# Natural Language Processing: Seq2seq and Attention

HSE Faculty of Computer Science

Machine Learning and Data-Intensive Systems

Murat Khazhgeriev

# Table of Content

- **Sequence to Sequence (seq2seq)**

- Attention

- Practical tips

# Machine translation requires training parameters to provide results

Human Translation

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

The "probability" is intuitive and is given by a human translator's expertise

Machine Translation

model      parameters

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

Questions we need to answer

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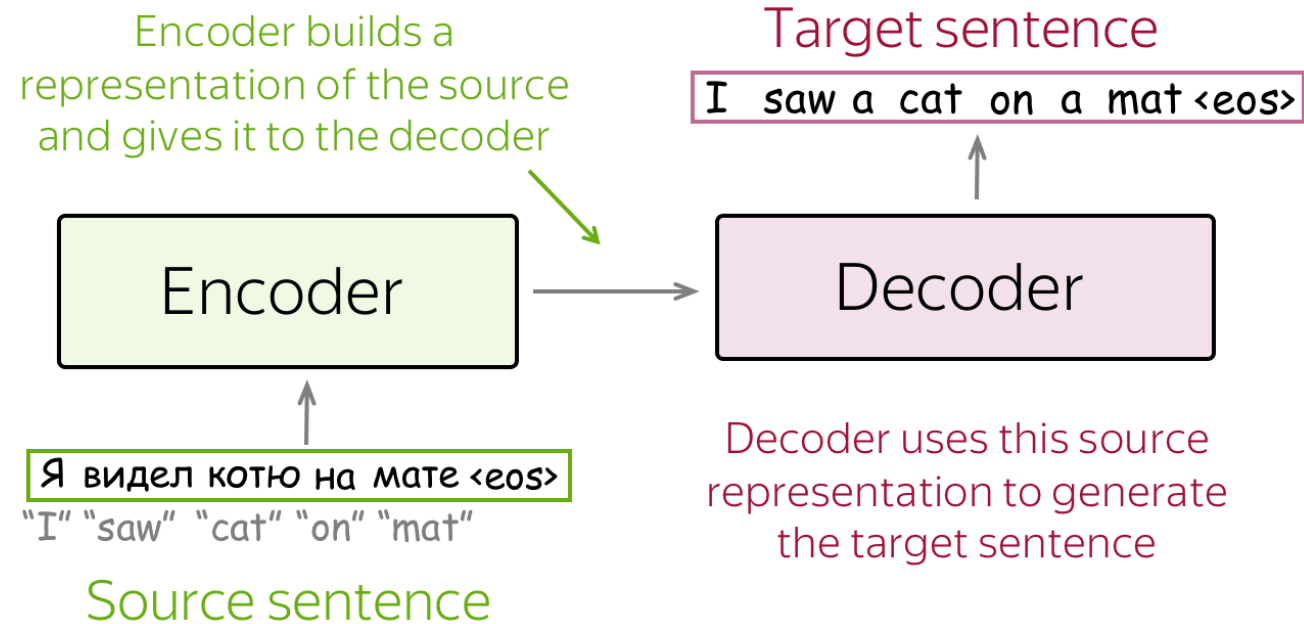
  How does the model for $p(y|x,\theta)$ look like?

- **learning**

  How to find $\theta$?

- **search**

  How to find the argmax?

