

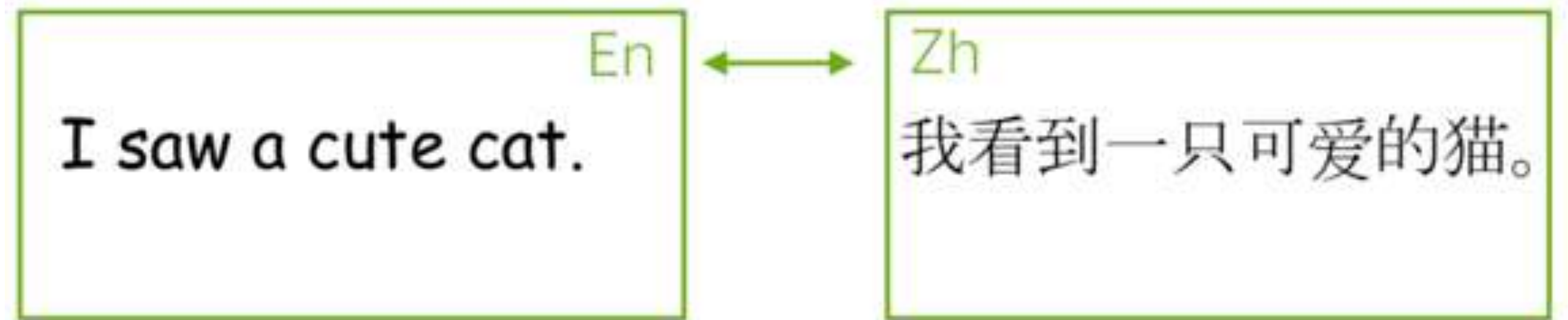


Seq2seq and Attention


Lena Voita

Sequence to Sequence Task


- Translation between natural languages
- More generally, translation between any sequences



What is going to happen:

- Seq2seq Basics
- Attention
- Transformer
- Subword Segmentation: BPE
-  Analysis and Interpretability

What is going to happen:

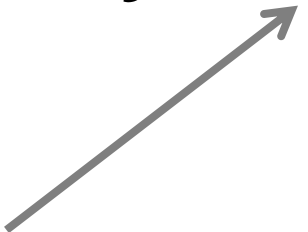
- Seq2seq Basics →
 - Machine Translation Task
 - Encoder-Decoder Framework
 - Conditional LMs
 - The Simplest RNN Model
 - Training
 - Inference
- Attention
- Transformer
- Subword Segmentation: BPE
-  Analysis and Interpretability

Translation

Human Translation

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

The “probability” is
intuitive and is given
by a human
translator’s expertise

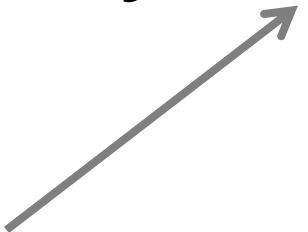


Translation


Human Translation

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

The “probability” is
intuitive and is given
by a human
translator’s expertise



Machine Translation

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


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

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

model parameters

Questions we need to answer

- **modeling**

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

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

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

model parameters

Questions we need to answer

- modeling

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

- learning

How to find θ ?

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

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

model parameters

Questions we need to answer

- **modeling**

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

- **learning**

How to find θ ?

- **search**

How to find the argmax?

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

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

model parameters

Questions we need to answer

- **modeling**

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

- **learning**

How to find θ ?

- **search**

How to find the argmax?

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 θ ?

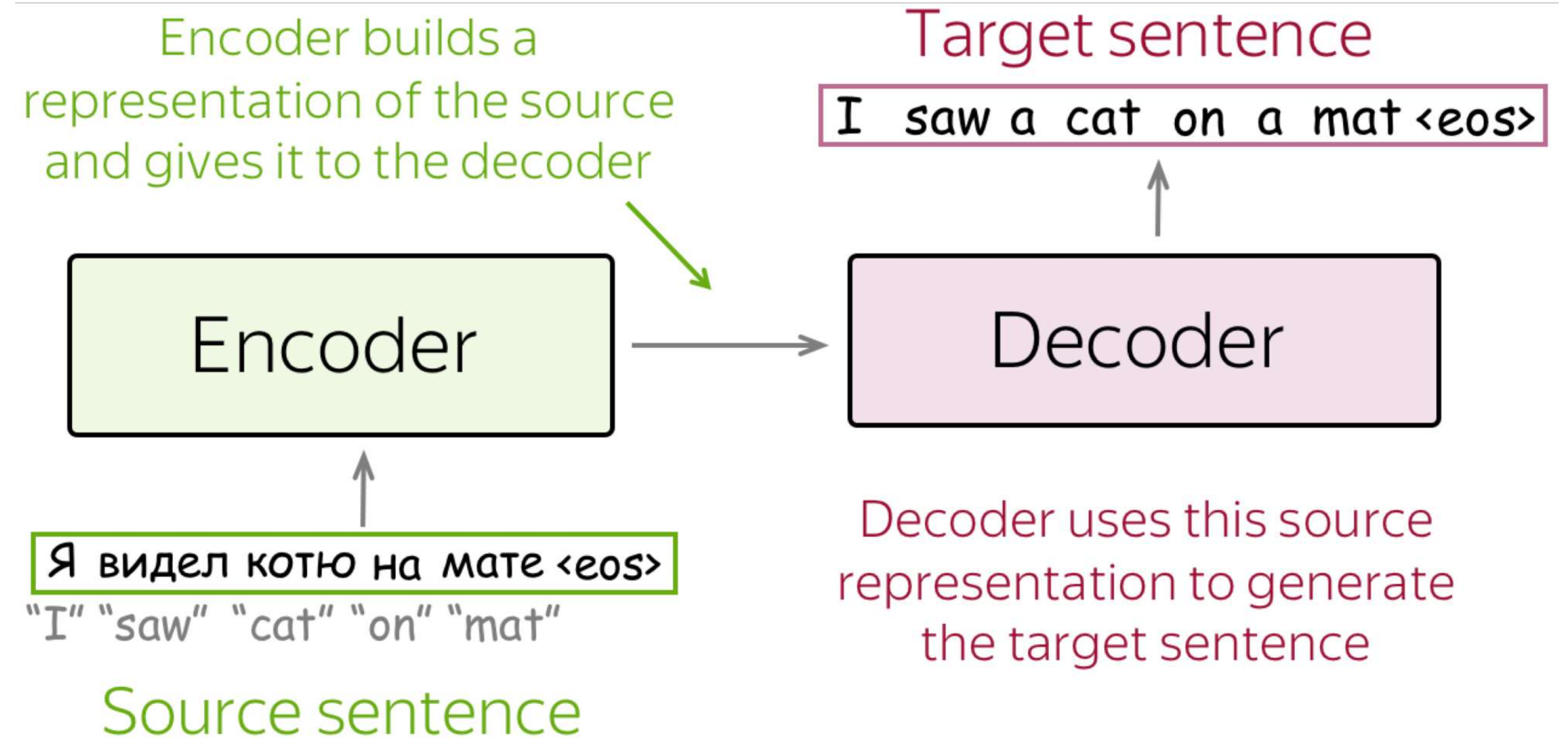
- **search**

How to find the argmax?

Encoder-Decoder Framework

The standard modeling paradigm:

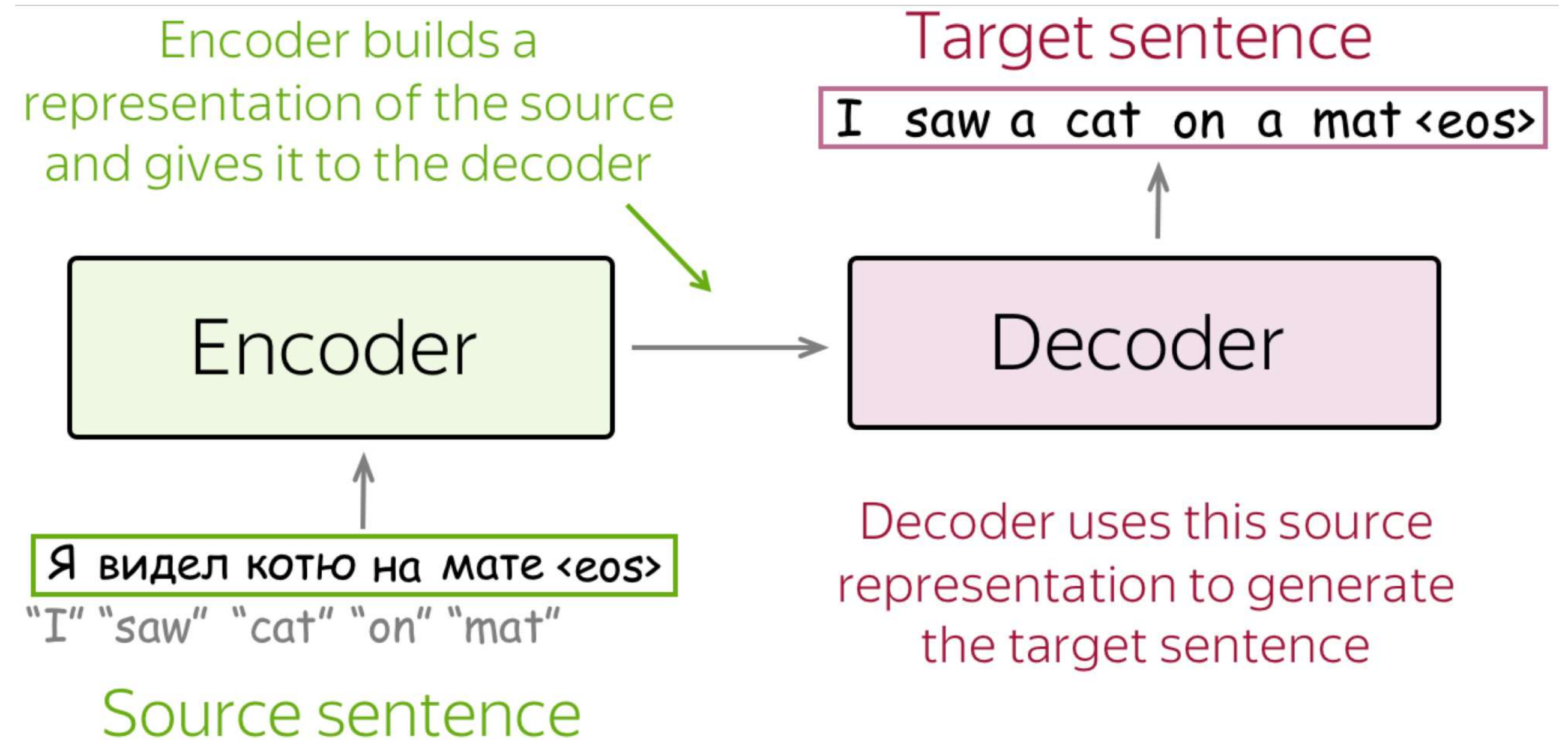
- **Encoder** – reads the source sentence and produces its representation



Encoder-Decoder Framework

The standard modeling paradigm:

- **Encoder** – reads the source sentence and produces its representation
- **Decoder** - uses source representation from the encoder to generate the target sequence.




Conditional Language Models

Language Models:
(left-to-right)

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

Conditional Language Models

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

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

condition on source x

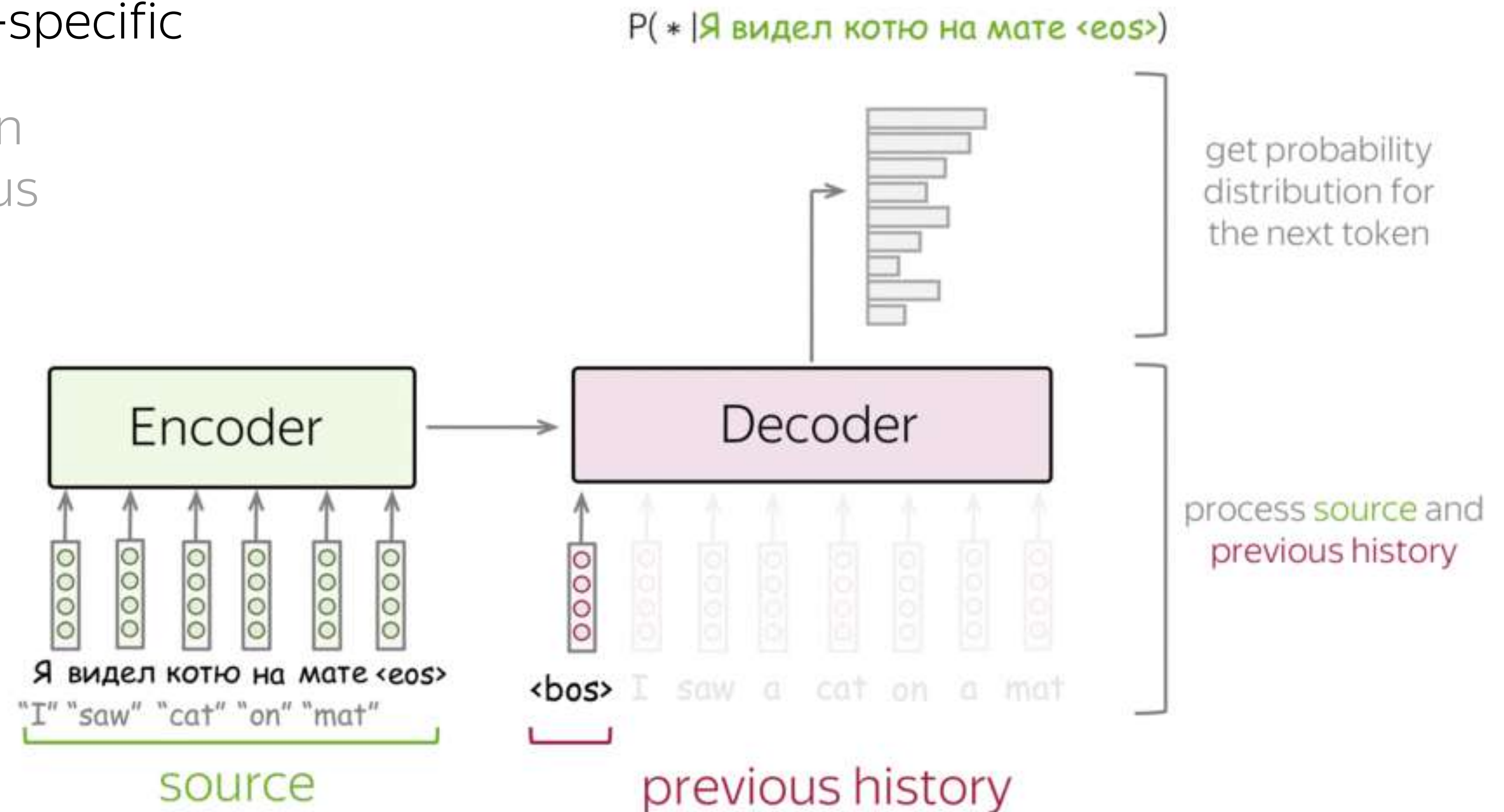
General View

- process context – model-specific

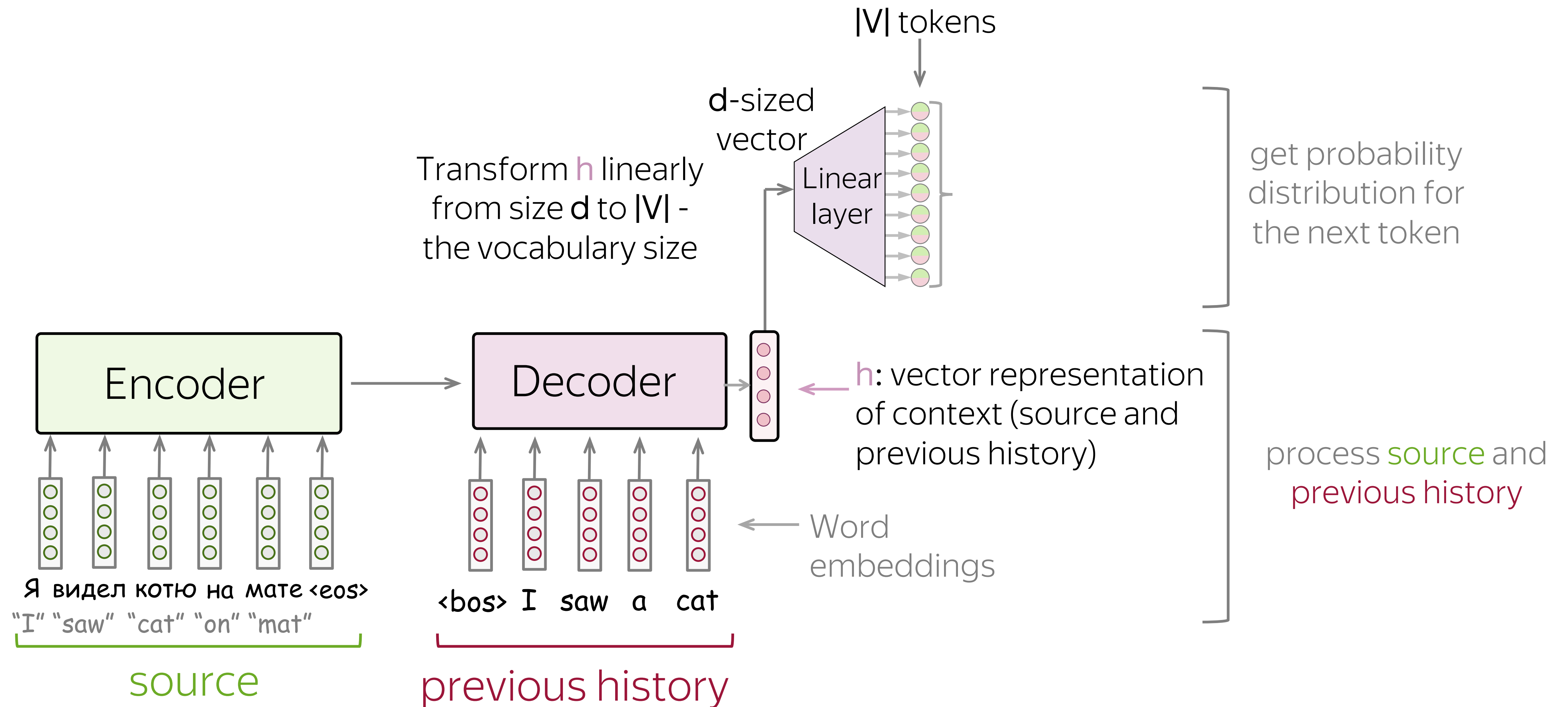
Get vector representation of the source and previous target tokens

- evaluate probabilities – model-agnostic

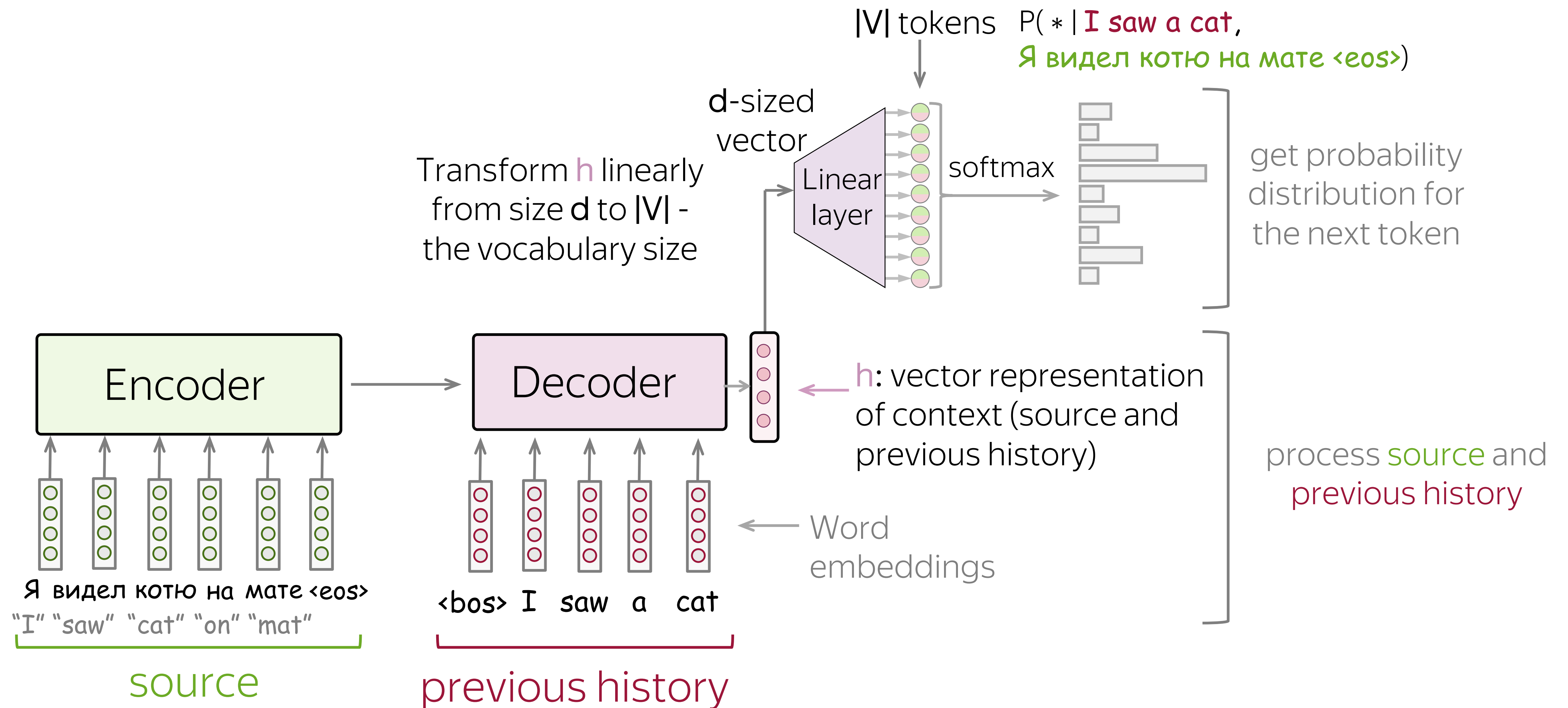
Predict probability distribution for the next target token



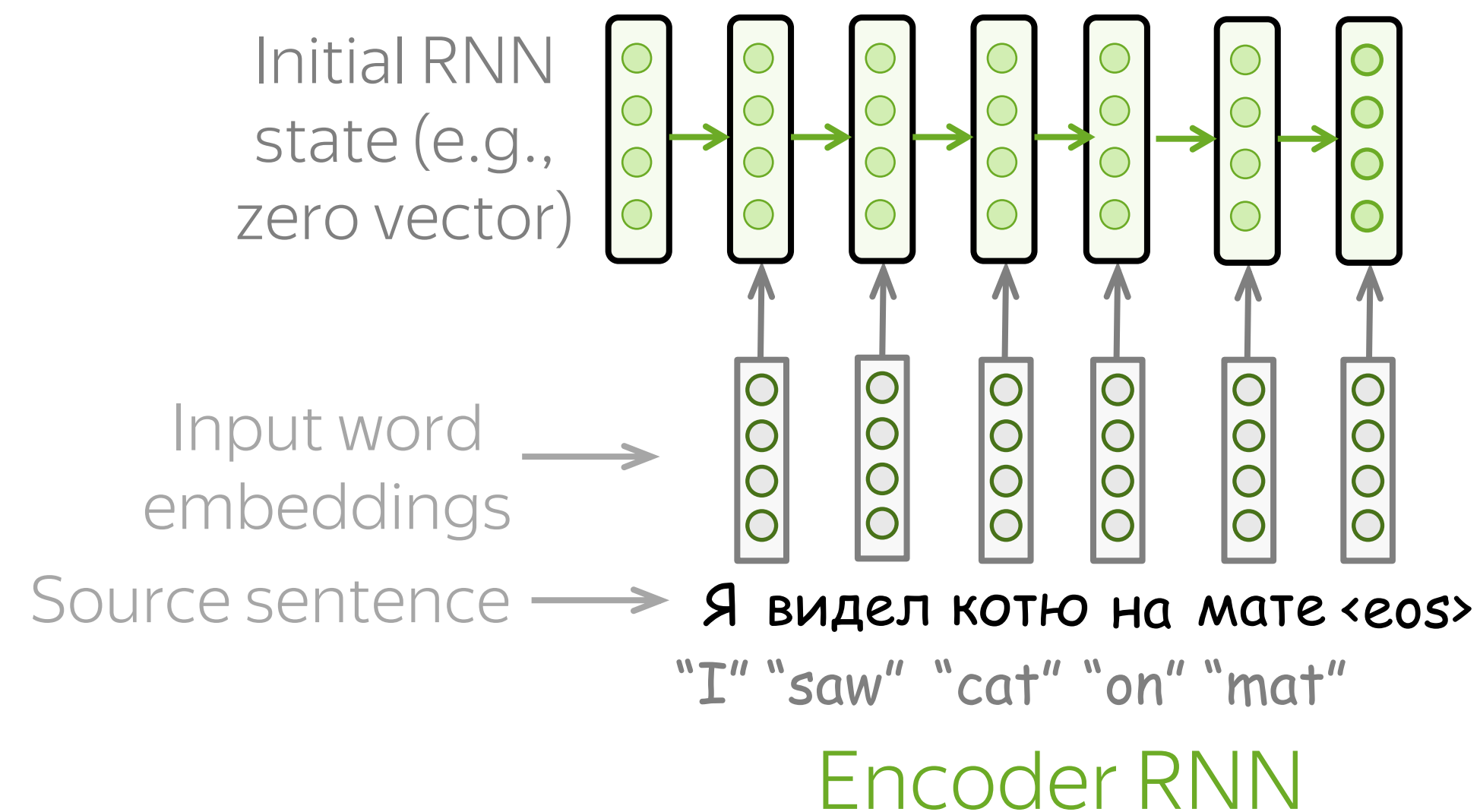
High-Level Pipeline



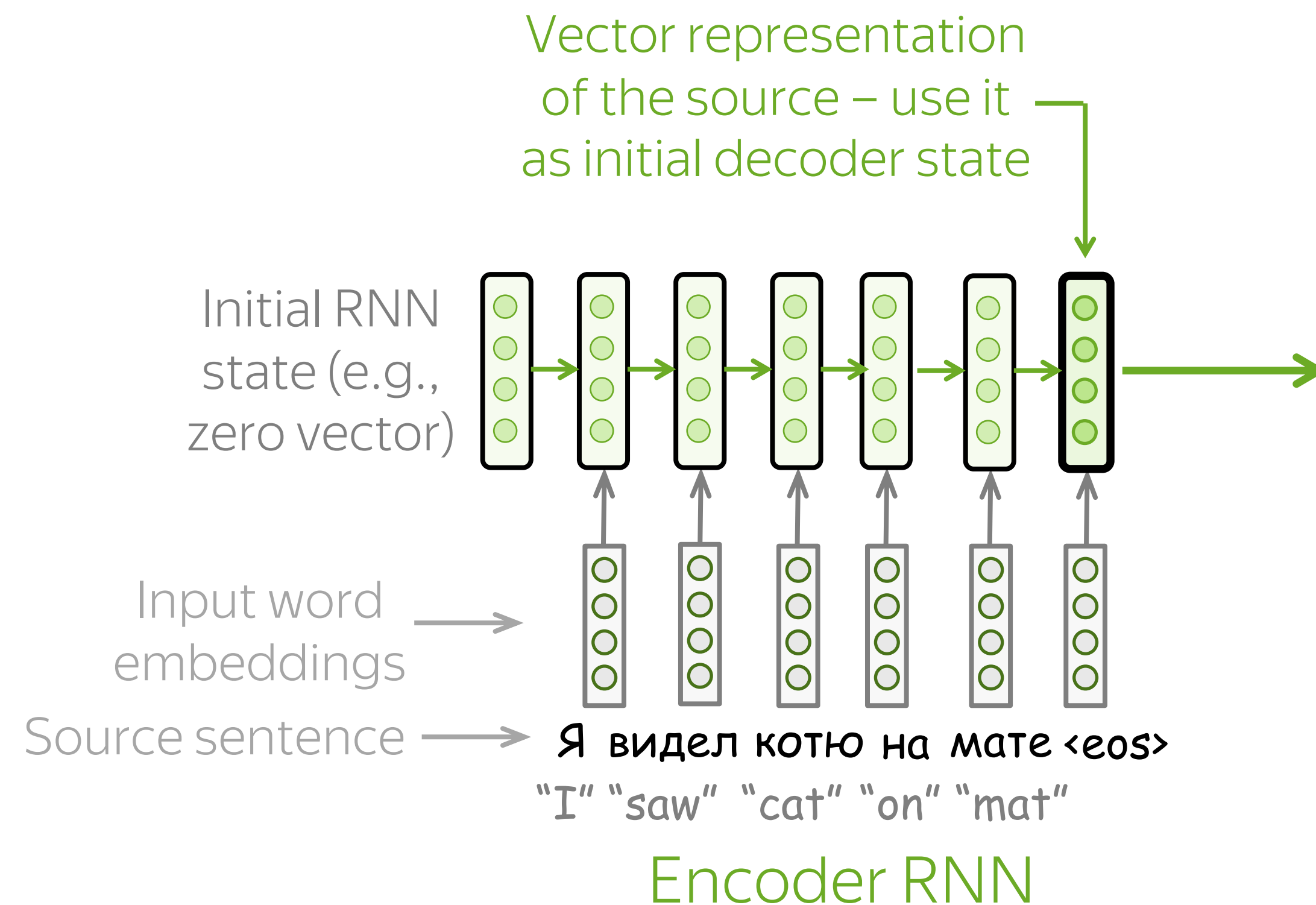
High-Level Pipeline



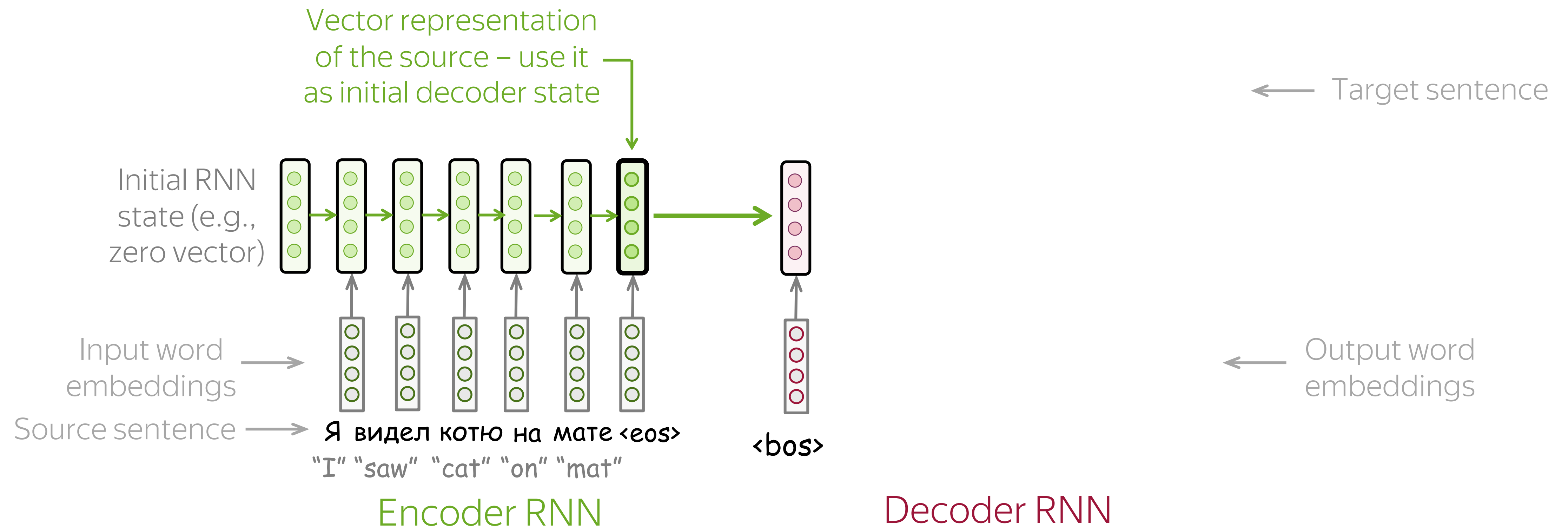
The Simplest Model: RNN Encoder and Decoder



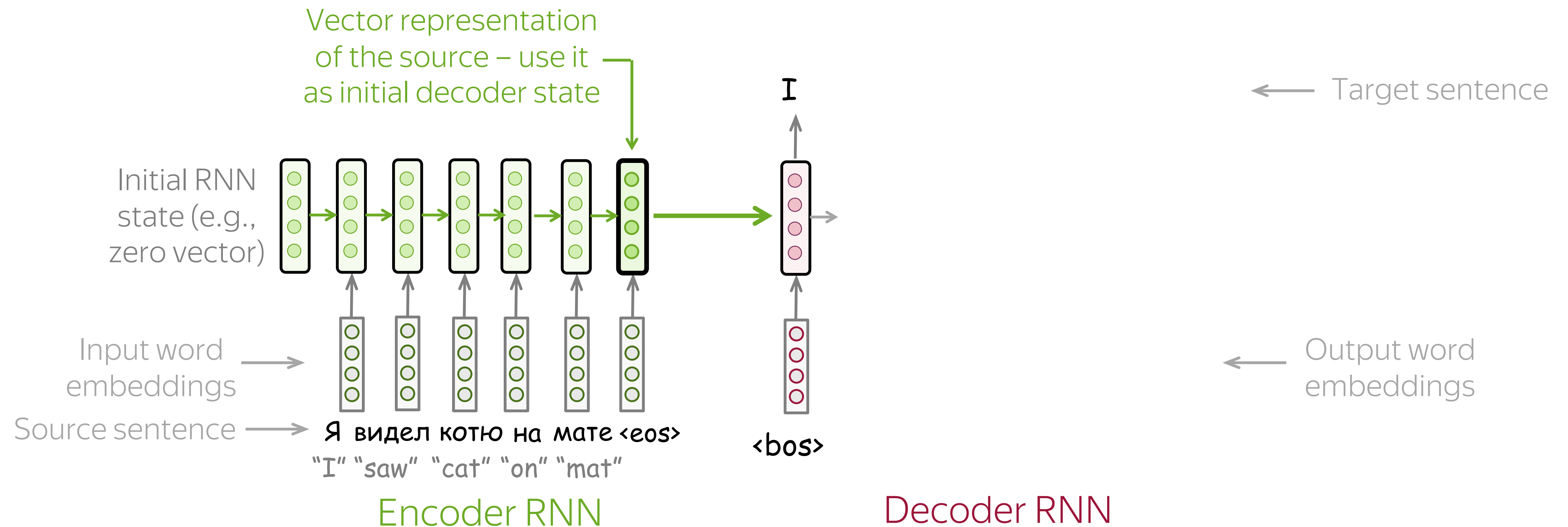
The Simplest Model: RNN Encoder and Decoder



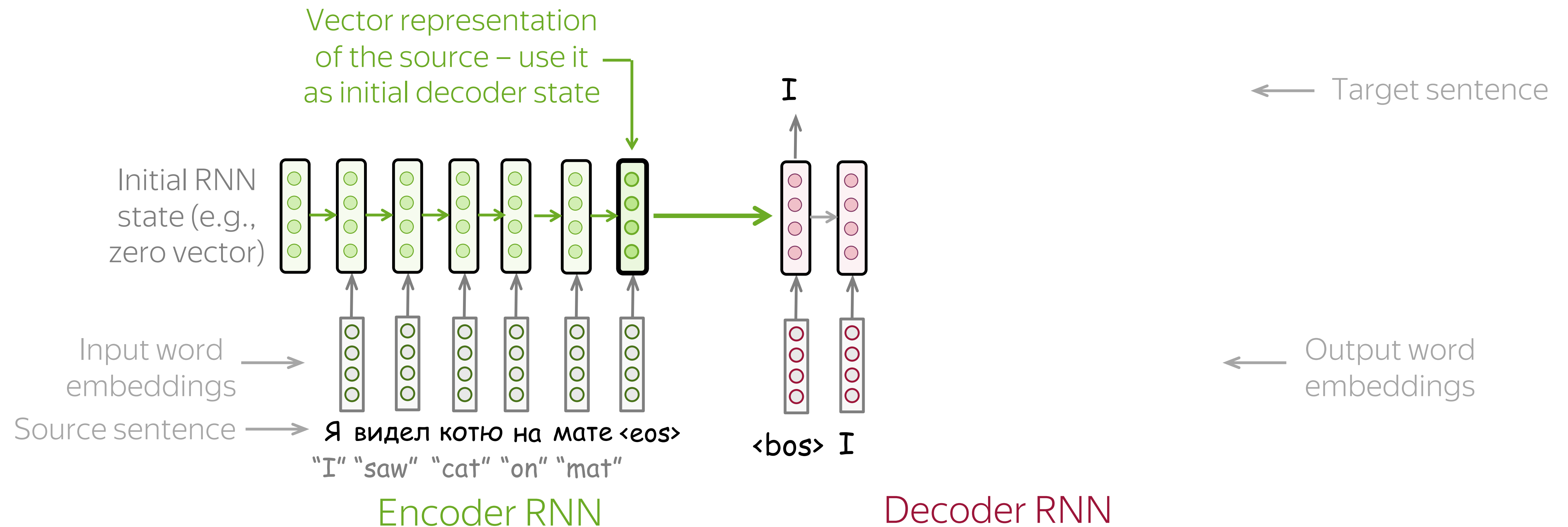
The Simplest Model: RNN Encoder and Decoder



