# Deep-Learning:

# Recurrent Neural Networks (RNN)

**Pr. Fabien MOUTARDE**
**Center for Robotics**
**MINES ParisTech**
**PSL Université Paris**

`Fabien.Moutarde@mines-paristech.fr`
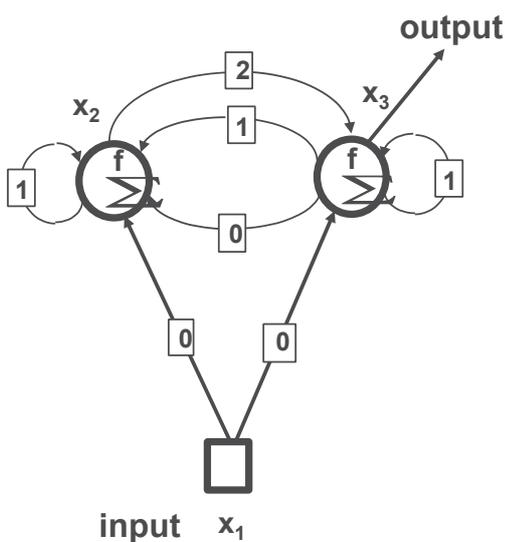
`http://people.mines-paristech.fr/fabien.moutarde`

# Acknowledgements

**During preparation of these slides, I got inspiration and borrowed some slide content from several sources, in particular:**
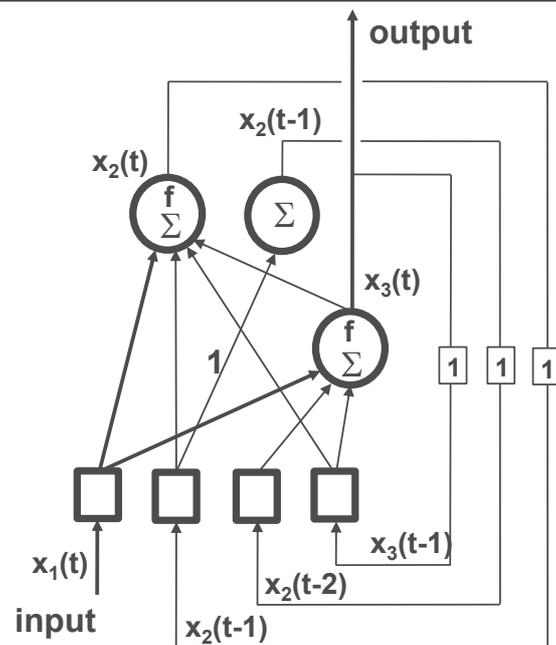
- **Fei-Fei Li + J.Johnson + S.Yeung: slides on "*Recurrent Neural Networks*" from the** "*Convolutional Neural Networks for Visual Recognition*" **course at Stanford**
http://cs231n.stanford.edu/slides/2019/cs231n_2019_lecture10.pdf

- **Yingyu Liang: slides on "*Recurrent Neural Networks*" from the** "*Deep Learning Basics*" **course at Princeton**
https://www.cs.princeton.edu/courses/archive/spring16/cos495/slides/DL_lecture9_RNN.pdf

- **Arun Mallya: slides "*Introduction to RNNs*" from the** "Trends in Deep Learning and Recognition" **course of Svetlana LAZEBNIK at University of Illinois at Urbana-Champaign**
http://slazebni.cs.illinois.edu/spring17/lec02_rnn.pdf

- **Tingwu Wang: slides on "*Recurrent Neural Network*" for a course at University of Toronto**
https://www.cs.toronto.edu/%7Etingwuwang/rnn_tutorial.pdf

- **Christopher Olah: online tutorial "*Understanding LSTM Networks*"**
*https://colah.github.io/posts/2015-08-Understanding-LSTMs/*

# Outline

- **Standard Recurrent Neural Networks**
- Training RNN: BackPropagation Through Time
- LSTM and GRU
- Applications of RNNs

