

# *Attention in Machine Translation*

Nikolay Arefyev

*CMC MSU Department of Algorithmic Languages &  
Samsung Moscow Research Center*

# What is Machine Translation?

- Machine Translation (MT): translate a sentence in one language to a sentence in another language.
  - Special case of **sequence transduction**

***Source language:*** ru

x: Мороз и солнце - день чудесный!

***Target language:*** en

y: Frost and the sun - a wonderful day!

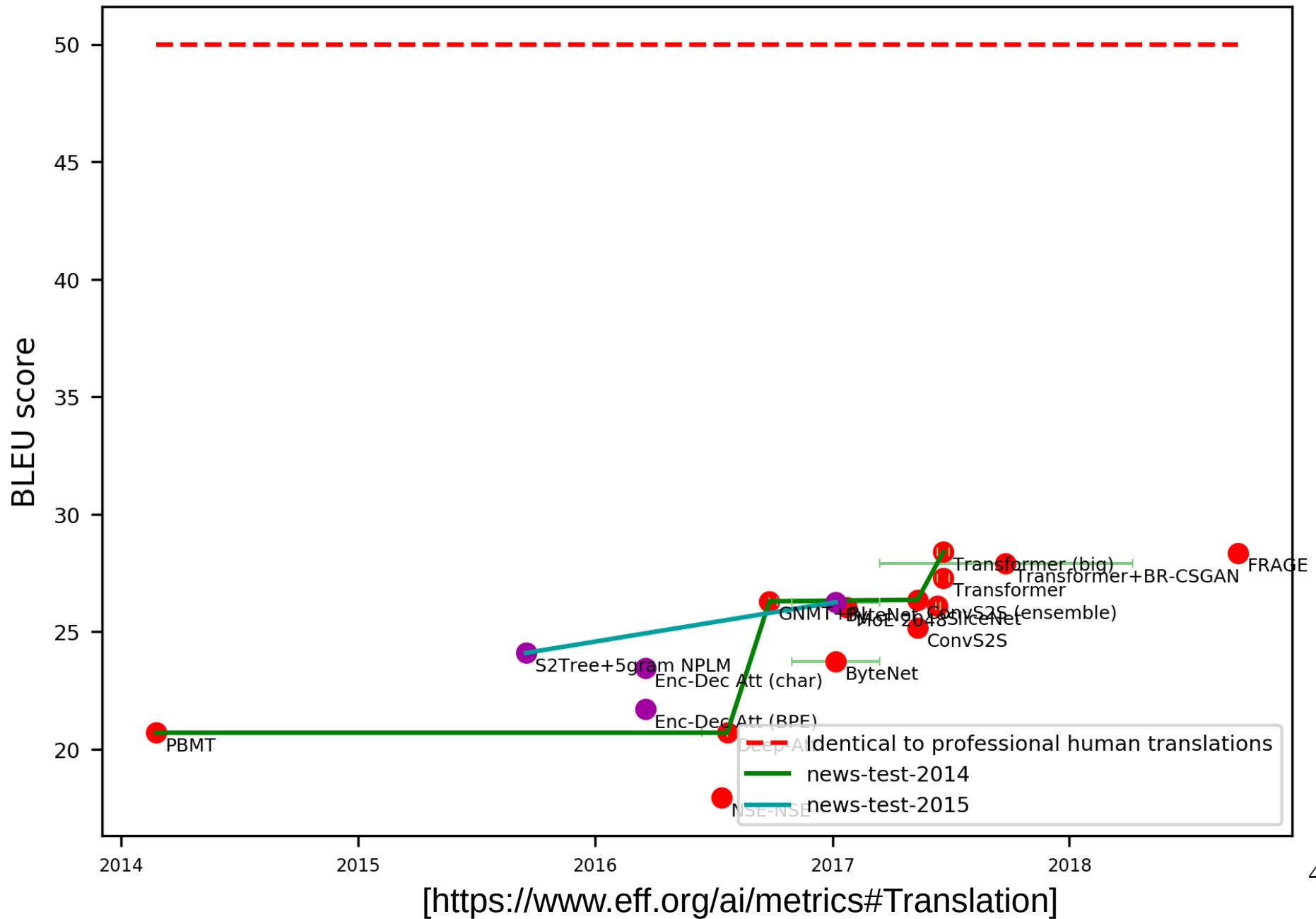
# Approaches

- 1950s: Russian → English MT (Cold War)
  - Rule-based systems using bilingual dictionaries
- 1990s-2010s: **Statistical** MT (SMT)
  - Mainly **Phrase-Based** MT (PBMT)
  - Learns from (large) **sentence-aligned** bilingual corpora
    - Millions of sentence pairs
  - Consists of separate very complex components, learnt separately
  - developed by large groups for decades, doesn't generalize to new language pairs
- Since 2014: **Neural** MT (NMT)
  - Single NN learnt **end-to-end**
  - Learns from (large) **sentence-aligned** bilingual corpora
  - A few (good) student / month to implement
  - better quality of translation than SMT systems developed for decades

Google Translate: NMT, Yandex Translate: NMT+PBMT

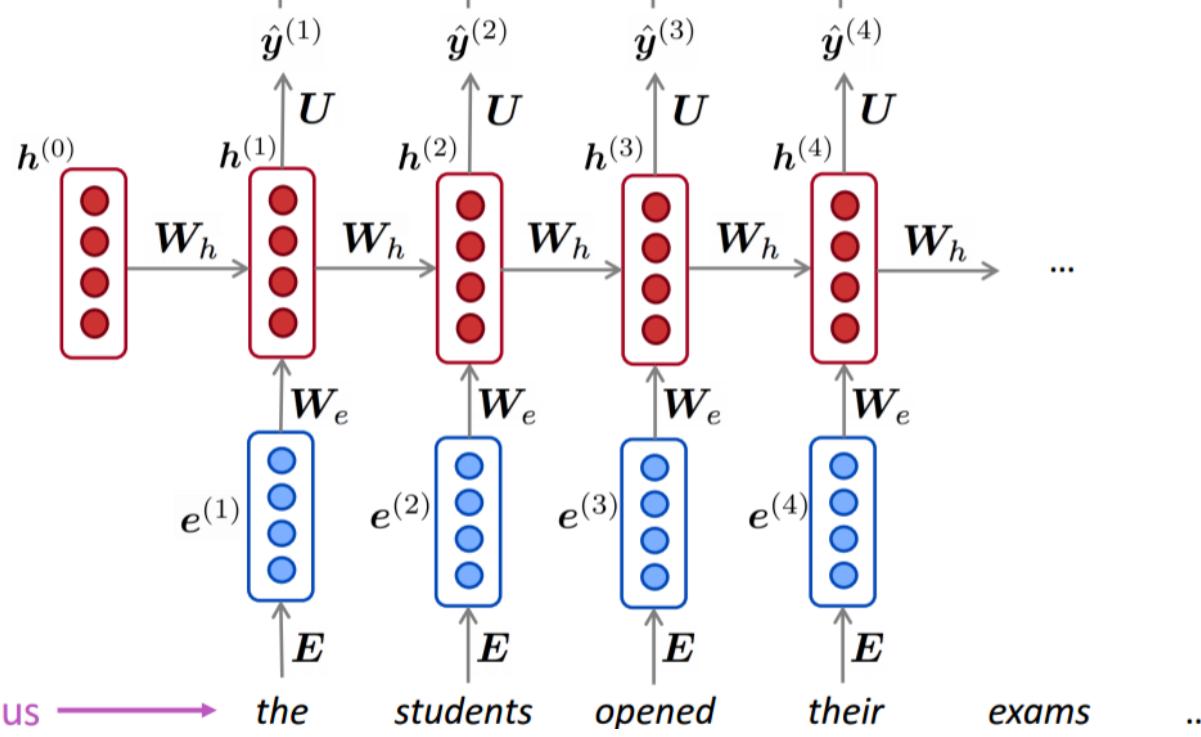
# MT vs. Human translation

En-De Translation BLEU scores



# RNN LM reminder

Loss  $\longrightarrow J^{(1)}(\theta) + J^{(2)}(\theta) + J^{(3)}(\theta) + J^{(4)}(\theta) + \dots = J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta)$



output distribution

$$\hat{y}^{(t)} = \text{softmax} (\mathbf{U} \mathbf{h}^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

hidden states

$$\mathbf{h}^{(t)} = \sigma (\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + \mathbf{b}_1)$$

$\mathbf{h}^{(0)}$  is the initial hidden state

word embeddings

$$\mathbf{e}^{(t)} = \mathbf{E} \mathbf{x}^{(t)}$$

words / one-hot vectors

$$\mathbf{x}^{(t)} \in \mathbb{R}^{|V|}$$

# RNN/LSTM LM reminder

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

- estimate  $P(y)$
- or generate new  $y$ :

$y_1$  = random word / <START>

for  $i$  from 2 to  $T$ :

$y_i = \operatorname{argmax}_y P(y|y_1 \dots y_{i-1})$  / sample from  $P(y|y_1 \dots y_{i-1})$

# LSTM Reminder

- Solves (partially) vanishing gradient problem of vanilla RNNs  
=> better models long range dependencies

$$i_t = \sigma(W^{(i)}x_t + U^{(i)}h_{t-1})$$

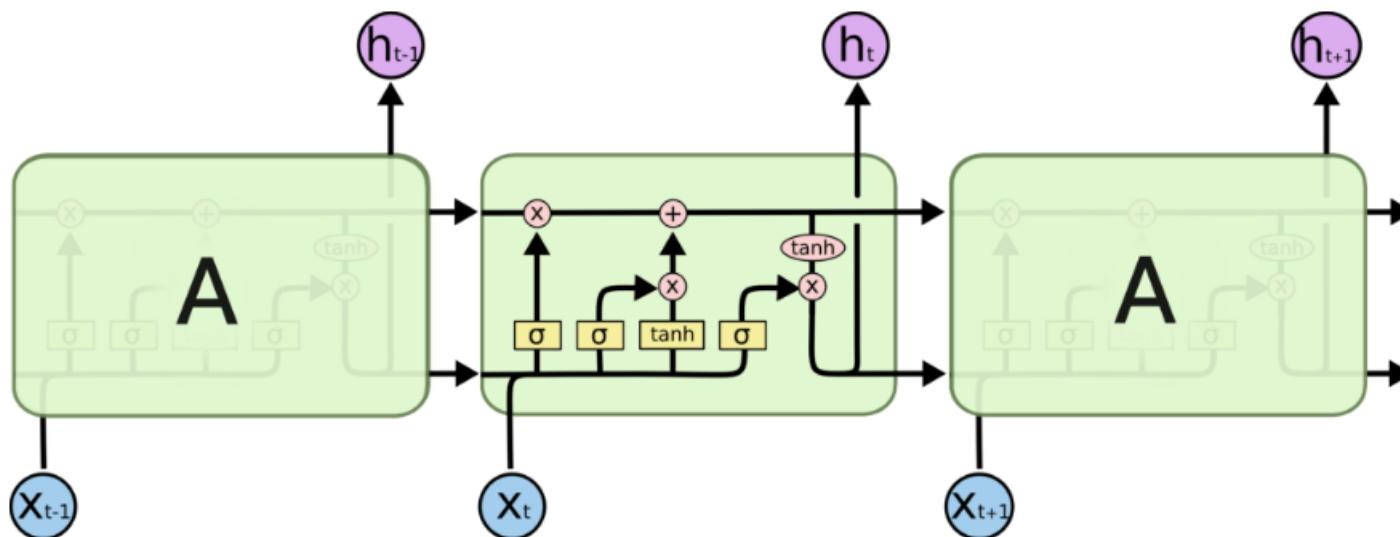
$$f_t = \sigma(W^{(f)}x_t + U^{(f)}h_{t-1})$$

$$o_t = \sigma(W^{(o)}x_t + U^{(o)}h_{t-1})$$

$$\tilde{c}_t = \tanh(W^{(c)}x_t + U^{(c)}h_{t-1})$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$h_t = o_t \circ \tanh(c_t)$$

