# Sequence to Sequence Models

Mirko Bronzi
Applied Research Scientist, Mila
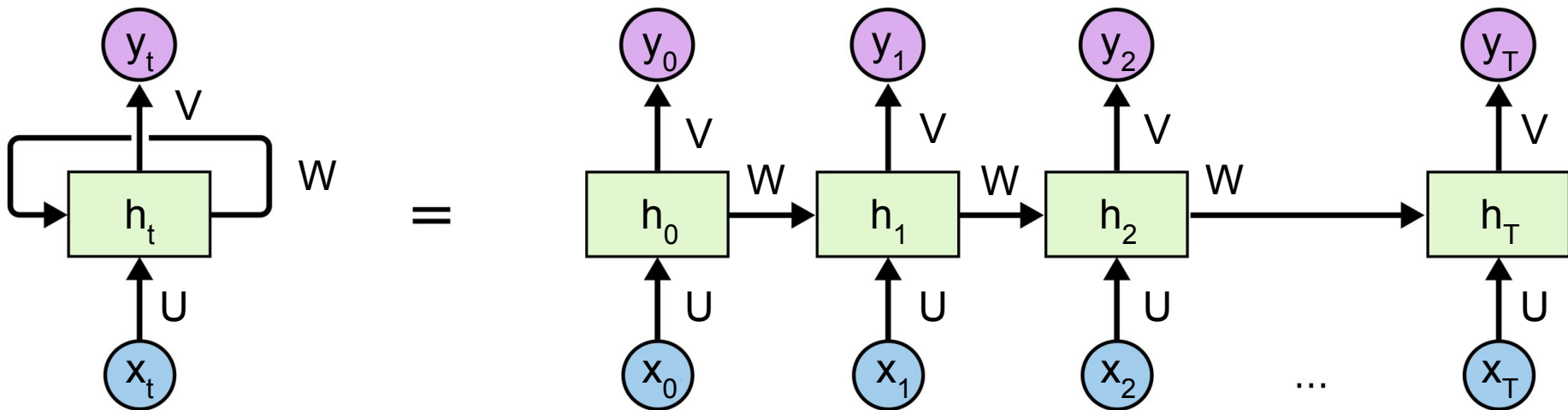mirko.bronzi@mila.quebec

# Plan

- RNN Recap
- Sequence to Sequence Models
- Attention Mechanism
- Transformer
- Libraries and References

Mila

# Plan

- **RNN Recap**
- Sequence to Sequence
- Attention Mechanism
- Transformer
- Libraries and References

Mila

# Recurrent Neural Networks

- The parameters of the model are **shared** over time.
- The internal state ($h_t$) is updated at each time step.

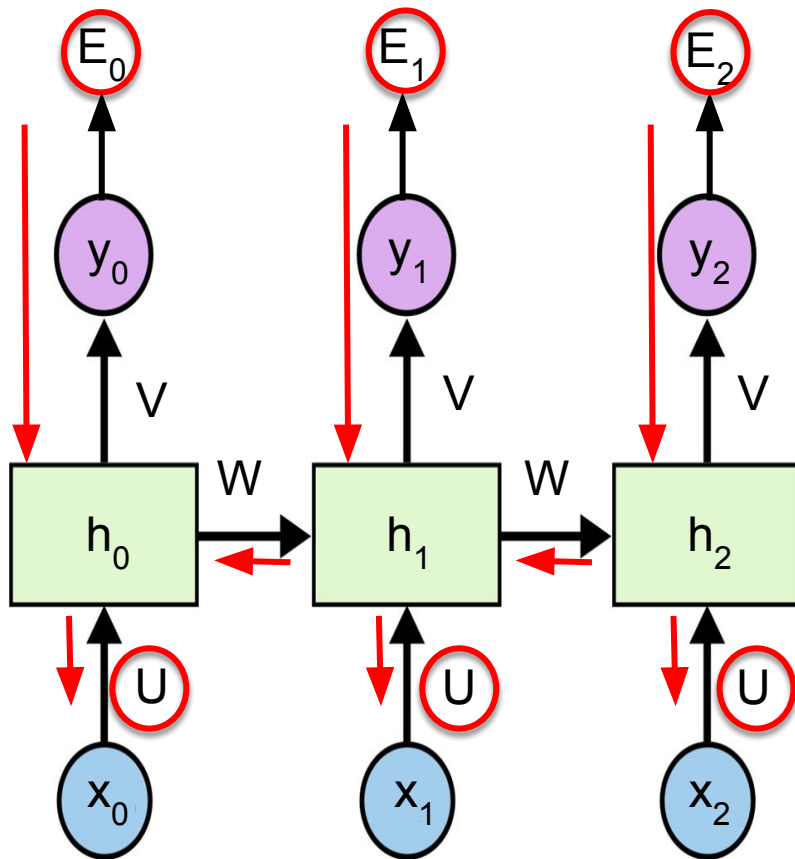The initial internal state ($h_{-1}$) is dropped for simplicity

# Backpropagation Through Time

- The global error is:

$$E = \sum_{t=0}^{T} E_t$$

- To compute the gradient of the global error with respect to a parameter, we compute the gradient of the individual error at each time step, and then sum all those values. For example:

$$\frac{\partial E}{\partial U} = \sum_{t=0}^{T} \frac{\partial E_t}{\partial U}$$

Image from Christopher Olah's blog

# Long-Term Dependencies

- Long-term dependencies are difficult to learn due to the long chain of gradients $\frac{\partial h_T}{\partial h_{T-1}} \ldots \frac{\partial h_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial h_0}$ that can lead to vanishing gradients
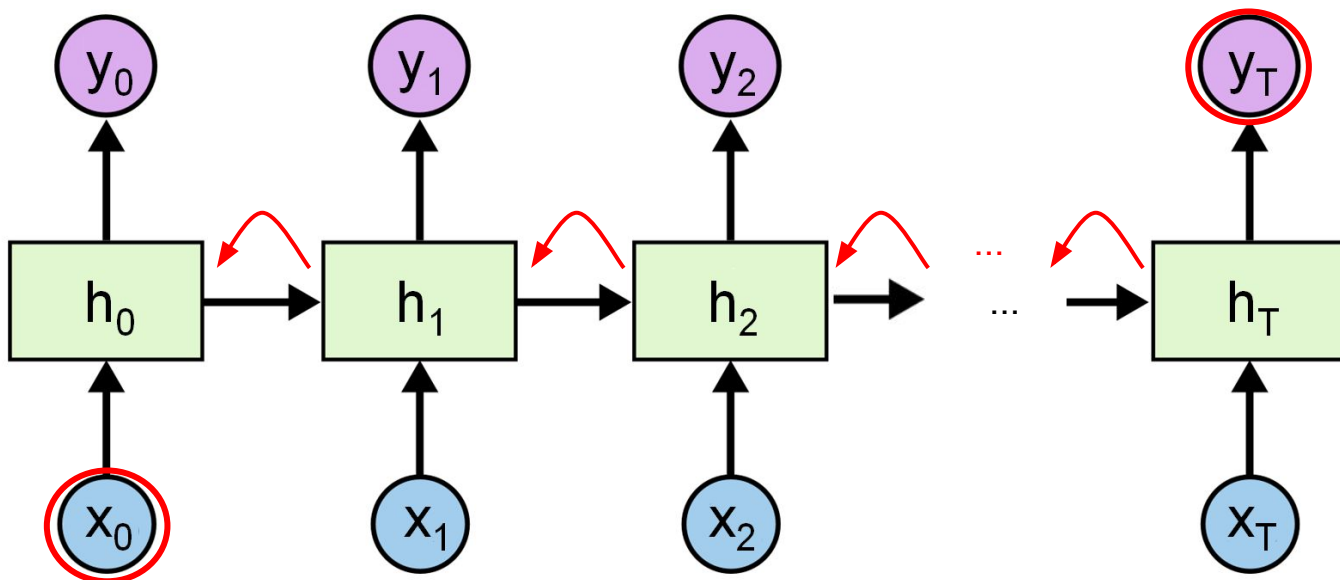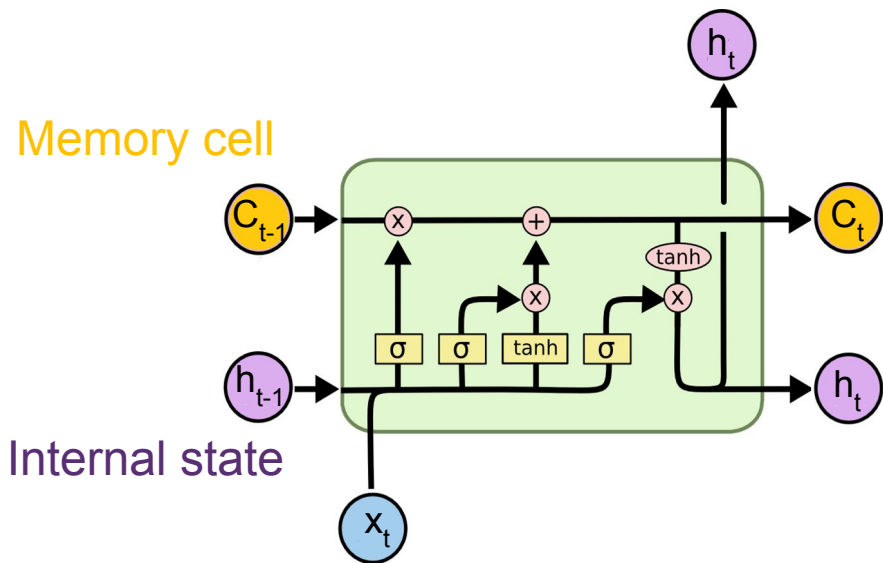


Image from Christopher Olah's blog

# Long Short-Term Memory (LSTM)

- Reduce the vanishing gradient problem using a gate mechanism and adding a memory cell.

