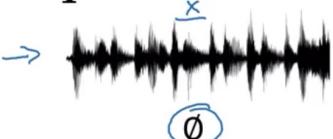


# Course 5: Sequence Models

## Week 1: Recurrent Neural Networks

### Why sequence models

#### Examples of sequence data

Speech recognition		→	"The quick brown fox jumped over the lazy dog."
Music generation		→	
Sentiment classification	"There is nothing to like in this movie."	→	
DNA sequence analysis	→ AGCCCCCTGTGAGGAACCTAG	→	AGCCCCTGTGAGGAAC <color>ACTAG</color>
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

### Notation

say we want a sequence model to automatically tell the name in this sentence

#### Motivating example

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

index position

$x^{(1)}$	$x^{(2)}$	- ..	$x^{(t)}$	- -	$x^{(n)}$
1	1	0	1	1	0 0 0 0
$y^{(1)}$	$y^{(2)}$	- ..	$y^{(t)}$	- ..	$y^{(n)}$

$T_x = \text{length of input sequence} = 9$   $\downarrow$  can be different  
 $T_y = \text{length of output sequence} = 9$

$x^{(i)l}$ :  $t^{\text{th}}$  element in the sequence of training examples i

$T_x^{(i)}$ : length of input sequence of training example i

$y^{(i)l}$ :  $t^{\text{th}}$  element in output sequence of i<sup>th</sup> example

$T_y^{(i)}$ : length of the output sequence in the i<sup>th</sup> training sample

For NLP, we need to decide how to represent individual words in the sequence.

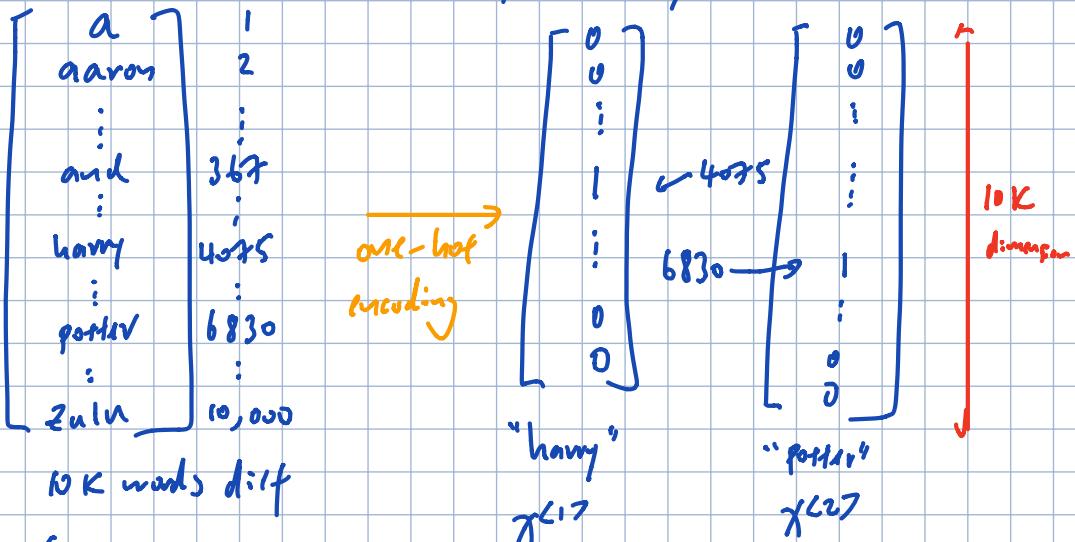
i.e. "Harry"  $\leftrightarrow x^{(1)}$  : how to represent

## Representing words

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

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad \dots \quad x^{<9>}$

come out with a vocabulary / dictionary, i.e. all the words we will be using



(get this by

① choose top 10K words  
in training set

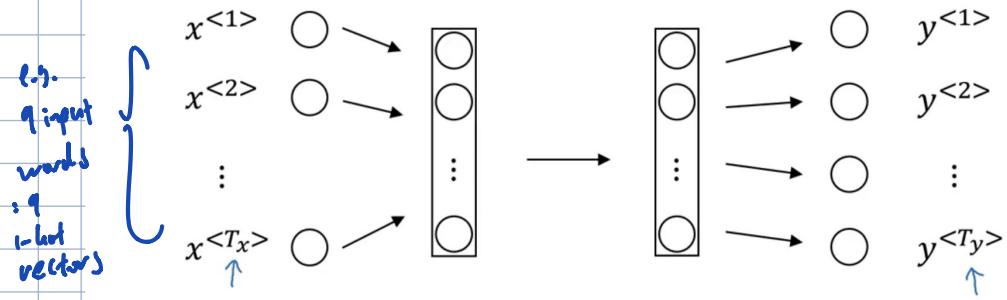
② ordine

so we have 9  $x^{<t>}$ , each a 10K dimension vector

goal:  $x \rightarrow y$

# Recurrent Neural Network Model

Why not a standard network?



$$\begin{matrix} T_x & \omega \\ T_y & \alpha \end{matrix} \xrightarrow{\rho} \begin{matrix} T_x & \omega + T_x(c_{i+1}) \\ T_y & \alpha + T_y(c_{i+1}) \end{matrix}$$

Problems:

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

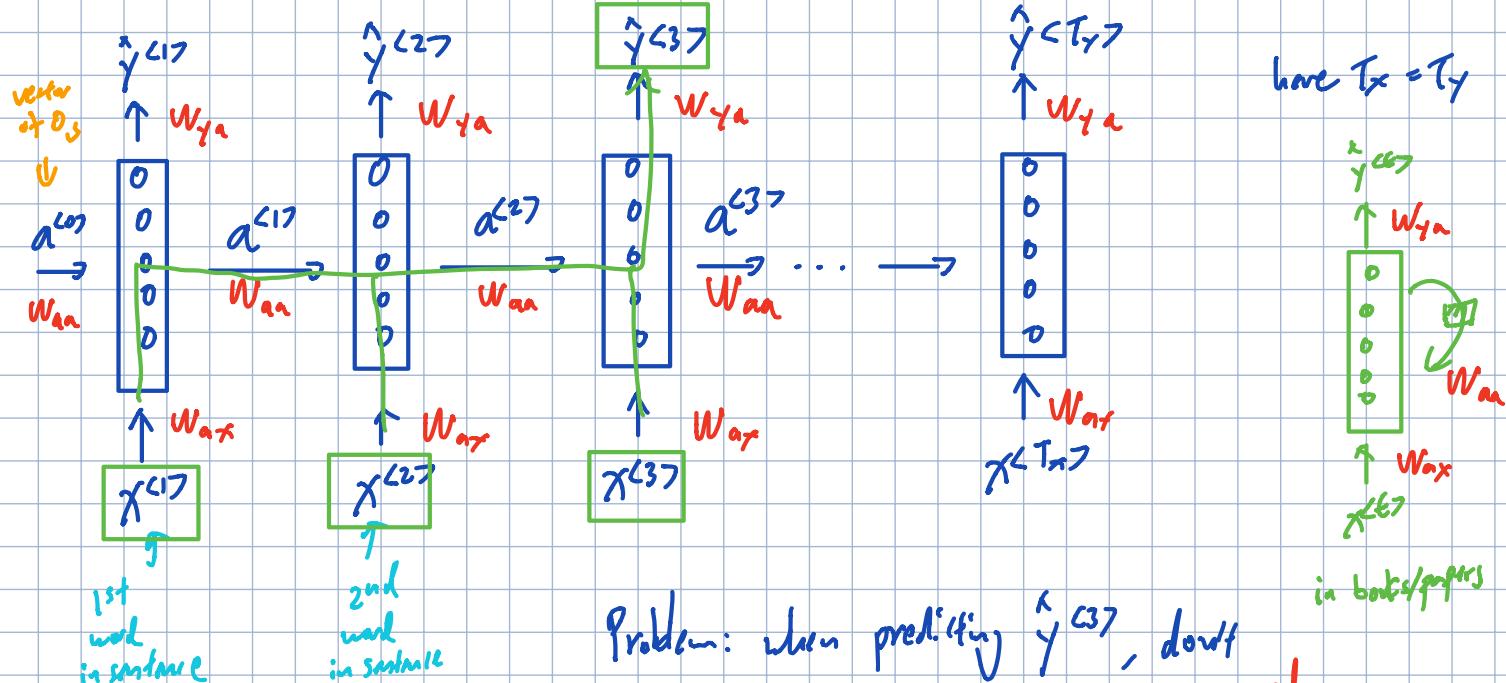
↳ e.g. "harry" at position  $1: x^{<1>}$ , but can be transferred

if "harry" at position  $t: x^{<t>}$ .

We want similar effect as CNN, where features

learned are generalized quickly to other parts of the image

## RNN



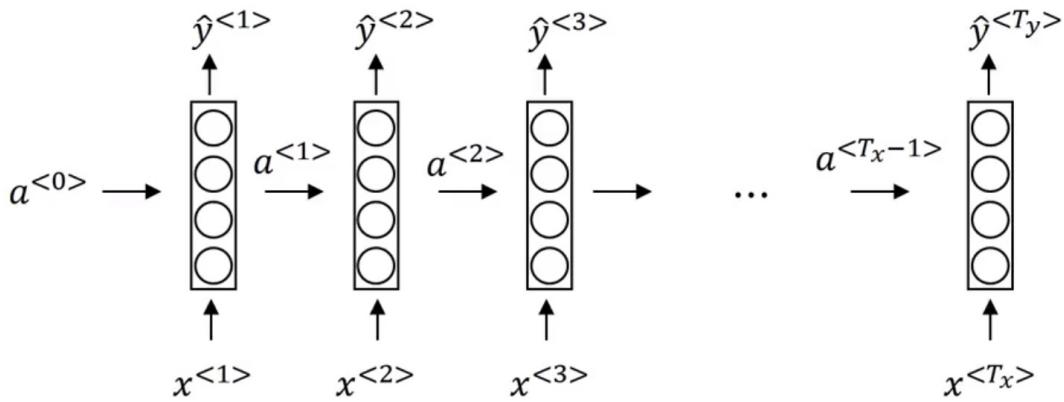
in books/papers

need  
Bi-directional  
RNNs  
(BRNNs)

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

He said, "Teddy bears are on sale!"

# Forward Propagation



$$a^{<0>} = \vec{0}$$

$$a^{<1>} = g(W_{aa} a^{<0>} + W_{ax} x^{<1>} + b_a) \leftarrow \text{tanh/ReLU}$$

vanishing gradient  
↳ have ways to solve

$$y^{<1>} = g_2(W_{ya} a^{<1>} + b_y) \leftarrow \begin{matrix} \text{Sigmoid/Softmax} \\ \text{binary multiple} \end{matrix}$$

activation (compute diff)

$$a^{<t>} = g(W_{aa} a^{<t-1>} + W_{ax} x^{<t>} + b_a) \quad \text{stack horizontally}$$

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

Simplified RNN notation

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

$(100, 100)$   $\xrightarrow{100}$   $(100, 100)$   $\xrightarrow{10,000}$

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

|||

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

$$[W_{aa}; W_{ax}] = W_a \quad (100, 100)$$

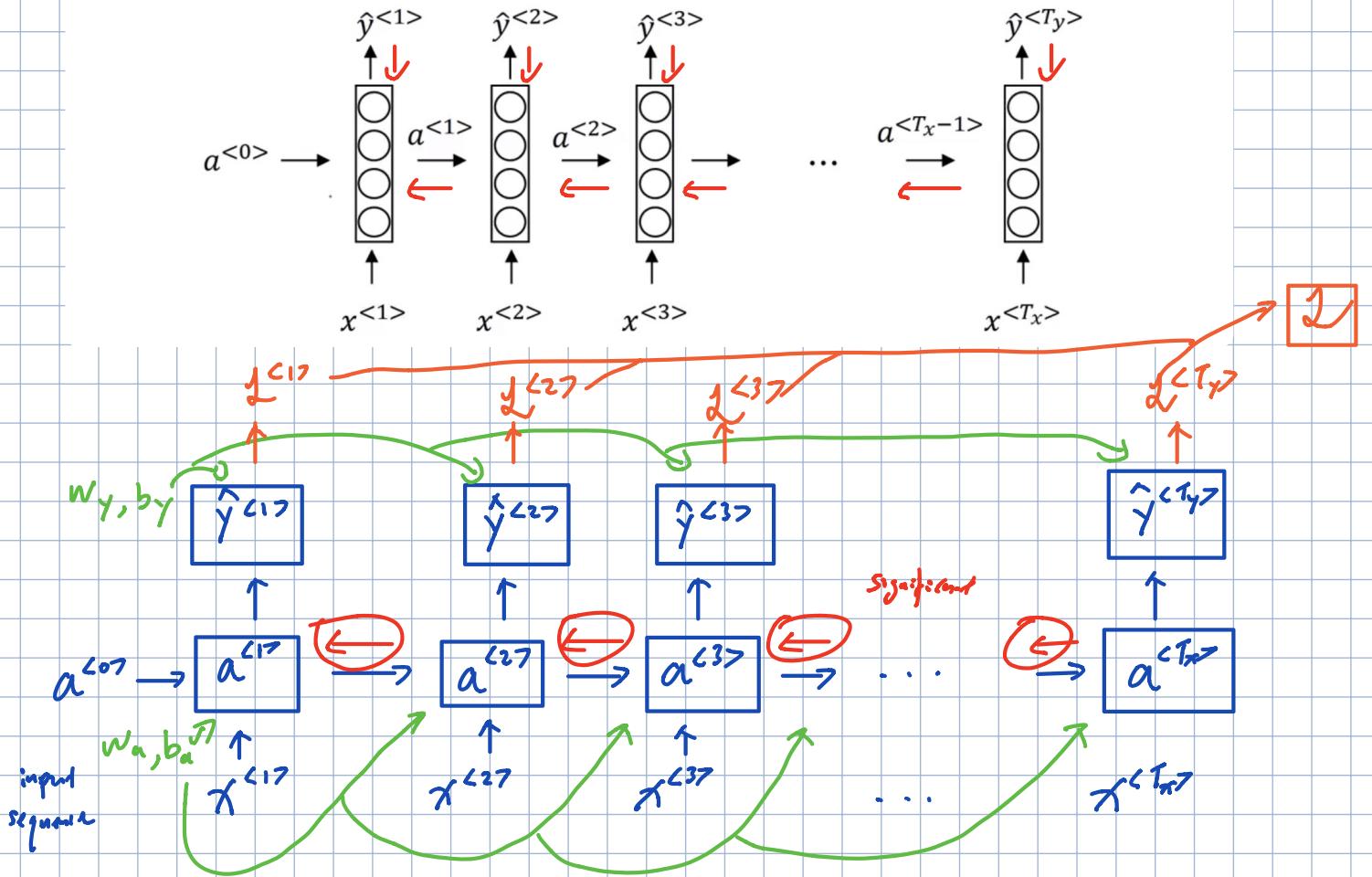
$$a^{<t>} = g(W_a [a^{<t-1>} \mid x^{<t>}] + b_a)$$

$$[a^{<t-1>} \mid x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} \in \mathbb{R}^{100 \times 1000}$$

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

## Backpropagation through time

# Forward propagation and backpropagation



To get back prop, need loss function. We define an element-wise loss first, which is supposed for a certain word in the sequence.

$$L^{(t)}(\hat{y}^{(t)}, y^{(t)}) = -y^{(t)} \log \hat{y}^{(t)} - (1-y^{(t)}) \log (1-\hat{y}^{(t)})$$

standard logistic  
regression loss /  
cross-entropy loss

This is the loss associated with a single prediction at a single position or at a single time step.

Overall loss for the entire sequence:

$$\mathcal{L}(\hat{y}, y) = \sum_{t=1}^{T_y} \mathcal{L}^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

## Different types of RNNs

So far we have seen RNN architecture where # of inputs ( $T_x$ ) = # of outputs ( $T_y$ ). For some applications,  $T_x$  may not equal to  $T_y$ .