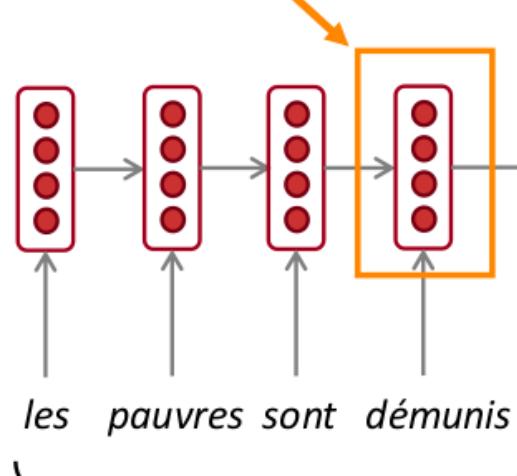


# NMT: seq2seq

The sequence-to-sequence model

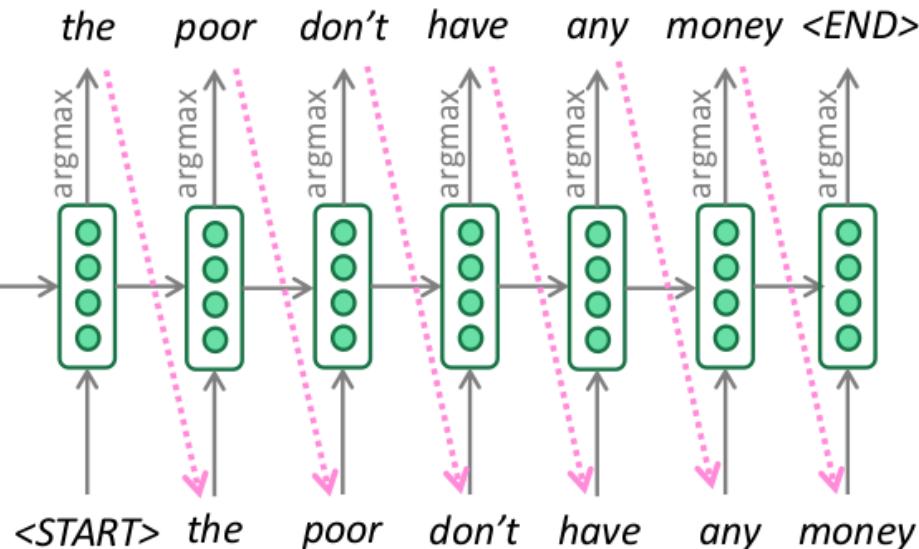
Encoding of the source sentence.  
Provides initial hidden state  
for Decoder RNN.

Encoder RNN



Encoder RNN produces  
an encoding of the  
source sentence.

Target sentence (output)



Decoder RNN

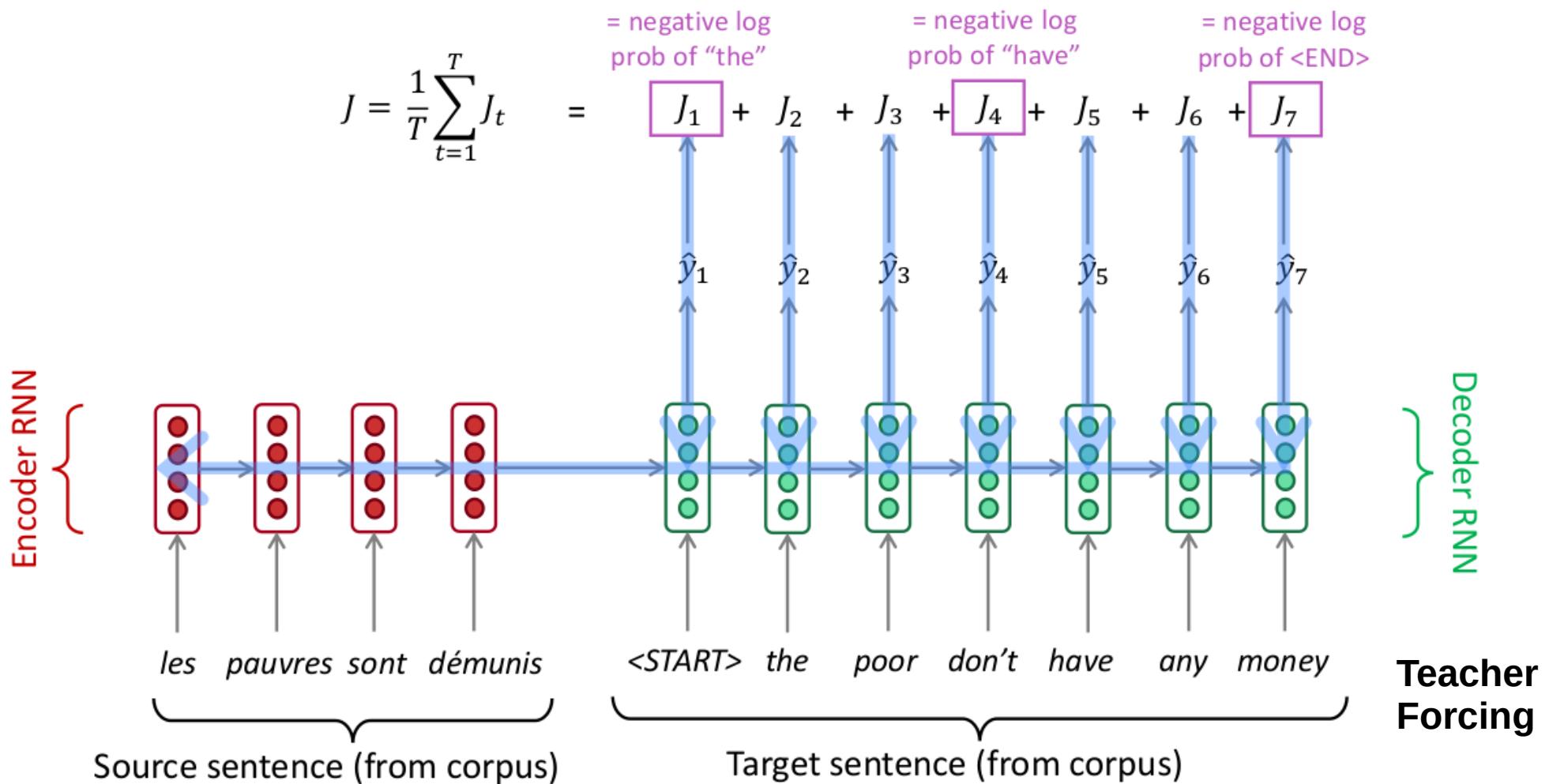
Greedy  
Decoding

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

# NMT: seq2seq loss

## Training a Neural Machine Translation system



# Conditional LM

- Decoder - conditional language model:

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

- Machine translation
  - Image captioning
  - Speech recognition
- The whole system (encoder + decoder) is trained “**end-to-end**”
  - Calculate gradients of loss w.r.t. parameters of both encoder and decoder & do SGD

# Sutskever 2014. Tricks

Sutskever, Vinyals, Le. Sequence to sequence learning with neural networks, 2014

- **reversed input:** shorter path btw.  $y_1$  and  $x_1$
- **stacked LSTMs: 4 layers,** 1000 cells per layer
- different LSTMs for encoder and decoder (different weights)
- **beam-search decoder**
- 1000 dim embeddings

# Sutskever 2014. Training

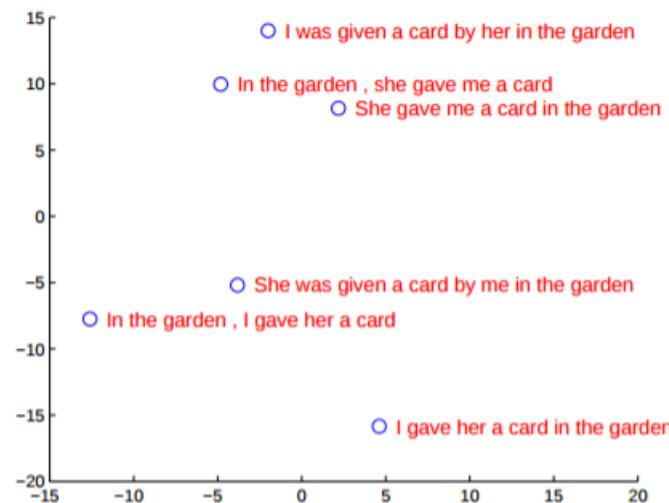
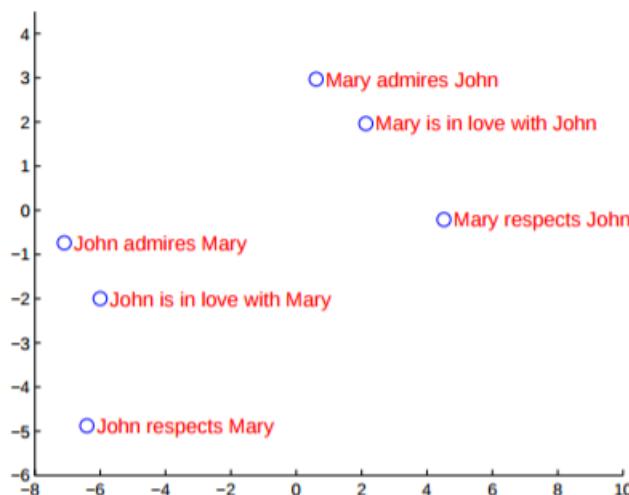
Sutskever, Vinyals, Le. Sequence to sequence learning with neural networks, 2014

