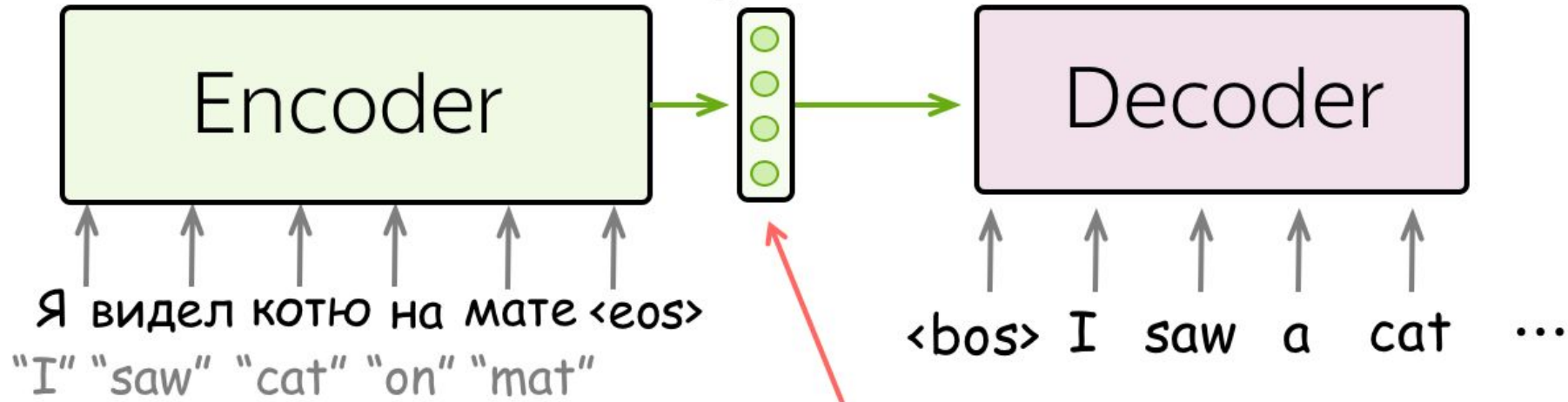


Transformer

Марк Блуменау, магистратура ИИ

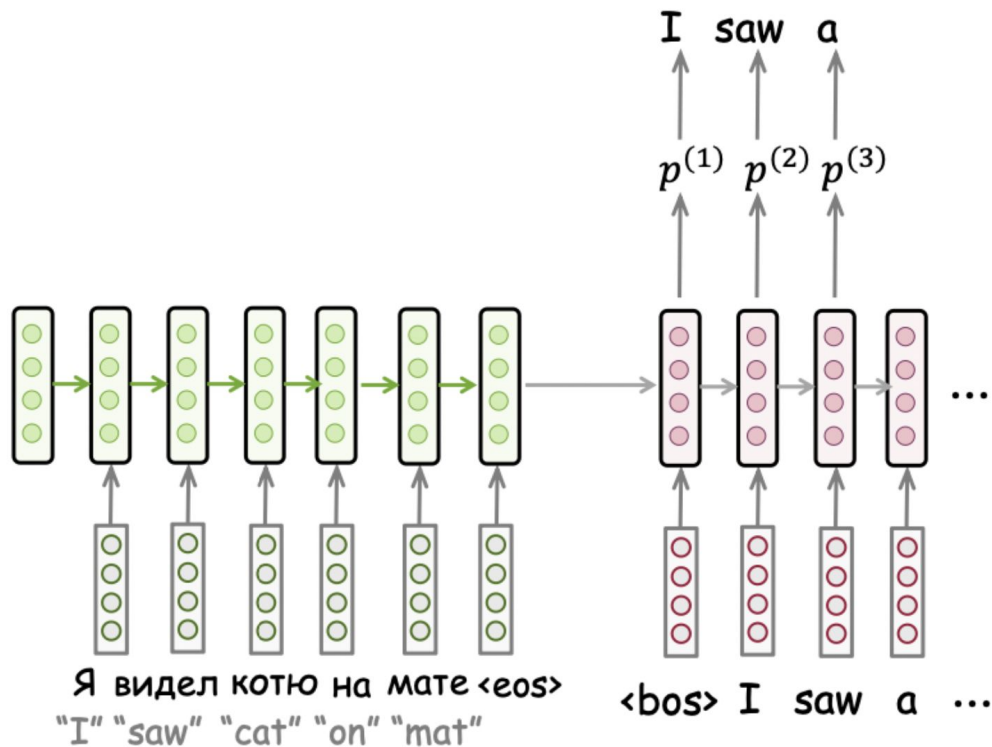
Почему RNN плохо?

We saw: encoder compresses the source into a single vector

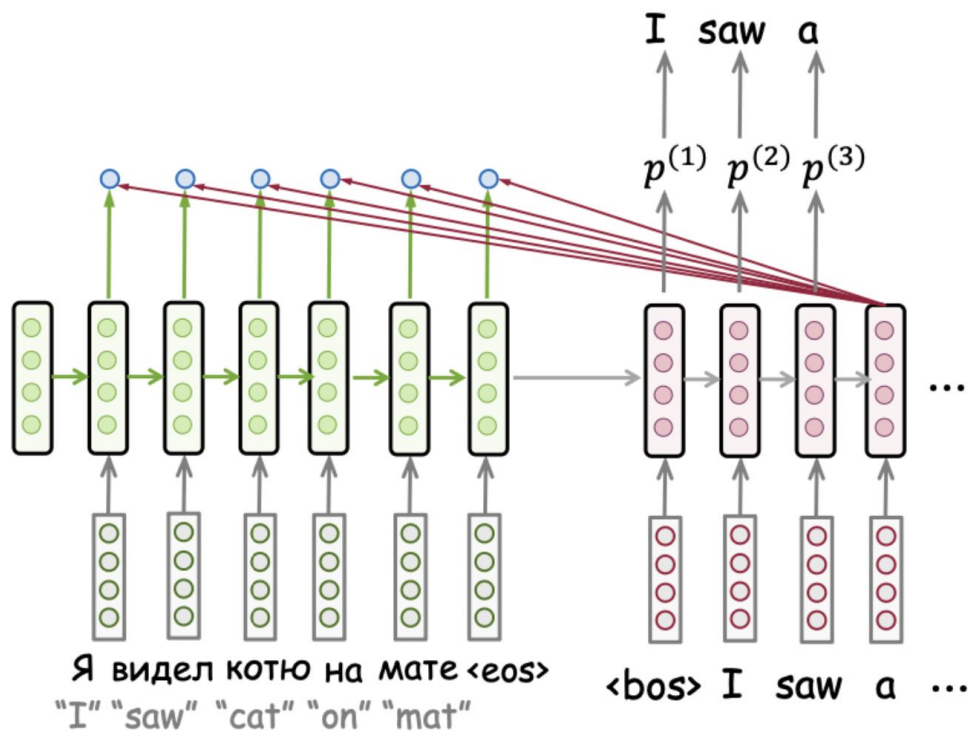


Problem: this is a bottleneck!

