

# Encoder Decoder model (Sequence to Sequence Model)

## Machine Translation

- Given a sentence in a source language, translate into a target language
- These two sequences may have different lengths

English ▼



Hindi ▼

UTD is a good university



यूटीडी एक अच्छा विश्वविद्यालय है

yooteedee ek achchha vishvavidyaalay  
hai

# Machine Translation

- Given a sentence in a source language, translate into a target language
- These two sequences may have different lengths

We cannot translate sentences word by word

Can you me help this sentence to translate  
↑ ↑ ↑ ↑ ↑ ↑ ↑ ↑  
Kannst du mir helfen diesen Satz zu uebersetzen ?

Can you help me to translate this sentence  
↑ ↑ × × × × × ×  
Kannst du mir helfen diesen Satz zu uebersetzen ?

# Neural Machine Translation



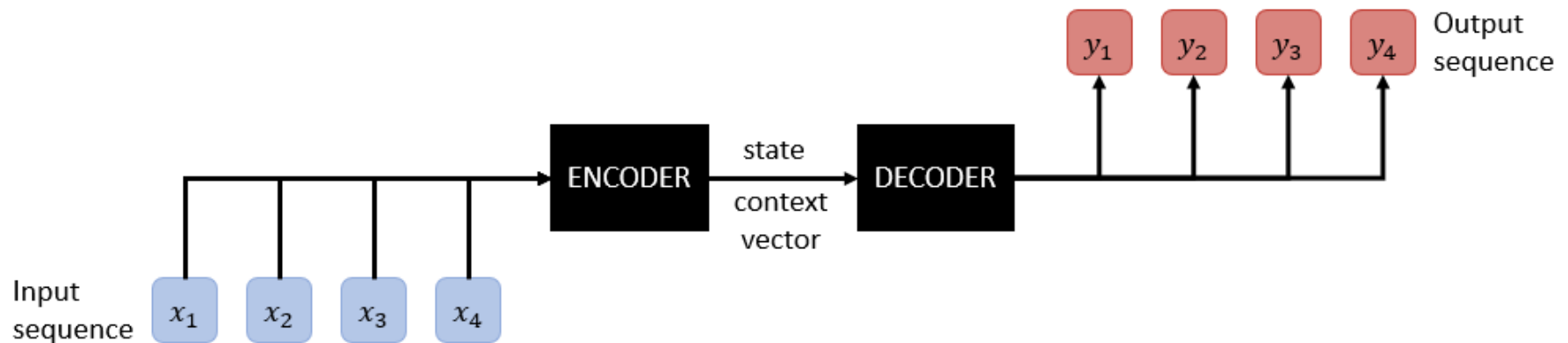
Neural machine Translation (NMT) – Machine translation with a single neural network



The neural network architecture typically involves two RNNs – this is called sequence-to-sequence

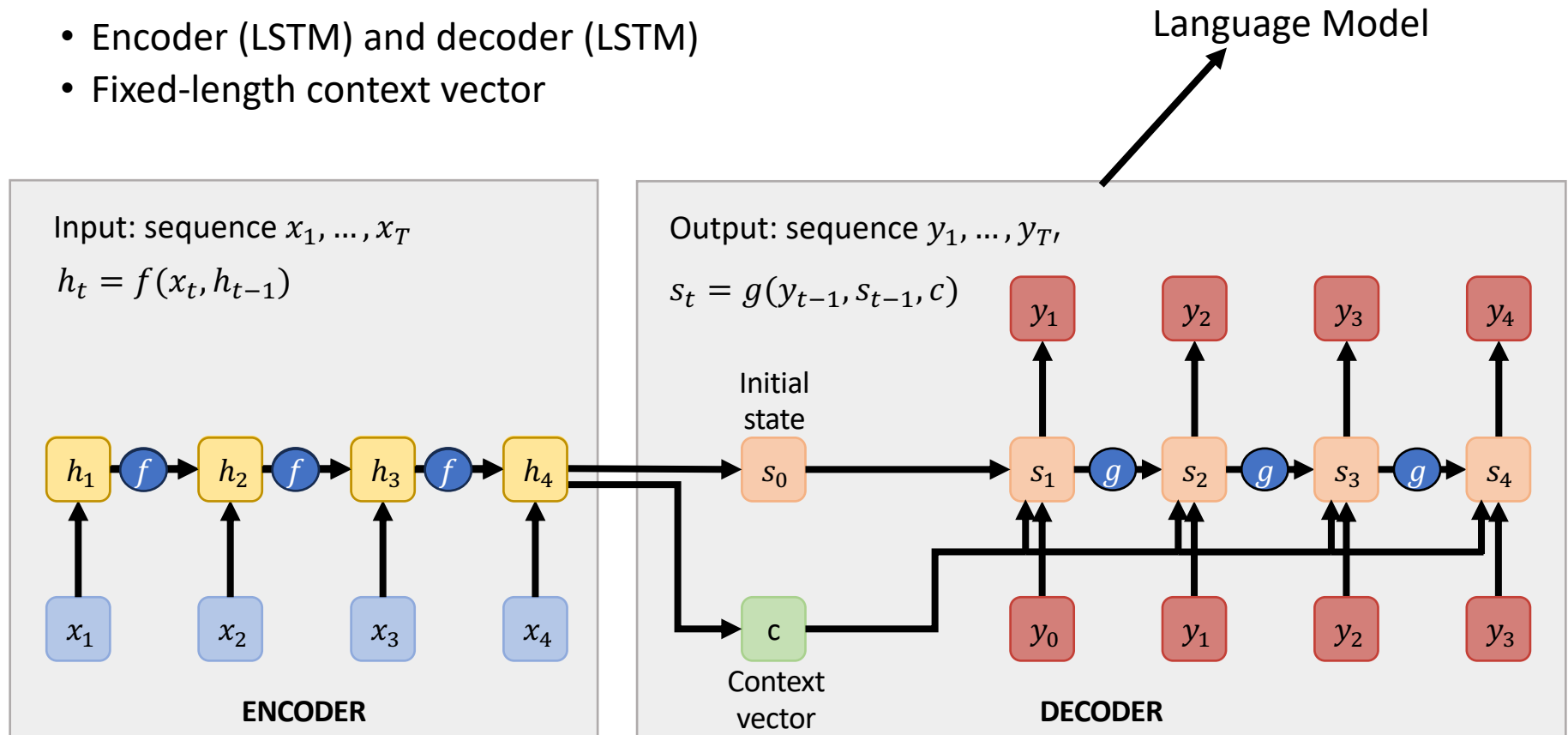
# The Encoder-decoder Architecture

- A model is partitioned into two parts
  - The encoder process inputs
  - The decoder generates outputs



# Sequence to Sequence with RNNs

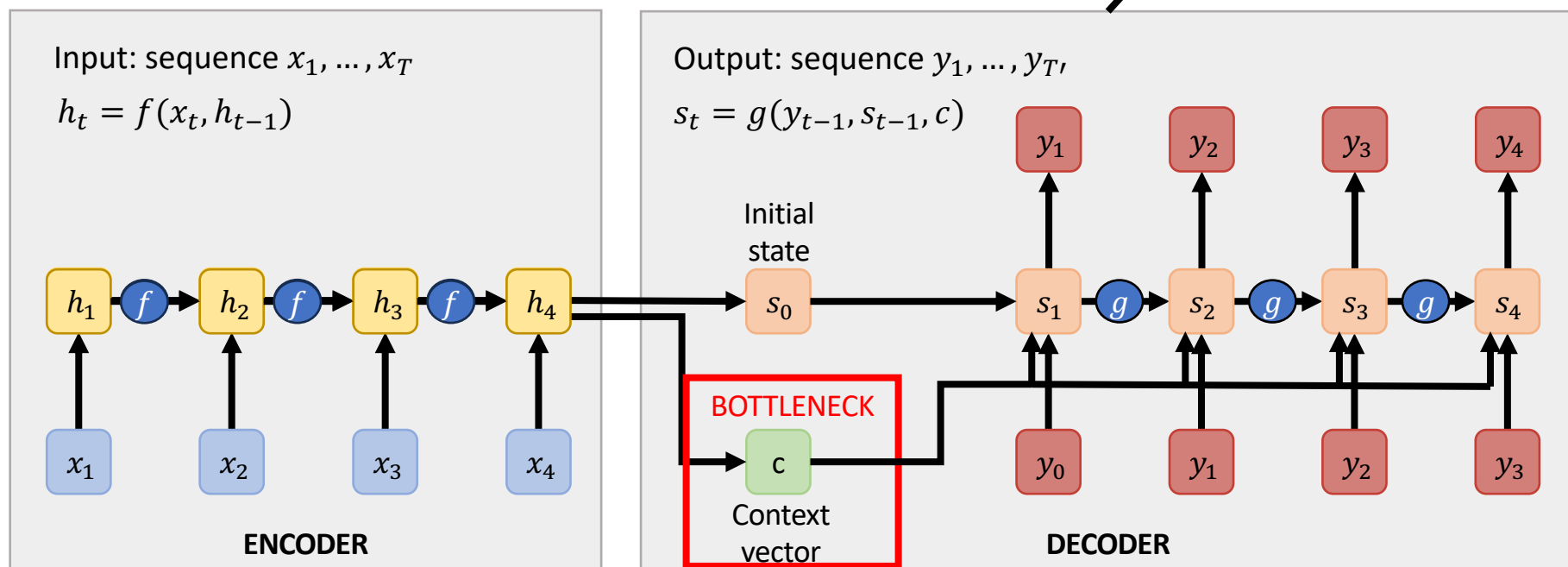
- Encoder (LSTM) and decoder (LSTM)
- Fixed-length context vector



I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Proceedings of the 27th International Conference on Neural Information Processing Systems (NIPS)*, 2014, pp. 3104–3112.

