

Deep Learning

Recurrent Neural Network



Geoffroy Peeters

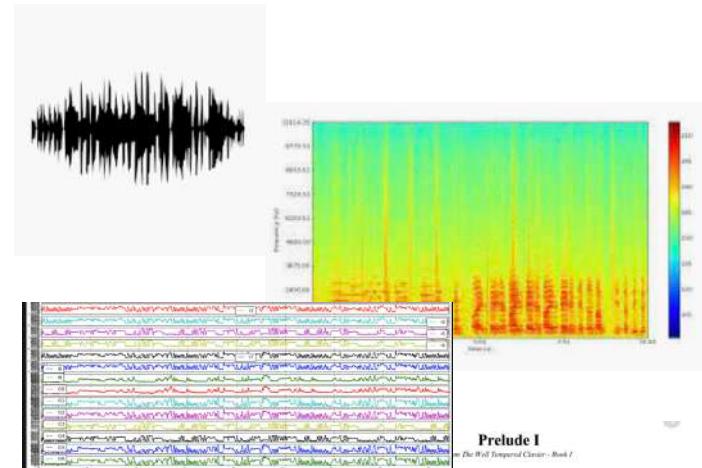
contact: geoffroy.peeters@telecom-paris.fr

Télécom-Paris, IP-Paris, France

Recurrent Neural Network for Sequential Data

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
 - Video (sequence of images)
 - Sensors (heart-beat, plane altitude)
 - Stock market
 - **Ordered:** (but not with time)
 - Text/Words
 - Genes/ADN
- **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 ?
Renie ton père et abdique ton nom; ou, si tu ne le veux pas, jure de m'aimer, et je ne serai plus une Capulet.

Roméo, à part : Dois-je l'écouter encore ou lui répondre ?

Juliette : Ton nom est mon ennemi. Tu n'es pas un Montague, tu es moi-même. Qu'est-ce qu'un Montague ? Ce n'est ni une main, ni un pied, ni un bras, si un visage, ni rien qui fasse partie d'un homme... Oh ! Sois quelque autre nom ! Qu'y a-t-il dans un nom ? Ce que nous appelons une rose embaumeraut tout auant un autre nom. Ainsi, quand Roméo ne s'appellerait plus Roméo, il conserverait encore les chères perfections qu'il

```
/* Simple HelloButton() method.
 * @version 1.0
 * @author John Doe <doe.j@example.com>
 */
HelloButton()
{
    JButton hello = new JButton( "Hello, world" );
    hello.addActionListener( new HelloButtonList
    // use the JFrame type until support for it
    // new component is finished
    JFrame frame = new JFrame( "Hello Button"
        pane = frame.getContentPane();
        hello );
    pane.add( hello );
    frame.setVisible( true );
    // display the frame
}
```



Recurrent Neural Network for Sequential Data

Notations

- Inputs:
 - $\mathbf{x}^{(t)}$: input vector \mathbf{x} at time t (of dimension $n^{[0]}$)
 - $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(T_x)}\}$: a sequence of inputs of length T_x
- Outputs:
 - $\mathbf{y}^{(t)}$: output vector \mathbf{y} at time t
 - $\{\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(T_y)}\}$: a 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

Recurrent Neural Network for Sequential Data

Notations

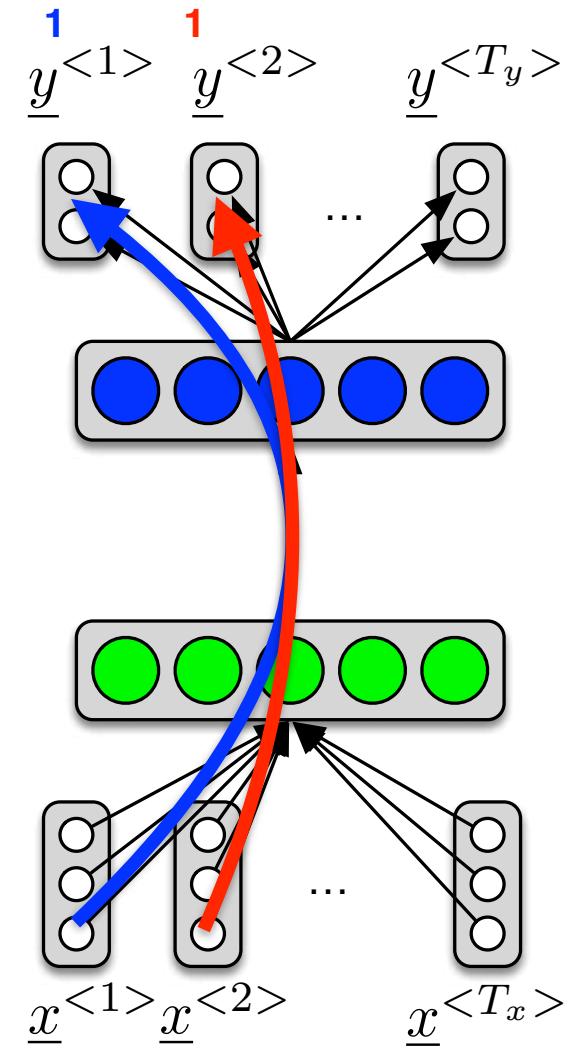
- Different types of output \mathbf{y} :
 - Continuous values $\mathbf{y} \in \mathbb{R}^{n_y}$ (Activation: linear, Loss: MSE)
 - Binary-classes $\mathbf{y} \in \{0,1\}$ (Activation: sigmoid, Loss: Binary-Cross-Entropy)
 - Multi-classes $y \in \{1\dots K\} \Rightarrow \mathbf{y} \in \{0,1\}^K$ (Activation: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

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

Recurrent Neural Network for Sequential Data

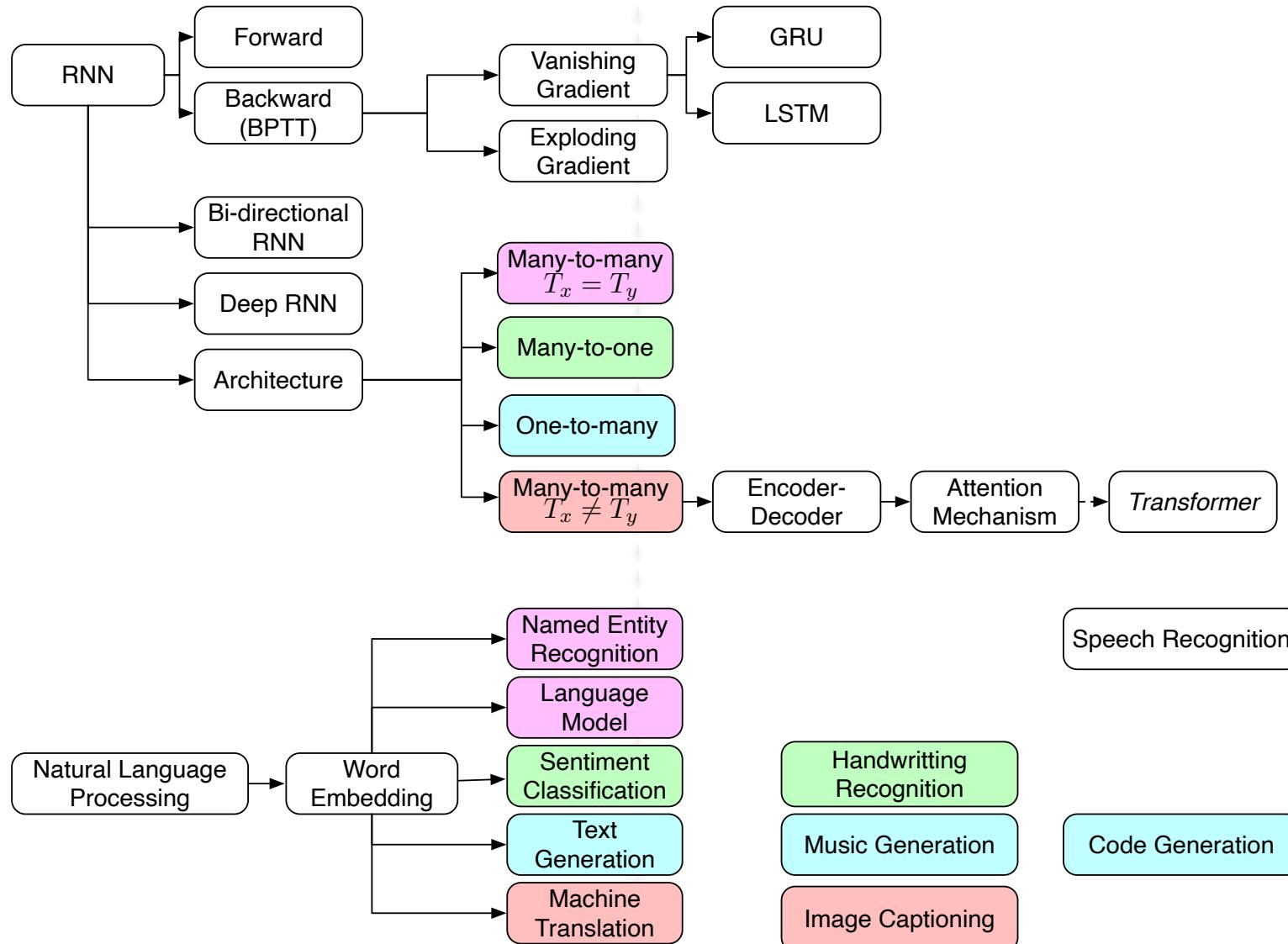
- **Why not using Multi-Layer-Perceptron ?**

- We could indeed represent an input sequence $\{\underline{x}^{(1)}, \dots, \underline{x}^{(T_x)}\}$ as a large-vector of dimension $(n^{[0]} T_x)$
 - but requires a huge input dimension
 - but $T_x^{(i)}$ can be different for each input sequence,
⇒ the dimension of the input would be different for each $(n^{[0]} T_x^{(i)})$
- Solution: zero-padding ? $(n^{[0]} \max_i T_x^{(i)})$
 - still requires a huge input dimension
 - 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)

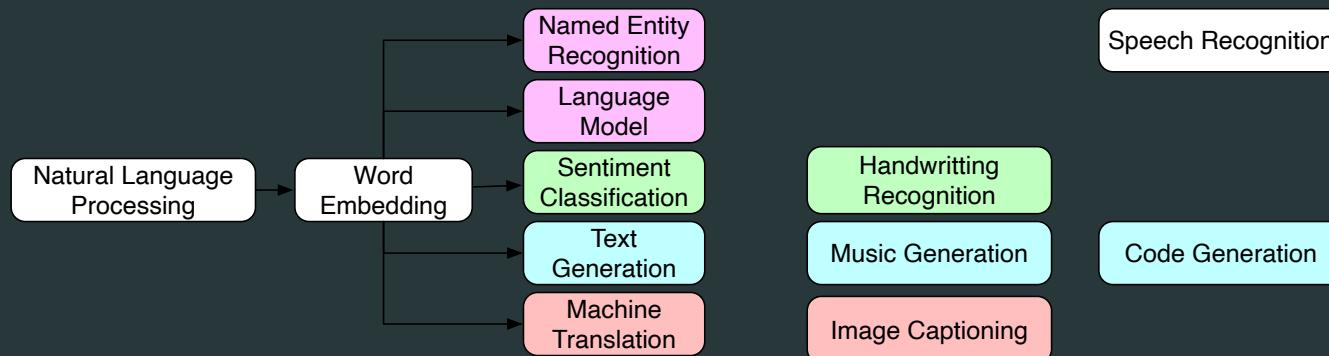
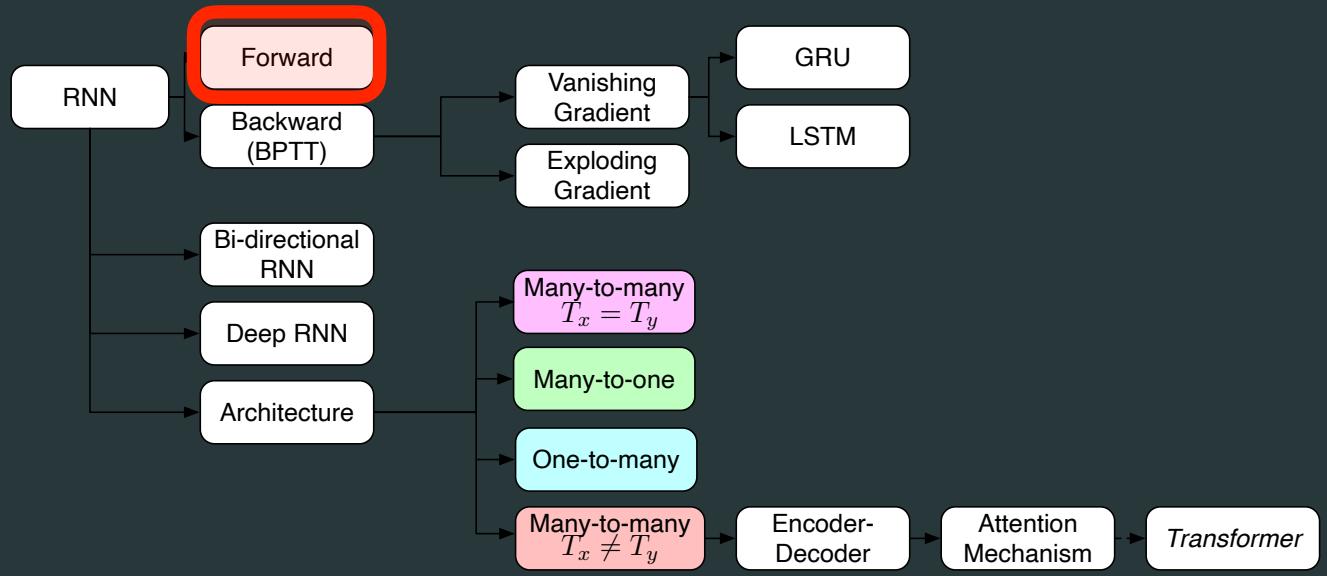


EmmanueEmmanuel

Overview

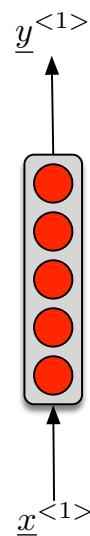


Recurrent Neural Network for Sequential Data



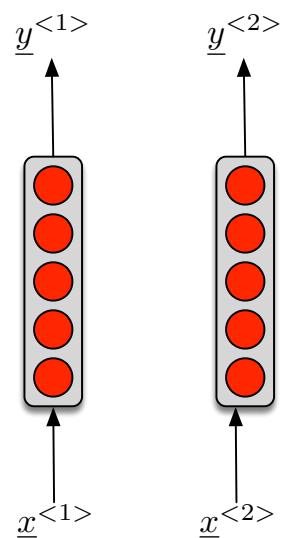
Recurrent Neural Network for Sequential Data

What is an RNN ?



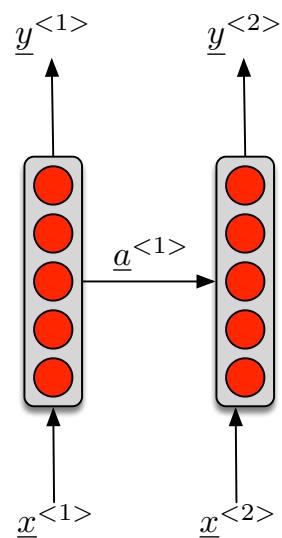
Recurrent Neural Network for Sequential Data

What is an RNN ?



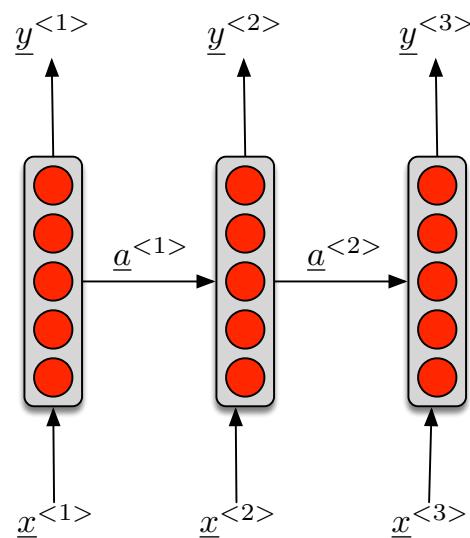
Recurrent Neural Network for Sequential Data

What is an RNN ?



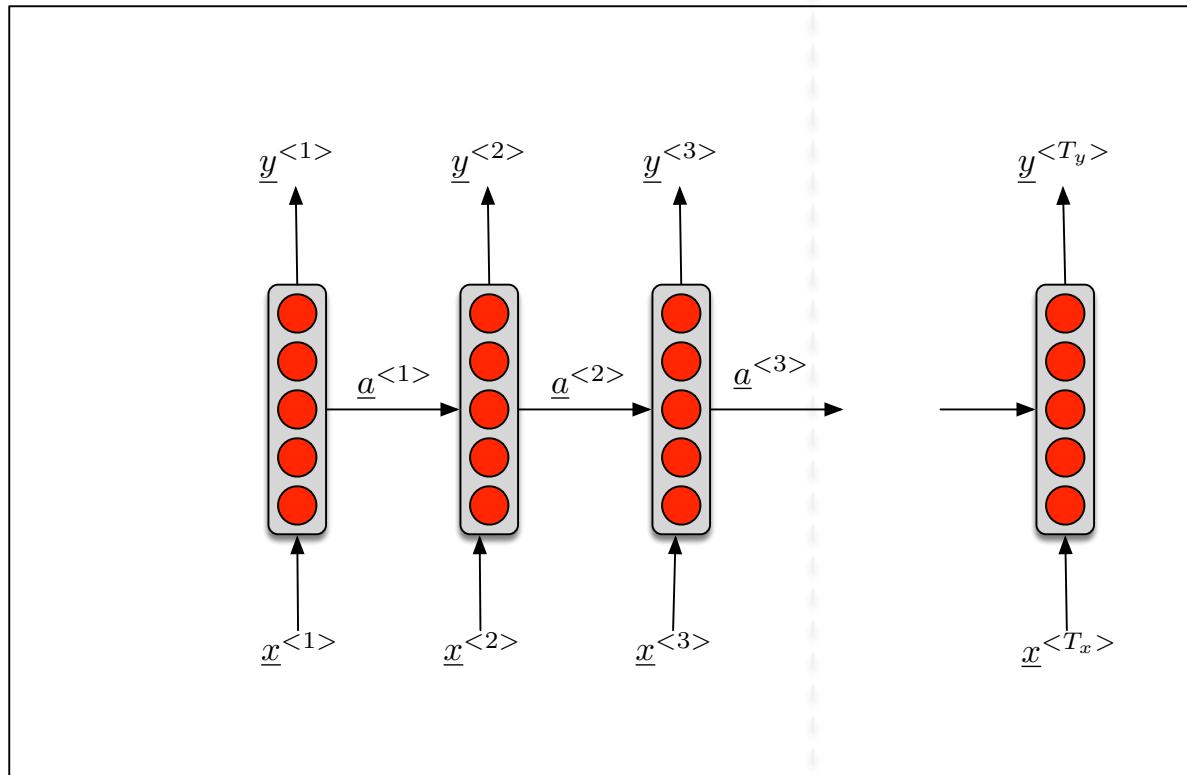
Recurrent Neural Network for Sequential Data

What is an RNN ?



Recurrent Neural Network for Sequential Data

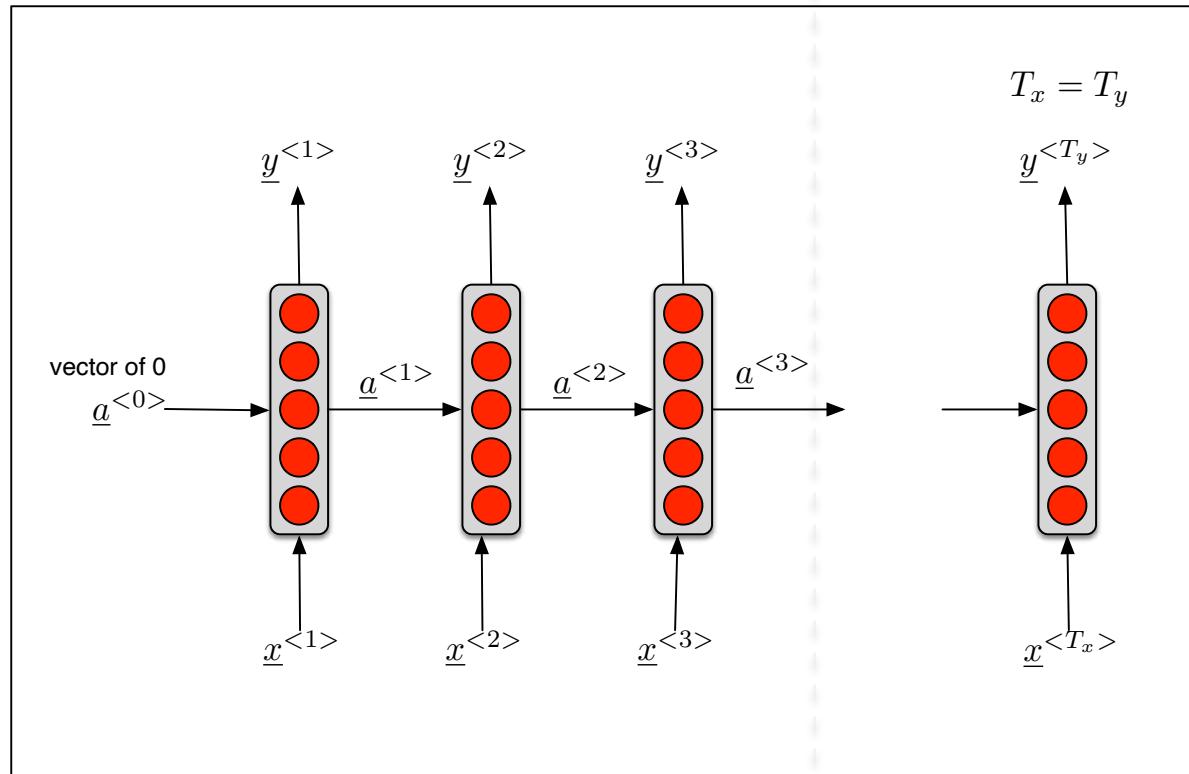
What is an RNN ?



- At each time step $\langle t \rangle$, the RNN gets the activation $a^{\langle t-1 \rangle}$ from the previous time $\langle t-1 \rangle$

Recurrent Neural Network for Sequential Data

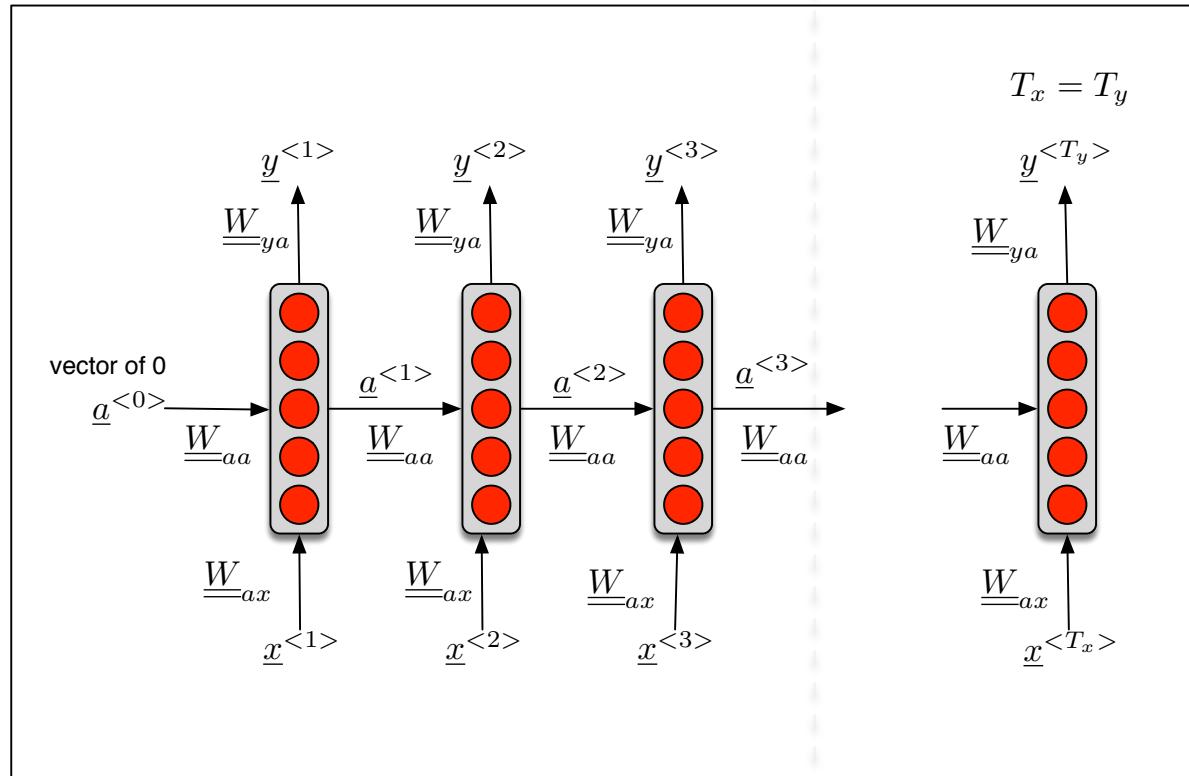
What is an RNN ?



- Initialise the first activation with $\underline{a}^{<0>} = 0$

Recurrent Neural Network for Sequential Data

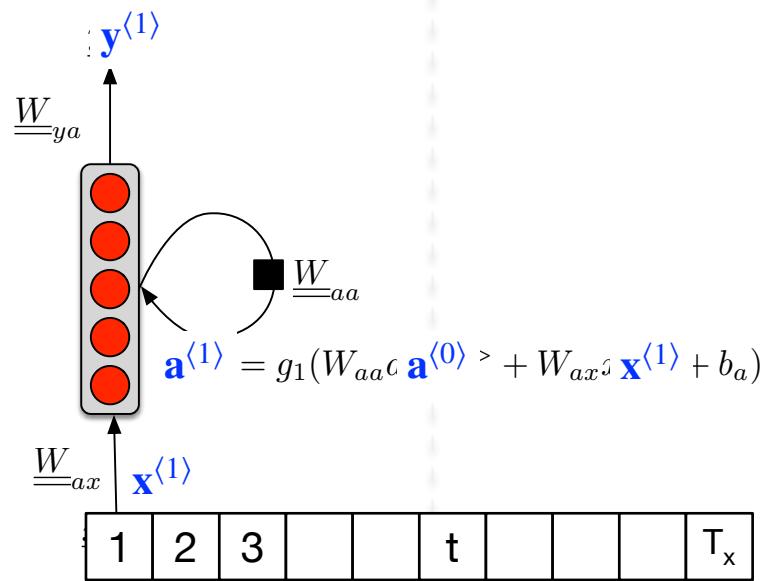
What is an RNN ?



- Weight matrices: \mathbf{W}_{ax} from $\mathbf{x}^{(t)}$ $\rightarrow \mathbf{a}^{(t)}$, \mathbf{W}_{aa} from $\mathbf{a}^{(t-1)}$ $\rightarrow \mathbf{a}^{(t)}$, \mathbf{W}_{ay} from $\mathbf{a}^{(t)}$ $\rightarrow \mathbf{y}^{(t)}$
- They are shared from all time-steps

Recurrent Neural Network for Sequential Data

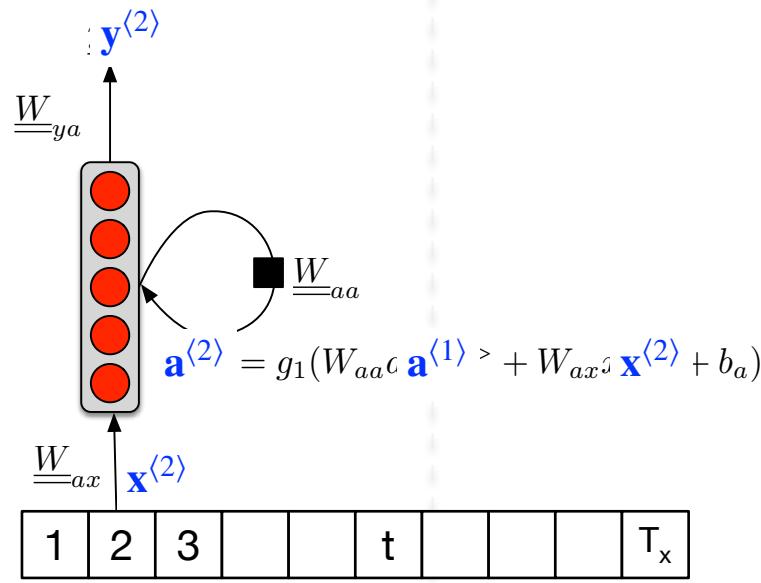
What is an RNN ?



- "**Rolled**" version of RNN (the real one, only one network)

Recurrent Neural Network for Sequential Data

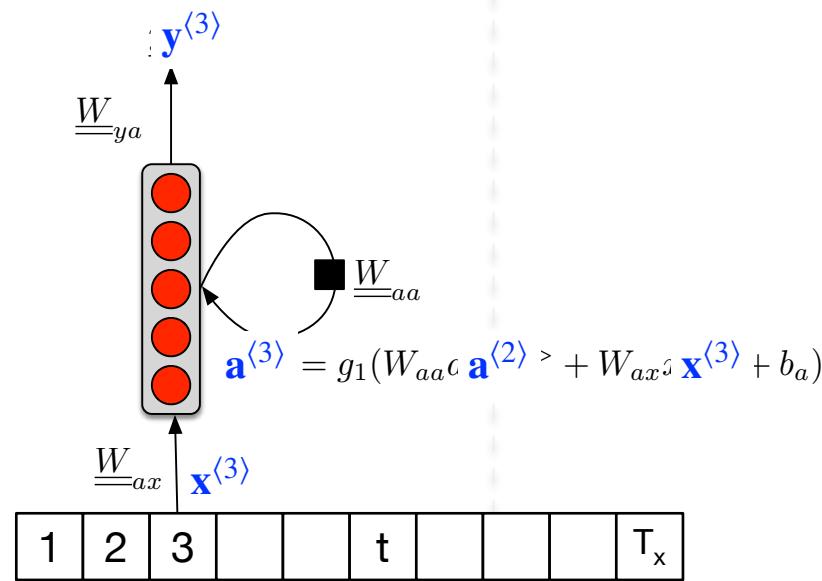
What is an RNN ?



- "**Rolled**" version of RNN (the real one, only one network)

Recurrent Neural Network for Sequential Data

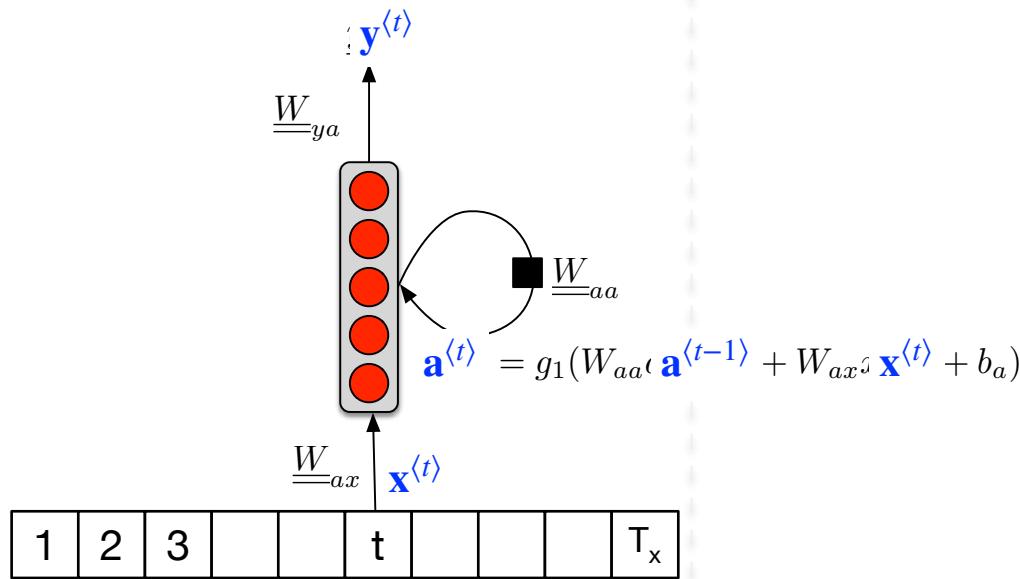
What is an RNN ?



