

Seq2Seq, Attention

# Machine translation task

# Machine translation task

*source* =  $(x_1, x_2, \dots, x_n)$

The cat sits on the floor

*target* =  $(y_1, y_2, \dots, y_m)$

Кошка сидит на полу

Machine translation task is to find the most probable translation given source:

$$\widehat{target} = \underset{target}{argmax} P(target \mid source, \theta)$$

# Machine translation task


*source* =  $(x_1, x_2, \dots, x_n)$

The cat sits on the floor

*target* =  $(y_1, y_2, \dots, y_m)$

Кошка сидит на полу

$$\widehat{target} = \underset{target}{argmax} \quad P(target \mid source, \theta)$$


$$P(target \mid source) = P(y_1 \mid source) \cdot P(y_2 \mid y_1, source) \dots P(y_m \mid y_1, \dots, y_{m-1}, source)$$

# Machine translation task

*source* =  $(x_1, x_2, \dots, x_n)$

The cat sits on the floor

*target* =  $(y_1, y_2, \dots, y_m)$

Кошка сидит на полу

$$P(\textit{target} \mid \textit{source}) = P(y_1 \mid \textit{source}) \cdot P(y_2 \mid y_1, \textit{source}) \dots P(y_m \mid y_1, \dots, y_{m-1}, \textit{source})$$

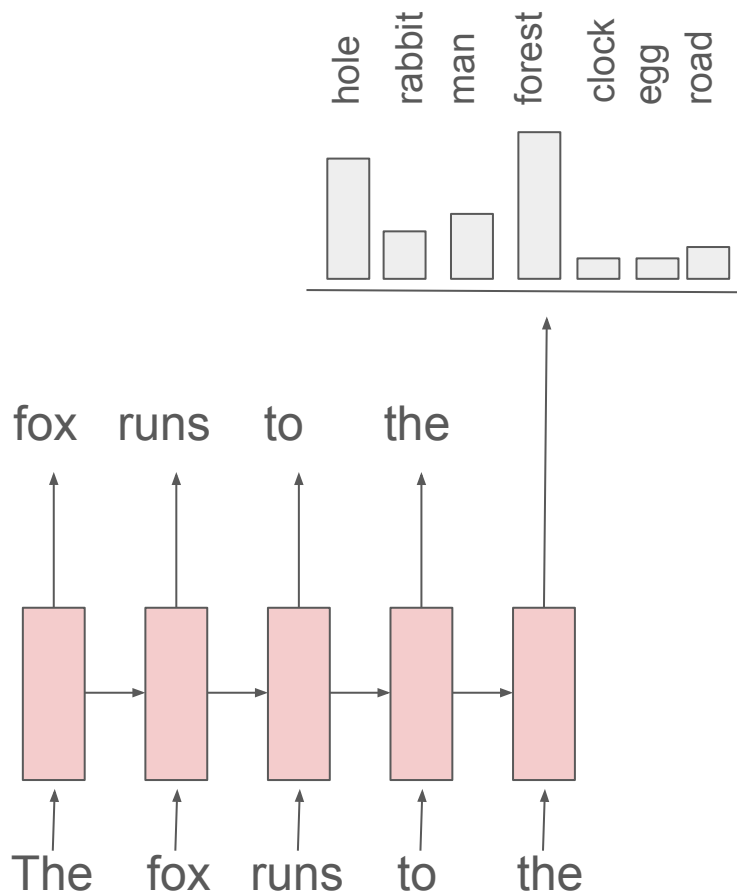
Conditional language model (conditioned on source)

# Machine translation techniques

- rule-based (1950th)
- statistics-based
- neural-based (2010th)