# Sequence to Sequence with RNNs

- Encoder (LSTM) and decoder (LSTM)
- Fixed-length context vector (**bottleneck**)

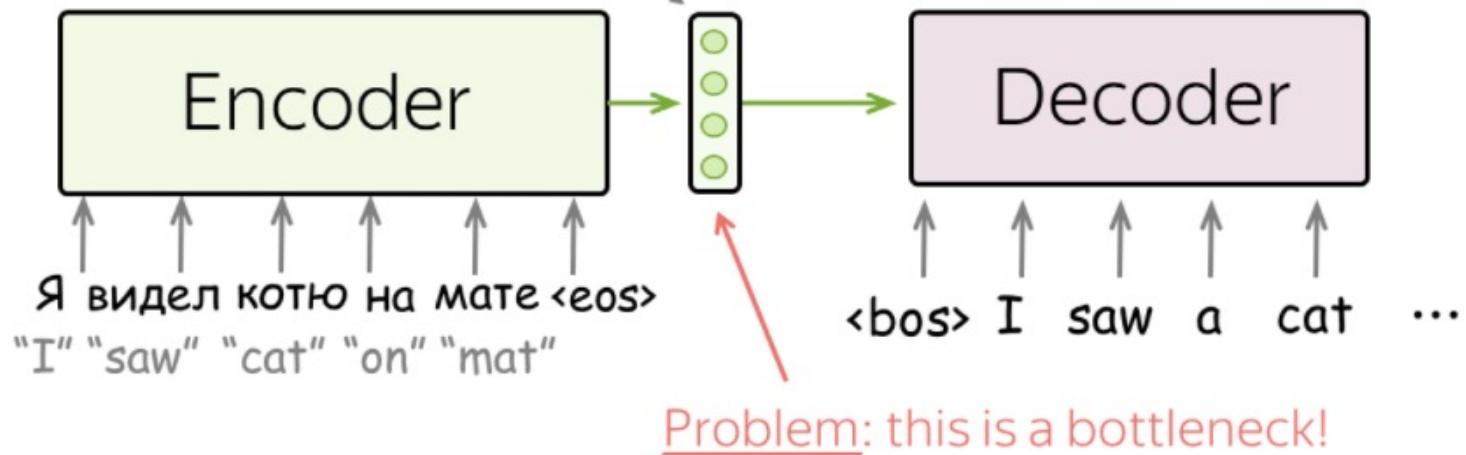


I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Proceedings of the 27th International Conference on Neural Information Processing Systems (NIPS)*, 2014, pp. 3104–3112.

## sequence-to-sequence – bottleneck problem

- Fixed source representation is suboptimal for the encoder, it is hard to compress the sentence
- For the decoder, different information may be relevant at different steps

We saw: encoder compresses the source into a single vector





+

•

○



# ATTENTION



# 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



# Sequence to Sequence with RNNs + Attention

**Idea!** Use a different context vector for each timestep in the decoder

$$s_t = g(y_{t-1}, s_{t-1}, c_t)$$

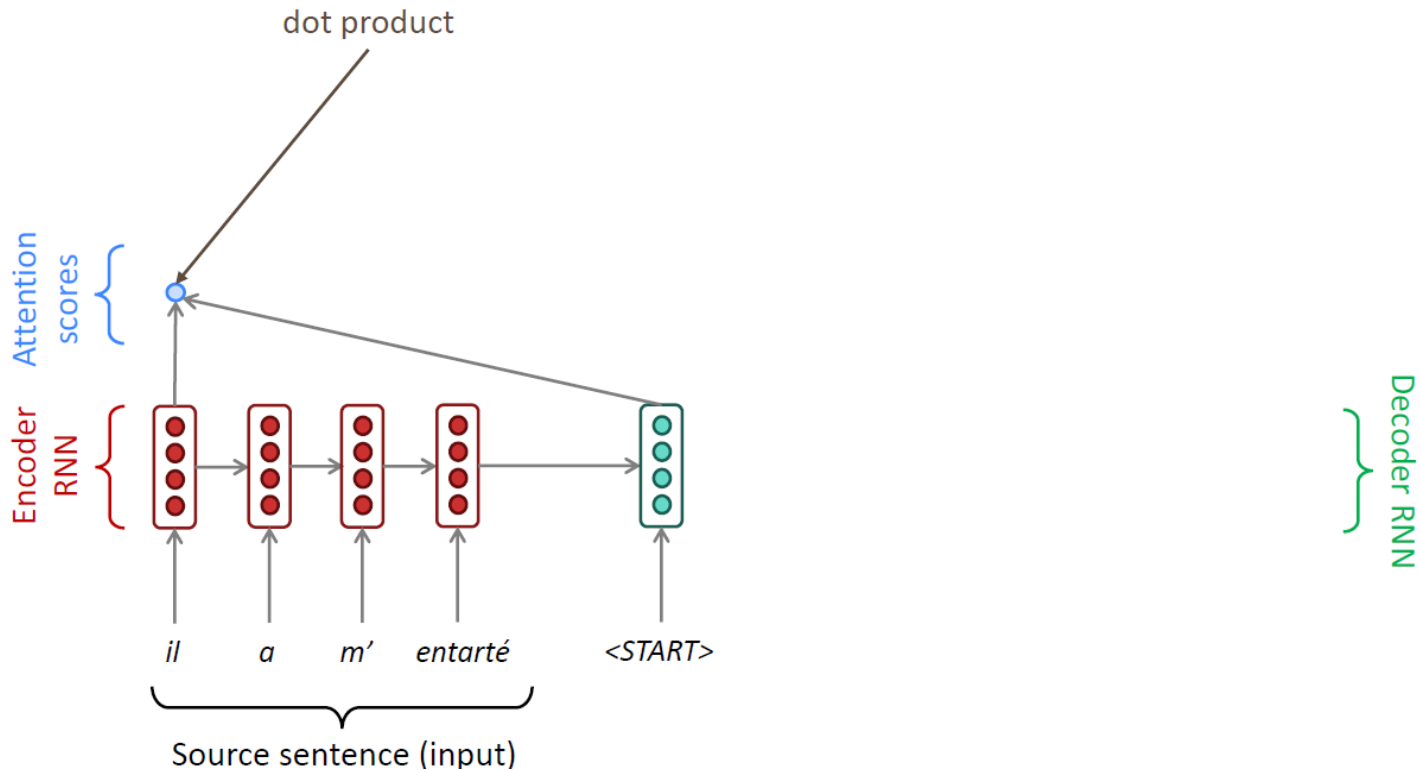
No more bottleneck through a single vector

Craft the context vector so that it “looks at” different parts of the input sequence for each decoder timestep

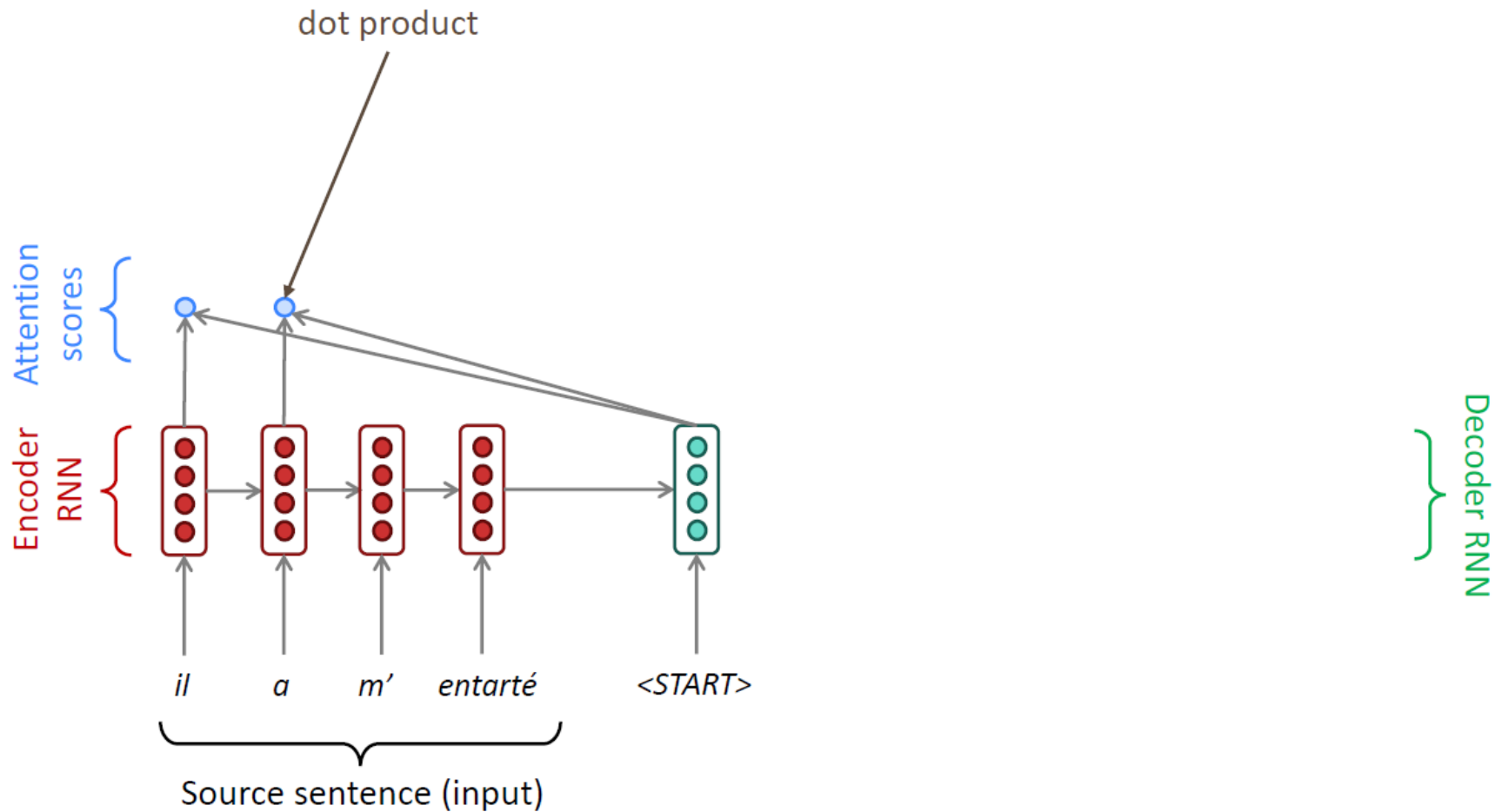
D. Bahdanau, K. Cho, and Y. Bengio, “Neural Machine Translation by Jointly Learning to Align and Translate,” in *3rd International Conference on Learning Representations (ICLR)*, 2015.