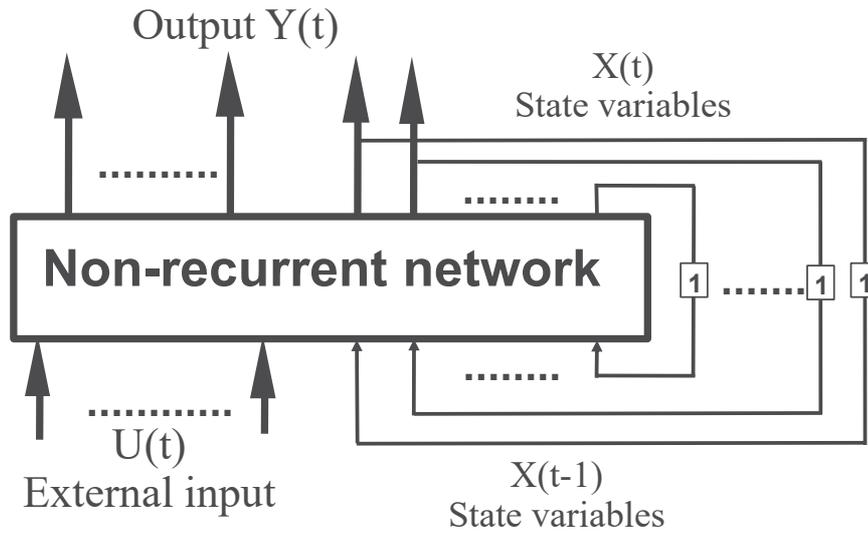# <u>R</u>ecurrent <u>N</u>eural <u>N</u>etworks (RNN)
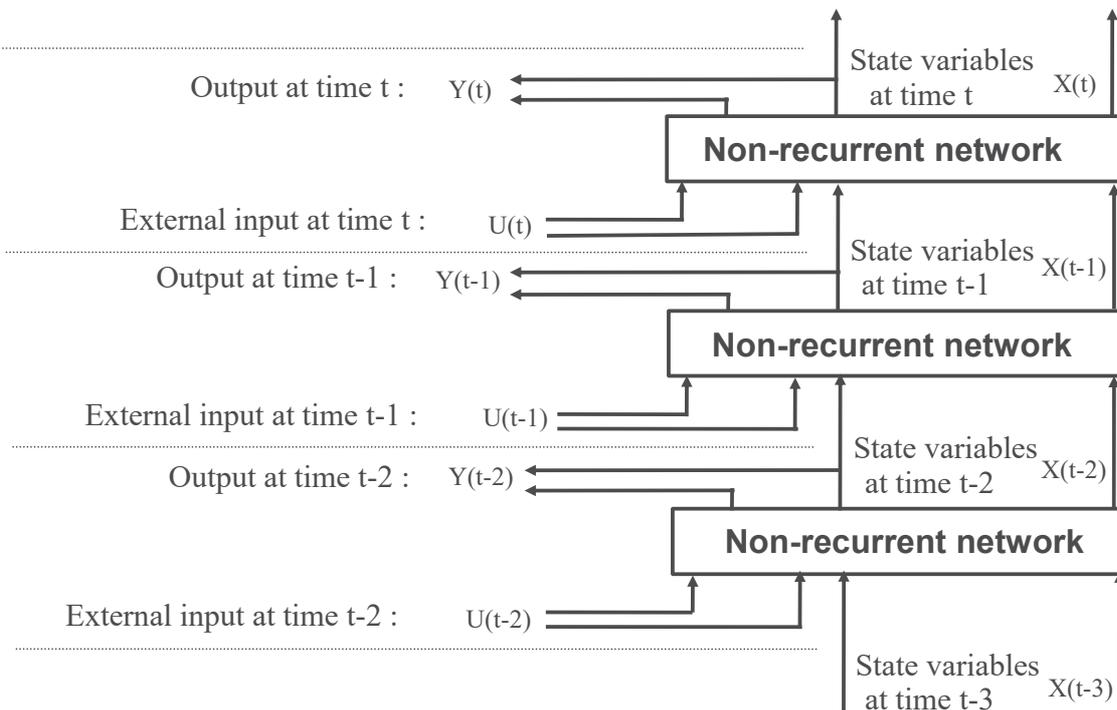


**Time-delay for each connection**
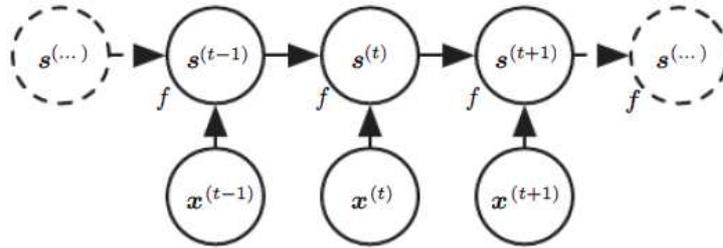
**Equivalent form**

Output Y(t)

X(t)
State variables

**Non-recurrent network**

1 ....... 1 1

U(t)
External input

X(t-1)
State variables

Output at time t :    Y(t)

State variables
at time t    X(t)

**Non-recurrent network**

External input at time t :    U(t)

State variables    X(t-1)
at time t-1

Output at time t-1 :    Y(t-1)

**Non-recurrent network**

External input at time t-1 :    U(t-1)

State variables
at time t-2    X(t-2)

Output at time t-2 :    Y(t-2)

**Non-recurrent network**

External input at time t-2 :    U(t-2)

