

# NLP

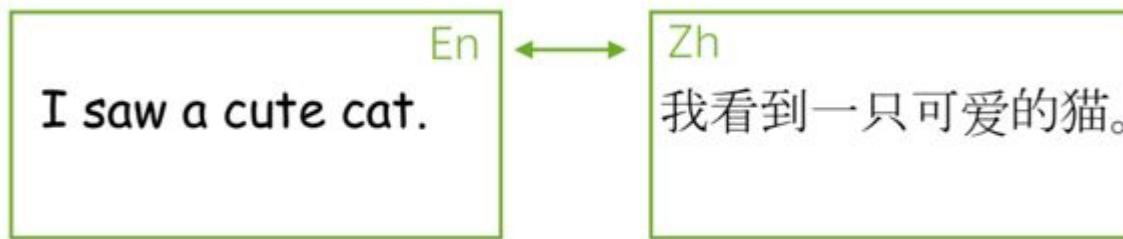
## Transformer

Основано на [NLP Course | For You](#)

# План

- seq2seq
- attention
- transformer
- BERT
- attention performance
- finetuning

# Sequence-to-sequence



# Задача

## Human Translation

$$y^* = \arg \max_y p(y|x)$$

The “probability” is  
intuitive and is given  
by a human  
translator’s expertise

## Machine Translation

$$y' = \arg \max_y p(y|x, \theta)$$

model

parameters

Questions we need to answer

- **modeling**

How does the model  
for  $p(y|x, \theta)$  look like?

- **learning**

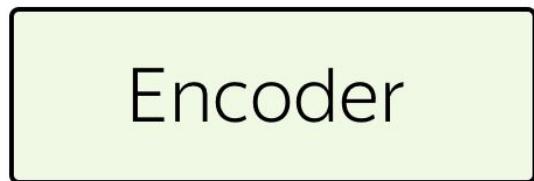
How to find  $\theta$ ?

- **search**

How to find  
the argmax?

# Encoder-decoder

Encoder builds a representation of the source and gives it to the decoder



Я видел котю на мате <eos>  
"I" "saw" "cat" "on" "mat"

Source sentence

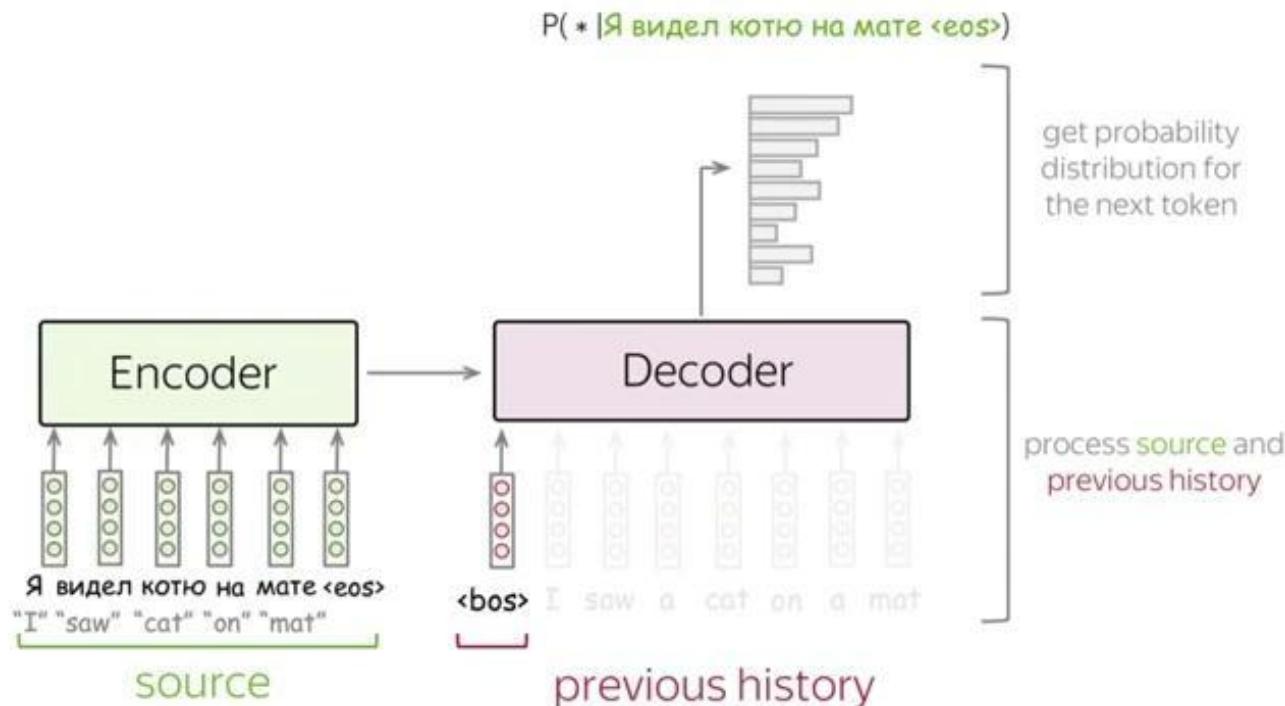
Target sentence

I saw a cat on a mat <eos>

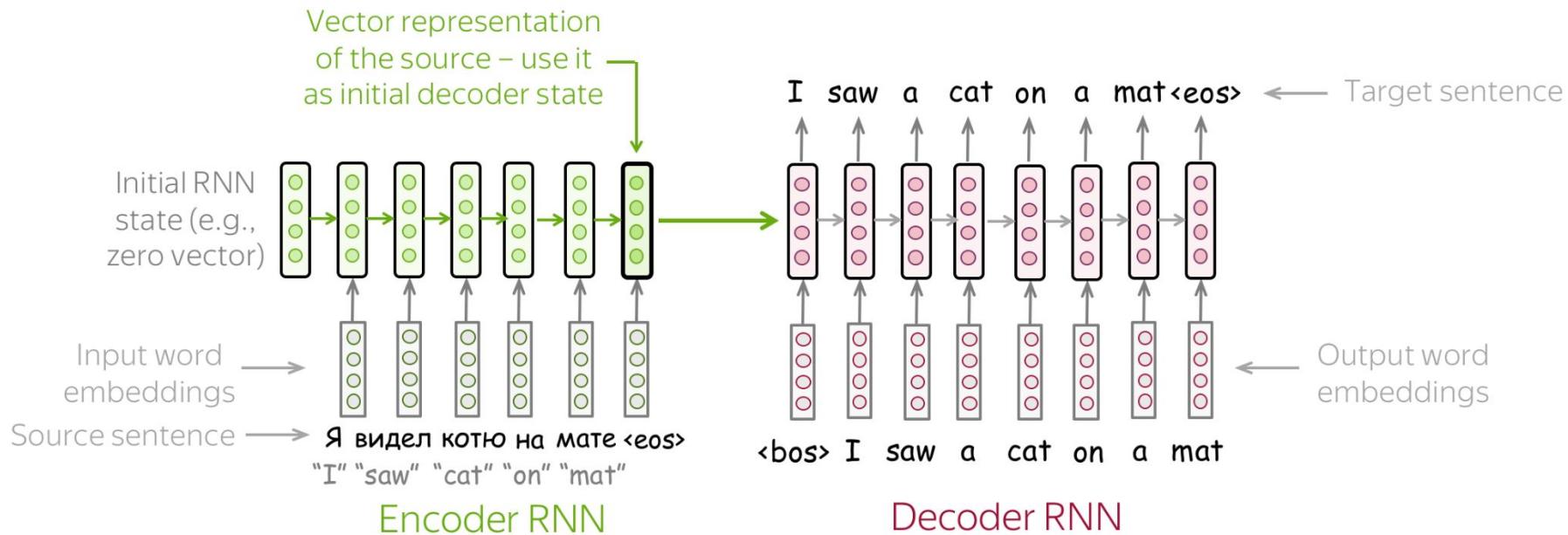


Decoder uses this source representation to generate the target sentence

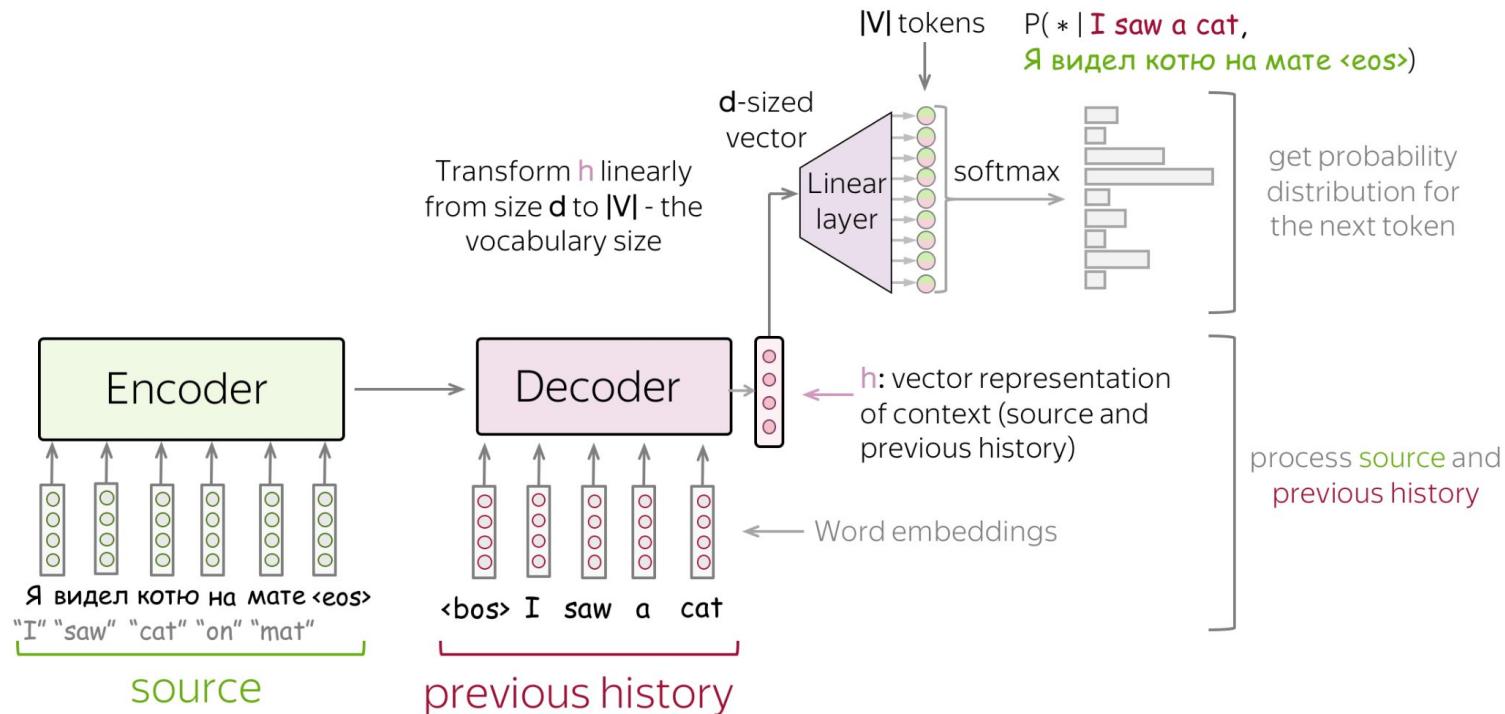
# Encoder-decoder



# RNN



# Как генерируем?



# Как генерируем?

Source sequence:

Я видел котю на мате <eos>  
"I" "saw" "cat" "on" "mat"

Target sequence:

I saw a **cat** on a mat <eos>

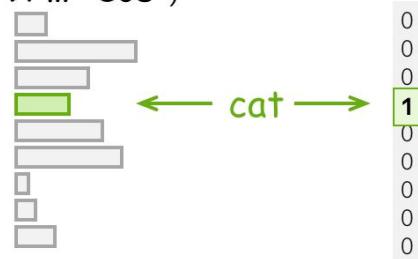
← one training example

previous tokens

we want the model  
to predict this

← one step for this example

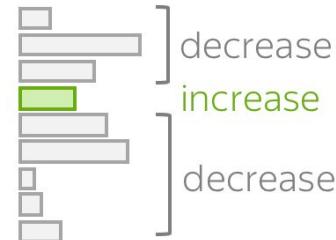
Model prediction:  $p(* | I \text{ saw } a,$   
 $\text{я} \dots \text{<eos>})$



Target



$$\text{Loss} = -\log(p(\text{cat})) \rightarrow \min$$



# Как найти argmax?

$$y' = \arg \max_y p(y|x) = \arg \max_y \prod_{t=1}^n p(y_t|y_{<t}, x)$$

How to find the argmax?

