

Lecture 5

- attention history behind
- dot-product attention
- transformer architecture
- positional embeddings
- bert and gpt

Shkhanukova Milana, 2023

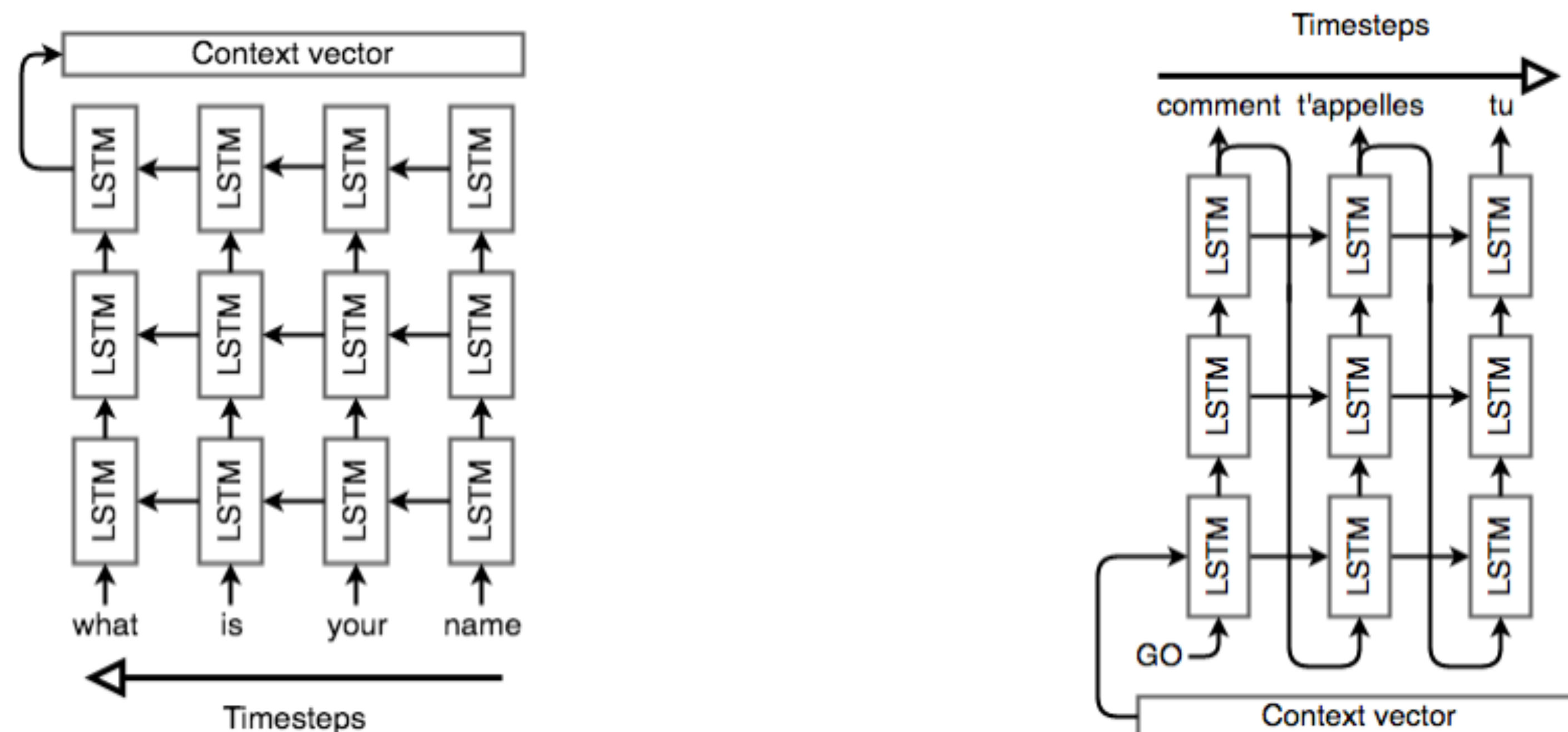
Seq2seq

Откуда все пошло?

Sequence-to-sequence, or «Seq2Seq» (2014 for English-French translation)

At a high level, a sequence-to-sequence model is an **end-to-end model** made up of two recurrent neural networks:

- an **encoder**, which takes the model's input sequence as input and encodes it into a fixed-size "context vector", and
- a **decoder**, which uses the context vector from above as a "seed" from which to generate an output sequence

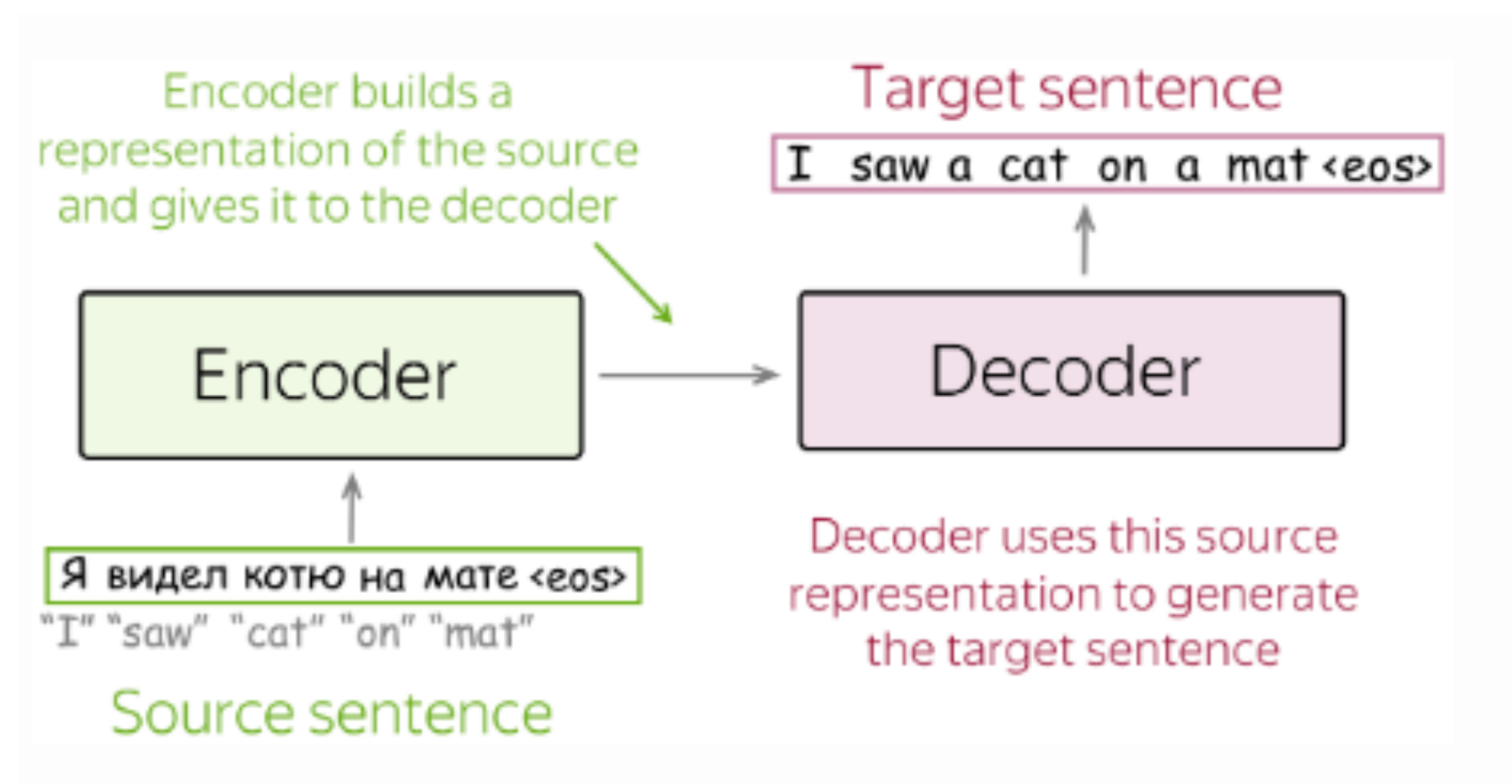


Откуда все пошло?

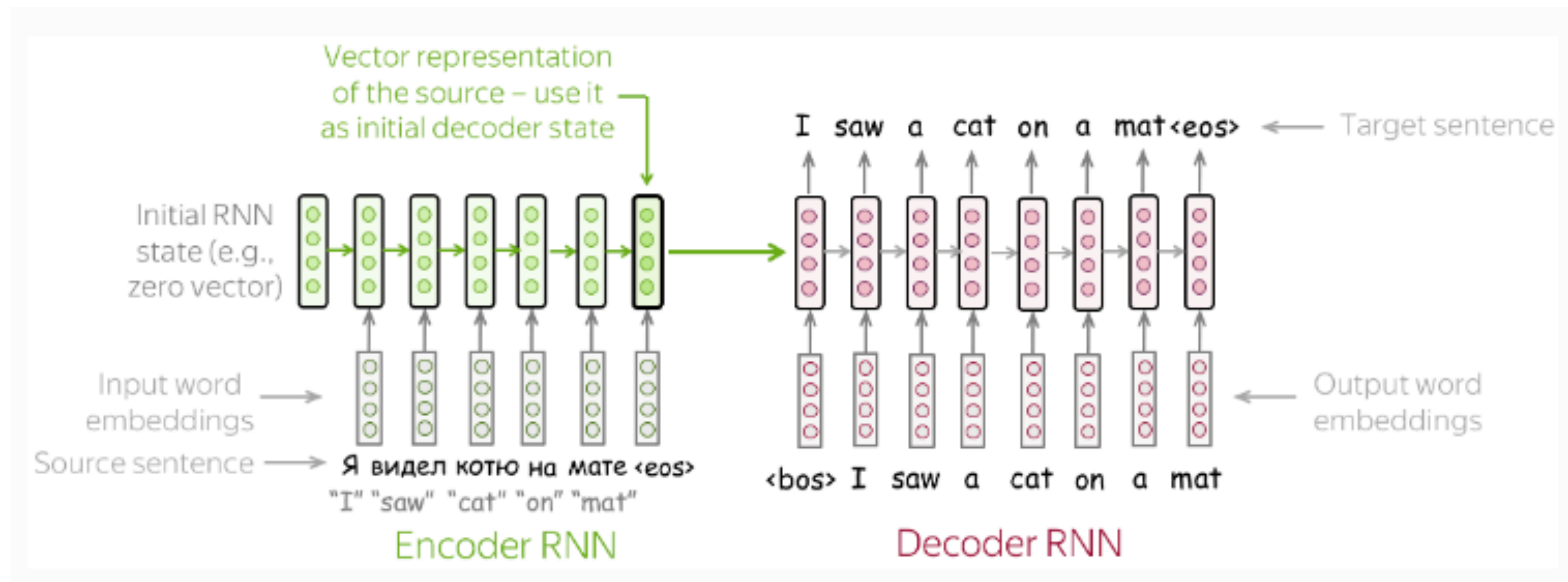
Sequence-to-sequence, or «Seq2Seq» (2014 for English-French translation)

At a high level, a sequence-to-sequence model is an **end-to-end model** made up of two recurrent neural networks:

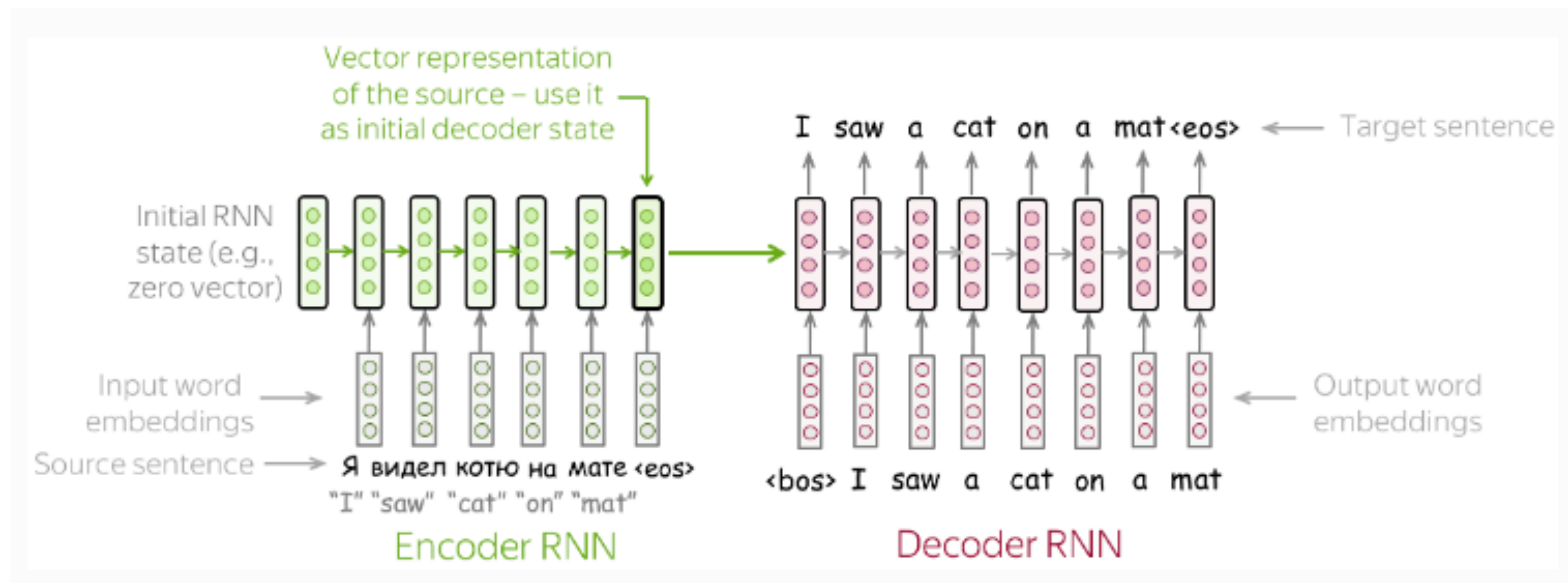
- an **encoder**, which takes the model's input sequence as input and encodes it into a fixed-size "context vector", and
- a **decoder**, which uses the context vector from above as a "seed" from which to generate an output sequence



Какие тут могут быть две проблемы?



