



deeplearning.ai

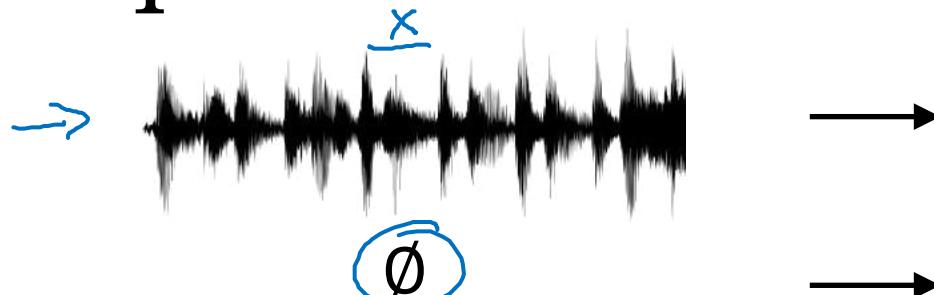
# Recurrent Neural Networks

---

Why sequence  
models?

# Examples of sequence data

Speech recognition



$y$   
“The quick brown fox jumped  
over the lazy dog.”

Music generation



Sentiment classification

“There is nothing to like  
in this movie.”



DNA sequence analysis → AGCCCCTGTGAGGAAC TAG



AG $\textcolor{red}{CCCCTGTGAGGAAC}$  TAG

Machine translation

Voulez-vous chanter avec  
moi?



Do you want to sing with  
me?

Video activity recognition



Running

Name entity recognition

Yesterday, Harry Potter  
met Hermione Granger.



Yesterday, **Harry Potter**  
met **Hermione Granger**.

Andrew Ng



deeplearning.ai

# Recurrent Neural Networks

---

## Notation

# Motivating example

NLP ~ Notation ~

Input sentence:

x:

(Harry Potter) and (Hermione Granger) invented a new spell.

Sequence

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad \dots \quad x^{<t>} \quad \dots \quad x^{<q>}$

$$T_x = q$$

len. de la seq

Output sequence:

y:

1 1 0 1 1 0 0 0 0  
 $y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad \dots \quad y^{<q>}$

El ej. es de "Name Entity Recognition." Busca id. donde hay Nombres en la secuencia.

$$\begin{matrix} x^{(i)<t>} \\ y^{(i)<t>} \end{matrix}$$

Notación p/ el elemento i de u seq (Training example)

$$\begin{matrix} T_x^{(i)} = q \\ T_y^{(i)} \end{matrix}$$

$$T_y = q$$

→  $T_y$  Puede ser  $\neq$  de  $T_x$  y Ademas c/ Training example Puede ser  $\neq$   $\rightarrow T_x^{(i)} + T_x^{(j)}$

Andrew Ng

# Representing words

$$x^{<\leftrightarrow>} \quad x \rightarrow y \quad (x, y)$$

x: Harry Potter and Hermione Granger invented a new spell.

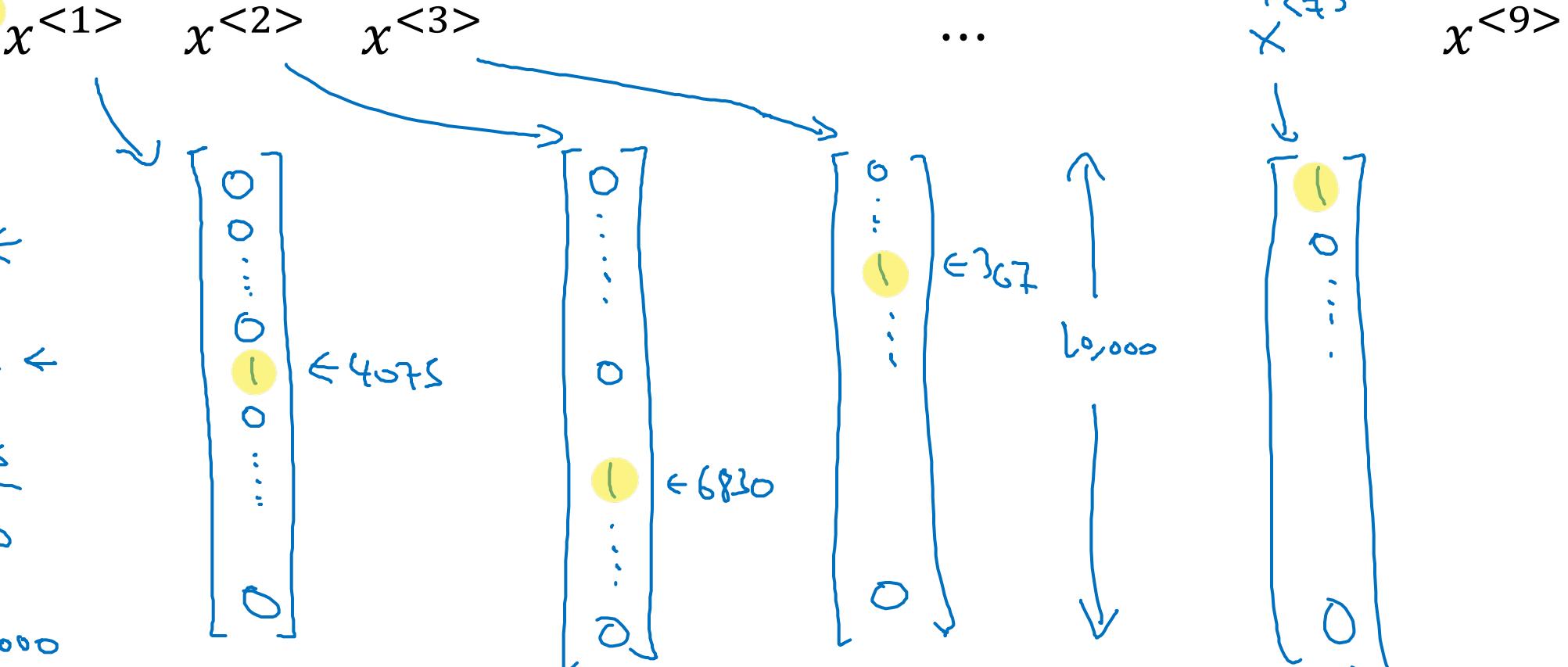
Ok, Pero como representamos words??

Así

Vocabulary  $\otimes$

|        |        |
|--------|--------|
| a      | 1      |
| aaron  | 2      |
| :      | :      |
| and    | 367    |
| :      | :      |
| harry  | 4075   |
| :      | :      |
| potter | 6830   |
| :      | :      |
| zulu   | 10,000 |

UNK 10,000  $\textcircled{I}$



One-hot

$\otimes$  lista de words q venios a usar, e del idioma.

Andrew Ng

II los posos van de ~30k - 50k y las Big tech llegan a USAR ~1M.

# Representing words

x: Harry Potter and Hermione Granger invented a new spell.  
 $x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad \dots \quad x^{<9>}$

\* Si HAY una word q no est. en vocab.  
Agrega una nueva "word" en vocab.  
word.

And = 367  
Invented = 4700  
A = 1  
New = 5976  
Spell = 8376  
Harry = 4075  
Potter = 6830  
Hermione = 4200  
Gran... = 4000



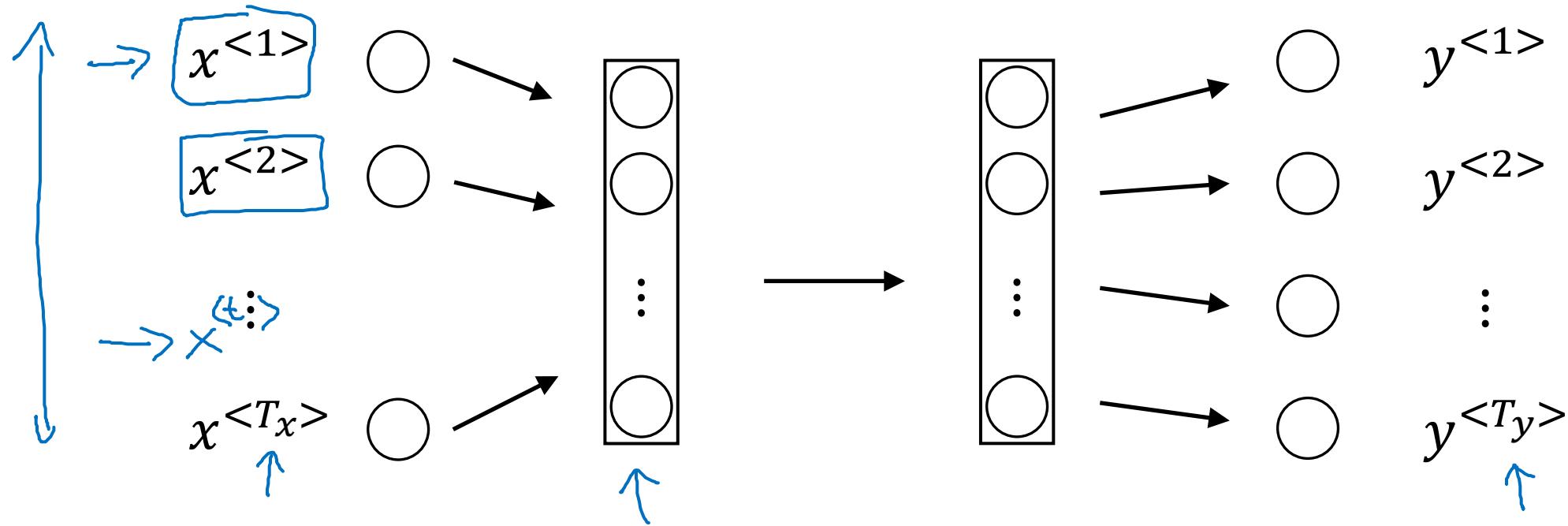
deeplearning.ai

# Recurrent Neural Networks

---

## Recurrent Neural Network Model

# Why not a standard network? feed forward NLP

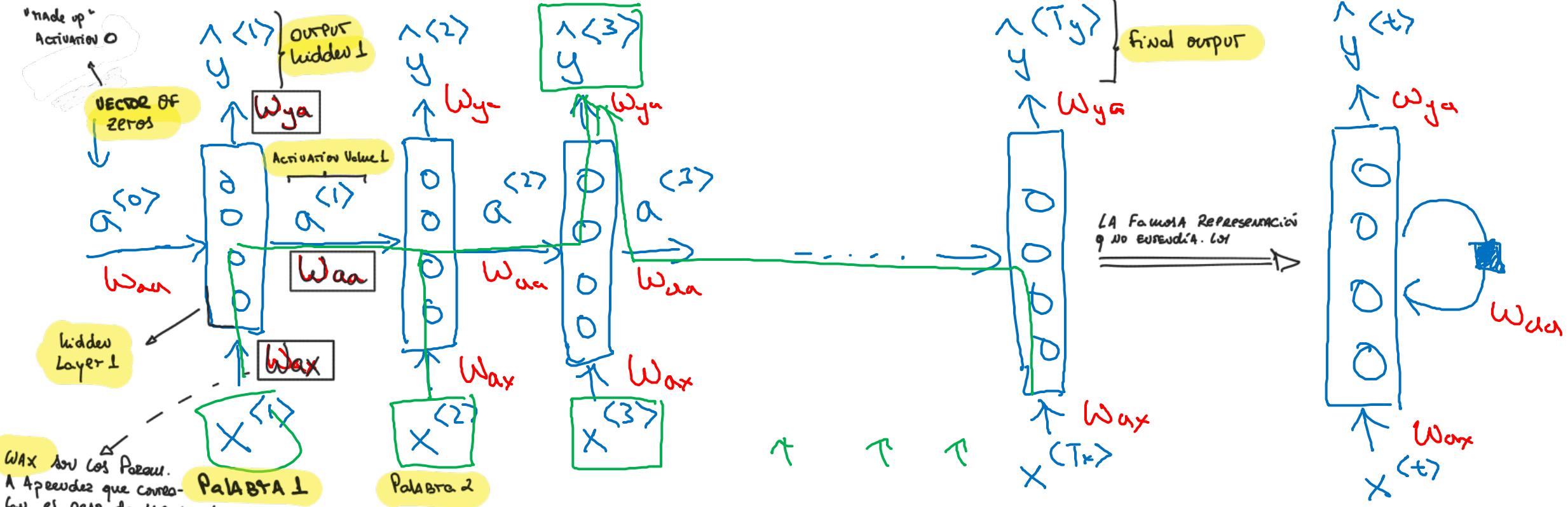


## Problems:

- - Inputs, outputs can be different lengths in different examples.
- - Doesn't share features learned across different positions of text.

→ Analogia a la Convolución. No explora la estructura secuencial y NO generaliza a tareas de EMA.

# Recurrent Neural Networks



W<sub>ax</sub> son los pesos.  
Aprenden que corresponde el peso de las palabras en el paso. Son los mismos, por lo tanto.

Una vez, las Activaciones serán solo por pesos \$W\_{aa}\$, más tarde para todas las capas.

W<sub>ya</sub> es el análogo a esas pesos, pero para las salidas de cada capa \$y^{(i)}

Parámetros a aprender.

He said, "Teddy Roosevelt was a great President."

He said, "Teddy bears are on sale!"

RNN a secas. Unidireccional.  
Los Bidireccionales usan info del futuro y los valores a ver + Adelante.

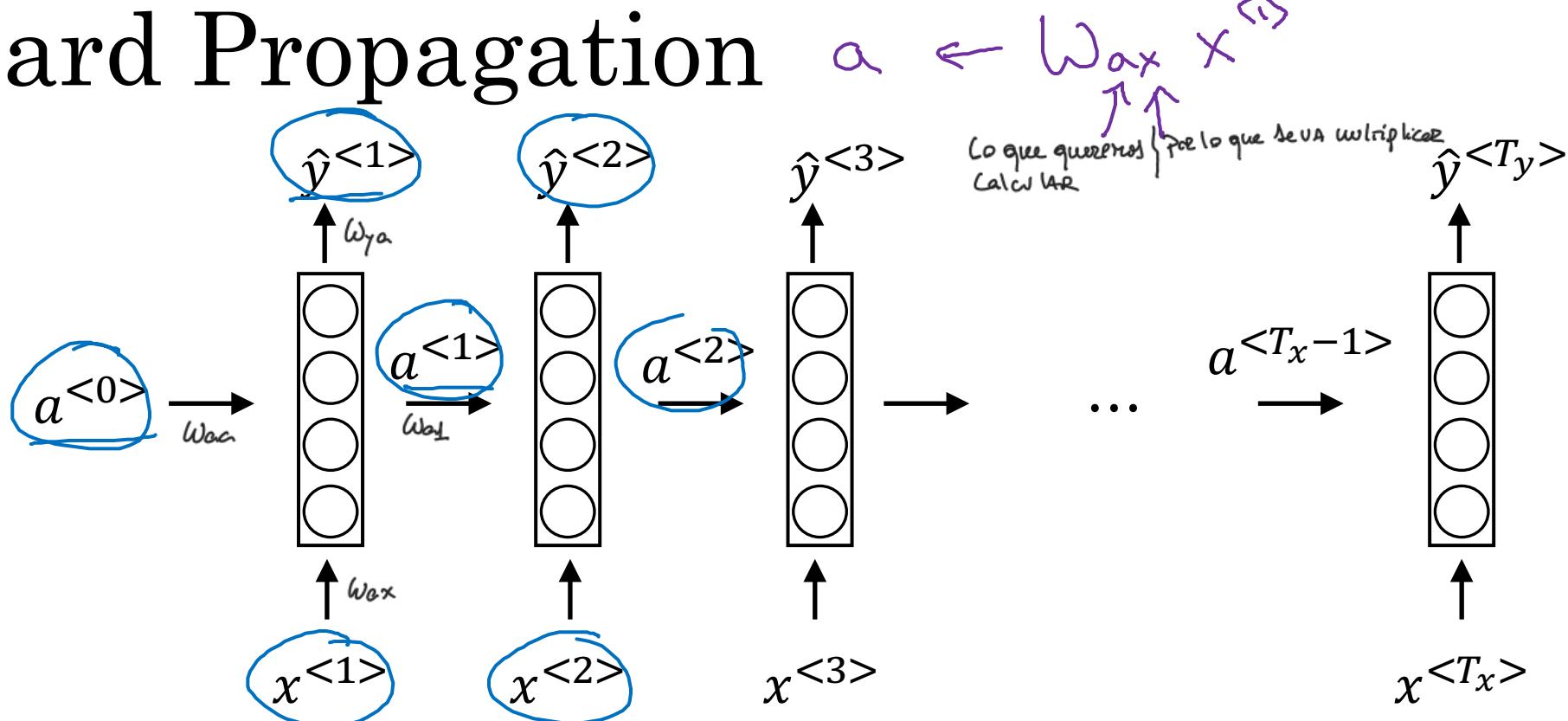
¿Qué es la Activación Value? \$y^{(i)}\$ Activado con una ReLU, Sigmoid, etc?

⇒ Usando info de las etapas anteriores. Lo único es q es una Arg. NO usa info de las etapas siguientes. Problema: EN EL EJ. PI saber si Teddy es un Noun, sería def. Saber de "president". Solución: Bidireccional RNN.

Andrew Ng

# Forward Propagation

Notación de los Weights: (I)



$$a^{(0)} = \vec{0}.$$

$$a^{(t)} = g_1(W_{aa} a^{(t-1)} + W_{ax} x^{(t)} + b_a) \leftarrow \underbrace{\tanh \text{ / ReLU}}$$

$$\hat{y}^{(t)} = g_2(W_{ya} a^{(t)} + b_y) \leftarrow \begin{array}{l} \text{Sigmoid} \\ \xrightarrow{g_2} \text{puedes ver +} \end{array}$$

VER (I)

|  |                             |
|--|-----------------------------|
| $a^{(t)} = g(W_{aa} a^{(t-1)} + W_{ax} x^{(t)} + b_a)$ | GENERALIZANDO               |
| $\hat{y}^{(t)} = g(W_{ya} a^{(t)} + b_y)$              | → Los Valores a Simplificar |

RNN  
FORWARD  
PROP.

Andrew Ng

# Simplified RNN notation

Matricialmente y Simplificando  
con Notación.

Partiendo de (Notación Matricial)



$$a^{(t)} = g(W_{aa}a^{(t-1)} + W_{ax}x^{(t)} + b_a)$$

↑ 100      ↑ 10,000

$\in \mathbb{R}^{(100, 100)}$      $\in \mathbb{R}^{(100, 10,000)}$

$$\hat{y}^{(t)} = g(W_{ya}a^{(t)} + b_y)$$

Análogamente, Simplificando:

$$y^{(t)} = g(W_y a^{(t)} + b_y)$$

Finalmente:

FORWARD Prop. RNN Simplificado

$$a^{(t)} = g(W_a [a^{(t-1)}, x^{(t)}] + b_a)$$

$$\hat{y}^{(t)} = g(W_y a^{(t)} + b_y)$$

Reescribimos

$$a^{(t)} = g(W_a [a^{(t-1)}, x^{(t)}] + b_a)$$

donde:

①  $[W_{aa}; W_{ax}]$  =  $W_a$  → eje x.

$\in \mathbb{R}^{(100, 10100)}$

②  $[a^{(t-1)}, x^{(t)}]$  =   
 A su vez, Esta Not. Es poner los v<sup>o</sup> uno  
 arriba del otro.

la matriz  $W_a$  es  $W_{ax}$   
y  $W_{aa}$  "pegadas" horizontalmente  
anteriormente a lo largo del eje x.

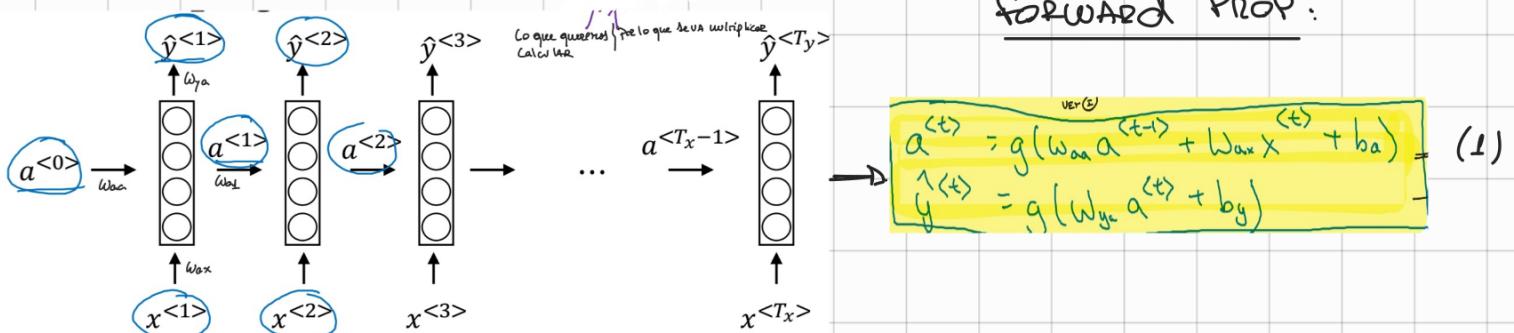
$[a^{(t-1)}]$        $x^{(t)}$

$\in \mathbb{R}^{(100, 10100)}$

③  $[W_{aa}; W_{ax}] \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} = W_{aa}a^{(t-1)} + W_{ax}x^{(t)}$

# Lo Hacemos más Ordenado:

Si Perdimos de Nuestra Notación Original p/ la BRNN, tenemos:



Si escribimos (1) en forma matricial, obtenemos:

$$\begin{cases} a^{(t)} = g(W_{aa}a^{(t-1)} + W_{ax}x^{(t)} + b_a) \\ \hat{y}^{(t)} = g(W_{ya}a^{(t)} + b_y) \end{cases}$$

y Con la Notación "STACKADA" definida en ① y ② de la Página anterior, obtenemos Por Último:

RNN

FORWARD Prop. ~~Stackada~~ Simplificado

$$a^{(t)} = g(W_a [a^{(t-1)}, x^{(t)}] + b_a)$$

$$\hat{y}^{(t)} = g(W_y a^{(t)} + b_y)$$



deeplearning.ai

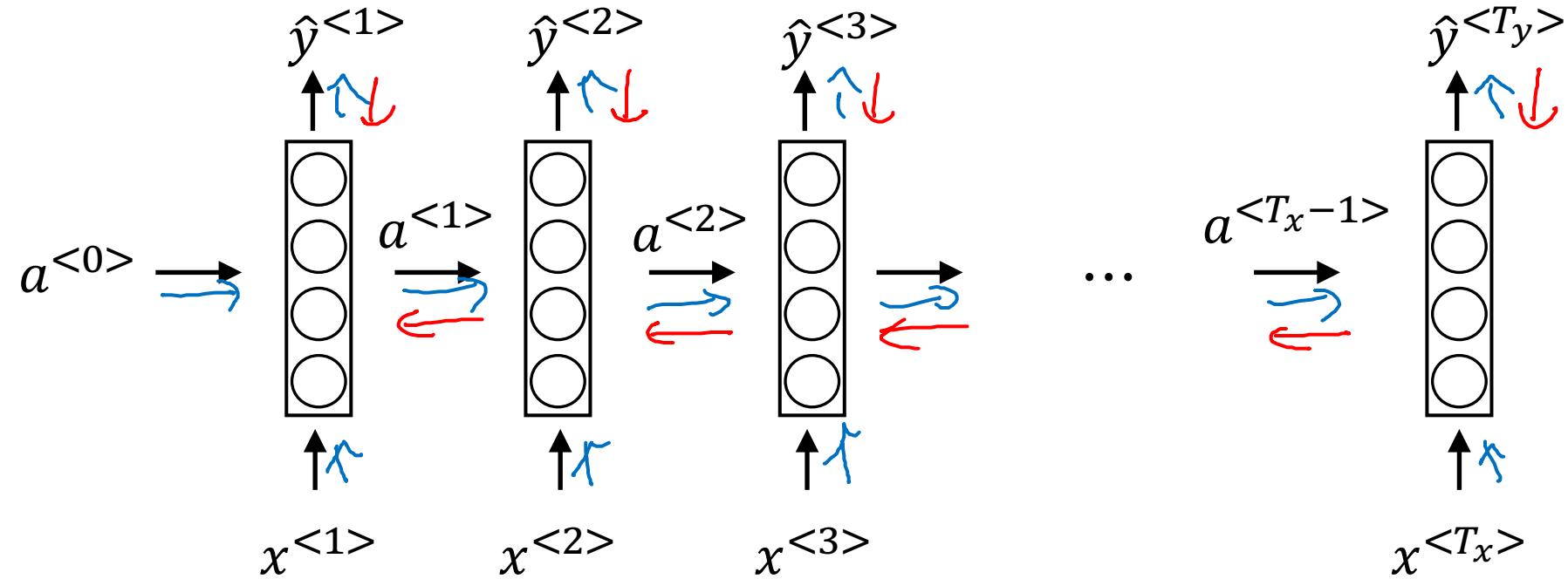
# Recurrent Neural Networks

---

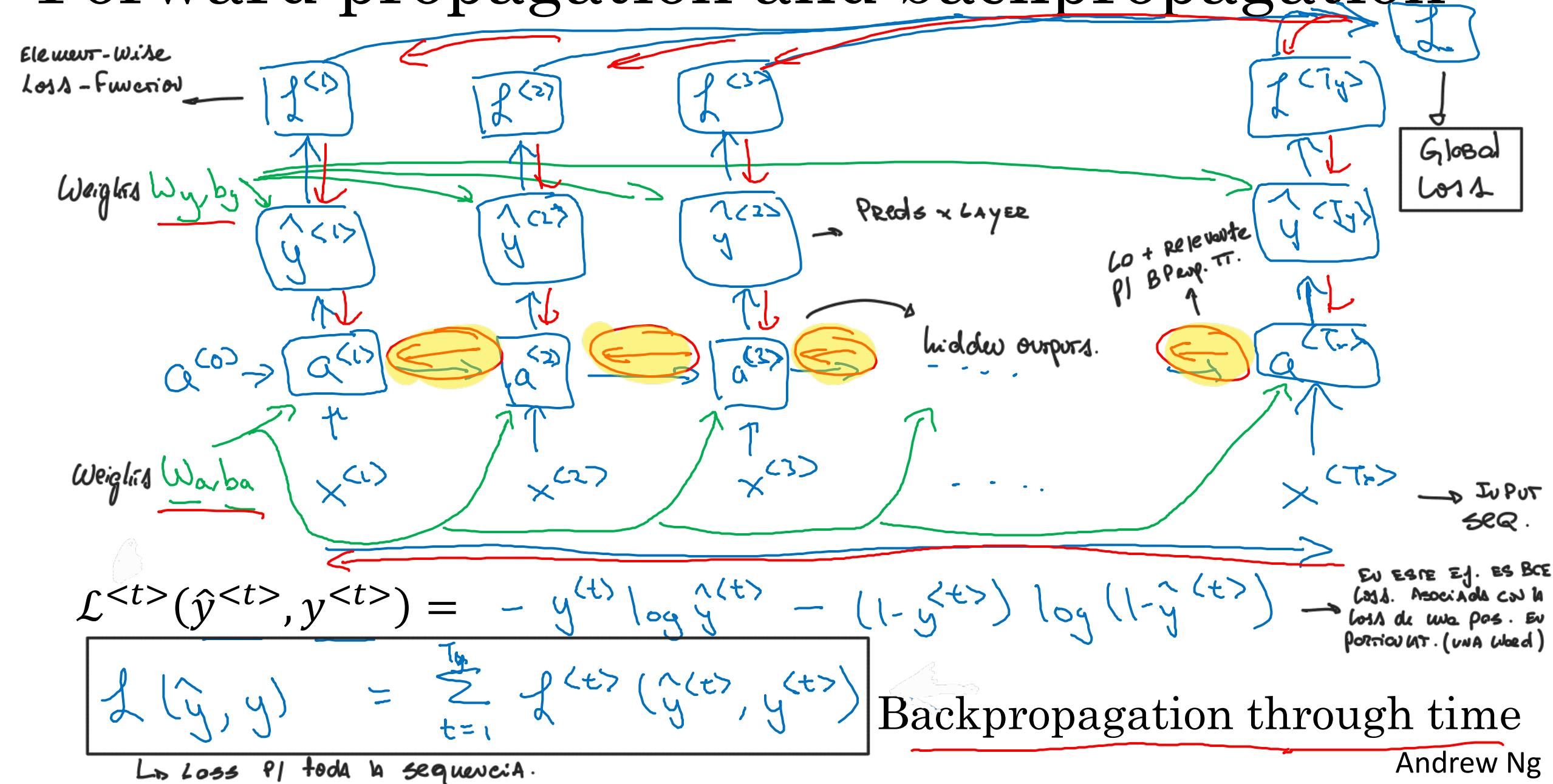
## Backpropagation through time

# Forward propagation and backpropagation

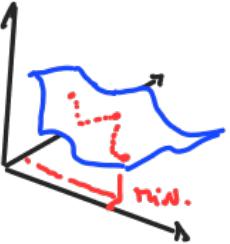
■ FORWARD  
✖ BACKWARD



# Forward propagation and backpropagation



↳ Calculando la Loss, y a veces calcular las derivadas de  $\frac{\partial h}{\partial w_i}$  el parámetros y con el  $\vec{\nabla}$  gradiente de derivadas parciales de Perceptrón, minimizar  $h(w_i)$  en el ESPACIO de PARÁMETROS con GRADIENT DESCENT !!.



# Recurrent Neural Networks

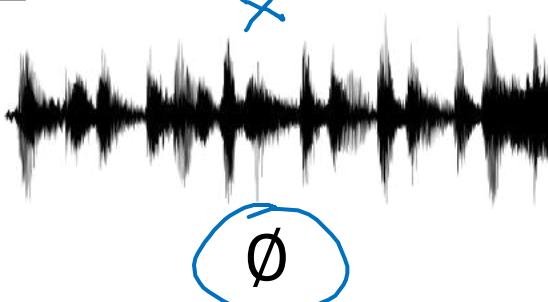
## Different types of RNNs

deeplearning.ai

Hasta ahora solo vimos casos que inputs = outputs ( $T_x = T_y$ ). pero esto puede no ser siempre así:

# Examples of sequence data

Speech recognition



$$T_x \quad T_y$$

y

“The quick brown fox jumped over the lazy dog.”

Music generation



Sentiment classification

“There is nothing to like  
in this movie.”

$$\text{εj de } \tau_x f \tau_y$$

→



DNA sequence analysis

AGCCCCTGTGAGGAAC TAG

$$\rightarrow$$

AGCCCCTGTGAGGAAC **TAG**

Machine translation

Voulez-vous chanter avec  
moi?

$$\text{εj de } \tau_x f \tau_y$$

→

Do you want to sing with  
me?

Video activity recognition



$$\rightarrow$$

Running

Name entity recognition

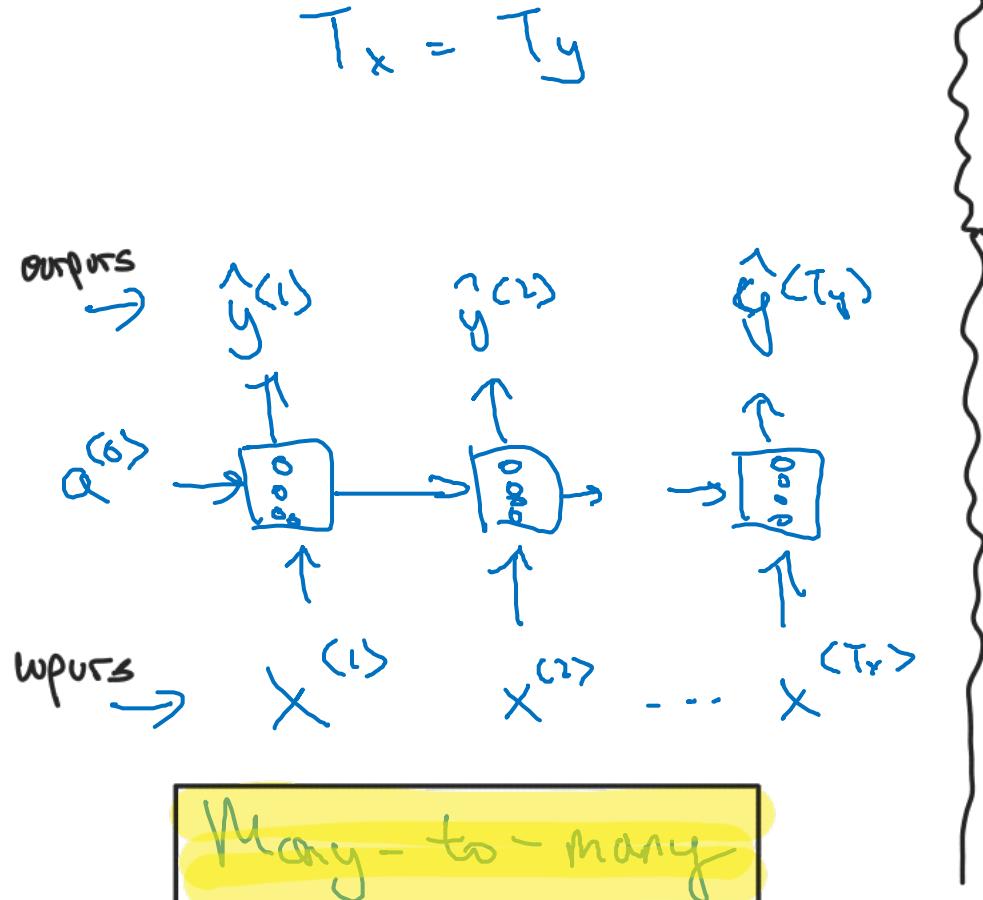
Yesterday, Harry Potter  
met Hermione Granger.

$$\rightarrow$$

Yesterday, **Harry Potter**  
met **Hermione Granger**.

