

# CIS 6930 Special Topics in Large Language Models

## Sequence-to-Sequence LM

# Outline

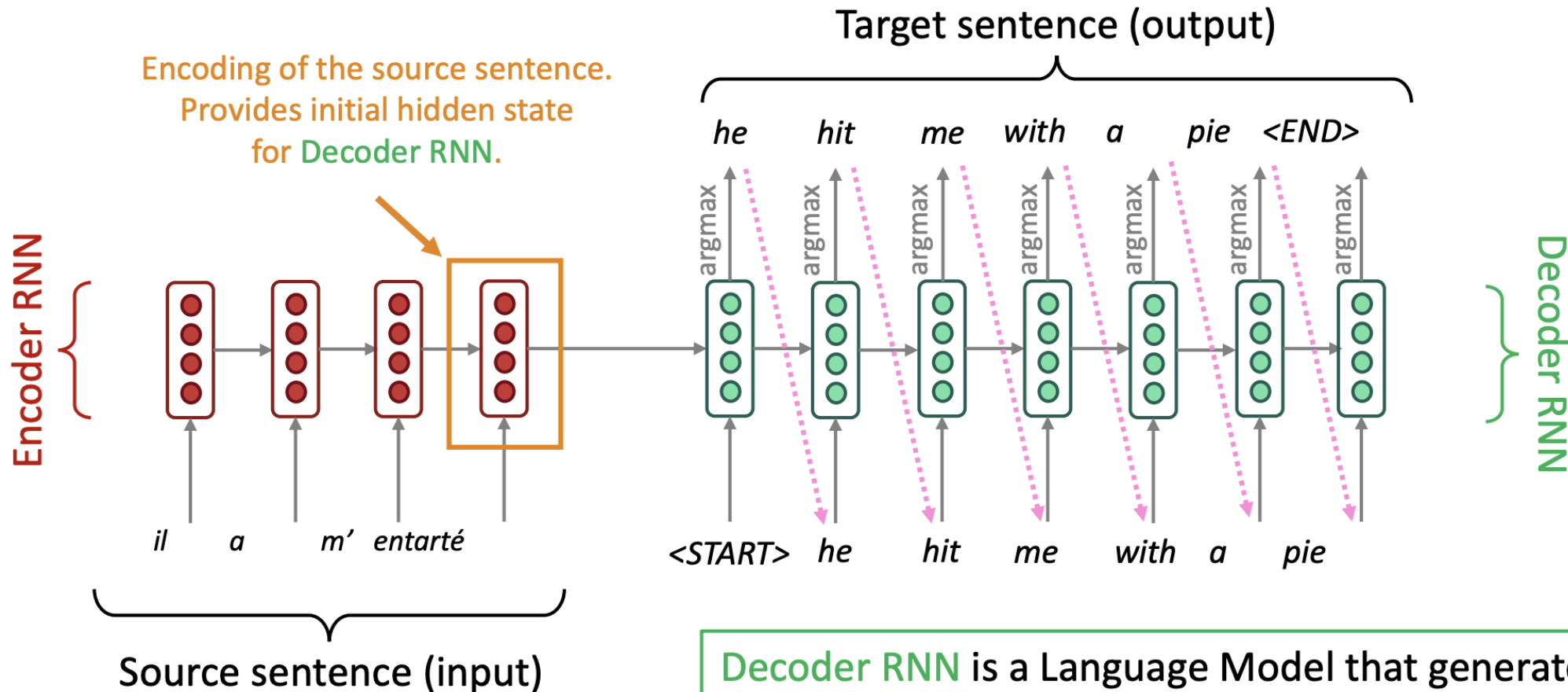
---

- Sequence-to-Sequence Language Model
- Sequence-to-Sequence with Attention

# Sequence-to-Sequence LM

# Sequence-to-Sequence

## Machine Translation



Encoder RNN produces an **encoding** of the source sentence.

Decoder RNN is a Language Model that generates target sentence, *conditioned on encoding*.

Note: This diagram shows **test time** behavior: decoder output is fed in ..... as next step's input

# Sequence-to-Sequence

---

- 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
- 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)

