

# Neural Attention

Makarand Tapaswi  
CS7.505 Spring 2024  
14<sup>th</sup> February 2024

# Attention is All You Need!

- Almost there ... but, before that
- Neural Machine Translation by Jointly Learning to Align and Translate
- Dzmitry Bahdanau, Kyunghyun Cho, Yoshua Bengio
- ICLR 2015 Oral
- <https://arxiv.org/abs/1409.0473>

# Recurrent Neural Networks

$$h_t = f(h_{t-1}, x_t)$$

- Why recurrent?
- Different variants of function  $f$  (GRU, LSTM, etc.)
- Long Short-Term Memory Units
- Gated Recurrent Units

# Examples of Sequence-to-Sequence tasks?

- Image Captioning
  - (strictly: image input to sequence)
- Sentiment Classification
  - (strictly: sequence to single output)
- Machine Translation
  - yes! sequence to sequence

# Very nice material that we will follow!

- [https://lennavoita.github.io/nlp\\_course/seq2seq\\_and\\_attention.html#main\\_content](https://lennavoita.github.io/nlp_course/seq2seq_and_attention.html#main_content)



# Formulating machine translation

## 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**

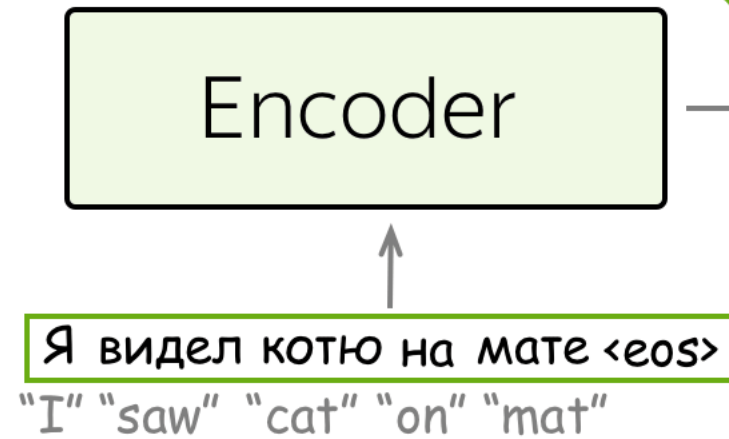
How to find  $\theta$ ?

- **search**

How to find the argmax?

# Encoder - decoder

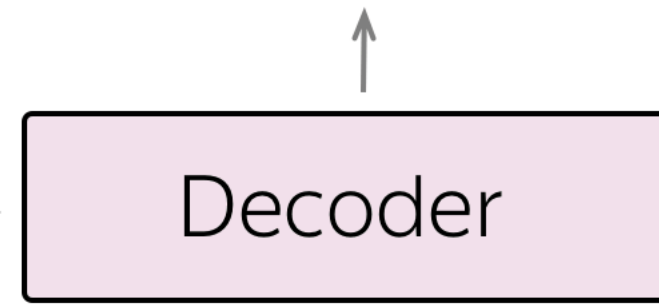
Encoder builds a representation of the source and gives it to the decoder



Source sentence

Target sentence

I saw a cat on a mat <eos>




Decoder uses this source representation to generate the target sentence

# Conditional Language Models

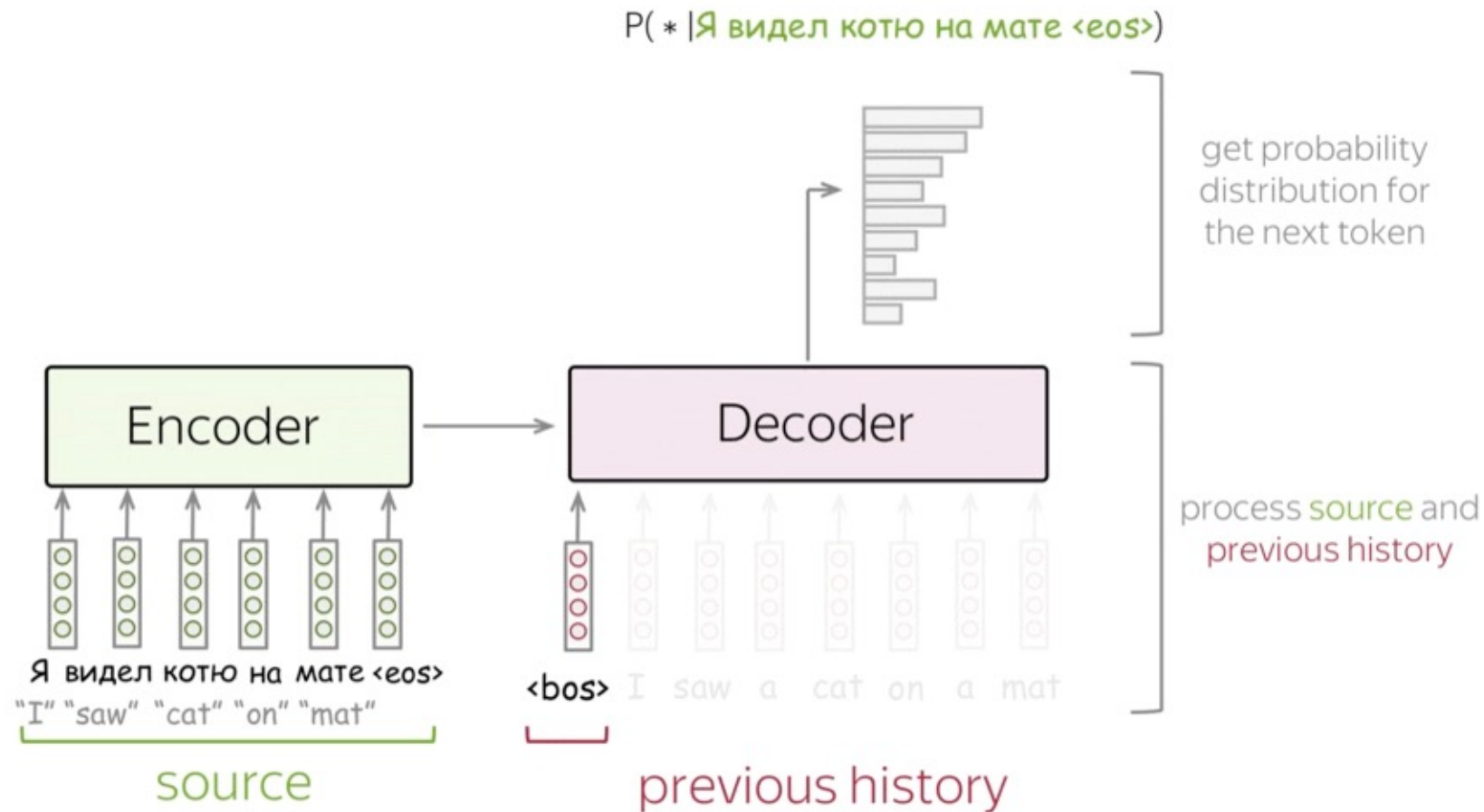
Language Models:  $P(y_1, y_2, \dots, y_n) = \prod_{t=1}^n p(y_t | y_{<t})$