### Examples of sequence data

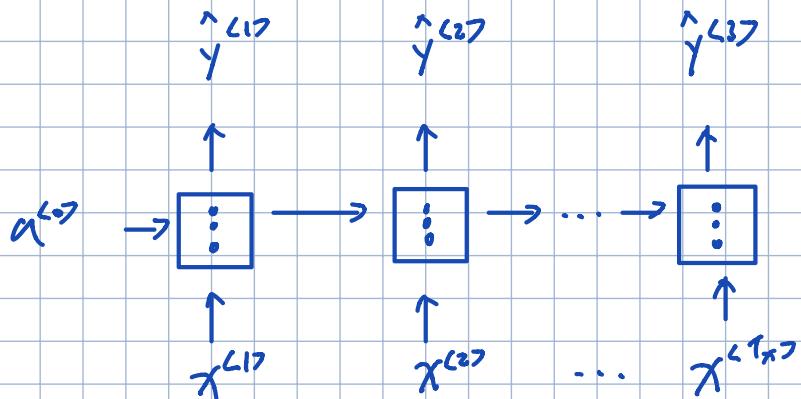
Speech recognition		$\rightarrow$	"The quick brown fox jumped over the lazy dog"
Music generation		$\rightarrow$	$\text{y}$
Sentiment classification		$\rightarrow$	
DNA sequence analysis	$\rightarrow$ AGCCCTGTGAGGAACCTAG	$\rightarrow$	AGCCCTGTGAGGAACCTAG
Machine translation	Voulez-vous chanter avec moi?	$\rightarrow$	Do you want to sing with me?
Video activity recognition		$\rightarrow$	Running
Name entity recognition	$\rightarrow$ Yesterday, Harry Potter met Hermione Granger.	$\rightarrow$	Yesterday, Harry Potter met Hermione Granger. Andrew Ng

not the same as input!

some type sequence,  
different lengths

## Examples of RNN architectures

$$T_x = T_y$$

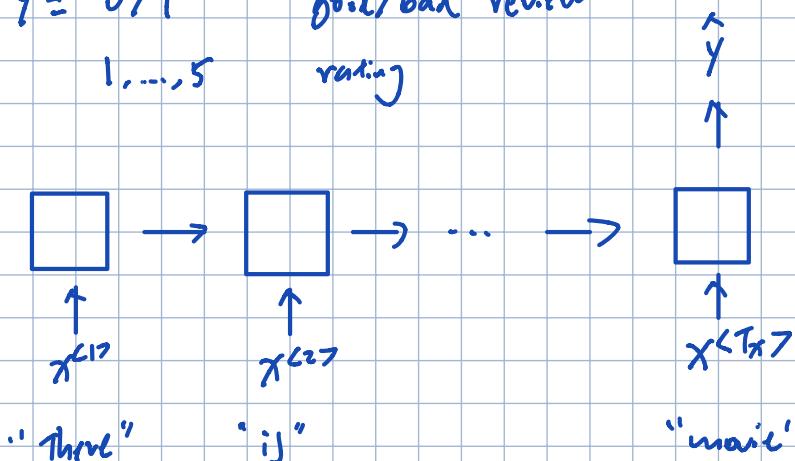


Many-to-Many (same input/output lengths)

## Sentiment classification

$x = \text{text}$  e.g. movie review

$y = 0/1$  good/bad review  
 $1, \dots, 5$  rating

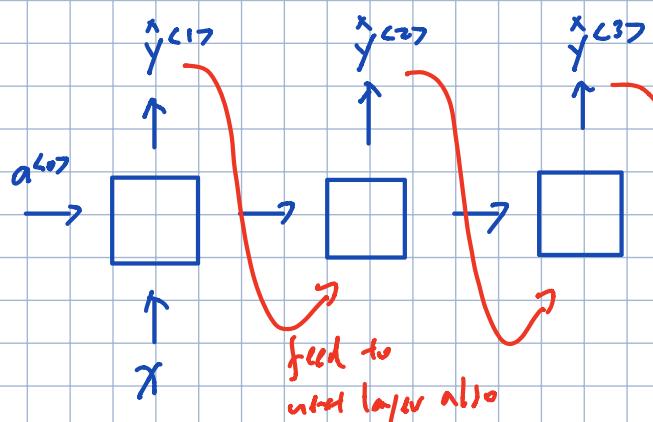


many - to - one

also have one - to - one  
 ↳ just the word SN

## Music Generation

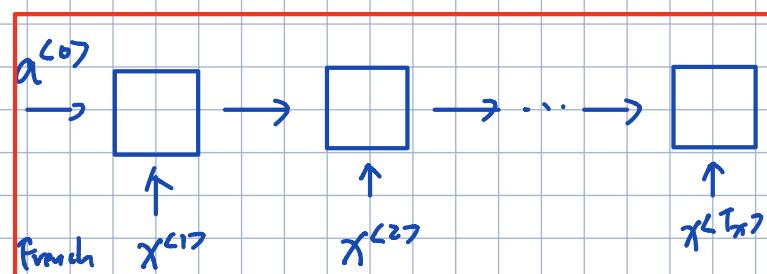
$x \rightarrow \hat{y}^{(1)} \hat{y}^{(2)} \dots \hat{y}^{(T_y)}$



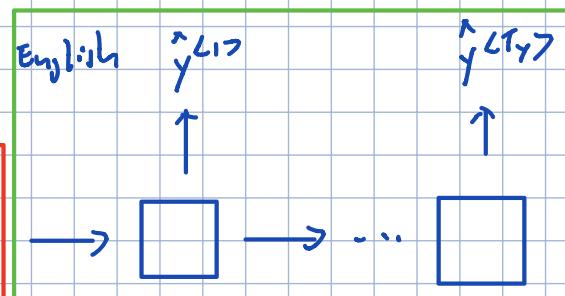
music note

one - to - many

## Machine Translation

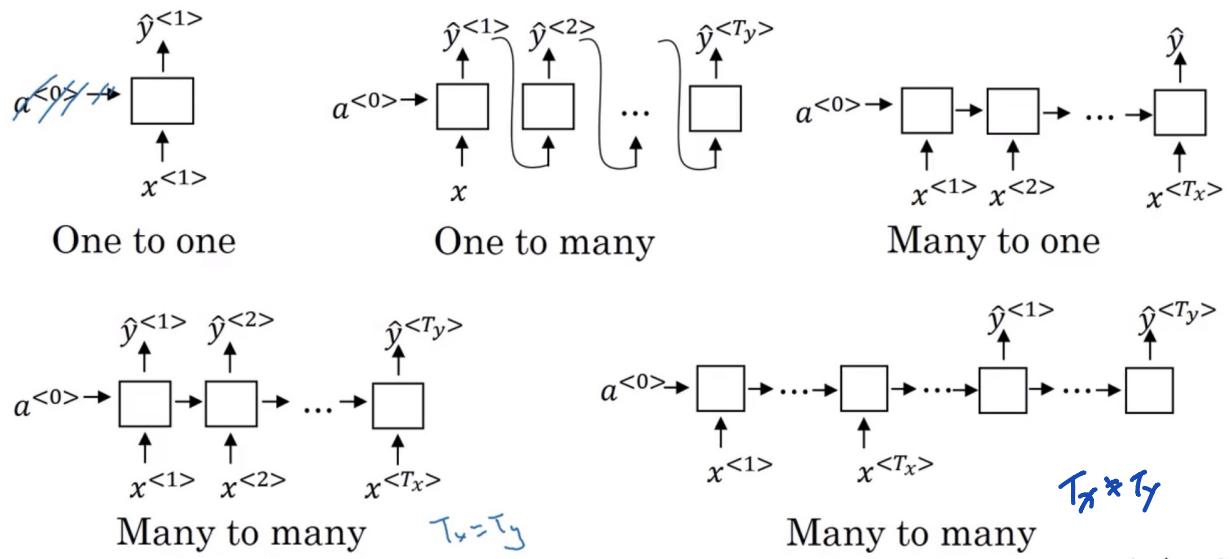


encoder



many - to - many  
 (different input/output lengths)

# Summary of RNN types



Andrew Ng

## Language model and sequence generation

### What is language modelling?

Speech recognition

The apple and pair salad. } sound the same

→ The apple and pear salad.

$$\begin{array}{l} \text{model} \\ \text{output} \end{array} \left\{ \begin{array}{l} P(\text{The apple and pair salad}) = 3.2 \times 10^{-13} \\ P(\text{The apple and pear salad}) = 5.8 \times 10^{-10} \end{array} \right. \xrightarrow{\text{much more likely}}$$

Given any sentence, language modelling tells us the probability of sentence

$$\text{i.e. } P(\text{sentence}) = ?$$

$$\hookrightarrow P(y^{<1>} \rightarrow y^{<2>} \rightarrow \dots \rightarrow y^{<T_y>})$$

where  $y$  rather than  $x$   
output                    input

# Language modelling with an RNN

Training set: large corpus of english text

e.g. a sentence in training set

Cats average 15 hours of sleep a day. (EOS)

$y^{(1)}$   $y^{(2)}$   $y^{(3)}$  ...

and -sentence

(EOS)

tokens

we ignore punctuation

what if words are not in our vocab?

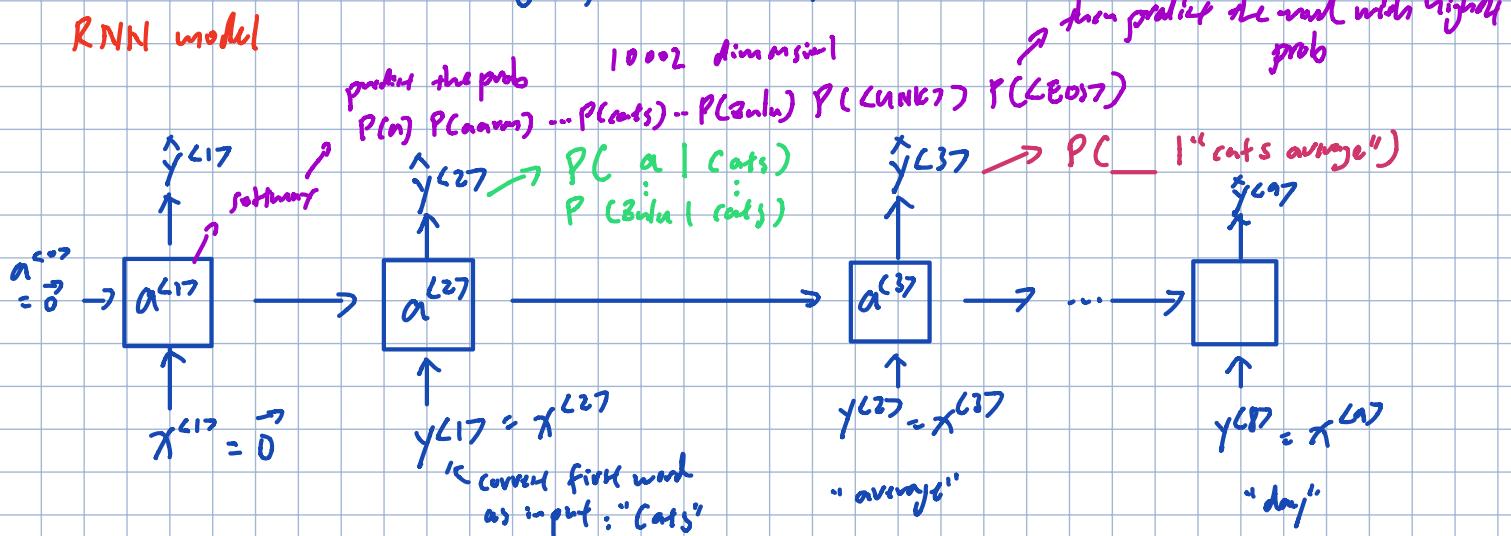
The Egyptian ~~Xan~~ is a bread of cat. (EOS)

(UNK) for unknown words

↳ model the chance of unknown words instead of the specific word

In RNN model, we setting input  $\pi^{(t-1)} = y^{(t-1)}$

RNN model



Cats average 15 hours of sleep a day. (EOS)

Each step in the RNN look at one set of preceding words, and ask what is the distribution over the next word?

↳ predict 1 word at a time, left to right.

To train the RNN, need loss function:

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

overall loss:

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

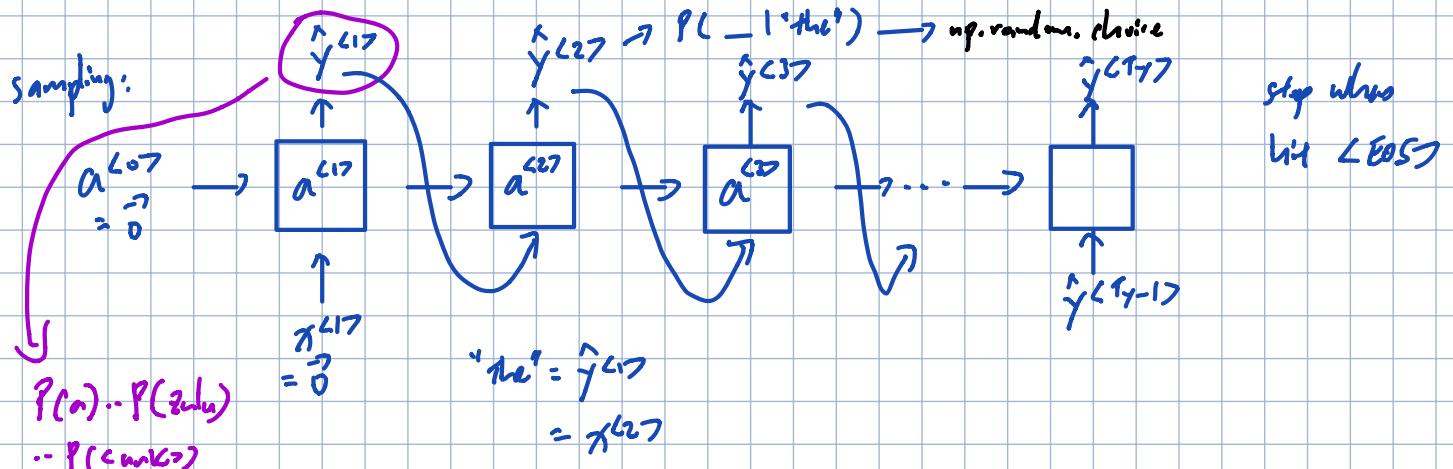
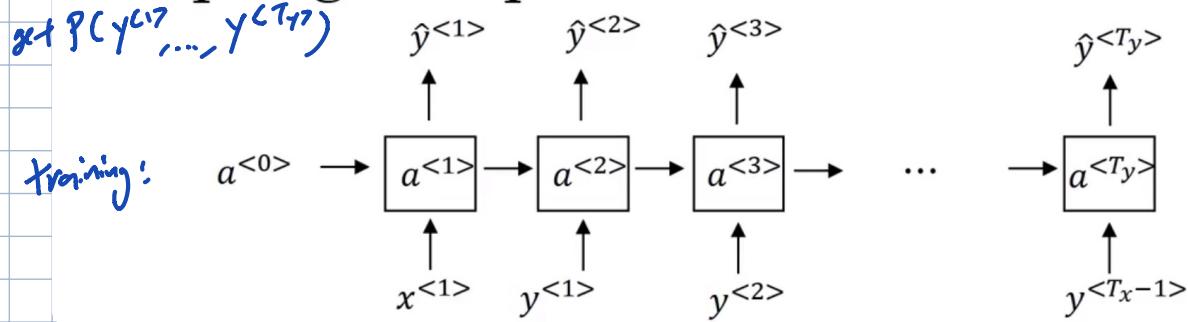
1. q. given a new sentence with 3 words:

can compute  $P(y^{(1)}, y^{(2)}, y^{(3)}) \rightarrow$  probability of the 3-word sentence

$$= P(y^{(1)}) P(y^{(2)} | y^{(1)}) P(y^{(3)} | y^{(1)}, y^{(2)})$$

## Sampling novel Sequences

### Sampling a sequence from a trained RNN



by np. random. choice

i.e. "The" chosen

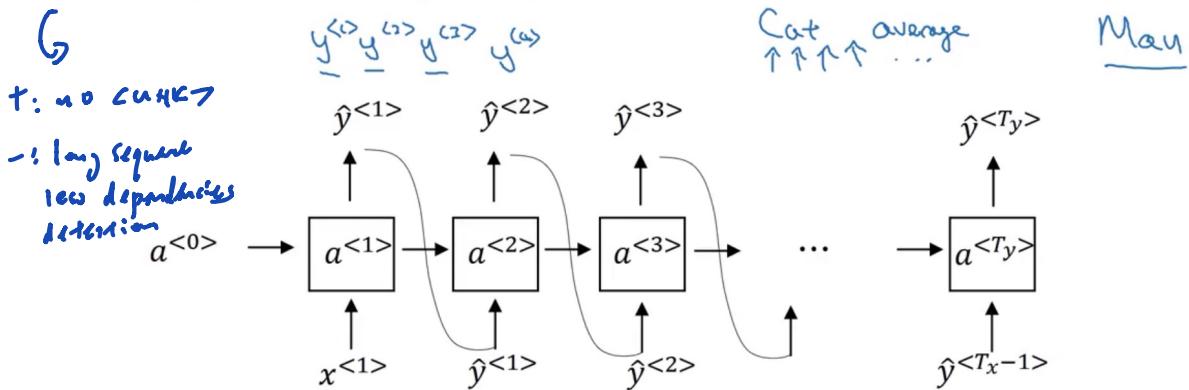
this is how we generate a randomly chosen sentence from a trained RNN language model!

