

DLAI – Marta R. Costa-jussà

# **Recurrent Neural Networks**

# Outline

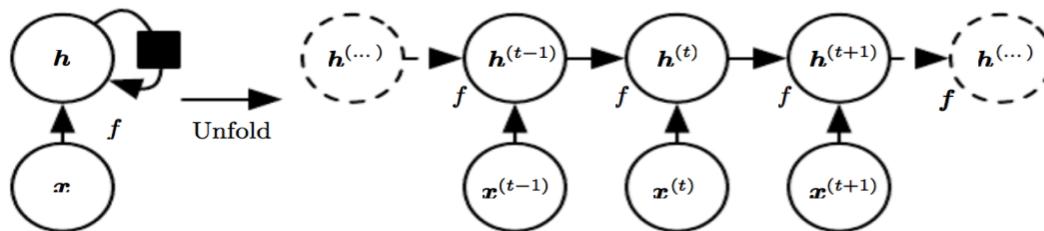
- Introduction
- Recurrent Neural Networks
- Vanishing gradient
- Gated Units
- Variants
- Example

# Motivation

- In artificial intelligence applications we often deal with SEQUENCES (e.g. speech, video, machine translation...)
- We need models that take information of SEQUENCE
- HOW?
- By recurrency

# Recurrent computational graph

- Regardless of the sequence length, the learned model always has **the same input size**, because it is specified in terms of transition from one state to another state, rather than specified in terms of a variable-length history of states.
- It is possible to use the **same transition function  $f$**  with **the same parameters** at every time step

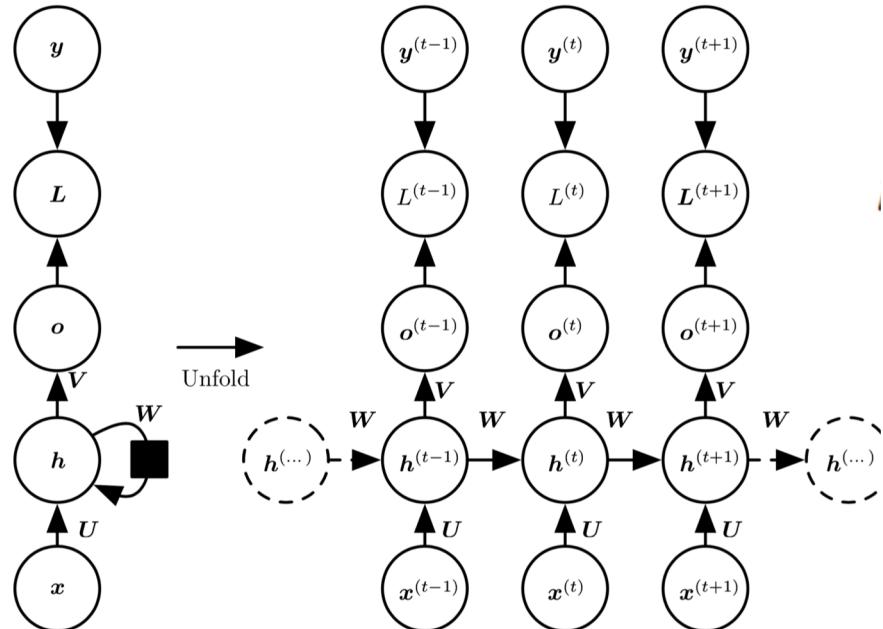
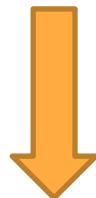


# Vanilla RNN

# Recurrent Neural Network (RNN) adding the “temporal” evolution

Allow to build specific  
connections capturing “history”

$$\begin{aligned} \mathbf{a}^{(t)} &= \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)}, \\ \mathbf{h}^{(t)} &= \tanh(\mathbf{a}^{(t)}), \\ \mathbf{o}^{(t)} &= \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)}, \\ \hat{\mathbf{y}}^{(t)} &= \text{softmax}(\mathbf{o}^{(t)}), \end{aligned}$$



TIME-STEPs: the memory do you want to include in your network.  
If you want your network to have memory of 60 characters, this number should be 60.

# Recurrent Neural Network (RNN) : sharing Weights

Hence we have two data flows:

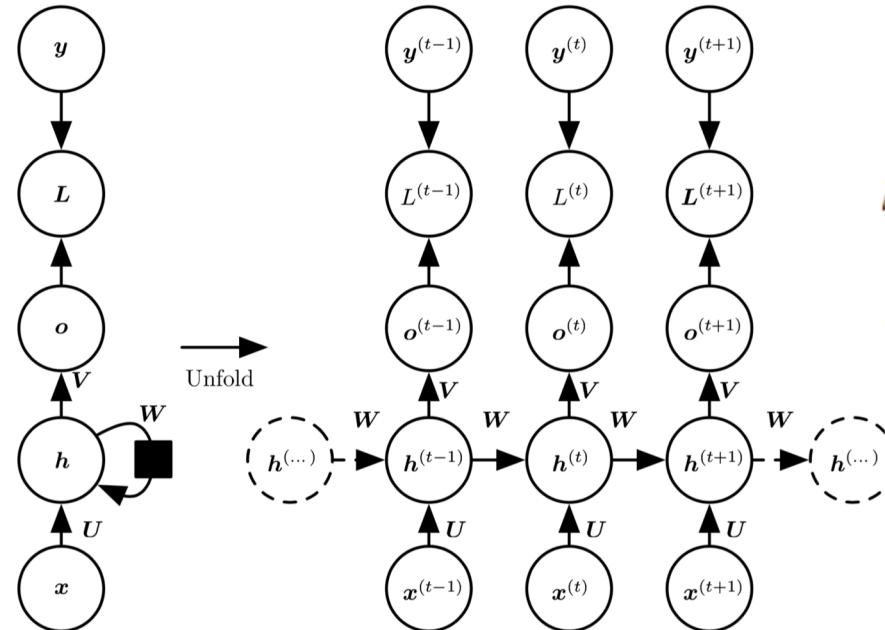
**Forward in space + time**

propagation: **2 projections per layer**

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

W, U, V shared across all steps

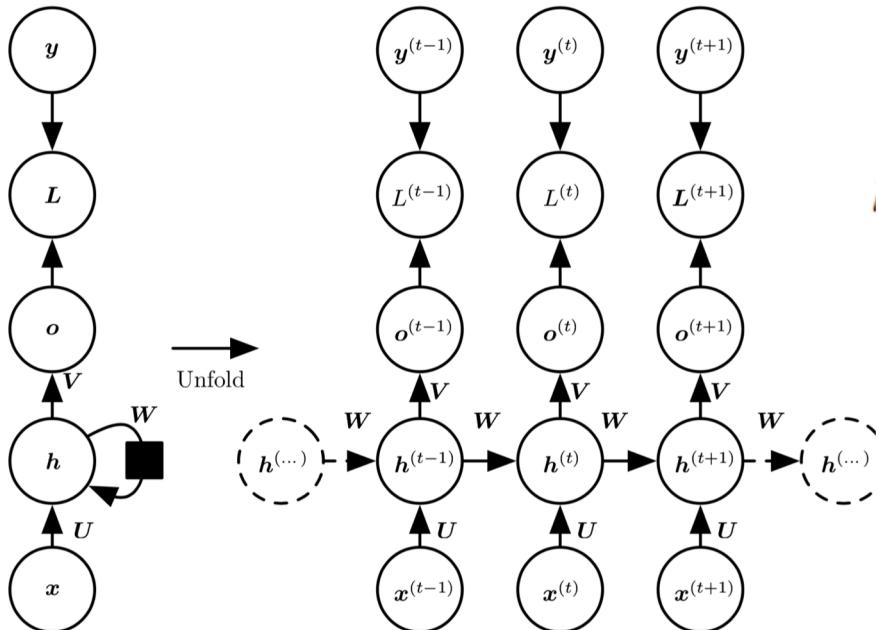
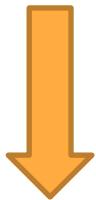
$$\begin{aligned} \mathbf{a}^{(t)} &= \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)}, \\ \mathbf{h}^{(t)} &= \tanh(\mathbf{a}^{(t)}), \\ \mathbf{o}^{(t)} &= \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)}, \\ \hat{\mathbf{y}}^{(t)} &= \text{softmax}(\mathbf{o}^{(t)}), \end{aligned}$$



