# Foundations for Natural Language Processing

## Improving Encoder-Decoder, Attention

Ivan Titov

(with graphics/materials from Elena Voita)

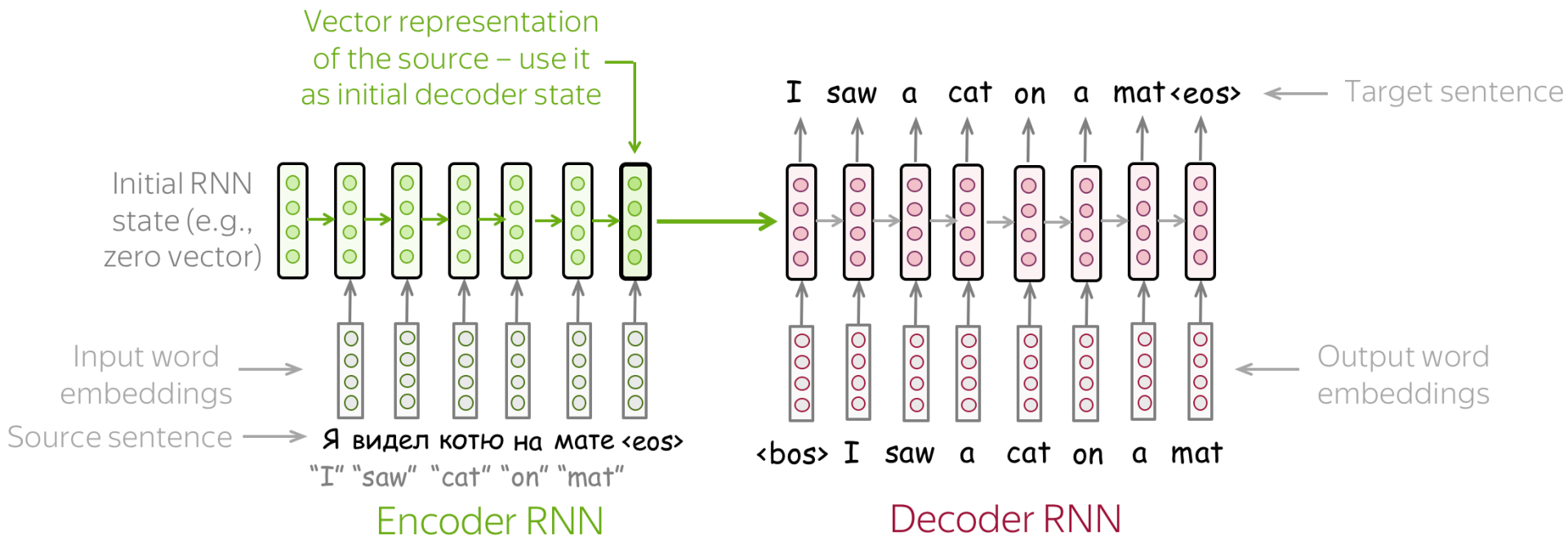**School of informatics**

# Plan for today

Last time:
- Vanilla Encoder-Decoder
- Text Generation (training, inference, evaluation)

Today

- Subword segmentation

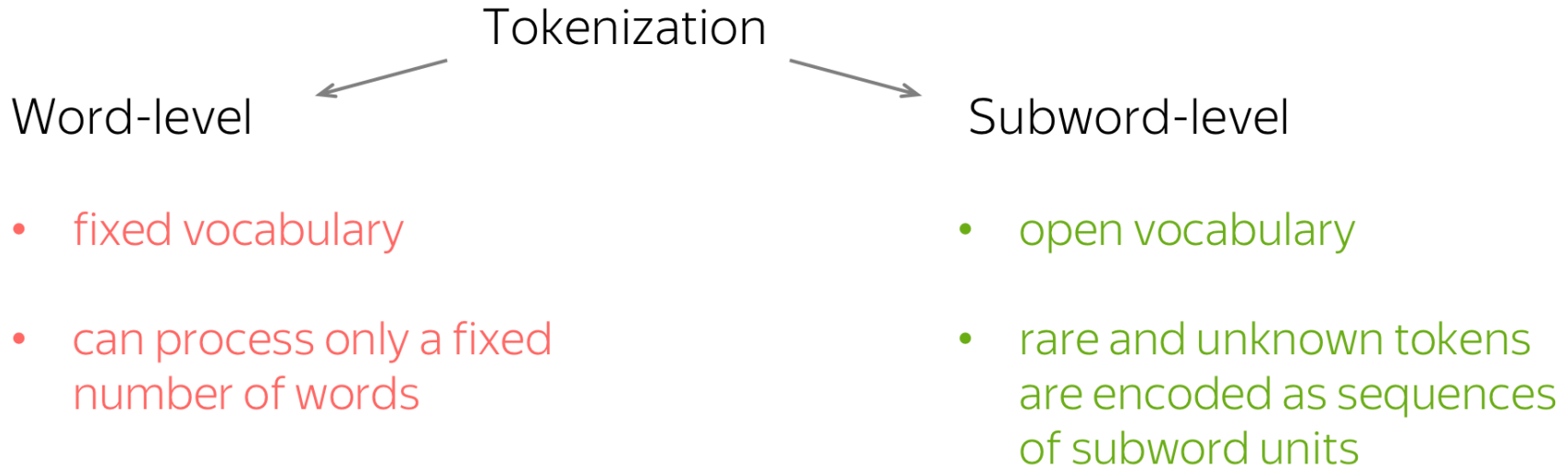- Improving Encoder-Decoder models

- Modeling Attention

# Tokenization

Vector representation
of the source – use it
as initial decoder state

I saw a cat on a mat<eos> ← Target sentence

Initial RNN
state (e.g.,
zero vector)

Input word
embeddings

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

Encoder RNN

<bos> I saw a cat on a mat

Output word
embeddings

Decoder RNN

We considered tokenization of sentences into 'words'
(whatever we mean by a 'word')

# Tokenization

Tokenization

Word-level

Subword-level

- fixed vocabulary

- can process only a fixed number of words

- open vocabulary

- rare and unknown tokens are encoded as sequences of subword units

Instead of '*unrelated*', we get two tokens   'un@@' '*related*'