Какие тут могут быть две проблемы?

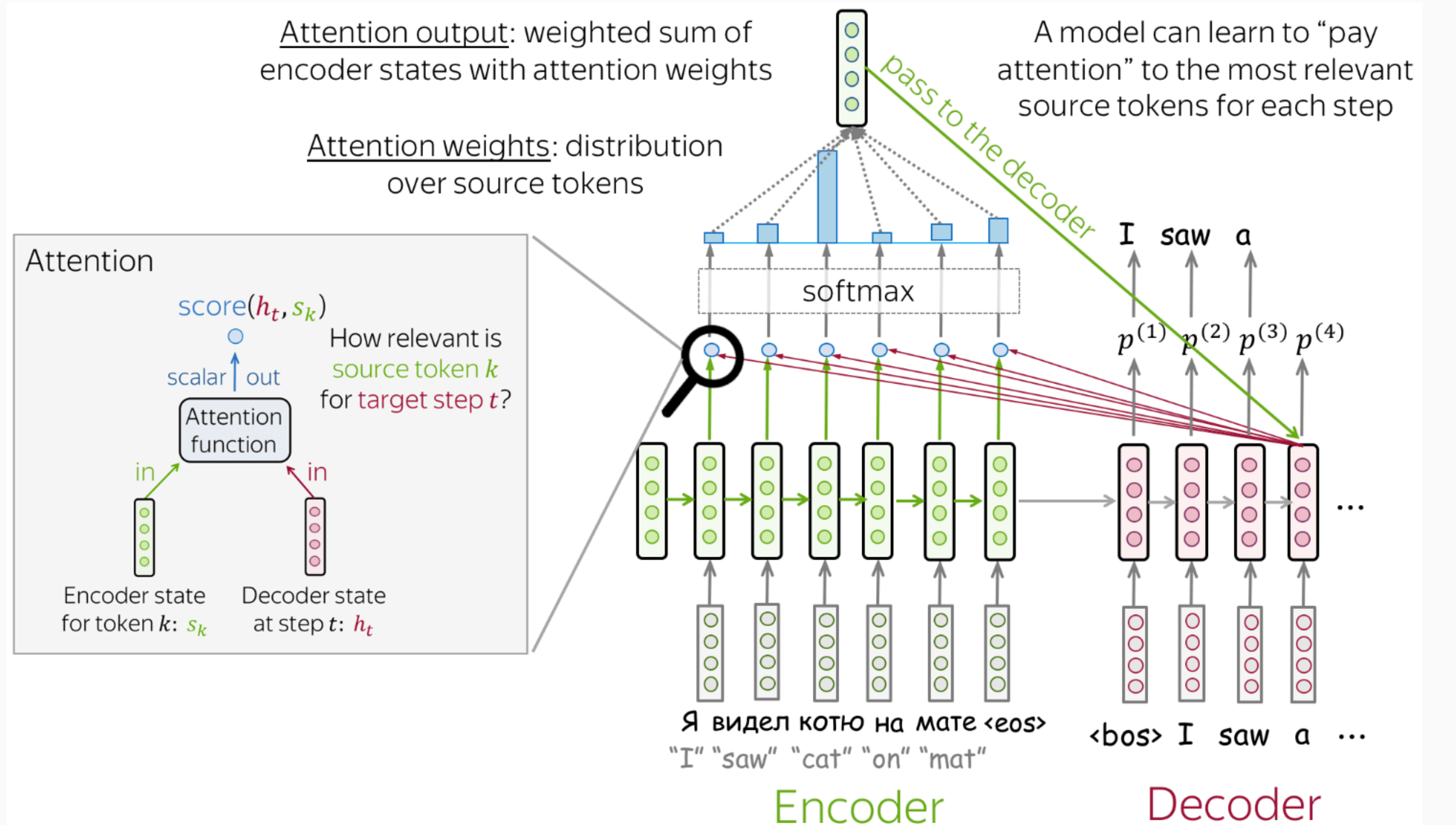


for the encoder, it is hard to compress the sentence

for the decoder, different information may be relevant at different steps.

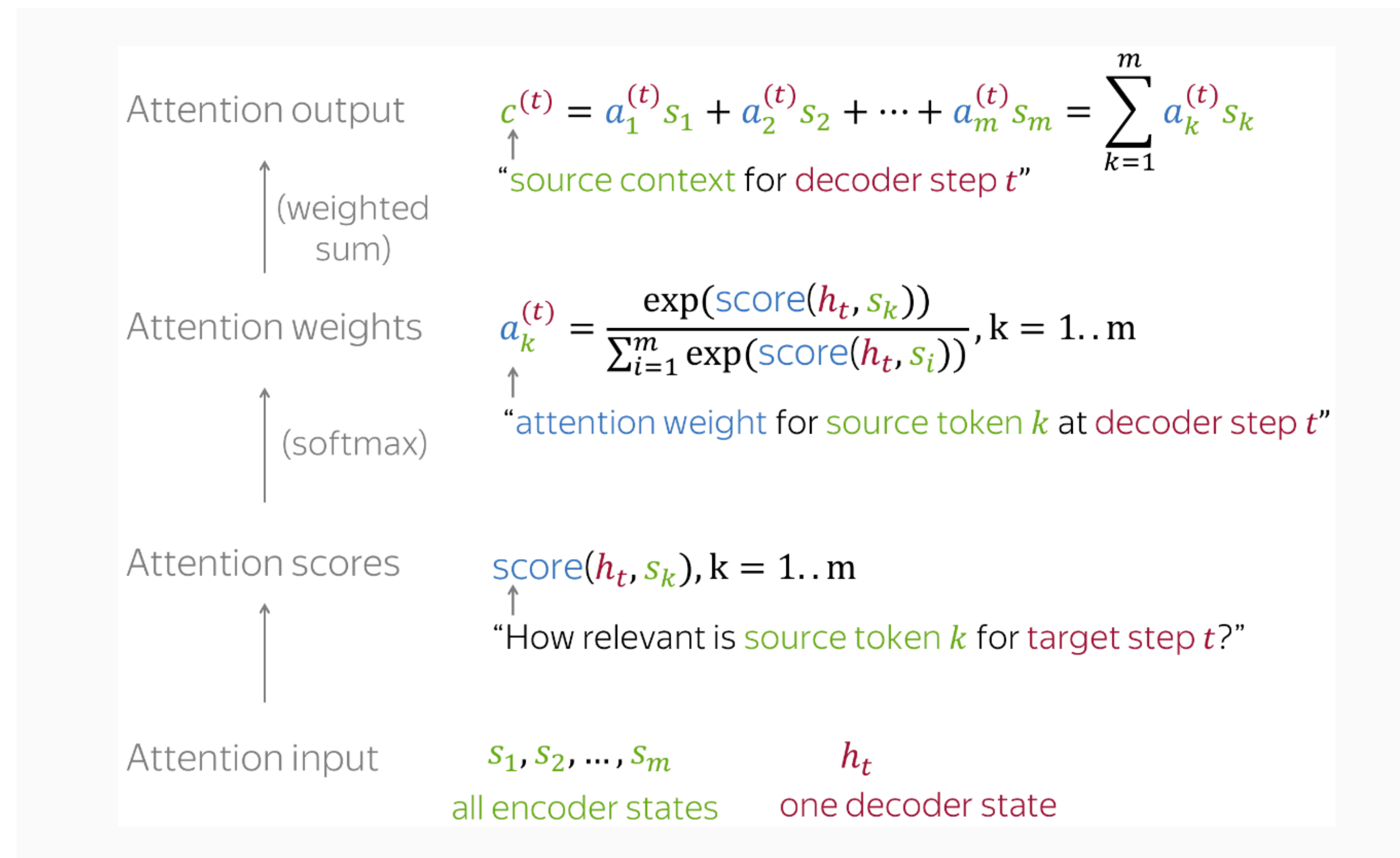
Simple idea behind

At different steps, let a model "focus" on different parts of the input.



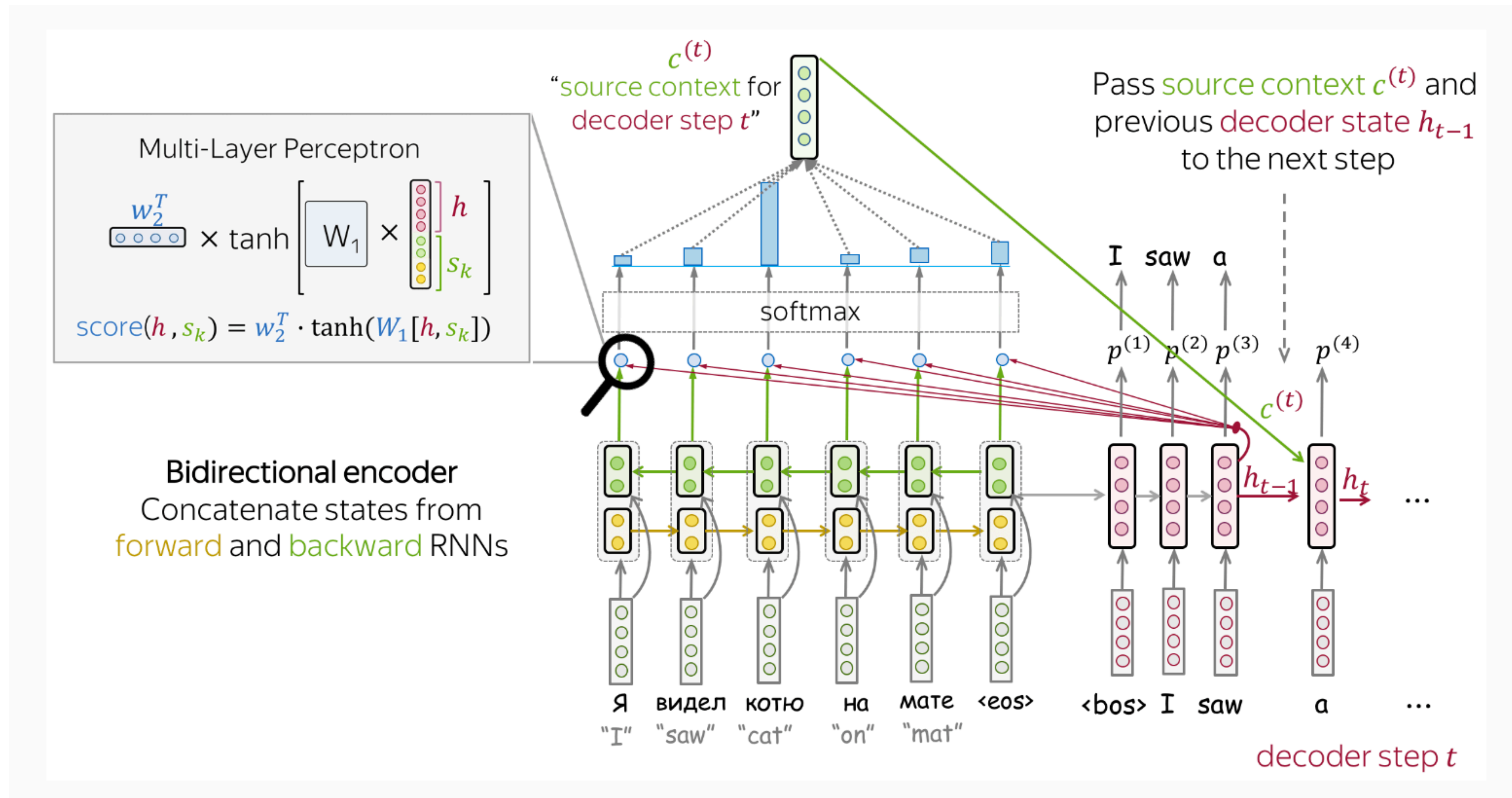
Simple idea behind

At different steps, let a model "focus" on different parts of the input.

