

Seq2seq model **and application for machine** **translation**

Nguyen Van Vinh

UET, VNU-Hanoi

Content

- Introduction to Machine Translation
- The seq2seq model
- Attention mechanism
- **Practice:** Machine translation với mô hình seq2seq

Machine Translation

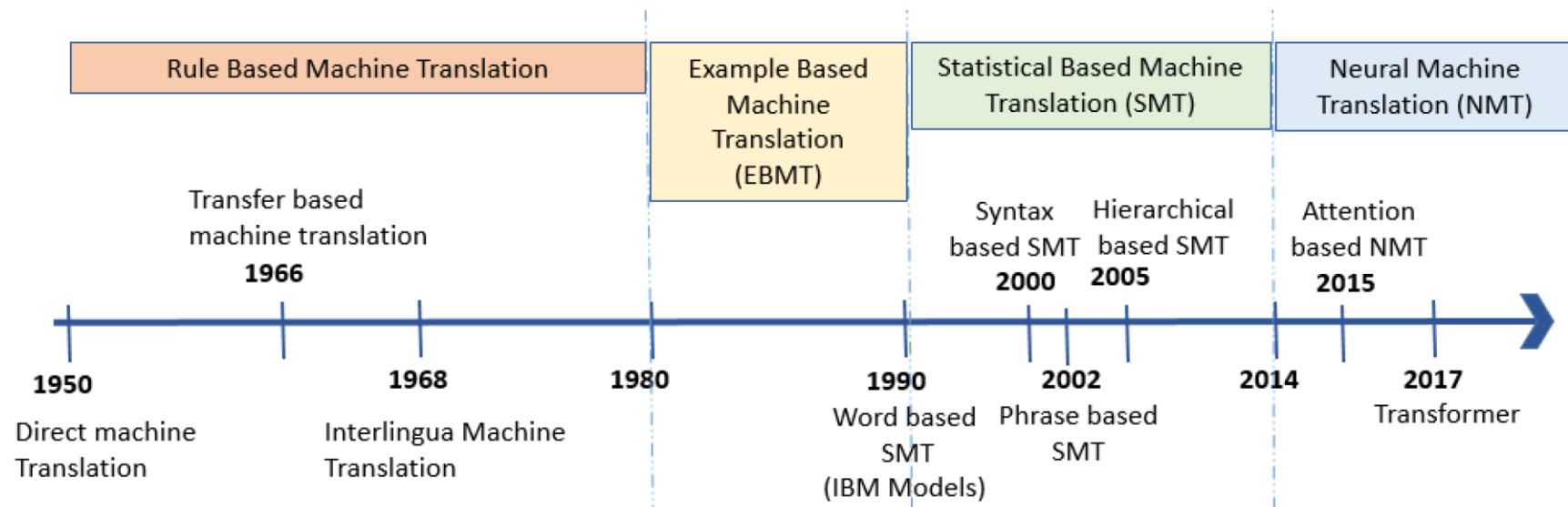
Machine Translation (MT) is the task of translating a sentence x from one language (the **source language**) to a sentence y in another language (the **target language**).

$x:$ *L'homme est né libre, et partout il est dans les fers*



$y:$ *Man is born free, but everywhere he is in chains*

History of MT



What is Neural Machine Translation?

- Neural Machine Translation (NMT) is a way to do Machine Translation with a *single neural network*

Sequence to Sequence Learning with Neural Networks

Ilya Sutskever
Google
ilyasu@google.com

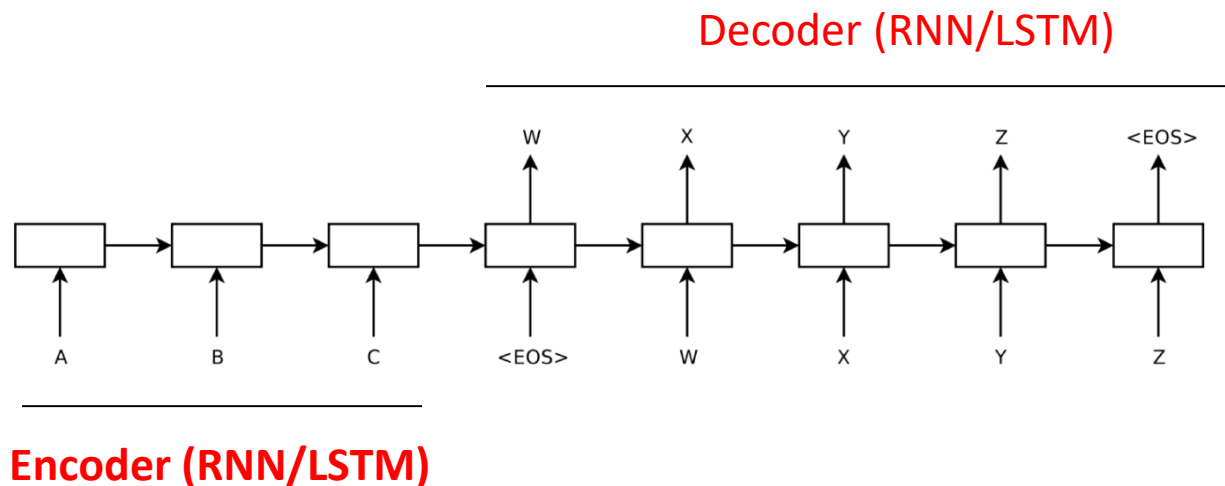
Oriol Vinyals
Google
vinyals@google.com

Quoc V. Le
Google
qvl@google.com



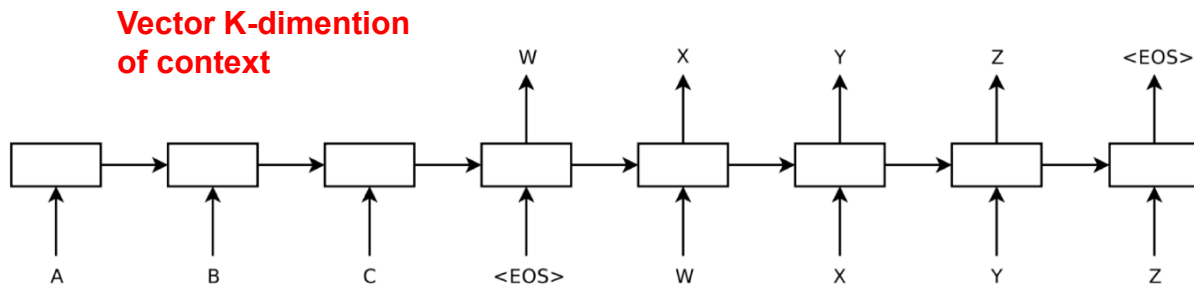
- The neural network architecture is called *sequence-to-sequence* (aka *seq2seq*) and it involves *two* RNNs (LSTMs).

Encoder-decoder Framework



Ilya Sutskever, Oriol Vinyals, Quoc V. Le, "Sequence to Sequence Learning with Neural Networks", NIPS 2014

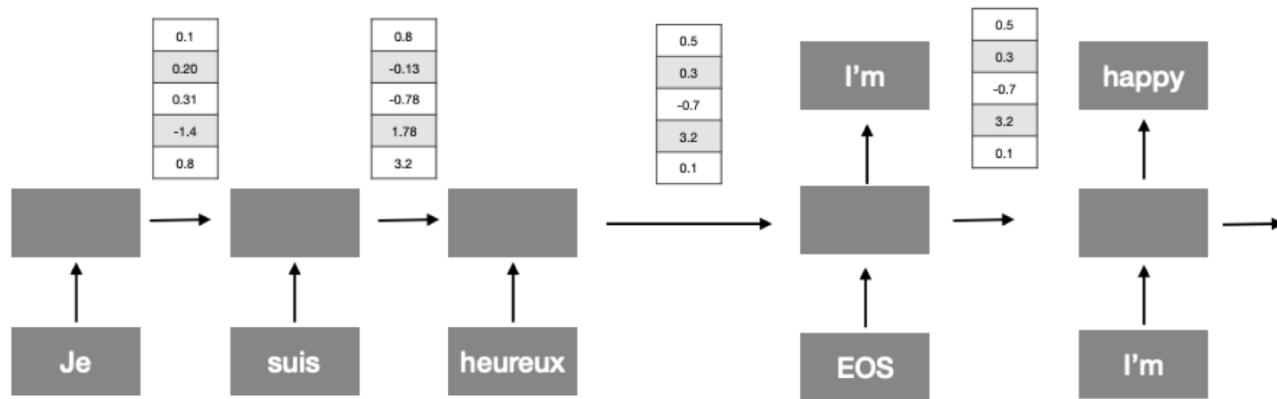
Encoder-decoder Framework



**Condition of the word to be generate from
the translation system**

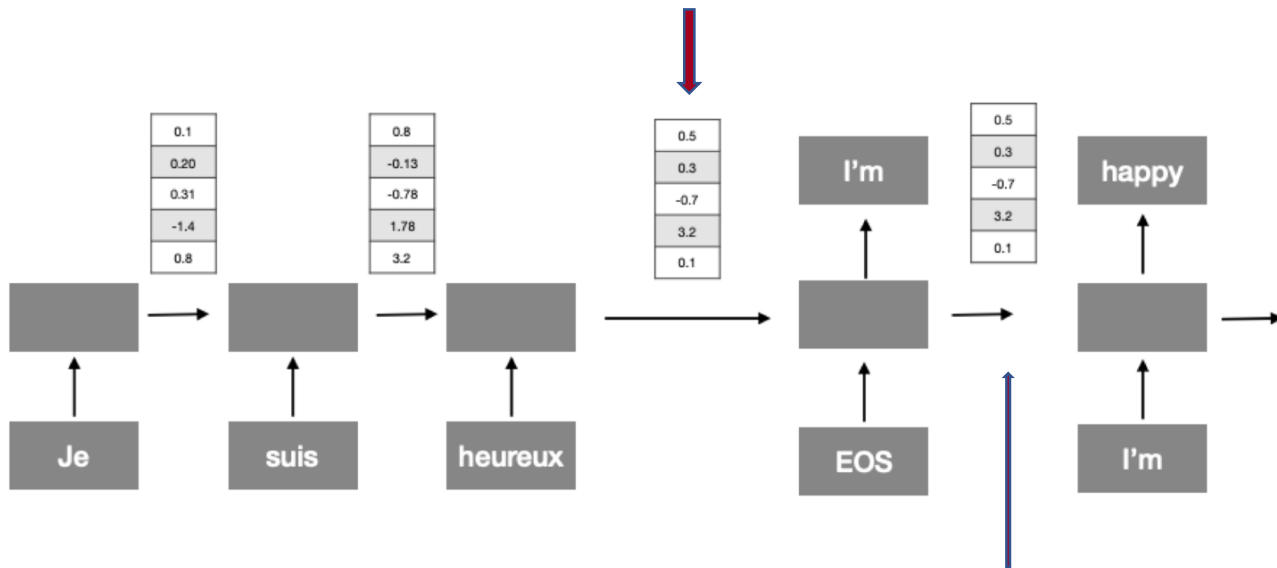
Ilya Sutskever, Oriol Vinyals, Quoc V. Le, "Sequence to Sequence Learning with Neural Networks", NIPS 2014

Encoder-decoder Framework



Encoder-decoder Framework

All inputs are aggregated in this single vector

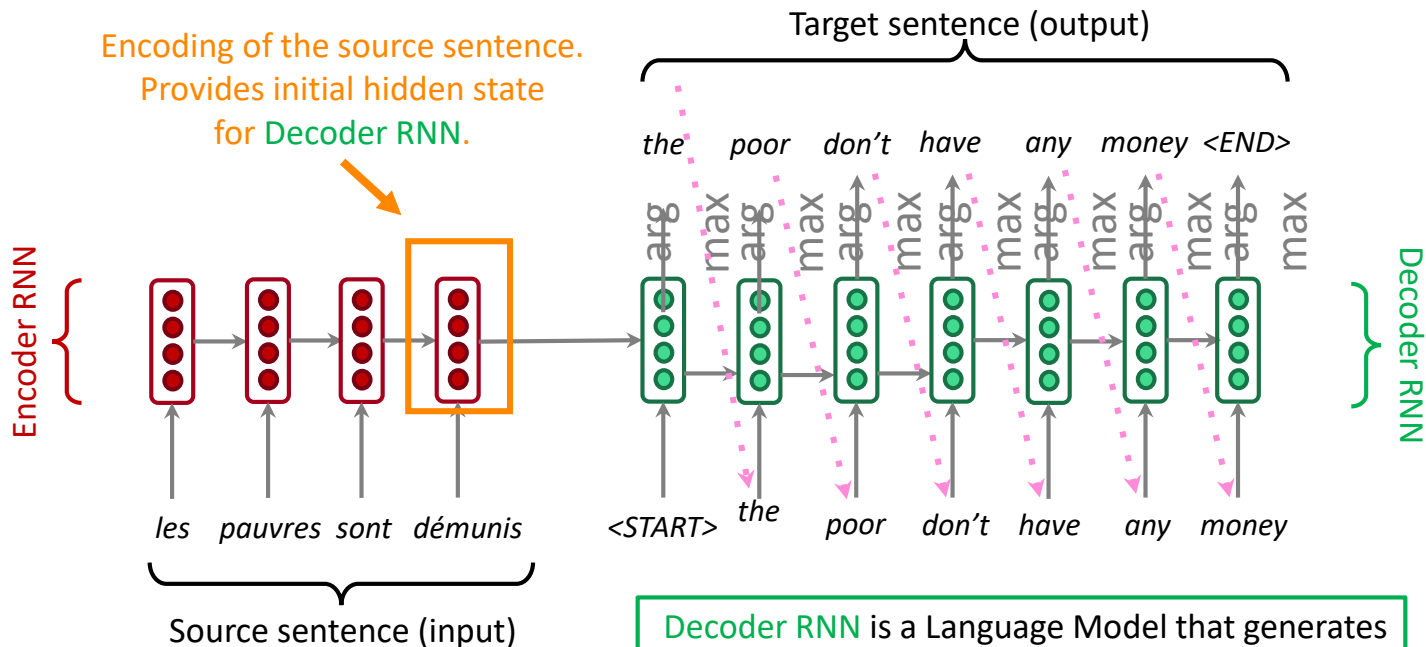


In the seq2seq, the decoder's state depends only on the previous state and the previous output

