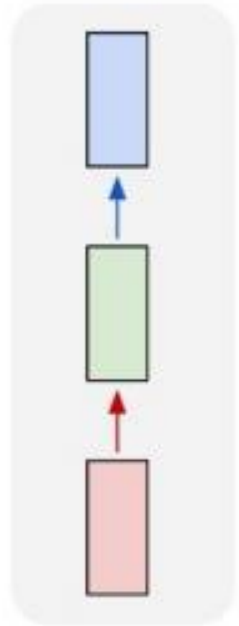


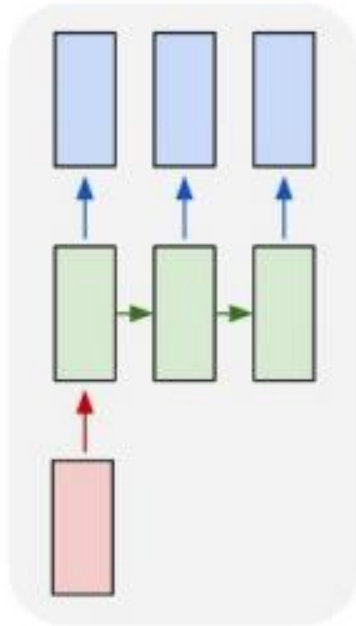
Types of RNNs

- many-to-one

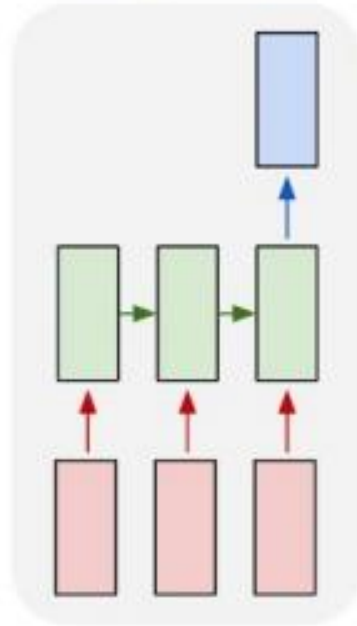
one to one



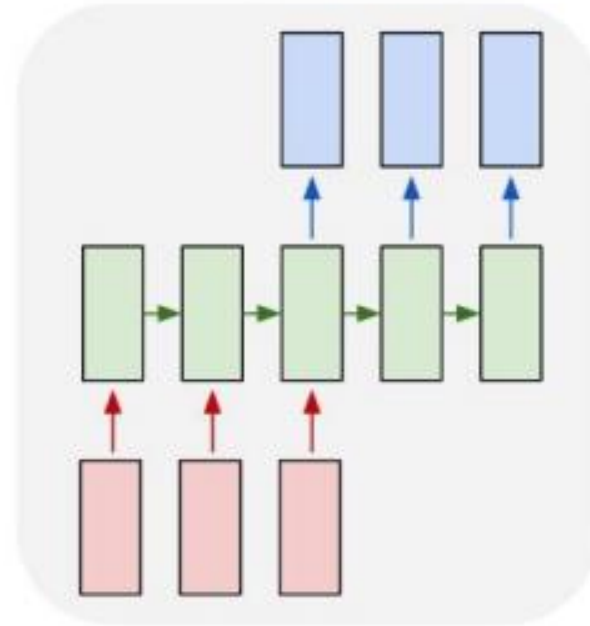
one to many



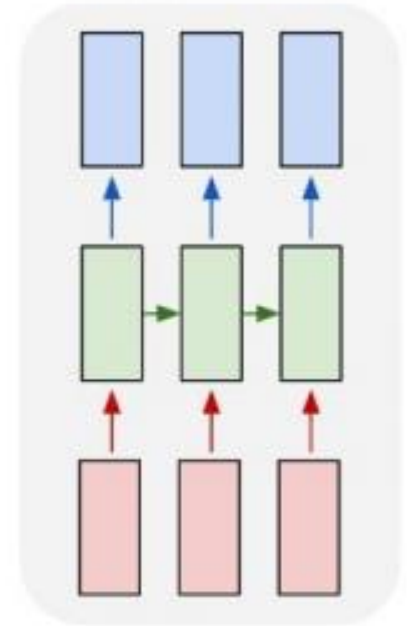
many to one



many to many



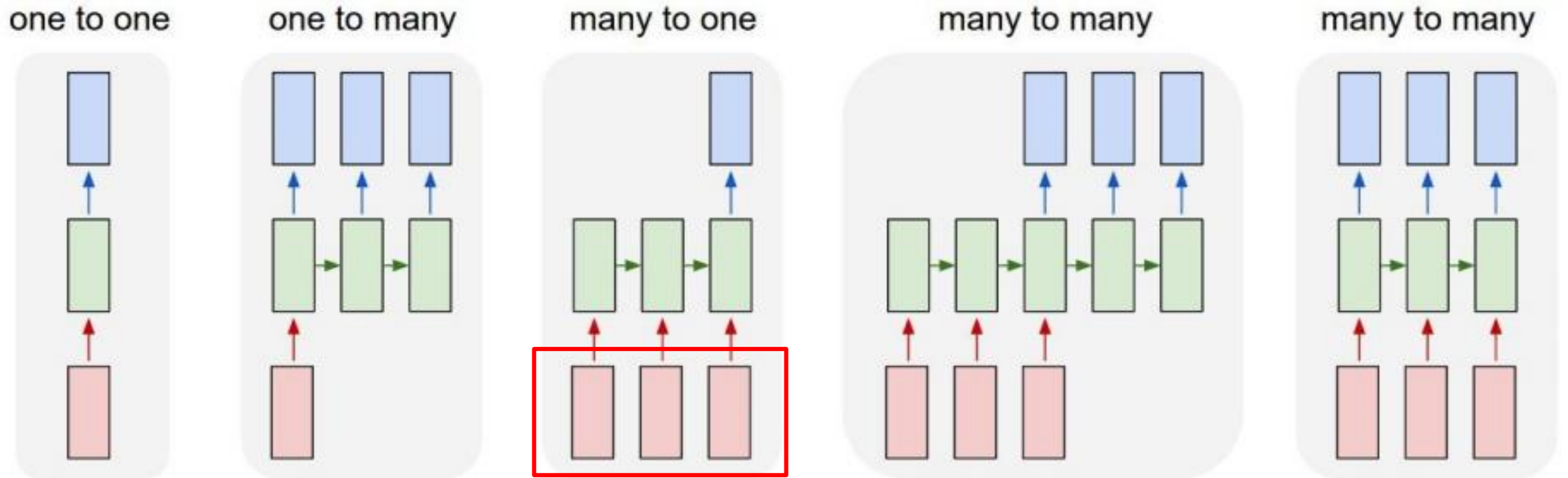
many to many



↖ e.g. **Sentiment Classification**
sequence of words -> sentiment

Types of RNNs

- many-to-one

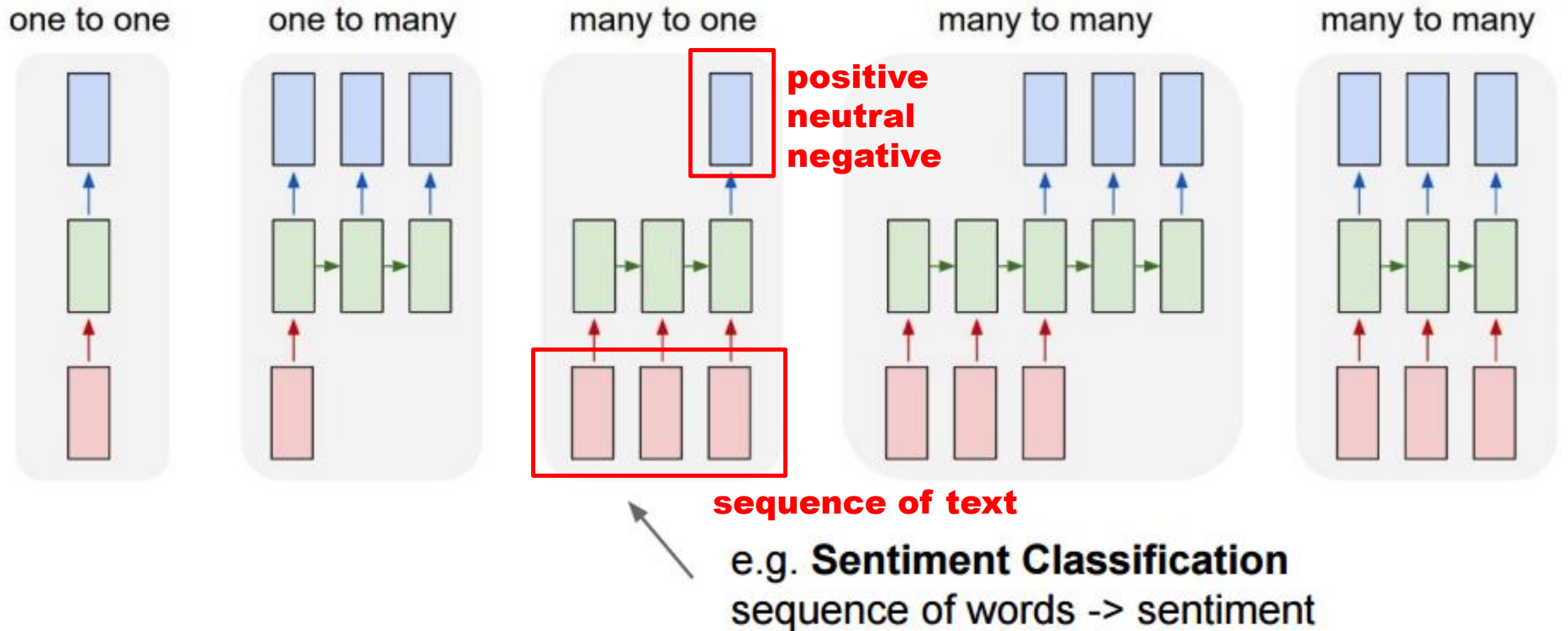


sequence of text

↖
e.g. Sentiment Classification
sequence of words -> sentiment

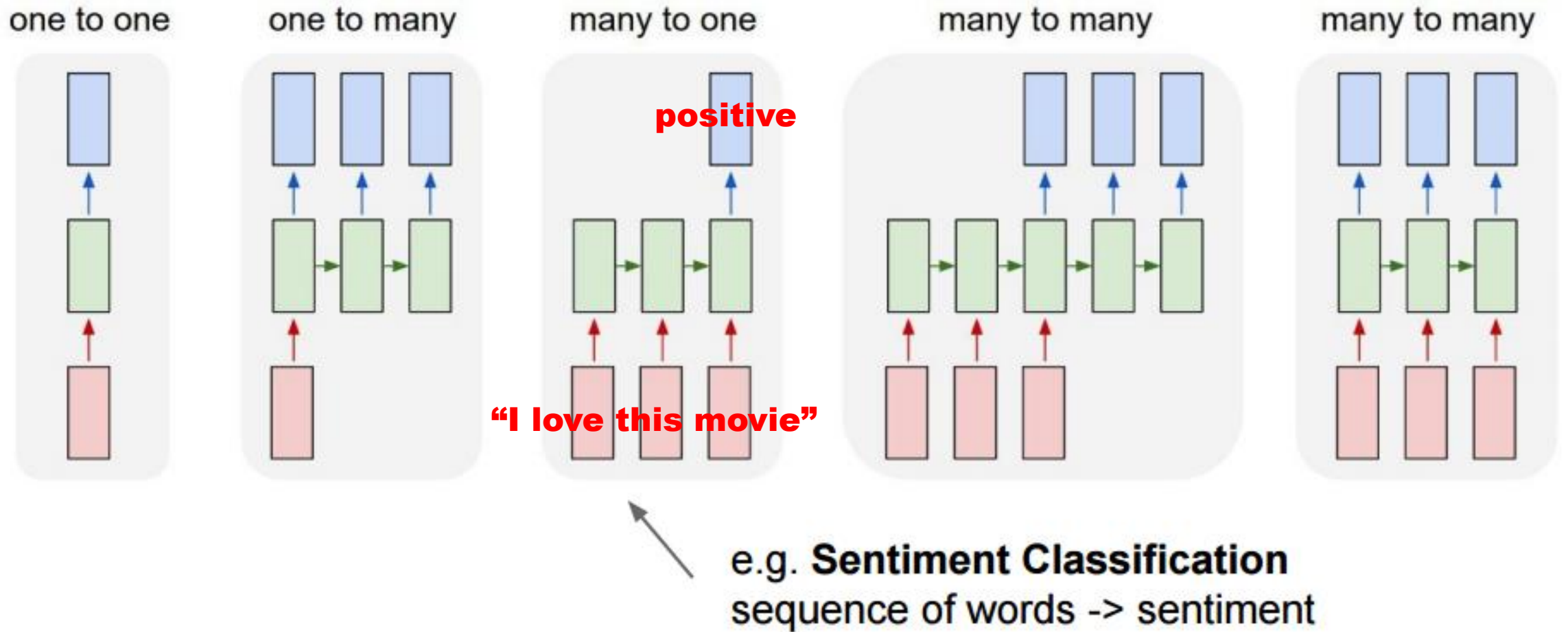
Types of RNNs

- many-to-one



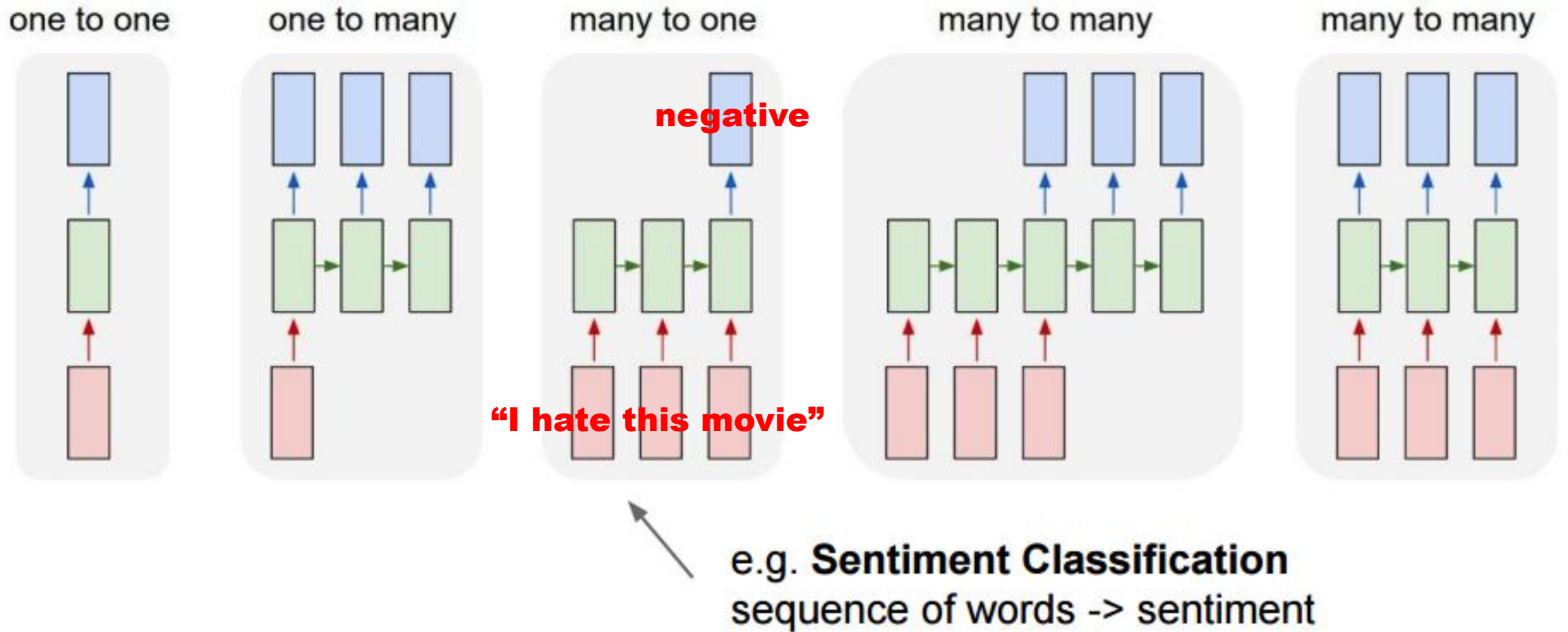
Types of RNNs

- many-to-one



Types of RNNs

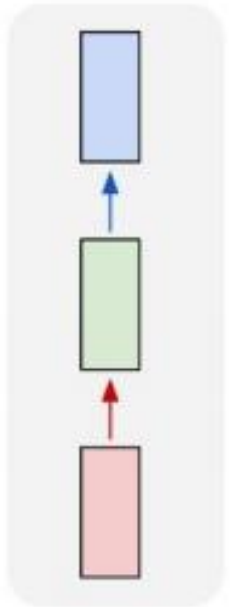
- many-to-one



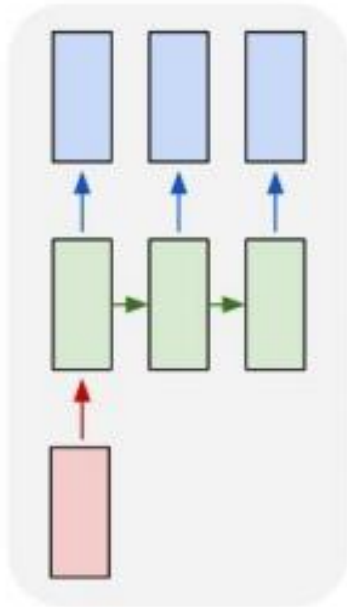
Types of RNNs

- Sequence-to sequence

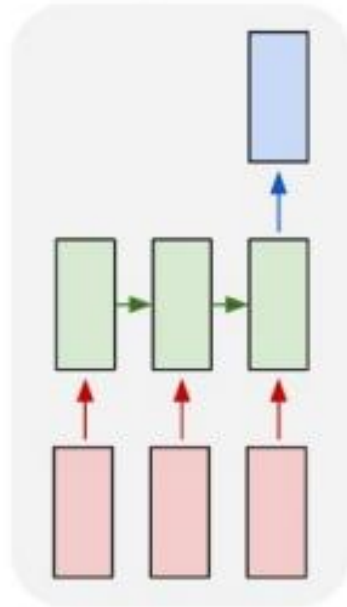
one to one



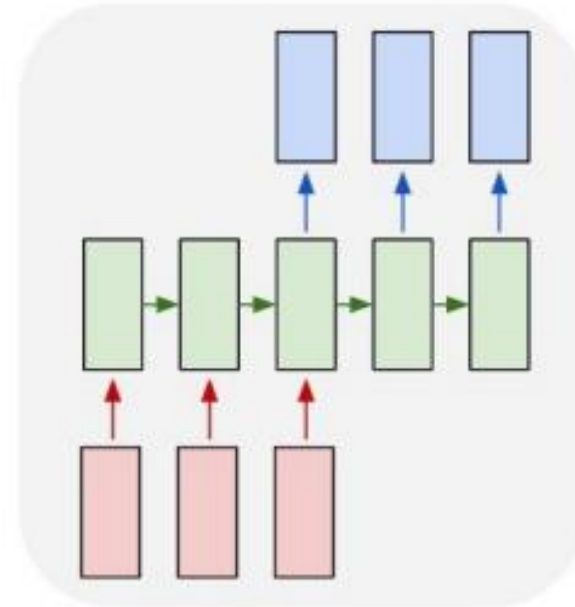
one to many



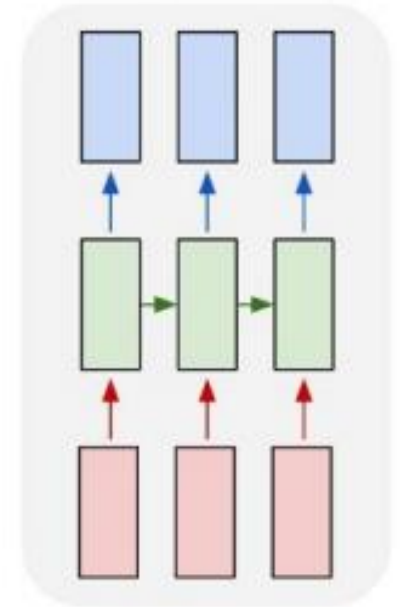
many to one



many to many



many to many

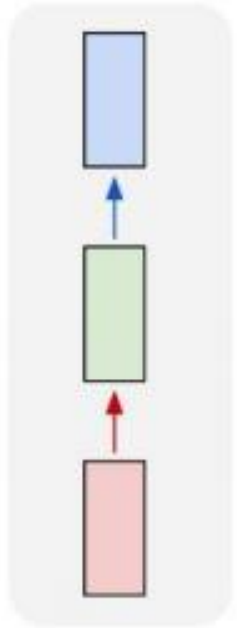


↖ e.g. **Machine Translation**
seq of words -> seq of words