- "Rolled" version of RNN (the real one, only one network)

Recurrent Neural Network for Sequential Data

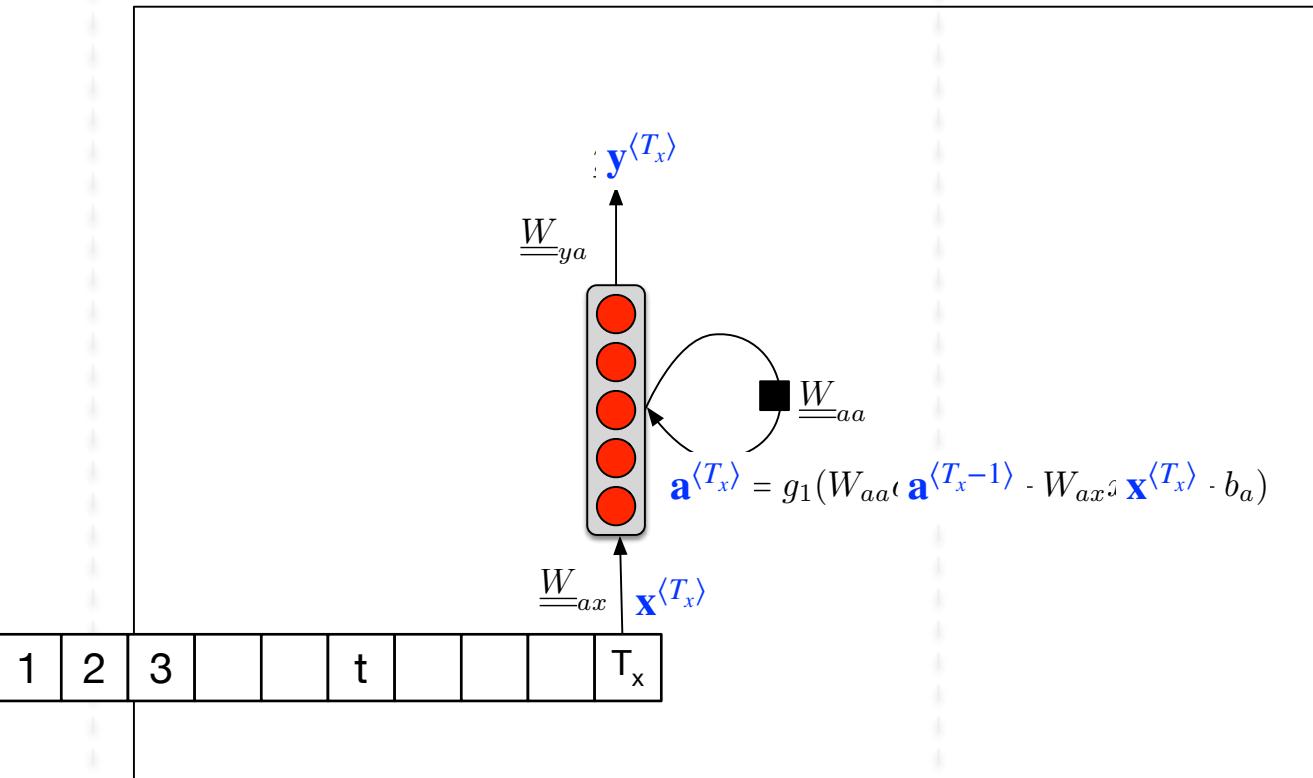
What is an RNN ?



- "Rolled" version of RNN (the real one, only one network)

Recurrent Neural Network for Sequential Data

What is an RNN ?



- "Rolled" version of RNN (the real one, only one network)

Recurrent Neural Network for Sequential Data

What is an RNN ?

- **Forward propagation**

$$\mathbf{a}^{(0)} = 0$$

$$\mathbf{a}^{(1)} = g_1 \left(\mathbf{W}_{aa} \mathbf{a}^{(0)} + \mathbf{W}_{ax} \mathbf{x}^{(1)} + \mathbf{b}_a \right)$$

$$\hat{\mathbf{y}}^{(1)} = g_2 \left(\mathbf{W}_{ya} \mathbf{a}^{(1)} + \mathbf{b}_y \right)$$

...

$$\mathbf{a}^{(t)} = g_1 \left(\mathbf{W}_{aa} \mathbf{a}^{(t-1)} + \mathbf{W}_{ax} \mathbf{x}^{(t)} + \mathbf{b}_a \right)$$

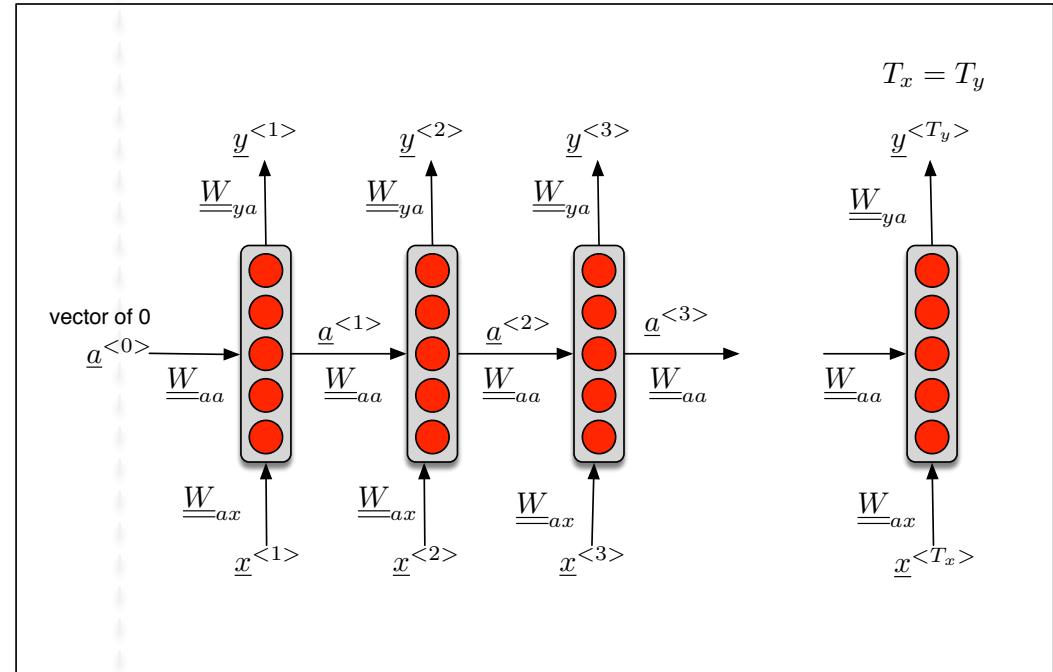
$$\hat{\mathbf{y}}^{(t)} = g_2 \left(\mathbf{W}_{ya} \mathbf{a}^{(t)} + \mathbf{b}_y \right)$$

- **Compact notations**

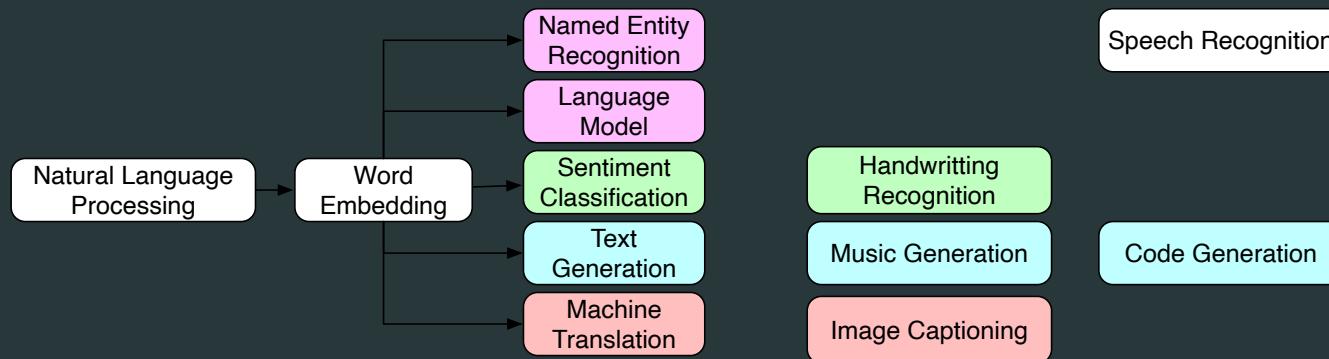
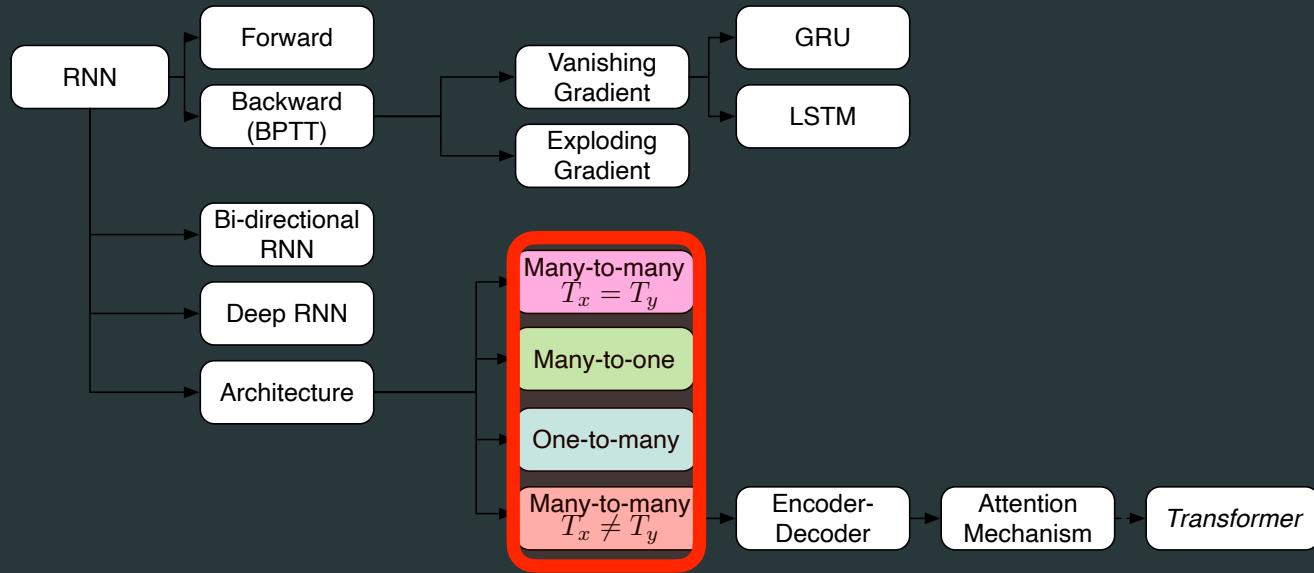
$$[\mathbf{W}_{aa} \mid \mathbf{W}_{ax}] \begin{bmatrix} \mathbf{a}^{(t-1)} \\ \mathbf{x}^{(t)} \end{bmatrix}$$

$$\mathbf{a}^{(t)} = g_1 \left(\mathbf{W}_a [\mathbf{a}^{(t-1)}; \mathbf{x}^{(t)}] + \mathbf{b}_a \right)$$

$$\hat{\mathbf{y}}^{(t)} = g_2 \left(\mathbf{W}_y \mathbf{a}^{(t)} + \mathbf{b}_y \right)$$



Various types of sequential data



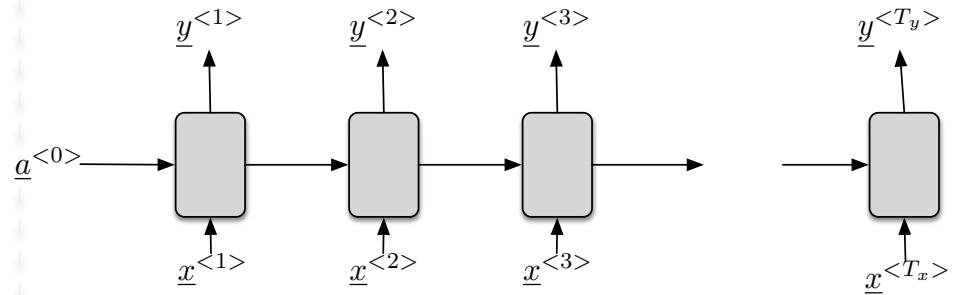
Various types of sequential data

Many-to-Many : $T_x = T_y$.

Input \underline{x}	Output \underline{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.

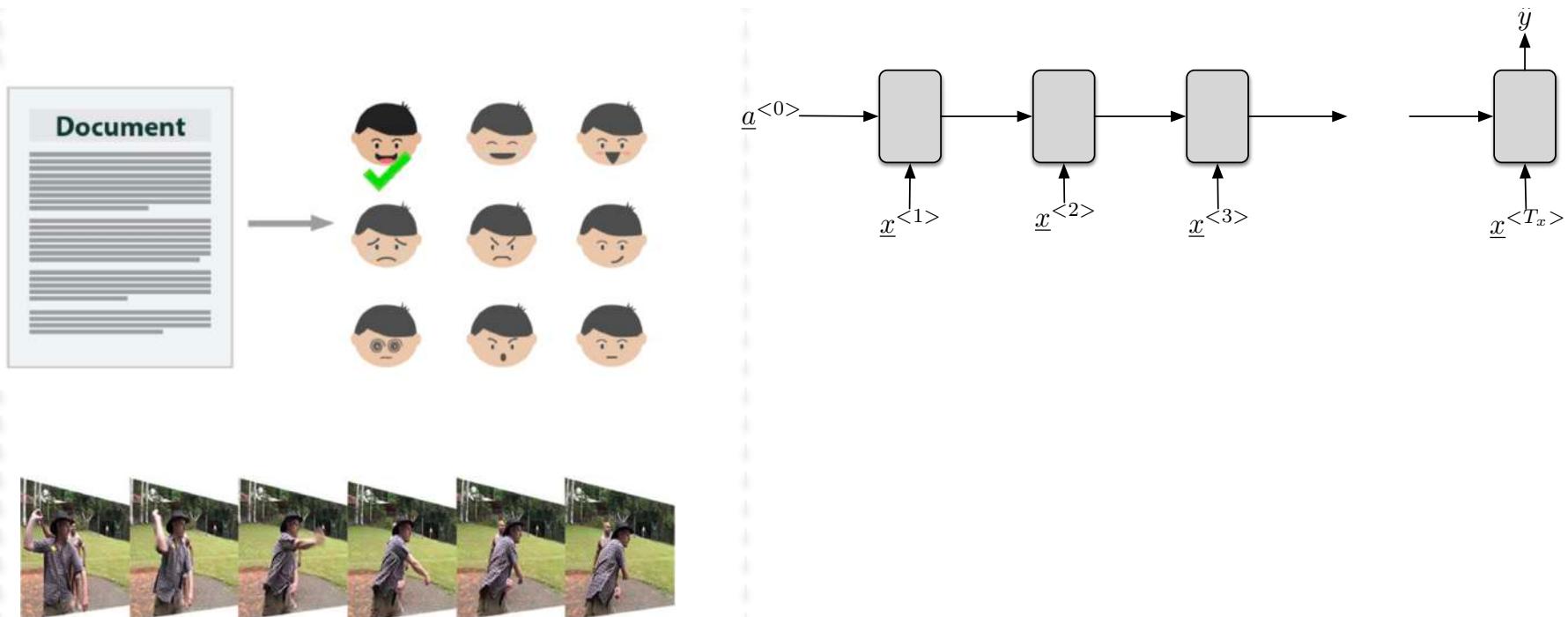
Tag colours:
LOCATION TIME PERSON ORGANIZATION MONEY PERCENT DATE



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

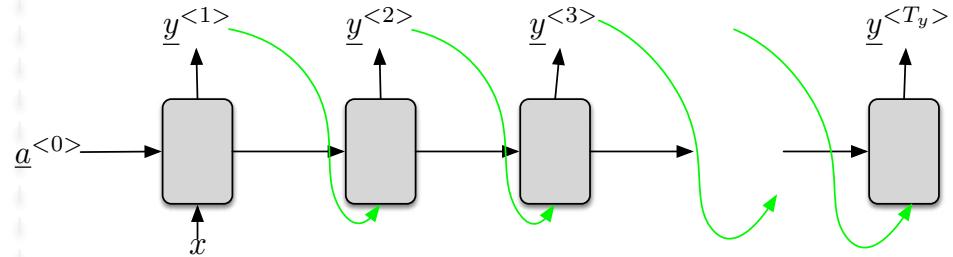


Various types of sequential data

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 Guangzhang's sovereignty. His generals were the powerful ruler of the Portugal in the [[Protestant Imminenes]], 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' 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 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.)



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

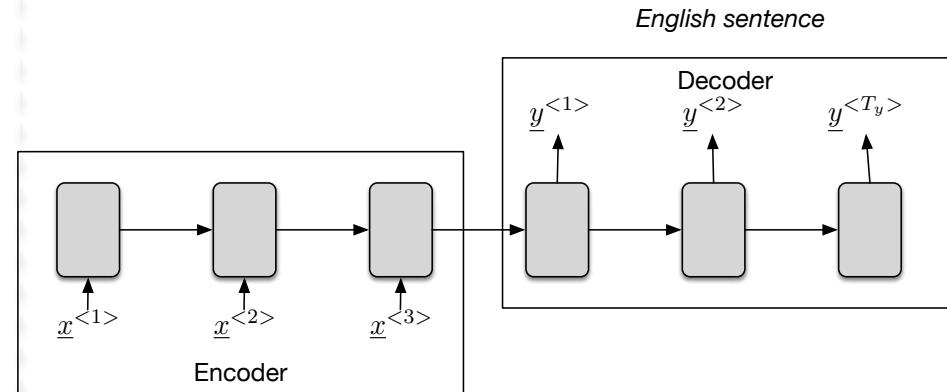
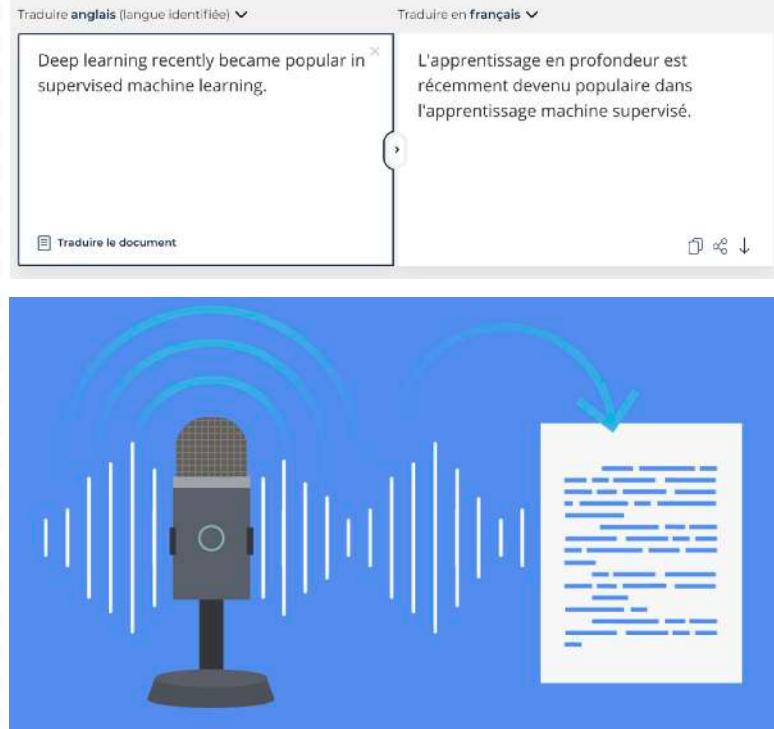
Russell E Glaue



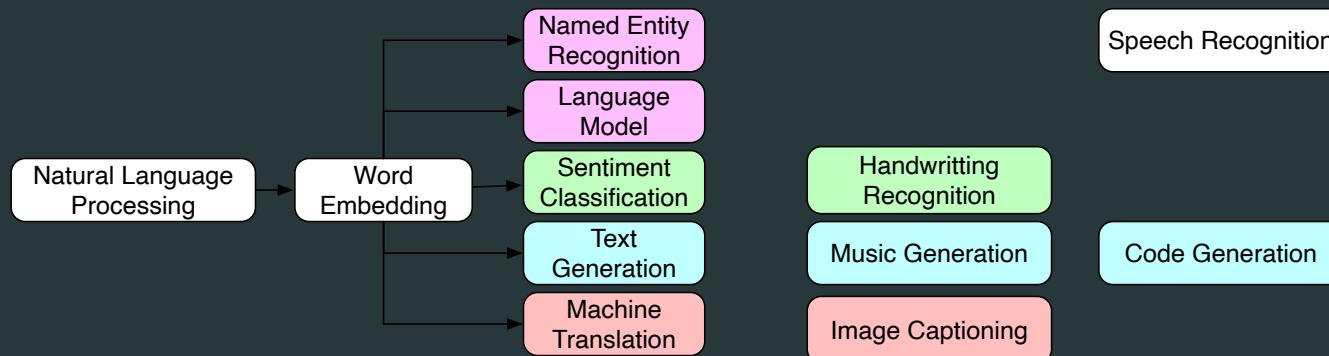
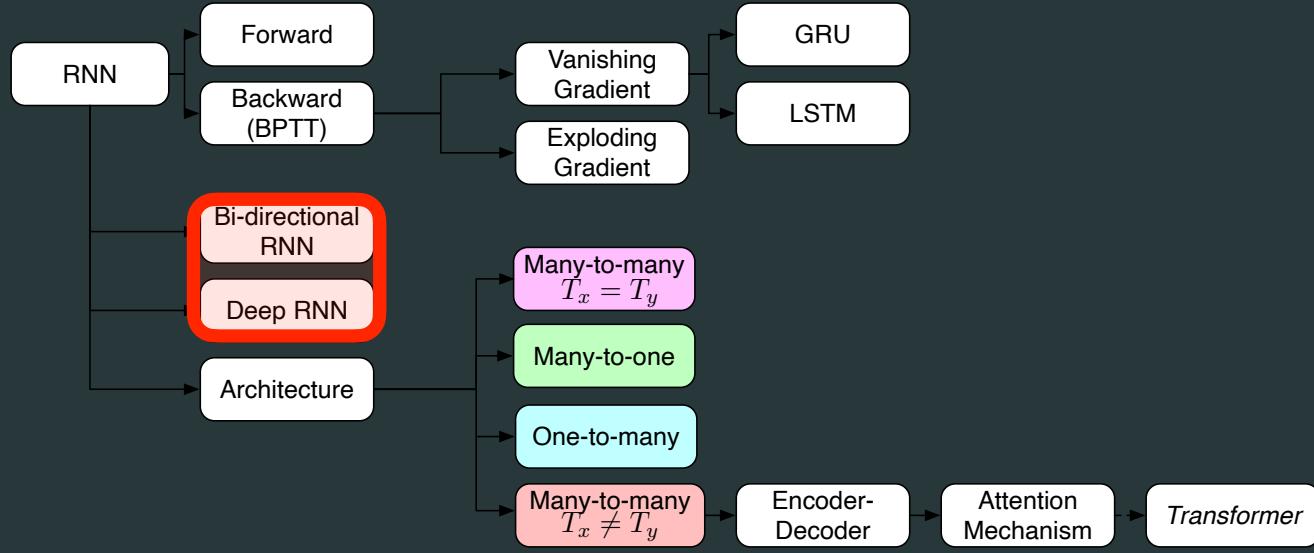
Various types of sequential data

Many-to-Many : $T_x \neq T_y$.

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

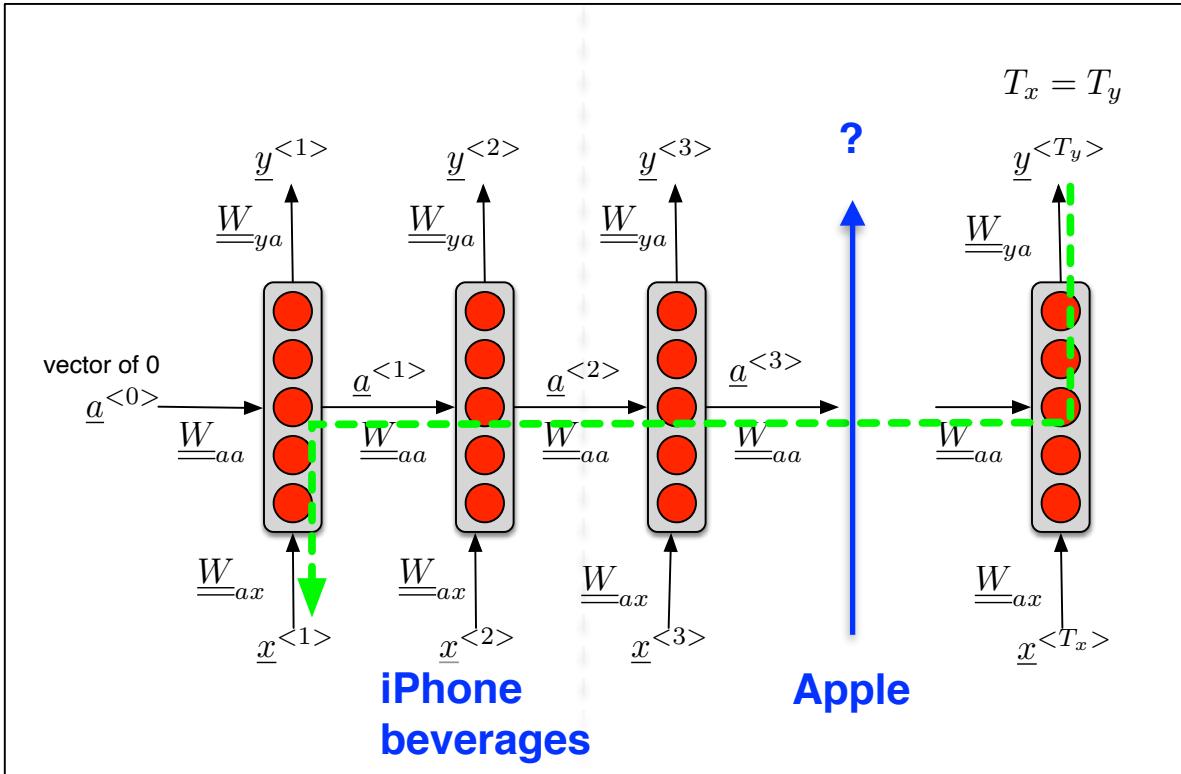


Different architectures



Different architectures

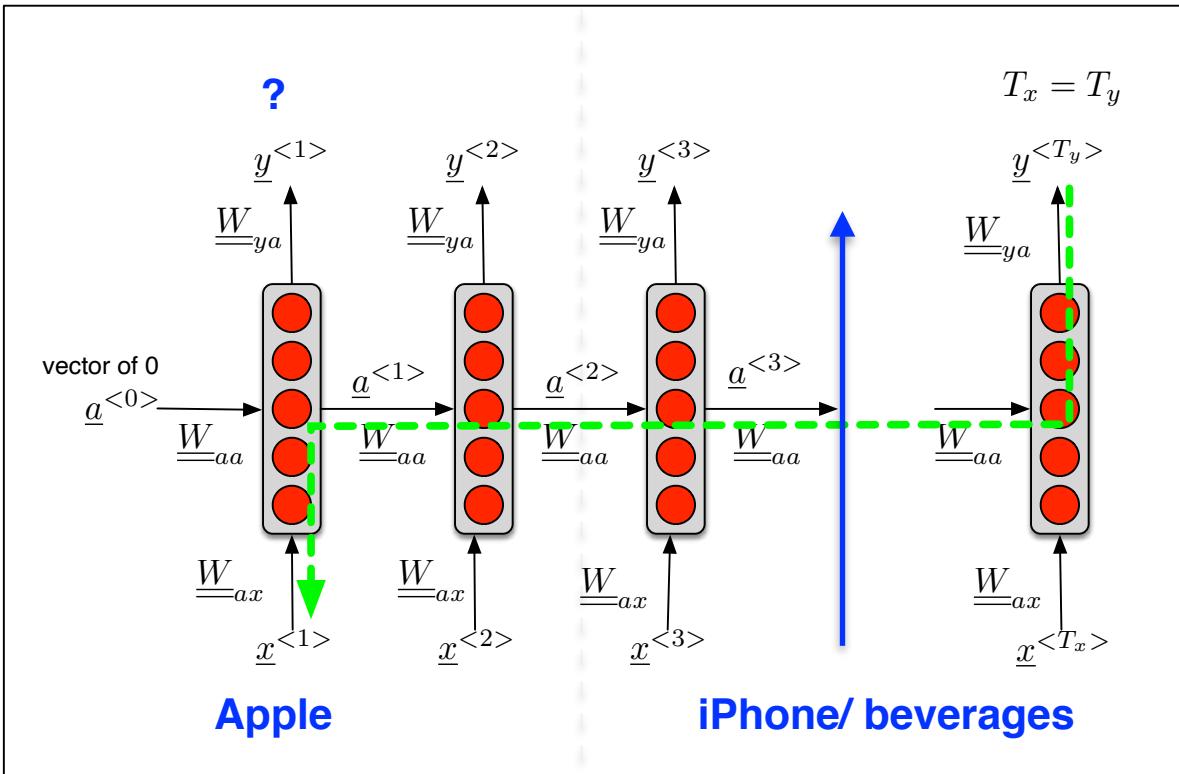
RNN



- RNN scans data from left to right \Rightarrow can do skip connections but only 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)

Different architectures

RNN



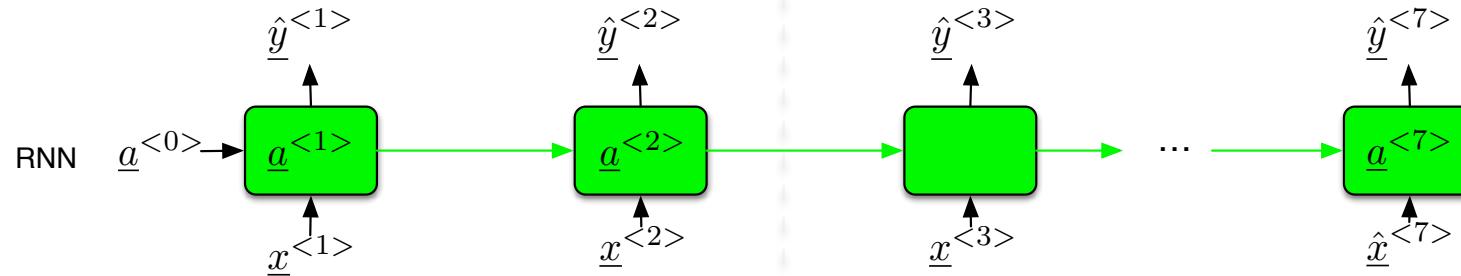
- RNN scans data from left to right \Rightarrow can do skip connections but only in the past
- Question: is "**Apple**" a named entity ?
 - "**Apple** stores sale the new **iPhone**."
 - "**Apple** juice is a refreshing and healthy **beverages**."
- \Rightarrow we cannot predict (since RNN does not depend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

Different architectures

Bi-directional RNN

– Standard RNN:

- Left-Right hidden states: $a^{(t)}$

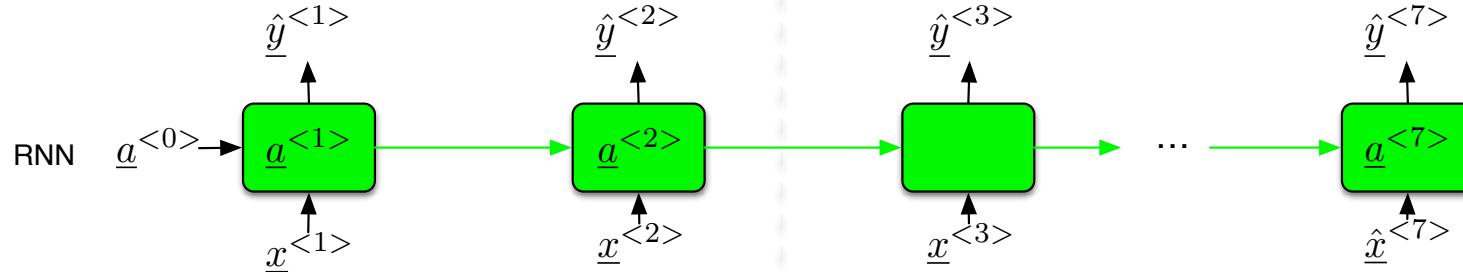


Different architectures

Bi-directional RNN

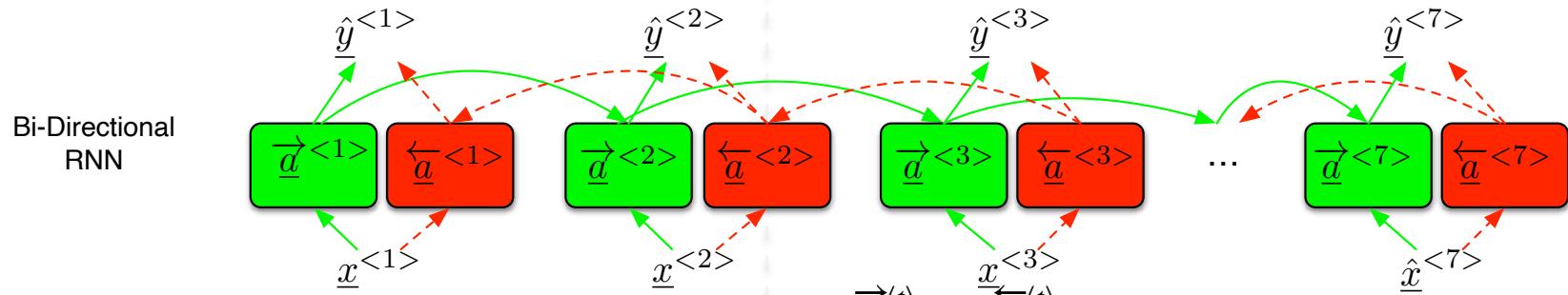
– Standard RNN:

- Left-Right hidden states: $\underline{a}^{(t)}$



– Bi-directional RNN:

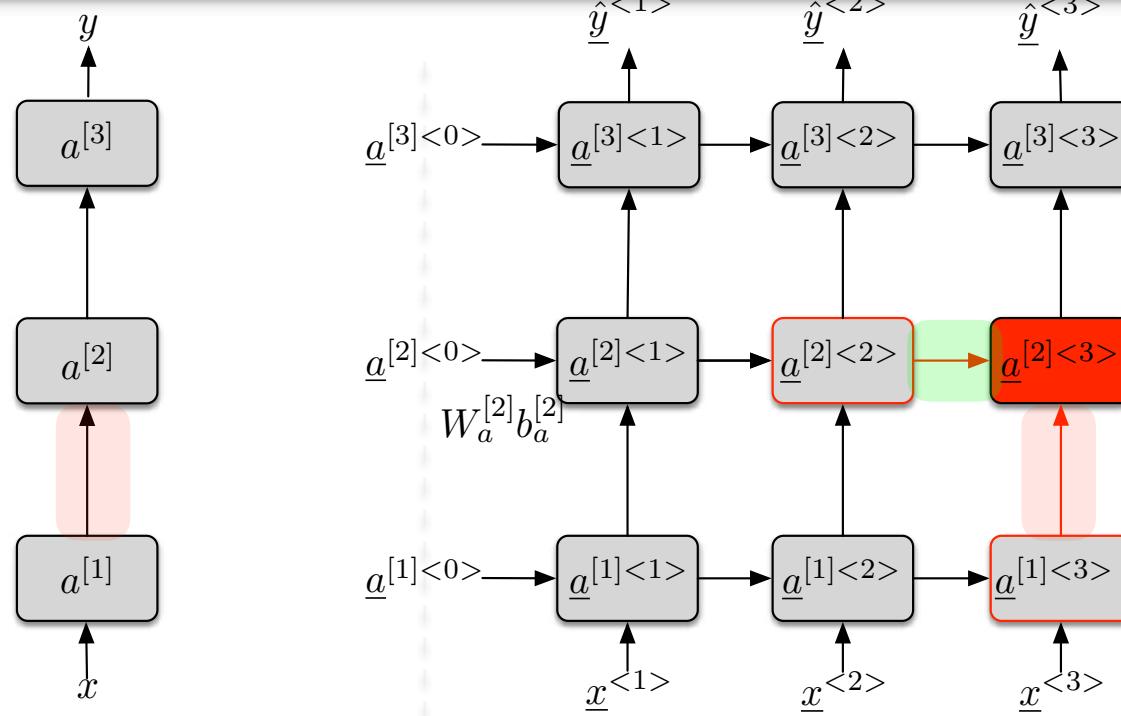
- Left-Right hidden states: $\overrightarrow{a}^{(t)}$
- Right-Left hidden states: $\overleftarrow{a}^{(t)}$



- The prediction is done from the concatenation of $\overrightarrow{a}^{(t)}$ and $\overleftarrow{a}^{(t)}$
- $\hat{y}^{(t)} = g(W_y[\overrightarrow{a}^{(t)}, \overleftarrow{a}^{(t)}] + b_y)$

Different architectures

Deep-RNN

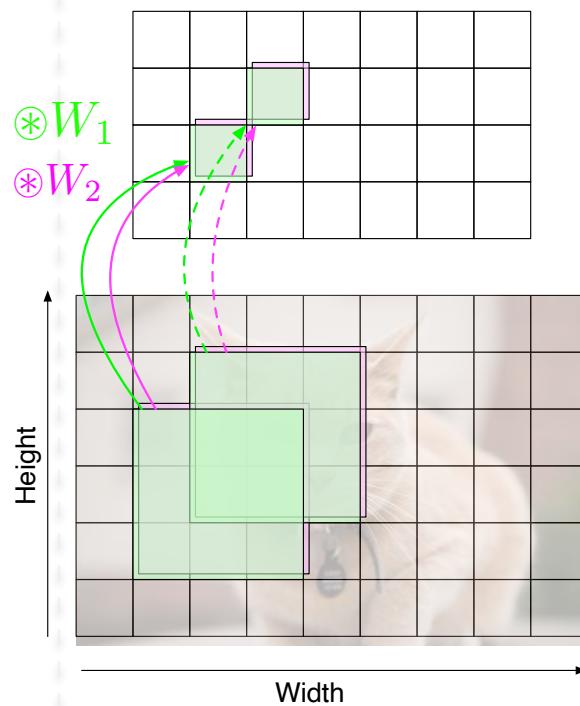


- For the layer $[1]$, the inputs of the cell at time $\langle t \rangle$ are
 - the input at the current time $\langle t \rangle$: $x^{(t)} = a^{[0]\langle t \rangle}$
 - the value of the cell at the previous time $\langle t - 1 \rangle$: $a^{[1]\langle t-1 \rangle}$
 - For the layer $[l]$, the inputs of the cell at time $\langle t \rangle$ are
 - the value of the cell of the previous layer $[l - 1]$ at the current time $\langle t \rangle$: $a^{[l-1]\langle t \rangle}$
 - the value of the cell of the current layer $[l]$ at the previous time $\langle t - 1 \rangle$: $a^{[l]\langle t-1 \rangle}$
- $$a^{[l]\langle t \rangle} = g \left(W_a^{[l]} [a^{[l]\langle t-1 \rangle}, a^{[l-1]\langle t \rangle}] + b_a^{[l]} \right)$$

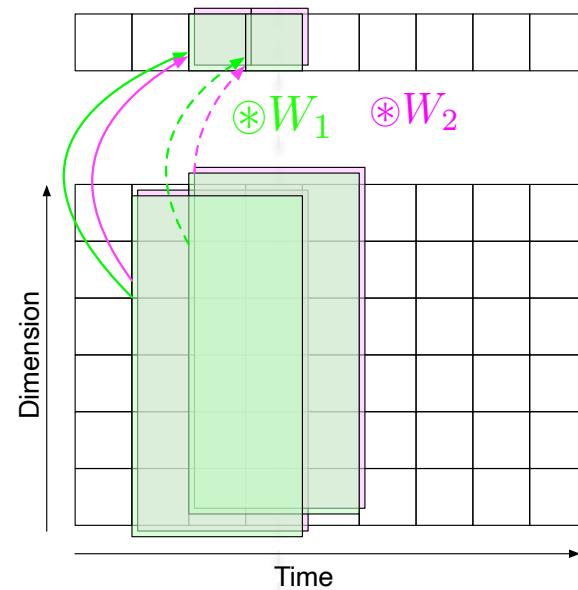
Different architectures

Using 1D convolution instead of RNN

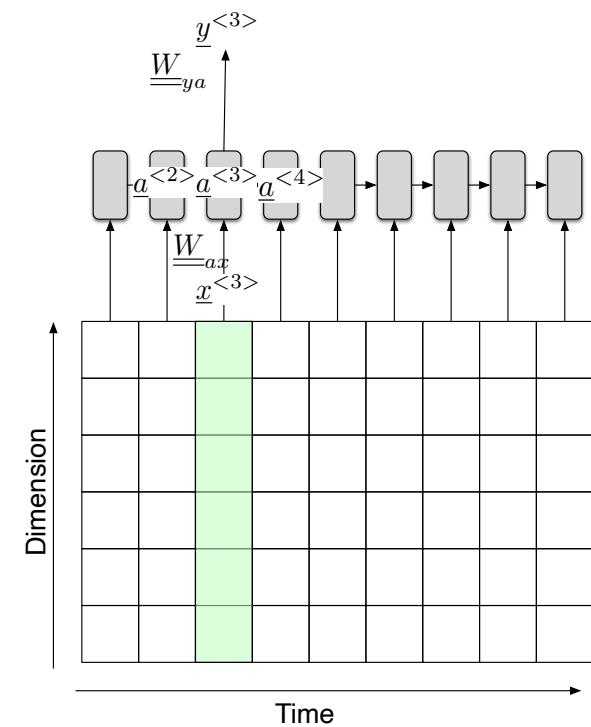
Conv2D



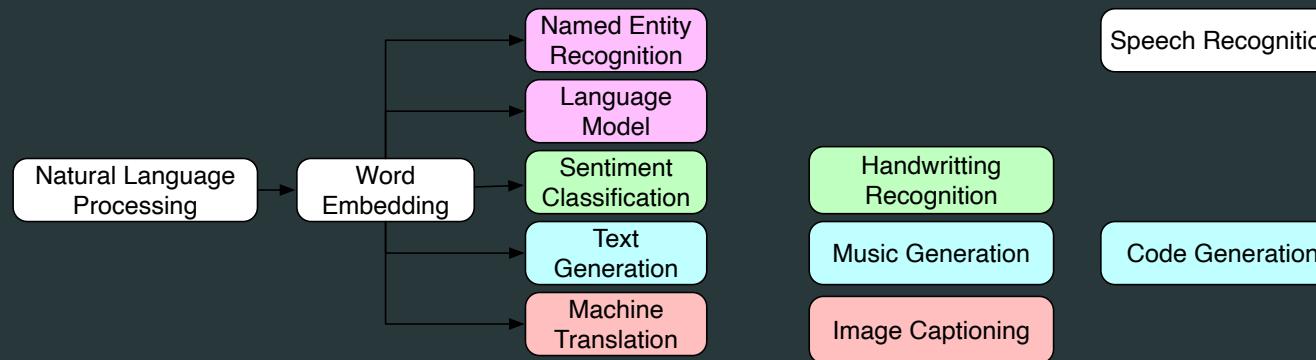
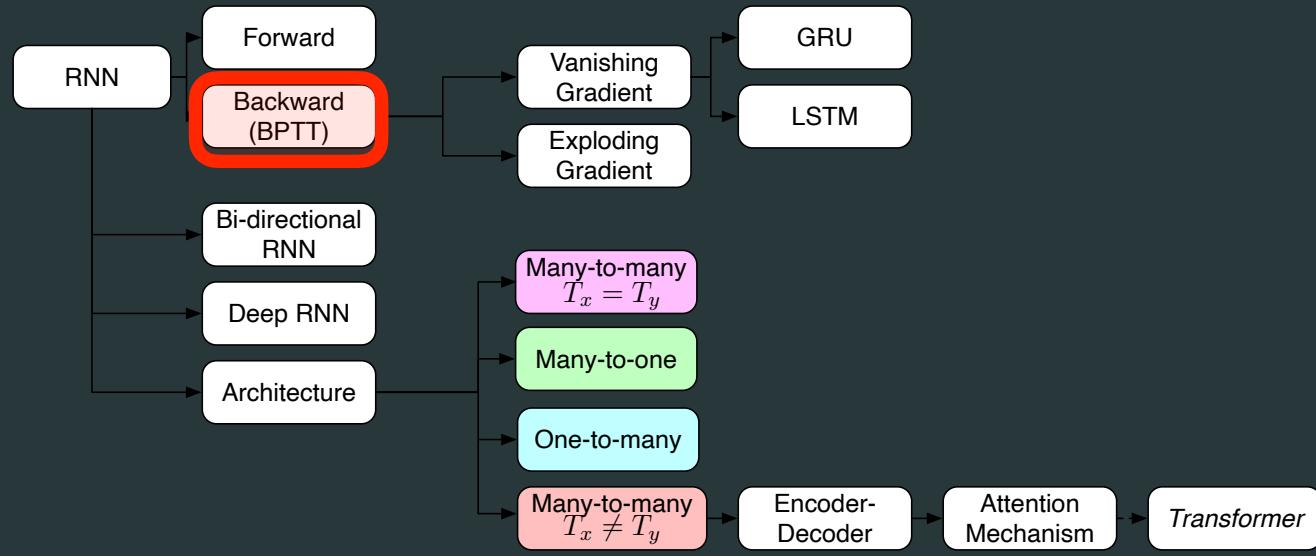
Conv1D



RNN



Back Propagation Through Time (BPTT)



Back Propagation Through Time (BPTT)

Forward pass

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

– Forward

$$\mathbf{a}^{(t)} = g_a \left(\mathbf{W}_{ax} \mathbf{x}^{(t)} + \mathbf{W}_{aa} \mathbf{a}^{(t-1)} + \mathbf{b}_a \right)$$

$$\hat{y}^{(t)} = g_y \left(\mathbf{W}_{ya} \mathbf{a}^{(t)} + \mathbf{b}_y \right)$$

– Loss

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

- $y^{(t)}$: ground-truth

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

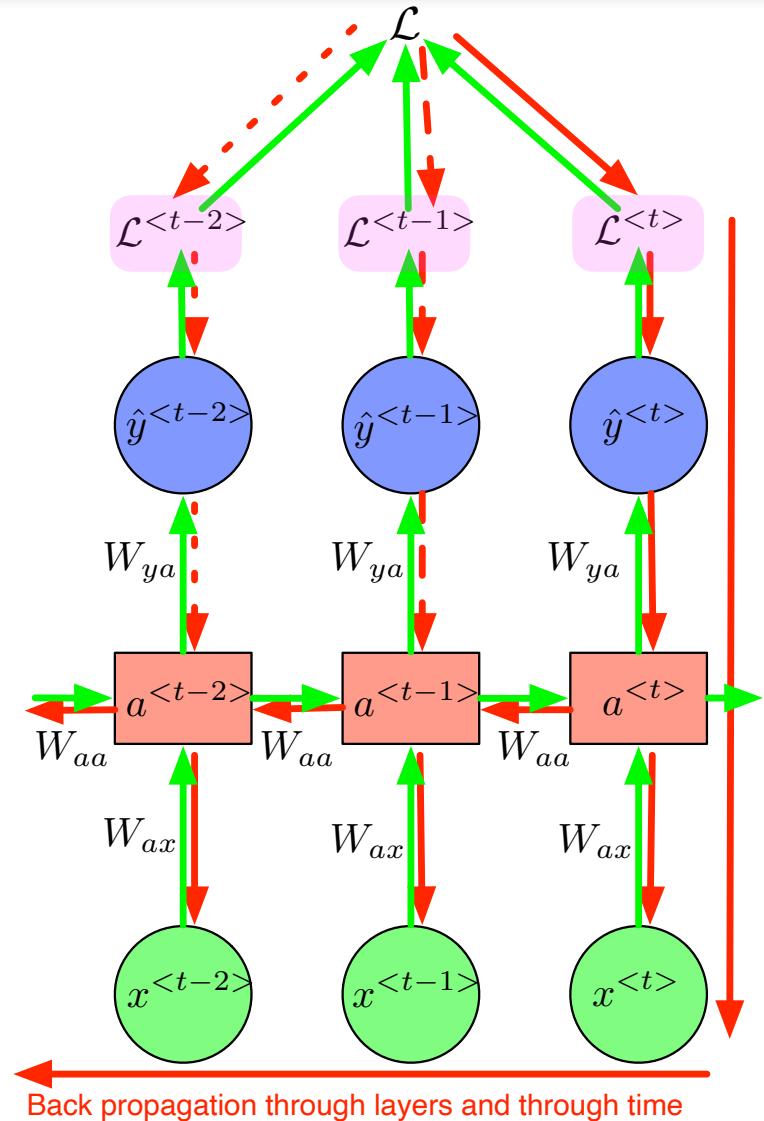
- Total loss

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

– Gradients ?

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

- All weights are shared across time steps !!!



Back Propagation Through Time (BPTT)

Backward: $\frac{\partial \mathcal{L}}{\partial W_{ya}}$

- Linearity

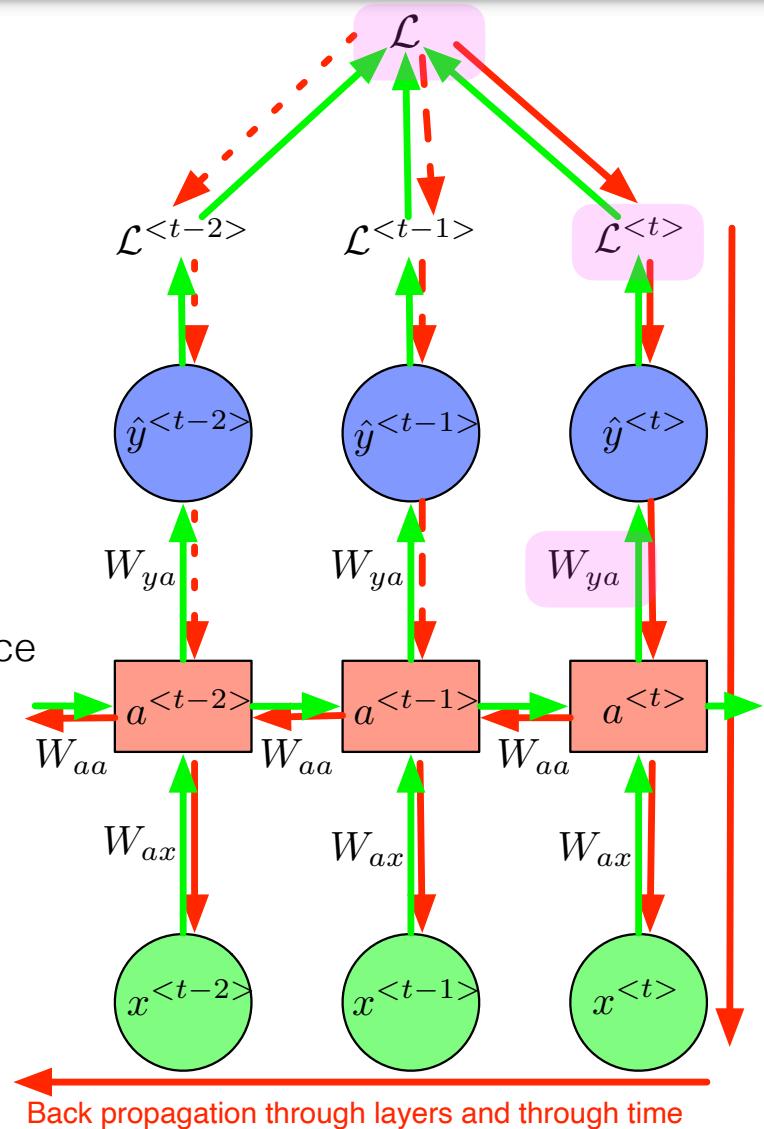
$$\frac{\partial \mathcal{L}}{\partial W_{ya}} = \sum_{t=1}^T \frac{\partial \mathcal{L}^{(t)}}{\partial W_{ya}}$$

- How much varying W_{ya} affect $\mathcal{L}^{(t)}$? : chain rule

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

- Very simple since $\hat{y}^{(t)}$ depends on W_{ya} on only place
(indeed $a^{(t)}$ does not depend on W_{ya})

$$\hat{y}^{(t)} = g_y \left(\mathbf{W}_{ya} \mathbf{a}^{(t)} + \mathbf{b}_y \right)$$



Back Propagation Through Time (BPTT)

Backward: $\frac{\partial \mathcal{L}}{\partial W_{aa}} .$

- Linearity

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

- How much varying W_{aa} affect $\mathcal{L}^{(t)}$? no chain rule

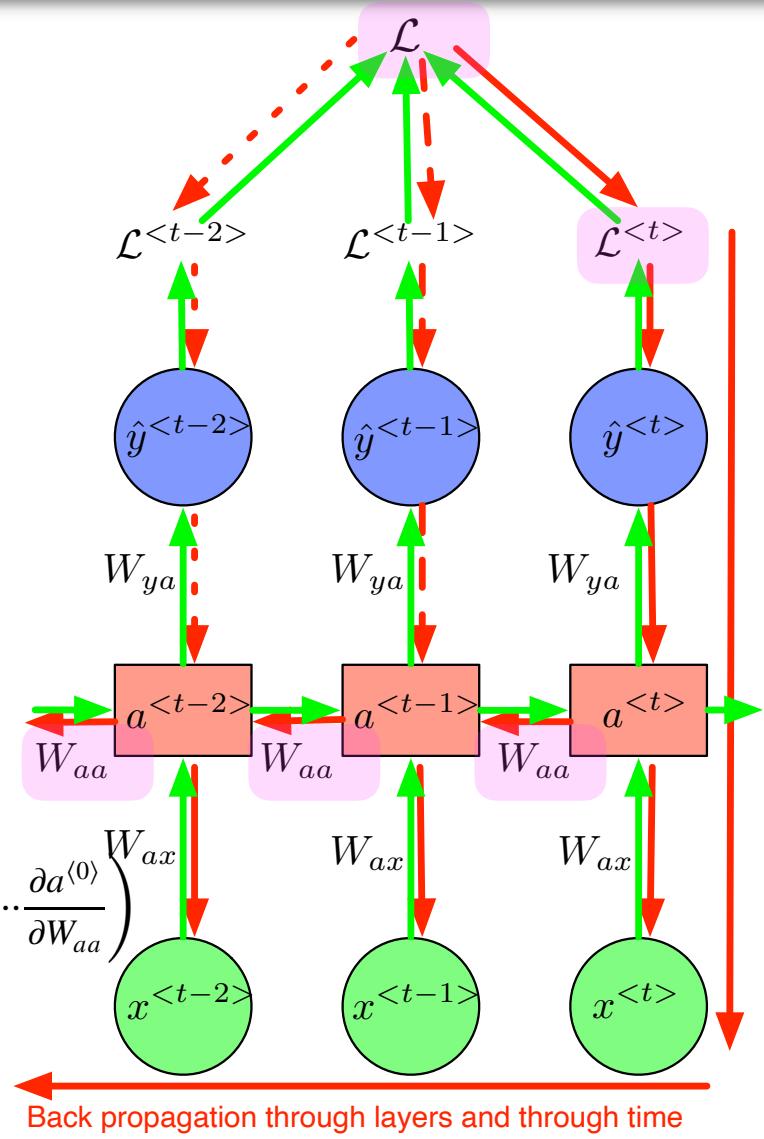
$$\frac{\partial \mathcal{L}^{(t)}}{\partial W_{aa}} \neq \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial a^{(t)}} \frac{\partial a^{(t)}}{\partial W_{aa}}$$

- because $a^{(t)}$ also depends on W_{aa} through $a^{(t-1)}$ which itself depends on ...

$$\mathbf{a}^{(t)} = g_a \left(\mathbf{W}_{ax} \mathbf{x}^{(t)} + \mathbf{W}_{aa} \mathbf{a}^{(t-1)} + \mathbf{b}_a \right)$$

- We use a formula of total derivative

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



Back Propagation Through Time (BPTT)

Backward: $\frac{\partial \mathcal{L}}{\partial W_{ax}}$.

- Linearity

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

- How much varying W_{ax} affect $\mathcal{L}^{(t)}$? no chain rule

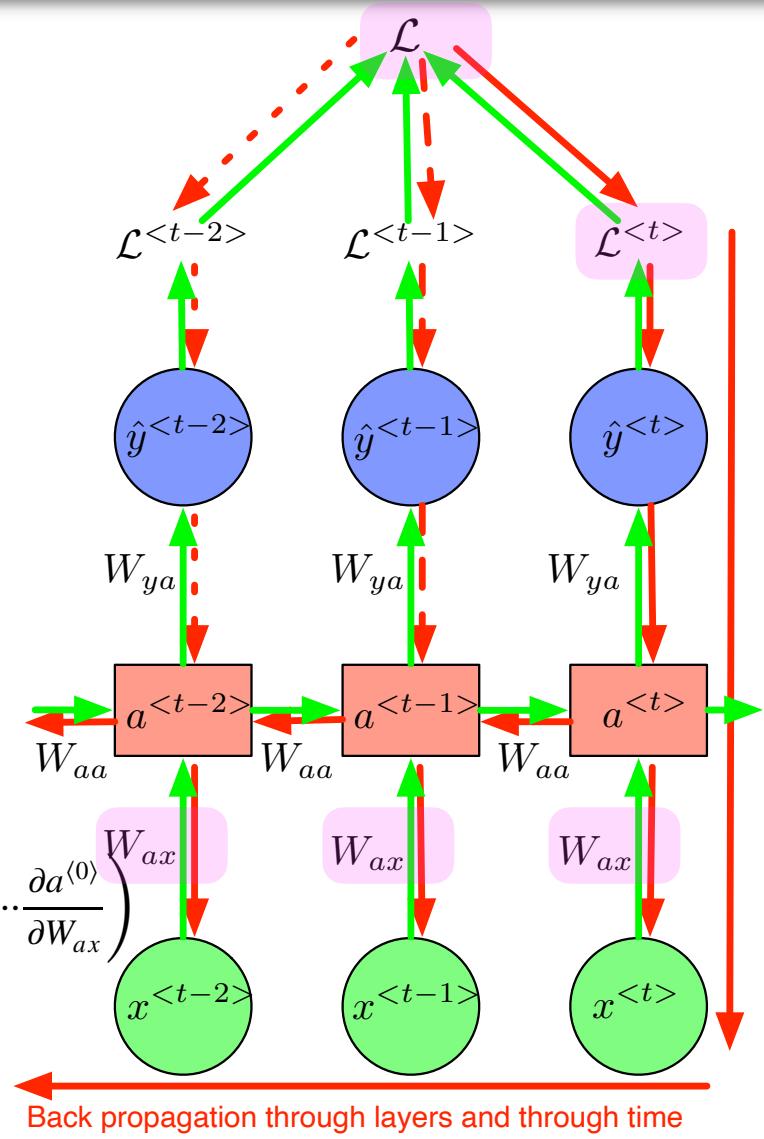
$$\frac{\partial \mathcal{L}^{(t)}}{\partial W_{ax}} \neq \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial a^{(t)}} \frac{\partial a^{(t)}}{\partial W_{ax}}$$

- because $a^{(t)}$ also depends on W_{ax} through $a^{(t-1)}$ which itself depends on ...

$$\mathbf{a}^{(t)} = g_a \left(\mathbf{W}_{ax} \mathbf{x}^{(t)} + \mathbf{W}_{aa} \mathbf{a}^{(t-1)} + \mathbf{b}_a \right)$$

- We use a formula of total derivative

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



Vanishing and exploding gradients

- Vanishing gradient
 - In MLP: very deep NN
 - very difficult to propagate back the gradient to affect the weights of the early Layers
 - In RNN: very difficult to learn 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.

Back Propagation Through Time (BPTT)

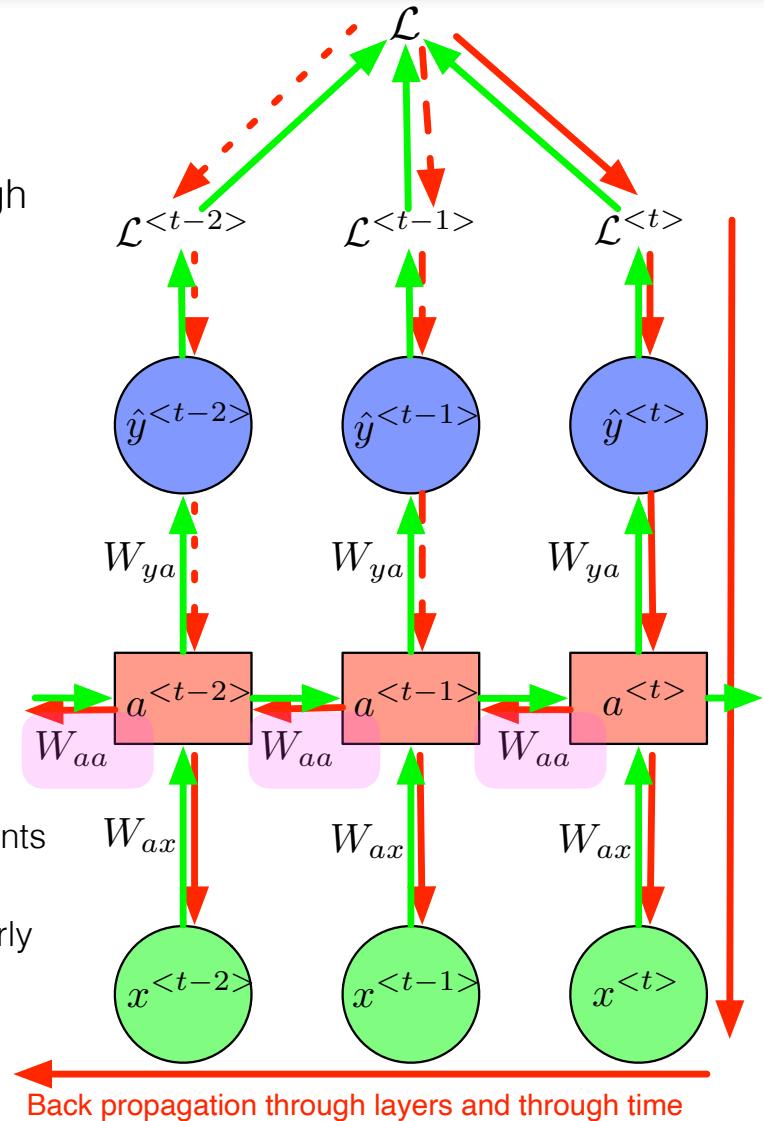
Vanishing and exploding gradients

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

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

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

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



Back Propagation Through Time (BPTT)

Vanishing and exploding gradients

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

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

_ Suppose only one hidden units, then $a^{(t)}$ is a scalar and consequently $\frac{\partial a^{(t)}}{\partial a^{(t-1)}}$ is also a scalar

- if $\left| \frac{\partial a^{(t)}}{\partial a^{(t-1)}} \right| < 1$, then the product goes to **0** exponentially fast
- if $\left| \frac{\partial a^{(t)}}{\partial a^{(t-1)}} \right| > 1$, then the product goes to **∞** exponentially fast

_ Vanishing gradients: $\left| \frac{\partial a^{(t)}}{\partial a^{(t-1)}} \right| < 1$

- contributions from faraway steps vanish and don't affect the training
- difficult to learn long-range dependencies

_ Exploding gradients: $\left| \frac{\partial a^{(t)}}{\partial a^{(t-1)}} \right| > 1$

- make the learning process unstable
- gradient could even become **NaN**

Back Propagation Through Time (BPTT)

Dealing with the exploding gradient

– Solution 1: gradient clipping

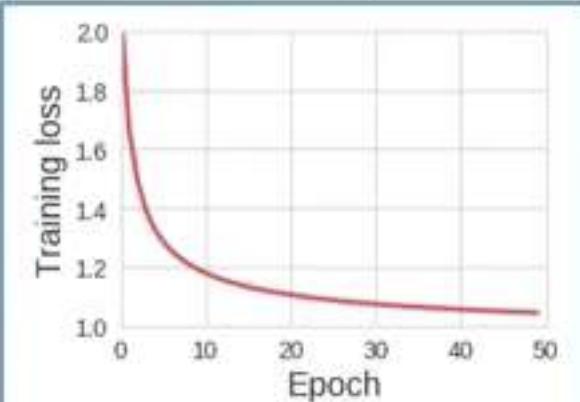
– gradient $d\theta = \frac{\partial \mathcal{L}}{\partial \theta}$ where θ are all the parameters

- if $\|\theta\| > \text{threshold}$

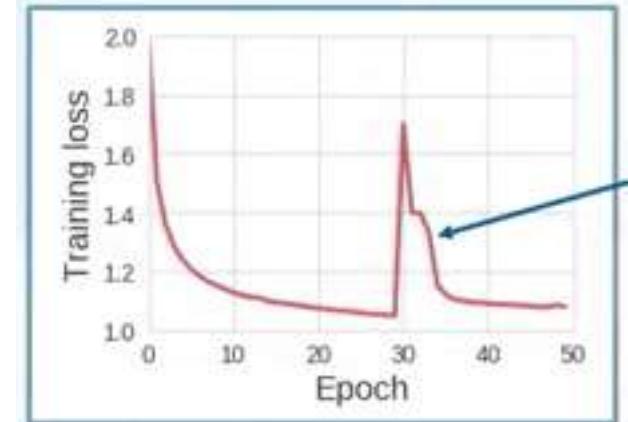
$$d\theta \leftarrow \frac{\text{threshold}}{\|\theta\|} d\theta$$

- 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

Stable learning curve



Unstable learning curve

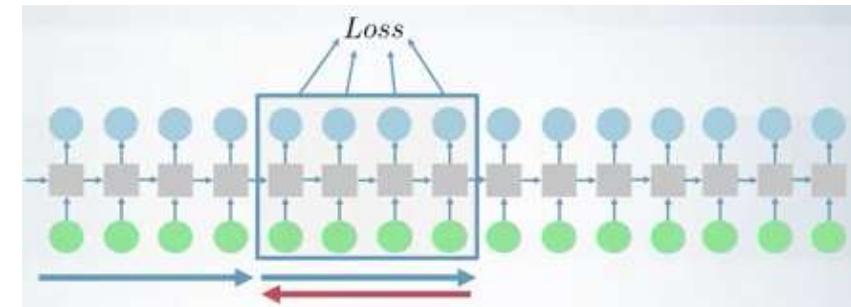
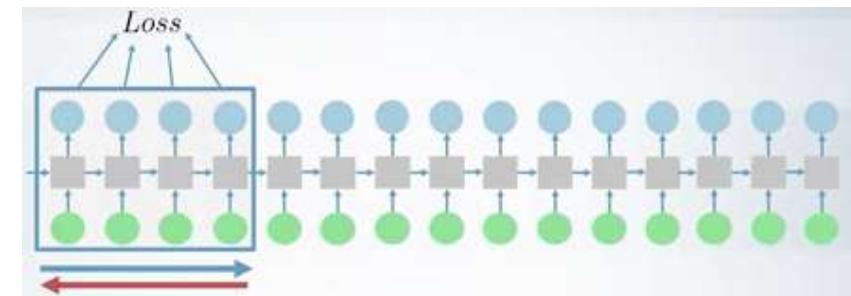
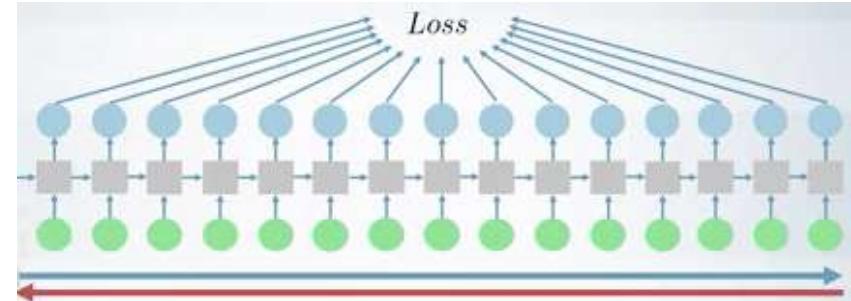


Back Propagation Through Time (BPTT)

