



# Master in Computer Vision Barcelona

---

[<http://pagines.uab.cat/mcv/>]

Xavier Giro-i-Nieto



[@DocXavi](https://twitter.com/DocXavi)



[xavier.giro@upc.edu](mailto:xavier.giro@upc.edu)

Associate Professor  
Universitat Politècnica de Catalunya  
Institut de Robòtica i Informàtica Industrial



## Module 6 - Day 6 - Lecture 1 Recurrent Neural Networks RNN

24th March 2022

# Acknowledgments



Santiago Pascual



Marta R. Costa-jussà

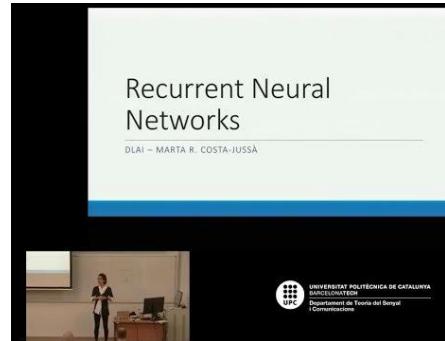


UNIVERSITAT POLITÈCNICA DE CATALUNYA  
BARCELONATECH  
Language and Speech Technologies



# Video-Lectures

Santiago Pascual  
UPC TelecomBCN  
DLSL 2017



Marta R. Costa-Jussà  
UPC TelecomBCN  
DLAI 2017



Xavier Giró  
UPC TelecomBCN  
DLAI 2020

# Outline

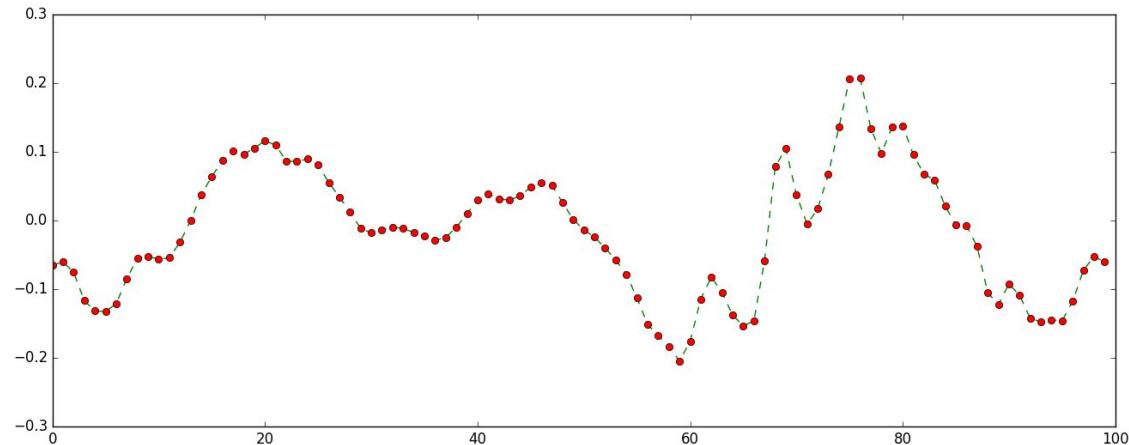
## 1. Motivation



Source: ([AI Memes](#)) AI & Deep Learning Memes For Back-propagated Poets

# Sequences (a.k.a time series)

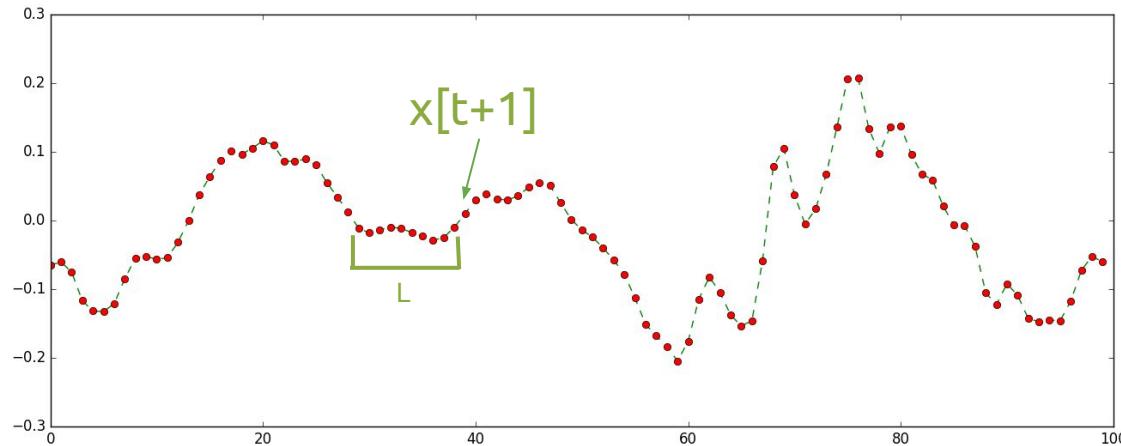
Consider a sequence of samples (eg. of temperature):



# Sequences with Feed Forward MLP

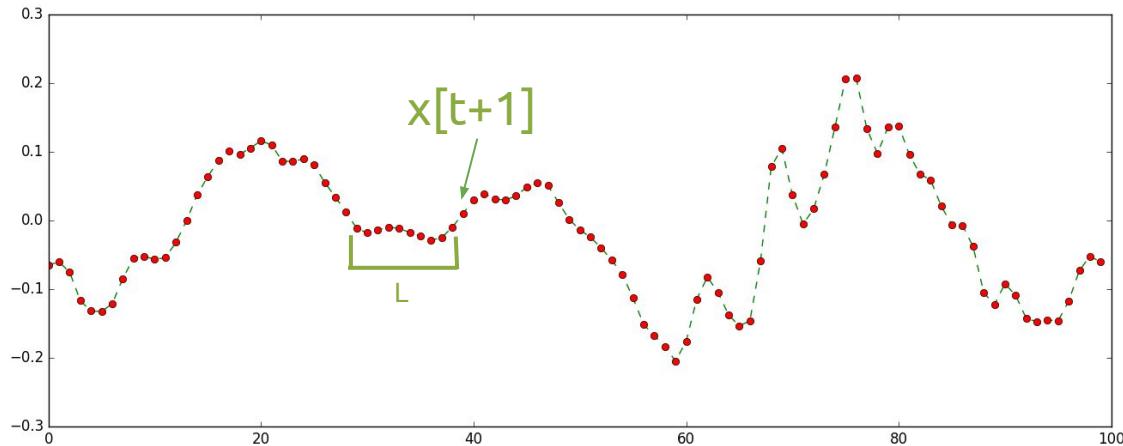
## Task:

Predict sample  $x[t+1]$  knowing previous values  $\{x[t], x[t-1], x[t-2], \dots, x[t-L]\}$



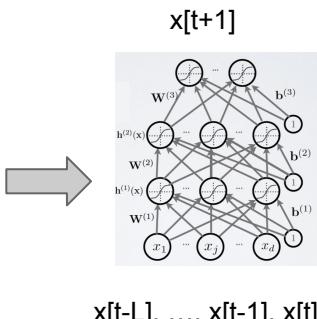
How could you solve this task with a naive feed forward MLP ?

# Sequences with Feed Forward MLP

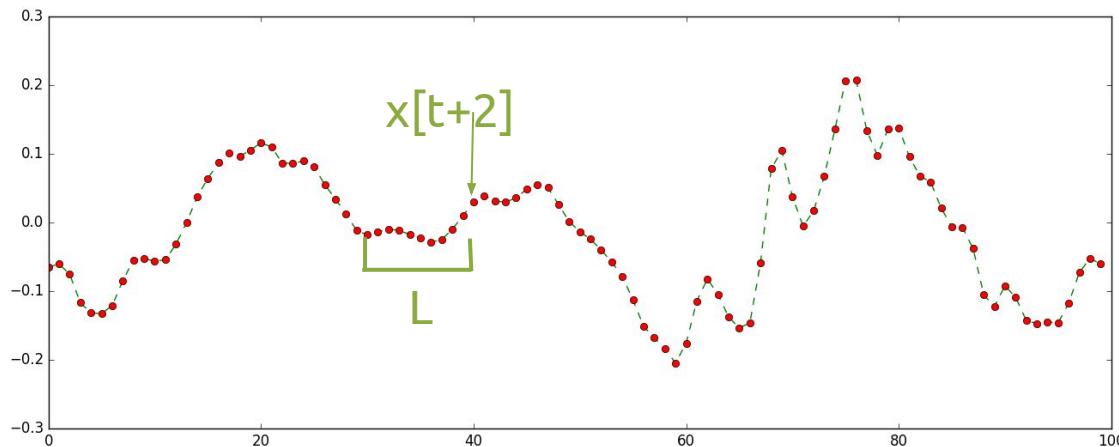


Feed Forward approach:

- static window of size  $L$
- slide the window time-step wise

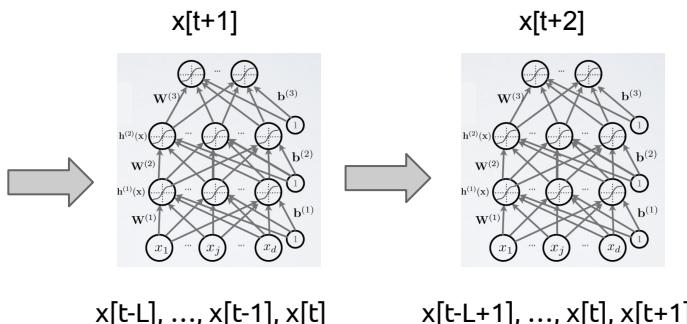


# Sequences with Feed Forward MLP

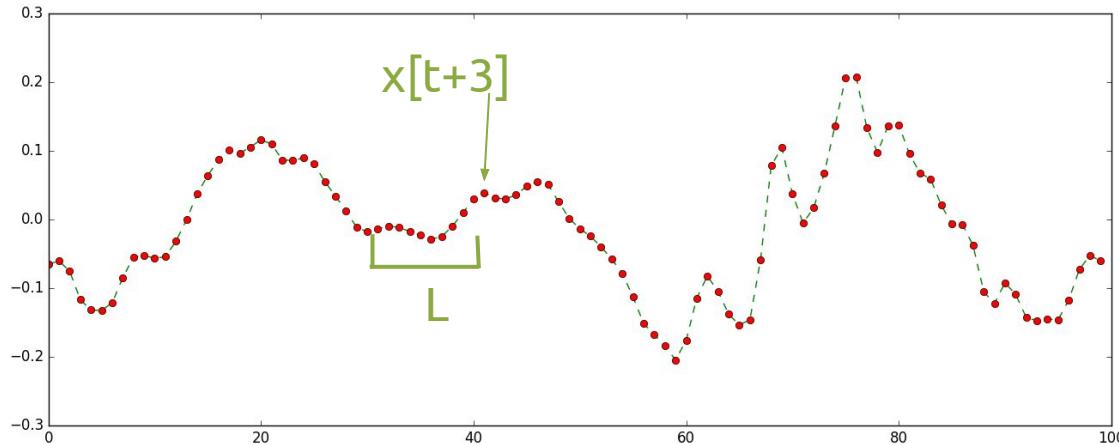


Feed Forward approach:

- static window of size L
- slide the window time-step wise

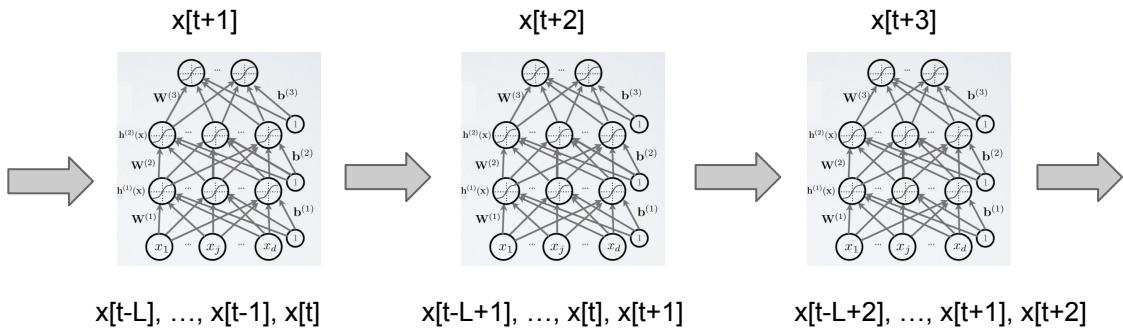


# Sequences with Feed Forward MLP



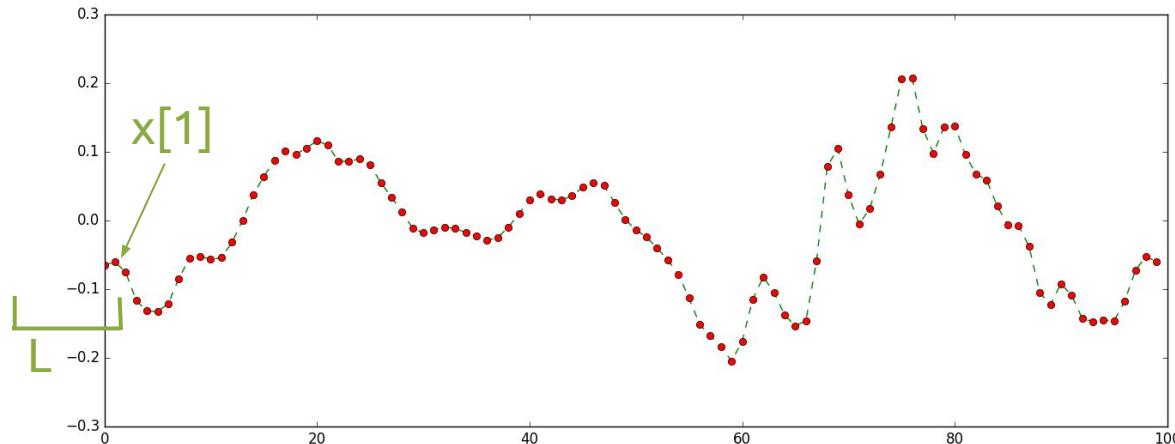
Feed Forward approach:

- static window of size  $L$
- slide the window time-step wise



# Sequences with Feed Forward MLP

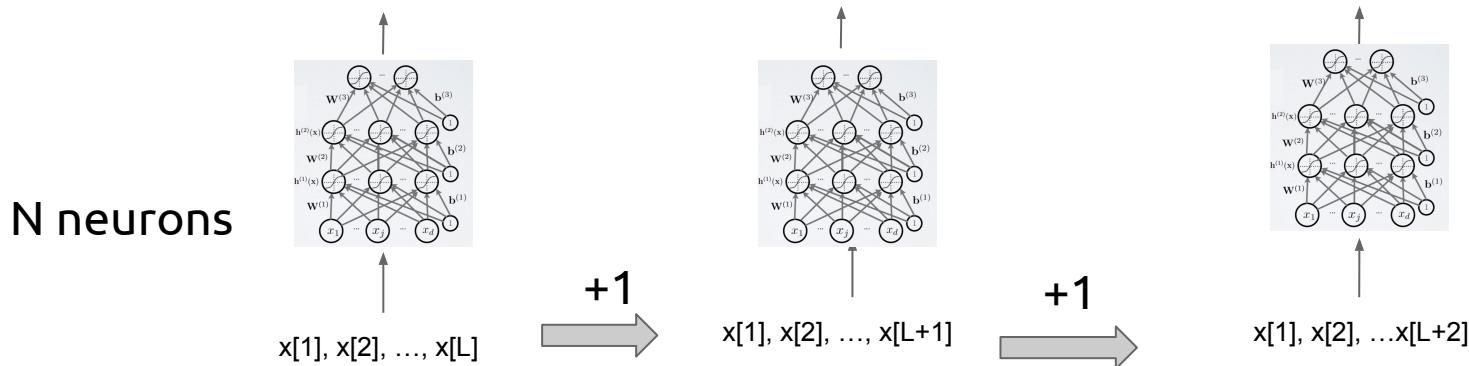
How would you deal with the prediction at the beginning of the sequence ?



By adding (zero) padding for the samples necessary to fill the input layer of the MLP.

# Sequences with Feed Forward MLP

In the first layer, how does an increase of the receptive field ( $L$ ) affect the amount of parameters to learn ?



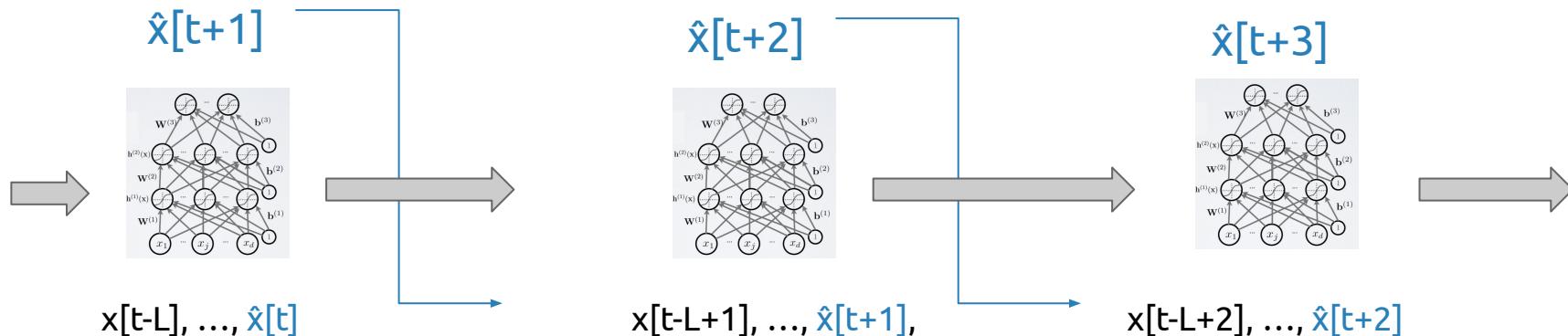
Receptive field	$L$	$L+1$	$L+2$
#params in 1st layer	$N \times L + N$	$N \times (L+1) + N$	$N \times (L+2) + N$

A +1 increase of the receptive field increases the parameters with a factor  $\times N$ .

# Sequences with Feed Forward MLP

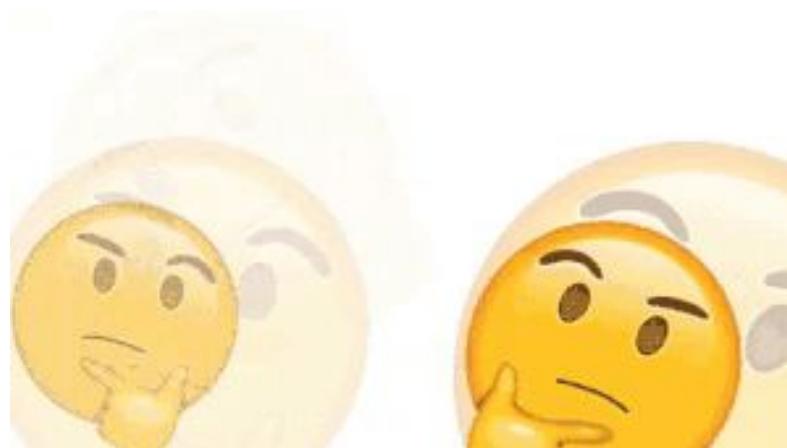
How could the previous decisions condition the future predictions ?

Adopting an **auto-regressive** strategy, by which the previous outputs become inputs in future steps.



Could we find an architecture that:

- adapts to any input length to avoid padding.
- does not increase the complexity with the receptive field.
- can learn what to remember from the past ?



# Outline

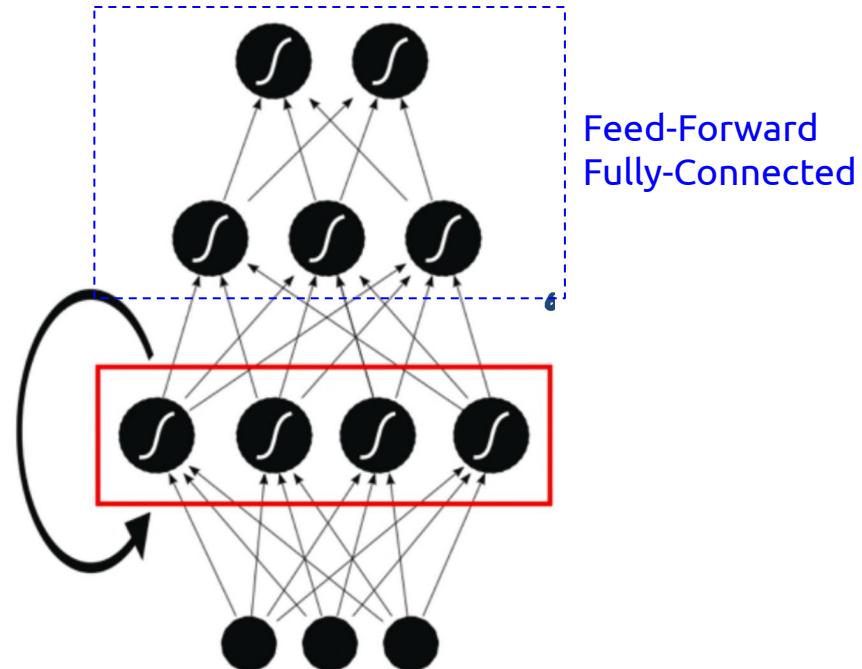
1. Motivation
2. Recurrent Neural Layer

# Recurrent Neural Layer

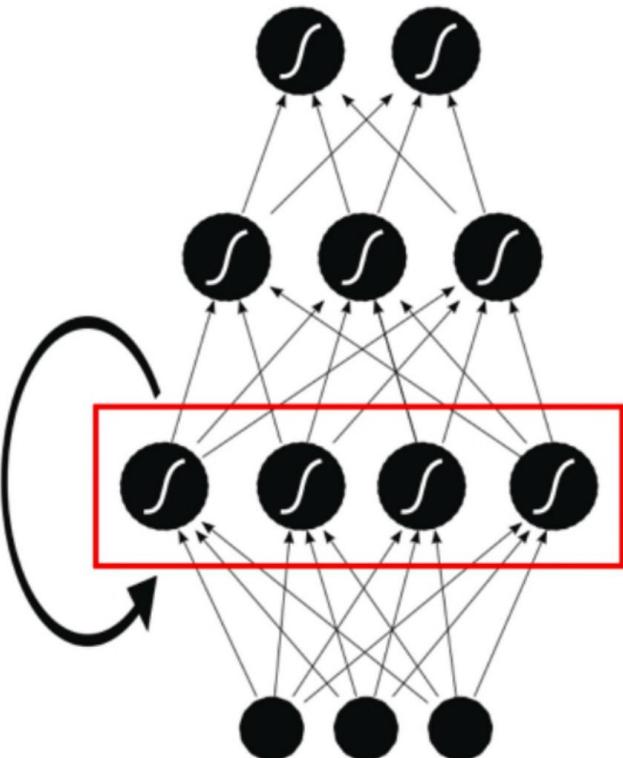


Recurrent layer (RNN)

The hidden state of a recurrent layer  $h_t$  also depends from its previous state  $h_{t-1}$ .