## Жадный подход

$$\arg \max_y \prod_{t=1}^n p(y_t | y_{<t}, x) \neq \prod_{t=1}^n \arg \max_{y_t} p(y_t | y_{<t}, x)$$

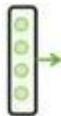
# Beam search

<bos>

Start with the begin of sentence token or with an empty sequence

# Once again генерация

Encoder: read source



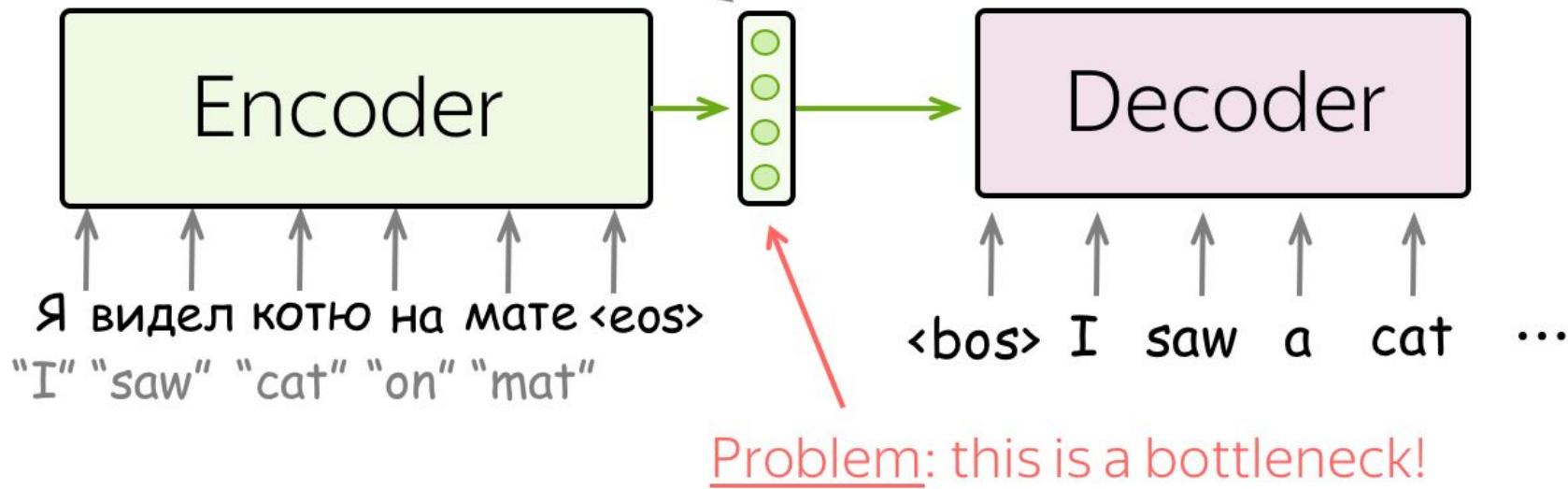
we are here

Source: Я видел котю на мате <eos>  
"I" "saw" "cat" "on" "mat"

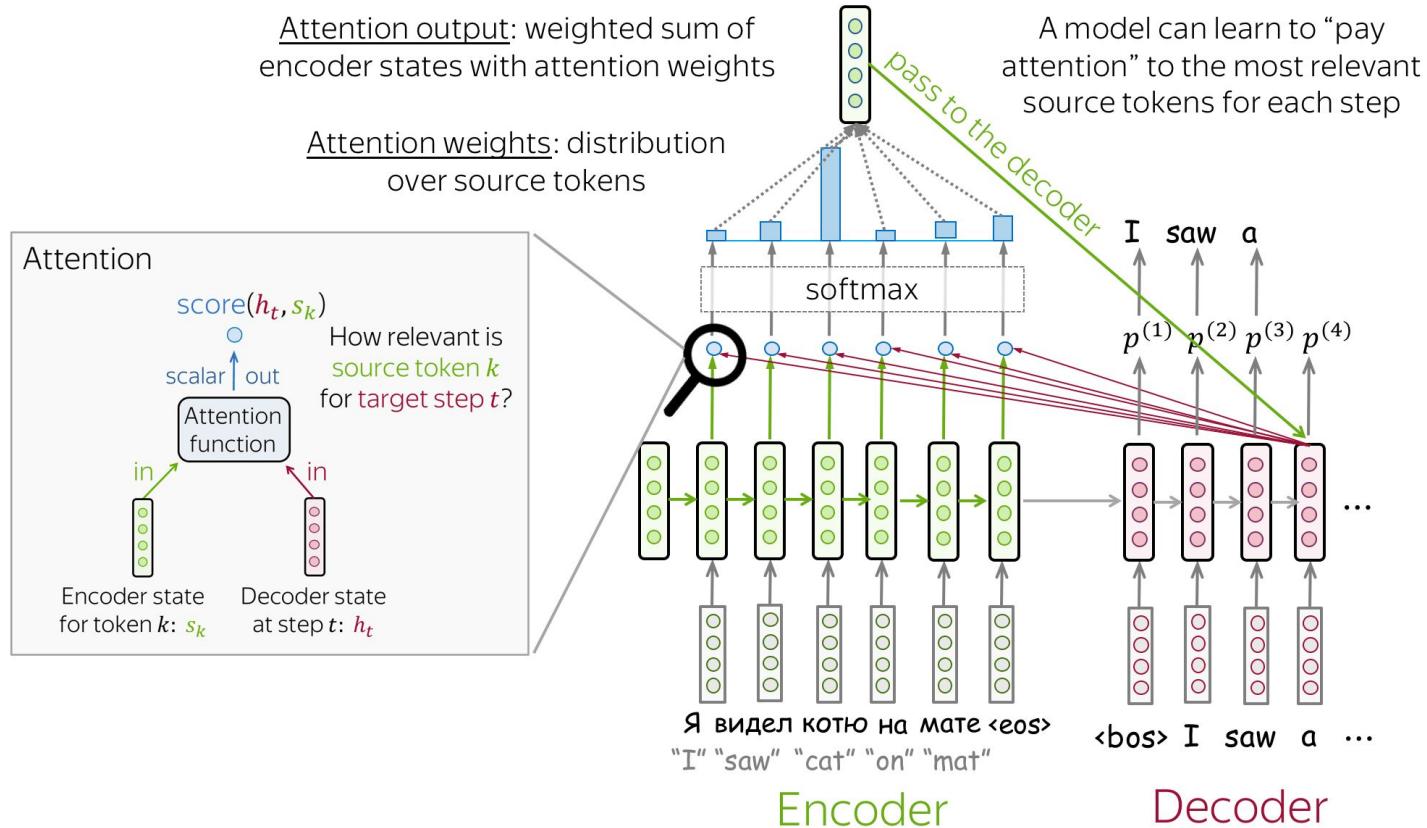
Target: I saw a cat on a mat <eos>

# Но есть нюанс

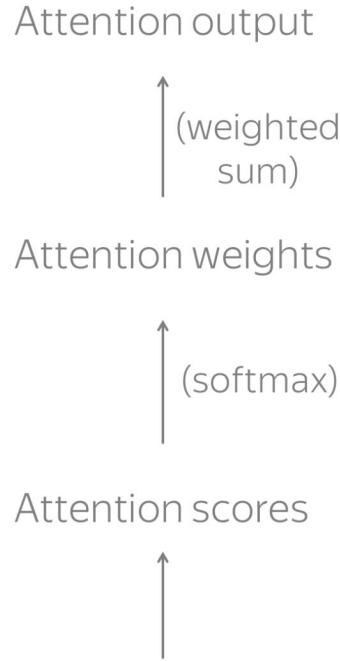
We saw: encoder compresses the source into a single vector



# Attention



# Подсчёт



$$c^{(t)} = a_1^{(t)} s_1 + a_2^{(t)} s_2 + \dots + a_m^{(t)} s_m = \sum_{k=1}^m a_k^{(t)} s_k$$

“source context for decoder step  $t$ ”

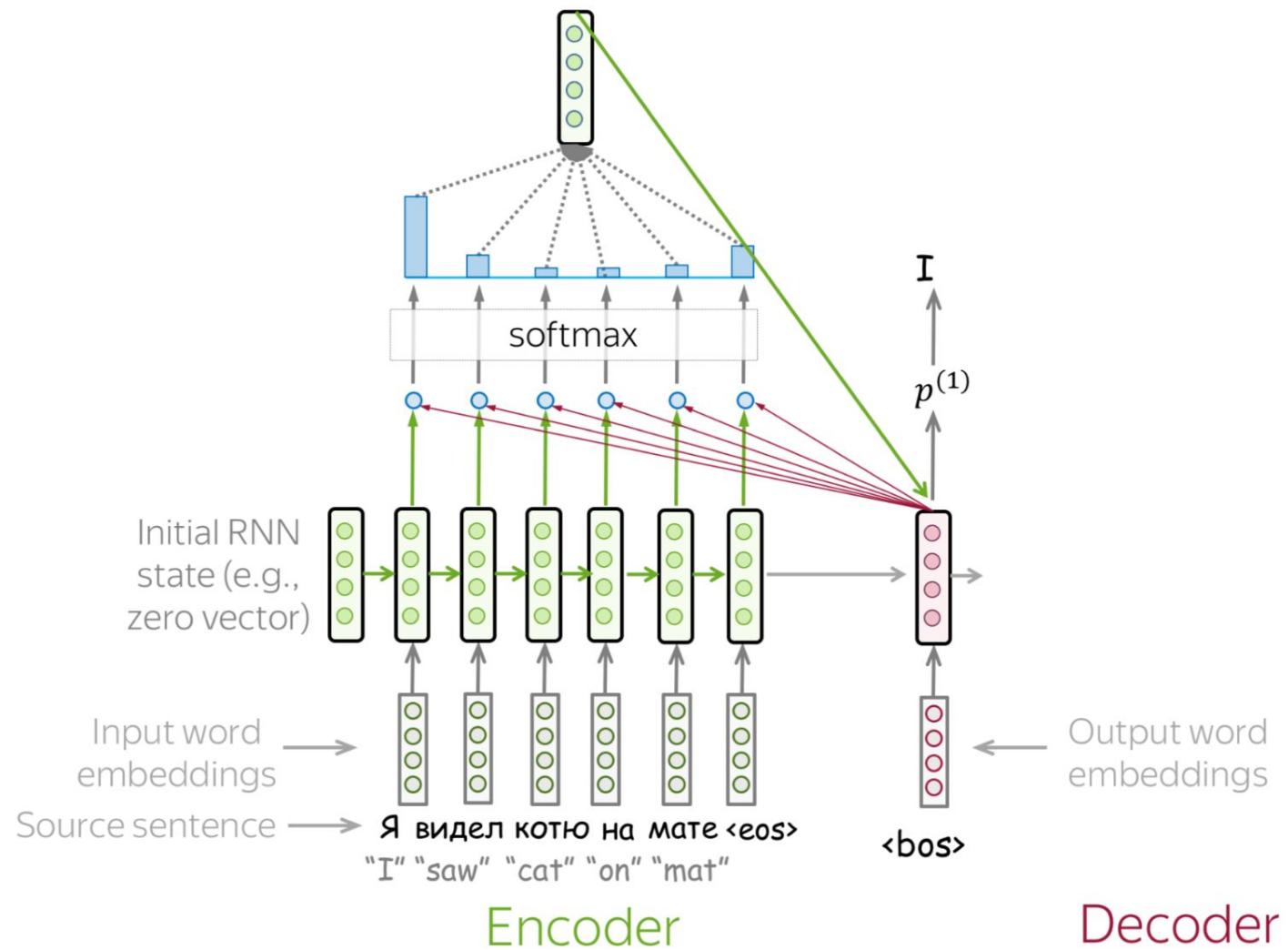
$$a_k^{(t)} = \frac{\exp(\text{score}(h_t, s_k))}{\sum_{i=1}^m \exp(\text{score}(h_t, s_i))}, k = 1..m$$