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

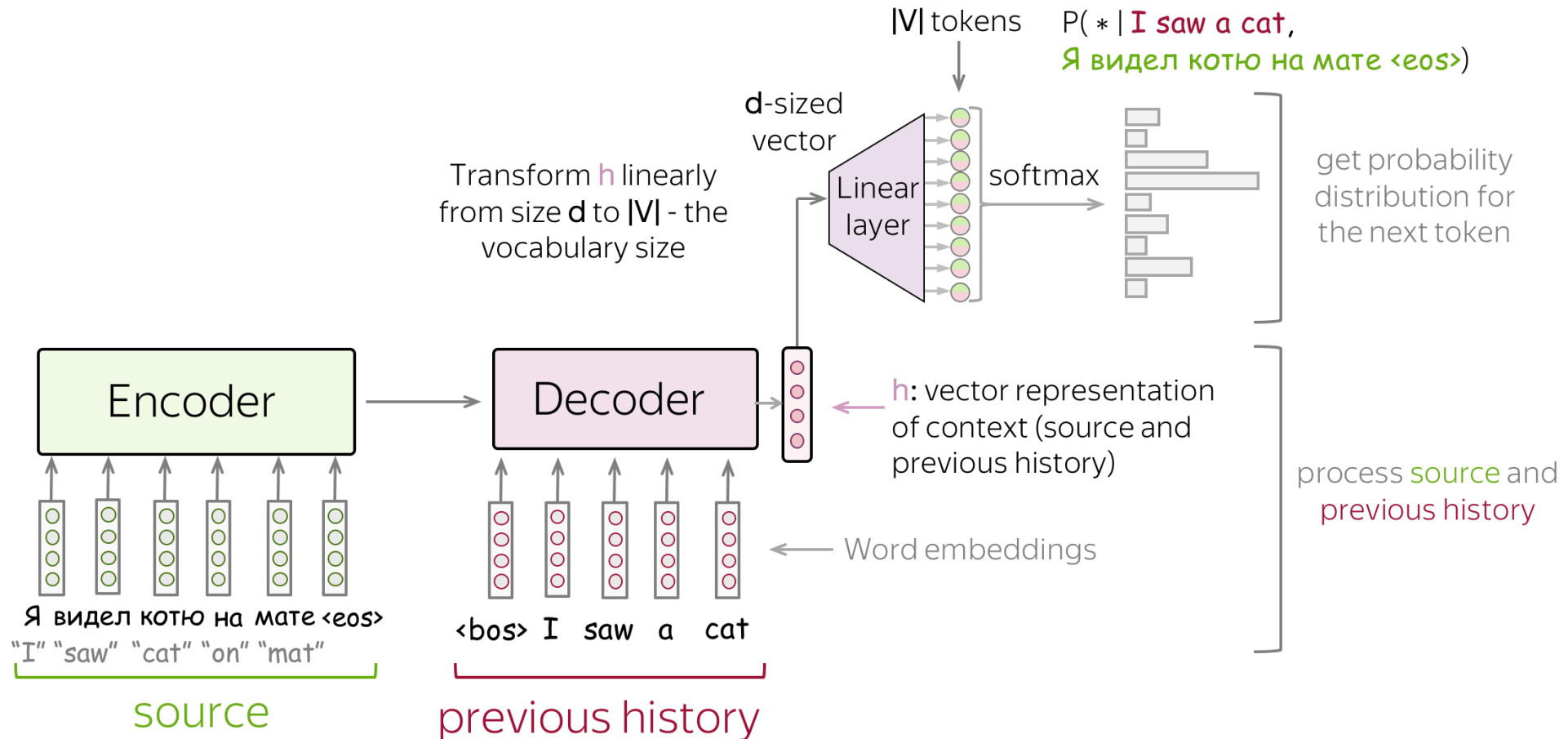
 condition on source  $x$



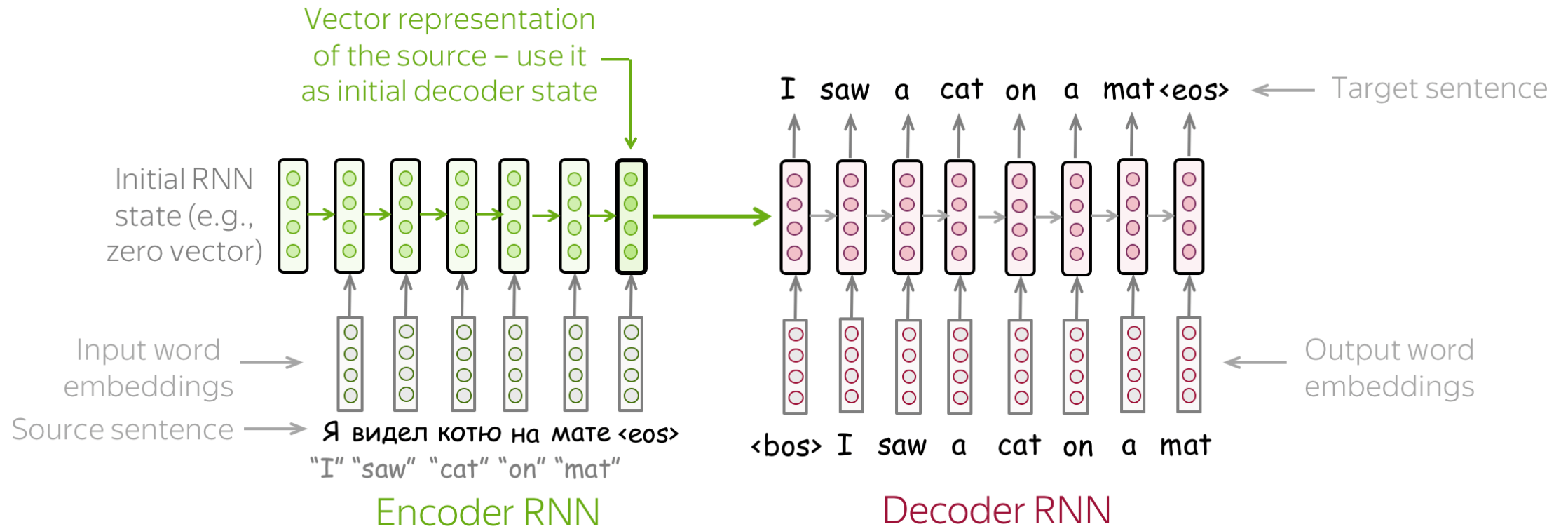
# Encoder Decoder explained



# Using