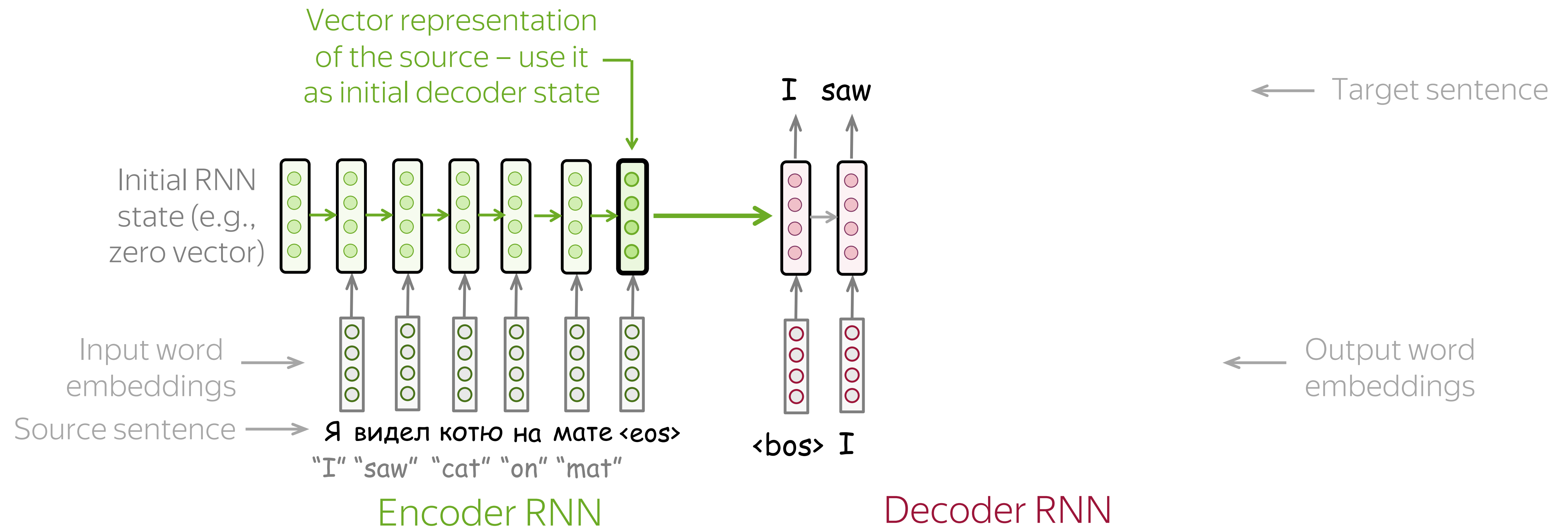
The Simplest Model: RNN Encoder and Decoder



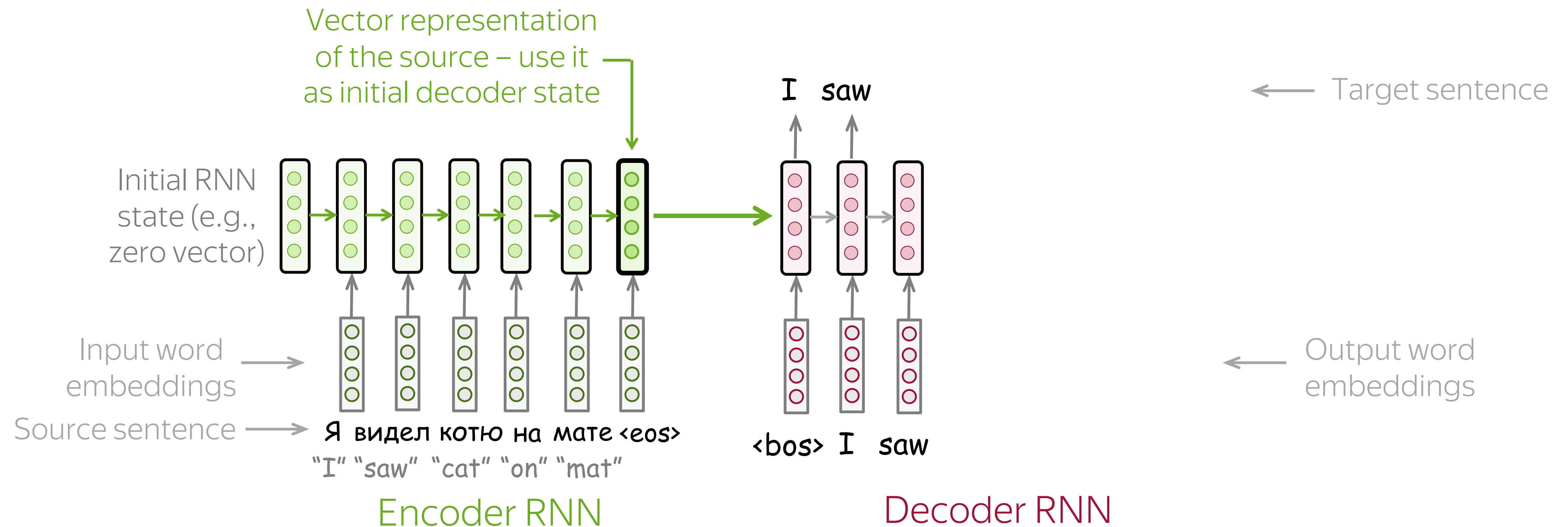
The Simplest Model: RNN Encoder and Decoder



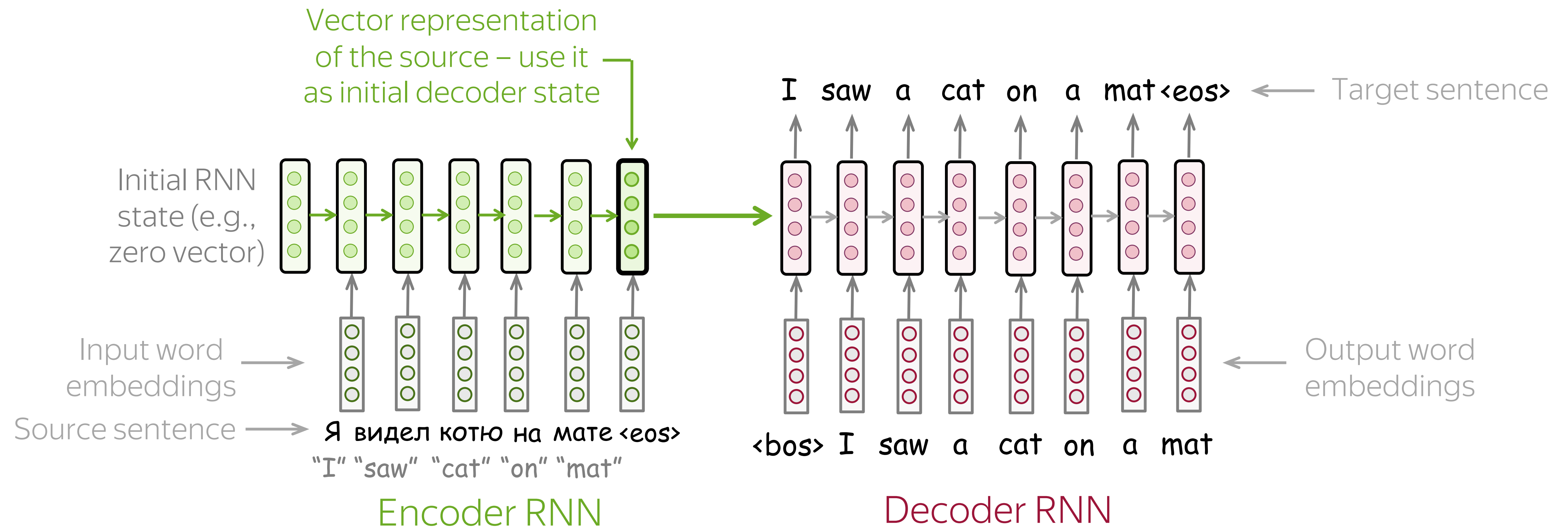
The Simplest Model: RNN Encoder and Decoder



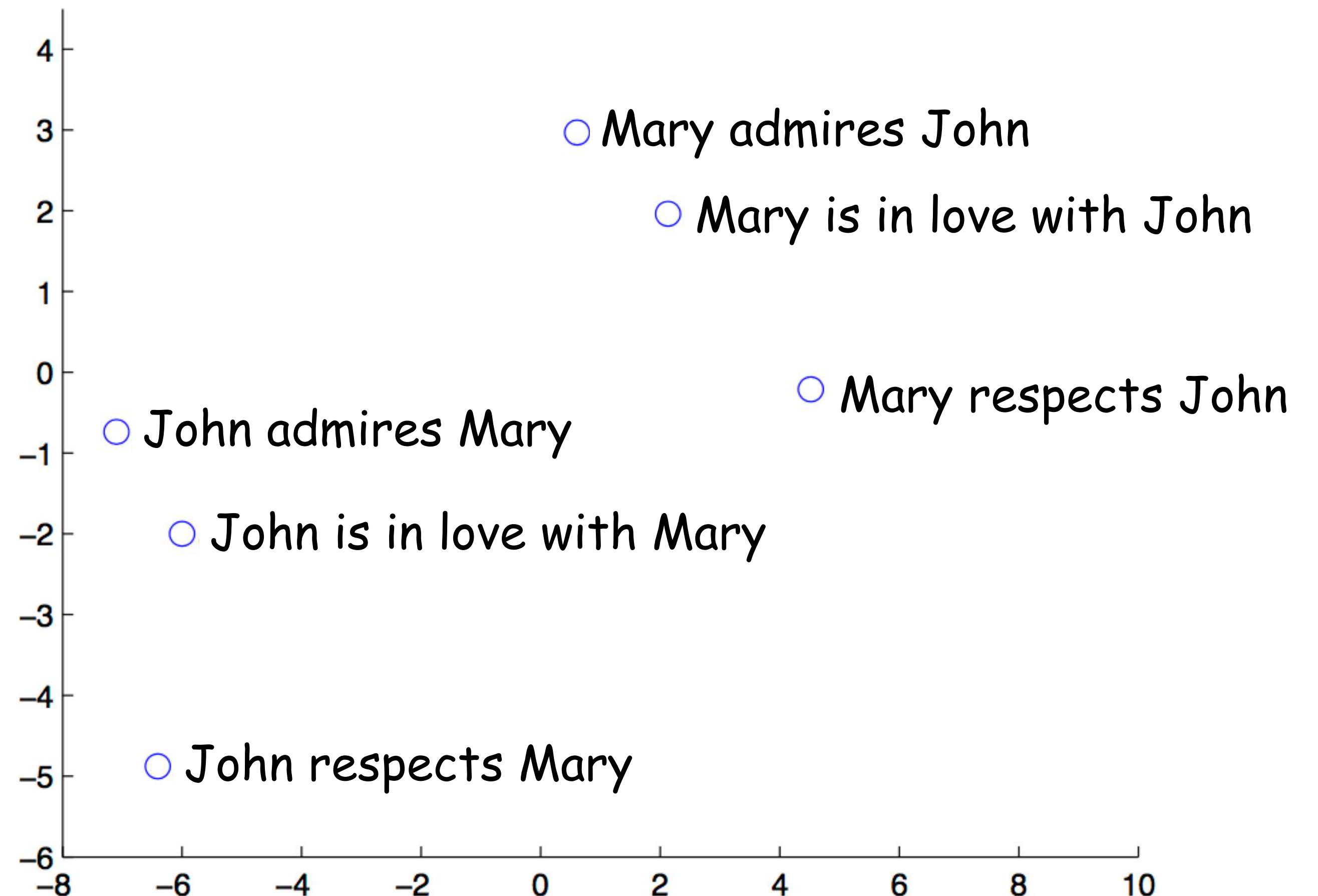
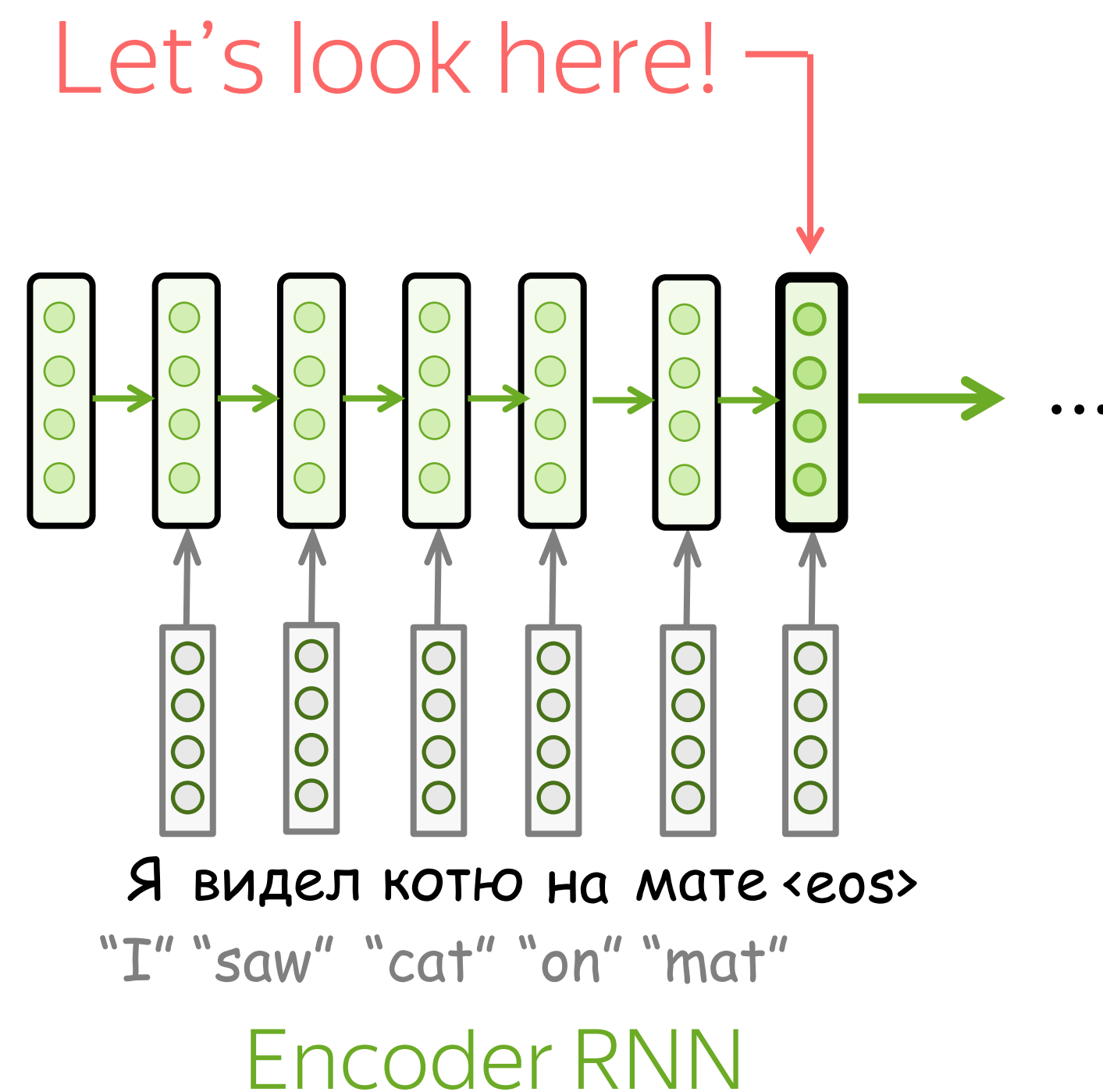
The Simplest Model: RNN Encoder and Decoder



The Simplest Model: RNN Encoder and Decoder

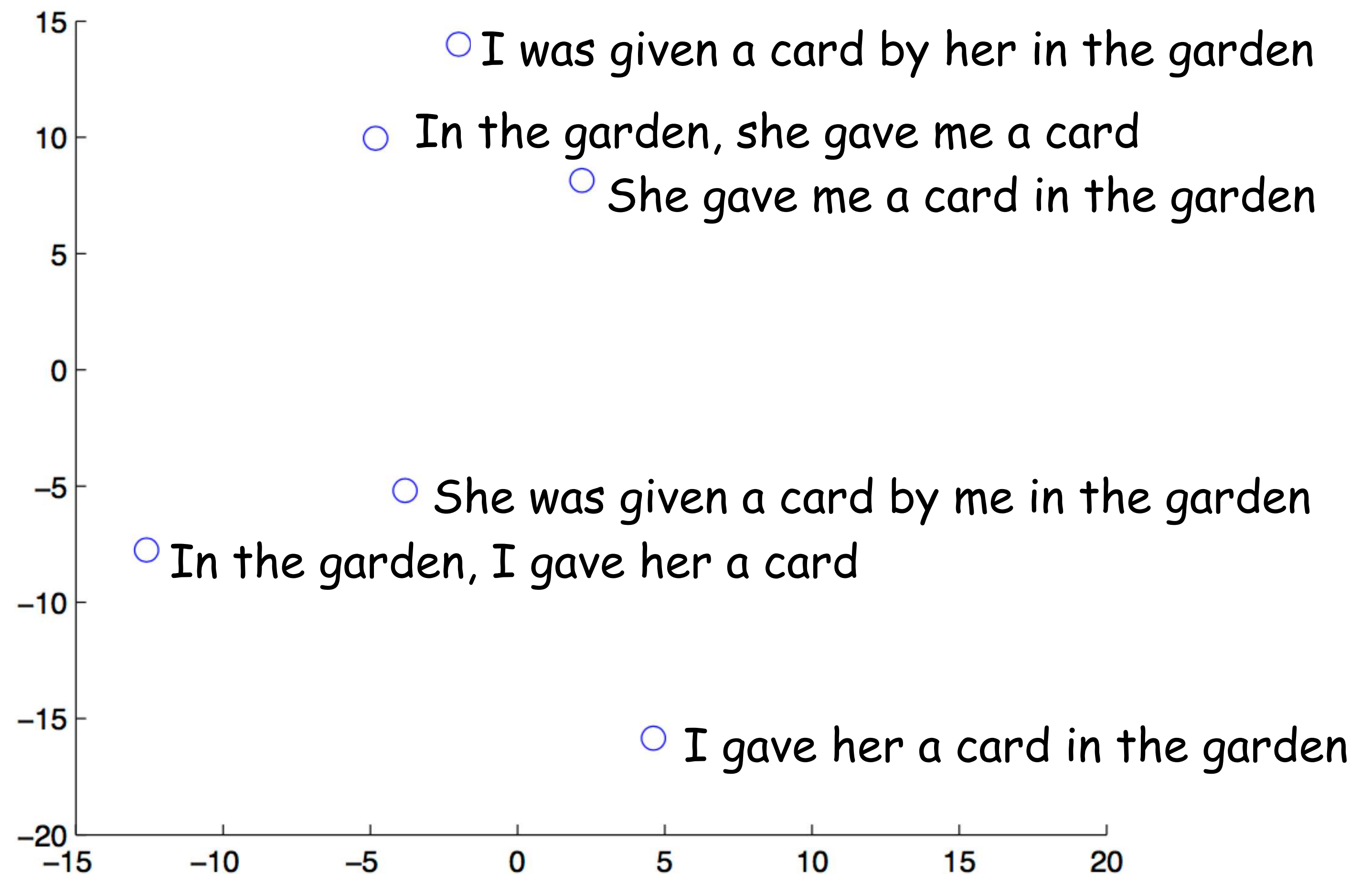
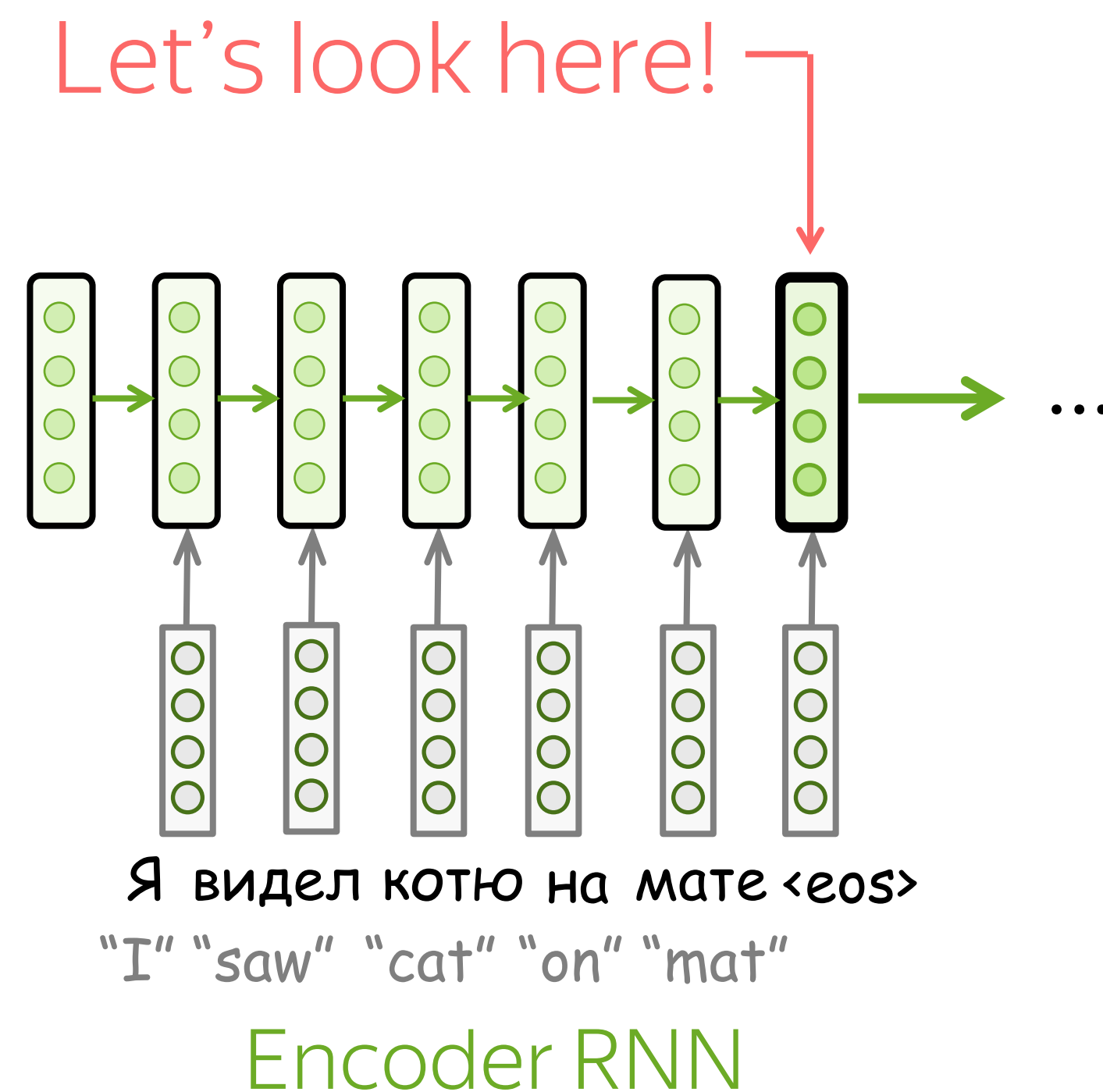


What does final encoder state represent?



The examples are from the paper Sequence to Sequence Learning with Neural Networks

What does final encoder state represent?



The examples are from the paper Sequence to Sequence Learning with Neural Networks

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

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

model parameters

Questions we need to answer

- **modeling**

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

- **learning**

How to find θ ?

- **search**

How to find the argmax?

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 θ ?

- **search**

How to find the argmax?

Cross-Entropy – again!

Source sequence:

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

Target sequence:

I saw a cat on a mat <eos>
previous tokens we want the model to predict this

← one training example
← one step for this example

Model prediction: $p(* | \text{I saw a, Я ... <eos>})$



cat

Target



Loss = $-\log(p(\text{cat})) \rightarrow \min$



decrease

increase

decrease

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 θ ?

- **search**

How to find the argmax?

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 θ ?

- **search**

How to find the argmax?

How To Generate a Translation?

$$y' = \arg \max_y p(y|x) = \arg \max_y \prod_{t=1}^n p(y_t | y_{<t}, x)$$

Greedy Decoding

$$y' = \arg \max_y p(y|x) = \arg \max_y \prod_{t=1}^n p(y_t | y_{<t}, x)$$

Straightforward:

- **greedy** - at each step, pick token with the highest probability

Greedy Decoding

$$y' = \arg \max_y p(y|x) = \arg \max_y \prod_{t=1}^n p(y_t | y_{<t}, x)$$

Straightforward:

- **greedy** - at each step, pick token with the highest probability

Wait a minute...

Greedy Decoding

$$y' = \arg \max_y p(y|x) = \arg \max_y \prod_{t=1}^n p(y_t | y_{<t}, x)$$

Straightforward:

- **greedy** - at each step, pick token with the highest probability

$$\arg \max_y \prod_{t=1}^n p(y_t | y_{<t}, x) \neq \prod_{t=1}^n \arg \max_{y_t} p(y_t | y_{<t}, x) \quad \text{- this is bad!}$$

Beam Search

- At each step, keep several best hypotheses

Beam Search

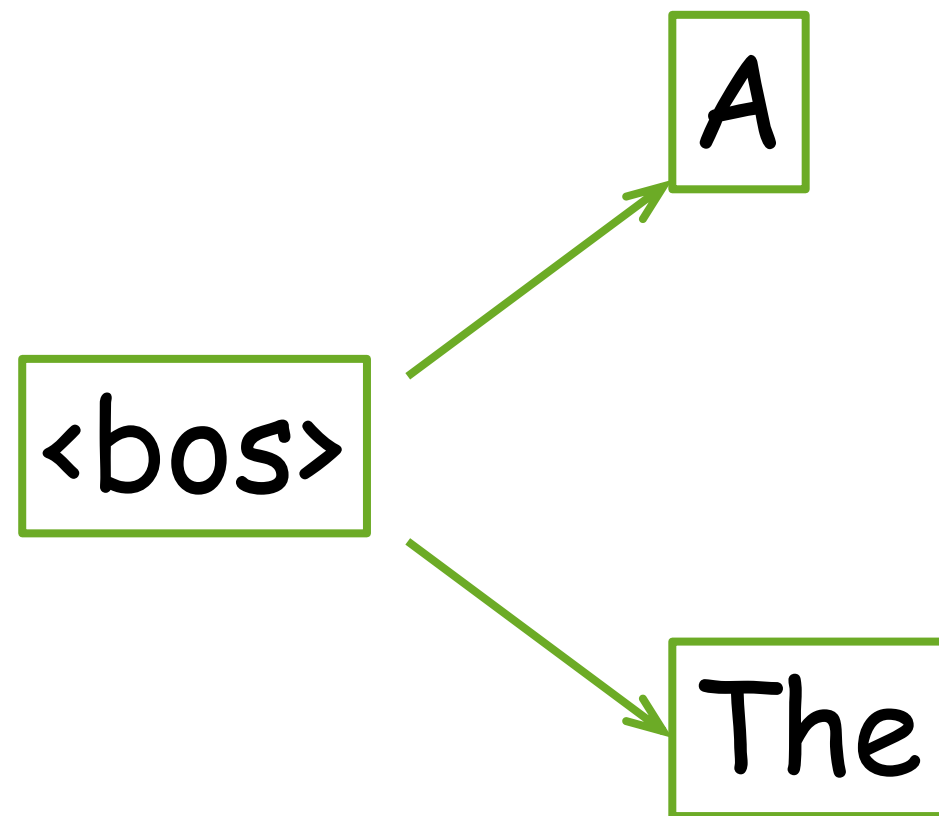
- At each step, keep several best hypotheses

<bos>

Start with the begin of sentence token or with an empty sequence

Beam Search

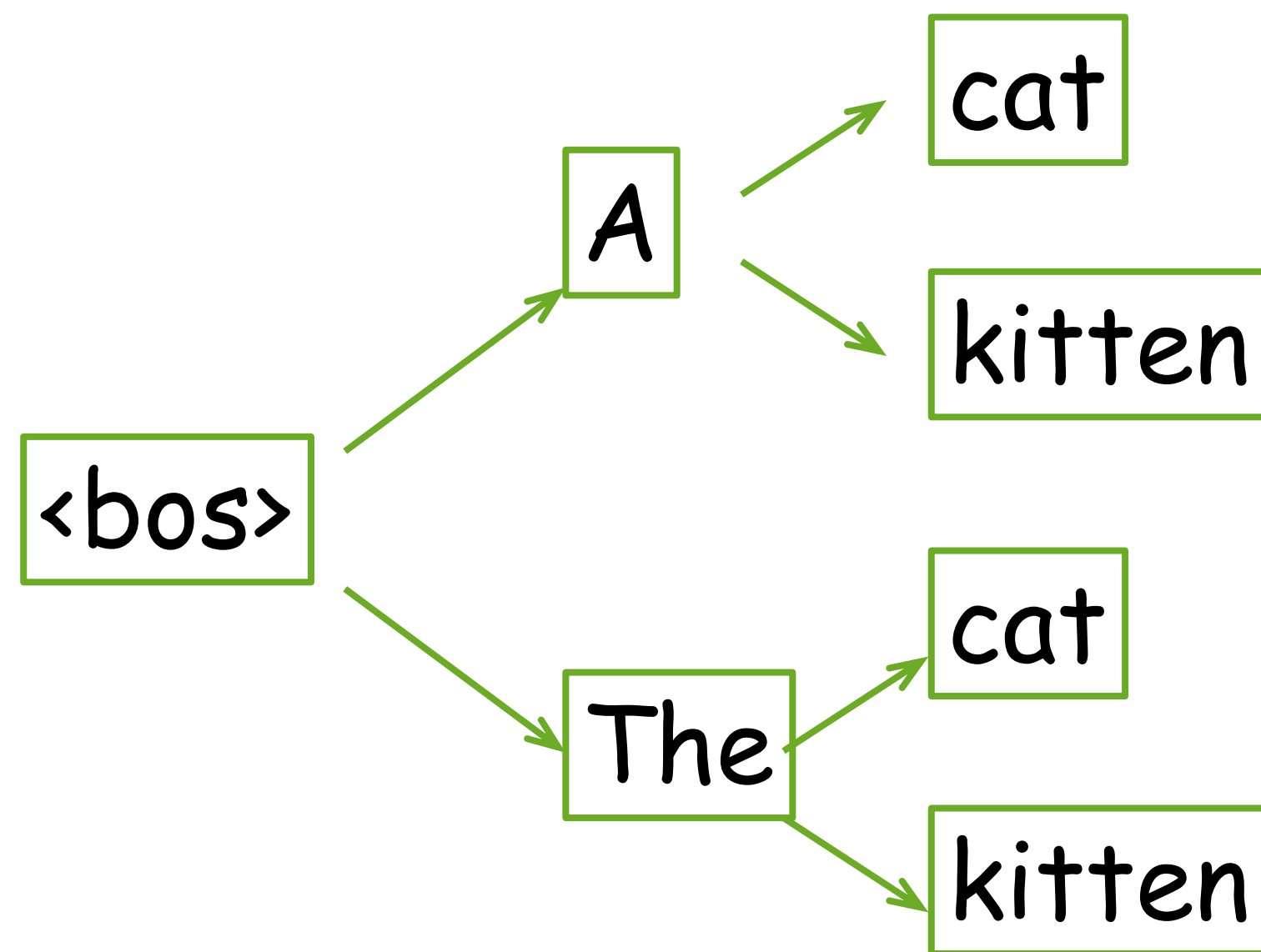
- At each step, keep several best hypotheses



Generate: `beam_size` most probable tokens

Beam Search

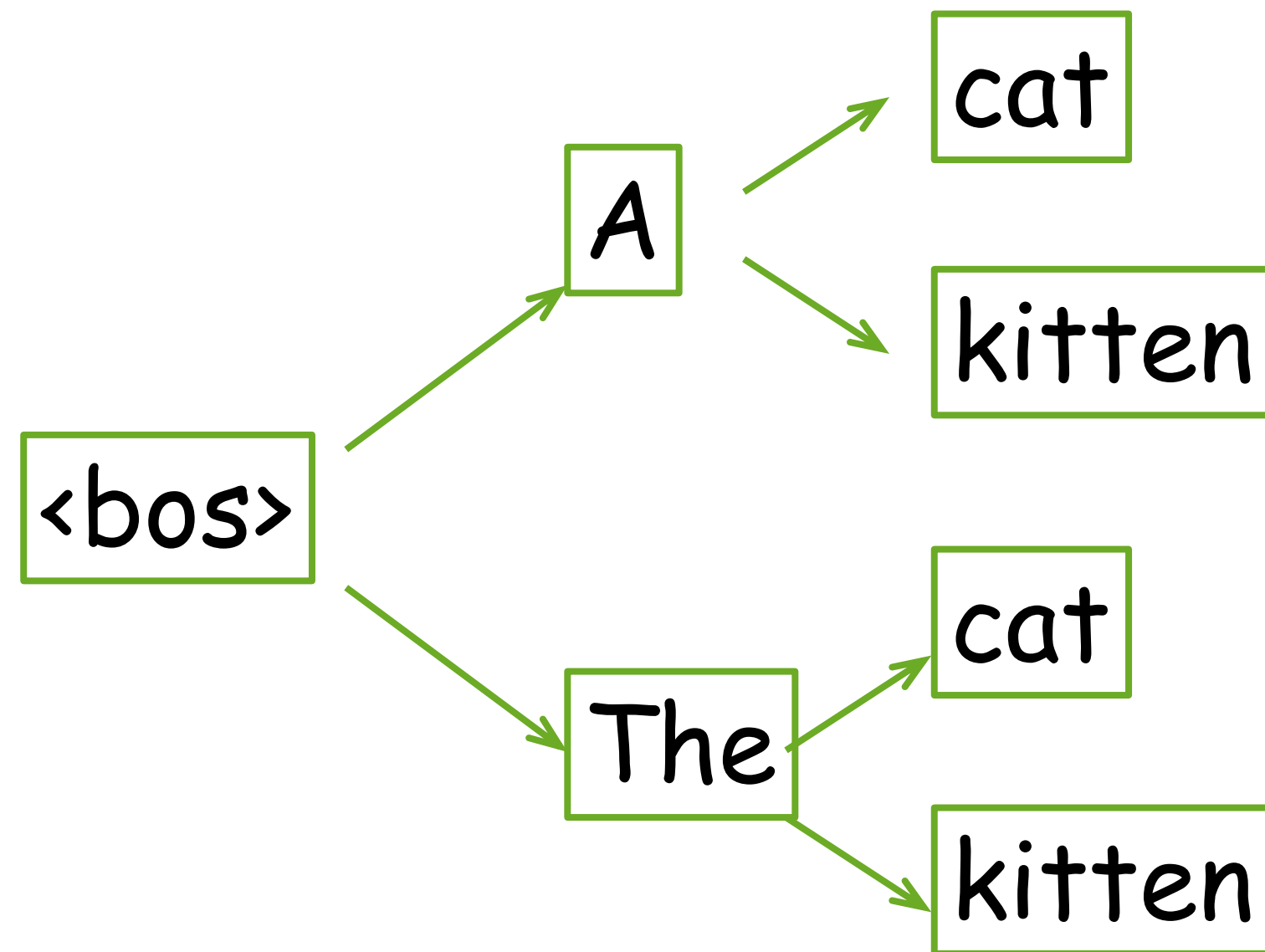
- At each step, keep several best hypotheses



