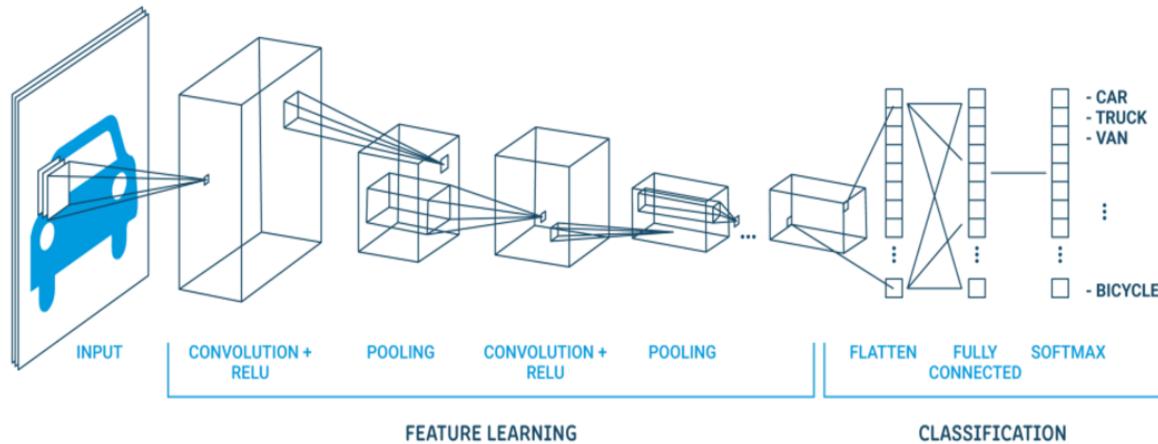


# RNNs and Attention

Arsalan Mosenia (am12546@nyu.edu)

# Limitations of CNNs



- Fixed size inputs (example: images of a given resolution)
  - Say I design a CNN, can you use it for inputs of a different dimension?
- Inherent spatial structure: spatially localized kernels/filters

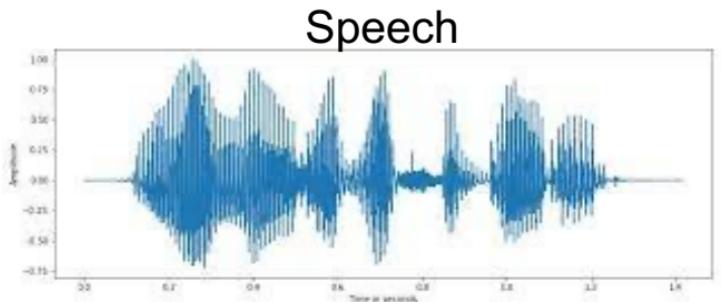
固有的空间结构

空间小范围的

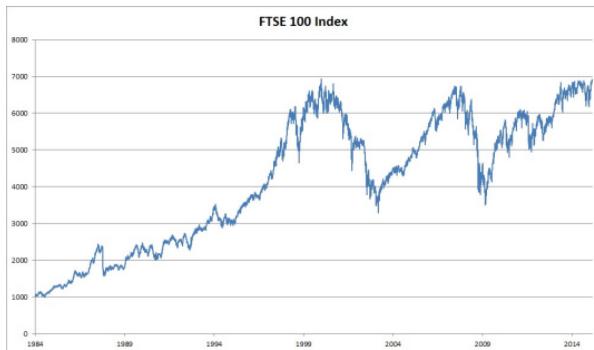
# Sequential Inputs



Video Frames



# Properties of Sequential Inputs



- Can be viewed as an ordered set of many inputs
  - Example: a video is a sequence of *many* images

固有因果模型

- Inherent causal structure (not always)
  - Future depends only on the past -> no time travel!
  - Example: the value of a stock at time  $t$  depends only what happened in the past

# Stock Prediction

- Say you would like to predict the value of a stock  $x_t$  at time  $t$
- Start by building a probabilistic model of the form:

$$x_t \sim P(x_t \mid x_{t-1}, \dots, x_1).$$



Function  
of  $(t-1)$   
inputs

- What's the problem here?
  - Variable length function
  - You need a new function (*and more complex function*) for every time step!

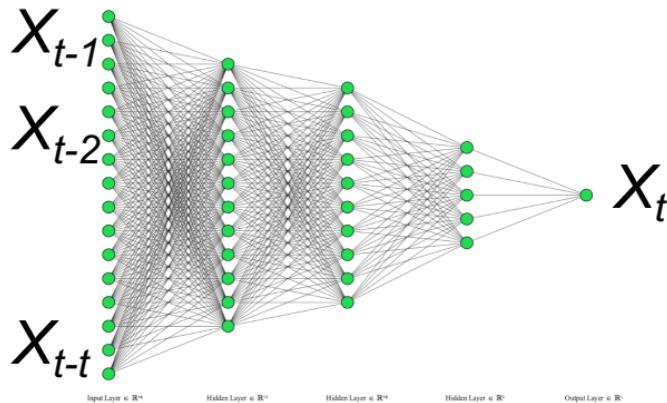
# Autoregressive Models

- Strategy 1: only look at  $t$  time-steps in the past

$$x_t \sim P(x_t \mid x_{t-1}, \dots, x_1).$$



$$x_t \sim P(x_t \mid x_{t-1}, \dots, x_{t-\tau}).$$



Drawbacks?

- Models can still be quite large *like image*
- Old history forgotten

*lose important value*

*explosively model  
x\_t: image  $\rightarrow$  large*

# Latent Autoregressive Models

- **Strategy 2:** use a “state” variable  $h_t$  to store a summary of the history
  - For example, the state could simply count and keep track of time. Of course, in practice, the state will do much more complex things.
  - How do we model state?

Summary of data  
we've seen so far

$$h_t = g(h_{t-1}, x_{t-1})$$

Current state      Past state      Past input

$$\hat{x}_t = P(x_t | h_t)$$

Prediction  
depends on  
state

# Latent Autoregressive Models

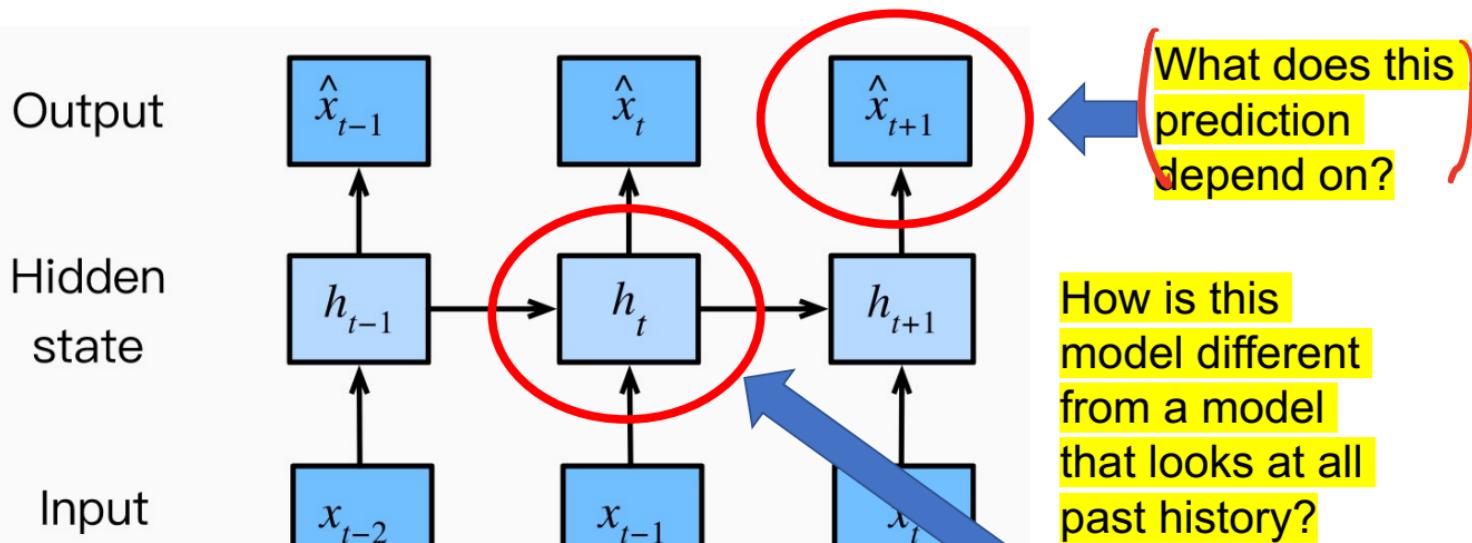


Fig. 8.1.2 A latent autoregressive model.

RNN

# Recurrent Neural Networks

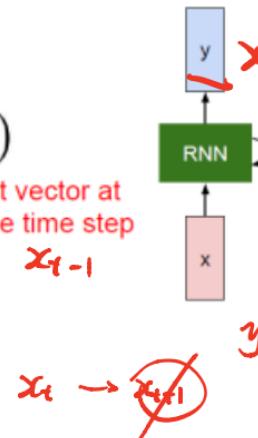
Small mod: we'll use past state and current input to output current state

We can process a sequence of vectors  $x$  by applying a recurrence formula at every time step:

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

new state      old state      input vector at some time step  
some function with parameters  $W$       *not*  $x_{t-1}$

