

From Recurrent Neural Networks to Transformers

Advanced Machine Learning

Master degree in Computer Science

Nicoletta Noceti

Machine Learning Genoa Center – University of Genoa

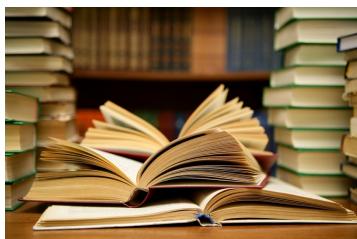
Sequential data and problem formulation

Dealing with sequential data

- Today we consider data in the form of sequences



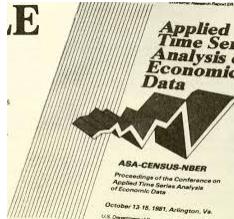
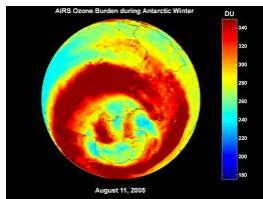
Audio



Text

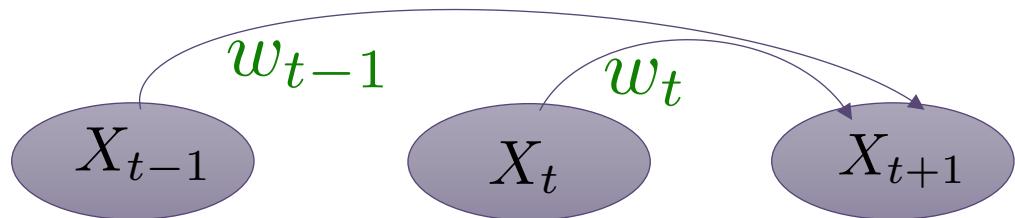


IoT

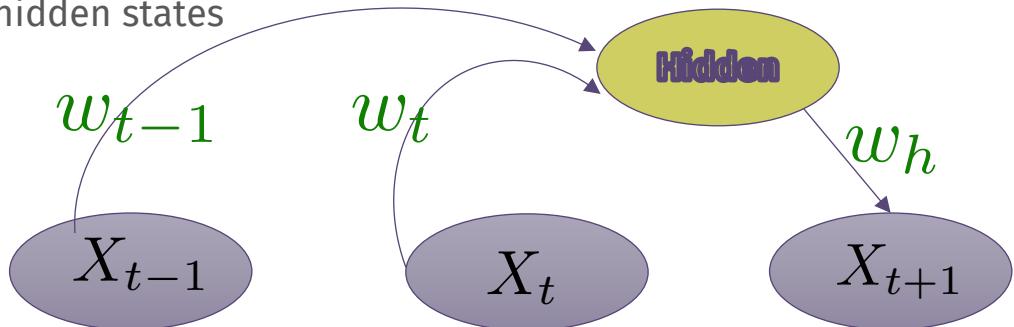


Dealing with sequences: a first summary

- Autoregressive models



- “Using” hidden states



Different formulations of the problem

- Let us introduce a variety of problem formulations involving a special kind of sequential data, i.e. text

“This morning I am happy since I am taking my dog for a walk”

- What are the problems that we may want to address?

“This morning I am happy since I am taking my dog for a walk”

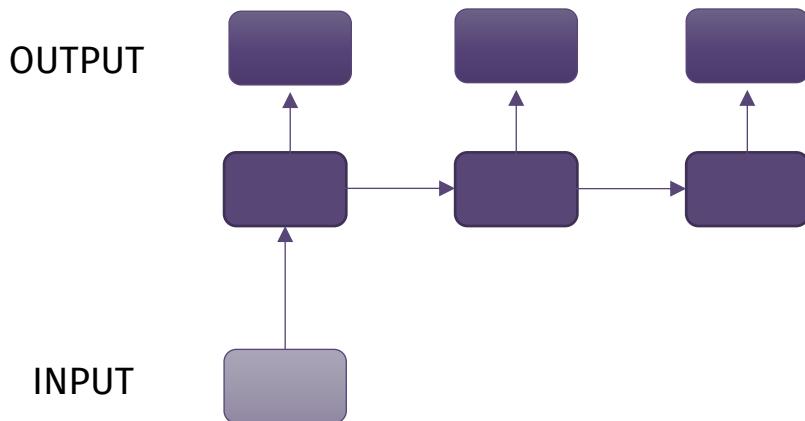


“This morning I am happy since I am taking my dog for a walk”

Different formulations of the problem

One-to-many: the input is in a standard format (not a sequence!), the output is a sequence

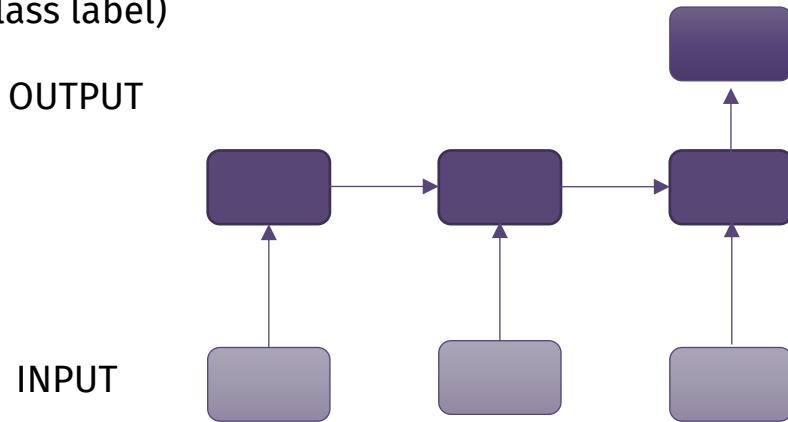
Example of applications: image captioning (input: image, output: text describing the image content)



Different formulations of the problem

Many-to-one: the input is a sequence, the output is a fixed-size vector (not a sequence!)

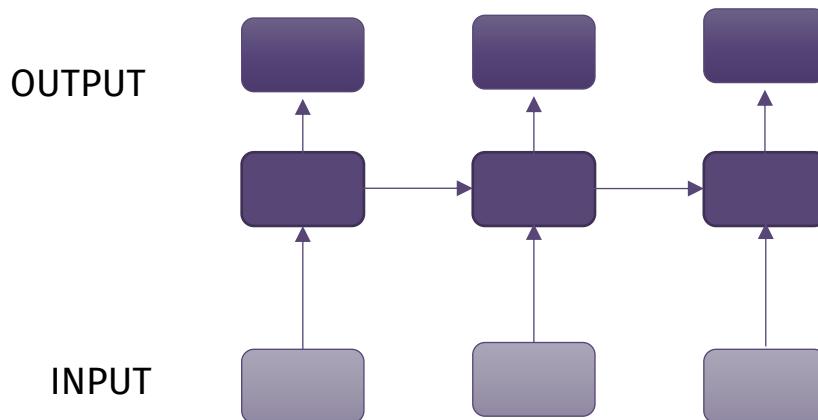
Example of applications: sentiment analysis (input: text, output: class label)



Different formulations of the problem

Direct Many-to-many: input and output are both sequences

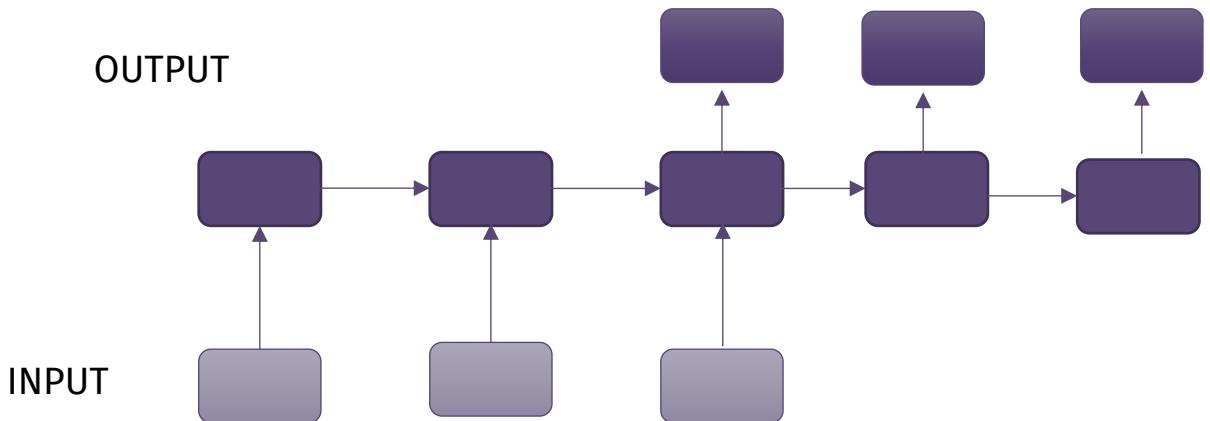
Example of applications: video captioning (input is a sequence of images, output is text)



Different formulations of the problem

Delayed Many-to-many: input and output are both sequences

Example of applications: language translation (input is a text, output is a text)



A special case: text data

- ML can not (directly) handle text data... we need a numerical descriptor
- Word embeddings can do the job

Vocabulary

my this I
walk taking
a the since
 am dog

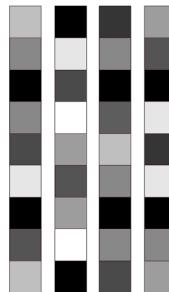
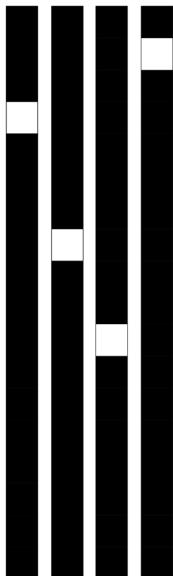
Indexing

a → 0
am → 1
dog → 2
...
walk → N

Embeddings

$A = [1000\dots 0]$	happy	walk
$Am = [0100\dots 0]$	sad	run
$dog = [0010\dots 0]$	dog	day
$Walk = [0000\dots 1]$	cat	night

One-hot encoding vs word embeddings



One-hot encoding:

- Very sparse
- High dimensional
- Hard-coded

Word embeddings:

- Dense
- Lower dimensional
- Learned from data

An example

An example: predicting the next word

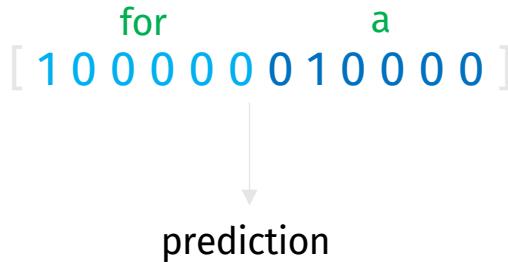
“This morning I took my dog for a walk”

- Predicting the last word in the sentence given all or some of the previous words
- Let's reason on possible solutions

Solution 1

*"This morning I took my dog **for a** walk"*

- Using small fixed windows and one-hot encoding



Problem: what if the few words are not informative?

e.g. "Italy is where I grew up but now I live in London. I speak fluent _____"

Longer-term dependences may be needed!