Generate: `beam_size` most probable tokens

Beam Search

- At each step, keep several best hypotheses

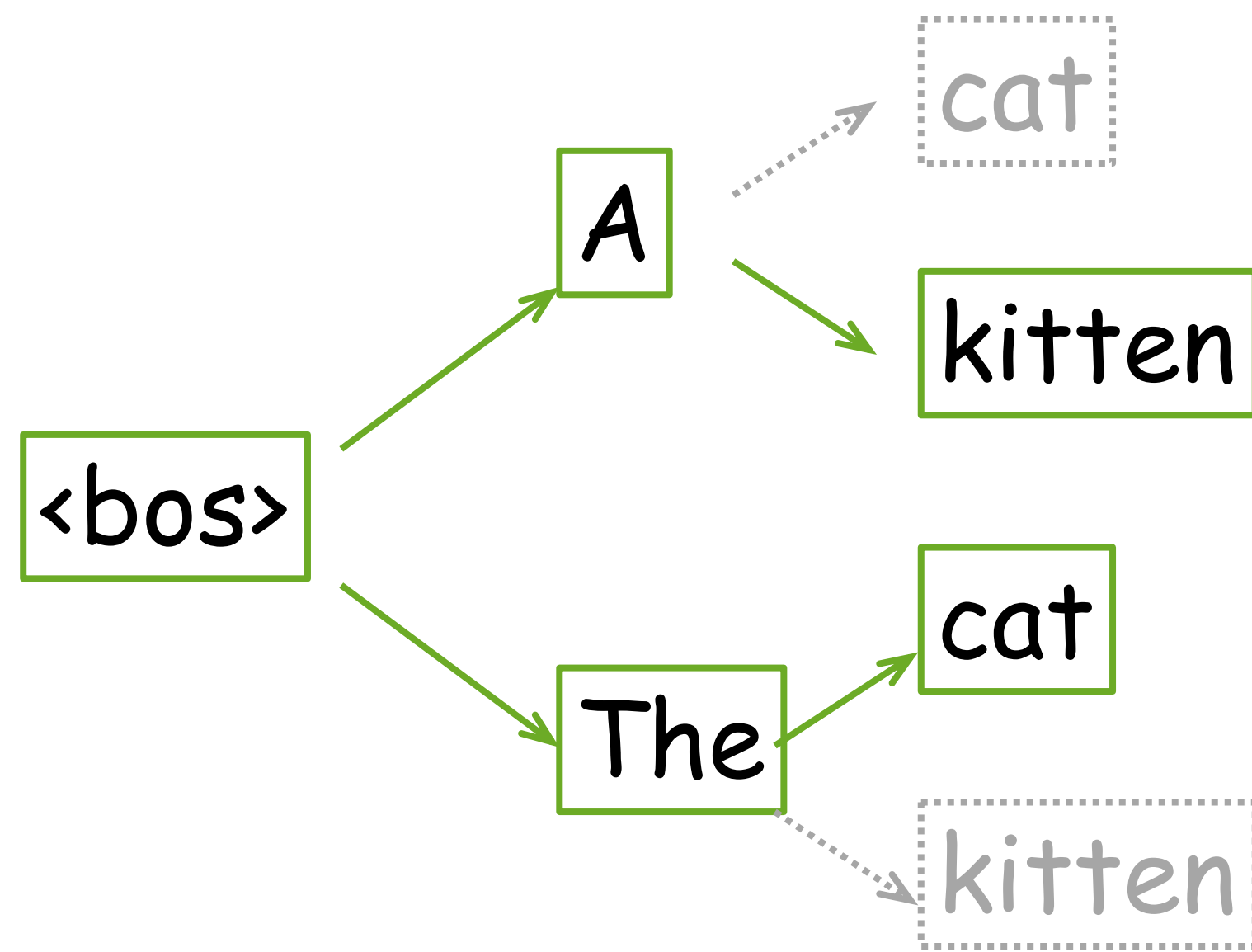


Look at probabilities: $P(\text{The cat}) > P(\text{A kitten}) > P(\text{A cat}) > P(\text{The kitten})$

Top 2 hypos

Beam Search

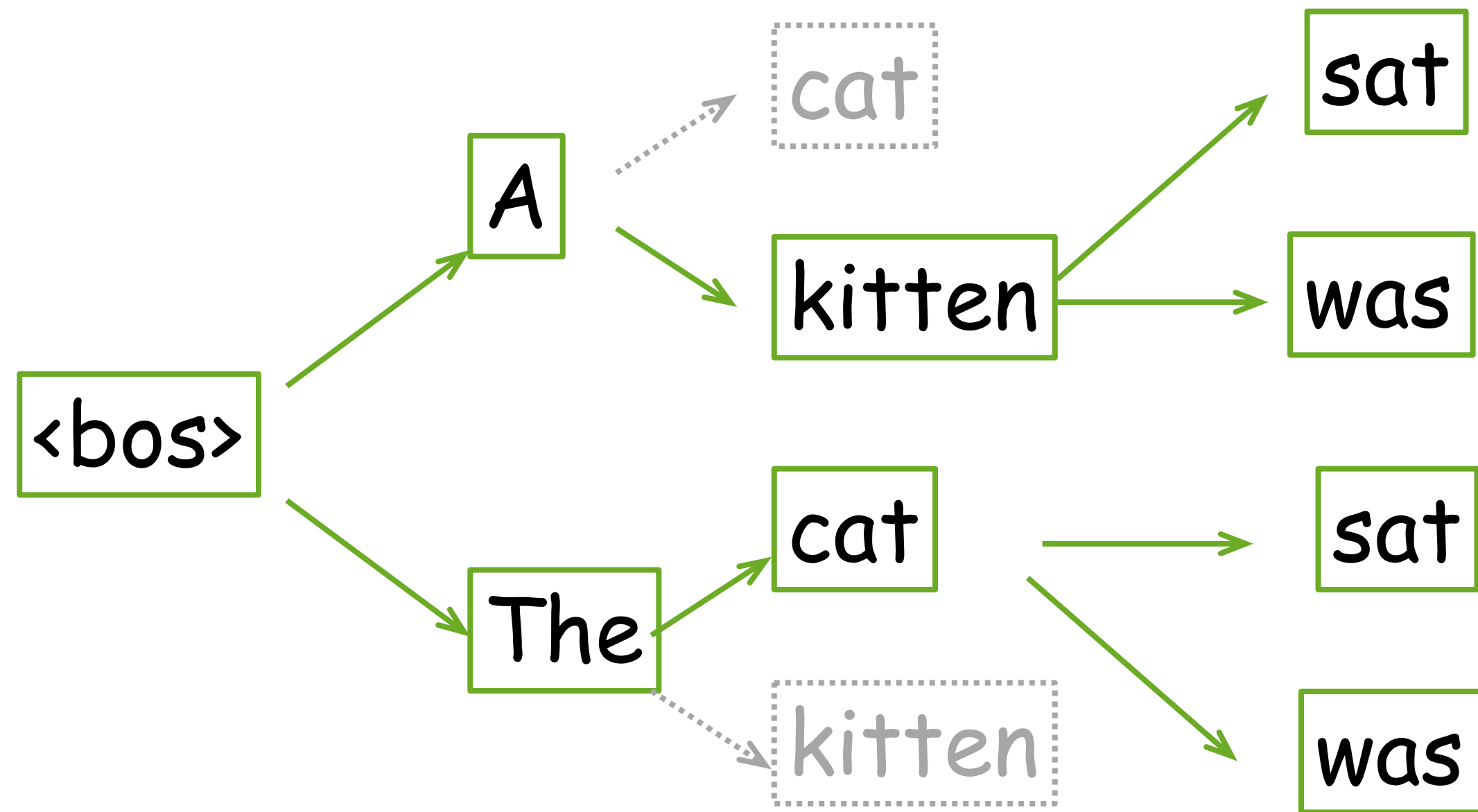
- At each step, keep several best hypotheses



Pick top `beam_size` hypos, terminate the rest

Beam Search

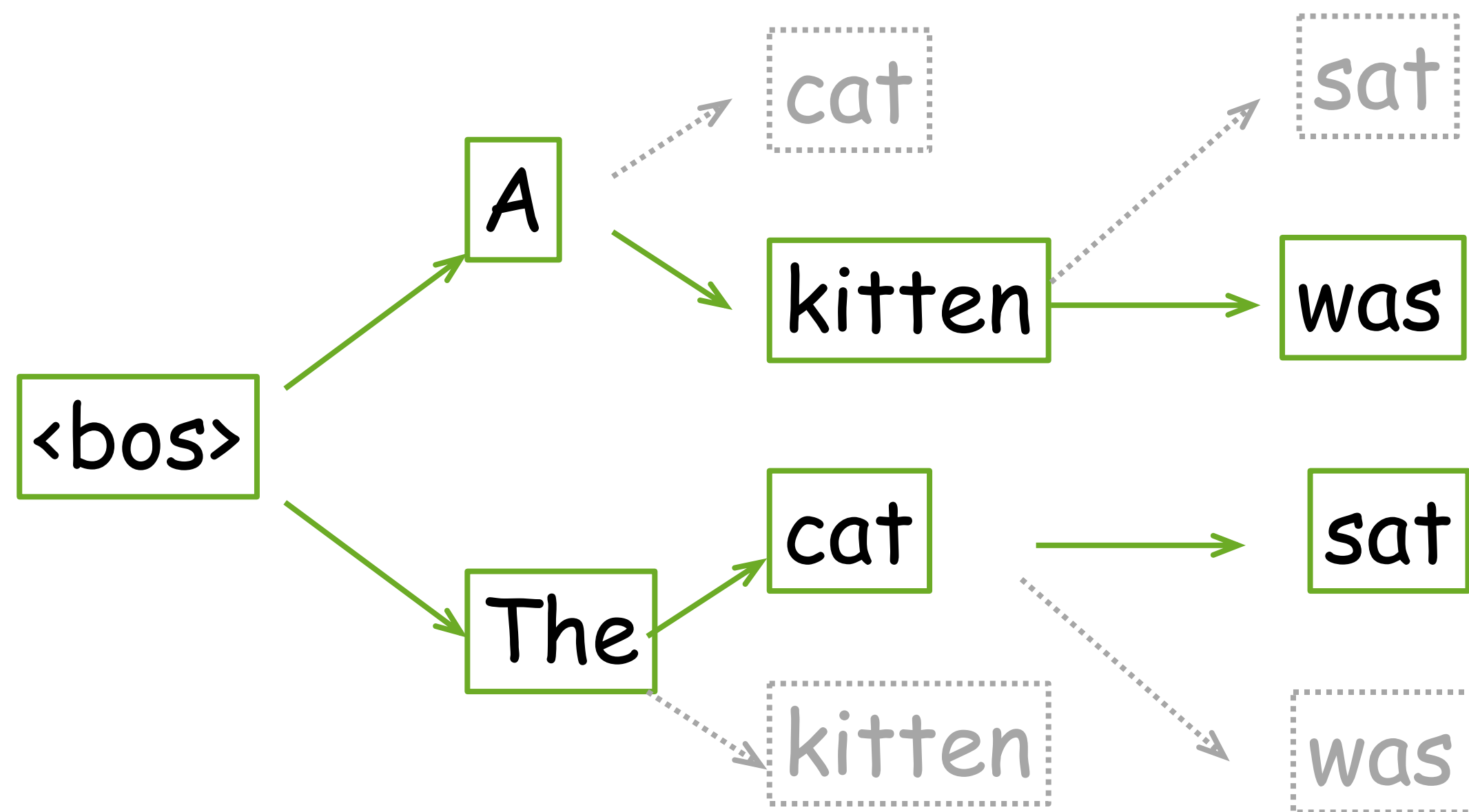
- At each step, keep several best hypotheses



Generate: `beam_size` most probable tokens

Beam Search

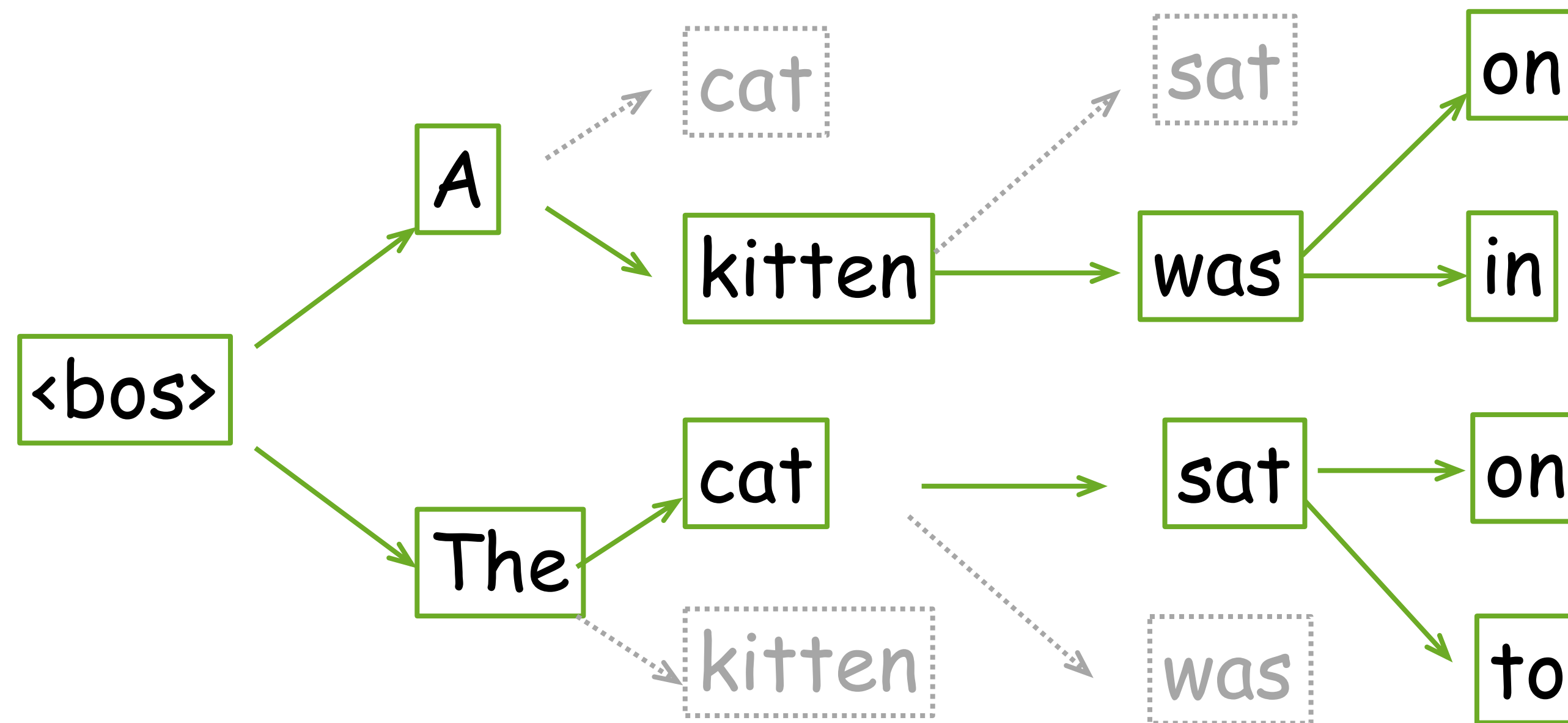
- At each step, keep several best hypotheses



Pick top `beam_size` hypos, terminate the rest

Beam Search

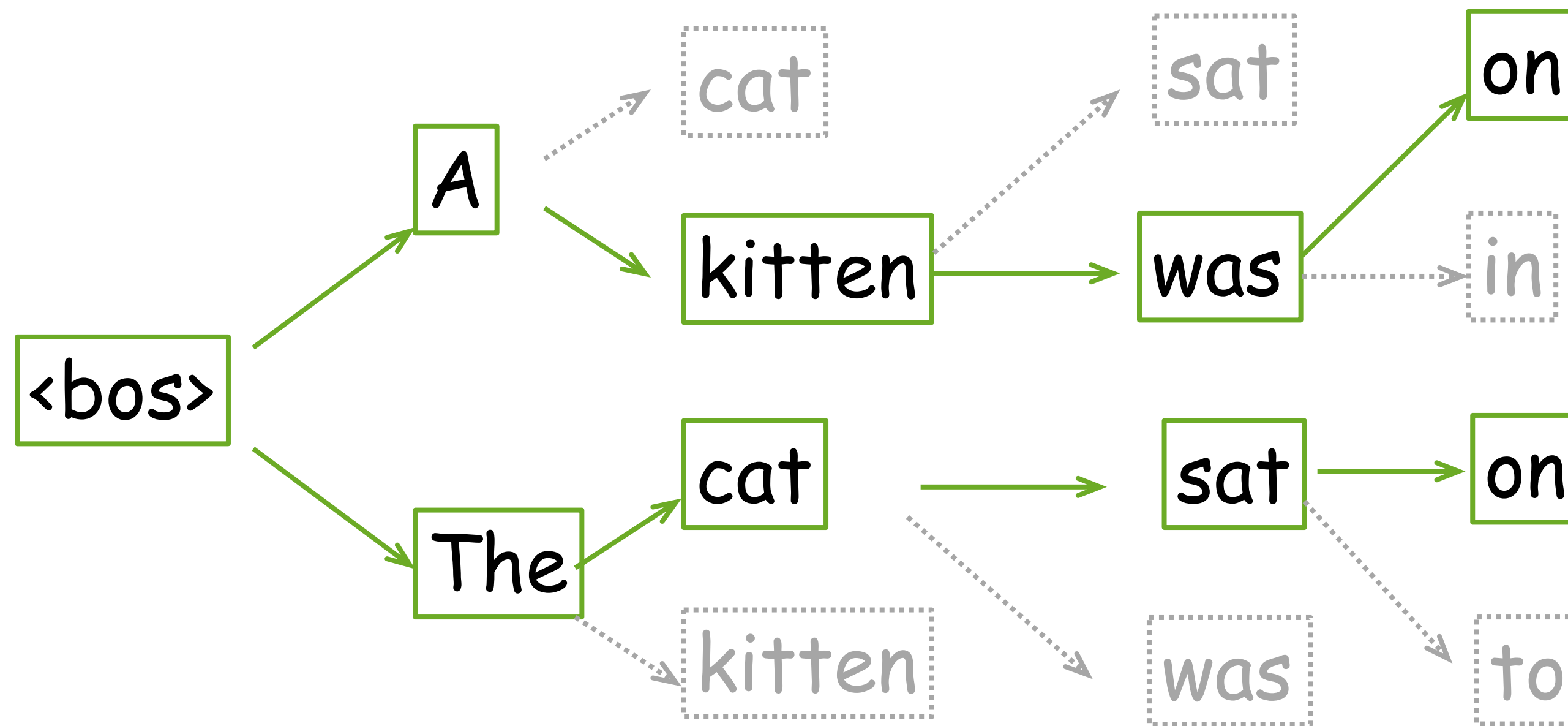
- At each step, keep several best hypotheses



Generate: `beam_size` most probable tokens

Beam Search

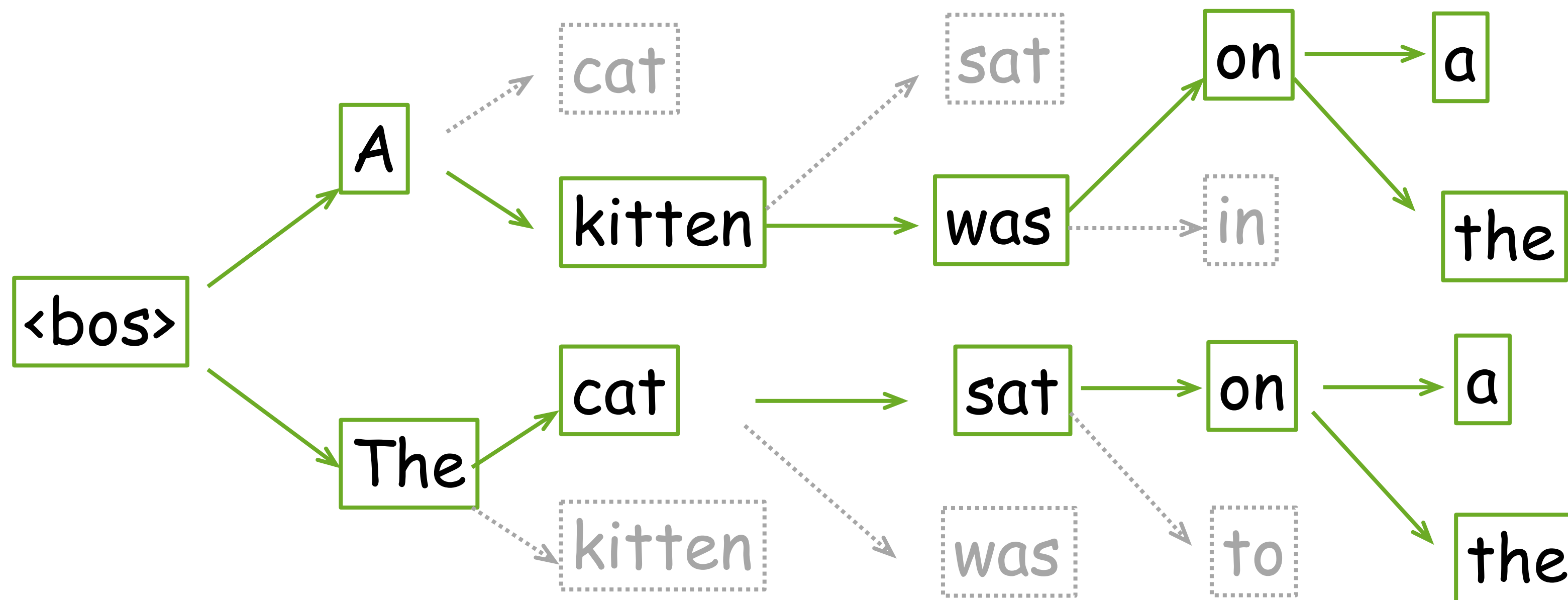
- At each step, keep several best hypotheses



Pick top `beam_size` hypos, terminate the rest

Beam Search

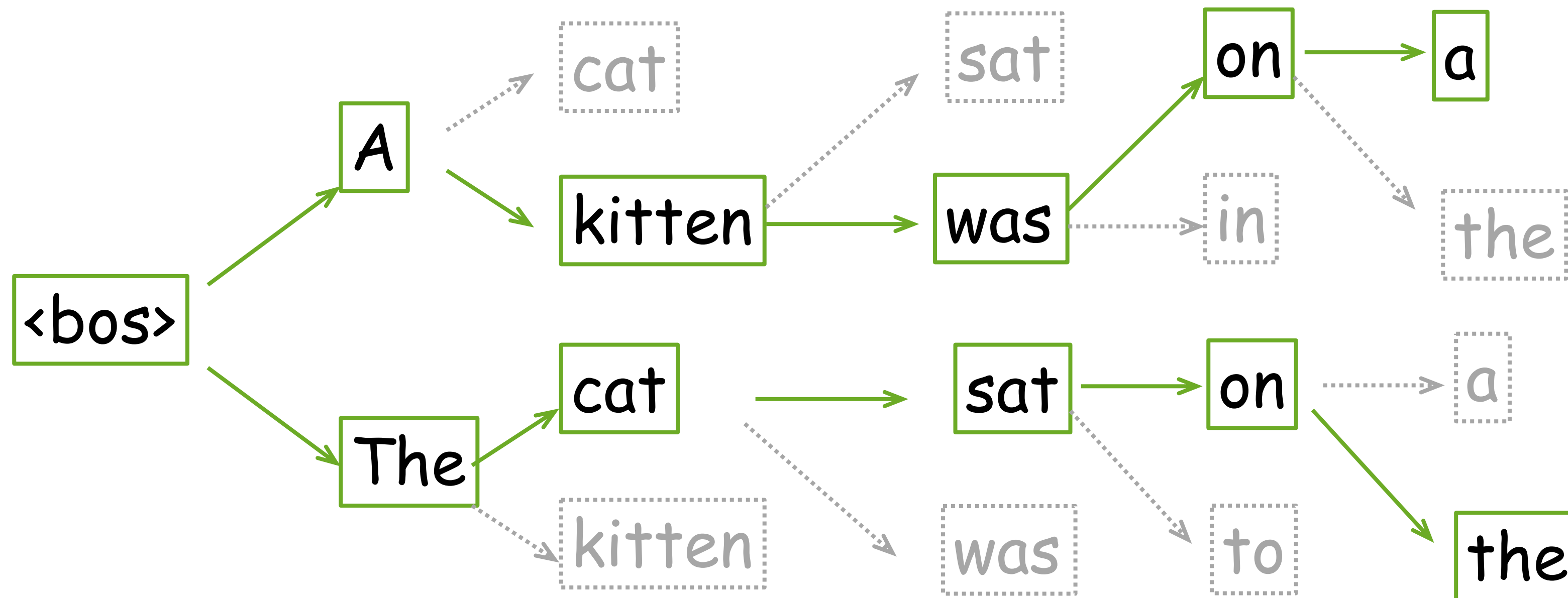
- At each step, keep several best hypotheses



Generate: `beam_size` most probable tokens

Beam Search

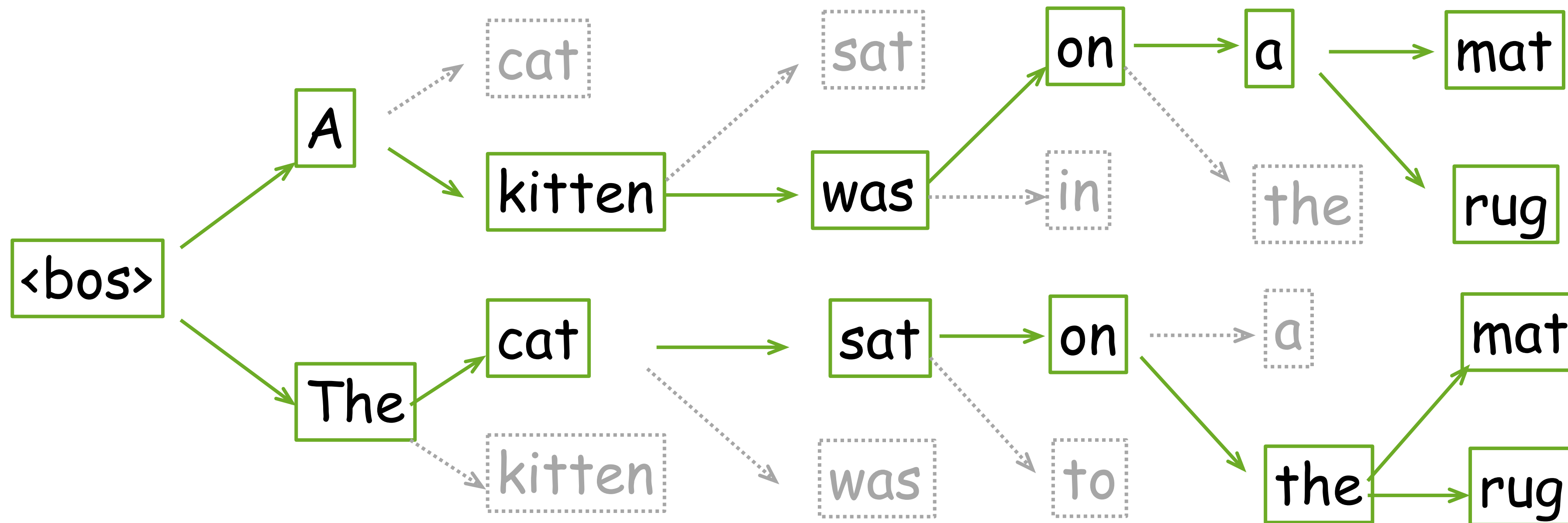
- At each step, keep several best hypotheses



Pick top `beam_size` hypos, terminate the rest

Beam Search

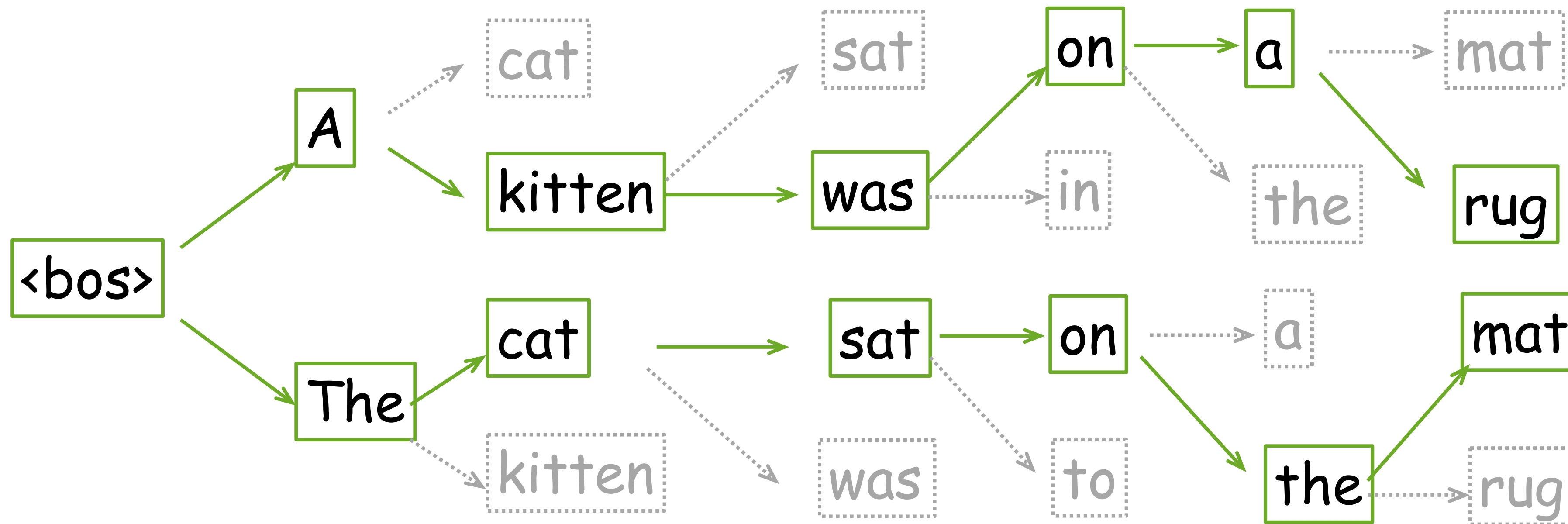
- At each step, keep several best hypotheses



Generate: `beam_size` most probable tokens

Beam Search

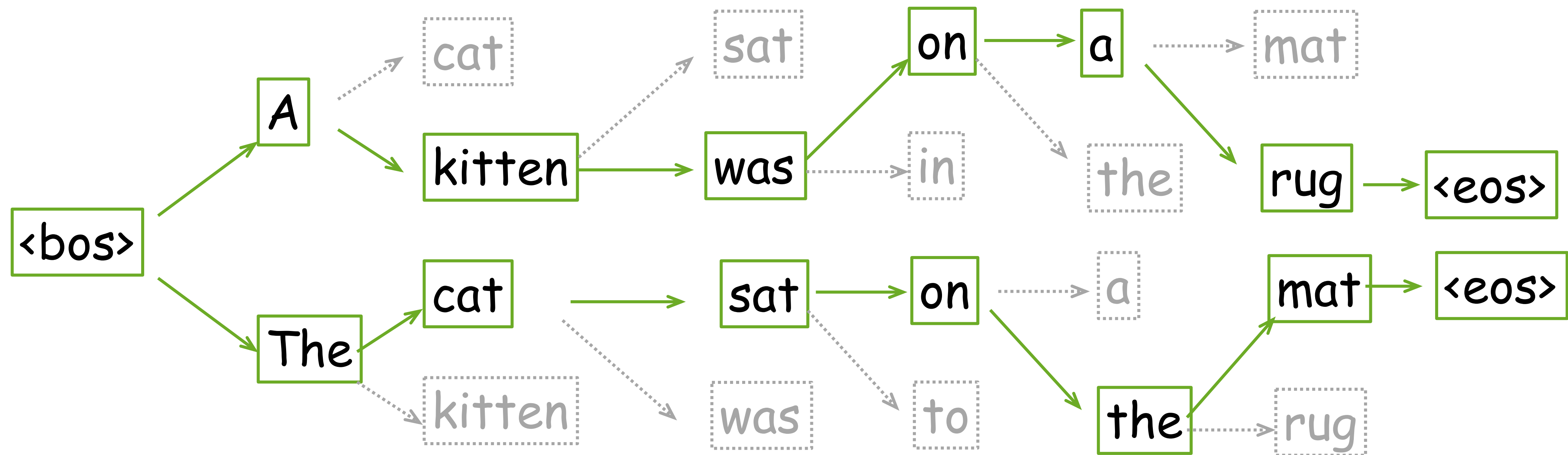
- At each step, keep several best hypotheses



Pick top `beam_size` hypos, terminate the rest

Beam Search

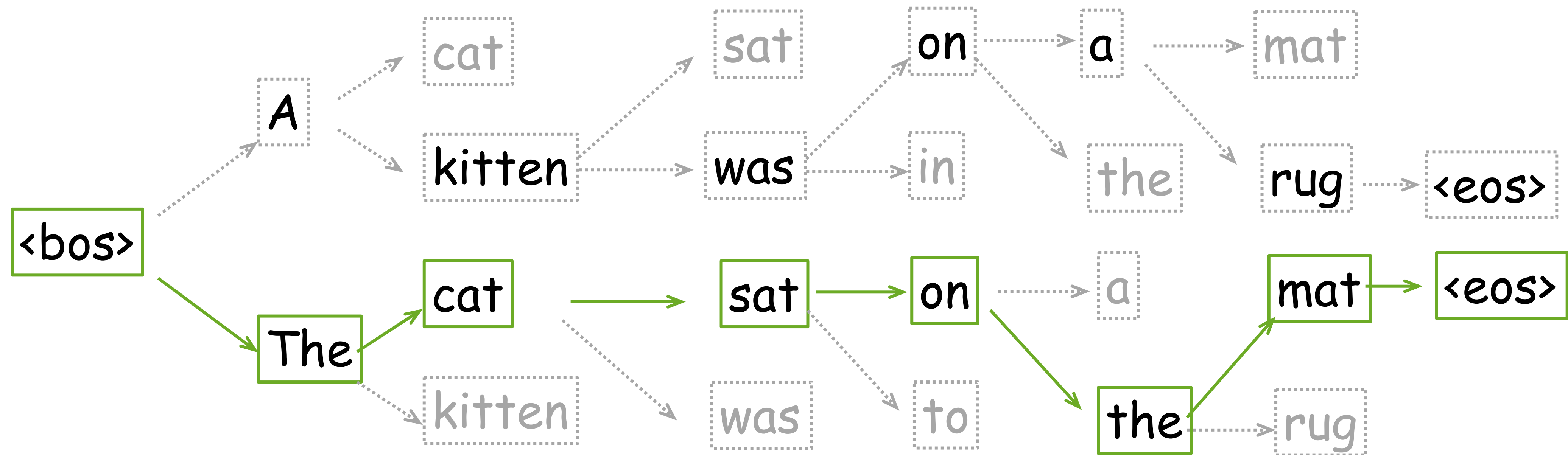
- At each step, keep several best hypotheses



All hypotheses are complete - generation ended

Beam Search

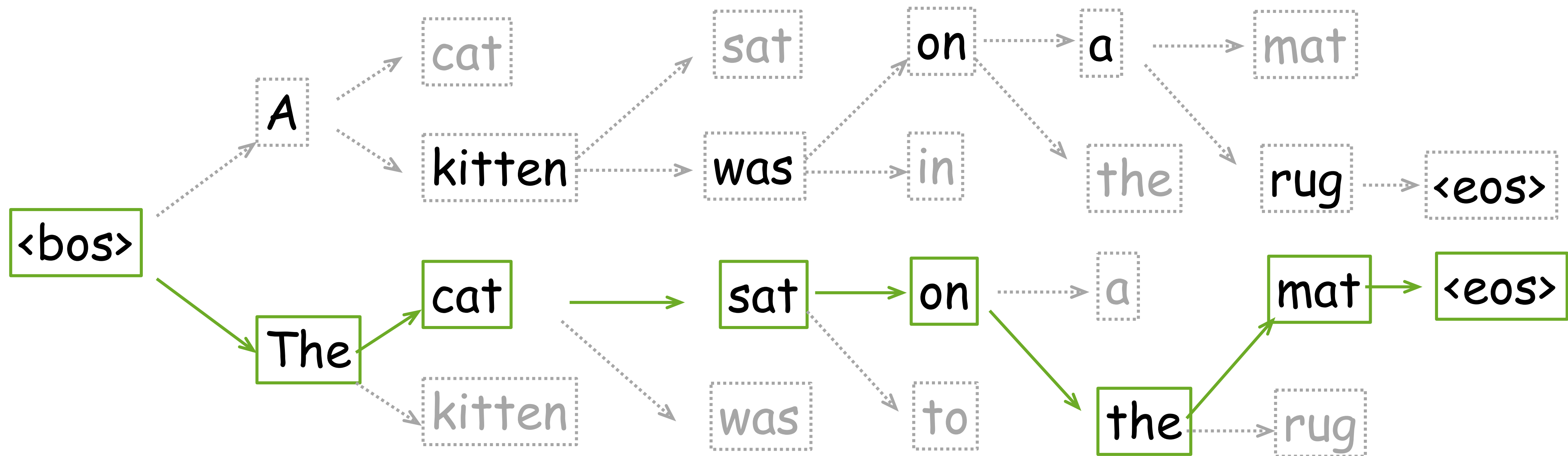
- At each step, keep several best hypotheses



Pick the hypothesis with the highest probability


Beam Search

- At each step, keep several best hypotheses




Result: The cat sat on the mat

What is going to happen:

- Seq2seq Basics
- Attention
- Transformer
- Subword Segmentation: BPE
-  Analysis and Interpretability