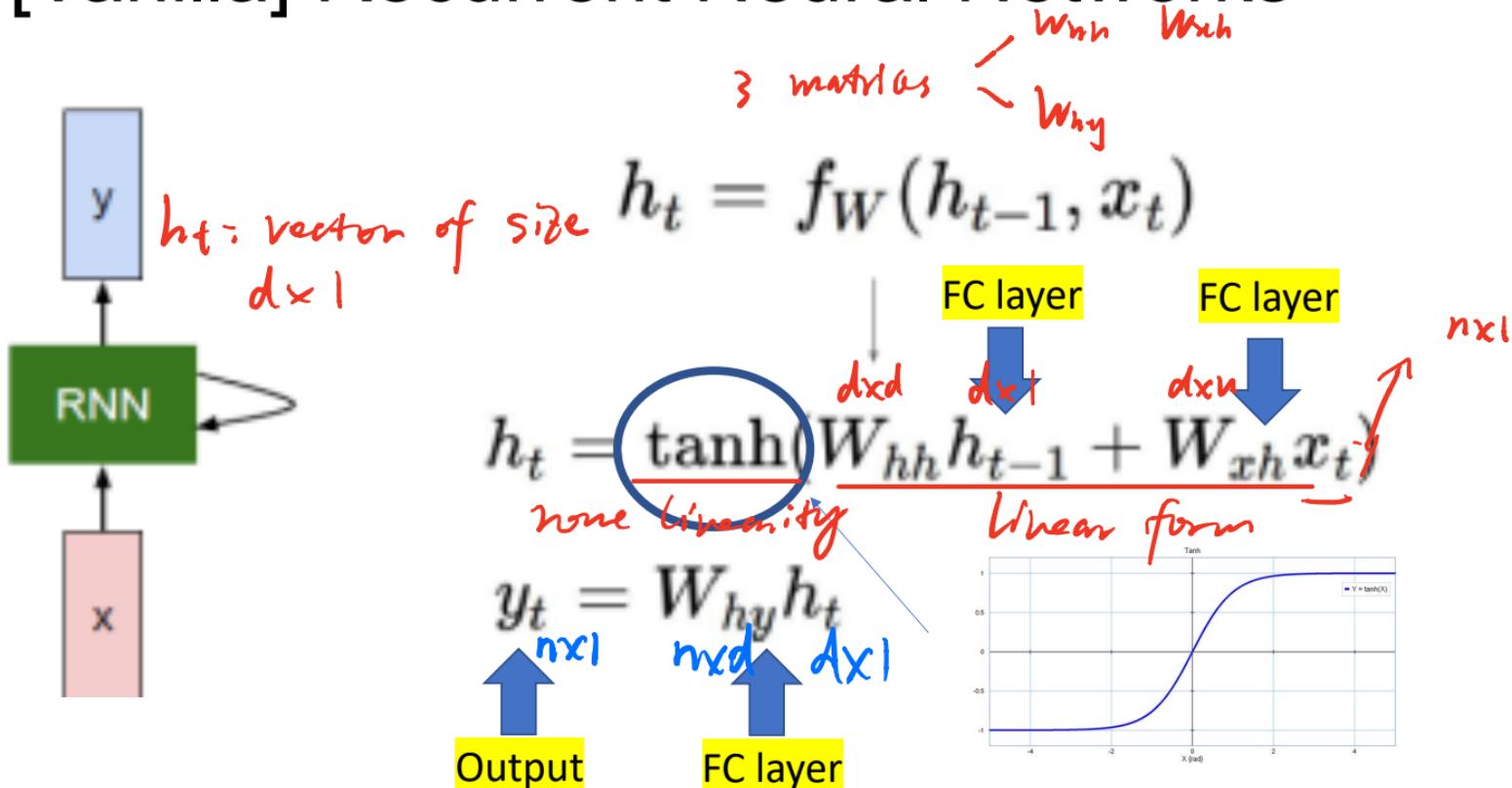
$g = f_W = \text{Neural network}$



What does  $f_W$  look like?

A neural network with  $N \geq 1$  layers!

# [Vanilla] Recurrent Neural Networks



## Character-level language model example

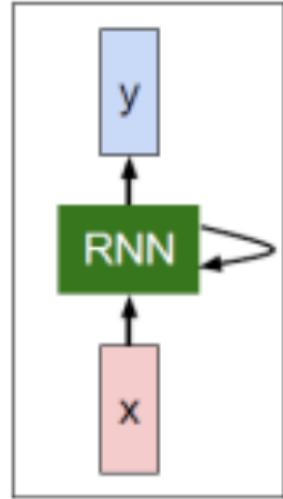
Vocabulary:  
[h,e,l,o]

Example training sequence:  
“hello”

**Goal:** Train a model to auto-complete a word, assuming a vocabulary of 4 letters [h,e,l,o]

**One-hot coding:**

h = [1,0,0,0]  
e = [0,1,0,0]  
l = [0,0,1,0]  
o = [0,0,0,1]



# [Vanilla] Recurrent Neural Networks

Assume a model has been trained! (we'll see later how to train the model)

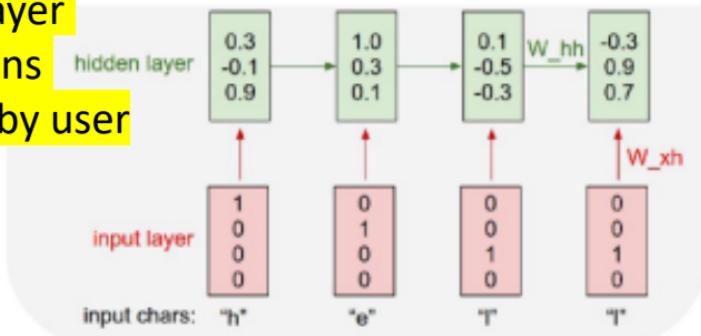
Character-level  
language model  
example

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
“hello”

Hidden layer  
dimensions  
selected by user

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

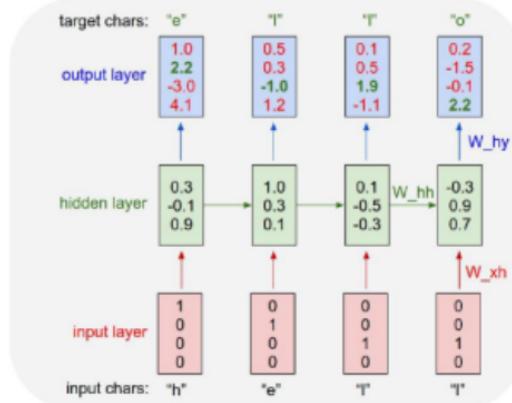


# [Vanilla] Recurrent Neural Networks

Character-level  
language model  
example

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
“hello”



What happens if we input “h”, “e”, “l”, “l”?

## Examples

open source textbook on algebraic geometry

## Latex source

<p><b>Proof.</b> Quotiated. <math>\square</math></p> <p><b>Lemma 0.1.</b> Let <math>C</math> be a set of the construction.  <math>\text{Let } C \text{ be a geer covering. Let } F \text{ be a quasi-coherent sheaf of } \mathcal{O}\text{-modules. We have to show that}</math></p> $\mathcal{O}_X(C) = \mathcal{O}_X(C)$ <p><b>Proof.</b> This is an algebraic space with the composition of datum <math>\mathcal{F}</math> on <math>X_{\text{max}}</math> we have</p> $\mathcal{O}_X(C) = (\text{morph}, s_{C,F}((\mathcal{F}))$ <p>where <math>\mathcal{G}</math> defines an isomorphism <math>\mathcal{F} \rightarrow \mathcal{G}</math> of <math>\mathcal{O}</math>-modules.</p> <p><b>Lemma 0.2.</b> This is an affine <math>Z</math> is injective.</p> <p><b>Proof.</b> See Lemma 0.1.</p> <p><b>Lemma 0.3.</b> Let <math>S</math> be a scheme. Let <math>X</math> be a scheme and <math>Z</math> is an affine open subscheme. Let <math>\mathcal{F}</math> be a <math>\mathcal{O}_Z</math>-module of finite type. Let <math>X</math> be a scheme. Let <math>\mathcal{F}'</math> be a <math>\mathcal{O}_X</math>-module is equal to the formal completion.</p> <p>The following to the construction of the lemma follows.</p> <p>Let <math>X</math> be a scheme. Let <math>Z</math> be a scheme covering. Let</p> $h: X \rightarrow Y \rightarrow Y' = Y \times^Z Y \rightarrow X,$ <p>be a morphism of algebras spaces over <math>S</math> and <math>Y</math>.</p> <p><b>Proof.</b> Let <math>X</math> be a noetherian scheme of <math>X</math>. Let <math>Z</math> be an algebraic space. Let <math>\mathcal{F}</math> be a quasi-coherent sheaf of <math>\mathcal{O}</math>-modules. The following are equivalent</p> <ul style="list-style-type: none"> <li>(1) <math>\mathcal{F}</math> is an algebraic space over <math>S</math>.</li> <li>(2) <math>\mathcal{F}</math> is an open covering.</li> </ul> <p>Consider a coarsening structure on <math>X</math> and <math>Z</math> the functor <math>\mathcal{O}_X(U)</math> which is locally of finite type</p>	<p>This time <math>\mathcal{F} = F</math> and <math>c = G</math> the diagram</p> <pre>     \begin{array}{ccc}     \mathcal{O}_X(C) &amp; \xrightarrow{\quad} &amp; \mathcal{O}_X(C) \\     \downarrow &amp; &amp; \downarrow \\     \mathcal{O}_Z(G) &amp; \xrightarrow{\quad} &amp; \mathcal{O}_X(C)     \end{array}   </pre> <p>and <math>\mathcal{O}_Z(G)</math></p> <p>is a <math>\mathcal{O}</math>-algebra. Then <math>\mathcal{O}_Z(G)</math> is a field and <math>\mathcal{O}_X(C)</math> is a set. If <math>G</math> is a set and <math>F</math> is a finite type, <math>Z</math> is of finite type, always. Then</p> <ul style="list-style-type: none"> <li>• the composition <math>\mathcal{F} \rightarrow \mathcal{O}_Z(G)</math> is a regular sequence;</li> </ul> <p><b>Proof.</b> We know that <math>X = Z</math> and <math>\mathcal{F}</math> is a finite type representable by algebraic space. The property <math>\mathcal{F}</math> is a finite cover of a regular scheme. Then <math>\mathcal{O}_X(C)</math> is a <math>\mathcal{O}</math>-algebra. Moreover <math>Z</math> is an open neighborhood of <math>Z</math>. Then <math>\mathcal{O}_Z(G)</math> is a <math>\mathcal{O}</math>-algebra. Then we know <math>\mathcal{F}</math>. A refined scheme where we know that <math>\mathcal{F}</math> is an open covering of <math>C</math>. The functor <math>\mathcal{O}_X</math> is a "full"</p> $\mathcal{O}_X \rightarrow \mathcal{F} = \mathcal{O}_{Z(\mathcal{O}_X(C))} \rightarrow \mathcal{O}_{Z(\mathcal{O}_X(C))}$ <p>is a presentation of covering of <math>Z</math> by <math>\mathcal{O}_X(C)</math> the unique element of <math>Z</math> such that <math>X</math> is an open neighborhood of <math>Z</math>.</p> <p>The property <math>\mathcal{F}</math> is a regular space of <math>\mathcal{O}_{Z(\mathcal{O}_X(C))}</math> and we can refine this covering. Then <math>\mathcal{O}_Z(G)</math> is a <math>\mathcal{O}</math>-algebra. Then <math>\mathcal{F}</math> is a presentation of <math>\mathcal{O}_Z(G)</math> if <math>\mathcal{F}</math> is a regular space of <math>\mathcal{O}_{Z(\mathcal{O}_X(C))}</math>.</p> <p>If <math>\mathcal{F}</math> is a finite étale sheaf, <math>\mathcal{F}</math> is a closed immersion, see Lemma 0.2. This is a</p>
---	--

## Sonnet 116 – Let me not ...

*by William Shakespeare*

Let me not to the marriage of true minds  
Admit impediments. Love is not love  
Which alters when it alteration finds,  
Or bends with the remover to remove:  
O no! it is an ever-fixed mark  
That looks on tempests and is never shaken;  
It is the star to every wandering bark,  
Whose worth's unknown, although his height be taken.  
Love's not Time's fool, though rosy lips and cheeks  
Within his bending sickle's compass come:  
Love alters not with his brief hours and weeks,  
But bears it out even to the edge of doom.  
If this be error and upon me proved,  
I never writ, nor no man ever loved.

at first:

tyntd-iafhatawiaoahrdemot lytdws e ,tfti, astai f ogoh eoase rrwanbyne 'mhthnee  
nja tklrod t o idoo ns smtt b re etia h brents pictike acanpas lnp

**train more**

"Mont thithey" fomescerlundi  
Keushay. Thos here  
sheulke, ammerenith ol svih I lalterthend Bleipile shuiwy fil on aseterlome  
coaniogennc Phe lish thond hon at. MeDimorion in ther thize."