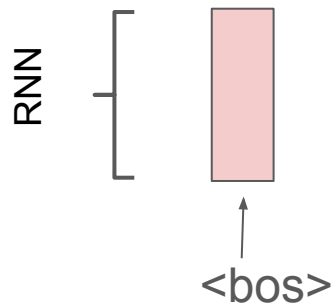
Seq2Seq

# LM (reminder)

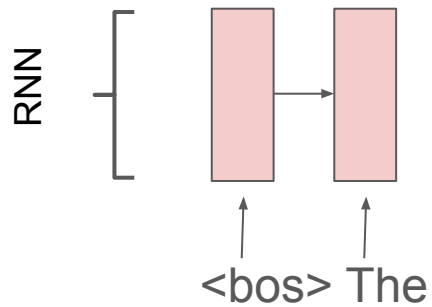




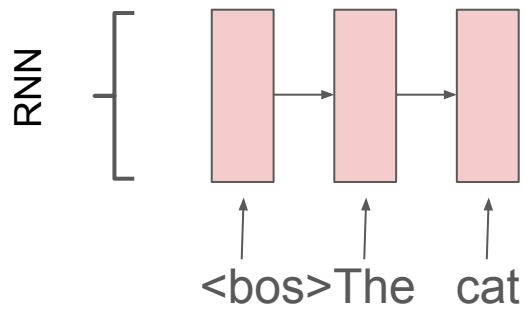
# Encoder-Decoder



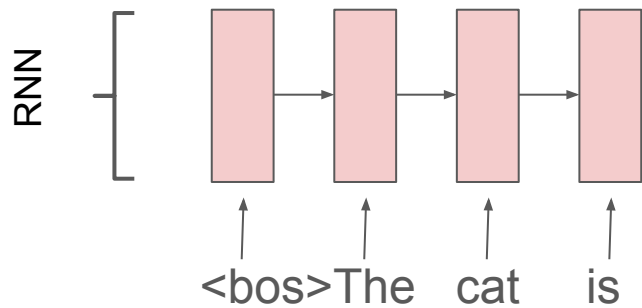
# Encoder-Decoder



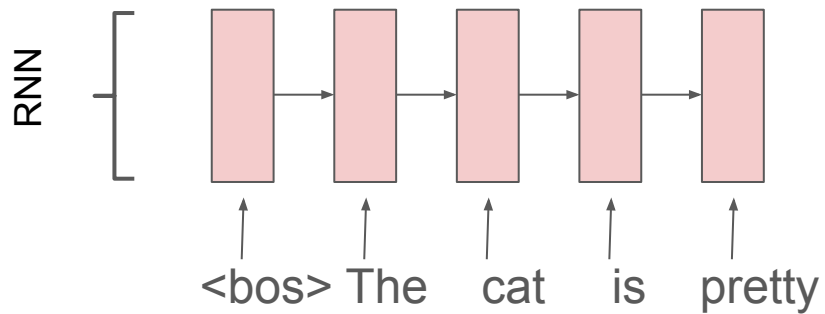
# Encoder-Decoder



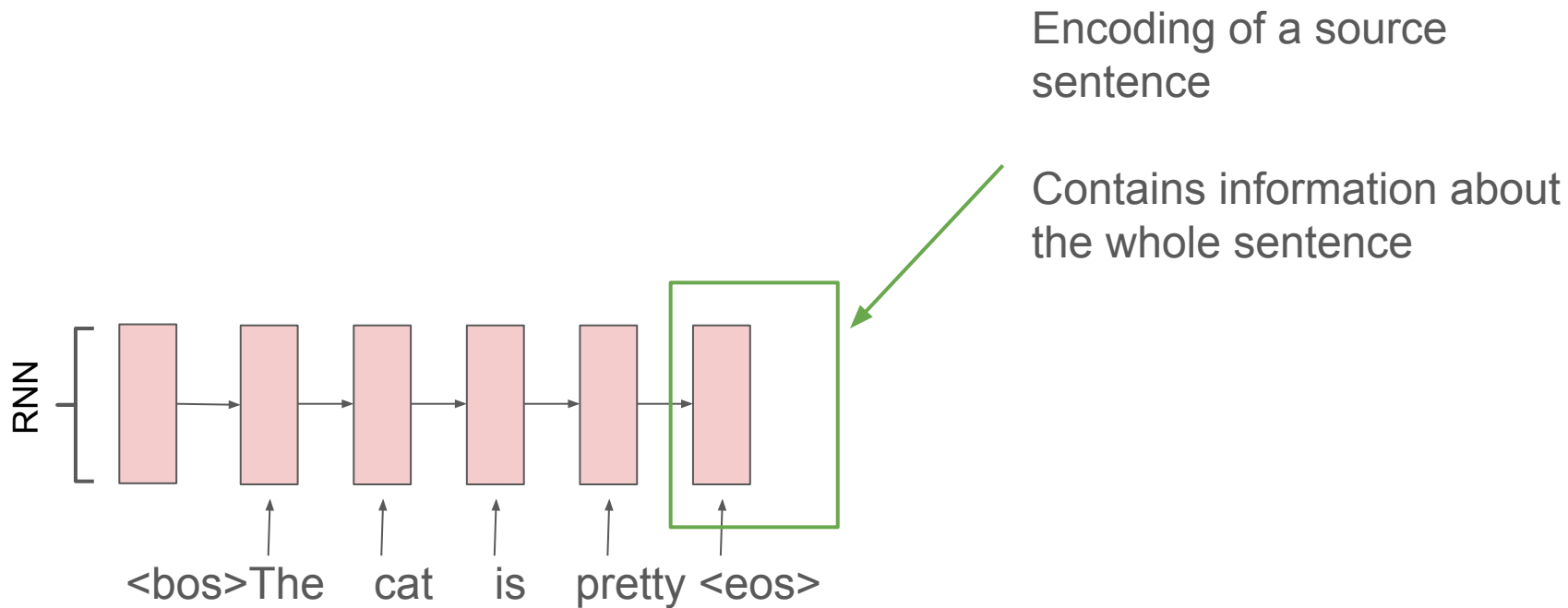
# Encoder-Decoder



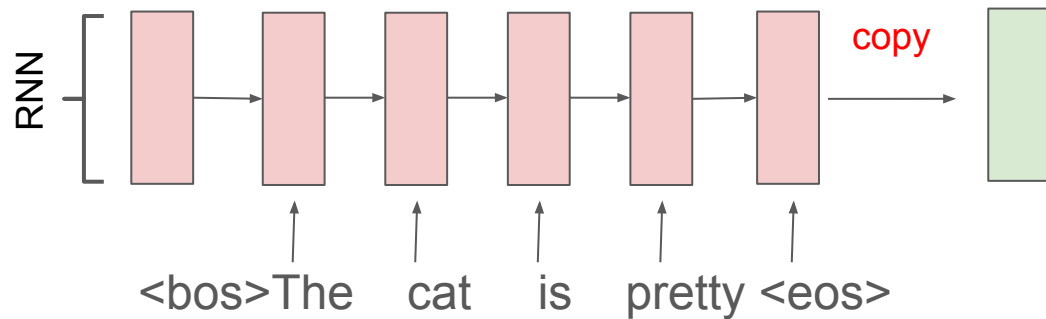
# Encoder-Decoder



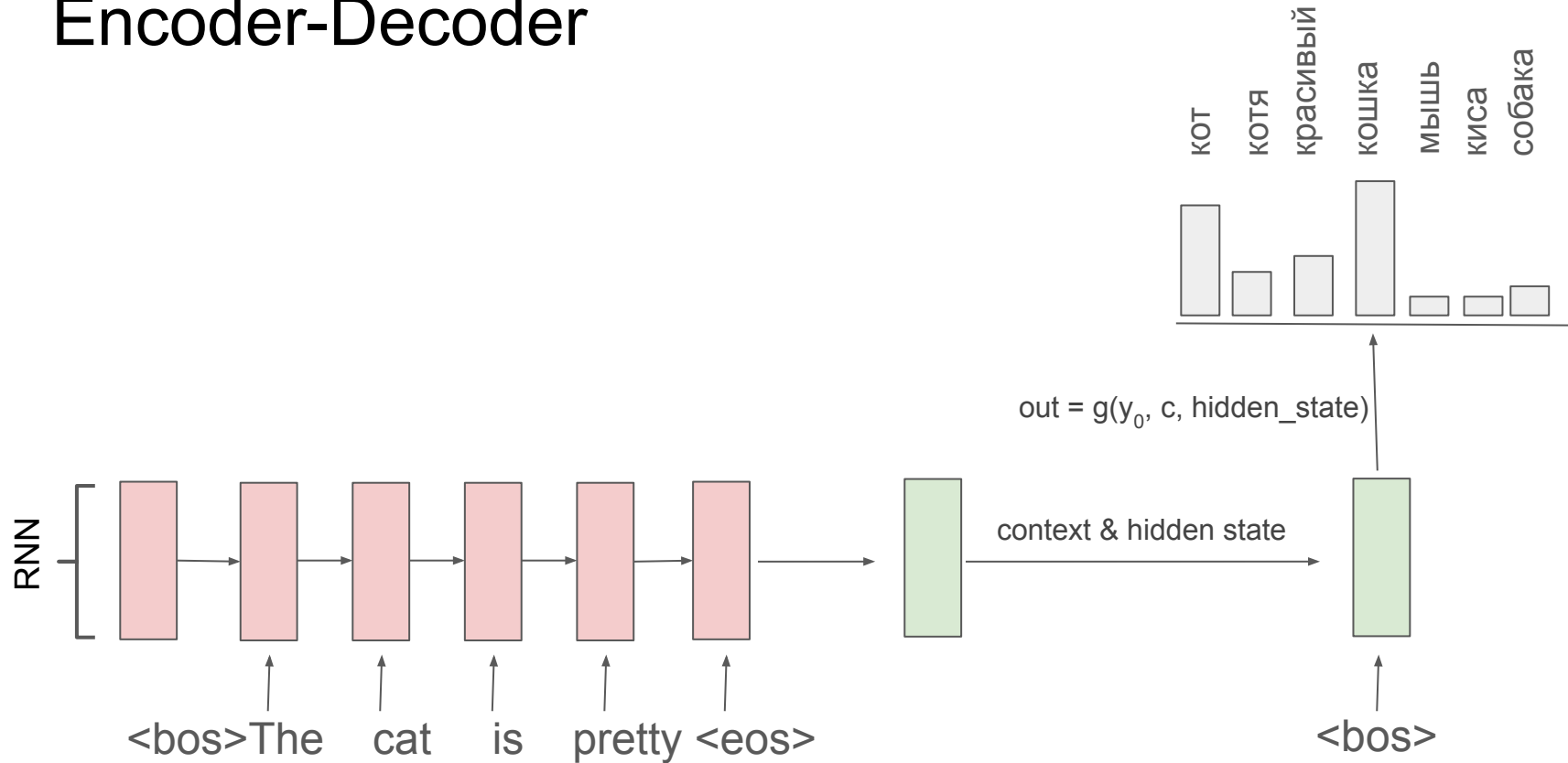
# Encoder-Decoder



# Encoder-Decoder

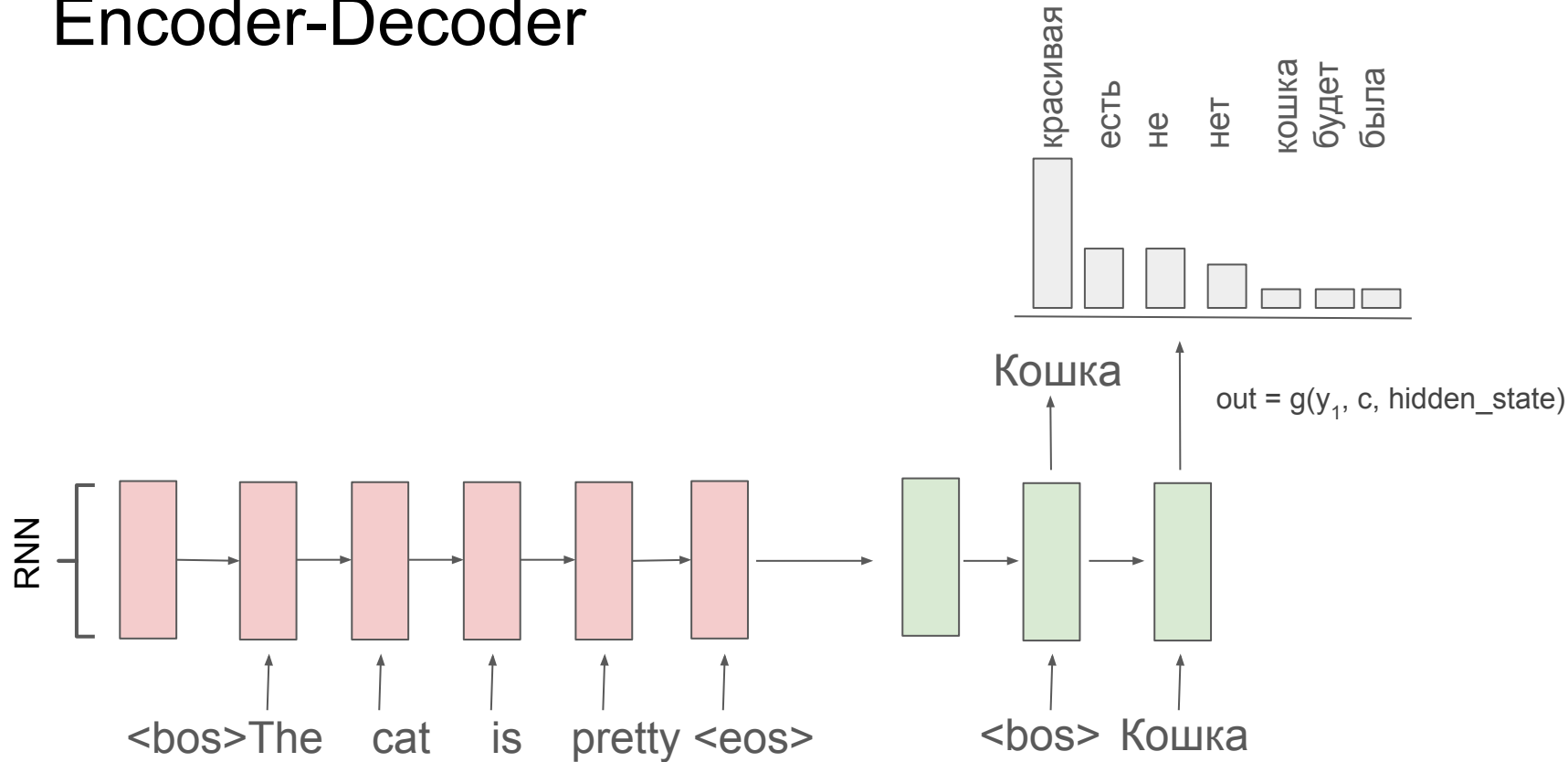


# Encoder-Decoder

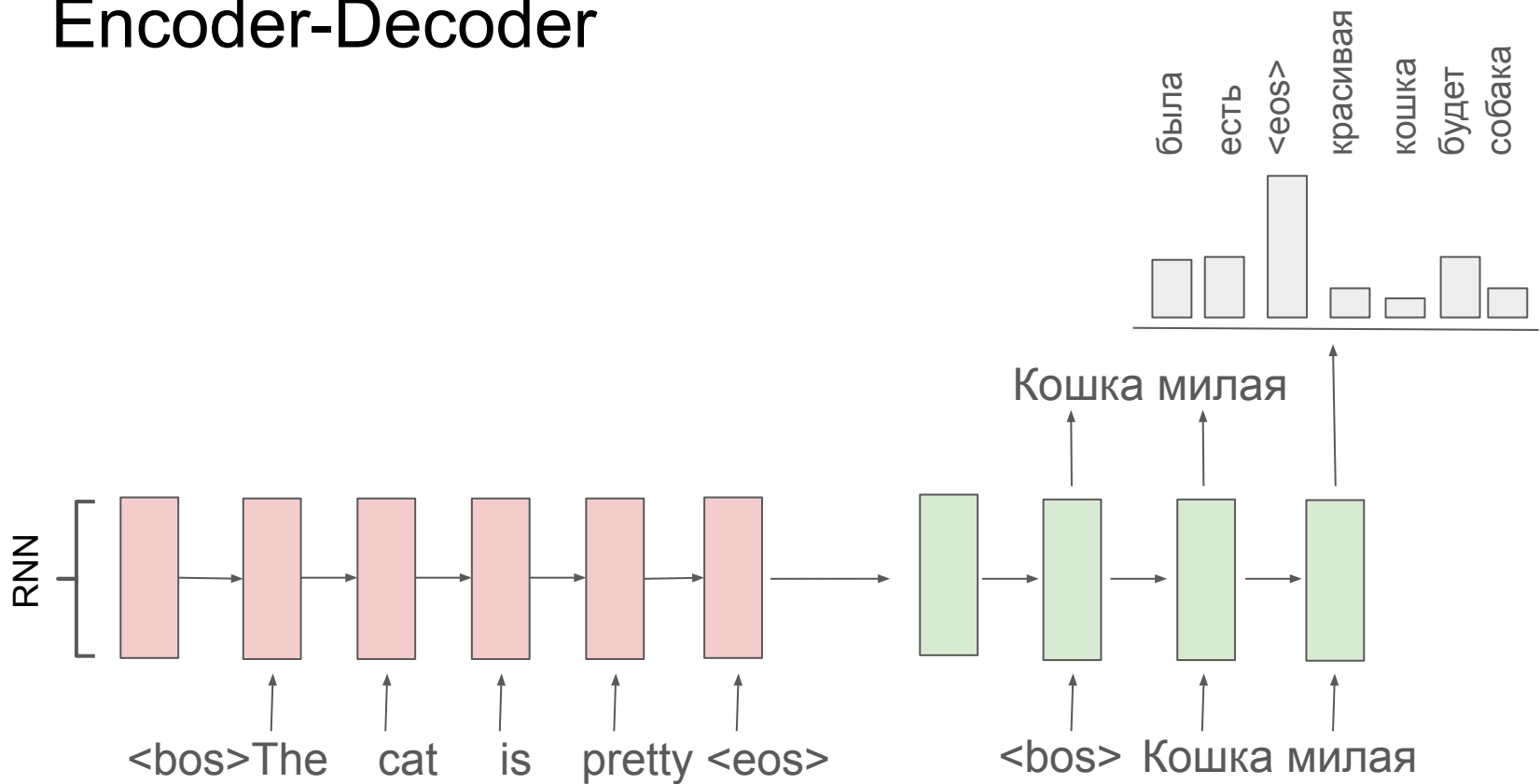




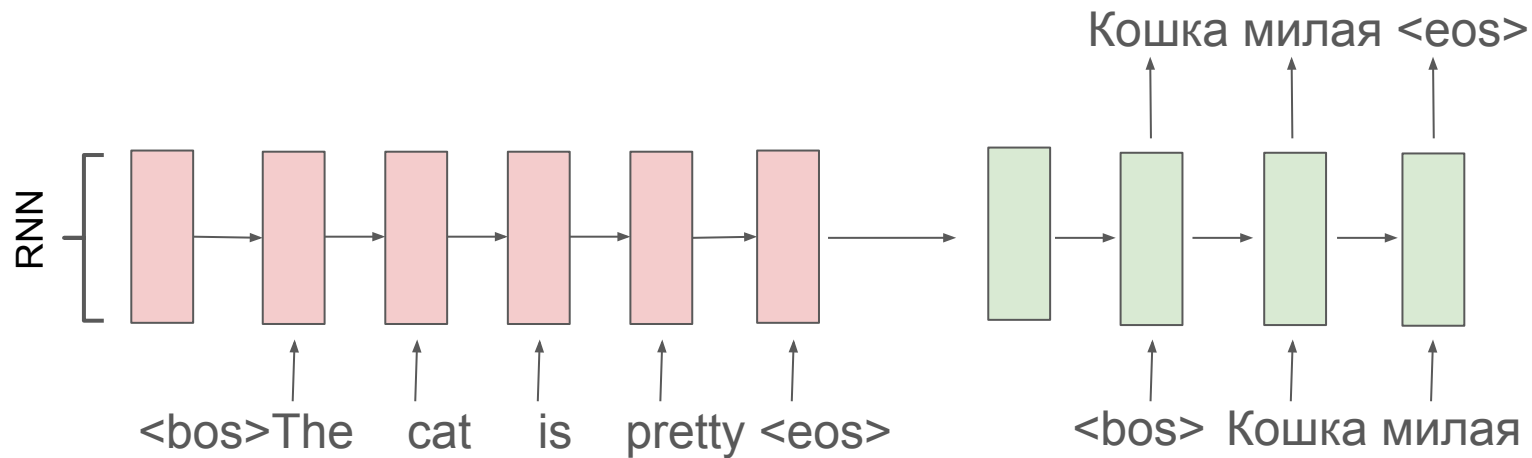
# Encoder-Decoder



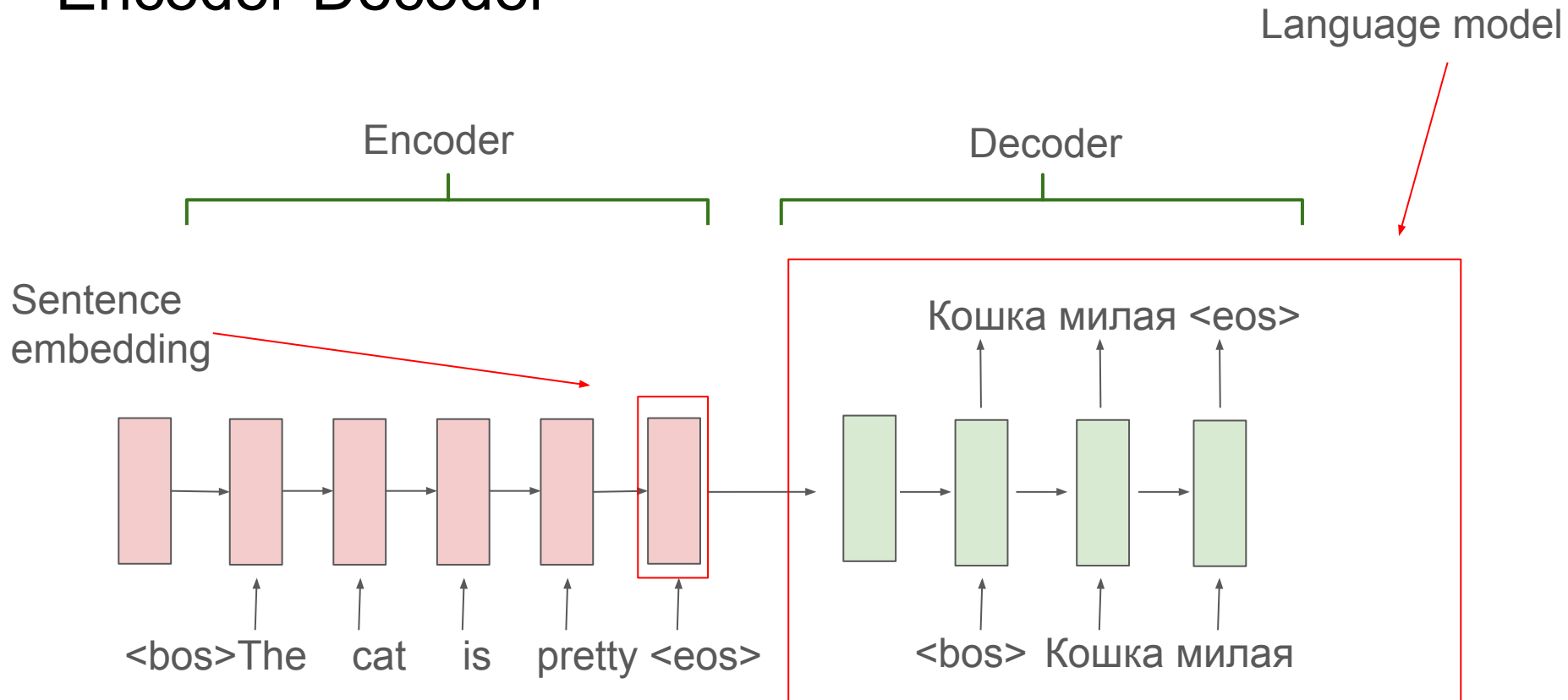
# Encoder-Decoder



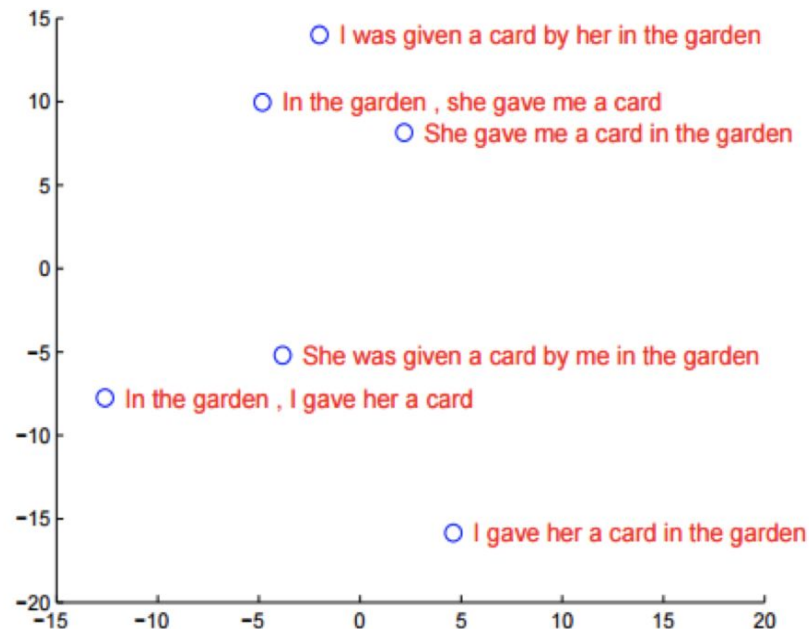
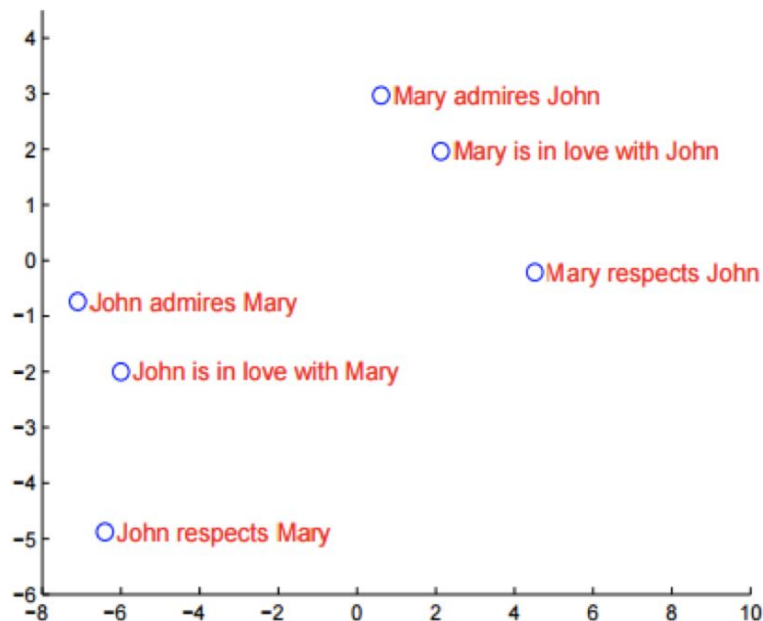
# Encoder-Decoder