Neural Machine Translation

The sequence-to-sequence model

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



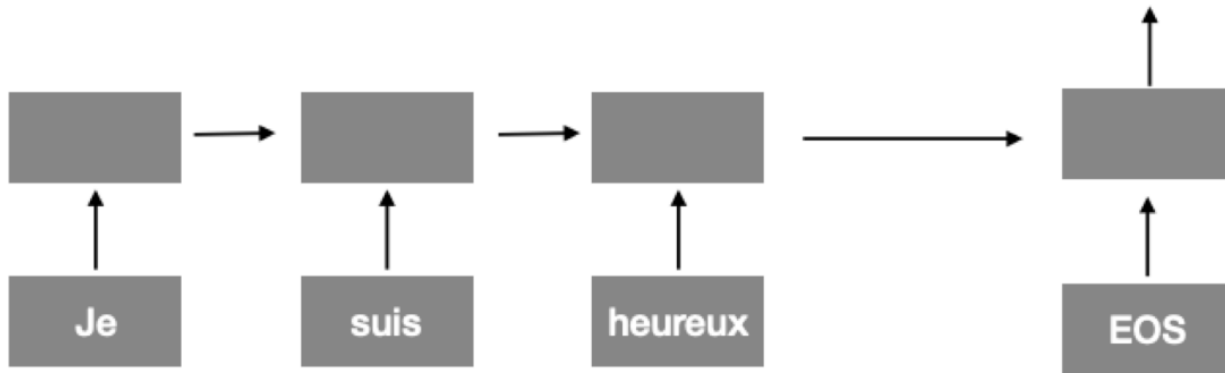
Encoder RNN produces
an **encoding** of the
source sentence.

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

Training of NMT system

- As in other RNN models, we can train by minimizing the loss function between what we predict at each step and its ground true value.



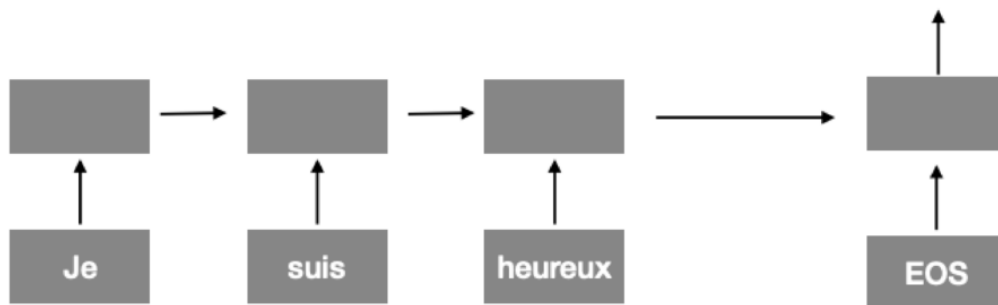
Training of NMT system

Truth

I'm	you	are	the	...
1	0	0	0	0

Predicted

I'm	you	are	the	...
0.03	0.05	0.02	0.01	0.009

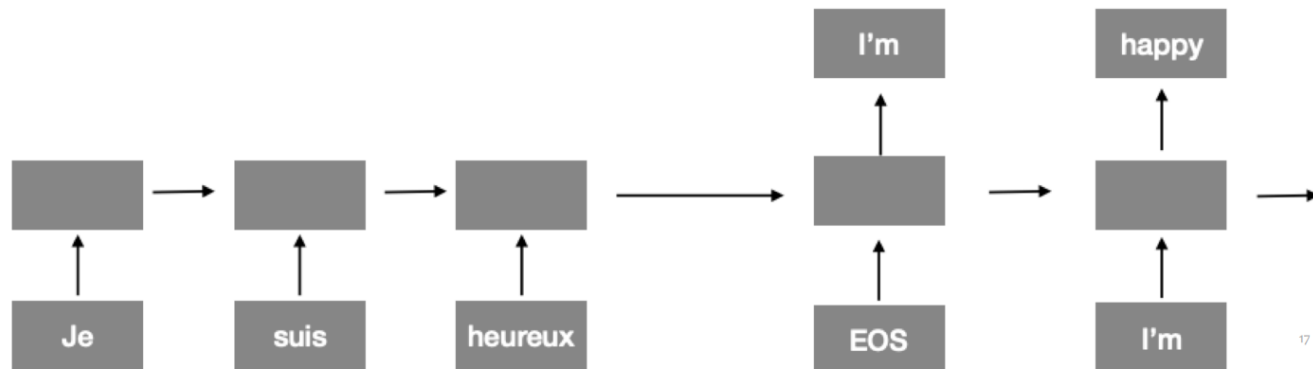


Truth

happy	great	bad	ok	...
1	0	0	0	0

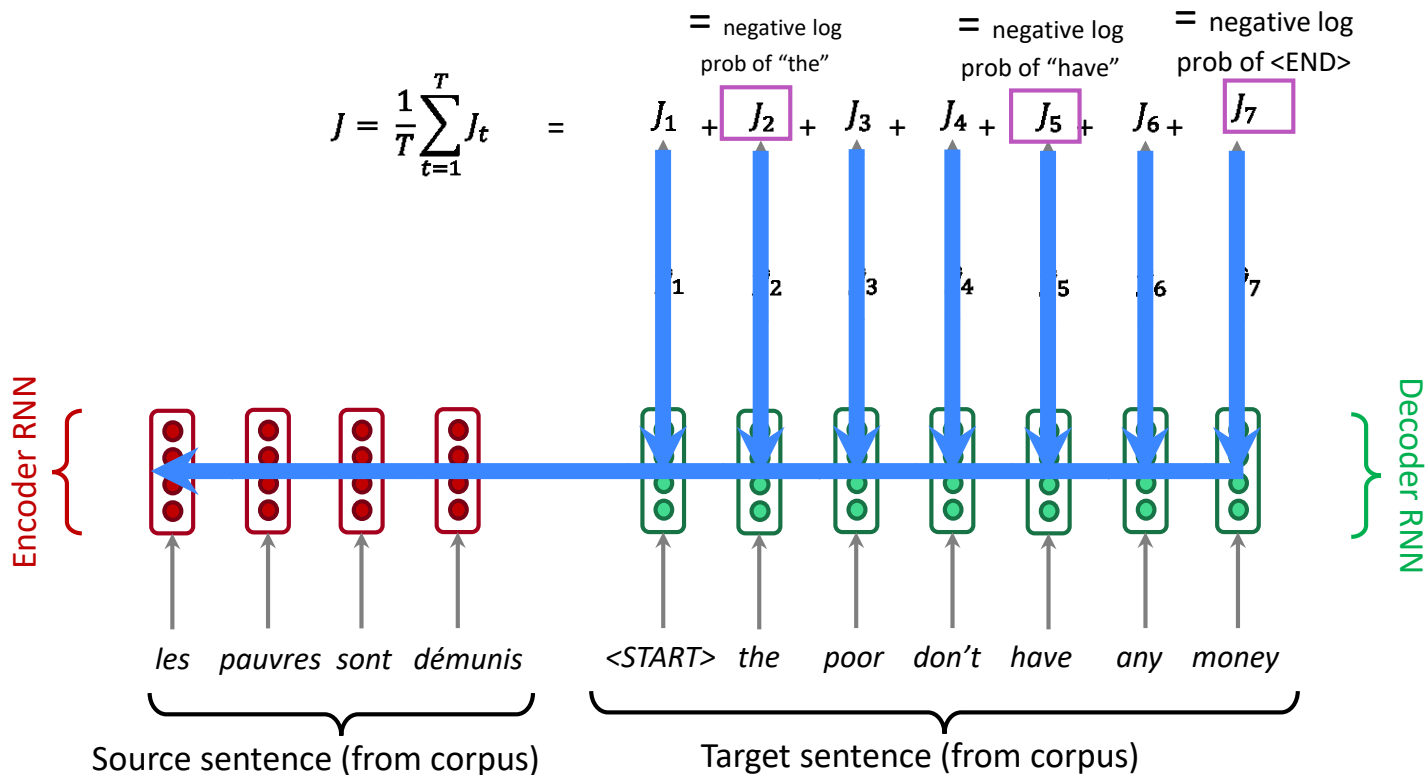
Predicted

happy	great	bad	ok	...
0.13	0.08	0.01	0.03	0.009



17

Training of NMT system



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

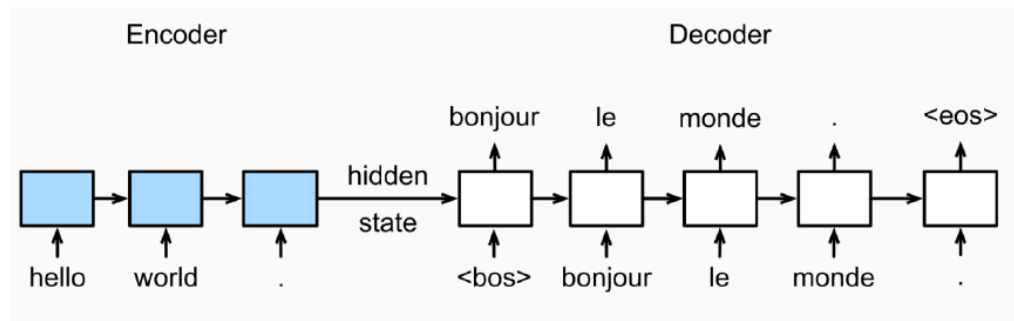
Seq2seq: Decoder

- A **conditional** language model
 - It is a **language model** because the decoder is predicting the next word of the target sentence
 - **Conditional** because the predictions are also conditioned on the source sentence through h_{src}
- NMT directly calculates $P(\mathbf{w}(t) \mid \mathbf{w}(s))$
 - Denote $\mathbf{w}(t) = y_1, \dots, y_T$

$$P(\mathbf{w}^{(t)} \mid \mathbf{w}^{(s)}) = P(y_1 \mid \mathbf{w}^{(s)})P(y_2 \mid y_1, \mathbf{w}^{(s)})P(y_3 \mid y_1, y_2, \mathbf{w}^{(s)}) \dots P(y_T \mid y_1, \dots, y_{T-1}, \mathbf{w}^{(s)})$$

$$\hat{\mathbf{y}}_t = \text{softmax}(\mathbf{W}_o \mathbf{h}_t) \quad P(y_{t+1} \mid y_1, \dots, y_t, \mathbf{w}^{(s)}) = \hat{\mathbf{y}}_t(y_{t+1})$$

Understanding seq2seq

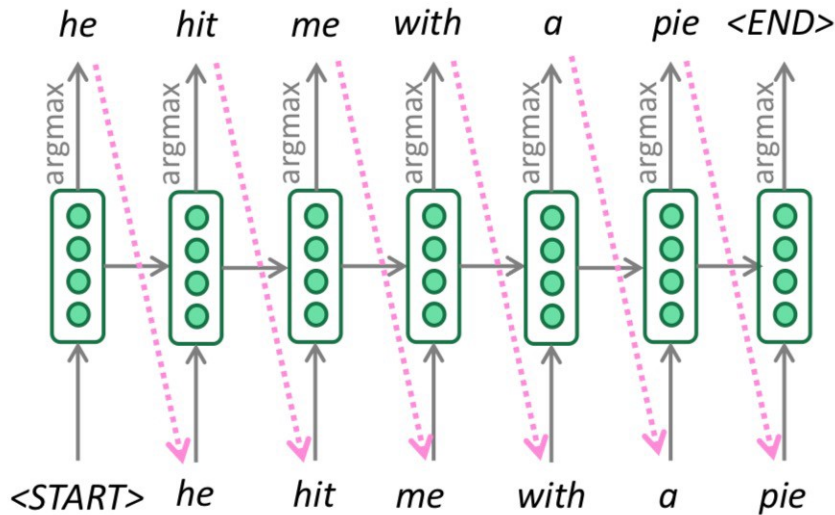


Which of the following is correct?

- (A) We can use bidirectional RNNs for both encoder and decoder
- (B) The decoder has more parameters because of the output matrix \mathbf{W}_o
- (C) The encoder and decoder have separate word embeddings
- (D) The encoder and decoder's parameters are optimized together

Decoding seq2seq models

- Greedy decoding = Compute argmax at every step of decoder to generate word



- Exhaustive search is very expensive: $\arg \max_{y_1, \dots, y_T} P(y_1, \dots, y_T | \mathbf{w}^{(s)})$ we even don't know what T is

Better-than-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

Decoder based on 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!

