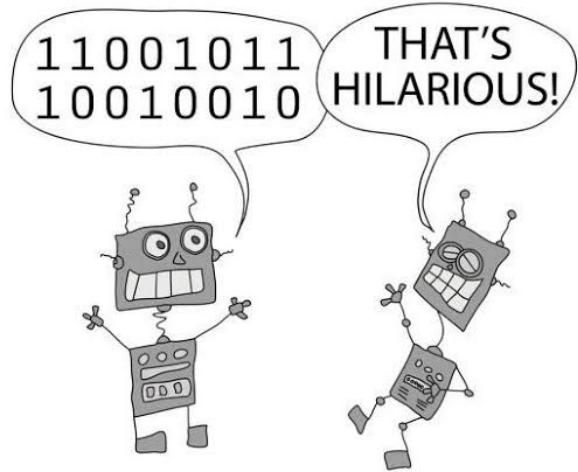


# THIS IS CS5045!

GCR: ioc7cdl

# MACHINE TRANSLATION

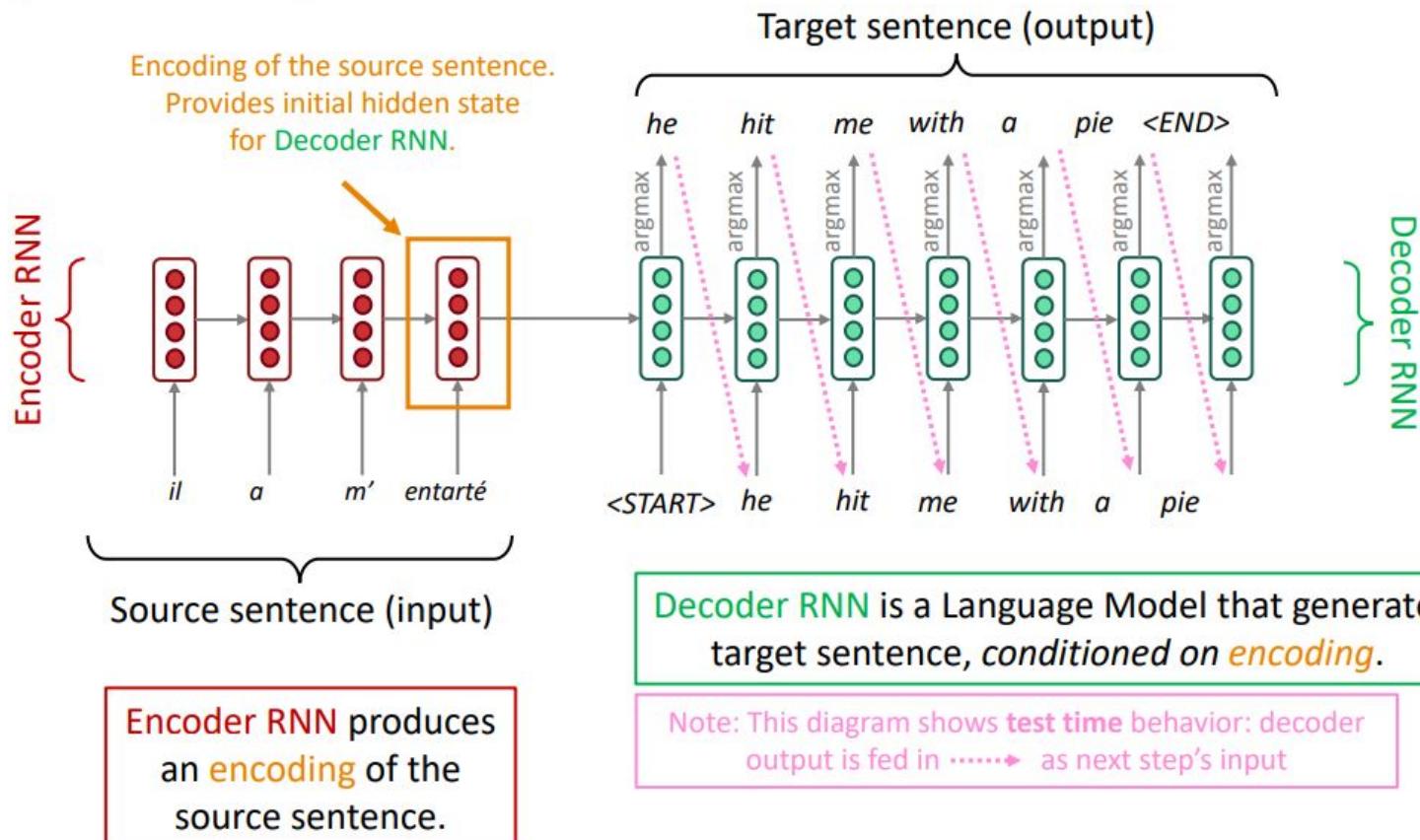


# What is Neural Machine Translation?

- Neural Machine Translation (NMT) is a way to do Machine Translation with a *single end-to-end neural network*
- The neural network architecture is called a sequence-to-sequence model (aka seq2seq) and it involves *two RNNs*

# Neural Machine Translation (NMT)

## The sequence-to-sequence model



## Sequence-to-sequence is versatile!

- The general notion here is an **encoder-decoder** model
  - One neural network takes input and produces a neural representation
  - Another network produces output based on that neural representation
  - If the input and output are sequences, we call it a seq2seq model
- Sequence-to-sequence is useful for *more than just MT*
- Many NLP tasks can be phrased as sequence-to-sequence:
  - **Summarization** (long text → short text)
  - **Dialogue** (previous utterances → next utterance)
  - **Parsing** (input text → output parse as sequence)
  - **Code generation** (natural language → Python code)

# Neural Machine Translation (NMT)

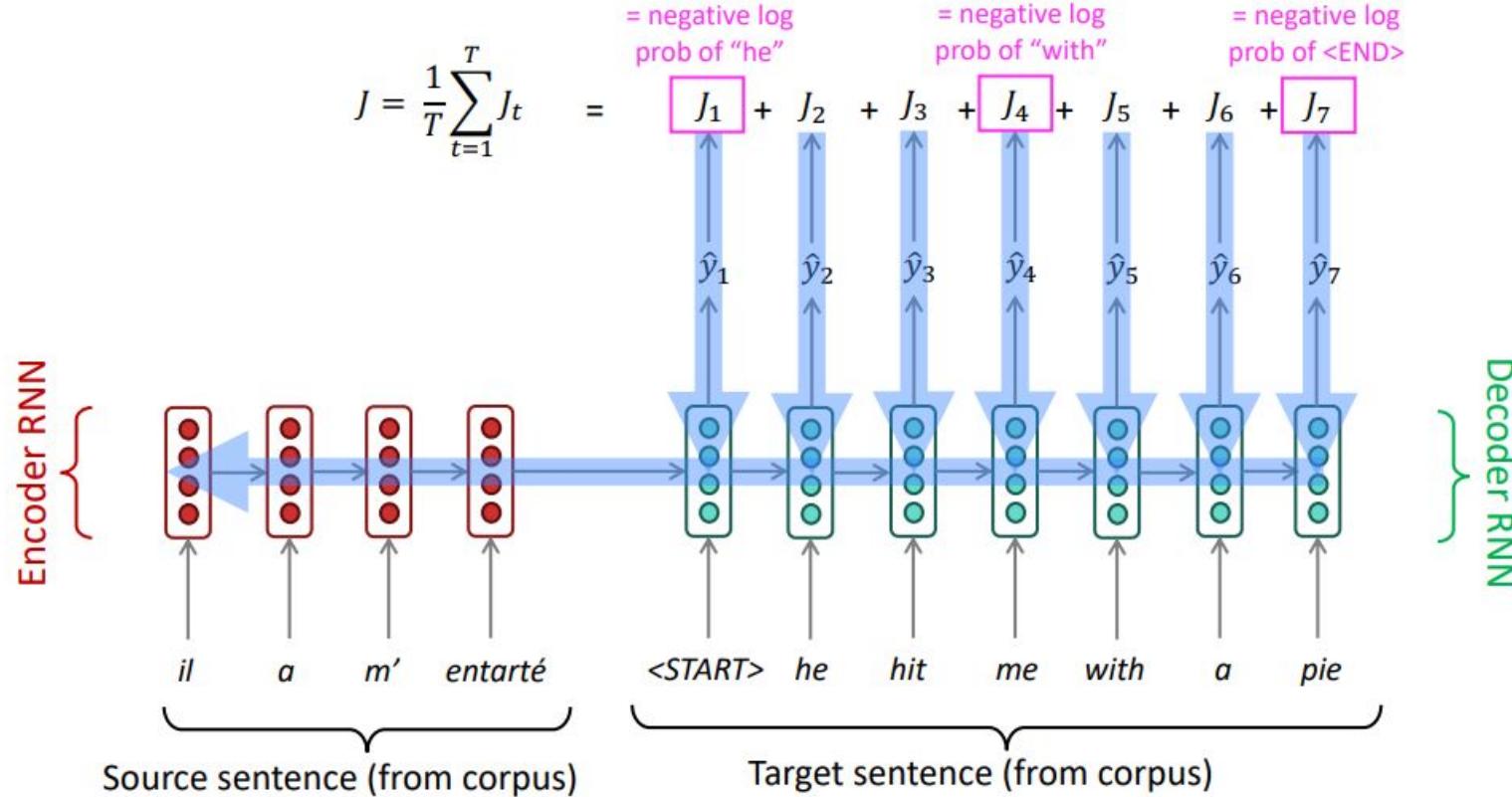
- The sequence-to-sequence model is an example of a **Conditional Language Model**
  - **Language Model** because the decoder is predicting the next word of the target sentence  $y$
  - **Conditional** because its predictions are *also* conditioned on the source sentence  $x$
- NMT directly calculates  $P(y|x)$ :

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

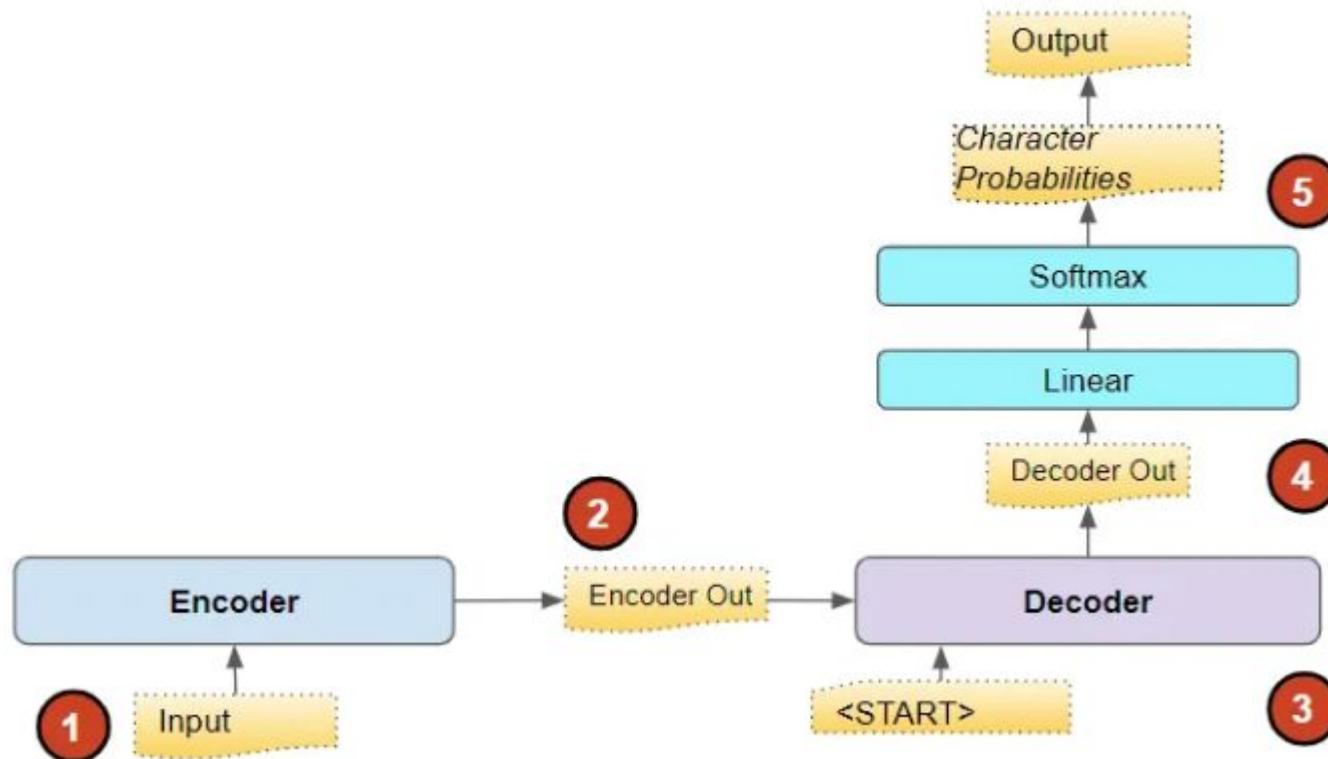
  
Probability of next target word, given  
target words so far and source sentence  $x$

- **Question:** How to train an NMT system?
- **(Easy) Answer:** Get a big parallel corpus...
  - But there is now exciting work on “unsupervised NMT”, data augmentation, etc.

# Training a Neural Machine Translation system



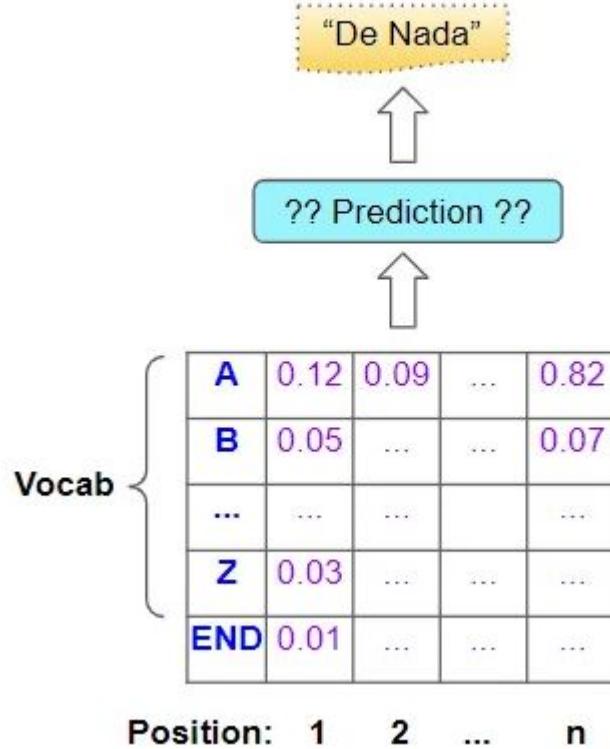
Seq2seq is optimized as a **single system**. Backpropagation operates “end-to-end”.



Vocab

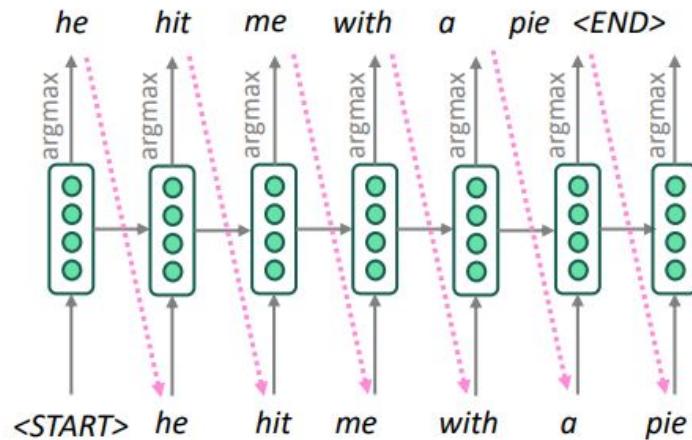
<b>A</b>	0.12	0.09	...	0.82
<b>B</b>	0.05	...	...	0.07
...	...	...		...
<b>Z</b>	0.03	...	...	...
<b>END</b>	0.01	...	...	...

