

$x = \text{Harry Potter and Hermione Granger invented a new spell.}$

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

$\rightarrow y: \quad 1 \quad 1 \quad 0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0$

$$T_x = 9$$

$T_x \rightarrow \text{length of the input sequence}$

$T_y \rightarrow \text{length of the output sequence}$

$x^{(i)} \rightarrow i^{\text{th}} \text{ training example}$

$x^{(i)<t>} \rightarrow t^{\text{th}} \text{ element in the sequence of training example } i$

$T_x^{(i)} \rightarrow \text{input sequence length for training example } i$

$y^{(i)<t>} \rightarrow t^{\text{th}} \text{ element in the output sequence of the } i^{\text{th}} \text{ training example}$

$T_y^{(i)} \rightarrow \text{length of the output sequence in the } i^{\text{th}} \text{ training example}$

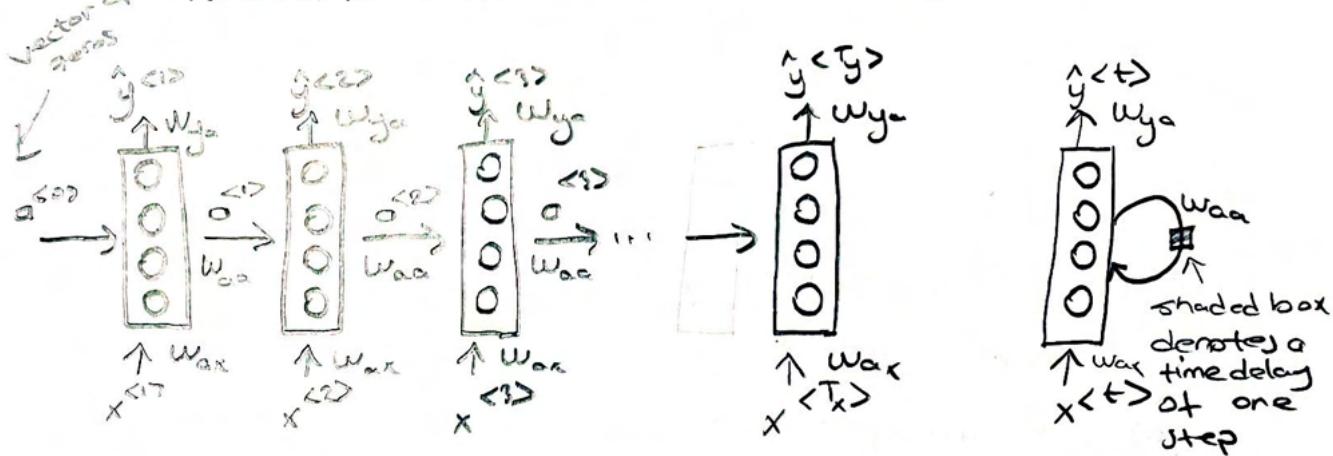
Vocabulary/
Dictionary

	$x^{<1>}$	$x^{<2>}$	$x^{<3>}$	one-hot encoding
aaron	1	[0]	[0]	
:	2	[0]	[0]	
and	367	[0]	[0]	
:	367	[0]	[0]	
harry	4075	[0]	[0]	
:	4075	[0]	[0]	
potter	8880	[0]	[0]	
:	8880	[0]	[0]	
Aulu	10000	[0]	[0]	

10000 D Vector

The goal is given this representation for X to learn a mapping using a sequence model to then target output y

RECURRENT NEURAL NETWORKS



Why not a standard network?

- Inputs, outputs can be different length in different examples
- A naive NN architecture doesn't share features learned across different positions of texts.

The RNN scans through the data from left to right. The parameters it uses for each time step are shared

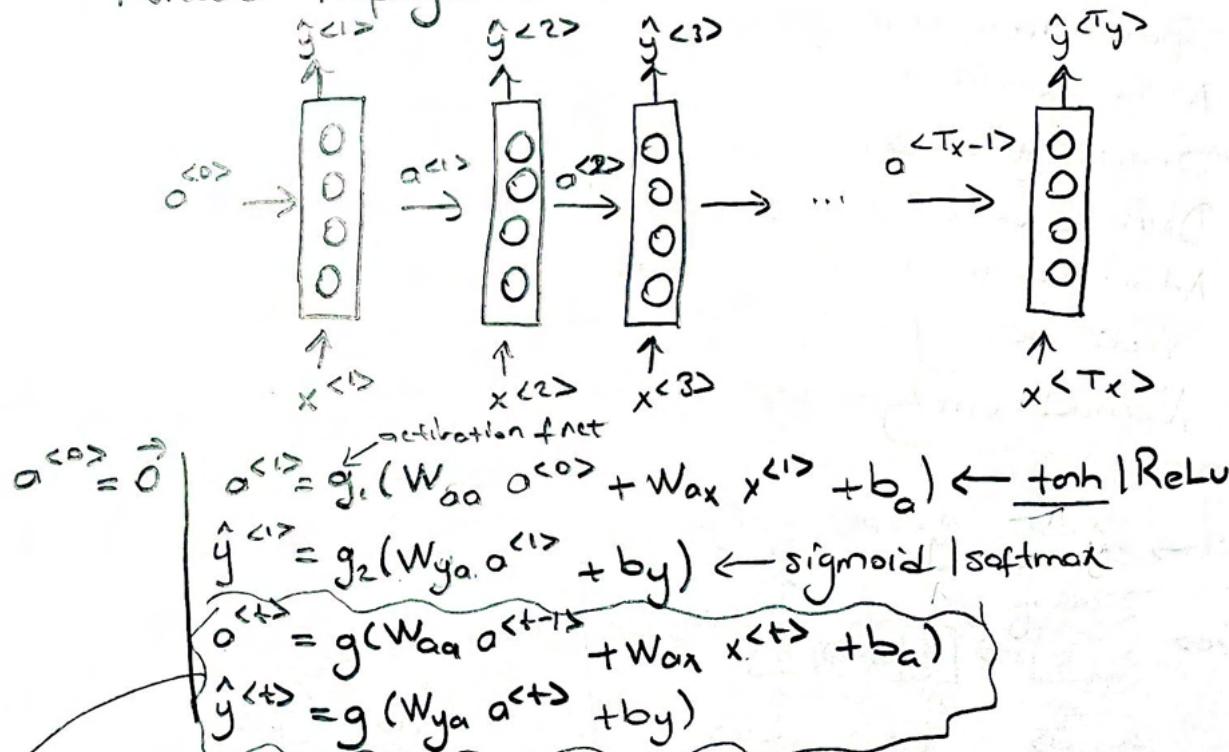
When making the prediction for $\hat{y}^{(3)}$, it gets the information not only from $x^{(3)}$ but also the information from $x^{(1)}$ and $x^{(2)}$ because the information on $x^{(1)}$ can pass through this way to help to prediction with $\hat{y}^{(3)}$

One weakness of this RNN is that it only used the info that is earlier in this sequence to make a prediction. When predicting $\hat{y}^{(3)}$, it doesn't use info about the words $x^{(4)}, x^{(5)}, x^{(6)}$, and so on. If you're given a sentence "He said, 'Teddy Roosevelt was a great President'" in order to decide whether or not the word "Teddy" is a part of a person's name it would be really useful to know not just info from the first two words but to know info from the later words in the sentence as well. The sentence could also have been "He said "Teddy bears are on sale!"

Bidirectional RNN (BRNN)

Unidirectional RNN

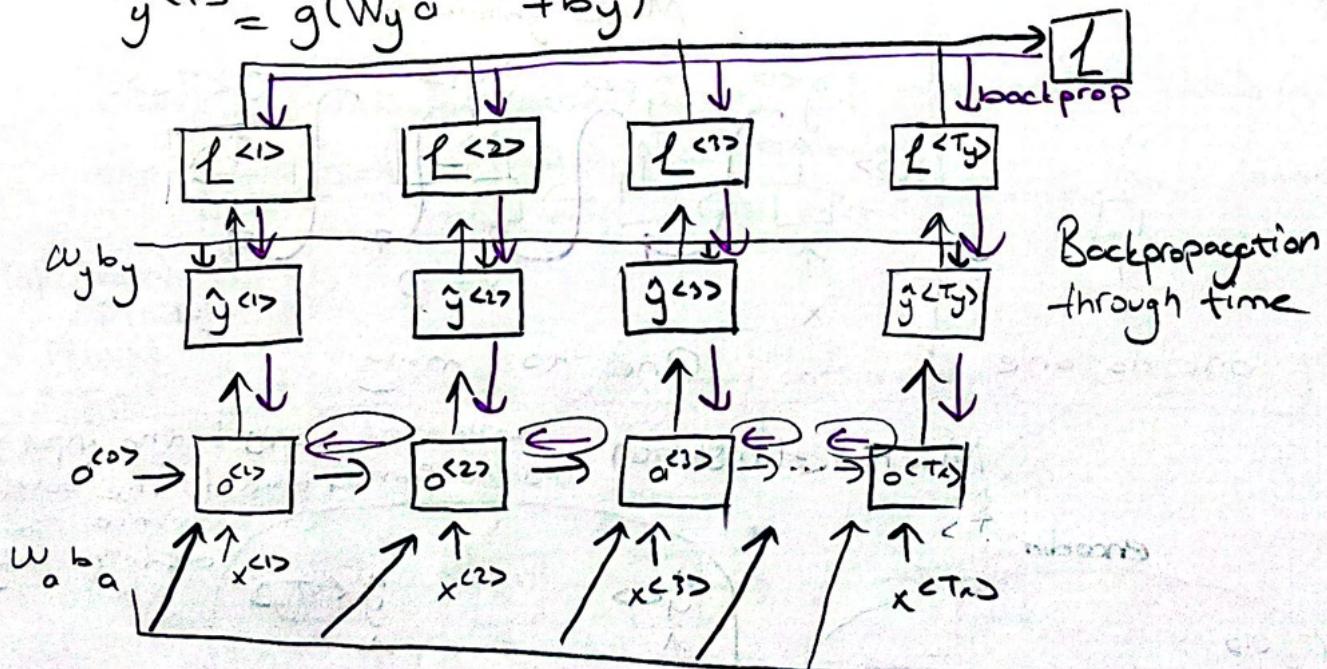
Forward Propagation



tanh prevents vanishing gradient

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

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



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

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

Examples of sequence data

Speech recognition

Music generation

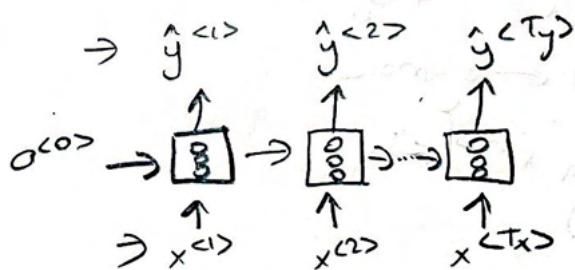
Sentiment classification

DNA sequence analysis

Machine translation

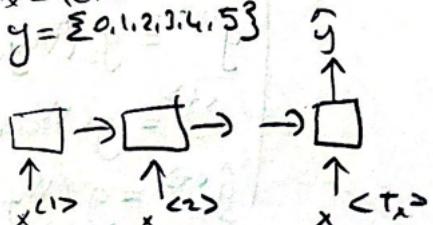
Video activity recognition

Named entity recognition

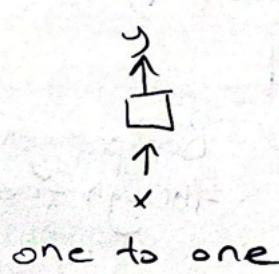


Many to many
(many inputs many outputs)