# Encoder-Decoder



# Sentence embedding



(Sutskever et al., 2014)

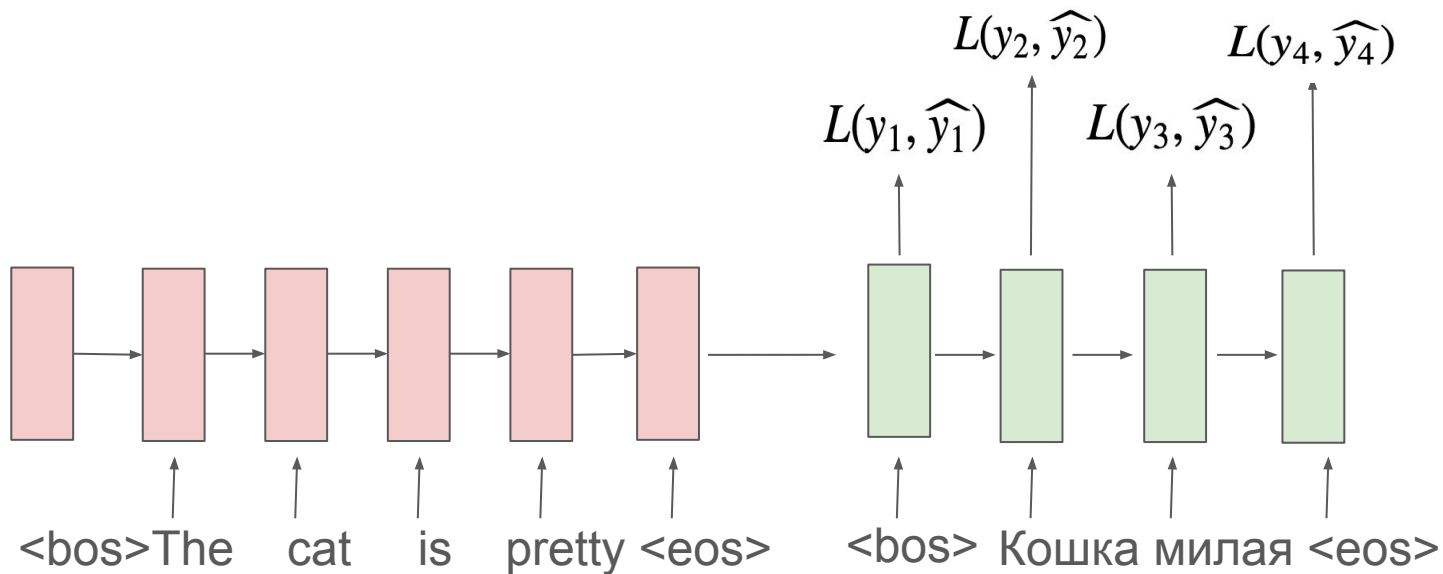
# Encoder-Decoder

- consists of two parts: encoder and decoder
- encoder gathers information about source sentence

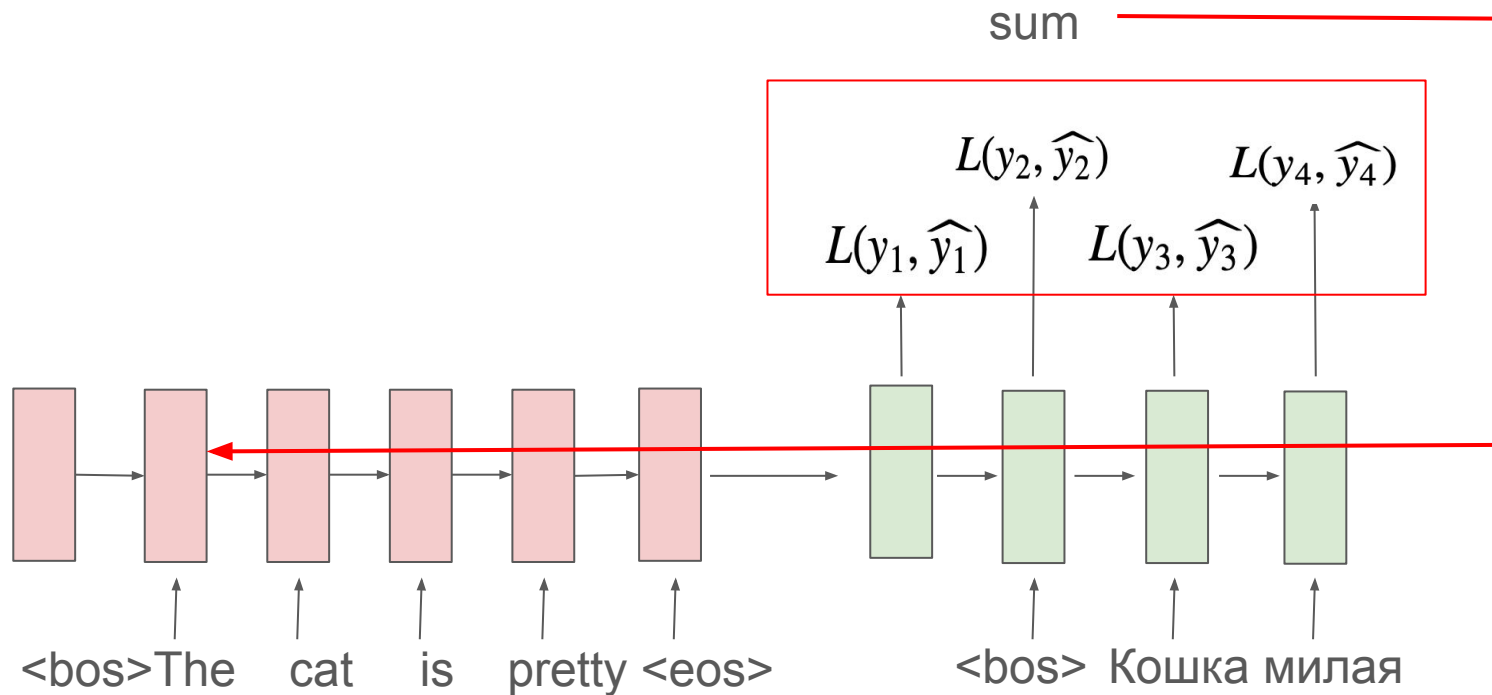
(may be not only RNN, but anything)

- decoder uses this information to generate new sequence

# Encoder-Decoder training



# Encoder-Decoder training





# Encoder-Decoder

- consists of two parts: encoder and decoder
- encoder gathers information about source sentence

(may be not only RNN, but anything)

- decoder uses this information to generate new sequence

Any improvement ideas?

Any problems seen?

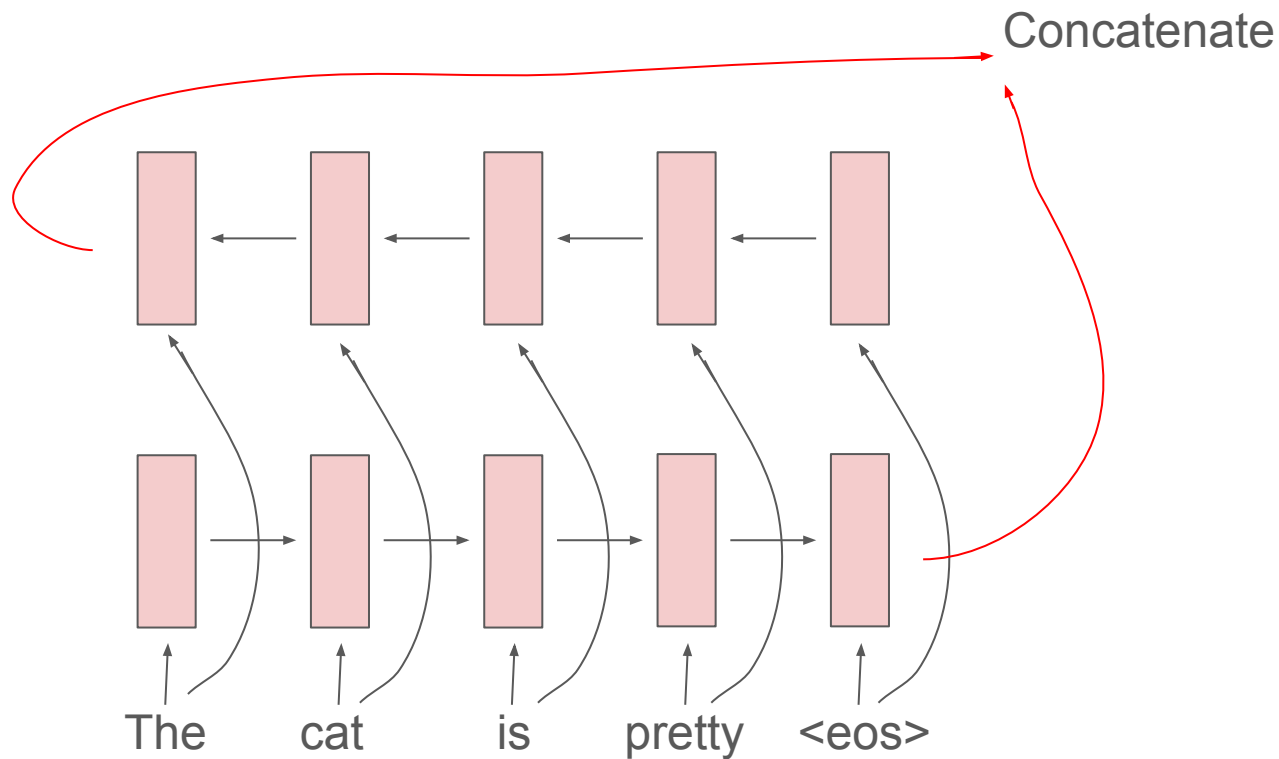
# Encoder-Decoder: problems

- Encoder can forget beginning of long sentences
- Greedy decoding
- Teacher-forcing
- All the complicated human language is put into single vector!

# Encoder-Decoder: problems

- Encoder can forget beginning of long sentences
- Greedy decoding
- Teacher-forcing
- All the complicated human language is put into single vector!

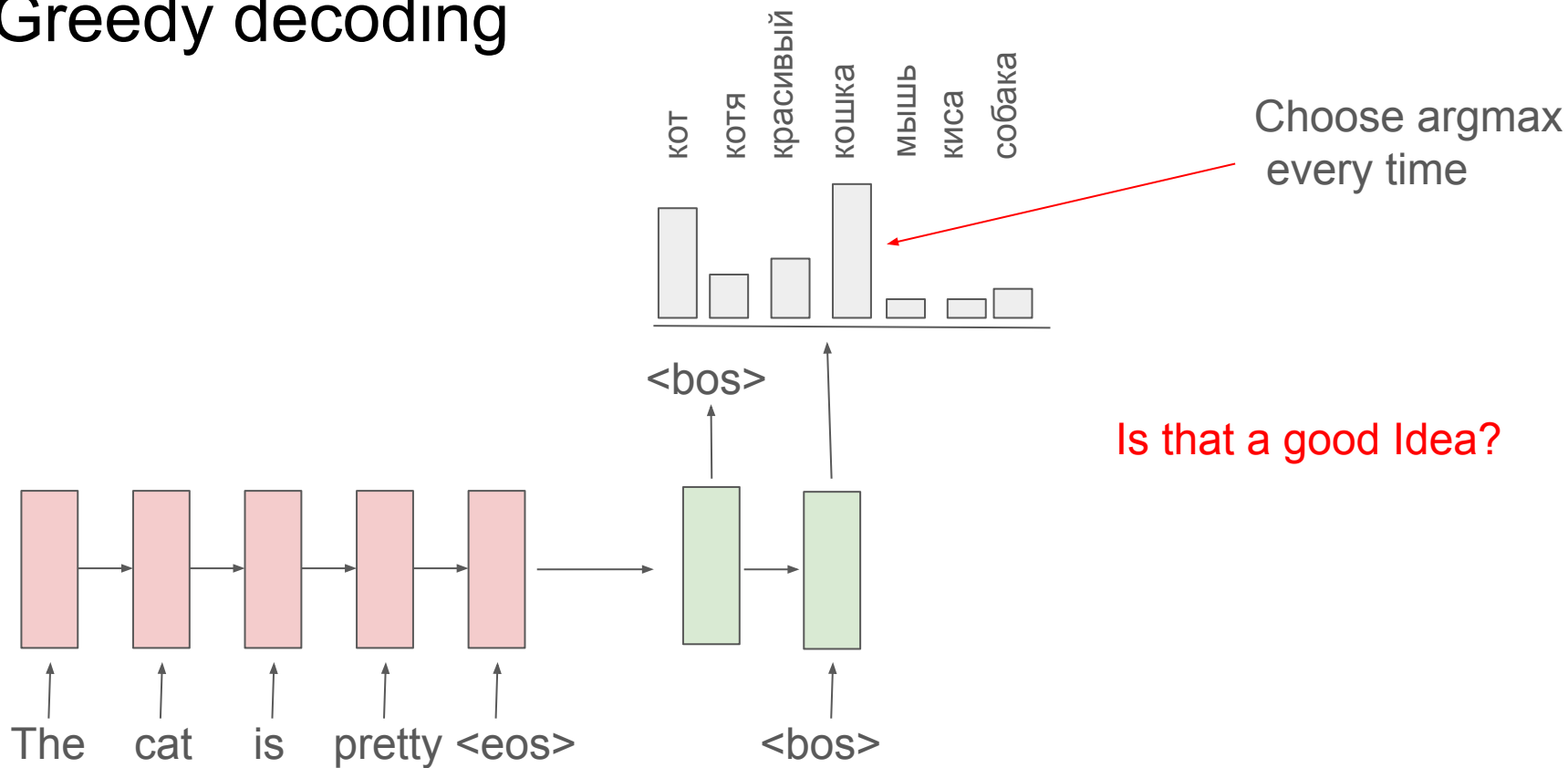
# Bidirectional RNN



# Encoder-Decoder: problems

- Encoder can forget beginning of long sentences
- Greedy decoding
- Teacher-forcing
- All the complicated human language is put into single vector!

