

# SEQUENCE MODELS

①

## Recurrent Neural Networks

### Why Sequence Models

#### Examples of Sequence Data

Speech Recognition     

"The quick brown fox jumped over the lazy dog"

→ Both input and output are sequence data

Music Generation     

$\emptyset$   
can be nothing or integer

→ Sequence data

Sentiment Classification

"There's nothing to like in this movie"

DNA Sequence Analysis

AGCCCTG... AGCCC...

Machine Translation

Valuez-vous chanter avec moi?

Do you want to sing with me?

Video Activity Recognition

Running

Name Entity Recognition

Yesterday, Harry Potter met Hermione Granger

Harry Potter  
Hermione Granger

### Notation

#### Motivating Example:

$X$ : Harry Potter and Hermione Granger invented a new spell.

$x^{(1)} \ x^{(2)} \ x^{(3)} \ x^{(4)} \ x^{(5)} \ x^{(t)} \dots \ x^{(9)}$

$y^{(1)} \ y^{(2)} \ y^{(3)} \dots \ y^{(9)}$

$X^{(i)}$  →  $i$ th training example

$T_x^{(i)}$  → length of input in example  $i$

$X^{(i)(t)}$  →  $i$ th training example  $t^{\text{th}}$  element

$T_y^{(i)}$  → " " output in example  $i$

#### Representing words:

##### Vocabulary

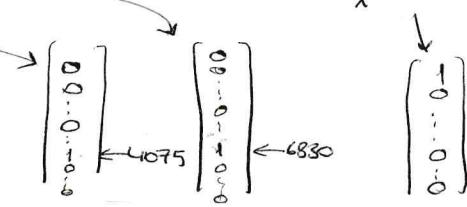
aeron	1
and	367
harry	4075
potter	6830
zulu	10000

$X$ : Harry Potter and Hermione Granger invented a new spell.

$x^{(1)} \ x^{(2)} \ x^{(3)}$

$x^{(7)} \ x^{(9)}$

one-hot representation

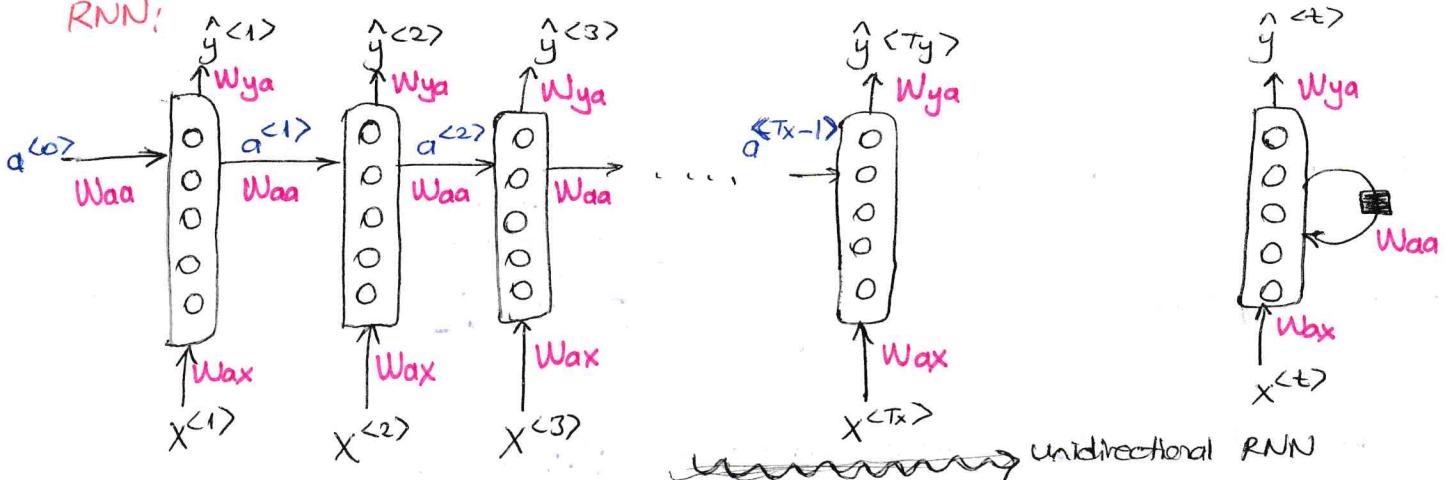


## Recurrent Neural Network Model

Why not a standard network?

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

RNN:



He said, "Teddy Roosevelt was a great President"

He said, "Teddy bears are on sale!"

} In order to detect the names, it necessitates bidirectional RNN

$$\begin{aligned}
 a^{<0>} &= \vec{0} \\
 a^{<1>} &= g(W_{aa} \cdot a^{<0>} + W_{ax} \cdot x^{<1>} + b_a) && \leftarrow \text{tanh/ReLU} \\
 \hat{y}^{<1>} &= g_2(W_{ya} \cdot a^{<1>} + b_y) && \leftarrow \text{sigmoid softmax} \\
 a^{<t>} &= g(W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>} + b_a) \\
 \hat{y}^{<t>} &= g_2(W_{ya} \cdot a^{<t>} + b_y)
 \end{aligned}$$

Simplified RNN:

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

$$\begin{array}{c}
 \boxed{W_{aa} \quad | \quad W_{ax}} \\
 \downarrow 100 \quad \quad \quad \downarrow 10000 \\
 \boxed{W_a} = W_a
 \end{array}$$

$(100, 10100)$

For example:

$$\begin{array}{c}
 a^{<t>} = g(W_{aa} a^{<t-1>} + W_{ax} x^{<t>} + b_a) \\
 \uparrow 100 \quad \quad \quad \uparrow 10000 \\
 (100, 100) \quad \quad \quad (100, 10000)
 \end{array}$$

$$\begin{array}{c}
 \boxed{[a^{<t-1>}, x^{<t>}]} = \boxed{\frac{a^{<t-1>}}{x^{<t>}}} \\
 \boxed{W_{aa} \quad | \quad W_{ax}} \boxed{[a^{<t-1>} \quad x^{<t>}]} = W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>}
 \end{array}$$