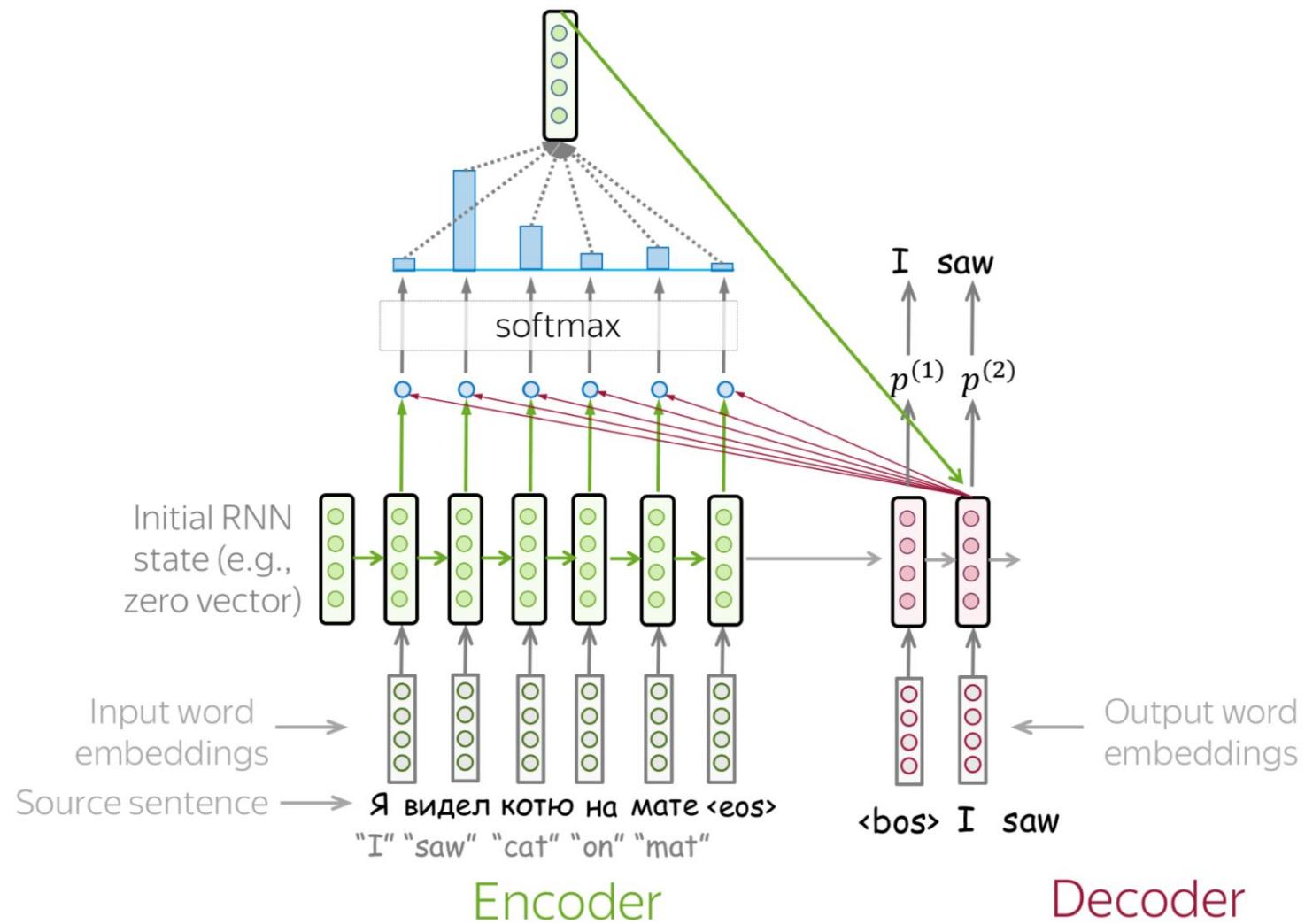
“attention weight for source token  $k$  at decoder step  $t$ ”

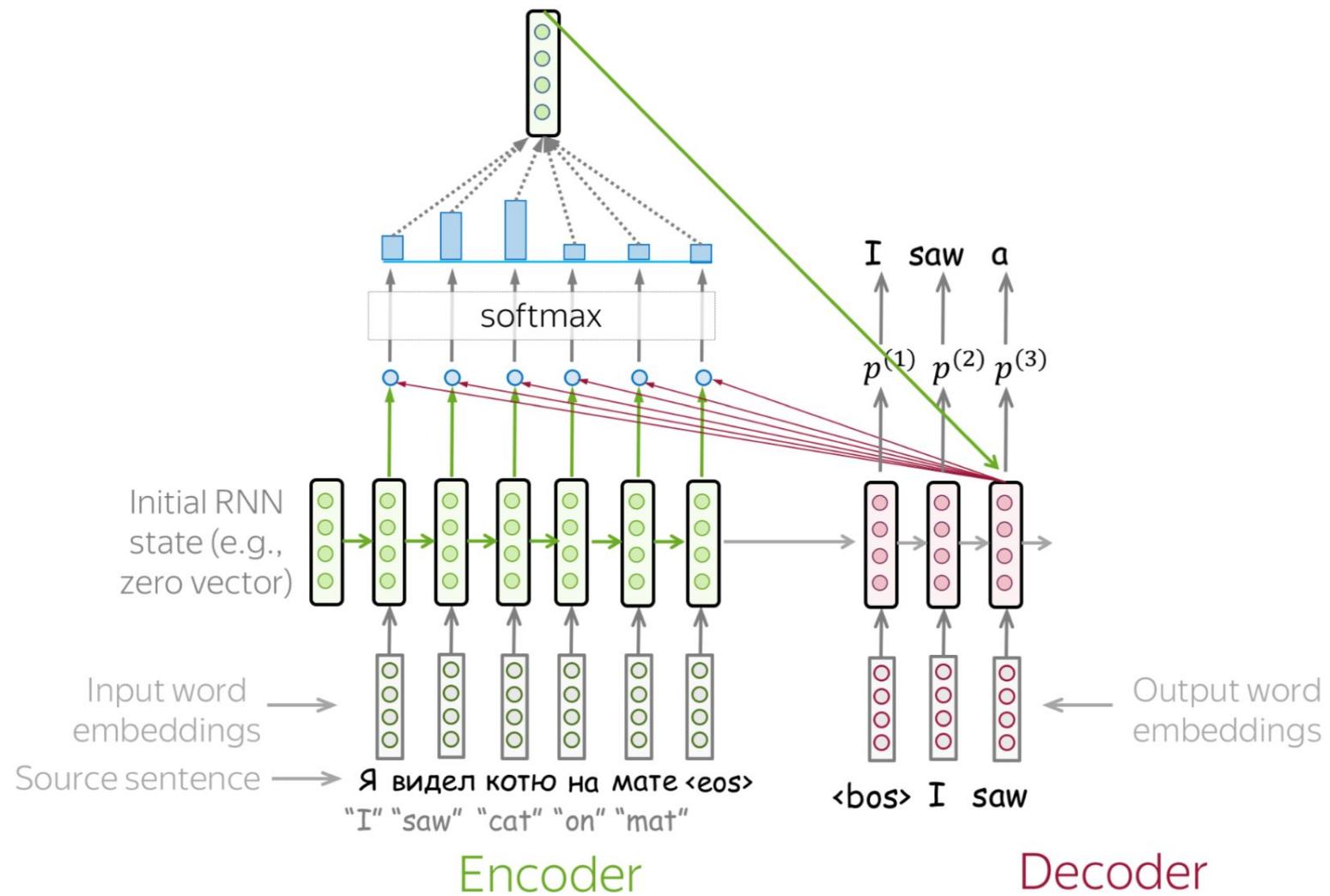
$$\text{score}(h_t, s_k), k = 1..m$$

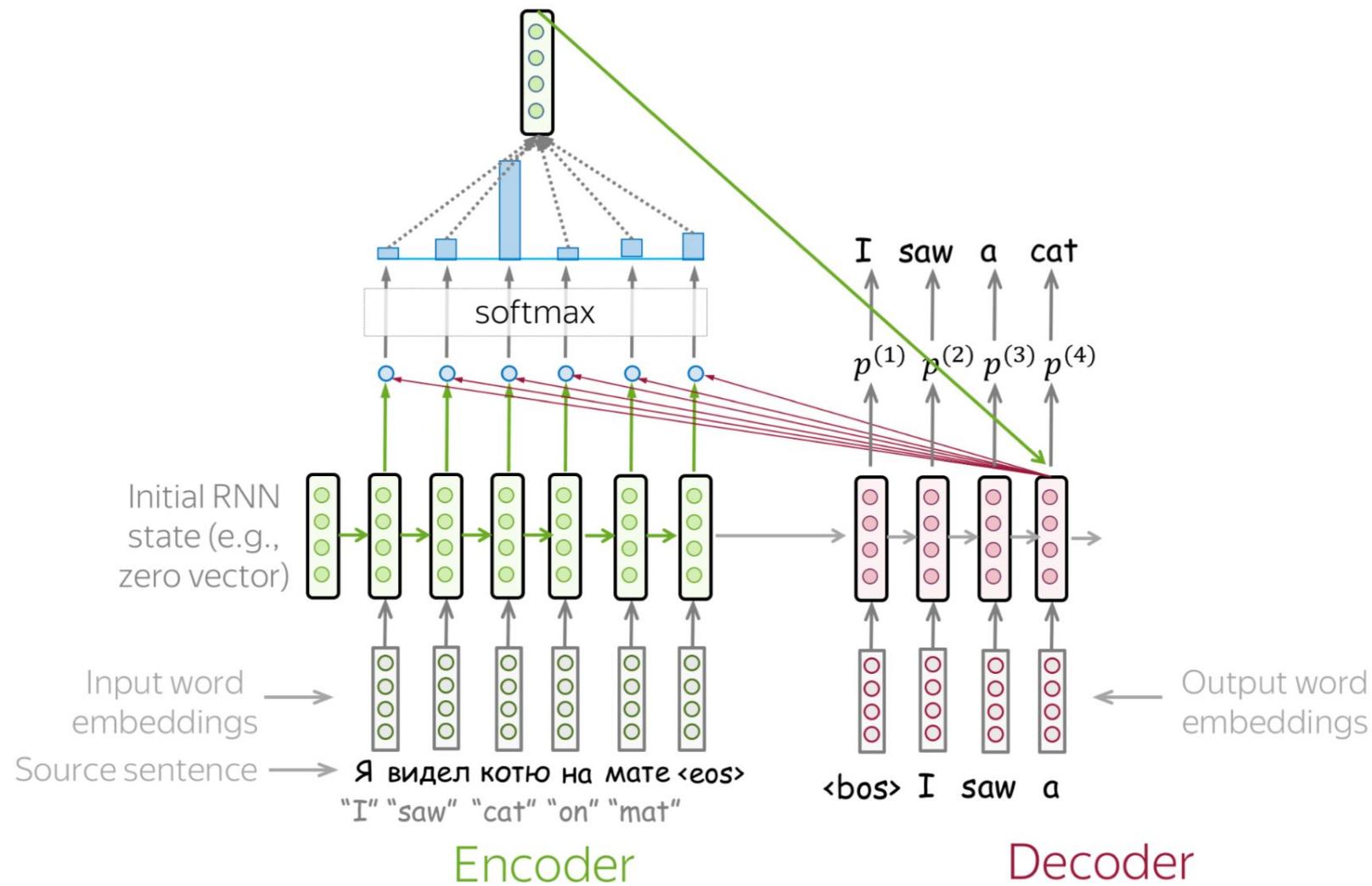
“How relevant is source token  $k$  for target step  $t$ ? ”

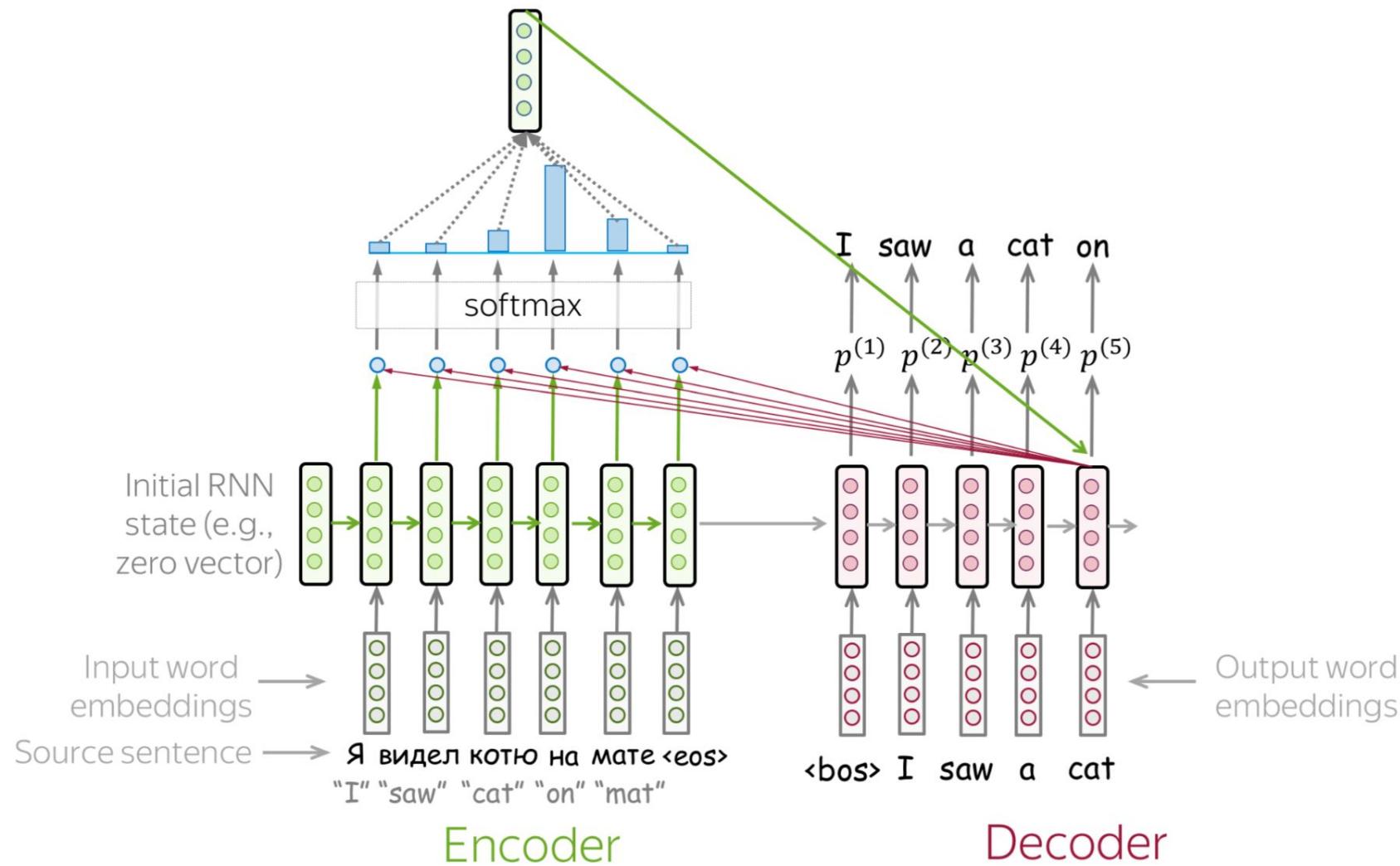
$s_1, s_2, \dots, s_m$   
all encoder states       $h_t$   
                            one decoder state