Sentiment classification
 $x = \text{text}$
 $y = \{0, 1, 2, 3, 4, 5\}$

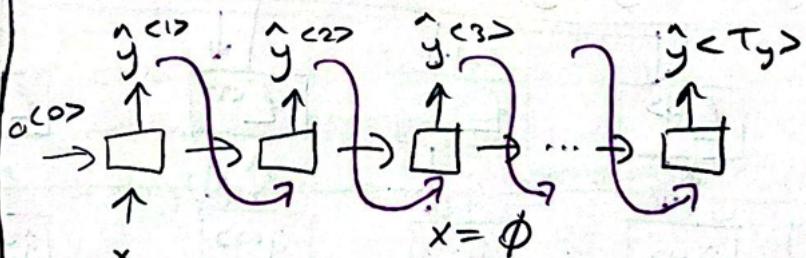


There is ... movie
Many to one



one to one

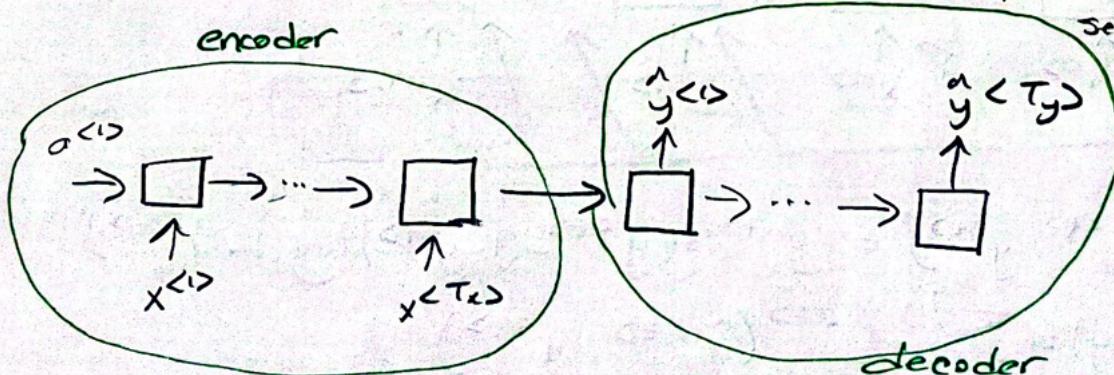
Music generation



One to many

Machine translation

Number of words in the input sentence
and the number of words in the output
sentence could be
different lengths



Many to many

LANGUAGE MODELING

5

Speech Recognition

The apple and pair salad.

The apple and pear salad.

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

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

$$P(y^{<1>} y^{<2>} \dots y^{<T>})$$

Language Modeling With an RNN

Training set: large corpus of English text.

* Cats average 15 hours of sleep a day. <EOS>

$$\begin{array}{ccccccc} -y^{<1>} & y^{<2>} & y^{<3>} & \dots & & y^{<8>} & y^{<9>} \\ x^{<t>} = y^{<t-1>} & & & & & & \end{array}$$

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

<UNK> stands for unknown words. If some of the words in the training set are not in the vocabulary

1,0002 way softmax output (Vocabulary length = 10000)

$$P(a) P(aaron) \dots P(cats) \rightarrow$$

$$\dots P(gulu)$$

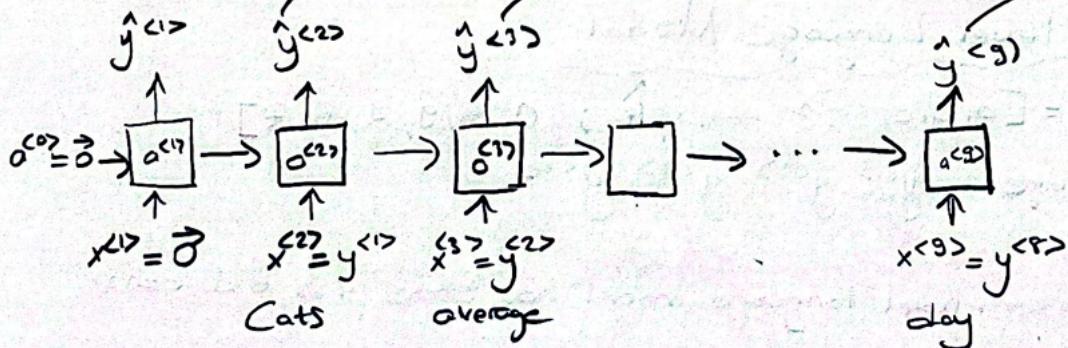
$$P(UNK)$$

$$P(EOS)$$

$$P(\text{average|cats})$$

$$P(__ | \text{"cats average"})$$

$$P(\underline{\text{EOS}} | \dots)$$



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

What $a^{<1>}$ does? It will make a softmax prediction to try to figure out what is the probability of the first word y .

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

\uparrow \uparrow
 NH true word
 softmax predicted

$$L = \sum_t L_t$$

$$L = \sum_t L^{<\leftrightarrow>} (y^{<\leftrightarrow>}, y'^{<\leftrightarrow>})$$

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

$$= p(y^{<1>}) p(y^{<2>} | y^{<1>})$$

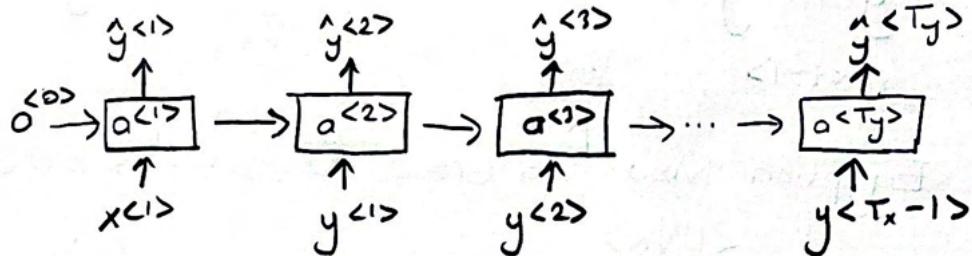
What's the chance of $y^{(2)}$ given $y^{(1)}$

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

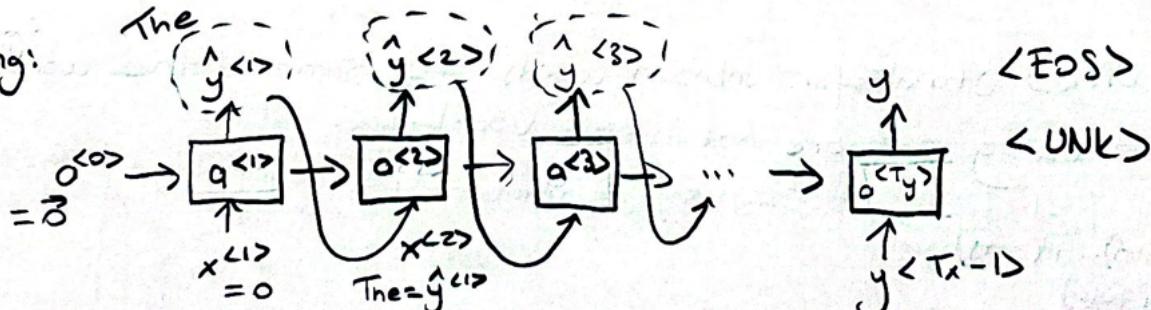
→ What's the chance of $y^{(s)}$ given $y^{(1)}$ and $y^{(2)}$

Sampling A Sequence from A Trained RNN

Training!



Sampling:



$P(a) P(aaron) \dots P(zulu) P(<\text{UNK}>)$ np.random.choice $P(__ | \text{the})$

Character-level Language Model

Vocabulary = [a, b, c, ..., f, —, ., , ;, 0, ..., 9, A, ..., F]

cat average. $y^{(1)}$ $y^{(2)}$ $y^{(3)}$ $y^{(4)}$

Using character-level language model has some pros and cons.

Character-level language model;

+ no `<UNK>` tokens

- The main disadvantage is much longer sequences.

- Computationally expensive to train

Vanishing gradient \rightarrow GRU, LSTM

Exploding gradient \rightarrow Gradient clipping can be used

Gated Recurrent Unit (GRU)

RNN Unit

$y \leftrightarrow$

Softmax

a^{t-1}

$x \leftrightarrow$

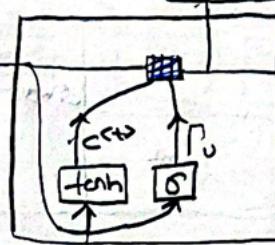
$a^t \leftrightarrow$

tanh

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

GRU

$$C^{t-1} = o^{t-1}$$



$$C^{t-1} = o^{t-1}$$

$$C^t = a^{t-1}$$

$$a^t = C^t$$

$$x \leftrightarrow$$

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

$$\Gamma_u = 0.000001$$

Γ : gate

\tilde{C} : candidate value it doesn't suffer from vanishing gradient problem

$$\Gamma_u = 1 \quad \Gamma_u = 0 \quad \Gamma_u = 0 \quad \Gamma_u = 0 \dots$$

$$C^{t-1} = 1 \quad \dots \quad \dots \quad \dots = 1$$

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

Singular or plural
↓
↓
0

C^{t-1} can be a vector

(8)

Full GRU

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

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

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

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

Γ_r : tells you how relevant is $c^{<t-1>}$ to computing the next candidate for $c^{<t>}$.

GRU

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

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

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

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

(update)

(forget)

(output)

LSTM

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

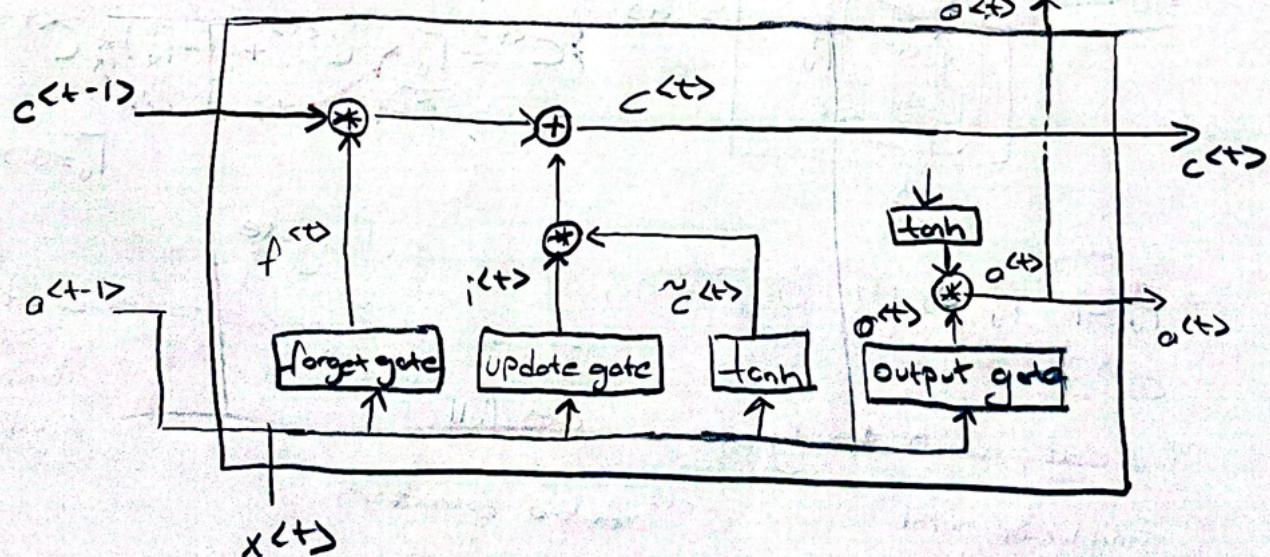
$$\Gamma_u = \sigma(W_u [\sigma^{<t-1>} x^{<t>}] + b_u)$$

