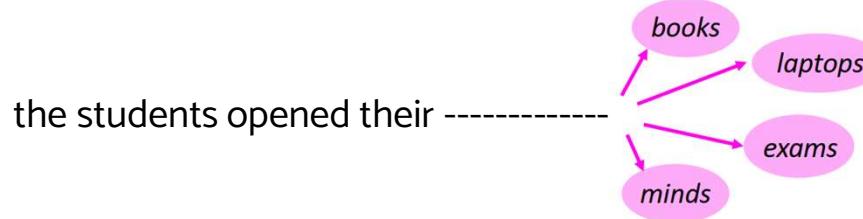


# RNNs and Language Models

Mohammad Taher Pilehvar

# Language Modeling

- Language Modeling is the task of predicting what word comes next



More formally: given a sequence of words  $x_1, x_2, \dots, x_t$ , compute the probability distribution of the next word  $x_{t+1}$ :

$$P(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})$$

where  $x_{t+1}$  can be any word in the vocabulary  $V = \{w_1, \dots, w_{|V|}\}$

# Language Modeling

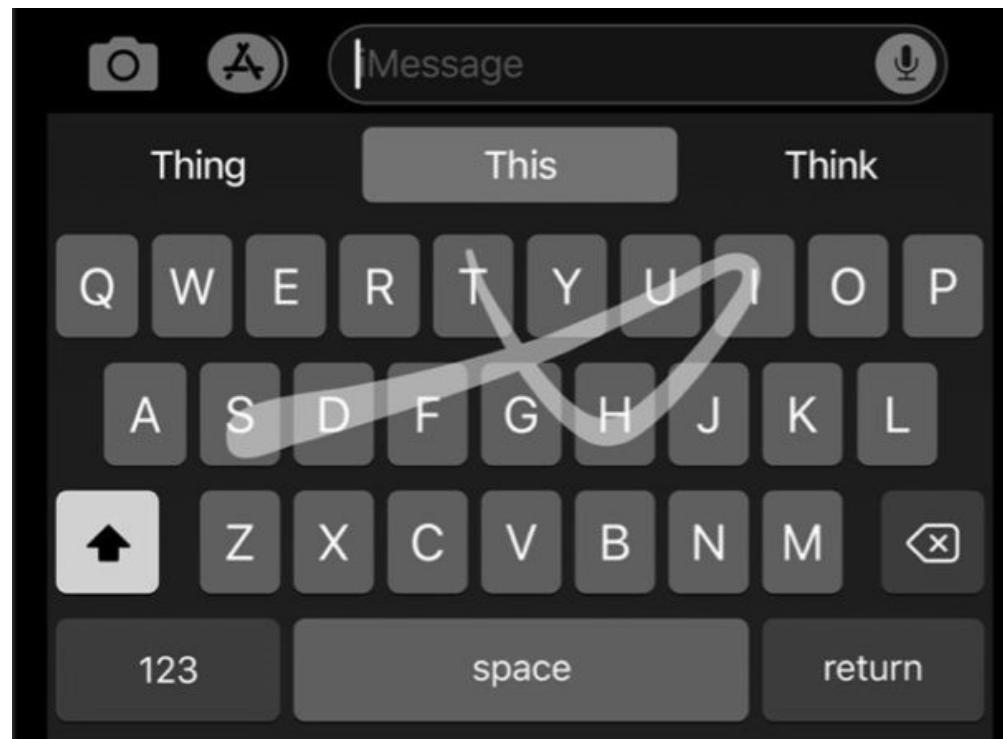
- You can also think of Language Model as a system that assigns a probability to a piece of text
- For example, if we have some text  $x_1, x_2, \dots, x_t$ , then the probability of this text (according to the Language Model) is:

$$P(x_1, \dots, x_T) = P(x_1) \times P(x_2|x_1) \times \cdots \times P(x_T|x_{T-1}, \dots, x_1)$$

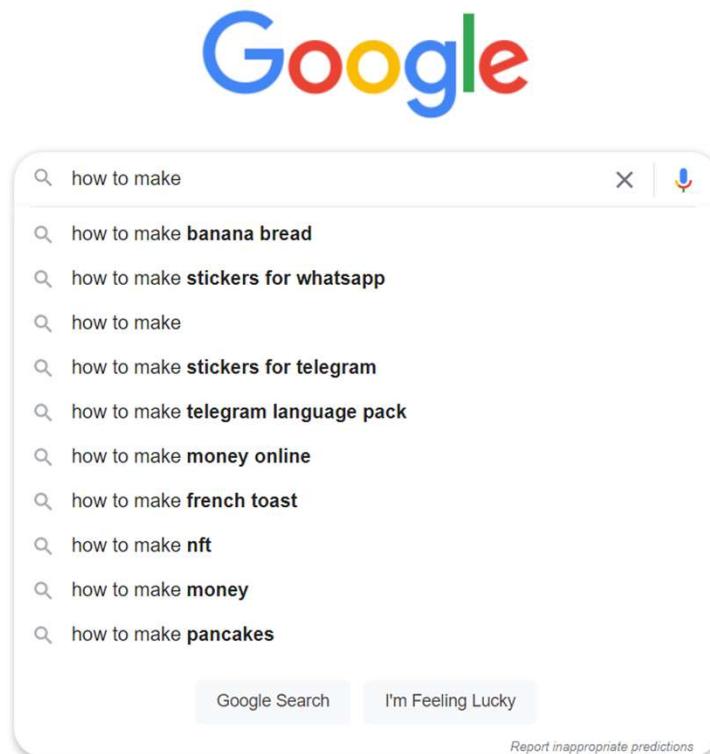
$$= \prod_{t=1}^T P(x_t | \underbrace{x_{t-1}, \dots, x_1}_{\text{This is what our LM provides}})$$

This is what our LM provides

# You use LMs on the daily basis!



# You use LMs on the daily basis!



# n-gram Language Model

- Question: How to learn a Language Model?
- Answer (pre- Deep Learning): learn an n-gram language model!
- Definition: An *n-gram* is a chunk of  $n$  consecutive words.
  - **uni**grams: “the”, “students”, “opened”, “their”
  - **bi**grams: “the students”, “students opened”, “opened their”
  - **tri**grams: “the students opened”, “students opened their”
  - **four**-grams: “the students opened their”
- Idea: Collect statistics about how frequent different n-grams are and use these to predict next word

# n-gram Language Model

- First we make a **Markov assumption**:  $x^{(t+1)}$  depends only on the preceding  $n-1$  words

$$P(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) = P(\mathbf{x}^{(t+1)} | \underbrace{\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)}}_{n-1 \text{ words}})$$

(assumption)

$$\begin{aligned} & \xrightarrow{\text{prob of a n-gram}} P(\mathbf{x}^{(t+1)}, \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)}) \\ & \xrightarrow{\text{prob of a (n-1)-gram}} P(\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)}) \end{aligned}$$

(definition of  
conditional prob)

- Question: How do we get these n-gram and (n-1)-gram probabilities?
- Answer: By counting them in some large corpus of text!

$$\approx \frac{\text{count}(\mathbf{x}^{(t+1)}, \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})}{\text{count}(\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})}$$

(statistical  
approximation)

# n-gram Language Model: example

- Suppose we are learning a 4-gram Language Model.

*condition on this*  
~~discard~~  
~~as the proctor started the clock, the students opened their \_\_\_\_\_~~

$$P(\mathbf{w}|\text{students opened their}) = \frac{\text{count(students opened their } \mathbf{w})}{\text{count(students opened their)}}$$

- For example, suppose that in the corpus:

- “students opened their” occurred 1000 times
- “students opened their books” occurred 400 times  
 $P(\text{books} | \text{students opened their}) = 0.4$
- “students opened their exams” occurred 100 times  
 $P(\text{exams} | \text{students opened their}) = 0.1$

Was it good that we discarded the context of “proctor”?

# Sparsity Problems with n-gram Language Models

## Sparsity Problem 1

**Problem:** What if “students opened their  $w$ ” never occurred in data? Then  $w$  has probability 0!

**(Partial) Solution:** Add small  $\delta$  to the count for every  $w \in V$ . This is called *smoothing*.

$$P(w|\text{students opened their}) = \frac{\text{count(students opened their } w\text{)}}{\text{count(students opened their)}}$$

## Sparsity Problem 2

**Problem:** What if “students opened their” never occurred in data? Then we can’t calculate probability for *any*  $w$ !

**(Partial) Solution:** Just condition on “opened their” instead. This is called *backoff*.

**Note:** Increasing  $n$  makes sparsity problems *worse*. Typically, we can’t have  $n$  bigger than 5.

# Sparsity Problems with n-gram Language Models

**Storage:** Need to store count for all  $n$ -grams you saw in the corpus.

$$P(\mathbf{w}|\text{students opened their}) = \frac{\text{count(students opened their } \mathbf{w})}{\text{count(students opened their)}}$$

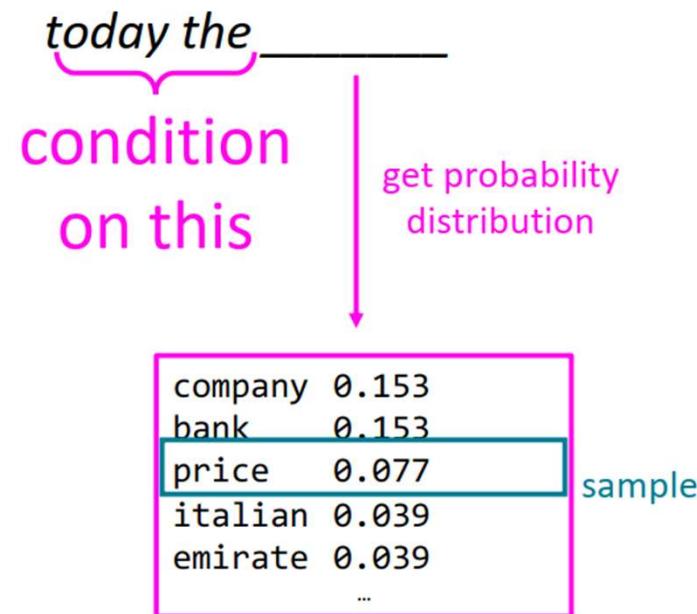
Increasing  $n$  or increasing corpus increases model size!

# n-gram LMs in practice

- You can build a simple trigram Language Model over a 1.7 million word corpus (Reuters) in a few seconds on your laptop\*

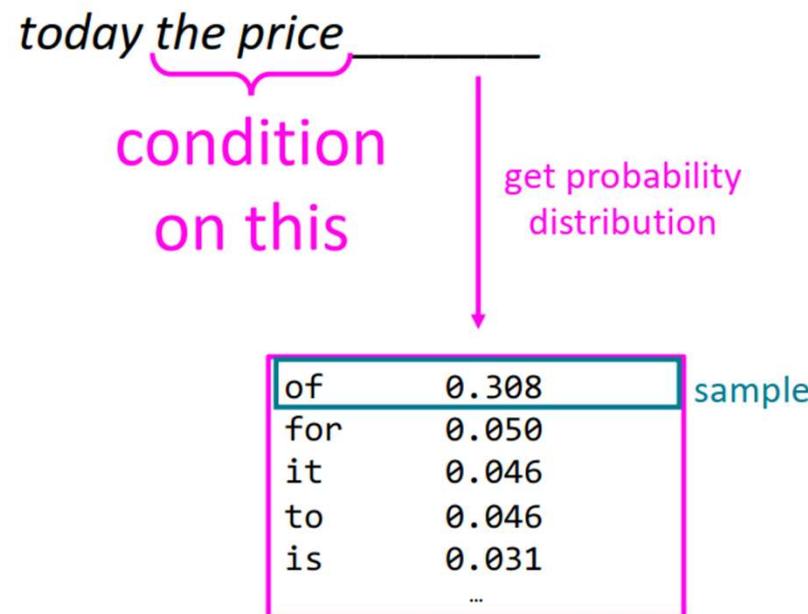
# Generating text with a n-gram Language Model

You can also use a Language Model to generate text



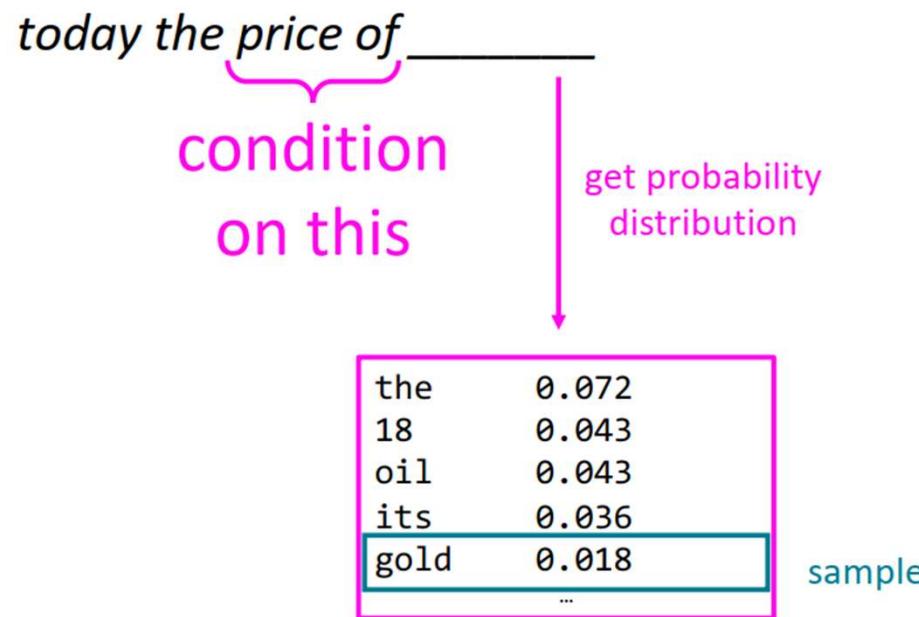
# Generating text with a n-gram Language Model

You can also use a Language Model to generate text



# Generating text with a n-gram Language Model

You can also use a Language Model to generate text



# Generating text with a n-gram Language Model

You can also use a Language Model to generate text

*today the price of gold per ton , while production of shoe  
lasts and shoe industry , the bank intervened just after it  
considered and rejected an imf demand to rebuild depleted  
european stocks , sept 30 end primary 76 cts a share .*

Surprisingly grammatical!

...but **incoherent**. We need to consider more than three words at a time if we want to model language well.

But increasing  $n$  worsens sparsity problem,  
and increases model size...

# Neural Networks

# How to build a neural Language Model?

- Recall the Language Modeling task:
  - Input: sequence of words  $x_1, x_2, \dots, x_t$
  - Output: prob distribution of the next word  $P(x_{t+1}|x_t, \dots, x_1)$
- How about a window-based neural model?

# A fixed-window neural Language Model

# A fixed-window neural Language Model

output distribution

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

hidden layer

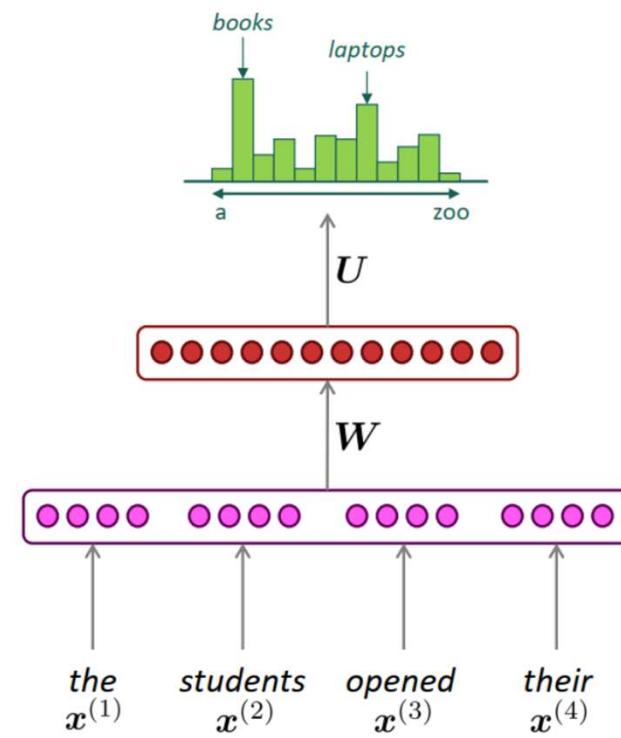
$$\mathbf{h} = f(\mathbf{W}\mathbf{e} + \mathbf{b}_1)$$

concatenated word embeddings

$$\mathbf{e} = [\mathbf{e}^{(1)}; \mathbf{e}^{(2)}; \mathbf{e}^{(3)}; \mathbf{e}^{(4)}]$$

words / one-hot vectors

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \mathbf{x}^{(4)}$$



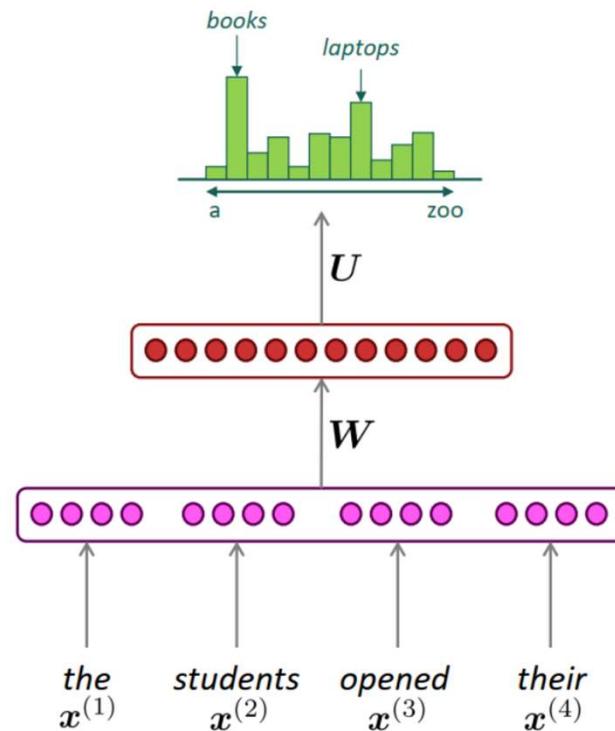
# A fixed-window neural Language Model

**Improvements** over n-gram LM:

- No sparsity problem
- Don't need to store all observed  $n$ -grams

Remaining **problems**:

- Fixed window is **too small**
- Enlarging window enlarges  $W$
- Window can never be large enough!
- $x(0)$  and  $x(1)$  are multiplied by completely different weights in  $W$ . **No symmetry** in how the inputs are processed.

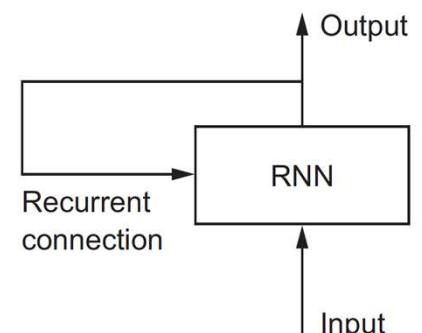


# Understanding RNNs

- In Dense and ConvNet, each input is processed independently, with no state kept in between inputs:
  - To process a sequence or a temporal series of data points, you have to show the entire sequence to the network at once: turn it into a single data point.
- In contrast, human processes a text word by word, while keeping memories of what came before.
  - Biological intelligence processes information incrementally while maintaining an internal model of what it's processing, built from past information and constantly updated as new information comes in.

# Understanding RNNs

- RNN adopts the same principle, albeit in an extremely simplified version:
  - It processes sequences by iterating through the sequence elements and maintaining a state containing information relative to what it has seen so far.
- The state of the RNN is reset between processing two different, independent sequences.
  - We can still consider one sequence a single data point: a single input to the network.



**Figure 6.9** A recurrent network: a network with a loop

# Understanding RNNs

To make it clear, let's implement the forward pass of a toy RNN in Numpy.

```
state_t = 0           ← The state at t
for input_t in input_sequence:   ← Iterates over sequence elements
    output_t = f(input_t, state_t)
    state_t = output_t   ← The previous output becomes the state for the next iteration.
```

# Understanding RNNs

- You can even flesh out the function  $f$ : the transformation of the input and state into an output will be parameterized by two matrices,  $W$  and  $U$ , and a bias vector.

```
state_t = 0
for input_t in input_sequence:
    output_t = activation(dot(W, input_t) + dot(U, state_t) + b)
    state_t = output_t
```

# Understanding RNNs

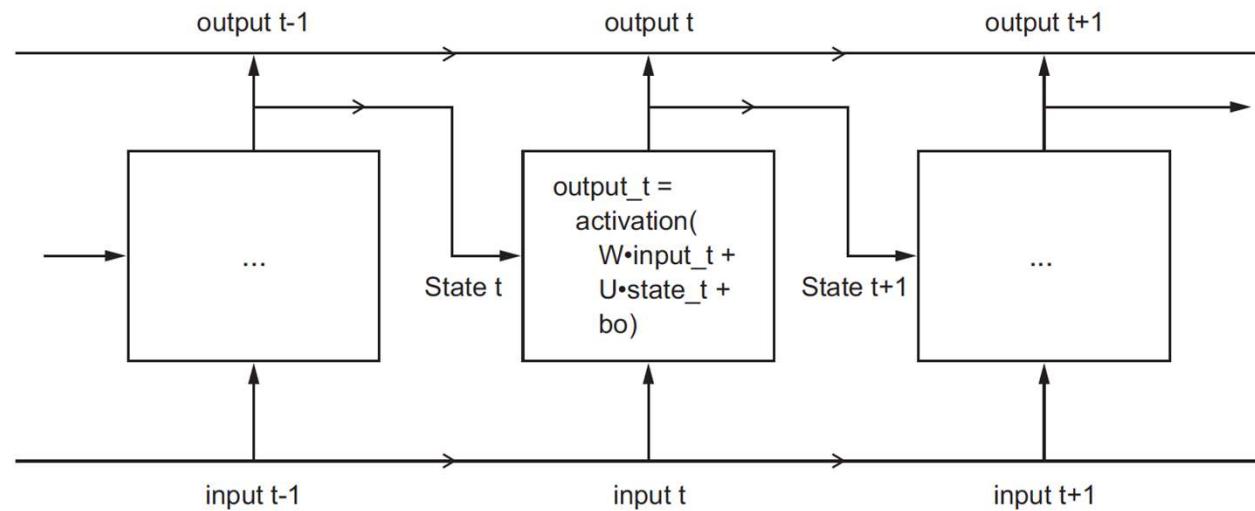
```
Number of timesteps in  
the input sequence  
Dimensionality of the  
input feature space  
Dimensionality of the  
output feature space  
Input data: random  
noise for the sake of  
the example  
Initial state: an  
all-zero vector  
Creates random  
weight matrices  
input_t is a vector of  
shape (input_features,).  
The final output is a 2D tensor of  
shape (timesteps, output_features).  
Updates the state of the  
network for the next timestep  
Stores this output in a list  
Combines the input with the current  
state (the previous output) to obtain  
the current output
```

```
import numpy as np  
timesteps = 100  
input_features = 32  
output_features = 64  
  
inputs = np.random.random((timesteps, input_features))  
state_t = np.zeros((output_features, ))  
  
W = np.random.random((output_features, input_features))  
U = np.random.random((output_features, output_features))  
b = np.random.random((output_features, ))  
  
successive_outputs = []  
for input_t in inputs:  
    output_t = np.tanh(np.dot(W, input_t) + np.dot(U, state_t) + b)  
    successive_outputs.append(output_t)  
  
    state_t = output_t  
  
final_output_sequence = np.concatenate(successive_outputs, axis=0)
```

# Understanding RNNs

- RNNs are characterized by their step function, such as the following function in this case:

```
output_t = np.tanh(np.dot(W, input_t) + np.dot(U, state_t) + b)
```

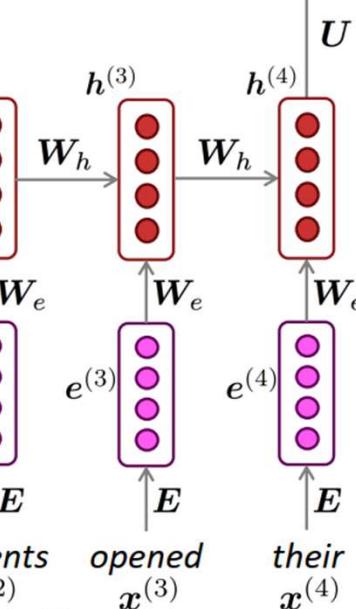
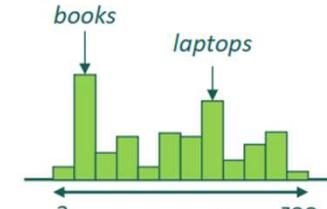


# Understanding RNNs

output distribution

$$\hat{y}^{(t)} = \text{softmax}(\mathbf{U}h^{(t)} + b_2) \in \mathbb{R}^{|V|}$$

$$\hat{y}^{(4)} = P(x^{(5)} | \text{the students opened their})$$



hidden states

$$h^{(t)} = \sigma(\mathbf{W}_h h^{(t-1)} + \mathbf{W}_e e^{(t)} + b_1)$$

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

word embeddings

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

words / one-hot vectors

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

*Note: this input sequence could be much longer now!*

# RNN Language Models

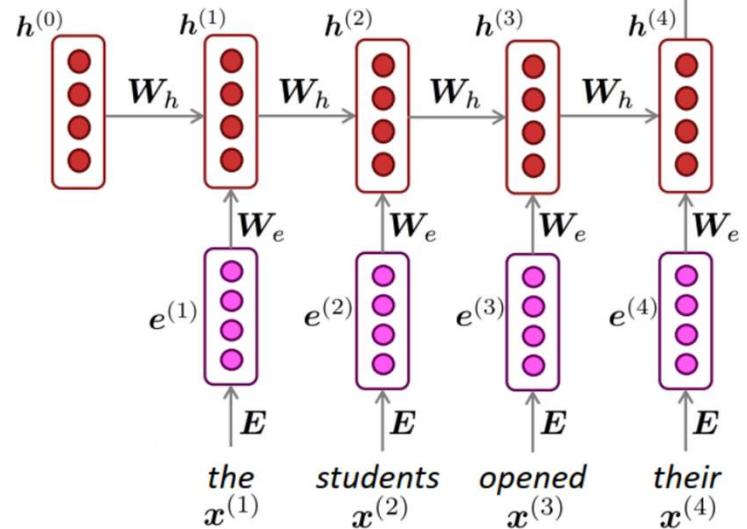
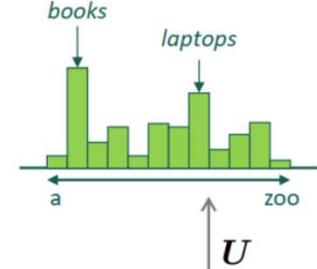
## RNN Advantages:

- Can process **any length** input
- Computation for step  $t$  can (in theory) use information from **many steps back**
- **Model size doesn't increase** for longer input context
- Same weights applied on every timestep, so there is **symmetry** in how inputs are processed.