Andrew Ng

# Examples of RNN architectures



$\epsilon_j =$  Sentiment classification

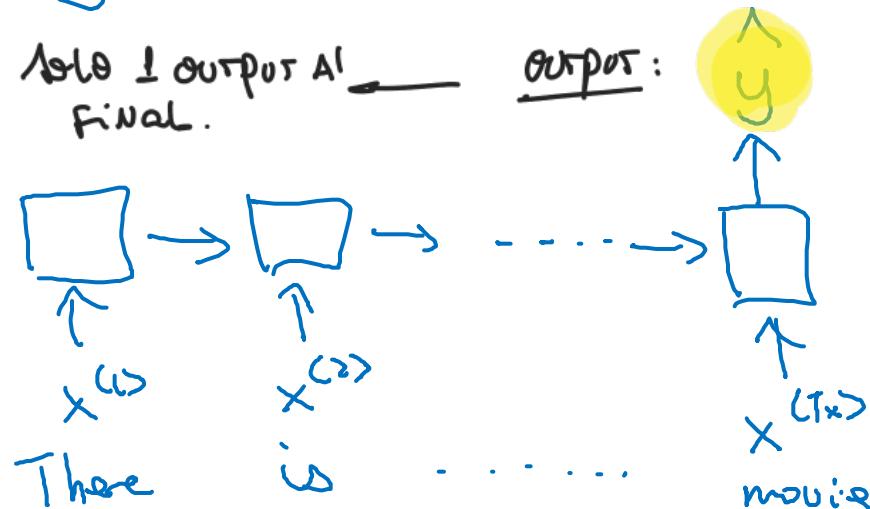
$x = \text{text}$

$y = 0/1 \quad 1 \dots 5$

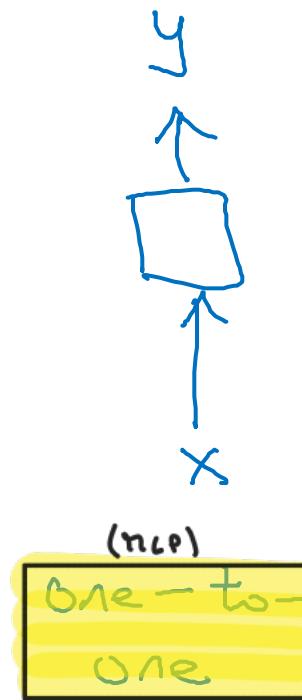
one 1 output at final.

Input: There is ...

Output:  $\hat{y}$

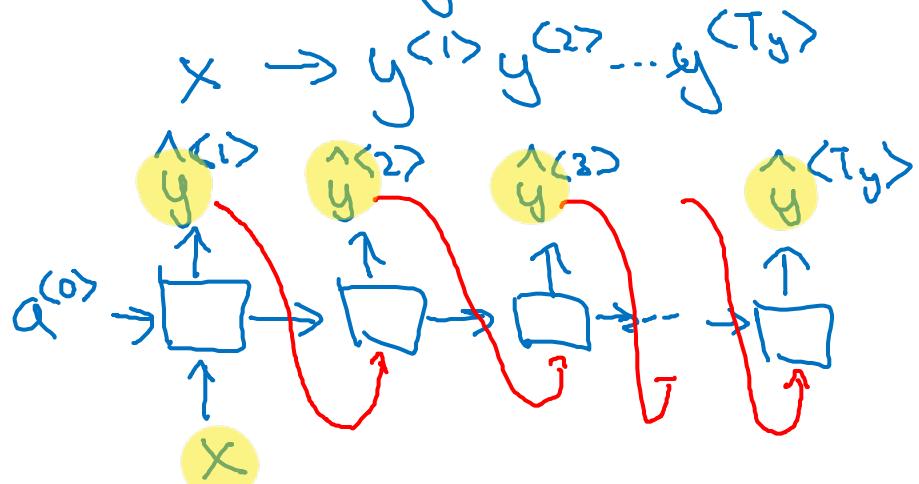


Many-to-one



# Examples of RNN architectures

Ej: Music generation

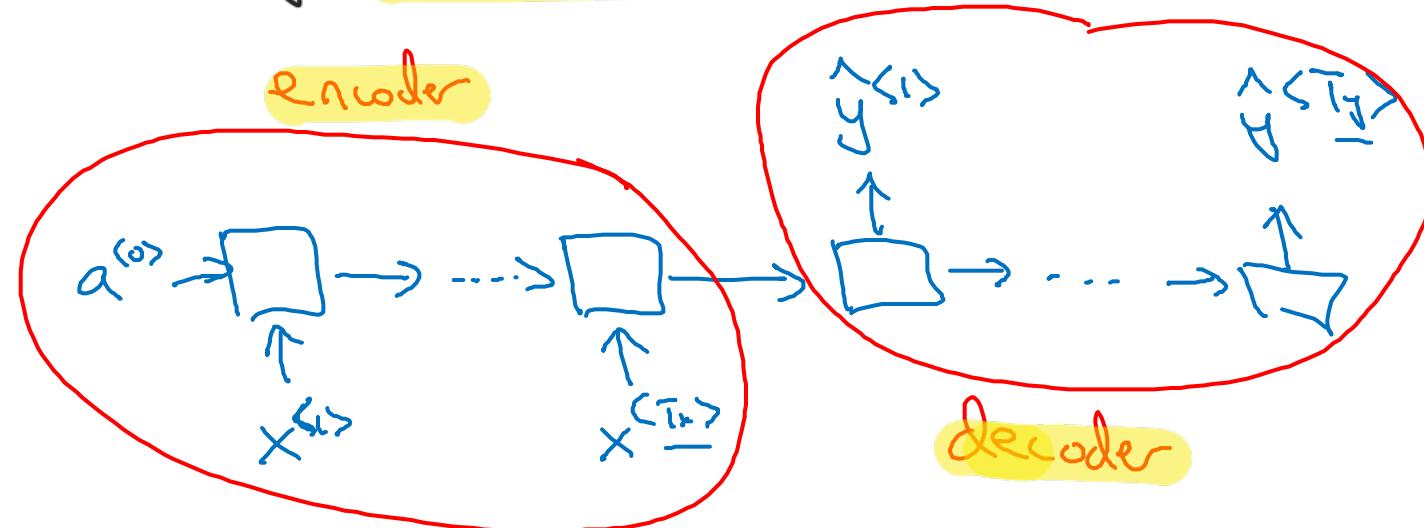


One-to-many

$$x = \phi$$

Ej: Machine translation

encoder



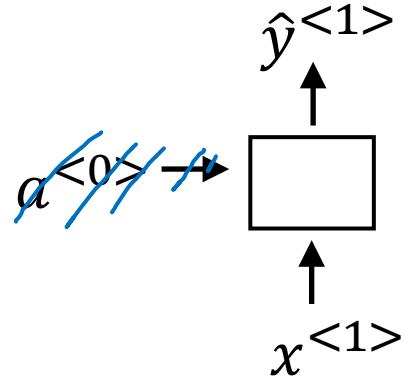
decoder

Many-to-many

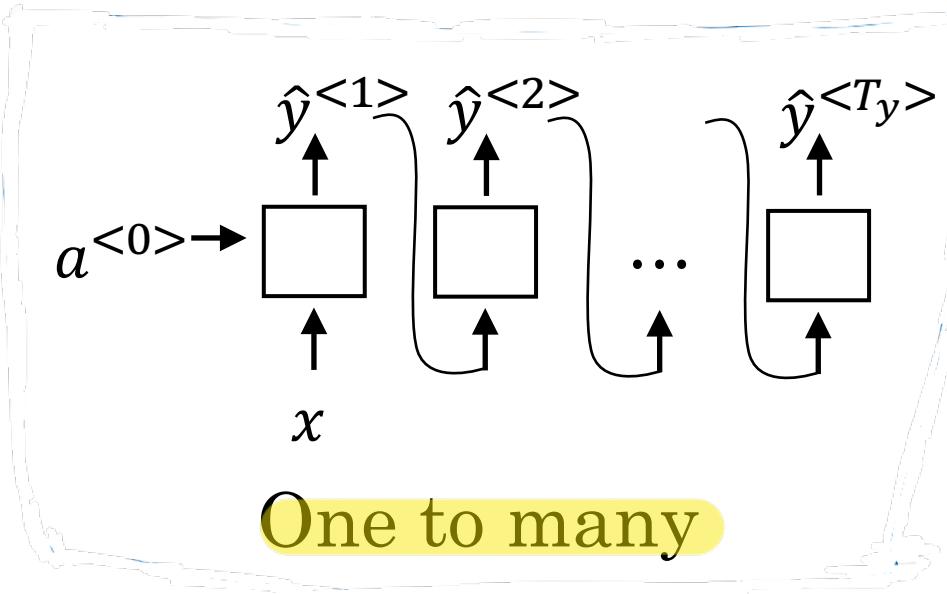
→ Zero cov input length  
≠ output length!

Encoder - Decoder

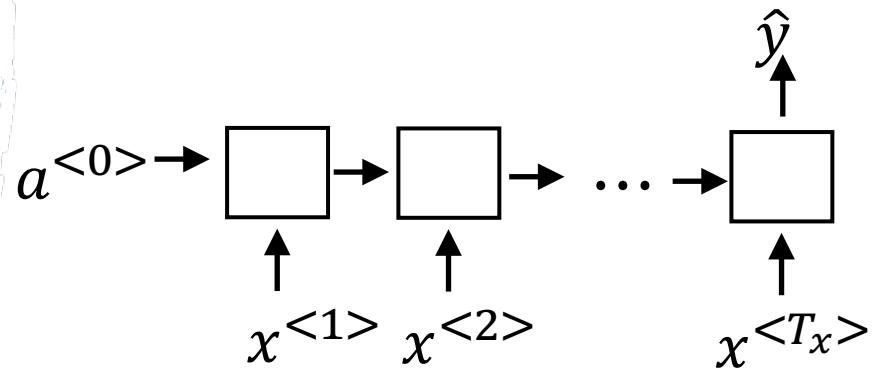
# Summary of RNN types



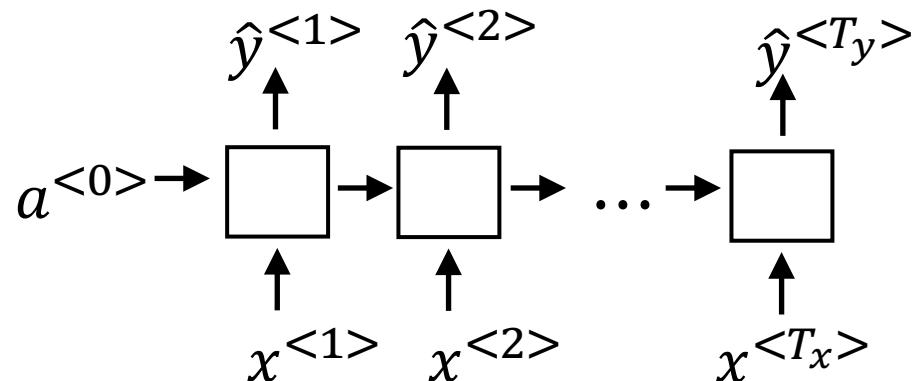
One to one



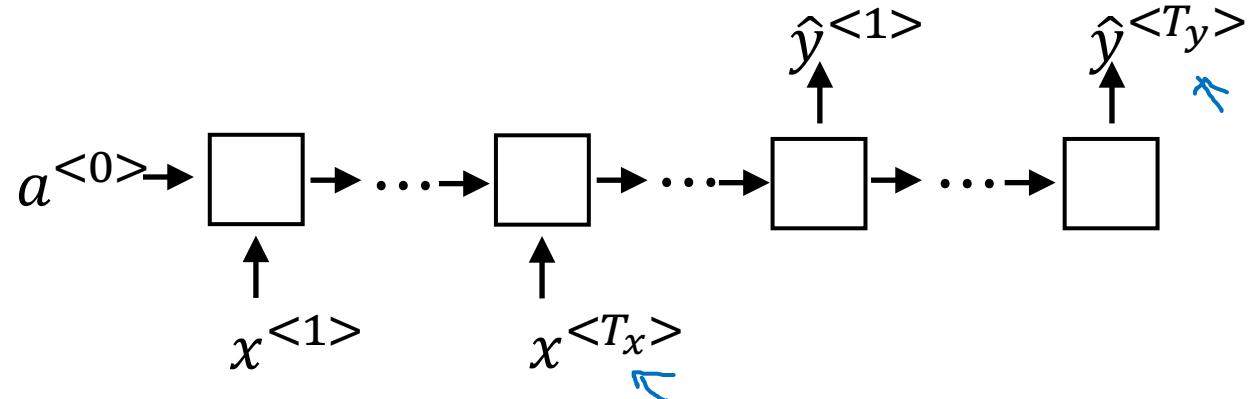
One to many



Many to one



Many to many  $T_x = T_y$



Many to many  $T_x \neq T_y$



deeplearning.ai

god. sequence gen.  
music, text, etc.  
god.

# Recurrent Neural Networks

---

Language model and  
sequence generation

# What is language modelling?

## Speech recognition

The apple and pair salad.

→ The apple and pear salad.

Si dije esto (Audio), cual es más probable que haya dicho pair?

$$\left\{ \begin{array}{l} P(\text{The apple and } \underline{\text{pair}} \text{ salad}) = 3.2 \times 10^{-3} \\ P(\text{The apple and } \underline{\text{pear}} \text{ salad}) = 5.7 \times 10^{-10} \end{array} \right. \Rightarrow P(\text{PEAR}) > P(\text{PAIR})$$

$$\Rightarrow \boxed{P(\text{Sentence}) = ?} \longrightarrow \boxed{P(y^{(1)}, y^{(2)}, \dots, y^{(T_y)})}$$

# Language modelling with an RNN

y... como buildemos  
un Language model?

Training set: large corpus of english text.

Tokenize

→ Doe Andrew solo habla de one-hot encoding. le dig.  
semeja europeo & jugar los embeddings.

Cats average 15 hours of sleep a day.  $\downarrow \langle \text{EOS} \rangle$

$y^{(1)}$        $y^{(2)}$   
 $x^{(t)} = y^{(t-1)}$

$y^{(3)}$       ...  
 $y^{(8)}$        $y^{(9)}$

avdeew dice q ueamos a terminar setteando el N° de inputs  
 $x^{(t)} = y^{(t-1)}$ , que despues nos explica pq.

The Egyptian ~~Mau~~ is a bread of cat.  $\langle \text{EOS} \rangle$

10,000

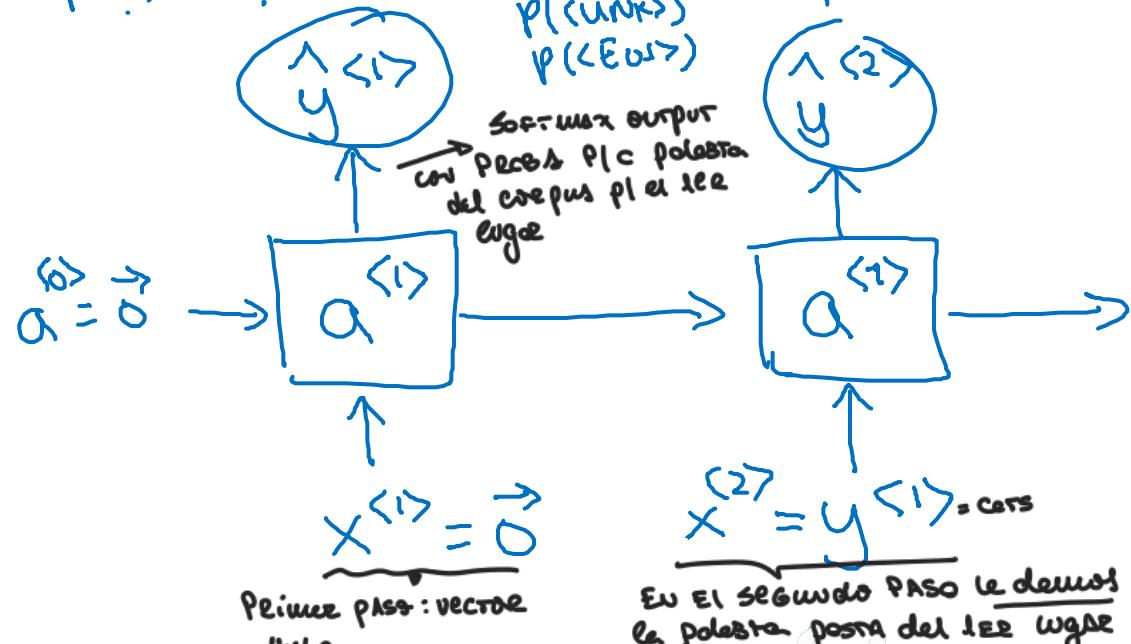
$\langle \text{UNK} \rangle$

tokens q no se ven en el corpus.

# RNN model

$p(a) p(aaron) \dots p(\text{cats}) \dots p(\text{tulu})$

$p(\text{cure})$   
 $p(\langle \text{EOS} \rangle)$

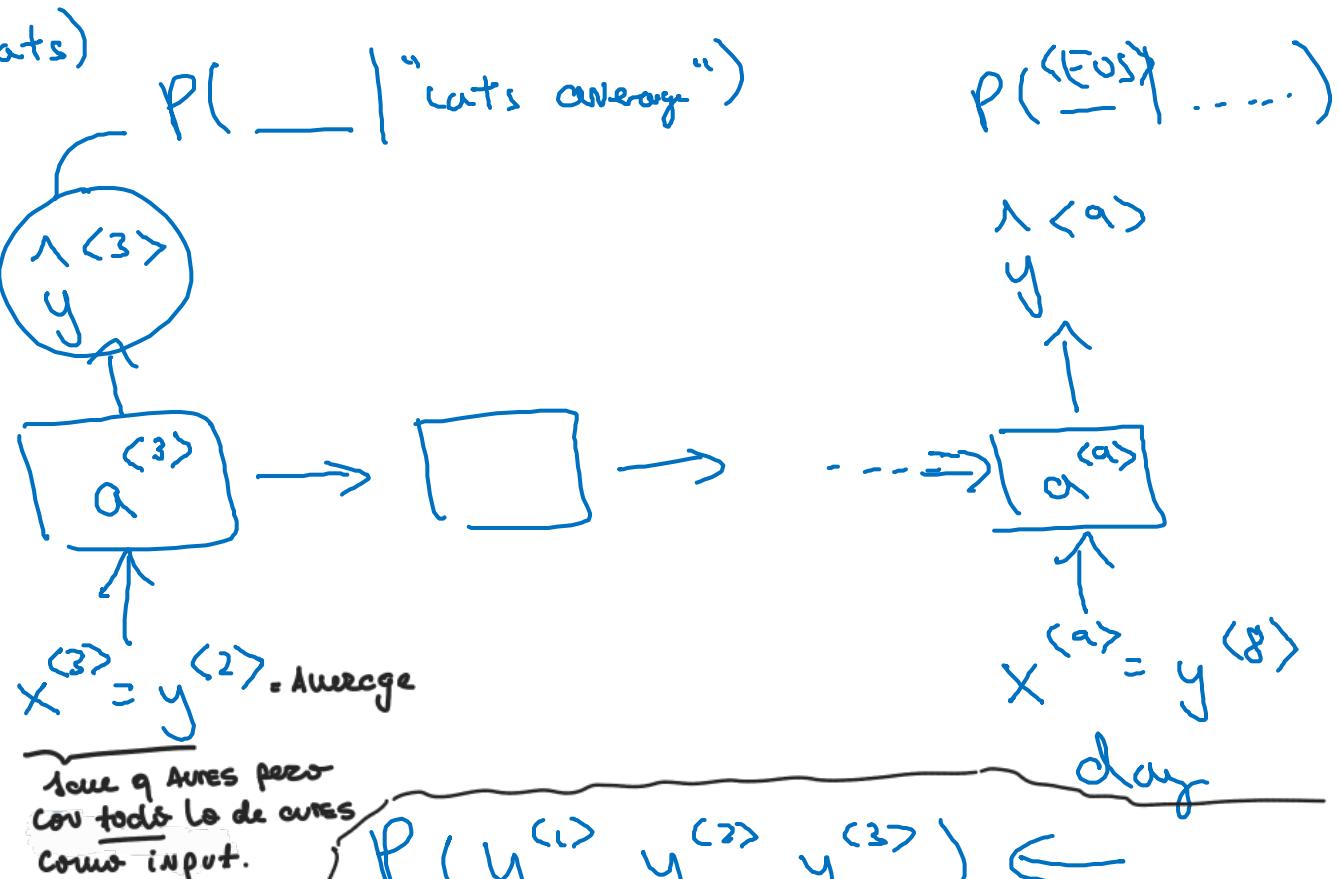


→ Cats average 15 hours of sleep a day.  $\langle \text{EOS} \rangle$

$$\mathcal{L}(\hat{y}^{<t>}, y^{<t>}) = - \sum_i y_i^{<t>} \log \hat{y}_i^{<t>}$$

$$\mathcal{L} = \sum_t \mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

Loss function. Con esto evalúa la prob. de cada output contra ground truth



- 1 = chance de  $y^{(1)}$
- 2 = chance de  $y^{(2)}$  dado  $y^{(1)}$
- 3 = chance de  $y^{(3)}$  dado  $[y^{(1)}, y^{(2)}]$

$\prod$  = Prob de todos la seguencia.

Andrew Ng



deeplearning.ai

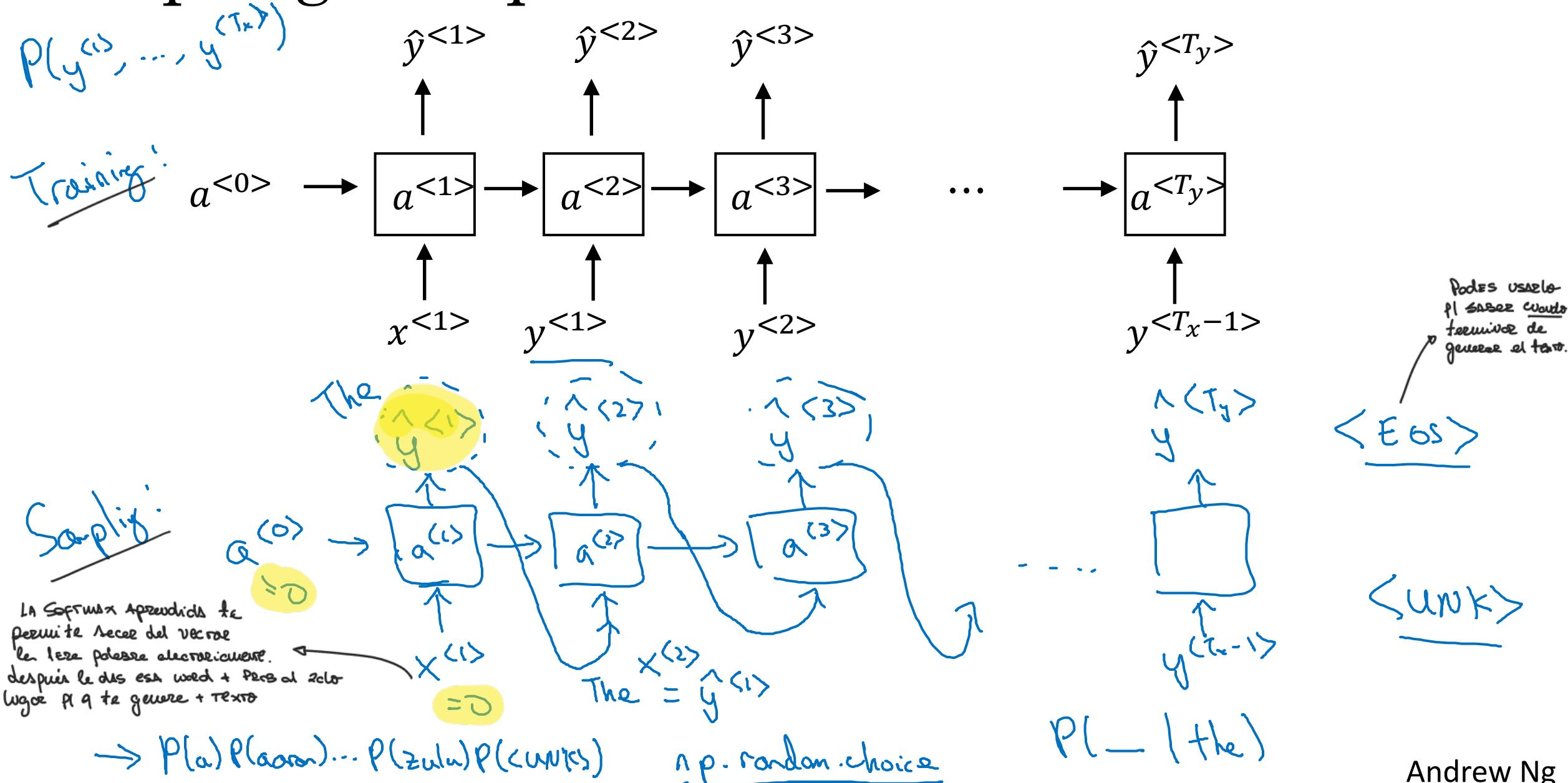
# Recurrent Neural Networks

---

Sampling novel  
sequences

Una vez que entrenamos  
el modelo, podemos pedirle  
que haga eso.

# Sampling a sequence from a trained RNN



# Character-level language model

- ④ more computationaly  
Expensive pero se  
dejar modelos polinomiales  
que no  $\Delta$ .

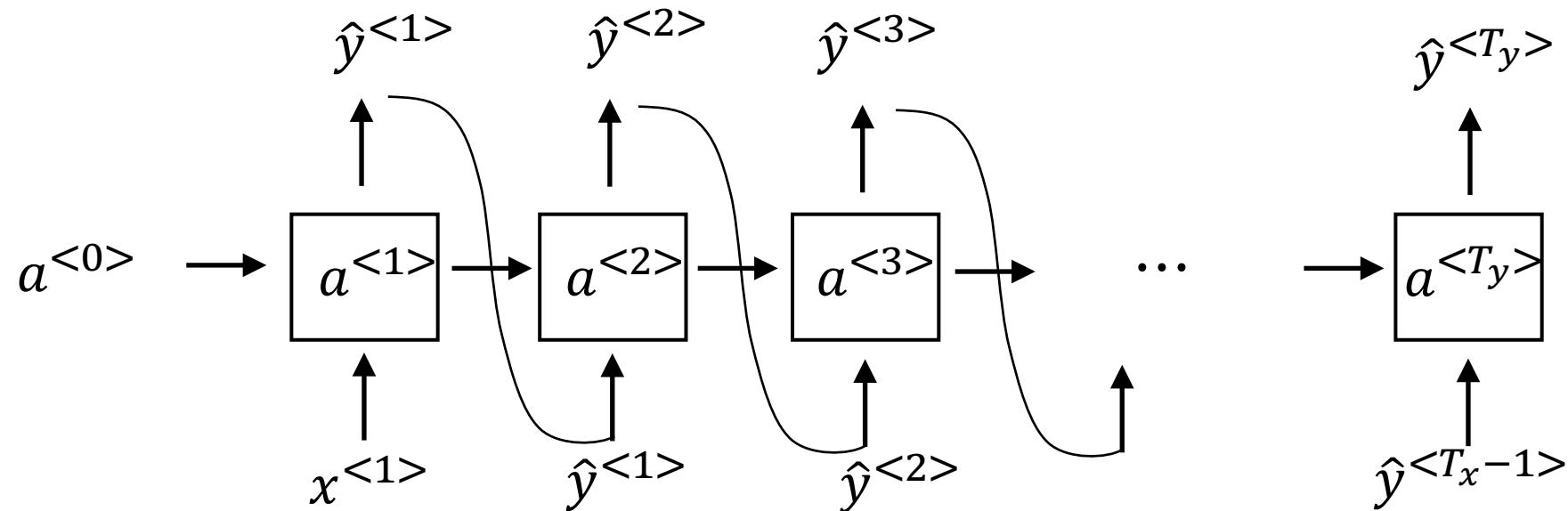
→ Vocabulary = [a, aaron, ..., zulu, <UNK>] → word level

$\rightarrow$  Vocabulary = [a, b, c, ..., z, \cup, o, , , ;, O, ..., Q, A, ..., Z]  $\rightarrow$  CHARACTER LEVEL

$$y^{(0)} = y^{(1)} = y^{(2)} = y^{(3)}$$

Cat average  
↑ ↑ ↑ ↑ . . .

Man



# Sequence generation

## News

President enrique peña nieto, announced  
sench's sulk former coming football langston  
paring.

“I was not at all surprised,” said hich langston.

“Concussion epidemic”, to be examined. ←

The gray football the told some and this has on  
the uefa icon, should money as.

## Shakespeare

The mortal moon hath her eclipse in love.  
And subject of this thou art another this fold.

When lesser be my love to me see sabl’s.  
For whose are ruse of mine eyes heaves.



deeplearning.ai

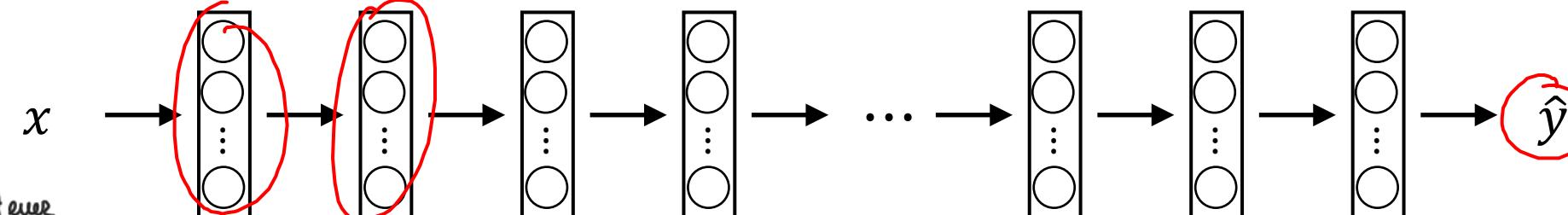
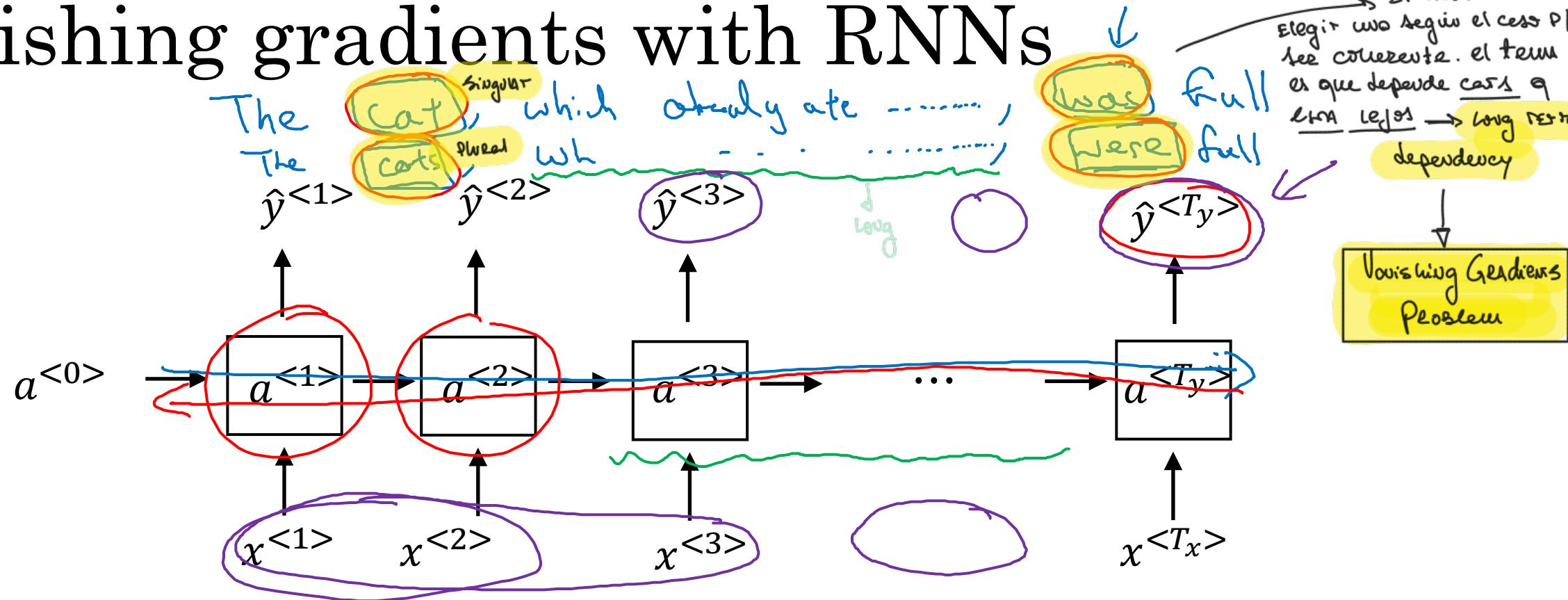
# Recurrent Neural Networks

---

[GRU's and LSTM's]

Vanishing gradients  
with RNNs

# Vanishing gradients with RNNs



También podemos tener este problema, aunque es menos frecuente, es más fácil de ver por el overflow y q se vea a la gente los pesos.

Exploding gradients.