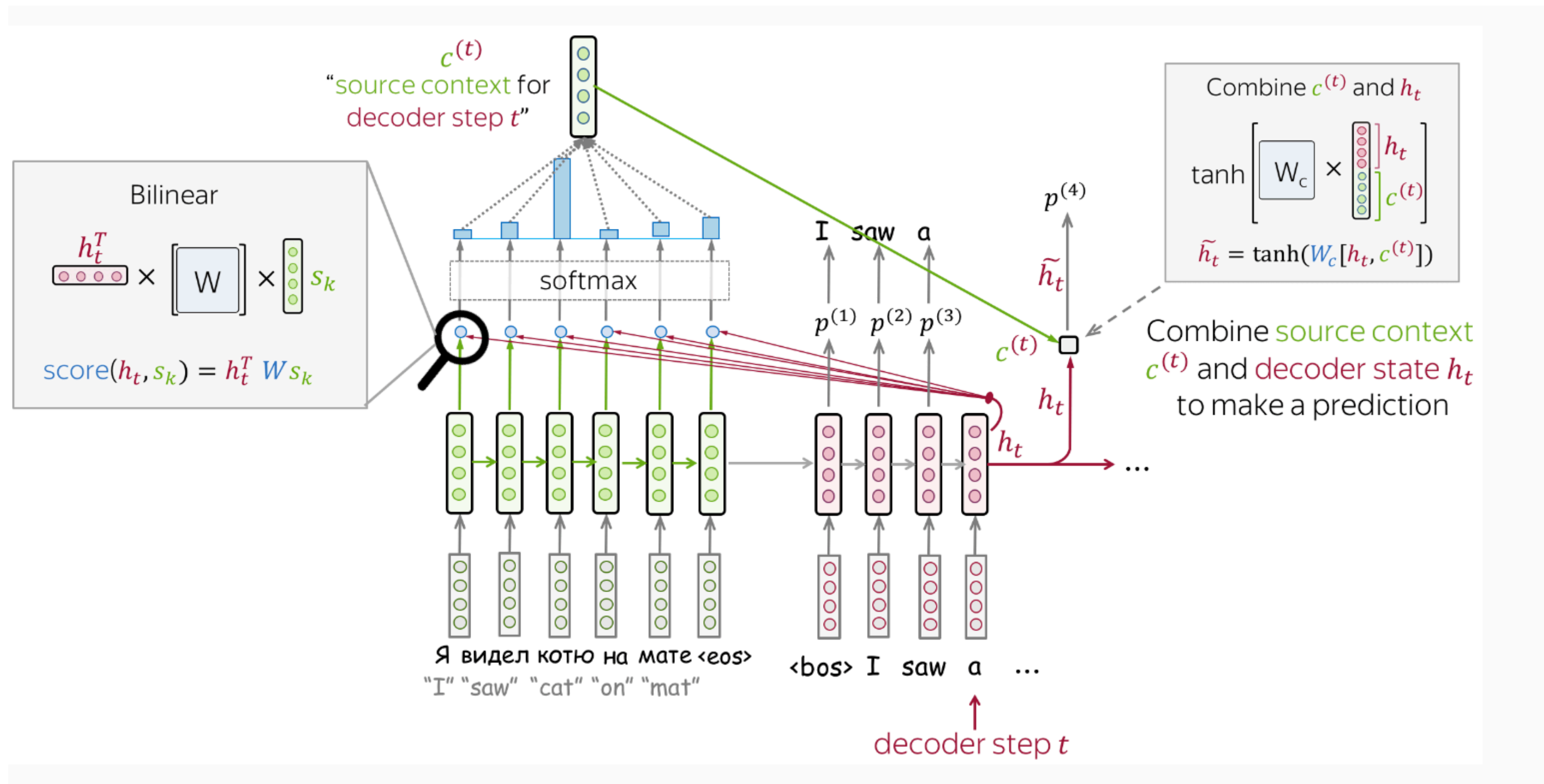


**Attention
before 2017**

Bahdanau

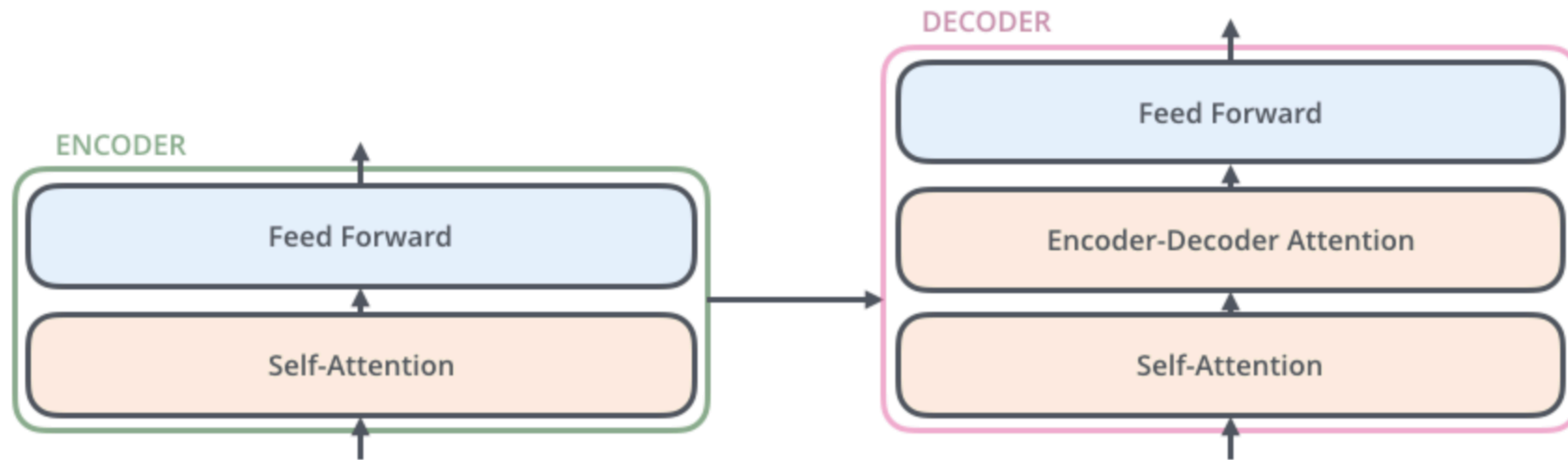


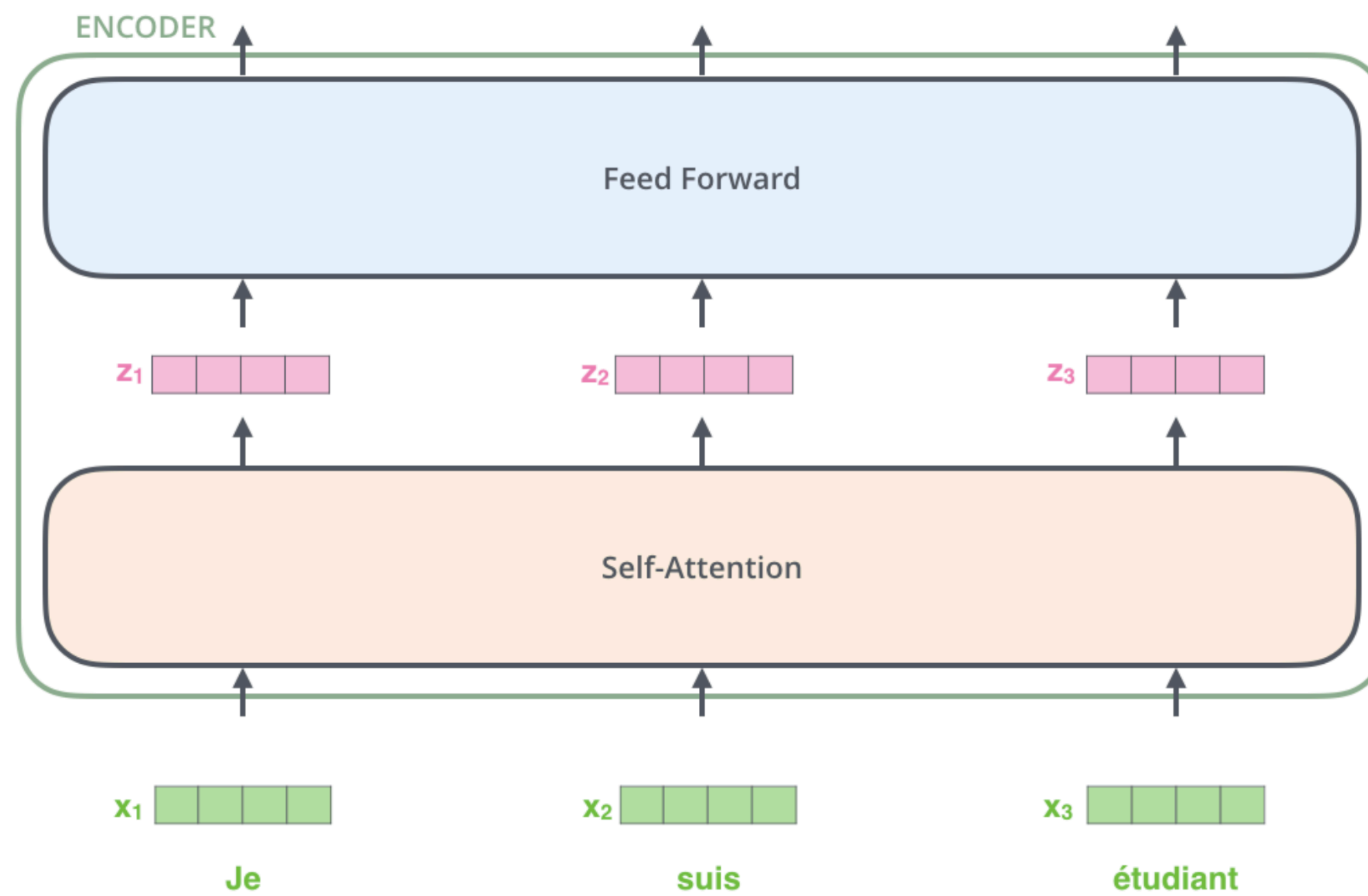
Luong



Attention In transformer step-by-step

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





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

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

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

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

1. Создаем три вектора - queries, key, values.
2. Взвешиваем каждое слово относительно всех остальных.
(Dot-product, query and key)
3. Softmax - теперь все значения будут positive и суммироваться в 1
4. Умножаем value на softmax score

The main idea behind self-attention is that instead of using a fixed embedding for each token, we can use the whole sequence to compute a *weighted average* of each embedding.

Given a sequence of token embeddings x_1, \dots, x_n , self-attention produces a sequence of **new embeddings** x'_1, \dots, x'_n where each x'_i is a **linear combination of all the x_j** :

The coefficients w_{ji} are called *attention weights* and are normalized so that $\sum_j w_{ji} = 1$.

$$x'_i = \sum_{j=1}^N w_{ji} x_j \qquad \sum_j w_{ji} = 1.$$

Query, keys, values

you compare the query with the keys to get scores/weights for the values
(Сравниваем query с keys, чтобы получить веса для обновления values)

- key - относительно кого считаем, исходный текст на языке оригинала
- query - для кого считаем, переведенный текст на целевой язык
- value - на чем мы считаем, снова исходный текст

Супермаркет. Каждый ингредиент для окрошки - это **query** - то, что нужно купить. На полках много товаров.

Они могут не также называться, для нас это **key**.