Memory cell

Internal state



$$i_t = \sigma(U_i x_t + W_i h_{t-1} + b_i)$$

$$f_t = \sigma(U_f x_t + W_f h_{t-1} + b_f)$$
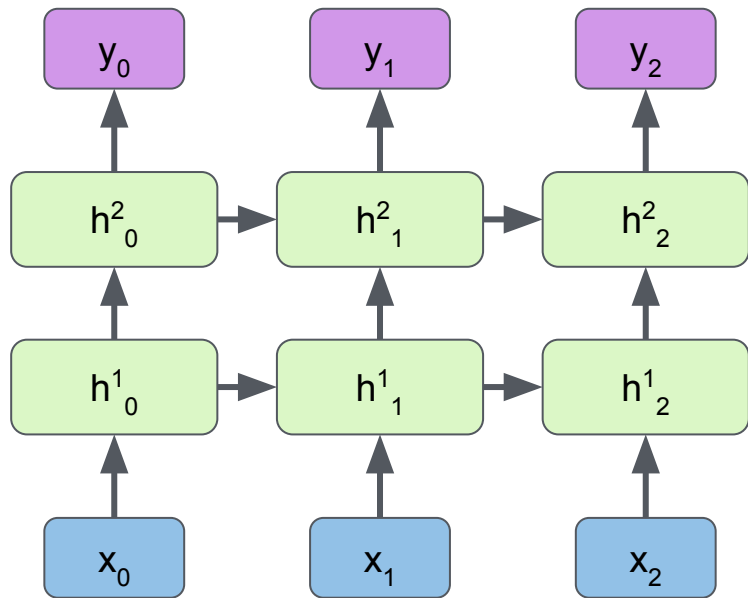
$$o_t = \sigma(U_o x_t + W_o h_{t-1} + b_o)$$

$$g_t = tanh(U_g x_t + W_g h_{t-1} + b_g)$$
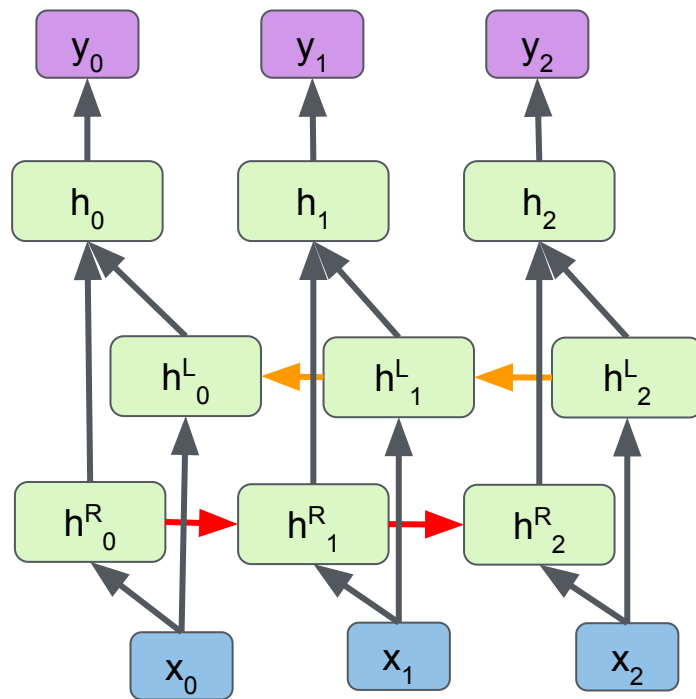
$$C_t = i_t \times g_t + f_t \times C_{t-1}$$

$$h_t = o_t \times tanh(C_t)$$

Image from Christopher Olah's blog

Hochreiter et al., Long short-term memory, Neural Computation 1997

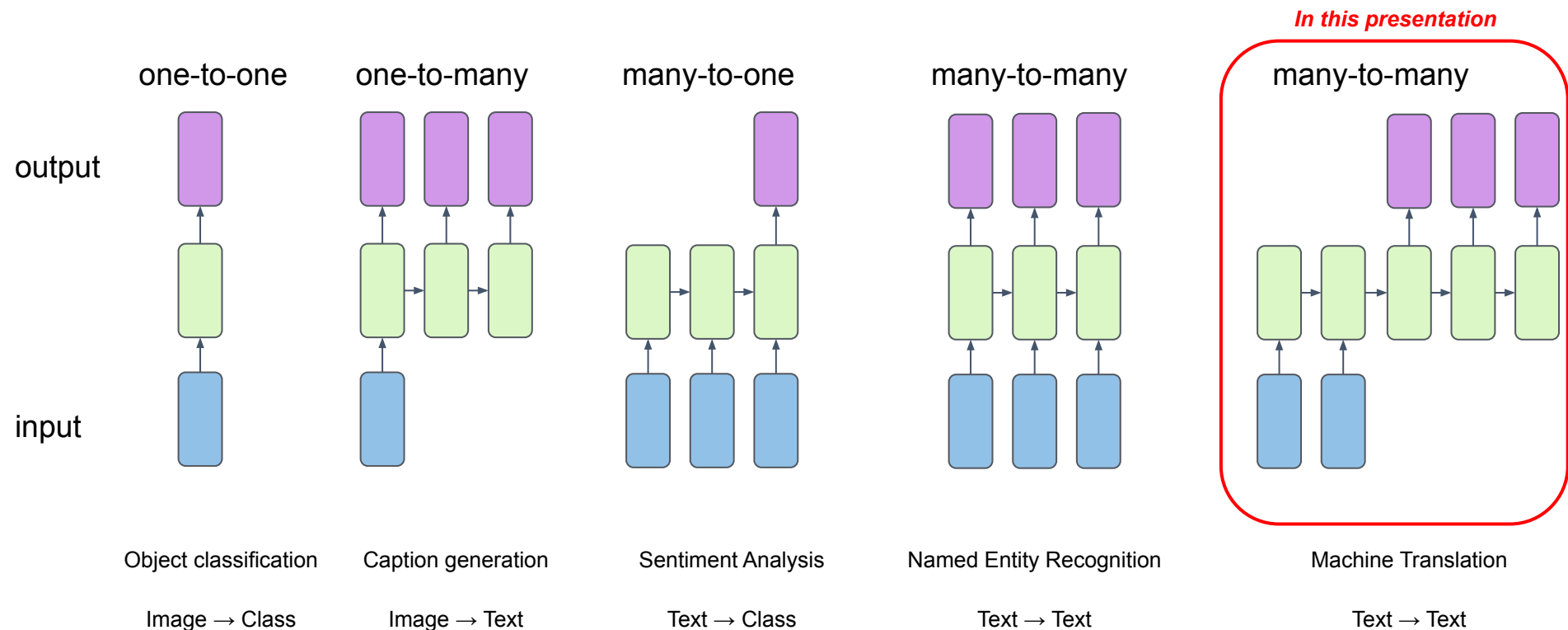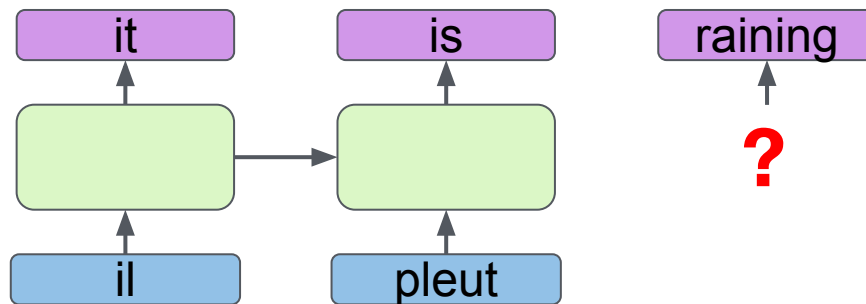# Multi-Layer and Bidirectional RNNs



Layers of RNNs

Bidirectional RNNs

# Plan

- RNN Recap
- **Sequence to Sequence Models**
- Attention Mechanism
- Transformer
- Libraries and References

# Modeling Sequences



*In this presentation*

| one-to-one | one-to-many | many-to-one | many-to-many | many-to-many |

output

input

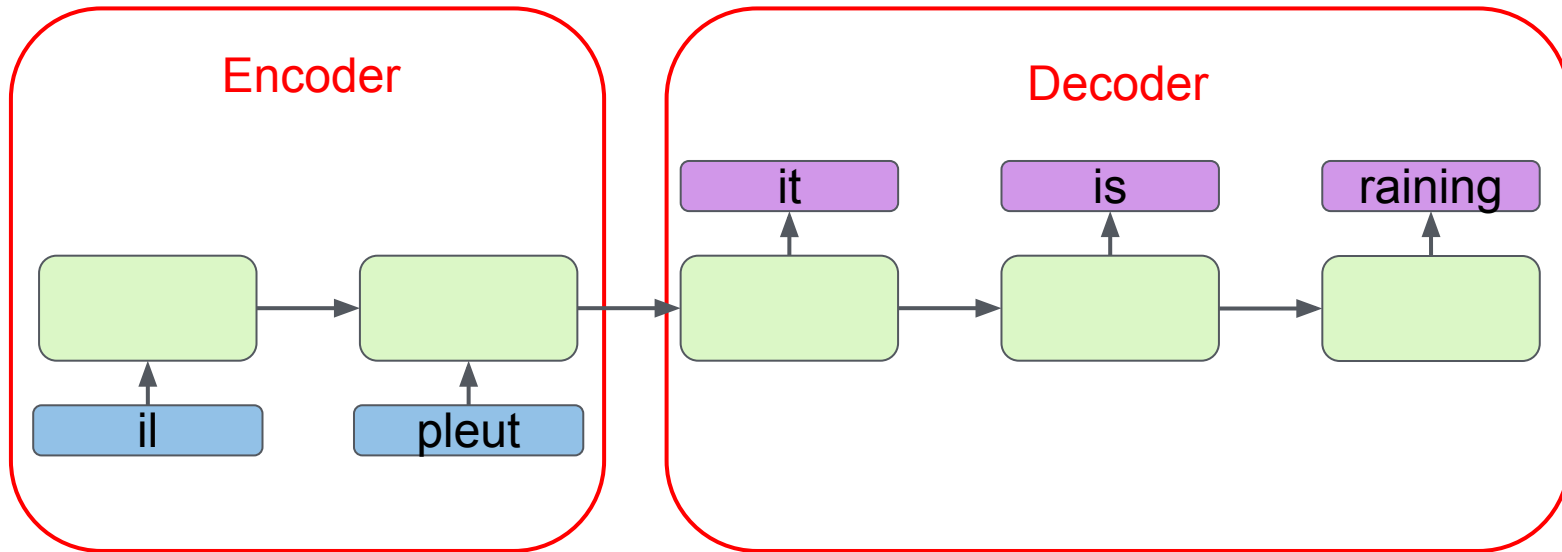| Object classification | Caption generation | Sentiment Analysis | Named Entity Recognition | Machine Translation |
| Image → Class | Image → Text | Text → Class | Text → Text | Text → Text |