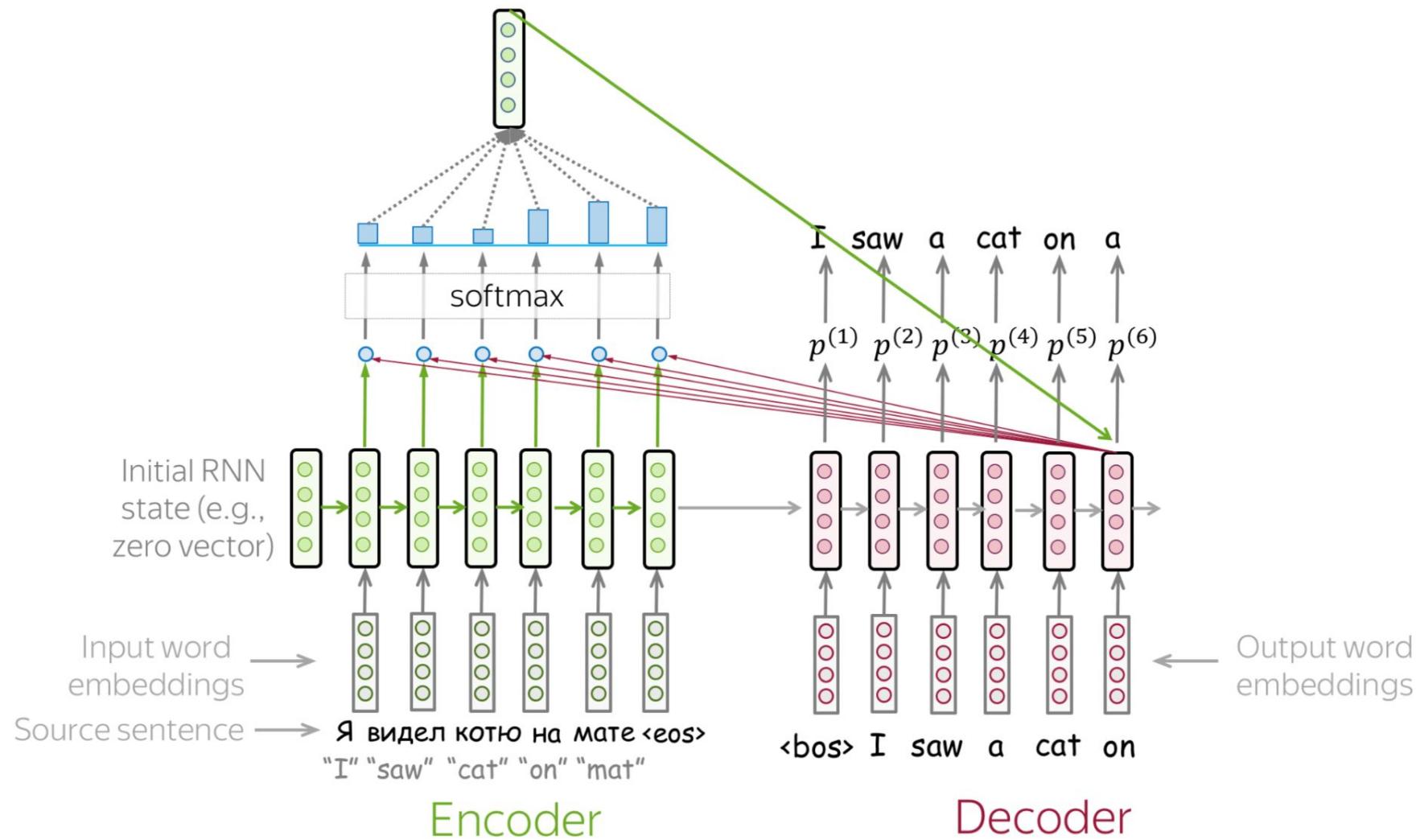


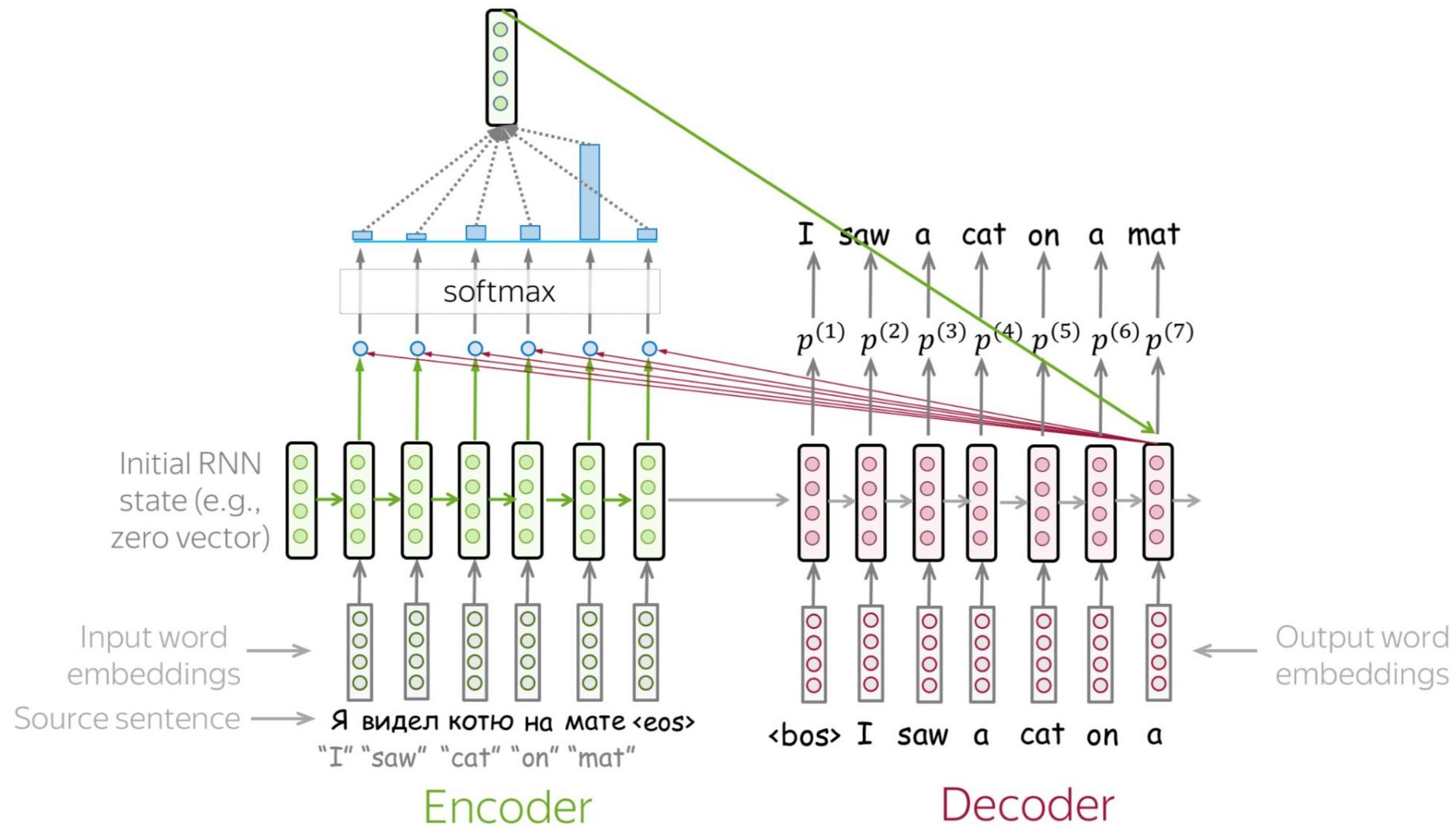


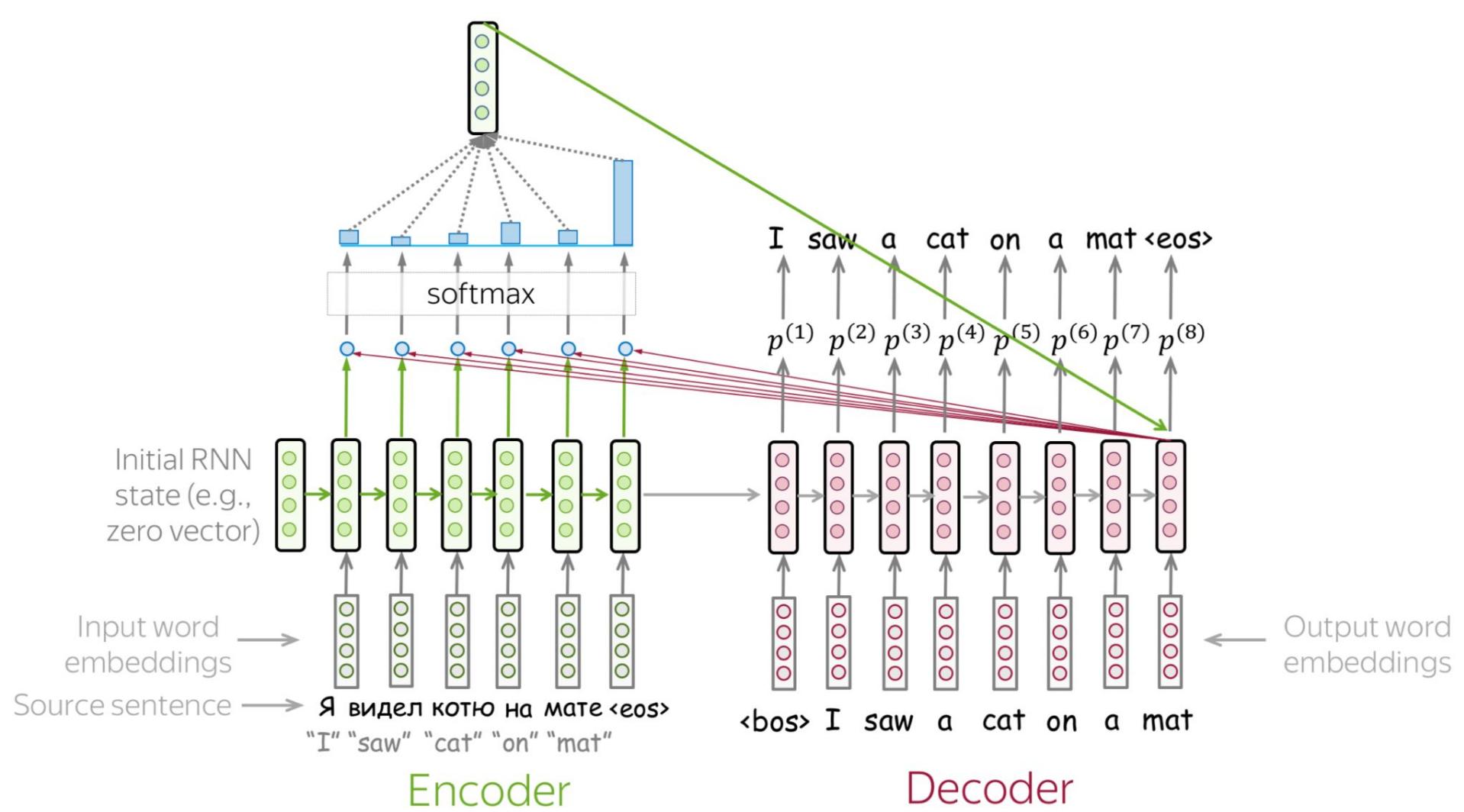




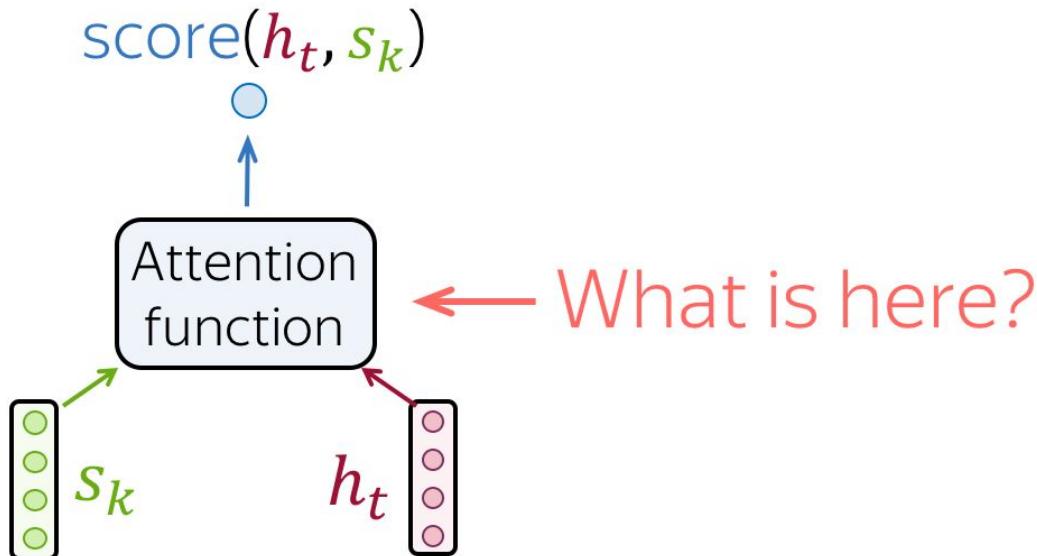






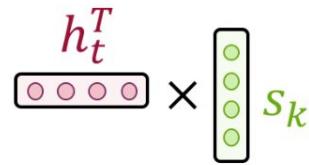


# А что за функция?



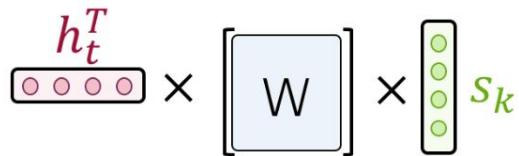
# В целом любая

Dot-product



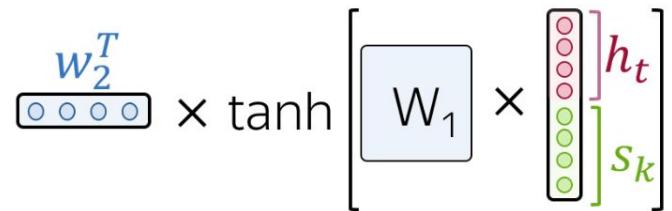
$$\text{score}(h_t, s_k) = h_t^T s_k$$

Bilinear



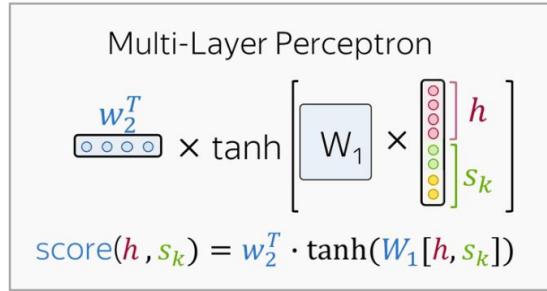
$$\text{score}(h_t, s_k) = h_t^T W s_k$$

Multi-Layer Perceptron

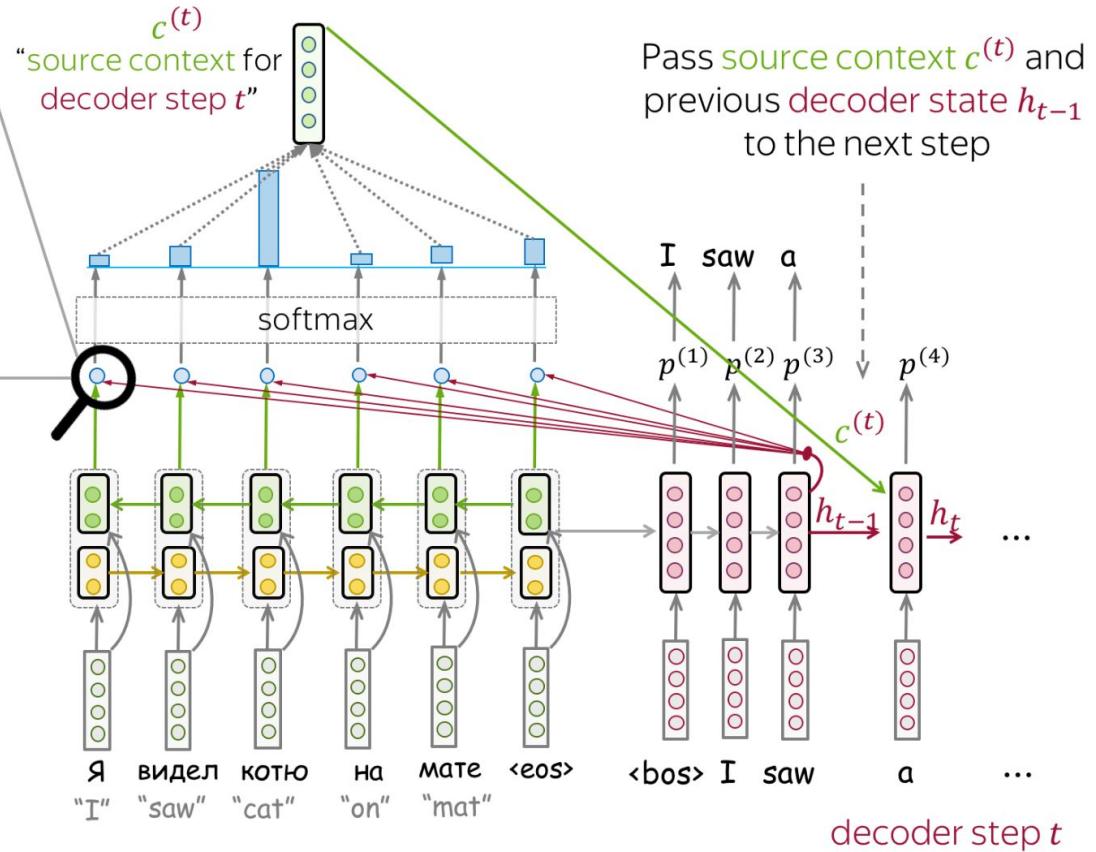


$$\text{score}(h_t, s_k) = w_2^T \cdot \tanh(W_1[h_t, s_k])$$

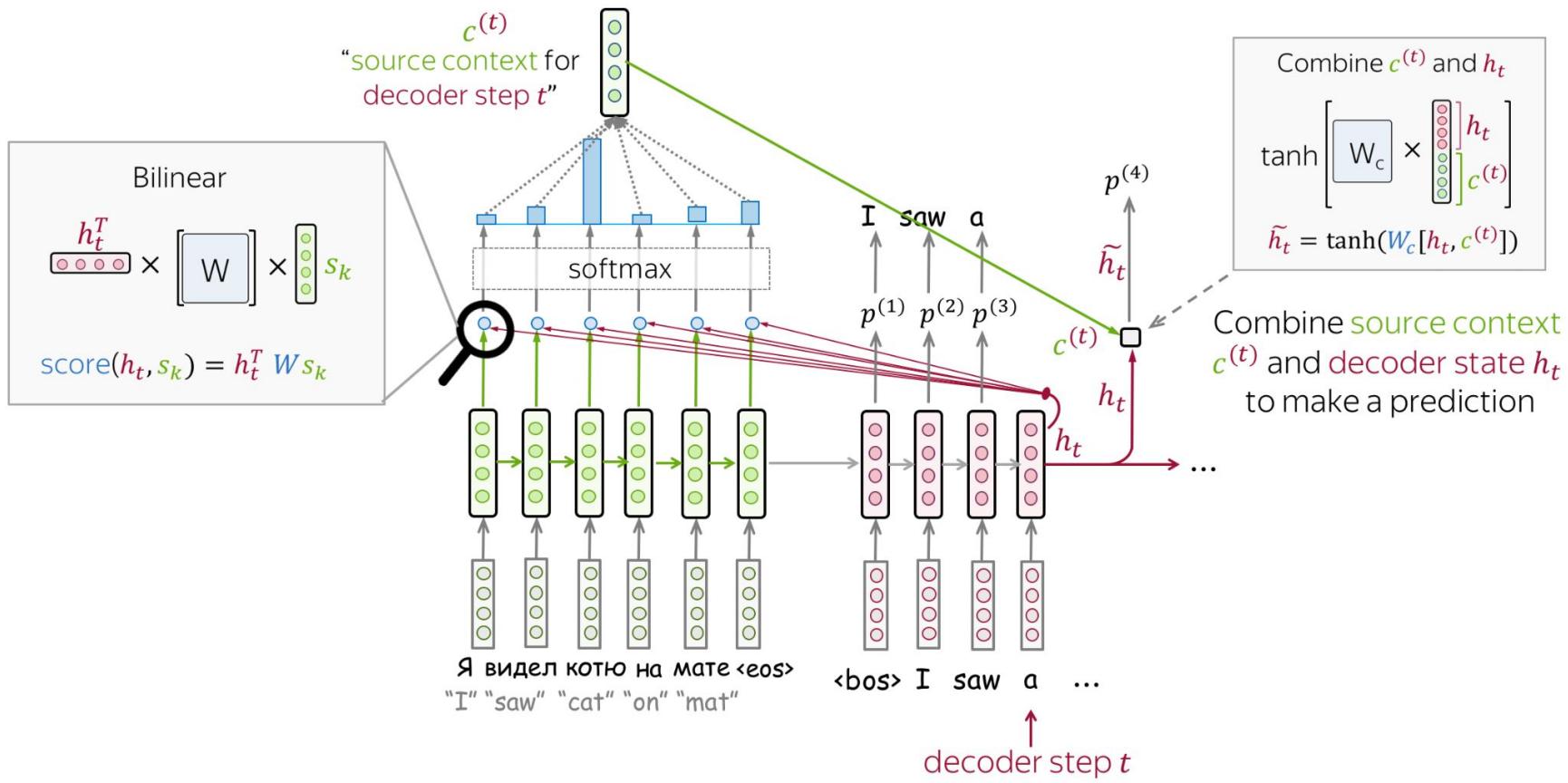
# Bahdanau model



Bidirectional encoder  
Concatenate states from  
forward and backward RNNs



# Luong model



# Alignment

L'  
accord  
sur  
la  
zone  
économique  
européenne  
a  
été  
signé  
en  
août  
1992  
.  
