# Recurrent Neural Networks (RNNs)

Michel RIVEILL

michel.riveill@univ-cotedazur.fr

▸ Until now, the order of the entries was not taken into account...

*→ Order of words is important :*

▸ Except that this one is important:

▸ word order:

▸ I like chocolate but not beer ≠ I like beer but not chocolate

▸ Music ?



▸ Signal ?

# Motivation

- Humans don't start their thinking from scratch every second
  - Thoughts have persistence
- Traditional neural networks can't characterize this phenomena
  - Ex: classify what is happening at every point in a movie
  - How a neural network can inform later events about the previous ones
- Recurrent neural networks address this issue
  - Some applications
    - NLP: Same word may have a different label depending on the context.
      - NER - Apple CEO Tim Cook eat an apple
      - POS - The bank is located on the bank (La banque est située sur la berge)
    - Forcasting - Time-series Prediction
- How?
  - Add state to artificial neurons

# What are RNNs?

- Main idea is to make use of sequential information
- How RNN is different from neural network?
  - Vanilla neural networks (MLP) assume all inputs are independent of each other
    - Features independence
  - But for many tasks, that's a very bad idea
- What RNN does?
  - Perform the same task for every element of a sequence
    - That's what recurrent stands for
  - Output depends on the previous computations!
- Another way of interpretation – RNNs have a "memory"
  - To store previous computations

# Some applications (not recent)

- RNN Generated TED Talks
  - YouTube Link - https://youtu.be/-OodHtJIsaY?t=31s
- RNN Generated Eminem rapper
  - RNN Shady - https://soundcloud.com/mrchrisjohnson/recurrent-neural-shady
- RNN Generated Music
  - Music Link - http://www.hexahedria.com/2015/08/03/composing-music-with-recurrent-neural-networks/

# From vanilla NN to recurrent NN

▸ Vanilla cell

  ▸ $y = F(U.X)$

y

$F$

U

X

One i/p
1 weight

# From vanilla NN to recurrent NN

▸ Vanilla cell

   ▸ $y = F(U.X)$

▸ Recurrent cell → use 2 weights matrix

   ▸ Add an internal variable: $h$

   ▸ The output depends to the current entry and the previous internal variable:

      ▸ $h_t = F(W.h_{t-1} + U.X_t)$

      ▸ Could be rewritten on $h_t = F(V.[h_{t-1}, X_t])$

$h_t$

$h$ ■   F

$x_t$

# From vanilla NN to recurrent NN

- Vanilla cell
  - $y = F(U.X)$
- Recurrent cell → use for each time step the same weights matrix
  - $h_t = F(W.[h_{t\_1}, Xt])$
- Recurrent layer, step by step
  - at each time step
    - A new entry is being supplied
    - And a new output $(h_t)$ is calculated using:
      - The new input $X_t$
      - The output of the previous step $h_{t\_1}$

- $h_1 = F(W.[h_0, X_1])$
- $h_2 = F(W.h_1, X_2])$
- $h_3 = F(W.h_2, X_3])$

- ...

# From vanilla NN to recurrent NN

▸ Vanilla cell

  ▸ $y = F(U.X)$

▸ Recurrent cell

  ▸ $h_t = F(W.[h_{t-1}, X_t])$

▸ Recurrent neural networks are "unrolled" programmatically during training and prediction

  ▸ All neurons share the same weight matrix

# Remember

| From | To |
|------|-----|
|  |  |

# Remember

| From | To |
|------|-----|
|  | |

# RNN in action

# Problems with naive RNN

- RNNs do not learn easily
- Unfolding the network for learning leads to vanishing gradient problems!

# Learning process
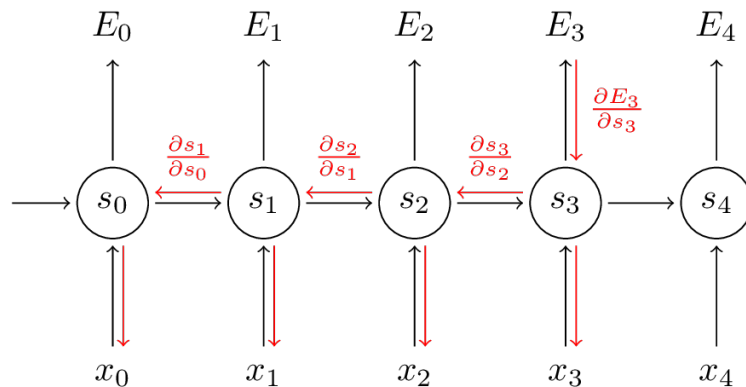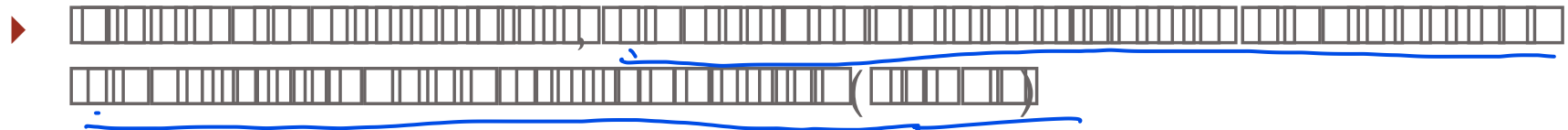## Vanishing gradient problem

https://www.youtube.com/watch?v=8z3DFk4VxRo

For rigorous proofs and derivations, please refer to

On the difficulty of training recurrent neural networks, Pascanu *et al.*, 2013

Long Short-Term Memory, Hochreiter *et al.*, 1997

# Main learning problem

▸ Remember

    ▸ In the backpropagation step we adjust weights with the following rules

        ▸ $W_{new} = W_{old} - \lambda \, gradient$

    ▸ if there are several layers, the gradient is proportional to the product of the derivative of the transfer function (sigmoid)



▸ $\dfrac{\partial E_3}{\partial W} = \sum_{k=0}^{3} \dfrac{\partial E_3}{\partial \widehat{y_3}} \dfrac{\partial \widehat{y_3}}{\partial s_3} \dfrac{\partial s_3}{\partial s_k} \dfrac{\partial s_k}{\partial W}$

▸ $\dfrac{\partial E_3}{\partial W} = \dfrac{\partial E_3}{\partial \widehat{y_3}} \dfrac{\partial \widehat{y_3}}{\partial s_3} \dfrac{\partial s_3}{\partial W} + \dfrac{\partial E_3}{\partial \widehat{y_3}} \dfrac{\partial \widehat{y_3}}{\partial s_3} \dfrac{\partial s_3}{\partial s_2} \dfrac{\partial s_2}{\partial W} + \ldots + \dfrac{\partial E_3}{\partial \widehat{y_3}} \dfrac{\partial \widehat{y_3}}{\partial s_3} \dfrac{\partial s_3}{\partial s_2} \dfrac{\partial s_2}{\partial s_1} \dfrac{\partial s_1}{\partial s_0} \dfrac{\partial s_0}{\partial W}$