# Modeling Sequences

- How to handle input and output sequences of different lengths?
  - Machine translation.
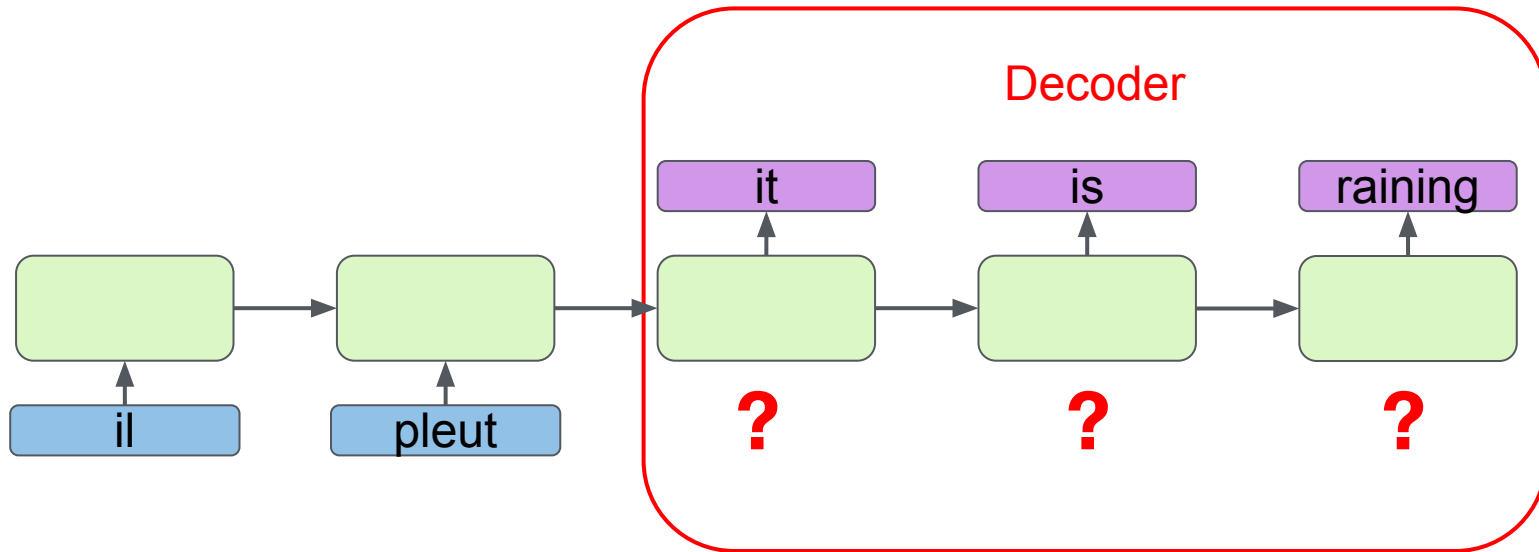  - Text summarization.
  - …

# Different input-output sequence sizes

- Create an architecture composed of two components (e.g. two different RNNs):
    - Encoder that processes the input sequence.
    - Decoder that generates the output sequence based on the encoded input.
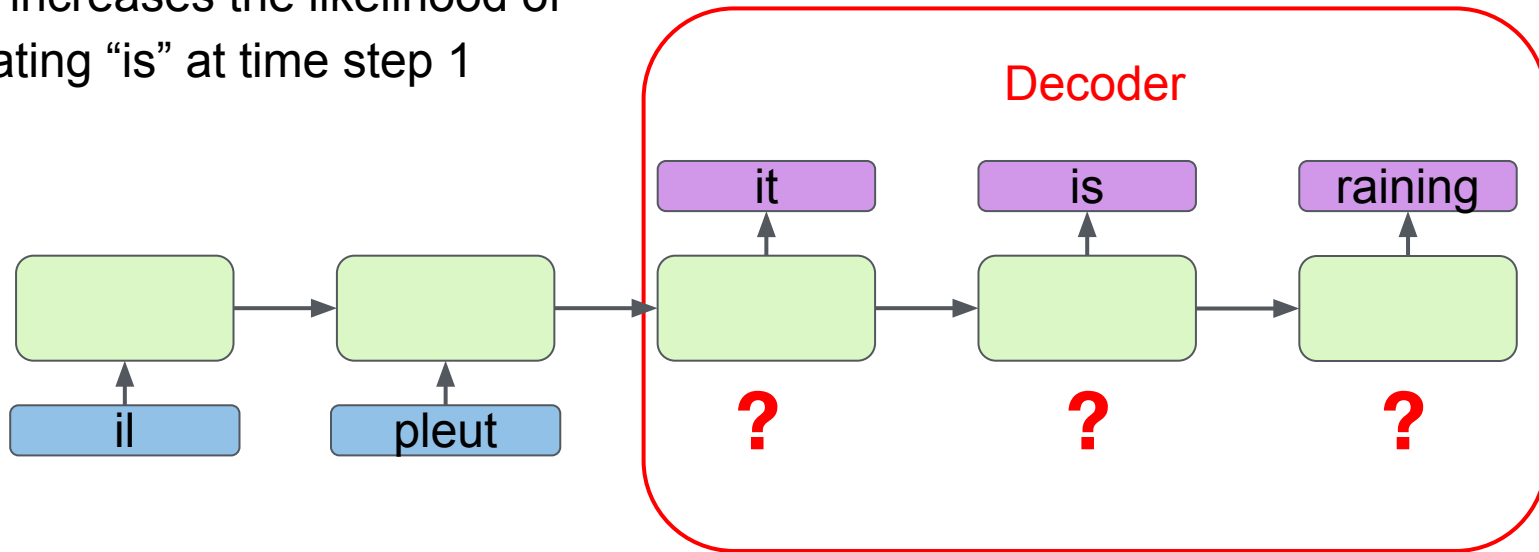
# Different input-output sequence sizes

- How to implement the decoder?
- Note the missing input. We need a mechanism that will allow the decoder to generate consistent outputs across time.

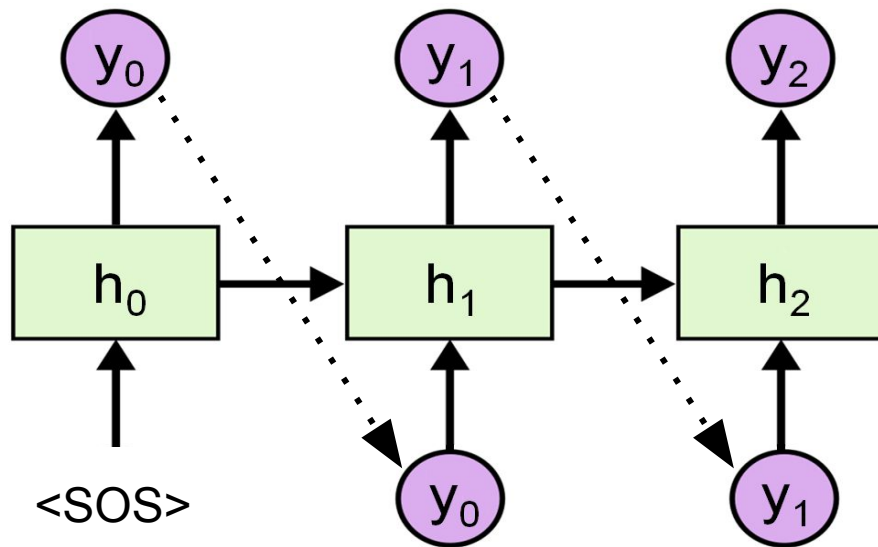# Different input-output sequence sizes

Example: knowing that the decoder generated "it" at time step 0 increases the likelihood of generating "is" at time step 1

