

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

AGCCCCCTGTGAGGAAC TAG



AGCCCCCTGTGAGGAAC **T**AG

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.

Motivating example

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

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

$$T_x = 9$$

$\rightarrow y:$

$1 \quad 1 \quad 0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0$

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

$$T_y = 9$$

$x^{(i)<t>}$

$$T_x^{(i)} = 9$$

15

$y^{(i)<t>}$

$$T_y^{(i)}$$

We want to do named entity recognition to identify where the names are in a sentence --
Training data is sentence mapped to 1s and 0s.

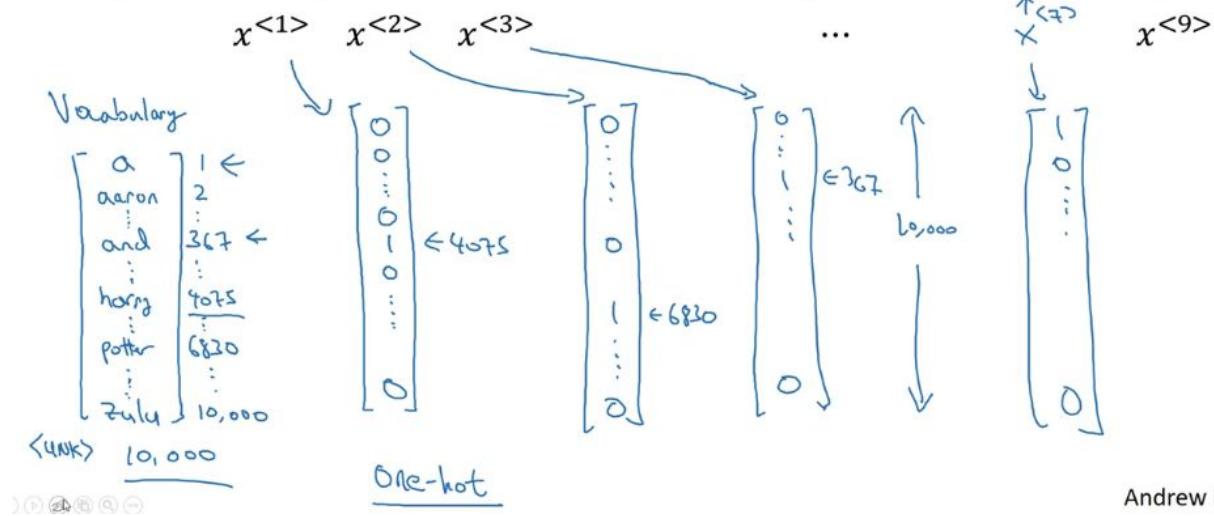
ONE HOT VECTOR

Representing words

$$x \xrightarrow{\leftrightarrow} (x, y)$$

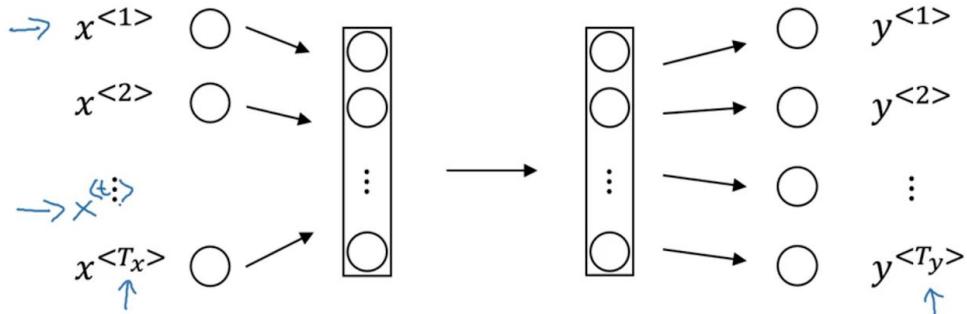
$$x \rightarrow y$$

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



Andrew Ng

Why not a standard network?