Dealing with the exploding gradient

– Solution 1: 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 back-propagate 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)

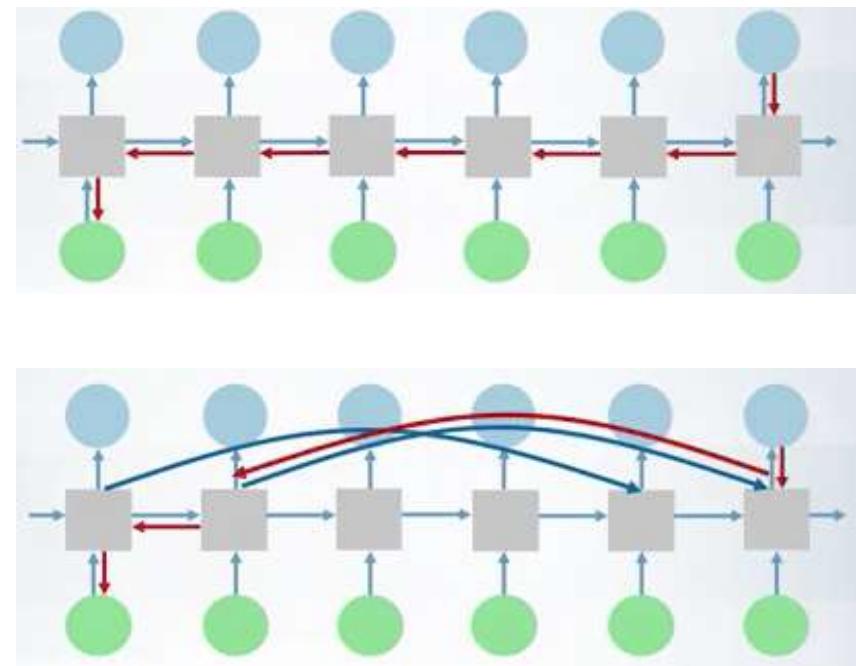


Back Propagation Through Time (BPTT)

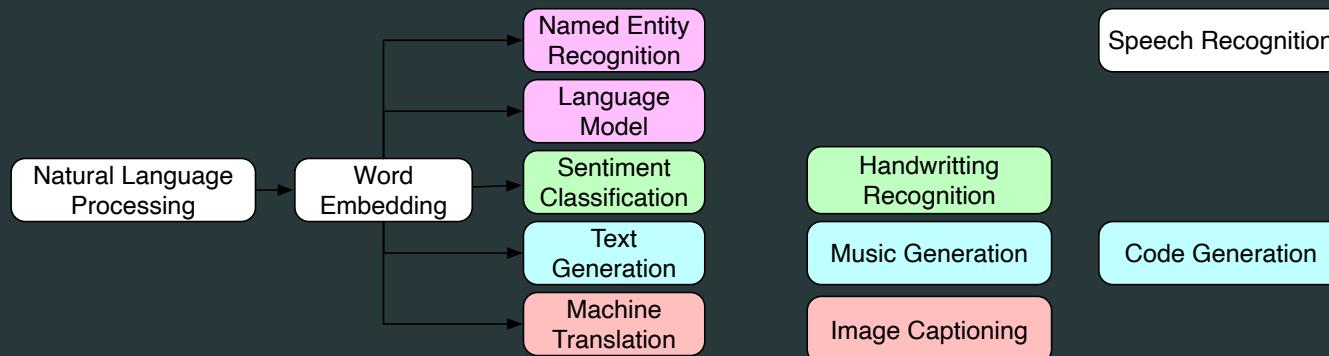
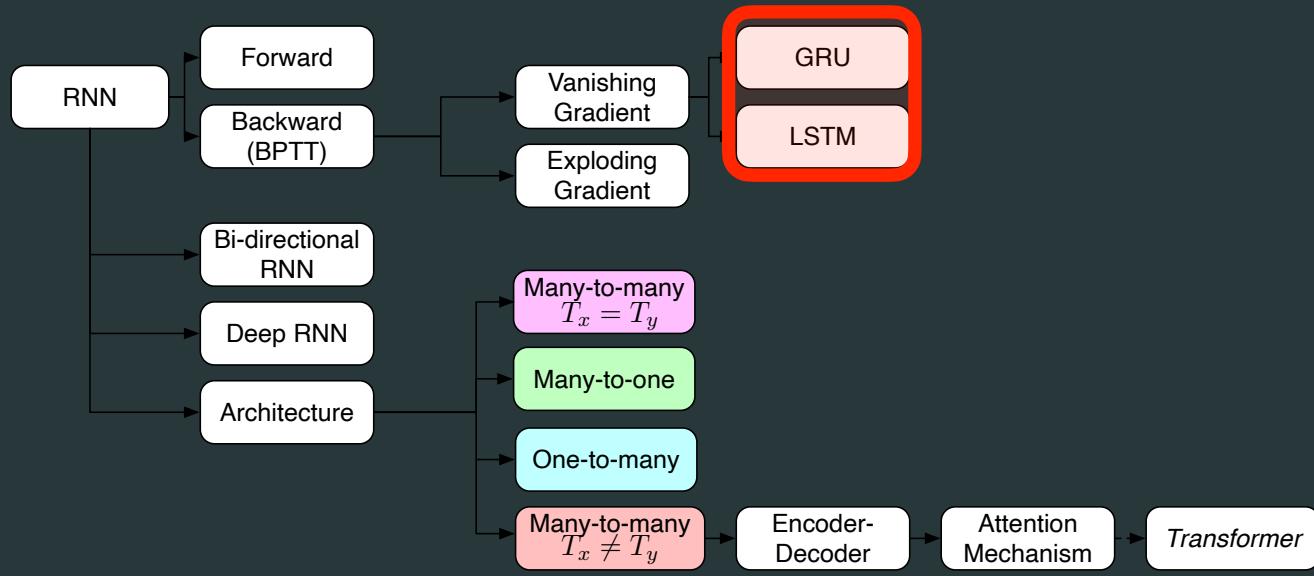
Dealing with the vanishing gradient

– 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
- Back-propagation through the shortcuts: the gradients vanish slower
 - we can learn long-term dependencies
- not exclusive to RNN (see ResNet)



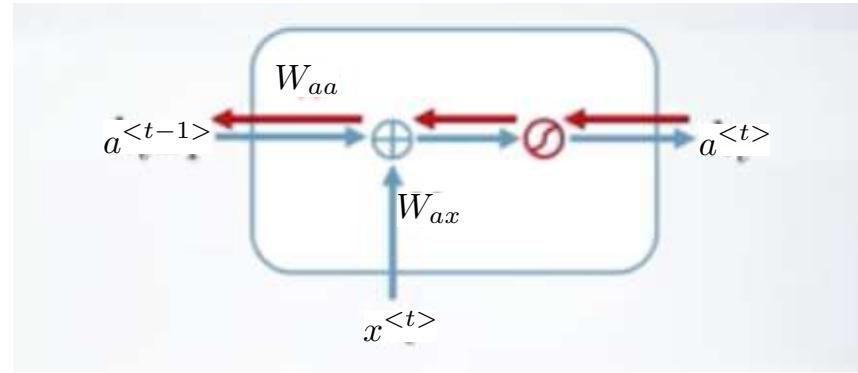
Gated units (LSTM and GRU)



Main ideas : Recurrent Neural Network (RNN)

$$\mathbf{a}^{(t)} = g_a \left(\mathbf{W}_{ax} \mathbf{x}^{(t)} + \mathbf{W}_{aa} \mathbf{a}^{(t-1)} + \mathbf{b}_a \right)$$

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



Gated units (LSTM and GRU)

Main ideas : Long Short Term Memory Units (LSTM)

- Add a new path: the **memory cell** $c^{(t)}$
 - same dimension as the hidden layer $a^{(t)}$
- **Simplified LSTM** (no possibility to erase)
 - Candidate cell value

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

- Gates (update, output)

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

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

- Output memory cell

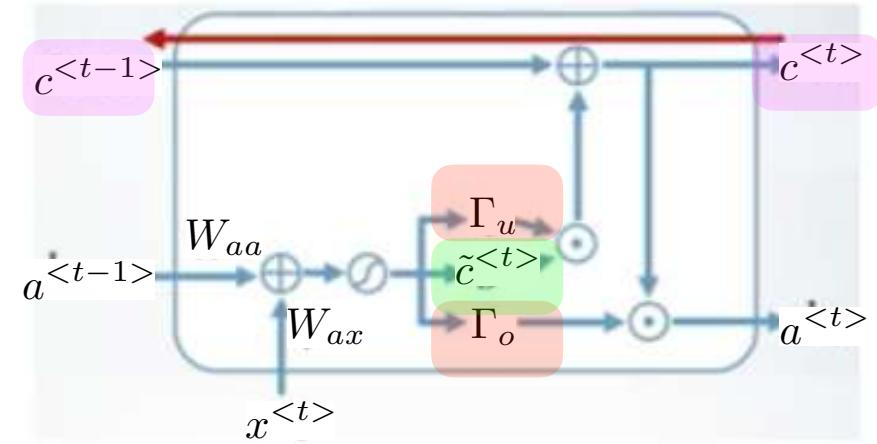
$$c^{(t)} = c^{(t-1)} + \Gamma_u \odot \tilde{c}^{(t)}$$

- Output hidden state

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

- no non-linearity or multiplication in the memory cell-path

- $\frac{\partial c^{(t)}}{\partial c^{(t-1)}} = \text{diag}(1) \Rightarrow$ no vanishing gradients



Gated units (LSTM and GRU)

Main ideas : Long Short Term Memory Units (LSTM)

– Forget Gate LSTM (no possibility to erase)

- Gates (forget)

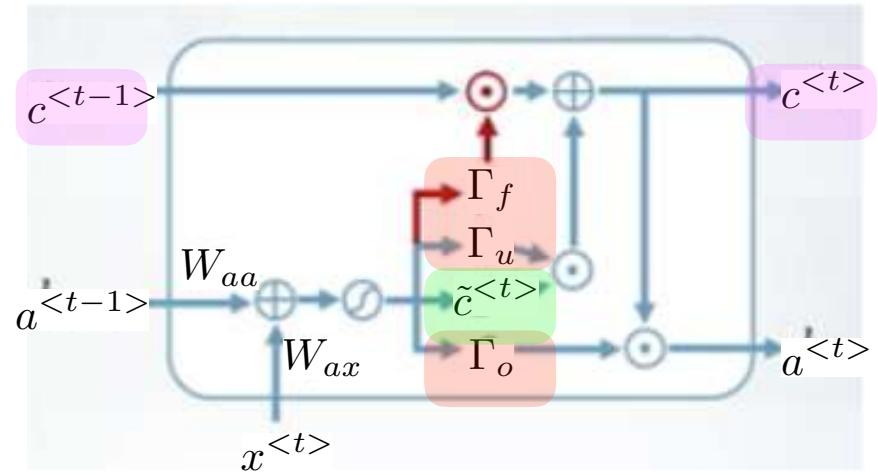
$$\Gamma_f = \sigma(W_{ax}^f x^{(t)} + W_{aa}^f a^{(t-1)} + b^f)$$

- Output memory cell

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

- Output hidden state

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



– We now have a multiplication on the short-way

- $\frac{\partial c^{(t)}}{\partial c^{(t-1)}} = \text{diag}(\Gamma_f)$

- with $\Gamma_f \in [0,1]$ the forget gate: $\Gamma_f = \sigma(W_{ax}^f x^{(t)} + W_{aa}^f a^{(t-1)} + b^f)$

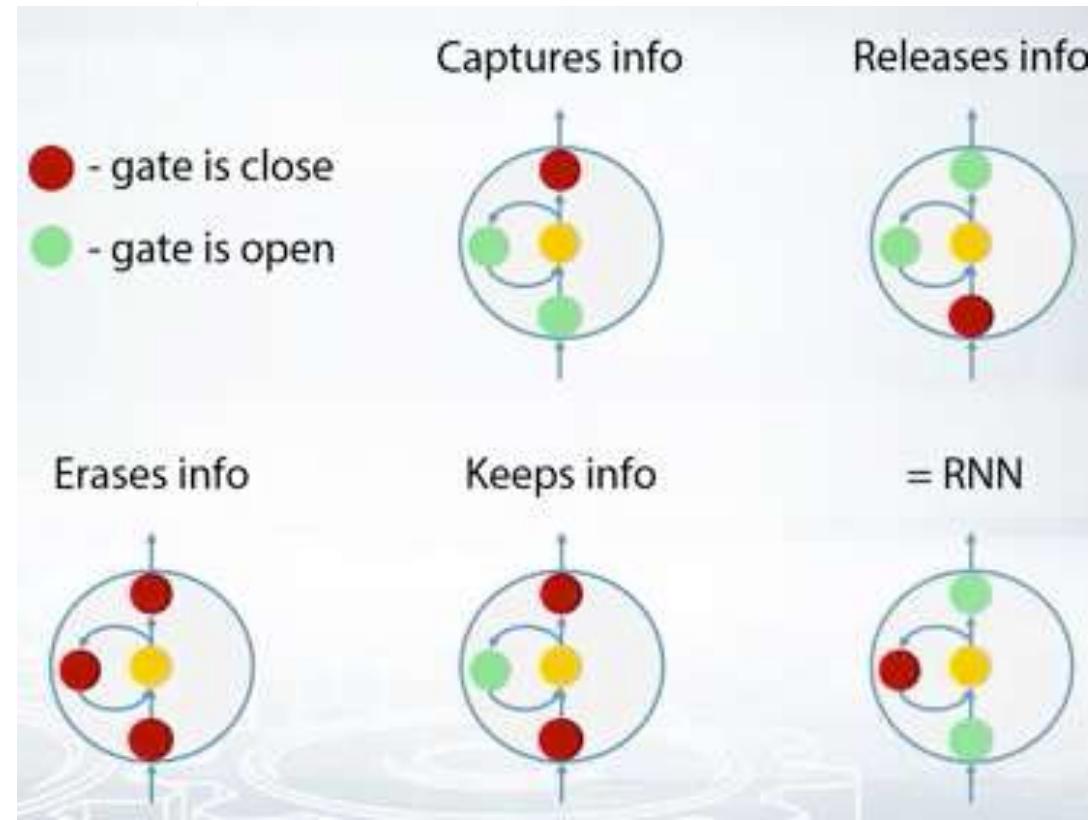
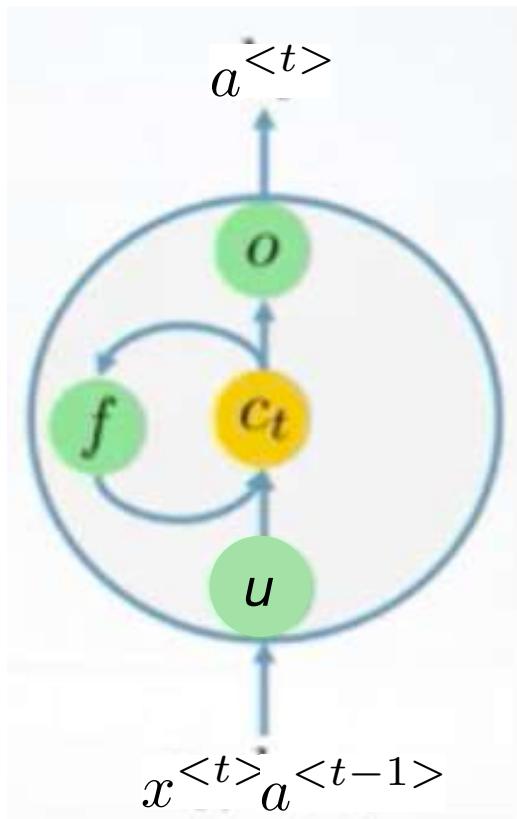
- set a high initial value for $b^f \Rightarrow$ at the beginning of training the gate is open, the LSTM doesn't forget and can learn long-term dependencies in the data

– **Note:** LSTM has four times more parameters to be trained than RNN

- $W_{ax}^f, W_{aa}^f, W_{ax}^u, W_{aa}^u, W_{ax}^o, W_{aa}^o, W_{ax}^c, W_{aa}^c, b^f, b^u, b^o, b^c$

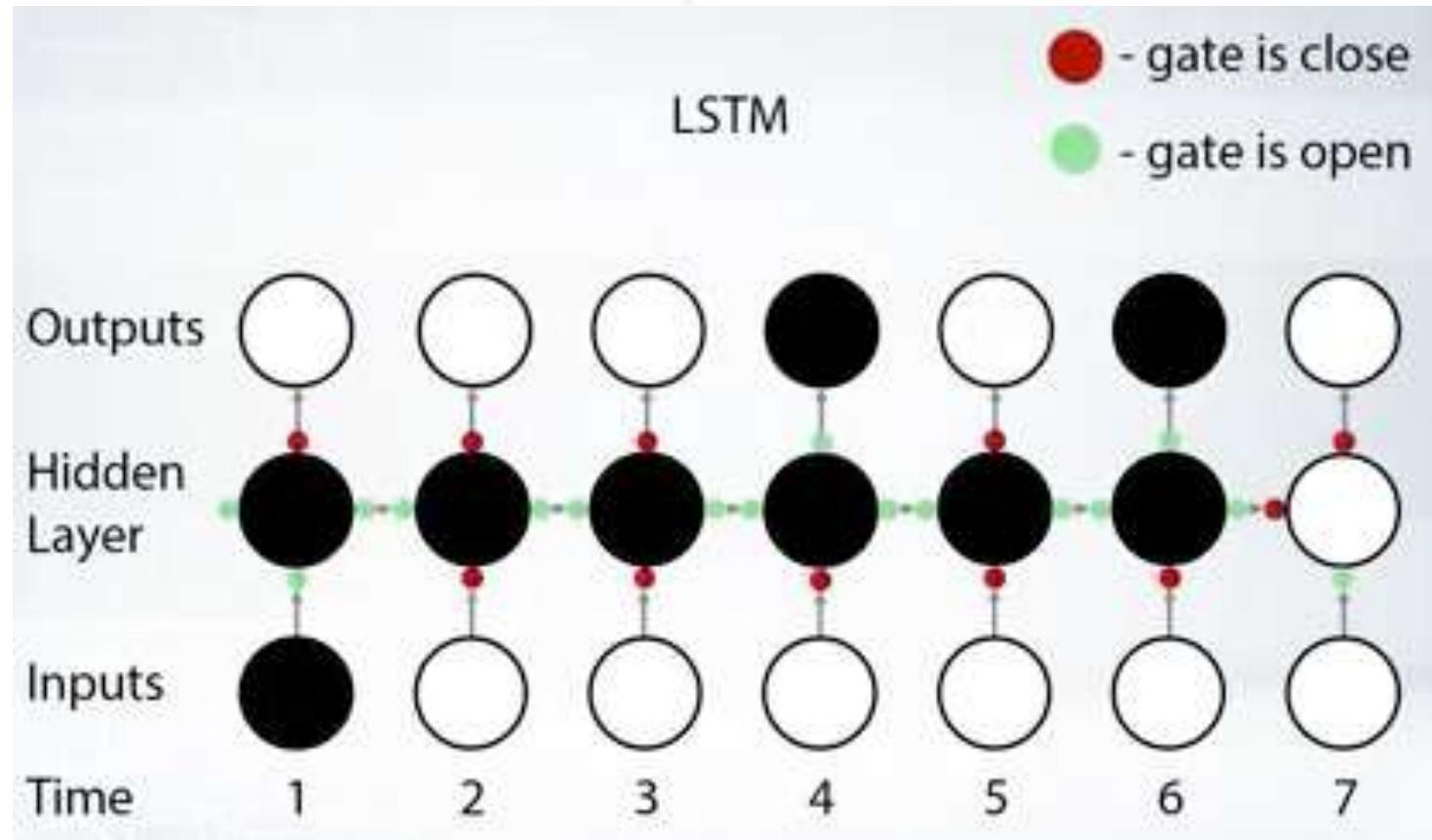
Gated units (LSTM and GRU)

LSTM regimes



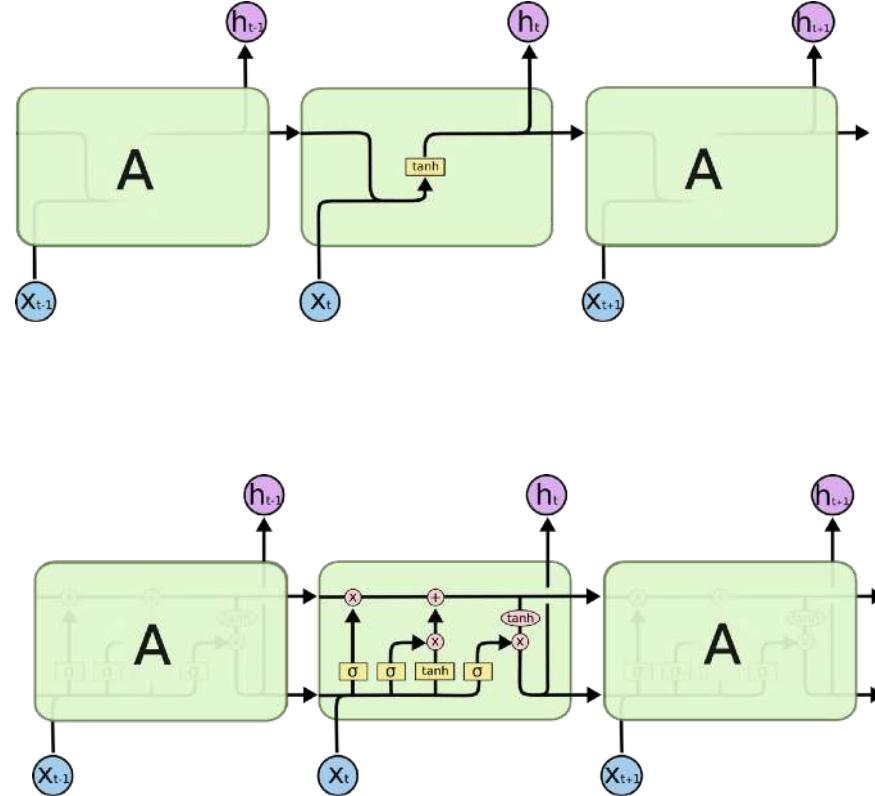
Gated units (LSTM and GRU)

LSTM regimes



LSTM example usage

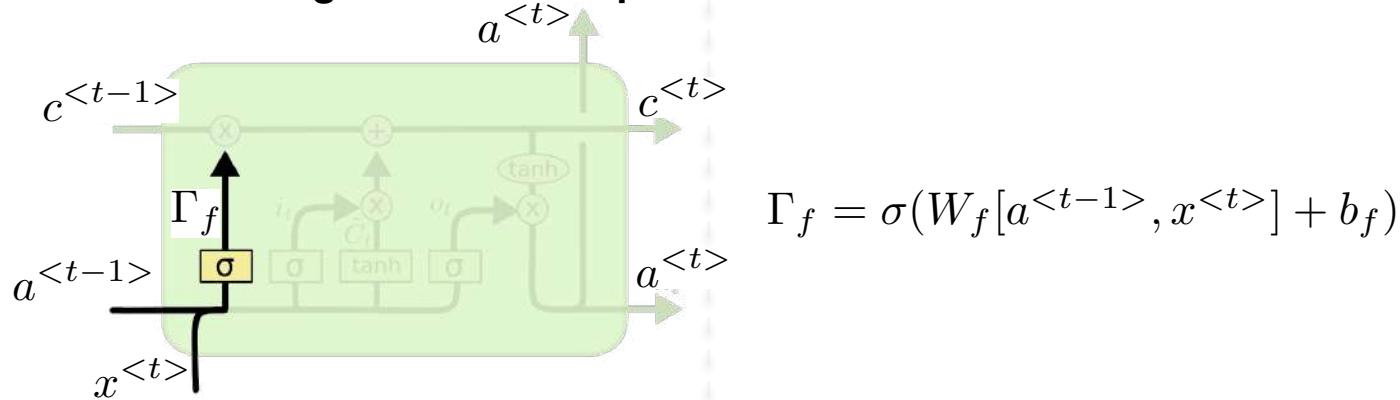
- Instead of having a single neural network layer, there are four, interacting in a very special way.
 - Lines carry a vector
 - A recurrent net transmits its hidden states
 - LSTM introduces a second channel : the cell state
 - It acts as a memory
 - Gates control the memory



Gated units (LSTM and GRU)

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

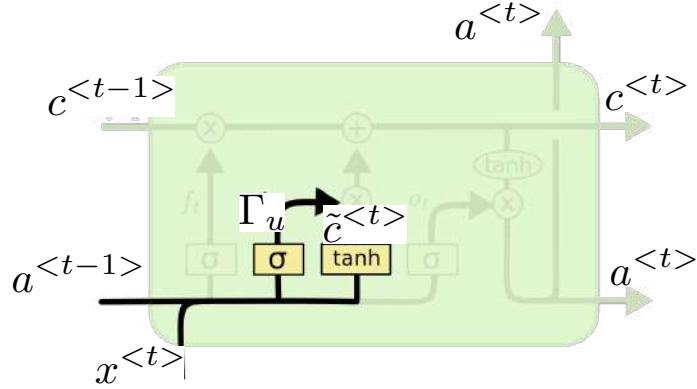
- Intuition for language modelling

- 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

Gated units (LSTM and GRU)

LSTM example usage

- What should be taken into account ?



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

$$\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

- Action

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

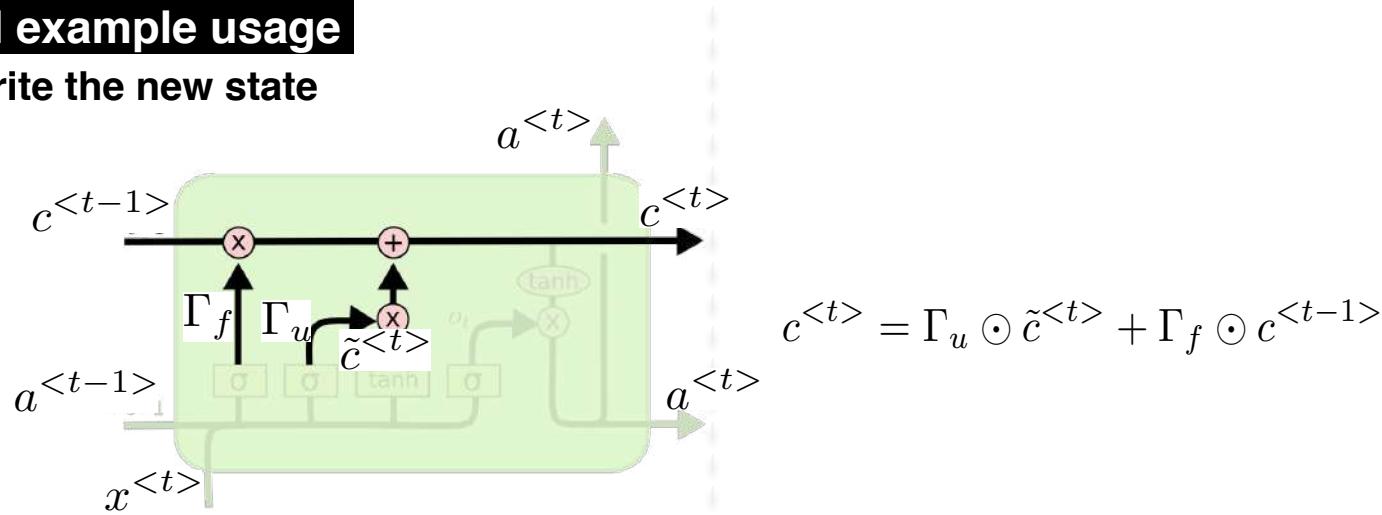
- Intuition for language modelling

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

Gated units (LSTM and GRU)

LSTM example usage

- Write the new state



- Action

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

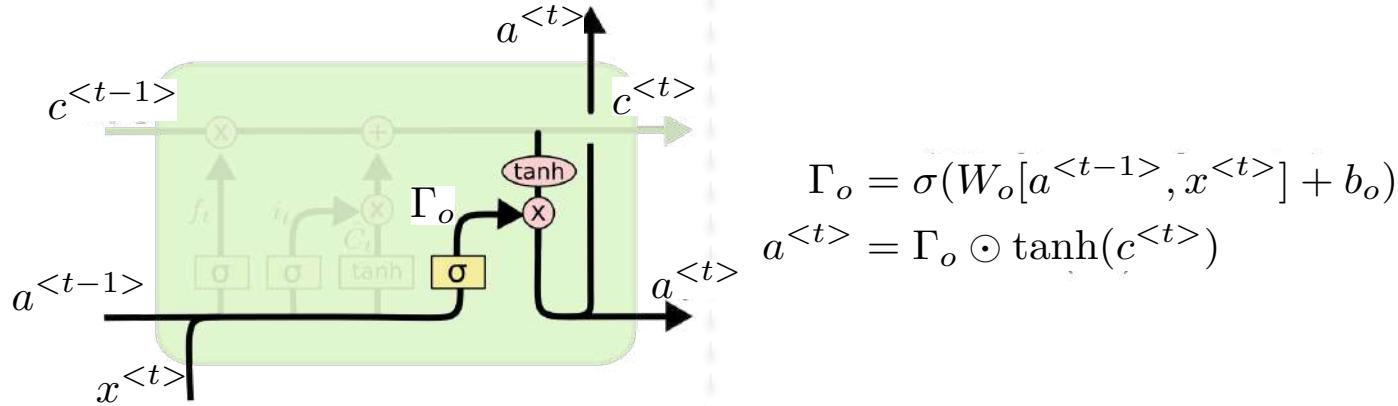
- Intuition for language modelling

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

Gated units (LSTM and GRU)

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 modelling

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

Gated units (LSTM and GRU)

Gated Recurrent Unit (GRU)

– Simple GRU

- Candidate cell value

$$\tilde{c}^{(t)} = \tanh \left(W^c [c^{(t-1)}, x^{(t)}] + b^c \right)$$

$$c^{(t-1)} = a^{(t-1)}$$

- Gates (update)

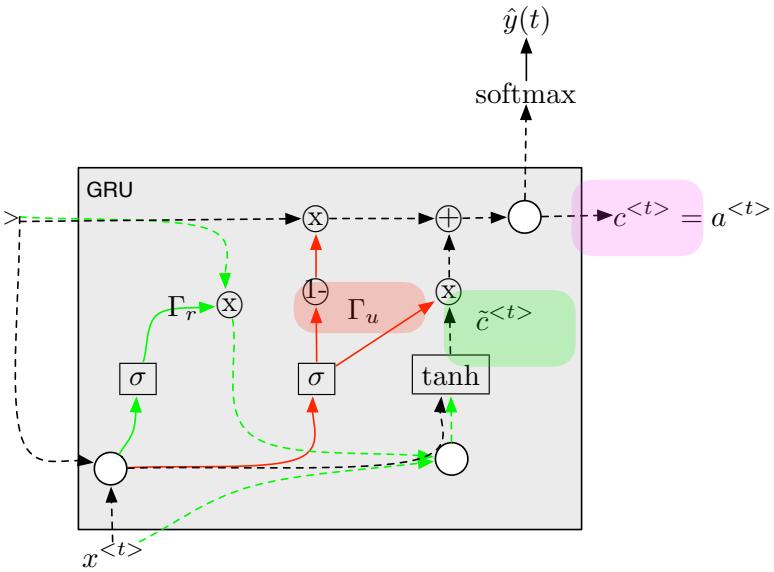
$$\Gamma_u = \sigma \left(W^u [c^{(t-1)}, x^{(t)}] + b^u \right)$$

- Output memory cell

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

- Output hidden state

$$a^{(t)} = c^{(t)}$$



Gated units (LSTM and GRU)

Gated Recurrent Unit (GRU)

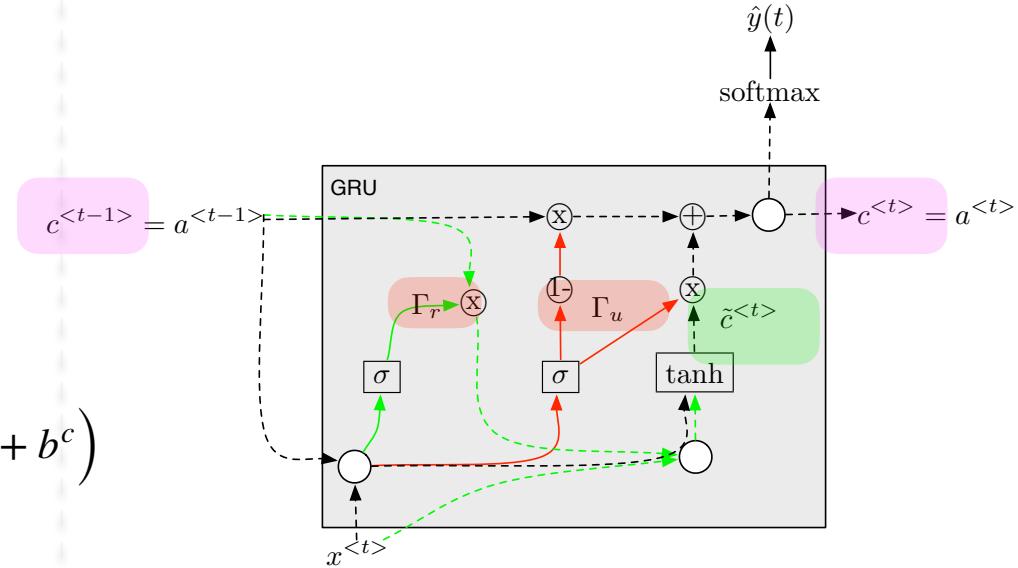
– Full GRU

- add a relevance gate Γ_r
- instead of using $c^{(t-1)}$ to compute $\tilde{c}^{(t)}$ we first filter it using a relevance gate Γ_r