- SGD without momentum, lr=0.7, train for 7.5 epochs, after 5 epochs halve lr every 0.5 epoch
- Bs=128, balanced by length => 2x speedup!
- Gradient clipping:  $g := 5 * g / \|g\|$
- parallel En → Fr corpora: 12M sentence pairs
  - vocabulary: 160K/80K most frequent words + UNK
- Custom C++ impl on 8 GPUs – 10 days
  - 1GPU per LTSM layer + 4GPUs for softmax

# Sutskever 2014. Results

Sutskever, Vinyals, Le. Sequence to sequence learning with neural networks, 2014

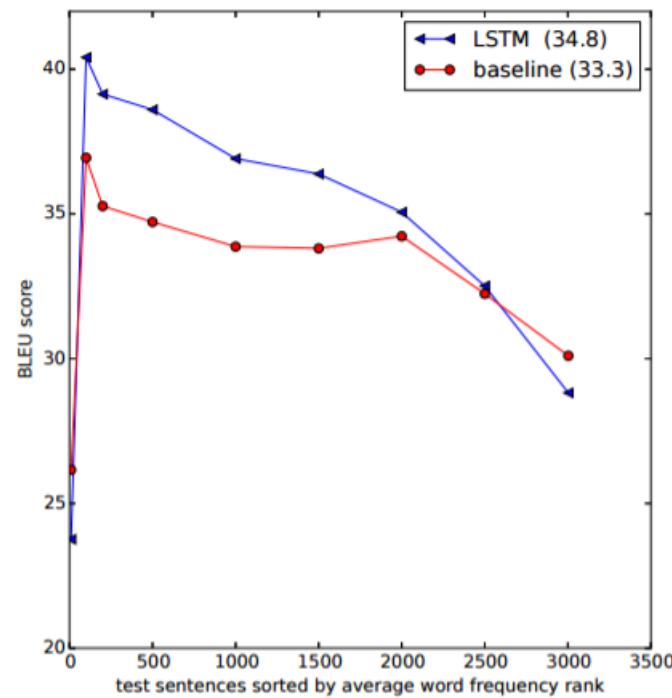
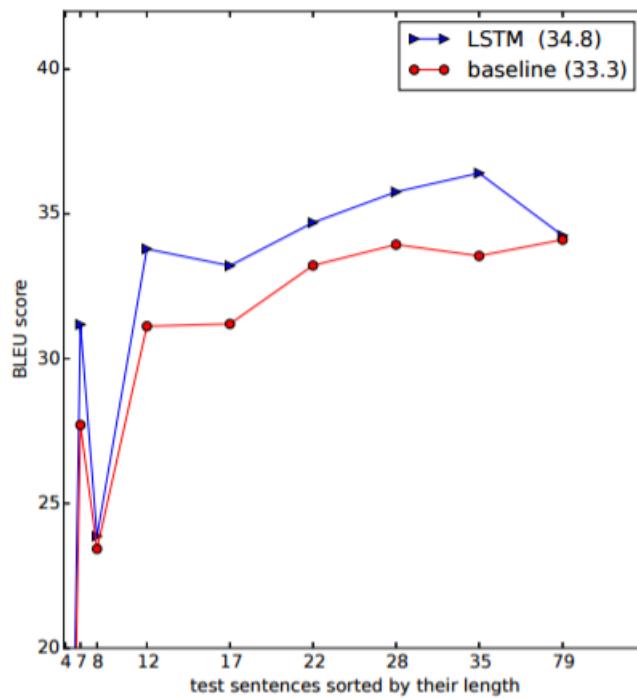
- Little worse than SOTA SMT system, but much simpler and easier to develop
- LSTM hidden state represents “thought”?
  - PCA of LSTM hidden states:



# Results

Sutskever, Vinyals, Le. Sequence to sequence learning with neural networks, 2014

- no problems with long sentences (!)
- but problems with rare words



# Decoding

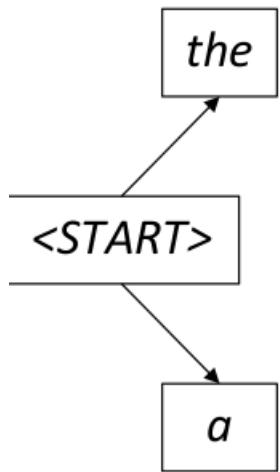
- Decoding task: search for sequence  $y$  such as  
$$y_1 y_2 \dots y_T = \operatorname{argmax}_y P(y|x) = \operatorname{argmax}_y P(y_1|x)P(y_2|y_1, x)\dots P(y_T|y_1, \dots, y_{T-1}, x)$$
- Cannot enumerate all possible  $y$  (do full search)
  - $V^T$ , where  $V=100K-1M$ ,  $T=10-100$
- Use heuristic search (approximate argmax):
  - Ancestral / forward sampling:  $y_i \sim P(y_i | y_1, \dots, y_{i-1})$
  - Greedy search:  $y_i = \operatorname{argmax}_y P(y_i | y_1, \dots, y_{i-1})$
  - Beam search  $\leftarrow$  most common for MT

# Greedy Decoding

- Greedy decoding has no way to undo decisions!
  - *les pauvres sont démunis (the poor don't have any money)*
  - → *the* \_\_\_\_
  - → *the poor* \_\_\_\_
  - → *the poor are* \_\_\_\_
- Better option: use **beam search** (a search algorithm) to explore *several* hypotheses and select the best one

