

# Introduction to NLP



# Credits

- › YSDA NLP course (Elena Voita et al.):  
[https://github.com/yandexdataschool/nlp\\_course/tree/master](https://github.com/yandexdataschool/nlp_course/tree/master)
- › CS224N: Stanford NLP course (Cristopher Manning, Abigail See et al.):  
<https://web.stanford.edu/class/cs224n>
- › <https://github.com/bentrevett/pytorch-sentiment-analysis>
- › <https://github.com/bentrevett/pytorch-seq2seq>

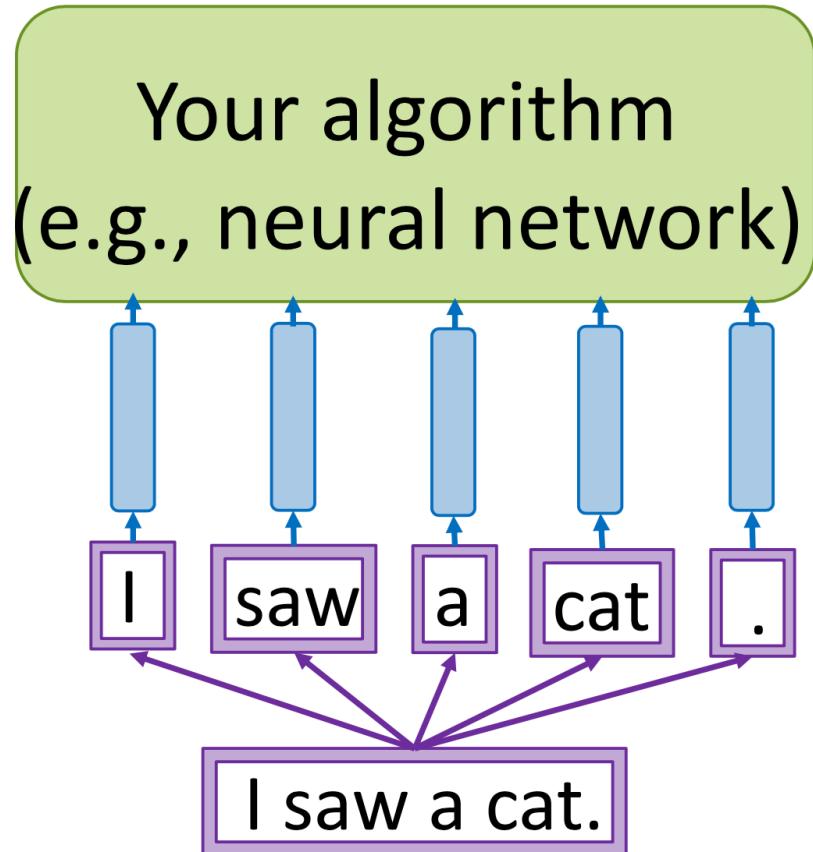
# NLP Problems

- › Text classification
- › Named Entity Recognition
- › Language models
- › Text Summarization
- › Machine Translation
- › Speech-to-Text
- › Question Answering
- › Dialogue Systems
- › and many more...

# NLP Problems

- › Text classification
- › Named Entity Recognition
- › Language models
- › Text Summarization
- › Machine Translation
- › Speech-to-Text
- › Question Answering
- › Dialogue Systems
- › and many more...

# Word representations



Any algorithm for solving any task

Word representation - vector  
(word embedding)

Sequence of tokens

Text

# Representing words as discrete symbols

- › We regard words as discrete symbols: `hotel`, `conference`, `motel`
- › Words can be represented by **one-hot vectors**:
  - › `motel=[000000000010000]`
  - › `hotel=[000000010000000]`
- › Vector dimension = number of words in vocabulary (e.g., 500,000)

# What is word meaning?

What is the meaning of the word “**bardiwac**”?

# What is word meaning?

What is the meaning of the word “**bardiwac**”?

- › "He handed her her glass of **bardiwac**"
- › "Beef dishes are made to complement the **bardiwacs**"
- › "Nigel staggered to his feet, face flushed from too much  
**bardiwac**"
- › "Malbec, one of the lesser-known **bardiwac** grapes, responds  
well to Australia's sunshine"
- › "I dined off bread and cheese and this excellent **bardiwac**"
- › "The drinks were delicious: blood-red **bardiwac** as well as  
light, sweet Rhenish"

# What is word meaning?

What is the meaning of the word “**bardiwac**”?

- › "He handed her her glass of **bardiwac**"
- › "Beef dishes are made to complement the **bardiwacs**"
- › "Nigel staggered to his feet, face flushed from too much **bardiwac**"
- › "Malbec, one of the lesser-known **bardiwac** grapes, responds well to Australia's sunshine"
- › "I dined off bread and cheese and this excellent **bardiwac**"
- › "The drinks were delicious: blood-red **bardiwac** as well as light, sweet Rhenish"



*"You shall know a word by the company it keeps."* - Firth (1957)

**Bardiwac** is a heavy red alcoholic beverage made from grapes

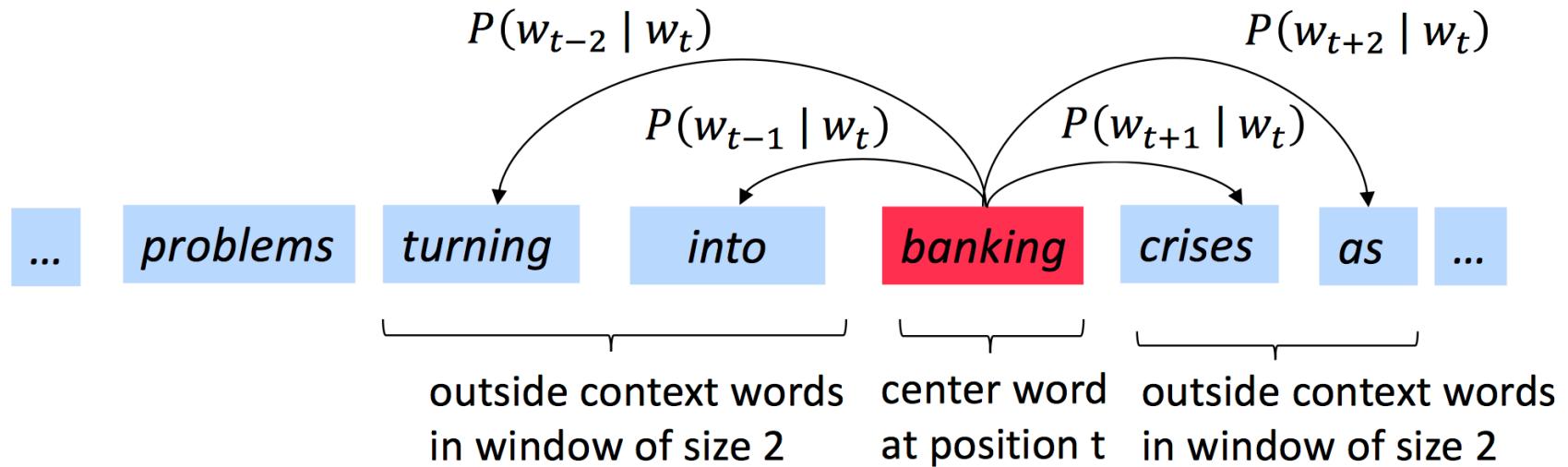
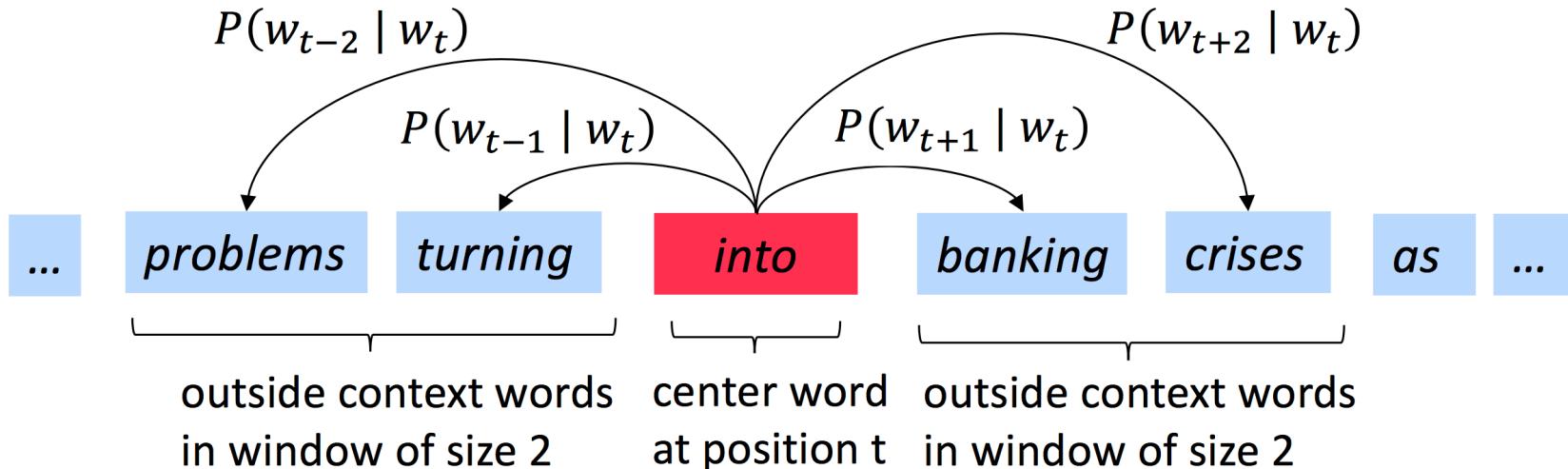


# Word2Vec

We will learn a dense vector for each word, chosen so that it is similar to vectors of words that appear in similar contexts.

- › Large text corpus
- › Every word in a fixed vocabulary is represented by a vector (an embedding-to-be)
- › Go through each position  $t$  in the text, which has a center word  $c$  and context (“outside”) words  $o$
- › Use the similarity of the word vectors for  $c$  and  $o$  to calculate the probability of  $o$  given  $c$  (or vice versa)
- › Keep adjusting the word vectors to maximize this probability

# Word2Vec



# Word2Vec Objective

For each position  $t = 1, \dots, T$ , predict context words within a window of fixed size  $m$ , given center word  $w_t$ .

$$\text{Likelihood} = L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_{t+j} | w_t; \theta)$$

$\theta$  is all variables  
to be optimized

The **objective function** is the (average) negative log likelihood:

$$J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_t; \theta)$$

# Word2Vec Objective

We want to minimize the objective function:

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_t; \theta)$$

Question: How to calculate  $P(w_{t+j} | w_t; \theta)$  ?

Answer: We will use two vectors per word  $w$ :

- $v_w$  when  $w$  is a center word
- $u_w$  when  $w$  is a context word

Exponentiation makes anything positive

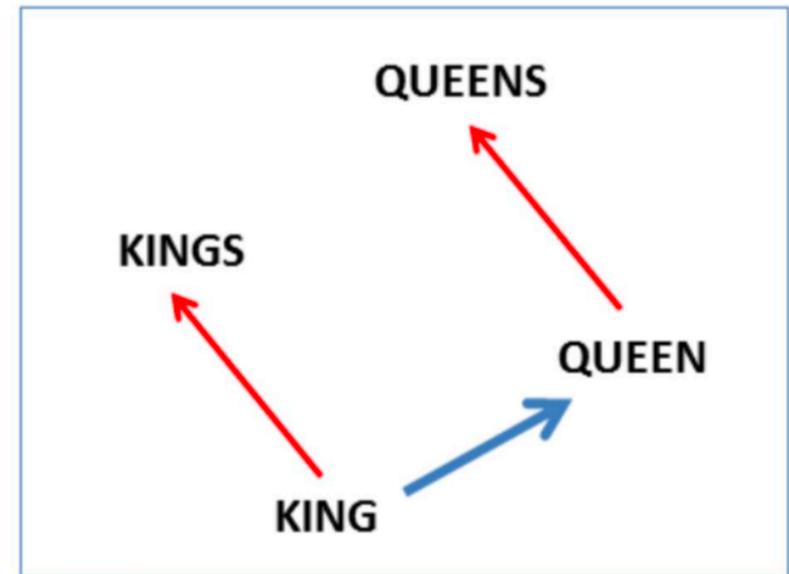
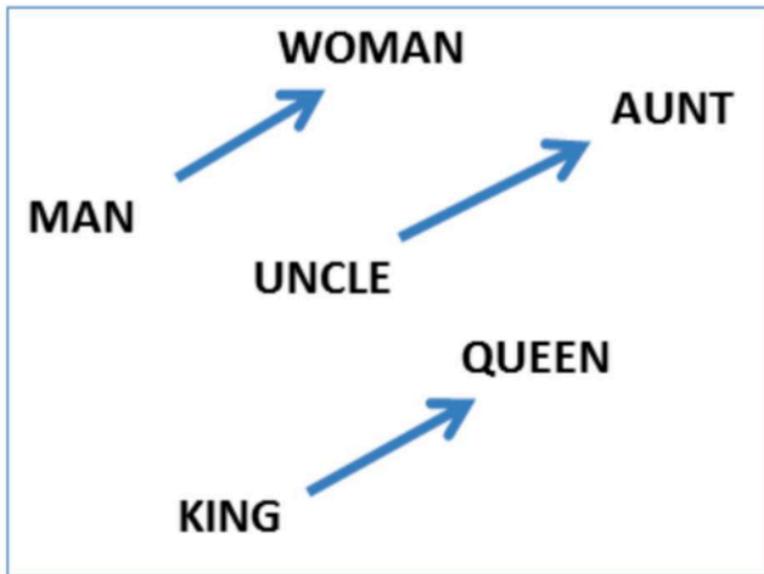
$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

Dot product compares similarity of  $o$  and  $c$ .  
 $u^T v = u \cdot v = \sum_{i=1}^n u_i v_i$   
Larger dot product = larger probability

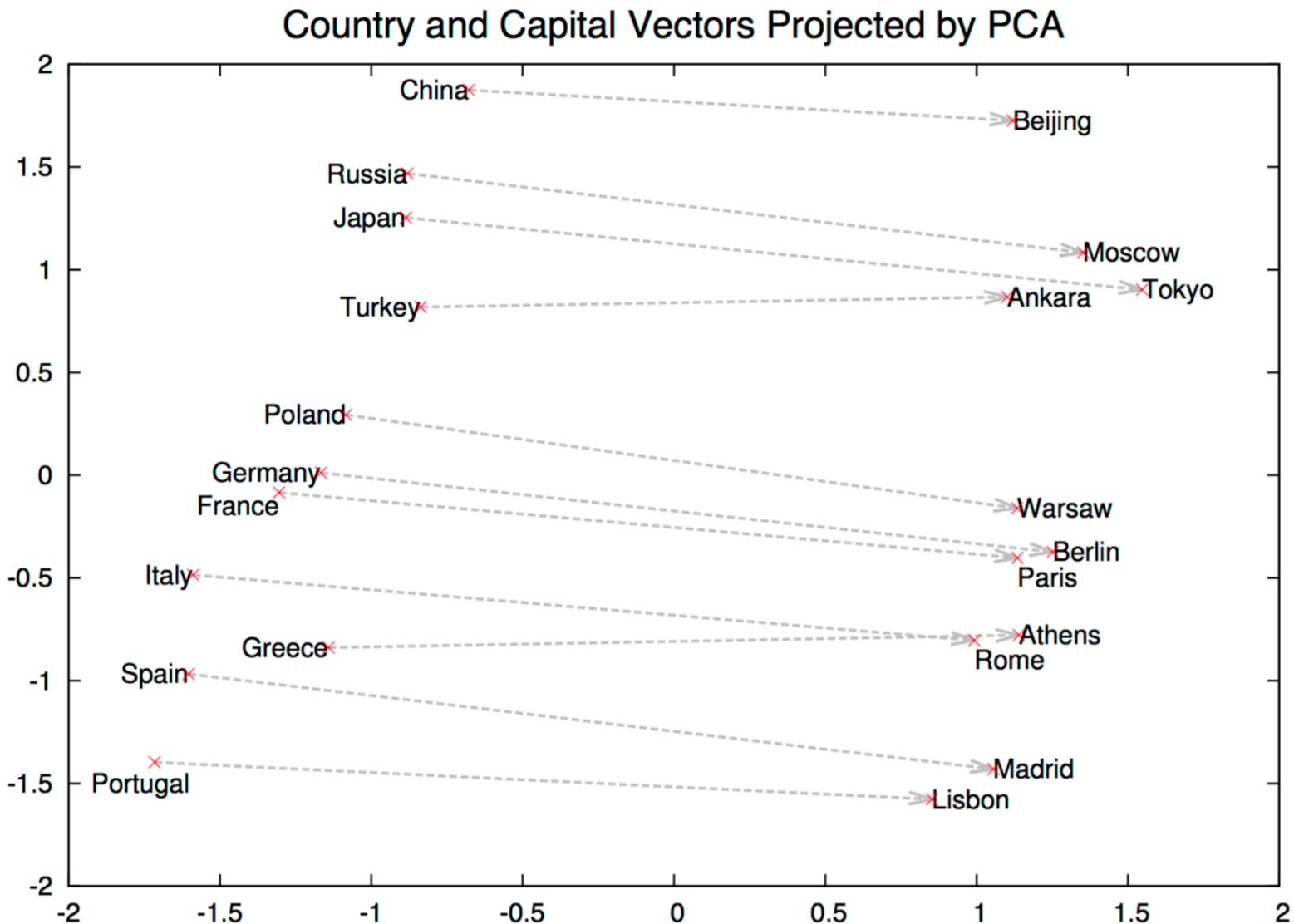
Normalize over entire vocabulary  
to give probability distribution

# Word2Vec: Relations between vectors

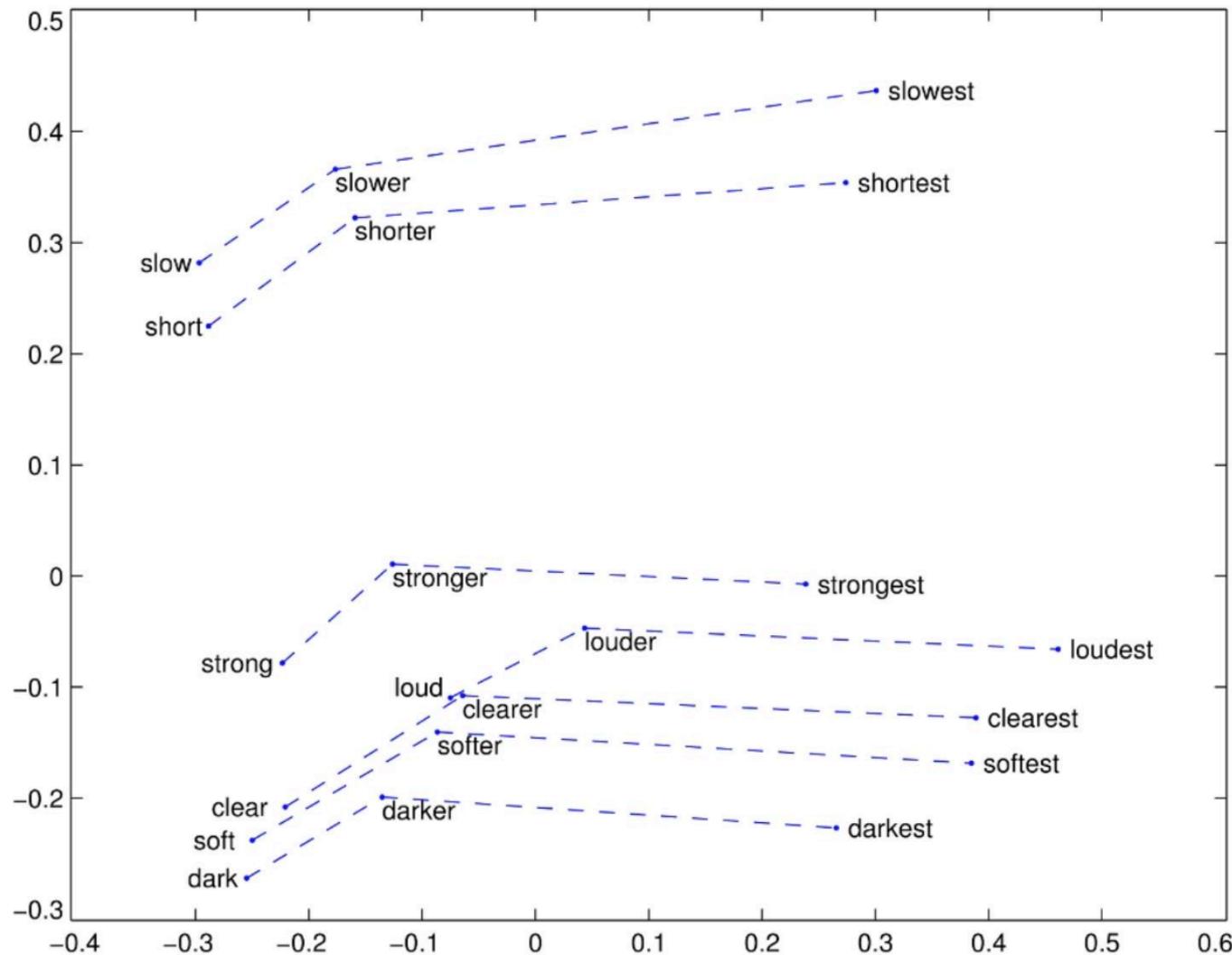
$$v(\text{king}) - v(\text{man}) + v(\text{woman}) \approx v(\text{queen})$$



# Word2Vec: Relations between vectors



# Word2Vec: Relations between vectors



# Text Encoders in NLP

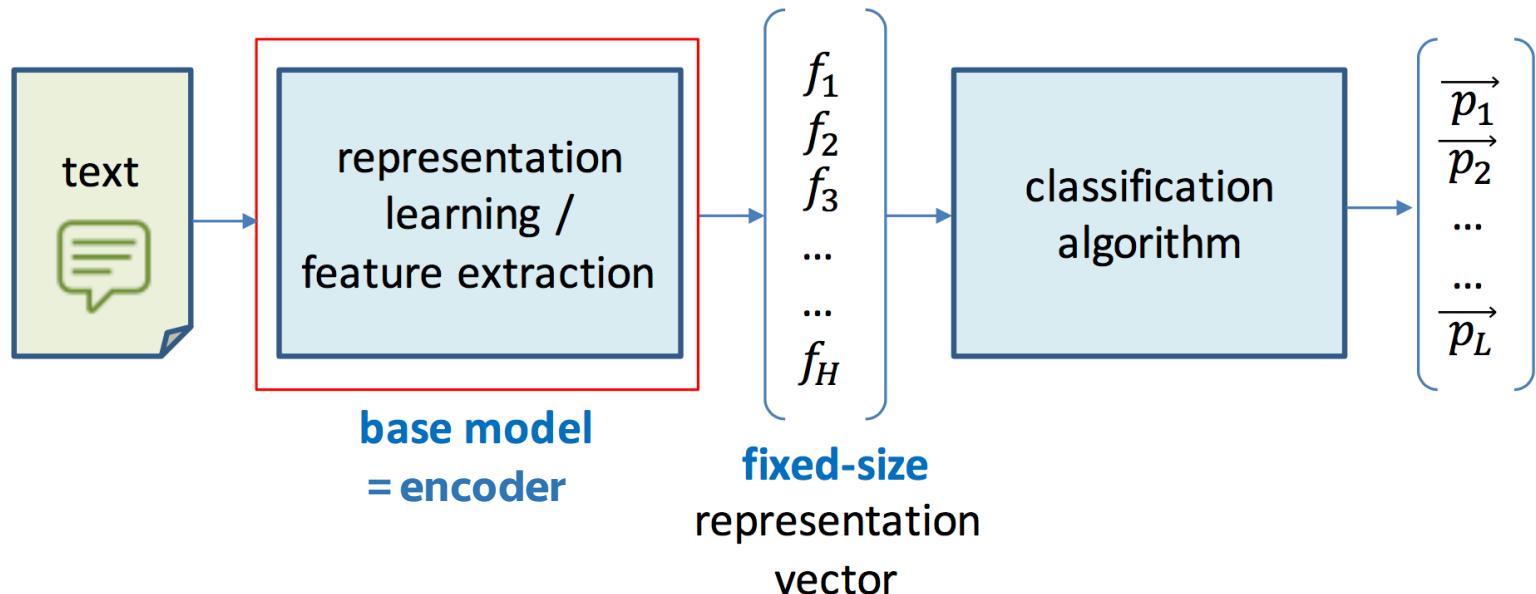
Common case: initializing the embedding layer with pre-trained embeddings

- › Vanilla [RNN](#)
- › Fancy [LSTM](#), [BiLSTM](#), [GRU](#)
- › Transformer-based encoders: [ELMo](#), [BERT](#), [RoBERTa](#), [XLNet](#), [T5](#), ...



# Text classification problem

- › Dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  with  $x \in X, y \in Y$
- › Label set  $\mathcal{L} = \{1, 2, \dots, L\}$
- › The pairs  $(x_n, y_n)$  are i.i.d. following  $P(\mathbf{X}, \mathbf{Y})$  over a sample space  $X \times Y$
- › Multi-class classification:  $Y = \mathcal{L}$
- › Multi-label classification:  $Y = 2^{\mathcal{L}}$



# Language models

- › Language model: estimates the “probability” of a text

$$P(\mathbf{y}) = P(y_1)P(y_2|y_1)P(y_3|y_1, y_2) \dots P(y_T|y_1, y_2, \dots, y_{T-1})$$

The screenshot shows the Yandex search engine interface. In the search bar, the query "a fat cat sat on" is typed, followed by a small blue link labeled "bn". To the right of the search bar is a yellow "Search" button with a magnifying glass icon. Below the search bar, five search results are displayed:

- a **fat cat sat** on a mat
- a **fat cat sat** on a mat and ate a fat rat
- a **fat cat sat** on a mat.
- a **fat cat sat** on a mat перевод на русский
- a **fat cat sat** on a mat and ate a fat rat перевод на русский

# Vanilla RNN language model

Output distribution

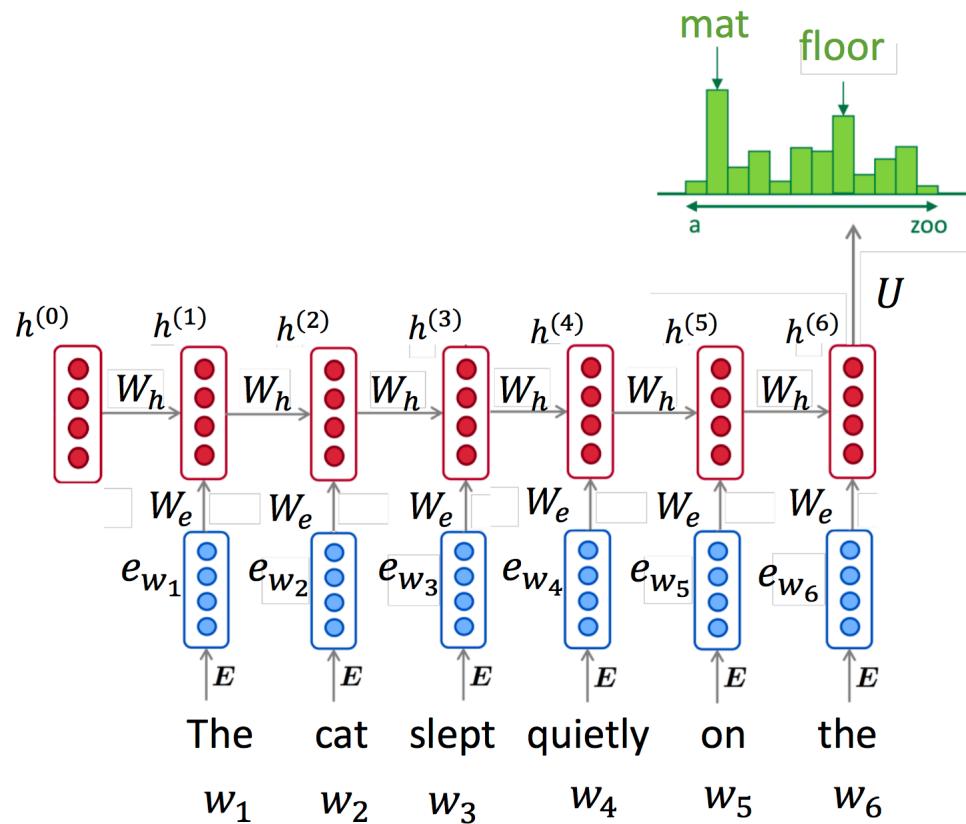
$$\hat{y}_t = \text{softmax}(Uh^{(t)} + b_2) \in \mathbb{R}^{|V|}$$

