

A close-up photograph of a person's hands holding a pair of black-framed glasses. The lenses are reflecting a blurred image of a large industrial port with numerous blue cranes and shipping containers. The background is a soft-focus view of the same port through the glasses.

# Attention Mechanism

# Agenda

- RNN
- Seq2Seq
- Attention
- Implement



### 2014 Seq2seq、GRU

- 1. Learning phrase representations
- 2. Seq2Seq Learning with NN

### 2015

4. Effective Approaches to Attention-based  
Attention 概念延伸

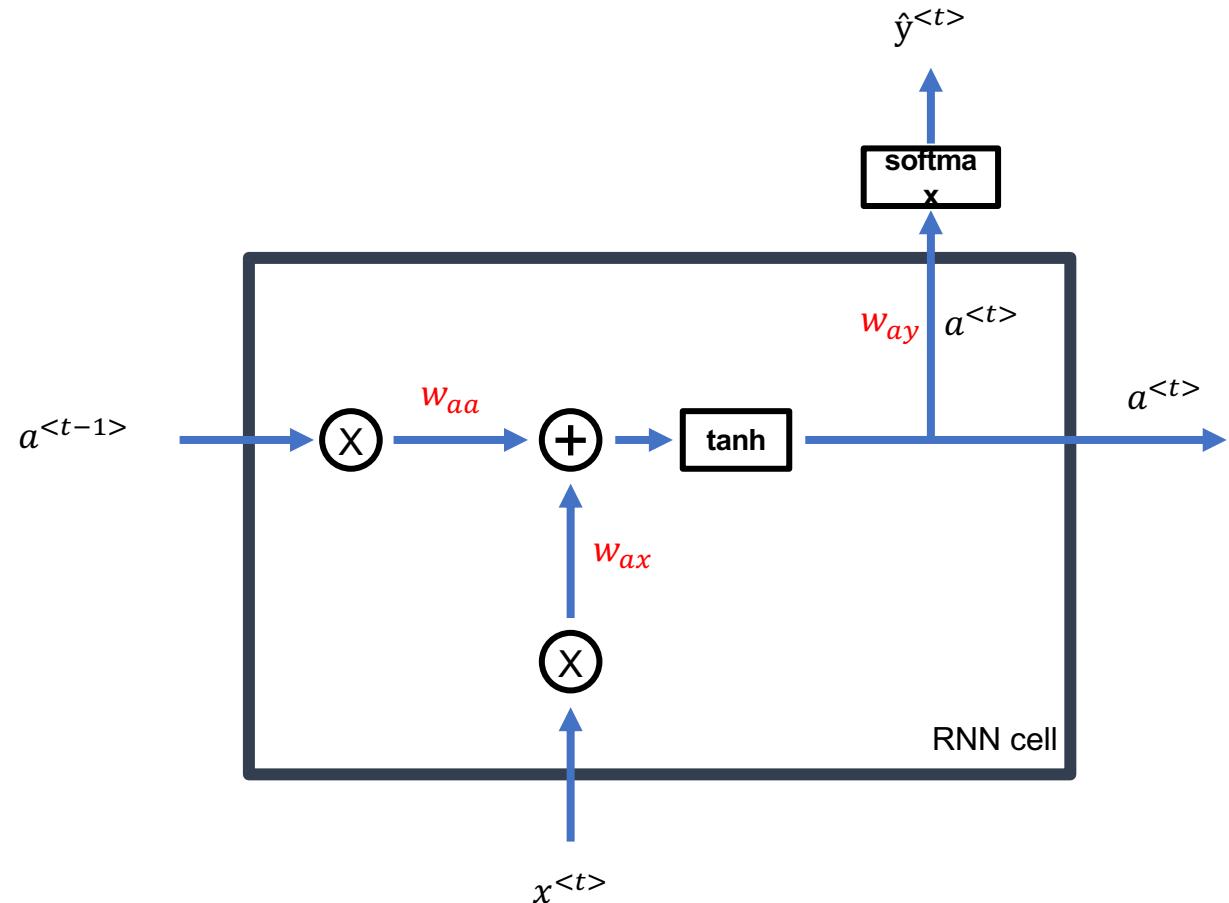
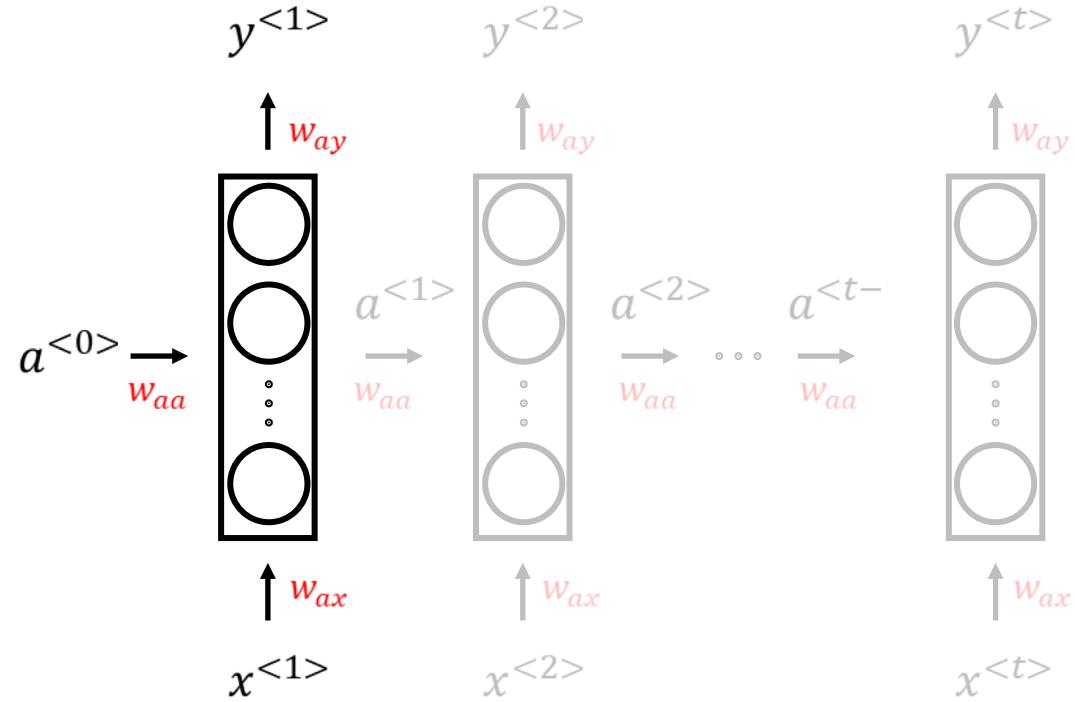
### 2017 Self-Attention取代Cnn、Rnn

- 6. Attention is all you need
- 7. Atrank: An Attention-Based User Behavior



# Close Look at RNN Cell

## Recurrent Neural Network

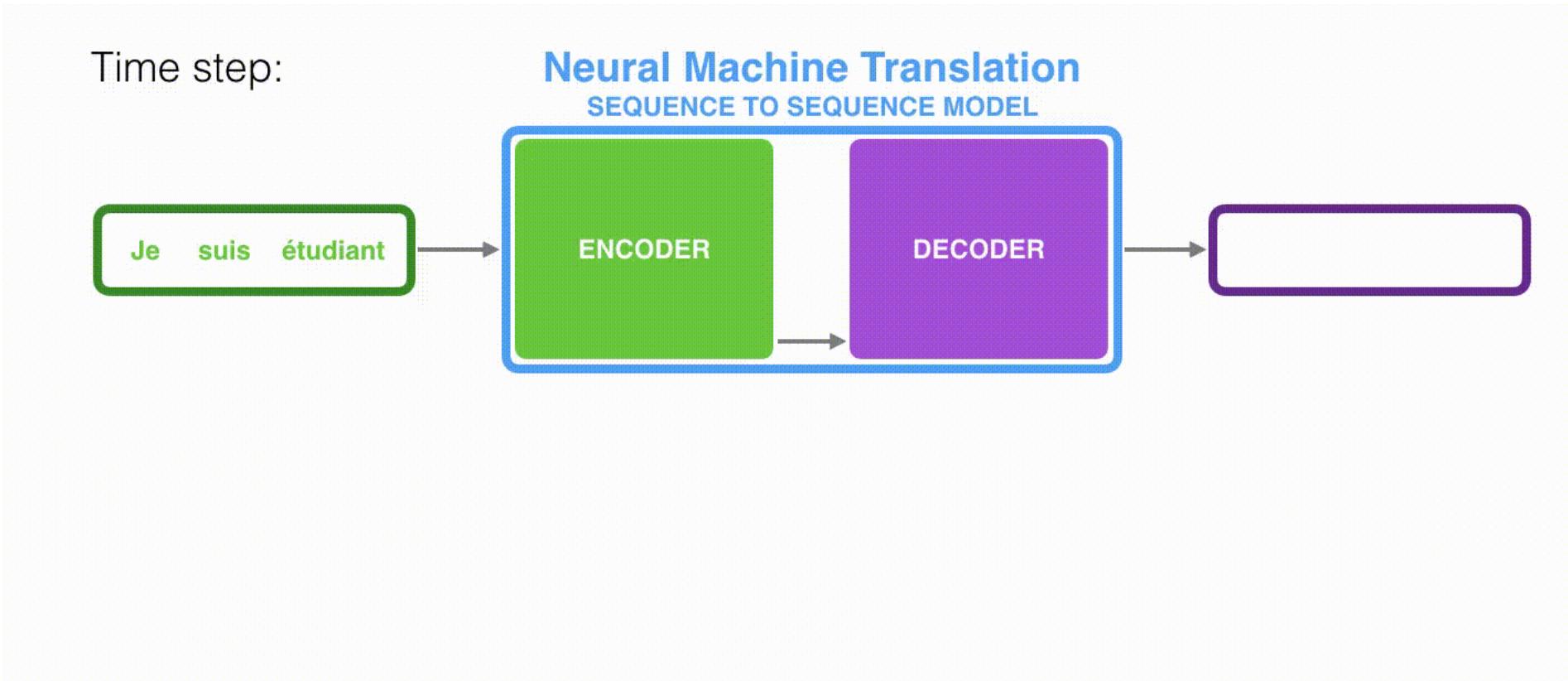


$$a^{<t>} = \tanh(W_{ax}x^{<t>} + W_{aa}a^{<t-1>} + b_a)$$

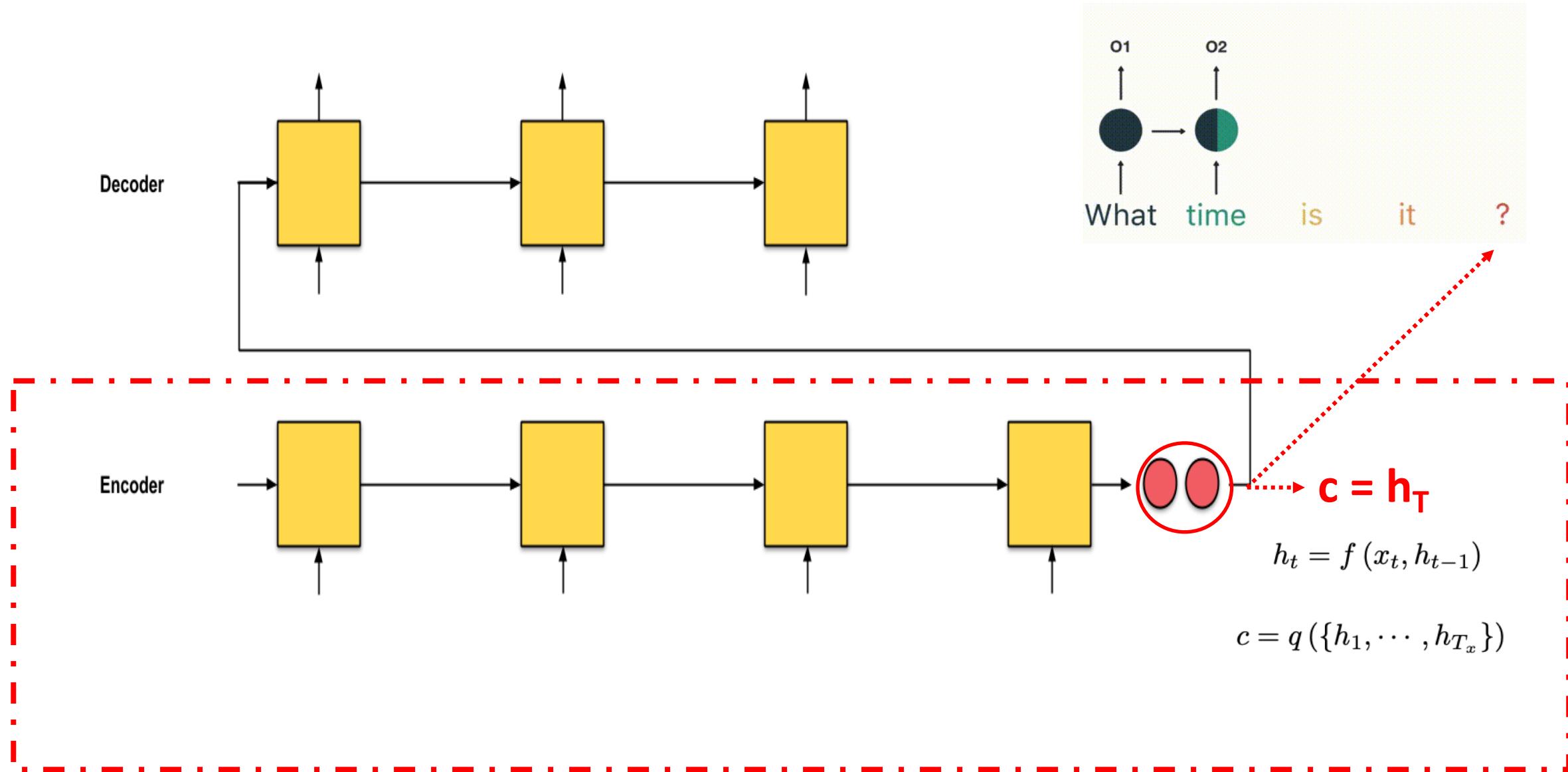
$$\hat{y}^{<t>} = \text{softmax}(W_{ay}a^{<t>} + b_y)$$