a decoder hidden representation

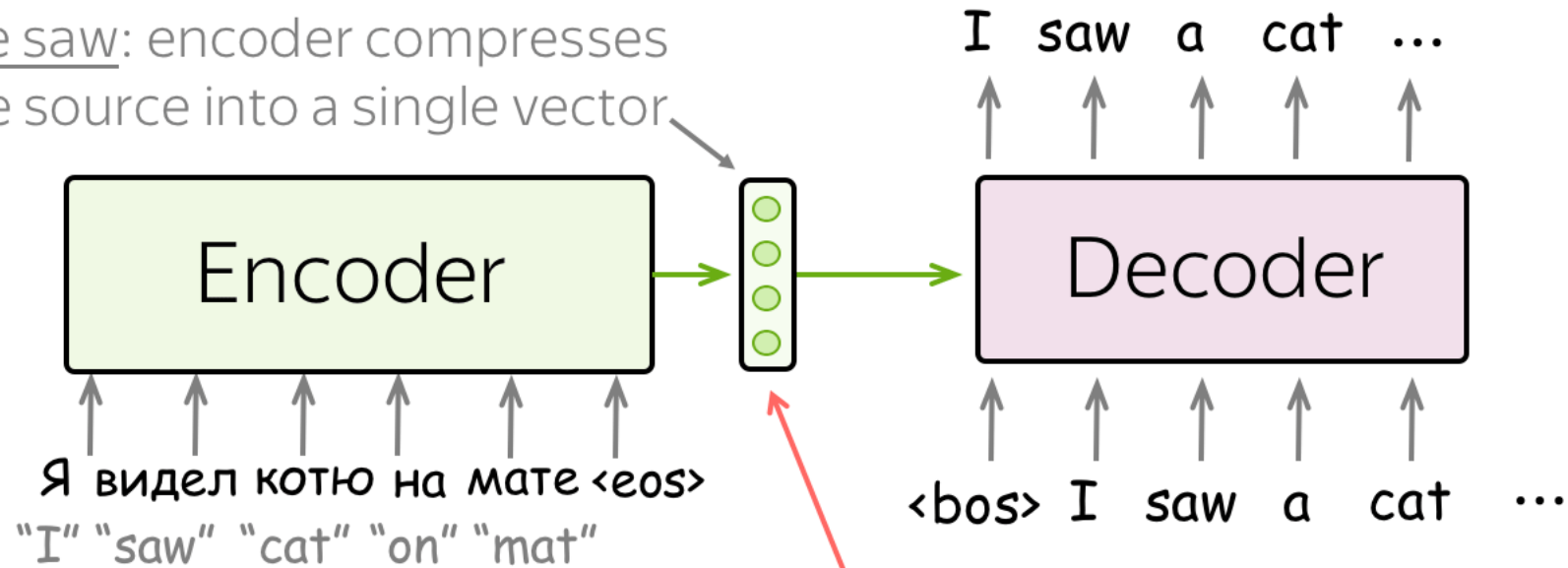


# Seq2Seq RNNs



# Fixed encoder representations are an issue

We saw: encoder compresses the source into a single vector