Subword segmentations reduces sparsity and results in a speeds-up (recall: softmax involves summation over all token types, few token typs -> faster computation)
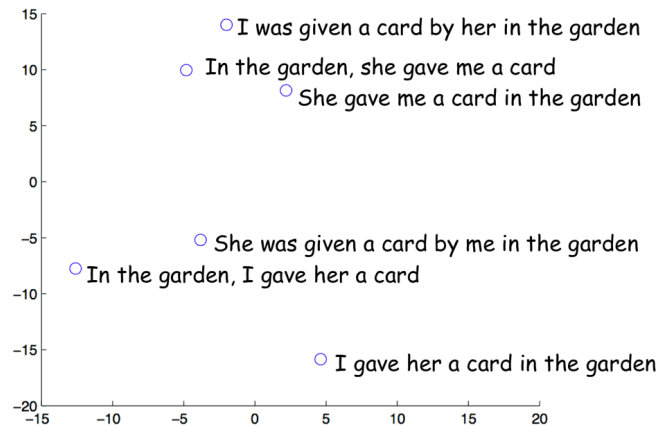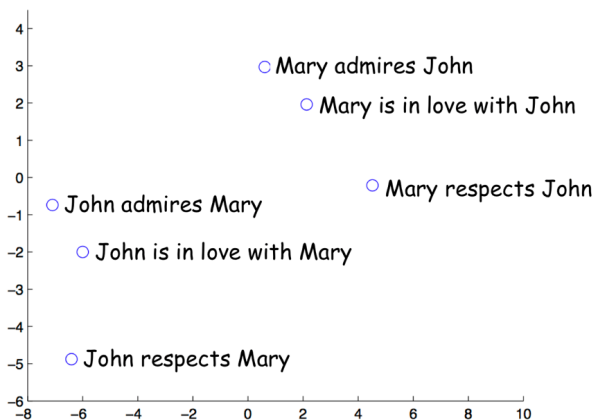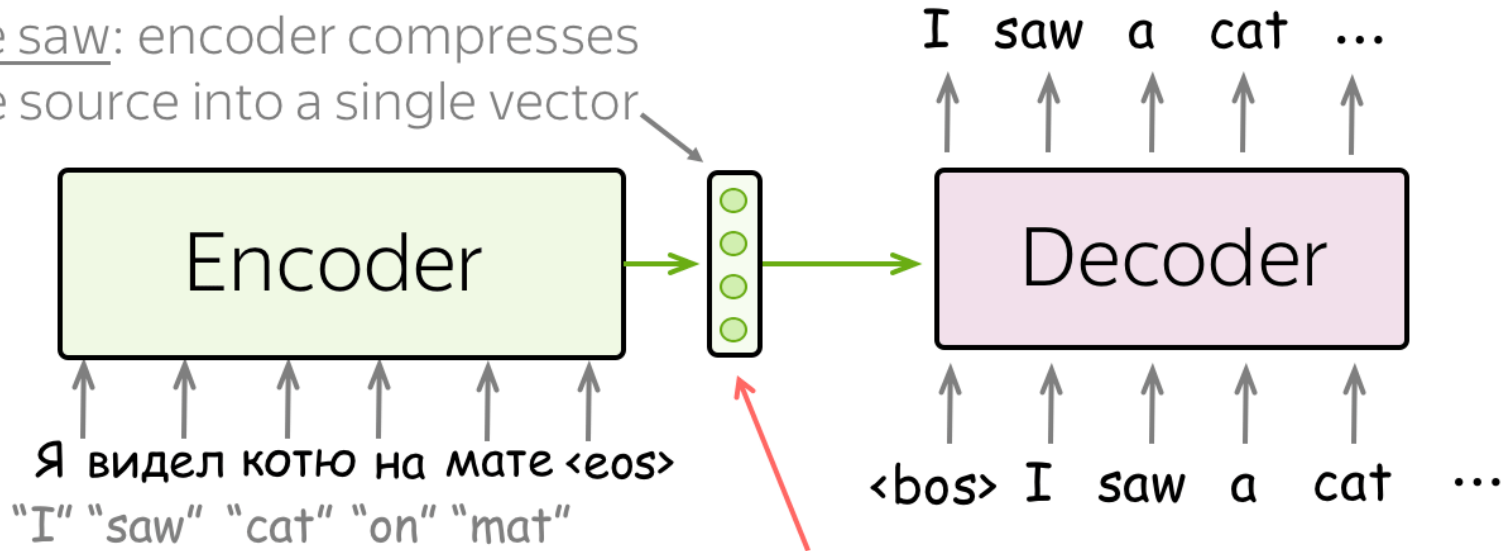
Crucial for morphologically-rich languages

Standard segmentation algorithms rely on character ngram frequency (not on morphology (e.g., Byte-Pair Encoding)

Used in virtually any modern neural model
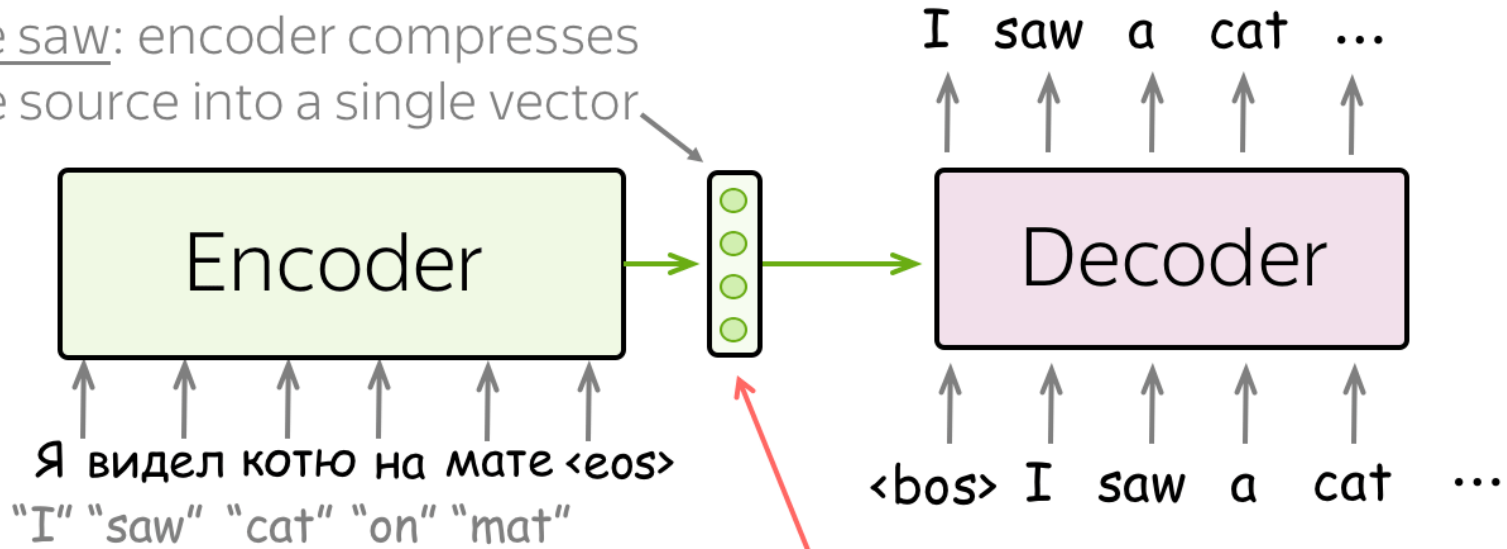
# The key problem with this approach

We saw: encoder compresses the source into a single vector.

**Problem**: this is a bottleneck!

I saw a cat ...

Encoder

Decoder

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

<bos> I saw a cat ...

# The key problem with this approach

We saw: encoder compresses the source into a single vector.



Encoder

Decoder

Я видел котю на мате <eos>

"I" "saw" "cat" "on" "mat"

I saw a cat ...

<bos> I saw a cat ...

Problem: this is a bottleneck!

