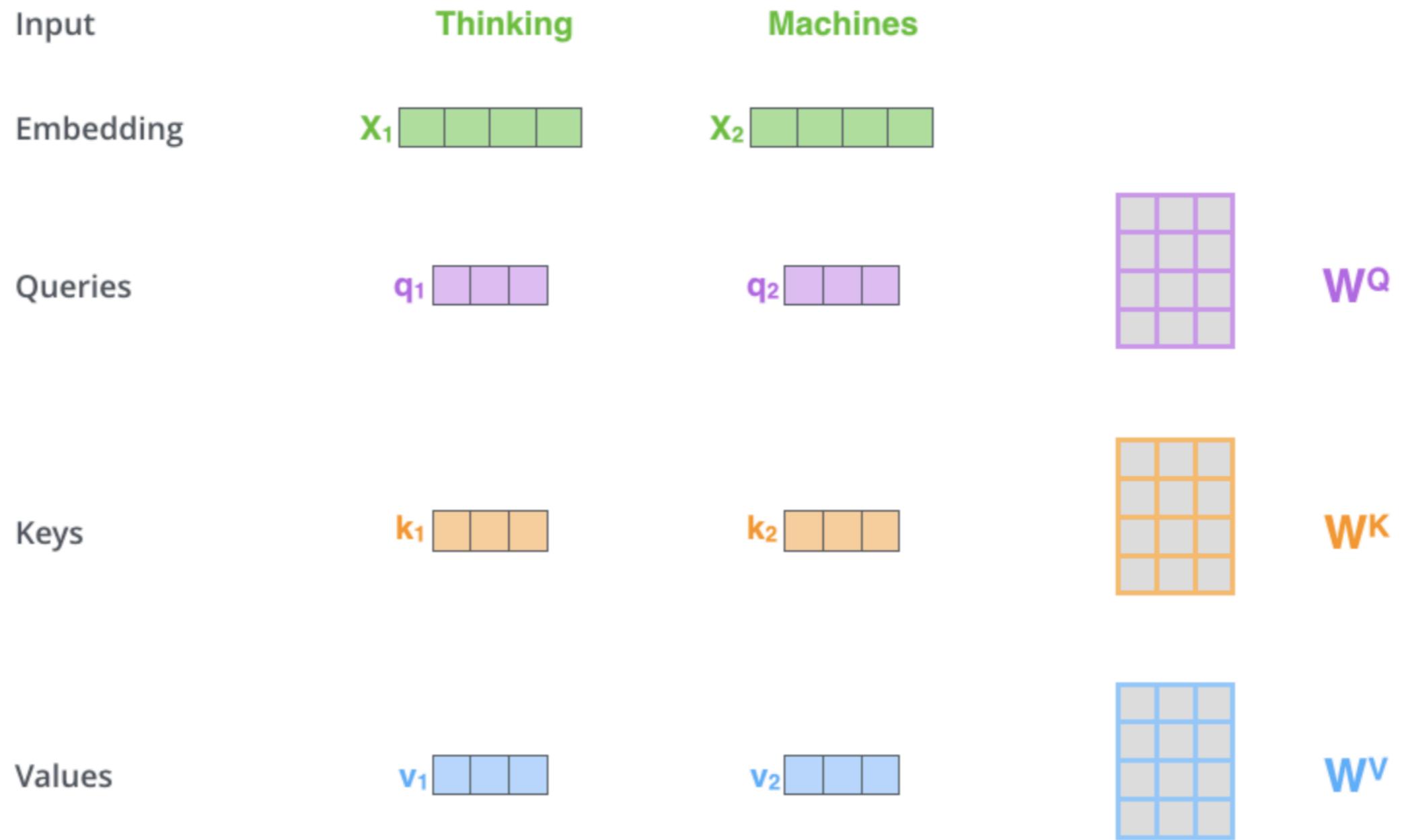


*The animal didn't cross the street because **it** was too tired*

Self-Attention



Шаг 1 - Векторы Queries, Keys, Values



Шаг 1 - Векторы Queries, Keys, Values

$$\mathbf{X} \times \mathbf{W}^Q = \mathbf{Q}$$

Diagram illustrating the computation of Query vectors (\mathbf{Q}). An input matrix \mathbf{X} (green, 2x4) is multiplied by a weight matrix \mathbf{W}^Q (purple, 4x4). The result is a query matrix \mathbf{Q} (purple, 2x4).