# Autoregressive RNNs

<SOS> *Start of sequence*

- We can use a RNN to **generate** a sequence.
- In order to generate consistent outputs across time, we can condition each output on previously generated outputs.
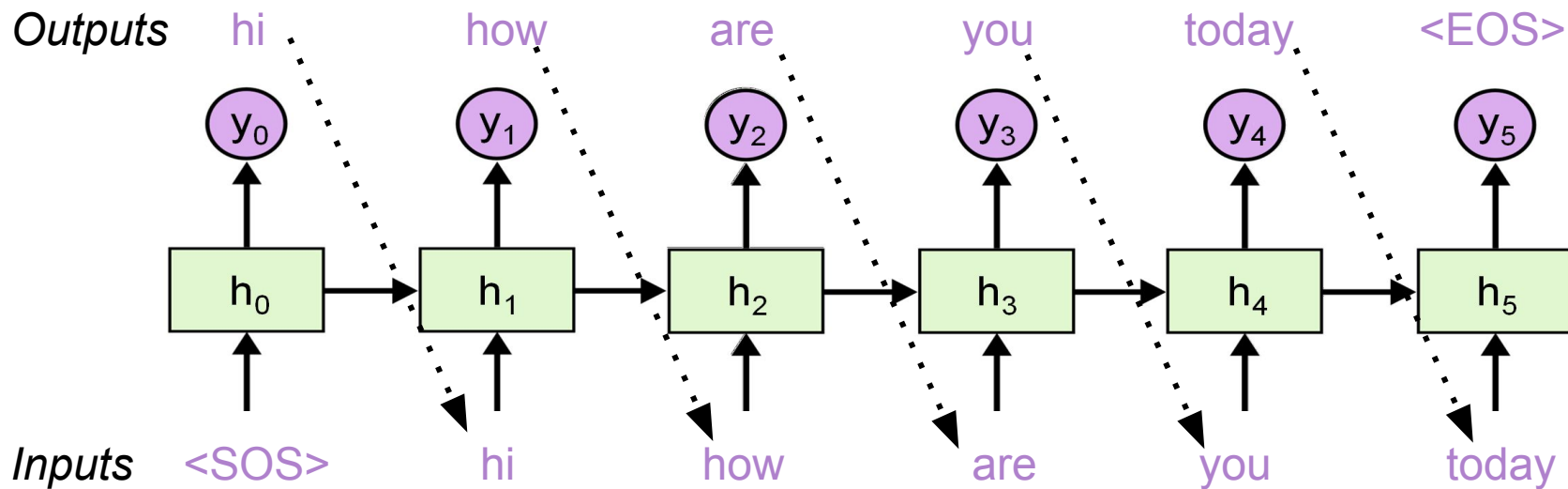- Such a model is called **autoregressive**.

# Autoregressive RNNs

<SOS> *Start of sequence*
<EOS> *End of sequence*

# Sequence-to-Sequence Models

*<SOS> Start of sequence*

*<EOS> End of sequence*

Sutskever et al., Sequence to Sequence Learning with Neural Networks

Cho et al., Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation

Mila

# Sequence-to-Sequence Models - Example

<SOS> Start of sequence

<EOS> End of sequence



| | | | |
|---|---|---|---|
| $h_0$ | $h_1$ | $h_2$ | |
| Je | suis | malade | <SOS> |

Encoder ︸ Decoder (autoregressive RNN)

# Sequence-to-Sequence Models - Example

<SOS> Start of sequence

<EOS> End of sequence



Encoder

Decoder (autoregressive RNN)

# Sequence-to-Sequence Models - Example

\<SOS\> Start of sequence

\<EOS\> End of sequence

# Sequence-to-Sequence Models - Example

<SOS> Start of sequence

<EOS> End of sequence

# Sequence-to-Sequence Models - Example

<SOS> Start of sequence

<EOS> End of sequence



Encoder

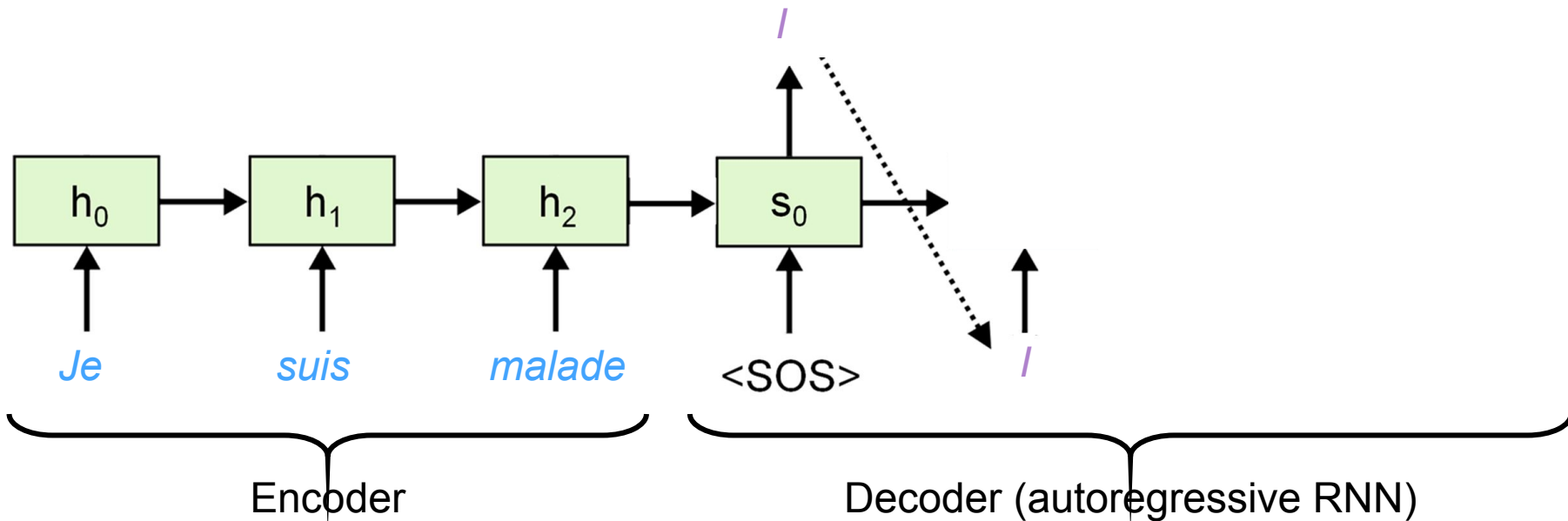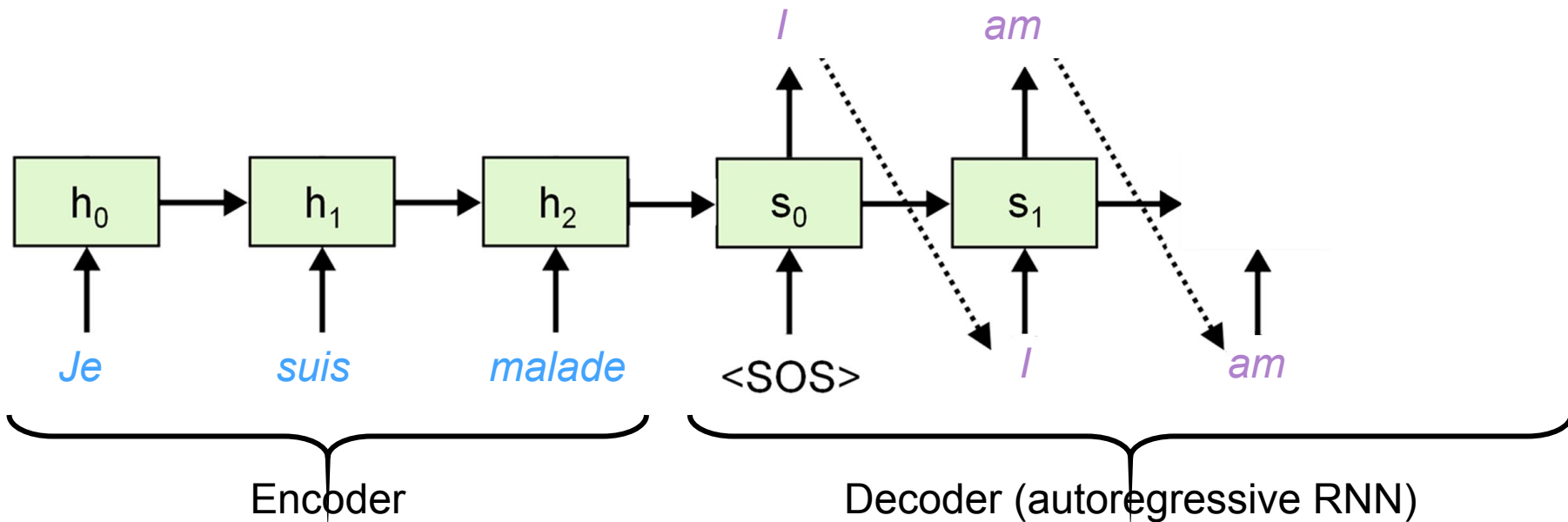Decoder (autoregressive RNN)

# Plan

- RNN Recap
- Sequence to Sequence Models
- **Attention Mechanism**
- Transformer
- Libraries and References

# Sequence-to-Sequence Models - Bottleneck

- The encoder has to store/compress all the information from the input into a **fixed size** vector ($h_2$ in this example).

# Sequence-to-Sequence Models - Bottleneck

- This is not easy to do with very long input sequences.
- $h_n$ is a bottleneck.

fixed size!

| $h_0$ | $h_1$ | $h_2$ | $h_{...}$ | $h_{...}$ | $h_n$ | Decoder |

| je | pense | que | ... | ... | merci |

Encoder

Mila

# Attention Mechanism

- Problem: it is not easy to store all the necessary information from an **arbitrary long** sequence into a **fixed-size** vector.



Bahdanau et al., Neural Machine Translation by Jointly Learning to Align and Translate

# Attention Mechanism

- Problem: it is not easy to store all the necessary information from an **arbitrary long** sequence into a **fixed-size** vector.
- A possible solution can be to allow the decoder to "selectively look back" at the encoded input sequence.