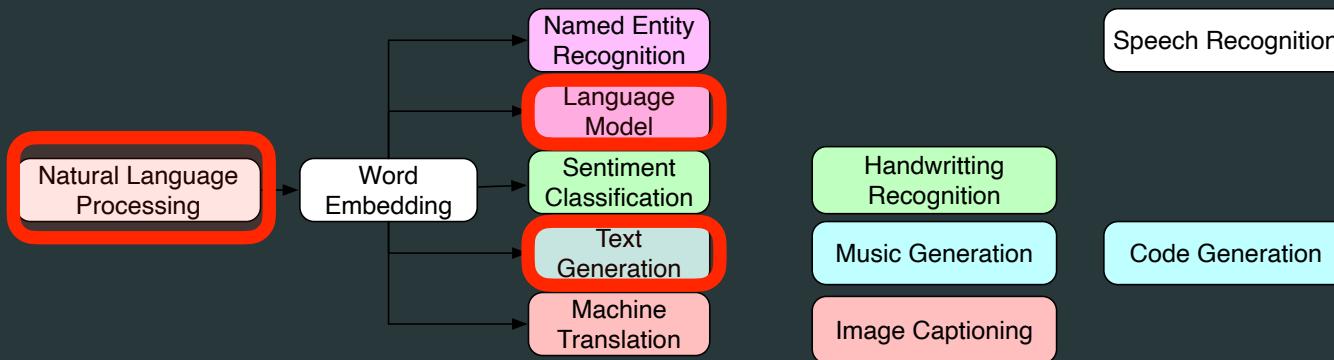
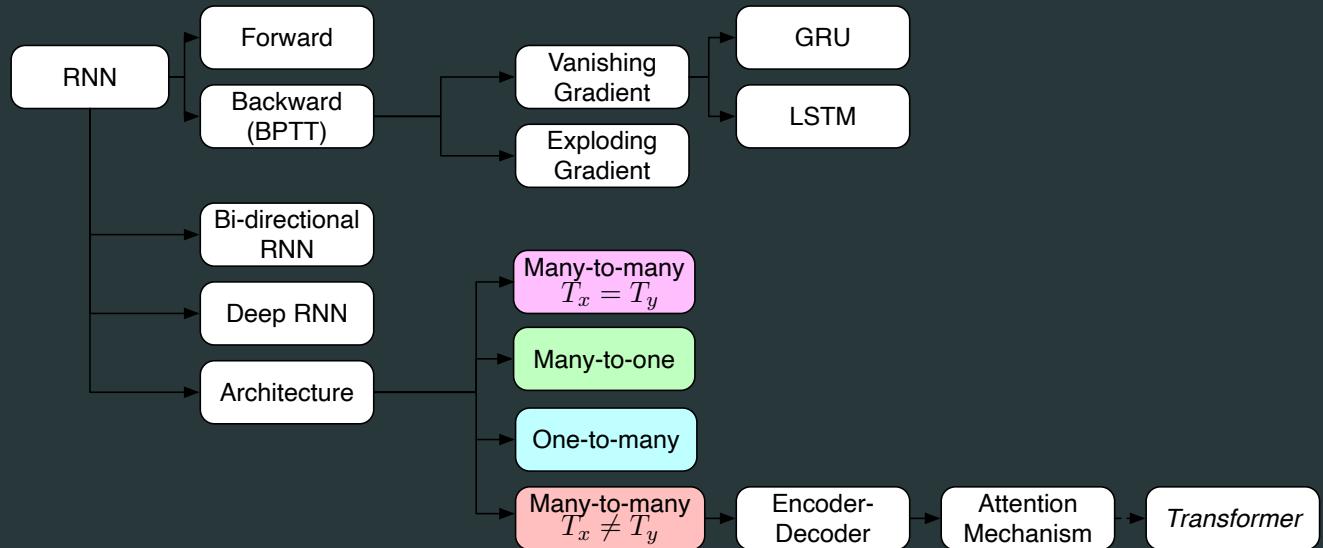
$$\Gamma_r = \sigma \left(W^r [c^{(t-1)}, x^{(t)}] + b^r \right)$$

- Candidate cell value

$$\tilde{c}^{(t)} = \tanh \left(W^c [\Gamma_r \odot c^{(t-1)}, x^{(t)}] + b^c \right)$$



Natural Language Processing



Language model

- **Applications**
 - Automatic Speech Recognition, Machine Translation, Optical Character Recognition, ...
- **Automatic Speech Recognition**
 - We need a language model to disambiguate what the acoustic model has eared
 - Choice between "This is an **example** of an homophone" or "This is an **egg-sample** of an homophone"
 - $p(\text{"example"}) = 5.7 \cdot 10^{-10} \Rightarrow$ is more likely in this context
 - $p(\text{"egg-sample"}) = 3.2 \cdot 10^{-13}$
- **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

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

- Simplification: the **N-gram assumption** :

$$P(w_1^L) = \prod_{i=1}^L P(w_i | w_{i-N+1}^{i-1}), \quad \forall i, w_i \in V$$

– Bi-gram (N=2) assumption

$$P(w_1^L) = P(w_1)P(w_2 | w_1)P(w_3 | w_2)P(w_4 | w_3)\dots P(w_L | w_{L-1})$$

- using **Recurrent Neural Network (RNN)**

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

Language model

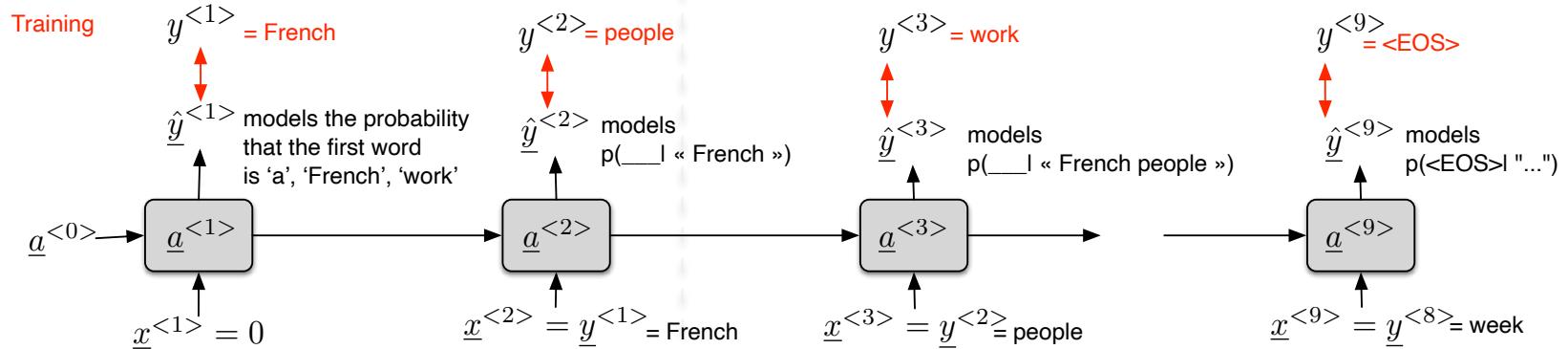
– How to build a language model ?

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

Natural Language Processing

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 (given by the hidden cell $a^{(t)}$) predict (maximise the probability to observe) the next word $y^{(t)}$ as output $\hat{y}^{(t)}$
- **Example:** "French people work an average of 35 hours a week."



- **Training:**
 - $\hat{y}^{(t)}$: output with a softmax with K neurons, gives the probability of each of the K words in the dictionary
 - $y^{(t)}$: ground-truth, one-hot-encoding with K classes
 - Minimise the cross-entropy loss

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

Natural Language Processing

Language model: (2) Test using a RNN

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

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

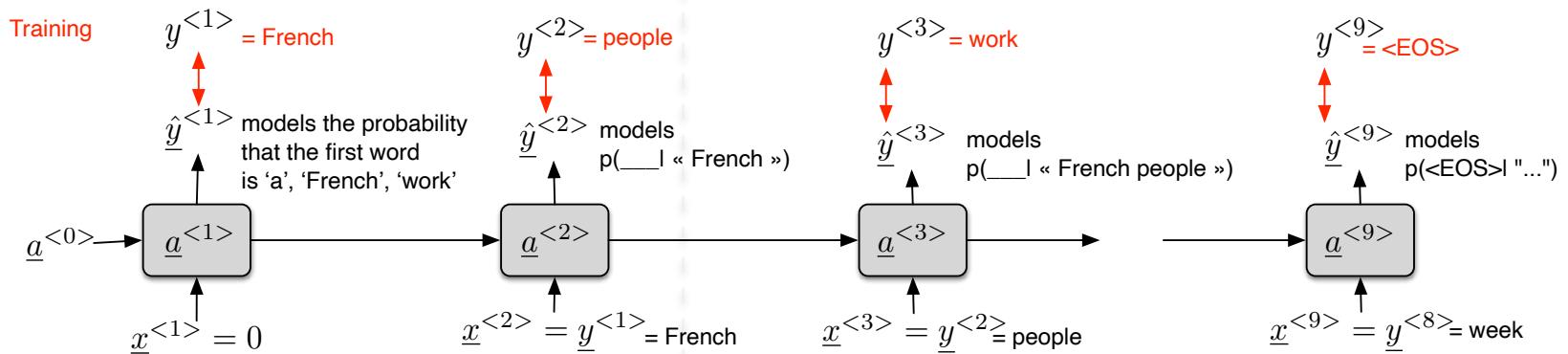
$$p(y^{(2)} | y^{(1)})$$

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

given by the first softmax

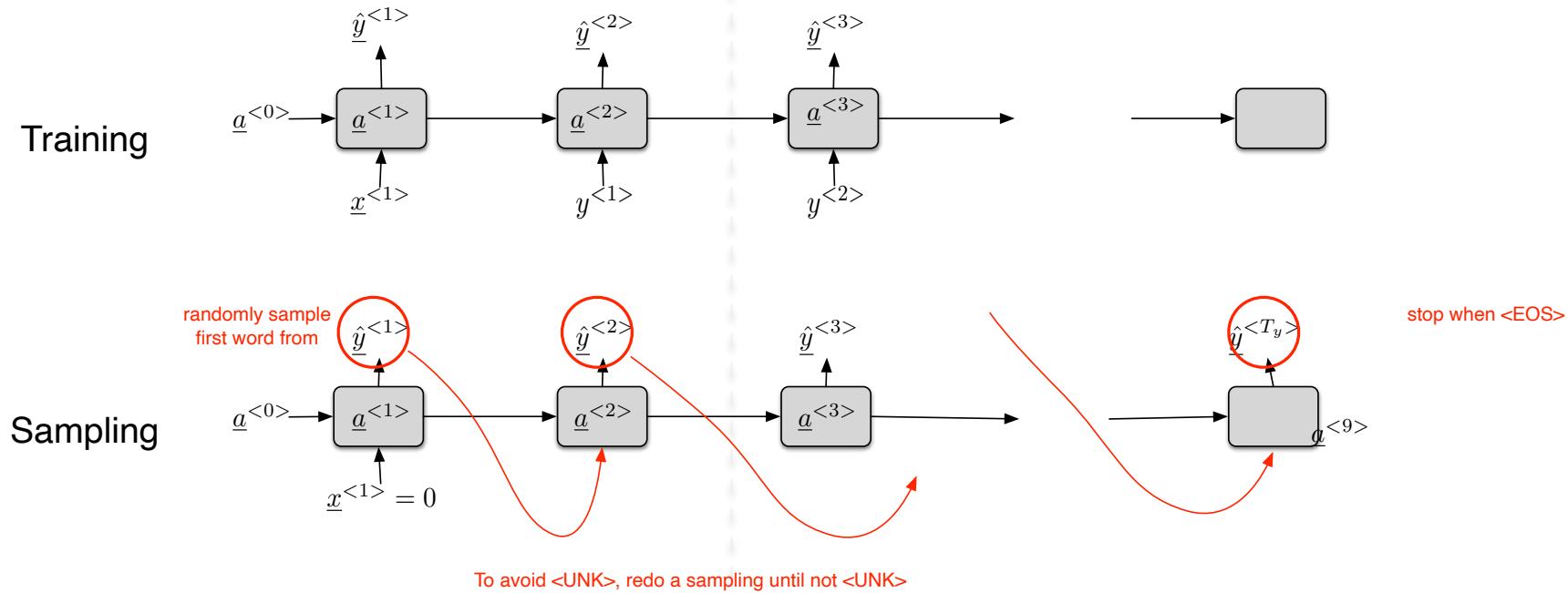
given by the second softmax

given by the third softmax



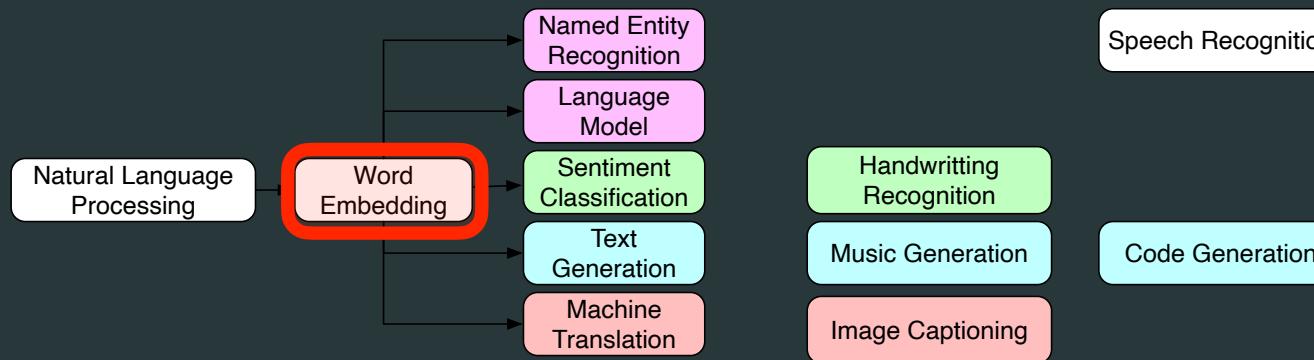
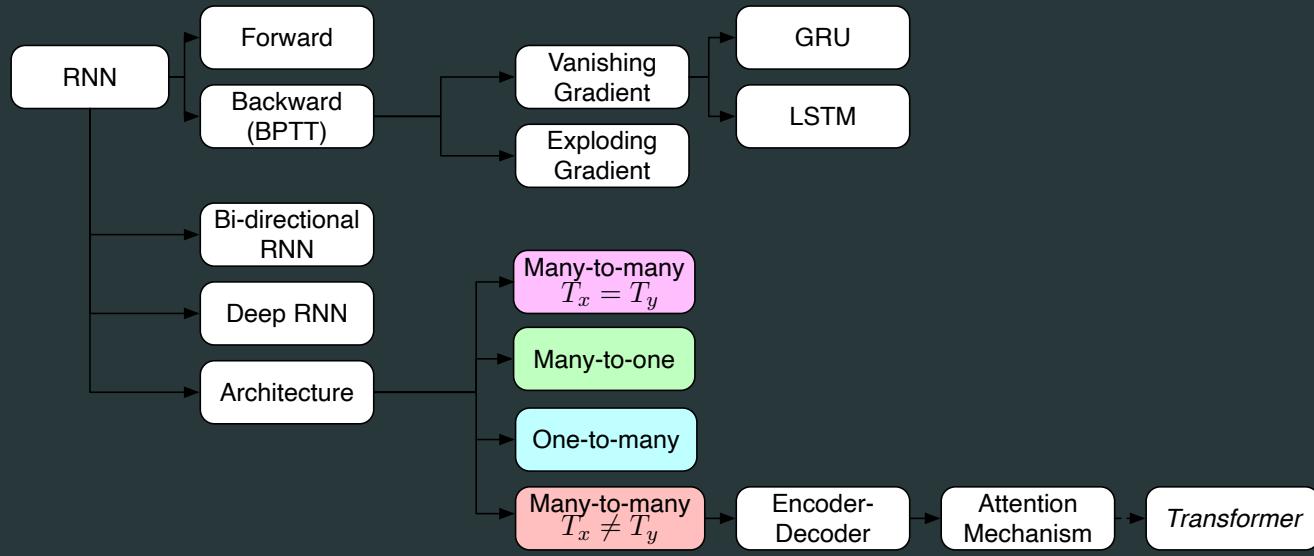
Natural Language Processing

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



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
...	chancelor
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

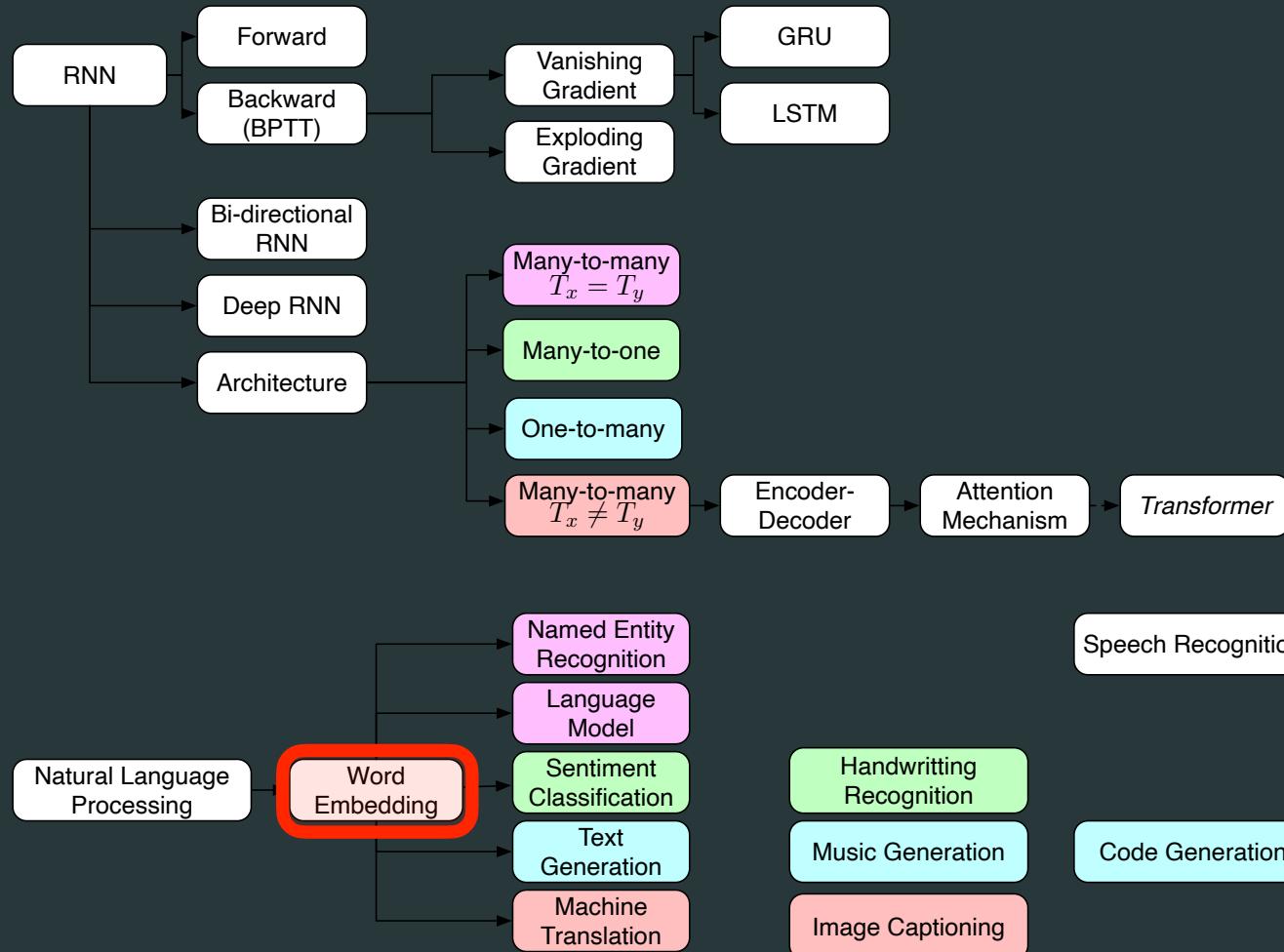
– Weaknesses of one-hot-encoding

- 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 generalise

– 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



Natural Language Processing

Word embedding

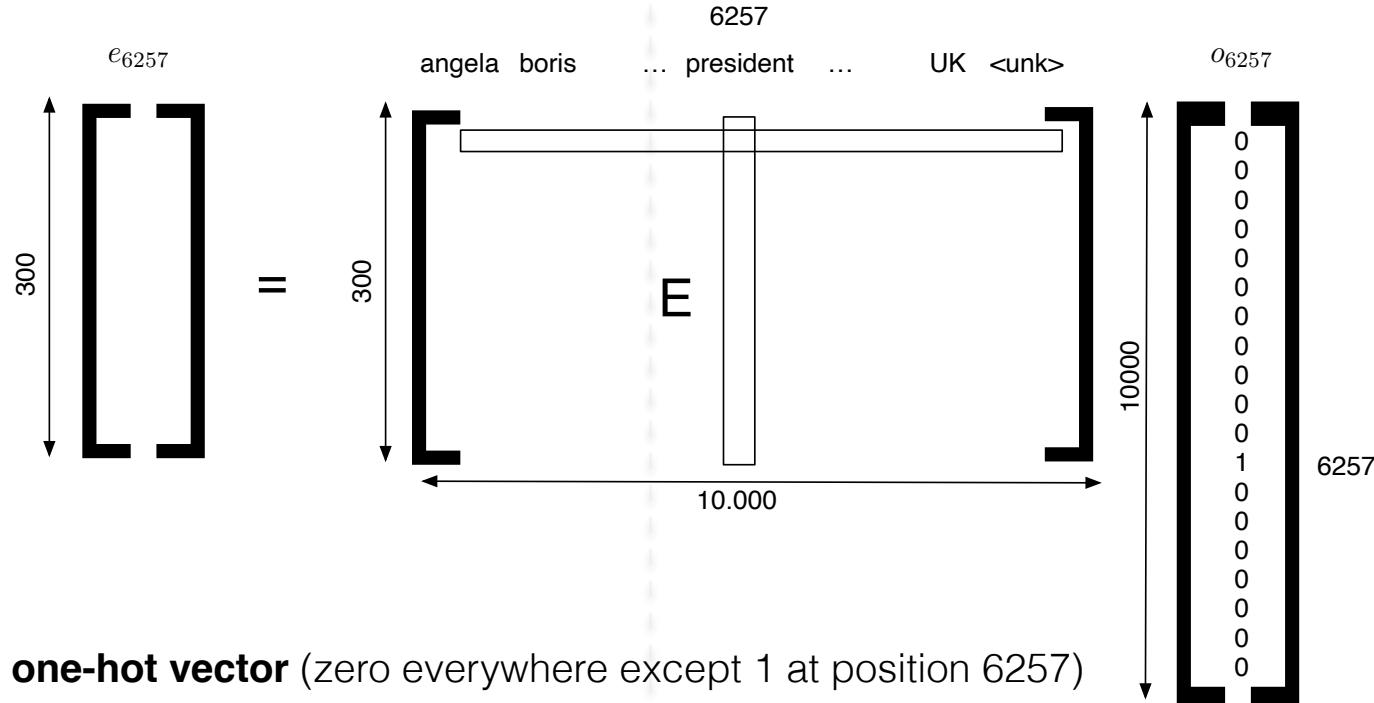
- learn a continuous representation of words
- each word w is associated with a real-valued vector e_{ID}
- if embedding size is 300, e_{5391} is a 300 dimensional vector
- we would like the distance $\|e_{ID1} - e_{ID2}\|$ to reflect the meaningful similarities between words

	Man 5391	Woman 9853	King 4914	Queen 7157	Uncle 456	Aunt 6257	President 7124	Chancellor 3212	France 6789	Germany 1234	Boat 5923	Plane 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 generalisation) that
 - The president is on a **boat** ____?____ \Rightarrow "trip"
 - The president is on a **plane** ____?____ \Rightarrow "trip"

Natural Language Processing

Word embedding/ Matrix



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

Word embedding/ Learning

- Various existing methods to learn word embedding:
 - Classic neural language model [Bengio et al., 2003]
 - Collobert and Weston model [Collobert, Weston, 2008]
 - Word2Vec [Mikolov et al. 2013]
 - Continuous bag-of-words (CBOW)
 - Skip-gram
 - Glove model [Pennington et al. 2014]

[Y. Bengio, R. Ducharme, P. Vincent, and C. Jauvin. A neural probabilistic language model. *Journal of machine learning research*, 3(Feb):1137–1155, 2003.]
[R. Collobert and J. Weston. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, pages 160–167. ACM, 2008.]

[T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.]

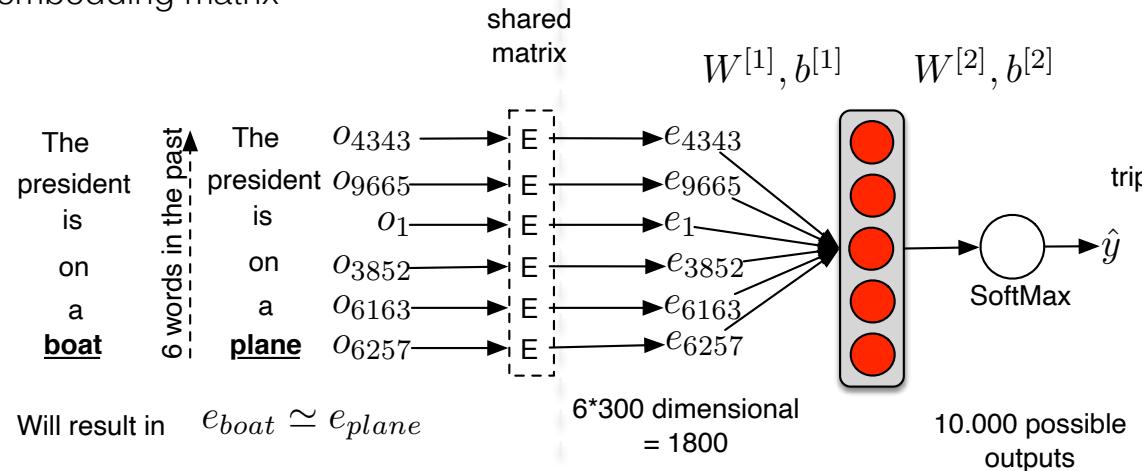
[T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.]

[J. Pennington, R. Socher, and C. Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.]

Natural Language Processing

Word embedding/ Classic neural language model

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



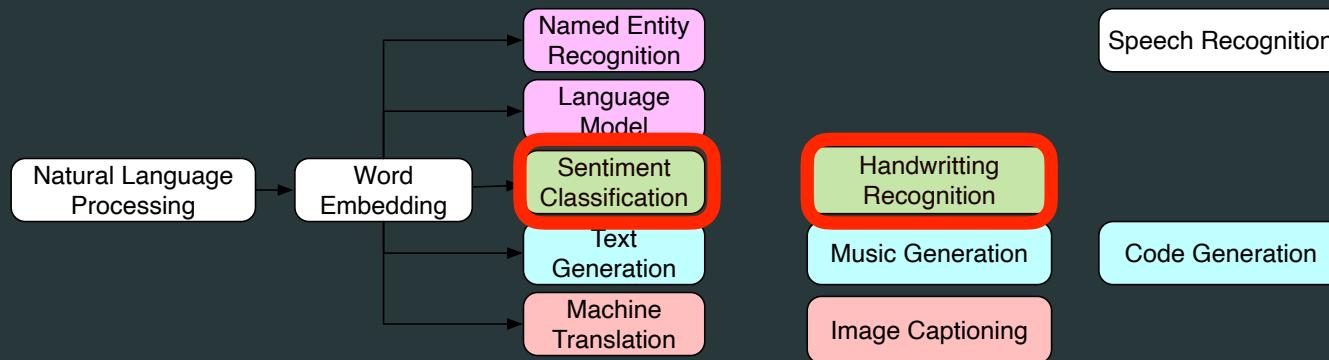
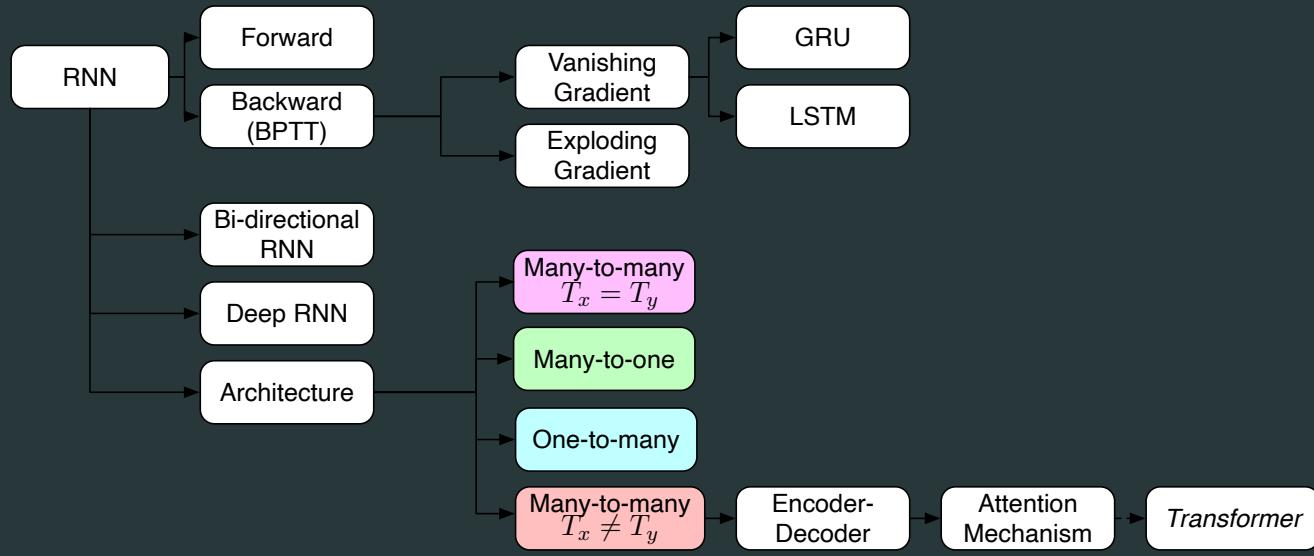
– Example:

- "The president is on a boat trip to the United States."
target

– Other possible definition of the context

- last 4 words
- 4 words on the left, 4 words on the right: "is on a boat ____?____ to the United States"
- last 1 words: "boat ____?____"
- nearby 1 word (skip-gram): "president ____?____"

Application: Sentiment Classification



Application: Sentiment Classification

x

y

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



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

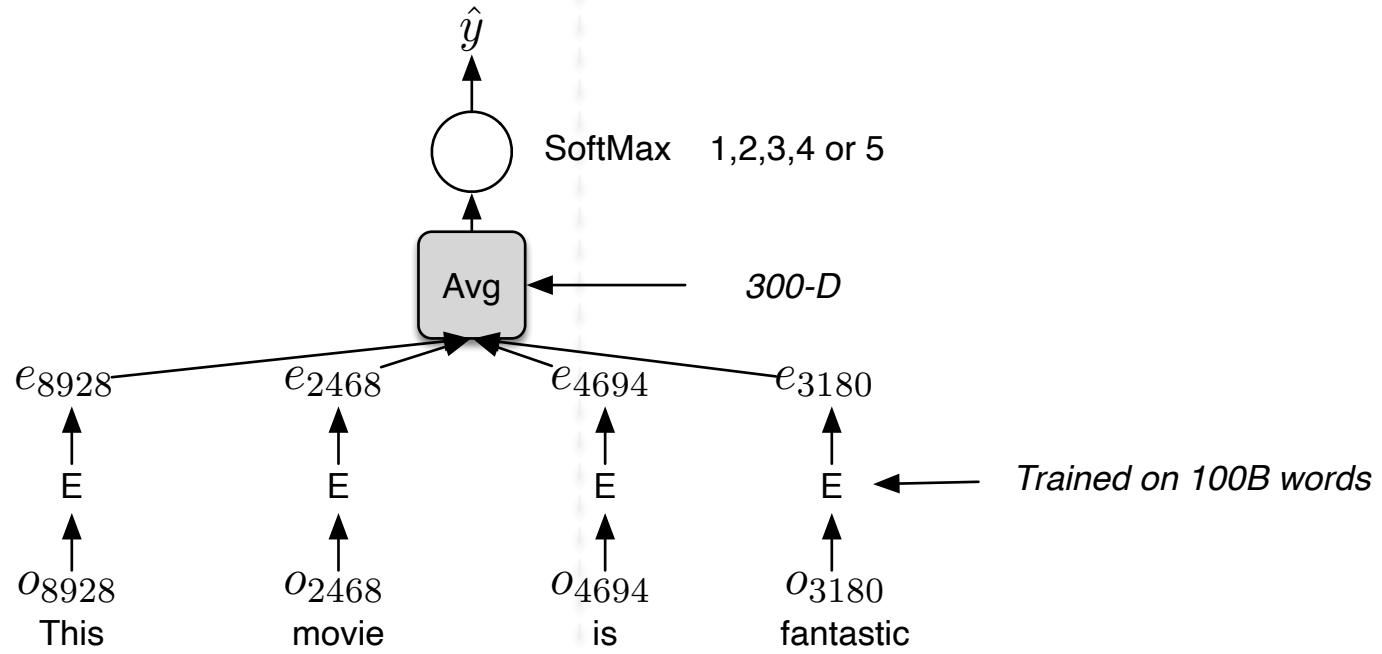


This movie was completely lacking a good script,
good actors and a good director.