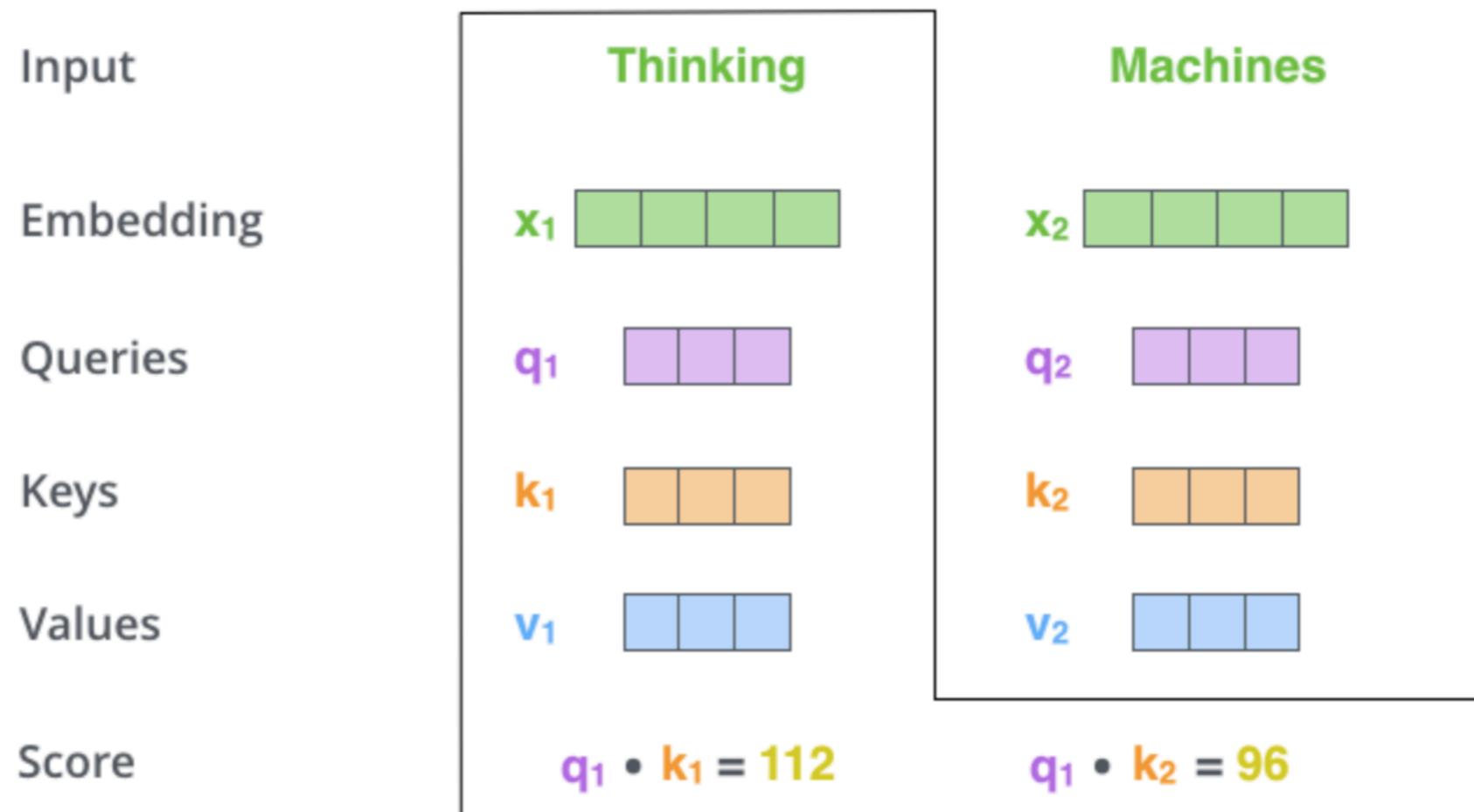
$$\mathbf{X} \times \mathbf{W}^K = \mathbf{K}$$

Diagram illustrating the computation of Key vectors (\mathbf{K}). An input matrix \mathbf{X} (green, 2x4) is multiplied by a weight matrix \mathbf{W}^K (orange, 4x4). The result is a key matrix \mathbf{K} (orange, 2x4).

$$\mathbf{X} \times \mathbf{W}^V = \mathbf{V}$$

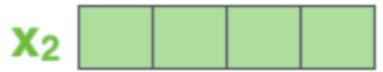
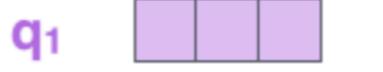
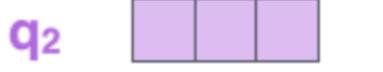
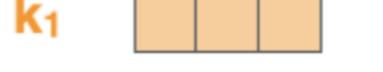
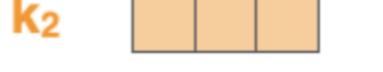
Diagram illustrating the computation of Value vectors (\mathbf{V}). An input matrix \mathbf{X} (green, 2x4) is multiplied by a weight matrix \mathbf{W}^V (blue, 4x4). The result is a value matrix \mathbf{V} (blue, 2x4).