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

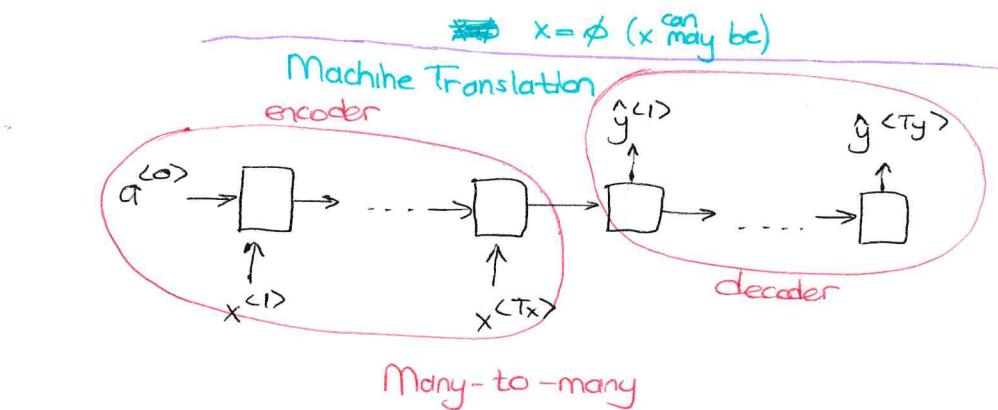
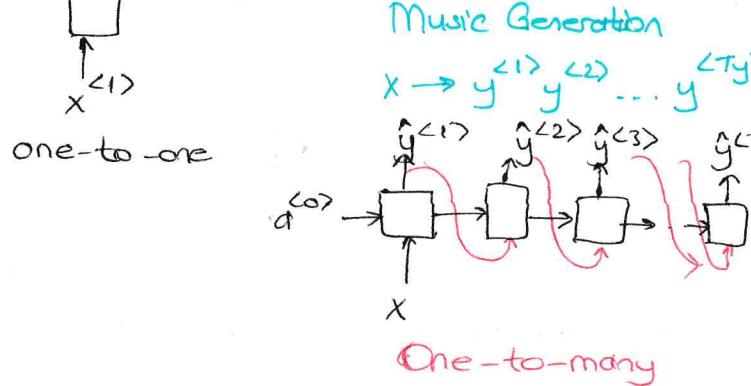
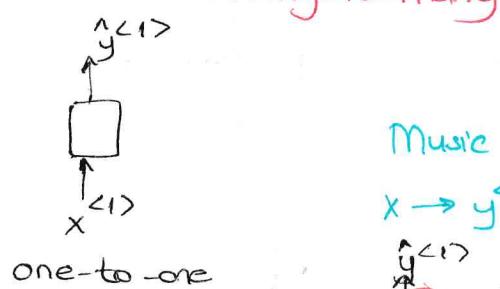
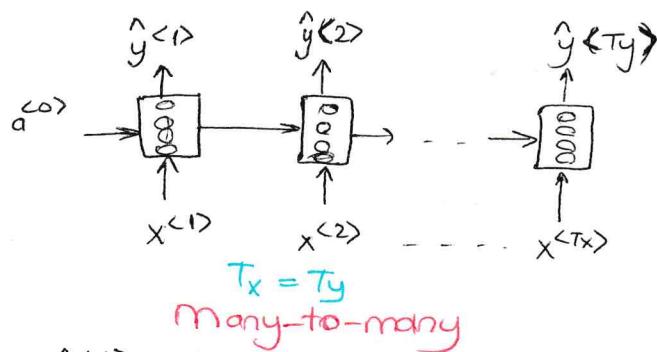
Backpropagation Through Time

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

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

Different Types of RNNs

The presentation in this video was inspired by a blog post by Andrej Karpathy, titled, The Unreasonable Effectiveness of Recurrent Neural Networks.



Language Model and Sequence Generation

What is language modelling?

Speech Recognition

$$P(\text{The apple and pair salad}) = 3.2 \times 10^{-13}$$

$$P(\text{The apple and pear salad}) = 5.7 \times 10^{-10}$$

$$P(\text{sentence}) = ? \quad P(y^{(1)}, y^{(2)}, \dots, y^{(T)})$$

Language modelling with an RNN

Training set: Large corpus of English text.

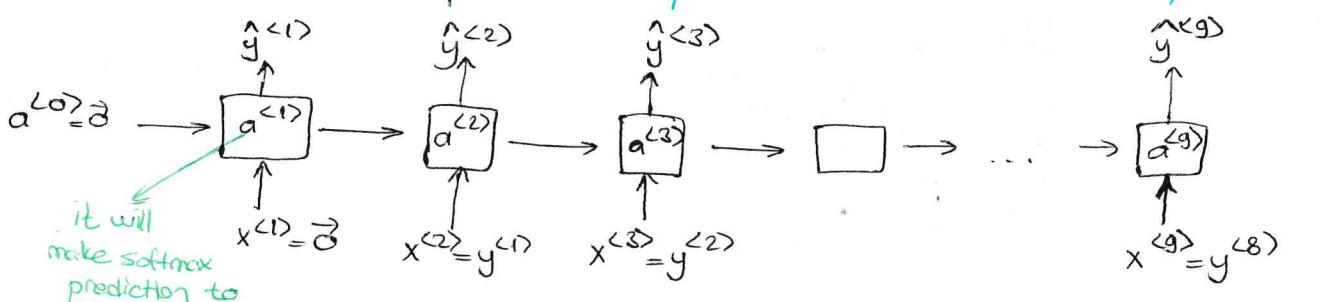
"Cats ~~avarage~~ 15 hours of sleep a day." <EOS>

$$y^{(1)} \quad y^{(2)} \quad y^{(3)} \quad \dots \quad y^{(8)} \quad y^{(9)}$$

"The Egyptian Mau is a breed of cat." <EOS>  
<UNK>

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

RNN model:



Cats avarage 15 hours of sleep a day. <EOS>

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

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

Sampling Novel Sequences

Word-level language model:

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

Character-level language model:

Vocabulary = [a, b, c, ..., z, Ȑ, ., , , , 0, ..., 9, A, ..., Z]

Pros:

- You don't ever have to worry about unknown word tokens,
- also used in more specialized applications where you have a more specialized vocabulary.

Cons:

- You end up with much longer sequences.
- Not as good as word level lang. models at capturing long range dependencies between how the earlier parts of the sentences also affect the later part of the sentence.
- more computationally expensive to train.
- They are not in widespread use today.

Vanishing Gradients with RNNs

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

The cats, which already ate ...., were full.

↓  
This is right or wrong, it's just very difficult for the area to backpropagate all the way to the beginning of the sequences and therefore to modify how the neural network is doing computations earlier in the sequence.

This is weakness of RNN algorithm.