# What is the value of the derivative 'chain'? $\frac{\partial s_3}{\partial s_2}\frac{\partial s_2}{\partial s_1}\frac{\partial s_1}{\partial s_0}$

▸ For sigmoid:
  ▸ Max = 0.25
  ▸ $0.25^2 = 0.0625$
  ▸ $0.25^4 = 0.00391$
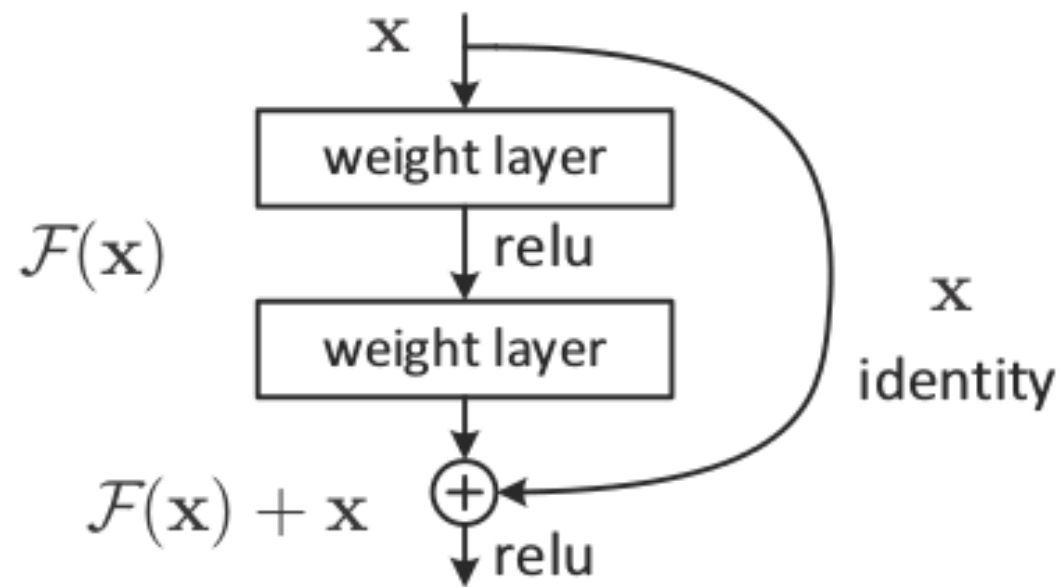  ▸ $0.25^8 = 0.0000152$
  ▸ $0.25^{16} = 0.00000000233$

▸ With RNN, we often have sequences longer than **50**
  ▸ **G**radient is near: $\mathbf{7 * 10^{-31}}$

▸ Usually, how can this be corrected?
  ▸ For short chains: change the activation function → relu vs. sigmoid
  ▸ **F**or long chains: use residuals

# Residual ?

# How to introduce residual in RNN

- The context is close to the word to be predicted
  - Few iterations separate them.
  - No problem
    - Few steps :
      - no vanishing gradient problem



$h_0$    $h_1$    $h_2$    $h_3$    $h_4$

A    A    A    A    A

$x_0$    $x_1$    $x_2$    $x_3$    $x_4$

clouds    in    the    blue    sky

# From vanilla Sort Term Memory...

- Traditionnal implementation of RNN cells
  - $h_t = tanh(W \cdot [h_{t-1}, xt\,])$
  - Involves a single level of processing
  - No control of hidden state
  - Creating the risk of the evanescent gradient.

Hidden state $h_{t-1}$         $h_t$         $h_{t+1}$

A    tanh    A

$X_{t-1}$ Input     $X_t$     $X_{t+1}$

# ... to LSTM (Long Short Term Memory)

- But if the context is far from the word to predict
  - Many iterations separate them!
  - Pb **1:** Possible gradient problem
  - Pb**2:** We want to be able to control the information to be stored in memory: to preserve the past or, on the contrary, to renew the state



I grew up in France...                                    I speak French...

# Dealing with the vanishing gradient problem → LSTM cell

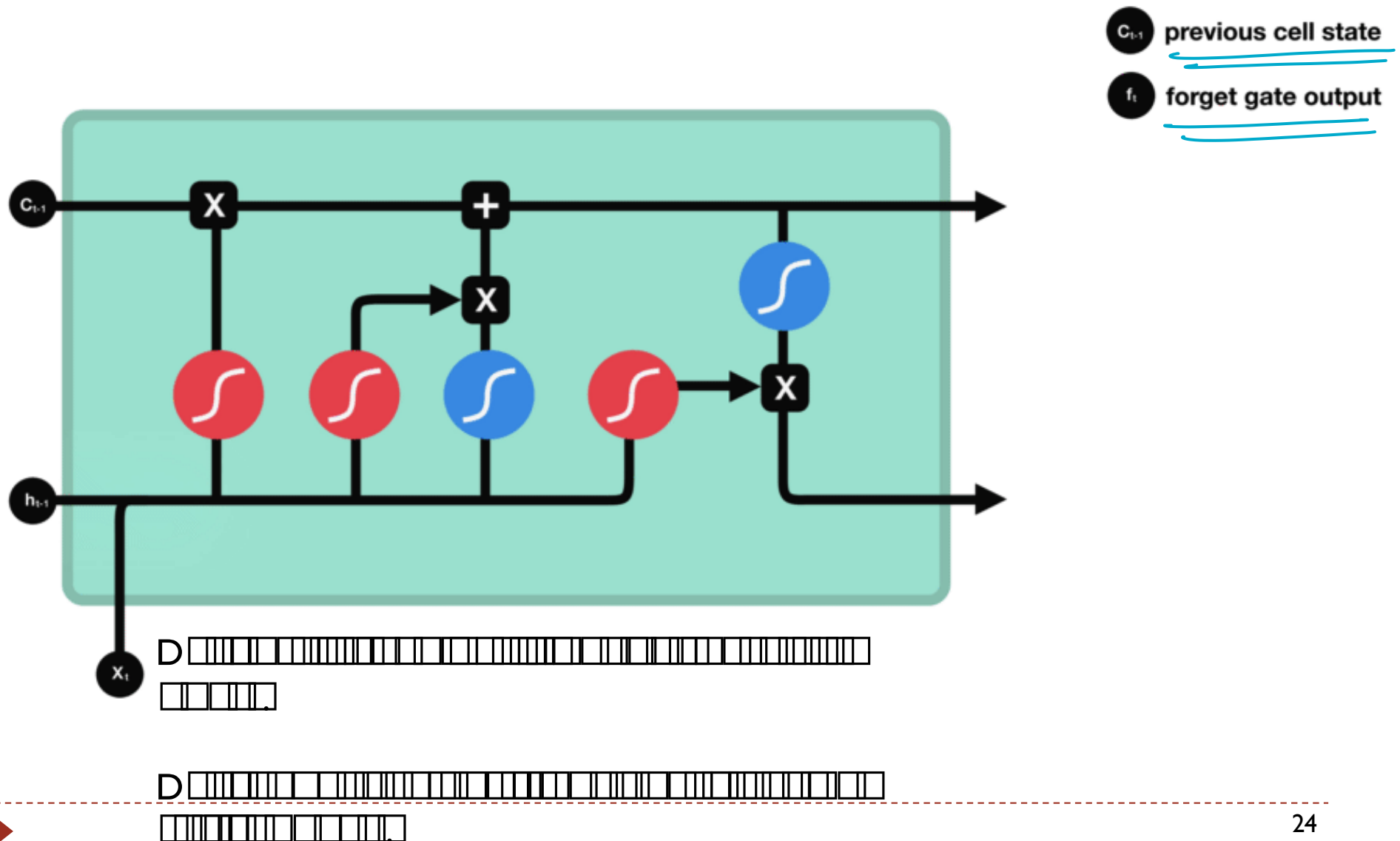Resolve pb1 → add a kind of residual chain