Attention



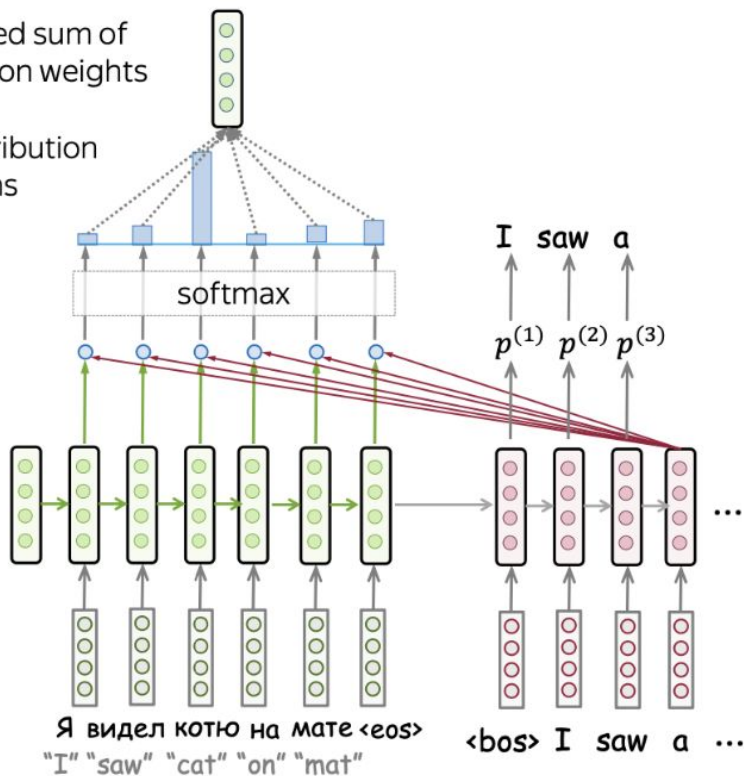
А хочется ведь так



Собираем

Attention output: weighted sum of encoder states with attention weights

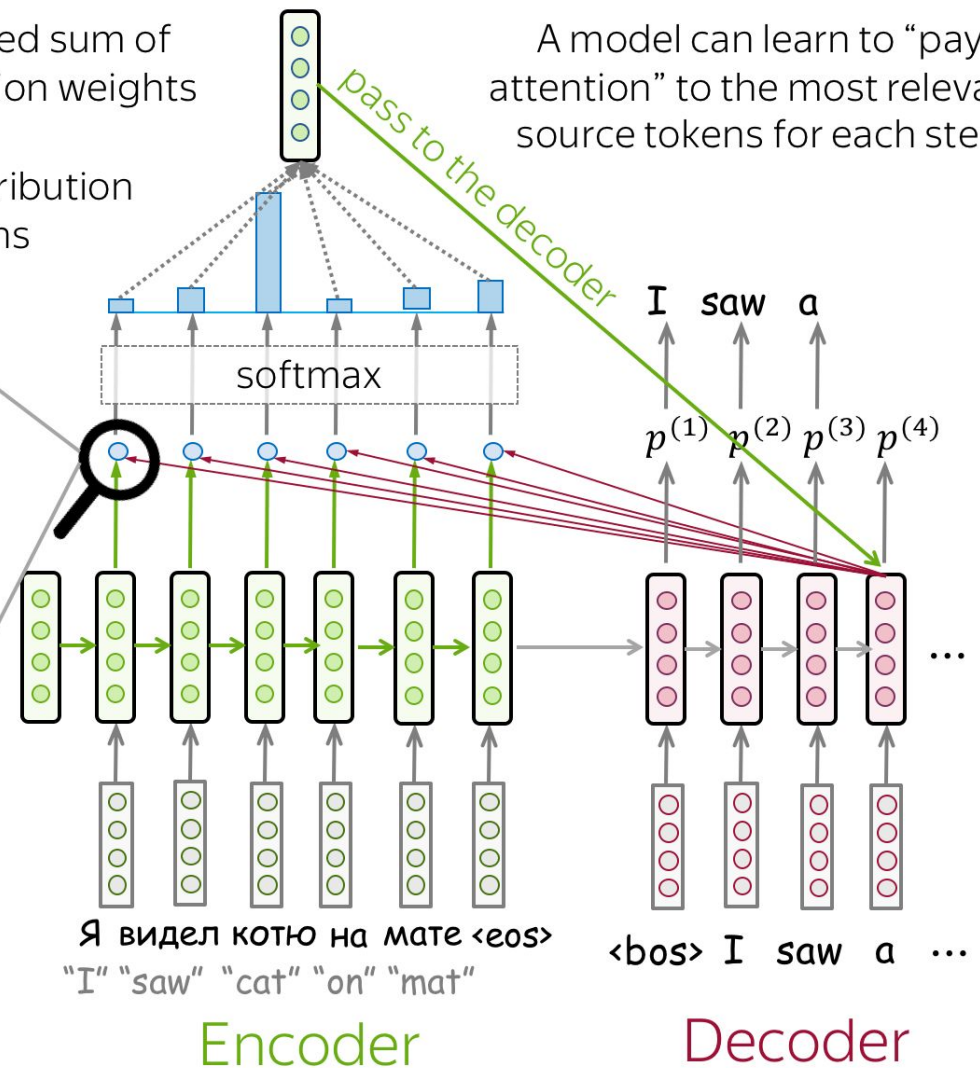
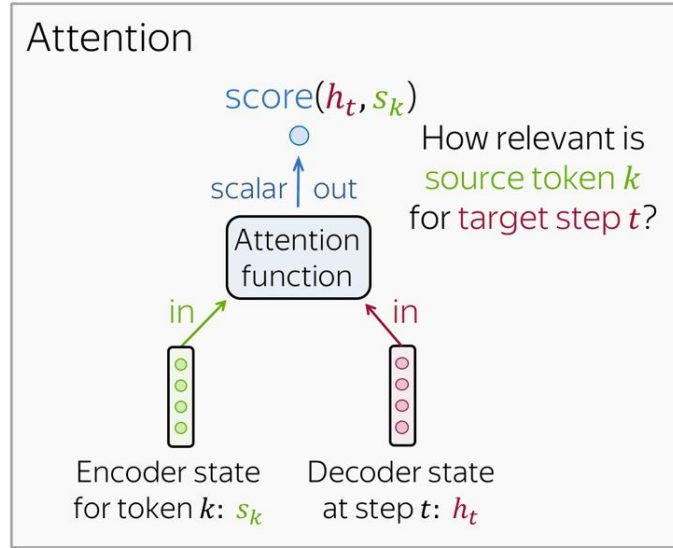
Attention weights: distribution over source tokens



