

Deep Learning

Topics: Recurrent Neural Networks for Sequential Data

Geoffroy Peeters

LTCI, Télécom Paris, IP Paris



1. Three main types of Nets

1.1 Multi Layers Perceptron (MLP), Fully-Connected (FC), Feed-Forward

1.2 Convolutional Neural Networks (CNN)

1.3 Recurrent Neural Networks (RNN)

2. Recurrent Neural Networks for Sequential Data

2.1 What are sequential data ?

2.2 Notations

2.3 How can we process Sequential data with a Neural Network

2.4 What is an RNN ?

2.5 Forward Propagation

3. Various types of sequential data

3.1 Many-To-Many $T_y = T_x$

3.2 Many-To-One $T_x > 1, T_y = 1$

3.3 One-To-Many $T_x = 1, T_y > 1$

3.4 Many-To-Many $T_y \neq T_x$

4. Different architectures

4.1 RNN

4.2 Bi-directional RNN

4.3 Deep RNN

4.4 Using 1D Convolution instead of RNN

5. Back Propagation Through Time (BPTT)

5.1 Forward pass

5.2 Backward pass : 1) compute $\frac{\partial L}{\partial W_{ya}}$

5.3 Backward pass : 2) compute $\frac{\partial \mathcal{L}}{\partial W_{aa}}$

5.4 Backward pass : 3) compute $\frac{\partial \mathcal{L}}{\partial W_{ax}}$

5.5 Vanishing and exploding gradients

5.6 Dealing with the exploding gradient problem

5.7 Dealing with the vanishing gradient problem

6. Gated units (LSTM and GRU)

6.1 Main ideas : Recurrent Neural Network (RNN)

6.2 Main ideas : Long Short Term Memory Units (LSTM)

6.3 Main ideas : Gated Recurrent Unit (GRU)

6.4 Recurrent Neural Network (RNN)

6.5 Long Short Term Memory Units (LSTM)

6.6 Gated Recurrent Unit (GRU)

6.7 LSTM Example usage

7. Natural Language Processing

7.1 Notations

7.2 Language model

7.3 One-hot-encoding

7.4 Word embedding

7.5 Application : Sentiment Classification

7.6 Application : Handwritten Character Recognition

8. Sequence to sequence

8.1 Introduction

8.2 Machine Translation

8.3 Attention model

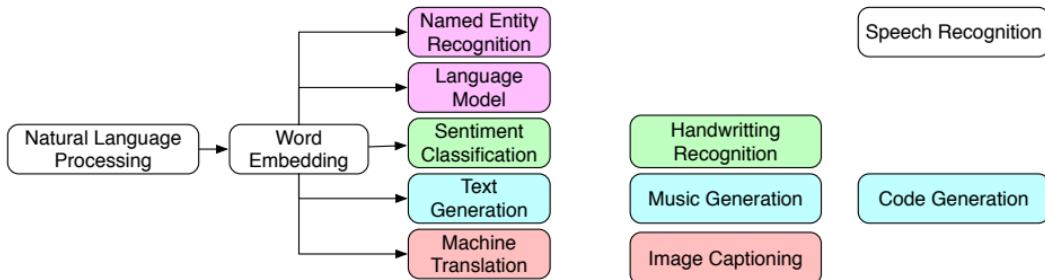
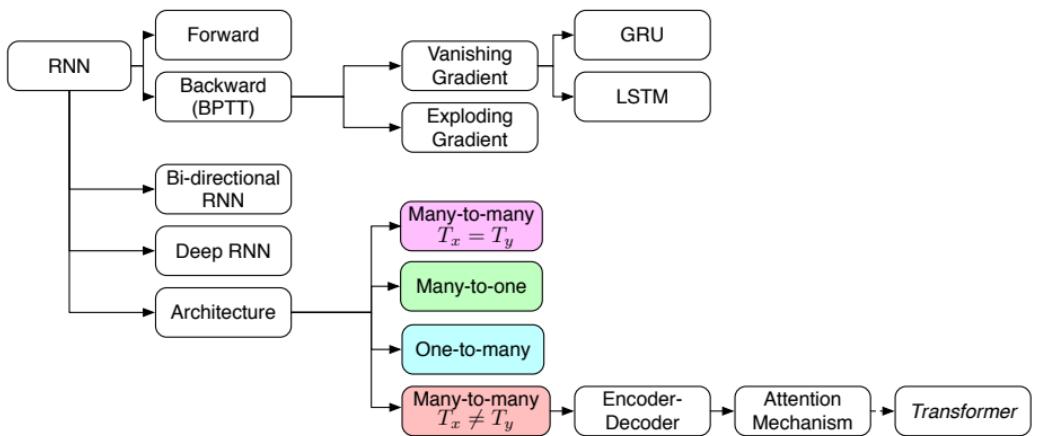
8.4 Application : Image captioning

8.5 Transformer

8.6 BERT

9. Implementations in keras

Overview



What are sequential data ?

- **Sequential data are any type of data for which the order matters**

- Examples :

- **Time series :**

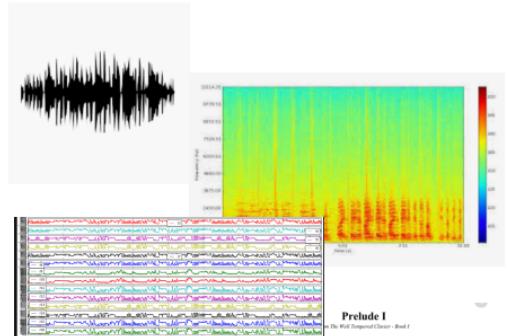
- Audio (sequence of sound pressures or sequence of Fourier transforms)
 - Sequence of musical notes (C)
 - Video (sequence of images)
 - Sensors (heart-beat, plane altitude)
 - Stock market

- **Ordered :** (but not with time)

- Text/Words (C)
 - Genes/ADN (C)

- Common models :

- Vector linear autoregression (VAR)
 - Markov model (Discrete-State HMMs)
 - Kalman filters (Continuous-State HMMs)
 - ...



ROMEO & JULIETTE

WILLIAM SHAKESPEARE

Juliette : O Roméo ! Roméo ! Pourquoi es-tu Roméo ?
Réire ton père et abdiquer ton nom, ou, si tu ne le veux pas, jure de t'arrêter, et je ne semerai plus une Capulet.

Roméo, à part : Dois-je t'écouter encore ou lui répondre ?

Juliette : Ton nom est mon ennemi. Tu n'es pas un Montage, tu es loi-même. Qu'est-ce qu'un Montage ? Ce n'est ni une main, ni un pied, ni un bras, ni un visage, ni un cœur, ni une tête, ni rien du tout... Oh ! Sois mon autre nom. Que y a-t-il dans un nom ? Ce que nous appelons une rose embrassera aussi bien un autre nom. Ainsi, quand Roméo ne s'appelleraient plus Roméo, il conserverait encore les choses précieuses qu'il



Notations

- Inputs :
 - $\underline{x}^{(t)}$: input vectors \underline{x} at time t (of dimension $n^{[0]}$)
 - $\{\underline{x}^{(1)}, \underline{x}^{(2)}, \dots, \underline{x}^{(T_x)}\}$ the sequence of inputs of length T_x
- Outputs :
 - $\underline{y}^{(t)}$: output vectors \underline{y} at time t
 - $\{\underline{y}^{(1)}, \underline{y}^{(2)}, \dots, \underline{y}^{(T_y)}\}$ the sequence of outputs of length T_y
- T_x and T_y can be of different lengths
- $x^{(i)<t>}$ is the i^{th} example in the training-set
 - $T_x^{(i)}$ the length of the i^{th} input sequence
 - $T_y^{(i)}$ the length of the i^{th} output sequence

Notations

- Different types of output \underline{y} :
 - Continuous values $\underline{y} \in \mathbb{R}^{n_y}$ (Activation : linear, Loss : MSE)
 - Binary-classes $\underline{y} \in \{0, 1\}$ (Activation : sigmoid, Loss : Binary-Cross-Entropy)
 - Multi-classes $y \in \{1 \dots K\} \Rightarrow \underline{y} \in \{0, 1\}^K$ (Activation : softmax, Loss : Cross-Entropy)
- Different types of input \underline{x}
 - Continuous values $\underline{x} \in \mathbb{R}^{n_x}$
 - Categorical values $x \in \{1 \dots K\} \Rightarrow \underline{x} \in \{0, 1\}^K$
- **Example** : categorical inputs/ binary outputs
 - "Named Entity Recognition"

{x}	$x^{<1>}$	$x^{<2>}$...				$x^{<7>}$
	Emmanuel	Macron	is	the	president	of	France
{y}	$y^{<1>}$	$y^{<2>}$...				$y^{<7>}$
	1	1	0	0	0	0	1

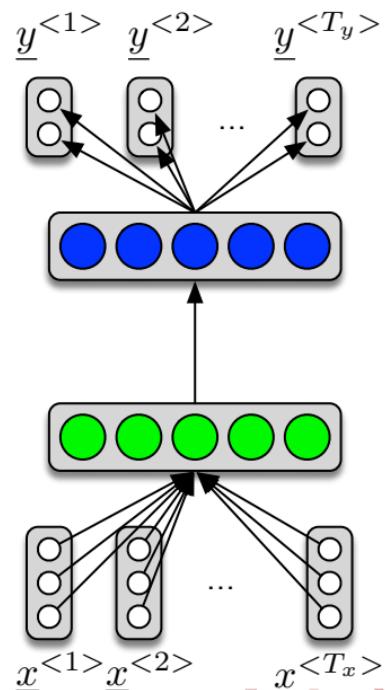
- How do we process categorical \circledC values as input ?

- One-hot-encoding
- Embedding

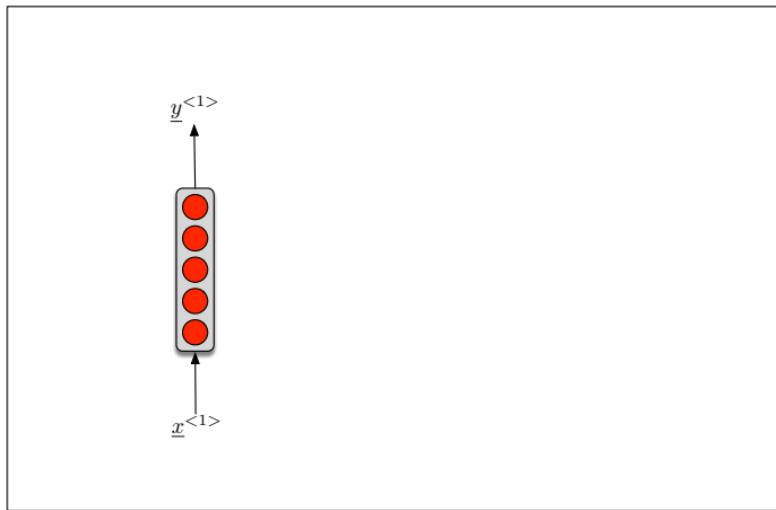
How can we process Sequential data with a Neural Network

- **Why Multi-Layer-Perceptron/Fully-connected networks are not appropriate ?**

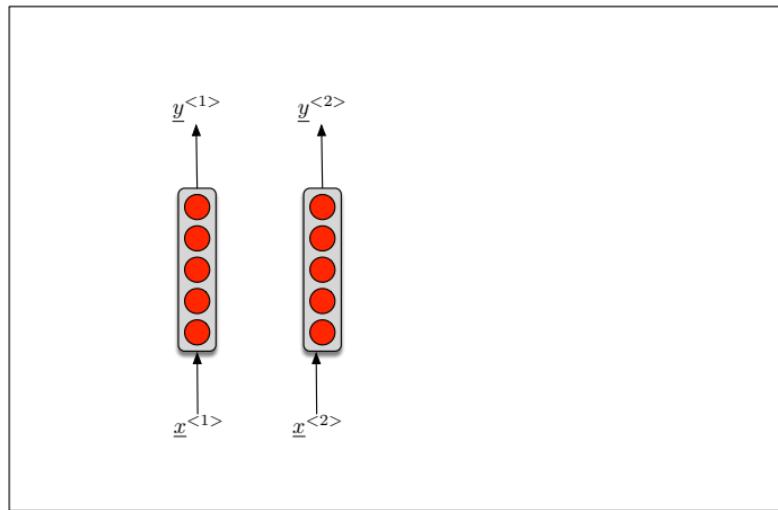
- We could represent an input sequence $\{\underline{x}^{<1>} \cdots \underline{x}^{<T_x>}\}$ as a large-vector of dimension $(n^{[0]} T_x)$
- but (because $T_x^{(i)}$ can be different for each input sequence) the dimension of this large-vector would change for each
 - zero-padding ?
- still require a huge input dimension ($n^{[0]} T_x$)
- the network will have to learn different weights for the same dimension $n^{[0]}$ at various time $\langle t \rangle$ in the sequence
 - no weight sharing !
- **Solutions :**
 - Recurrent Neural Network
 - 1D Convolutional Neural Network (convolution only over time)



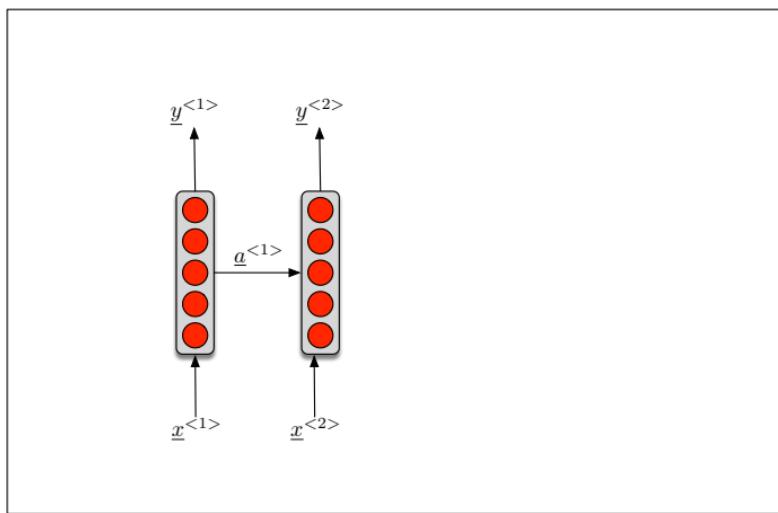
What is an RNN ?



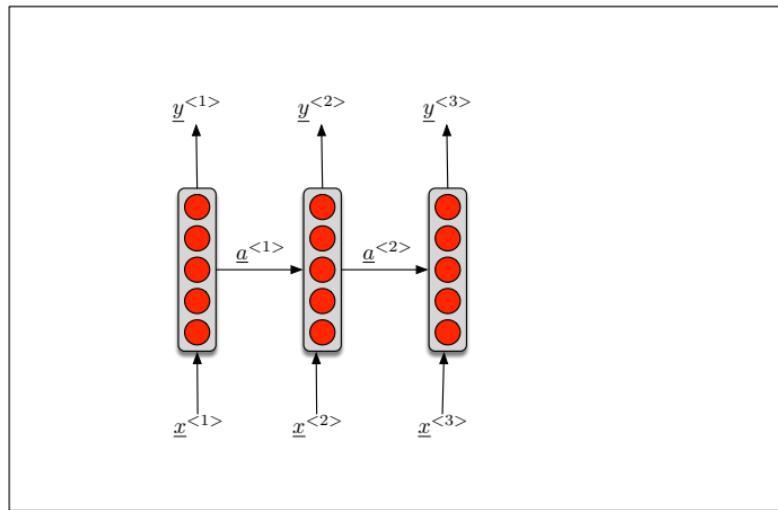
What is an RNN ?