Solution 2

- Using all the available words to form an embedding

“This morning I took my dog for a walk”

[1 0 1 0 1 1 0 1 0 0 0 0 0 1 1 0 0 1 0 0]



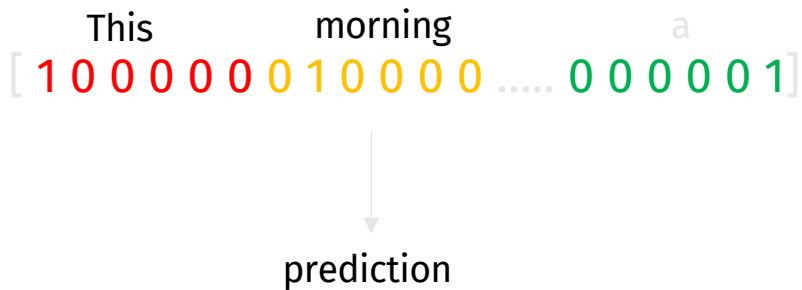
prediction

Problem: order is not preserved

e.g. “The food was good, not bad at all” vs “The food was bad, not good at all”

Solution 3

- Going back to solution 1 with larger windows



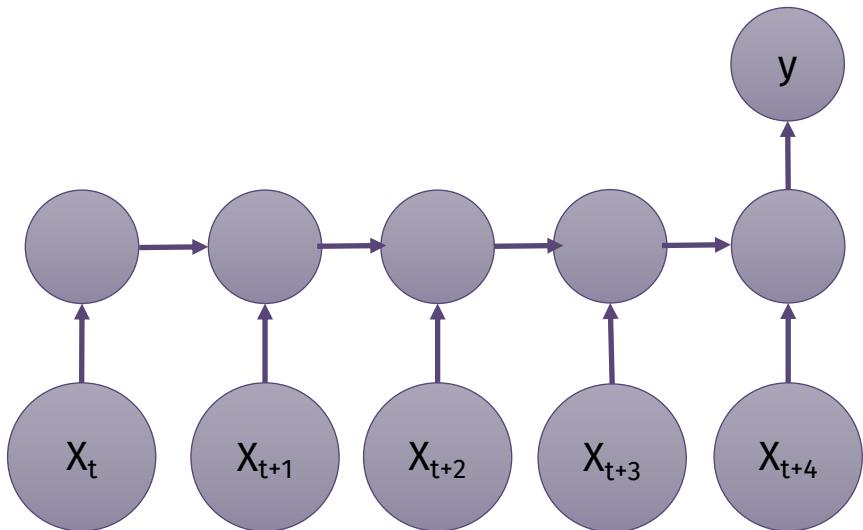
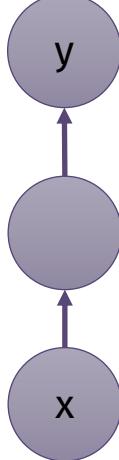
Problem: no parameter sharing, and the knowledge we may derive does not transfer if words appear elsewhere

How to approach sequence modelling

- Handling sequences of **different lengths**
- Taking into account **short** and **long term** dependences
- Considering **order** between elements
- **Sharing parameters** across the sequence

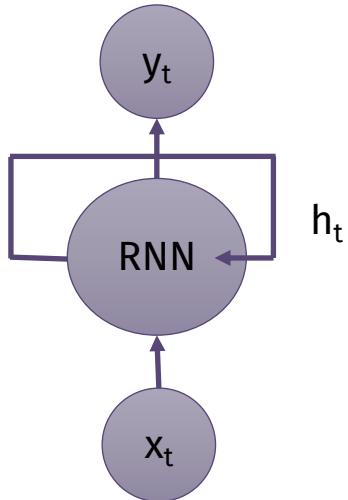
Recurrent Neural Networks

Modelling sequences



Standard one-to-one vanilla network

Recurrent Neural Networks



A recurrence relation is applied at each time step to model the sequence

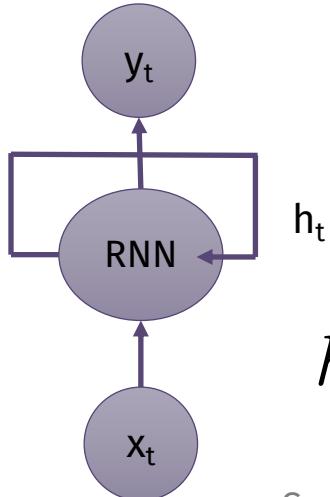
$$h_t = f_W(h_{t-1}, x_t)$$

Same function and weights are used for each time step, so they are a way to share weights over time

Recurrent Neural Networks (1986)

- Recurrence adds memory to the NN
- It also provides a way to model causal relationships between observations: the decision a recurrent net reached at time step $t-1$ affects the decision it will reach at time step t
- RNNs have two sources of input: the present and the recent past, which combine to determine how they respond to new data
- It is finding correlations between events separated by many moments, and these correlations are called “*long-term dependencies*”

RNNs: Forward propagation



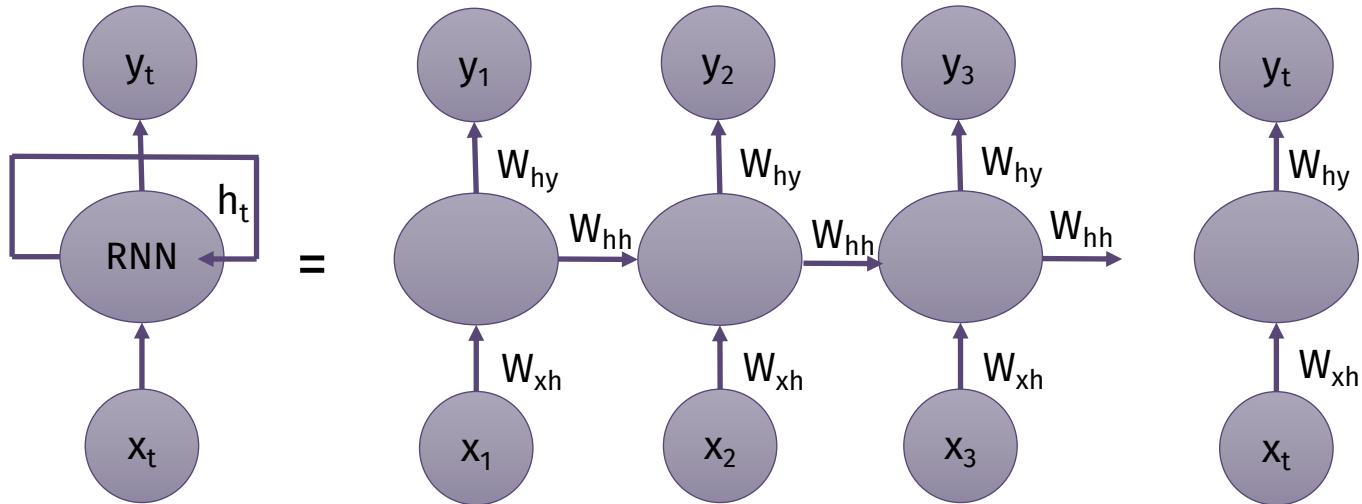
$$y_t = W_{hy} h_t$$

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

Same function and weights are used for each time step, so they are a way to share weights over time

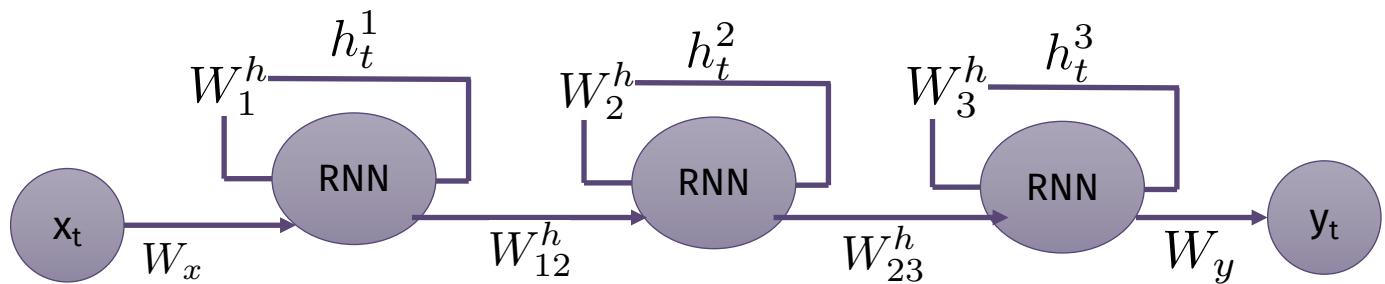
Unrolling RNNs over time

A RNN can be seen as a sequence of multiple, communicating copies of the same network

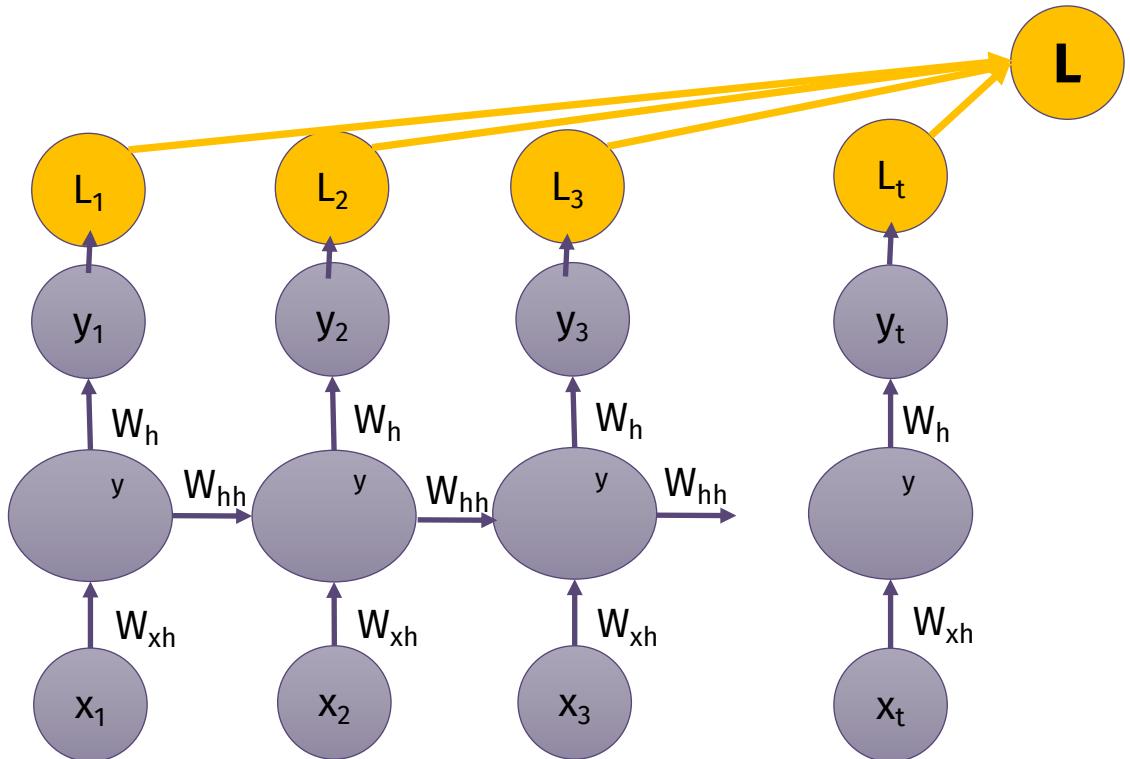


The weight matrices are filters that determine **how much importance to give to both the present input and the past hidden state**

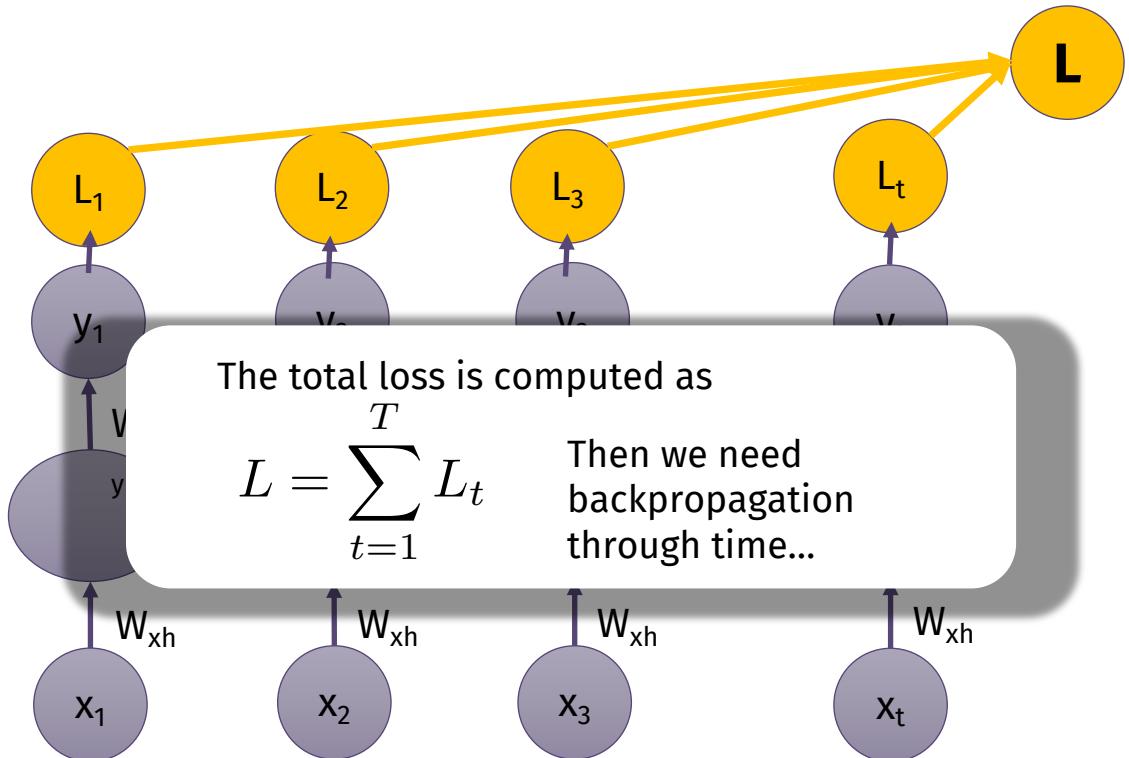
What about having multiple hidden layers?