Nan

Gradient clipping

→ max threshold q pone pl. Evitar Exploding Gradients. Andrew Ng



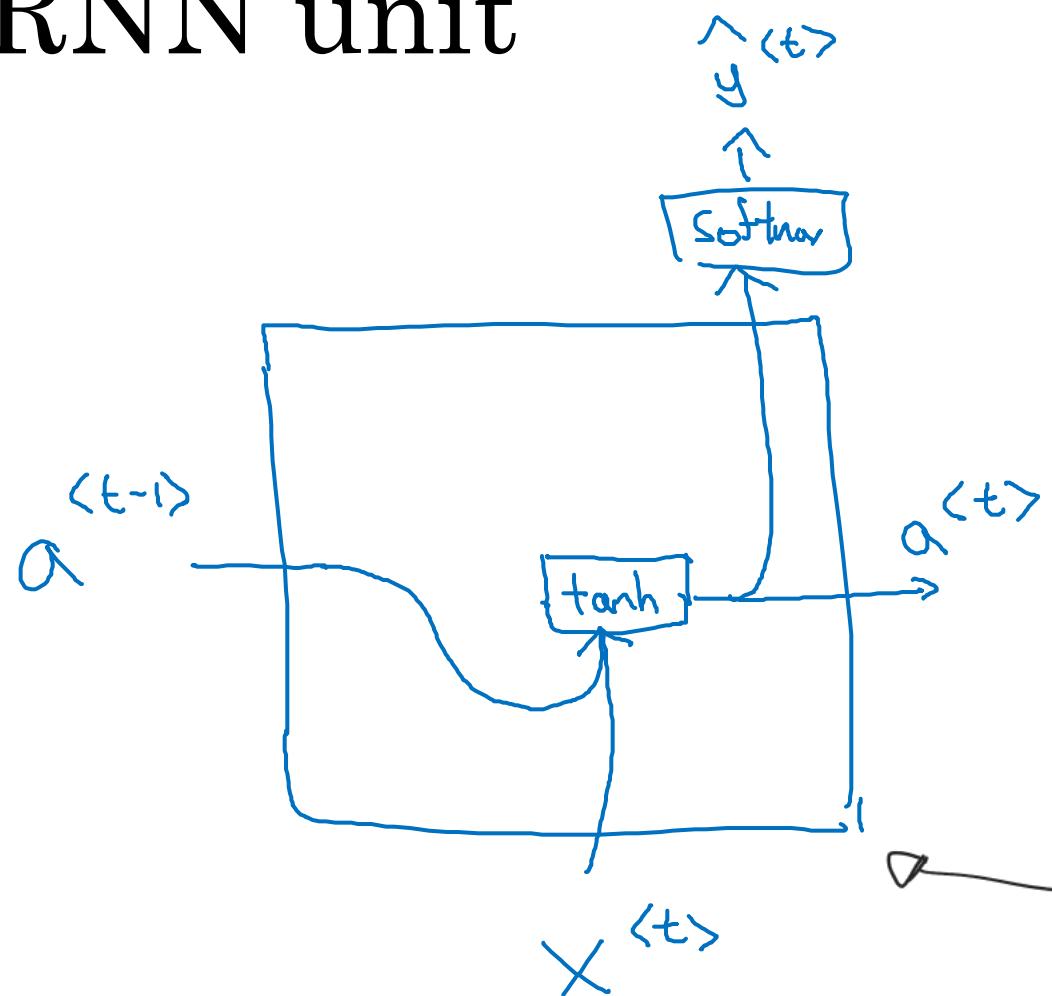
deeplearning.ai

# Recurrent Neural Networks

---

Gated Recurrent  
Unit (GRU)

# RNN unit

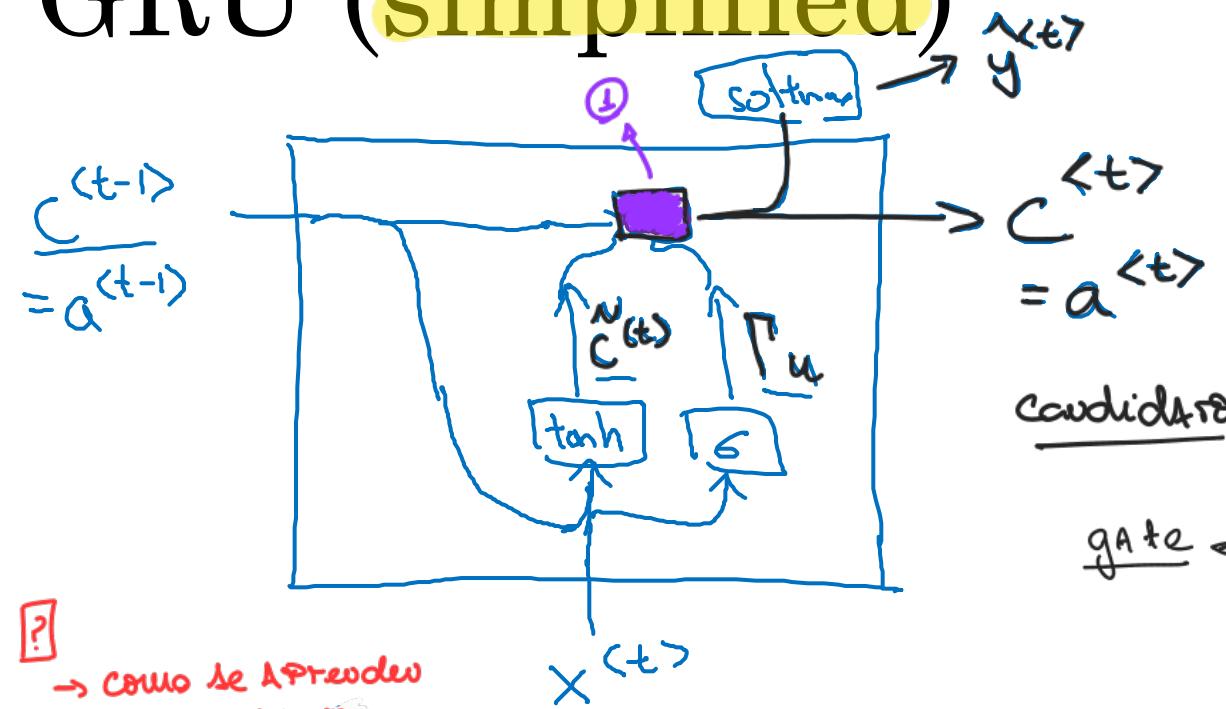


→ Formula de la Activación de una RNN, segun lo visto:

$$a^{(t)} = g(W_a[a^{(t-1)}, x^{(t)}] + b_a)$$

graficamente, la hidden layer sería algo así.

# GRU (simplified)



→ como se aprenden  $\Gamma_u$  y  $C$ ? O solo  $\Gamma$  capa.

$$\Gamma_u = 1$$

$$C^t = 1$$

$$\Gamma_u = 0 \quad \Gamma_u = 0 \quad \Gamma_u = 0 \quad \dots$$

→ The cat, which already ate ..., was full.

Candidato que en el timestamp va a INTRODUCIR ejemplo de memoria "C<sup>t</sup>"

$C = \text{Memory cell}$

$C^t = a^t$

$\text{Value previous memory} = a^{t-1}$

$\text{Value Accrued seq.} = a^t$

$N^t = \tanh(W_c [C^{t-1}, x^t] + b_c)$

$\Gamma_u = \sigma(W_u [C^{t-1}, x^t] + b_u)$

①  $C^t = \frac{\Gamma_u}{\Gamma_u + \hat{C}} * \hat{C} + (1 - \frac{\Gamma_u}{\Gamma_u + \hat{C}}) * C^{t-1}$

→ Recalculamos  $C^{t+1} = 1$   
Memoria en base a  $\Gamma$  y  $\hat{C}$ .  
el gate filtra al candidato.

element-wise

$$\Gamma_u = 0.000001$$

Mantiene el valor  
de  $C$  por ser <<<  
previene los Vanishing

Andrew Ng

[Cho et al., 2014. On the properties of neural machine translation: Encoder-decoder approaches]

[Chung et al., 2014. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling]

# Full GRU

$\tilde{h}$

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

$u$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$r$

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$h$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

↳ Notación de Andrew

↳ Notación más técnica.

Aprende la Relevancia de  $c^{<t-1>}$  p/ comprender el alg. candidato p/  $c^{<t>}$

→ Porque Agregamos  $\Gamma_r$  y no usamos la versión simple de antes? el research indica que tenerla es mejor

→ En cuanto a como se aprenden los pesos, creo que la Rta esta en los  $W_r$  y  $W_u$  q serian los weights (matrices) de estas unidades.

The cat, which ate already, was full.



deeplearning.ai

Andrew dice que son + power  
que los GRU.

# Recurrent Neural Networks

---

LSTM (long short term memory) unit

# GRU and LSTM

## GRU

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$a^{<t>} = c^{<t>}$$



## LSTM

$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

$$\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \text{Tanh}(c^{<t>})$$

# LSTM units

## GRU

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

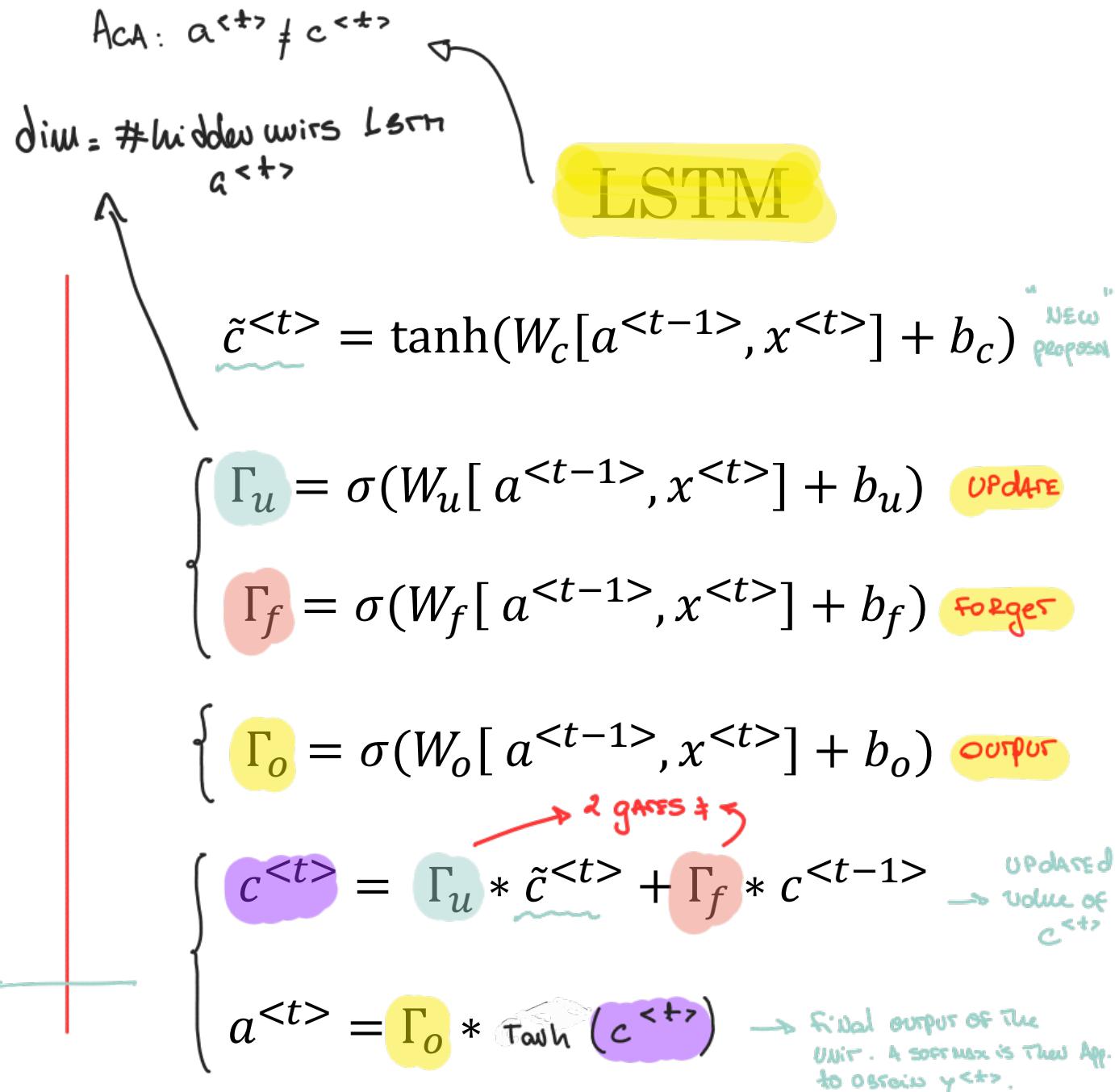
$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$a^{<t>} = \tilde{c}^{<t>} \quad \text{UNA SOLA GATE ACTIVA COMO UPDATE + FORGET}$$

Is, é a  $c^{<t>}$  the "hidden state" or dim in LSTM's?



# LSTM in pictures

$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

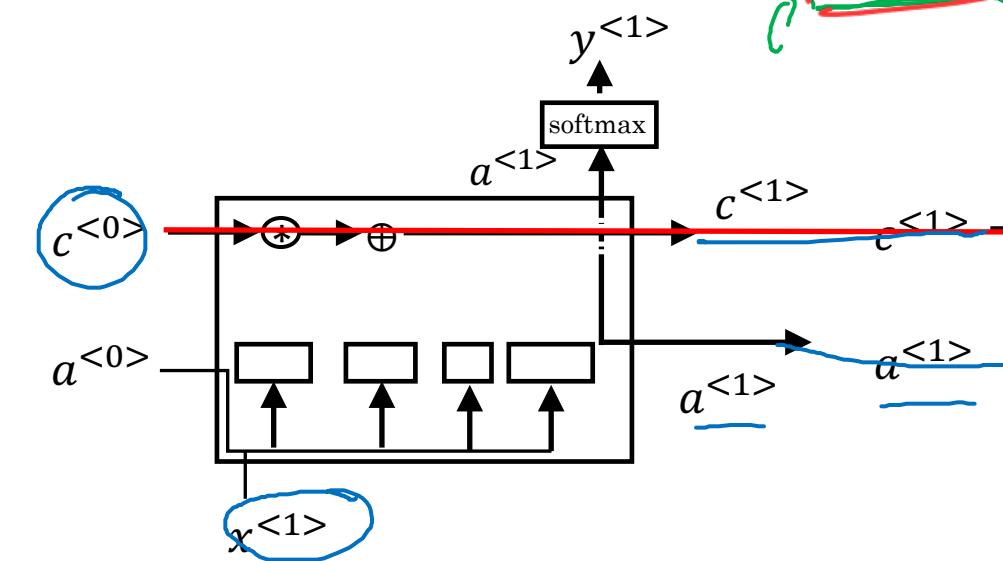
$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

$$\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$

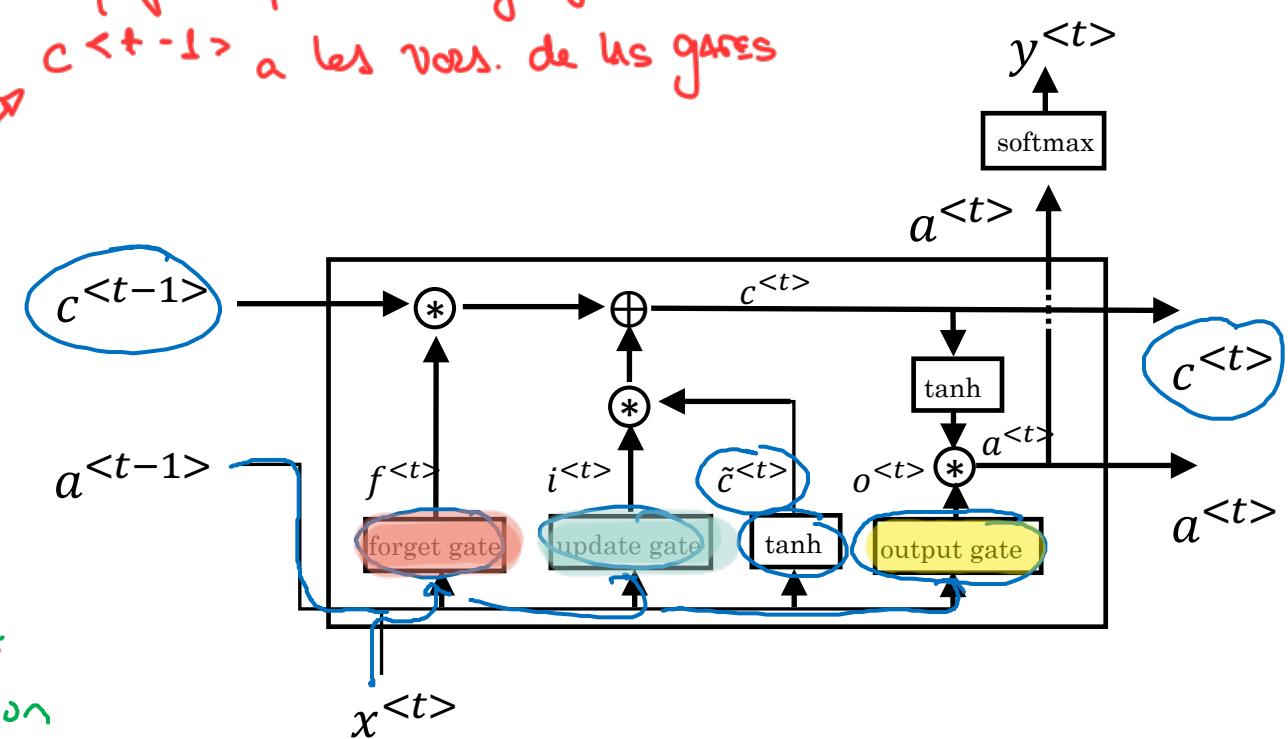
$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * c^{<t>}$$

*peephole connection*



*Hay gente q ademas Agrega el  $c^{<t-1>}$  a los vrs. de las gatres*



→ No hay concurso de cuando usar GRU's y  
LSTMs. conviene probar los dos.



deeplearning.ai

# Recurrent Neural Networks

---

out for lunch

## Bidirectional RNN

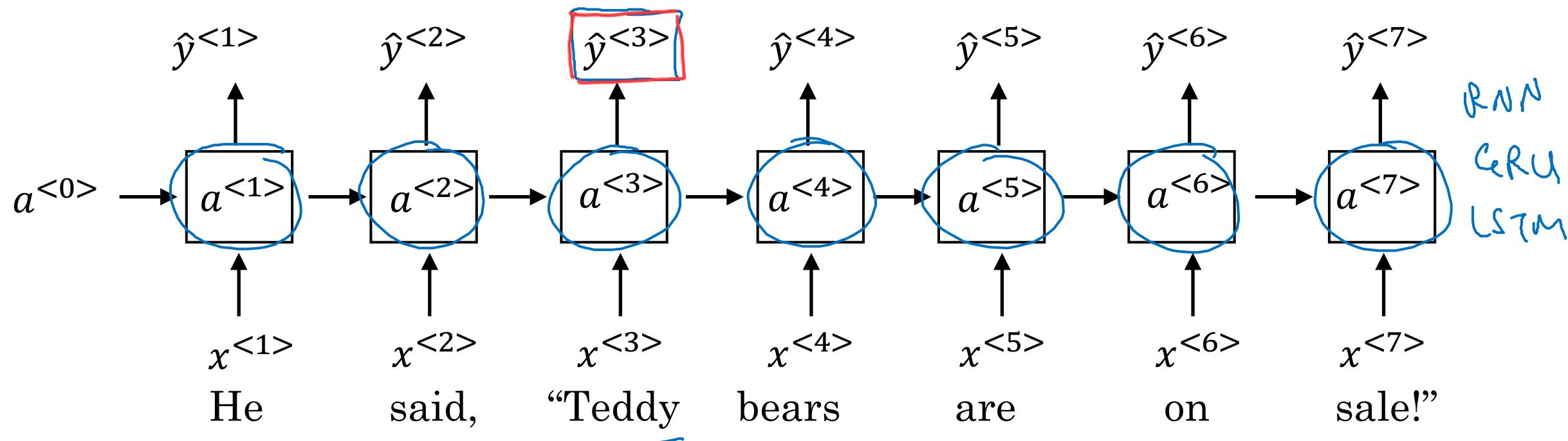
Getting info from the  
future.

# Getting information from the future

He said, "Teddy bears are on sale!"

He said, "Teddy Roosevelt was a great President!"

Bear or President? The answer is in the future.



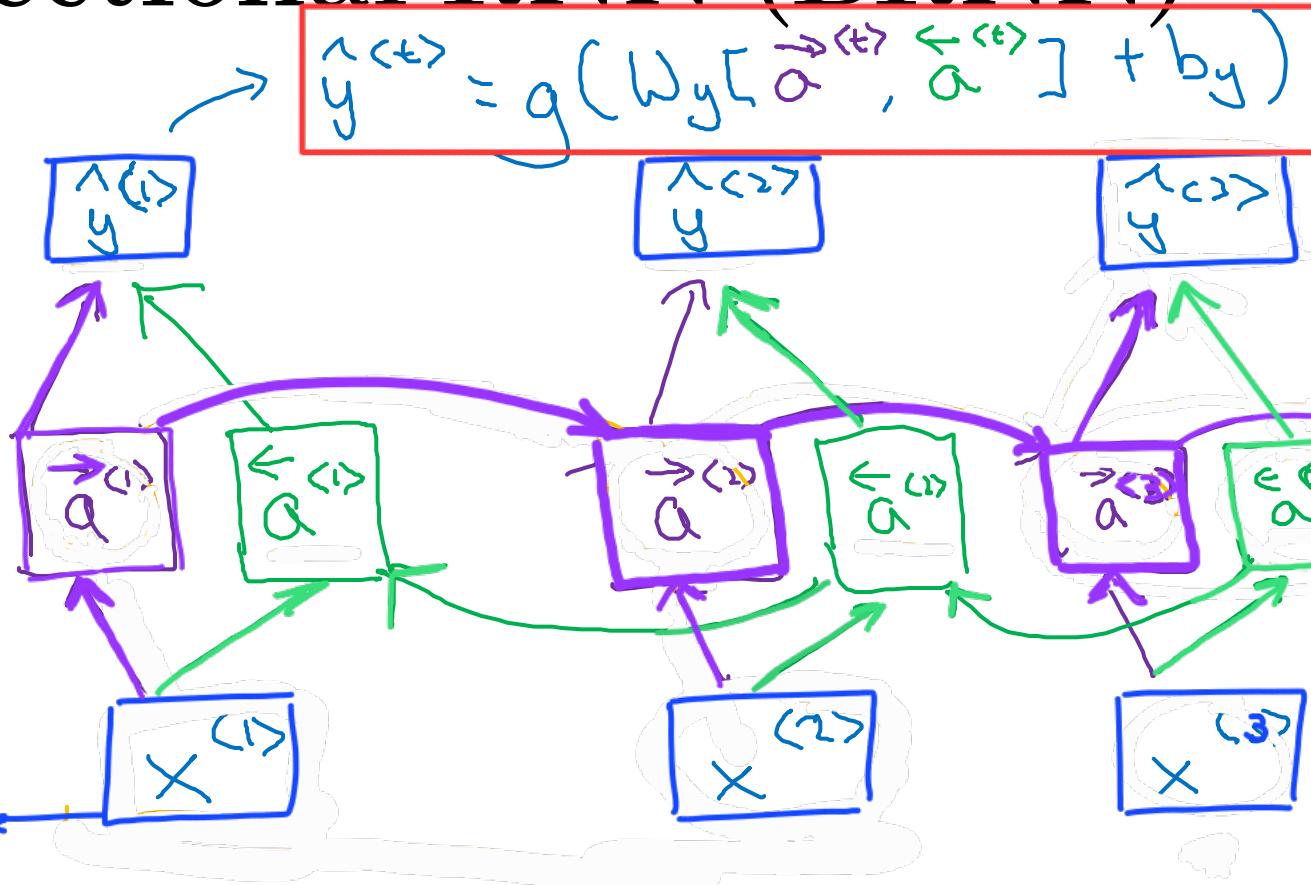
# Bidirectional RNN (BRNN)

los bloques  
puedes ser:

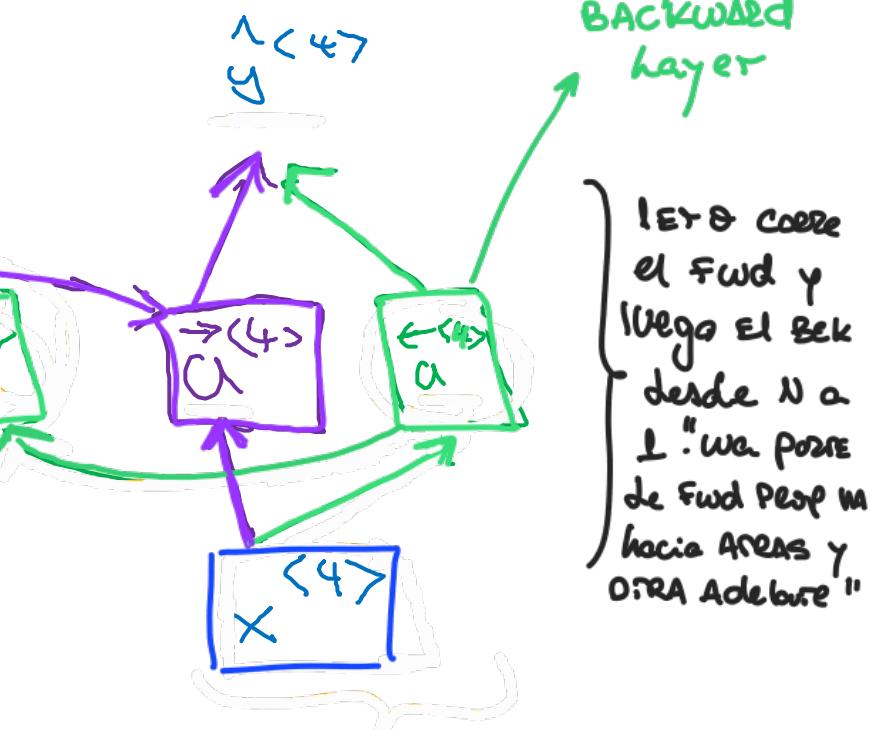
GRU  
LSTM

FORWARD  
Layer  
Componente de  
siempre

Words



En el paso 3, pidej. te llega  
info de  $x_4 \rightarrow a^+$  e info  
de  $x_1 \leftarrow a^-$



Acyclic graph

desventaja  $\Rightarrow$  NECESITAMOS TODA LA SEQ de una pl para la predicción. Pl real tiene que  
esperar pq tiene q esperar a que termine  
de leer la persona.

BRNN w/ LSTM

He said,

"Teddy Roosevelt ..."



deeplearning.ai

# Recurrent Neural Networks

---

Deep RNNs

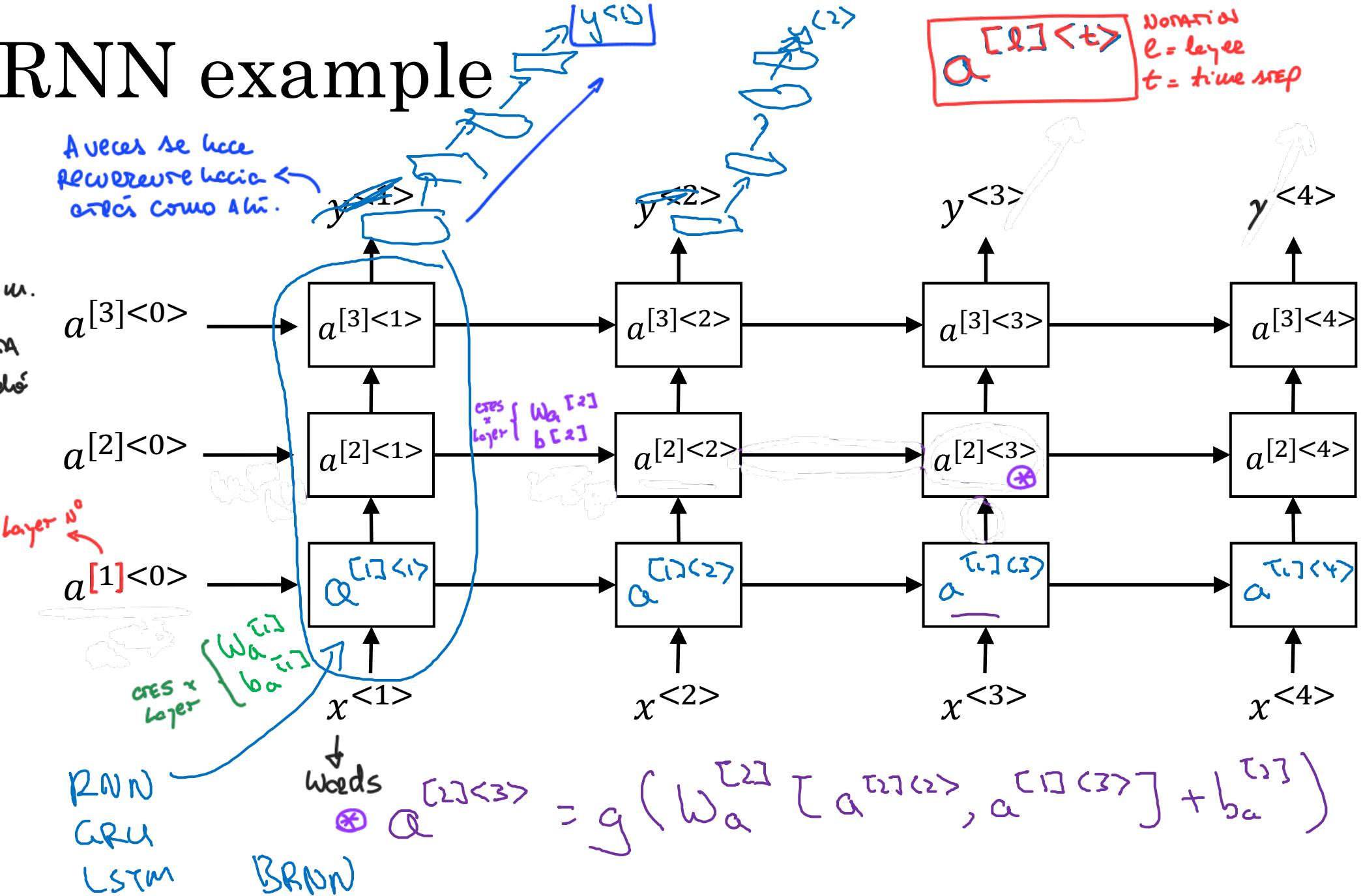
Stack multiple layers of RNN's together

# Deep RNN example

→ NO suelen verse  
muchos más layers  
stacked q 3. pq  
escalon a la gama  
funcion en k t. dim.

→ La VERDAD q esca  
Arch. NO me quedó  
muy close.

A veces se hace  
recurrente hacia  
atrás como ahí.



# Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



deeplearning.ai

# NLP and Word Embeddings

---

## Word representation

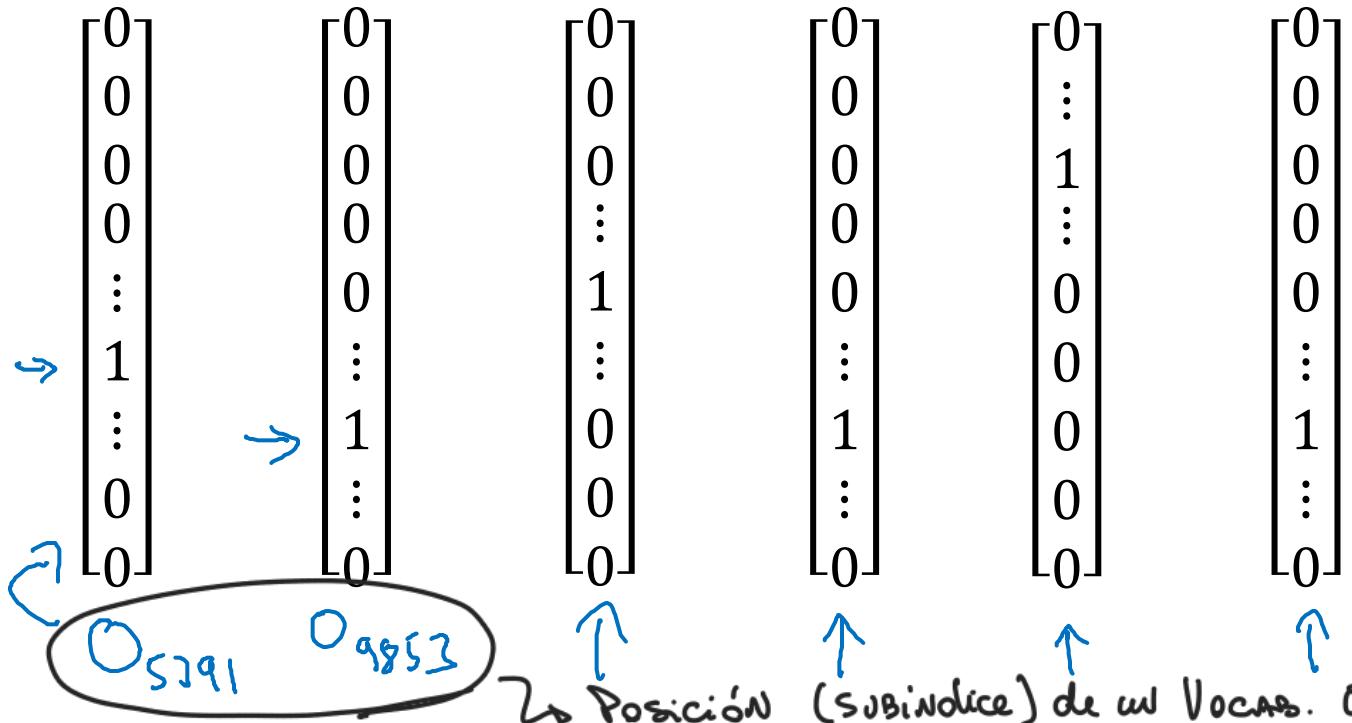
# Word representation

$$V = [a, \text{aaron}, \dots, \text{zulu}, \text{<UNK>}]$$

$$|V| = 10,000$$

## 1-hot representation

Man (5391)   Woman (9853)   King (4914)   Queen (7157)   Apple (456)   Orange (6257)



I want a glass of orange juice.