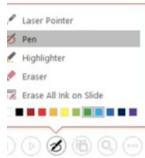
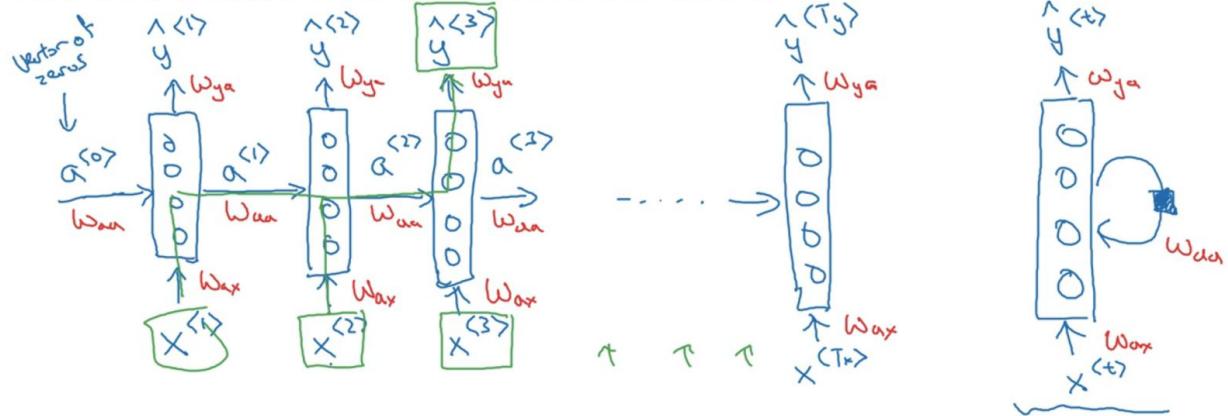


Problems:

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

Need a better way to generalize for knowledge learnt in one part to be used elsewhere. Like CNNs for images to generalize the angles etc.

Recurrent Neural Networks



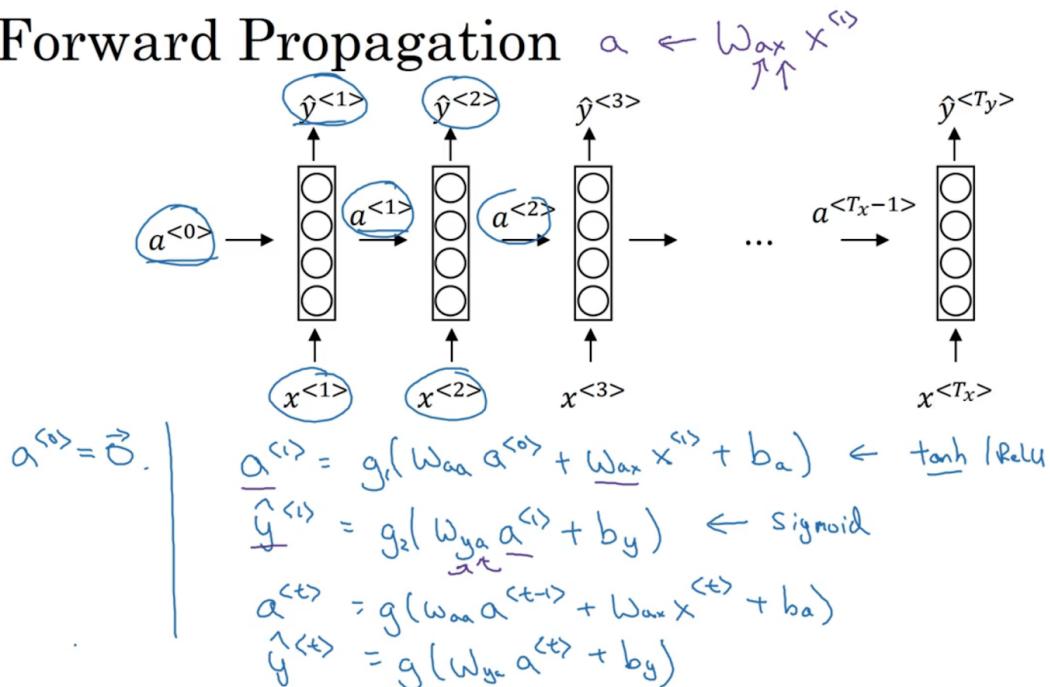
He said, "Teddy Roosevelt was a great President."