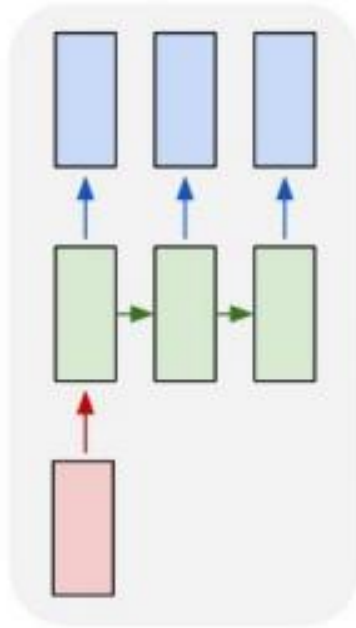
Types of RNNs

- Sequence-to sequence

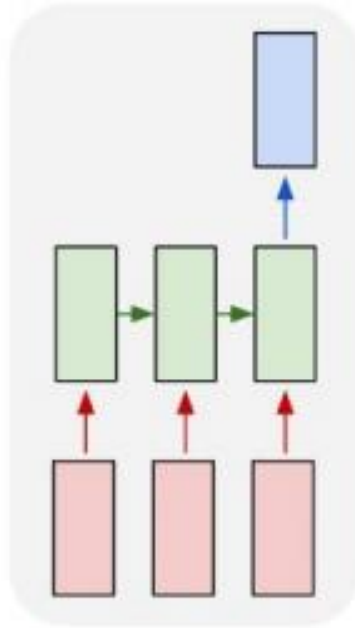
one to one



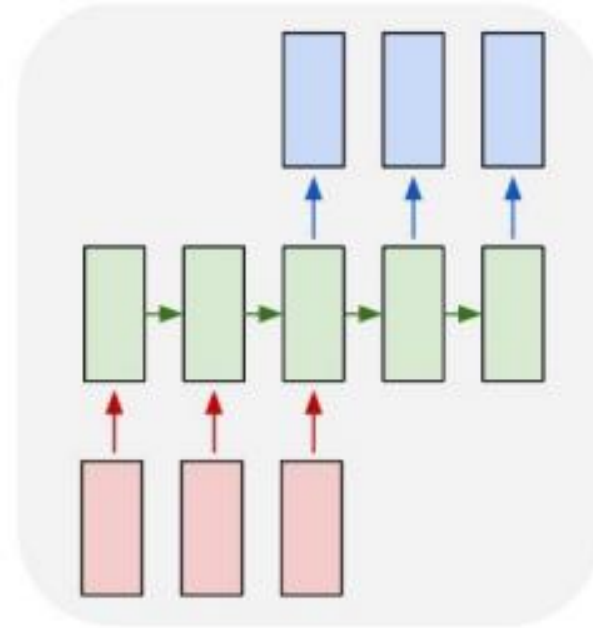
one to many



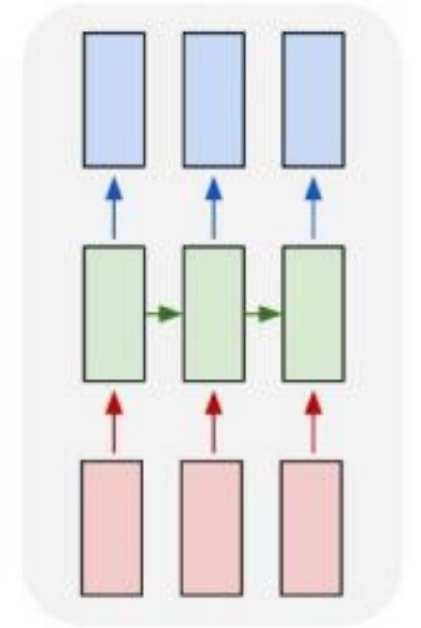
many to one



many to many



many to many

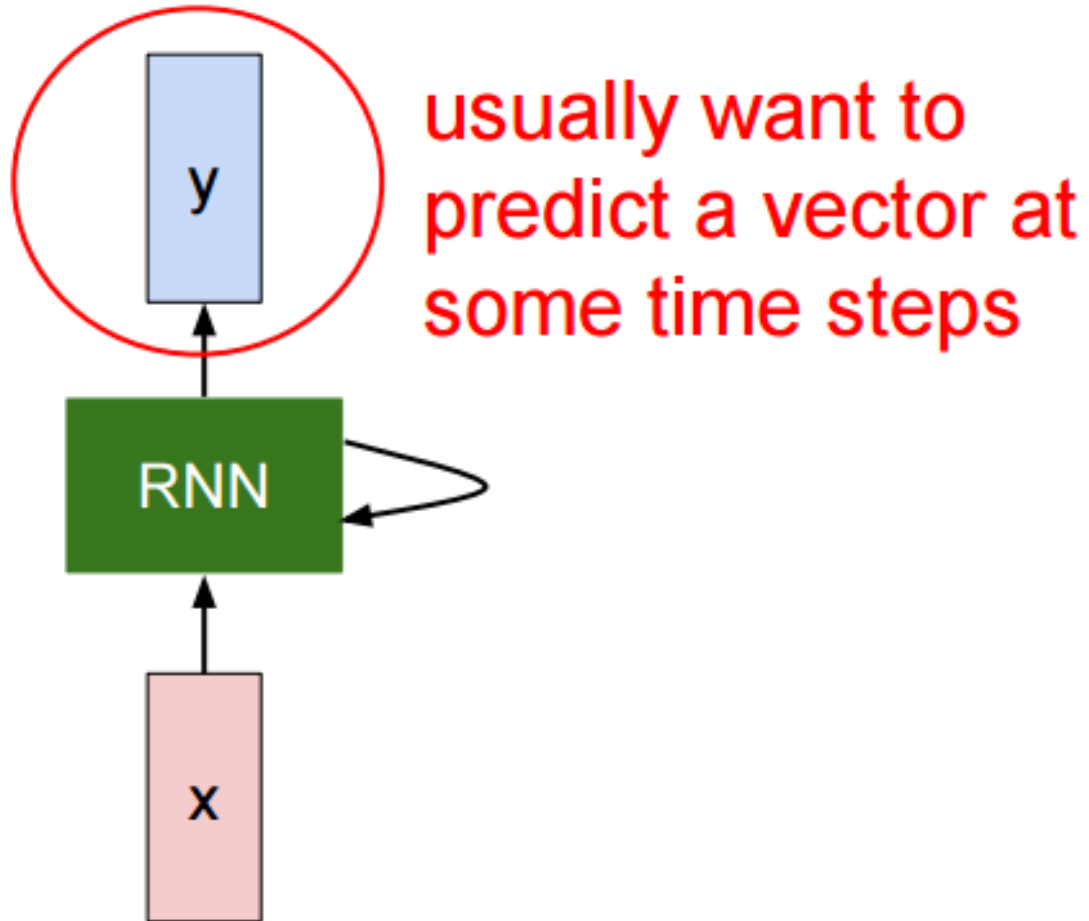


e.g. **Video classification on frame level**



Recurrent Neural Network

- Inputs and outputs of RNNs (rolled version)



Recurrent Neural Network

- How to calculate the hidden state of RNNs

We can process a sequence of vectors \mathbf{x} by applying a recurrence formula at every time step:

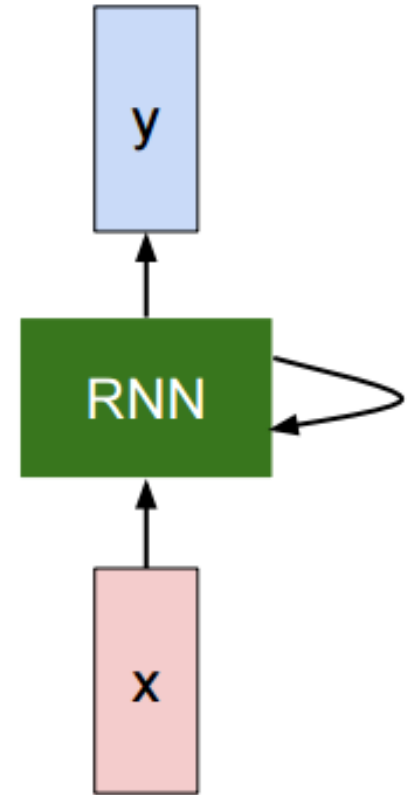
$$\boxed{h_t} = \boxed{f_W}(\boxed{h_{t-1}}, \boxed{x_t})$$

new state

some function with parameters W

old state

input vector at some time step



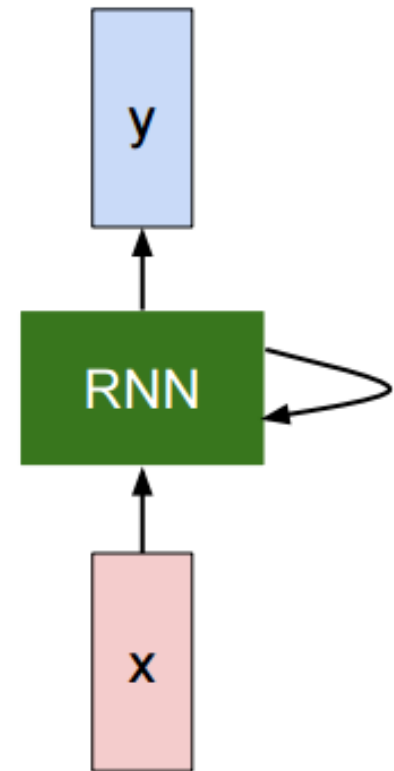
Recurrent Neural Network

- How to calculate the hidden state of RNNs

We can process a sequence of vectors \mathbf{x} by applying a recurrence formula at every time step:

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

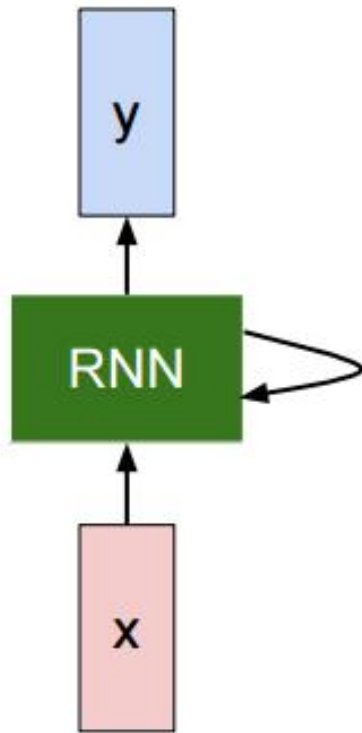
Notice: the same function and the same set of parameters are used at every time step.