Bahdanau et al., Neural Machine Translation by Jointly Learning to Align and Translate
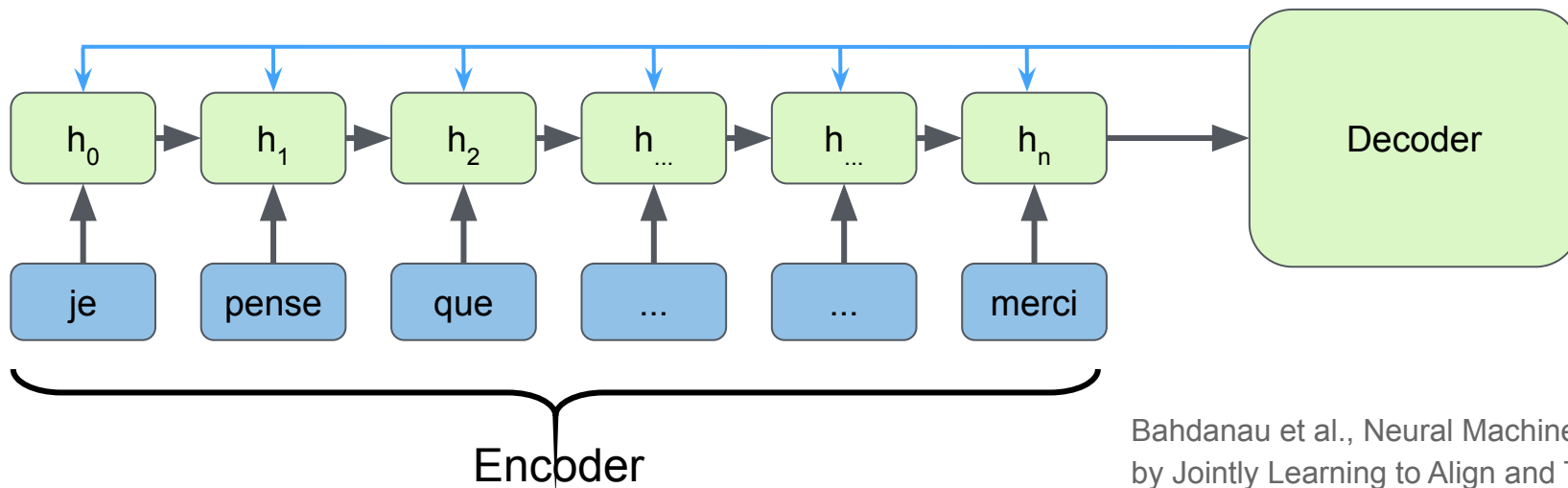
# Attention Mechanism

- This can be done with **attention**:
  - At any decoding time step, the decoder can use attention to fetch the relevant information for that step from the encoded input sequence.
- E.g., when producing the output word "think" (in a machine translation task), the decoder can focus on the encoding of the input word "pense".



Bahdanau et al., Neural Machine Translation by Jointly Learning to Align and Translate

Mila

# Attention Mechanism - Formalization

- Attention is a function **A** that, given a decoder state **s** and an encoded input sequence **h**, identifies the elements in **h** that are important at the current decoding time step.



Bahdanau et al., Neural Machine Translation by Jointly Learning to Align and Translate

# Attention Mechanism - Formalization

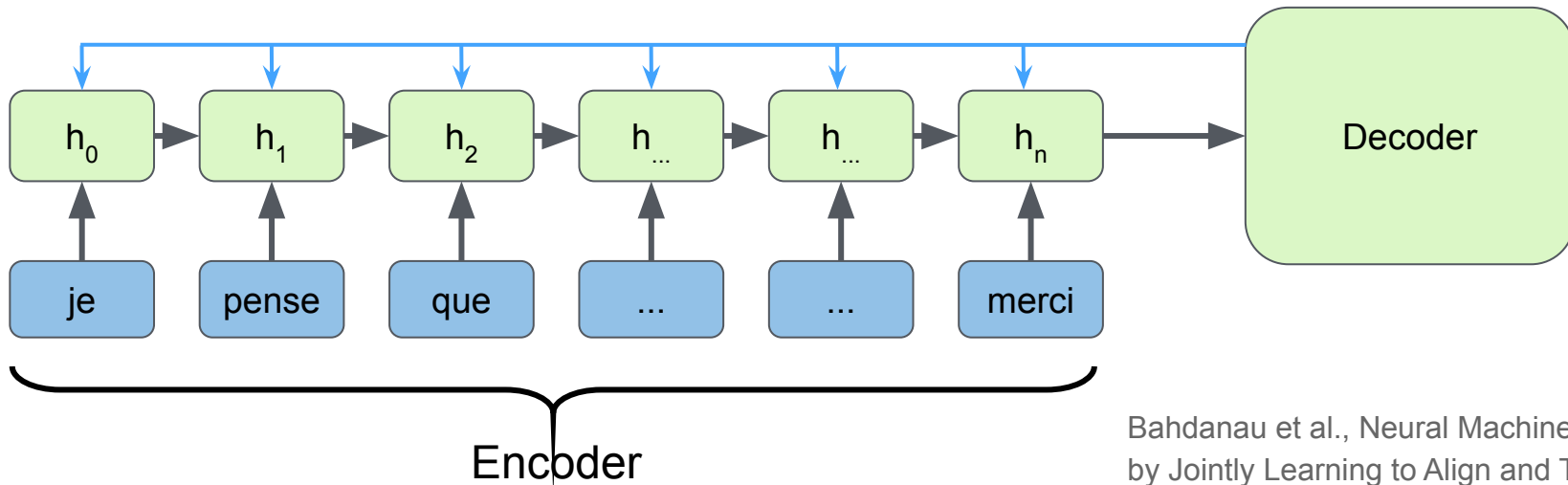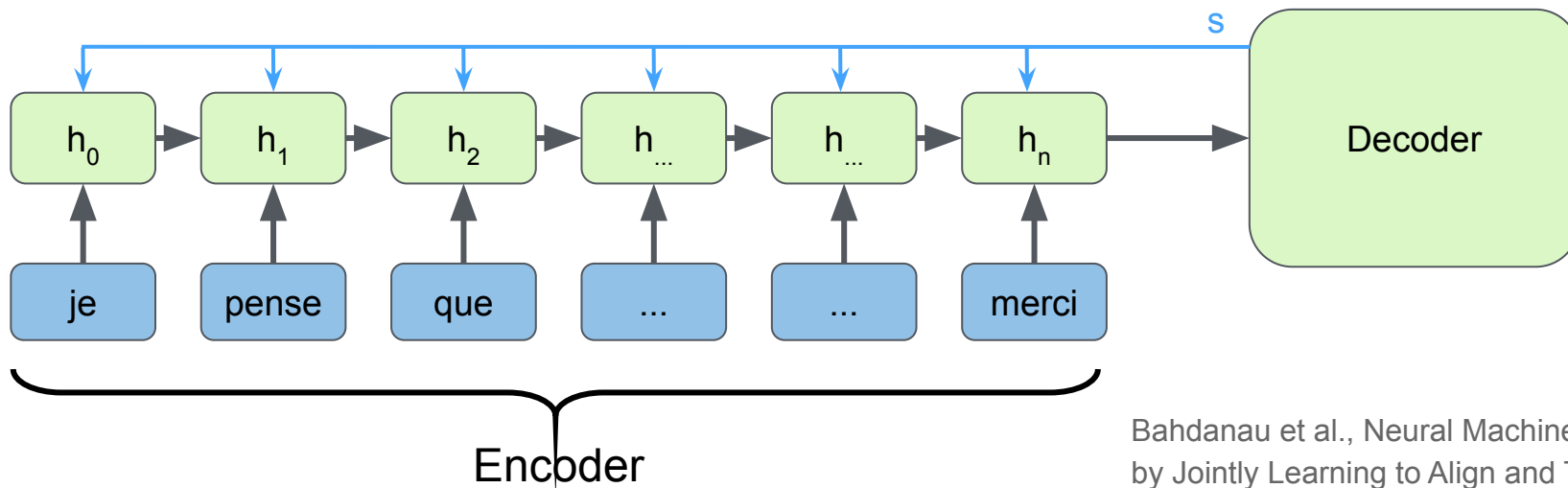- Attention is a function A that, given a decoder state s and an encoded input sequence h, identifies the elements in h that are important at the current decoding time step.
    - **A** assigns weights **w** to the elements in **h**.
    - Those weights are normalized (to sum to 1).