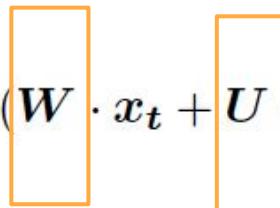
# Recurrent Neural Layer



Feed-forward  
Weights (W)

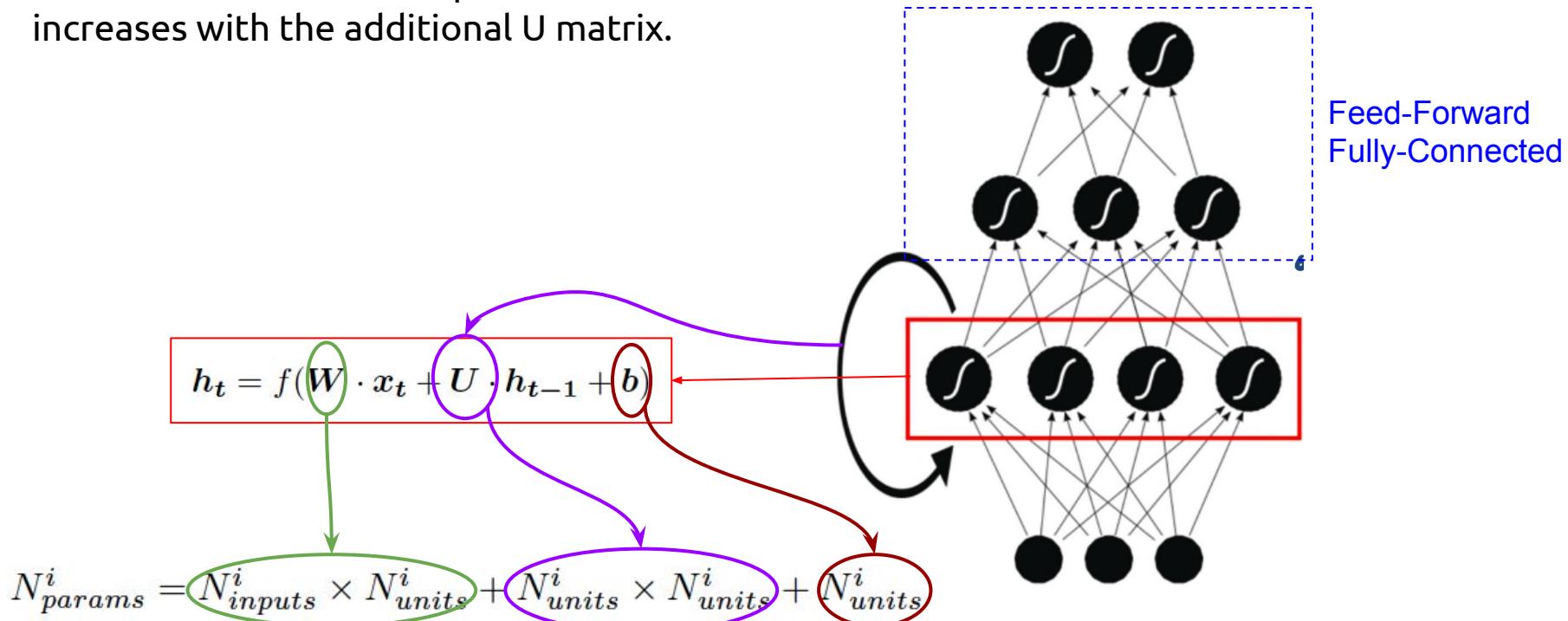
$$h_t = f(W \cdot x_t + U \cdot h_{t-1} + b)$$

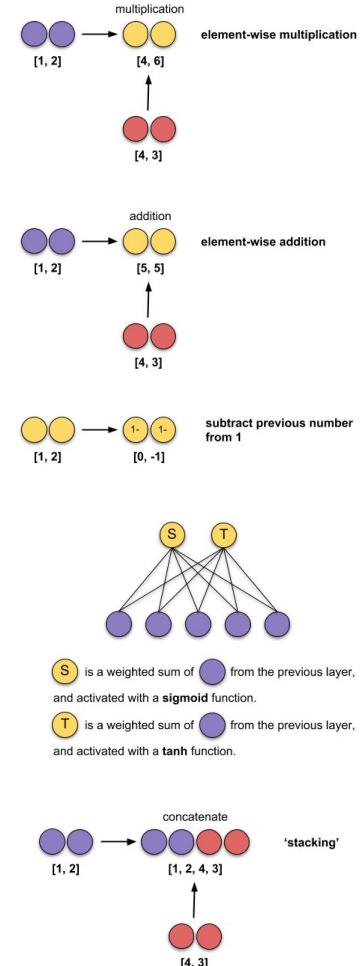
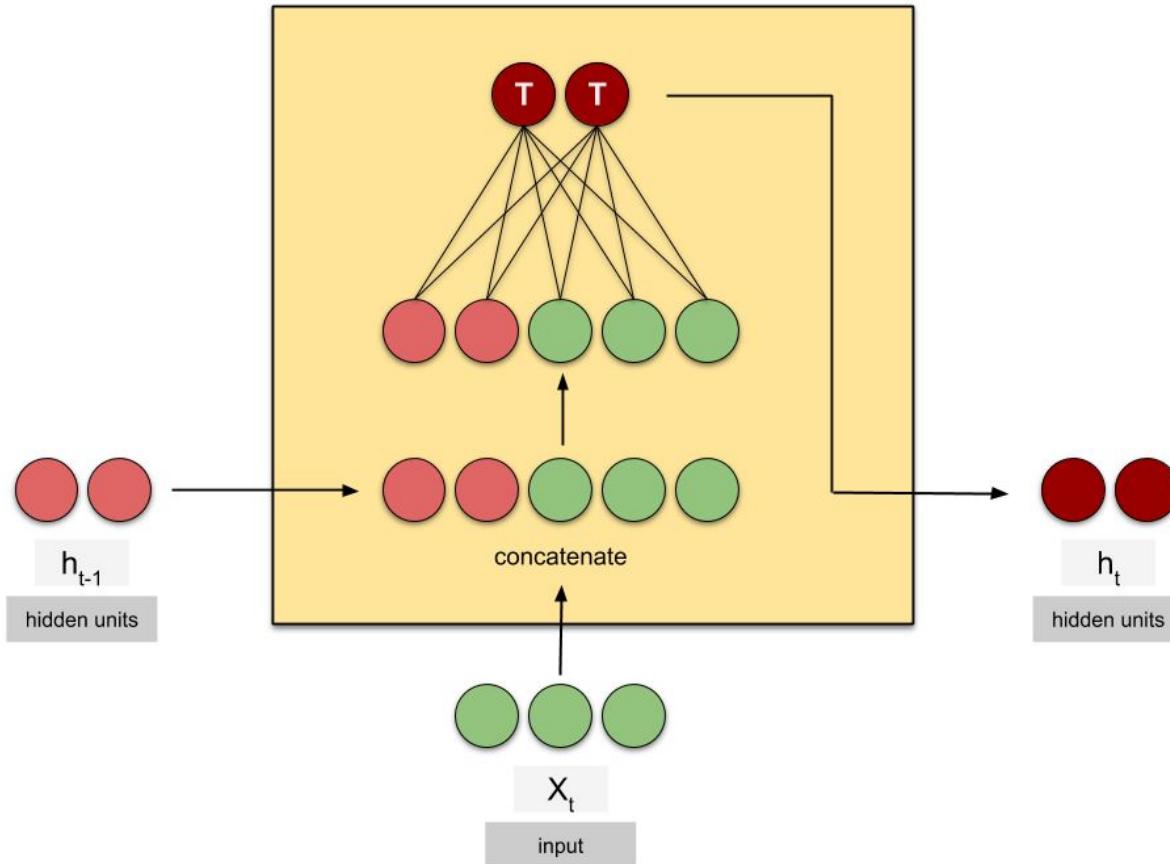
Recurrent  
Weights (U)

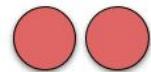


# Recurrent Neural Layer

The amount of learnable parameters increases with the additional U matrix.





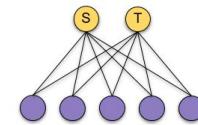
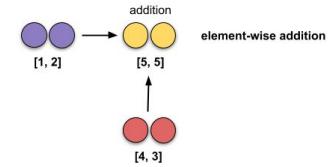
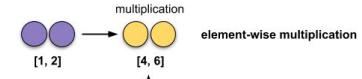
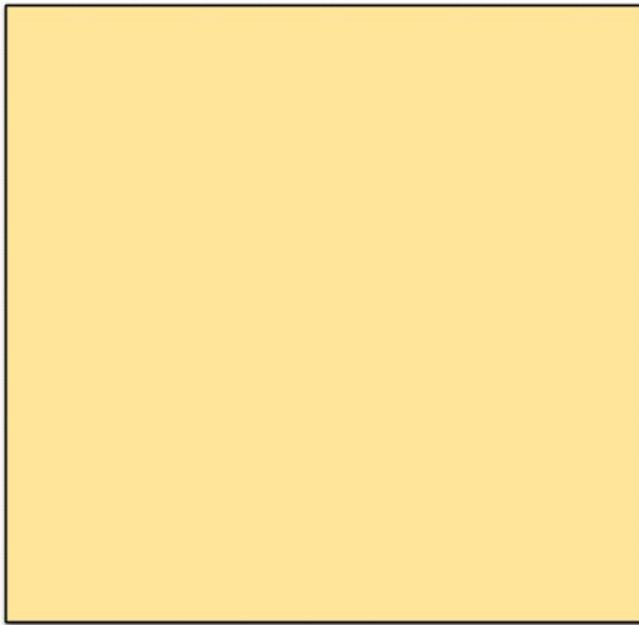


$h_{t-1}$   
hidden units



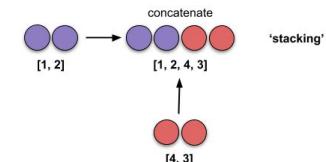
$X_t$

input



$S$  is a weighted sum of  $h_{t-1}$  and activated with a sigmoid function.

$T$  is a weighted sum of  $h_{t-1}$  and activated with a tanh function.



# Recurrent Neural Layer

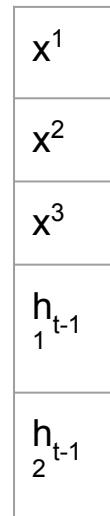
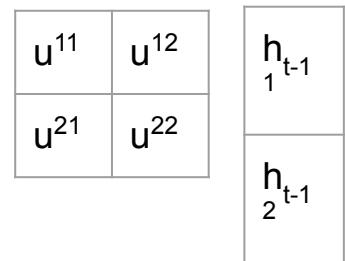
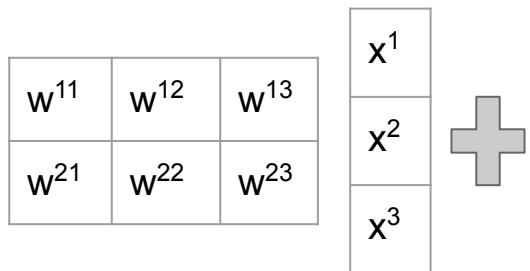
Prove that a concatenation of vectors  $h_{t-1}$  and  $x_t$  is equivalent to the following expression:

$$h_t = f(W \cdot x_t + U \cdot h_{t-1} + b)$$

# Recurrent Neural Layer

Prove that a concatenation of vectors  $h_{t-1}$  and  $x_t$  is equivalent to the following expression:

$$h_t = f(W \cdot x_t + U \cdot h_{t-1} + b)$$

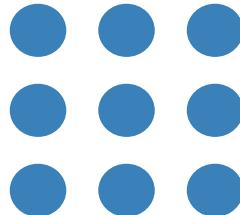


# Outline

1. Motivation
2. Recurrent Neural Layer
3. Forward & Backward Passes Through Time

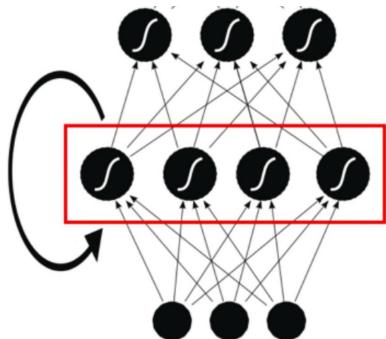
# Forward Pass in RNN

Front View



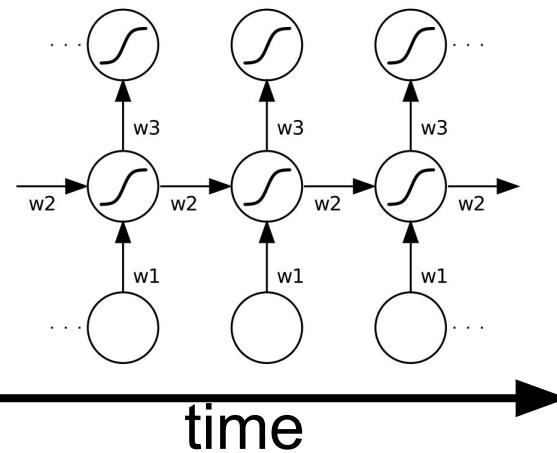
Rotation  
90°

Side View



time

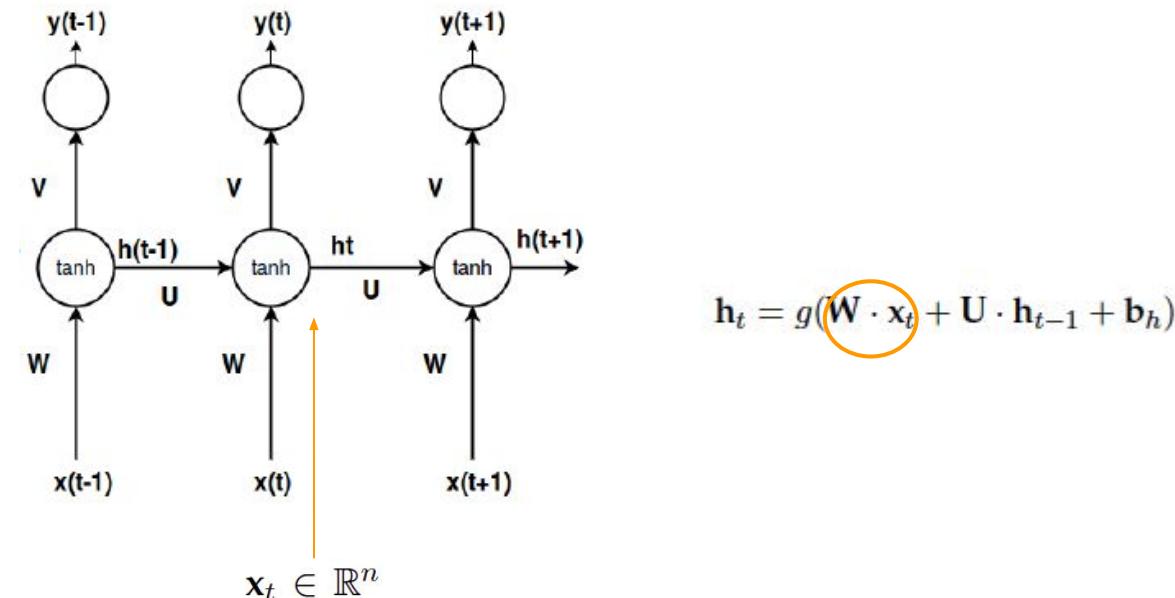
Unfold  
(Rotation  
90°)



# Forward Pass in RNN

Hence we have two data flows: **Forward in neural layers + time** propagation

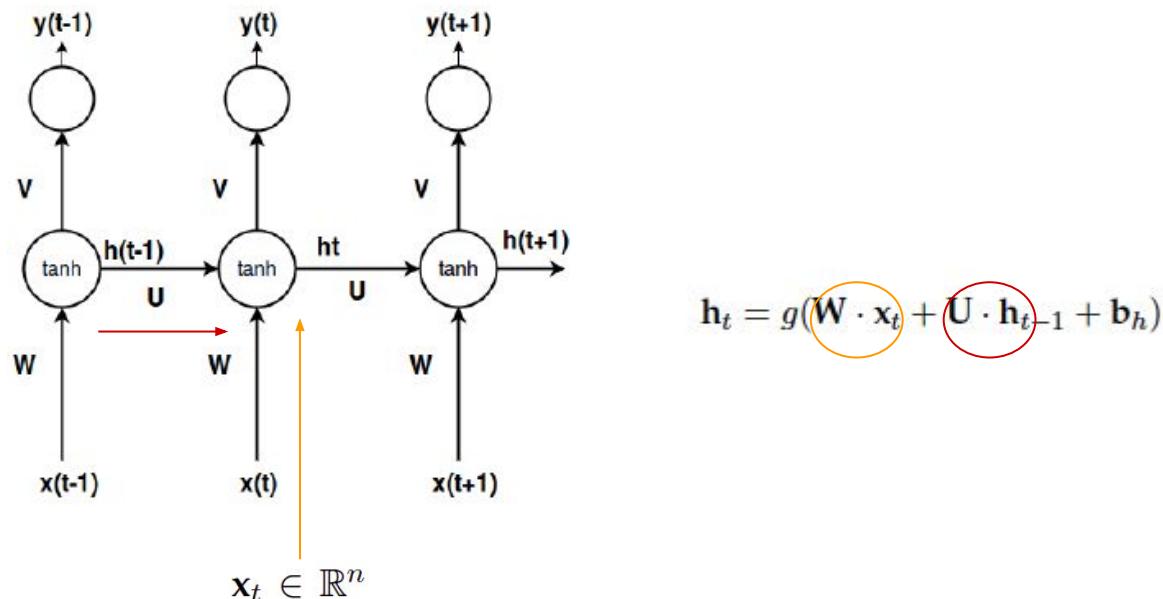
**BEWARE:** We have extra temporal depth now! Every time-step is an extra level of depth (as a deeper stack of layers in a feed-forward fashion!)



# Forward Pass in RNN

Hence we have two data flows: **Forward in layers + time** propagation

- Last time-step includes the context of our decisions recursively

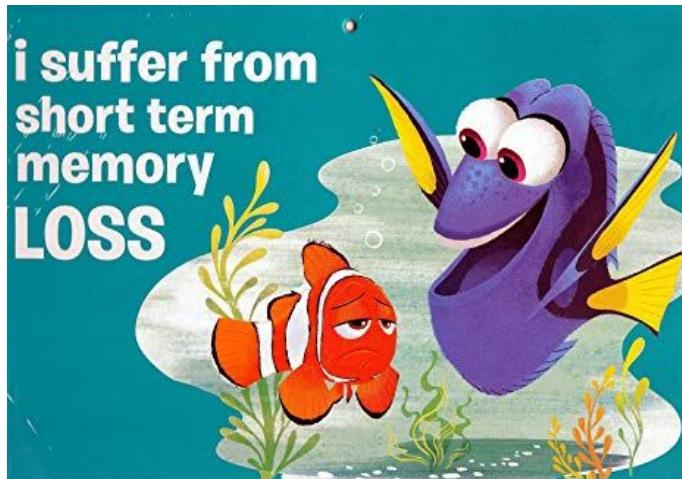


# Forward Pass in RNN

Main problem:

- **Long-term memory** (remembering quite far time-steps) **vanishes quickly** because of the recursive operation with non-linearities  $\mathbf{g}(\cdot)$  and  $\mathbf{U}$ .

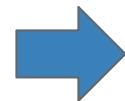
$$\mathbf{h}_t = \mathbf{g}(\mathbf{W} \cdot \mathbf{x}_t + \mathbf{U} \cdot \mathbf{g}(\cdots \mathbf{g}(\mathbf{W} \cdot \mathbf{x}_{t-T} + \mathbf{U} \cdot \mathbf{h}_{t-T} + \mathbf{b}_h) \cdots) + \mathbf{b}_h)$$



# Backpropagation Through Time (BPTT)

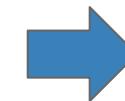
**Back Propagation Through Time (BPTT):** The training method has to take into account the time operations:

Total error at the output is the sum of errors at each time-step t



$$E(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{t=1}^T E_t(\mathbf{y}_t, \hat{\mathbf{y}}_t)$$

Total gradient is the sum of gradients at each time-step t

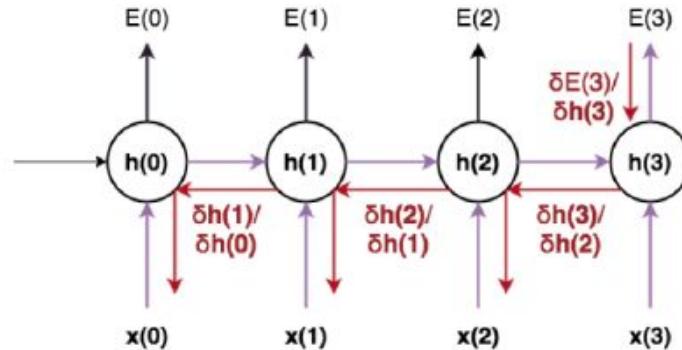


$$\frac{\partial E}{\partial \mathbf{W}} = \sum_{t=0}^{T-1} \frac{\partial E_t}{\partial \mathbf{W}}$$

**T: max amount of time-steps to do back-prop.** In Keras this is specified when defining the “input shape” to the RNN layer, by means of:  
(batch size, sequence length (T), input\_dim)

Input shape

3D tensor with shape `(nb_samples, timesteps, input_dim)`.



Example back-prop in time with 3 time-steps

# Backpropagation Through Time (BPTT)

Main problems:

- **Exploding / vanishing gradients:** During training gradients explode/vanish easily because of depth-in-time.