Position: 1    2    ...    n



# Greedy decoding

- We saw how to generate (or “decode”) the target sentence by taking argmax on each step of the decoder



- This is **greedy decoding** (take most probable word on each step)
- Problems with this method?**

## Problems with greedy decoding

- Greedy decoding has no way to undo decisions!

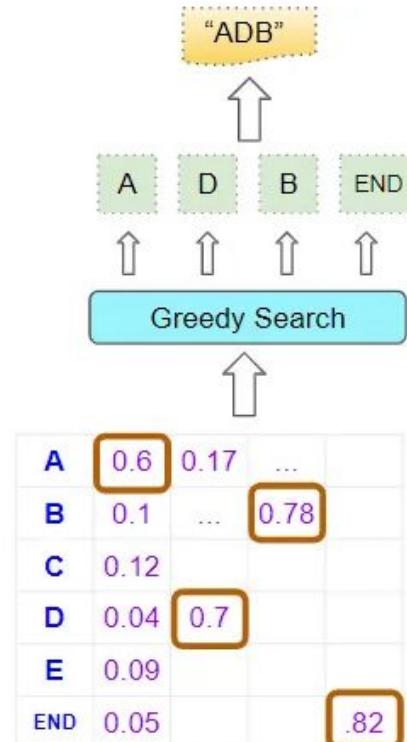
- Input: *il a m'entarté*      (*he hit me with a pie*)

- → *he* \_\_\_\_\_

- → *he hit* \_\_\_\_\_

- → *he hit a* \_\_\_\_\_

(*whoops! no going back now...*)



Greedy Search (Image by Author)

# How do we evaluate Machine Translation?

## BLEU (Bilingual Evaluation Understudy)

- BLEU compares the machine-written translation to one or several human-written translation(s), and computes a **similarity score** based on:
  - *n*-gram precision (usually for 1, 2, 3 and 4-grams)
  - Plus a penalty for too-short system translations
- BLEU is **useful** but **imperfect**
  - There are many valid ways to translate a sentence
  - So a **good** translation can get a **poor** BLEU score because it has low *n*-gram overlap with the human translation 😞

# BLEU score against 4 reference translations

## Reference translation 1:

The U.S. island of Guam is maintaining a high state of alert after the Guam airport and its offices both received an e-mail from someone calling himself the Saudi Arabian Osama bin Laden and threatening a biological/chemical attack against public places such as the airport.

## Reference translation 2:

Guam International Airport and its offices are maintaining a high state of alert after receiving an e-mail that was from a person claiming to be the wealthy Saudi Arabian businessman Bin Laden and that threatened to launch a biological and chemical attack on the airport and other public places .

## Machine translation:

The American [?] international airport and its office all receives one calls self the sand Arab rich business [?] and so on electronic mail , which sends out ; The threat will be able after public place and so on the airport to start the biochemical attack . [?] highly alerts after the maintenance.

## Reference translation 3:

The US International Airport of Guam and its office has received an email from a self-claimed Arabian millionaire named Laden , which threatens to launch a biochemical attack on such public places as airport . Guam authority has been on alert .

## Reference translation 4:

US Guam International Airport and its office received an email from Mr. Bin Laden and other rich businessman from Saudi Arabia . They said there would be biochemical air raid to Guam Airport and other public places . Guam needs to be in high precaution about this matter .

[Papineni et al. 2002]

# BLEU SCORE

$$\text{BP} = \begin{cases} 1 & \text{if } c > r \\ e^{(1-r/c)} & \text{if } c \leq r \end{cases} .$$

**Reference Text:** It was raining heavily today

Then,

**Predicted Text:** It It It is raining heavily

Unigram:

Clipped Precision count= 3/6=1/2

Bigrams for reference text: ["It was", "was raining", "raining heavily", "heavily today"]

Bigrams for predicted text: ["It It", "It It", "It is", "is raining", "raining heavily"]

Clipped Precision: 1/5

$$\text{BLEU} = \text{BP} \cdot \exp \left( \sum_{n=1}^N w_n \log p_n \right) .$$

# BLEU SCORE

$$\text{Global Average Precision} = \exp\left(\sum_{n=1}^N w_n \log p_n\right) = \prod_{n=1}^N p_n^{w_n} = (p_1)^{1/2} \cdot (p_2)^{1/2}$$

Usually we use **N=4** and  $w_n = 1/4$ , but in this case we will use **N=2** and  $w_n = 1/2$

$$\text{Brevity Penalty} = \begin{cases} 1, & \text{if } c > r \\ e^{(1-r/c)}, & \text{if } c \leq r \end{cases}$$

r = number of words in the reference text

c = number of words in the predicted text

Brevity penalty cannot be larger than 1, even if the predicted text is larger than the reference text.

r = 5 c = 6

Bleu = 0.316

# IMPLEMENTATION FOR BLEU SCORE WITH PYTHON

```
from nltk.translate.bleu_score import sentence_bleu  
  
reference = ['It was raining heavily today'].split()  
  
candidate = 'It It It is raining heavily'.split()  
  
print(sentence_bleu(reference, candidate, weights=(0.5, 0.5, 0, 0)))
```

<https://www.scaler.com/topics/nlp/bleu-score-in-nlp/>

# NMT: perhaps the biggest success story of NLP Deep Learning?

Neural Machine Translation went from a **fringe research attempt** in **2014** to the **leading standard method** in **2016**

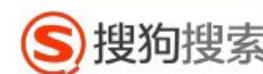
- **2014:** First seq2seq paper published
- **2016:** Google Translate switches from SMT to NMT – and by 2018 everyone has



Microsoft



Tencent 腾讯



- This is amazing!
  - **SMT** systems, built by **hundreds** of engineers over many **years**, outperformed by NMT systems trained by a **small group** of engineers in a few **months**

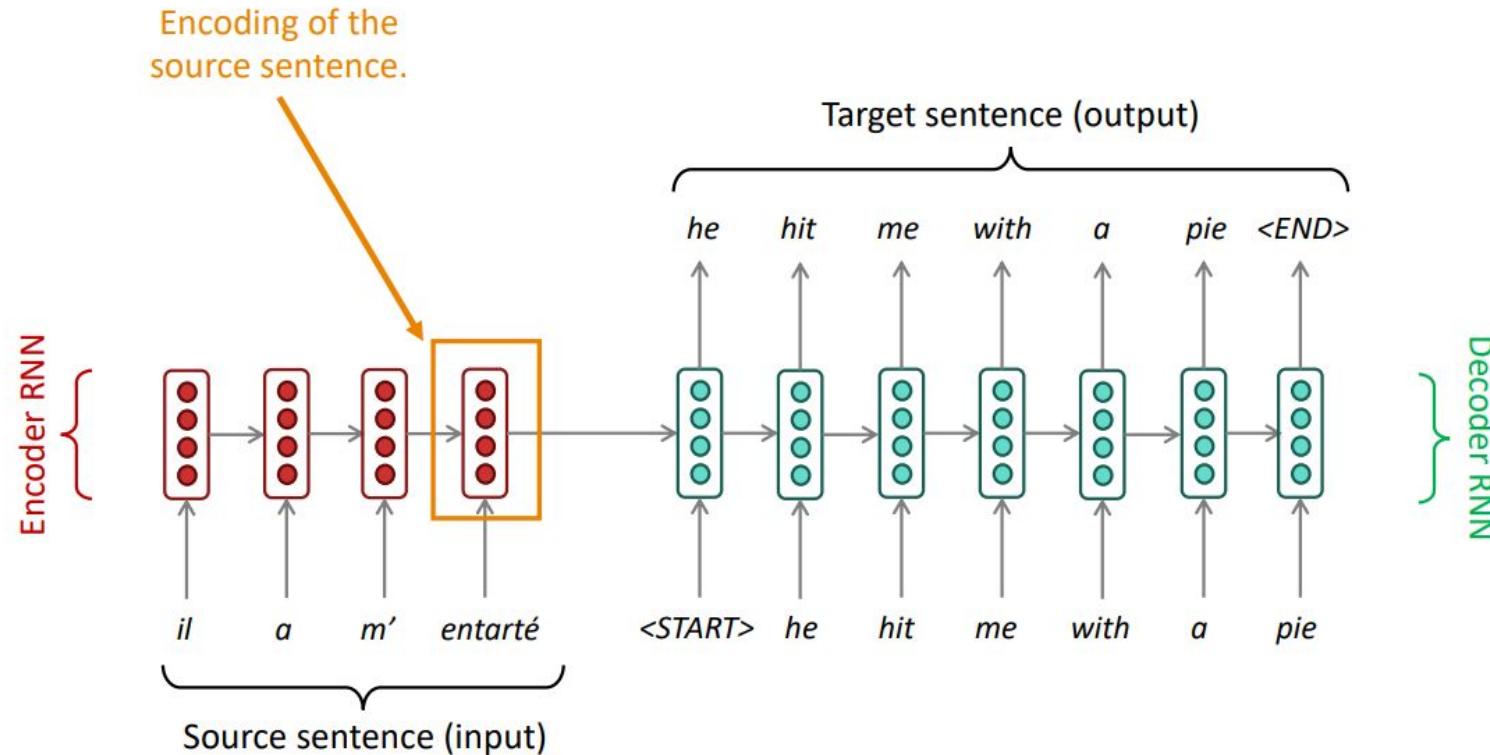
## NMT research continues

NMT is a **flagship task** for NLP Deep Learning

- NMT research has **pioneered** many of the recent **innovations** of NLP Deep Learning
- NMT research continues to **thrive**
  - Researchers have found **many, many improvements** to the “vanilla” seq2seq NMT system we’ve just presented
  - But we’ll present next **one improvement** so integral that it is the new vanilla...

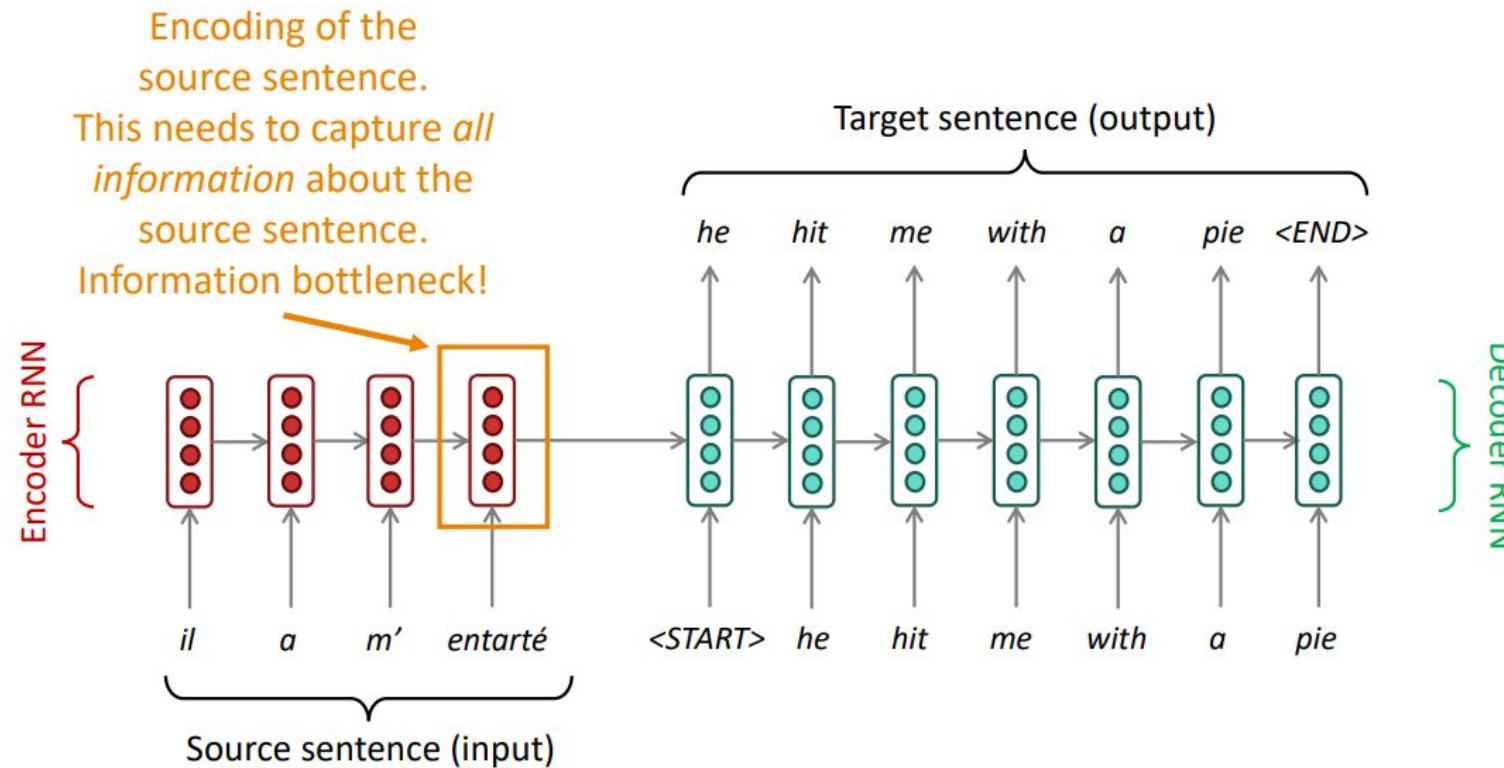
# ATTENTION

# 1. Why attention? Sequence-to-sequence: the bottleneck problem



Problems with this architecture?