# Recurrent Neural Network (RNN) : parameters

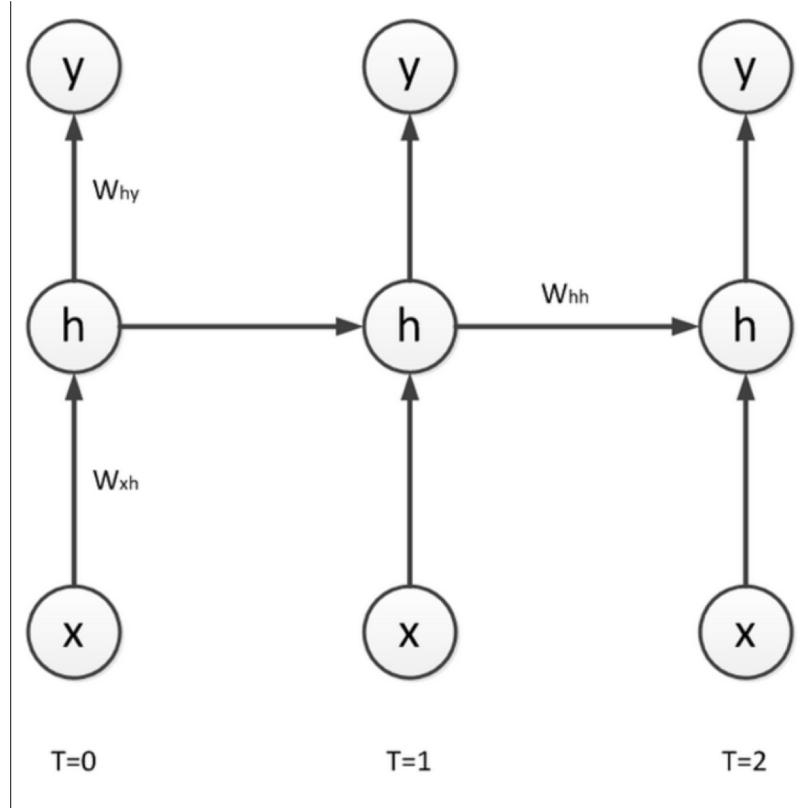
Allow to build specific  
connections capturing "history"

$$\begin{aligned} \mathbf{a}^{(t)} &= \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)}, \\ \mathbf{h}^{(t)} &= \tanh(\mathbf{a}^{(t)}), \\ \mathbf{o}^{(t)} &= \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)}, \\ \hat{\mathbf{y}}^{(t)} &= \text{softmax}(\mathbf{o}^{(t)}), \end{aligned}$$

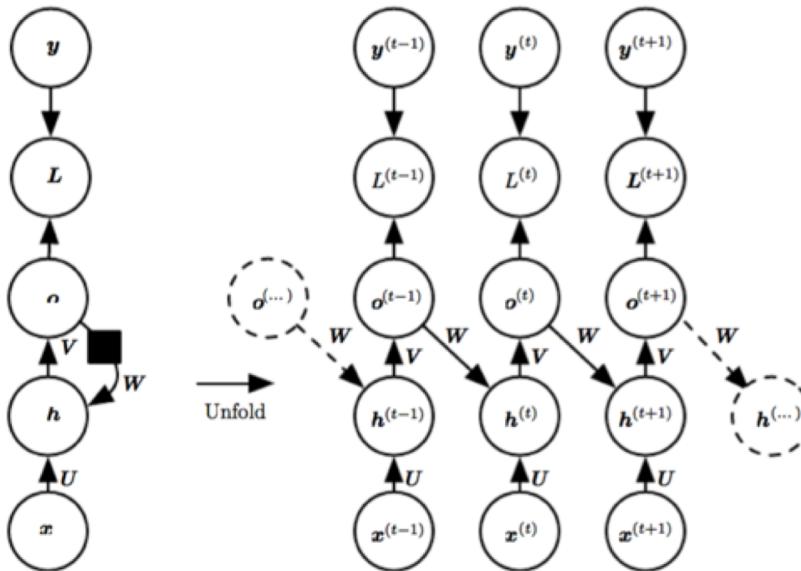


# Quick Exercise

- The figure shows a RNN with one input unit  $x$ , one logistic hidden unit  $h$ , and one linear output unit  $y$ . The Network parameters are  $W_{xh}=-0.1$ ,  $W_{hh}=0.5$  and  $W_{hy}=0.25$ ,  $hbias=0.4$  and  $ybias=0.0$ . The input takes the values 18, 9, -8 at time steps 0,1 and 2.
- 1-Compute the hidden value  $h_0$
- 2-Compute the output value  $y_1$
- 3-Compute the output value  $y_2$



# Alternatives of recurrence



# Training a RNN I: BPTT

- Backpropagation through time (BPTT): The training algorithm for updating network weights to minimize error including time

# Remember BackPropagation

1. Present a training input pattern and propagate it through the network to get an output.
2. Compare the predicted outputs to the expected outputs and calculate the error.
3. Calculate the derivatives of the error with respect to the network weights.
4. Adjust the weights to minimize the error.
5. Repeat.

# BPTT I: Loss

$$\begin{aligned} L & \left( \{\boldsymbol{x}^{(1)}, \dots, \boldsymbol{x}^{(\tau)}\}, \{\boldsymbol{y}^{(1)}, \dots, \boldsymbol{y}^{(\tau)}\} \right) \\ &= \sum_t L^{(t)} \\ &= - \sum_t \log p_{\text{model}} \left( y^{(t)} \mid \{\boldsymbol{x}^{(1)}, \dots, \boldsymbol{x}^{(t)}\} \right), \end{aligned}$$

The total loss for a given sequence  $x(t)$

$$\frac{\partial L}{\partial L^{(t)}} = 1.$$

Our goal is to calculate the gradients of the error with respect to our parameters U, W and V and then learn good parameters using Stochastic Gradient Descent.

# BPTT II: backward in time and in space

$$(\nabla_{\boldsymbol{o}^{(t)}} L)_i = \frac{\partial L}{\partial o_i^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial o_i^{(t)}} = \hat{y}_i^{(t)} - \mathbf{1}_{i=y^{(t)}}.$$

Gradient on the outputs at time step t

Backward in time

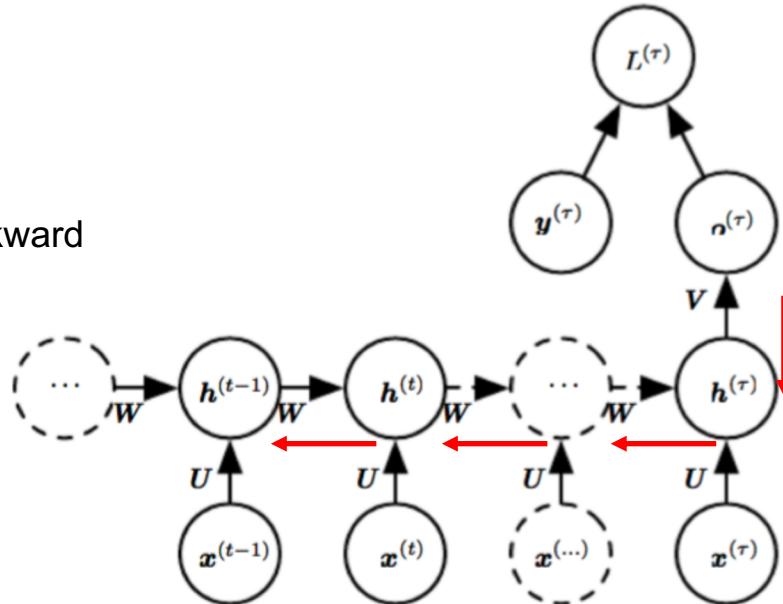
Backward in space

$$\begin{aligned}\nabla_{\boldsymbol{h}^{(t)}} L &= \left( \frac{\partial \boldsymbol{h}^{(t+1)}}{\partial \boldsymbol{h}^{(t)}} \right)^\top (\nabla_{\boldsymbol{h}^{(t+1)}} L) + \left( \frac{\partial \boldsymbol{o}^{(t)}}{\partial \boldsymbol{h}^{(t)}} \right)^\top (\nabla_{\boldsymbol{o}^{(t)}} L) \\ &= \boldsymbol{W}^\top \text{diag} \left( 1 - \left( \boldsymbol{h}^{(t+1)} \right)^2 \right) (\nabla_{\boldsymbol{h}^{(t+1)}} L) + \boldsymbol{V}^\top (\nabla_{\boldsymbol{o}^{(t)}} L),\end{aligned}$$

# BPTT III: gradient of the Loss

$$\nabla_{h^{(\tau)}} L = V^\top \nabla_{o^{(\tau)}} L.$$

Start at the end of the sequence, going backward



# BPTT IV: gradient on the remaining parameters

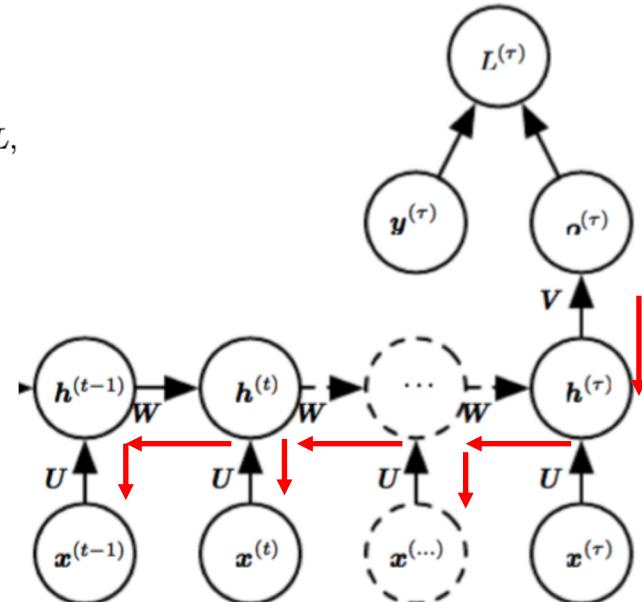
$$\nabla_{\mathbf{c}} L = \sum_t \left( \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{c}} \right)^\top \nabla_{\mathbf{o}^{(t)}} L = \sum_t \nabla_{\mathbf{o}^{(t)}} L,$$

$$\nabla_{\mathbf{b}} L = \sum_t \left( \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{b}^{(t)}} \right)^\top \nabla_{\mathbf{h}^{(t)}} L = \sum_t \text{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) \nabla_{\mathbf{h}^{(t)}} L,$$