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

3

# Encoder-decoder architecture maps data semantics



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

Target sentence

I saw a cat on a mat <eos>

Encoder

Decoder

Я видел котю на мате <eos>

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

Source sentence

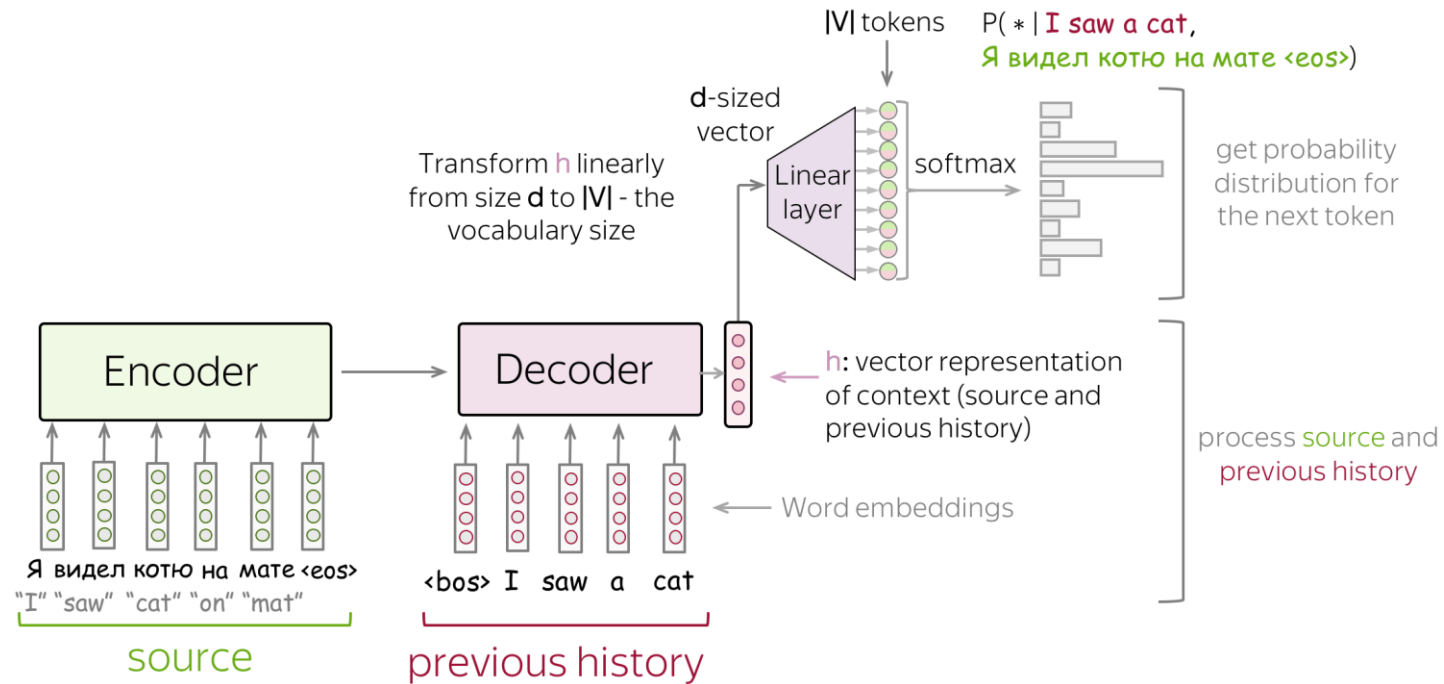Decoder uses this source representation to generate the target sentence

# Encoder encapsulates a condition for a decoder Language Model

Language Models: $P(y_1, y_2, \ldots, y_n) = \prod_{t=1}^{n} p(y_t | y_{<t})$

Conditional

Language Models: $P(y_1, y_2, \ldots, y_n, |x) = \prod_{t=1}^{n} p(y_t | y_{<t}, x)$

condition on source $x$

# A helicopter view on the Encoder-Decoder architecture

# Model Loss is a well-known Cross-Entropy

Source sequence:

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

Target sequence:        ←——— one training example

I saw a cat on a mat <eos>   ←——— one step for this example

previous tokens    we want the model
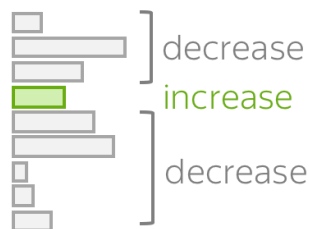                   to predict this
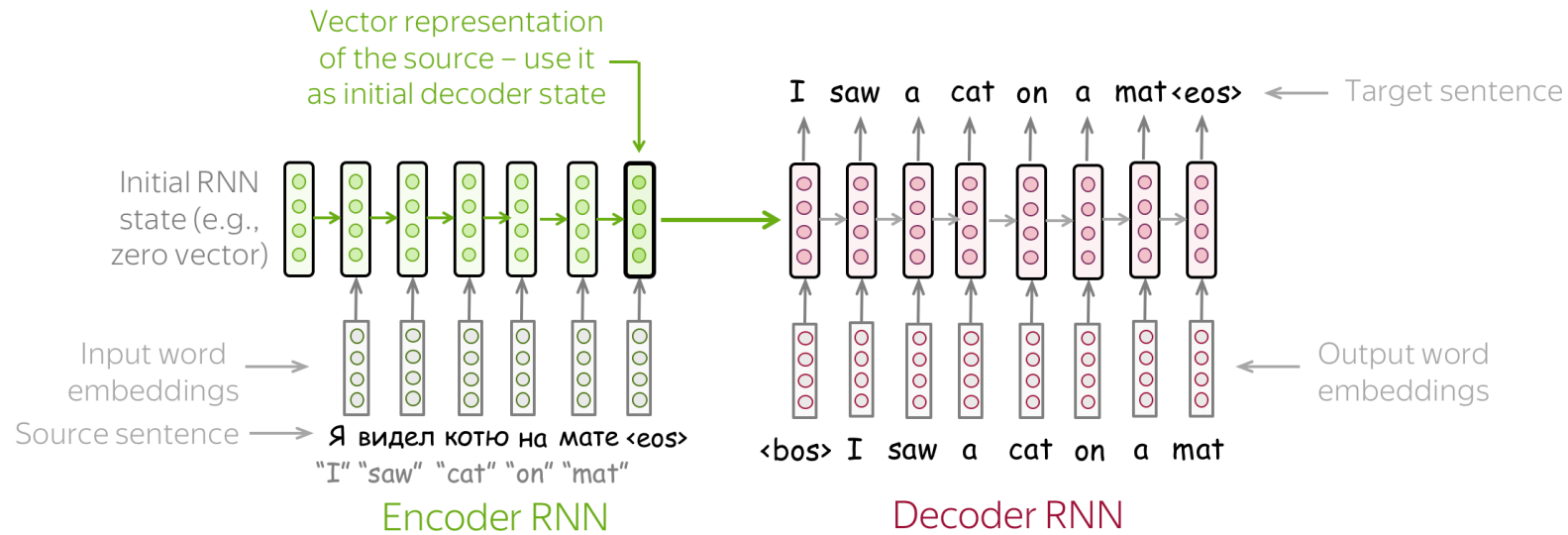
Model prediction: p( * |I saw a,   Target    Loss = -log (p(cat)) → min
Я ... <eos>)

0
0
0
←— cat —→  1    ] decrease
0            increase
0
0          ] decrease
0
0
0