# Seq2Seq

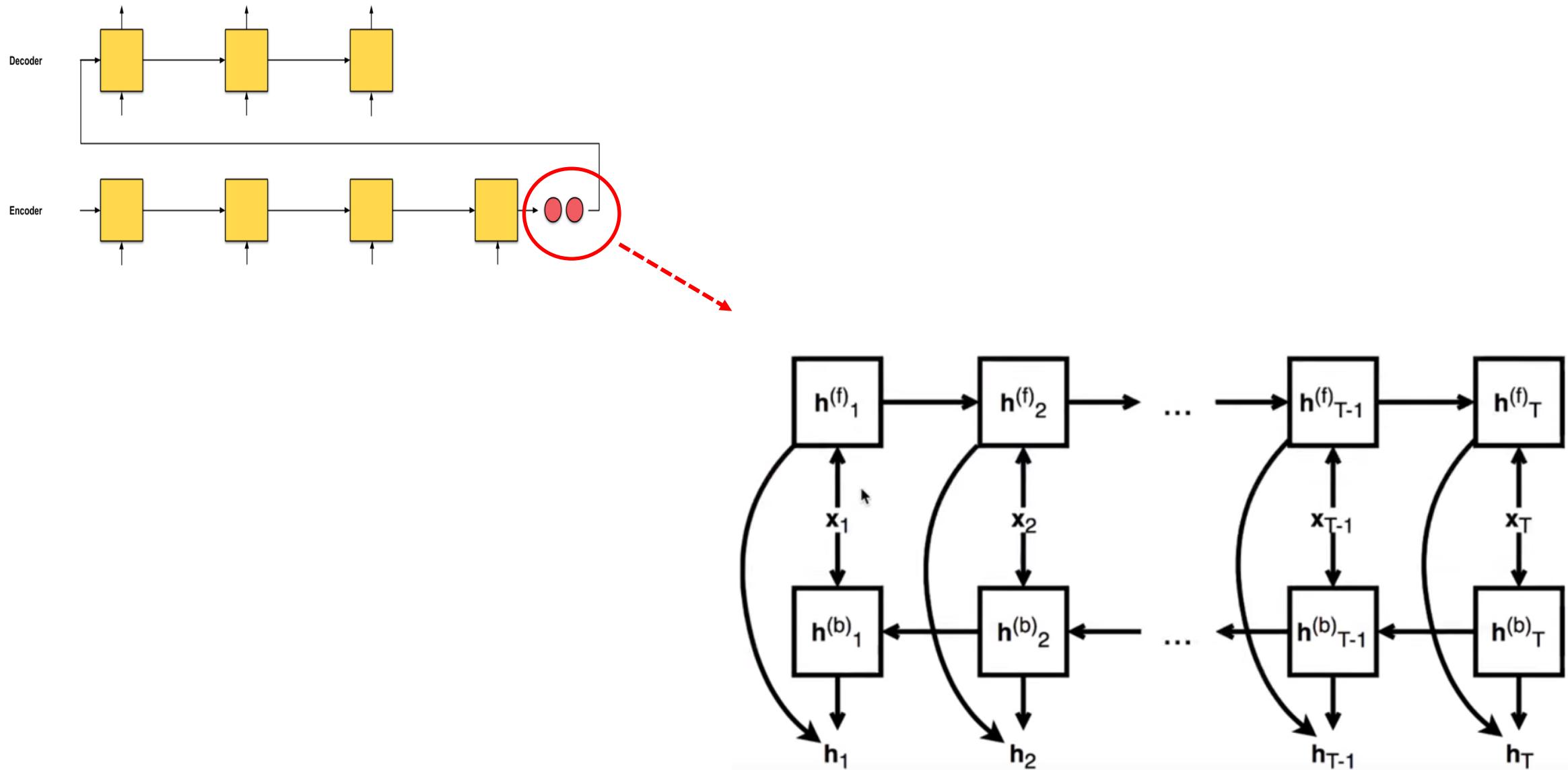
RNN input:長度不固定 ouput:固定



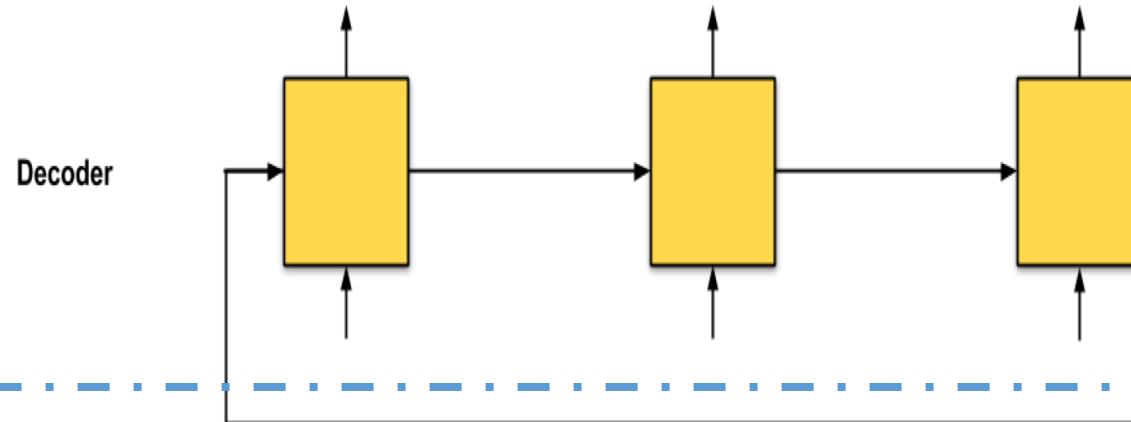
# Seq2Seq



# Seq2Seq

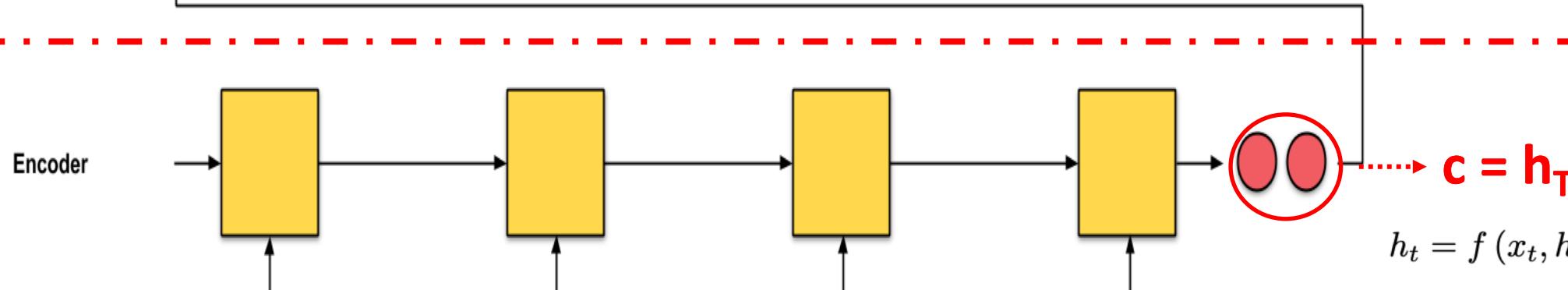


# Seq2Seq



$$p(\mathbf{y}) = \prod_{t=1}^T p(y_t | \{y_1, \dots, y_{t-1}\}, c),$$

$$p(y_t | \{y_1, \dots, y_{t-1}\}, c) = g(y_{t-1}, s_t, c),$$

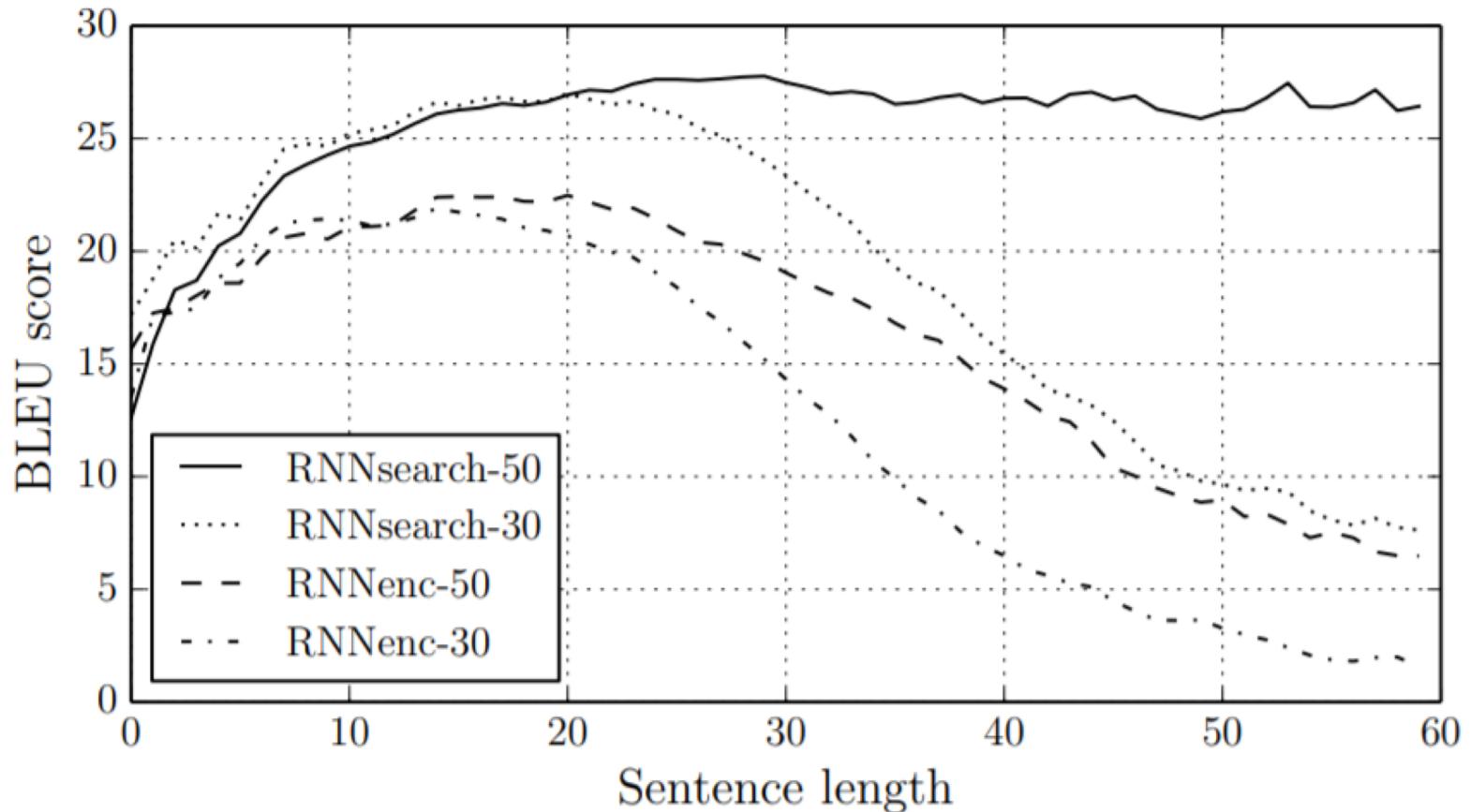