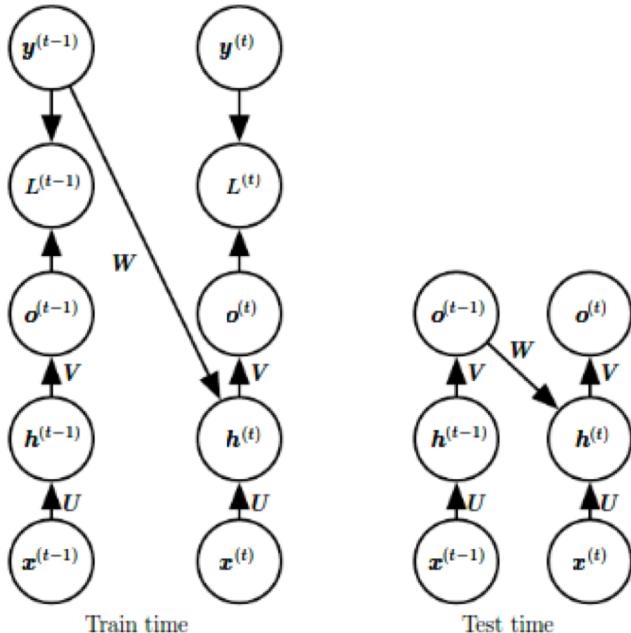
$$\nabla_{\mathbf{V}} L = \sum_t \sum_i \left( \frac{\partial L}{\partial o_i^{(t)}} \right) \nabla_{\mathbf{V}^{(t)}} o_i^{(t)} = \sum_t (\nabla_{\mathbf{o}^{(t)}} L) \mathbf{h}^{(t)\top},$$

$$\begin{aligned} \nabla_{\mathbf{W}} L &= \sum_t \sum_i \left( \frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{W}^{(t)}} h_i^{(t)} \\ &= \sum_t \text{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{h}^{(t-1)\top}, \end{aligned}$$

$$\begin{aligned} \nabla_{\mathbf{U}} L &= \sum_t \sum_i \left( \frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{U}^{(t)}} h_i^{(t)} \\ &= \sum_t \text{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{x}^{(t)\top}, \end{aligned}$$



# Note on teacher forcing



Pro's

No need BPTT  
Parallel Training

Con's

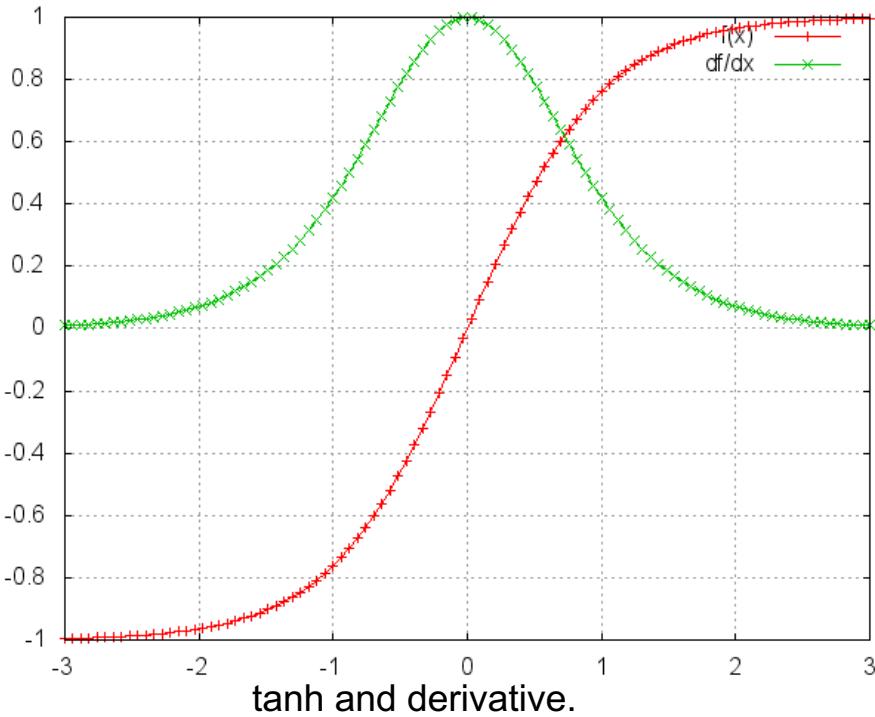
Exposure bias (but in practice is not so relevant)  
Recurrence from the output, not internal state

# Vanishing gradient

# Vanishing gradient I

- During training gradients explode/vanish easily because of depth-in-time → Exploding/Vanishing gradients!

## Vanishing gradient II



- Traditional activation functions such as hyperbolic tangent have gradients in the range (-1,1) and backpropagation computes gradients by the chain rule

- The vanishing (and exploding) gradient problem is caused by the repeated use of the recurrent weight matrix in RNN
  - This has the effect of multiplying n of these small numbers to compute gradients
  - This means that the gradient (error signal) decreases exponentially with n

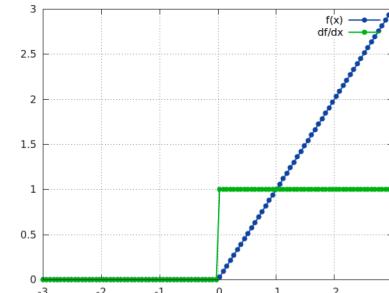
$$\begin{aligned}\nabla_{\mathbf{W}} L &= \sum_t \sum_i \left( \frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{W}^{(t)}} h_i^{(t)} \\ &= \sum_t \text{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{h}^{(t-1)\top},\end{aligned}$$

# Question

- Do FNNs with several hidden layers suffer from vanishing gradient problem?

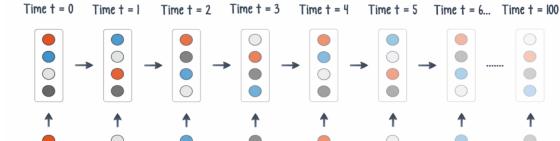
# Standard Solutions

- Proper initialization of Weight Matrix
- Regularization of outputs or Dropout
- Use of ReLU Activations as it's derivative is either 0 or 1