Resolve pb2 → control the memory

Residual



LSTM = *Long Short-Term Memory*

*(crédit : image modifiée de Michaël Nguyen)*

# LSTM cell

▸ **Contain three "gates":**

  ▸ Calculation zones that regulate the flow of information (by performing specific actions).

  ▸ Forget gate (porte d'oubli)

  ▸ Input gate (porte d'entrée)

  ▸ Output gate (porte de sortie)

▸ **$h_t$: Hidden state (état caché)**

  ▸ The eventual output

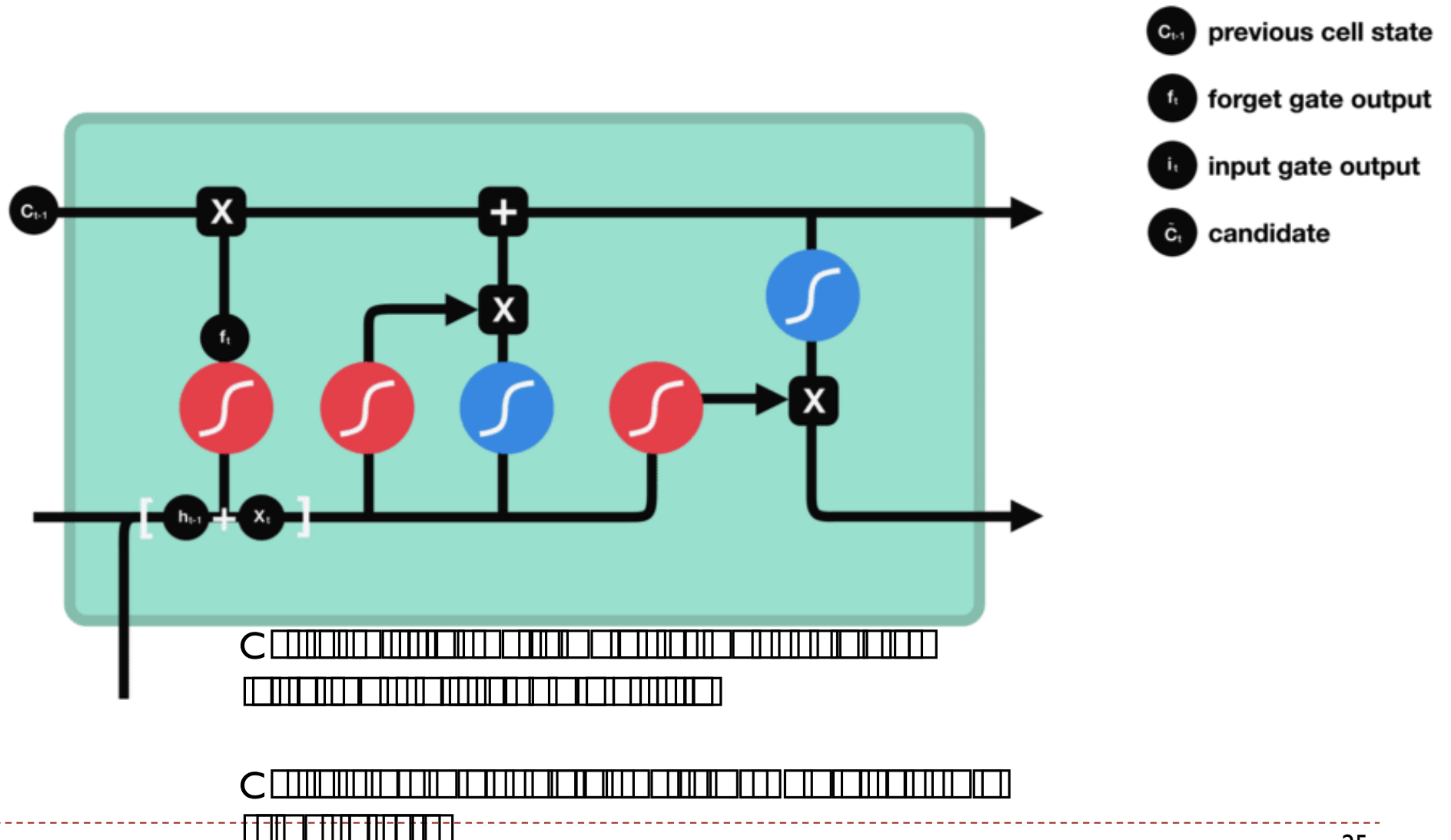▸ **$c_t$: Cell state (état de la cellule**