Шаг 2 - Вычисление Score



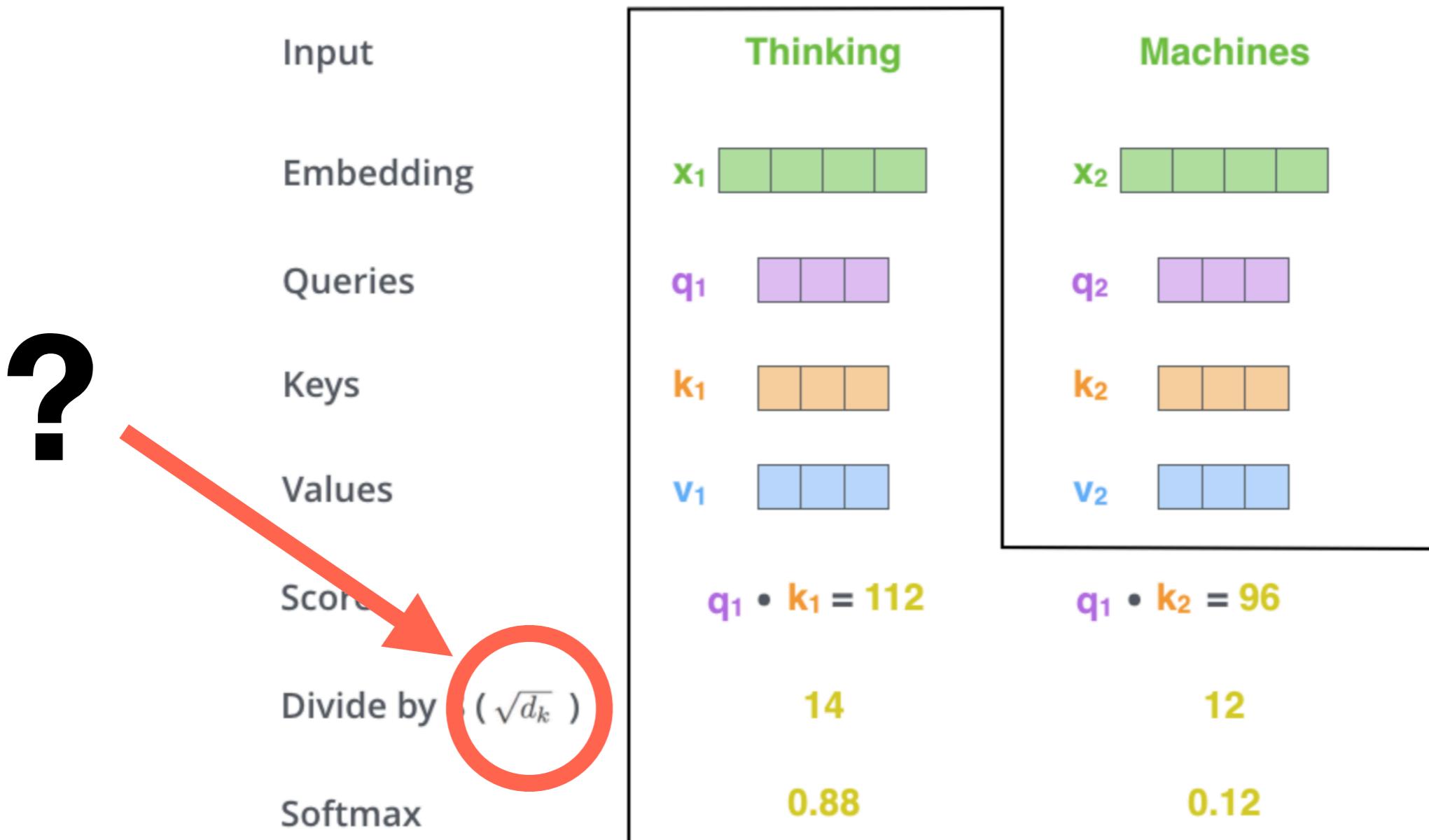
Шаг 3: Делим Скор на $\sqrt{d_k}$

Шаг 4: Softmax

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

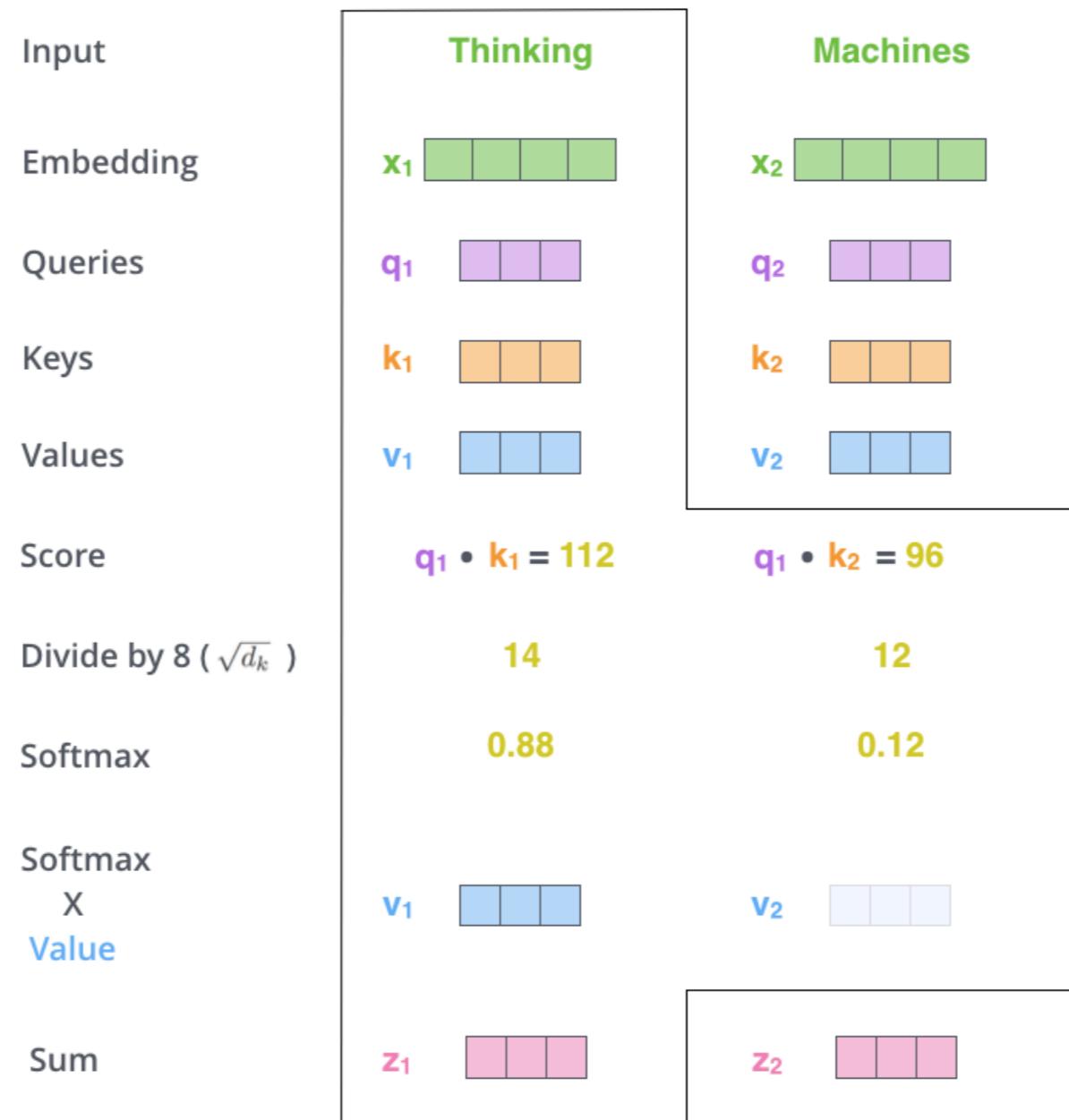