$w_0=0.1$   $w_1=0.6$   $w_2=0.05$   ...   …   $w_n=0.03$   s

| $h_0$ | $h_1$ | $h_2$ | $h_{...}$ | $h_{...}$ | $h_n$ | Decoder |

| je | pense | que | ... | ... | merci |

Encoder

Bahdanau et al., Neural Machine Translation by Jointly Learning to Align and Translate

# Attention Mechanism - Formalization

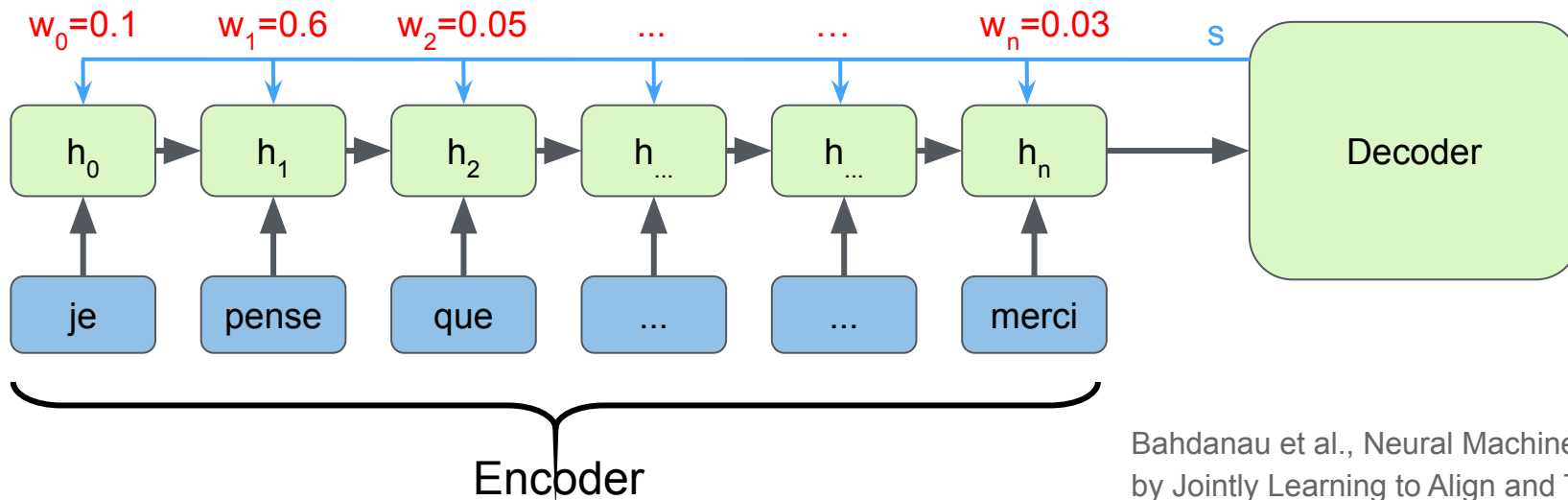- The weights **w** are used to compute a weighted sum **c** of the elements in the sequence **h**. **c** is called **context vector**.

$$c = \sum_{i=0}^{n} w_i \cdot h_i$$



c = $w_0 \cdot h_0$ + $w_1 \cdot h_1$ + $w_2 \cdot h_2$ + ... + ... + $w_n \cdot h_n$    s

| $h_0$ | $h_1$ | $h_2$ | $h_{...}$ | $h_{...}$ | $h_n$ | Decoder |

| je | pense | que | ... | ... | merci |

Encoder

# Attention Mechanism - Formalization

- The context vector **c** is then fed to the decoder.



$$c = w_0 \cdot h_0 \ + \ w_1 \cdot h_1 \ + \ w_2 \cdot h_2 \ + \ \ldots \ + \ \ldots \ + \ w_n \cdot h_n$$

s

| $h_0$ | $h_1$ | $h_2$ | $h_{\ldots}$ | $h_{\ldots}$ | $h_n$ | Decoder |

| je | pense | que | ... | ... | merci |

Encoder

Mila

# Attention Mechanism - Formalization

- Let's see a full step-by-step example of the attention mechanism.
- We will then look at how to implement the function **A**.

$$c = w_0 \cdot h_0 + w_1 \cdot h_1 + w_2 \cdot h_2 + \dots + \dots + w_n \cdot h_n$$

| $h_0$ | $h_1$ | $h_2$ | $h_{\dots}$ | $h_{\dots}$ | $h_n$ | Decoder |

| je | pense | que | ... | ... | merci |

Encoder

Mila

# Example: Sequence-to-Sequence + Attention

$x_0$    $x_1$    $x_2$

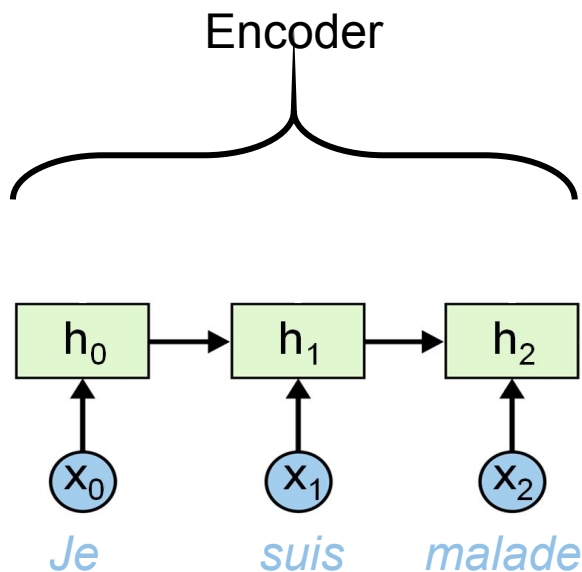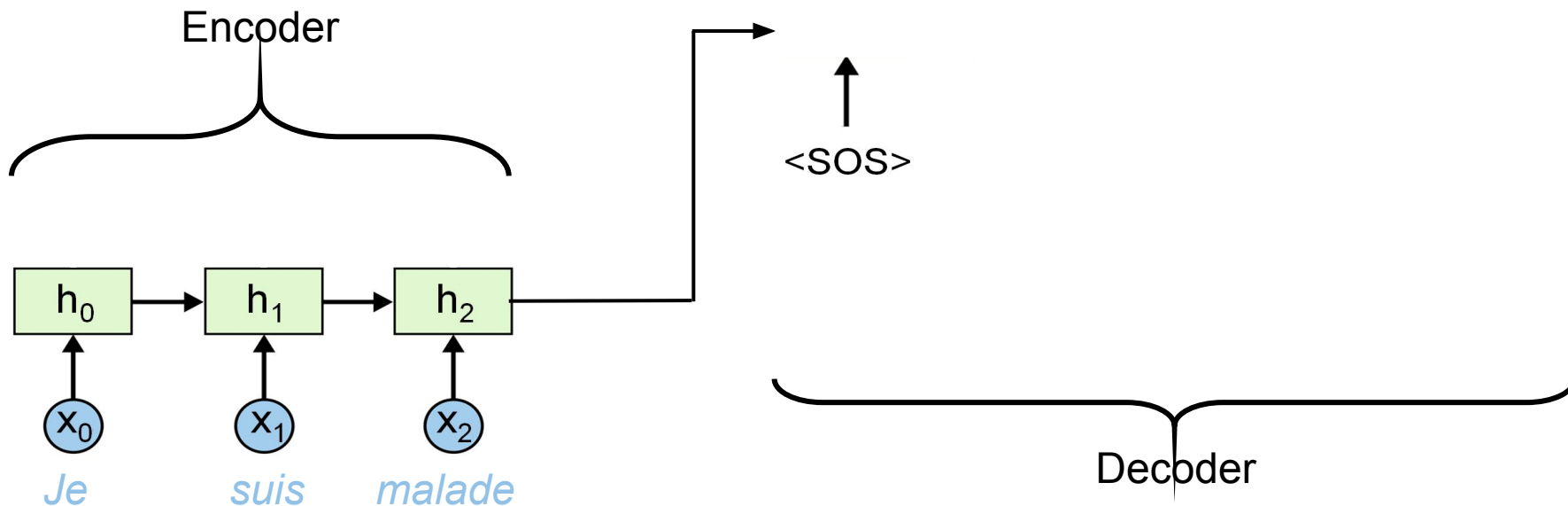*Je*    *suis*    *malade*

# Example: Sequence-to-Sequence + Attention

All the x sequence is encoded into a vector of **fixed size** (h2).

Encoder

| $h_0$ | $h_1$ | $h_2$ |

$x_0$  $x_1$  $x_2$
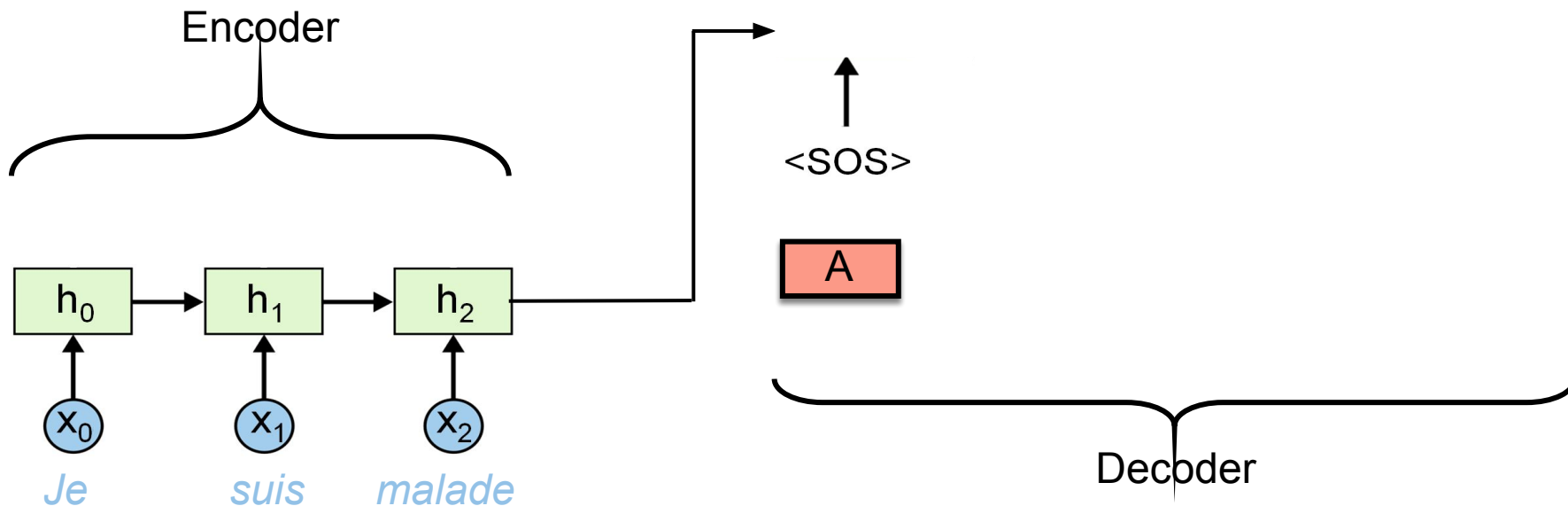
*Je*  *suis*  *malade*

