



PONTIFICIA  
UNIVERSIDAD  
CATÓLICA  
DE CHILE

# IIC 3800 Tópicos en CC NLP

<https://github.com/marcelomendoza/IIC3800>

# Conceptos de NLP

## POS tagging

- ▶ Etiquetar cada término de acuerdo a la función que este cumple en el texto.
- ▶ Puede ayudarnos en tareas como detección de estilo, parsing, detección de colocaciones.
- ▶ Tarea importante en NLP.

Text:

John likes the blue house at the end of the street .

Adjective	Determiner	Preposition
Adverb	Noun	Pronoun
Conjunction	Number	Verb

# Conceptos de NLP

## POS tagging

	Tag	Description	Example
Open Class	<b>ADJ</b>	Adjective: noun modifiers describing properties	<i>red, young, awesome</i>
	<b>ADV</b>	Adverb: verb modifiers of time, place, manner	<i>very, slowly, home, yesterday</i>
	<b>NOUN</b>	words for persons, places, things, etc.	<i>algorithm, cat, mango, beauty</i>
	<b>VERB</b>	words for actions and processes	<i>draw, provide, go</i>
	<b>PROPN</b>	Proper noun: name of a person, organization, place, etc..	<i>Regina, IBM, Colorado</i>
	<b>INTJ</b>	Interjection: exclamation, greeting, yes/no response, etc.	<i>oh, um, yes, hello</i>
Closed Class Words	<b>ADP</b>	Adposition (Preposition/Postposition): marks a noun's spacial, temporal, or other relation	<i>in, on, by under</i>
	<b>AUX</b>	Auxiliary: helping verb marking tense, aspect, mood, etc.,	<i>can, may, should, are</i>
	<b>CCONJ</b>	Coordinating Conjunction: joins two phrases/clauses	<i>and, or, but</i>
	<b>DET</b>	Determiner: marks noun phrase properties	<i>a, an, the, this</i>
	<b>NUM</b>	Numeral	<i>one, two, first, second</i>
	<b>PART</b>	Particle: a preposition-like form used together with a verb	<i>up, down, on, off, in, out, at, by</i>
	<b>PRON</b>	Pronoun: a shorthand for referring to an entity or event	<i>she, who, I, others</i>
	<b>SCONJ</b>	Subordinating Conjunction: joins a main clause with a subordinate clause such as a sentential complement	<i>that, which</i>
Other	<b>PUNCT</b>	Punctuation	<i>! , ()</i>
	<b>SYM</b>	Symbols like \$ or emoji	<i>\$, %</i>
	<b>X</b>	Other	<i>asdf, qwfg</i>

## Conceptos de NLP

### POS tagging en NLTK

```
>>> text = word_tokenize("And now for something completely different")
>>> nltk.pos_tag(text)
[('And', 'CC'), ('now', 'RB'), ('for', 'IN'), ('something', 'NN'),
 ('completely', 'RB'), ('different', 'JJ')]
```

```
>>> text = word_tokenize("They refuse to permit us to obtain the refuse permit")
>>> nltk.pos_tag(text)
[('They', 'PRP'), ('refuse', 'VBP'), ('to', 'TO'), ('permit', 'VB'), ('us', 'PRP'),
 ('to', 'TO'), ('obtain', 'VB'), ('the', 'DT'), ('refuse', 'NN'), ('permit', 'NN')]
```

# Conceptos de NLP

## POS tagging en Spacy (Español)

```
!python -m spacy download es_core_news_sm
!python -m spacy download es_core_news_md
```

```
import spacy
import es_core_news_sm
import es_core_news_md

nlp = es_core_news_sm.load()
doc = nlp('''Un desastroso espíritu posee tu tierra:
            donde la tribu unida blandió sus mazas,
            hoy se enciende entre hermanos perpetua guerra,
            se hieren y destrozan las mismas razas.''' )

for token in doc: print(token.text, "|", token.lemma_, '|', token.pos_)
```

Un | Un | DET  
desastroso | desastroso | NOUN  
espíritu | espíritu | PROPN  
posee | poseer | VERB  
tu | tu | DET  
tierra | tierra | NOUN  
: | : | PUNCT