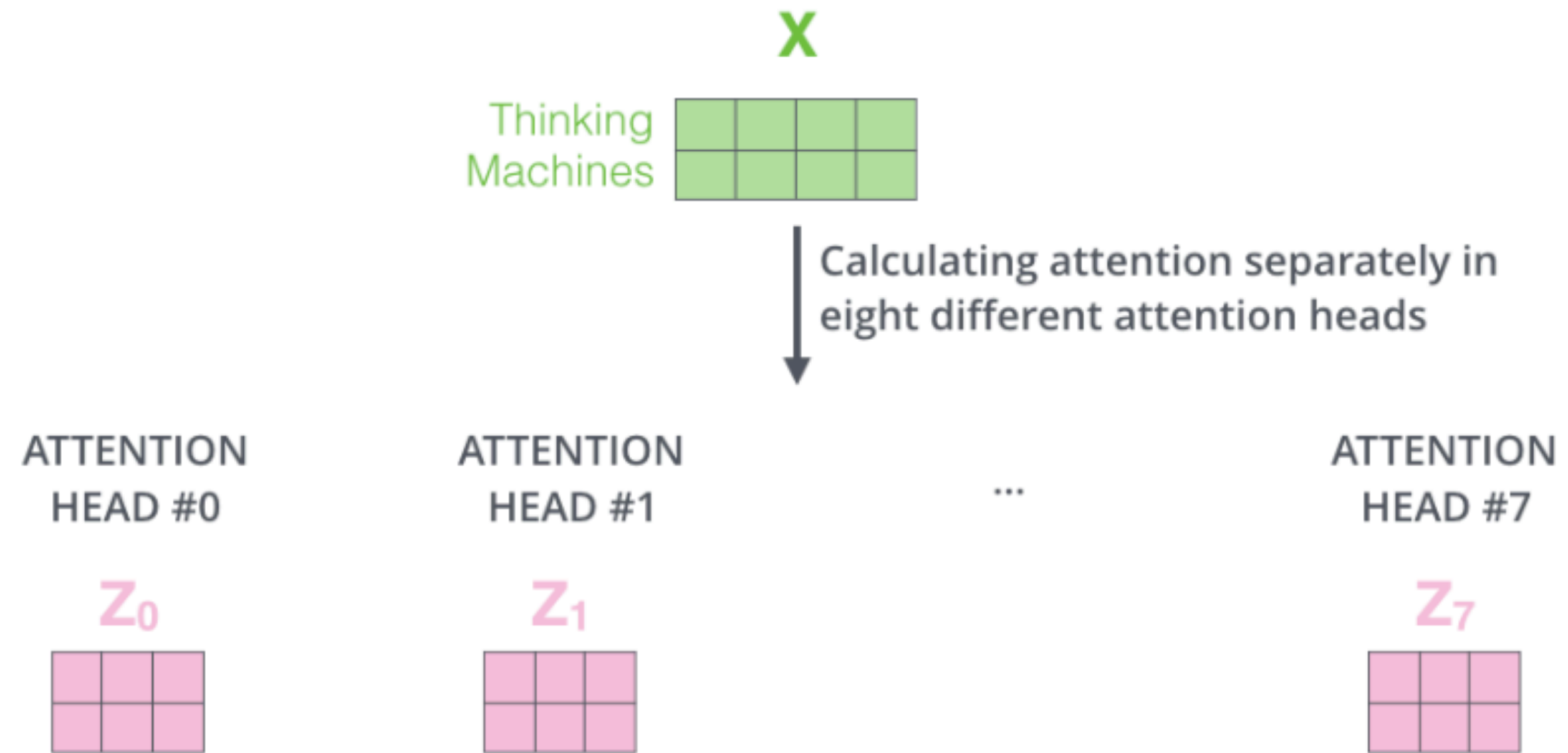
Сравниваем, как сильно похоже описание в списке ингредиентов с товарами.

Для self-attention значения будут немного более smooth.

Let's write from scratch

<https://colab.research.google.com/drive/10h66zyXkEyQCLc-DCyhdhrOf3qOzRtZ7?usp=sharing>

Multi-head attention



1. Give multiple representation subspaces
2. Focus on different positions

1) This is our input sentence*

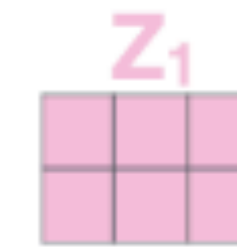
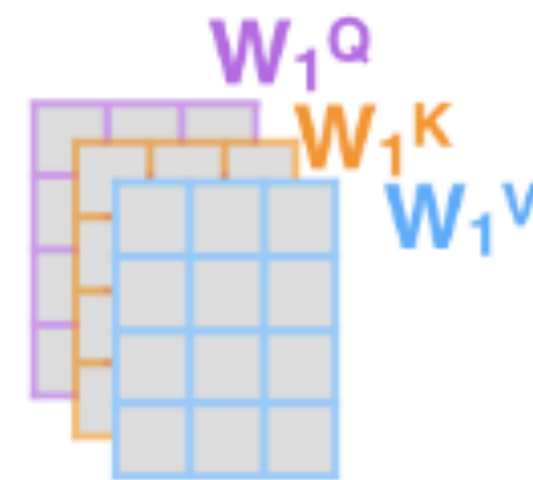
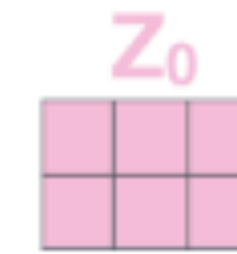
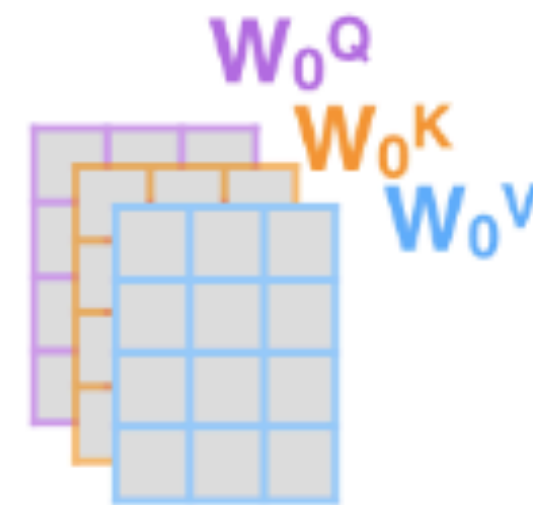
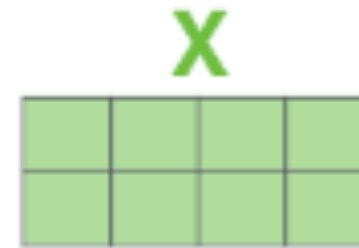
2) We embed each word*

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

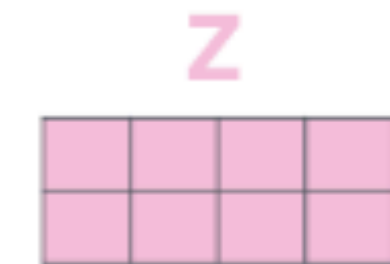
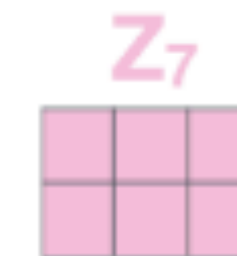
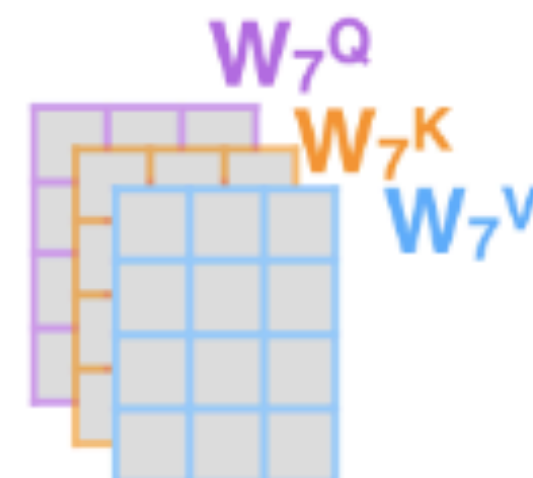
Thinking
Machines



...

...

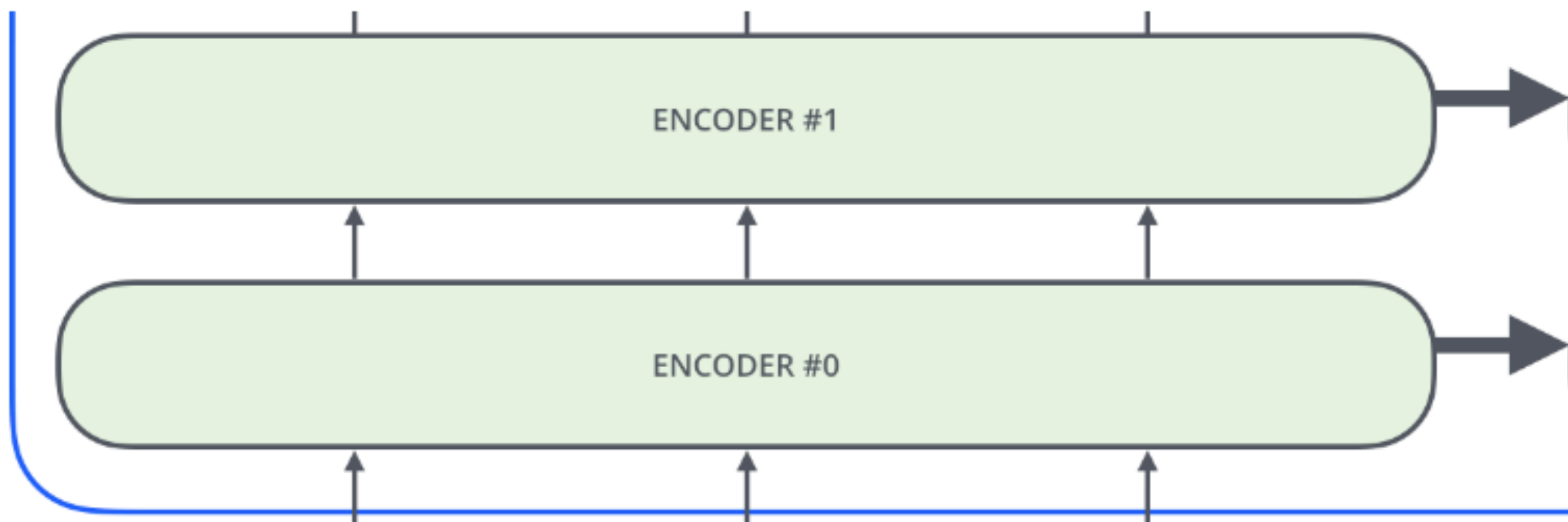
...



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

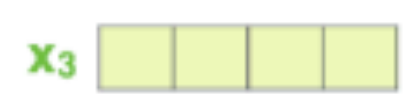
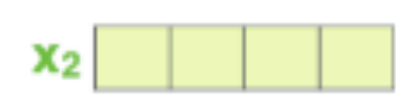
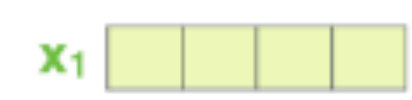


Positional encoding



1. Every position should have the same identifier, irrespective of the sequence length or input
2. Should not be huge to push the meaning in other subspace

EMBEDDING
WITH TIME
SIGNAL

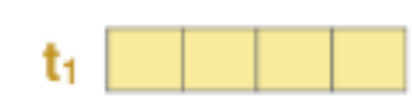


=

=

=

POSITIONAL
ENCODING

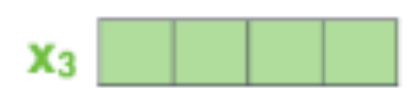
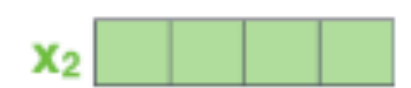
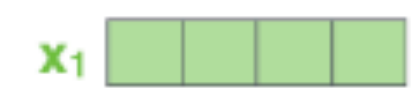


+

+

+

EMBEDDINGS



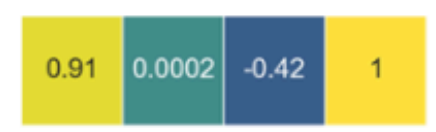
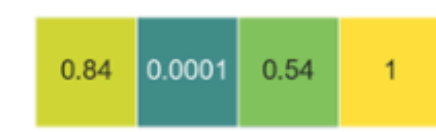
INPUT

Je

suis

étudiant

POSITIONAL
ENCODING

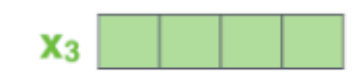
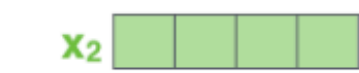
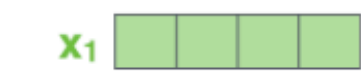