[train more](#)

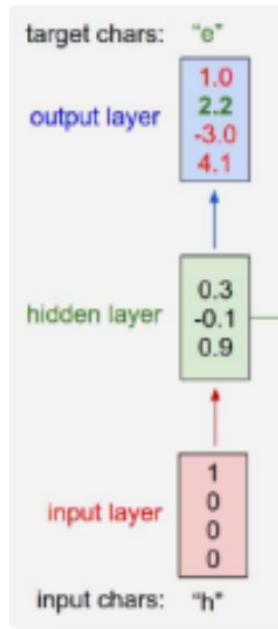
Aftair fall unsuch that the hall for Prince Velzonski's that me of her heartly, and behs to so arwage fiving were to it beloge, pavu say falling misfort how, and Gogitton is so overleical and after.

train more

"Why do what that day," replied Natasha, and wishing to himself the fact the princess, Princess Mary was easier, fed in had oftened him.  
Pierre aking his soul came to the packs and drove up his father-in-law women

# Generating New Words

- Once we've trained the model, how do we use it?



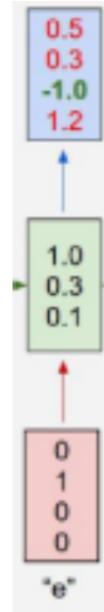
$$p_i = \frac{e^{y_i}}{\sum e^{y_i}}$$

RNN output  
defines a  
probability  
distribution  
over characters



Sample from  
distribution

'e'



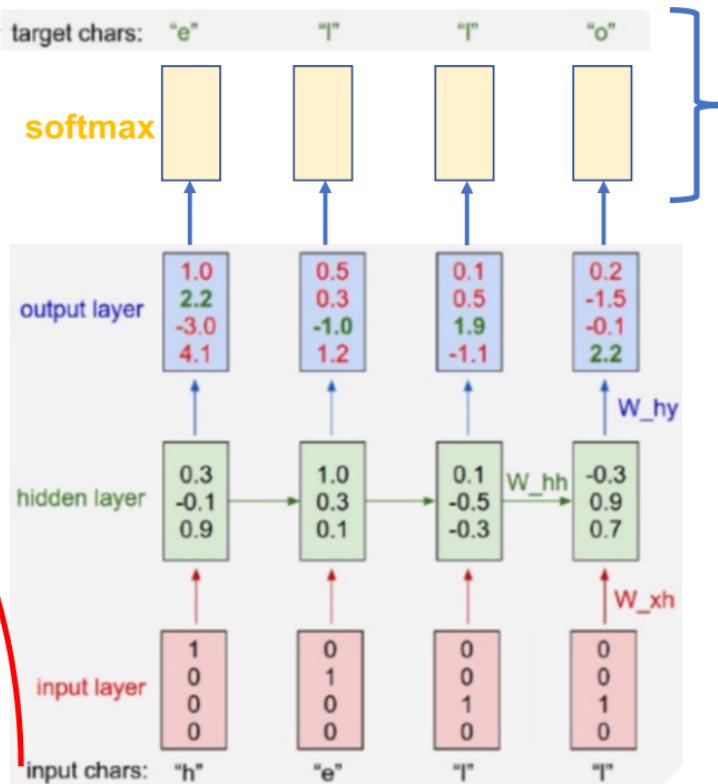
Sample from  
distribution

'h'

Note: with some  
probability you  
might pick every  
other letter also

# RNN Training

Input sequence shifted right by 1



Cross-entropy loss between target chars and predictions, averaged over the T predictions.

Select T (# of time-steps) for training.  
Let's select T=4.  
The RNN is  
“unrolled” T times....

# Training Data Creation

- In practice, you have a loooo...ooong sequence of text. For example, an entire book. We will use HG Wells' book, *The Time Machine* as an example and assume T=5.

the time machine by h g wells

The entire book can be partitioned into  $T=5$  length character strings starting with an offset 0.

Different offsets give different sequences. How should we create our training data and batches?

Fig. 8.3.1 Different offsets lead to different subsequences when splitting up text.

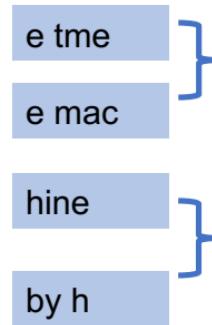
# Sequential Sampling

- In each epoch

```
the time machine by h g wells  
the time machine by h g wells  
the time machine by h g wells  
the time machine by h g wells ← blue arrow  
the time machine by h g wells  
the time machine by h g wells
```

Fig. 8.3.1 Different offsets lead to different subsequences when splitting up text.

Start with a random offset (ex: O=2)



We set T=5.  
Mini-batch size B=2

Samples from adjacent mini-batches are adjacent in the original sequence

# Random Sampling

Recall, we set  $T=5$ .  
Assume mini-batch size  $B=2$

- In each epoch

the time machine by h g wells  
the time machine by h g wells

Start with a  
random  
offset (ex:  
 $O=2$ )

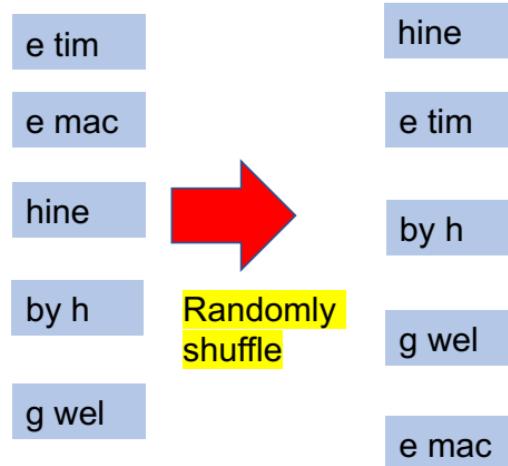


Fig. 8.3.1 Different offsets lead to different subsequences when splitting up text.



1. avoid overfitting
2. having a slightly different distribution in your dataset, more variance

Samples from adjacent mini-batches are NOT adjacent in the original sequence

# Real-world Problems are Complex

```
Source code of C ++ 'Bubble sort'  
// This procedure realizes Bubble sort  
void BubbleSort(int a[], int n)  
{  
    int right = n-1;  
    int temp = 0;  
    while(right > 0)  
    {  
        int k = 0;  
        for(int i = 0; i < right; i++)  
            if(a[i+1] < a[i])  
            {  
                temp = a[i+1];  
                a[i+1] = a[i];  
                a[i] = temp;  
                k = i;  
            }  
        right = k;  
    }  
}
```

What does this code do?

Sometimes the earliest tokens might be the most informative!

I am writing about topic X...

On a different note...

Back to topic X....

What is the topic of this email?

Sometimes intermediate tokens might be irrelevant!

GRU  
LSTM } extensions of RNN

How are you doing?  
Hope you've been keeping well...

I am writing because...

What is the topic of this email?

Sometimes you might want to “reset” your RNN

# Going Beyond Character Level Modeling

- In practice, character level models are not very successful. We would ideally like to encode more semantically meaningful inputs
- **Tokenization**: partitions inputs into meaningful tokens (example: characters, words, combinations of words, or sub-words)