$$h_t = f(x_t, h_{t-1})$$

$$c = q(\{h_1, \dots, h_{T_x}\})$$

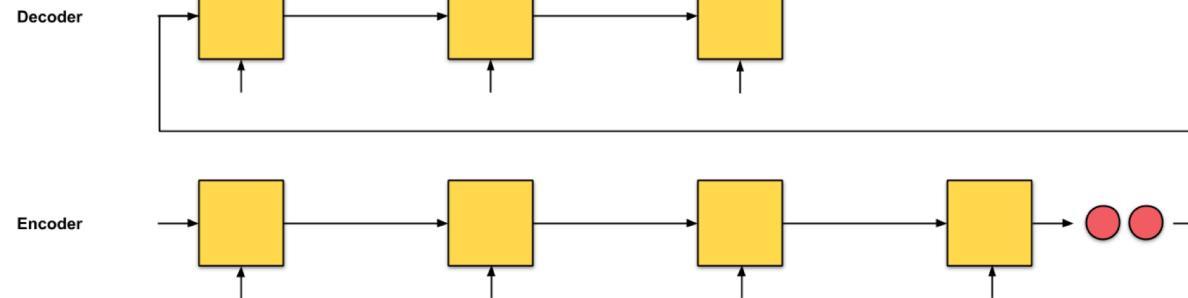
# Seq2Seq bottleneck

- 1、把輸入 $X$ 所有訊息壓縮到一個固定長度的hidden state  $C$ 。當輸入句子長度很長，模型的效能急劇下降。
- 2、把輸入 $X$  encoder成一個固定的長度，對於句子中每個詞都賦予相同的權重，這樣做是不合理。

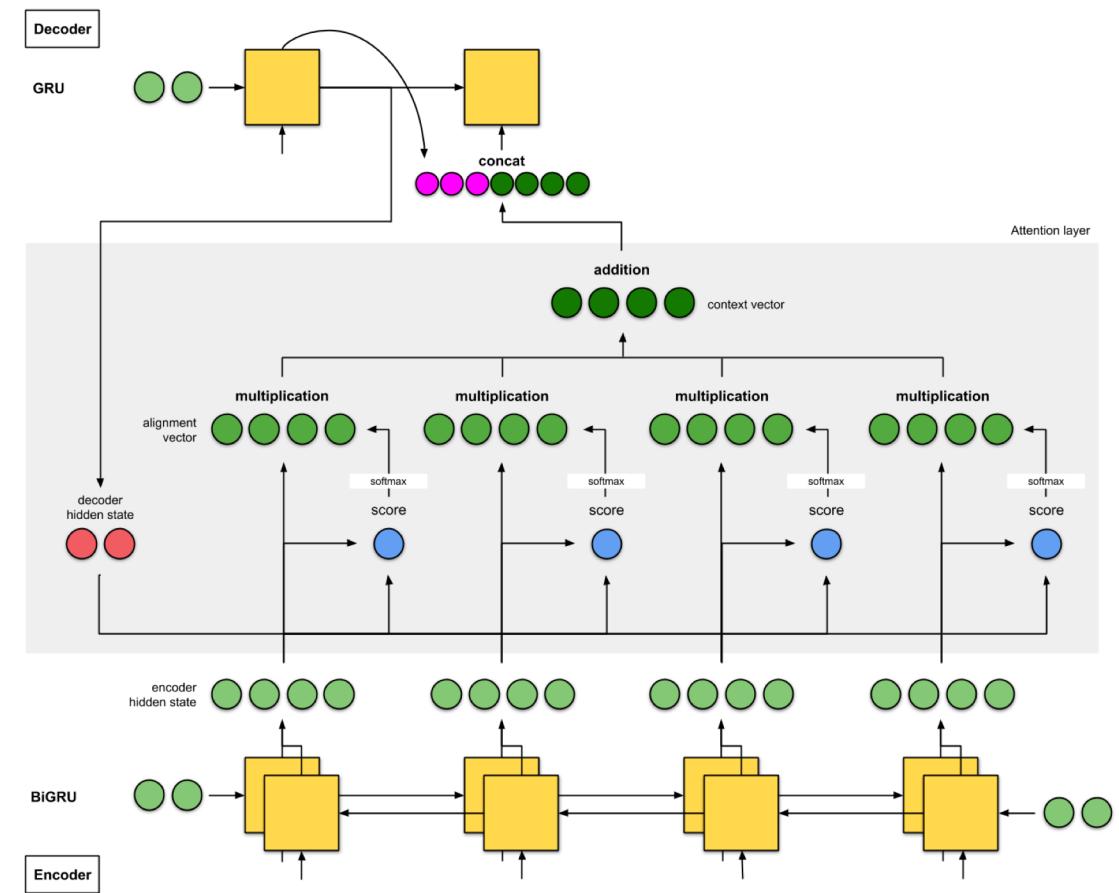


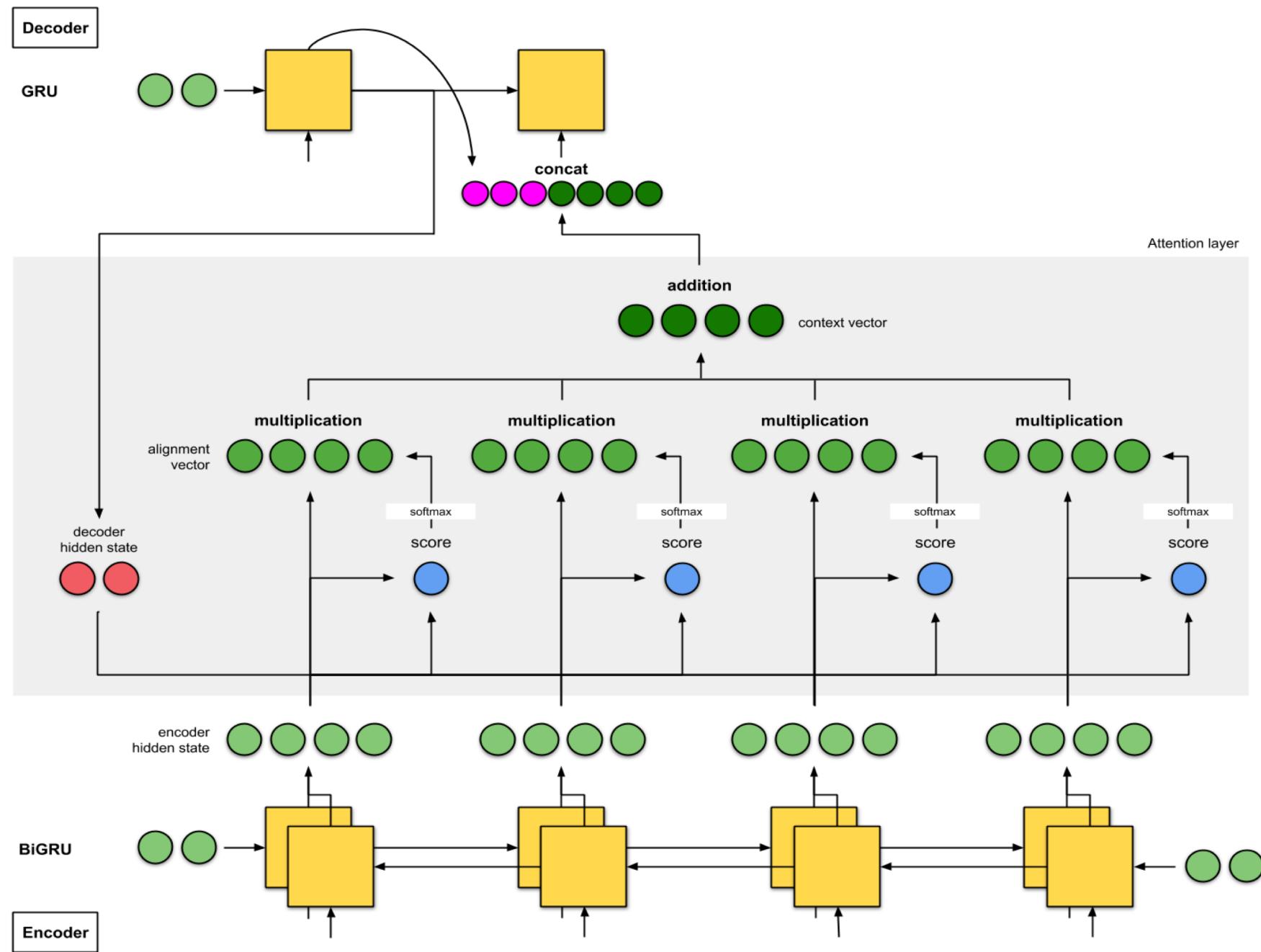
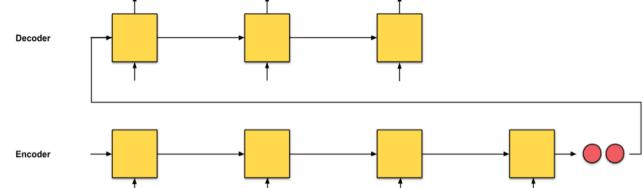
# Attention