# Greedy decoding



# Greedy decoding

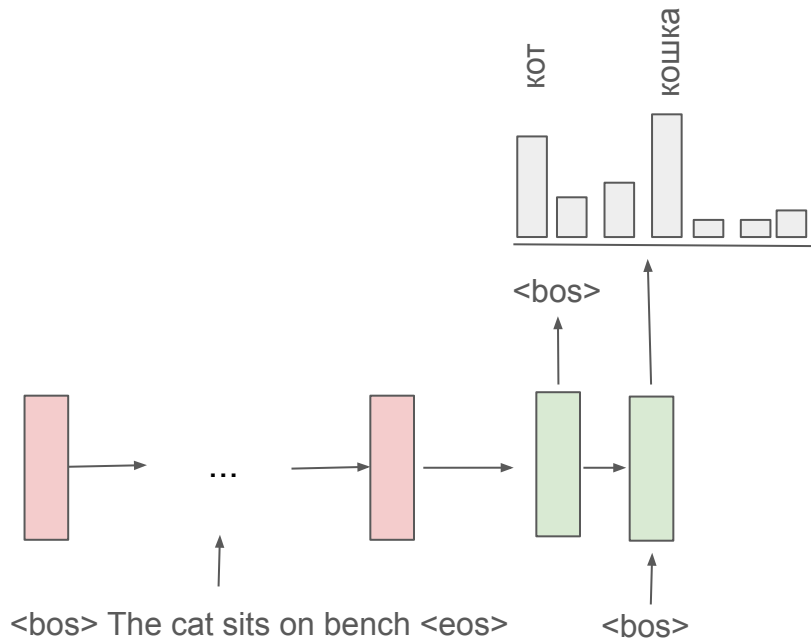
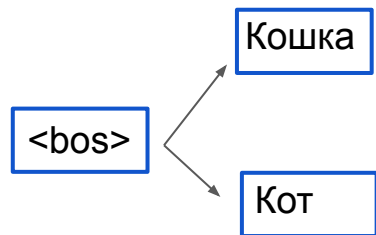
Had there been a vaccine, would he be alive.

Была тут вакцина ...

Если бы там была вакцина ...

# Beam search

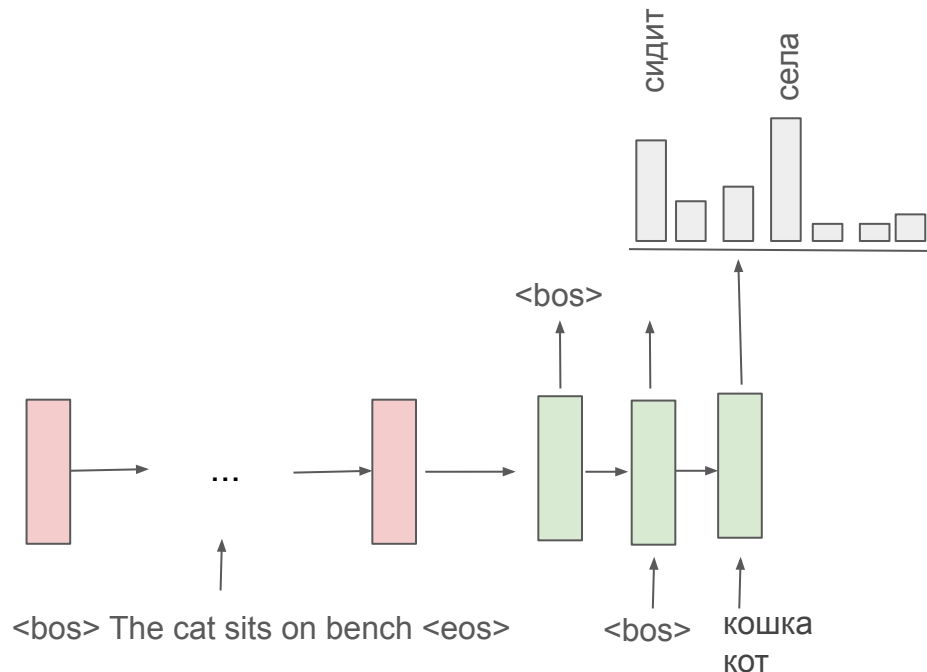
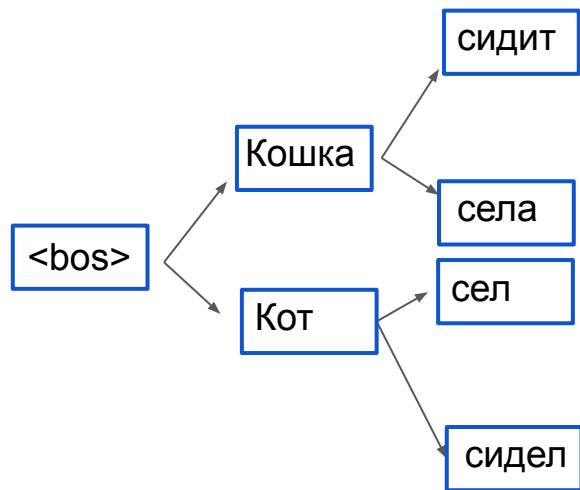
Maintain fixed number of hypotheses, extend them and choose most probable ones





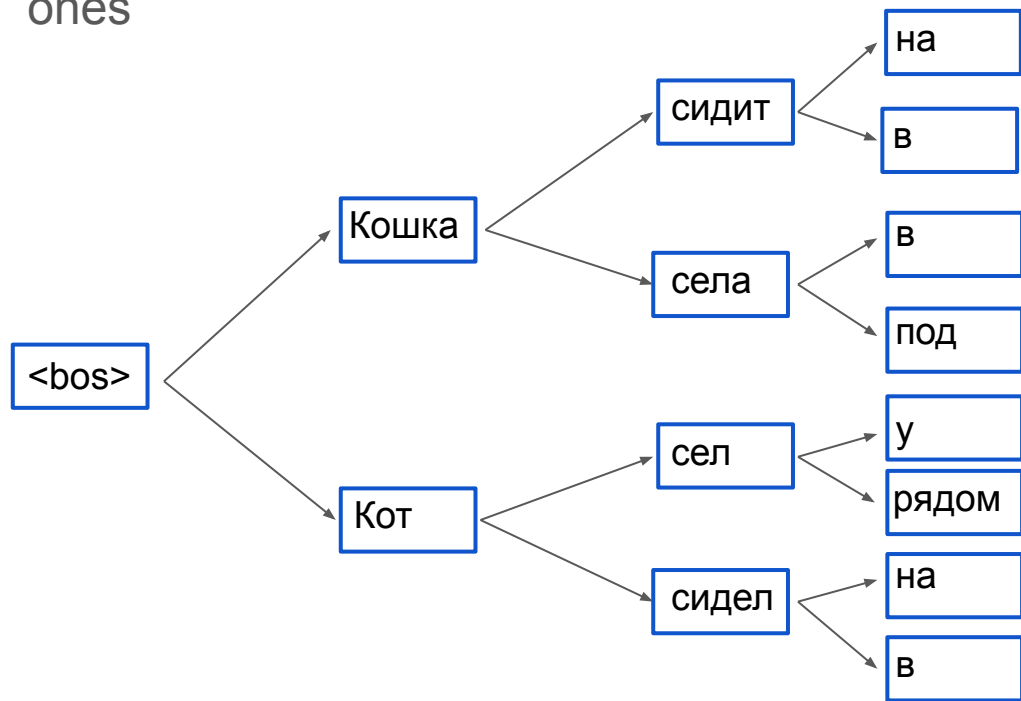
# Beam search

Maintain fixed number of hypotheses, extend them and choose most probable ones



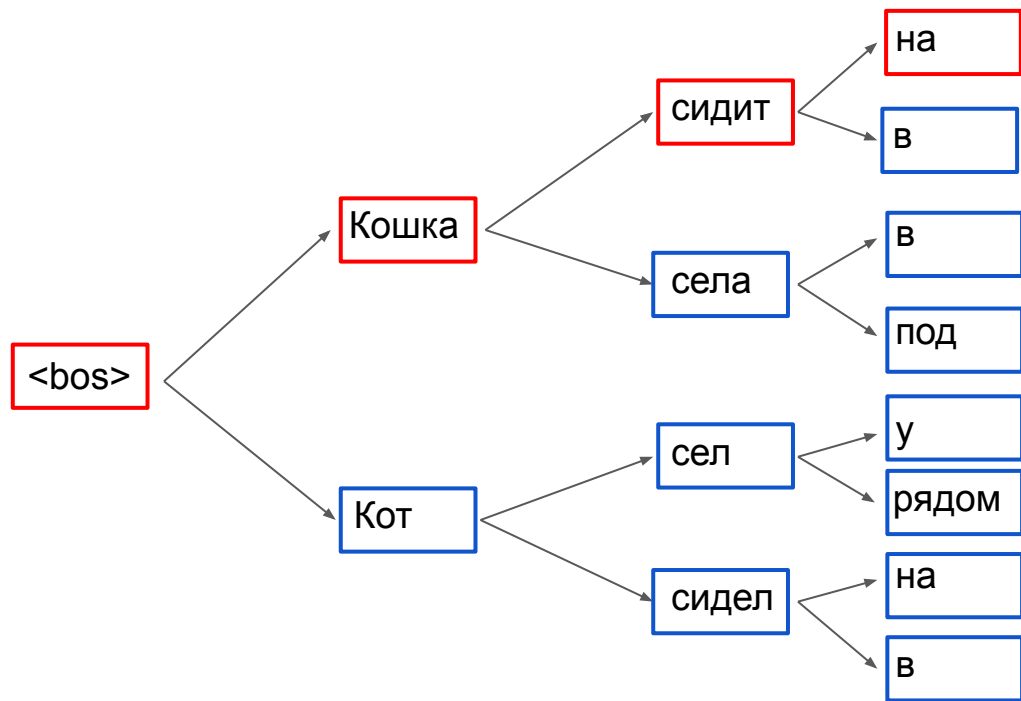
# Beam search

Maintain fixed number of hypotheses, extend them and choose most probable ones



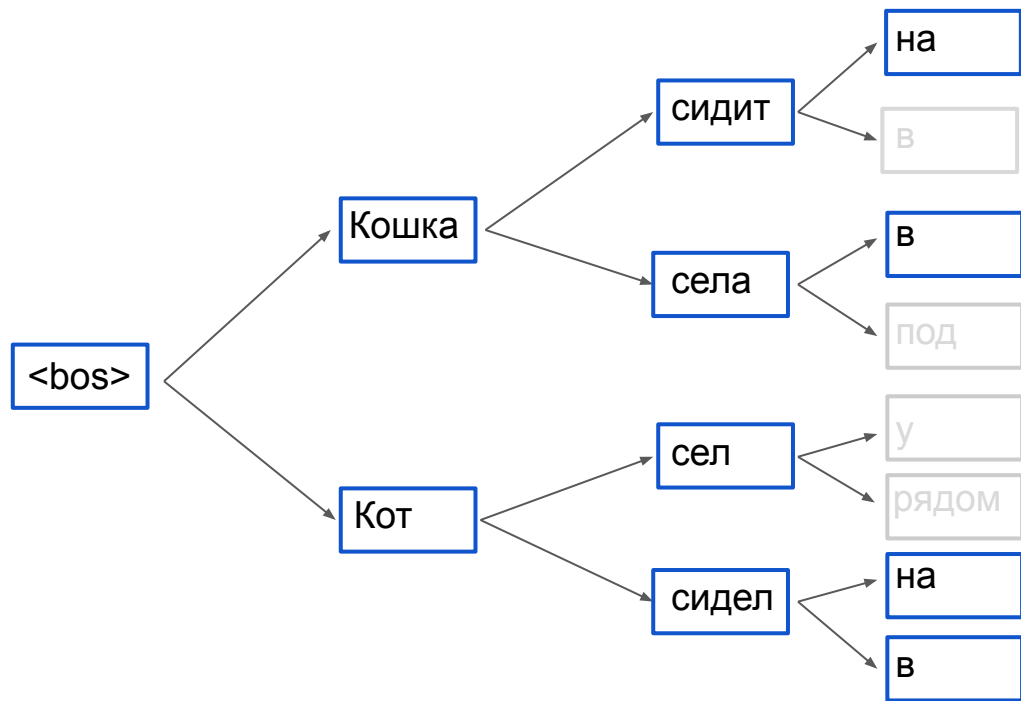
# Beam search

$$P(\text{Кошка сидит на}) = P(\text{на} \mid \text{кошка сидит}) * P(\text{сидит} \mid \text{кошка}) * P(\text{кошка})$$