+

+

+

EMBEDDINGS



INPUT

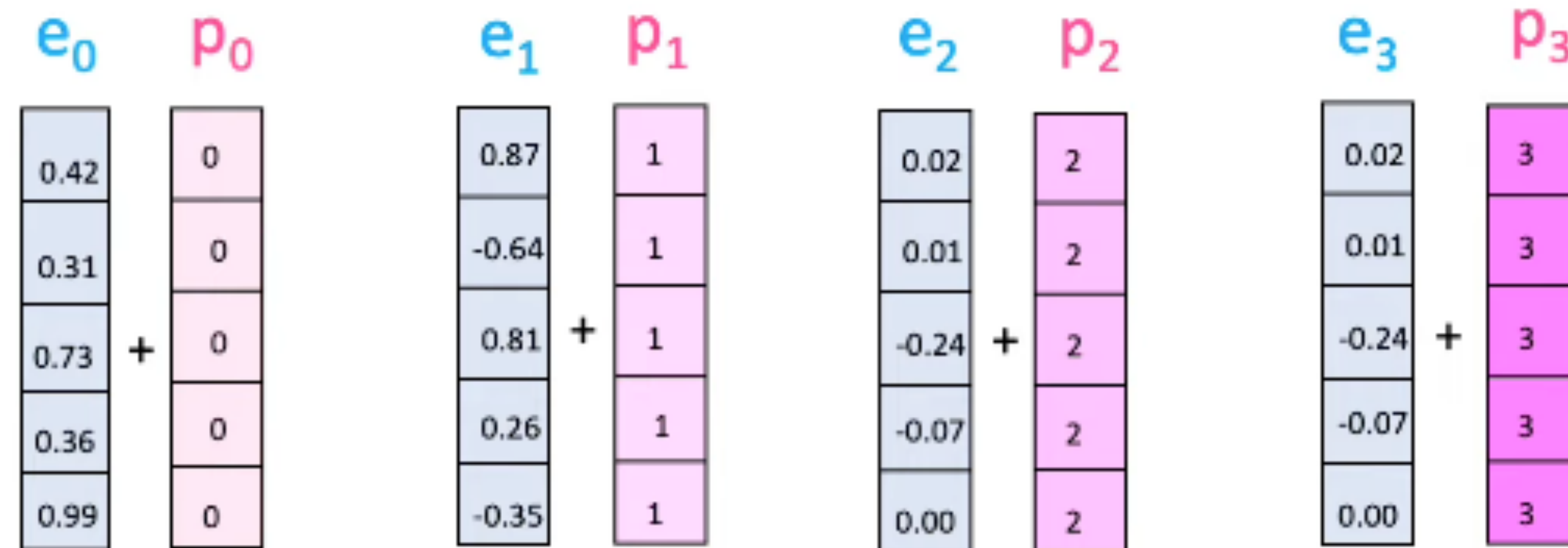
Je

suis

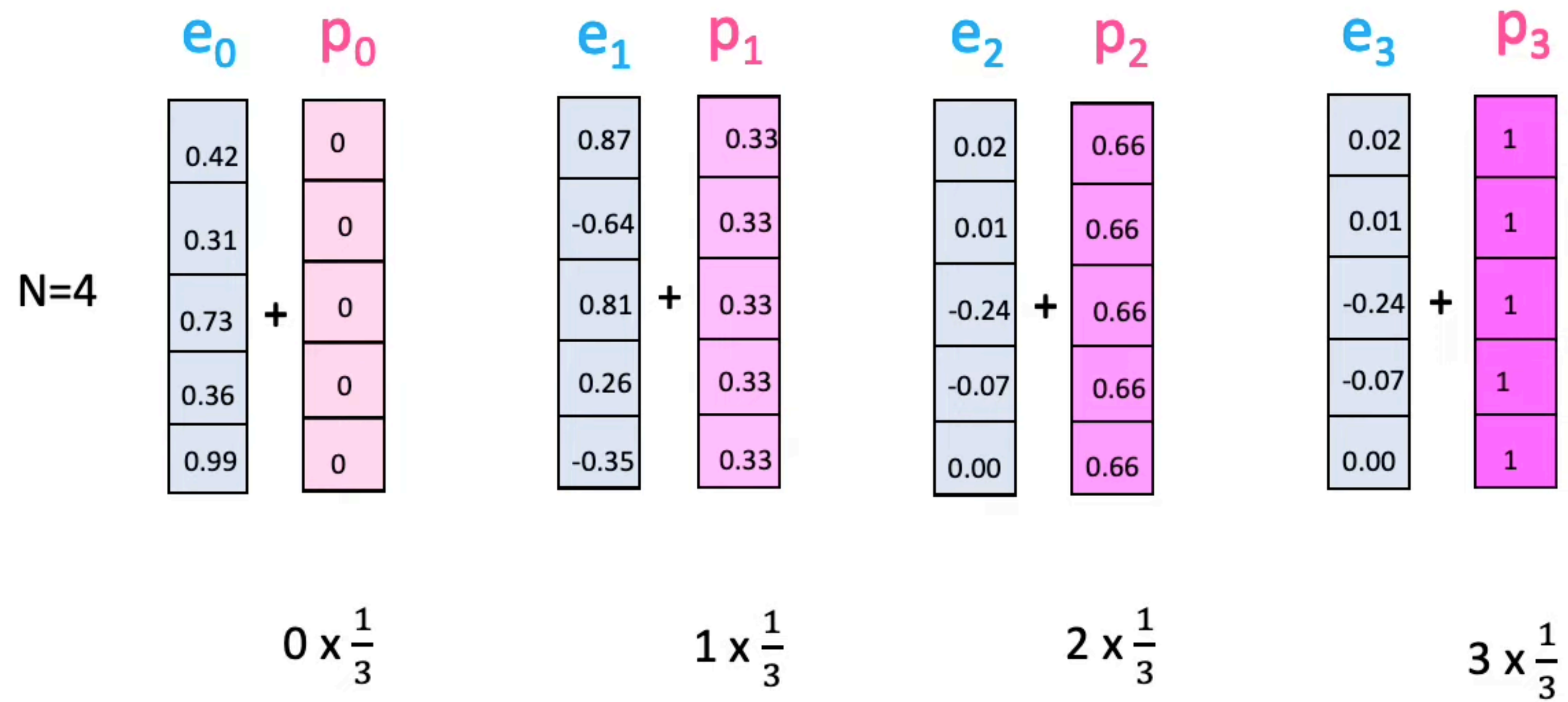
étudiant

Positional encoding

Even though she did **not** win the award, she was satisfied.
Even though she did win the award, she was **not** satisfied.