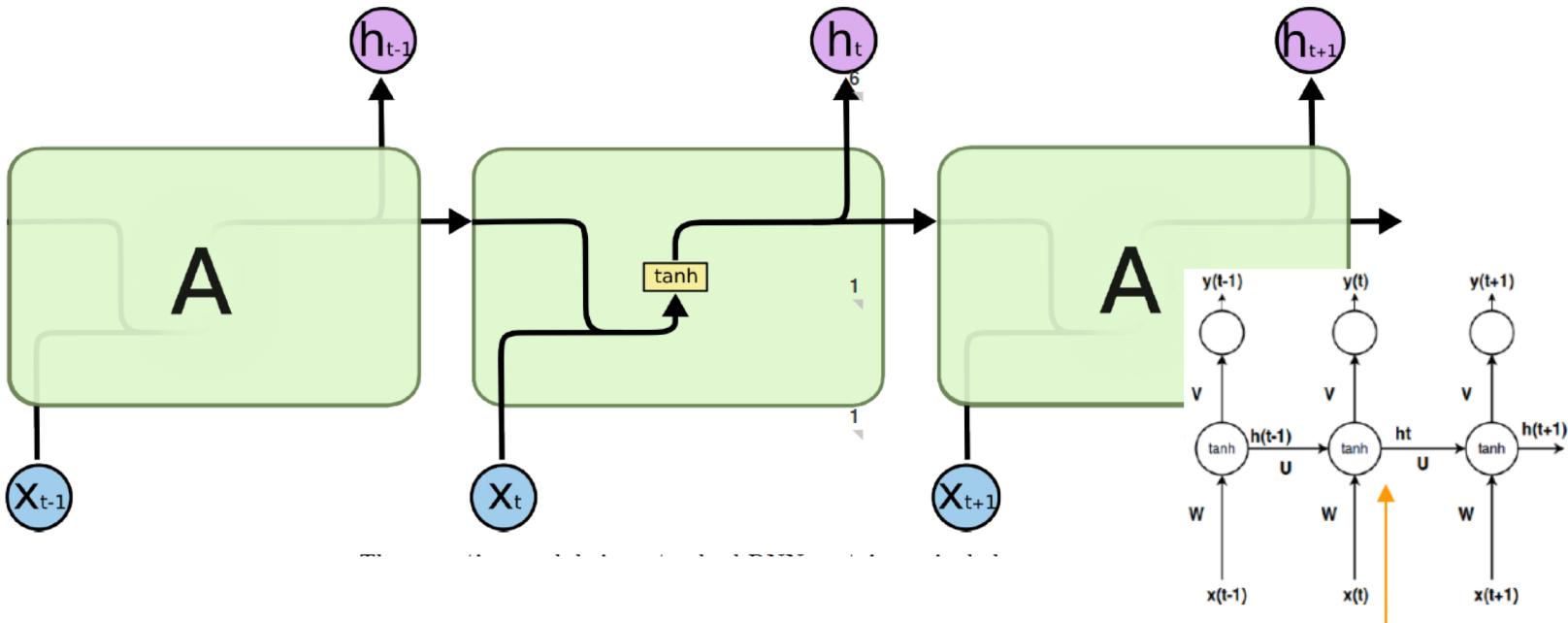


# Gated Units

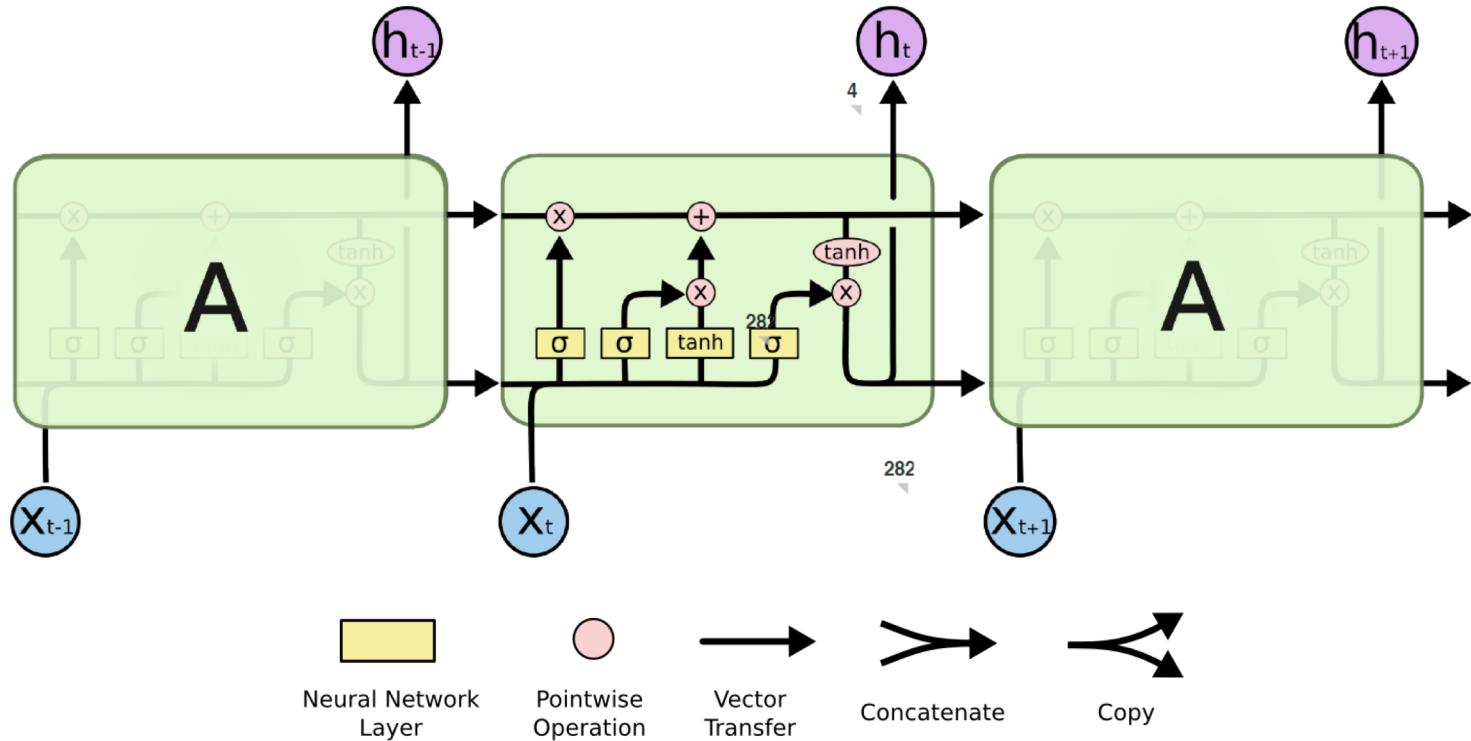
# Standard RNN



<https://www.nextbigfuture.com/2016/03/recurrent-neural-nets.html>



# Long-Short Term Memory (LSTM)



# Gating method

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 decide whether the current time-step information matters or not**, to accept or discard (optimized reading mechanism)
3. **Make every RNN unit able to forget whatever may not be useful anymore** by clearing that info from the cell state (optimized clearing 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..

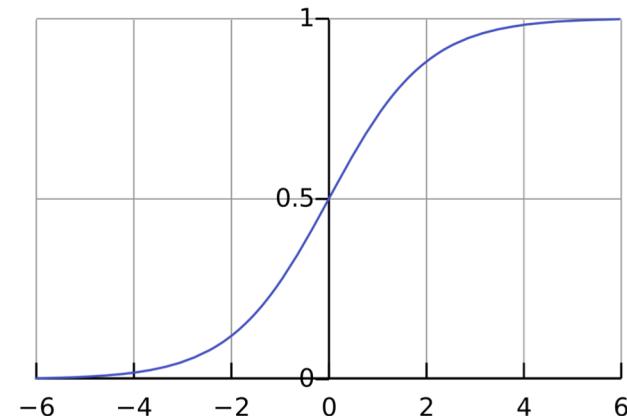
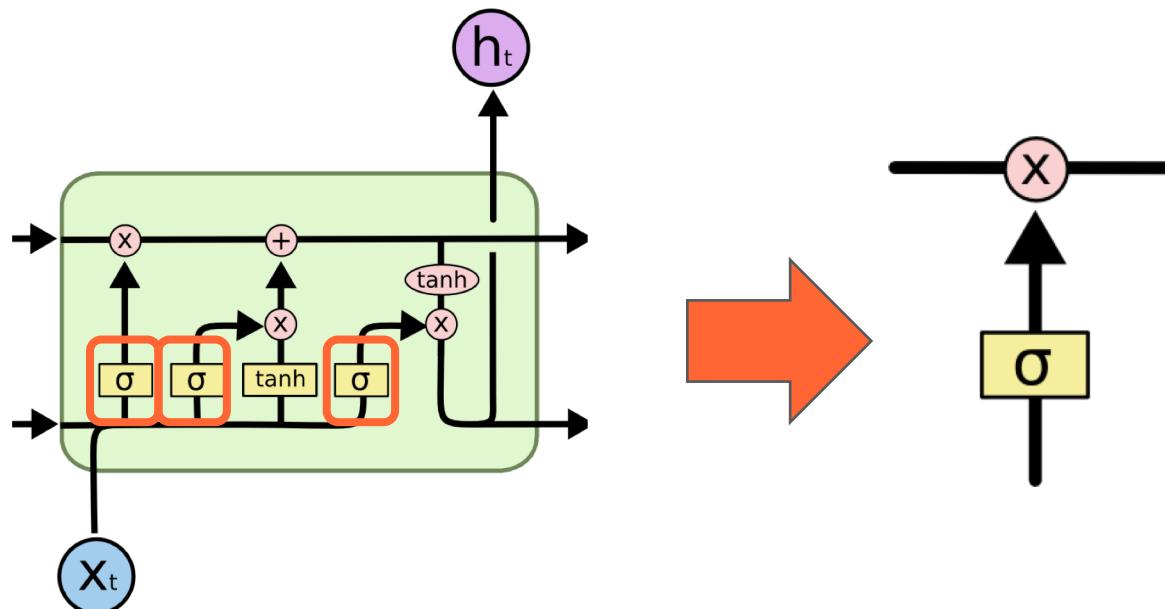
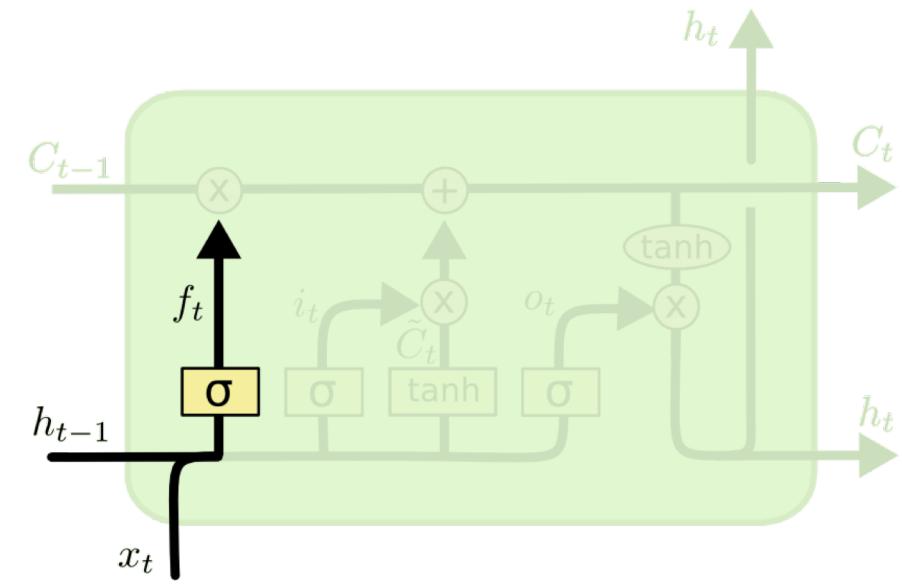