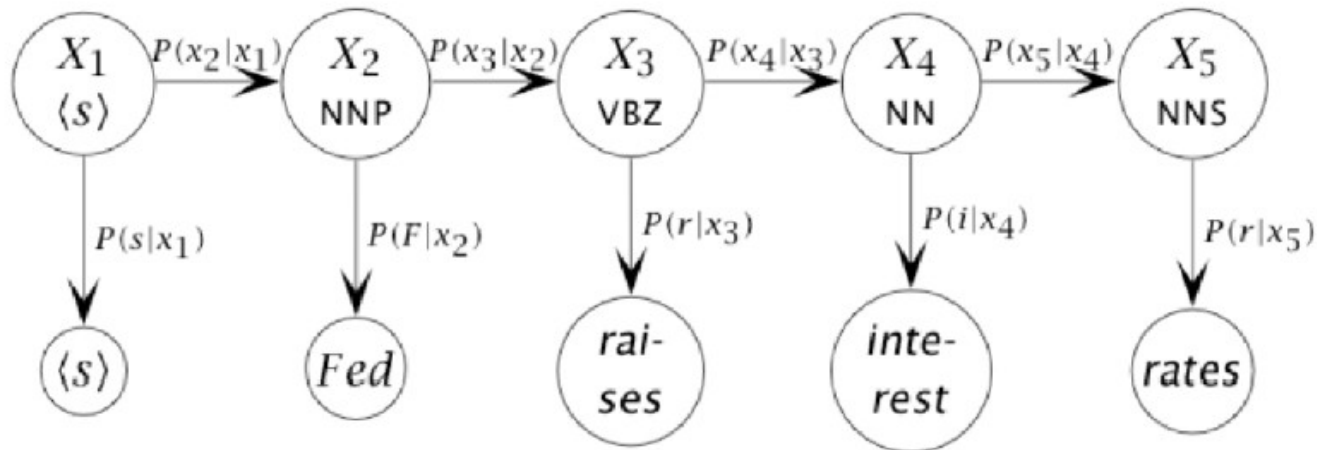
## Conceptos de NLP

### POS tagging con HMM (modelo clásico)

- ▶ Se dispone de un corpus etiquetado.
- ▶ La secuencia de tags es interpretada como una cadena de Markov:  
 $P(x_{t+1} \mid x_t, \dots, x_1) = P(x_{t+1} \mid x_t)$ ,  $x_1, \dots, x_{t+1}$  representan tags
- ▶ Usamos un modelo generativo para términos, con tags como estados ocultos:  $P(t \mid x_1, \dots, x_{t+1}) = P(t \mid x_{t+1})$

## Conceptos de NLP

### POS tagging con HMM (modelo clásico)



- En general muestran buena precisión (sobre 90 %).