MXNET

PYTORCH

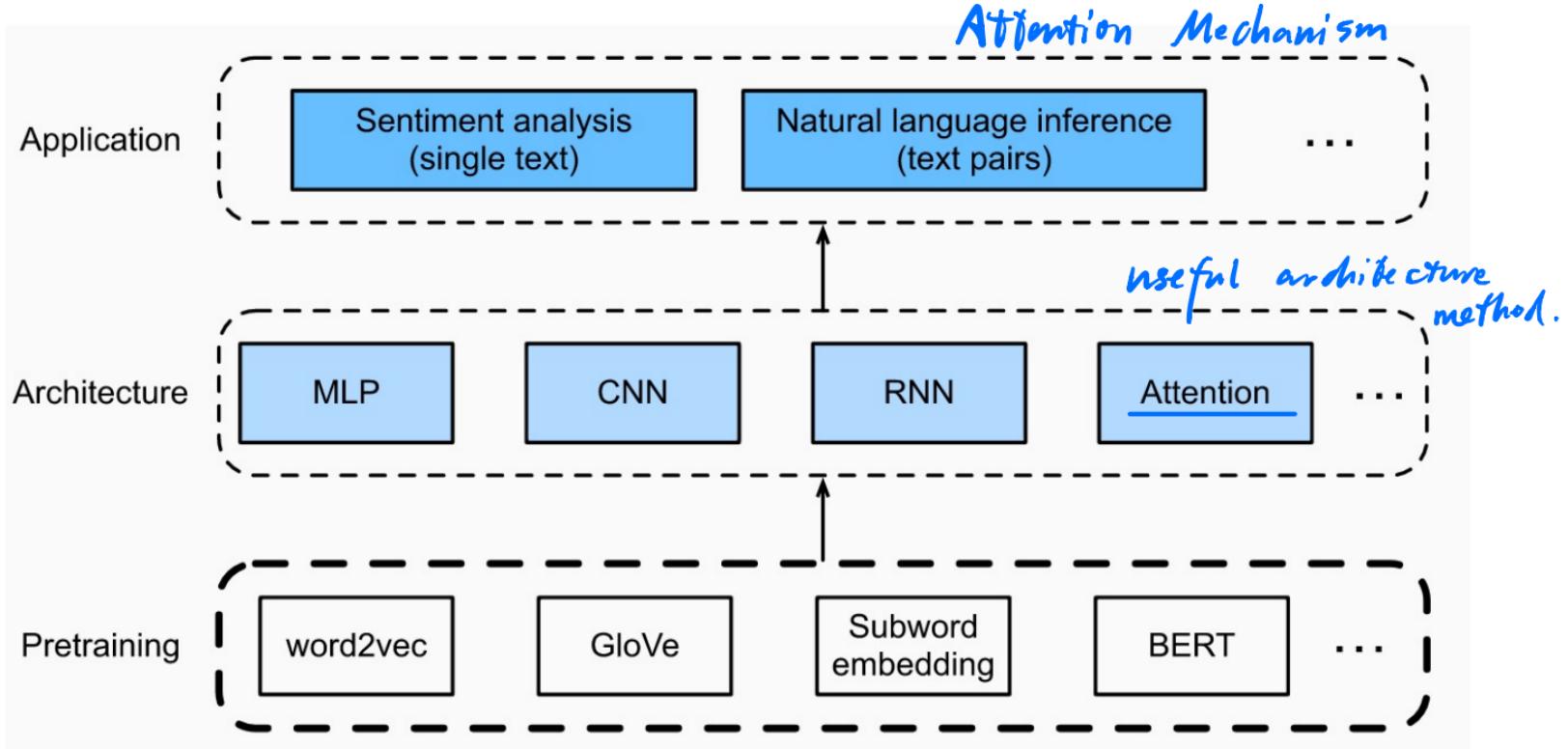
TENSORFLOW

```
def tokenize(lines, token='word'):    #@save
    """Split text Lines into word or character tokens."""
    if token == 'word':
        return [line.split() for line in lines]
    elif token == 'char':
        return [list(line) for line in lines]
    else:
        print('ERROR: unknown token type: ' + token)

tokens = tokenize(lines)
for i in range(11):
    print(tokens[i])
```

```
['the', 'time', 'machine', 'by', 'h', 'g', 'wells']
[]
[]
[]
[]
[]
['i']
[]
[]
['the', 'time', 'traveller', 'for', 'so', 'it', 'will', 'be', 'convenient', 'to', 'spe
['was', 'expounding', 'a', 'recondite', 'matter', 'to', 'us', 'his', 'grey', 'eyes', '
['twinkled', 'and', 'his', 'usually', 'pale', 'face', 'was', 'flushed', 'and', 'animat
```

# NLP Overview



# NLP Applications

- Machine Translation
- Sentiment Analysis (e.g., Social Media Monitoring)
- Question-answering (e.g., Chatbots)
- Hiring and Recruitment
- Automatic Summarization
- Spam Filtering
- Voice Assistants

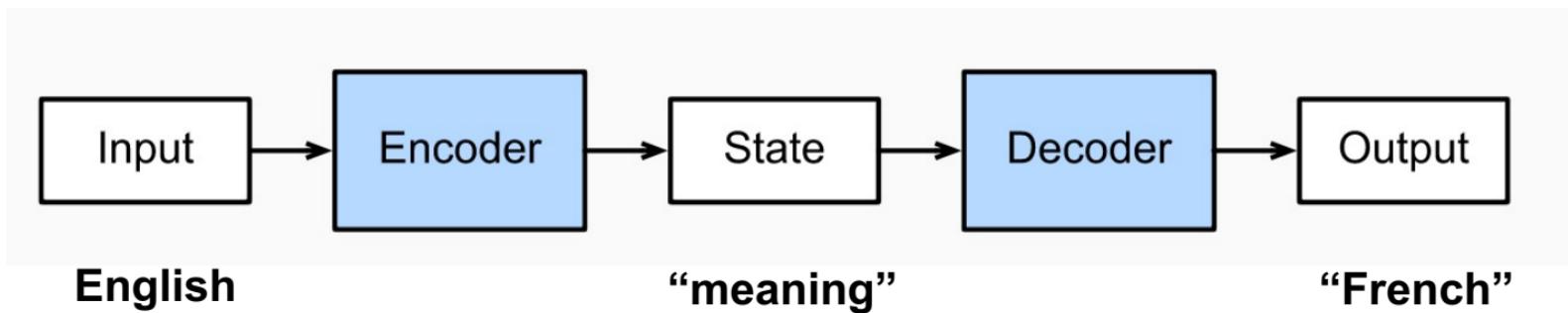
# Machine Translation

- Statistical Machine Translation (SMT)
- Rule-based Machine Translation (RBMT)
- Hybrid Machine Translation (HMT)
- Neural Machine Translation (NMT)

Classic Sequence-to-Sequence Problem

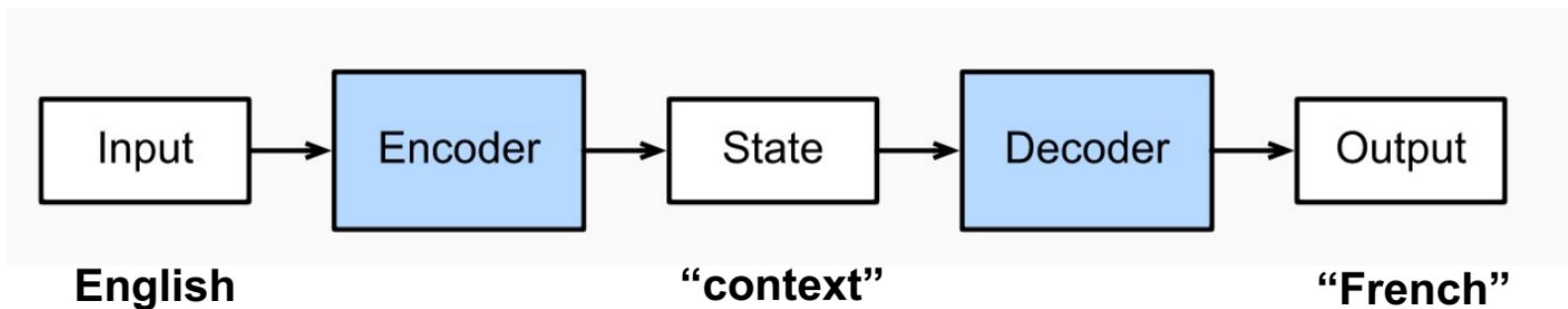
# Neural Machine Translation (NMT)

- Language translation

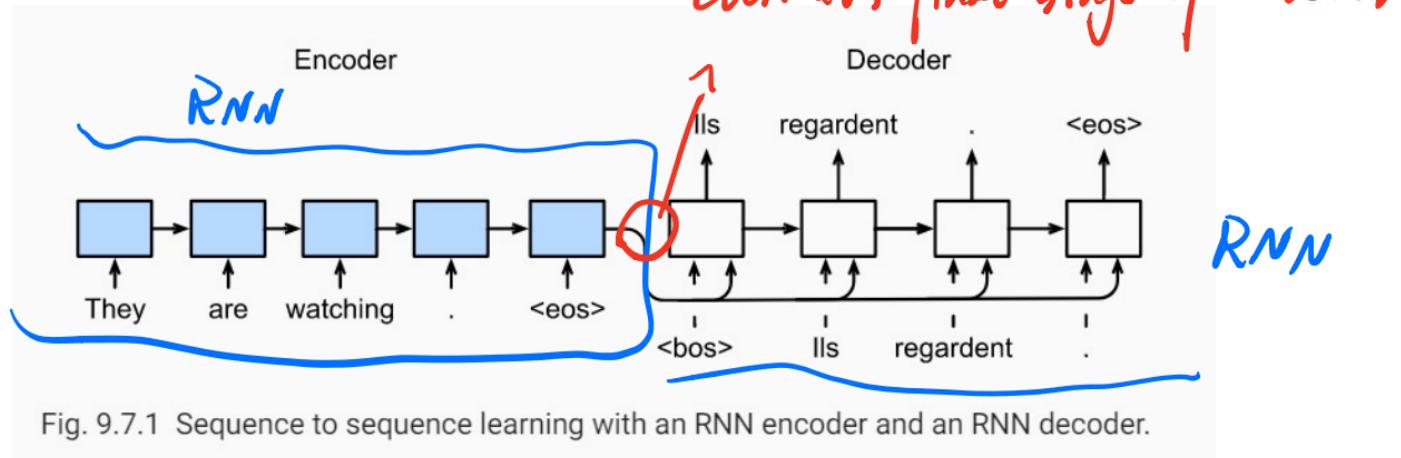


# Neural Machine Translation (NMT)

- Language translation



# Implementation



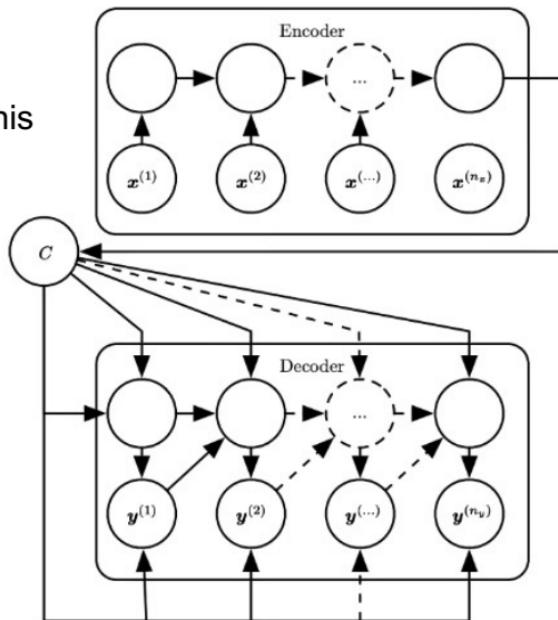
**Context:** final encoder state, or a combination of all intermediate states