→ this is word-level RNN, also can have character-level RNN

# Character-level language model

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

Vocabulary = [a, b, c, ..., z, Ȑ, ., , ;, ȏ, ..., ȑ, A, ..., Z] characters



## 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.

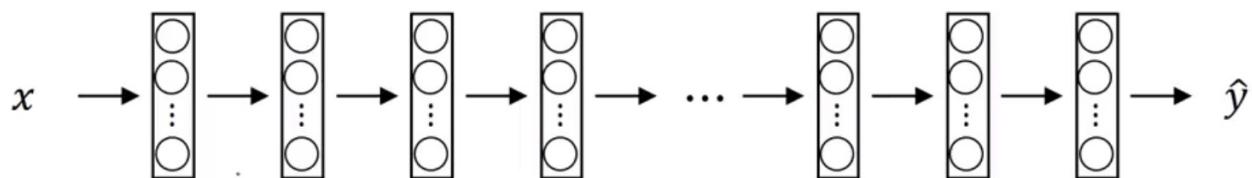
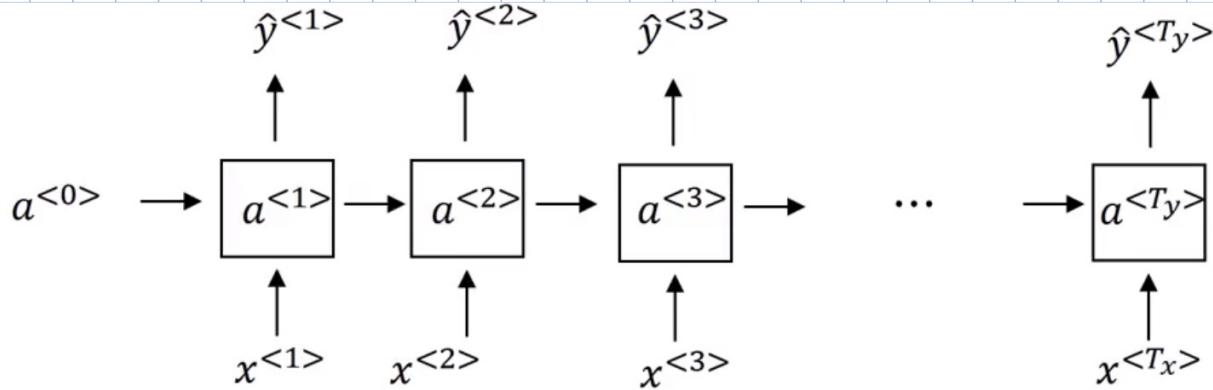
And subject of this thou art another this fold.

When lesser be my love to me see sabl's.

For whose are ruse of mine eyes heaves.

## Vanishing gradients with RNNs

e.g. The cat, which ate ... , was full. } language has long-term dependency  
 The cats, which ate ... , were full. }



backprop are hard for deep NN, to affect the weights of earlier layers.

In RNN there are a lot of local influences, i.e.  $\hat{y}^{<3>}$  is mainly influenced by  $x^{<1>} \rightarrow x^{<2>} \rightarrow x^{<3>}$  (i.e. values that are close to  $\hat{y}^{<3>}$ )

Exploding gradients also happens in RNN, but it is less common.

↳ apply gradient clipping

## Gated Recurrent Unit (GRU)

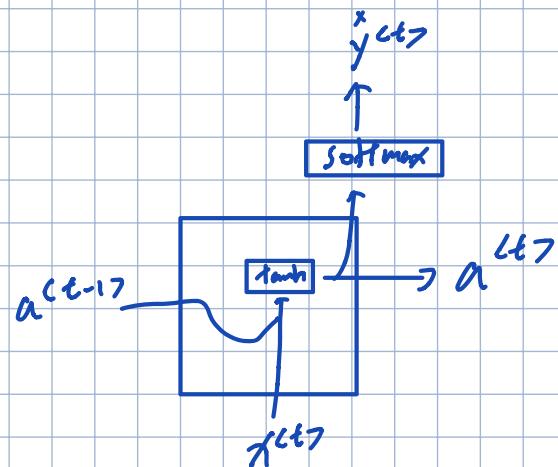
GRU is a modification to the RNN hidden layer, makes it better to capture long range connections and helps with the vanishing gradient.

RNN unit:

activations at time  $t$  of RNN

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

$\nwarrow \text{tanh}$



### GRU (simplified)

e.g. the cat, which already ate ..., was full.

as we read the sentence, left to right there is a new variable,

$C$  = memory cell  $\rightarrow$  provides memory if cat is singular/plural

at time  $t$ ,  $C^{(t)}$ . the GRU will output  $a^{(t)} = C^{(t)}$  (for GRU, they are

same, but for LSTM, they are different)

at every time step, we consider  $\tilde{C}^{(t)}$ , i.e. a candidate for replacing  $C^{(t)}$

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

Important! GRU have a gate:  $\Gamma_u$ , "u" for "update", value between 0 and 1.

we can think of  $\Gamma_u$  always 0 or 1 for intuition (most of the time)

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

sigmoid

update  $C^{(t)}$

/ don't update  $C^{(t)}$



The  $\tilde{C}^{(t)}$  is a candidate,  $\Gamma_u$  will decide if we actually update it.

$$\Gamma_u = 1$$

$$\Gamma_u = 0$$

$$\Gamma_u = 0$$

$$C^{(t)} = \dots \dashrightarrow \dots$$

$$= 1 \quad \text{remember!}$$

The rat, which already ate ..., was full.

$\leftarrow C^{(t)}$  is set 0 or 1, depend on whether the word are singular/plural.

GRU will memorise the values at the  $C^{(t)}$  all the way until "word"

$\Gamma_u$  decides when do we update this value

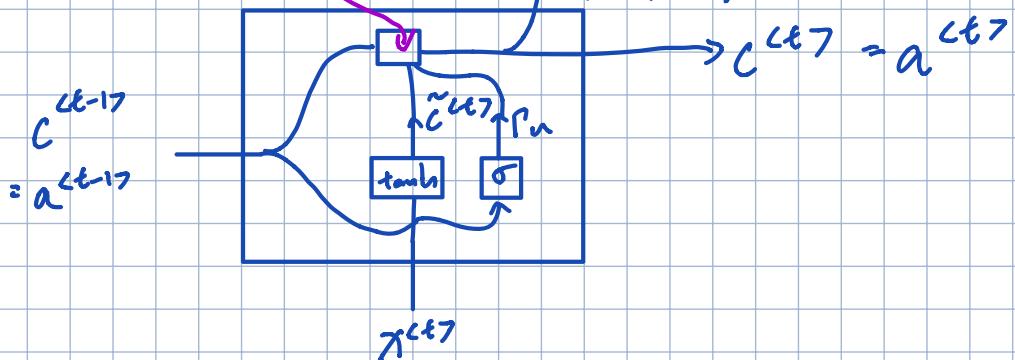
e.g. update @ "cat", and carry till "word". then decide job done and no need to memorize anymore.

element-wise

$$L7 \quad C^{(t)} = \Gamma_u * \tilde{C}^{(t)} + (1 - \Gamma_u) * C^{(t-1)}$$

if  $= 1$ , set  $C^{(t)} := \tilde{C}^{(t)}$

$\rightarrow t+1 \text{ max} \rightarrow y$



$\because \Gamma_u$  is quite easy to set to 0

$\therefore C^{(t)}$  is good at maintaining the value for the cell

$\hookrightarrow$  even across many time steps

$\hookrightarrow$  helps with vanishing gradient problems.

$\hookrightarrow$  long range dependencies

$C^{(t)}$  and  $\tilde{C}^{(t)}$  have the same dimension (i.e. a vector of hidden activation units)

$\Gamma_u$  also have the same dimension

## Full GRU

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

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

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

[another gate, "r" stands for "relevance"]

[how relevant is  $c^{t-1}$  to compute next candidate for  $c^{t+1}$ ]

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

[alternative symbol in literature]

## Long Short Term Memory (LSTM) unit

### GRU and LSTM

#### GRU

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

$$\begin{aligned} \Gamma_u &= \sigma(W_u[c^{t-1}, x^{t+1}] + b_u) \\ \Gamma_r &= \sigma(W_r[c^{t-1}, x^{t+1}] + b_r) \end{aligned} \quad \left. \begin{array}{l} \text{2 gates} \\ \text{value for LSTM} \end{array} \right\}$$

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

$$\underline{a^{t+1}} = \underline{c^{t+1}}$$

gives the option of keeping  
the old values and add  
new value in the memory cell

more powerful / generic

$\downarrow$

LSTM

gate  $a^{t+1}$  explicitly:  
 $\rightarrow$  no  $\Gamma_r$       output  $a^{t+1}$

$$\begin{aligned} \tilde{c}^{t+1} &= \tanh(W_c[a^{t-1}, x^{t+1}] + b_c) \\ \Gamma_u &= \sigma(W_u[a^{t-1}, x^{t+1}] + b_u) \\ \Gamma_f &= \sigma(W_f[a^{t-1}, x^{t+1}] + b_f) \\ \Gamma_o &= \sigma(W_o[a^{t-1}, x^{t+1}] + b_o) \end{aligned} \quad \left. \begin{array}{l} \text{3 gates} \\ \text{instead of } (1 - \Gamma_u), we have "forget gate" \\ \text{output gate} \end{array} \right\}$$

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

↑  
element-wise  
vector × vector

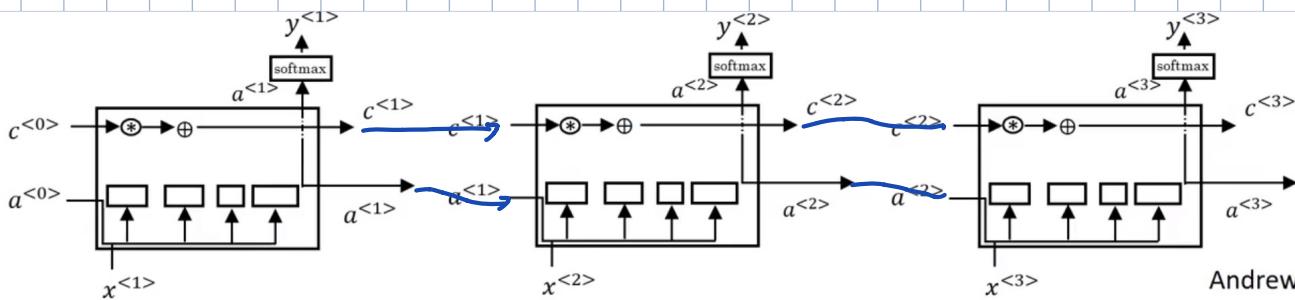
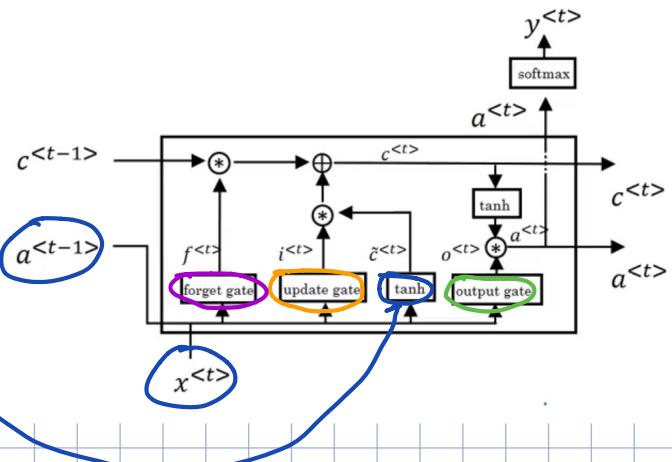
$$a^{t+1} = \Gamma_o * \tanh(c^{t+1})$$

# LSTM in pictures

$$\begin{aligned}\tilde{c}^{<t>} &= \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c) \\ \Gamma_u &= \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u) \\ \Gamma_f &= \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f) \\ \Gamma_o &= \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o) \\ c^{<t>} &= \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>} \\ a^{<t>} &= \Gamma_o * \tanh c^{<t>}\end{aligned}$$

add  
c <t-1>  
people  
connection

If we have them in parallel



$c^{<t>}$ 's values can be passed down easily, for some forget and update gate

$$\therefore C^{LST} = C^{LST}$$

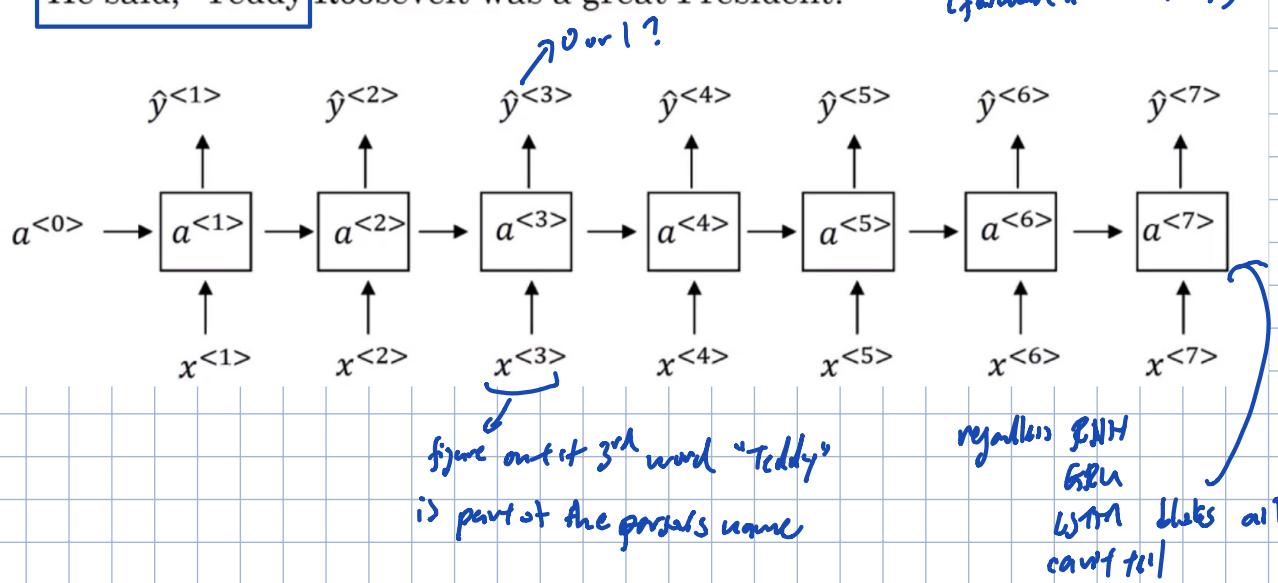
## Bidirectional RNN (BRNN)

### Getting information from the future

He said, "Teddy bears are on sale!"