## Seq2seq

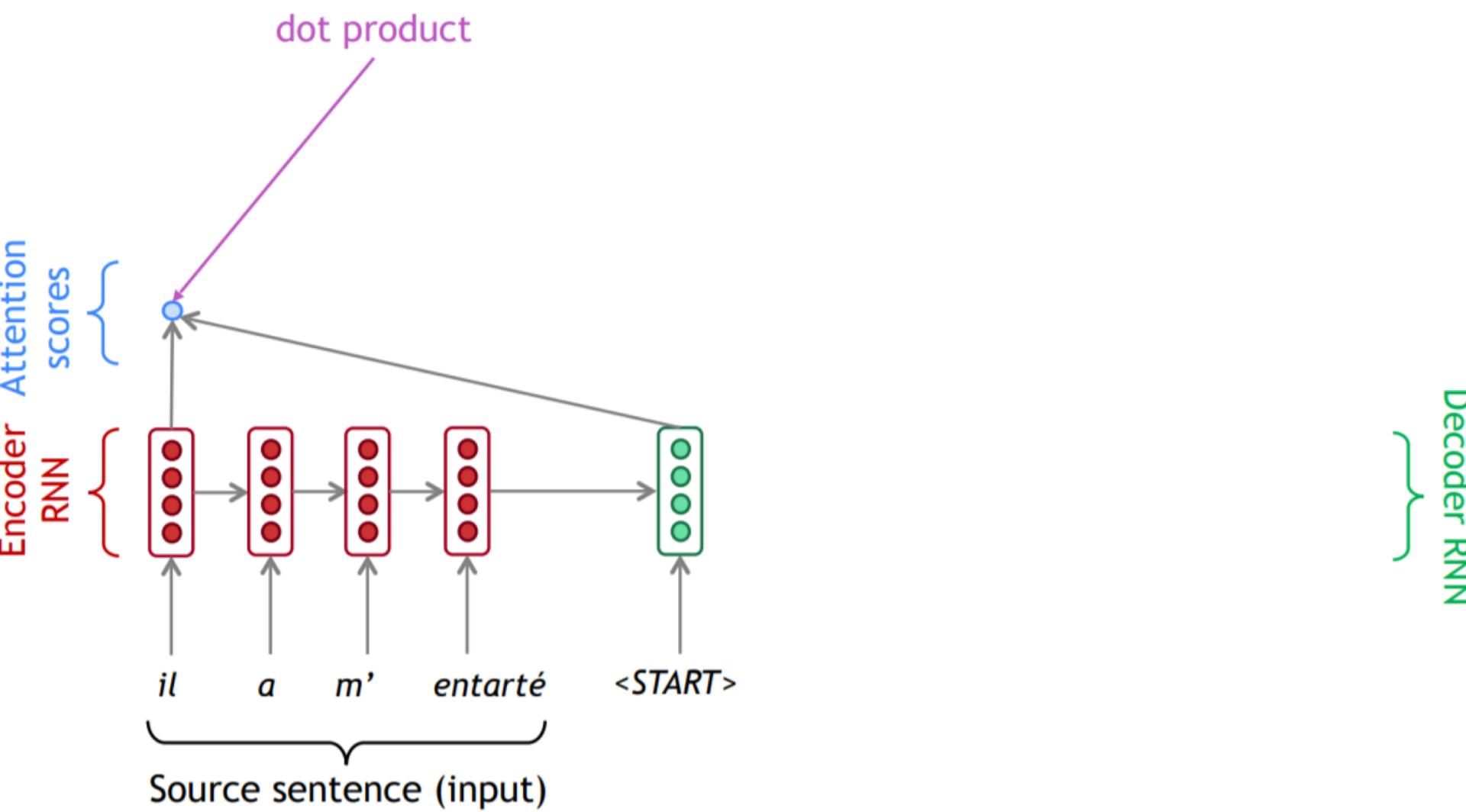


## Seq2seq+Attention

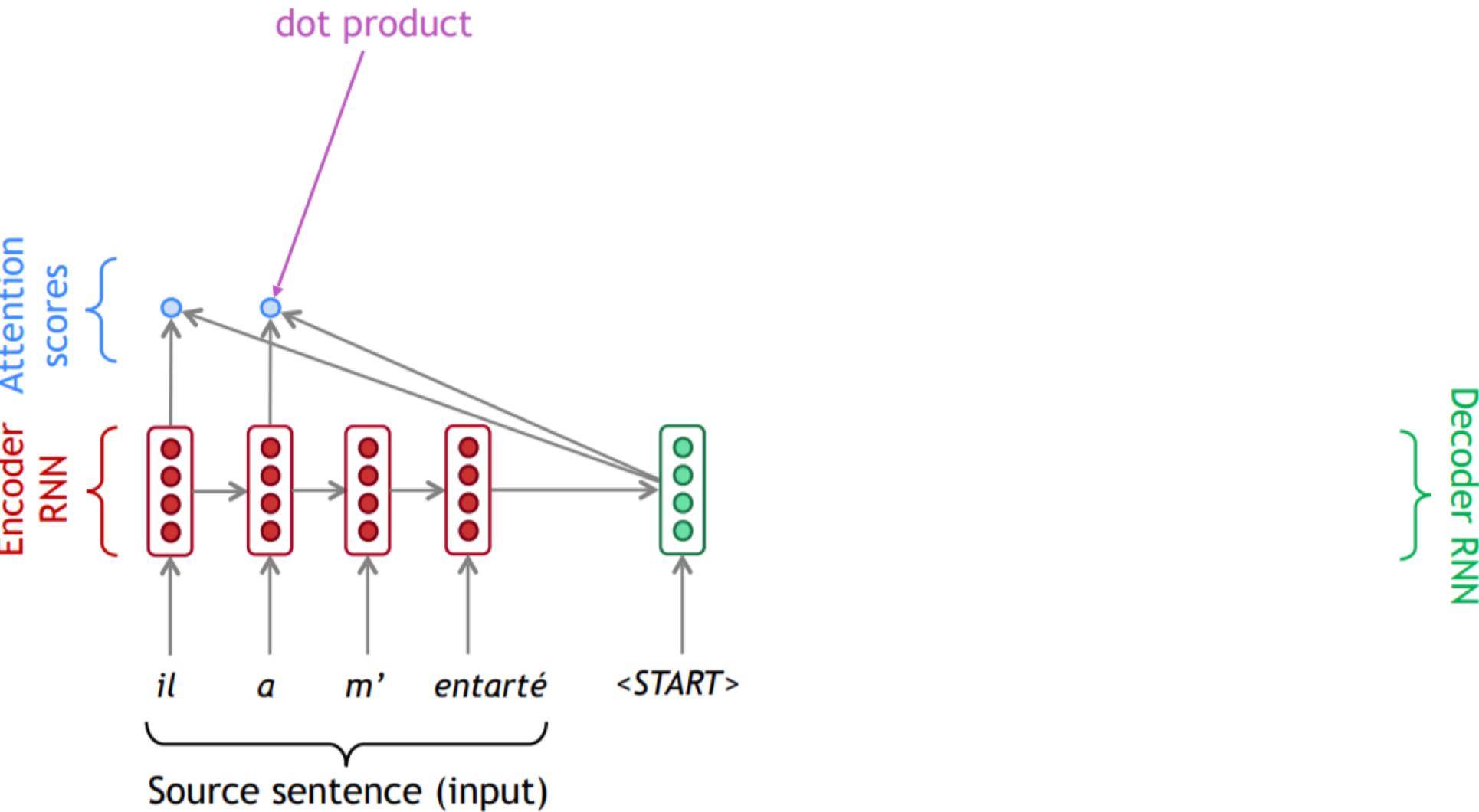




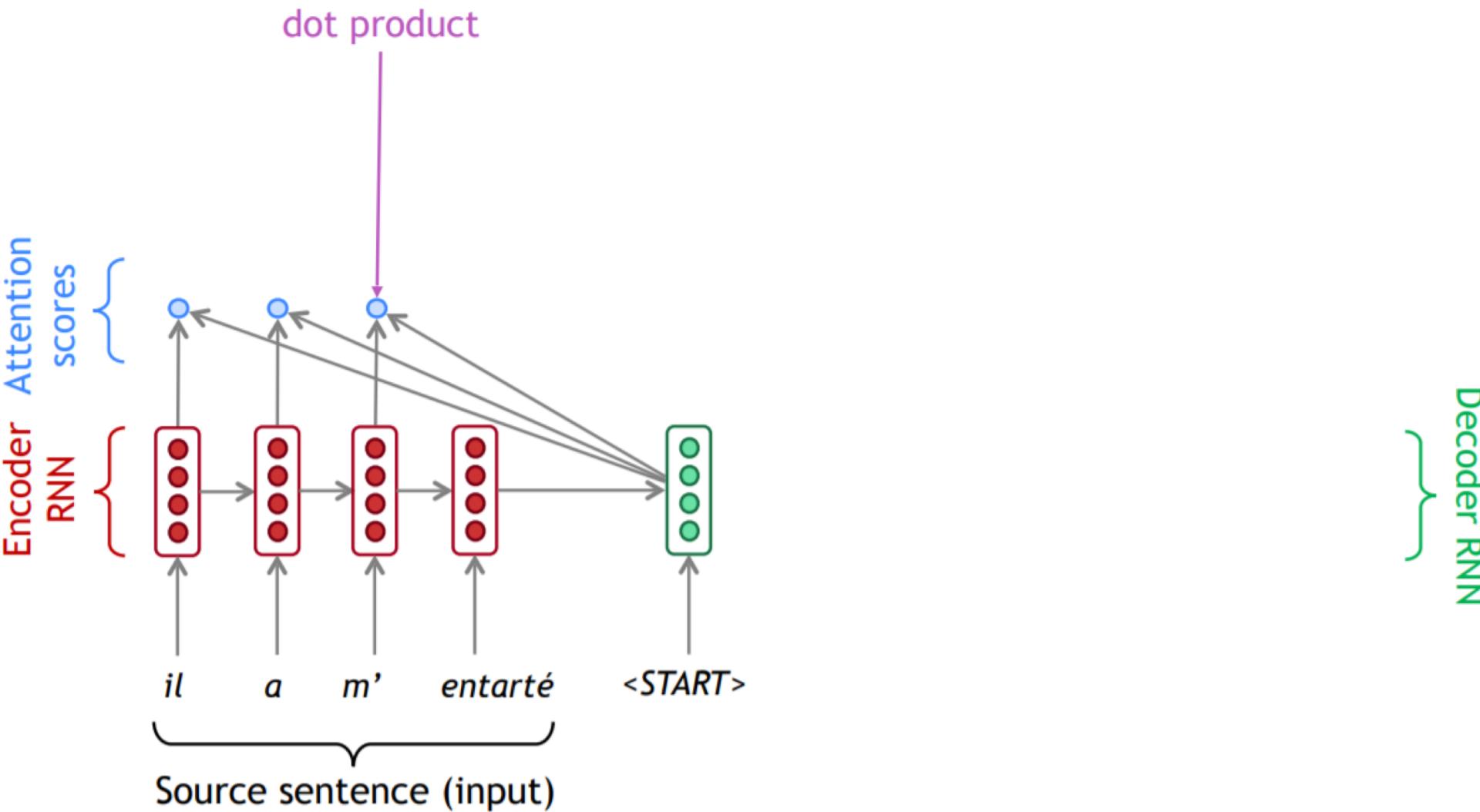
# Attention



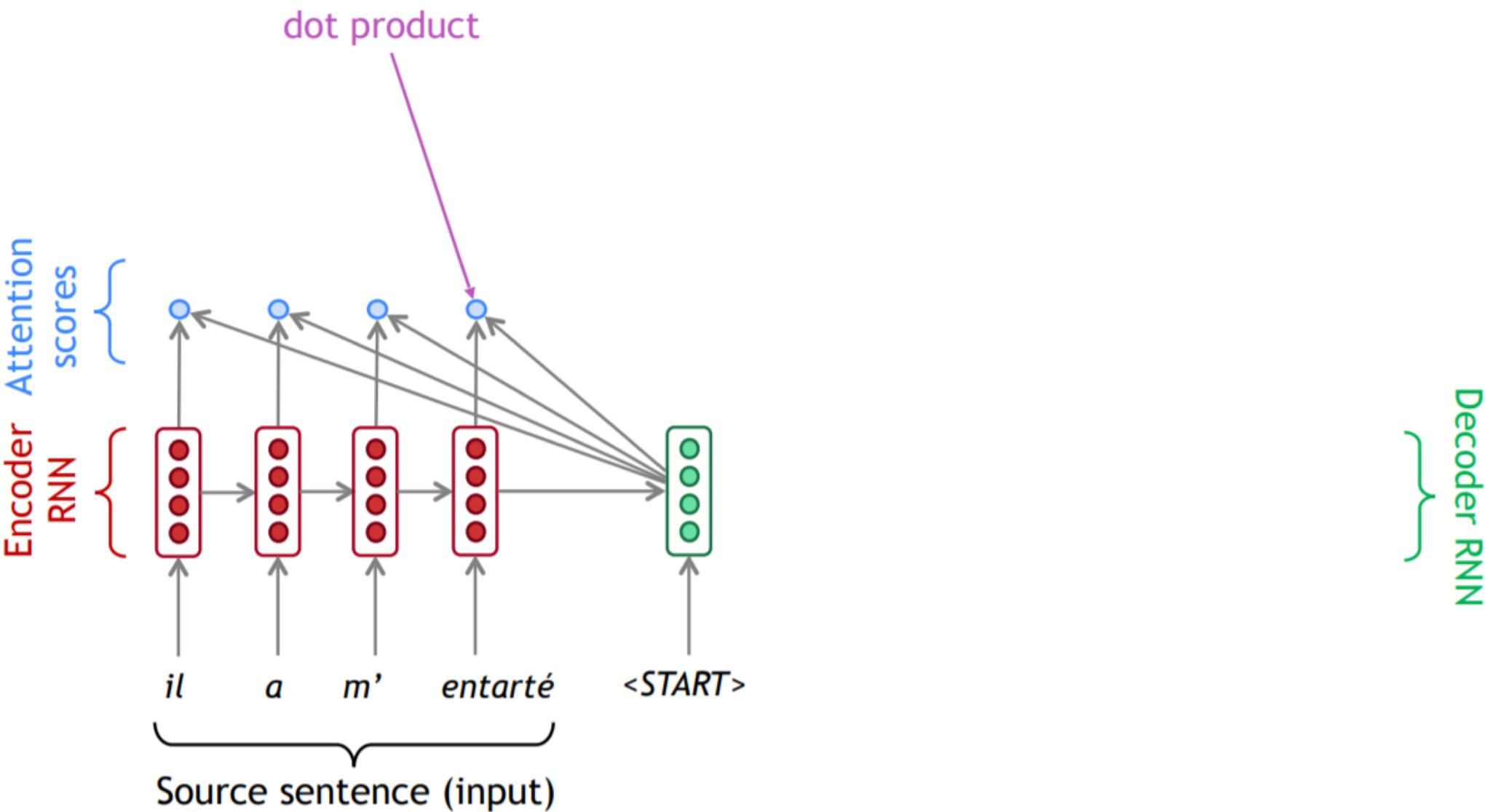
# Attention



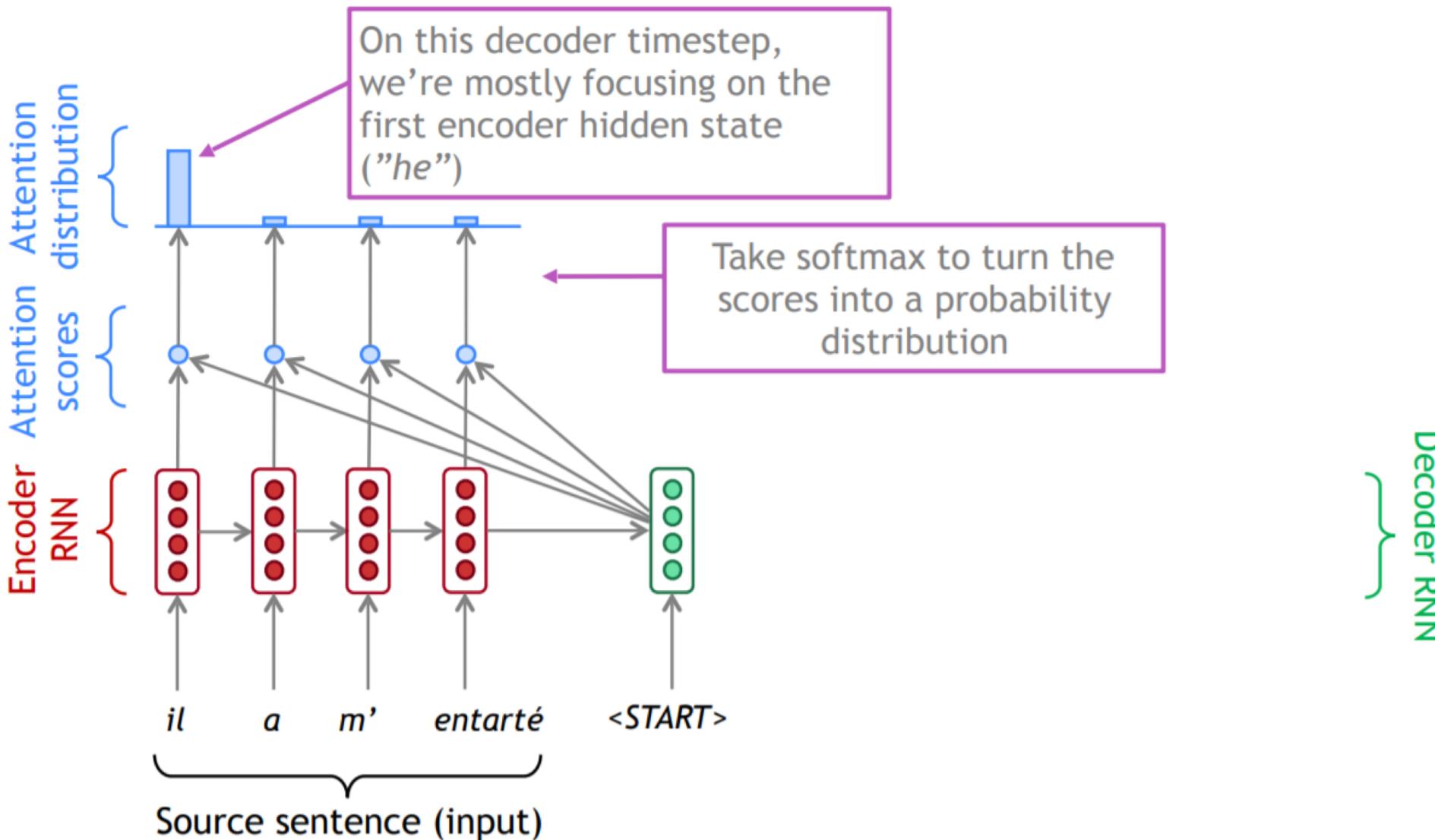
# Attention



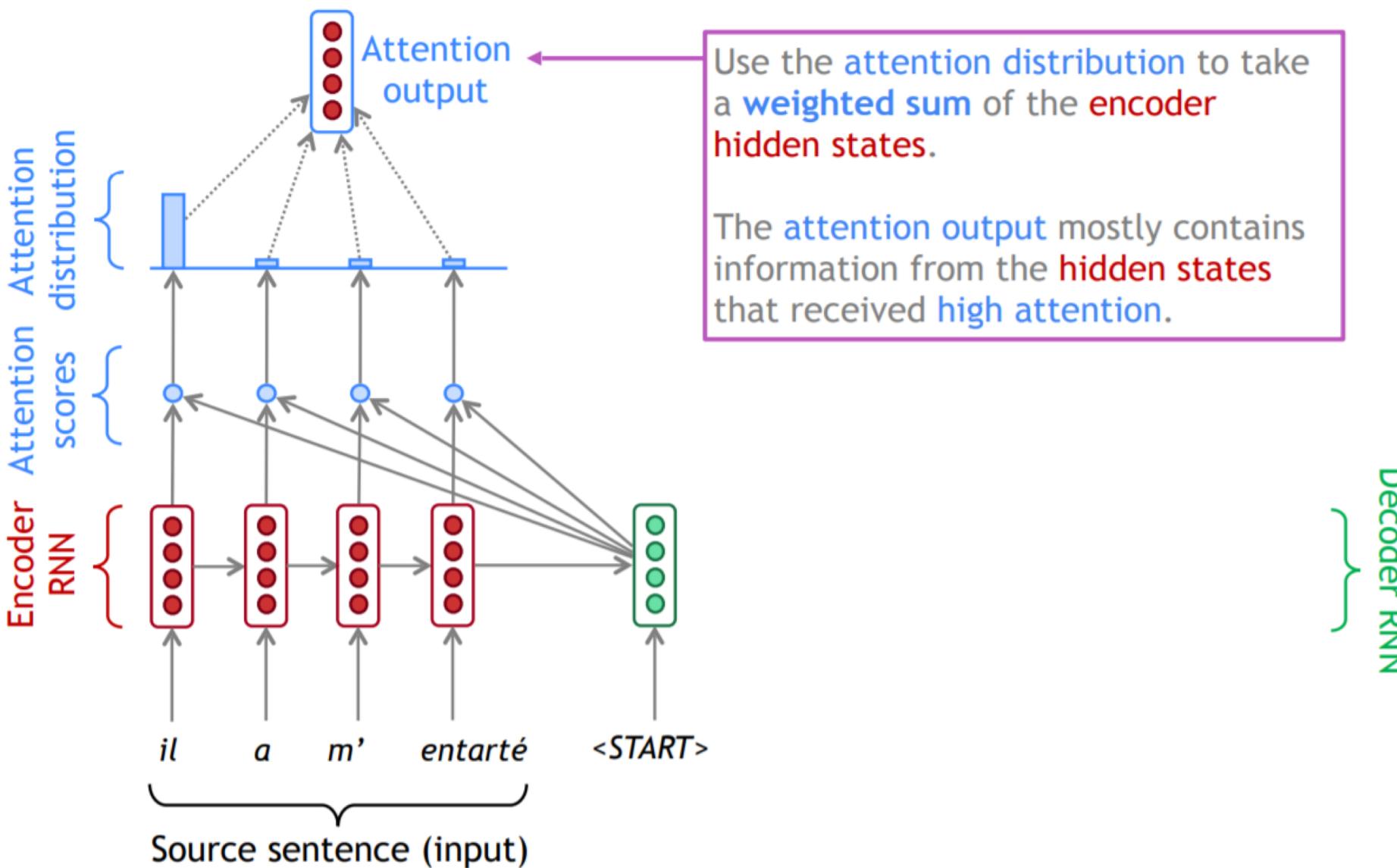