Attention output: weighted sum of encoder states with attention weights

Attention weights: distribution over source tokens

A model can learn to “pay attention” to the most relevant source tokens for each step



Attention output

↑
(weighted
sum)

Attention weights

↑
(softmax)

Attention scores

↑

Attention input

$$\underset{\substack{\uparrow \\ \text{"source context for decoder step } t"}}}{c^{(t)}} = a_1^{(t)} s_1 + a_2^{(t)} s_2 + \dots + a_m^{(t)} s_m = \sum_{k=1}^m a_k^{(t)} s_k$$

$$\underset{\substack{\uparrow \\ \text{"attention weight for source token } k \text{ at decoder step } t"}}}{a_k^{(t)}} = \frac{\exp(\text{score}(h_t, s_k))}{\sum_{i=1}^m \exp(\text{score}(h_t, s_i))}, k = 1..m$$

$$\underset{\substack{\uparrow \\ \text{"How relevant is source token } k \text{ for target step } t?"}}}{\text{score}(h_t, s_k)}, k = 1..m$$

s_1, s_2, \dots, s_m

all encoder states

h_t

one decoder state

Badhau

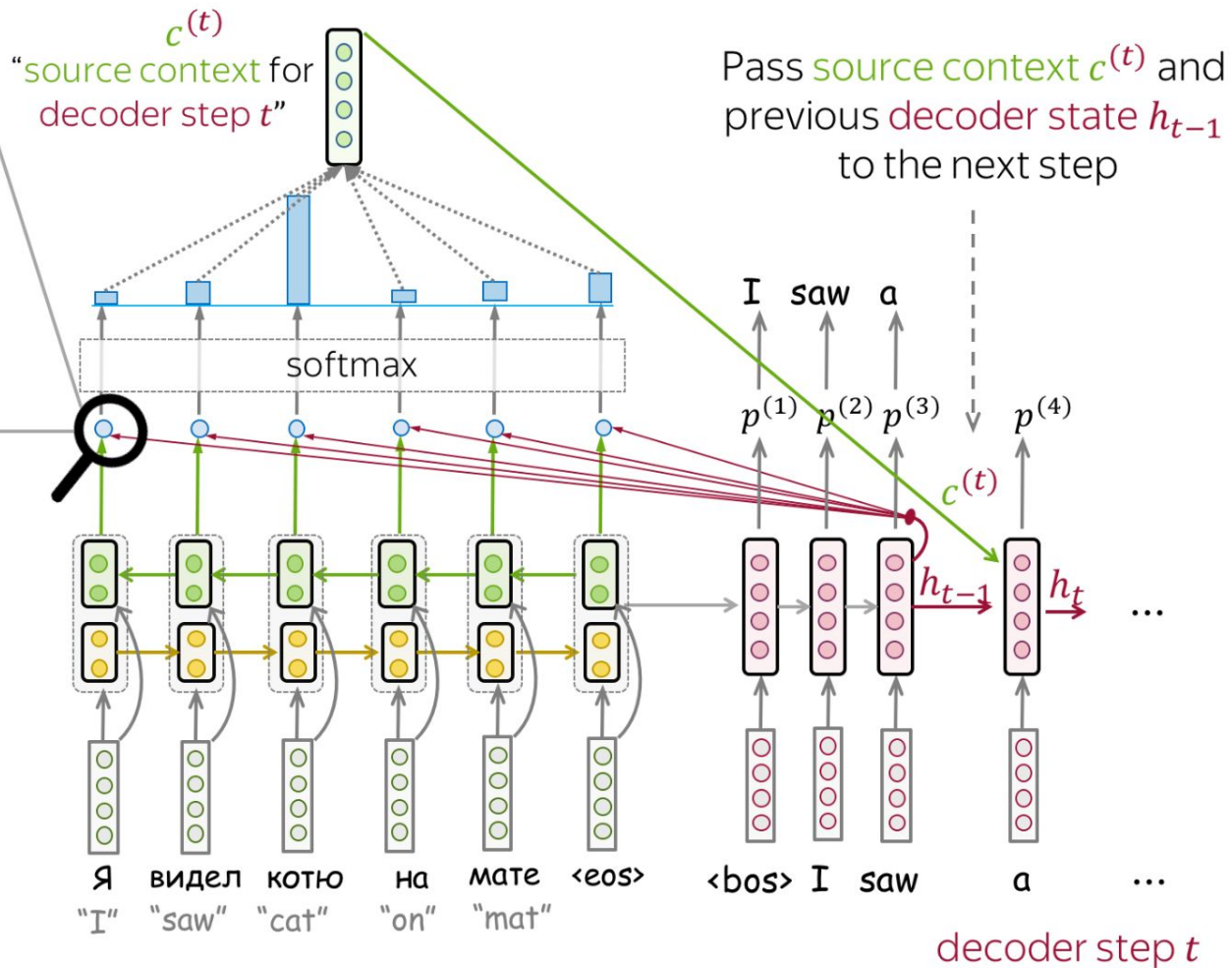
Multi-Layer Perceptron

$$w_2^T \times \tanh \left[W_1 \times \begin{bmatrix} h \\ s_k \end{bmatrix} \right]$$