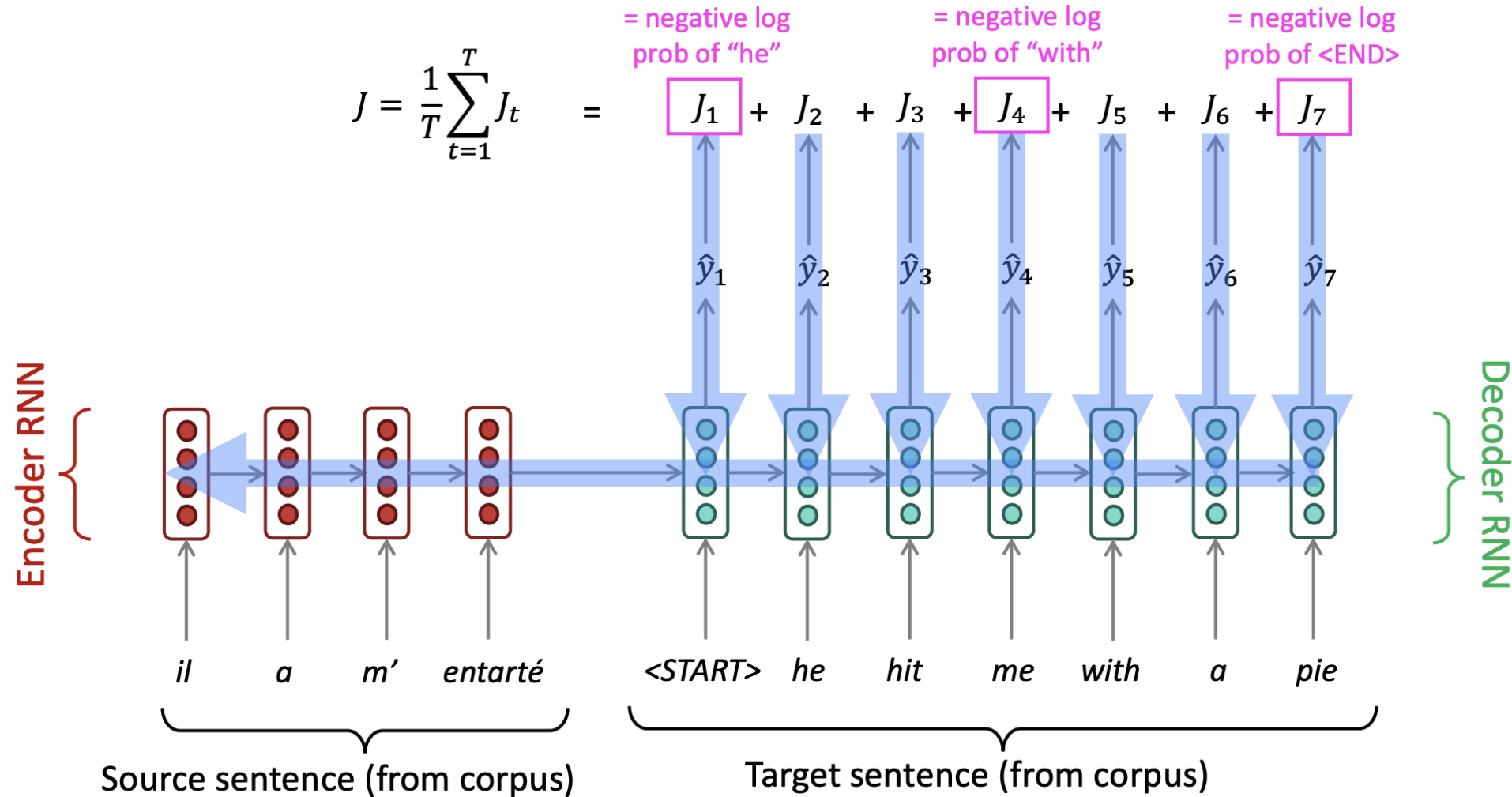
# Sequence-to-Sequence

- 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  $x$

# Training Sequence-to-Sequence



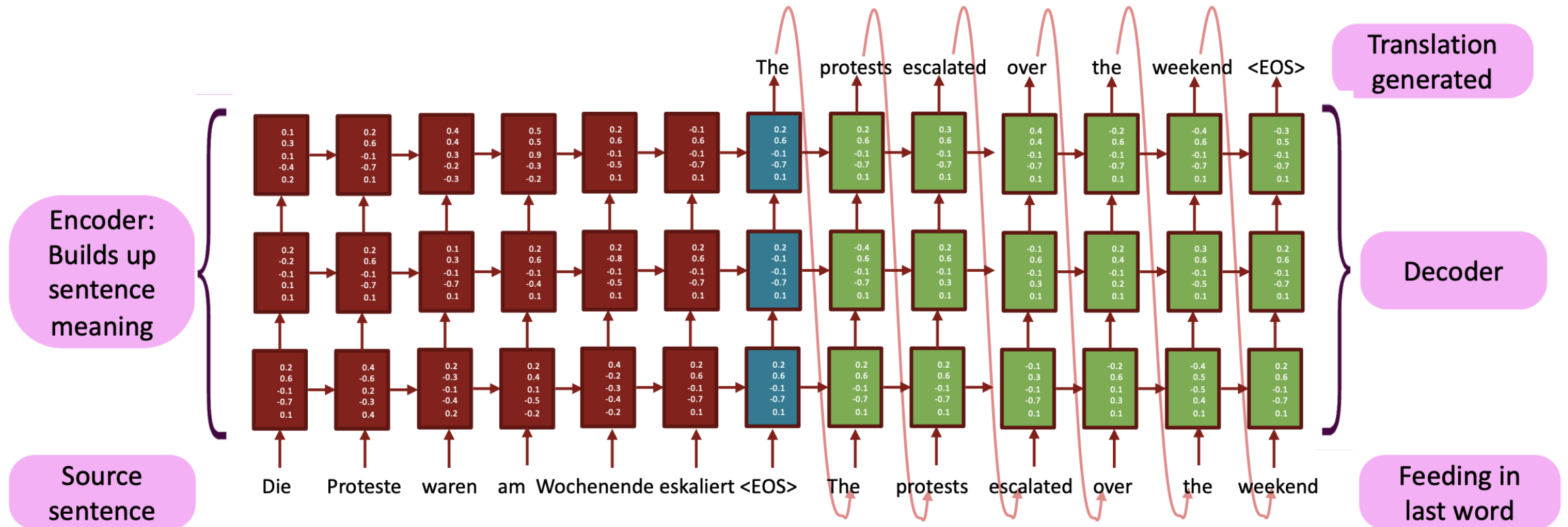
Seq2seq is optimized as a **single system**. Backpropagation operates "*end-to-end*".