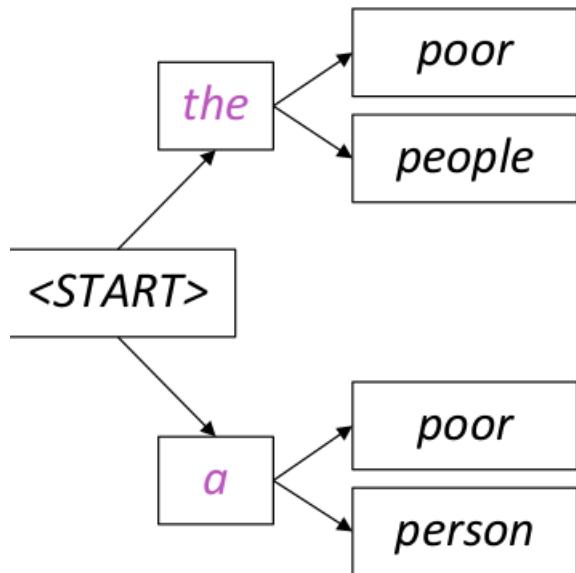
# Beam search

Beam size = 2



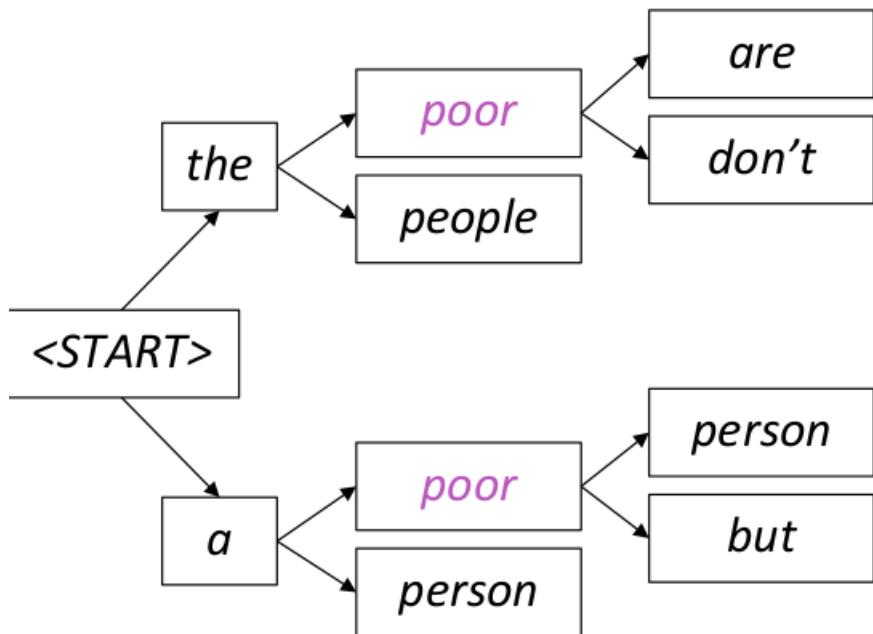
# Beam search

Beam size = 2



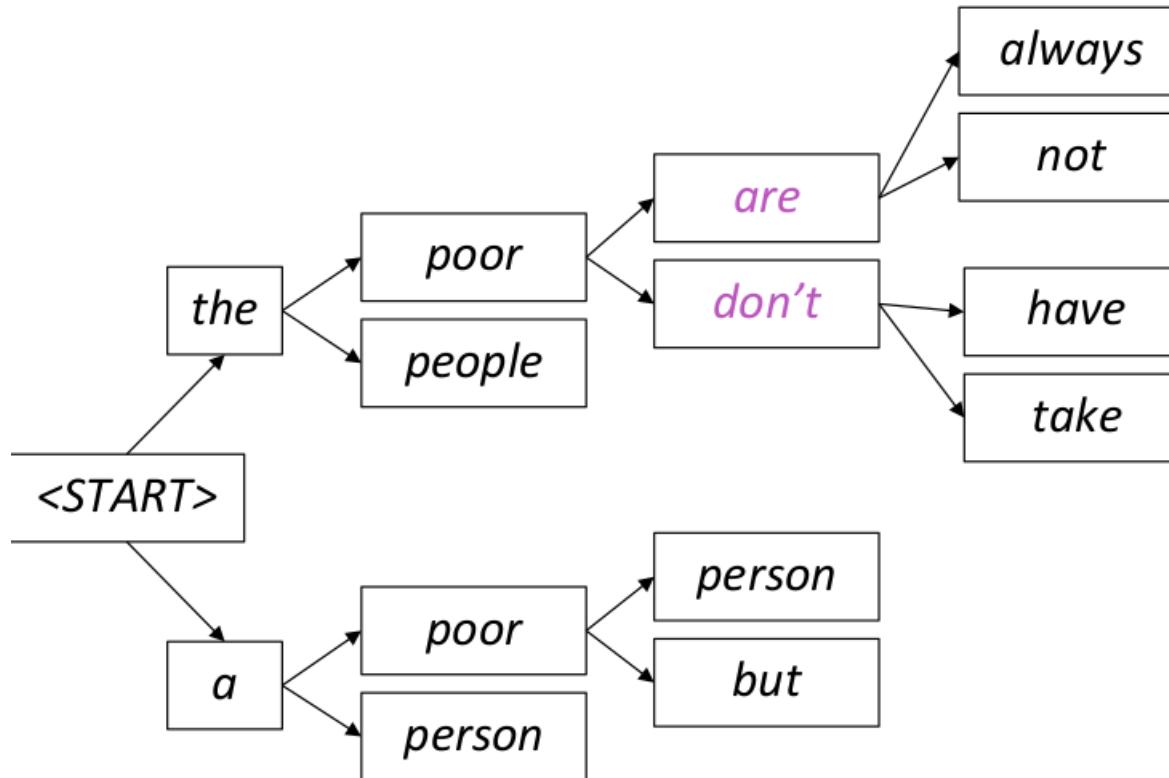
# Beam search

Beam size = 2



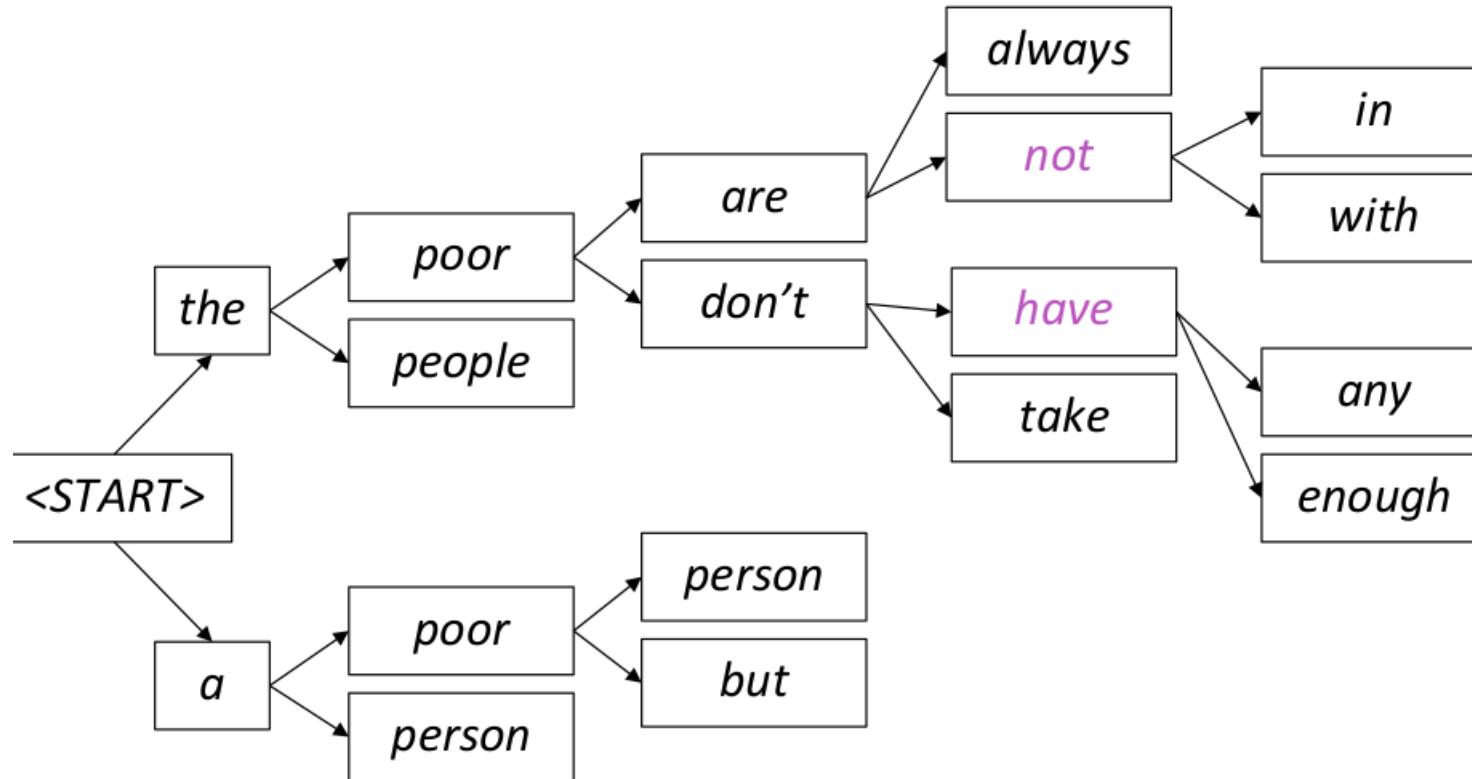
# Beam search

Beam size = 2



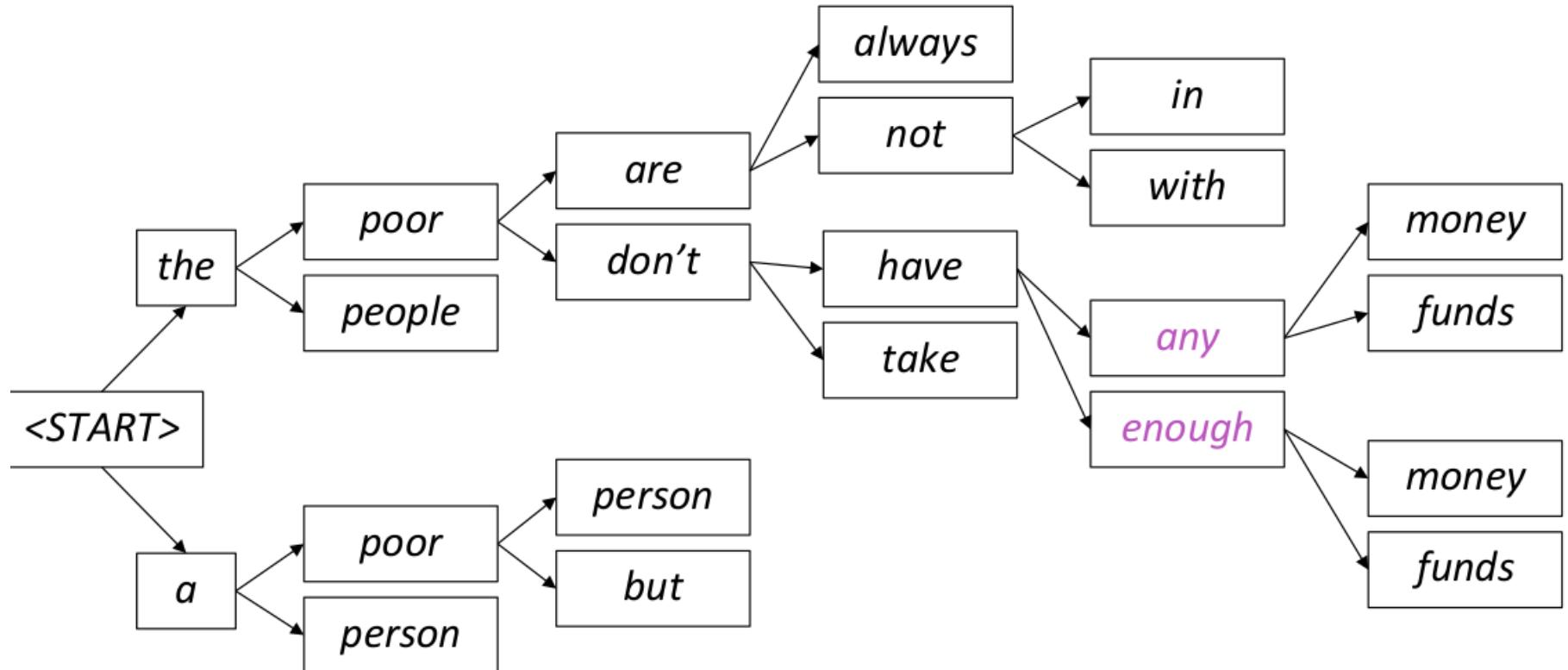
# Beam search

Beam size = 2



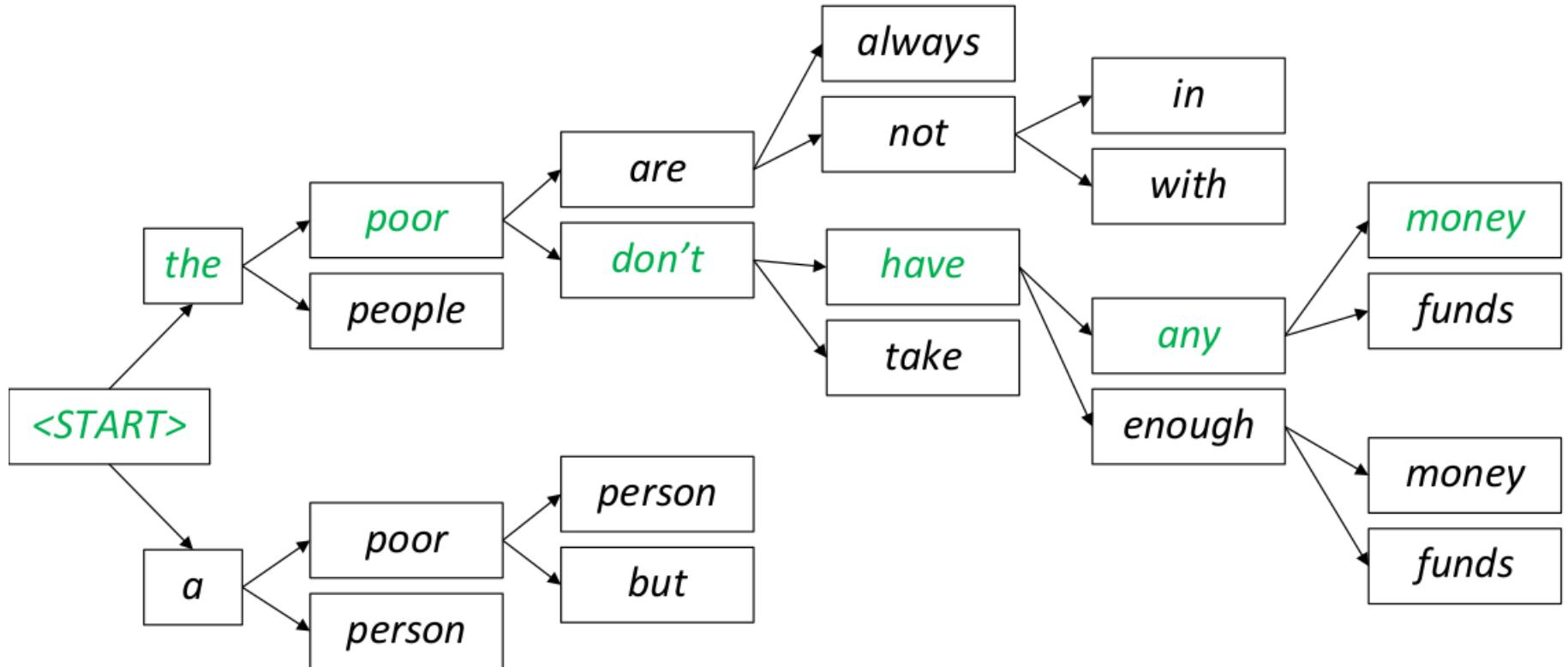
# Beam search

Beam size = 2



# Beam search

Beam size = 2



# Beam search

- Ideally we want to find  $y$  that maximizes

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

- We could try enumerating all  $y \rightarrow$  too expensive!
  - Complexity  $O(V^T)$  where  $V$  is vocab size and  $T$  is target sequence length
- Beam search: On each step of decoder, keep track of the  $k$  most probable partial translations
  - $k$  is the beam size (in practice around 5 to 10)
  - Not guaranteed to find optimal solution
  - But much more efficient!