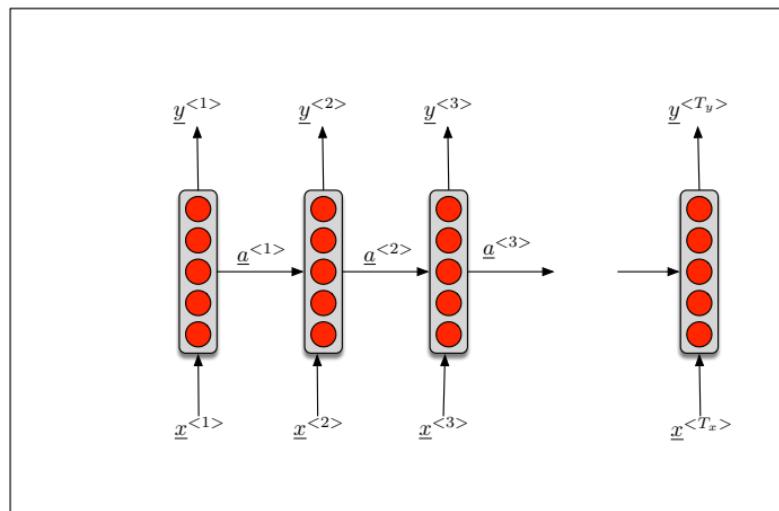
What is an RNN ?



What is an RNN ?

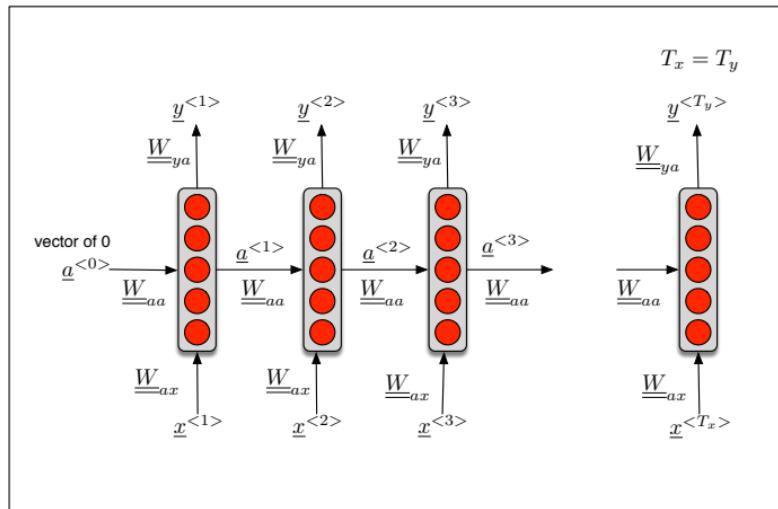


What is an RNN ?



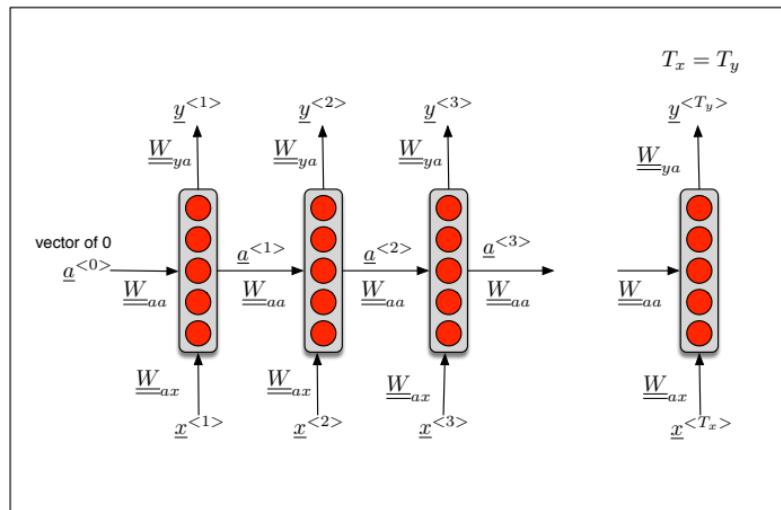
- At each time step $< t >$, the RNN passes its activation $a^{<t>}$ to the next time step $< t + 1 >$

What is an RNN ?



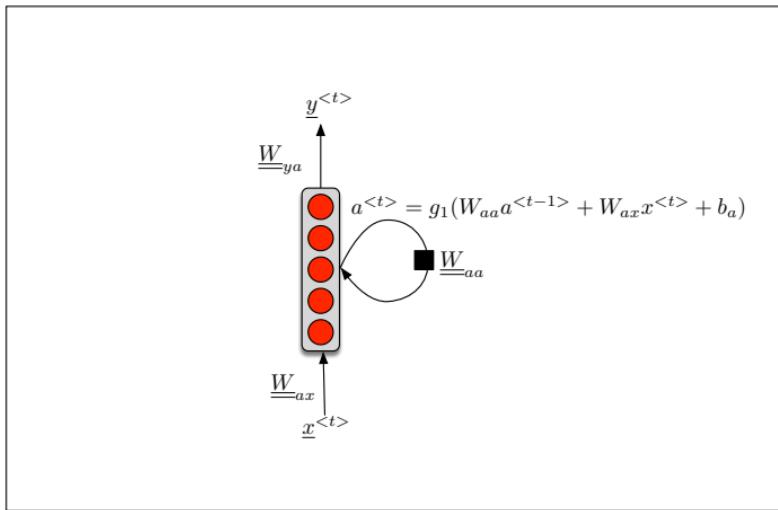
- Initialize the first activation $a^{<0>}$ as zero

What is an RNN ?



- \underline{W}_{ax} ($x^{<t>} \rightarrow a^{<t>}$) is the same for every time step
- \underline{W}_{aa} ($a^{<t-1>} \rightarrow a^{<t>}$) is the same for every time step
- \underline{W}_{ya} ($a^{<t>} \rightarrow y^{<t>}$) is the same for every time step

Other possible representation



- \underline{W}_{ax} ($x^{<t>} \rightarrow a^{<t>}$) is the same for every time step
- \underline{W}_{aa} ($a^{<t-1>} \rightarrow a^{<t>}$) is the same for every time step
- \underline{W}_{ya} ($a^{<t>} \rightarrow y^{<t>}$) is the same for every time step

Forward Propagation

- Forward propagation :

$$a^{<0>} = 0$$

$$a^{<1>} = g_1 \left(W_{aa} a^{<0>} + W_{ax} x^{<1>} + b_a \right)$$

$$\hat{y}^{<1>} = g_2 \left(W_{ya} a^{<1>} + b_y \right)$$

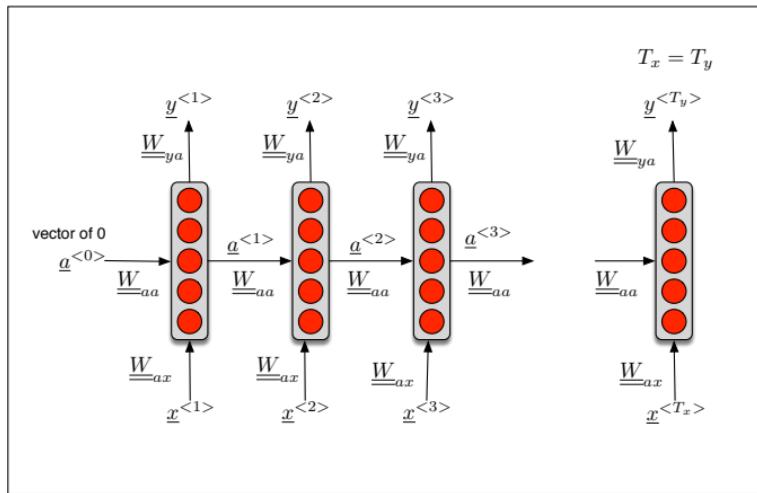
...

$$a^{<t>} = g_1 \left(W_{aa} a^{<t-1>} + W_{ax} x^{<t>} + b_a \right)$$

$$\hat{y}^{<t>} = g_2 \left(W_{ya} a^{<t>} + b_y \right)$$

- Compact notations :

$$\begin{bmatrix} W_{aa} & | & W_{ax} \\ (n_a, n_a) & & (n_a, n_x) \end{bmatrix} \begin{bmatrix} a^{<t-1>} \\ (n_a, m) \\ X^{<t>} \\ (n_x, m) \end{bmatrix}$$



$$a^{<t>} = g_1 \left(W_a [a^{<t-1>} ; x^{<t>}] + b_a \right)$$

$$\hat{y}^{<t>} = g_2 \left(W_y a^{<t>} + b_y \right)$$

Various types of sequential data

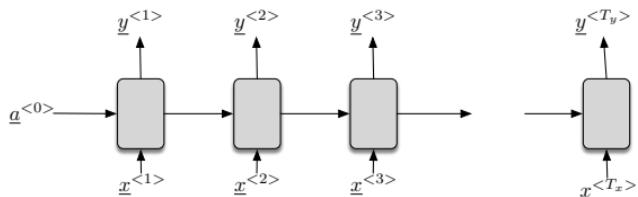
Various types of sequential data/ Many-To-Many $T_y = T_x$

Input x	Output y	Type	Examples
sequence $T_x > 1$	sequence $T_y = T_x$	Many-To-Many	Named entity recognition

In 1917, Einstein applied the general theory of relativity to model the large-scale structure of the universe. He was visiting the United States when Adolf Hitler came to power in 1933 and did not go back to Germany, where he had been a professor at the Berlin Academy of Sciences. He settled in the U.S., becoming an American citizen in 1940. On the eve of World War II, he endorsed a letter to President Franklin D. Roosevelt alerting him to the potential development of "extremely powerful bombs of a new type" and recommending that the U.S. begin similar research. This eventually led to what would become the Manhattan Project. Einstein supported defending the Allied forces, but largely denounced using the new discovery of nuclear fission as a weapon. Later, with the British philosopher Bertrand Russell, Einstein signed the Russell-Einstein Manifesto, which highlighted the danger of nuclear weapons. Einstein was affiliated with the Institute for Advanced Study in Princeton, New Jersey, until his death in 1955.

Tag colours:

LOCATION TIME PERSON ORGANIZATION MONEY PERCENT DATE

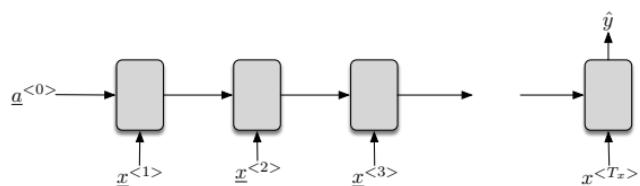
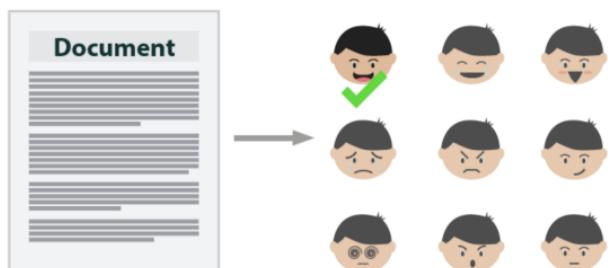


- **Many-to-many** : $T_x = T_y$

Various types of sequential data

Various types of sequential data/ Many-To-One $T_x > 1, T_y = 1$

Input x	Output y	Type	Examples
sequence $T_x > 1$	single $T_y = 1$	Many-To-One	Sentiment analysis, Video activity detection



- **Many-to-one** : sentiment classification
 - x = text (movie review)
 - y = 0/1 or 1 ⋯ 5 rating

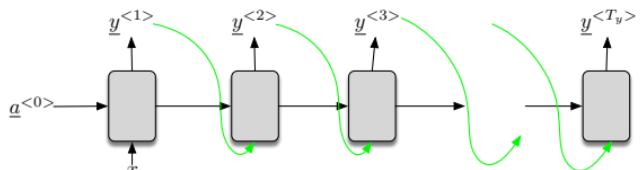


Various types of sequential data

Various types of sequential data/ One-To-Many $T_x = 1, T_y > 1$

Input x	Output y	Type	Examples
single $T_x = 1$	sequence, $T_y > 1$	One-To-Many	Text generation, music generation

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 slot 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 the [[Punjab Resolution]] (PJS) [<http://www.humanah.yahoo.com/guardian.cfm/7754800786d17551963a89.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 tripod of aid exile.]]



- **One-to-many :** text, code, music generation

AI-music, 192x96x576, 193 epochs, 805 seconds, 7 training sets

Russell E Glaue



Various types of sequential data

Various types of sequential data/ Many-To-Many $T_y \neq T_x$

Input x	Output y	Type	Examples
sequence $T_x > 1$	sequence $T_y > 1, T_y \neq T_x$	Many-To-Many	Automatic Speech Recognition, Machine translation

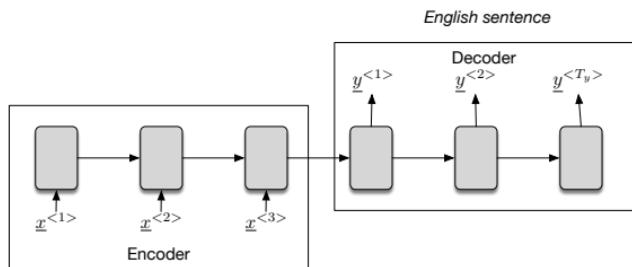
Traduire **anglais** (langue identifiée) ▾

Deep learning recently became popular in supervised machine learning.

Traduire en **français** ▾

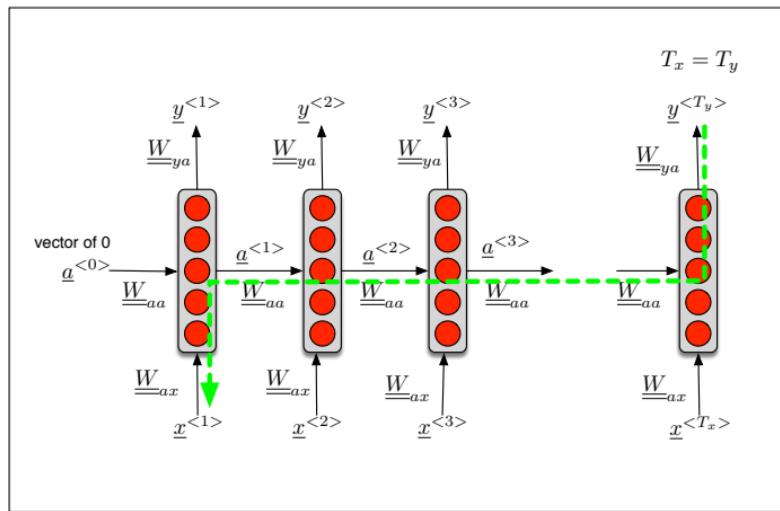
L'apprentissage en profondeur est récemment devenu populaire dans l'apprentissage machine supervisé.