# Multi-layer Sequence-to-Sequence

[Sutskever et al. 2014; Luong et al. 2015]

The hidden states from RNN layer  $i$  are the inputs to RNN layer  $i+1$





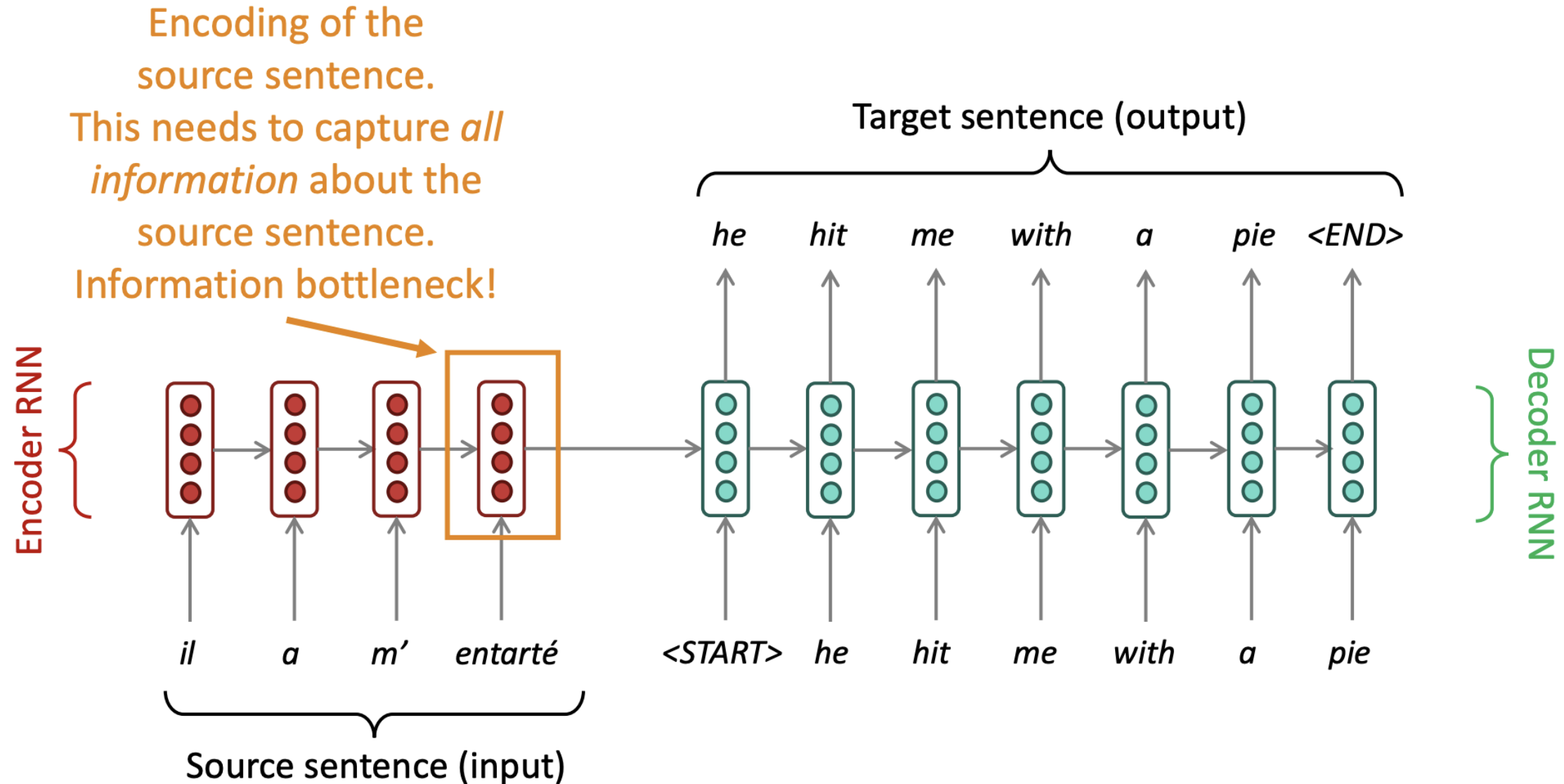
# Multi-layer Sequence-to-Sequence

---

- Multi-layer or stacked RNNs allow the network to compute **more complex representations**
  - The **lower RNNs** should **compute lower-level features** and the **higher RNNs** should compute **higher-level features**.
- **High-performing RNNs are usually multi-layer** (but aren't as deep as convolutional or feed-forward networks)
- For example: In a 2017 paper, Britz et al. find that for Neural Machine Translation, **2 to 4 layers** is best for the encoder RNN, and **4 layers** is best for the decoder RNN
  - Often 2 layers is a lot better than 1, and 3 might be a little better than 2
  - Usually, **skip-connections/dense-connections** are needed to train deeper RNNs (e.g., **8 layers**)