State variables
at time t-3    X(t-3)

$$s^{(t+1)} = f\left(s^{(t)}, x^{(t+1)}\right)$$



**If using a Neural Net for f, this is EXACTLY a RNN!**
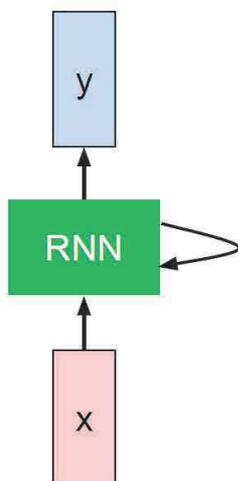


Figures from *Deep Learning*, Goodfellow, Bengio and Courville

---

**State vector s ←→ vector h of hidden neurons**



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

new state          old state  input vector at
                                  some time step

some function
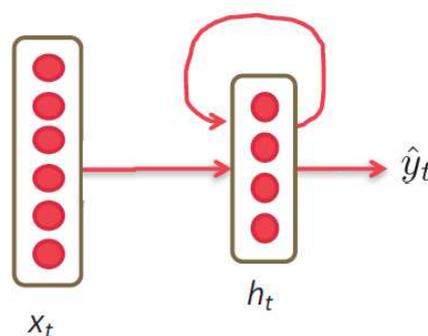with parameters W

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

$$y_t = W_{hy} h_t$$

ou $y_t = \mathtt{softMax}(W_{hy} h_t)$

# Advantages of RNN

The *hidden state* s of the RNN builds a kind of <u>lossy summary of the past</u>

RNN totally <u>adapted to processing SEQUENTIAL data</u> (same computation formula applied at each time step, but modulated by the evolving "memory" contained in state s)
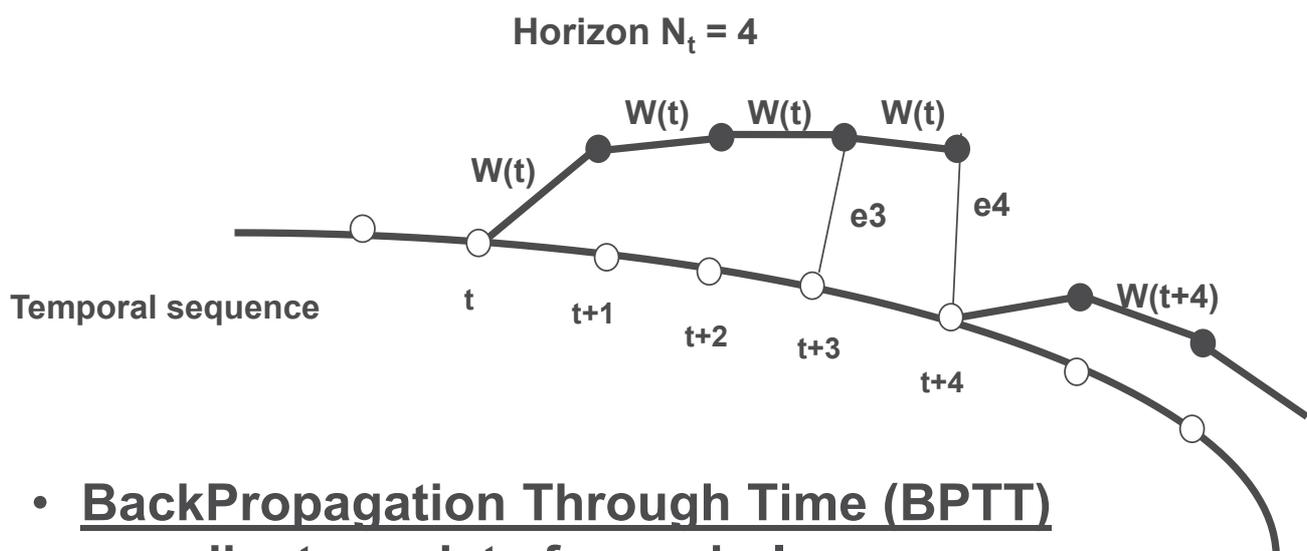
<u>Universality of RNNs</u>: any function computable by a Turing Machine can be computed by a finite-size RNN (Siegelmann and Sontag, 1995)