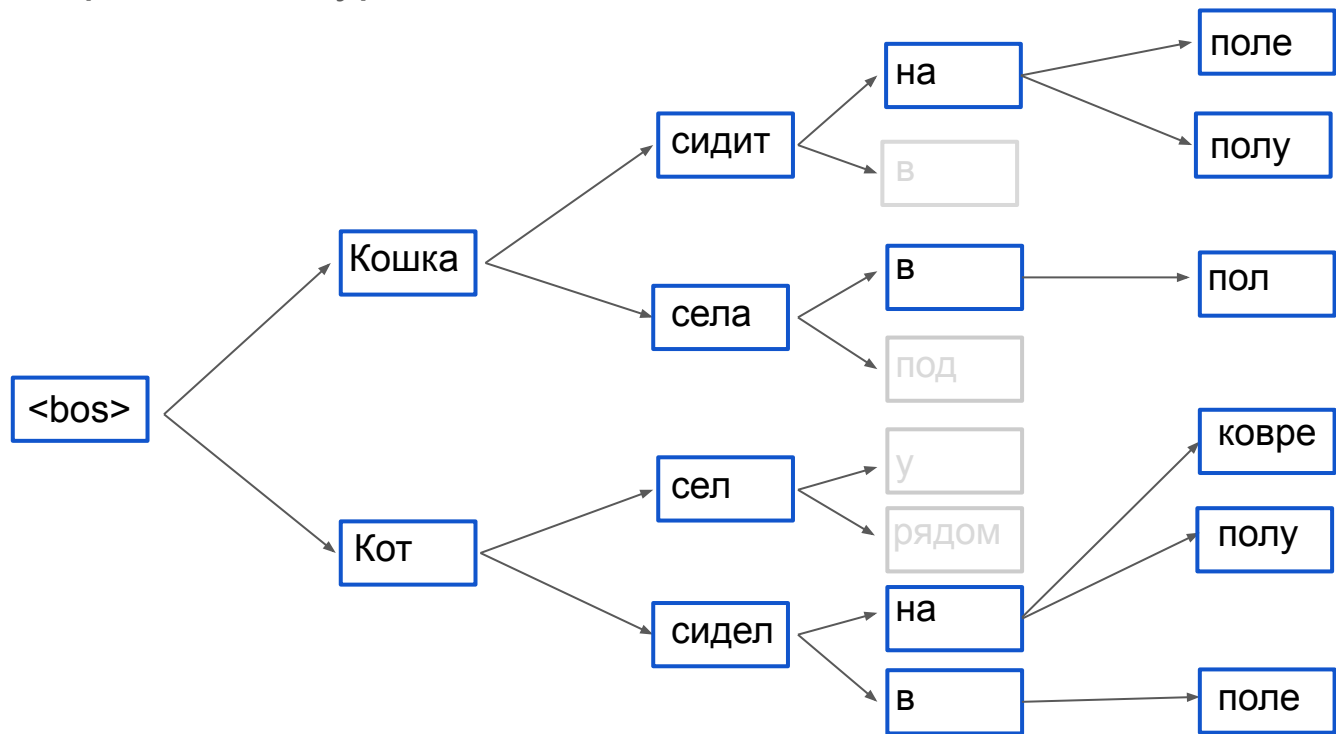
# Beam search

Leave only fixed number of hypotheses with best scores

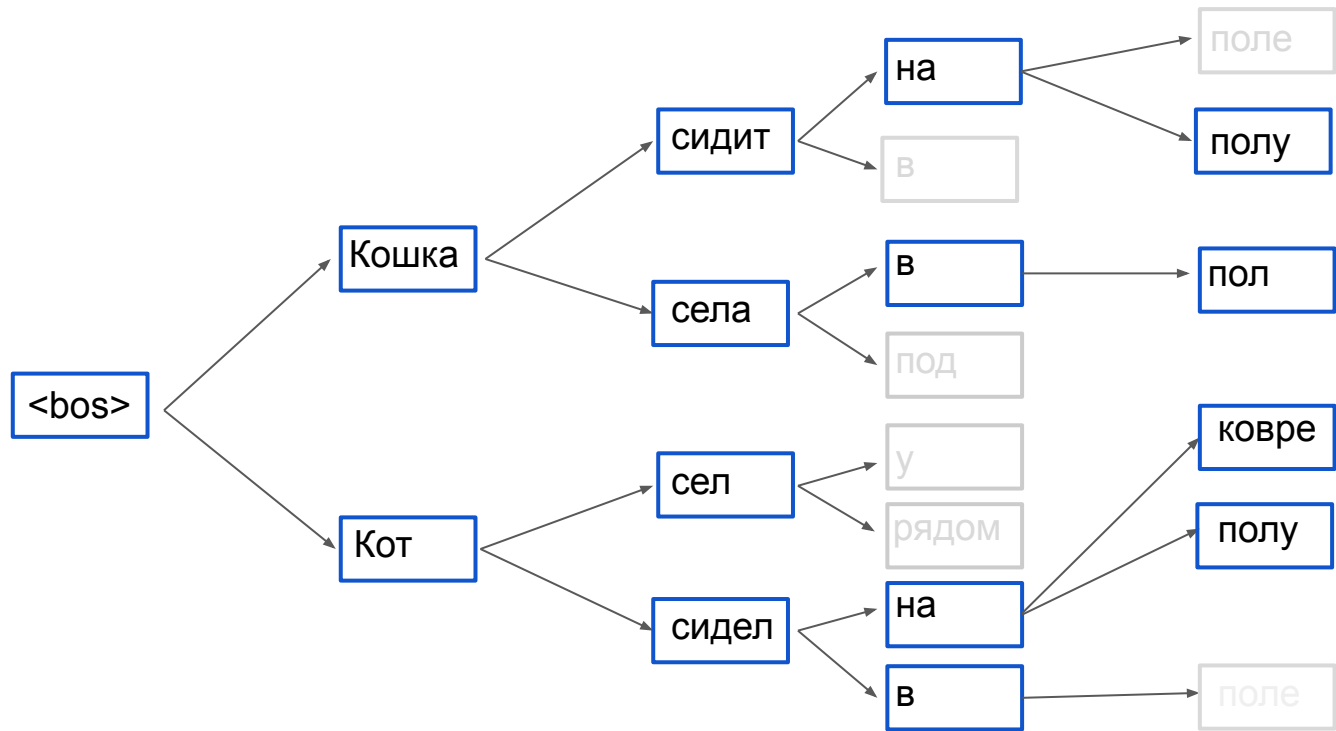


# Beam search

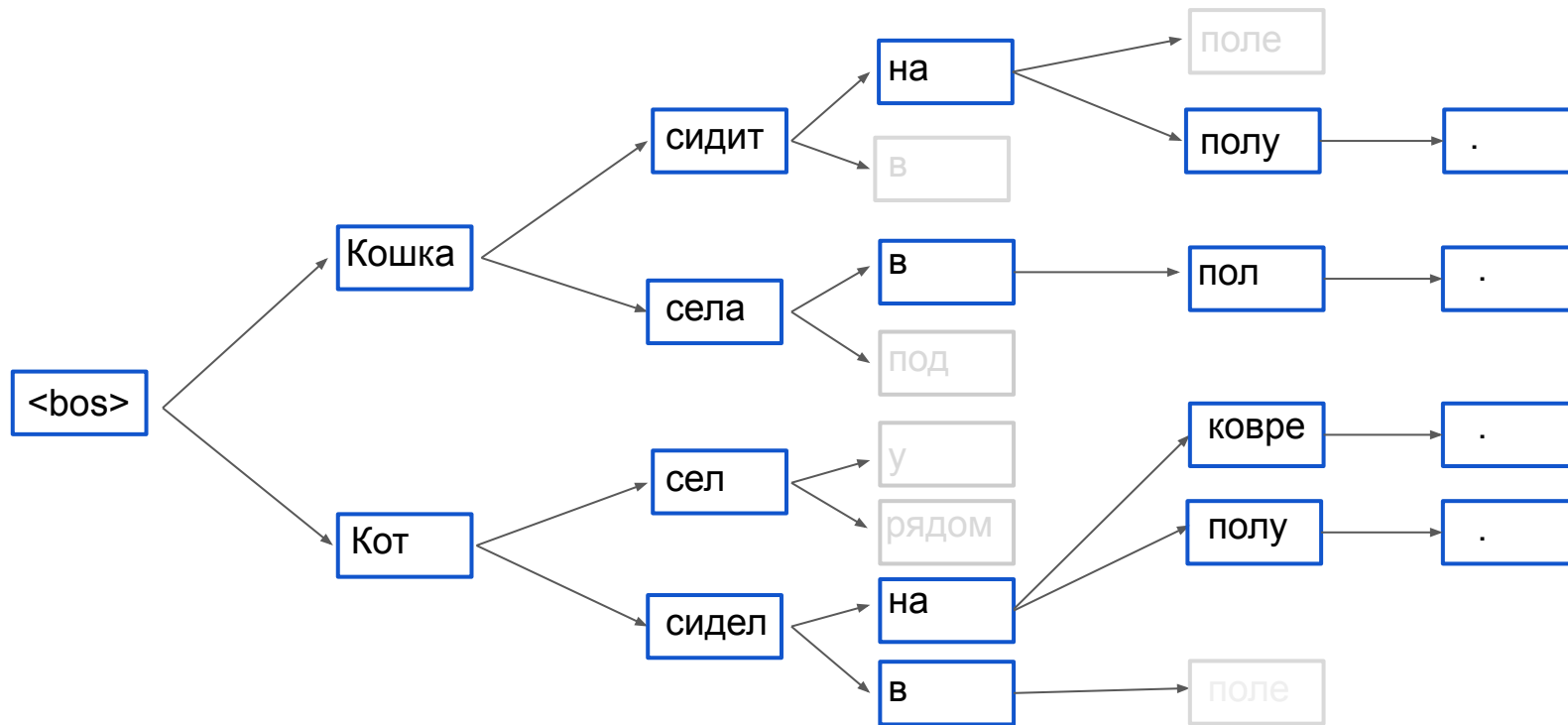
Expand best hypotheses further



# Beam search

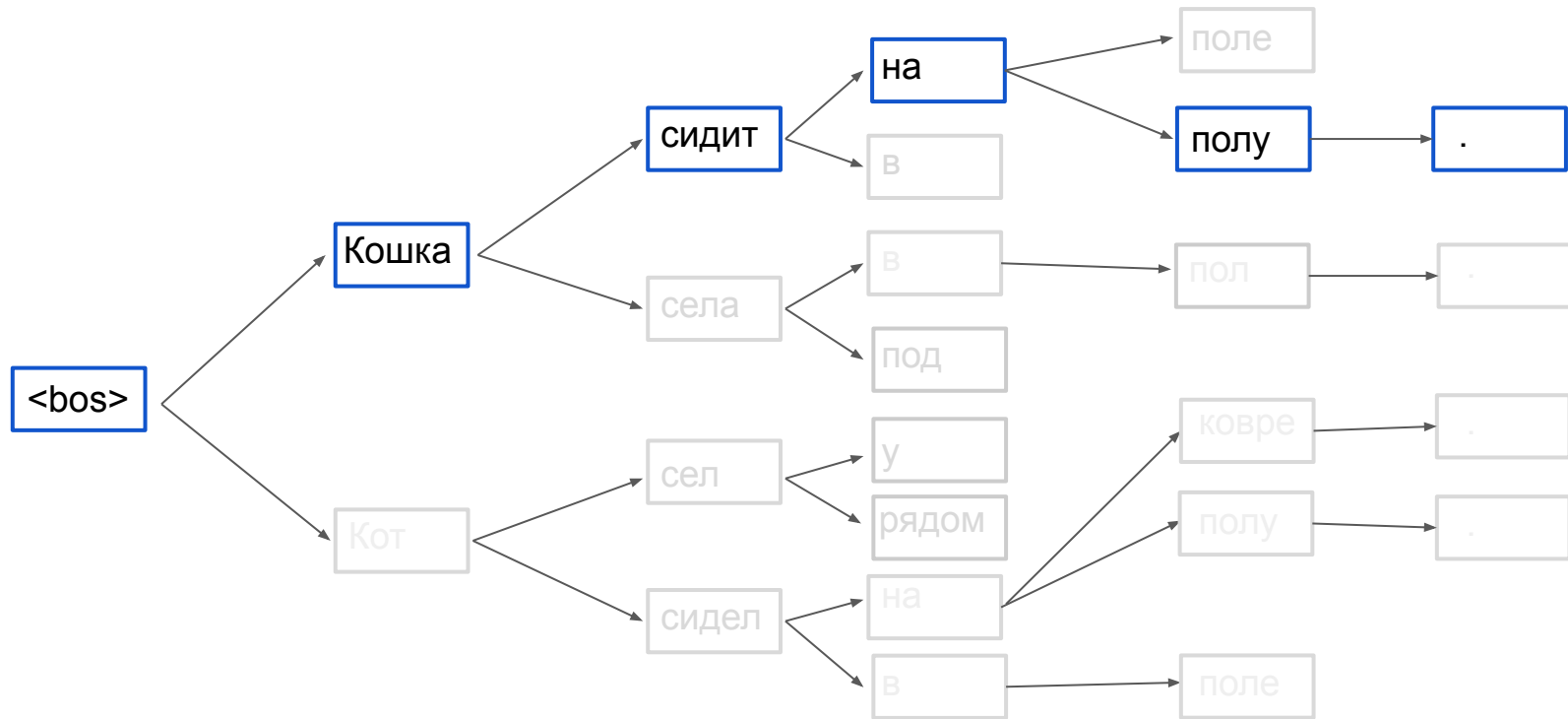


# Beam search



# Beam search

Choose the best overall hypothesis.





# Beam search

Maintain fixed number of hypotheses, extend them and choose most probable ones

Optimal beam size  $\sim 4-8$

If beam size is too big, the translation model behaves badly (quality decreases)

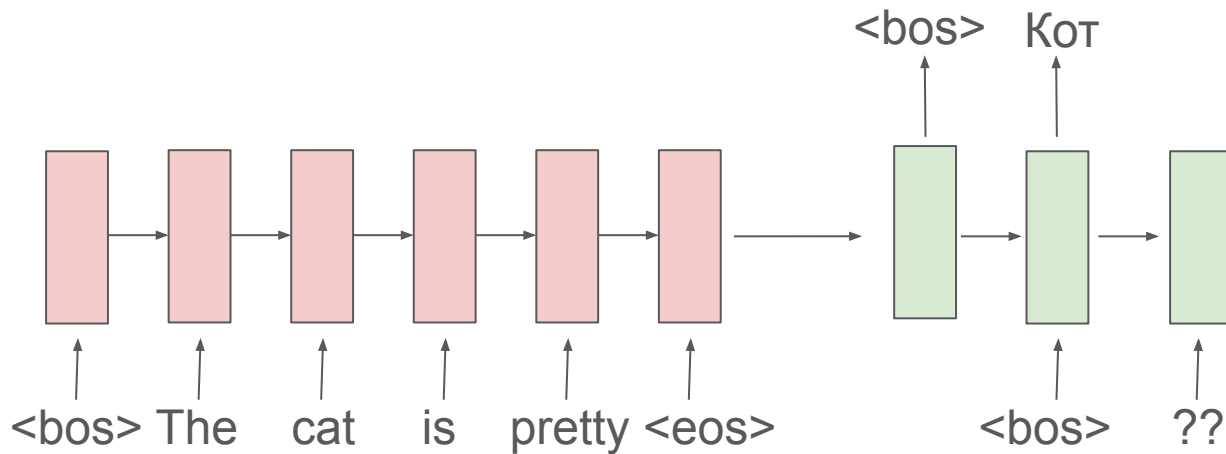
# Encoder-Decoder: problems

- Encoder can forget beginning of long sentences
- Greedy decoding
- **Teacher-forcing**
- All the complicated human language is put into single vector!