Figure: Cristopher Olah, [“Understanding LSTM Networks”](#) (2015)

slide credit: Xavi Giro

# Forget Gate



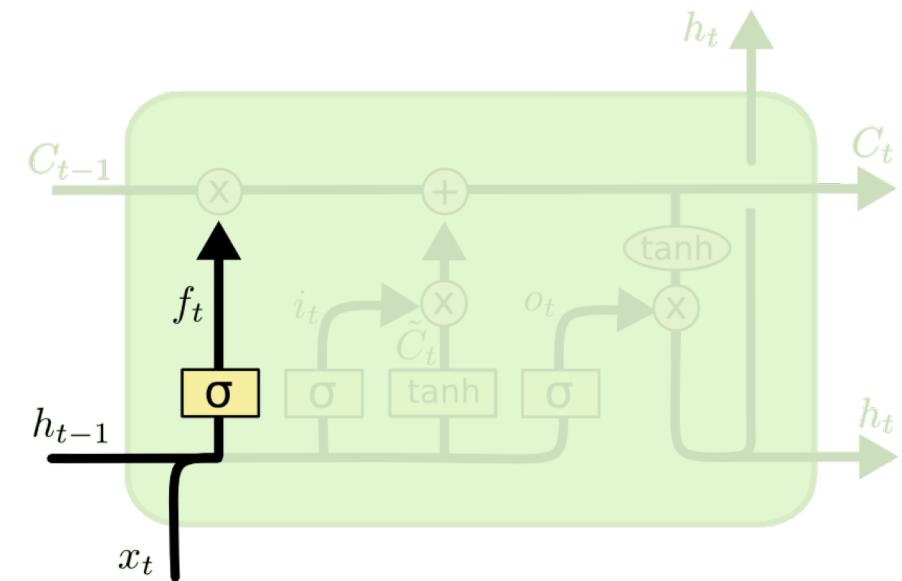
**Forget Gate:**

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

Concatenate

What part of memory to  
“forget” – zero means  
forget this bit

# Forget Gate: Example



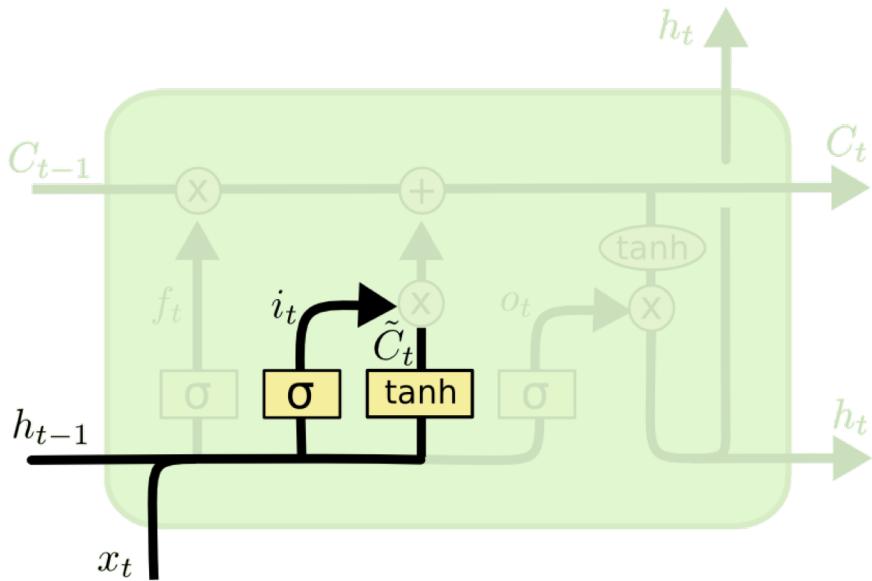
$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

LANGUAGE MODELING

Joan es un chico activo y **Anna** es una chica  
calmada

Forget about "male" gender

# Input Gate



What bits to insert into the next states

## Input Gate Layer

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

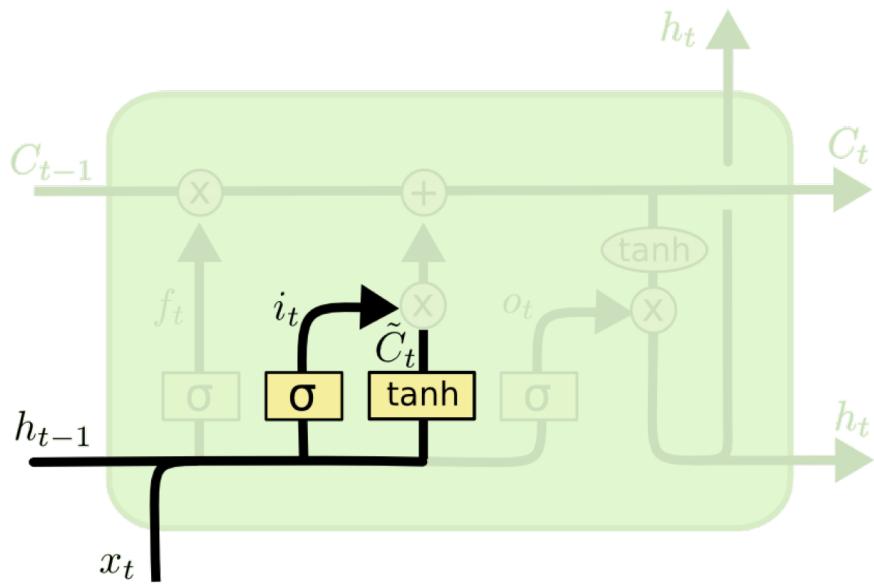
## New contribution to cell state

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

Classic neuron

What content to store into the next state

# Input Gate: Example



## Input Gate Layer

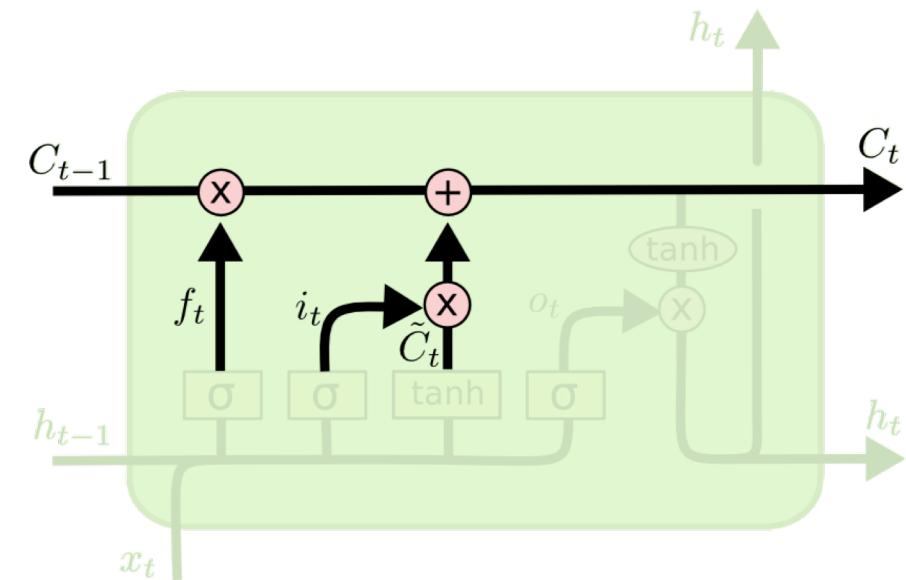
$$i_t = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$$

LANGUAGE MODELING

Joan es un chico activo y **Anna** es una chica  
calmada

Input about "female" gender

# Update Cell State

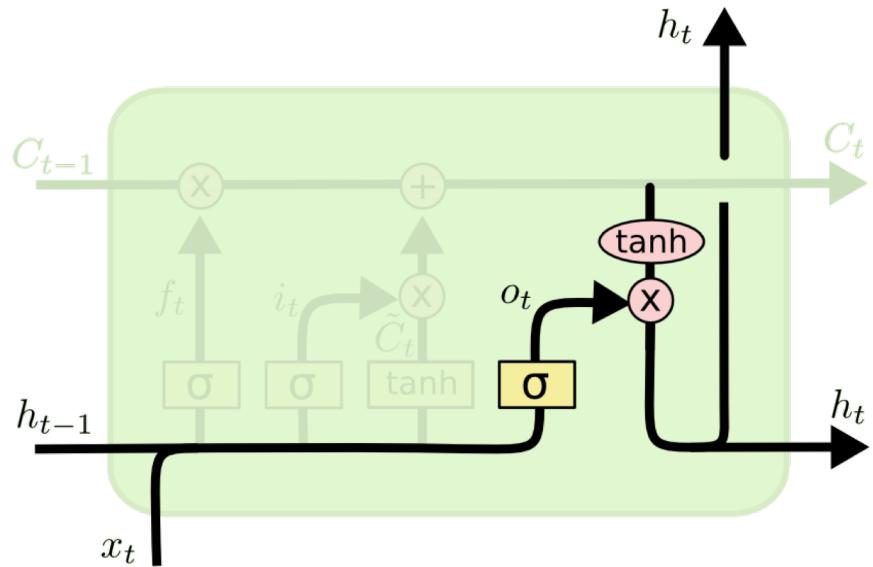


Next memory cell content – mixture of not-forgotten part of previous cell and insertion

## Update Cell State (memory):

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

# Output Gate



What part of cell to output

## Output Gate Layer

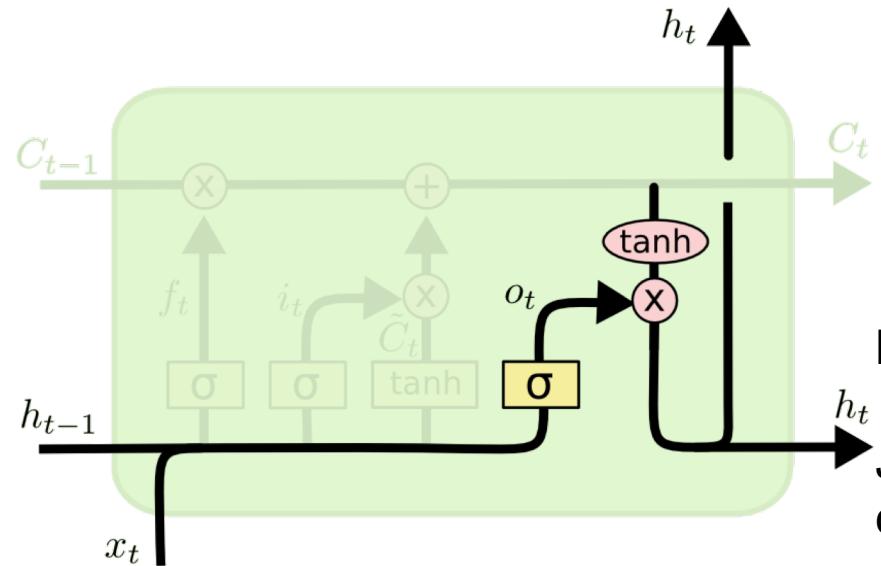
$$o_t = \sigma (W_o [ h_{t-1}, x_t ] + b_o)$$