# Sequence to Sequence with RNNs + Attention

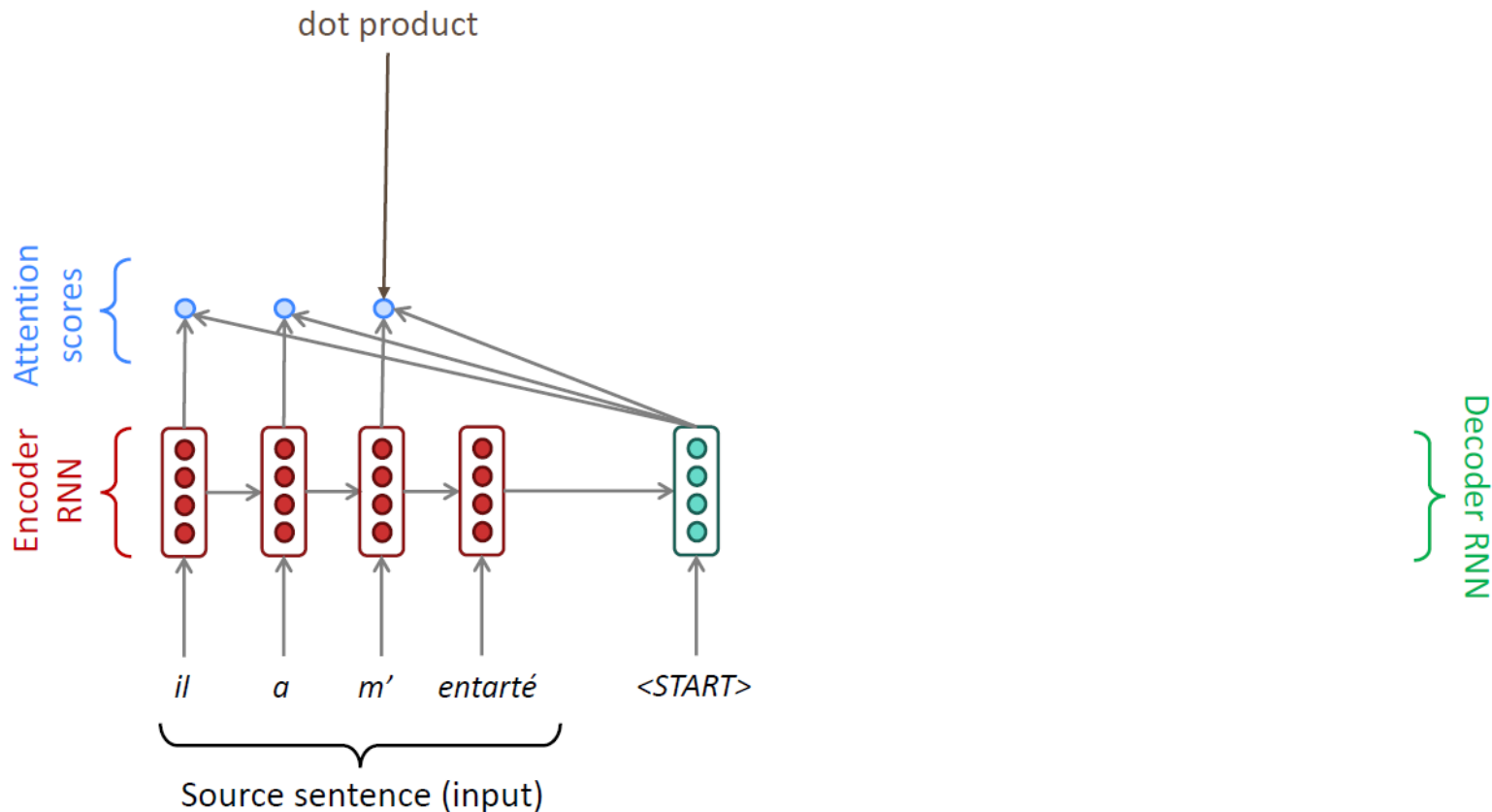
**Core idea:** on each step of the decoder, *use direct connection to the encoder* to *focus on a particular part* of the source sequence



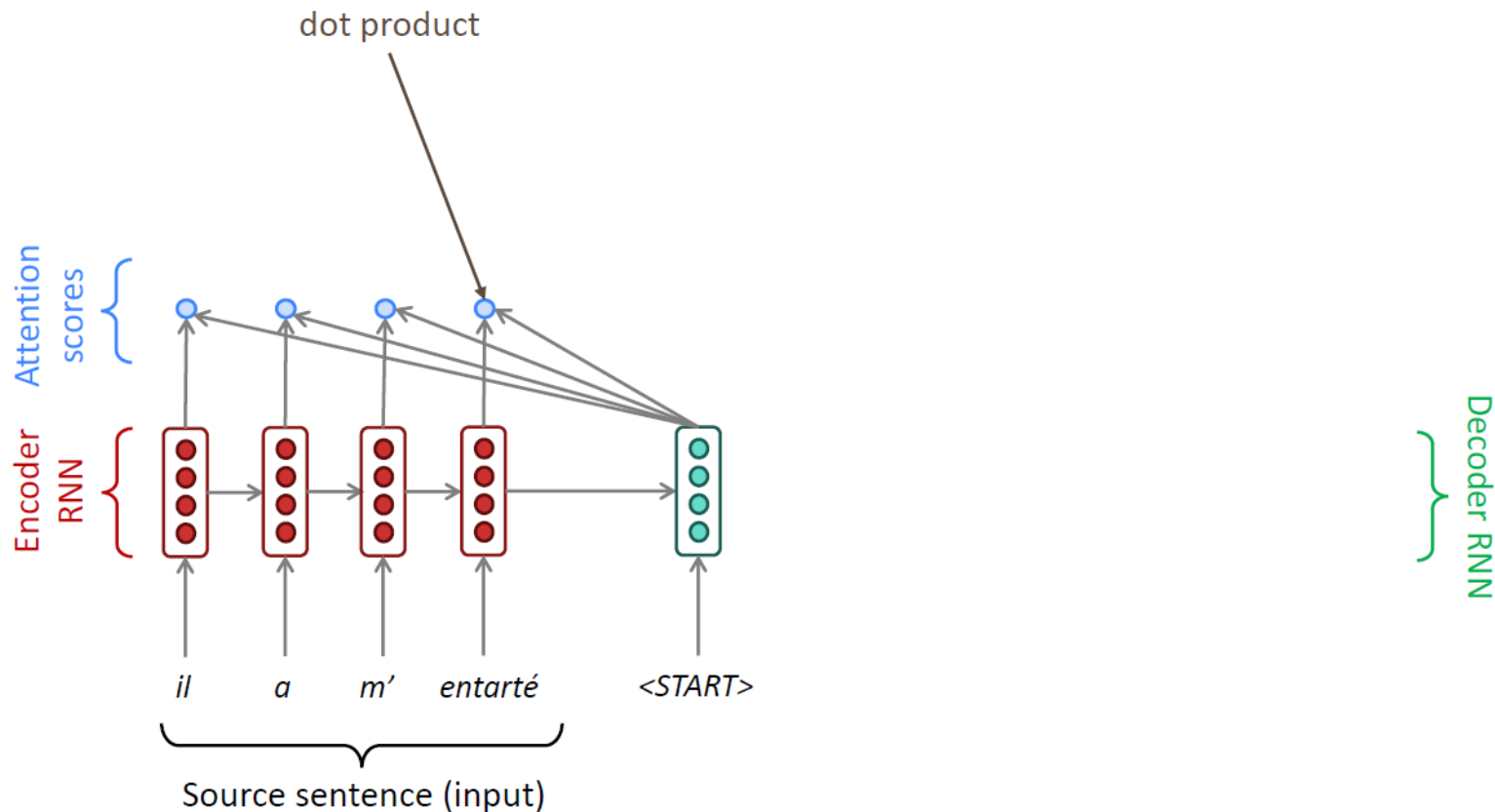
# Sequence to Sequence with RNNs + Attention