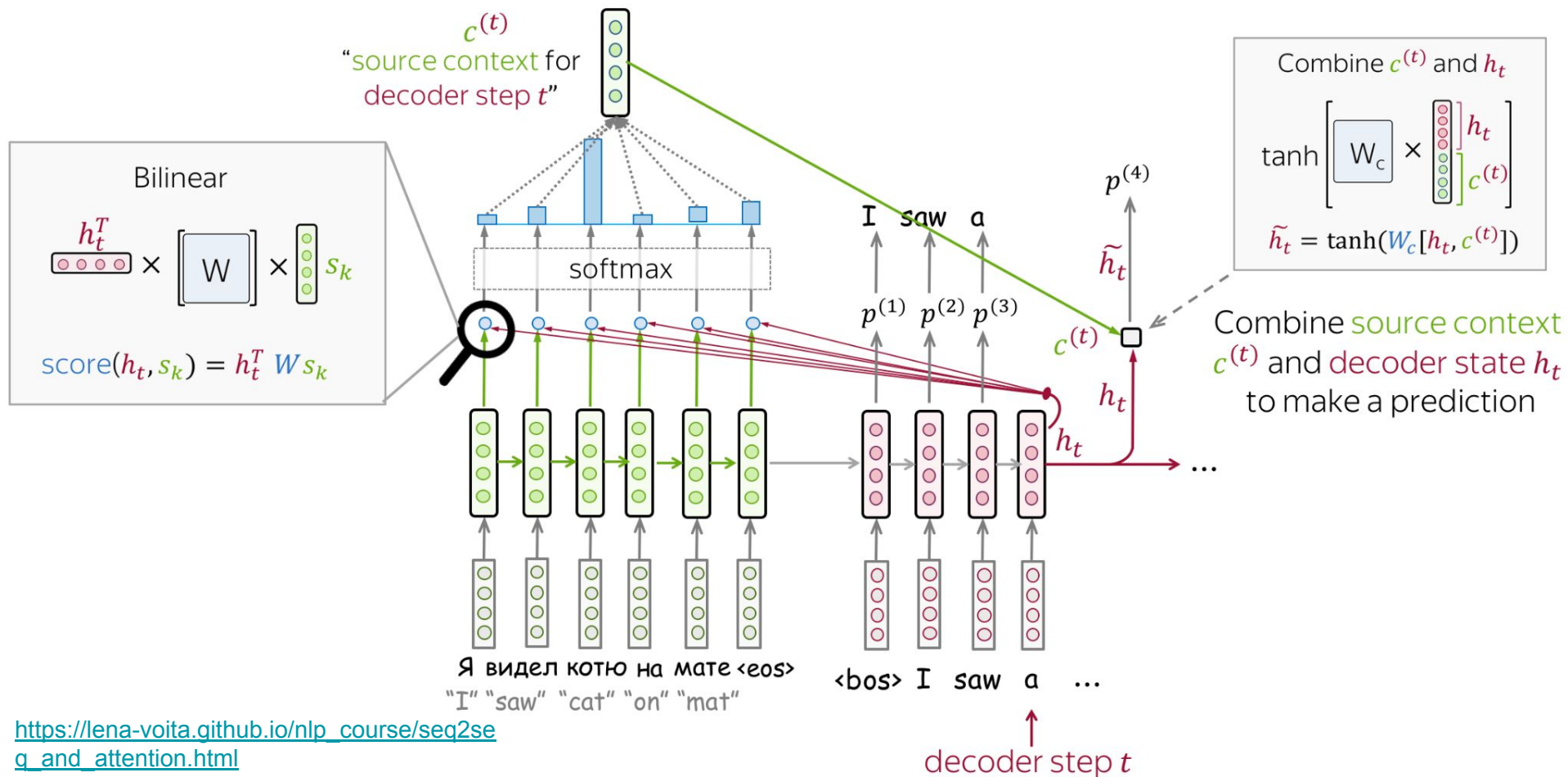
$$\text{score}(h, s_k) = w_2^T \cdot \tanh(W_1[h, s_k])$$

Bidirectional encoder

Concatenate states from
forward and backward RNNs

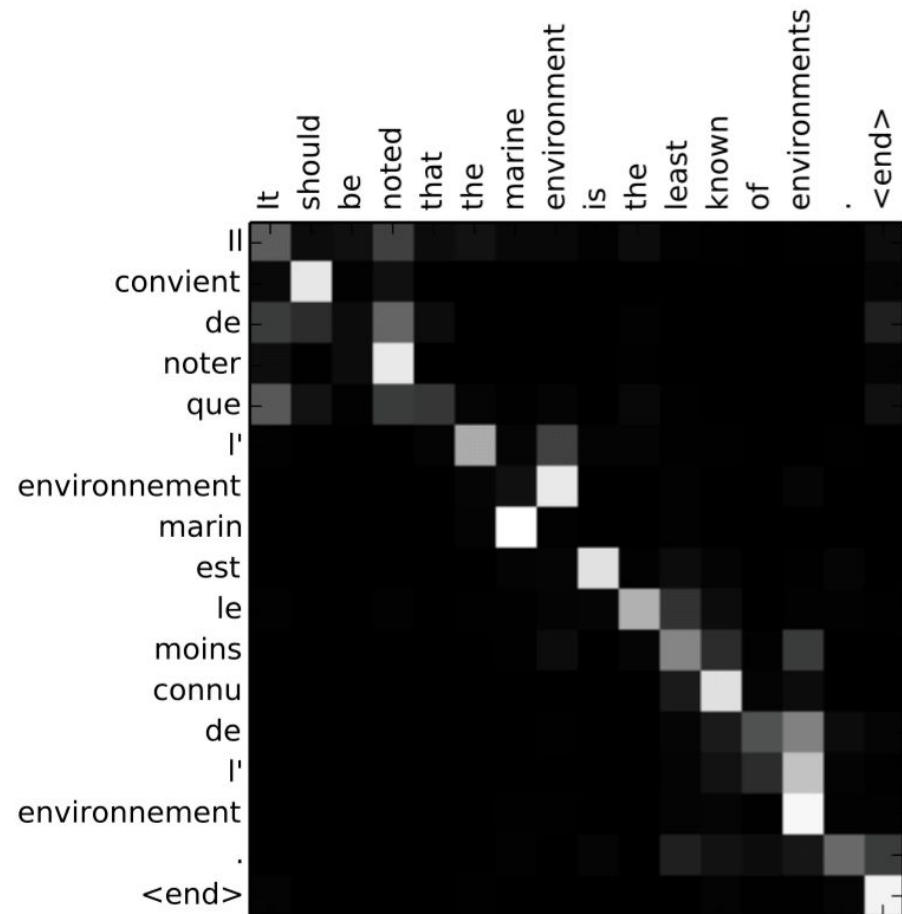
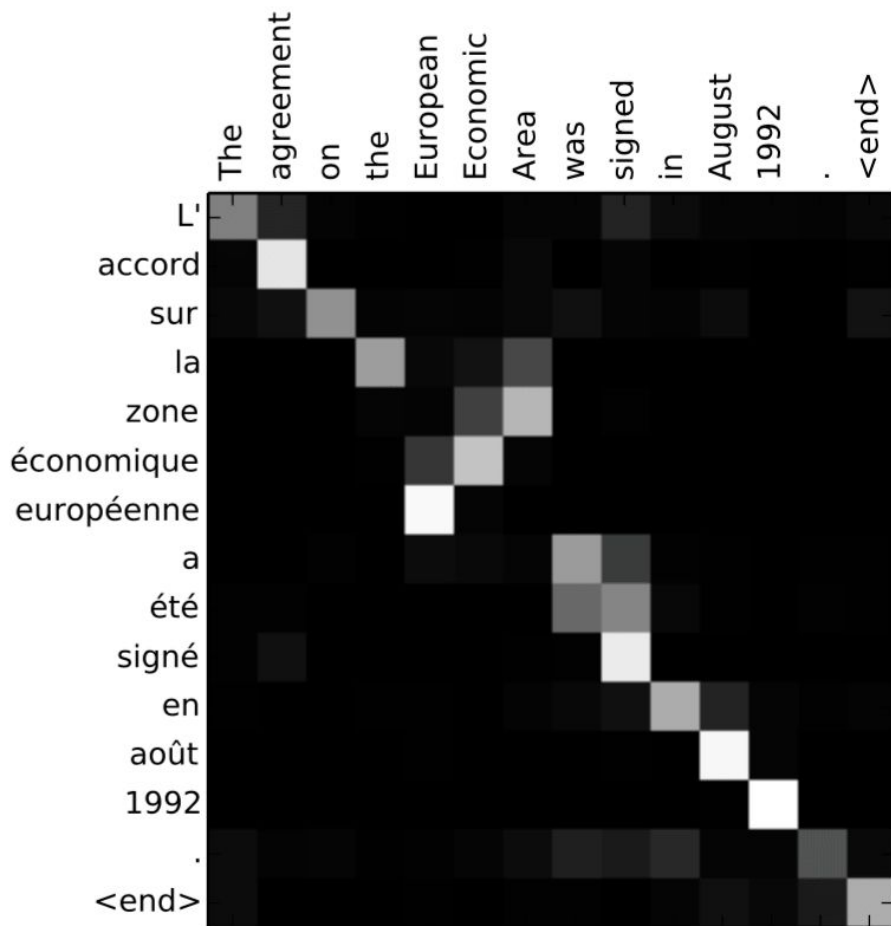