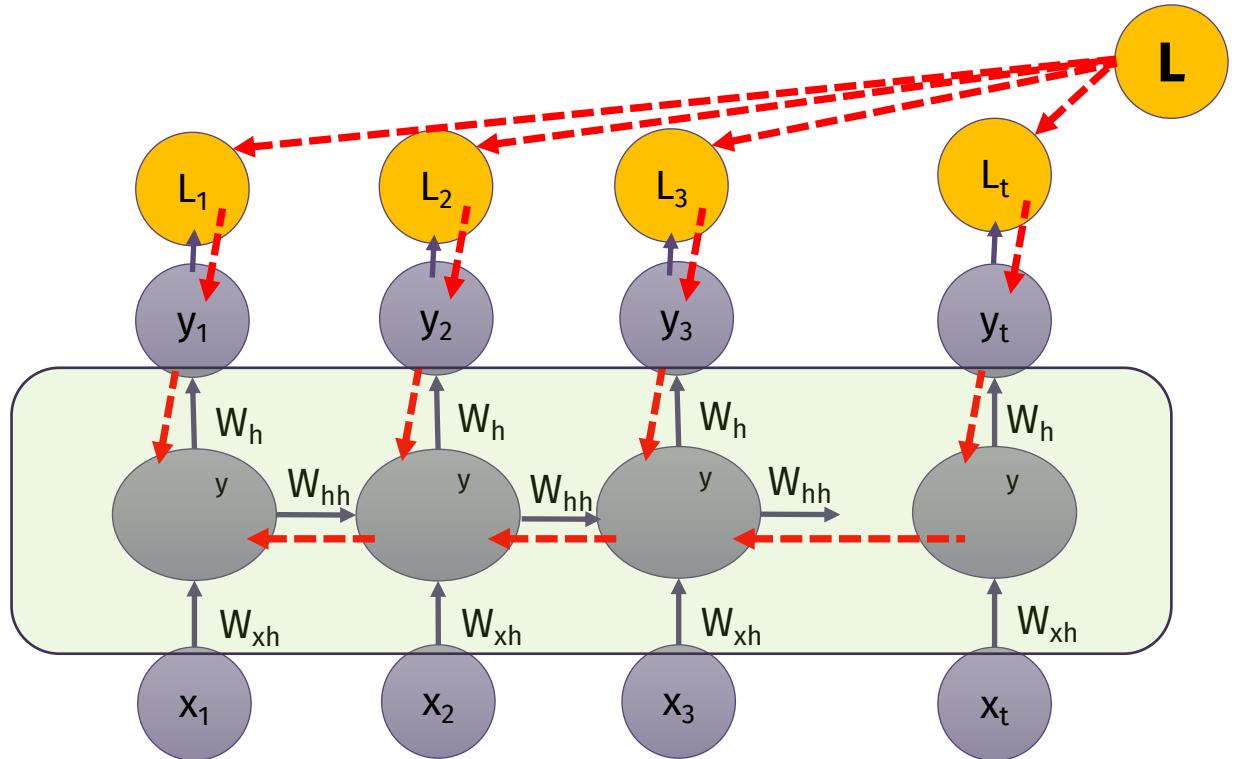
Again on forward propagation



Again on forward propagation



Backward propagation



Loss and backpropagation

$$L = \sum_{t=1}^T L_t \quad \frac{\partial L_t}{\partial W_{hh}} = \frac{\partial L_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \left(\sum_{k=1}^t \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial W_{hh}} \right)$$

$$\frac{\partial h_t}{\partial h_k} = \prod_{i=k+1}^t \frac{\partial h_i}{\partial h_{i-1}}$$

Computed by multiplying consecutive time steps

Gradients-related issues for long-term dependences

- The computation of the loss gradient as successive multiplication leads to instability of the gradient and may take very long training times
- Many values < 1 lead to **vanishing** gradient problems
- Many values > 1 lead to **exploding** gradient problems

Exploding gradient

- Many values > 1 lead to **exploding** gradient problems: **the update with SGD is done with very large steps, leading to bad results**

$$\Theta_{new} = \Theta_{old} - \alpha \nabla_{\Theta} L(\Theta)$$

- A possible solution is gradient clipping: if the gradient is greater than some threshold, scale it down before applying SGD update
- You make a step in the same direction but with a smaller step

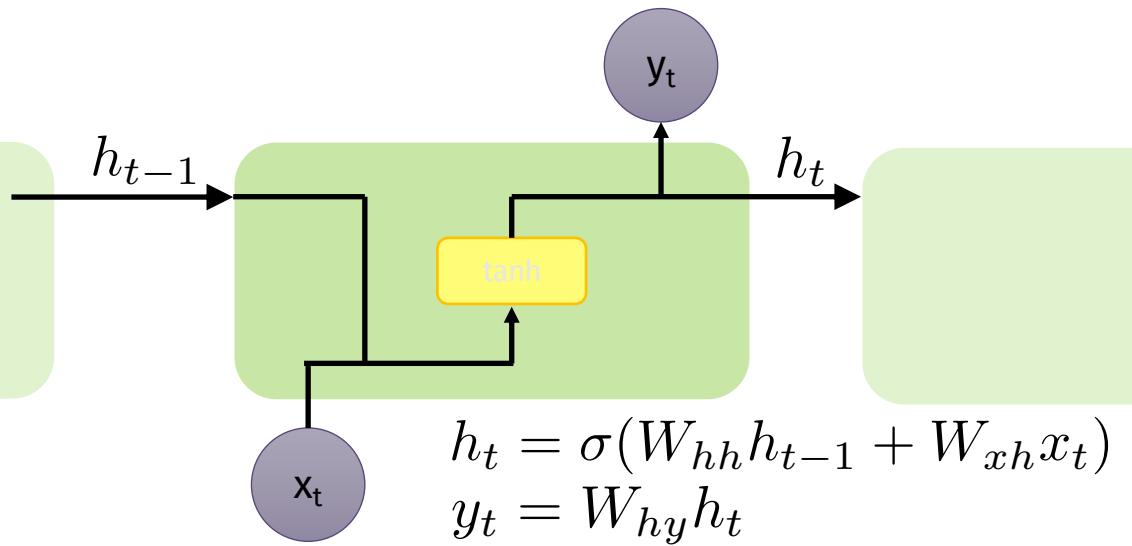
Vanishing gradient

- Many values < 1 lead to **vanishing** gradient problems: **gradient signal far over time is lost because it's much smaller than gradient signal from closer times**
- Model weights are updated only with respect to near effects, not long-term effects
- A possible solution to learn long-term dependences in the data is to use **gated cells**

Long-Short Term Memory

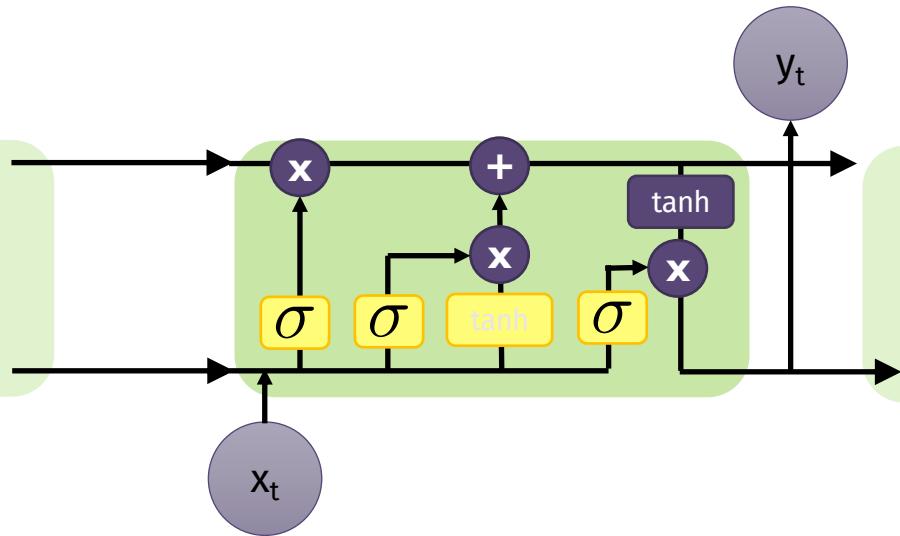
RNNs cells

- In standard RNNs, the cells contain a simple computation and their state is constantly re-written



Gated cells

- Gated cells contain computational blocks that control information flow



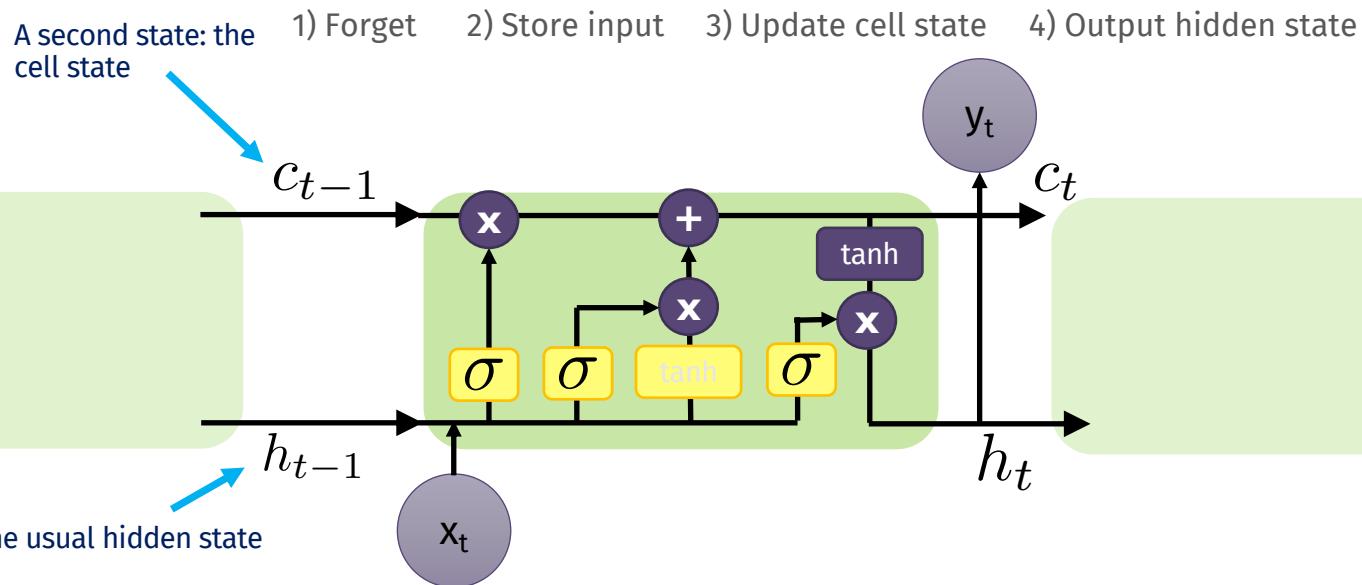
Long-short term memory (LSTM, 1997)

- They consider **connection weights that may change at each time step**
- Information is accumulated over a long duration
- Once the information has been used, it may be useful for the RNN to forget the old state or keep the information
- The LSTM learns how to decide when to do that (this is in fact the role of gated units)