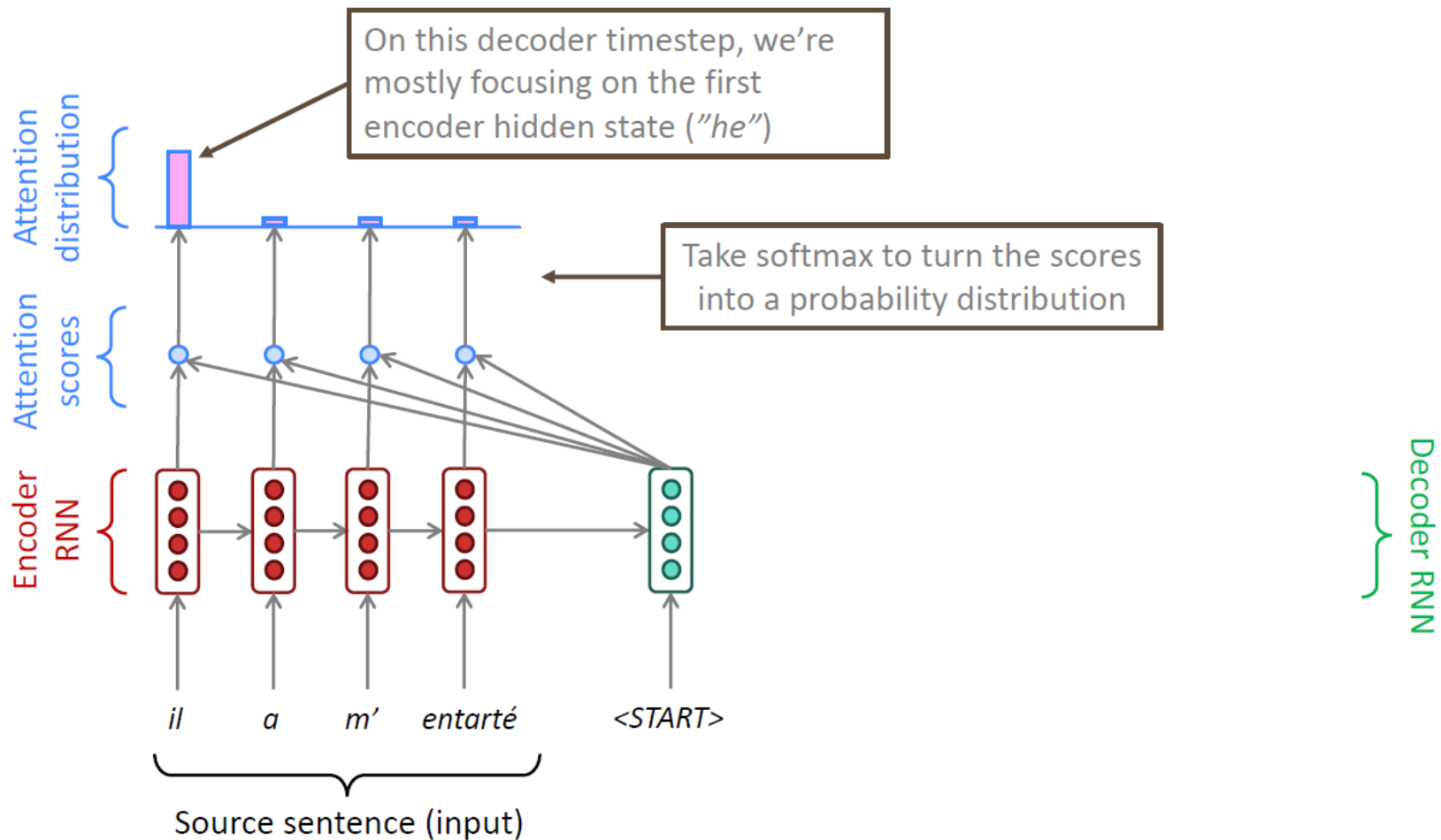
# Sequence to Sequence with RNNs + Attention