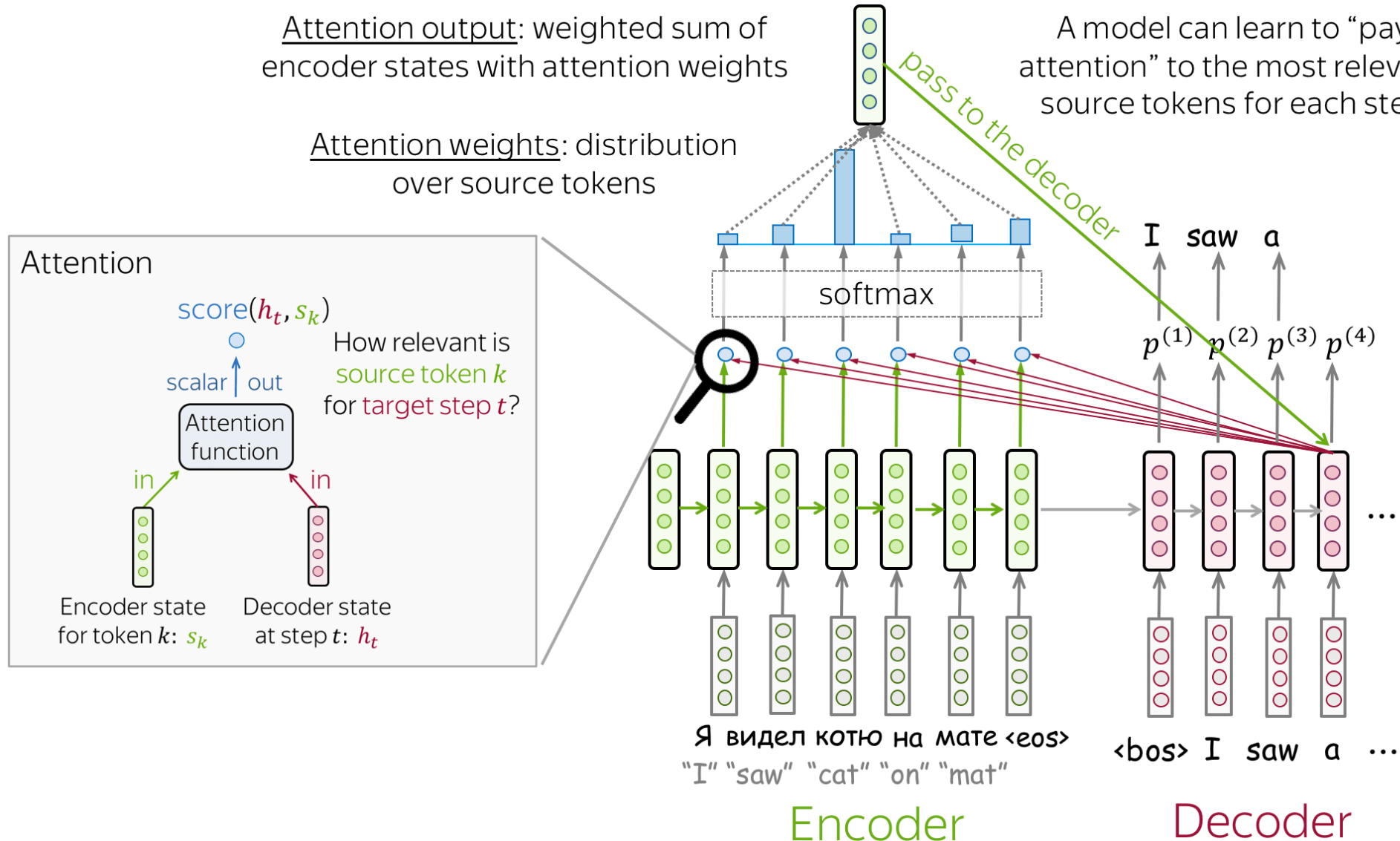
Problem: this is a bottleneck!

# Attention

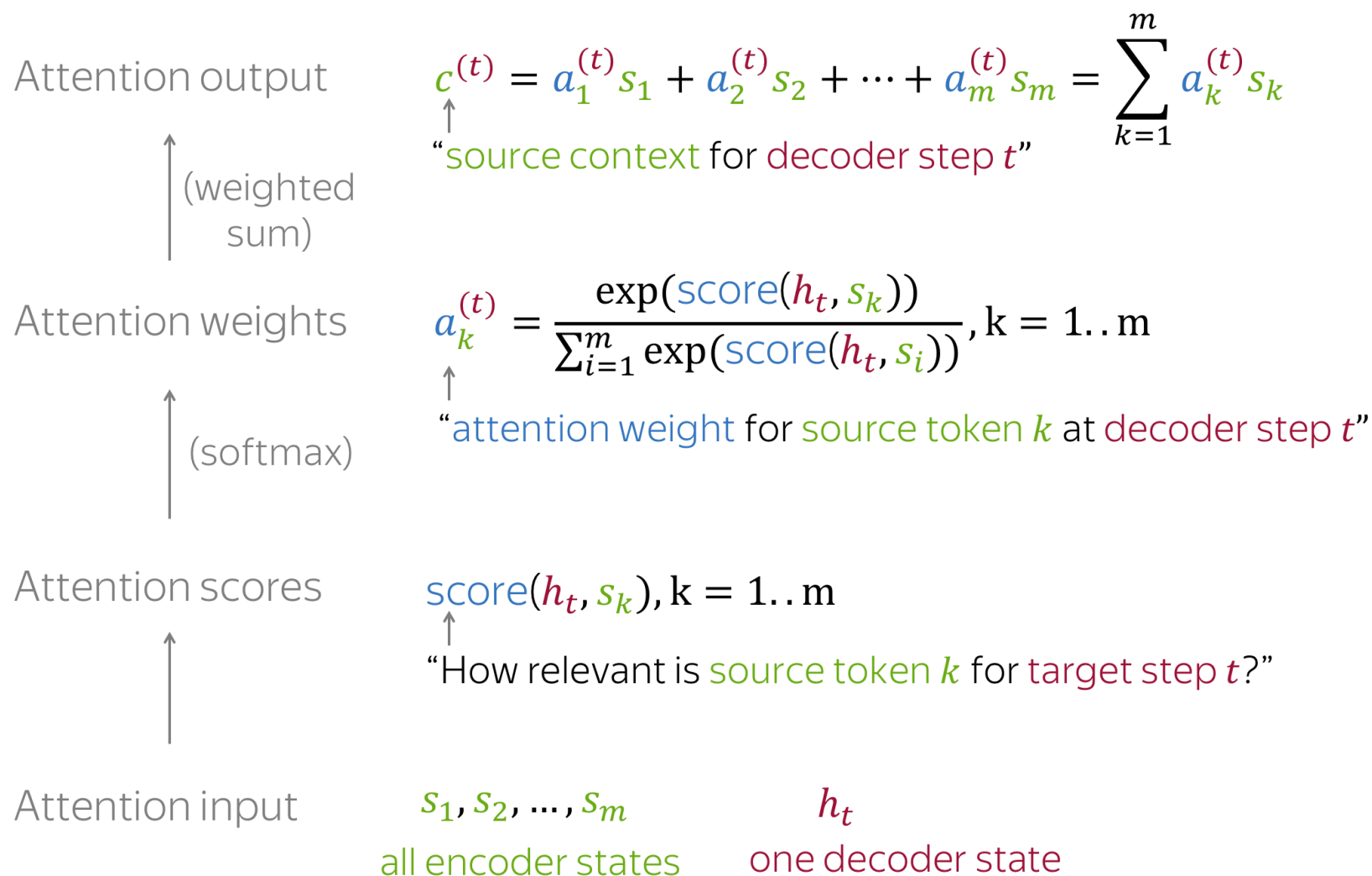
Attention output: weighted sum of encoder states with attention weights