Luong



А что получается?

https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html



Ну давайте его везде теперь засунем

	Seq2seq without attention	Seq2seq with attention	Transformer
processing within encoder	RNN/CNN	RNN/CNN	attention
processing within decoder	RNN/CNN	RNN/CNN	attention
decoder-encoder interaction	static fixed- sized vector	attention	attention

Who is who

Encoder

Who is doing:

- all source tokens

What they are doing:

- look at each other
- update representations

repeat
N times

Decoder

Who is doing:

- target token at the current step

What they are doing:

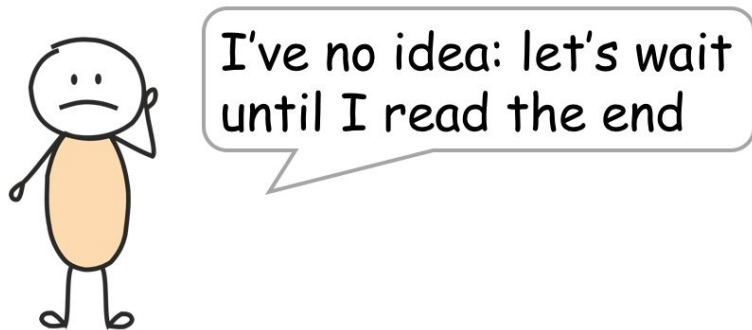
- looks at previous target tokens
- looks at source representations
- update representation

repeat
N times

Напоминание: RNN медленный

I arrived at the **bank** after crossing thestreet? ...river?

What does **bank** mean in this sentence?



I've no idea: let's wait
until I read the end

RNNs

$O(N)$ steps to process a
sentence with length N



I don't need to wait - I
see all words at once!

Transformer

Constant number of steps
to process any sentence

Стоп, а как посмотреть самим на себя???



[Зеркала > Зеркало в раме в античном
стиле золотое купить в интернет-магазине](#)

Each vector receives three representations (“roles”)

$$\begin{bmatrix} W_Q \end{bmatrix} \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{blue} \\ \text{blue} \\ \text{blue} \end{bmatrix}$$

Query: vector **from** which the attention is looking

“Hey there, do you have this information?”

$$\begin{bmatrix} W_K \end{bmatrix} \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{yellow} \\ \text{yellow} \\ \text{yellow} \end{bmatrix}$$

Key: vector **at** which the query looks to compute weights

“Hi, I have this information – give me a large weight!”

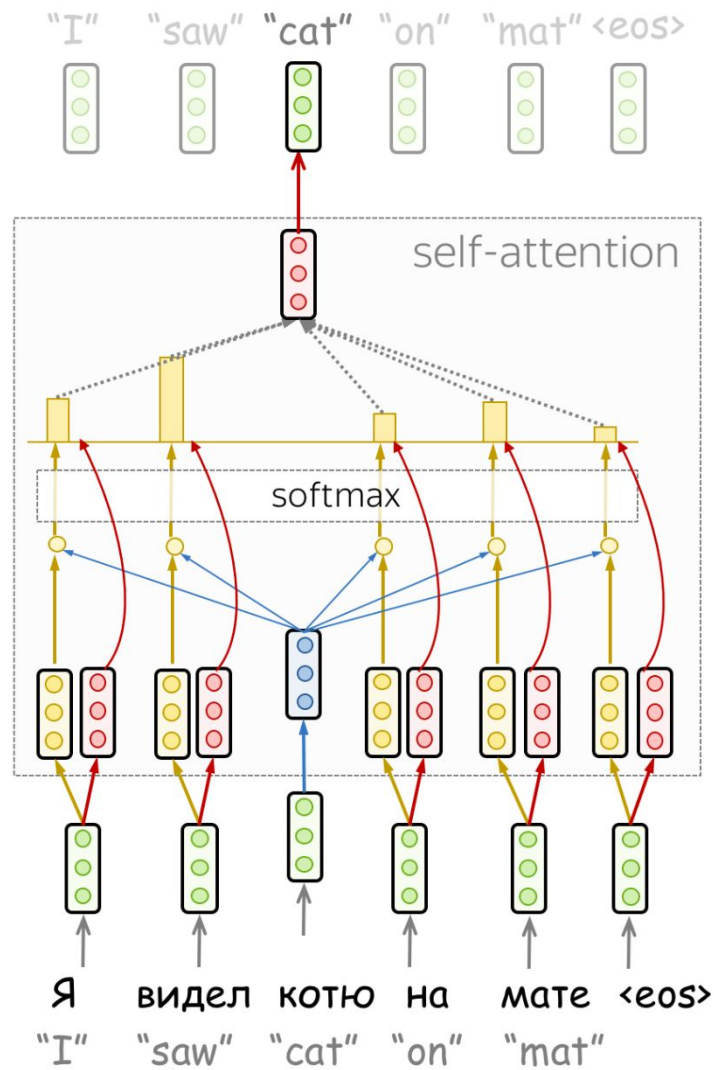
$$\begin{bmatrix} W_V \end{bmatrix} \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{red} \\ \text{red} \\ \text{red} \end{bmatrix}$$