Hidden states

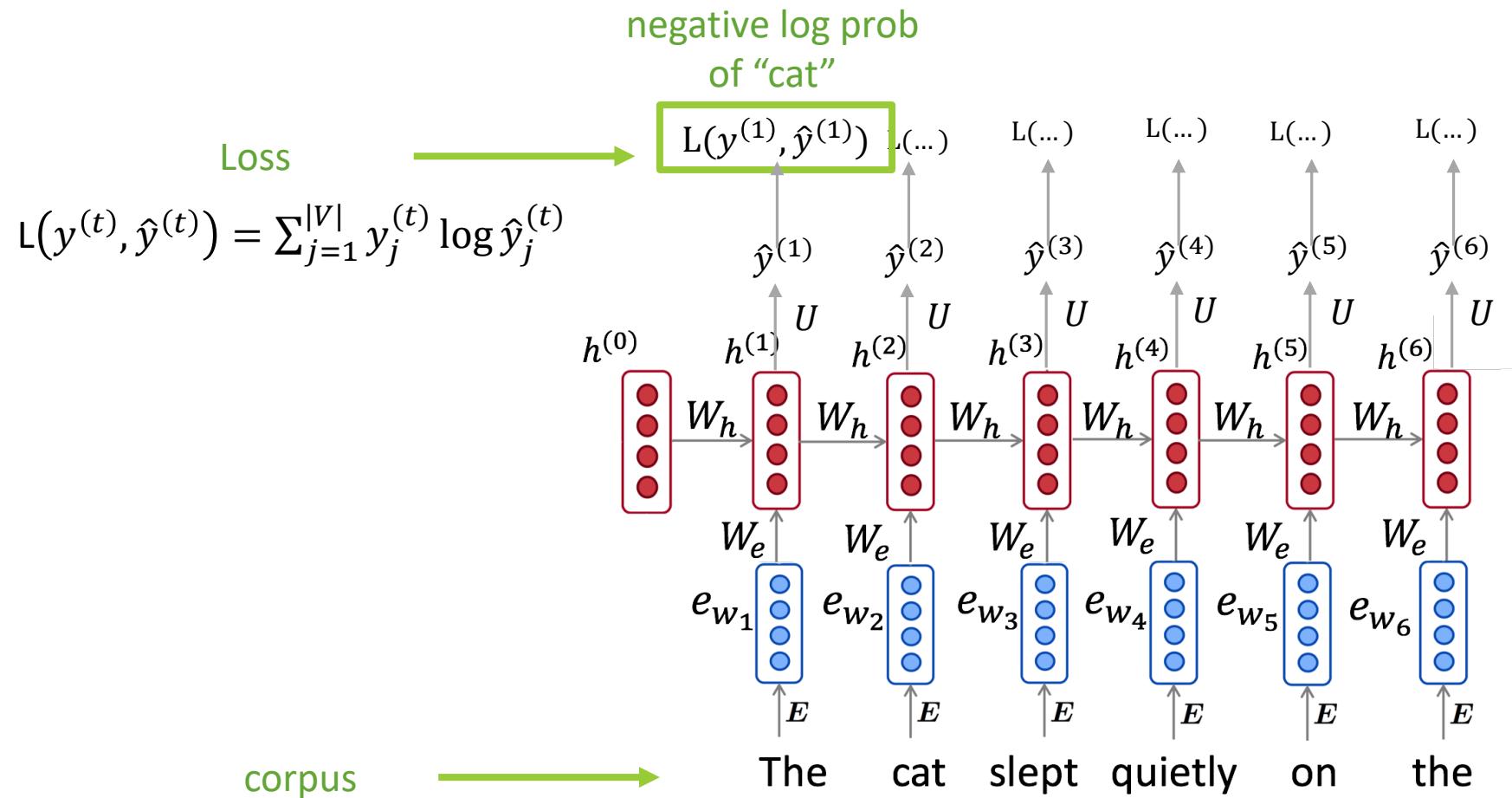
$$h^{(t)} = \sigma(W_h h^{(t-1)} + W_e e_{w_{t-1}} + b_1)$$

Word embeddings

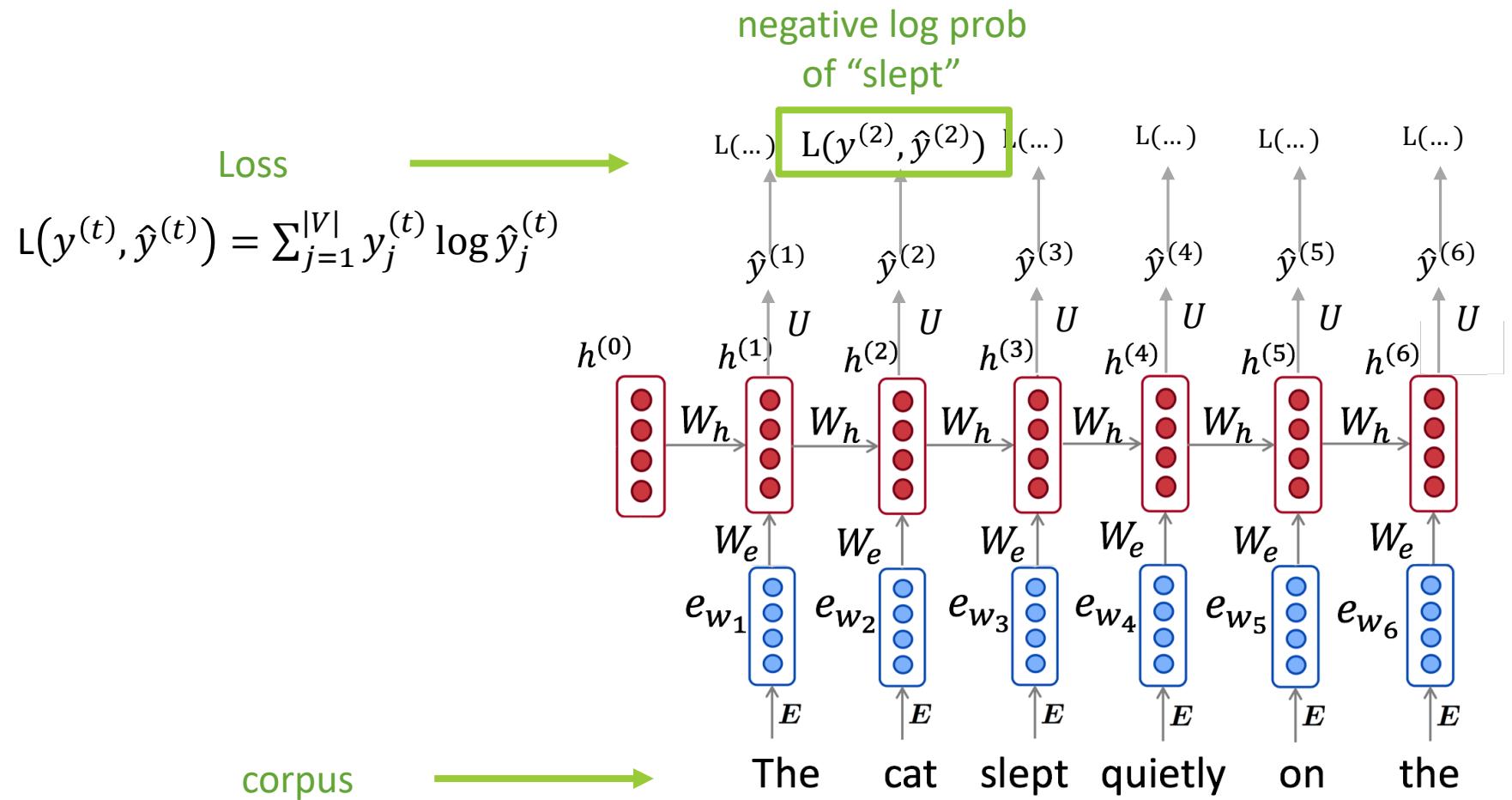
Words/tokens



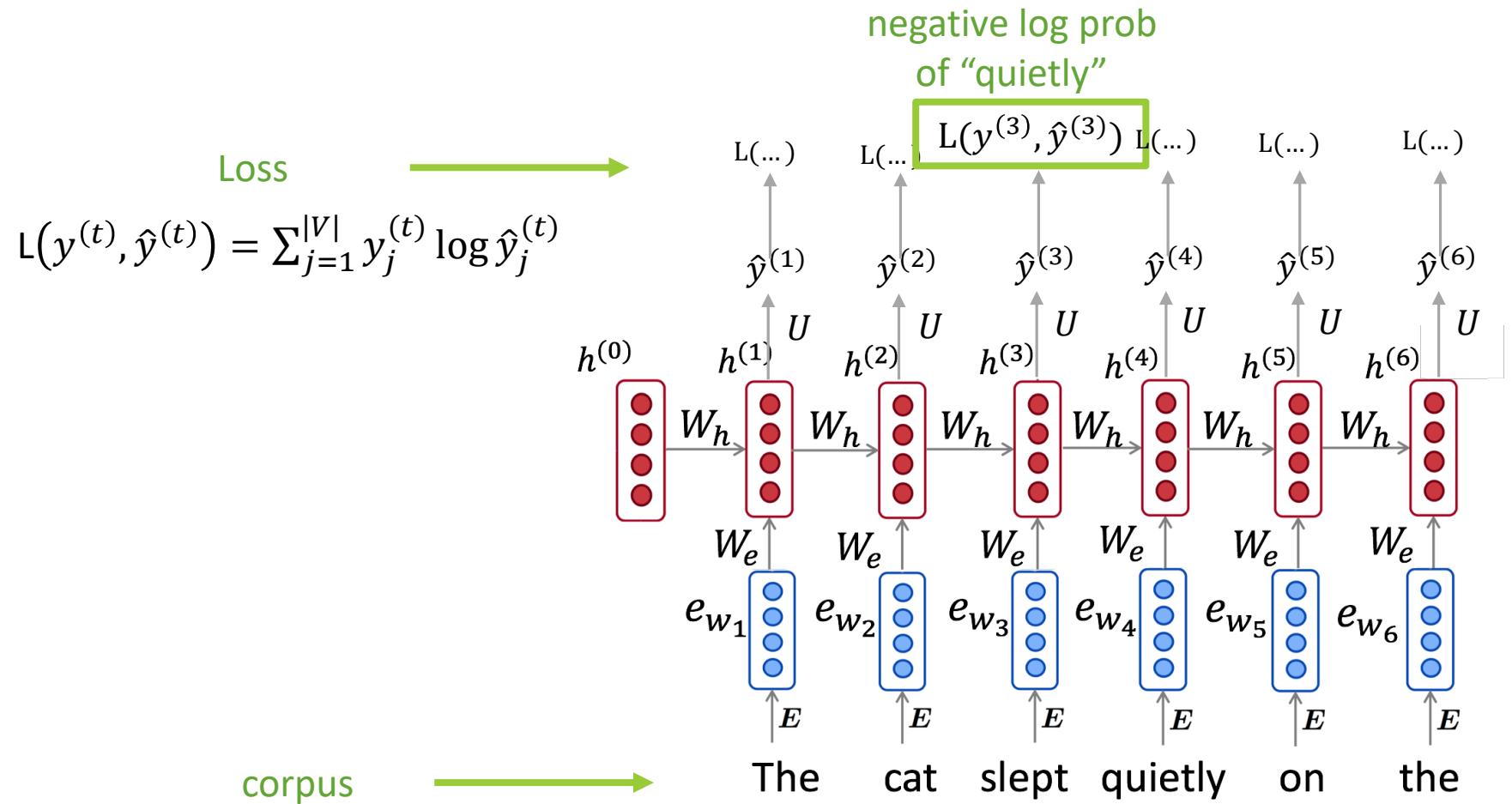
# Training RNN-LM



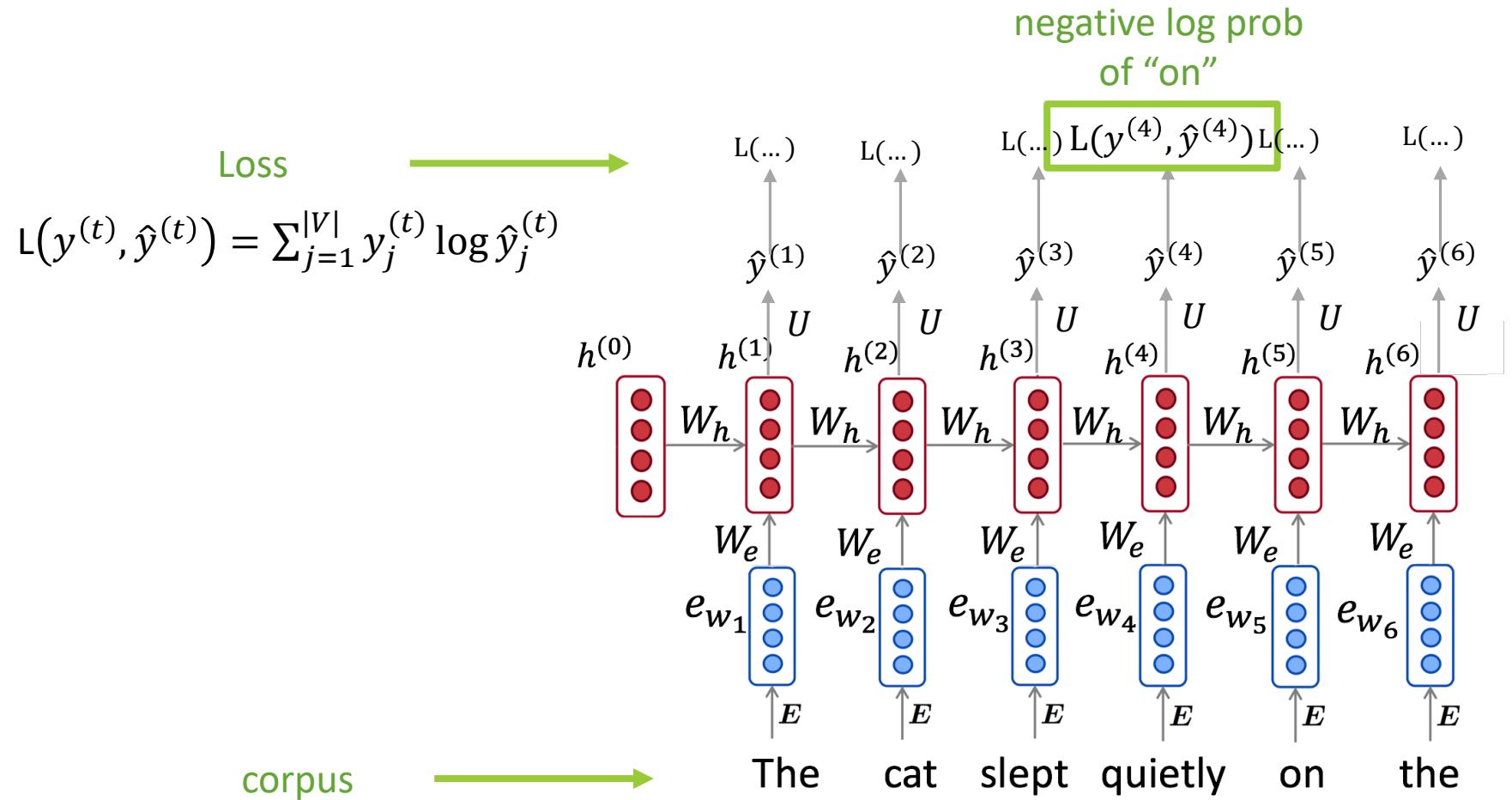
# Training RNN-LM



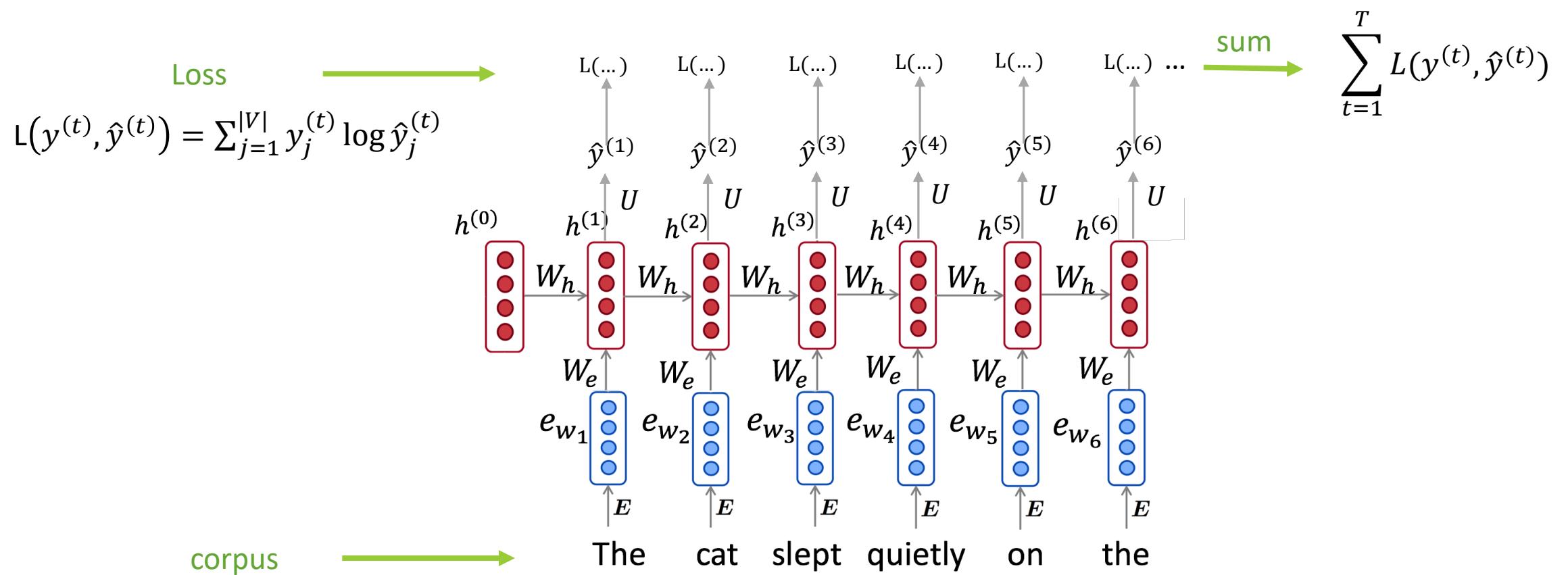
# Training RNN-LM



# Training RNN-LM



# Training RNN-LM



# Visualizing and Understanding Recurrent Networks

Char-based LM trained  
on Latex of book on  
algebraic geometry

*Proof.* Omitted.  $\square$

**Lemma 0.1.** Let  $\mathcal{C}$  be a set of the construction.

Let  $\mathcal{C}$  be a gerber covering. Let  $\mathcal{F}$  be a quasi-coherent sheaves of  $\mathcal{O}$ -modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

*Proof.* This is an algebraic space with the composition of sheaves  $\mathcal{F}$  on  $X_{\text{étale}}$  we have

$$\mathcal{O}_X(\mathcal{F}) = \{\text{morph}_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where  $\mathcal{G}$  defines an isomorphism  $\mathcal{F} \rightarrow \mathcal{F}$  of  $\mathcal{O}$ -modules.  $\square$

**Lemma 0.2.** This is an integer  $\mathcal{Z}$  is injective.

*Proof.* See Spaces, Lemma ???.  $\square$

**Lemma 0.3.** Let  $S$  be a scheme. Let  $X$  be a scheme and  $X$  is an affine open covering. Let  $\mathcal{U} \subset X$  be a canonical and locally of finite type. Let  $X$  be a scheme. Let  $X$  be a scheme which is equal to the formal complex.

The following to the construction of the lemma follows.

Let  $X$  be a scheme. Let  $X$  be a scheme covering. Let

$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

be a morphism of algebraic spaces over  $S$  and  $Y$ .

*Proof.* Let  $X$  be a nonzero scheme of  $X$ . Let  $X$  be an algebraic space. Let  $\mathcal{F}$  be a quasi-coherent sheaf of  $\mathcal{O}_X$ -modules. The following are equivalent

- (1)  $\mathcal{F}$  is an algebraic space over  $S$ .
- (2) If  $X$  is an affine open covering.

Consider a common structure on  $X$  and  $X$  the functor  $\mathcal{O}_X(U)$  which is locally of finite type.  $\square$

This since  $\mathcal{F} \in \mathcal{F}$  and  $x \in \mathcal{G}$  the diagram

$$\begin{array}{ccccc}
 S & \xrightarrow{\quad} & & & \\
 \downarrow & & & & \\
 \xi & \longrightarrow & \mathcal{O}_{X'} & \nearrow & \\
 \text{gor}_s & & \uparrow & & \\
 & & = \alpha' & \longrightarrow & \\
 & \uparrow & = \alpha' & \longrightarrow & \alpha \\
 \text{Spec}(K_\psi) & & \text{Mor}_{\text{Sets}} & & X \\
 & & & & \downarrow \\
 & & & & d(\mathcal{O}_{X_{X/k}}, \mathcal{G})
 \end{array}$$

is a limit. Then  $\mathcal{G}$  is a finite type and assume  $S$  is a flat and  $\mathcal{F}$  and  $\mathcal{G}$  is a finite type  $f_*$ . This is of finite type diagrams, and

- the composition of  $\mathcal{G}$  is a regular sequence,
- $\mathcal{O}_{X'}$  is a sheaf of rings.

$\square$

*Proof.* We have see that  $X = \text{Spec}(R)$  and  $\mathcal{F}$  is a finite type representable by algebraic space. The property  $\mathcal{F}$  is a finite morphism of algebraic stacks. Then the cohomology of  $X$  is an open neighbourhood of  $U$ .  $\square$

*Proof.* This is clear that  $\mathcal{G}$  is a finite presentation, see Lemmas ??.

A reduced above we conclude that  $U$  is an open covering of  $\mathcal{C}$ . The functor  $\mathcal{F}$  is a “field”

$$\mathcal{O}_{X,x} \rightarrow \mathcal{F}_{\mathbb{F}} \dashv (\mathcal{O}_{X_{\text{étale}}} \rightarrow \mathcal{O}_{X_\ell}^{-1} \mathcal{O}_{X_\lambda}(\mathcal{O}_{X_\eta}^\pi))$$

is an isomorphism of covering of  $\mathcal{O}_{X_\ell}$ . If  $\mathcal{F}$  is the unique element of  $\mathcal{F}$  such that  $X$  is an isomorphism.

The property  $\mathcal{F}$  is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme  $\mathcal{O}_X$ -algebra with  $\mathcal{F}$  are opens of finite type over  $S$ . If  $\mathcal{F}$  is a scheme theoretic image points.  $\square$

If  $\mathcal{F}$  is a finite direct sum  $\mathcal{O}_{X_\lambda}$  is a closed immersion, see Lemma ???. This is a sequence of  $\mathcal{F}$  is a similar morphism.

More hallucinated algebraic geometry. Nice try on the diagram (right).

Karpathy et al, ICLR workshop, 2016 <https://arxiv.org/pdf/1506.02078.pdf>

# Visualizing and Understanding Recurrent Networks

---

Char-based LM trained on  
Linux Source Code

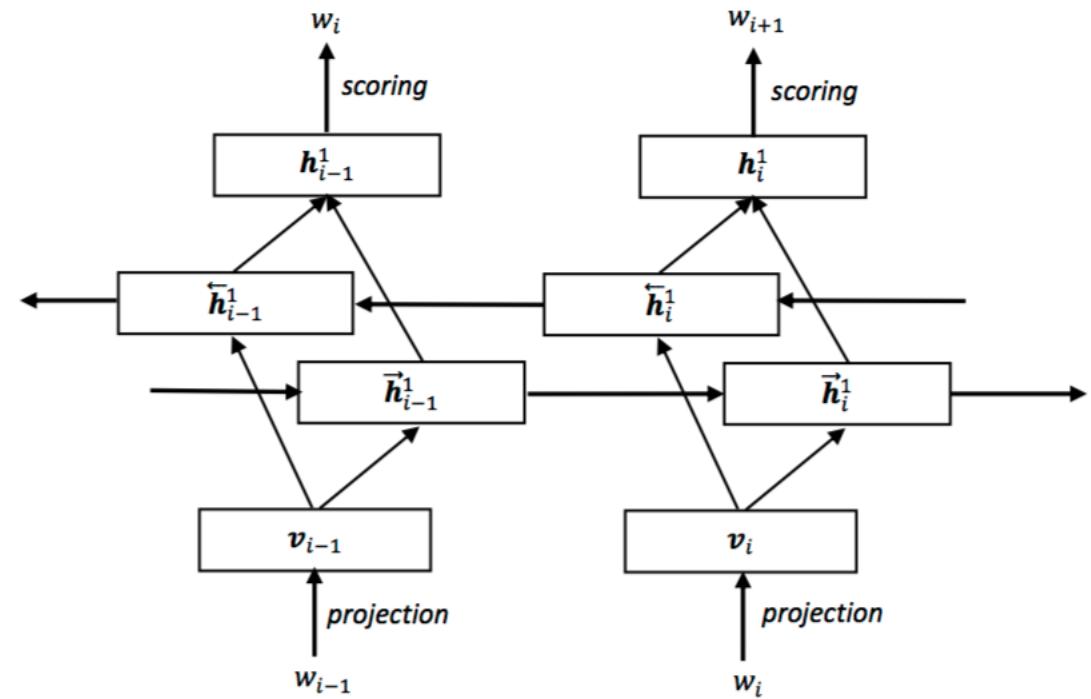
```
/*
 * Increment the size file of the new incorrect UI_FILTER group information
 * of the size generatively.
 */
static int indicate_policy(void)
{
    int error;
    if (fd == MARN_EPT) {
        /*
         * The kernel blank will coeld it to userspace.
         */
        if (ss->segment < mem_total)
            unblock_graph_and_set_blocked();
        else
            ret = 1;
        goto bail;
    }
    segaddr = in_SB(in.addr);
    selector = seg / 16;
    setup_works = true;
    for (i = 0; i < blocks; i++) {
        seq = buf[i++];
        bpf = bd->bd.next + i * search;
        if (fd) {
            current = blocked;
        }
    }
    rw->name = "Getjbbregs";
    bprm_self_clearl(&iv->version);
    regs->new = blocks[(BPF_STATS << info->historidac)] | PFMR_CLOBATHINC_SECONDS << 12;
    return segtable;
}
```

Karpathy et al, ICLR workshop, 2016 <https://arxiv.org/pdf/1506.02078.pdf>

# Bi-RNN Language Models

---

- Can take into account both left and right contexts
- Can not be used for generation
- Then, where can it be used?