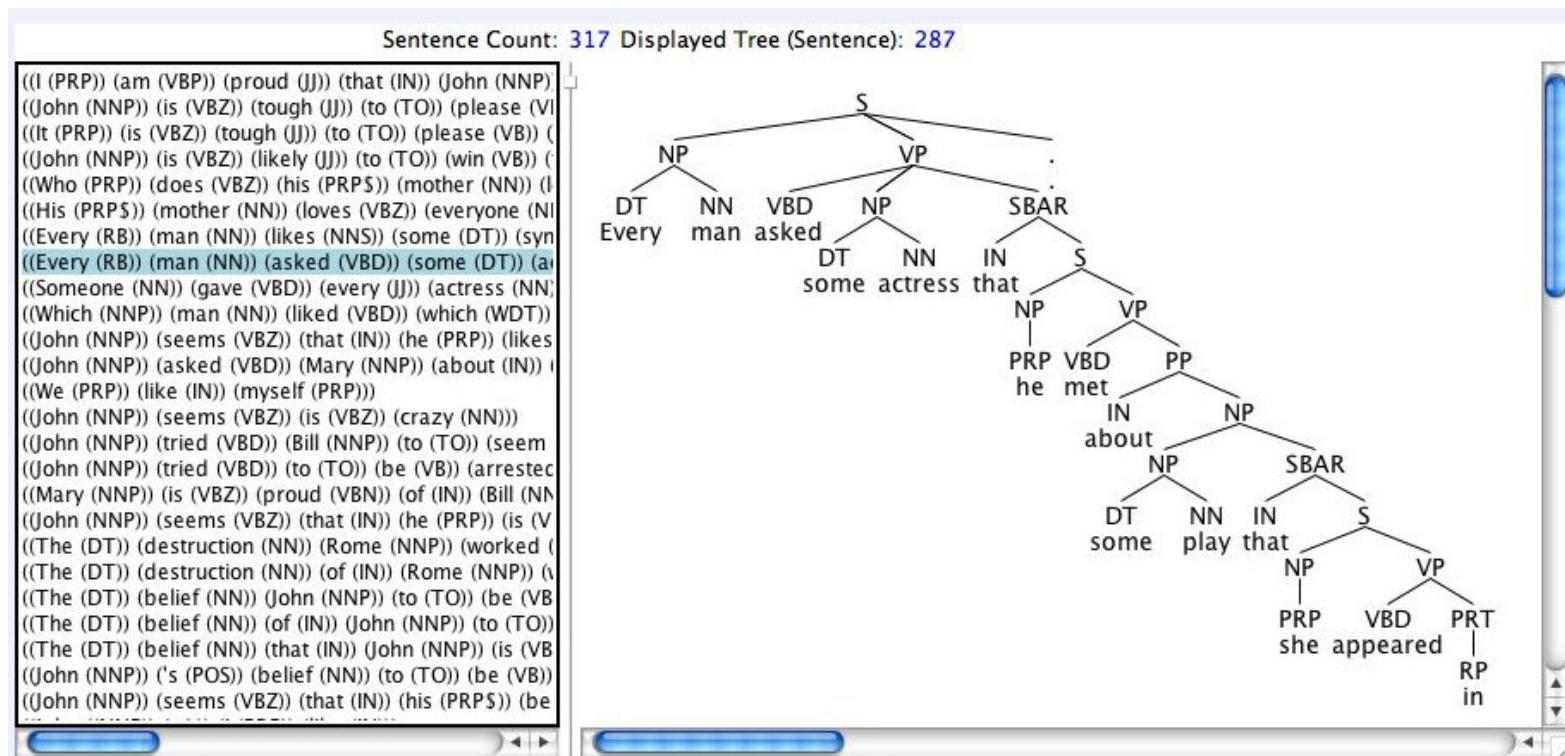
El algoritmo de entrenamiento se llama Viterbi.

# Conceptos de NLP

## POS tagging ¿Cuáles datos usan?

Treebanks: [Penn treebank](#) (más famoso), [UAM Spanish Treebank](#), ...

[Treebank viewer](#):

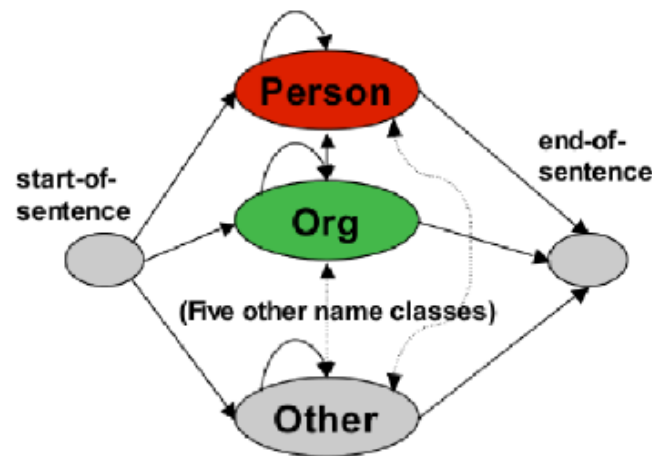




# Conceptos de NLP

## Named Entity Recognition

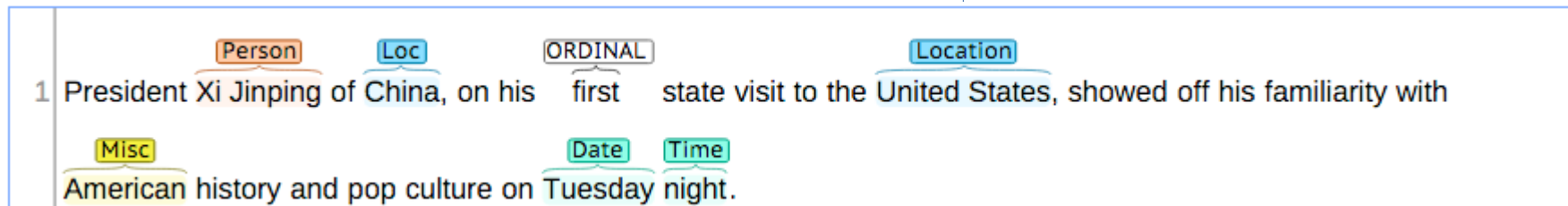
- ▶ Tarea: Identificar entidades en texto (personas, organizaciones, etc.)
- ▶ Separa el text en chunks, y para cada cual asocia una NE. Opera sobre texto tagged.
- ▶ NER types: organization, person, location, date, time, money, percent, facility (human made artifacts), gpe (geo-political ents).
- ▶ POS tagging puede ayudar, agregando *entity* como un estado mas.