# Attention



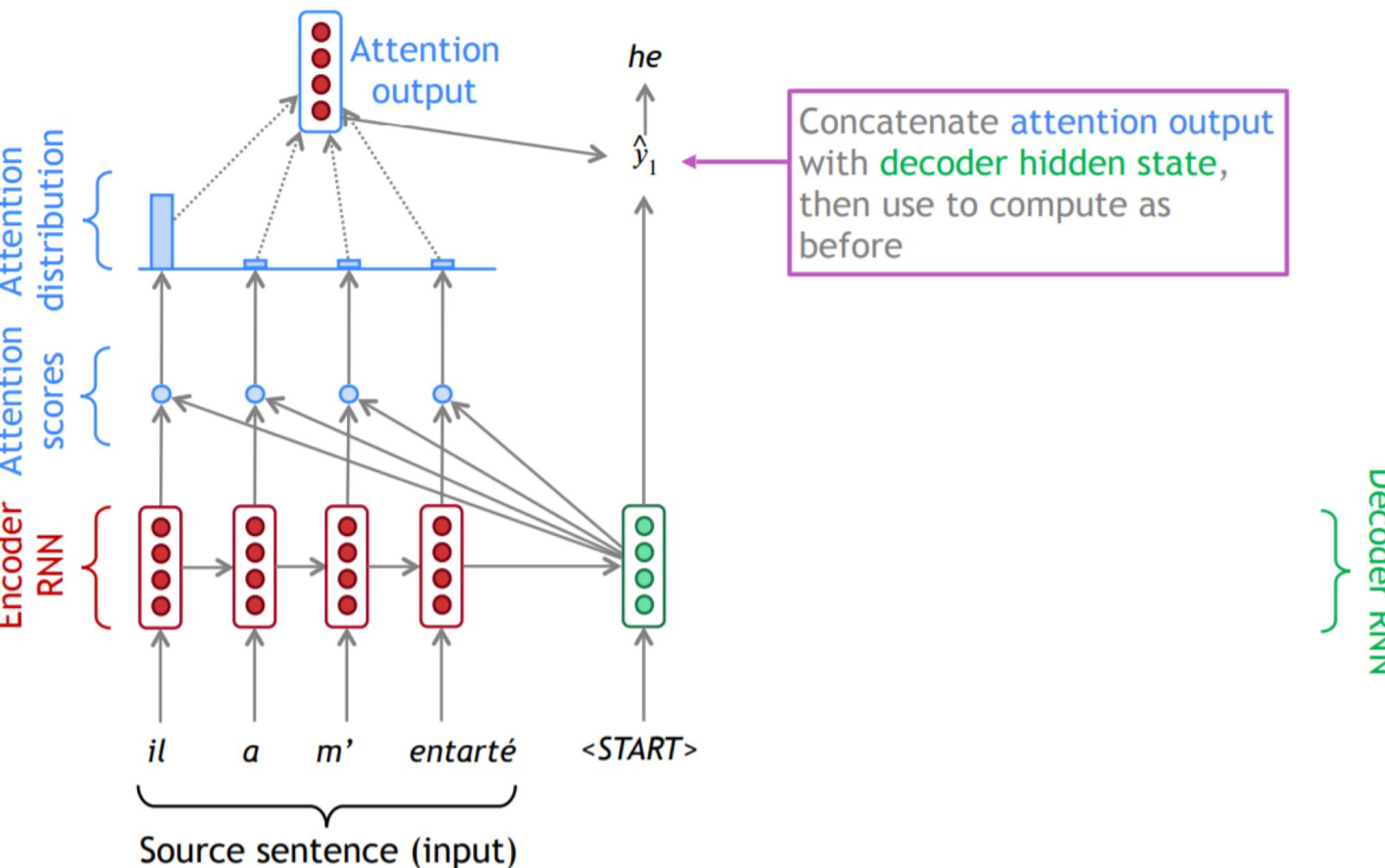
# Attention



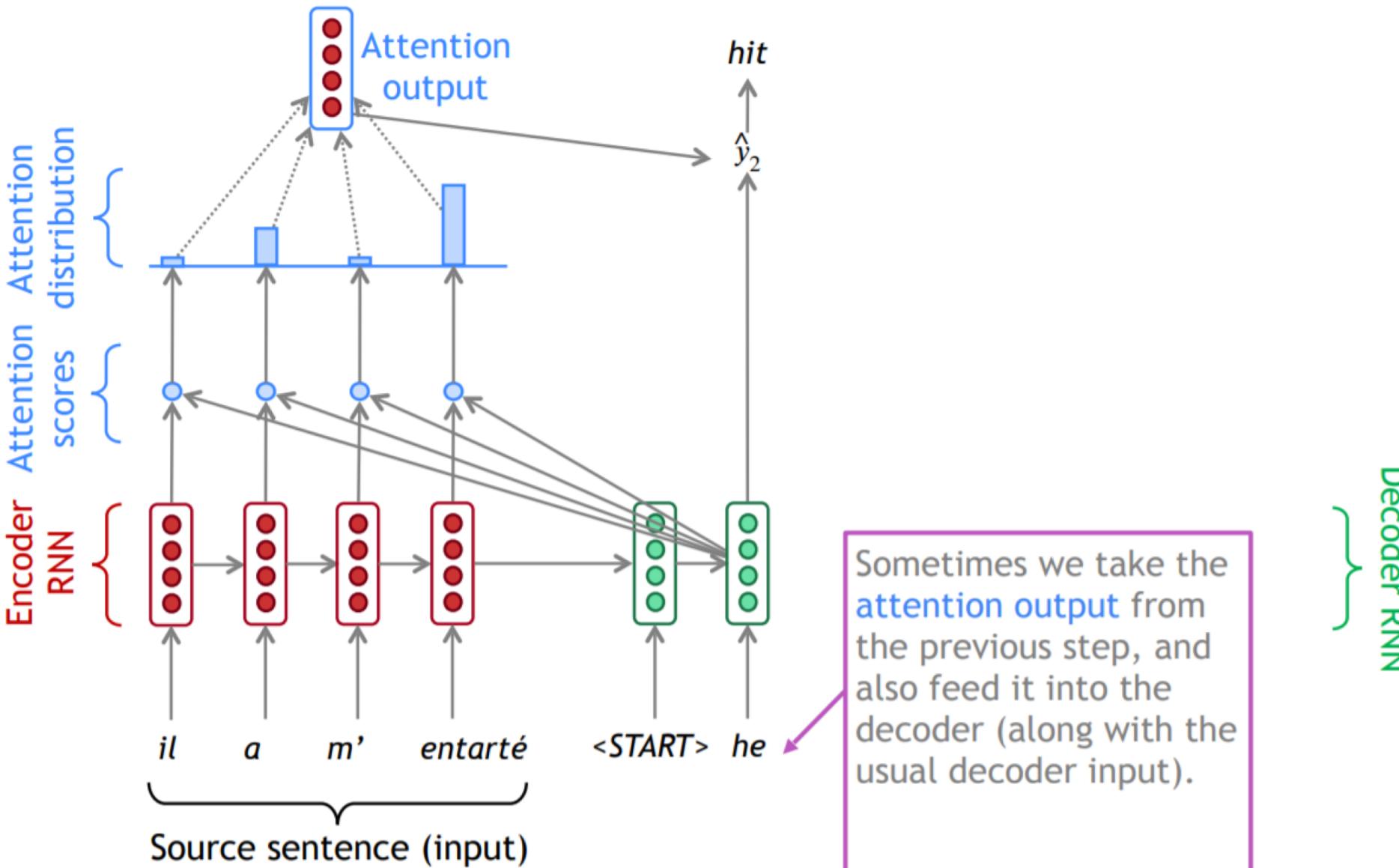
# Attention



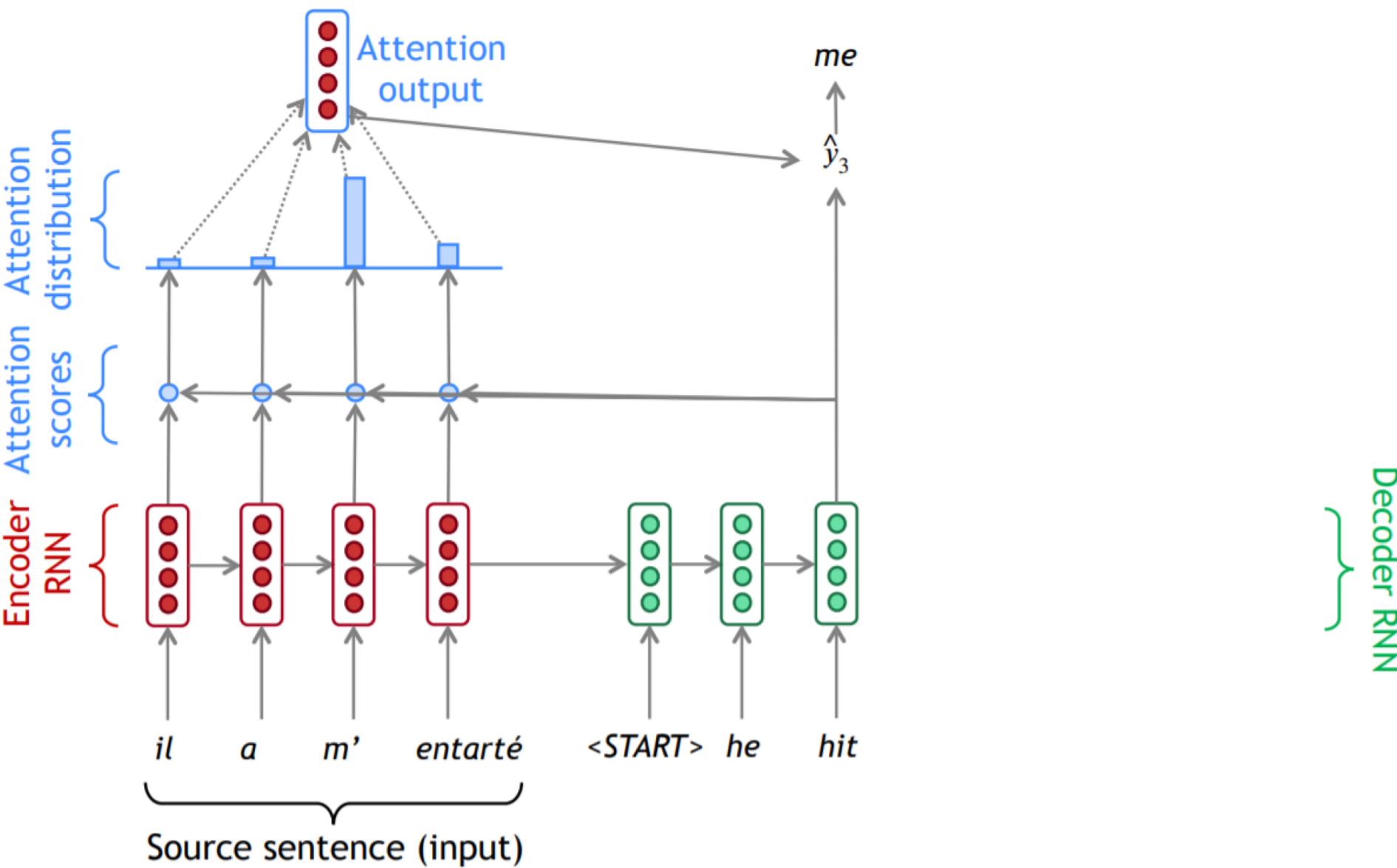
# Attention



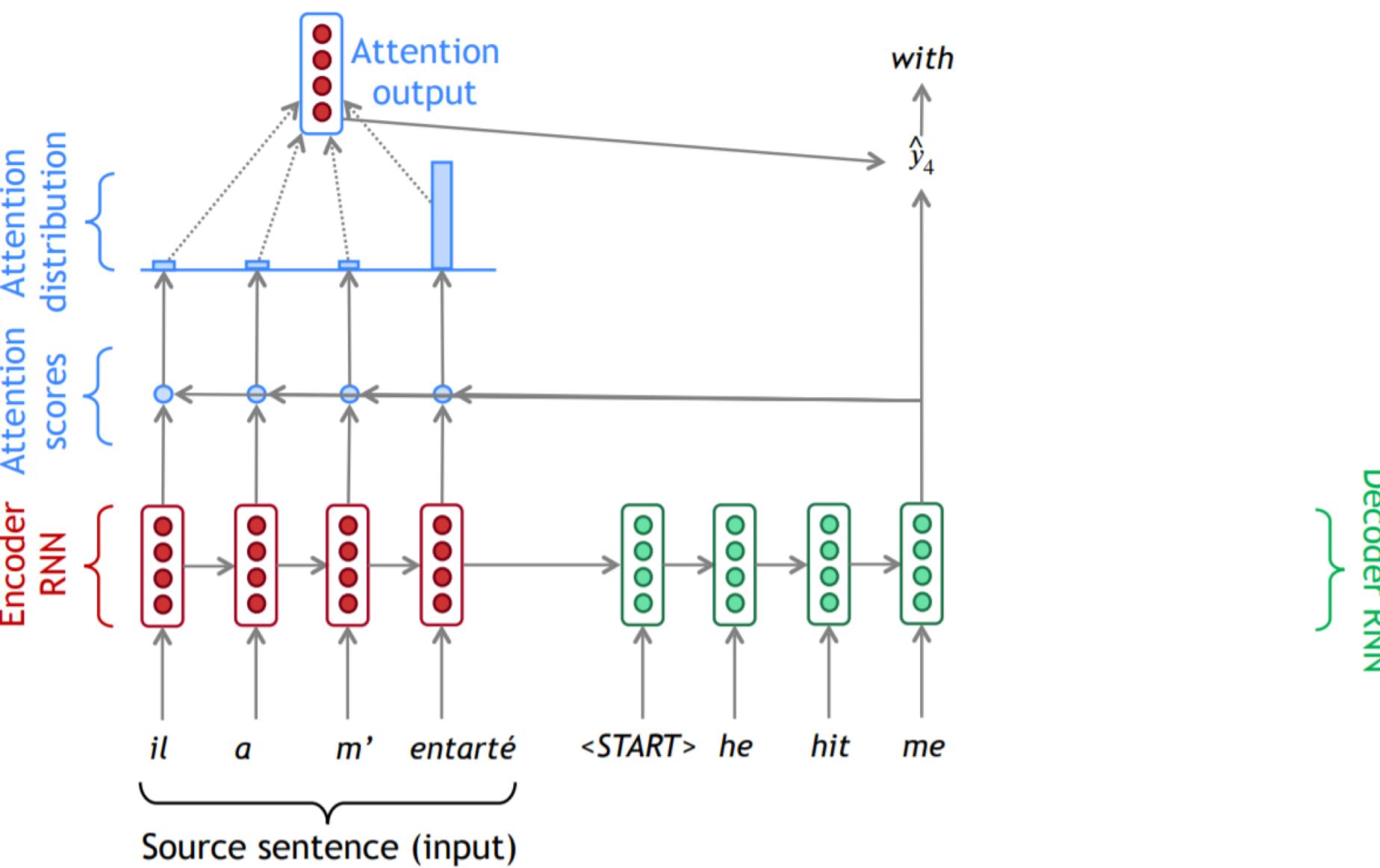
# Attention



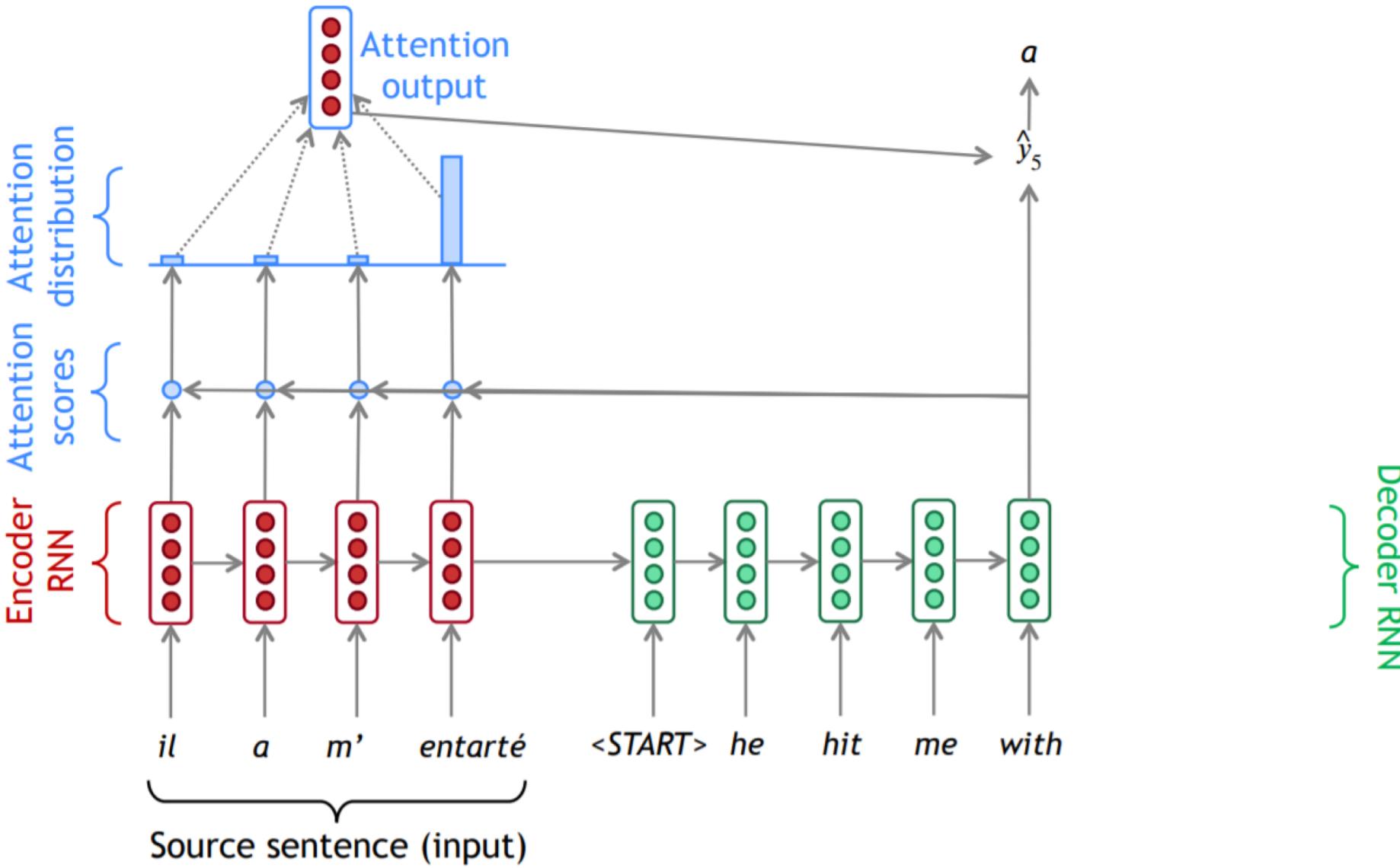
# Attention



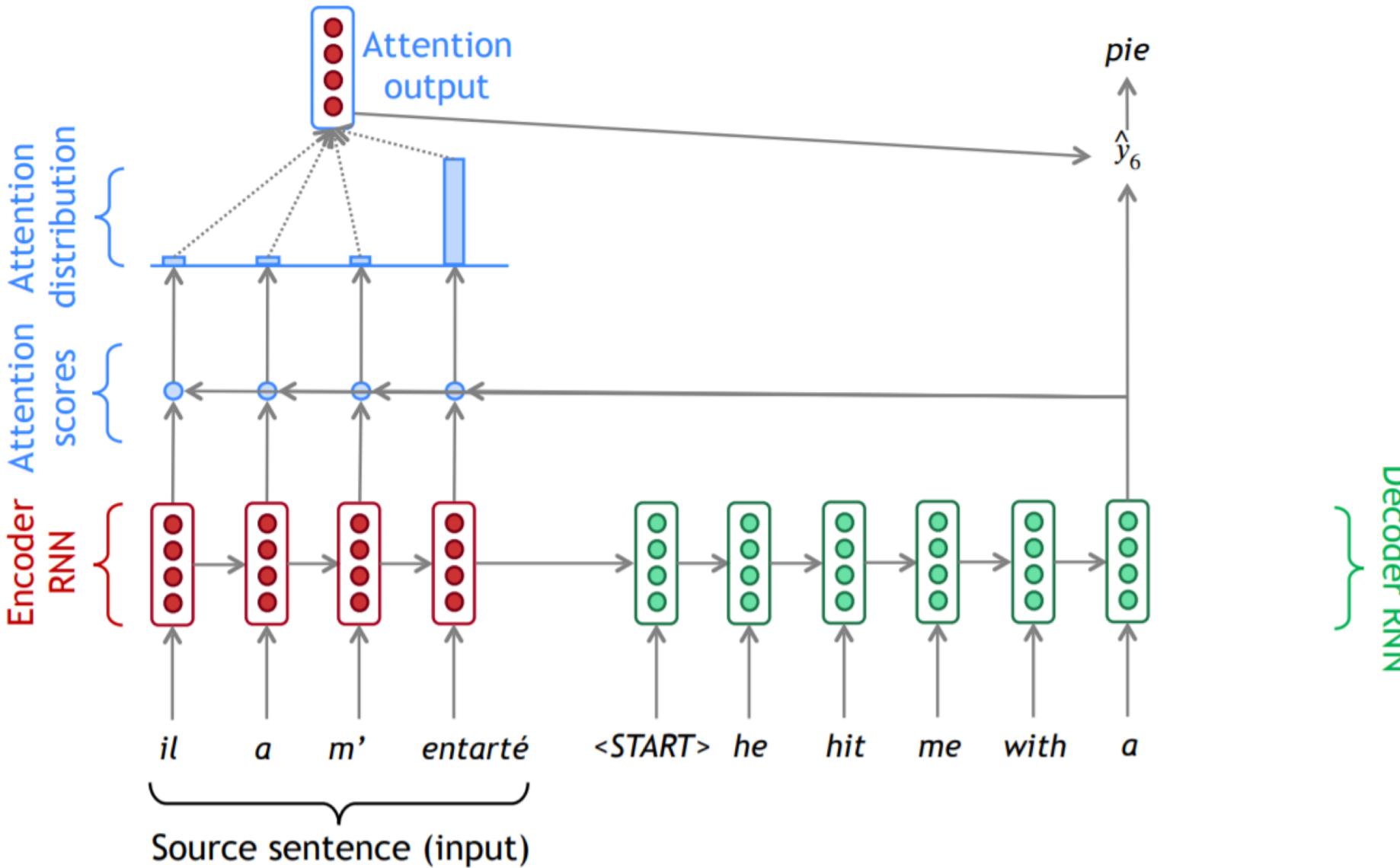