# Length bias in NMT

- Maximizing  $s(e) = P(e|source)$ :  $s(e) = \sum_{i=1}^m \log P(e_i | e_{1:i})$   
results in the shorter translations be preferred to longer translations.
- Length normalization:  $s'(e) = s(e) / m$
- Google NMT:  $s'(e) = s(e) \left/ \frac{(5+m)^\alpha}{(5+1)^\alpha} \right.$
- Word reward  $s'(e) = s(e) + \gamma m$

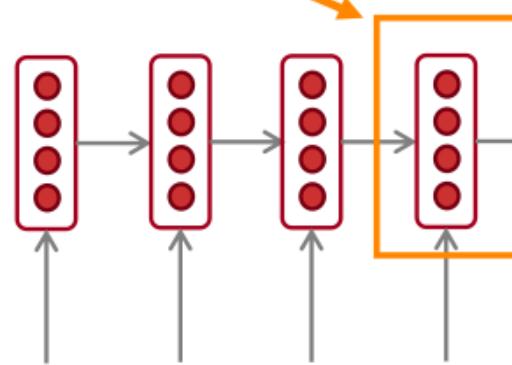
# Seq2seq bottleneck

Encoding of the source sentence.

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

Information bottleneck!

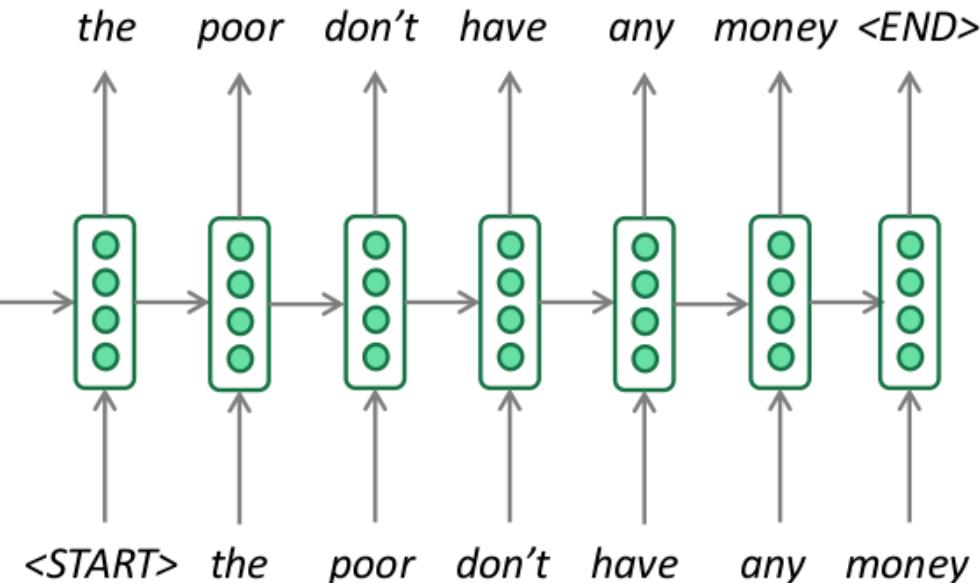
Encoder RNN



*les pauvres sont démunis*

Source sentence (input)

Target sentence (output)

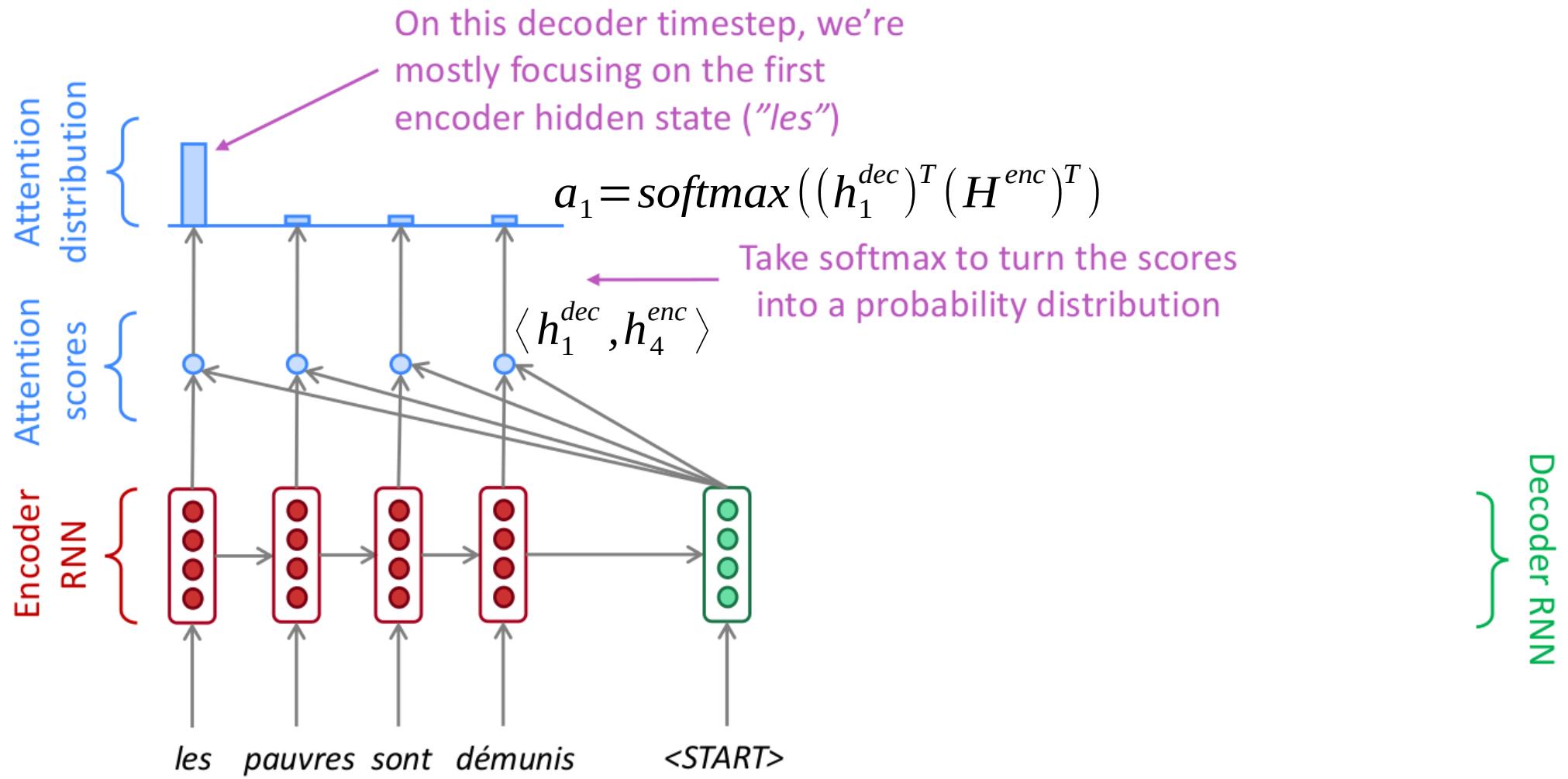


Decoder RNN

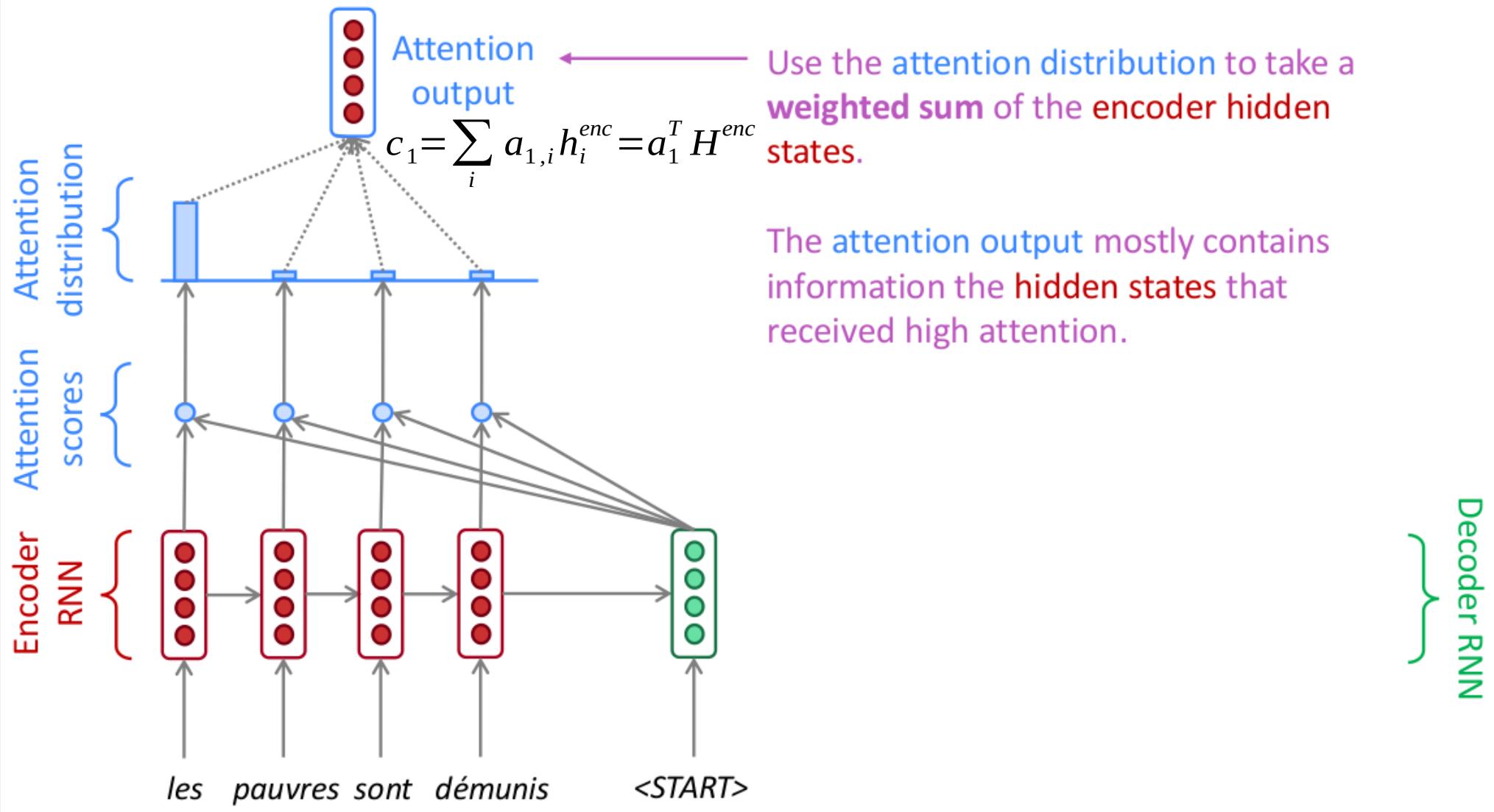
# Attention: Idea

- At each timestep (when generating next word) the decoder should focus on the **most relevant (to this timestep) part** of input sequence
- Hidden state of the decoder at each timestep represents a “rough idea” of what it wants to say next
  - Help it to “put this idea into words”: find a part of input sequence relevant to this “idea” and “remind” this part to him