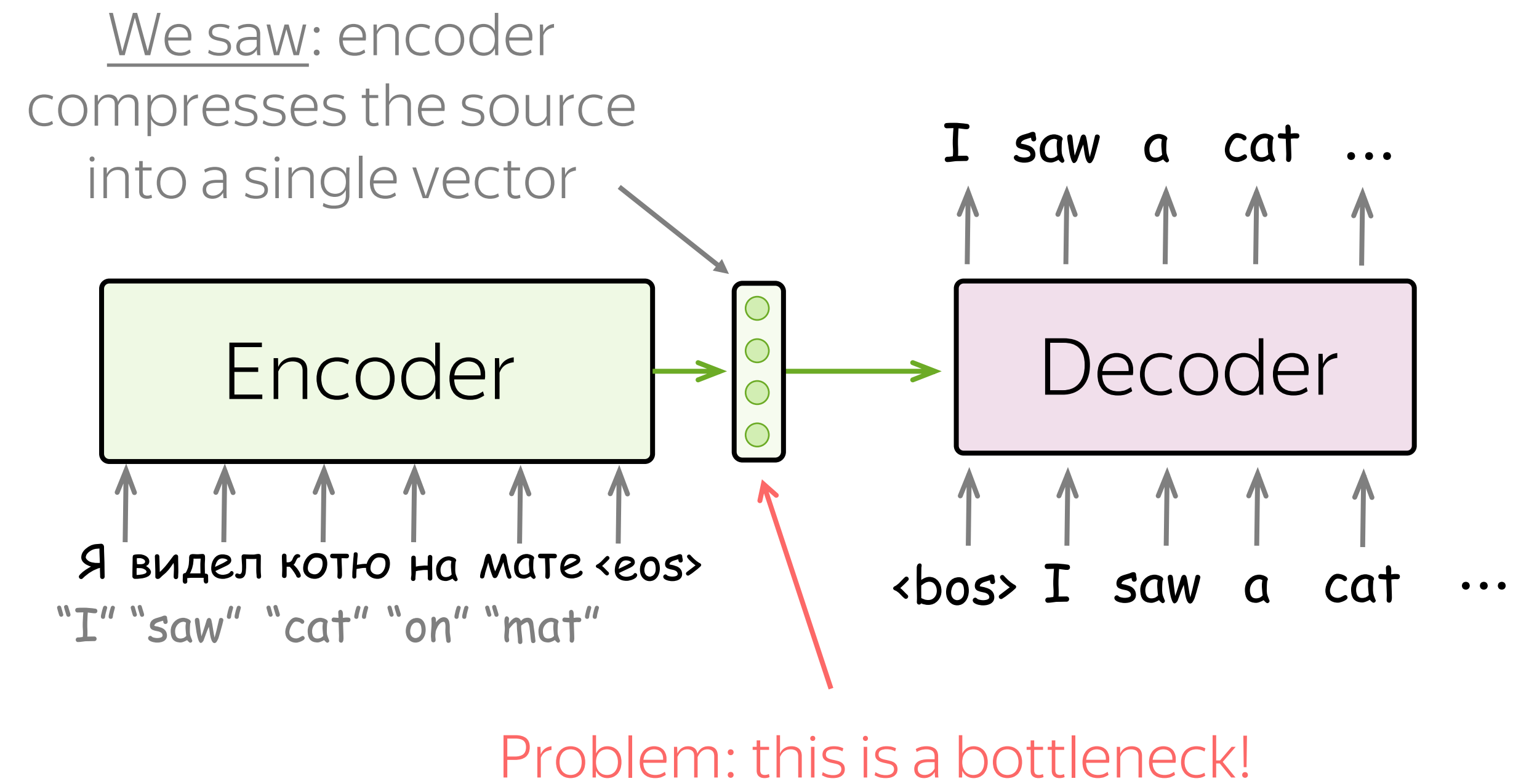
What is going to happen:

- Seq2seq Basics
- Attention →
 - Why do we need it?
 - Attention: High-Level
 - Attention Score Functions
 - Models: Bahdanau vs Luong
- Transformer
- Subword Segmentation: BPE
-  Analysis and Interpretability

The Problem of Fixed Encoder Representation

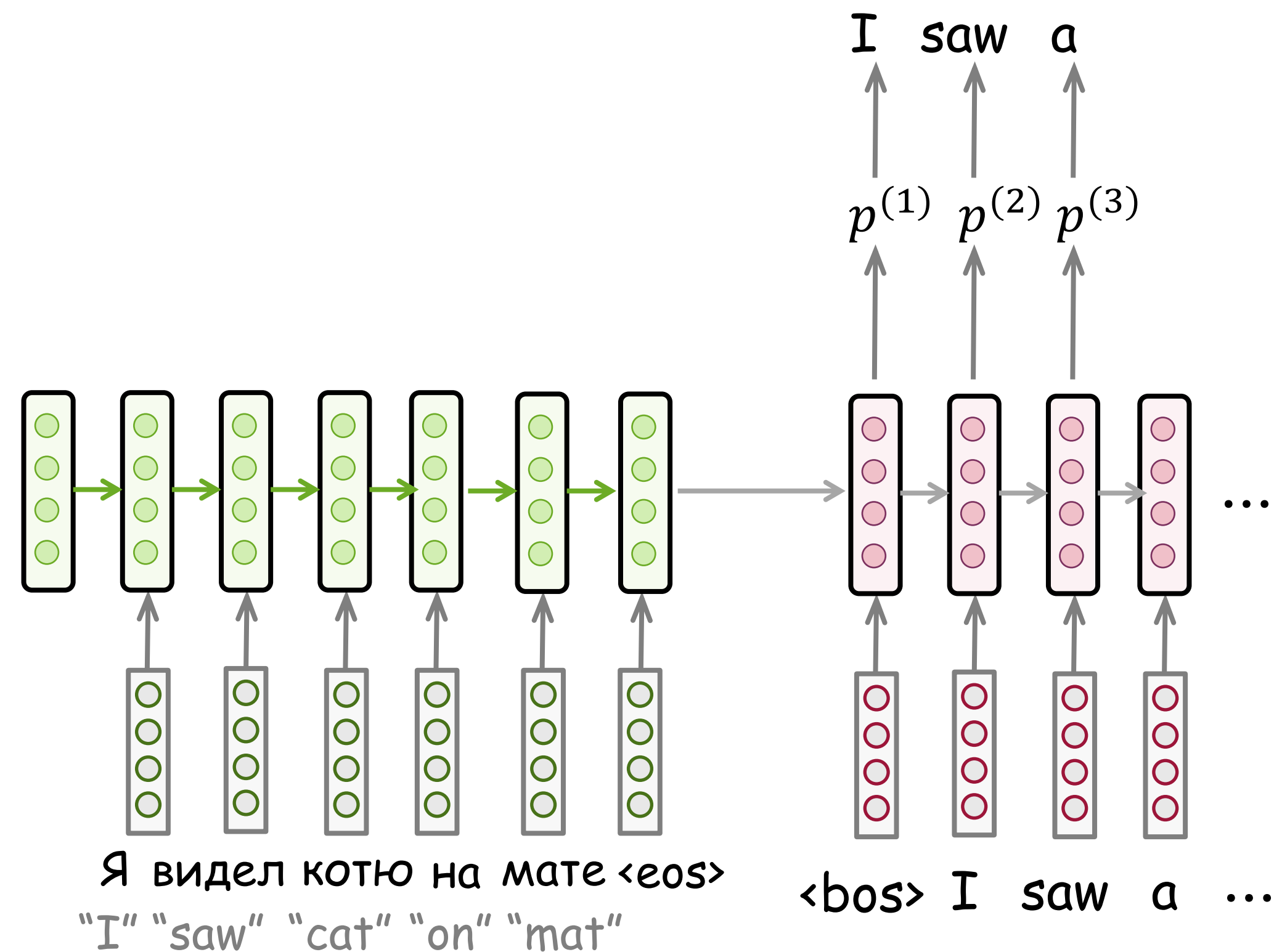
Fixed source representation is bad:

- for **encoder**, it is hard to compress a sentence
- for **decoder**, different information may be needed at different steps



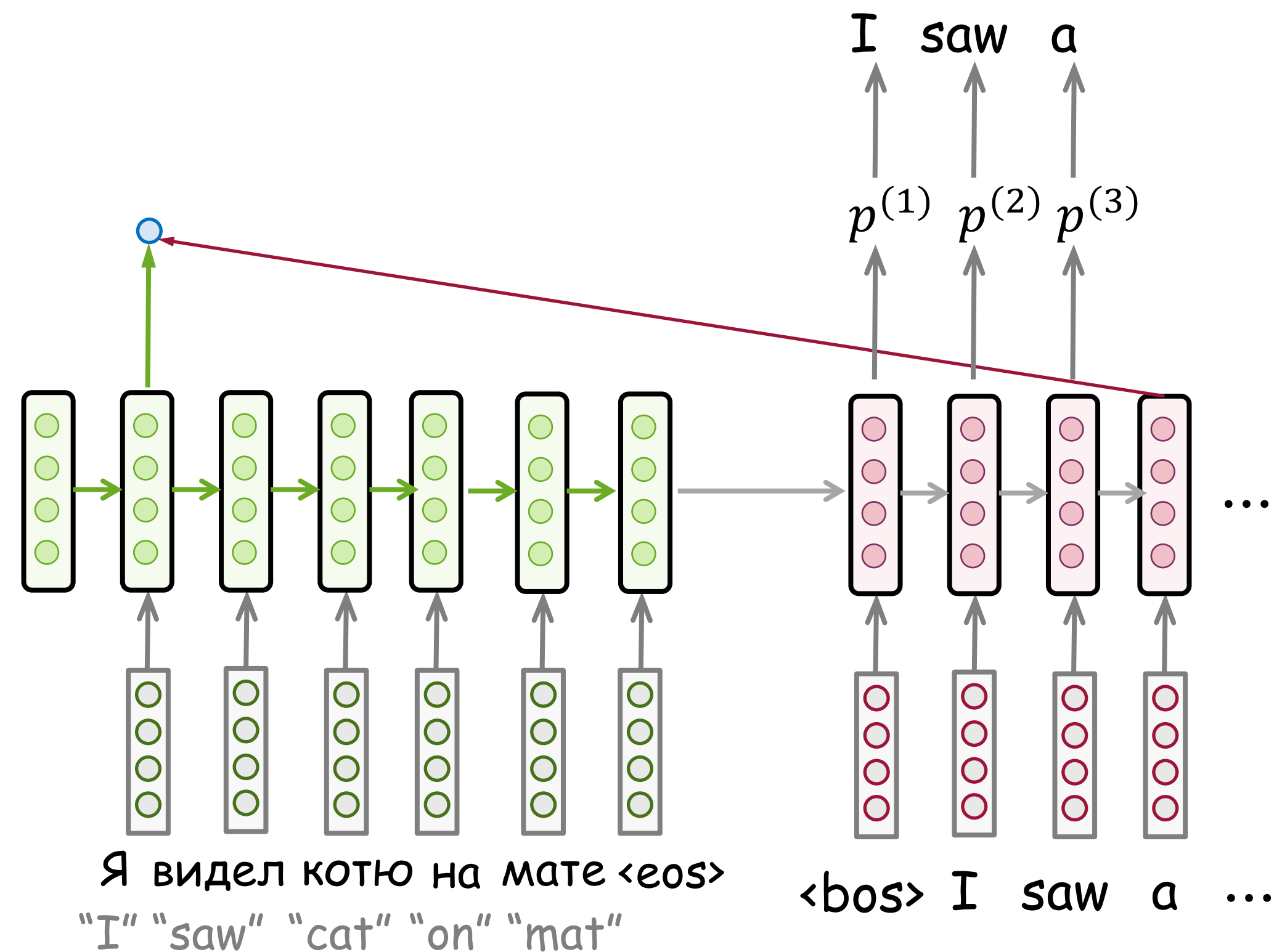
Attention: High-Level View

Attention: At different steps, let a model "focus" on different parts of the input.



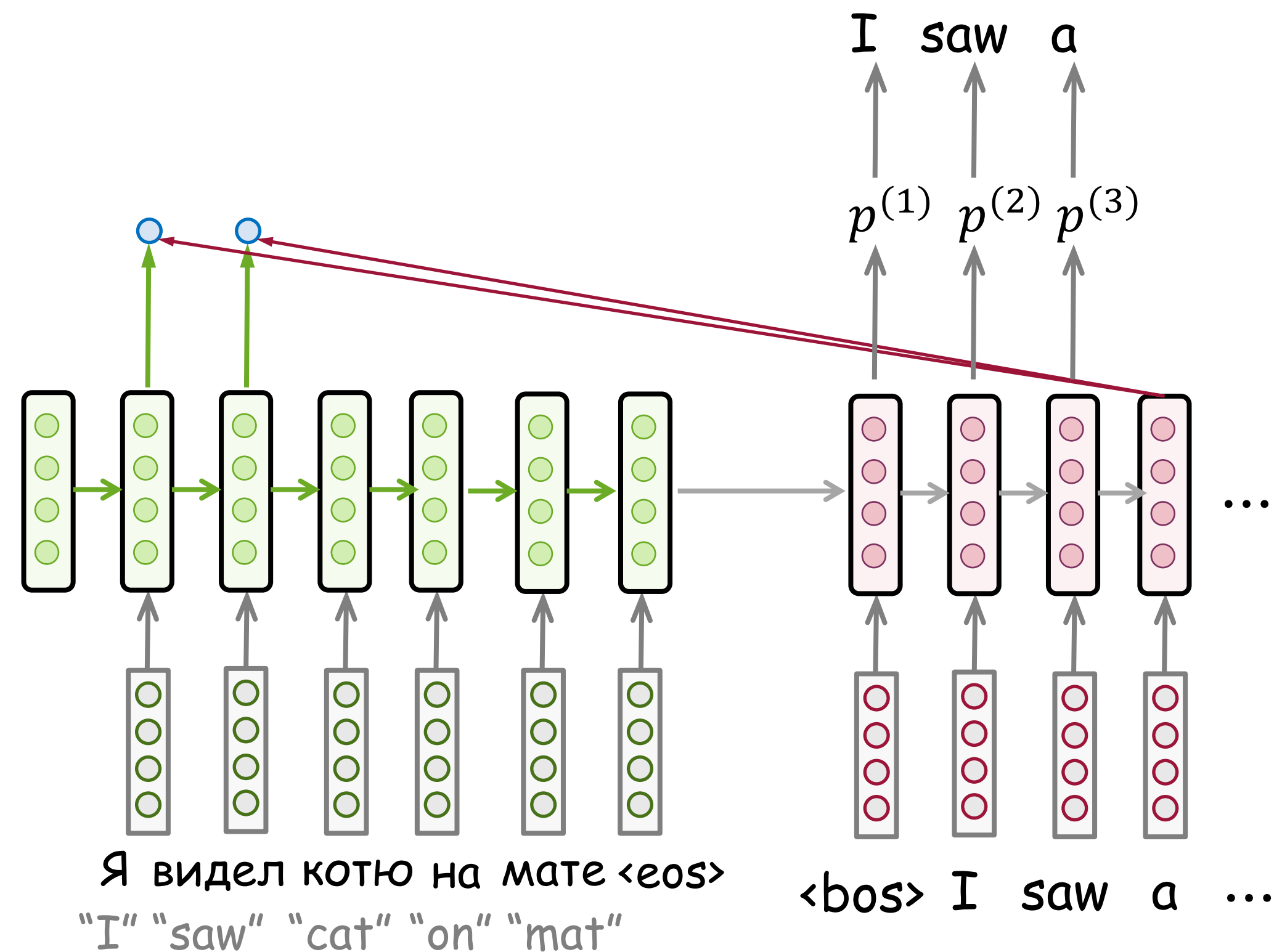
Attention: High-Level View

Attention: At different steps, let a model "focus" on different parts of the input.



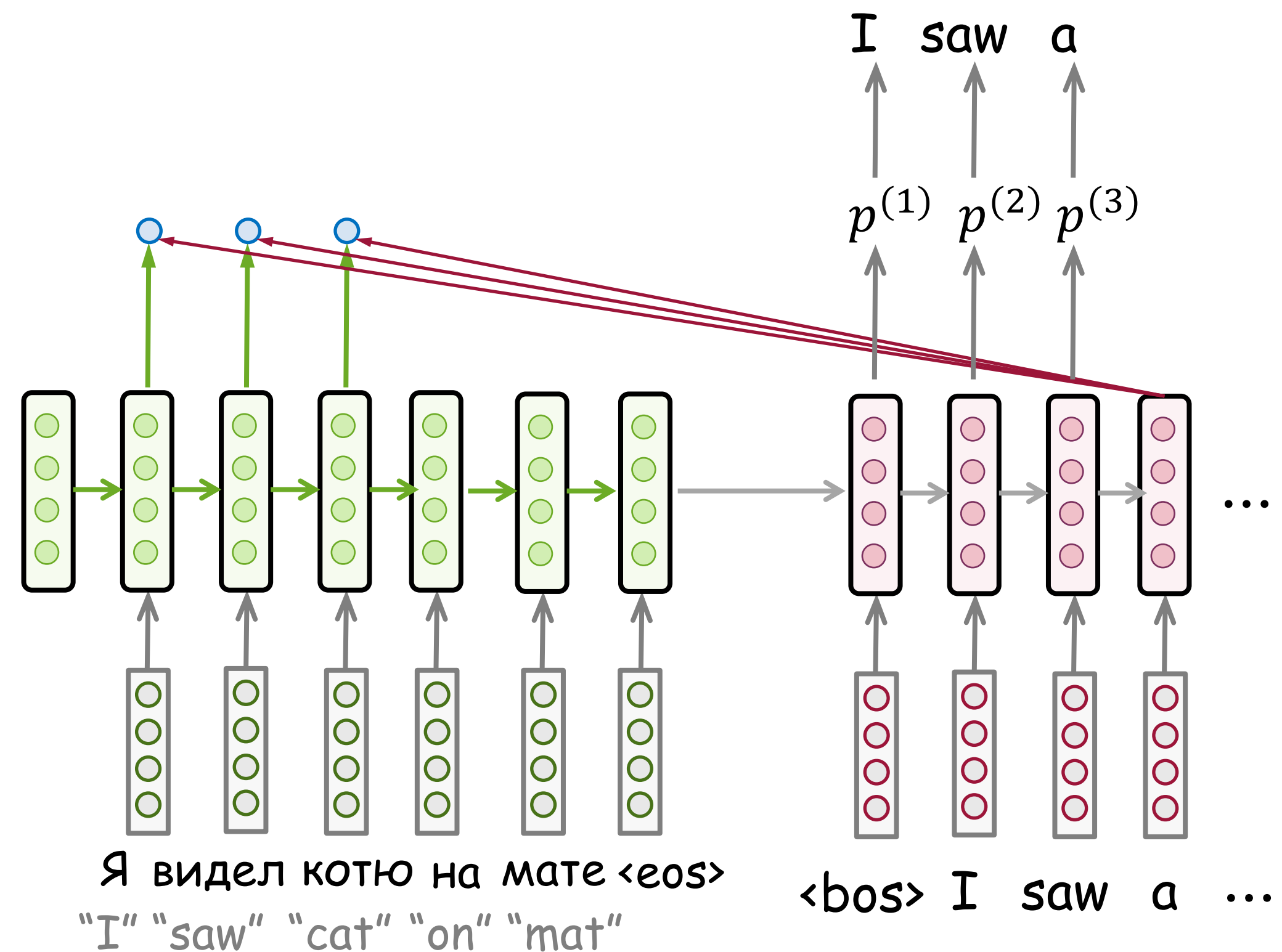
Attention: High-Level View

Attention: At different steps, let a model "focus" on different parts of the input.



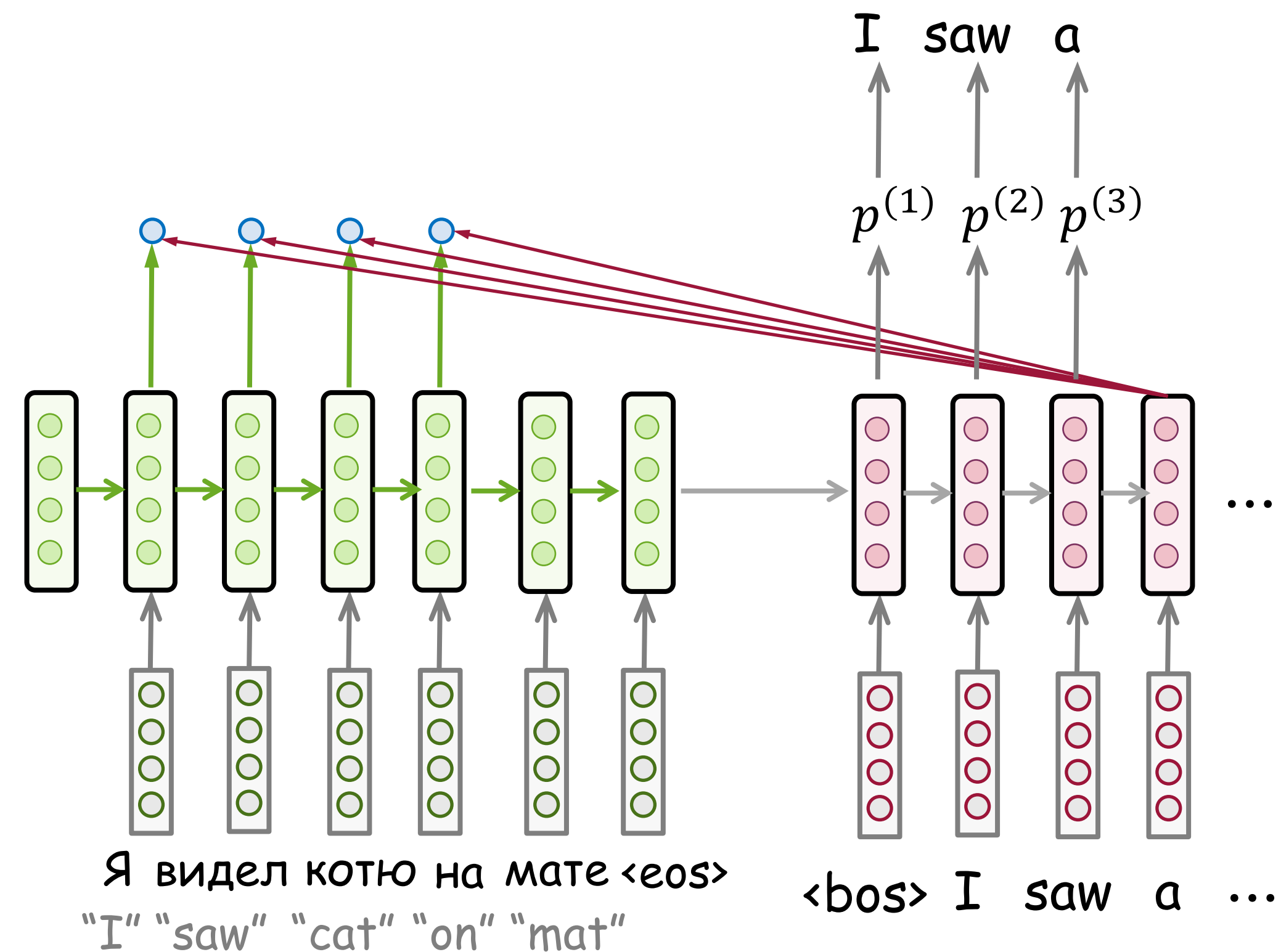
Attention: High-Level View

Attention: At different steps, let a model "focus" on different parts of the input.



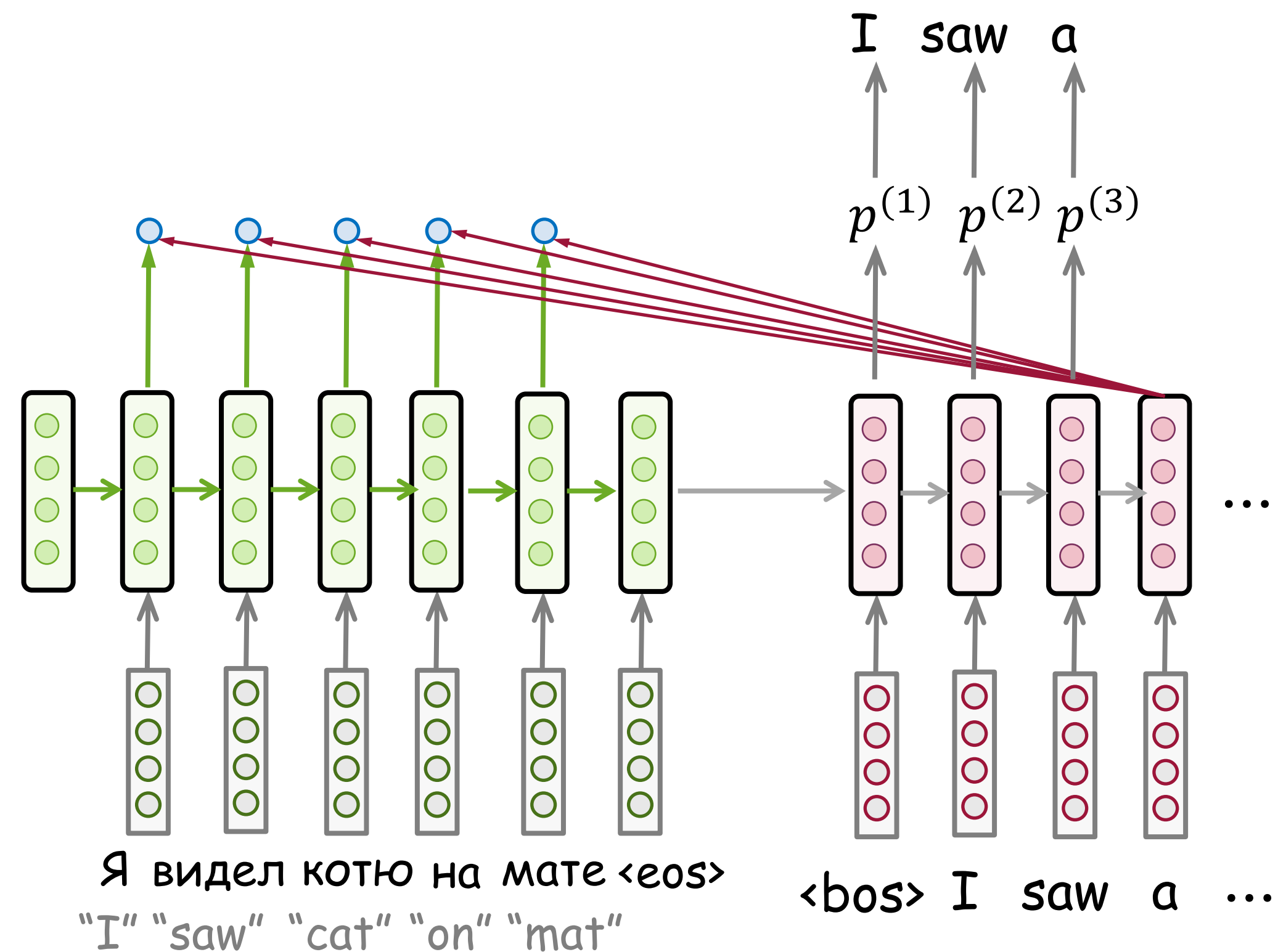
Attention: High-Level View

Attention: At different steps, let a model "focus" on different parts of the input.



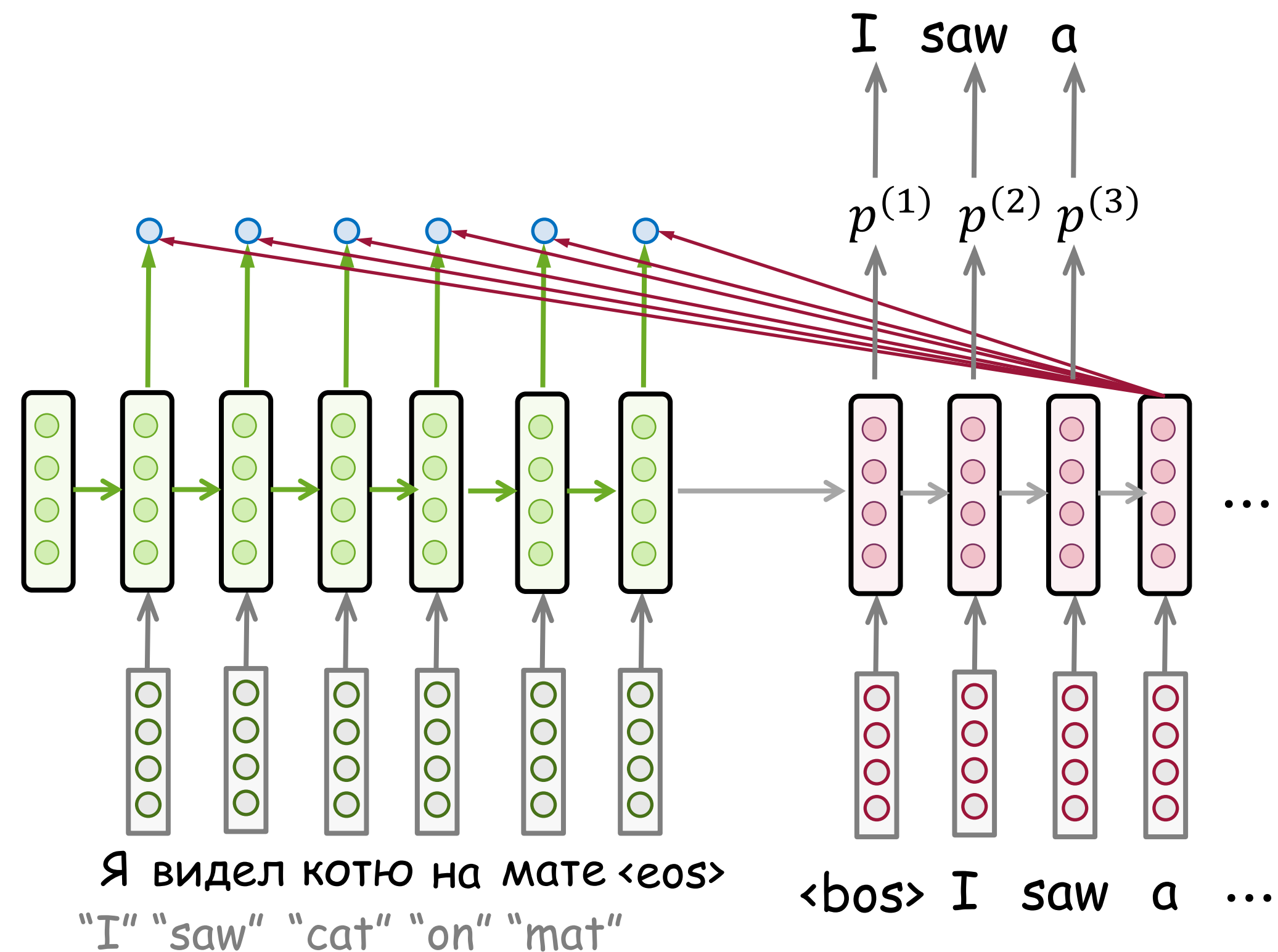
Attention: High-Level View

Attention: At different steps, let a model "focus" on different parts of the input.



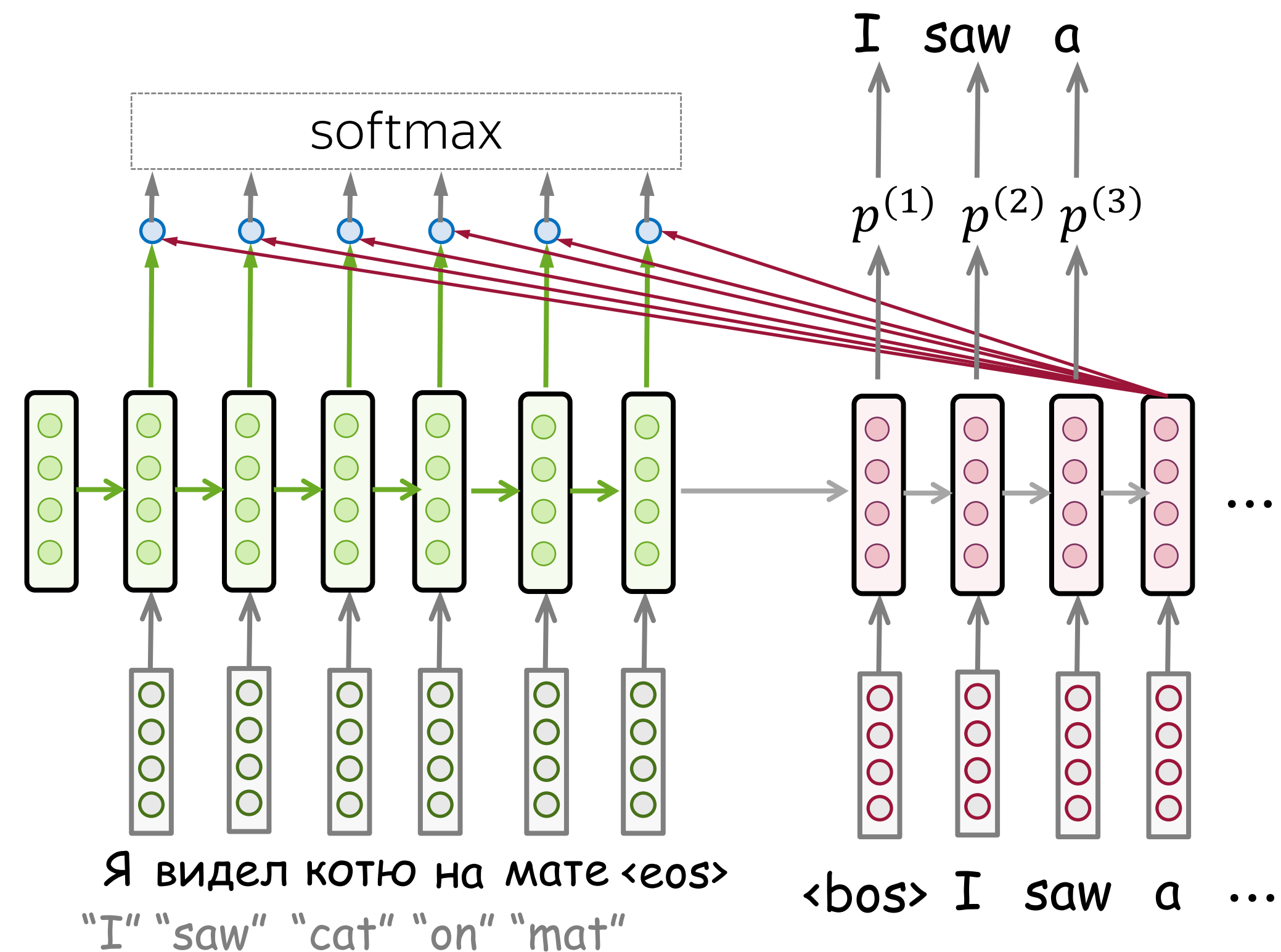
Attention: High-Level View

Attention: At different steps, let a model "focus" on different parts of the input.



Attention: High-Level View

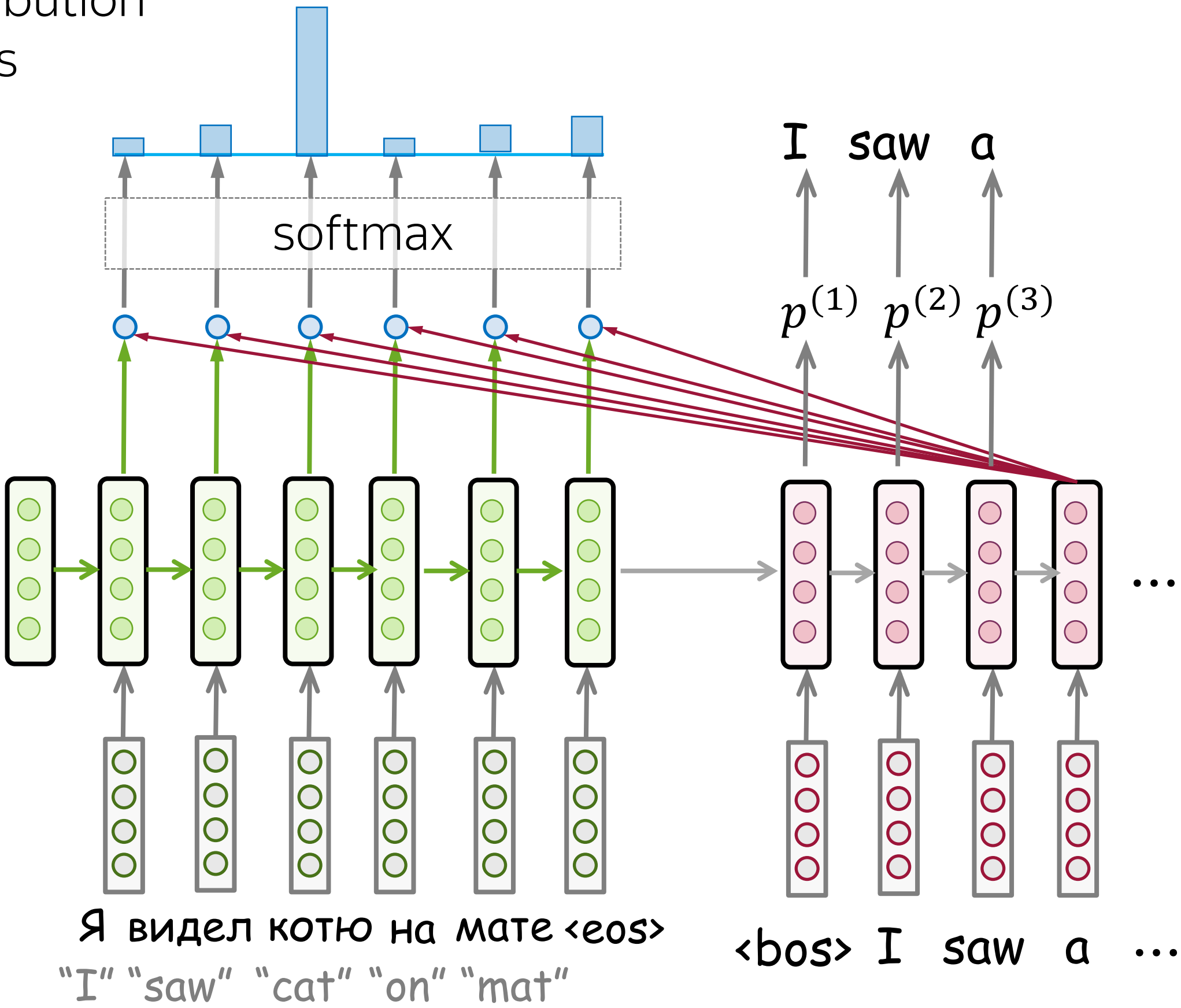
Attention: At different steps, let a model "focus" on different parts of the input.