# Attention



# Attention



# Attention



# 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)

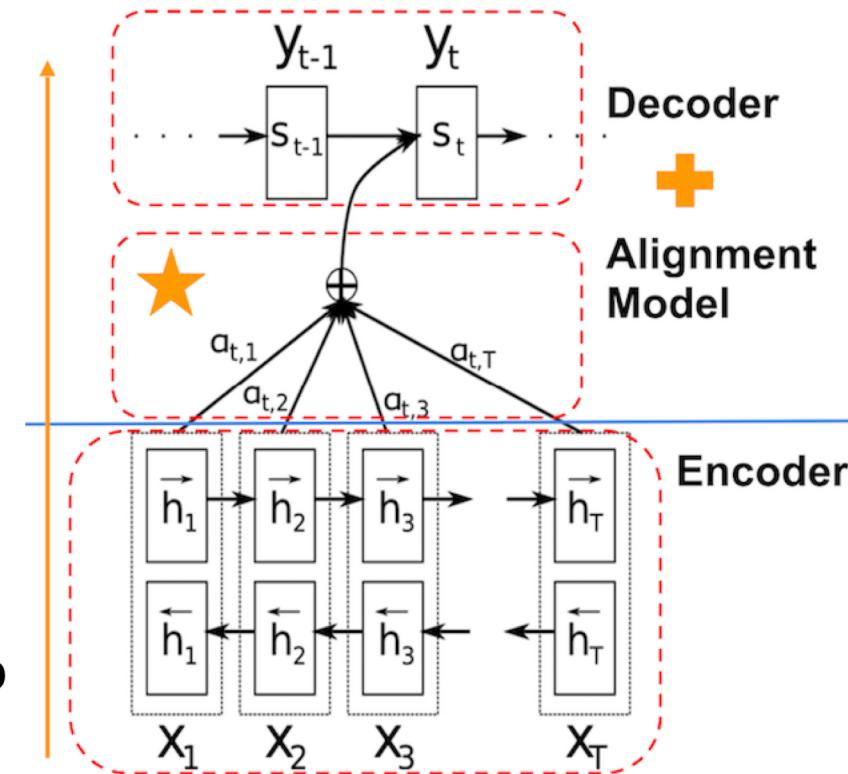
$$\alpha^t = \text{softmax}(e^t) \in \mathbb{R}^N$$

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

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