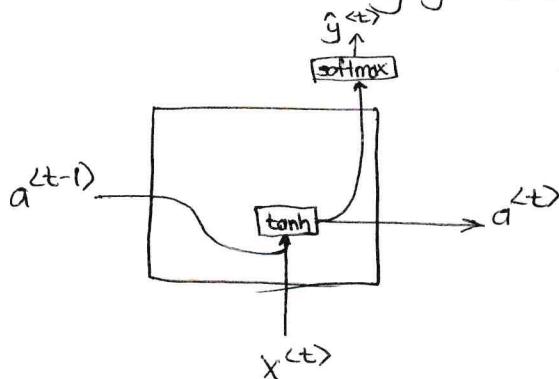
When exploding gradients happened, you see NaN's. Use gradient clipping.

robust solution

looks at gradient vectors, and if it is bigger than some threshold, rescale some of your gradient vector so that it is not too big.

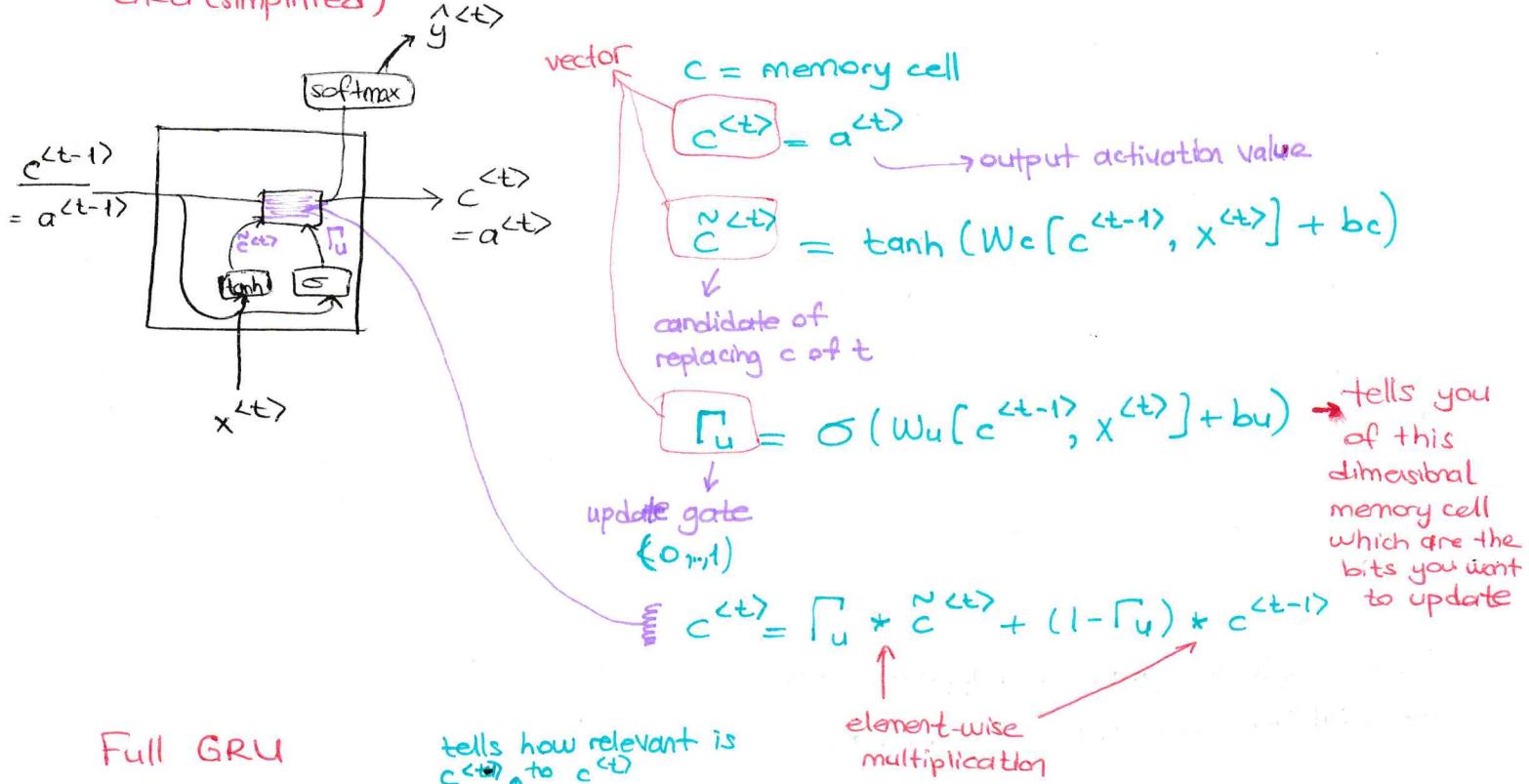
Gated Recurrent Unit (GRU)

It's a solution for vanishing gradient problems.



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

GRU (simplified)



Full GRU

tells how relevant is  $c^{t-1}$  to  $c^t$

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

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

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

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

↓

in literature

## Long Short Term Memory (LSTM),

COLLEGE 5  
WEEK 1

7

It is more powerful and more general version of GRU

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

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

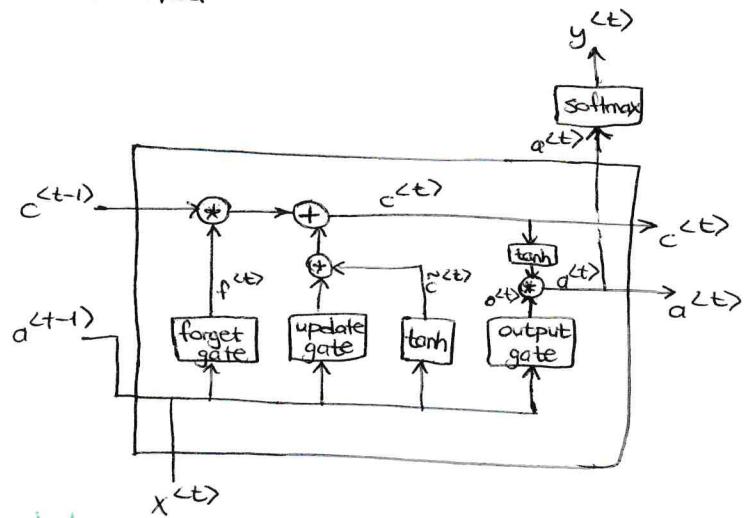
$$f_f = \sigma(W_f[a^{(t-1)}, x^{(t)}] + b_f)$$