He said, "Teddy bears are on sale!"

Andrew N

Only uses past to predict future words. Bi directional RNNs are better for this.

Forward Propagation



Andrew Ng

Activations or rather RNN has 100 elements in it, there are 10000 words each word having its own weight vector -- this results in current activation which goes into next activations.

Simplified RNN notation

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

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

$$y^{<t>} = g(W_ya^{<t>} + b_y)$$

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

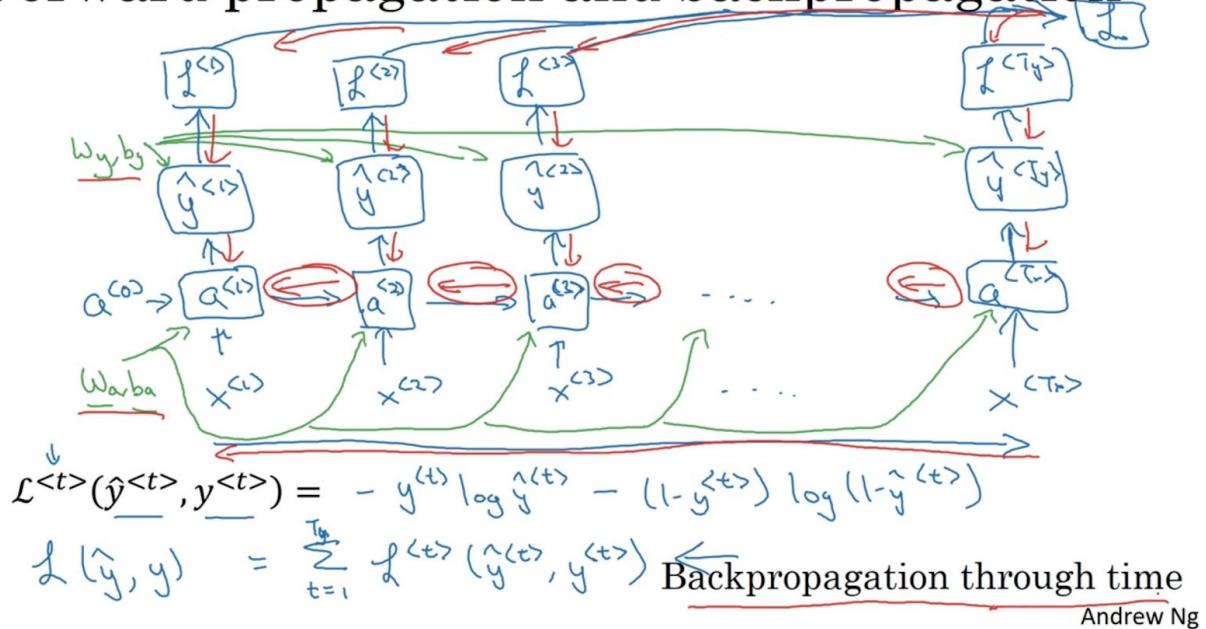
$$W_a = \begin{bmatrix} W_{aa} & W_{ax} \end{bmatrix} \quad (100, 100) \times (100, 10, 000) = (100, 10100)$$

$$[a^{<t-1>}, x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} \quad (100, 10000) \times (100, 10100) = (100, 10100)$$