## Output to next layer

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

tanh maps bits to  $[-1, +1]$  range

# Output Gate: Example



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

LANGUAGE MODELING

Joan es un chico activo y **Anna** es una chica  
calmada

Info relevant for a verb? : "female", "singular",  
"3<sup>rd</sup> person"

# Implementing an LSTM

For  $t = 1, \dots, T$ :

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

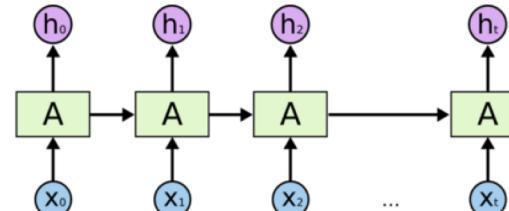
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

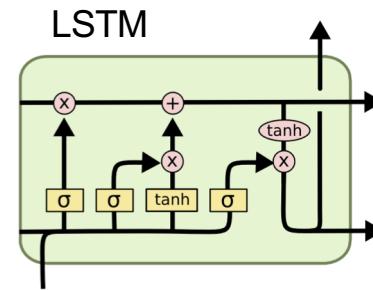
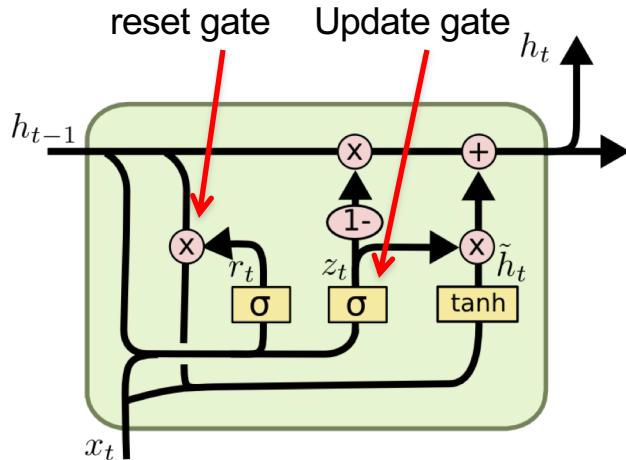
$$(2) \quad C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

$$(3) \quad o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

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



# GRU – gated recurrent unit (more compression)



$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

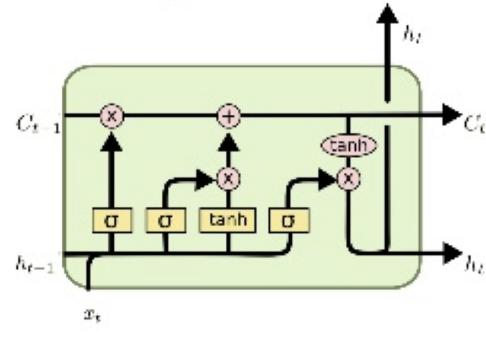
It combines the **forget** and **input** into a single **update gate**.

X,\*: element-wise multiply

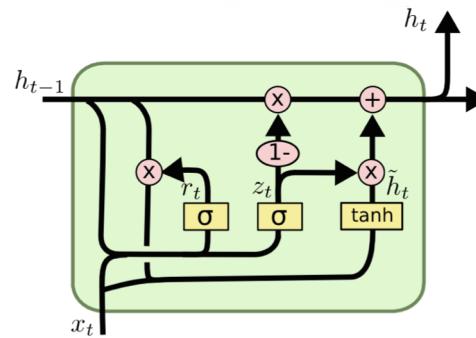
It also merges the cell state and hidden state. This is simpler than LSTM. There are many other variants too.

# LSTM and GRU

- LSTM [Hochreiter&Schmidhuber97]



- GRU [Cho+14]



GRUs also takes  $x_t$  and  $h_{t-1}$  as inputs. They perform some calculations and then pass along  $h_t$ . What makes them different from LSTMs is that GRUs don't need the cell layer to pass values along. The calculations within each iteration insure that the  $h_t$  values being passed along either retain a high amount of old information or are jump-started with a high amount of new information.