<end>

The  
agreement  
on  
the  
European  
Economic  
Area  
was  
signed  
in  
August  
1992  
. <end>

Il  
convient  
de  
noter  
que  
l'  
environnement  
marin  
est  
le  
moins  
connu  
de  
l'  
environnement  
.  
<end>

It  
should  
be  
noted  
that  
the  
marine  
environment  
is  
the  
least  
known  
of  
environments  
. <end>

# Attention is ALL you need

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

# Что хотим от трансформера?

## 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

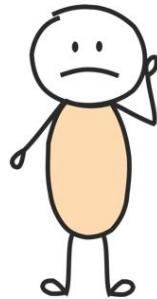
repeat  
N times



# Почему RNN будет работать хуже?

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

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



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

RNNs

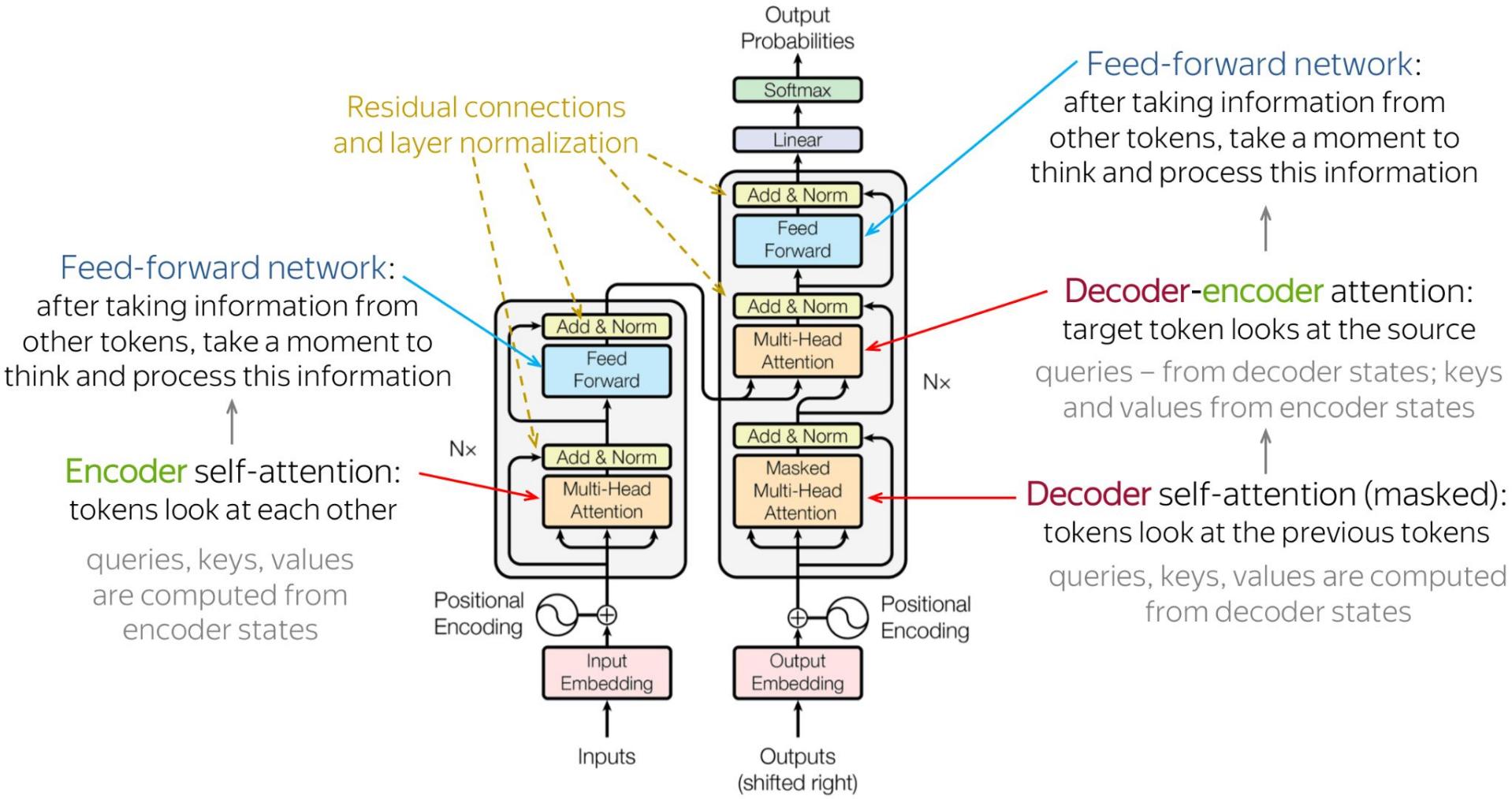
$O(N)$  steps to process a  
sentence with length  $N$