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

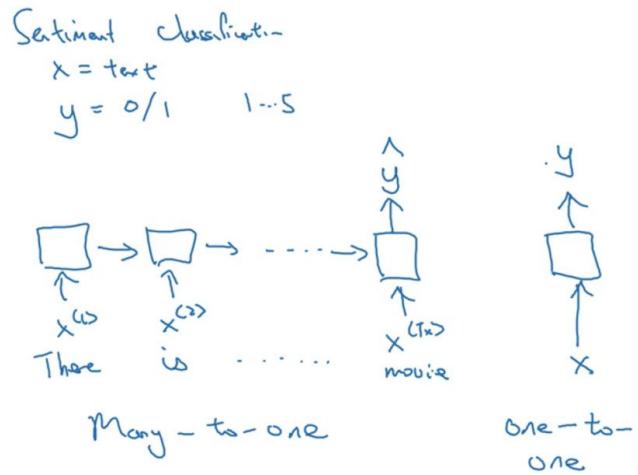
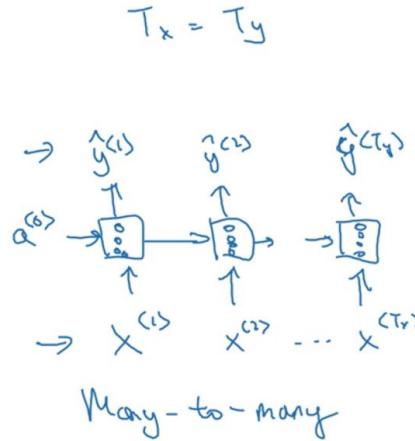
LOSS : Compute loss at every time step and combine for one single loss. Basically at every time step, there is a sigmoid 0 or 1. So you use $-\text{ylog}(\hat{y})$ at every point and try to minimize it so that you can best predict where Harry Potter is in the sentence.

Forward propagation and backpropagation

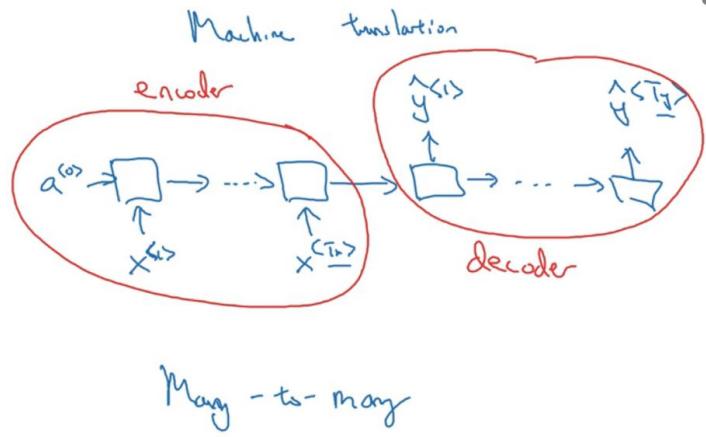
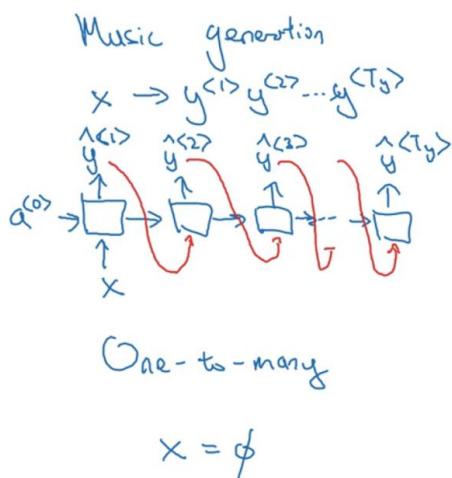