Attention: High-Level View

Attention: At different steps, let a model "focus" on different parts of the input.

Attention weights: distribution over source tokens

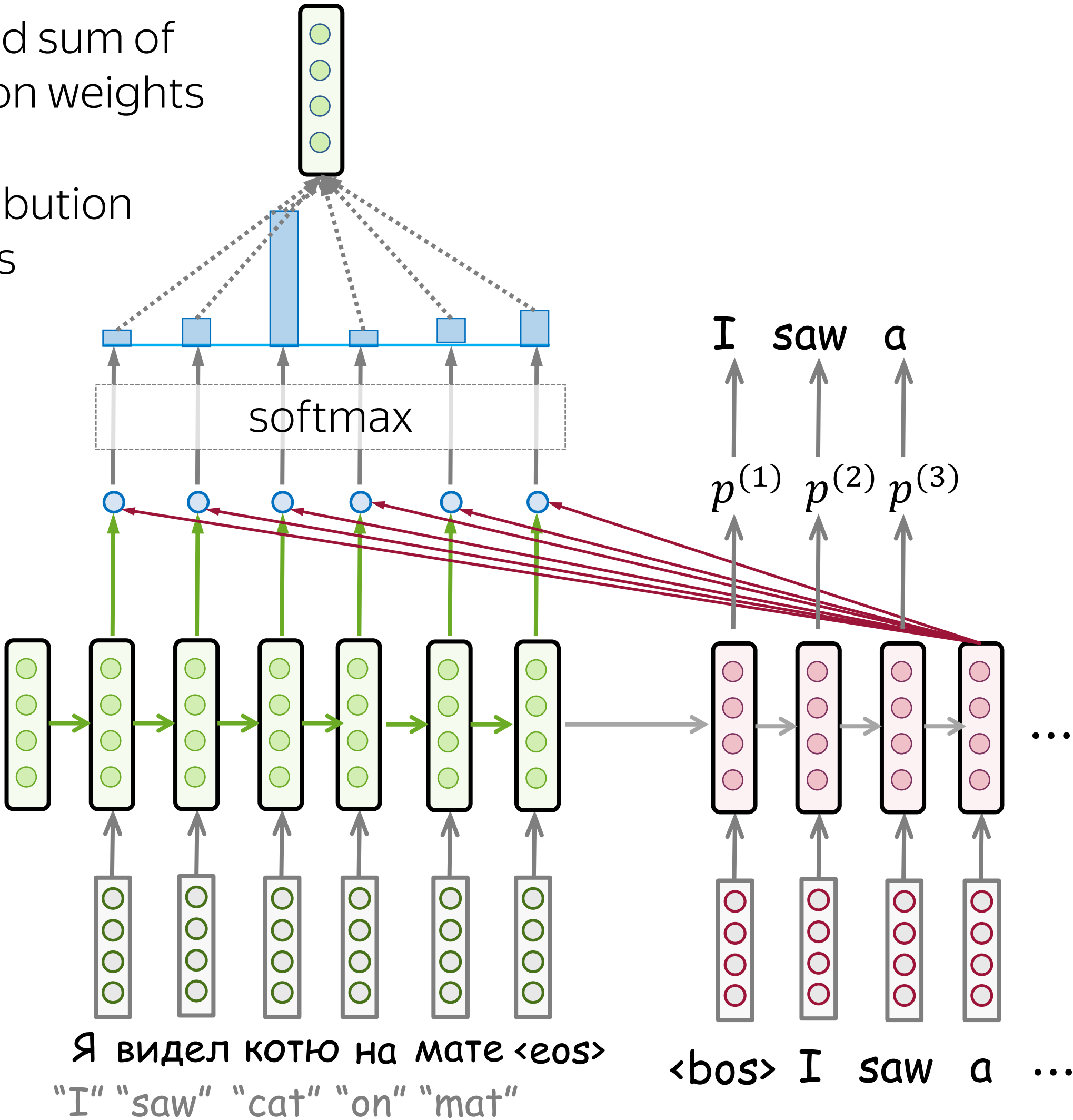


Attention: High-Level View

Attention: At different steps, let a model "focus" on different parts of the input.

Attention output: weighted sum of encoder states with attention weights

Attention weights: distribution over source tokens

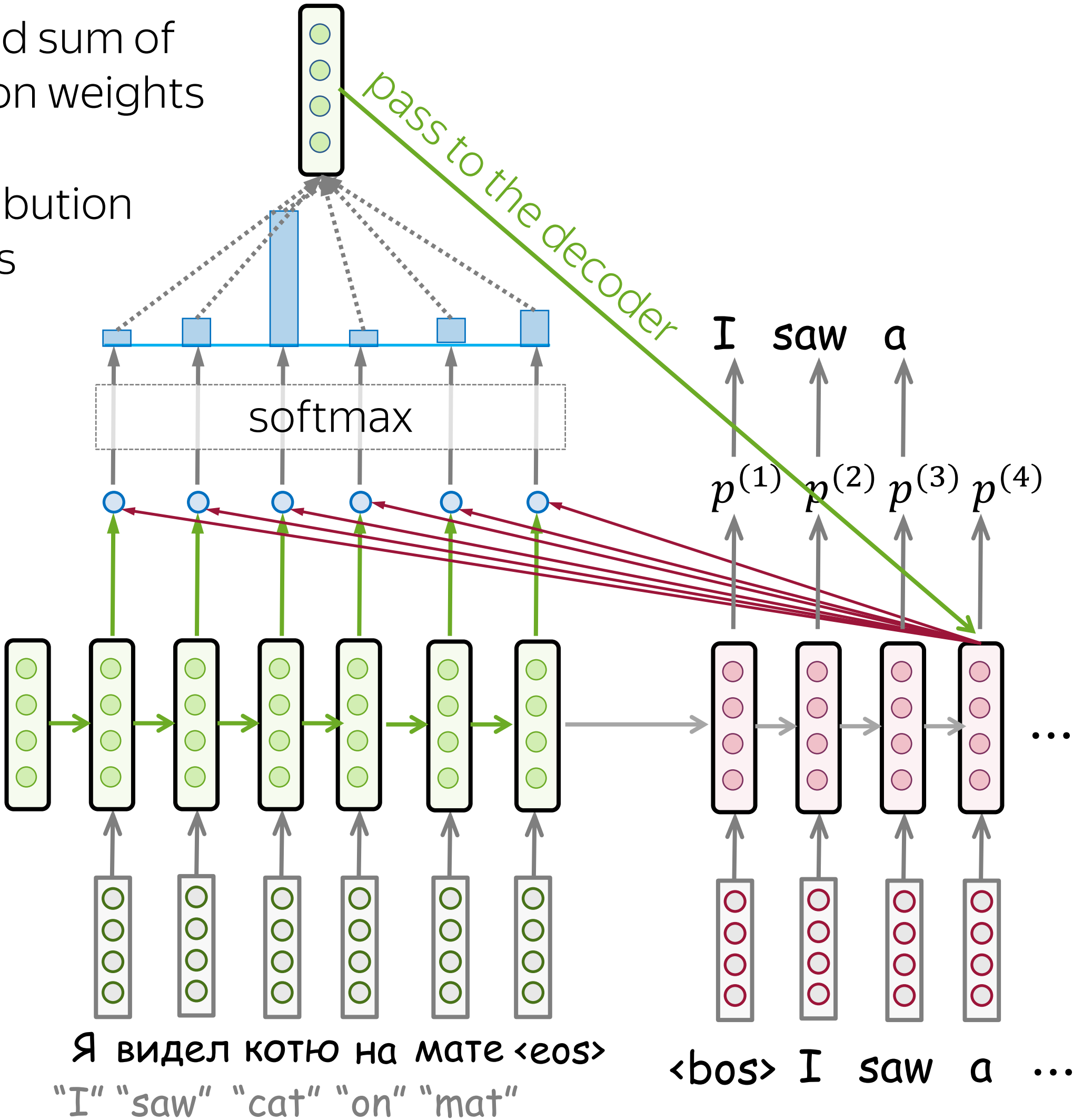


Attention: High-Level View

Attention: At different steps, let a model "focus" on different parts of the input.

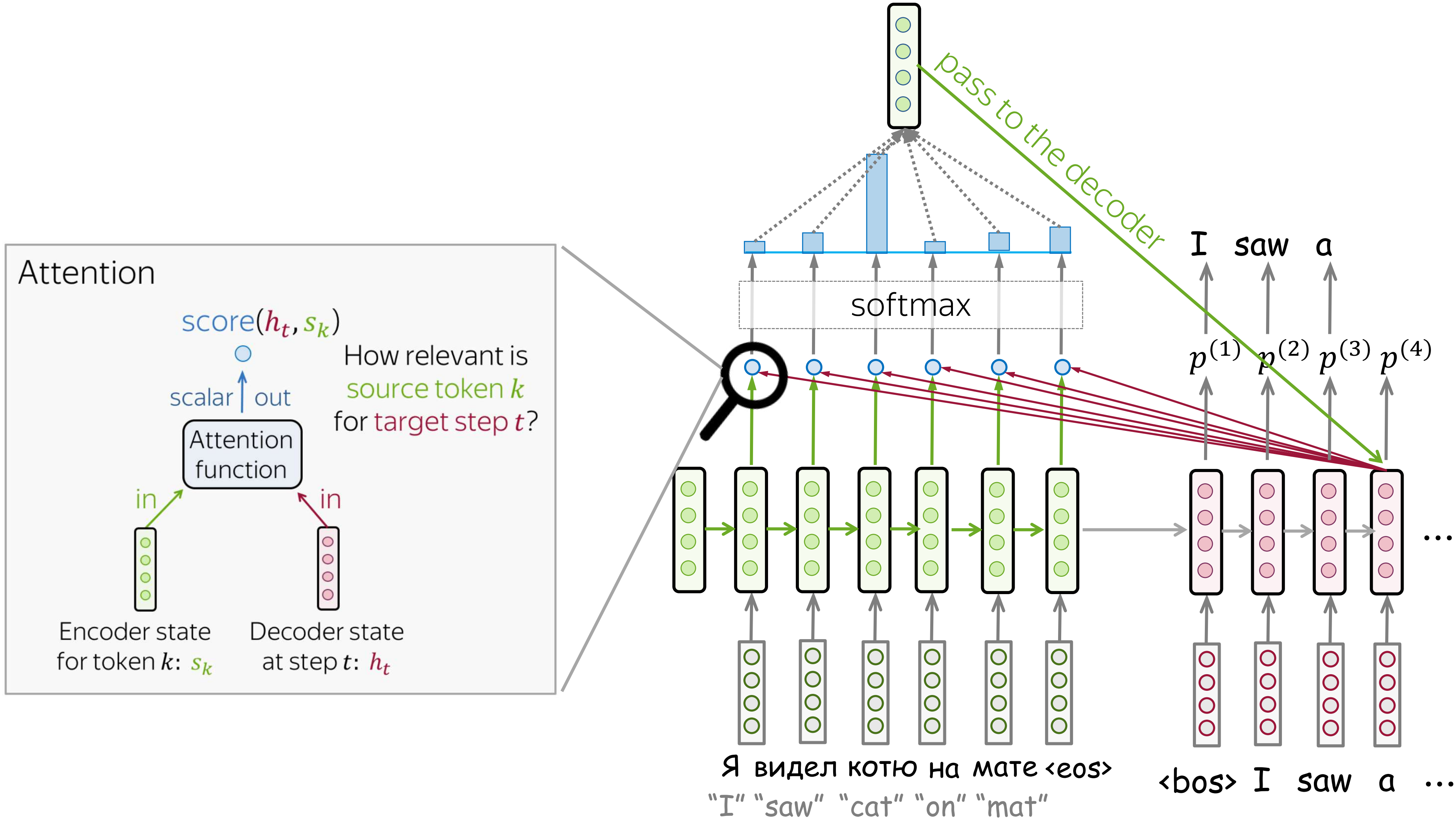
Attention output: weighted sum of encoder states with attention weights

Attention weights: distribution over source tokens



Attention: High-Level View

Attention: At different steps, let a model "focus" on different parts of the input.



Computation Pipeline

Attention input

s_1, s_2, \dots, s_m
all encoder states

h_t
one decoder state

Computation Pipeline

Attention scores

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



Attention input

s_1, s_2, \dots, s_m
all encoder states

h_t
one decoder state

Computation Pipeline

Attention scores



Attention input

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



“How relevant is source token k for target step t ?”