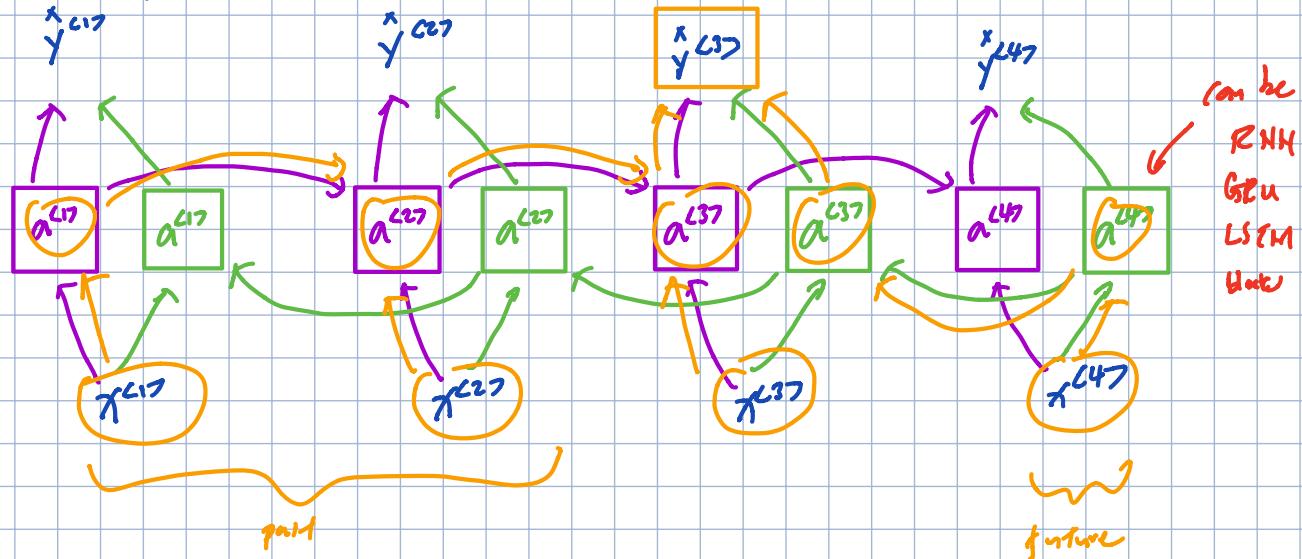
He said, "Teddy Roosevelt was a great President!"

unidirectional RNN  
(forward direction only)



## BRNN

$$\text{Eq. } \hat{y}^{[t]} = g(W_y [a^{[t]}, a^{[t]}] + b_y)$$



forward recurrent layer: start from  $a^{[1]}$

Acyclic graph

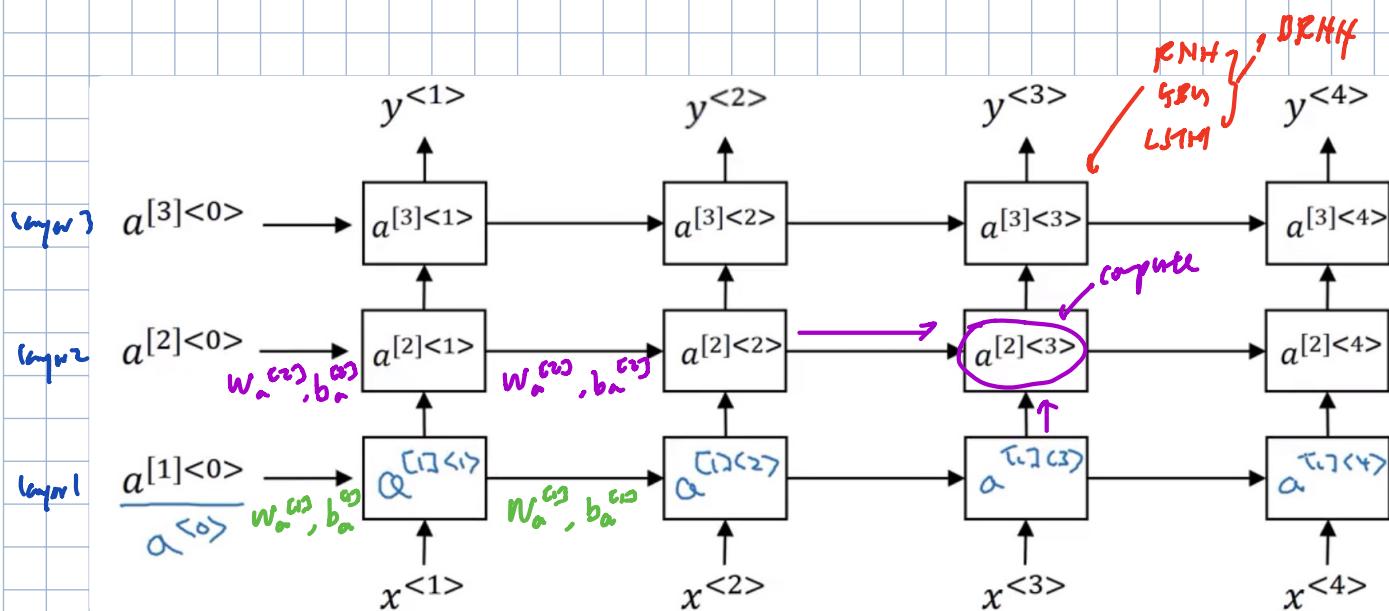
backward recurrent layer: start from  $\bar{a}^{[4]}$

\* both are forward prop, just different direction in time!

\* Disadvantage of BRNN: need entire sequence before prediction.

## Deep RNN

$a^{[1,2,3,4]}$  layer  
time



$$a^{[2,3,4]} = g(W_a^{[2]} [a^{[2,3,4]}, a^{[2,3,4]}] + b_a^{[2,3,4]})$$

## Week 2 : Natural Language Processing & Word Embeddings

### Word Representation

#### Word representation

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

vocab

$$|V| = 10,000$$

1-hot representation

Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ 0 \end{bmatrix}$
$\hookrightarrow O_{5391}$	$\hookrightarrow O_{9853}$				

same closeness

as other pair of words

what's the next word

I want a glass of orange juice.

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

Hard for algorithm to generalize from "orange juice" to "apple juice"

"0" stands for one-hot

bag of words

treats each word as a thing on its own  $\rightarrow$  weakness

Inner product of any 1-hot vector = 0

$\hookrightarrow$  can't distinguish any pair of these vectors, same distance

$\hookrightarrow$  algorithm doesn't know "apple" and "orange" are much more similar than "King" or "Queen"- "orange".

Andrew Ng

#### Featurized representation: word embedding

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	+1	-0.95	0.97	0.00	0.01
Royal	0.00	0.02	0.93	0.95	-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						
Alive						

↑  
300  
for  
tiny

↓

300 dimension  
vector for "man"  $\rightarrow l_{5391}$

It would be nice for each word in the dictionary to learn some features (with associated values)

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

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

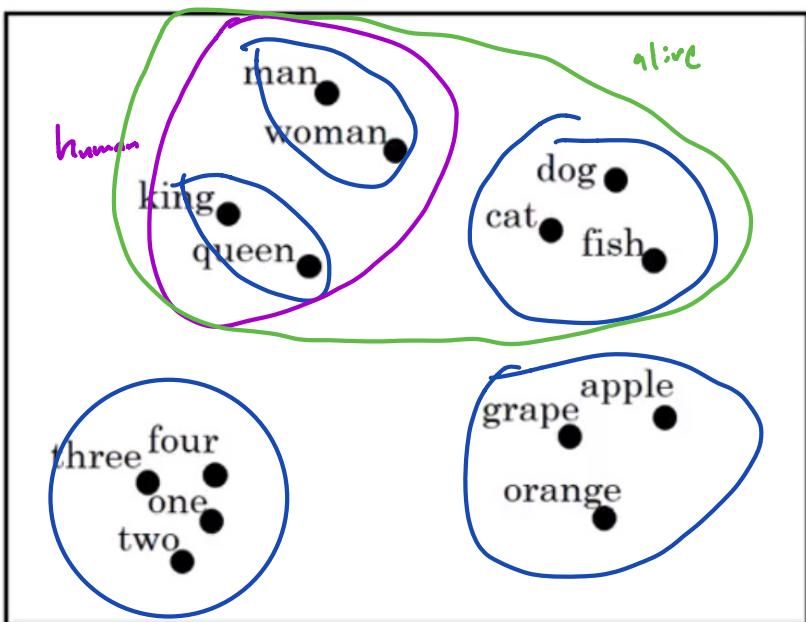
Some of their features are similar

↳ increase the odds the learning algorithm has figure out orange juice is a thing

↳ quick figure out apple juice is also a thing.

## visualizing word embeddings

From 300 dimensions (for each word)

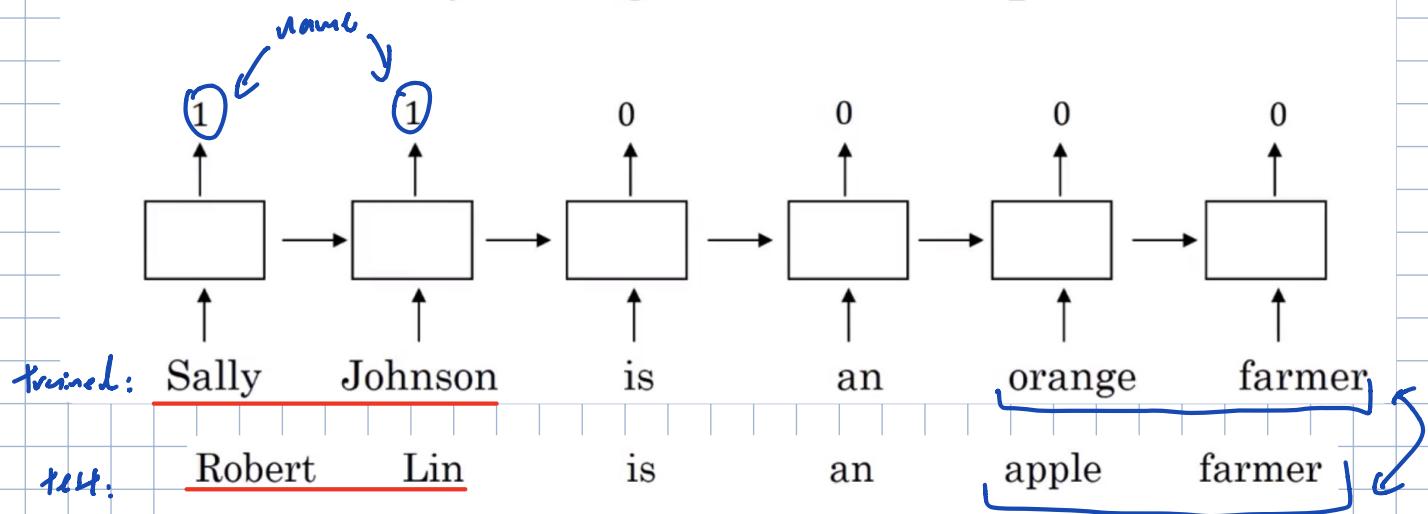


2-Dimensional

t-SNE

## Using word embeddings

### Named entity recognition example



↳ "orange farmer" and "apple farmer" are similar

↳ helps to identify "Robert Lin" is also a name

↳ how about "durian Cultivator"?

↳ need to learn "durian" is a fruit

Transfer learning!

1B - 100B words unlabeled data

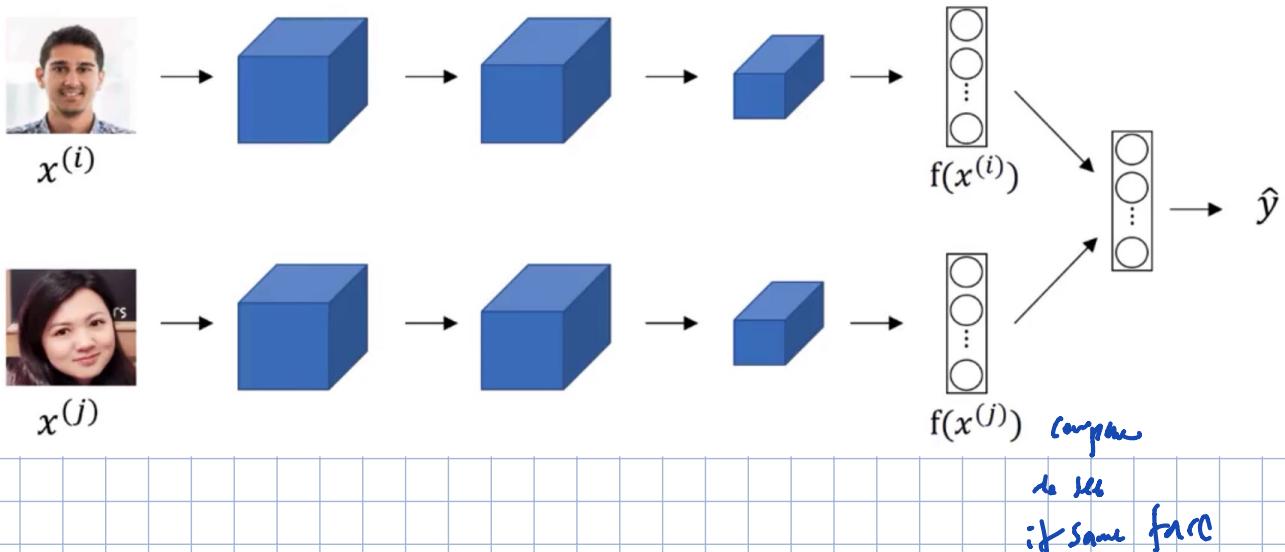
↳ figure out "durian" is a fruit

↳ apply 100K training set → for named recognition

### Transfer learning and word embeddings

1. Learn word embeddings from large text corpus. (1-100B words)  
(Or download pre-trained embedding online.)
2. Transfer embedding to new task with smaller training set.  
(say, 100k words)  $\rightarrow 10,000 \rightarrow 300$
3. Optional: Continue to finetune the word embeddings with new data.