# Example: Sequence-to-Sequence + Attention

The decoder starts with the **<SOS>** symbol.
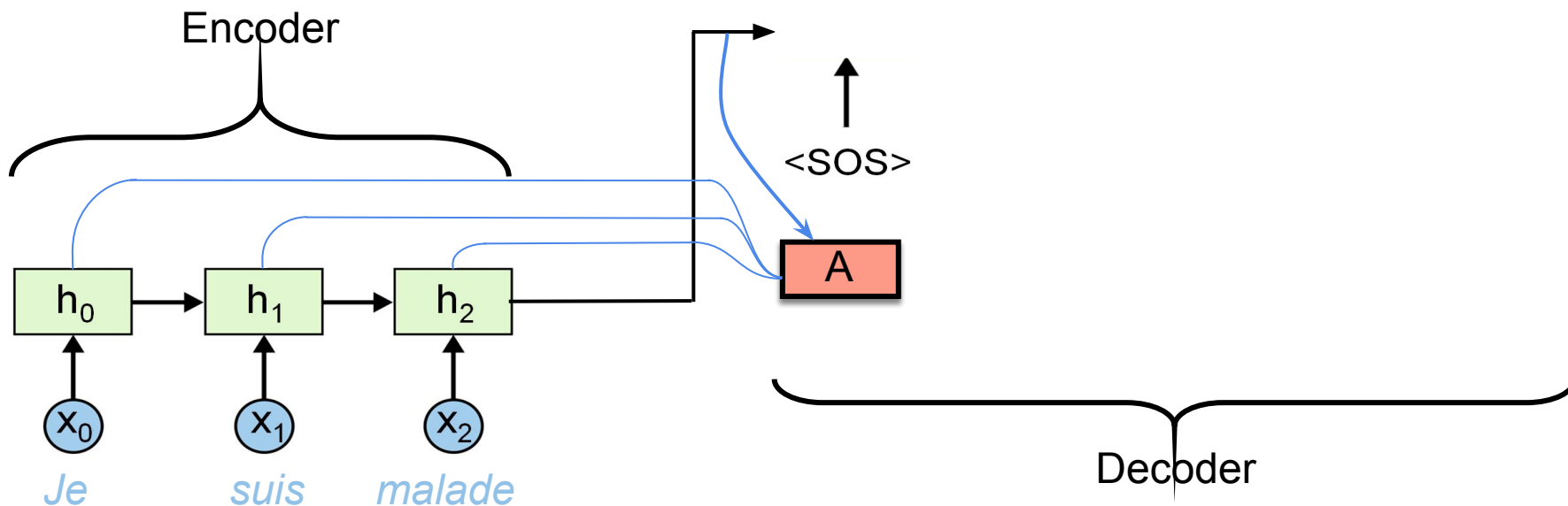
# Example: Sequence-to-Sequence + Attention

The attention model **A** is added to the decoder.

# Example: Sequence-to-Sequence + Attention

The decoder's previous state and the encoded input sequence **h** are fed as inputs to the attention model (**A**).

# Example: Sequence-to-Sequence + Attention

The context vector (output of the attention) is fed as an input to the decoder.

# Example: Sequence-to-Sequence + Attention

The internal state $s_0$ is computed.

# Example: Sequence-to-Sequence + Attention

The output **y₀** is computed and used as the next input.

# Example: Sequence-to-Sequence + Attention

# Example: Sequence-to-Sequence + Attention

# Example: Sequence-to-Sequence + Attention

# Attention Function

- There are several possible implementations for **A**.
- The most simple version is based on a dot product, i.e., $e_i = s_{t-1} \cdot h_i$ .

# Attention Function

- The dot product results are passed through a Softmax to get normalized weights $\mathbf{w}=[\mathbf{w}_0, ..., \mathbf{w}_n]$, which indicate how "important" the various elements are.
- The final result is the weighted sum of $\mathbf{h}$.

$$c = \sum_{i=0}^{n} w_i \cdot h_i$$

# Visualizing Attention



- The thick lines show where the decoder is focusing its attention when analyzing the encoded input sequence.

Image from Christopher Olah's blog

# Other Examples

**Image caption generation**

# Other Examples

A woman is throwing a frisbee in a park .



A     woman     is

throwing     a     frisbee     in

a     park     .

Xu et al., Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, 2015

Mila

# Plan

- RNN Recap
- Sequence to Sequence Models
- Attention Mechanism
- **Transformer**
- Libraries and References

# Beyond RNNs

- Sequence-to-Sequence + Attention systems perform well, but they are based on RNNs (simple RNNs / LSTMs / GRUs).
- RNNs suffer from two problems:
  - Not easy to parallelize.
  - Even in the more "complex" implementations (e.g. LSTMs), they struggle to capture (very) long-term dependencies.

Mila

# Beyond RNNs

- Every state in an RNN depends on the previous internal state.
- This creates a chain of computation which prevents parallelization.

# Beyond RNNs

- Long-term dependencies are hard to capture with RNNs.
- This problem is strongly mitigated using LSTMs / GRUs, but it's still there for very long sequences.

# Beyond RNNs

- There is no easy solution to deal with those RNN-related problems.

Mila

# Beyond RNNs

- To improve seq2seq systems, we need to find a replacement for RNNs.

# Beyond RNNs

- Note that we keep:
  - the encoder-decoder architecture,
  - the attention mechanism,
  - the autoregressive nature of the decoder.

# Transformer

- The Transformer architecture was introduced in the paper "**Attention is all you need**".

- Note: in the next slides we will focus on providing the intuition, thus simplifying some aspects of the architecture.



Vaswani et al, "Attention Is All You Need", 2017

# Transformer

- Several key points:
  - Recurrence replaced with self-attention and multi-head attention.
  - Positional encodings.
  - Residual connections.
  - Layer normalization.
  - Position-wise feed-forward networks.
- We will focus on the self-attention and the multi-head attention.



Vaswani et al, "Attention Is All You Need", 2017

Mila

# Self-Attention

- Before introducing Self-Attention, let's recap how a RNN works.
- At each time step, a RNN encodes the current input taking into consideration the past context (or the future context for right-to-left models).

- Example: the hidden state $h_2$ is encoding the information from the current input $x_2$ as well as the previous context.

# Self-Attention

- We want to do something similar with self attention:
  - encode the current input taking into consideration the surrounding context.
- The attention mechanism is used to identify the elements in a sequence which are "relevant" to encode the current one.

# Self-Attention

- Let's consider time step 0 with its element $x_0$.
- We will identify all elements in the sequence which are "relevant" to encode $x_0$.

# Self-Attention

- This is done by assigning a weight to each element (by computing a dot product between the element and $x_0$)…

| $w_{0,0}=0.2$ | $w_{0,1}=0.4$ | $w_{0,2}=0.3$ | $w_{0,3}=0.1$ | $w_{0,4}=0.0$ | $w_{0,5}=0.0$ |

$x_0$

| $x_0$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |

| t=0 | t=1 | t=2 | t=3 | t=4 | t=5 |

Mila

# Self-Attention

- This is done by assigning a weight to each element (by computing a dot product between the element and x0)…
- ... and then computing a normalized weighted sum.

$$h_0$$

$$x_0$$

$$h_0 = w_{0,0} \cdot x_0 \quad + \quad w_{0,1} \cdot x_1 \quad + \quad w_{0,2} \cdot x_2 \quad + \quad w_{0,3} \cdot x_3 \quad + \quad w_{0,4} \cdot x_4 \quad + \quad w_{0,5} \cdot x_5$$

| $x_0$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
|-------|-------|-------|-------|-------|-------|
| t=0 | t=1 | t=2 | t=3 | t=4 | t=5 |

Mila

# Self-Attention

- This is repeated for every step… (note that the weights vary across steps)



$h_1$

$h_1 = w_{1,0} \cdot x_0 \quad + \quad w_{1,1} \cdot x_1 \quad + \quad w_{1,2} \cdot x_2 \quad + \quad w_{1,3} \cdot x_3 \quad + \quad w_{1,4} \cdot x_4 \quad + \quad w_{1,5} \cdot x_5$

$x_1$

| $x_0$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |

t=0    t=1    t=2    t=3    t=4    t=5

# Self-Attention

- This is repeated for every step… (note that the weights vary across steps)
- … until all the steps are completed.

# Self-Attention - Multiple Heads

- The self-attention is meant to identify all elements in a sequence which are "relevant" to encode $x_t$.
- Given that there can be several types of relevant information, we can have several attention mechanisms.

Mila

# Self-Attention - Multiple Heads

- The self-attention is meant to identify all elements in a sequence which are "relevant" to encode $x_t$.
- Given that there can be several types of relevant information, we can have several attention mechanisms.
- Each attention mechanism is called a **head**, leading to a **multi-head self-attention**.
  - In this example, there are head#0 and head#1, each with different weights.

$$h^1_0 = w^1_{0,0} \cdot x_0 \quad + \quad w^1_{0,1} \cdot x_1 \quad + \quad w^1_{0,2} \cdot x_2 \quad + \quad w^1_{0,3} \cdot x_3 \quad + \quad w^1_{0,4} \cdot x_4 \quad + \quad w^1_{0,5} \cdot x_5$$

$$h^0_0 = w^0_{0,0} \cdot x_0 \quad + \quad w^0_{0,1} \cdot x_1 \quad + \quad w^0_{0,2} \cdot x_2 \quad + \quad w^0_{0,3} \cdot x_3 \quad + \quad w^0_{0,4} \cdot x_4 \quad + \quad w^0_{0,5} \cdot x_5$$

| $x_0$ |

| $x_0$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
| :---: | :---: | :---: | :---: | :---: | :---: |
| t=0 | t=1 | t=2 | t=3 | t=4 | t=5 |

Mila

# Self-Attention - Multiple Heads

- The various heads are then merged together.
  - E.g., they are concatenated, $\mathbf{h_0} = [\textcolor{red}{\mathbf{h^0_0}}, \textcolor{green}{\mathbf{h^1_0}}]$.