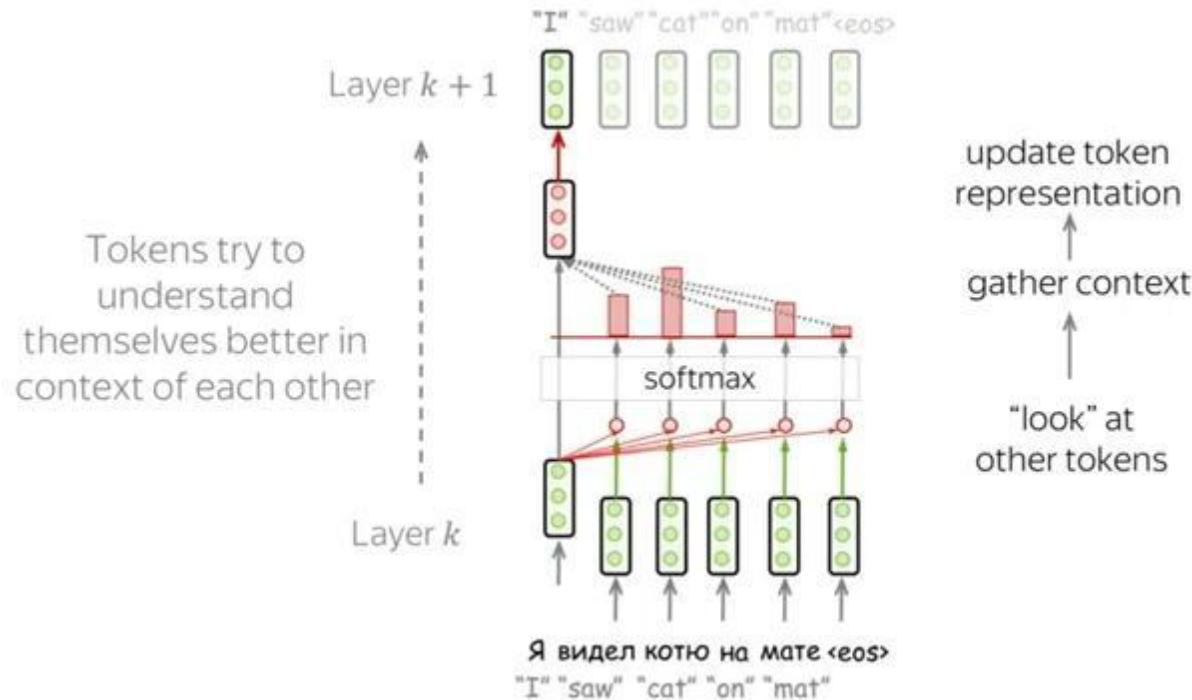
I don't need to wait - I  
see all words at once!

Transformer

Constant number of steps



# Self-attention



# Self-attention

Each vector receives three representations (“roles”)

$$[W_Q] \times \begin{array}{|c|c|c|}\hline & \textcolor{green}{\bullet} & \textcolor{green}{\bullet} \\ \hline & \textcolor{green}{\bullet} & \textcolor{green}{\bullet} \\ \hline & \textcolor{blue}{\bullet} & \textcolor{blue}{\bullet} \\ \hline \end{array} = \begin{array}{|c|c|c|}\hline & \textcolor{blue}{\bullet} & \textcolor{blue}{\bullet} \\ \hline & \textcolor{blue}{\bullet} & \textcolor{blue}{\bullet} \\ \hline & \textcolor{blue}{\bullet} & \textcolor{blue}{\bullet} \\ \hline \end{array}$$

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

“Hey there, do you have this information?”

$$[W_K] \times \begin{array}{|c|c|c|}\hline & \textcolor{green}{\bullet} & \textcolor{green}{\bullet} \\ \hline & \textcolor{green}{\bullet} & \textcolor{green}{\bullet} \\ \hline & \textcolor{yellow}{\bullet} & \textcolor{yellow}{\bullet} \\ \hline \end{array} = \begin{array}{|c|c|c|}\hline & \textcolor{yellow}{\bullet} & \textcolor{yellow}{\bullet} \\ \hline & \textcolor{yellow}{\bullet} & \textcolor{yellow}{\bullet} \\ \hline & \textcolor{yellow}{\bullet} & \textcolor{yellow}{\bullet} \\ \hline \end{array}$$

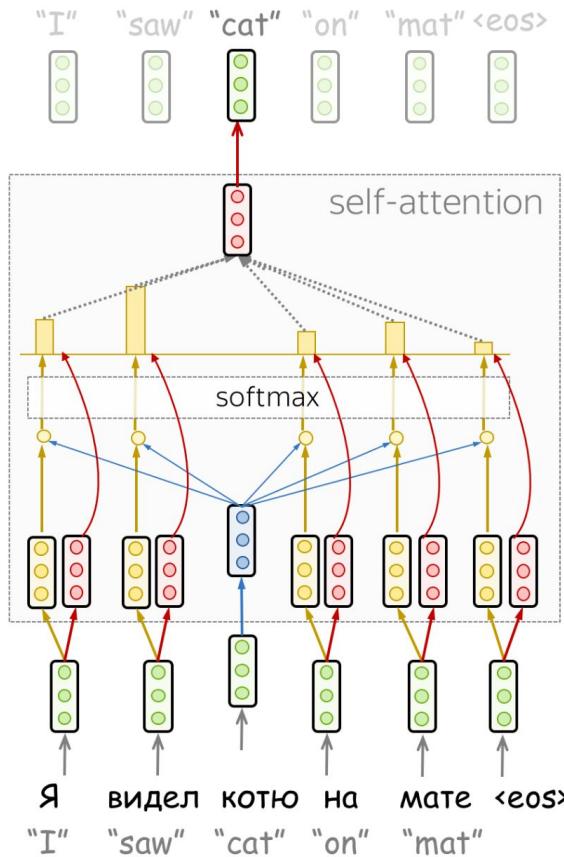
**Key:** vector **at** which the query  
looks to compute weights

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

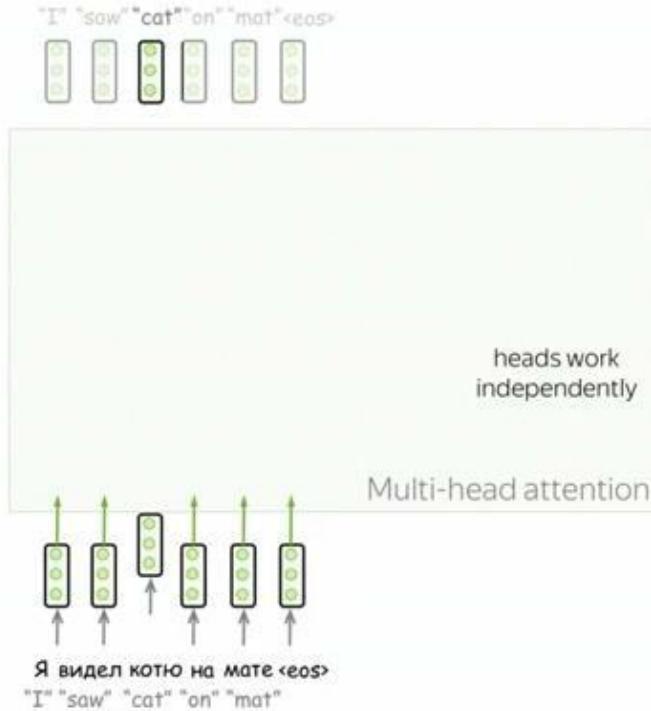
$$[W_V] \times \begin{array}{|c|c|c|}\hline & \textcolor{green}{\bullet} & \textcolor{green}{\bullet} \\ \hline & \textcolor{green}{\bullet} & \textcolor{green}{\bullet} \\ \hline & \textcolor{red}{\bullet} & \textcolor{red}{\bullet} \\ \hline \end{array} = \begin{array}{|c|c|c|}\hline & \textcolor{red}{\bullet} & \textcolor{red}{\bullet} \\ \hline & \textcolor{red}{\bullet} & \textcolor{red}{\bullet} \\ \hline & \textcolor{red}{\bullet} & \textcolor{red}{\bullet} \\ \hline \end{array}$$

**Value:** their weighted sum is  
attention output

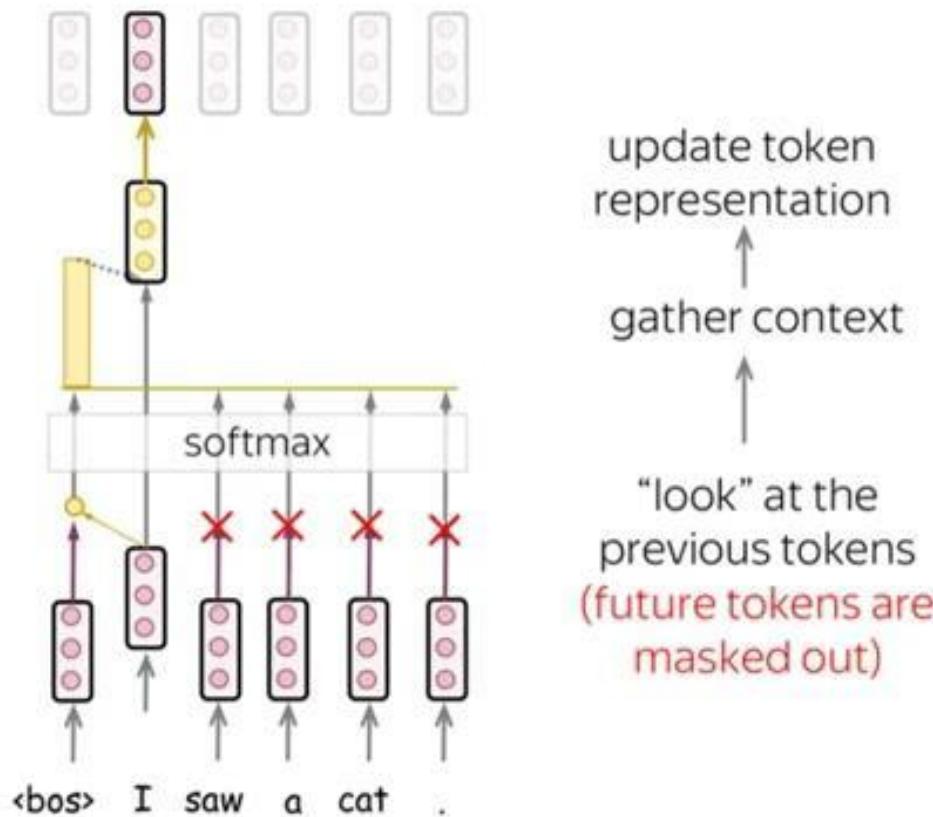
“Here’s the information I have!”