Traduire le document



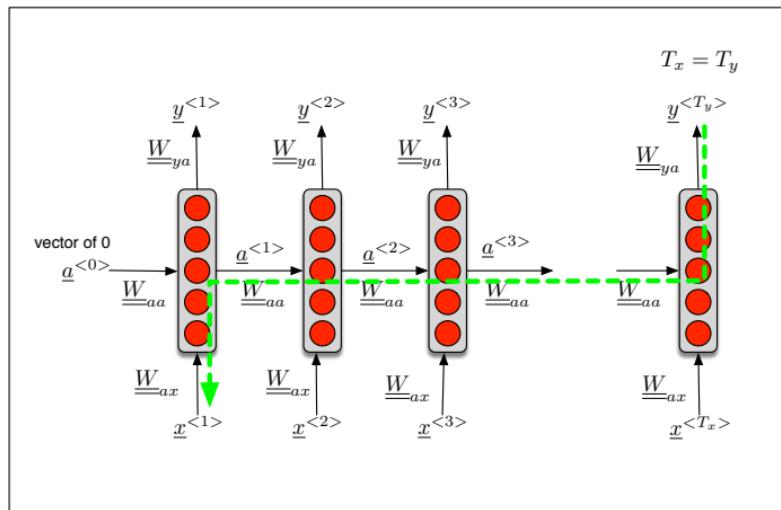
- **Many-to-many** : $T_x \neq T_y$: Machine translation
 - Encoder/Decoder
 - **Attention** mechanism
 - **Transformer** architecture

RNN



- RNN scans through the data from left to right \Rightarrow can do skip connections in the past
- Question : is "Apple" a named entity ?
 - "The new **iPhone** is sold in **Apple** stores."
 - "A refreshing and healthy **beverages** is an **Apple** juice."
- \Rightarrow we can predict (since RNN depends on the past)

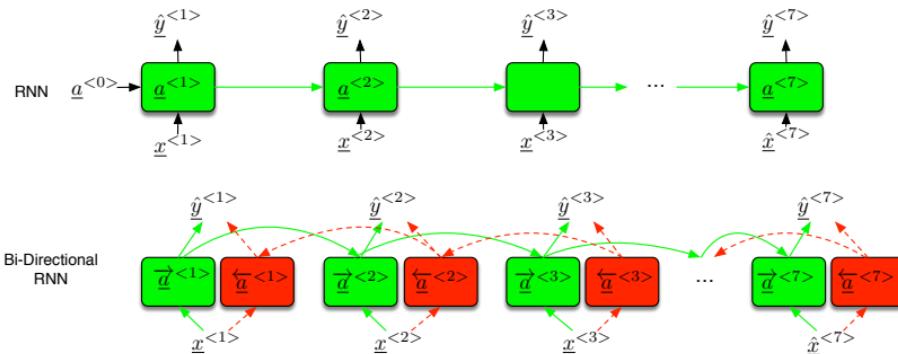
Bi-directional RNN



- RNN scans through the data from left to right \Rightarrow can do skip connections **only in the past**
- Question : is "Apple" a named entity ?
 - "Apple stores sale the new **iPhone.**"
 - "Apple juice is a refreshing and healthy **beverages.**"
- \Rightarrow we cannot predict (since RNN does not depend on the future)

Bi-directional RNN

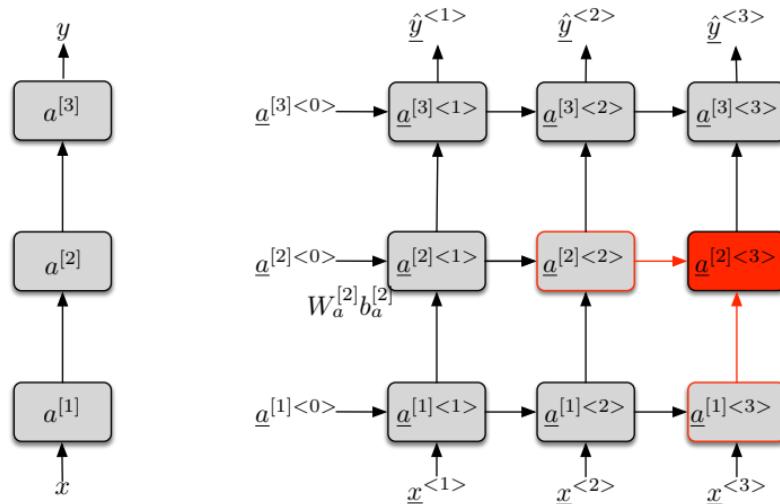
- Solution ?
 - Bi-directional RNN (BRNN)
- Question : is "Apple" a named entity ?
 - "Apple stores sale the new iPhone."
 - "Apple juice is a refreshing and healthy beverages."
- \Rightarrow we can predict since we now depend on the past and the future



- The prediction $\hat{y}^{<t>}$ is done from the concatenation of
 - the **left-to-right** $\vec{a}^{<t>}$ and
 - the **right-to-left** $\underline{a}^{<t>}$ hidden states of RNN

$$\hat{y}^{<t>} = g(W_y[\vec{a}^{<t>} \text{, } \underline{a}^{<t>}] + b_y)$$

Deep RNN

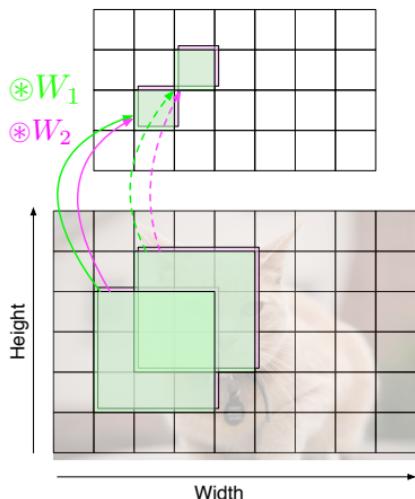


- For the first layer, the inputs of the cell at time $< t >$ are
 - the input x at the current time $< t >$: $x^{<t>}$
 - the value of the cell at the previous time $< t - 1 >$: $a^{<t-1>}$
- For the layer $[l]$, the inputs of the cell at time $< t >$ are
 - the value of the cell of the previous layer $[l - 1]$ at the current time $< t >$: $a^{[l-1]} < t >$
 - the value of the cell of the current layer $[l]$ at the previous time $< t - 1 >$: $a^{[l]} < t-1 >$

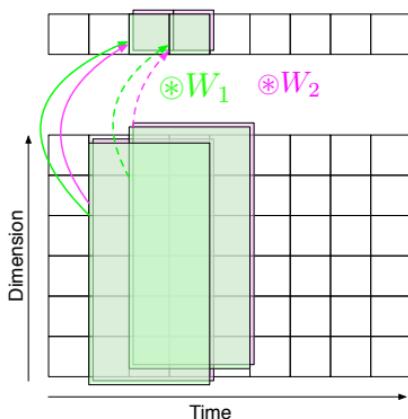
$$a^{[l]} < t > = g \left(W_a^{[l]} [a^{[l]} < t-1 >, a^{[l-1]} < t >] + b_a^{[l]} \right)$$

Using 1D Convolution instead of RNN

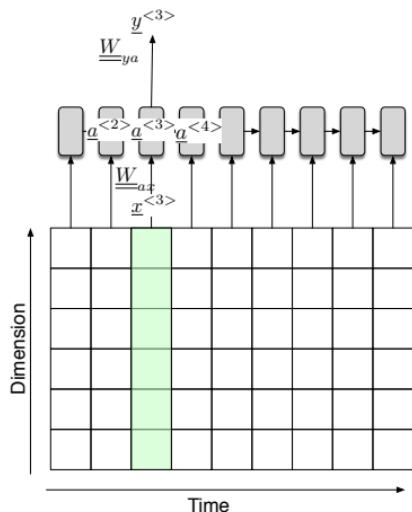
Conv2D



Conv1D



RNN



Back Propagation Through Time (BPTT)

Back Propagation Through Time (BPTT)

Forward pass

- Compute $a^{(t)}$, $\hat{y}^{(t)}$, $\mathcal{L}^{(t)}$, \mathcal{L}

$$a^{(t)} = g_a(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b_a)$$

$$\hat{y}^{(t)} = g_y(W_{ya}a^{(t)} + b_y)$$

- Compute the Loss

- Loss for time $< t >$: $\mathcal{L}^{(t)}(y^{(t)}, \hat{y}^{(t)})$

- $y^{(t)}$: true label

- $\hat{y}^{(t)}$: predicted label

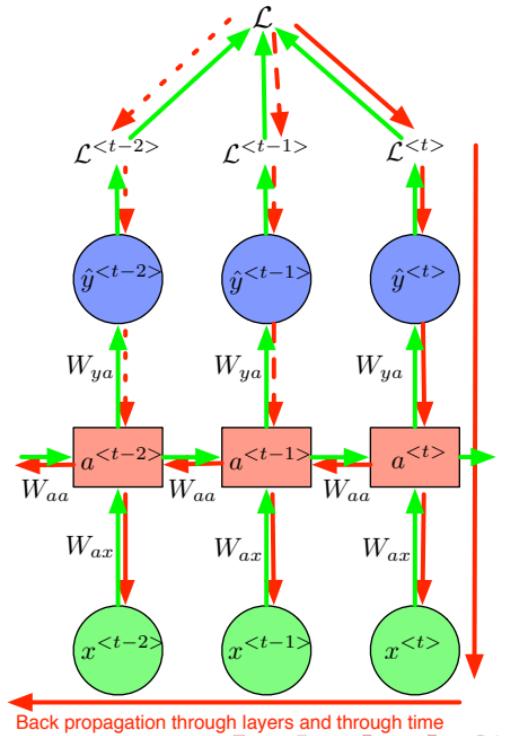
- Total loss :

$$\mathcal{L} = \sum_t \mathcal{L}^{(t)}(y^{(t)}, \hat{y}^{(t)})$$

- Backward pass

- compute $\frac{\partial \mathcal{L}}{\partial W_{ya}}$, $\frac{\partial \mathcal{L}}{\partial W_{ax}}$, $\frac{\partial \mathcal{L}}{\partial W_{aa}}$, $\frac{\partial \mathcal{L}}{\partial b_x}$, $\frac{\partial \mathcal{L}}{\partial b_a}$

- All weights are shared across time steps!!!



Back Propagation Through Time (BPTT)

Backward pass : 1) compute $\frac{\partial \mathcal{L}}{\partial W_{ya}}$

- We measure how much varying W_{ya} affect $\mathcal{L}^{(t)}$.

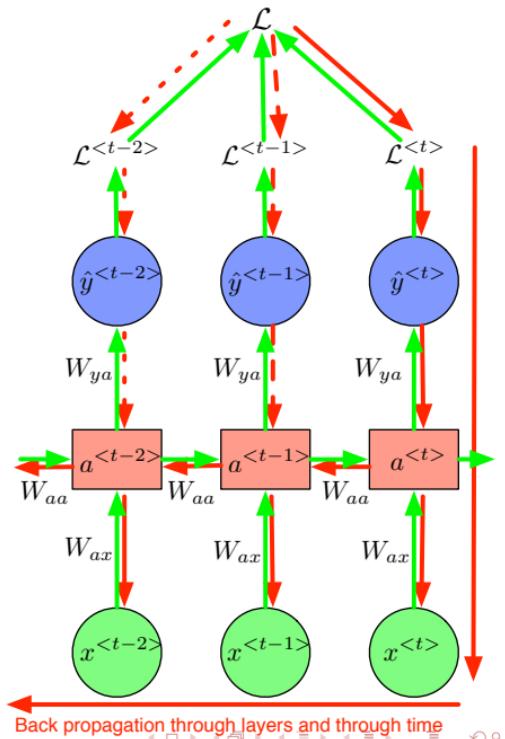
$$\frac{\partial \mathcal{L}}{\partial W_{ya}} = \sum_{t=1}^T \frac{\partial \mathcal{L}^{(t)}}{\partial W_{ya}}$$

- We can use the chain rule.

$$\frac{\partial \mathcal{L}^{(t)}}{\partial W_{ya}} = \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial W_{ya}}$$

- This is very simple since \hat{y}_t depends on W_{ya} on one place (indeed $a^{(t)}$ does not depend on W_{ya}).

$$\hat{y}^{(t)} = g_y(W_{ya}a^{(t)} + b_y)$$



Back Propagation Through Time (BPTT)

Backward pass : 2) compute $\frac{\partial \mathcal{L}}{\partial W_{aa}}$

- We measure how much varying W_{aa} affect $\mathcal{L}^{(t)}$.

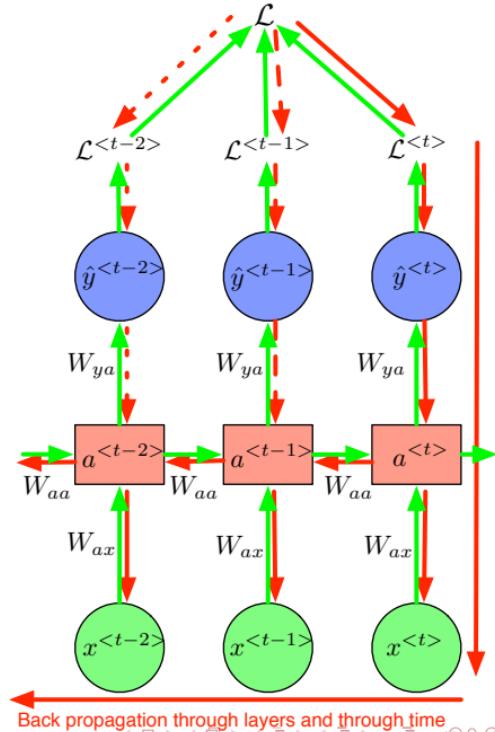
$$\frac{\partial \mathcal{L}}{\partial W_{aa}} = \sum_{t=1}^T \frac{\partial \mathcal{L}^{(t)}}{\partial W_{aa}}$$

- We cannot use the chain rule ($\frac{\partial \mathcal{L}^{(t)}}{\partial W_{aa}} = \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial a^{(t)}} \frac{\partial a^{(t)}}{\partial W_{aa}}$) because a_t depends on W_{aa} also through $a^{(t-1)}$ which itself

$$\dots a^{(t)} = g_a(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b_a)$$

- We use a formula of total derivative

$$\begin{aligned} \frac{d\mathcal{L}^{(t)}}{dW_{aa}} &= \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial a^{(t)}} \left(\frac{\partial a^{(t)}}{\partial W_{aa}} + \frac{\partial a^{(t)}}{\partial a^{(t-1)}} \frac{\partial a^{(t-1)}}{\partial W_{aa}} + \dots + \frac{\partial a^{(t)}}{\partial a^{(t-1)}} \dots \frac{\partial a^{(0)}}{\partial W_{aa}} \right) \\ &= \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\hat{y}^{(t)}}{a^{(t)}} \sum_{k=0}^t \frac{\partial a^{(t)}}{\partial a^{(t-1)}} \dots \frac{\partial a^{(k+1)}}{\partial a^{(k)}} \frac{\partial a^{(k)}}{\partial W_{aa}} \\ &= \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial a^{(t)}} \sum_{k=0}^t \left(\prod_{i=k+1}^t \frac{\partial a^{(i)}}{\partial a^{(i-1)}} \right) \frac{\partial a^{(k)}}{\partial W_{aa}} \end{aligned}$$



Back Propagation Through Time (BPTT)

Backward pass : 3) compute $\frac{\partial \mathcal{L}}{\partial W_{ax}}$

- We measure how much varying W_{ax} affect $\mathcal{L}^{(t)}$.

$$\frac{\partial \mathcal{L}}{\partial W_{ax}} = \sum_{t=1}^T \frac{\partial \mathcal{L}^{(t)}}{\partial W_{ax}}$$

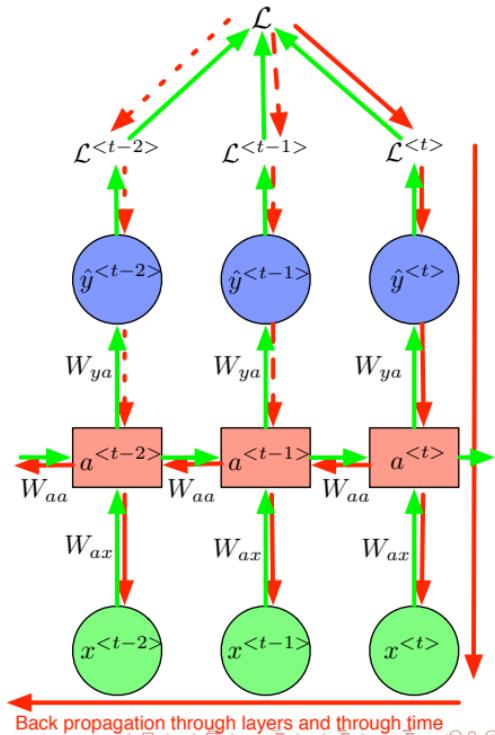
- It is also necessary to go backwards in time to calculate $\frac{\partial \mathcal{L}}{\partial W_{ax}}$

- same situation as with W_{aa} :
- $a^{(t)}$ depends on W_{ax} also through $a^{(t-1)}$ which itself ...