Recurrent Neural Network

- How to calculate output of RNN

The state consists of a single “*hidden*” vector \mathbf{h} :



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



$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$

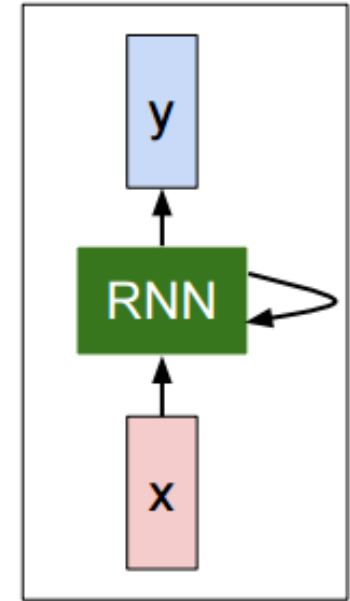
Character-level Language Model

- Example of training sequence “hello”

Character-level language model example

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”



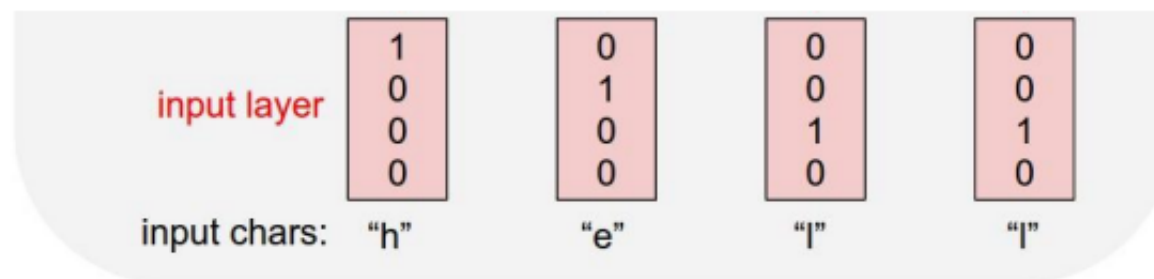
Character-level Language Model

- Example of training sequence “hello”

Character-level language model example

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”



Character-level Language Model

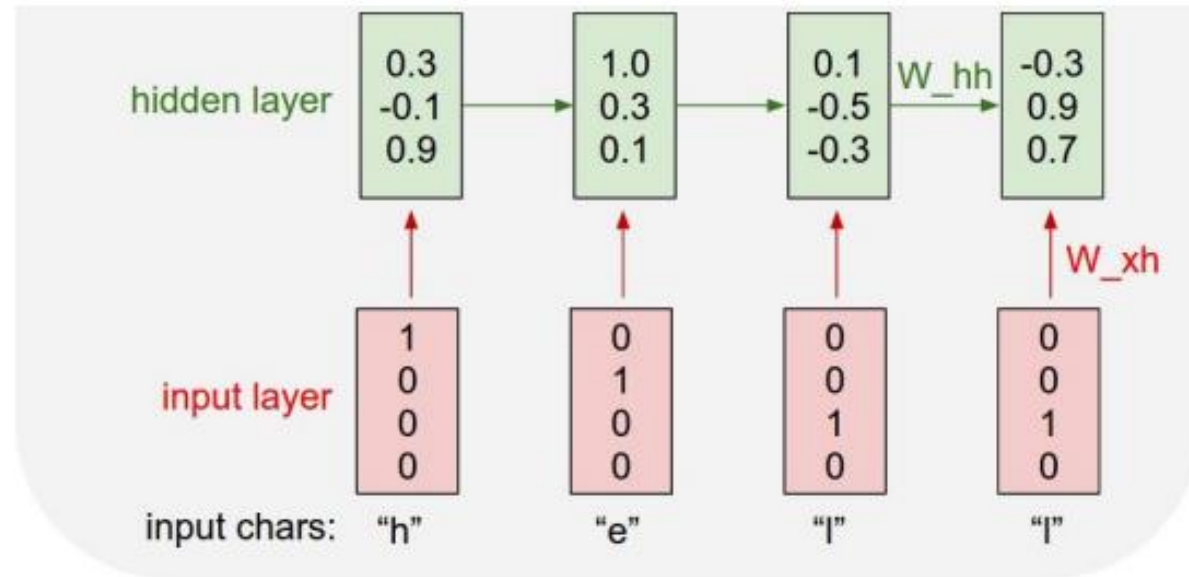
- Example of training sequence “hello”

Character-level language model example

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$



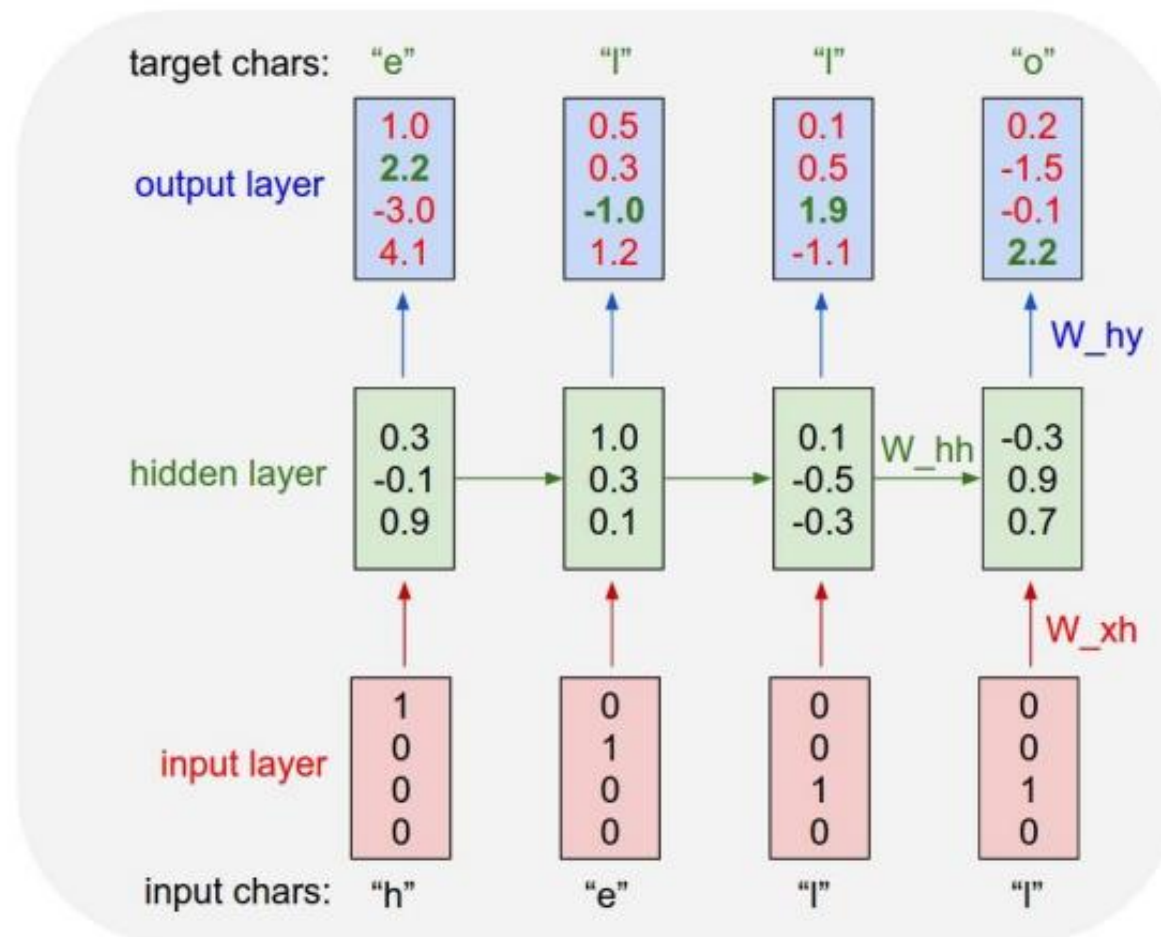
Character-level Language Model

- Example of training sequence “hello”

Character-level language model example

Vocabulary:
[h,e,l,o]

Example training sequence:
“hello”



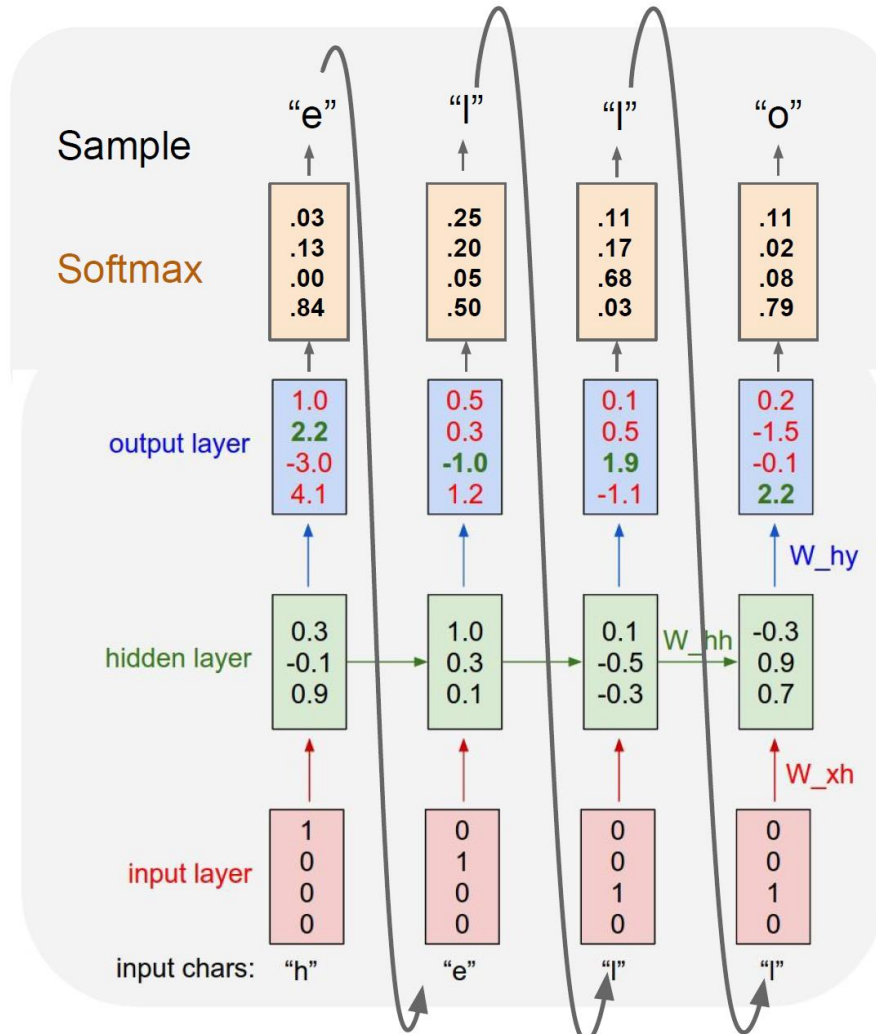
Character-level Language Model

- Example of training sequence “hello”

**Example:
Character-level
Language Model
Sampling**

Vocabulary:
[h,e,l,o]

At test-time sample
characters one at a time,
feed back to model



min-char-rnn.py

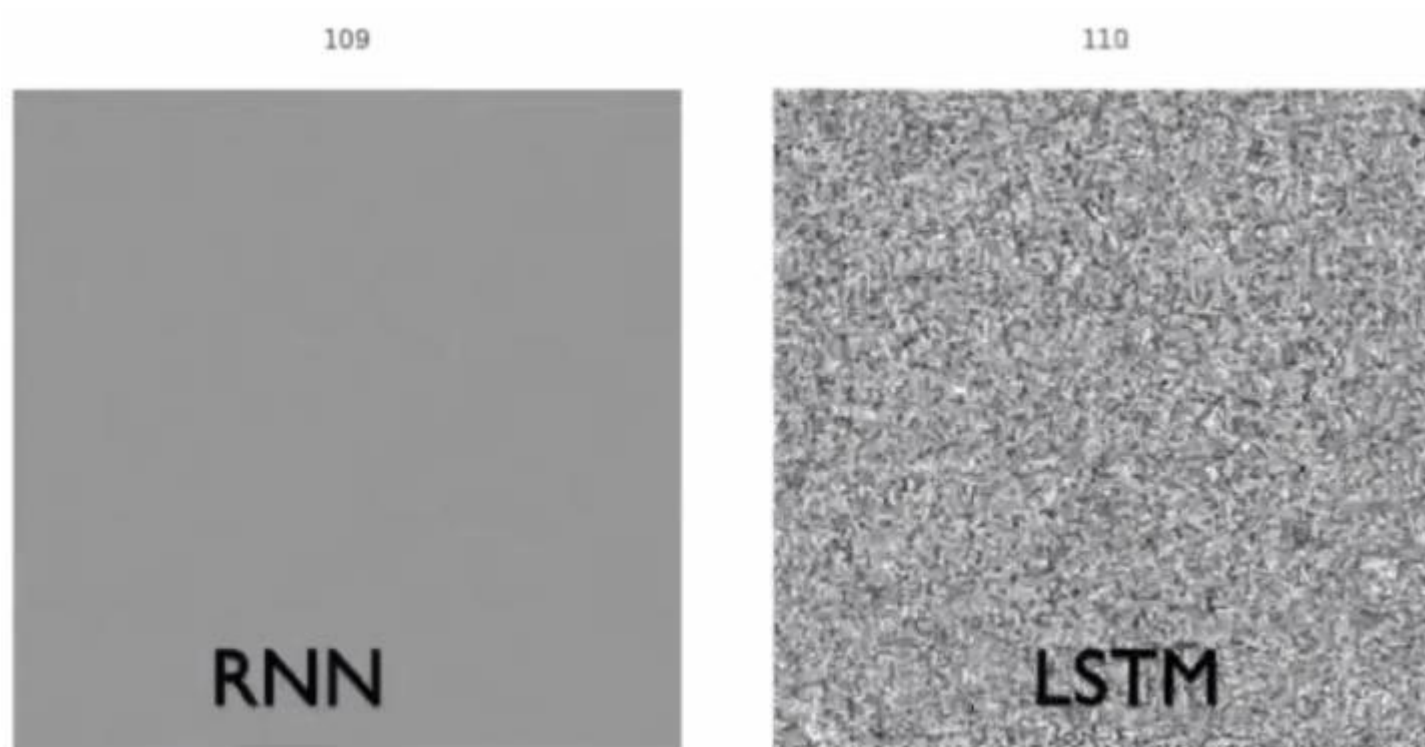
min-char-rnn.py gist: 112 lines of Python

```
1  """
2  Minimal character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)
3  BSD license
4  """
5  import numpy as np
6
7  # data I/O
8  data = open('input.txt', 'r').read() # should be simple plain text file
9  chars = list(set(data))
10 data_size, vocab_size = len(data), len(chars)
11 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
12 char_to_ix = { ch:i for i,ch in enumerate(chars) }
13 ix_to_char = { i:ch for i,ch in enumerate(chars) }
14
15 # hyperparameters
16 hidden_size = 100 # size of hidden layer of neurons
17 seq_length = 25 # number of steps to unroll the RNN for
18 learning_rate = 1e-1
19
20 # model parameters
21 wxh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
22 whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
23 why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
24 bh = np.zeros((hidden_size, 1)) # hidden bias
25 by = np.zeros((vocab_size, 1)) # output bias
26
27 def lossFun(inputs, targets, hprev):
28     """
29     inputs, targets are both list of integers.
30     hprev is Hx1 array of initial hidden state
31     returns the loss, gradients on model parameters, and last hidden state
32     """
33     xs, hs, ys, ps = {}, {}, {}, {}
34     hs[-1] = np.copy(hprev)
35     loss = 0
36     # forward pass
37     for t in xrange(len(inputs)):
38         xs[t] = np.zeros((vocab_size,1)) # encode in 1-of-k representation
39         xs[t][inputs[t]] = 1
40         hs[t] = np.tanh(np.dot(wxh, xs[t]) + np.dot(whh, hs[t-1]) + bh) # hidden state
41         ys[t] = np.dot(why, hs[t]) + by # unnormalized log probabilities for next chars
42         ps[t] = np.exp(ys[t]) / np.sum(np.exp(ys[t])) # probabilities for next chars
43         loss += -np.log(ps[t][targets[t],0]) # softmax (cross-entropy loss)
44     # backward pass: compute gradients going backwards
45     dwxh, dwhh, dwhy = np.zeros_like(wxh), np.zeros_like(whh), np.zeros_like(why)
46     dbh, dby = np.zeros_like(bh), np.zeros_like(by)
47     dhnext = np.zeros_like(hs[0])
48     for t in reversed(xrange(len(inputs))):
49         dy = np.copy(ps[t])
50         dy[targets[t]] -= 1 # backprop into y
51         dwhy += np.dot(dy, hs[t].T)
52         dby += dy
53         dh = np.dot(why.T, dy) + dhnext # backprop into h
54         dhraw = (1 - hs[t] * hs[t]) * dh # backprop through tanh nonlinearity
55         dbh += dhraw
56         dwxh += np.dot(dhraw, xs[t].T)
57         dwhh += np.dot(dhraw, hs[t-1].T)
58         dhnext = np.dot(whh.T, dhraw)
59     for dparam in [dwxh, dwhh, dwhy, dbh, dby]:
60         np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
61     return loss, dwxh, dwhh, dwhy, dbh, dby, hs[len(inputs)-1]
62
63 def sample(h, seed_ix, n):
64     """
65     sample a sequence of integers from the model
66     h is memory state, seed_ix is seed letter for first time step
67     """
68     x = np.zeros((vocab_size, 1))
69     x[seed_ix] = 1
70     ixes = []
71     for t in xrange(n):
72         h = np.tanh(np.dot(wxh, x) + np.dot(whh, h) + bh)
73         y = np.dot(why, h) + by
74         p = np.exp(y) / np.sum(np.exp(y))
75         ix = np.random.choice(range(vocab_size), p=p.ravel())
76         x = np.zeros((vocab_size, 1))
77         x[ix] = 1
78         ixes.append(ix)
79     return ixes
80
81 n, p = 0, 0
82 mwxh, mwhh, mwhy = np.zeros_like(wxh), np.zeros_like(whh), np.zeros_like(why)
83 mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
84 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
85 while True:
86     # prepare inputs (we're sweeping from left to right in steps seq_length long)
87     if p+seq_length+1 >= len(data) or n == 0:
88         hprev = np.zeros((hidden_size,1)) # reset RNN memory
89         p = 0 # go from start of data
90     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
91     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
92
93     # sample from the model now and then
94     if n % 100 == 0:
95         sample_ix = sample(hprev, inputs[0], 200)
96         txt = ''.join(ix_to_char[ix] for ix in sample_ix)
97         print '----\n %s \n----' % (txt, )
98
99     # forward seq_length characters through the net and fetch gradient
100    loss, dwxh, dwhh, dwhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
101    smooth_loss = smooth_loss * 0.999 + loss * 0.001
102    if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
103
104    # perform parameter update with Adagrad
105    for param, dparam, mem in zip([dwxh, dwhh, dwhy, dbh, dby],
106                                  [dwxh, dwhh, dwhy, dbh, dby],
107                                  [mwxh, mwhh, mwhy, mbh, mby]):
108        mem += dparam * dparam
109        param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
110
111    p += seq_length # move data pointer
112    n += 1 # iteration counter
```