# 1. Why attention? Sequence-to-sequence: the bottleneck problem



# Attention

- **Attention** provides a solution to the bottleneck problem.
- Core idea: on each step of the decoder, *use direct connection to the encoder to focus on a particular part* of the source sequence

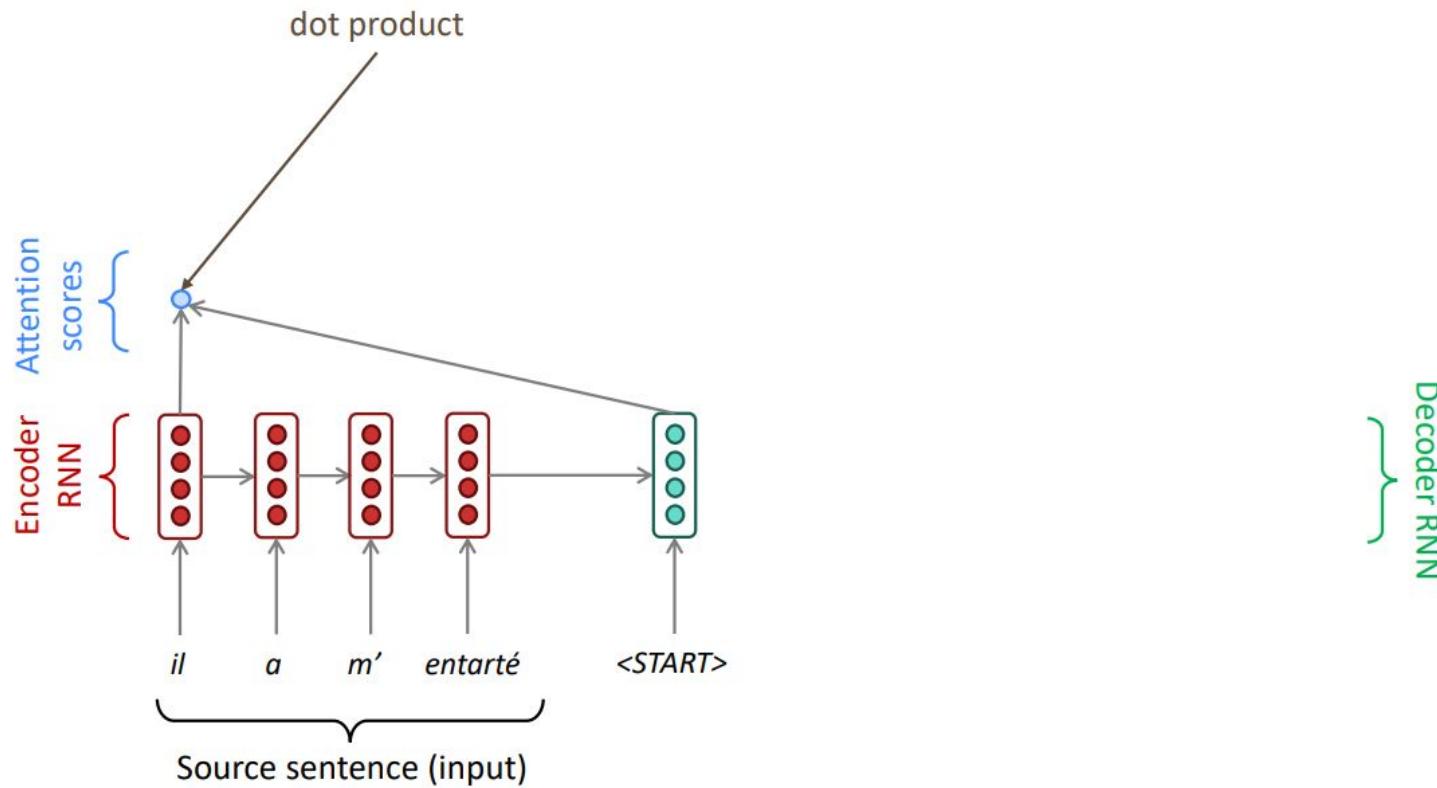


- First, we will show via diagram (no equations), then we will show with equations