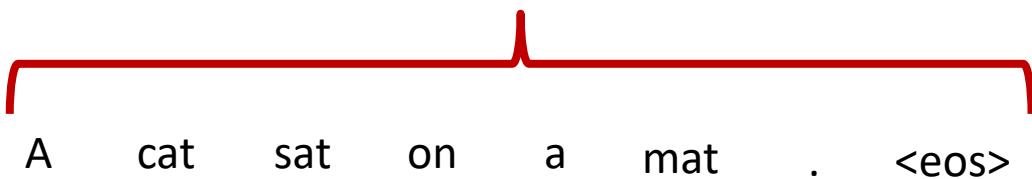
<https://arxiv.org/pdf/1602.06064.pdf>

# RNN Encoder-Decoder (Vanilla)

---

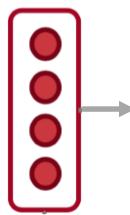


**Source sentence**



# RNN Encoder-Decoder (Vanilla)

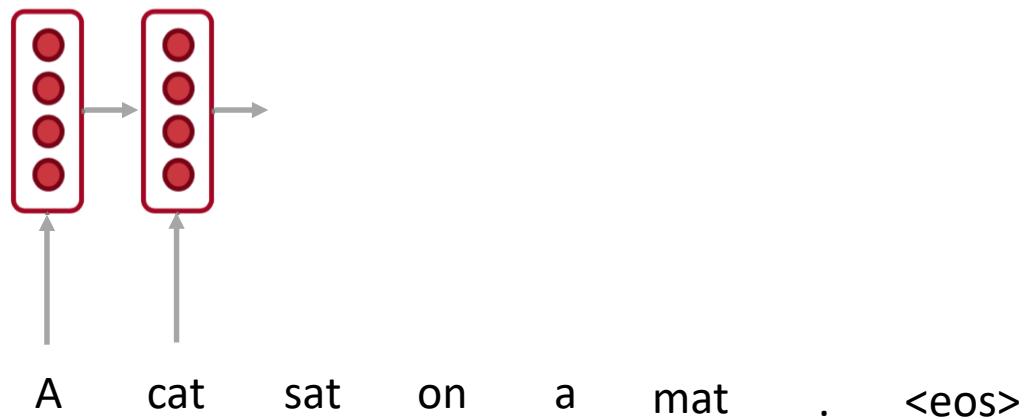
---



A    cat    sat    on    a    mat    .    <eos>

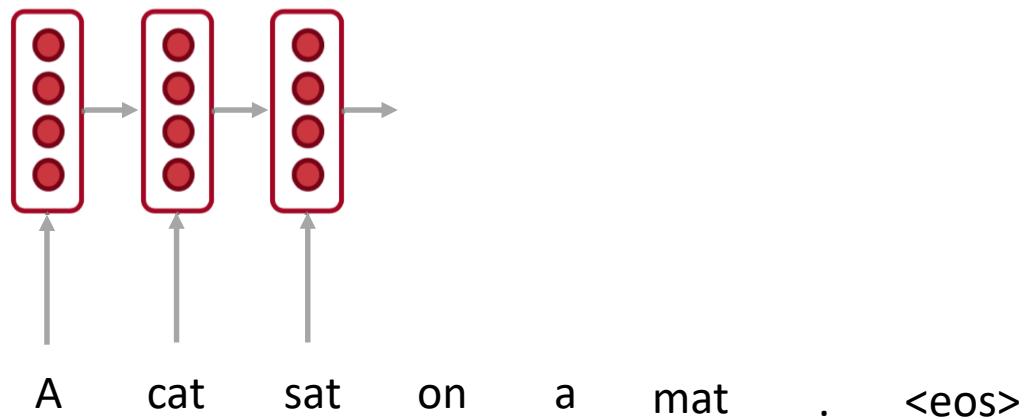
# RNN Encoder-Decoder (Vanilla)

---



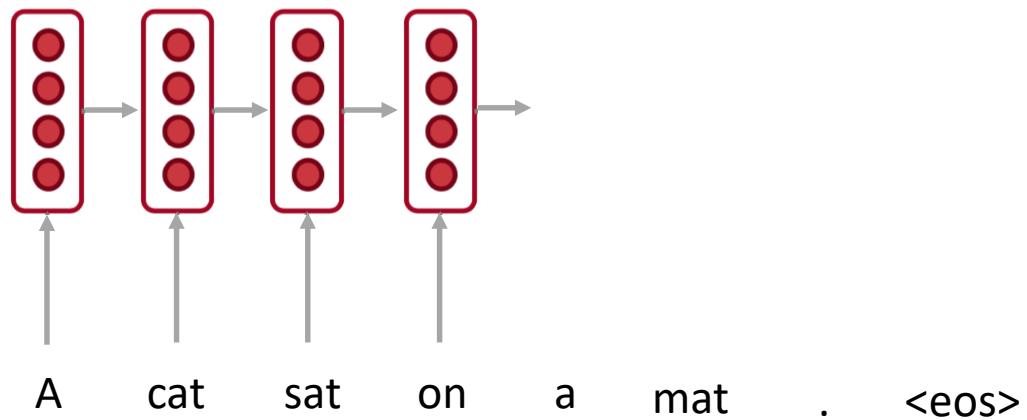
# RNN Encoder-Decoder (Vanilla)

---



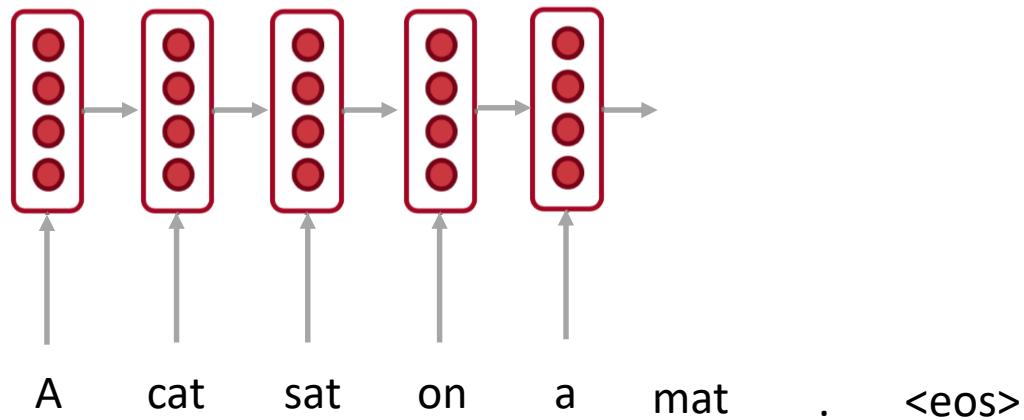
# RNN Encoder-Decoder (Vanilla)

---



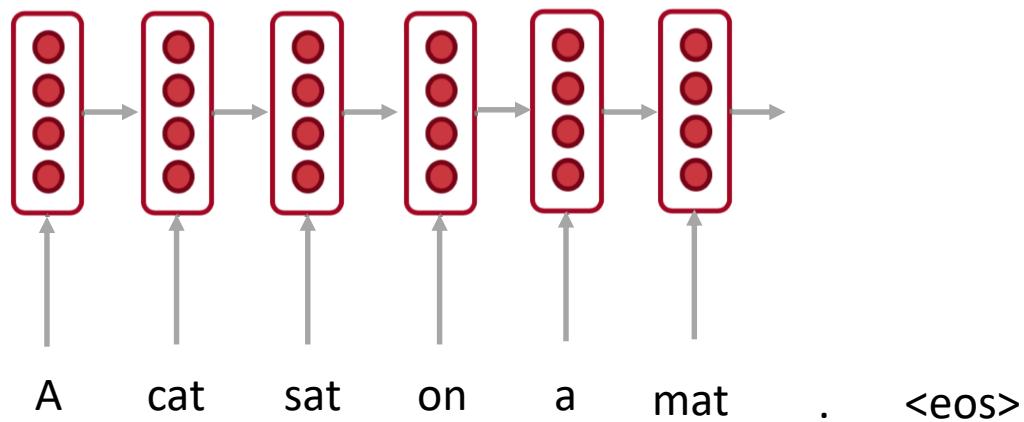
# RNN Encoder-Decoder (Vanilla)

---



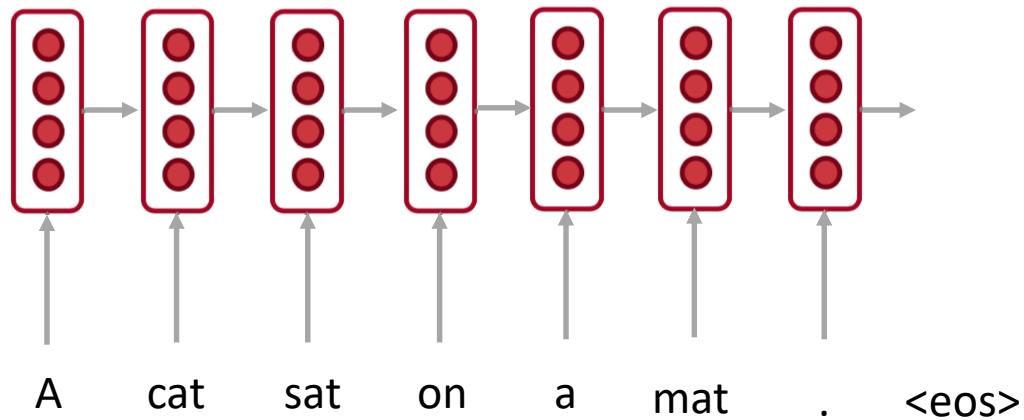
# RNN Encoder-Decoder (Vanilla)

---



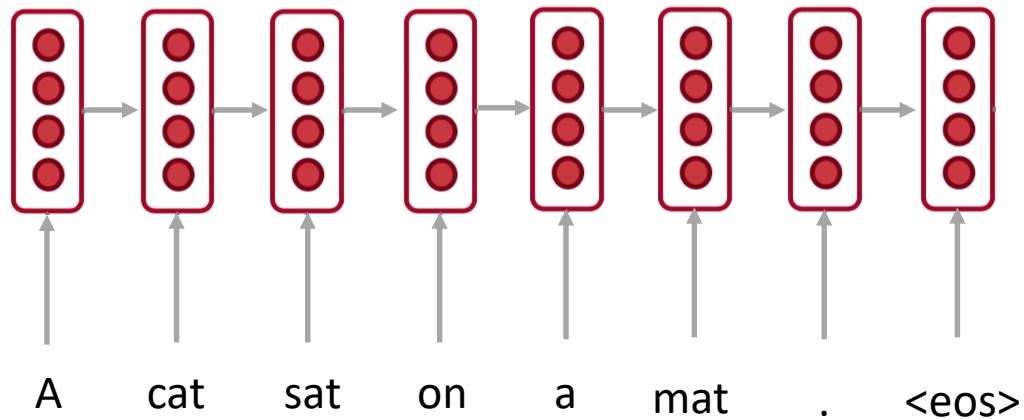
# RNN Encoder-Decoder (Vanilla)

---

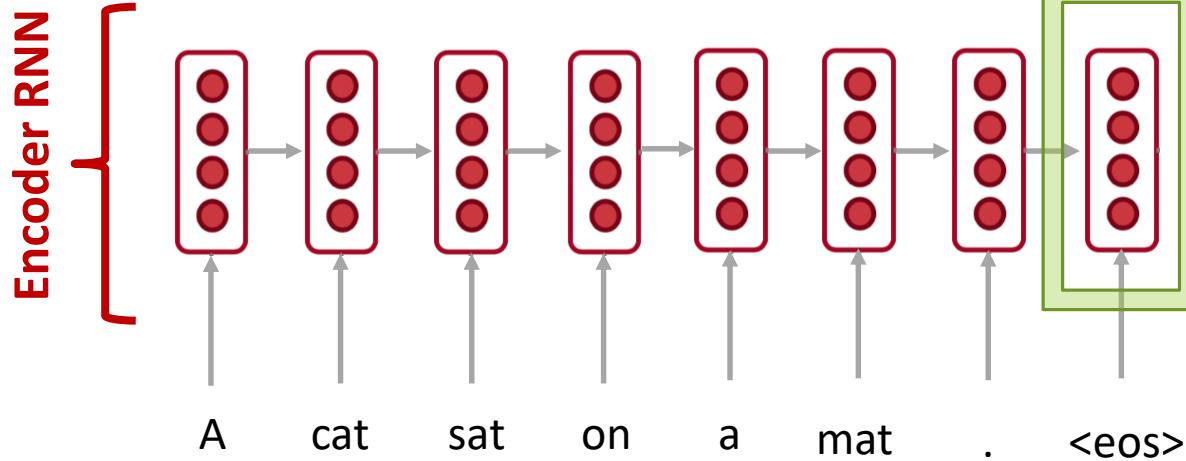


# RNN Encoder-Decoder (Vanilla)

---



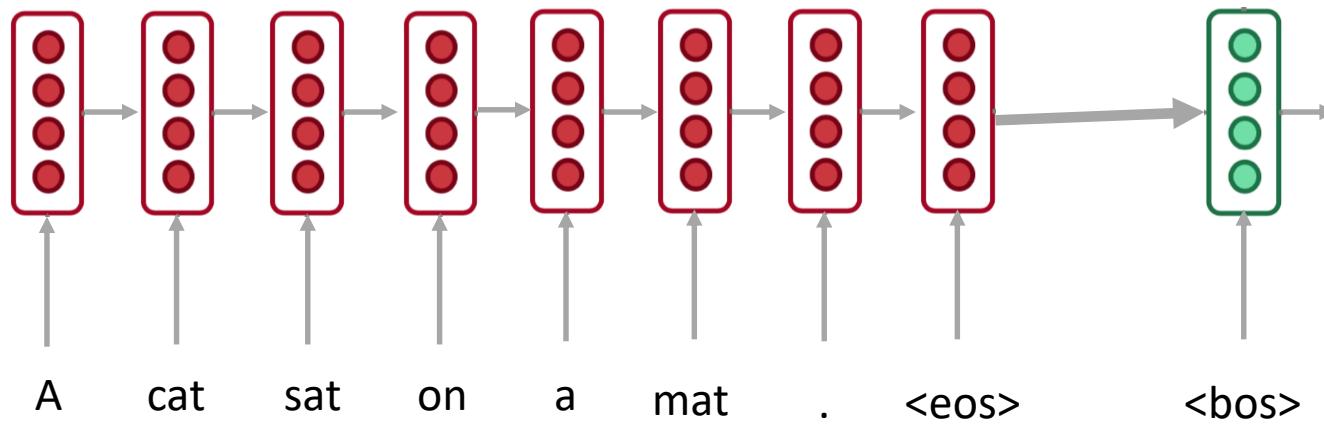
# RNN Encoder-Decoder (Vanilla)



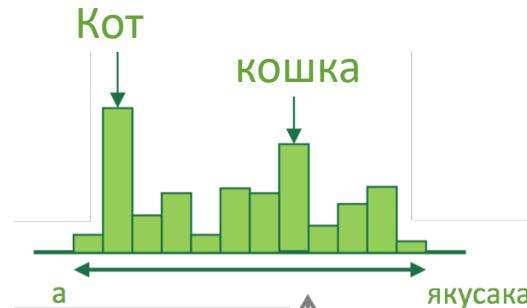
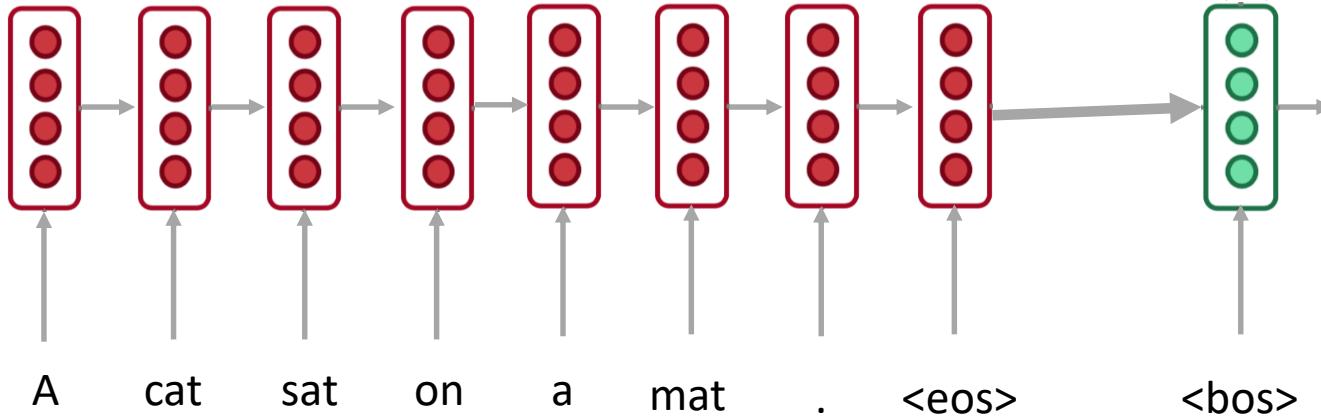
Encoding of a source sentence.  
Provides an initial state for decoder RNN

# RNN Encoder-Decoder (Vanilla)

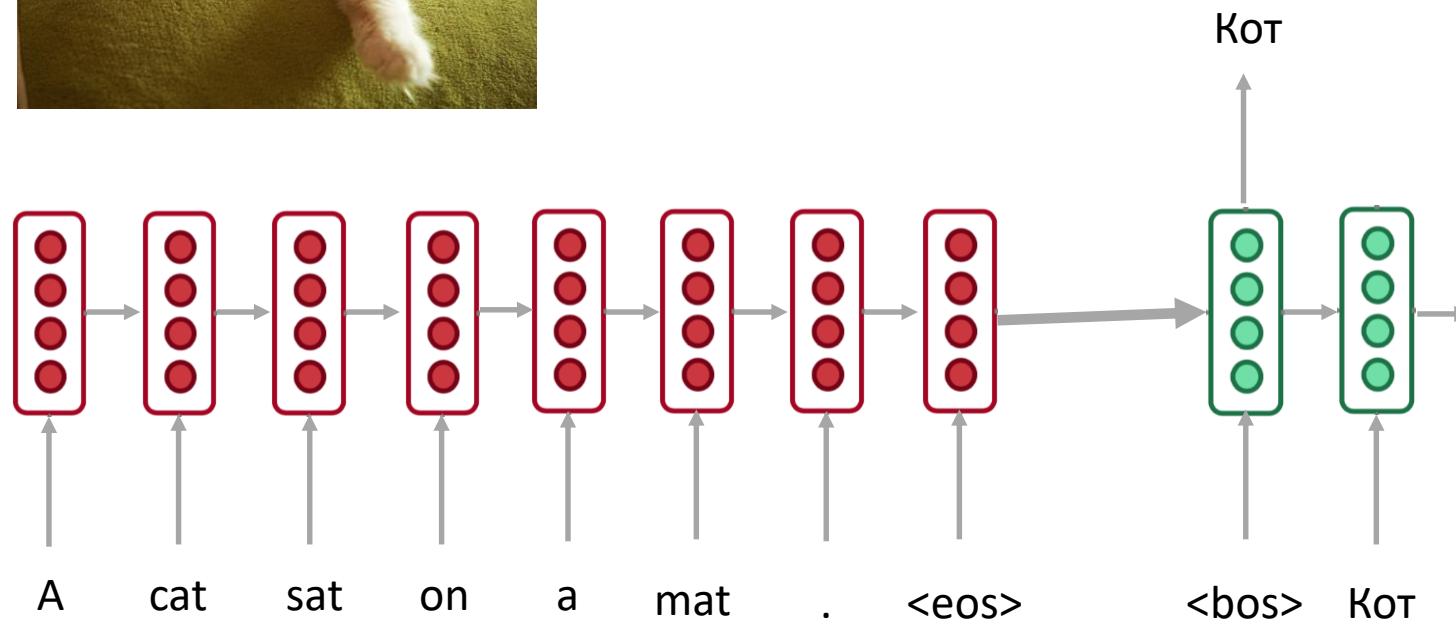
---



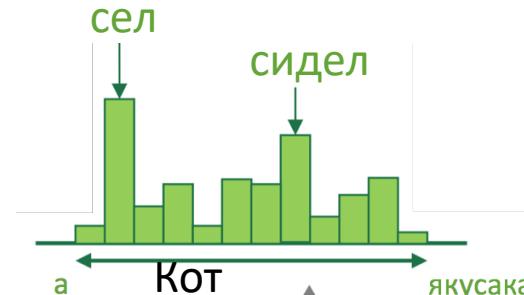
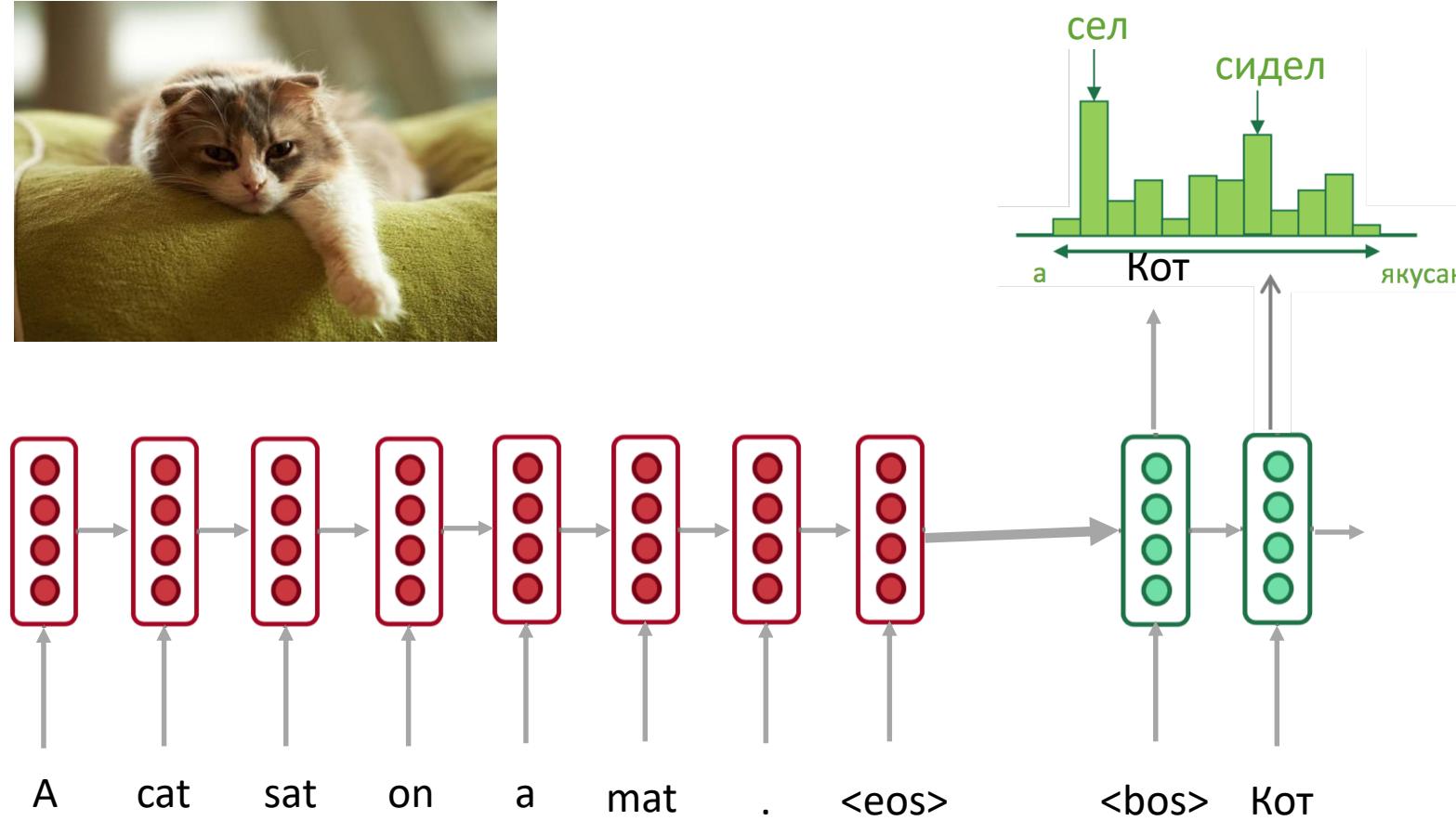
# RNN Encoder-Decoder (Vanilla)



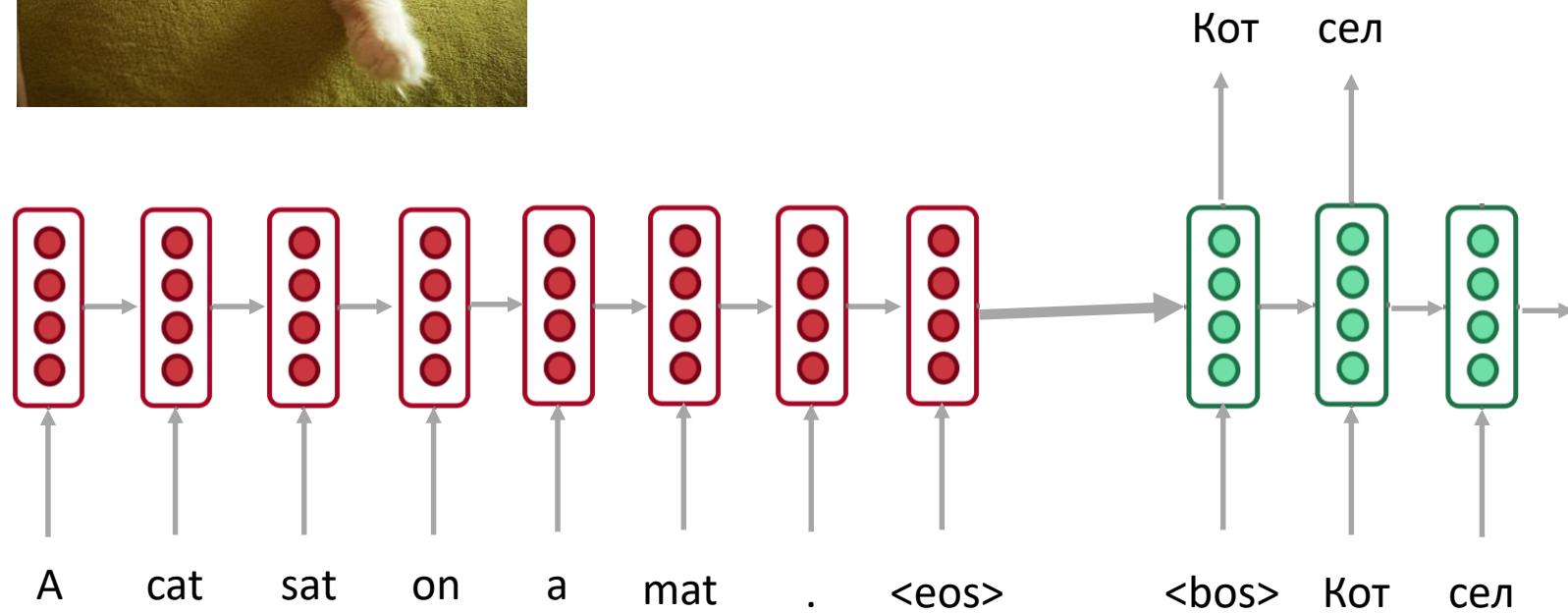
# RNN Encoder-Decoder (Vanilla)



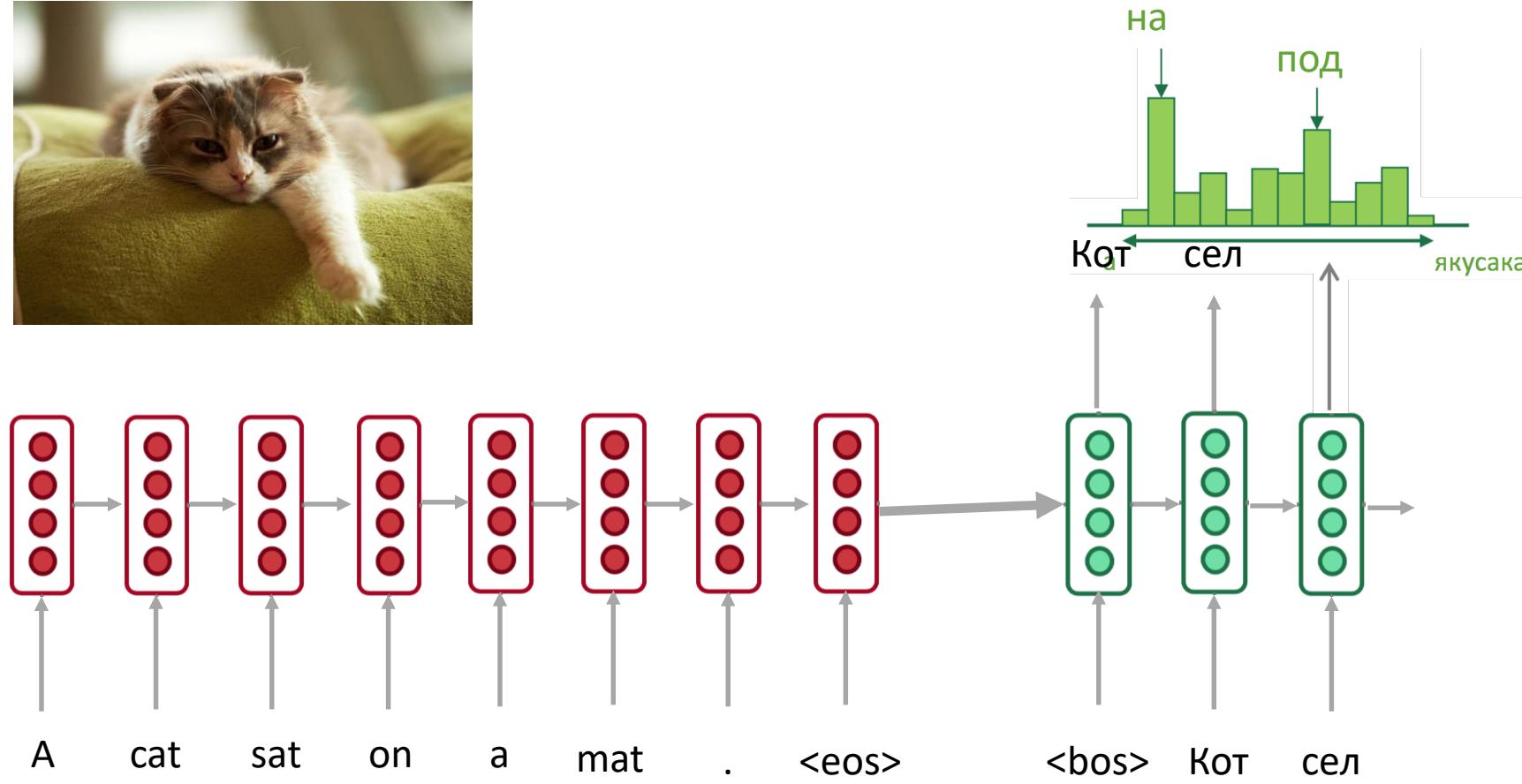
# RNN Encoder-Decoder (Vanilla)



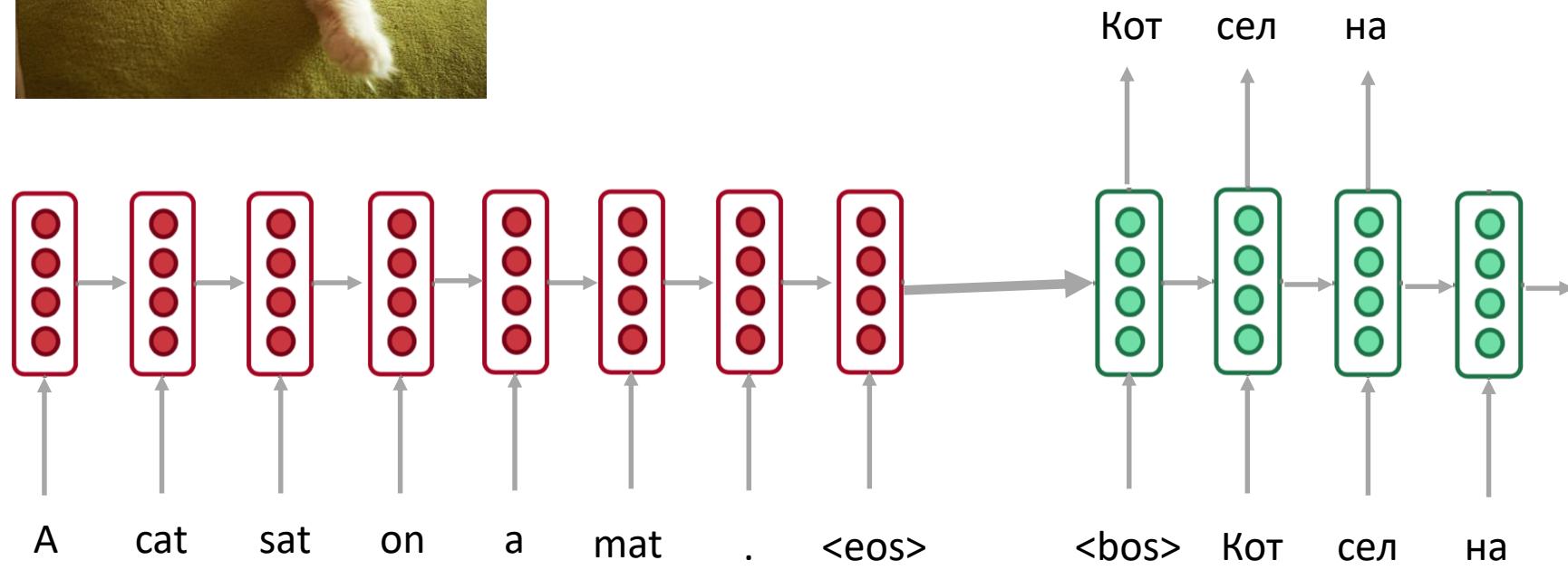
# RNN Encoder-Decoder (Vanilla)