## RNN Disadvantages:

- Recurrent computation is **slow**
- In practice, difficult to access information from **many steps back**

$$\hat{y}^{(4)} = P(\mathbf{x}^{(5)} | \text{the students opened their})$$



# Training an RNN LM

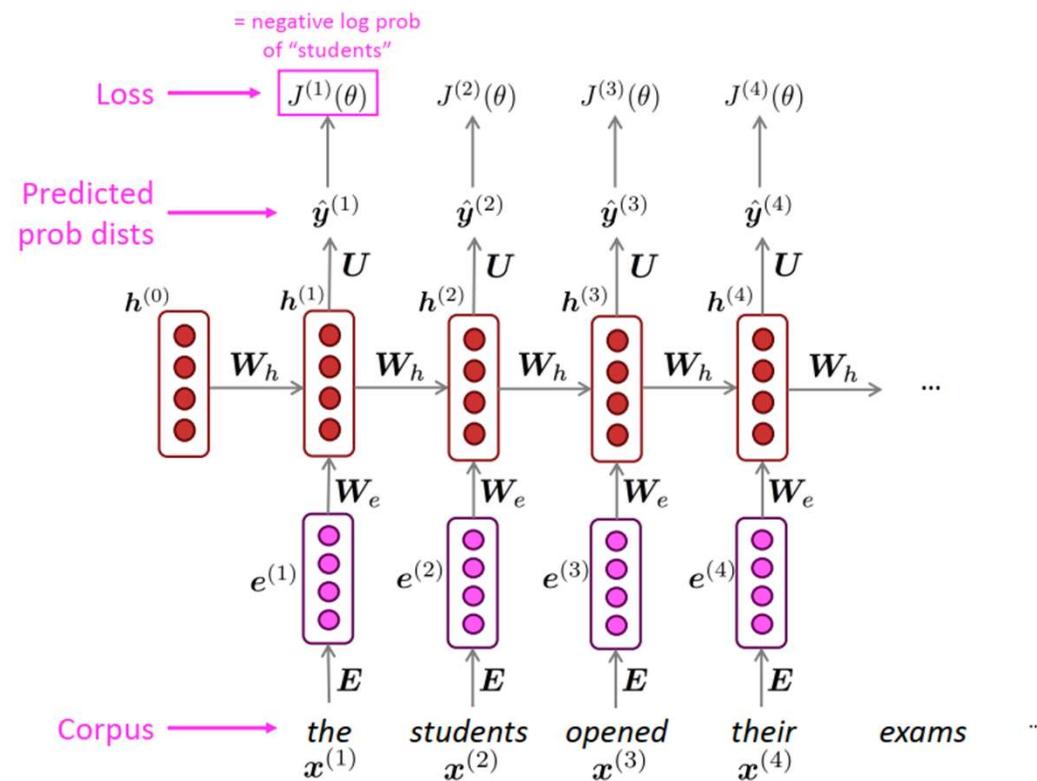
- Get a big corpus of text which is a sequence of words  $x^1, \dots, x^T$
- Feed into RNN-LM; compute output distribution  $\hat{y}^{(t)}$  for every step  $t$ .
  - i.e. predict probability dist of every word, given words so far
- Loss function on step  $t$  is **cross-entropy** between predicted probability distribution  $\hat{y}^{(t)}$ , and the true next word  $y^{(t)}$  (one-hot for  $x^{(t+1)}$ ):

$$J^{(t)}(\theta) = CE(y^{(t)}, \hat{y}^{(t)}) = - \sum_{w \in V} y_w^{(t)} \log \hat{y}_w^{(t)} = - \log \hat{y}_{x_{t+1}}^{(t)}$$

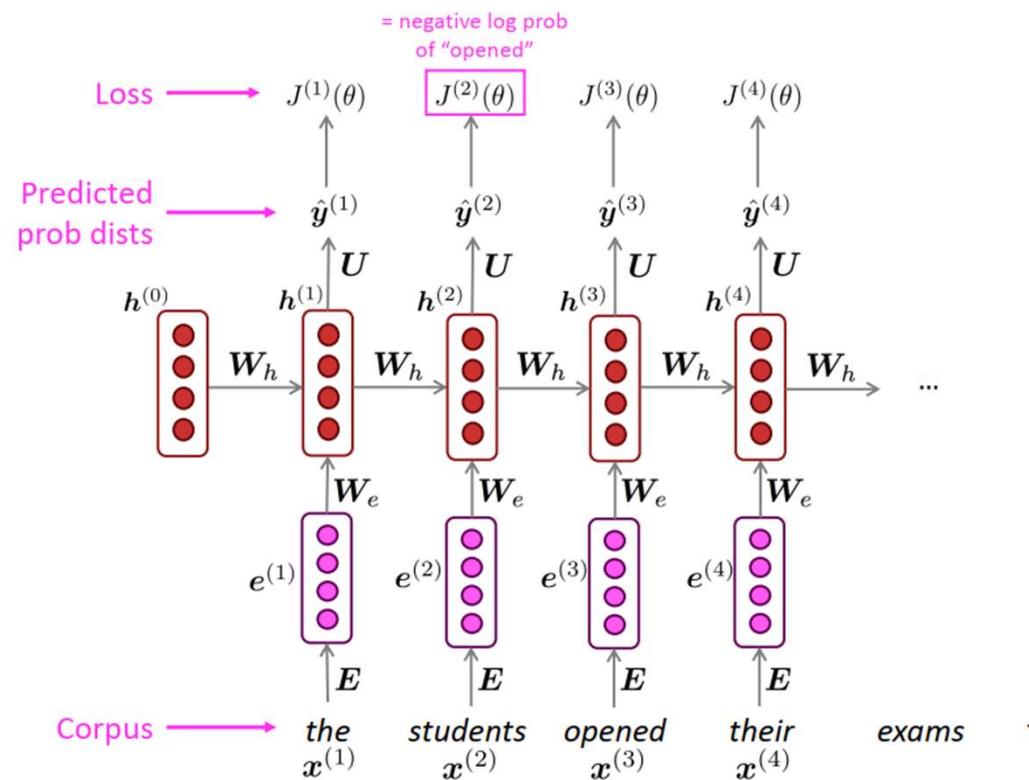
- Average this to get overall loss for entire training set:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^T - \log \hat{y}_{x_{t+1}}^{(t)}$$

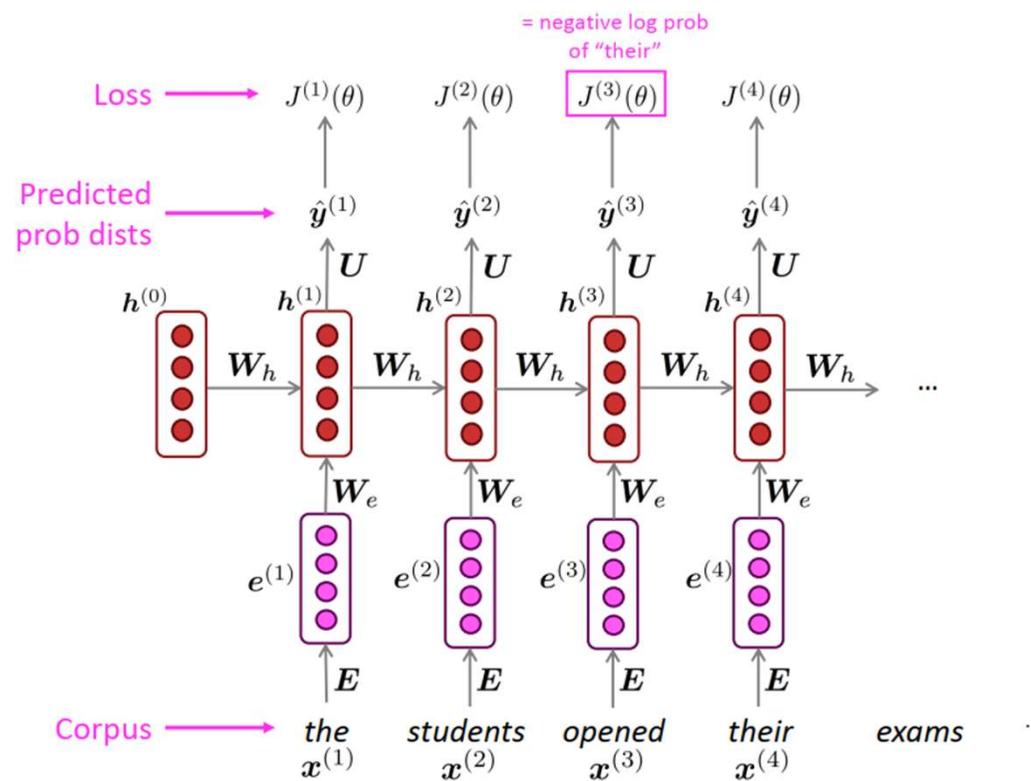
# Training an RNN LM



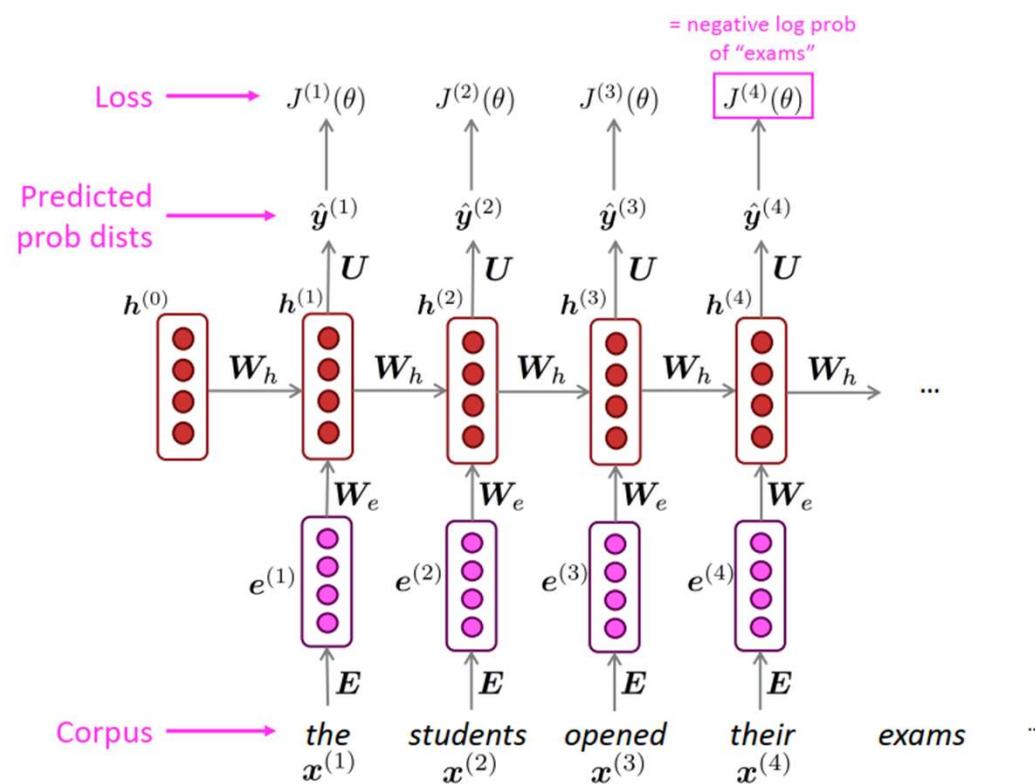
# Training an RNN LM



# Training an RNN LM

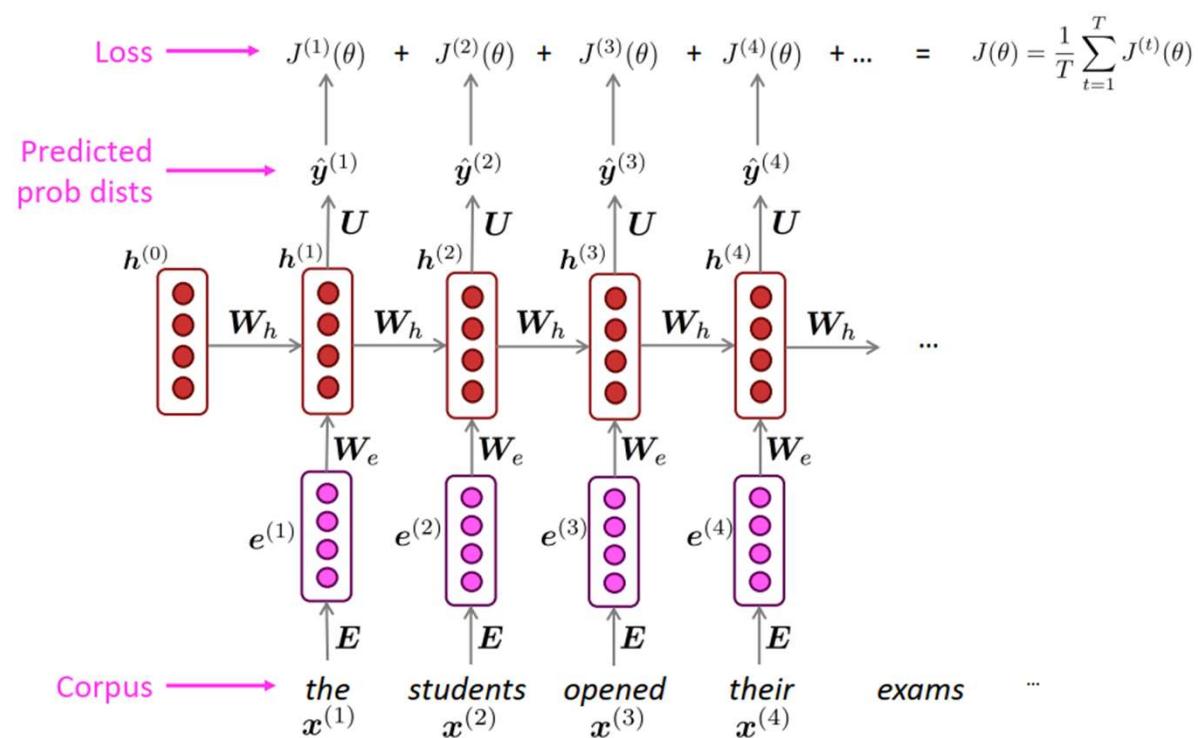


# Training an RNN LM

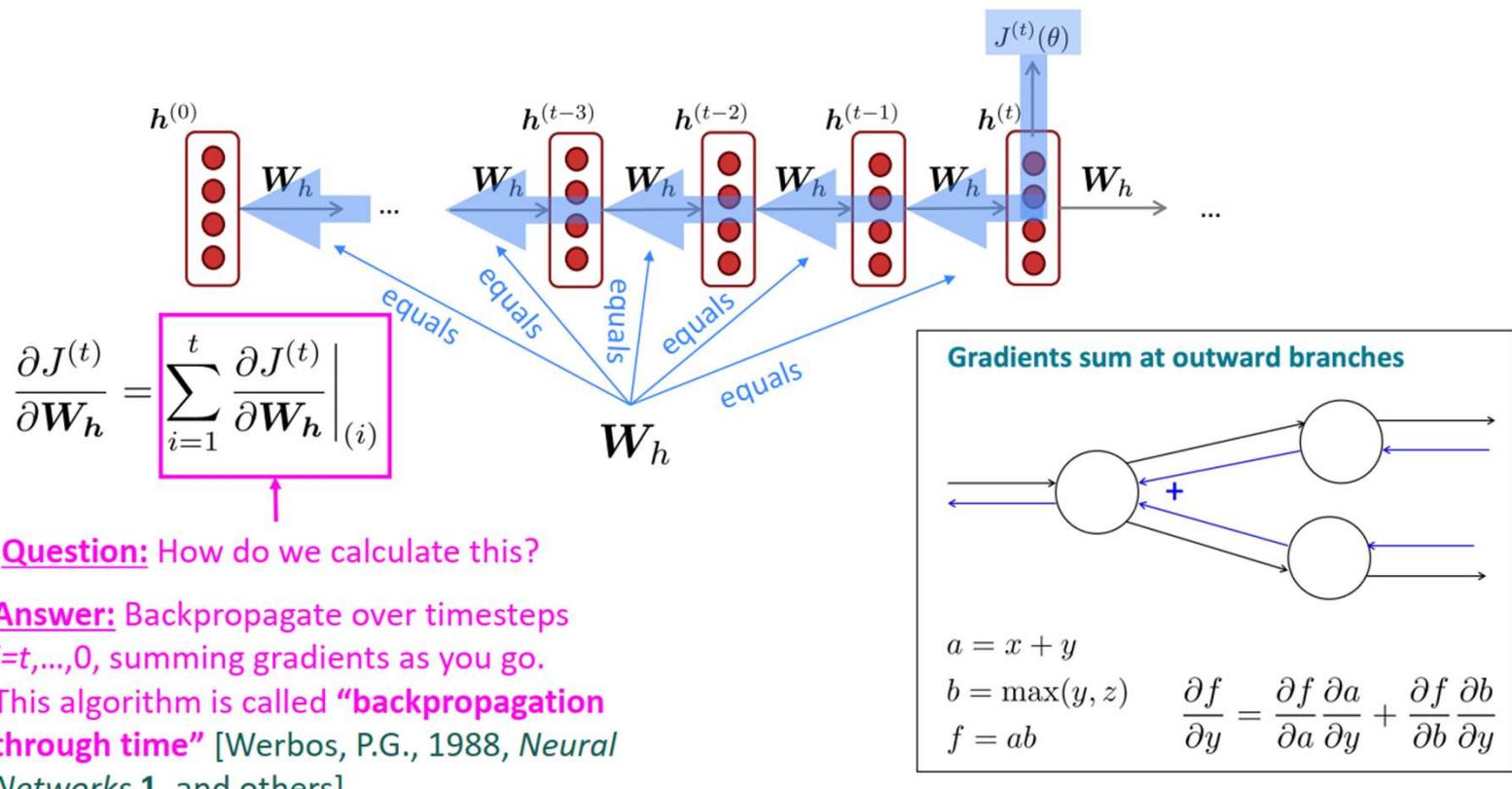


# Training an RNN LM

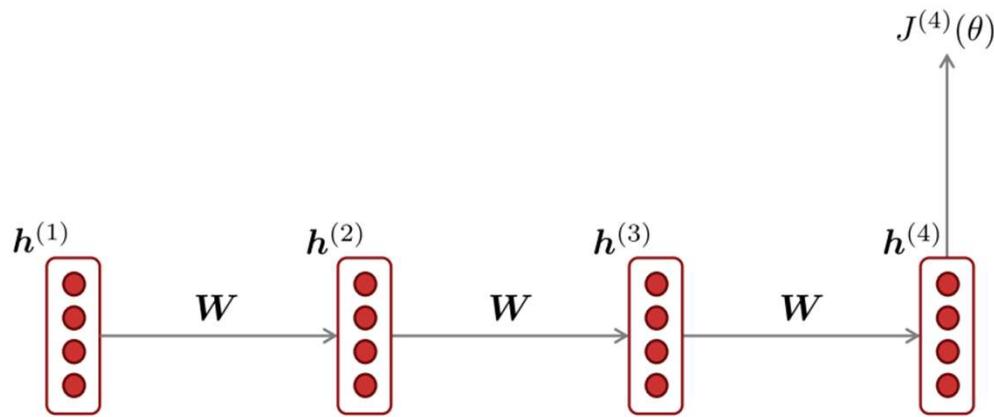
Teacher Forcing



# Backpropagation for RNNs

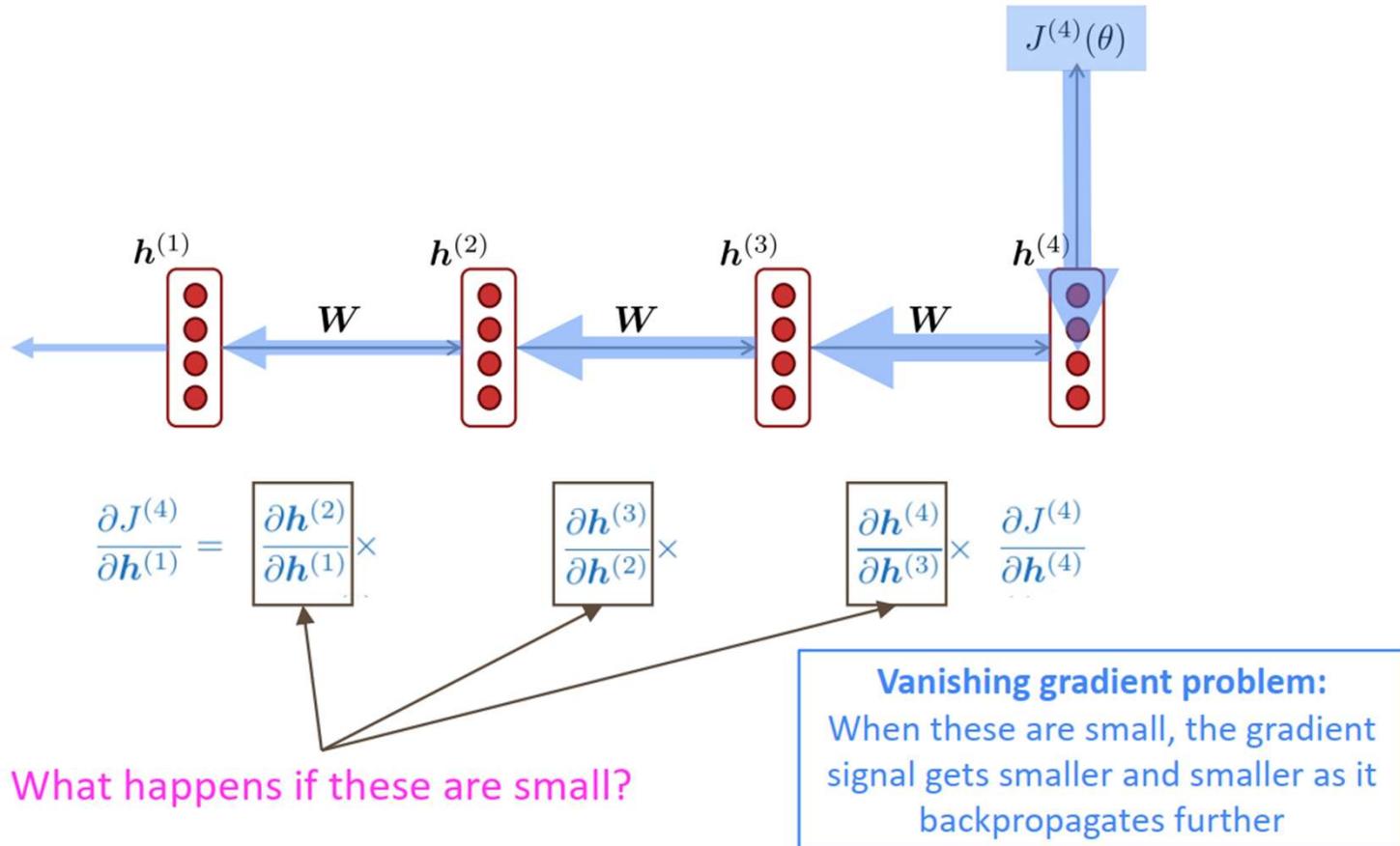


# Backpropagation for RNNs



$$\frac{\partial J^{(4)}}{\partial \mathbf{h}^{(1)}} = \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}} \times \quad \frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(2)}} \times \quad \frac{\partial \mathbf{h}^{(4)}}{\partial \mathbf{h}^{(3)}} \times \frac{\partial J^{(4)}}{\partial \mathbf{h}^{(4)}}$$