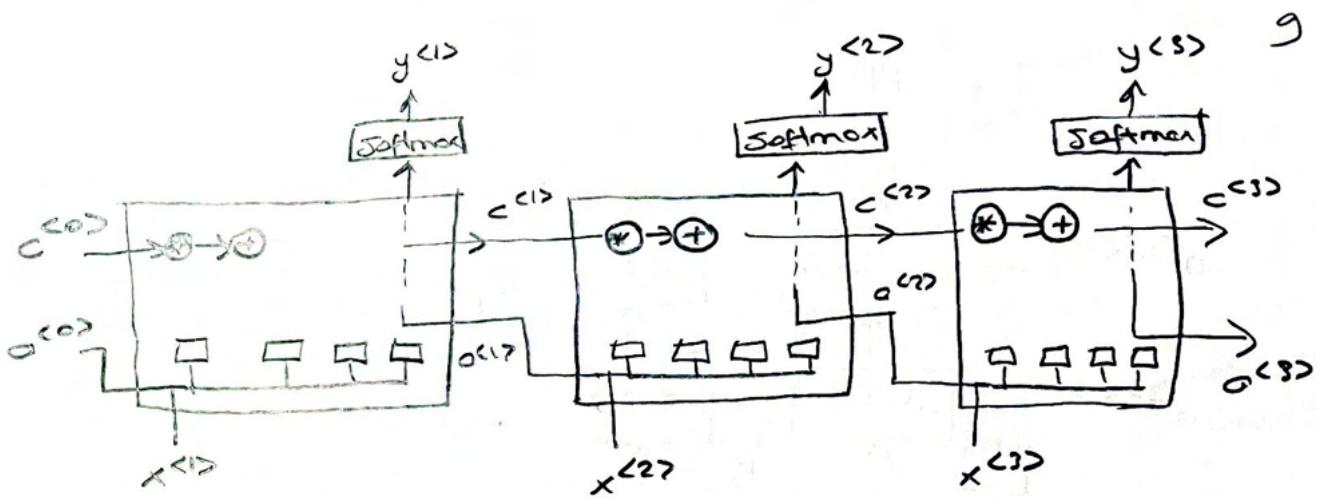
$$\Gamma_f = \sigma(W_f [\sigma^{<t-1>} x^{<t>}] + b_f)$$

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

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

$$o^{<t>} = \Gamma_o * \tanh(c^{<t>}) \quad \text{gate connection}$$

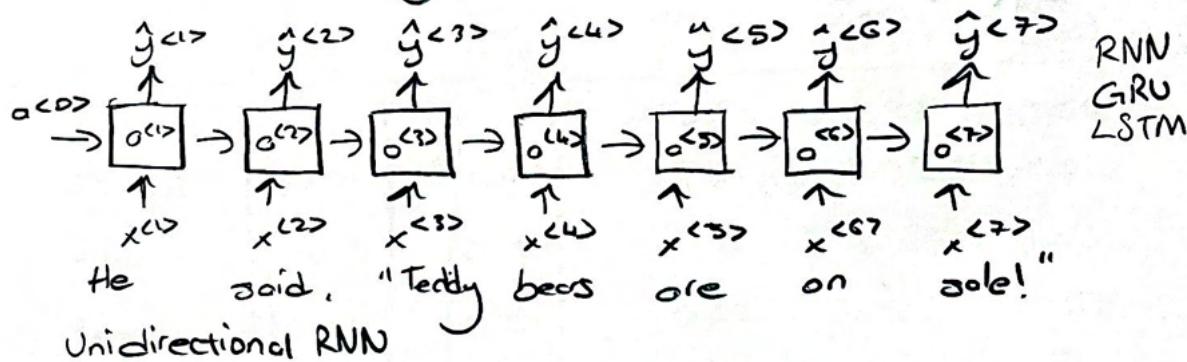




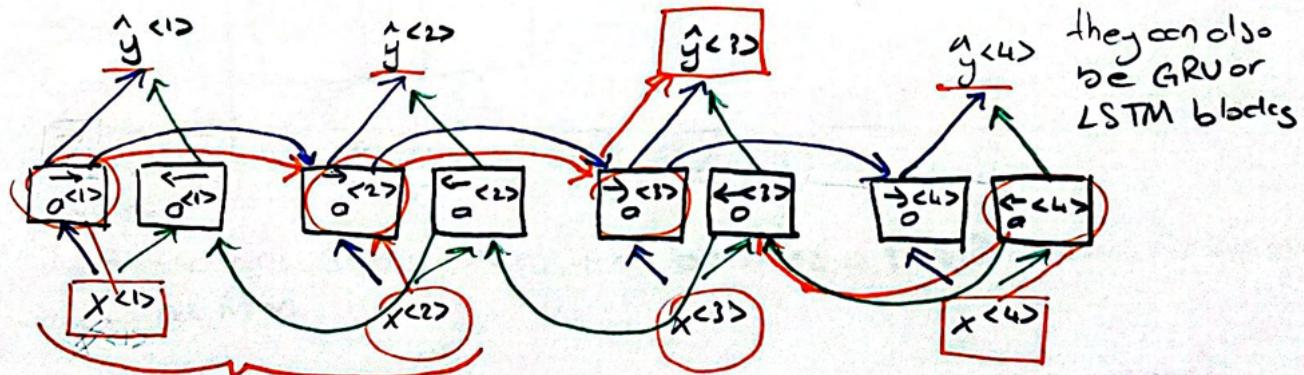
Bidirectional RNNs

He said, "Teddy bears are on sale!"

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



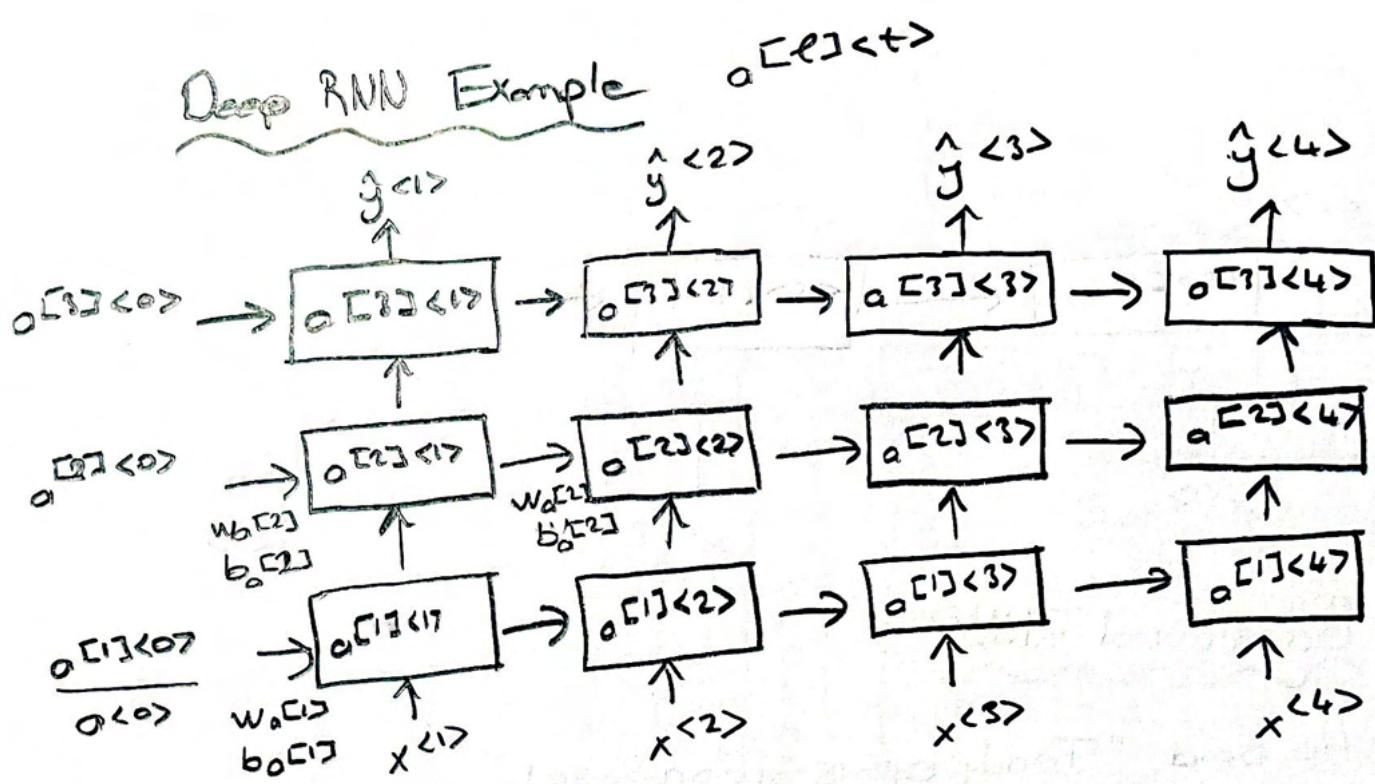
Unidirectional RNN



Acyclic graph (bipartite, no graph cycles)

$$y^{<t>} = g(W_y [o^{<t>} \rightarrow, o^{<t>} \leftarrow] + b_y)$$

Bidirectional LSTMs are commonly used in NLP



$$o^{[2]}_{[3]} = g(w_o^{[2]} [o^{[2]}_{[2]}, o^{[1]}_{[3]}] + b_o^{[2]})$$

Word Representation

$$V = [v_1, v_2, \dots, v_{100}, \text{UNK}] \quad |V| = 10,000$$

one-hot representation

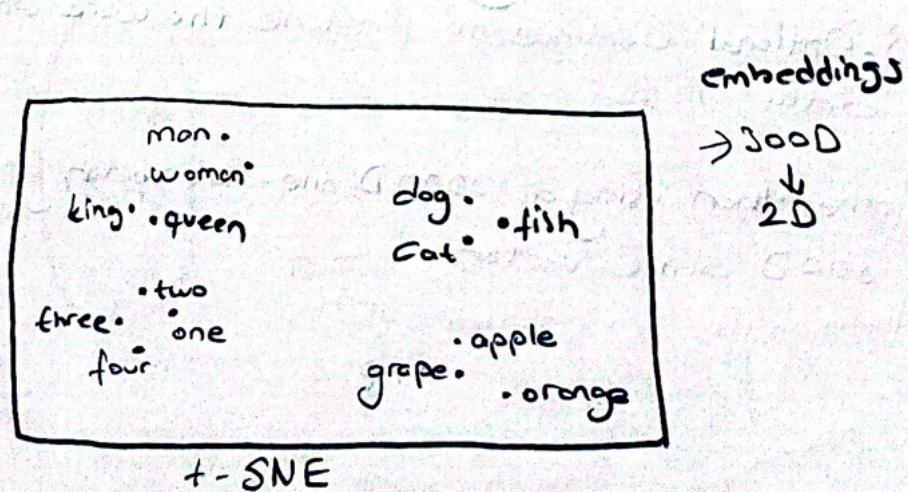
Mon (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
0	0	0	0	0	0
0	00	00	000	-	000
0	000	00	000	000	000
..
-	0	-	0	000	0
0	00	..
0	-	0	0	0	0
0	0	0	0	0	0

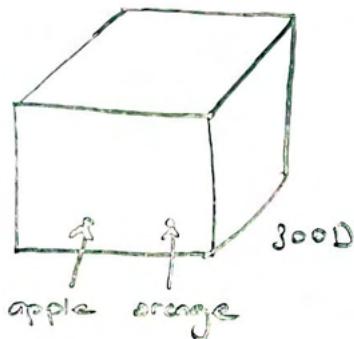
I want a glass of orange juice
I want a glass of apple juice

Featured representation: word embedding

	Men (5391)	women (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	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
Site	:	:	:	:	:	:
Cost						
Alive						
Verb						
None						

If there are 300 features, then becomes a 300 D vector for representing the word Man.