$$a^{(t)} = g_a(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b_a)$$

- We use a formula of total derivative

$$\begin{aligned} \frac{d\mathcal{L}^{(t)}}{dW_{ax}} &= \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial a^{(t)}} \left(\frac{\partial a^{(t)}}{\partial W_{ax}} + \frac{\partial a^{(t)}}{\partial a^{(t-1)}} \frac{\partial a^{(t-1)}}{\partial W_{ax}} + \dots + \frac{\partial a^{(t)}}{\partial a^{(t-1)}} \cdots \frac{\partial a^{(0)}}{\partial W_{ax}} \right) \\ &= \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial a^{(t)}} \sum_{k=0}^t \frac{\partial a^{(t)}}{\partial a^{(t-1)}} \cdots \frac{\partial a^{(k+1)}}{\partial a^{(k)}} \frac{\partial a^{(k)}}{\partial W_{ax}} \\ &= \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial a^{(t)}} \sum_{k=0}^t \left(\prod_{i=k+1}^t \frac{\partial a^{(i)}}{\partial a^{(i-1)}} \right) \frac{\partial a^{(k)}}{\partial W_{ax}} \end{aligned}$$



Vanishing and exploding gradients

Vanishing gradient

- Long-term dependency
 - The **student**, which already followed 6 hours of lessons, **was** tired.
 - The **students**, which already followed 6 hours of lessons, **were** tired.
- Very deep NN
 - very difficult to propagate back the gradient to affect the weights of the early Layers

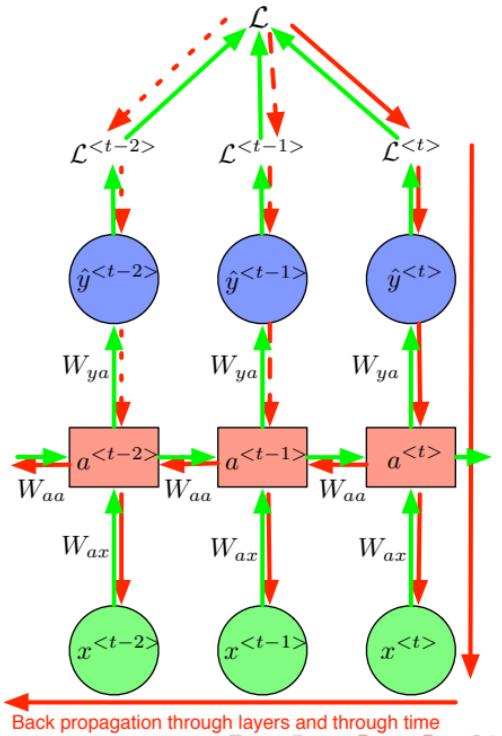
Vanishing and exploding gradients

- To train an RNN we need to backpropagate through layers and through time

$$\frac{\partial \mathcal{L}}{\partial W_{aa}} = \sum_{t=0}^T \frac{\partial \mathcal{L}^{(t)}}{\partial W_{aa}}$$

$$\frac{\partial \mathcal{L}^{(t)}}{\partial W_{aa}} \propto \sum_{k=0}^t \left(\prod_{i=k+1}^t \frac{\partial a^{(i)}}{\partial a^{(i-1)}} \right) \frac{\partial a^{(k)}}{\partial W_{aa}}$$

- each term in the sum is the contribution of
 - a state at time k
 - to the gradient of the loss at time step t
- the more steps between k and t , the more elements in this product
- the values of these Jacobian matrices have particularly severe impact on the contributions of faraway steps



Vanishing and exploding gradients

$$\frac{\partial \mathcal{L}}{\partial W_{aa}} = \sum_{t=0}^T \frac{\partial \mathcal{L}^{(t)}}{\partial W_{aa}}$$

$$\frac{\partial \mathcal{L}^{(t)}}{\partial W_{aa}} \propto \sum_{k=0}^t \left(\prod_{i=k+1}^t \frac{\partial a^{(i)}}{\partial a^{(i-1)}} \right) \frac{\partial a^{(k)}}{\partial W_{aa}}$$

- Suppose only one hidden units, then a_i is a scalar and consequently $\frac{\partial a^{(i)}}{\partial a^{(i-1)}}$ is also a scalar
 - If $\left| \frac{\partial a^{(i)}}{\partial a^{(i-1)}} \right| < 1$ then the product goes to 0 exponentially fast
 - If $\left| \frac{\partial a^{(i)}}{\partial a^{(i-1)}} \right| > 1$ then the product goes to infinity exponentially fast
- Vanishing gradients
 - $\left| \frac{\partial a^{(i)}}{\partial a^{(i-1)}} \right| < 1$
 - contributions from faraway steps vanish and don't affect the training
 - difficult to learn long-range dependencies
- Exploding gradients
 - $\left| \frac{\partial a^{(i)}}{\partial a^{(i-1)}} \right| > 1$
 - make the learning process unstable
 - gradient could even become NaN

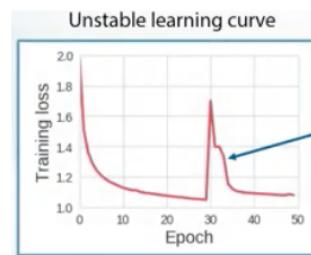
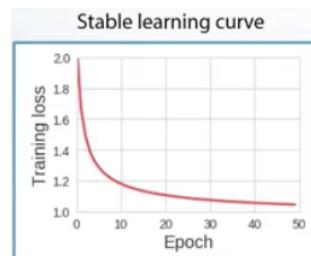
Dealing with the exploding gradient problem

- **Solution 1 : gradient clipping**

- gradient $d\theta = \frac{\partial \mathcal{L}}{\partial \theta}$, where θ are all the parameters of the network
- if $\|d\theta\| > \text{threshold}$

$$d\theta \leftarrow \frac{\text{threshold}}{\|d\theta\|} d\theta$$

- clipping doesn't change the direction of the gradient but change its length
- we can clip only the norm of the part which causes the problem
- we choose the threshold manually : start with a large threshold and then reduce it

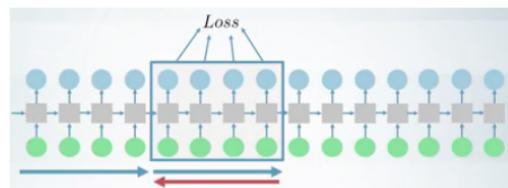
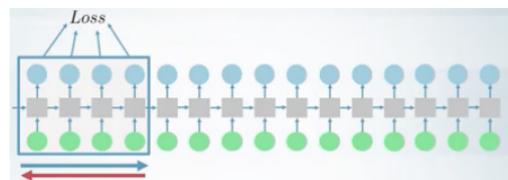
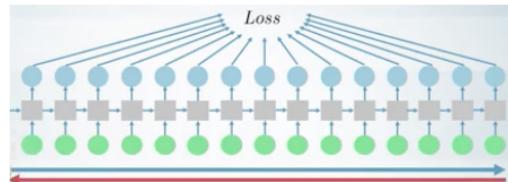


Back Propagation Through Time (BPTT)

Dealing with the exploding gradient problem

- **Solution 2 : truncated BPTT**

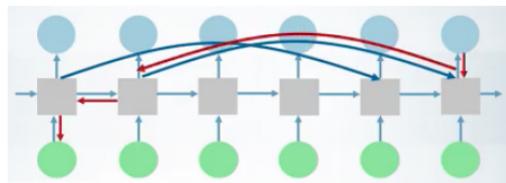
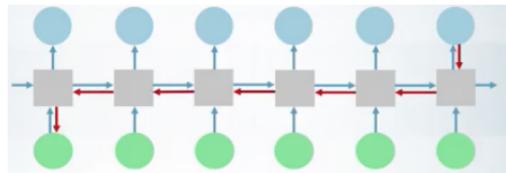
- Training very long sequences
 - time consuming !
 - exploding gradients !
- Truncated BPTT :
 - run forward and backward passes through the chunks of the sequence (instead of the whole sequence)
- Forward
 - We carry hidden states forward in time forever
- Backward
 - only backpropagate in the chunks (small number of steps)
- Much faster but dependencies longer than the chunk size don't affect the training (at least they still work at Forward pass)



Dealing with the vanishing gradient problem

- **Solution 2 : use skip-connections**

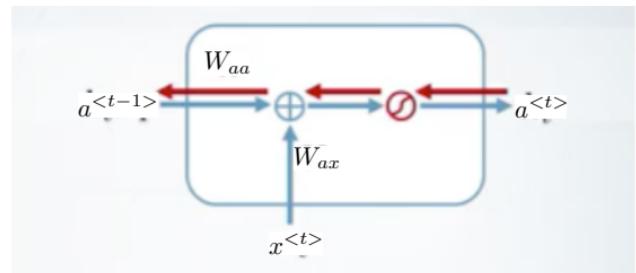
- in RNN : because at each step we multiply the Jacobian matrices → vanishing gradient
 - we cannot learn long-term dependencies
- Solution : add short-cuts between hidden states that are separated by more than one time step
 - are usual connections (with their own parameter matrices)
 - create much shorter paths between far away time-steps
- Backpropagation through the shortcuts : the gradients vanish slower
 - we can learn long-term dependencies
- not exclusive to RNN (see ResNet)



Main ideas : Recurrent Neural Network (RNN)

$$a^{(t)} = g \left(W_{ax} x^{(t)} + W_{aa} a^{(t-1)} + b_a \right)$$

- non-linearities and/or multiplication creates the vanishing gradient problem
- Solution ?
 - create a short-way for backpropagation without any non-linearities or multiplication



Gated units (LSTM and GRU)

Main ideas : Long Short Term Memory Units (LSTM)

[Hochreiter and Schmidhuber 1997. "Long short-term memory", Neural Computation Volume 9 | Issue 8 | November 15, 1997 p.1735-1780]

- Add a **memory cell**
 - the same dimension as hidden layer
- **Simplified LSTM** (no possibility to erase anything) :
 - Candidate cell value

$$\tilde{c}^{(t)} = g \left(W_{ax}^c x^{(t)} + W_{aa}^c a^{(t-1)} + b_g \right)$$

Update gate

$$\Gamma_u = \sigma \left(W_{ax}^u x^{(t)} + W_{aa}^u a^{(t-1)} + b_u \right)$$

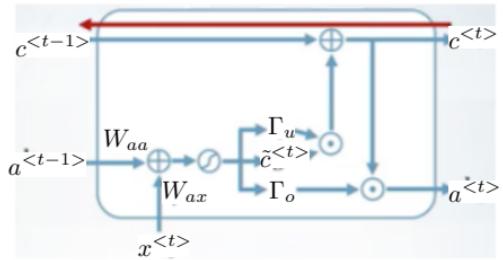
Output gate

$$\Gamma_o = \sigma \left(W_{ax}^o x^{(t)} + W_{aa}^o a^{(t-1)} + b_o \right)$$

$$c^{(t)} = c^{(t-1)} + \Gamma_u \odot \tilde{c}^{(t)}$$

$$a^{(t)} = \Gamma_o \odot f(c^{(t)})$$

- no non-linearity or multiplication in the memory cell-path
- $\frac{\partial c^{(t)}}{\partial c^{(t-1)}} = \text{diag}(1)$
- no vanishing gradients!!!



Main ideas : Long Short Term Memory Units (LSTM)

[Hochreiter and Schmidhuber 1997. "Long short-term memory", Neural Computation Volume 9 | Issue 8 | November 15, 1997 p.1735-1780]

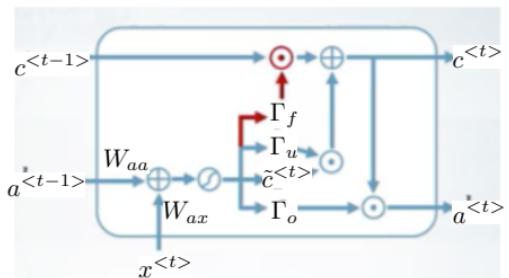
- **Forget gate LSTM :**

Forget gate

$$\Gamma_f = \sigma(W_{ax}^f x^{(t)} + W_{aa}^f a^{(t-1)} + b_f)$$

$$c^{(t)} = \Gamma_f \odot c^{(t-1)} + \Gamma_u \odot \tilde{c}^{(t)}$$

$$a^{(t)} = \Gamma_o \odot f(c^{(t)})$$

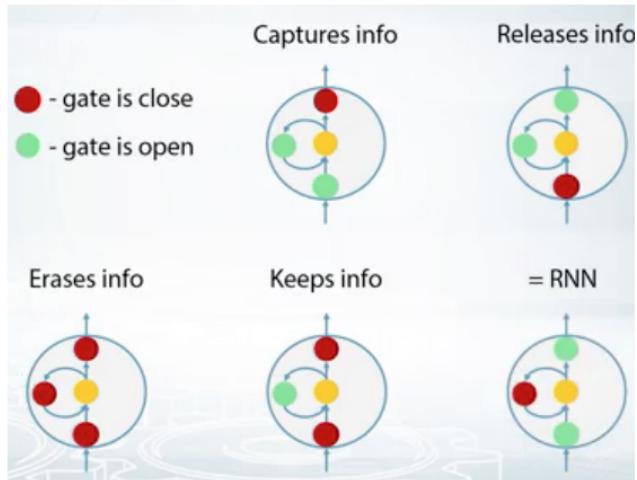
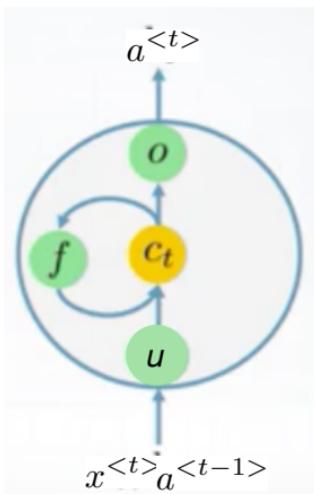


- We now have a multiplication on the short-way

- $\frac{\partial c^{(t)}}{\partial c^{(t-1)}} = \text{diag}(\Gamma_f)$
- with Γ_f the forget gate : $\Gamma_f = \sigma(W_{ax}^f x^{(t)} + W_{aa}^f a^{(t-1)} + b_f)$
- Γ_f takes value between 0 and 1
 - \Rightarrow set a high initial value for b_f
 - at the beginning of the training, the LSTM doesn't forget and can find long-term dependencies in the data

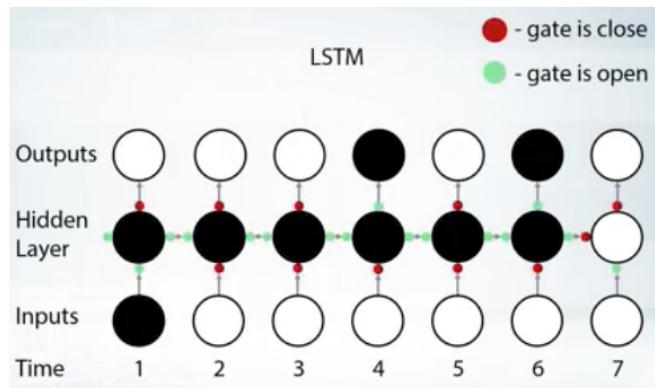
- LSTM are very nice but has four times more parameters than RNN

LSTM regimes



Main ideas : Long Short Term Memory Units (LSTM)

- LSTM regimes

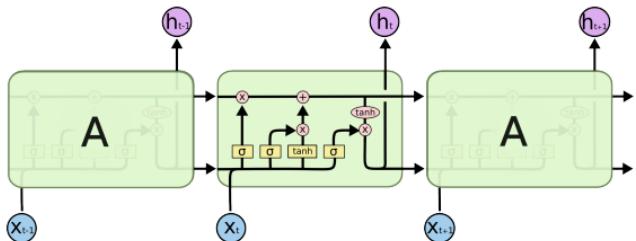
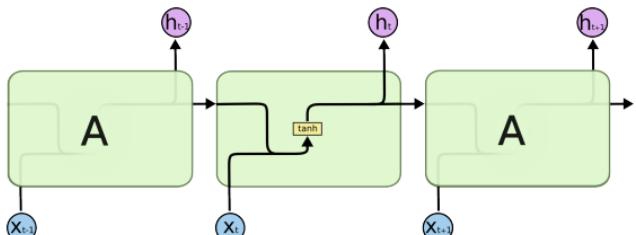


LSTM Example usage

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

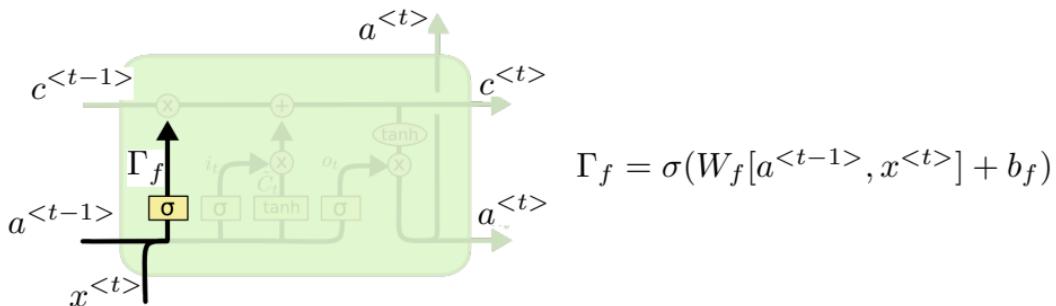
Instead of having a single neural network layer, there are four, interacting in a very special way.