Problem: fixed source representation is suboptimal:
• for the encoder, it is hard to compress the sentence;
• for the decoder, different information may be relevant at different steps.
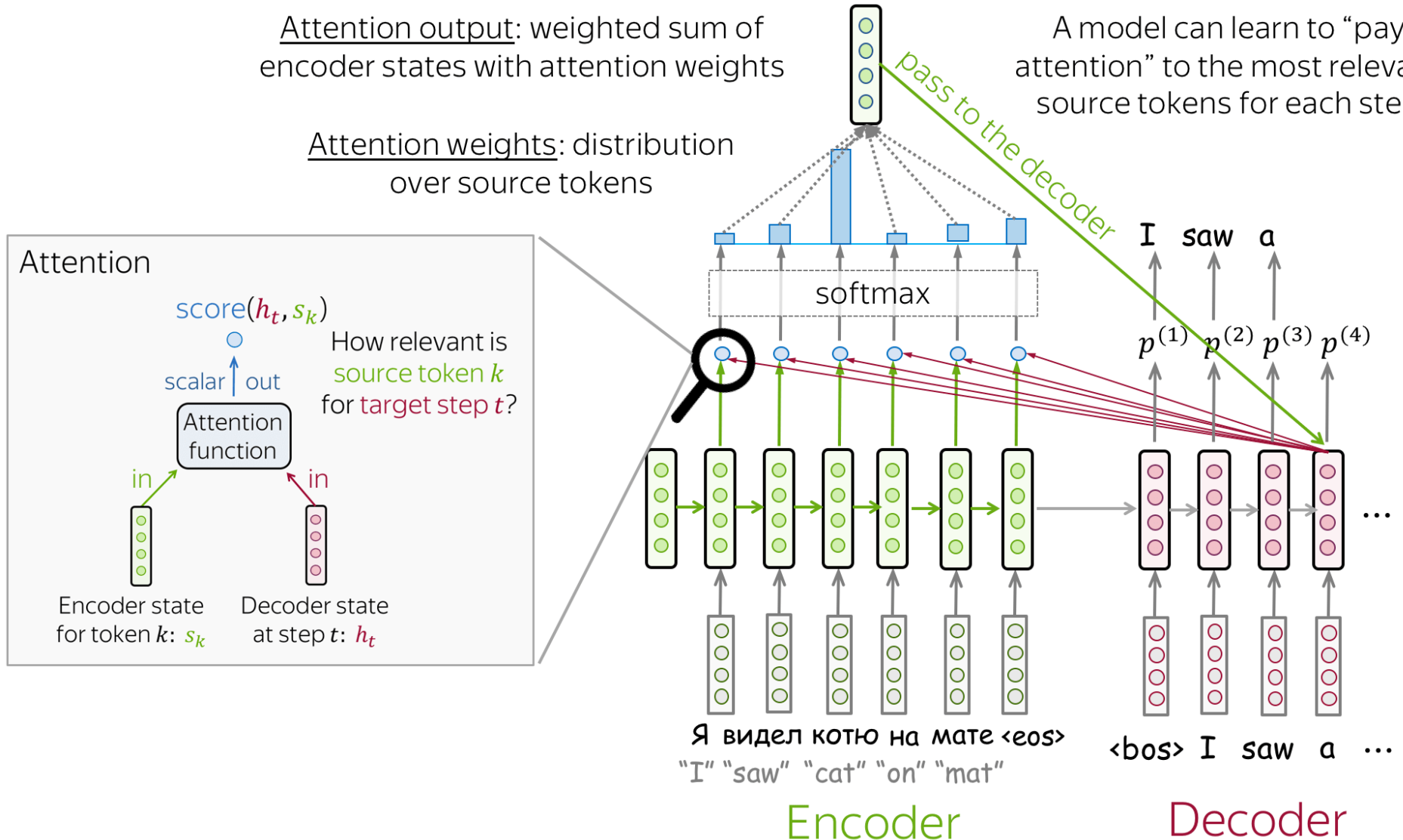
Solution: modeling "attention"

# Attention: Intuition

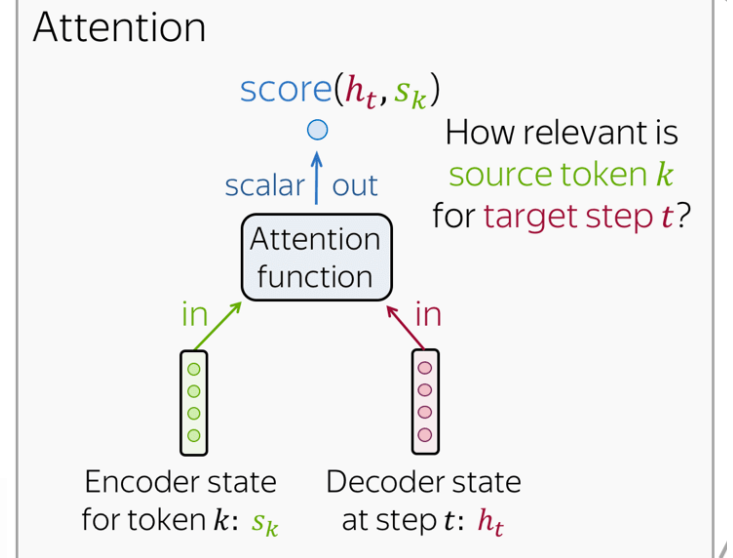## At every step, the decoder decide on which input tokens to focus

Attention output: weighted sum of encoder states with attention weights

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

_pass to the decoder_

Attention weights: distribution over source tokens

### Attention

$\text{score}(h_t, s_k)$

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

scalar $\uparrow$ out

Attention function

in $\qquad$ in

Encoder state for token $k$: $s_k$ $\qquad$ Decoder state at step $t$: $h_t$

softmax

I saw a

$p^{(1)}\ p^{(2)}\ p^{(3)}\ p^{(4)}$

Я видел котю на мате <eos>

"I" "saw" "cat" "on" "mat"

<bos> I saw a ...

**Encoder** $\qquad$ **Decoder**

# Attention



Attention

$\text{score}(h_t, s_k)$

○

scalar out

Attention function

in

in

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

Encoder state for token $k$: $s_k$

Decoder state at step $t$: $h_t$

At each decoder step, attention

- receives **attention input**: a decoder state $h_t$ and all encoder states $s_1$, $s_2$, ..., $s_m$;
- computes **attention scores**
  For each encoder state $s_k$, attention computes its "relevance" for this decoder state $h_t$. Formally, it applies an attention function which receives one decoder state and one encoder state and returns a scalar value $score(h_t, s_k)$;
- computes **attention weights**: a probability distribution - softmax applied to attention scores;
- computes **attention output**: the weighted sum of encoder states with attention weights.

# Attention

Attention output

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

$\uparrow$
"source context for decoder step $t$"

(weighted sum)

Attention weights

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

$\uparrow$
"attention weight for source token $k$ at decoder step $t$"

(softmax)

Attention scores

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

$\uparrow$
"How relevant is source token $k$ for target step $t$?"