Шаг 3: Делим Скор на $\sqrt{d_k}$

Шаг 4: Softmax

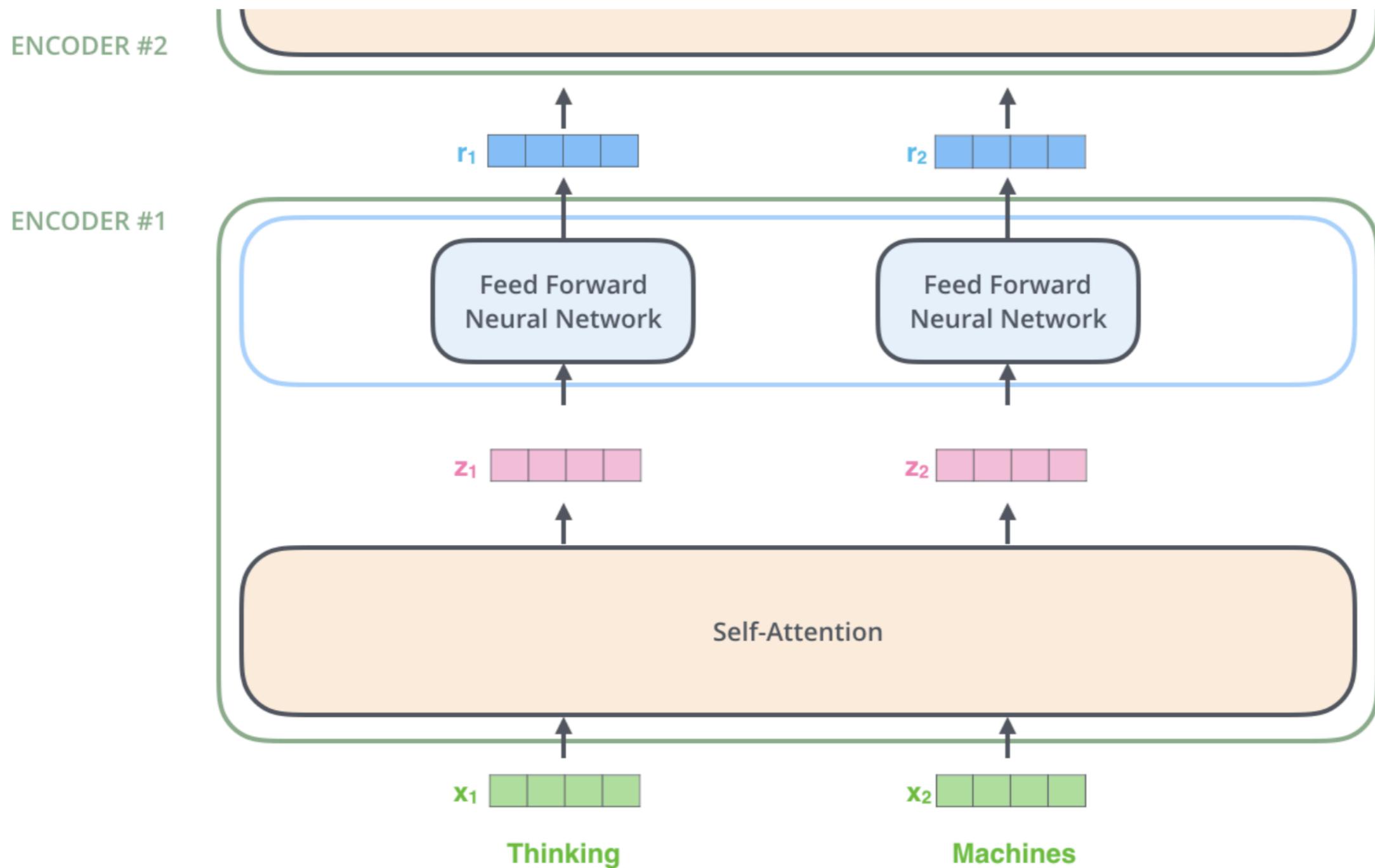


Шаг 5: Умножение каждого Values вектора на значение softmax

Шаг 6: Суммирование векторов



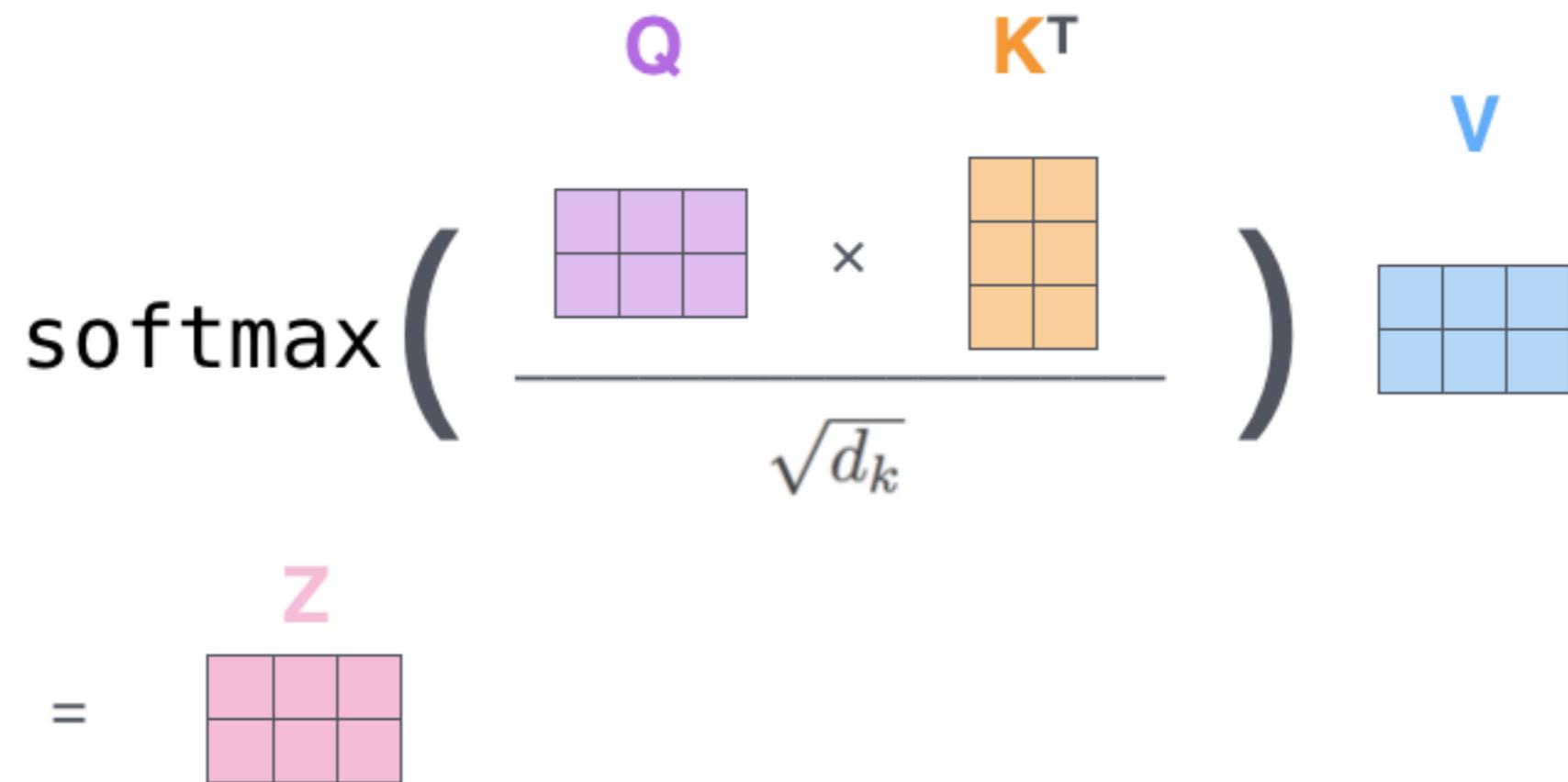