# Relation to face encoding



## Properties of word embeddings

### Analogy

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

Qn:

$$e_{5391} - e_{9853} \approx$$

Man  $\rightarrow$  Woman as King  $\rightarrow$  ?



$$e_{\text{Man}} - e_{\text{Woman}} \approx e_{\text{King}} - e_{?}$$

$$e_{\text{Man}} - e_{\text{Woman}} \approx$$

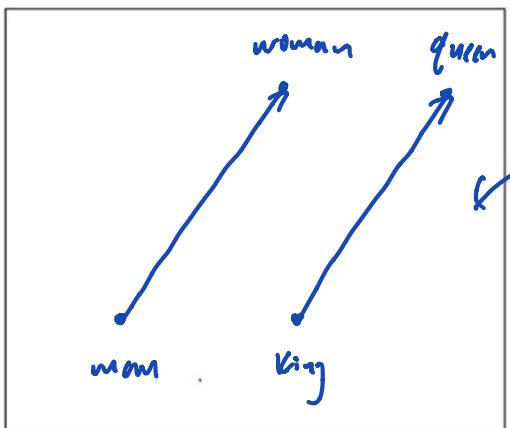
$$\begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

gender

$$e_{\text{King}} - e_{\text{Queen}} \approx$$

$$\begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

# Analogy using word vectors



Not t-SNE  
↳ highly nonlinear process  
multiple visual feature space

$$e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\text{queen}}$$

300 dimensional

Solve:

Find word  $w$  so that:  $\underset{w}{\operatorname{argmax}} \sim(\ell_w, \ell_{\text{king}} - \ell_{\text{man}} + \ell_{\text{woman}})$

some similarity function

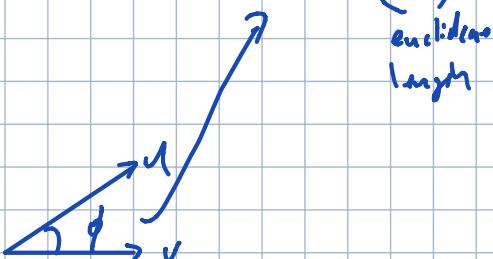
Most commonly used similarity function: Cosine Similarity

Cosine Similarity

$$\sim(\ell_u, \ell_v) = \sim(\ell_u, \ell_v - \ell_u + \ell_v)$$

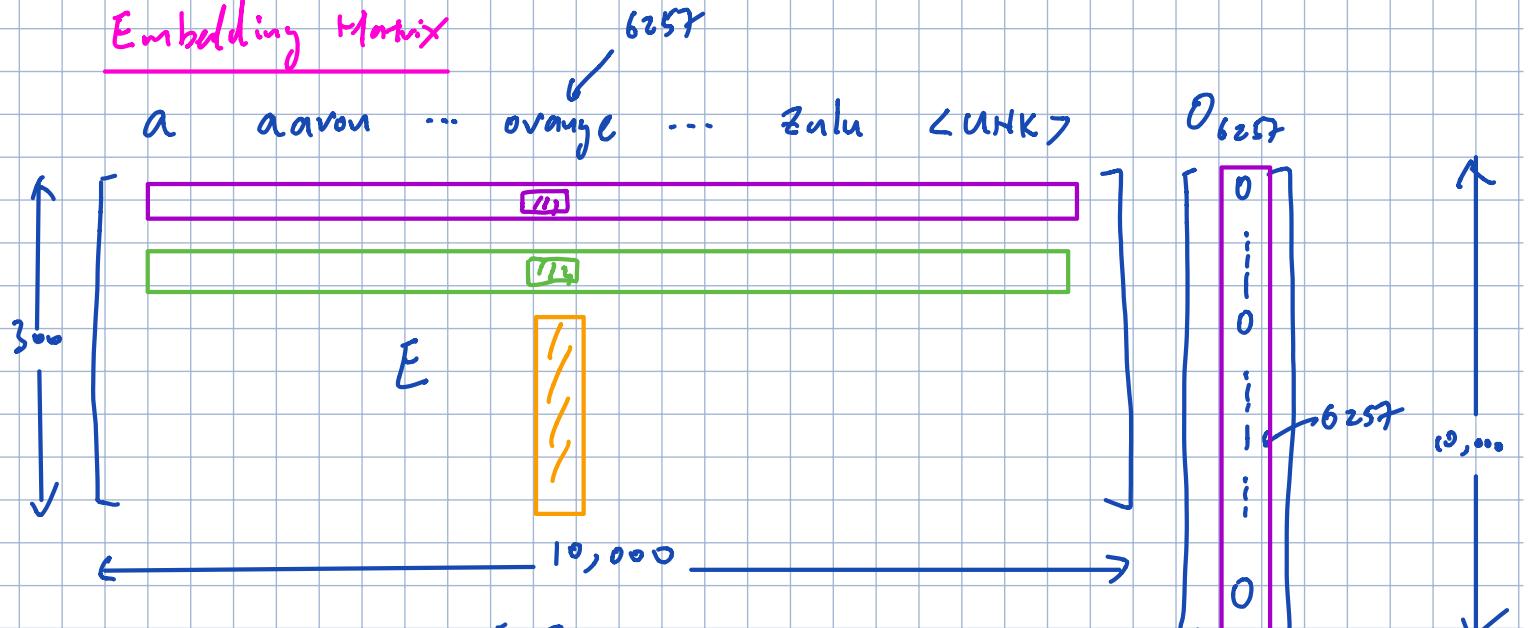
$$\sim(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$

$u^T v$  cosine prod of  $u$  and  $v$   
 $\|u\|_2$  /  $\|v\|_2$   
euclidean length



$$\text{can also use } \|u - v\|$$

## Embedding Matrix



$$E \cdot O_{6257} = \begin{bmatrix} \text{orange embedding} \\ \text{other embeddings} \end{bmatrix} \quad \begin{matrix} \leftarrow \\ \text{select out "orange" embedding} \end{matrix}$$

$(300 \times 10,000) \quad (10,000 \times 1)$

$$= l_{6257}$$

$\hookrightarrow l_w$  = embeddings for word  $w$

$$\hookrightarrow E \cdot O_j = l_j$$

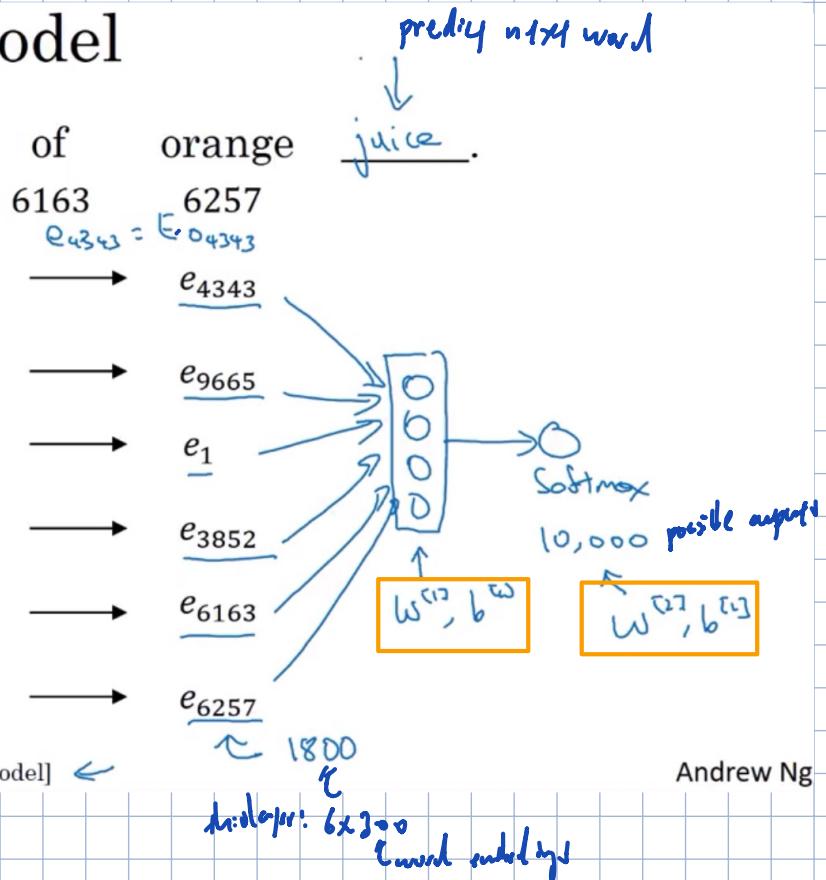
In practice, we specialised function to look up embedding

$\hookrightarrow$  :: word are  $O_s$ , not efficient.

# Learning word embeddings

## Neural language model

I	want	a	glass	of	orange
4343	9665	1	3852	6163	6257
	<sup>10^6 words</sup> → one-hot vector dim=1000 vector			$e_{4343} = E \cdot o_{4343}$	
I	$o_{4343}$			$e_{4343}$	
want	$o_{9665}$			$e_{9665}$	
a	$o_1$			$e_1$	
glass	$o_{3852}$			$e_{3852}$	
of	$o_{6163}$			$e_{6163}$	
orange	$o_{6257}$			$e_{6257}$	



We commonly use a fixed historical window: e.g. always want to predict the next word given by 4 previous words.

↳ can deal with arbitrary long sentence.

Other context/target pairs

I want a glass of orange juice to go along with my cereal.  
 context: Last 4 words  
 target: juice

context: Last 4 words

4 words on left & right

Last 1 word

Nearly 1 word

## Word2Vec

### Skip-grams

Come up with a few context to target errors to create our supervised learning problem

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

<u>Context</u>	<u>target</u>
orange	juice
orange	glass
orange	my

① randomly pick a word to be context word

② randomly pick another word, within some window, to be target word.

③ Supervised problem: given context, predict what is this randomly chosen words within the window.

### Model

Vocab size = 10,000 K

$x \longrightarrow y$   
 context C ("orange")  $\longrightarrow$  Target t ("juice")  
 6253  $\qquad\qquad\qquad$  4834

$O_c \rightarrow E$   $\rightarrow e_c = E \cdot O_c \longrightarrow$   
 param  $\downarrow$  (embedding vectors for input context)  $\qquad\qquad\qquad$   $O \rightarrow \hat{y}$   
 softmax  $\qquad\qquad\qquad$  parameters

$$\text{softmax: } P(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$$

$\theta_t$  = parameters associated with output t

$$\text{Loss: } \mathcal{L}(\hat{y}, y) = - \sum_{i=1}^{w,000} y_i \log \hat{y}_i$$

method

$$y = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \text{ 4834}$$

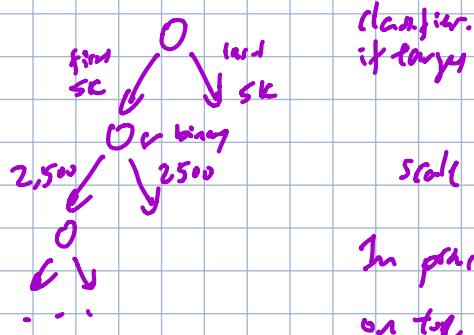
Problems with softmax classification → computational grid

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

every time evaluate this probability, need to sum over 10k.

Soln:

Hierarchical softmax



classifier.  
it starts in fam sc or land sc in the words

scale with  $\log |V|$

In practice, the tree is not balanced, i.e. common words on top, rare words are deep below.

How to sample context C?

(one context is sampled, target can be sampled within a window)

→ sample uniformly, randomly.

↪ the, of, a ... appear extremely frequently

## Negative Sampling

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

context	word	target?
orange	juice	1
orange	king	0
orange	book	0
	x	y
	← pick randomly from dictionary	

problem formulation:

given a pair of words like "orange" and "juice",  
predict is this a context-target pair?

To generate dataset, pick a context word and then target word → we example

↳ for K time:

pick some context word  
pick a random word from dictionary

Model for supervised learning  $x \rightarrow y$

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

	context	word	target?
1	orange	juice	1
2	orange	king	0
3	orange	book	0
4	orange	the	0
5	orange	of	0

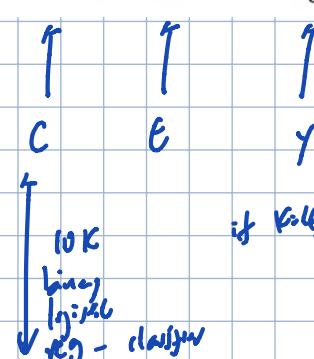
define logistic regression model:

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

orange  
6257

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

↑  
0 : juice?  
1 : king?  
0 : ?



So instead of having 1 giant  $10^k$ -way softmax, (expensive to compute)  
we have  $10^k$  binary classification problems. (cheap to compute)

[ negative sampling : have 1 positive example, then deliberately generate -ve examples.]

## Selecting negative examples

context	word	target?
orange	juice	1
orange	king	0
orange	book	0
orange	the	0
orange	of	0

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

- (1) sample according to how often different words appear.  
↳ high "the", "an", ...  
(2)  $\frac{1}{|U|}$

## GloVe word vectors

GloVe (global vectors for word representation)

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

picking C, t

$X_{ij} = 1$  if word j occurs in the context of word i  
 $\begin{matrix} t \\ C \end{matrix}$

define context and target as whether 2 words appear in close proximity

$$X_{ij} = X_{ji}$$

$X_{ij}$  is a count that captures how often do words  $i$  and  $j$  appear with each other

### Model

minimize

$$\sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) (\theta_i^T e_j + b_i + b'_j - \log X_{ij})^2$$

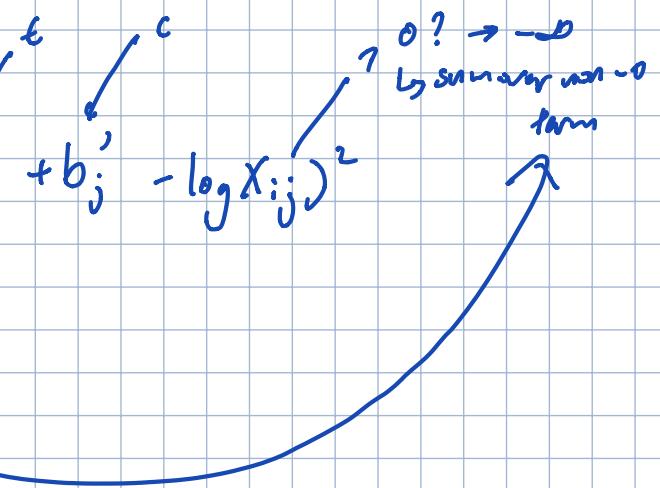
$\theta_i^T e_j$   
"weighting"

$$f(X_{ij}) = 0 \text{ if } X_{ij} = 0$$

"this" "is" ...  $\rightarrow$  stop words  
 "domain" ..  $\rightarrow$  rare words

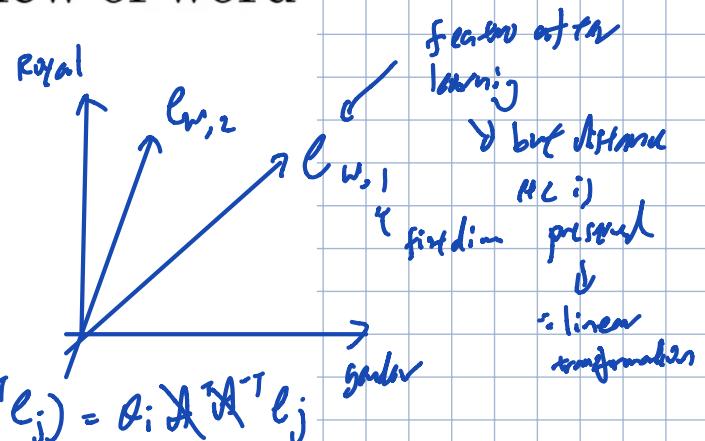
$\theta_i, e_j$  are symmetric

$$e_w^{(sim)} = \frac{e_w + \theta_w}{2}$$



### A note on the featurization view of word embeddings

	Man	Woman	King	Queen	
	(5391)	(9853)	(4914)	(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	←



$$(\theta_i \theta_i^T) (\theta_i^T e_j) = \theta_i \theta_i^T e_j$$

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

You cannot guarantee that individual components of the embeddings are interpretable

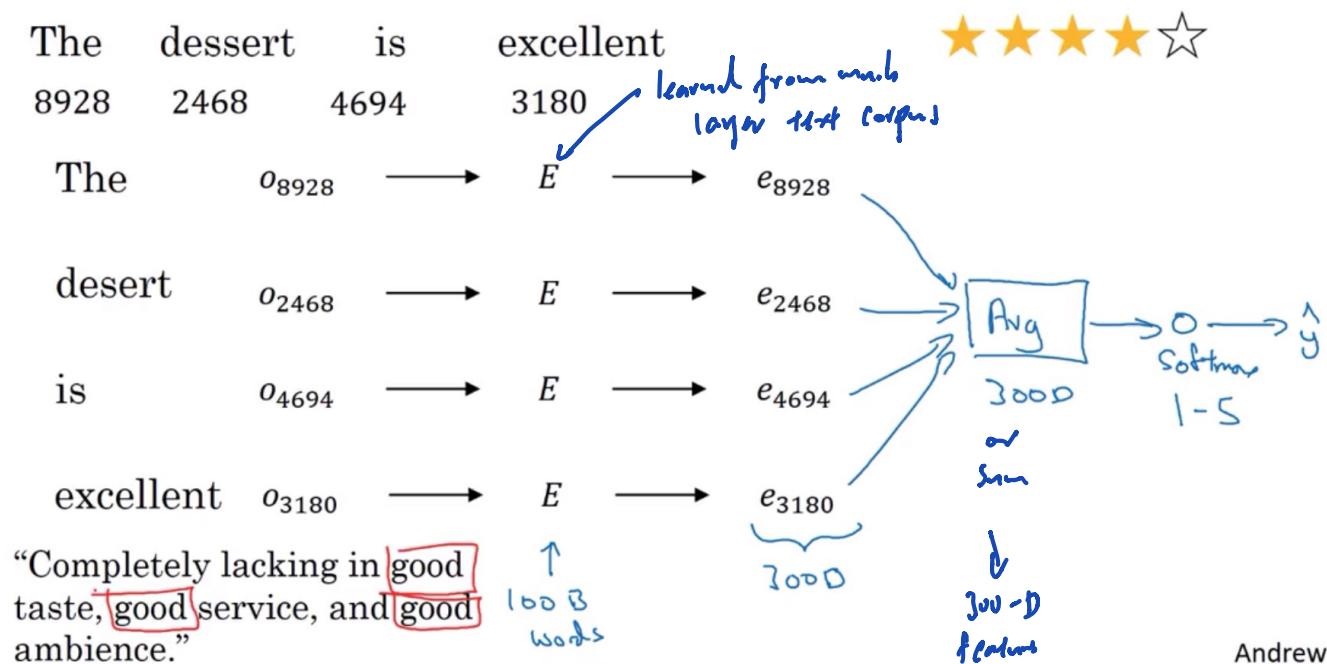
## Sentiment Classification

It's the task of looking at a piece of text and telling if someone likes or dislikes the thing they talking about.