# A simple RNN is a valid Encoder-Decoder model



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

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

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

Output word embeddings

<bos> I saw a cat on a mat
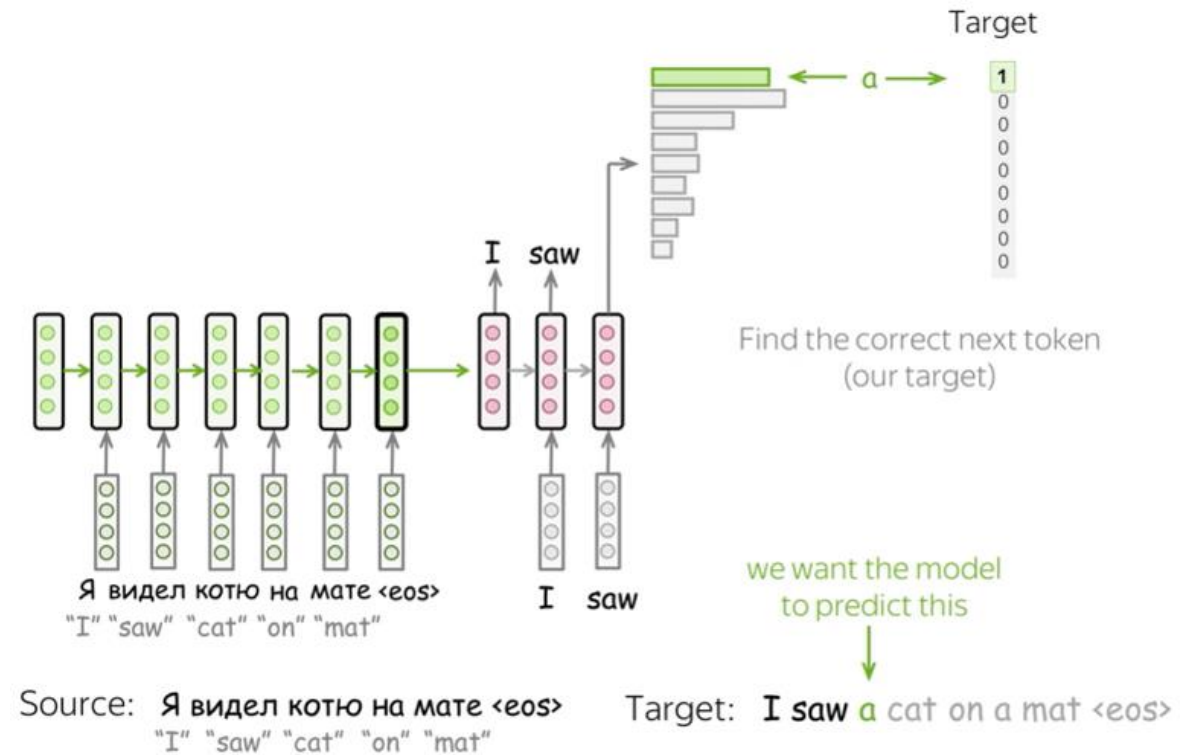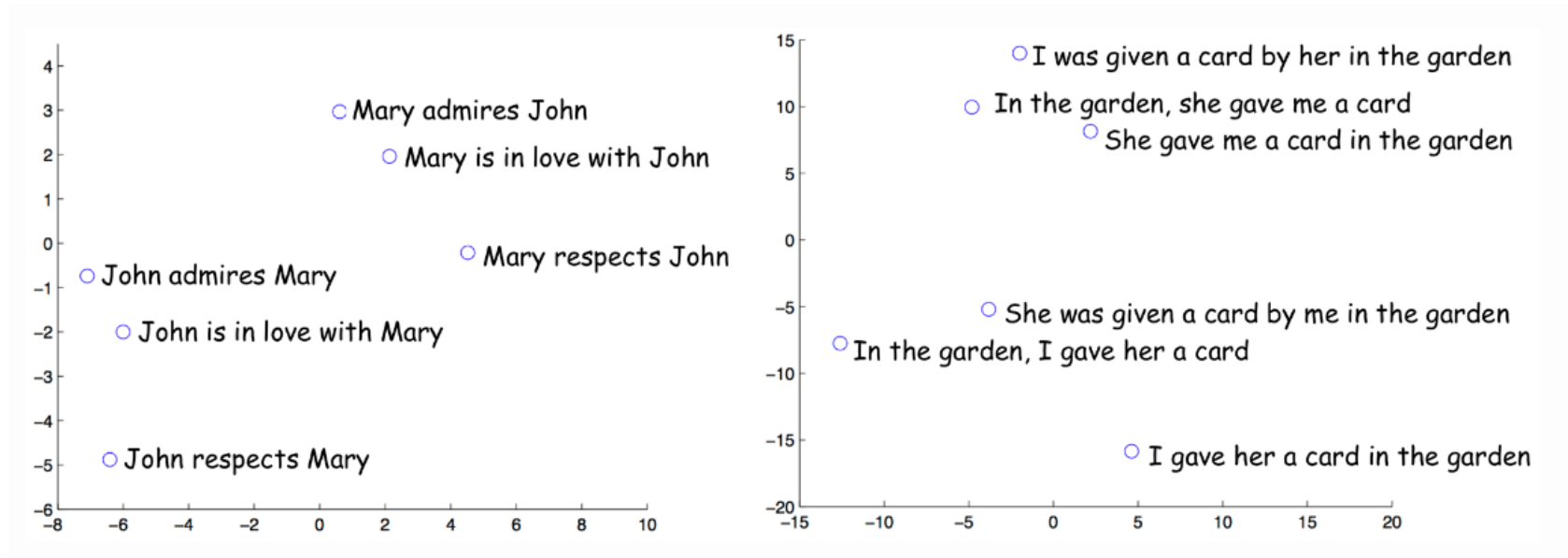
Decoder RNN