# Backpropagation for RNNs



## Vanishing gradient proof sketch (linear case)

- Recall:

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

- What if  $\sigma$  were the identity function,  $\sigma(x) = x$  ?

$$\begin{aligned} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} &= \text{diag}\left(\sigma'(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1)\right) \mathbf{W}_h && \text{(chain rule)} \\ &= \mathbf{I} \mathbf{W}_h = \mathbf{W}_h \end{aligned}$$

- Consider the gradient of the loss  $J^{(i)}(\theta)$  on step  $i$ , with respect to the hidden state  $\mathbf{h}^{(j)}$  on some previous step  $j$ . Let  $\ell = i - j$

$$\begin{aligned} \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(j)}} &= \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \prod_{j < t \leq i} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} && \text{(chain rule)} \\ &= \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \prod_{j < t \leq i} \mathbf{W}_h = \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \boxed{\mathbf{W}_h^{\ell}} && \text{(value of } \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \text{ )} \end{aligned}$$

If  $\mathbf{W}_h$  is “small”, then this term gets exponentially problematic as  $\ell$  becomes large

**Source:** “On the difficulty of training recurrent neural networks”, Pascanu et al, 2013. <http://proceedings.mlr.press/v28/pascanu13.pdf>  
 (and supplemental materials), at <http://proceedings.mlr.press/v28/pascanu13-supp.pdf>

## Vanishing gradient proof sketch (linear case)

- What's wrong with  $\mathbf{W}_h^\ell$  ?
- Consider if the eigenvalues of  $\mathbf{W}_h$  are all less than 1:
  - $\lambda_1, \lambda_2, \dots, \lambda_n < 1$
  - $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n$  (eigenvectors)
- We can write  $\frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \mathbf{W}_h^\ell$  using the eigenvectors of  $\mathbf{W}_h$  as a basis:

$$\frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \mathbf{W}_h^\ell = \sum_{i=1}^n c_i \boxed{\lambda_i^\ell} \mathbf{q}_i \approx \mathbf{0} \text{ (for large } \ell\text{)}$$

↑  
Approaches 0 as  $\ell$  grows, so gradient vanishes

- What about nonlinear activations  $\sigma$  (i.e., what we use?)
  - Pretty much the same thing, except the proof requires  $\lambda_i < \gamma$  for some  $\gamma$  dependent on dimensionality and  $\sigma$

# Understanding LSTM and GRU layers

SimpleRNN is generally too simplistic to be of real use.

- Although it should theoretically be able to retain at time  $t$  information about inputs seen many timesteps before, in practice, such long-term dependencies are impossible to learn.

*Vanishing gradient problem*

# Vanishing gradient problem

An effect that is similar to what is observed with non-recurrent networks (feedforward networks) that are many layers deep: as you keep adding layers to a network, the network eventually becomes untrainable.

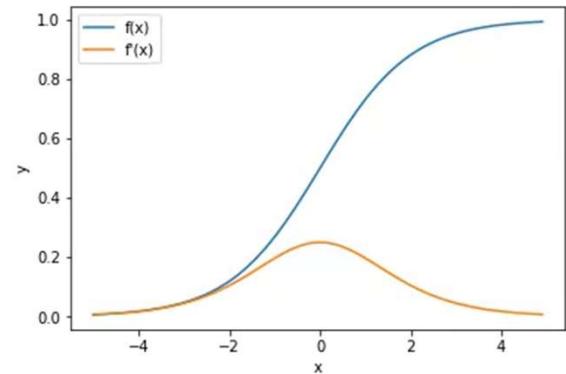
# Vanishing gradients

Arises due to the nature of the back-propagation optimization

- The weight and bias values in the various layers within a neural network are updated each optimization iteration by stepping in the direction of the *gradient* of the weight/bias values with respect to the loss function.

Problematic for deep feedforward networks

- Particularly with sigmoid activation functions.
- When the sigmoid function value is either too high or too low, the derivative (orange line) becomes very small i.e.  $\ll 1$ .



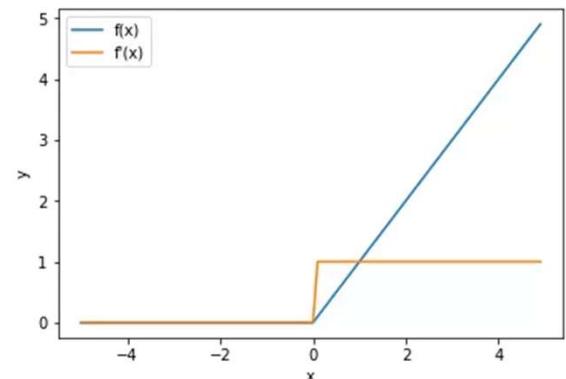
Read more here: <https://adventuresinmachinelearning.com/vanishing-gradient-problem-tensorflow/>

# Vanishing gradients

ReLU activation can partly solve the problem

No degradation of the error signal

But, certain weights can be cancelled out whenever there is a negative input into a given neuron



# Vanishing gradients

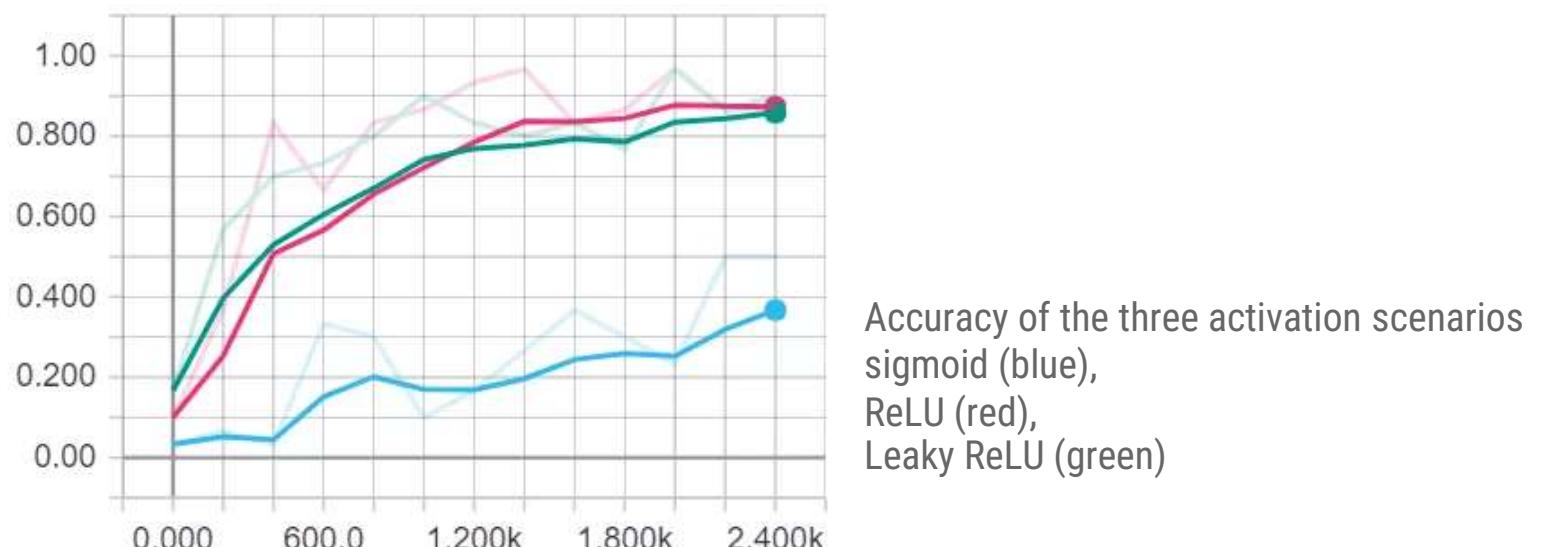
Leaky ReLU

$$f(x) = \max(\alpha x, x)$$

The good thing about the Leaky ReLU activation function is that the derivative when  $x$  is below zero is alpha – i.e. it is a small but no longer 0.

# Vanishing gradients (in practice)

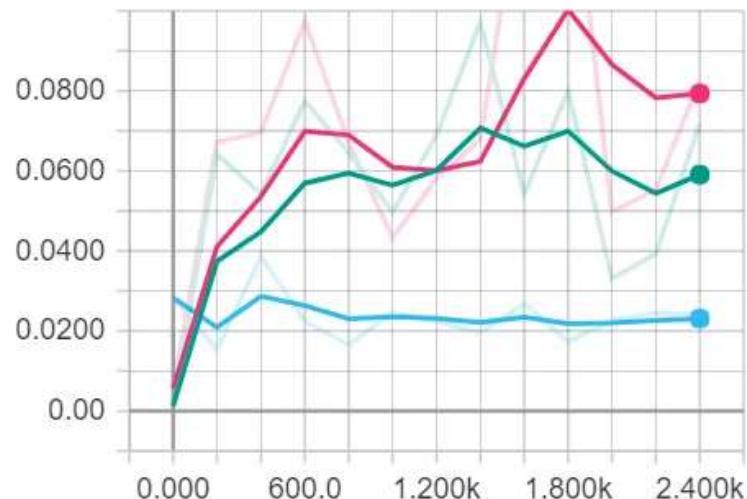
MNIST using a 7-layer densely connected network



# Vanishing gradients (in practice)

MNIST using a 7-layer densely connected network

Mean absolute gradient logs during training

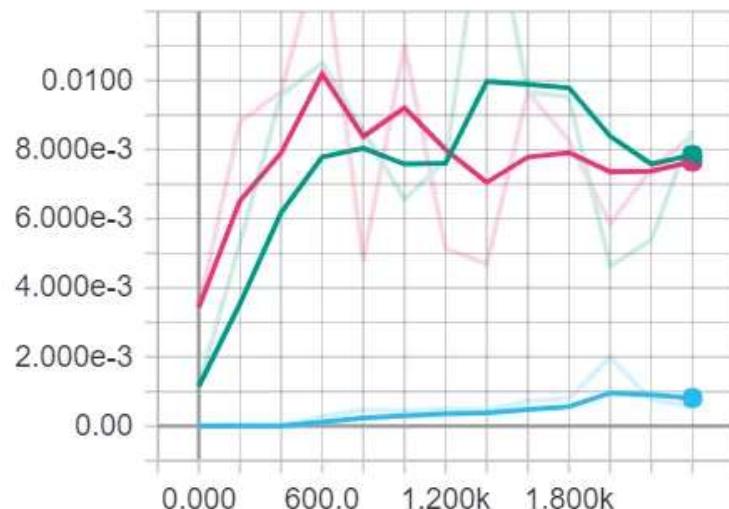


Mean absolute gradients  
Output layer (6th layer)  
sigmoid (blue),  
ReLU (red),  
Leaky ReLU (green)

# Vanishing gradients (in practice)

MNIST using a 7-layer densely connected network

Mean absolute gradient logs during training



Mean absolute gradients  
1st layer  
sigmoid (blue),  
ReLU (red),  
Leaky ReLU (green)

# From SimpleRNN to LSTM

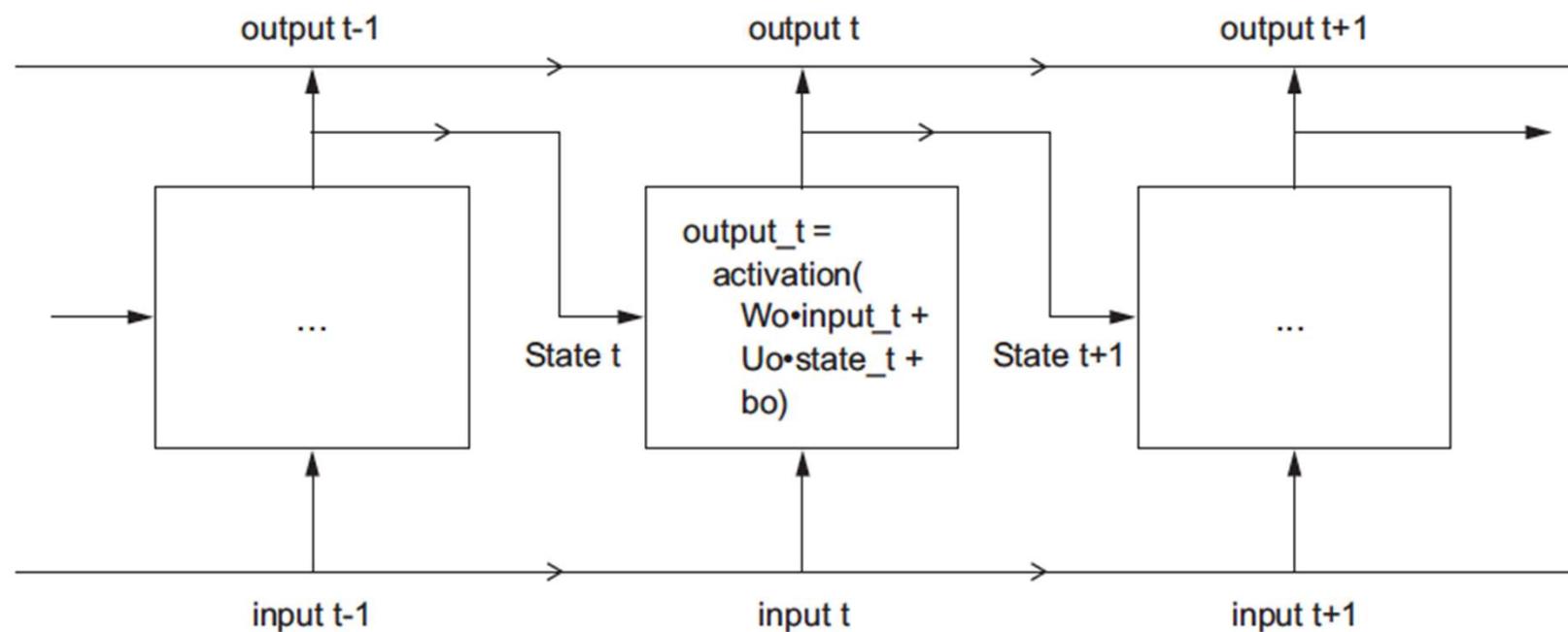


Figure 6.13 The starting point of an LSTM layer: a SimpleRNN

# From SimpleRNN to LSTM

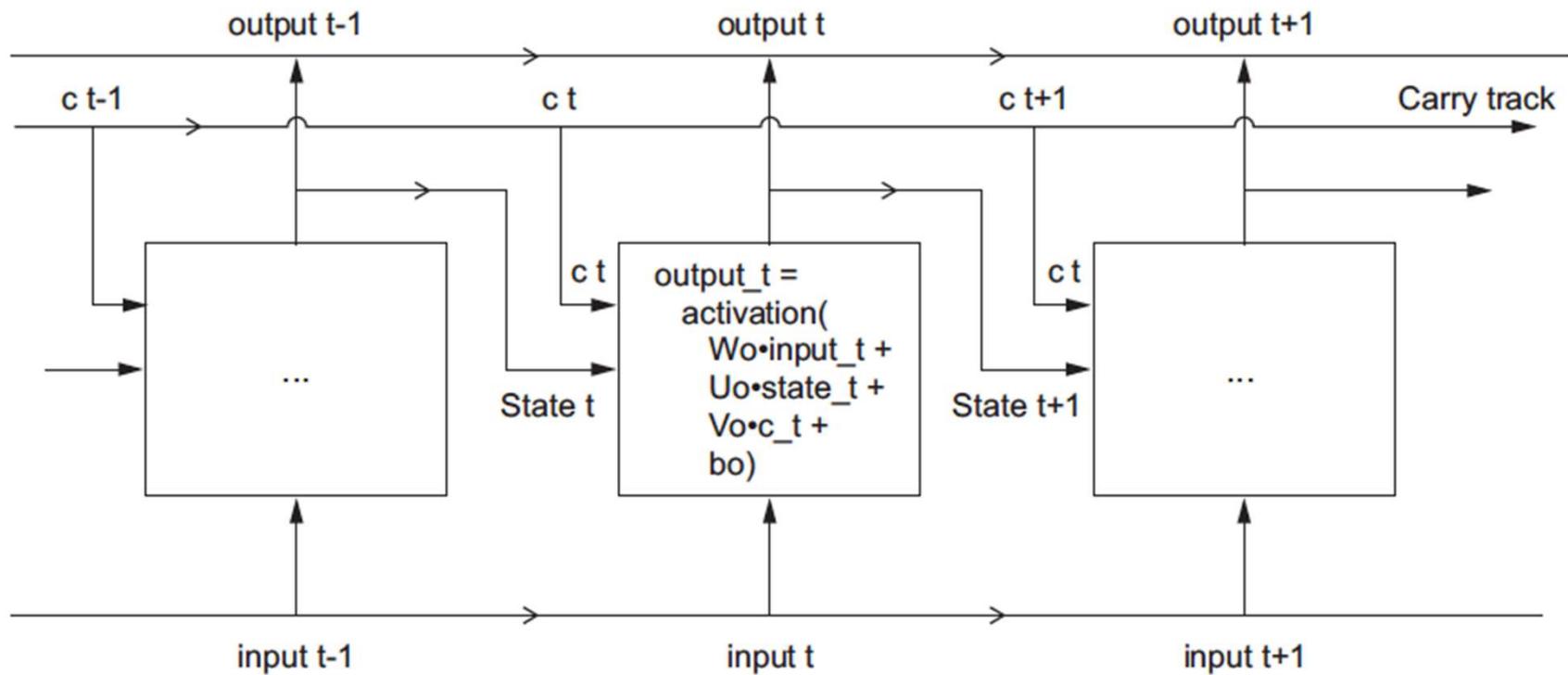


Figure 6.14 Going from a SimpleRNN to an LSTM: adding a carry track

# From SimpleRNN to LSTM

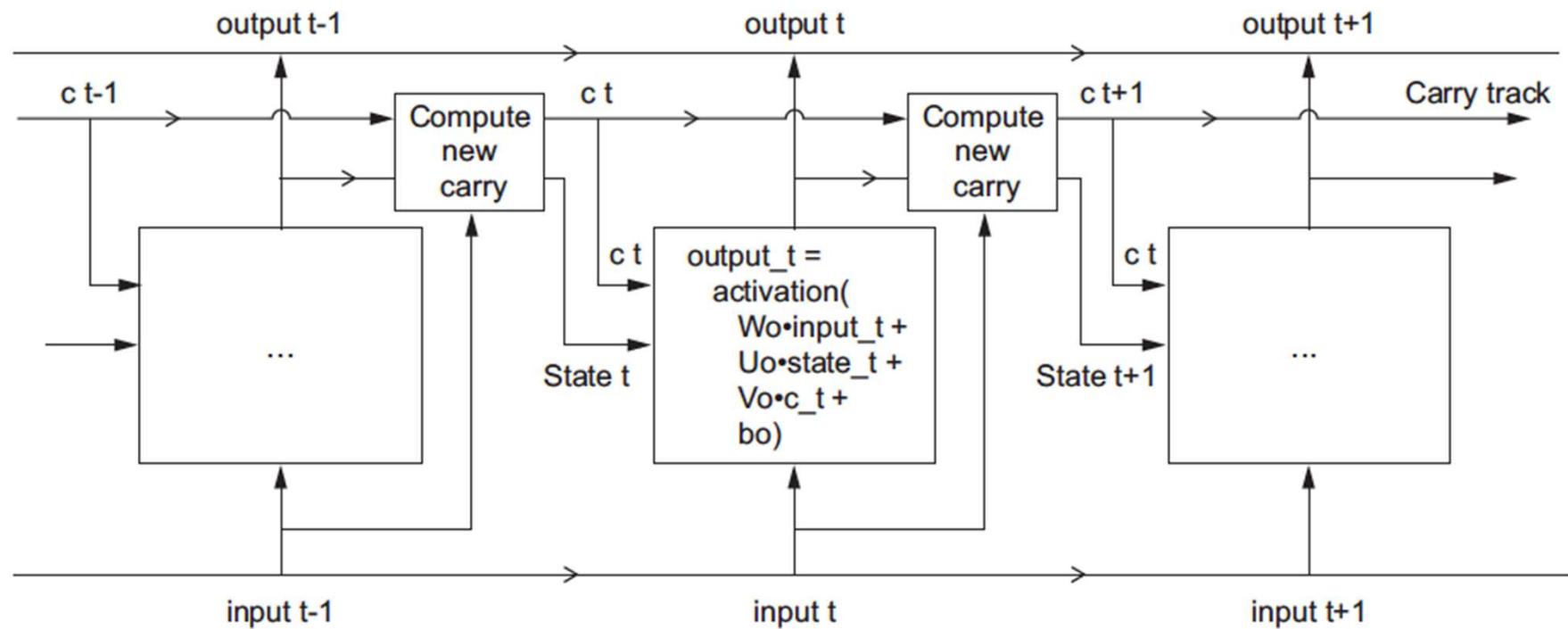


Figure 6.15 Anatomy of an LSTM

# From SimpleRNN to LSTM

## **Listing 6.25 Pseudocode details of the LSTM architecture (1/2)**

```
output_t = activation(dot(state_t, Uo) + dot(input_t, Wo) + dot(C_t, Vo) + bo)  
i_t = activation(dot(state_t, Ui) + dot(input_t, Wi) + bi)  
f_t = activation(dot(state_t, Uf) + dot(input_t, Wf) + bf)  
k_t = activation(dot(state_t, Uk) + dot(input_t, Wk) + bk)
```

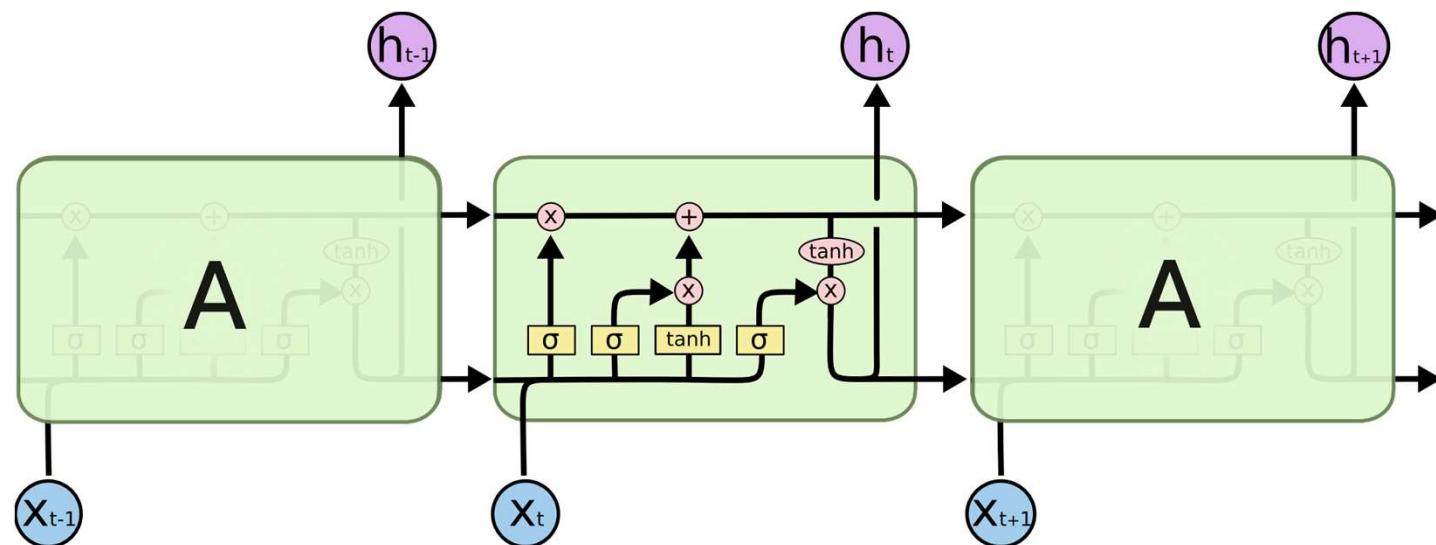
You obtain the new carry state (the next  $c_{t+1}$ ) by combining  $i_t$ ,  $f_t$ , and  $k_t$ .