I want a glass of apple \_\_\_\_.

→ El problema de representar one-hot es q el Algoritmo No tiene forma de saber que Apple y Orange tienen algo q ver. Es q el producto interno de cualquier por one-hot es cero y → su diss. es max.

→ Buscemos otra representación, con Features

# Featurized representation: word embedding

|        | Man<br>(5391) | Woman<br>(9853) | King<br>(4914) | Queen<br>(7157) | Apple<br>(456) | Orange<br>(6257) |
|--------|---------------|-----------------|----------------|-----------------|----------------|------------------|
| Gender | -1            | 1               | -0.95          | 0.97            | 0.00           | 0.01             |
| Royal  | 0.01          | 0.62            | <u>0.93</u>    | <u>0.95</u>     | -0.01          | 0.00             |
| Age    | 0.03          | 0.02            | 0.7            | 0.69            | 0.03           | -0.02            |
| Food   | 0.04          | 0.01            | 0.02           | 0.01            | 0.95           | 0.97             |
| Size   | ⋮             | ⋮               | ⋮              | ⋮               | ⋮              | ⋮                |
| Cost   | ⋮             | ⋮               | ⋮              | ⋮               | ⋮              | ⋮                |
| Verb   | ⋮             | ⋮               | ⋮              | ⋮               | ⋮              | ⋮                |

↑ Features

↑ Weeds

↑ ↗ e<sub>5391</sub>

↑ ↗ e<sub>9853</sub>

Word Embedding ↗

Como los Aprendemos a los Features de C/Palabra? ↗

I want a glass of orange juice ↗

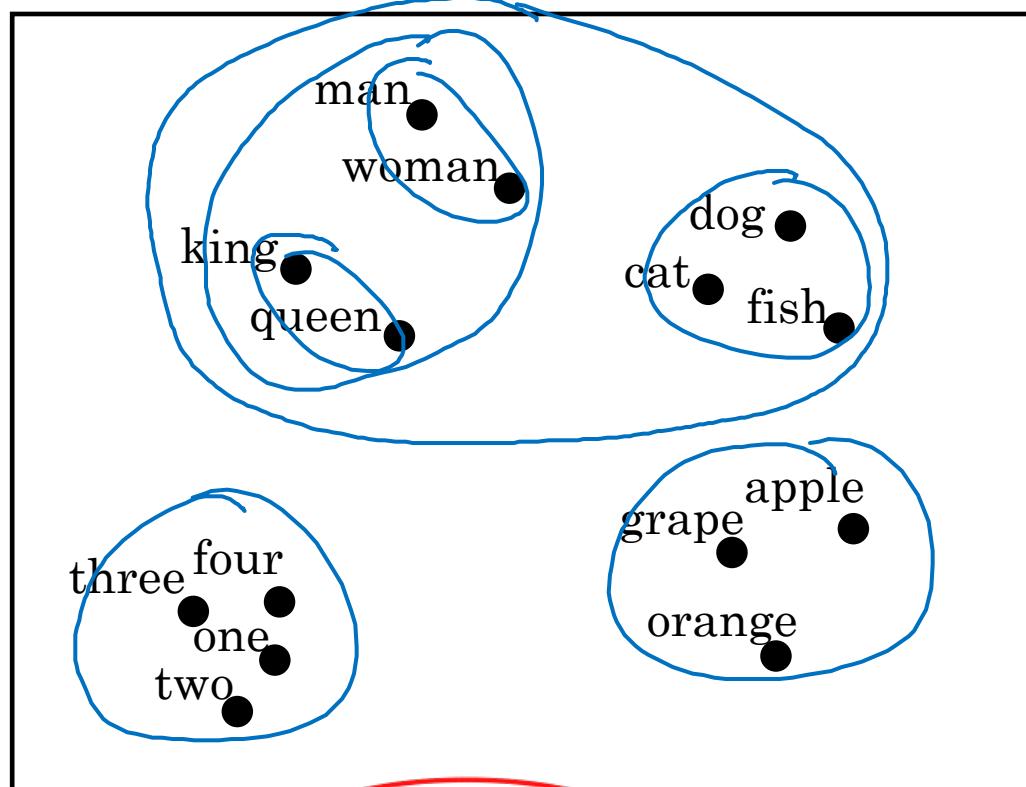
I want a glass of apple juice ↗

Andrew ↗

AHORA SU REPRESENTACIÓN EN ESTE ESPACIO VECTORIAL ES PARECIDA "

# Visualizing word embeddings

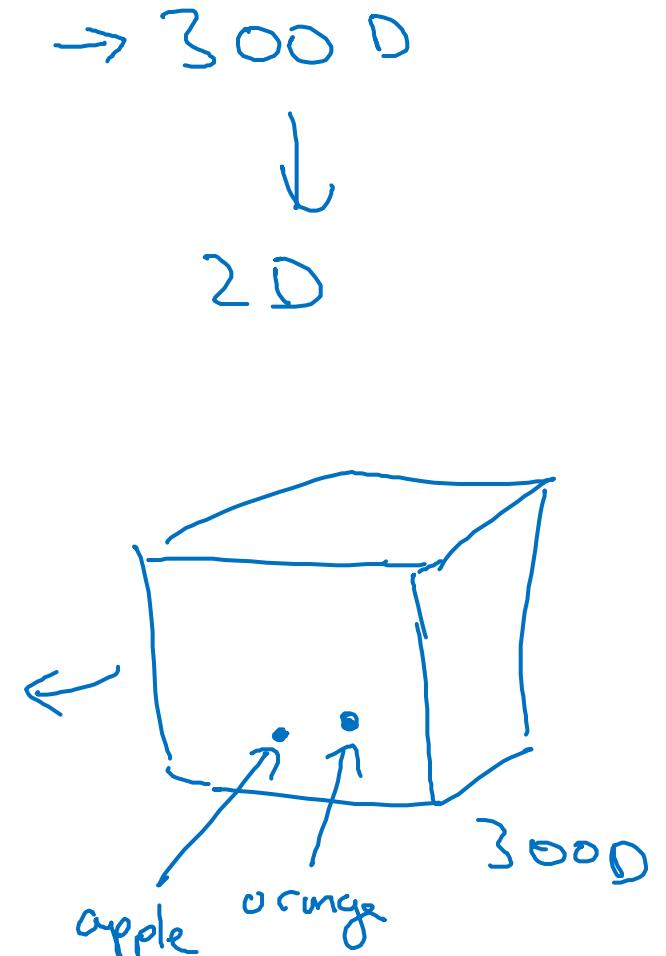
Reducir los vectores  
de embeddings para  
poder interpretarlos



dimensionality  
reduction



t-SNE





deeplearning.ai

# NLP and Word Embeddings

---

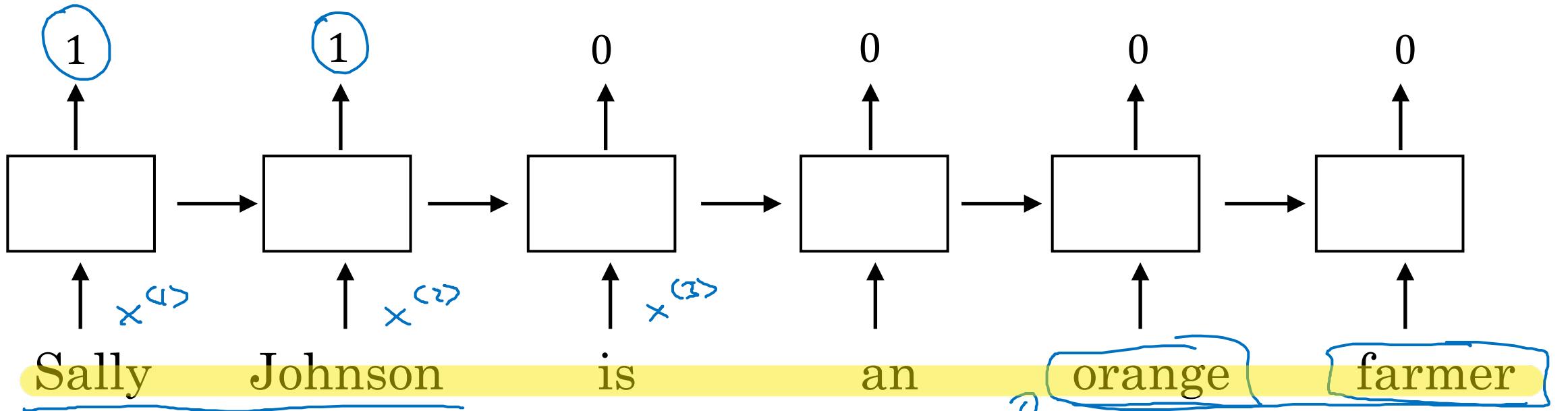
Using word embeddings

\* Diccionario y cultivo de palabras raras q prob.  
NO ESTEN en el training set → pero si q  
W. emb. sea q sou, puede inferir q va  
a faltar q esto, ya que la palabra en si COZGARA'CA  
INFO & FRUSTRADA y q FALTARÉE. Y

# Named entity recognition example

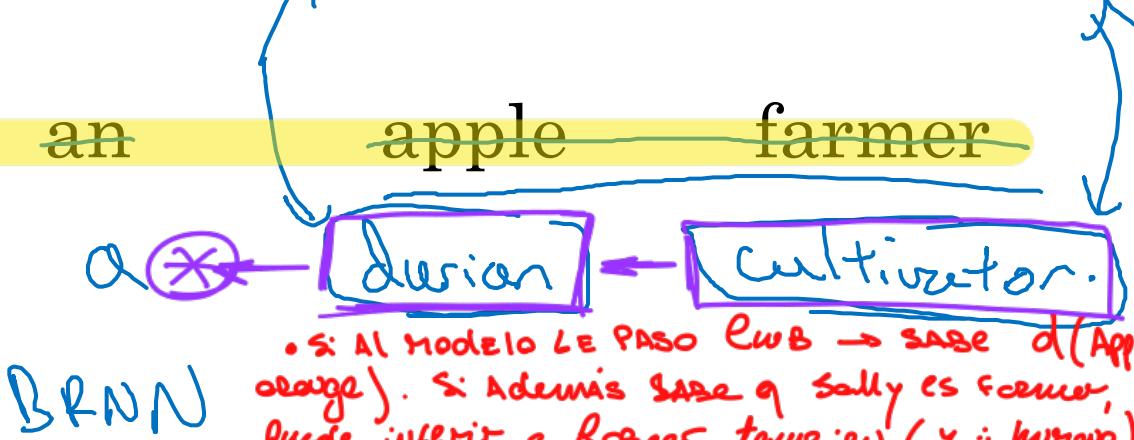
→ Problema que  
TUVE EU Ford

→ DETECTAR NOMBRES de Personas



Robert Lin is an

apple farmer



BRNN

DATASET SIZE  
embedding → 1B words - 100B words  
Named Entity → 100K words  
"transfer learning"

Andrew Ng

# Transfer learning and word embeddings

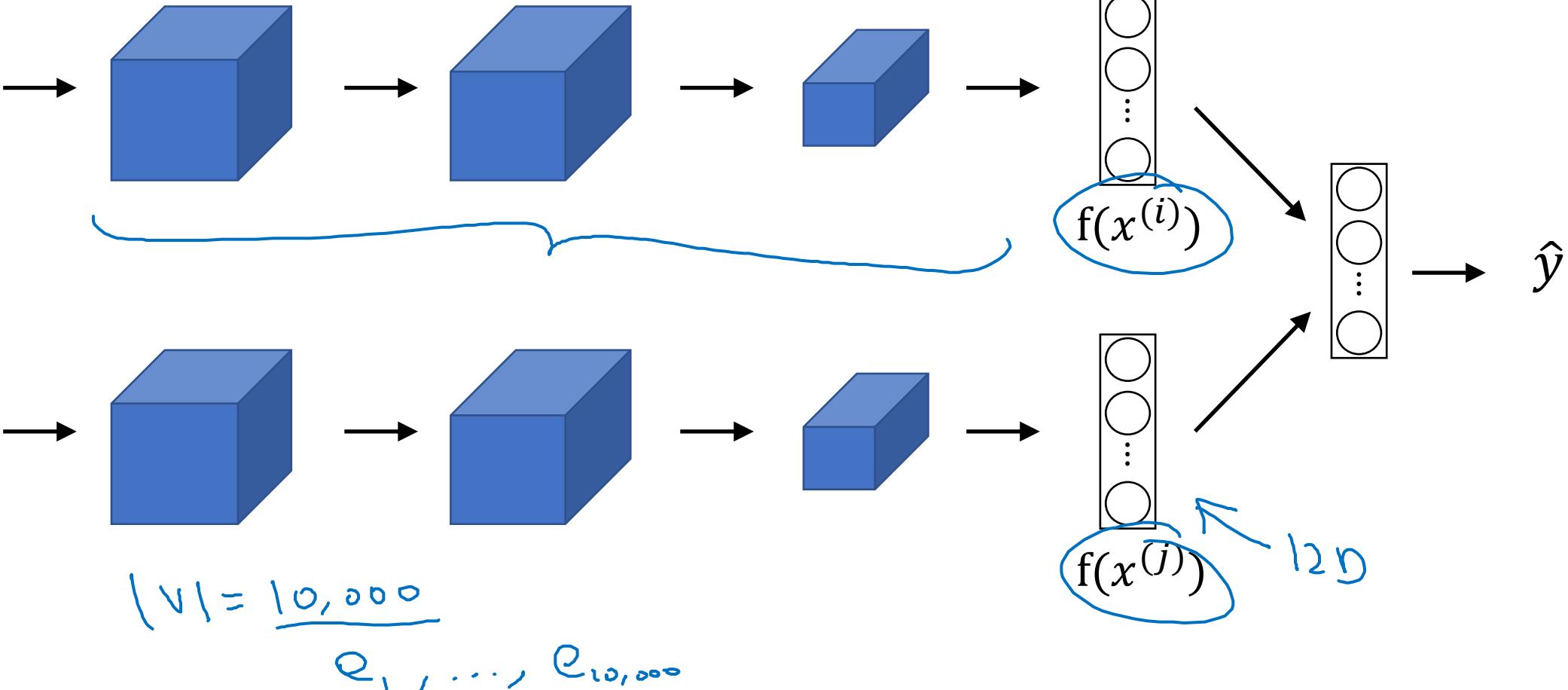
- A 1. Learn word embeddings from large text corpus. (1-100B words)  
(Or download pre-trained embedding online.)
- B 2. Transfer embedding to new task with smaller training set.  
(say, 100k words)



3. Optional: Continue to finetune the word embeddings with new data.

→ Para lo que, los embeddings son  
claves de tareas que no tenemos datasets  
grandes.

# Relation to face encoding (embedding) 128D





deeplearning.ai

# NLP and Word Embeddings

---

## Properties of word embeddings

# Analogies

|        | Man<br>(5391) | Woman<br>(9853) | King<br>(4914) | Queen<br>(7157) | Apple<br>(456) | Orange<br>(6257) |
|--------|---------------|-----------------|----------------|-----------------|----------------|------------------|
| Gender | -1            | 1               | -0.95          | 0.97            | 0.00           | 0.01             |
| Royal  | 0.01          | 0.02            | 0.93           | 0.95            | -0.01          | 0.00             |
| Age    | 0.03          | 0.02            | 0.70           | 0.69            | 0.03           | -0.02            |
| Food   | 0.09          | 0.01            | 0.02           | 0.01            | 0.95           | 0.97             |

$$\begin{matrix} e_{5391} \\ e_{\text{man}} \end{matrix}$$

$$\underline{\text{Man} \rightarrow \text{Woman}}$$

$$e_{\text{man}} - e_{\text{woman}}$$

$$e_{\text{woman}}$$

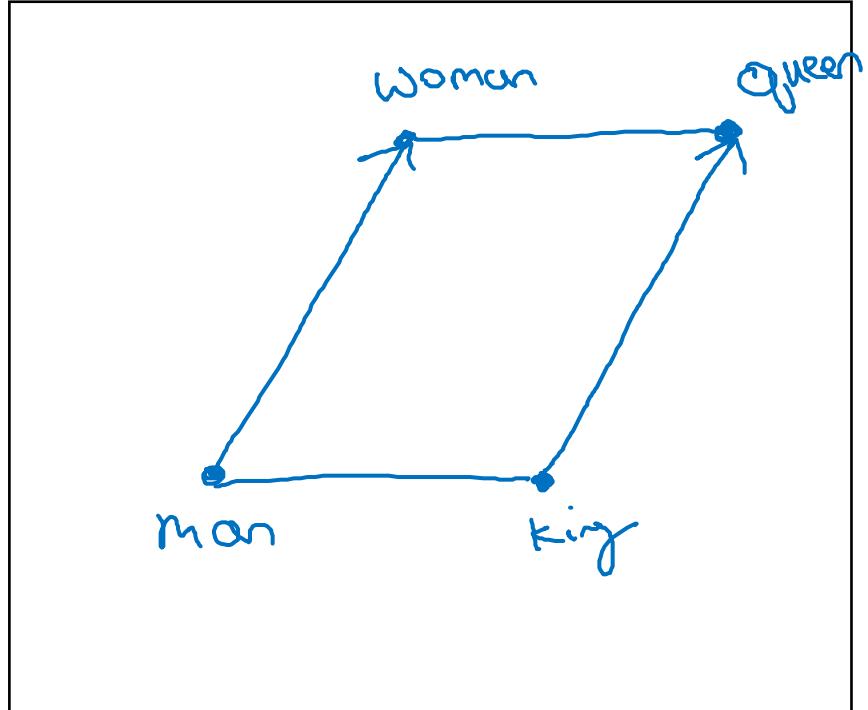
$$\underline{\text{King} \rightarrow ? \text{ Queen}}$$

$$e_{\text{king}} - e_{? \text{ Queen}}$$

$$\underline{e_{\text{man}} - e_{\text{woman}}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\underline{e_{\text{king}} - e_{\text{queen}}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

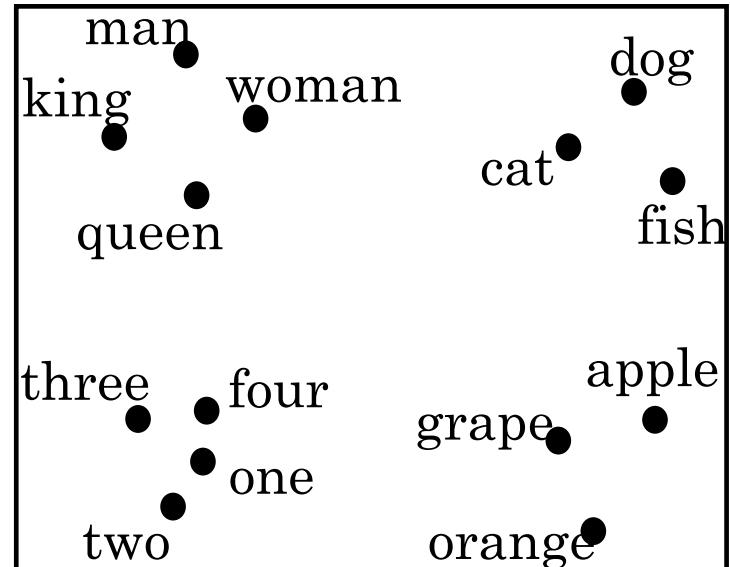
# Analogies using word vectors



300 D

Find word  $w_i: \arg \max_w$

$300D \rightarrow 2D$



$$e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\underline{\text{?}}} \quad e_w$$

$$\underbrace{\hspace{10cm}}$$

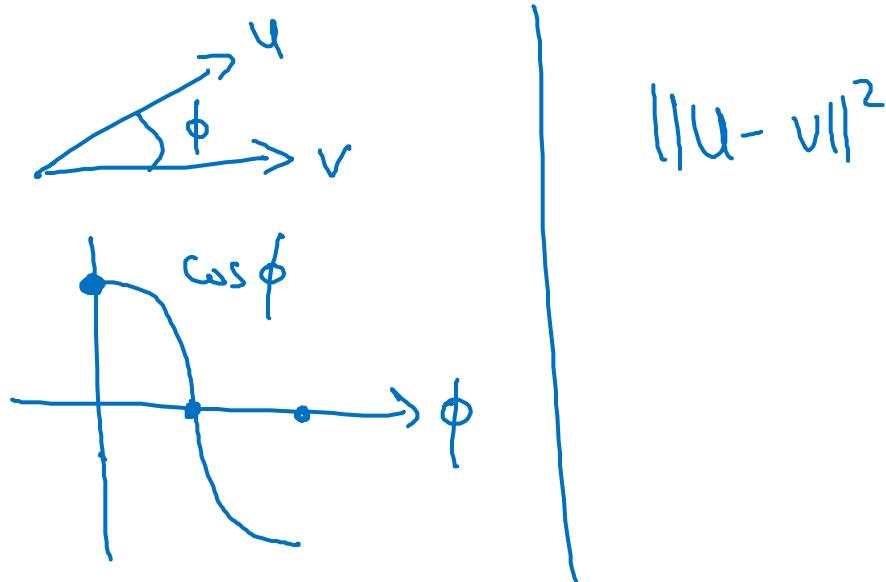
$$\boxed{\text{Sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})}$$

30 - 75%

# Cosine similarity

$$\rightarrow \boxed{\text{sim}(e_w, e_{king} - e_{man} + e_{woman})}$$

$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$



Man:Woman as Boy:Girl  
Ottawa:Canada as Nairobi:Kenya  
Big:Bigger as Tall:Taller  
Yen:Japan as Ruble:Russia



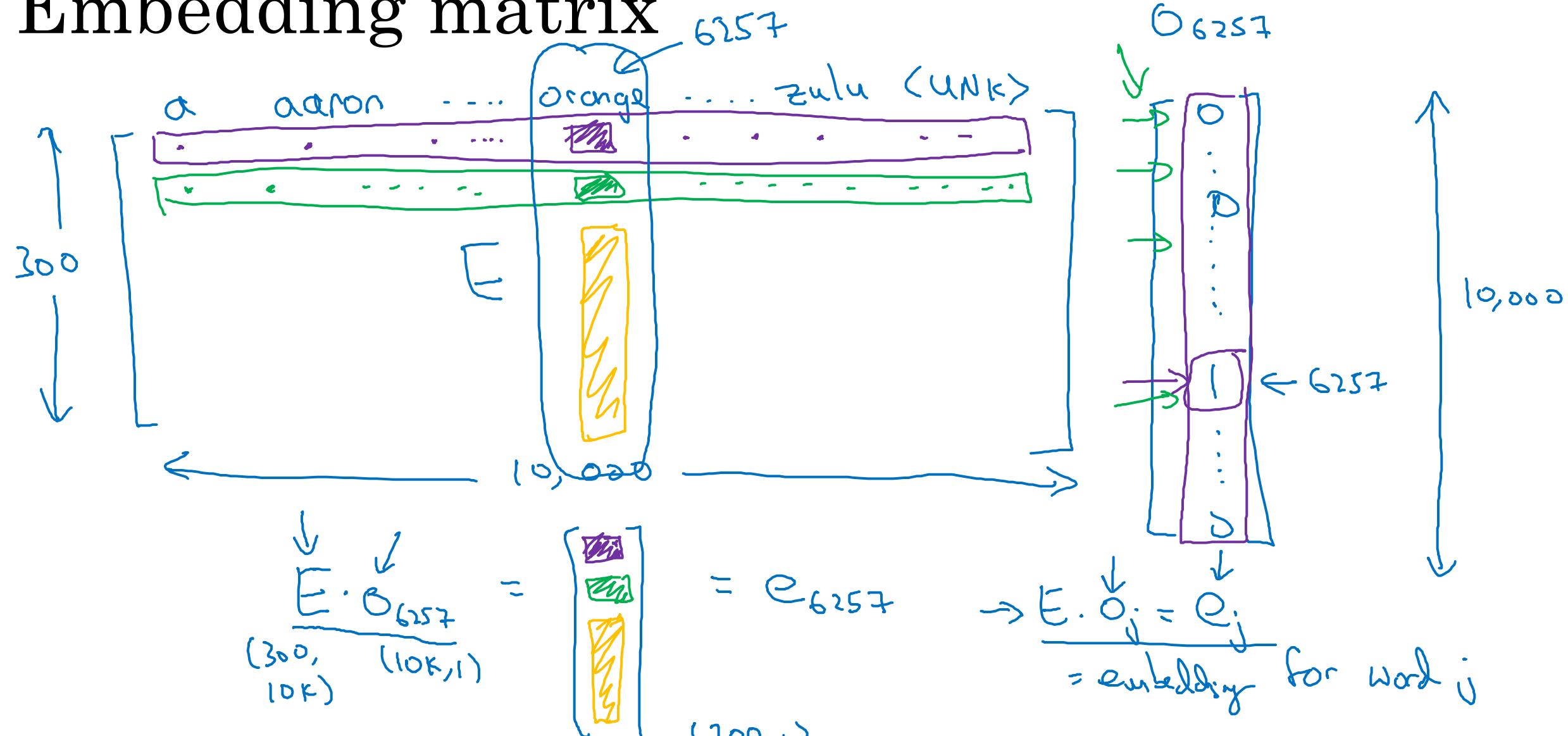
deeplearning.ai

# NLP and Word Embeddings

---

## Embedding matrix

# Embedding matrix



In practice, use specialized function to look up an embedding.  
→ Embedding



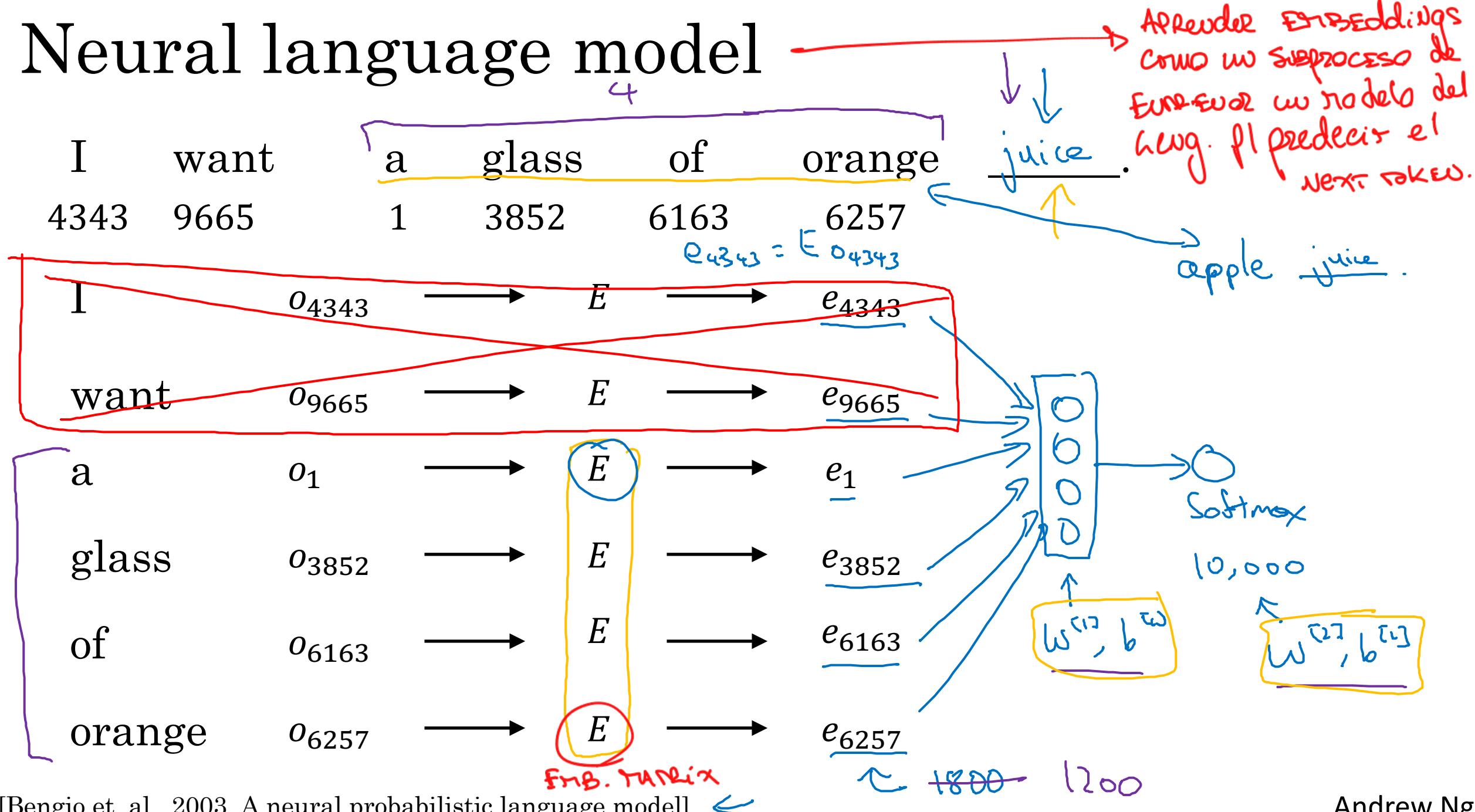
deeplearning.ai

# NLP and Word Embeddings

---

## Learning word embeddings

# Neural language model



# Other context/target pairs

I want a **glass** of **orange** juice to go along with my cereal.

Context: Last 4 words.

4 words on left & right

Last 1 word

Nearby 1 word

skip gram

a glass of orange ? to go along with

orange ?

glass . ?



deeplearning.ai

# NLP and Word Embeddings

---

## Word2Vec

# Skip-grams

I want a glass of orange juice to go along with my cereal.

new. ANTES  
pero con w CONTEXTOS  
nacho mas chico  $\rightarrow$  1  
otra palabra q una  
dissancia skipped!

Lo para Aprender embeddings ve joya.

Context

orange

orange

orange

Target

juice

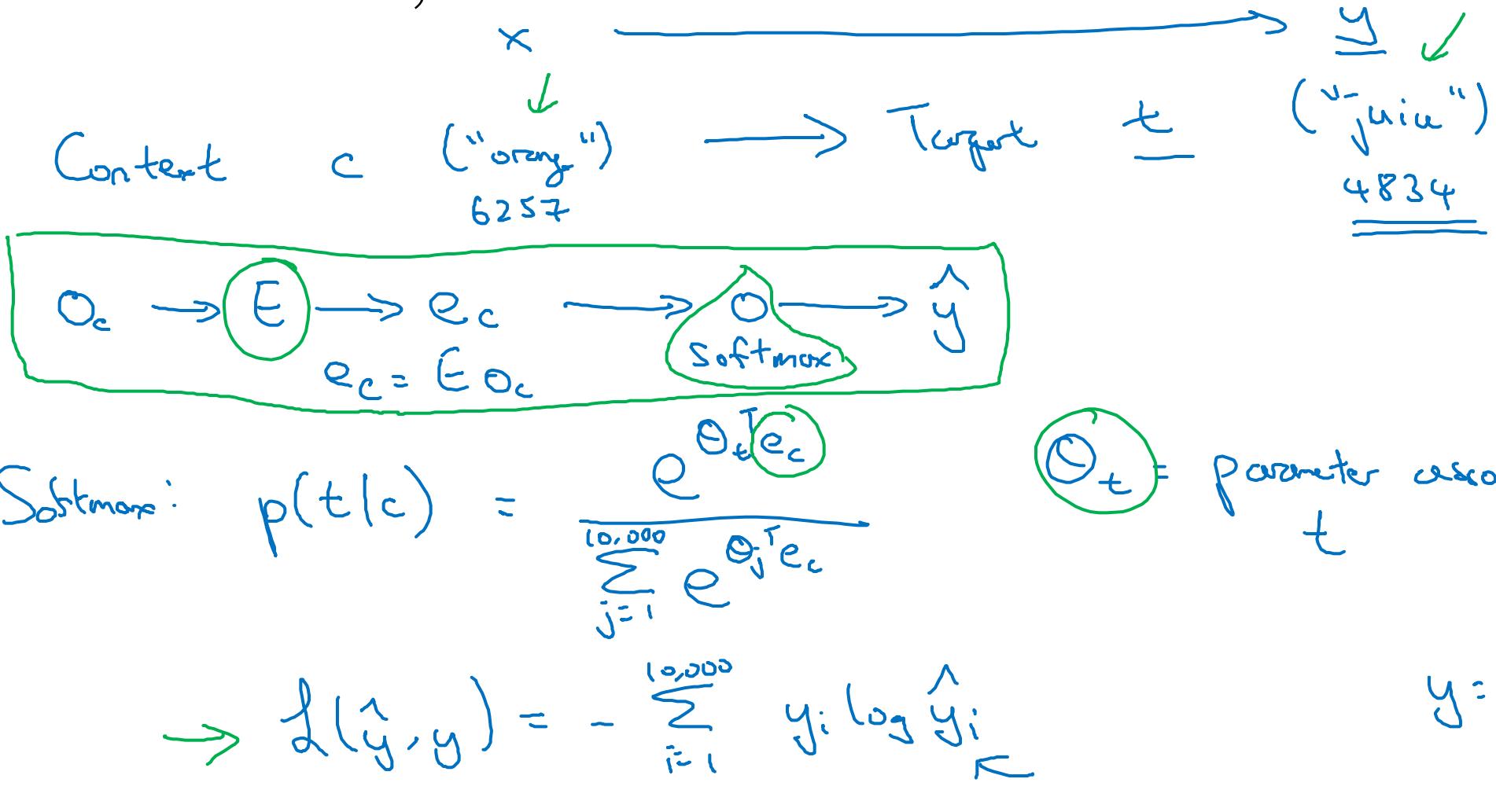
glass

my



# Model

Vocab size = 10,000k

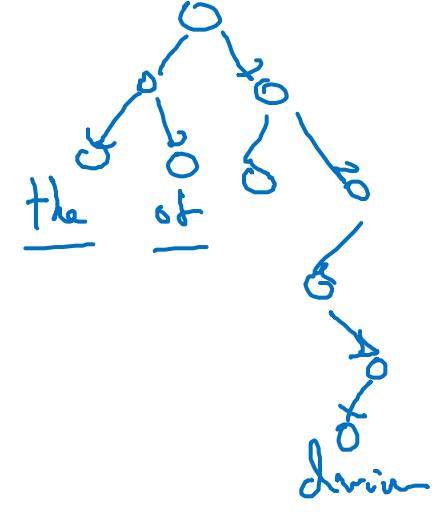
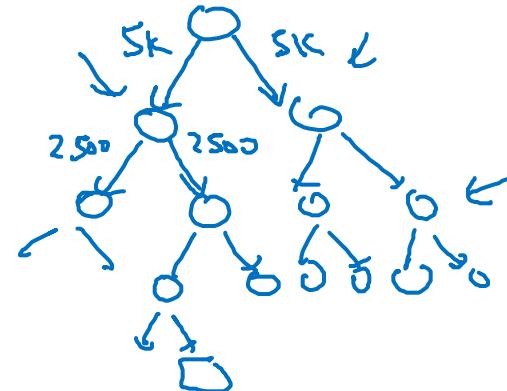


# Problems with softmax classification

$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

$\log |V|$

Hierarchical softmax.