# Teacher-forcing

**Source:** The cat is pretty

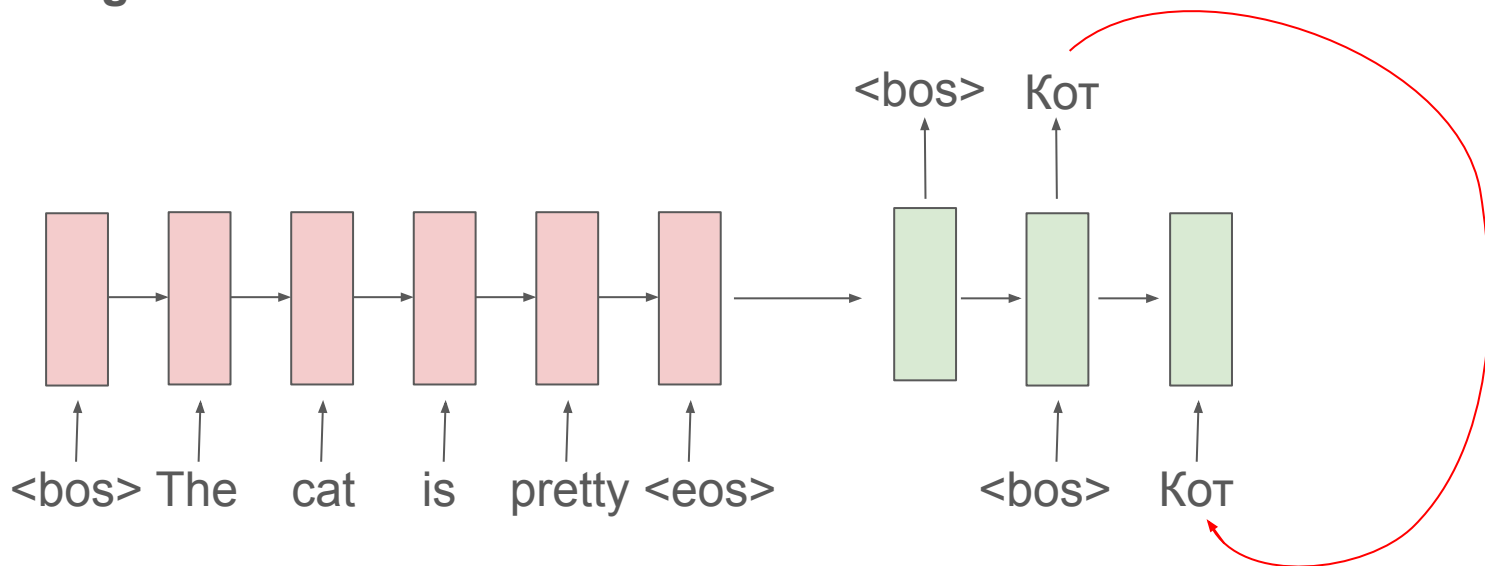
**Target:** Кошка милая



# Teacher-forcing

**Source:** The cat is pretty

**Target:** Кошка милая

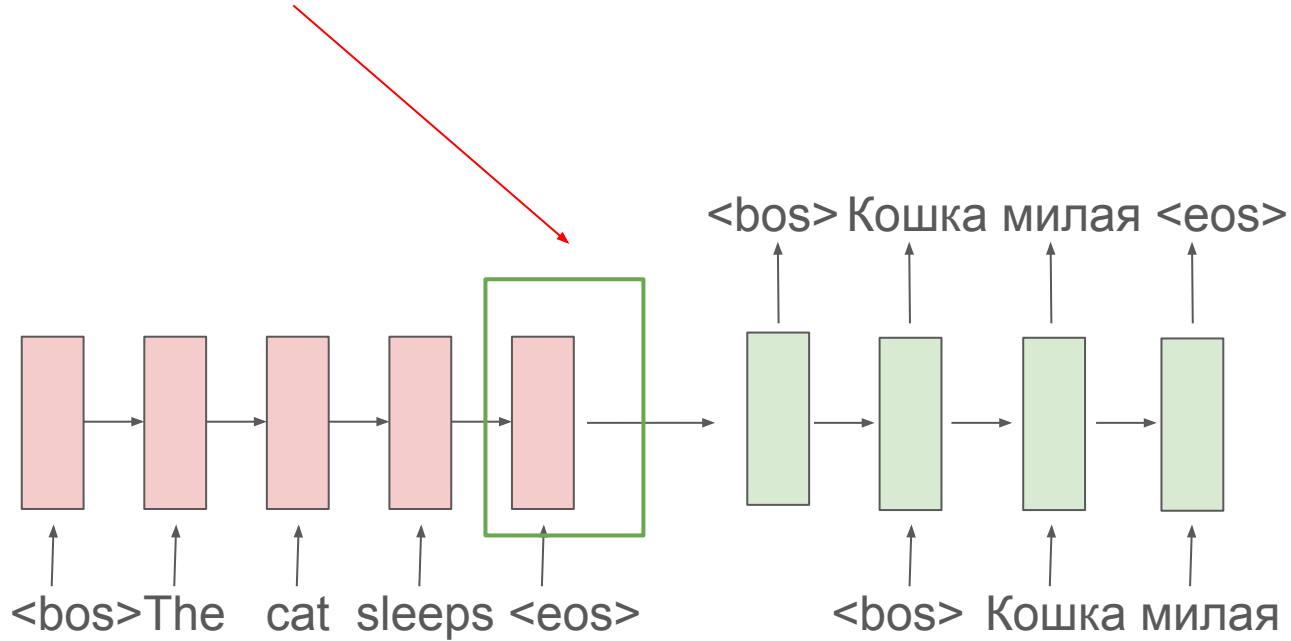


# Encoder-Decoder: problems

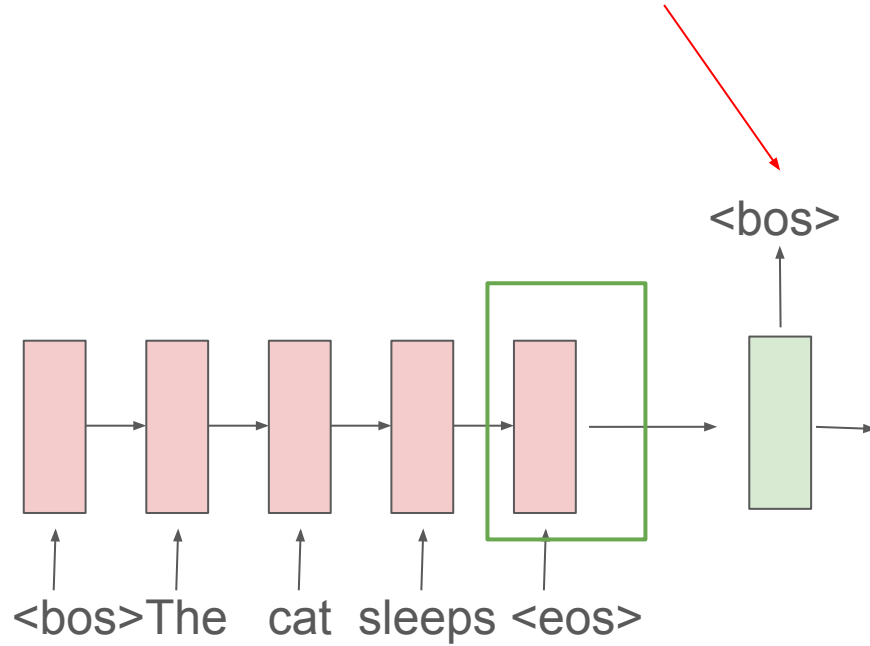
- Encoder can forget beginning of long sentences
- Greedy decoding
- Teacher-forcing
- All the complicated human language is put into single vector!

# Attention

A poor single vector responsible for whole sentence

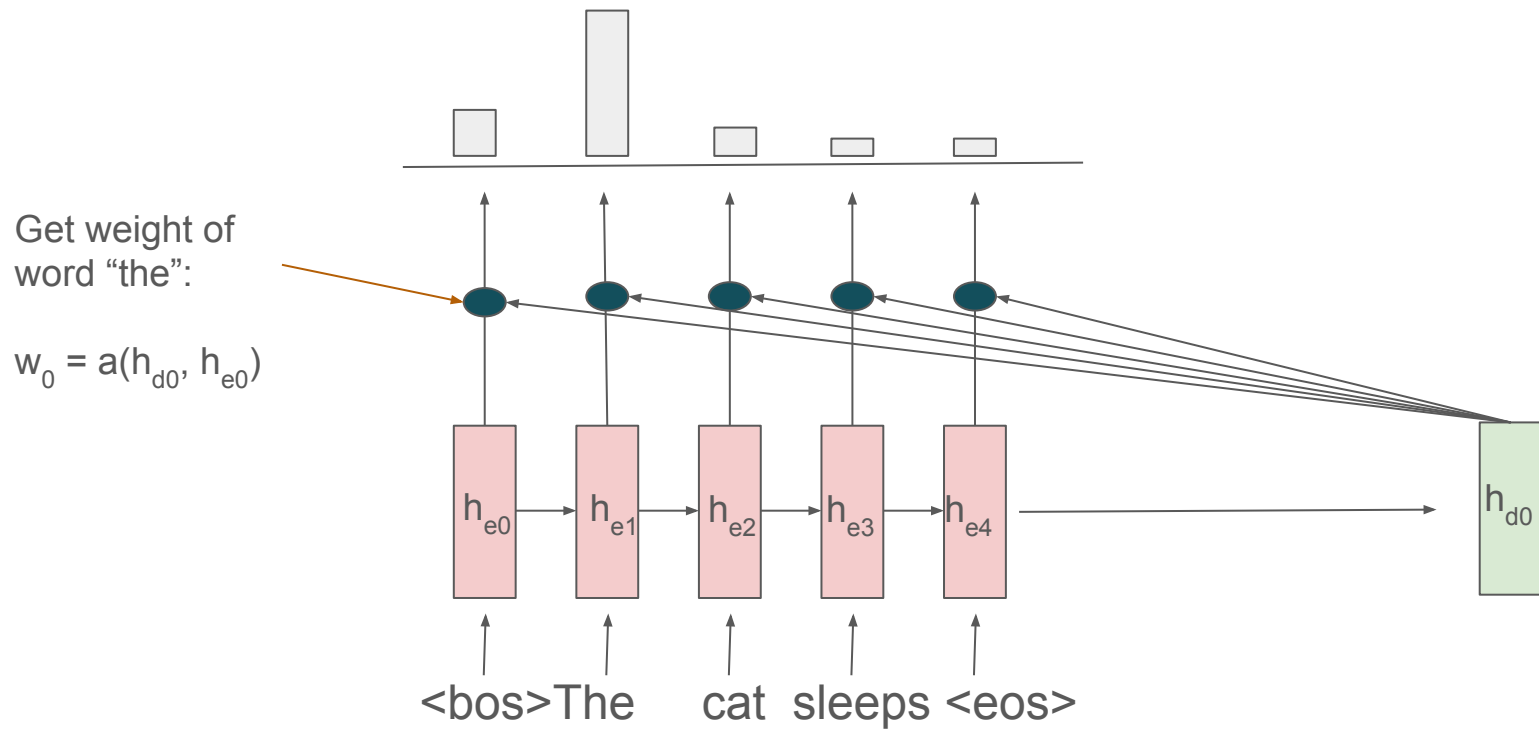


Get prediction based **only** on a hidden state of RNN & context (embedding)