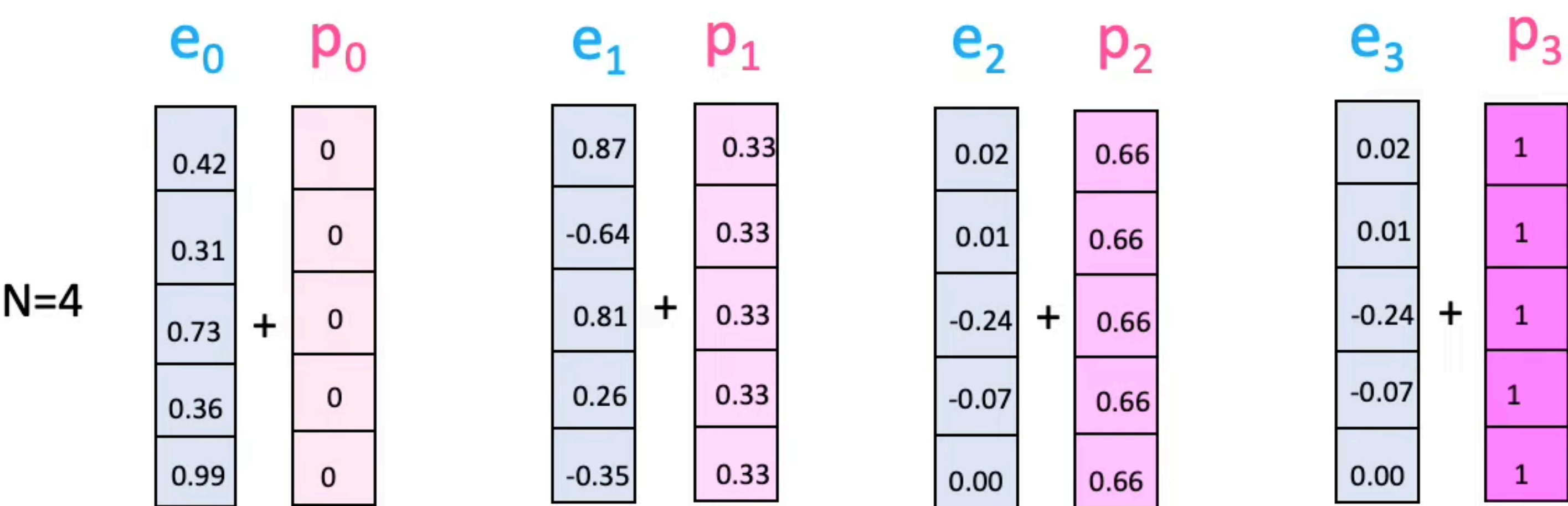
Positional encoding



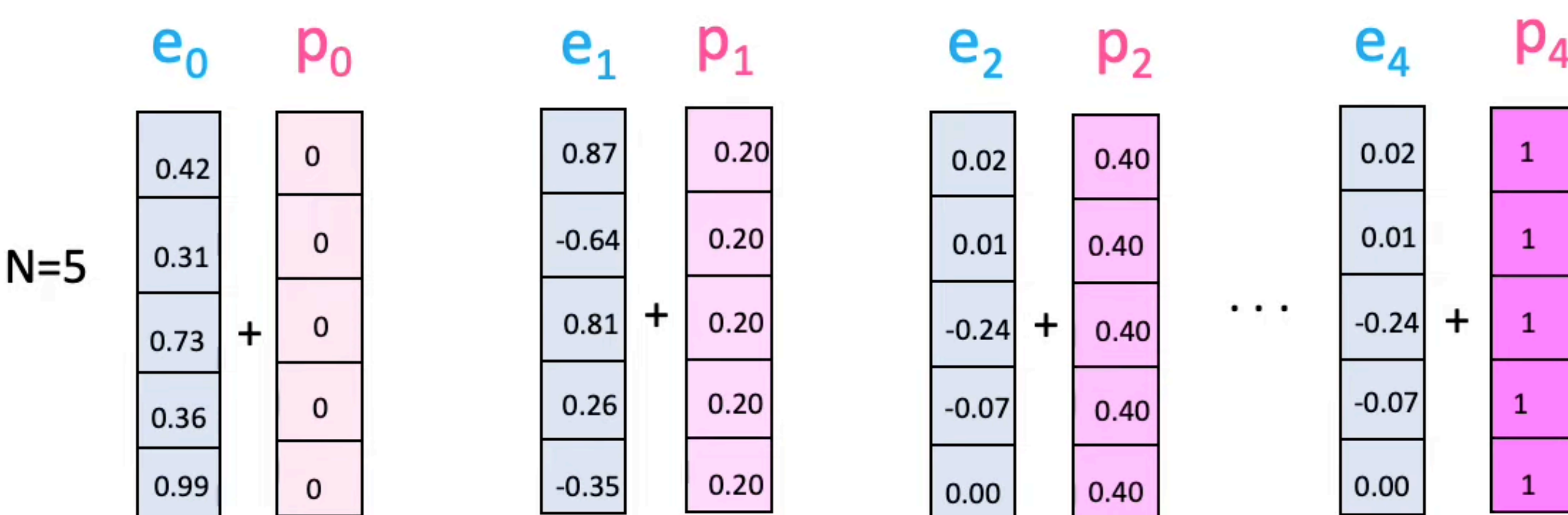
$$\frac{1}{N-1} = \frac{1}{3}$$

Positional encoding

Sentence 1



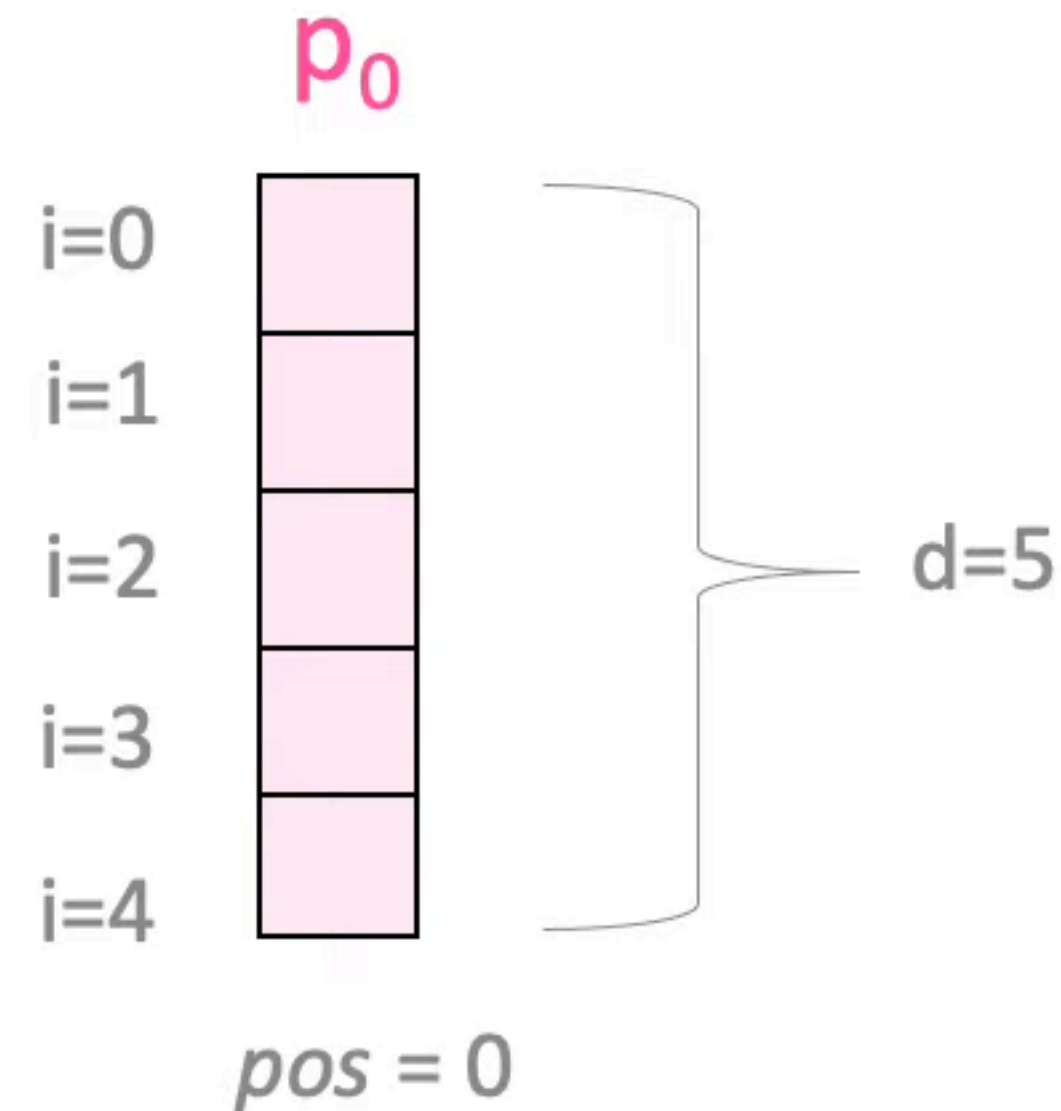
Sentence 2



Positional encoding

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$



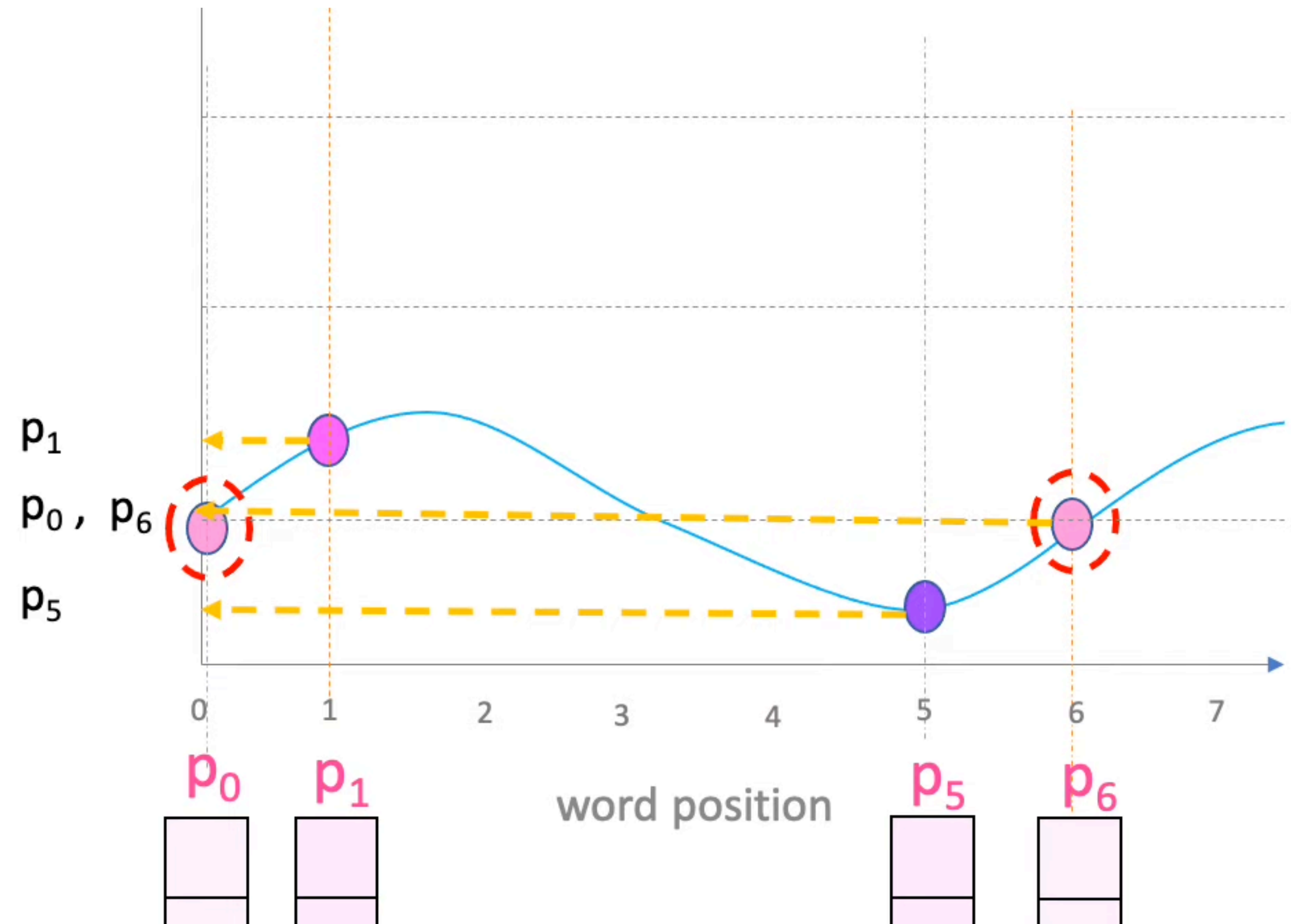
$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

The diagram shows the formula for the positional encoding. The term pos in the numerator and the denominator $10000^{\frac{2i}{d}}$ are circled with a red dashed line. A red arrow points to the denominator.

Positional encoding

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

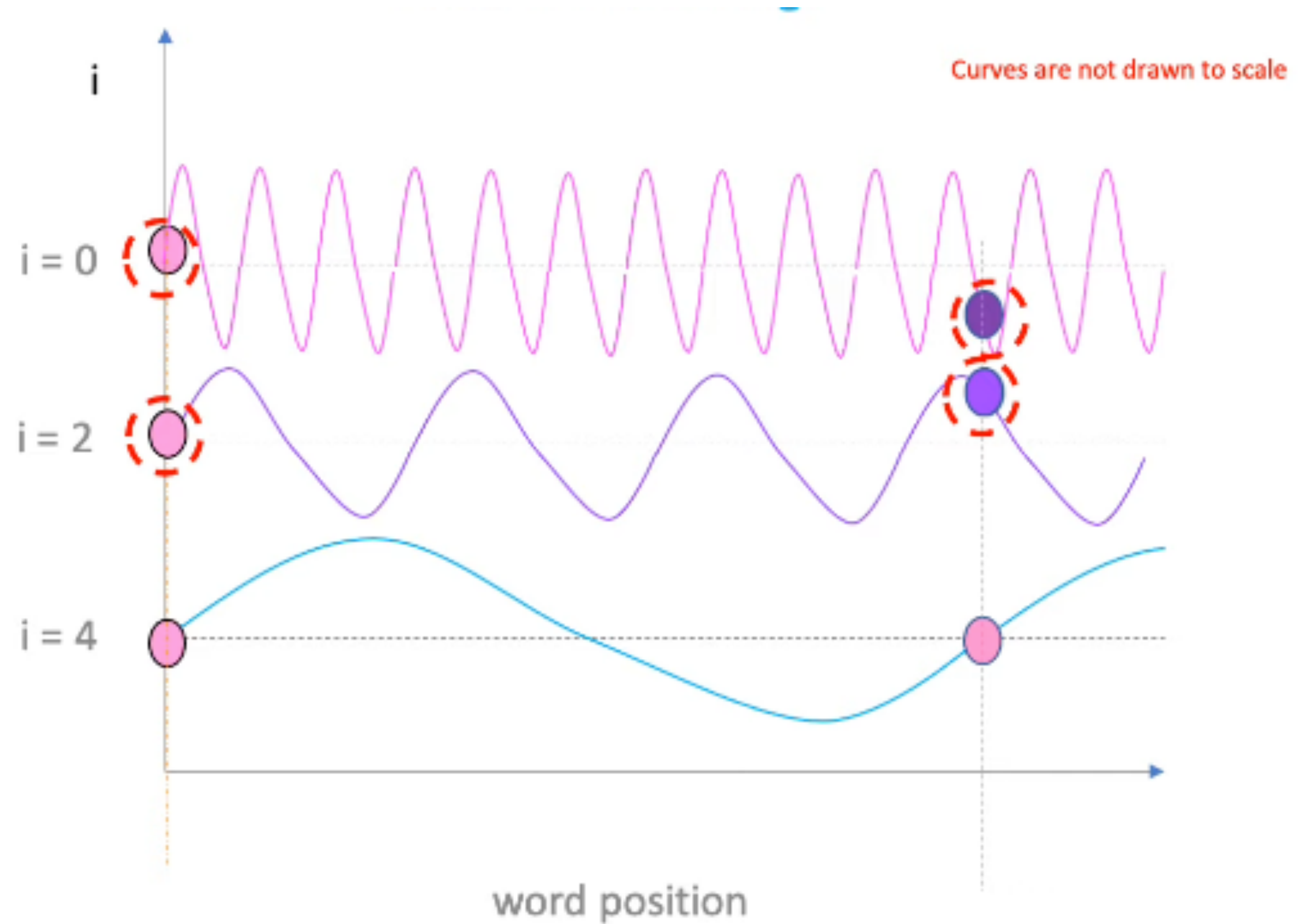
$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$