- Lines carry a vector
- A recurrent net transmits its hidden states
- LSTM introduces a second channel : the cell state
- It acts as a memory
- Gates control the memory



LSTM Example usage

- What should be forgotten from the previous cell state ?



Action

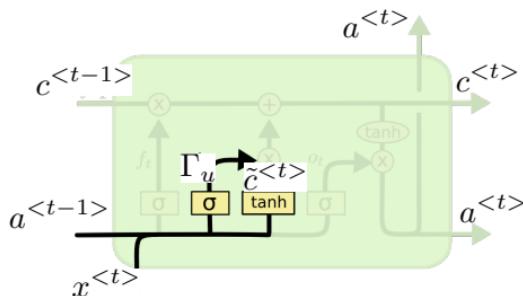
- The sigmoid (forget gate) answers for each component :
 - 1 : to keep it,
 - 0 to forget it, or
 - a value in-between to mitigate its influence

Intuition for language modeling

- The cell state might embed the gender of the present subject :
 - keep it to predict the correct pronouns,
 - or forget it, when a new subject appears

LSTM Example usage

- What should be taken into account ?



$$\begin{aligned}\tilde{c}^{<t>} &= \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c) \\ \Gamma_u &= \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)\end{aligned}$$

Action

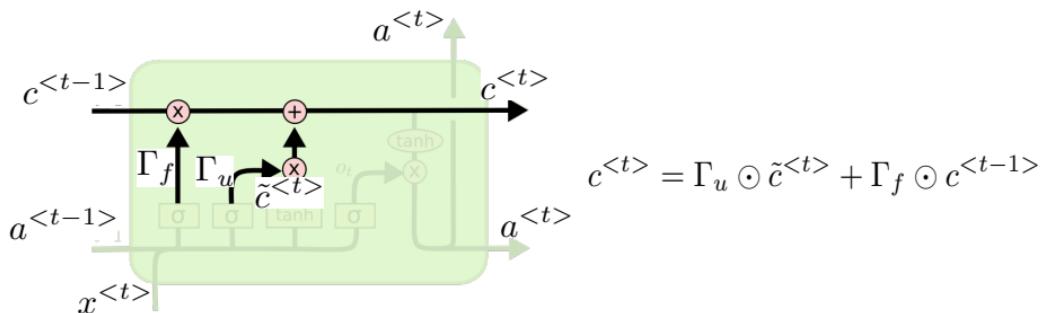
- Create the update $\tilde{c}^{<t>}$ of the cell state
- and its contribution Γ_u (the update gate with a sigmoid activation)

Intuition for language modeling

- Add the gender of the new subject to the cell state,
- to replace the old one we're forgetting. appears

LSTM Example usage

- Write the new state



Action

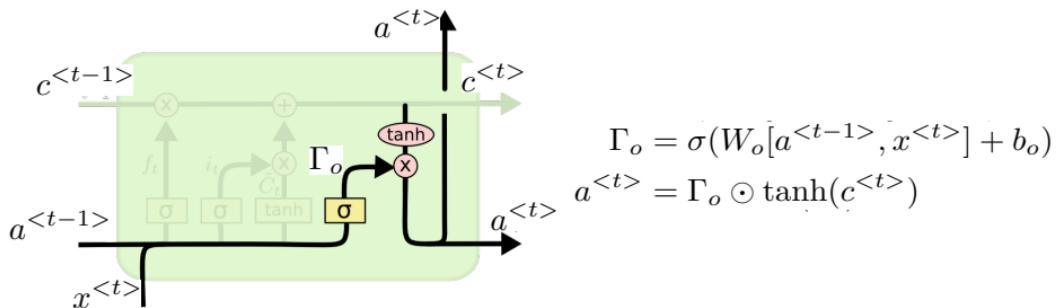
- Merge the old cell state modified by the forget gate
- with the new input

Intuition for language modeling

- Decide to drop the information about the old subject
- Refresh the memory appears

LSTM Example usage

- Write the new hidden state



Action

- Decide what parts of the (filtered) cell state to output $a^{<t>}$
- Compute the hidden state

Intuition for language modeling

- Since we just saw a subject,
- output the relevant information for the future (gender, number) appears

Gated Recurrent Unit (GRU)

[Cho et al. 2014. "On the Properties of Neural Machine Translation: Encoder–Decoder Approaches"]

[Chung et al. 2014 "Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling"]

- Simple GRU

Candidate cell value

$$\tilde{c}^{<t>} = \tanh(W_c[c^{<t-1>}, x^{<t>}] + b_c)$$

Update gate

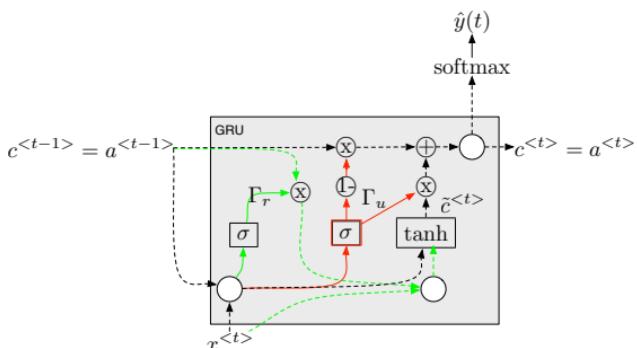
$$\Gamma_u = \sigma(W_c[c^{<t-1>}, x^{<t>}] + b_c)$$

Cell update

$$c^{<t>} = \Gamma_u \odot \tilde{c}^{<t>} + (1 - \Gamma_u) \odot c^{<t-1>}$$

Output

$$c^{<t>} = a^{<t>}$$



Gated Recurrent Unit (GRU)

[Cho et al. 2014. "On the Properties of Neural Machine Translation: Encoder–Decoder Approaches"]

[Chung et al. 2014 "Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling"]

• Full GRU

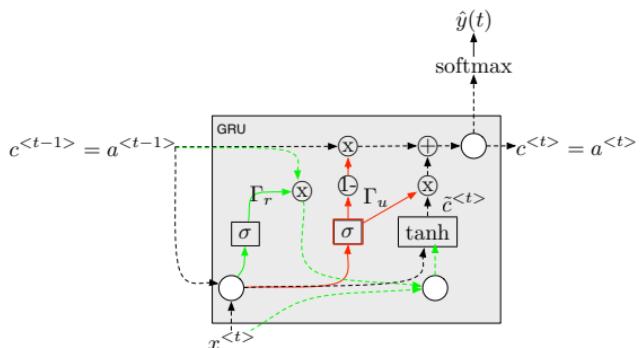
- adding a relevance gate Γ_r
- instead of using all $c^{}$ to compute $\tilde{c}^{}$, we will first filter it using a "Relevant gate" Γ_r .
- The "Relevant gate" Γ_r is computed using $c^{}$ and $x^{}$.

Relevant gate

$$\Gamma_r = \sigma(W_r[c^{}, x^{}] + b_r)$$

Cell update

$$\tilde{c}^{} = \tanh(W_c[\Gamma_r \odot c^{}, x^{}] + b_c)$$



Language model

- **Applications**

- Automatic Speech Recognition, Machine Translation, OCR, ...

- **Automatic Speech Recognition :**

- "This is an **example** of an homophone." or "This is an **egg-sample** of an homophone."
 - $p(\text{"example"}) = 5.7 \cdot 10^{-10}$
 - $p(\text{"egg - sample"}) = 3.2 \cdot 10^{-13}$

- **Language model : $p(\text{sentence})$**

- $p(y^{(1)}, y^{(2)}, \dots, y^{(T_y)})$

Language model

- Notation :

$$\mathbf{w} = w_1^L : w_1, w_2, \dots, w_L$$

- Goal

- Estimate the (non-zero) probability of a word sequence for a given vocabulary

$$\begin{aligned} P(w_1^L) &= P(w_1, w_2, \dots, w_L) = \prod_{i=1}^L P(w_i | w_1^{i-1}) \quad \forall i, w_i \in V \\ &= P(w_1)P(w_2|w_1)P(w_3|w_1, w_2)P(w_4|w_1, w_2, w_3) \dots P(w_L|w_1 \dots w_{L-1}) \end{aligned}$$

- with the N-gram assumption :

$$\begin{aligned} P(w_1^L) &= \prod_{i=1}^L P(w_i | w_{i-n+1}^{i-1}), \forall i, w_i \in V \\ &= P(w_1)P(w_2|w_1)P(w_3|w_2)P(w_4|w_2) \dots P(w_L|w_{L-1}) \end{aligned}$$

- in the recurrent way :

$$P(w_i | w_1^{i-1})$$

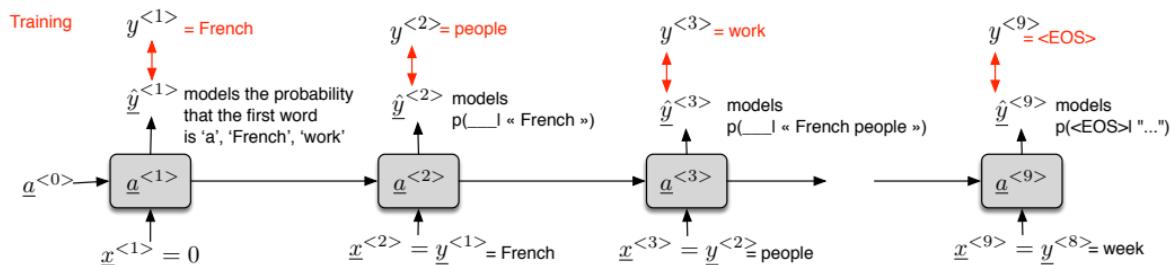
- How to build a language model ?

- need a large corpus of English text
- define a dictionary
- tokenize the text
 - convert to one-hot-vector
 - + <EOS> (end of sentence)
 - + <UNK> (unknown word, very uncommon vocalular) : Gif-sur-Yvette

Language model/ ① Training using a RNN

- **Algorithm :**

- start with empty context, empty input word
- given the previous word as input $x^{(t)}$ and the context (in the hidden cell $a^{(t)}$) predict (maximize the probability to observe) the next word $y^{(t)}$ as output $\hat{y}^{(t)}$
- "French people work an average of 35 hours a week."



- **Training :**

- output of the network (softmax) is $\hat{y}^{(t)}$ (probability of each word)
- ground-truth is $y^{(t)}$ (one-hot-encoding)
- minimize the loss

$$\mathcal{L}_t(\hat{y}^{(t)}, y^{(t)}) = - \sum_{c=1}^K y_c^{(t)} \log(\hat{y}_c^{(t)})$$

$$\mathcal{L} = \sum_t \mathcal{L}_t(\hat{y}^{(t)}, y^{(t)})$$

Language model/ ② Testing using a RNN

- Given a new sentence of three words ($y^{(1)}, y^{(2)}, y^{(3)}$)
- what is the probability of this sentence?

$$P(y^{(1)}, y^{(2)}, y^{(3)}) = p(y^{(1)})$$

given by the first softmax

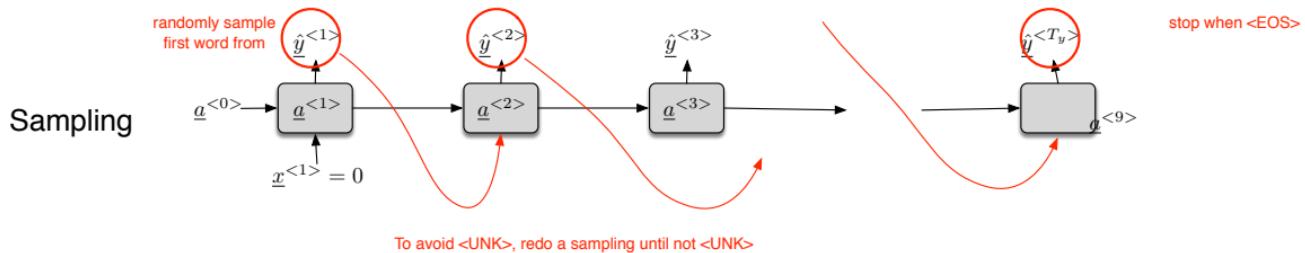
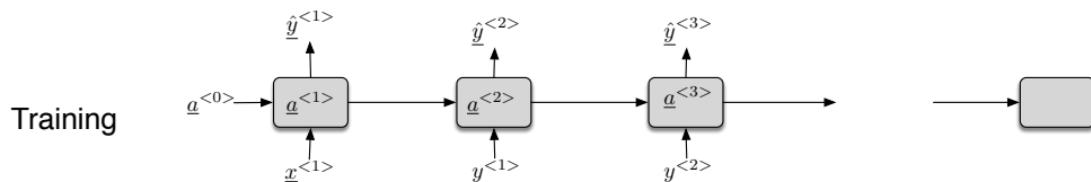
$$p(y^{(2)}|y^{(1)})$$

given by the second softmax

$$p(y^{(3)}|y^{(1)}, y^{(2)})$$

given by the third softmax

Language model/ ③ Sampling novel sequences using RNN



- Word-level RNN
- Character-level RNN :
 - Vocabulary= ['a', 'b', 'c', ..., '.', ',', ';', 'A', 'B', 'C', ...]
- Examples :
 - Shakespeare, wikipedia, source-code generation
 - <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

One-hot-encoding

- How to represent the individual words in a sentence ?
- What will be $x^{(t)}$?
- Construct the list of the N most used categories (words) in the corpus
 - = "vocabulary/dictionary"
 - N is usually in the range 10.000 words - 50.000.
 - add a <UNK> (unknown word) for the words that are not in vocabulary/dictionary
- Each word is now associated with an ID

1	Angela
...	Boris
...	chancellor
367	Emmanuel
...	France
...	Germany
4075	is
...	Johnson
6830	Macron
...	Merkel
...	of
...	president
...	prime-minister
...	the
10.000	UK

One-hot-encoding

- **One-hot-encoding :**

- the one-hot vector of an ID is a vector filled with 0s, except for a 1 at the position associated with the ID
- each word w is associated with a one-hot-vector vector o_{ID}
- If $N = 10.000$, $x^{(t)}$ has dimension 10.000
- a one-hot encoding makes no assumption about word similarity
 - all words are equally different from each other

$\{x\}$	$x^{<1>}$	$x^{<2>}$	$x^{<3>}$	$x^{<4>}$	$x^{<5>}$	$x^{<6>}$	$x^{<7>}$
	Emmanuel	Macron	is	the	president	of	France
1	0	0	0				
2	0	0	0				
...				
367	1	0	0				
...				
4075	0	0	1				
...				
6830	0	1	0				
...				
10.000	0	0	0				

One-hot-encoding

- **Weaknesses of one-hot-encoding**
 - it is very high-dimensional
 - vulnerability to overfitting, computationally expensive
 - all words are equally different from each other
 - the inner product between any two one-hot vectors is zero
 - as a consequence, we cannot generalize
- Example :
 - we have trained a language model to complete the sentence
 - The president is on a **boat** ? " ⇒ "trip"
 - since there is no specific relationship between "boat" and "plane", the algorithm cannot complete the sentence
 - The president is on a **plane** ? " ⇒ ?