Attention input

$$s_1, s_2, \ldots, s_m \qquad h_t$$

all encoder states          one decoder state

# Attention

Attention output

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

↑
"source context for decoder step $t$"

(weighted sum)

Attention weights

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

↑
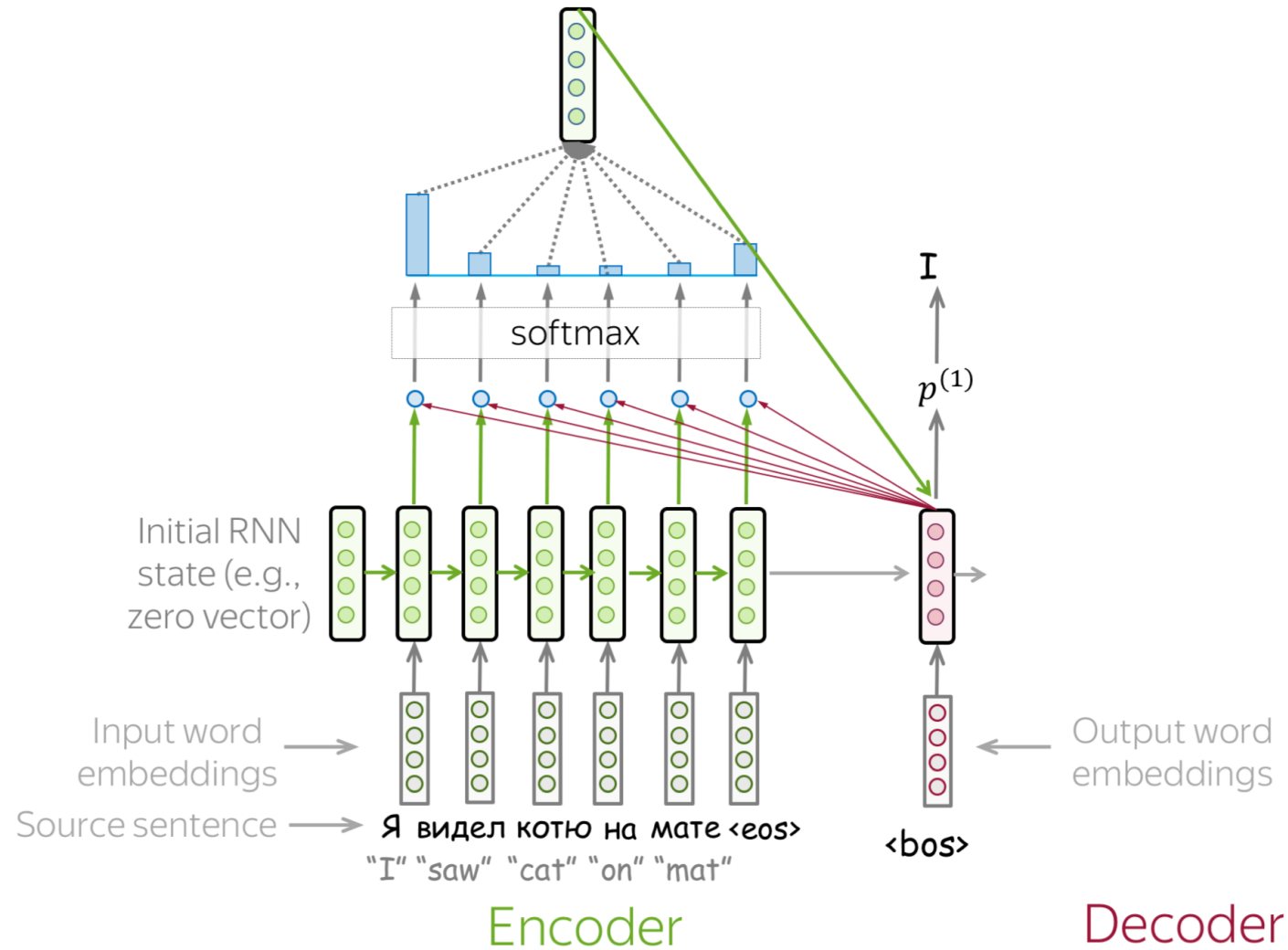"attention weight for source token $k$ at decoder step $t$"

Atte

The model learns to choose relevant input tokens for every step

the computation is differentiable so
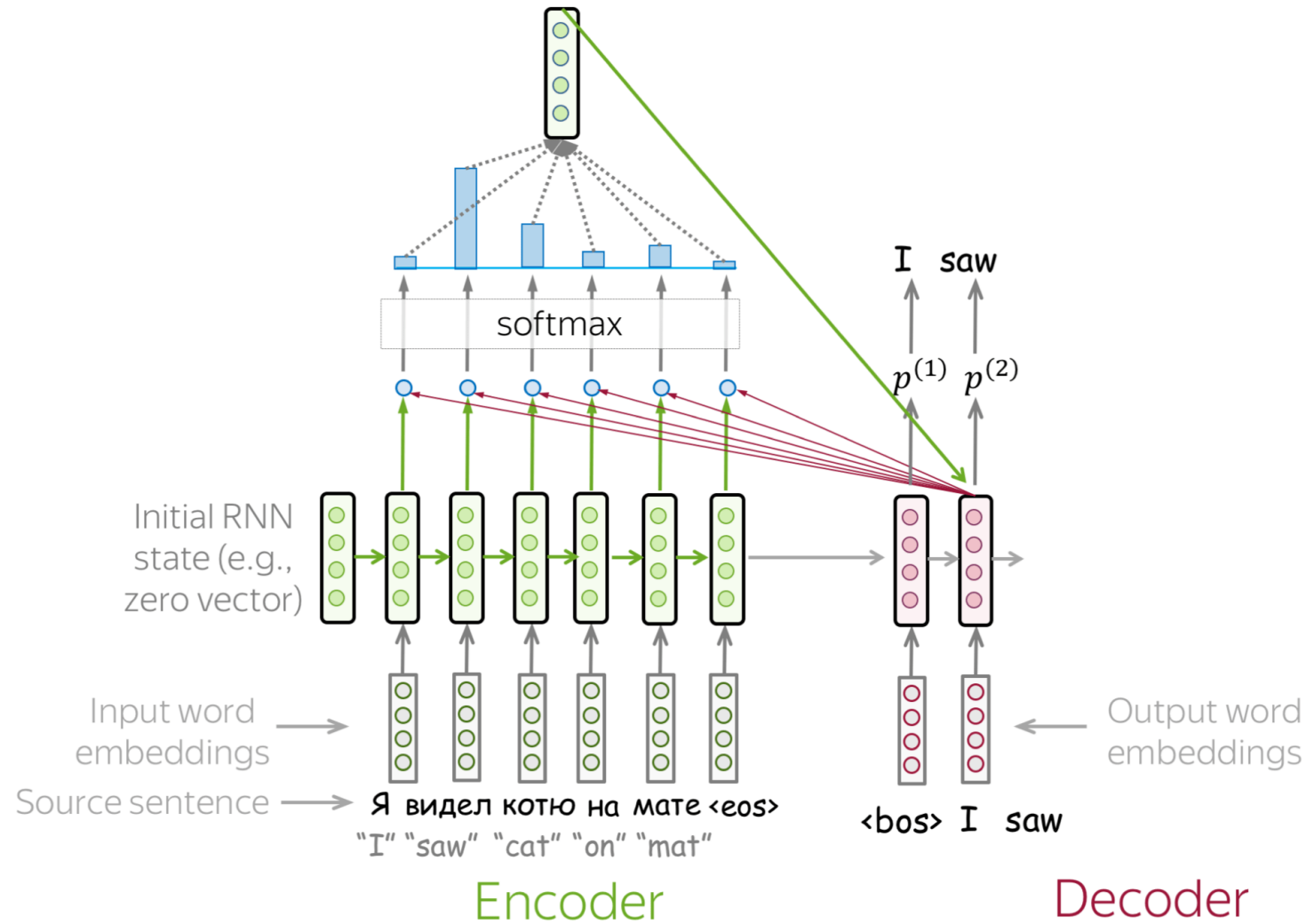everything here is learned end-to-end
with backpropagation

?"