# RNN hyper-parameters



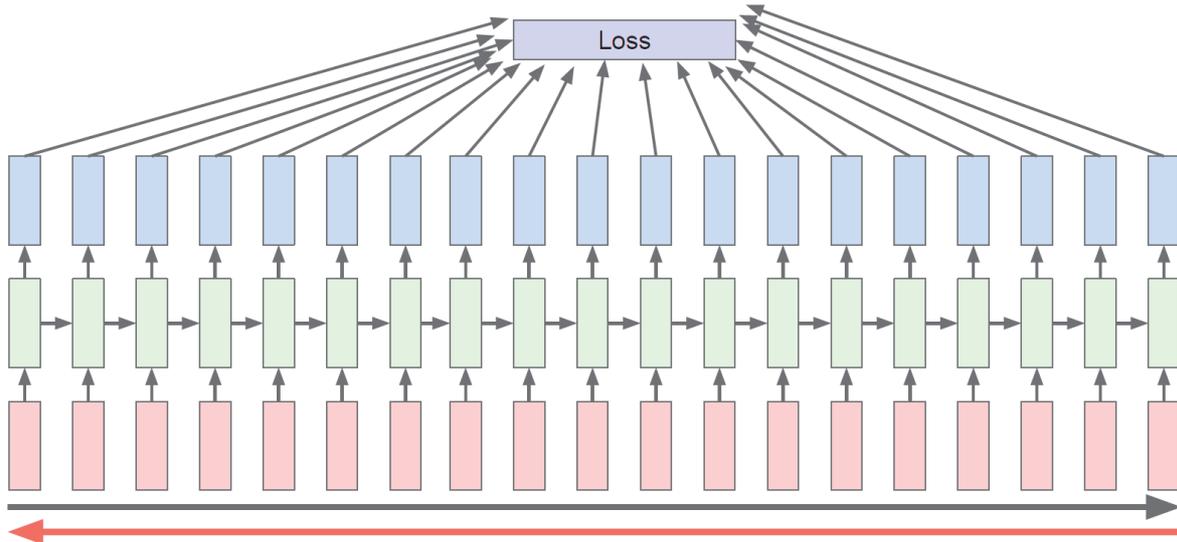- **As for MLP, main hyperparameter = <u>size of hidden layer</u> (=size of vector h)**

- Standard Recurrent Neural Networks
- **Training RNN: BackPropagation Through Time**
- LSTM and GRU
- Applications of RNNs

---

# RNN training
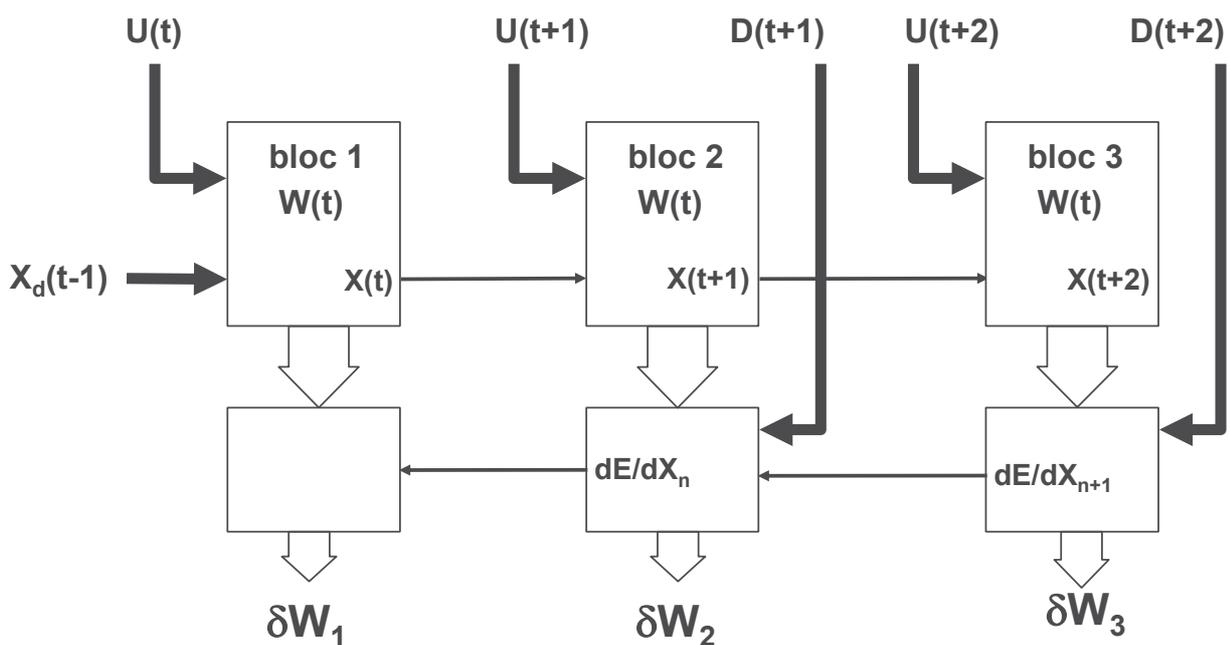
**Horizon $N_t$ = 4**



Temporal sequence

- **<u>BackPropagation Through Time (BPTT)</u>** gradients update for a whole sequence
- **or Real Time Recurrent Learning (RTRL)** gradients update for each frame in a sequence