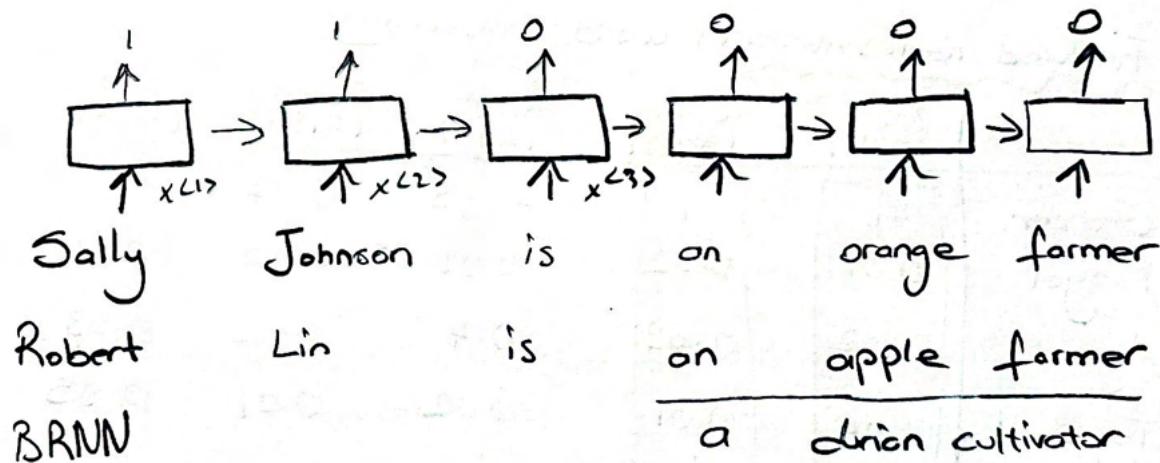


Orange gets embedded to a point in this
300D space

featureized representations
one hot representations

Using Word Embeddings

Named entity recognition example:



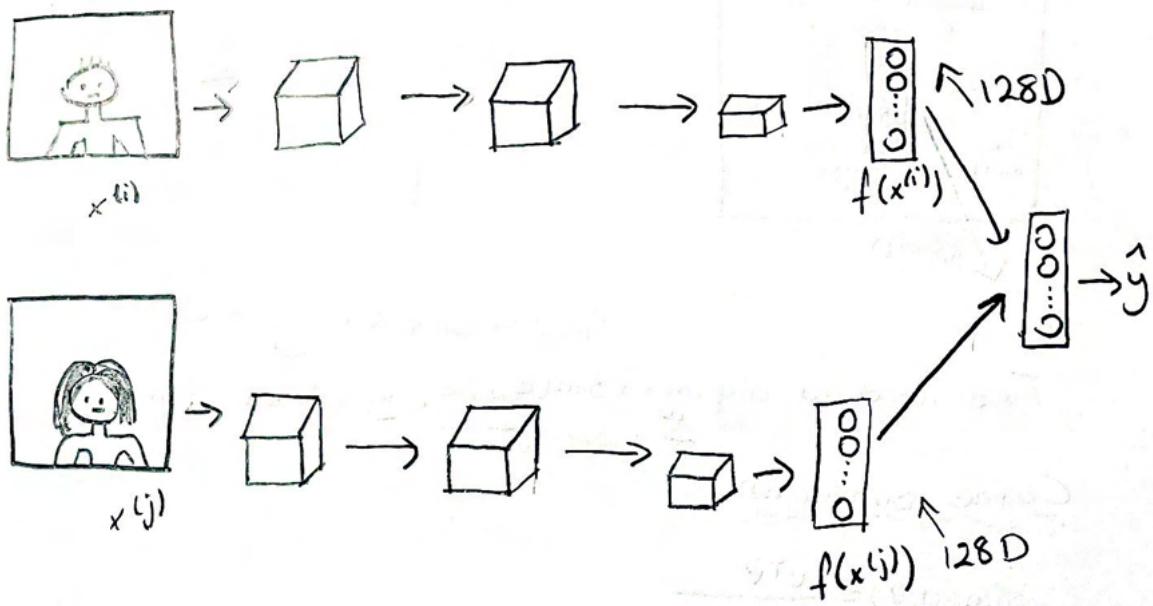
Transfer learning and word embeddings

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

Rather than using a 10000 D one-hot vector, you can instead use a 800 D dense vector.

through the one-hot vector

Relation to face encoding (embedding)



For face recognition we train a siamese network architecture that would learn a 128 D representation for different faces and compare these encodings in order to figure out if these two pictures are of the same face.

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	0.7	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97

↓
 e_{man}
 e_{woman}

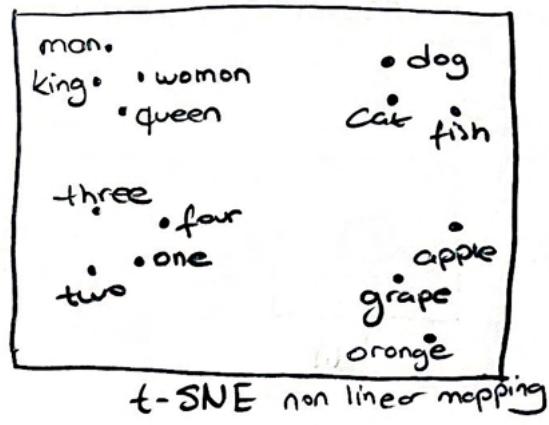
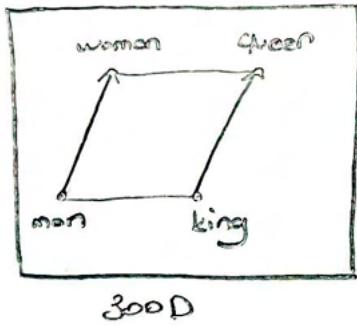
4D embeddings

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

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

$\text{Man} \rightarrow \text{Woman}$ as $\text{King} \rightarrow ?\text{Queen}$
 $e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} \rightarrow e_{\text{?Queen}}$

Analogies using word vectors

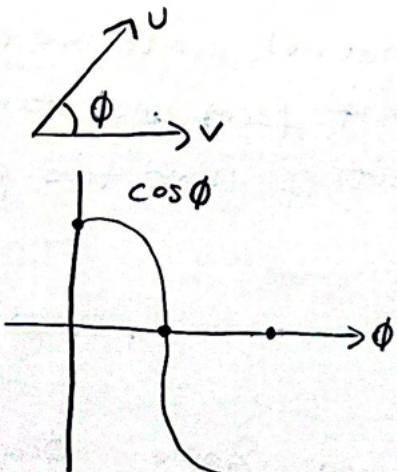


$$e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\text{?}}$$

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

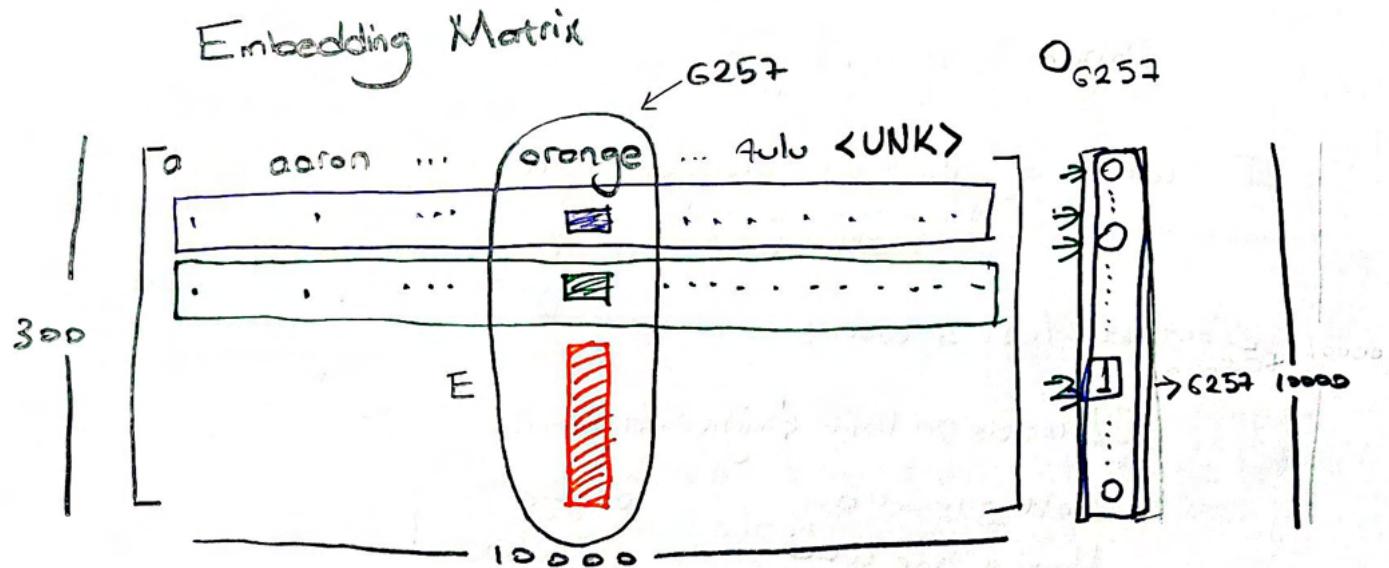
Cosine similarity

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



Squared distance or euclidean distance can also be used
 $\|u - v\|^2$

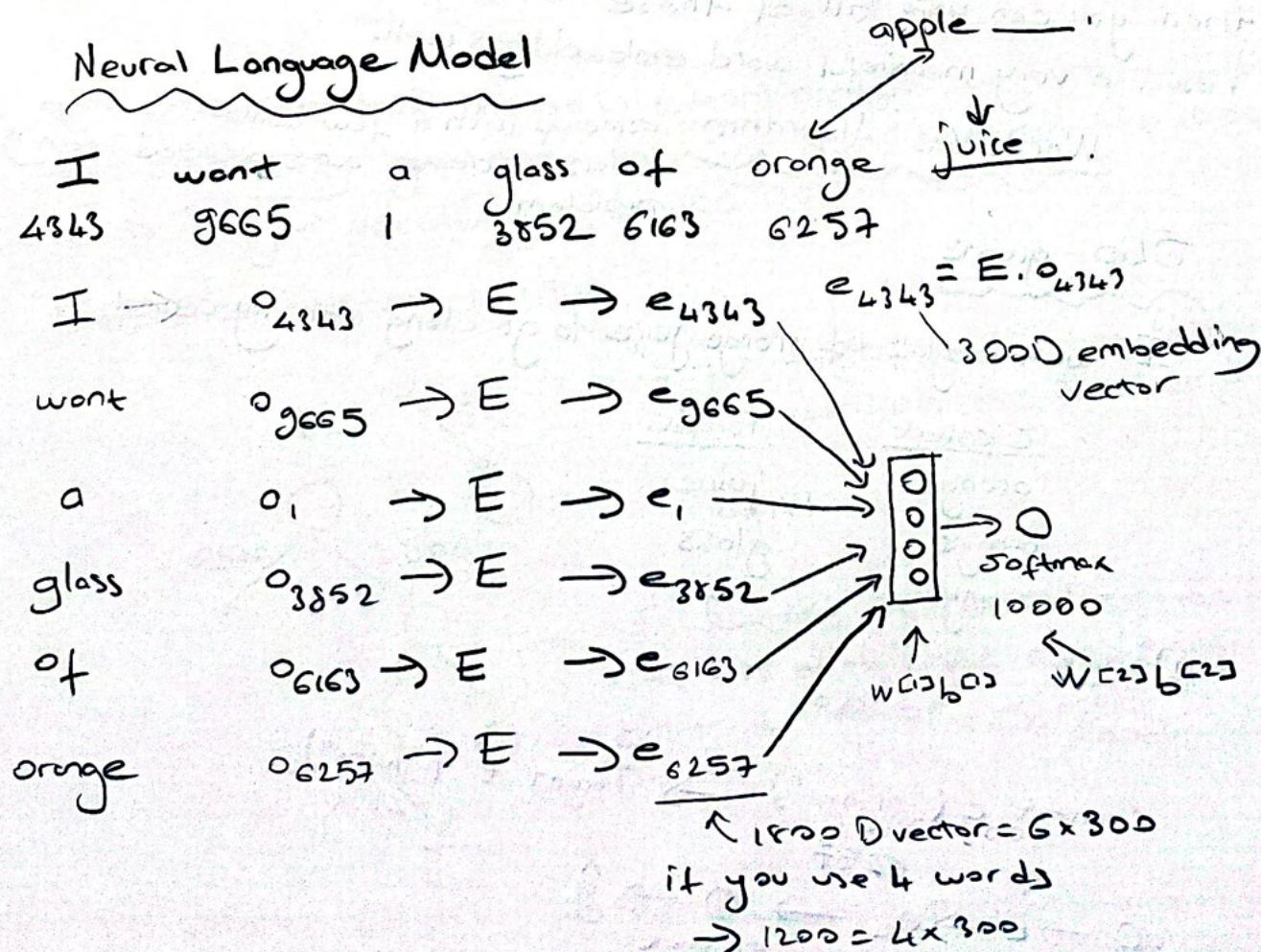
- Mon: Woman as Boy: Girl
- Ottawa: Canada as Nairobi: Kenya
- Big: Bigger as Tall: Taller
- Yen: Japan as Ruble: Russia



$$E \rightarrow (300, 10000)$$

$$O_{G257} \rightarrow (10000, 1)$$

Neural Language Model



Other context/target pairs

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

Context



target ... try to predict



Context: last 4 words

4 words on left & right a glass of orange ? to go along with

Last one word

orange ?