# Conceptos de NLP

## Named Entity Recognition

bi-grama



Named Entity Recognition puede ser muy desafiante:

Back in 2000 , People Magazine PUBLISHER highlighted Prince Williams' PERSON style who at the time was a little more fashion-conscious , even making fashion statements at times .

Now-a-days the prince mainly wears navy COLOR suits ITEM ( sometimes double-breasted DESIGN ) , light blue COLOR button-ups ITEM with classic LOOK pointed DESIGN collars PART , and burgundy COLOR ties ITEM .

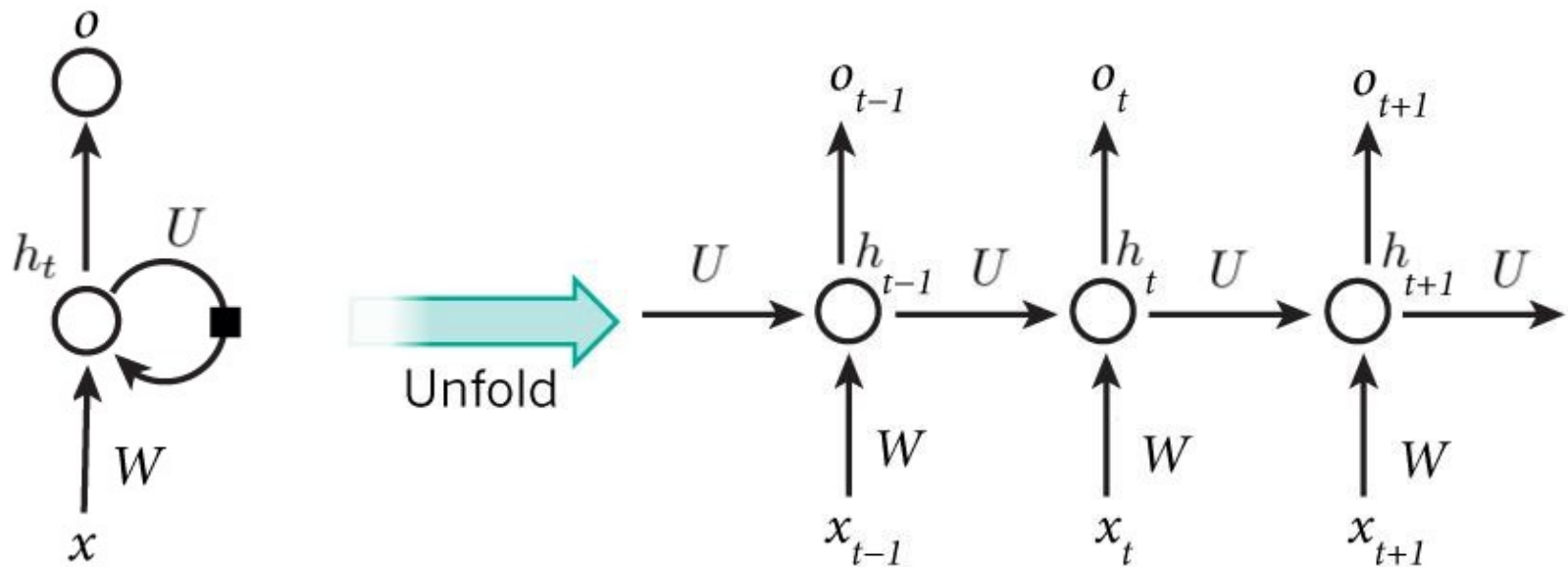
But who knows what the future holds ...

Duchess Kate PERSON did wear an Alexander McQueen BRAND dress ITEM to the wedding OCCASION in the fall of 2017 SEASON .