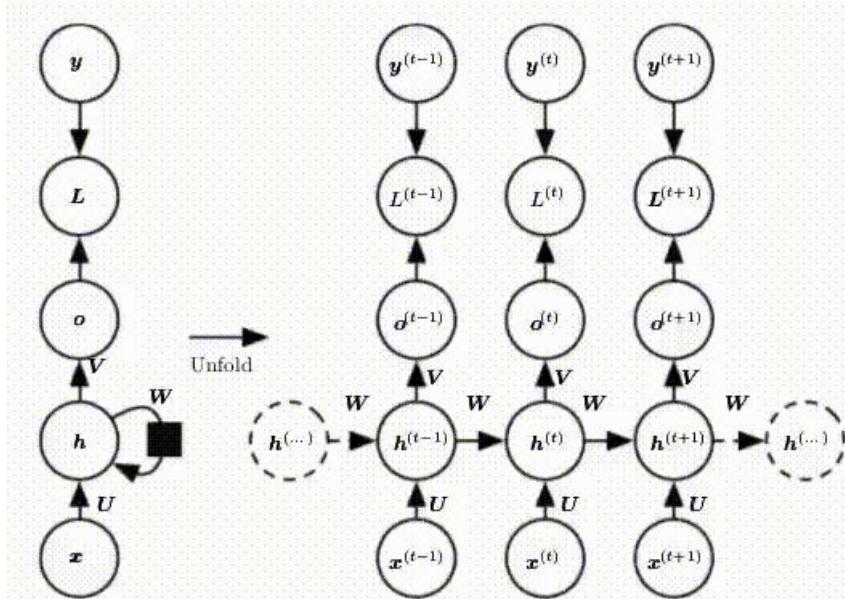


Figure: [Jordi Pons](#)

# Outline

1. Motivation
2. Recurrent Neural Layer
3. Forward & Backward Passes Through Time
4. Gated Recurrent Neurons
  - a. LSTM
  - b. GRU

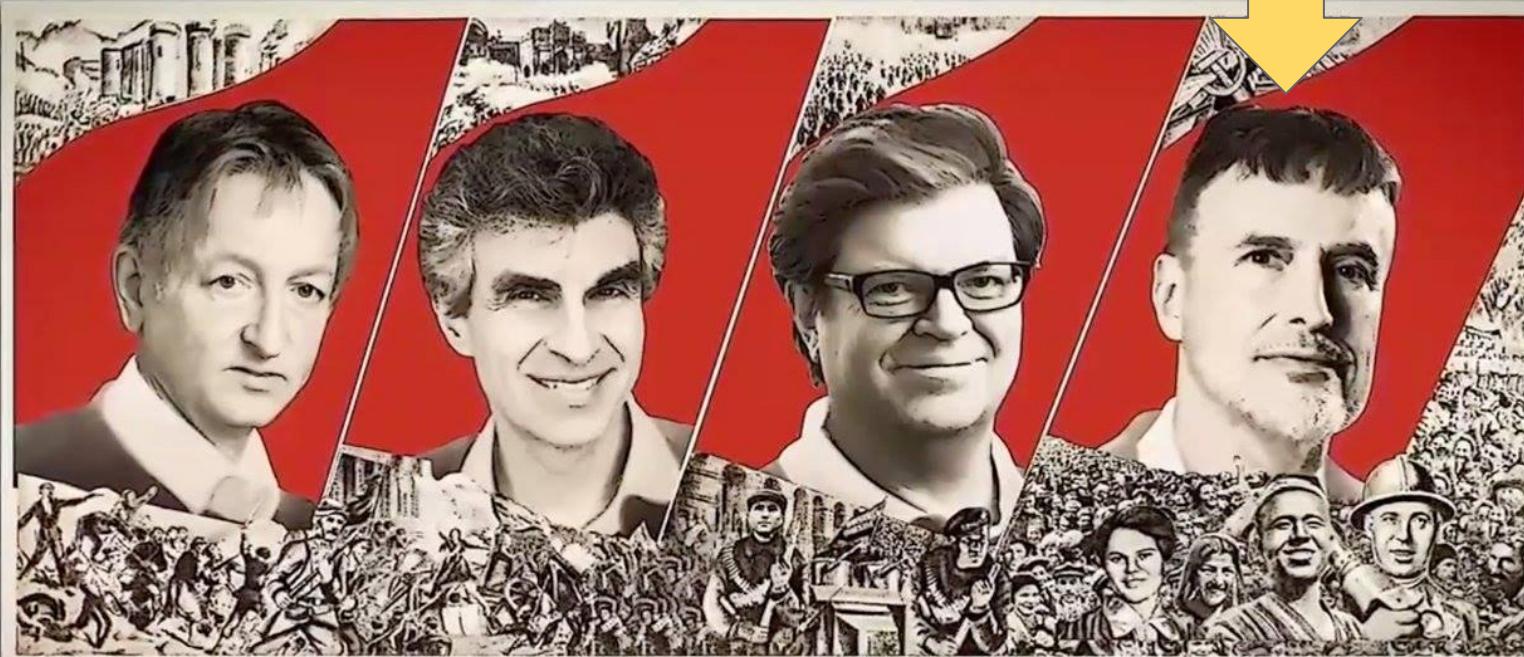


Figure: Michael Bronstein, "[Geometric Deep Learning](#)" (ICLR 2021)

# Long Short-Term Memory (LSTM)

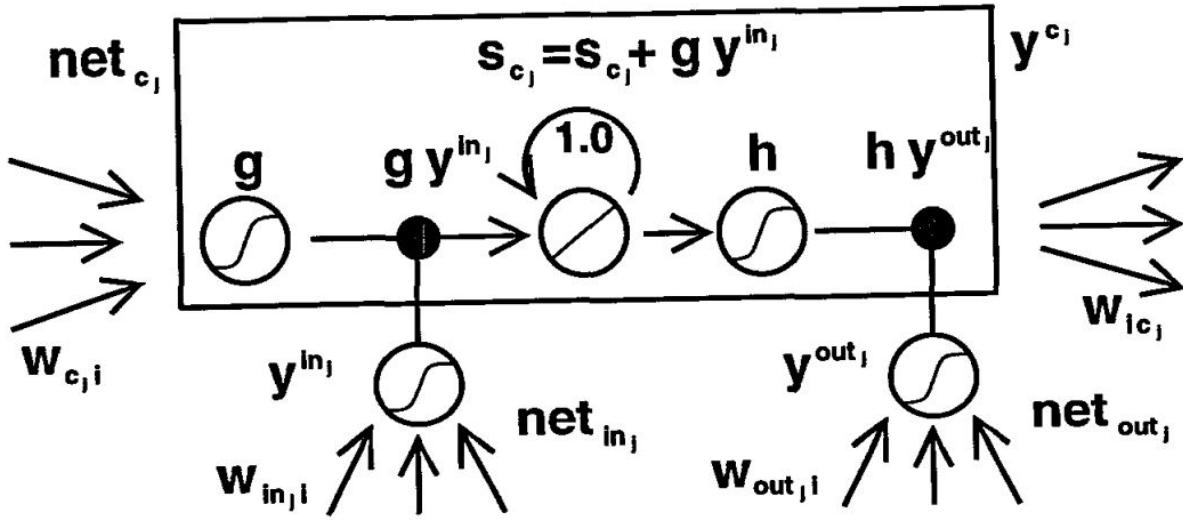
The New York Times, ["When A.I. Matures, It May Call Jürgen Schmidhuber 'Dad'"](#)  
(November 2016)



# Long Short-Term Memory (LSTM)

1744

Sepp Hochreiter and Jürgen Schmidhuber



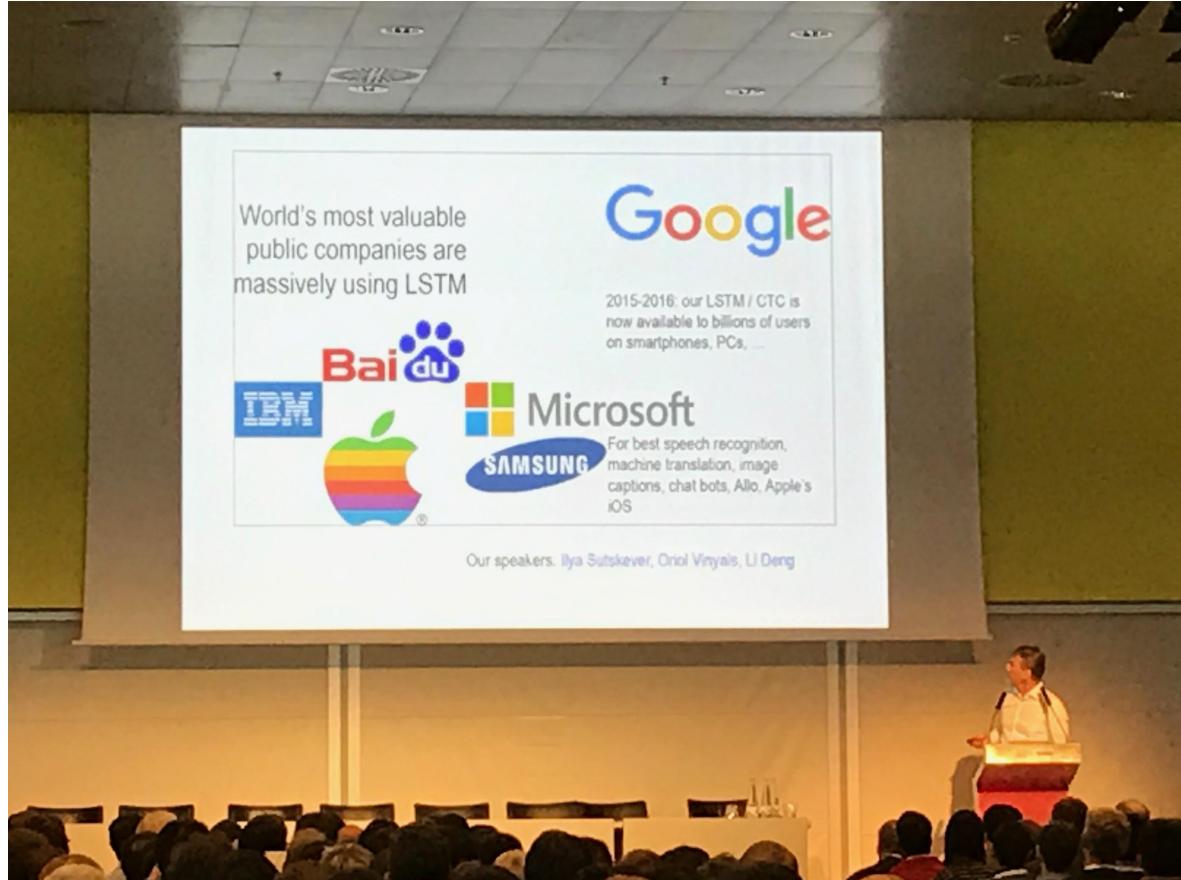
Hochreiter, Sepp, and Jürgen Schmidhuber. ["Long short-term memory."](#) Neural computation 9, no. 8 (1997): 1735-1780.

# Long Short-Term Memory (LSTM)

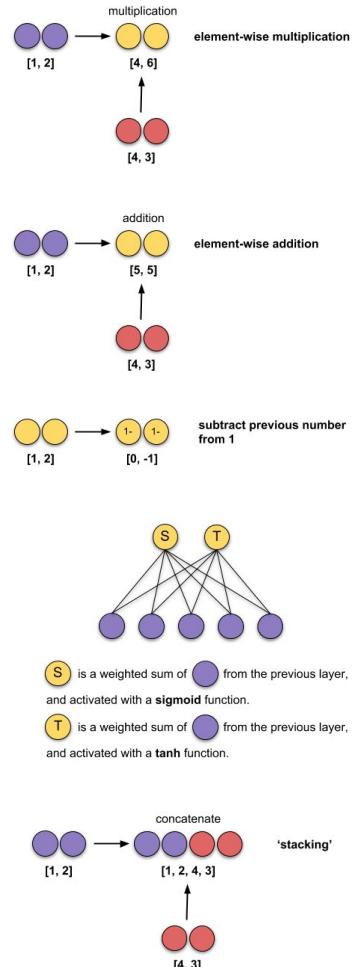
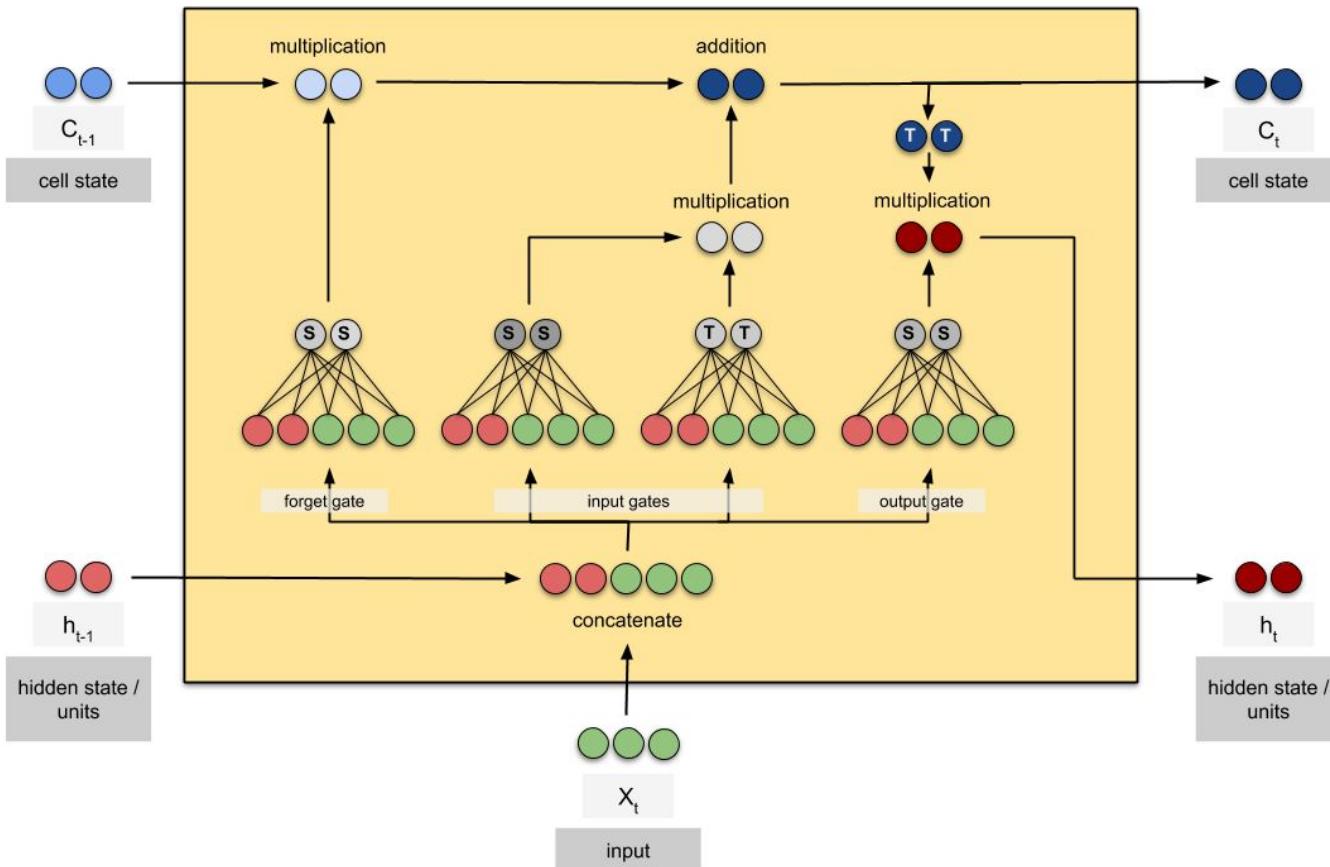


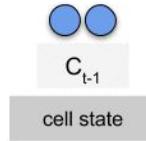
Jürgen  
Schmidhuber @  
NIPS 2016  
Barcelona

# Long Short-Term Memory (LSTM)



Jürgen Schmidhuber  
@ NIPS 2016 Barcelona





$C_{t-1}$

cell state



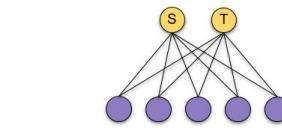
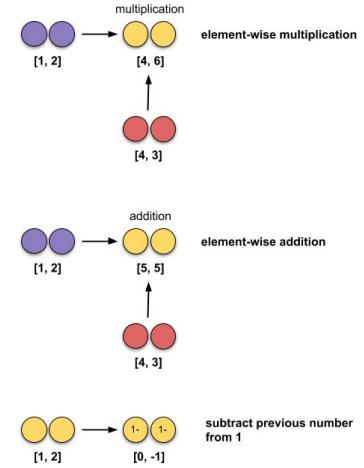
$h_{t-1}$

hidden state / units



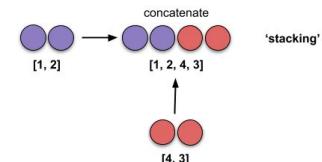
input

$X_t$



$S$  is a weighted sum of  $h_{t-1}$  from the previous layer, and activated with a **sigmoid** function.

$T$  is a weighted sum of  $C_{t-1}$  from the previous layer, and activated with a **tanh** function.



# Gating method



Solutions:

1. Change the way in which past information is kept → create the notion of **cell state**: a memory unit that keeps long-term information in a safer way by protecting it from recursive operations.
2. Make every RNN unit able to **forget whatever may not be useful anymore** by clearing that info from the cell state (optimized clearing mechanism)
3. Make every RNN unit able to decide whether **the current time-step information matters or not**, to accept or discard (optimized reading mechanism)
4. Make every RNN unit able to **output the decisions whenever it is ready to do so** (optimized output mechanism)

# Long Short-Term Memory (LSTM)

Three **gates** are governed by *sigmoid* units (btw [0,1]) define the control of in & out information with a product..

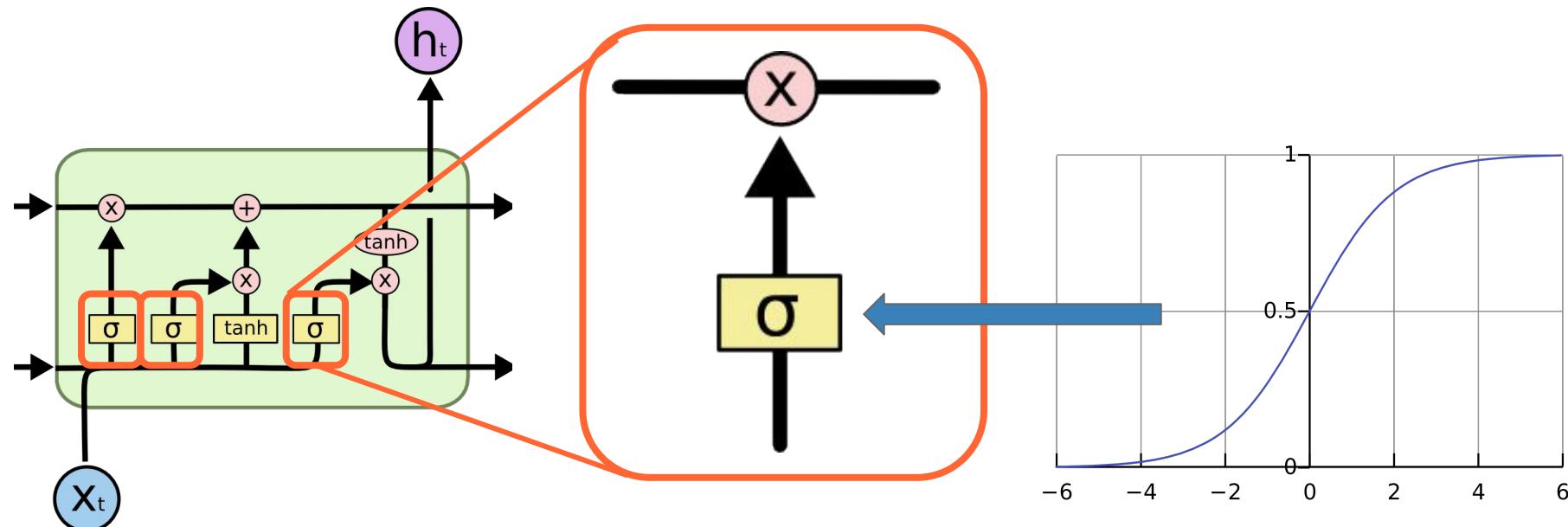
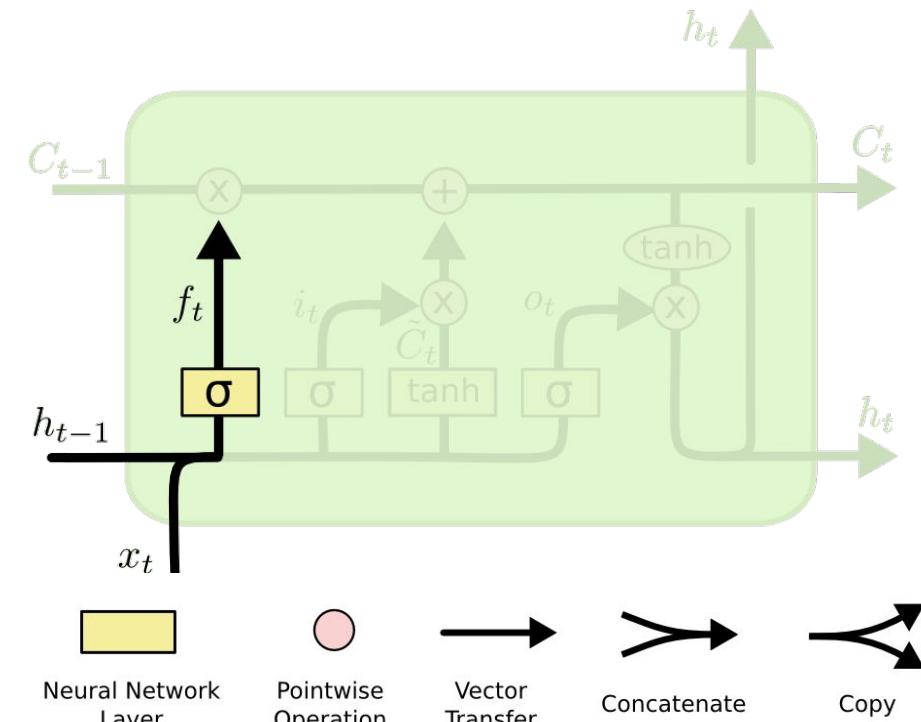


Figure: Cristopher Olah, ["Understanding LSTM Networks"](#) (2015)

# Long Short-Term Memory (LSTM)

Make every RNN unit able to **forget whatever may not be useful anymore** by clearing that info from the cell state (optimized clearing mechanism)



**Forget Gate:**

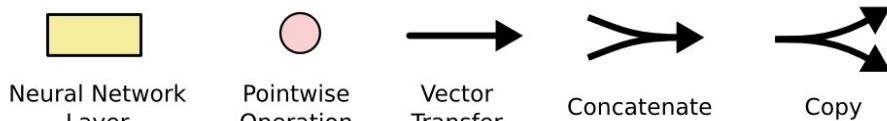
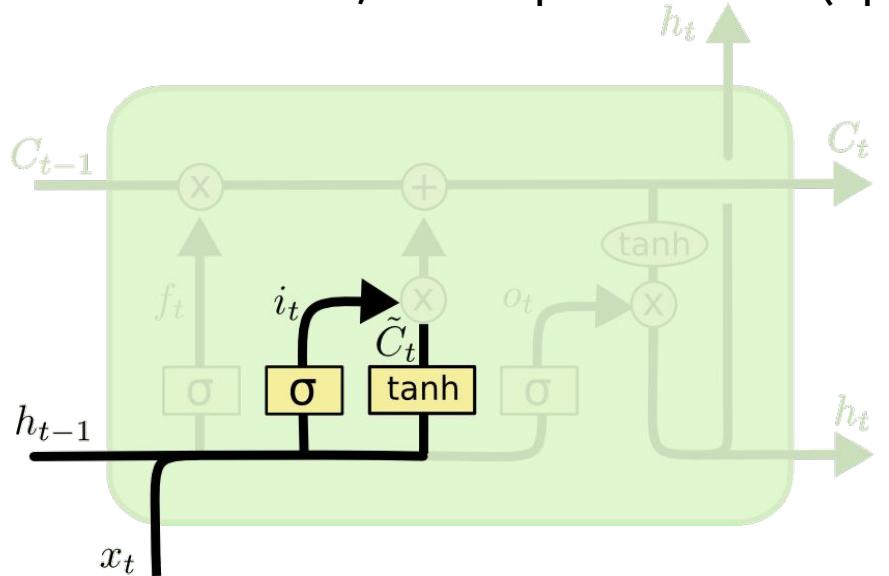
$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

Concatenate

Concatenate

# Long Short-Term Memory (LSTM)

Make every RNN unit able to decide whether **the current time-step information matters or not**, to accept or discard (optimized reading mechanism)



## Input Gate Layer

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

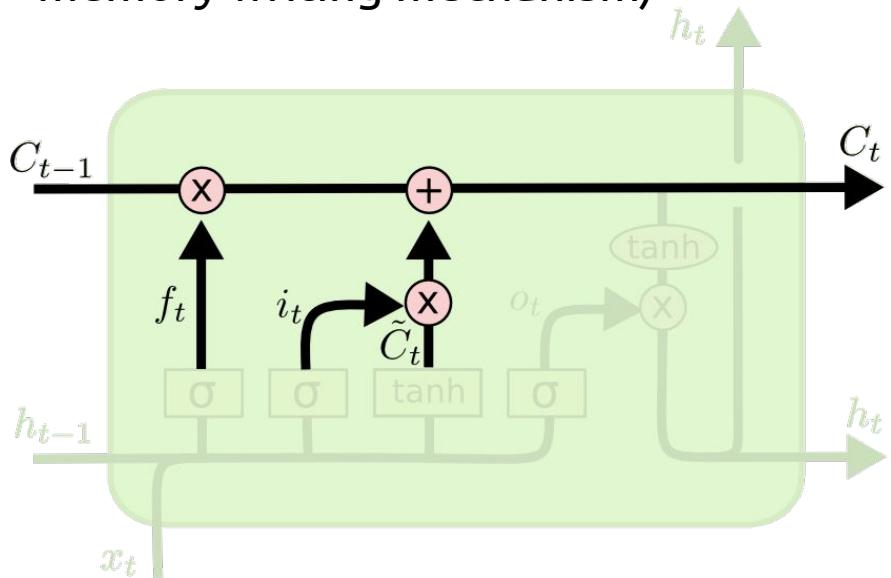
## New contribution to cell state

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

Classic neuron

# Long Short-Term Memory (LSTM)

Make every RNN unit able to decide how to **to update the cell state** (optimized memory writing mechanism)