(<https://gist.github.com/karpathy/d4dee566867f8291f086>)

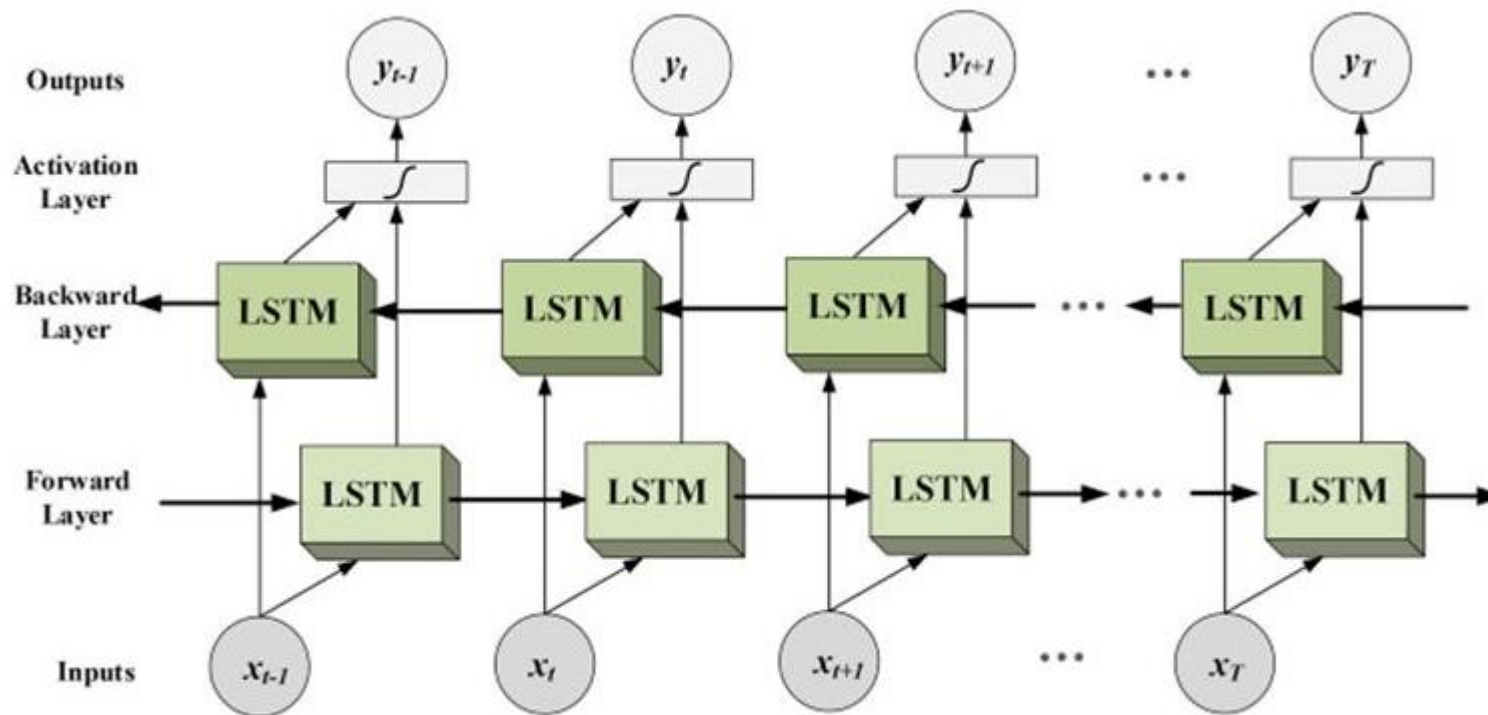
Vanishing Gradient Problem

- The reason why the vanishing gradient problem is important:



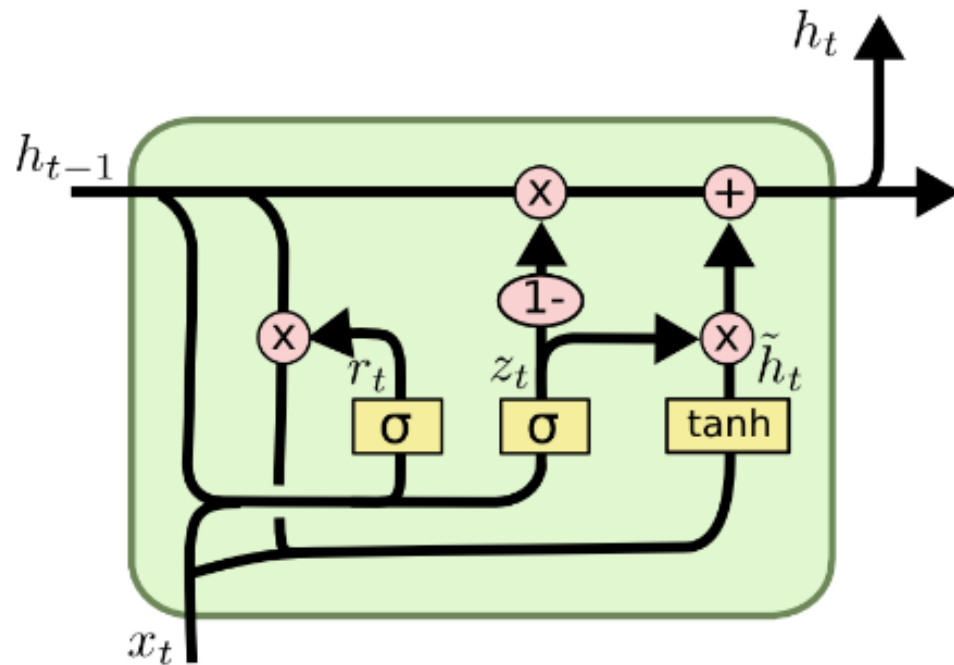
[RNN vs LSTM: Vanishing Gradients - GIF on Imgur](#)

Bidirectional RNN



GRU

- What is GRU(Gated Recurrent Unit)?



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

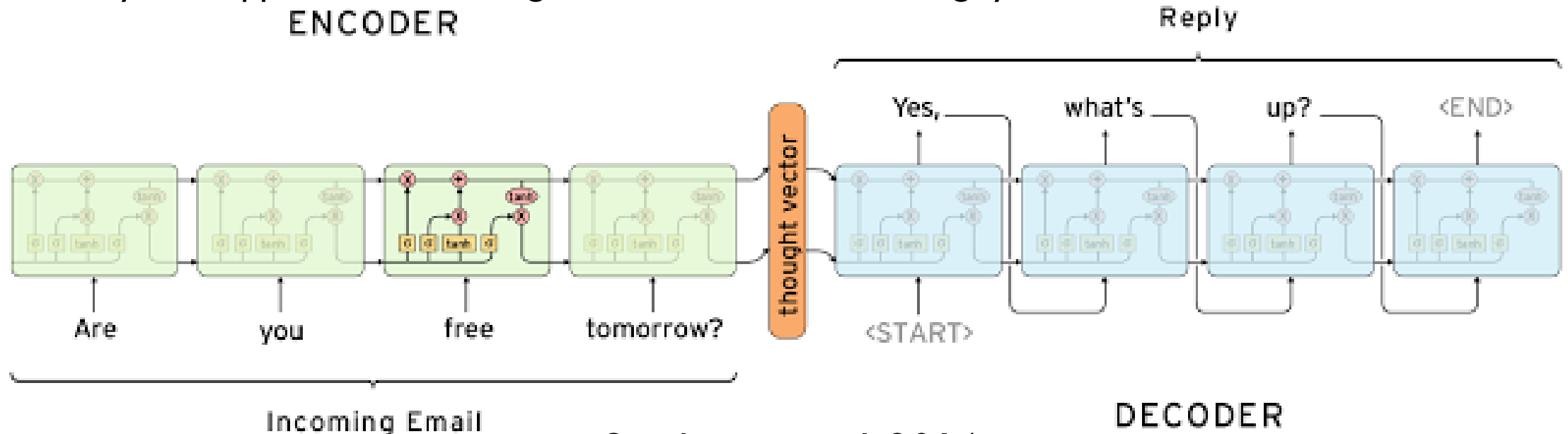
$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

[Understanding LSTM Networks -- colah's blog](#)

Seq2Seq Model

- It takes a sequence of words as input and gives a sequence of words as output.
- It composed of an encoder and a decoder.
- Many NLP applications exist, e.g., machine translation, dialog systems, and so on.

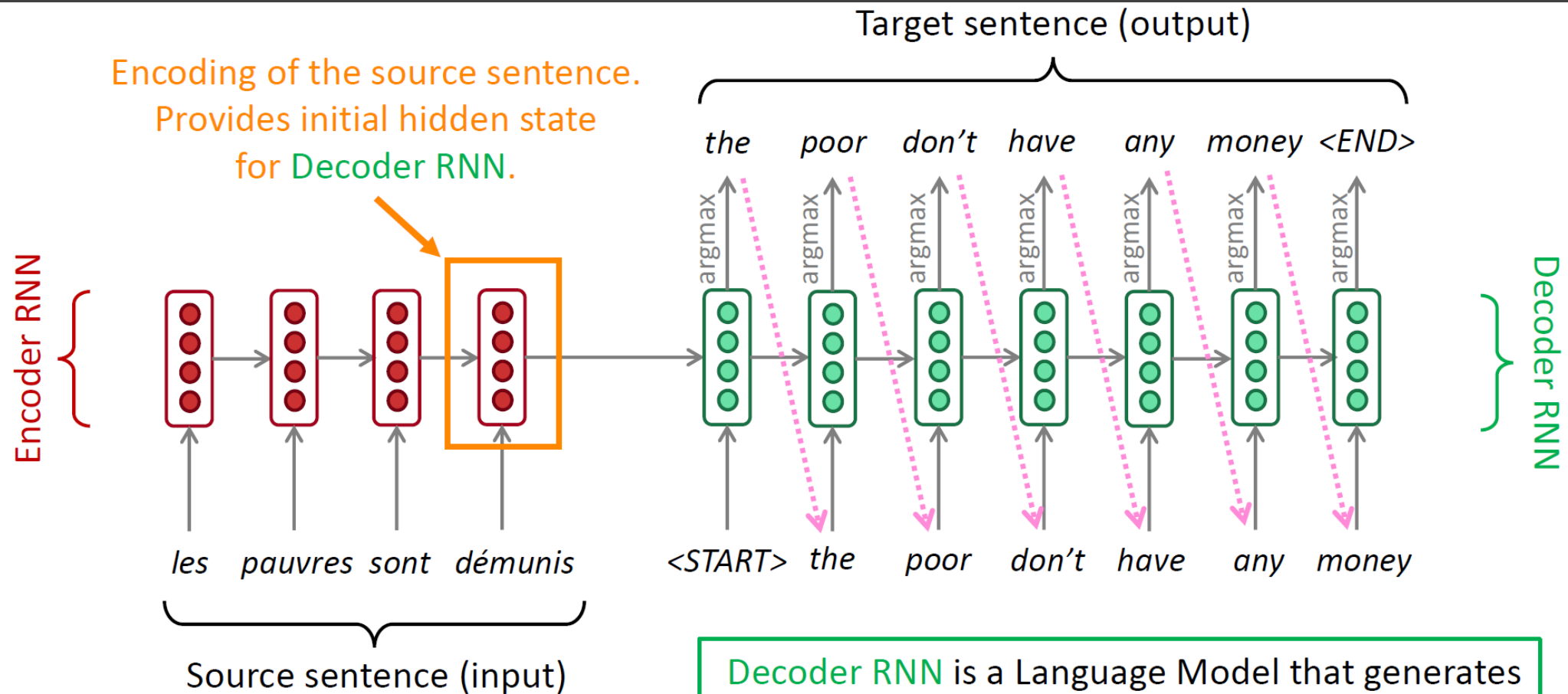


Sutskever et al. 2014

“Sequence to Sequence Learning with Neural Networks”

Encode source into fixed length vector, use it as
initial recurrent state for target decoder model

Seq2Seq for Machine Translation

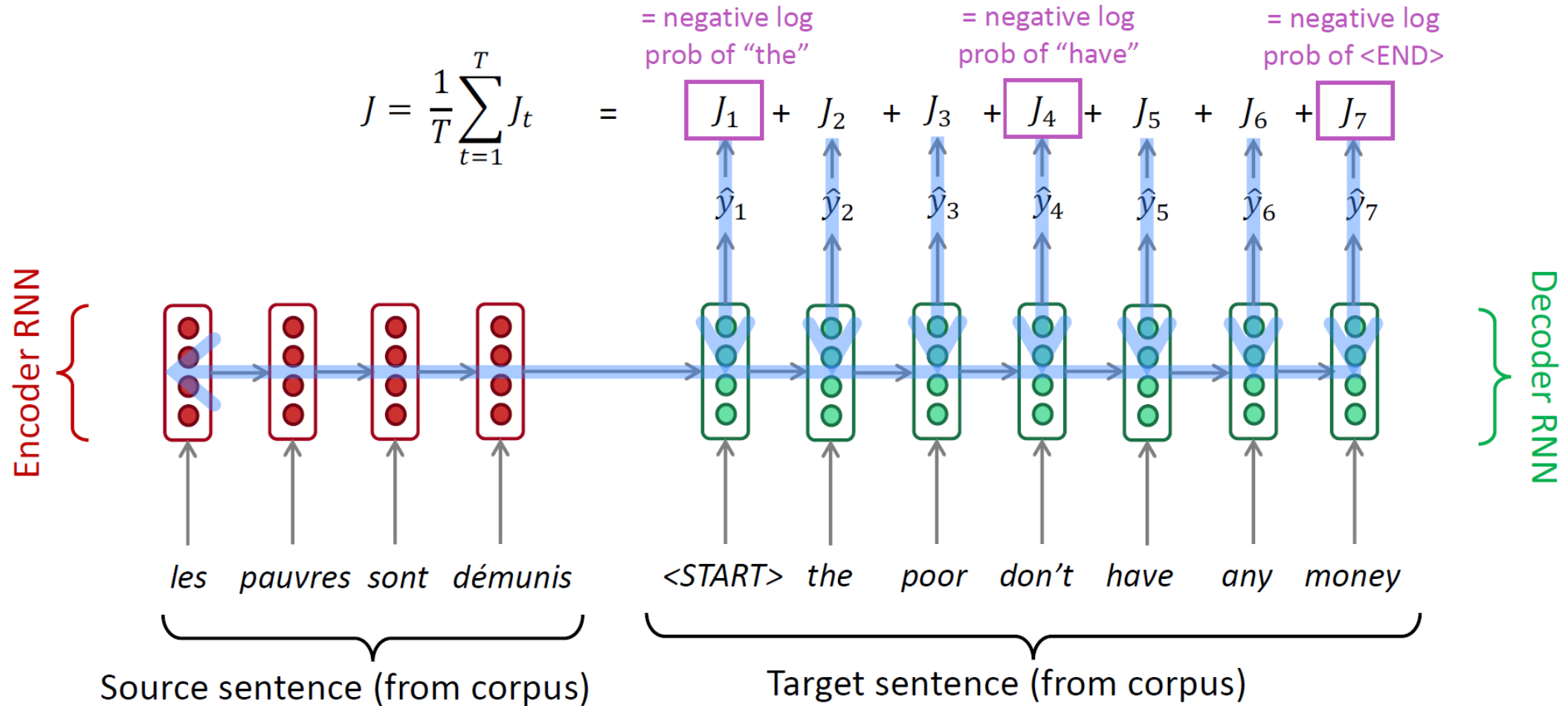


Encoder RNN produces an **encoding** of the source sentence.