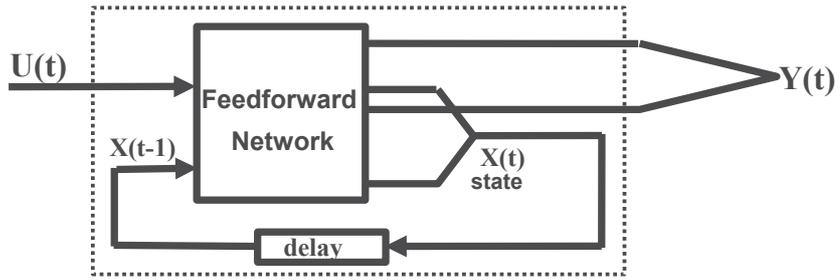# BackPropagation THROUGH TIME (BPTT)



- **Forward through entire sequence to compute SUM of losses at ALL (or part of) time steps**
- **Then backprop through ENTIRE sequence to compute gradients**

# BPTT computation principle



$$\delta W = \delta W_1 + \delta W_2 + \delta W_3$$

# BPTT algorithm

$$W(t+N_t) = W(t) - \lambda \ \text{grad}_W(E) \ \text{avec} \ E = \sum_\tau (Y_\tau - D_\tau)^2$$

$$\frac{\partial E}{\partial W} = \sum_{t=1}^{T} \frac{\partial E_t}{\partial W} \quad \textbf{and} \quad \forall t, \frac{\partial E_t}{\partial W} = \frac{\partial E_t}{\partial Y_t} \frac{\partial Y_t}{\partial X_{t-1}} \frac{\partial X_{t-1}}{\partial W} \quad \textbf{(chain rule)}$$

$$\frac{\partial X_t}{\partial W} = \sum_{k=1}^{t-1} \frac{\partial X_t}{\partial X_{t-k}} \frac{\partial X_{t-k}}{\partial W} \qquad \frac{\partial X_t}{\partial X_{t-k}} = \prod_{j=1}^{t} \boxed{\frac{\partial X_j}{\partial X_{j-1}}} \quad \textbf{\textcolor{purple}{Jacobian matrix of the Feedforward net}}$$

---

# Vanishing/exploding gradient problem

- **If eigenvalues of Jacobian matrix >1, then <u>gradients tend to EXPLODE</u>**
  **➔ Learning will never converge.**

- **Conversely, if eigenvalues of Jacobian matrix <1, then <u>gradients tend to VANISH</u>**
  **➔ Error signals can only affect small time lags**
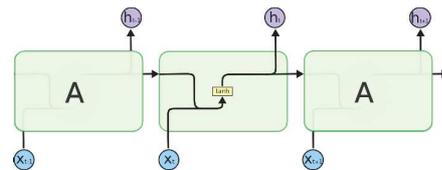  **➔ short-term memory.**

**➔Possible solutions for exploding gradient:**
  **_CLIPPING_ trick**

**➔ Possible solutions for vanishing gradient:**
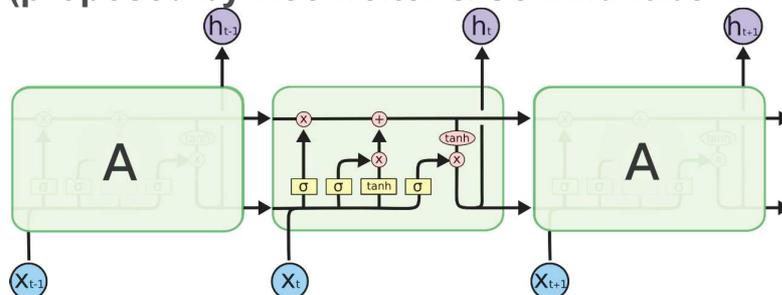  **– use _ReLU_ instead of tanh**
  **– _change what is inside the RNN!_**

# Outline

---

# Long Short-Term Memory (LSTM)

**Problem of *standard* RNNs =
no actual LONG-TERM memory**



## LSTM = RNN variant for solving this issue
**(proposed by Hochreiter & Schmidhuber in 1997)**



*[Figures from https://colah.github.io/posts/2015-08-Understanding-LSTMs/]*

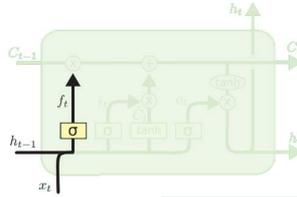- **Key idea = use "gates" that modulate respective influences of input and memory**

# LSTM gates

Gate = *pointwise* multiplication by $\sigma$ in $]0;1[$
➔ modulate between "*let nothing through*"
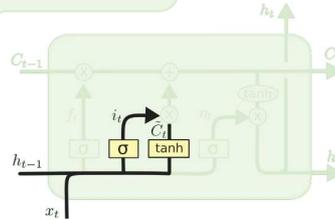and "*let everything through*"