Nearby one word

glass ?

Skip gram

What researchers found was that if you really want to build a language model, it's natural to use the last few words as a context. But if your main goal is to learn a word embedding, then you can use all of these other contexts and they will result in very meaningful word embeddings well.

Word2Vec Algorithm

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

Skip-grams

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

Model

Vocab size = 10000

Context c ("orange") → Target t ("juice")
6257 4834

6257

$$O_c \rightarrow E \rightarrow e_c \rightarrow \odot \rightarrow \hat{y}$$

$$e_c = E O_c$$

one-hot vector

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

↑ target context ↑ transpose
↑ transpose ↑ transpose

θ_t = parameter associated with output t

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

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

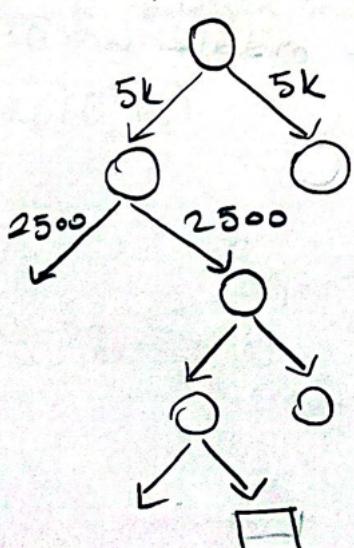
This is called skip-gram model because it is taking as input one word like orange and then trying to predict some words skipping a few words from the left or the right side.

There are a couple problems with using this algorithm. The primary problem is computational speed.

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

For the softmax model, everytime you want to evaluate this probability, you need to carry out a sum over all 10000 words in your vocabulary. (or 100,000-1,000,000)

Hierarchical softmax classifier.



is the target word in the first 5000 words in the vocabulary?

$$\log |V|$$

In practice, hierarchical softmax classifier doesn't use a perfectly balanced tree or perfectly symmetric tree.

Hierarchical softmax classifier can be developed so that the common words tend to be on top where as the less common words can be buried much deeper in the tree.

How to sample the context c ?

How do you choose the context c ? One thing you can do is just sample uniformly, at random, from your training corpus. When you do that, you find that there are some words like 'the', 'of', 'a', 'and', 'to'... appear extremely frequently. Your context to target mapping pairs just get these types of words very frequently, whereas there are other words like 'apple', 'orange', 'durian' that don't appear that often. You spend almost all the effort updating etc.

So in practice, the distribution of words P_c isn't taken just entirely uniformly at random for the training set purpose. But instead there are some heuristics that you could use in order to balance out something from the common words together with the less common words.

Cbow (continuous backwards) model takes the surrounding contexts from middle word, and uses the surrounding words to try to predict the middle word.

Negative Sampling

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

x	mapping		y
context	word	target	
orange	juice	L	?
orange	king	0	?
orange	book	0	?
orange	the	0	?
orange	of	0	?

k = 5-20 for small datasets

k = 2-5 for large datasets

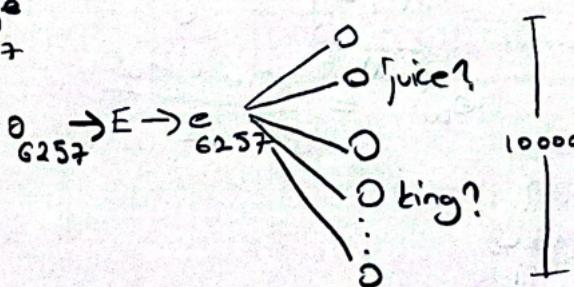
Softmax: $p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$

10000 way softmax

Define a logistic regression model:

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

Orange
6257



10000 binary logistic regression classifiers

10000 binary classification problems

On every iteration we're only going to train five of them ($k+1$). Computation cost of this algorithm is much lower

context	word	target
orange	juice	L
orange	king	0
orange	book	0
orange	the	0
orange	of	0

k-to-1 ratio of - to + examples

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

Glove (global vectors for word representation)

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

c, t

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

\uparrow \uparrow \uparrow
 c + t c

$$x_{ij} = x_{ji} \quad \text{depending on the definition of context}$$

Symmetric relationship

For the purposes of the Glove algorithm, we can define context and target as whether or not the two words appear in close proximity.

Model

$$\text{minimize} \quad \sum_{i=1}^{10000} \sum_{j=1}^{10000} f(x_{ij}) (\theta_i^T e_j + b_i + b_j - \log x_{ij})^2$$

this, is, of, a... \rightarrow stop words

θ_i, e_j are symmetric

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

How related are those two words?
words + one
j i

End product is a good predictor for how often the two words occur together

Sentiment Classification

Simple Sentiment Classification Model

The dessert is excellent.

8928 2468 4694 3180

★★★★★

The 8928 → E → e₈₉₂₈

dessert 2468 → E → e₂₄₆₈

is 4694 → E → e₄₆₉₄

excellent 3180 → E → e₃₁₈₀

100B
words

300D

Average
Avg

300D

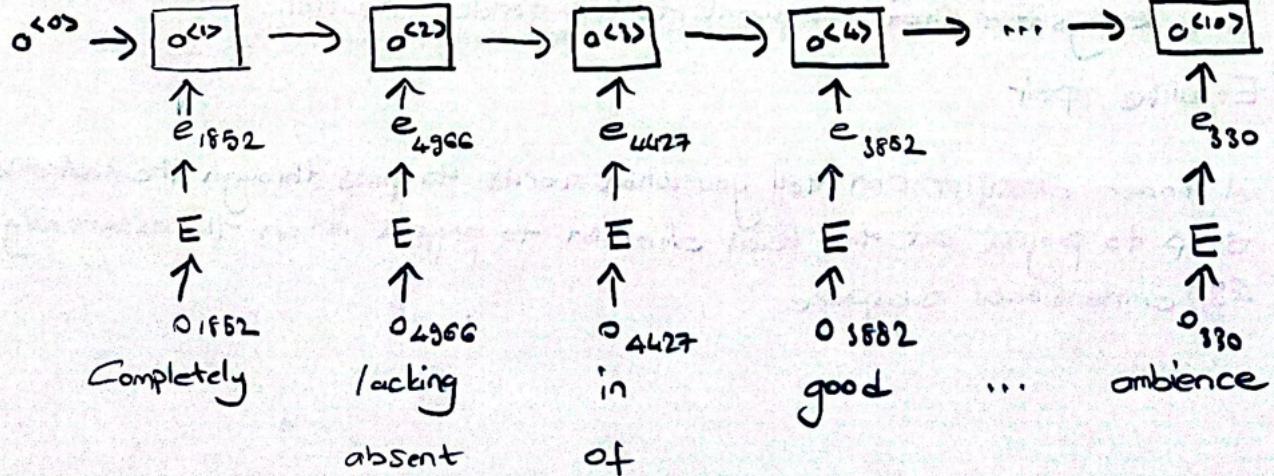
1-5
stars

One of the problems with this algorithm is it ignores word order.

"Completely lacking in good taste, good service, and good ambience" this is a very negative review.

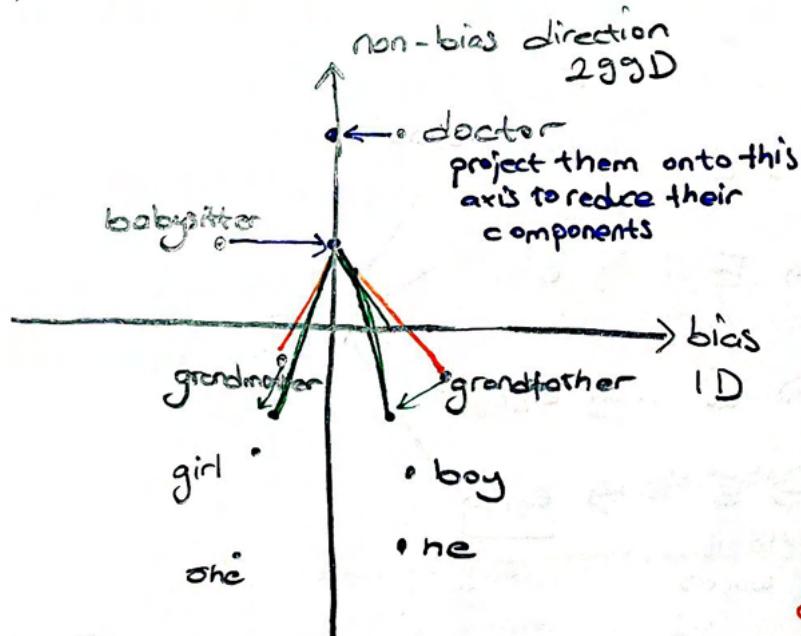
RNN for Sentiment Classification

many-to-one



Debiasing Word Embeddings

Addressing Bias in Word Embeddings



Singular Value Decomposition (SVD)
uses ideas similar to PCA

in this example, the distance (similarity) between babysitter and grandmother is smaller than the ~~the~~ distance between babysitter and grandfather

1- Identify bias direction

$$\begin{aligned} & e_{\text{he}} - e_{\text{she}} \\ & e_{\text{male}} - e_{\text{female}} \end{aligned} \quad \left. \begin{array}{l} \text{average} \\ \text{overage} \end{array} \right.$$

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

There are some words that intrinsically capture gender. (grandmother, grandfather, girl, boy, she, he) Whereas there are other words like doctor and babysitter that we want to be gender neutral.

3- Equalize pair

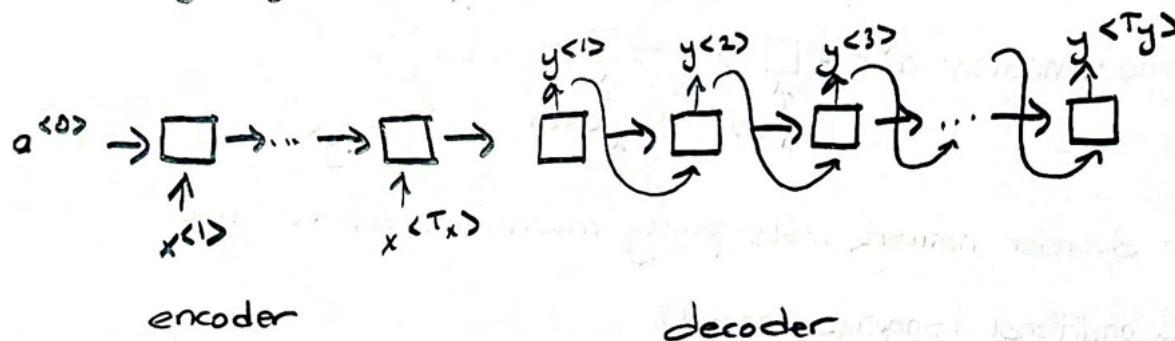
A linear classifier can tell you what words to pass through the neutralization step to project out this bias direction to project it on this essentially 2gg dimensional subspace

Sequence to Sequence 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.

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

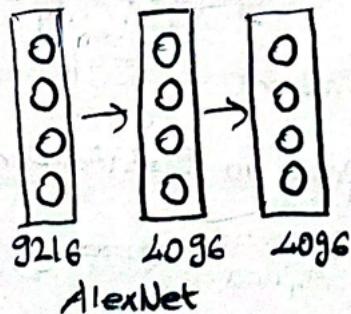


After ingesting the input sequence the RNN outputs a vector that represents the input sequence.