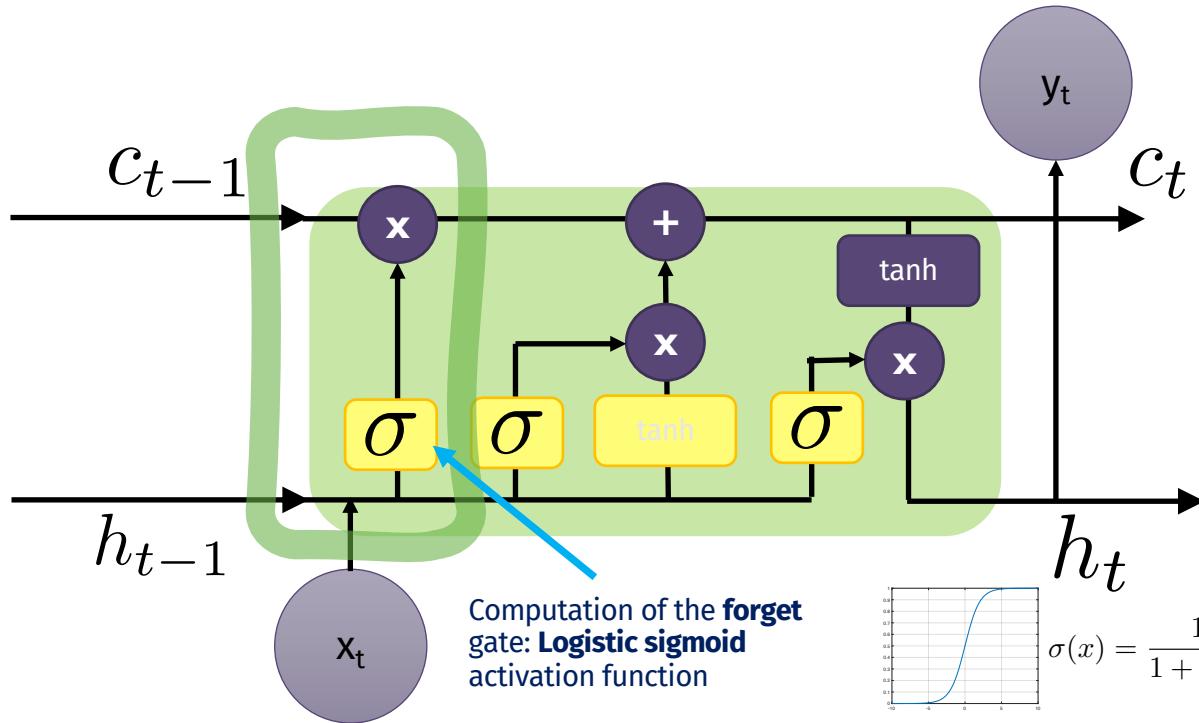
LSTM cell

- It includes two states: a **hidden state** and a **cell state**, both vectors of length n
- The cell stores long-term information. The LSTM can **erase**, **write** and **read** information from the cell
- The selection of which information is erased/written/read is controlled by three corresponding **gates**, vectors again of length n
- Each element of the gates can be open (1), closed (0), or somewhere in-between
- The gates are dynamic: their value is computed based on the current context

LSTM cell: how does it work

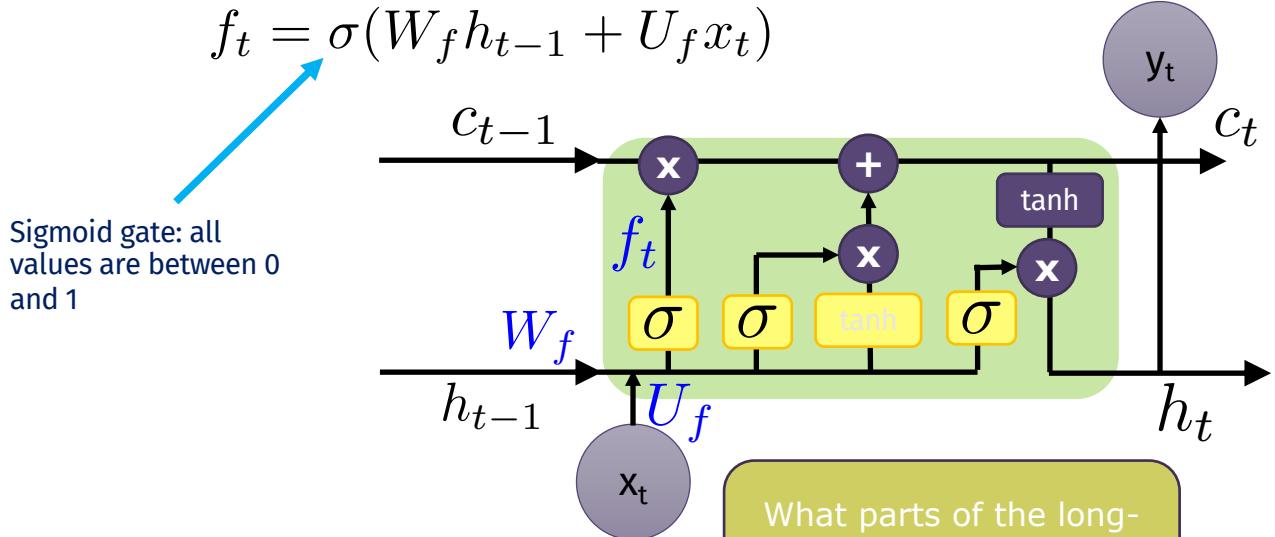


LSTM cell – 1) Forget



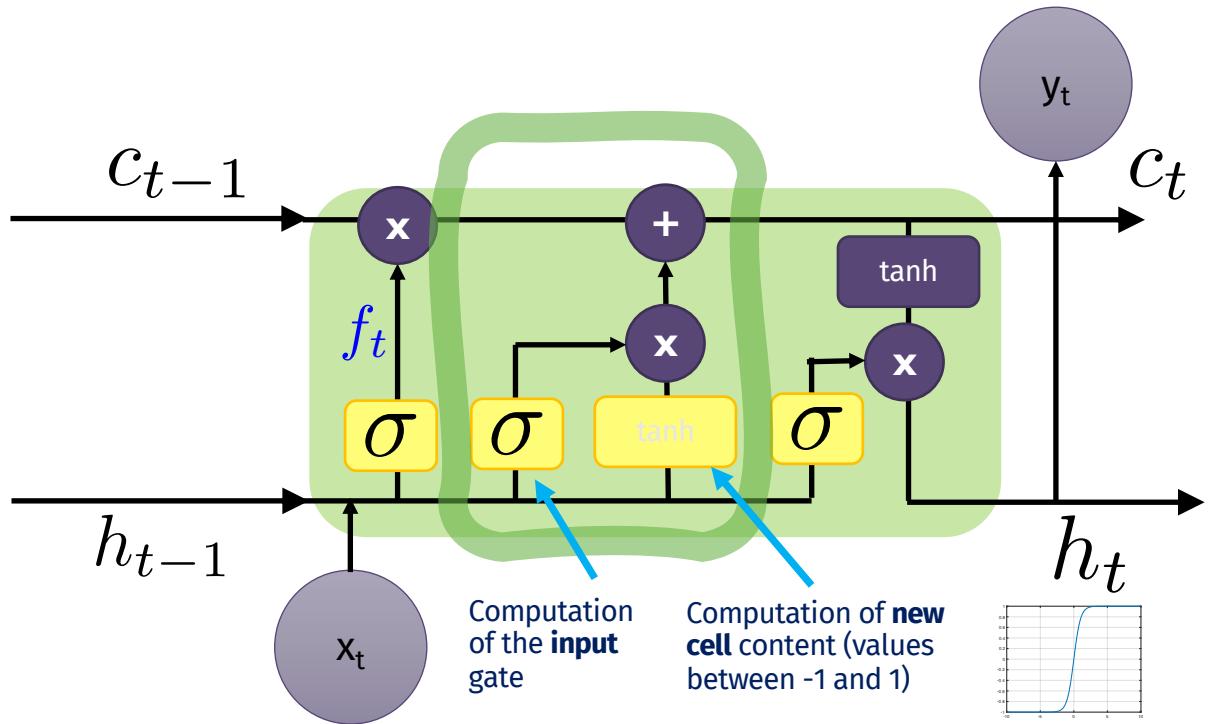
Forget gate

- It decides which information to keep and what to forget from the previous **cell state**



What parts of the long-term memory can now be forgotten?

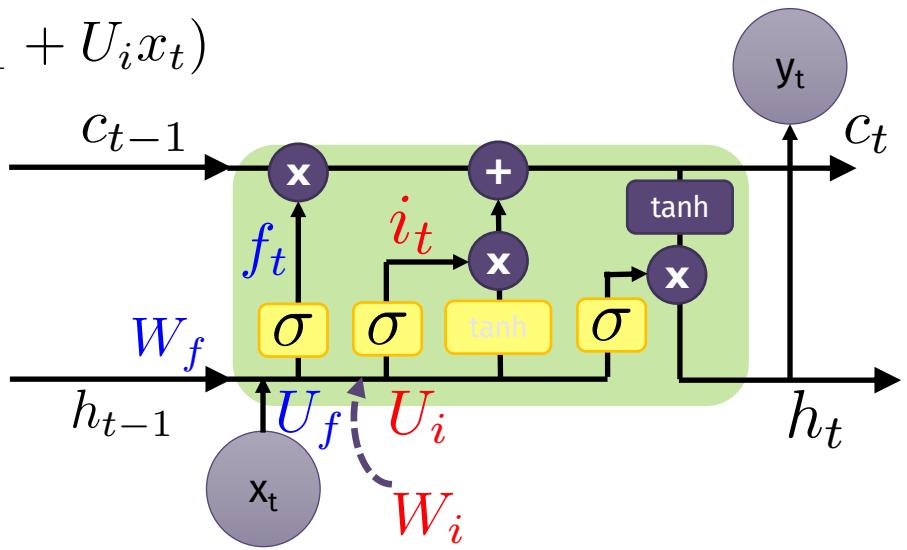
LSTM cell – 2) Store input



Input gate

- It decides what new information (from previous hidden state and new input data) storing in the cell states

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

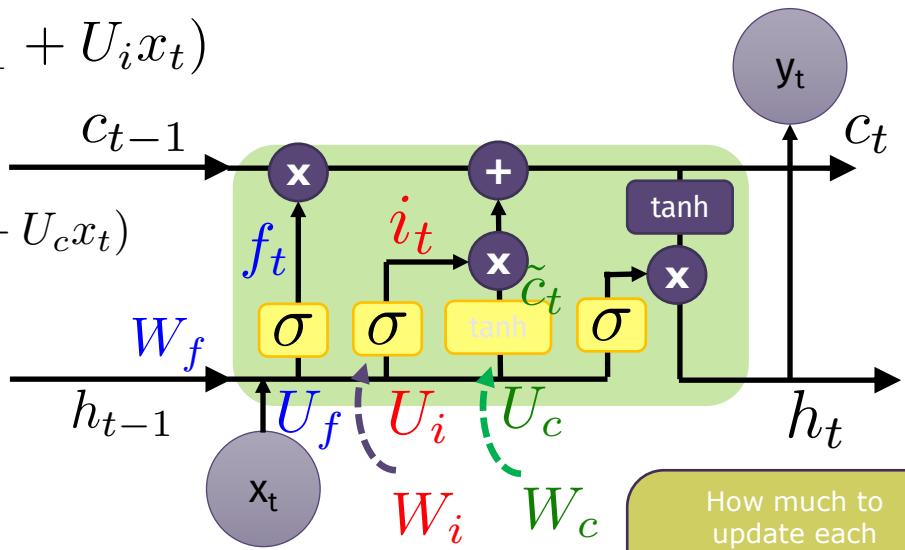


Input gate

- It decides what new information (from previous hidden state and new input data) storing in the cell states