  ▸ Like residual



*(crédit : image modifiée de Michaël Nguyen)*
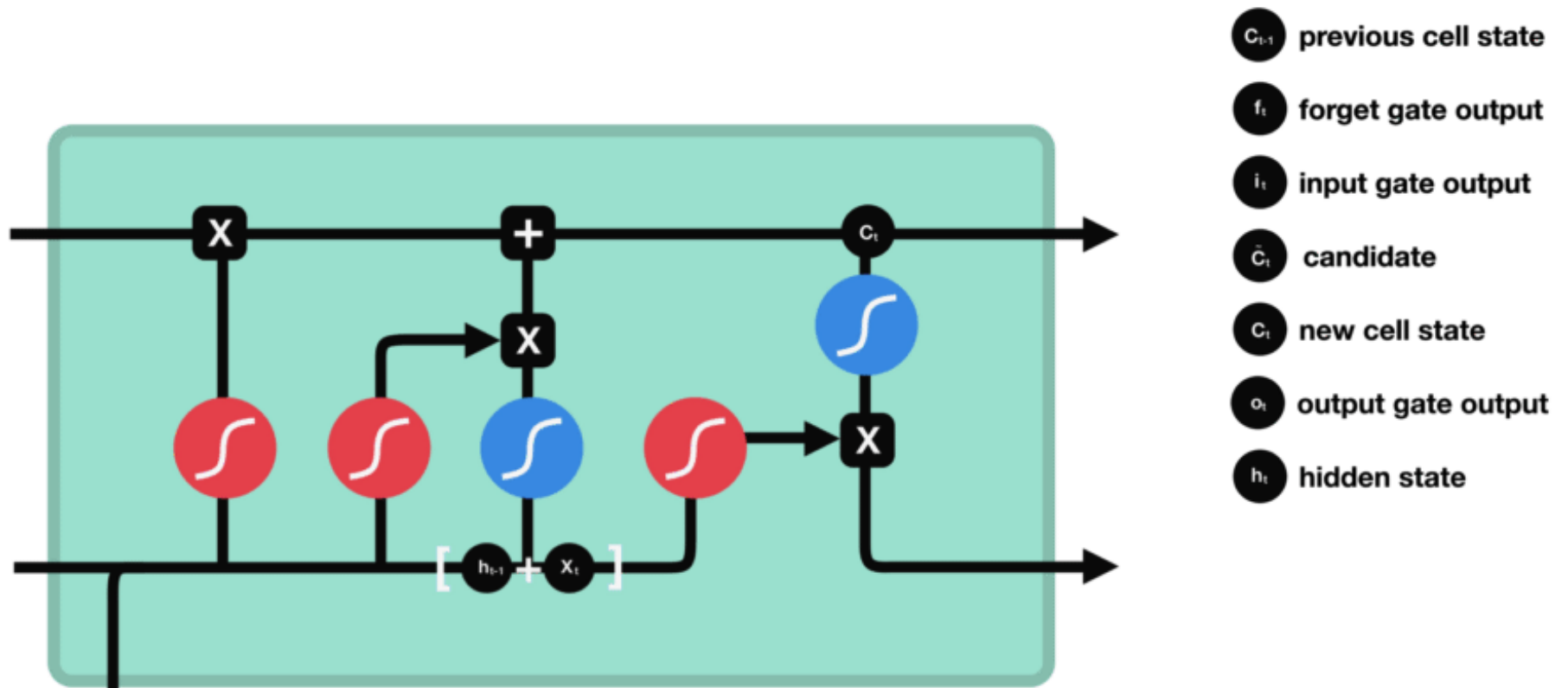
# LSTM cell (porte oubli / forget get)



Décide si l'état de la cellule doit être conservé ou non.

Decides whether the state of the cell should be retained or not.

# LSTM cell (porte entrée / input get)



Capture l'information d'entrée qui doit être incluse dans l'état de la cellule

Captures the input information to be included in the cell state

# LSTM cell (état de la cellule / cell state)



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

L'état de la cellule se calcule assez simplement à partir de la porte d'oubli et de la porte d'entrée.

The state of the cell is calculated quite simply from the forget gate and the input gate
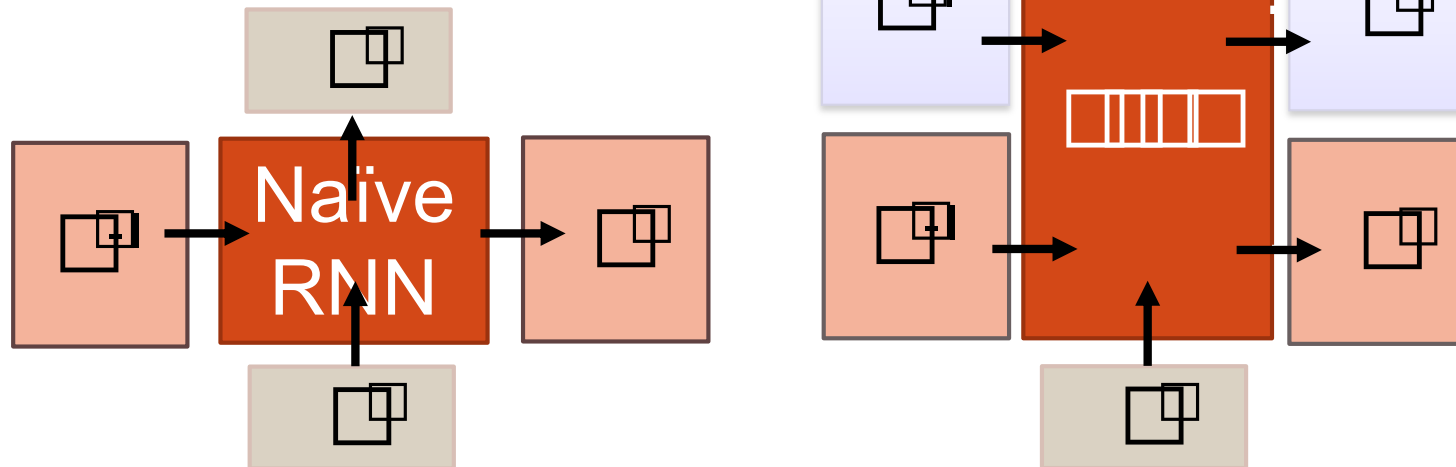
26

# LSTM cell (porte de sortie / output gate)



La porte de sortie décide quel sera le prochain état caché. Il contient des informations sur les entrées précédentes du réseau et sert aux prédictions.

The output gate decides what the next hidden state will be. It contains information about previous inputs to the network and is used for predictions.

27

# Naïve RNN vs LSTM

▸ Naïve RNN
  ▸ Reuse at each step the previous Output

▸ LSTM
  ▸ At each step 3 gate control the use use of Input value, Cell state and previous Output