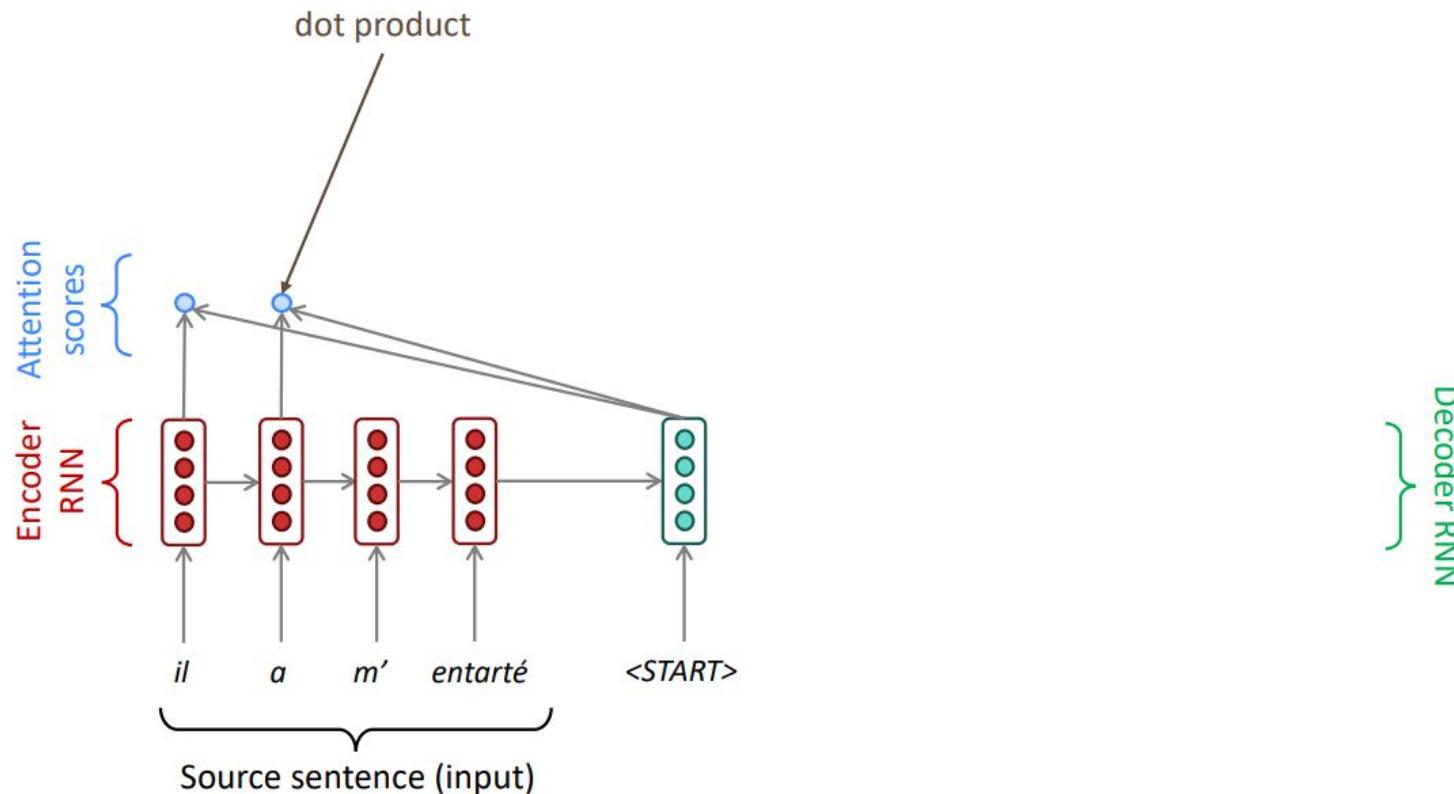
## **Attention:**

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

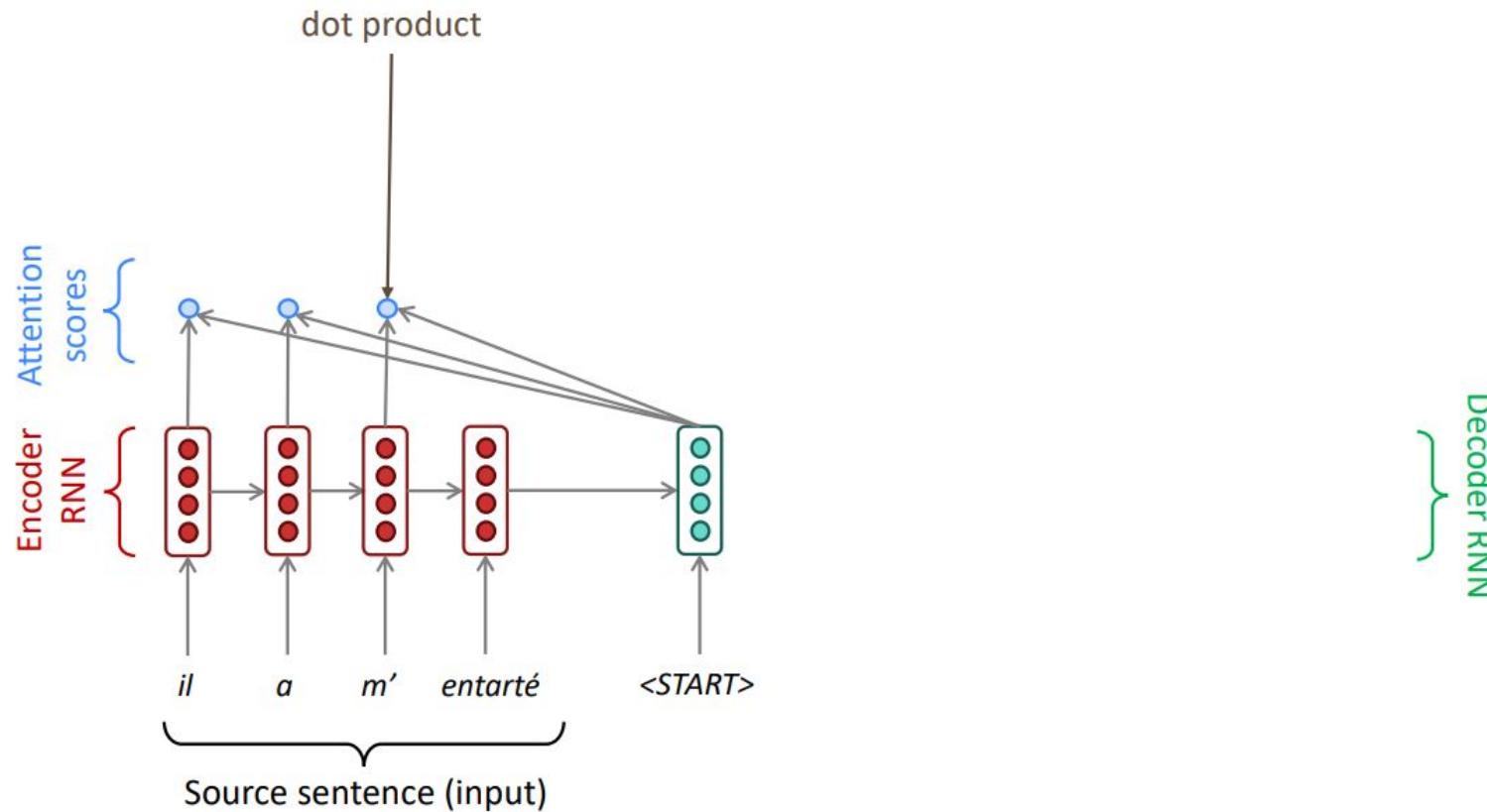
# Sequence-to-sequence with attention



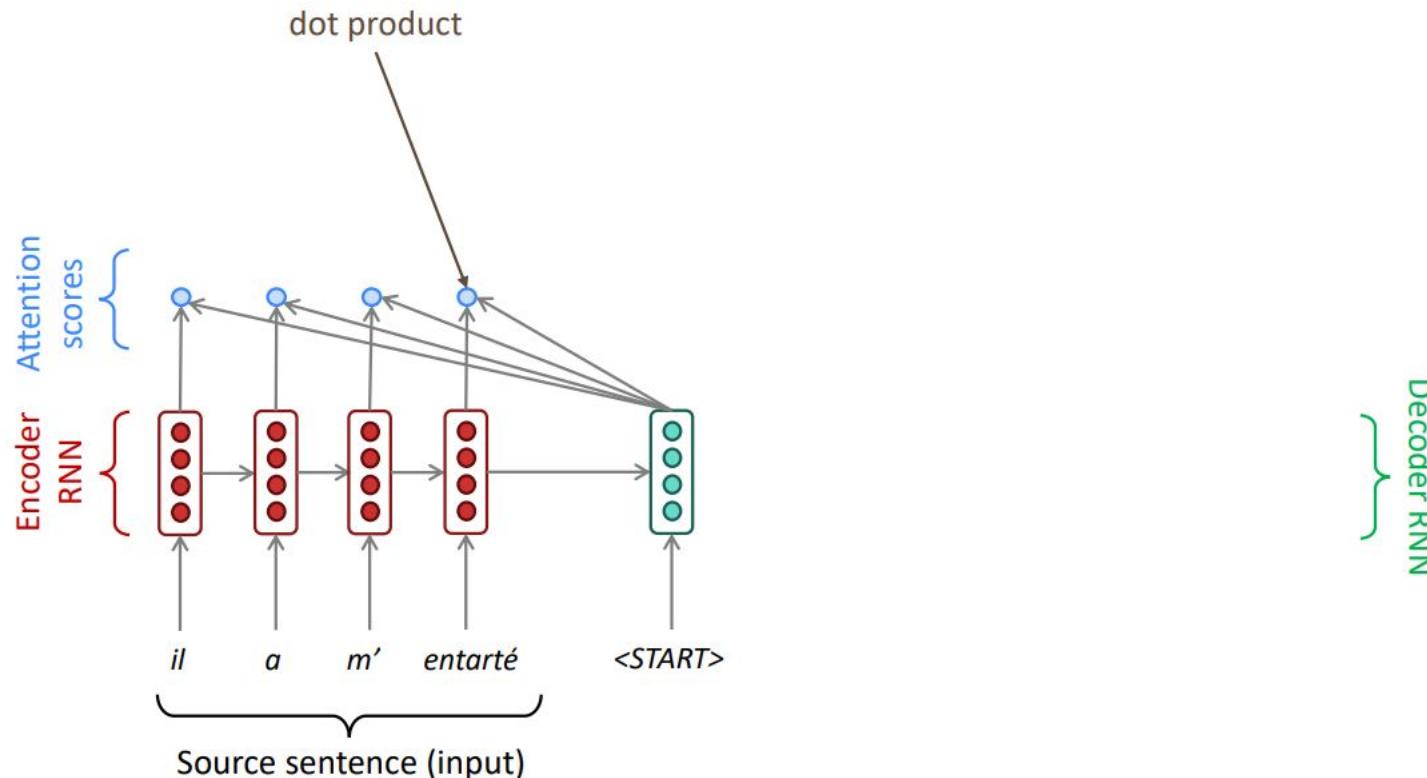
# Sequence-to-sequence with attention



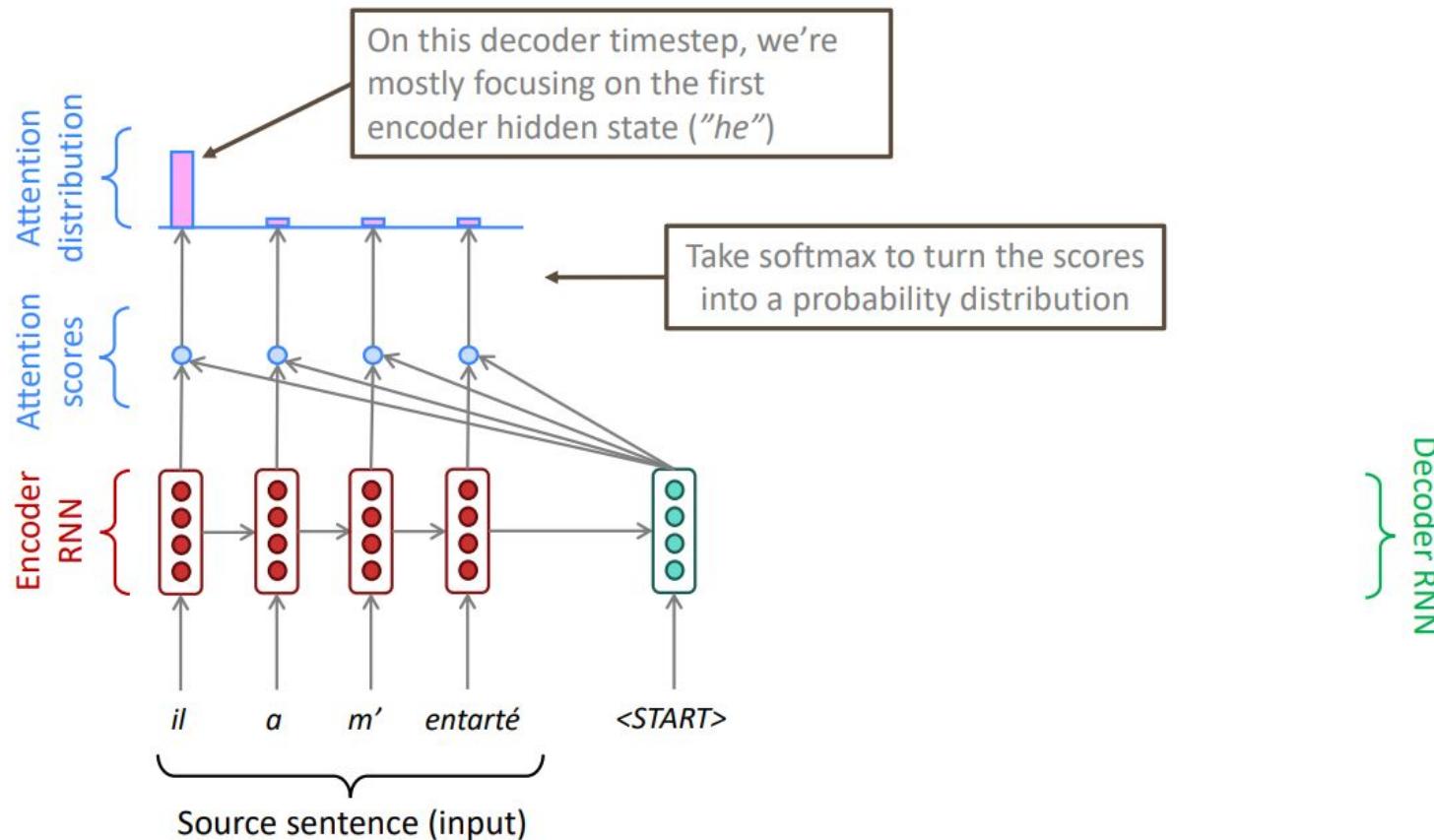
# Sequence-to-sequence with attention



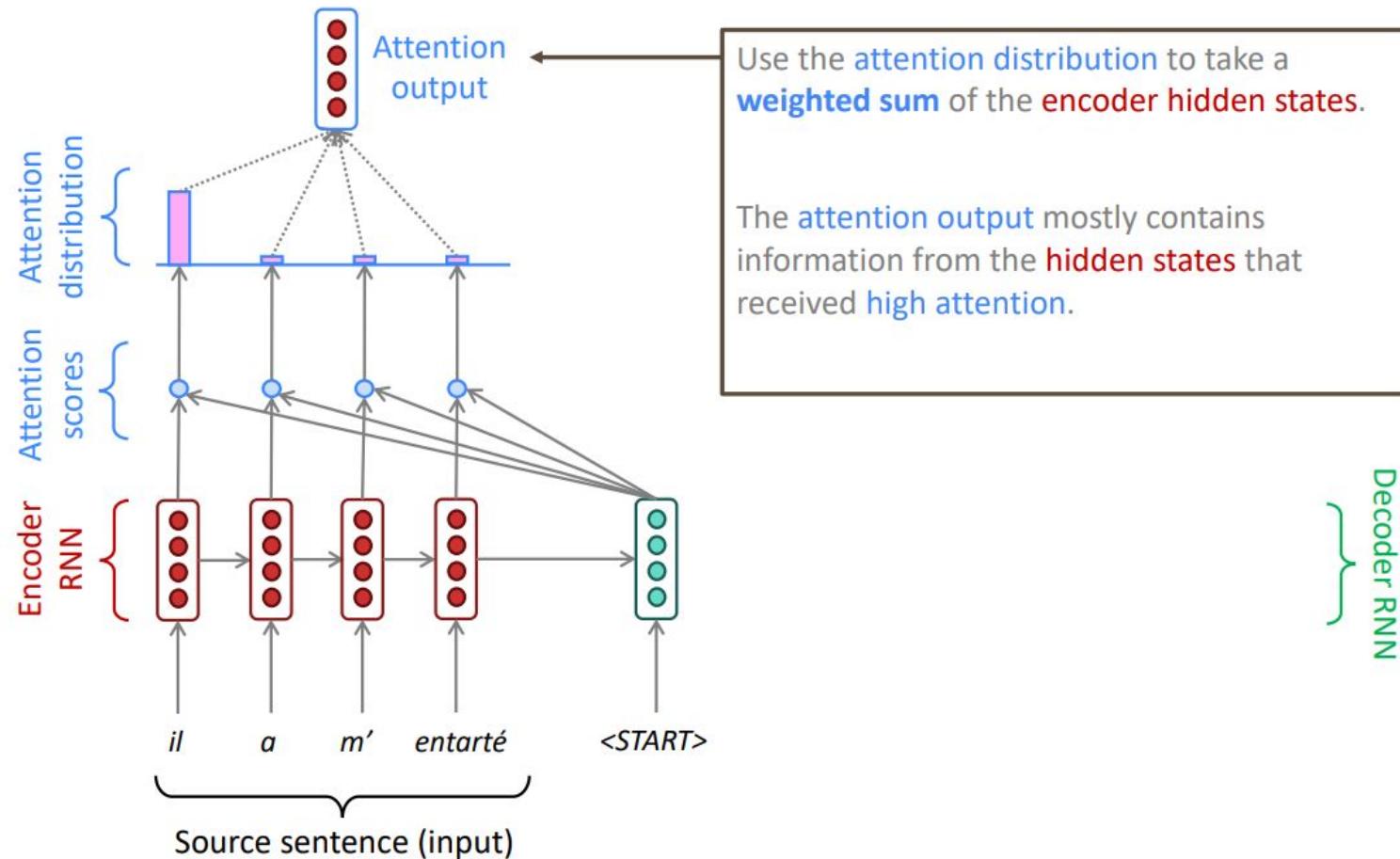
# Sequence-to-sequence with attention



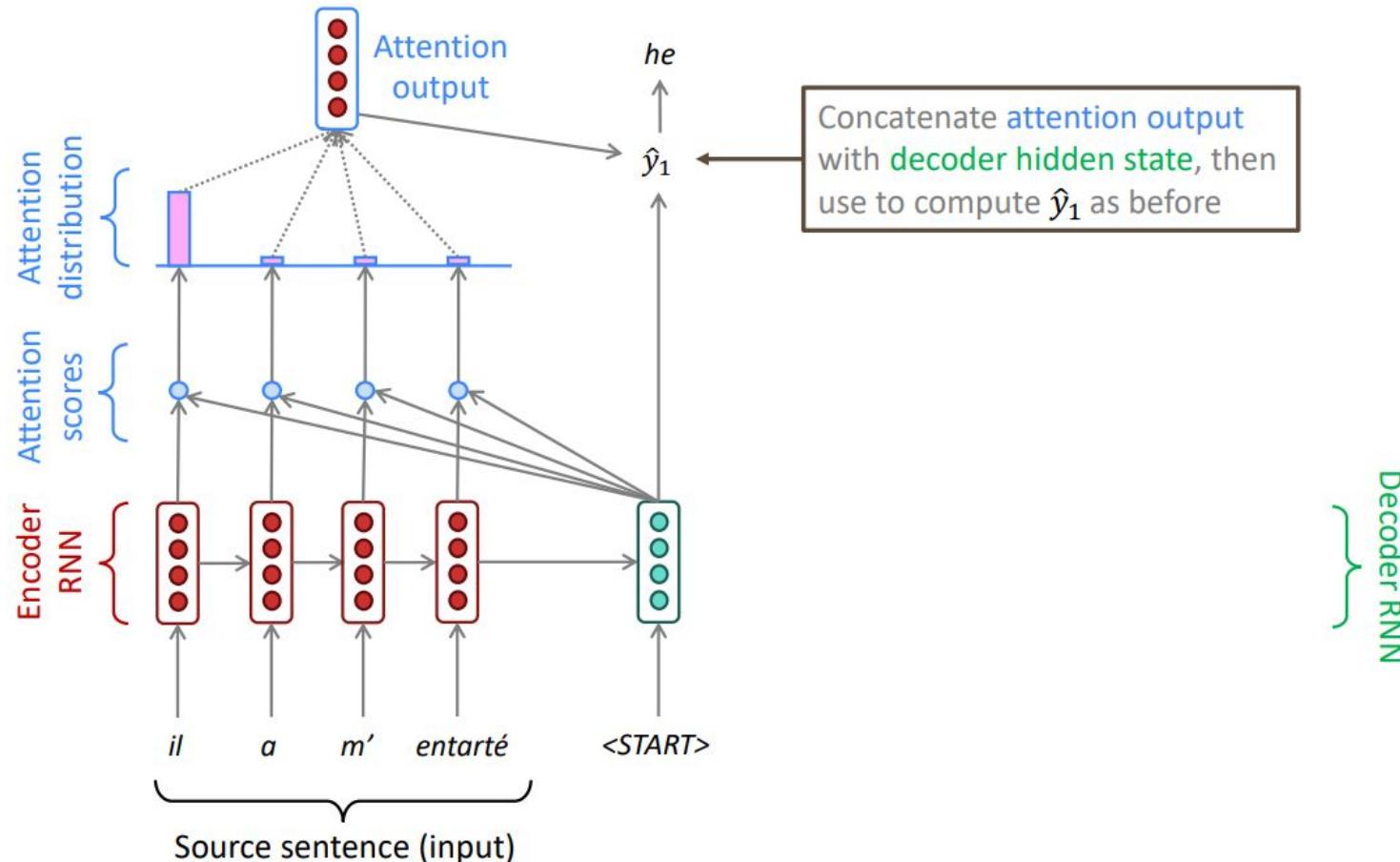
# Sequence-to-sequence with attention



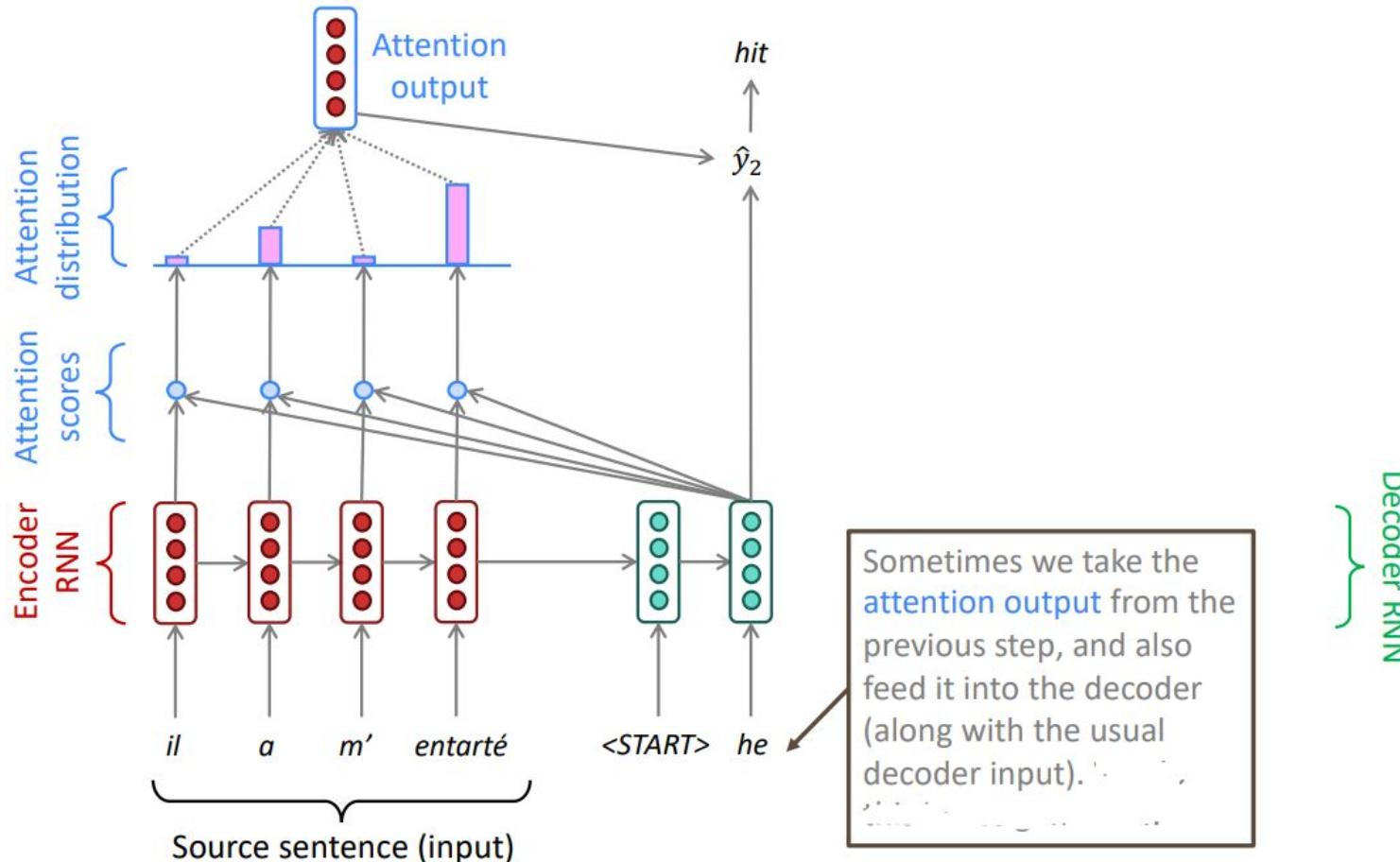
# Sequence-to-sequence with attention



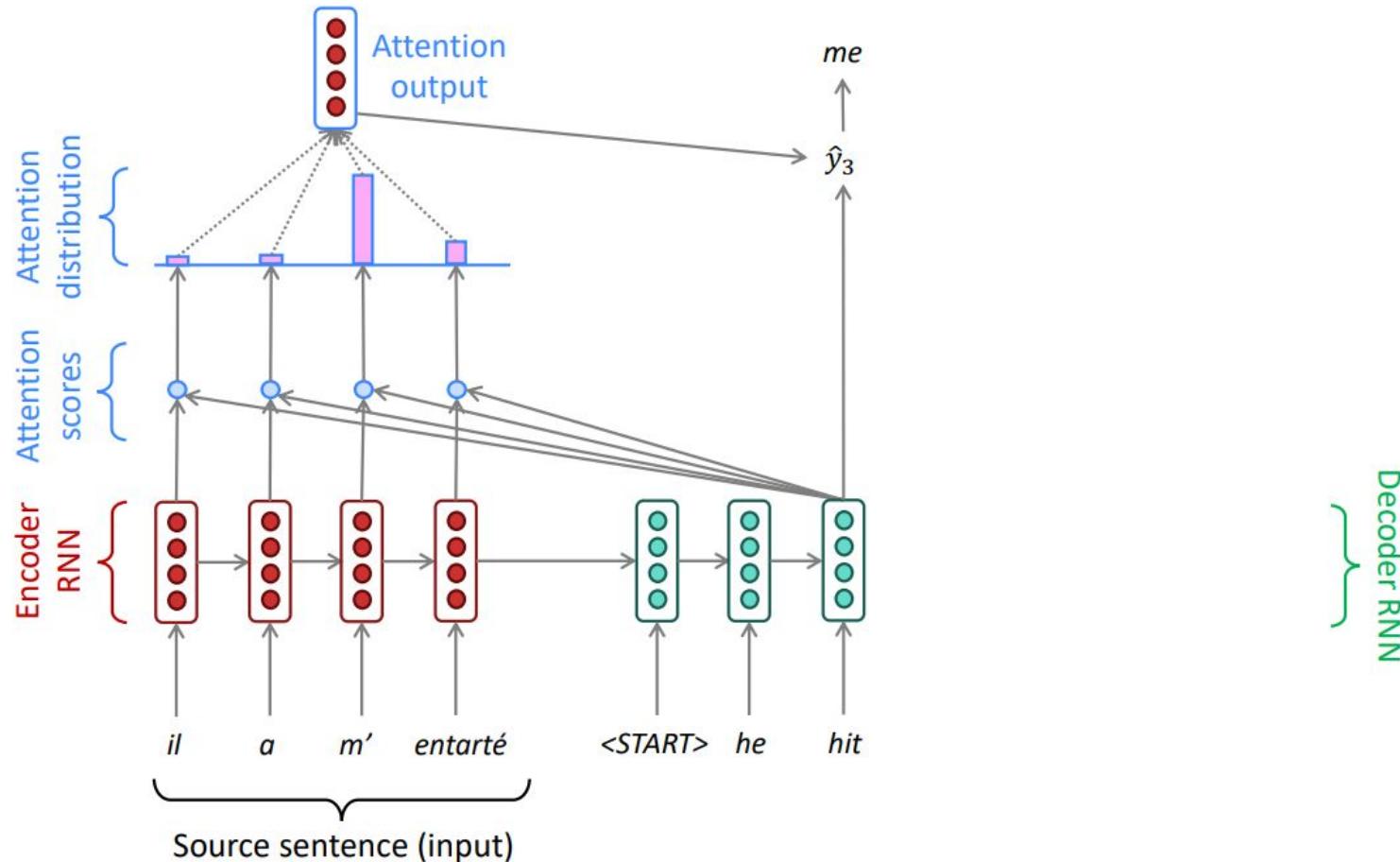
# Sequence-to-sequence with attention



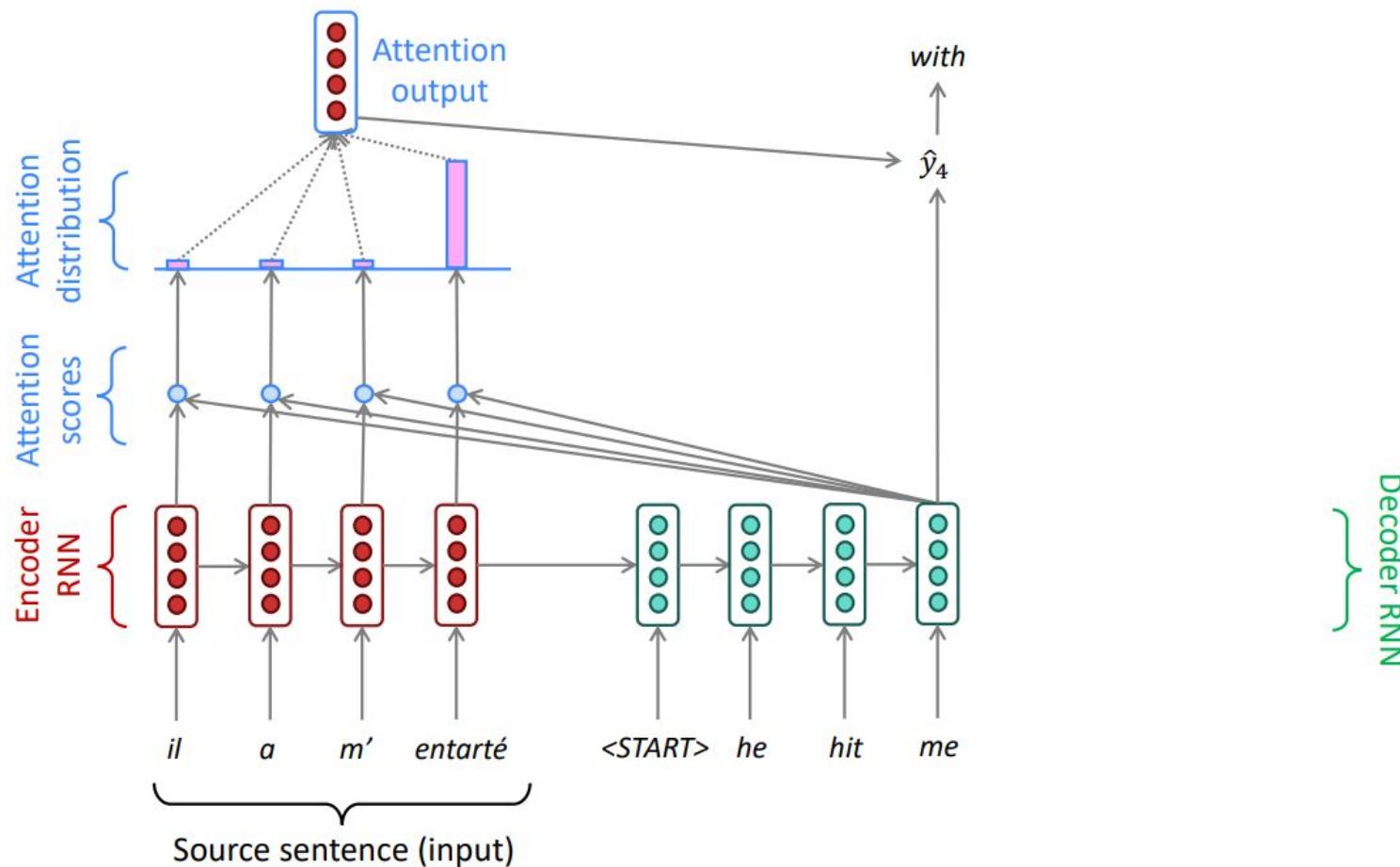
# Sequence-to-sequence with attention



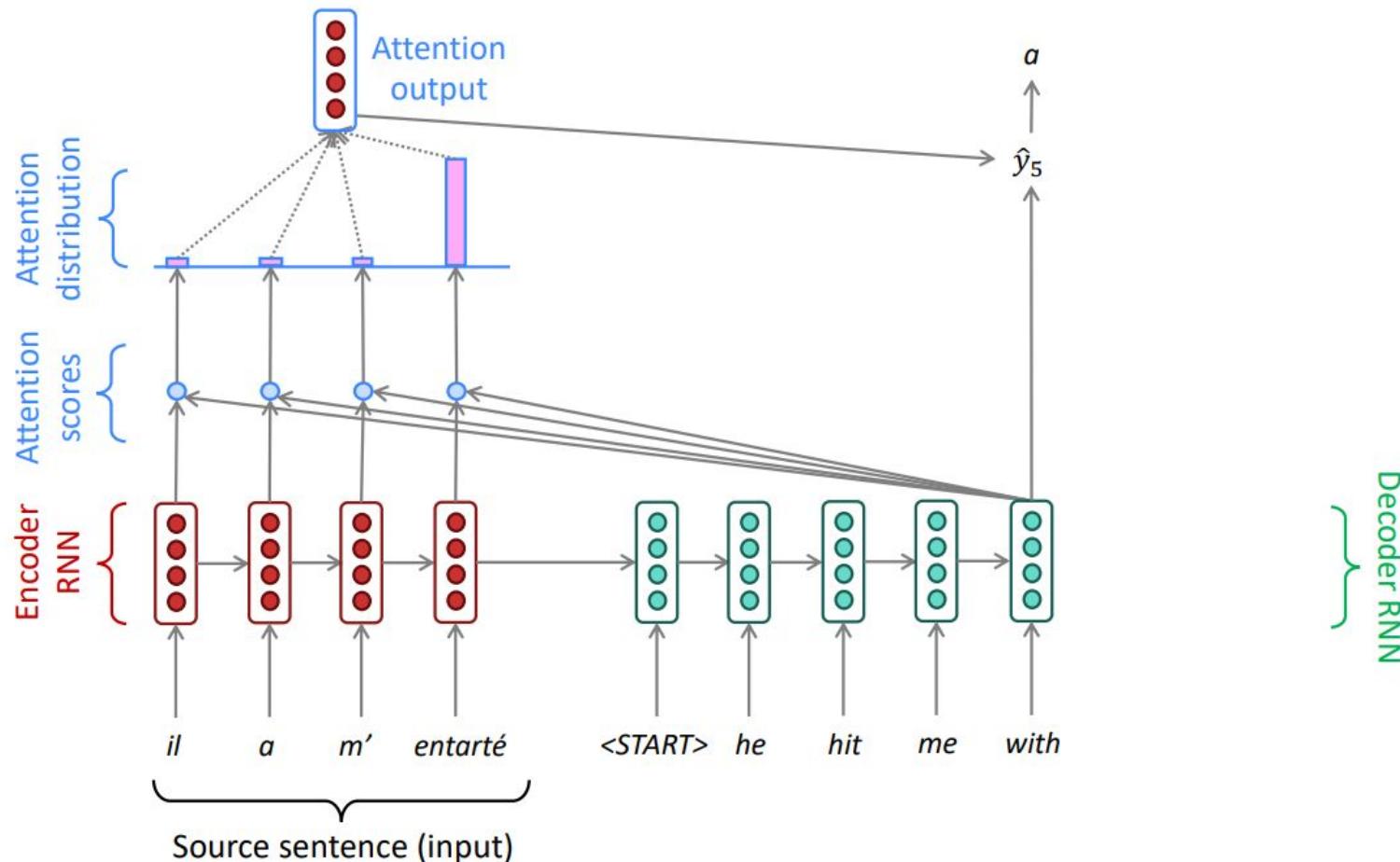
# Sequence-to-sequence with attention



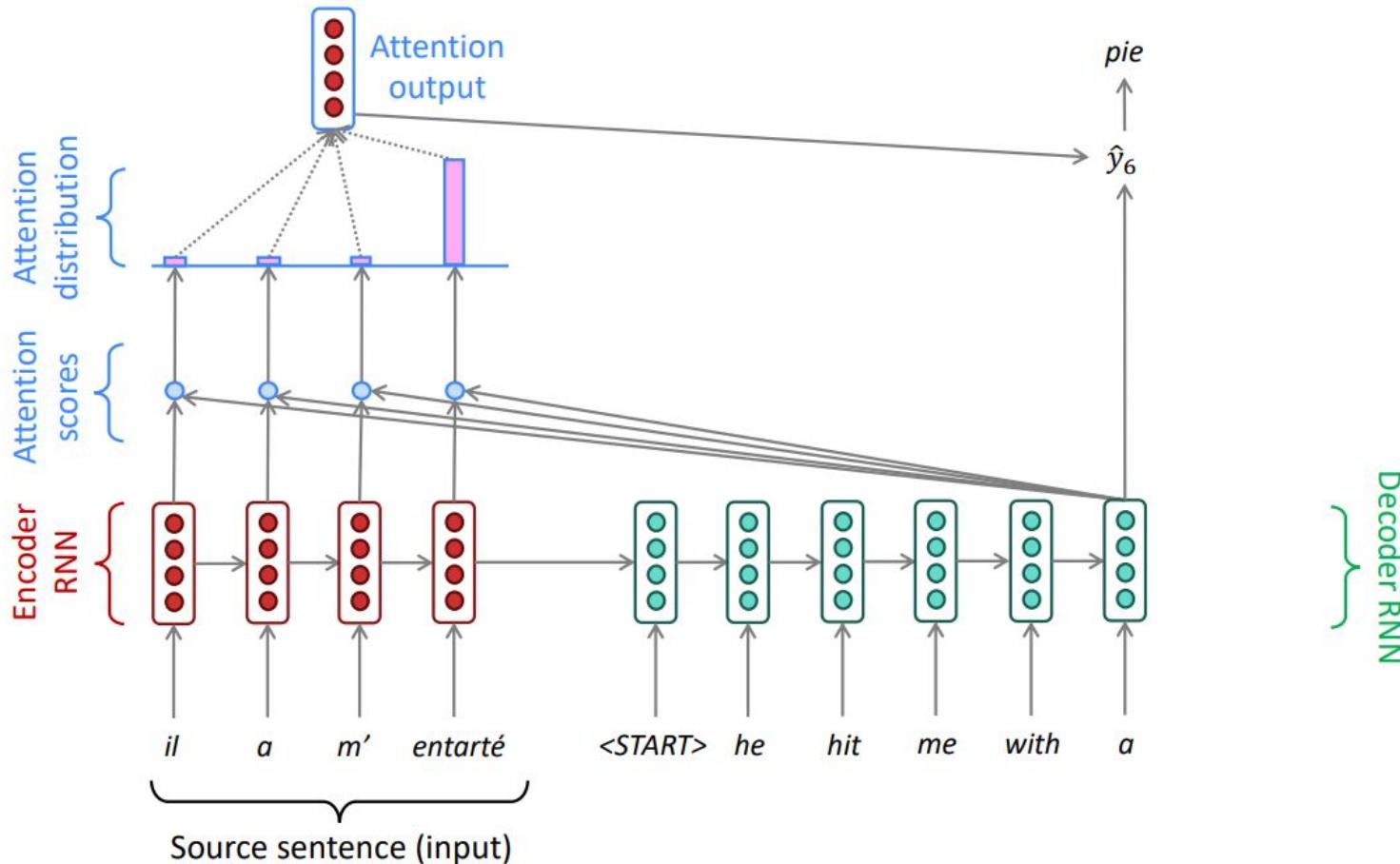