$h_0$

$$\textcolor{green}{h^1_0 = w^1_{0,0} \cdot x_0 \quad + \quad w^1_{0,1} \cdot x_1 \quad + \quad w^1_{0,2} \cdot x_2 \quad + \quad w^1_{0,3} \cdot x_3 \quad + \quad w^1_{0,4} \cdot x_4 \quad + \quad w^1_{0,5} \cdot x_5}$$

$$\textcolor{red}{h^0_0 = w^0_{0,0} \cdot x_0 \quad + \quad w^0_{0,1} \cdot x_1 \quad + \quad w^0_{0,2} \cdot x_2 \quad + \quad w^0_{0,3} \cdot x_3 \quad + \quad w^0_{0,4} \cdot x_4 \quad + \quad w^0_{0,5} \cdot x_5}$$

$\mathbf{x_0}$

| $x_0$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
|-------|-------|-------|-------|-------|-------|
| t=0 | t=1 | t=2 | t=3 | t=4 | t=5 |

# Self-Attention - Multiple Heads

- The weights for a given head are based on a dot-product.
  - E.g., $\mathbf{w}^0_{3,4} = \mathbf{x}^0_3 \cdot \mathbf{x}^0_4$

- The dot-product is computed in a different space for each attention head. This space is obtained by learning a projection from the original space to the one dedicated to a particular head.

**Original** space

$x_3$

$x_4$

**head#0** space

$x^0_4$

$x^0_3$

**head#1** space

$x^1_4$

$x^1_3$

# Self-Attention - Visualization



The Law will never be perfect, but its application should be just - this is what we are missing, in my opinion. <EOS> <pad>

The Law will never be perfect, but its application should be just - this is what we are missing, in my opinion. <EOS> <pad>

Vaswani et al, "Attention Is All You Need", 2017

Mila

# Self-Attention - Advantages

- The multi-head self-attention can be computed in parallel at all time steps.
- There are no dependencies between time steps.

# Self-Attention - Advantages

- The RNN chain of computation is **not** there anymore.
- The information does **not** need to flow over a long chain of elements.
- E.g., $x_5$ has direct access to both $x_0$ and $x_4$.



| $x_0$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
|-------|-------|-------|-------|-------|-------|
| **t=0** | **t=1** | **t=2** | **t=3** | **t=4** | **t=5** |

# Plan

- RNN Recap
- Sequence to Sequence models
- Attention Mechanism
- Transformer
- **Libraries and References**

# Libraries

- RNNs are included in the main DL frameworks:
  - PyTorch : https://pytorch.org/docs/stable/nn.html#recurrent-layers
  - Tensorflow: https://www.tensorflow.org/tutorials/recurrent
- There are several Transformer implementations:
  - in Tensorflow: https://github.com/tensorflow/tensor2tensor
  - in PyTorch: https://github.com/huggingface/pytorch-transformers

Mila

# References

- Christopher Olah's blog about LSTMs:
  http://colah.github.io/posts/2015-08-Understanding-LSTMs/
- Christopher Olah's publications on the attention mechanism:
  https://distill.pub/2016/augmented-rnns/
- Andrej Karpathy's blog about RNNs:
  http://karpathy.github.io/2015/05/21/rnn-effectiveness/
- The Deep Learning Book (Goodfellow et al.): http://www.deeplearningbook.org/

Mila

Quebec
Artificial
Intelligence
Institute

Mila

Questions?

# Self-Attention

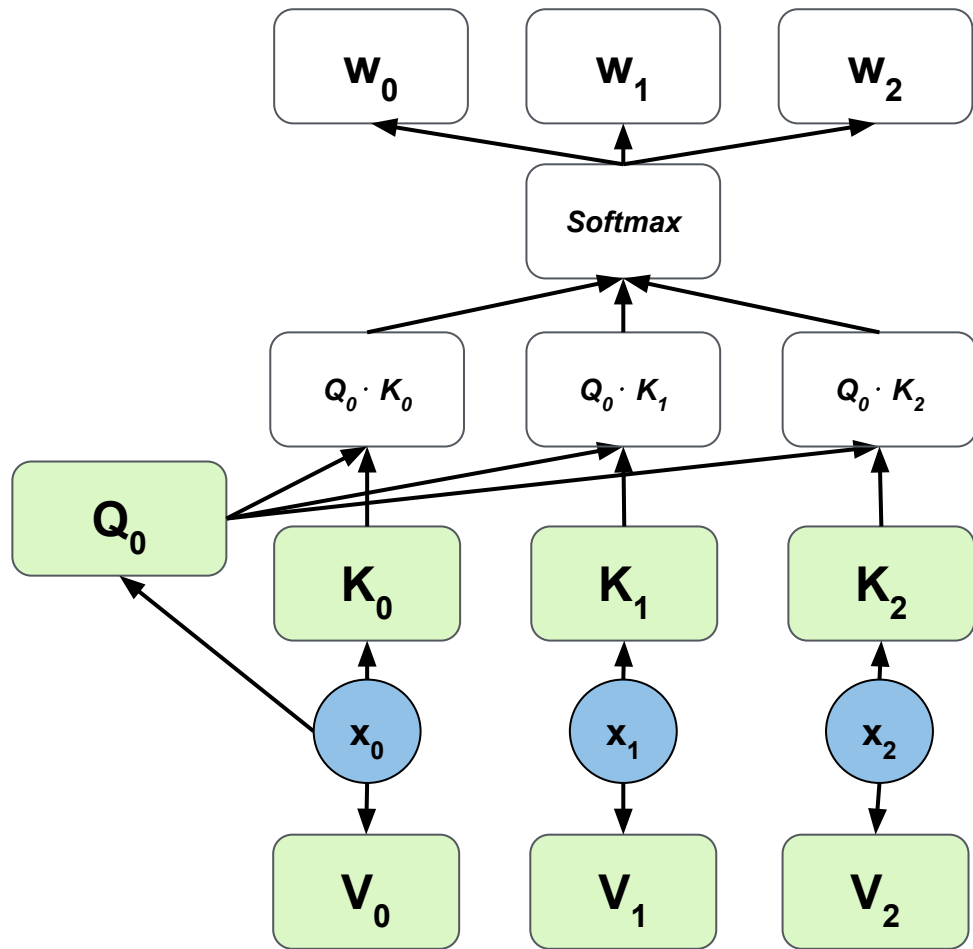- In reality, every input (x) has 3 different "views": Q, K, V.
  - Q: query - used to represent the current state (called **s** before).
  - K: keys - used for the dot-product (to compute the weights - called **h** before).
  - V: values - used in the final weighted sum.
- In this example, we are focusing on $x_0$.

$$x_0' = \sum_{i=0}^{n} w_i \cdot V_i$$

# Decoder Probabilities

- Want to maximize:

$$p(y_1, y_2, \ldots, y_n | x)$$

- We can decompose as:

$$p(y_1, y_2, \ldots, y_n | x) = p(y_1 | x) p(y_2, \ldots, y_n | x, y_1)$$

and again:

$$p(y_1, y_2, \ldots, y_n | x) = p(y_1 | x) p(y_2 | x, y_1) p(y_3, \ldots, y_n | x, y_1, y_2)$$

until:

$$p(y_1, y_2, \ldots, y_n | x) = p(y_1 | x) p(y_2 | x, y_1) \ldots p(y_n | x, y_1, \ldots y_{n-1})$$

Mila

# Results - Seq2Seq W/ and W/O Attention

| Model | Max sequence length | BLEU score |
|---|---|---|
| Seq2seq | 30 | 13.9 |
| Seq2seq + attention | 30 | **21.5** |
| Seq2seq | 50 | 17.8 |
| Seq2seq + attention | 50 | **26.7** |

- The BLEU score (BiLingual Evaluation Understudy) measures the quality of the translation (the higher the score the better).
  - BLEU is a modified form of the Precision metric based on the overlap between output and target.
- Attention improves the score significantly.

Bahdanau et al., Neural Machine Translation by Jointly Learning to Align and Translate

Mila

# Results - Transformer

- Better results than previous systems (left column)...
- … with a fraction of the training cost (right column).

| Model | BLEU | | Training cost (FLOP) | |
|-------|------|------|------|------|
| | **EN-DE** | **EN_FR** | **EN-DE** | **EN_FR** |
| *GNMT (RNNs)* | 24.6 | 39.92 | $2.3*10^{19}$ | $1.4*10^{20}$ |
| *Transformer (base)* | 27.3 | 38.1 | $3.3*10^{18}$ | |
| *Transformer (big)* | 28.4 | 41.0 | $2.3*10^{19}$ | |

- BLEU score (BiLingual Evaluation Understudy): the higher the score the better.

Mila

# Self-Attention - Word Order

- The (weighted) sum leads to the loss of the elements' order.
  - E.g., $w_0 \cdot x_0 + w_1 \cdot x_1 = w_1 \cdot x_1 + w_0 \cdot x_0$

- Solution: attach the position information to every word.
  - E.g., $x'_0 = [x_0 + p_0]$, $x'_1 = [x_1 + p_1]$

- Now the order of the elements has an impact on the results:

  $w_0 \cdot [x_0 + p_0] + w_1 \cdot [x_1 + p_1] \neq w_1 \cdot [x_1 + p_0] + w_0 \cdot [x_0 + p_1]$

Mila