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

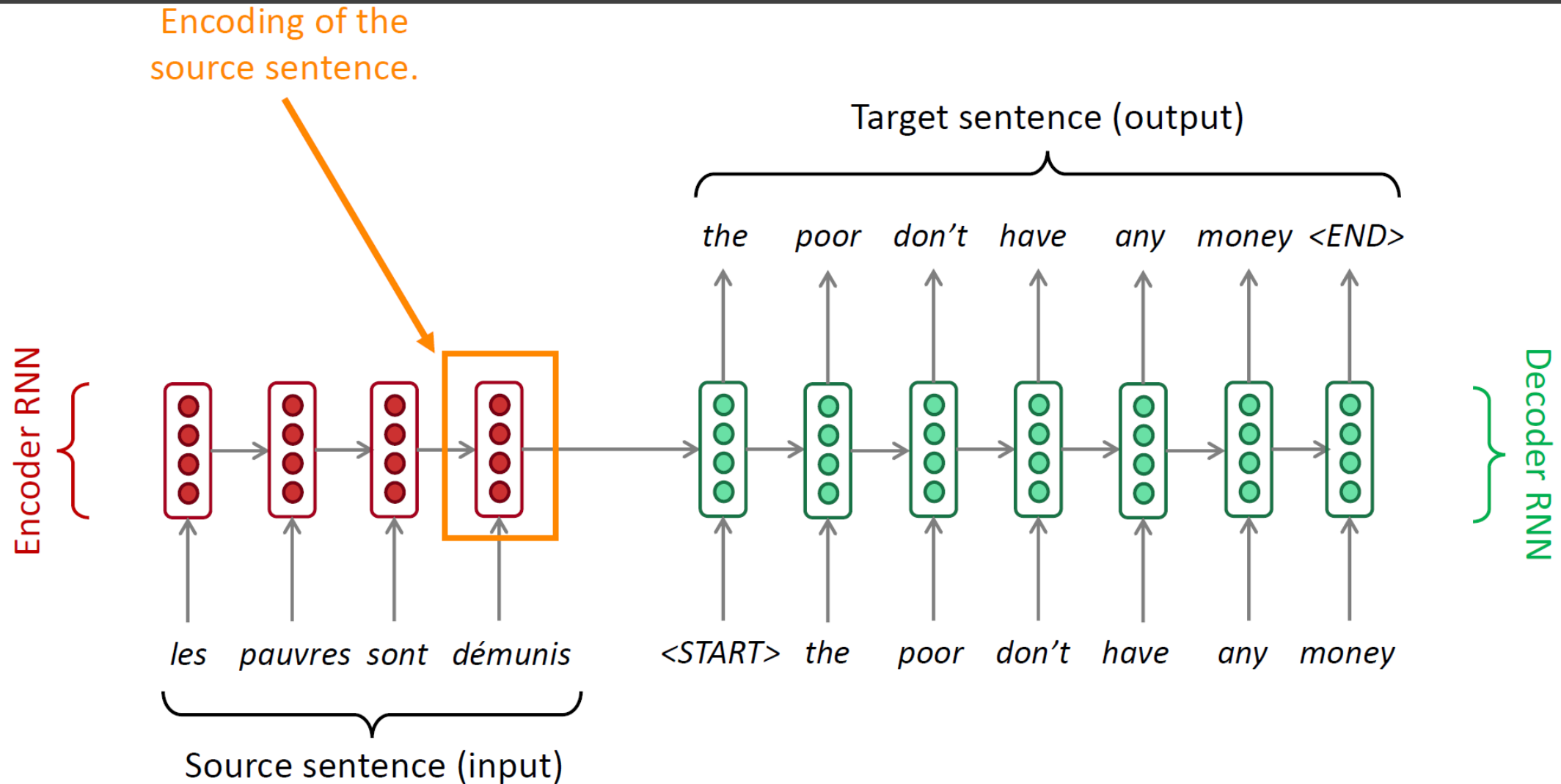
Note: This diagram shows **test time** behavior: decoder output is fed in➤ as next step's input

Seq2Seq Model: Training



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

Seq2Seq Model: Bottleneck Problem



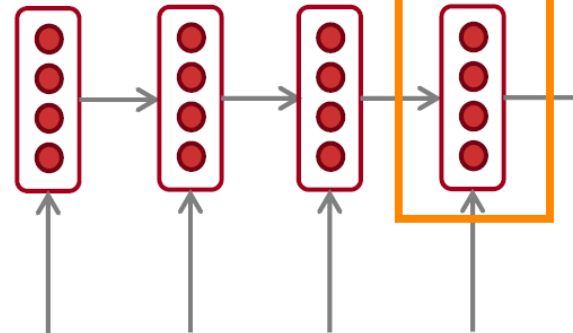
Problems with this architecture?

Seq2Seq Model: Bottleneck Problem

Encoding of the
source sentence.

This needs to capture *all*
information about the
source sentence.
Information bottleneck!

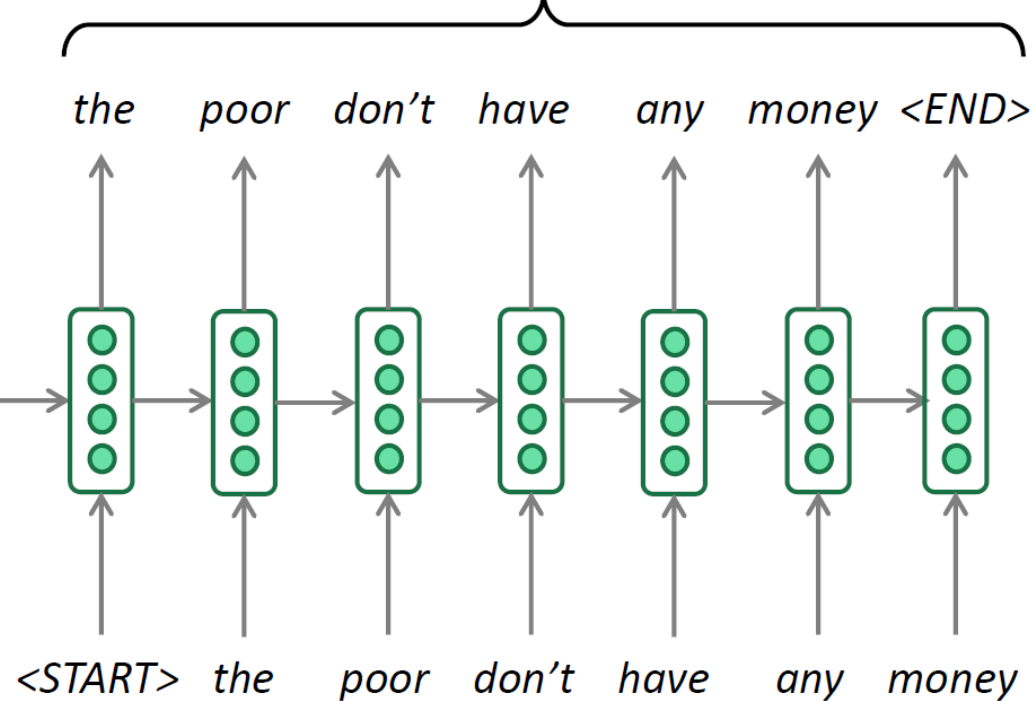
Encoder RNN



les pauvres sont démunis

Source sentence (input)

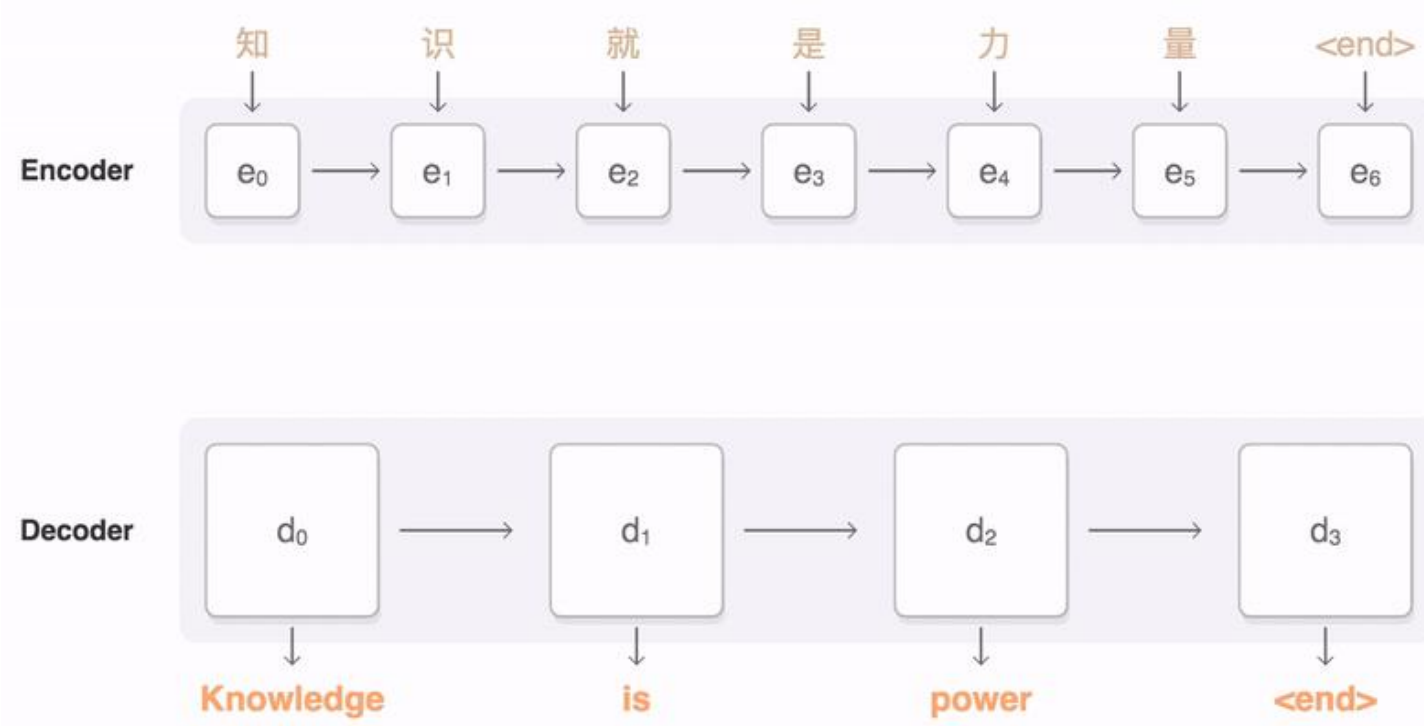
Target sentence (output)