$h^t$

$c^{t-l}$  LSTM  $c^t$

$h^{t-l}$  $h^t$

$x^t$

$h^t$

$h^{t-l}$  Naïve RNN  $h^t$

$x^t$

c changes slowly ➡ $c^t$ is $c^{t-1}$ added by something

h changes faster ➡ $h^t$ and $h^{t-1}$ can be very different

# GRU – gated recurrent unit

▸ GRU = a light LSTM Cell



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

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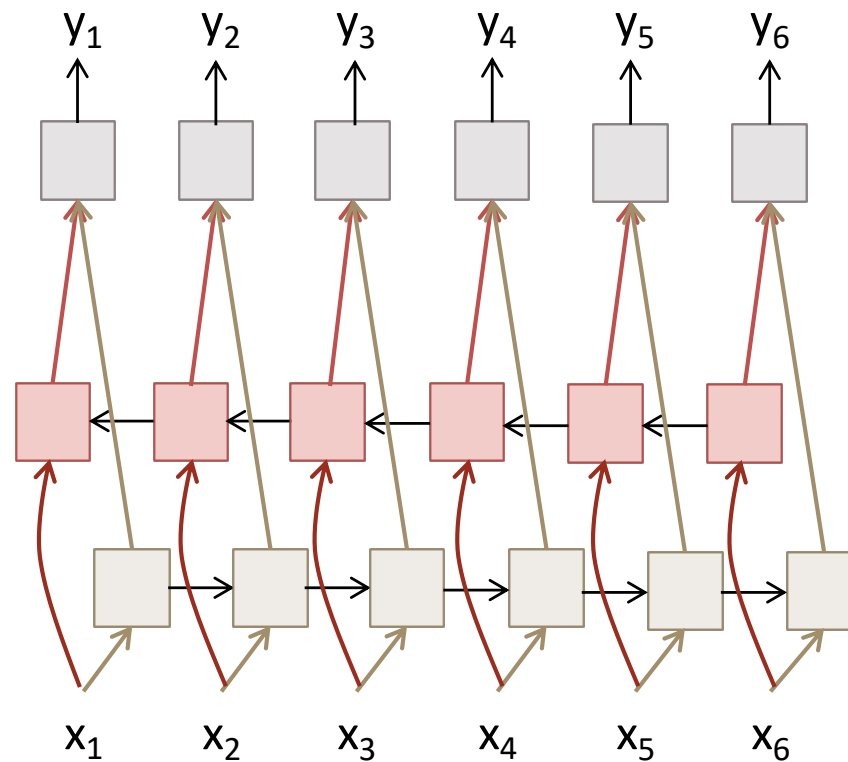
$$\tilde{h}_t = \tanh\left(W \cdot [r_t * h_{t-1}, x_t]\right)$$

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

- It combines the forget and input into a single update gate.
- It also merges the cell state and hidden state.
→ This is simpler/faster than LSTM.

# Bi-directional RNNs

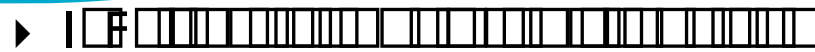- RNNs can process the input sequence in forward and in the reverse direction



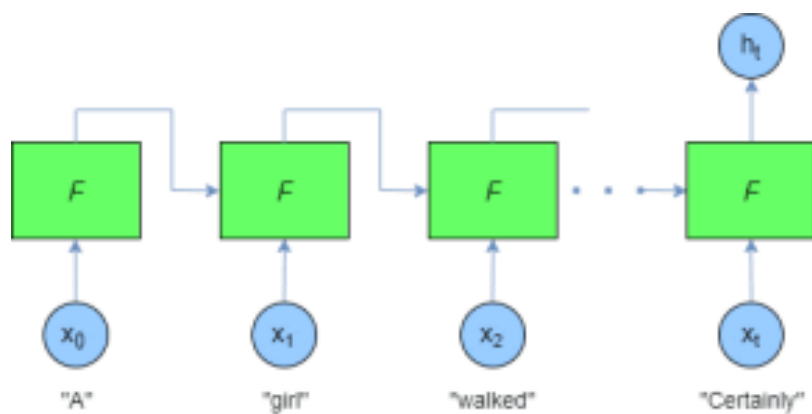- Popular in speech recognition, could be used also with text
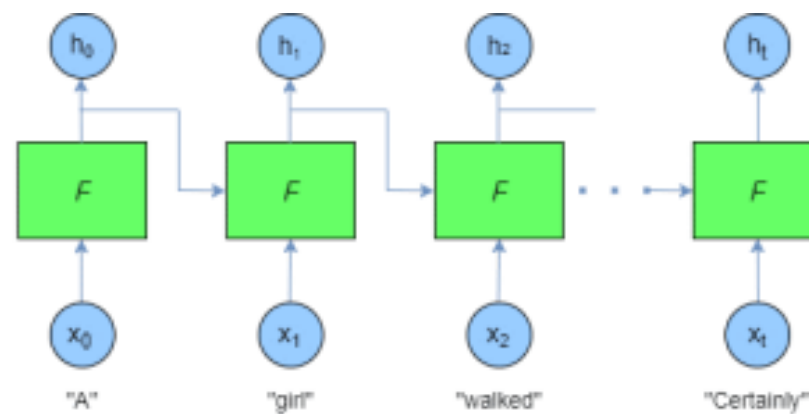
# RNN cell in Keras

# Keras Long Short-Term Memory Cell

▸ from tf.keras.layers import LSTM

▸ Main params
  ▸ Units: dimension of output space

  ▸ return_sequences: True or False
    ▸ If False return only the last output
    ▸ If True return the full sequence of the output sequence
      ▫ Output sequence = hidden state (the vocabulary change regardind documentation)
  ▸ return_state:True or False
    ▸ If True return 3 values
      ▫ The full output sequence or only the last one (depend on return_sequences)
      ▫ The last output sequence
      ▫ The cell state
    ▸ If False return nothing

return_sequences = False

return_sequences = True

# If there's a Dense after LSTM. How many dense cells are used?

▸ X = LSTM(500, return sequence = False or True)

▸ Output = Dense(…)(X)

With return_sequences=False,
        Dense layer is applied only once at the last cell

With return_sequences=True
        Dense layer is applied to every timestep

# A basic example: forcasting

inputs = Input(shape=(None,I))   # Input size are unknow
  ▸ Output shape=[None, None, I] → *One Features.*
output = LSTM(16, return_sequences=False)(embedding)
  ▸ Output shape=[None, 16]
predictions = Dense(I, activation=linear')(output)
  ▸ Output shape=[None, I] → *Regression.*

▸ **Fit by batch**
  ▸ Model.fit(X, y, ….). ← all item have the same length

▸ **Fit by item**
  ▸ For  i in range(len(X)): ← could be different length
    ▸ Model.fit(X[i], y[i], …)

# RNN for forcasting

# RNN for forecasting

# RNN for forecasting

- A model that makes a prediction
  - one hour into the future,
  - given six hours of history

# To make a single prediction

- 24 hours into the future,
- given 24 hours of history

# Dense model



```
multi_step_dense = tf.keras.Sequential([
    tf.keras.layers.Flatten(),
    # Shape: (time, features) => (time*features)

    tf.keras.layers.Dense(units=32, activation='relu'),
    tf.keras.layers.Dense(units=1),
])
```

# Recurrent model: return_state=False



```
lstm_model = tf.keras.models.Sequential([

    tf.keras.layers.LSTM(32, return_sequences=False),
    # Shape [batch, time, features] => [batch, lstm_units]

    tf.keras.layers.Dense(units=1)
])
```