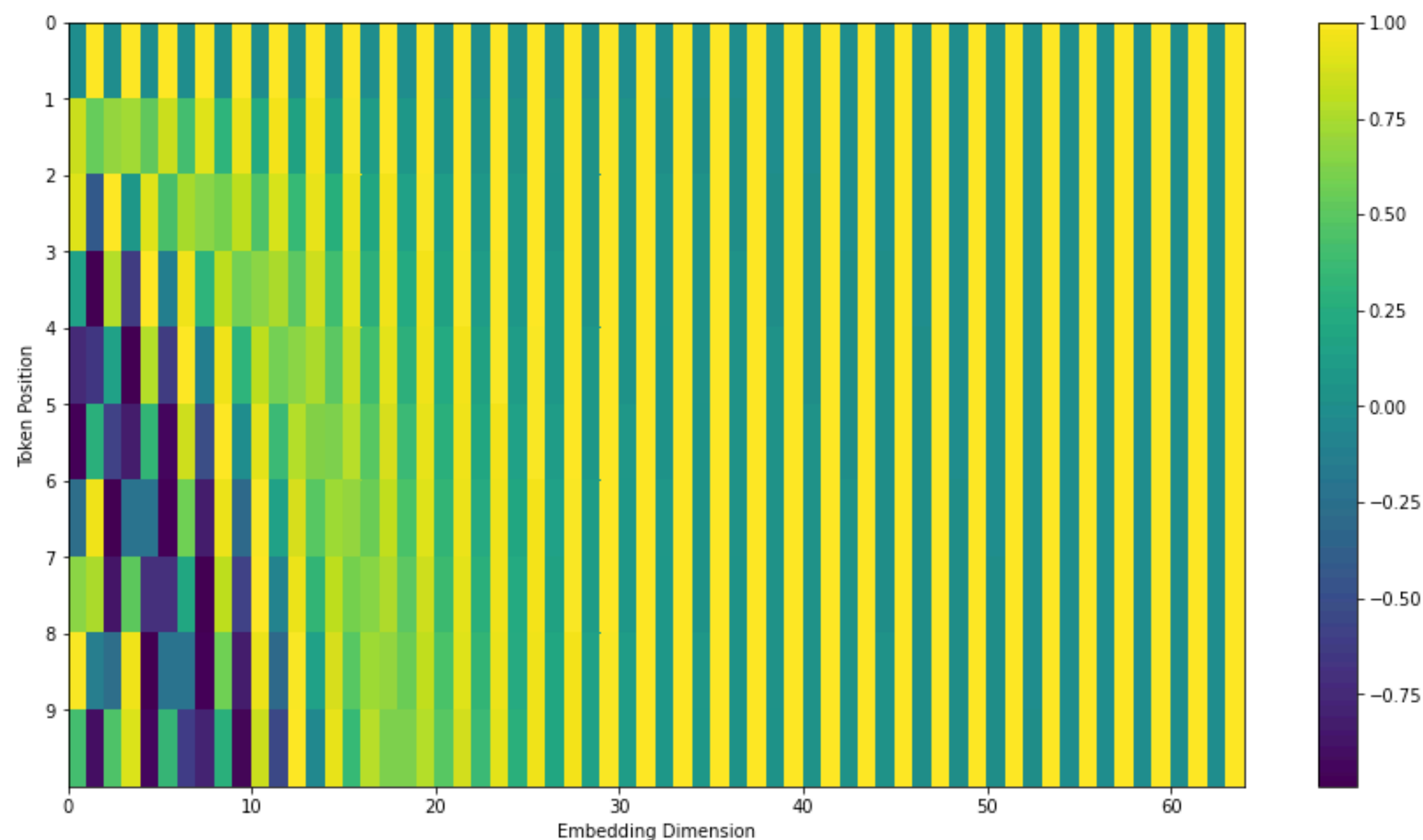
Positional encoding

$$PE_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

$$PE_{(pos, 2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$



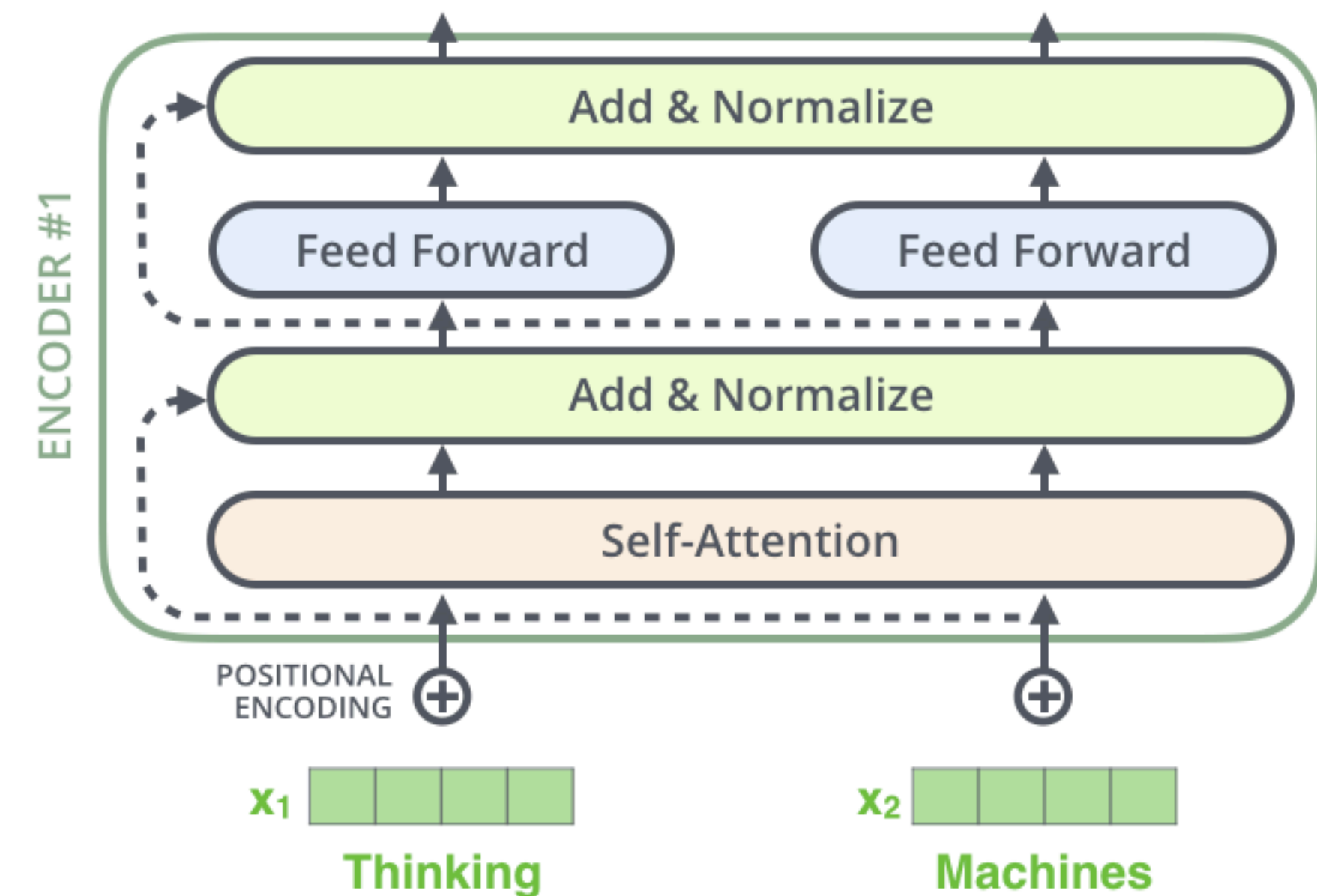
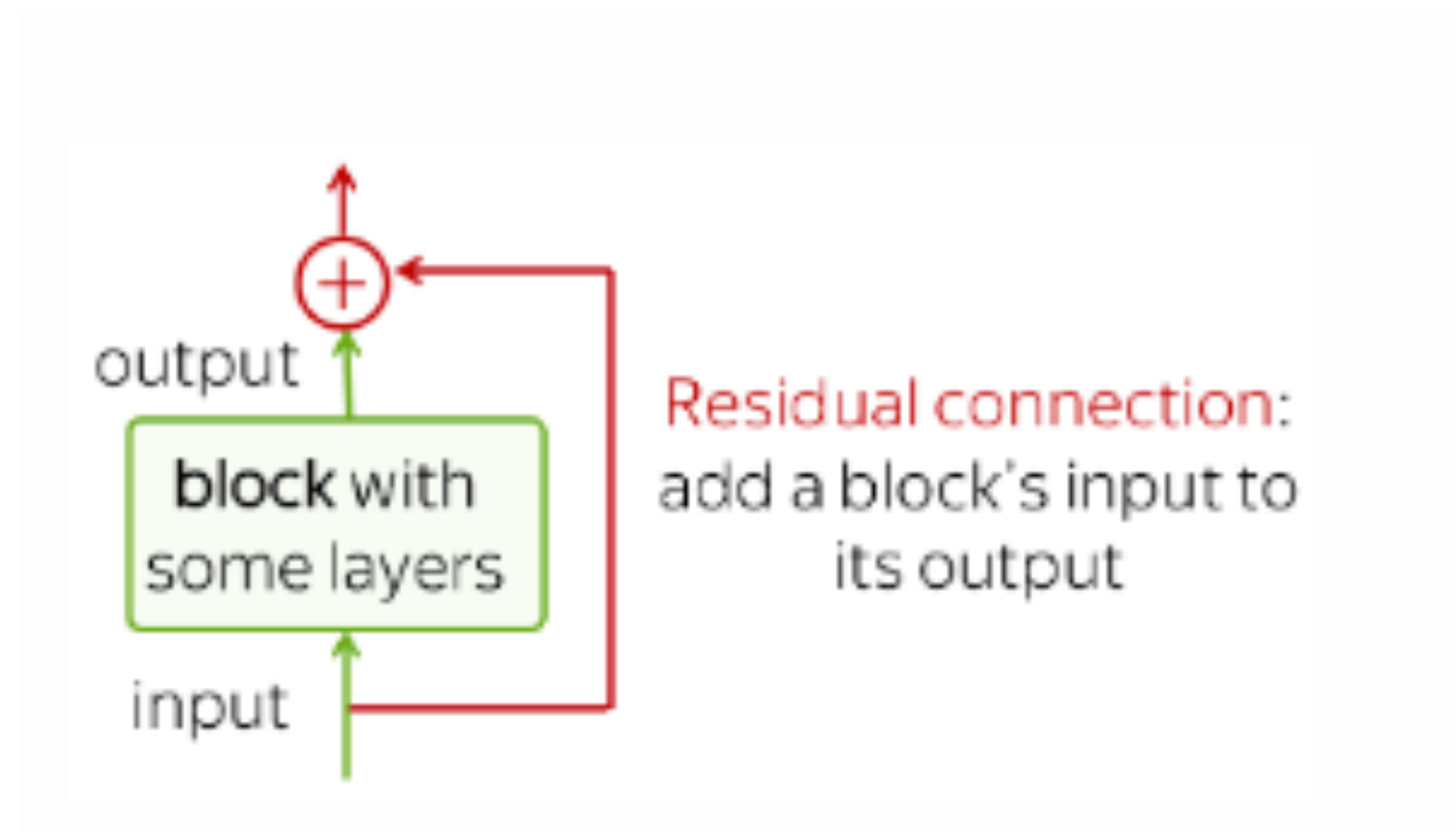
Positional encoding



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

Residual connection

directly passing "raw" embeddings to the next layer can actually be very helpful!



Layer Norm and add

Normalize each input in the batch to have zero mean and unity variance.

Scale and bias = trainable parameters and used after normalization to rescale output.

The difference between the two arrangements is illustrated in Figure 3-6.