Self-Attention



Self-Attention

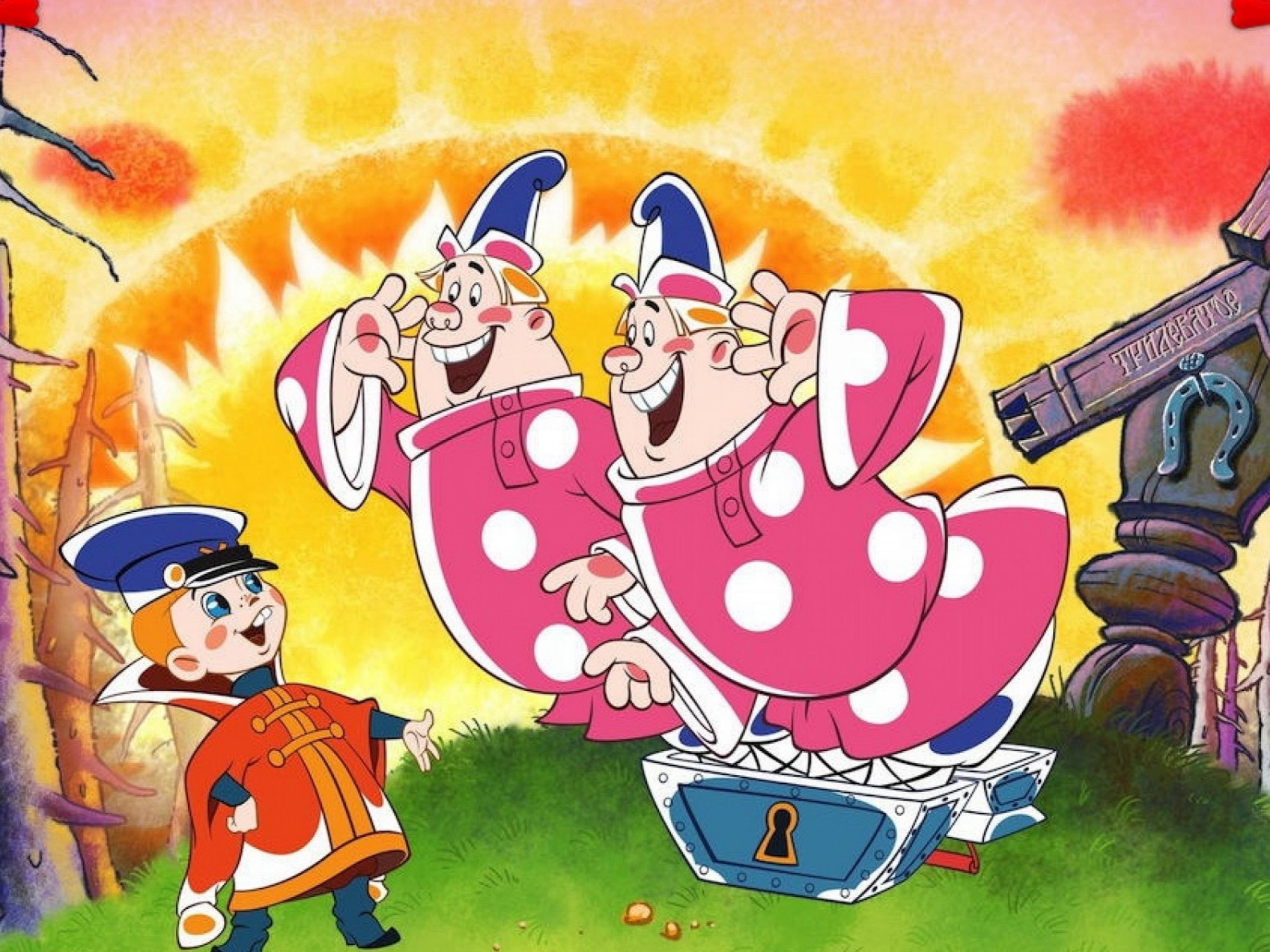
$$\text{softmax} \left(\frac{\begin{matrix} \mathbf{Q} & \mathbf{K}^T \\ \times & \end{matrix}}{\sqrt{d_k}} \right) \mathbf{V}$$