Atte
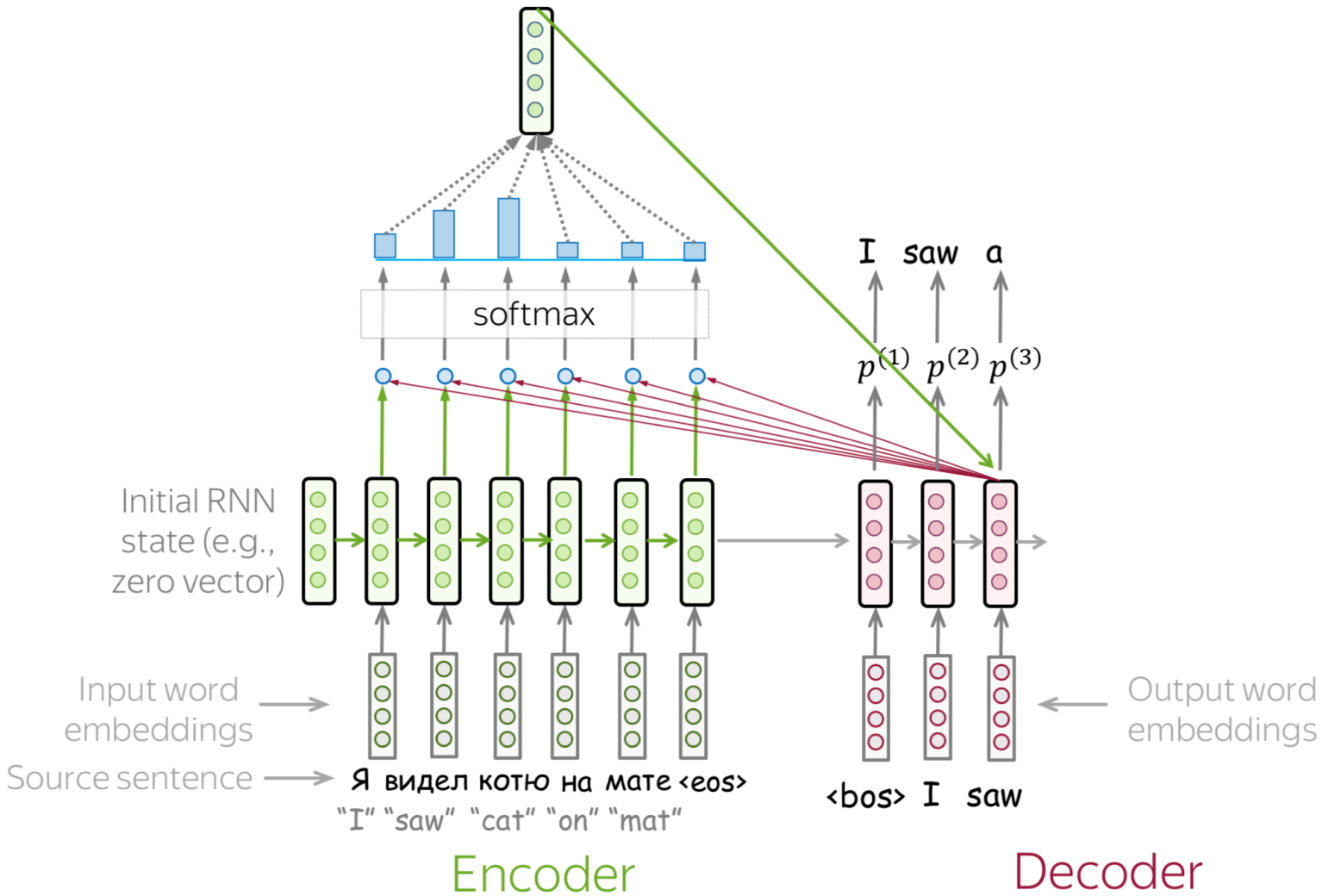
# Attention: step by step



softmax

I

$p^{(1)}$

Initial RNN state (e.g., zero vector)

Input word embeddings

Source sentence

Я видел котю на мате <eos>

"I" "saw" "cat" "on" "mat"

<bos>

Output word embeddings

Encoder

Decoder

# Attention: step by step



softmax

$p^{(1)}$  $p^{(2)}$

I  saw

Initial RNN
state (e.g.,
zero vector)

Input word
embeddings

Source sentence

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

<bos> I saw

Output word
embeddings

Encoder

Decoder

# Attention: step by step



I saw a

$p^{(1)}$ $p^{(2)}$ $p^{(3)}$

softmax

Initial RNN state (e.g., zero vector)

Input word embeddings

Source sentence

Я видел котю на мате ‹eos›
"I" "saw" "cat" "on" "mat"

‹bos› I saw

Output word embeddings

Encoder

Decoder

# Attention: step by step



Encoder

Decoder

Initial RNN state (e.g., zero vector)

Input word embeddings

Source sentence

softmax

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

<bos> I saw a

I saw a cat

$p^{(1)}$ $p^{(2)}$ $p^{(3)}$ $p^{(4)}$

Output word embeddings

# Attention: step by step



Initial RNN state (e.g., zero vector)

Input word embeddings

Source sentence

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

softmax

Encoder

I saw a cat on

$p^{(1)}$ $p^{(2)}$ $p^{(3)}$ $p^{(4)}$ $p^{(5)}$

<bos> I saw a cat

Output word embeddings

Decoder

# Attention: step by step



Encoder

Decoder

# Attention: step by step



Initial RNN state (e.g., zero vector)

Input word embeddings

Source sentence

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

softmax

$p^{(1)}$ $p^{(2)}$ $p^{(3)}$ $p^{(4)}$ $p^{(5)}$ $p^{(6)}$ $p^{(7)}$

I saw a cat on a mat

<bos> I saw a cat on a

Output word embeddings

Encoder

Decoder

# Attention: step by step



softmax

$p^{(1)}$ $p^{(2)}$ $p^{(3)}$ $p^{(4)}$ $p^{(5)}$ $p^{(6)}$ $p^{(7)}$ $p^{(8)}$

I saw a cat on a mat <eos>

Initial RNN state (e.g., zero vector)