- **FORGET gate**

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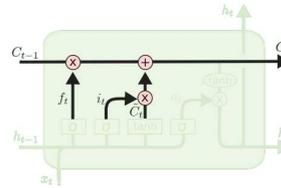
  - **INPUT gate**

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

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

➔ **next state = mix between
pure memory or pure new**
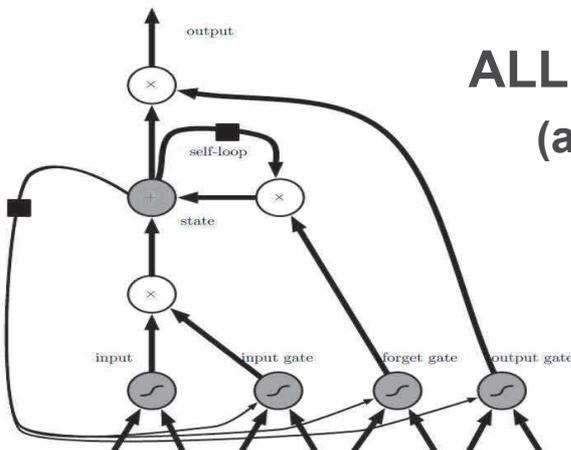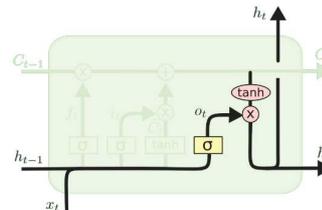
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

*[Figures from https://colah.github.io/posts/2015-08-Understanding-LSTMs/]*

---

# LSTM summary

- **OUTPUT gate**

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

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

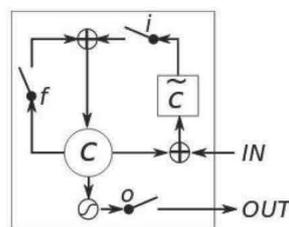**ALL weigths $\mathtt{W_f}$, $\mathtt{W_i}$, $\mathtt{W_c}$ and $\mathtt{W_o}$
(and biases) are LEARNT**

*[Figure from Deep Learning book by I. Goodfellow, Y. Bengio & A. Courville]*
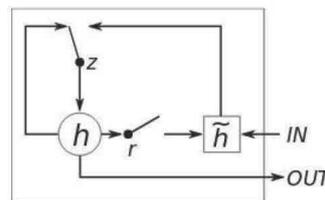
# Why LSTM avoids vanishing gradients?

## Uninterrupted gradient flow!

# Gated Recurrent Unit (GRU)

## Simplified variant of LSTM, with only 2 gates:
## a RESET gate & an UPDATE gate
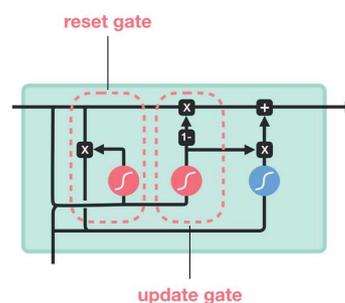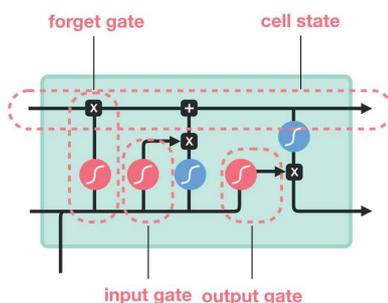### (proposed by Cho, et al. in 2014)
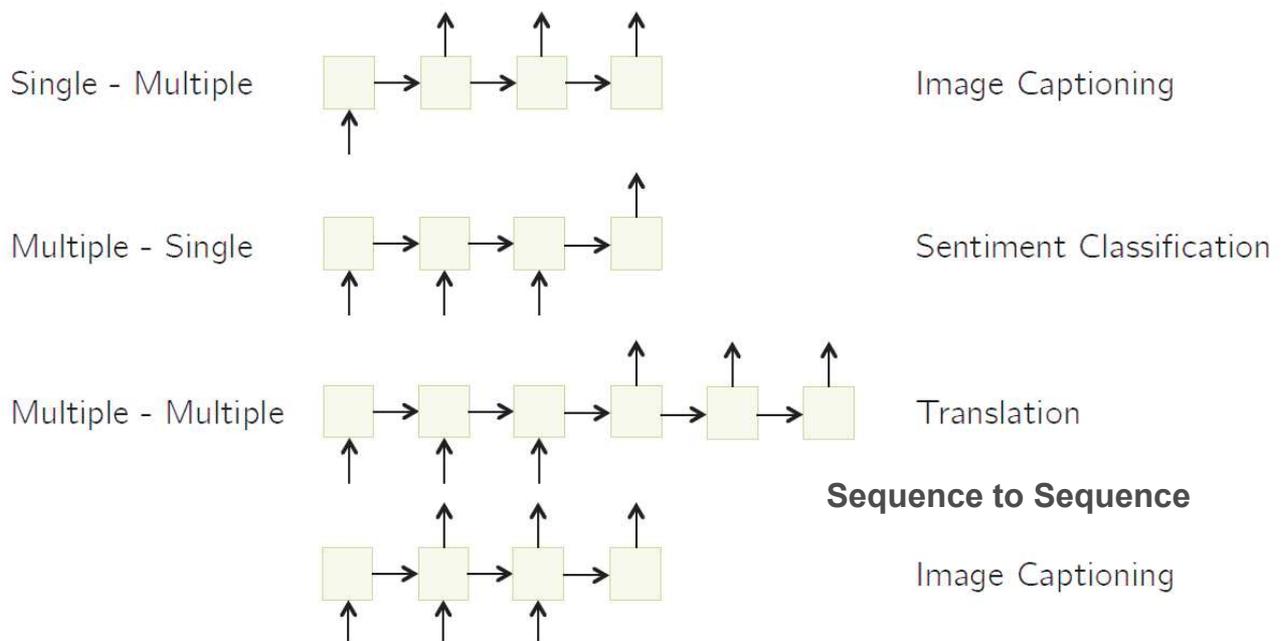