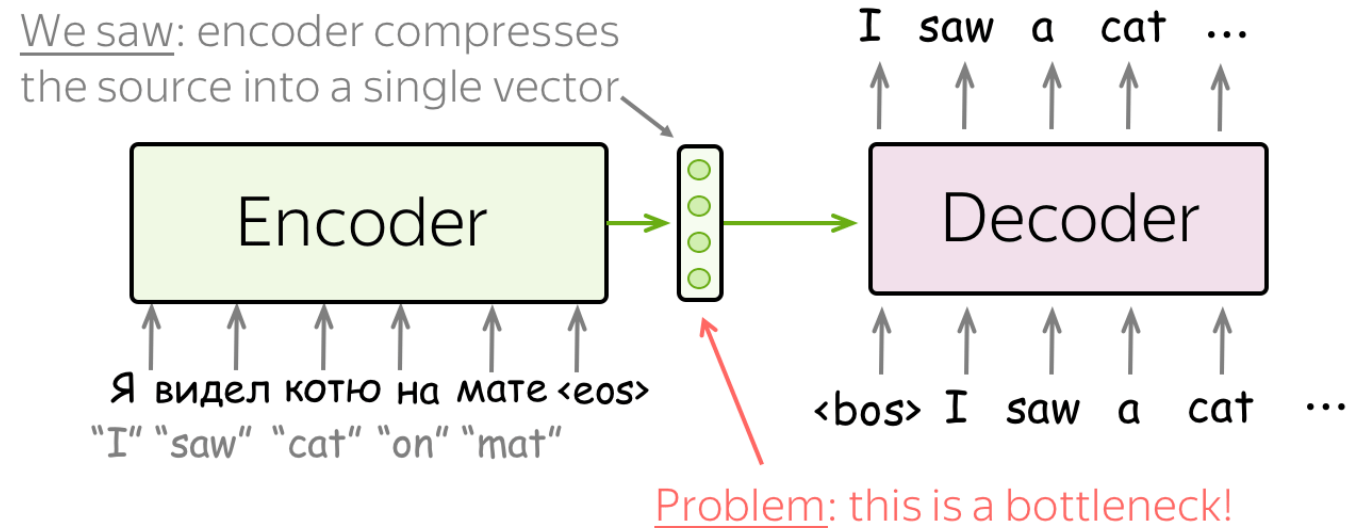# A simple RNN is a valid Encoder-Decoder model