Input word embeddings

Source sentence

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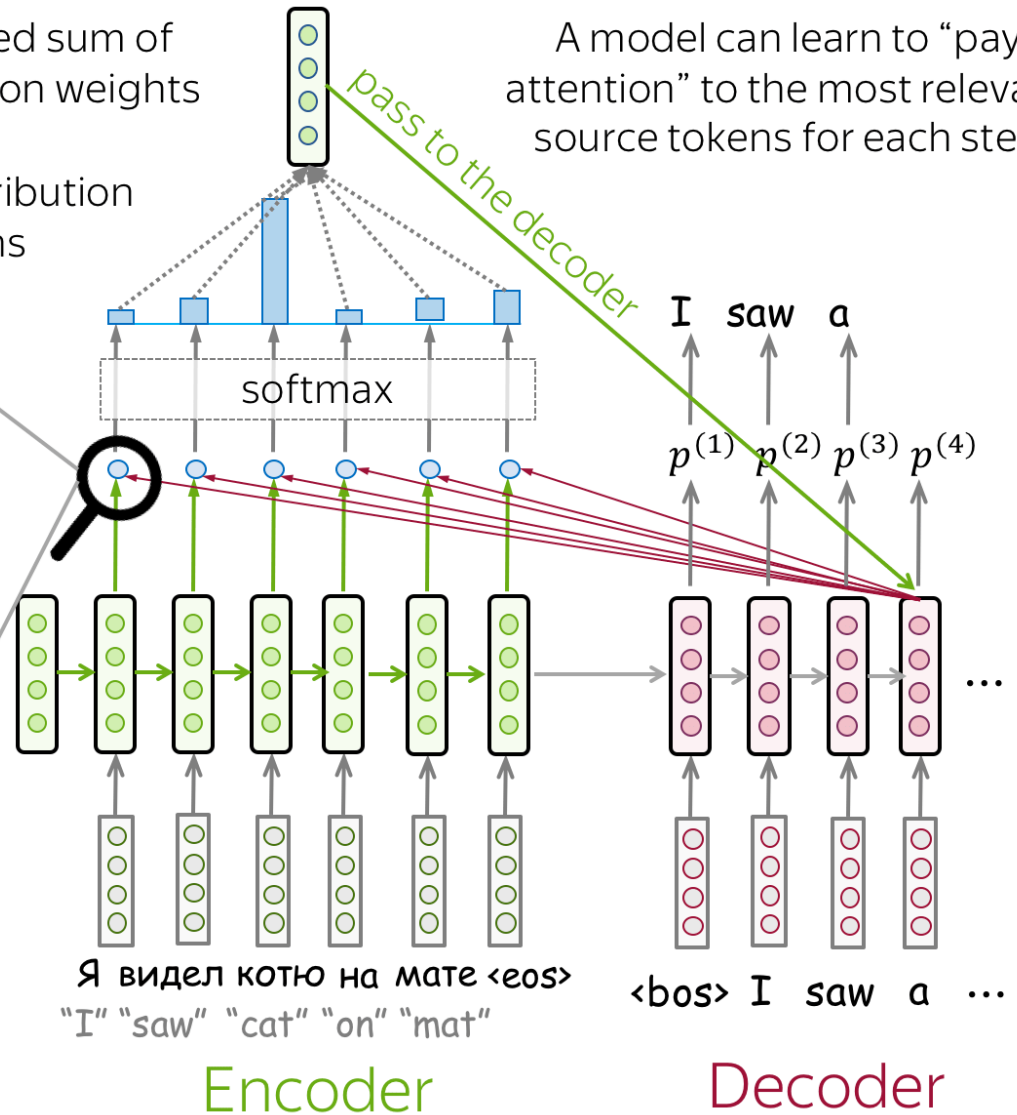"I" "saw" "cat" "on" "mat"

<bos> I saw a cat on a mat

Output word embeddings
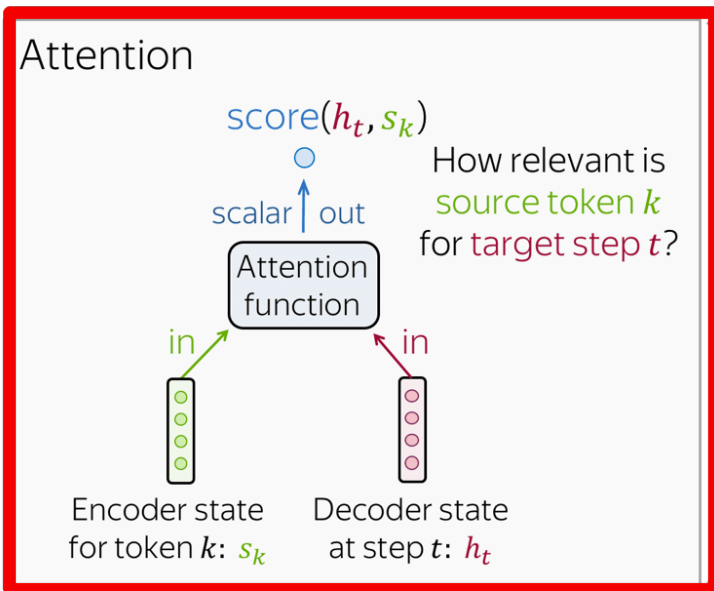
Encoder

Decoder

# Attention

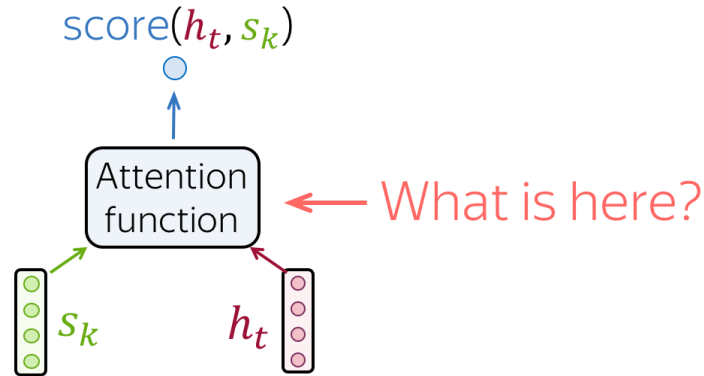Attention output: weighted sum of encoder states with attention weights

Attention weights: distribution over source tokens

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

pass to the decoder

softmax

I  saw  a

$p^{(1)}$  $p^{(2)}$  $p^{(3)}$  $p^{(4)}$

## Attention

$score(h_t, s_k)$

scalar ↑ out

Attention function

in        in

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

Encoder state for token $k$: $s_k$

Decoder state at step $t$: $h_t$

**What goes in here?**

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

<bos>  I  saw  a  ...

Encoder

Decoder

# How do we compute attention scores?  Alternatives

$\text{score}(h_t, s_k)$

Attention function

← What is here?

$s_k$     $h_t$

## Dot-product

$h_t^T$     $s_k$

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

## Bilinear

$h_t^T$  ×  W  ×  $s_k$

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

aka 'Luong attention'

## Multi-Layer Perceptron

$w_2^T$  ×  tanh  $W_1$  ×  $\begin{bmatrix} h_t \\ s_k \end{bmatrix}$

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

aka 'Bahdanau' attention
(from the original paper)