(a) Long Short-Term Memory            (b) Gated Recurrent Unit

LSTM                                   GRU

forget gate          cell state        reset gate

input gate   output gate               update gate
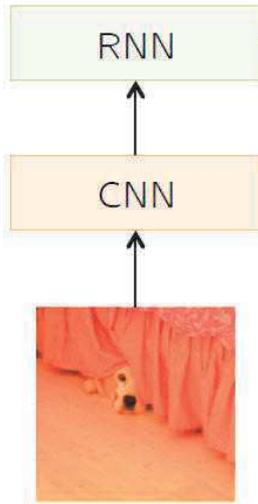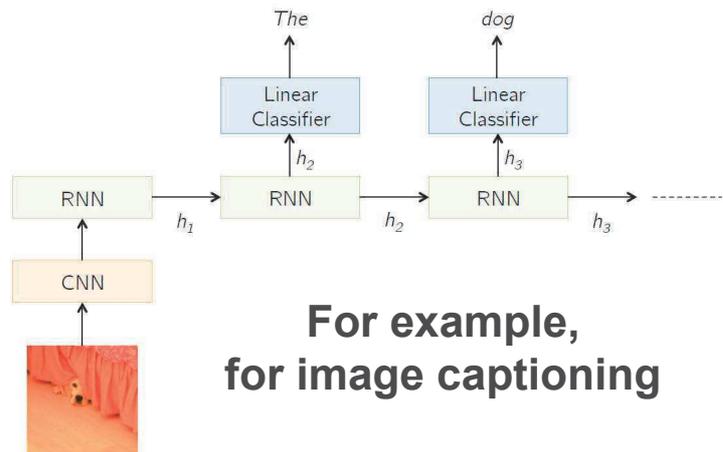
- Standard Recurrent Neural Networks
- Training RNN: BackPropagation Through Time
- LSTM and GRU
- **Applications of RNNs**
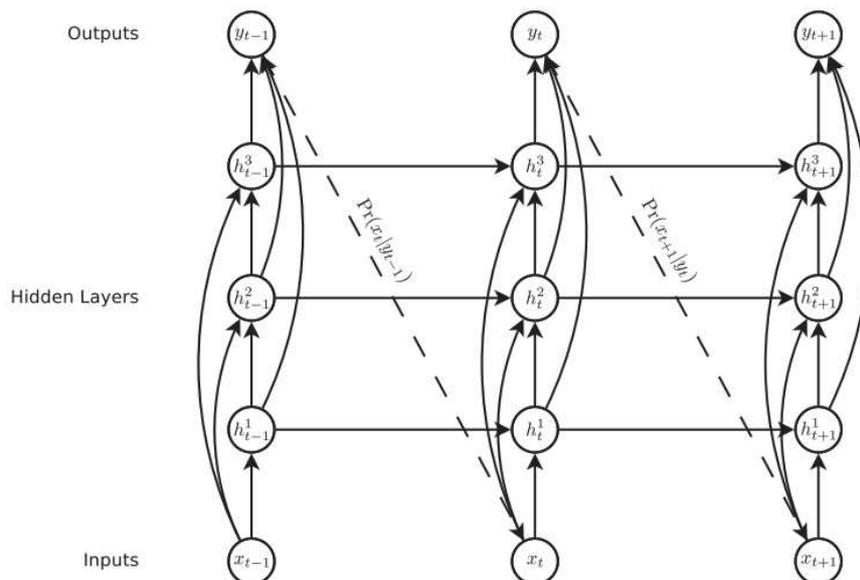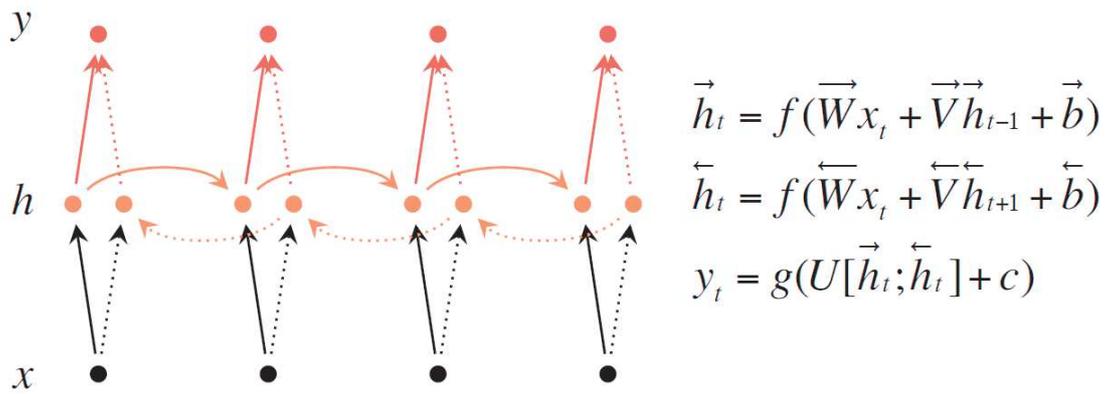
---