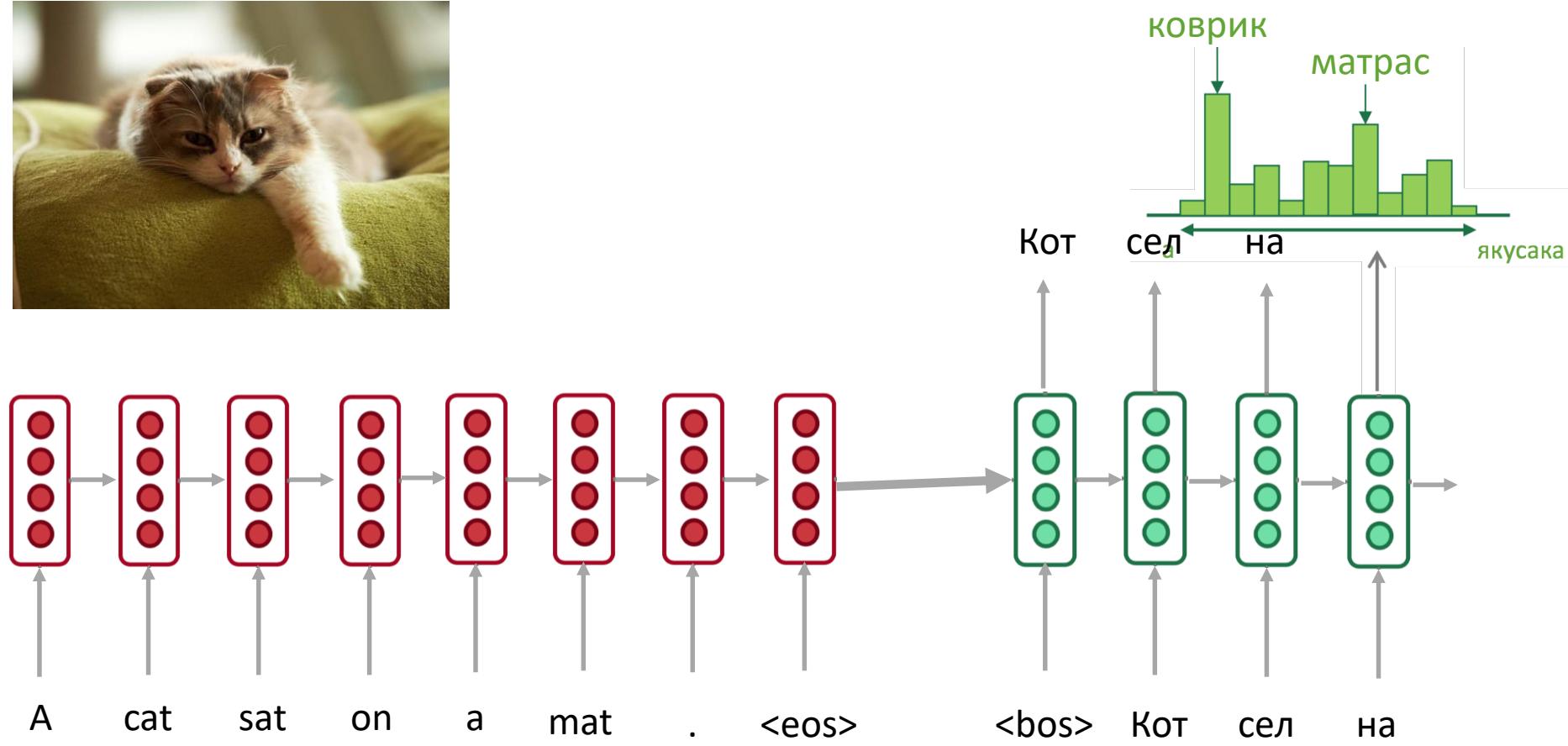
# RNN Encoder-Decoder (Vanilla)



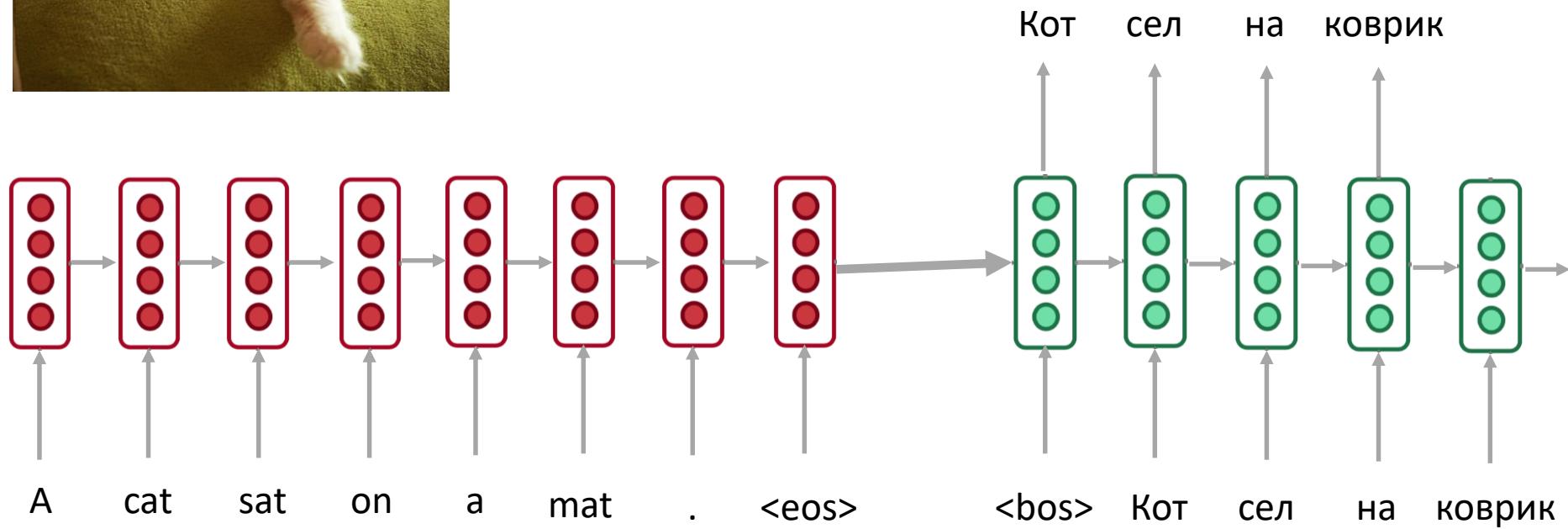
# RNN Encoder-Decoder (Vanilla)



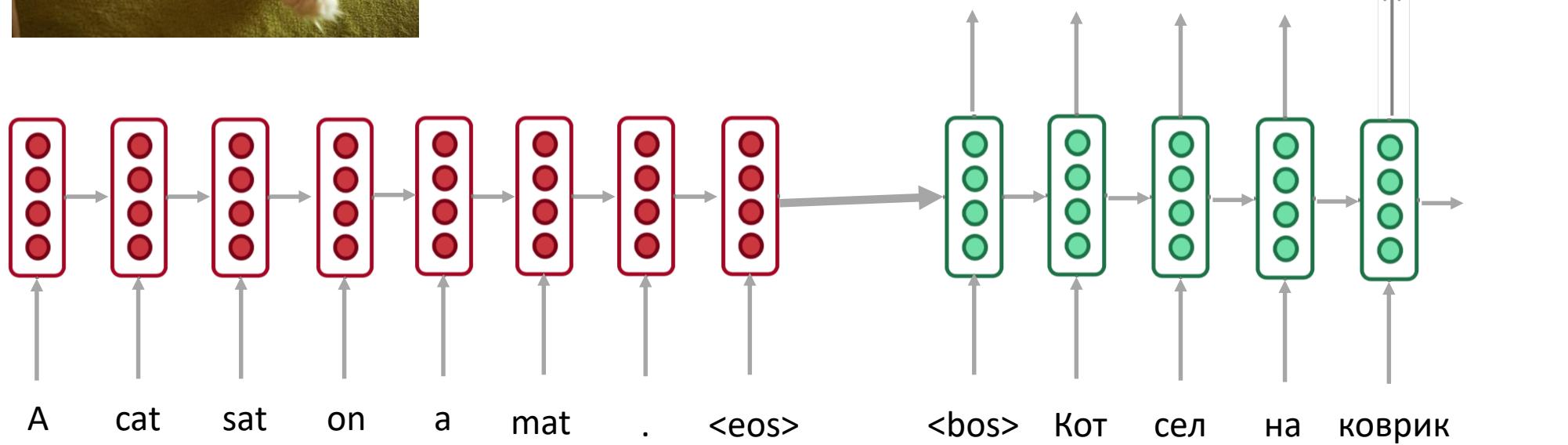
# RNN Encoder-Decoder (Vanilla)



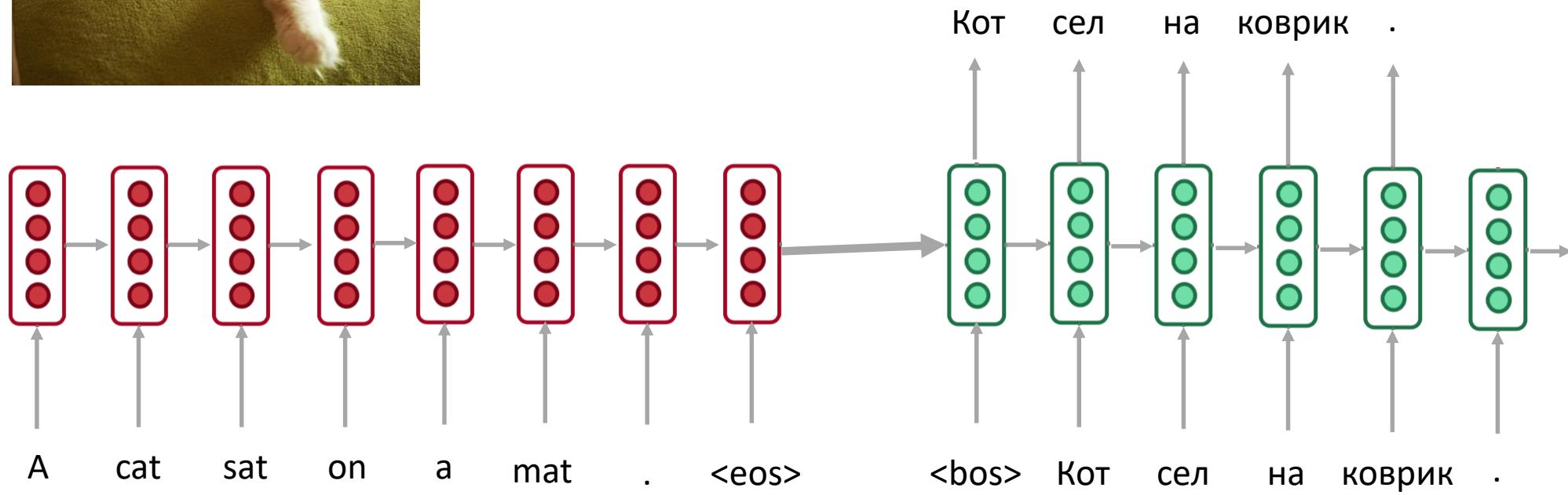
# RNN Encoder-Decoder (Vanilla)



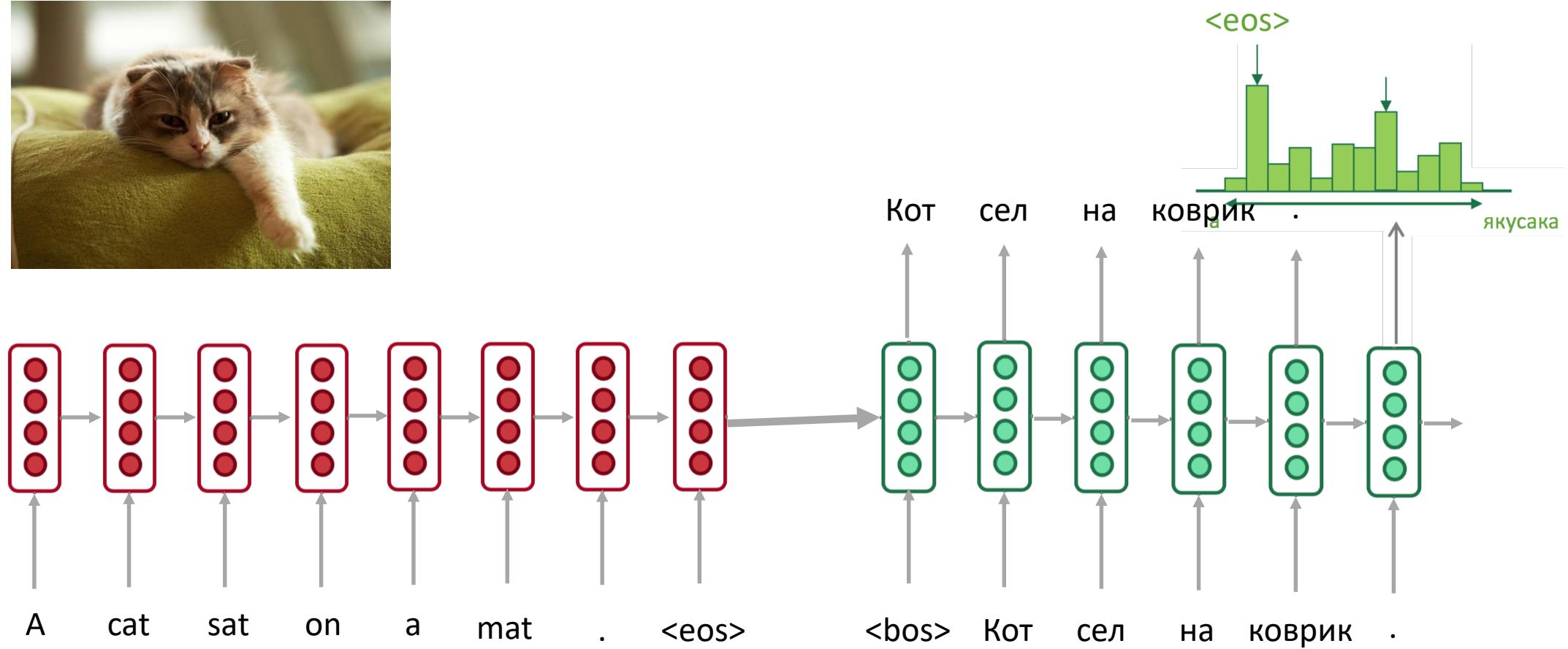
# RNN Encoder-Decoder (Vanilla)



# RNN Encoder-Decoder (Vanilla)



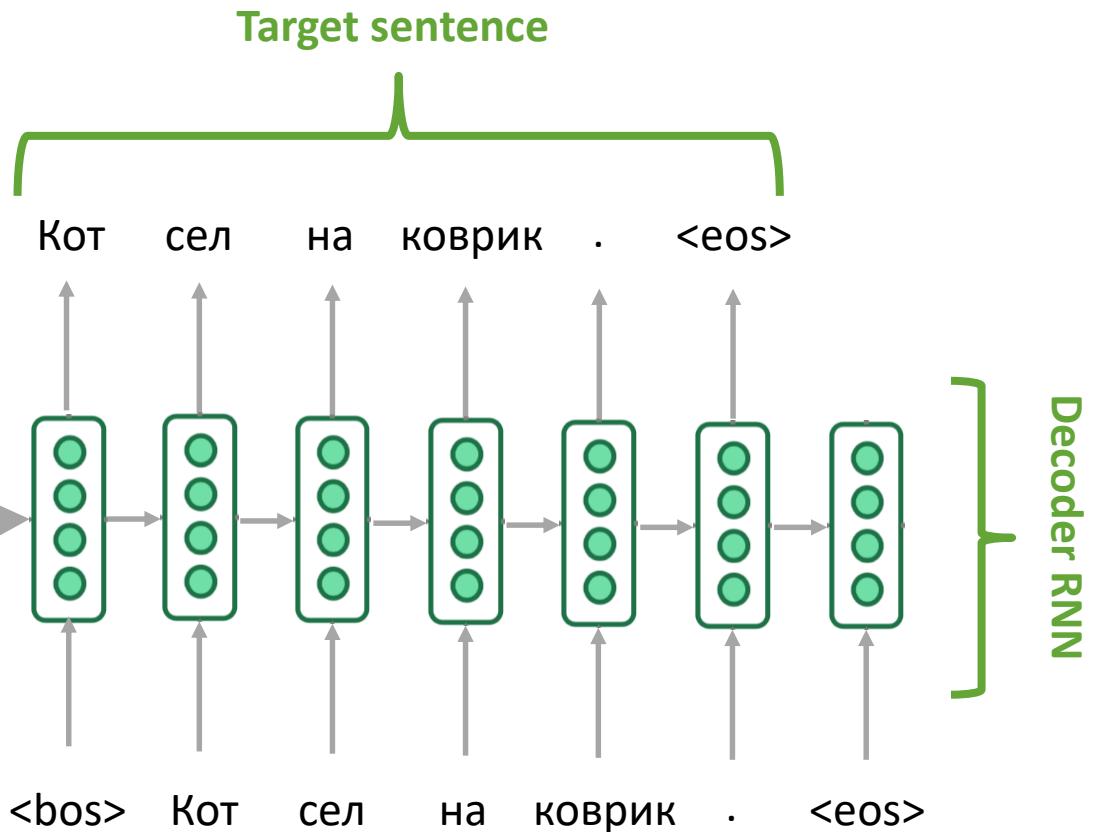
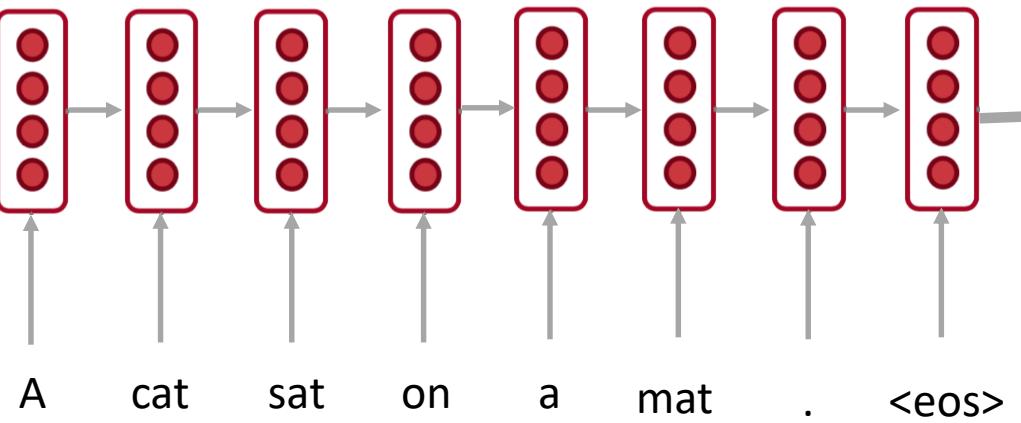
# RNN Encoder-Decoder (Vanilla)



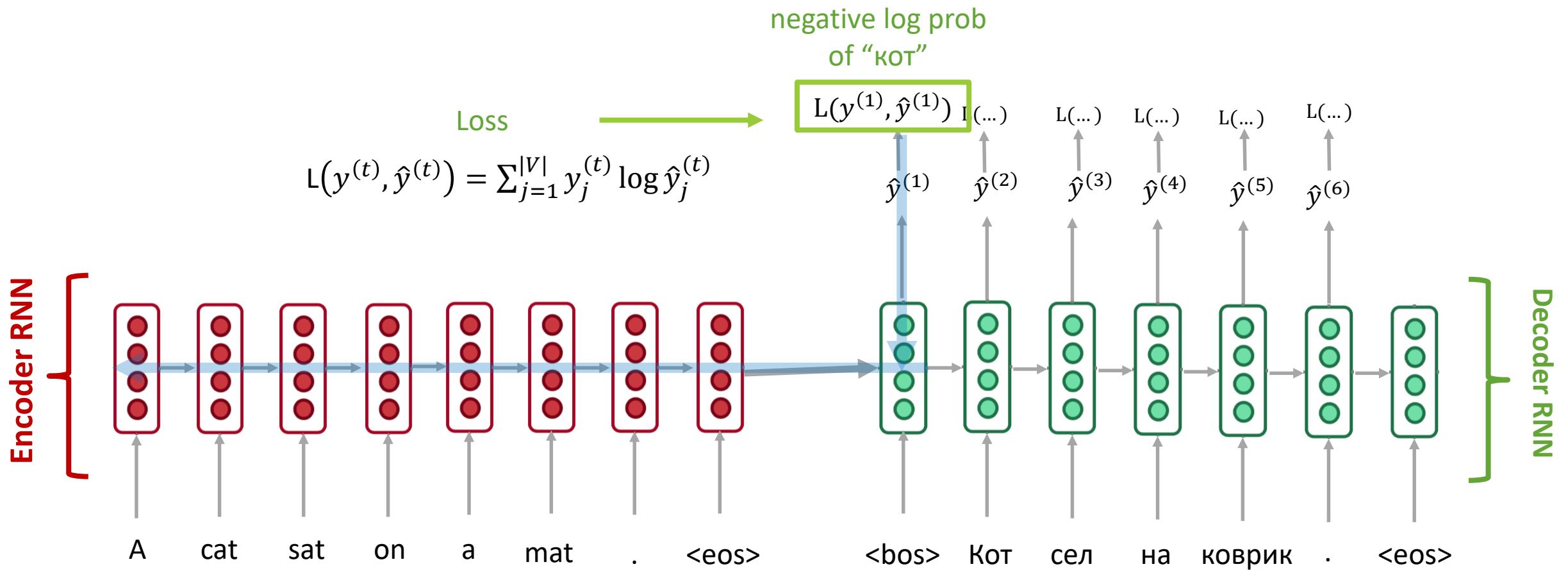
# RNN Encoder-Decoder (Vanilla)



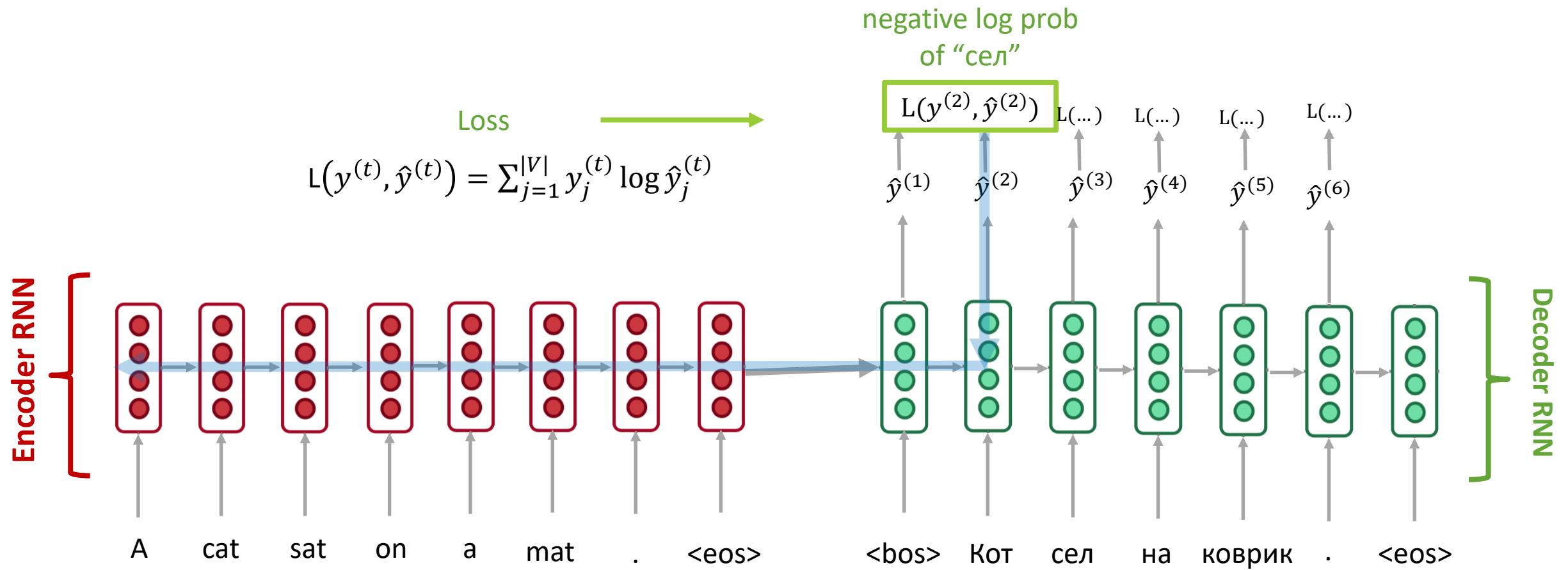
Encoder RNN



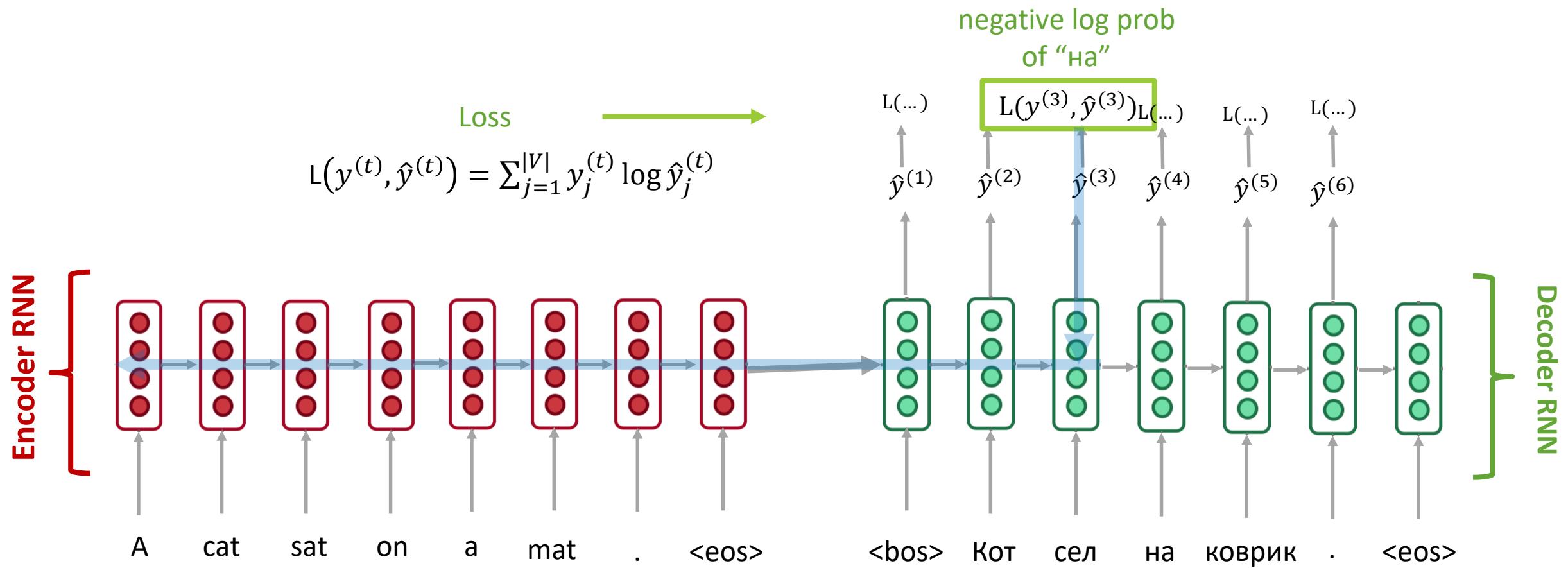
# RNN Encoder-Decoder (Vanilla)