Word embedding

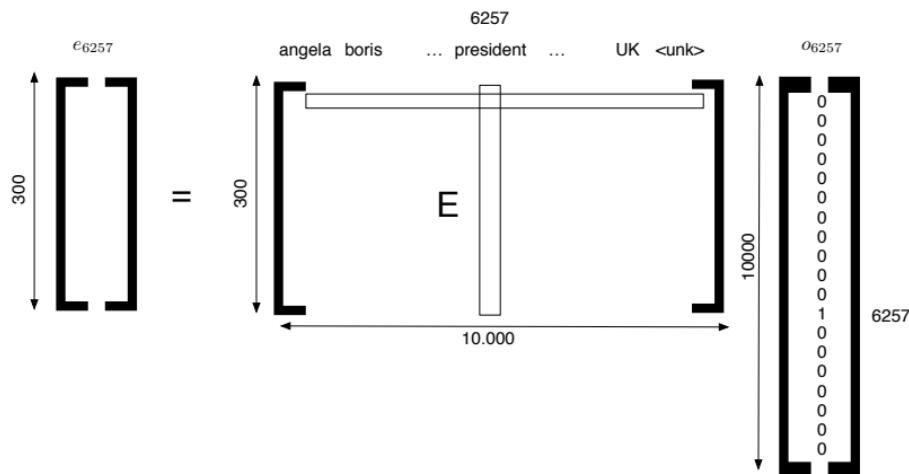
Word embedding

- **Word Embedding** : learn a continuous representation of words
 - each word w is associated with a real-valued vector e_{ID}
 - if embedding size is 300, e_{5391} is a 300 dimensional vector
 - we would like the distance $\|e_{ID1} - e_{ID2}\|$ to reflect the meaningful similarities between words

	Man	Woman	King	Queen	Uncle	Aunt	President	Chancellor	France	Germany	Boat	Plane
	5391	9853	4914	7157	456	6257	7124	3212	6789	1234	5923	8871
Gender	-1	1	-1	1	-1	1	-0,7	-0,3	0	0	0	0
Royal	0	0	1	1	0	0	0,5	0,5	0	0	0	0
Age	0	0	0,7	0,7	0,5	0,5	0,7	0,7	0	0	0	0
Country	0	0	0,5	0,5	0	0	0,5	0,2	1	1	0	0
Transportation	0	0	0	0	0	0	0	0	0	0	1	1
...

- In this representation "boat" and "plane" are quite similar
 - some of the features may differ but most features are similar
- Therefore the algorithm can find (by generalization) that
 - The president is on a boat ? " \Rightarrow "trip"
 - The president is on a plane ? " \Rightarrow "trip"

Word embedding/ Matrix



- \underline{o}_{6257} = **one-hot vector** (zero everywhere except 1 at position 6257)
 - dimension 10,000
- $\underline{e}_{6257} = \frac{\underline{\underline{I}}}{(300)} \cdot \underline{o}_{6257} = \text{embedding of } \underline{o}_{6257}$
- $\underline{e}_j = \underline{\underline{I}} \cdot \underline{o}_j = \text{embedding for word } j$
- In practice, use **look-up table** (take the column vector of $\underline{\underline{I}}$ corresponding to 6257)

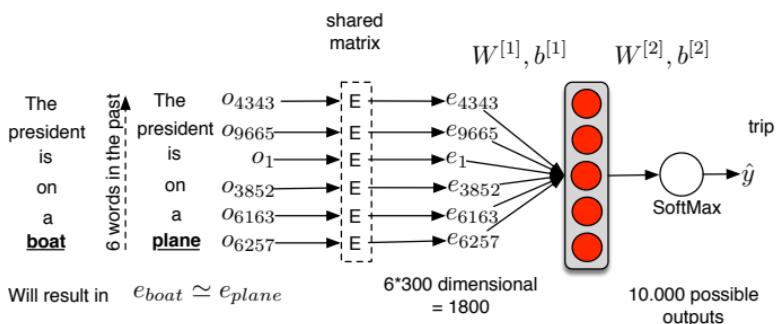
Word embedding/ Learning using

- Various existing methods to learn word embedding :
 - Classic neural language model (Bengio)
 - Collobert and Weston model
 - Word2Vec (Mikolov)
 - Continuous bag-of-words (CBOW)
 - Skip-gram
 - Glove model (Pennington)

Word embedding/ Classic neural language model

[Bengio et al., 2003 "A Neural probabilistic model"]

- **Method** : build a language model (learn to predict the following word) using a MLP
 - Use backpropagation to learn the parameters : $E, W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$
 - E is the **embedding matrix**



- **Example :**
 - "The president is on a boat $\underline{\text{trip}_{\text{target}}}$ to the United States."

- **Other possible definition of the context**
 - last 4 words
 - 4 words on the left, 4 words on the right : "is on a boat $\underline{\text{?}}$ to the United States"
 - last 1 words : "boat $\underline{\text{?}}$ "
 - nearby 1 word (skip-gram) : "president $\underline{\text{?}}$ "

Application : Sentiment Classification

X

This movie is fantastic ! I really like it because it is so good !

y



Not to my taste, will skip and watch another movie.

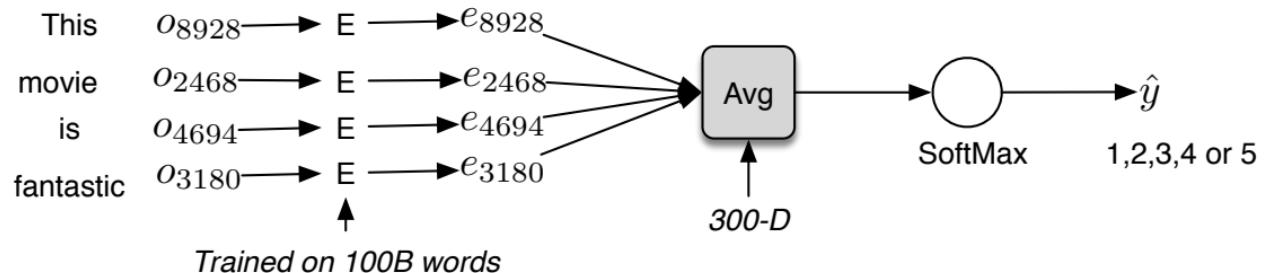


This movie was completely lacking a good script,
good actors and a good director.



Application : Sentiment Classification/ using a simple model

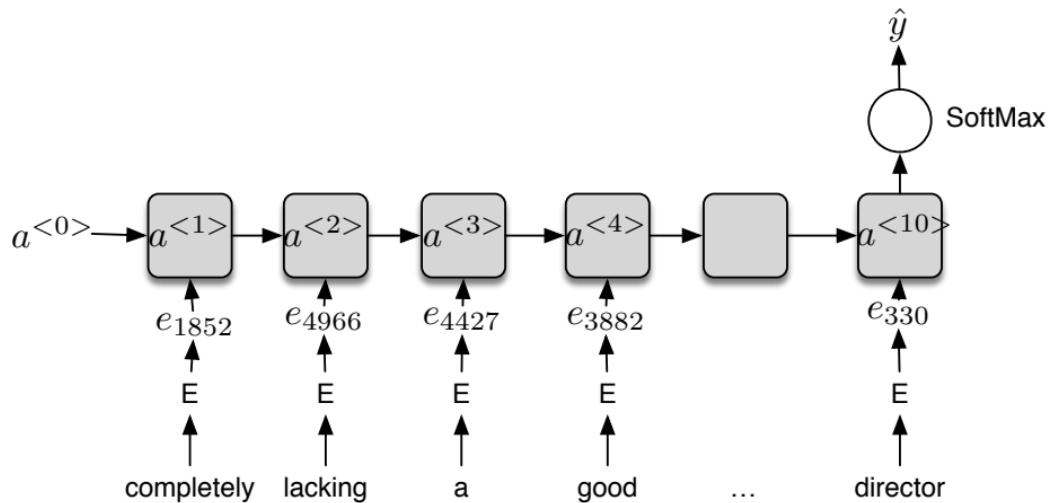
The	movie	is	fantastic
8928	2468	4694	3180



- Does not work for "This movie was completely lacking a **good** script, **good** actors and a **good** director."
 - Avg= good \Rightarrow does not understand

Application : Sentiment Classification/ using RNN

- **Many-to-one** architecture

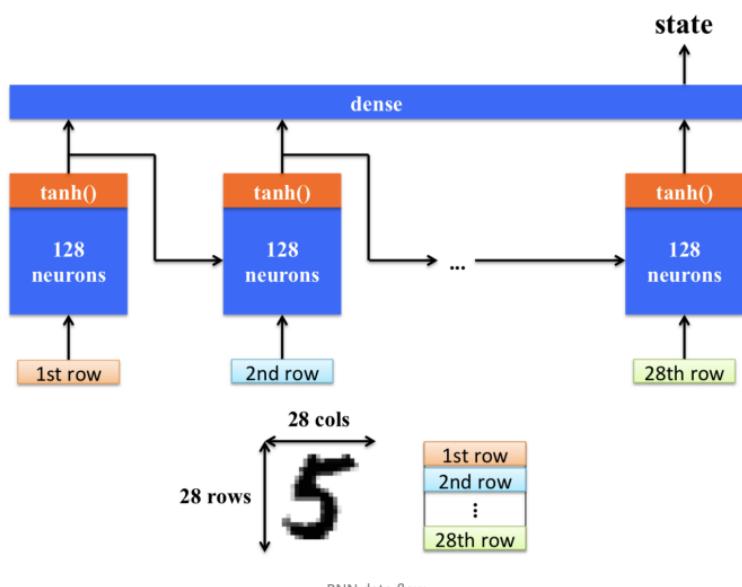


Application : Handwritten Character Recognition

[Graves et al. 2008 "Unconstrained online handwriting recognition with recurrent neural networks"]
[Graves et al. 2009 "Offline Handwriting Recognition with Multidimensional Recurrent Neural Networks"]

- source

<https://medium.com/machine-learning-algorithms/mnist-using-recurrent-neural-network-2d070a5915a2>



Introduction

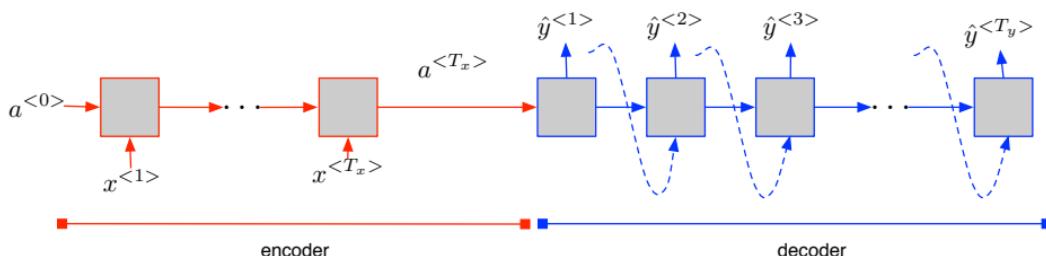
[Sutskever et al. 2014, "Sequence to sequence learning with neural networks"]

[Cho et al. 2014 "Learning phrase representations using rnn encoder-decoder for statistical machine translation"]

- What happens if $T_x \neq T_y$ (many-to-many architecture but with different lengths)
 - Example : Machine Translation

{x}	x ^{<1>}	x ^{<2>}	...			x ^{<6>}		
	Macron	voyage	aux	Etats-Unis à		Noël		
{y}	y ^{<1>}	y ^{<2>}	...				y ^{<8>}	
	Macron	is	travelling	to	the	United-Statfor	Christmas	

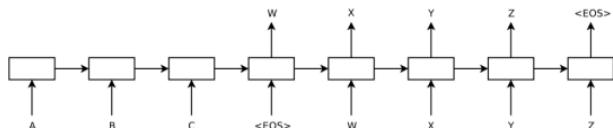
- Solution :
 - Sequence to sequence [Sutskever] or Encoder-Decoder [Cho] architecture



Sequence to sequence

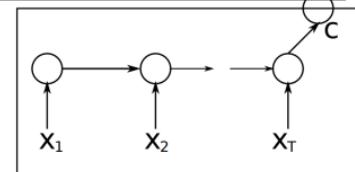
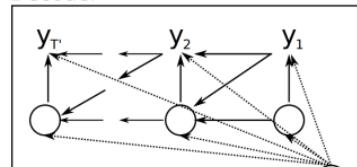
Introduction

[Sutskever et al. 2014, "Sequence to sequence learning with neural networks"]



[Cho et al. 2014 "Learning phrase representations using rnn encoder-decoder for statistical machine translation"]

Decoder

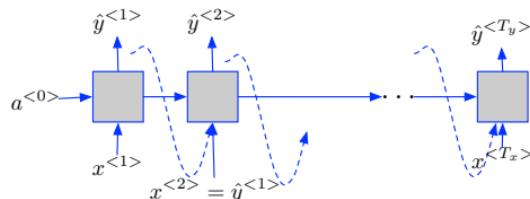


Encoder

Machine Translation

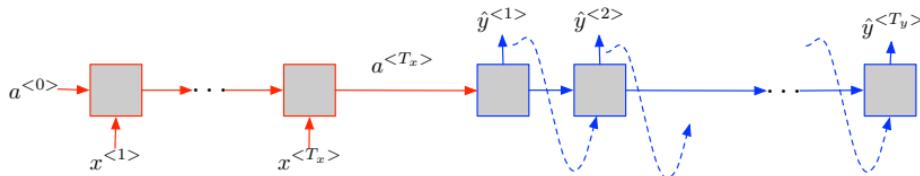
- **Language model :**

- $p(y^{<1>}, \dots, y^{<T_y>})$



- **Machine Translation :**

- Conditional language model : $p(y^{<1>}, \dots, y^{<T_y>} | x^{<1>}, \dots, x^{<T_x>})$
 - French sentence : $x^{<1>}, \dots, x^{<T_x>}$
 - English sentence : $y^{<1>}, \dots, y^{<T_y>}$



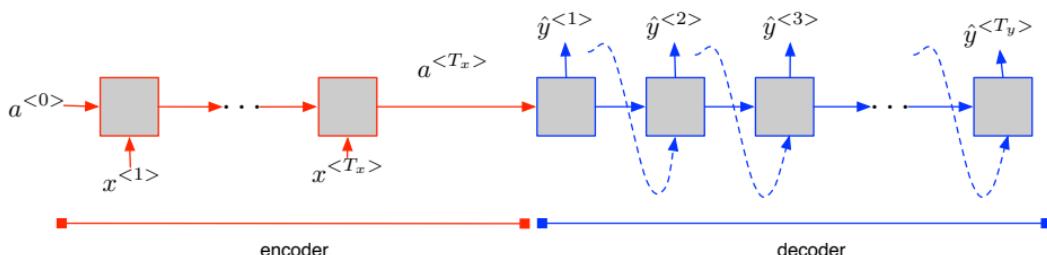
Attention model

[Bahdanau et al. 2015 "Neural machine translation by jointly learning to align and translate"]

- **The problem of long sequences**

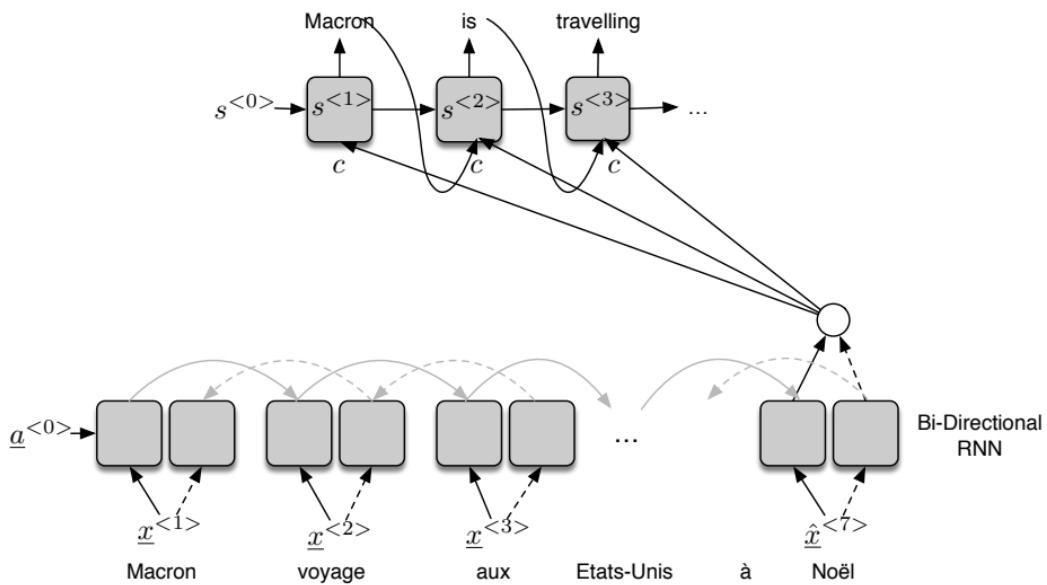
This kind of experience is part of Disney's efforts to "extend the lifetime of its series and build new relationships with audiences via digital platforms that are becoming ever more important," he added.