# Recurrent model: return_state=True



lstm_model = tf.keras.models.Sequential([

   tf.keras.layers.LSTM(32, return_sequences=True),
   # Shape [batch, time, features] => [batch, time, lstm_units]

   tf.keras.layers.Dense(units=1)
   # One Dense by time step
])

# How to predict futur multiple values



Not possible to build the network in few line of code :

Seq2Seq model in another lecture
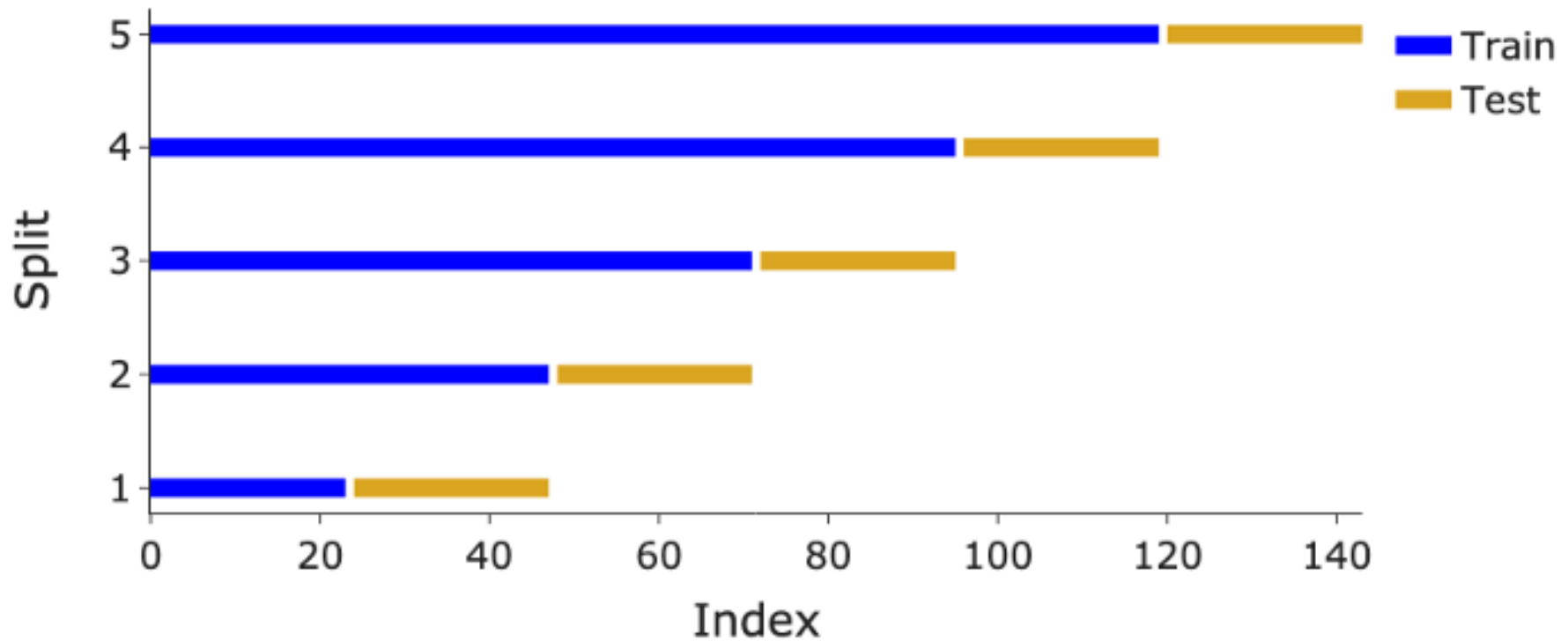
# Time series cross-validation

from sklearn.model_selection import TimeSeriesSplit

# RNN in NLP task

# Some use of RNN
## → Text Classification / Sentiment analysis

**A**ffect a label to a text

▶ Classify a
  ▶ restaurant review from Yelp!
  ▶ movie review from IMDB
    …
    as positive or negative

▶ Inputs:
  ▶ Multiple words, one or more sentences

▶ Outputs:
  ▶ Positive / Negative classification

▶ "The food was really good"

▶ "The chicken crossed the road because it was uncooked"

# Sentiment analysis - solution 1

▶ retrieve only the last state

# Sentiment analysis – solution2

▸ Other possible architecture

  ▸ Average all the internal state

# Some use of RNN
## → Named Entity Recognition / Part of Speech Tagging

▸ Affect a label to each word

  ▸ find and classify names in text

    ▸ Could bean entity : number, country, person, … (NER)

    ▸ Could be a function : noun, verb, adverbs, … (POS)

In fact, the [Chinese NORP] market has the [three CARDINAL] most influential names of the retail and tech space – [Alibaba GPE] , [Baidu ORG] , and [Tencent PERSON] (collectively touted as [BAT ORG] ), and is betting big in the global [AI GPE] in retail industry space . The [three CARDINAL] giants which are claimed to have a cut-throat competition with the [U.S. GPE] (in terms of resources and capital) are positioning themselves to become the 'future [AI PERSON] platforms'. The trio is also expanding in other [Asian NORP] countries and investing heavily in the [U.S. GPE] based [AI GPE] startups to leverage the power of [AI GPE] . Backed by such powerful initiatives and presence of these conglomerates, the market in APAC AI is forecast to be the fastest-growing [one CARDINAL] , with an anticipated [CAGR PERSON] of [45% PERCENT] over [2018 - 2024 DATE] .

To further elaborate on the geographical trends, [North America LOC] has procured [more than 50% PERCENT] of the global share in [2017 DATE] and has been leading the regional landscape of [AI GPE] in the retail market. The [U.S. GPE] has a significant credit in the regional trends with [over 65% PERCENT] of investments (including M&As, private equity, and venture capital) in artificial intelligence technology. Additionally, the region is a huge hub for startups in tandem with the presence of tech titans, such as [Google ORG] , [IBM ORG] , and [Microsoft ORG] .