# Sequence-to-sequence with attention



# Sequence-to-sequence with attention



# Sequence-to-sequence with attention



# Attention is great!



- Attention significantly improves NMT performance
  - It's very useful to allow decoder to focus on certain parts of the source
- Attention provides more “human-like” model of the MT process
  - You can look back at the source sentence while translating, rather than needing to remember it all
- Attention solves the bottleneck problem
  - Attention allows decoder to look directly at source; bypass bottleneck
- Attention helps with the vanishing gradient problem
  - Provides shortcut to faraway states
- Attention provides some interpretability
  - By inspecting attention distribution, we see what the decoder was focusing on
  - We get (soft) alignment for free!
  - This is cool because we never explicitly trained an alignment system
  - The network just learned alignment by itself

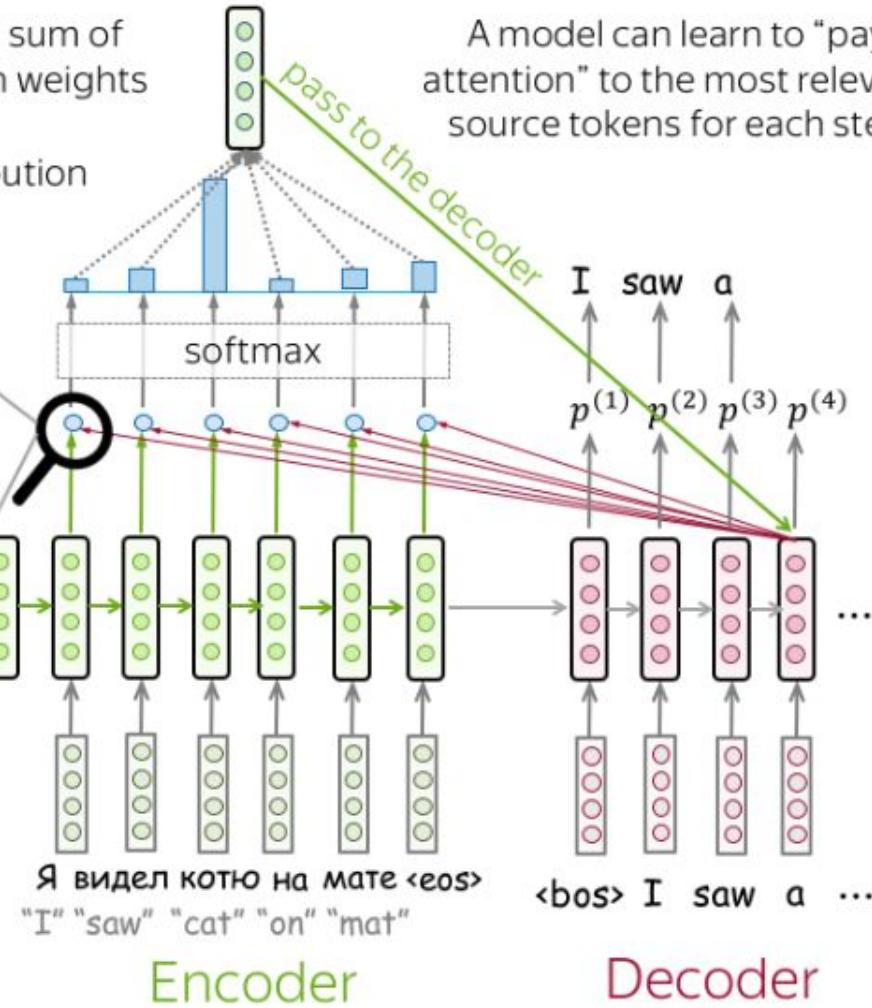
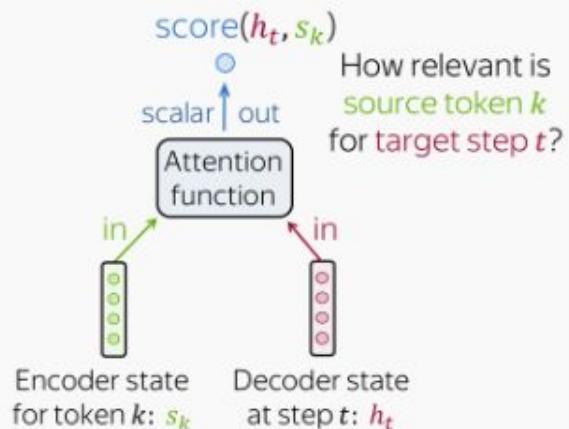
	he	hit	me	with	a	pie
il						
a						
m'						
entarté						