# Semantic space of text embeddings

# Table of Content

- Sequence to Sequence (seq2seq)
- **Attention**
- Practical tips
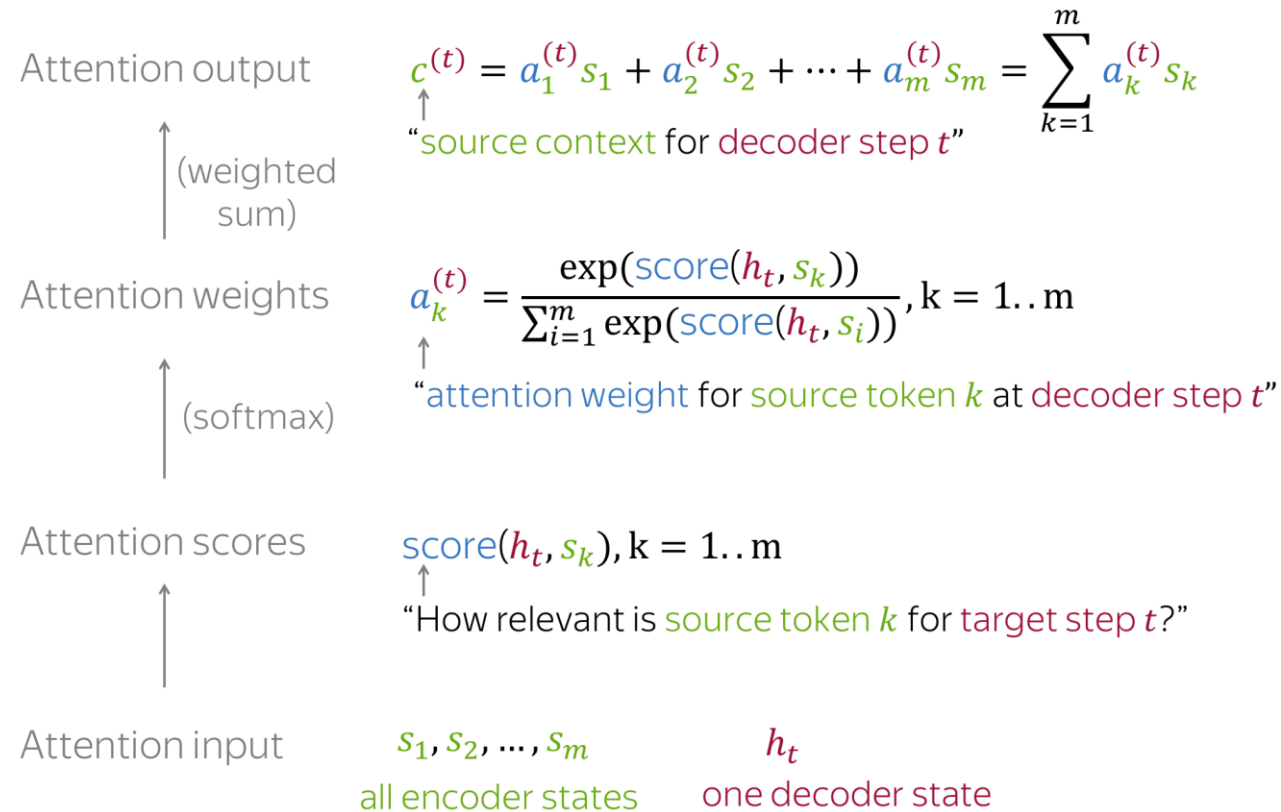
# The final hidden state is a bottleneck

We saw: encoder compresses the source into a single vector

I saw a cat ...

Encoder → Decoder

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

<bos> I saw a cat ...

Problem: this is a bottleneck!