**Decoder inputs:** target sequences and in some cases, context

**Decoder outputs:** shifted target sequence

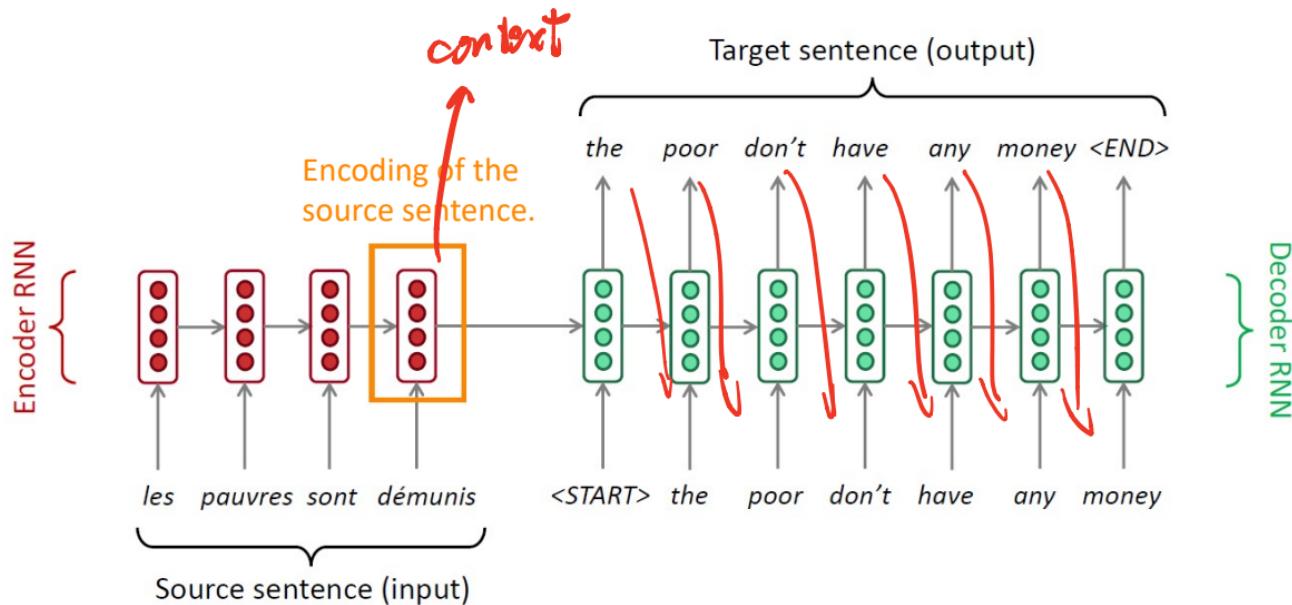
# Decoder-Encoder Model

Les pauvres sont démunis



The poor do not have any money

# Information Bottleneck



**Information Bottleneck:** All information about the source sentence captured as a single context

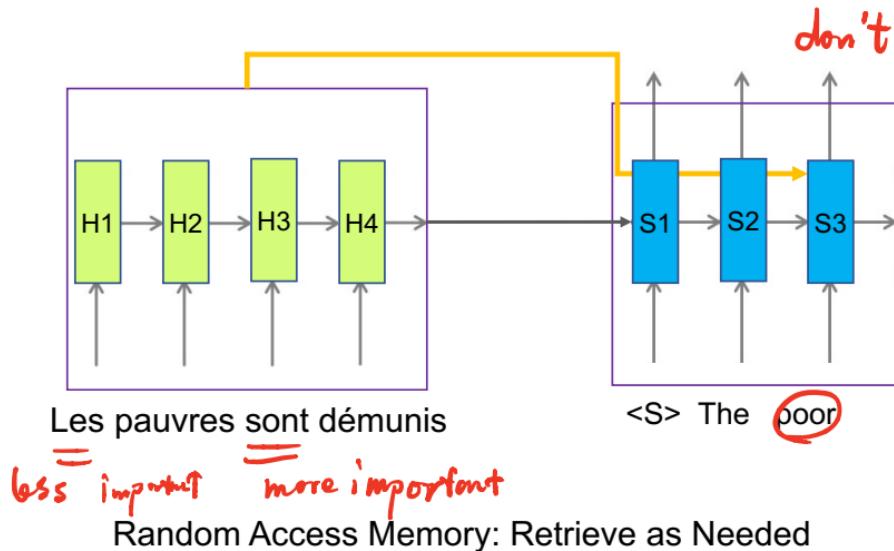
only pass h<sub>4</sub>      h<sub>4</sub> has less information about previous

# Classic Attention Mechanism

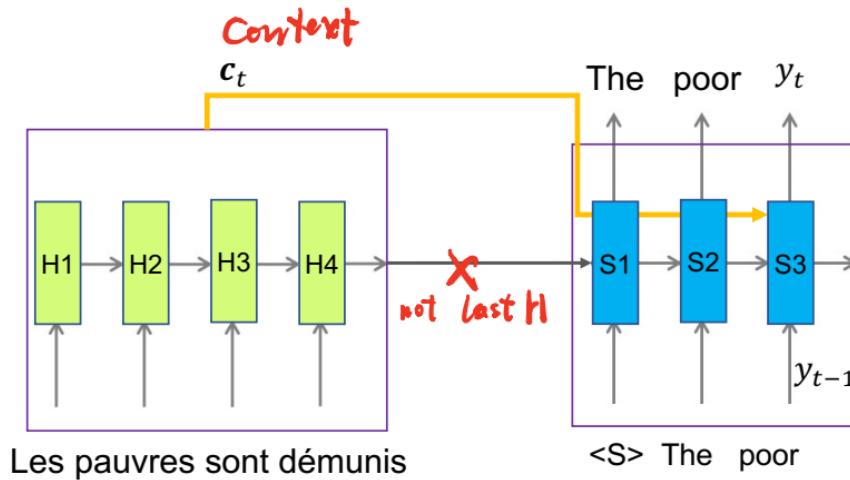
Solution to the bottleneck problem: *Attention*

**Key idea:** On each step of the decoder, *focus on a particular part* of the source sequence

# Classic Attention Mechanism

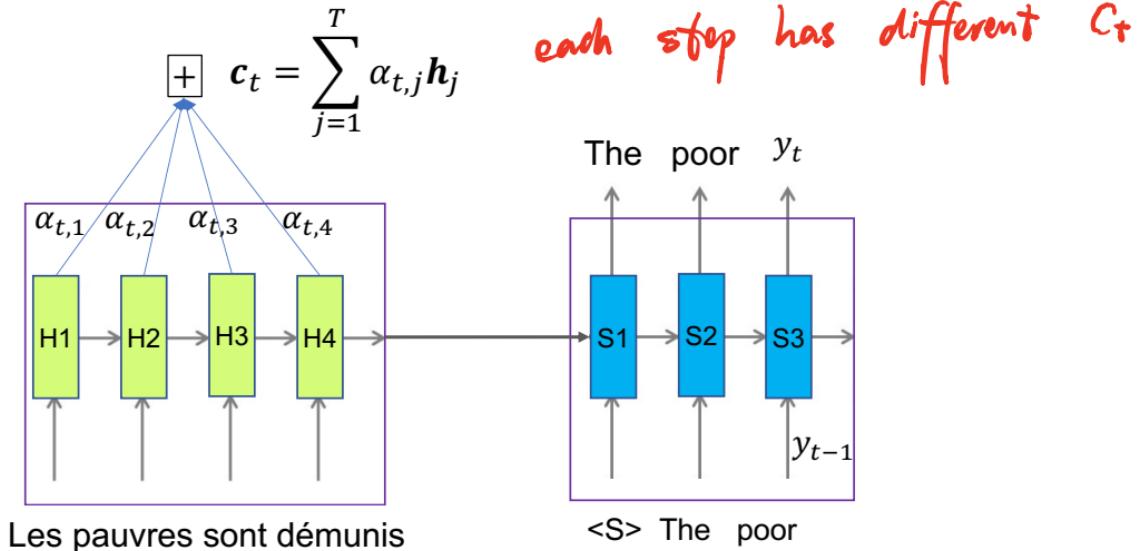


# Classic Attention Mechanism



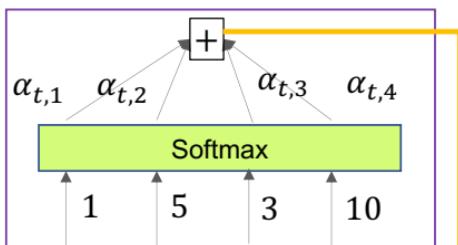
$$p(y_t | y_{t-1}, \dots, y_1, x) = f(\underbrace{y_{t-1}}_{\text{Hidden state}}, \underbrace{s_t}_{\text{Hidden state}}, \underbrace{c_t}_{\text{Hidden state}})$$

# Classic Attention Mechanism

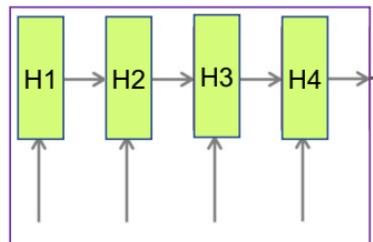


# Classic Attention Mechanism

$$\alpha_{t,j} = \frac{\exp(Score_{t,j})}{\sum_{k=1}^T \exp(Score_{t,k})} = \text{Softmax}(\text{scores})$$

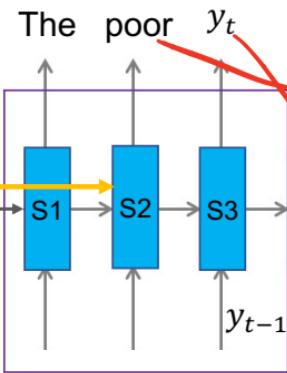


$c_t = [0 \ 0 \ 0 \ 1]$   
for all cases  
 $c_t$  is different  
for different steps