= \mathbf{Z}

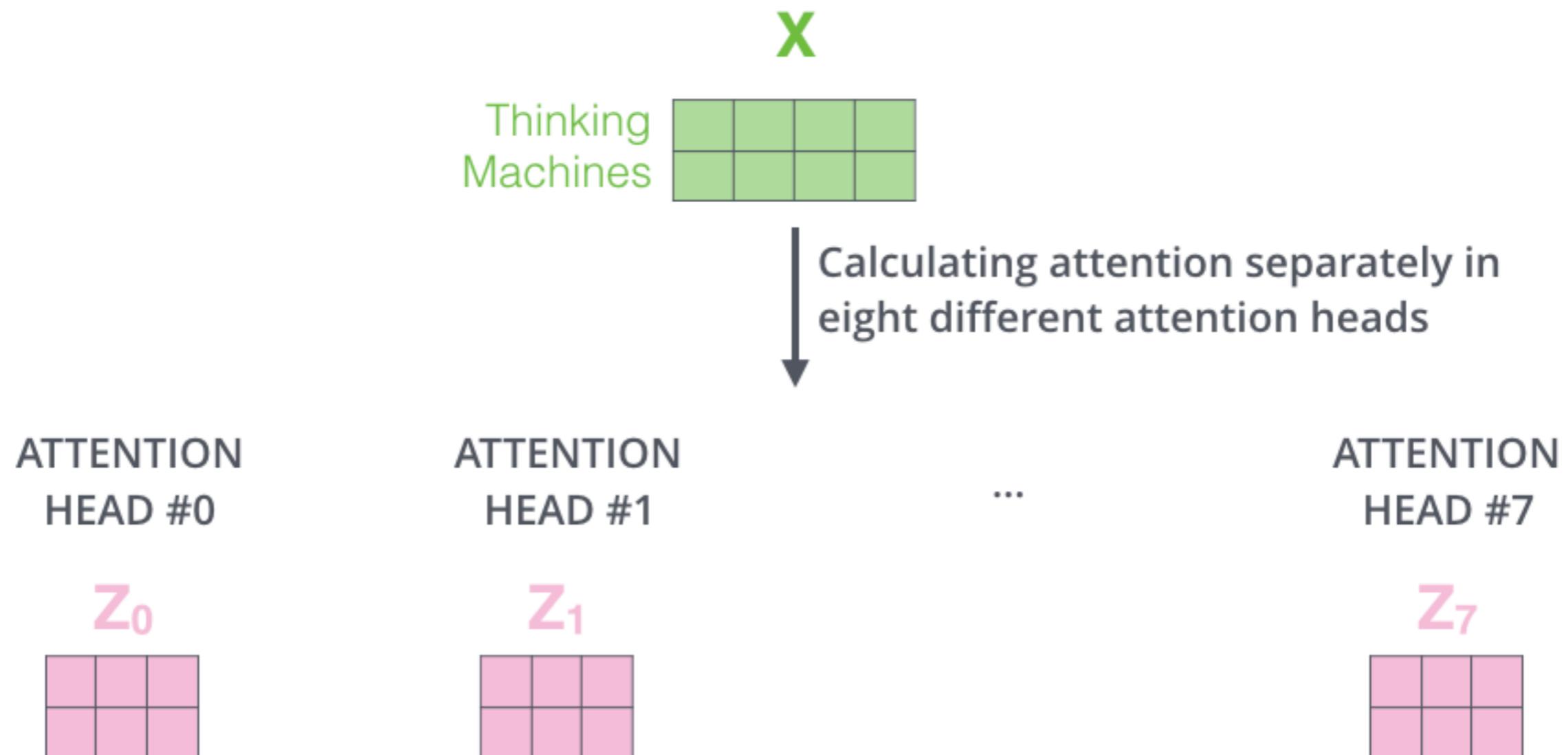


The diagram illustrates the Self-Attention mechanism. It shows the computation of attention weights from query matrix \mathbf{Q} and key matrix \mathbf{K}^T , scaled by the square root of d_k , followed by a multiplication with value matrix \mathbf{V} . The result is labeled \mathbf{Z} .

The matrices \mathbf{Q} , \mathbf{K}^T , and \mathbf{V} are represented as 3x3 grids. The query matrix \mathbf{Q} is purple, the key matrix \mathbf{K}^T is orange, and the value matrix \mathbf{V} is blue. The resulting matrix \mathbf{Z} is pink.



Multihead Self-Attention

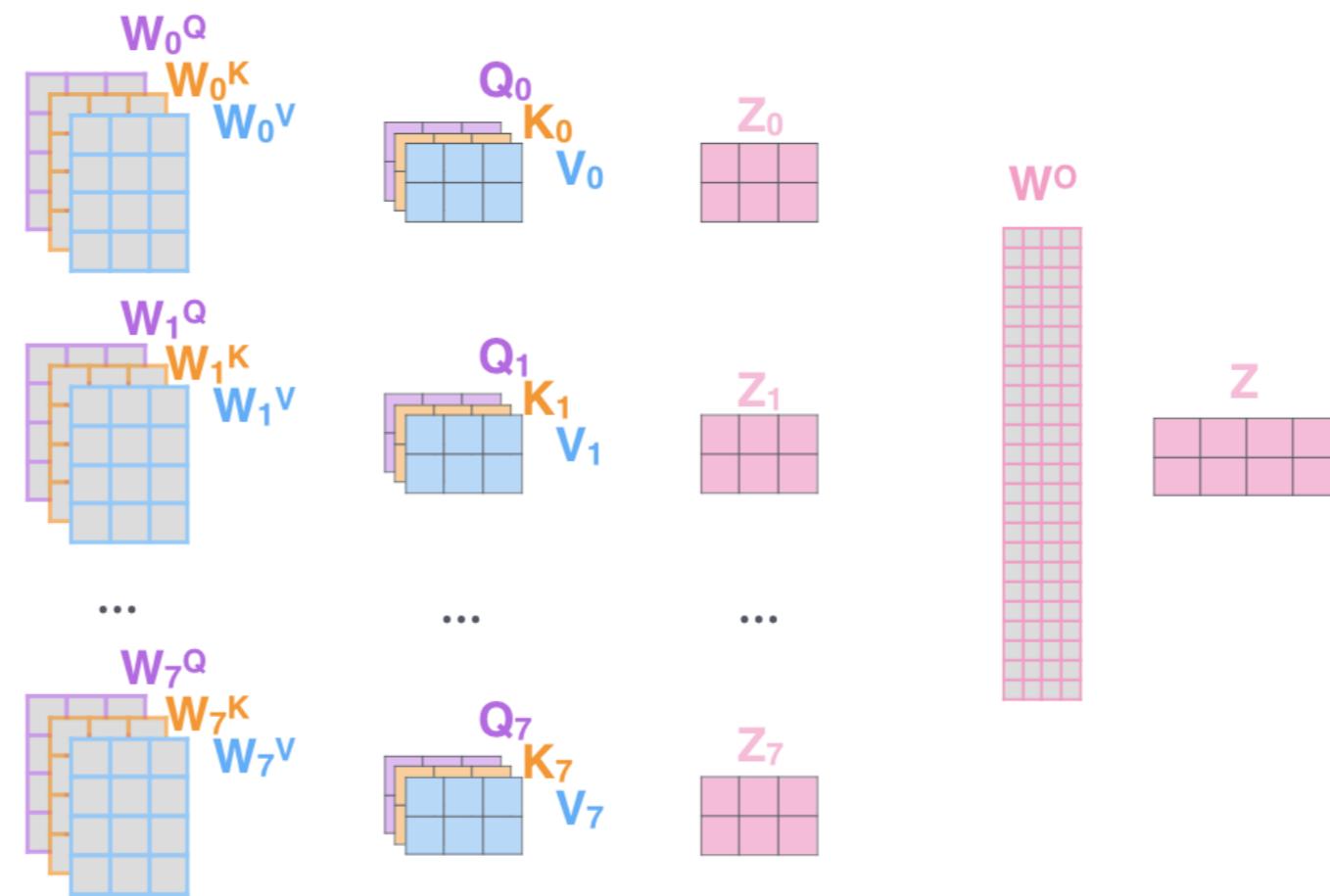


Multihead Self-Attention

- 1) This is our input sentence* X
- 2) We embed each word* R
- 3) Split into 8 heads. We multiply X or R with weight matrices W_0^Q, W_0^K, W_0^V
- 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