Attention weights: distribution over source tokens

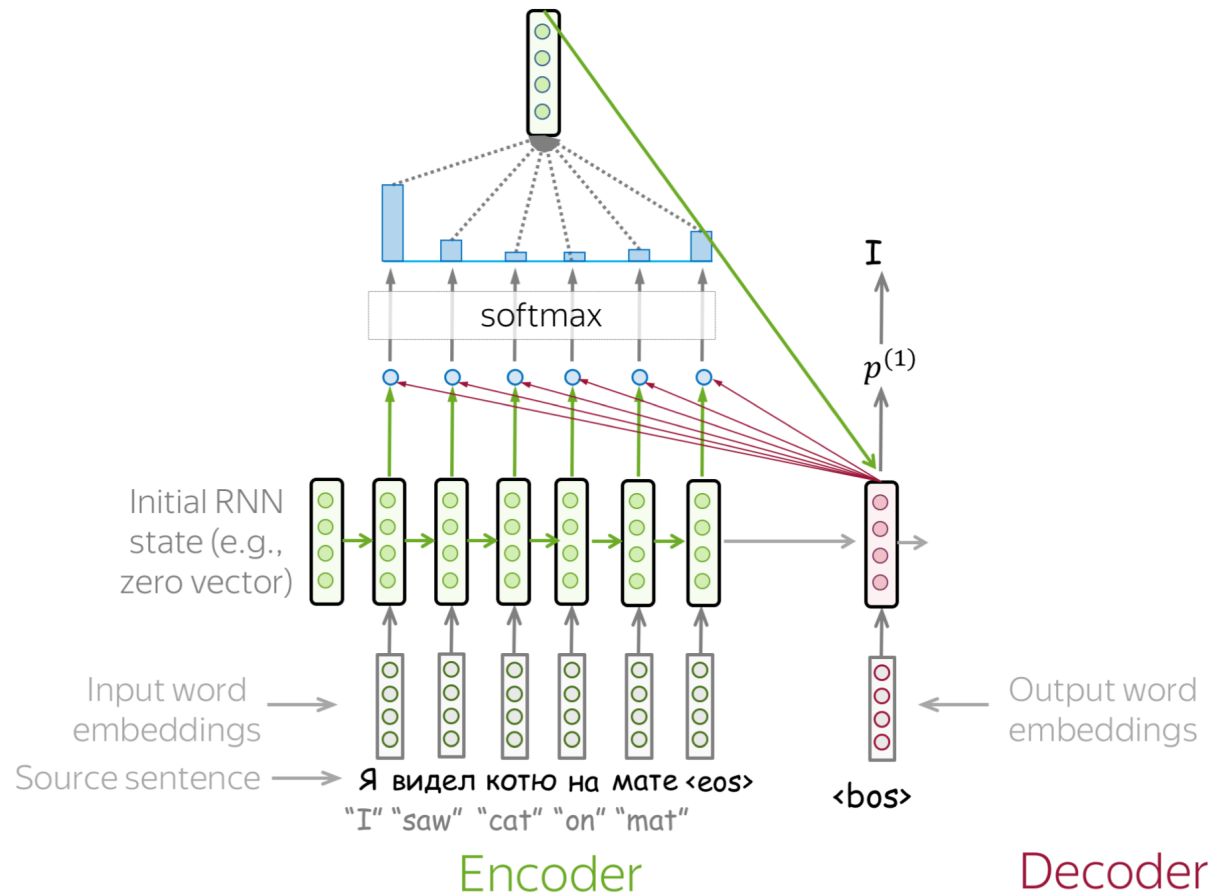
A model can learn to “pay attention” to the most relevant source tokens for each step