# Visual Comparison FNN, Vanilla RNNs and LSTMs

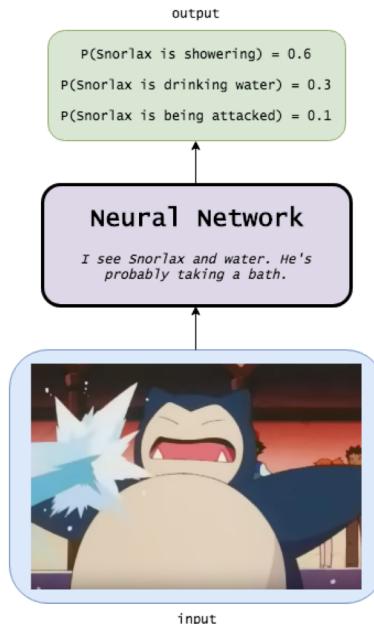


Image src <http://blog.echen.me/2017/05/30/exploring-lstms/>

# Visual Comparison FNN, Vanilla RNNs and LSTMs

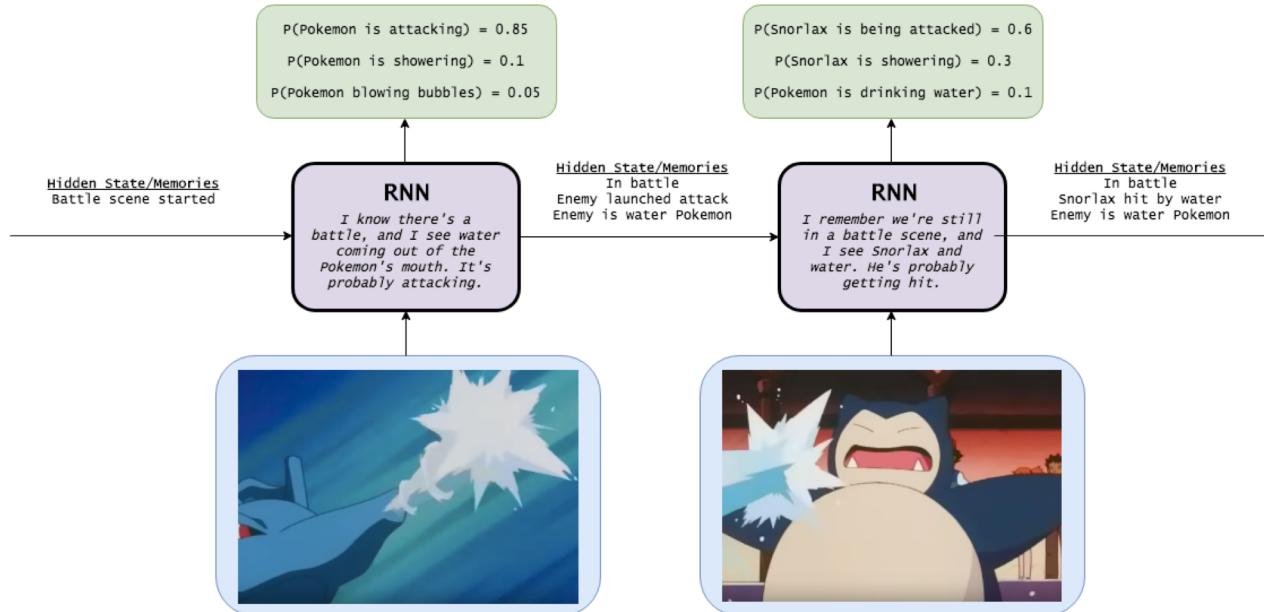


Image src <http://blog.echen.me/2017/05/30/exploring-lstms/>

# Visual Comparison FNN, Vanilla RNNs and LSTMs

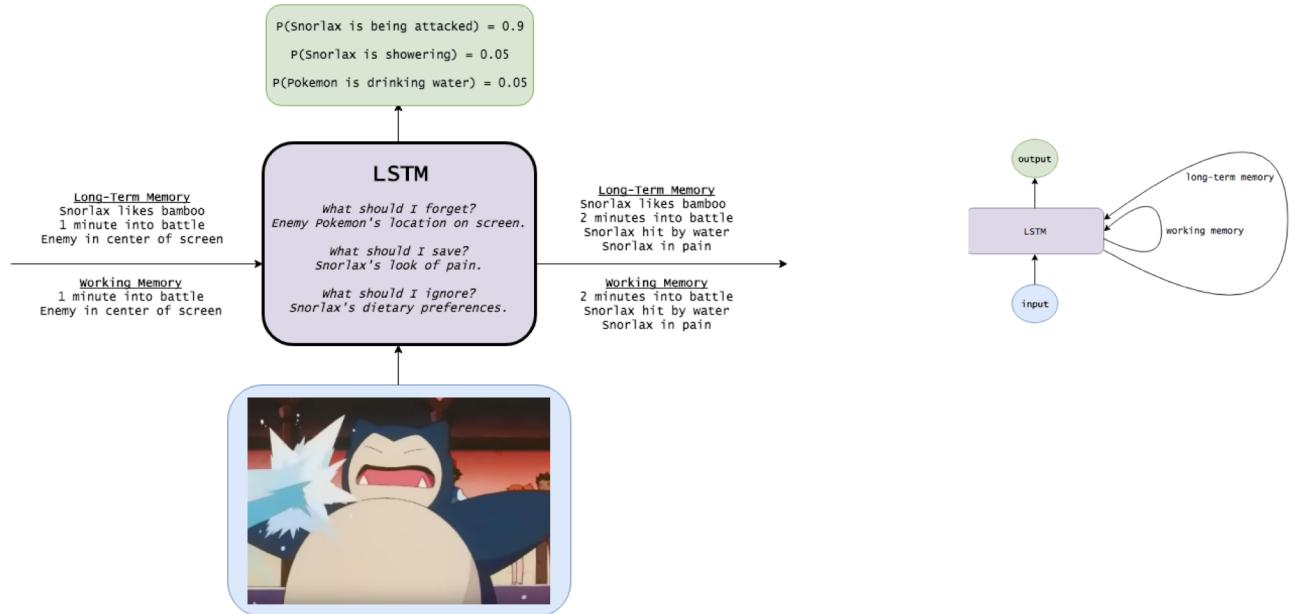
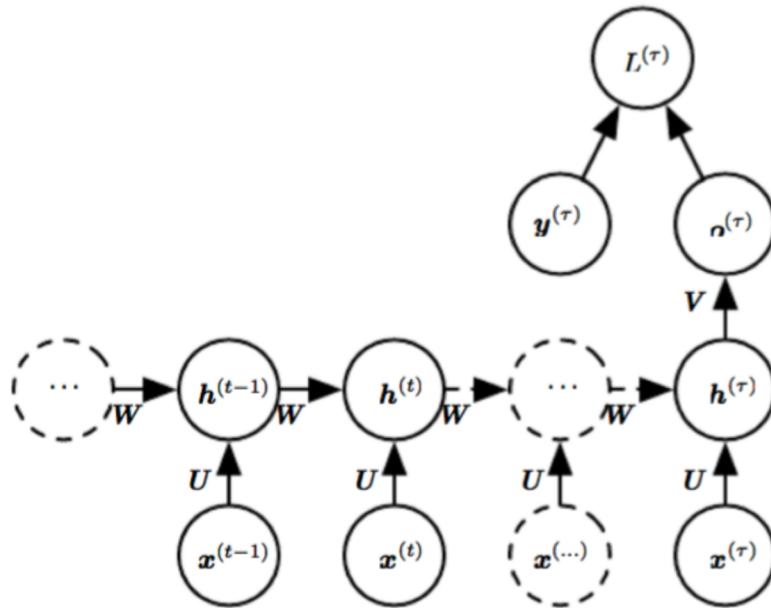


Image src <http://blog.echen.me/2017/05/30/exploring-lstms/>

# **Variants of RNNs**

# Sequence to vector



# Vector to sequence

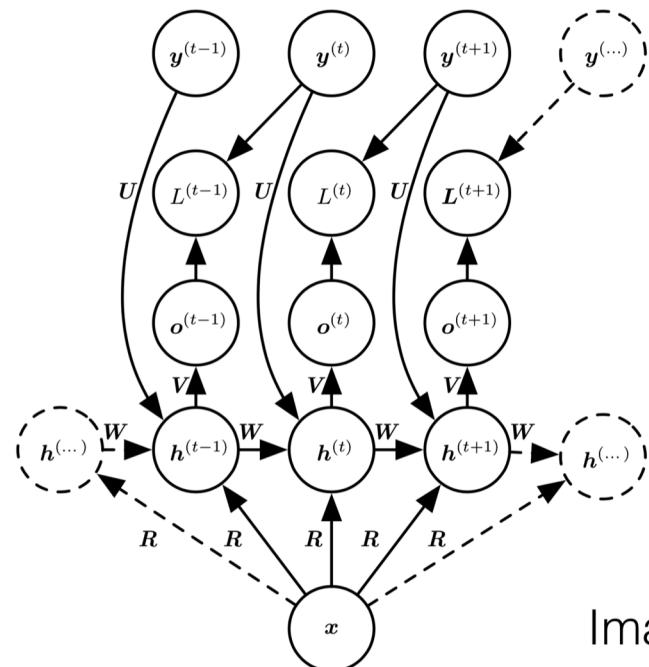
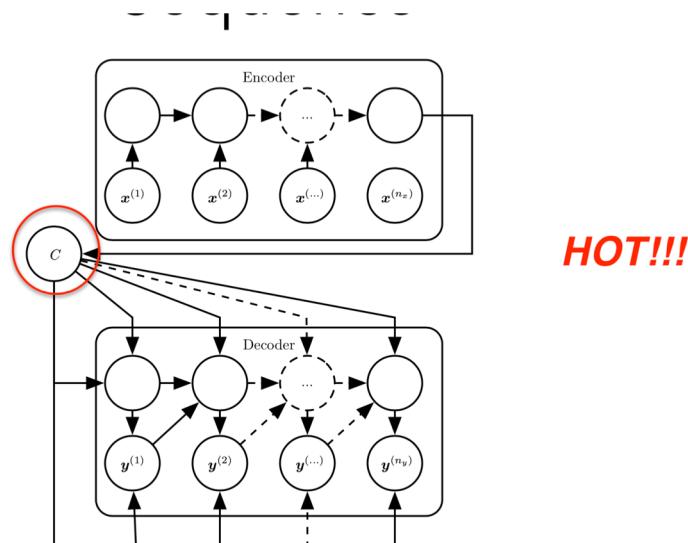


Image as extra input

# Exercise

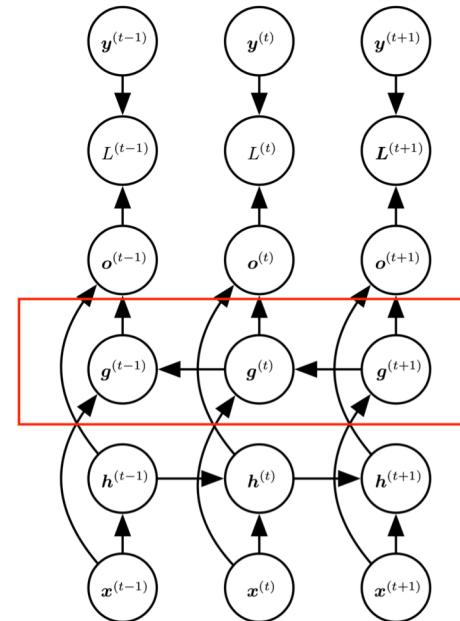
- Draw the computational graph of a sequence to sequence model

# Sequence to sequence



To learn the context Variable C which represents a semantic summary of the input sequence, and for later decoder RNNN