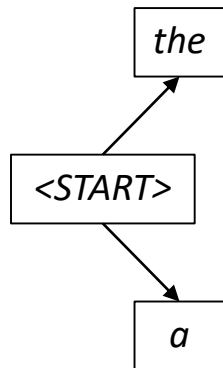
Decoder based on Beam search: Example

Beam size = 2

<START>

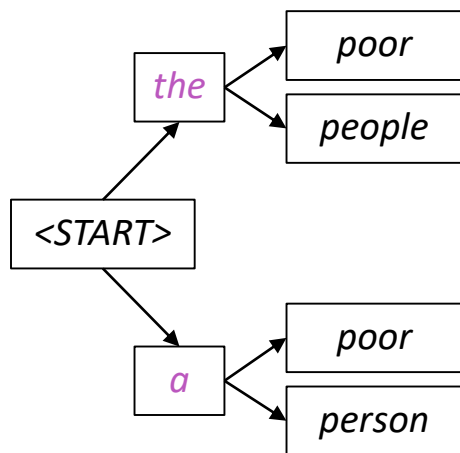
Decoder based on Beam search: Example

Beam size = 2



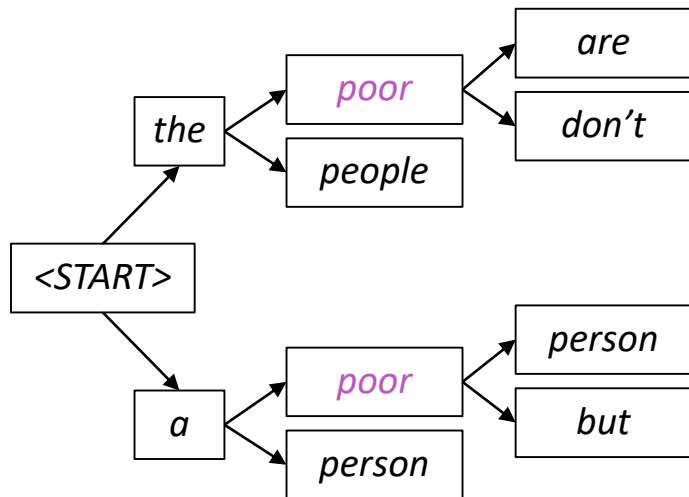
Decoder based on Beam search: Example

Beam size = 2



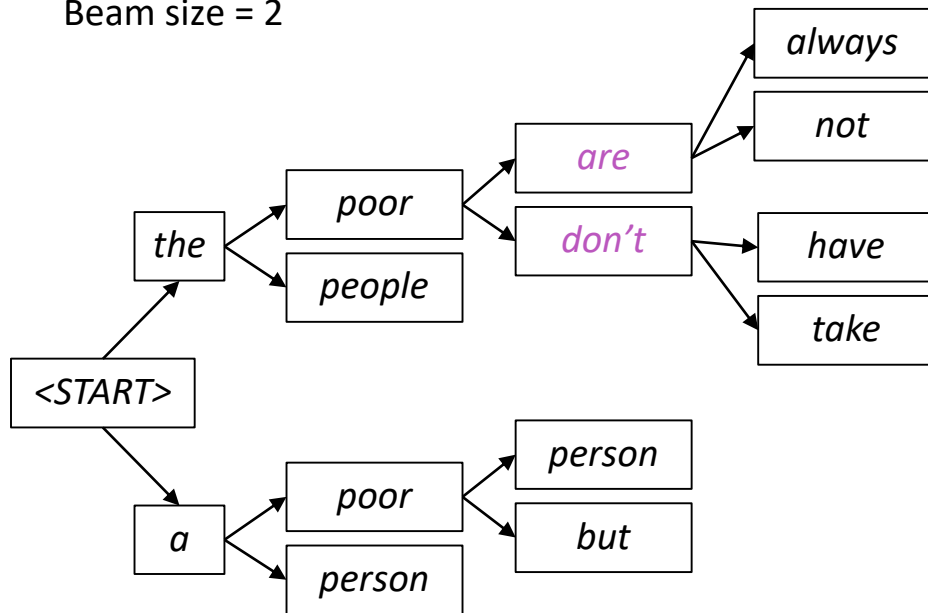
Decoder based on Beam search: Example

Beam size = 2



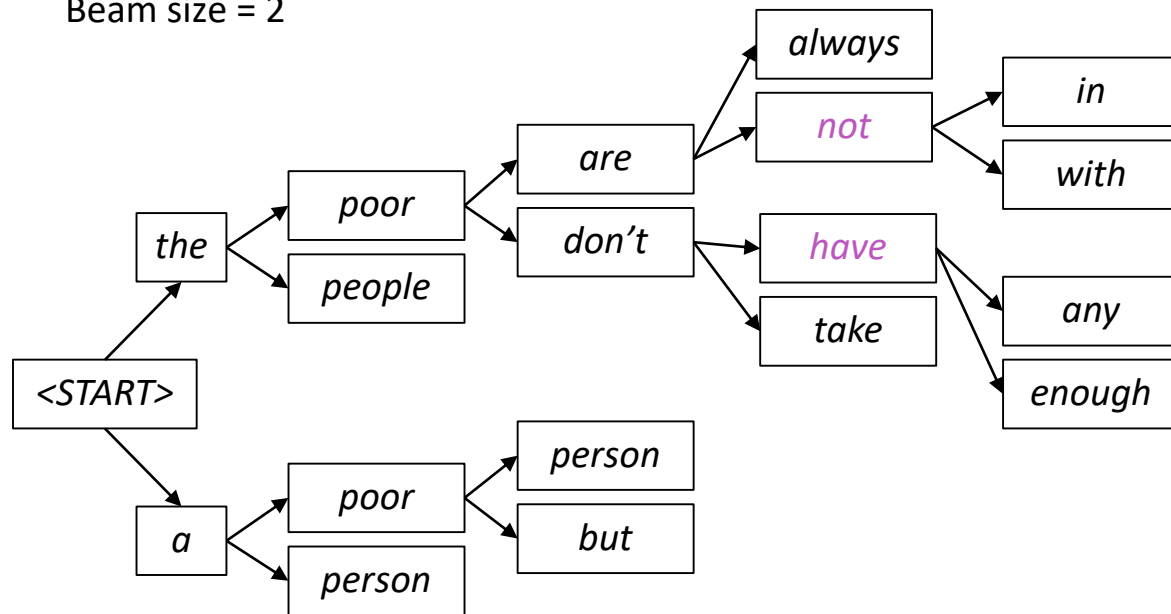
Decoder based on Beam search: Example

Beam size = 2



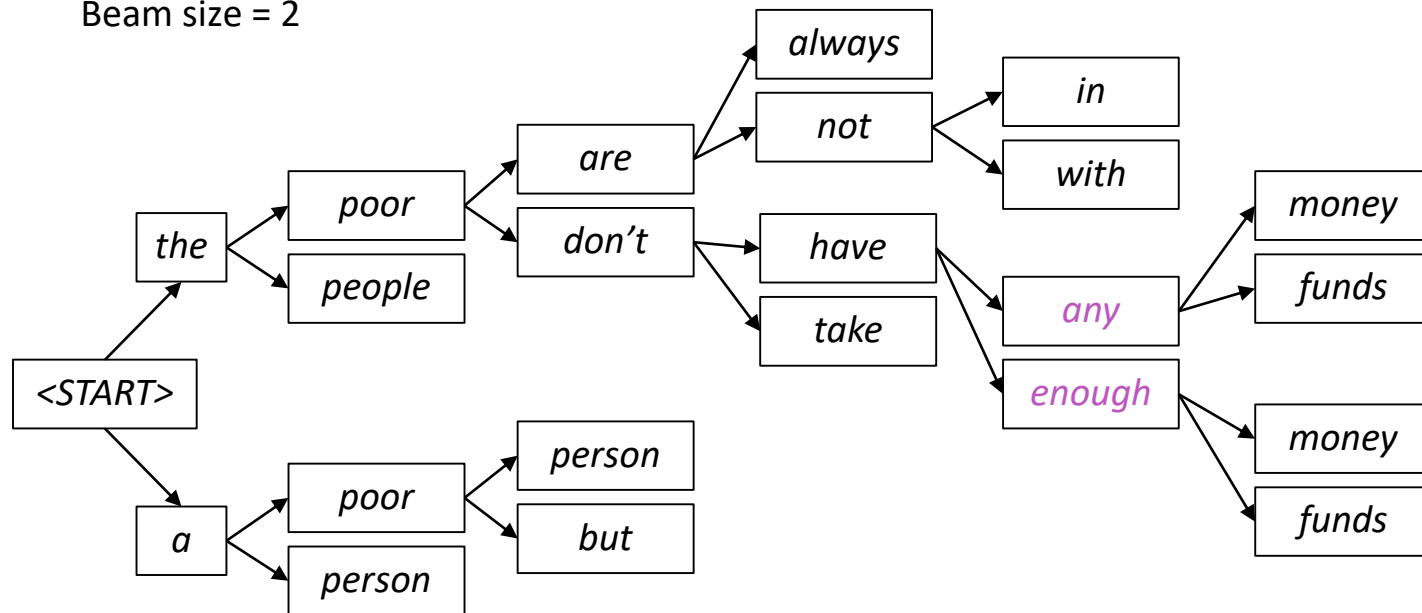
Decoder based on Beam search: Example

Beam size = 2



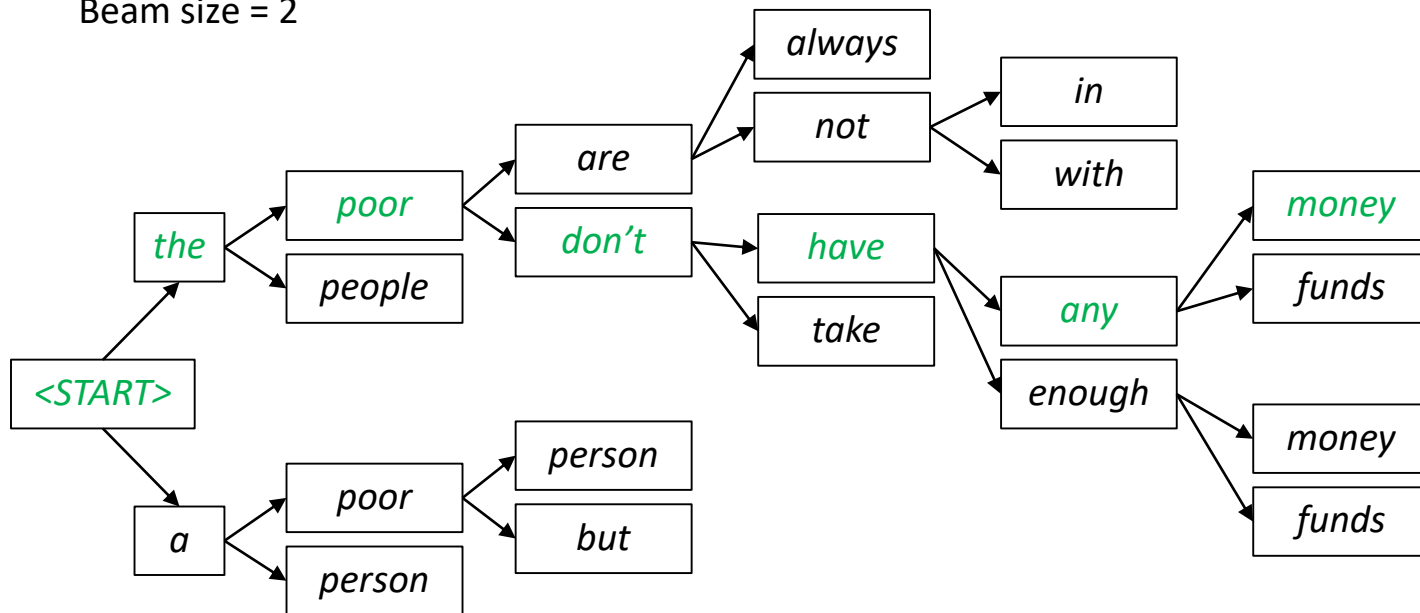
Decoder based on Beam search: Example

Beam size = 2



Decoder based on Beam search: Example

Beam size = 2



Beam search: stopping criterion

- In greedy decoding, usually we decode until the model produces an `<END>`

Ví dụ : `<START>` *he hit me with a pie* `<END>`

- In **beam search decoding**, different hypotheses may produce `<END>` tokens on different timesteps
 - When a hypothesis produces `<END>`, that hypothesis is **complete**.
 - **Place it aside** and continue exploring other hypotheses via beam search.
- Usually we continue beam search until:
 - We reach timestep T (where T is some pre-defined cutoff), or
 - We have at least n completed hypotheses (where n is pre-defined cutoff)

Advantage of NMT

Compared to SMT, NMT has many **advantages**:

Better **performance**

- More **fluent**
 - Better use of **context**
 - Better use of **phrase similarities**
-
- A **single neural network** to be optimized end-to-end
 - No subcomponents to be individually optimized
-
- Requires much **less human engineering effort**
 - No feature engineering
 - Same method for all language pairs

Weakness of NMT?