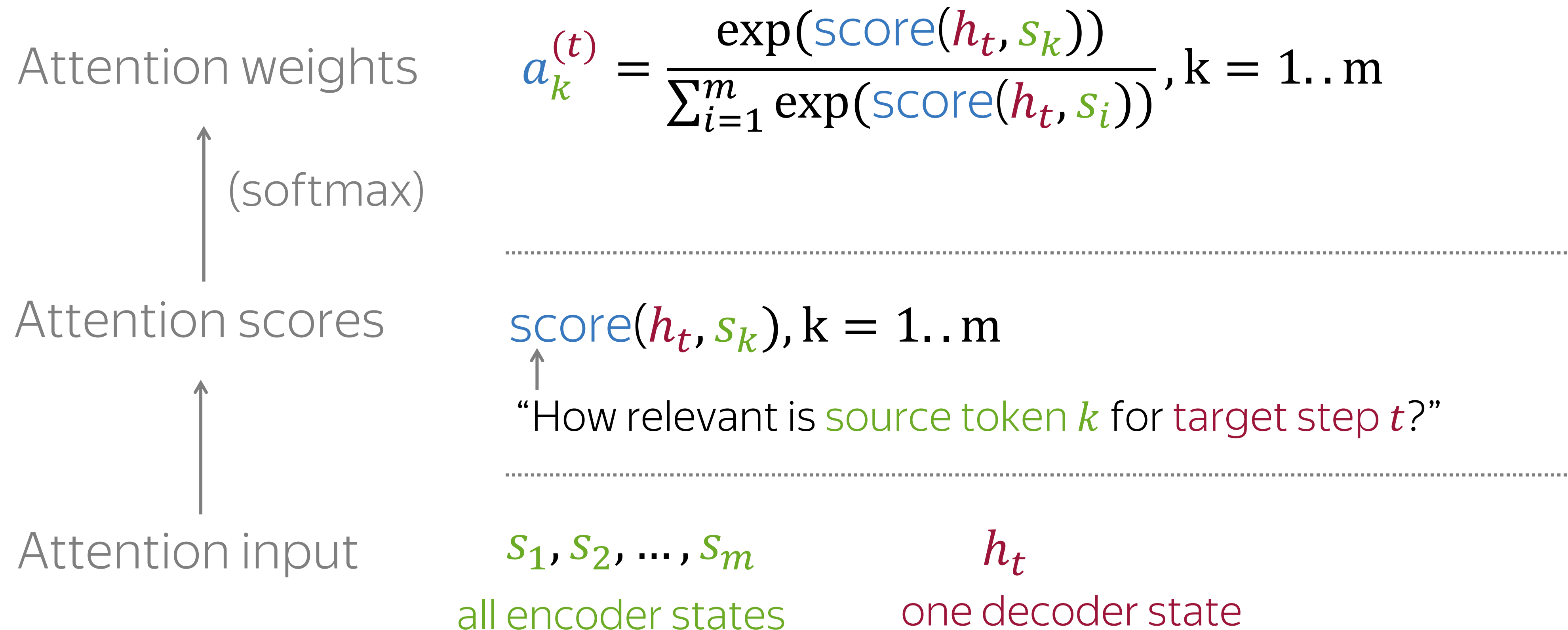
s_1, s_2, \dots, s_m

all encoder states

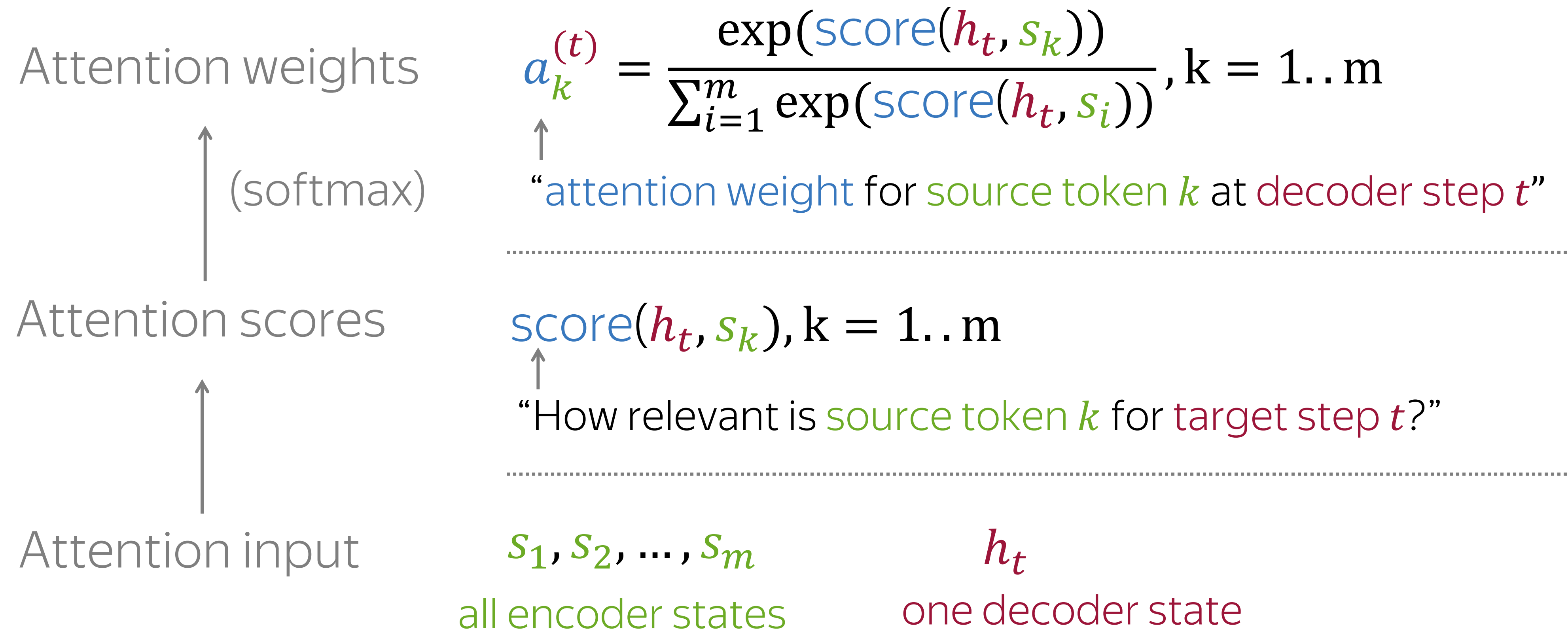
h_t

one decoder state

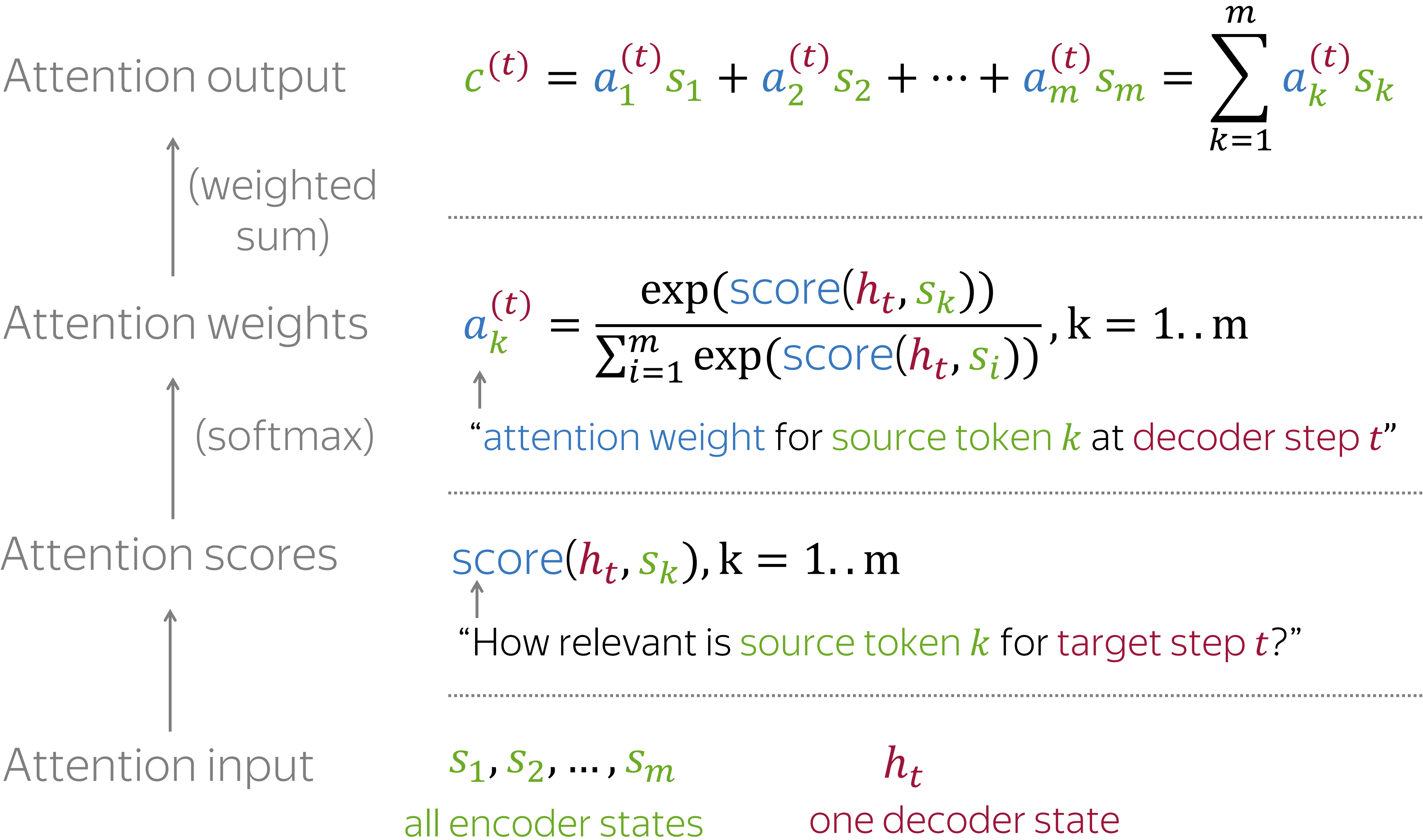
Computation Pipeline



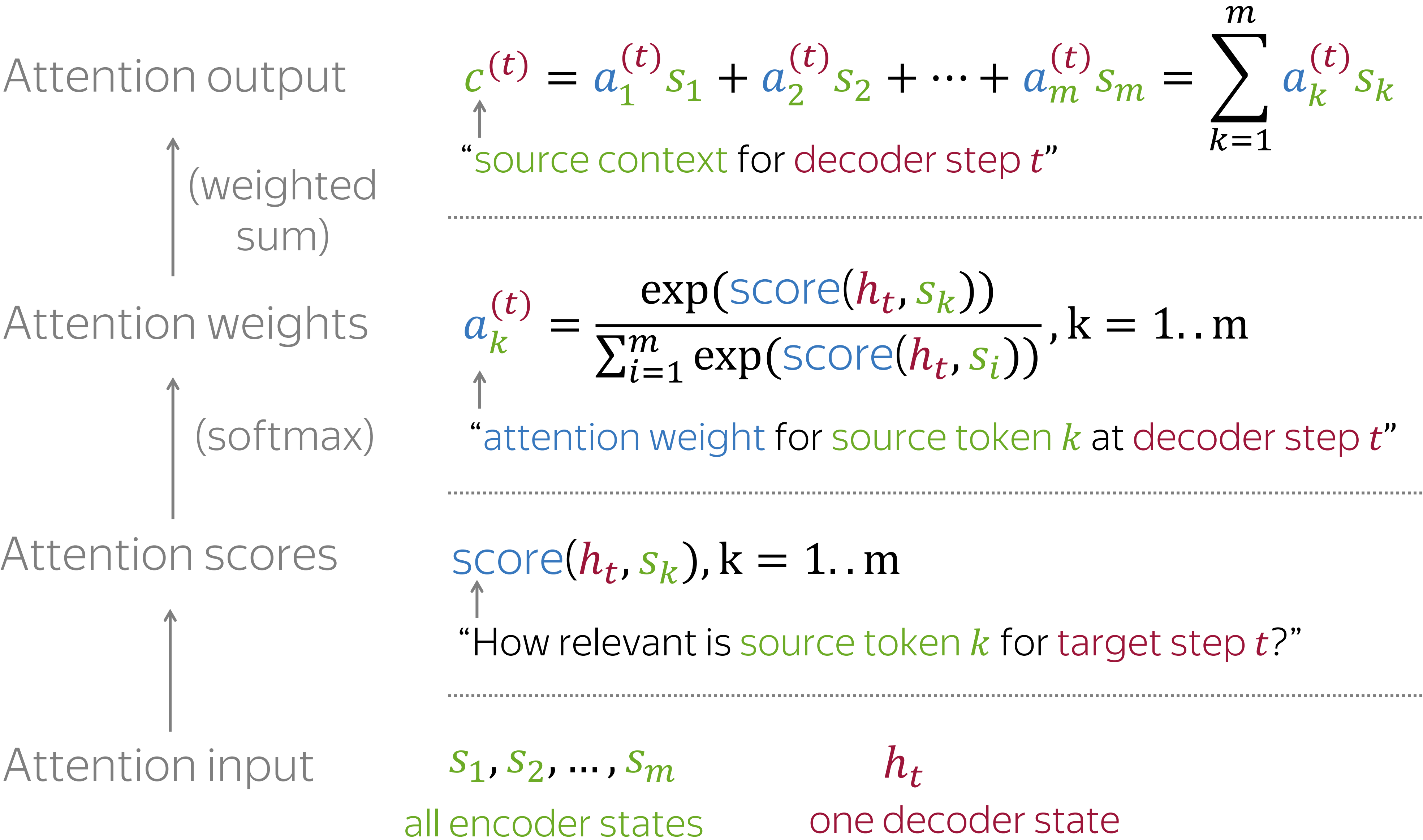
Computation Pipeline



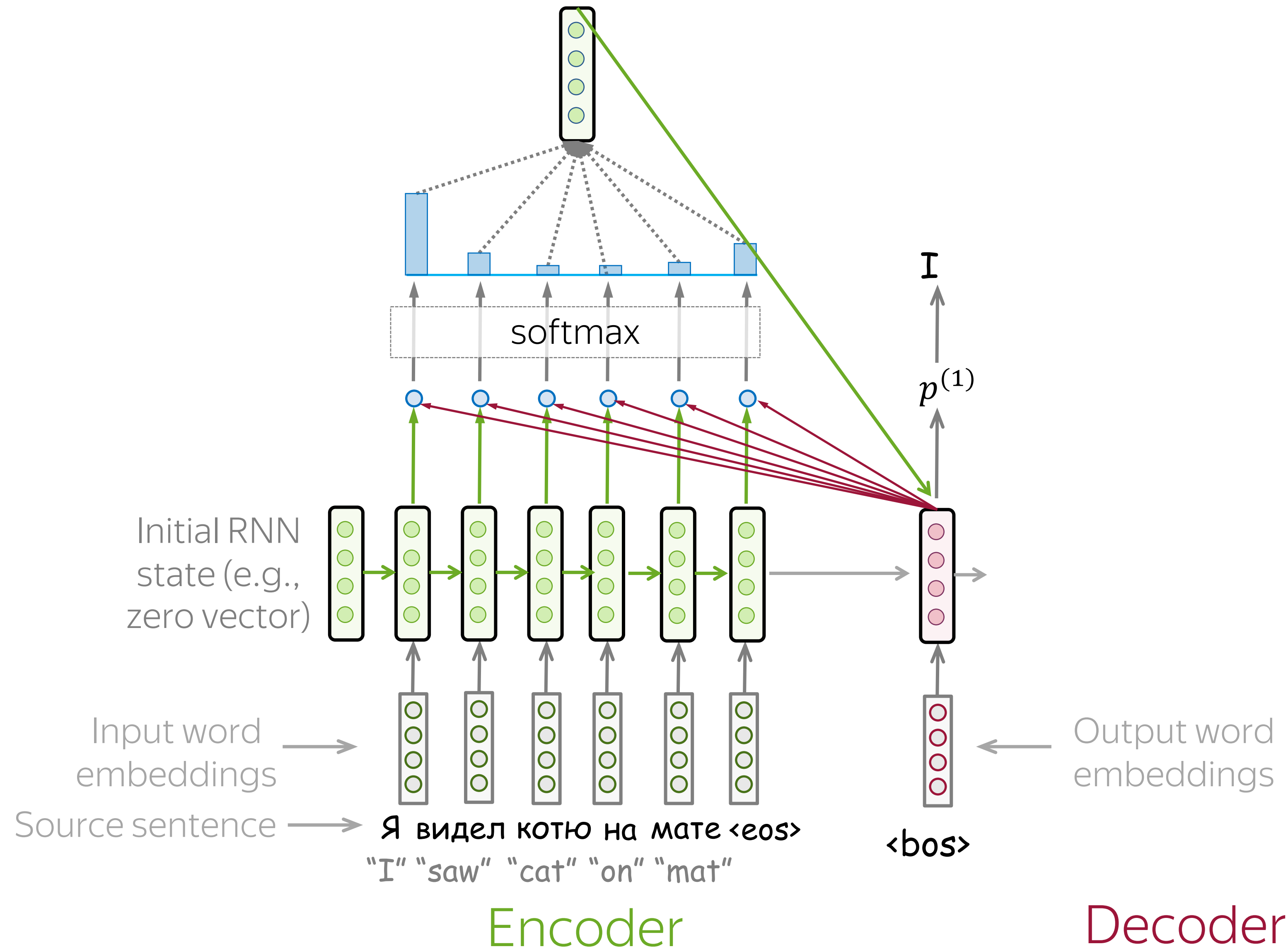
Computation Pipeline



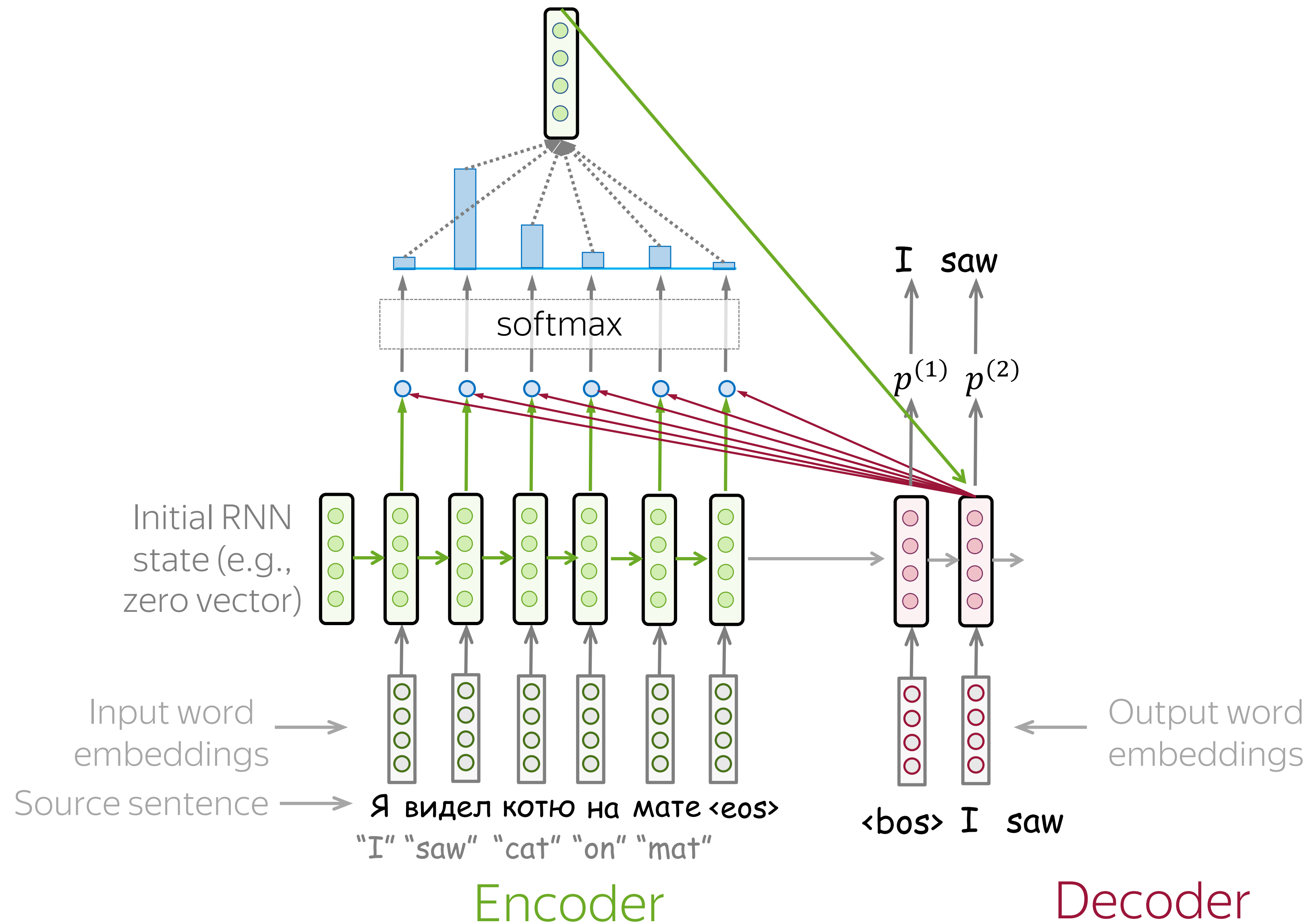
Computation Pipeline



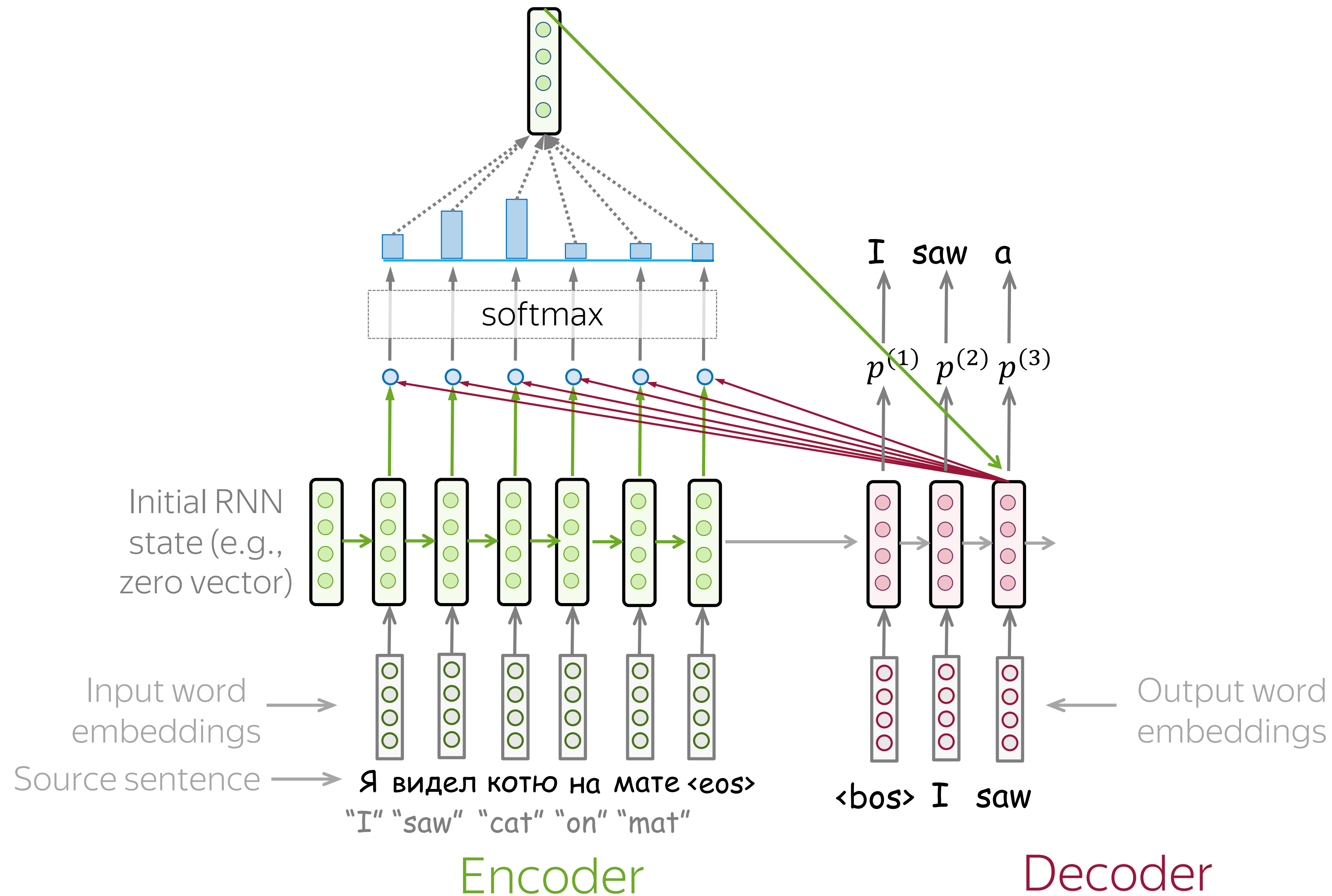
Model Learns to Pick Relevant Tokens



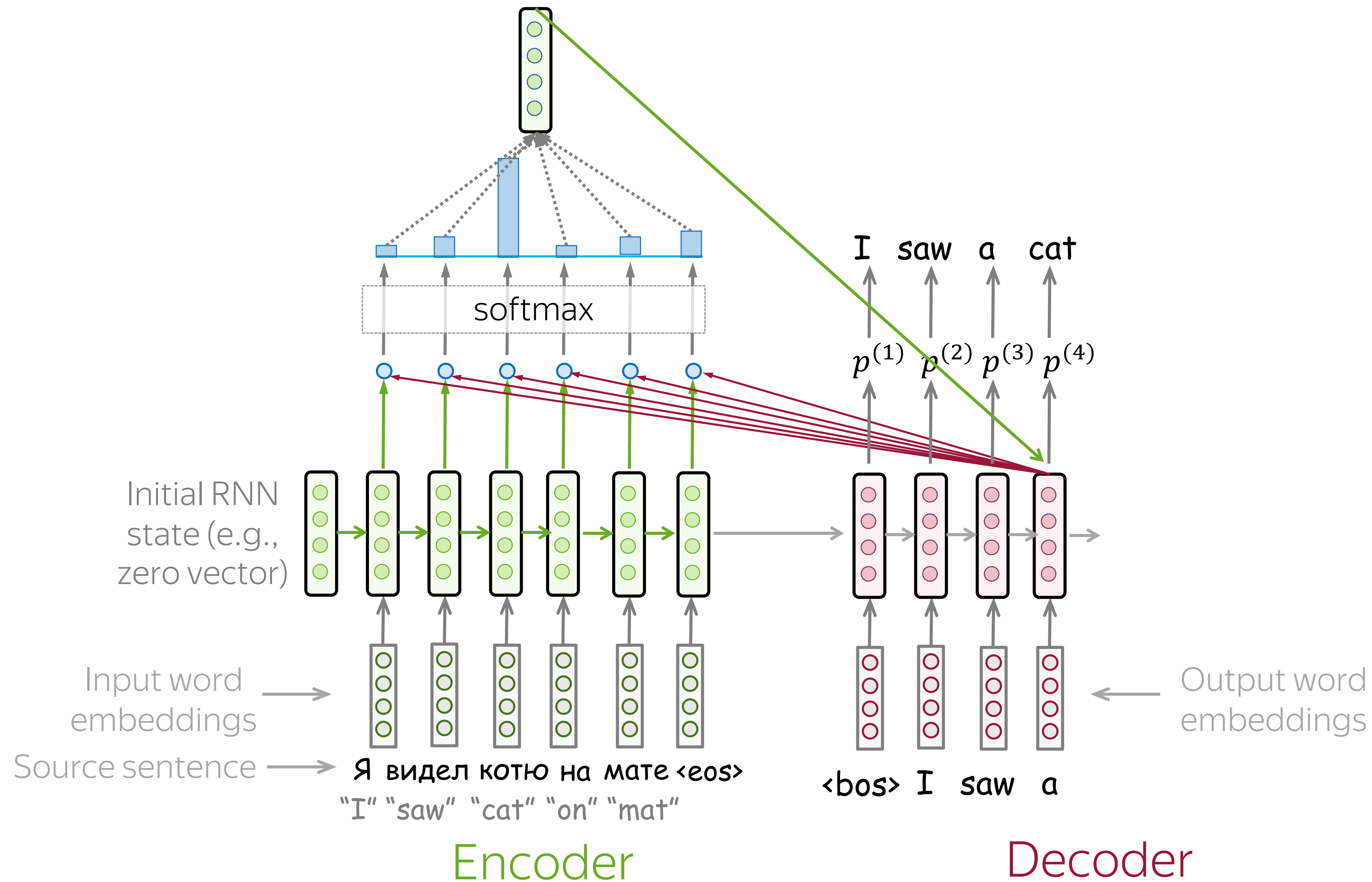
Model Learns to Pick Relevant Tokens



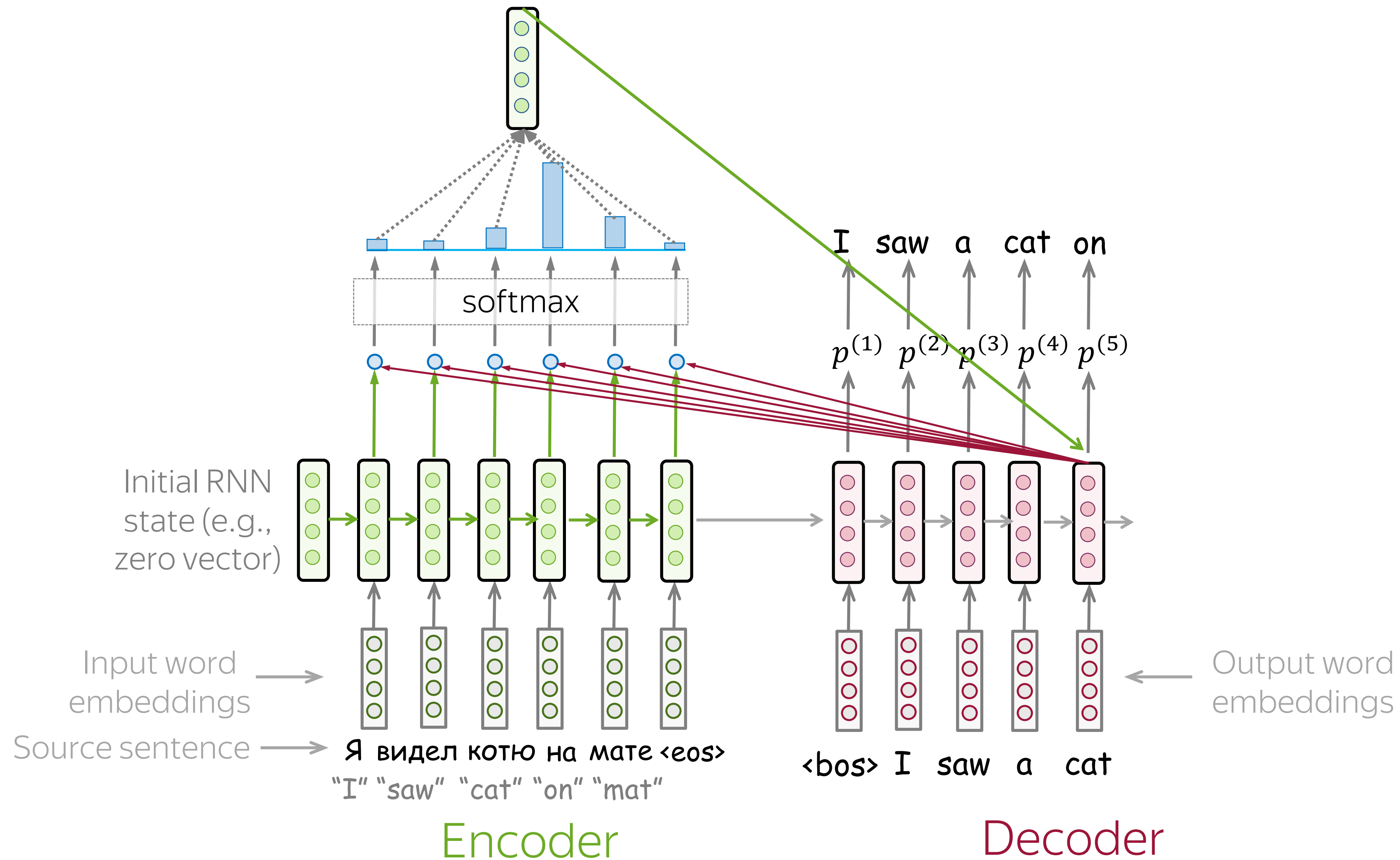
Model Learns to Pick Relevant Tokens



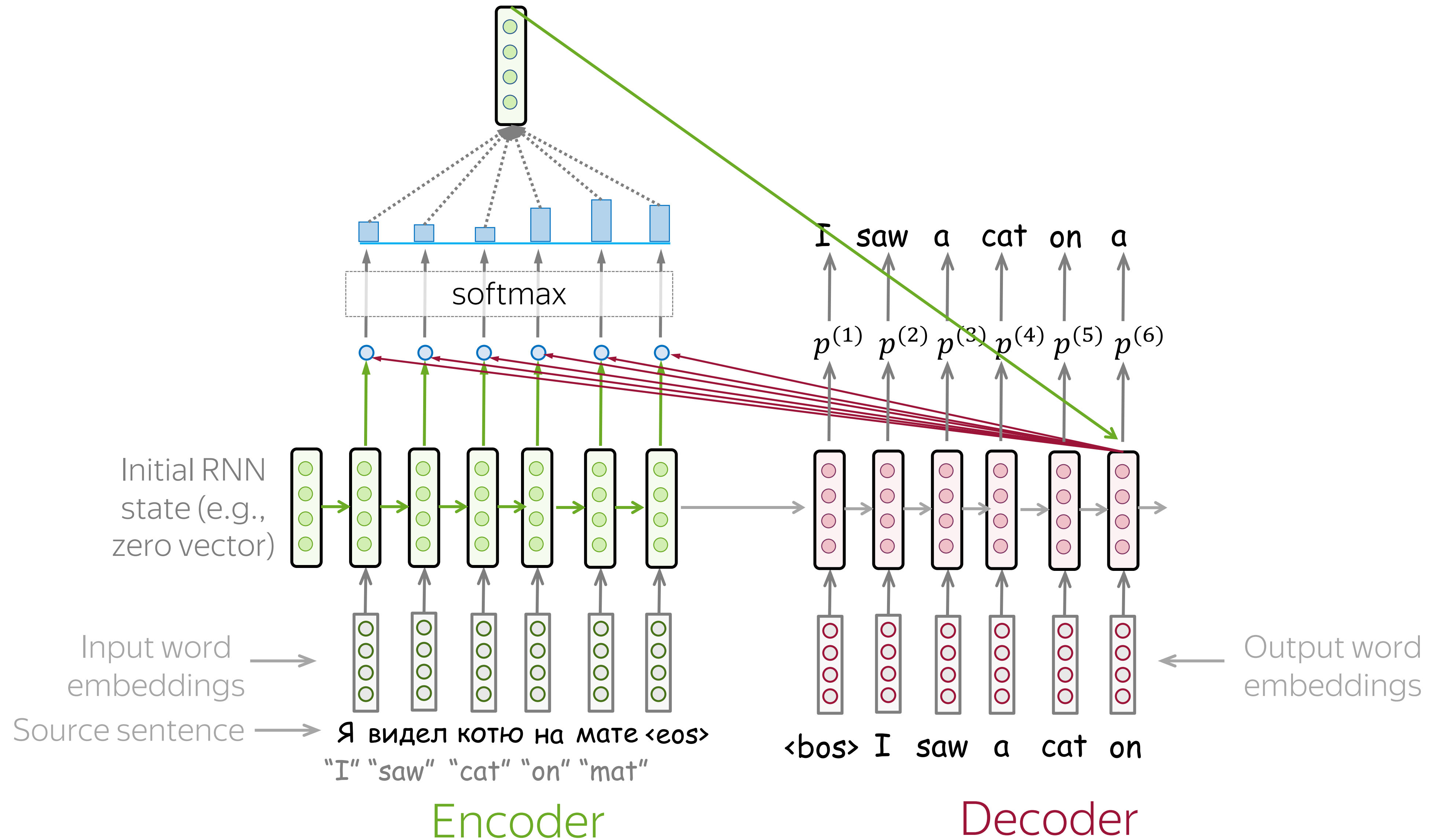
Model Learns to Pick Relevant Tokens



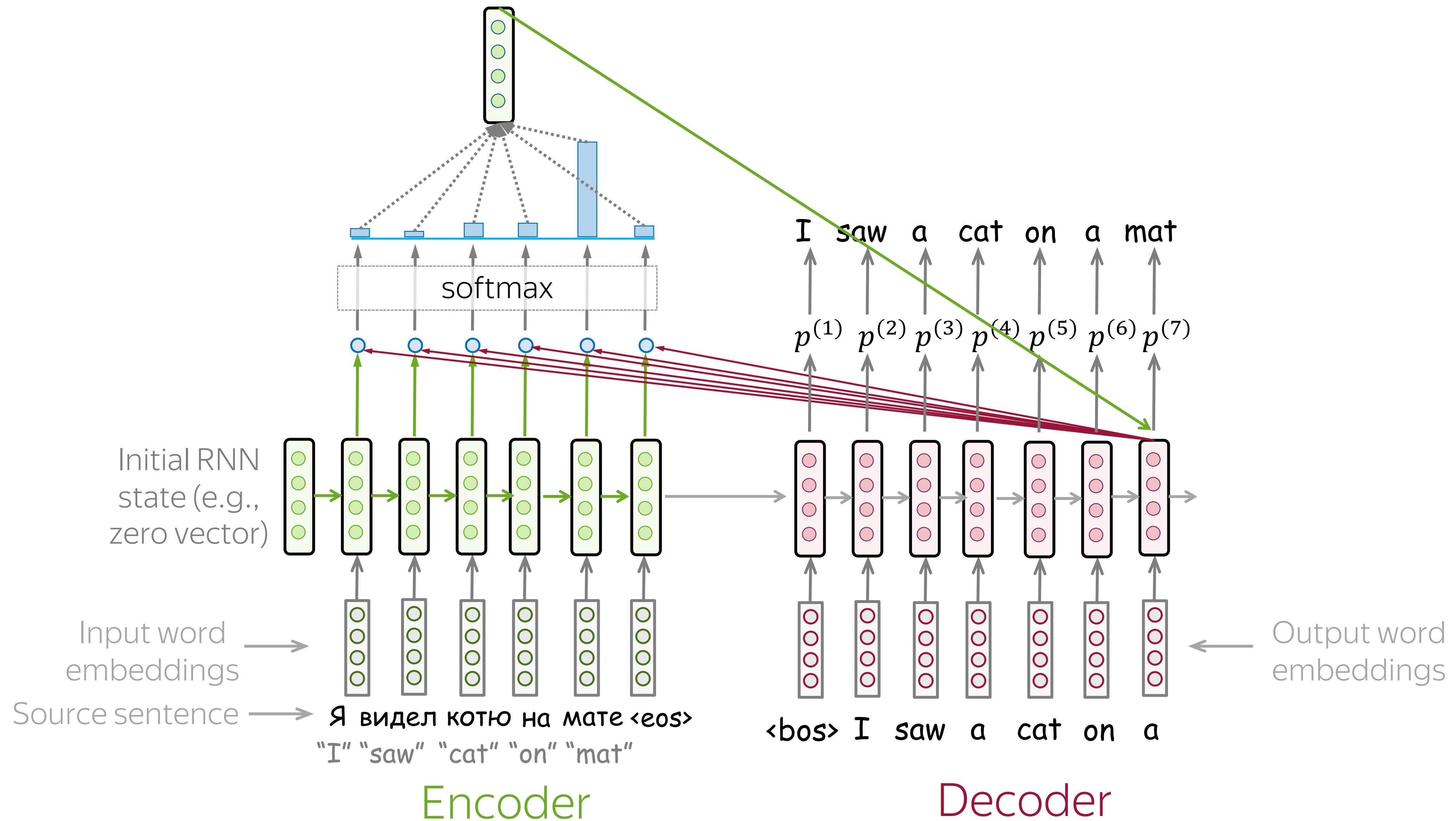
Model Learns to Pick Relevant Tokens



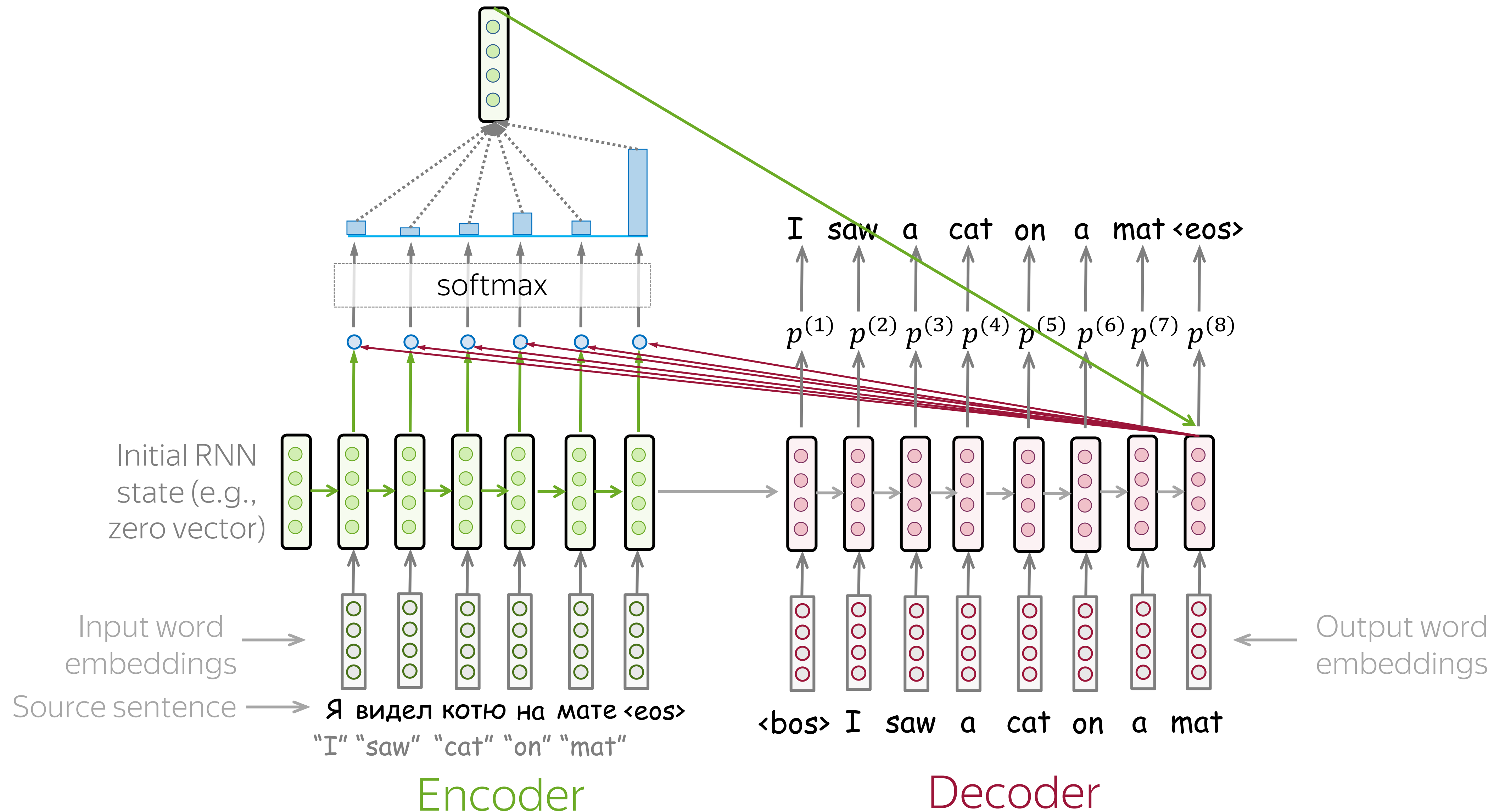
Model Learns to Pick Relevant Tokens



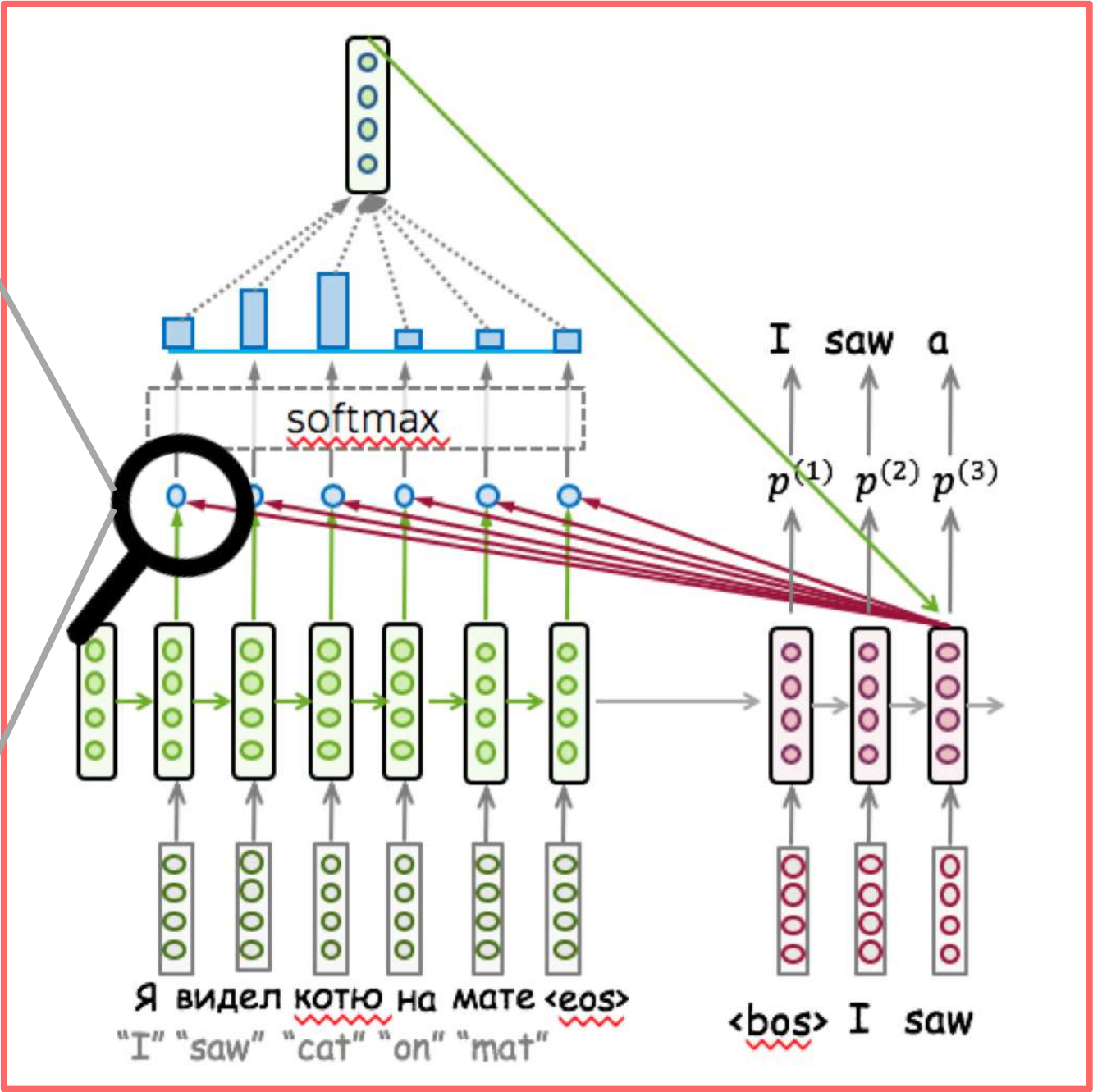
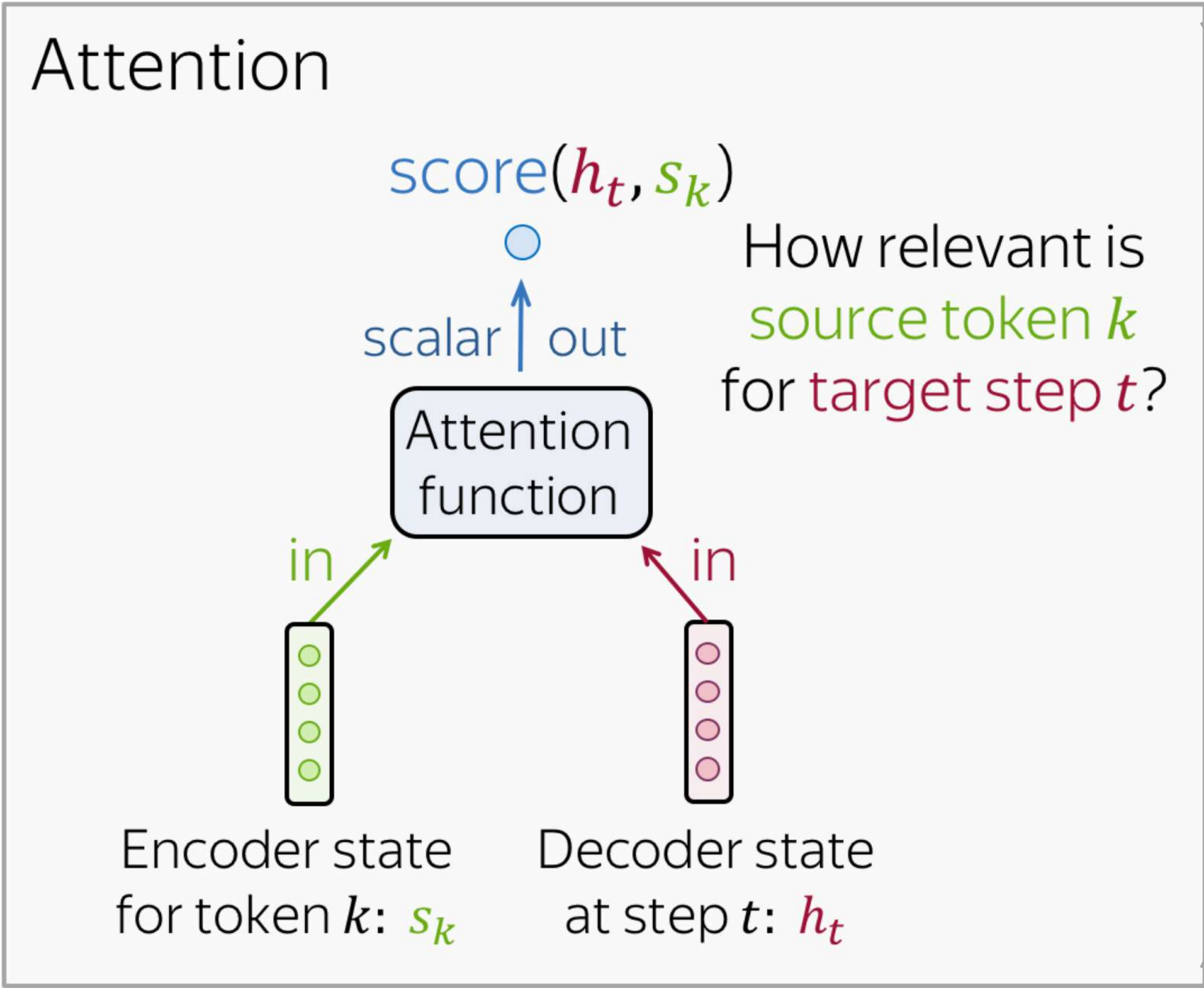
Model Learns to Pick Relevant Tokens



Model Learns to Pick Relevant Tokens



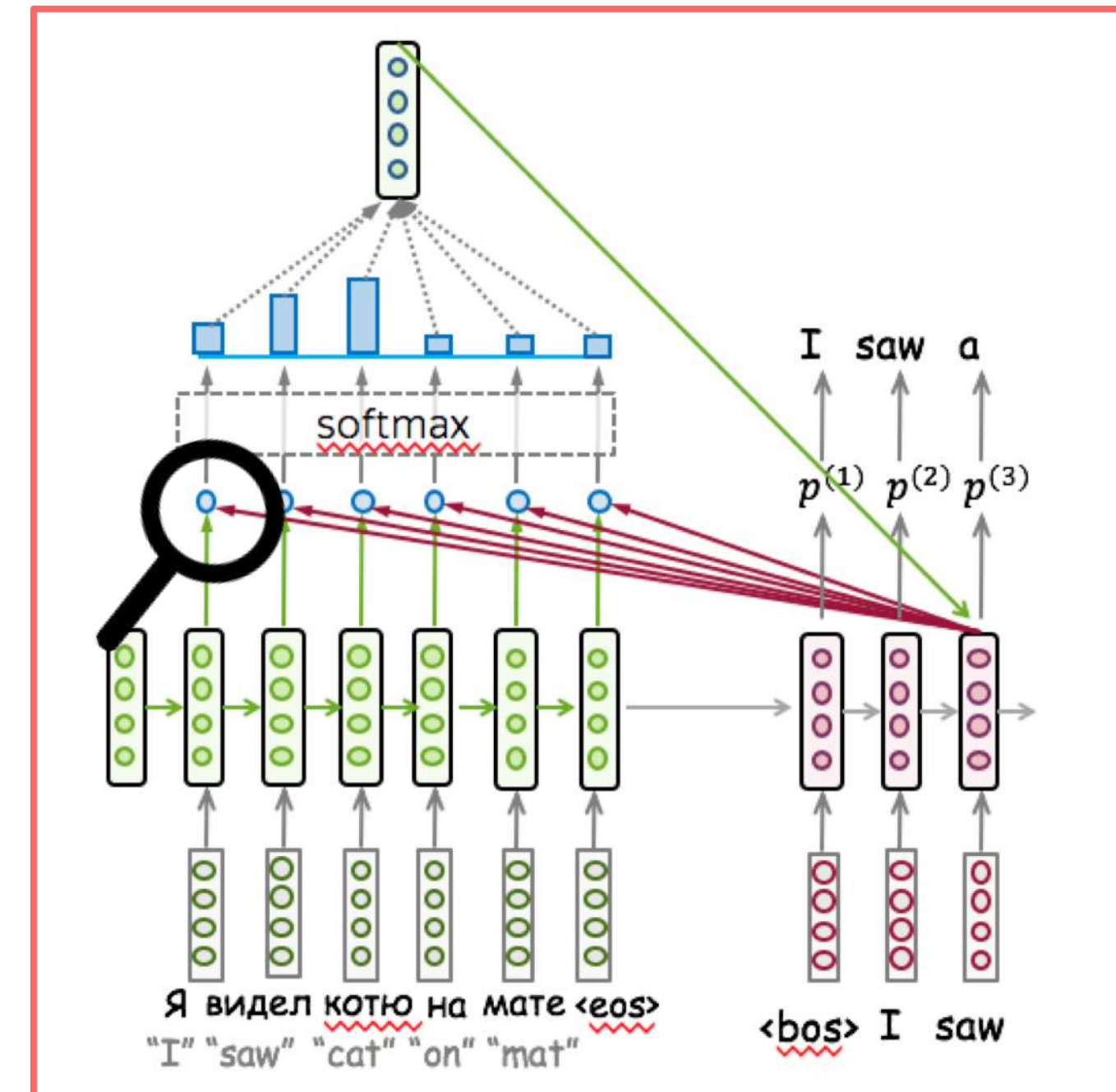
Attention Score Functions



Attention Score Functions

- Dot-product: $\text{score}(h_t, s_k) = h_t^T s_k$

$$\begin{matrix} h_t^T \\ \text{---} \end{matrix} \times \begin{matrix} \text{---} \\ s_k \end{matrix}$$



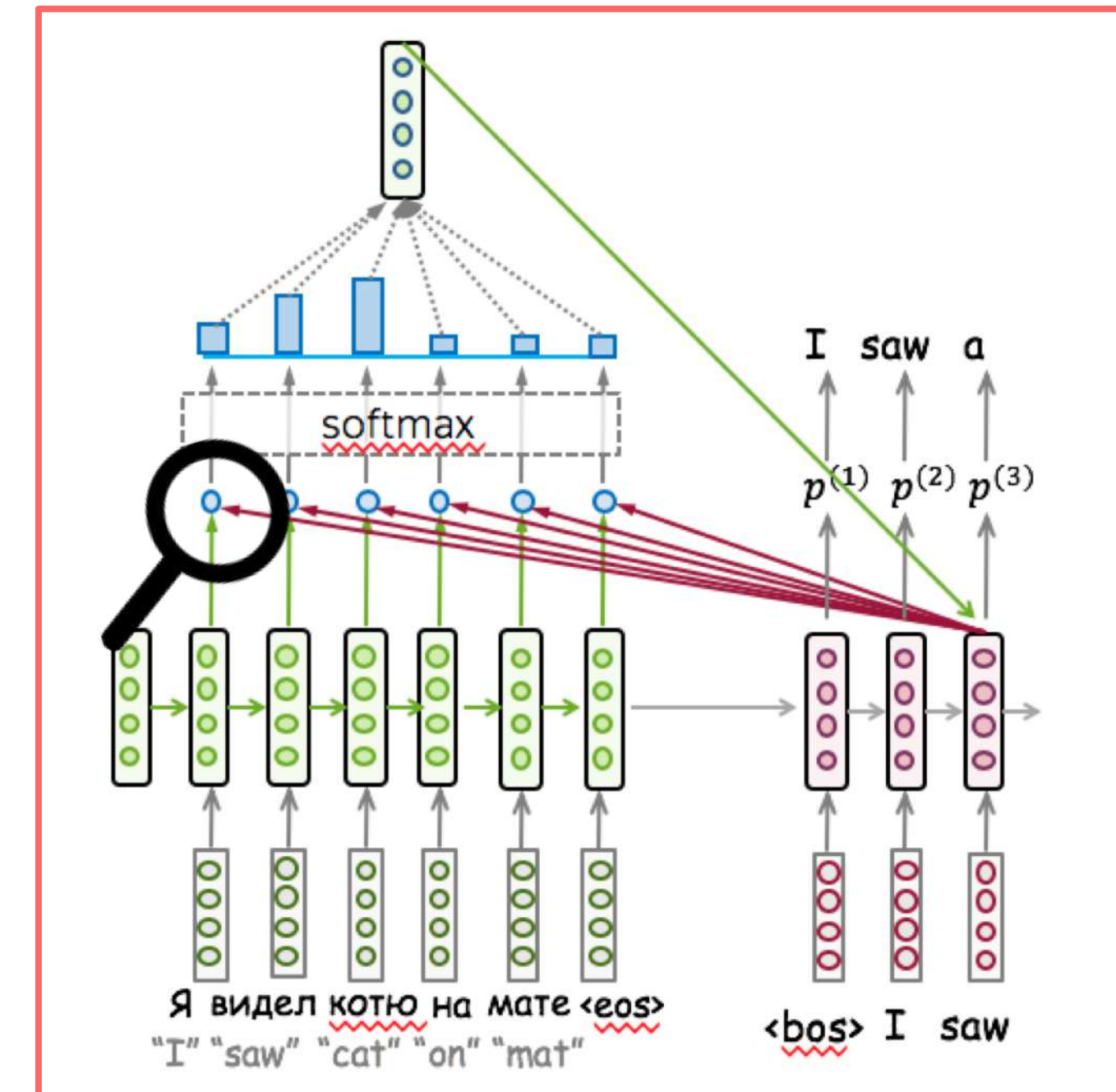
Attention Score Functions

- Dot-product: $\text{score}(h_t, s_k) = h_t^T s_k$

$$\begin{matrix} h_t^T \\ \text{---} \end{matrix} \times \begin{matrix} \text{---} \\ s_k \end{matrix}$$

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

$$\begin{matrix} h_t^T \\ \text{---} \end{matrix} \times \begin{bmatrix} W \end{bmatrix} \times \begin{matrix} \text{---} \\ s_k \end{matrix}$$



Attention Score Functions

- Dot-product: $\text{score}(h_t, s_k) = h_t^T s_k$

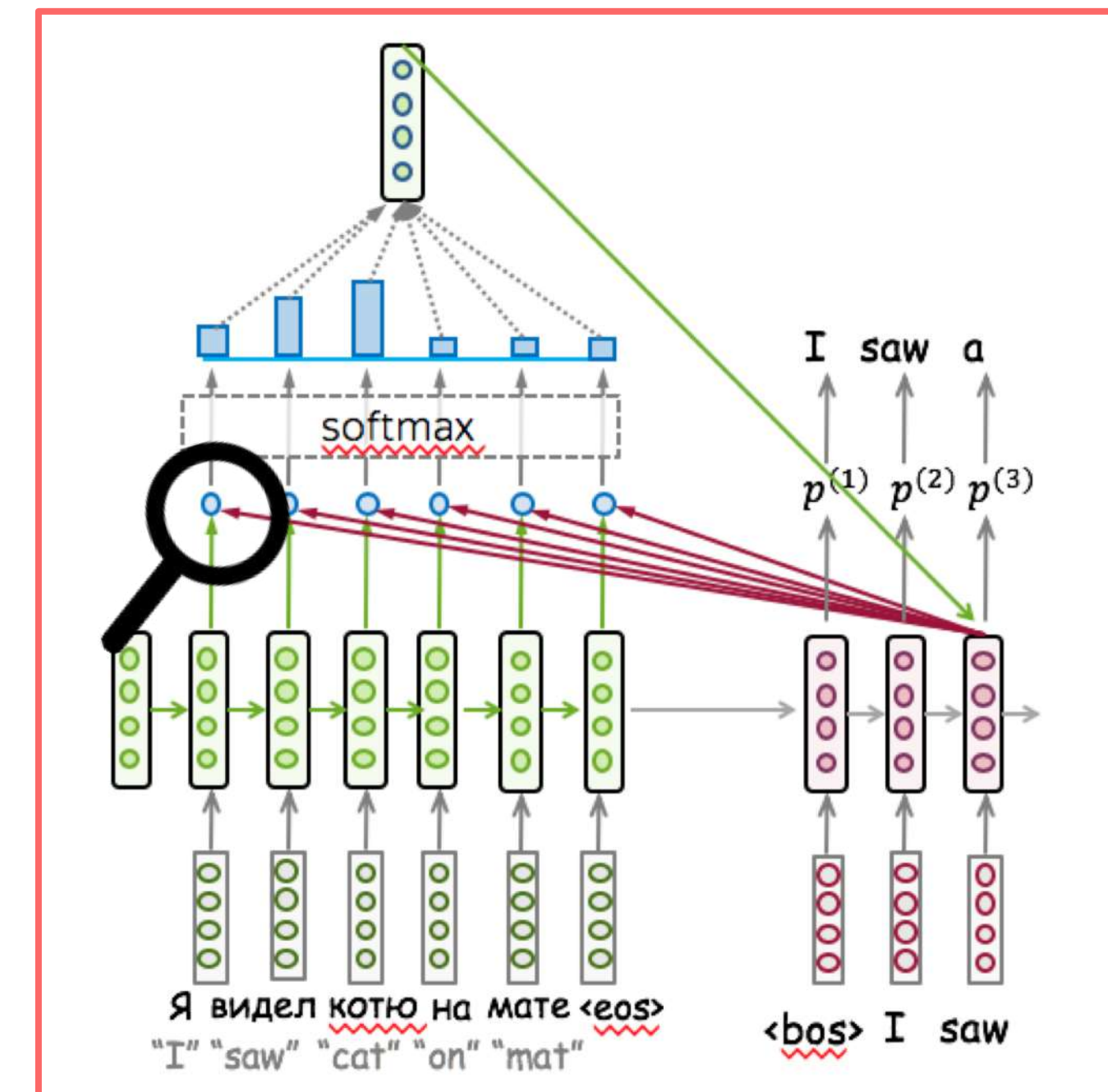
$$\begin{matrix} h_t^T \\ \text{---} \end{matrix} \times \begin{matrix} \text{---} \\ s_k \end{matrix}$$

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

$$\begin{matrix} h_t^T \\ \text{---} \end{matrix} \times \begin{bmatrix} W \end{bmatrix} \times \begin{matrix} \text{---} \\ s_k \end{matrix}$$

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

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



What is going to happen:

- Seq2seq Basics

- Attention 

- Why do we need it?