Value: their weighted sum is attention output

“Here’s the information I have!”

$$\underset{\substack{\text{from} \\ \text{to}}}{\text{Attention}(q, k, v)} = \overbrace{\text{softmax}\left(\frac{qk^T}{\sqrt{d_k}}\right)}^{\text{Attention weights}} v$$

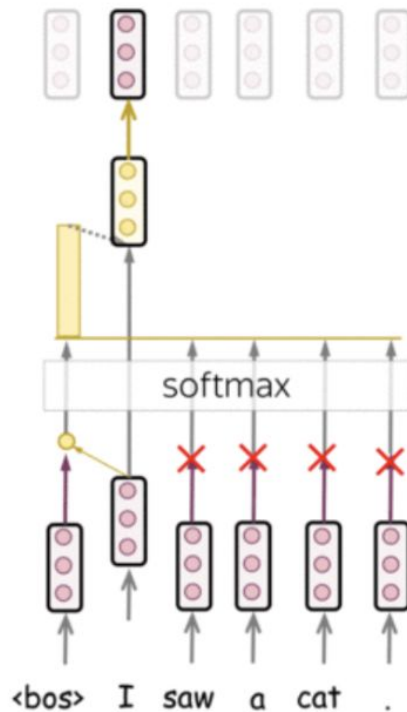
vector dimensionality of K, V



Смотреть в будущее оставим Ванге

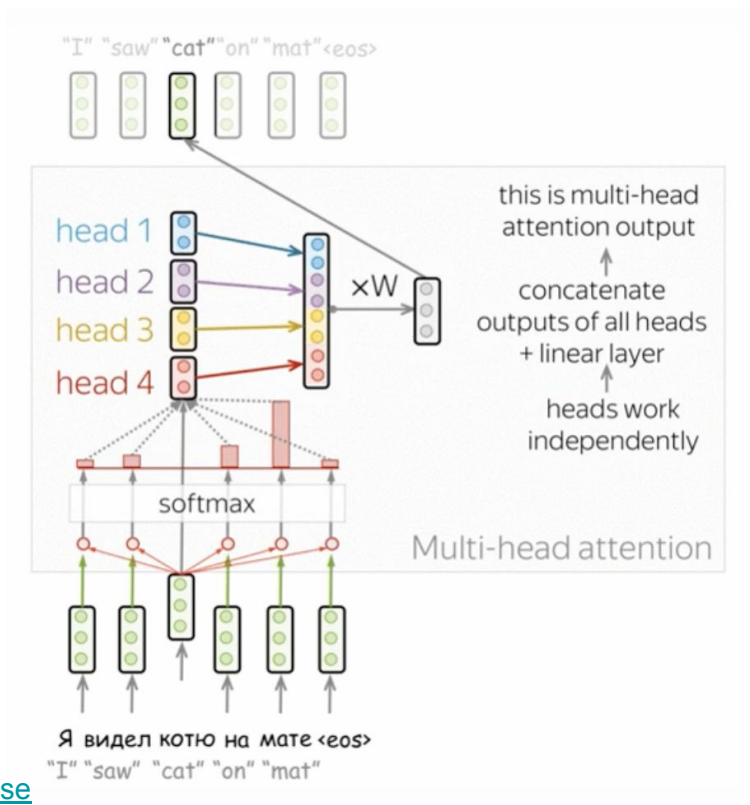
In the decoder, we forbid looking at future tokens – we don't know them

Note: in training, decoder processes all target tokens at once – without masks, it would see future