### Sentiment classification problem



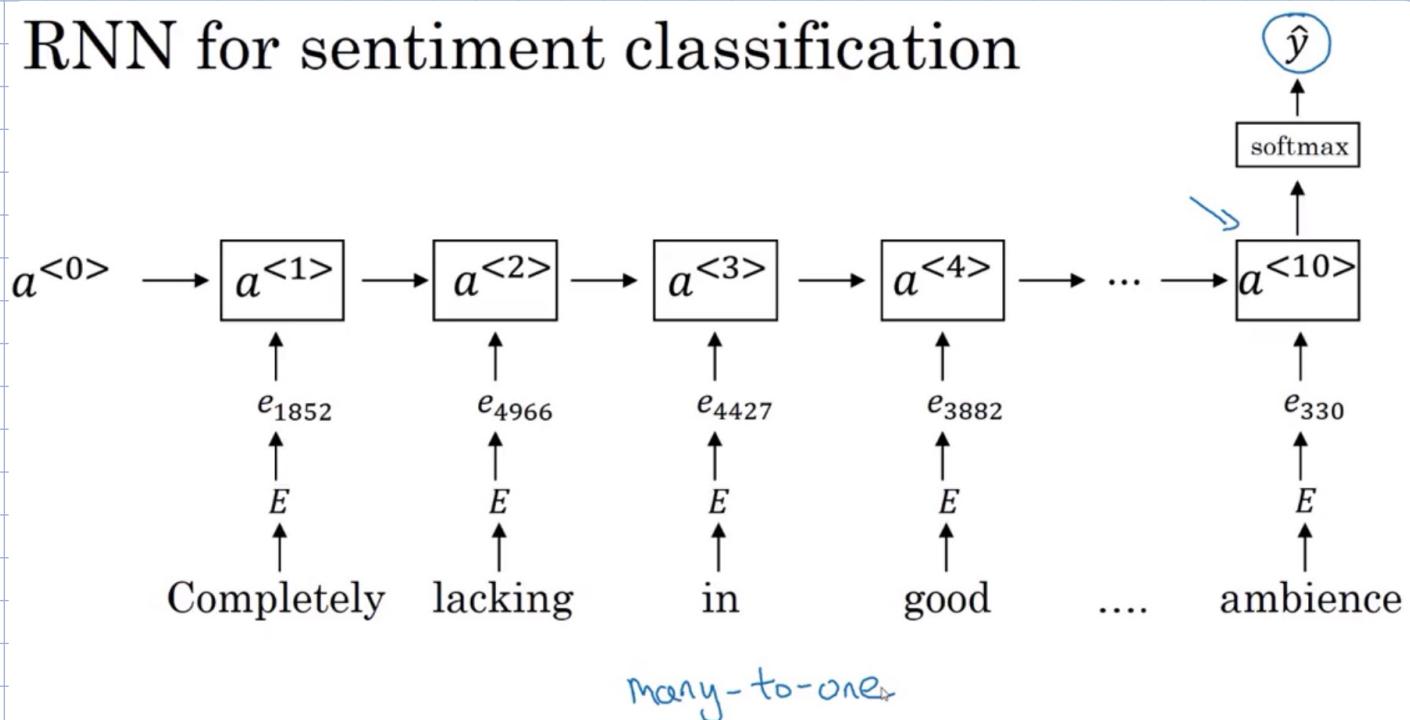
### Simple sentiment classification model



↑ problem: ignore word orders

↓ use RNN

# RNN for sentiment classification



↳ take care of word order

## Debiasing word embeddings

### The problem of bias in word embeddings

Man:Woman as King:Queen

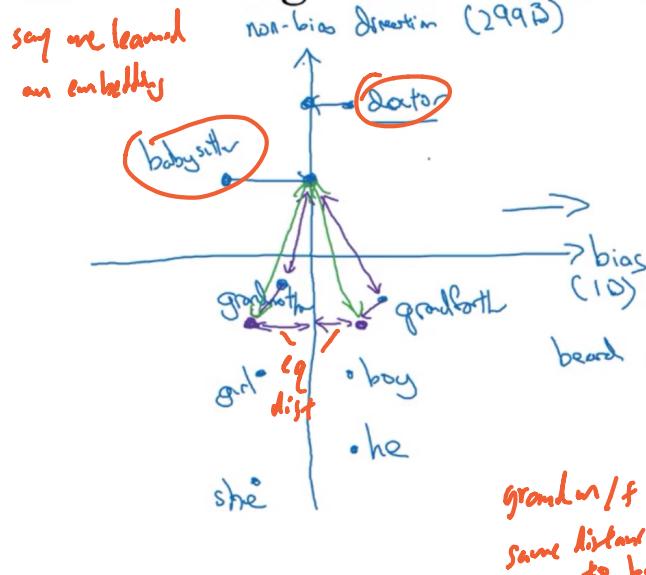
→ gender stereotype

Man:Computer\_Programmer as Woman:Homemaker X

Father:Doctor as Mother:Nurse X

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



1. Identify bias direction. (say girl)

$$\begin{cases} \text{she} - \text{she} \\ \text{male} - \text{female} \end{cases} \rightarrow \text{average}$$

2. Neutralize: For every word that is not definitional, project to get rid of bias. → Not grandfather etc  
e.g. doctor etc  
intrinsically F/M

3. Equalize pairs.

$$\rightarrow \text{grandmother} - \text{grandfather} = \text{boy} - \text{girl}$$

## Week 3: Sequence model & Attention mechanism

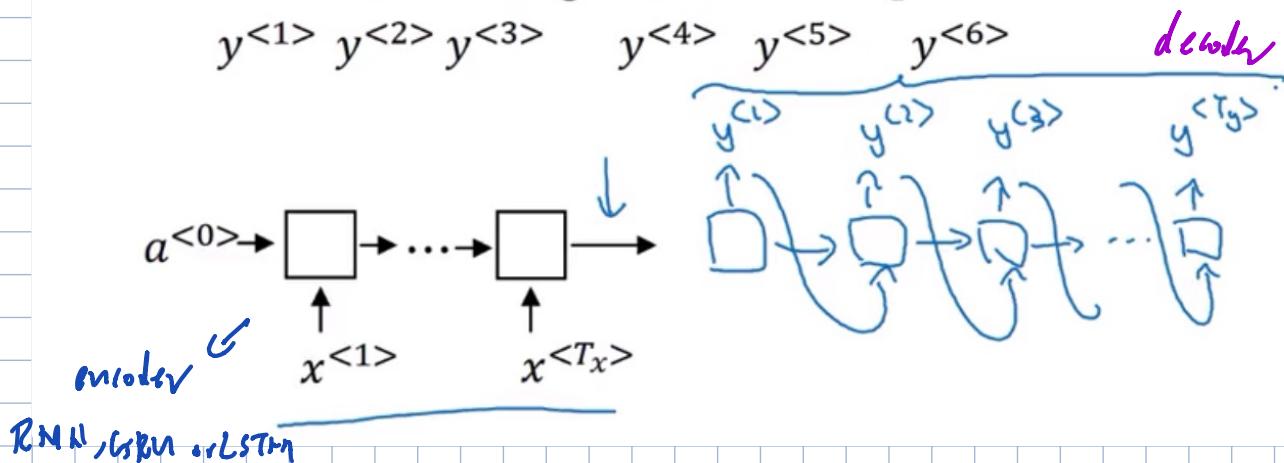
### Basic models

## Sequence to sequence model

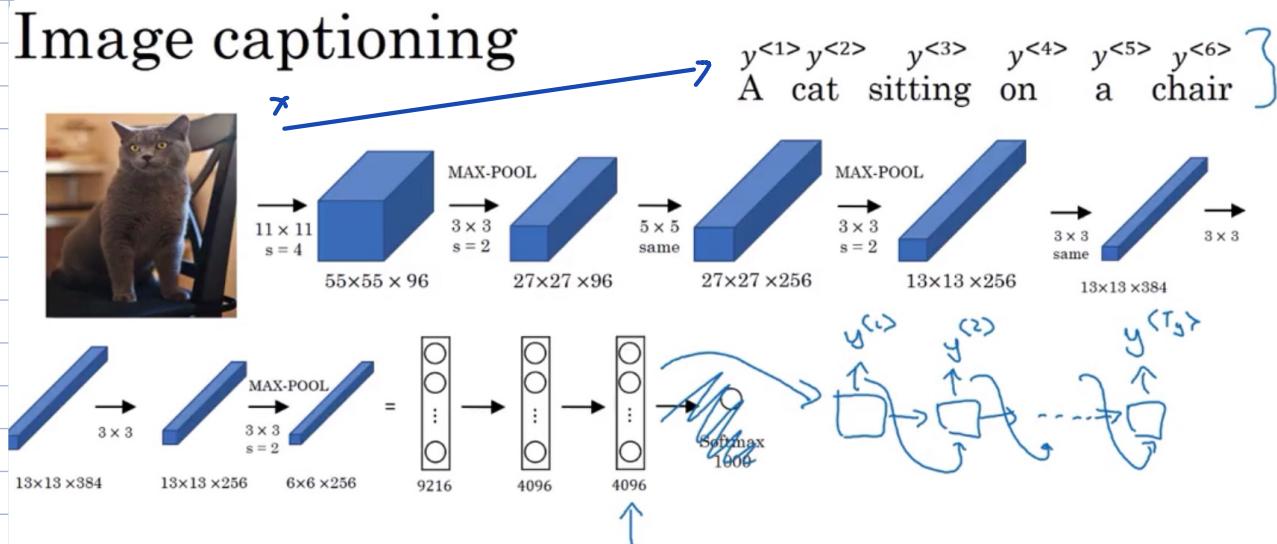
$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$

Jane visite l'Afrique en septembre

→ Jane is visiting Africa in September.



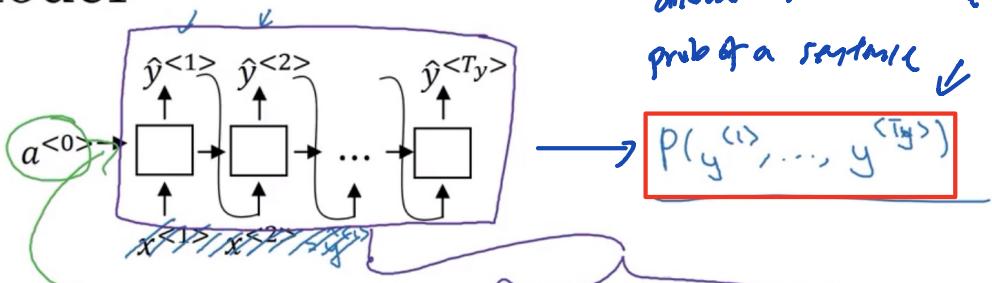
## Image captioning



## Picking the most likely sentence

# Machine translation as building a conditional language model

Language model:

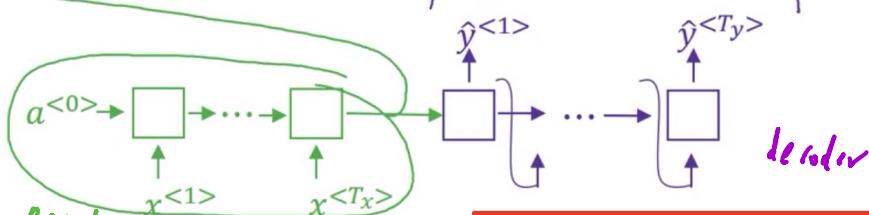


allows us to estimate  
prob of a sentence

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

Machine translation:

instead of ~~verbs~~,  
we consider networks  
that figure out representation  
for the input sentence



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

Andrew Ng

J. French

Chinese

Machine translation as building a conditional language model!

Finding the most likely translation

Jane visite l'Afrique en septembre.

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

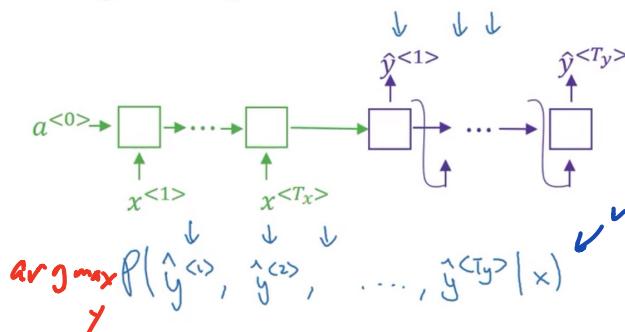
we don't want to sample output  
at random

- 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.

} different translation output each time  
algorithm to maximize conditional prob (... find the right y)  
beam search

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

Why not a greedy search?