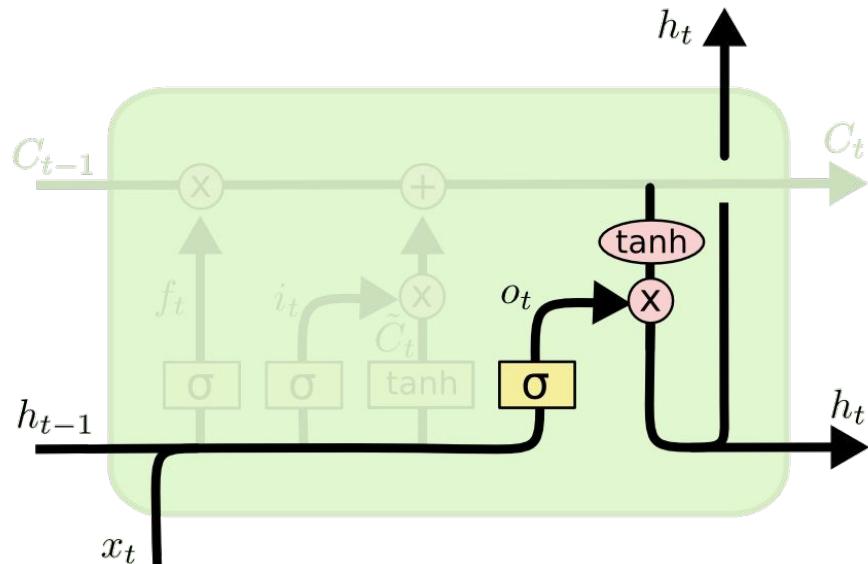


**Forget + Input Gates =  
Update Cell State (memory):**

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

# Long Short-Term Memory (LSTM)

Make every RNN unit able to **output the decisions whenever it is ready to do so** (optimized output mechanism)



## Output Gate Layer

$$o_t = \sigma (W_o [ h_{t-1}, x_t ] + b_o)$$

**Output to next layer & timestep**

$$h_t = o_t * \tanh (C_t)$$

# Long Short-Term Memory (LSTM)

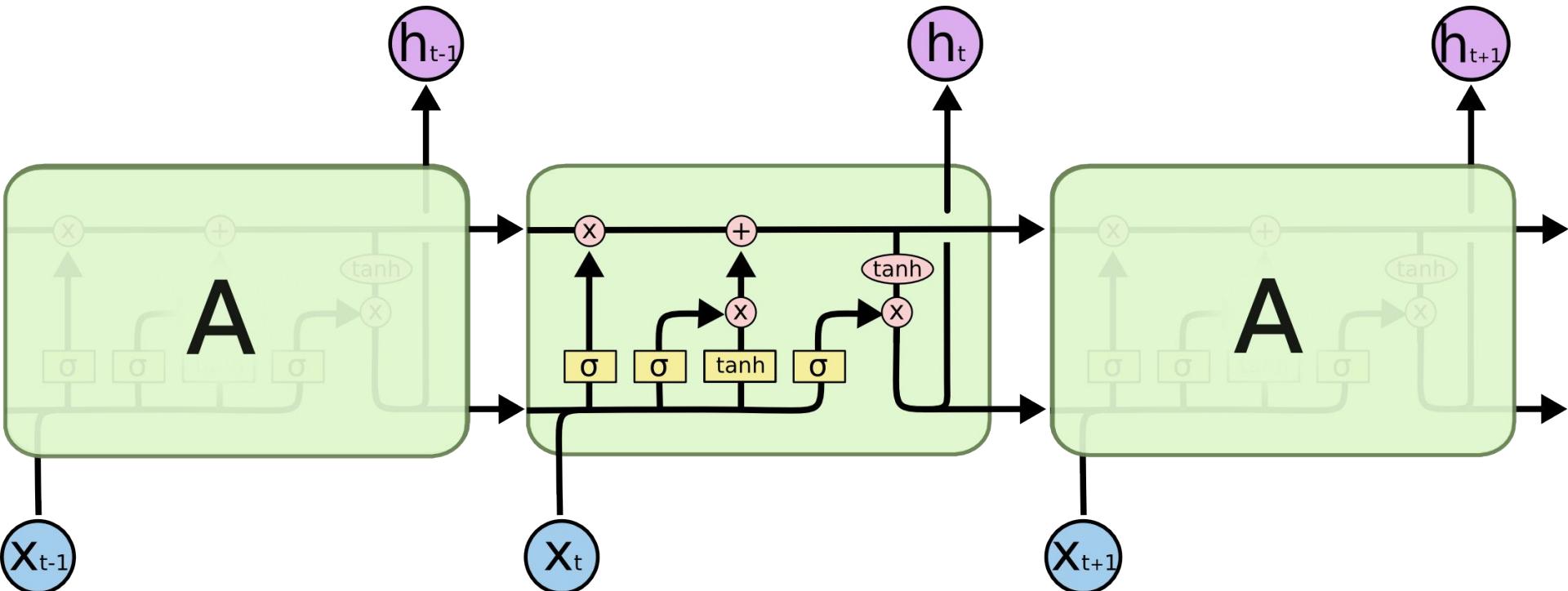
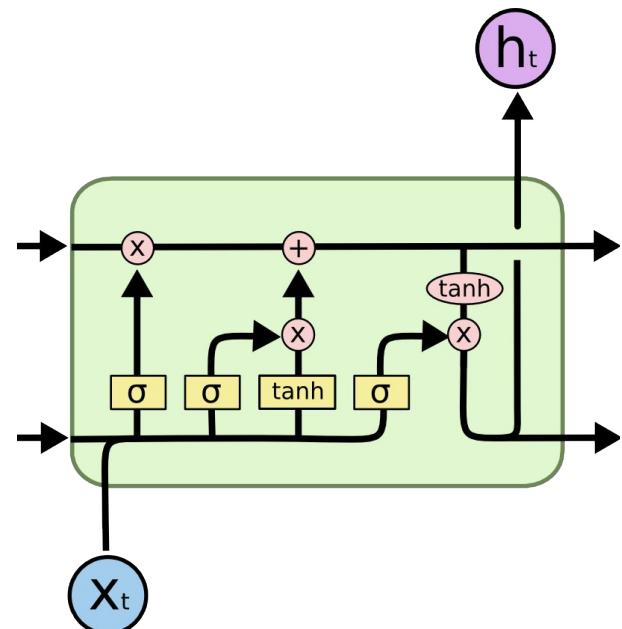


Figure: Cristopher Olah, ["Understanding LSTM Networks"](#) (2015)

# Long Short-Term Memory (LSTM)



Compared to a non-gated RNN, an LSTM has four times more parameters because of the additional neurons that govern the gates:

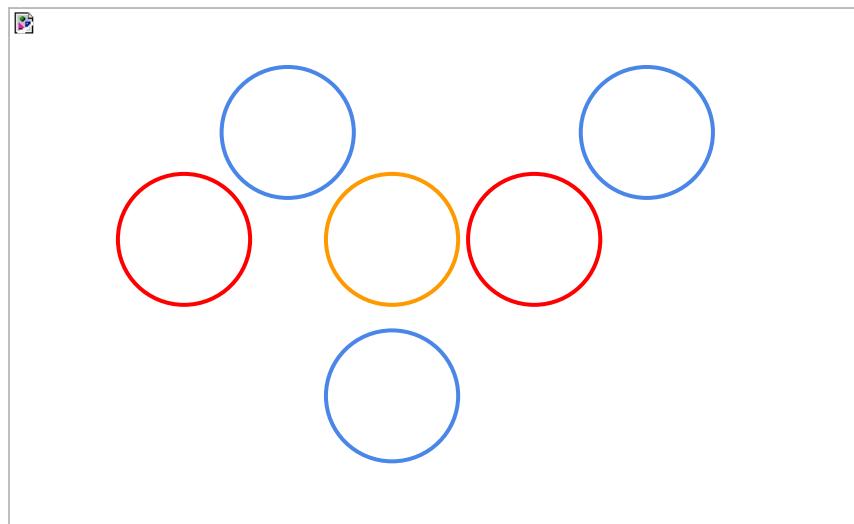
$$N_{params}^i = 4 \times (N_{inputs}^i \times N_{units}^i + N_{units}^i \times N_{units}^i + N_{units}^i)$$

3 sigmoid gates +  
input activation ( $\tanh$  in the figure)

# Long Short Term Memory (LSTM) cell

Updating an LSTM cell requires 6 computations:

1. Gates
2. Activation units
3. Cell state



Computation Flow

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

$\hat{C}_t = \tanh(W_c x_t + U_c h_{t-1} + b_c)$

$f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f)$

$C_t = i_t \odot \hat{C}_t + f_t \odot C_{t-1}$

$o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o)$

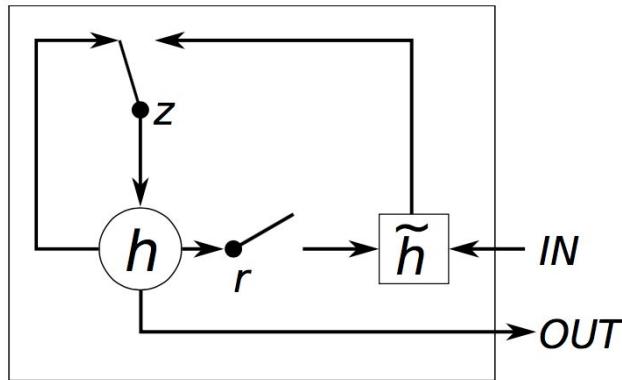
$h_t = o_t \odot \tanh(C_t)$

# Outline

1. Motivation
2. Recurrent Neural Layer
3. Forward & Backward Passes Through Time
4. Gated Recurrent Neurons
  - a. LSTM
  - b. GRU

# Gated Recurrent Unit (GRU)

GRU obtain a similar performance as LSTM with one gate less.



$$u_i = \sigma(W^{(u)}x_i + U^{(u)}h_{i-1} + b^{(u)}) \quad (1)$$

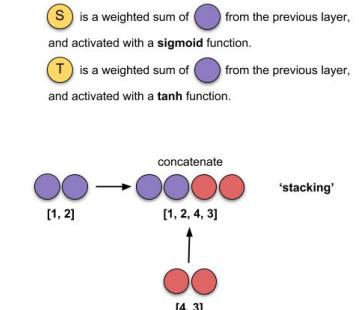
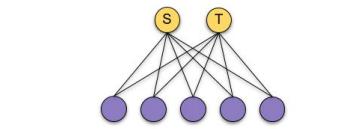
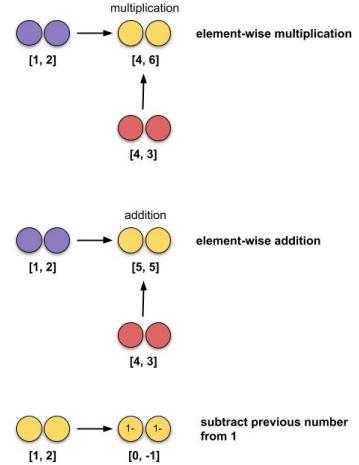
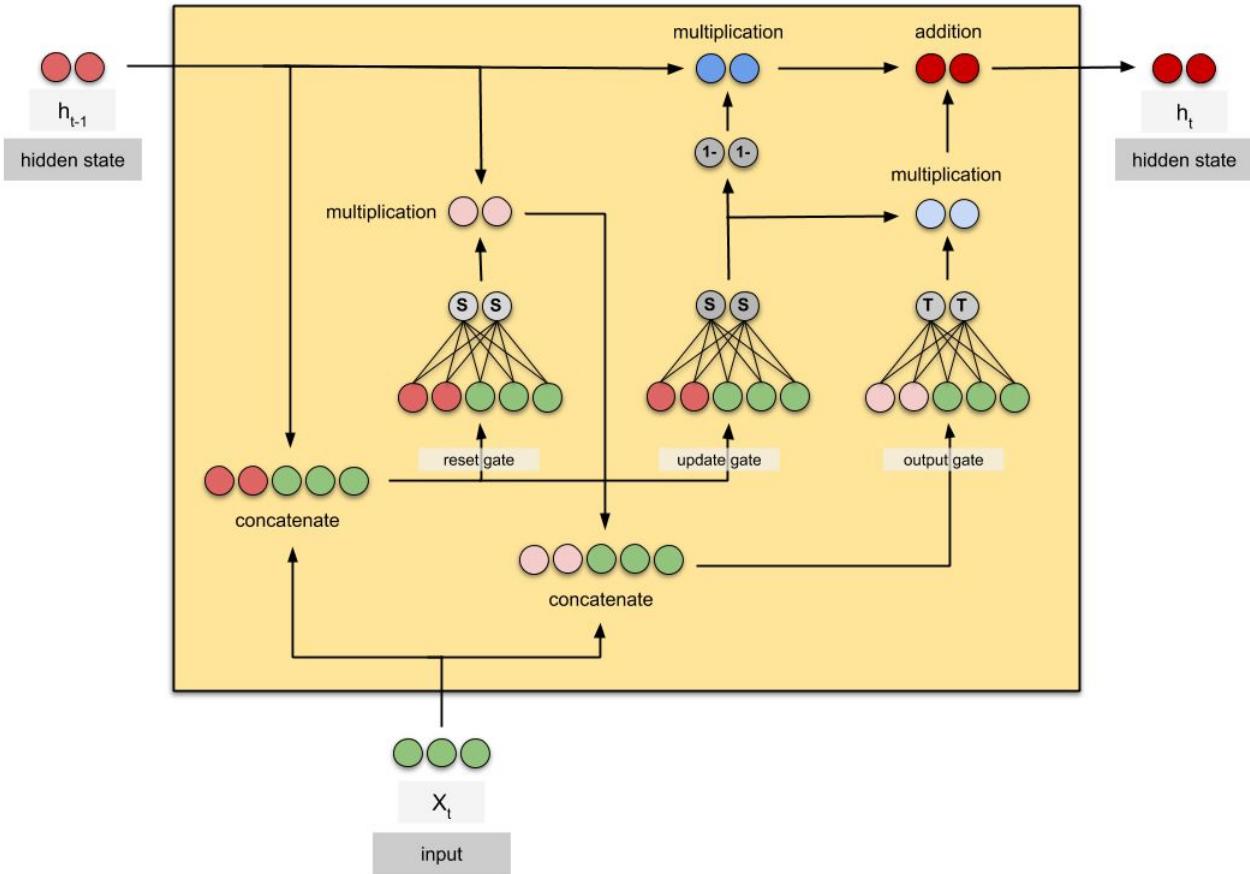
$$r_i = \sigma(W^{(r)}x_i + U^{(r)}h_{i-1} + b^{(r)}) \quad (2)$$

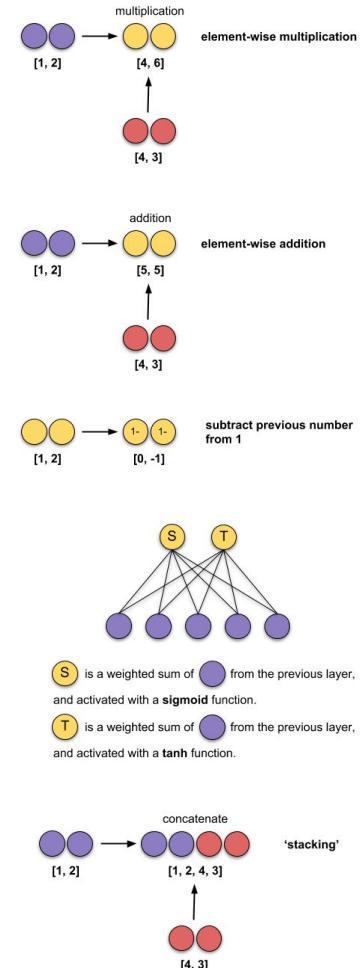
$$\tilde{h}_i = \tanh(Wx_i + r_i \circ Uh_{i-1} + b^{(h)}) \quad (3)$$

$$h_i = u_i \circ \tilde{h}_i + (1 - u_i) \circ h_{i-1} \quad (4)$$

$$N_{params}^i = 3 \times (N_{inputs}^i \times N_{units}^i + N_{units}^i \times N_{units}^i + N_{units}^i)$$

#GRU Cho, Kyunghyun, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. ["Learning phrase representations using RNN encoder-decoder for statistical machine translation."](#) EMNLP 2014.





# Outline

1. Motivation
2. Recurrent Neural Layer
3. Forward & Backward Passes Through Time
4. Gated Recurrent Neurons
  - a. LSTM
  - b. GRU
5. Stacked RNNs

# Deep (stacked) RNNs

Recurrent layers can be stacked as any non-recurrent layer.

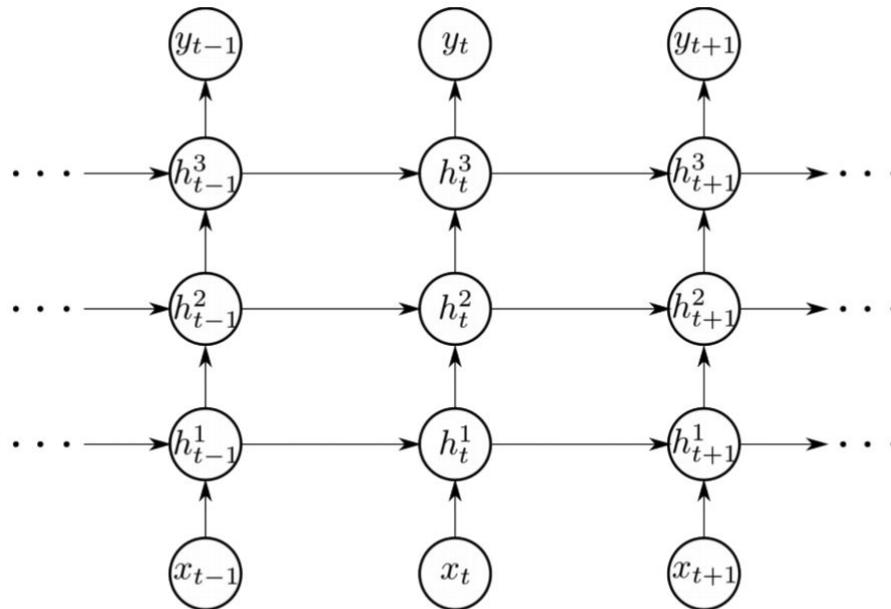


Figure:  
[Lambert](#)  
(Stanford University).

# Bidirectional RNN (BRNN)

Sequences can be process in both directions by stacking two RNN layers.

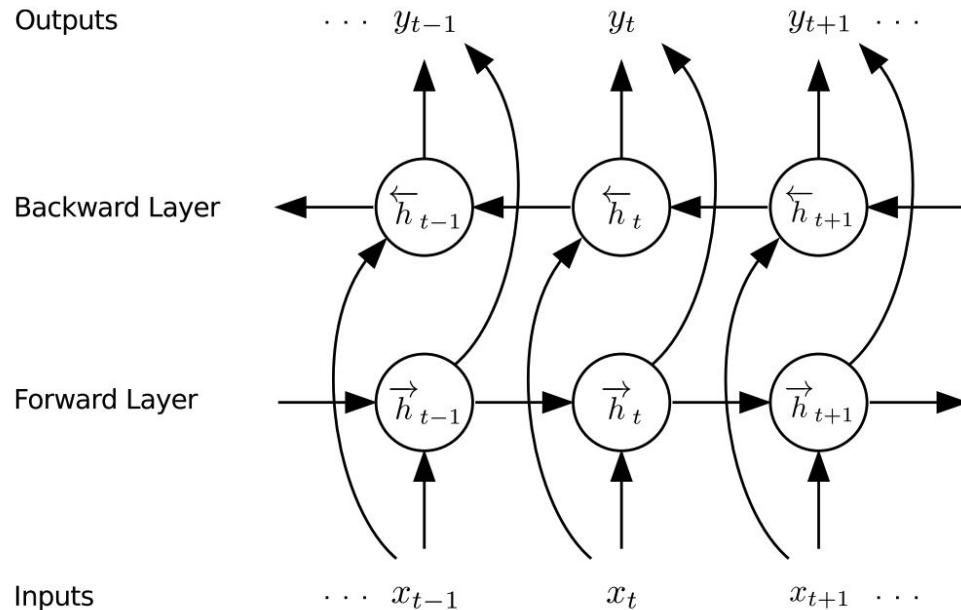


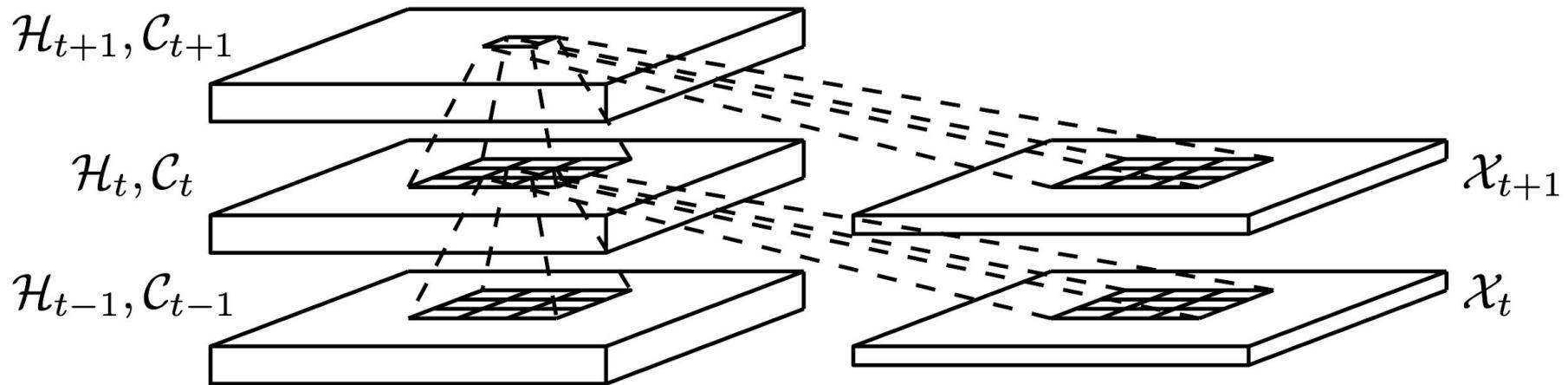
Figure:  
Graves et al. 2013.

#BRNN Schuster, M., & Paliwal, K. K. (1997). [Bidirectional recurrent neural networks](#). IEEE transactions on Signal Processing, 45(11), 2673-2681.