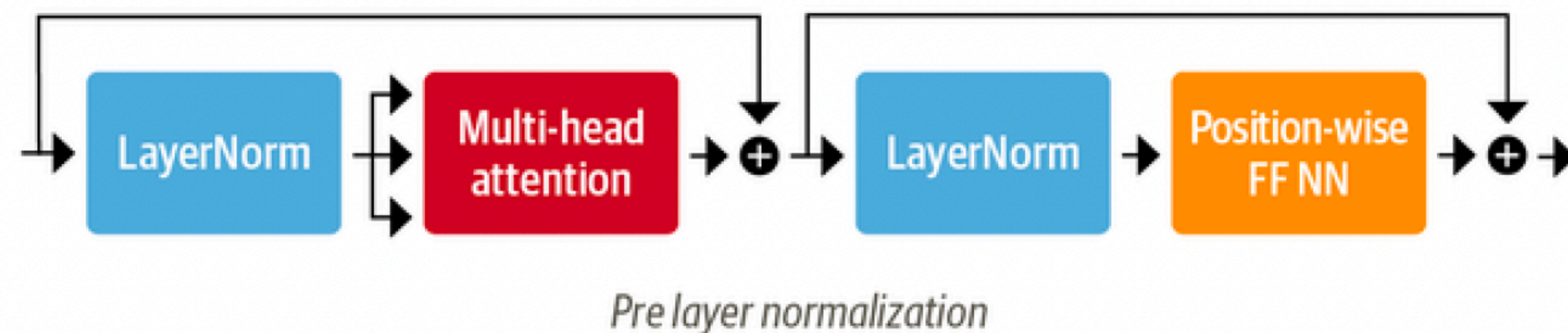
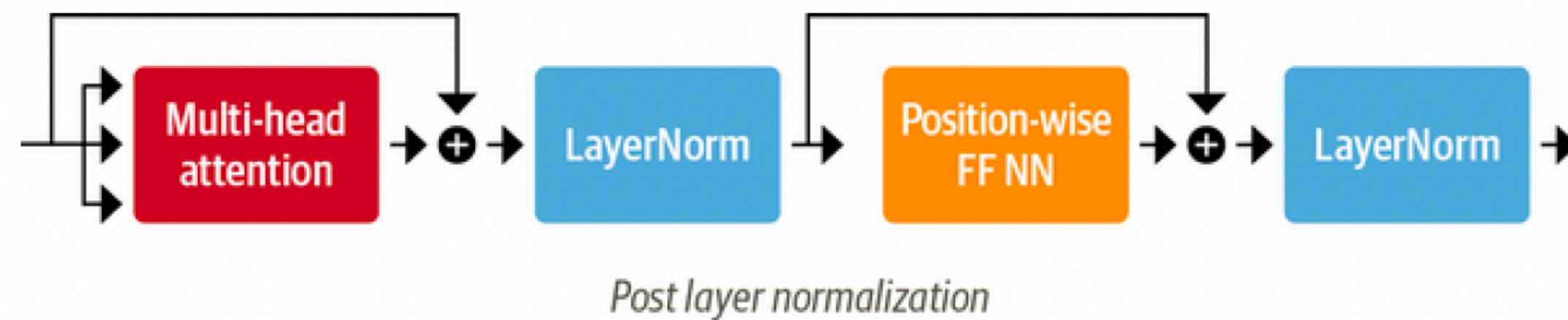
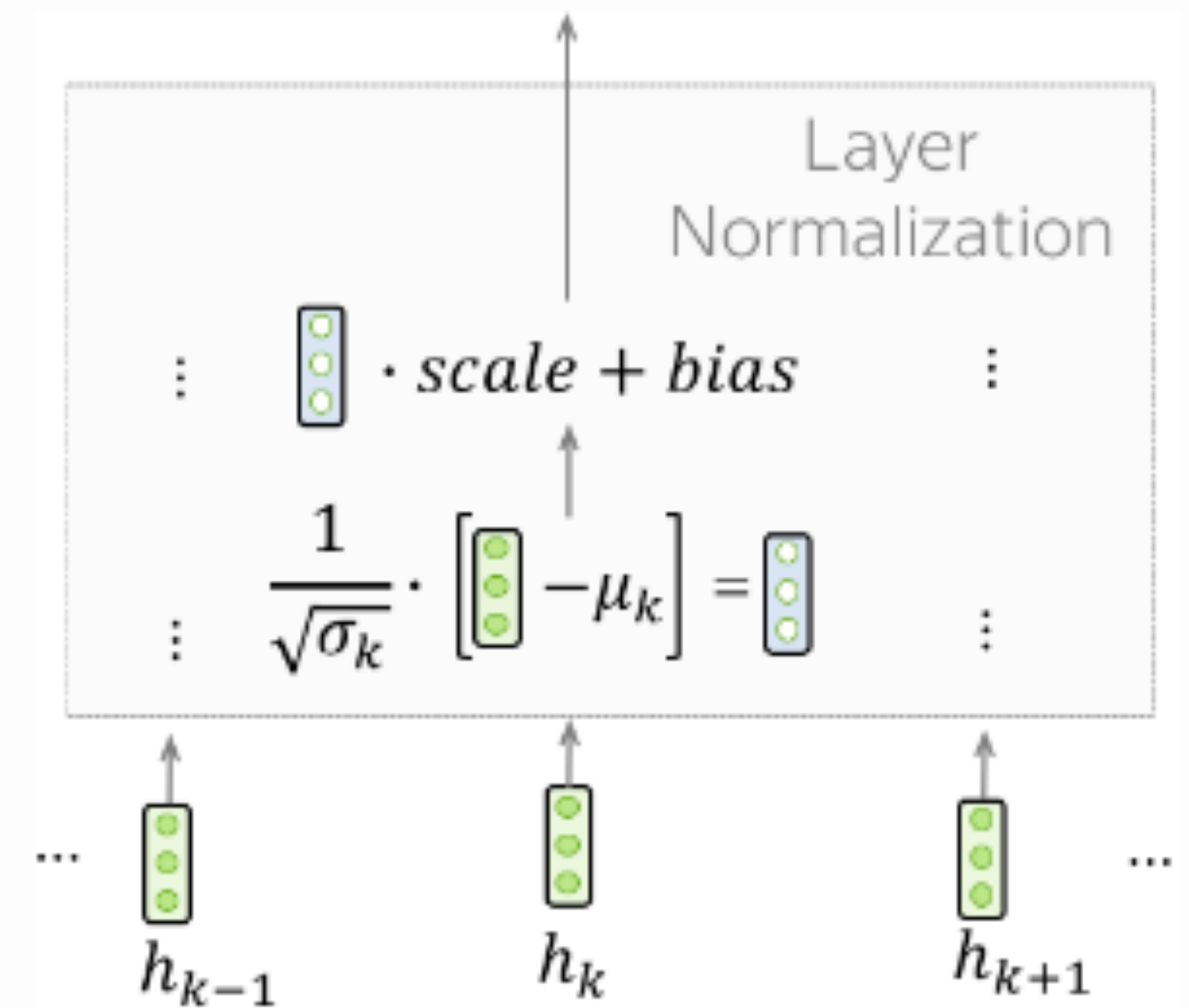


Figure 3-6. Different arrangements of layer normalization in a transformer encoder layer



Local vs global

Local VS global attention

global VS local - в глобал мы берем все стейты от encoder, в local предсказываем позиции токенов, на которые будем смотреть.

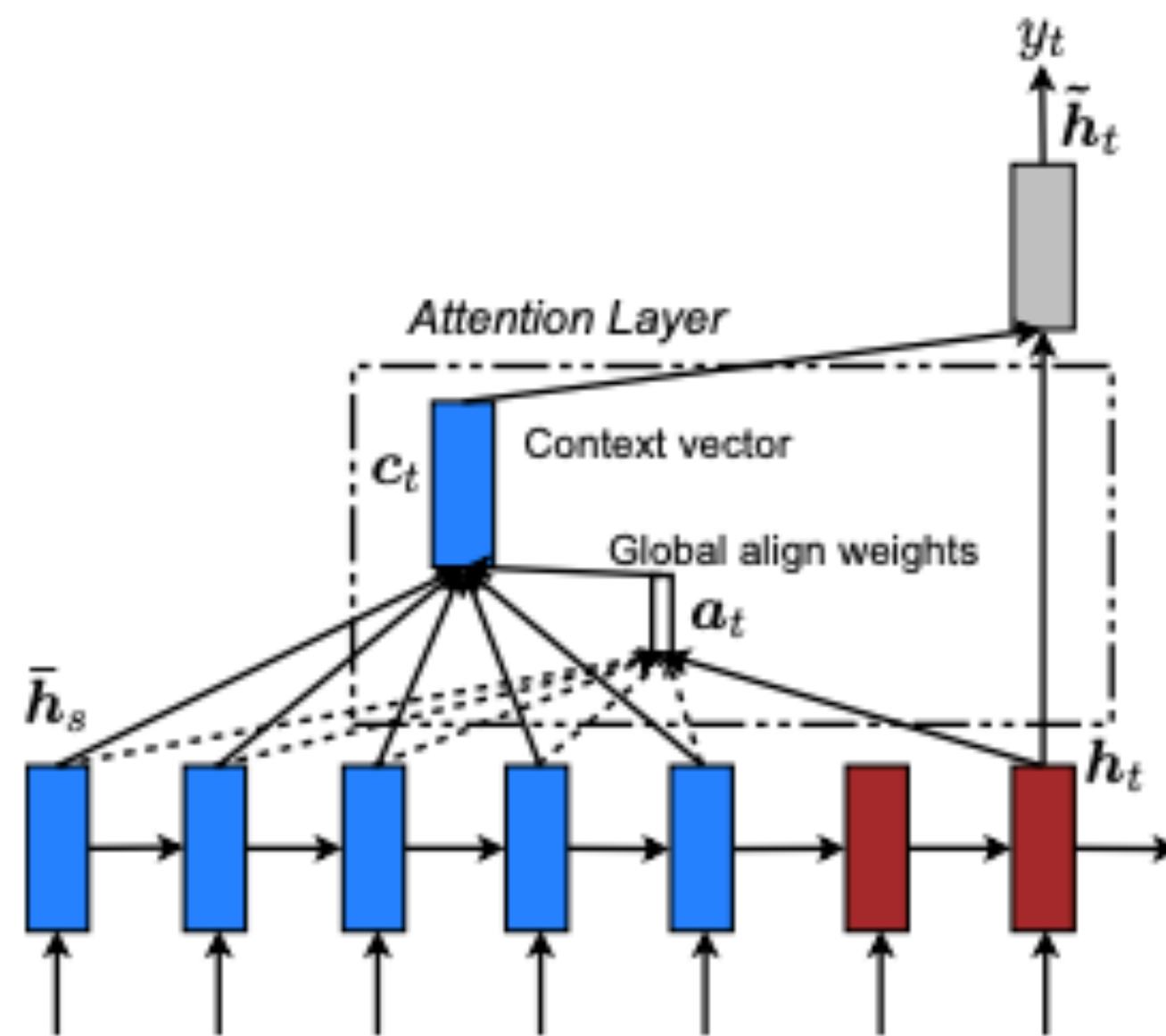


Figure 2: **Global attentional model** – at each time step t , the model infers a *variable-length* alignment weight vector a_t based on the current target state h_t and all source states \bar{h}_s . A global context vector c_t is then computed as the weighted average, according to a_t , over all the source states.

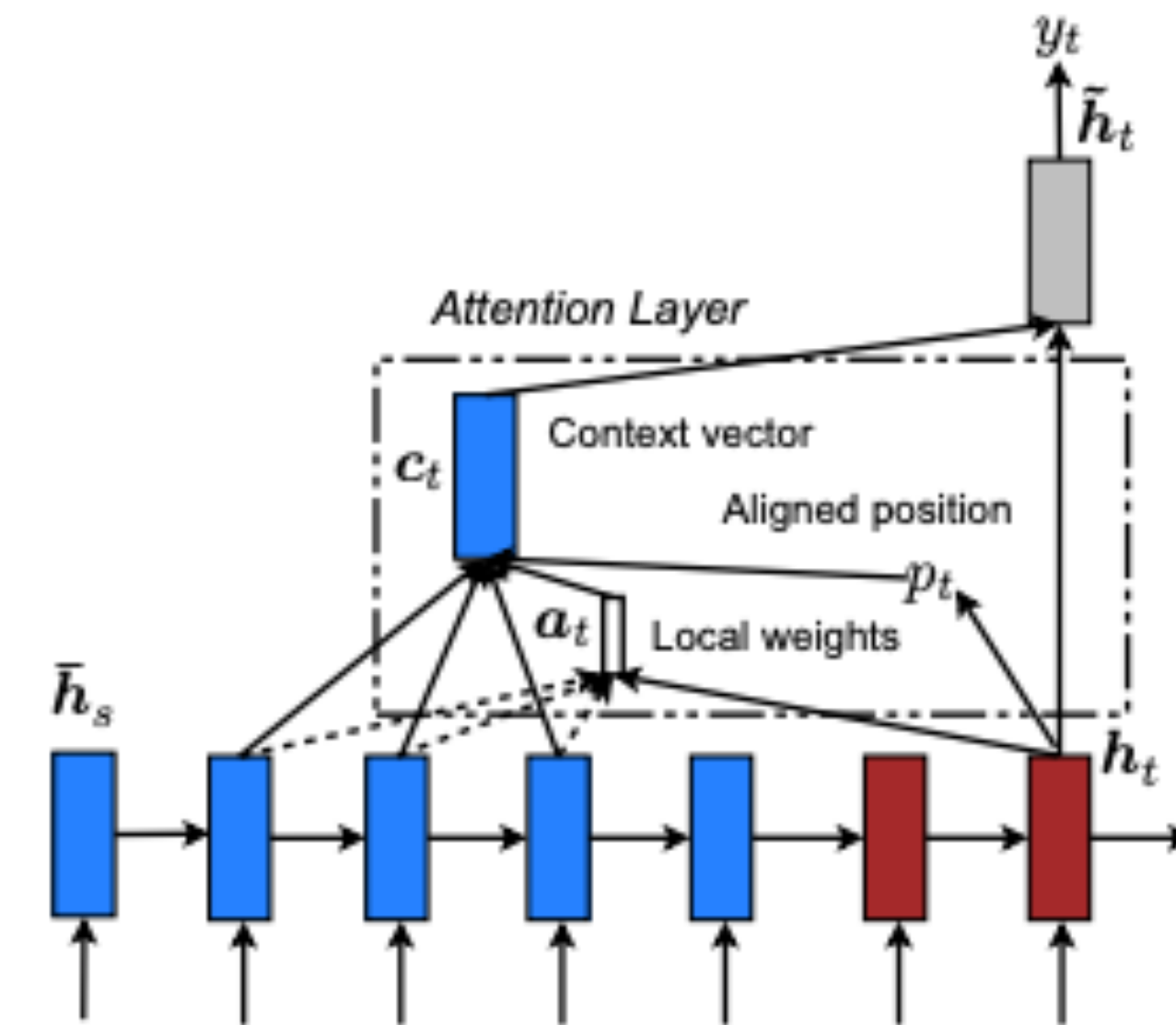


Figure 3: **Local attention model** – the model first predicts a single aligned position p_t for the current target word. A window centered around the source position p_t is then used to compute a context vector c_t , a weighted average of the source hidden states in the window. The weights a_t are inferred from the current target state h_t and those source states \bar{h}_s in the window.

Transformers