# Sequence to Sequence with RNNs + Attention

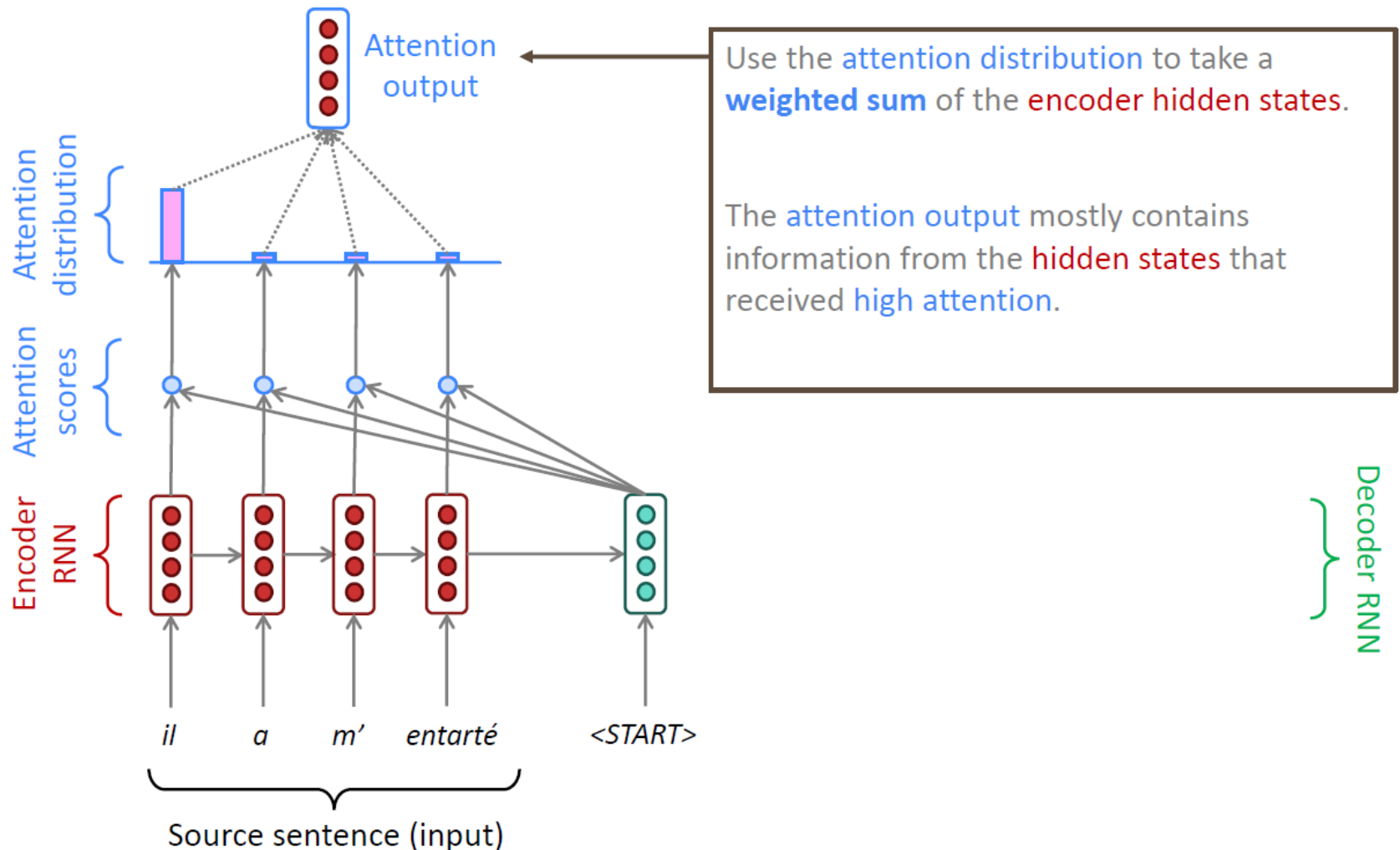


# Sequence to Sequence with RNNs + Attention

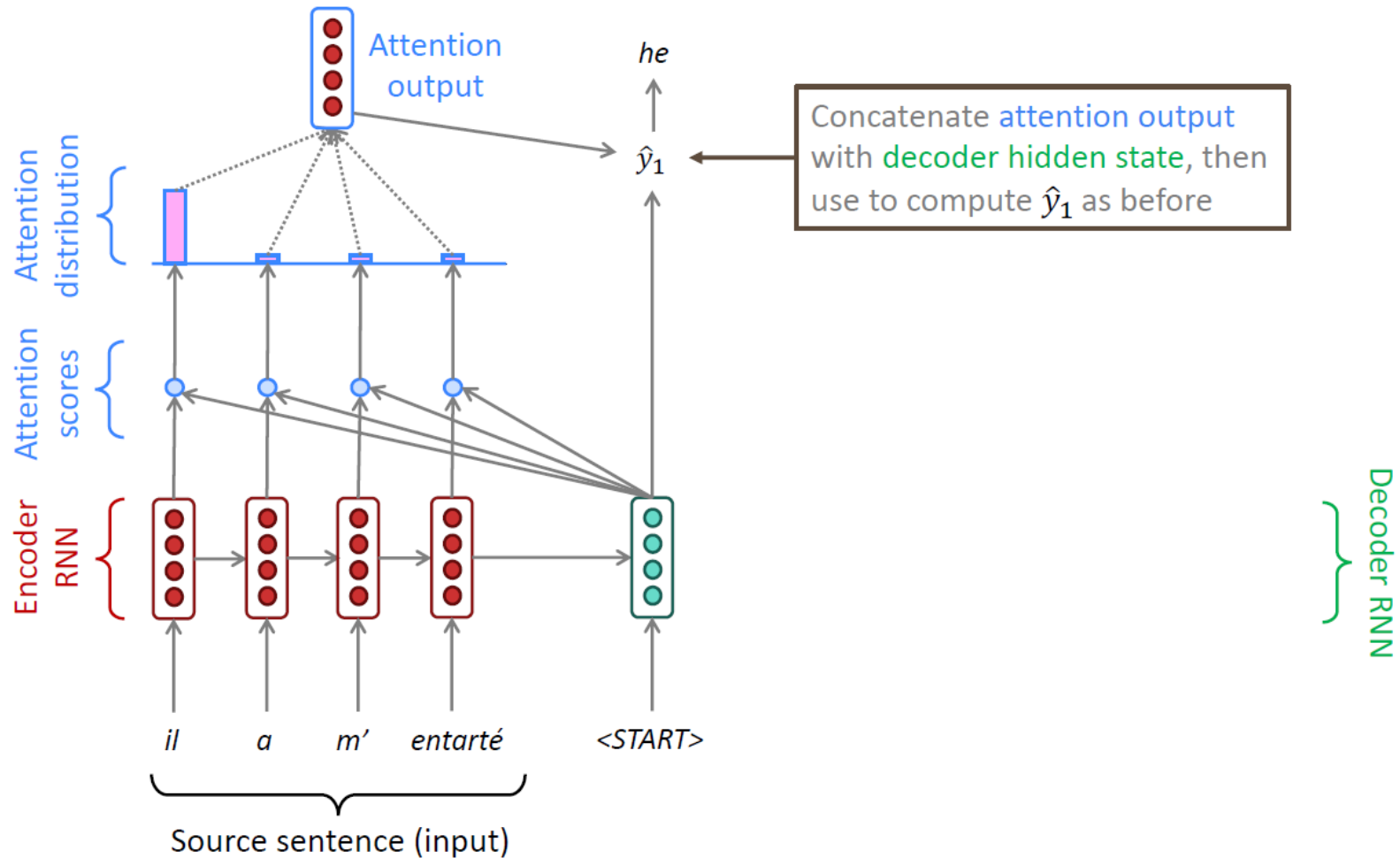




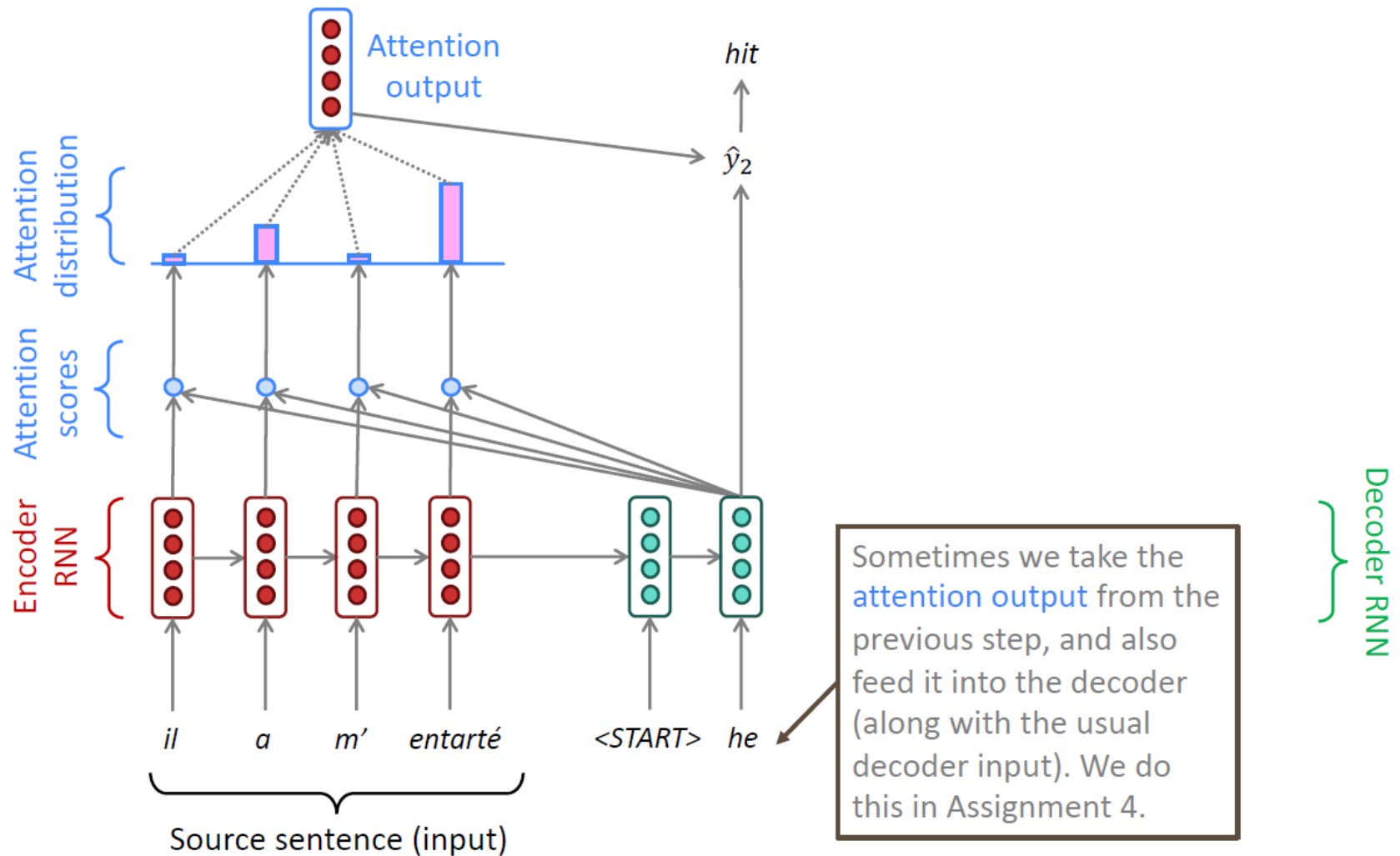
# Sequence to Sequence with RNNs + Attention



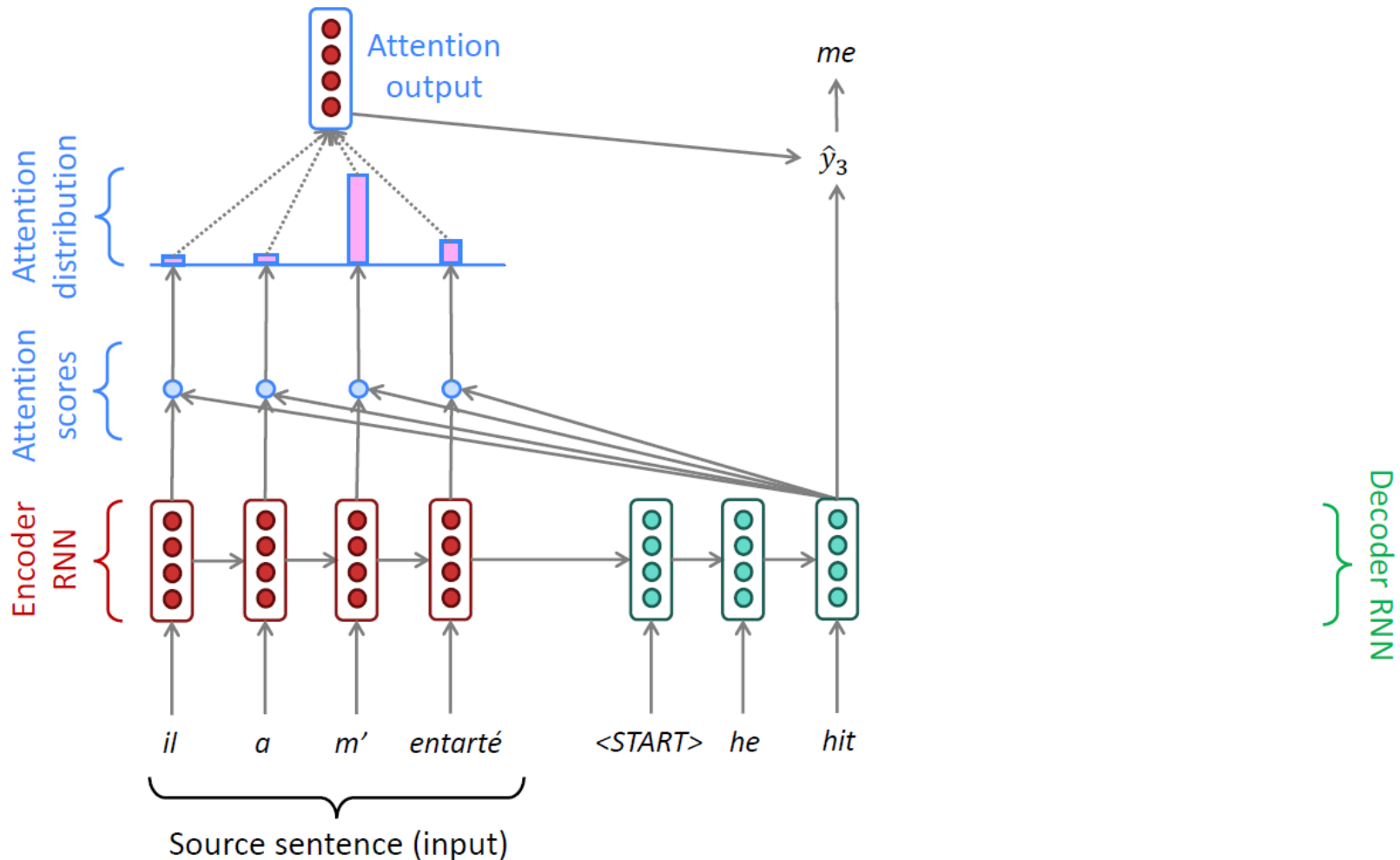
# Sequence to Sequence with RNNs + Attention