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

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

$$a^{(t)} = f_o * \tanh(c^{(t)})$$

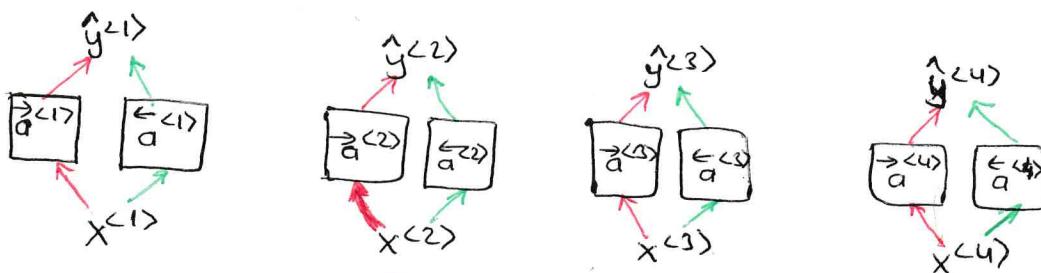
(peephole connection)



- LSTM came much earlier. GRUs were relatively recent invention
  - GRU is simpler model, it is actually easier to build or scale
  - GRU has 2 gates, so computationally, it runs a bit faster.
  - LSTM is more powerful and more effective since it has 3 gates instead of 2.
- Andrew NG thinks LSTM has been the historically more proven choice.

## Bidirectional RNN,

also can be  
GRU  
LSTM



$$\hat{y}^{(t)} = g(W_y[\vec{a}^{(t)}, \vec{a}^{\leftarrow(t)}] + b_y)$$

Disadvantage: You need the entire sequence of data before you can make predictions anywhere.

# Introduction to Word Embeddings

11

## Word Representation

COURSES  
WEEK 2

$$V = \{a, aaron, \dots, zulu, \text{UNK}\} \quad |V|=10\,000$$

1-hot representation

Man  
(5391)

$$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

05391

One hot

Feature representation ; word embedding

	Man (5391)	Woman (9853)	King (ugly)	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	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
;						
size						
cost						
alive						
noun						
verb						
;						
e <sub>5391</sub> (300,1)						
e <sub>9853</sub> (300,1)						

→ This allows it to generalize better across different words.

## Using Word Embeddings

COURSE 5  
WEEK 2



### Named Entity Recognition Example:

Sally Johnson is an orange farmer.

Robert Lih is an apple farmer.  
a durian cultivator.

~~Transfer Learning:~~

→ 1B words - 100B words

Bidirectional RNN

→ 100k words

### Transfer Learning and Word Embeddings

1. Learn 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)

→ 10000

→ 300

sparse vector      dense vector

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

In practice, you would do this only if task 2 has pretty big data set.

If your label data set for step 2 is quite small, then usually, I would not bother to continue to finetune the word embeddings.

Word embeddings has been useful for:

- named entity recognition
- text summarization
- co-reference resolution
- parsing

less useful for:

- language modeling
- machine translation

## Properties of Word Embeddings

	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

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

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

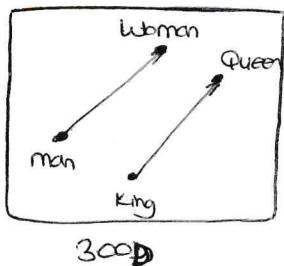
Man → Woman vs King → Queen

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

## Analogies Using Word Vectors

COURSE 5  
WEEK 2

3



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

Find word  $w$ :  $\arg \max \sim(e_w, e_{\text{King}} - e_{\text{man}} + e_{\text{woman}})$

Most commonly used similarity function is called Cosine Similarity

## Cosine Similarity

$$\sim(e_w, e_{\text{King}} - e_{\text{man}} + e_{\text{woman}})$$

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

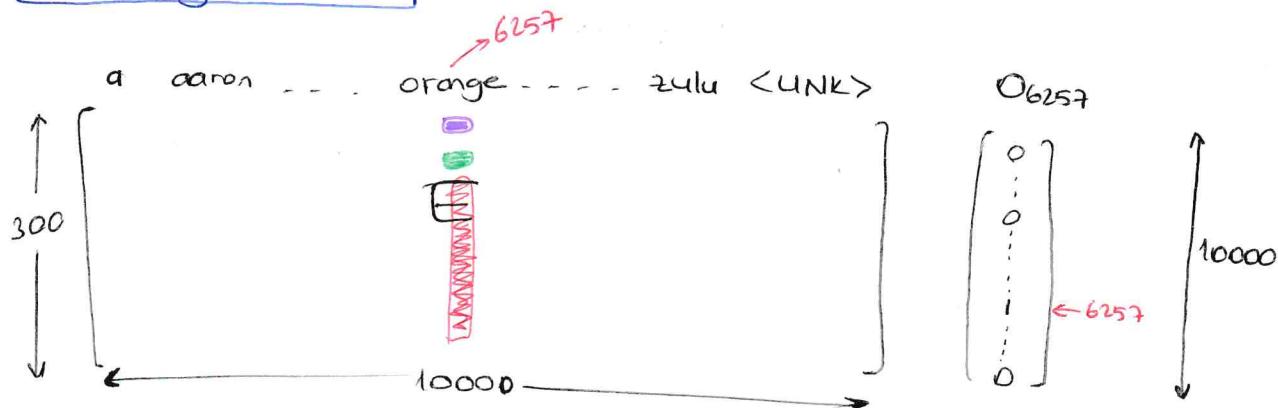
if  $u$  and  $v$  are very similar, their inner product will tend to be large.

Square Distance (Euclidean Distance):

$$\|u - v\|^2$$

Main difference between these is how it normalizes the lengths of the vectors  $u$  and  $v$ .

## Embedding Matrix



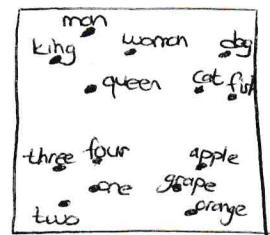
$$E \cdot O_{6257} = \begin{bmatrix} \text{row 1} \\ \text{row 2} \\ \vdots \\ \text{row 6257} \end{bmatrix} = e_{6257}$$

$(300, 10k) \quad (10k, 1)$

$$E \cdot O_J = e_J$$

$= \text{embedding for word } J$

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



t-SNE

→ non-linear mapping

(It takes 300-D data and maps it in a very non-linear way to a 2D space.)