Natural Language Processing

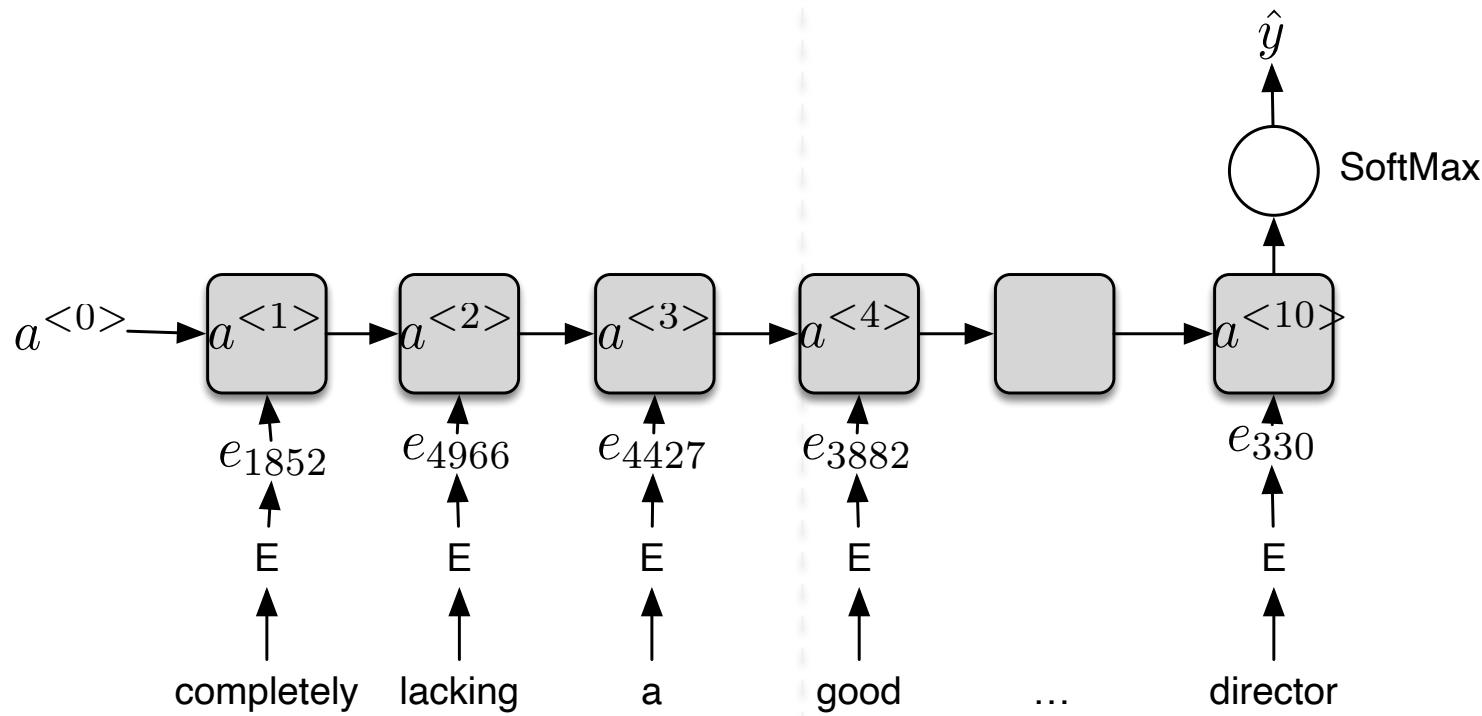
Application: Sentiment Classification / using a simple model



- Does not work for
"This movie was completely lacking a **good** script, **good** actors and a **good** director"
 - This is because the Avg of "good" is positive
 - The model can not understand that we mean the opposite ("lacking")
 - We need a Sequence model \Rightarrow RNN

Application: Sentiment Classification / using a RNN

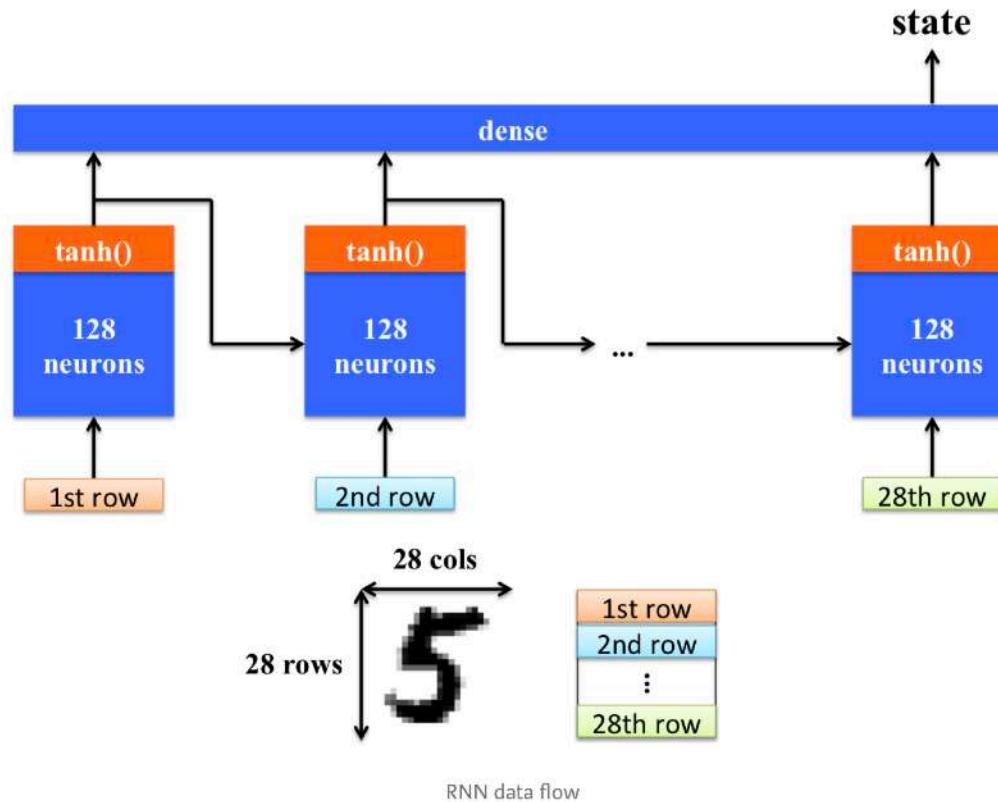
- **Many-to-one** architecture



Natural Language Processing

Application: Handwritten Character Recognition

- **Many-to-one** architecture

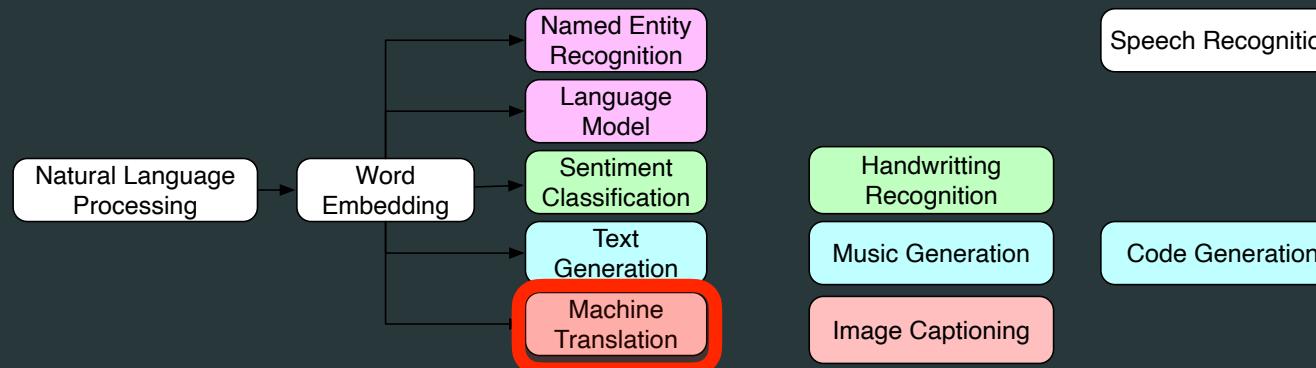
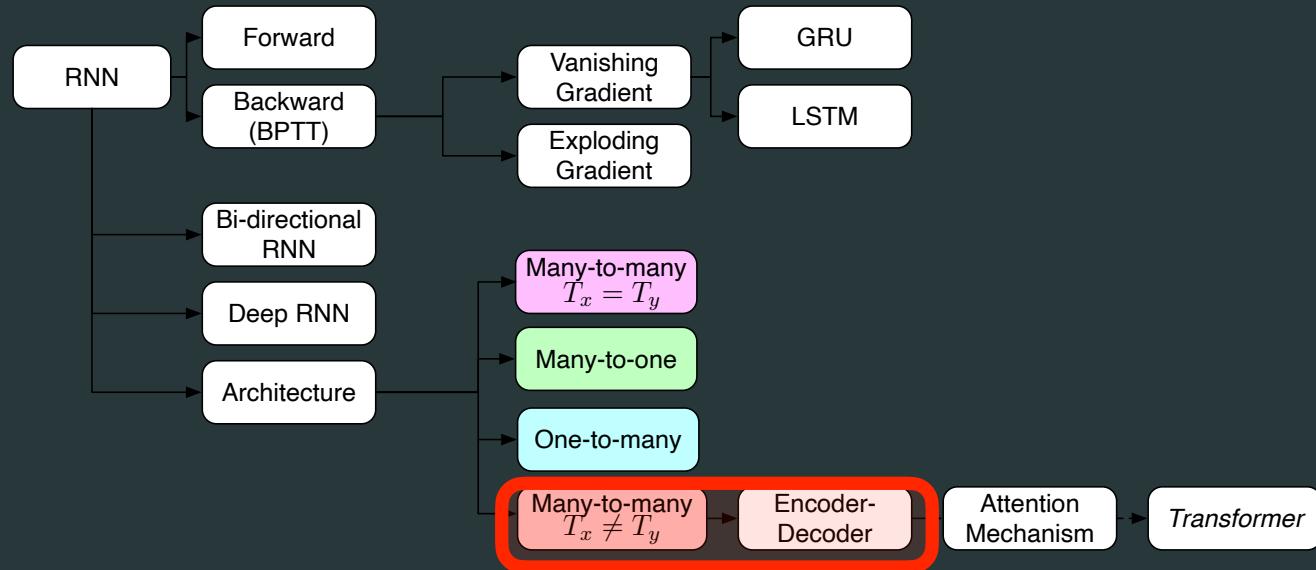


source <https://medium.com/machine-learning-algorithms/mnist-using-recurrent-neural-network-2d070a5915a2>

[Graves et al. 2008 "Unconstrained online handwriting recognition with recurrent neural networks"]

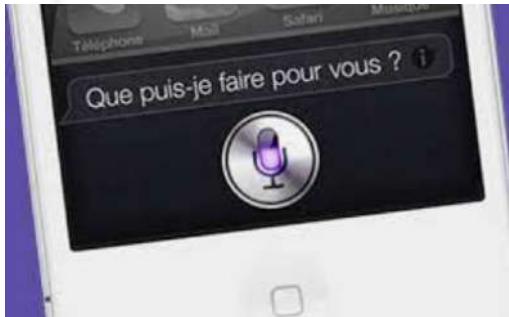
[Graves et al. 2009 "Offline Handwriting Recognition with Multidimensional Recurrent Neural Networks"]

Sequence to Sequence

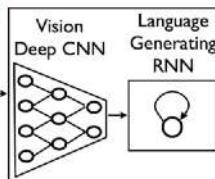
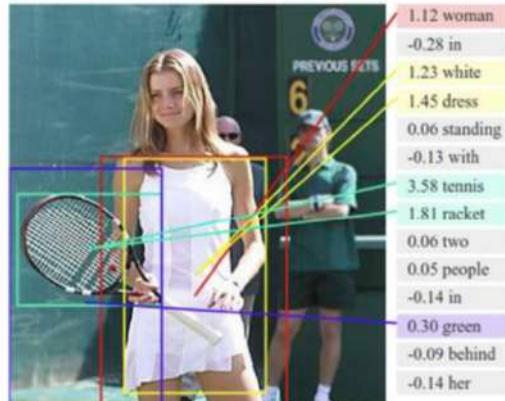


Applications

Automatic Speech Recognition



Automatic picture captioning



A group shopping at an outdoor market.
There are many vegetables at the fruit stand.

Self Driving Cars



DALL-E



DEEP LEARNING HAS MASTERED GO Google Alpha Go

nature International weekly journal of science

Home | News & Comment | Research | Careers & Jobs | Current Issue | Authors & Ideas | Photo & Video | Previous Issues | About Us | Help

AlphaGo | Where We Stand | New York | News | Article

ARTICLE | REVIEW

Google reveals secret test of AI bot to beat top Go players

Updated version of DeepMind's AlphaGo program beats last mystery online competitor.

Deep River, Alip Heming, Chris J. Maddison, Arthur O'Dwyer, Laurent Bégin, Ongene van den Driessche, Jukka Suominen, Sørensen Amrhein, Wei Pan, Van Den Herrewyn, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Milner, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Massaketh Leon, Koen Keijzer, Thore Graepel & Demis Hassabis

ARTICLES | Contributors | Correspondence | Authors

Nature 509, 484–489 (28 January 2014) | doi:10.1038/nature14911
Received: 11 November 2013 | Accepted: 05 January 2014 | Published online: 27 January 2014



Style Transfer

Zebras ↪ Horses



zebra → horse



horse → zebra

Old → Young

Male → Female

Female → Male

Neural Machine Translation

DeepL Traducteur DeepL Pro Connexion

Découvrir DeepL Pro Connexion

Anglais (langue détectée) ▾ Fr... ▾ formelle/informel ▾ Glossaire

Deep learning (also known as deep structured learning) is part of a broader family of machine learning methods based on artificial neural networks with representation learning. Learning can be supervised, semi-supervised or unsupervised.

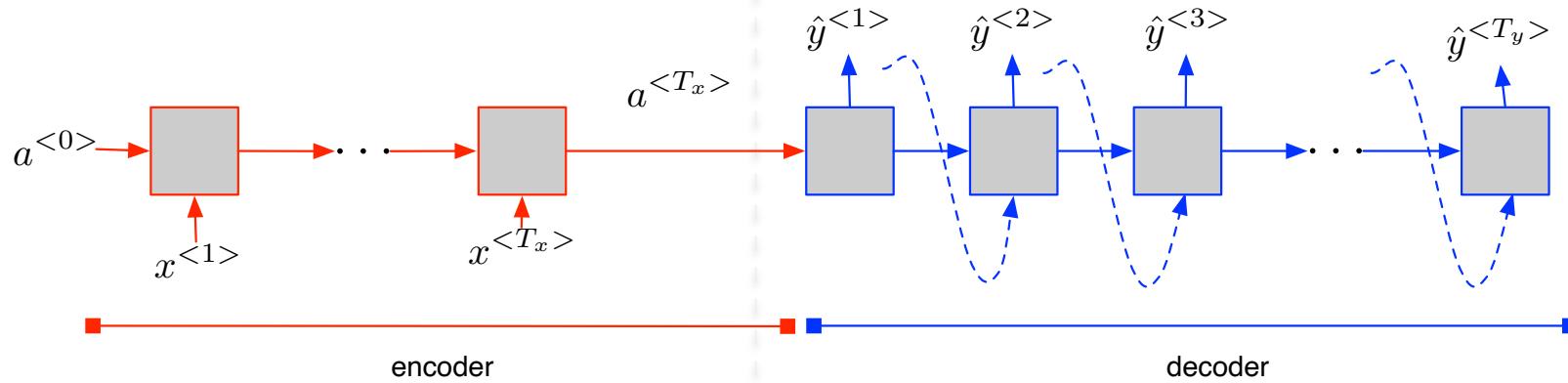
L'apprentissage profond (également appelé apprentissage structuré profond) fait partie d'une famille plus large de méthodes d'apprentissage automatique basées sur les réseaux neuronaux artificiels avec apprentissage par représentation. L'apprentissage peut être supervisé, semi-supervisé ou non supervisé.

Sequence to Sequence

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

{x}	x ^{<1>}	x ^{<2>}	...			x ^{<7>}		
	Deep	Learning	is	part	of	machine	learning	
{y}	y ^{<1>}	y ^{<2>}	...					y ^{<9>}
	L'	apprentissage	profond	fait	partie	de	l'	apprentissage machine

- Solution:
 - **Sequence to sequence** [Sutskever et al., 2014] or **Encoder-Decoder** [Cho et al., 2014] architecture



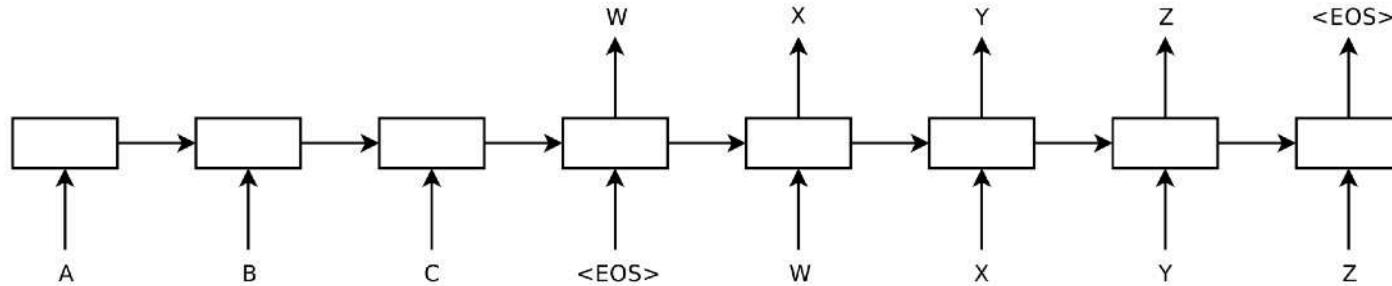
[Sutskever et al. 2014, "Sequence to sequence learning with neural networks"] [LINK](#)

[Cho et al. 2014 "Learning phrase representations using RNN encoder-decoder for statistical machine translation"] [LINK](#)

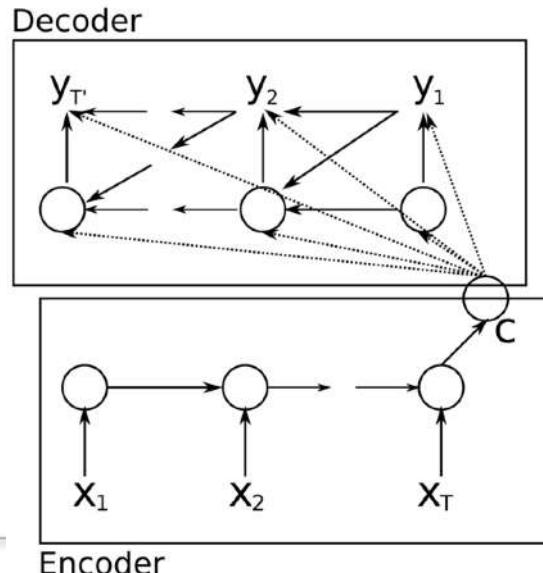
Sequence to Sequence

Introduction

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



[Cho et al. 2014 "Learning phrase representations using RNN encoder-decoder for statistical machine translation"] [LINK](#)

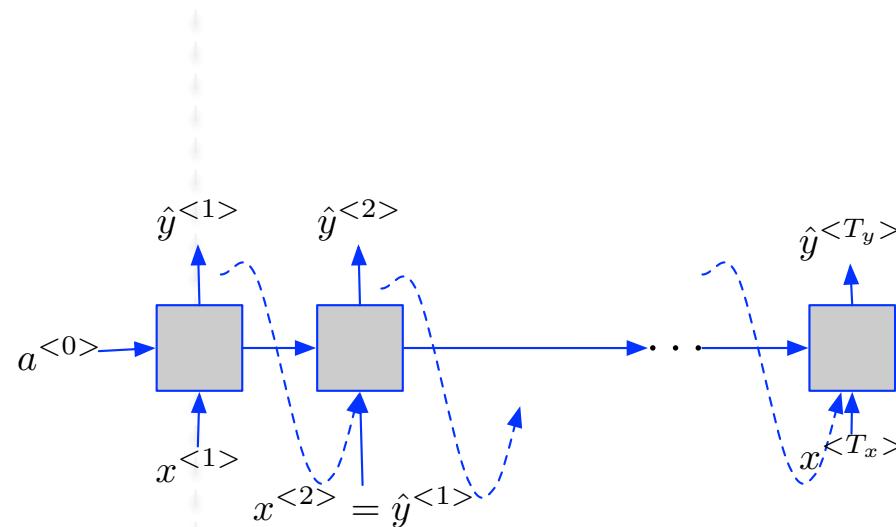


Sequence to Sequence

Machine Translation

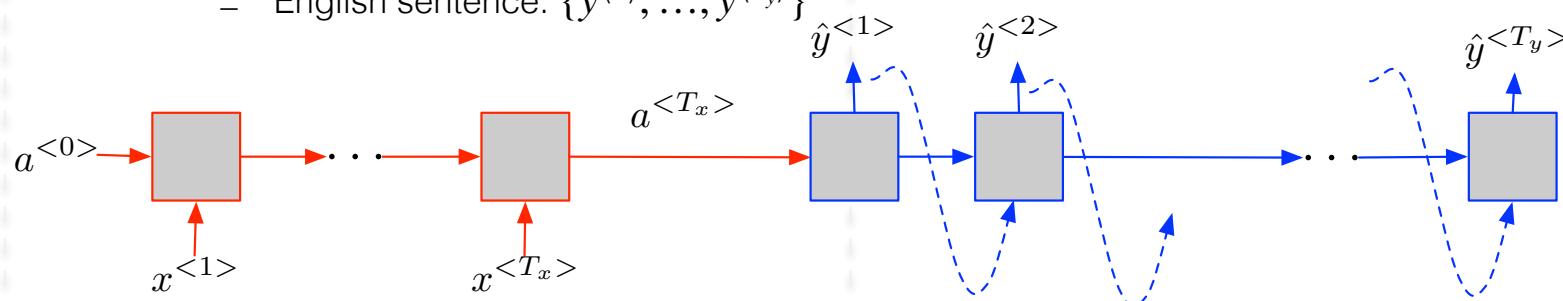
– Language model

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

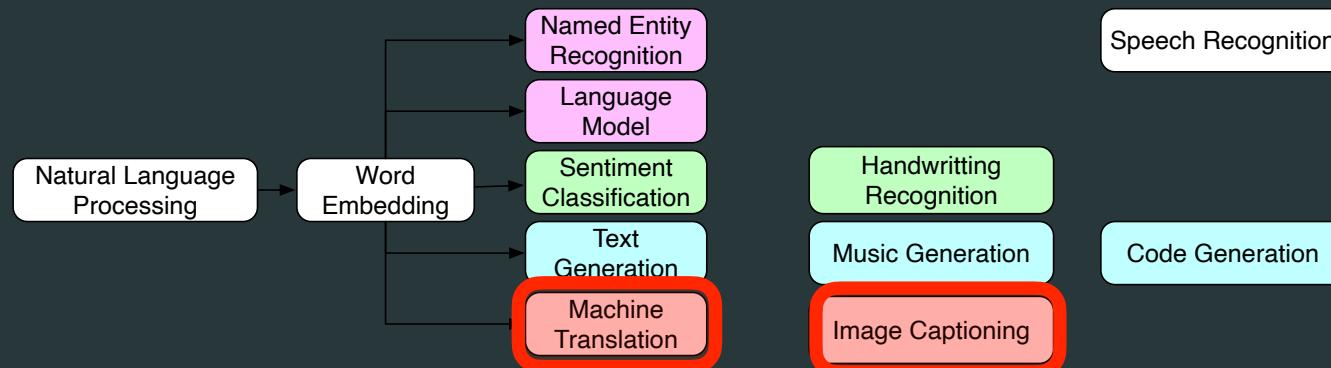
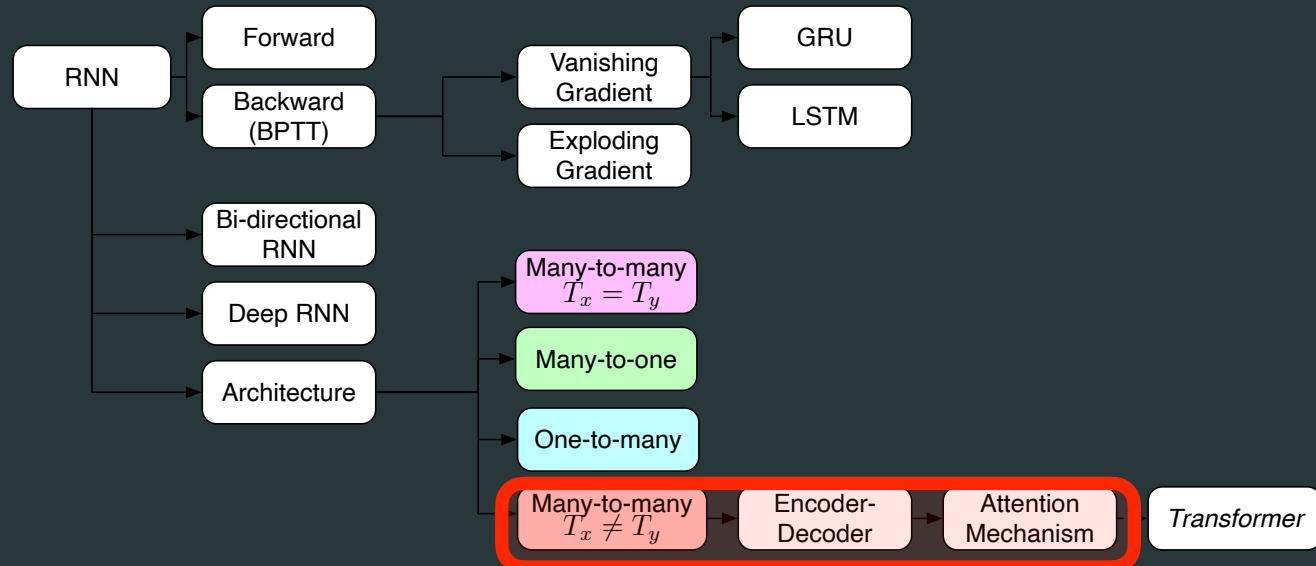


– Machine Translation:

- 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



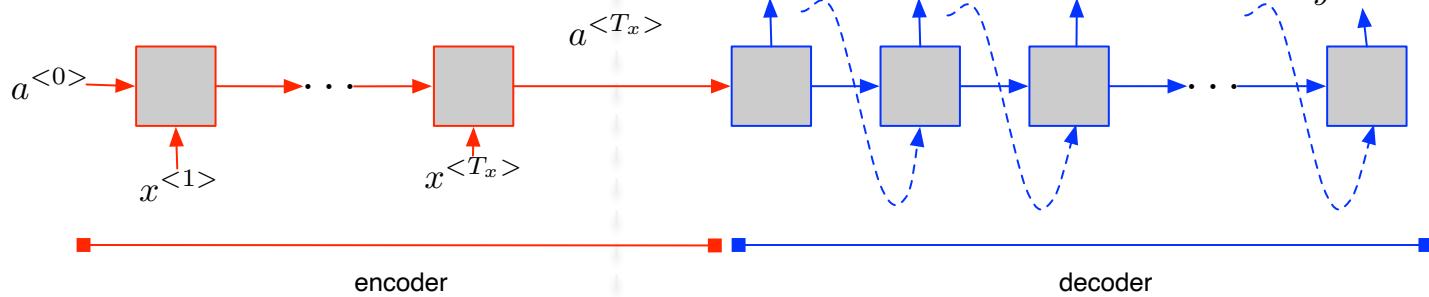
Sequence to sequence

Attention model

– 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.

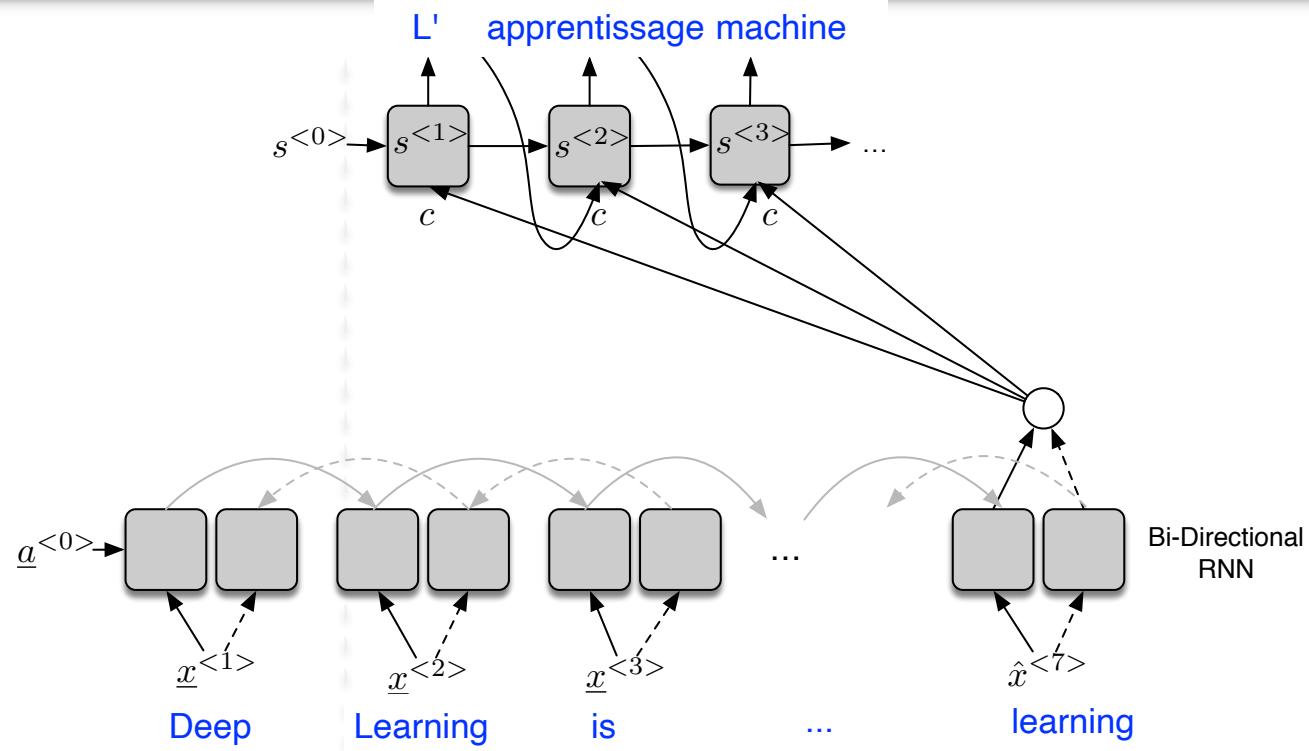
Ce type d'expérience fait partie des initiatives du Disney pour "prolonger la durée de vie de ses nouvelles et de développer des liens avec les lecteurs numériques qui deviennent plus complexes.