# Extract information from text



Vente Villa 4 pièces Nice (06000)
Réf. 12390: Sur les Hauteurs de Nice. Superbe villa moderne (190m2), 2 chambres et 1 suite parentale, 3 salles de bain. Très grand salon/salle à manger, cuisine américaine équipée. Prestations de haut standing. Vue panoramique sur la mer. Cette villa a été construite en 2005. 1 270 000 euros. Si vous êtes intéressés, contactez vite Mimi LASOURIS 06.43.43.43. 43

**REAL ESTATE TEMPLATE**
Reference: 12390
Prize: 1 270 000
Surface: 190 m2
Year Built: 2005
Rooms: 4
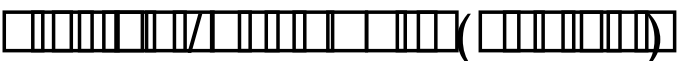Owner: Mimi LASOURIS
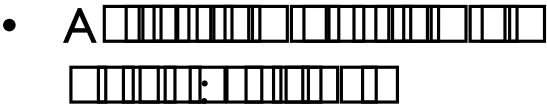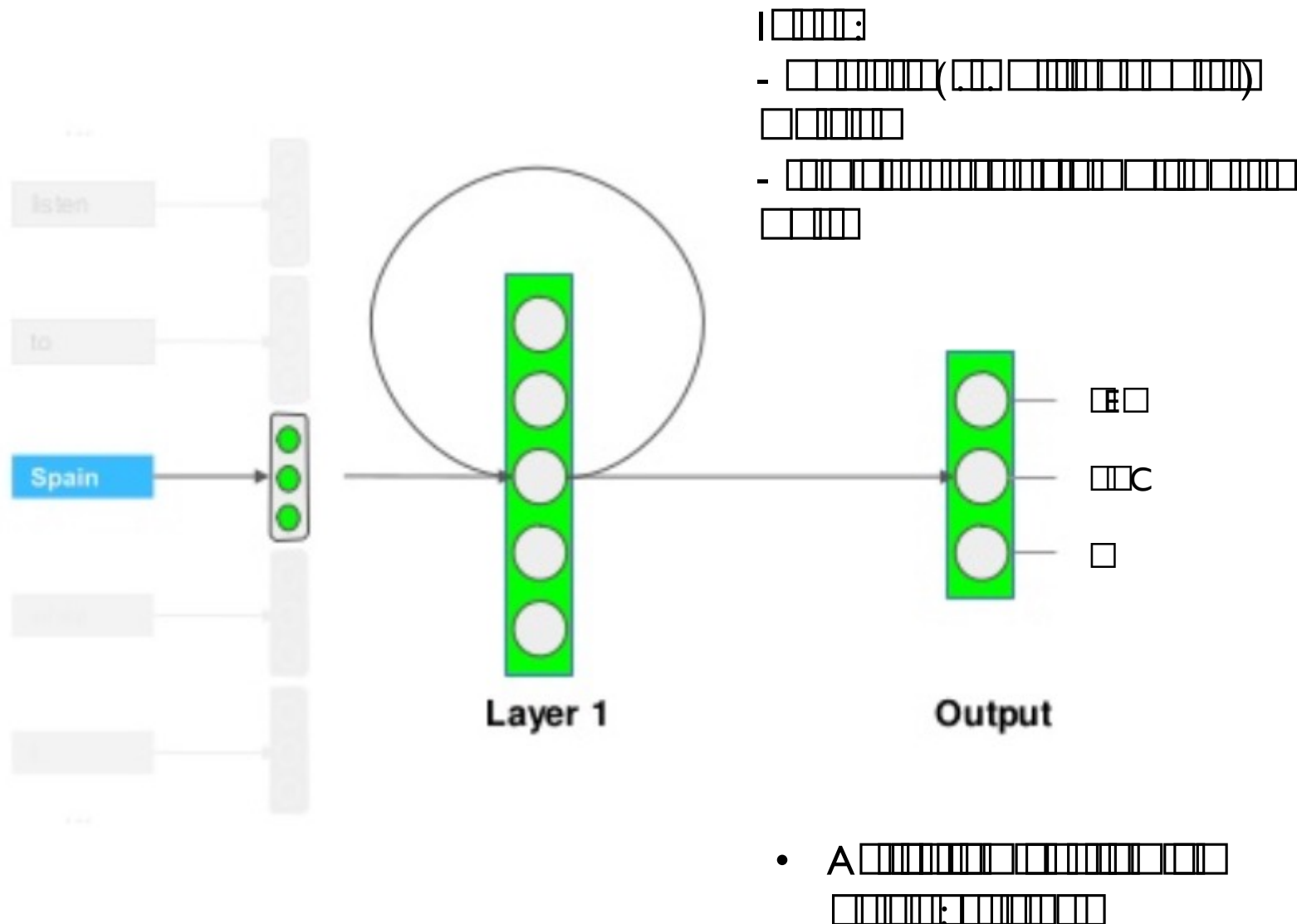Telephone: 06.43.43.43. 43

# Machine Learning approach

- Features for sequence labeling
  - Words
    - Current word (essentially like a learned dictionary)
    - Previous/next word (context)
    - Other kinds of inferred linguistic classification – Part-of-speech tags

  - Label context
    - Previous (and perhaps next) label

Input :
- a word and its context (the words before and after it)
Output
- the label of the central word



- Activation function for output: softmax

# Recurrent neural network for NER

Input :
- a phrase (i.e. a list of words)
Output
- the label associated with each word