How to sample the context  $c$ ?

→ the, of, a, and, to, ...

→ orange, apple, durian

$P_{\text{durian}}$

$t$   
 $c \rightarrow t$

$P(c)$



deeplearning.ai

# NLP and Word Embeddings

---

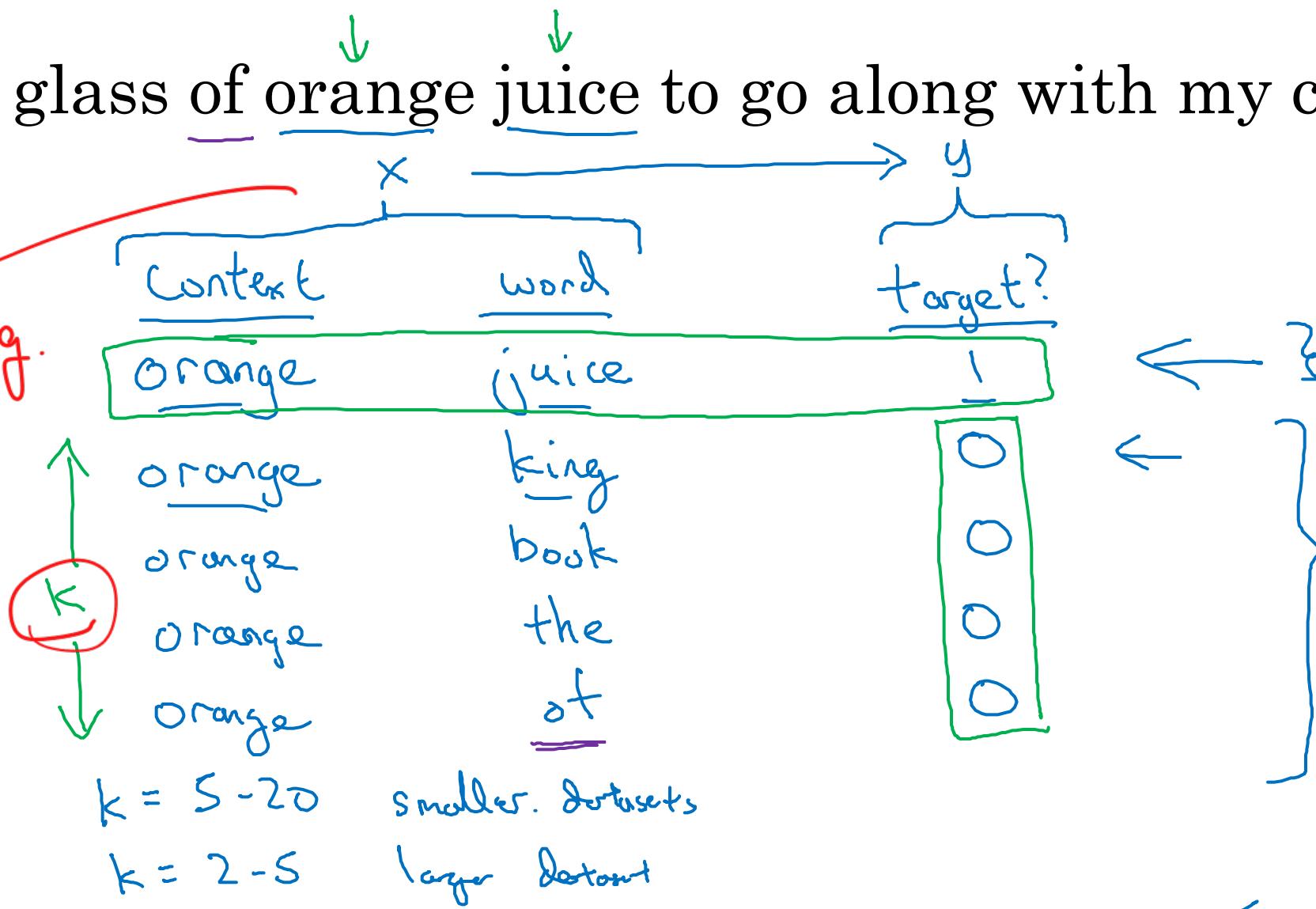
## Negative sampling

Para resolver la escala de la  
Softmax

# Defining a new learning problem

I want a glass of orange juice to go along with my cereal.

To convert this  
in Sup. learning.



## Model

# Softmax:

$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

*10,000 ways  
of forming*

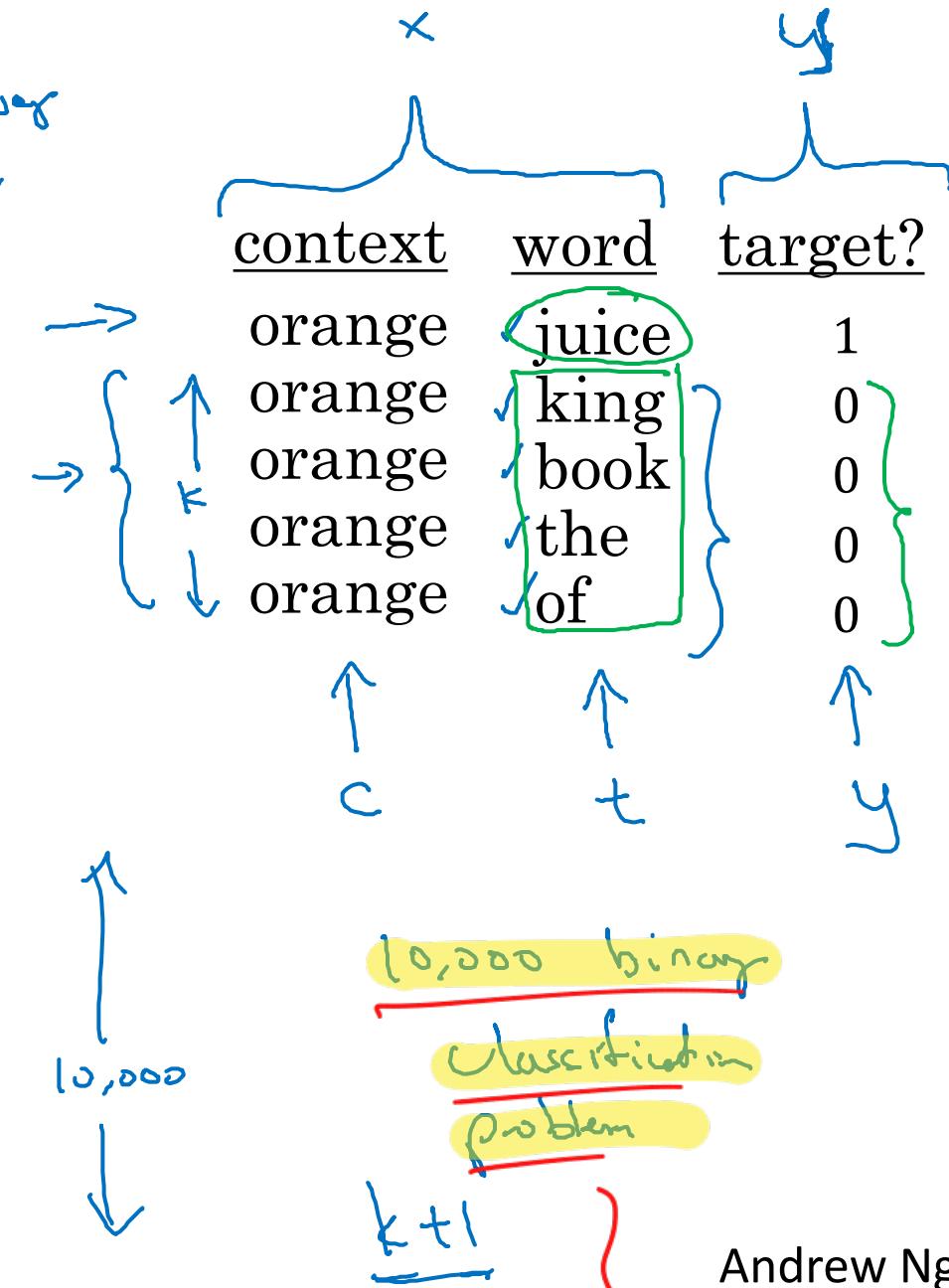
$$P(y=1 | c, t) = \sigma(\theta_t^T e_c) \leftarrow$$

Orange  
6257

$$O_{6257} \rightarrow E \rightarrow e_{6257}$$

juice

king



# Selecting negative examples

|     | <u>context</u> | <u>word</u> | <u>target?</u> |
|-----|----------------|-------------|----------------|
| pos | orange         | juice       | 1              |
| neg | orange         | king        | 0              |
|     | orange         | book        | 0              |
|     | orange         | the         | 0              |
|     | orange         | of          | 0              |

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10,000} f(w_j)^{3/4}}$$

Prob q recomienda el paper pl samplez palabras. de ultima leez del paper.

$$\frac{1}{|V|}$$

the , of, and, ... 10k, Espanol

( k+1 ) clasificadores binarios q usan  
funcion mas basico.

K neg Random p  
1 pos posita.

EN VEZ de usarlo  
usar softmax ?!



$k_{neg}$

$t$



deeplearning.ai

# NLP and Word Embeddings

---

## GloVe word vectors

FACE DETECTION

yolo face detector → embedder → Cosine Similarity?  
↓  
Vs. classifier?

# GloVe (global vectors for word representation)

I want a glass of orange juice to go along with my cereal.

$$c, t \rightarrow \begin{array}{l} t = \text{target} \\ c = \text{context} \end{array}$$

$$x_{ij} = \# \text{ times } i \text{ appears in context of } j.$$

$x_{ij}$  = # times  $i$  appears in context of  $j$ .  
 $c = t$        $i = t$        $j = c$

$x_{ij} = x_{ji}$   $\leftarrow$  Esto va a depender de como definimos el contexto.

Cuenta q tan seguido las palabras  $i, j$  aparecen juntas en el contexto

# Model

minimize

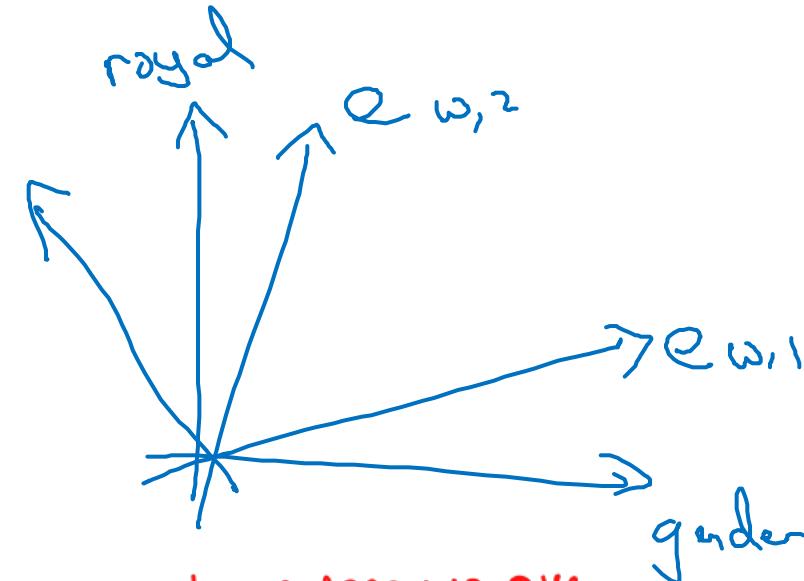
$$\sum_{i=1}^{10,000} \sum_{j=1}^{100,000} f(x_{ij}) (\theta_i^T e_j + b_i + b_j) - \log \frac{1}{x_{ij}} + \lambda \theta_i^T e_c$$

Annotations:

- $\lambda$ ? index. donde  $x_{ij}=0$  no sumar! por eso usamos  $\lambda$
- $\theta_i^T e_c$  weight $x_{ij}$   $\Rightarrow$  weighting factor term
- $f(x_{ij}) = 0$  of  $x_{ij} = 0$ . \theta \log 0 = 0 \Rightarrow No sumo si  $x_{ij} = 0$ .
- this is of a ...
- dorian \Rightarrow Palabra poco frec. Le da "peso" final.
- + detalles en el paper.
- $\theta_w$  (final) =  $\frac{\theta_w + \bar{\theta}_w}{2}$  \Rightarrow Random init  $\theta_w$  y  $\bar{\theta}_w$ , Entreno con  $\bar{\theta}_w$  y tomo el promedio al final como final.

# A note on the featurization view of word embeddings

|        | Man<br>(5391) | Woman<br>(9853) | King<br>(4914) | Queen<br>(7157) |
|--------|---------------|-----------------|----------------|-----------------|
| Gender | -1            | 1               | -0.95          | 0.97            |
| Royal  | 0.01          | 0.02            | 0.93           | 0.95            |
| Age    | 0.03          | 0.02            | 0.70           | 0.69            |
| Food   | 0.09          | 0.01            | 0.02           | 0.01            |



⇒ No podemos asegurarnos que  
esta dimensión tenga una interpretabilidad. Sin embargo,  
las analogías y similitudes se mantienen.

$$\text{minimize } \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) (\underbrace{\theta_i^T e_j}_{\text{feature vector}} + b_i - b_j' - \log X_{ij})^2$$

$$(A\theta_i)^T (A^T e_j) = \cancel{\theta_i^T A^T A} \cancel{\theta_i^T} e_j$$



deeplearning.ai

# NLP and Word Embeddings

---

## Sentiment classification

# Sentiment classification problem

→ god: Cn Embeddings poderos  
bridesor Seuivre classifiers Cn  
para labels!



The dessert is excellent.



Service was quite slow.



Good for a quick meal, but nothing special.



Completely lacking in good taste, good service, and good ambience.



10,000 → 100,000 words

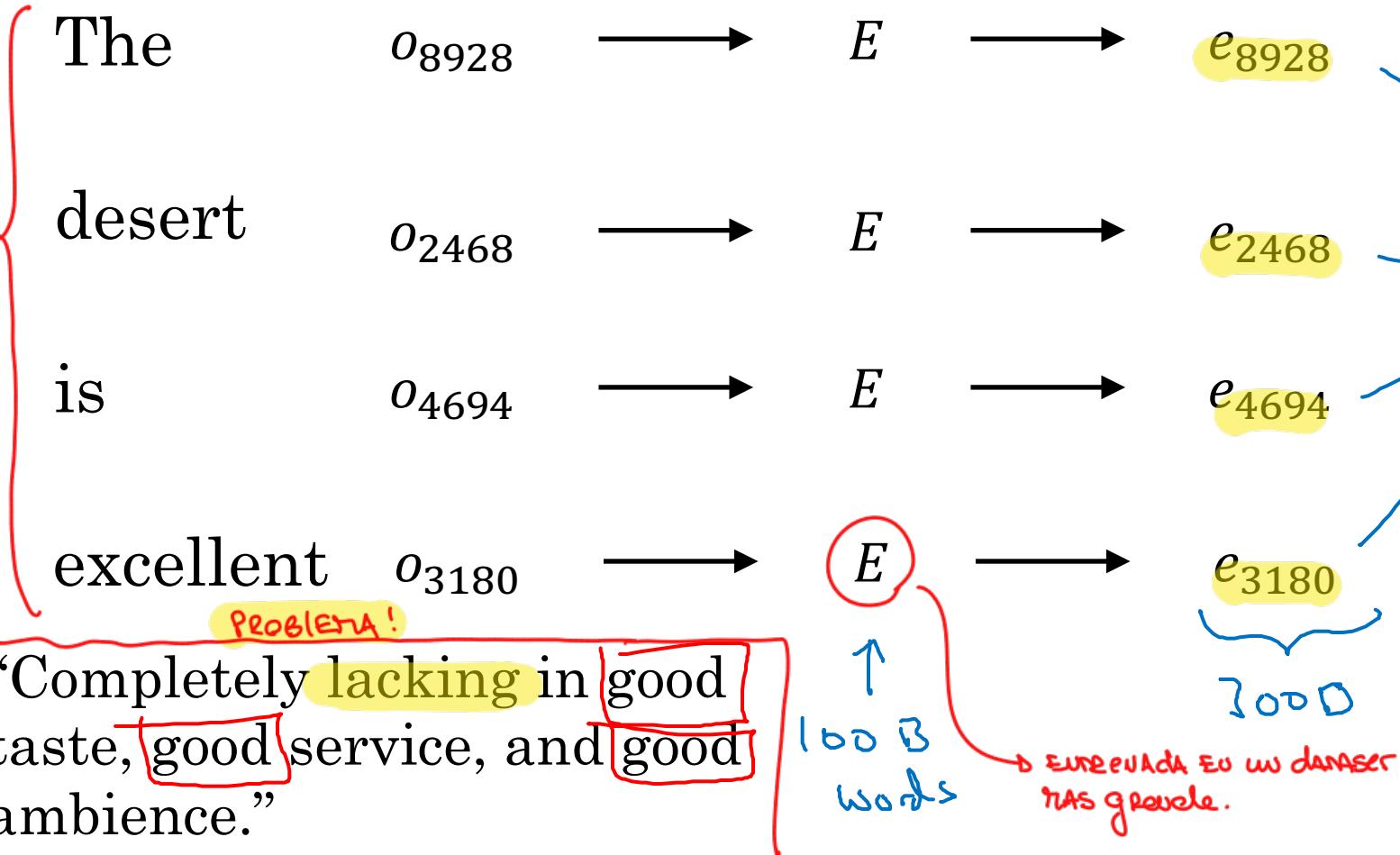
# Simple sentiment classification model

The dessert is excellent

8928 2468 4694 3180



Embeddings



$e_{8928}$

$e_{2468}$

$e_{4694}$

$e_{3180}$

USAR el promedio TE independiza del Largo.

Avg o Sum

Ang

300D

clasificador.

Softmax

1-5

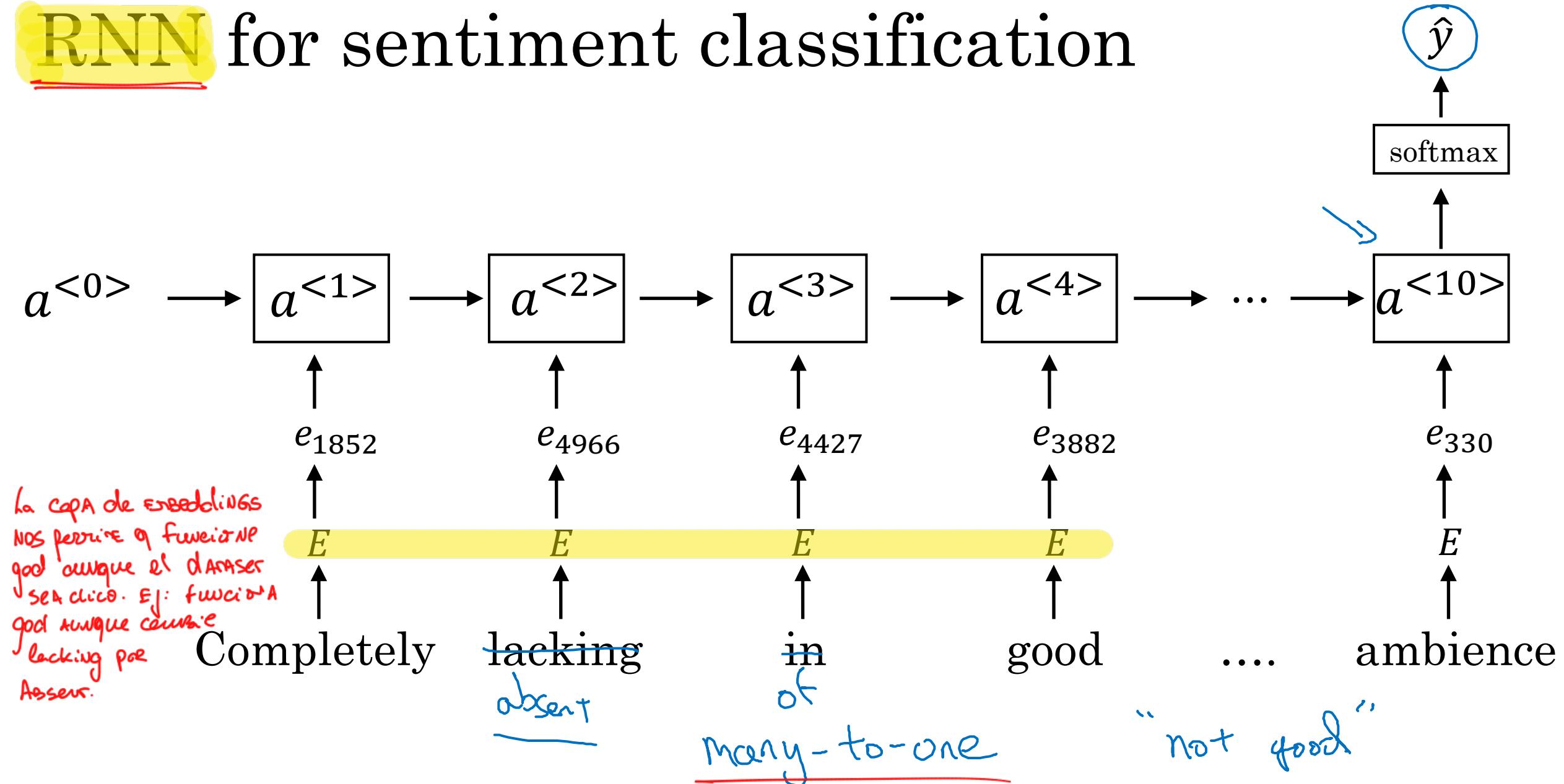
→ Lo malo es que ignora el orden!

Una review ~~4\*~~ Así sería la clasificación en ese modelo! ↴

Una solución es usar una RNN! ↴

Andrew Ng

# RNN for sentiment classification





deeplearning.ai

# NLP and Word Embeddings

---

## Debiasing word embeddings

# The problem of bias in word embeddings

gender, race, etc.

Man:Woman as King:Queen

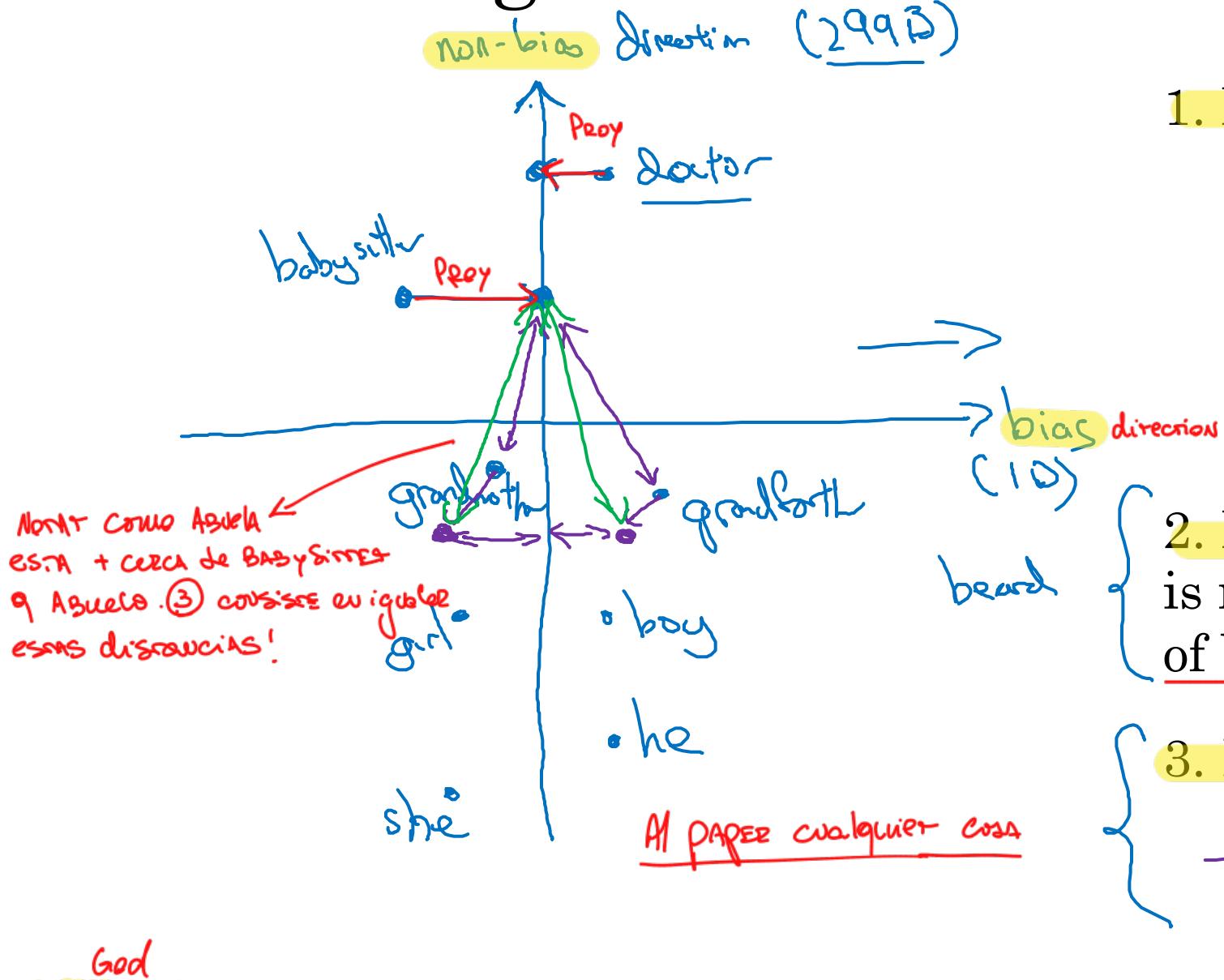
Man:Computer\_Programmer as Woman:Homemaker !

Father:Doctor as Mother:Nurse !

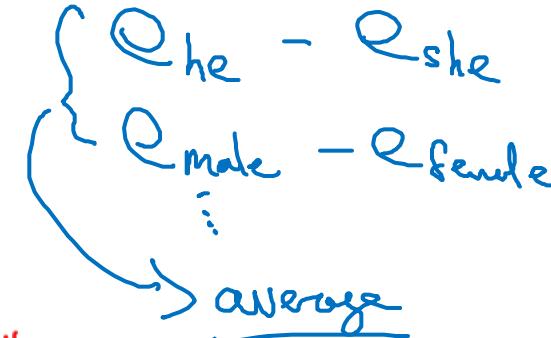
Word embeddings can reflect gender, ethnicity, age, sexual orientation, and other biases of the text used to train the model.