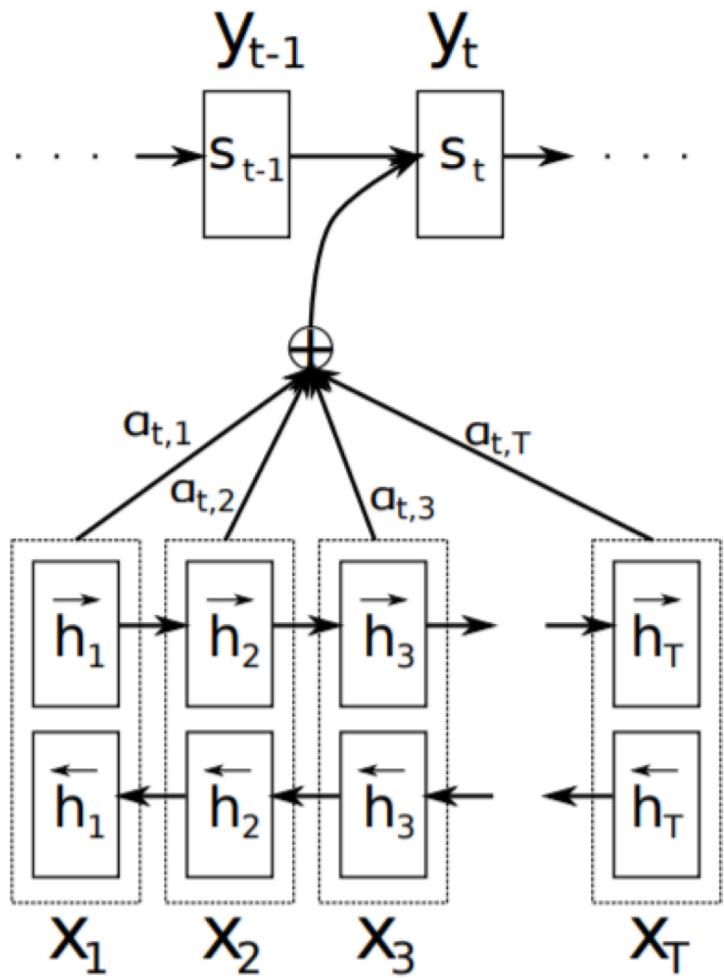
- Finally we concatenate the attention output  $a_t$  with the decoder hidden state  $s_t$  and proceed as in the non-attention seq2seq model

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

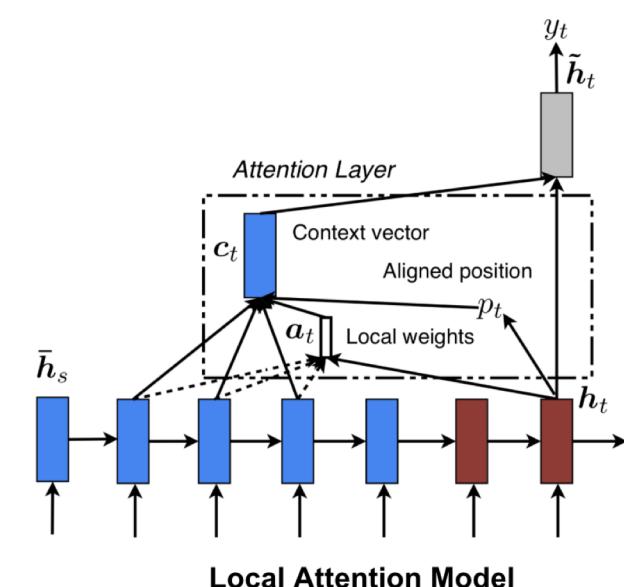
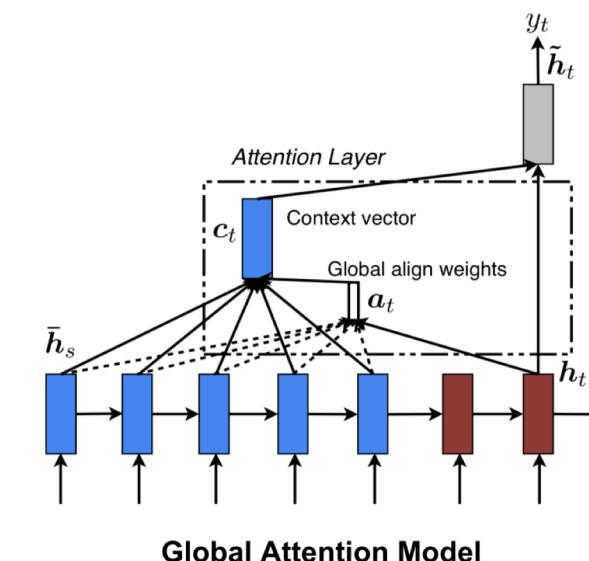