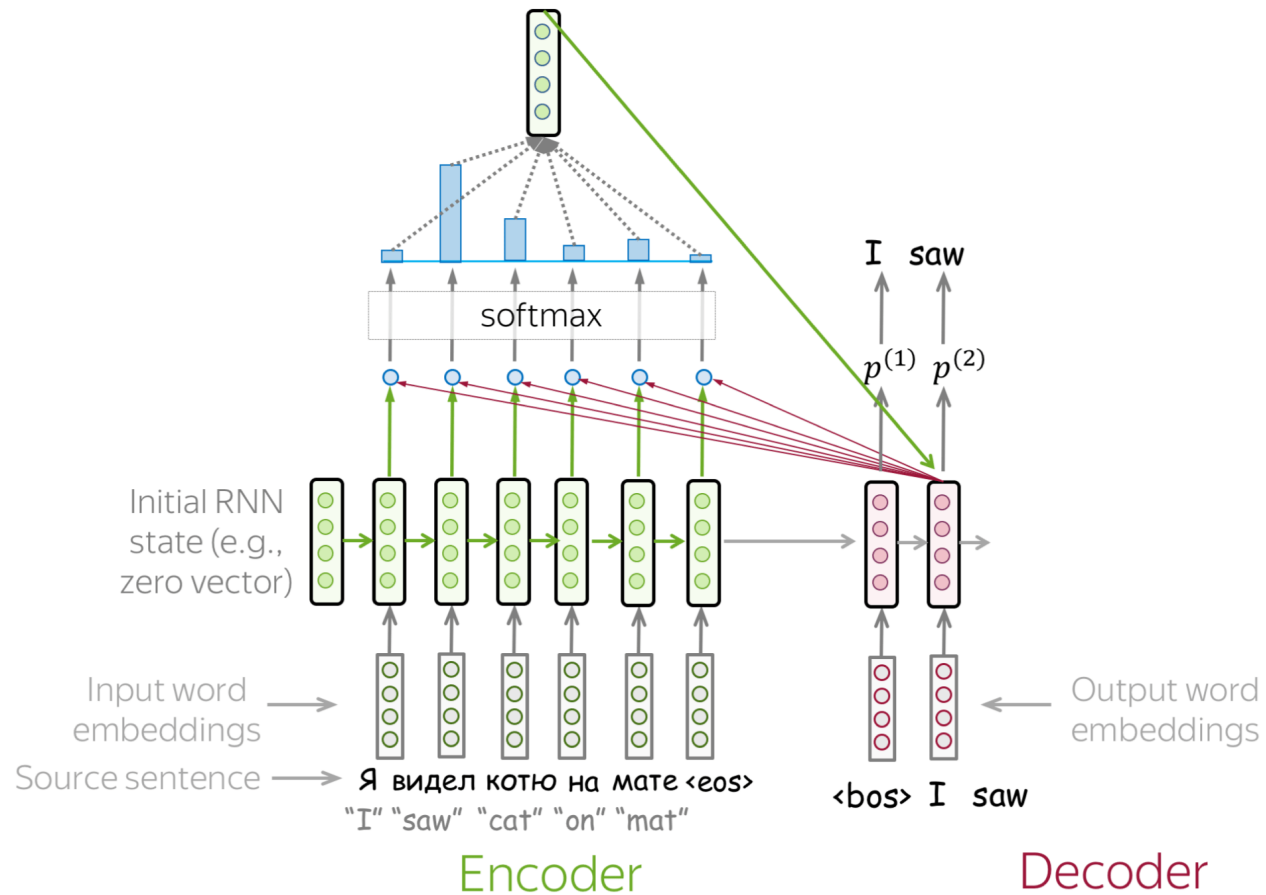
# How to compute attention?