## **Listing 6.26 Pseudocode details of the LSTM architecture (2/2)**

```
c_t+1 = i_t * k_t + c_t * f_t
```

# LSTMs: a deeper look

<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

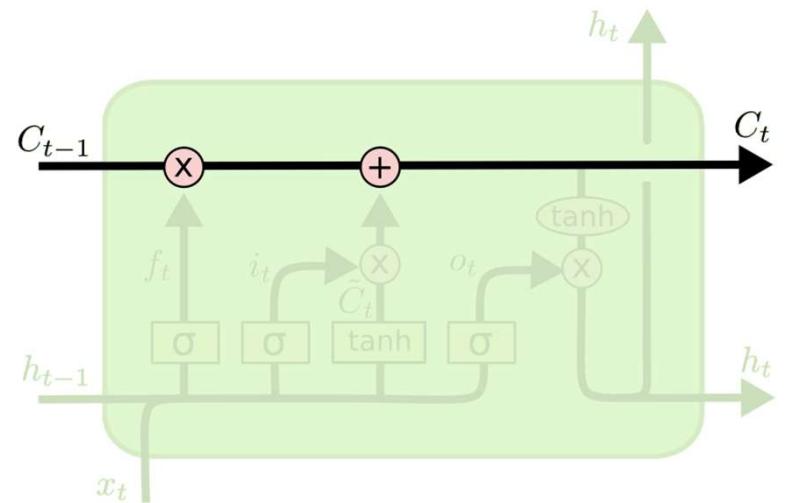


# LSTMs

## *Cell state*

- kind of like a conveyor belt

It runs straight down the entire chain,  
with only some minor linear interactions.

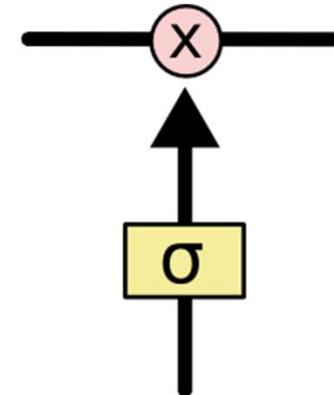


It's very easy for information to just flow along it unchanged.

# LSTMs

## Gates

A way to optionally let information through.



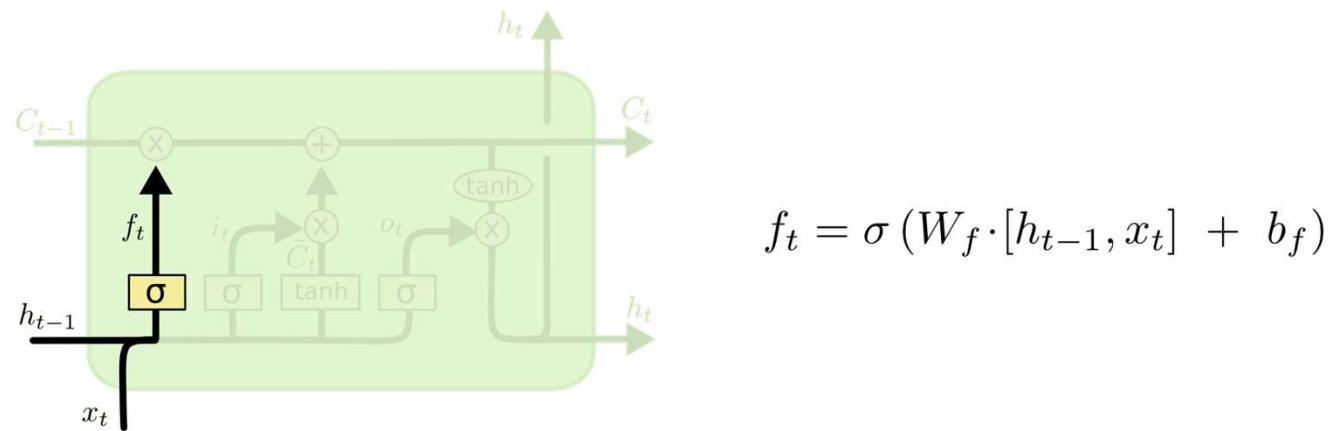
They are composed of a sigmoid neural net layer and a pointwise multiplication operation.

- A value of zero means “let nothing through,” while a value of one means “let everything through!”

# LSTMs

*Forget gate layer*

What information we're going to throw away from the cell state?

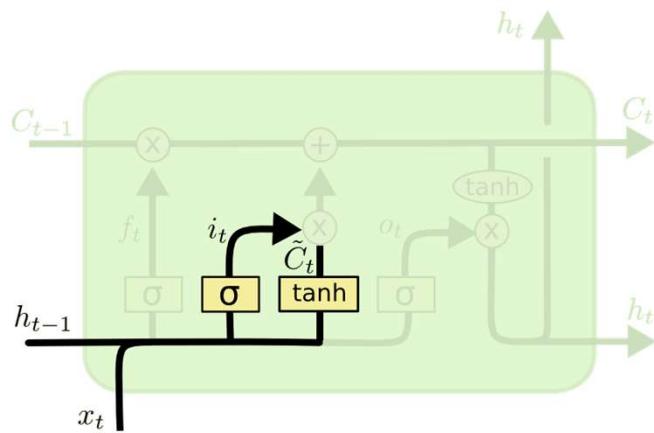


# LSTMs

What new information we're going to store in the cell state?

Two parts:

1. A sigmoid layer called the *input gate layer* decides which values we'll update.
2. A tanh layer creates a vector of new candidate values, that could be added to the state.

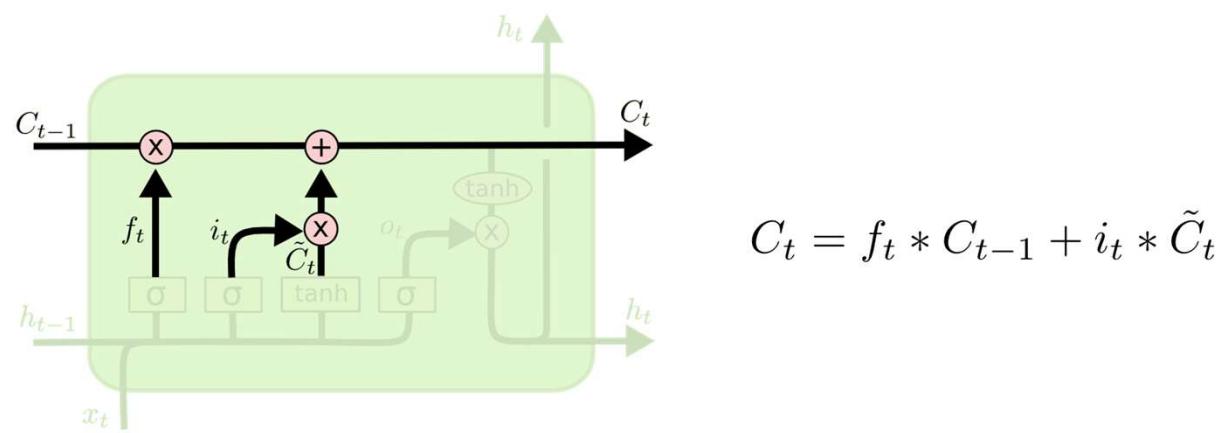


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

# LSTMs

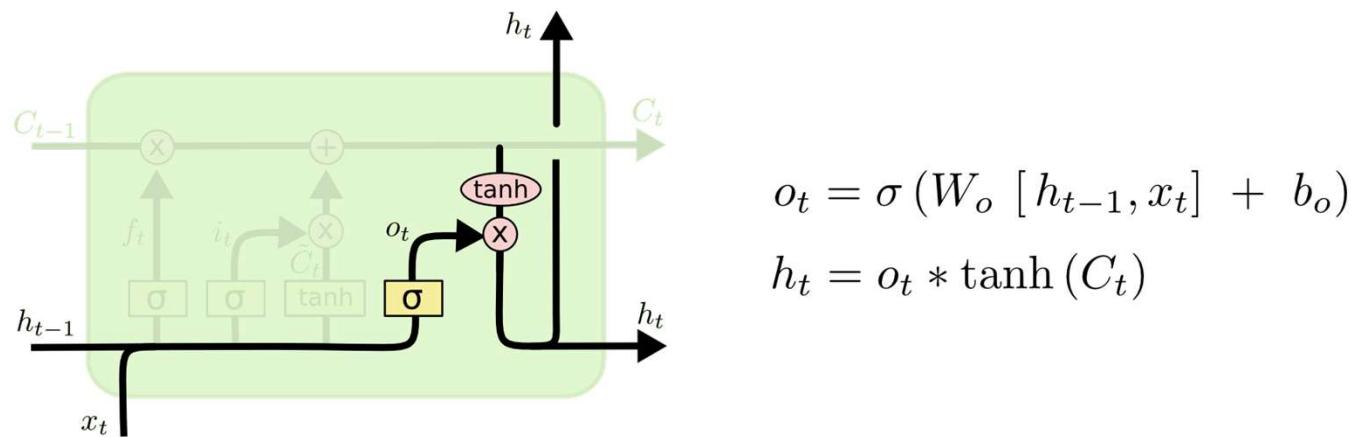
It's now time to update the old cell state  $C_{t-1}$  to  $C_t$



# LSTMs

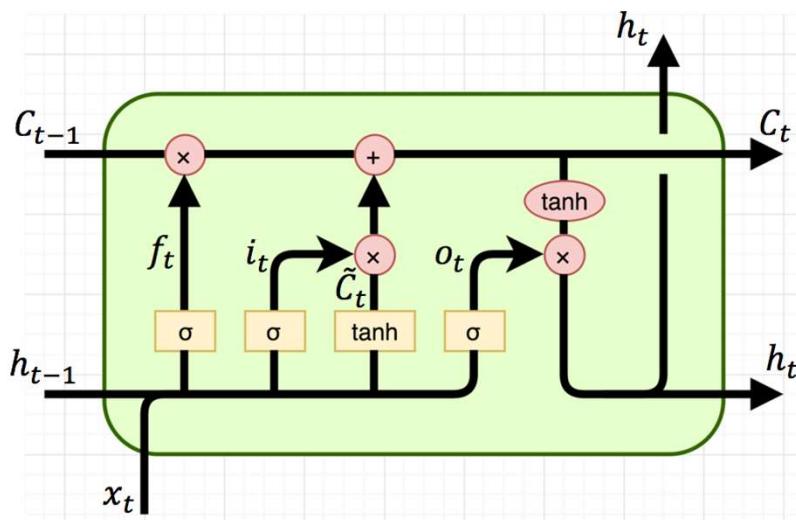
Finally, we need to decide what we're going to output. This output will be based on our cell state but will be a filtered version.

1. A sigmoid layer to decide what parts of the cell state to output.
2. Put the cell state through tanh (to push the values to be between  $-1$  and  $1$ ) and multiply it by the output of the sigmoid gate, so that we only output the parts we decided to.

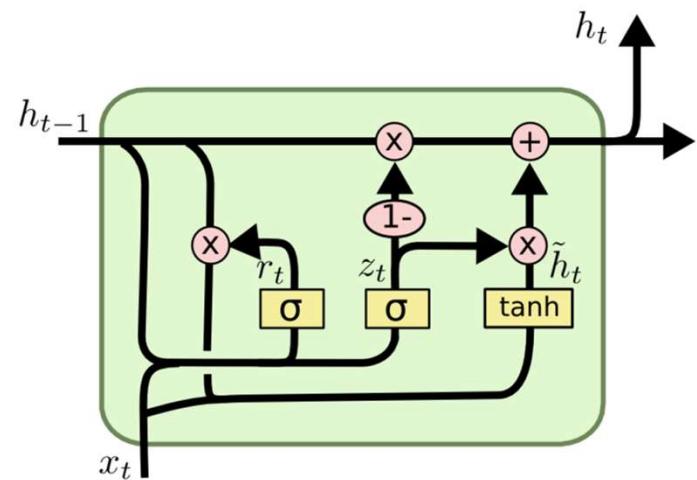


# GRU: Gated Recurrent Units

Kyunghyun Cho et al (2014)

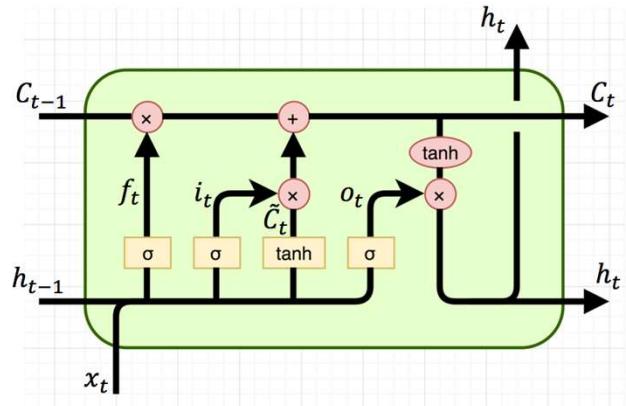


(a) Long Short-Term Memory

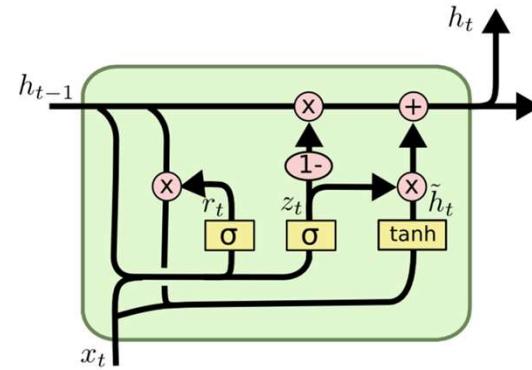


(b) Gated Recurrent Unit

# GRU: Gated Recurrent Units



(a) Long Short-Term Memory



(b) Gated Recurrent Unit

$$\begin{aligned}
 i_t &= \sigma(x_t U^i + h_{t-1} W^i) \\
 f_t &= \sigma(x_t U^f + h_{t-1} W^f) \\
 o_t &= \sigma(x_t U^o + h_{t-1} W^o) \\
 \tilde{C}_t &= \tanh(x_t U^g + h_{t-1} W^g) \\
 C_t &= \sigma(f_t * C_{t-1} + i_t * \tilde{C}_t) \\
 h_t &= \tanh(C_t) * o_t
 \end{aligned}$$

$$\begin{aligned}
 z_t &= \sigma(x_t U^z + h_{t-1} W^z) \\
 r_t &= \sigma(x_t U^r + h_{t-1} W^r) \\
 \tilde{h}_t &= \tanh(x_t U^h + (r_t * h_{t-1}) W^h) \\
 h_t &= (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t
 \end{aligned}$$

# Evaluating language models

- The standard evaluation metric for Language Models is **perplexity**

$$\text{perplexity} = \prod_{t=1}^T \left( \underbrace{\frac{1}{P_{\text{LM}}(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})}}_{\text{Inverse probability of corpus, according to Language Model}} \right)^{1/T}$$

Normalized by  
number of words

- This is equal to the exponential of the cross-entropy loss :

$$= \prod_{t=1}^T \left( \frac{1}{\hat{y}_{\mathbf{x}_{t+1}}^{(t)}} \right)^{1/T} = \exp \left( \frac{1}{T} \sum_{t=1}^T -\log \hat{y}_{\mathbf{x}_{t+1}}^{(t)} \right) = \exp(J(\theta))$$