# Learning Word Embeddings: Word2vec & GloVe

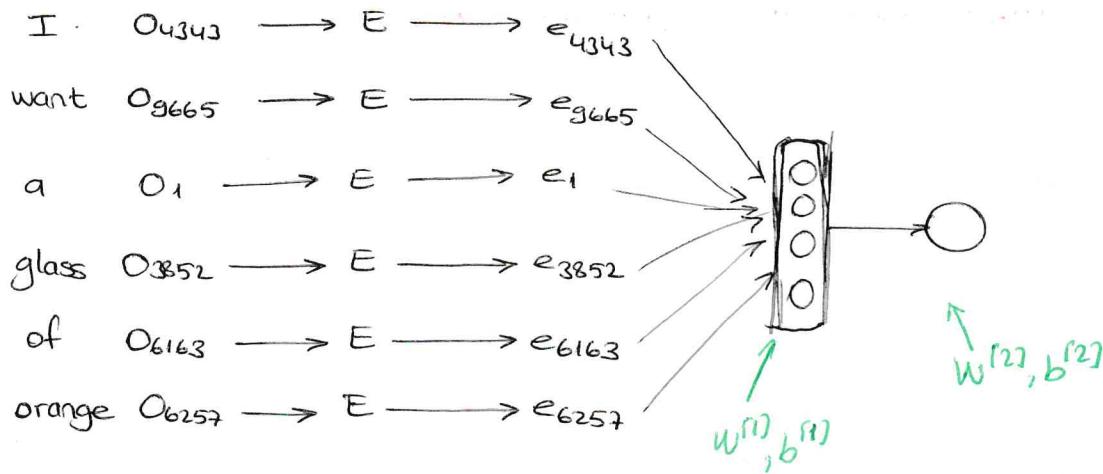
34

## Learning Word Embeddings

COURSE 5  
WEEK 2

Neural Language Model :

I want a glass of orange —  
4343 9665 1 3852 6163 6257



if we use 6 words, dimensions will be 1800 with their embeddings.

If 4 words are used, 1200 dimensions will be.

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      a glass of orange ? to go along with

Last 1 word      orange ?

Nearby 1 word → idea of a Skip-Gram model      glass ?

## Word2Vec

Skip-grams

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

<u>Context</u>	<u>Target</u>	→ (within ± 10 word window)
randomly ← pick a word	orange	juice
	orange	glass
	orange	my

Mikolov et al., 2013,  
Efficient estimation of word representations in vector space

⇒ supervised learning

# Model

COURSE 5  
WEEK 2

5

Vocab size = 10 000k

$x$  —————→  $y$   
 Context  $c$  ("orange") —————→ Target  $t$  ("juice")  
 6257 4834

$O_c \rightarrow E \rightarrow e_c \rightarrow O \rightarrow \hat{y}$   
 $e_c = E \cdot O_c$  softmax

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 a output  $t$

$$L(\hat{y}, y) = -\sum_{i=1}^{10000} y_i \log \hat{y}_i$$

$$y = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \leftarrow 4834$$

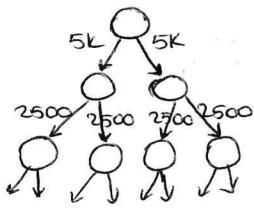
Problems with softmax Classification

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

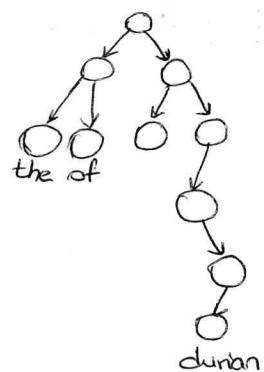
The primary problem is computational speed.

Solution: Hierarchical Softmax Classifier

computational size =  $(\log V)$



In practice:



How to sample the context  $c$ ?

→ the, of, a, and, to, ...

orange, apple, durian

$P(c)$



In practice, the distribution of words  $p(c)$  isn't taken just entirely random for the training set purpose, instead, there are different heuristics in order to balance out something from the common words together with less common words.

## Negative Sampling

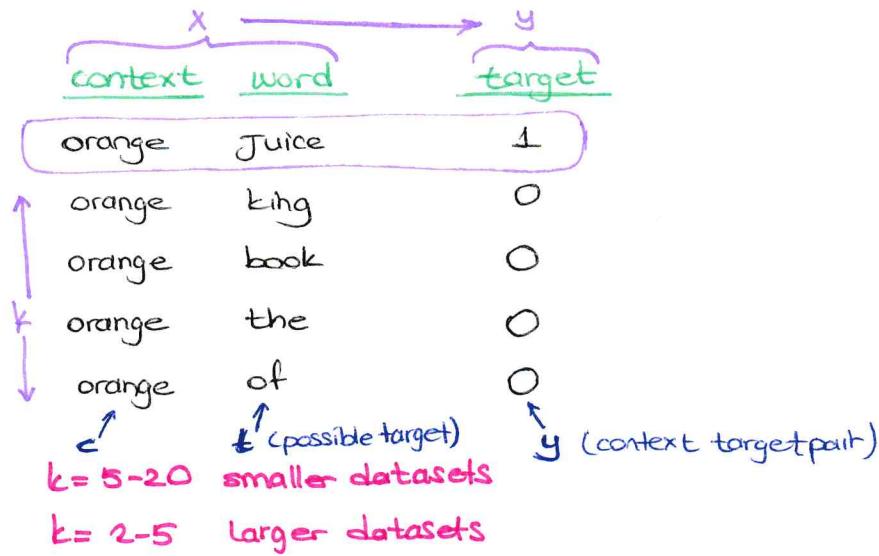
COURSE 5  
WEEK 2

E6

Defining a new Learning problem :

Mikolov et al. 2013, Distributed representation of words and phrases and their compositionality

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

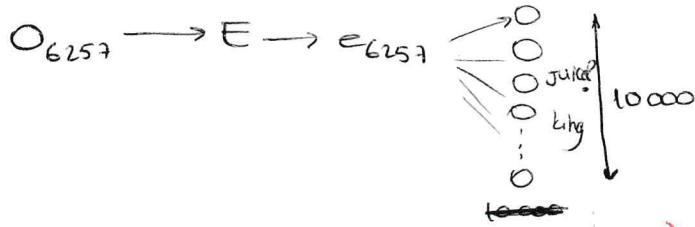


$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}} \rightarrow 10000 \text{ way softmax}$$

parameter vector  $\theta$  for each possible target word

$$p(y=1 | c, t) = \delta(\theta_t^T e_c)$$

separate parameter vector embedding vector for each possible context word



We are going to train one responding to the actual target word and then  $k$  randomly chosen negative examples.