# Encoder-Decoder variants: Bahdanau



Multi-Layer Perceptron

$$w_2^T \times \tanh \left[ W_1 \times \begin{bmatrix} h \\ s_k \end{bmatrix} \right]$$
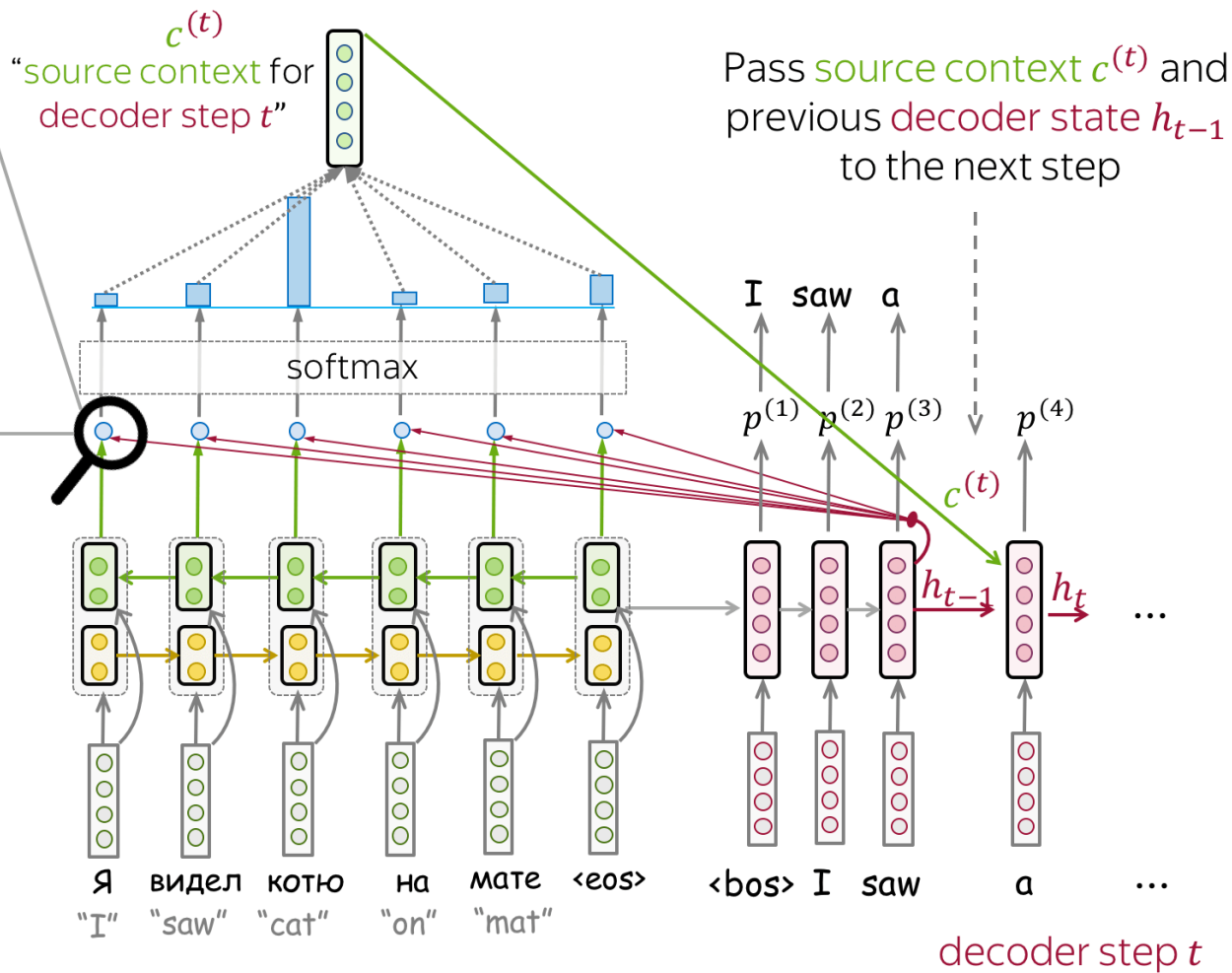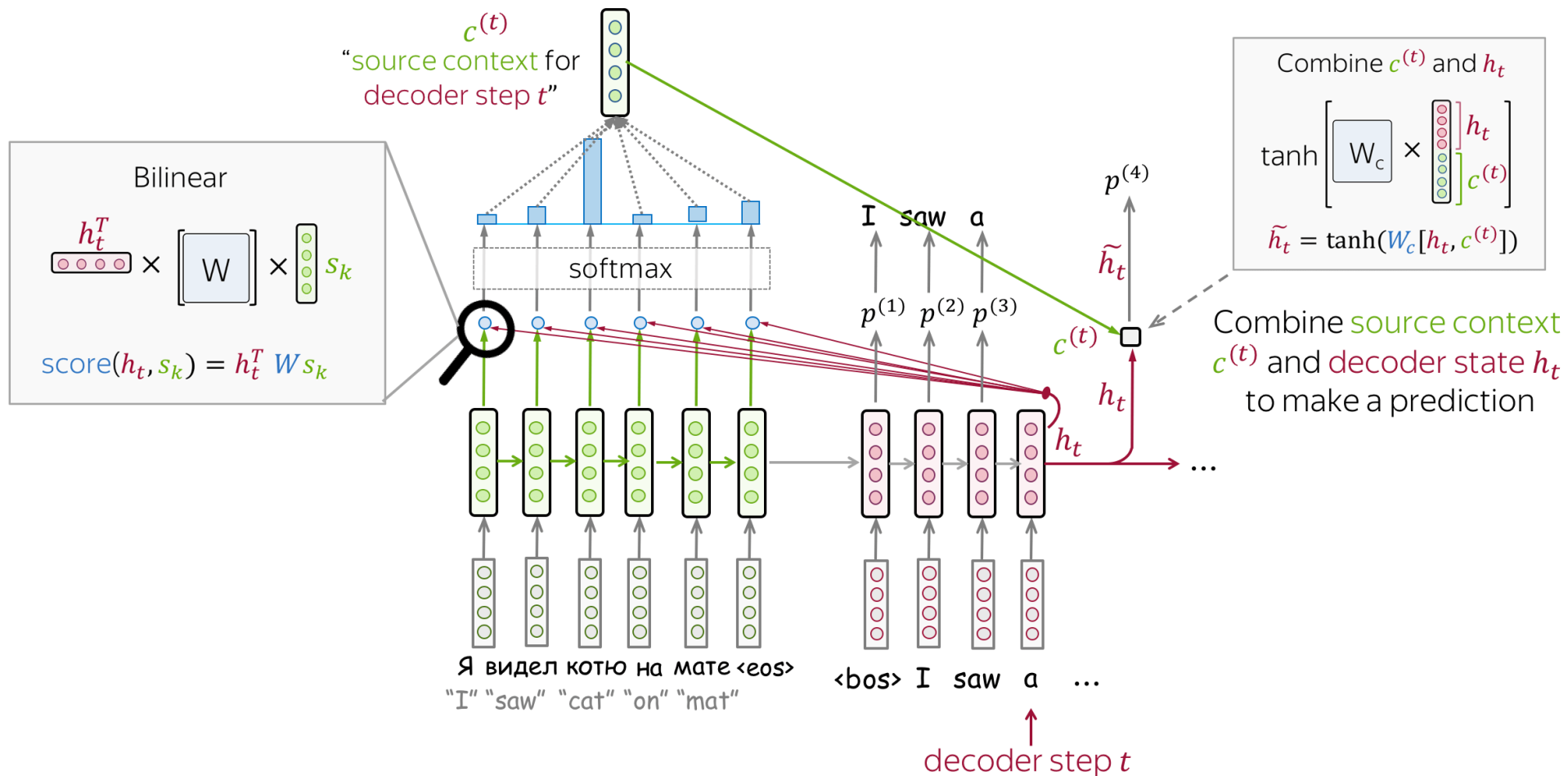
$$\text{score}(h, s_k) = w_2^T \cdot \tanh(W_1[h, s_k])$$

$c^{(t)}$
"source context for decoder step $t$"

Pass source context $c^{(t)}$ and previous decoder state $h_{t-1}$ to the next step

softmax

I saw a

$p^{(1)}$ $p^{(2)}$ $p^{(3)}$ $p^{(4)}$

$c^{(t)}$

**Bidirectional encoder**
Concatenate states from forward and backward RNNs

$h_{t-1}$ $h_t$ ...