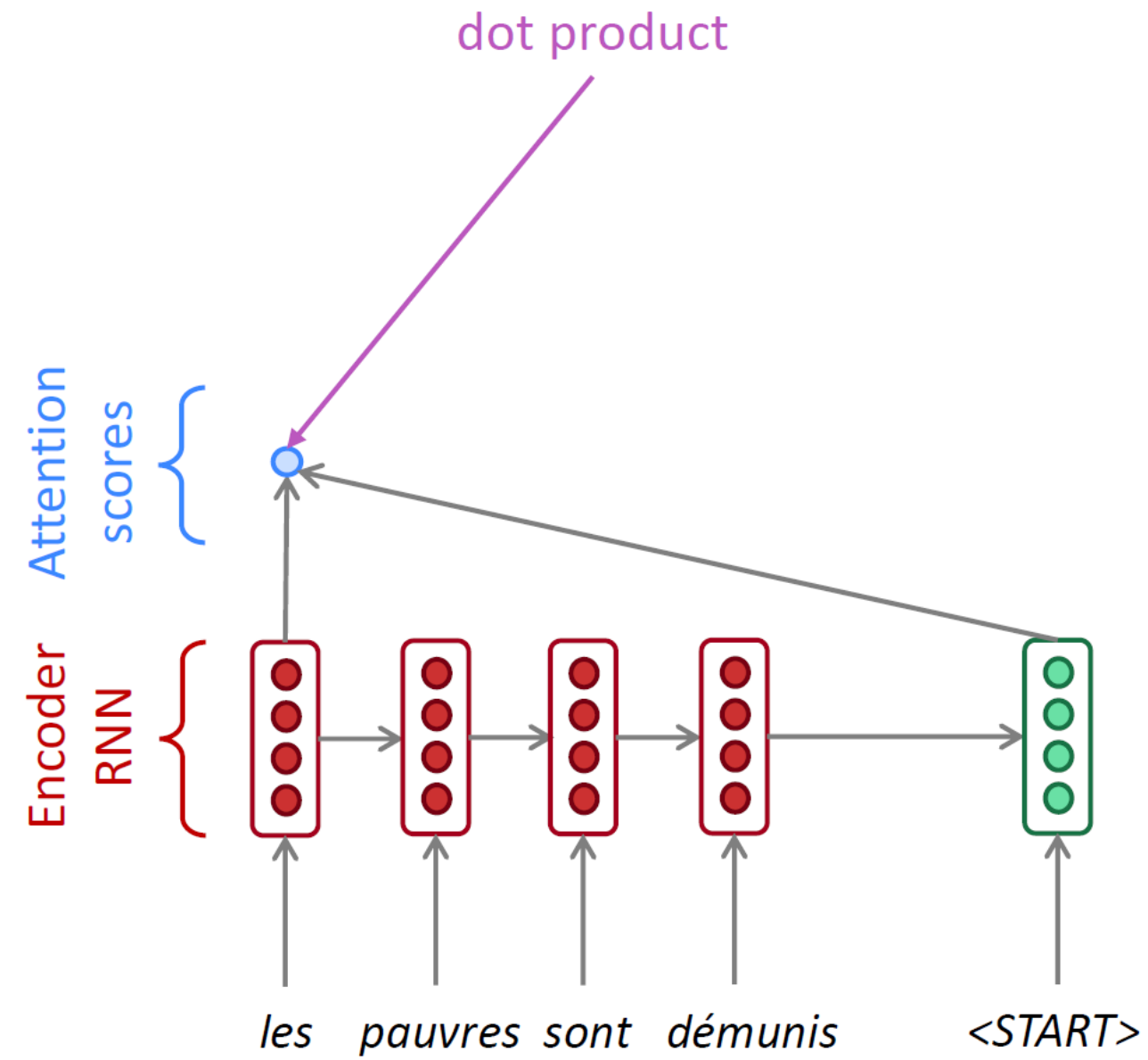


Decoder RNN

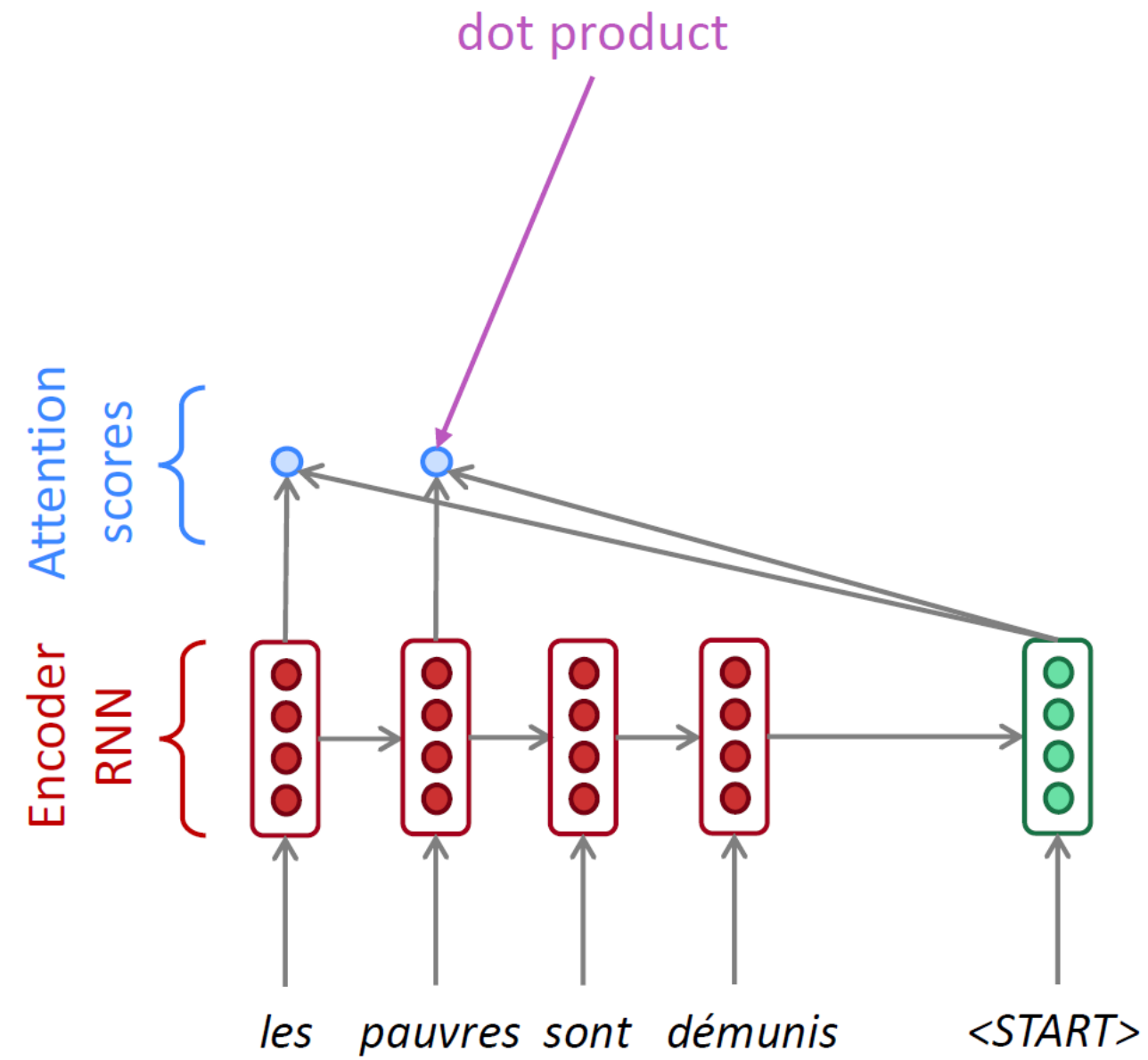
Seq2Seq with Attention

- **Attention** provides a solution to the bottleneck problem.
- Core idea: At each time step of the decoder, **focus on a particular part of the source sequence**
- Tensorflow official implementation: <https://github.com/google/seq2seq>