# Addressing bias in word embeddings $\Rightarrow \boxed{\text{God}}$



1. Identify bias direction.



\* PALABRAS q NO tienen  
pasado el género,  
(f Abuela, Abuelo  
x ej)

2. Neutralize: For every word that is not definitional, project to get rid of bias. \*

3. Equalize pairs. (3) con ALGEBRA basically

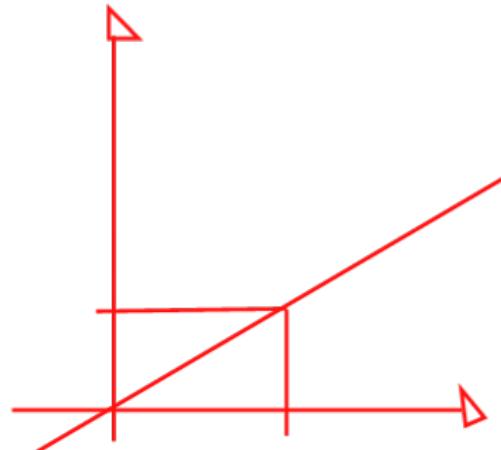
$$\begin{aligned} \rightarrow \text{grandmother} &\leftrightarrow \text{grandfather} \\ \text{girl} &\leftrightarrow \text{boy} \end{aligned}$$

# Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



U.  
Lucas  
Alejandro  
Argento  
DNI: 40540183

AÑO 2024  
siglo XXI.  
playas  
tierra



deeplearning.ai

# Sequence to sequence models

*; Translation!*

---

## Basic models

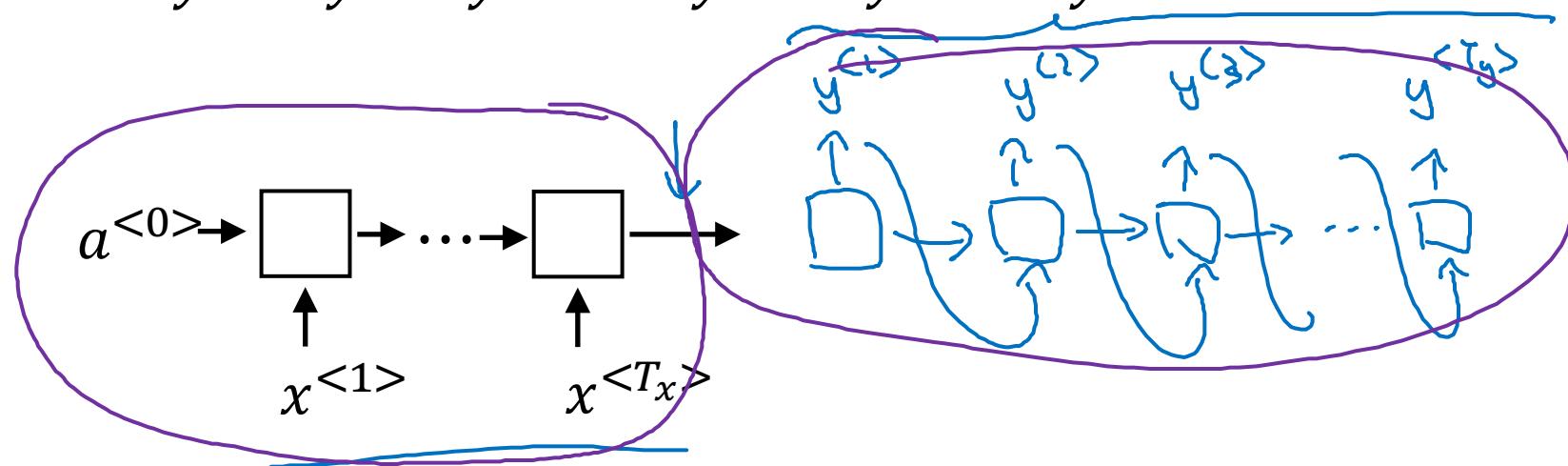
# Sequence to sequence model

$$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$$

Jane visite l'Afrique en septembre

→ Jane is visiting Africa in September.

$$y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad y^{<4>} \quad y^{<5>} \quad y^{<6>}$$

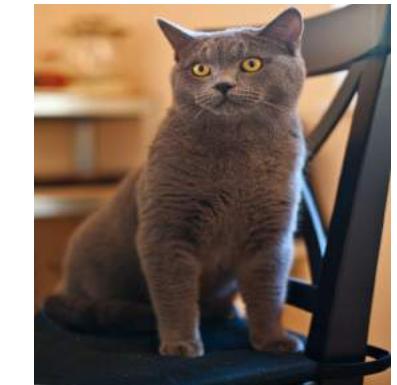


[Sutskever et al., 2014. Sequence to sequence learning with neural networks] ↪

[Cho et al., 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation] ↪

Andrew Ng

# Image captioning



AlexNet

$11 \times 11$   
s = 4

$55 \times 55 \times 96$

MAX-POOL  
 $3 \times 3$   
s = 2

$27 \times 27 \times 96$

$5 \times 5$   
same

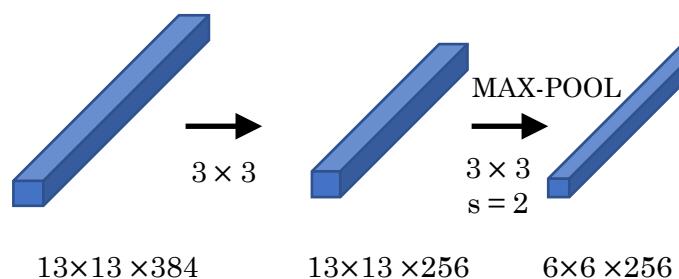
$27 \times 27 \times 256$

MAX-POOL  
 $3 \times 3$   
s = 2

$13 \times 13 \times 256$

$3 \times 3$   
same

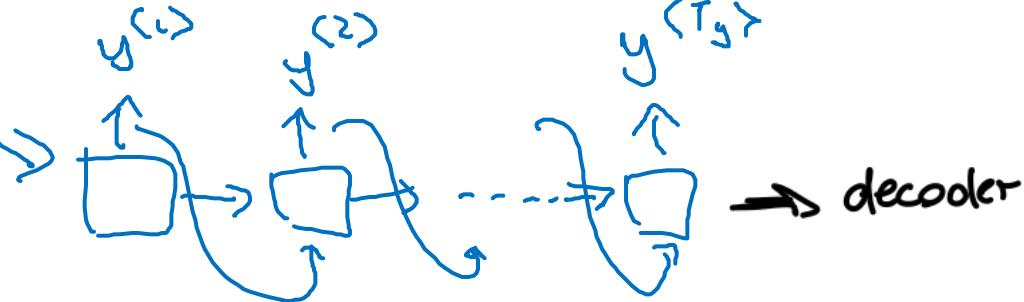
$13 \times 13 \times 384$



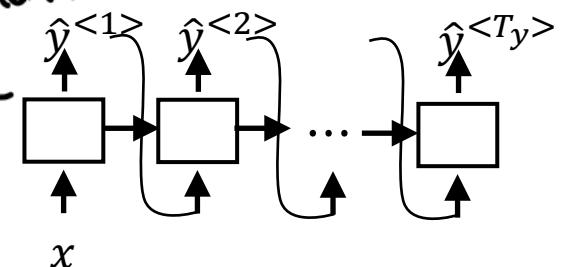
$9216$

$4096$

$4096$



esta en AlexNet  
Los sacamos p  
usarla de encoder



[Mao et. al., 2014. Deep captioning with multimodal recurrent neural networks]

[Vinyals et. al., 2014. Show and tell: Neural image caption generator]

[Karpathy and Li, 2015. Deep visual-semantic alignments for generating image descriptions]

Andrew Ng



deeplearning.ai

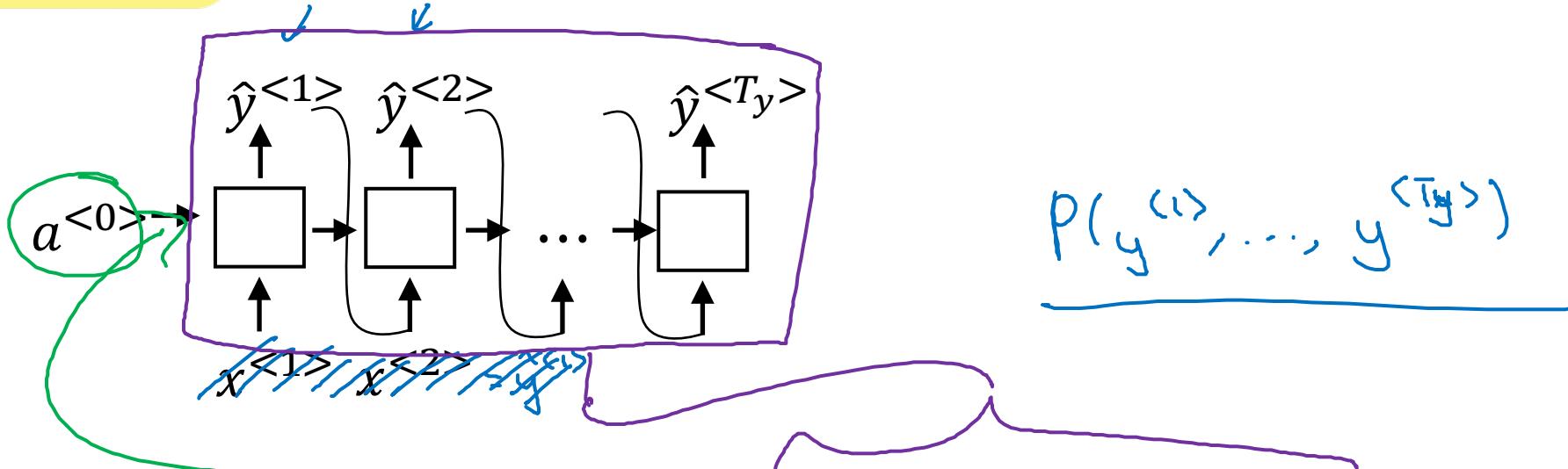
# Sequence to sequence models

---

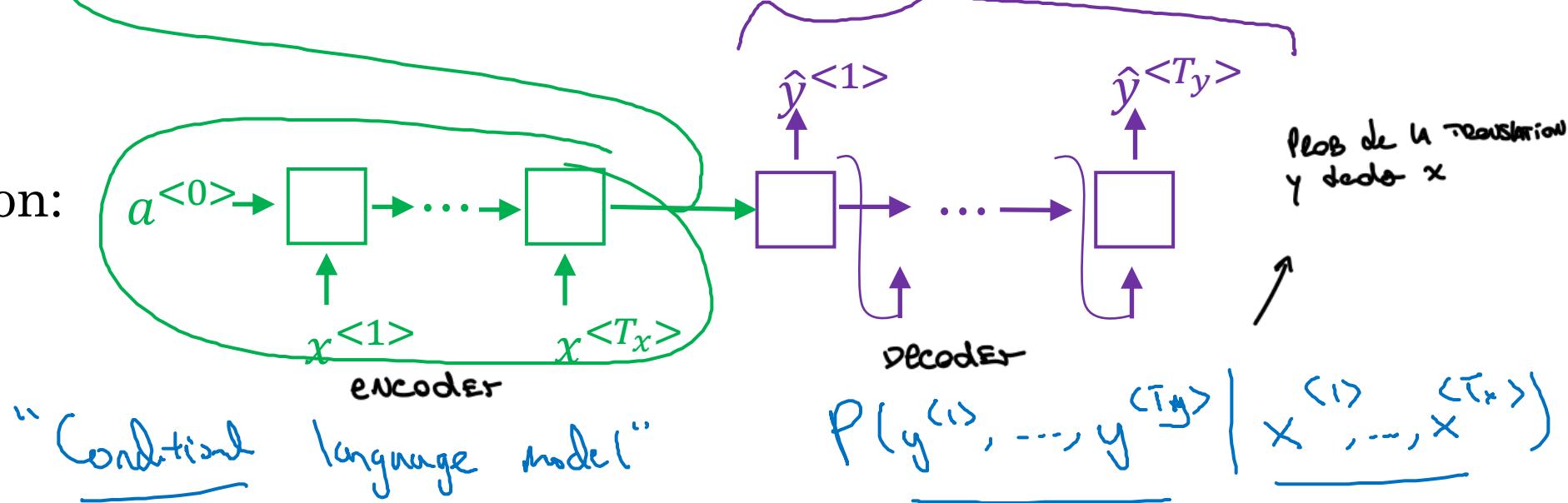
Picking the most  
likely sentence

# Machine translation as building a conditional language model

Language model:



Machine translation:



Andrew Ng

# Finding the most likely translation

Jane visite l'Afrique en septembre.

$$P(y^{<1>} , \dots , y^{<T_y>} | x)$$

French  
↓  
English

- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in September.
- In September, Jane will visit Africa.
- Her African friend welcomed Jane in September.

outputs del modelo  
Cuál elegirás y  
Cómo?

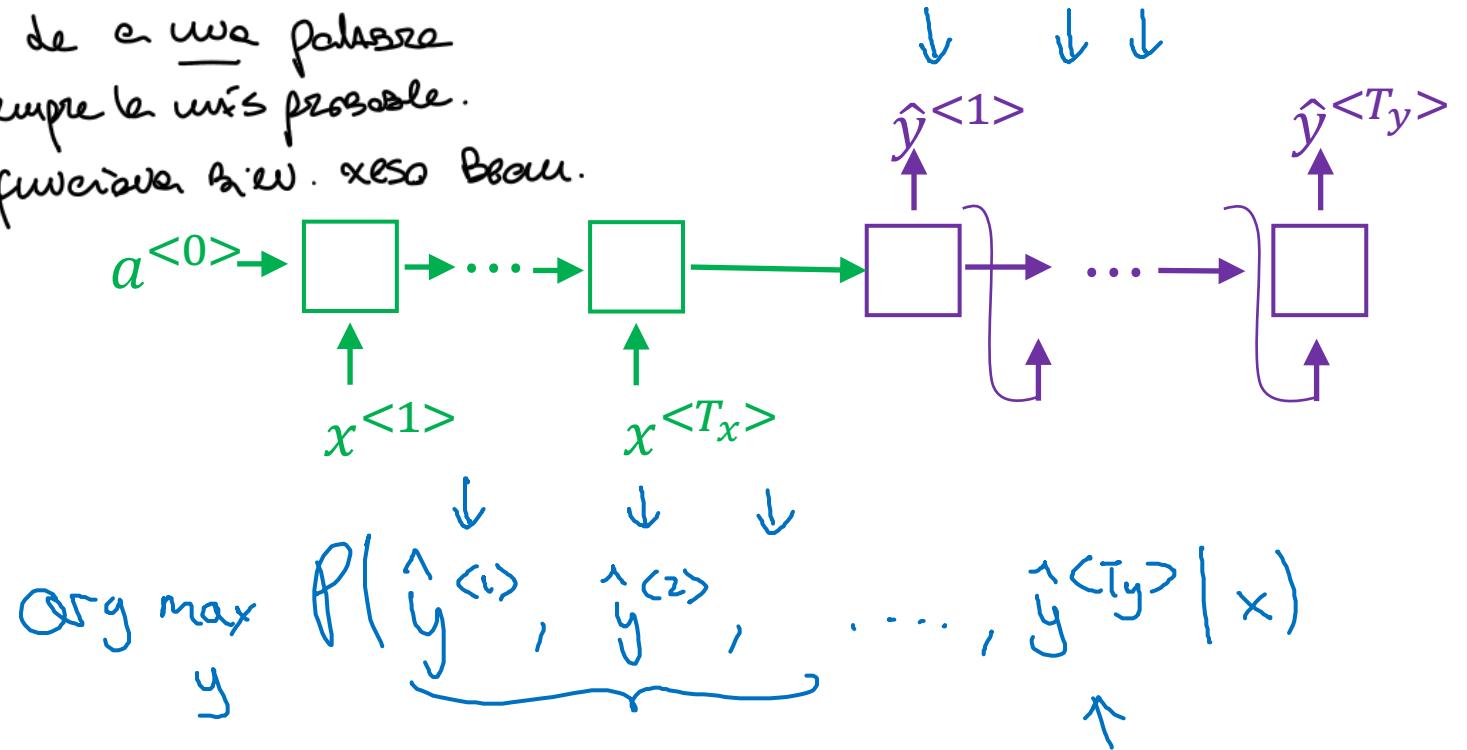
$$\arg \max_{y^{<1>} , \dots , y^{<T_y>}} P(y^{<1>} , \dots , y^{<T_y>} | x)$$

⇒ queremos maximizar la prob., buscamos ese  $y$ .  
Se suele hacer con Beam Search

# Why not a greedy search?

$$P(\hat{y}^{(1)} | x)$$

Estás eligiendo de entre una palabra.  
eligieras siempre la más probable.  
pero eso no funcionaría bien. como Beau.



$$\begin{aligned} & 10,000 \\ & 10 \\ & \underline{10,000^{10}} \\ & \underline{P(y|x)} \end{aligned}$$

→ Jane is visiting Africa in September.

→ Jane is going to be visiting Africa in September.

$$P(\text{Jane is going } | x) > P(\text{Jane is visiting } | x)$$



deeplearning.ai

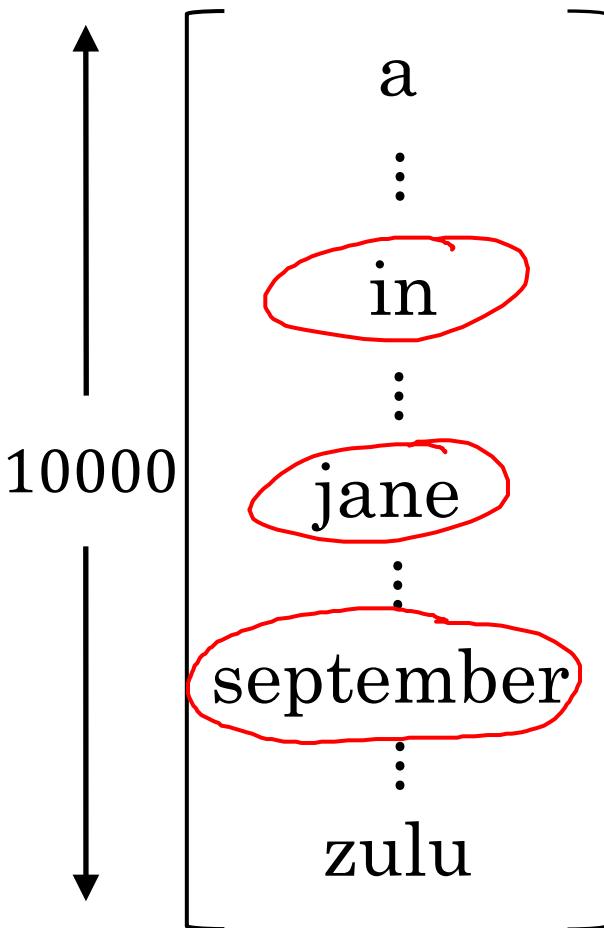
# Sequence to sequence models

---

Beam search

# Beam search algorithm

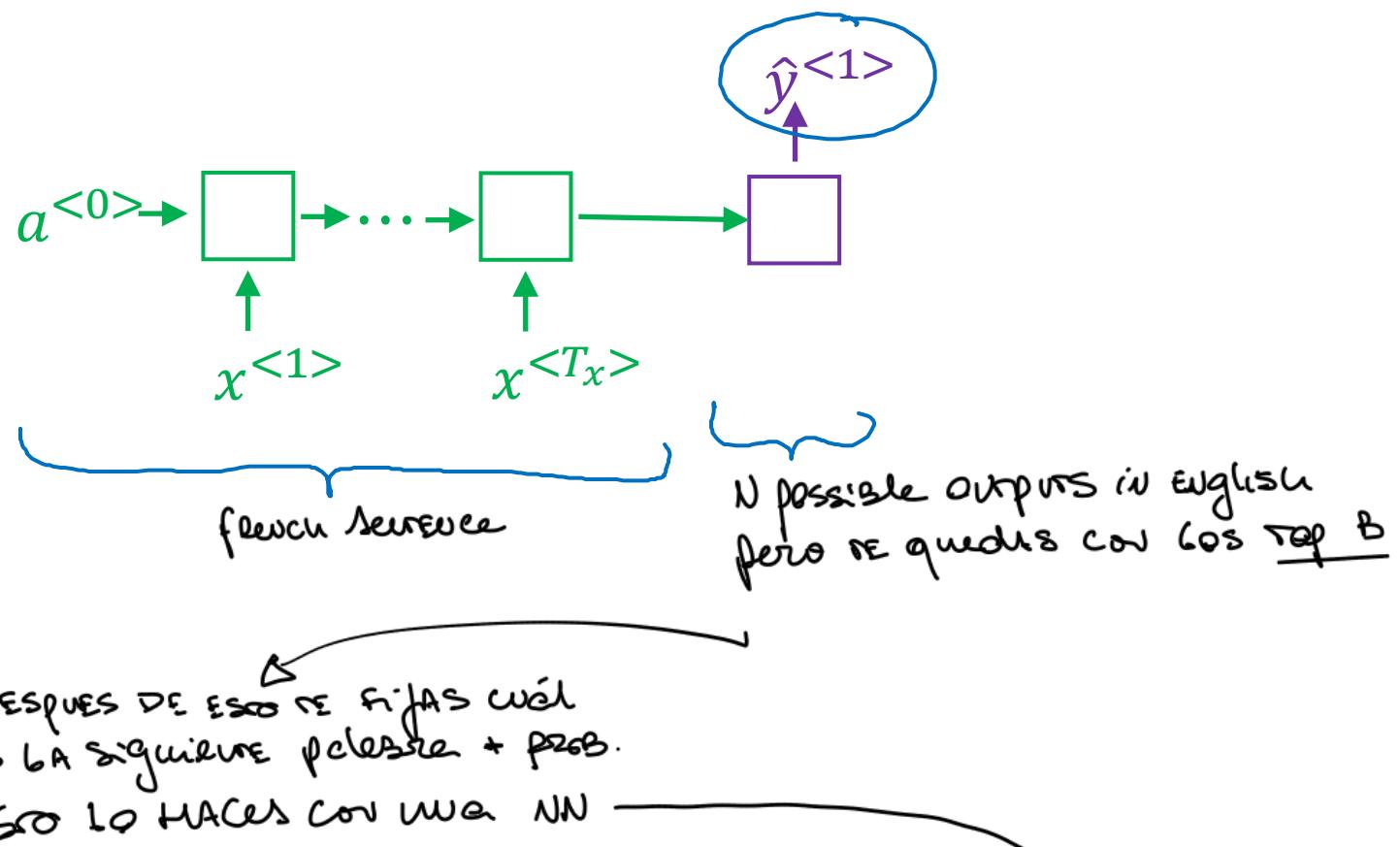
Step 1



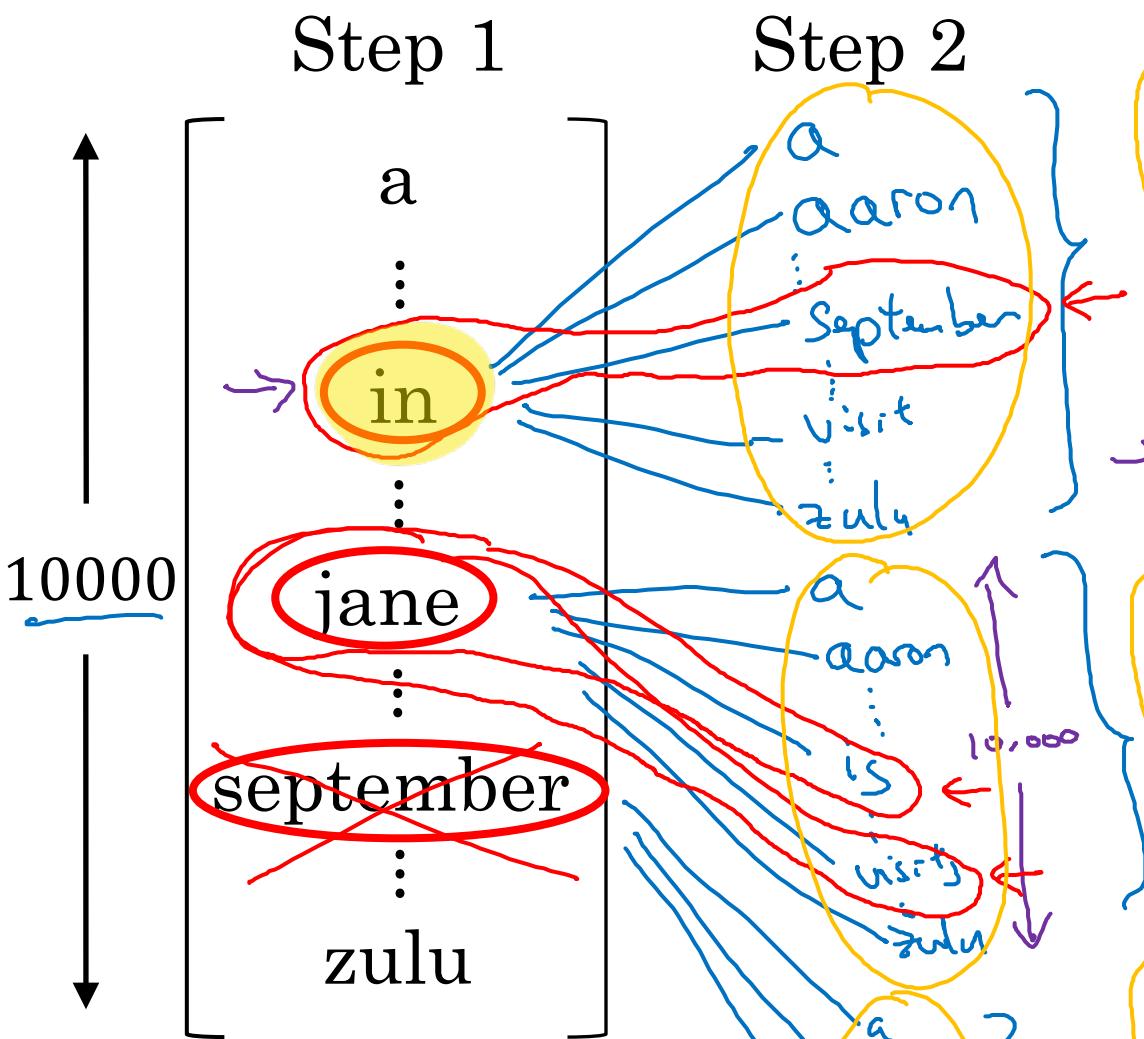
$$B = 3$$

Cantidad de palabras max  
A evaluar como most likely output.

$$\rightarrow P(y^{<1>} | x)$$

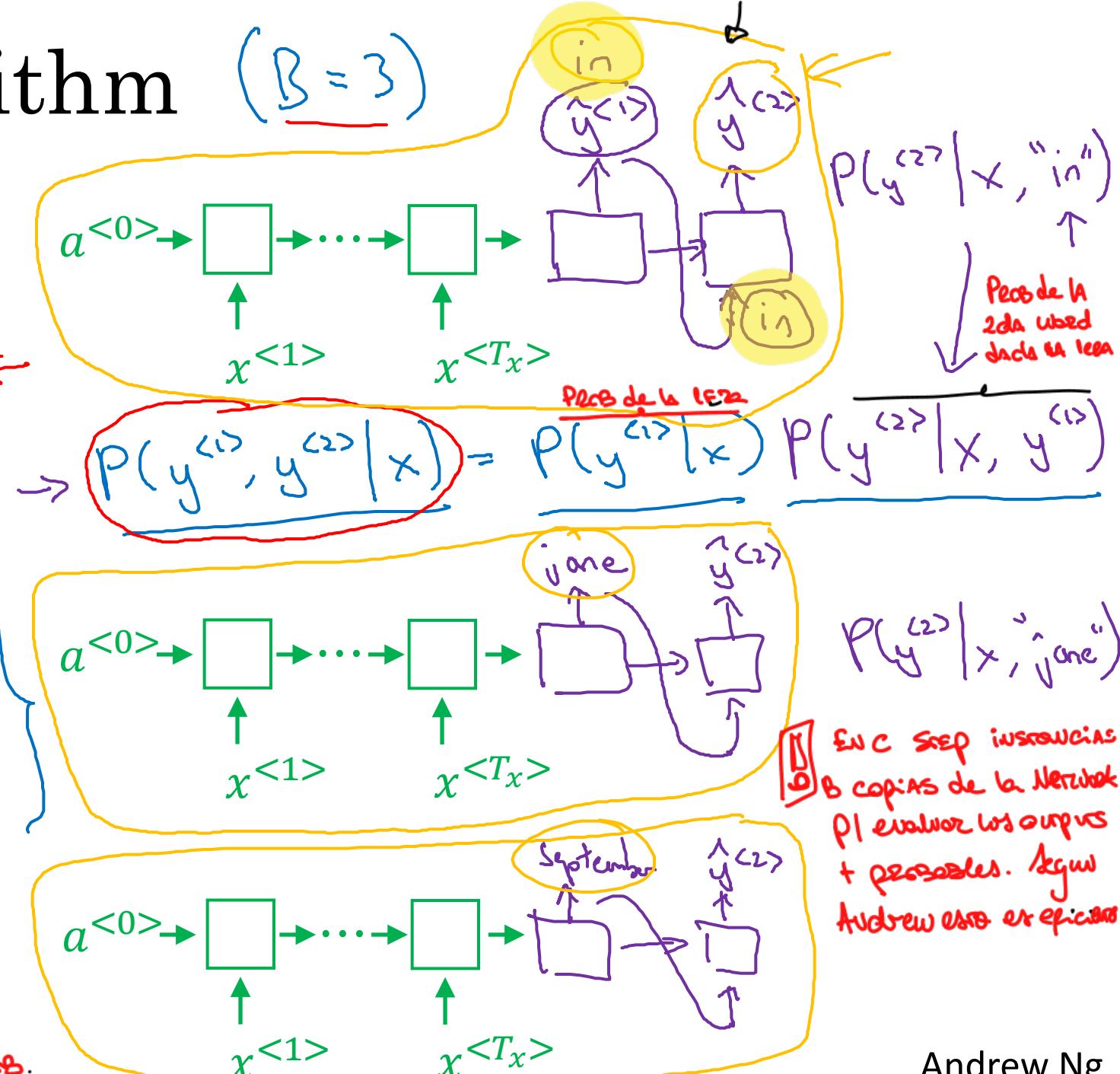


# Beam search algorithm



Cada vez q evalutas  $y^{(1)}, y^{(2)}$   
operarias, evalutas  
B.N Posibilidades  $y$  te quedas  $B + \text{Prob.}$

( $B = 3$ )



# Beam search ( $B = 3$ )

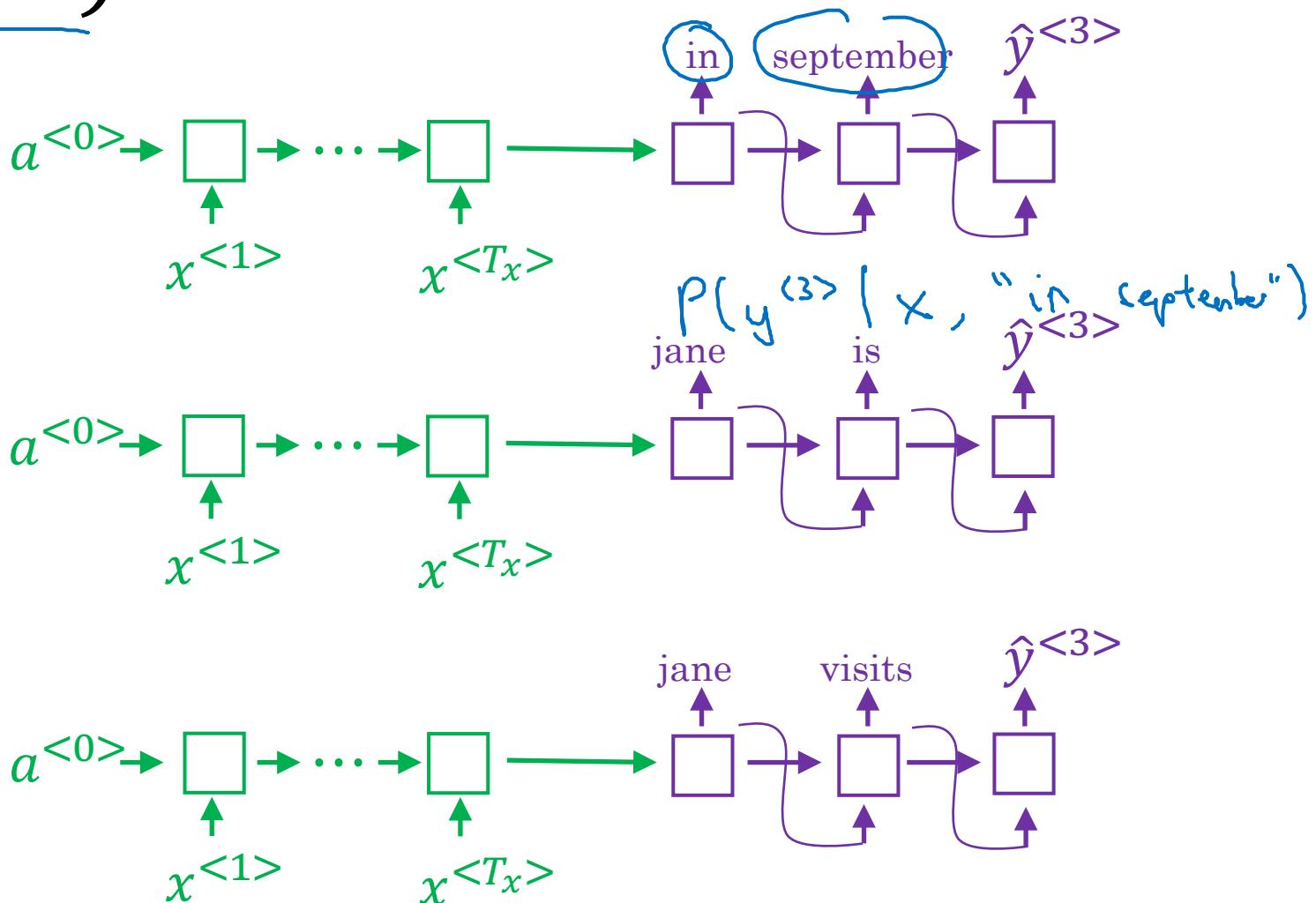
in september

jane is

jane visits

$$P(y^{<1>} , y^{<2>} | x)$$

$B=1 \rightsquigarrow$  greedy search



jane visits africa in september. <EOS>



deeplearning.ai

# Sequence to sequence models

---

## Refinements to beam search

# Length normalization

$$p(y^{(1)} \dots y^{(T_y)} | x) = P(y^{(1)} | x) P(y^{(2)} | x, y^{(1)}) \dots P(y^{(T_y)} | x, y^{(1)}, \dots, y^{(T_y-1)})$$

$$\arg \max_y \prod_{t=1}^{T_y} P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$

*Puede haber Underflow si los prob son muy chiquitos. Entonces tomar el log.*

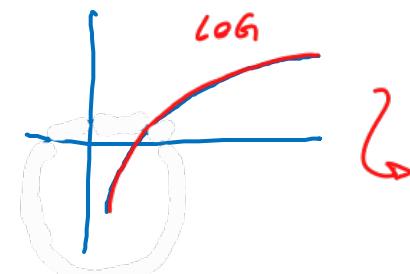
$$\arg \max_y \sum_{t=1}^{T_y} \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$

*→ ojo! Si tengo una oración muy larga, estoy multiplicando muchos números < 1 y eso me da un producto (prob) muy pequeño. Lo mismo es cierto con el log (de va haciendo + negativo). Para arreglarlo:*

$$T_y = 1, 2, 3, \dots, 30.$$

Elegis la sequence con el value + la grandeza de esa de todas las q vienes en beam search

$$\frac{1}{T_y^\alpha} \sum_{t=1}^{T_y} \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$



$$\begin{aligned} \log P(y|x) \\ P(y|x) \end{aligned}$$

→ ojo! Si tengo una oración muy larga, estoy multiplicando muchos números < 1 y eso me da un producto (prob) muy pequeño. Lo mismo es cierto con el log (de va haciendo + negativo). Para arreglarlo:

$$\alpha = 0.7$$

Límpiate para controlar esto →

$$\left\{ \begin{array}{l} \alpha = 1 \rightarrow \text{Normalizar} \\ \alpha = 0 \rightarrow \text{sin normalizar} \end{array} \right.$$

# Beam search discussion

Beam width B?

$1 \rightarrow 3 \rightarrow 10, \quad 100, \quad 1000 \rightarrow 3000$

{  
large B: better result, slower  
small B: worse result, faster

Unlike exact search algorithms like BFS (Breadth First Search) or DFS (Depth First Search), Beam Search runs faster but is not guaranteed to find exact maximum for  $\arg \max_y P(y|x)$ .

y



deeplearning.ai

# Sequence to sequence models

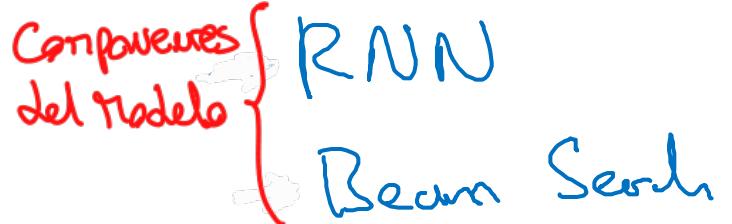
---

## Error analysis on beam search

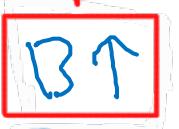
What if Beam is causing problems? how do you know its this?

# Example

Jane visite l'Afrique en septembre.