- Attention: High-Level

- Attention Score Functions

- Models: Bahdanau vs Luong

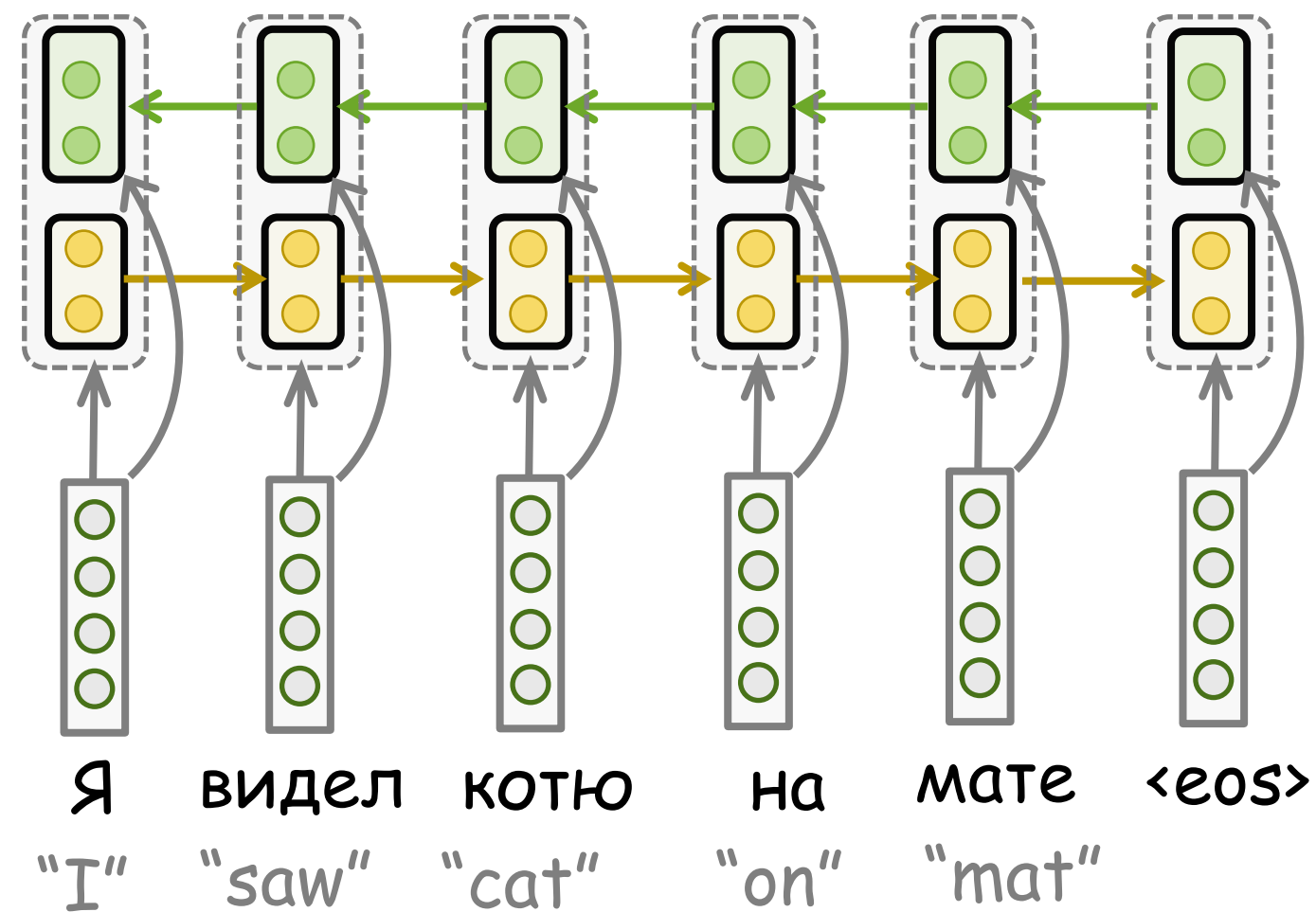
-  Analysis and Interpretability

What is going to happen:

- Seq2seq Basics
- Attention →
 - Why do we need it?
 - Attention: High-Level
 - Attention Score Functions
 - Models: Bahdanau vs Luong
- Transformer
- Subword Segmentation: BPE
-  Analysis and Interpretability

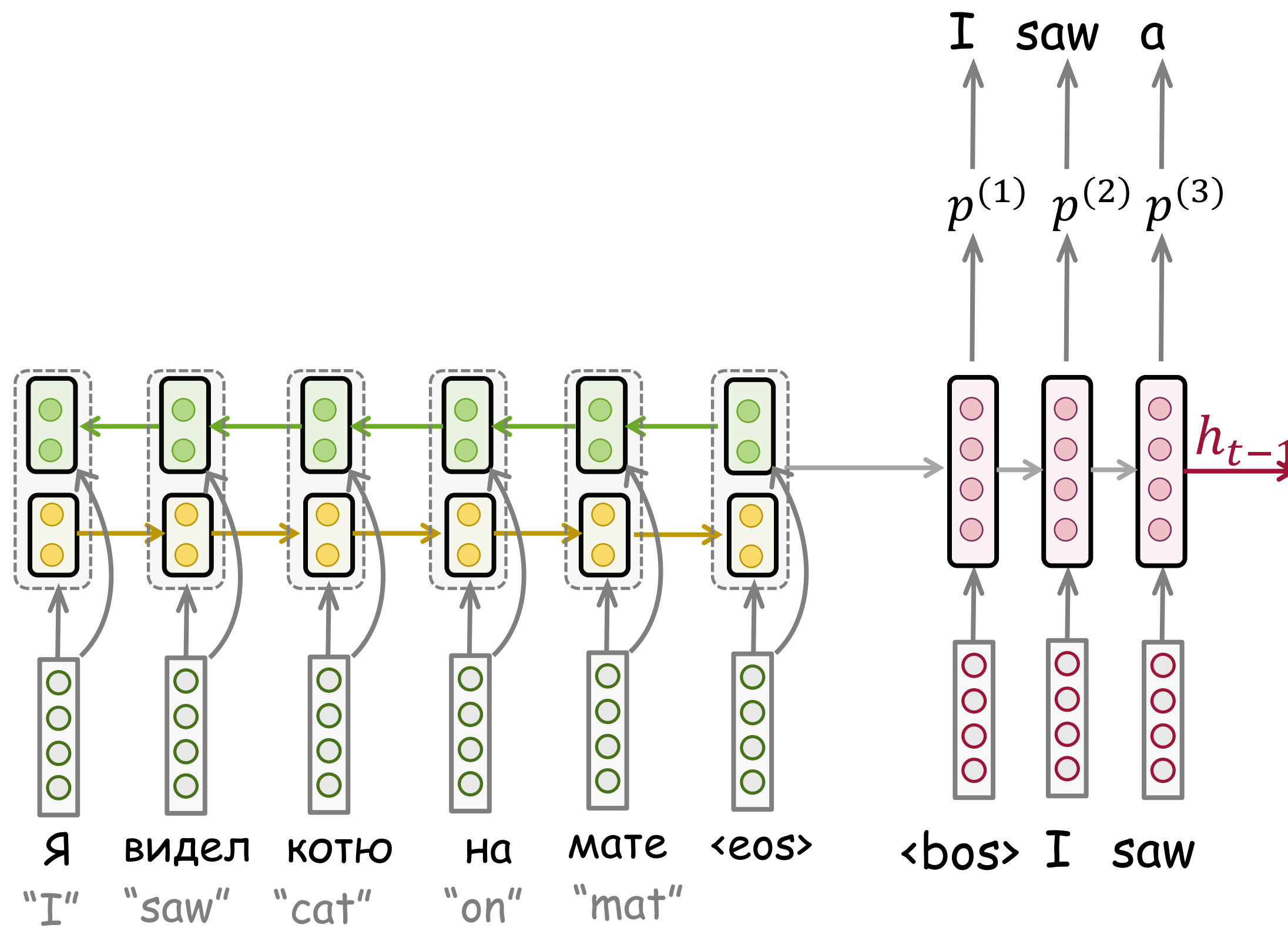
Bahdanau Model (the original attention model)

Bidirectional encoder
Concatenate states from
forward and backward RNNs



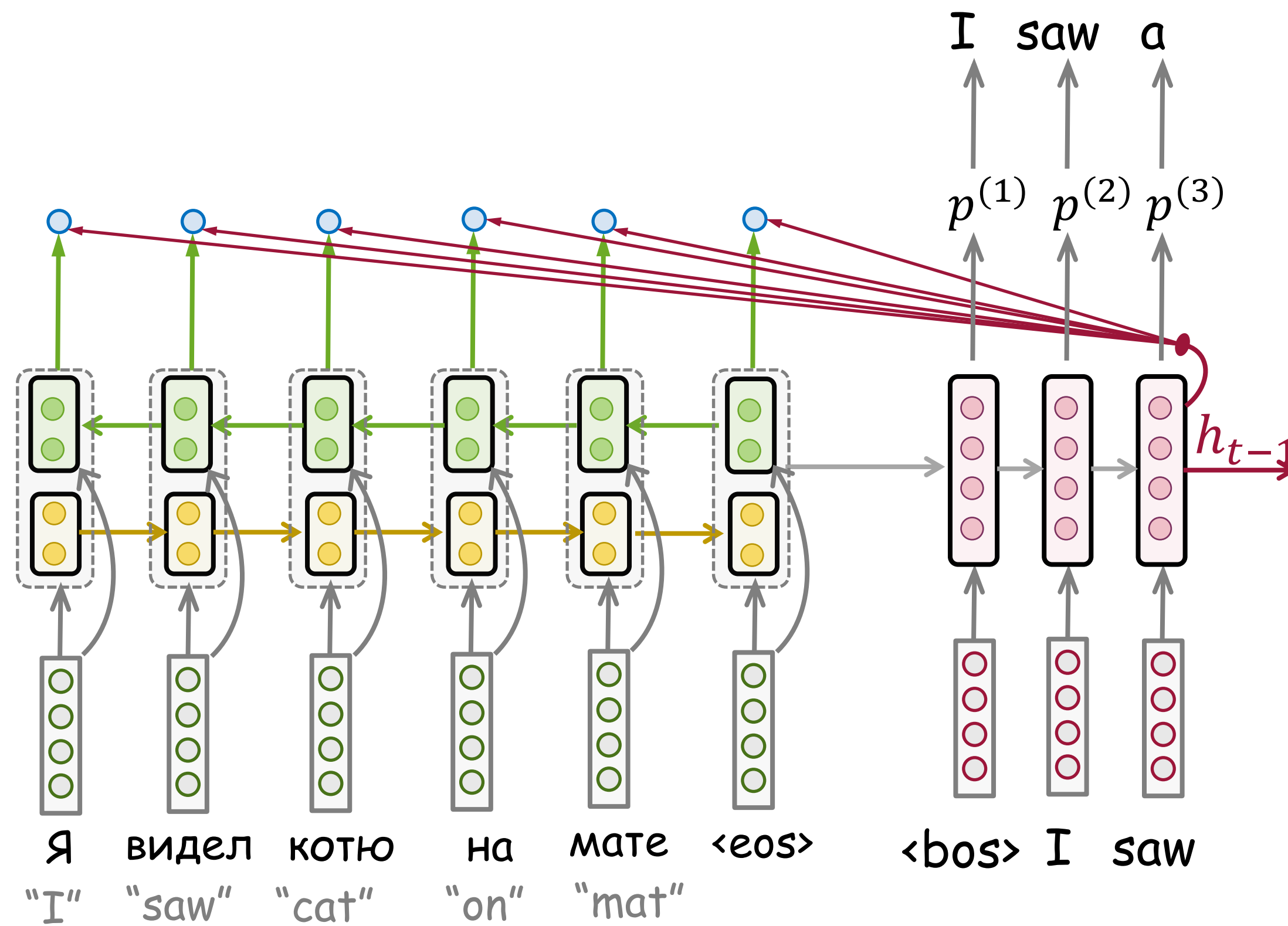
Bahdanau Model (the original attention model)

Bidirectional encoder
Concatenate states from
forward and backward RNNs

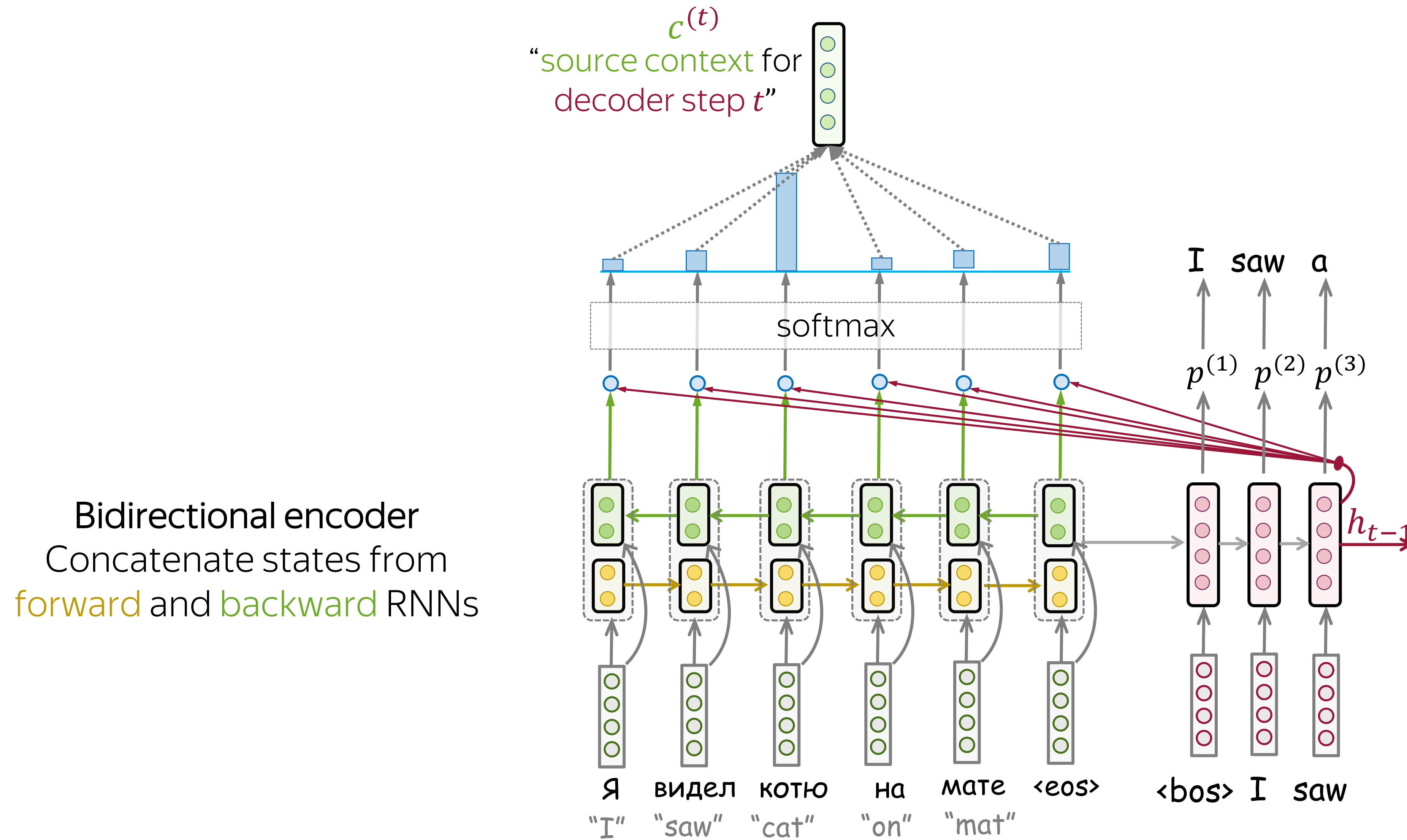


Bahdanau Model (the original attention model)

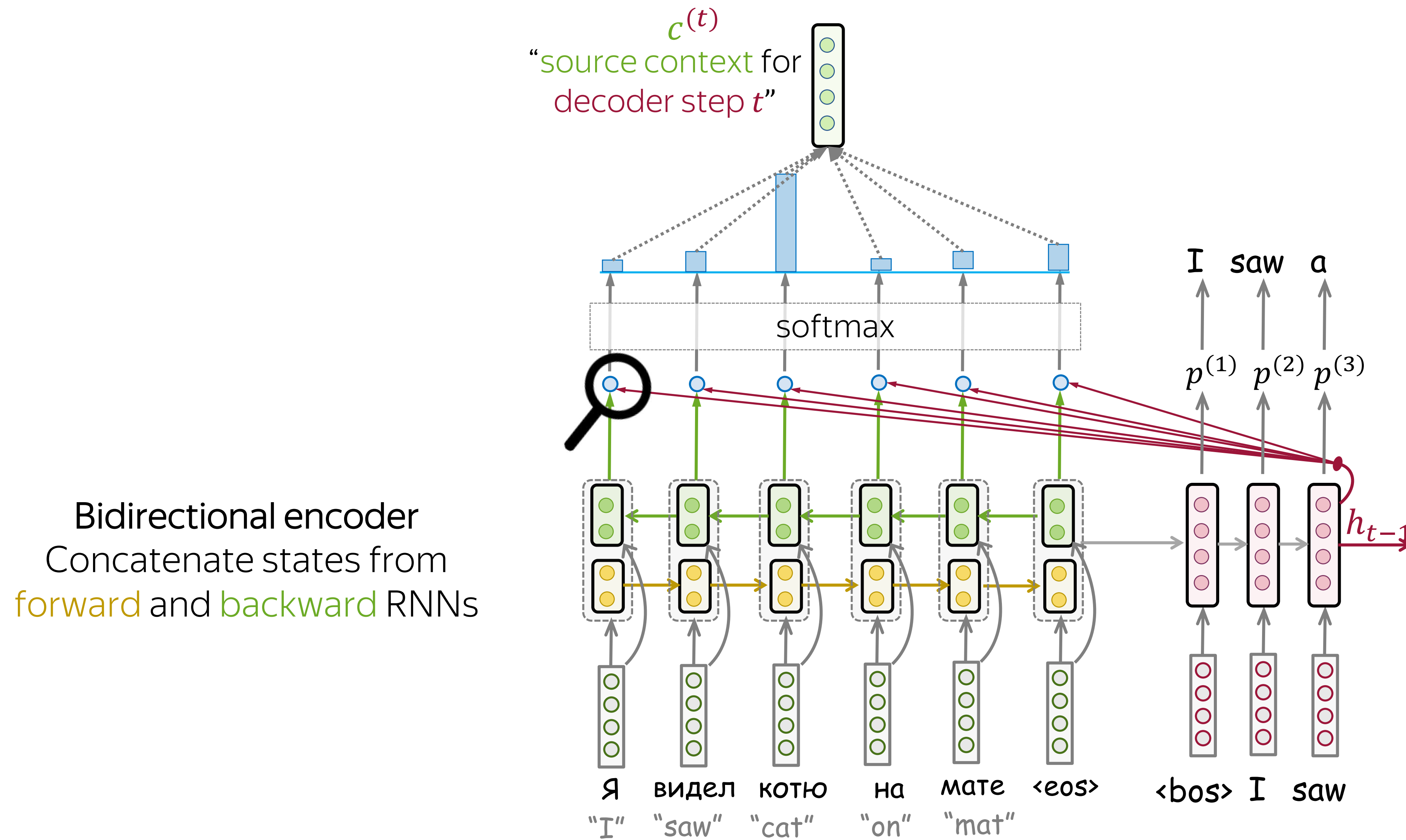
Bidirectional encoder
Concatenate states from
forward and backward RNNs



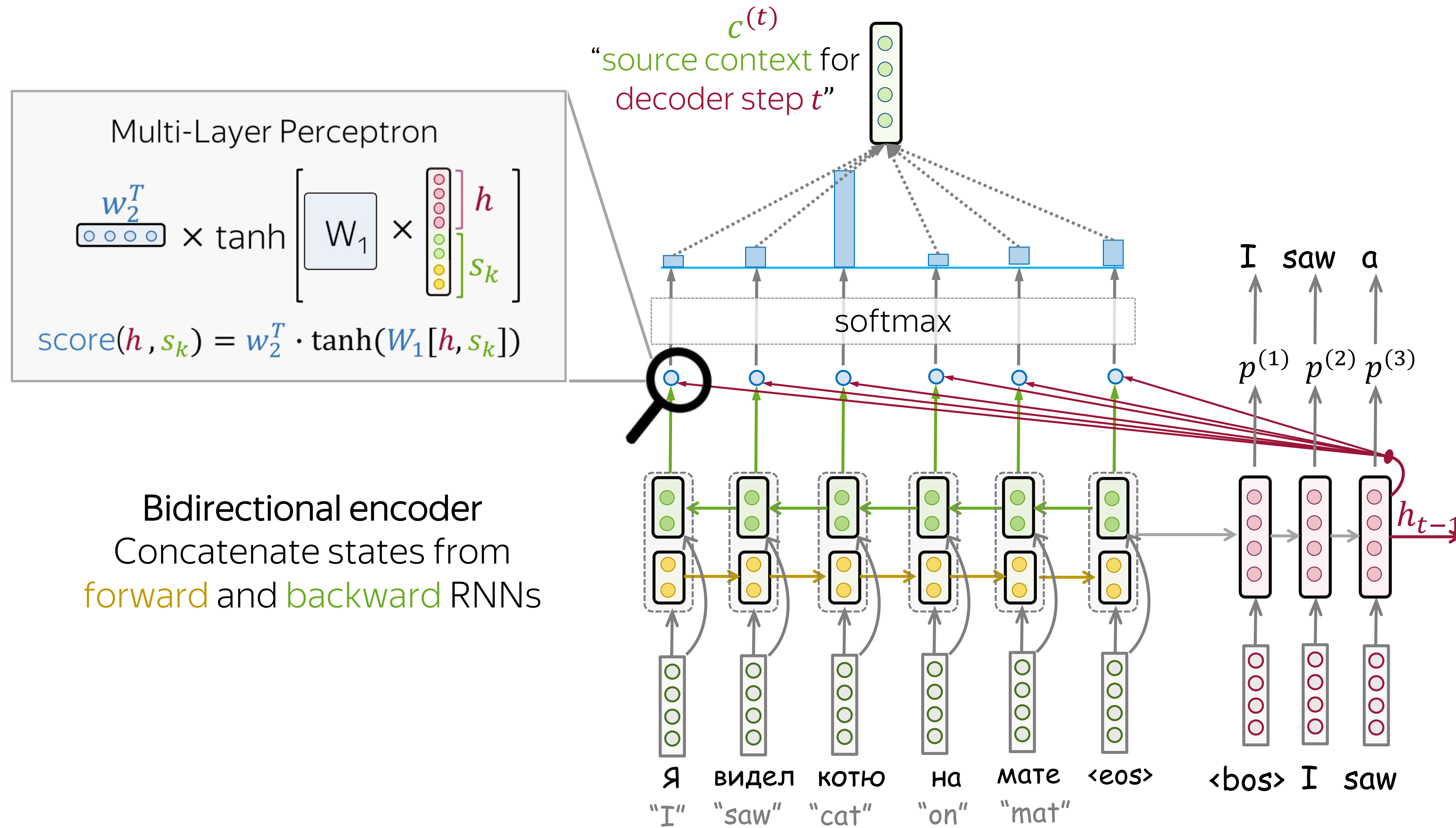
Bahdanau Model (the original attention model)



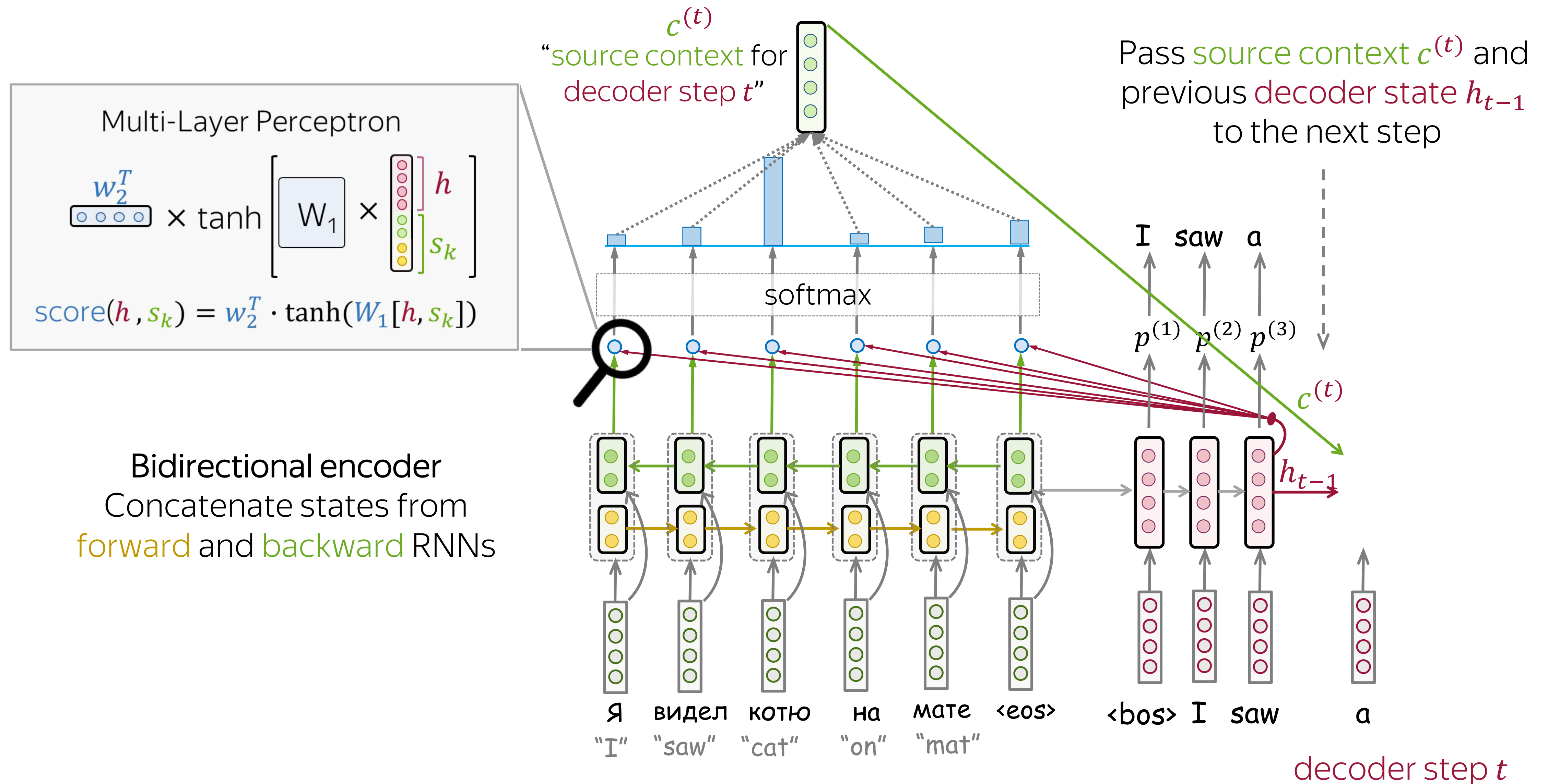
Bahdanau Model (the original attention model)



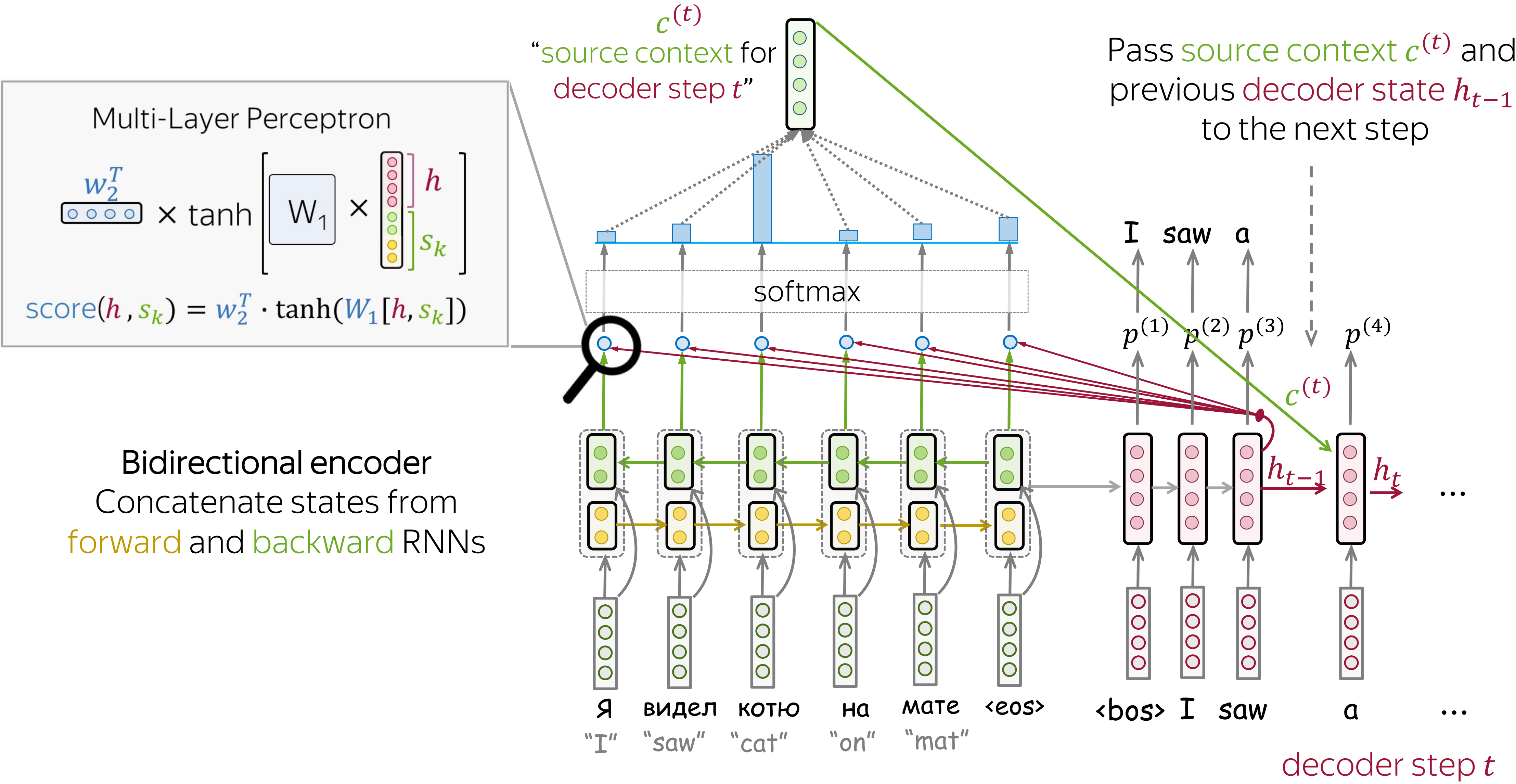
Bahdanau Model (the original attention model)