# RNNs for sequence labeling

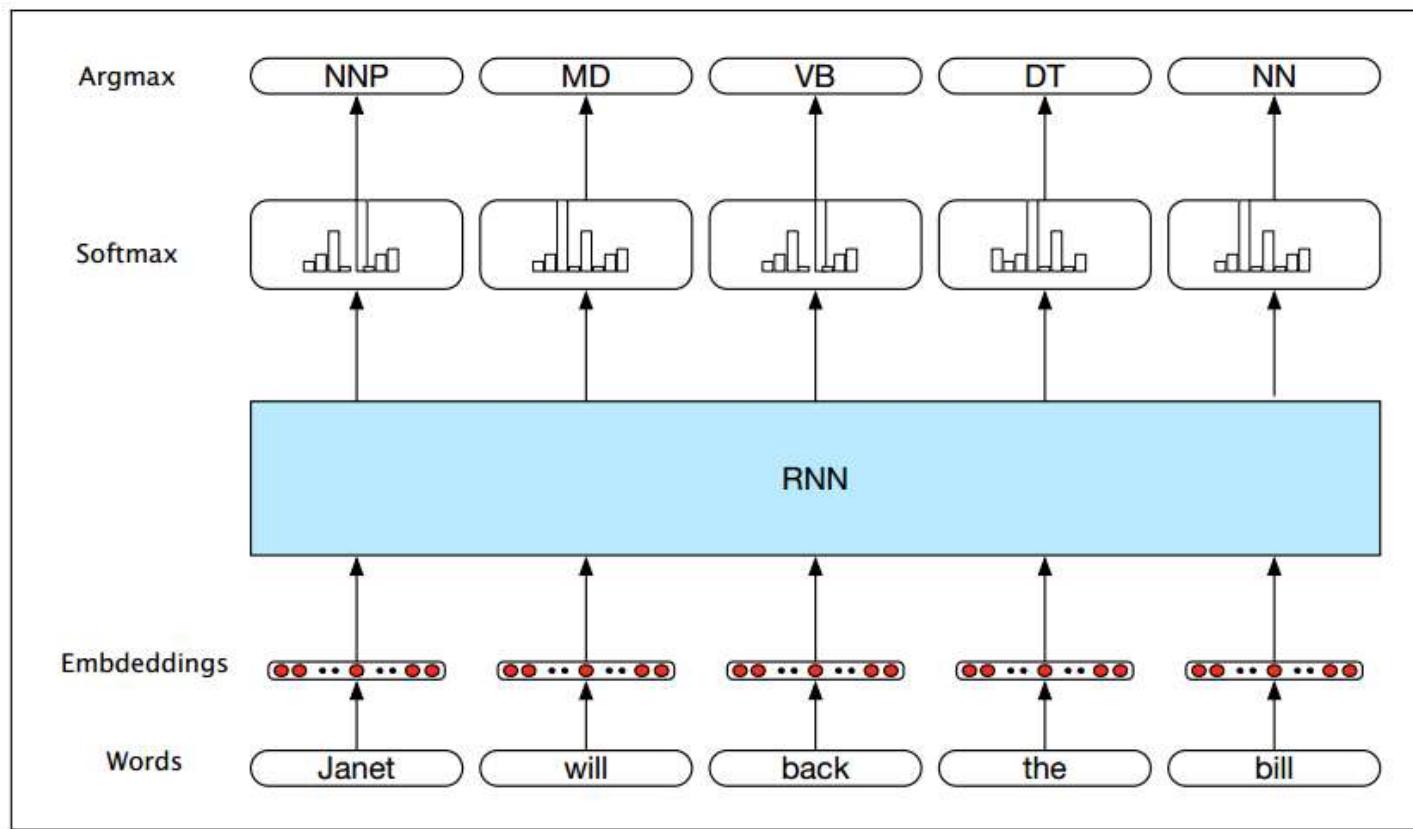


Image source

# RNNs for language modeling -- training

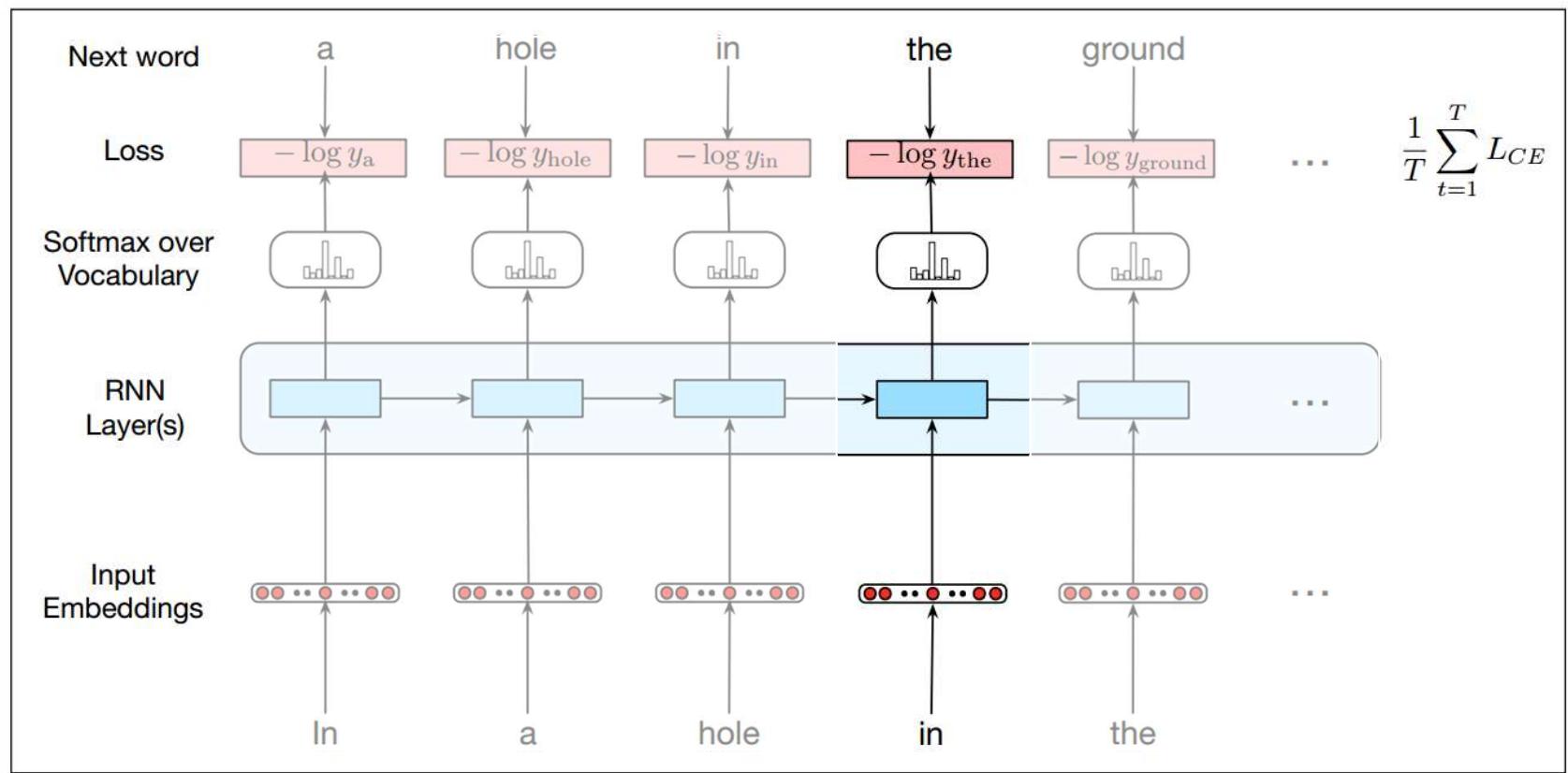


Image source

# RNNs for language modeling -- autoregressive generation

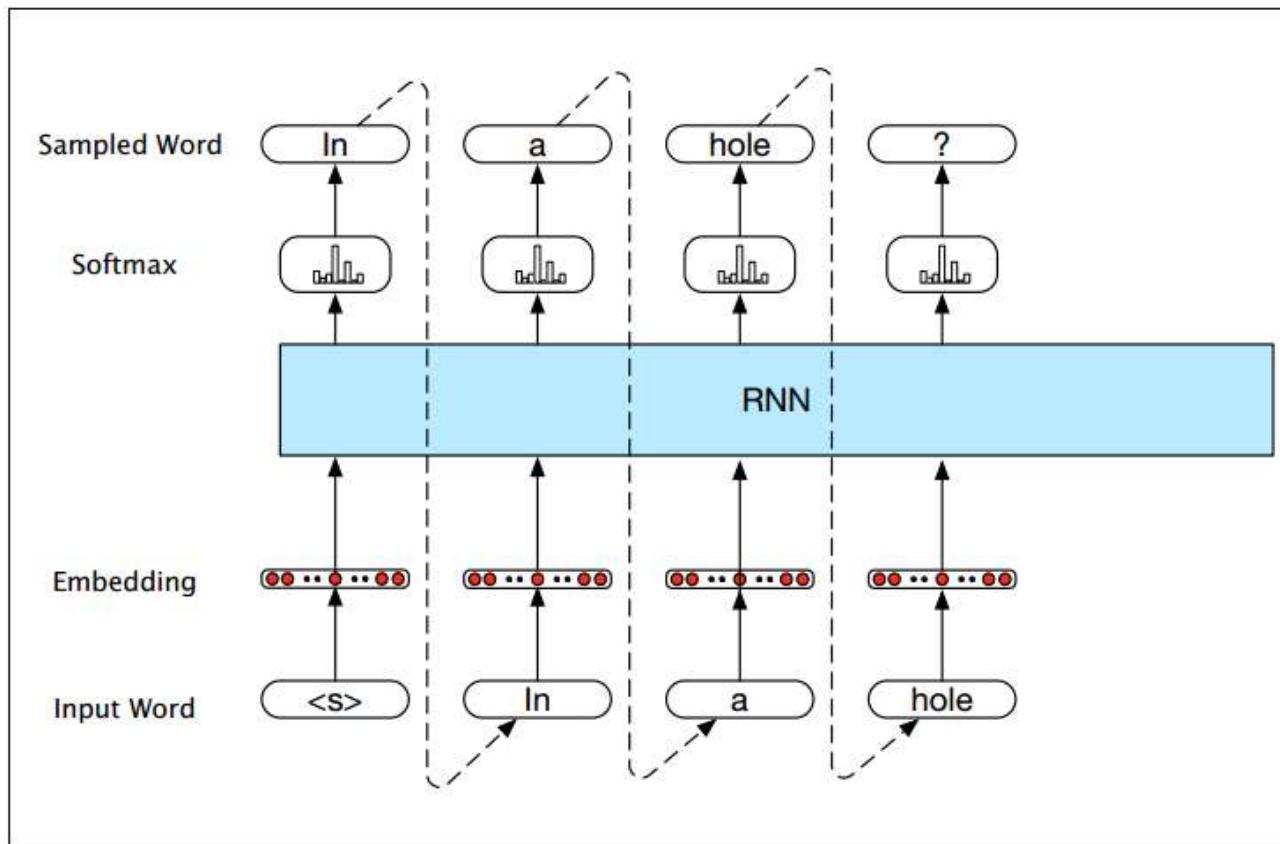
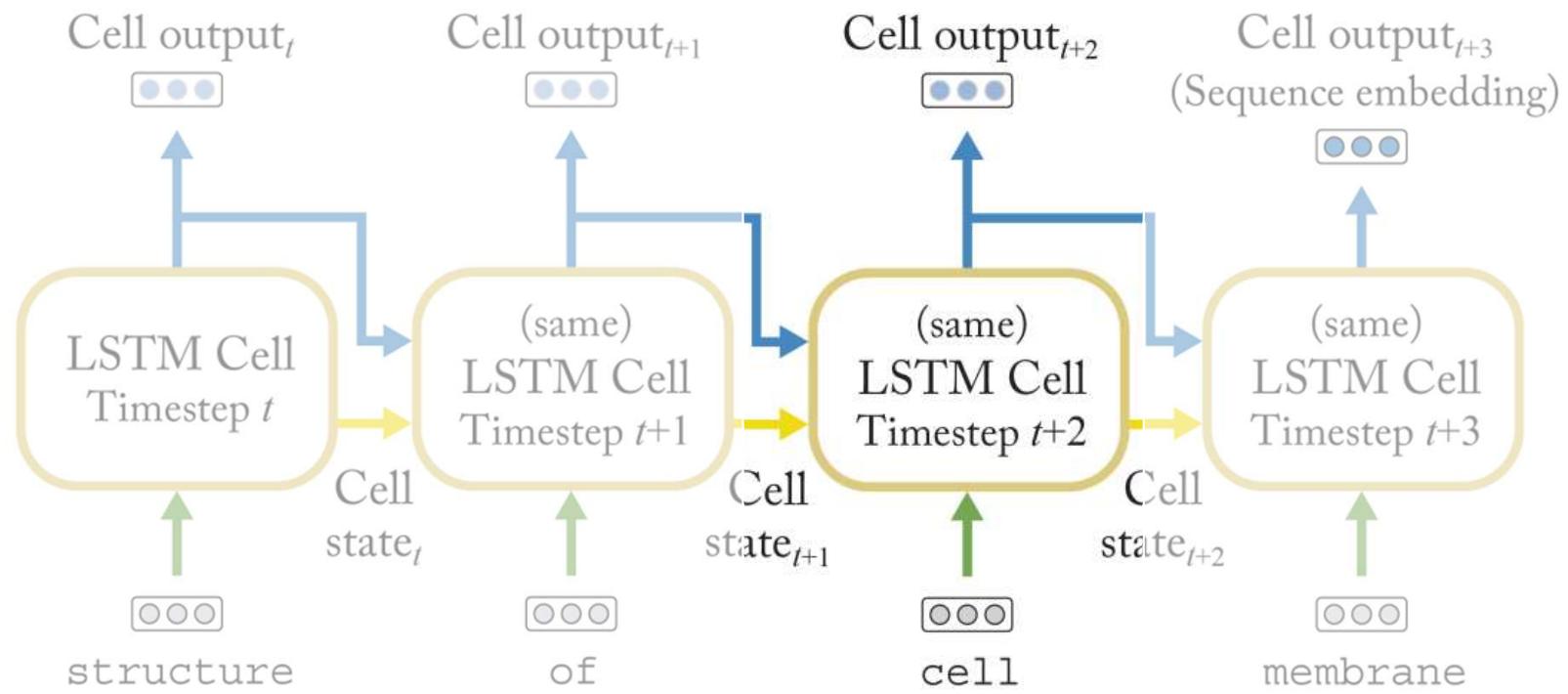


Image source

# RNNs for language modeling



# RNNs for language modeling

- Persian poetry!

The screenshot shows a user interface for generating Persian poetry. At the top, there is a text input field containing the word "study" and a "Generate" button. Below this, a blue header bar displays "[1s] Ready". Underneath the bar, the word "Poem" is shown in a dark box. To the right of "Poem" are five yellow stars. Below the stars, a poem is displayed in Persian:

The scientific curiosity!  
And learn to use an equal colour vision,  
Try and teach the public policy,  
A state of scientific discipline.

روان را به دانش به کار آوری

به توران به خواب اندر آرد سری

# RNNs for language modeling

- Persian poetry!

نیست کو معلین  
دامشت خانقه به من من از در میخانه بست  
در کار ما به در این خیال ما را ببرد  
ز شاه ساهد می و زلف تو را تا برسد  
تا به درد تو به میخانه در آب است  
که باز من آباد به در خسابت  
به خاک دل بر ت

دید او را اندر آمد ای جوان  
گفت می دامن امیران را به جان  
در دم و زشته که افزونتر بود  
چون تو این گفت این جهان خالی بود  
گر تو ما آن را کجا زید خداست  
گفت حقست آن دم آن هدیه هاست  
کیمیا بهتر از آن زاری شوی  
قبله و قدرت از آن علت شوی  
بازی زین بندها قربان شود  
گر نشاید خانه را پستی شود  
تا بود کار دگر خندان کند  
که پذیرای خدا پیدا کند  
یا چو او گردان همی گویید هین  
در پشیمانی در آید بر زمین

سواران گردنکشان دسته دید  
خردمند را او بدان خسته دید  
سکندر نگه کرد پس پهلوان  
به بدخواه شد شاد و روشن روان  
سپاه اندر آمد به پیش سوار  
خردمند و شایسته کارزار  
بفرمود تا بنده آگاه دید  
چنین تا بر شاه ایران کشید  
نهادند چیزی که پوشیده بود  
جهان را درم داد و دینار بود  
سران افسر از گوهر شاهوار  
نخست آفرین کرد بر کردگار  
چو بهرام بشنید گریان شدند

# RNNs for language modeling

- Wikipedia articles

Naturalism and decision for the majority of Arab countries' capitalide was grounded by the Irish language by [[John Clair]], [[An Imperial Japanese Revolt]], associated with Guangzham's sovereignty. His generals were the powerful ruler of the Portugal in the [[Protestant Immineners]], which could be said to be directly in Cantonese Communication, which followed a ceremony and set inspired prison, training. The emperor travelled back to [[Antioch, Perth, October 25|21]] to note, the Kingdom of Costa Rica, unsuccessful fashioned the [[Thrales]], [[Cynth's Dajoard]], known in western [[Scotland]], near Italy to the conquest of India with the conflict. Copyright was the succession of independence in the slop of Syrian influence that was a famous German movement based on a more popular servitious, non-doctrinal and sexual power post. Many governments recognize the military housing of the [[Civil Liberalization and Infantry Resolution 265 National Party in Hungary]], that is sympathetic to be to the [[Punjab Resolution]] (PJS)[<http://www.humah.yahoo.com/guardian.cfm/7754800786d17551963s89.htm>] Official economics Adjoint for the Nazism, Montgomery was swear to advance to the resources for those Socialism's rule, was starting to signing a major tripad of aid exile.]

# RNNs for language modeling

- Wikipedia articles

```
{ { cite journal | id=Cerling Nonforest Department|format=Newlymeslated|none } }
' 'www.e-complete''.

'''See also'''': [[List of ethical consent processing]]

== See also ==
*[[Iender dome of the ED]]
*[[Anti-autism]]

====[[Religion|Religion]]====
*[[French Writings]]
*[[Maria]]
*[[Revelation]]
*[[Mount Agamul]]
```

# RNNs for language modeling

- XML

```
<page>
    <title>Antichrist</title>
    <id>865</id>
    <revision>
        <id>15900676</id>
        <timestamp>2002-08-03T18:14:12Z</timestamp>
        <contributor>
            <username>Paris</username>
            <id>23</id>
        </contributor>
        <minor />
        <comment>Automated conversion</comment>
        <text xml:space="preserve">#REDIRECT [[Christianity]]</text>
    </revision>
</page>
```

# RNNs for language modeling

For  $\bigoplus_{n=1,\dots,m} \mathcal{L}_{m,n} = 0$ , hence we can find a closed subset  $\mathcal{H}$  in  $\mathcal{H}$  and any sets  $\mathcal{F}$  on  $X$ ,  $U$  is a closed immersion of  $S$ , then  $U \rightarrow T$  is a separated algebraic space.

*Proof.* Proof of (1). It also start we get

$$S = \text{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by  $\coprod Z \times_U U \rightarrow V$ . Consider the maps  $M$  along the set of points  $\text{Sch}_{fppf}$  and  $U \rightarrow U$  is the fibre category of  $S$  in  $U$  in Section, ?? and the fact that any  $U$  affine, see Morphisms, Lemma ???. Hence we obtain a scheme  $S$  and any open subset  $W \subset U$  in  $\text{Sh}(G)$  such that  $\text{Spec}(R') \rightarrow S$  is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that  $f_i$  is of finite presentation over  $S$ . We claim that  $\mathcal{O}_{X,x}$  is a scheme where  $x, x', s'' \in S'$  such that  $\mathcal{O}_{X,x'} \rightarrow \mathcal{O}'_{X',x'}$  is separated. By Algebra, Lemma ?? we can define a map of complexes  $\text{GL}_{S'}(x'/S'')$  and we win.  $\square$

To prove study we see that  $\mathcal{F}|_U$  is a covering of  $X'$ , and  $\mathcal{T}_i$  is an object of  $\mathcal{F}_{X/S}$  for  $i > 0$  and  $\mathcal{F}_p$  exists and let  $\mathcal{F}_i$  be a presheaf of  $\mathcal{O}_X$ -modules on  $C$  as a  $\mathcal{F}$ -module. In particular  $\mathcal{F} = U/\mathcal{F}$  we have to show that

$$\widetilde{\mathcal{I}}^\bullet = \mathcal{I}^\bullet \otimes_{\text{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1} \mathcal{F}$$

is a unique morphism of algebraic stacks. Note that

$$\text{Arrows} = (\text{Sch}/S)^{opp}_{fppf}, (\text{Sch}/S)_{fppf}$$

and

$$V = \Gamma(S, \mathcal{O}) \rightarrow (U, \text{Spec}(A))$$

is an open subset of  $X$ . Thus  $U$  is affine. This is a continuous map of  $X$  is the inverse, the groupoid scheme  $S$ .

*Proof.* See discussion of sheaves of sets.  $\square$

The result for prove any open covering follows from the less of Example ???. It may replace  $S$  by  $X_{\text{spaces},\text{étale}}$  which gives an open subspace of  $X$  and  $T$  equal to  $S_{\text{Zar}}$ , see Descent, Lemma ???. Namely, by Lemma ?? we see that  $R$  is geometrically regular over  $S$ .

**Lemma 0.1.** Assume (3) and (3) by the construction in the description.

Suppose  $X = \lim |X|$  (by the formal open covering  $X$  and a single map  $\underline{\text{Proj}}_X(\mathcal{A}) = \text{Spec}(B)$  over  $U$  compatible with the complex

$$\text{Set}(\mathcal{A}) = \Gamma(X, \mathcal{O}_{X,\mathcal{O}_X}).$$

When in this case of to show that  $\mathcal{Q} \rightarrow \mathcal{C}_{Z/X}$  is stable under the following result in the second conditions of (1), and (3). This finishes the proof. By Definition ?? (without element is when the closed subschemes are catenary. If  $T$  is surjective we may assume that  $T$  is connected with residue fields of  $S$ . Moreover there exists a closed subspace  $Z \subset X$  of  $X$  where  $U$  in  $X'$  is proper (some defining as a closed subset of the uniqueness it suffices to check that the following theorem

(1)  $f$  is locally of finite type. Since  $S = \text{Spec}(R)$  and  $Y = \text{Spec}(R)$ .

*Proof.* This is form all sheaves of sheaves on  $X$ . But given a scheme  $U$  and a surjective étale morphism  $U \rightarrow X$ . Let  $U \cap U = \coprod_{i=1,\dots,n} U_i$  be the scheme  $X$  over  $S$  at the schemes  $X_i \rightarrow X$  and  $U = \lim_i X_i$ .  $\square$

The following lemma surjective restrocomposes of this implies that  $\mathcal{F}_{x_0} = \mathcal{F}_{x_0} = \mathcal{F}_{X,\dots,0}$ .

**Lemma 0.2.** Let  $X$  be a locally Noetherian scheme over  $S$ ,  $E = \mathcal{F}_{X/S}$ . Set  $\mathcal{I} = \mathcal{J}_1 \subset \mathcal{I}'_n$ . Since  $\mathcal{I}'_n \subset \mathcal{I}^n$  are nonzero over  $i_0 \leq p$  is a subset of  $\mathcal{J}_{n,0} \circ \mathcal{A}_2$  works.

**Lemma 0.3.** In Situation ???. Hence we may assume  $q' = 0$ .

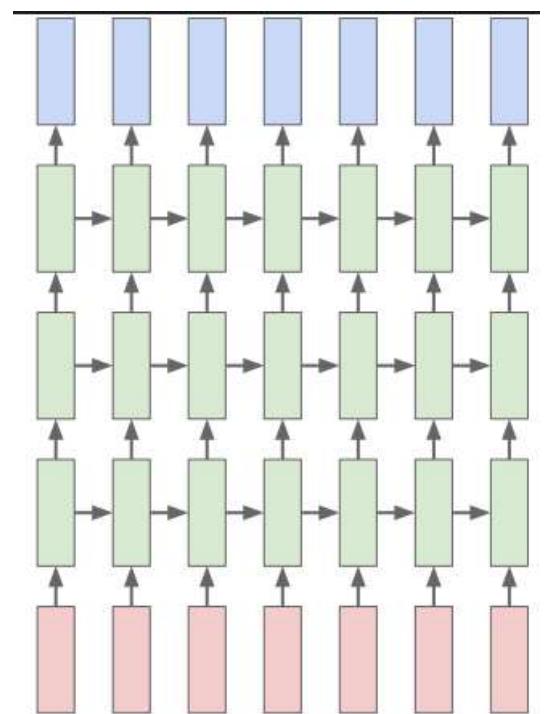
*Proof.* We will use the property we see that  $p$  is the next functor (??). On the other hand, by Lemma ?? we see that

$$D(\mathcal{O}_{X'}) = \mathcal{O}_X(D)$$

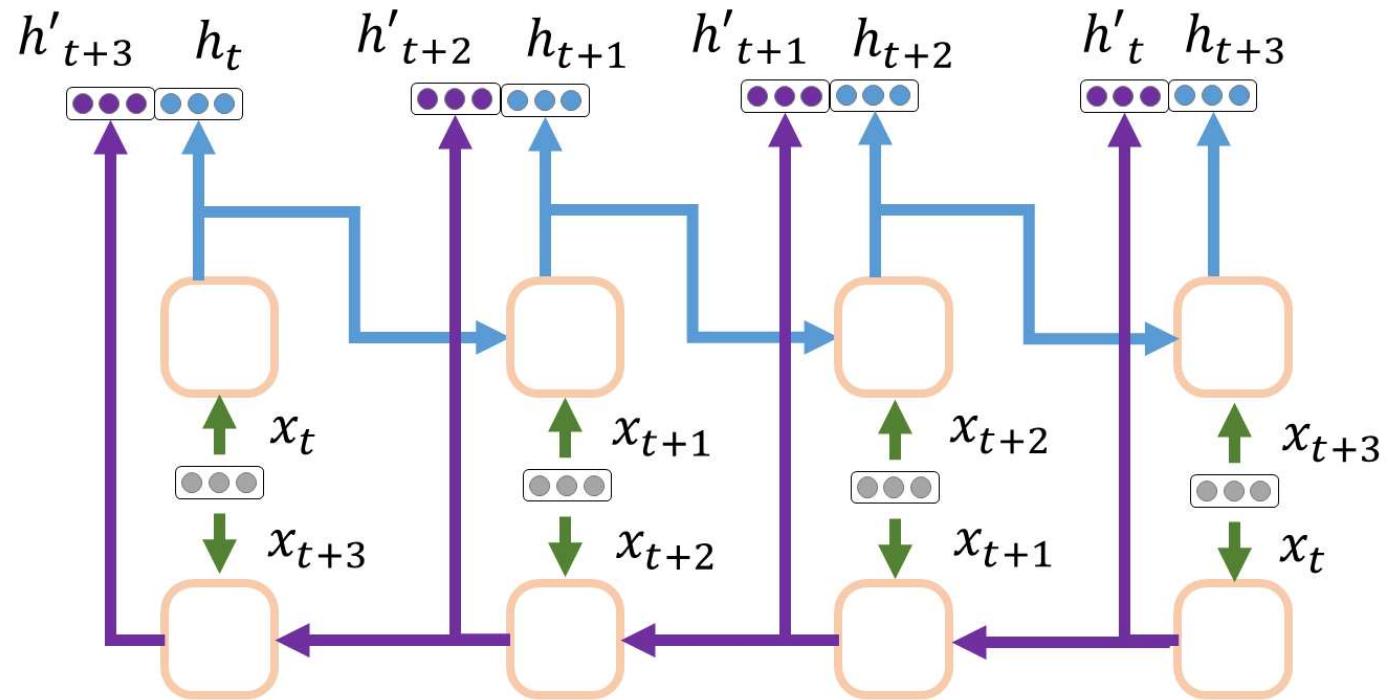
where  $K$  is an  $F$ -algebra where  $\delta_{n+1}$  is a scheme over  $S$ .  $\square$

# Stacking recurrent layers

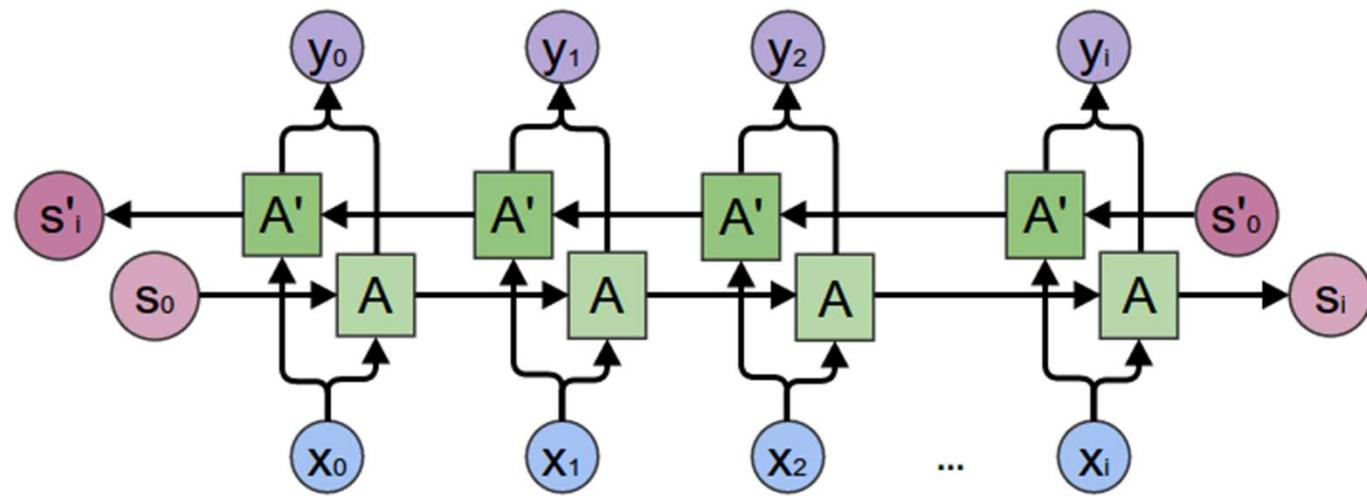
- It is generally a good idea to increase the capacity of your network until overfitting becomes your primary obstacle (assuming that you are already taking basic steps to mitigate overfitting, such as using dropout).
- Recurrent layer stacking is a classic way to build more powerful recurrent networks:
  - For instance, what currently powers the Google translate algorithm is a stack of seven large LSTM layers -- that's huge.