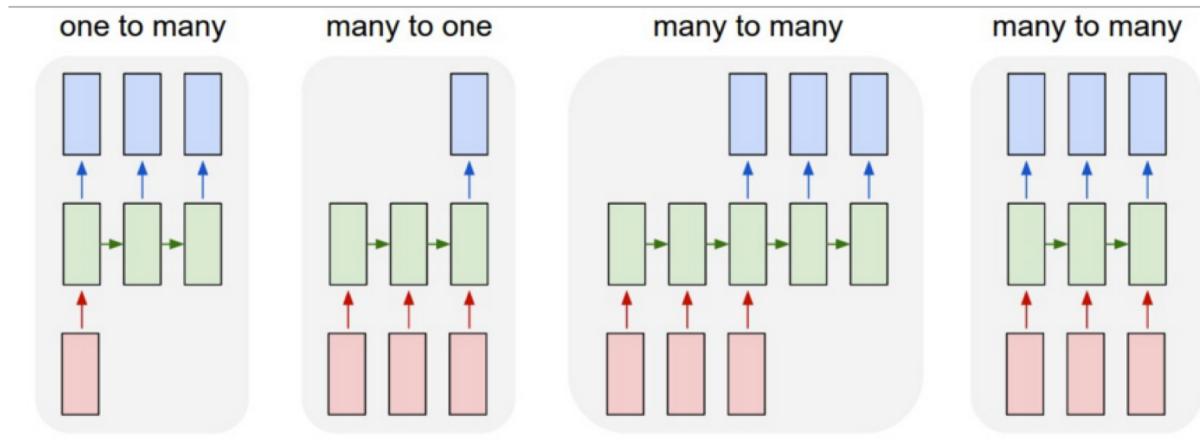
# Bidirectional RNNs



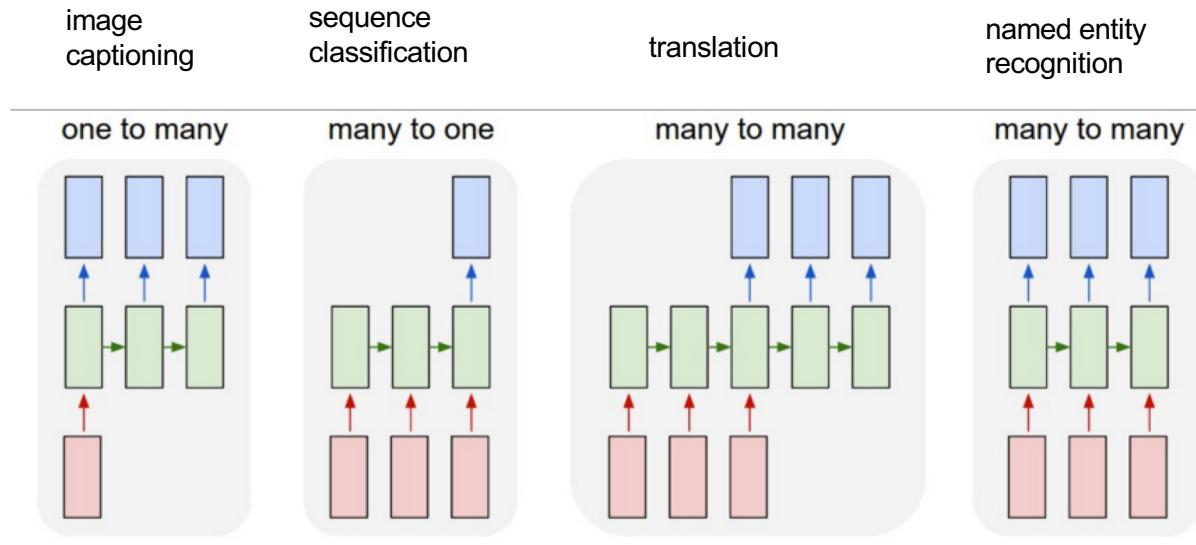
# Question

- You have a pet dog whose mood is heavily dependent on the current and past few days' weather. You've collected data for the past 365 days on the weather, which you represent as a sequence as  $x<1>, \dots, x<365>$ . You've also collected data on your dog's mood, which you represent as  $y<1>, \dots, y<365>$ . You'd like to build a model to map from  $x \rightarrow y$ . Should you use a Unidirectional RNN or Bidirectional RNN for this problem?

# Question: Can LSTMs be used for other sequence tasks: which ones?

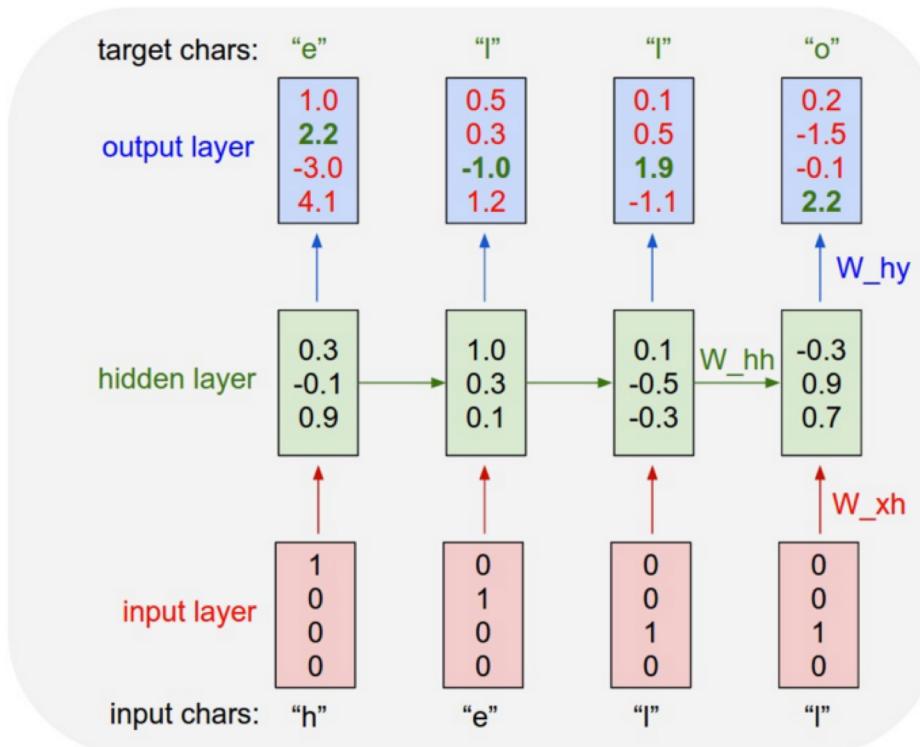


# LSTMs can be used for other sequence tasks: which ones?



# LSTM Example

# Character-level language model



Test time:

- pick a seed character sequence
- generate the next character
- then the next
- then the next ...

# Character-level language model

PANDARUS:

Alas, I think he shall be come approached and the day  
When little strain would be attain'd into being never fed,  
And who is but a chain and subjects of his death,  
I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul,  
Breaking and strongly should be buried, when I perish  
The earth and thoughts of many states.

DUKE VINCENTIO:

Well, your wit is in the care of side and that.

Second Lord:

They would be ruled after this chamber, and  
my fair nues begun out of the fact, to be conveyed,  
Whose noble souls I'll have the heart of the wars.

# Character-level language model

First Citizen:

Nay, then, that was hers,  
It speaks against your other service:  
But since the  
youth of the circumstance be spoken:  
Your uncle and one Baptista's daughter.

Yoav Goldberg:  
order-10 unsmoothed  
character n-grams

SEBASTIAN:

Do I stand till the break off.

BIRON:

Hide thy head.

VENTIDIUS:

He purposeth to Athens: whither, with the vow  
I made to handle you.

FALSTAFF:

My good knave.

# Character-level language model

MMMMM---- Recipe via Meal-Master (tm) v8.05

Title: BARBECUE RIBS

Categories: Chinese, Appetizers

Yield: 4 Servings

1 pk Seasoned rice

1 Beer -- cut into  
-cubes

1 ts Sugar

3/4 c Water  
Chopped finels,  
-up to 4 tblsp of chopped

2 pk Yeast Bread/over

MMMMM-----FILLING-

2 c Pineapple, chopped

1/3 c Milk

1/2 c Pecans  
Cream of each

2 tb Balsamic cocoa

2 tb Flour

2 ts Lemon juice  
Granulated sugar

2 tb Orange juice

1 c Sherry wheated curdup

1 Onion; sliced

1 ts Salt

2 c Sugar

1/4 ts Salt  
1/2 ts White pepper, freshly ground  
Sesame seeds

1 c Sugar

1/4 c Shredded coconut

1/4 ts Cumin seeds

Preheat oven to 350. In a medium bowl, combine milk, flour and water and then cornstarch. add tomatoes, or nutmeg; serve.

# Character-level language model

For  $\bigoplus_{n=1,\dots,m}$  where  $\mathcal{L}_{m,\bullet} = 0$ , hence we can find a closed subset  $\mathcal{H}$  in  $\mathcal{H}$  and any sets  $\mathcal{F}$  on  $X$ ,  $U$  is a closed immersion of  $S$ , then  $U \rightarrow T$  is a separated algebraic space.

*Proof.* Proof of (1). It also start we get

$$S = \text{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by  $\coprod Z \times_U U \rightarrow V$ . Consider the maps  $M$  along the set of points  $Sch_{fppf}$  and  $U \rightarrow U$  is the fibre category of  $S$  in  $U$  in Section, ?? and the fact that any  $U$  affine, see Morphisms, Lemma ???. Hence we obtain a scheme  $S$  and any open subset  $W \subset U$  in  $Sh(G)$  such that  $\text{Spec}(R') \rightarrow S$  is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that  $f_i$  is of finite presentation over  $S$ . We claim that  $\mathcal{O}_{X,x}$  is a scheme where  $x, x', s'' \in S'$  such that  $\mathcal{O}_{X,x'} \rightarrow \mathcal{O}'_{X',x'}$  is separated. By Algebra, Lemma ?? we can define a map of complexes  $\text{GL}_{S'}(x'/S'')$  and we win.  $\square$

To prove study we see that  $\mathcal{F}|_U$  is a covering of  $\mathcal{X}'$ , and  $\mathcal{T}_i$  is an object of  $\mathcal{F}_{X/S}$  for  $i > 0$  and  $\mathcal{F}_p$  exists and let  $\mathcal{F}_i$  be a presheaf of  $\mathcal{O}_X$ -modules on  $\mathcal{C}$  as a  $\mathcal{F}$ -module. In particular  $\mathcal{F} = U/\mathcal{F}$  we have to show that

$$\widetilde{M}^\bullet = \mathcal{I}^\bullet \otimes_{\text{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1} \mathcal{F}$$

is a unique morphism of algebraic stacks. Note that

$$\text{Arrows} = (Sch/S)_{fppf}^{opp}, (Sch/S)_{fppf}$$

and

$$V = \Gamma(S, \mathcal{O}) \longmapsto (U, \text{Spec}(A))$$

LaTeX “almost compiles”

# Character-level language model

```
/*
 * Increment the size file of the new incorrect UI_FILTER group information
 * of the size generatively.
 */
static int indicate_policy(void)
{
    int error;
    if (fd == MARN_EPT) {
        /*
         * The kernel blank will coeld it to userspace.
         */
        if (ss->segment < mem_total)
            unblock_graph_and_set_blocked();
        else
            ret = 1;
        goto bail;
    }
    segaddr = in_SB(in.addr);
    selector = seg / 16;
    setup_works = true;
    for (i = 0; i < blocks; i++) {
        seq = buf[i++];
        bpf = bd->bd.next + i * search;
        if (fd) {
            current = blocked;
        }
    }
}
```

# Get some practice after the class

- a simple RNN character model with PyTorch

<https://github.com/gabrielloye/RNN-walkthrough/blob/master/main.ipynb>