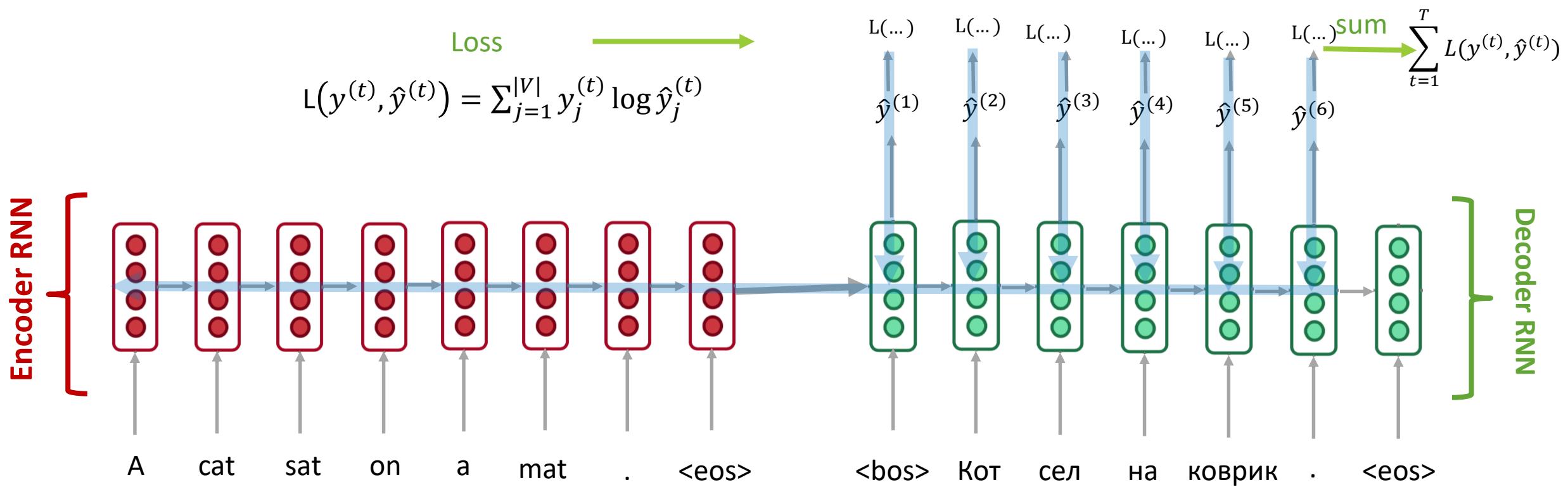
# RNN Encoder-Decoder (Vanilla)



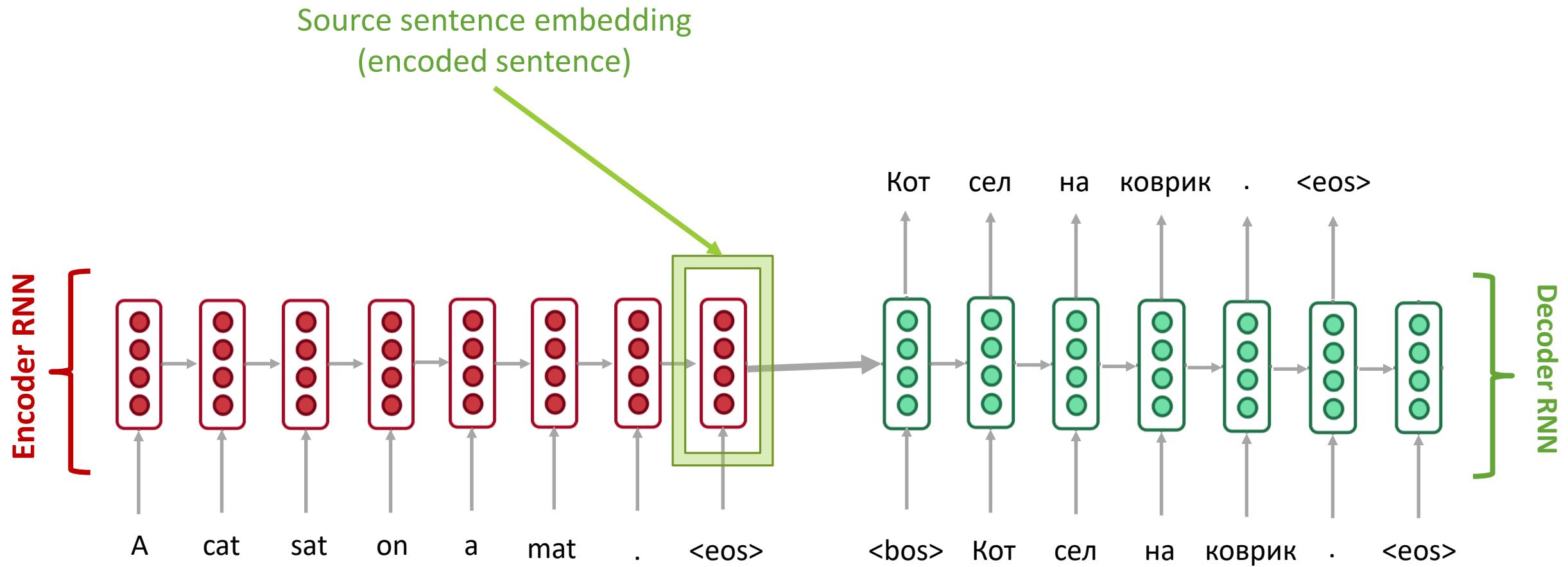
# RNN Encoder-Decoder (Vanilla)



# RNN Encoder-Decoder (Vanilla)

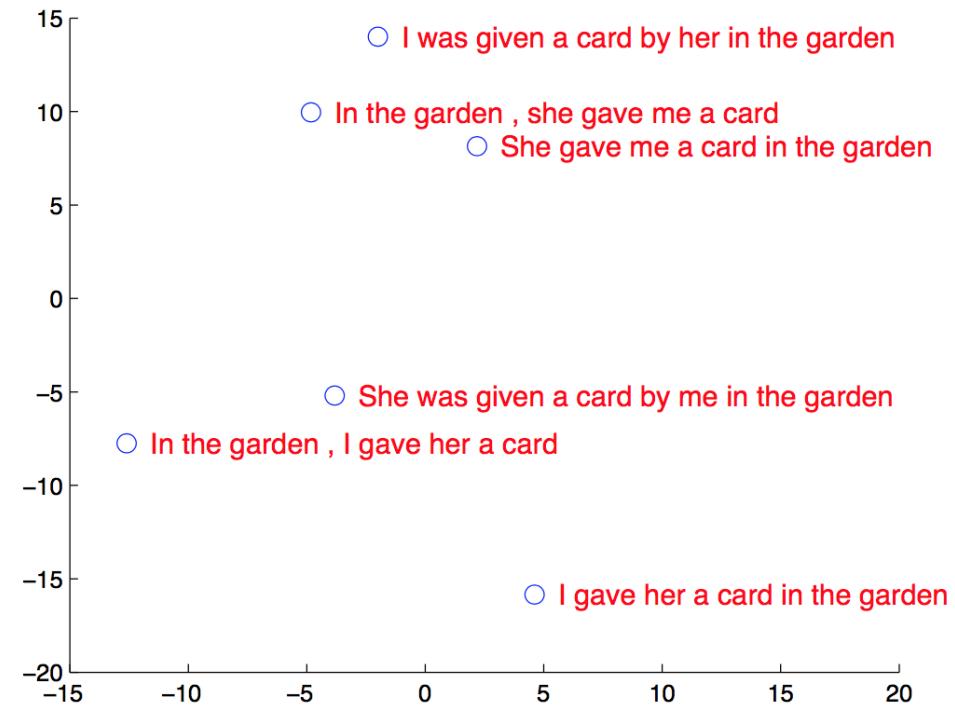
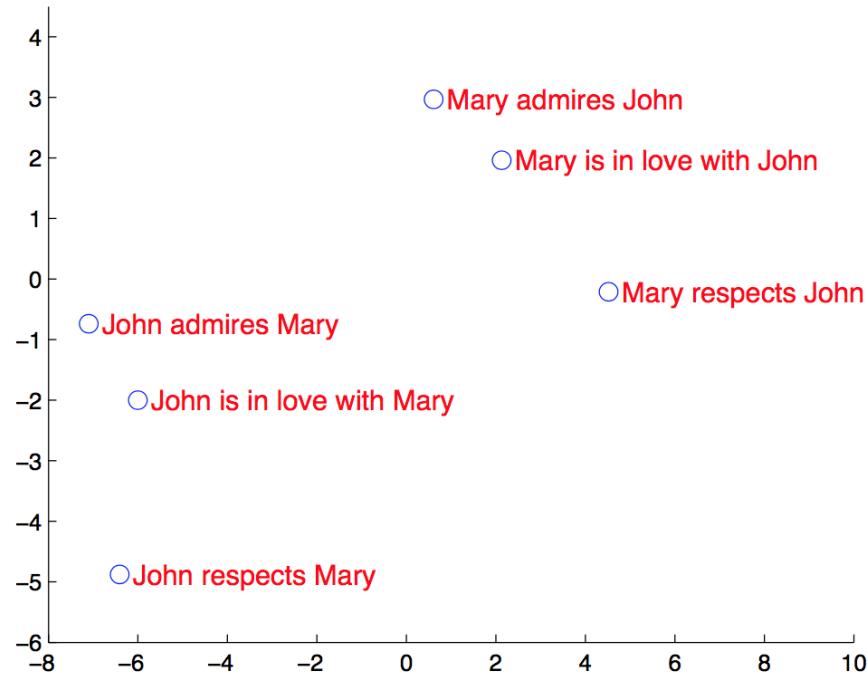


# Let's look at sentence embedding



# Let's look at sentence embedding

---

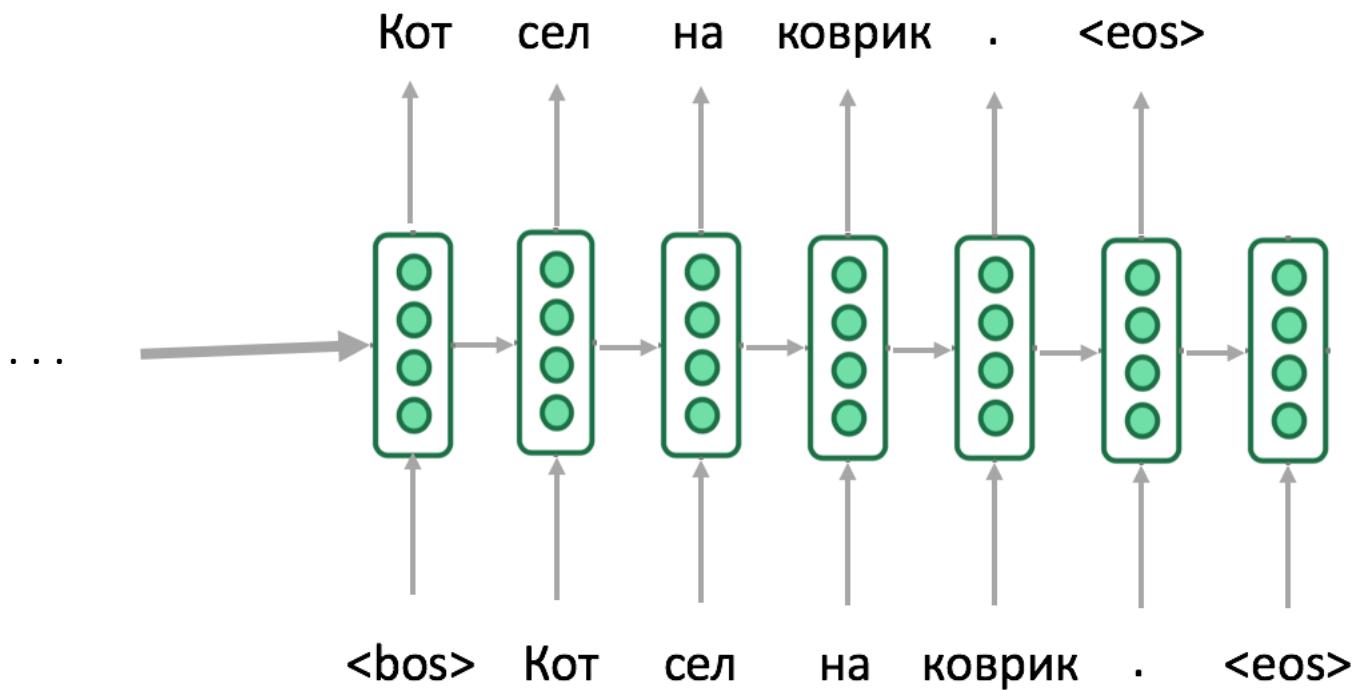


# Decoding

---

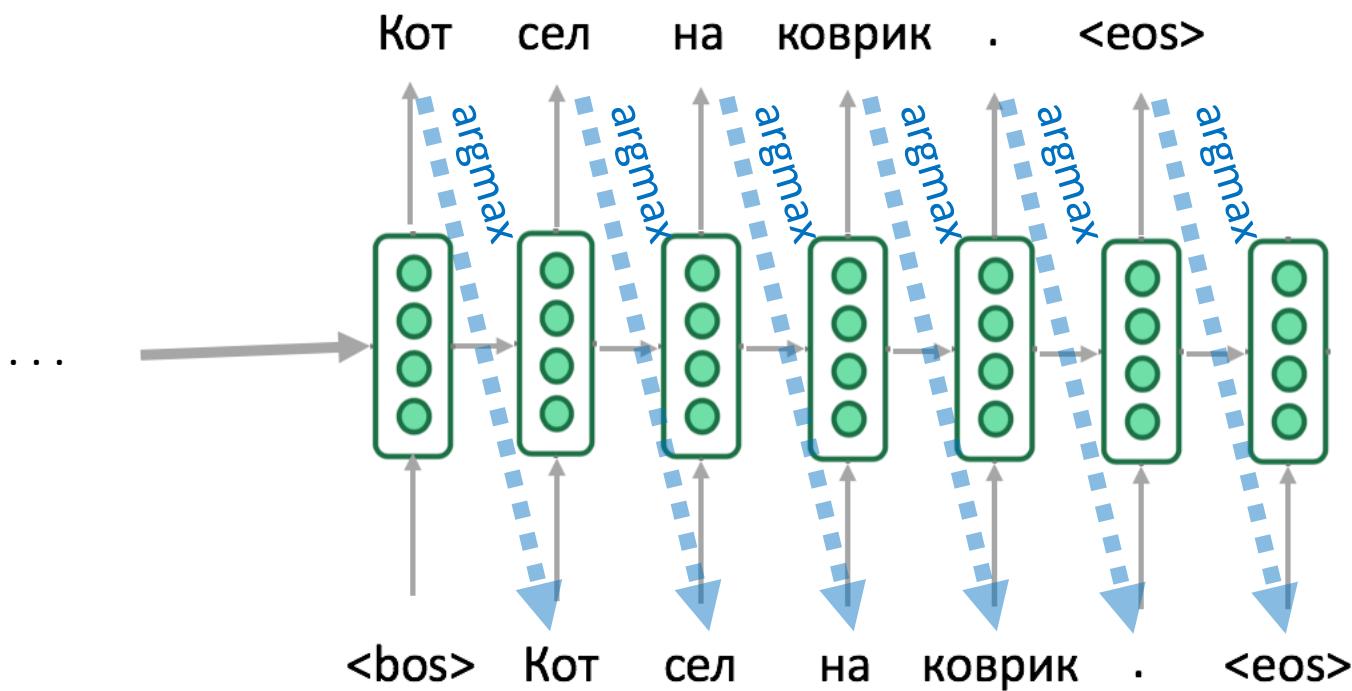
# Decoding: how to generate a sequence?

---



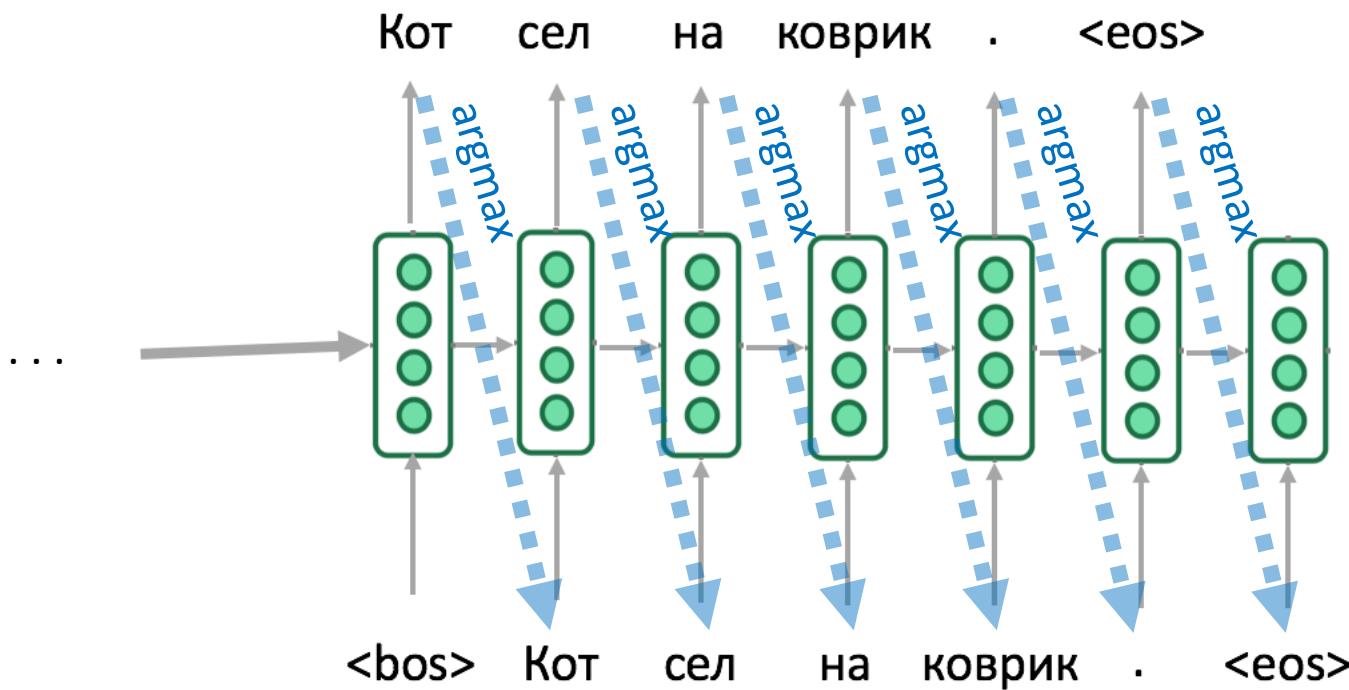
# Decoding: how to generate a sequence?

## 1. Greedy decoding

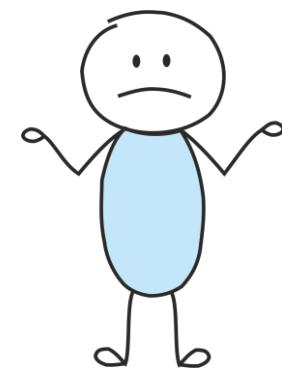


# Decoding: how to generate a sequence?

## 1. Greedy decoding



Is it really good?



# Decoding: how to generate a sequence?

---

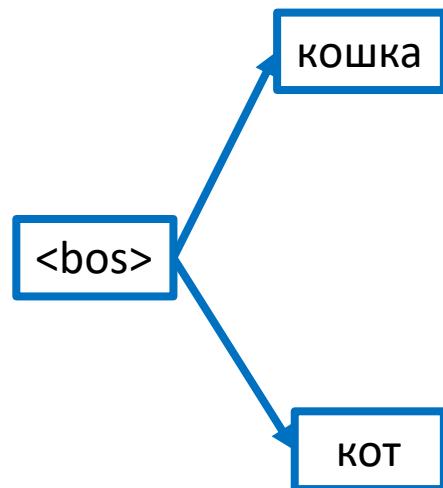
## 2. Beam search

<bos>

# Decoding: how to generate a sequence?

---

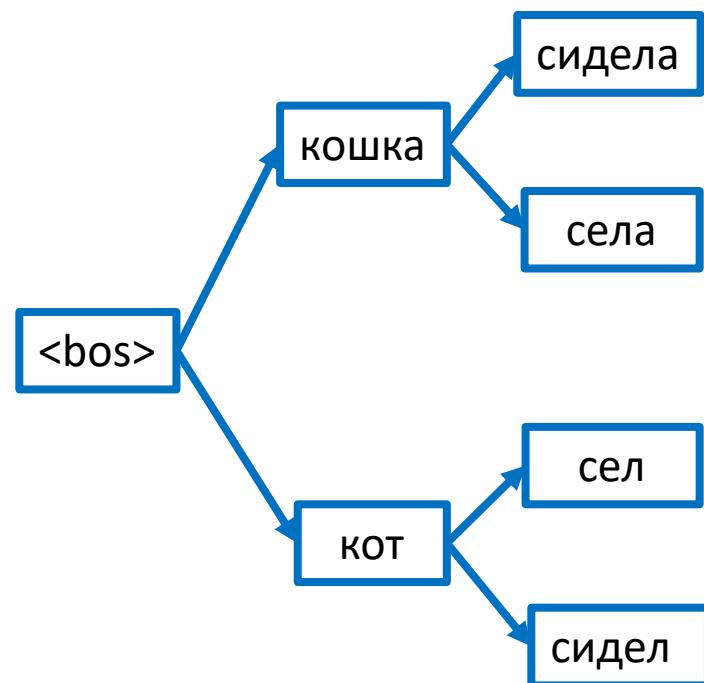
## 2. Beam search



# Decoding: how to generate a sequence?

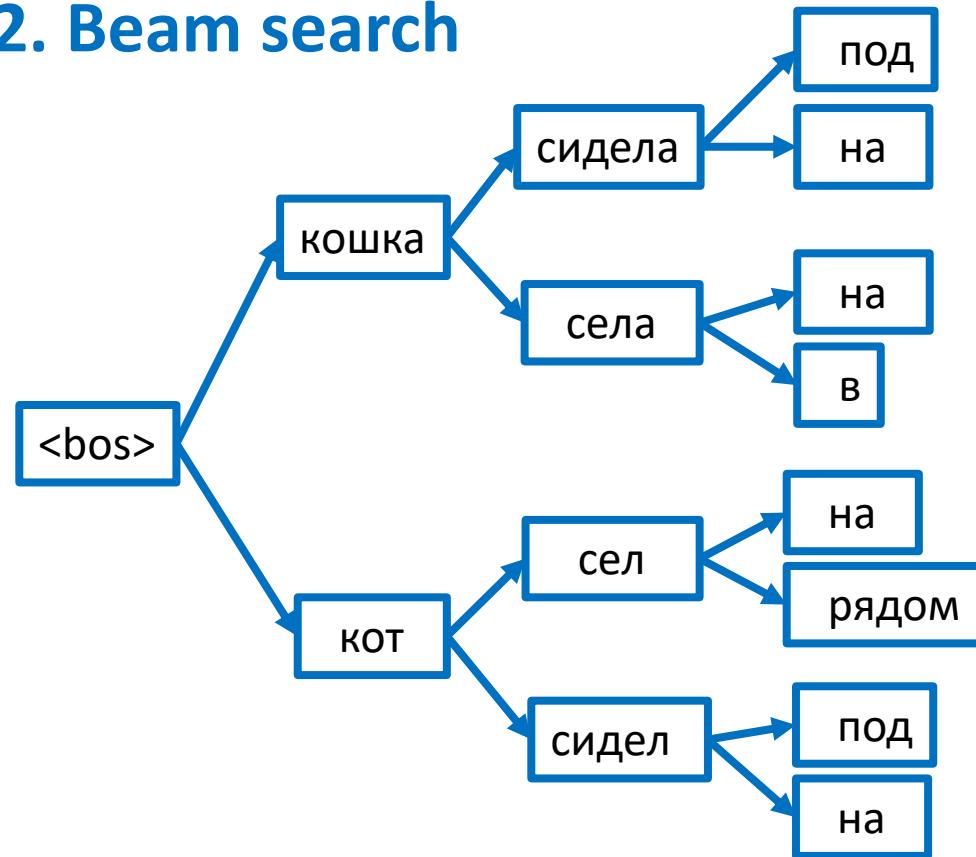
---

## 2. Beam search



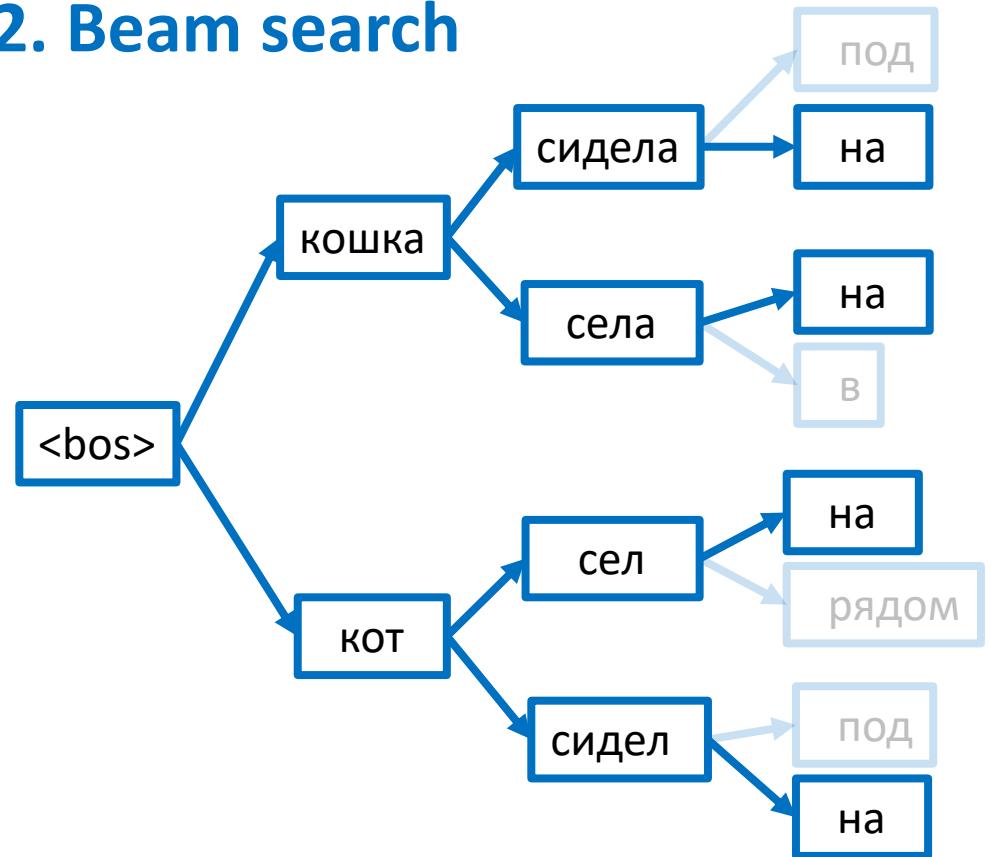
# Decoding: how to generate a sequence?

## 2. Beam search



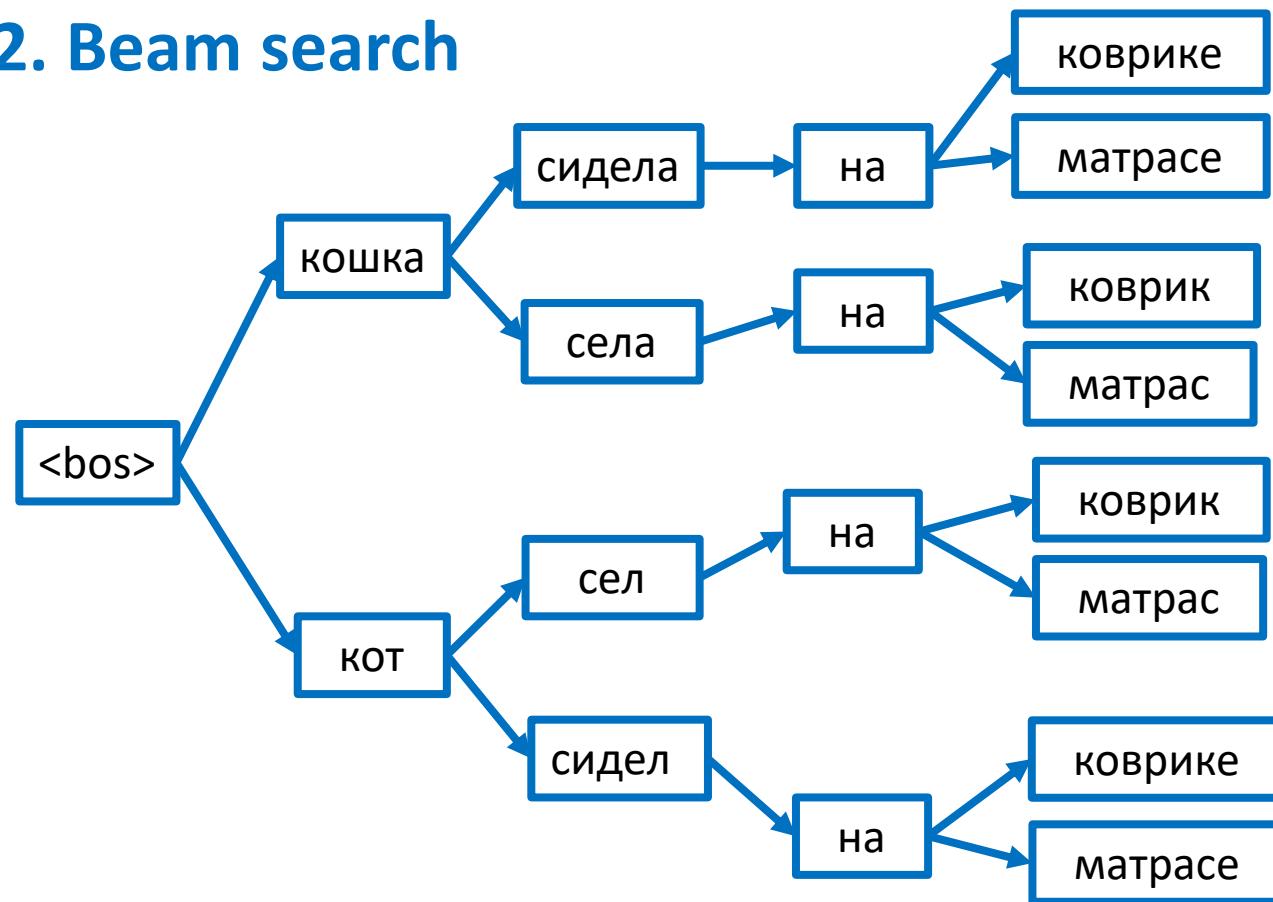
# Decoding: how to generate a sequence?