# Outline

1. Motivation
2. Recurrent Neural Layer
3. Forward & Backward Passes Through Time
4. Gated Recurrent Neurons
  - a. LSTM
  - b. GRU
5. Stacked RNNs
6. Advanced architectures

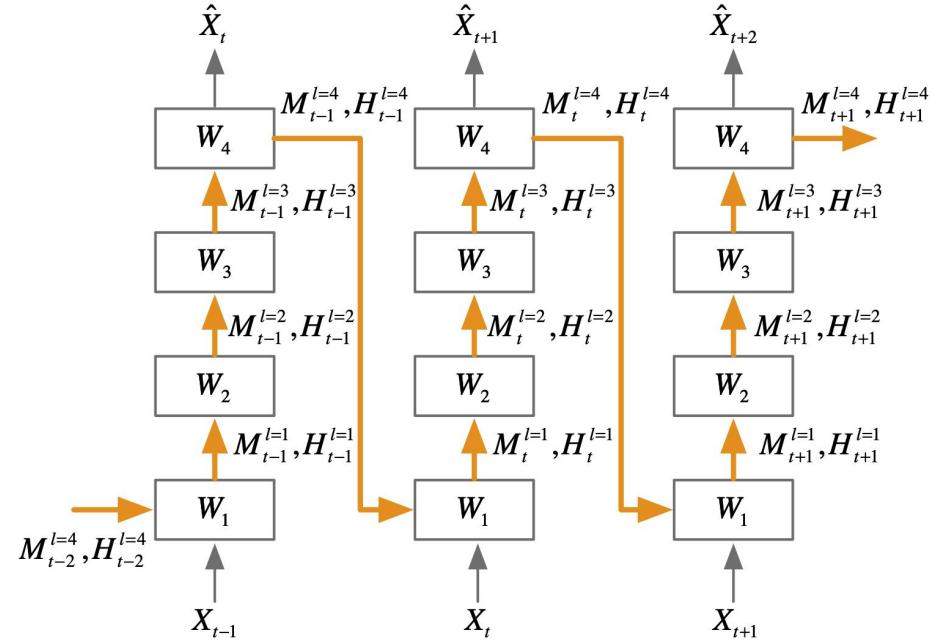
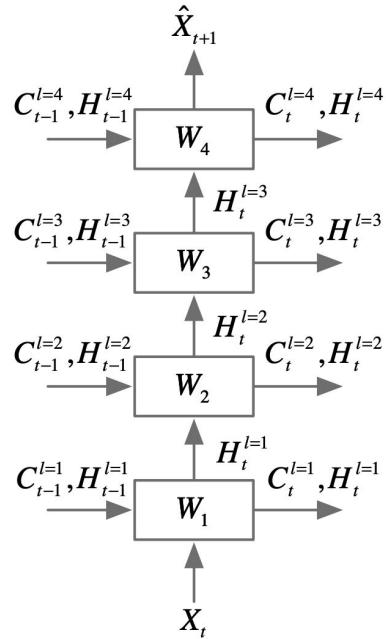
# Conv-LSTM



#ConvLSTM Xingjian, S. H. I., Chen, Z., Wang, H., Yeung, D. Y., Wong, W. K., & Woo, W. C. [Convolutional LSTM network: A machine learning approach for precipitation nowcasting](#). NIPS 2015.

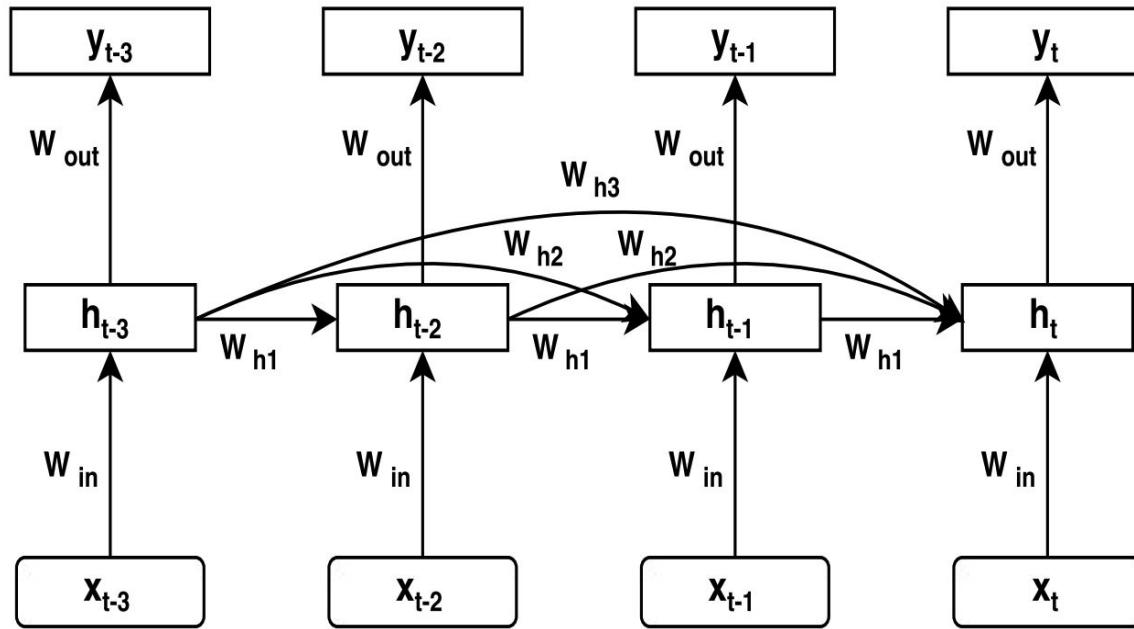
# Spatiotemporal LSTM (ST-LSTM)

Cell states are no longer constrained inside each LSTM unit. Instead, they are allowed to zigzag in two directions: across stacked RNN layers vertically and through all RNN states horizontally

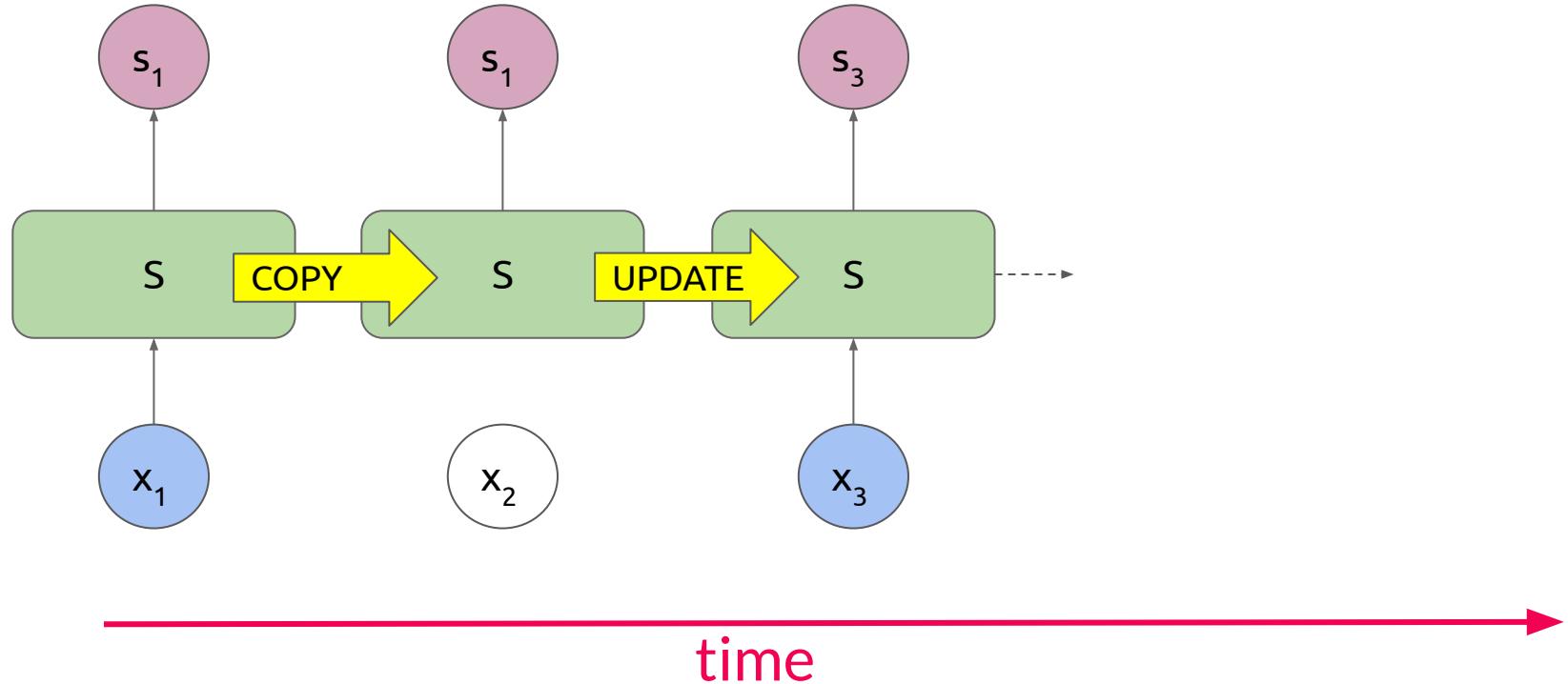


# Higher Order RNN

More memory units can be added to connect with further temporal states.

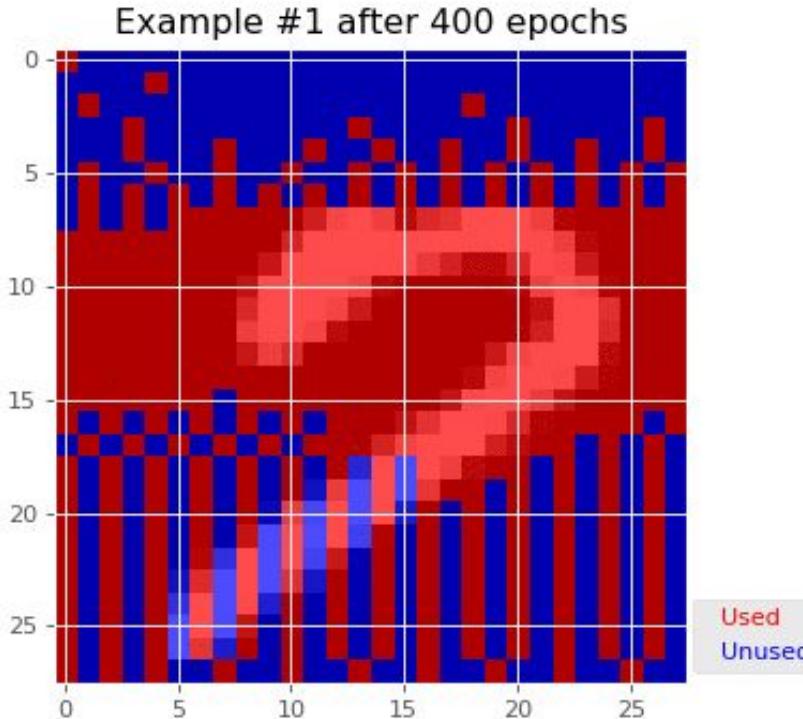


# SkipRNN



#SkipRNN Victor Campos, Brendan Jou, Xavier Giro-i-Nieto, Jordi Torres, and Shih-Fu Chang. ["Skip RNN: Learning to Skip State Updates in Recurrent Neural Networks"](#), ICLR 2018.

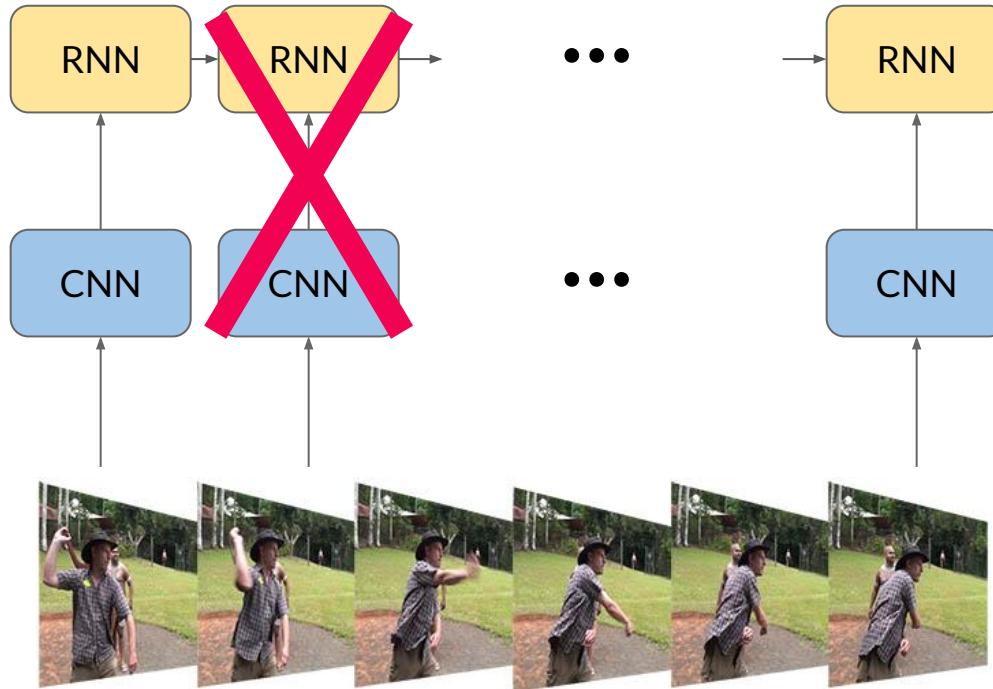
# SkipRNN



~95% acc

Used  
Unused

# SkipRNN



After processing a frame, let the RNN decide how many future frames can be skipped

In skipped frames, simply copy the output and state from the previous time step

There is no ground truth for which frames can be skipped. The RNN **learns** it by itself during training!

# SkipRNN

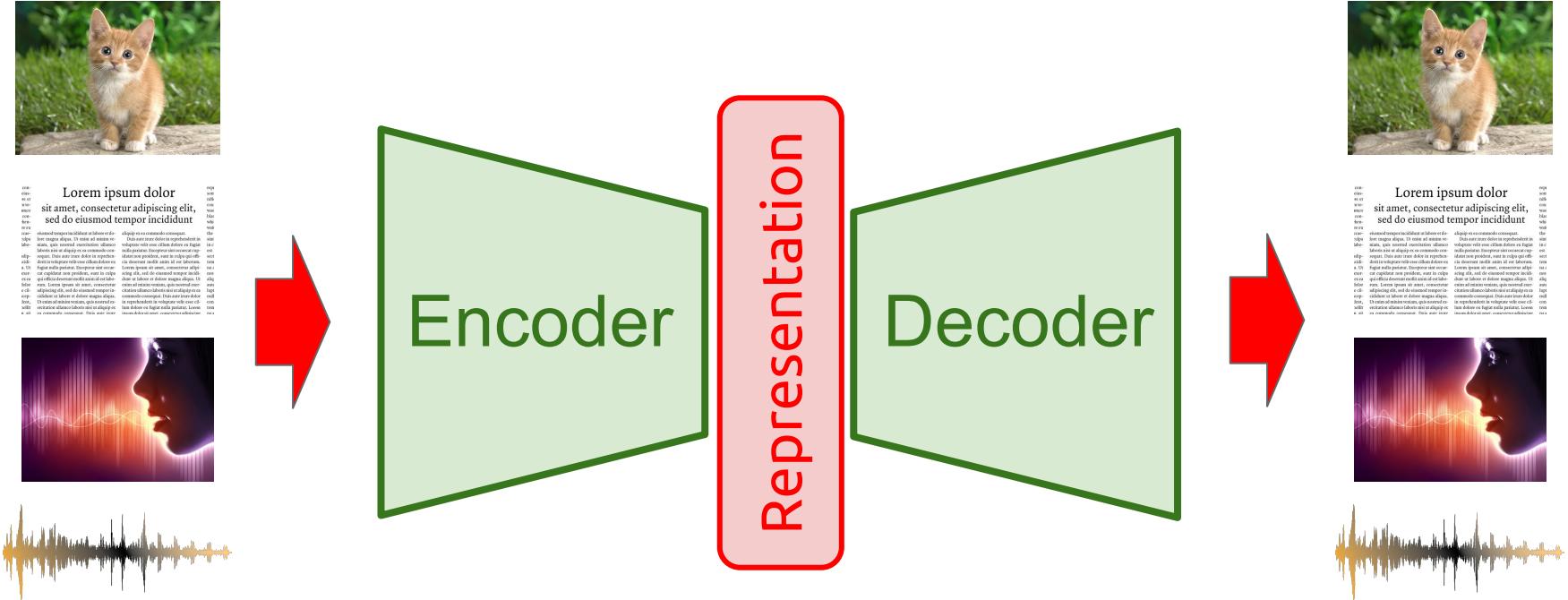


Used  
Unused

# Outline

1. Motivation
2. Recurrent Neural Layer
3. Forward & Backward Passes Through Time
4. Gated Recurrent Neurons
  - a. LSTM
  - b. GRU
5. Advanced Architectures
6. Encoder-Decoder Architectures

# Encoder-Decoder Applications



# Encoding of a Sequence of Words

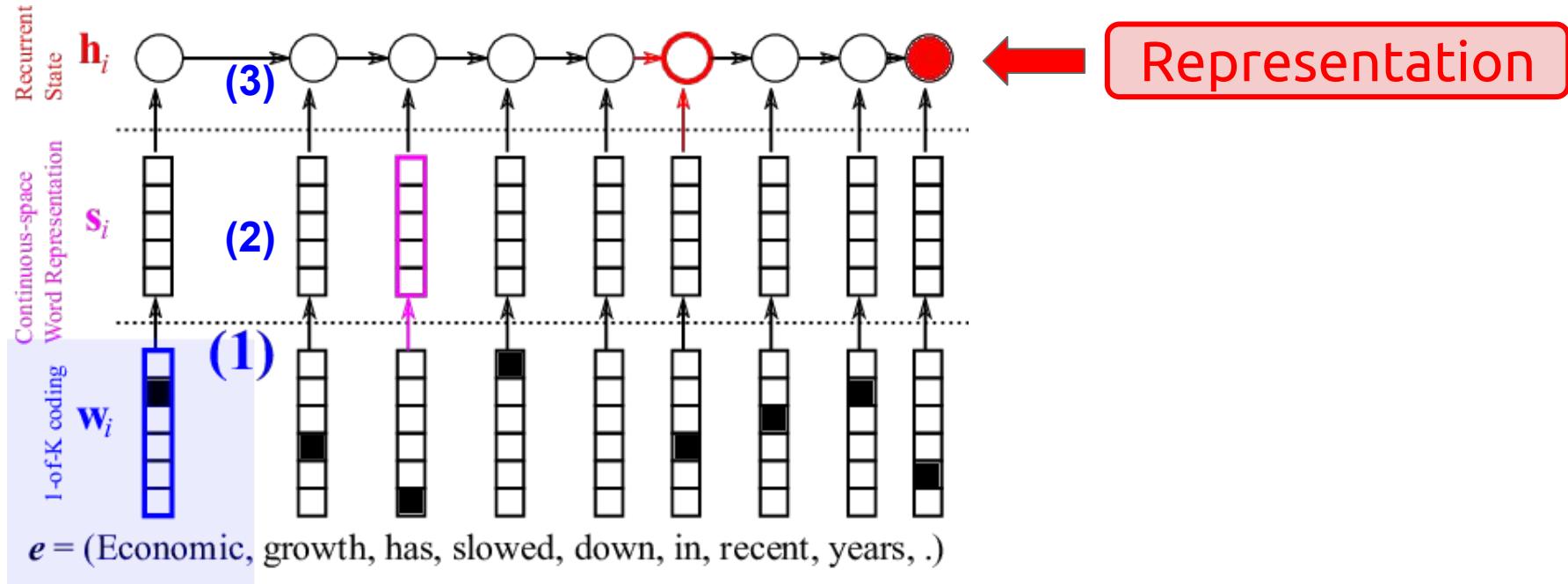
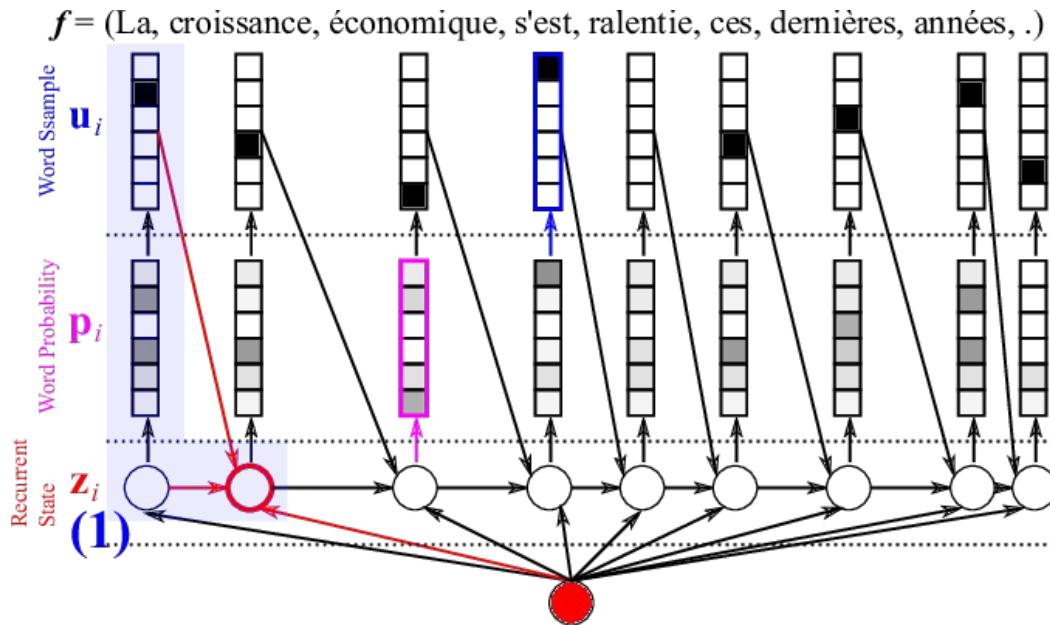


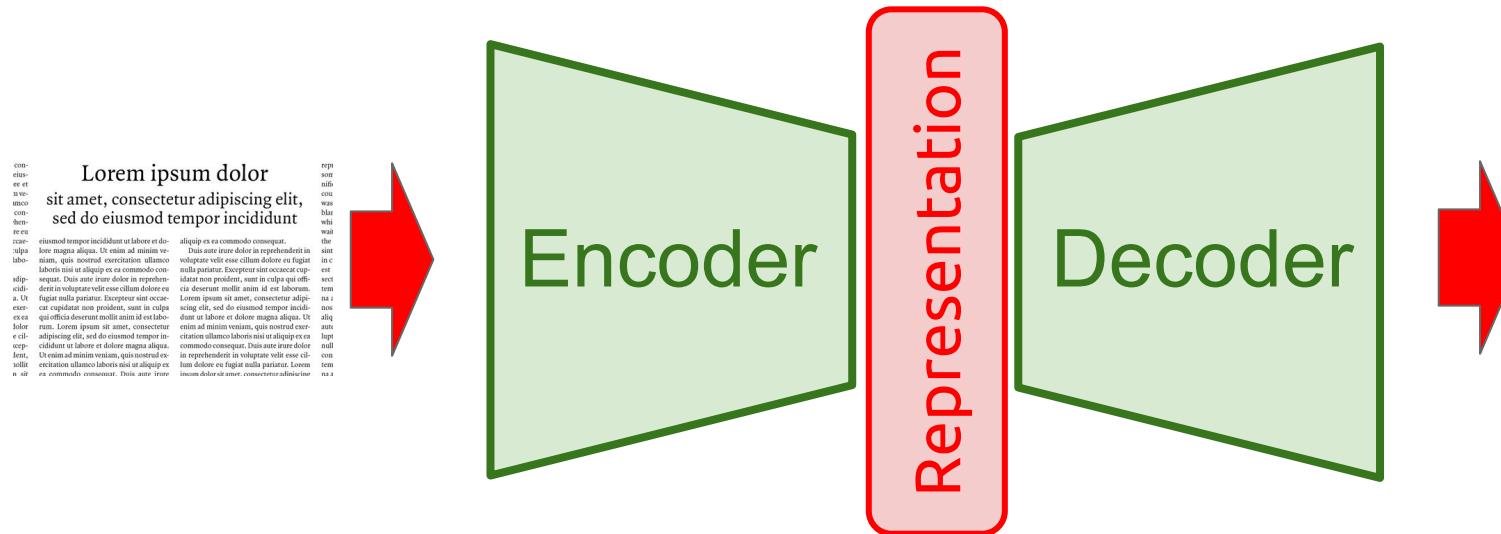
Fig: Kyunghyun Cho, "[Introduction to Neural Machine Translation with GPUs](#)" (2015)

Cho, Kyunghyun, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. "[Learning phrase representations using RNN encoder-decoder for statistical machine translation.](#)" EMNLP 2014.