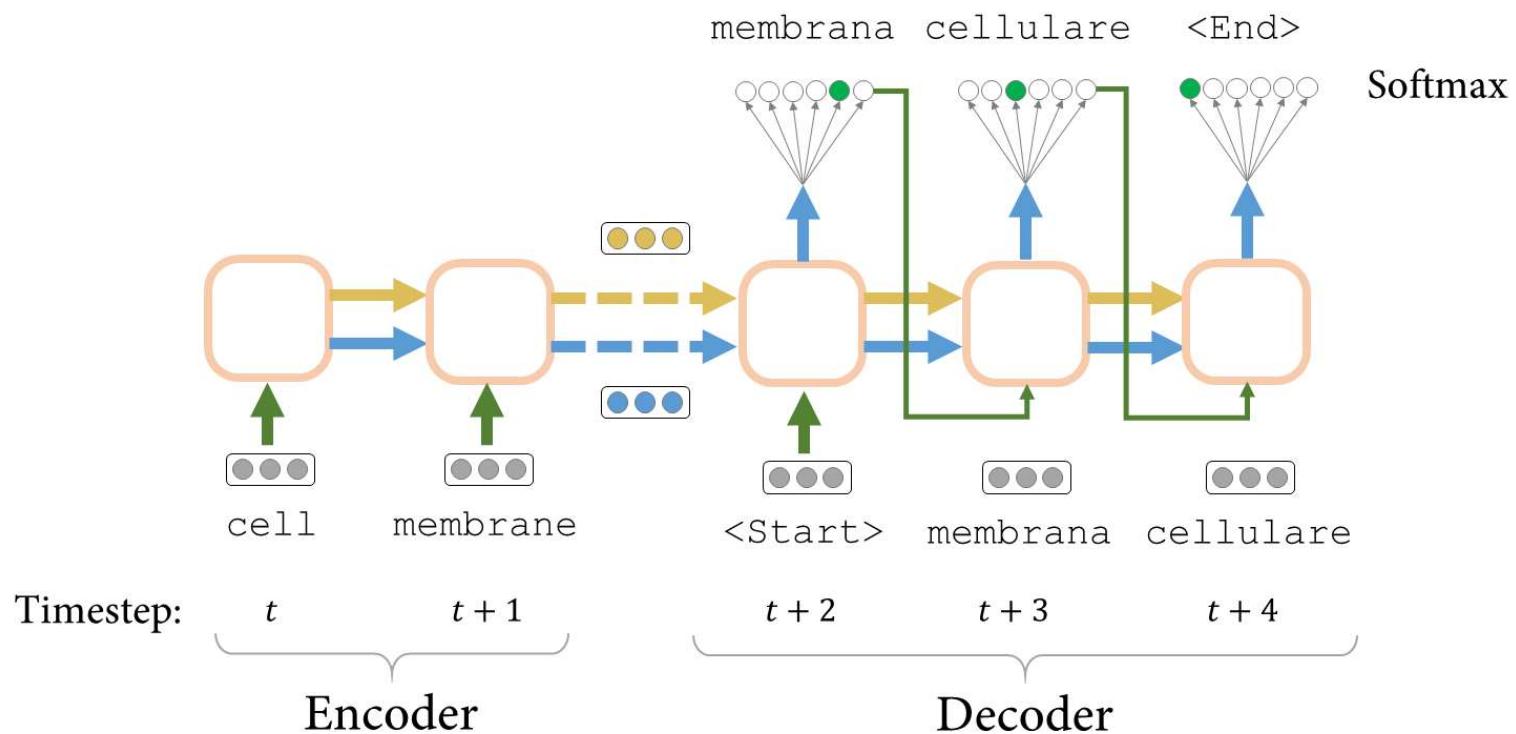
# Bidirectional RNNs



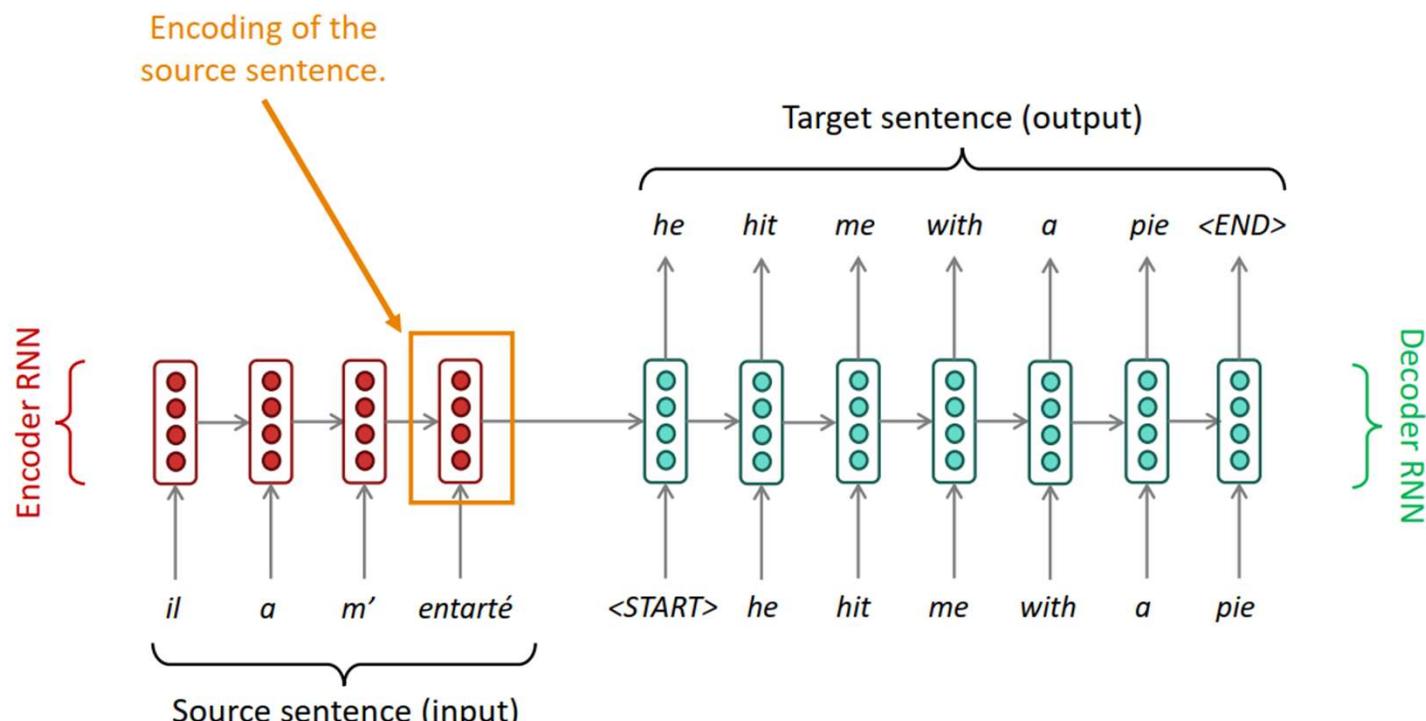
# Bidirectional RNNs



# Sequence transduction (Seq2seq)



# Seq2seq: the bottleneck problem



Problems with this architecture?

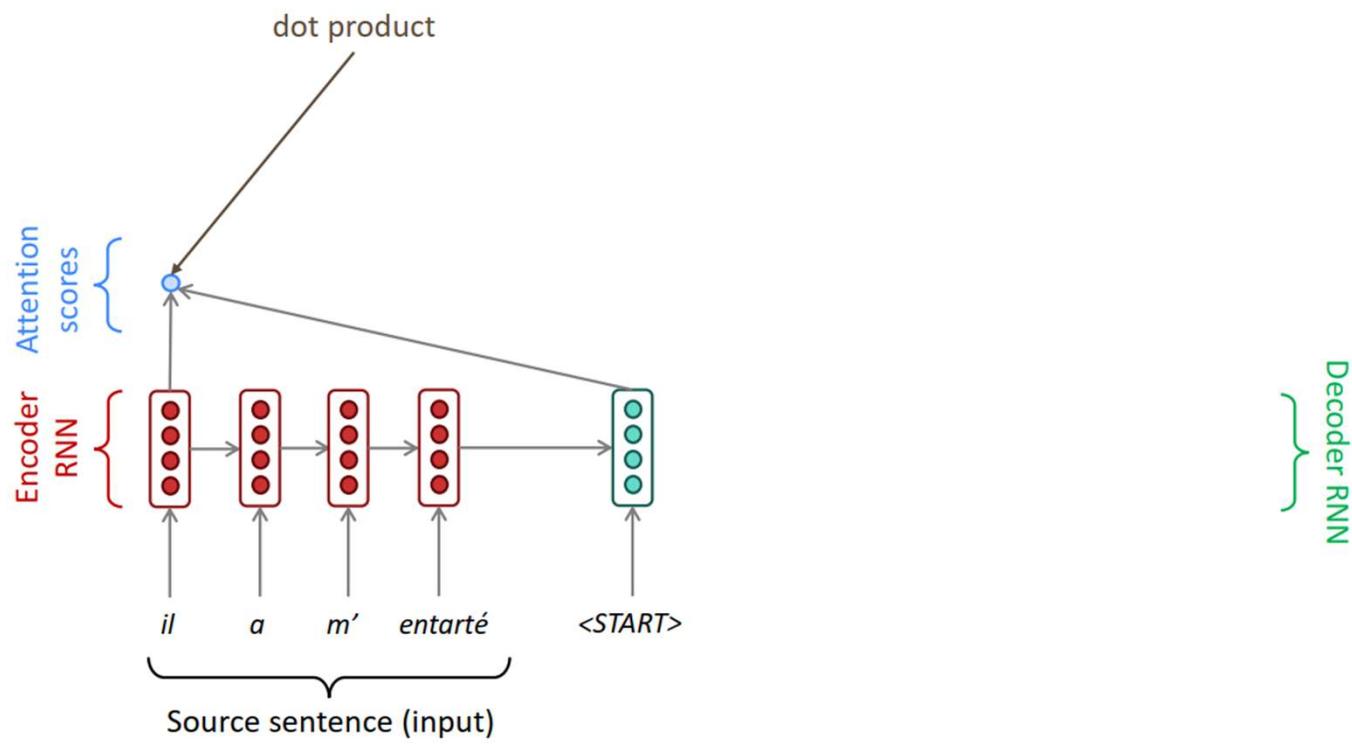
# Attention

Attention provides a solution to the bottleneck problem.

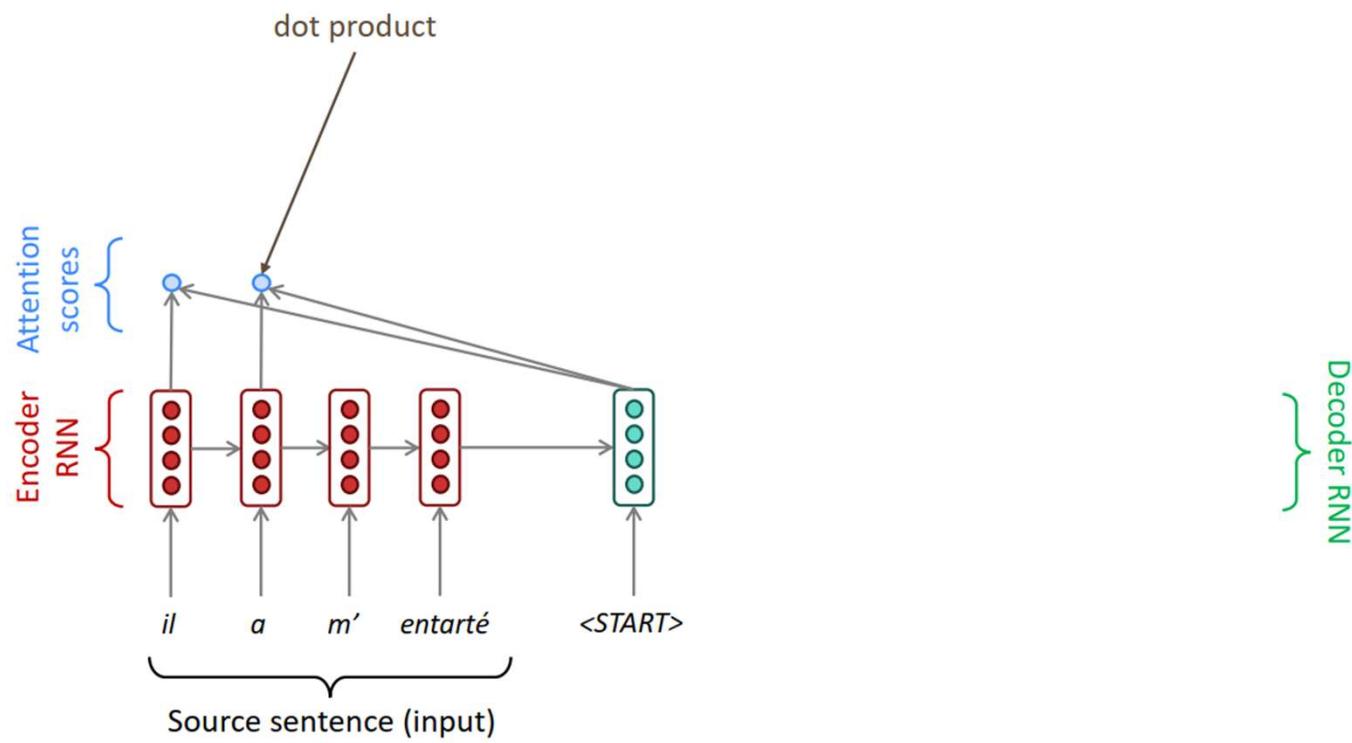
Core idea:

- on each step of the decoder, use direct connection to the encoder to focus on a particular part of the source sequence