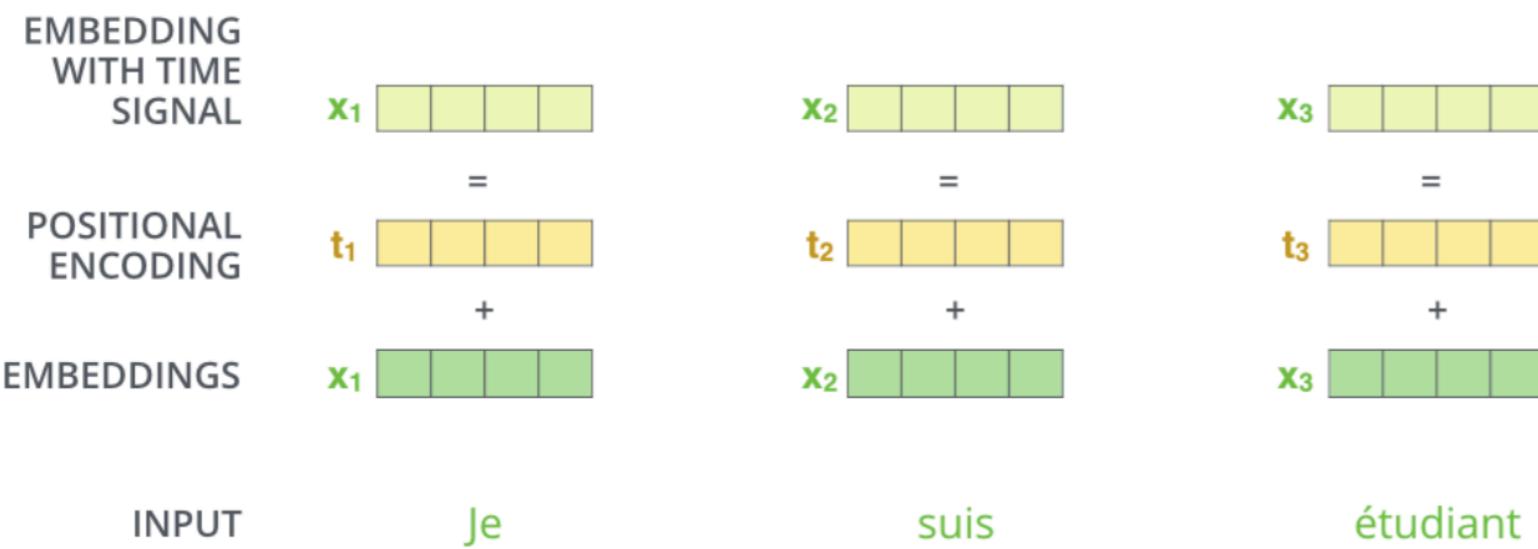
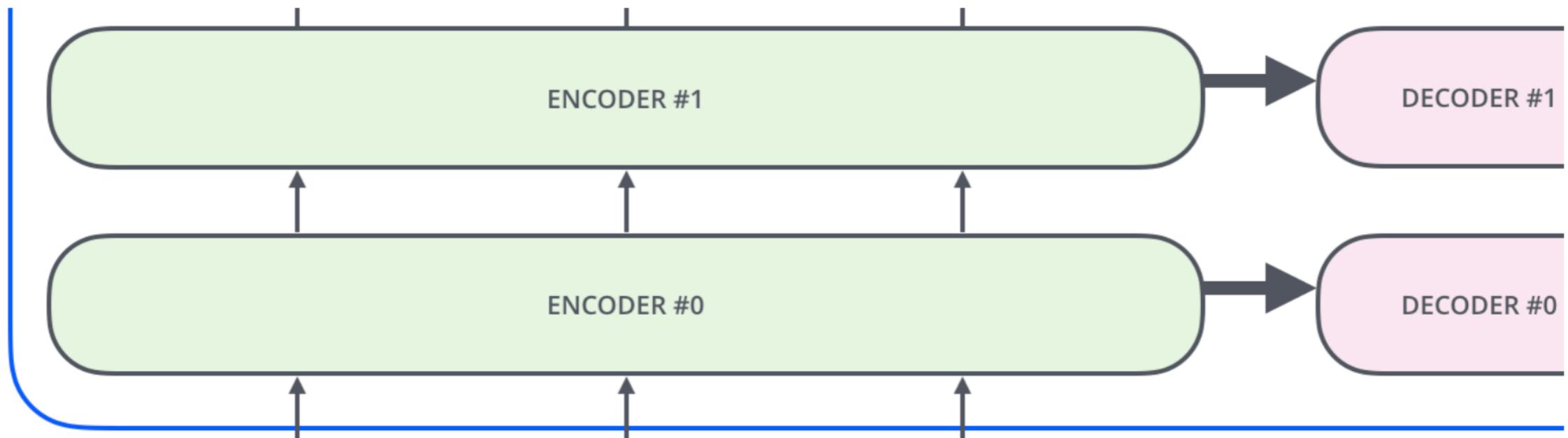
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



Но ведь нам нужно обрабатывать последовательности...