## 2. Beam search



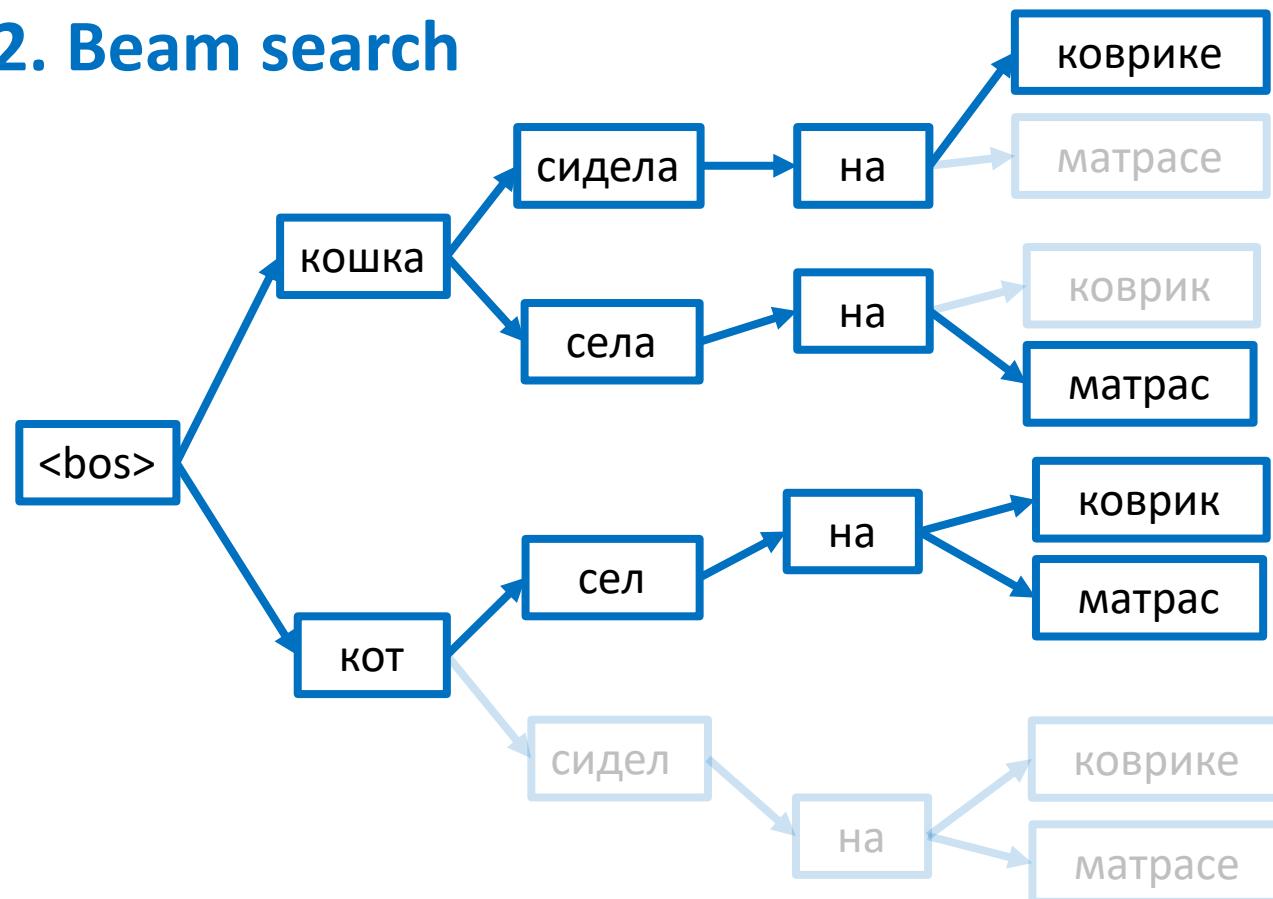
# Decoding: how to generate a sequence?

## 2. Beam search



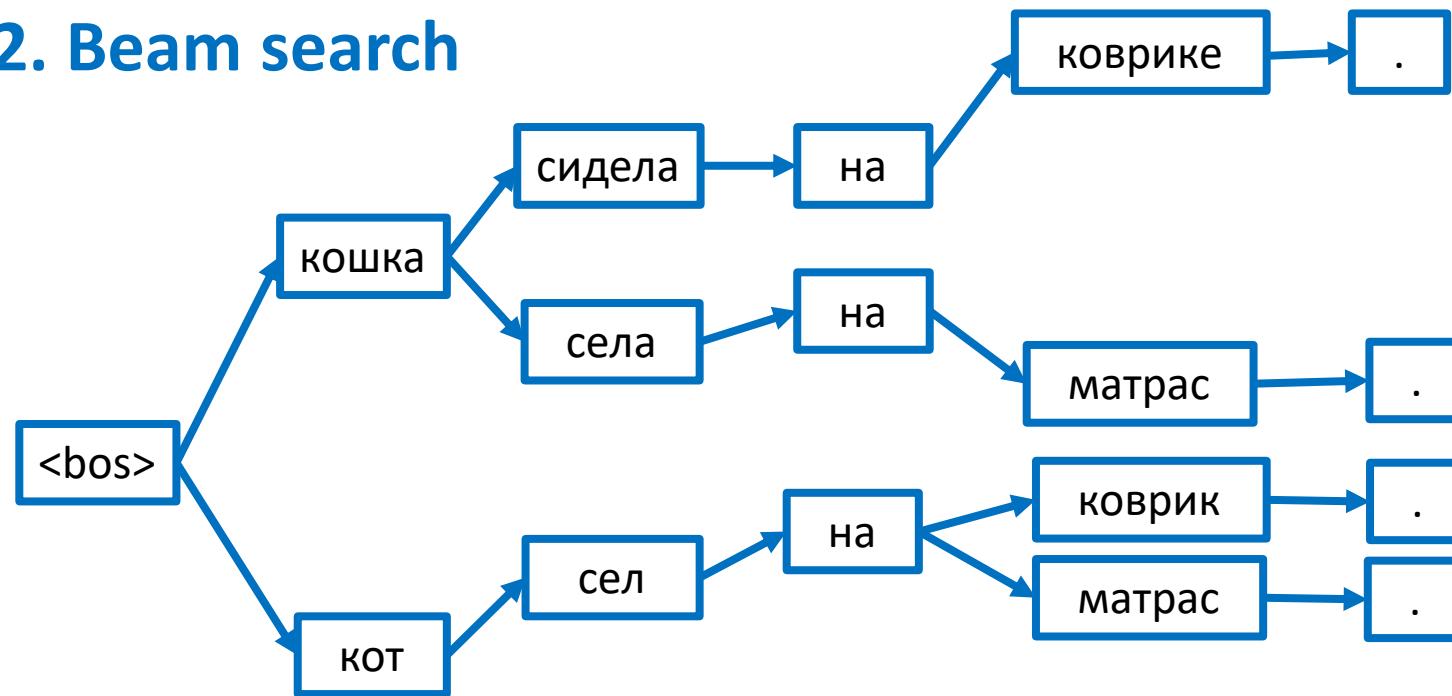
# Decoding: how to generate a sequence?

## 2. Beam search



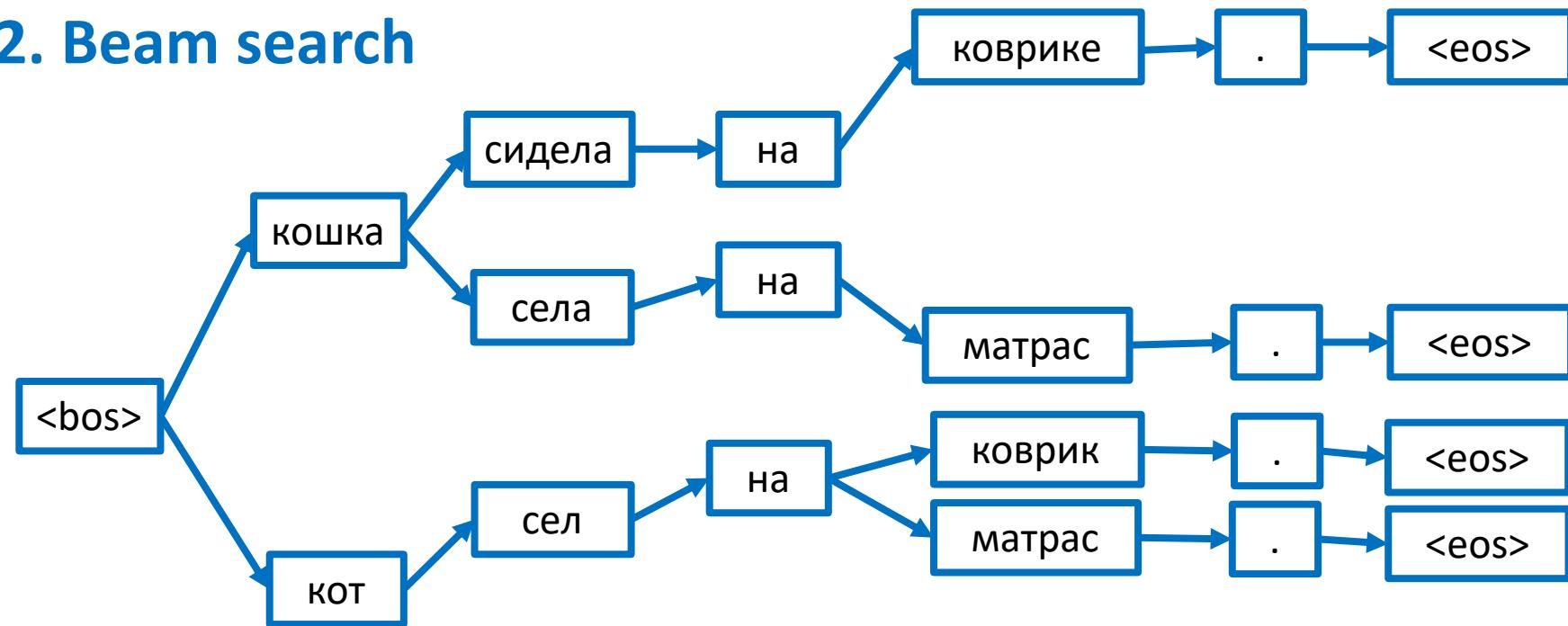
# Decoding: how to generate a sequence?

## 2. Beam search



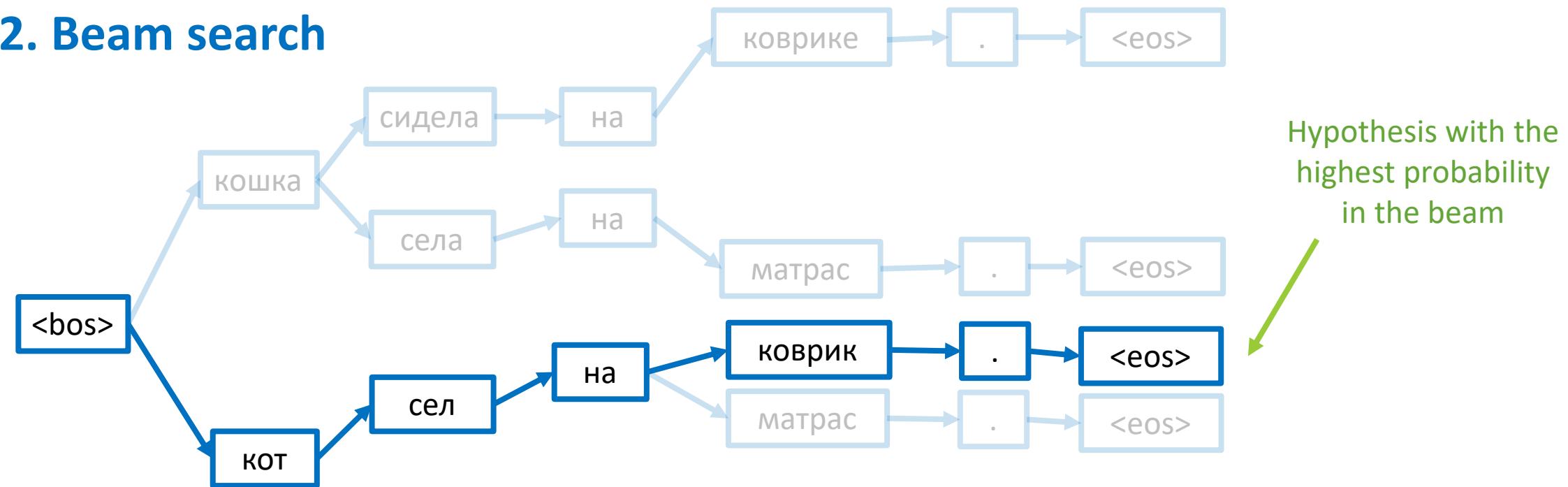
# Decoding: how to generate a sequence?

## 2. Beam search



# Decoding: how to generate a sequence?

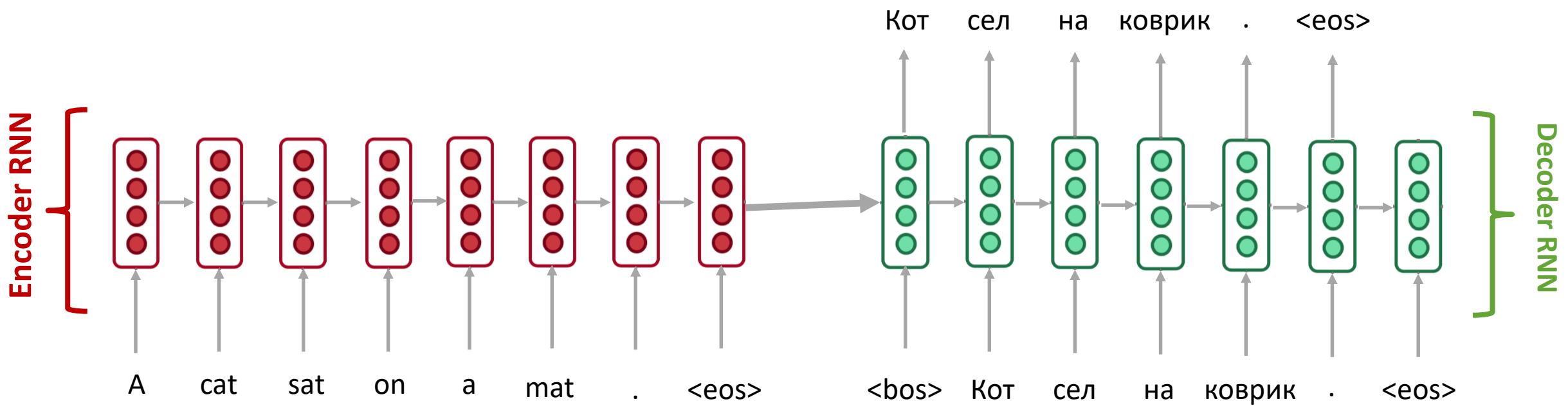
## 2. Beam search



# Attention-based models

---

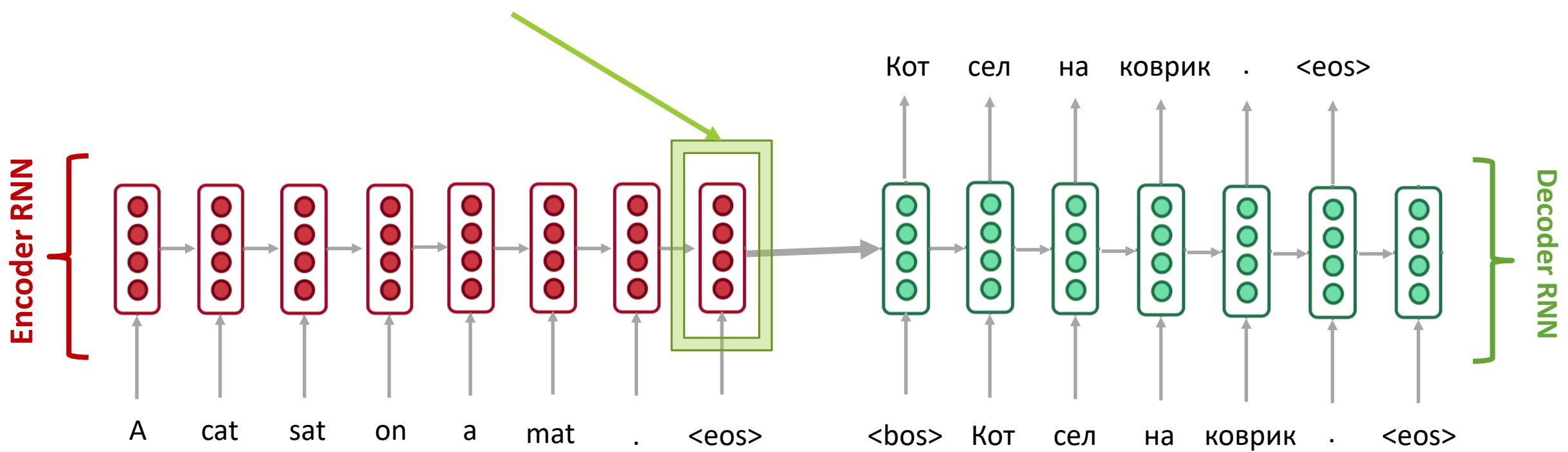
# What is the problem with such models?



# What is the problem with such models?

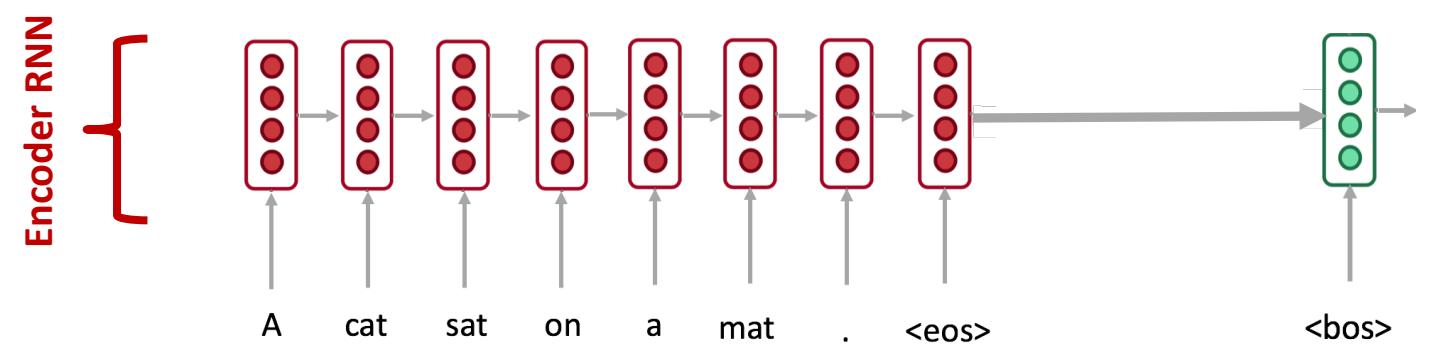
Information bottleneck:

the model has to put all information  
about sentence in one vector



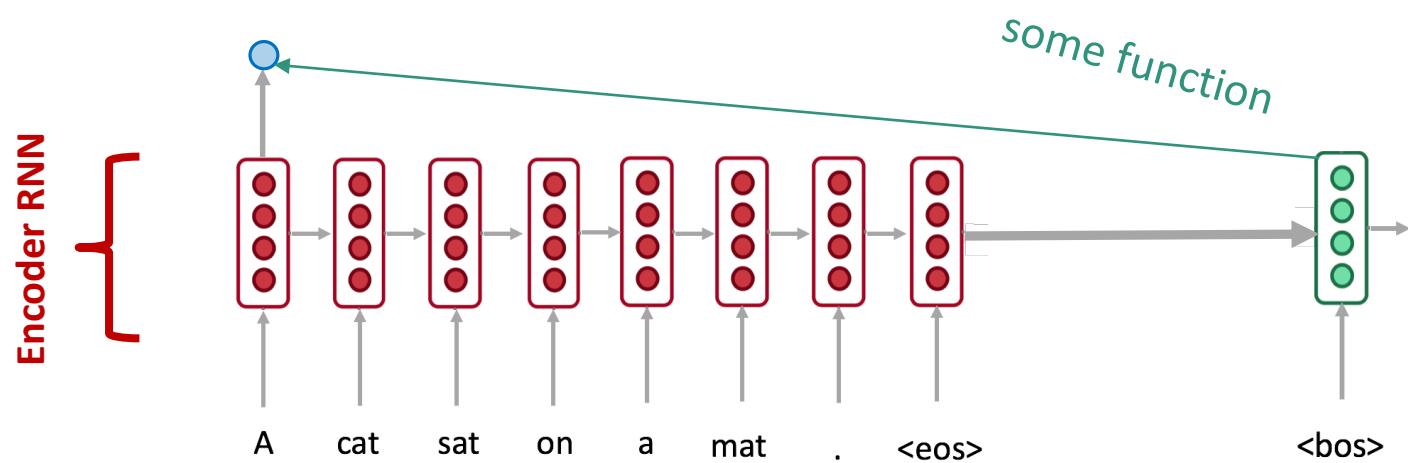
# Attention

---



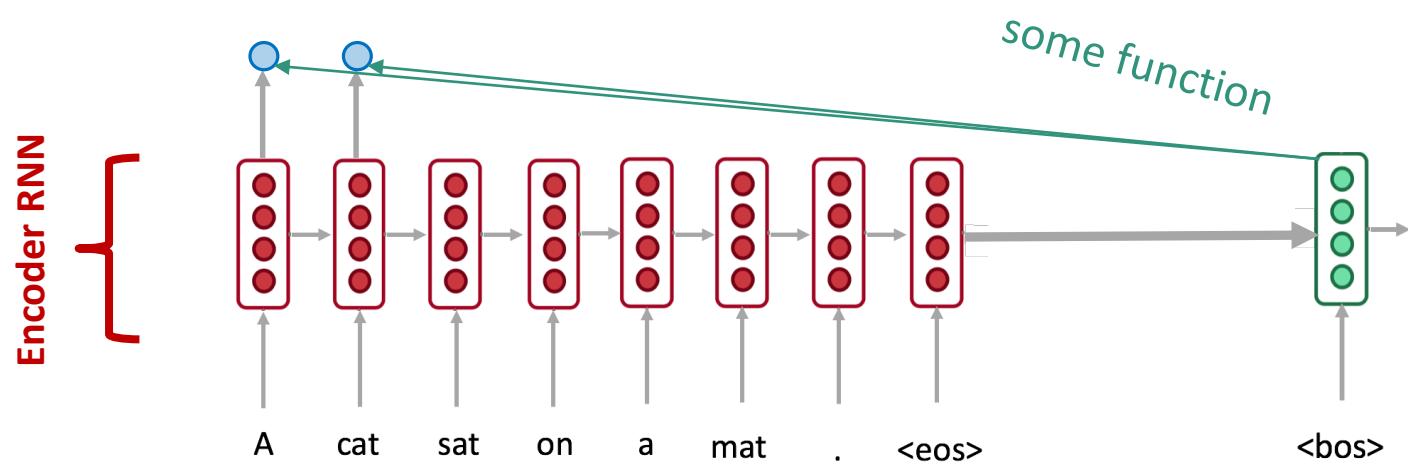
# Attention

---



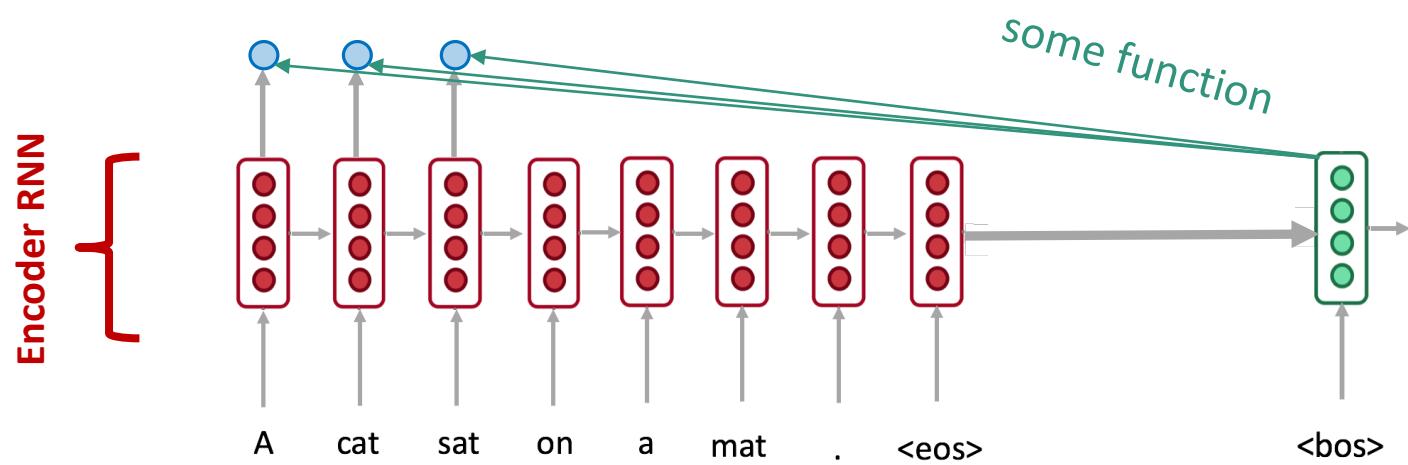
# Attention

---



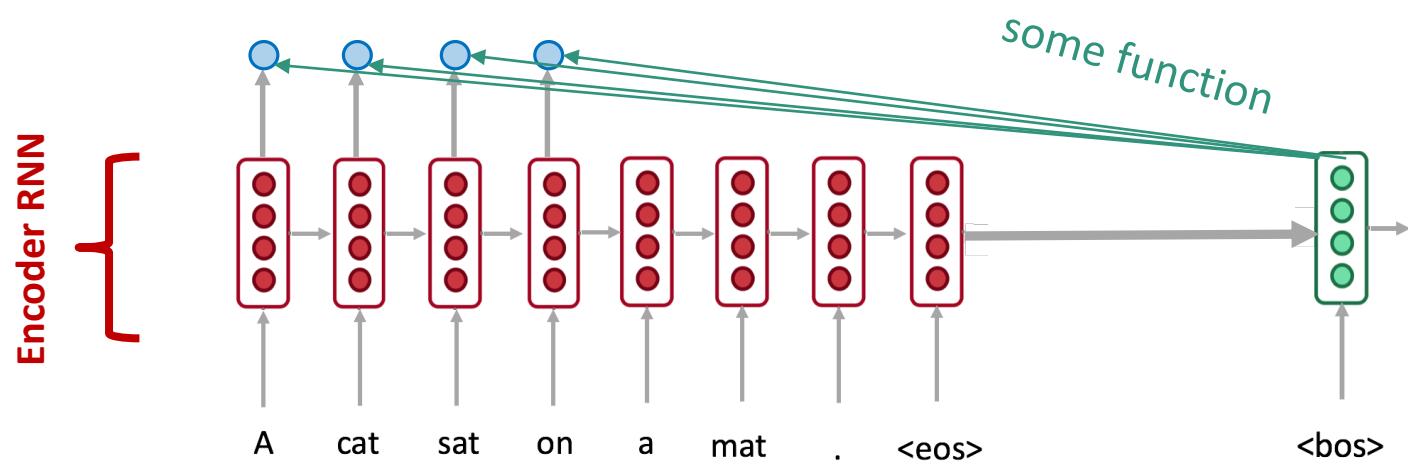
# Attention

---



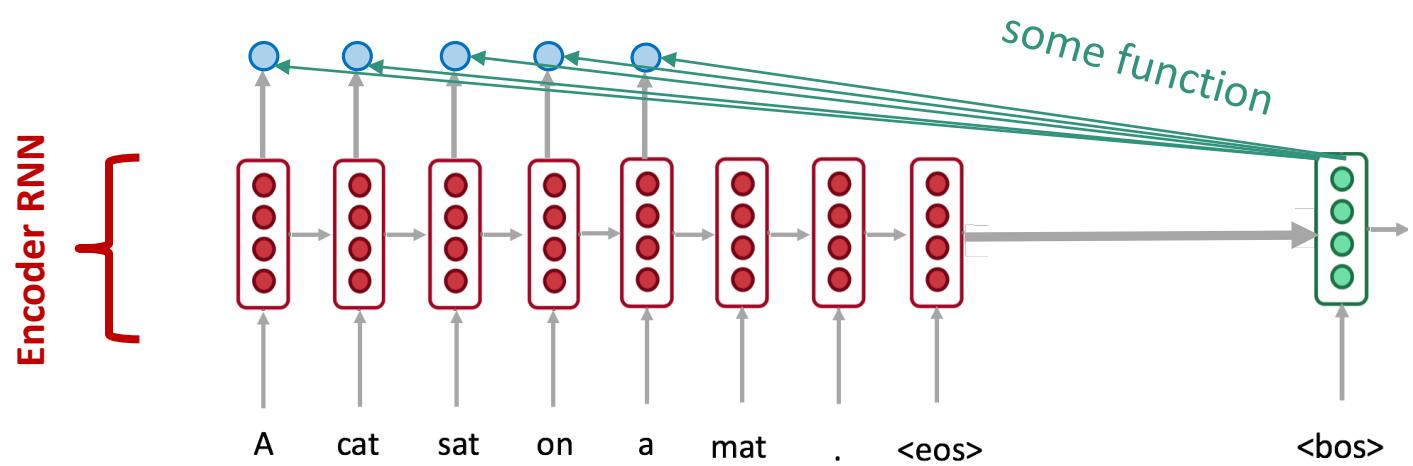
# Attention

---



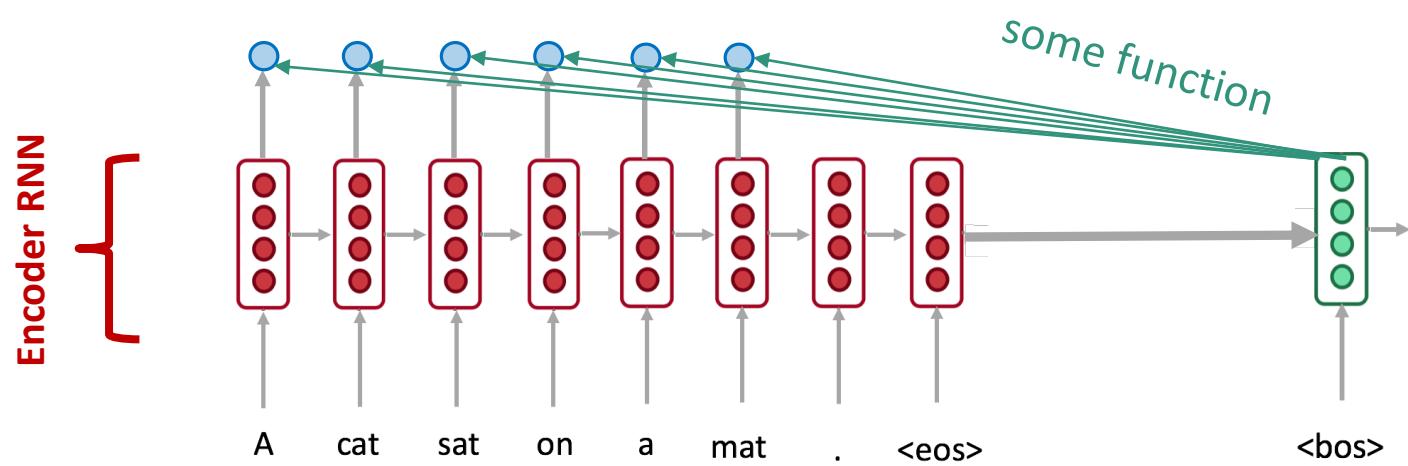
# Attention

---



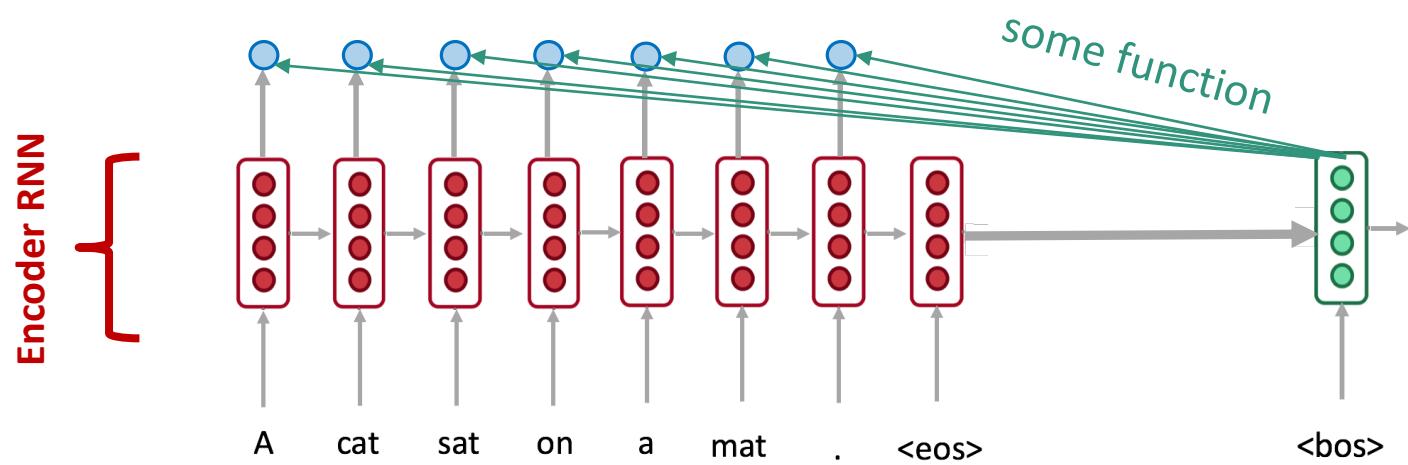
# Attention

---



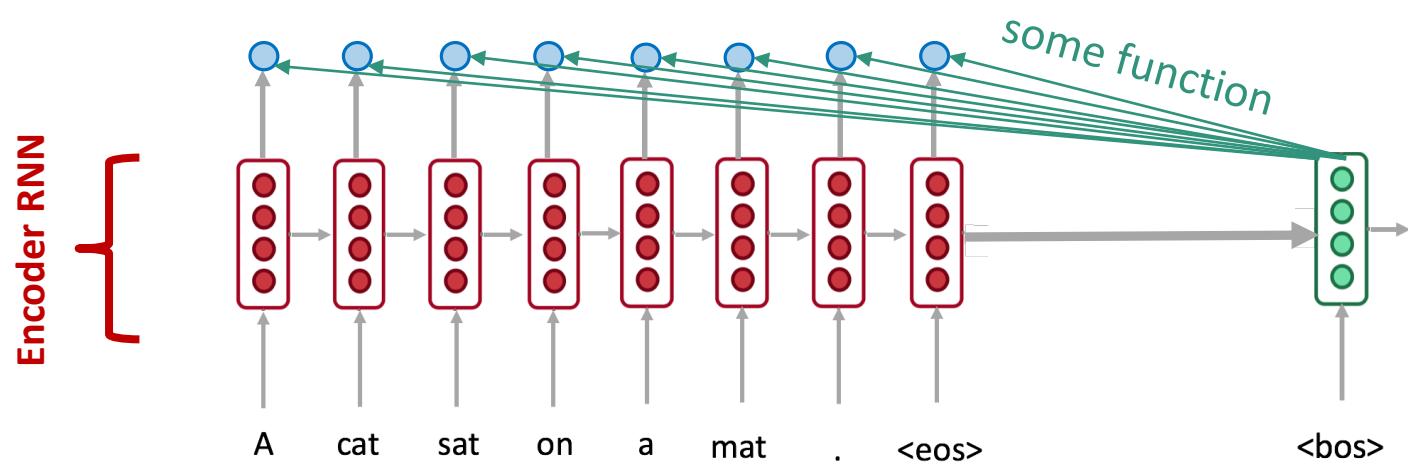
# Attention

---



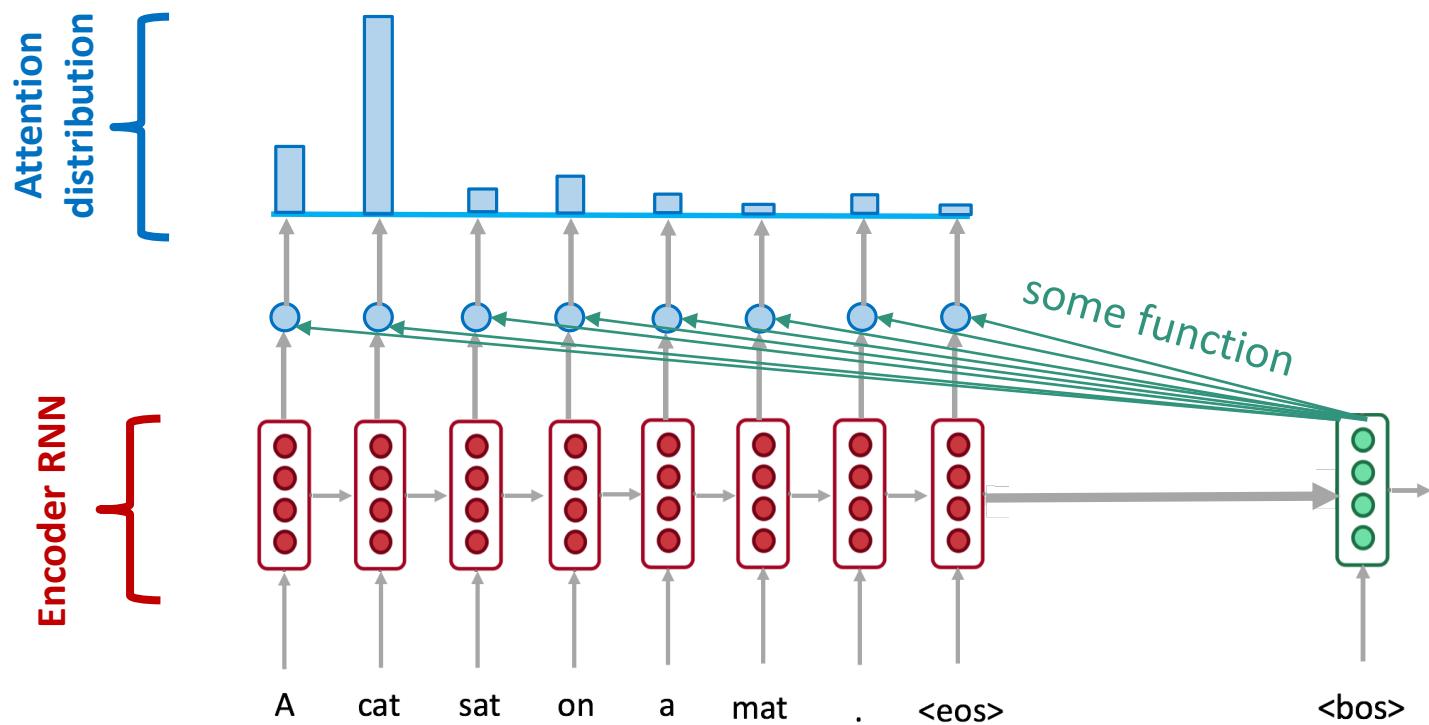
# Attention

---

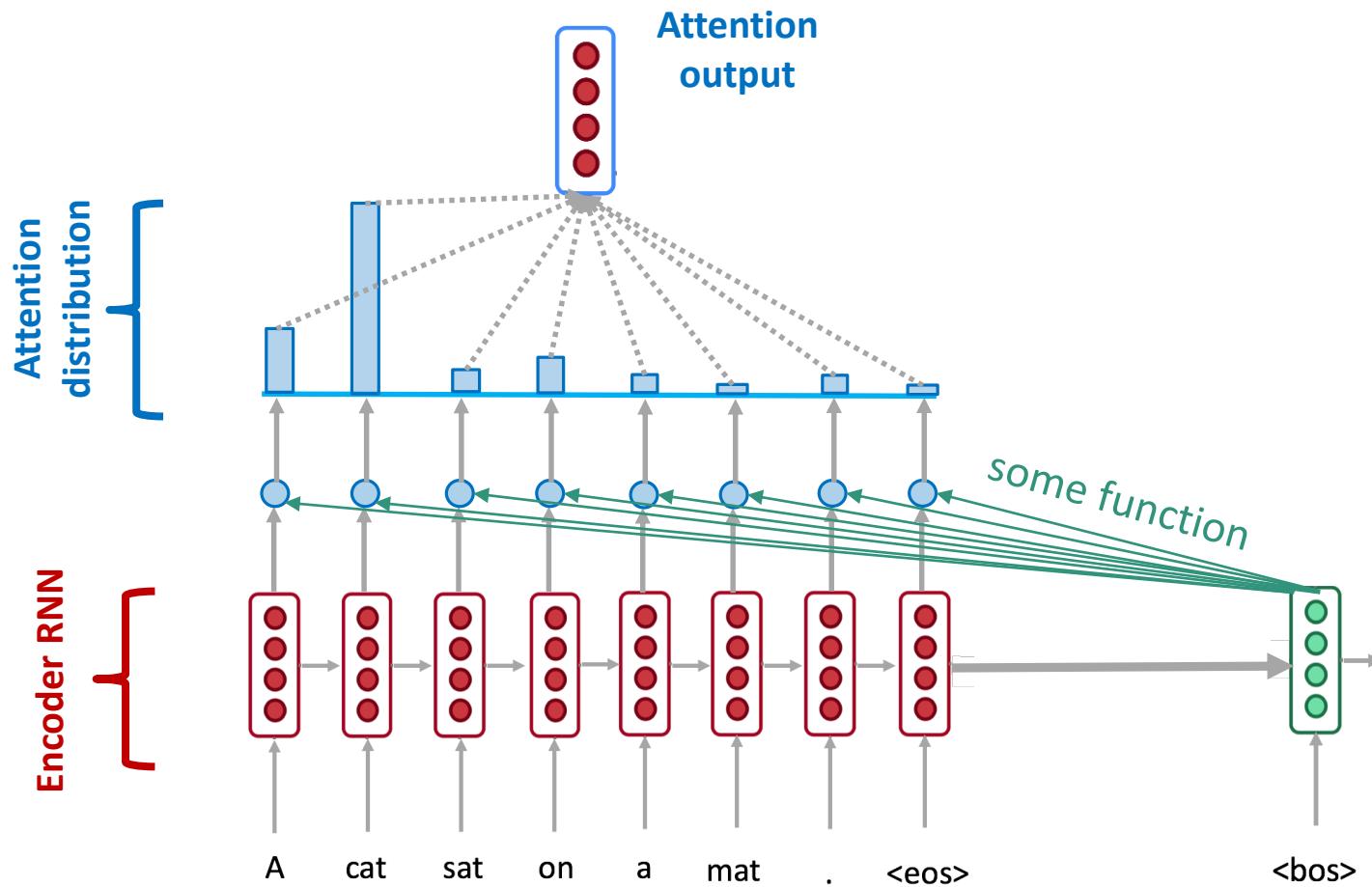


# Attention

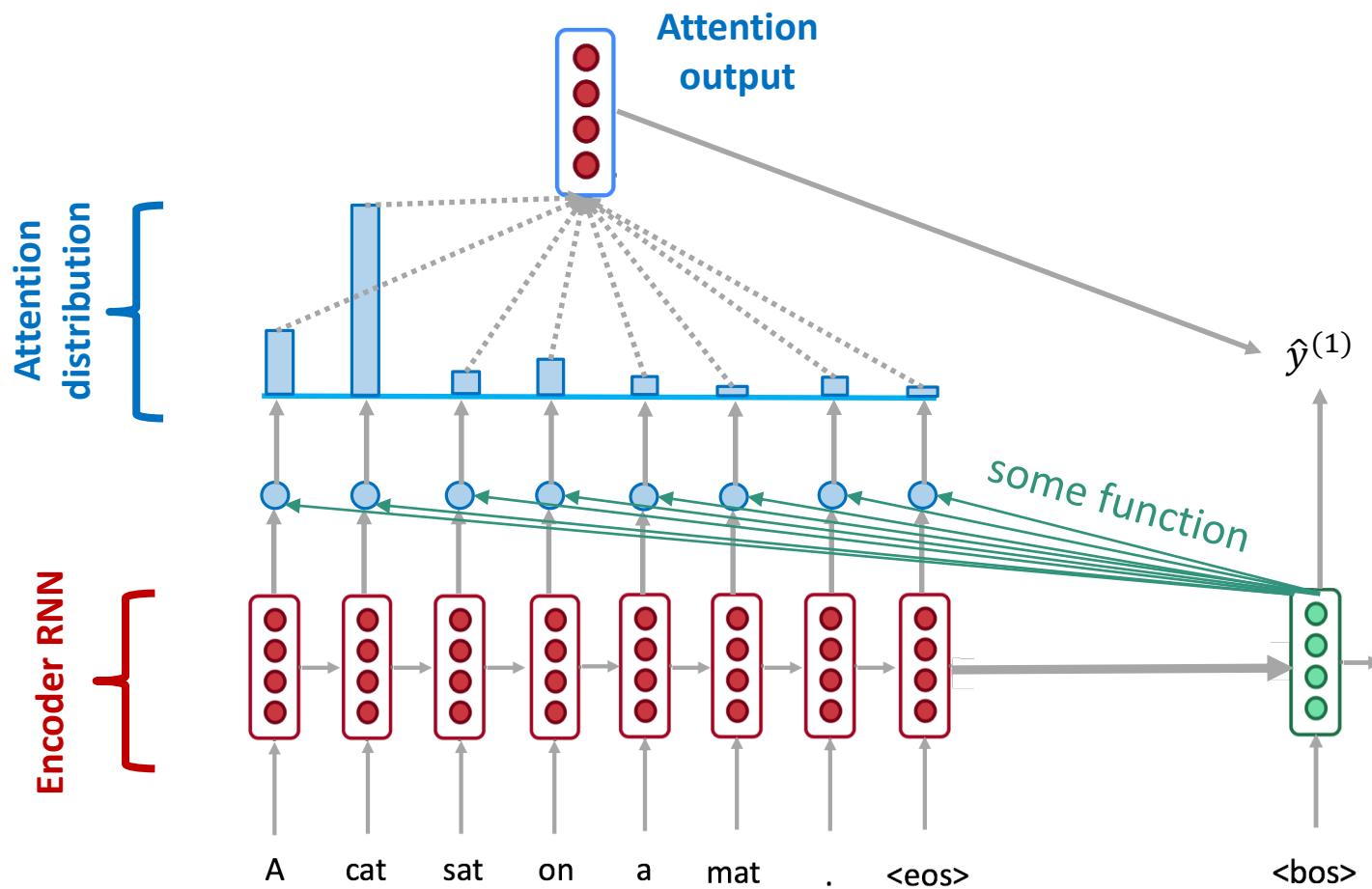
---



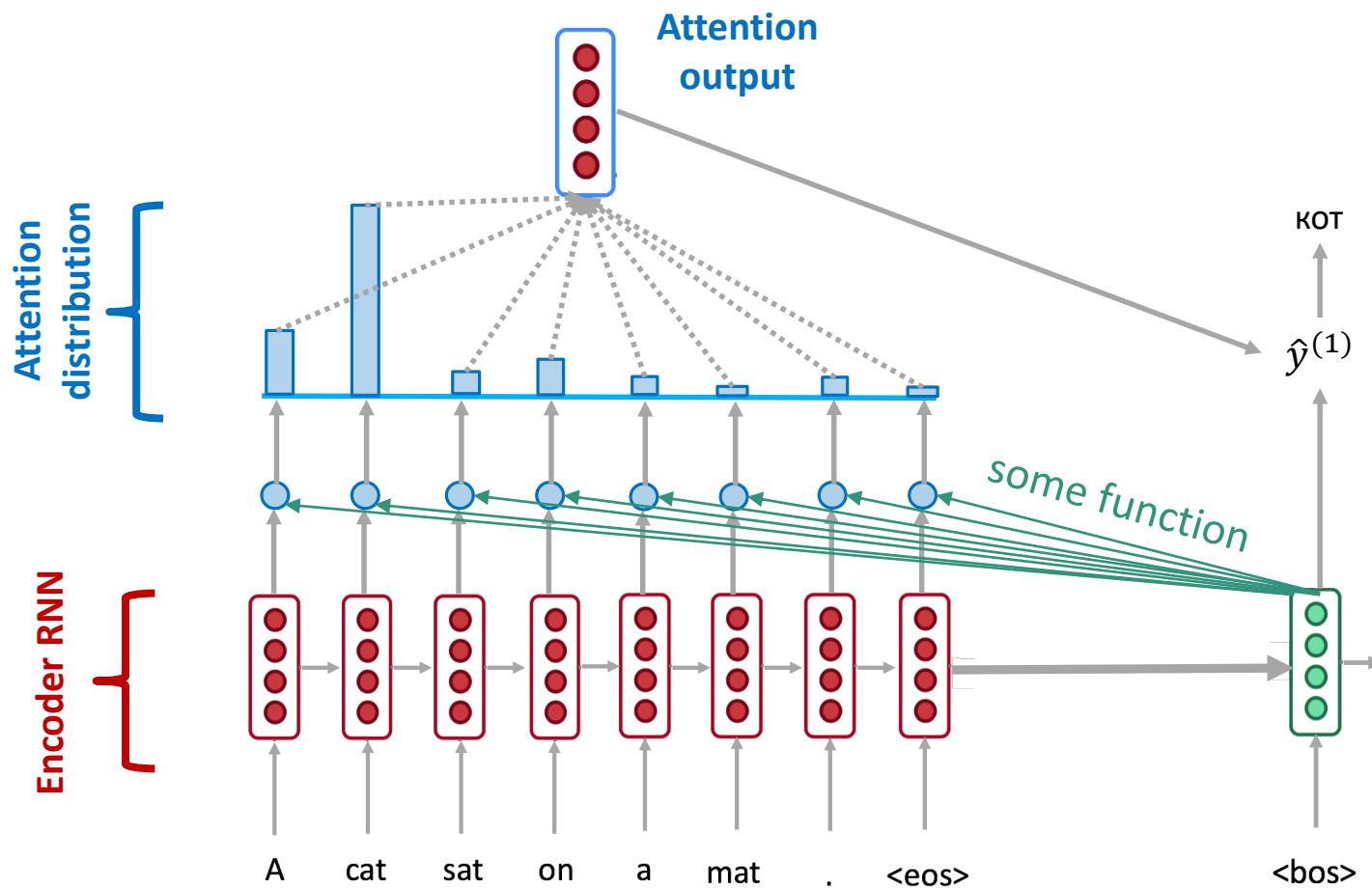
# Attention



# Attention

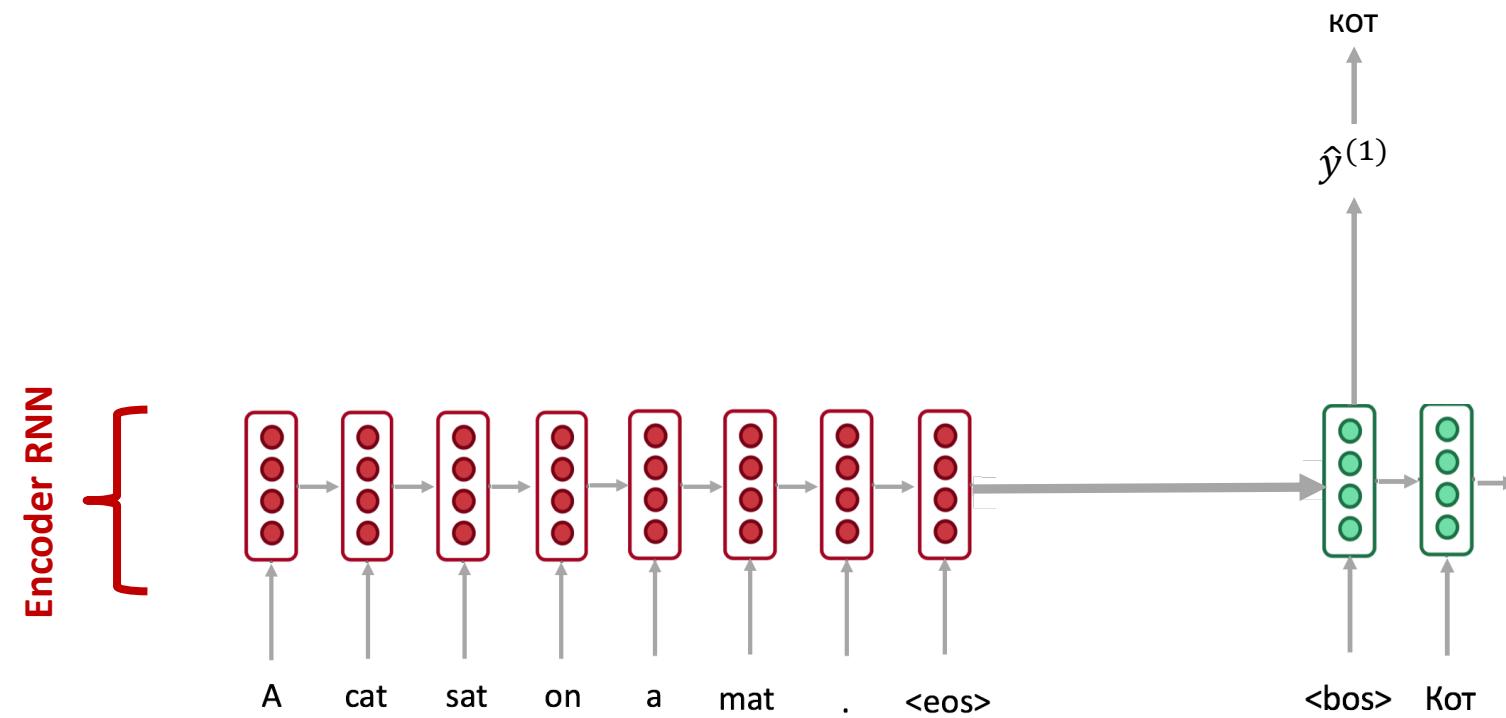


# Attention



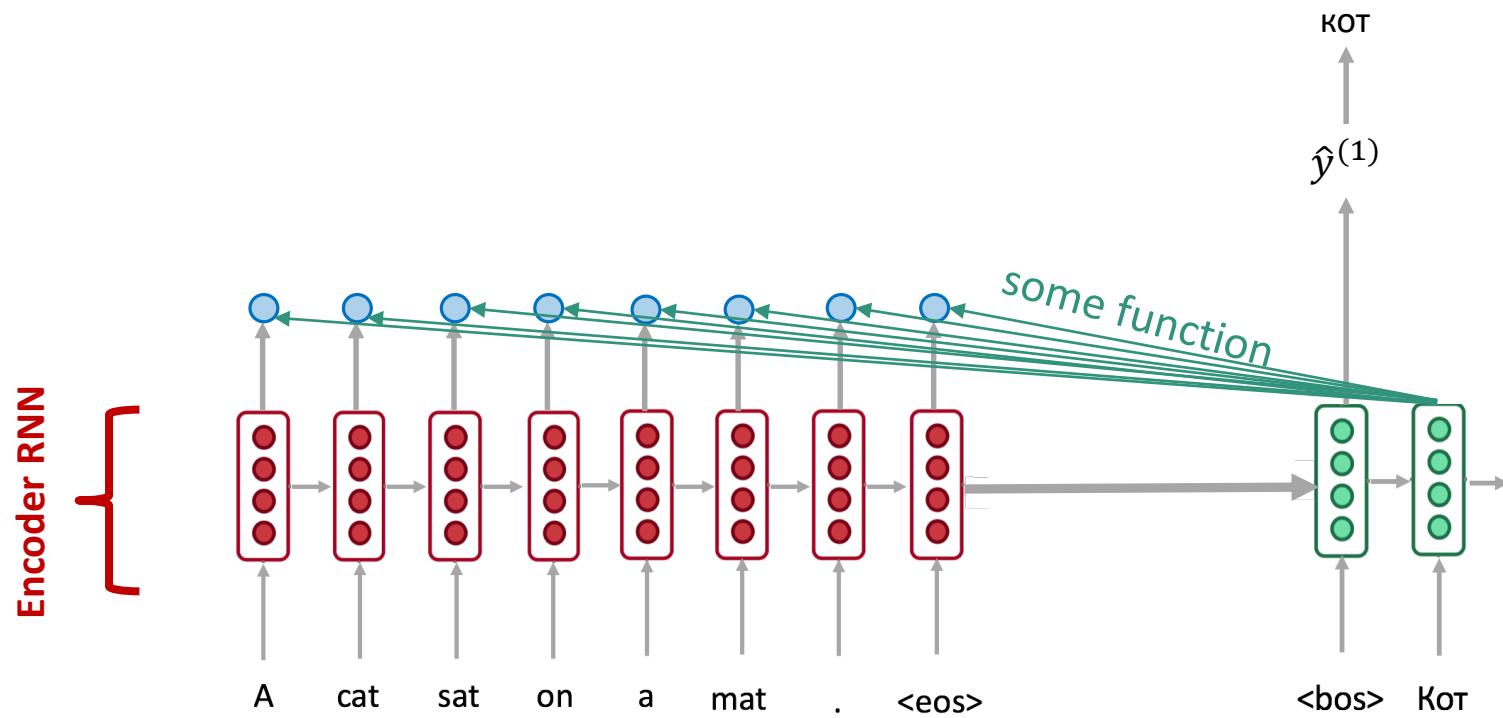
# Attention

---

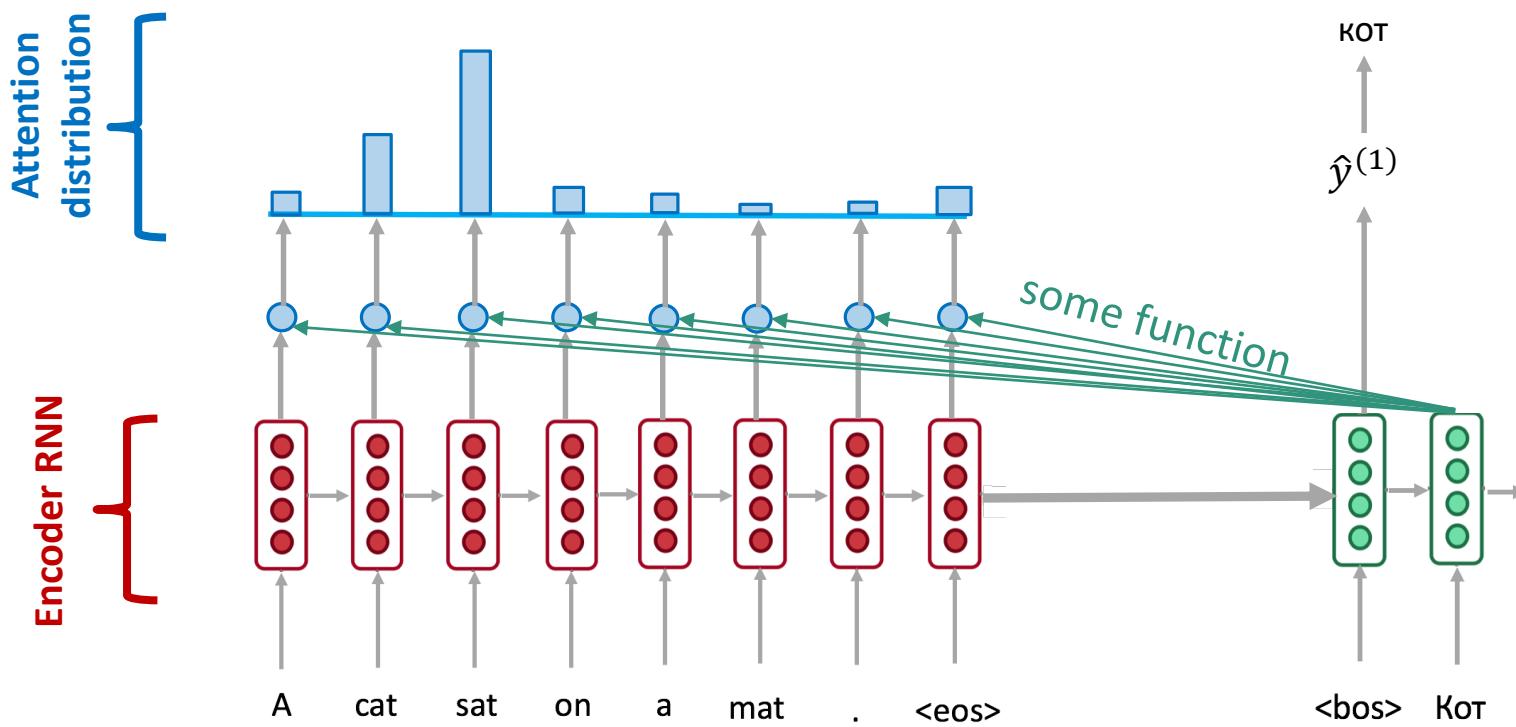


# Attention

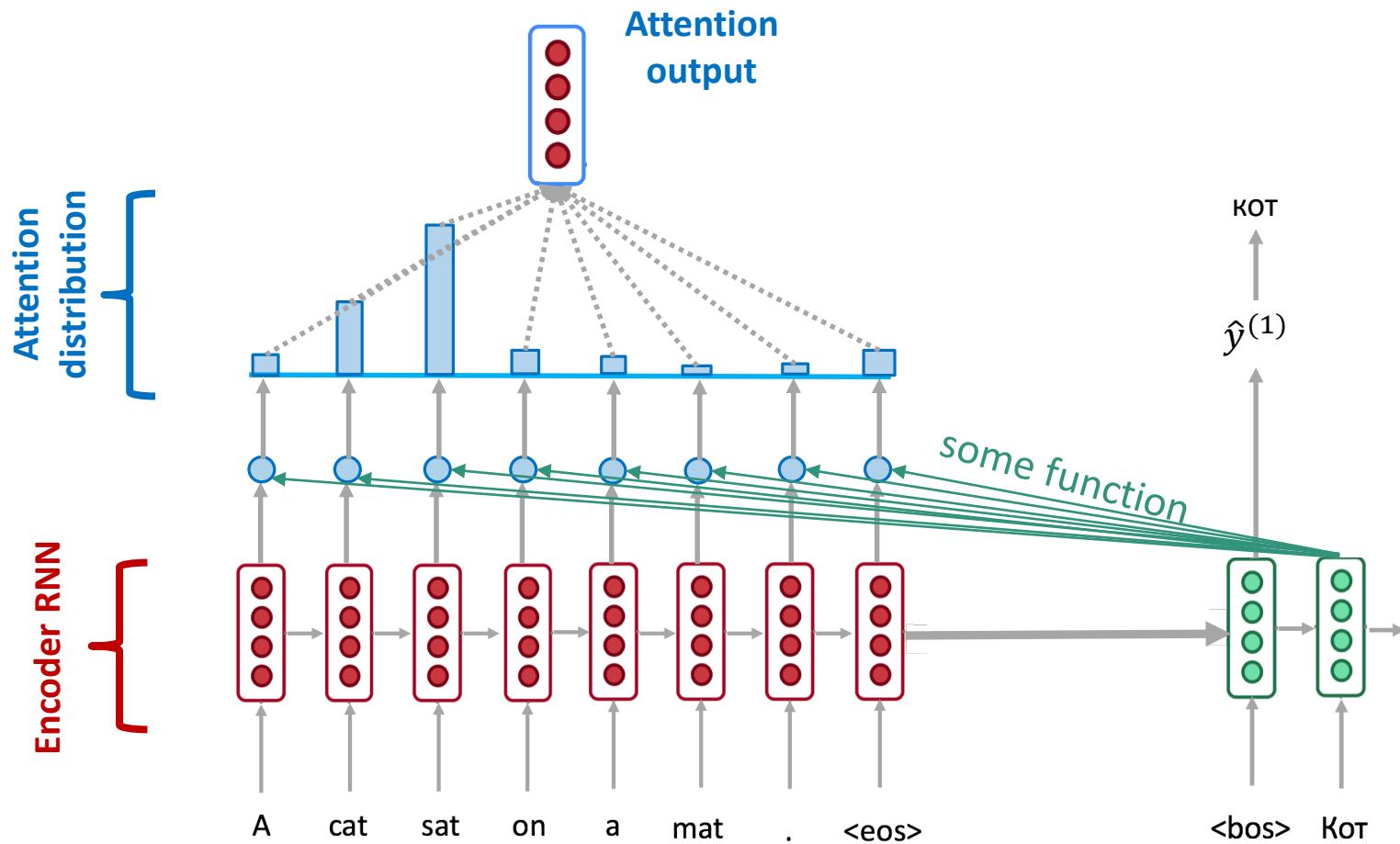
---



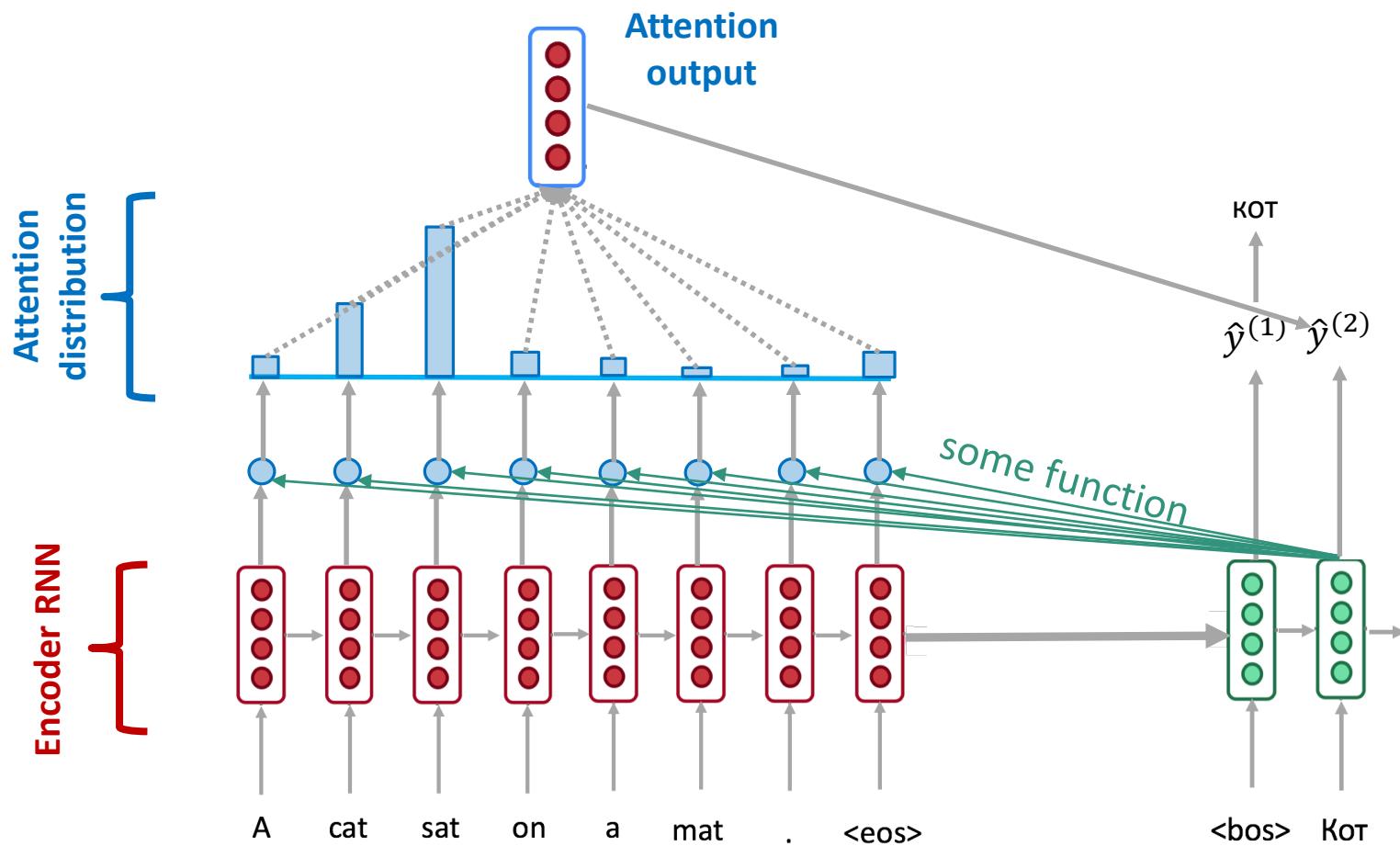
# Attention



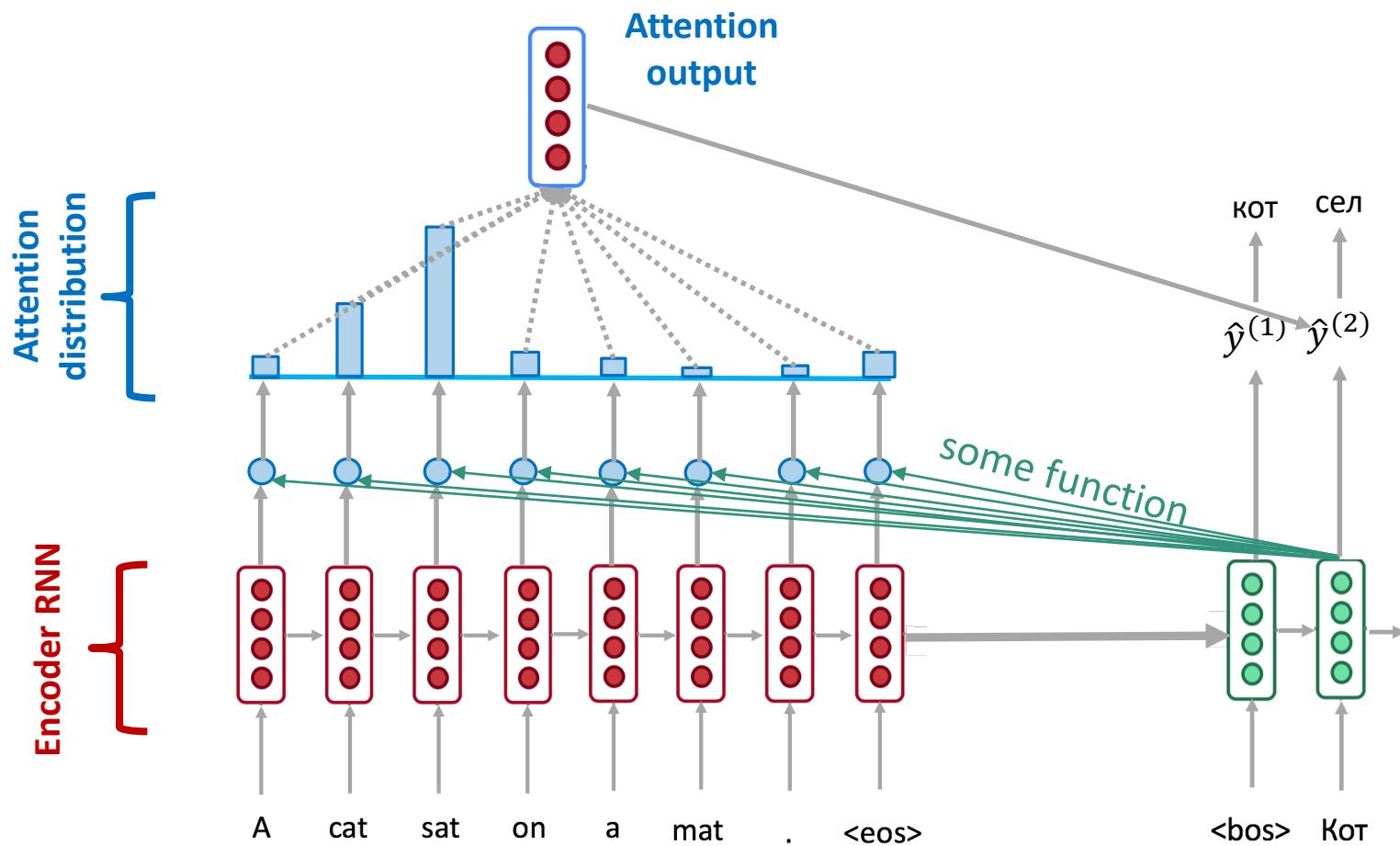