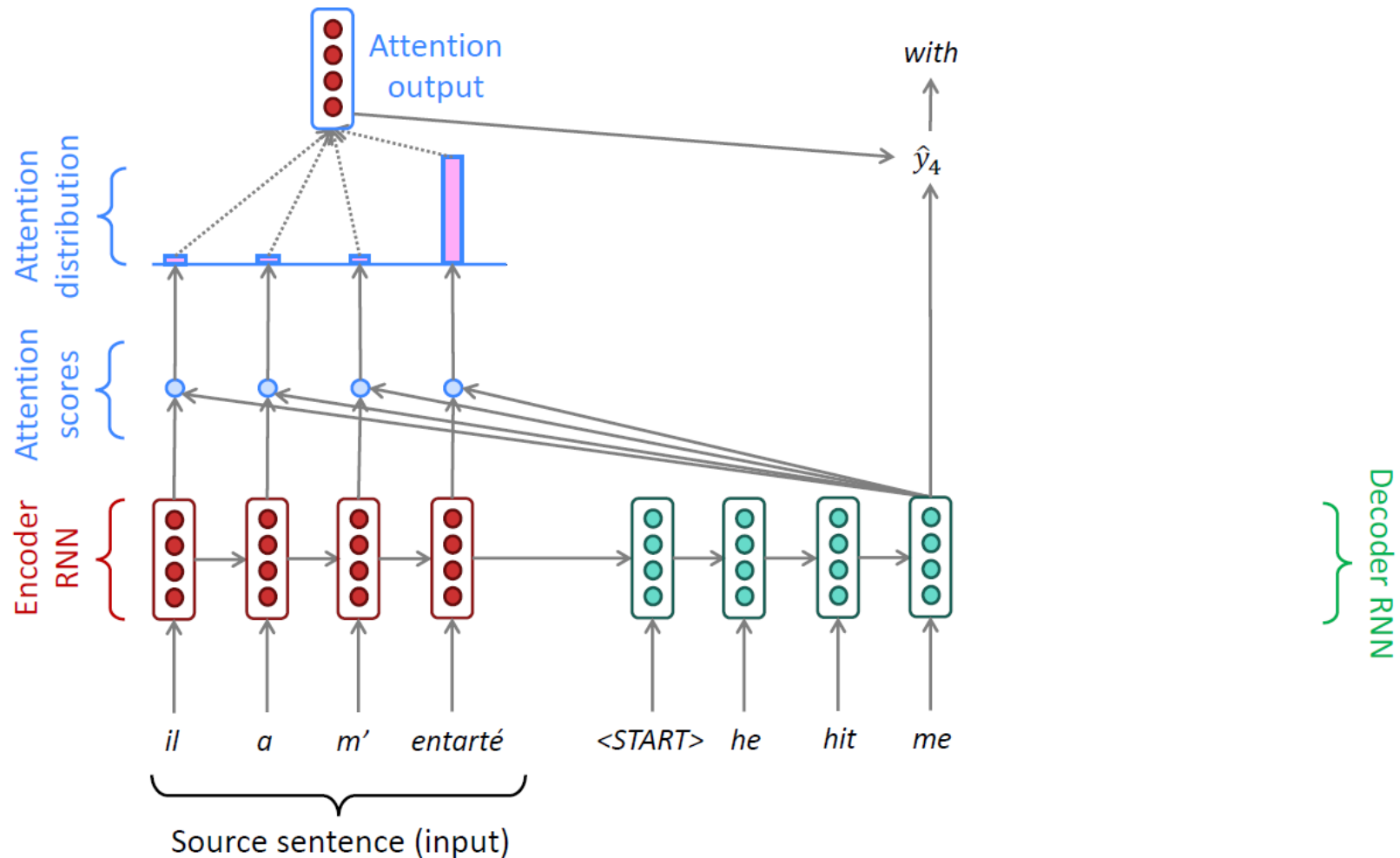
# Sequence to Sequence with RNNs + Attention



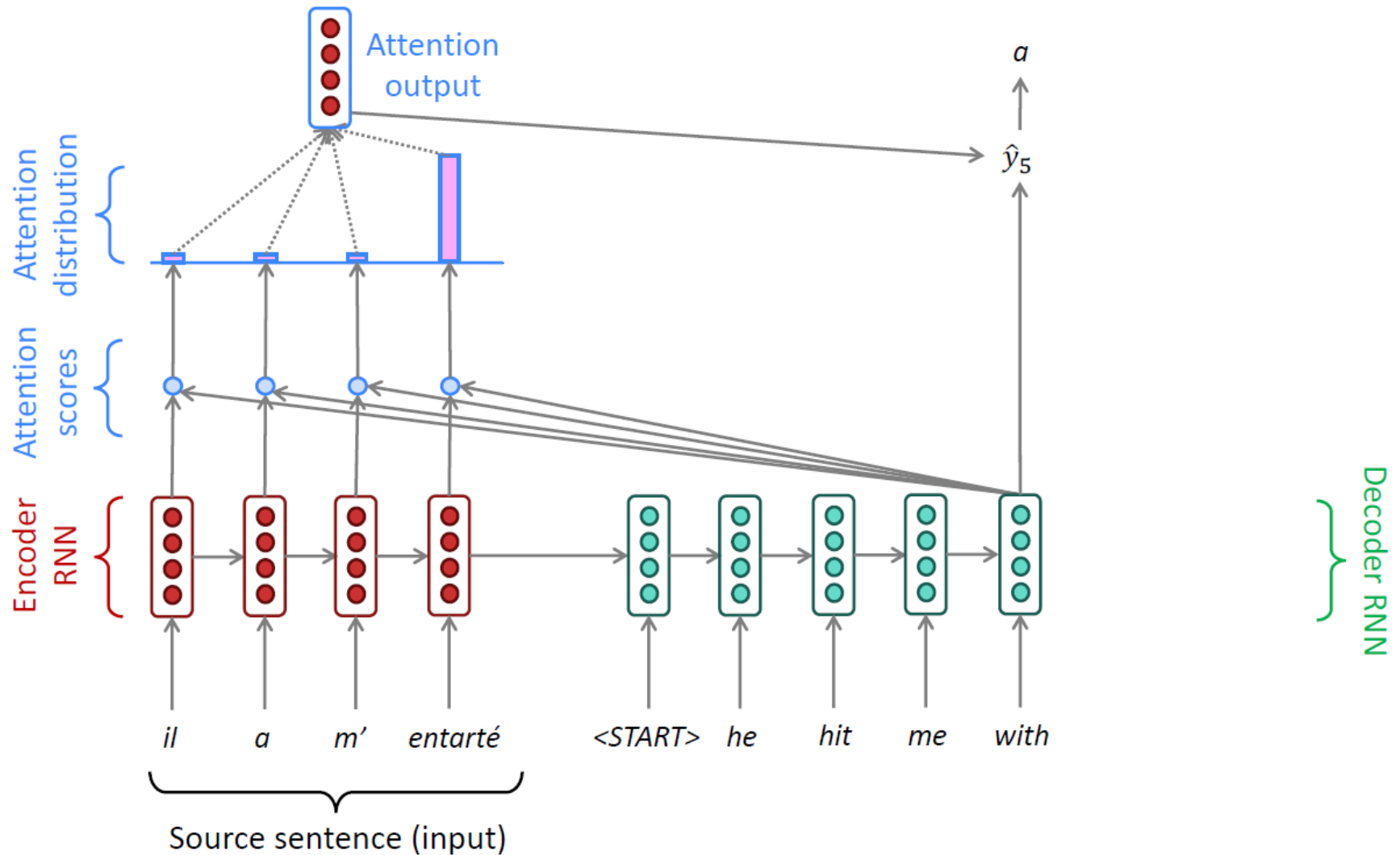
# Sequence to Sequence with RNNs + Attention