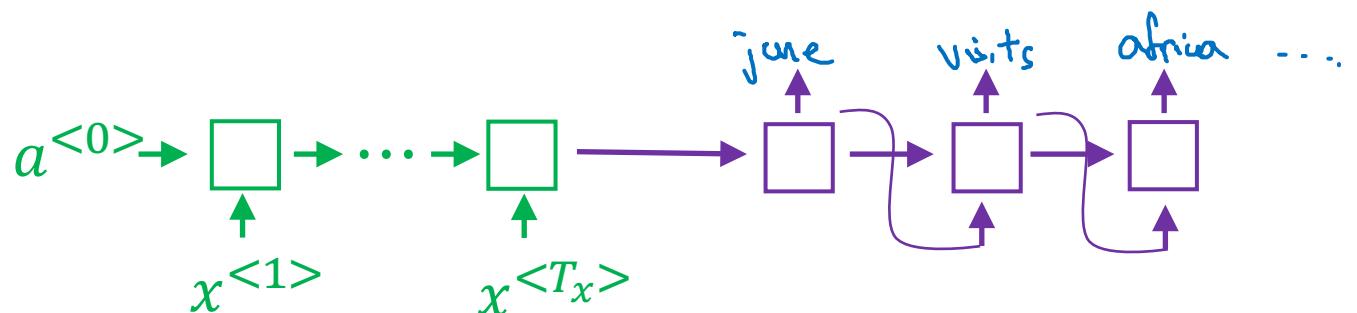


CASI Siempre  
funciona pero  
A veces no es suf.



- { Human: Jane visits Africa in September.  $(y^*)$
- Algorithm: Jane visited Africa last September.  $(\hat{y}) \leftarrow$

RNN computes:  $P(y^*|x) \leftarrow \text{RNN} \rightarrow P(\hat{y}|x)$



# Error analysis on beam search

$$P(y^*|x)$$

Human: Jane visits Africa in September. ( $y^*$ )

$$P(\hat{y}|x)$$

Algorithm: Jane visited Africa last September. ( $\hat{y}$ )

Case 1:  $P(y^*|x) > P(\hat{y}|x) \leftarrow$

$$\arg \max_y P(y|x)$$

Beam search chose  $\hat{y}$ . But  $y^*$  attains higher  $P(y|x)$ .

Conclusion: Beam search is at fault.

Case 2:  $\underline{P(y^*|x)} \leq \underline{P(\hat{y}|x)} \leftarrow$

$y^*$  is a better translation than  $\hat{y}$ . But RNN predicted  $P(y^*|x) < P(\hat{y}|x)$ .

Conclusion: RNN model is at fault.

# Error analysis process

| Human                            | Algorithm                           | <small>Human</small><br>$P(y^* x)$    | <small>Algorithm</small><br>$P(\hat{y} x)$ | At fault? |
|----------------------------------|-------------------------------------|---------------------------------------|--|-----------|
| Jane visits Africa in September. | Jane visited Africa last September. | <u><math>2 \times 10^{-10}</math></u> | <u><math>1 \times 10^{-10}</math></u>      | B         |
| ...                              | ...                                 | —                                     | —  | R         |
| ...                              | ...                                 | —                                     | —  | R         |
|                                  |                                     |                                       | —  | R         |
|                                  |                                     |                                       |  | :         |

Figures out what fraction of errors are “due to” beam search vs. RNN model



deeplearning.ai

# Sequence to sequence models

---

Bleu score  
(optional)

# Evaluating machine translation

French: Le chat est sur le tapis.

Reference 1: The cat is on the mat. ←

Reference 2: There is a cat on the mat. ←

MT output: the the the the the the. (?)

Precision:  $\frac{7}{7}$

de las 7 palabras, cuantas  
Aparecen en el Test?

No nos  
sirve  
→

Modified precision:

$$\frac{2}{7}$$

Max times q APARECE  
en un Ejemplo  
Count de la palabra que

Bilingual Evaluation Understudy

→ Inversion: Nos fijemos si las palabras  
que en el mt aparecen en el dexter  
correcto.

BLEU: Lo usamos

para medir Accuracy  
cuando hay más de  
una Rpta Correcta

# Bleu score on bigrams

→ calculate pls see us please

Example: Reference 1: The cat is on the mat.

Reference 2: There is a cat on the mat.

MT output: The cat the cat on the mat.

Count eu nr

the cat

2

cat the

1

cat on

1

on the

1

the mat

1

Count clip

1

1

1

The n° of times the  
Bigram Appears in At least  
one REF.

|   |                       |
|---|-----------------------|
| 4 | → $\Sigma$ count clip |
| 6 | → $\Sigma$ count nr   |

modified precision  
of Bigrams.

# Bleu score on unigrams

Example: Reference 1: The cat is on the mat.

$$P_1, P_2 = 1.0$$

Reference 2: There is a cat on the mat.

→ MT output: The cat the cat on the mat. (↑)

UNIGRAMS

$$P_1 = \frac{\sum_{\text{Unigrams } e^1} \text{Count clip (unigram)}}{\sum_{\text{Unigram } e^1} \text{Count (unigram)}}$$

N GRAMS

$$P_N = \frac{\sum_{N\text{grams } e^1} \text{Count clip (ngram)}}{\sum_{N\text{grams } e^1} \text{Count (ngram)}}$$

# Bleu details

$p_n$  = Bleu score on n-grams only

Combined Bleu score:

$$\text{Bleu} = \text{BP} \cdot \exp \left( \frac{1}{4} \sum_{n=1}^4 p_n \right)$$

BP = Brevity penalty  $\rightarrow$  extra traducciones muy cortas

$$\text{BP} = \begin{cases} 1 & \text{if MT\_output\_length} > \text{reference\_output\_length} \\ \exp(1 - \text{MT\_output\_length}/\text{reference\_output\_length}) & \text{otherwise} \end{cases}$$

Basicamente el Bleu es bajo si cuando hay multiples respuestas correctas.





deeplearning.ai

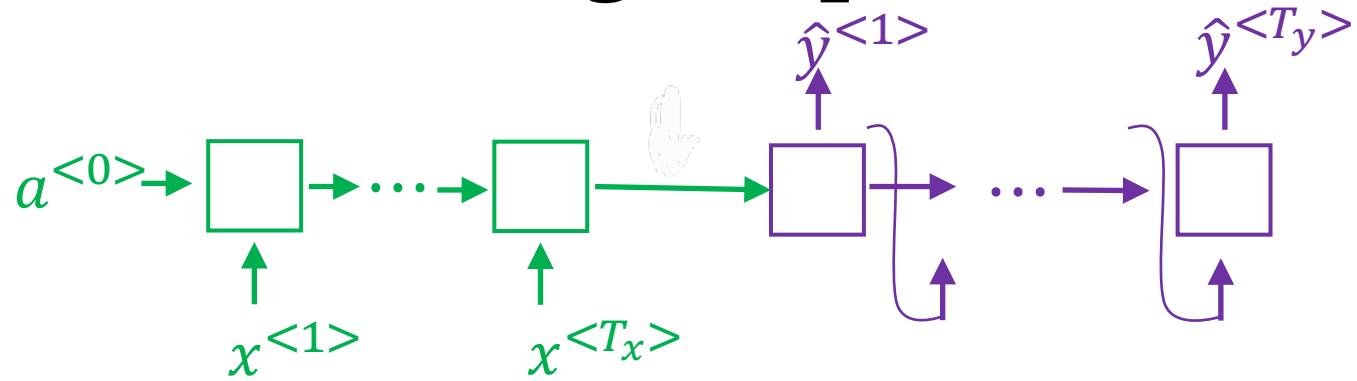
# Sequence to sequence models

---

Attention model  
intuition

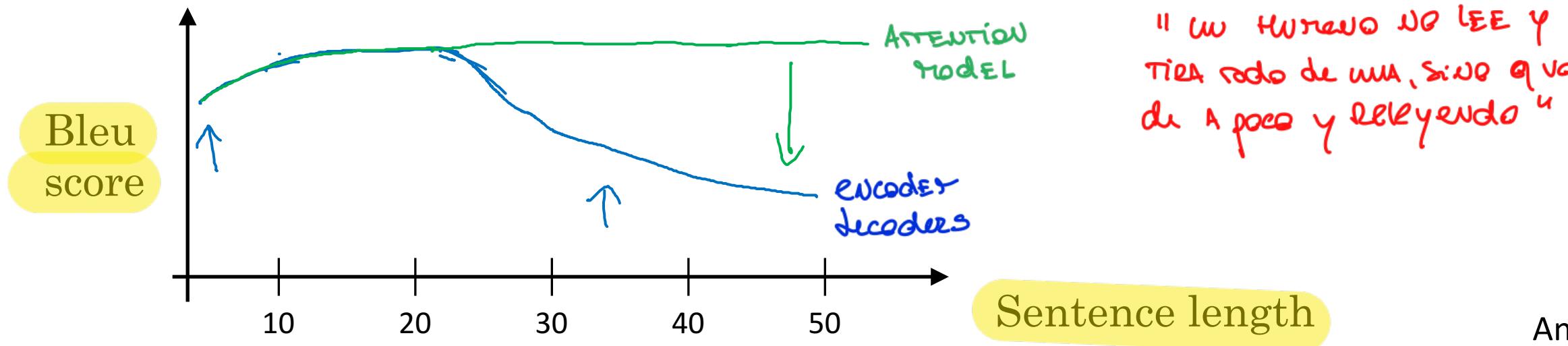
Road to Attention is all you need ☺

# The problem of long sequences

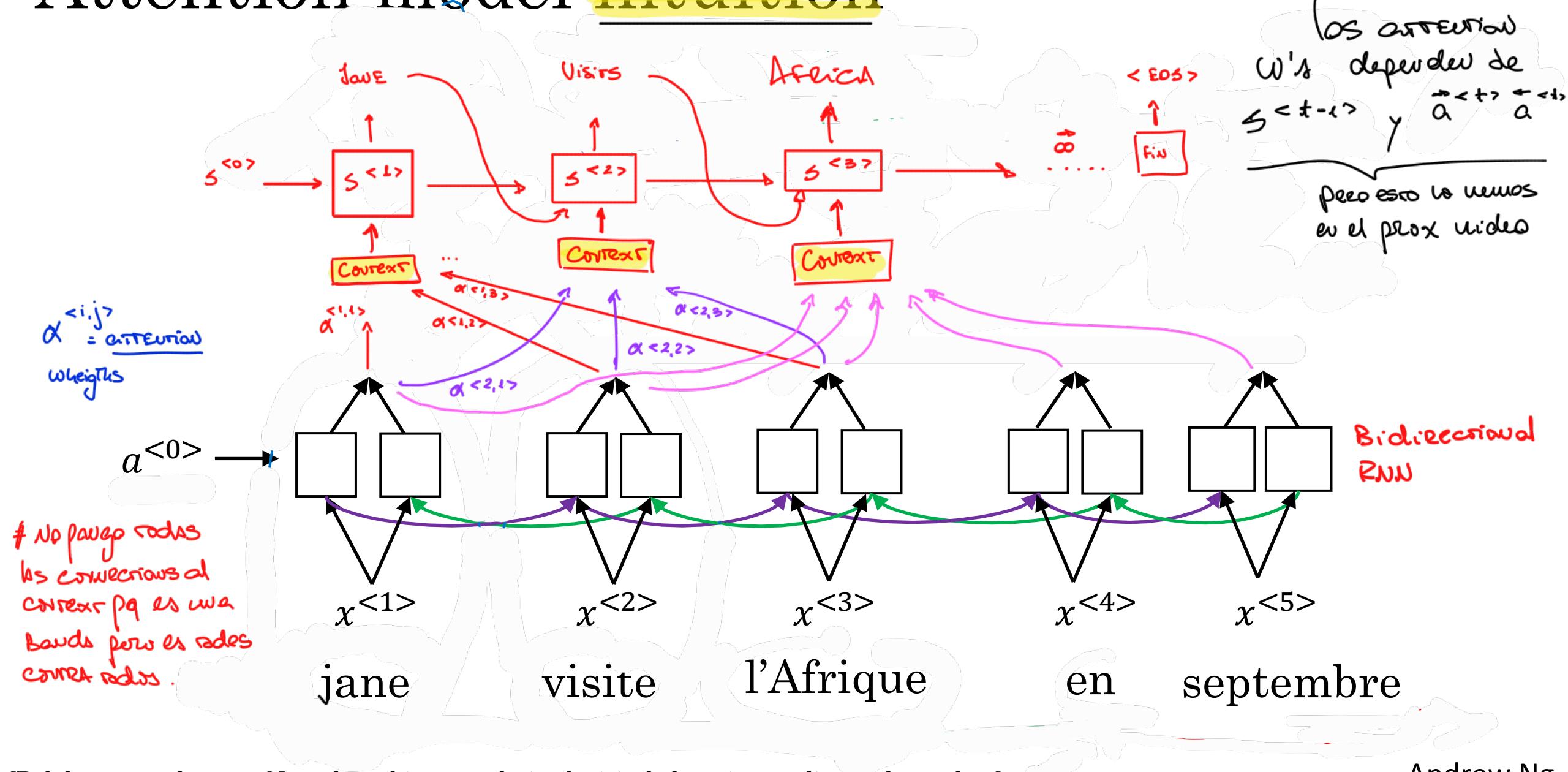


Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente d'y aller aussi.

Jane went to Africa last September, and enjoyed the culture and met many wonderful people; she came back raving about how wonderful her trip was, and is tempting me to go too.



# Attention model intuition





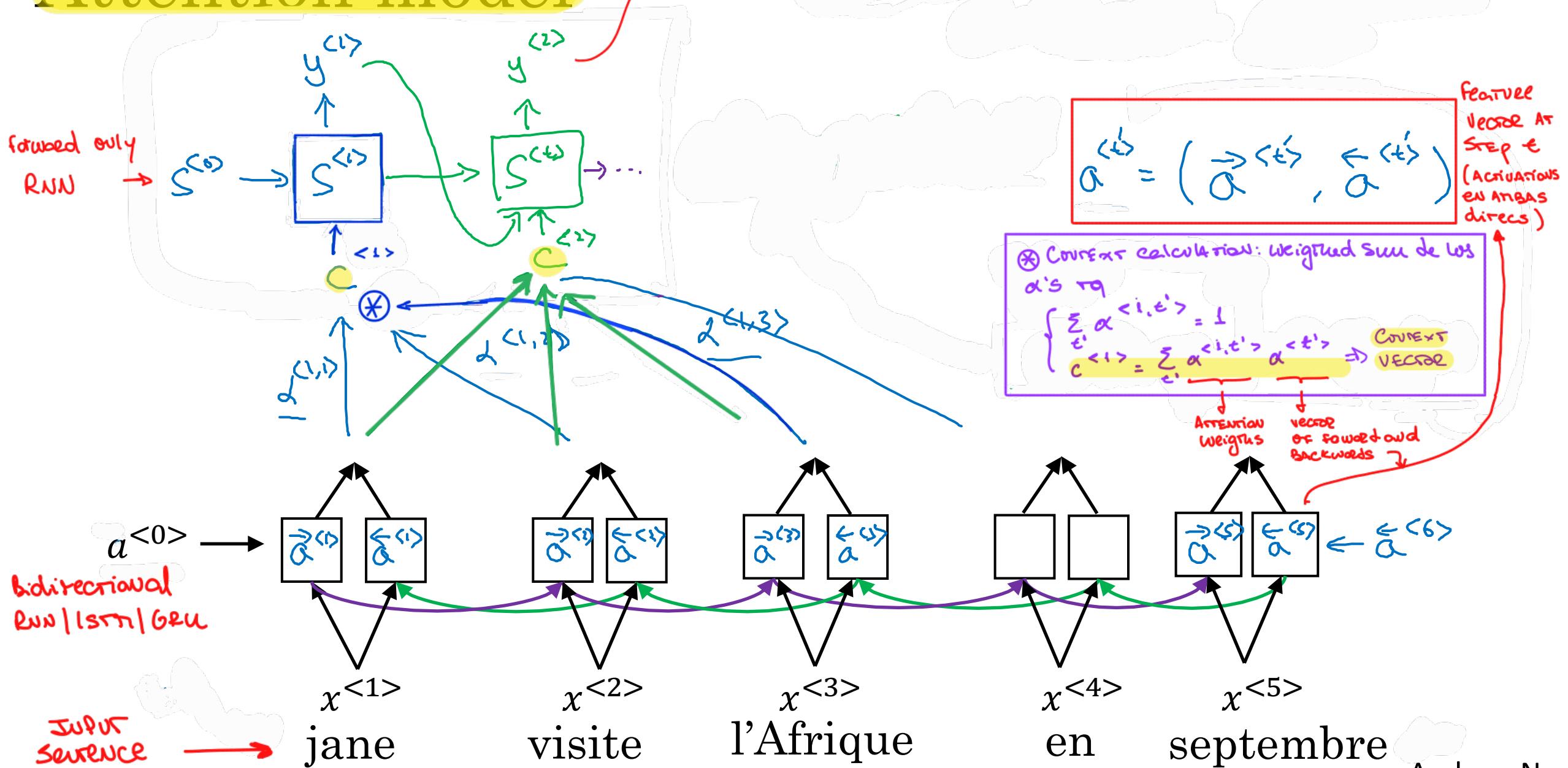
deeplearning.ai

# Sequence to sequence models

---

Attention model

# Attention model



# Computing attention $\alpha^{<t,t'>}$

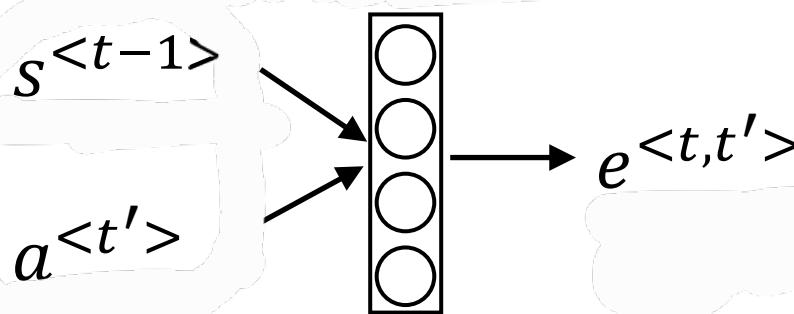
$\alpha^{<t,t'>} = \text{amount of attention } y^{<t>} \text{ should pay to } a^{<t'>}$

alfa! de attention

$$\alpha^{<t,t'>} = \frac{\exp(e^{<t,t'>})}{\sum_{t'=1}^{T_x} \exp(e^{<t,t'>})}$$

suman 1 plc t'

los e se  
calculan usando  
una mini  
red neuronal.

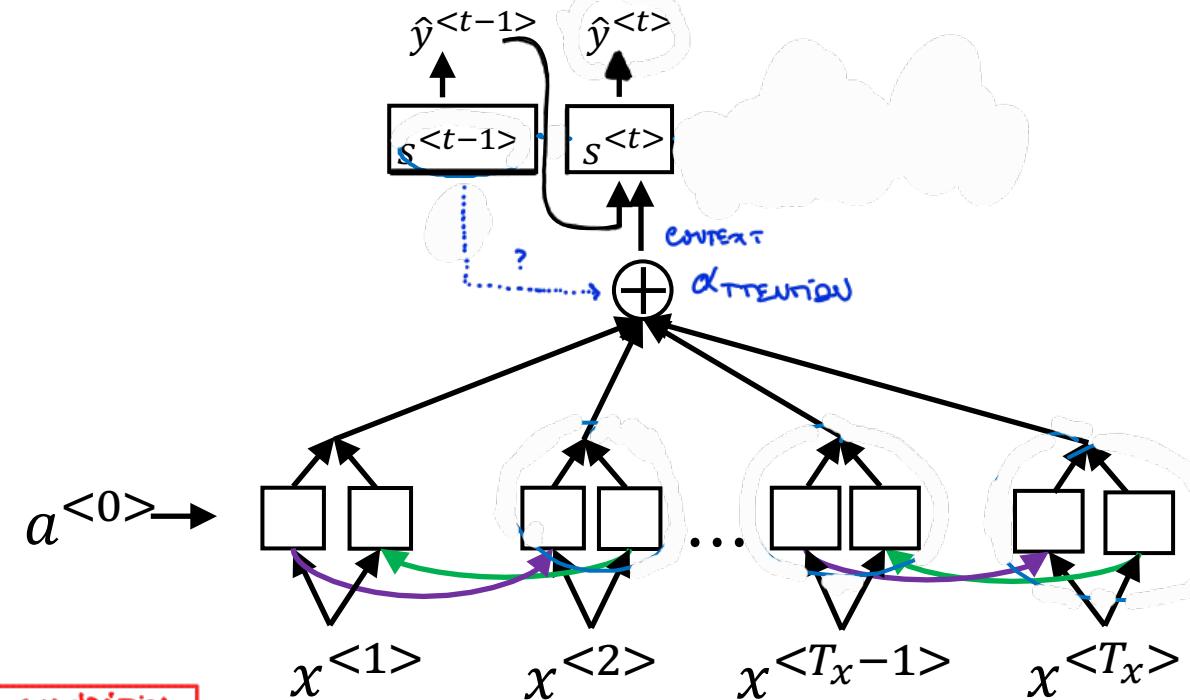


$s$  = la activación del paso anterior fijo de la  
Red de Atención. NO ENTENDÍ MUY BIEN ESTO,  
Si contexto necesita los  $a$  y los  $a$  necesitan  
los  $e$  q se calculan con los  $s$ , NO ESNA mal el dibujo?

Escala de forma cuadrática  
este Algoritmo. Hay Research  
para revertir de Bajarlo igual

Sí. Aquí esca ok

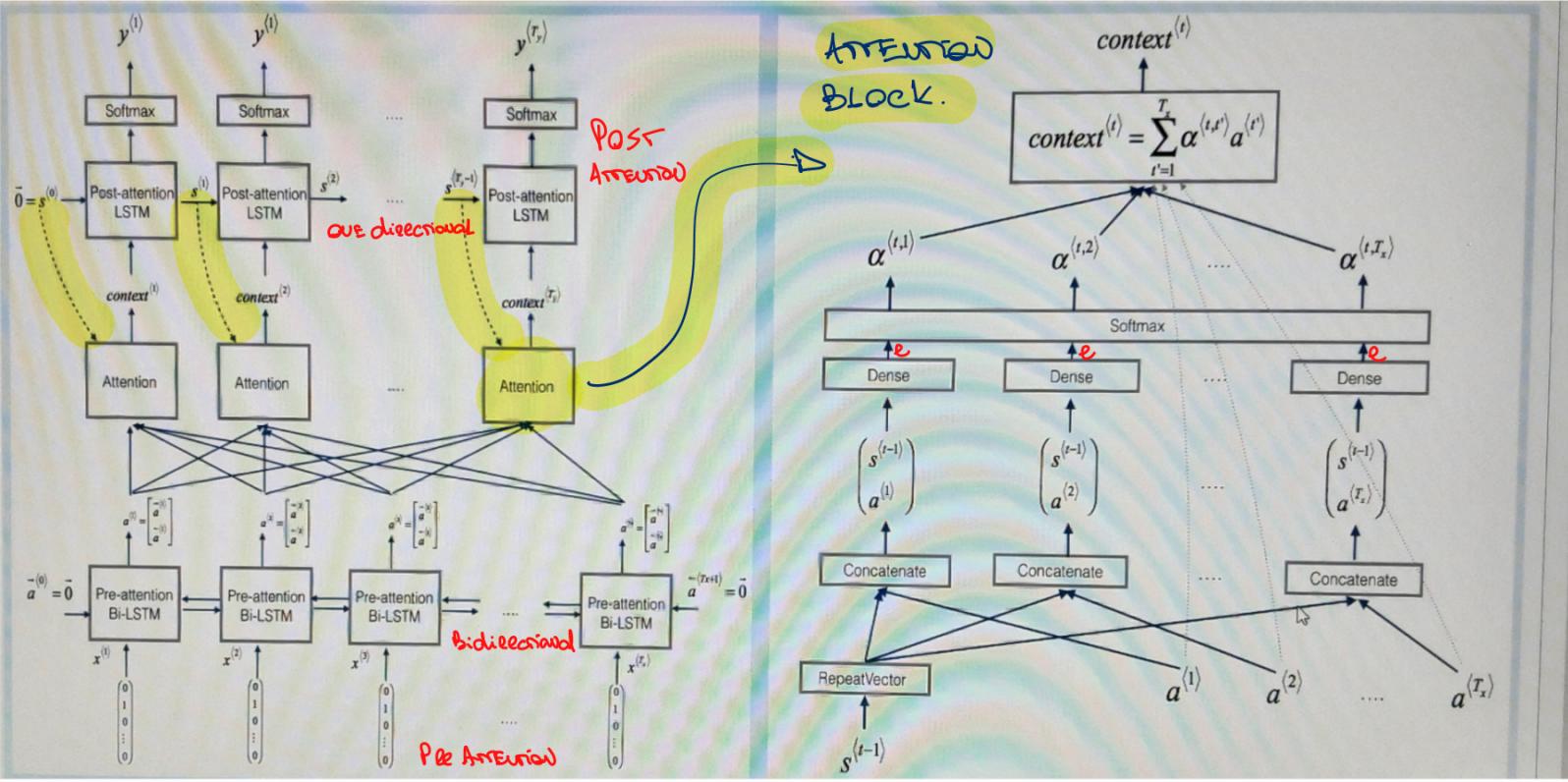
"a" de hidden state



[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

[Xu et. al., 2015. Show, attend and tell: Neural image caption generation with visual attention]

Andrew Ng



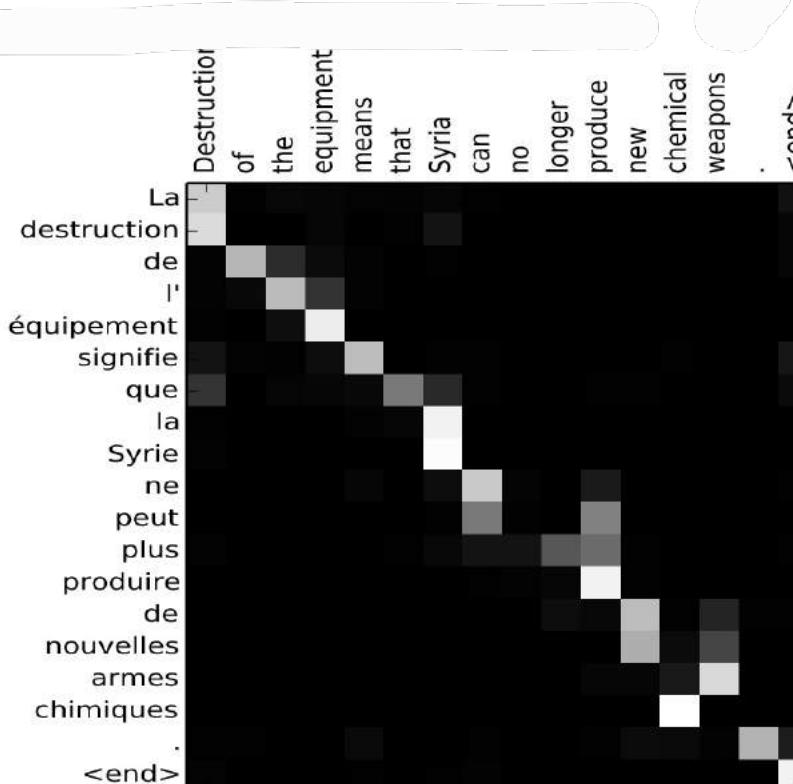
#  $e = \text{energies}$ . tended algo q vez con los EBN's de Yann LeCun?

# Attention examples

July 20th 1969 → 1969 – 07 – 20

23 April, 1564 → 1564 – 04 – 23

Visualization of  $\alpha^{<t,t'>}$ :



→ Attention Weights.



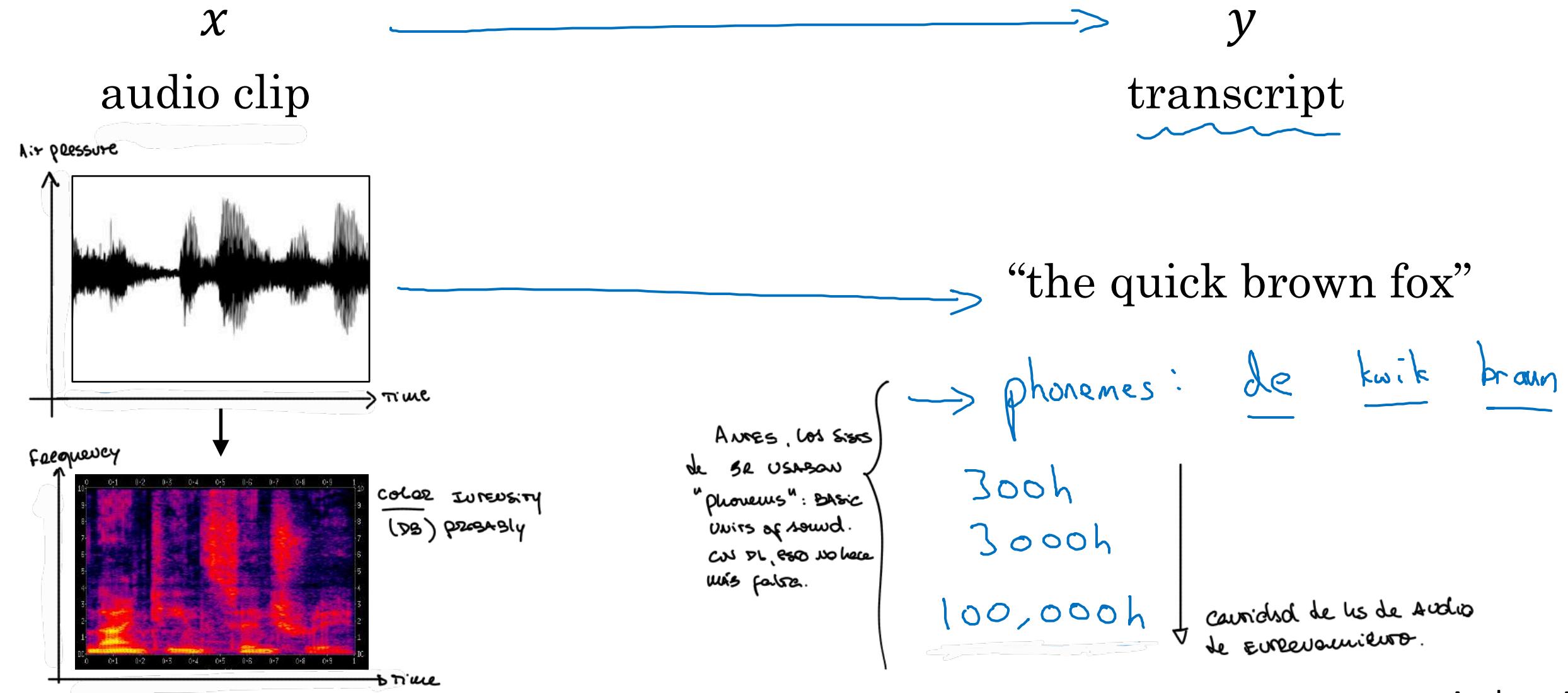
deeplearning.ai

Audio data

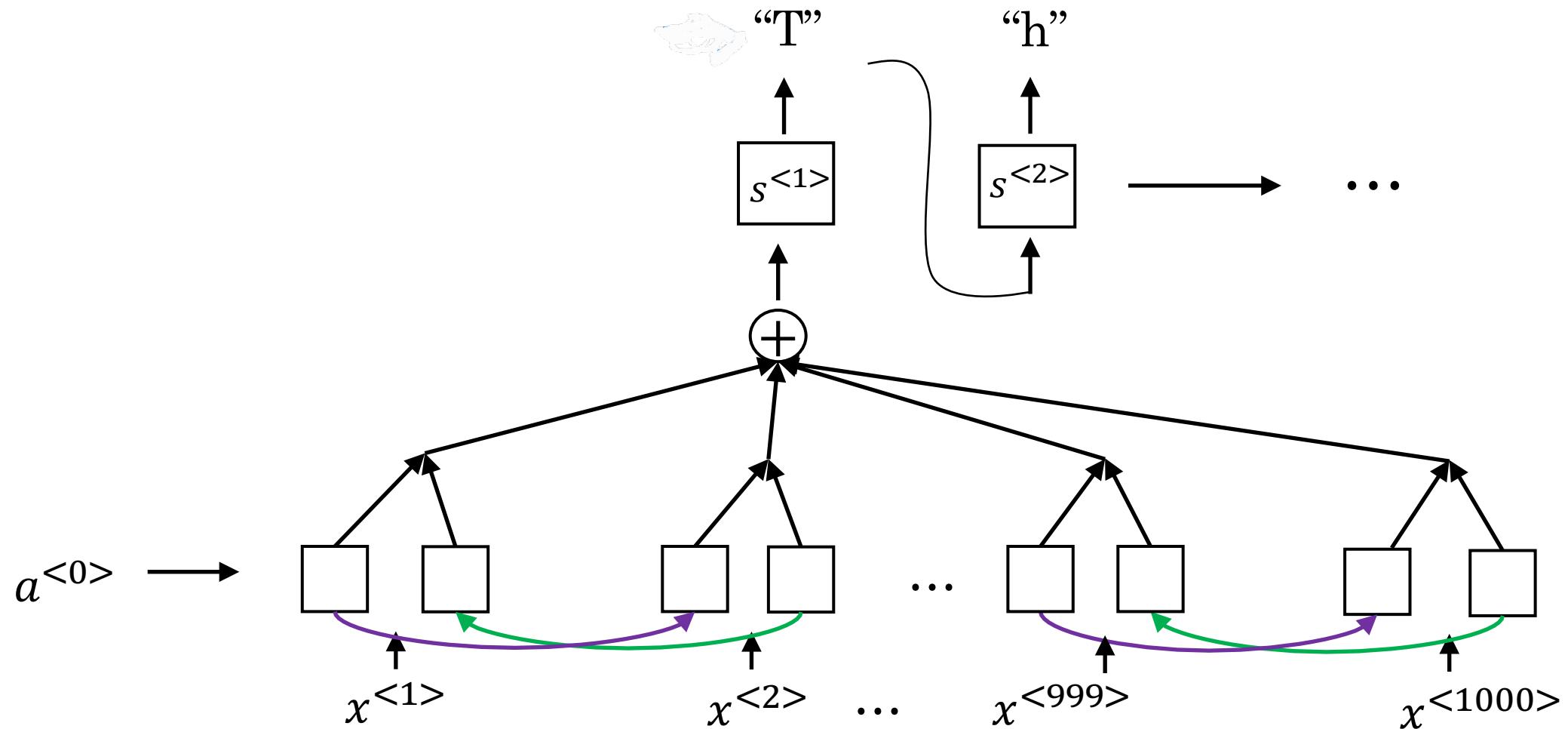
---

Speech recognition

# Speech recognition problem



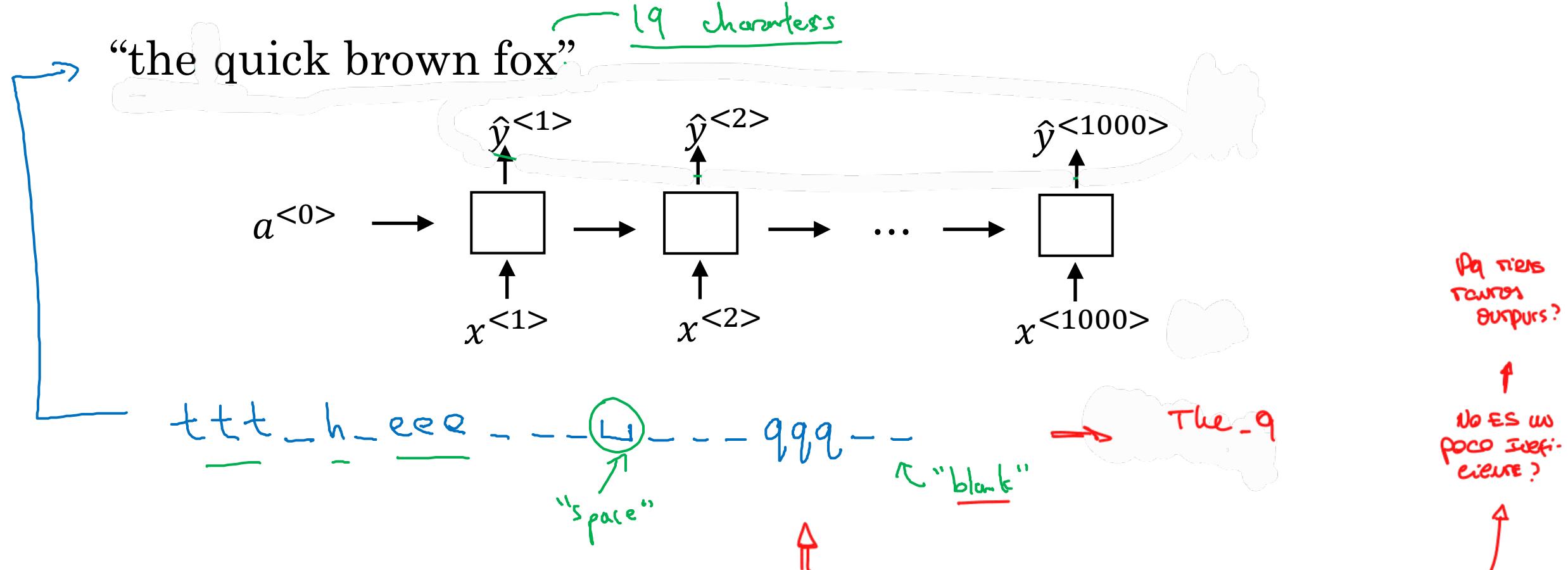
# Attention model for speech recognition