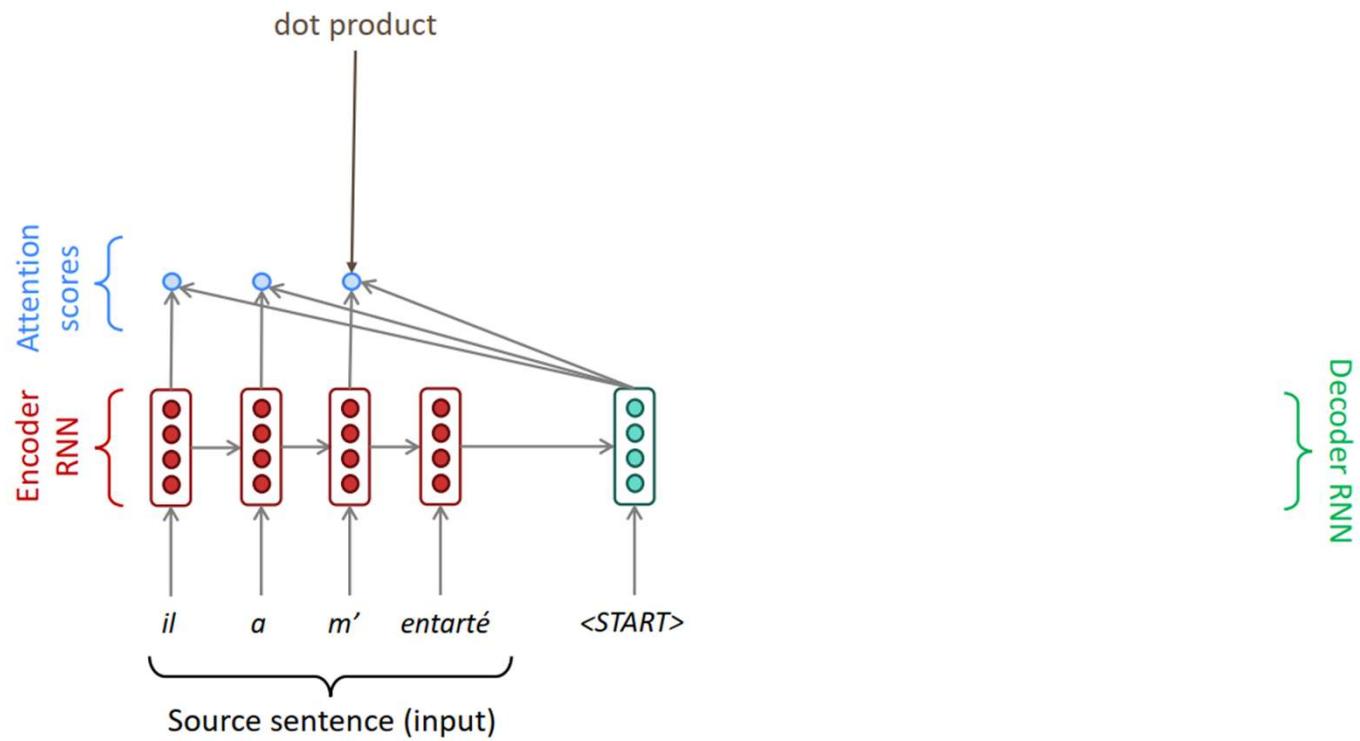
# Seq2seq with attention



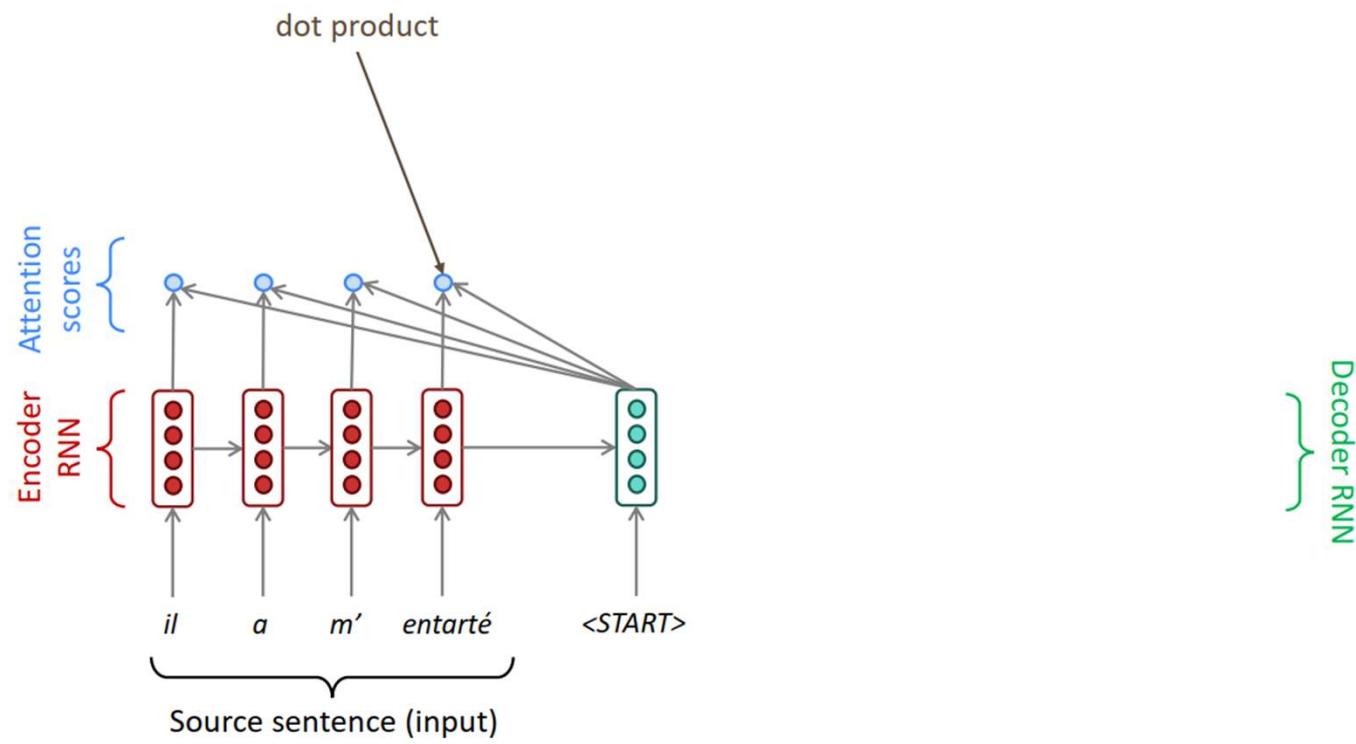
# Seq2seq with attention



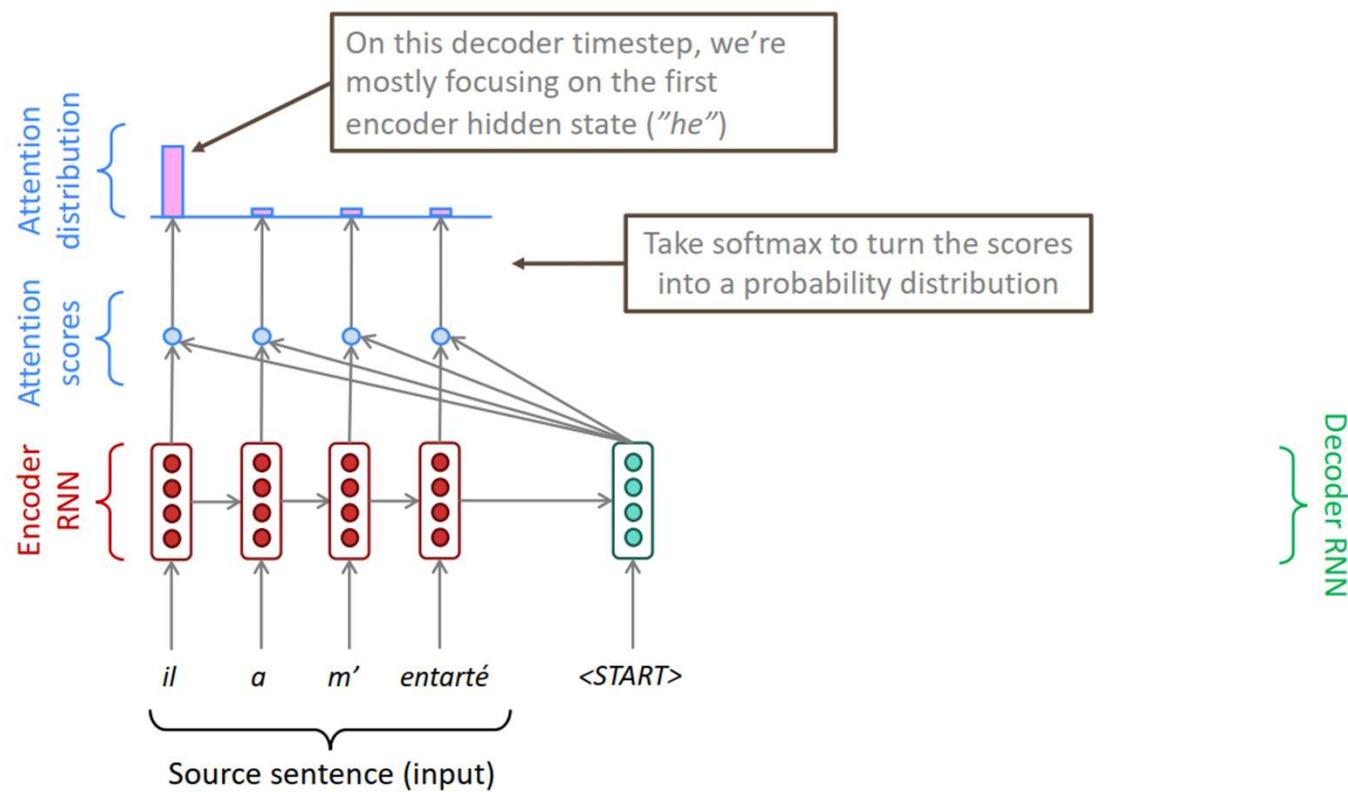
# Seq2seq with attention



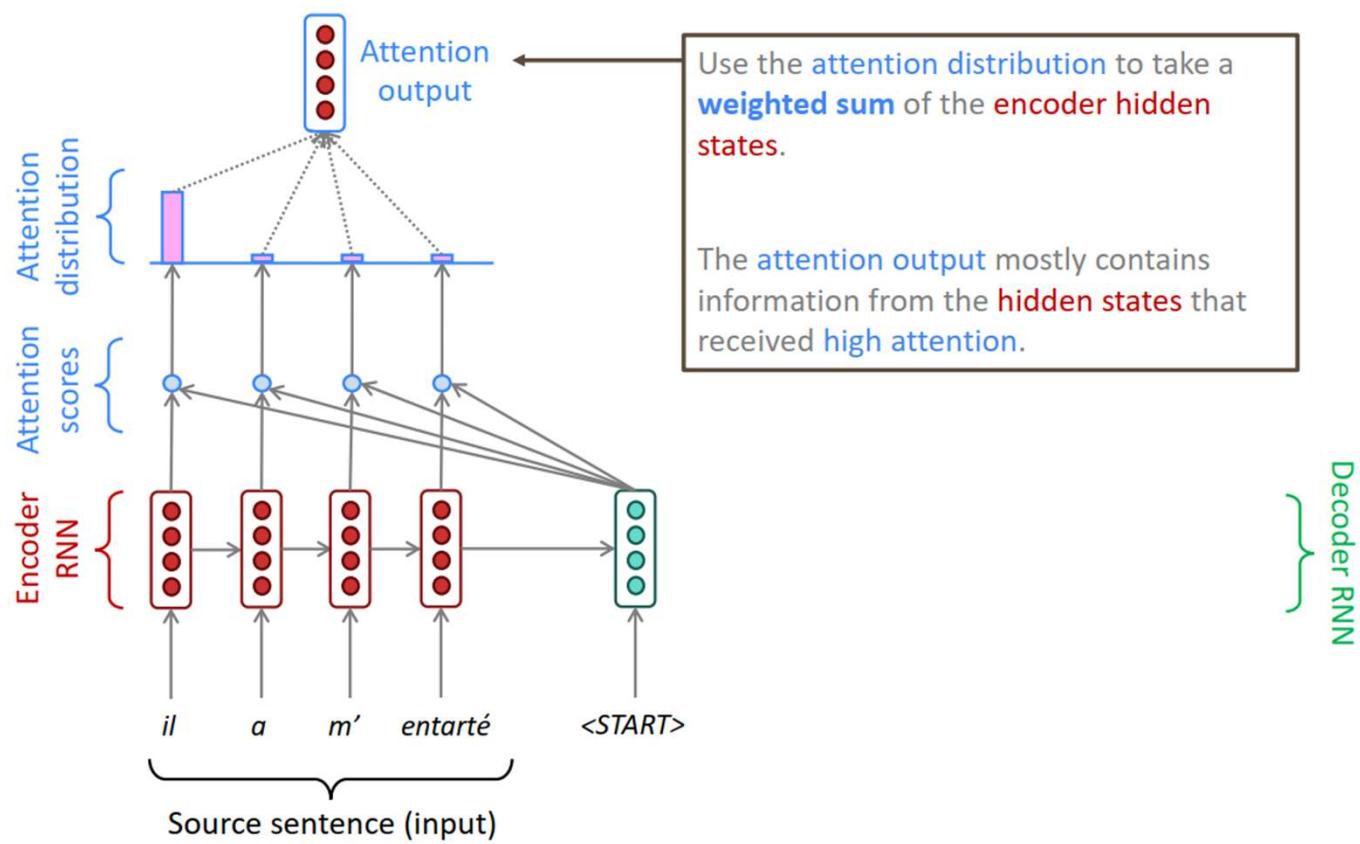
# Seq2seq with attention



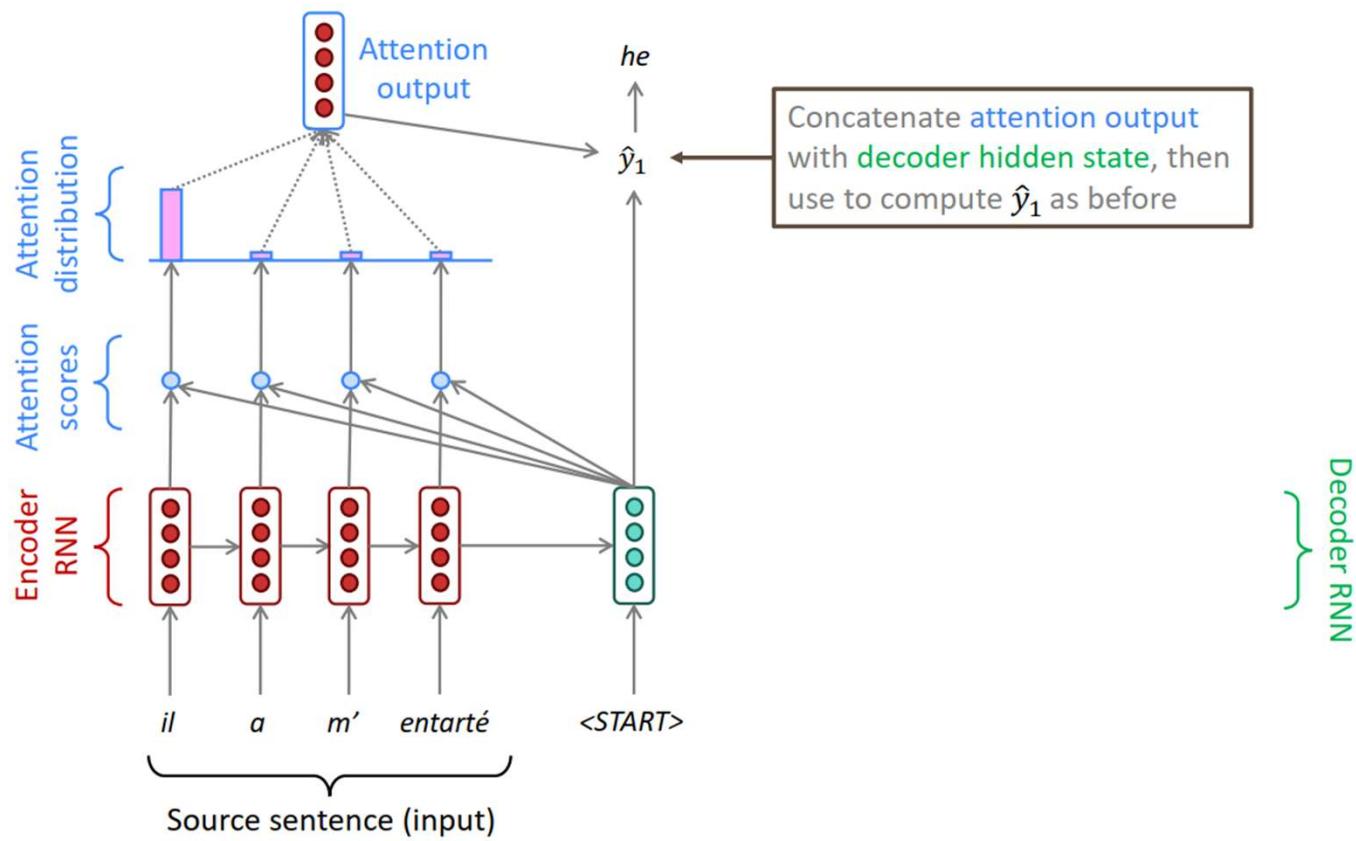
# Seq2seq with attention



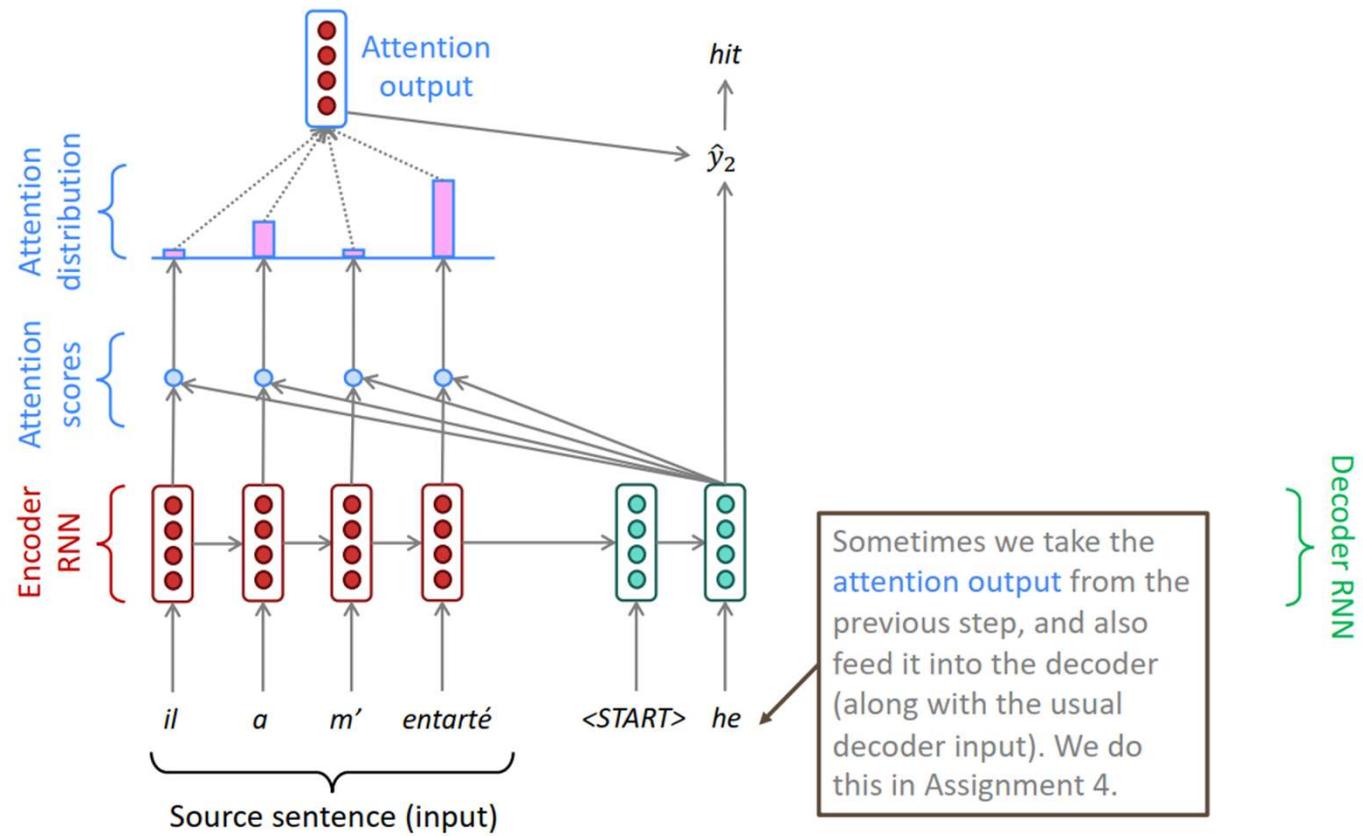
# Seq2seq with attention



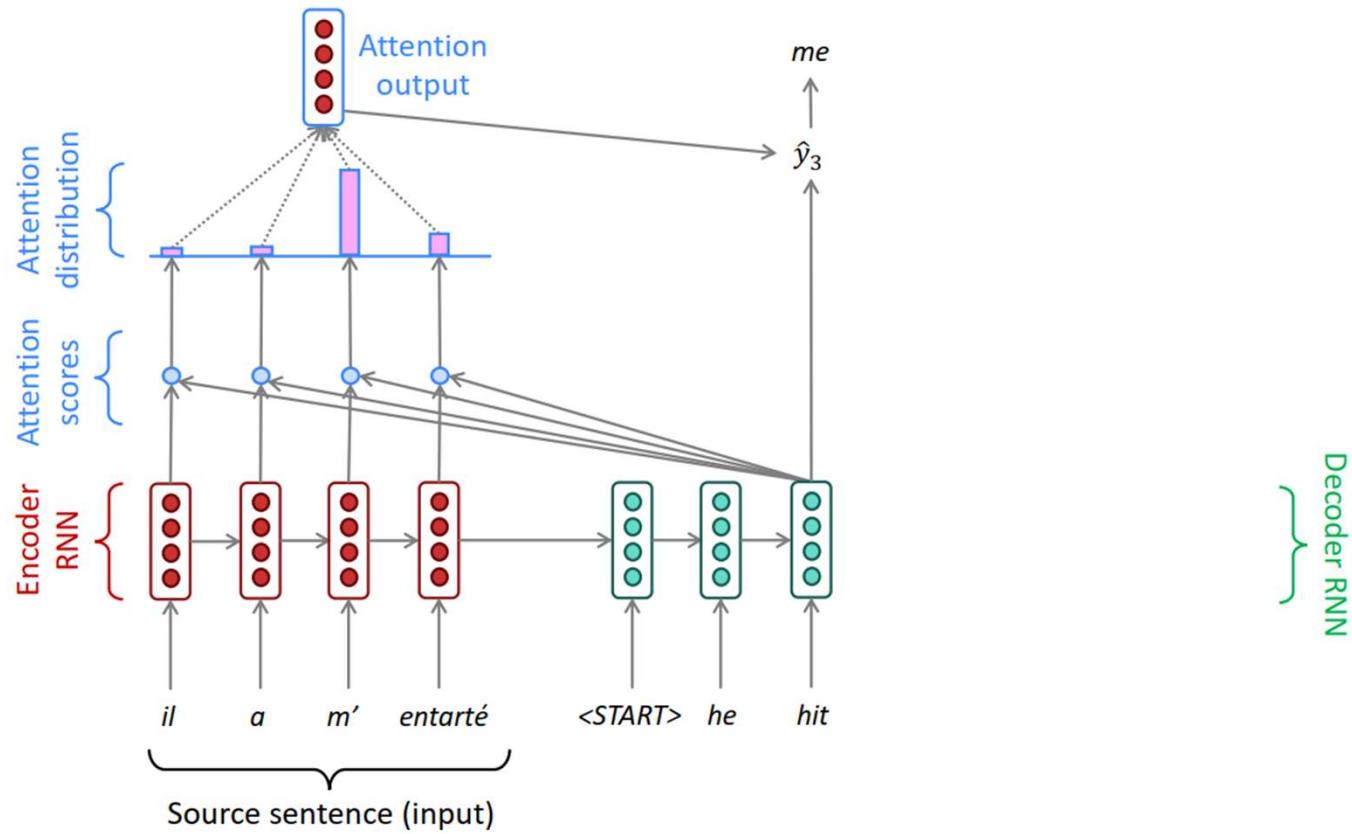
# Seq2seq with attention



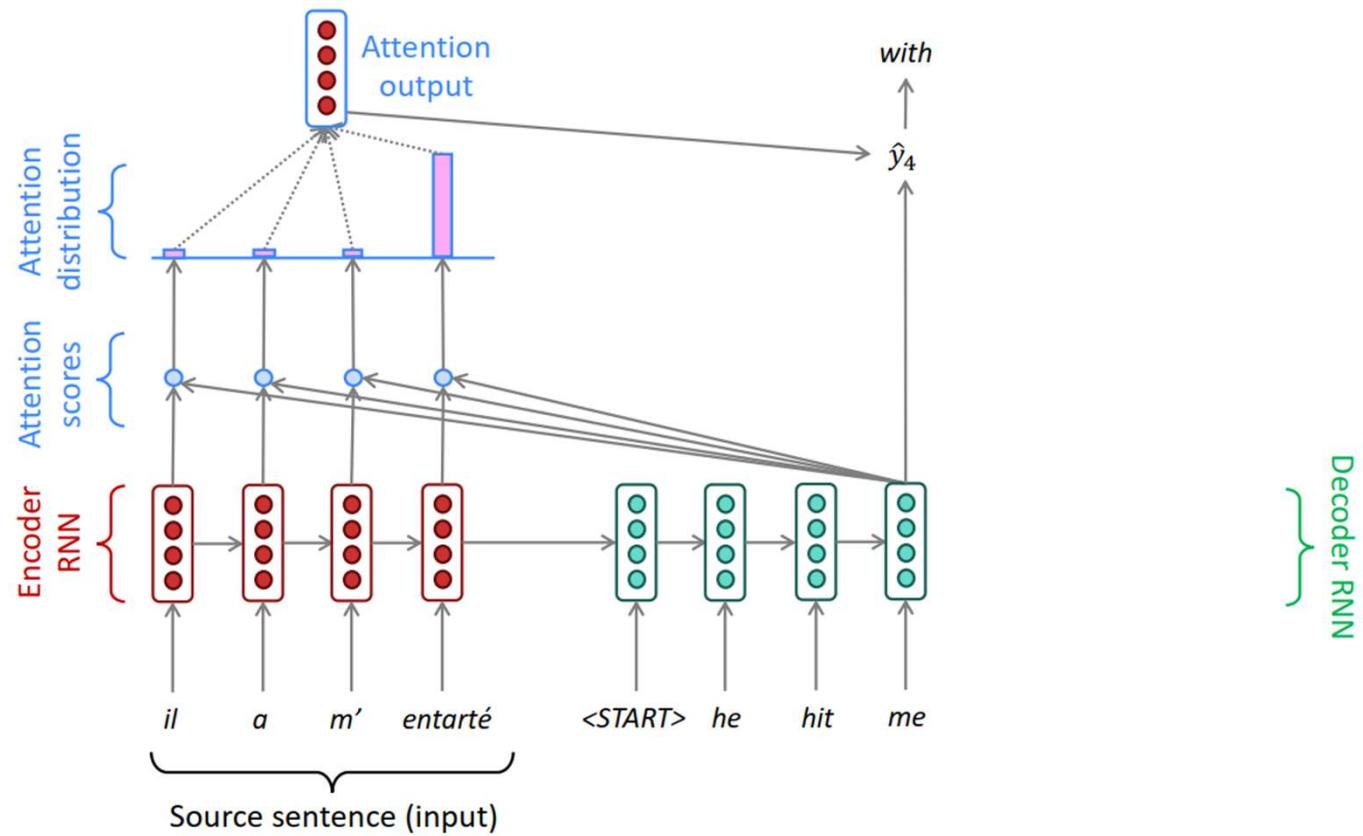
# Seq2seq with attention



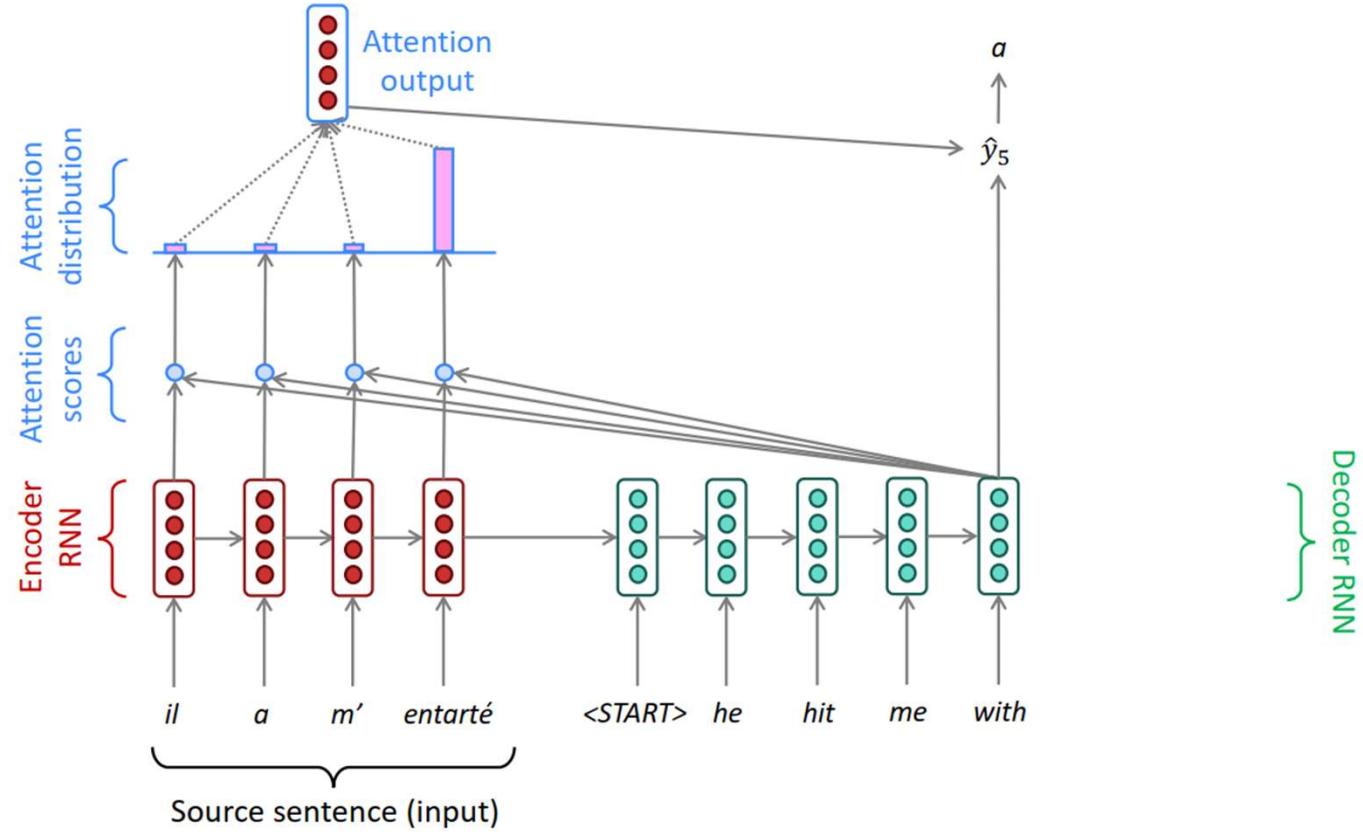
# Seq2seq with attention



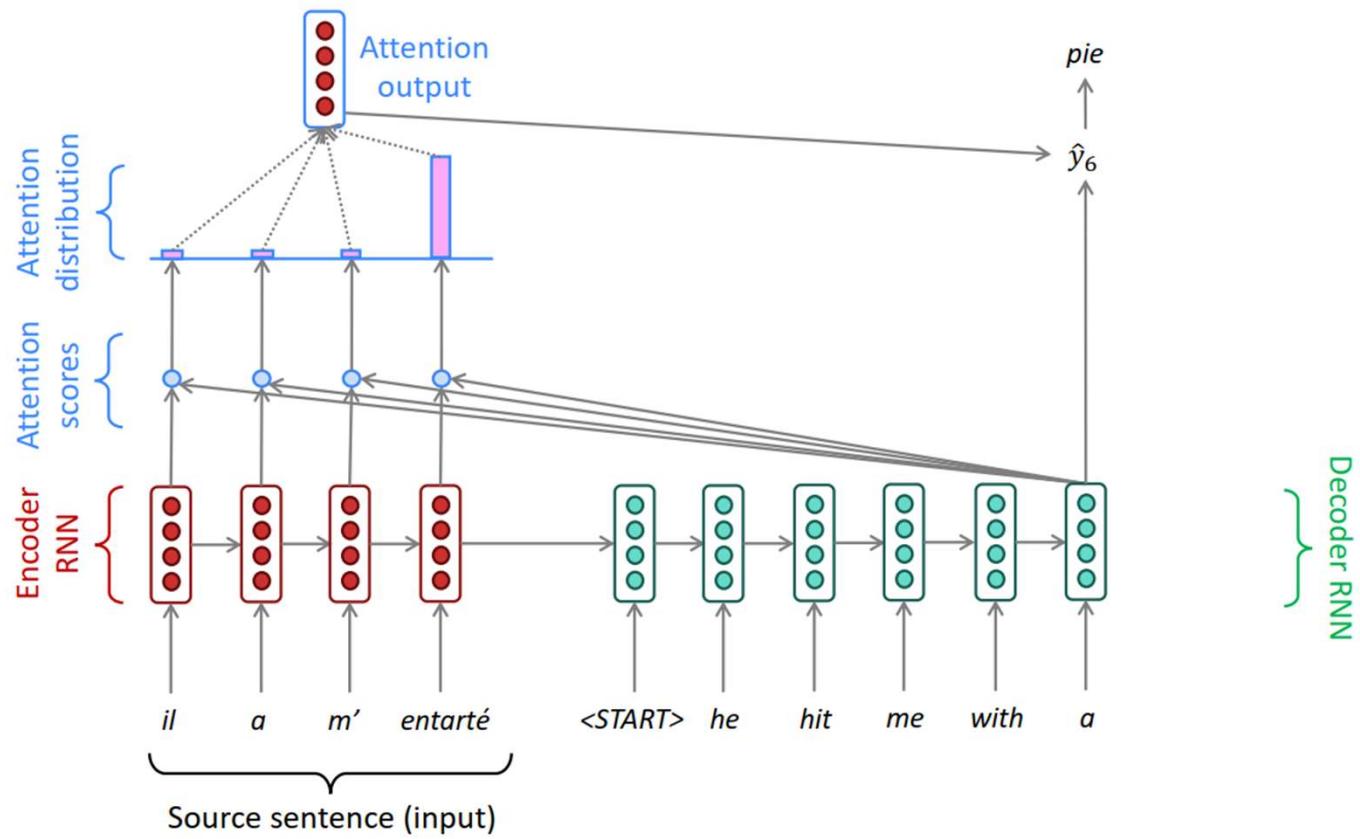
# Seq2seq with attention



# Seq2seq with attention



# Seq2seq with attention



# Attention in equations

- We have encoder hidden states  $h_1, \dots, h_N \in \mathbb{R}^h$
- On timestep t, we have decoder hidden state  $s_t \in \mathbb{R}^h$
- We get the attention scores  $e^t$  for this step:
- We take softmax to get the attention distribution  $\alpha^t$  for this step (this is a probability distribution)
- We use  $\alpha^t$  to take a weighted sum of the encoder hidden states to get the attention output  $a^t$
- Finally we concatenate the attention output with the decoder hidden state and proceed as in the non-attention seq2seq model