Les pauvres sont démunis

$$c_t = \sum_i \alpha_{t,i} \rightarrow \text{softmax}(\text{scores})$$



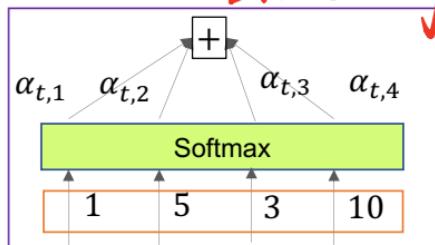
<S> The poor

$$c_t = [0.1, 0.2, 0.2] \rightarrow c_t = [0.5, 0.5, 0]$$

# Classic Attention Mechanism

$$\text{Score} = \sin(s_t, h_t) \xrightarrow{\text{softmax}} \begin{bmatrix} 0.5 \\ 0.1 \\ 0.3 \\ 0.1 \end{bmatrix} \Rightarrow C_t = 0.5H_1 + 0.1H_2 + 0.3H_3 + 0.1H_4$$

similarity function : dot product



$$Score_{t,j} = Score(s_t, h_j)$$

$y_{t-1}, s_t \rightarrow s_3$

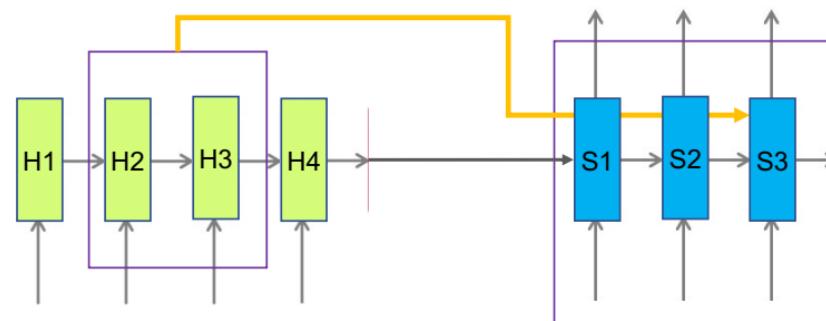
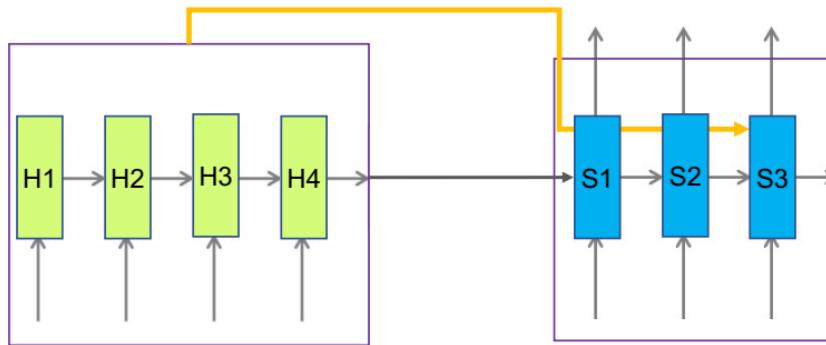
$$C_3 = \left\{ \begin{array}{l} \text{Sim}(s_3, H_1) \\ \text{Sim}(s_3, H_2) \\ \dots \\ \text{Sim}(H_4) \end{array} \right\}$$

Common choices

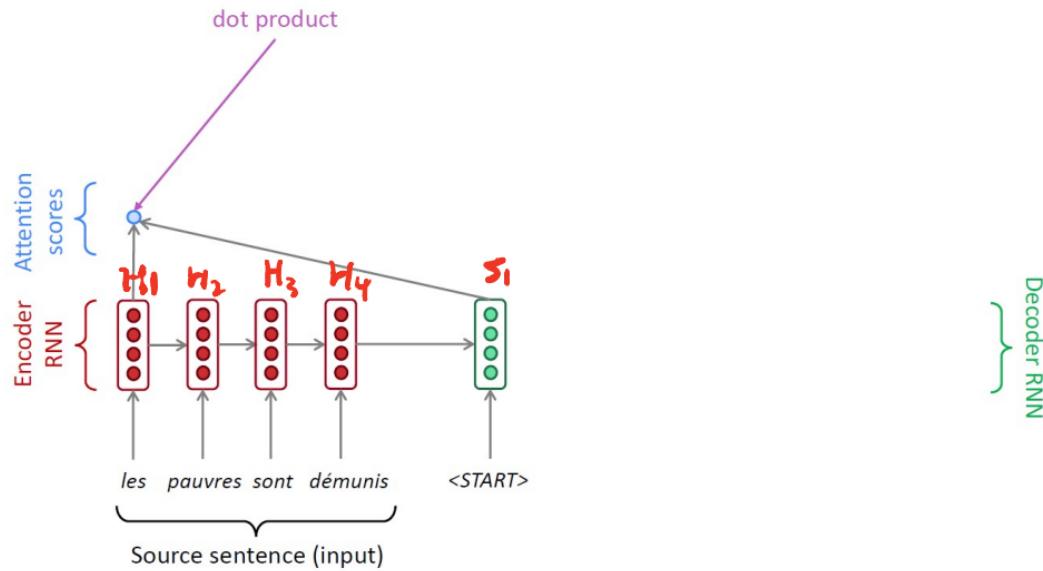
$$\begin{aligned} \text{Dot}(s_t, h_j) &= h_t^T s_t \\ &= h_t^T W_a s_t \end{aligned}$$