# CTC cost for speech recognition

(Connectionist temporal classification)

# Encontrar un modelo de SR hoy por hoy es muy costoso y difícil, necesitamos muchos datos. Un trigger word es + fácil y lo vemos en el big video



Basic rule: collapse repeated characters not separated by "blank"



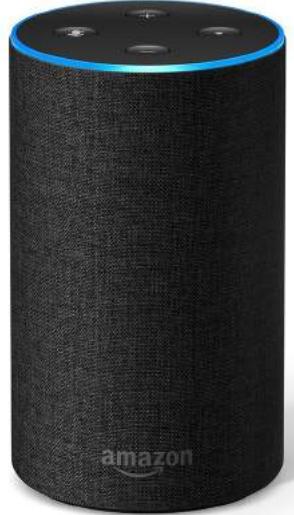
deeplearning.ai

Audio data

---

Trigger word  
detection

# What is trigger word detection?



Amazon Echo  
(Alexa)



Baidu DuerOS  
(xiaodunihao)

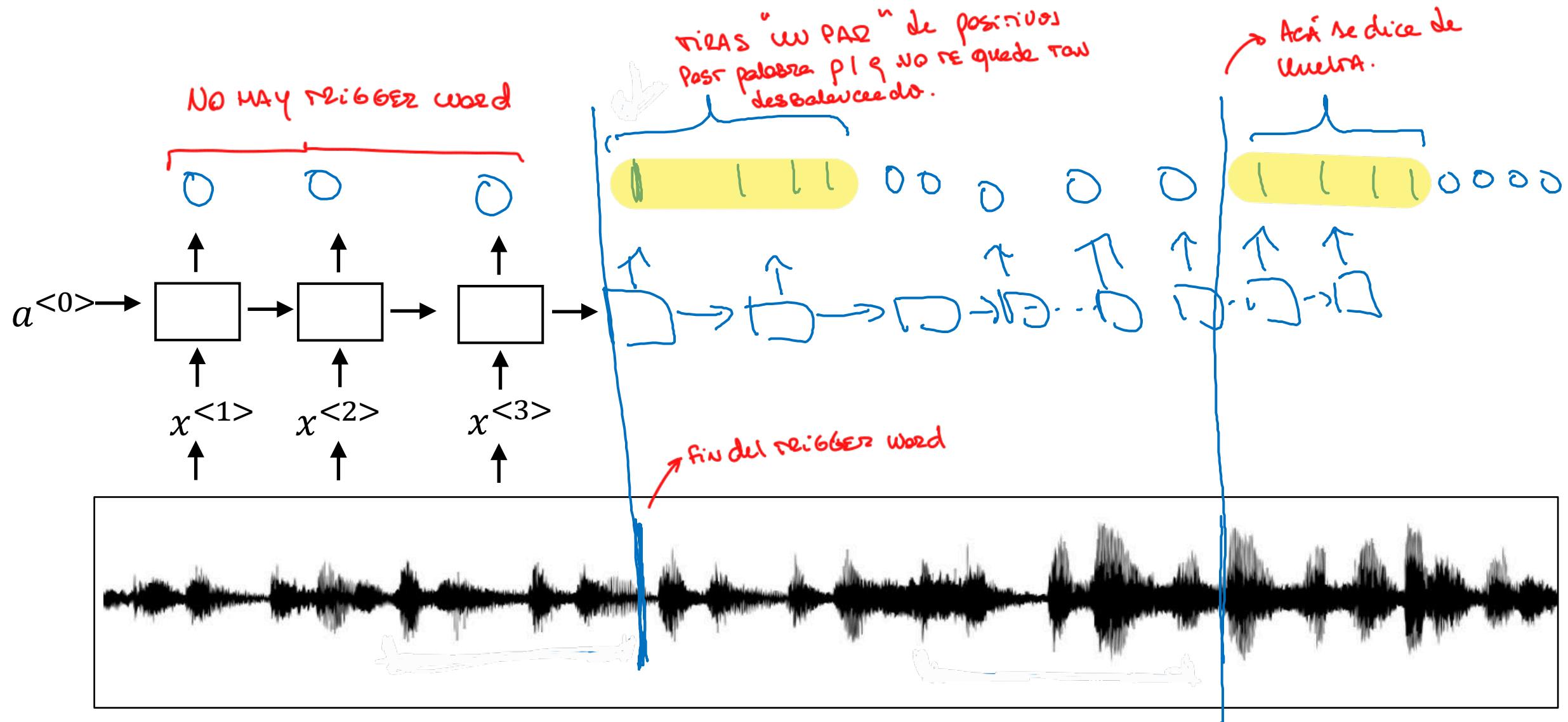


Apple Siri  
(Hey Siri)



Google Home  
(Okay Google)

# Trigger word detection algorithm





deeplearning.ai

# Conclusion

---

## Summary and thank you

# Specialization outline

1. Neural Networks and Deep Learning
2. Improving Deep Neural Networks: Hyperparameter tuning, Regularization and Optimization
3. Structuring Machine Learning Projects
4. Convolutional Neural Networks
5. Sequence Models

# Deep learning is a super power

Please buy this  
from shutterstock  
and replace in  
final video.



[www.shutterstock.com](http://www.shutterstock.com) • 331201091

Thank you.

- Andrew Ng

# Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



deeplearning.ai

# Sequence to sequence models

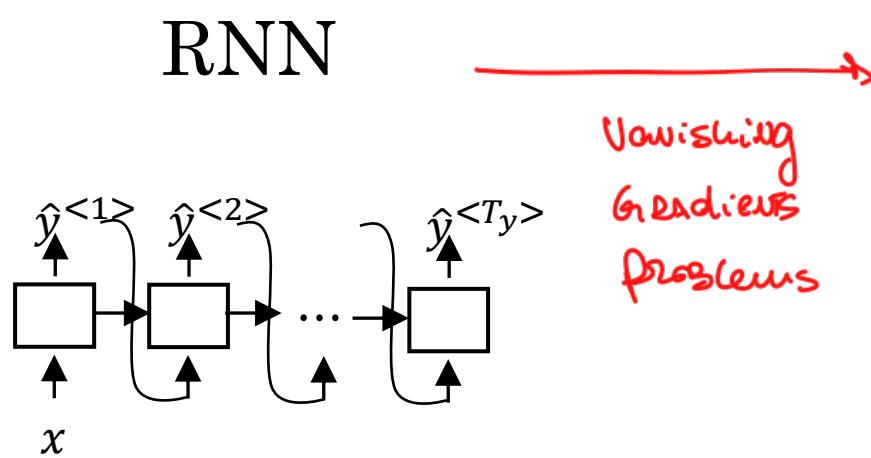
---

Transformers

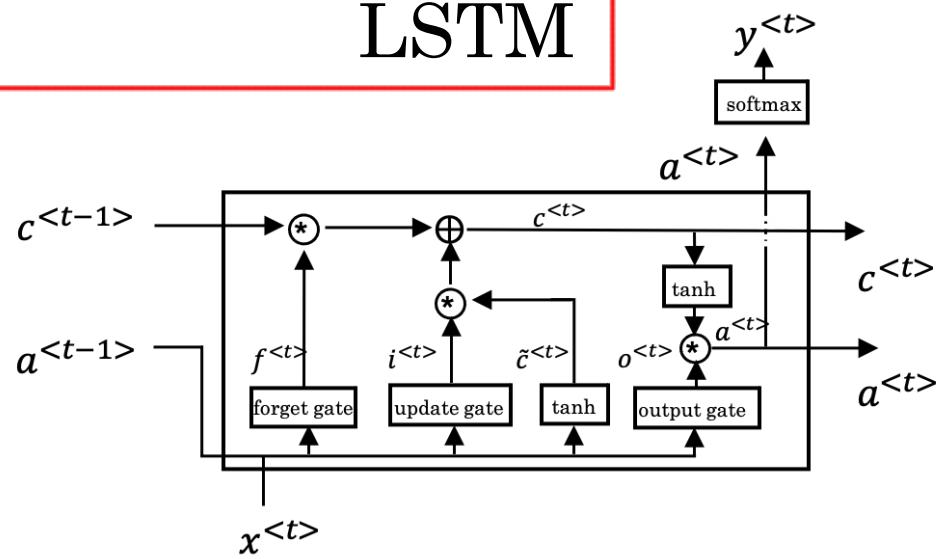
Intuition

# Transformers Motivation

Increased complexity,  
sequential



GRU      LSTM

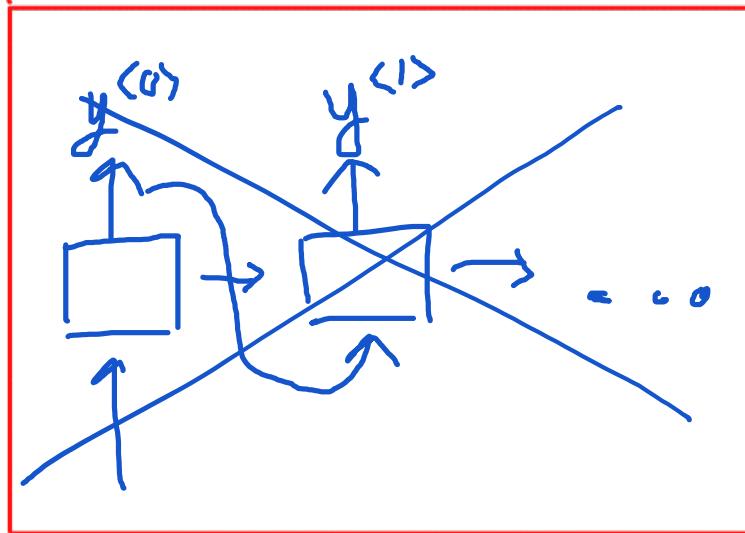


Mucho más complejos y todos sequential  
procesan un token at a time!

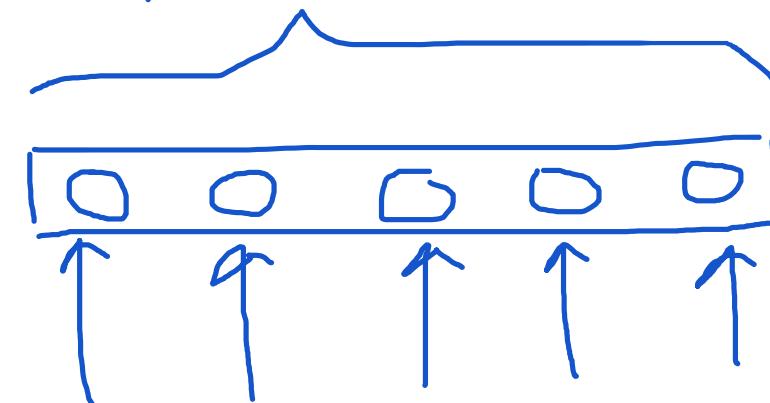
# Transformers Intuition

- Attention + CNN → CNN style
  - Self-Attention →  $A^{<1>} \dots A^{<N>}$  ⇒  $N$  representations in a sequence of  $N$  words  
q se puede Parallelizar
  - Multi-Head Attention →  $\frac{1}{N}$  for loop to see el proceso de self attention. → Muchas Rep de Self Attention  $\neq$ . ( $\frac{1}{N}$  copies)

Sequential way of processing is Recurrent



In parallel style of CNN's





deeplearning.ai

# Sequence to sequence models

---

## Self-Attention

# Self-Attention Intuition

→ Key, value, se explican luego

$A(q, K, V)$  = attention-based vector representation of a word

calcular para cada palabra

$A^{<1>} ; \dots ; A^{<N>}$

## RNN Attention

$$\alpha^{<t,t'>} = \frac{\exp(e^{<t,t'>})}{\sum_{t'=1}^{T_x} \exp(e^{<t,t'>})}$$

## Transformers Attention

$$A(q, K, V) = \sum_i \frac{\exp(e^{<q \cdot k^{<i>}>})}{\sum_j \exp(e^{<q \cdot k^{<j>}>})} v^{<i>}$$

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$   
Jane      visite      l'Afrique      en      septembre

↳ Para cada palabra {  
 $q$  = query  
 $K$  = key  
 $V$  = value}

# Self-Attention

AKA: scaled dot product attention.

Vectorized Rep

$$A(q, K, V) = \sum_i \frac{\exp(e^{q \cdot k^{<i>}})}{\sum_j \exp(e^{q \cdot k^{<j>}})} v^{<i>}$$

Vectorized Rep  $\rightarrow$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

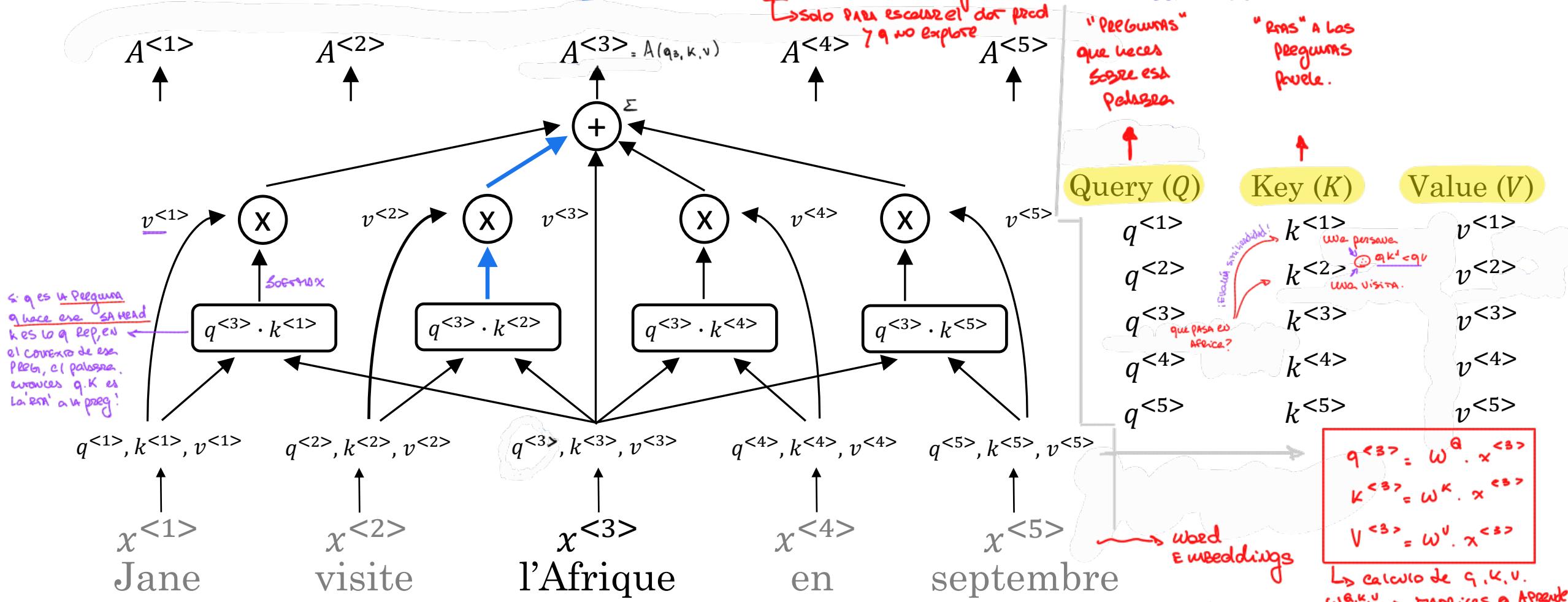
*Sí: ponemos todas las comparaciones de las palabras juntas*

*Solo para escuchar el otro lado q no explore*

ESO ES INTUITIVO MUCHO MÁS.

"Risas" a las preguntas

"Preguntas" que hacen sobre esa palabra



$$\begin{aligned} q^{<3>} &= W^Q \cdot x^{<3>} \\ k^{<3>} &= W^K \cdot x^{<3>} \\ v^{<3>} &= W^V \cdot x^{<3>} \end{aligned}$$

*L> calculo de q, k, v.  
W<sup>Q, K, V</sup> → matrices q aprende el Algoritmo.*



deeplearning.ai

# Sequence to sequence models

---

Multi-Head  
Attention

# Multi-Head Attention

BASICALLY, "A BIG FOR LOOP" SOBRE EL SLIDE ANTERIOR  
 ↴ PERO EN REALIDAD  
 MAY ALGORITMOS? ← ESO ES LO PARALELIZABLE !!

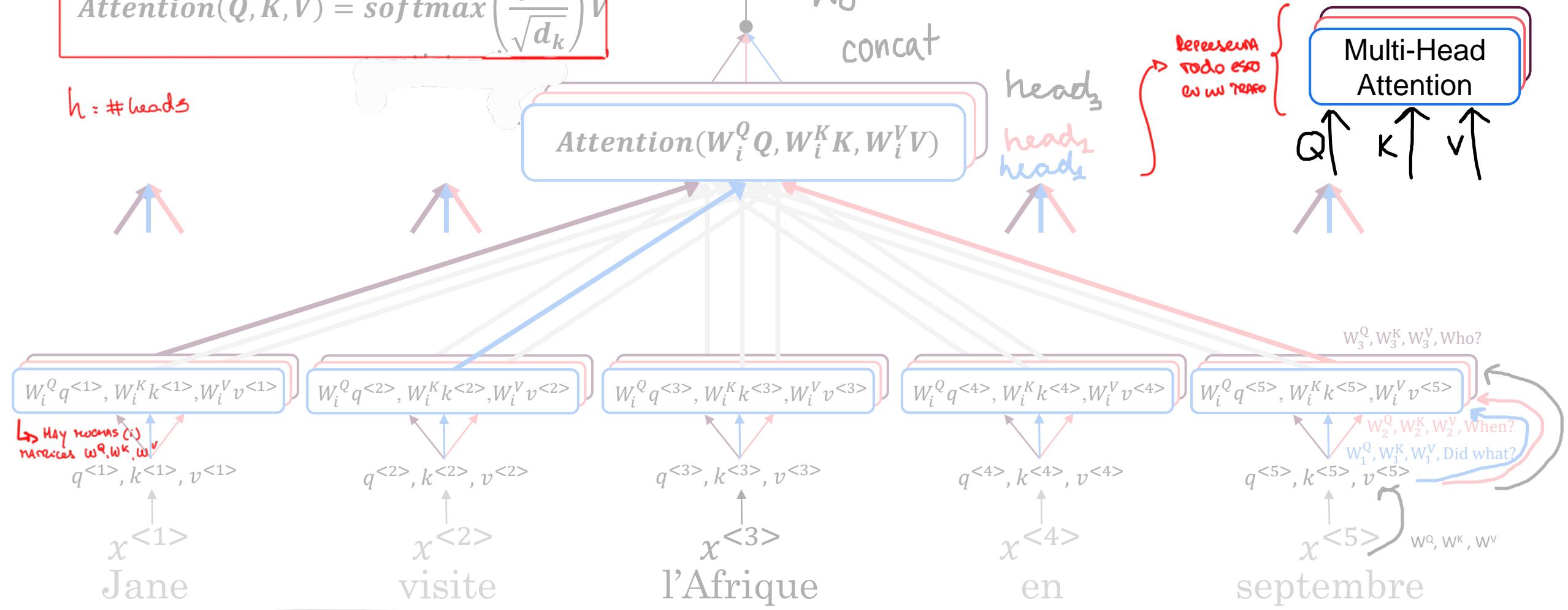
NU A UNO LE  
 PASO TODO.

$$\text{MultiHead}(Q, K, V) = \text{concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_n)W_o$$

$$\text{head}_i = \text{Attention}(W_i^Q Q, W_i^K K, W_i^V V)$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$h = \# \text{heads}$



•  $W_1^Q, W_1^K, W_1^V \rightarrow$  responder la pregunta "que es el pasado? p/c word/Reej".  $W_2^{Q,V}$  podrían preguntar "cuando es el pasado?". etc. Se aprenden esas "preguntas" → como si fueran Features!  
 q le vamos a dar a una NN ↴



deeplearning.ai

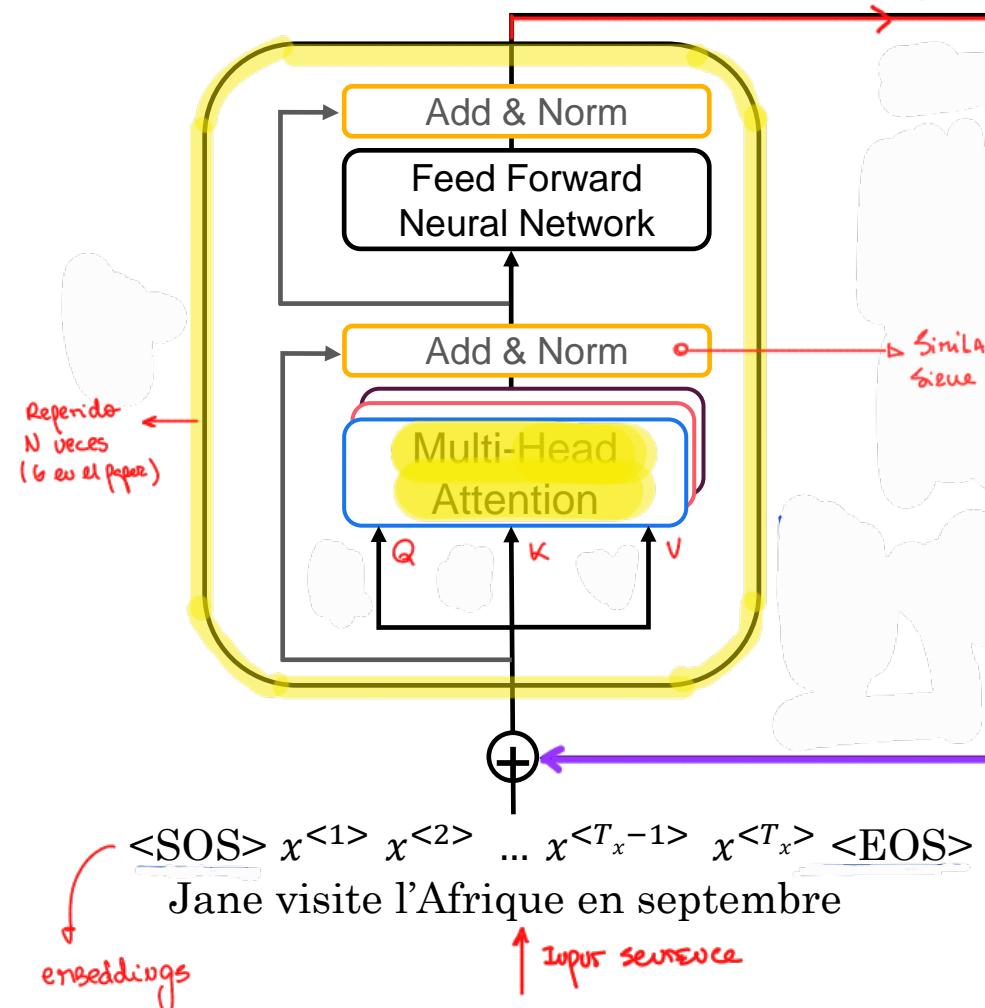
# Sequence to sequence models

---

Transformers

# Transformer Details

## Encoder



Contextual Information.

Similar to batch norm (?)  
Sigue para speedup training!

Aca puede ir: "Masked" Multi-Head Attention  
Importante para el training, cuando  
ya ATENDIO en Self Supervised.

$<\text{SOS}> x^{<1>} x^{<2>} \dots x^{<T_x-1>} x^{<T_x>} <\text{EOS}>$

Jane visite l'Afrique en septembre

↑ Input sequence

embeddings

Positional Encoding

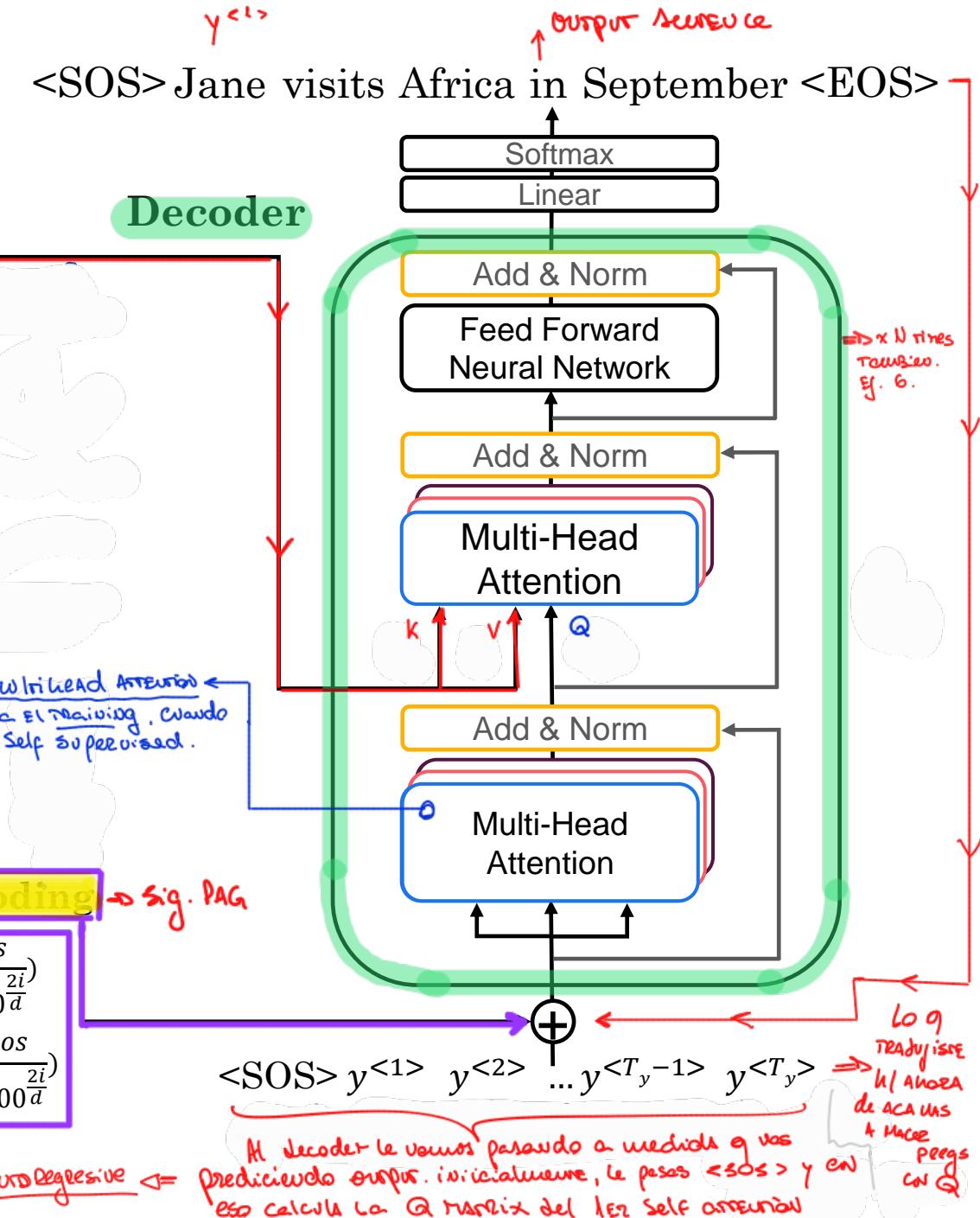
$$PE_{(pos,2i)} = \sin\left(\frac{pos}{1000^{\frac{2i}{d}}}\right)$$

$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{1000^{\frac{2i}{d}}}\right)$$

$y^{<1>} <\text{SOS}> \text{Jane visits Africa in September } <\text{EOS}>$

$y^{<1>}$

## Decoder



Softmax  
Linear

Output Sequence

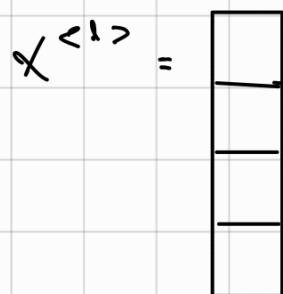
$y^{<1>}$

$\Rightarrow N$  times  
tambien.  
Ej. 6.

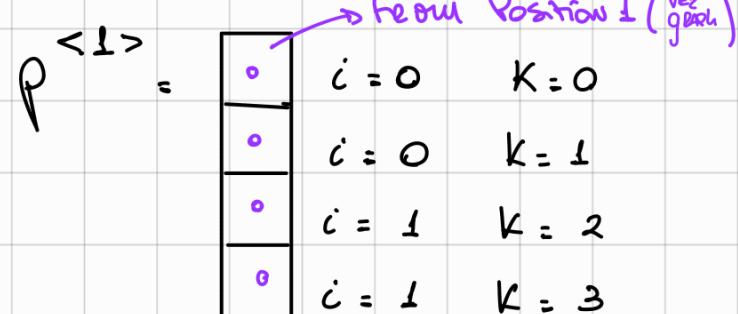
## Positional Encoding

$\Rightarrow$  anterior no toma en cuenta el orden. por eso usamos el PE.

Pense que tenemos un input embedding de tamaño  $d=4$



CREAMOS UN  
positional  
embedding



$$\left\{ \begin{aligned} p_E(\text{pos}; 2i) &= \sin \left[ \frac{\text{Pos}}{10000^{\frac{2i}{d}}} \right] \\ p_E(\text{pos}; 2i+1) &= \cos \left[ \frac{\text{Pos}}{10000^{\frac{2i}{d}}} \right] \end{aligned} \right.$$

$$\left\{ \begin{aligned} p_E(\text{pos}; 2i) &= \sin \left[ \frac{\text{Pos}}{10000^{\frac{2i}{d}}} \right] \\ p_E(\text{pos}; 2i+1) &= \cos \left[ \frac{\text{Pos}}{10000^{\frac{2i}{d}}} \right] \end{aligned} \right.$$

PE. Encoding único para  
palabra.

desde le sume  
a los inputs de los  
embeddings

- $K = 2i = 0 \rightarrow i = 0$



- $K = 2i + 1 = 0 \rightarrow i = 0$



- $K = 2i = 2 \rightarrow i = 1$



- $K = 2i + 1 = 3 \rightarrow i = 1$



$\Rightarrow$  Mayor freq

Position 1  $\Rightarrow$  Es un valor único para cada pos.