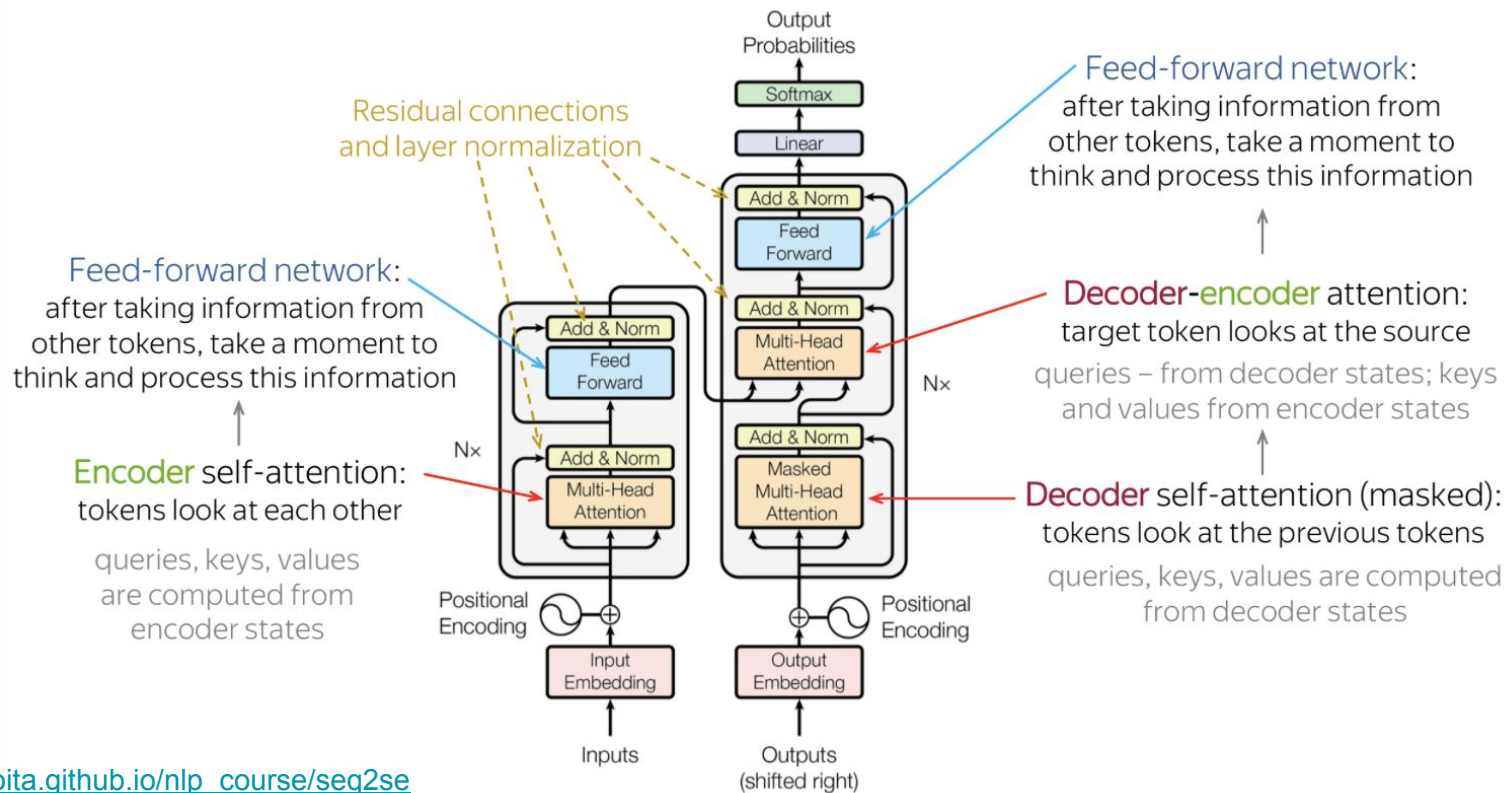


update token
representation
↑
gather context
↑
“look” at the
previous tokens
(future tokens are
masked out)