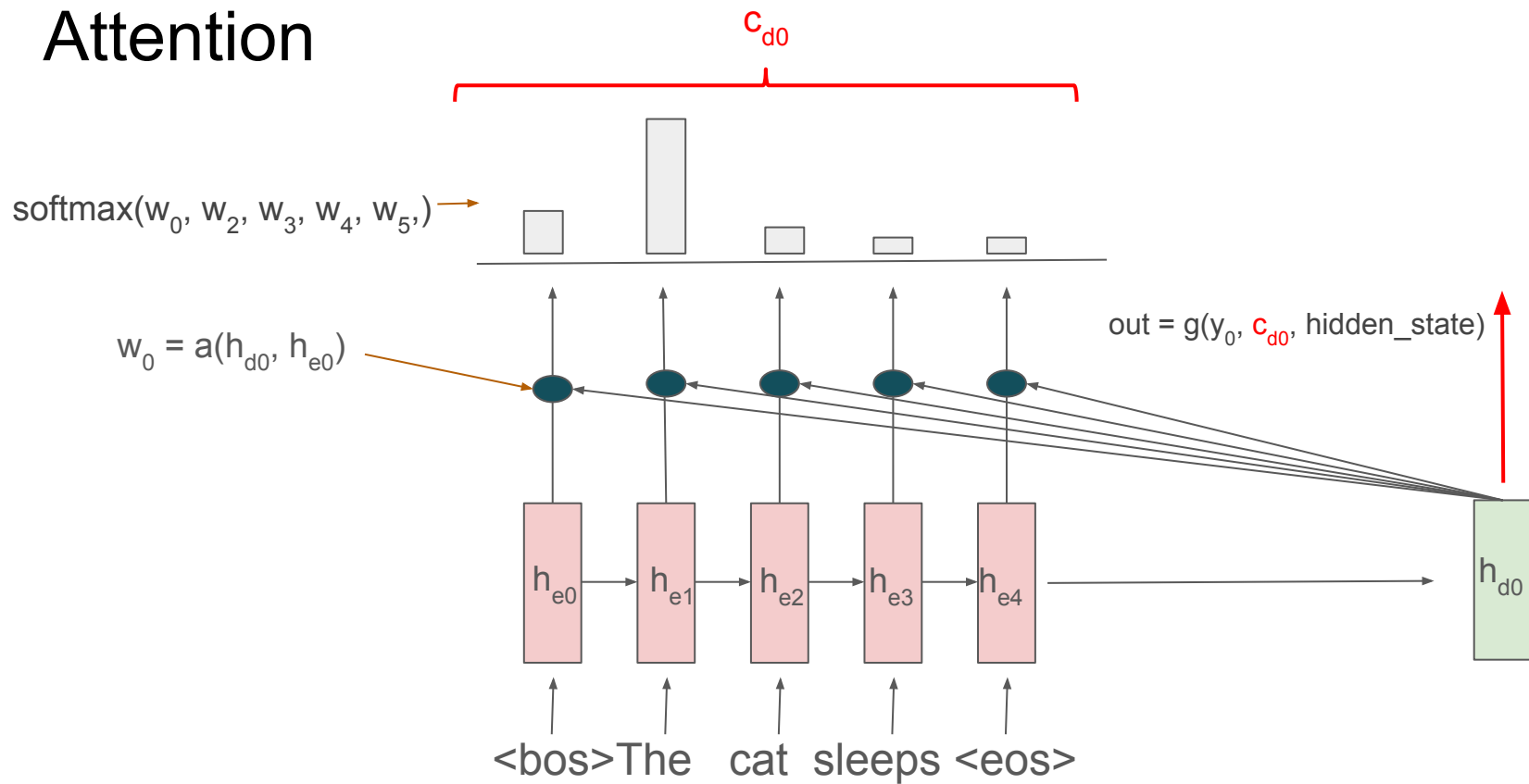




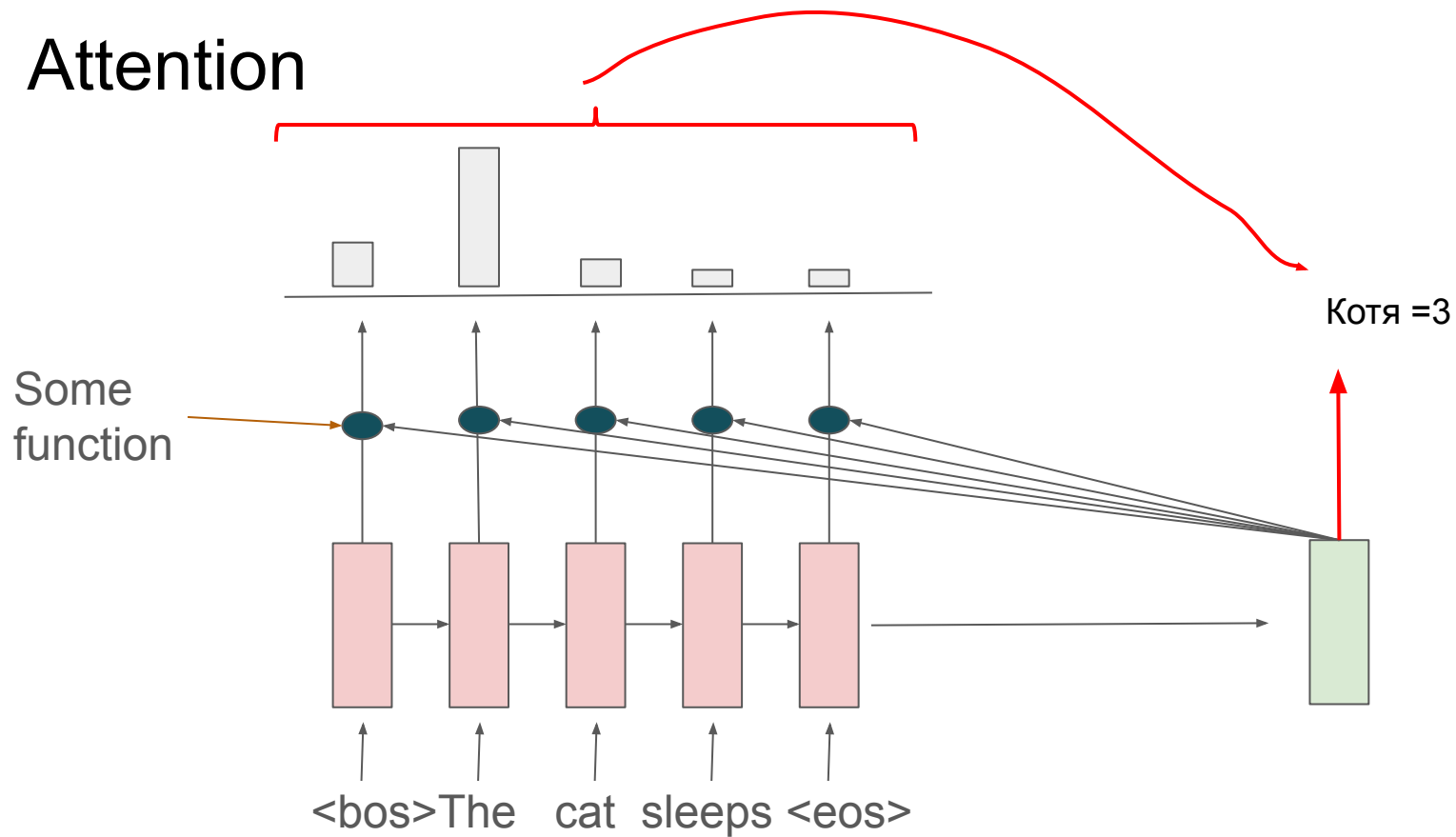
# Attention



# Attention



# Attention



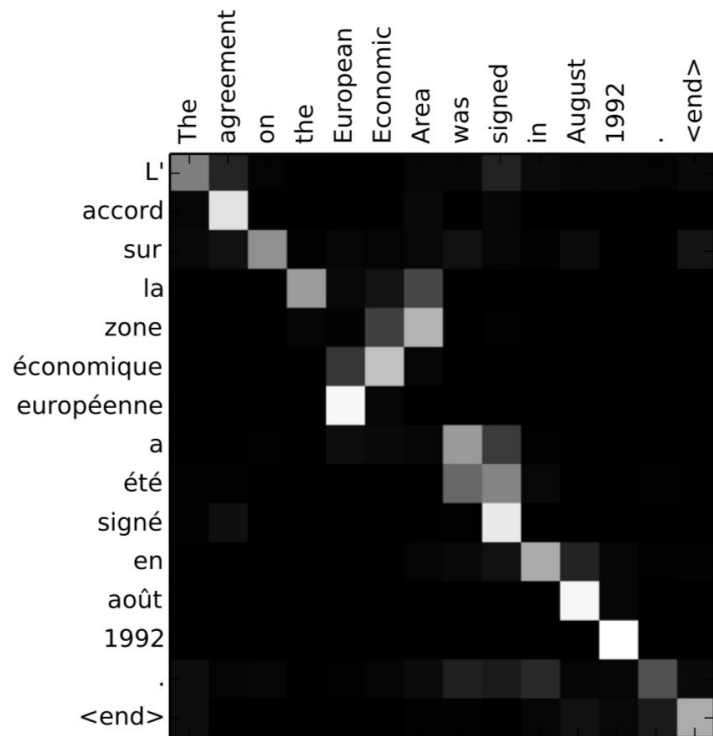
# Attention

Advantages:

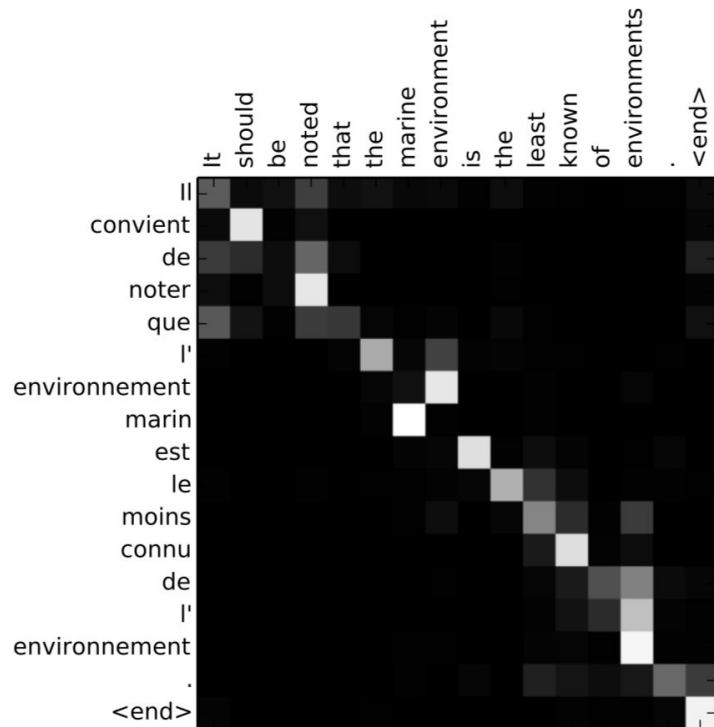
- behaves like a human
- we don't need to remember whole sentence

(don't need to put the whole sentence to one fixed-size vector)

# Attention



(a)



(b)

# Attention

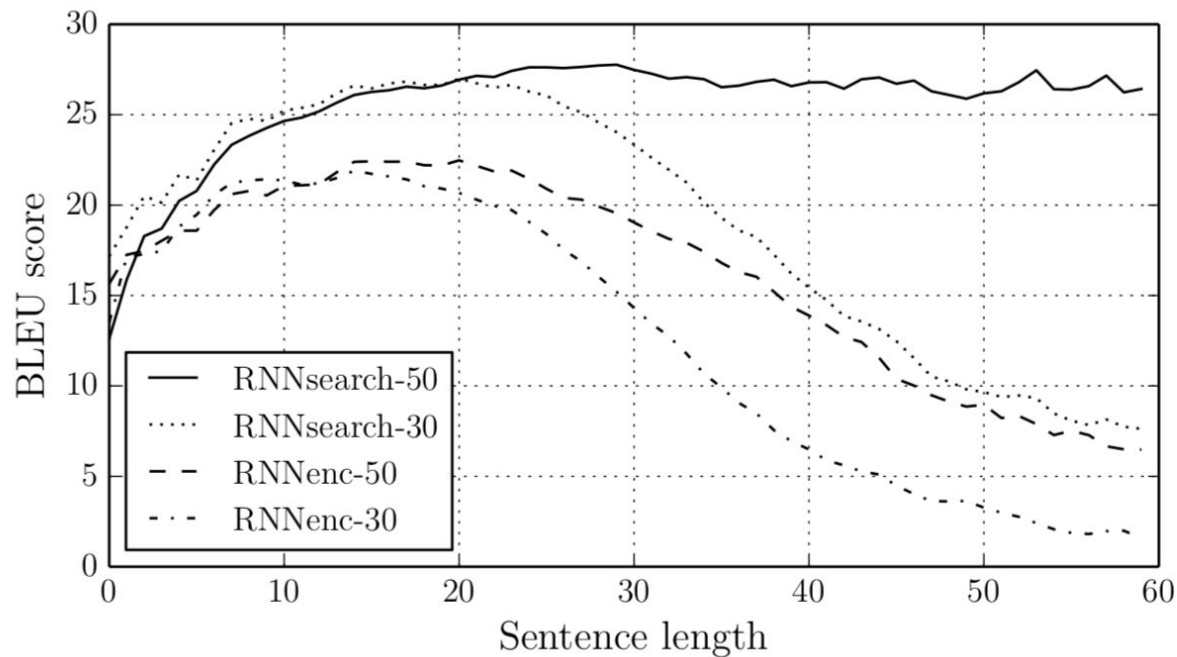
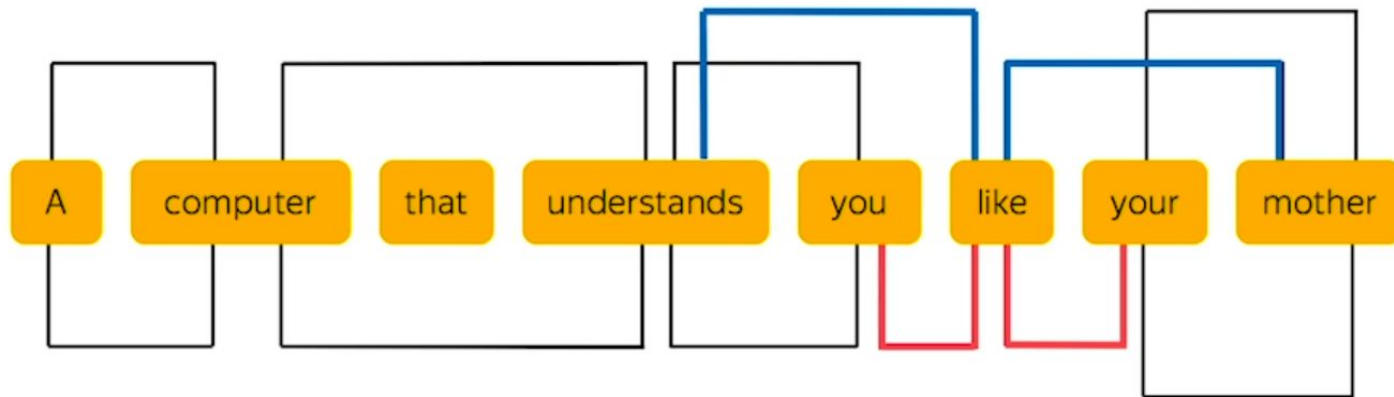


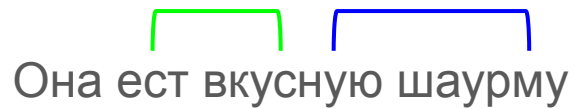
Figure 2: The BLEU scores of the generated translations on the test set with respect to the lengths of the sentences. The results are on the full test set which includes sentences having unknown words to the models.