# Sequence-to-Sequence with Attention

# Sequence-to-Sequence Bottleneck

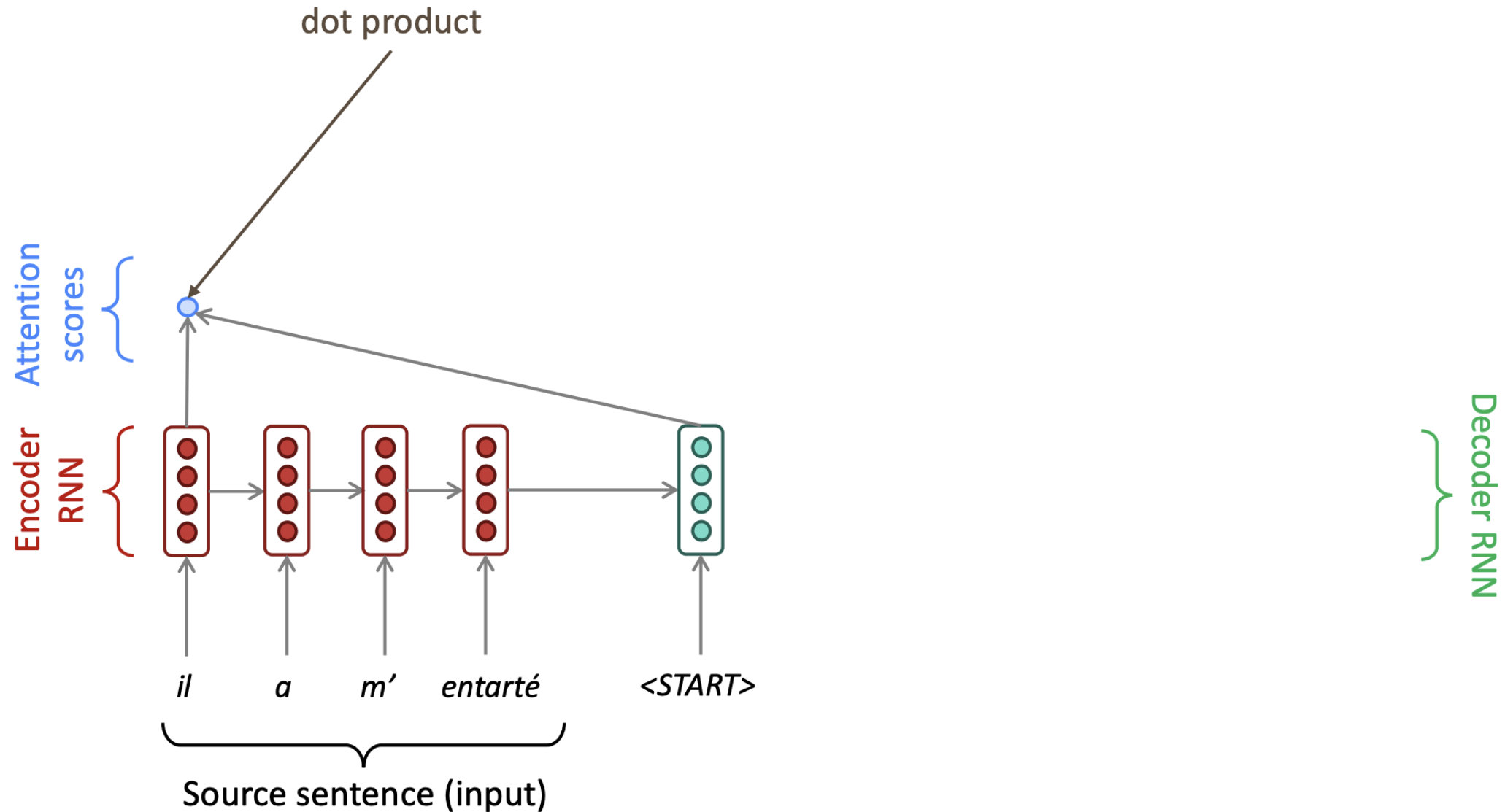


# Sequence-to-Sequence with Attention

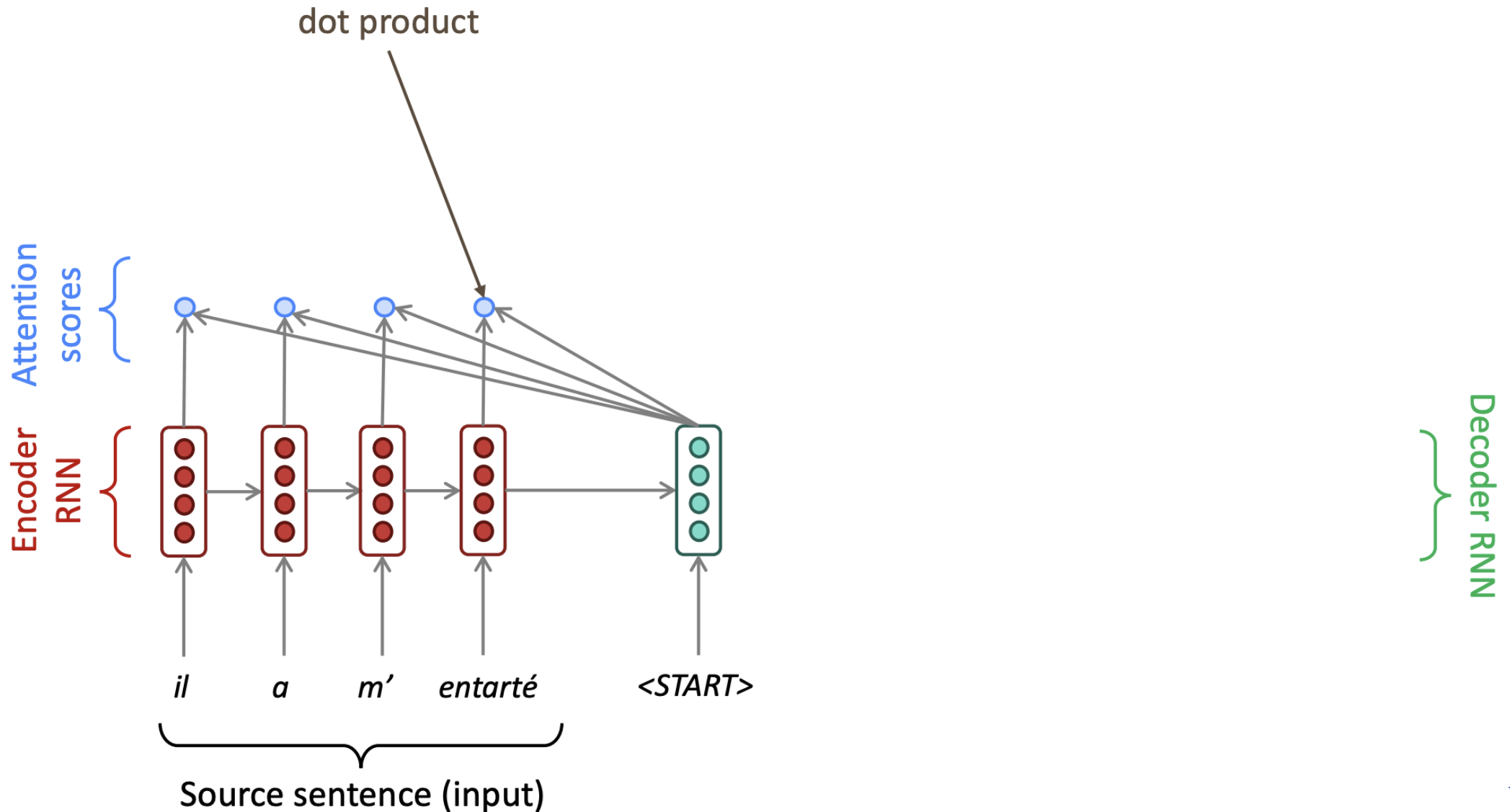
---

- **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

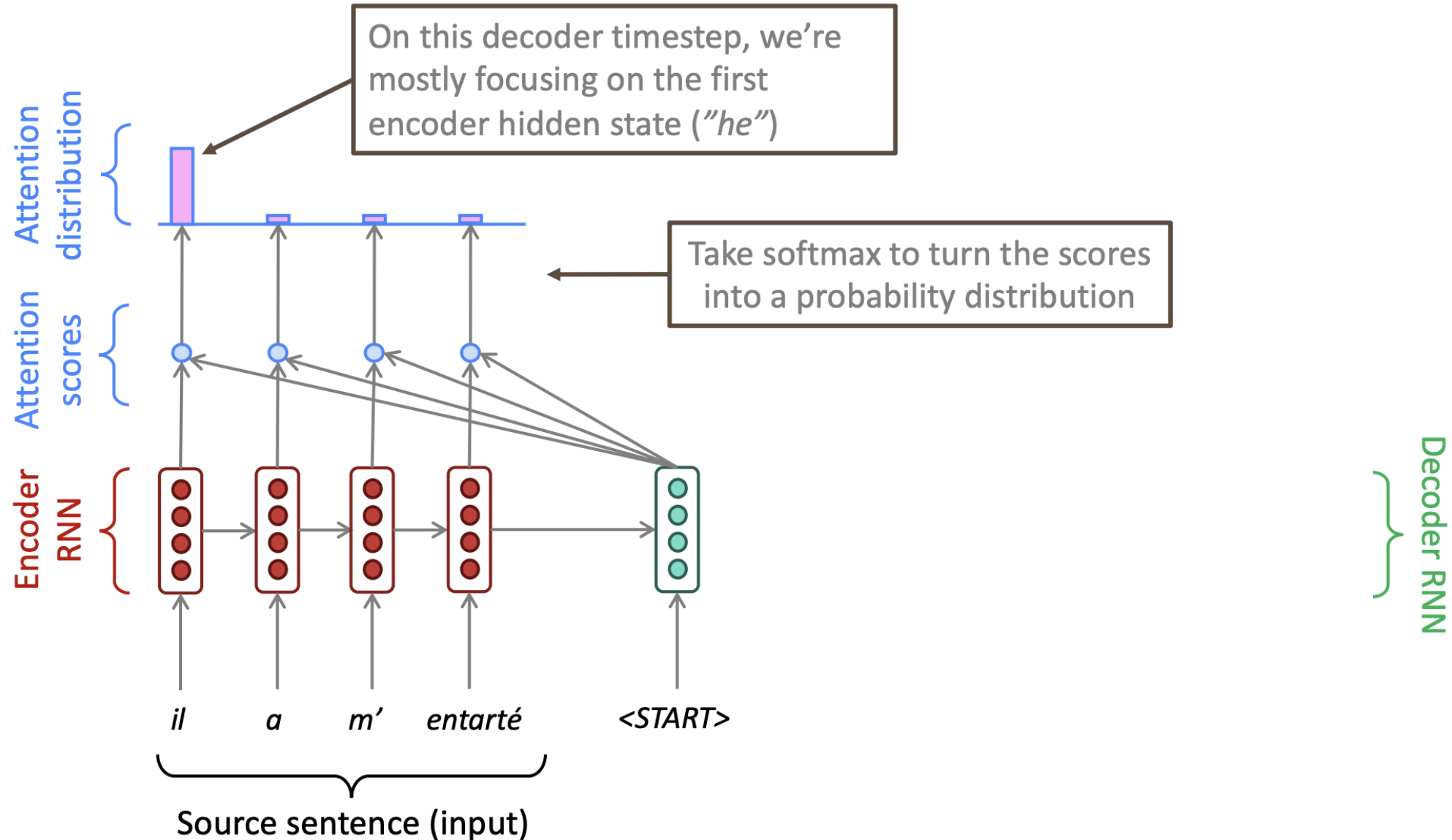
# Sequence-to-Sequence with Attention



# Sequence-to-Sequence with Attention

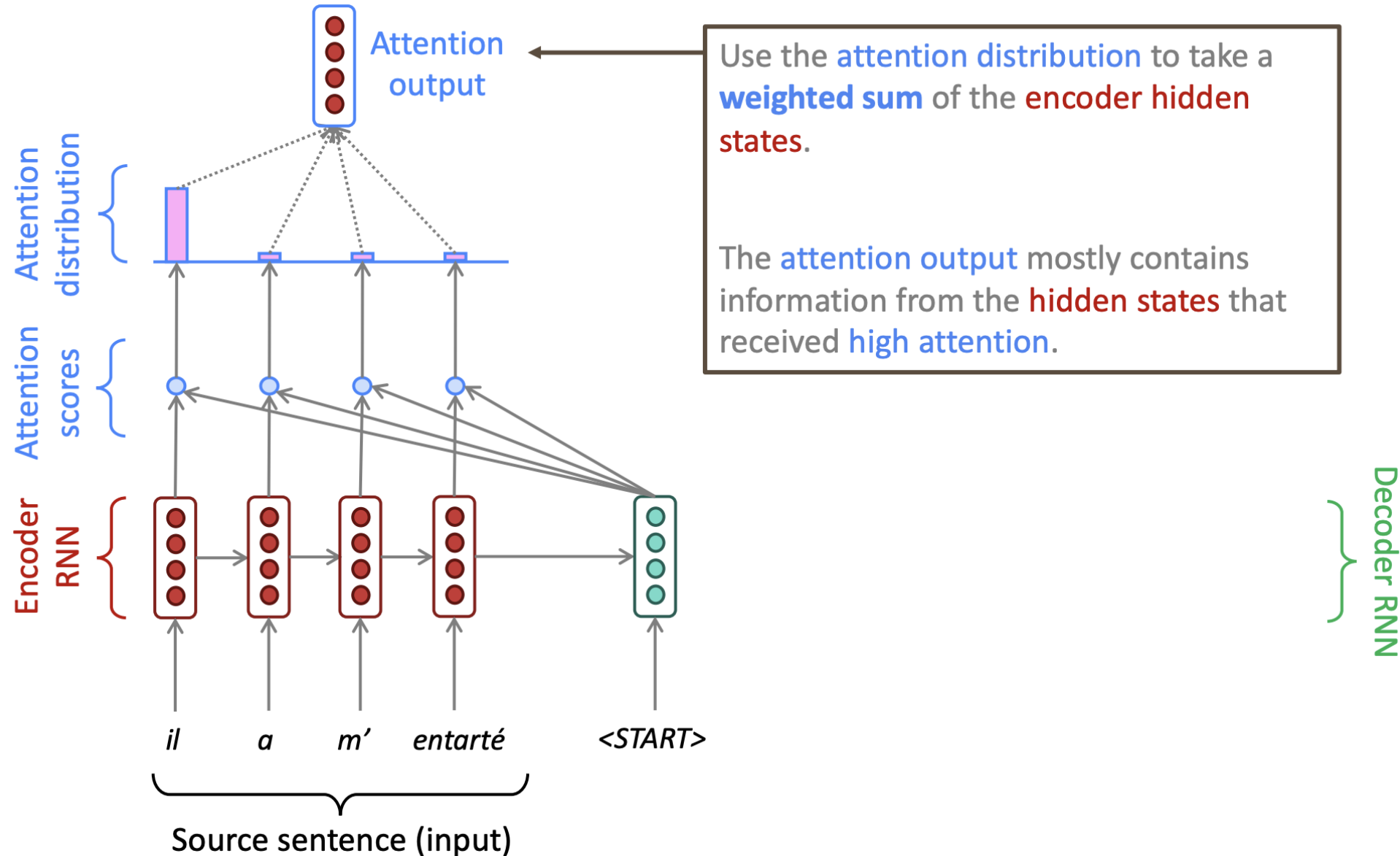


# Sequence-to-Sequence with Attention

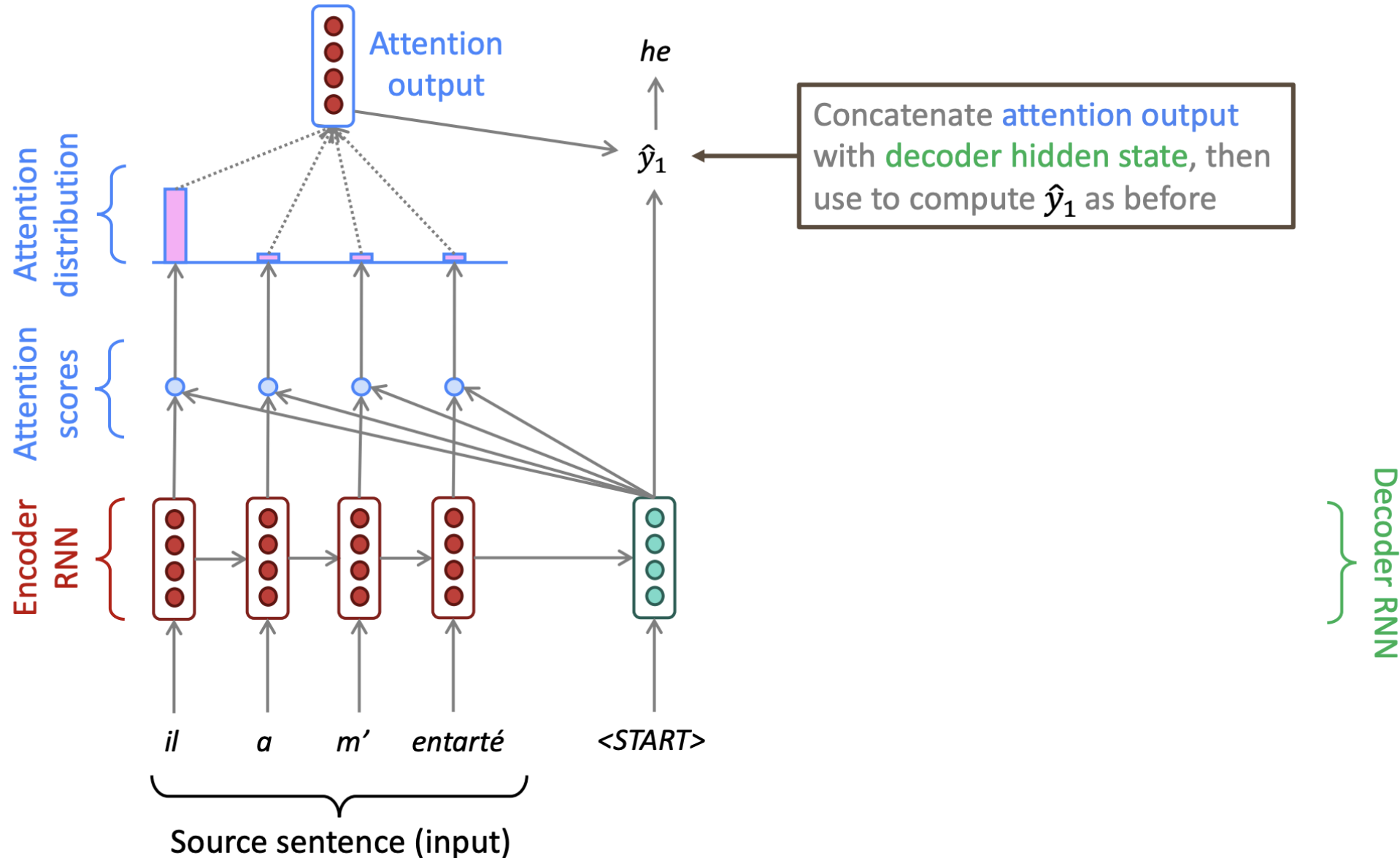




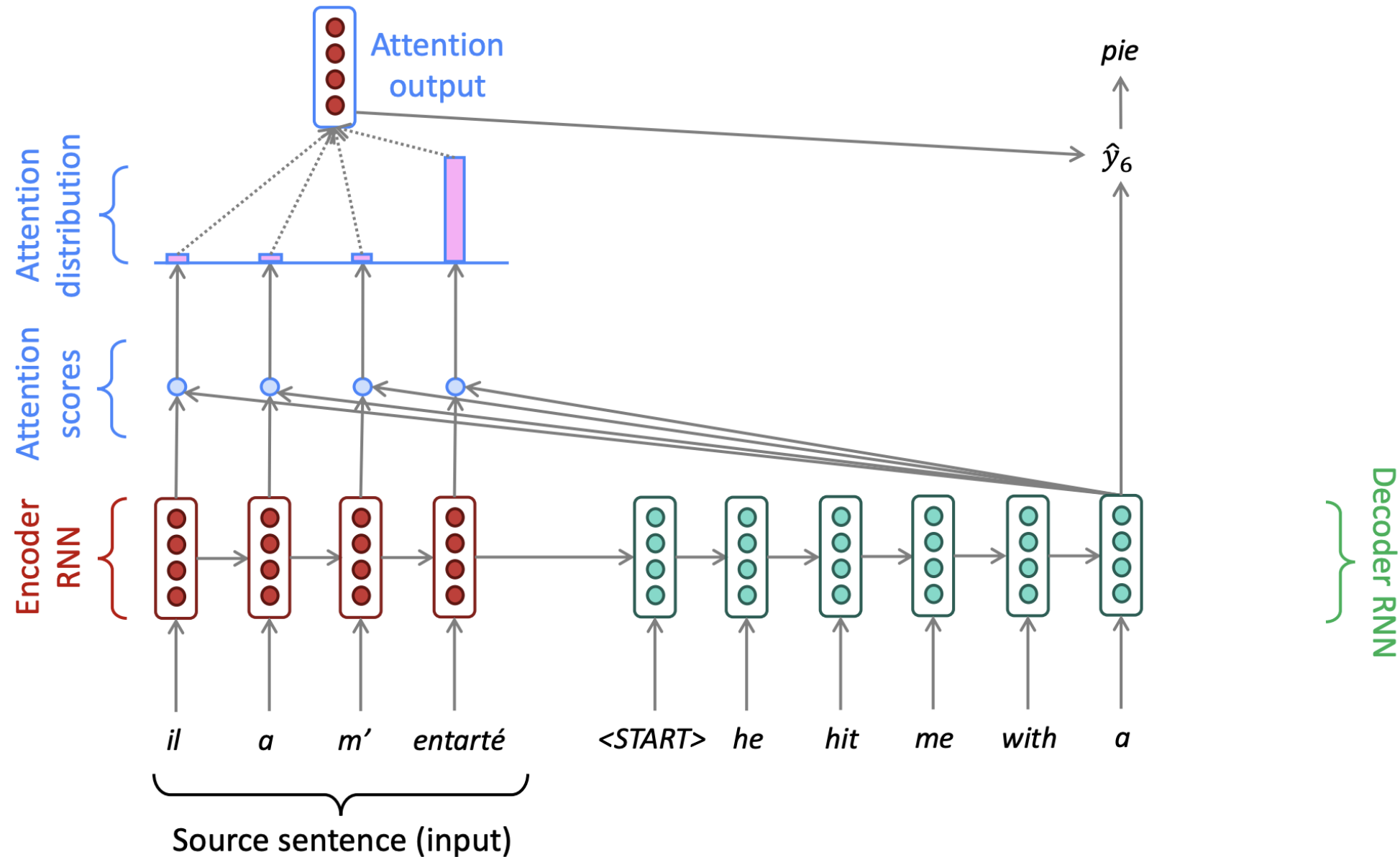
# Sequence-to-Sequence with Attention



# Sequence-to-Sequence with Attention



# Sequence-to-Sequence with Attention