# Decoding word 1

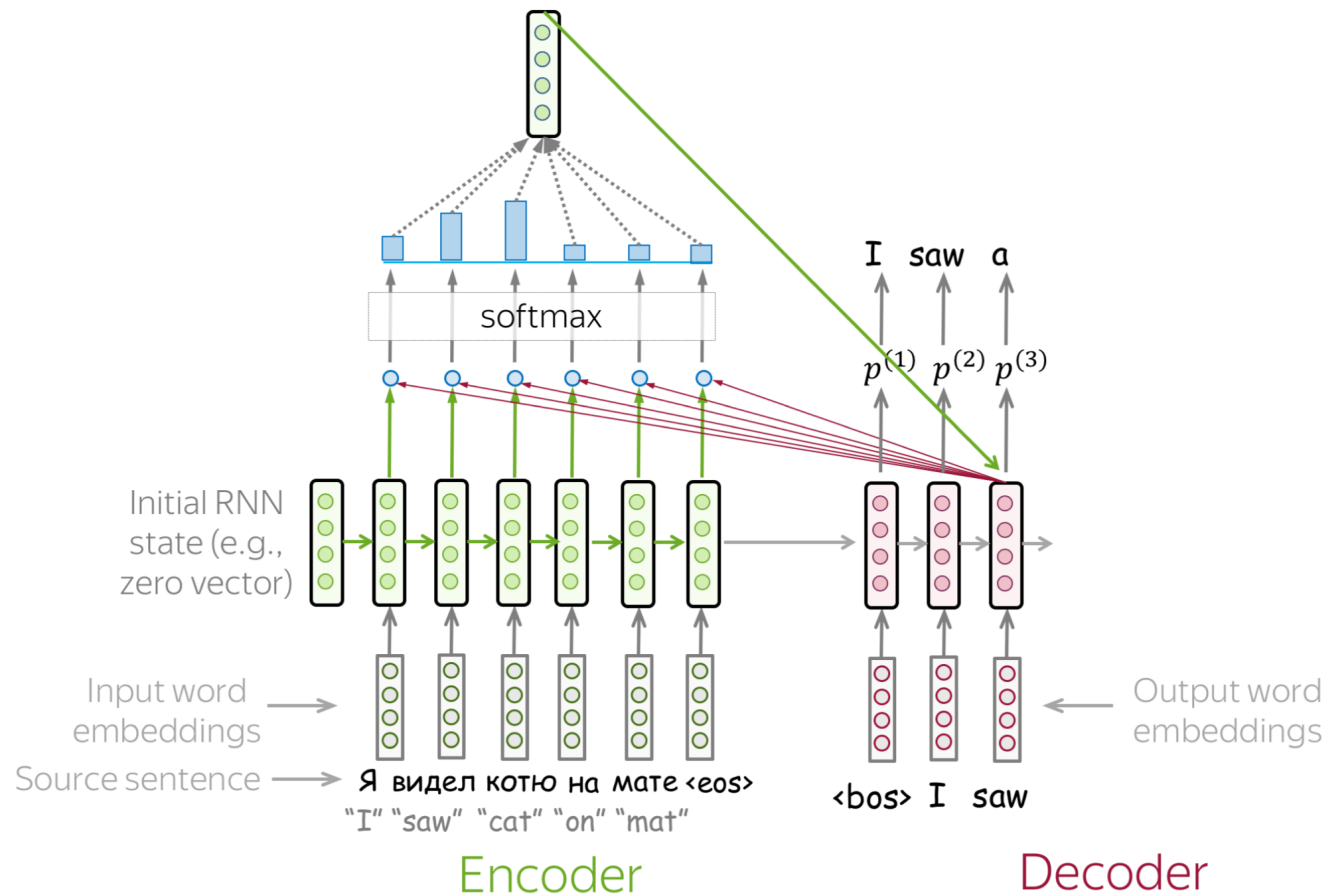


# Decoding word 2

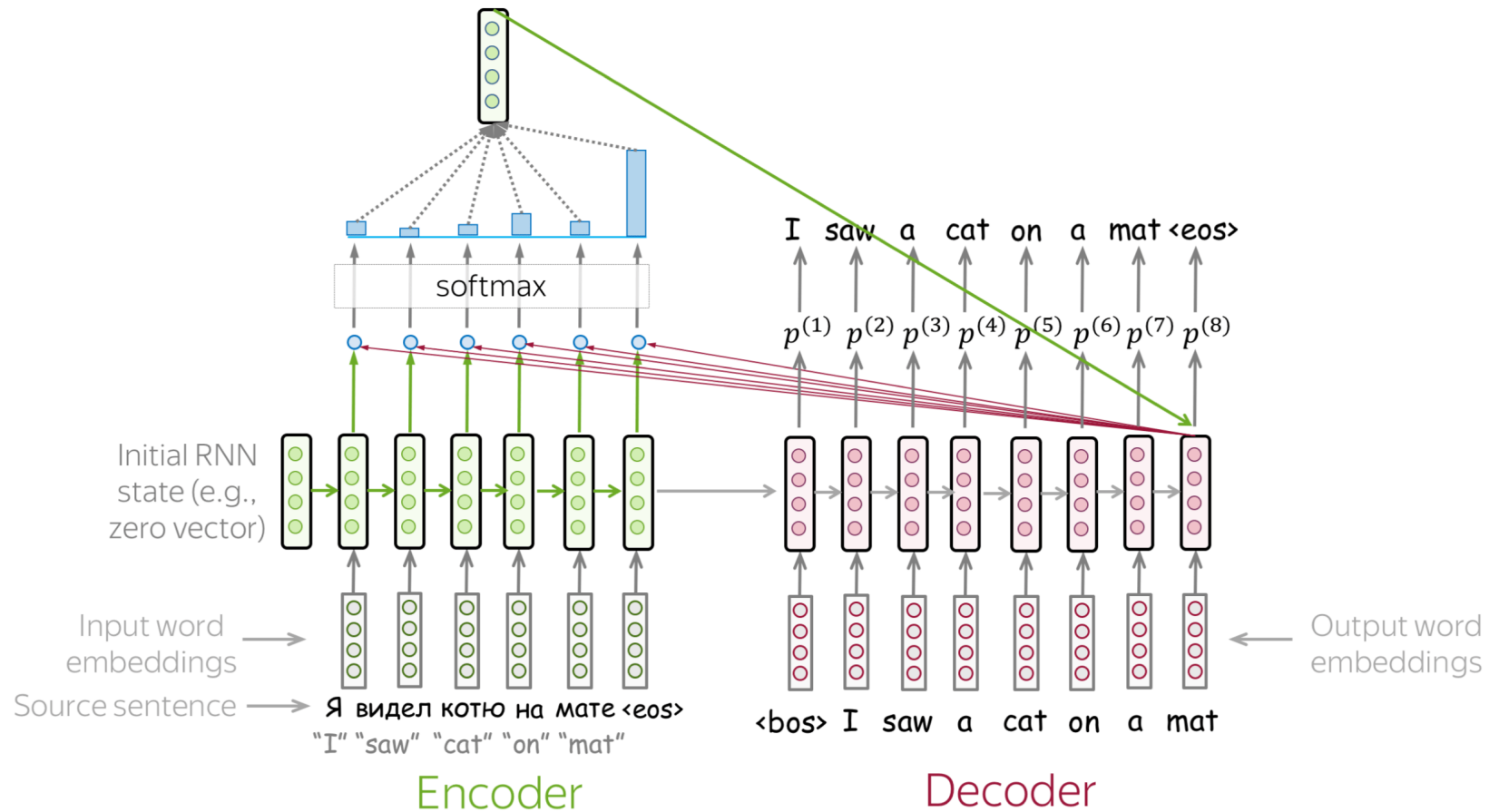




# Decoding word 3

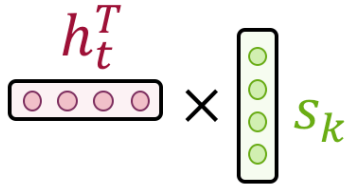


# Decoding word N



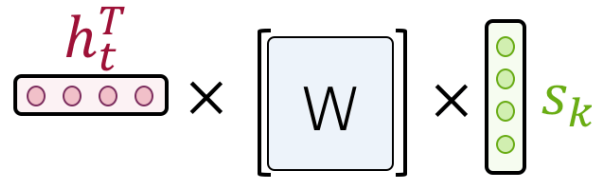
# Computing Attention

Dot-product


$$h_t^T \times s_k$$

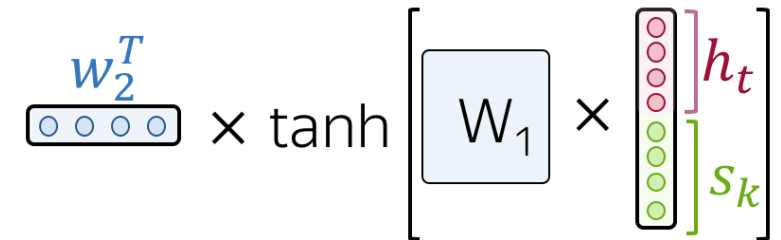
$$\text{score}(h_t, s_k) = h_t^T s_k$$