Compared to SMT:

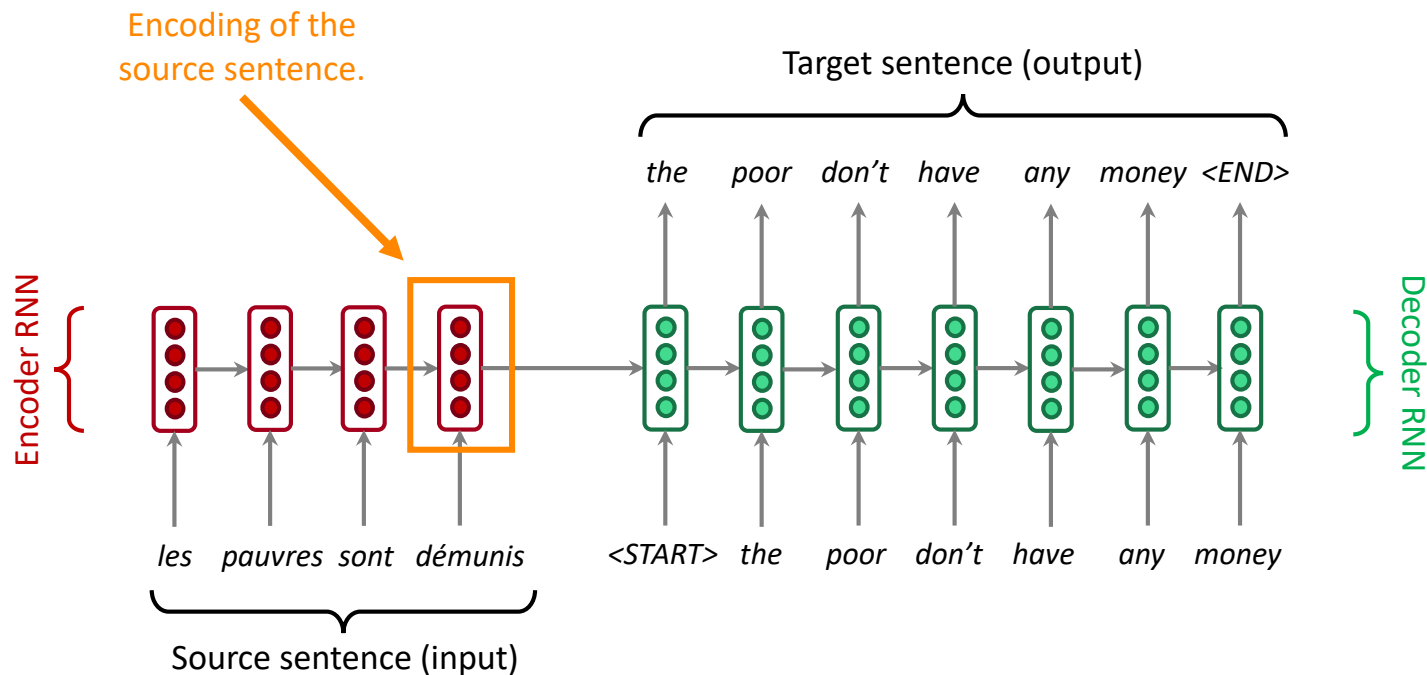
- NMT is **less interpretable**
 - Hard to debug
- NMT is **difficult to control**
 - For example, can't easily specify rules or guidelines for translation
 - Safety concerns!

How to evaluation MT system?

BLEU (Bilingual Evaluation Understudy)

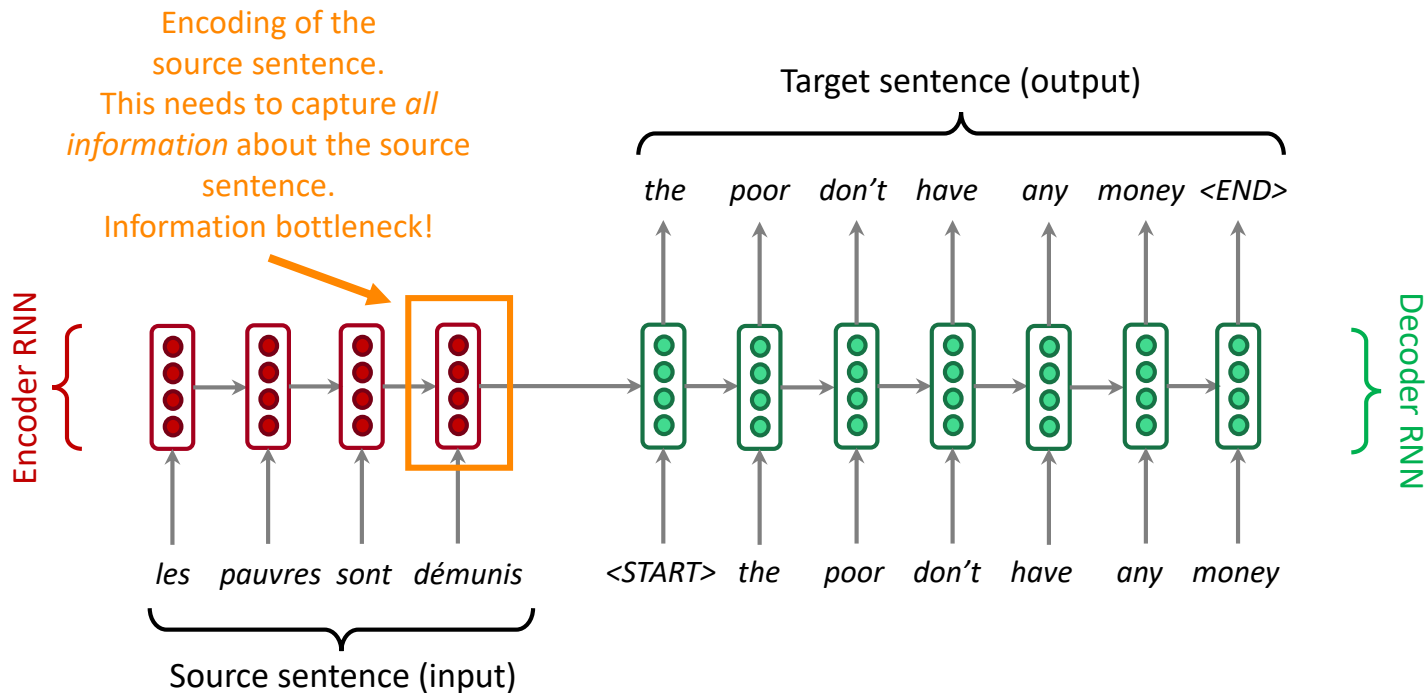
- BLEU compares the machine-written translation to one or several human-written translation(s), and computes a **similarity score** based on:
 - ***n*-gram precision** (usually up to 3 or 4-grams)
 - Penalty for too-short system translations
- BLEU is **useful** but **imperfect**
 - There are many valid ways to translate a sentence
 - So a **good** translation can get a **poor** BLEU score because it has low *n*-gram overlap with the human translation 😞
- SacreBLEU (2018)/COMET (2020)

Sequence-to-sequence: the bottleneck problem

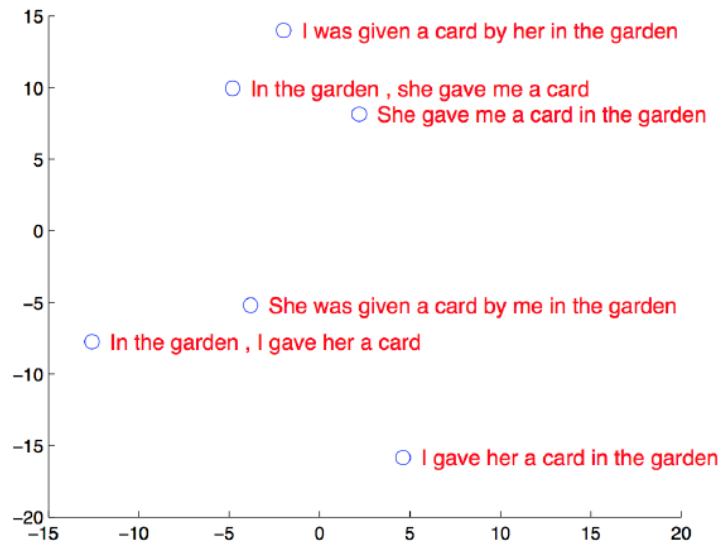
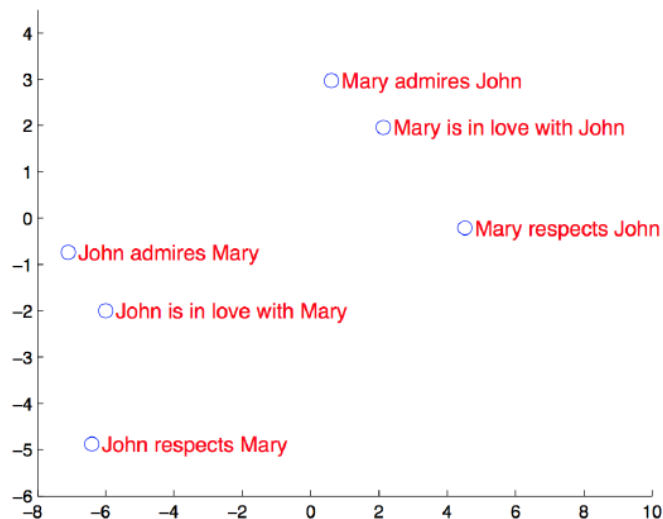


Problems with this architecture?

Sequence-to-sequence: the bottleneck problem



Watching sentence embedding



Attention (Chú ý)

- Attention provides a solution to the bottleneck problem

ICLR 2015

NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

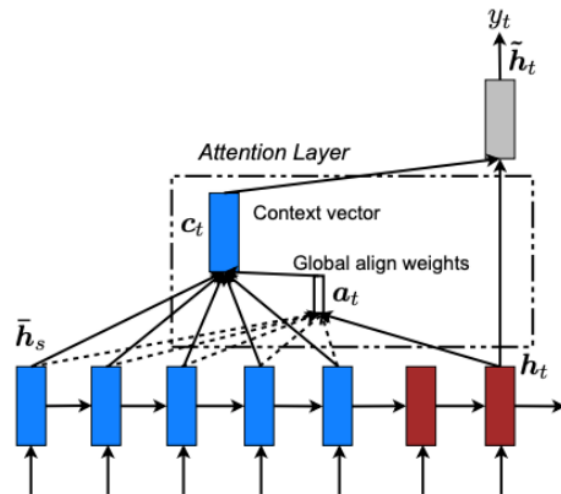
Dzmitry Bahdanau
Jacobs University Bremen, Germany

Kyunghyun Cho **Yoshua Bengio***
Université de Montréal

EMNLP 2015

Effective Approaches to Attention-based Neural Machine Translation

Minh-Thang Luong **Hieu Pham** **Christopher D. Manning**
Computer Science Department, Stanford University, Stanford, CA 94305
{lmthang, hyhieu, manning}@stanford.edu

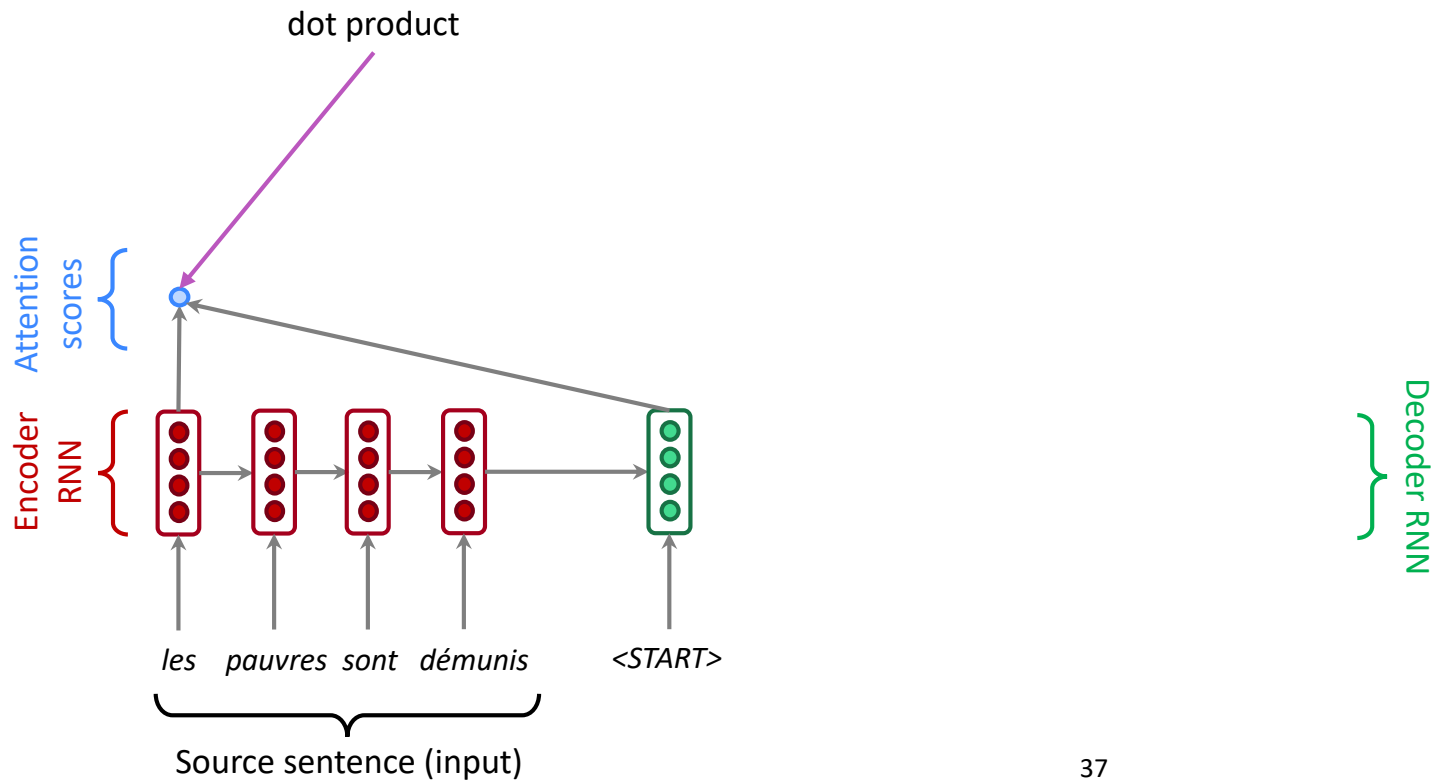


Attention (Chú ý)