state of target language      state of original language

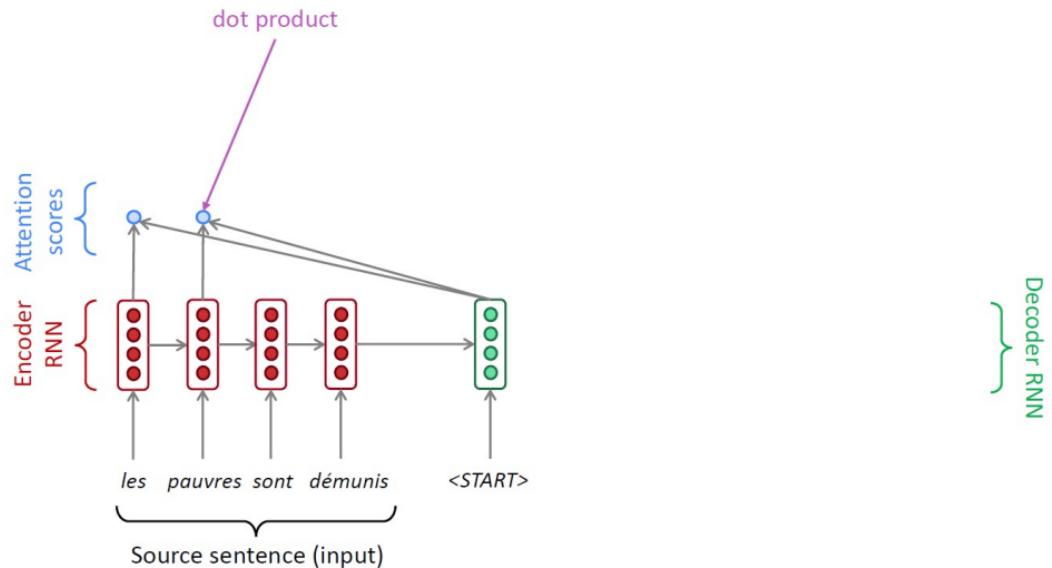
# Global vs Local Attention



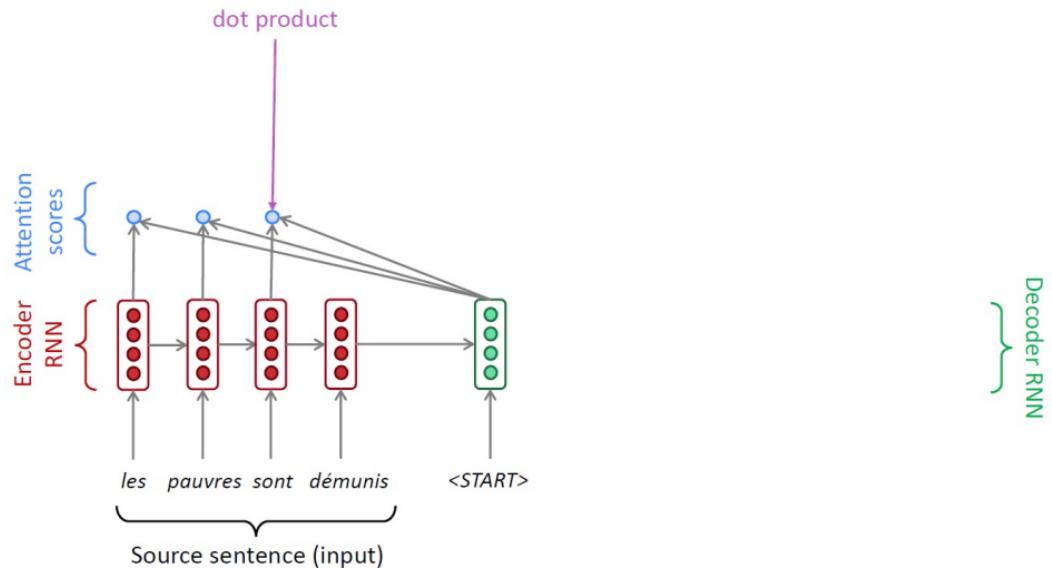
# Running Example



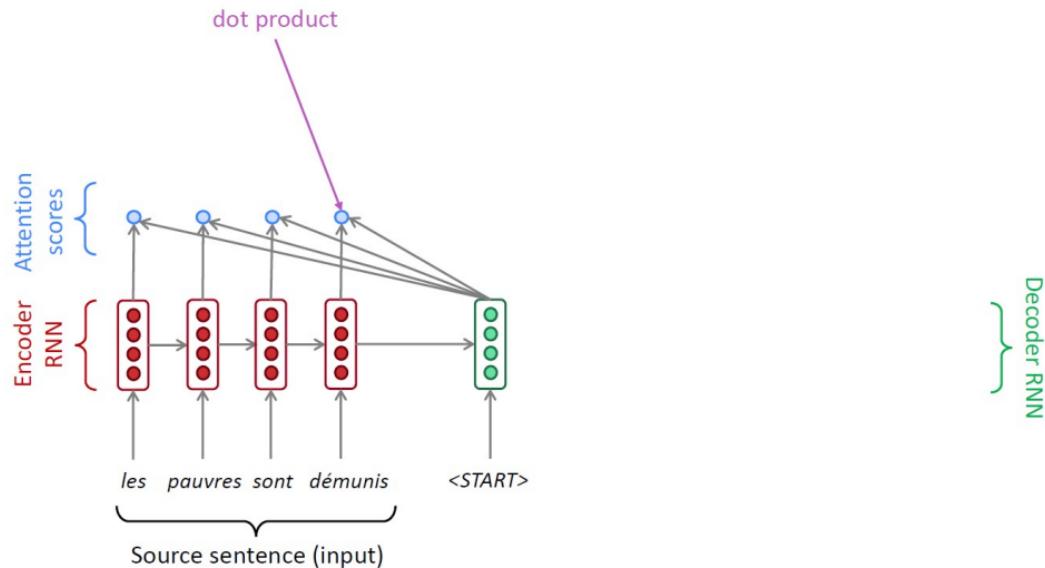
# Running Example



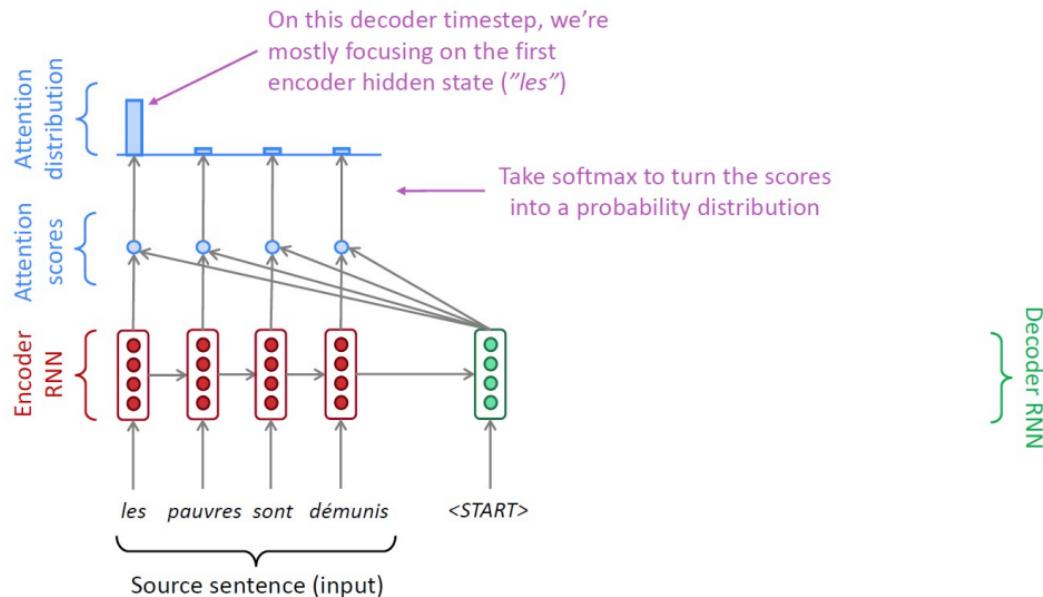
# Running Example



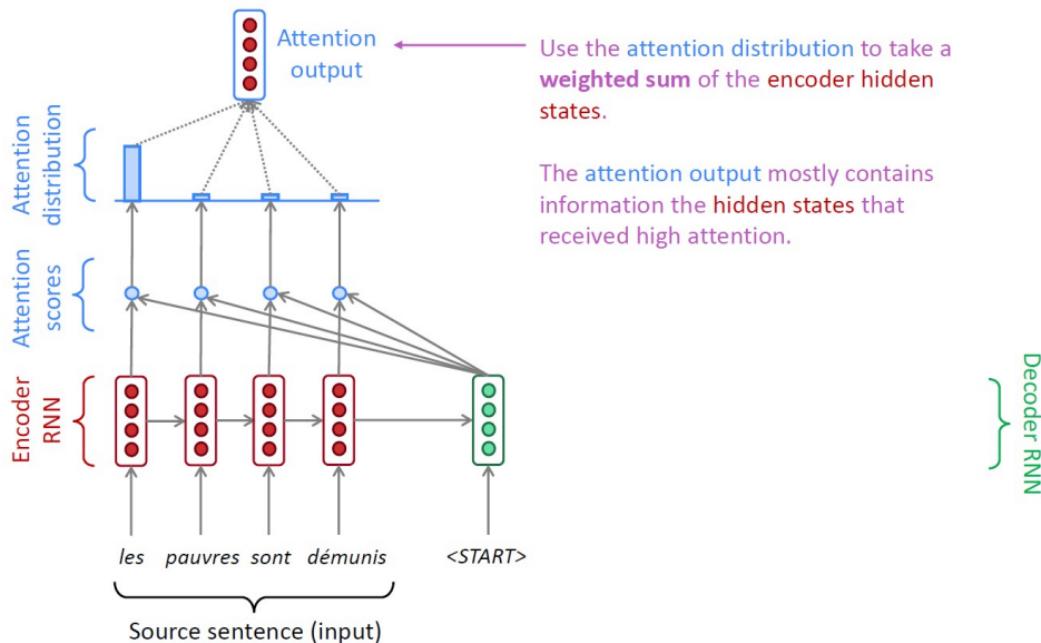
# Running Example



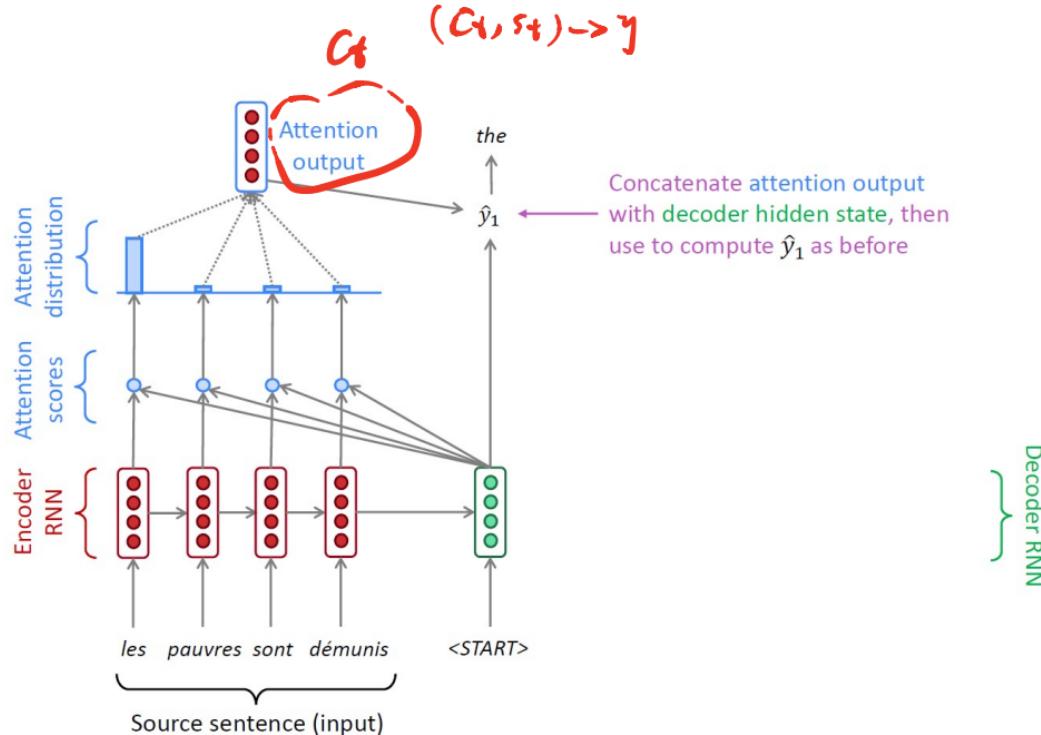
# Running Example



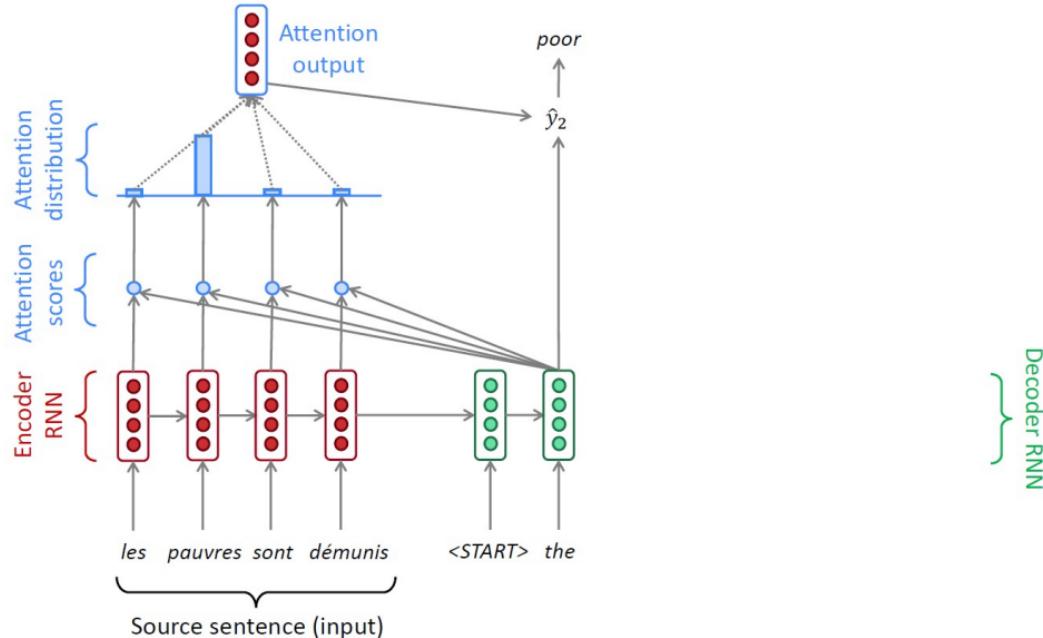