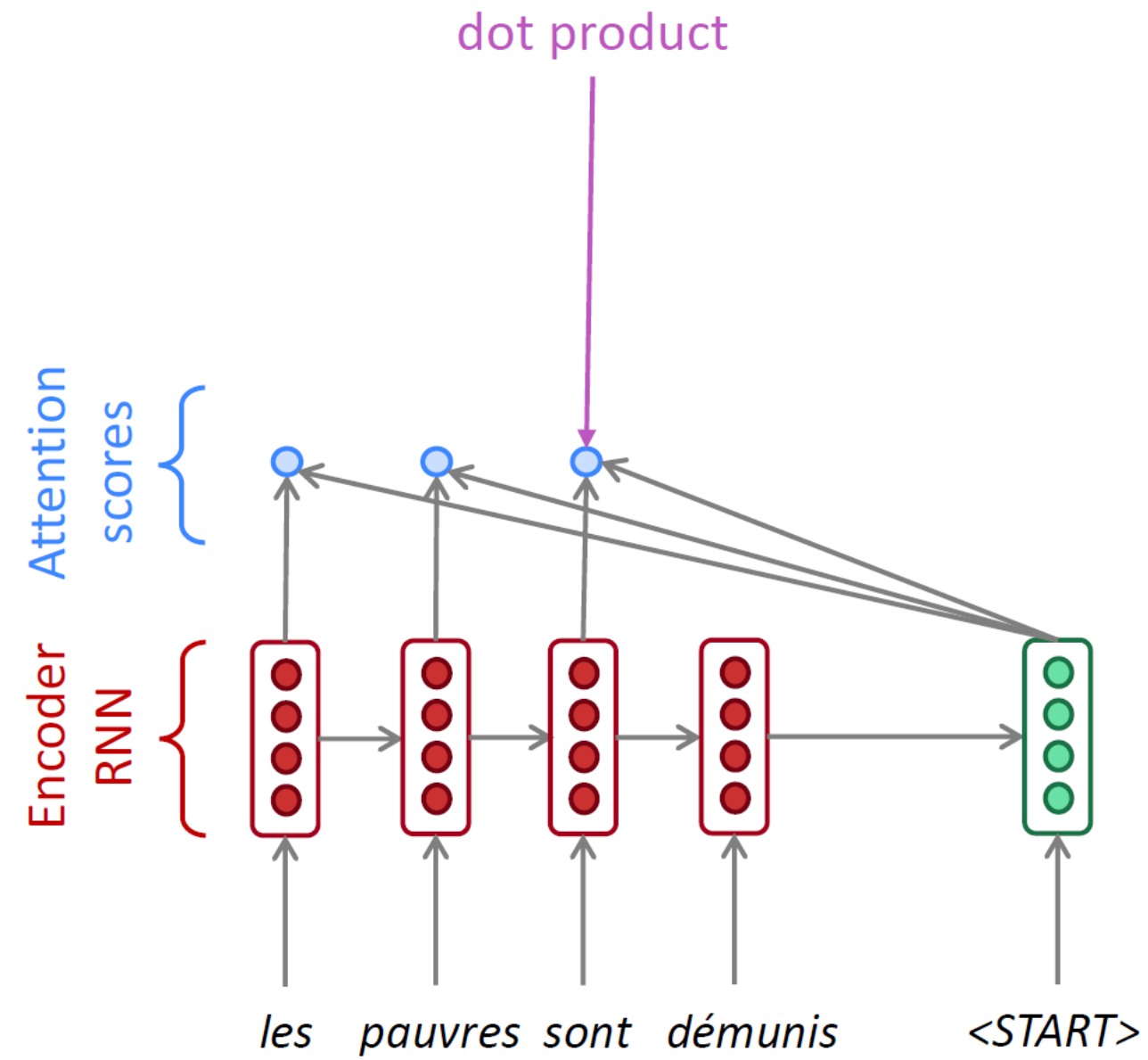




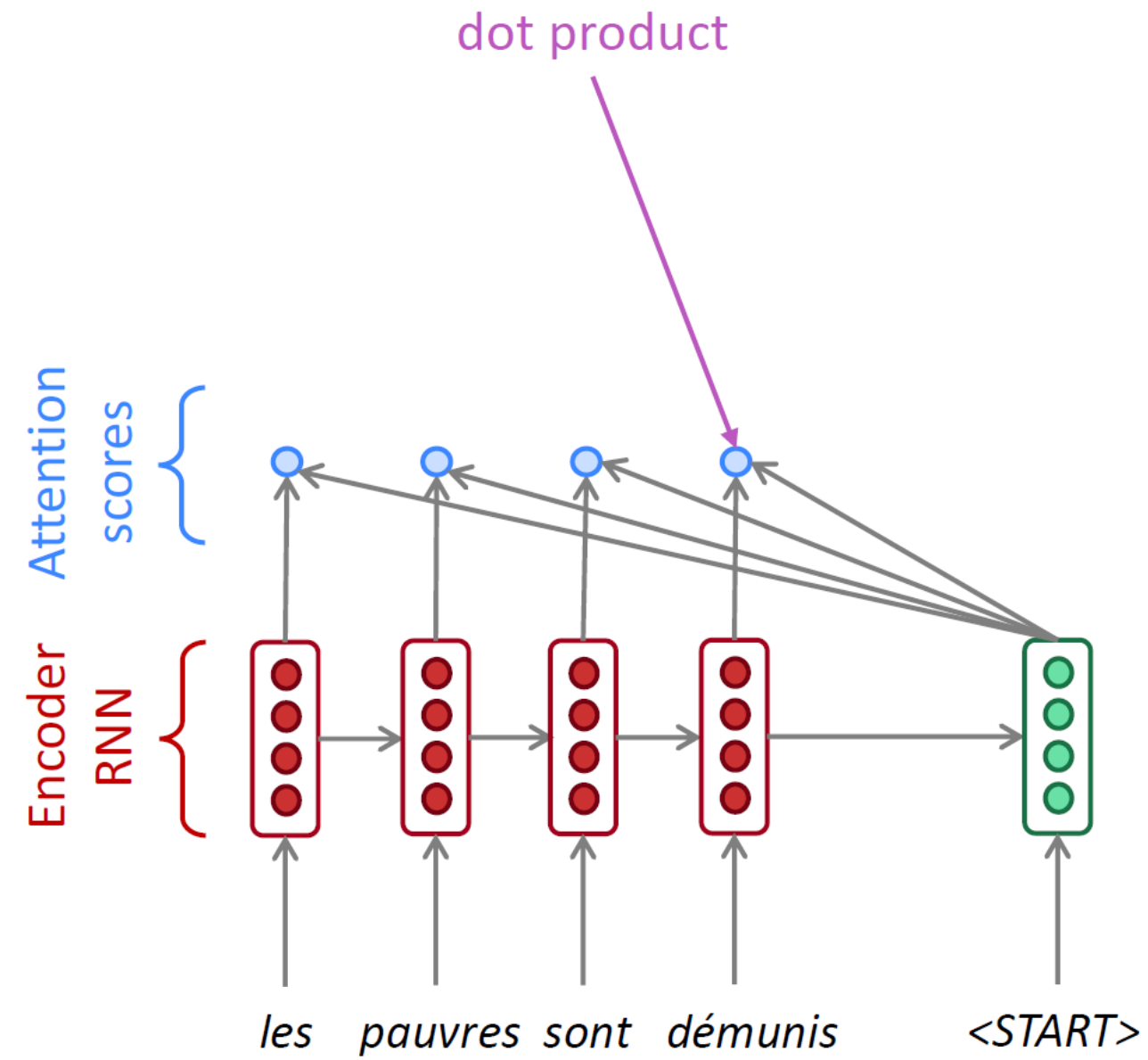
Decoder RNN



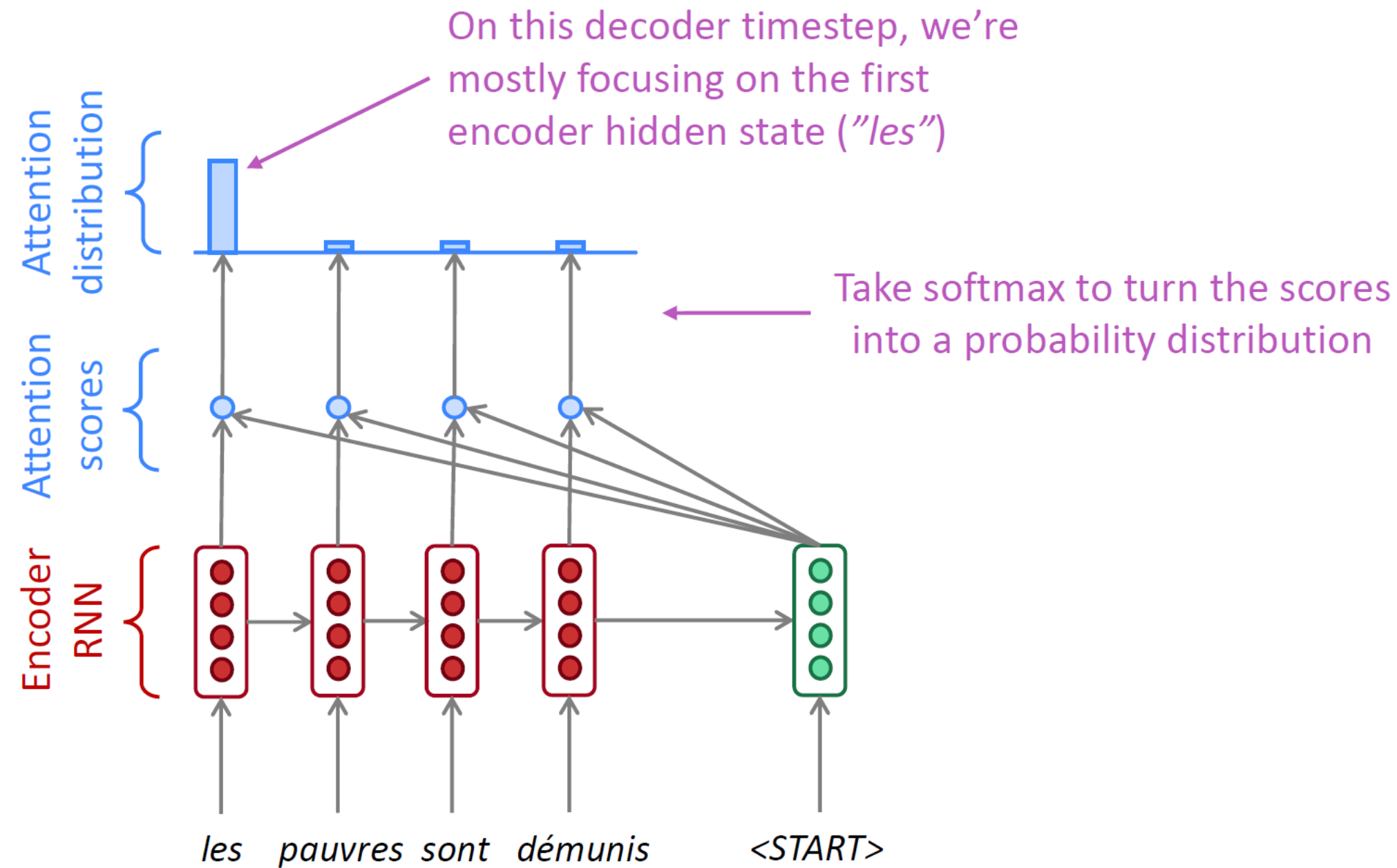
Decoder RNN

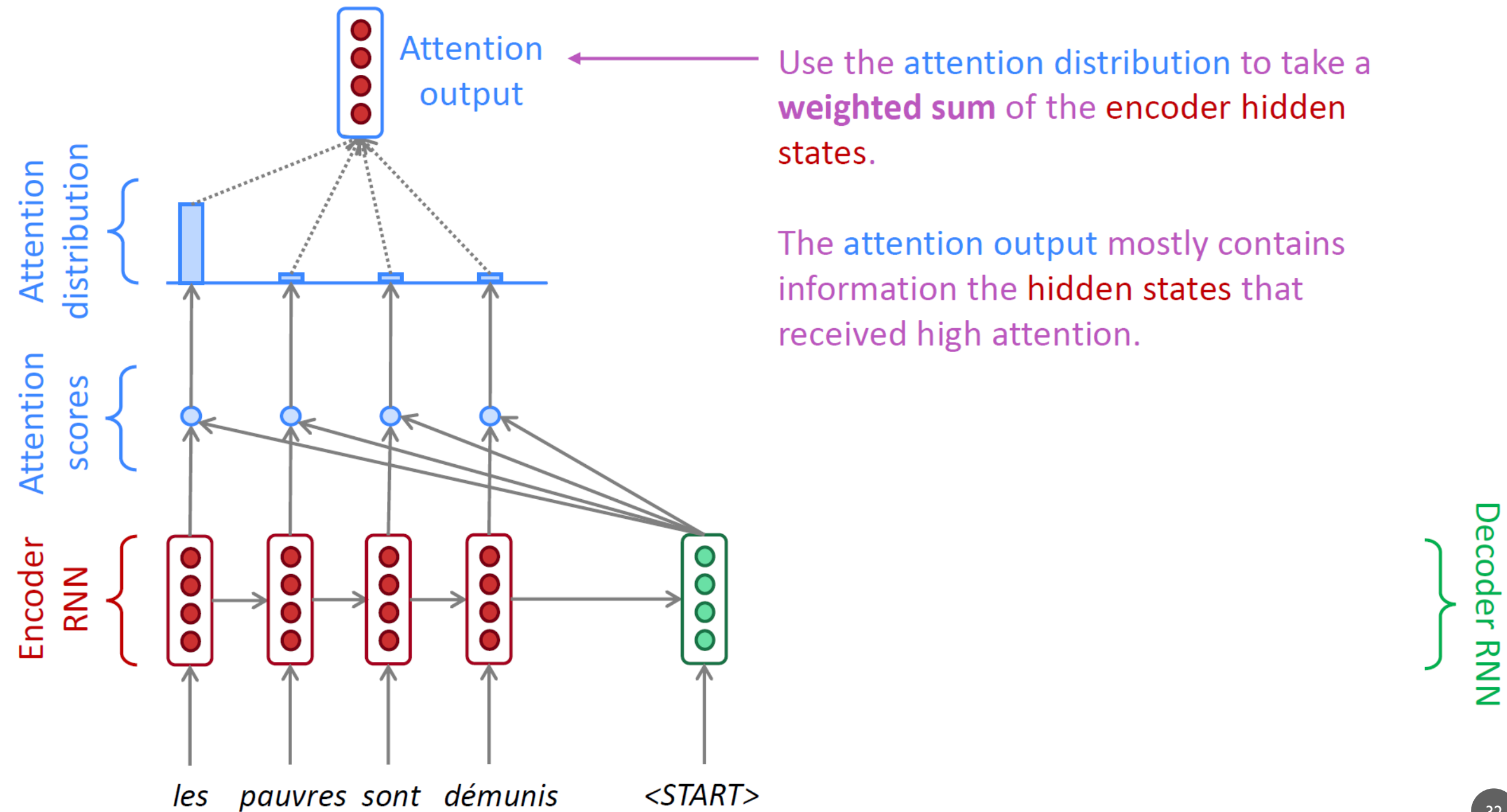


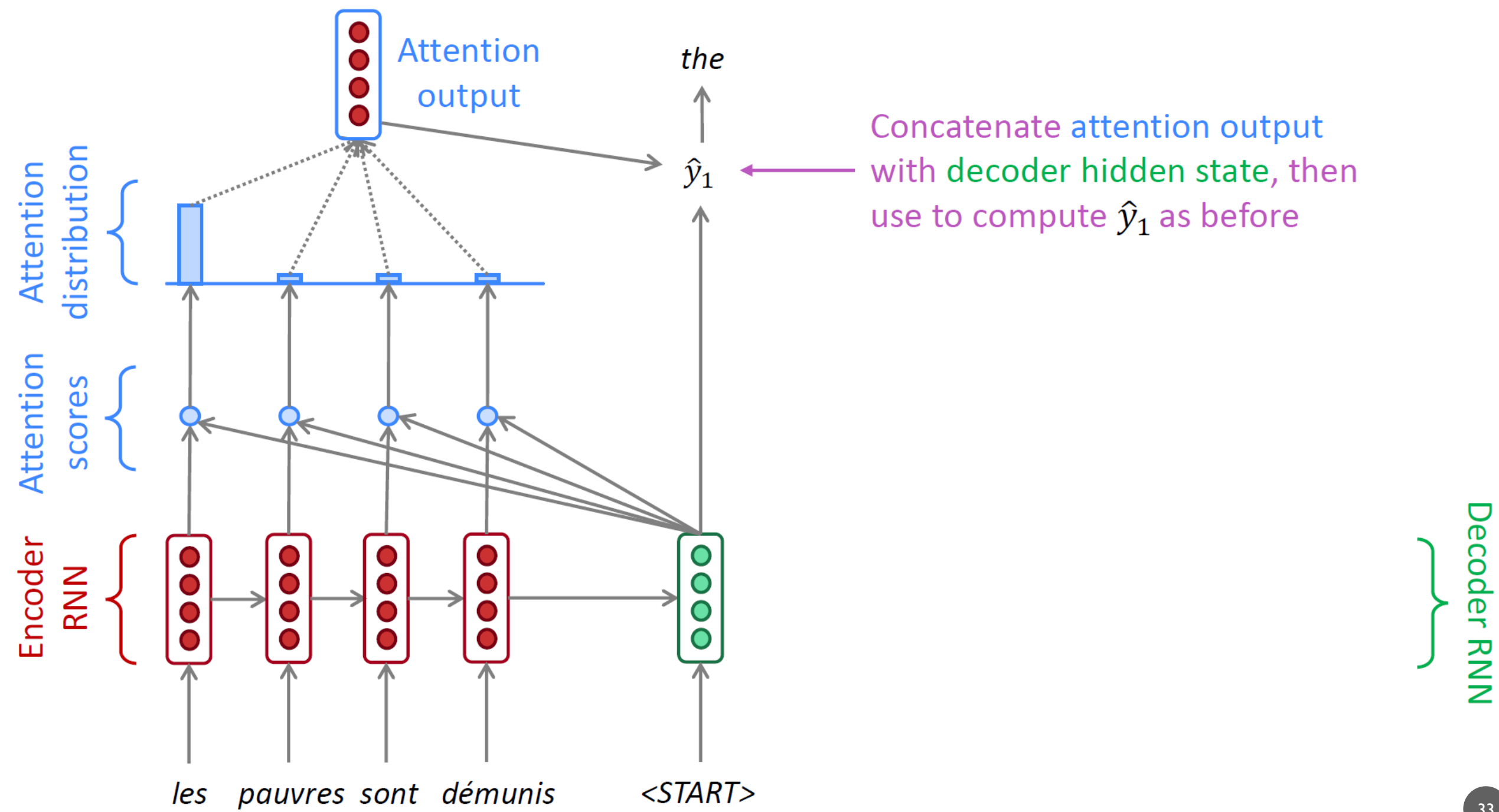
Decoder RNN

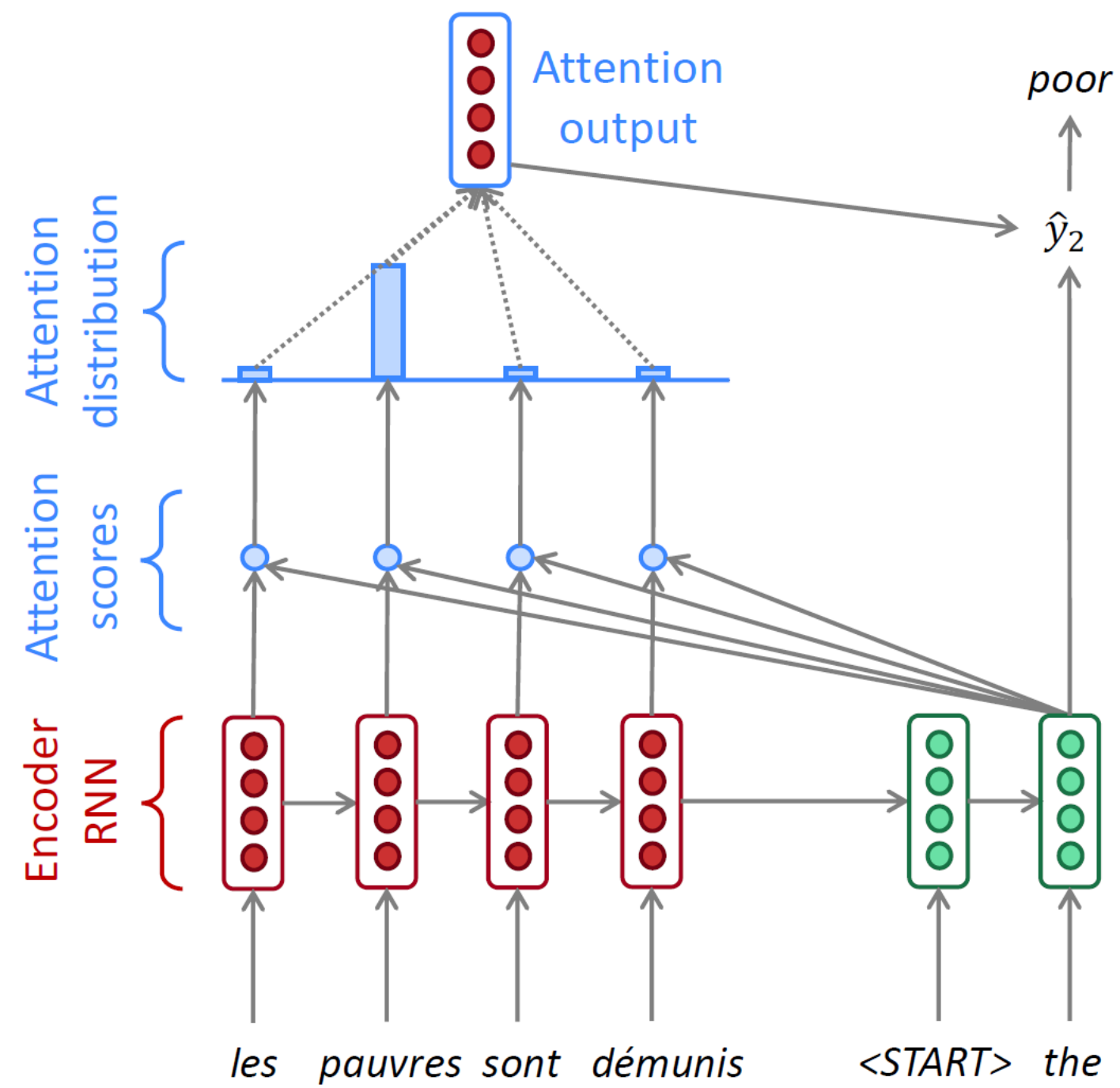


Decoder RNN

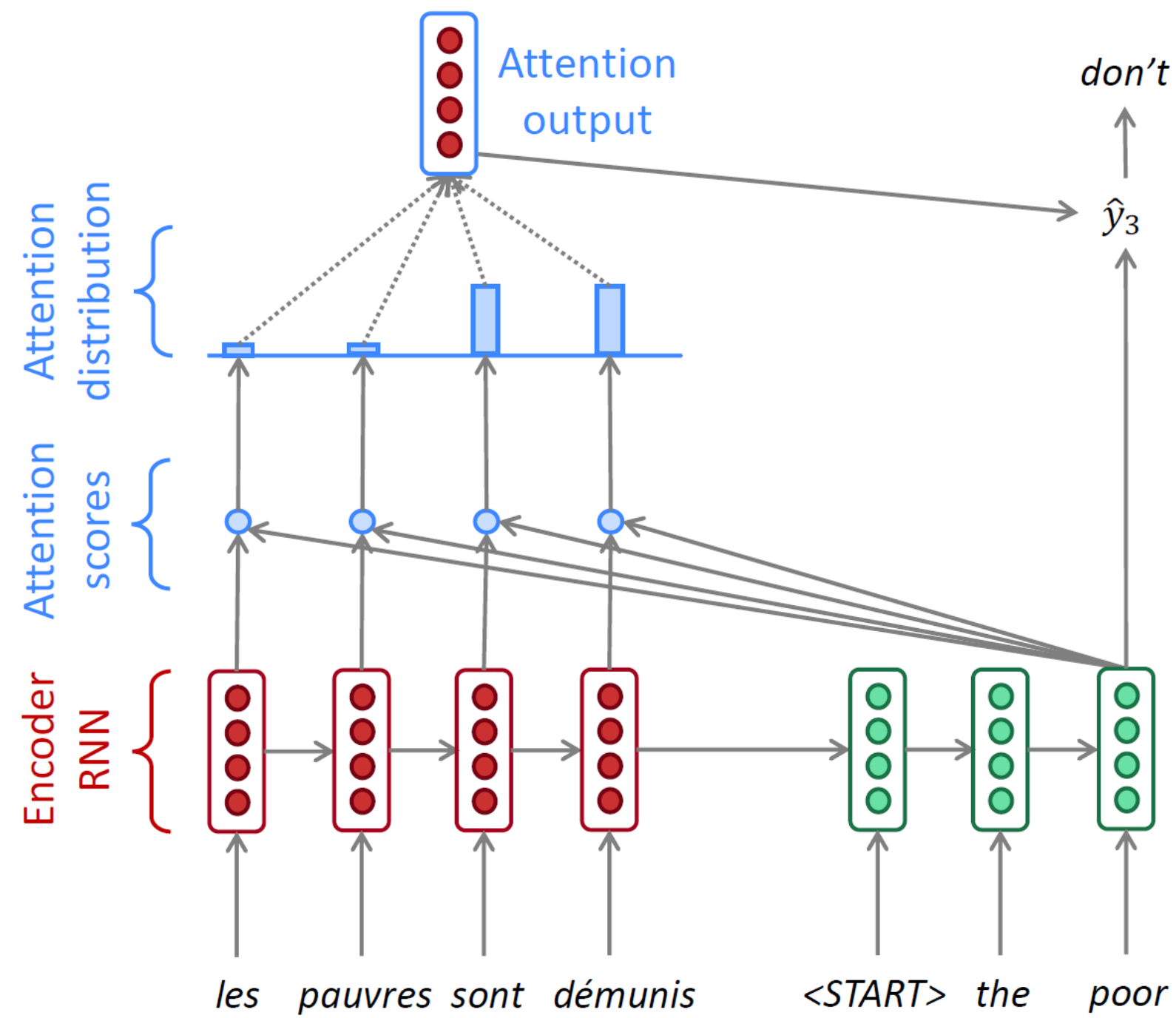




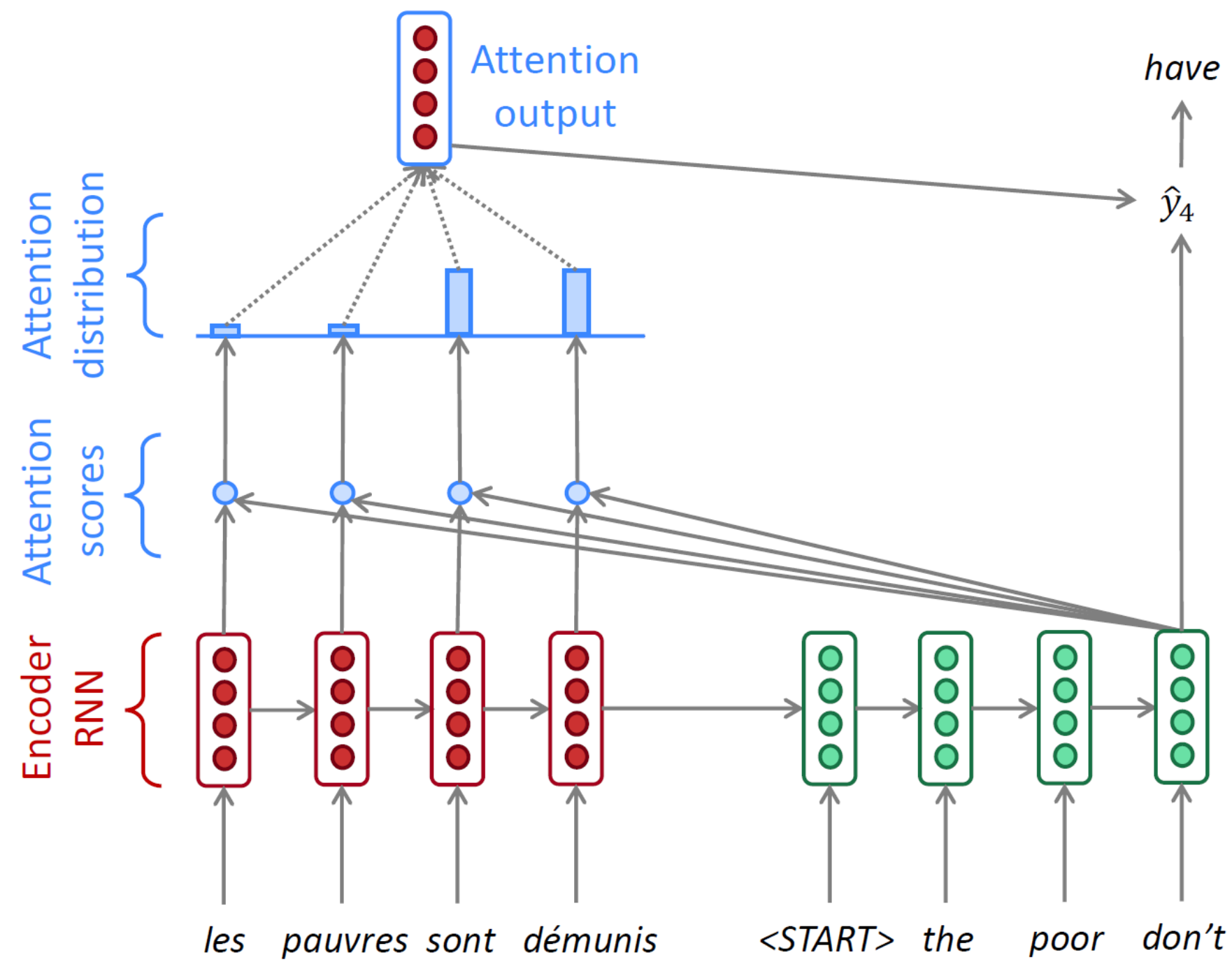




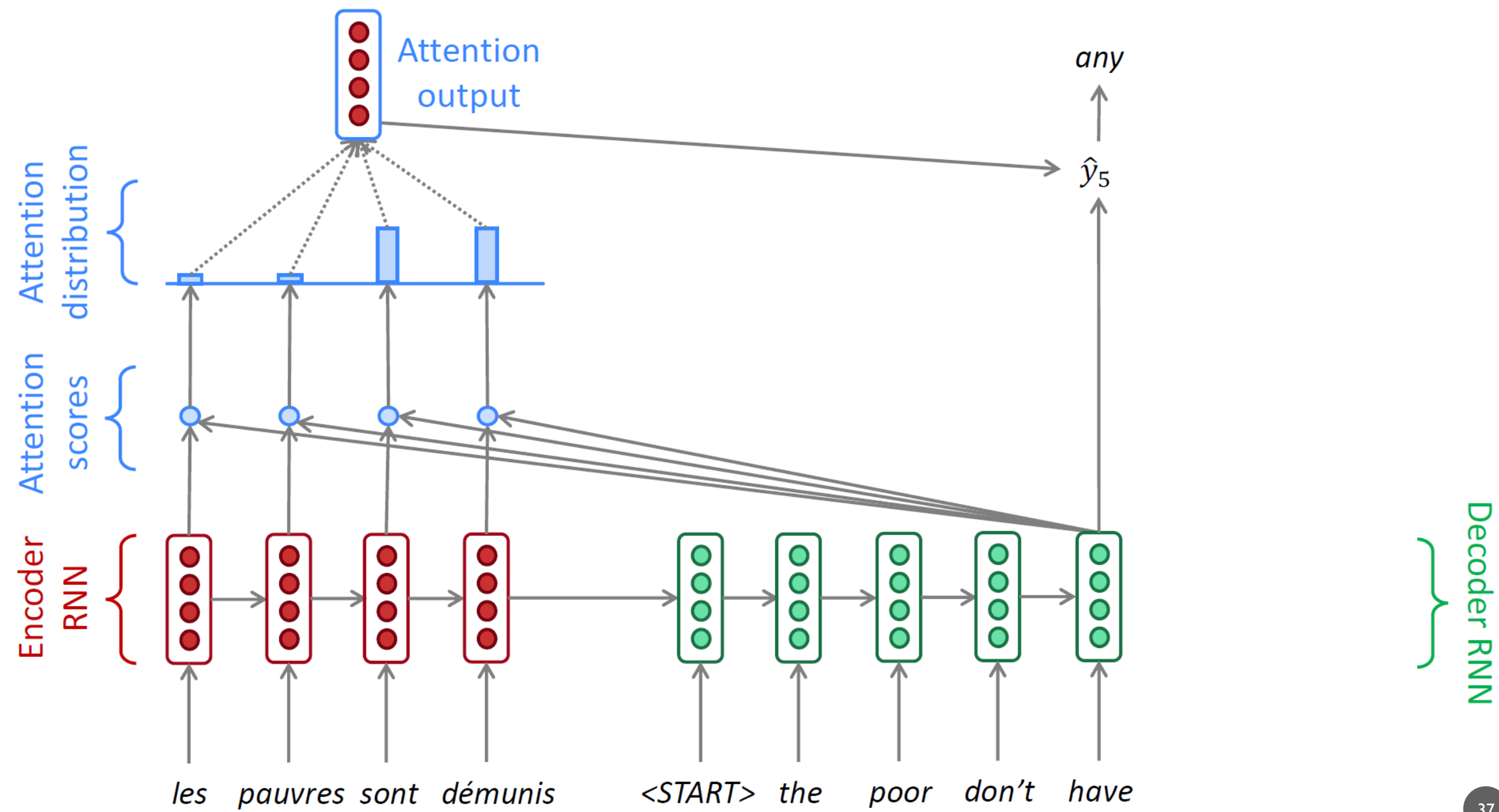
Decoder RNN

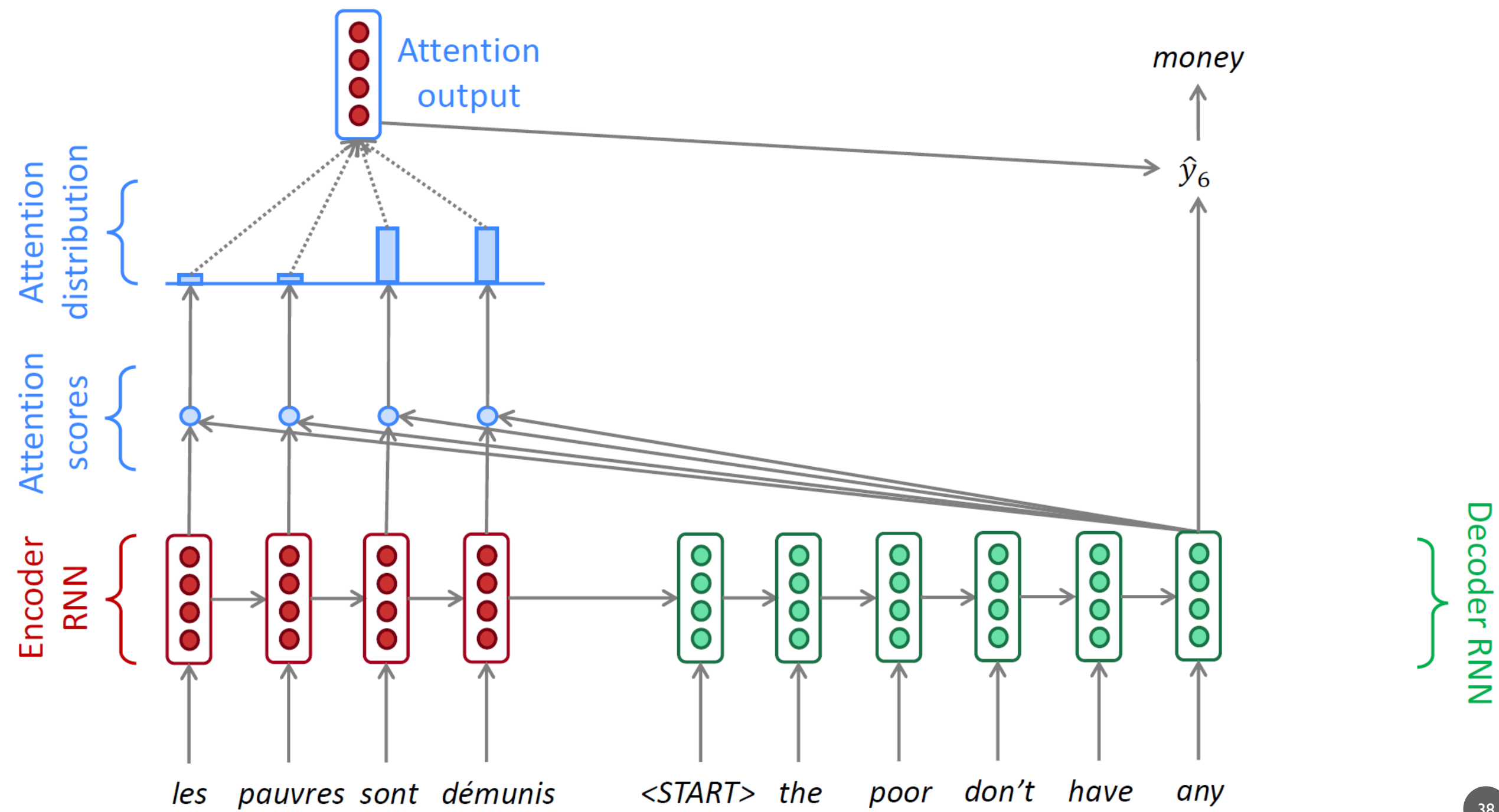


Decoder RNN




Decoder RNN





Attention is Great!

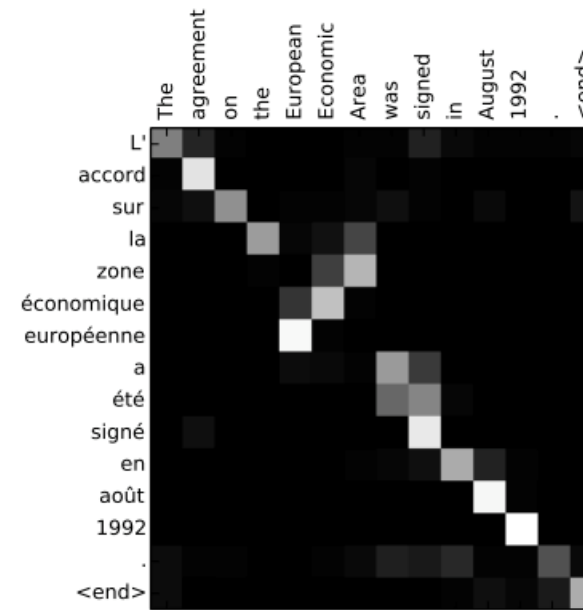
- Attention significantly improves NMT performance
 - It is useful to allow the decoder to focus on particular parts of the source
- Attention solves the bottleneck problem