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

$$\alpha^t = \text{softmax}(e^t) \in \mathbb{R}^h$$

$$a^t = \sum_{i=1}^N \alpha_i^t h_i \in \mathbb{R}^h$$

$$[a_t; s_t] \in \mathbb{R}^h$$

# Attention

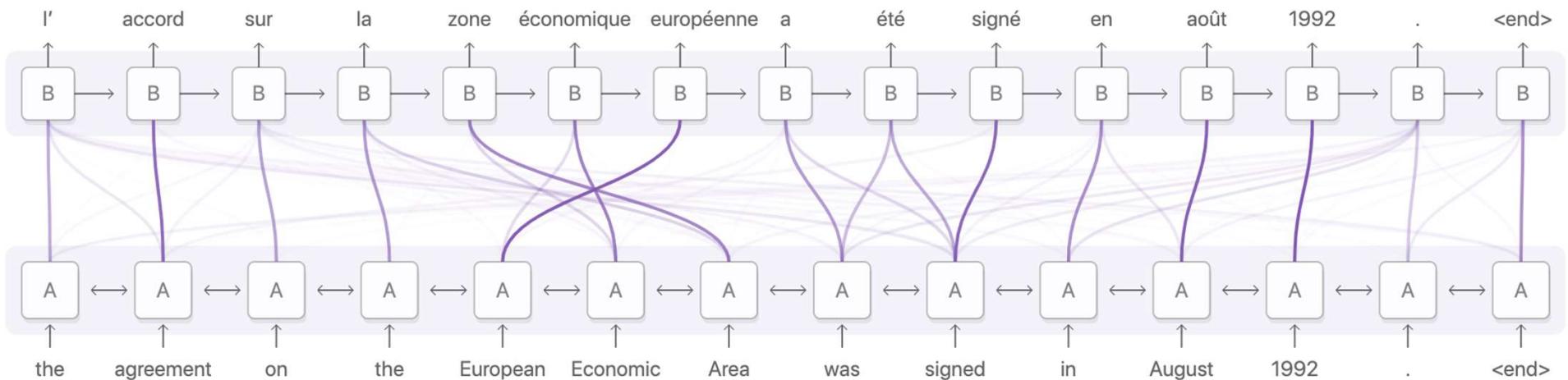
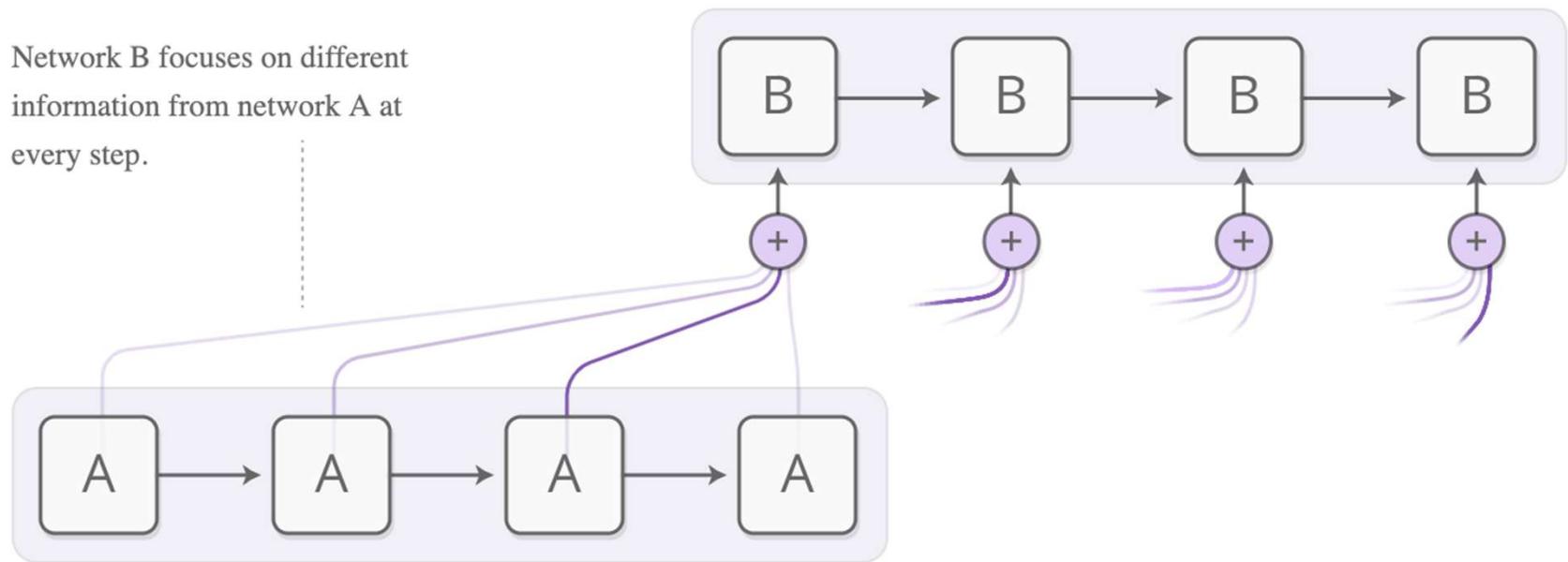


Diagram derived from Fig. 3 of [Bahdanau, et al. 2014](#)

# Attention

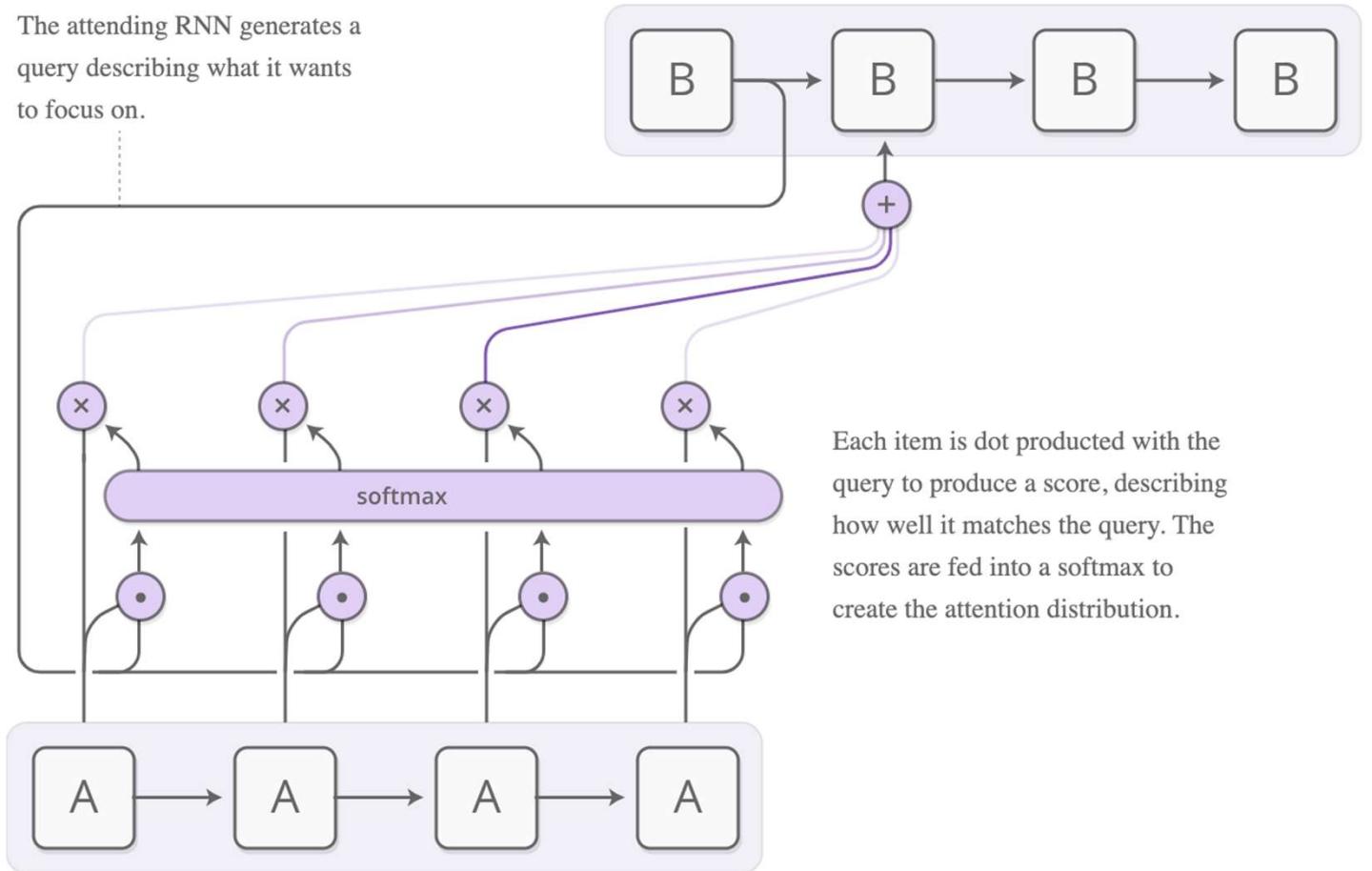
Idea: while generation focus on part of the information given

Network B focuses on different information from network A at every step.

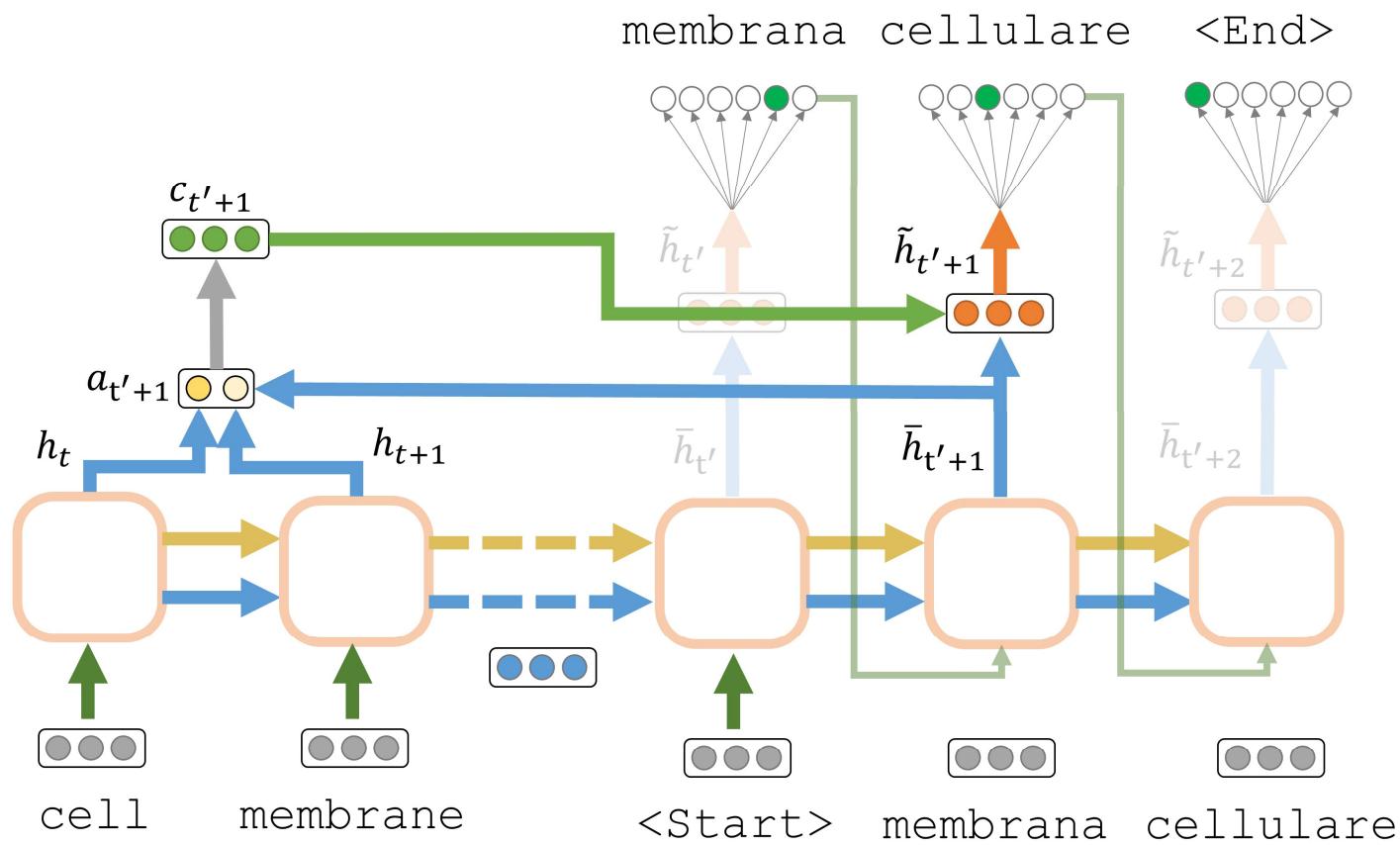


# Attention

The attending RNN generates a query describing what it wants to focus on.



# Attention

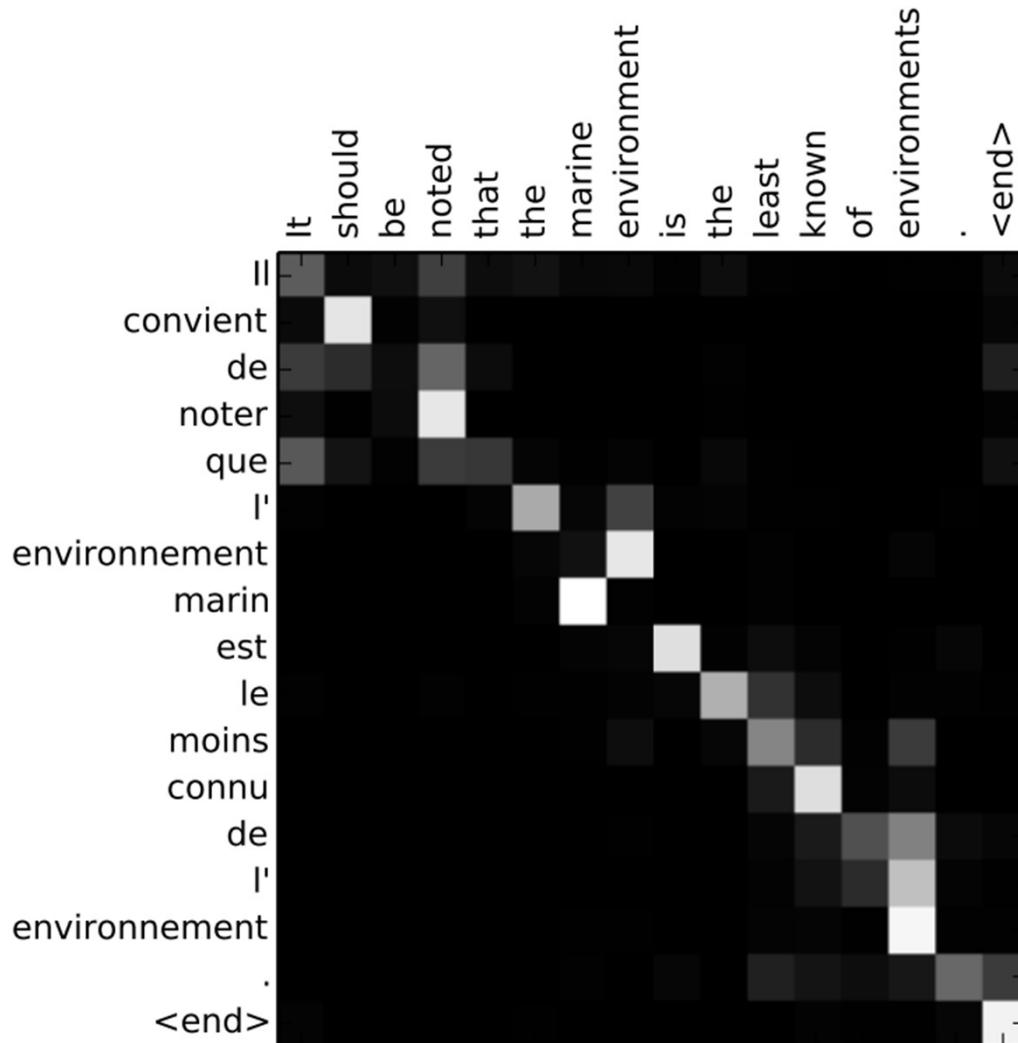


# Attention: advantages

- Attention significantly improves NMT performance
  - It's very useful to allow decoder to focus on certain parts of the source
- Attention provides more “human-like” model of the MT process
  - You can look back at 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
  - Provides shortcut to faraway states
- Attention provides some interpretability
  - By inspecting attention distribution, we see what the decoder was focusing on
  - We get (soft) alignment for free!
  - This is cool because we never explicitly trained an alignment system
  - The network just learned alignment by itself

# Attention

Inherent alignment



Dzmitry Bahdanau et al, 2014

# Attention variants

- Basic dot-product attention:  $e_i = s^T h_i \in \mathbb{R}$ 
  - Note: this assumes  $d_1 = d_2$ . This is the version we saw earlier.
- Multiplicative attention:  $e_i = s^T Wh_i \in \mathbb{R}$ 
  - Where  $W \in \mathbb{R}^{d_1 \times d_2}$  is a weight matrix. Perhaps better called “bilinear attention”
- Reduced-rank multiplicative attention:  $e_i = s^T (U^T V) h_i = (U_s)^T (V h_i)$ 
  - For low rank matrices  $U \in \mathbb{R}^{k \times d_2}$ ,  $V \in \mathbb{R}^{k \times d_1}$ ,  $k \ll d_1, d_2$
- Additive attention:  $e_i = v^T \tanh(W_1 h_i + W_2 s) \in \mathbb{R}$ 
  - Where  $W_1 \in \mathbb{R}^{d_3 \times d_1}$  and  $W_2 \in \mathbb{R}^{d_3 \times d_2}$  are weight matrices and  $v \in \mathbb{R}^{d_3}$  is a weight vector.
  - $d_3$  (the attention dimensionality) is a hyperparameter
  - “Additive” is a weird/bad name. It’s really using a feed-forward neural net layer.

# Attention is a general Deep Learning technique

- However: You can use attention in many architectures
  - (not just seq2seq) and many tasks (not just Seq2seq)
- 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
- We sometimes say that the query *attends* to the values.
- For example, in the seq2seq + attention model, each decoder hidden state (query) attends to all the encoder hidden states (values)

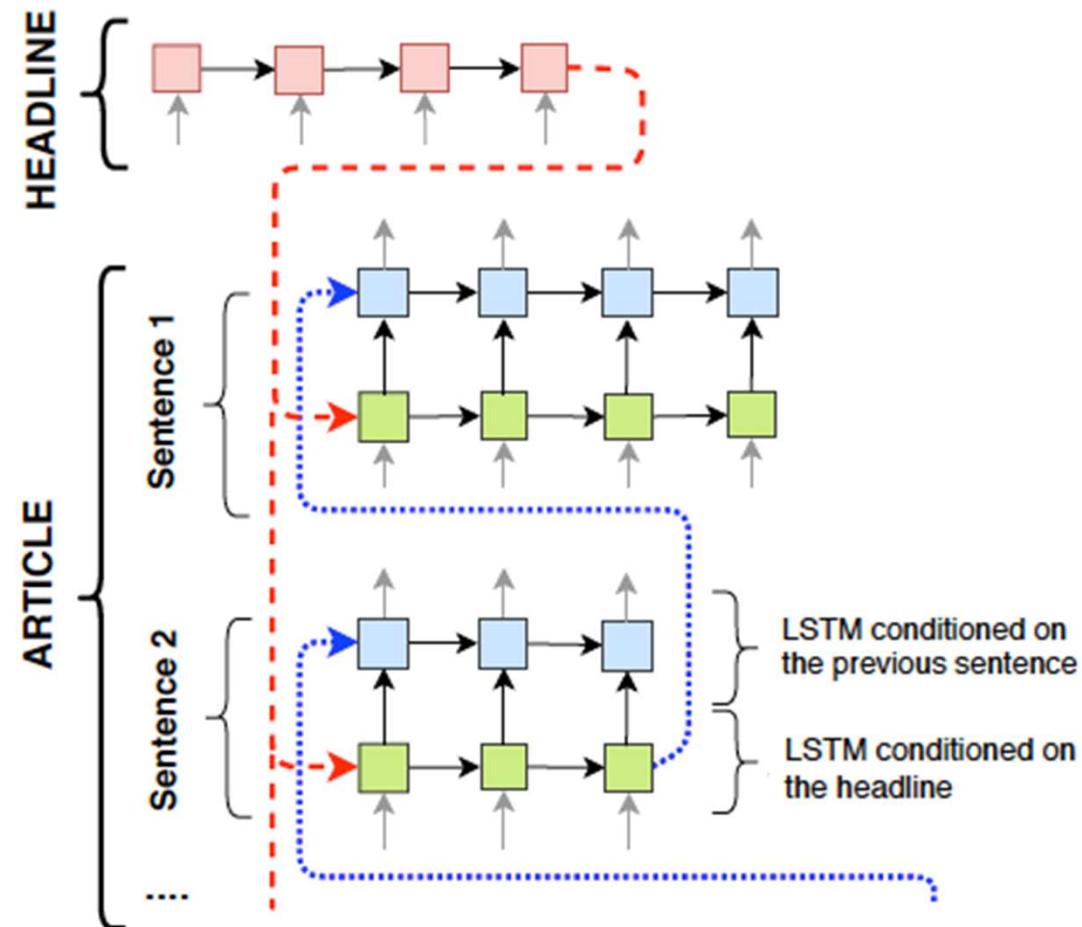
# Attention is a general Deep Learning technique

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!

# Attention

Veracity detection



# Attention



A stop sign is on a road with a mountain in the background.



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.

Show, attend, and tell, 2015.

Questions?

# RNNs: sequence bucketing

Too much padding can result in reduced performance

Place sequences with similar lengths in the same batch

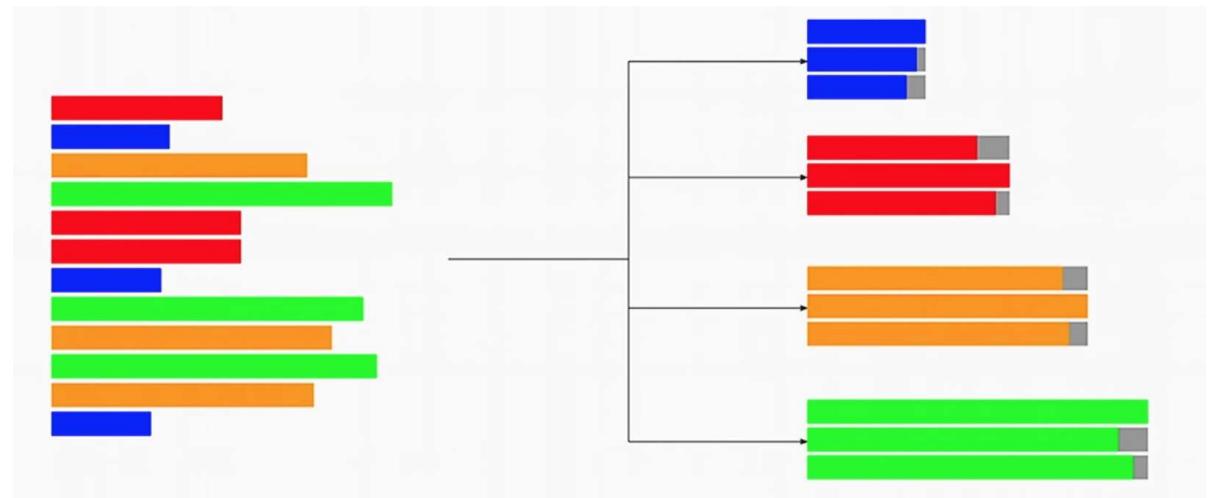


Image [source](#)