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

greedy search:

pick the next most ~~most~~ likely word sequentially.

we want this, joint probability maximization.

$$P(y|x)$$

- 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)\$  
half forward approach

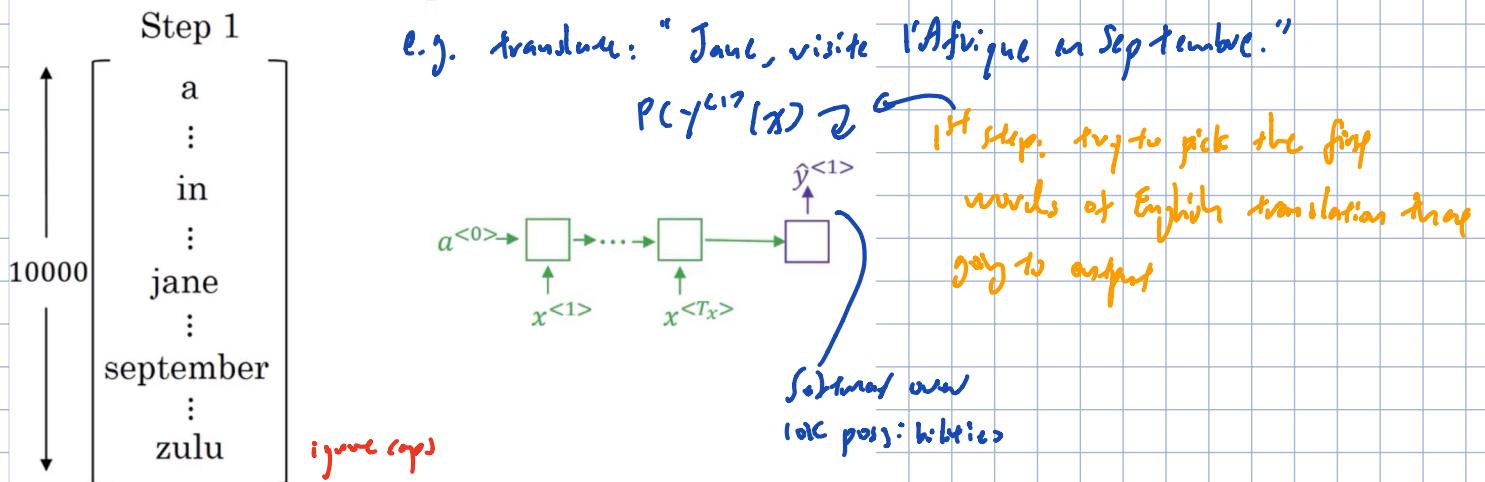
Rather than generate a sentence at random, we want to find the most likely English sentence for machine translation.

But the set of all English sentences of a certain length is too large to exhaustively enumerate → need an search algorithm.

## Beam Search

In machine translation, given a French input, we don't want to output a random English translation, but the best and most likely English translation.

### Beam search algorithm



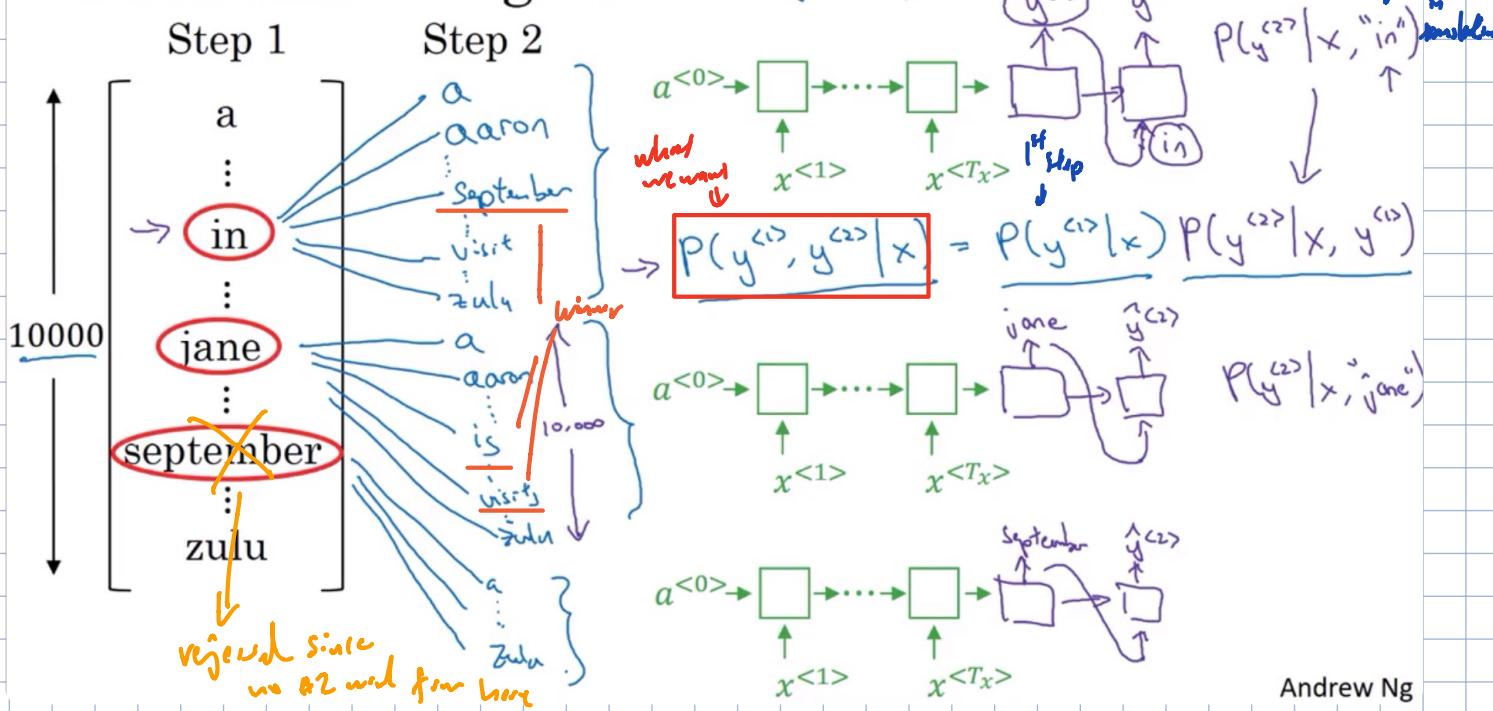
while greedy search will pick only the one most likely words and move on,  
beam search can consider multiple alternatives.

Beam width,  $B=3$  (3 possibilities)

e.g. evaluate  $P(y^{<1>} | x)$  over different choice of first words, highly  
i) "in", "Jane", "September"

Then beam search will store in memory that it will try all 3 of these words.

# Beam search algorithm ( $B=3$ )



In 2nd step, for each of these 3 choices, consider what is the 2nd word?

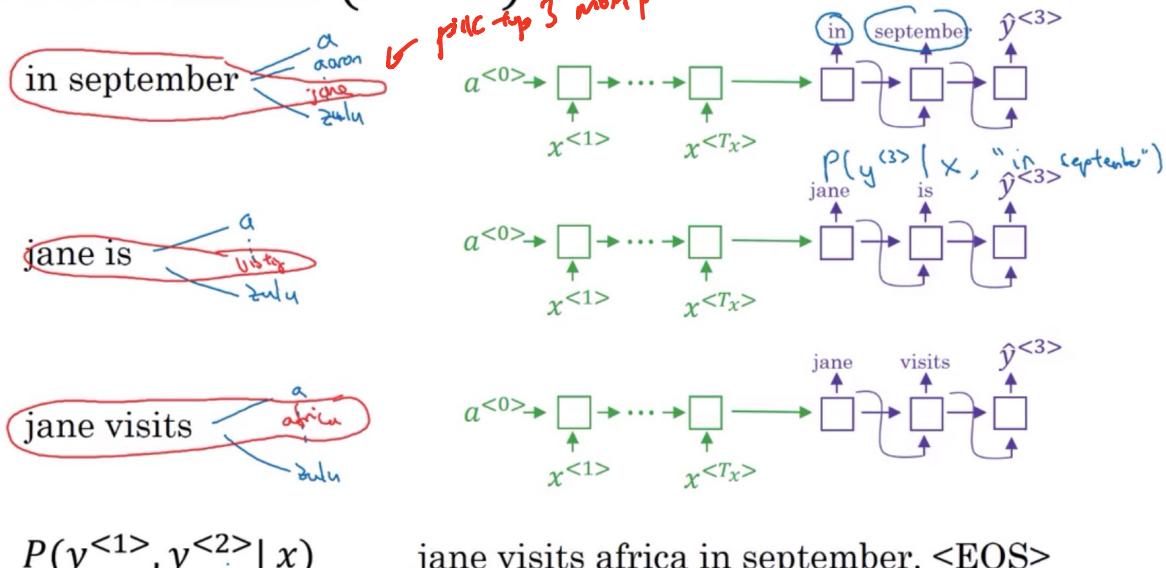
↳ In total, there will be 3x10,000 options.

↳ evaluate all these 30K options and pick the top 3.

Since we rejected 1 ("september"), we only keep track of 2 pairs  $y^{(1)}, y^{(2)}$   
(i.e. in september, jane is, Jane visits)

↳  $\therefore$  beam width = 3  $\therefore$  every step we instantiate three copies of the network  
to evaluate their parallel sentence fragments and output.

## Beam search ( $B = 3$ )



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

jane visits africa in september. <EOS>

Andrew Ng

## Refinements to Beam Search

Length normalisation is a small change to beam search to get much better results.

### Length Normalisation

beam search (max the probability)

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

⇒ log

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

all these numbers are very small, multiply them can result in numerical underflow

In practice, we take log, (so become sum)

$$\textcircled{2} \arg \max_y \sum_{t=1}^{T_y} \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$

$\max \log P(y|x)$

$\max P(y|x) \rightarrow$  same y gives max result

If the sentence is long, then the probability of that sentence is low, since multiply multiple terms together (most of them < 1)

tends towards very short translation

↳ same as the log version

$$\textcircled{3} \boxed{\frac{1}{T_y}} \sum_{t=1}^{T_y} \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$

normalize by the number of words in the translation

↓ static approach

$$\frac{1}{T_y^\alpha} (0 < \alpha < 1)$$

↑ complete normalisation by length

→ beam width = 'a'

↳ always keep track top 'a' sentences

↳ some of them are diag  
↳ final output

# Beam search discussion

Beam width B?

$1 \rightarrow 3 \rightarrow 10, \quad 100, \quad 1000, \rightarrow 3000$   
or      large      research

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

## Error analysis on beam search

Beam search is an approximate search algorithm (heuristic search algorithm)  $\rightarrow$  what it beam search make a mistake?

### Example

Jane visite l'Afrique en septembre.

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

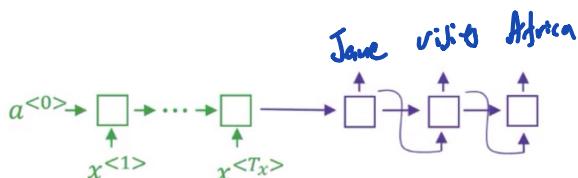
Algorithm: Jane visited Africa last September. ( $\hat{y}$ )

model has 2 components

① RNN: encoder - decoder

② Beam search

$\hookrightarrow$  can increase B ↑



RNN computes  $P(y^*|x)$  and  $P(\hat{y}|x)$ , and see which is larger

# Error analysis on beam search

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

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

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

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

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

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

Conclusion: Beam search is at fault.

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

$y^*$  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	$P(y^* x)$	$P(\hat{y} x)$	At fault?
Jane visits Africa in September.	Jane visited Africa last September.	$2 \times 10^{-10}$	$1 \times 10^{-10}$	B
...	...	—	—	R
...	...	—	—	Q
				R R R ...

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

Andrew N

## Bleu Score

Given a French sentence, there can be multiple English translations that are equally good → how to evaluate the machine translation system?  
↳ how to measure accuracy?

### Evaluating machine translation

French: Le chat est sur le tapis.

Bleu  
bilingual evaluation understudy

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.

Any machine translated sentence that is close to reference, gets a high score.  
these reference should be provided ~ point at dev set / test set.

Precision,  $\frac{7}{7}$   
 $\frac{7}{7}$   
fraction: # of ref words  
that appear  
in ref 1/2

Modified precision:

give each word  
credit up to a  
max of one  
if appears in ref

← count clip "the"  
 $\frac{2}{7} \approx$  (and sf "the")

since "the" appear  
at most 2 times  
we give limit = 2

We not only want single isolated word appearance, but also pairs of words as well.

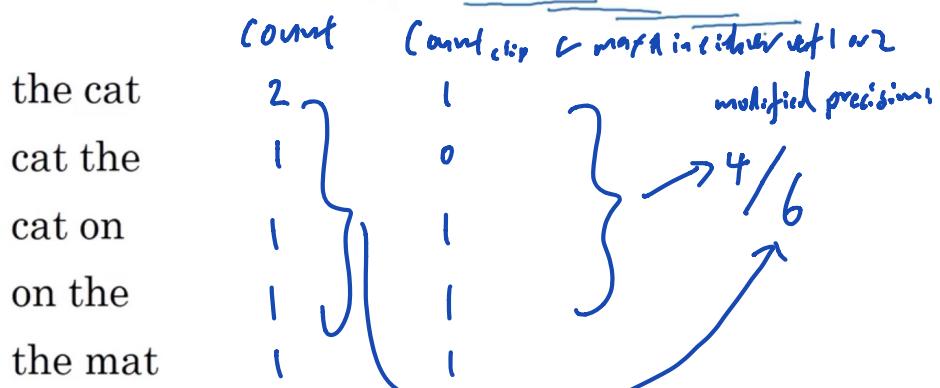
Bigram: pairs of words appear next to each other.

# Bleu score on bigrams

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. ←



[Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]

Andrew Ng

Formally:

unigram:

$$P_1 = \frac{\sum_{\text{unigrams } i} (\text{Count}_{\text{clip}}(\text{unigram}))}{\sum_{\text{unigrams } i} (\text{Count}(\text{unigram}))}$$

*modified precision for unigrams*

$$P_n = \frac{\sum_{\text{n-grams } i} (\text{Count}_{\text{clip}}(\text{n-gram}))}{\sum_{\text{n-grams } i} (\text{Count}(\text{n-gram}))}$$

*n-gram*

If MT output is exactly same as ref 1 or 2,  $P_1, P_2, \dots = 1.0$

Final Bleu score:

$P_n$  = Bleu score on n-grams only  $\rightarrow$  get  $P_1, P_2, \dots, P_4$

Log combined Bleu score:  $B_P \exp\left(\frac{1}{4} \sum_{n=1}^4 P_n\right)$

B\_P - brevity penalty  $\begin{cases} 1 & \text{if MT\_output\_length} \geq \text{reference\_output\_length} \\ \exp(1 - \text{reference\_output\_length}/\text{MT\_output\_length}) & \text{otherwise} \end{cases}$

! Shorthand translation

get easier Bleu score

Bleu usually for MT ! multiple good translation

not in speech recognition ! usually have ground truths to compare with.

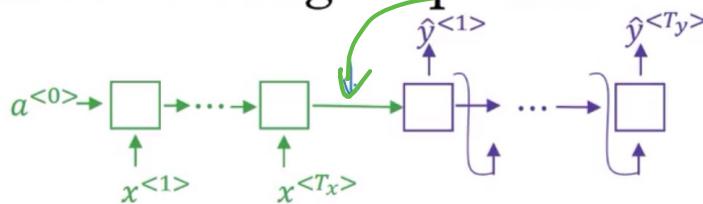
## Attention model intuition

Encoder - decoder architecture:

1 RNN reads in sentence, then 1 RNN outputs a sentence

A modification to this is called attention model

## The problem of long sequences

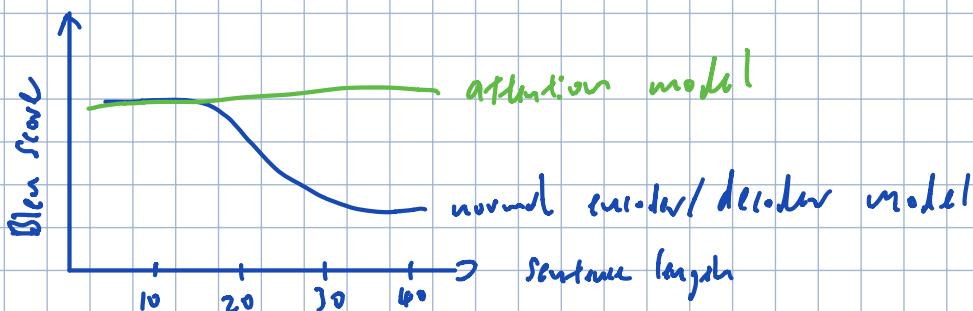


Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente d'y aller aussi.