## Sequence labeling

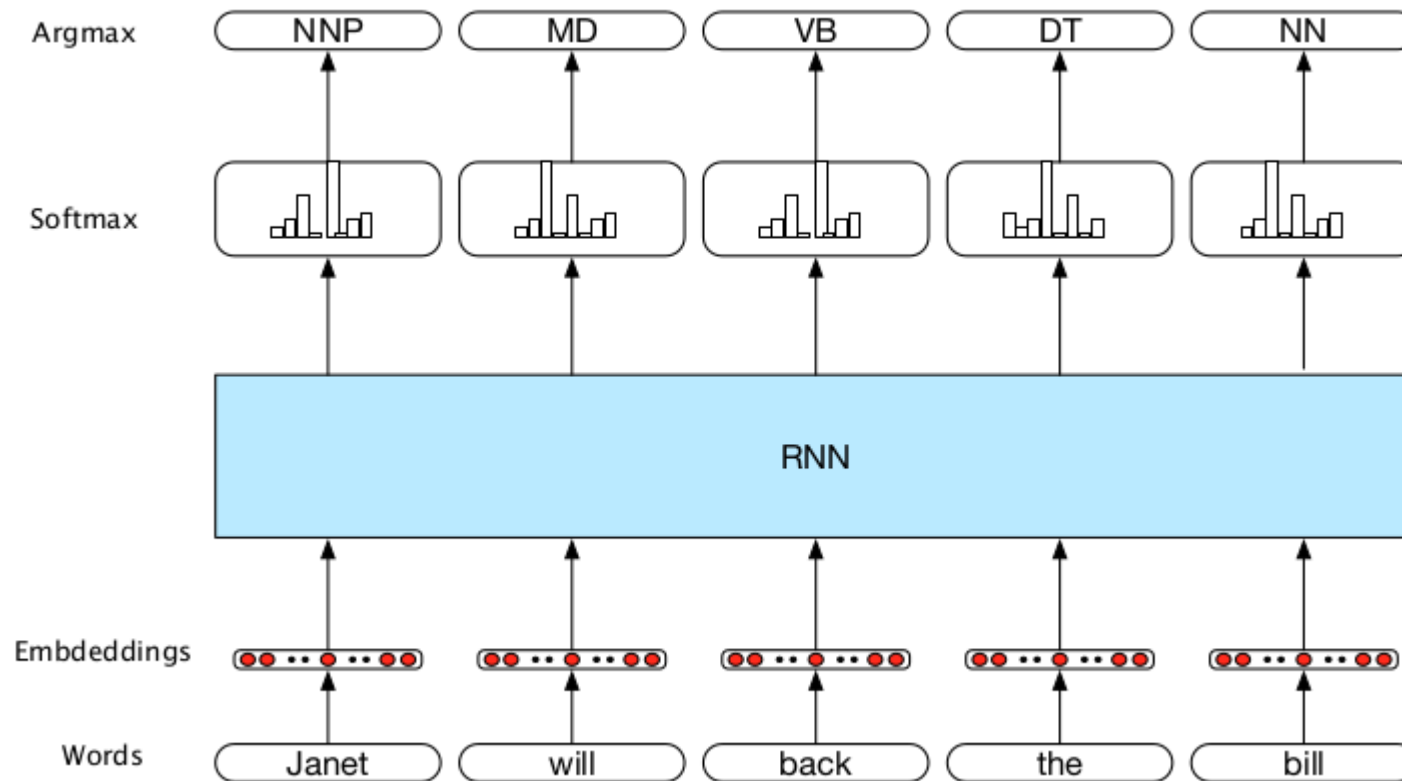
Presentamos los datos como secuencia:  $X = (x_1, x_2, \dots, x_T)$