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

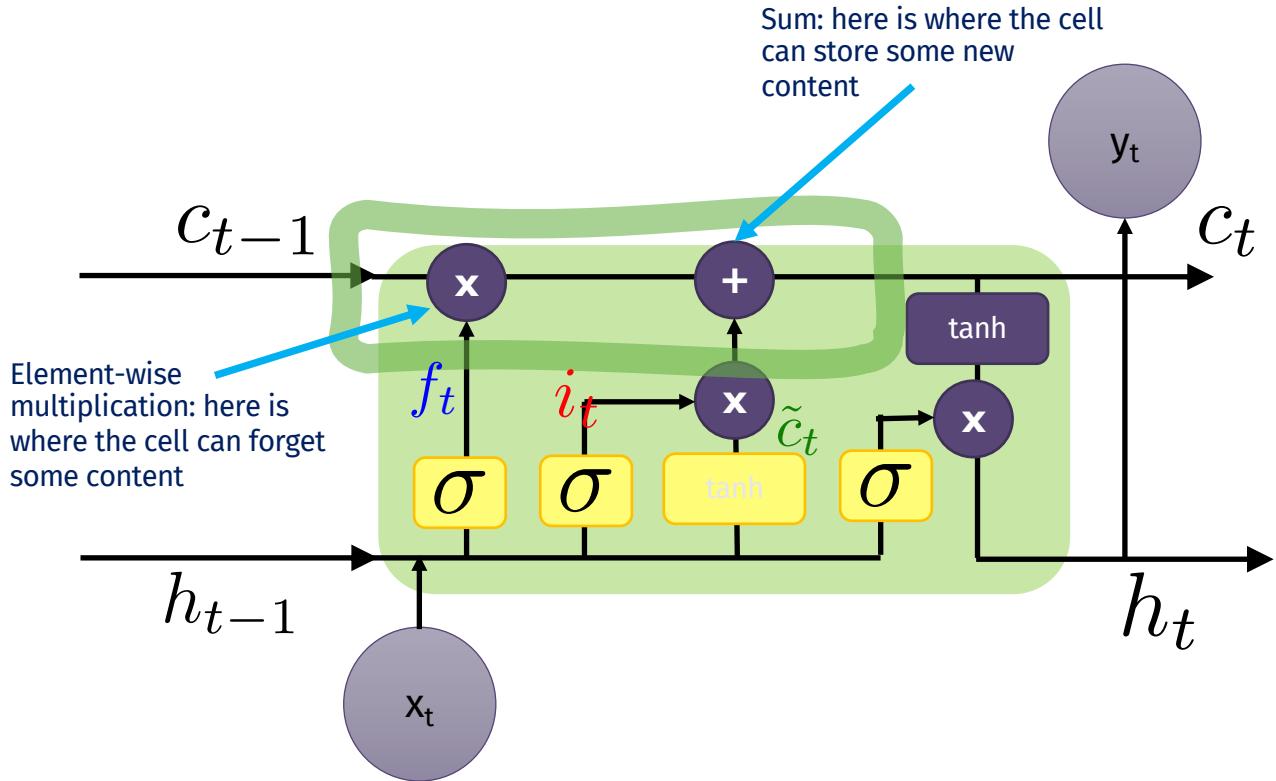
$$\tilde{c}_t = \tanh(W_c h_{t-1} + U_c x_t)$$



Is the new data worth remembering?

How much to update each component of the cell state given the new data and the hidden state

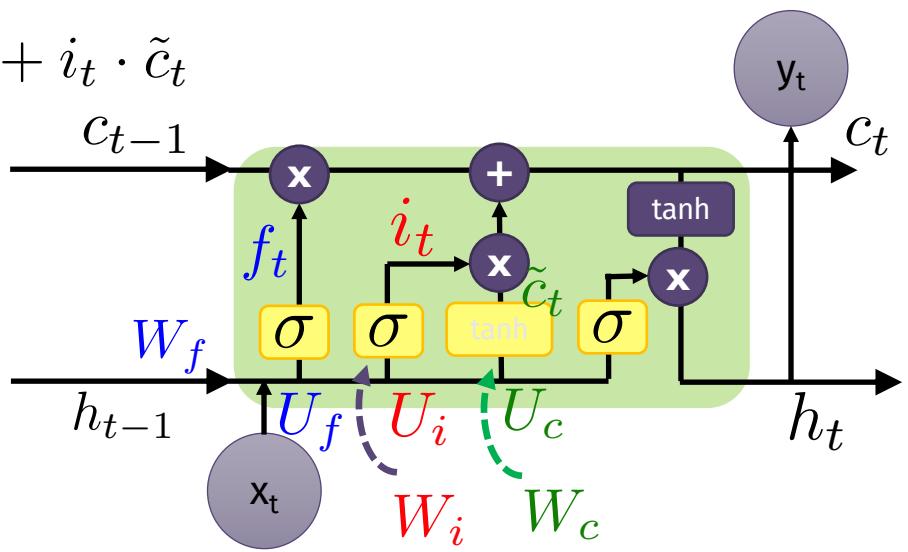
LSTM cell – 3) Update



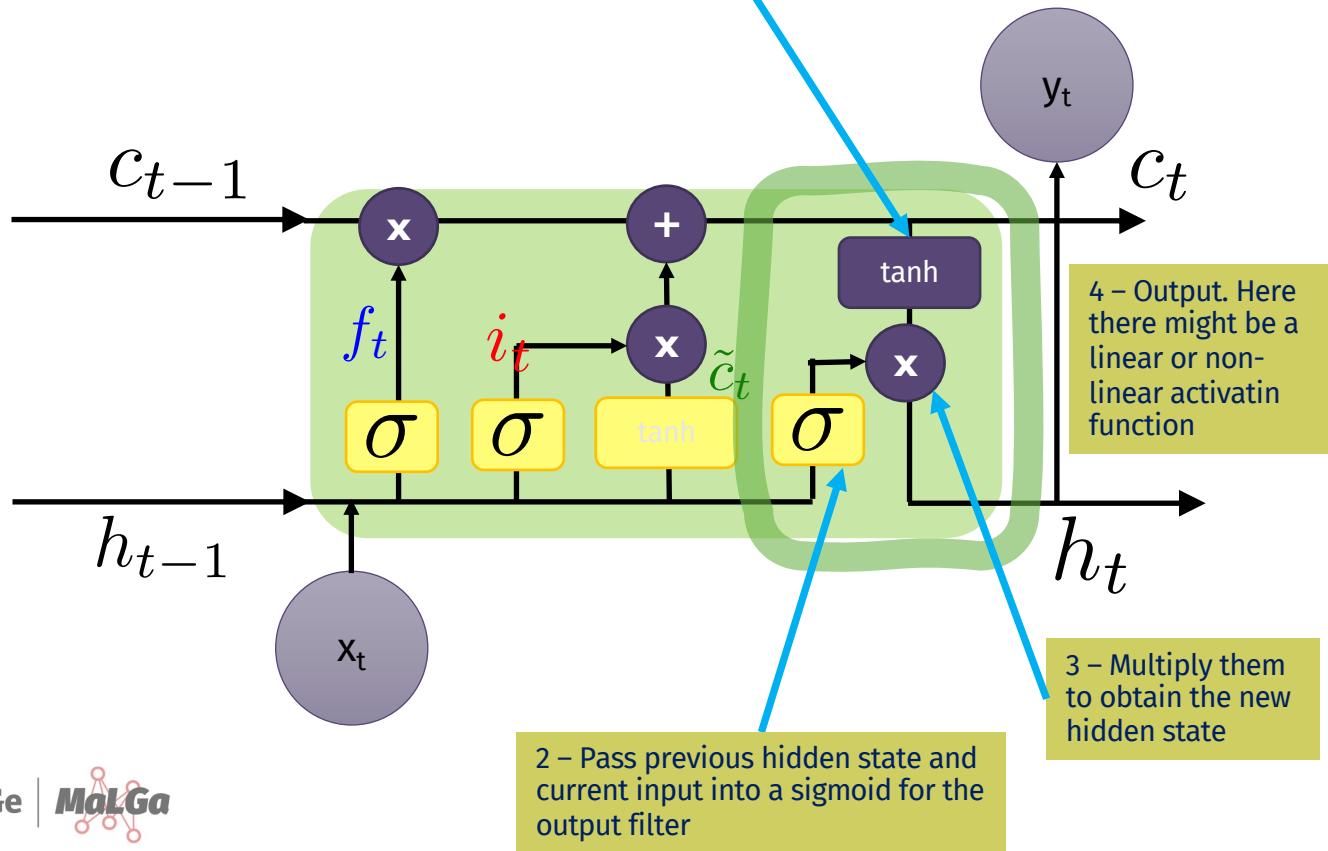
Update

- The cell state values are selectively updated

$$c_t = c_{t-1} \cdot f_t + i_t \cdot \tilde{c}_t$$



LSTM cell – 4) Output



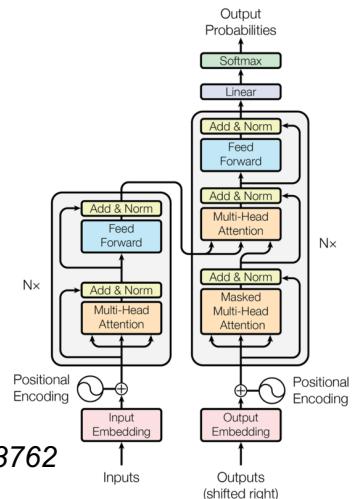
Variations on the theme and recent advances

GRU (Gated Recurrent Unit)

- It includes only the hidden state
- It also has two gates:
 - Update gate: it decides what information to keep and what to throw away
 - Reset gate: it decides how much past information to forget

Transformers

- It includes an attention mechanism that decides at each step which other part of a certain sequence is important

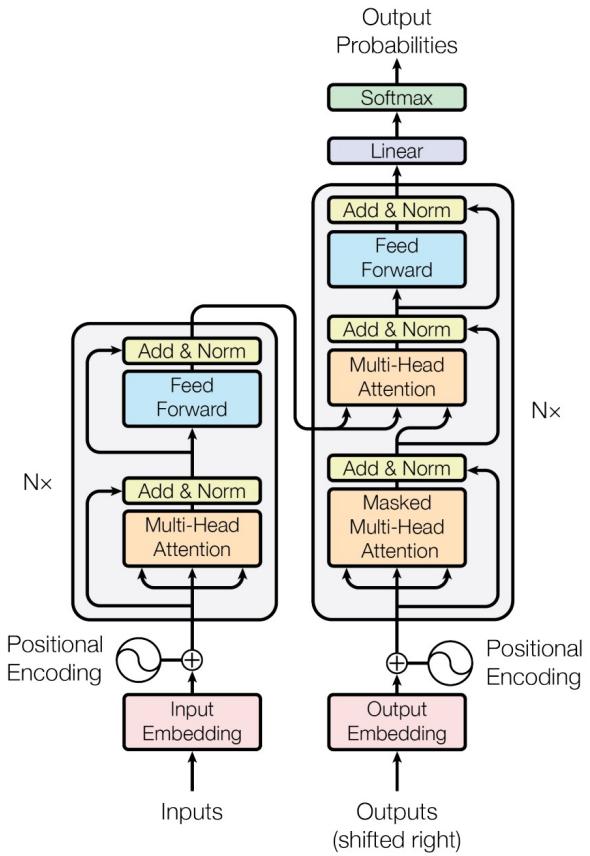


<https://arxiv.org/abs/1706.03762>

Self-attention and Transformers

Transformer (2017)

Originally proposed for language translation, they can be highly parallelized



From <https://arxiv.org/abs/1706.03762>

Figure 1: The Transformer - model architecture.

Transformer

- As other models for sequence transduction (e.g. language translation, they are based on an encoder-decoder structure)

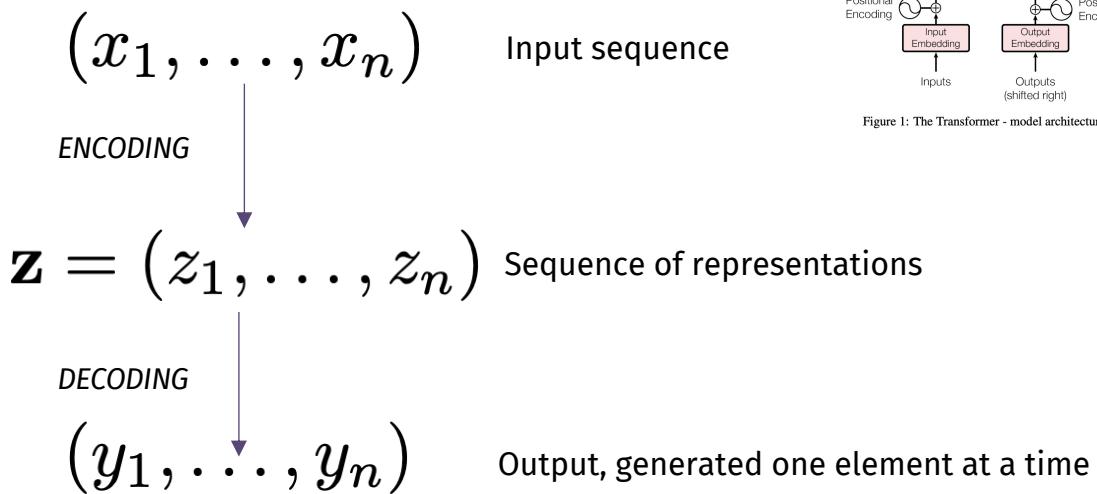


Figure 1: The Transformer - model architecture.