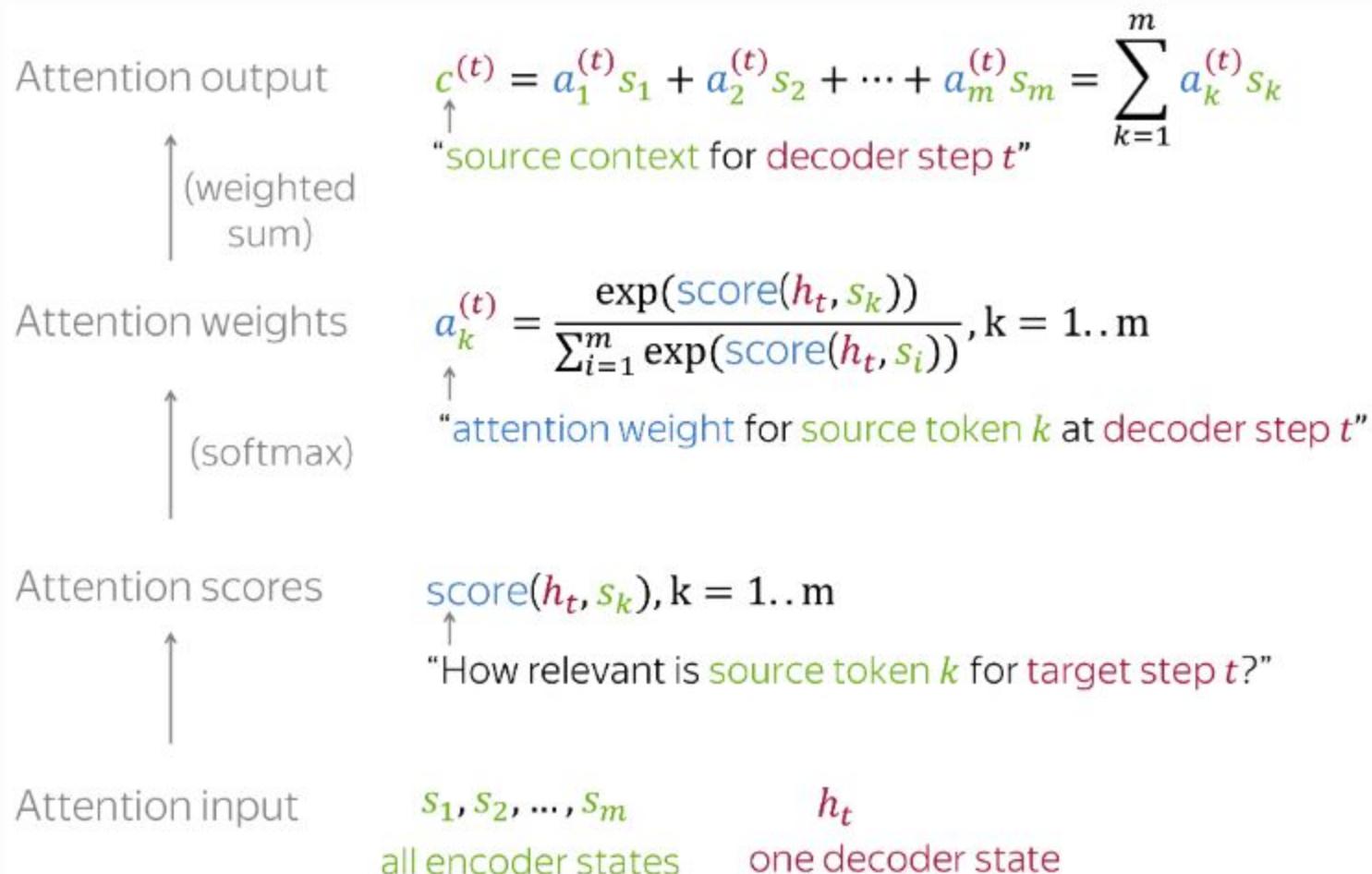
Attention output: weighted sum of encoder states with attention weights

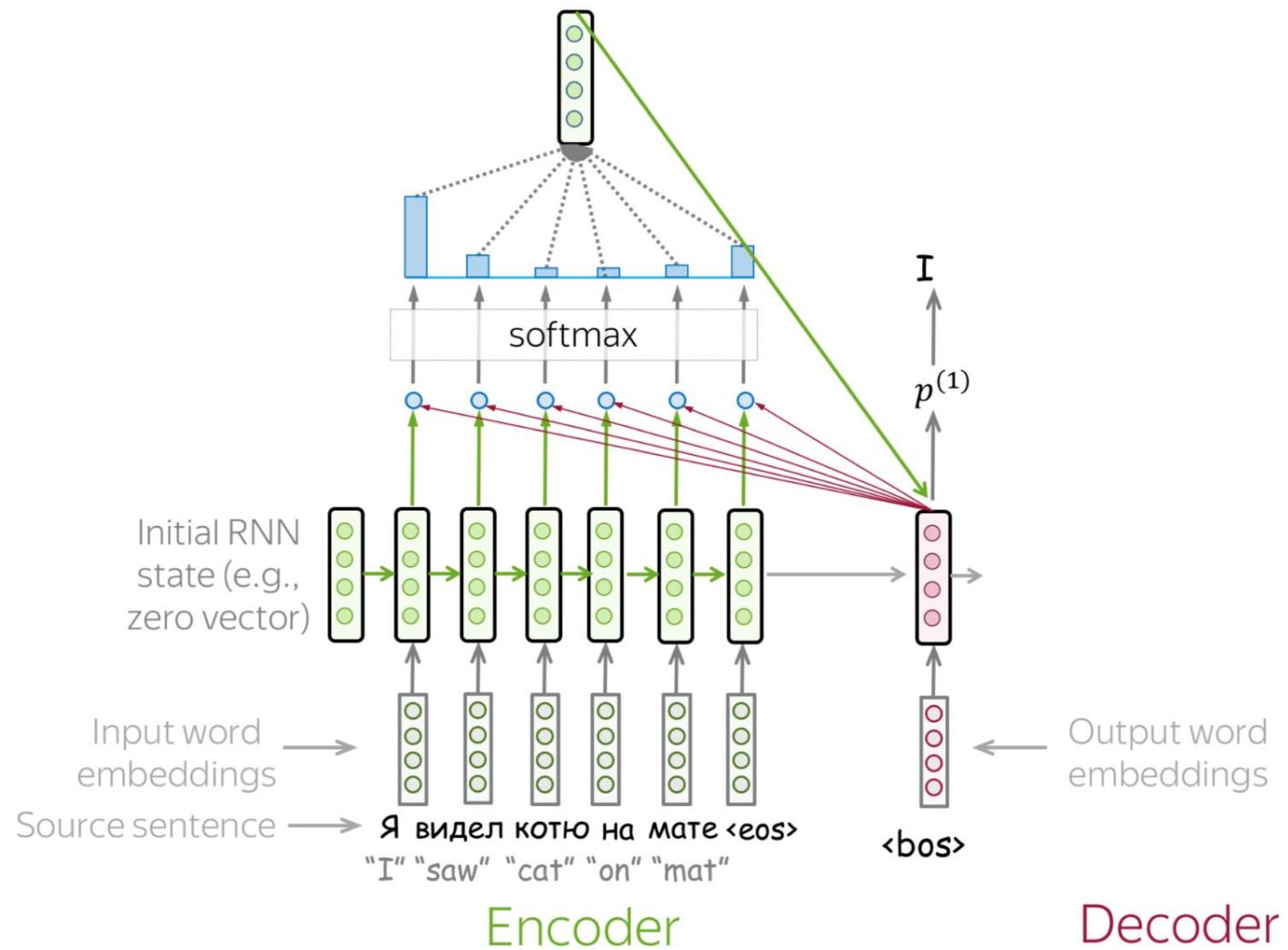
A model can learn to “pay attention” to the most relevant source tokens for each step

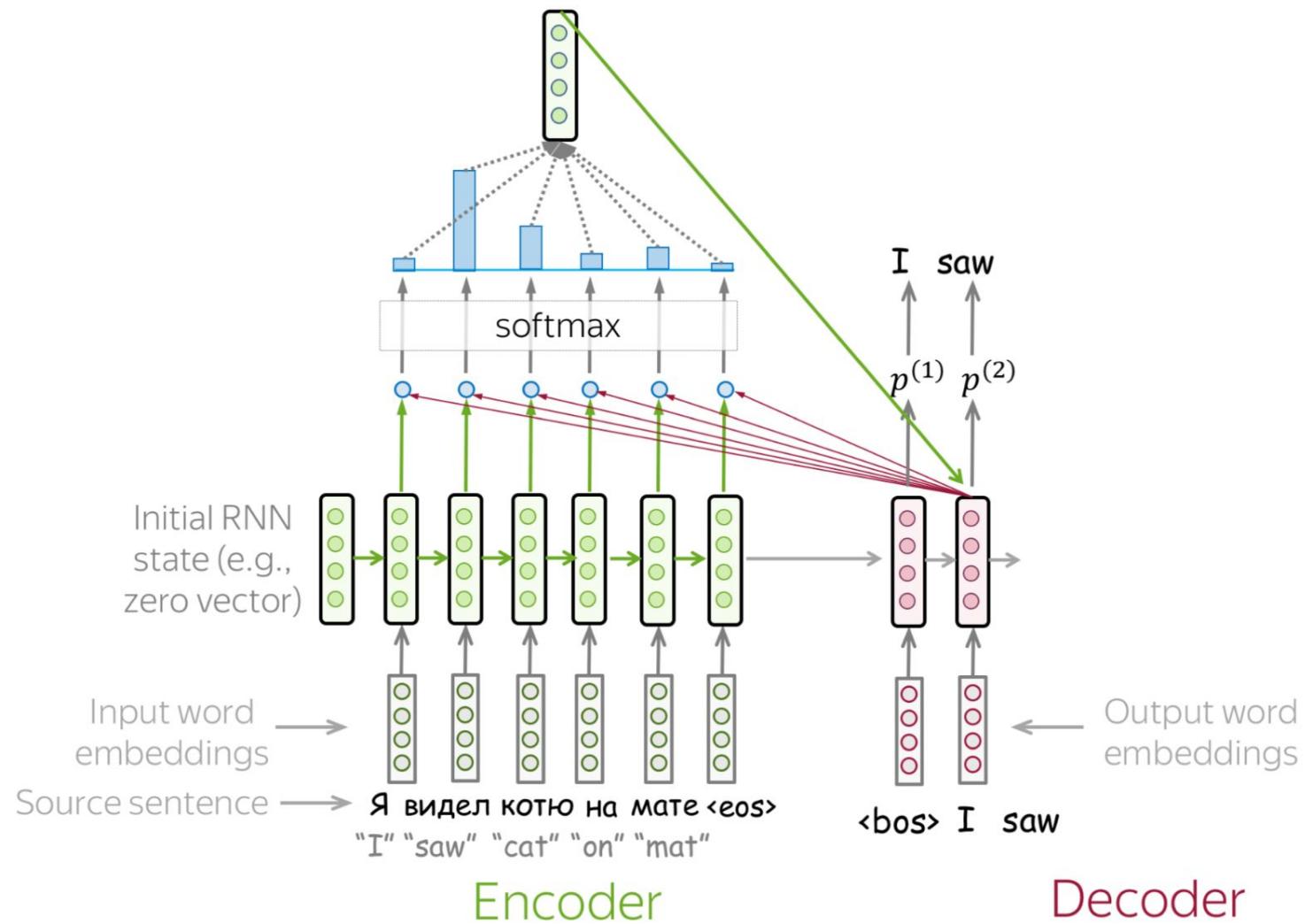
Attention weights: distribution over source tokens