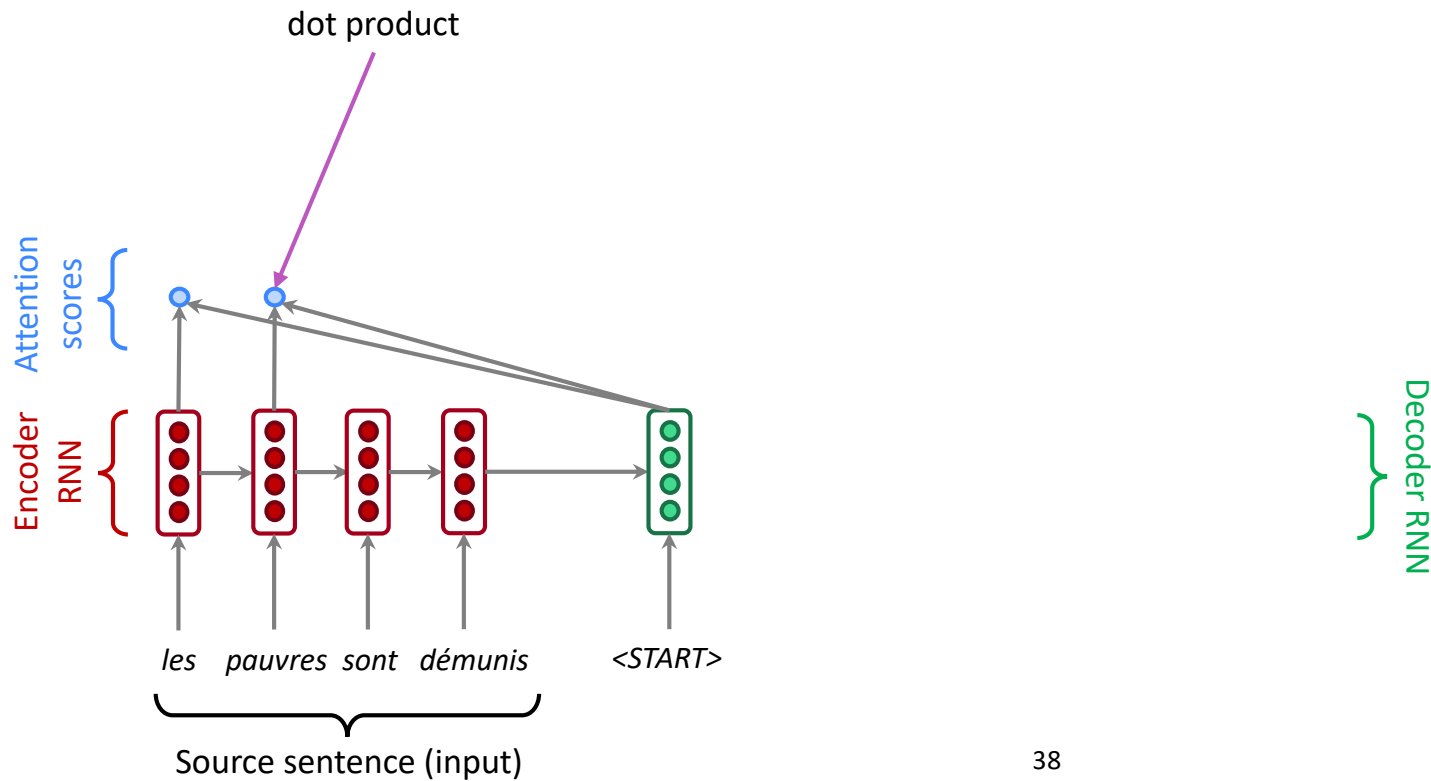


- **Attention: Solution** for bottleneck problem.
- **Main idea**: on each time step of the decoder, use *direct connection to the encoder* to *focus on a particular part* of the source sequence
 - This depends on the **decoder's** current hidden state ht_{dec} (i.e. an idea of what you are trying to decode)
 - Usually implemented as a probability distribution over the hidden states of the **encoder** (hi_{enc})