# Sequence-to-Sequence with Attention

- We have encoder hidden states  $h_1, \dots, h_N \in \mathbb{R}^h$
- On timestep  $t$ , we have decoder hidden state  $s_t \in \mathbb{R}^h$
- We get the attention scores  $e^t$  for this step:

$$e^t = [s_t^T h_1, \dots, s_t^T h_N] \in \mathbb{R}^N$$

- We take softmax to get the attention distribution  $\alpha^t$  for this step (this is a probability distribution and sums to 1)

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

- We use  $\alpha^t$  to take a weighted sum of the encoder hidden states to get the attention output  $a_t$

$$a_t = \sum_{i=1}^N \alpha_i^t h_i \in \mathbb{R}^h$$

- Finally we concatenate the attention output  $a_t$  with the decoder hidden state  $s_t$  and proceed as in the non-attention seq2seq model

$$[a_t; s_t] \in \mathbb{R}^{2h}$$

# Attention

---

- 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 far-away states

# 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*

$$\mathbf{e} \in \mathbb{R}^N$$

There are  
multiple ways  
to do this

2. Taking softmax to get *attention distribution*  $\alpha$ :

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

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 Variants

There are **several ways** you can compute  $e \in \mathbb{R}^N$  from  $\mathbf{h}_1, \dots, \mathbf{h}_N \in \mathbb{R}^{d_1}$  and  $\mathbf{s} \in \mathbb{R}^{d_2}$ :