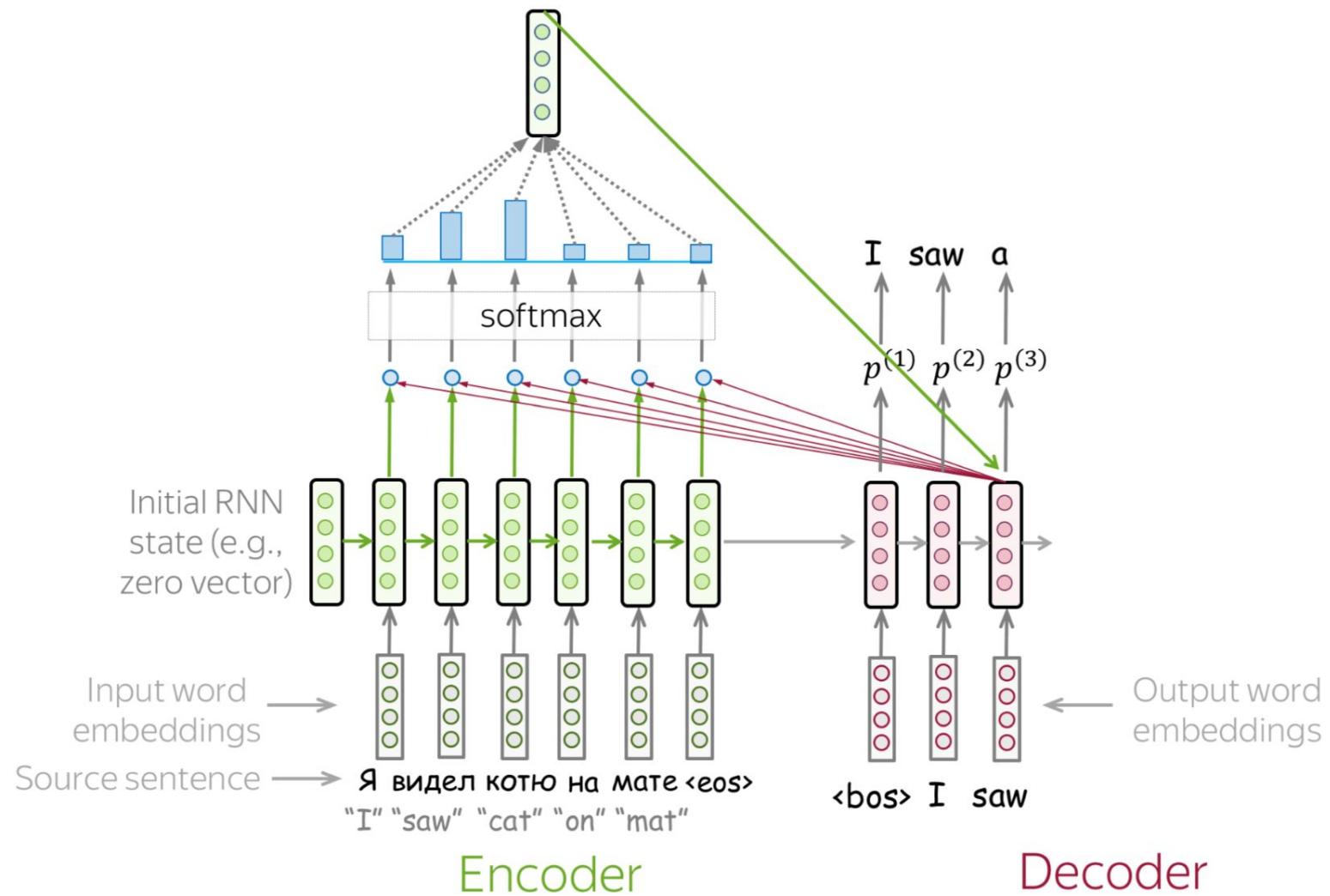
## Attention

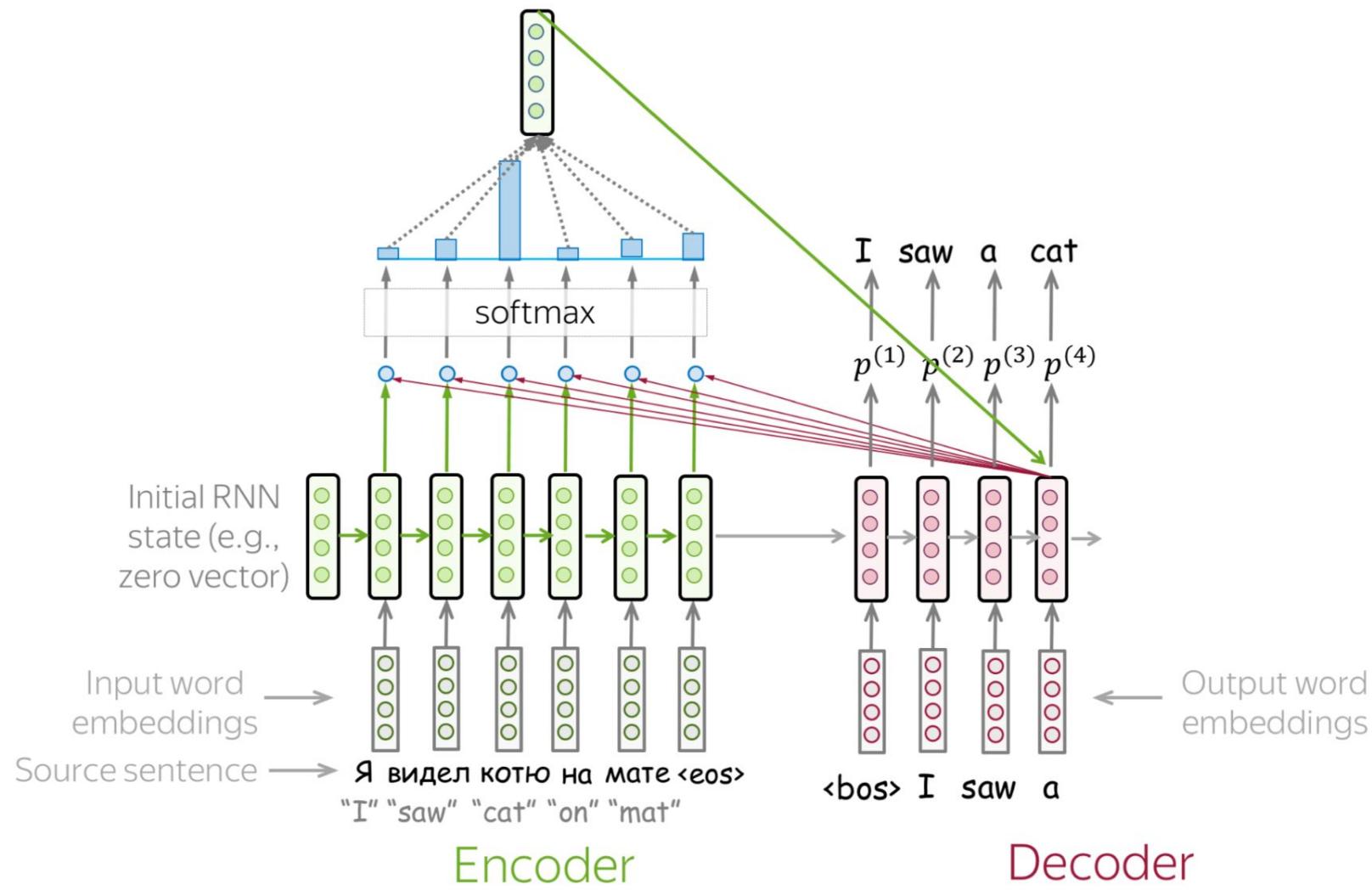


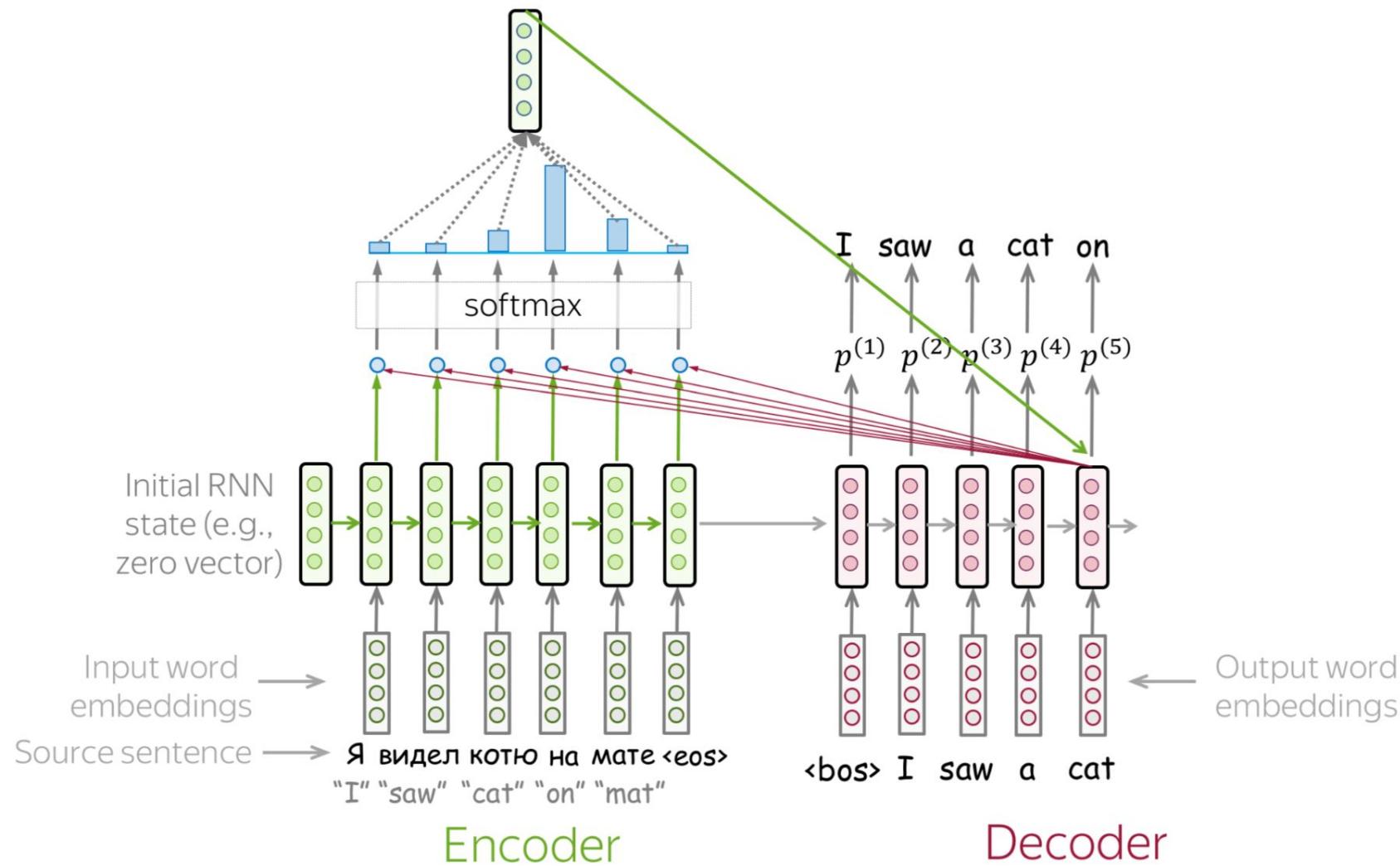


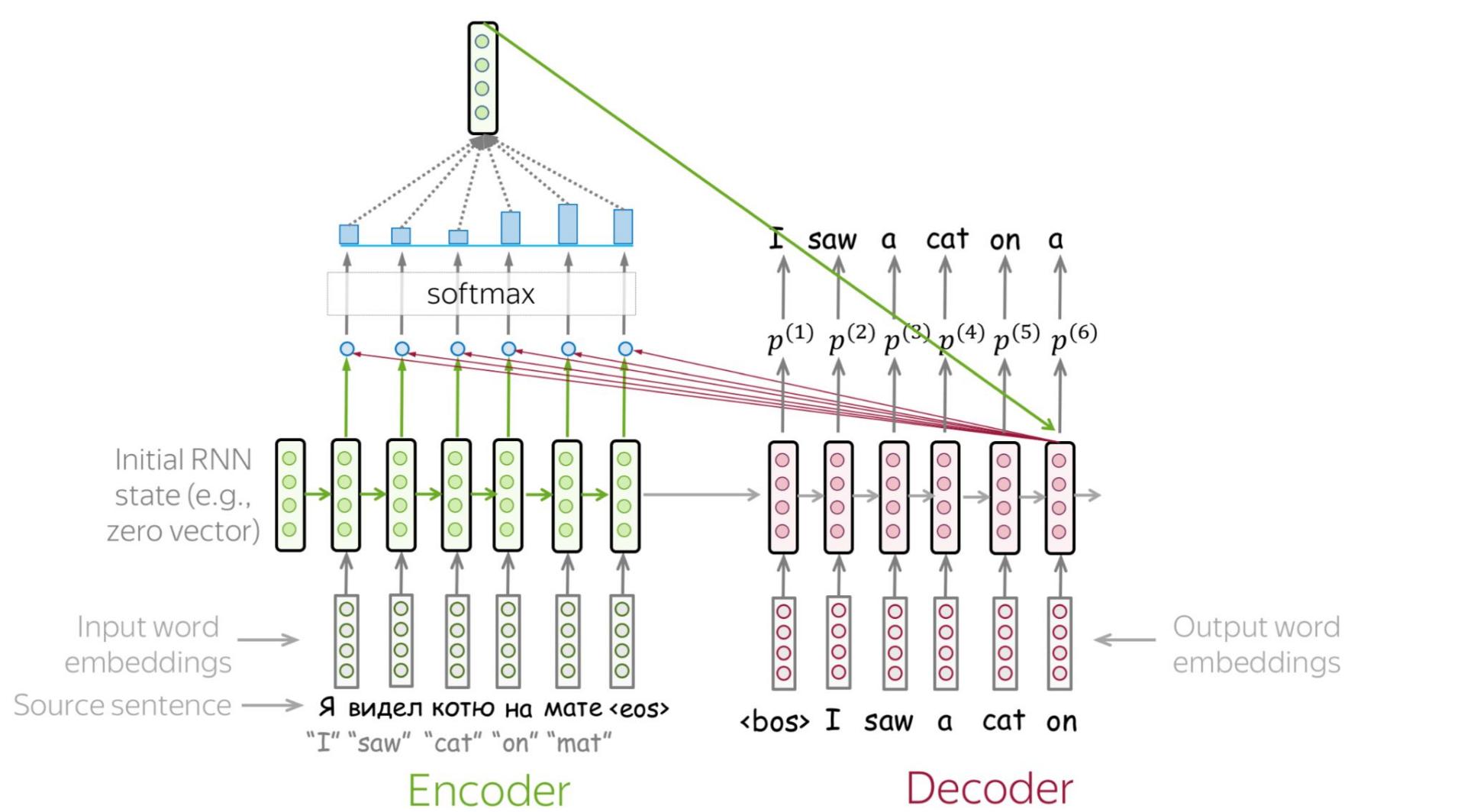


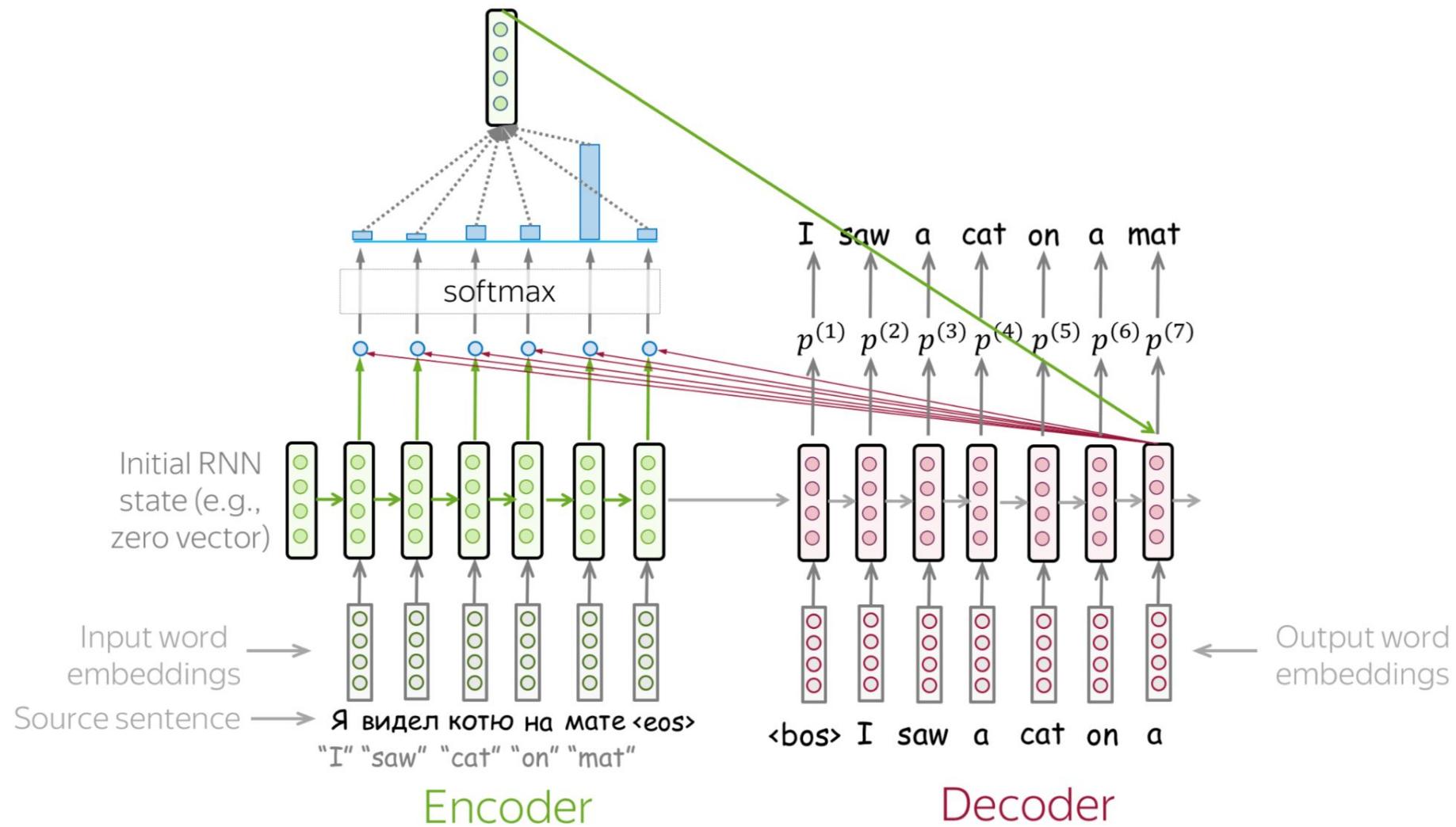


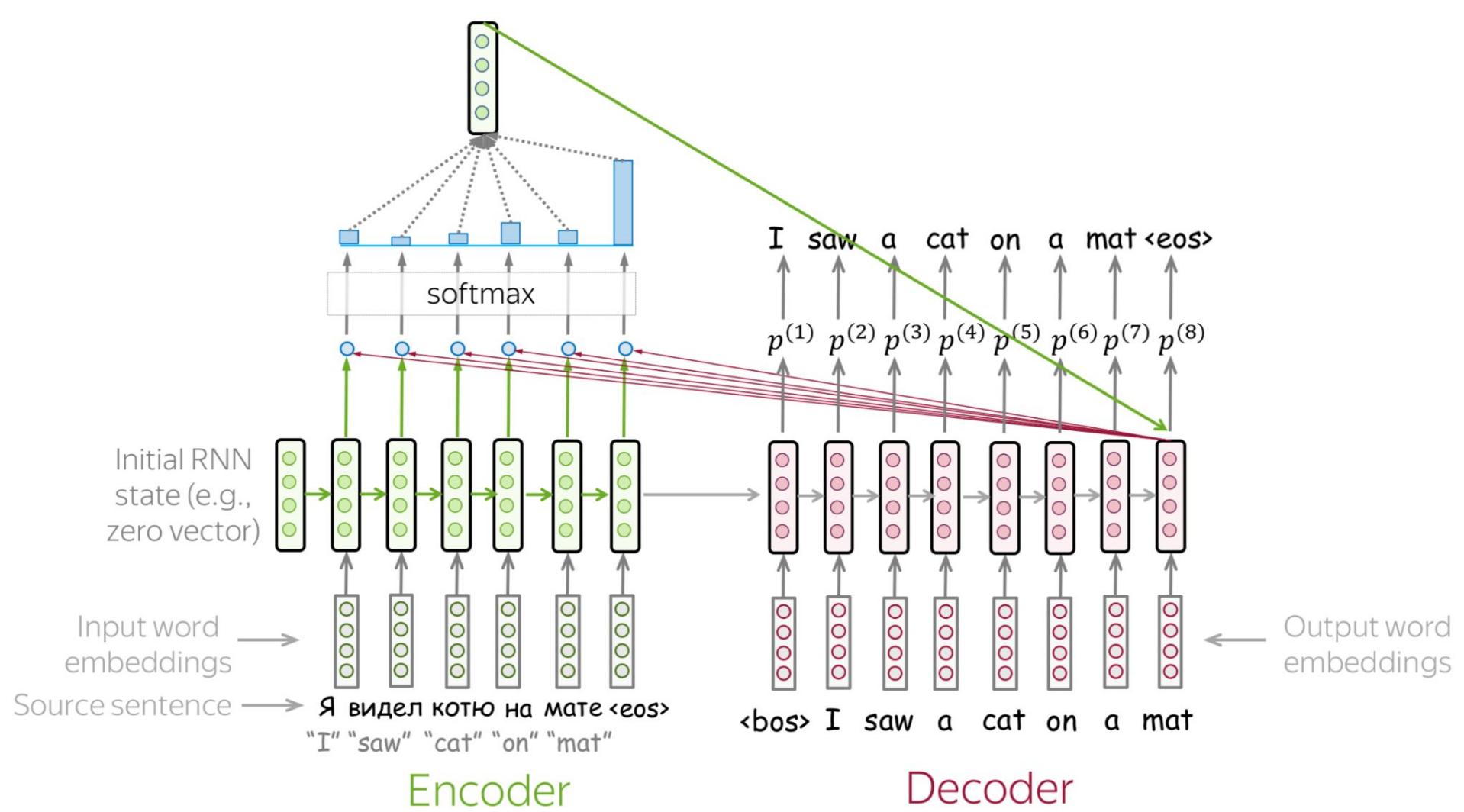




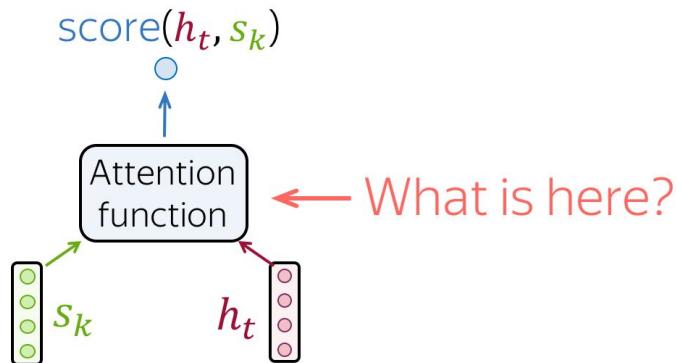








# HOW TO COMPUTE ATTENTION SCORE?



Dot-product

$$h_t^T \times \begin{bmatrix} \textcolor{green}{\circ} \\ \textcolor{green}{\circ} \\ \textcolor{green}{\circ} \\ \textcolor{green}{\circ} \end{bmatrix} s_k$$

Bilinear

$$h_t^T \times \begin{bmatrix} W \end{bmatrix} \times \begin{bmatrix} \textcolor{green}{\circ} \\ \textcolor{green}{\circ} \\ \textcolor{green}{\circ} \\ \textcolor{green}{\circ} \end{bmatrix} s_k$$

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

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

Multi-Layer Perceptron

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

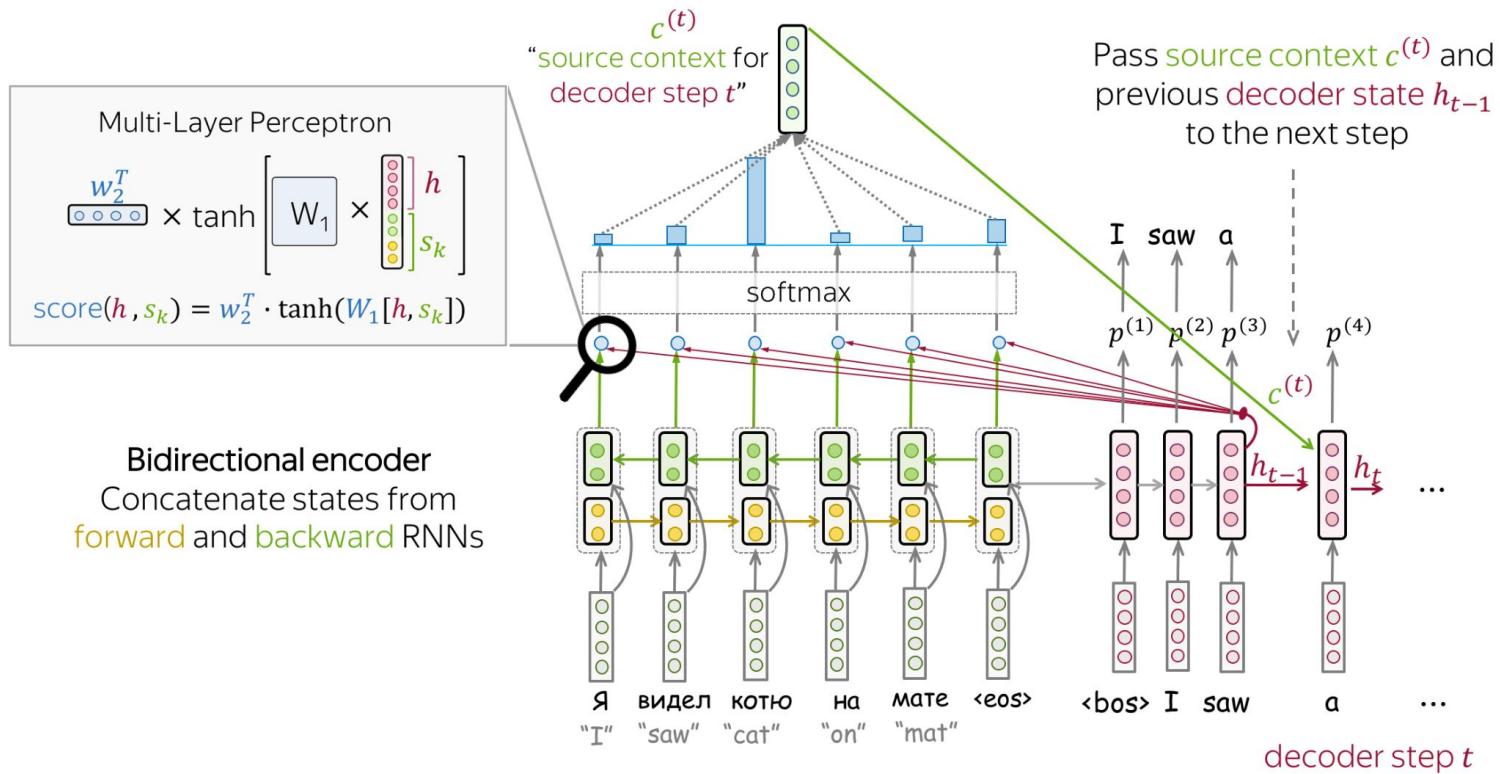
$$\text{score}(h_t, s_k) = w_2^T \cdot \tanh(W_1[h_t, s_k])$$

# HOW TO COMPUTE ATTENTION SCORE?

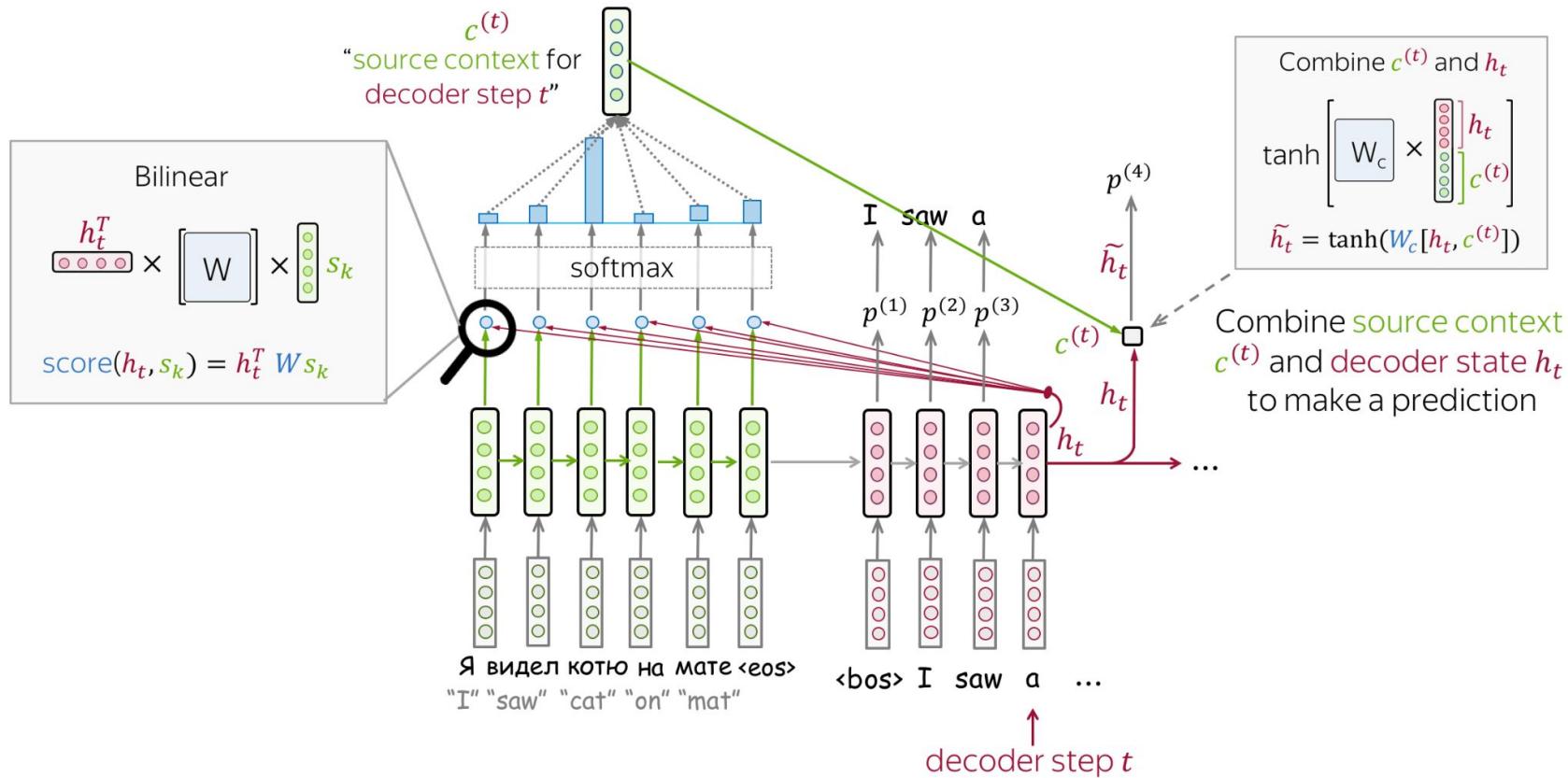
The most popular ways to compute attention scores are:

- dot-product – the simplest method
- bilinear function (aka "Luong attention") – used in the paper **Effective Approaches to Attention-based Neural Machine Translation**
- multi-layer perceptron (aka "Bahdanau attention") – the method proposed in the original paper (**NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE**).

# BAHDANAU MODEL



# LUONG MODEL



## There are *several* attention variants

- We have some *values*  $\mathbf{h}_1, \dots, \mathbf{h}_N \in \mathbb{R}^{d_1}$  and a *query*  $\mathbf{s} \in \mathbb{R}^{d_2}$
- Attention always involves:
  1. Computing the *attention scores*  $\mathbf{e} \in \mathbb{R}^N$
  2. Taking softmax to get *attention distribution*  $\alpha$ :

There are  
multiple ways  
to do this

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