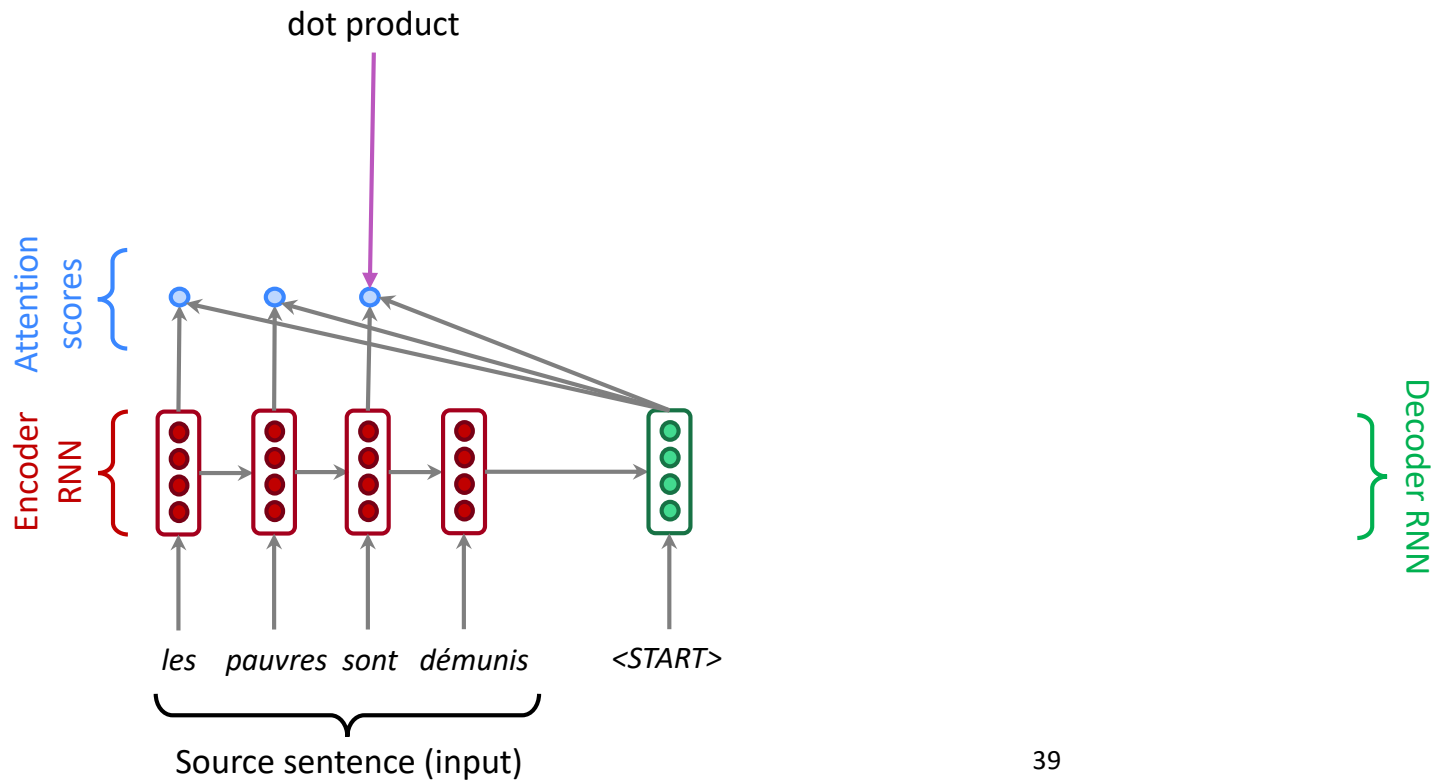
Encoder-Decoder with Attention



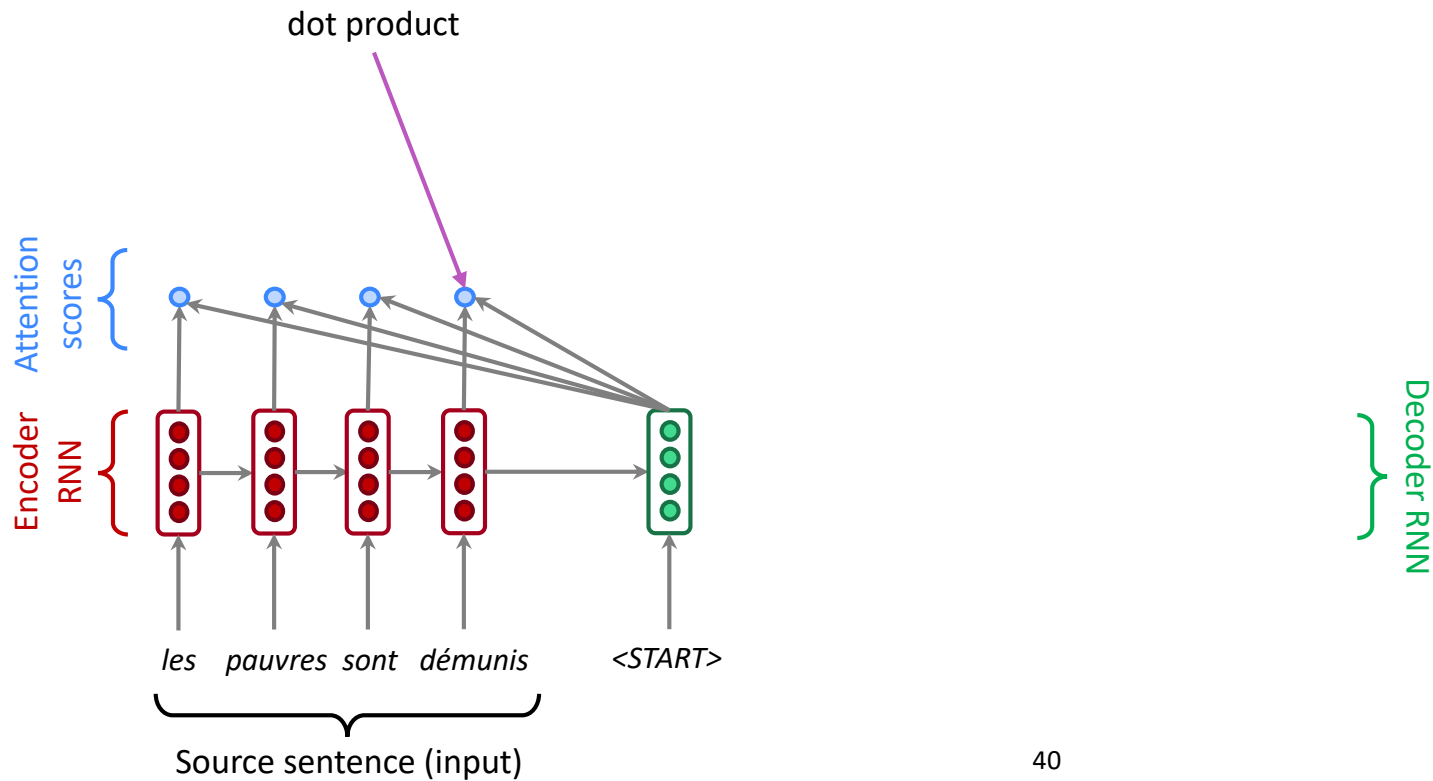
Seq2Seq with attention



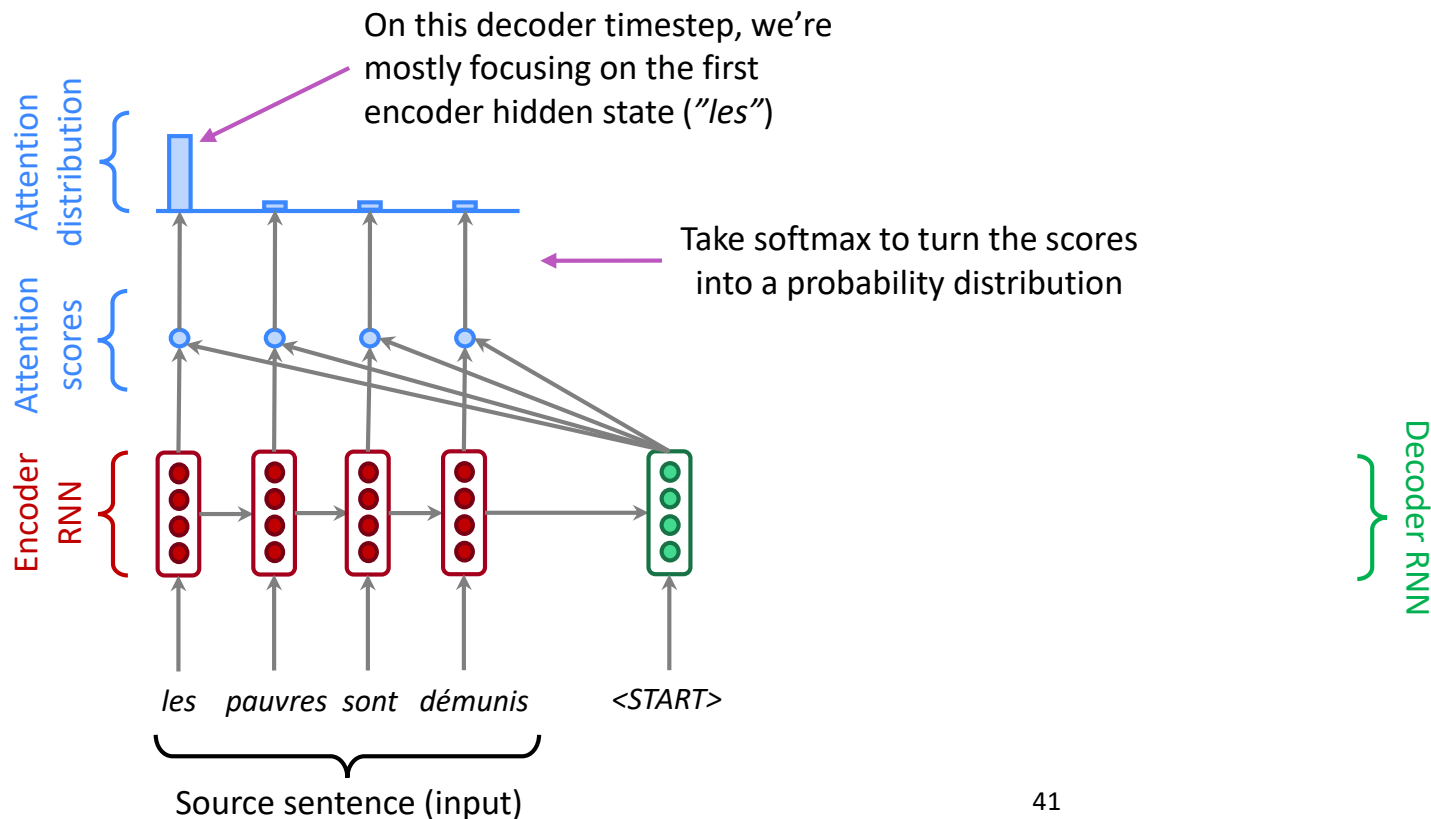
Sequence-to-sequence with attention



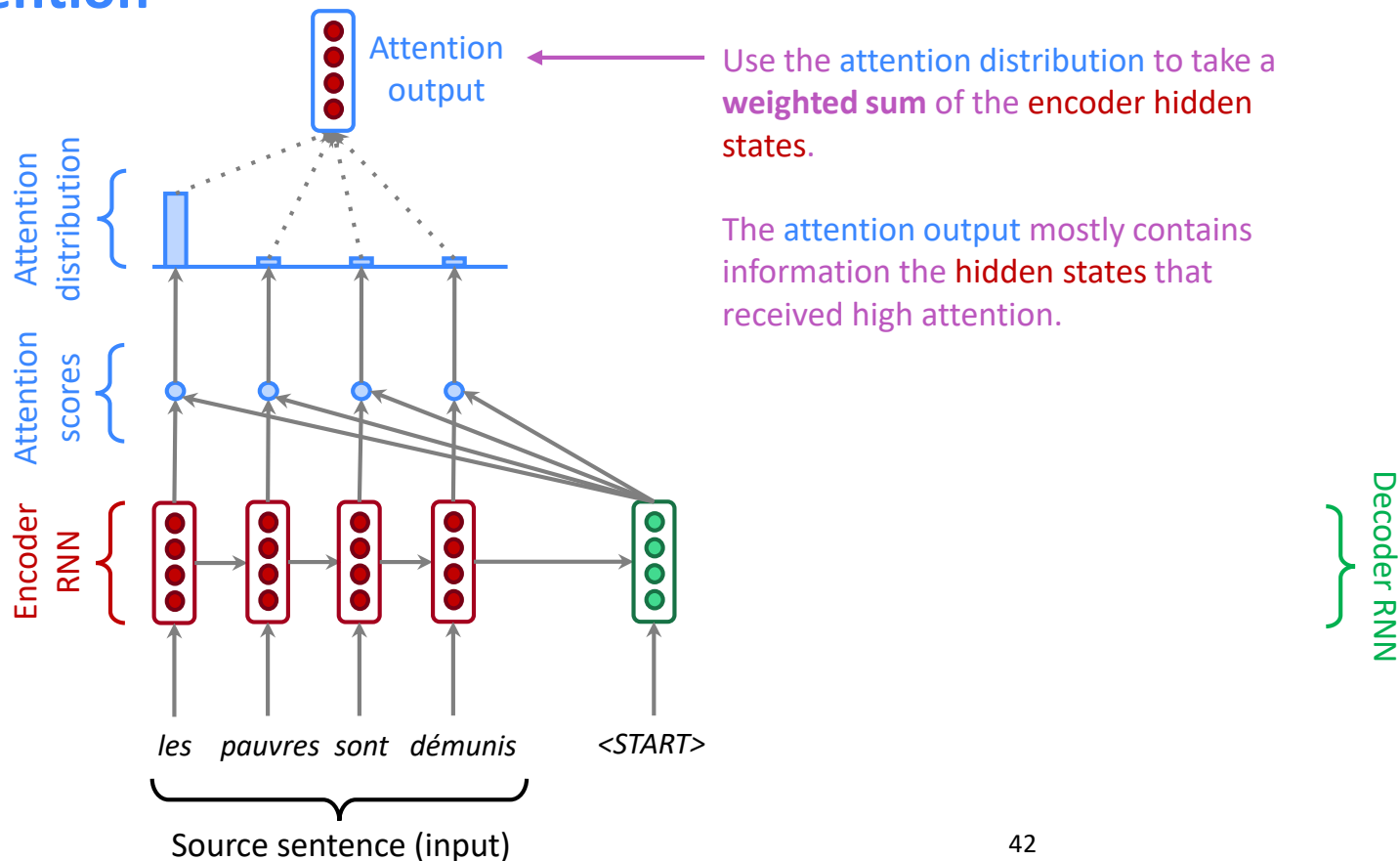
Sequence-to-sequence with attention



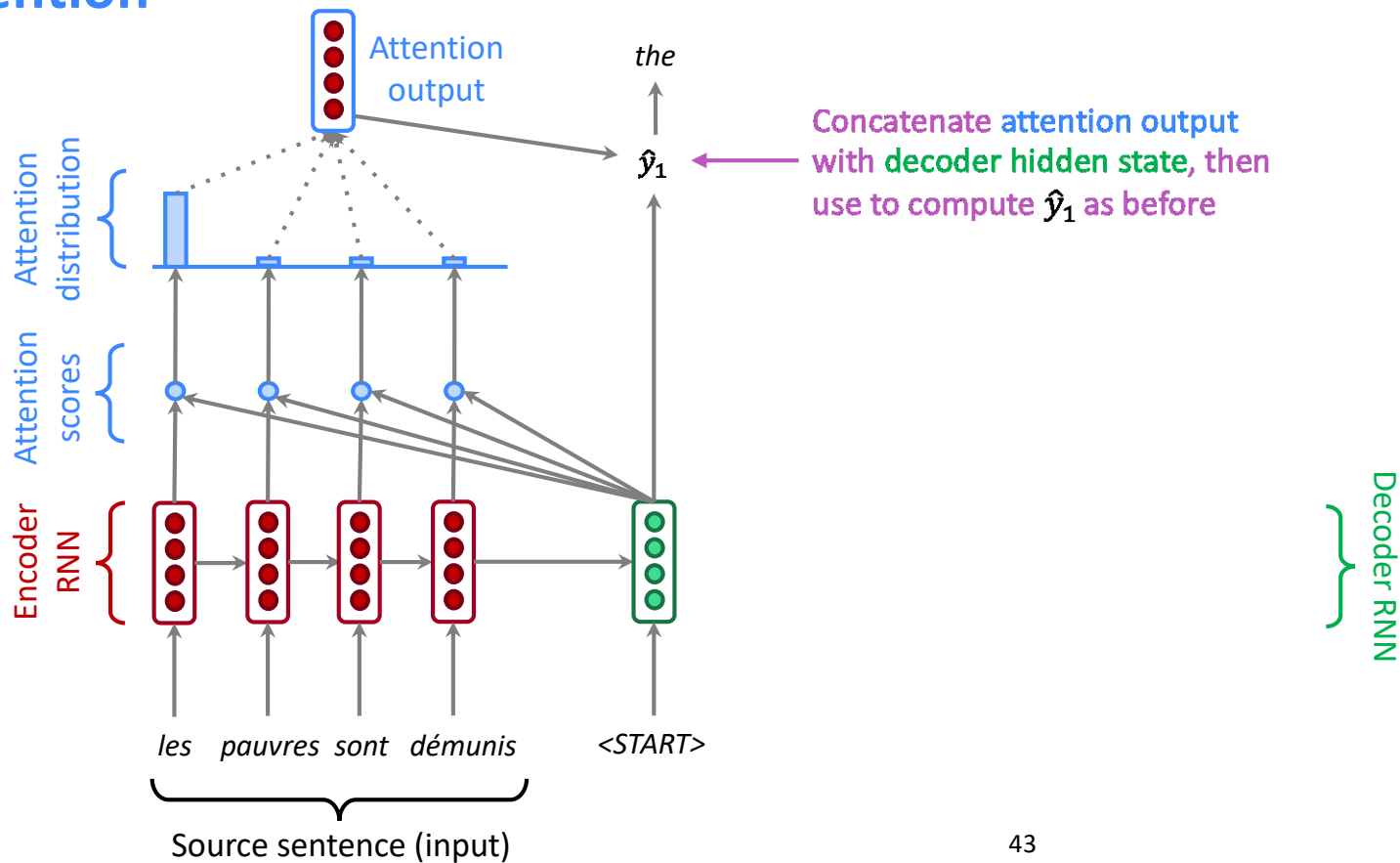
Sequence-to-sequence with attention



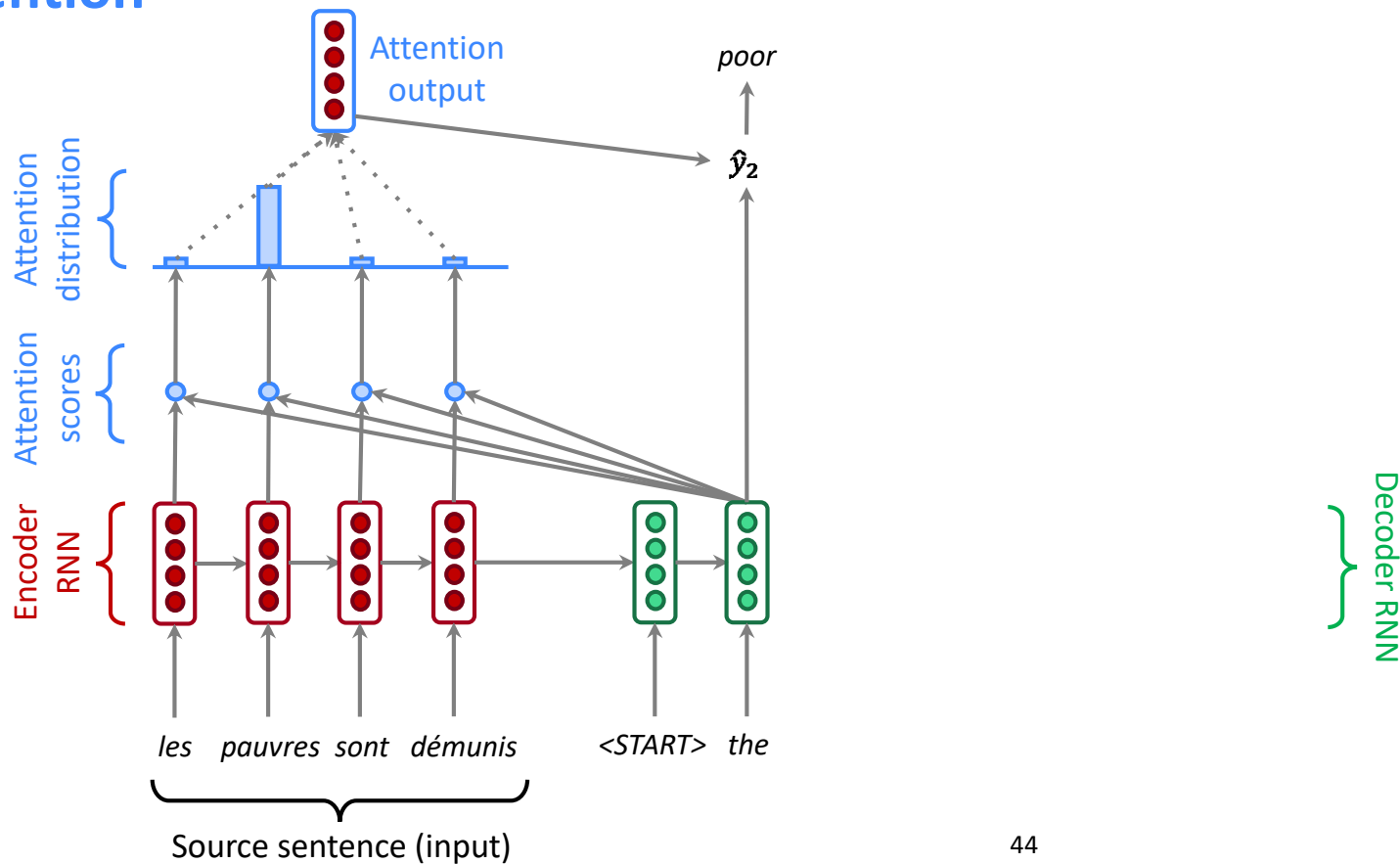
Sequence-to-sequence with attention



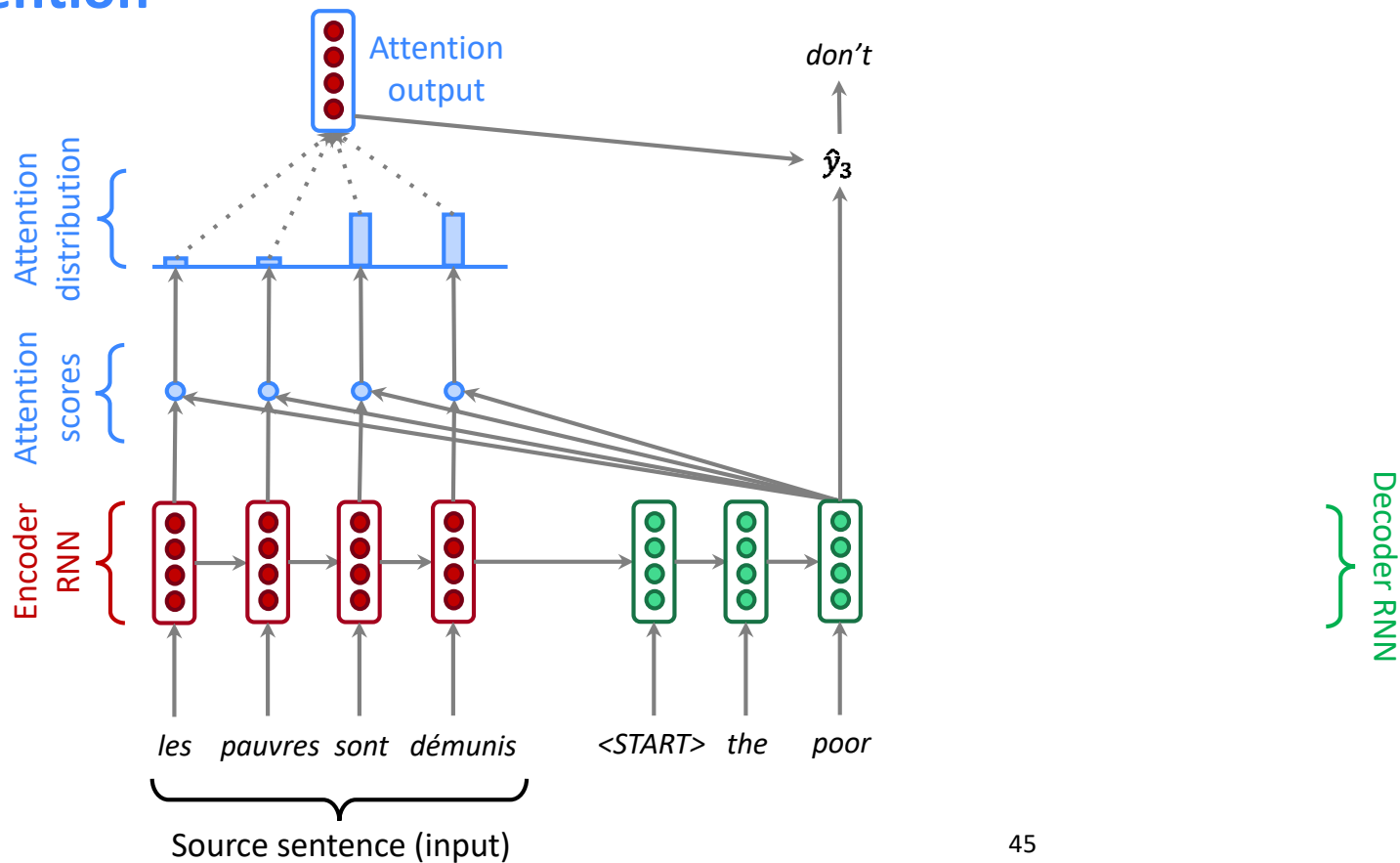
Sequence-to-sequence with attention



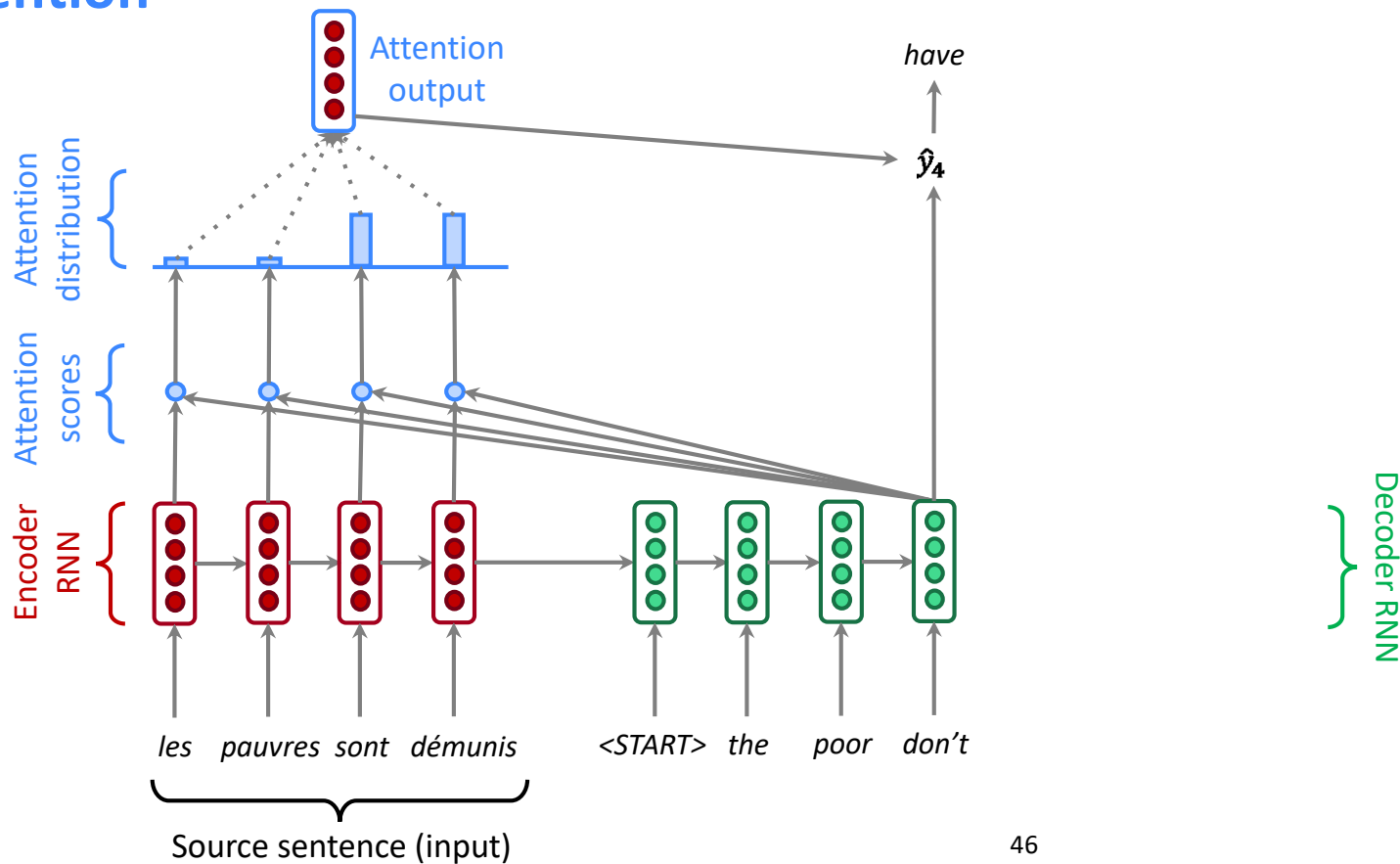
Sequence-to-sequence with attention



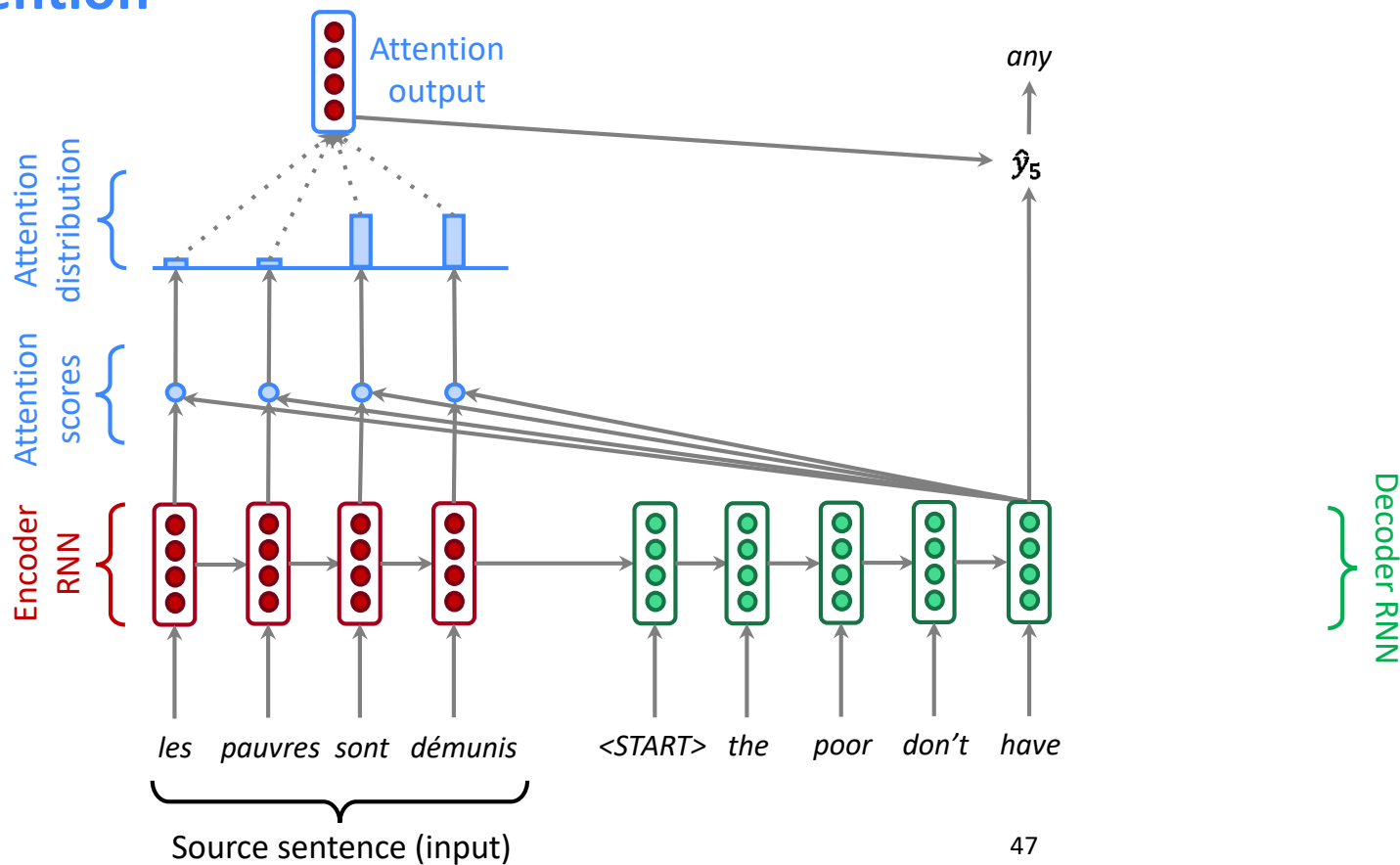
Sequence-to-sequence with attention



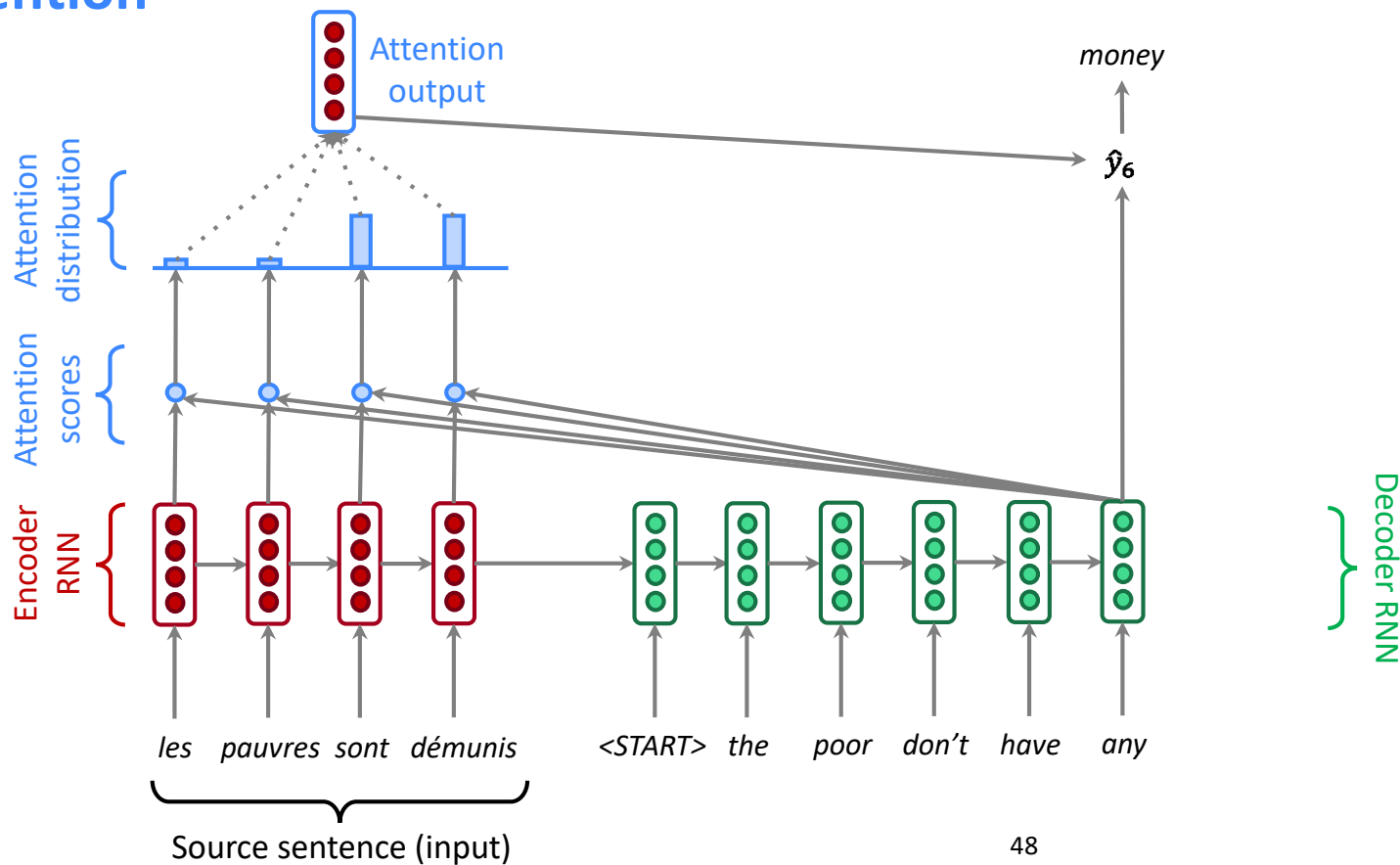
Sequence-to-sequence with attention



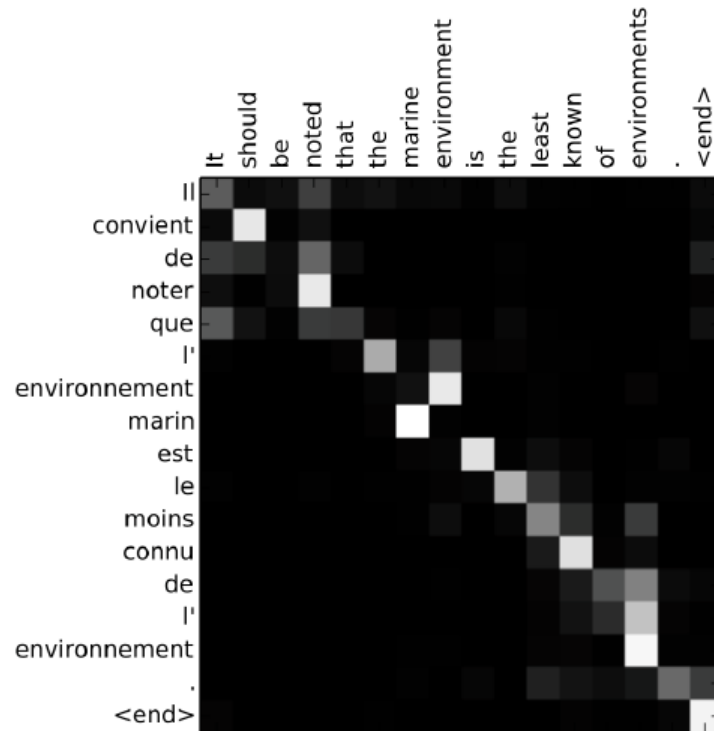
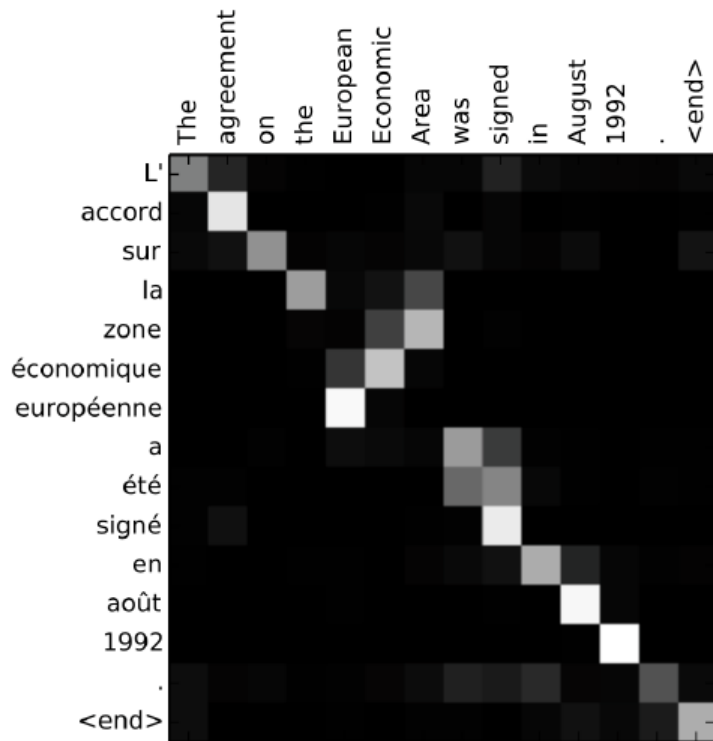
Sequence-to-sequence with attention



Sequence-to-sequence with attention

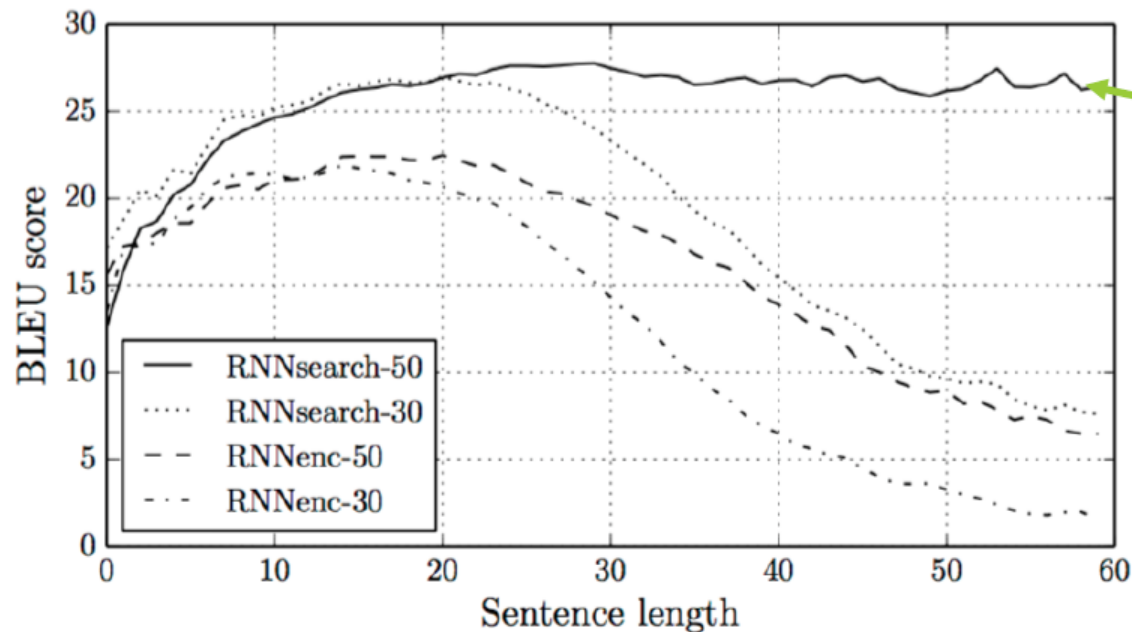


Neural Machine Translation by Jointly Learn to Align and Translate



Source: Bahdanau et al., ICLR 2015, <https://arxiv.org/abs/1409.0473>

Neural Machine Translation by Jointly Learn to Align and Translate



No performance drop
for long sentences!

Attention: Formula

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

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

There are multiple ways to do this

- We take softmax to get the attention distribution α^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 α^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 and proceed as in the non-attention: $[a_t; s_t] \in \mathbb{R}^{2h}$

Attention scoring function

- q is the query and k is the key
- **Multi-layer Perceptron** (Bahdanau et al. 2015)