- Encoder/Decoder:
 - $a^{<T_x>}$ is supposed to memorise 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

Sequence to sequence

Attention model

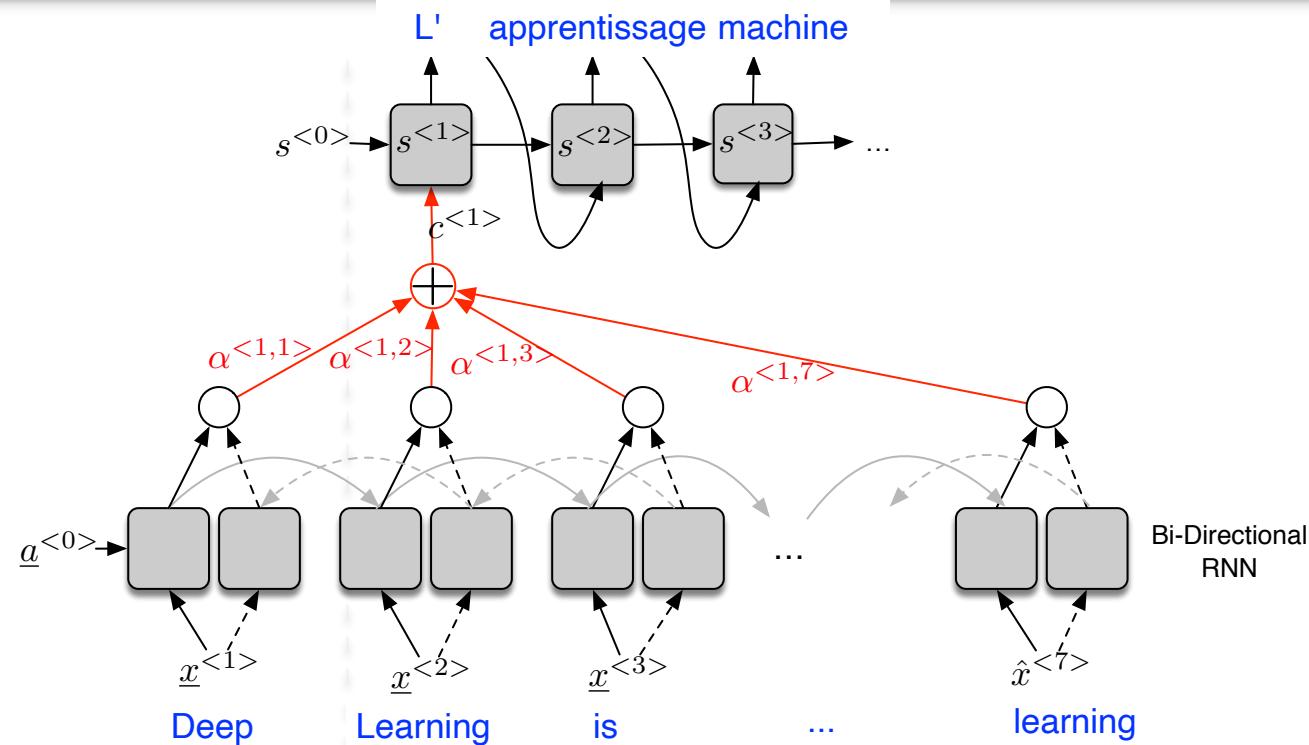


- In usual encoder-decoder,
 - information used for the decoding corresponds to the encoding at T_x : $c = a^{\langle T_x \rangle}$

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

Sequence to sequence

Attention model



– Attention model

- replace $c = a^{(T_x)}$ by a local version $c^{(\tau)}$
- $c^{(\tau)}$ 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 $t = \{1,2,3,\dots,7\}$

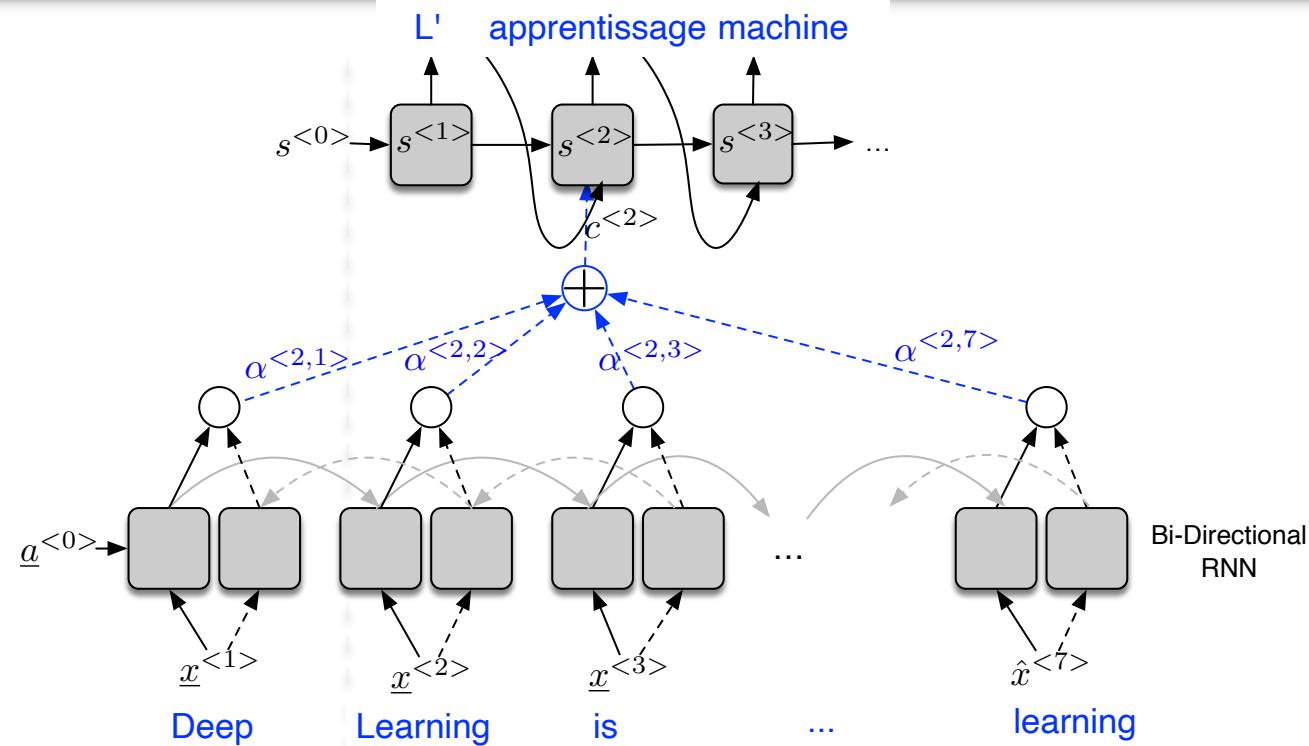
– Notation:

- $a^{(t)}, s^{(\tau)}$: hidden states of encoder/ decoder

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

Sequence to sequence

Attention model



– Attention model

- replace $c = a^{(T_x)}$ by a local version $c^{(\tau)}$
- $c^{(\tau)}$ 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 $t = \{1,2,3,\dots,7\}$

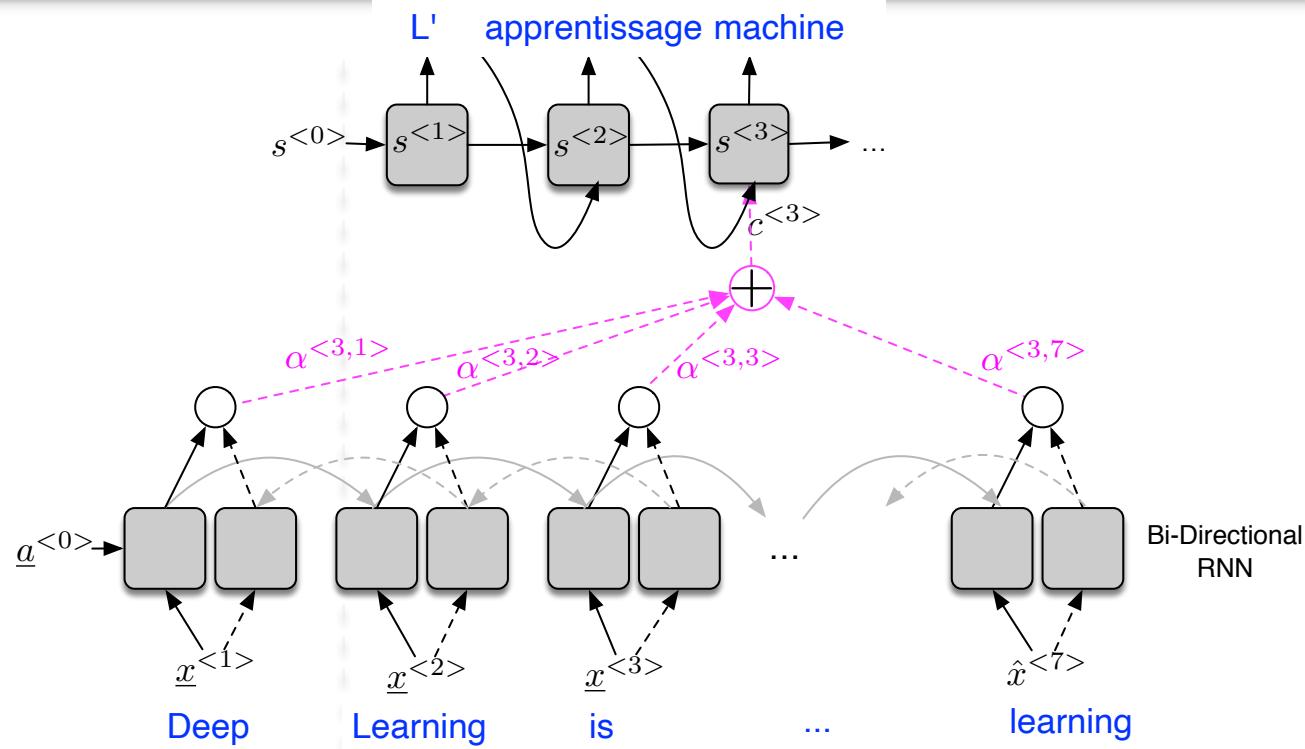
– Notation:

- $a^{(t)}, s^{(\tau)}$: hidden states of encoder/ decoder

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

Sequence to sequence

Attention model



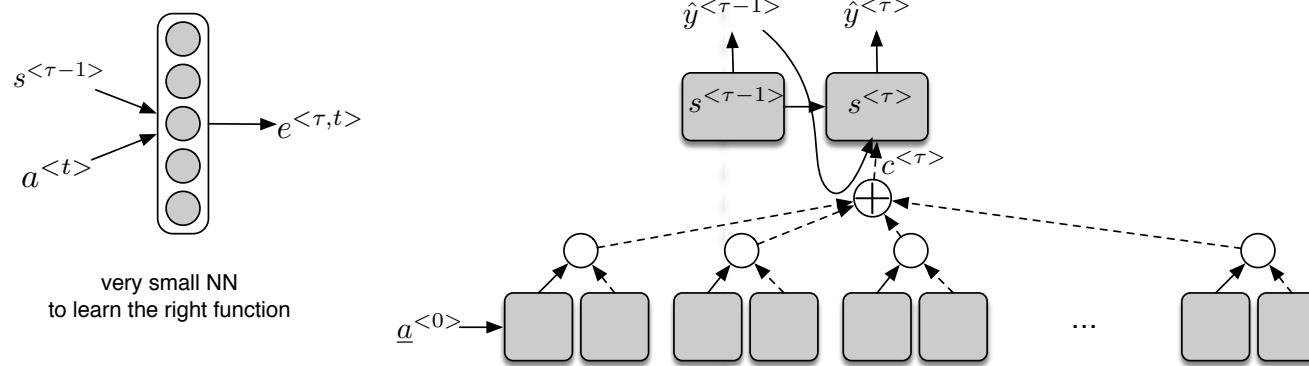
- Attention weights $\alpha^{(\tau=3,t)}$
 - describe an "alignment" between information at
 - **encoding time t** : computed using $\vec{a}^{(t)}$ and $\overleftarrow{a}^{(t)}$
 - we note $a^{(t)} = [\vec{a}^{(t)}, \overleftarrow{a}^{(t)}]$
 - **decoding time $\tau = 3$** : computed using $s^{(\tau=2)}$ (we do not yet observe $\tau = 3$)

Sequence to sequence

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** $\alpha^{(\tau,t)}$
- 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** $e^{(\tau,t)}$

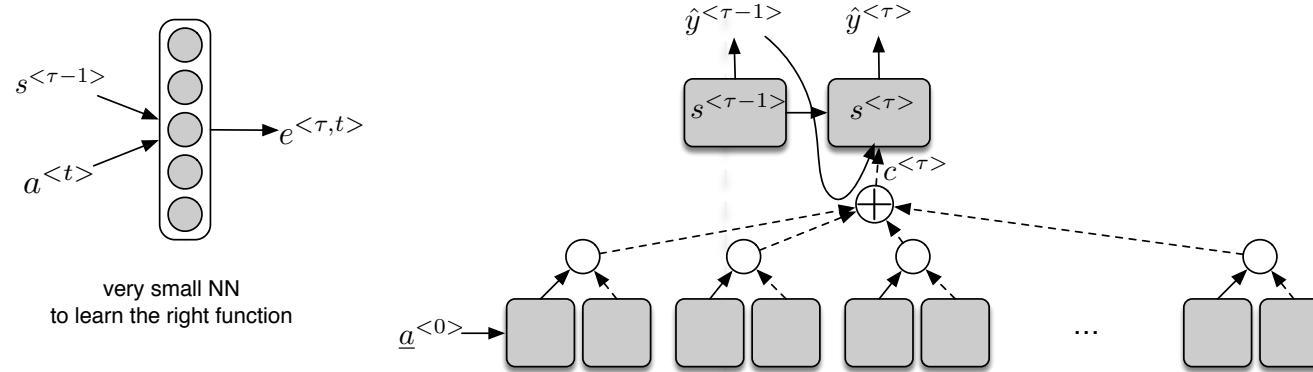
$$\alpha^{(\tau,t)} = \frac{\exp(e^{(\tau,t)})}{\sum_{t'=1}^{T_x} \exp(e^{(\tau,t')})} \text{ with } \sum_t \alpha^{(\tau,t)} = 1$$



Sequence to sequence

Attention model

- 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)



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

Sequence to sequence

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.

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

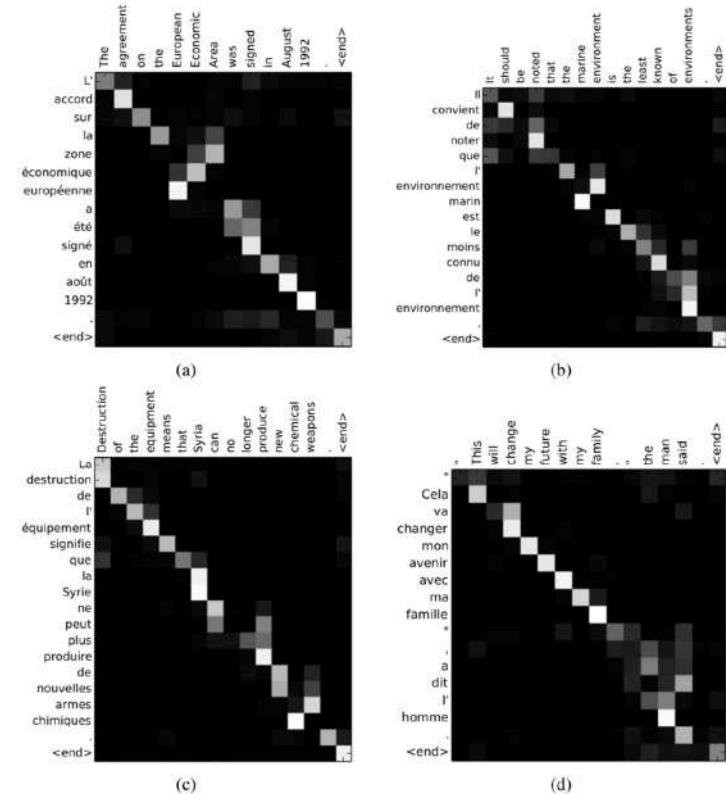
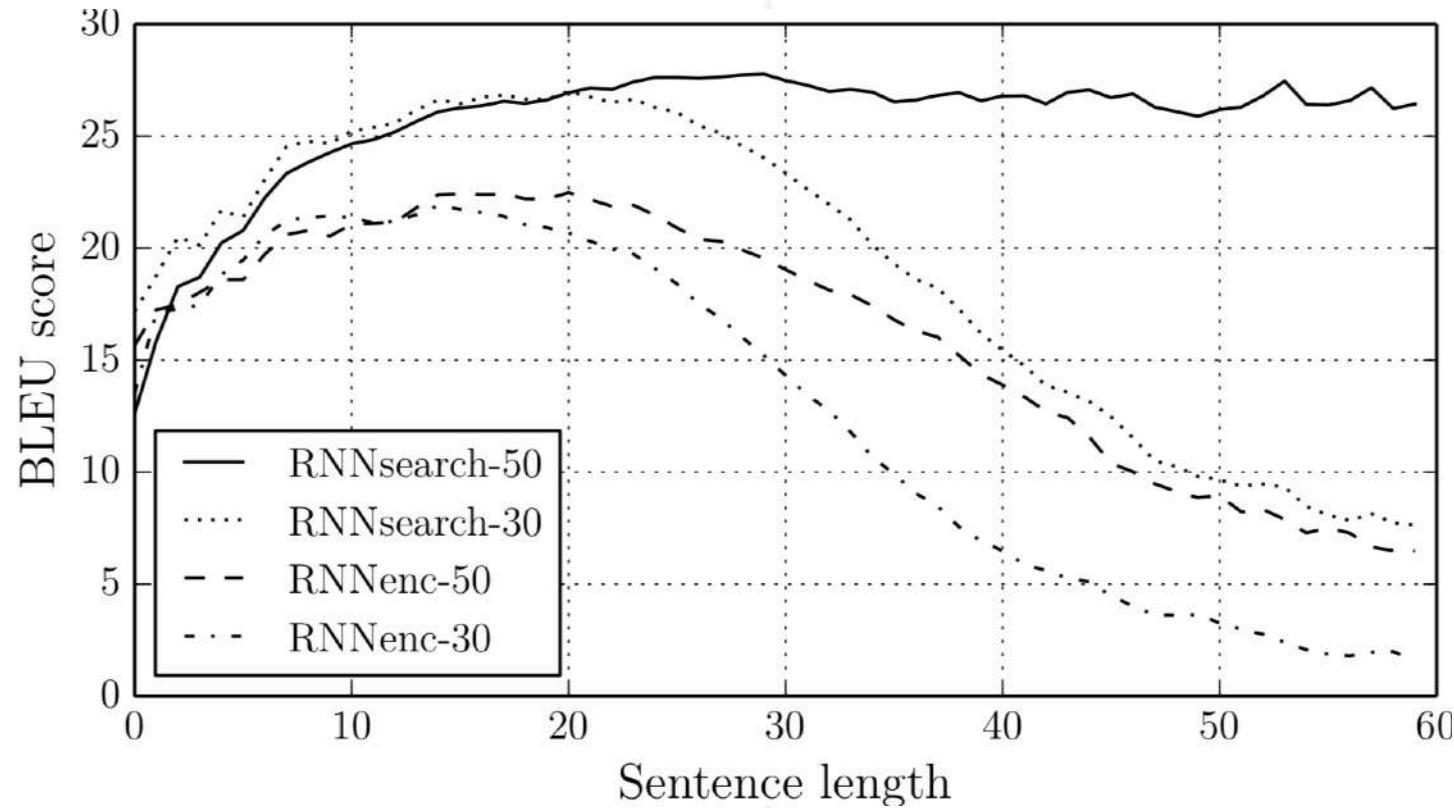


Figure 3: Four sample alignments found by RNNsearch-50. The x-axis and y-axis of each plot correspond to the words in the source sentence (English) and the generated translation (French), respectively. Each pixel shows the weight α_{ij} of the annotation of the j -th source word for the i -th target 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.

Sequence to sequence

Attention model

- Bleu score for long sentences



Applications

Automatic Speech Recognition



Automatic picture captioning



A group shopping at an outdoor market.

There are many vegetables at the fruit stand.

Self Driving Cars



DALL-E



DEEP LEARNING HAS MASTERED GO

Google Alpha Go



Style Transfer

Zebras ↡ Horses



Neural Machine Translation

Anglais (langue détectée) ▾

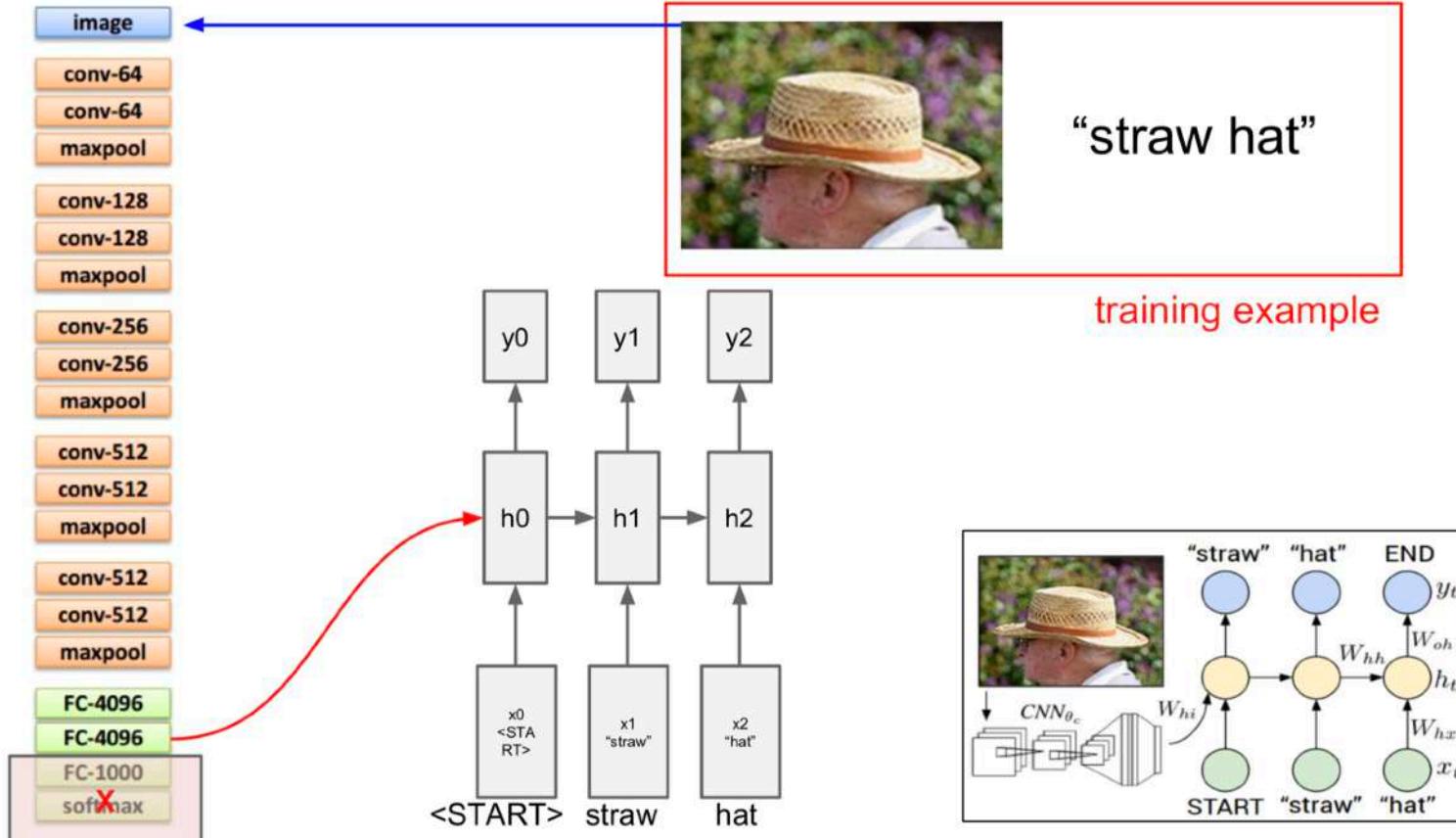
DeepL Traducteur DeepL Pro Connexion

Découvrir DeepL Pro

L'apprentissage profond (également appelé apprentissage structuré profond) fait partie d'une famille plus large de méthodes d'apprentissage automatique basées sur les réseaux neuronaux artificiels avec apprentissage par représentation. L'apprentissage peut être supervisé, semi-supervisé ou non supervisé.

Sequence to sequence

Application: Image captioning



SOURCE: http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture10.pdf

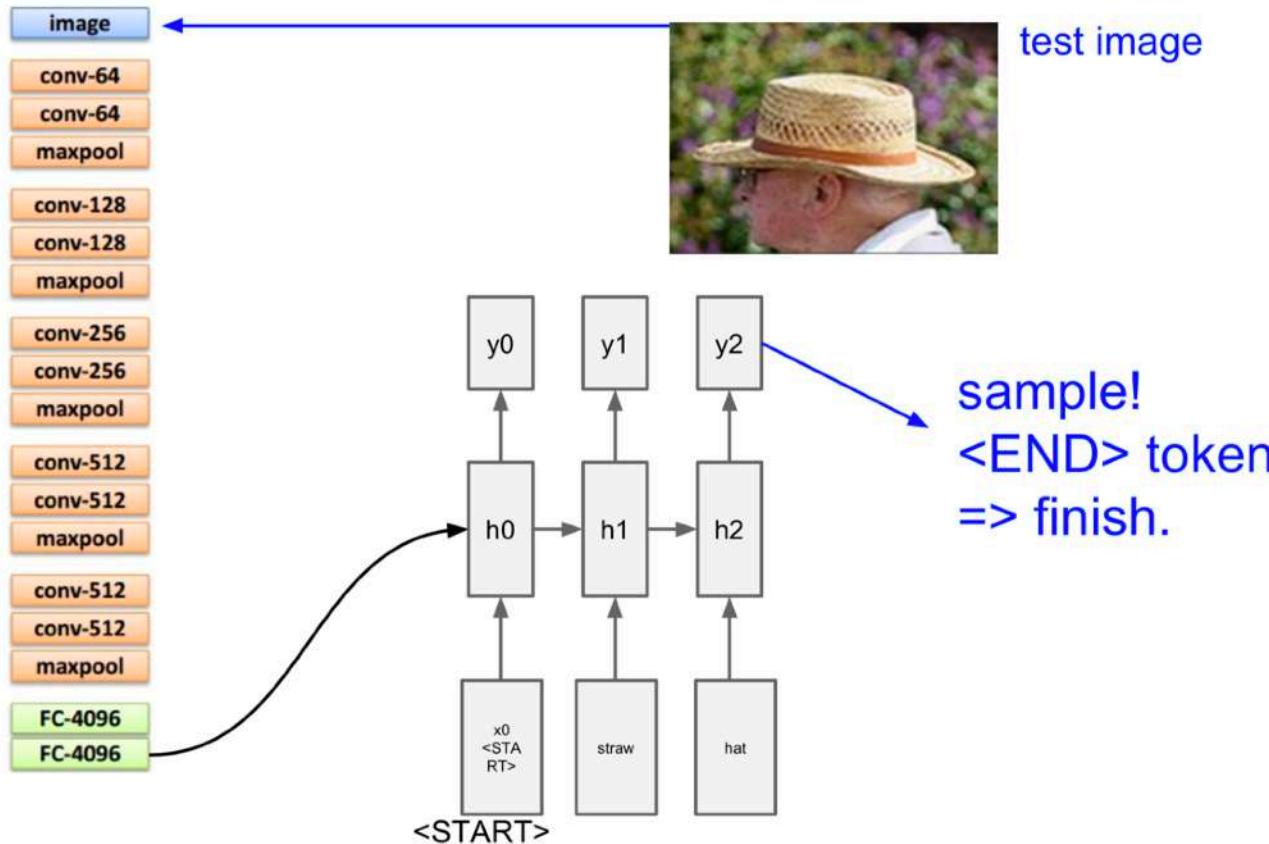
[Mao et al. 2014 "Deep captioning with multimodal recurrent neural networks (m-RNN)"] [LINK](#)

[Vinyals et al. 2015 "Show and tell: A neural image caption generator"] [LINK](#)

[Karpathy et al. 2015 "Deep visual-semantic alignments for generating image descriptions"] [LINK](#)

Sequence to sequence

Application: Image captioning



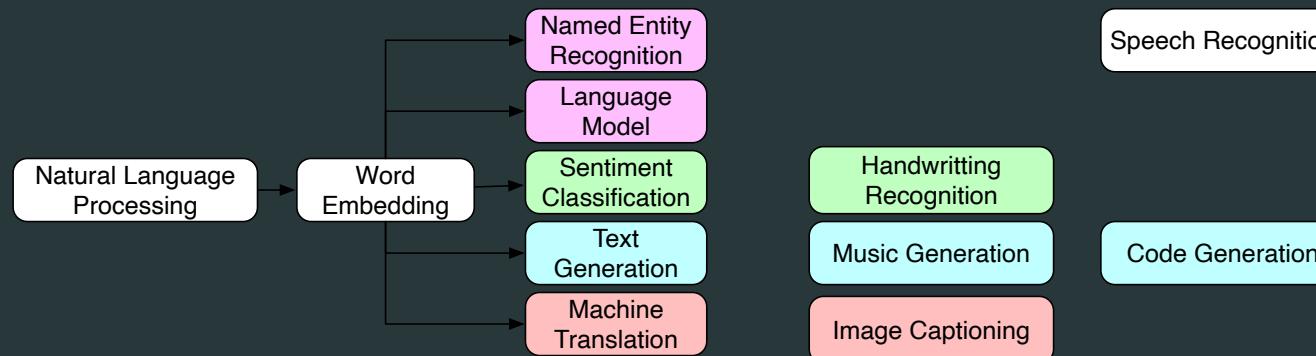
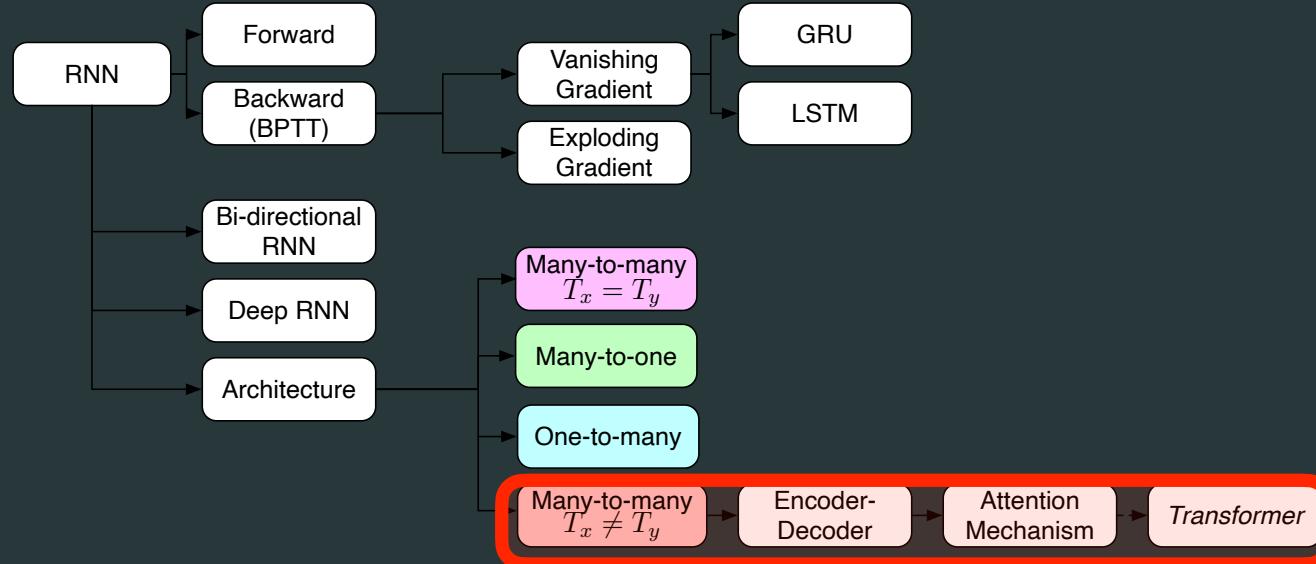
SOURCE: http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture10.pdf

[Mao et al. 2014 "Deep captioning with multimodal recurrent neural networks (m-RNN)"] [LINK](#)

[Vinyals et al. 2015 "Show and tell: A neural image caption generator"] [LINK](#)

[Karpathy et al. 2015 "Deep visual-semantic alignments for generating image descriptions"] [LINK](#)

Transformer

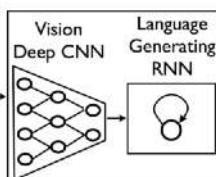


Applications

Automatic Speech Recognition



Automatic picture captioning



A group shopping at an outdoor market.
There are many vegetables at the fruit stand.

Self Driving Cars



DALL-E



DEEP LEARNING HAS MASTERED GO

Google Alpha Go

Google reveals secret test of AI bot to beat top Go players

Mastering the game of Go with deep neural networks and tree search

David Silver, Aljaž Petričević, Chris J. Maddison, Arthur Guez, Laurent Béjan, Oriol Vinyals, Julian Schrittwieser, Sophie Antiga, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Mnih, Noa Kischner, Ilya Sutskever, Thore Graepel & Demis Hassabis

Nature 529, 484–491 (28 January 2016) | doi:10.1038/nature16911
Received: 11 November 2015 | Accepted: 05 January 2016 | Published online: 27 January 2016



Neural Machine Translation

Anglais (langue détectée) ▾ Fr... ▾ forme/informel ▾ Glossaire

Deep learning (also known as deep structured learning) is part of a broader family of machine learning methods based on artificial neural networks with representation learning. Learning can be supervised, semi-supervised or unsupervised.

L'apprentissage profond (également appelé apprentissage structuré profond) fait partie d'une famille plus large de méthodes d'apprentissage automatique basées sur les réseaux neuronaux artificiels avec apprentissage par représentation. L'apprentissage peut être supervisé, semi-supervisé ou non supervisé.

Style Transfer



Sequence-to-Sequence

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

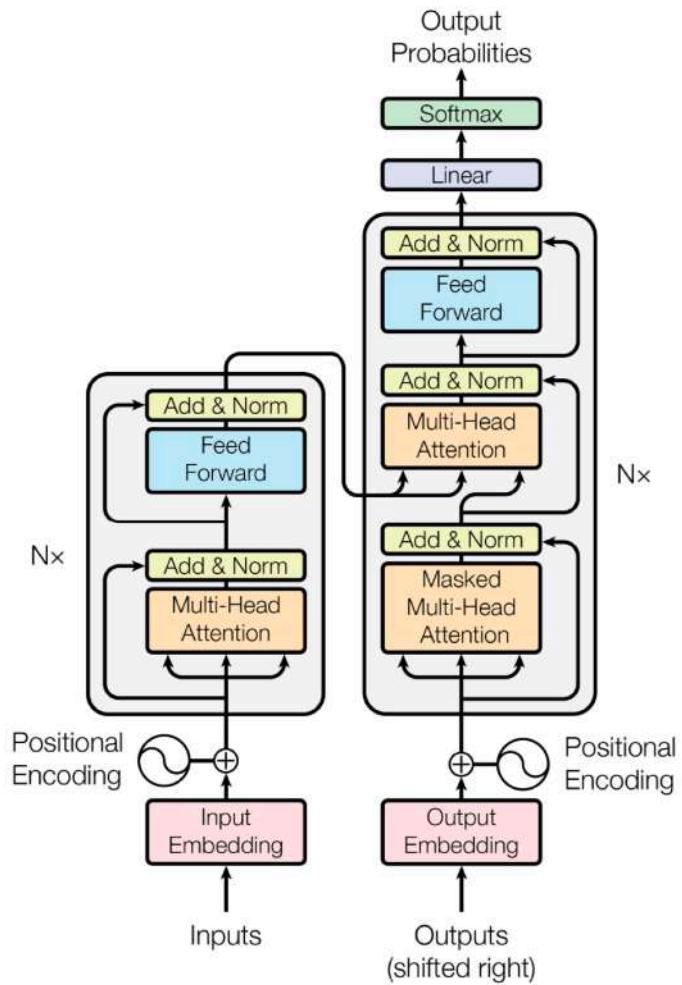


Figure 1: The Transformer - model architecture.