10000 binary classification problem

Selecting Negative Examples:

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10000} f(w_j)^{3/4}}$$

## GloVe Word

COURSE 5  
WEEK 2

7.

Pennington et al., 2014,  
GloVe : Global vectors  
for word representation

# Various Sequence to Sequence Architectures

13

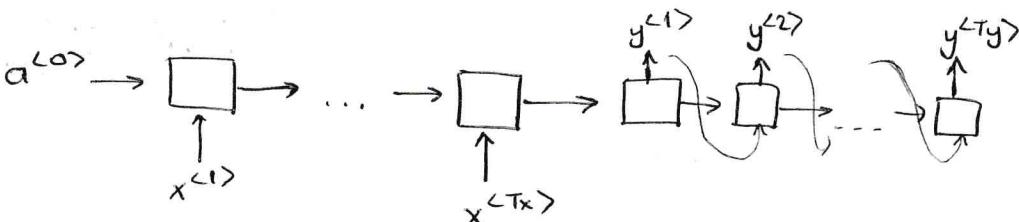
## Sequence To Sequence Model

### Basic Models

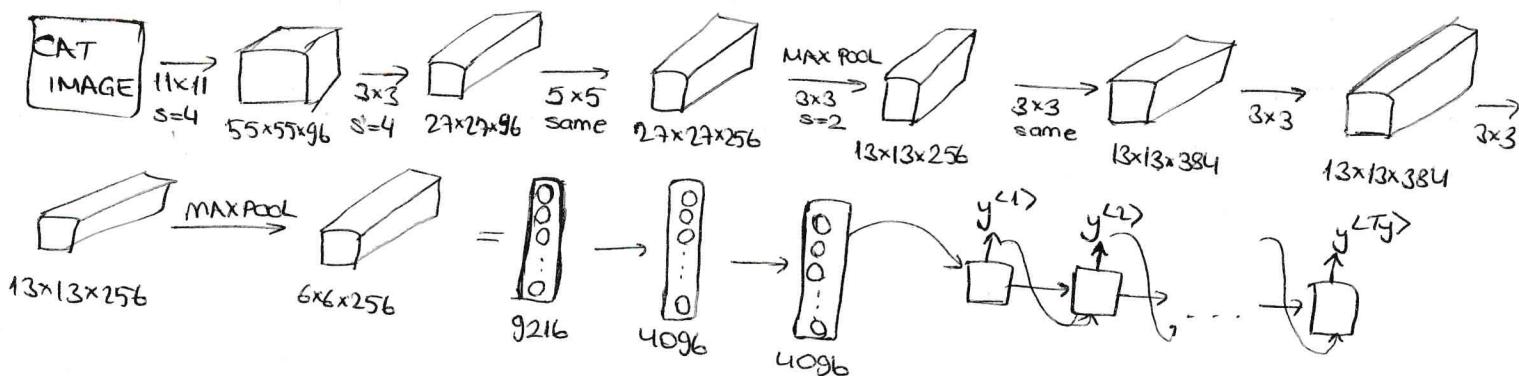
$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$   
 Jane visite l'Afrique en septembre.

→ Jane is visiting Africa in September.

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



## Image Captioning:



## Picking The Most Likely Sentence

### Finding the most likely translation

Jane visite l'Afrique en septembre.

→ 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.

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

$$\underset{y^{<1>} \dots y^{<Ty>}}{\operatorname{argmax}} P(y^{<1>} | x, \dots, y^{<Ty>} | x)$$

## Greedy Search :

$$P(\hat{y}^{<1>} | x, \hat{y}^{<2>} | x, \dots, \hat{y}^{<Ty>} | x)$$

↓  
pick-the  
best first  
word

After having  
the best first word,  
try to pick the best  
second word

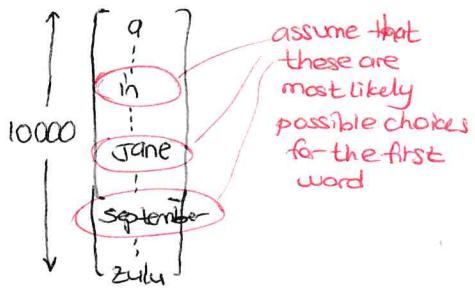
[ Sutskever et.al. 2014  
Sequence to sequence  
learning with neural  
networks.]

[ Cho et.al. 2014  
Learning phrase  
representation Using  
RNN encoder-decoder  
for statistical machine  
translation ]

## Beam Search

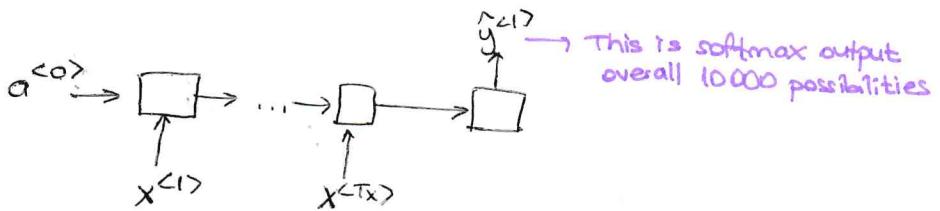
$$B = 3 \text{ (beam width)}$$

**Step 1 :**

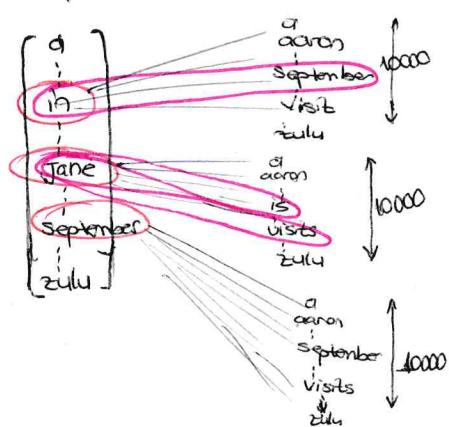


↳ means beam search will consider 3 possibility at a time

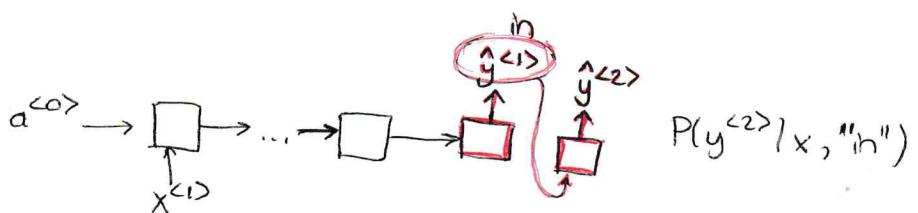
$$P(y^{<1>} | x)$$



**Step 2 :**

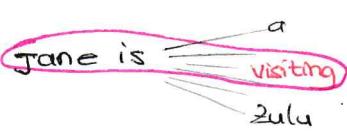
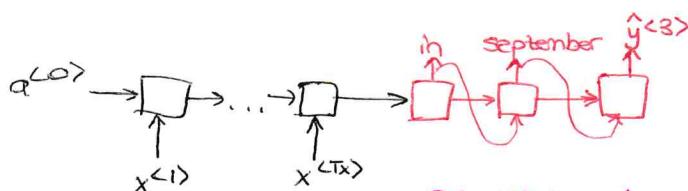