Bilinear


$$h_t^T \times W \times s_k$$

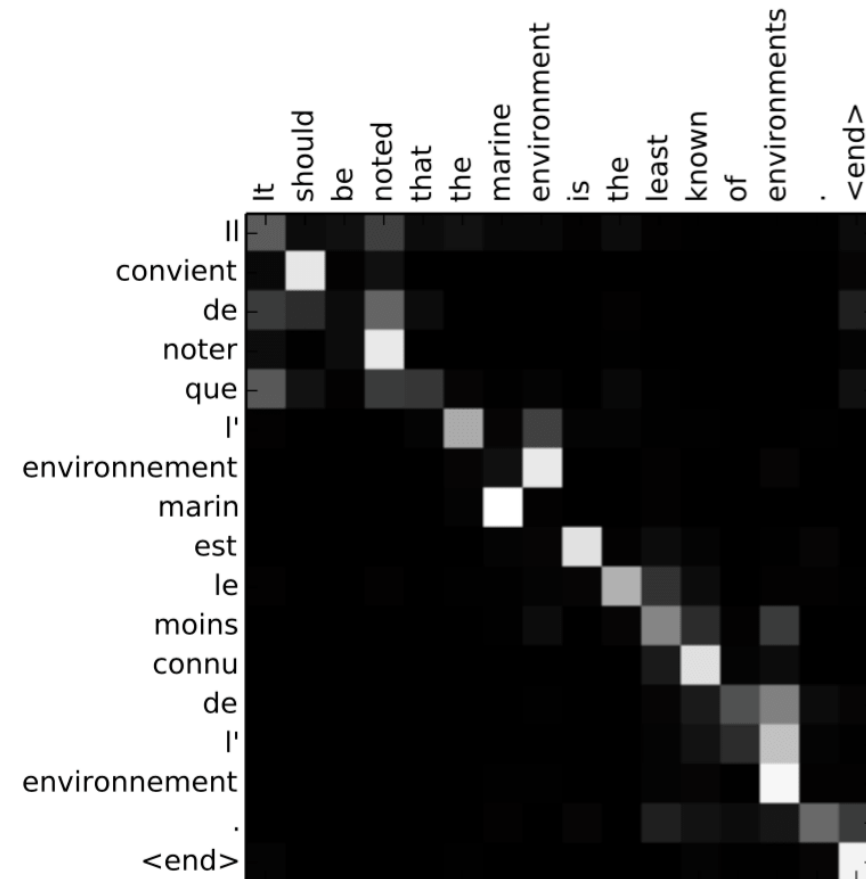
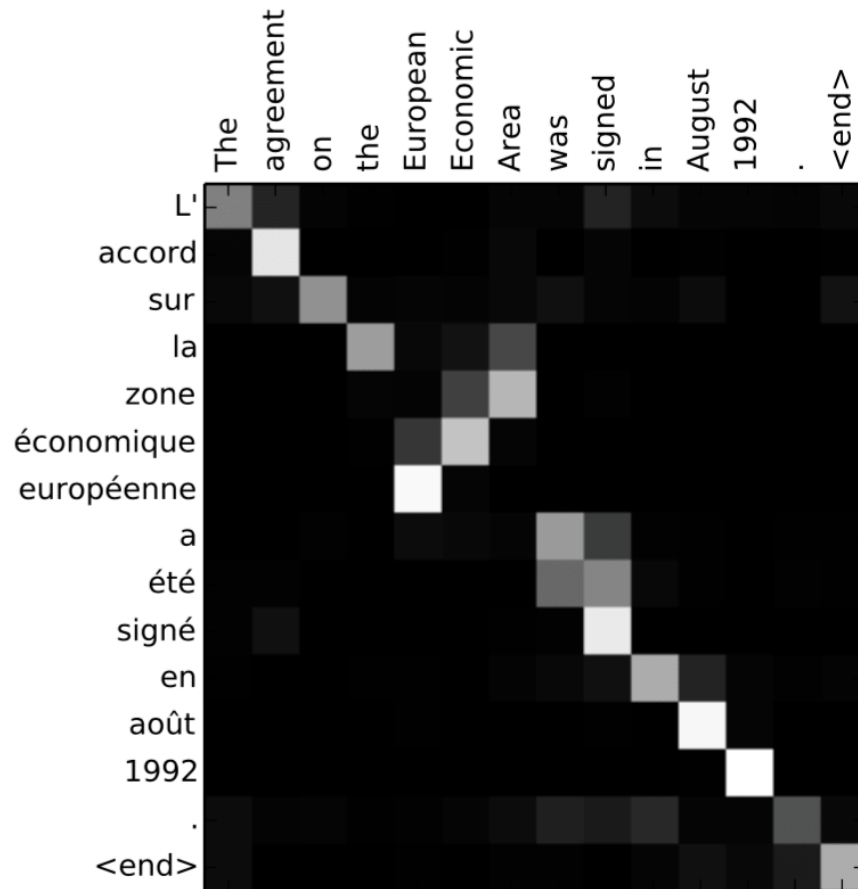
$$\text{score}(h_t, s_k) = h_t^T W s_k$$

Multi-Layer Perceptron


$$w_2^T \times \tanh \left[ W_1 \times \begin{bmatrix} h_t \\ s_k \end{bmatrix} \right]$$

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

# Attention Learns Alignment!



Thank you!