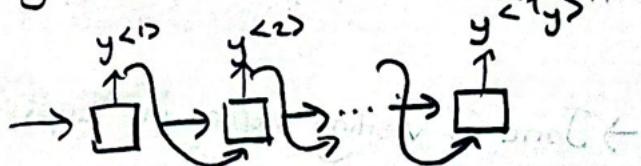
Image Captioning



11×11
 $s=4$

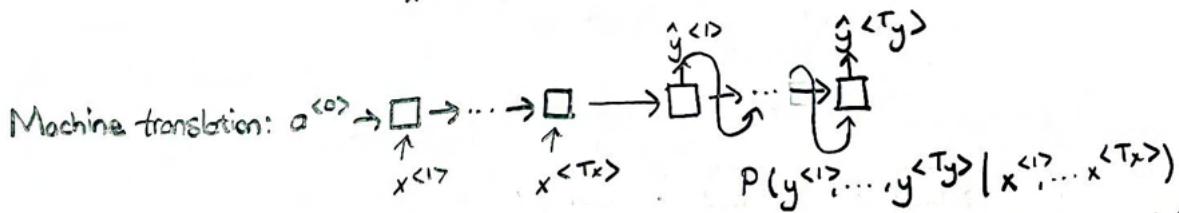
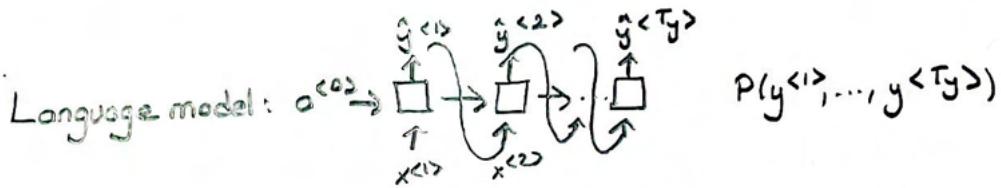


A cat sitting on a chair.
 $y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad y^{<4>} \quad y^{<5>} \quad y^{<6>}$



Pre-trained AlexNet can be the encoder network for the image. Now you have a 4096 dimensional vector that represents the image.

Machine Translation as Building a Conditional Language Model



The decoder network looks pretty much identical to the language model.

"Conditional language model".

Modeling the probability of the English translation, conditions on some input French sentence.

Finding the Most Likely Translation

Jane visite l'Afrique en septembre.

$$P(\underbrace{y^{(1)}, \dots, y^{(T_y)}}_{\text{english}} | \underbrace{x}_{\text{french}})$$

→ Jane is visiting Africa in September

→ Jane is going to be visiting Afrika in September

→ In September, Jane will visit Africa.

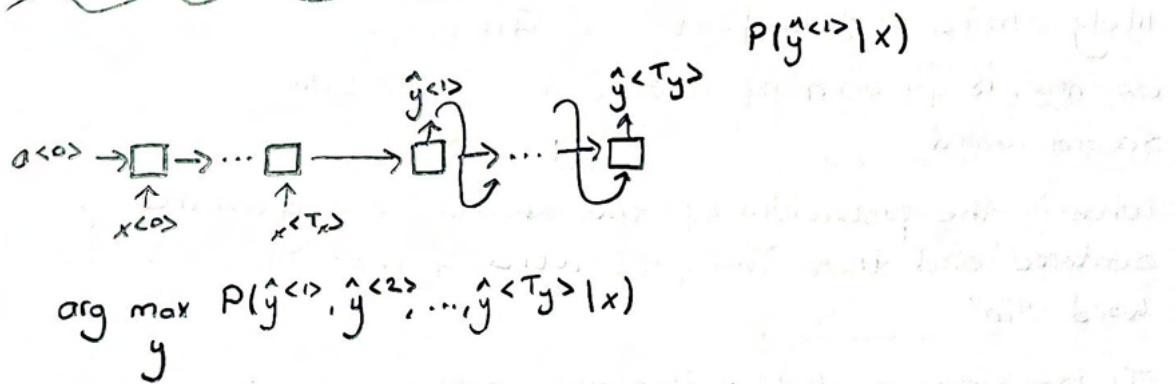
→ Her African friend welcomed Jane in September.

$$\underset{y^{(1)}, \dots, y^{(T_y)}}{\operatorname{arg\ max\ }} P(y^{(1)}, \dots, y^{(T_y)} | x)$$

When you're using this model for machine translation, you're not trying to sample at random from this distribution. Instead, what you would like to find the English sentence y , that maximizes that conditional probability. Developing a machine translation system, one of the things you need to do is come up with an algorithm that can actually find the value of y that maximizes this term.

Greedy search says to generate the first word just pick whatever is the most likely first word according to your conditional language model.

Why not a greedy search?



→ Jane is visiting Africa in September

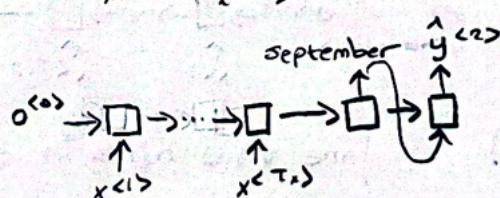
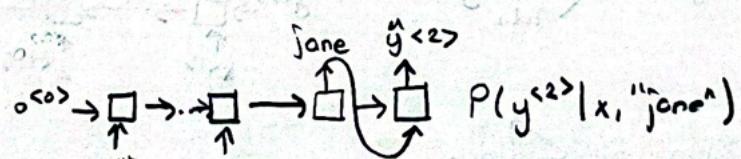
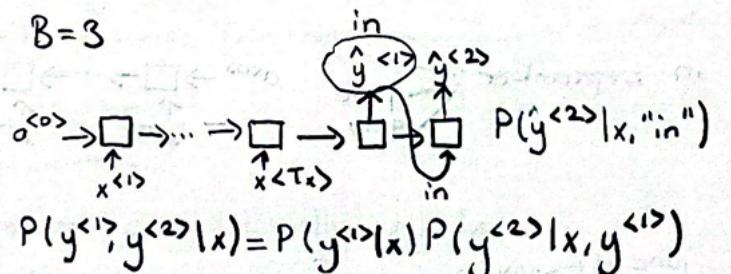
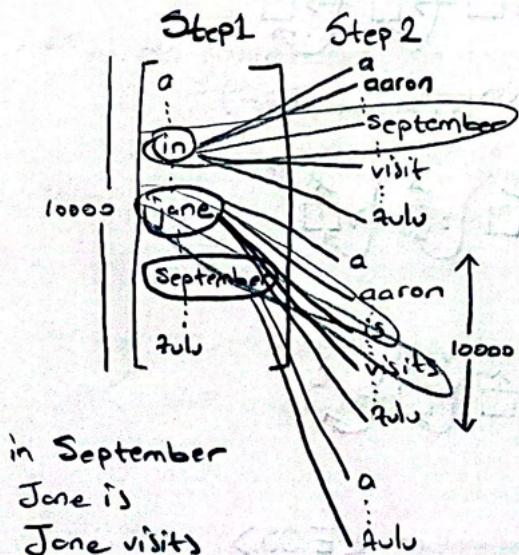
→ Jane is going to be visiting Africa in September

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

"going" is more common English word.

It's not always optimal to just pick one word at a time.

Beam Search Algorithm $B=3$

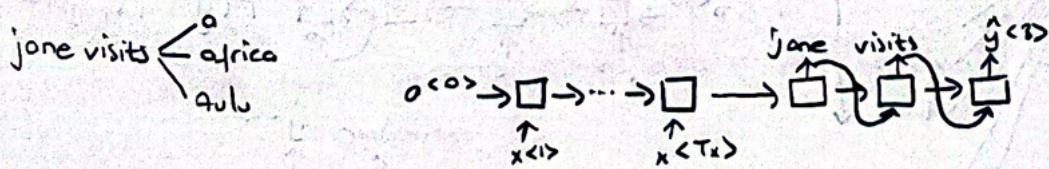
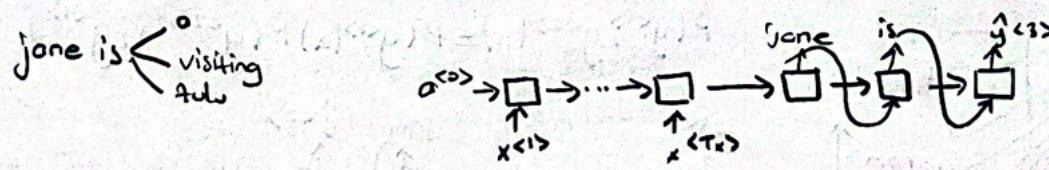
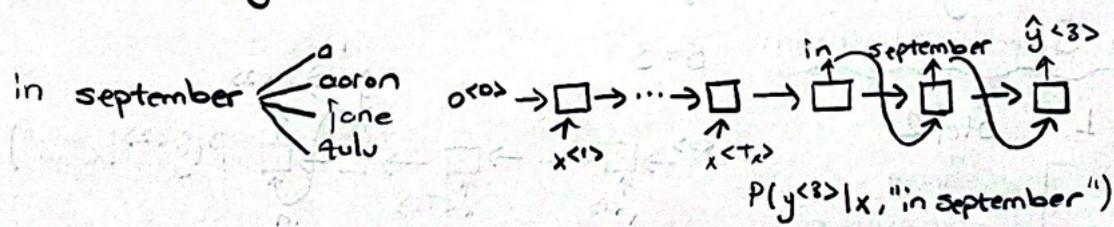


Whereas greedy search will pick only the one most likely word and move on, Beam search instead can consider multiple alternatives. So the Beam search algorithm has a parameter called B , which is called the beam width. This means beam search will consider not just one possibility but consider (for this example) 3 at a time.

In step 1, beam search picked Jane and September as the 3 most likely choice of the first word. In step 2, what beam search will do now, is for each of these 3 choices consider what should be the second word.

What is the probability of the second word given the input French sentence and that the first word of the translation has been the word "in"

If beam search decides that the most likely choices are in the first and second words are "in September," "Jane is", or "Jane visits" then what that means is that it is now rejecting September as a candidate for the first words of the output English translation. So we are now down to two possibilities for the first word but we still have a beam width of 3. Because of $B=3$, every step you instantiate 3 copies of the network to evaluate these partial sentence fragments and the output.



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

Jane visits africa in september. <EOS>

Refinements to Beam Search

Length Normalization

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

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

these probabilities are all numbers less than 1 result in a tiny number. which can result in numerical underflow. It's too small for the floating point representation in your computer to store accurately. So in practice, instead of maximizing this product, we will take logs

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

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

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

$$\alpha = 0.7 \quad \alpha = 1 \quad \alpha = 0$$

This takes the average of the log of the probability of each word. And this significantly reduces the penalty for outputting longer translations.

How do you choose the beam width B?

Large B: better result, slower, computationally expensive

Small B: worse result, faster

In production $\rightarrow 10, 100, 1000, 3000$

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

Error Analysis in Beam Search

Beam Search is an approximate search algorithm, also called a heuristic search algorithm. So it doesn't always output the most likely sentence.

Jane visite l'Afrique en septembre.

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

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

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

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

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

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

Conclusion: Beam search is at fault

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

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	2×10^{-10}	1×10^{-10}	B
...	...	—	—	R
...	...	—	—	B
				R
				R

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

Bleu Score (Bilingual Evaluation Understudy)

If there are multiple great answers, how do you measure accuracy?

French: Le chat est sur le tapis.

Reference 1: The cat is on the mat.

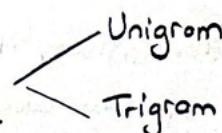
Reference 2: There is a cat on the mat.

Given a machine generated translation, bleu score allows you to automatically compute a score that measures how good is that machine translation.

MT output: the the the the the the the

Precision: $\frac{7}{7}$ Modified precision: $\frac{2}{7} \leftarrow \begin{matrix} \text{count clip ("the")} \\ \text{count ("the")} \end{matrix}$

In modified precision measure, we will give each word a credit up to the maximum number of times it appears in the reference sentences.

Bleu Score on Bigrams 

Bigrams means pairs of words appearing next to each other.