# Attention

## Bahdanau Attention

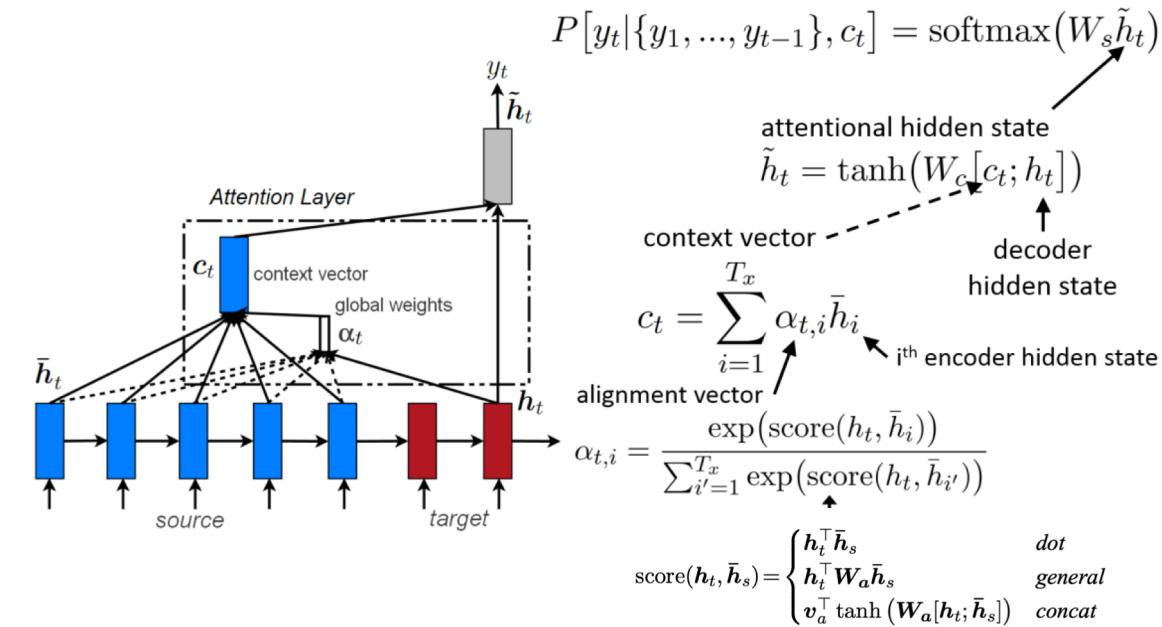


## Luong Attention

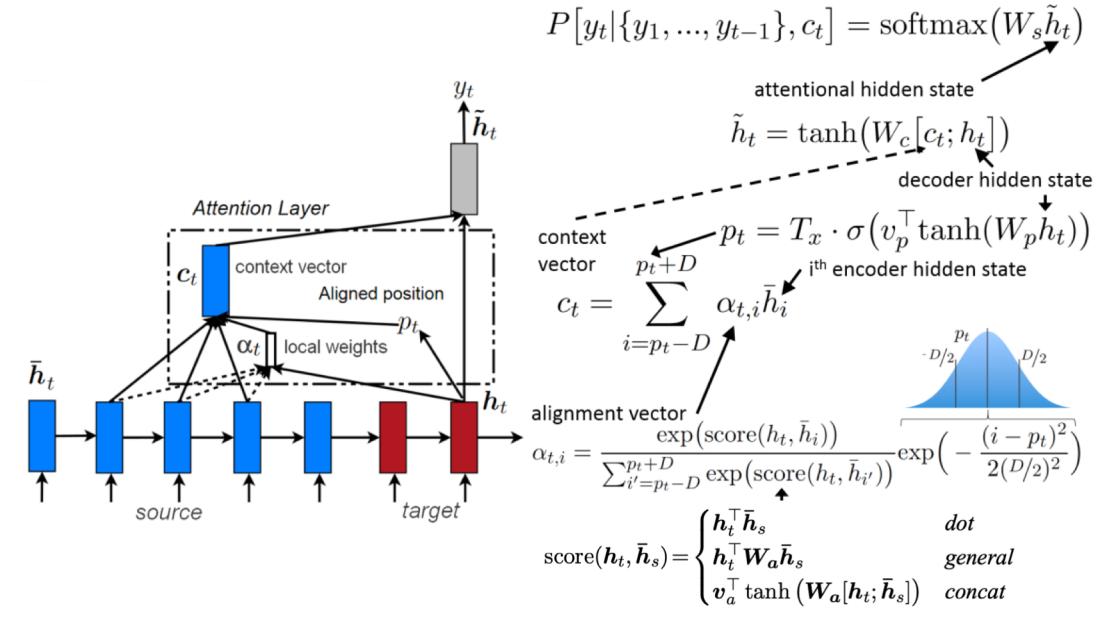


# Luong Attention

## Global Attention



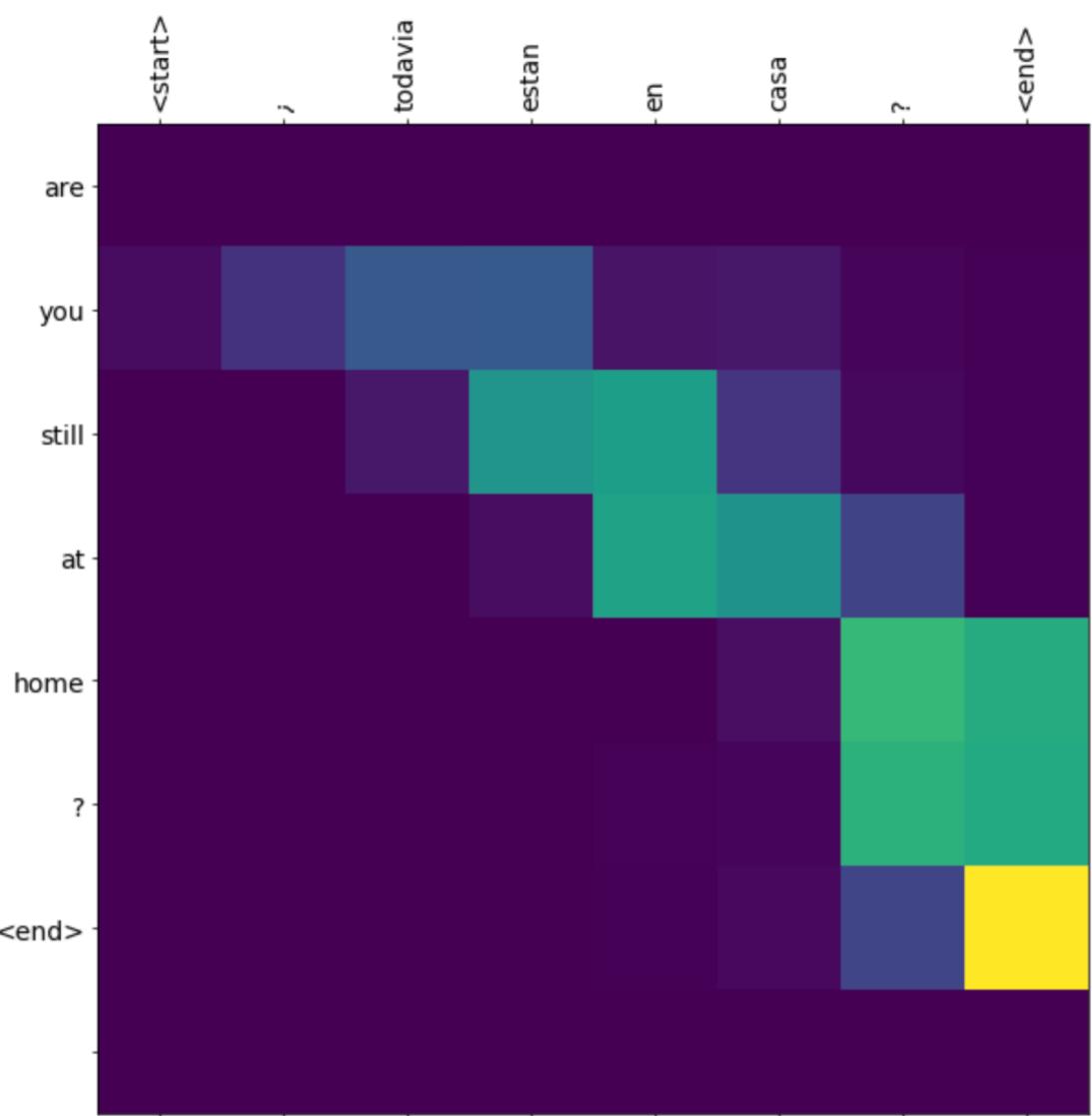
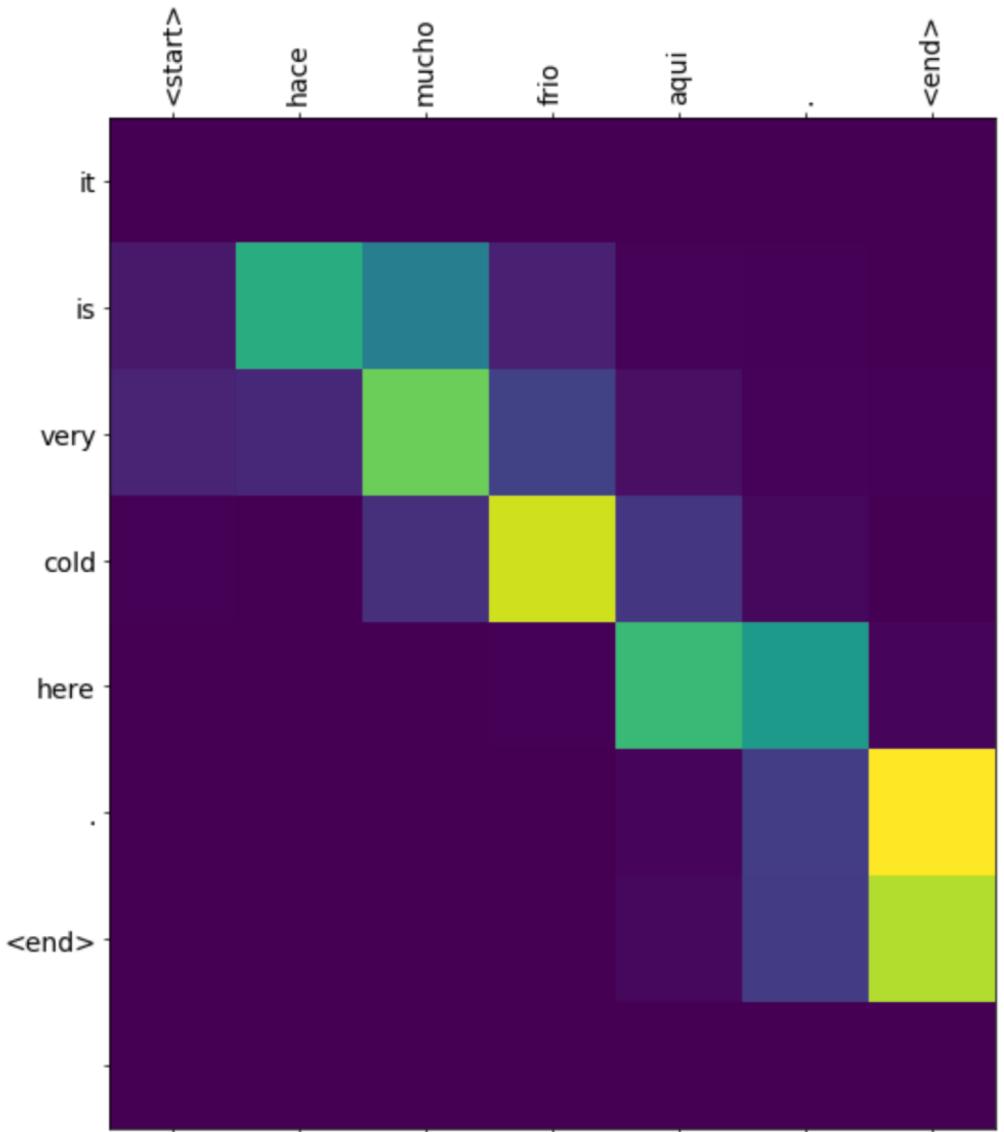
## local Attention



Monotonic alignment (local-m):  $p_t = t$

Predictive alignment (local-p):  $p_t = S \cdot \text{sigmoid}(\mathbf{v}_p^\top \tanh(\mathbf{W}_p \mathbf{h}_t))$

# Attention application



# Attention application



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



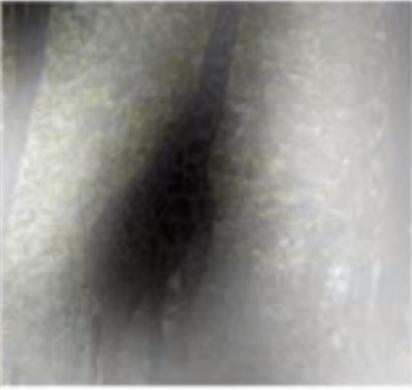
A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.



**To be continued**

---

**Transformer**

# References

---

- Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation
- Sequence to Sequence Learning with Neural Networks
- NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE
- Effective Approaches to Attention-based Neural Machine Translation
- Show, Attend and Tell: Neural Image Caption Generation with Visual Attention