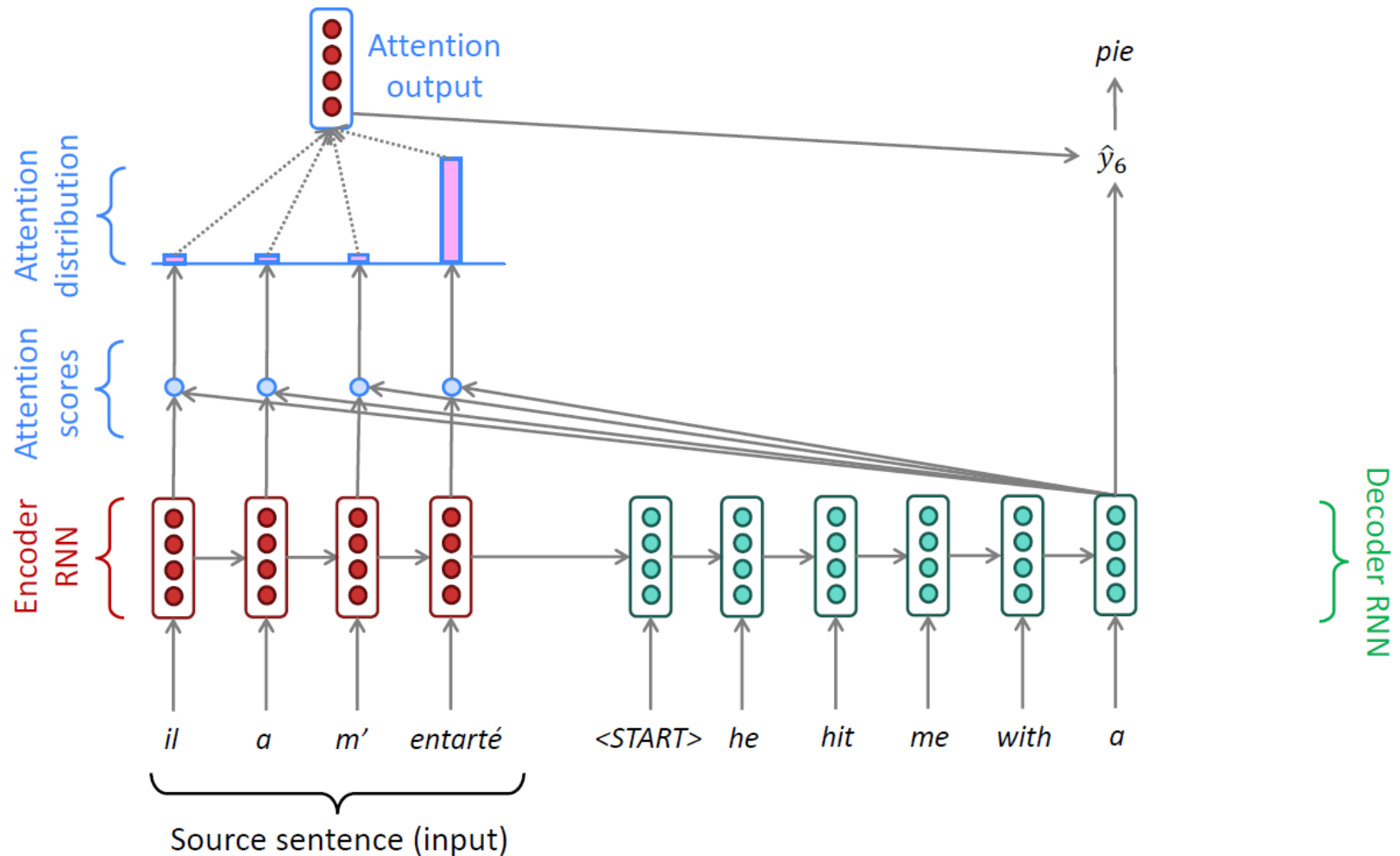
# Sequence to Sequence with RNNs + Attention



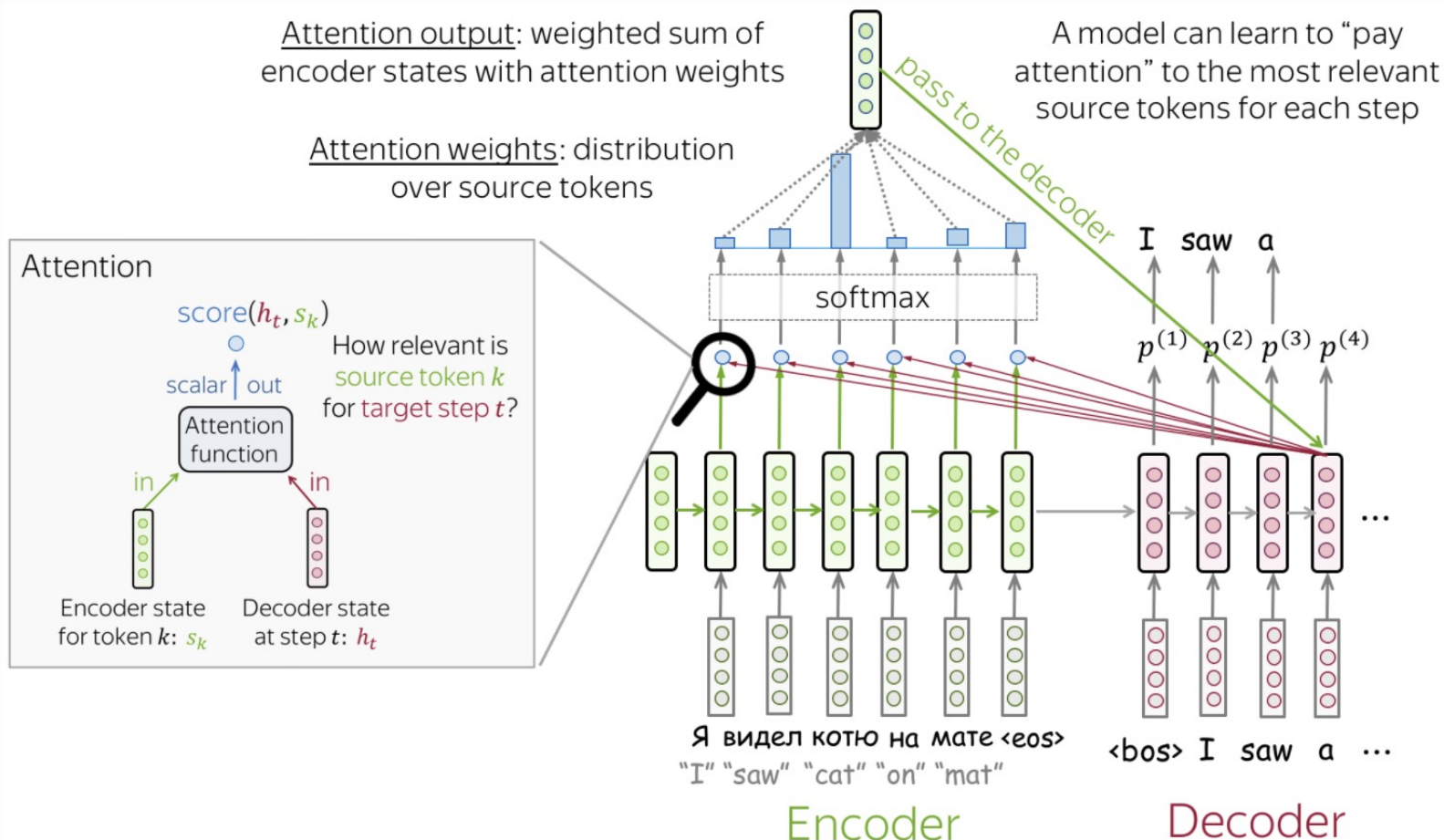
# Sequence to Sequence with RNNs + Attention