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.

what we ask the (green) encoder in our network do is read in the whole sentence and then we move the whole sentence and store it in the aspiration here the decoding network then generate the English translation

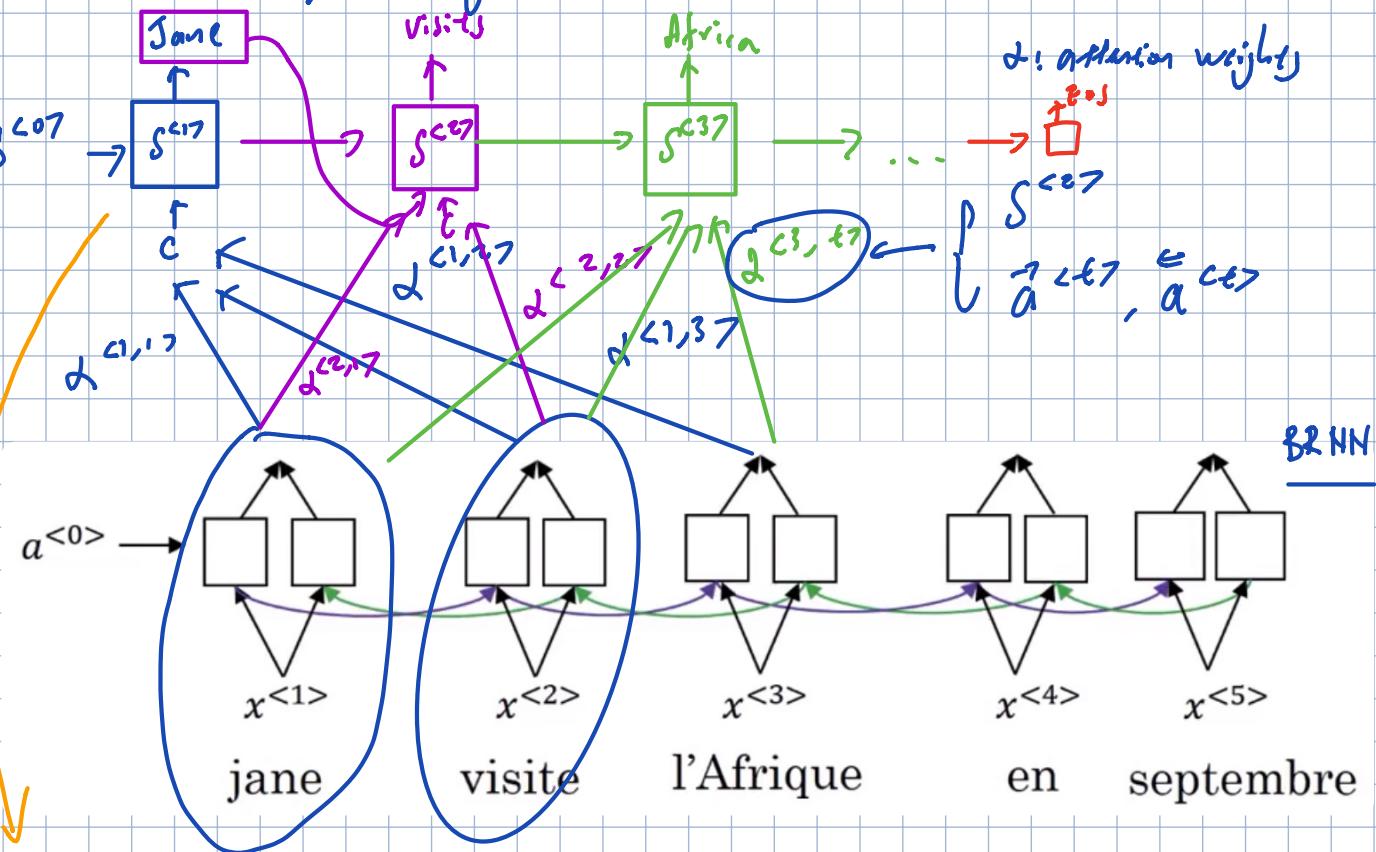
In human translation, we usually do it part by part through the sentence.



## Attention model intuition

use another RNN to generate English translations

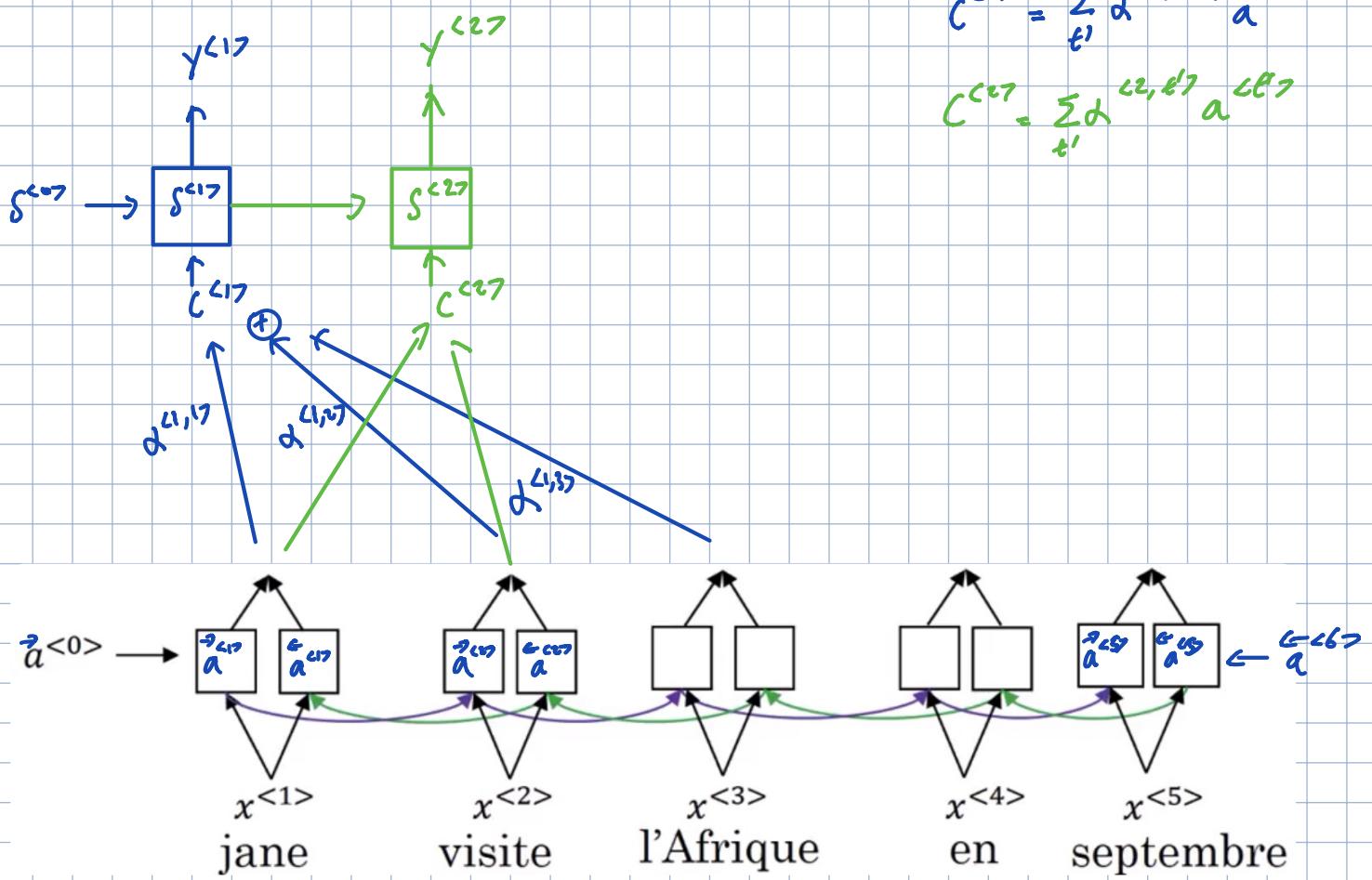
" $s^t$ " another notation for hidden unit here.



when translate the first English word, how much attention should we pay to the first word in the French sentence.

The RNN (above) at every step, there are these attention weights  $\alpha^{c,e,t}$ , when we are trying to generate  $t$ ' English word, how much should we pay attention to  $t'$  French word.

## Attention model



$$\sum_{t'} \alpha^{(1,t')} = 1$$

$$c^{<1>} = \sum_{t'} \alpha^{(1,t')} a^{<t'>}$$

$$c^{<2>} = \sum_{t'} \alpha^{(2,t')} a^{<t'>}$$

$$\alpha^{<t>} = (\vec{a}^{<t>}, \vec{e}^{<t>})$$

$\vec{e}$

both of the activations correlate together

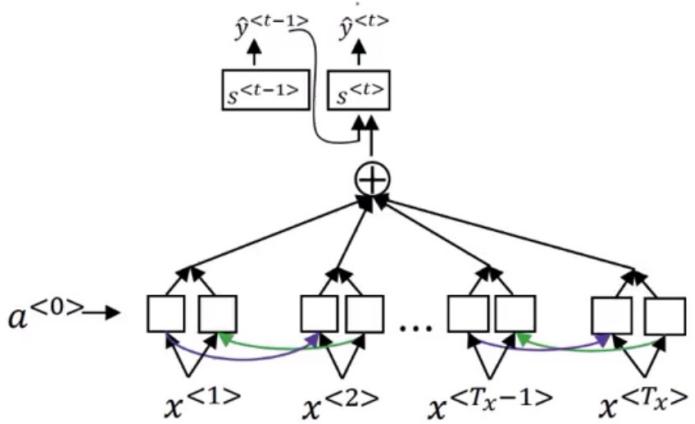
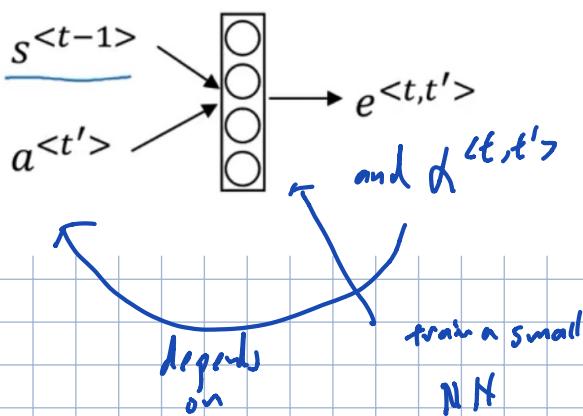
feature vector for time step  $t$

$\alpha^{(t,t')}$  = a way of attention that's  $y^{<t>}$  should pay to  $a^{<t'>}$

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

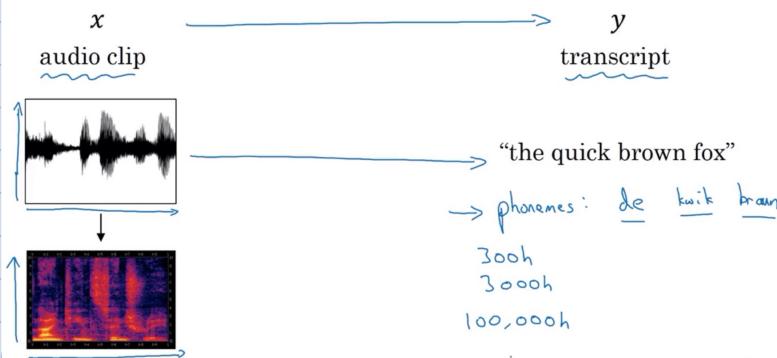
$\alpha^{<t,t'>} = \text{amount of attention } y^{<t>} \text{ should pay to } a^{<t'>}$

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

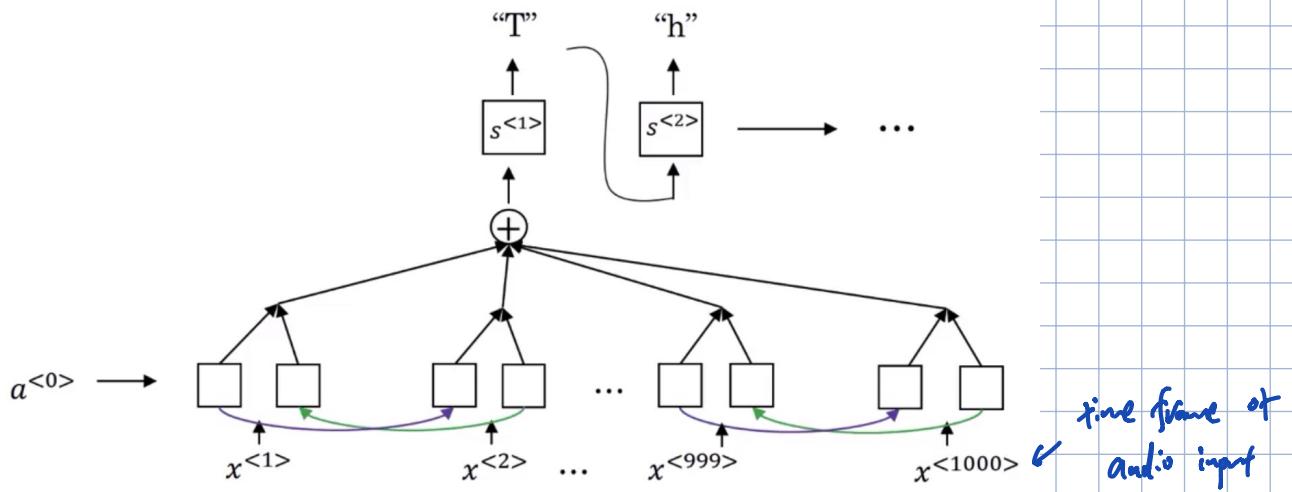


## Audio data: Speech Recognition

Speech recognition problem

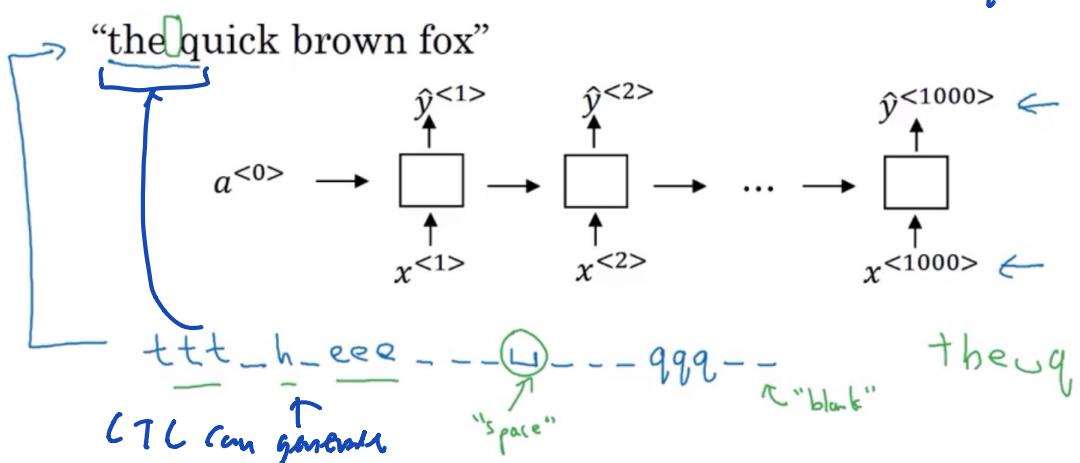


Attention model for speech recognition



# CTC cost for speech recognition

(Connectionist temporal classification)

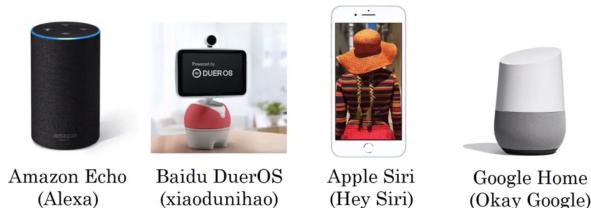


Basic rule: collapse repeated characters not separated by "blank" ↗

Graves et al., 2006. Connectionist Temporal Classification: Labeling unsegmented sequence data with recurrent neural networks] Andrew Ng

## Trigger word Detection

What is trigger word detection?



## Trigger word detection algorithm