Transformer

- **Self-Attention :**

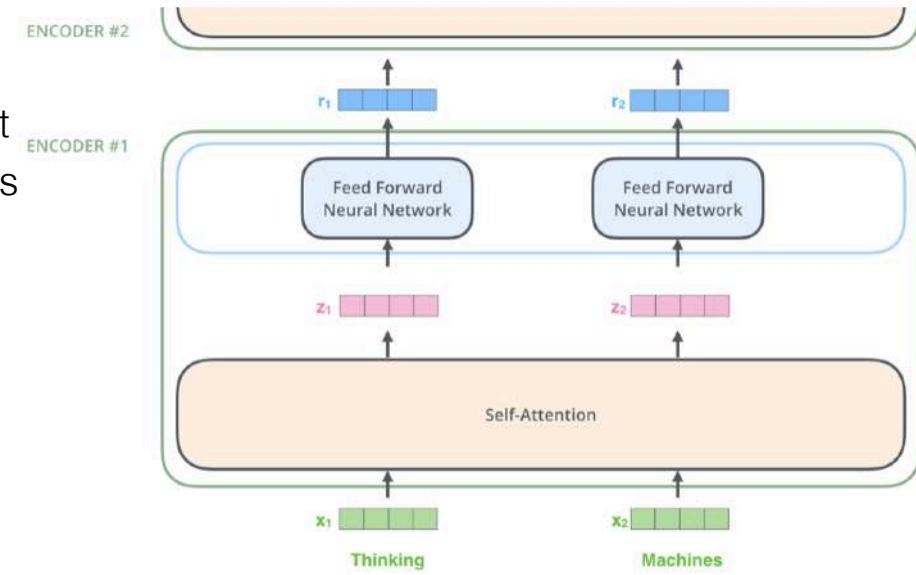
- Input sentence we want to translate : "*The animal didn't cross the street because it was too tired*"
 - What does "*it*" in this sentence refer to ? Is it referring to the "*street*" or to the "*animal*" ?
- When the model is processing the word "*it*", self-attention allows it to associate "*it*" 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

Sequence-to-Sequence

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

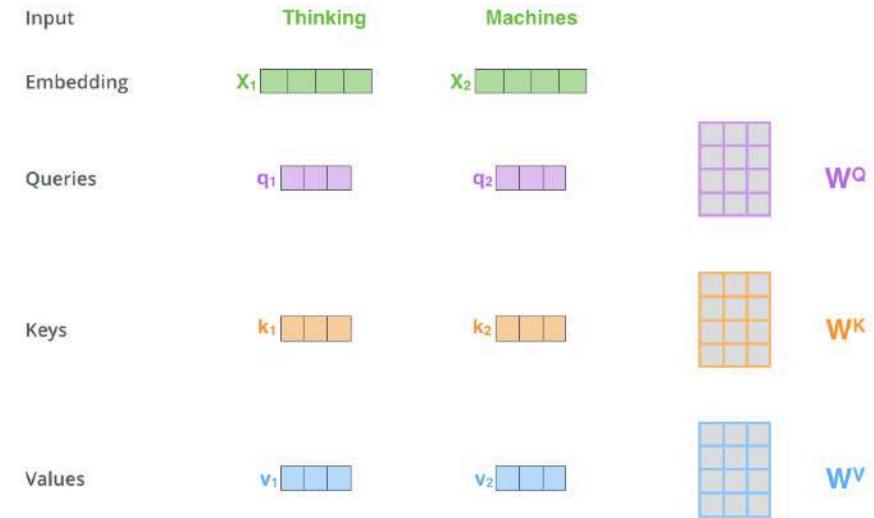


Sequence-to-Sequence

Transformer

- **First step :**

- create three vectors from each of the encoder's input vectors x_i :
 - a Query vector q_i ,
 - a Key vector k_i
 - a Value vector v_i
- the vectors are created by multiplying the input by three matrices W^Q, W^K, W^V that are trained during the training process
- the new vectors are smaller in dimension than the embedding vector



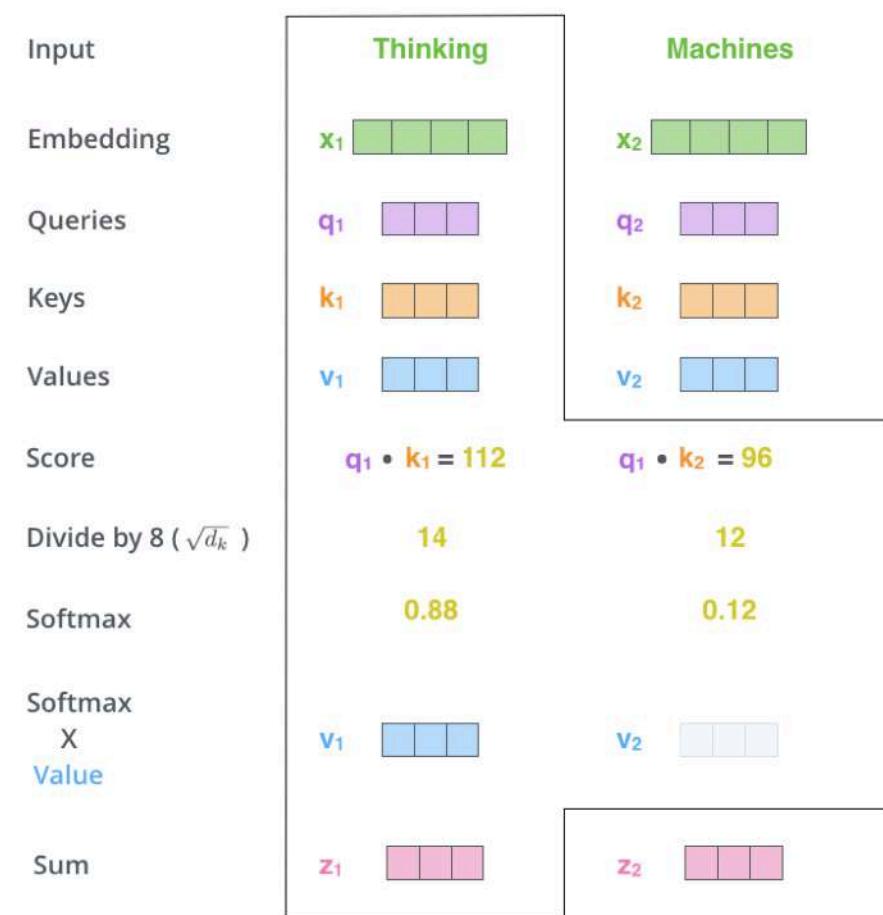
Multiplying x_1 by the W^Q weight matrix produces q_1 , the "query" vector associated with that word. We end up creating a "query", a "key", and a "value" projection of each word in the input sentence.

Sequence-to-Sequence

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 q_i with the key vector k_i 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

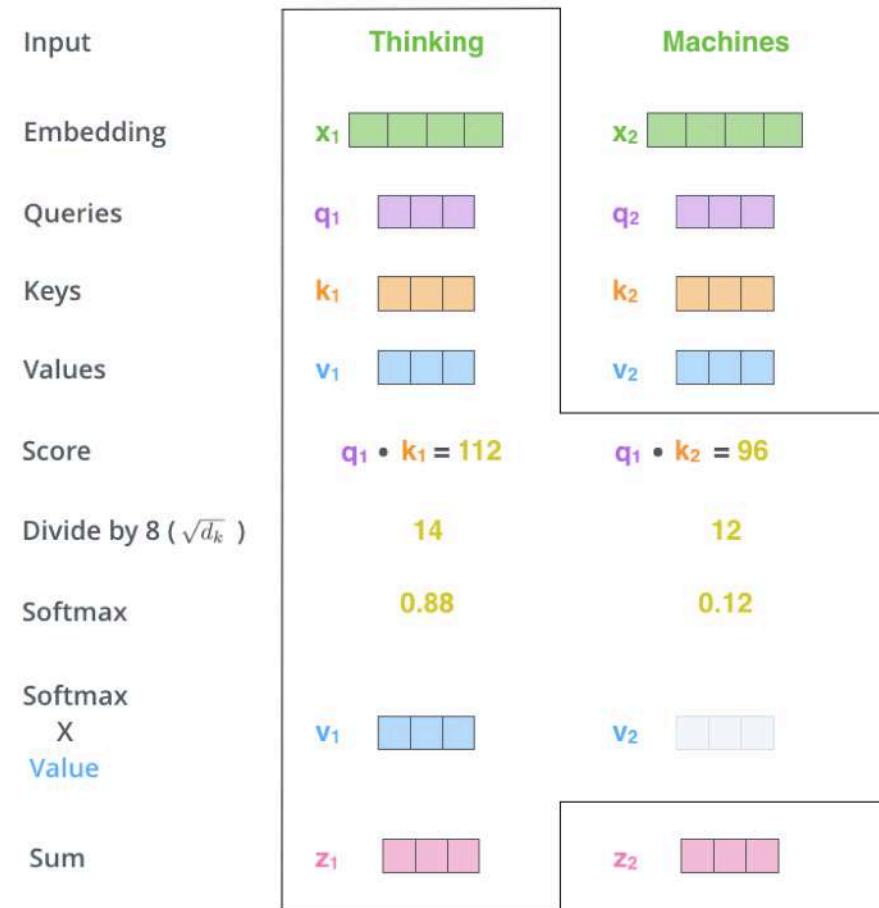


Sequence-to-Sequence

Transformer

- **Third step :**

- multiply each **value** v_i 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 z_i of the self-attention layer at this position (for the first word).

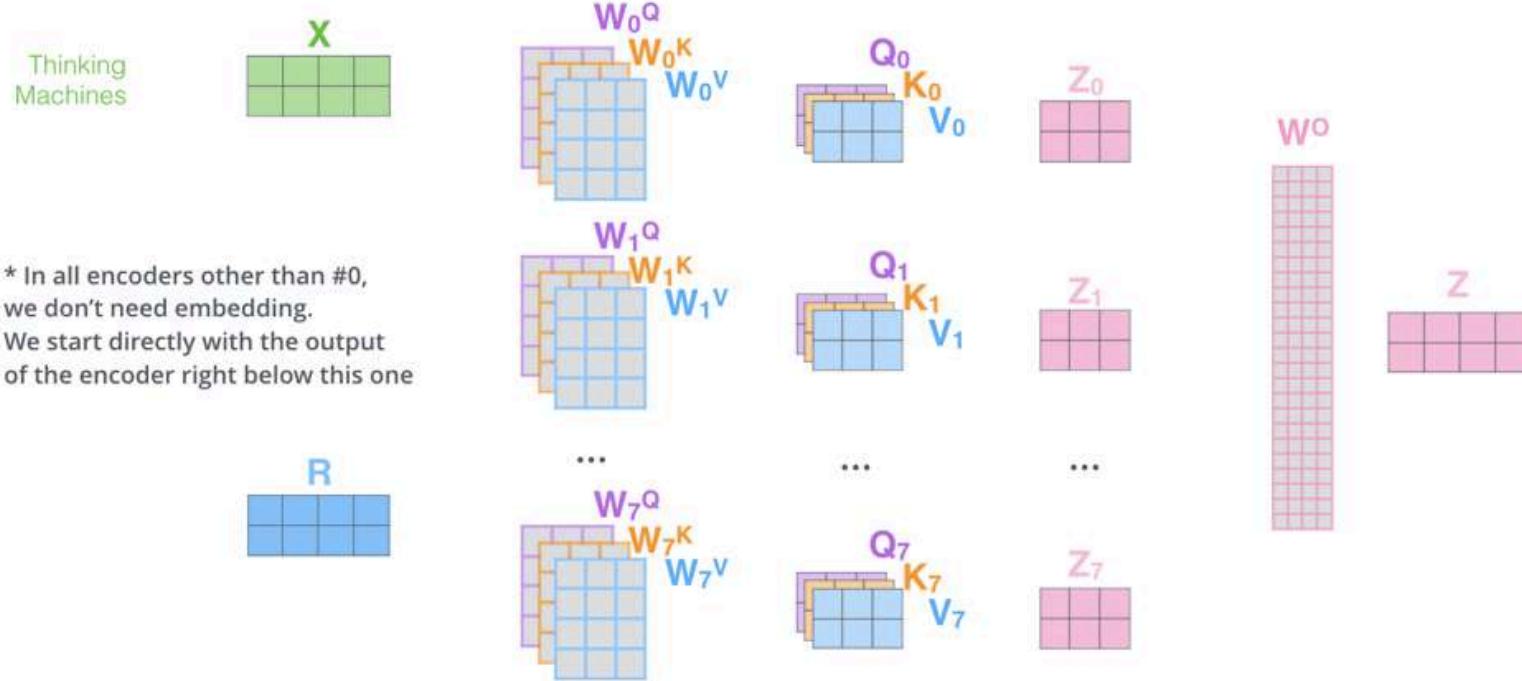


Sequence-to-Sequence

Transformer

- In summary

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

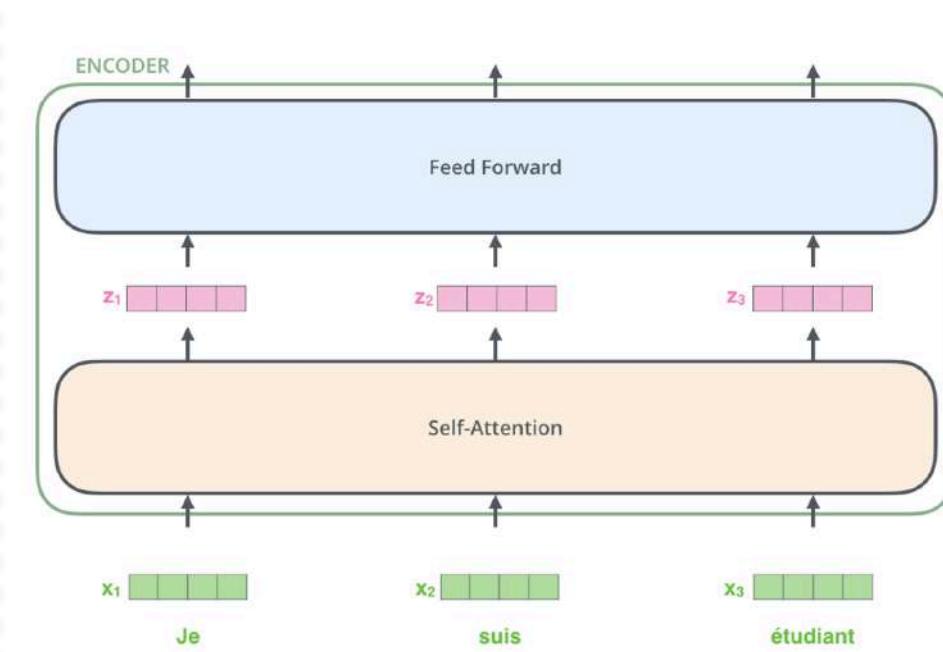


Sequence-to-Sequence

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

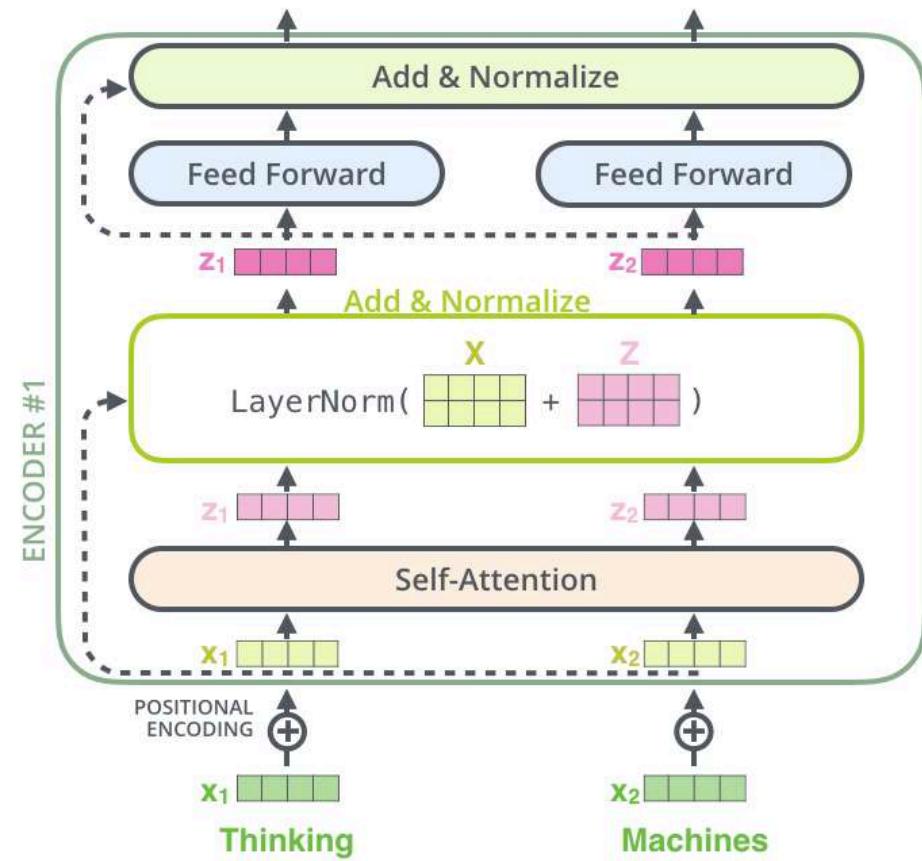


Sequence-to-Sequence

Transformer

- **Residuals + Normalisation**

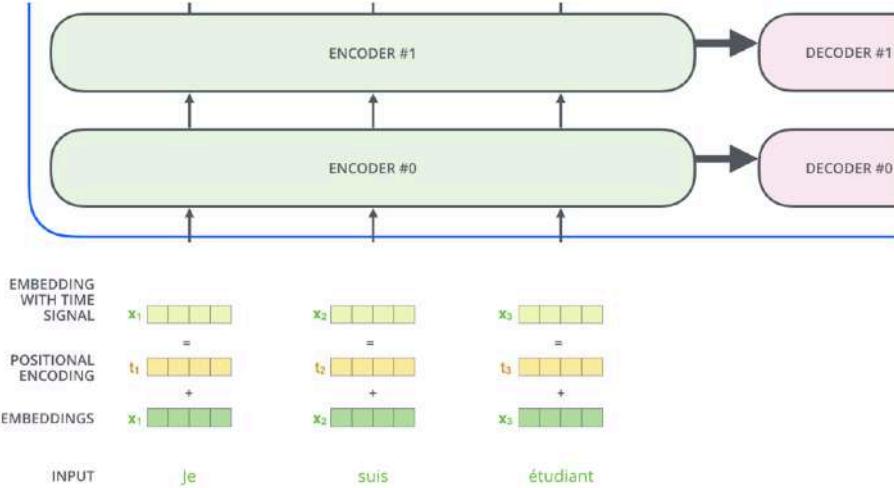
- each sub-layer (self-attention, feed-forward) in each encoder has
 - a residual connection around it
 - is followed by a layer-normalisation step



Sequence-to-Sequence

Transformer

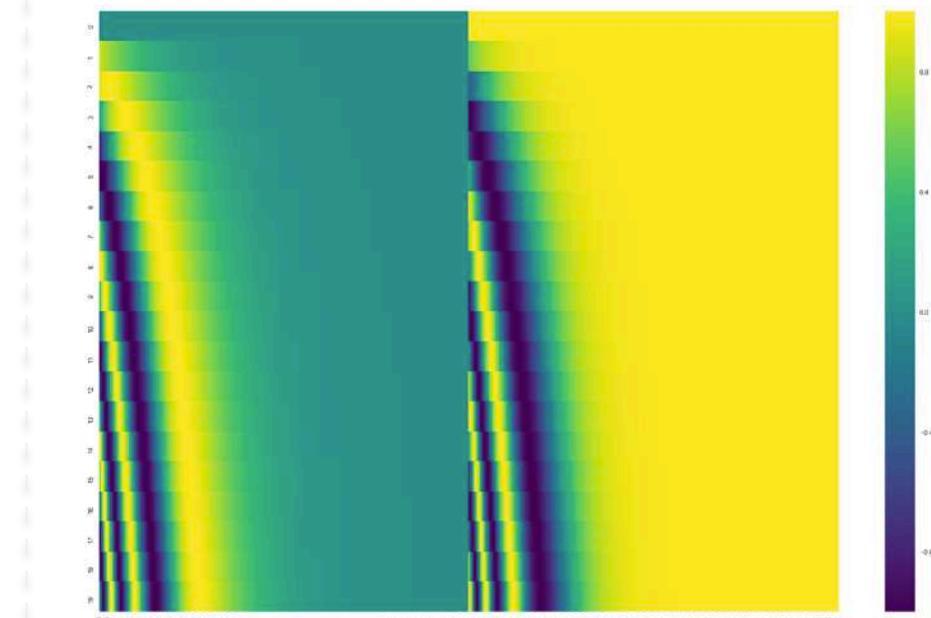
- **Representing The Order of The Sequence Using Positional Encoding**
 - Find a way to account for the order of the words in the input sequence
 - Transformer adds a vector to each input embedding
 - These vectors follow a specific pattern that the model learns, which helps it determine the position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



Sequence-to-Sequence

Transformer

- **Representing The Order of The Sequence Using Positional Encoding**
 - What might this pattern look like ?

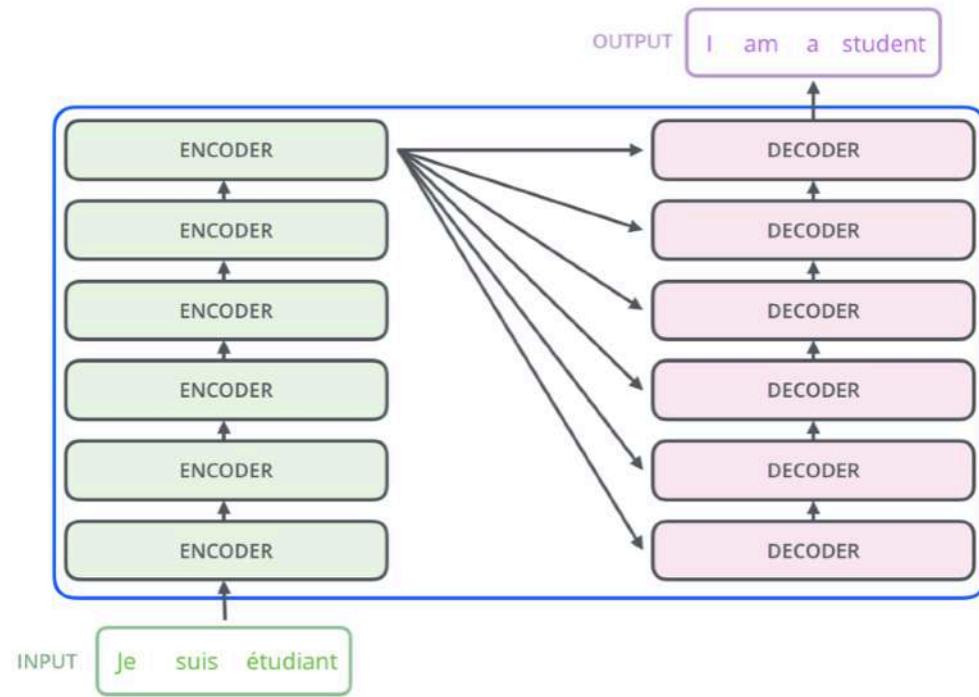


A real example of positional encoding for 20 words (rows) with an embedding size of 512 (columns). You can see that it appears split in half down the center. That's because the values of the left half are generated by one function (which uses sine), and the right half is generated by another function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

Sequence-to-Sequence

Transformer

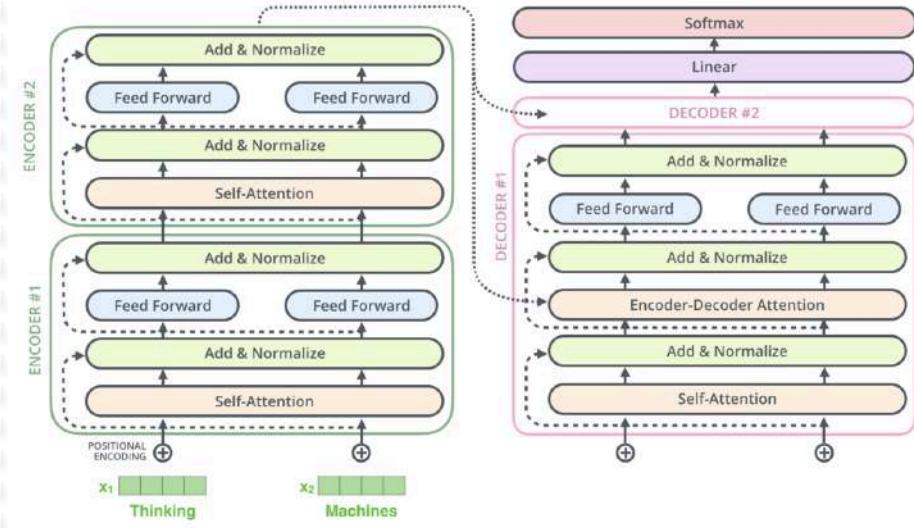
- The encoding component is a stack of 6 encoders on top of each other
- The decoding component is a stack of decoders of the same number



Sequence-to-Sequence

Transformer

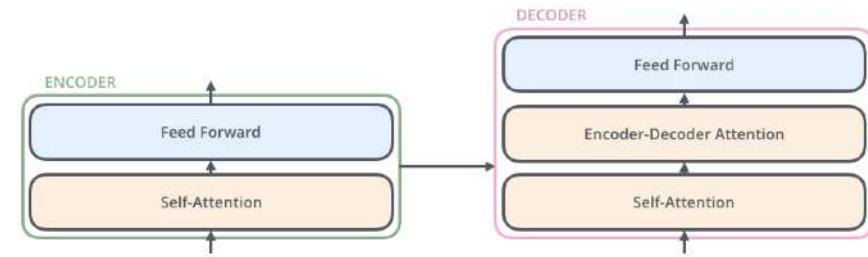
- **Plugin the encoder to the decoder**



Sequence-to-Sequence

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)
-

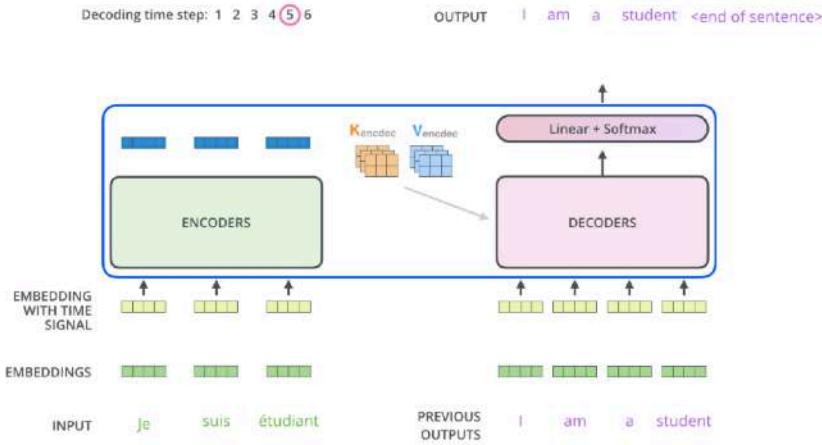


Sequence-to-Sequence

Transformer

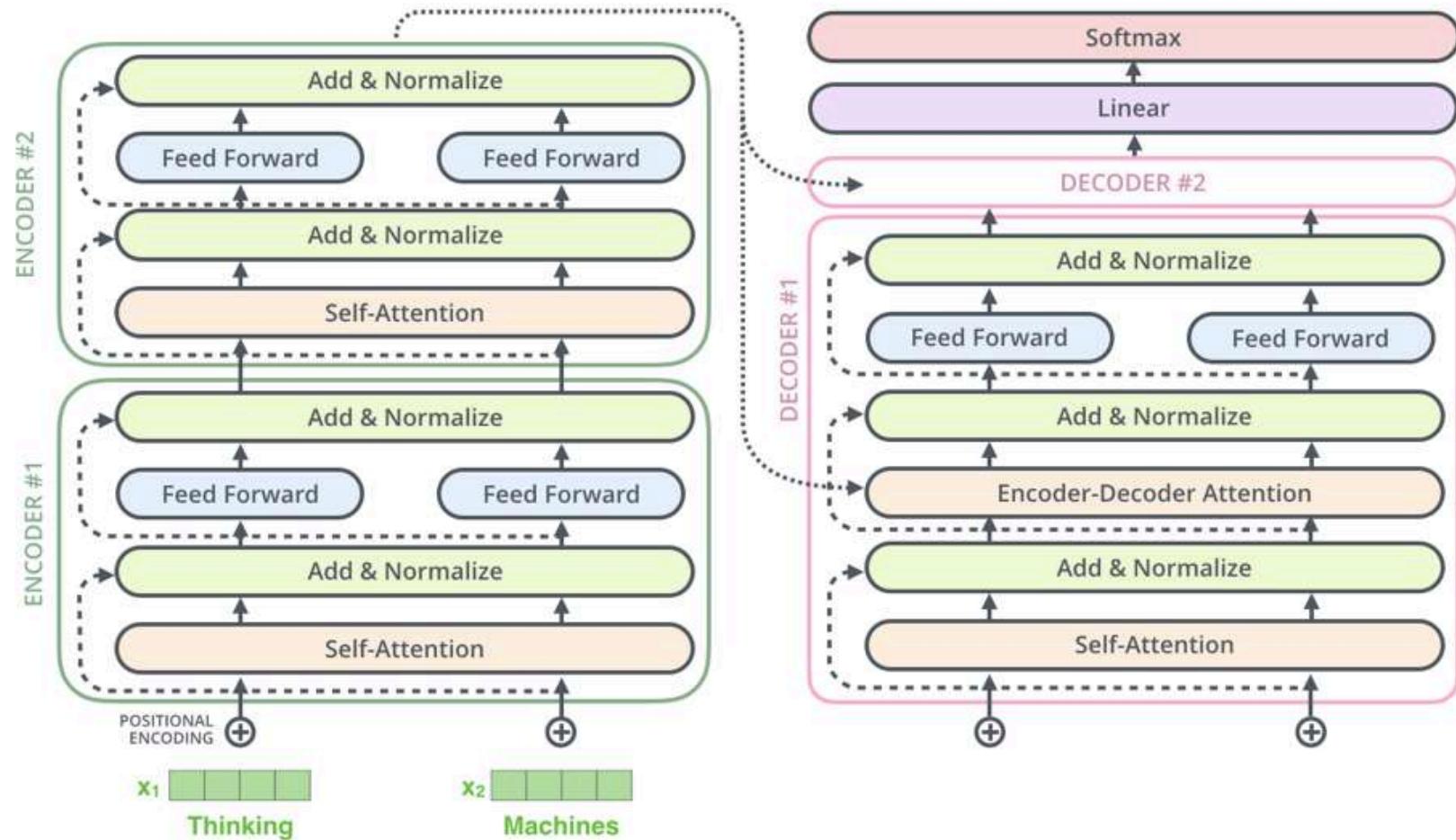
- **Decoder side**

- 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.



Sequence-to-Sequence

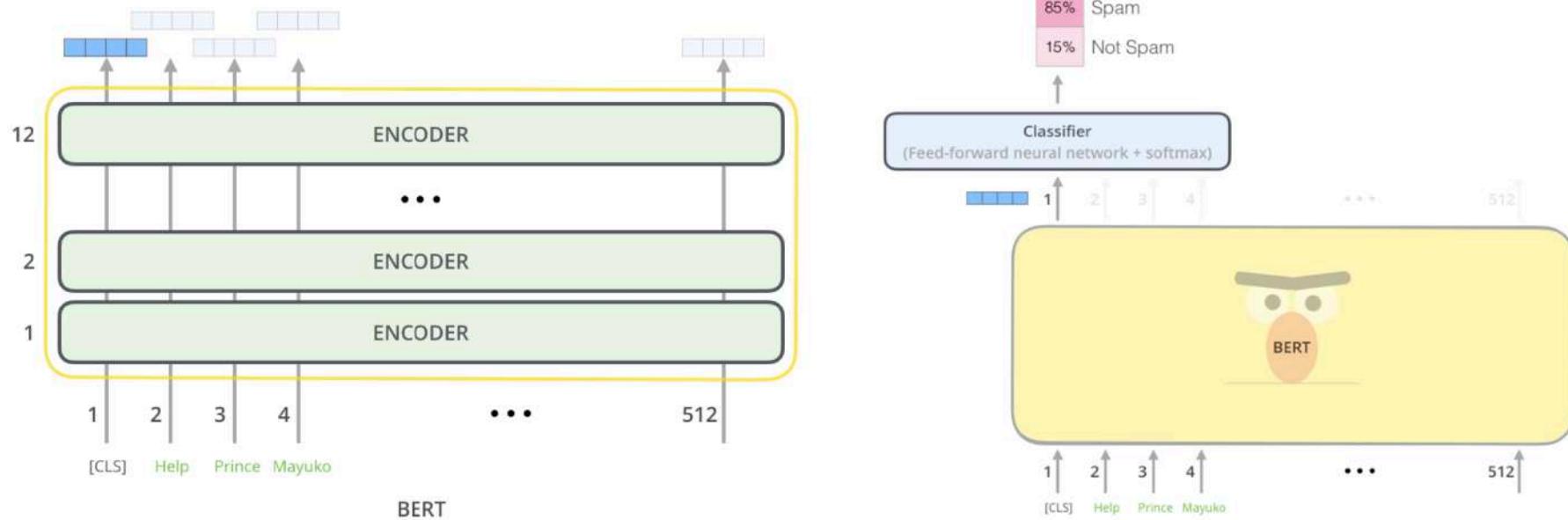
Transformer



Sequence-to-Sequence

BERT

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



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

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

Overview