Recurrente convencional:  $h_t = g(W \cdot x_t + U \cdot h_{t-1})$



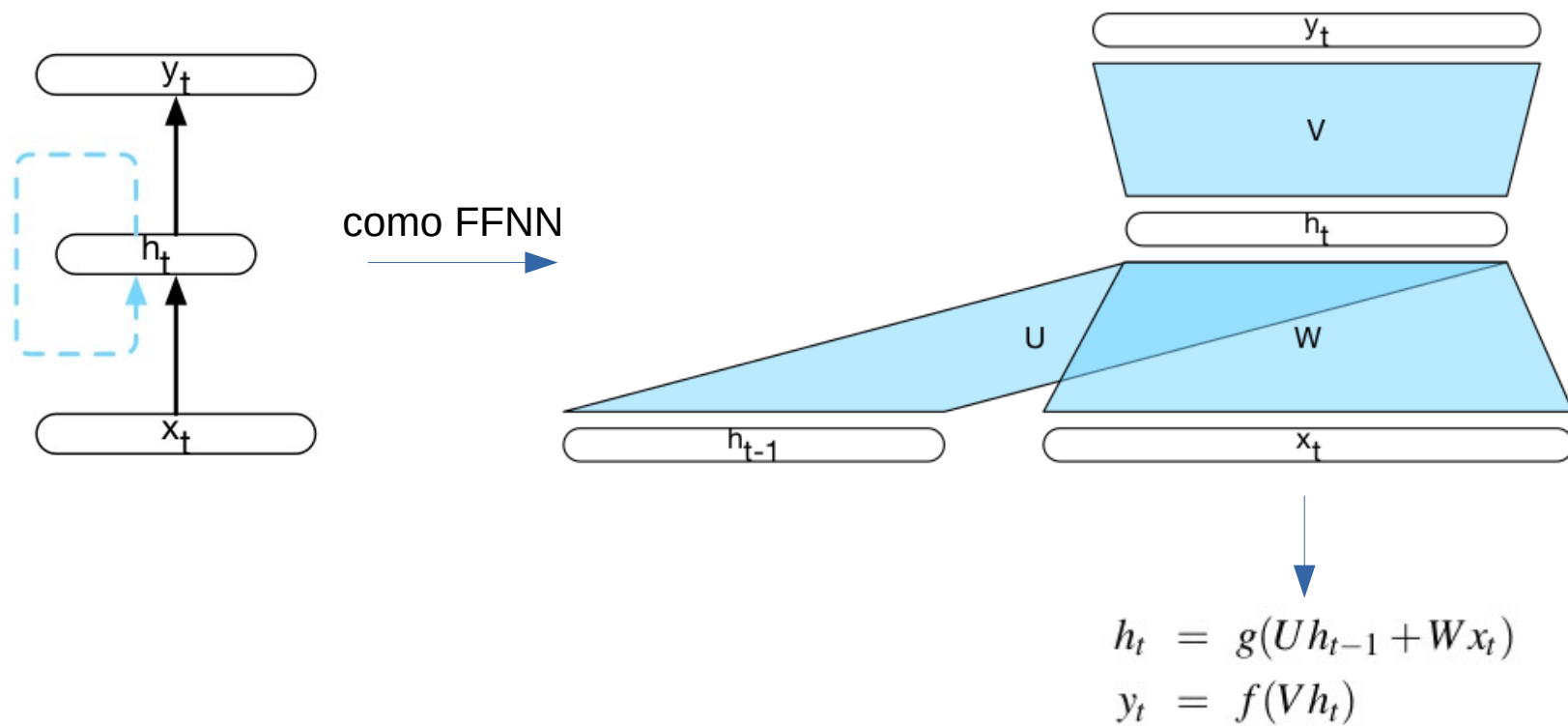
# Sequence labeling

Red recurrente (**sequence labeling**):