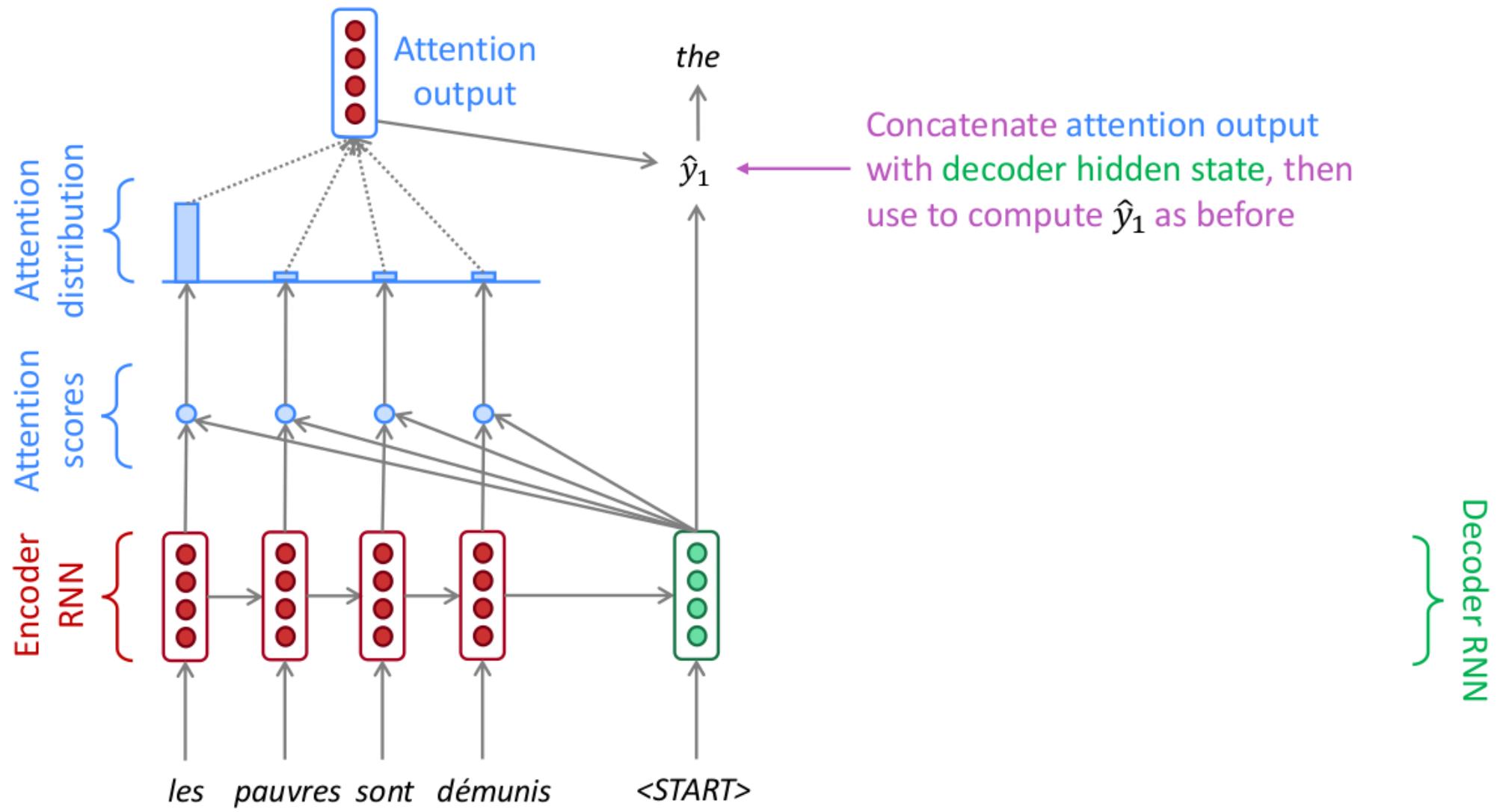
# Attention: Example



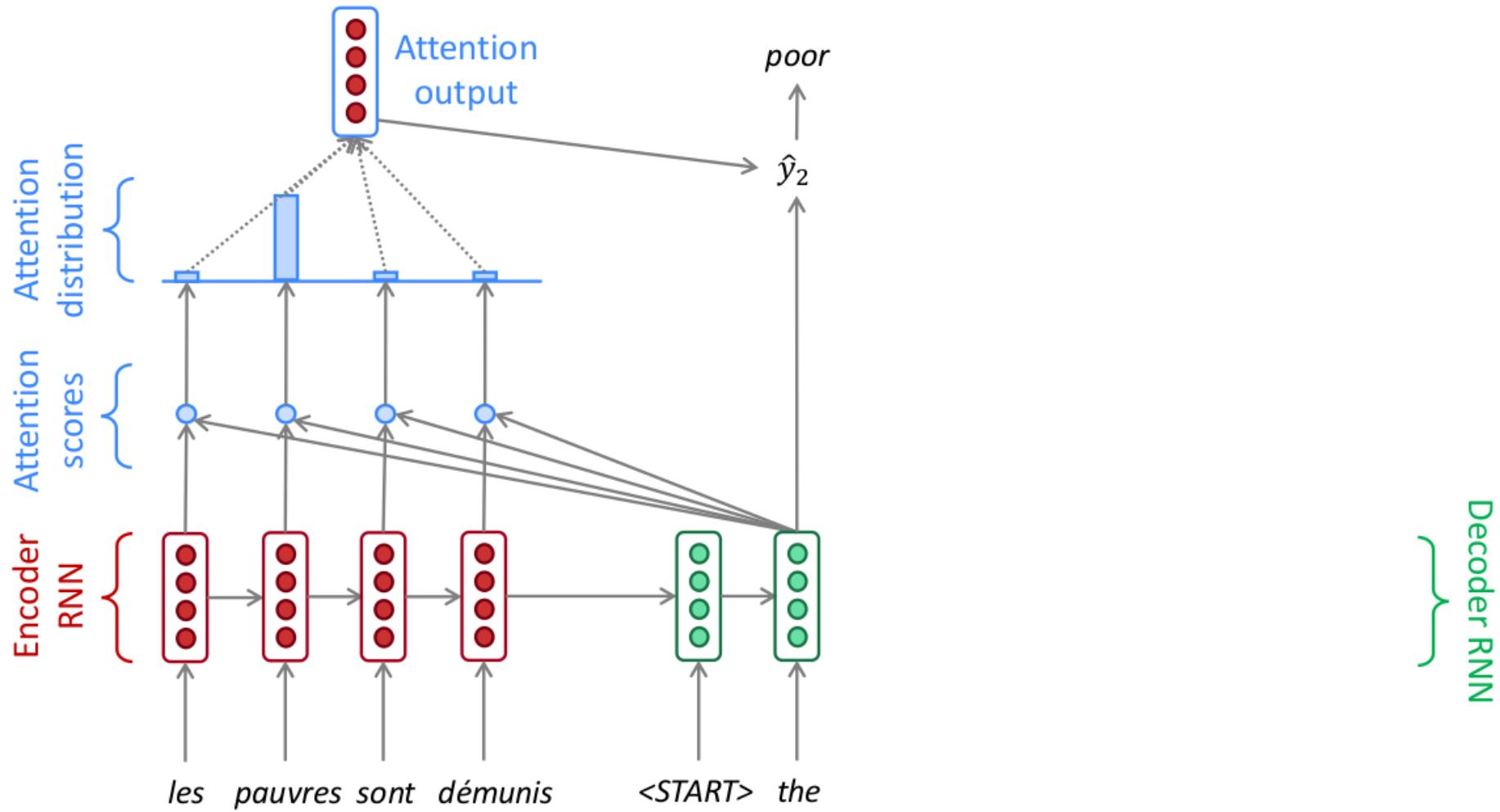
# Attention: Example



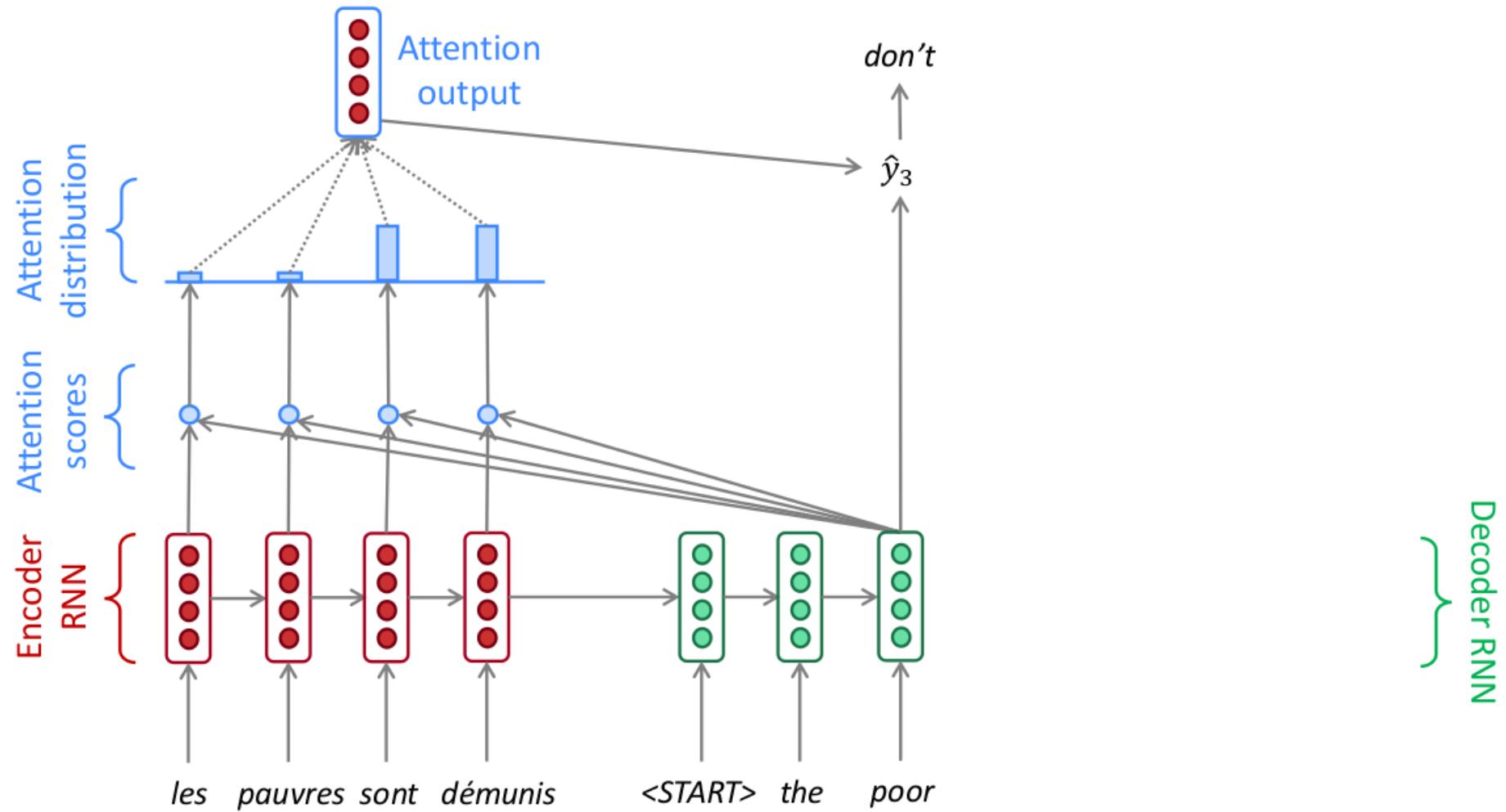
# Attention: Example



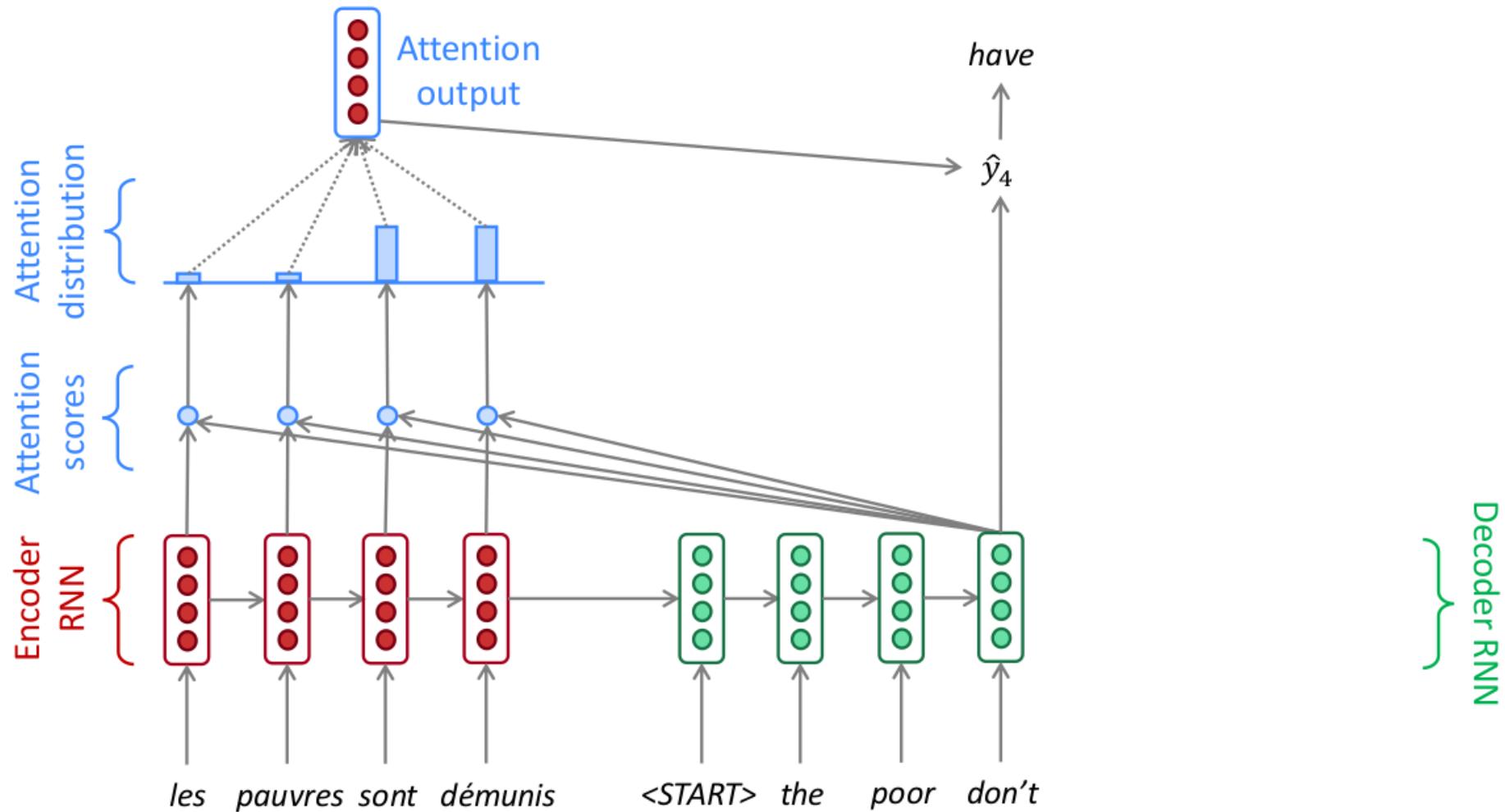
# Attention: Example



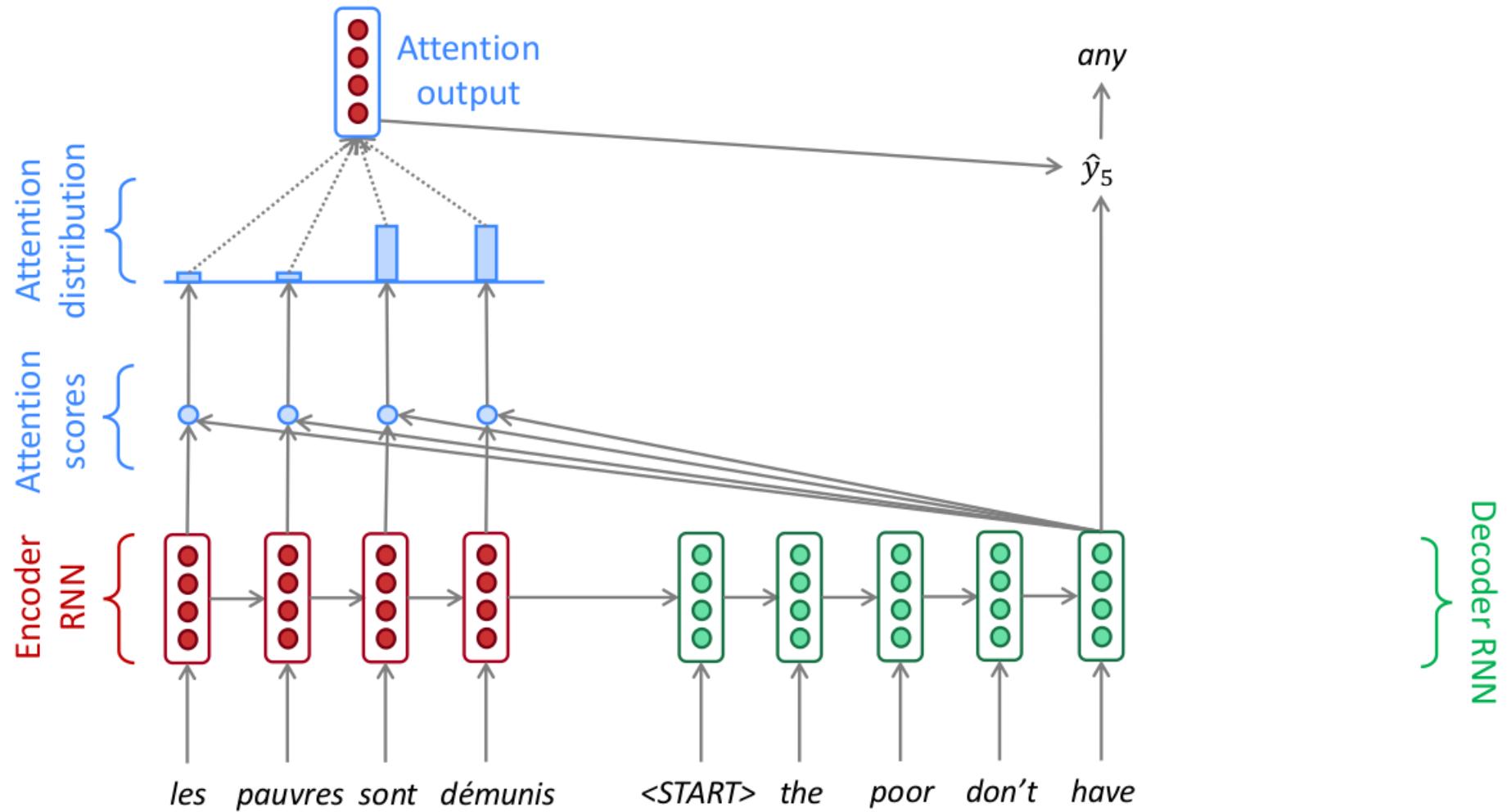
# Attention: Example



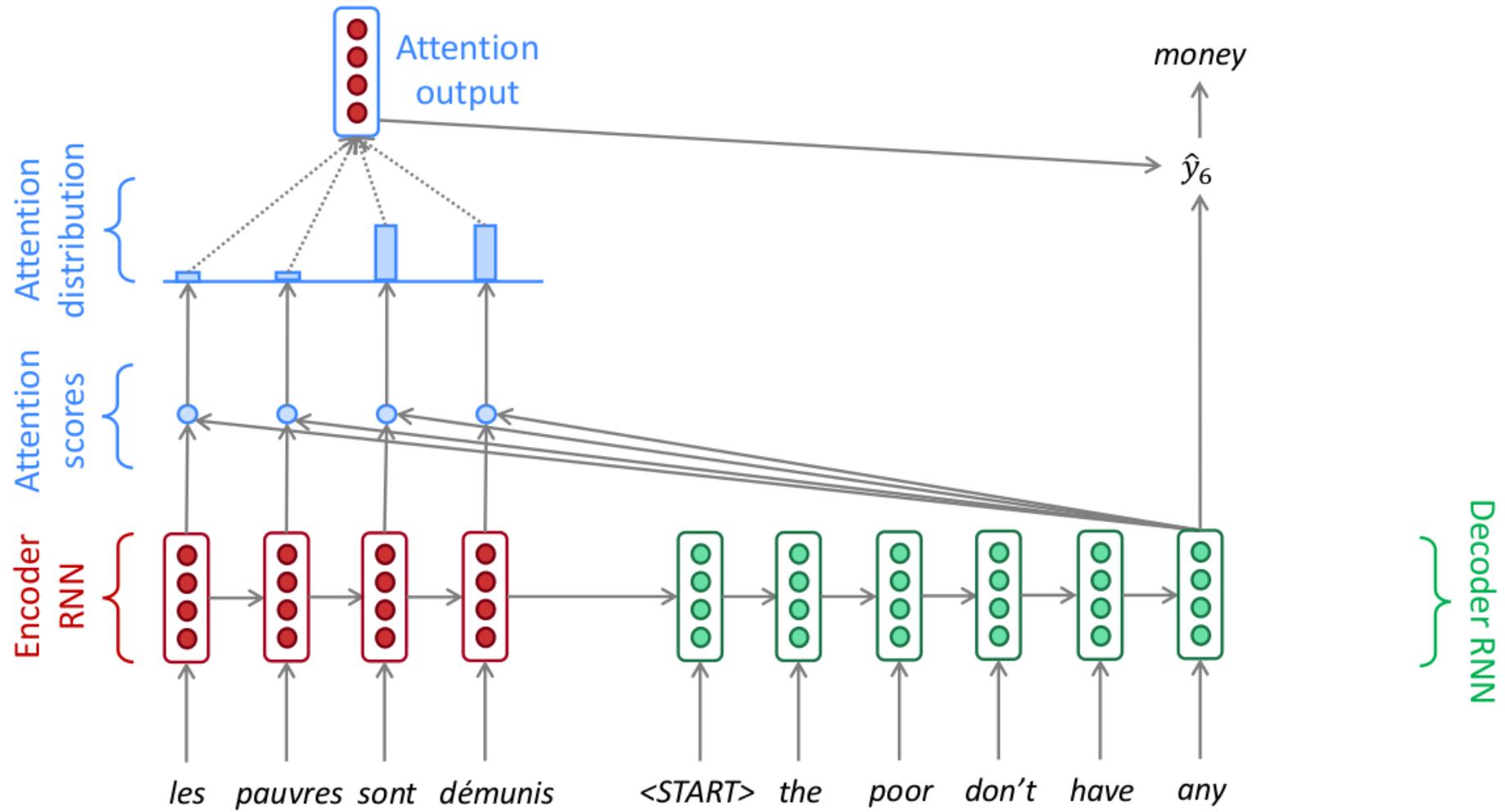
# Attention: Example



# Attention: Example



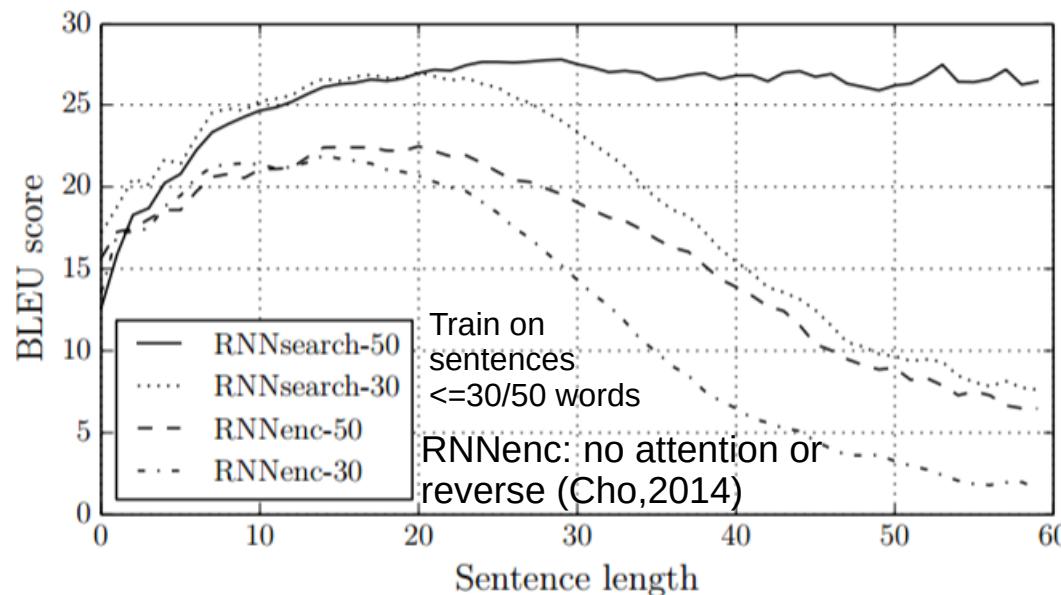
# Attention: Example



# Bahdanau 2015

Dzmitry Bahdanau, KyungHyun Cho, Yoshua Bengio. Neural Machine Translation by Jointly Learning to Align and Translate, 2015

- Solves the problem of low quality for long sequences
  - Alternative to reversing



# Bahdanau 2015

L'accord sur la zone économique européenne a été signé en août 1992.

The agreement on the European Economic Area was signed in August 1992.

<end>

(a)