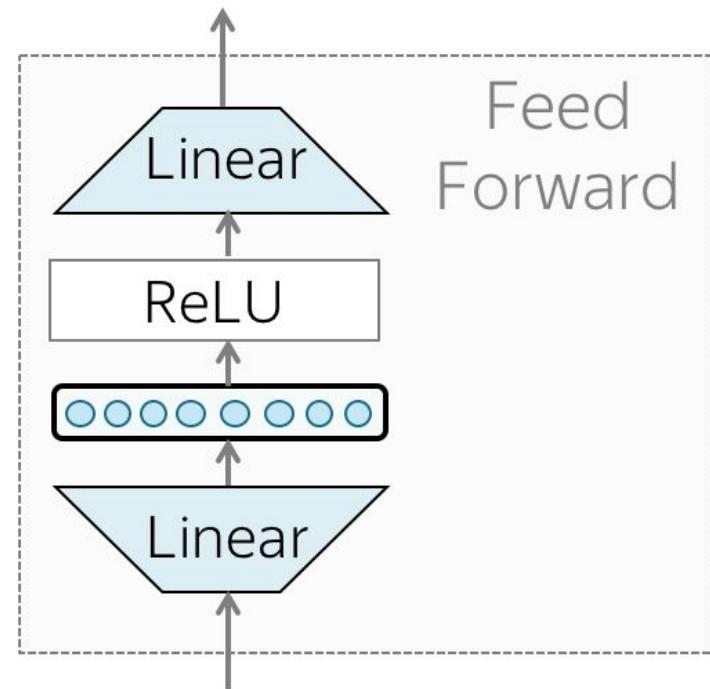
# Multi-head attention



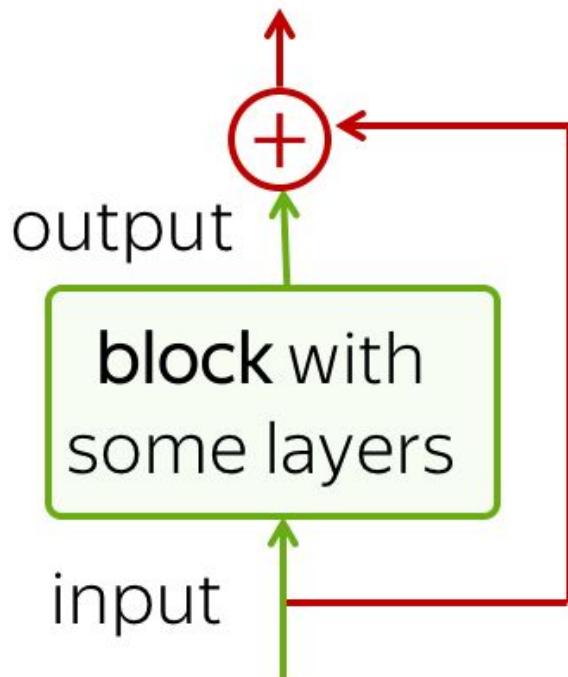
# Masked attention



# FFN

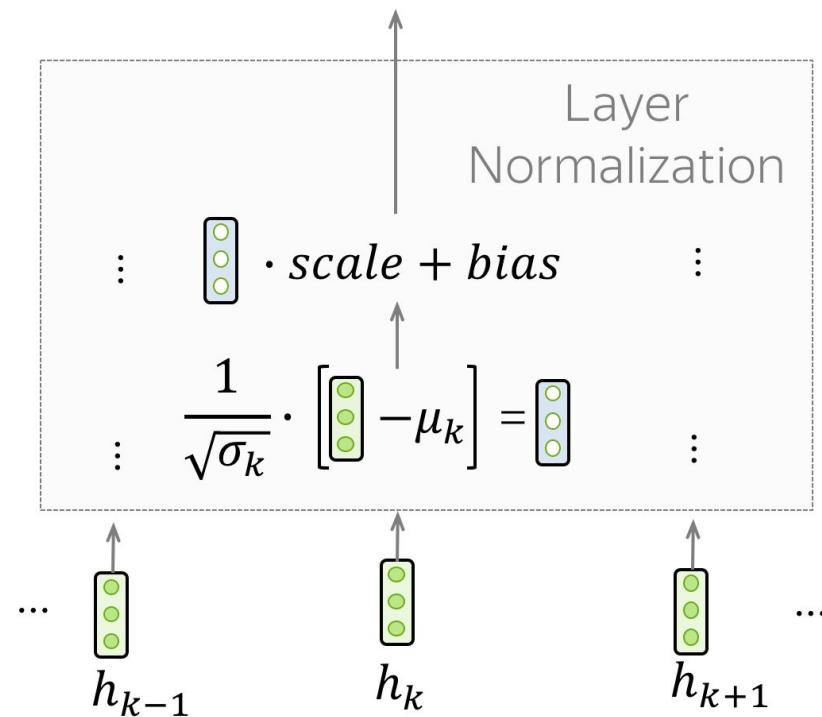


# Residual connection

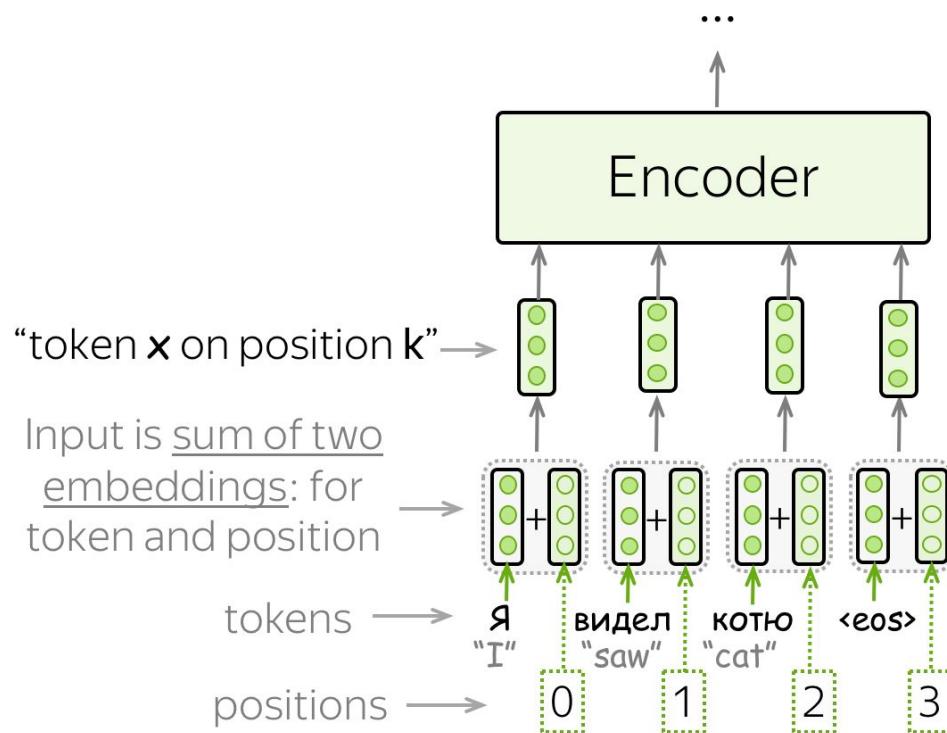


Residual connection:  
add a block's input to  
its output

# Layer norm



# Positional encoding



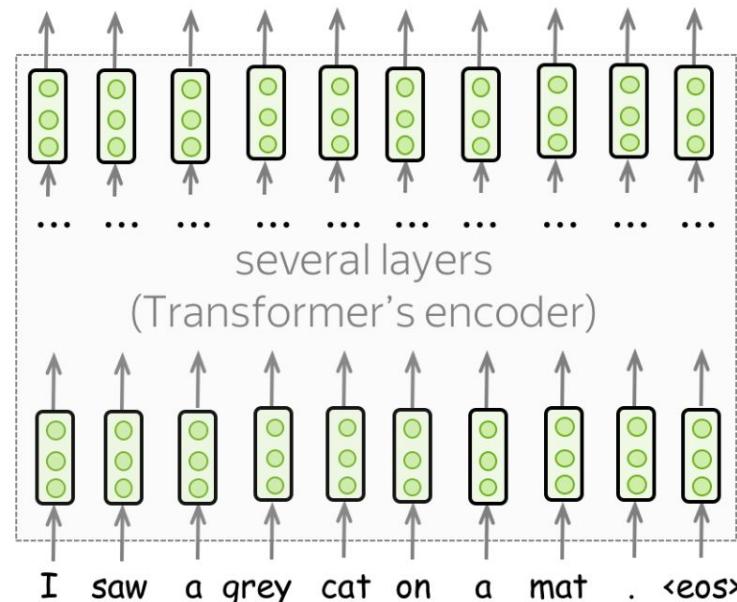
$$\text{PE}_{pos,2i} = \sin(pos/10000^{2i/d_{model}}),$$

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

# BPE (byte-pair encoding)

	Words in the data:		
Initial vocabulary: characters	word	count	Current merge table: (empty)
	c a t	4	
	m a t	5	
	m a t s	2	
	m a t e	3	
	a t e	3	
	e a t	2	
Split each word into characters			

# BERT



Model architecture:

- Transformer's encoder

What is special about it:

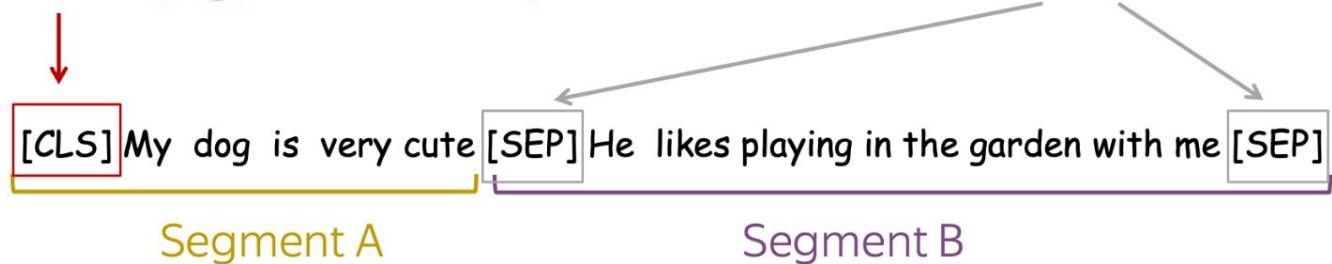
- Training objectives
  - MLM: Masked language modeling
  - NSP: Next sentence prediction
- The way it is used
  - No task-specific models

# Input

[CLS]: Special token

- Training time: predict if sentences are consecutive or not (Next Sentence Prediction /NSP objective)
- Test time: downstream tasks (e.g., classification)

[SEP]: Special token-separator



Training on pairs of sentences: either consecutive or random (50%/50%)

Training time: predict if sentences  
are consecutive (NSP objective)

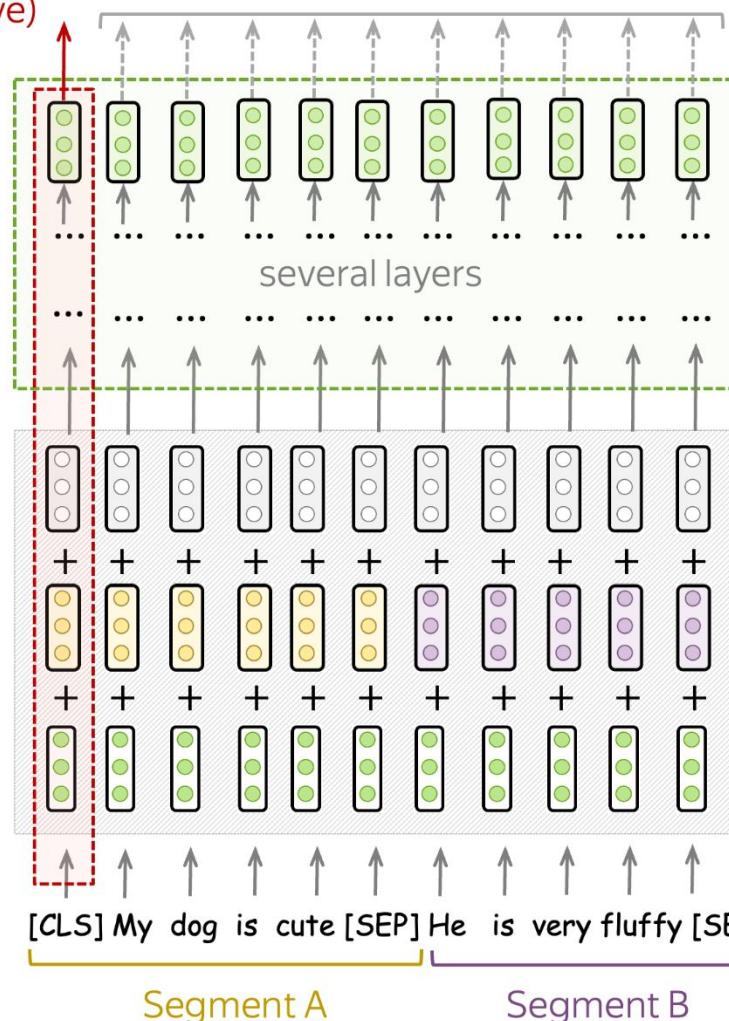
Test time: classification

Model  
(Transformer  
encoder)

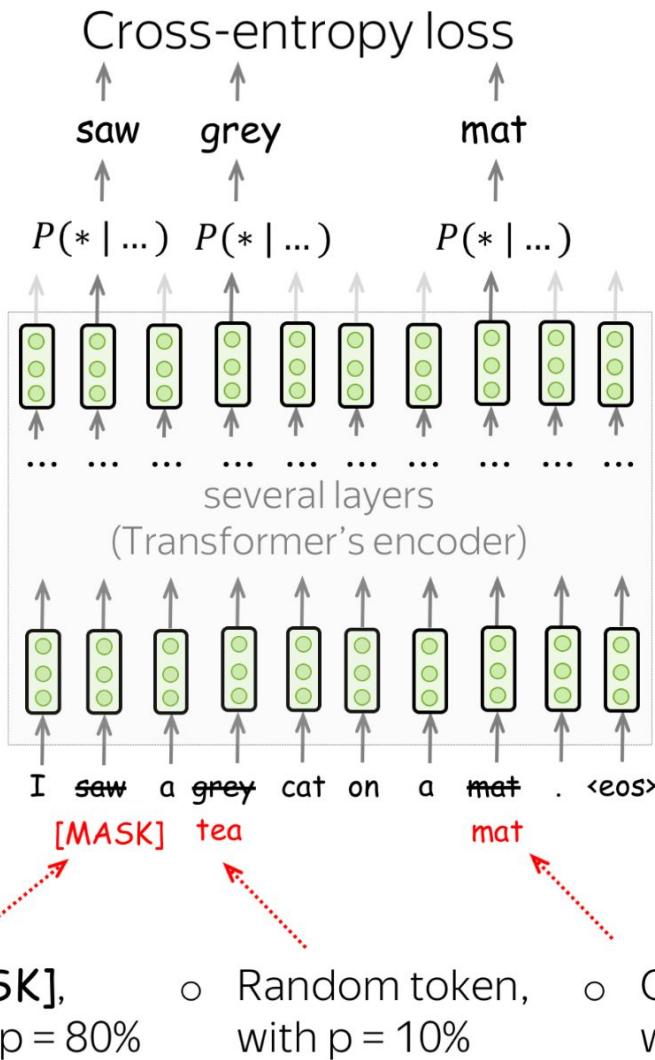
Input

Training on pairs of  
sentences: either  
consecutive or  
random (50%/50%)