# Sequence labeling

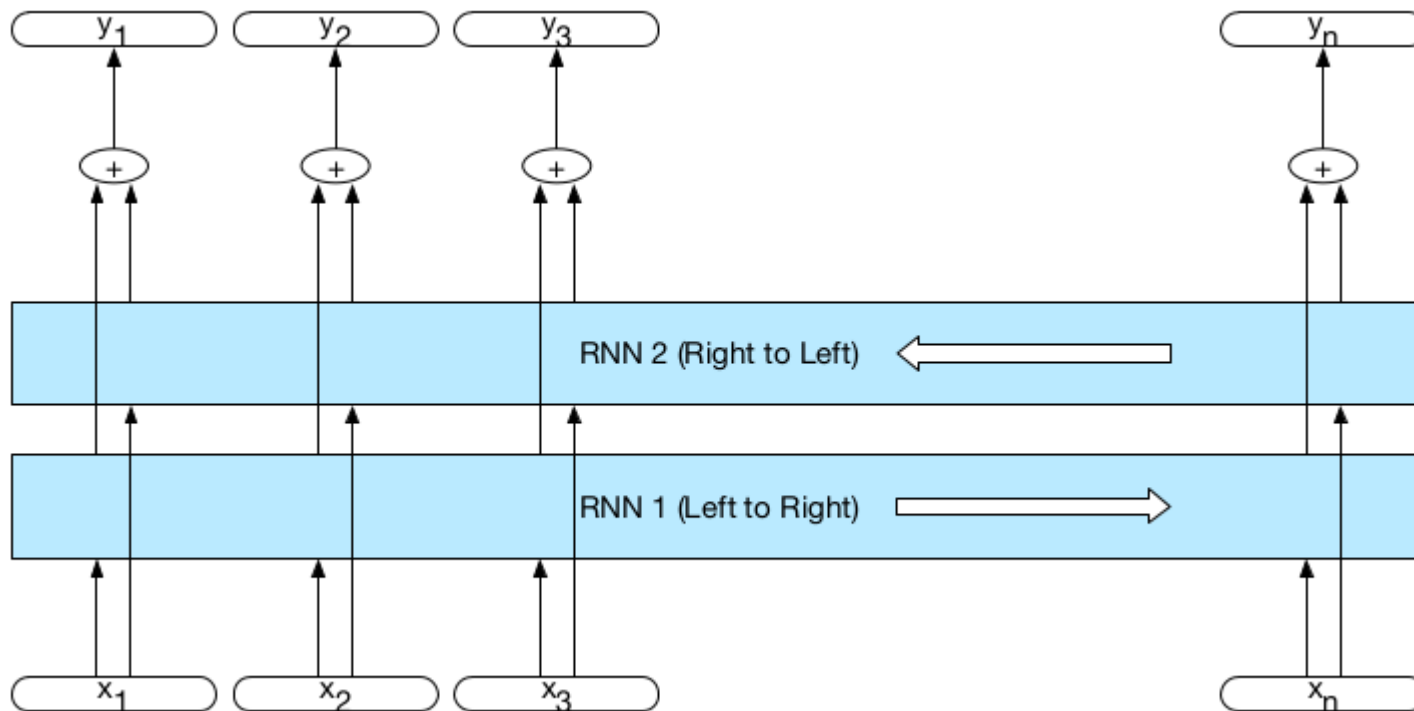
Red recurrente:



$$y_t = \text{softmax}(Vh_t)$$

## Sequence labeling

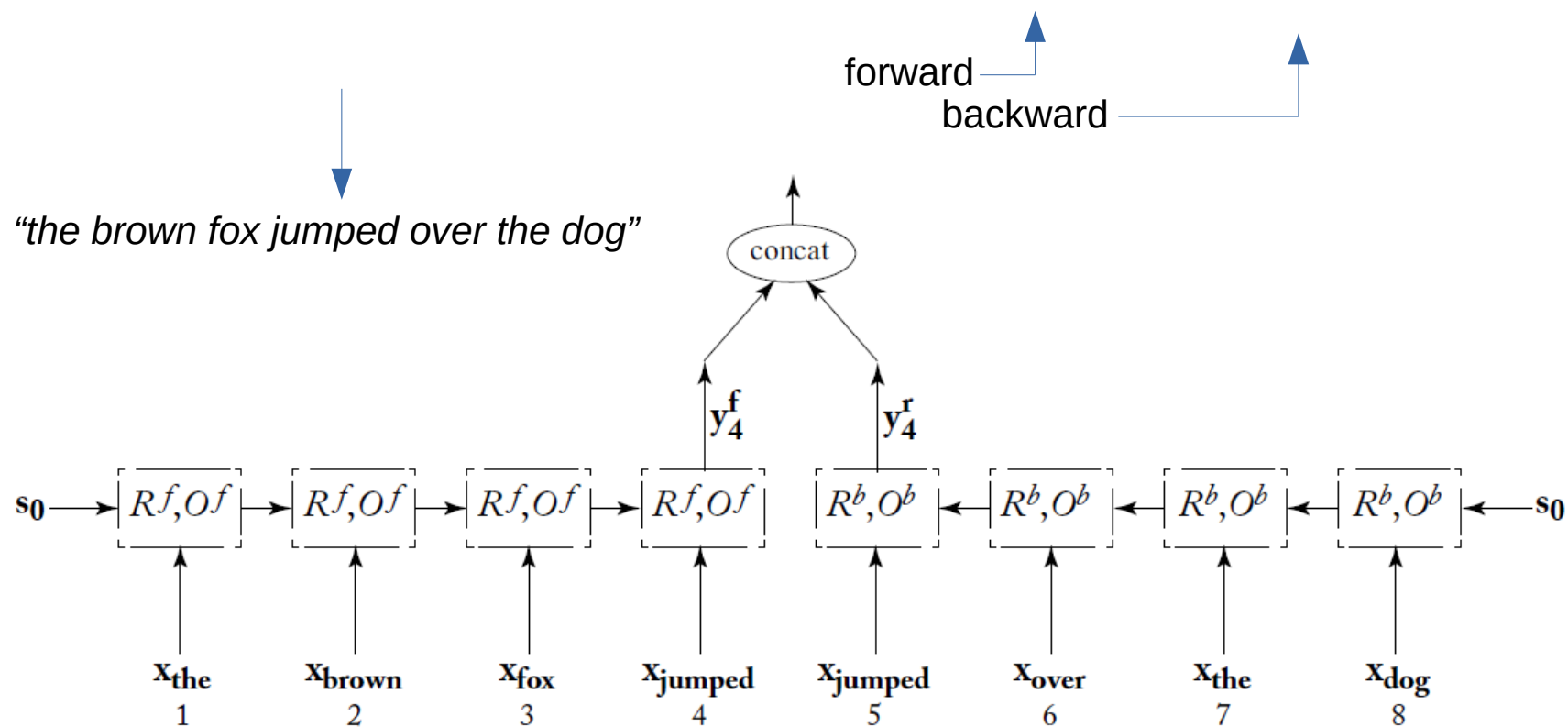
Bidirectional RNN (biRNN):