Ce type d'expérience fait partie des initiatives du Disney pour "prolonger la durée de vie de ses nouvelles et de développer des liens avec les lecteurs numériques qui deviennent plus complexes.



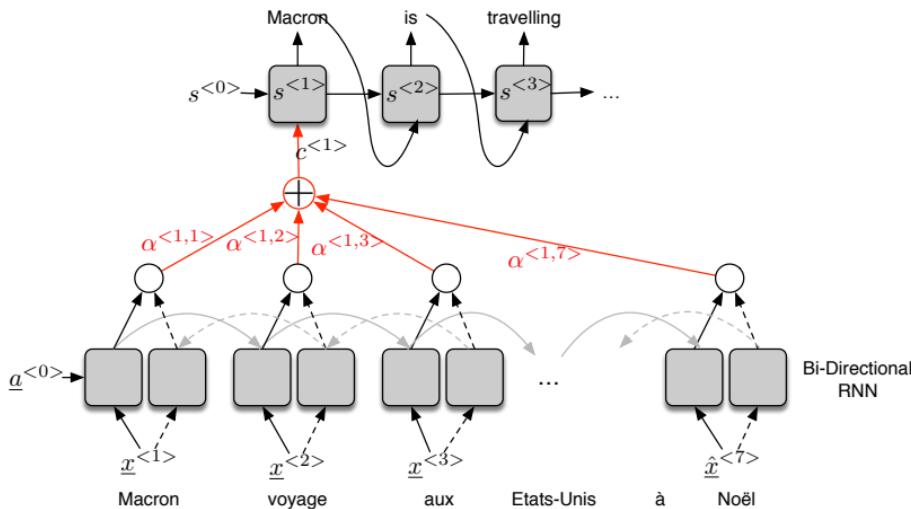
- Encoder/Decoder :
 - $a^{<T_x>}$ is supposed to memorize the whole sentence then translate it ;
 - human way : translate each part of a sentence at a time
- **Attention model ? (modification of the Encoder/Decoder)**
 - $a^{<T_x>}$ is replaced by a local version in the encoder \Rightarrow we pay attention on specific encoding

Attention model



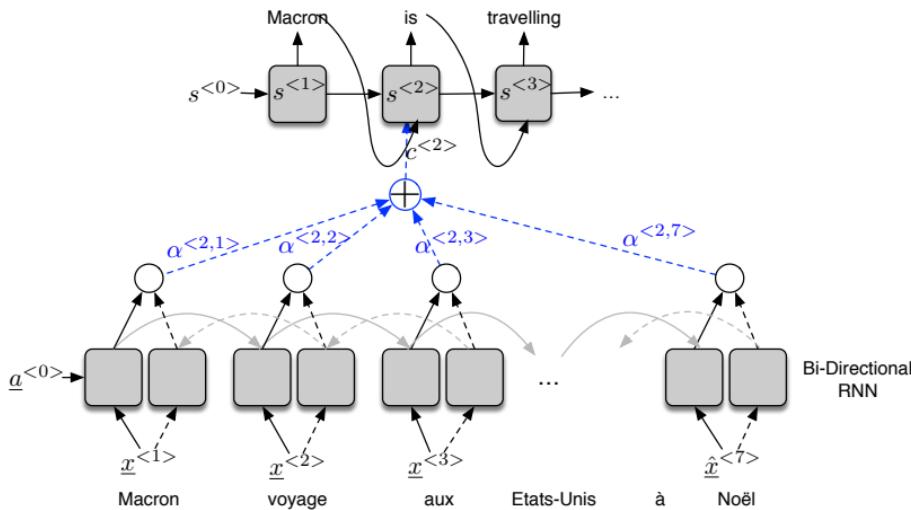
- In usual encoder-decoder,
 - information used for the decoding corresponds to the encoding at $T_x : c = a^{<T_x>}$

Attention model



- Attention model
 - replace $c = a^{<T_x>}$ by a local version $c^{<t>}$
 - $c^{<t>}$ is computed as a weighted sum of the encoding hidden states $a^{<t>}$
- **Attention weights $\alpha^{<\tau,t>}$:**
 - when generating information at time $\tau = 1$, how much, do we have to pay attention on information at time $t = 1, 2, 3 \dots 7$
- Notation :
 - $a^{<t>} / s^{<\tau>}$: hidden states of encoder/ decoder

Attention model



- Attention model
 - replace $c = a^{<T_x>}$ by a local version $c^{<t>}$
 - $c^{<t>}$ is computed as a weighted sum of the encoding hidden states $a^{<t>}$
- **Attention weights $\alpha^{<\tau,t>}$** :
 - when generating information at time $\tau = 1$, how much, do we have to pay attention on information at time $t = 1, 2, 3 \dots 7$
- Notation :
 - $a^{<t>} / s^{<\tau>}$: hidden states of encoder/ decoder

Attention model

- The **context vector** $c^{(\tau)}$ at decoding time τ
 - computed as the sum of the annotations $a^{(t)}$ weighted by their **attention weights** $\alpha^{(\tau,t)}$

$$c^{(\tau)} = \sum_t \alpha^{(\tau,t)} a^{(t)}$$

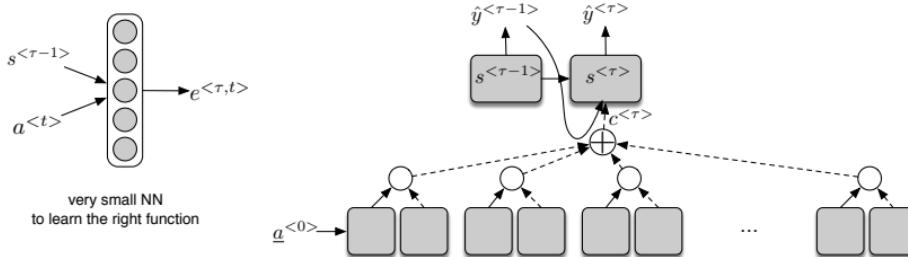
- The **attention weight** $\alpha^{(\tau,t)}$

- $\alpha^{(\tau,t)}$ = among of attention $y^{(\tau)}$ should pay to $a^{(t)}$
- computed using a softmax on the **alignment model** $e^{(\tau,t)}$

$$\alpha^{(\tau,t)} = \frac{\exp(e^{(\tau,t)})}{\sum_{t'=1}^{T_x} \exp(e^{(\tau,t')})} \quad \text{with } \sum_t \alpha^{(\tau,t)} = 1$$

- The **alignment model** $e^{(\tau,t)}$

- scores how well the inputs around position t match the output at position τ
- computed using two inputs
 - the RNN hidden state $s^{(\tau-1)}$ (just before emitting $y^{(\tau)}$)
 - the t -th annotation $a^{(t)}$ of the input sentence.
- computed using a feedforward neural network (jointly with all the other components)



Attention model

- Original sentence

This kind of experience is part of Disney's efforts to "extend the lifetime of its series and build new relationships with audiences via digital platforms that are becoming ever more important," he added.

- Encoder/Decoder without attention model

Ce type d'expérience fait partie des initiatives du Disney pour "prolonger la durée de vie de ses nouvelles et de développer des liens avec les lecteurs numériques qui deviennent plus complexes.

- Encoder/Decoder with attention model

Ce genre d'expérience fait partie des efforts de Disney pour "prolonger la durée de vie de ses séries et créer de nouvelles relations avec des publics via des plateformes numériques de plus en plus importantes", a-t-il ajouté.

- Visualization of $\alpha^{(\tau, t)}$ (English to French)

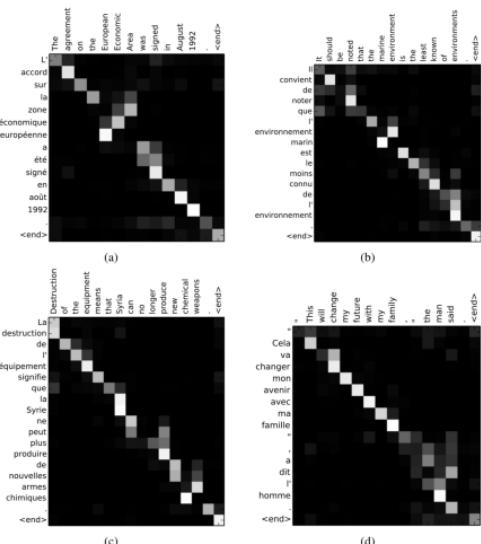
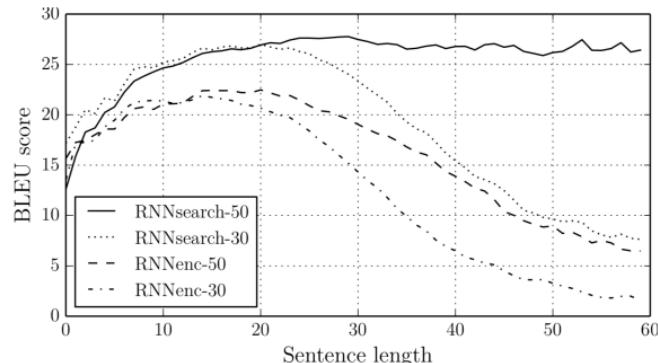


Figure 3: Four sample alignments found by RNNsearch-50. The x-axis and y-axis of each plot correspond to the words in the source sentence (English) and the generated translation (French), respectively. Each pixel shows the weight α_{ij} of the annotation of the j -th source word for the i -th target word (see Eq. (6)), in grayscale (0: black, 1: white). (a) an arbitrary sentence. (b-d) three randomly selected samples among the sentences without any unknown words and of length between 10 and 20 words from the test set.

Attention model

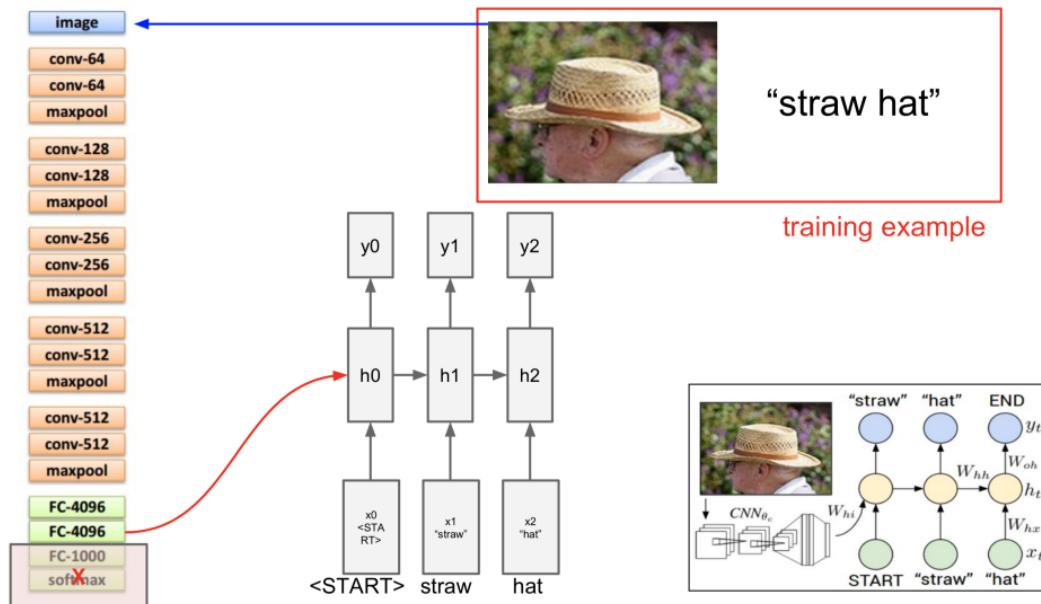
- Bleu score for long sentences



Application : Image captioning

[Mao et al. 2014 "Deep captioning with multimodal recurrent neural networks (m-rnn)"
[Vinyals et al. 2015 "Show and tell: A neural image caption generator"]
[Karpathy et al. 2015 "Deep visual-semantic alignments for generating image descriptions"]

- <http://cs231n.stanford.edu>



Transformer

[Vaswani et al. "Attention Is All You Need", NIPS, 2017]

[J. Alamar "The Illustrated Transformer"]

Transformer

- Global idea : get rid of complex recurrent or convolutional neural networks that include an encoder and a decoder
- a new simple network architecture based solely on attention mechanisms, dispensing with recurrence and convolutions entirely
- still an encoder decoder

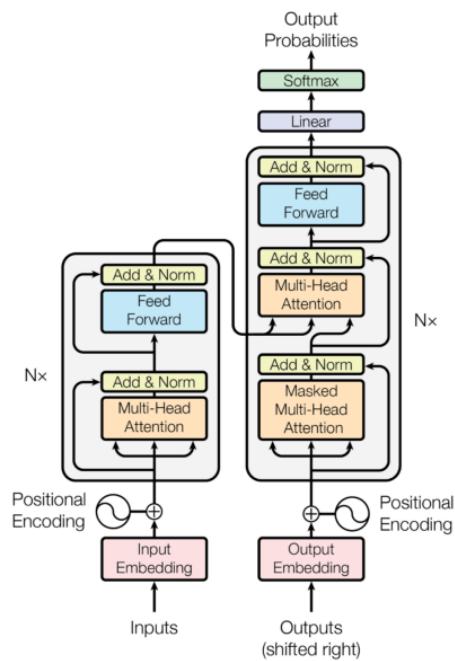


Figure 1: The Transformer - model architecture.

Transformer

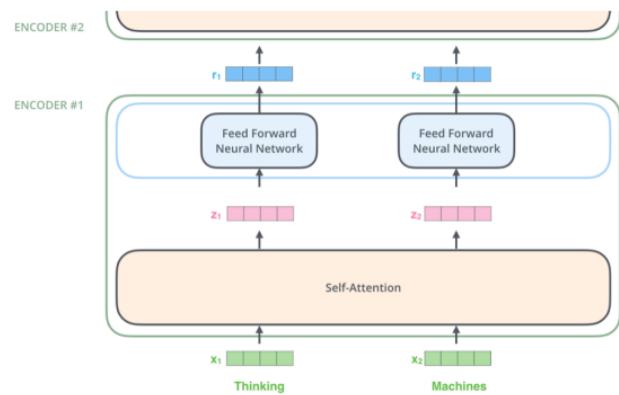
Self-Attention :

- Input sentence we want to translate : "The animal didn't cross the street because it was too tired"
 - What does "it" in this sentence refer to ? Is it referring to the "street" or to the "animal" ?
- When the model is processing the word "it", self-attention allows it to associate "it" with "animal"
- As the model processes each word (each position in the input sequence), self attention allows it to look at other positions in the input sequence for clues that can help lead to a better encoding for this word

Transformer

Now We're Encoding !