- Basic dot-product attention:  $e_i = \mathbf{s}^T \mathbf{h}_i \in \mathbb{R}$ 
  - Note: this assumes  $d_1 = d_2$
  - This is the version we saw earlier
- Multiplicative attention:  $e_i = \mathbf{s}^T \mathbf{W} \mathbf{h}_i \in \mathbb{R}$ 
  - Where  $\mathbf{W} \in \mathbb{R}^{d_2 \times d_1}$  is a weight matrix
- Additive attention:  $e_i = \mathbf{v}^T \tanh(\mathbf{W}_1 \mathbf{h}_i + \mathbf{W}_2 \mathbf{s}) \in \mathbb{R}$ 
  - Where  $\mathbf{W}_1 \in \mathbb{R}^{d_3 \times d_1}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{d_3 \times d_2}$  are weight matrices and  $\mathbf{v} \in \mathbb{R}^{d_3}$  is a weight vector.
  - $d_3$  (the attention dimensionality) is a hyperparameter



# Attention is a general concept

---

- You can use attention in many architectures and many tasks
- 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**.

- We sometimes say that the query attends to the values.

For example, in the seq2seq + attention model, each decoder hidden state (**query**) attends to all the encoder hidden states (**values**).

# Attention is a general concept

---

- 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).



Thank you!

**UF** | Herbert Wertheim  
College of Engineering  
UNIVERSITY *of* FLORIDA

---

LEADING THE CHARGE, CHARGING AHEAD