# Attention



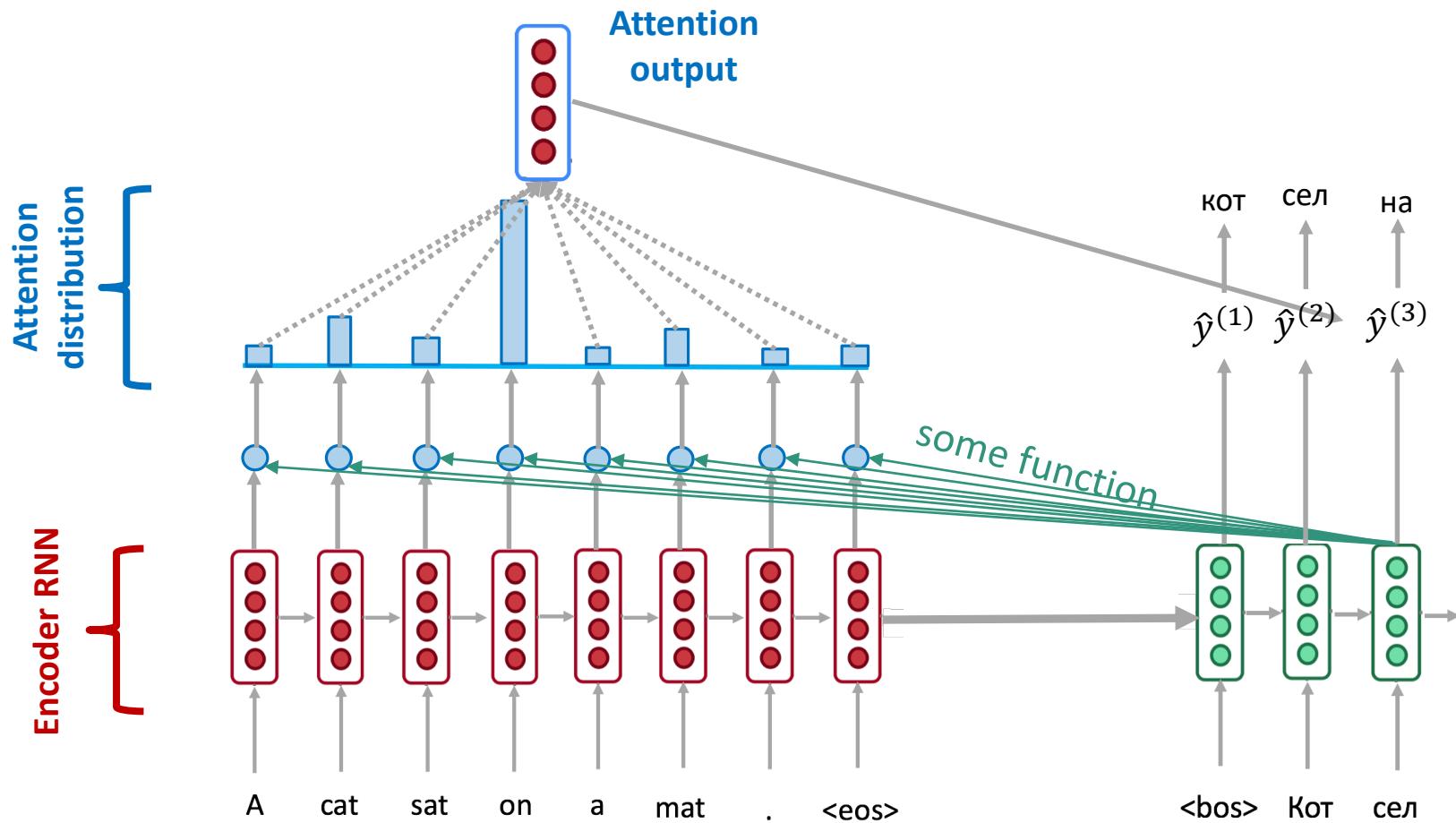
# Attention



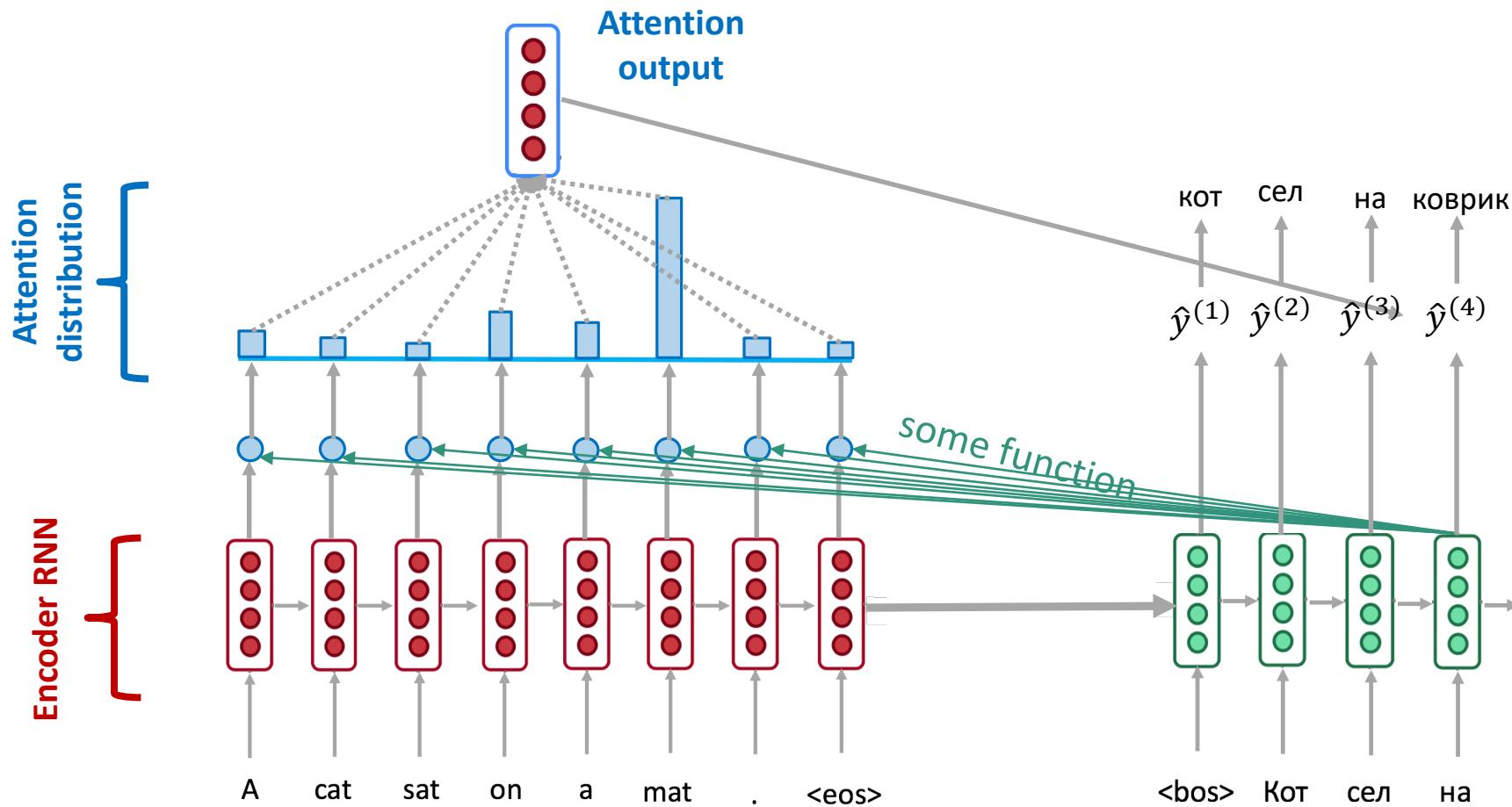
# Attention



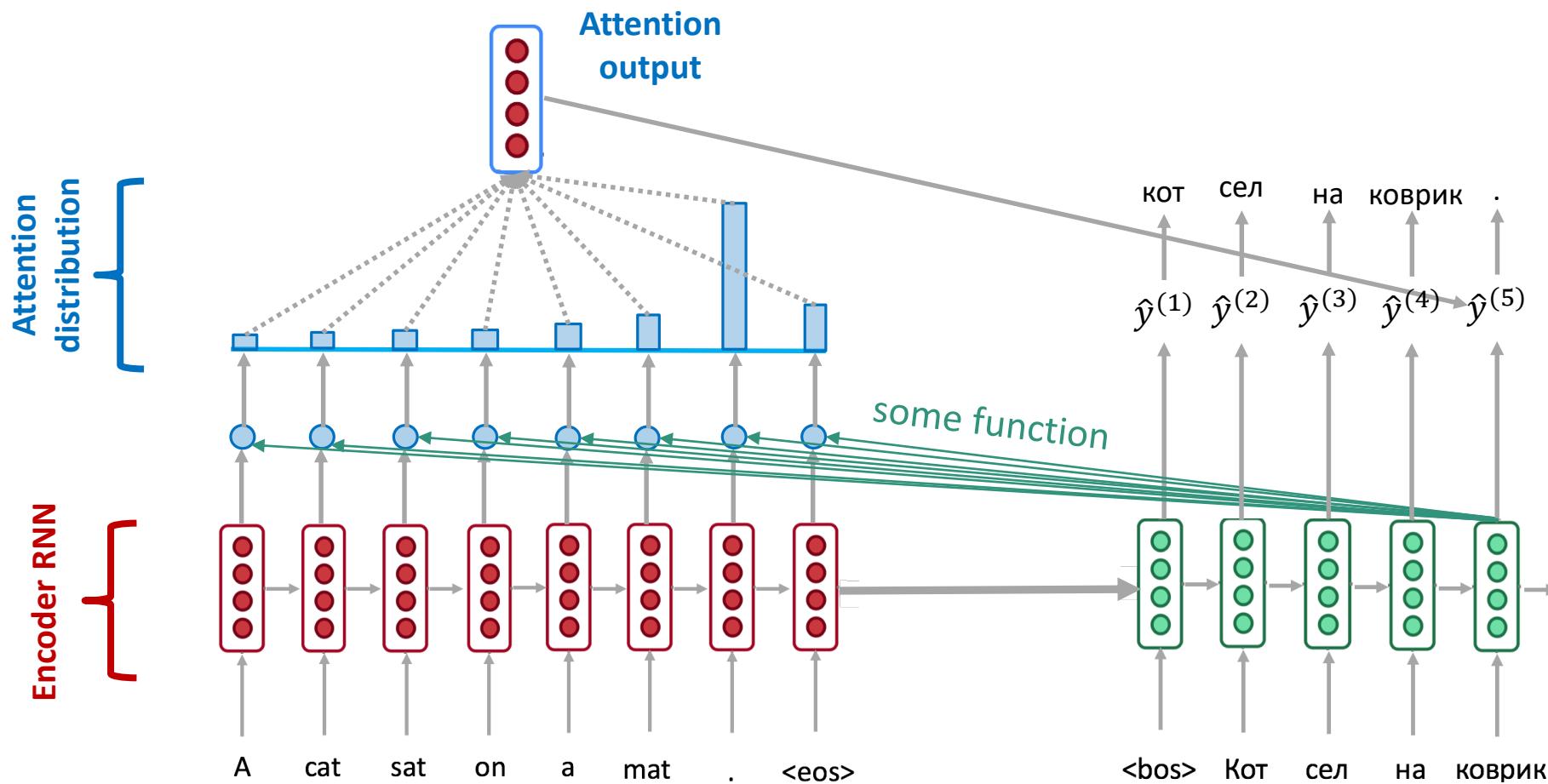
# Attention



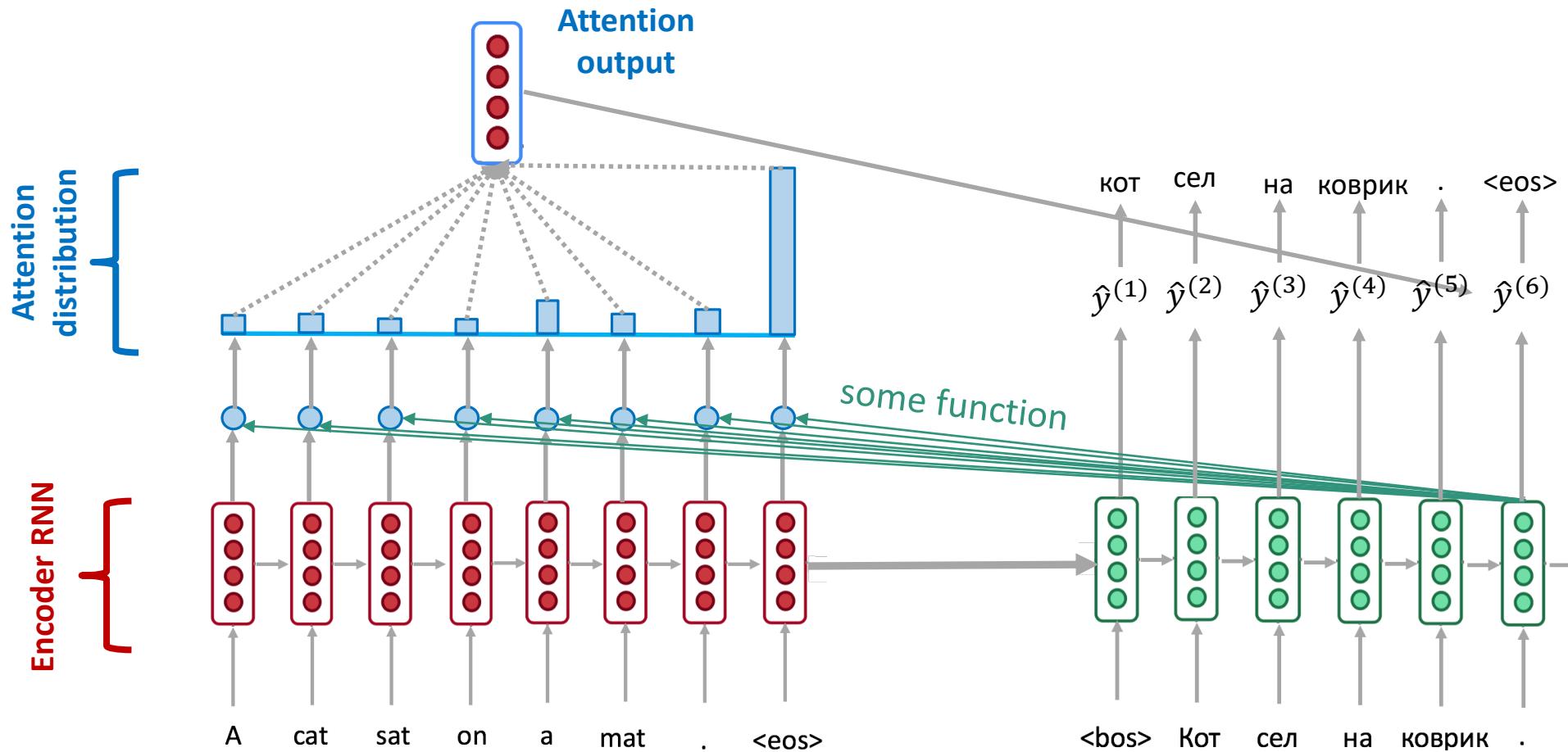
# Attention



# Attention



# Attention



# Attention: in equations

- We have encoder hidden states  $h_1, \dots, h_N \in \mathbb{R}^h$
- On timestep  $t$ , we have decoder hidden state  $s_t \in \mathbb{R}^h$
- We get the attention scores  $e^t$  for this step:

$$e^t = [s_t^T h_1, \dots, s_t^T h_N] \in \mathbb{R}^N$$

- We take softmax to get the attention distribution  $\alpha^t$  for this step (this is a probability distribution and sums to 1)

$$\alpha^t = \text{softmax}(e^t) \in \mathbb{R}^N$$

- We use  $\alpha^t$  to take a weighted sum of the encoder hidden states to get the attention output  $a_t$

$$a_t = \sum_{i=1}^N \alpha_i^t h_i \in \mathbb{R}^h$$

- Finally we concatenate the attention output  $a_t$  with the decoder hidden state  $s_t$  and proceed as in the non-attention seq2seq model

$$[a_t; s_t] \in \mathbb{R}^{2h}$$