Il convient de noter que l'environnement marin est le moins connu de l'environnement.

It should be noted that the marine environment is the least known of environments.

<end>

(b)

La destruction de l'équipement signifie que la Syrie ne peut plus produire de nouvelles armes chimiques.

Destruction of the equipment means that Syria can no longer produce new chemical weapons.

<end>

Cela va changer mon avenir avec ma famille.

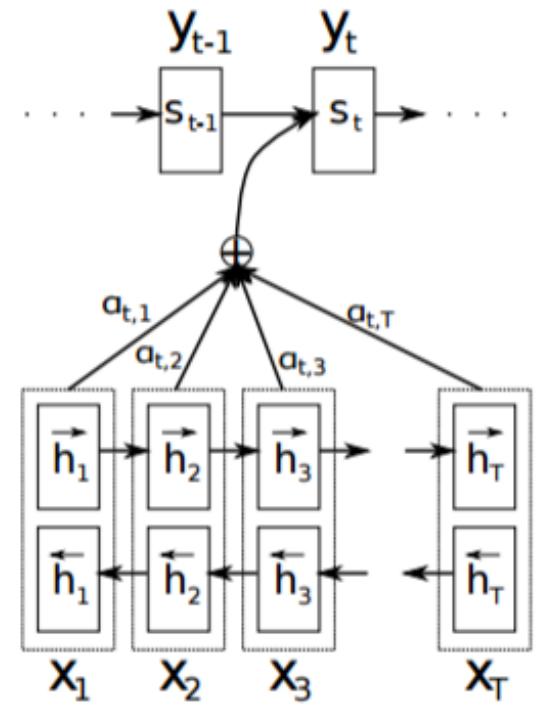
"This will change my future with my family."

the man said.

"

a dit l'homme

<end>



# Attention

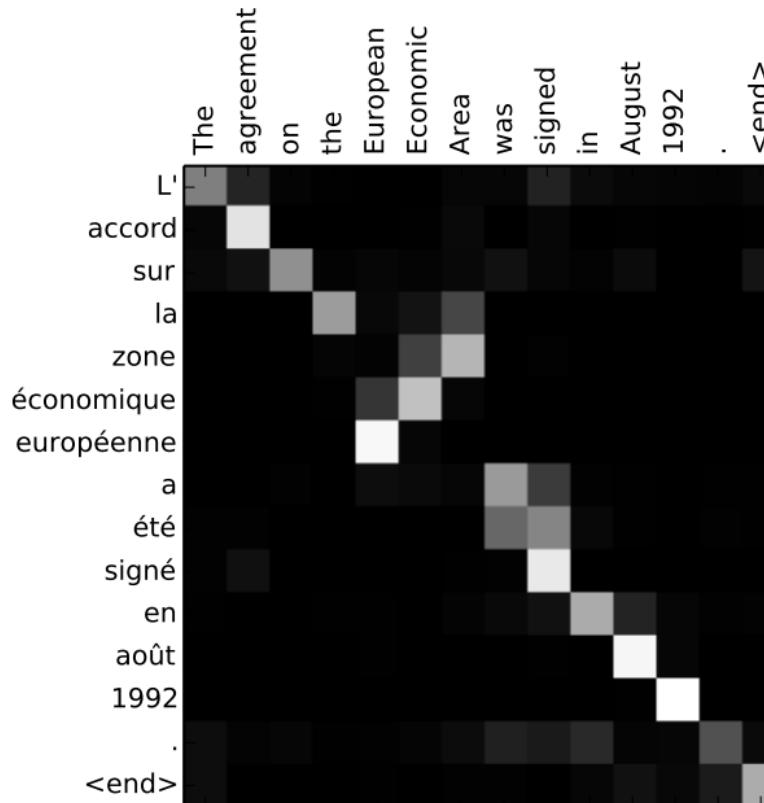
- Improves translation quality significantly
  - Solves **information bottleneck** problem, acts like an unlimited size memory
  - Helps **gradient propagation** from encoder to decoder, especially in long sequences
    - Remember vanishing gradient problem partially solved by LSTMs? But they have sigmoid forget gates! ( $<1$ )

Britz et al., 2017. Massive Exploration of Neural Machine Translation Architectures:

“we found that the attention-based models exhibited significantly larger gradient updates to decoder states throughout training. This suggests that the attention mechanism acts more like a “weighted skip connection” that optimizes gradient flow than like a “memory” that allows the encoder to access source states, as is commonly stated in the literature”

# Attention

- Improves interpretability
  - Networks learns alignment as byproduct of translation
  - We can look at which fragments were translated to which



# Attention

Attention is a new basic layer type (along with feed-forward, convolutional and recurrent)

- Works with variable length inputs (texts)
  - like RNNs, CNNs
  - unlike feed-forward
- Used in SOTA models in lots of tasks (question answering, image captioning, ...)

# Do I hold your attention?