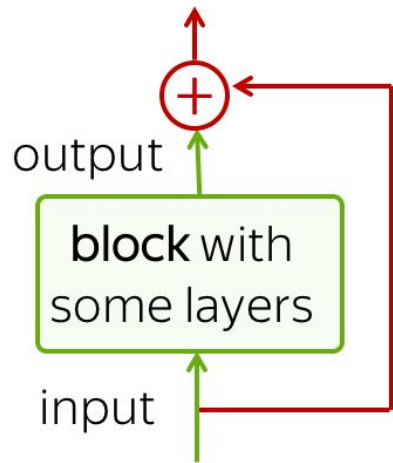
Одна голова хорошо, а две – лучше



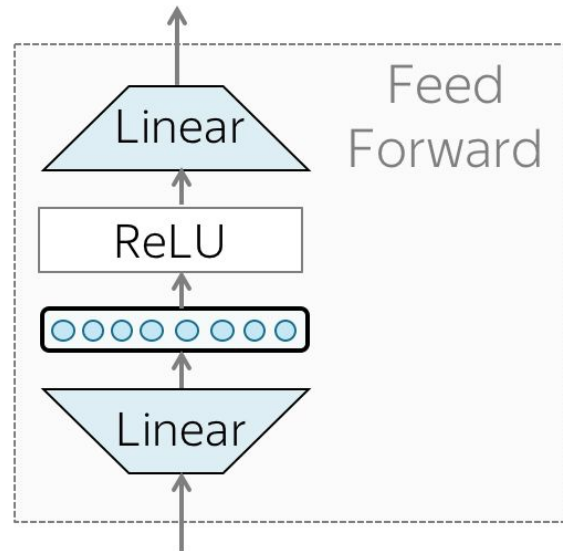
Ура давайте собрать лего



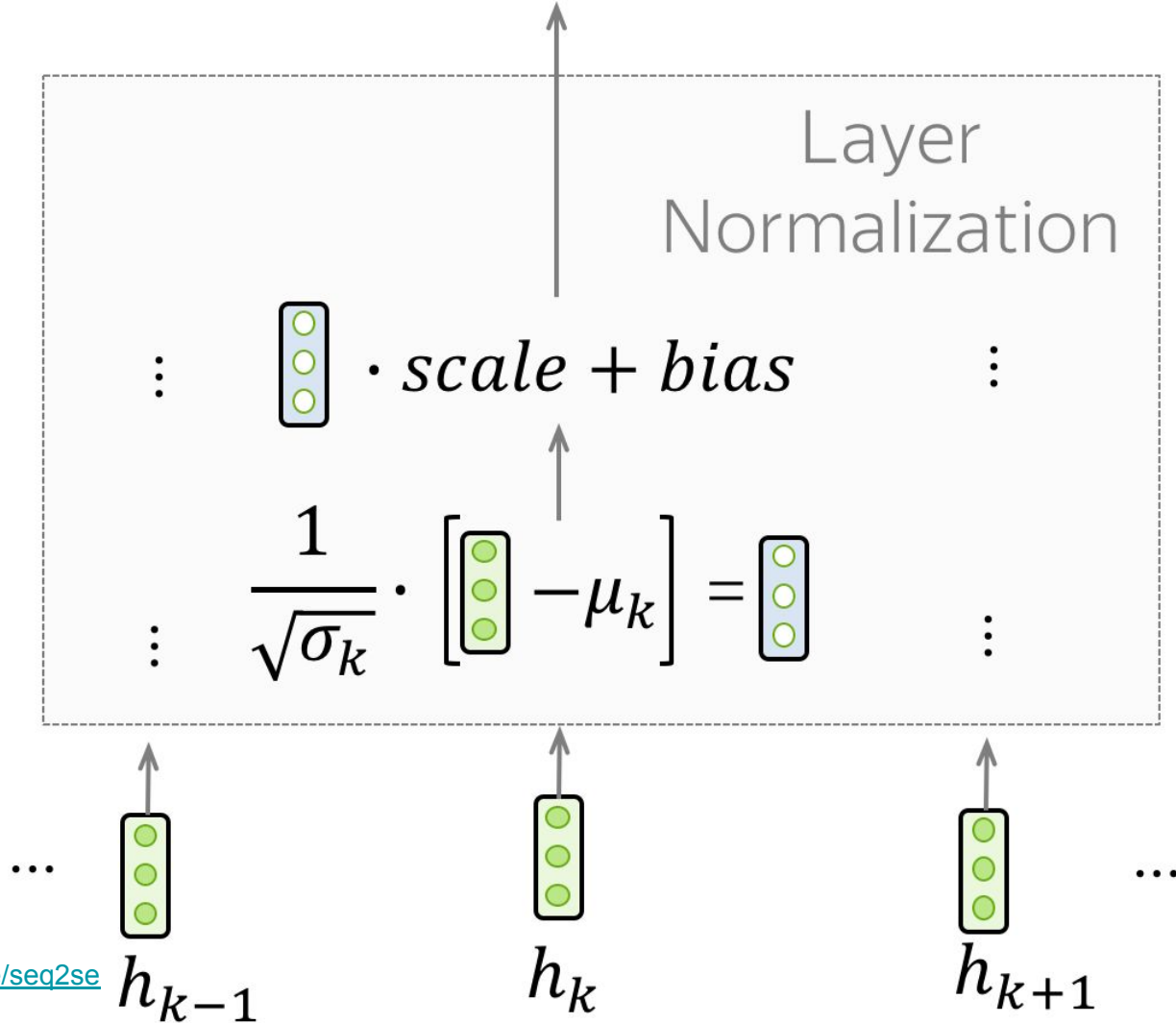
Эти покемоны нам известны



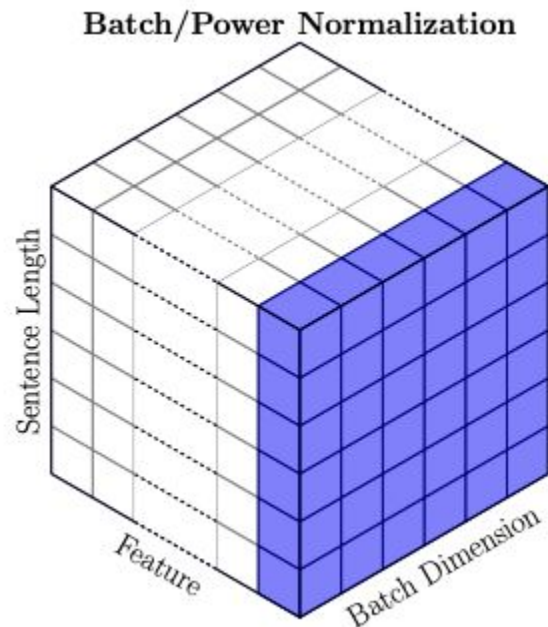
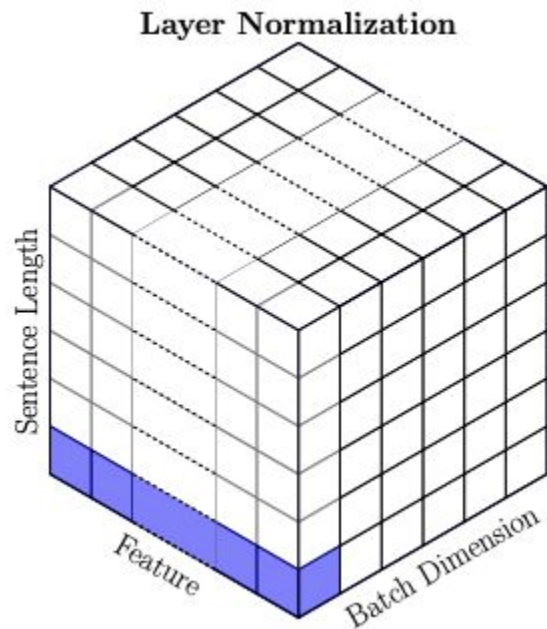
Residual connection:
add a block's input to
its output



А ЭТОТ НЕТ

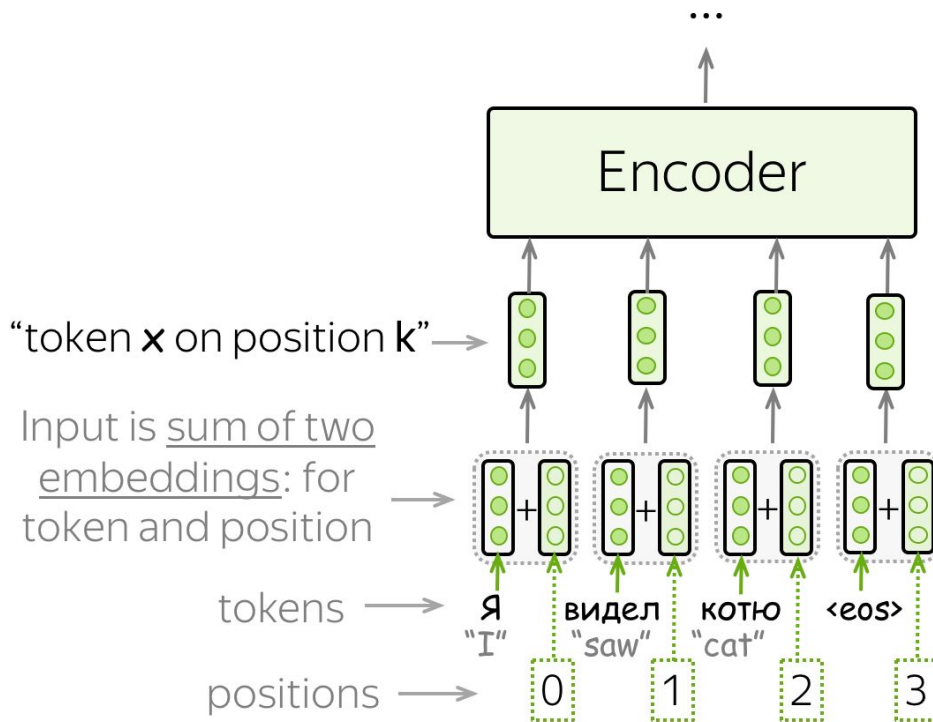


Просто ось поменяли



<https://stats.stackexchange.com/questions/474440/why-do-transformers-use-layer-norm-instead-of-batch-norm>

А позиция как работает?



Perplexity

$$L(y_{1:M}) = L(y_1, y_2, \dots, y_M) = \sum_{t=1}^M \log_2 p(y_t | y_{<t})$$

Log-likelihood of the text

$$Loss(y_{1:M}) = - \sum_{t=1}^M \log p(y_t | y_{<t})$$

Note: cross-entropy (our loss) is negative log-likelihood

$$Perplexity(y_{1:M}) = 2^{-\frac{1}{M} L(y_{1:M})}.$$

$$Perplexity(y_{1:M}) = 2^{-\frac{1}{M} L(y_{1:M})} = 2^{-\frac{1}{M} \sum_{t=1}^M \log_2 p(y_t | y_{1:t-1})} = 2^{-\frac{1}{M} \cdot M \cdot \log_2 \frac{1}{|V|}} = 2^{\log_2 |V|} = |V|.$$