# Decoding of a Sequence of Words



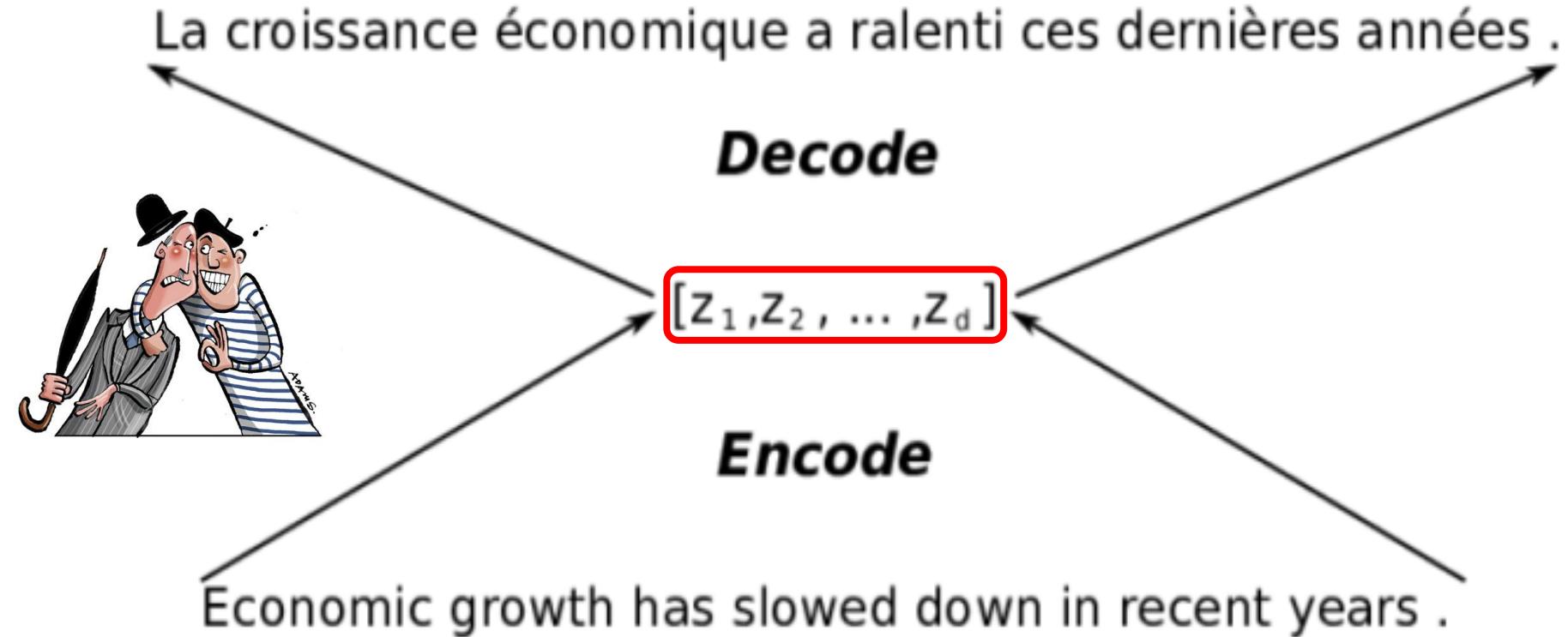
Representation



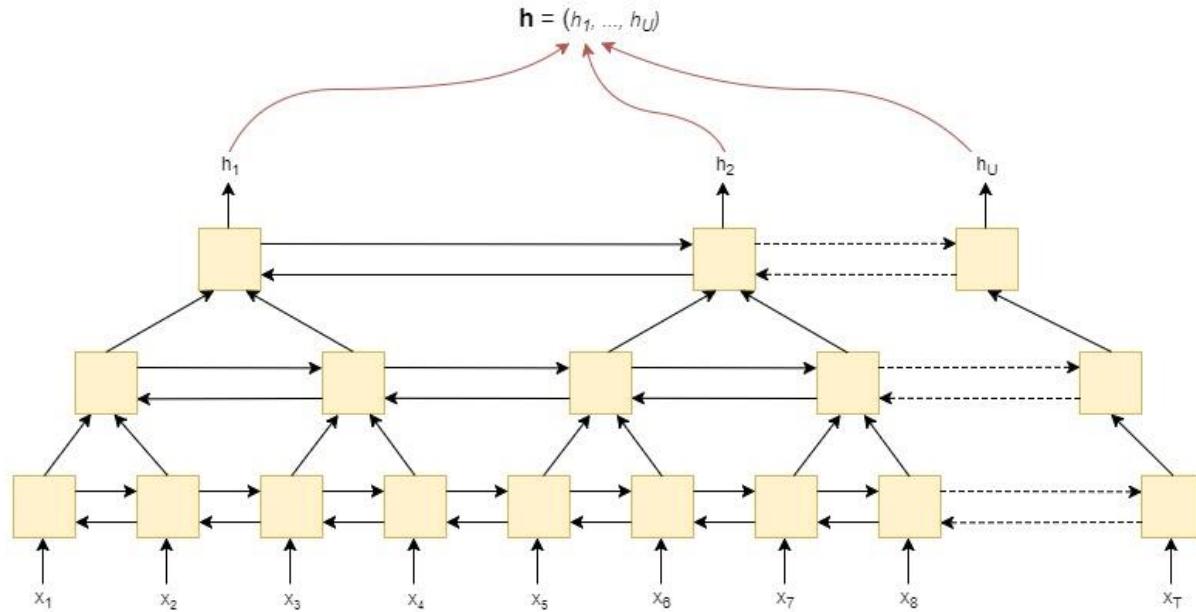
con-  
eius-  
et et  
n ve-  
unco  
con-  
ben-

**Lorem ipsum dolor**  
sit amet, consectetur adipiscing el-  
sed do eiusmod tempor incididunt

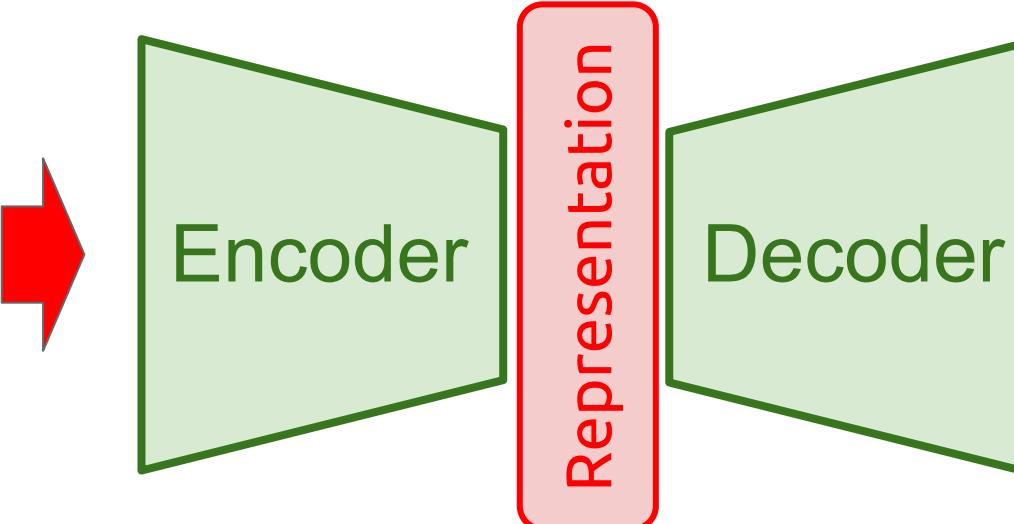
# Neural Machine Translation (NMT)



# Speech Encoding



Chan, William, Navdeep Jaitly, Quoc Le, and Oriol Vinyals. ["Listen, attend and spell: A neural network for large vocabulary conversational speech recognition."](#) ICASSP 2016.



**Encoder**  
Decoder

**Representation**

**Output Image**

**Input Image**

**Autoencoder Architecture**

Diagram illustrating an Autoencoder architecture. The process starts with an input image, which is processed by an Encoder to produce a compressed representation. This representation is then decoded back into the original image shape by a Decoder.

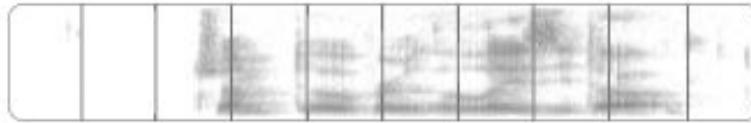
The diagram shows three main components: Encoder (green), Representation (red), and Decoder (green). The input image is transformed into a smaller representation, which is then expanded back into the original image size.

This architecture is commonly used for tasks like image denoising, compression, and generating new images based on existing ones.

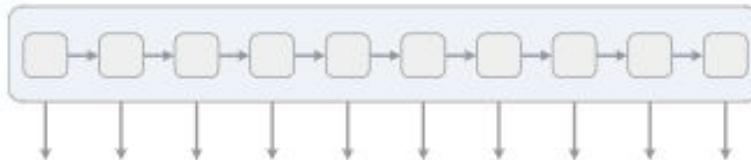
Key components of an autoencoder:

Common applications of autoencoders include:

# Automatic Speech Recognition (ASR)



We start with an input sequence,  
like a spectrogram of audio.



The input is fed into an RNN,  
for example.

h	h	h	h	h	h	h	h	h	h
e	e	e	e	e	e	e	e	e	e
l	l	l	l	l	l	l	l	l	l
o	o	o	o	o	o	o	o	o	o
$\epsilon$									

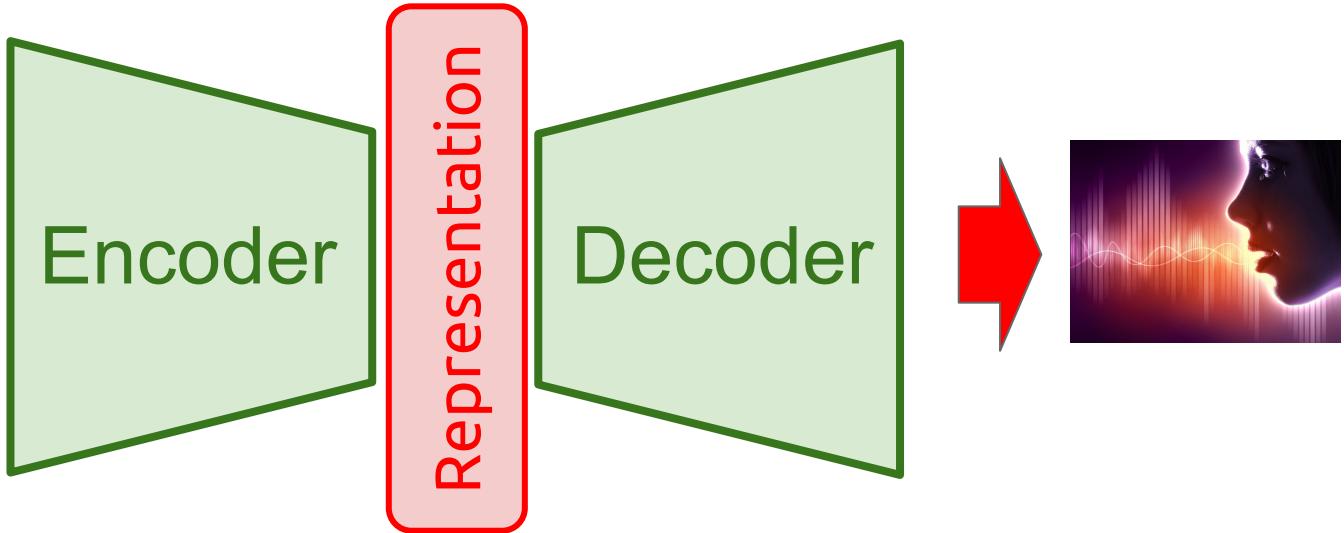
The network gives  $p_t(\alpha | X)$ ,  
a distribution over the outputs  
 $\{h, e, l, o, \epsilon\}$  for each input step.

con-  
eis-  
er et  
nve-  
mico  
mo-  
du-  
re eu  
nspa  
labo-  
ulip-  
cidic  
a. Ut  
exer-  
cito-  
fotor  
e cil-  
cep-  
solit  
n sit

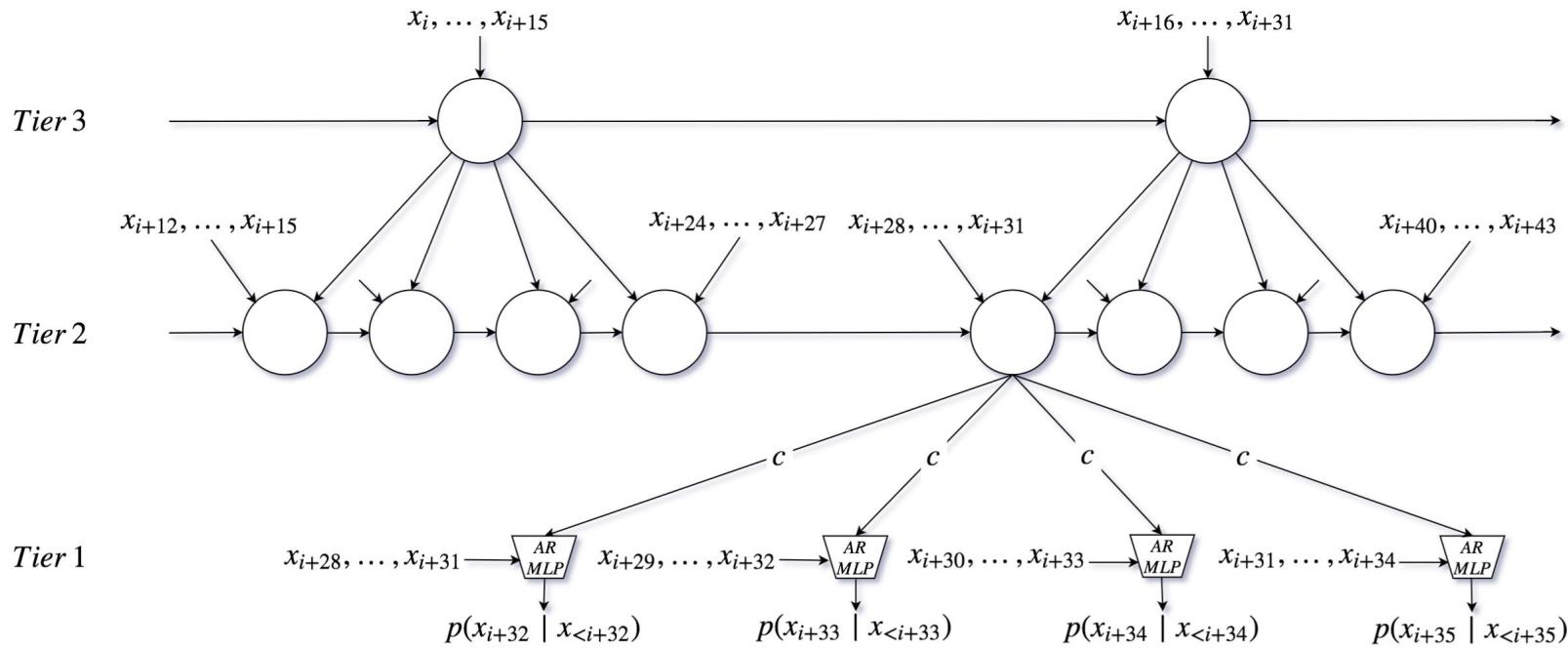
lopi  
ost  
nifb  
cou  
was  
tua  
whi  
wate  
the  
an  
in c  
est  
ses  
ma  
na  
nos  
hi  
auts  
lupt  
null  
tem  
na i

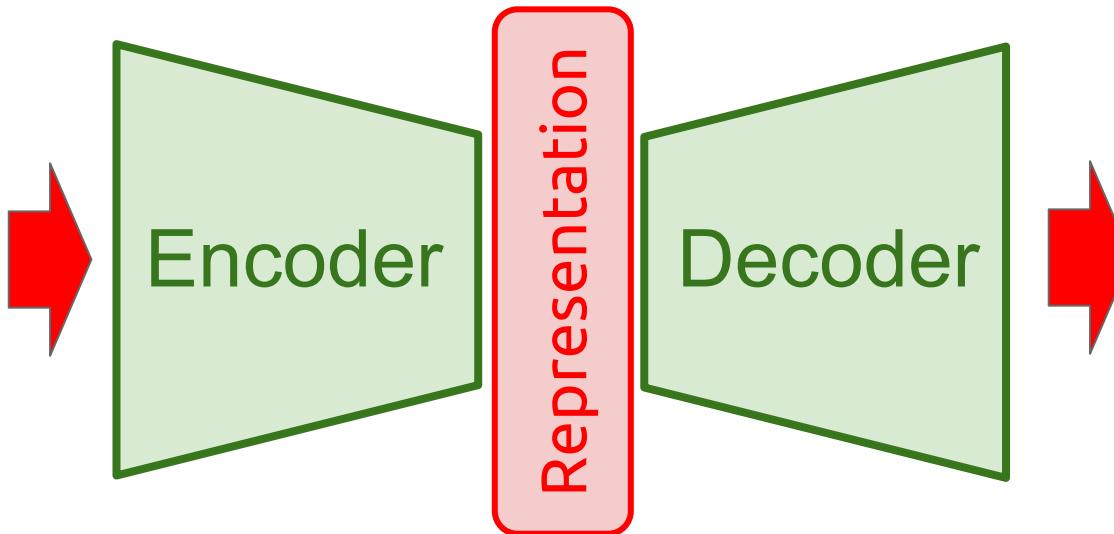
etiam  
tempor  
incididunt  
ut labore et  
dolore magna aliqua. Ut enim ad minim ve-  
niam, quis nostrud exercitation ullamco  
labore ut aliquip ex ea commodo conser-  
vatur. Duis aute irure dolor in reprehenderit  
in voluptate velit esse cillum dolore eu  
fugiat nulla pariatur. Excepteur sint occa-  
si cupidatat non proident, sunt in culpa  
qui officia deserunt mollit anim id est laborum.  
Lorem ipsum sit amet, consectetur adipi-  
scing elit, sed do eiusmod tempor inci-  
didunt ut labore et dolore magna aliqua.  
Ut enim ad minim veniam, quis nostrud exer-  
citation ullamco laboris nisi ut aliquip ex  
ea commodo consequat. Duis aute irure