Training time: MLM objective



Loss  
↑  
Target  
↑  
Prediction

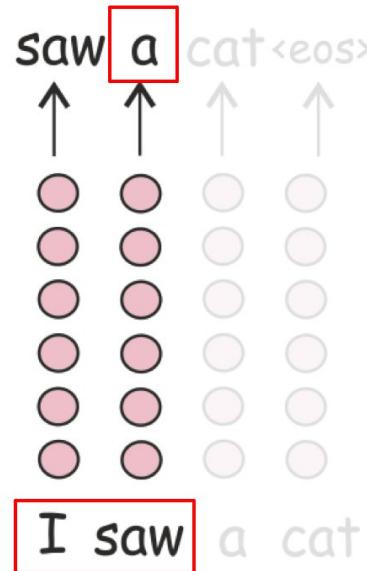


At each training step:

- pick randomly 15% of tokens
- replace each of the chosen tokens with something
- predict original chosen tokens

## Language Modeling

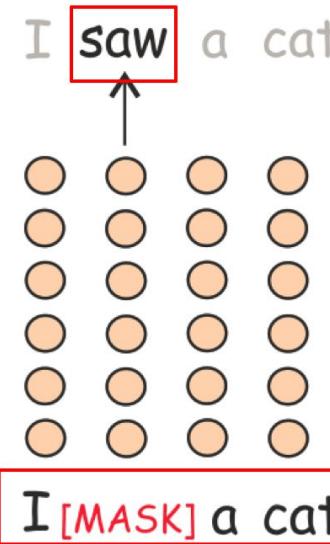
- Target: next token
- Prediction:  $P(*) | \text{I saw}$



left-to-right, does  
not see future

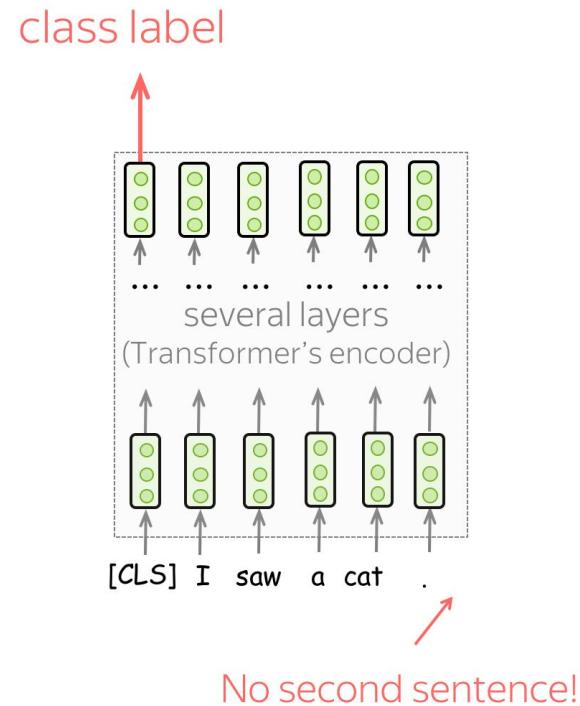
## Masked Language Modeling

- Target: current token (the true one)
- Prediction:  $P(*) | \text{I [MASK] a cat}$

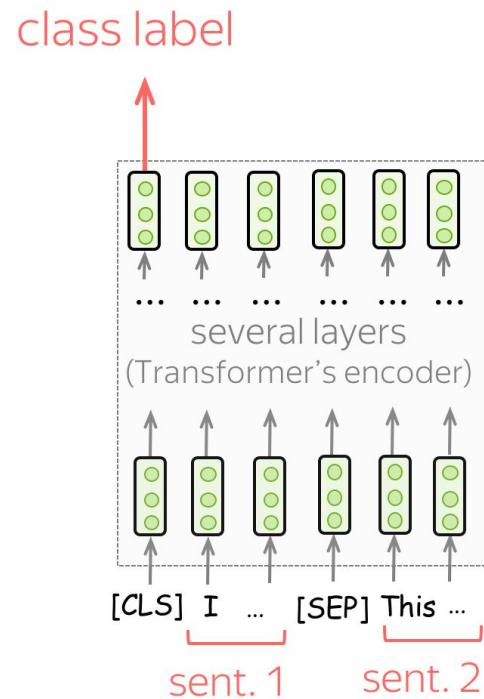


sees the whole text, but  
something is corrupted

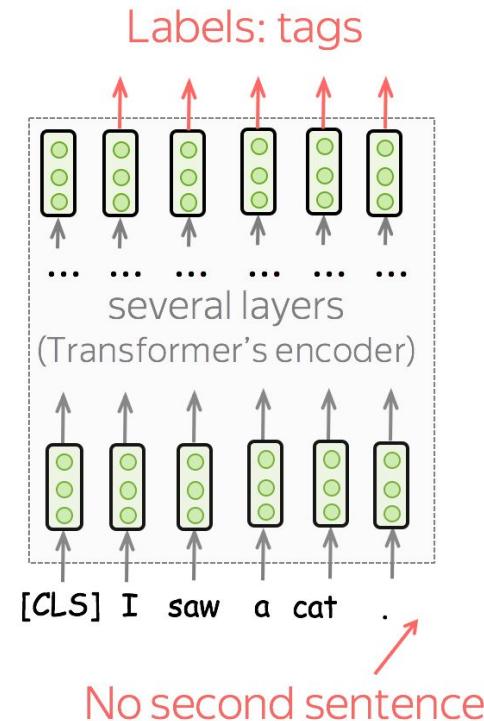
# Sequence classification



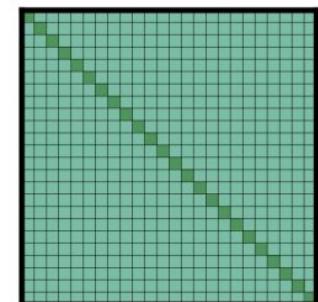
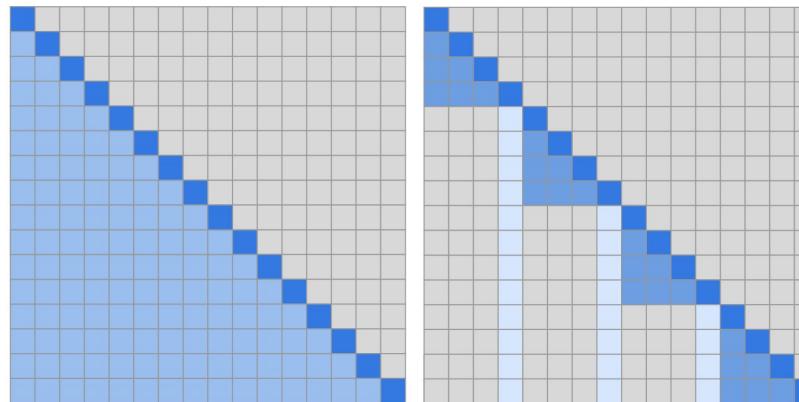
# Relevance (e.g.)



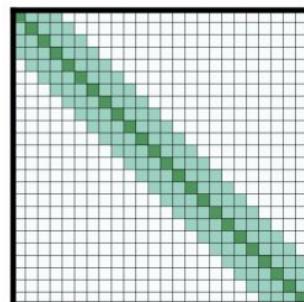
# NER (named entity recognition)



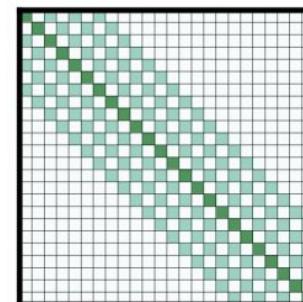
# Разные маски attention



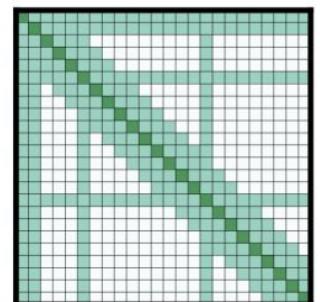
(a) Full  $n^2$  attention



(b) Sliding window attention

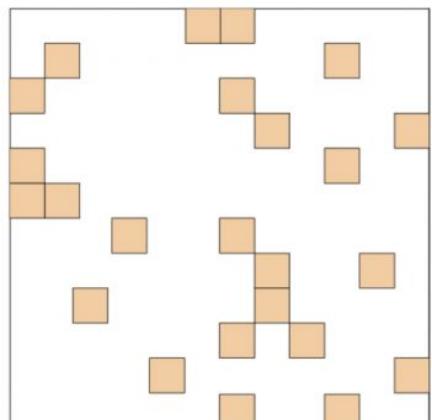


(c) Dilated sliding window

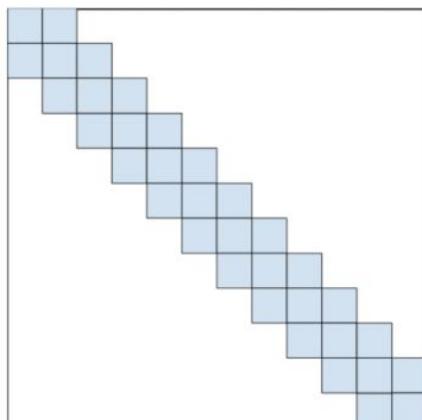


(d) Global+sliding window

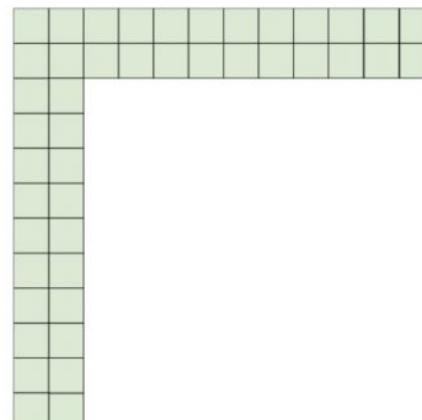
# Разные маски attention



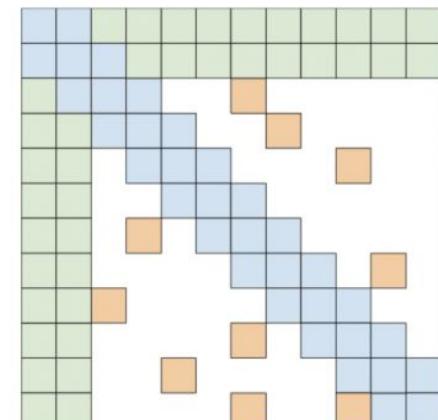
(a) Random attention



(b) Window attention



(c) Global Attention



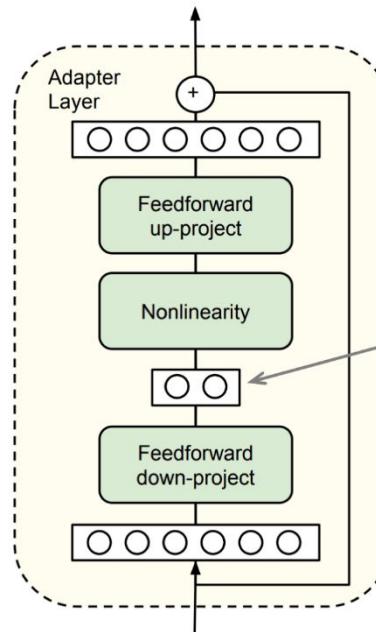
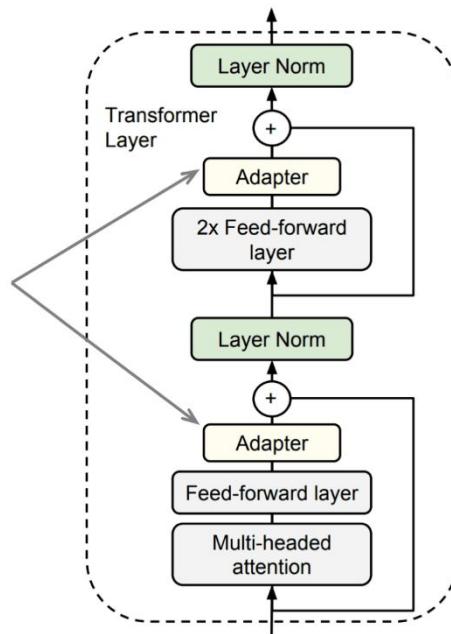
(d) BIGBIRD

# Модификации подсчёта внимания

$$\text{softmax}(Q_{n \times d} K_{n \times d}^T) V_{n \times d} \approx Q'_{n \times r} (K'^T_{n \times r} V_{n \times d})$$

# Адаптеры

Only these are trained,  
everything else is fixed and  
is the same for all tasks



# LoRA (Low Rank Adaptation)