# Typical usages of RNNs



| | | |
|---|---|---|
| Single - Multiple | | Image Captioning |
| Multiple - Single | | Sentiment Classification |
| Multiple - Multiple | | Translation |
| | **Sequence to Sequence** | |
| | | Image Captioning |

**Input into RNN the features from last convolutional layer**

RNN

CNN

The    dog

Linear Classifier    Linear Classifier

$h_2$    $h_3$

RNN — $h_1$ → RNN — $h_2$ → RNN — $h_3$ → - - - - -

CNN

**For example, for image captioning**

---

Outputs    $y_{t-1}$    $y_t$    $y_{t+1}$

$h^3_{t-1}$    $h^3_t$    $h^3_{t+1}$

Hidden Layers    $h^2_{t-1}$    $h^2_t$    $h^2_{t+1}$

$h^1_{t-1}$    $h^1_t$    $h^1_{t+1}$

Inputs    $x_{t-1}$    $x_t$    $x_{t+1}$

$\Pr(x_t|y_{t-1})$    $\Pr(x_{t+1}|y_t)$

**Several RNNs stacked (like layers in MLP)**

$$\vec{h}_t = f(\vec{W}x_t + \vec{V}\vec{h}_{t-1} + \vec{b})$$

$$\overleftarrow{h}_t = f(\overleftarrow{W}x_t + \overleftarrow{V}\overleftarrow{h}_{t+1} + \overleftarrow{b})$$

$$y_t = g(U[\vec{h}_t;\overleftarrow{h}_t] + c)$$
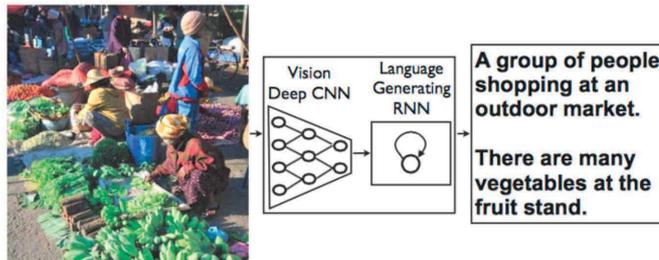
$h = [\vec{h};\overleftarrow{h}]$ now represents (summarizes) the past and future around a single token.

## (e.g. for offline classification of sequence of words)

## Wherever data is intrinsicly SEQUENTIAL

- **Speech recognition**
- ***Natural Language Processing (NLP)***
  - **Machine-Translation**
  - **Image caption generator**



- **Gesture recognition**
- **Music generation**
- ***Potentially any kind of time-series!!***

---

- **For SEQUENTIAL data (speech, text, …, gestures, …)**

- **Impressive results in Natural Language Processing (in particular Automated Real-Time Translation)**

- **Training of standard RNNs can be tricky (vanishing gradient…)**

- **LSTM / GRU now more used than standard RNNs**

# Any QUESTIONS ?