**Lorem ipsum dolor  
sit amet, consectetur adipiscing elit,  
sed do eiusmod tempor incididunt**



# Speech Synthesis

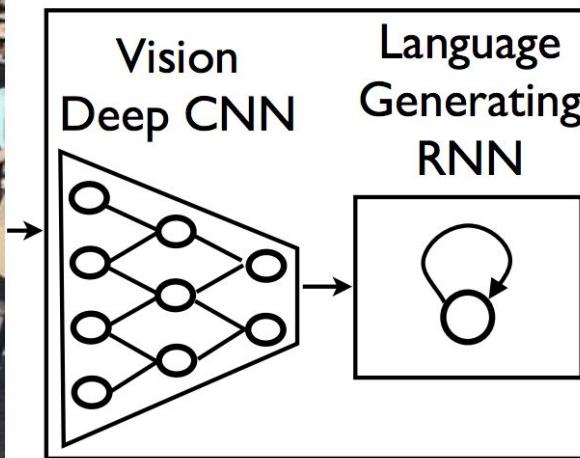




# **Lorem ipsum dolor**

it amet, consectetur adipiscing elit,  
sed do eiusmod tempor incididunt

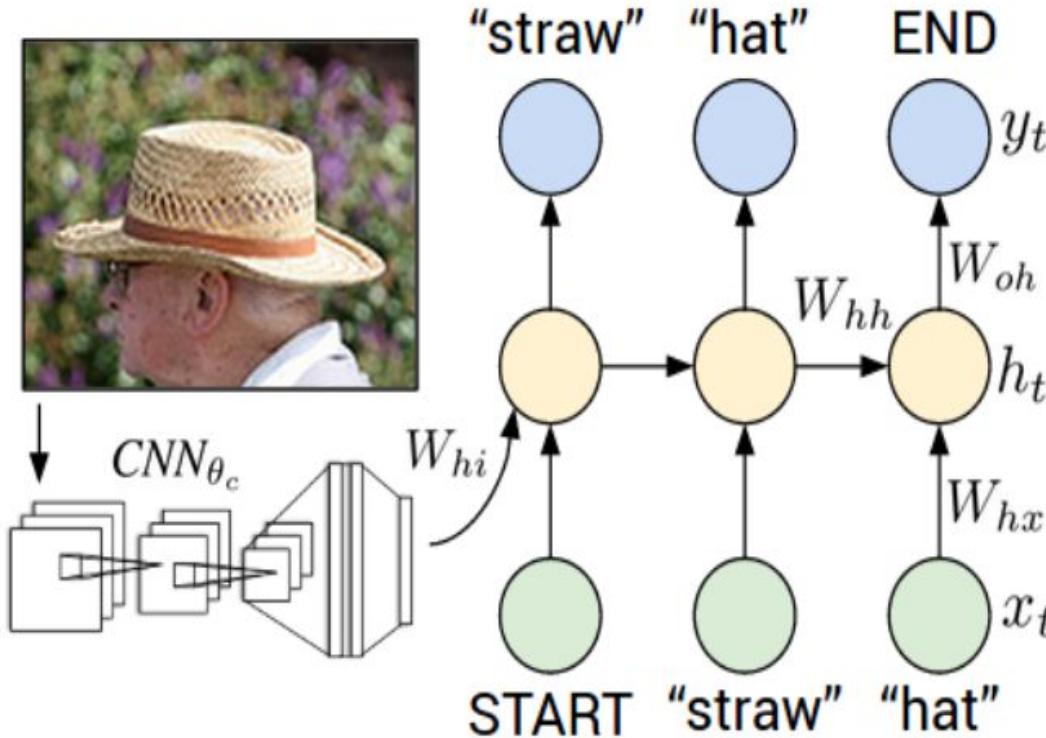
# Image Captioning with RNN

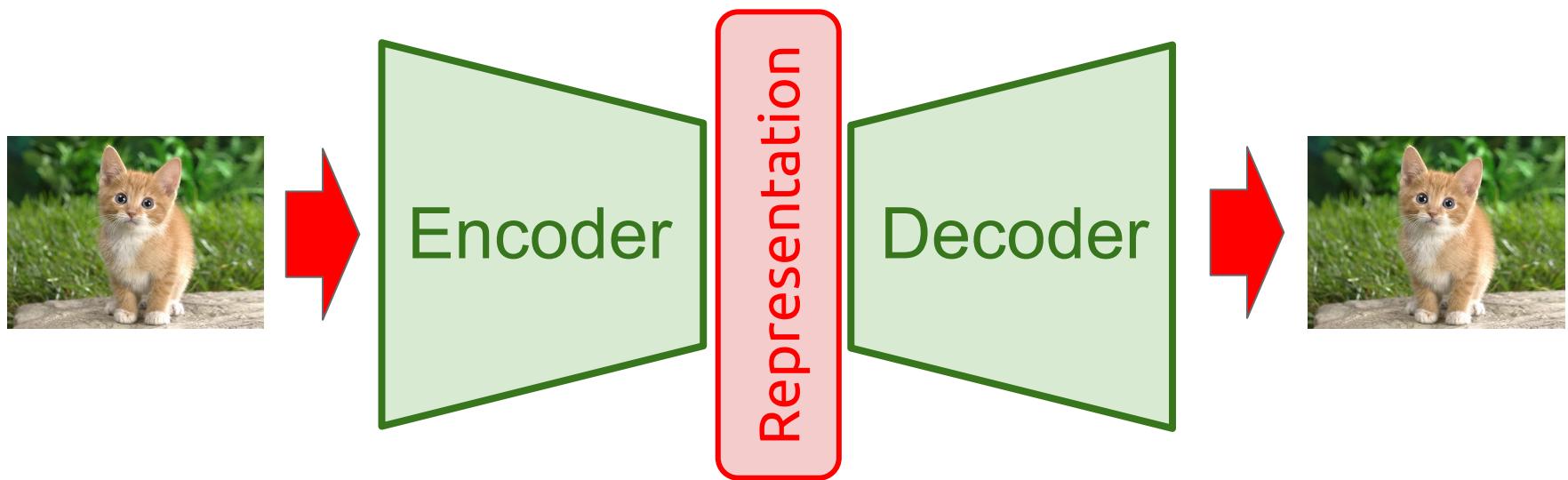


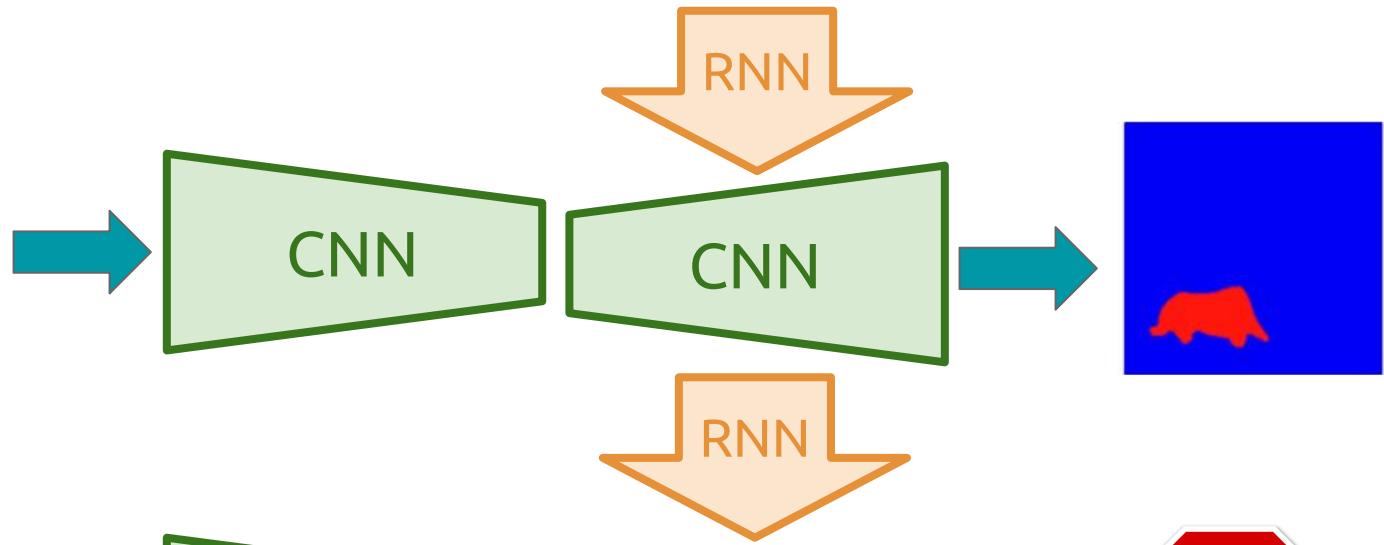
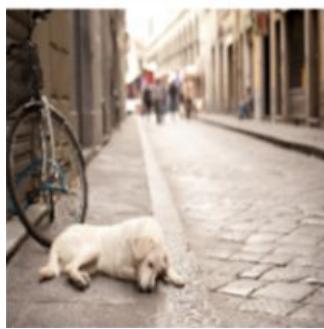
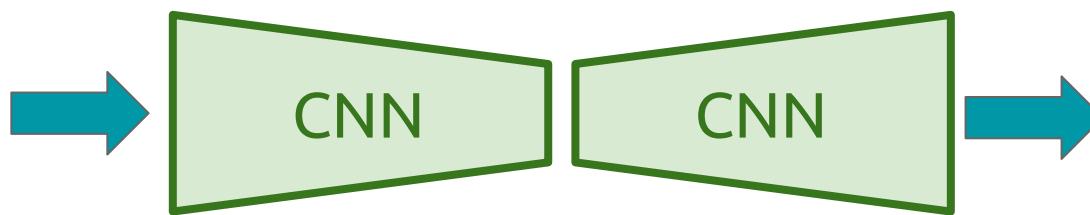
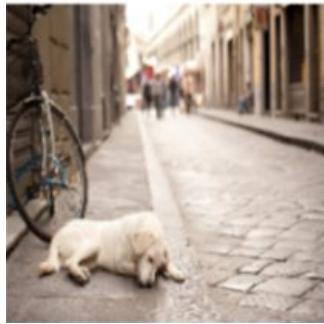
**A group of people shopping at an outdoor market.**

**There are many vegetables at the fruit stand.**

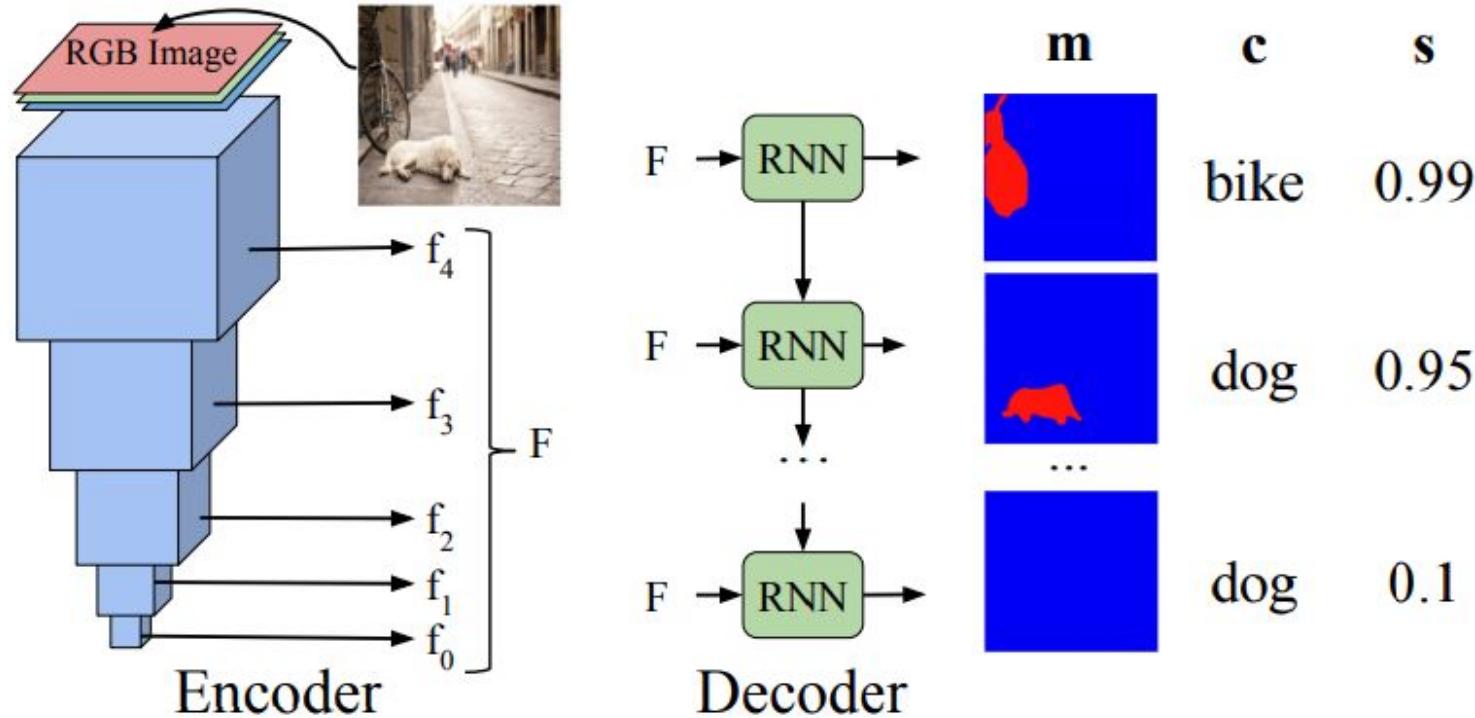
# Image Captioning with RNN



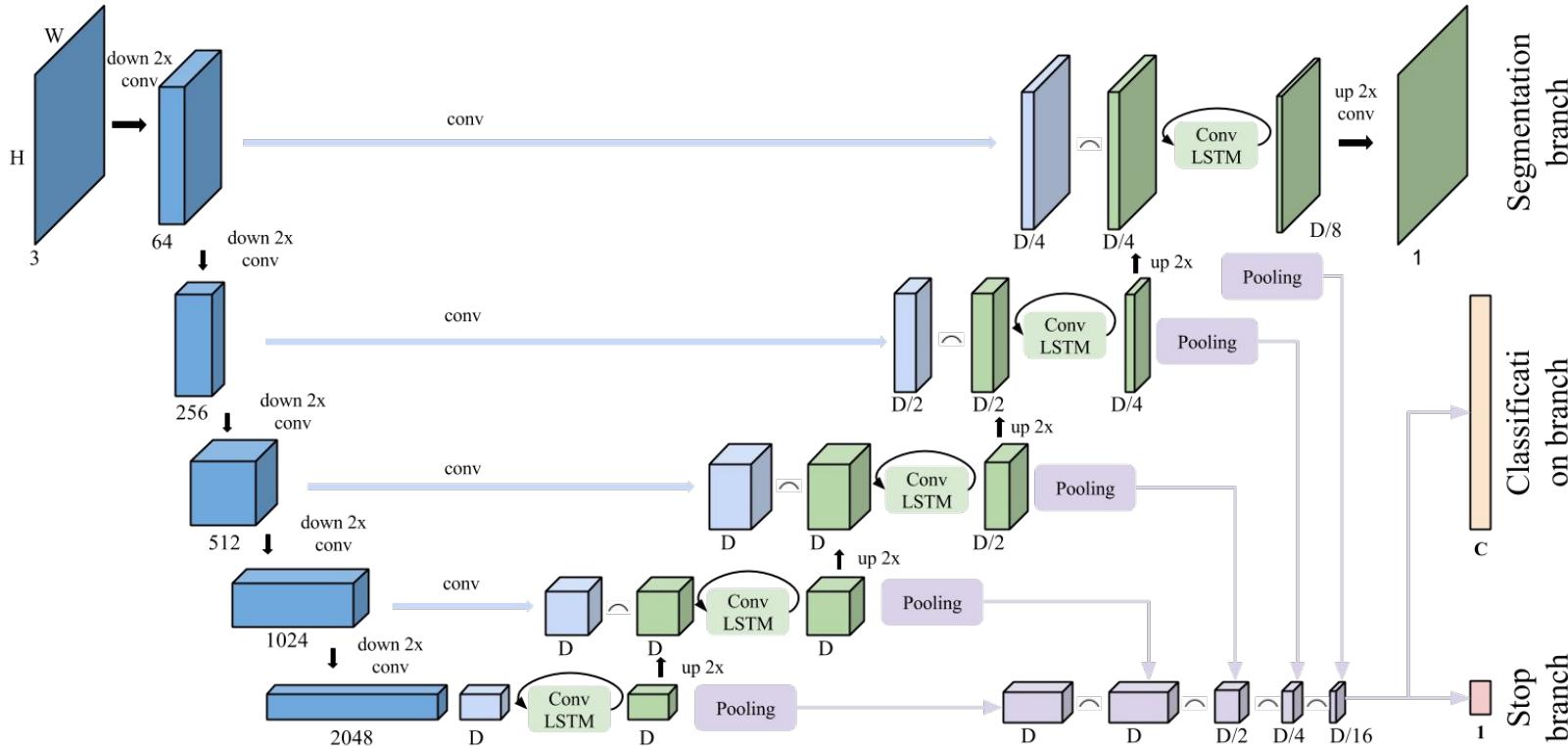




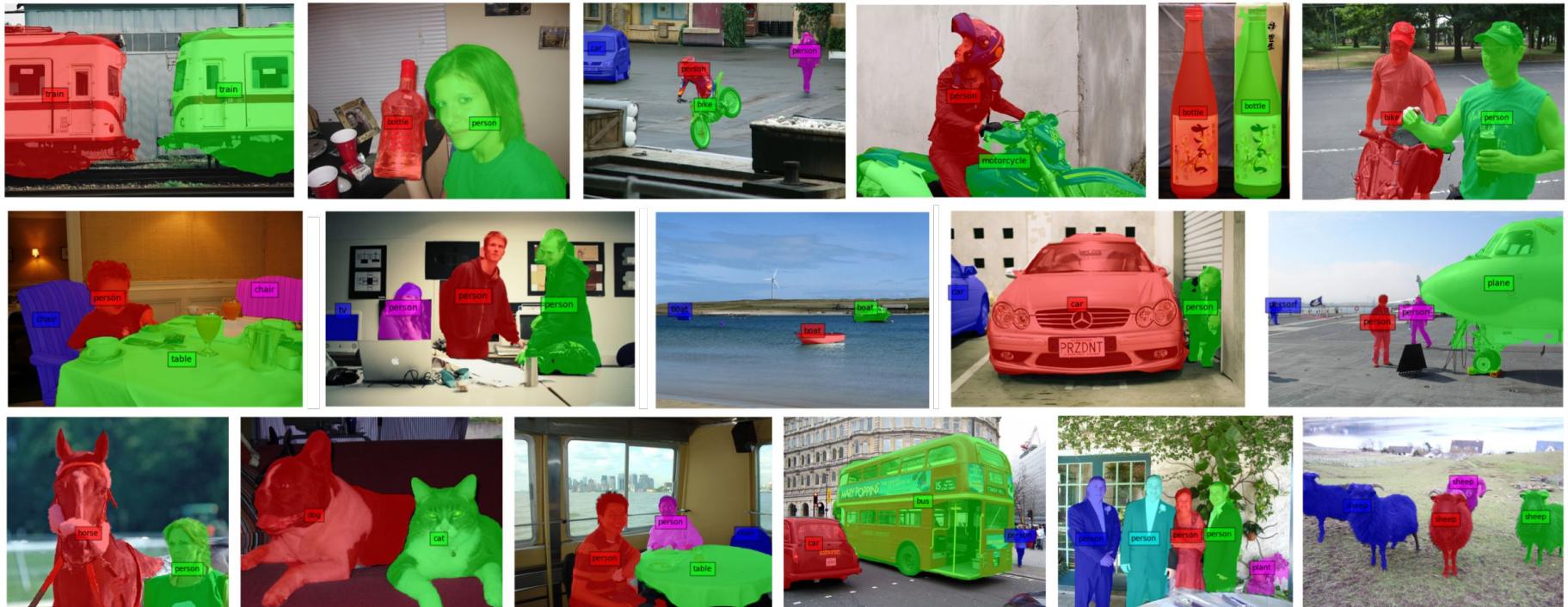
# Recurrent Instance Segmentation



# Recurrent Instance Segmentation

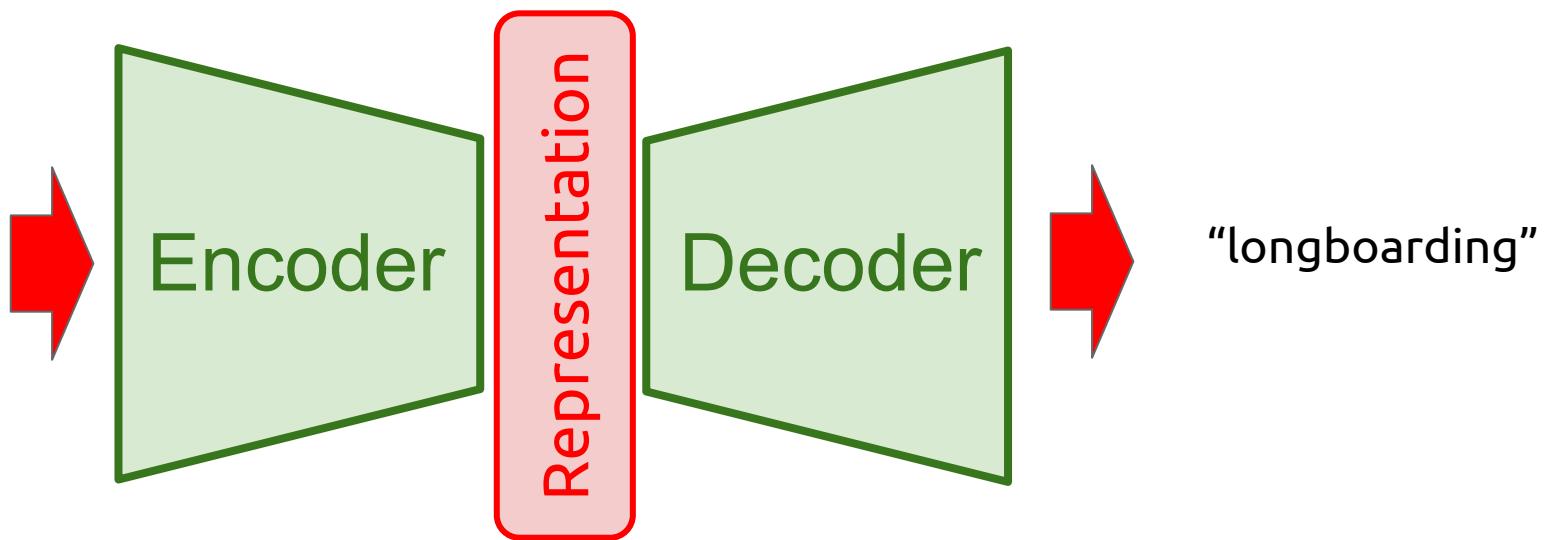


# Recurrent Instance Segmentation

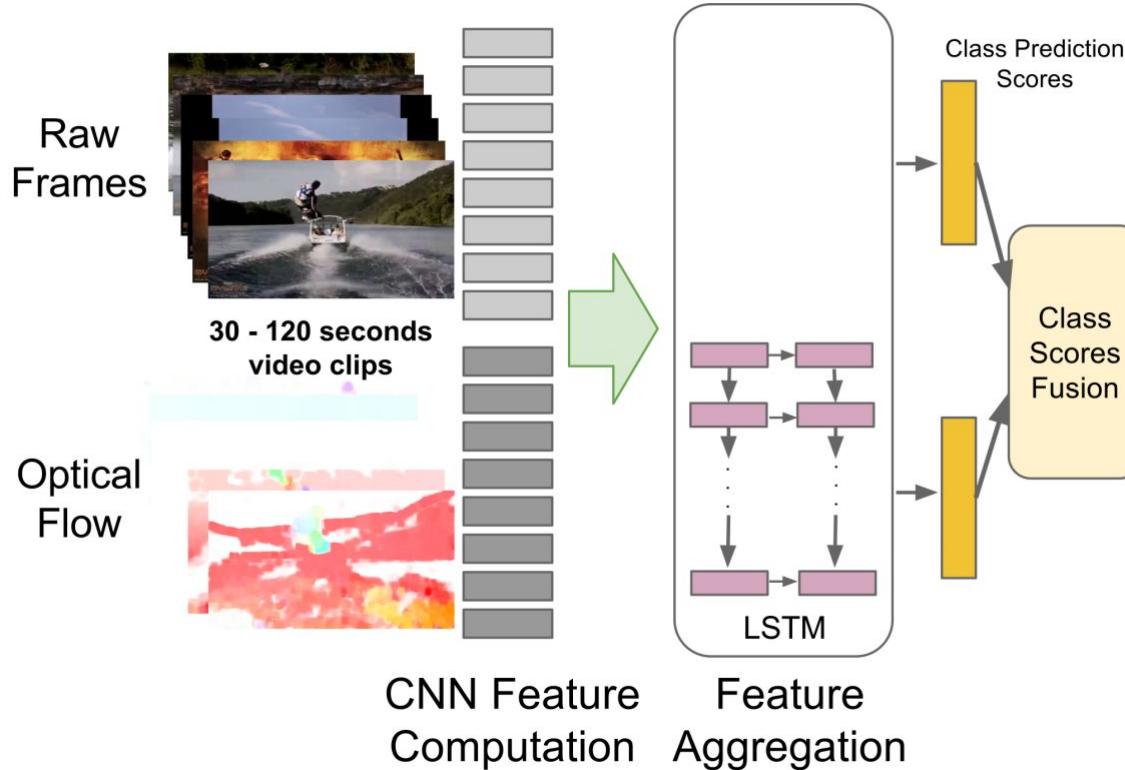


Color sequence:

#RSIS Salvador, Amaia, Miriam Bellver, Victor Campos, Manel Baradad, Ferran Marques, Jordi Torres, Xavier Giro-i-Nieto. ["Recurrent neural networks for semantic instance segmentation."](#) arXiv preprint arXiv:1712.00617 (2017). 80

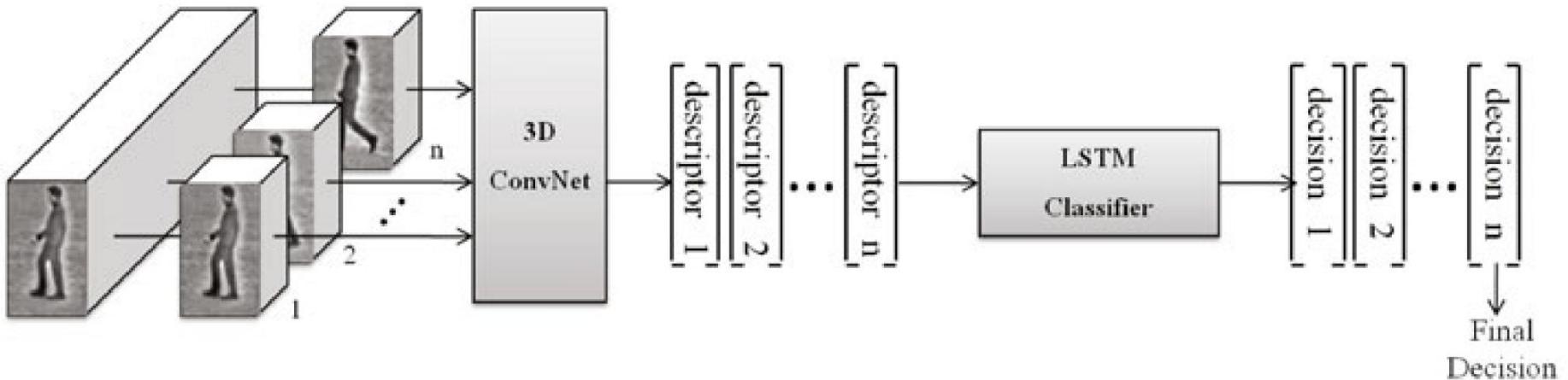


# Two-streams 2D CNNs + RNN



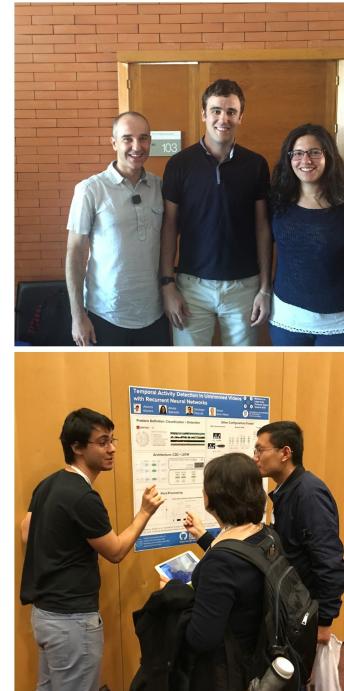
Yue-Hei Ng, Joe, Matthew Hausknecht, Sudheendra Vijayanarasimhan, Oriol Vinyals, Rajat Monga, and George Toderici.  
"Beyond short snippets: Deep networks for video classification." CVPR 2015

# 3D CNN + RNN



Baccouche, Moez, Franck Mamalet, Christian Wolf, Christophe Garcia, and Atilla Baskurt. ["Sequential deep learning for human action recognition."](#) HBU 2011.