Transformer: main structure

A stack of N identical layers each one composed by multi-head self-attention and a fully connected net

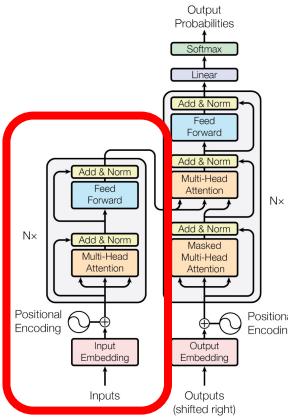
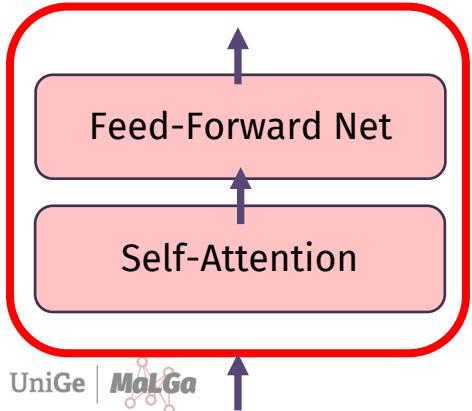
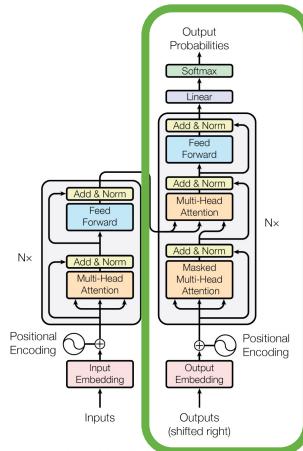


Figure 1: The Transformer - model architecture.

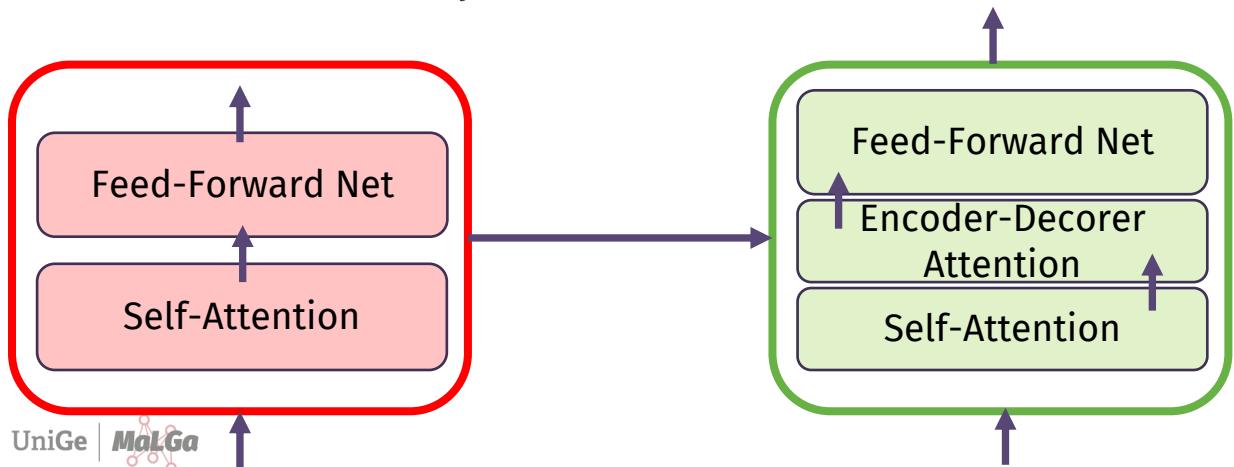


Transformer: main structure

A stack of N identical layers each one composed by multi-head self-attention and a fully connected net



A stack of N identical layers each one composed by masked multi-head self-attention, multi-head self-attention, and a fully connected net



Transformer: main structure

A stack of N identical layers each one composed by multi-head self-attention and a fully connected net

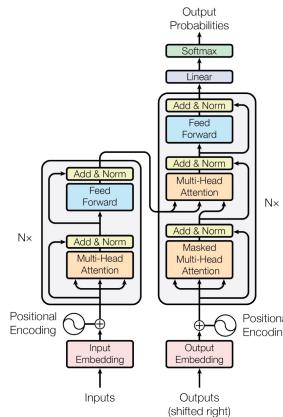


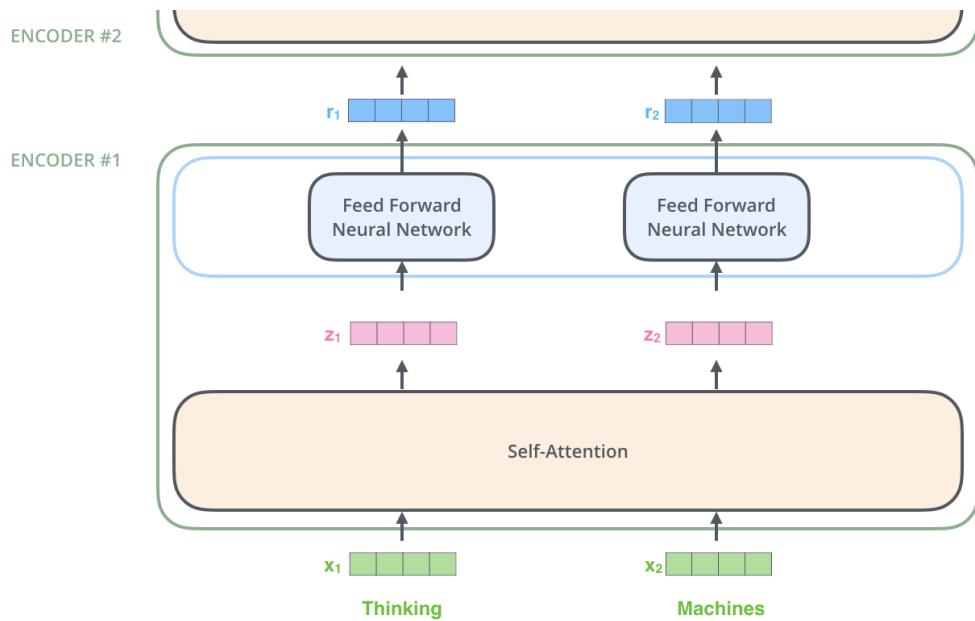
Figure 1: The Transformer - model architecture.

A stack of N identical layers each one composed by masked multi-head self-attention, multi-head self-attention, and a fully connected net

RESIDUAL CONNECTION: the output of each layer is

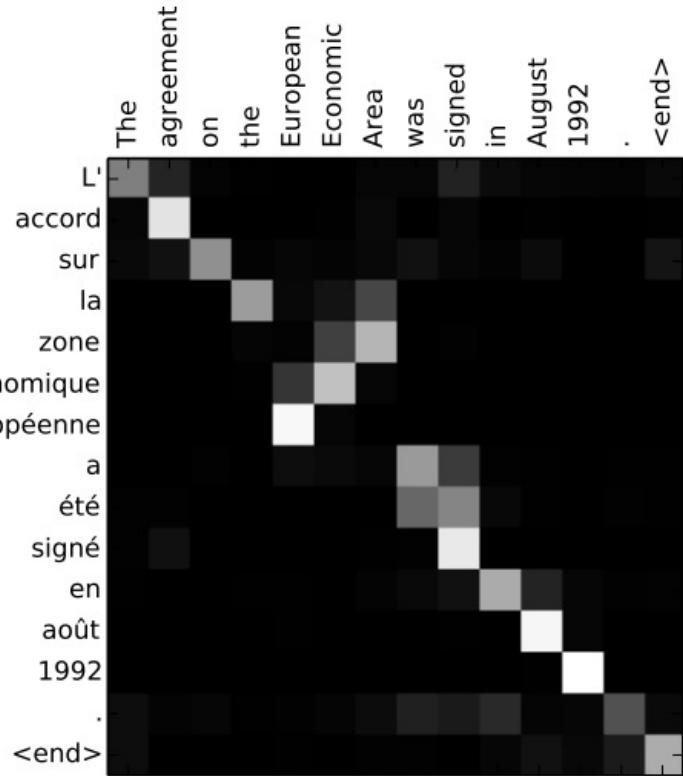
$$\sigma(x + f(x))$$

Where x is the input to the layer, while $f(x)$ is the function implemented by the layer itself



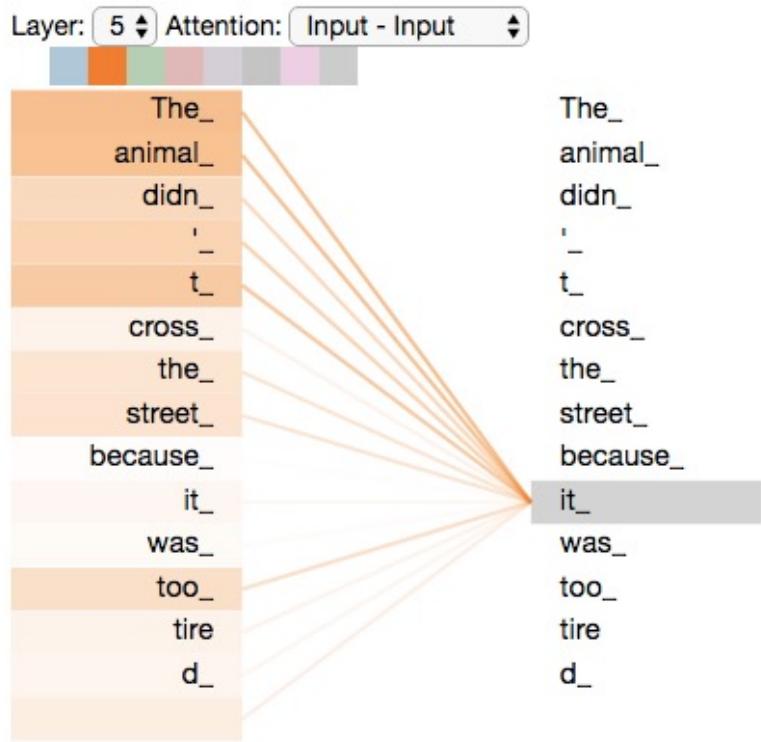
From <http://jalammar.github.io/illustrated-transformer/>

Why attention?



From «*Neural Machine Translation by Jointly Learning to Align and Translate*»
(2015)

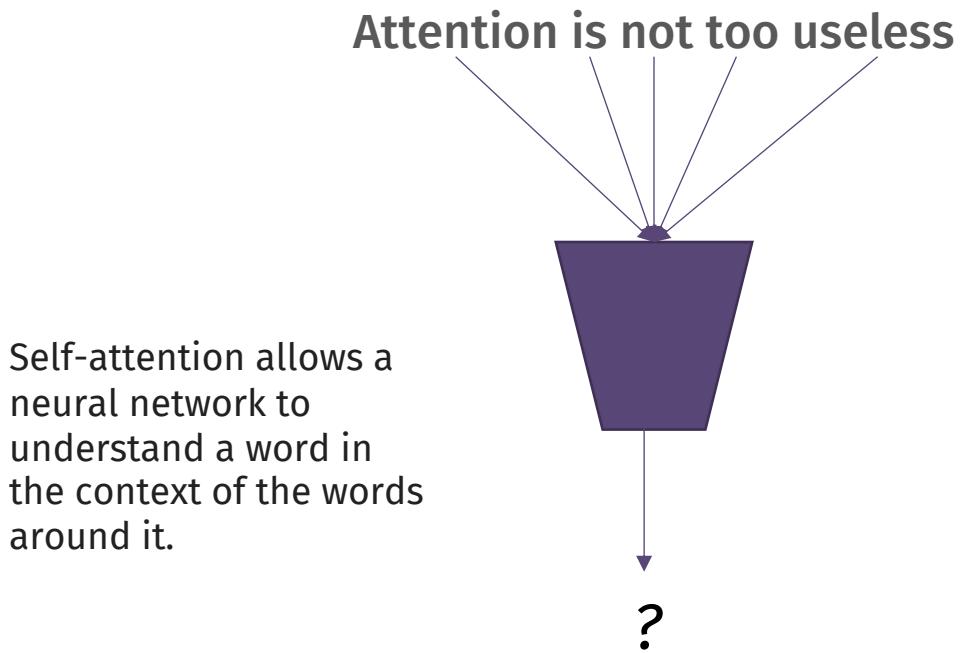
Why self-attention?