**Step 2 :**



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

**Step 3 :**



$$P(y^{<3>} | x, "in september")$$



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

→

The outcome of this process hopefully will be :

Jane visits Africa in September.  $\langle \text{EOS} \rangle$

$B=1 \rightarrow$  greedy search

$B=3 \dots 10 \rightarrow$  beam search will find a much better output sentence than greedy search.

## Refinements to Beam Search

Beam Search is a widely used algorithm in many production systems or in many commercial systems.

### Length Normalization:

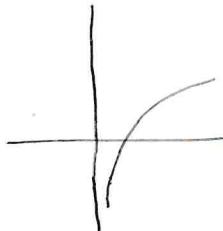
$$\arg \max_y \prod_{t=1}^{Ty} P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

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

$$P(y^{<1>} \dots y^{<Ty>} | x) = P(y^{<1>} | x) \cdot P(y^{<2>} | x, y^{<1>}) \dots P(y^{<Ty>} | x, y^{<1>} \dots y^{<Ty-1>})$$

it will be very tiny result.  
So we take log of this products

If you have long sentence, the probability of that sentence is going to be low, because you're multiplying as many terms here.



$\log P(y | x)$  should give the same result as maximizing  $P(y)$  given  $x$ .  
( $P(y | x)$ )

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

This significantly reduces the penalty for outputting longer translations

In practice:

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

$\alpha = 0.7$     $\alpha = 1 \rightarrow$  completely normalizing by length  
hyperparameter

$\alpha = 0 \rightarrow$  not normalizing at all.

### Beam Search Discussion:

Beam width  $B$ ? 10 100 1000 3000 (depends on the application)

The larger  $B$  is the more possibilities so you tend to get a better result.  
the more computationally expensive your algorithm is. (Because of keeping all possibilities around)  
It will be slow.

The memory requirements will grow, will also be computationally slower.

The small  $B$  is opposite of larger  $B$ .

Large  $B$ : better results, slower  
Small  $B$ : worse result, faster

So for many applications, from beam width 1 to 3 to 10,

there is huge gain as you go. But the gains as you go from 1000 to 3000 in beam width might not be as big,

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

## Error Analysis in Beam Search

4

Example:

Jane visite l'Afrique en Septembre.

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

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

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

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

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

Conclusion: Beam search is at fault.

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

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

Conclusion: RNN 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 <del>last September</del> <sup>model</sup>	$2 \times 10^{-10}$	$1 \times 10^{-10}$	B
-	-	-	-	R
-	-	-	-	B
-	-	-	-	R

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

RNN

Beam Search

always  
tempting to  
collect more  
training data  
that never hurts

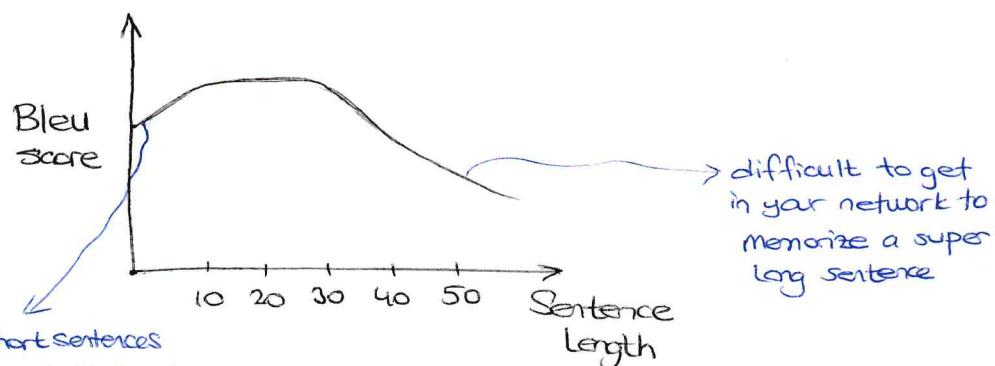
The way you  
got  $\hat{y}$  was  
you had RNN  
that was  
computing

$P(y|x)$  and  
beam search's  
job was to  
try to find a  
value of  
 $\text{argmax}_y P(y|x)$