# Self-attention



# Self-attention

Она ест вкусную шаурму

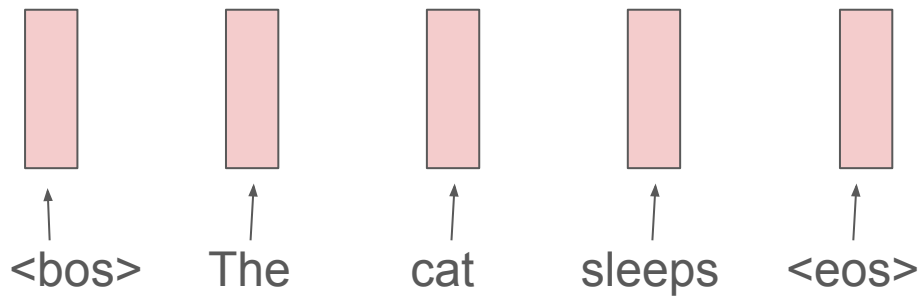


- Case agreement
- Gender agreement

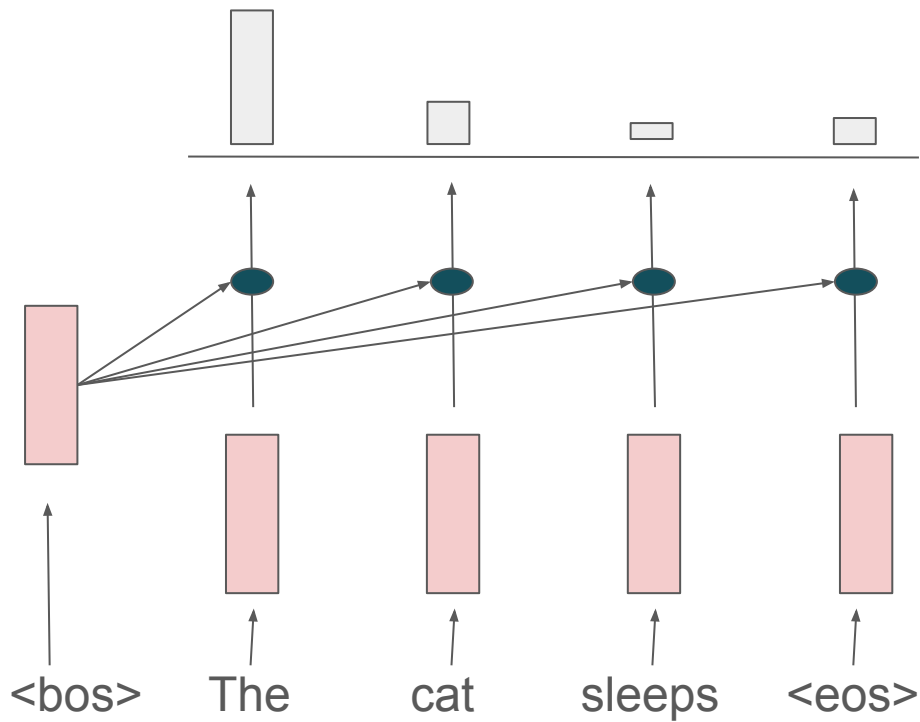


# Self-attention

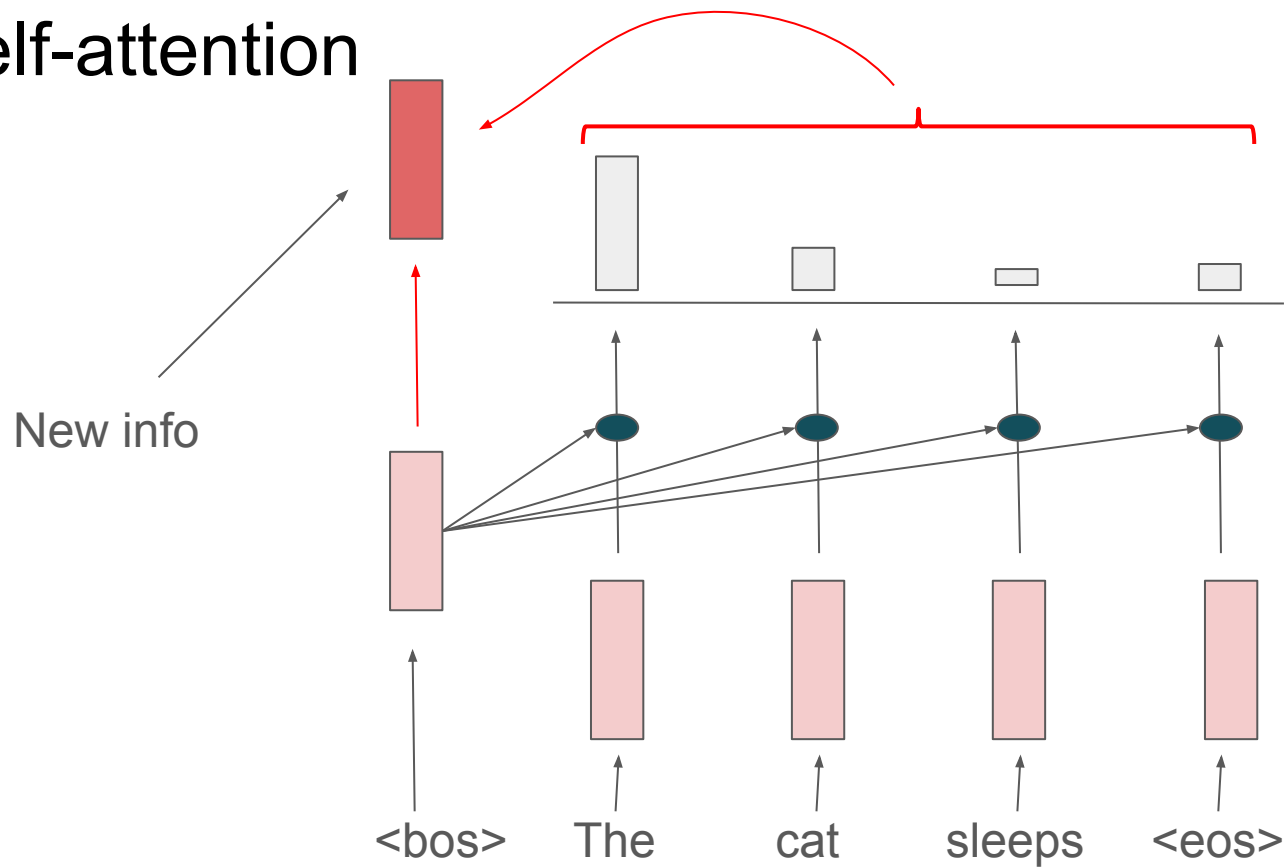
Embeddings



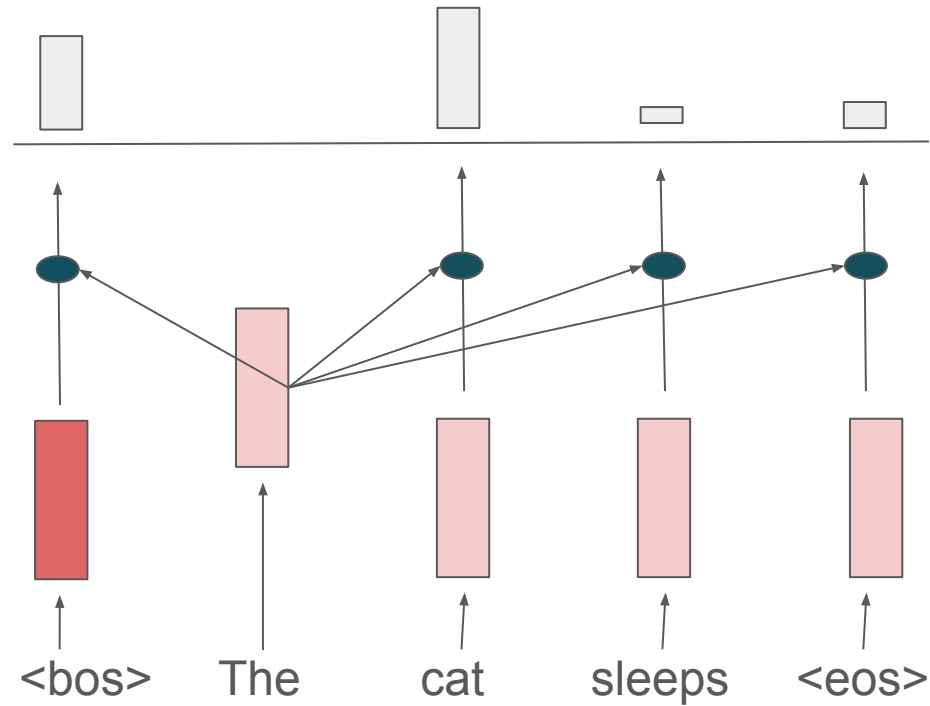
# Self-attention



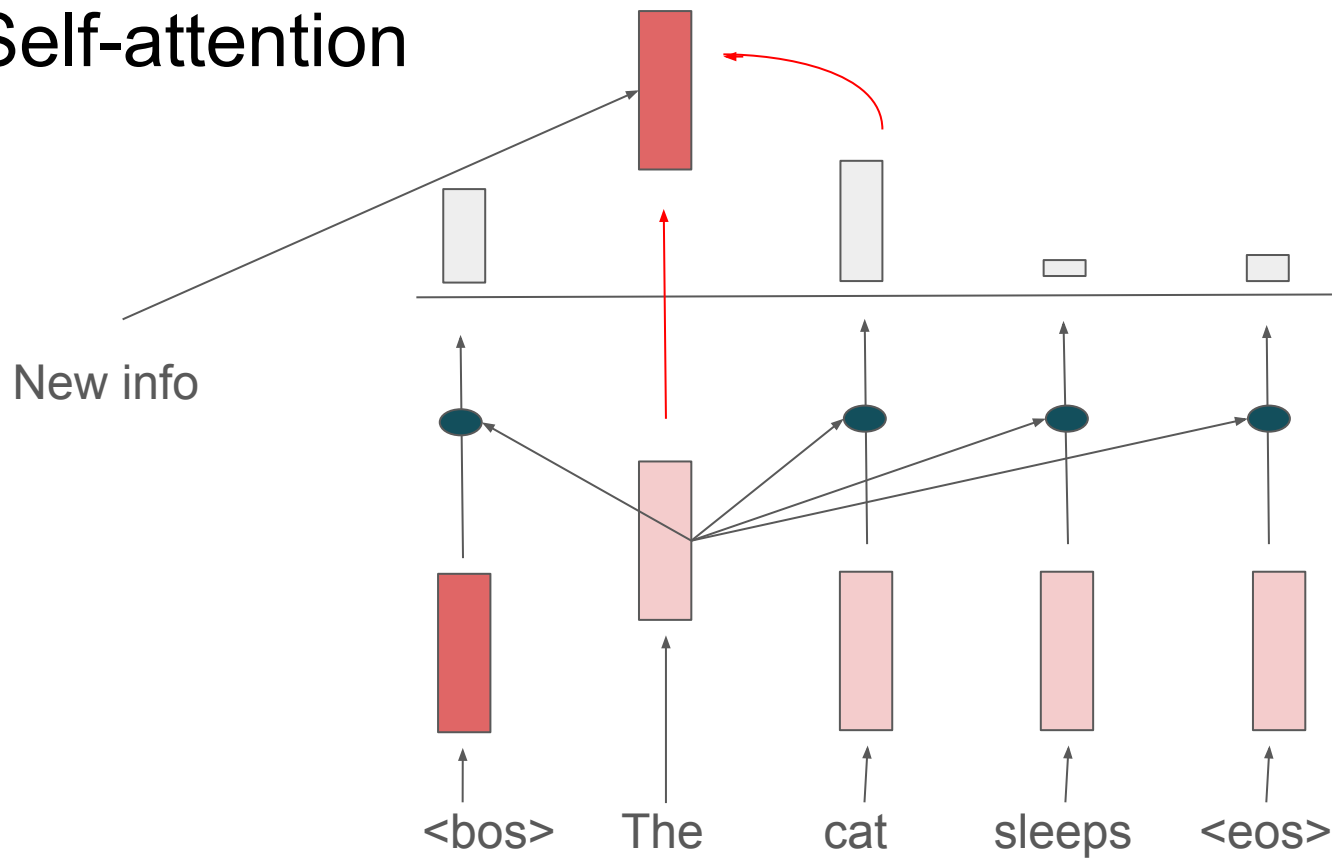
# Self-attention



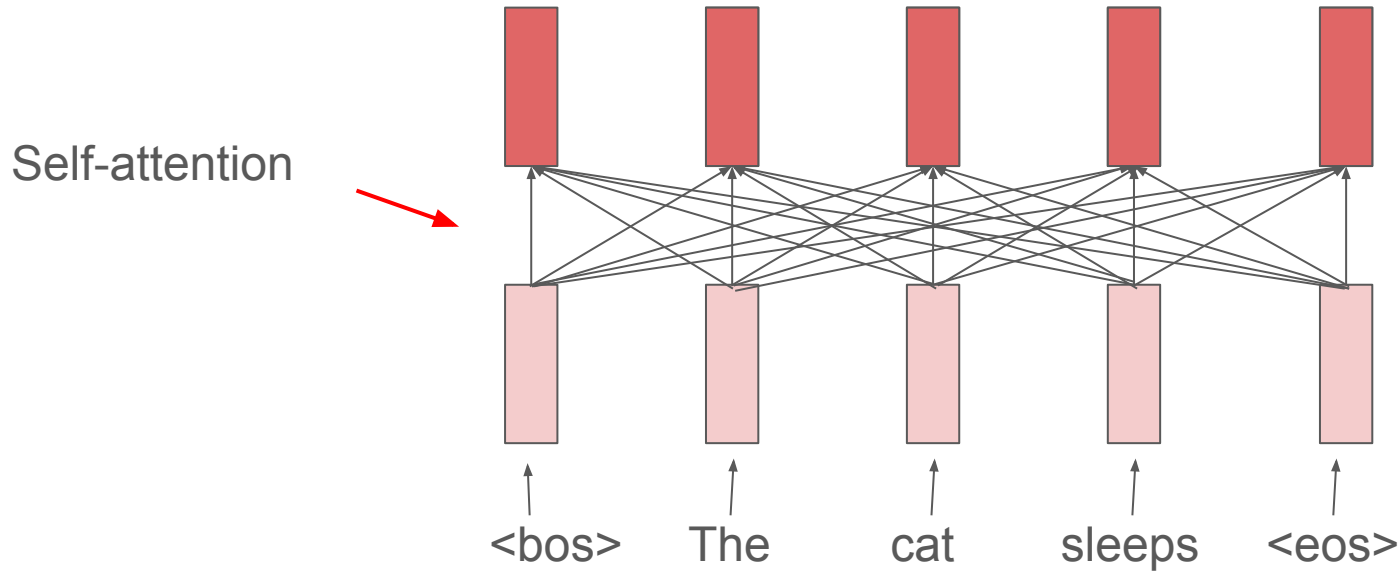
# Self-attention



# Self-attention




# Self-attention



# Multi-head attention

Она ест вкусную шаурму

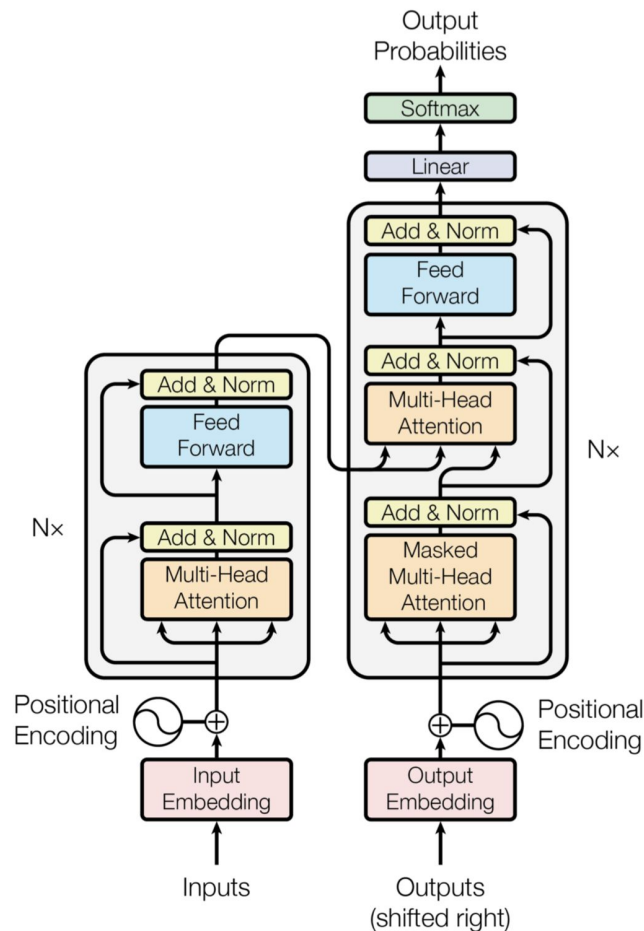


- Case agreement
- Gender agreement

$$Multi - head\ attn = Concat(attn_1, attn_2, \dots, attn_8)W$$

# Attention (~ Transformer)

- No recurrent -> parallel encoding  
-> faster
- Many attentions -> models does not have to remember much
- Multi-head attention -> able to pay attention to different aspects





# Transformer

# Transformer

*The animal didn't cross the street because it was too tired.*  
*L'animal n'a pas traversé la rue parce qu'il était trop fatigué.*

*The animal didn't cross the street because it was too wide.*  
*L'animal n'a pas traversé la rue parce qu'elle était trop large.*

# Transformer