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. (ŷ)

The possible bigrams are:

	Count	count clip	
the cat	2 ←	1 ←	
cat the	1 ←	0	$\frac{4}{6}$
cat on	1 ←	1 ←	
on the	1 ←	1 ←	
the mat	1 ←	1 ←	

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

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

Bleu Details

P_n = Bleu score on n-grams only

Combined Bleu score: $\text{BP} \exp\left(\frac{1}{4} \sum_{n=1}^4 P_n\right)$

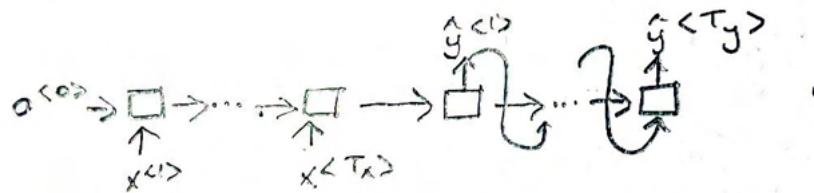
BP = Brevity Penalty

$$\text{BP} = \begin{cases} 1, & \text{if } \text{MT_output_length} > \text{reference_output_length} \\ \exp(1 - \text{reference_output_length}/\text{MT_output_length}), & \text{otherwise} \end{cases}$$

If you output very short translations, it's easier to get high precision. Because probably most of the words you output appear in the references. (We don't want this) BP is an adjustment factor that penalizes translation systems that output translations are too short.

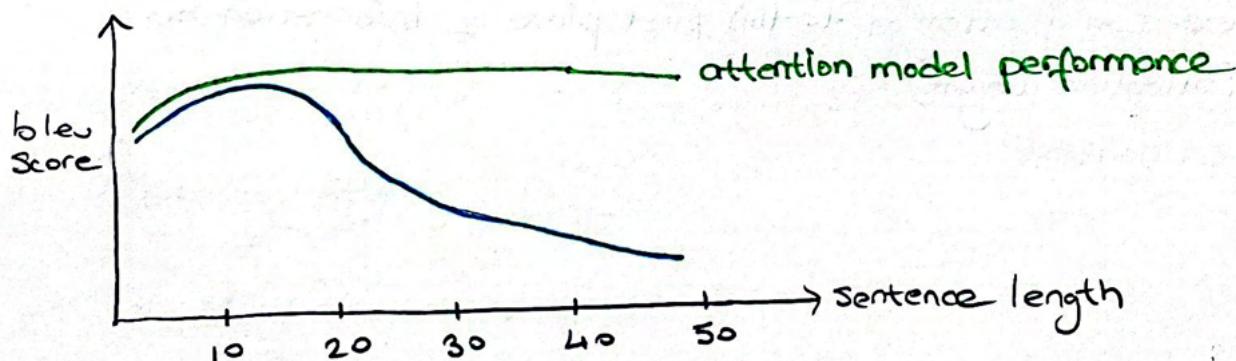
Attention Model Intuition

The Problem of Long Sequences



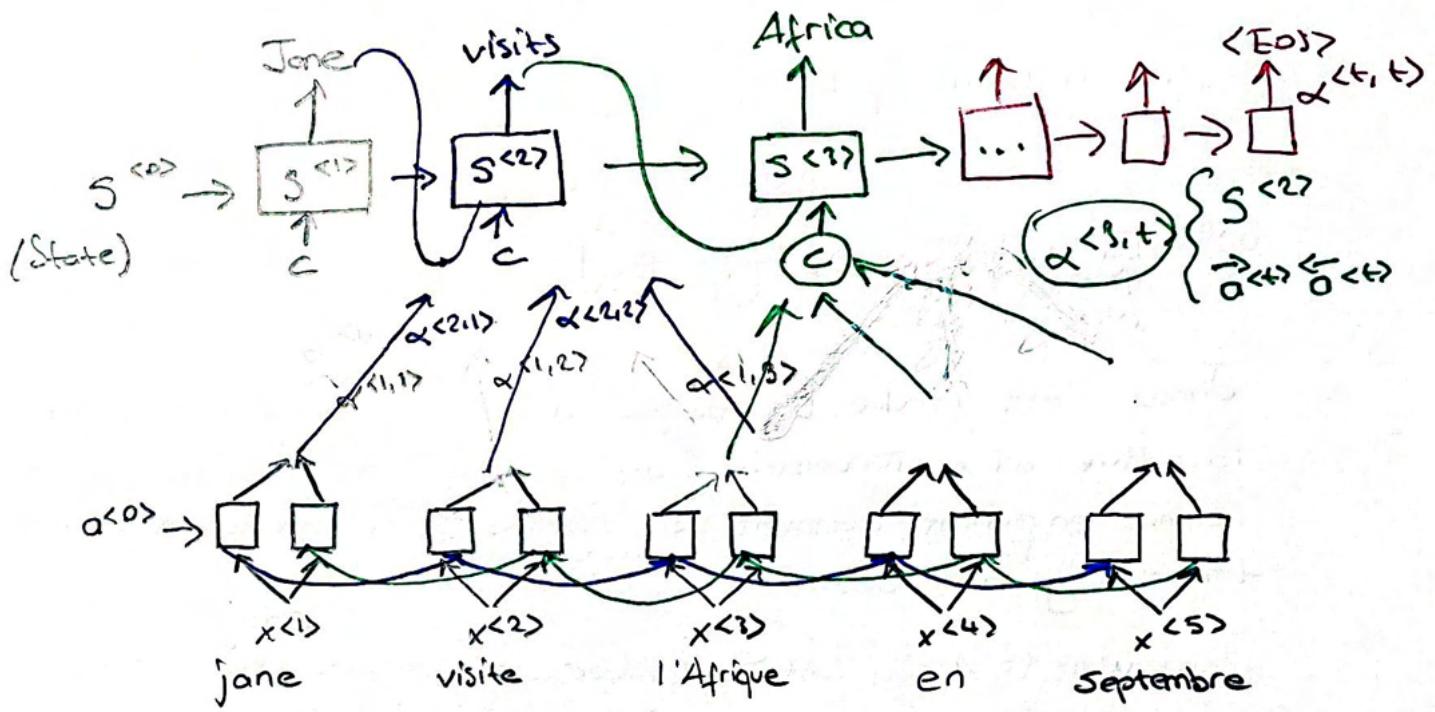
Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente d'y aller aussi.

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.



The encoder has to read the whole sentence and memorize it. A human translator would read the first part of the sentence, maybe generate part of the translation, look at the second part, generate a few more words and so on. (It's hard to memorize the whole sentence)

Attention model translates a bit more like humans might.



$\alpha^{(1,1)}$: when you're generating the first word, how much should you be paying attention to this first piece of information here.
 (attention weight)

c : context

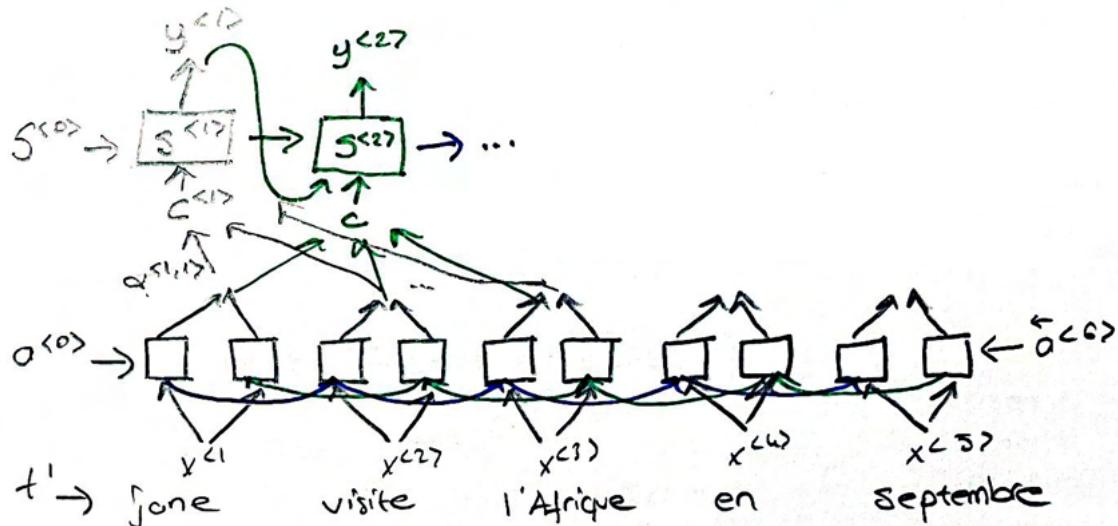
Attention Model

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

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

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

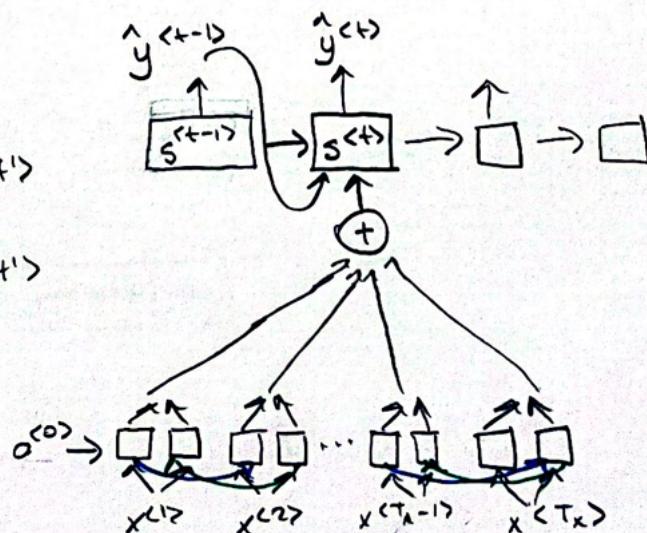
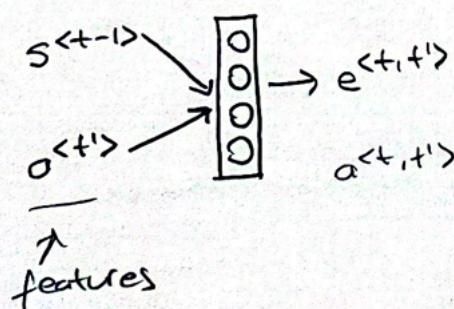
$$C^{(2)} = \sum_{t'} \alpha^{(2,t')} o^{(t')}$$



$$o^{(t)} = (\vec{o}^{(t)}, \overleftarrow{o}^{(t)}) \quad t' \rightarrow \text{French sentence}$$

Computing Attention $\alpha^{(t,t')}$

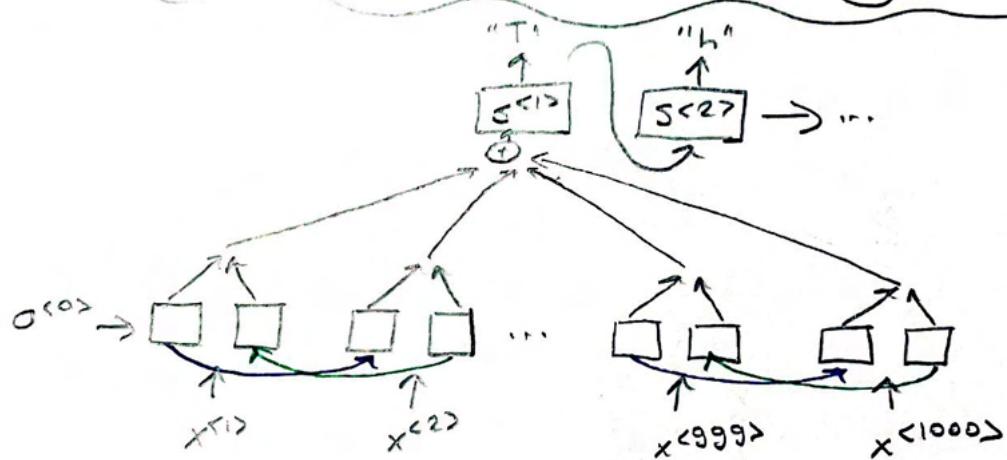
$$\alpha^{(t,t')} = \frac{\exp(e^{(t,t')})}{\sum_{t'=1}^{T_x} \exp(e^{(t,t')})}$$



total number of $\alpha^{(t,t')}$ parameters $T_x \cdot T_y$

Quadratic cost

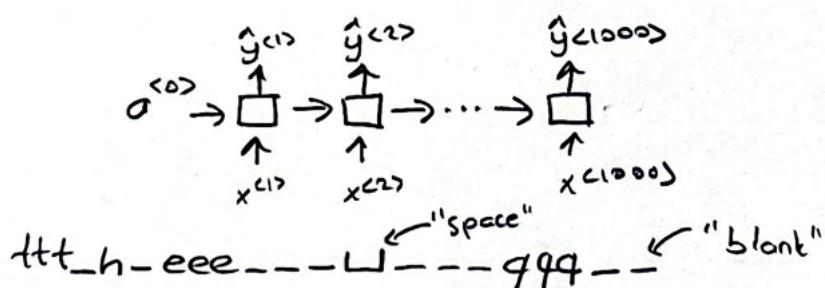
Attention Model for Speech Recognition



CTC Cost for Speech Recognition

(Connectionist Temporal Classification)