<http://jalammar.github.io/illustrated-transformer/>

What is self-attention?

An example on sentiment classification



What is self-attention?

An example on sentiment classification

Attention is **not** too **useless**



*To correctly classify
the sentence we
need to consider all
the words in it*

What is self-attention?

An example on sentiment classification

Attention is **not** too **useless**



*To correctly classify
the sentence we
need to consider all
the words in it*

*NOT ONLY: we also
need to understand
the relations
between them*

Let's proceed step by step

Input

Thinking

Machines

Embedding

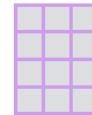
x_1

x_2

Queries

q_1

q_2



W^Q

Keys

k_1

k_2



W^K

Values

v_1

v_2

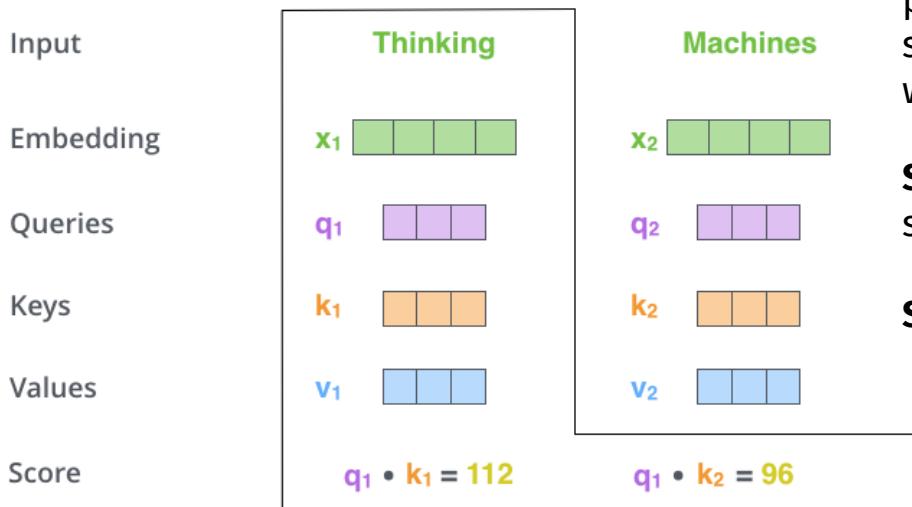


W^V

From <http://jalammar.github.io/illustrated-transformer/>

Step 1 - For each word, we create a Query vector, a Key vector, and a Value vector

Let's proceed step by step



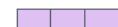
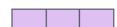
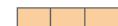
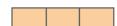
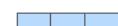
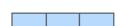
Step 2 - For each embedding, we compute the score that determines the importance of other parts of the input sentence as we encode a word at a certain position.

Step 3 – Normalize the scores

Step 4 – Apply softmax

From <http://jalammar.github.io/illustrated-transformer/>

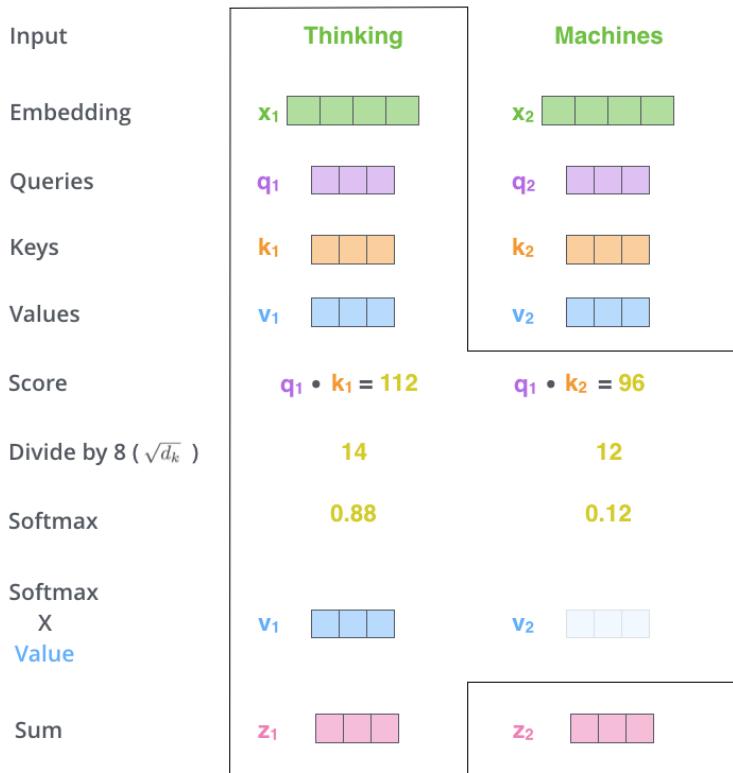
Let's proceed step by step

Input	Thinking		Machines	
Embedding	x_1		x_2	
Queries	q_1		q_2	
Keys	k_1		k_2	
Values	v_1		v_2	
Score	$q_1 \cdot k_1 = 112$		$q_1 \cdot k_2 = 96$	
Divide by 8 ($\sqrt{d_k}$)	14		12	
Softmax	0.88		0.12	

The softmax score determines the “level of expression” of the word at this position.

From <http://jalammar.github.io/illustrated-transformer/>

Let's proceed step by step



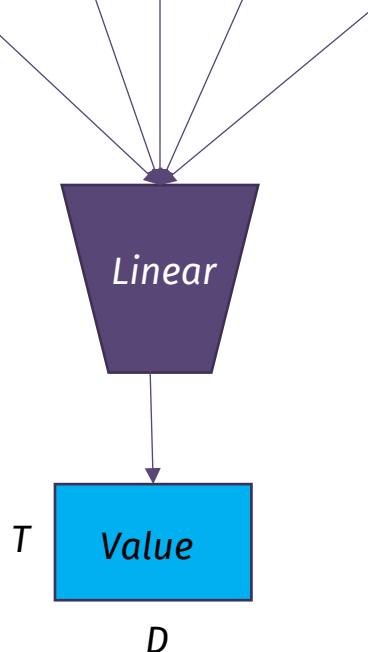
Step 5 – Multiply each value vector by the softmax score. This keeps intact the values of the word(s) we want to focus on, and attenuate the contribution of irrelevant words

Step 6 - Sum the weighted value vectors. This is the output of the self-attention layer at this position (for the first word).

From
<http://jalammar.github.io/illustrated-transformer/>

What is self-attention? An example on sentiment classification

Attention is **not** too **useless**



*We may encode
relationships between
values by summation*

Combining values/words

	Attention	is	not	too	useless
Attention	1	1	1	1	1
is	1	1	1	1	1
not	1	1	1	1	1
too	1	1	1	1	1
useless	1	1	1	1	1

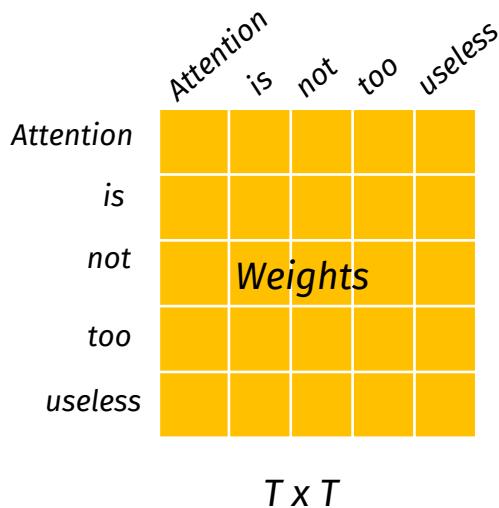
$T \times T$



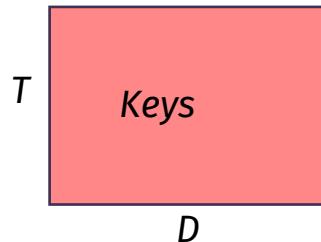
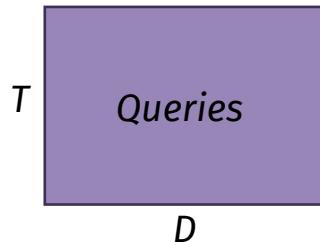
Here we would assume the same importance for all relationships... but we would like that, for instance, the relation between «not» and «useless» is more important than the one between «is» and «too»

Learning the attention weights

Inspired by
<https://twitter.com/MishaLaskin/status/1479246928454037508>

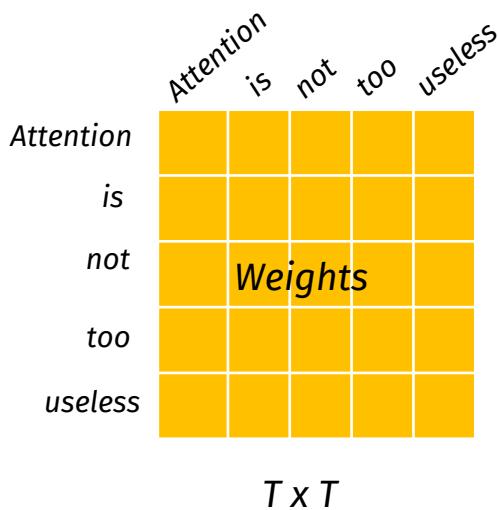


$$W = Q \cdot K^T$$



Learning the attention weights

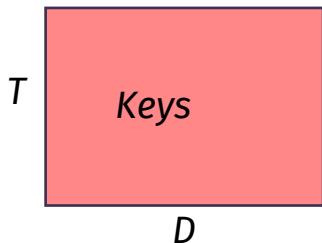
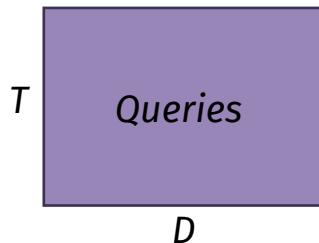
Inspired by
<https://twitter.com/MishaLaskin/status/1479246928454037508>



Intuition: we want the W matrix to weight the relationship between word_i as a context for word_j. We employ other two linear nets

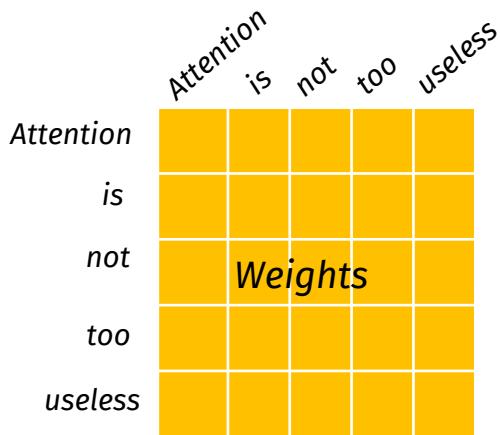


$$W = Q \cdot K^T \rightarrow \frac{Q \cdot K^T}{\sqrt{D}}$$

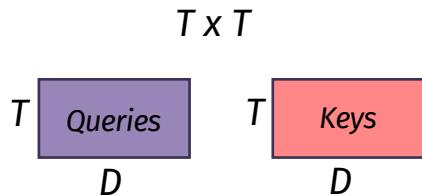
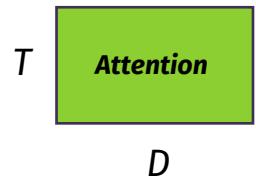
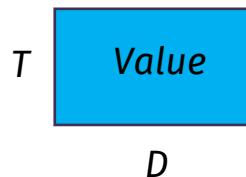