 - Flexible, often very good with large data

$$a(q, k) = w_2^T \tanh(W_1[q; k])$$

- **Bilinear** (Luong et al. 2015)

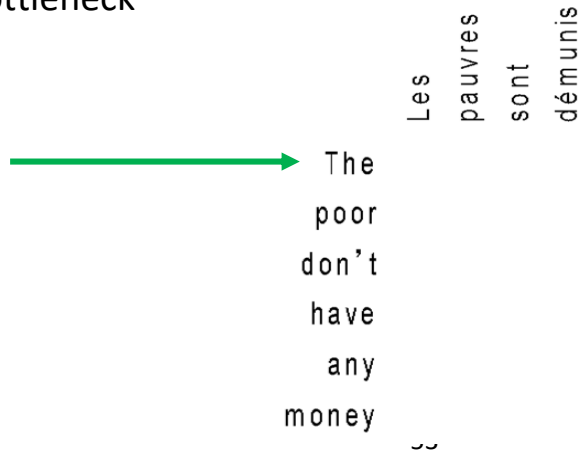
$$a(q, k) = q^T W k$$

- **Dot Product** (Luong et al. 2015)

$$a(q, k) = q^T k$$

Attention is so great (Bahdanau et al., 41067 Citations)

- 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**
- 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



Key developments in attention

Seq2Seq

- Cho et al. (2014)
- Sutskever et al. (2014)

Visual attention

- Xu et al. (2015)

Self attention

- Vaswani et al. (2017)

BERT

- Devlin et al. (2018)

ChatGpt

2014

2015

2016

2017

2018

2019

2020

Align & Translate

- Bahdanau et al. (2015)
- Luong et al. (2015)

Hierarchical attention

- Yang et al. (2016)

? GPT3

Attention is a *general* Deep Learning technique

- 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.

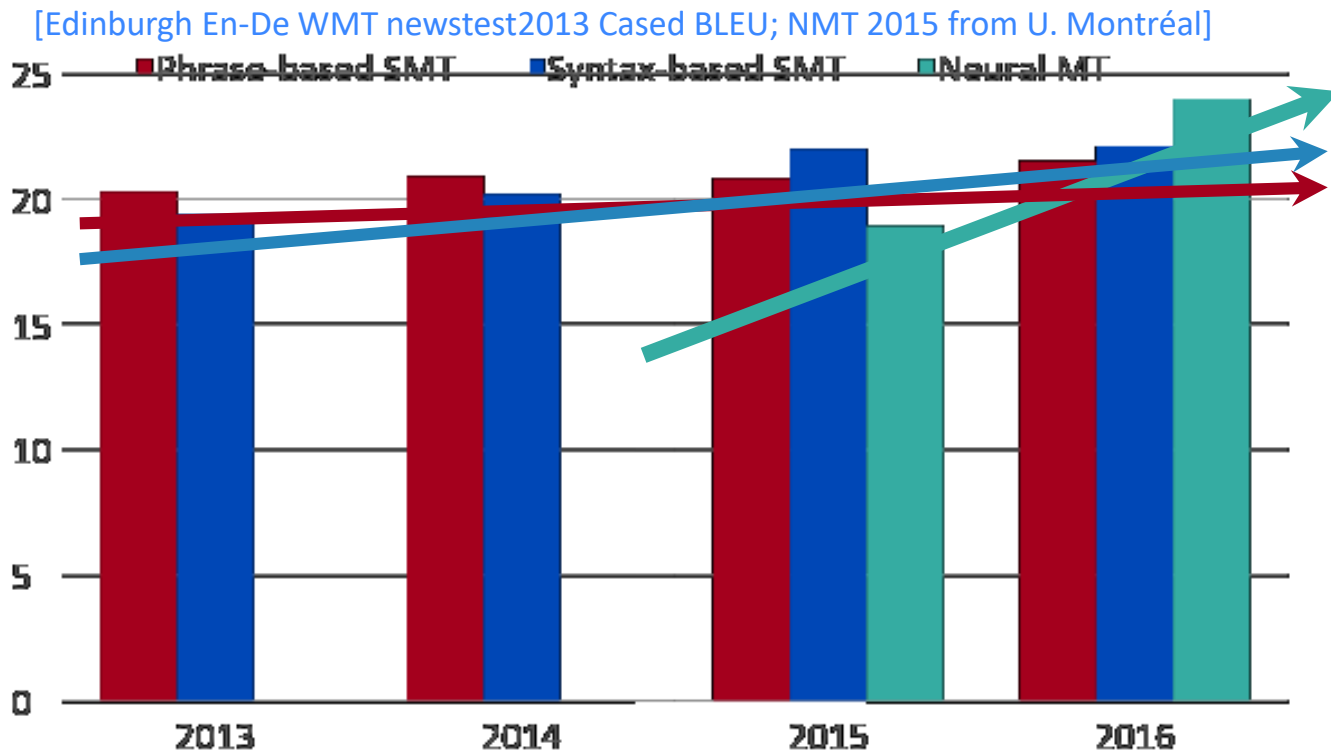
Intuition:

- 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).

Upshot:

- Attention has become the powerful, flexible, general way pointer and memory manipulation in all deep learning models. A new idea from after 2010! From NMT!

Evolution of the MT system over time



Source: http://www.meta-net.eu/events/meta-forum-2016/slides/09_sennrich.pdf

NMT: the first big story of NLP Deep Learning

Neural Machine Translation went from a **fringe research activity** in **2014** to the **leading standard method** in **2016**

- **2014**: First seq2seq paper published [[Sutskever et al., 2014](#)]
- **2016**: Google Translate switches from SMT to NMT and by 2018 everyone has



- **This is amazing!**
 - **SMT** systems, built by **hundreds** of engineers over many **years**, outperformed by NMT systems trained by a **small groups** of engineers in a few **months**