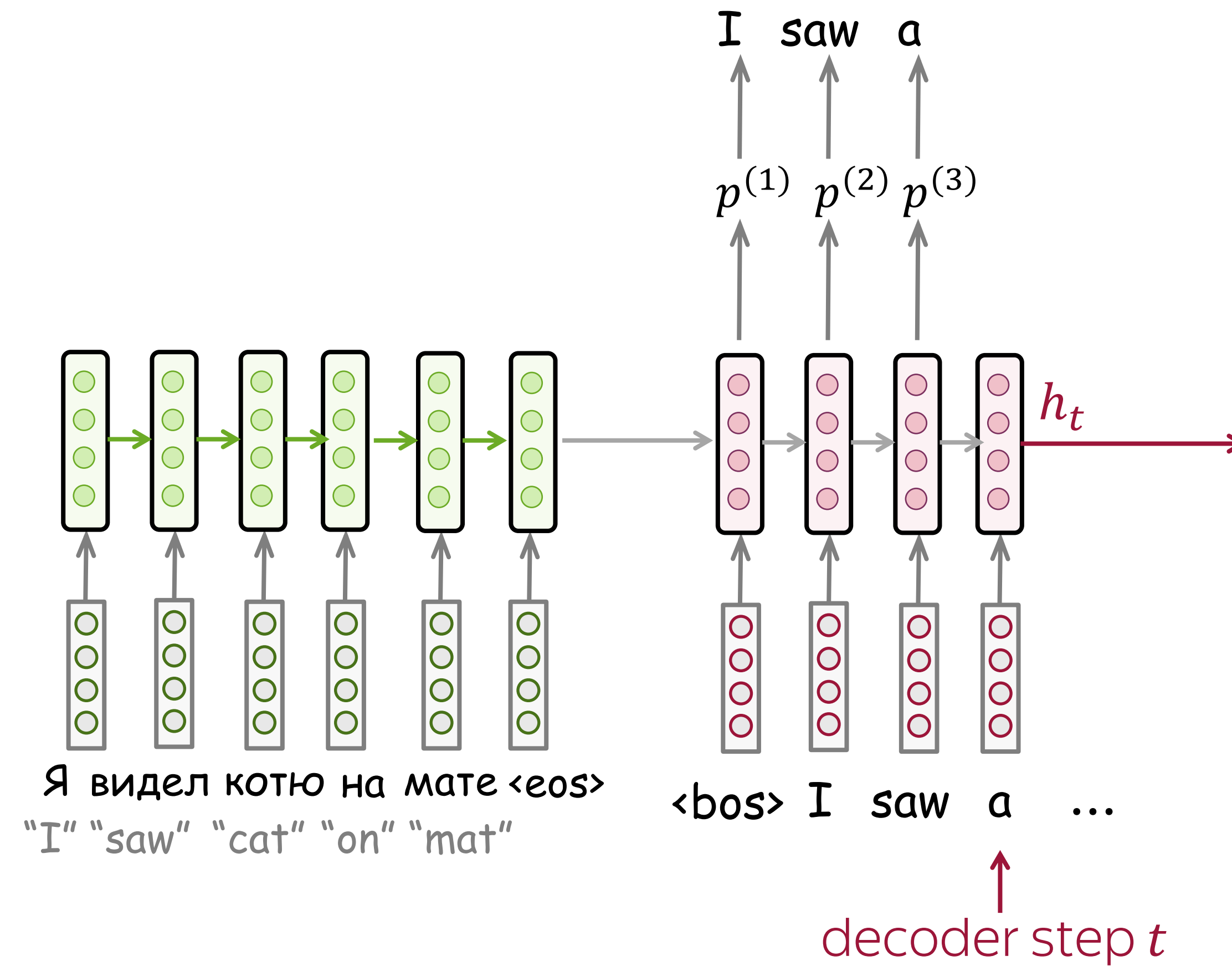
Bahdanau Model (the original attention model)



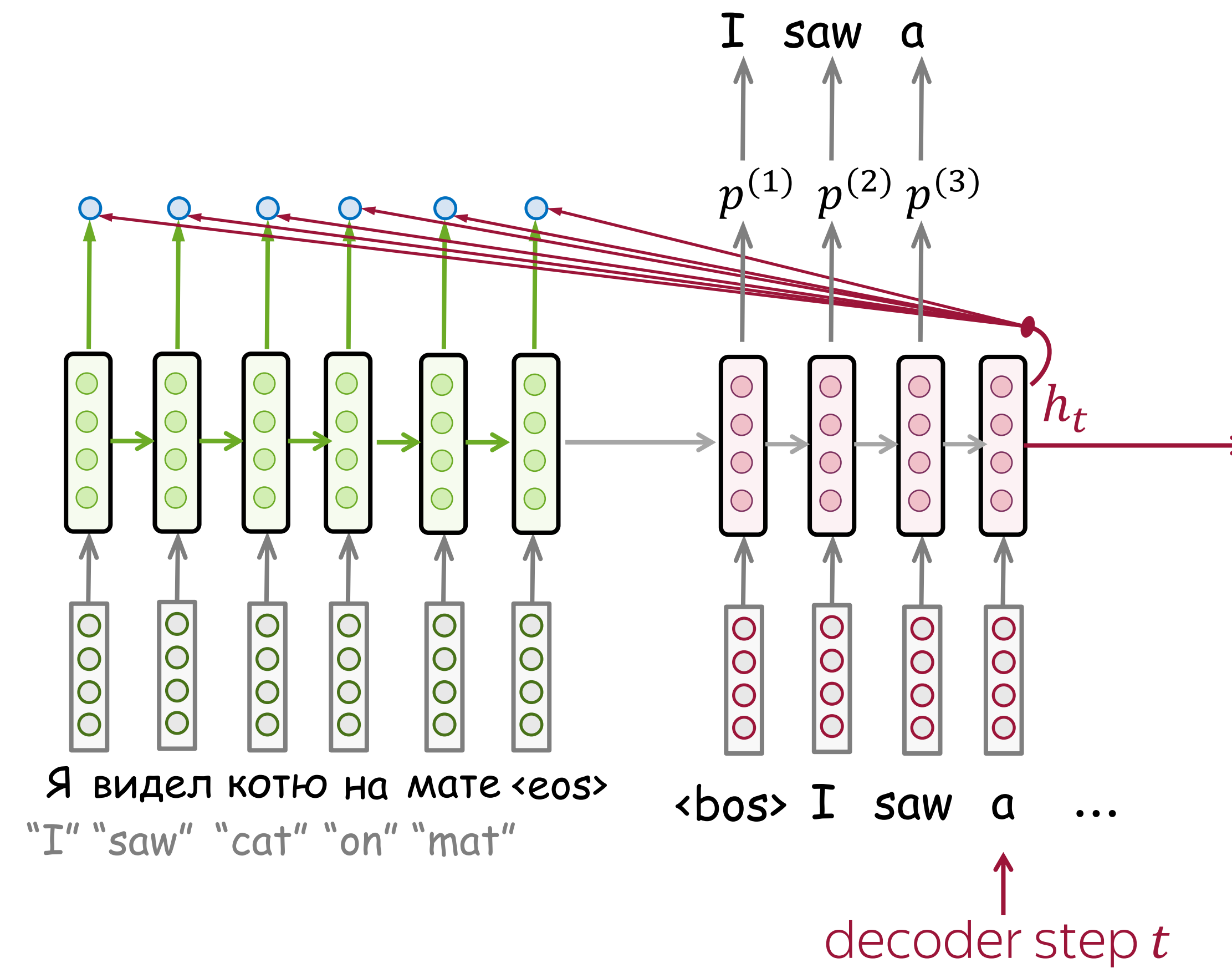
Bahdanau Model (the original attention model)



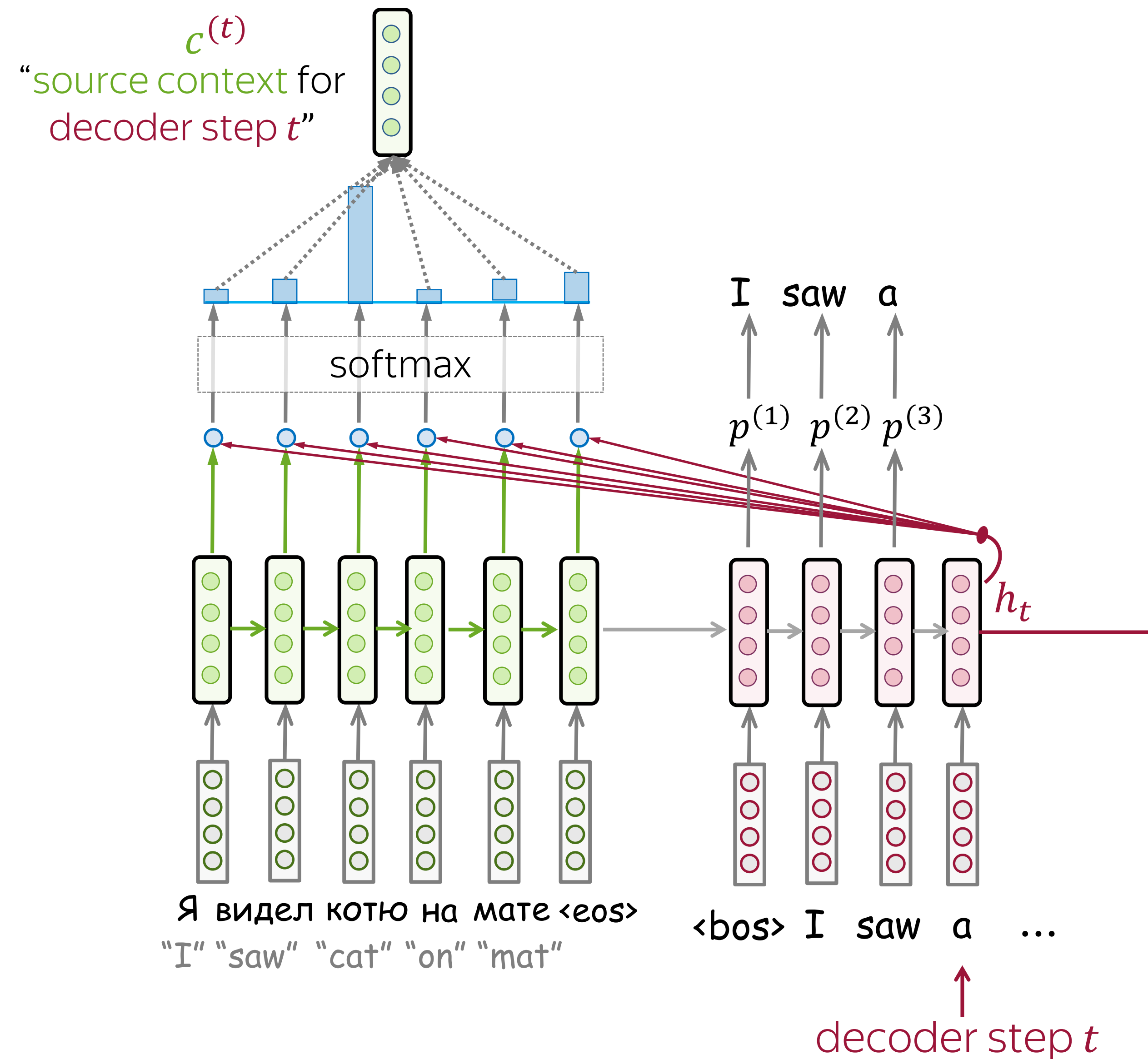
Luong Model



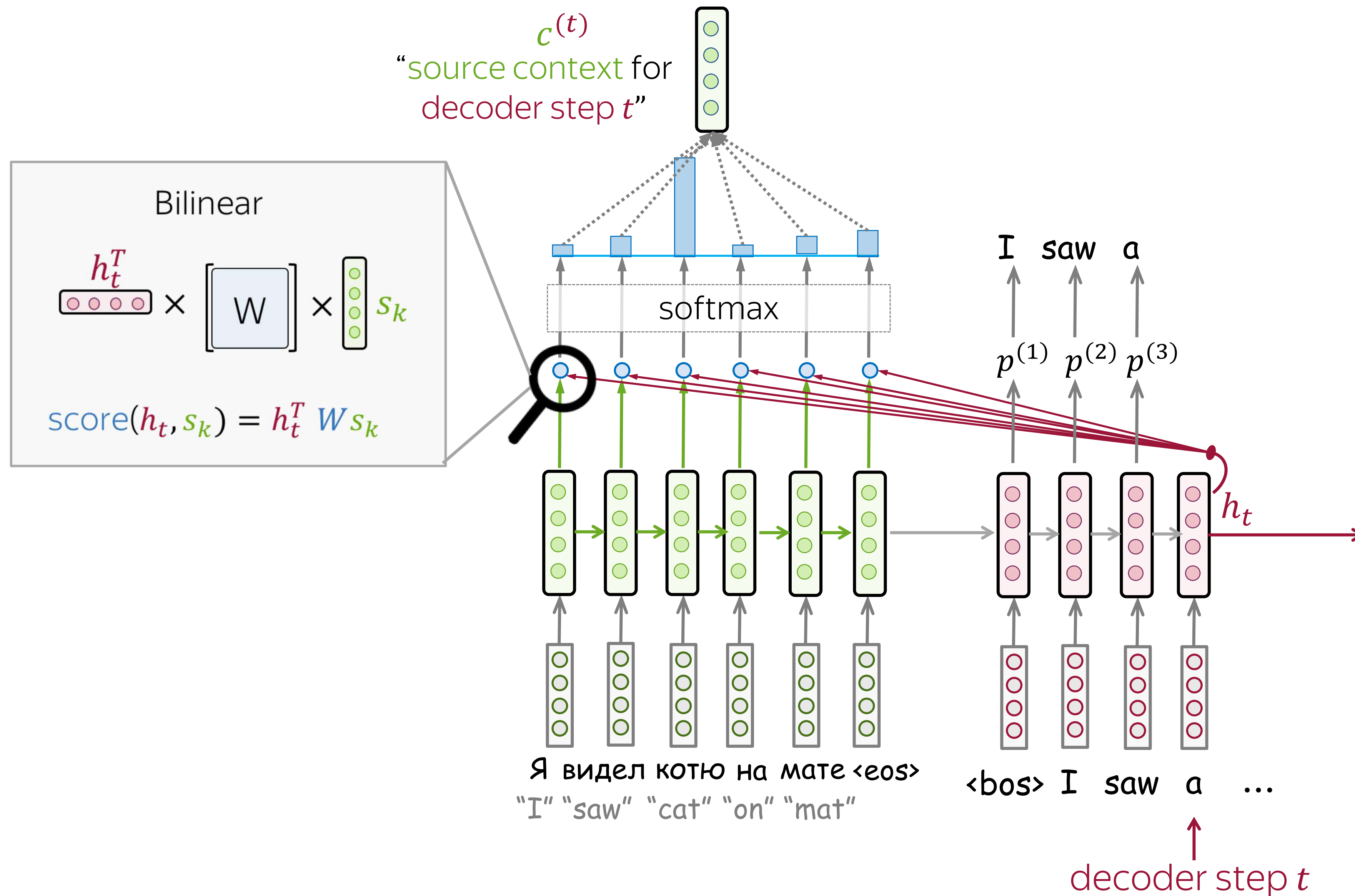
Luong Model



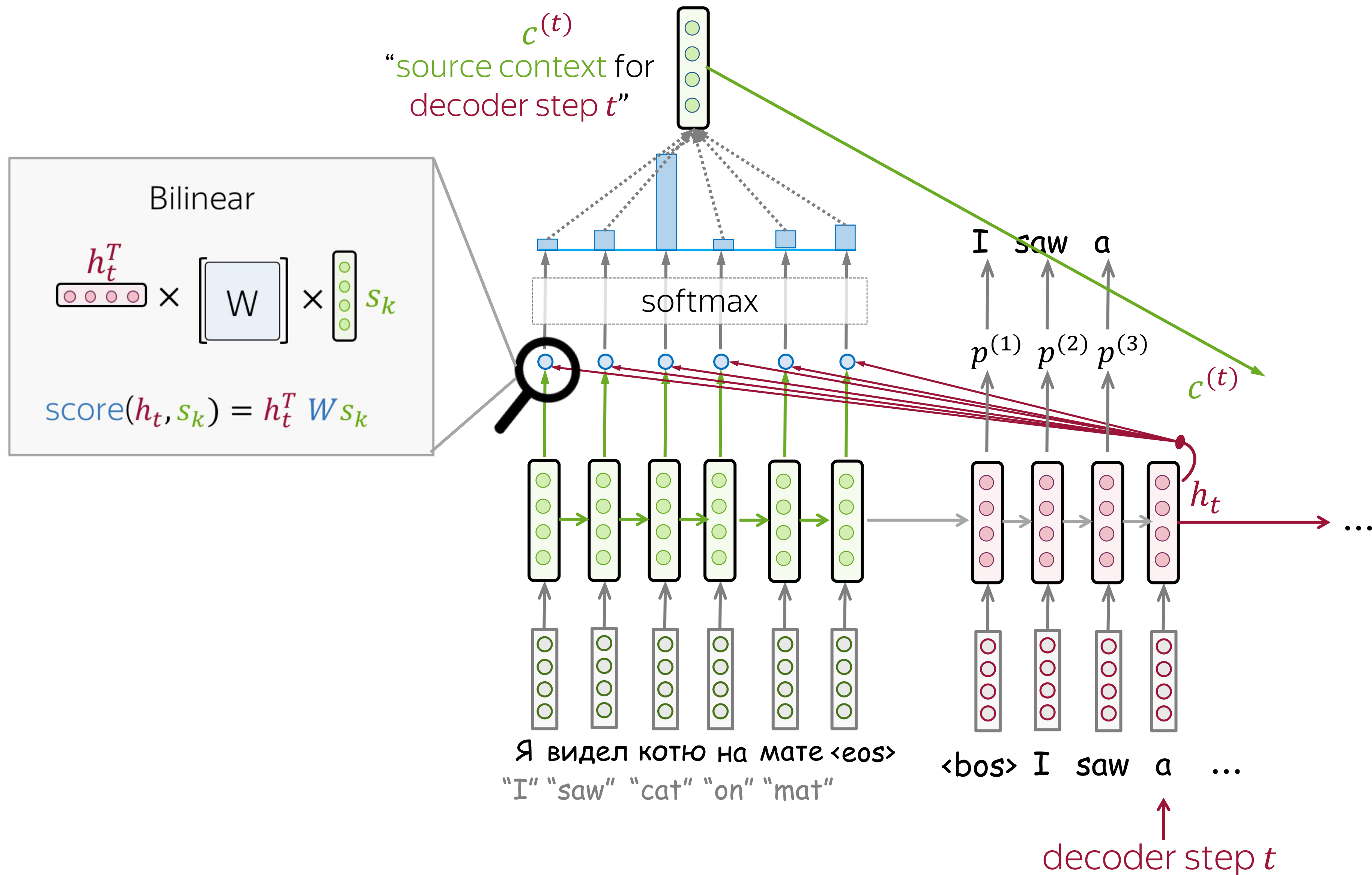
Luong Model



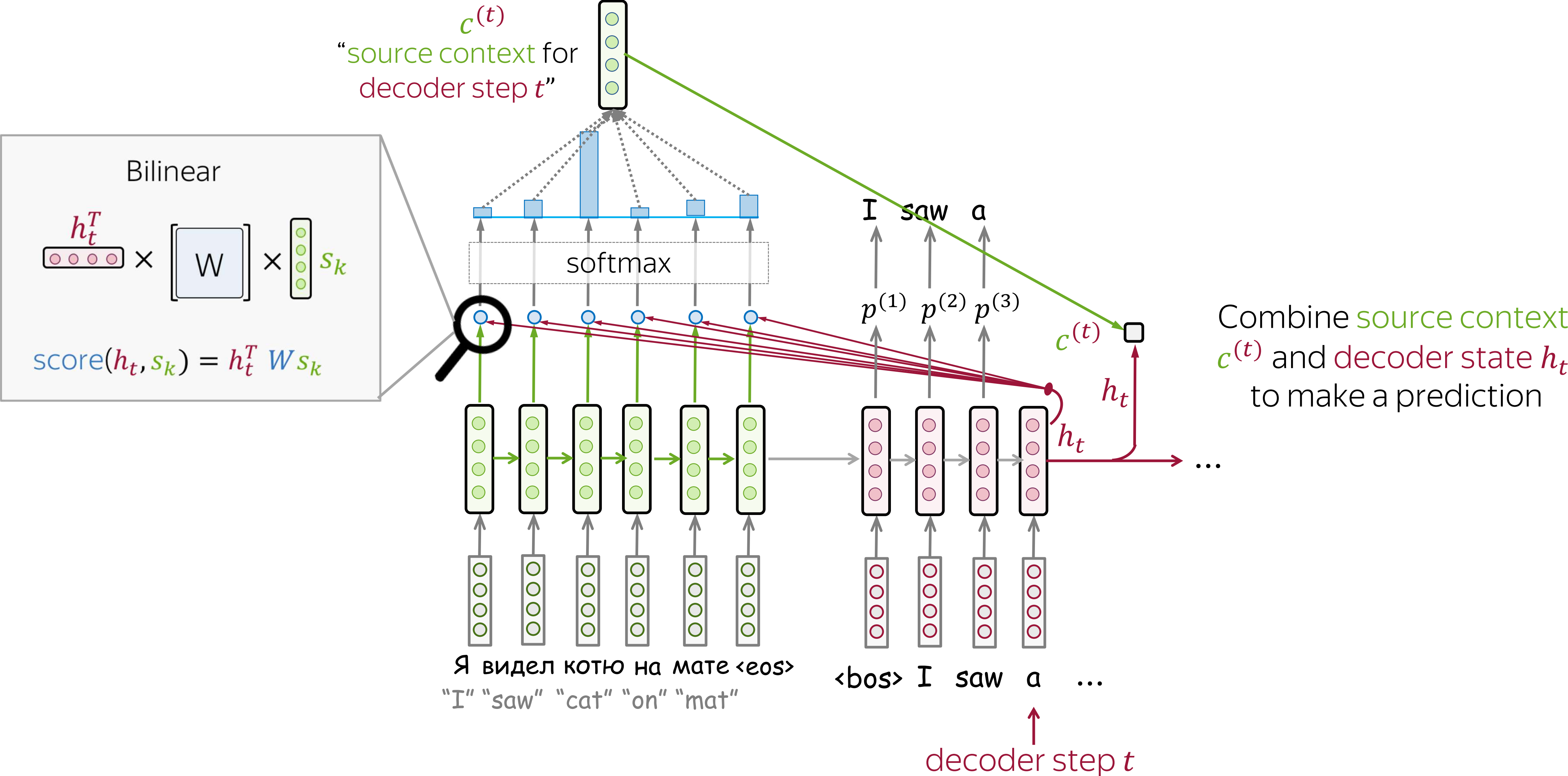
Luong Model



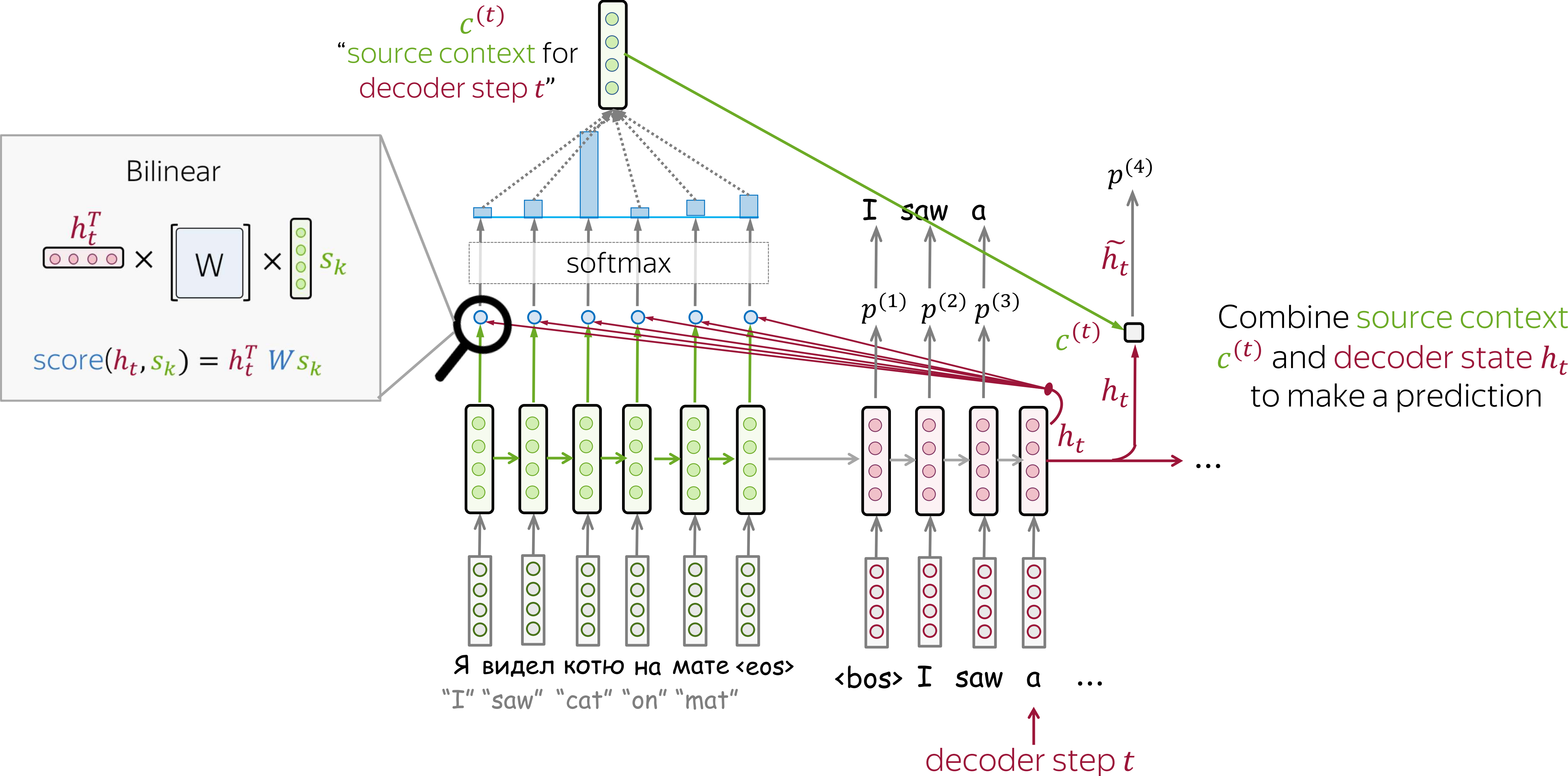
Luong Model



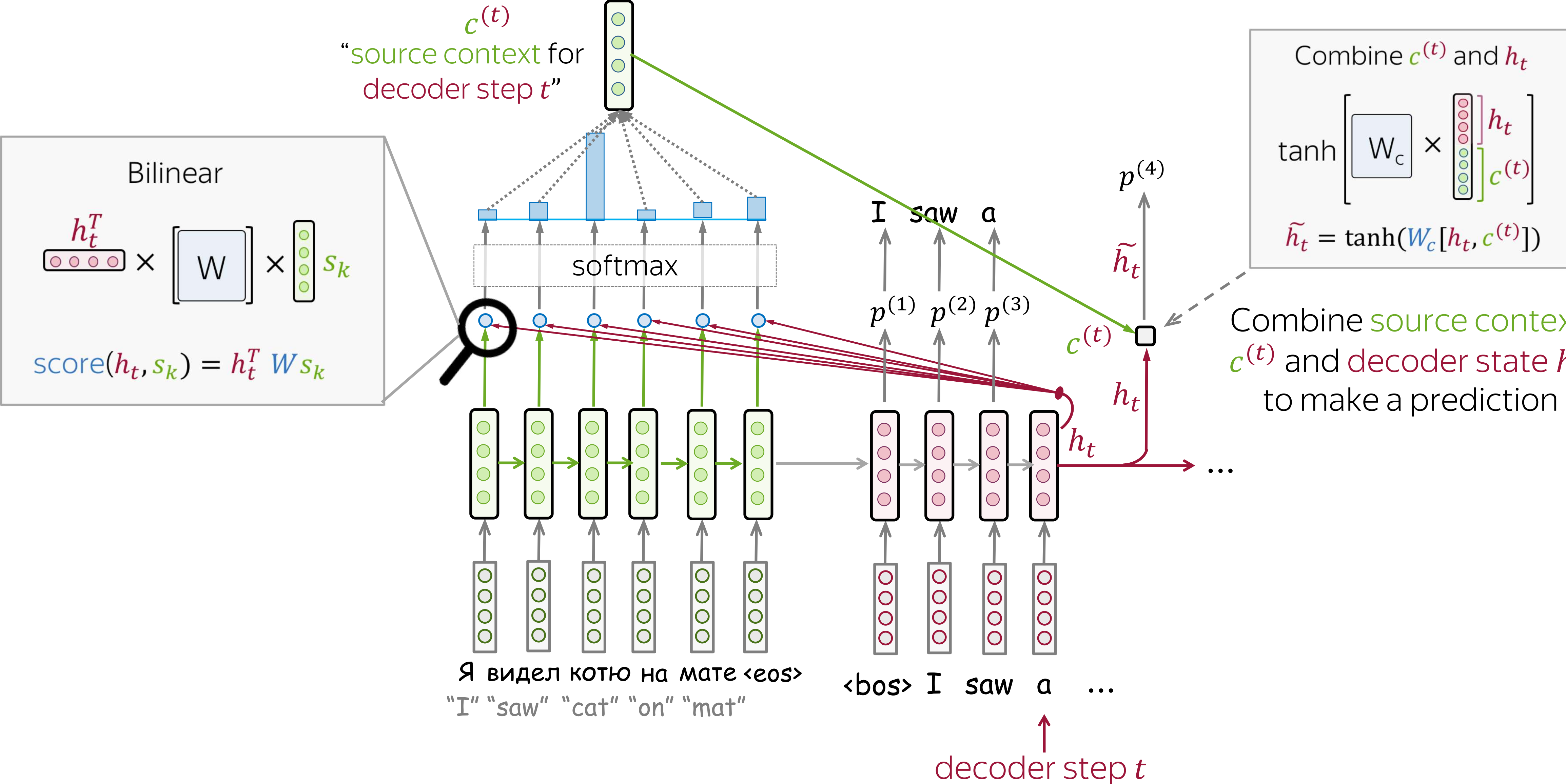
Luong Model



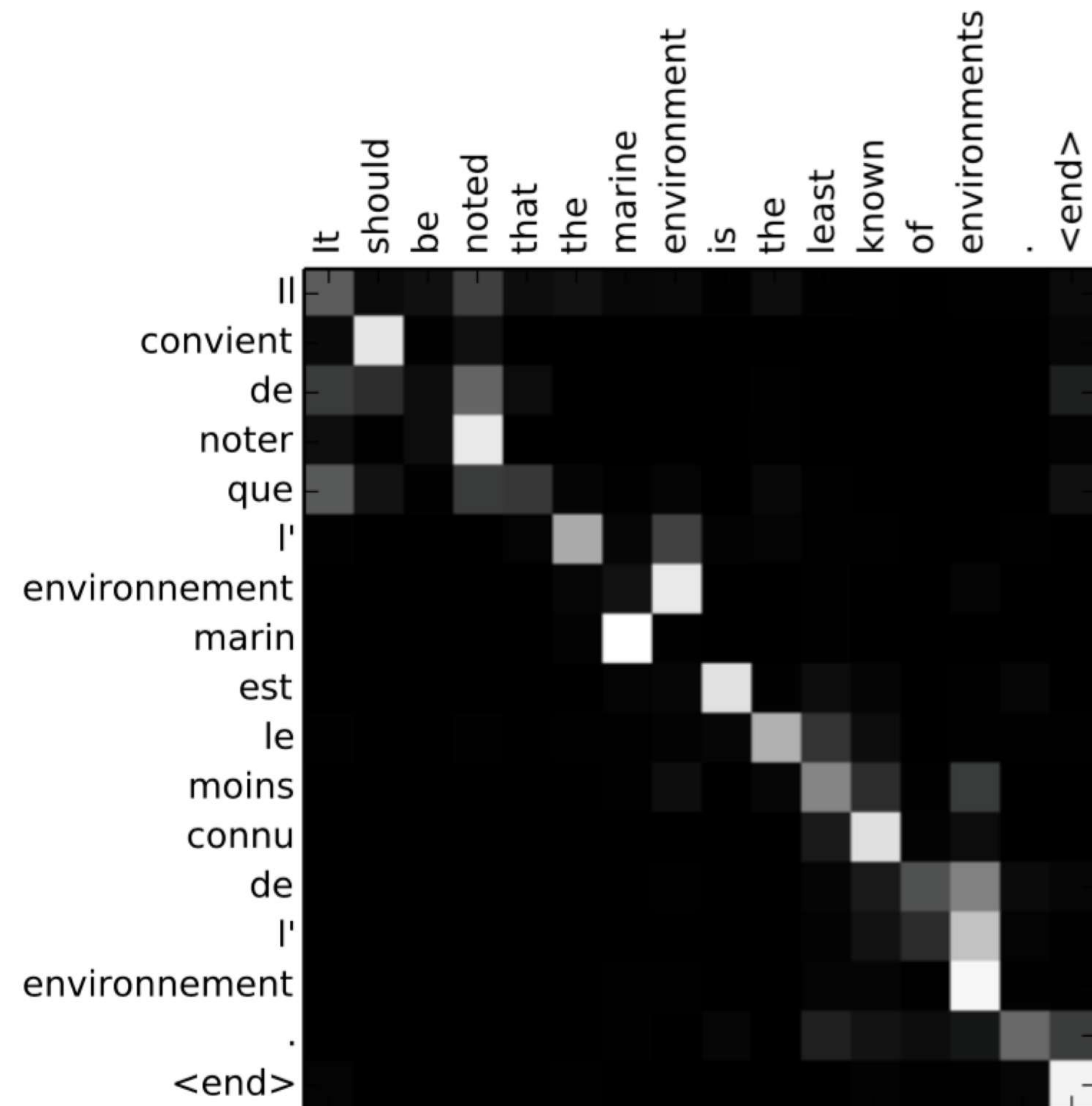
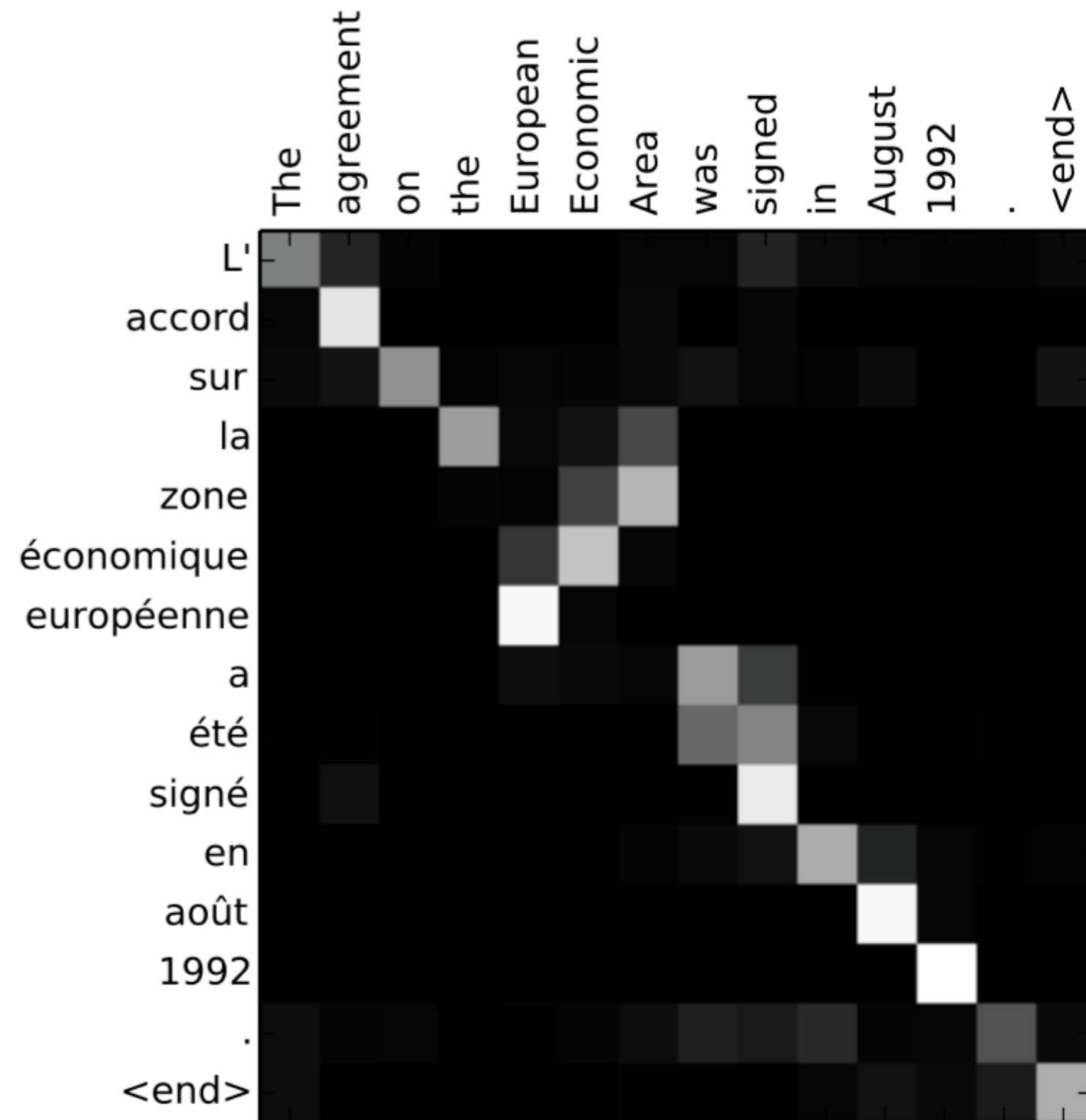
Luong Model



Luong Model



Attention Learned (almost) Alignment



The examples are from the paper [Neural Machine Translation by Jointly Learning to Align and Translate](#).