Single-head attention

W sees the relationships between words



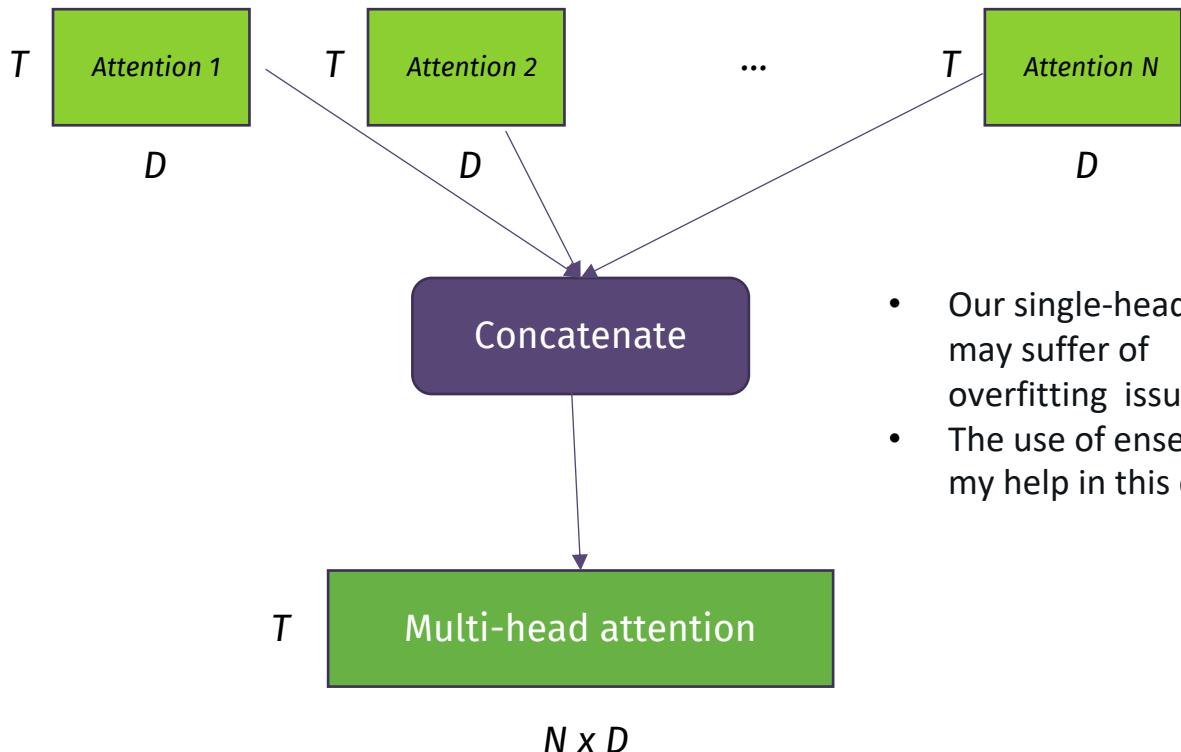
V sees all the words in the sentence



$$\text{softmax}\left(\frac{QK^T}{\sqrt{D}}\right)V$$

This is the (single head) self-attention

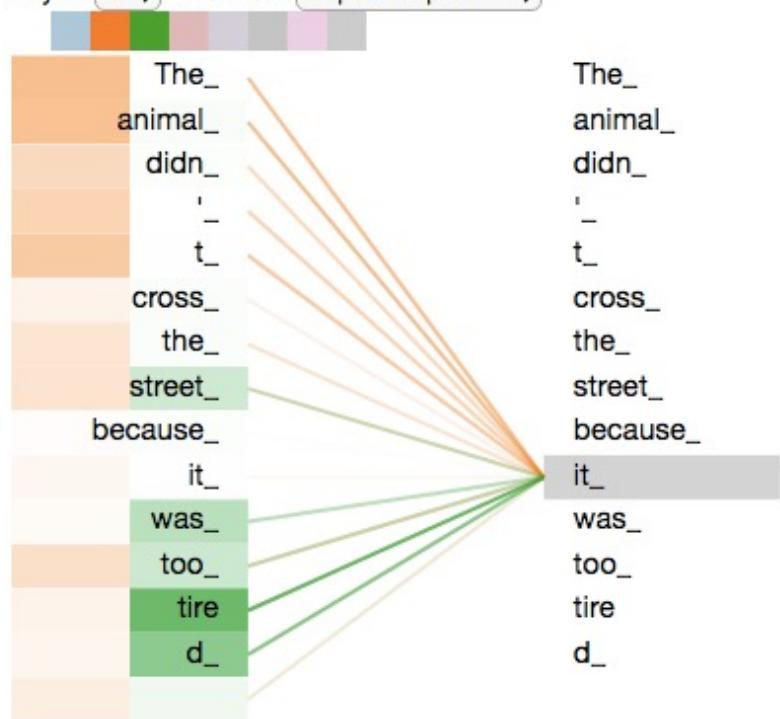
Multi-head attention



- Our single-head net may suffer of overfitting issues
- The use of ensembles my help in this case

Back to the example

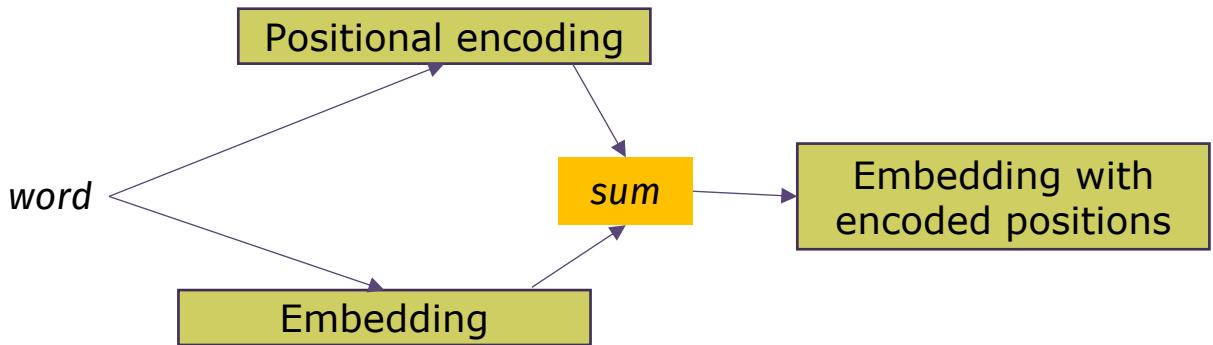
Layer: 5 Attention: Input - Input



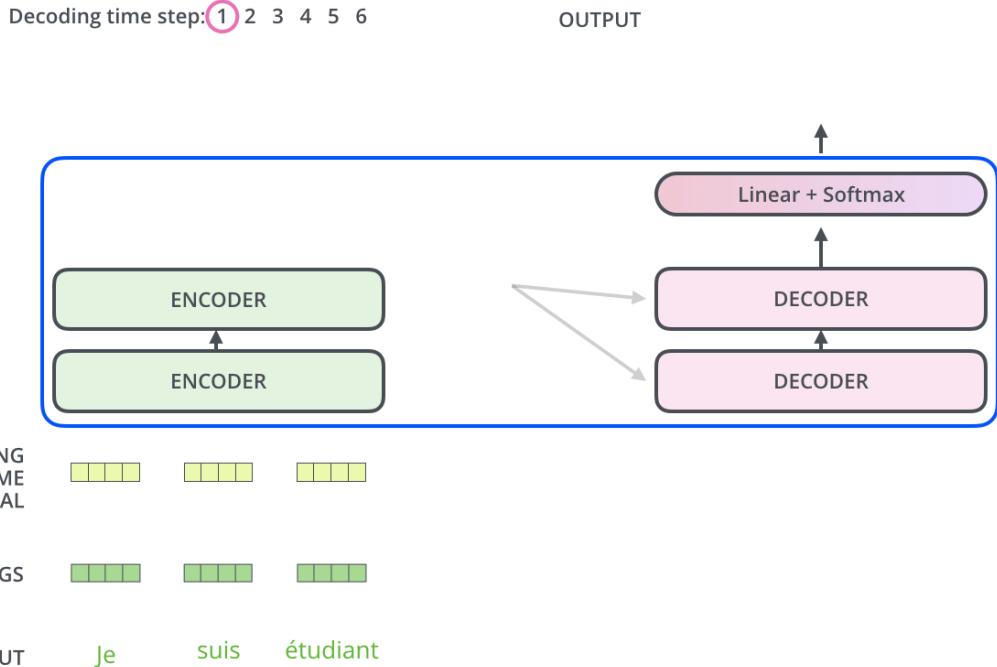
<http://jalammar.github.io/illustrated-transformer/>

What about the position of the words in the sequence?

- We need to «guide» the network to see the words in the correct order
- This is done by means of positional encoding



The decoder side

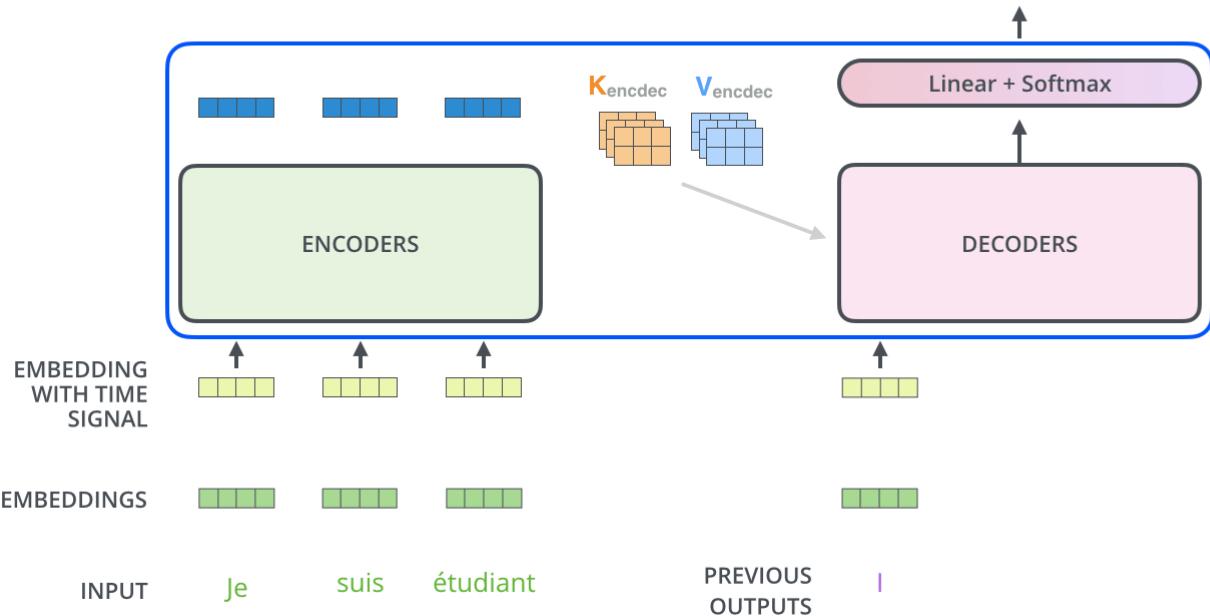


From
<http://jalammar.github.io/illustrated-transformer/>

The decoder side

Decoding time step: 1 2 3 4 5 6

OUTPUT |



From
<http://jalammar.github.io/illustrated-transformer/>

The decoder side: masked attention

- In the decoder, the self-attention layer is only allowed to consider earlier positions in the output sequence (it can not “see” the future)
- This is achieved using masked attention: future positions are set to $-\infty$ before the softmax step

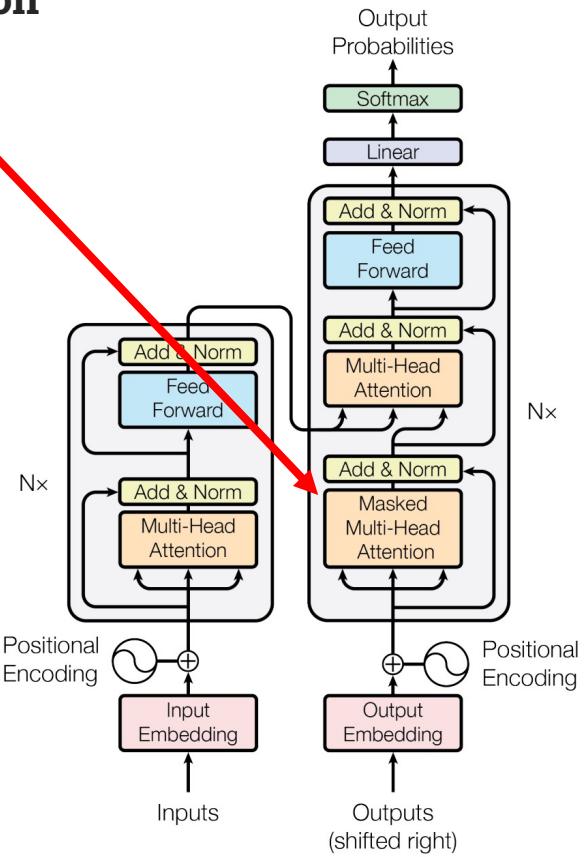


Figure 1: The Transformer - model architecture.

UniGe