Positional Encoding



Positional Encoding

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

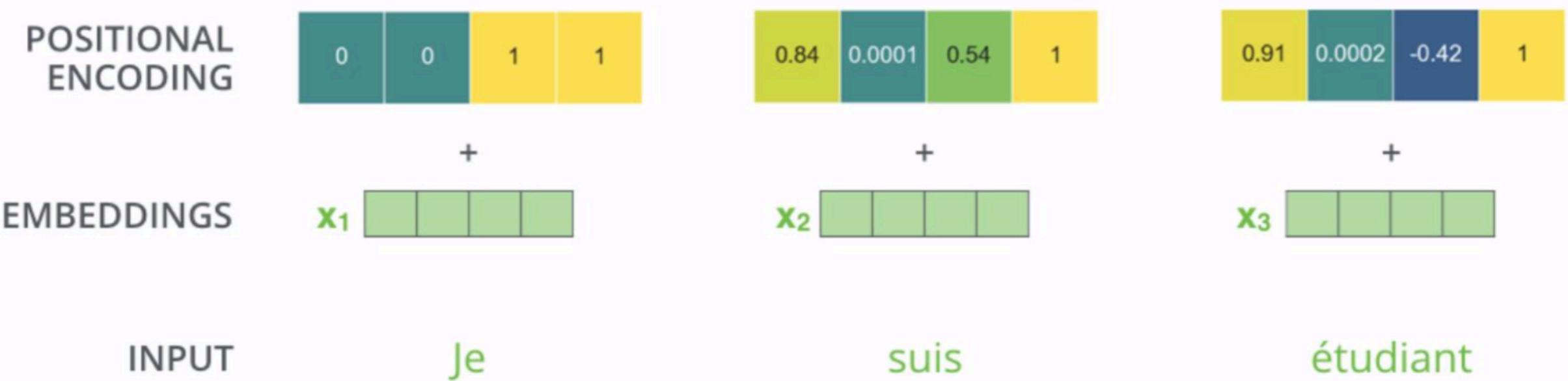
$$PE_{(pos, \ 2i + 1)} = \cos(pos / 10000^{2i / d_{\text{model}}})$$

- *pos* is the position
- *i* is the dimension.

Each dimension of the positional encoding corresponds to a sinusoid.
The wavelengths form a geometric progression from 2π to $2\pi * 10\ 000$

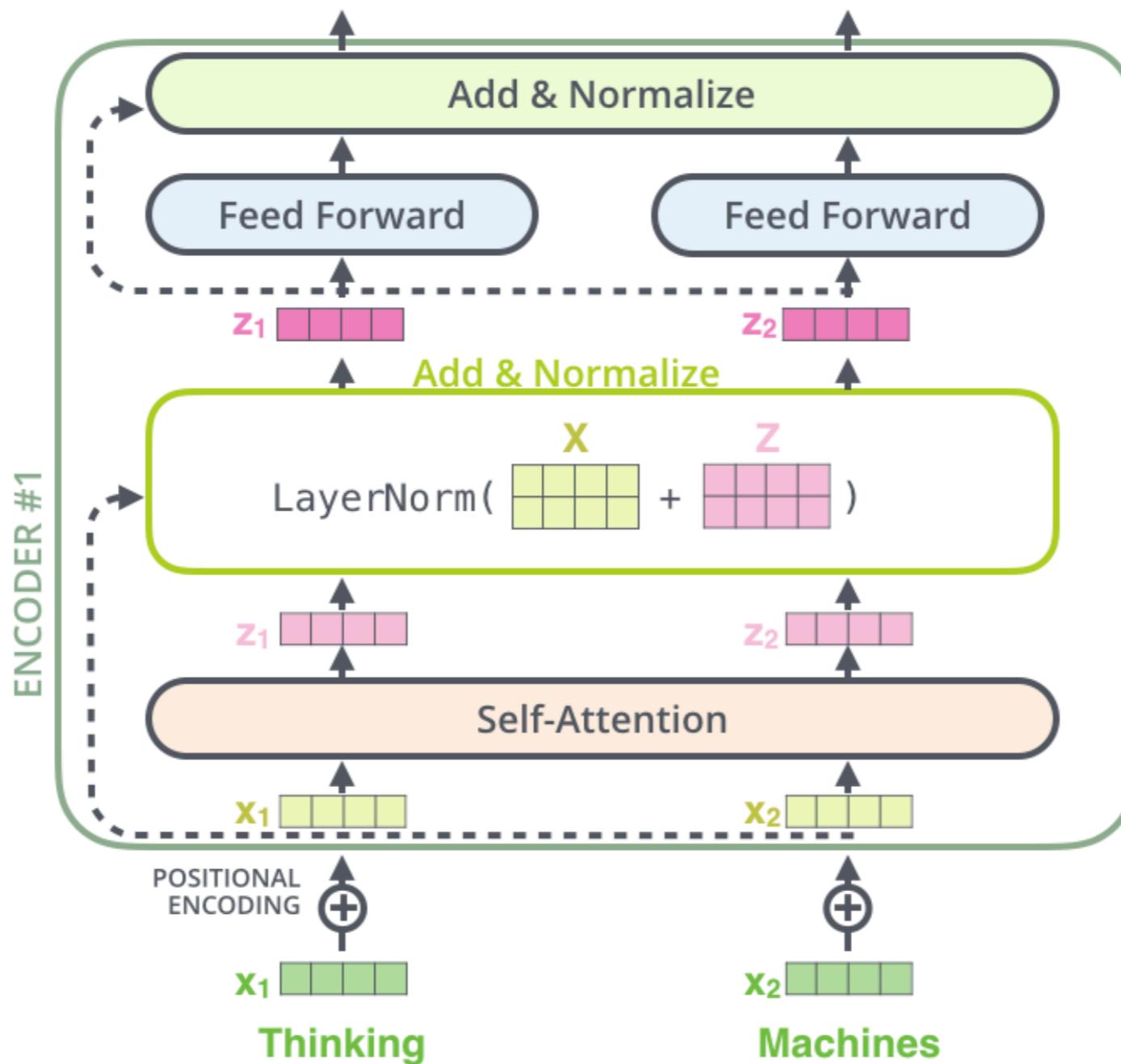


Positional Encoding



[https://github.com/ml-mipt/ml-mipt/blob/advanced/week04_Transformer/
week04_positional_encoding_carriers.ipynb](https://github.com/ml-mipt/ml-mipt/blob/advanced/week04_Transformer/week04_positional_encoding_carriers.ipynb)

The Residuals & layer-normalization



Self-Attention