# Running Example



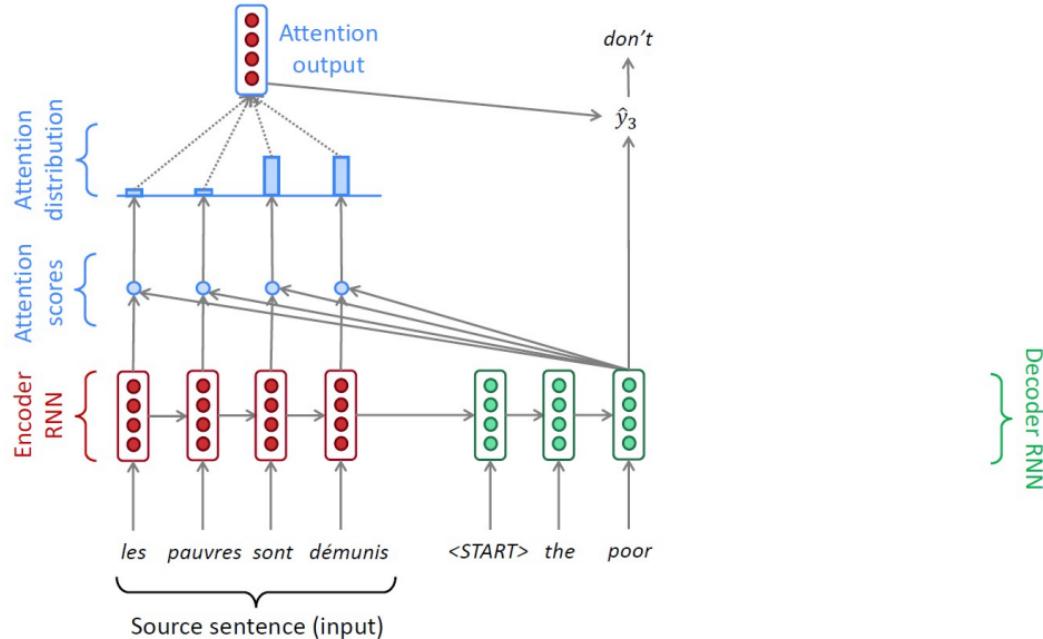
# Running Example



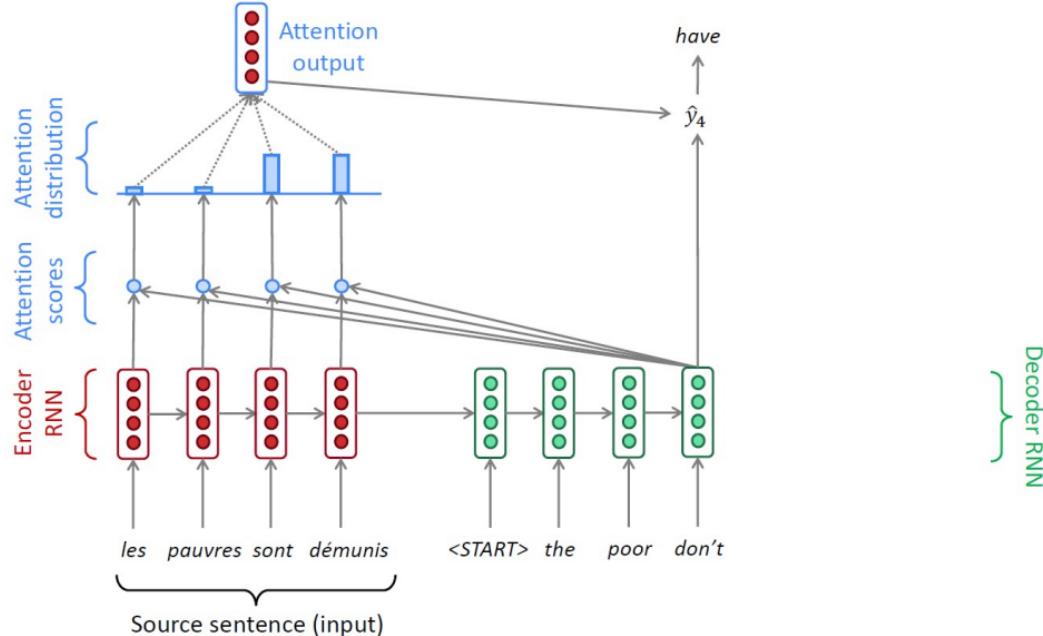
# Running Example



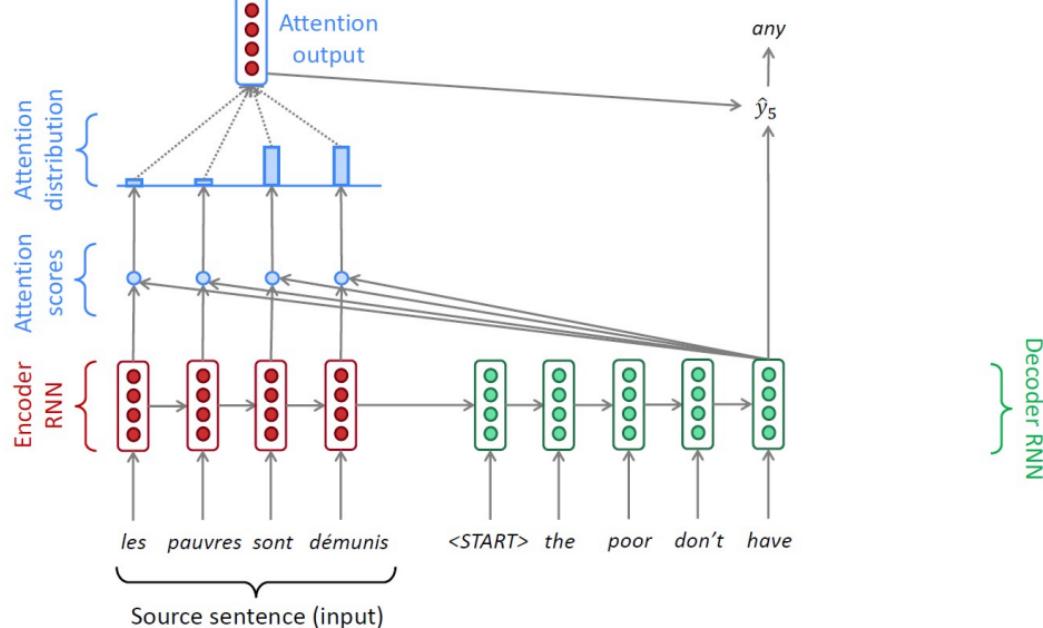
# Running Example



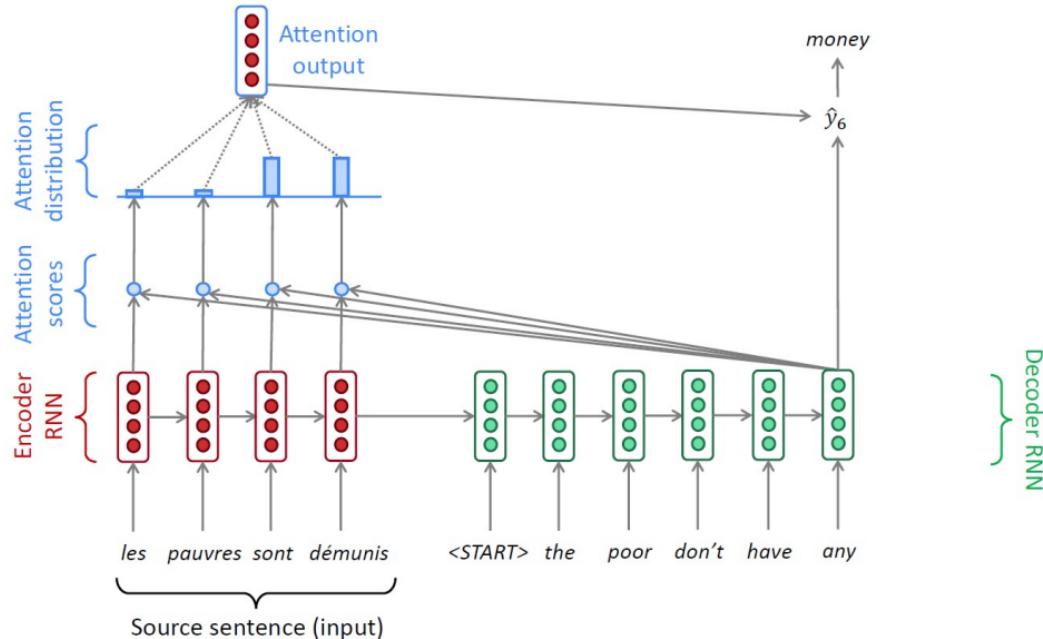
# Running Example



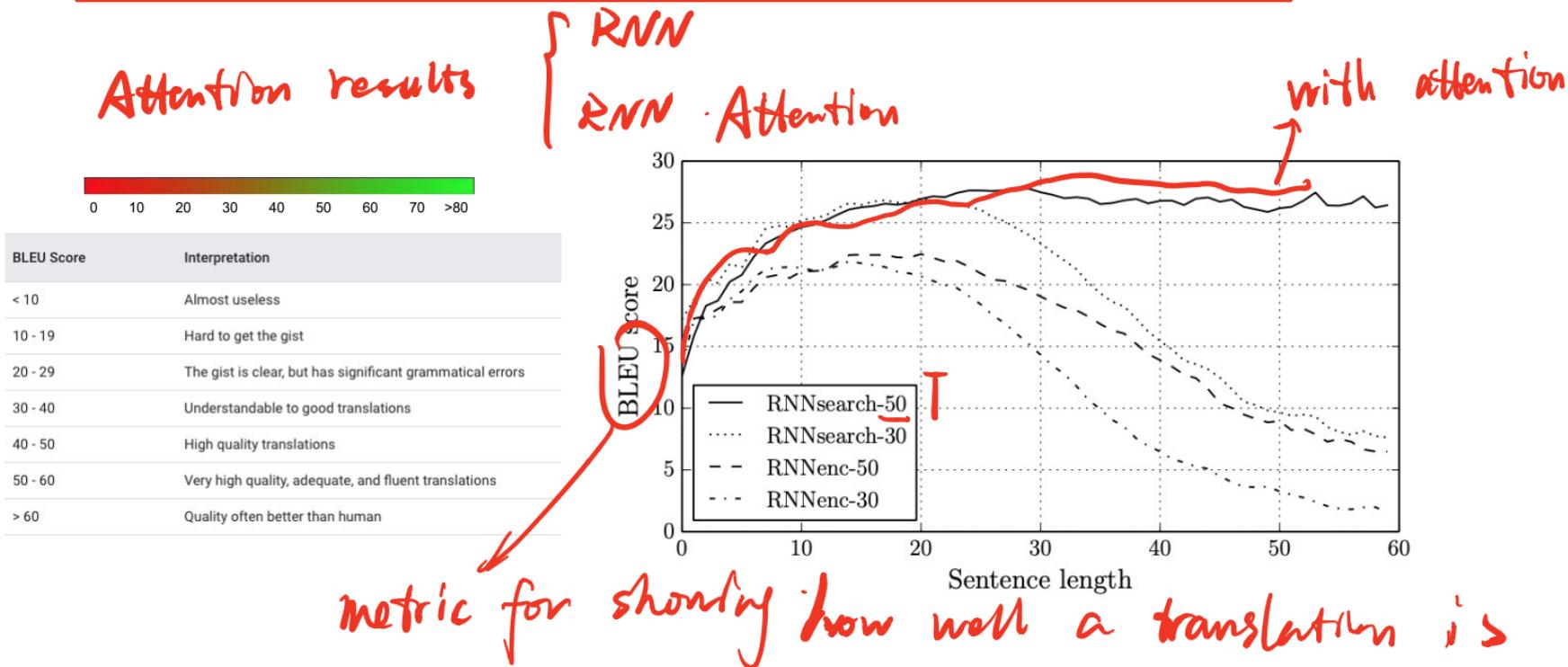
# Running Example



# Running Example



# Translation with Attention Mechanism



Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." *ICLR 2015*