# A simple Attention overview



**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

Attention

$score(h_t, s_k)$

scalar out

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

Attention function

in        in

Encoder state for token $k$: $s_k$        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        Decoder

# Encoder hidden states are weighed according to their attention score

**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**

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

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

(softmax)

**Attention scores**

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

"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

# Encoder hidden states are weighed according to their attention score



**Dot-product**

$$h_t^T \times s_k$$

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

**Bilinear**

$$h_t^T \times \boxed{W} \times s_k$$

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

**Multi-Layer Perceptron**

$$w_2^T \times \tanh\left(\boxed{W_1} \times \begin{bmatrix} h_t \\ s_k \end{bmatrix}\right)$$

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

# Table of Content

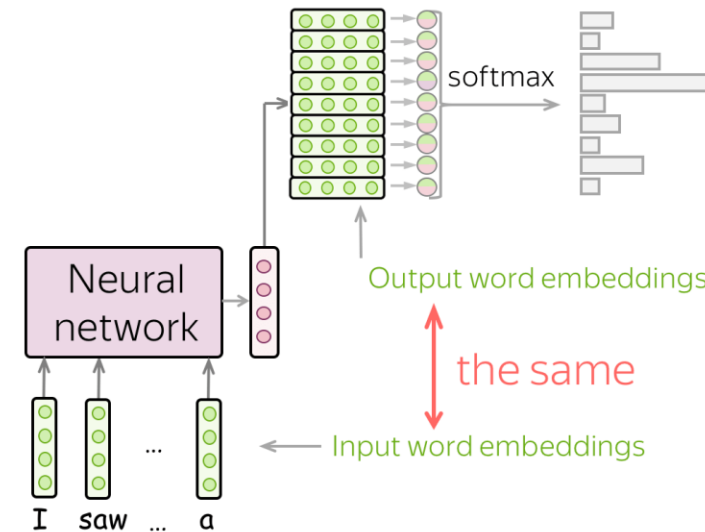- Sequence to Sequence (seq2seq)

- Attention

- **Practical tips**

# Weight tying is a way to significantly reduce the amount of parameters
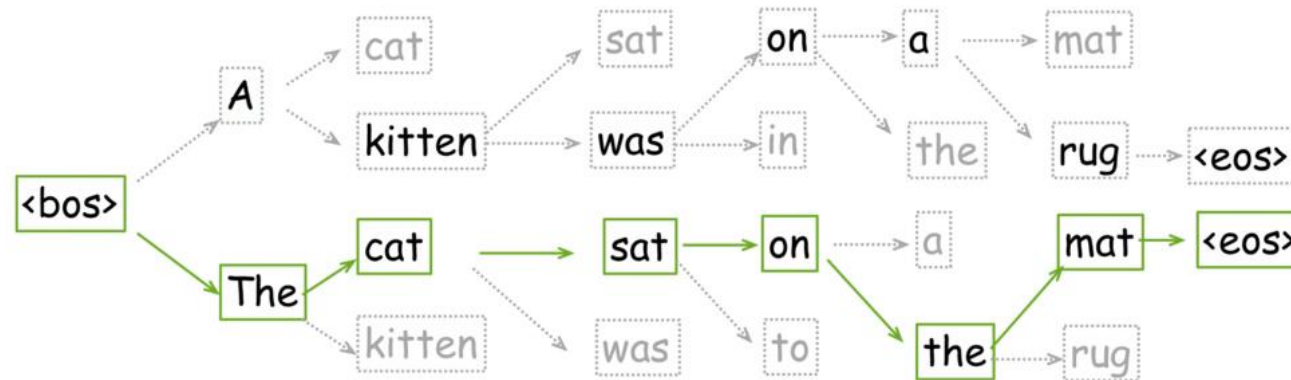
# Finding the next token is not that trivial

$$y' = \arg\max_y p(y|x) = \arg\max_y \prod_{t=1}^{n} p(y_t|y_{<t}, x)$$    How to find the argmax?

# Greedy Decoding: At each step, pick the most probable token

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

# A beam search illustration



Pick the hypothesis with the highest probability

# Temperature: the higher, the more chaotic the choice becomes

$$w_{next} \sim \frac{P(w_{next}|X)^{1/\tau}}{\sum_{\hat{w}} P(\hat{w}|X)^{1/\tau}}$$