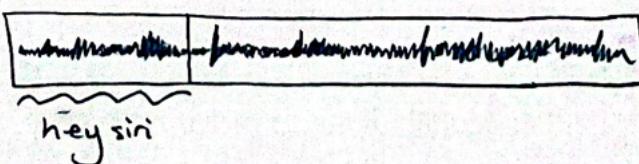
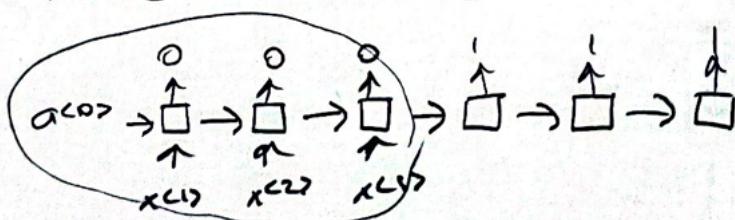
"the quick brown fox" \leftarrow 19 characters



10 secs of audio
100 hertz
100 samples per sec

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

Trigger Word Detection



Self-Attention Intuition

$A(q, K, V)$ = attention-based vector representation of a word
calculate for each word $A^{<1>} \dots A^{<5>}$

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

swani et al. 2017, Attention Is All You Need]

Andrew Ng

values
 $\{$
 $k^{<i>} : \text{key}$
 $v^{<i>} : \text{value}$

Self-Attention

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

Scaled dot product attention

The complexity of computing the softmax for SA block is quadratic with respect to the length of the input sequence.

$$\text{Softmax} \left(\frac{\exp(e^{q \cdot k^T})}{\sum_j \exp(e^{q \cdot k^T})} \right) v^{<i>}$$

$$A(q, K, V)$$

$$= \sum_i \exp(e^{q \cdot k^T}) v^{<i>}$$



Did what?

Query (Q)

Key (K)

Value (V)

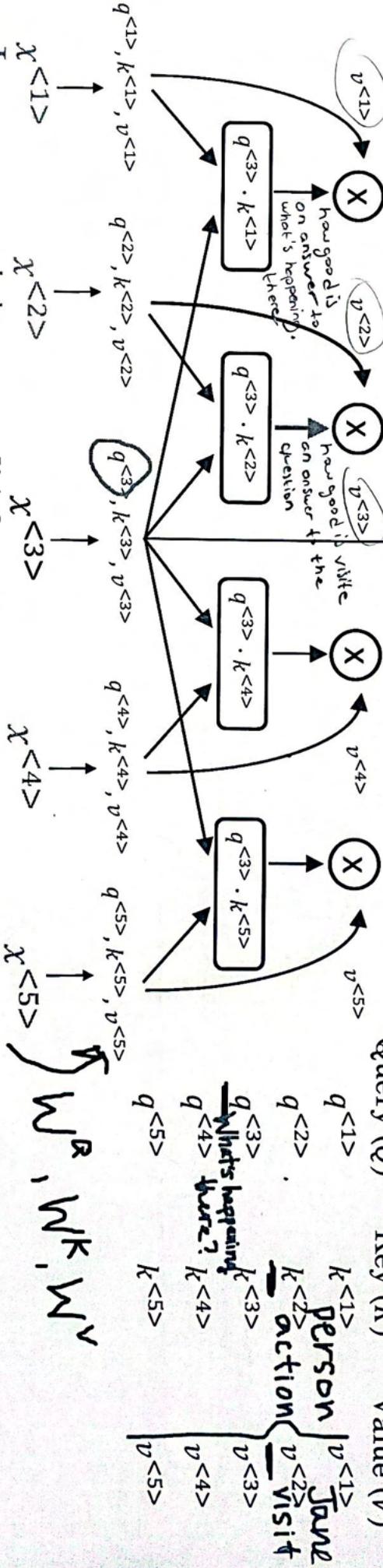
$q^{<1>} k^{<1>} v^{<1>} \text{Jane}$

$q^{<2>} k^{<2>} v^{<2>} \text{visit}$

$q^{<3>} k^{<3>} v^{<3>} \text{person}$

$q^{<4>} k^{<4>} v^{<4>} \text{action}$

$q^{<5>} k^{<5>} v^{<5>} \text{visit}$



Jane

visite

l'Afrique

en

septembre

Vaswani et al. 2017, Attention Is All You Need]

$q^{<3>} \in \mathbb{R}^{d_q}$ is a question that you get to ask about l'Afrique, represents a question like "what's happening there?"

$$q^{<3>} = W^Q x^{<3>} \\ k^{<3>} = W^K x^{<3>} \\ v^{<3>} = W^V x^{<3>}$$

Andrew Ng

Multi-Head Attention

consists of multiple Self-Attention blocks (heads)

$$\text{MultiHead}(Q, K, V) = \text{concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_n)W_0$$

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

$$h = \# \text{ heads}$$

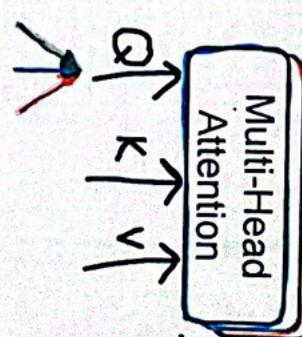
$$W_0^{\text{head}_i} = \text{Attention}(W_i^Q Q, W_i^K K, W_i^V V)$$

concat

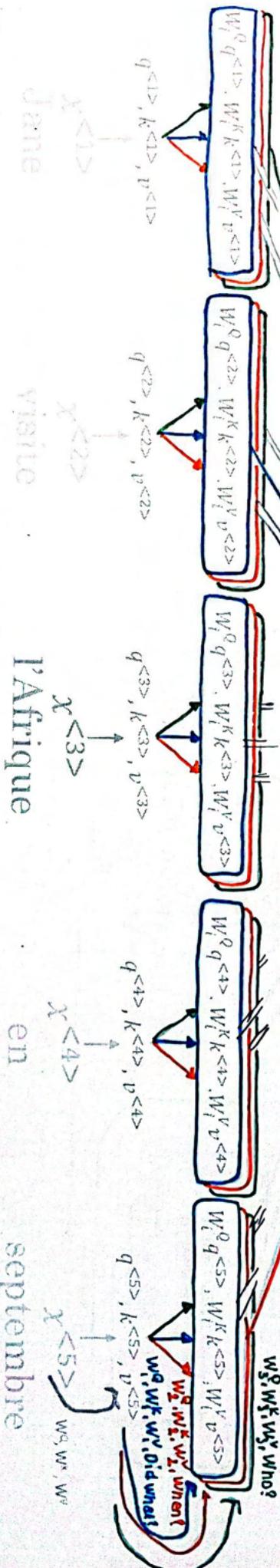
$$\text{Attention}(W_i^Q Q, W_i^K K, W_i^V V)$$

heads

$$\begin{matrix} Q \\ K \\ V \end{matrix}$$



you can compute different heads in a for loop
in practice you can calculate them in parallel because no one
has value depends on the value of any other head



Vaswani et al. 2017, Attention Is All You Need

$W_i^Q, W_i^K, W_i^V \rightarrow$ What's happening there?

$W_2^Q, W_2^K, W_2^V \rightarrow$ Who! Something happening?

$W_3^Q, W_3^K, W_3^V \rightarrow$ Who? The inner product between Jane's key vector and the l'Afrique query vector will be the highest

The complexity of computing the softmax for SA block Andrew Ng

is quadratic with respect to the length of the input sequence.

l'Afrique

en

septembre

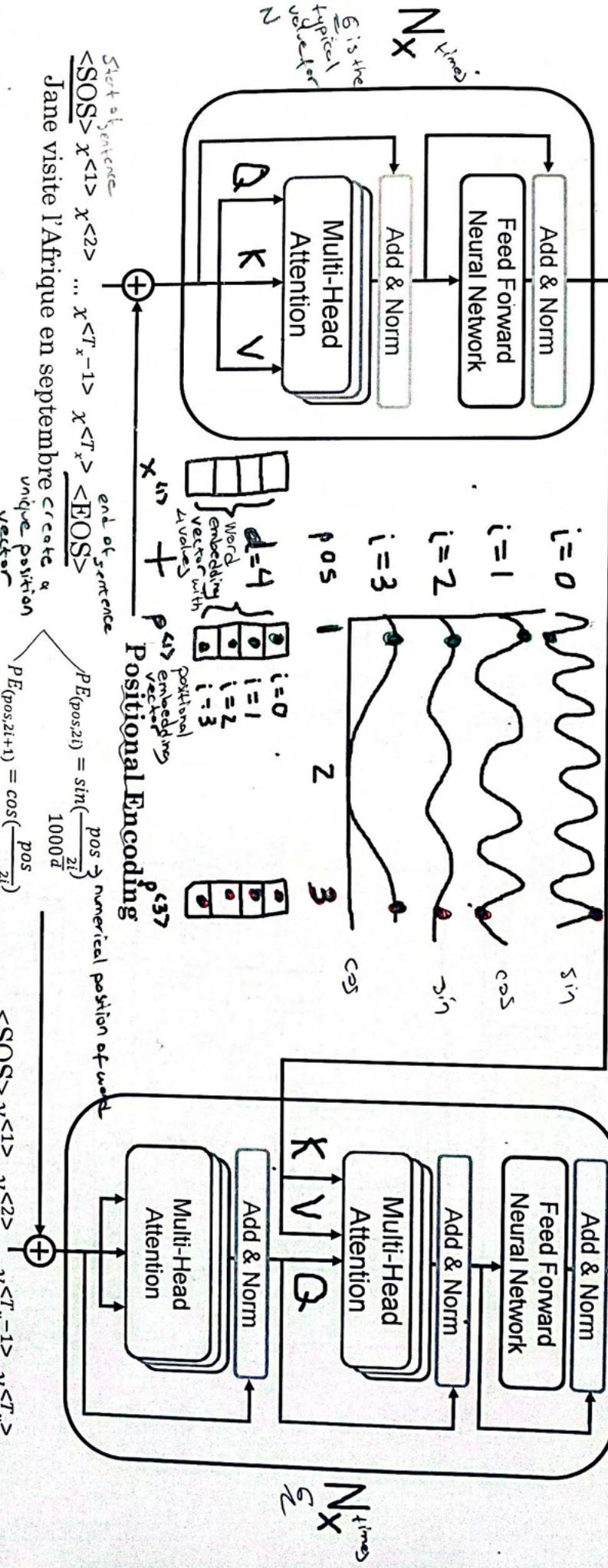
en

septembre

Transformer Details

<SOS> Jane visits Africa in September <EOS>

Encoder



Vaswani et al. 2017, Attention Is All You Need]

Andrew Ng

Outputs of the embedding layer is d (which in this case 4) the max length of sequence your model can take
In addition to adding these position encodings to the embeddings you'd also pass them through the network with residual connections.

Decoder