## Sequence labeling

Bidirectional RNN (biRNN):

$$\text{biRNN}(x_{1:n}, i) = y_i = [\text{RNN}^f(x_{1:i}); \text{RNN}^b(x_{n:i})].$$



## Sequence labeling

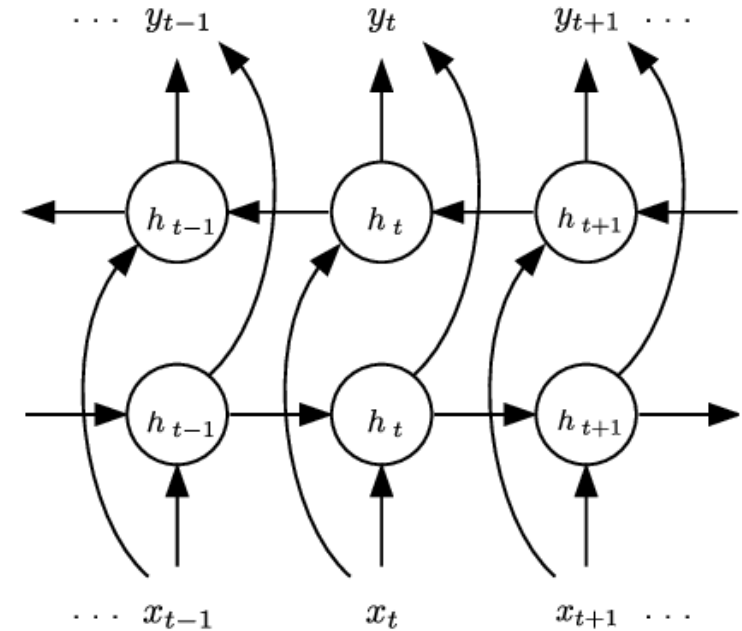
Bidirectional RNN (biRNN):

Dos capas:

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

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

$$y_t = g(V^{(f)} \cdot h_t^{(f)} + V^{(b)} \cdot h_t^{(b)})$$





## Sequence labeling

Bidirectional RNN para POS tagging con modelo preentrenado de subpalabras

Char embeddings:  $x_i = \phi(s, i) = [E_{[w_i]}; \text{RNN}^f(c_{1:\ell}); \text{RNN}^b(c_{\ell:1})]$ .

FastText  char n-grams

POS-tagging:  $p(t_i = j | w_1, \dots, w_n) = \text{softmax}(\text{MLP}(\text{biRNN}(x_{1:n}, i)))_{[j]}$

## Sequence labeling

Bidirectional RNN para POS tagging con modelo preentrenado de subpalabras