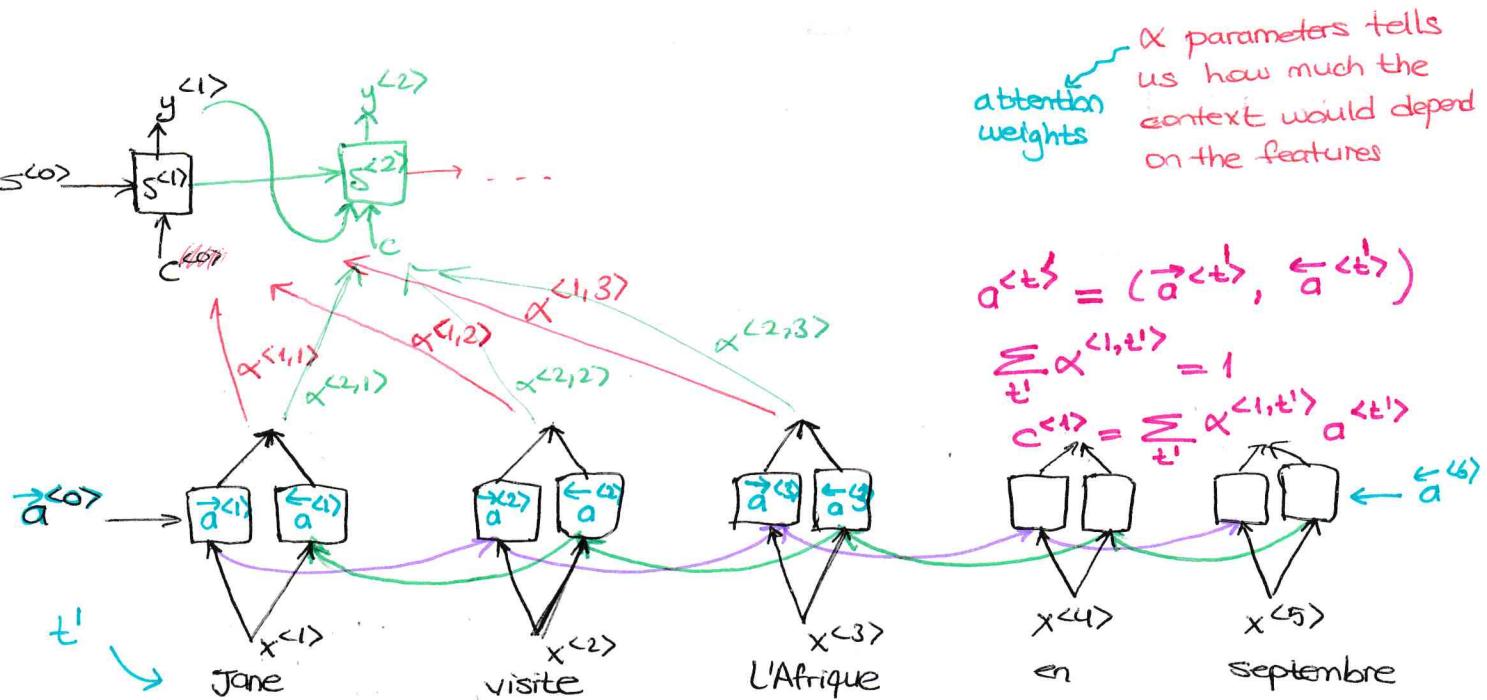
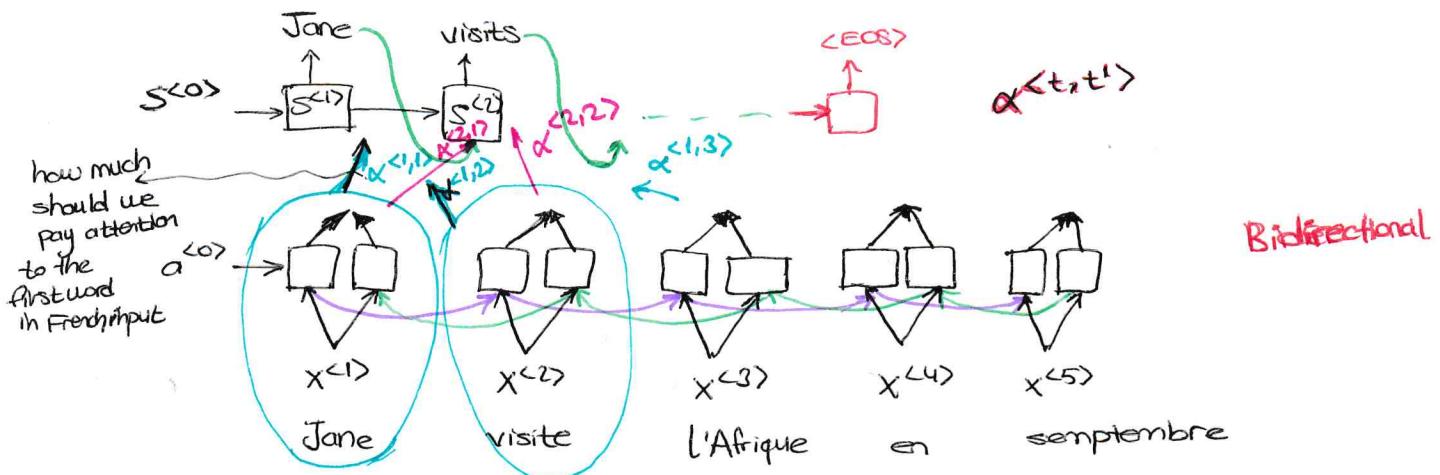
## ~~Attention~~ Attention Model

53

### The Problem of Long Sequences:



short sentences  
are just hard  
to translate,  
hard to get all  
the words



$$\alpha^{(t,t')} = (\vec{\alpha}^{(t,t')}, \overleftarrow{\alpha}^{(t,t')})$$

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

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

$$\alpha^{(t,t')} = \frac{1}{\sqrt{d}} \text{softmax}(\text{score}(x^{(t)}, h^{(t)}))$$

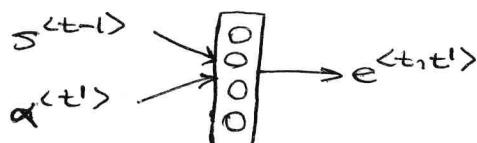
$\alpha^{(t,t')}$  = amount of "attention"  
 $y^{(t)}$  should pay to  $a^{(t')}$

$\alpha^{(t,t')}$  = amount of "attention"  
 $y^{(t)}$  should pay to  $a^{(t')}$

## Computing Attention $a^{(t+1)}$ :

$a^{(t+1)}$  = amount of attention  $y^{(t)}$  should pay to  $a^{(t)}$

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

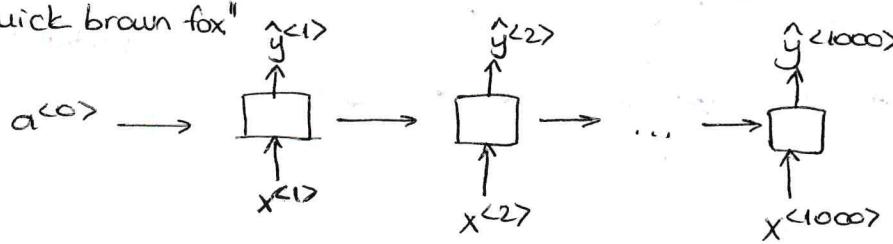


## Speech Recognition

CTS cost for speech recognition!

(Connectionist temporal classification)

"the quick brown fox" ^<1>



in practice,  
this network will  
usually be a bidirectional  
LSTM ~~or~~ bidirectional  
GRU and usually a deeper  
model.

ttt - h - eee - - ~~tt~~ - - - qqq - - "space" "blank"

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

Production skills speech recognition system is a pretty significant effort and requires a very large dataset.