Layer 1

Output

PER

LOC

…

- Activation function for output: softmax

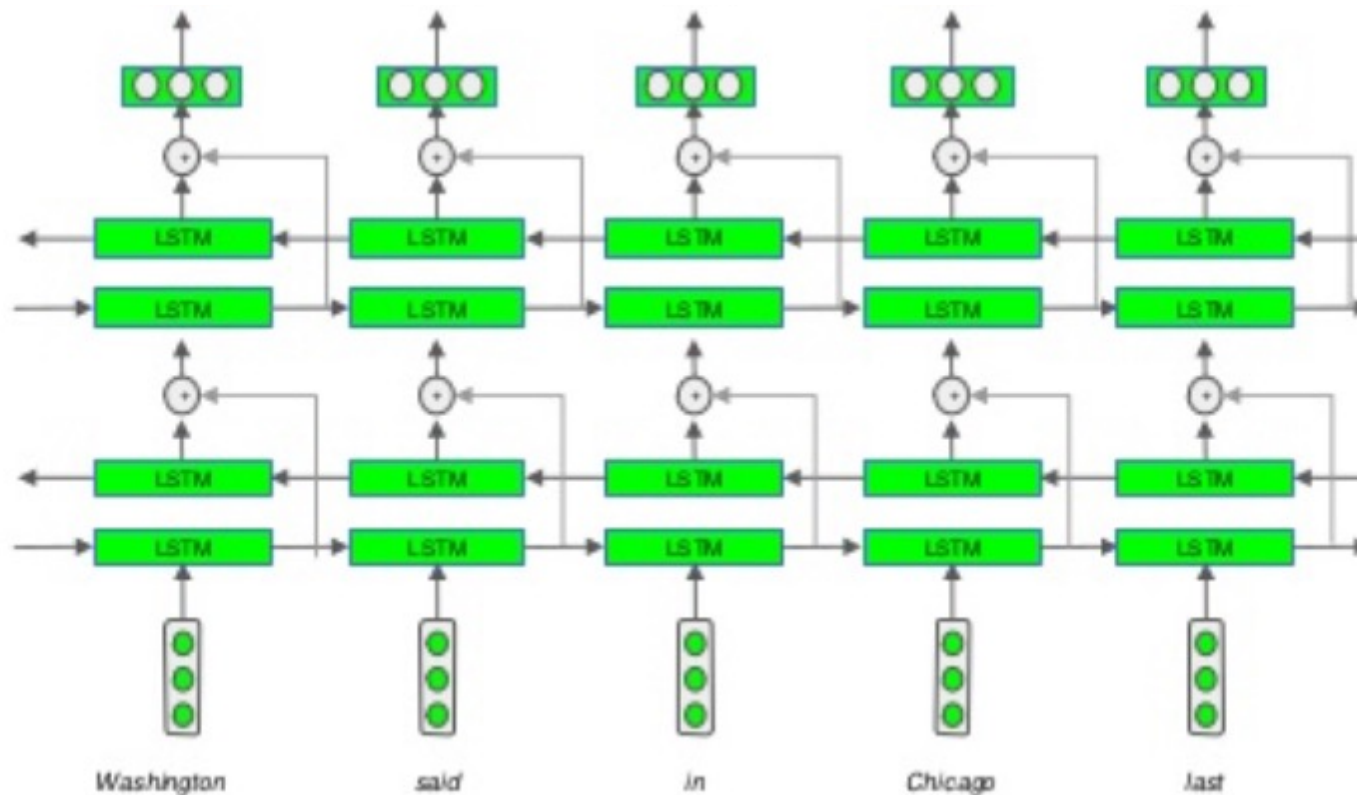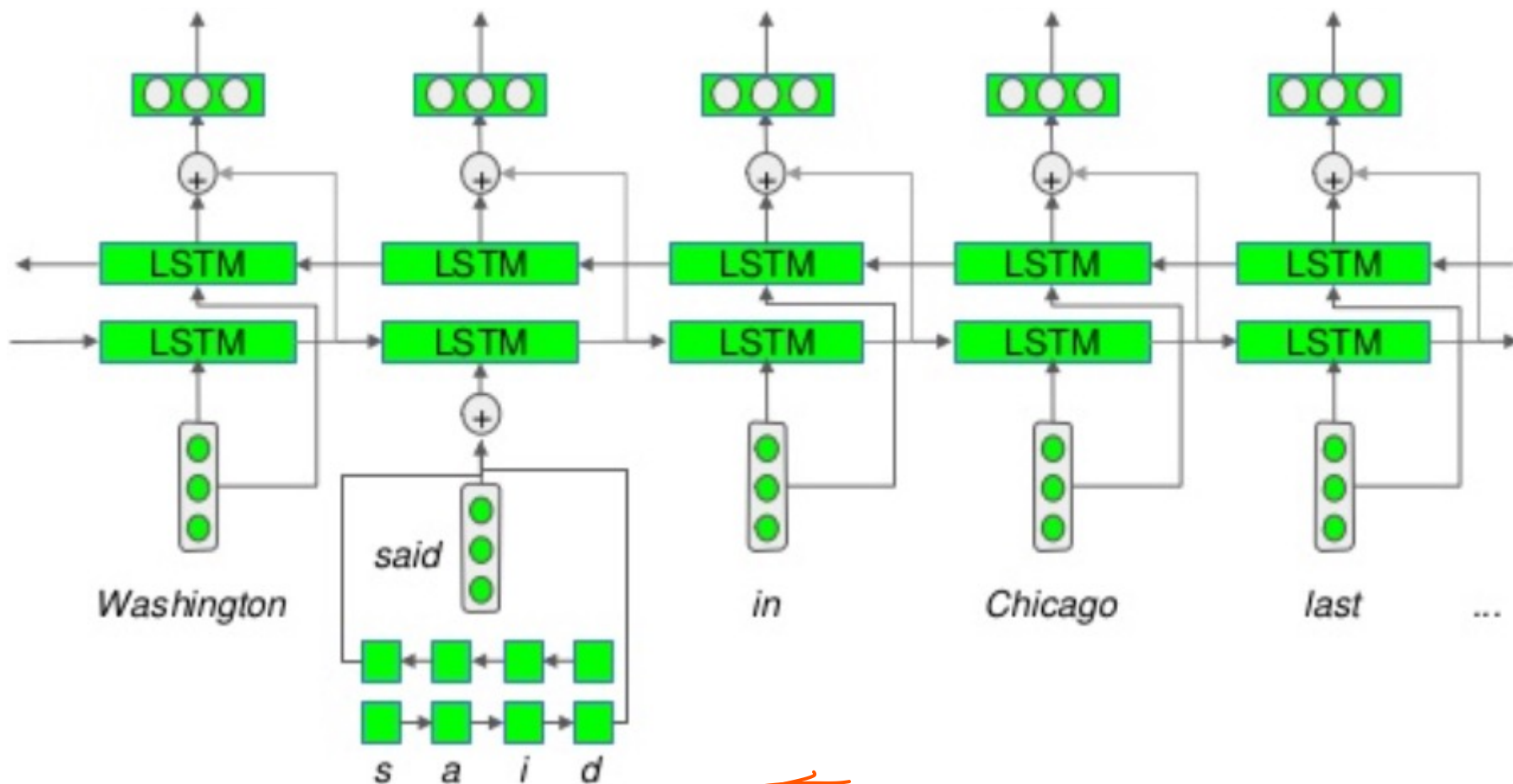# Recurrent neural network for NER (same network but unfolded)



- Activation function for output: softmax

# Bi directional recurrent neural network for NER



- Activation function for output: softmax

# Stacked Bi-RNN



- Activation function for output: softmax

# Multi-level encoding
## char encoding + word encoding

# Today Lab

For all students

▸ Building a parameterizable RNN for time series forecasting
  ▸ The length of the time series must not be known when the network is built
  ▸ The number of input features is known (n_in)
  ▸ The number of features to be predicted is known (n_out)

  ▸ Input shape of the network: (None, None, n_in)
  ▸ Output shape of the network: (None, I, n_out) or (None, n_out) ~~Build a Seq2Seq parameterizable model~~

▸ Fit and use the model for forecasting
  ▸ AvailableData set available
    ▸ filename = "http://www.i3s.unice.fr/~riveill/dataset/precipitation.csv.zip"
    ▸ df = pd.read_csv(filename, sep="\t", engine="python", on_bad_lines="skip")

  ▸ Use 24 months in the past to predict the next month
  ▸ Use cross validation with 5 splits (sklearn.model_selection.TimeSeriesSplit)
  ▸ Use MeanSquaredLogarithmicError to evaluate your network