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

decoder step $t$

# Encoder-Decoder variants: Luong



$c^{(t)}$
"source context for decoder step $t$"

Bilinear

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

softmax

I saw a

$p^{(4)}$

$\widetilde{h}_t$

$p^{(1)} \ p^{(2)} \ p^{(3)}$

$c^{(t)}$

$h_t$

$h_t$

Combine $c^{(t)}$ and $h_t$

$$\tanh \begin{bmatrix} W_c & \times & \begin{bmatrix} h_t \\ c^{(t)} \end{bmatrix} \end{bmatrix}$$

$$\widetilde{h}_t = \tanh(W_c[h_t, c^{(t)}])$$

Combine source context $c^{(t)}$ and decoder state $h_t$ to make a prediction

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

<bos> I saw a ...

decoder step $t$

https://arxiv.org/abs/1508.04025

# Is attention interpretable?

- The attention may be perceived as modelling alignment between input and output words, but
  - decoder computes attention **before** generating the target token (i.e. the choice of the token does not influence attention)
  - sometimes states encode something unexpected (e.g., <eos> may capture the general topic of the sentence)
  - attention is (?) not an explanation

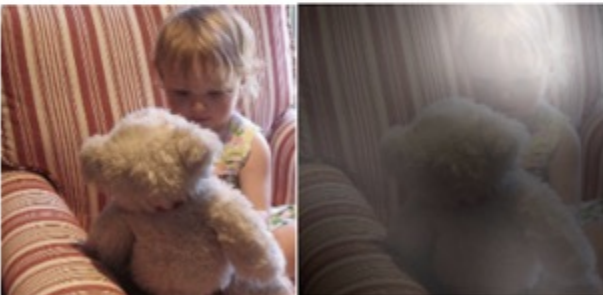# Enoder-decoders and attention are a very general idea

For example, caption generation



A woman is throwing a <u>frisbee</u> in a park.

A <u>dog</u> is standing on a hardwood floor.

A <u>stop</u> sign is on a road with a mountain in the background.

A little <u>girl</u> sitting on a bed with a teddy bear.

A group of <u>people</u> sitting on a boat in the water.

A giraffe standing in a forest with <u>trees</u> in the background.

From https://arxiv.org/pdf/1502.03044.pdf

# Attention is not necessarily faithful

Generally, you can make a change to model parameters such that attention score $a_k$ for a state $s_k$ decreases *N*-fold ($a'_k = s_k/N$) whereas $s_k$ magnitude increases *N*-fold ($s'_k = s_k N$)

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

↑
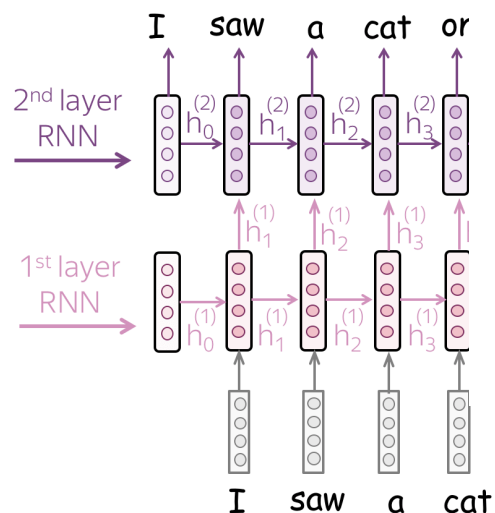"source context for decoder step $t$"

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

↑
"attention weight for source token $k$ at decoder step $t$"

Since the attention output is the weighted sum of embedding encoders states,  the attention output ($c$) will not change

# Attention is not necessarily faithful

Also, if we use a multilayer encoder (and we usually will), there is no guarantee that a state in a *n*-th layers (n > 1) encodes the n-th token
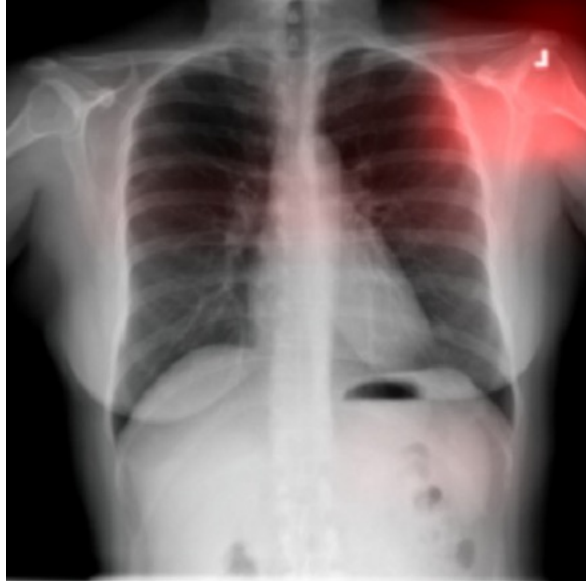


There are attribution methods (e.g., norm x gradient, layer-wise relevant propagation, integrated gradients, attention flow, zero valuing, …) which attempt to address these issues, but they also have pitfalls

Generally: attention cannot be trusted blindly but it can signal some issues with our models

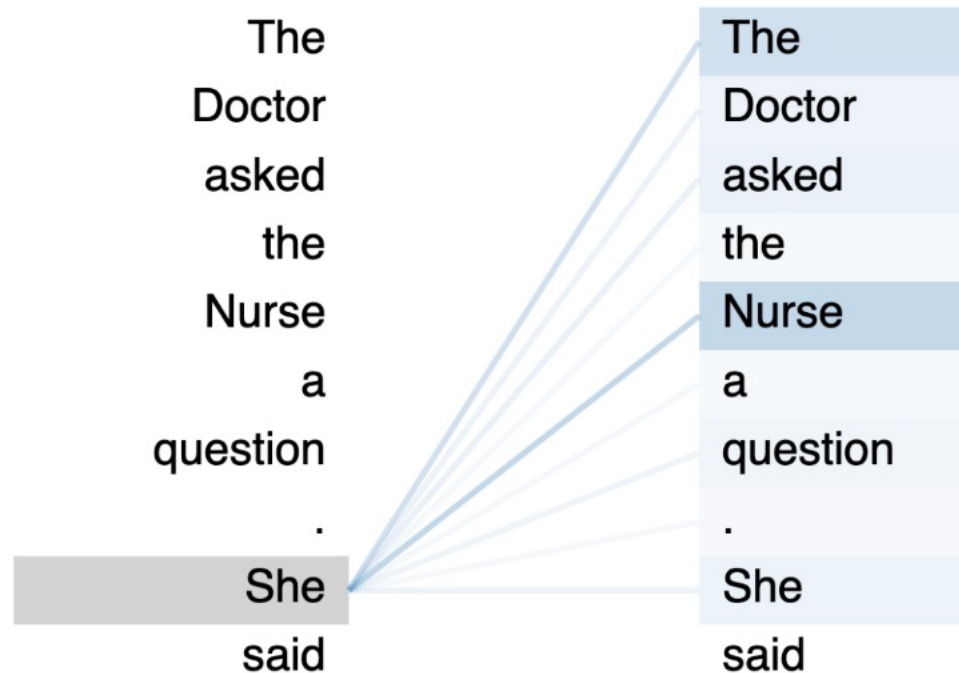# Attribution can help us detect issues with our models

Even if not perfectly faithful attention and attribution techniques help us detect issues with our models

Detecting pneumonia from an x-ray, the model 'looks' at the corner of an image, rather than at lungs.   Any thoughts on why?



Zech et al (PLOS Medecine 2018)

# Attribution can help us detect issues with our models

The model decides that pronoun 'she' refers to 'Nurse' rather than Doctor (gender bias)

It is not an encoder-decoder attention, it is attention with a language model (i.e. what we will consider on Friday)

# Transformer

Attention is all you need

| | 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 |

# Summary

- The general idea of attention
- Attention between encoder and decoder
- Is attention alignment / is it interpretable?
- Many applications in NLP and beyond