 - Attention allows the decoder to look directly at source; bypass the bottleneck
- Attention helps with vanishing gradient problem
 - Provides a shortcut to faraway states
- Attention provides some interpretability
 - By inspecting attention distribution, we can see what the decoder was focusing on
 - We get alignment for free!
 - This is cool because we never explicitly trained an alignment system
 - The network just learned alignment by itself



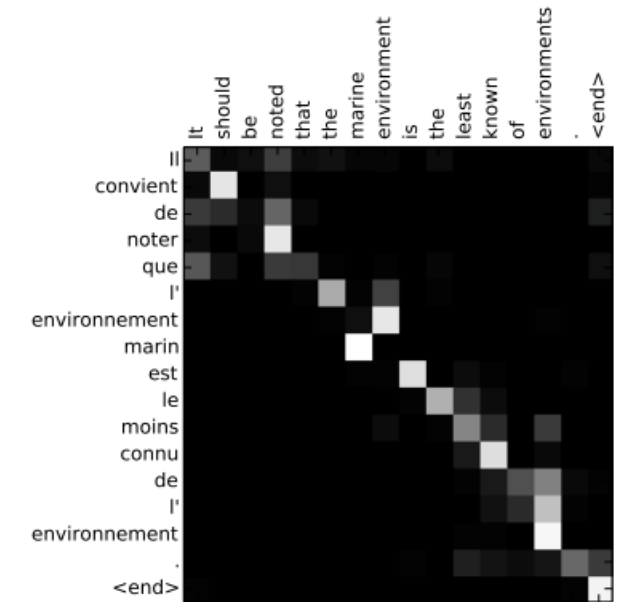
	Les	pauvres	sont	démunis
The	■			
poor		■		
don't			■	■
have			■	■
any			■	■
money			■	■

Attention Examples in Machine Translation

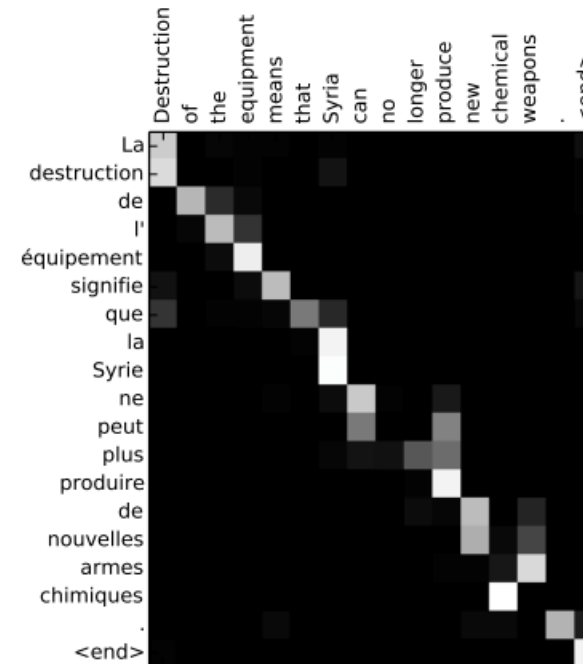
- It properly learns grammatical orders of words
- It skips unnecessary words such as an article



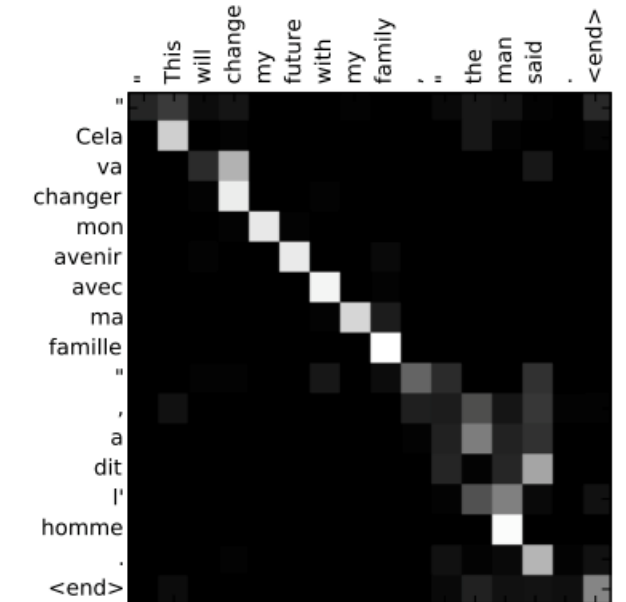
(a)



(b)



(c)



(d)