# Sequence to Sequence with RNNs + Attention



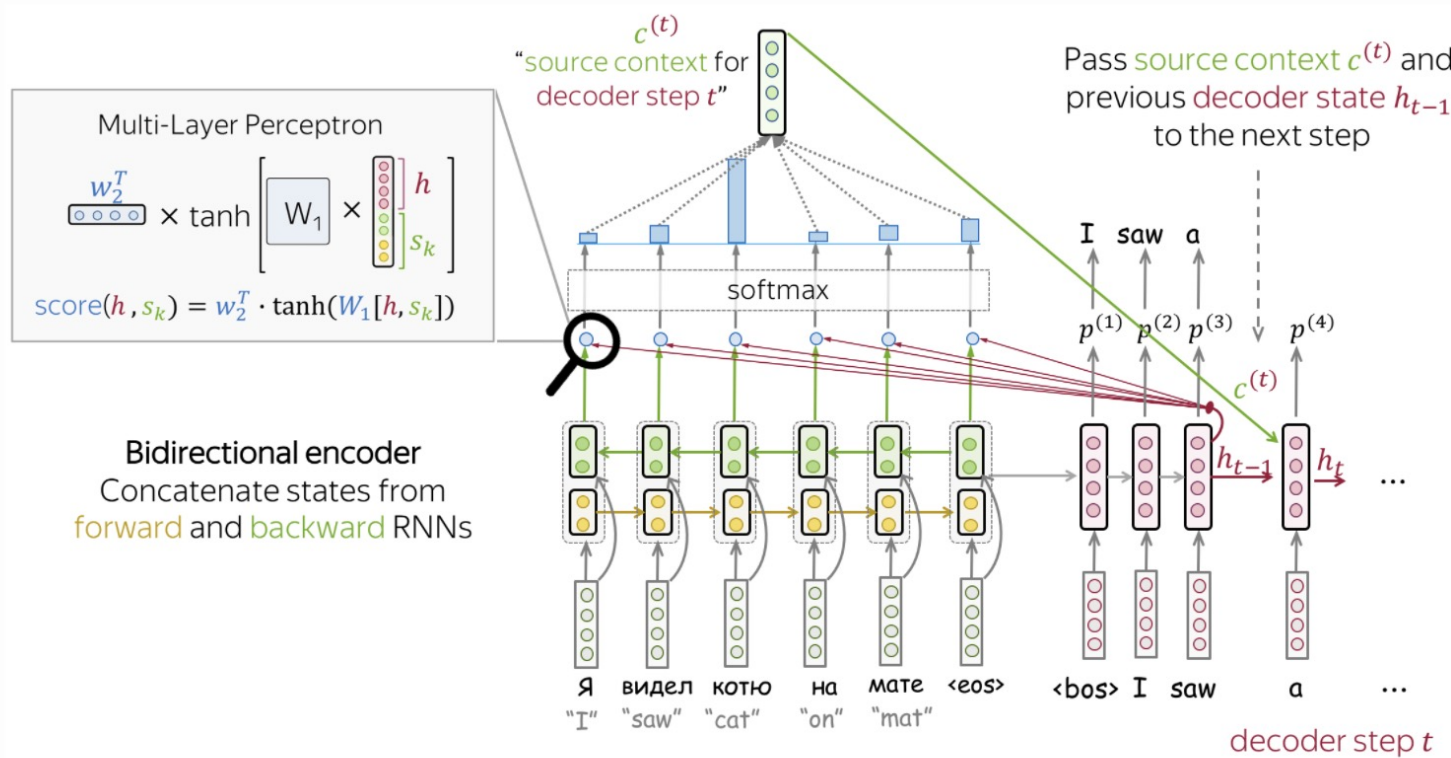
# Attention



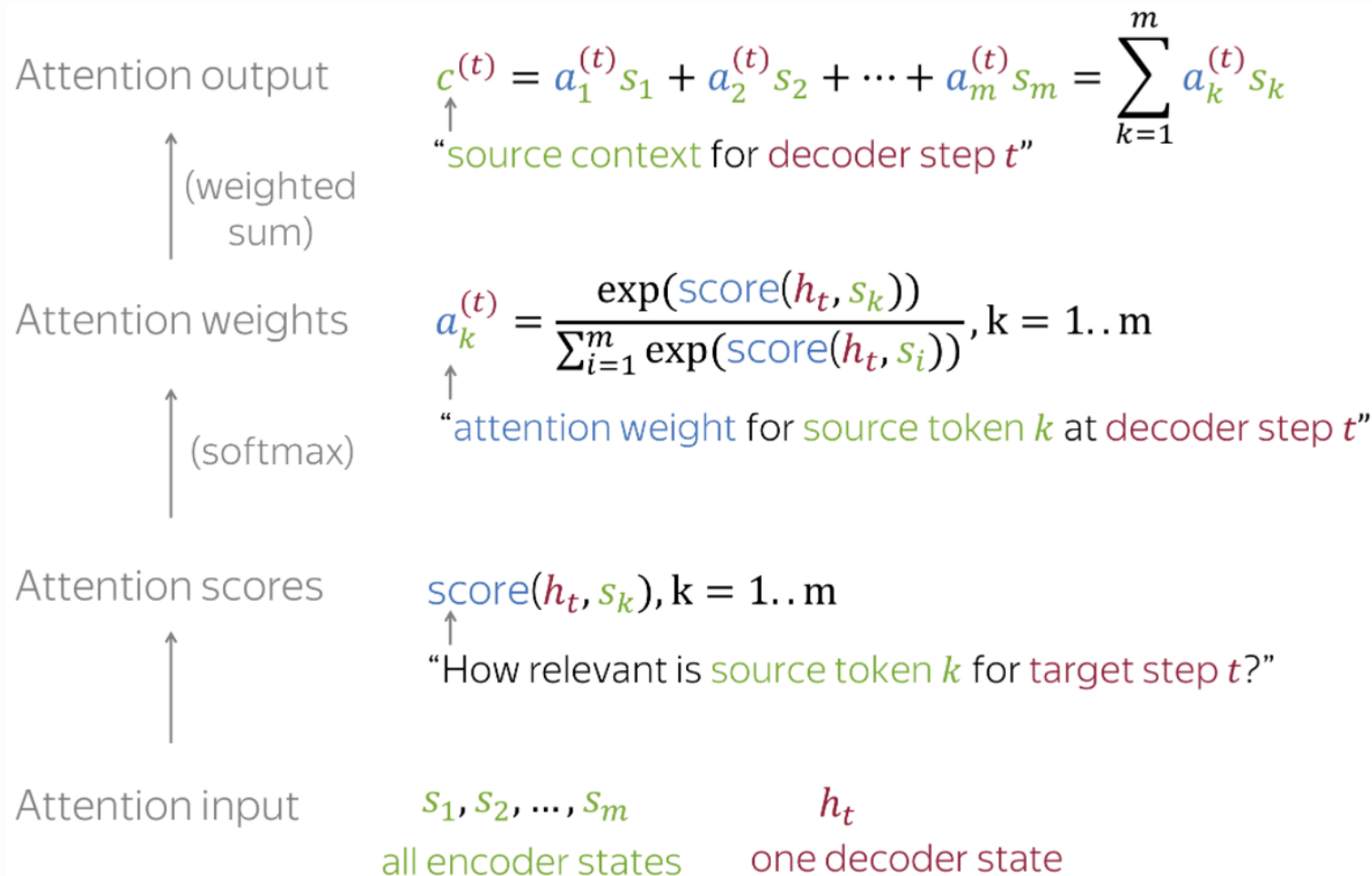


# Bahdanau Model

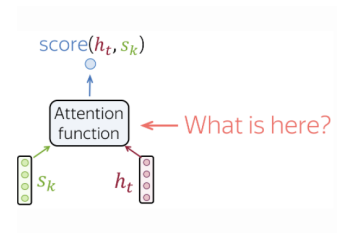
## Attention



# Attention Layer



# Ways to Calculate Attention scores



Dot-product

$$h_t^T \times s_k$$

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

**dot-product** - the simplest method

Bilinear

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

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

**bilinear function** (aka "Luong attention") - used in the paper [Effective Approaches to Attention-based Neural Machine Translation](#)

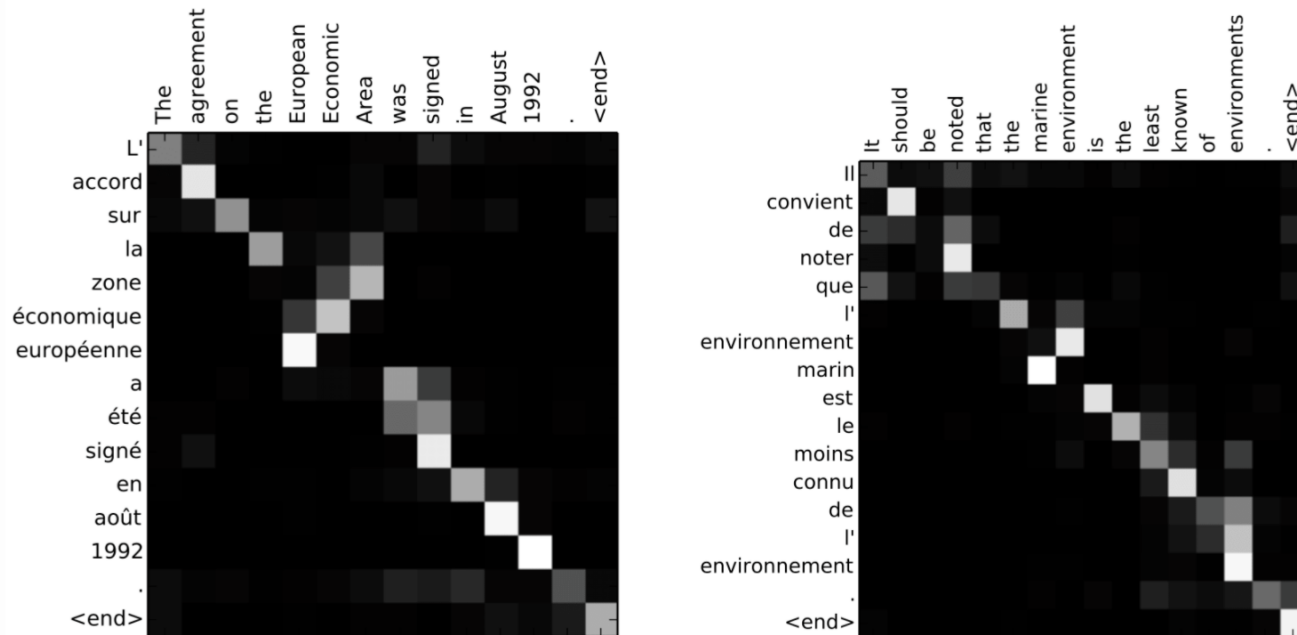
Multi-Layer Perceptron

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

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

**multi-layer perceptron** (aka "Bahdanau attention") - the method proposed in the [original paper](#).

# Attention Learns (Nearly) Alignment



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

# 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 a **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						