- 3. Using attention distribution to take weighted sum of values:

$$\mathbf{a} = \sum_{i=1}^N \alpha_i \mathbf{h}_i \in \mathbb{R}^{d_1}$$

thus obtaining the *attention output*  $\mathbf{a}$  (sometimes called the *context vector*)

# Attention is a *general* Deep Learning technique

- More general definition of attention:

- Given a set of vector *values*, and a vector *query*, attention is a technique to compute a weighted sum of the values, dependent on the query.

## Intuition:

- The weighted sum is a *selective summary* of the information contained in the values, where the query determines which values to focus on.
- Attention is a way to obtain a *fixed-size representation of an arbitrary set of representations* (the values), dependent on some other representation (the query).

## Upshot:

- Attention has become the powerful, flexible, general way pointer and memory manipulation in all deep learning models. A new idea from after 2010! From NMT!

---

# Attention Is All You Need

---

**Ashish Vaswani\***

Google Brain

[avaswani@google.com](mailto:avaswani@google.com)

**Noam Shazeer\***

Google Brain

[noam@google.com](mailto:noam@google.com)

**Niki Parmar\***

Google Research

[nikip@google.com](mailto:nikip@google.com)

**Jakob Uszkoreit\***

Google Research

[usz@google.com](mailto:usz@google.com)

**Llion Jones\***

Google Research

[llion@google.com](mailto:llion@google.com)

**Aidan N. Gomez\*** †

University of Toronto

[aidan@cs.toronto.edu](mailto:aidan@cs.toronto.edu)

**Lukasz Kaiser\***

Google Brain

[lukaszkaiser@google.com](mailto:lukaszkaiser@google.com)

**Illia Polosukhin\*** ‡

[illia.polosukhin@gmail.com](mailto:illia.polosukhin@gmail.com)

**ATTENTION IS ALL YOU NEED?**

**FALSE. YOU NEED WATER AND RATIONS.**

# REFERENCES

<https://web.stanford.edu/class/cs224n/slides/cs224n-2022-lecture08-final-project.pdf>

[https://lena-voita.github.io/nlp\\_course/seq2seq\\_and\\_attention.html](https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html)