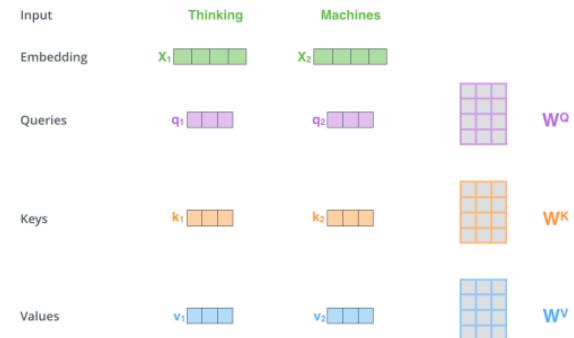
- An encoder receives a list of vectors as input
- It processes this list by passing these vectors into a 'self-attention' layer, then into a feed-forward neural network, then sends out the output upwards to the next encoder



Transformer

First step :

- create three vectors from each of the encoder's input vectors :
 - a Query vector,
 - a Key vector
 - a Value vector
- the vectors are created by multiplying the input by three matrices that we trained during the training process
- the new vectors are smaller in dimension than the embedding vector

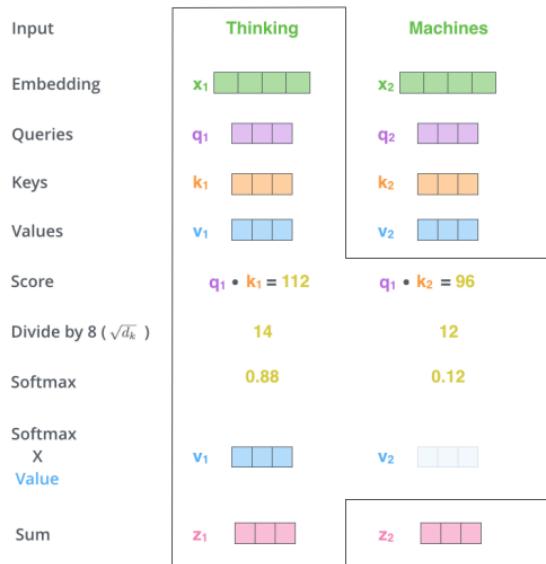


Multiplying x_1 by the W_Q weight matrix produces q_1 , the "query" vector associated with that word. We end up creating a "query", a "key", and a "value" projection of each word in the input sentence.

Transformer

Second step : calculate a score

- The score determines how much focus to place on other parts of the input sentence as we encode a word at a certain position.
- Example : calculate self-attention for the first word "Thinking" \Rightarrow we score each word of the input sentence against this word
- Score ?
 - dot product of the **query** vector with the **key** vector of the respective word we're scoring
 - divide by 8 (the square root of the dimension of the key vectors – 64)
 - pass the result through a softmax operation
 - softmax score determines how much each word will be expressed at this position



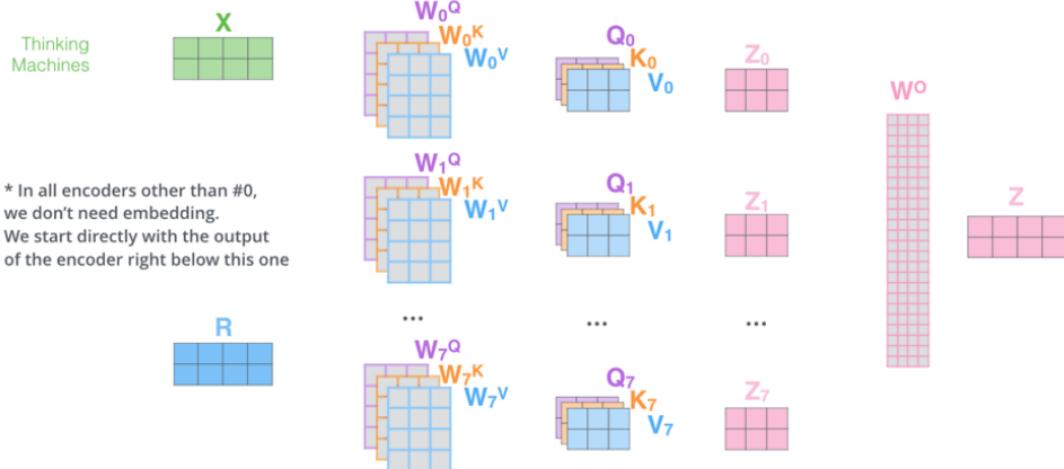
Third step :

- multiply each **value** vector by the **softmax** score
 - keep intact the values of the word(s) we want to focus on, and drown-out irrelevant words
- sum up the weighted value vectors
 - produces the output of the self-attention layer at this position (for the first word).

Transformer

In summary

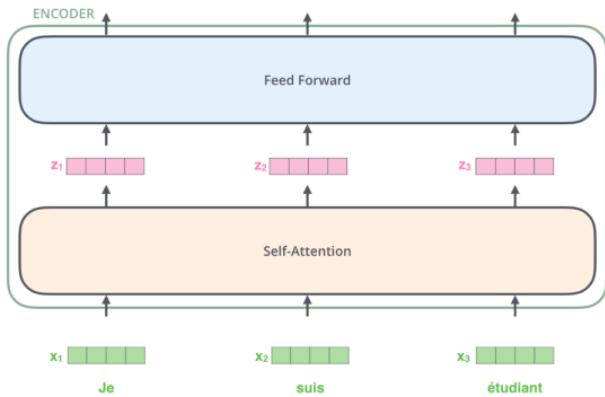
- 1) This is our input sentence* each word*
- 2) We embed
- 3) Split into 8 heads. We multiply \mathbf{X} or \mathbf{R} with weight matrices
- 4) Calculate attention using the resulting $\mathbf{Q}/\mathbf{K}/\mathbf{V}$ matrices
- 5) Concatenate the resulting \mathbf{Z} matrices, then multiply with weight matrix \mathbf{W}^o to produce the output of the layer



Transformer

Bringing The Tensors Into The Picture

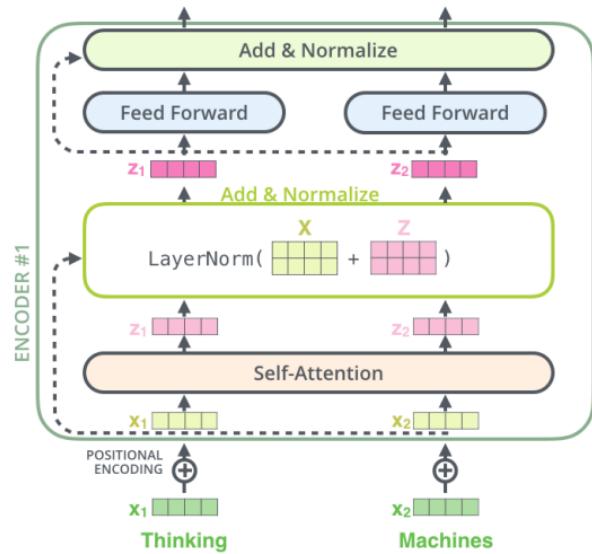
- first turn each input word into a vector using an embedding algorithm (only happens in the bottom-most encoder)
- after embedding the words in our input sequence, each of them flows through each of the two layers of the encoder
- word in each position flows through its own path in the encoder
- There are dependencies between these paths in the self-attention layer. The feed-forward layer does not have those dependencies, however, and thus the various paths can be executed in parallel while flowing through the feed-forward layer



Transformer

Residuals + Normalization

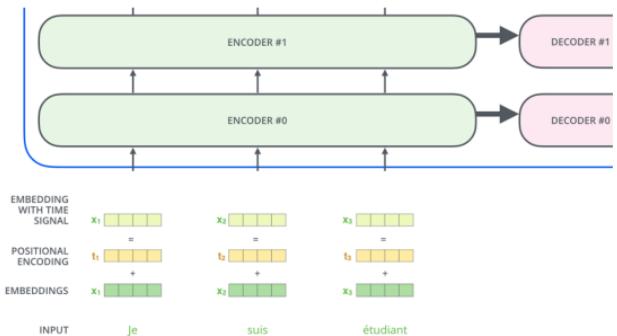
- each sub-layer (self-attention, feed-forward) in each encoder has
 - a residual connection around it
 - is followed by a layer-normalization step



Transformer

Representing The Order of The Sequence Using Positional Encoding

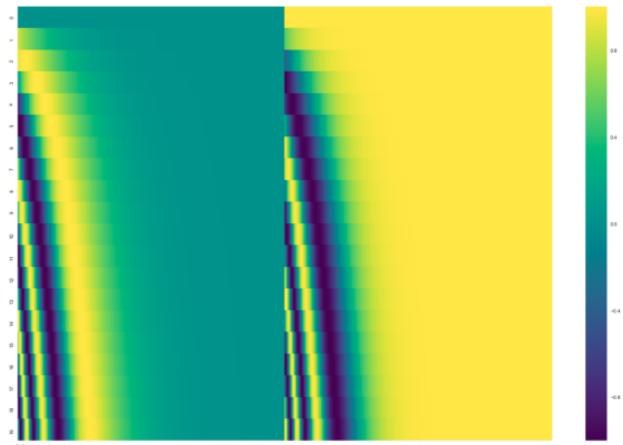
- Find a way to account for the order of the words in the input sequence
 - Transformer adds a vector to each input embedding
 - These vectors follow a specific pattern that the model learns, which helps it determine the position of each word, or the distance between different words in the sequence



To give the model a sense of the order of the words, we add positional encoding vectors -- the values of which follow a specific pattern.

Transformer

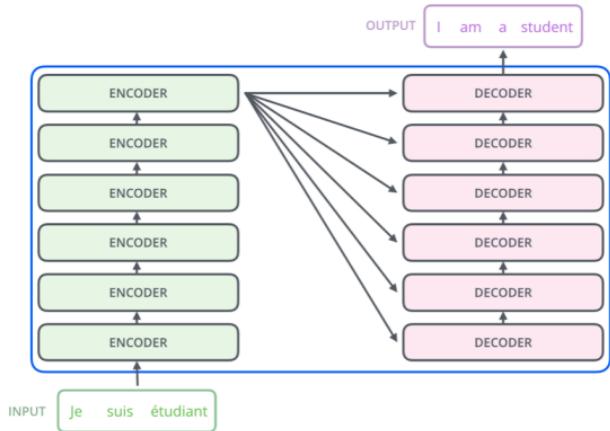
What might this pattern look like ?



A real example of positional encoding for 20 words (rows) with an embedding size of 512 (columns). You can see that it appears split in half down the center. That's because the values of the left half are generated by one function (which uses sine), and the right half is generated by another function (which uses cosine). They're then concatenated to form each of the positional encoding vectors.

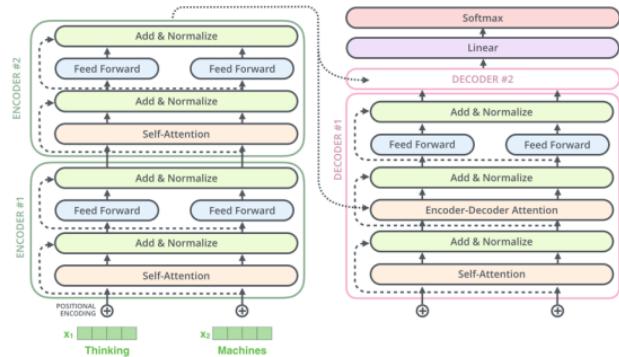
Transformer

- The encoding component is a stack of 6 encoders on top of each other
- The decoding component is a stack of decoders of the same number



Transformer

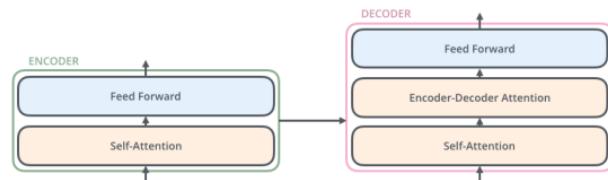
Plugin the encoder to the decoder



Transformer

The decoder has both those layers

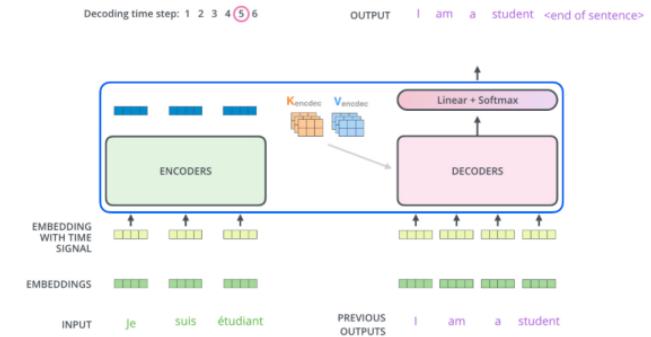
- but between them is an **attention layer** that helps the decoder focus on relevant parts of the input sentence (similar what attention does in seq2seq models)



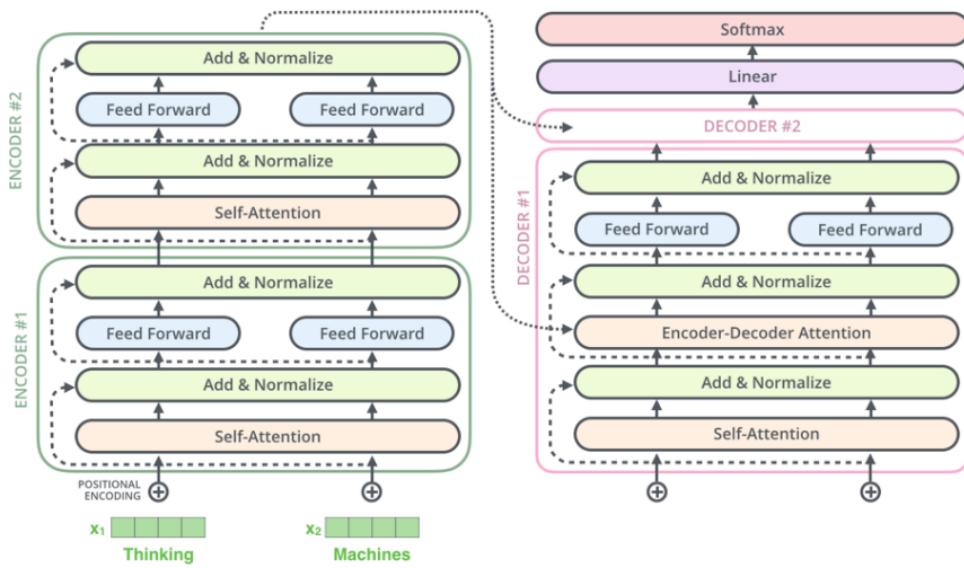
Transformer

Decoder side

- The output of the top encoder is then transformed into a set of attention vectors K and V
 - to be used by each decoder in its "encoder-decoder attention" layer which helps the decoder focus on appropriate places in the input sequence
- The output of the decoder of each step is fed to the bottom in the next step
- As for the encoder, we embed and add positional encoding to those decoder inputs
- In the decoder, the self-attention layer is only allowed to attend to earlier positions in the output sequence
- The "Encoder-Decoder Attention" layer works just like multiheaded self-attention, except it creates its Queries matrix from the layer below it, and takes the Keys and Values matrix from the output of the encoder stack.



Transformer



BERT

[J. Devlin et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", 2019]

[J. Alamar "The Illustrated BERT, ELMo, and co. (How NLP Cracked Transfer Learning)"]

- BERT ?

- a trained Transformer Encoder stack
- the first input token is supplied with a special [CLS] token. CLS here stands for Classification