# 3D CNN + RNN

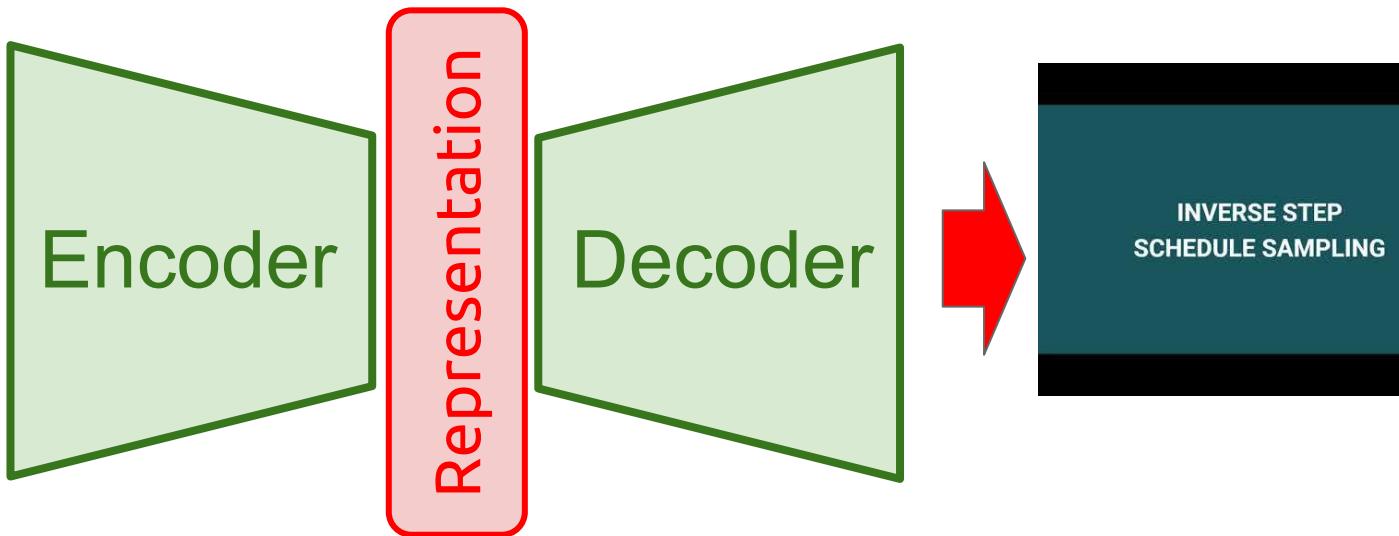


[A. Montes](#), Salvador, A., Pascual-deLaPuente, S., and Giró-i-Nieto, X., "Temporal Activity Detection in Untrimmed Videos with Recurrent Neural Networks", NIPS Workshop 2016 (best poster award)

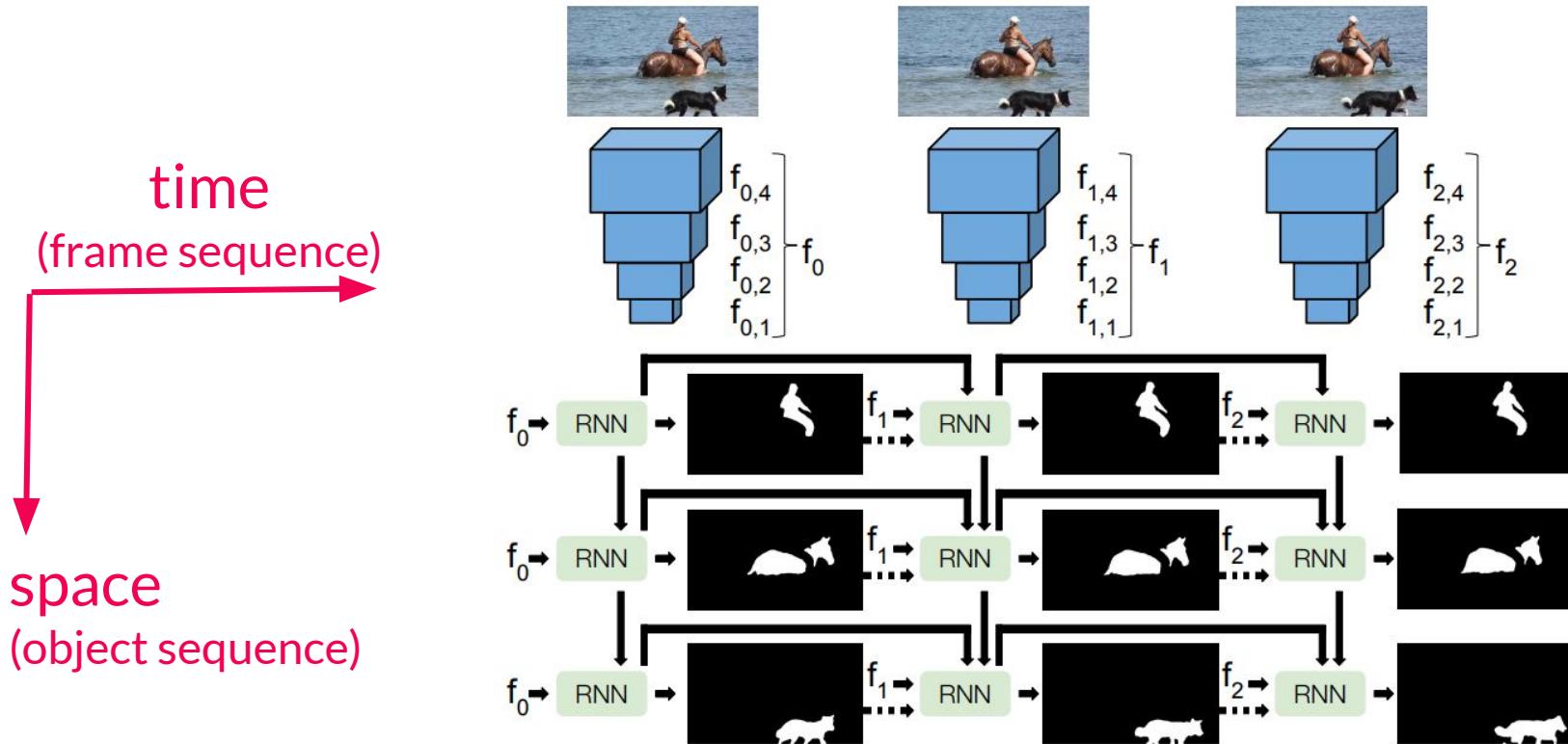


Timex 0.8689

0.8689 Longboarding



# RNN (Spatial + Temporal)

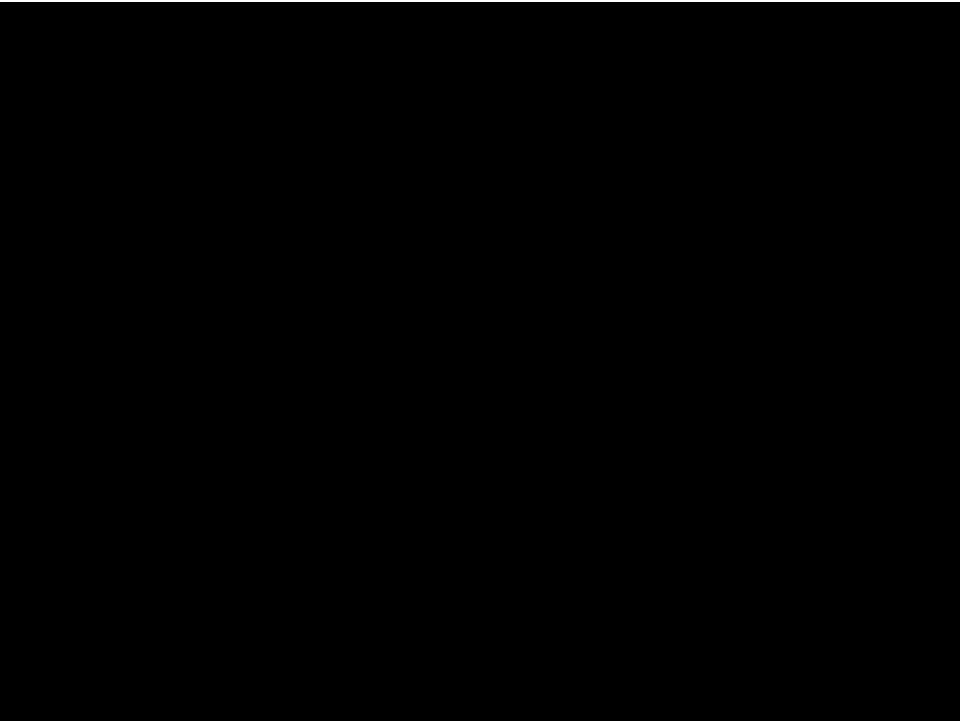


#RVOS Carles Ventura, Miriam Bellver, Andreu Girbau, Amaia Salvador, Ferran Marques and Xavier Giro-i-Nieto. ["RVOS: End-to-End Recurrent Network for Video Object Segmentation"](#), CVPR 2019.

# Outline

1. Motivation
2. Recurrent Neural Layer
3. Forward & Backward Passes Through Time
4. Gated Recurrent Neurons
  - a. LSTM
  - b. GRU
5. Advanced Architectures
6. Encoder-Decoder Architectures

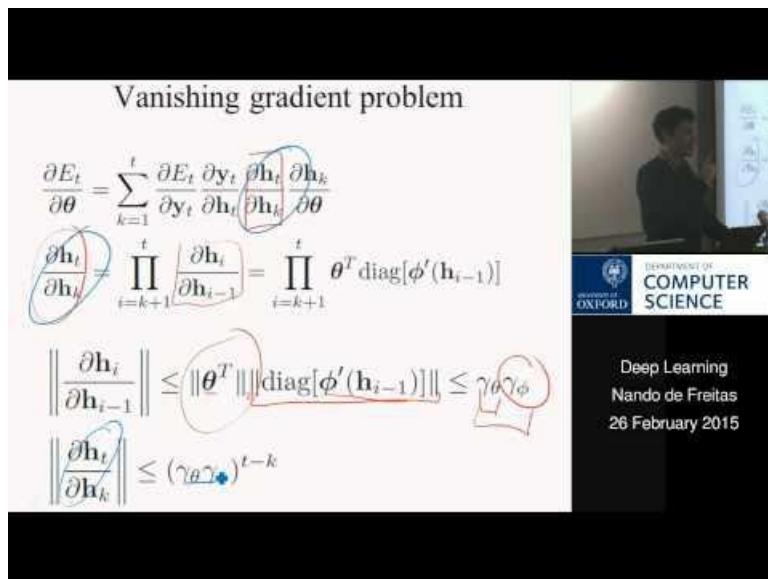
# Teaser of Kostas Derpanis lecture



[Kostas Derpanis, York University](#) (2022) - [@CSProfKGD](#)

# Learn more

- Chris Olah, Shan Carter, ["Attention and Augmented Recurrent Neural Networks"](#). distill.pub 2016.
- Jordi Pons (UPF): [Slides on Recurrent Neural Networks \[tweet\]](#)
- Ian Goodfellow, Yoshua Bengio, Aaron Courville, ["Deep Learning \(Chapter 10\)"](#). MIT Press 2016.
- Alex Graves, ["Supervised Sequence Labelling with Recurrent Neural Networks"](#)



Nando de Freitas, ["Recurrent Neural Nets and LSTMs"](#) (University of Oxford 2015).



Marta Garnelo, [Deepmind x UCL 2020 \[slides\]](#)

# (extra) PyTorch Lab on Google Colab

dai\_2019\_lab08\_rnn\_todo.ipynb

File Edit View Insert Runtime Tools Help Last edited on Dec 10, 2019

+ Code + Text

## Recurrent Neural Networks

### Lab credit

Lab created by [Santiago Pascual](#) and [Xavier Giro-i Nieto](#) for the [Postgraduate course in artificial intelligence with deep learning](#) in [UPC School \(2019\)](#).

### Slides

Related slides by [Marta R. Costa-jussà](#) from [Deep Learning for Artificial Intelligence](#) (UPC TelecomBCN 2018), and slides from [Xavier Giro-i Nieto](#) from [Deep Learning for Multimedia](#) (Dublin City University 2018).

### Video lectures

Related [video \(basic\)](#) & [video \(advanced\)](#) by [Santiago Pascual](#) from [Deep Learning for Speech and Language](#) (UPC TelecomBCN 2017), and [video](#) by [Marta R. Costa-jussà](#) from [Deep Learning for Artificial Intelligence](#) (UPC TelecomBCN 2019)



Santiago Pascual de la Puente  
santi.pascual@upc.edu

PhD 2019

Universitat Politècnica de Catalunya  
Technical University of Catalonia



# DEEP LEARNING FOR ARTIFICIAL INTELLIGENCE

3rd Master Course UPC ETSETB TelecomBCN Barcelona. Autumn 2019.

## Instructors

Xavier Giro-i Nieto    Marta R. Costa-jussà    Noé Casas    Verónica Vilaplana    Ramon Morros    Javier Ruiz    Albert Pumarola    Jordi Torres

## Organizers

UNIVERSITAT POLITÈCNICA  
DE CATALUNYA  
BARCELONATECH

telecom  
BCN

## Supporters

Google Cloud

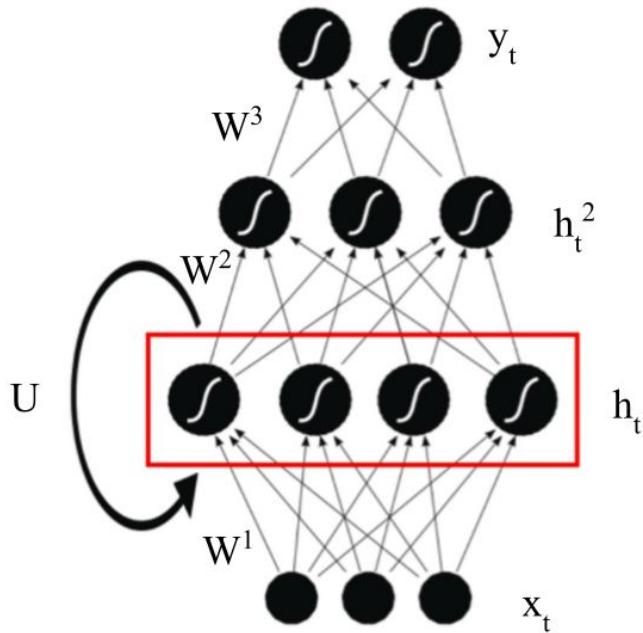
GitHub Education

+ info: <http://bit.ly/dlai2019>

- [Lectures](#) (with Slides & Videos)
- [Labs](#)

# Exercise

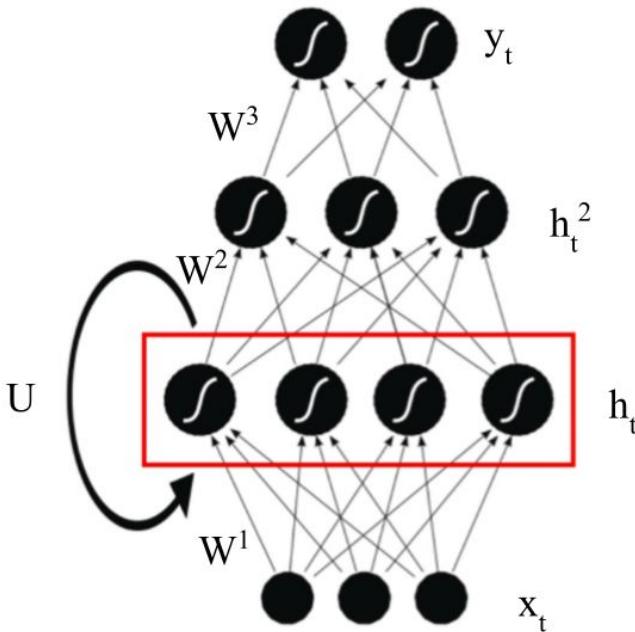
Consider a the neural network depicted in the figure, that represents the processing of a data sequence at timestep t. The layer in the box is a recurrent one, with a set of weights U.



- Formulate how is the output of the recurrent layer,  $h_t^1$ , computed considering a  $\tanh(\cdot)$  activation.
- How many parameters does this architecture contain ? Together with the depicted weight matrices  $W^i$ , also consider the biases  $b^i$ . Develop the calculations you made to reach the provided results.
- Draw the unfolded representation of this architecture for timesteps  $t-1$ ,  $t$  and  $t+1$ , labeling the edges with the appropriate weight matrices.

# Solution

Consider a the neural network depicted in the figure, that represents the processing of a data sequence at timestep t. The layer in the box is a recurrent one, with a set of weights U.



- a) Formulate how is the output of the recurrent layer,  $h_t^1$ , computed considering a  $\tanh(\cdot)$  activation.

$$h_t^1 = \tanh(U h_{t-1}^1 + W^1 x_t + b^1)$$

- b) How many parameters does this architecture contain ? Together with the depicted weight matrices  $W^i$ , also consider the biases  $b^i$ . Develop the calculations you made to reach the provided results.

$$\text{Layer } \#1: 3 \times 4 + 4 + 4 \times 4 = 32$$

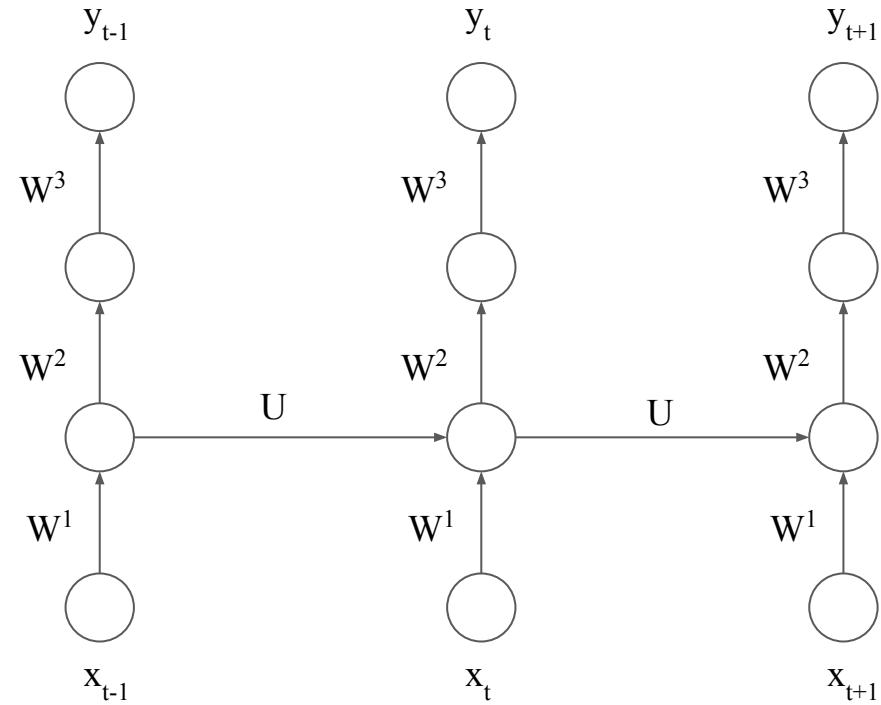
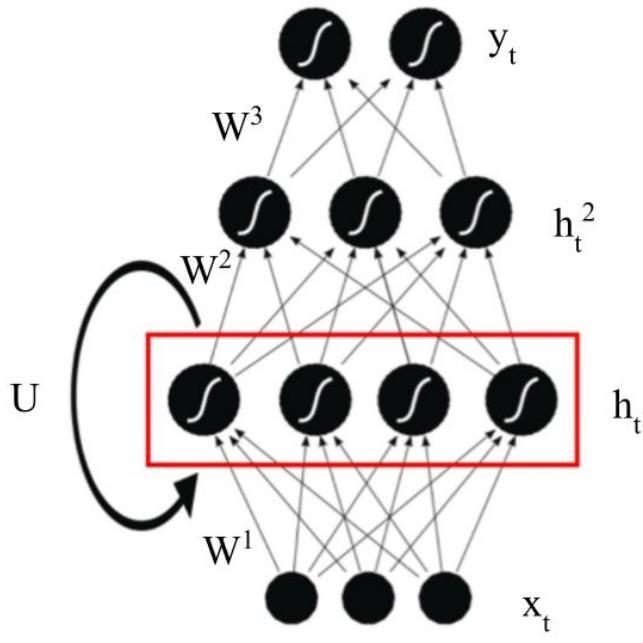
$$\text{Layer } \#2: 4 \times 3 + 3 = 15$$

$$\text{Layer } \#3: 3 \times 2 + 2 = 8$$

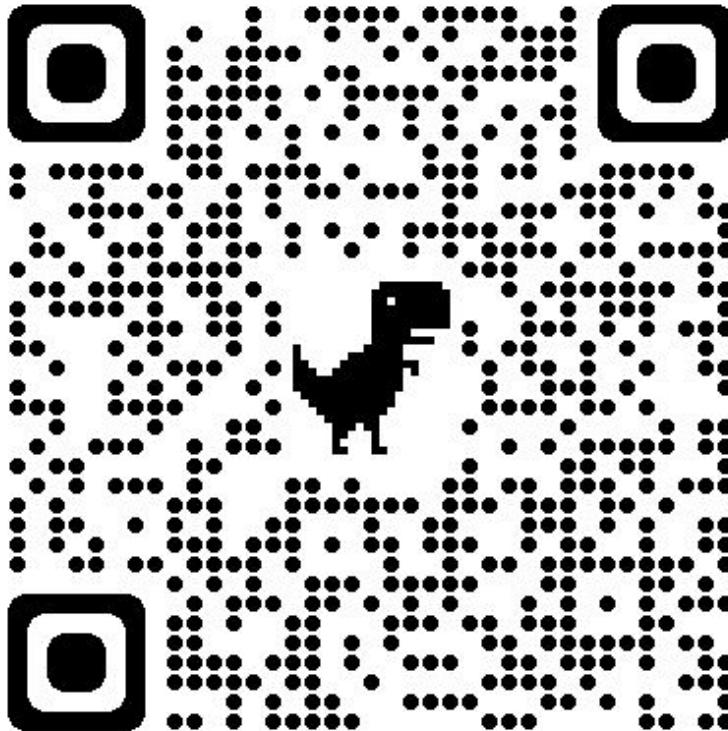
$$\mathbf{TOTAL: 20 + 15 + 8 = 55}$$

# Solution

c) Draw the unfolded representation of this architecture for timestep t-1, t and t+1, labeling the edges with the appropriate weight matrices.



# Quiz



```
import math
```

Define the data we need to solve the problem.

```
[ ] x0=12
Wxh=-0.1
Whh=0.5
hbias=0.1
Why=0.7
ybias=0.0
```

Compute the output of the hidden layer at the first time step, in which there is no previous hidden state.

```
[ ] dot = Wxh*x0 + hbias
h0 = 1.0 / (1 + math.exp(-dot) )
print(h0)
```

0.24973989440488234

```
[ ] yl=Why*h0+ybias
print( yl )
```

0.17481792608341762

```
[ ] xl=8
dot = xl * Wxh + h0*Whh + hbias
h1 = 1.0 / (1 + math.exp(-dot) )
print(h1)
```

0.36005393658863494

# Questions ?

## Undergradese

What undergrads ask vs. what they're REALLY asking

"Is it going to be an open book exam?"

Translation: "I don't have to actually memorize anything, do I?"

"Hmm, what do you mean by that?"

Translation: "What's the answer so we can all go home."

"Are you going to have office hours today?"

Translation: "Can I do my homework in your office?"

"Can i get an extension?"

Translation: "Can you re-arrange your life around mine?"

"Is this going to be on the test?"

Translation: "Tell us what's going to be on the test."

"Is grading going to be curved?"

Translation: "Can I do a mediocre job and still get an A?"