MT solved?

- **Nope!**
- Many difficulties remain:
 - Out-of-vocabulary words (unknown words)
 - Domain mismatch between train and test data
 - Maintaining context over longer text
 - Low-resource language pairs (Hallucinations)

Seq2seq is flexible and efficient!

- Seq2Seq is useful not only in the Machine Translation
- 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)
 - **OCR** (image of character → text sequence)
 - **ASR** (sequence of Acoustic → text sequence)

Machine Translation Problem

- Automatic translation from English sentence to Vietnamese sentence
- Training và testing data: IWSLT 2015

Example:

Input: I like a blue book

Output: Tôi thích quyển sách màu xanh

Apply seq2seq + attention for NMT

Goals to be achieved:

1. Training seq2seq model for English-Vietnamese translation problem.
2. Evaluating the translation system by BLEU score.
3. Seeing Attention score matrix to better understand the model

The tasks to be implemented:

1. Data preprocessing
2. Create training data
3. Write **encoder, decoder, attention** module
4. Training model
5. Translate new source sentence
6. Show Attention, BLEU score.

Conclusion

- **Sequence-to-sequence** is the architecture of most current NLP problems such as NMT, text generation, ...
- **Attention great** is a way to *focus on particular parts* of the input
 - Improved Seq2Seq model **a lot!**
 - As **the foundation** of the Transformer model (now dominated!)

References

- Speech and Language Processing 2025 (Chapter 13)
(<https://web.stanford.edu/~jurafsky/slp3/>)
- Machine Translation and Sequence-to-Sequence Models,
Neubig, 2019, CMU
- Some slides for Universities of Stanford 2024, MIT, Princeton,
...