ONE to ONE is vanilla neural network. asd

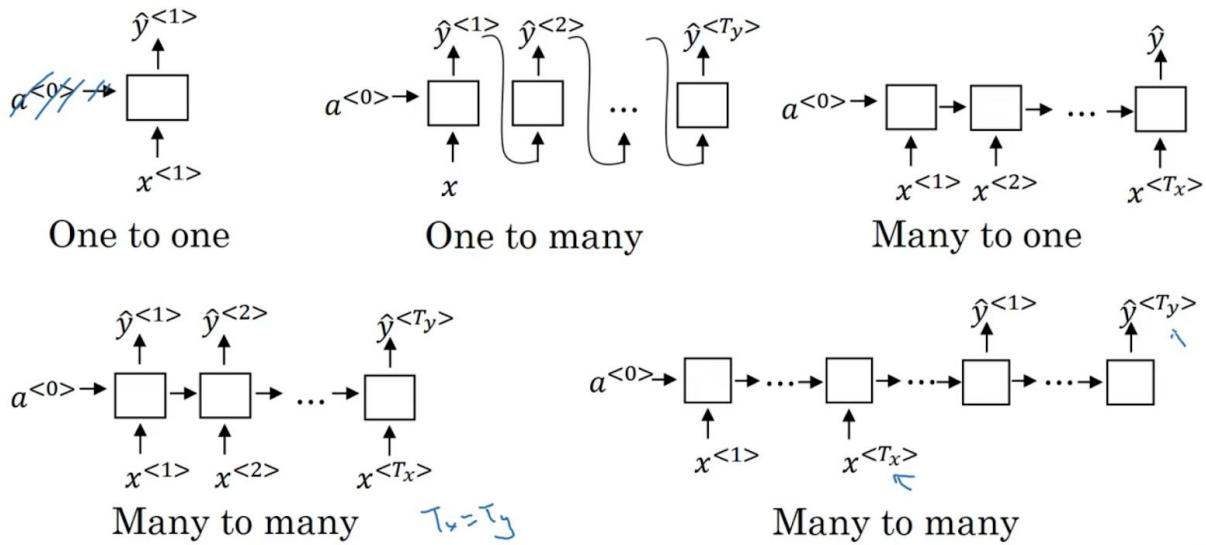
Examples of RNN architectures



Examples of RNN architectures



Summary of RNN types



LANGUAGE MODELLING:

Identical sentences, which one is more likely?

What is language modelling?

Speech recognition

The apple and pair salad.

→ The apple and pear salad.

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

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

$$P(\text{sentence}) = ?$$

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

Language modelling with an RNN

Training set: large corpus of english text.

Tokenize

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

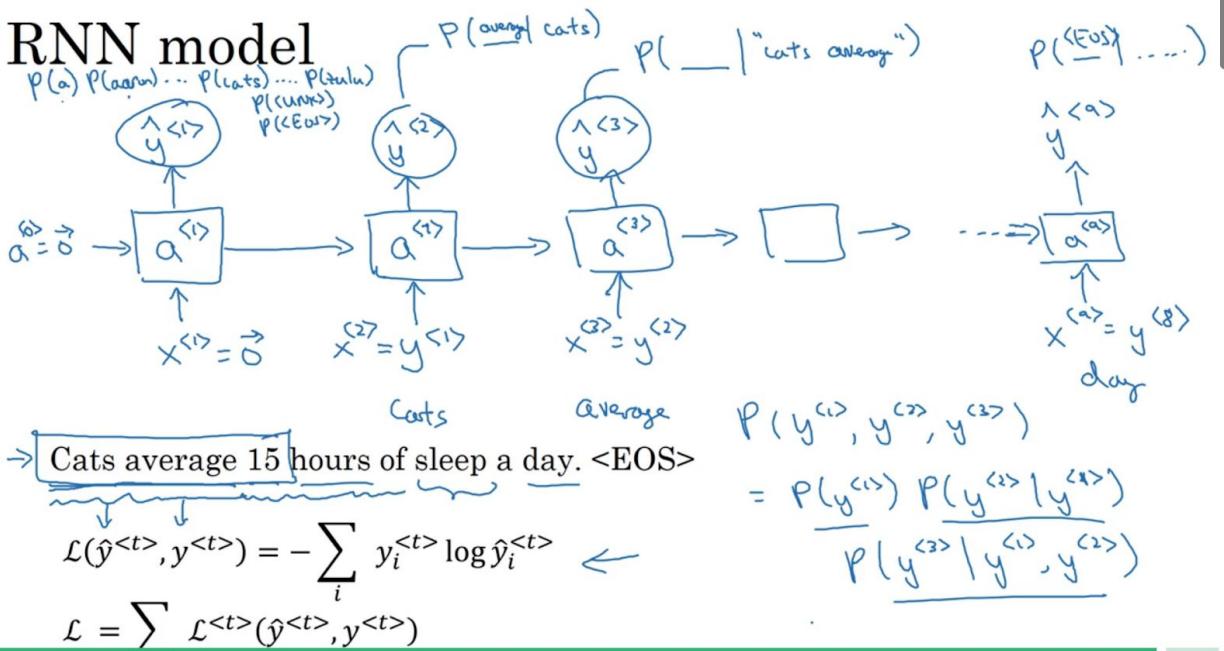
$y^{(1)}$ $y^{(2)}$ $y^{(3)}$... $y^{(8)}$ $y^{(9)}$

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

10,000

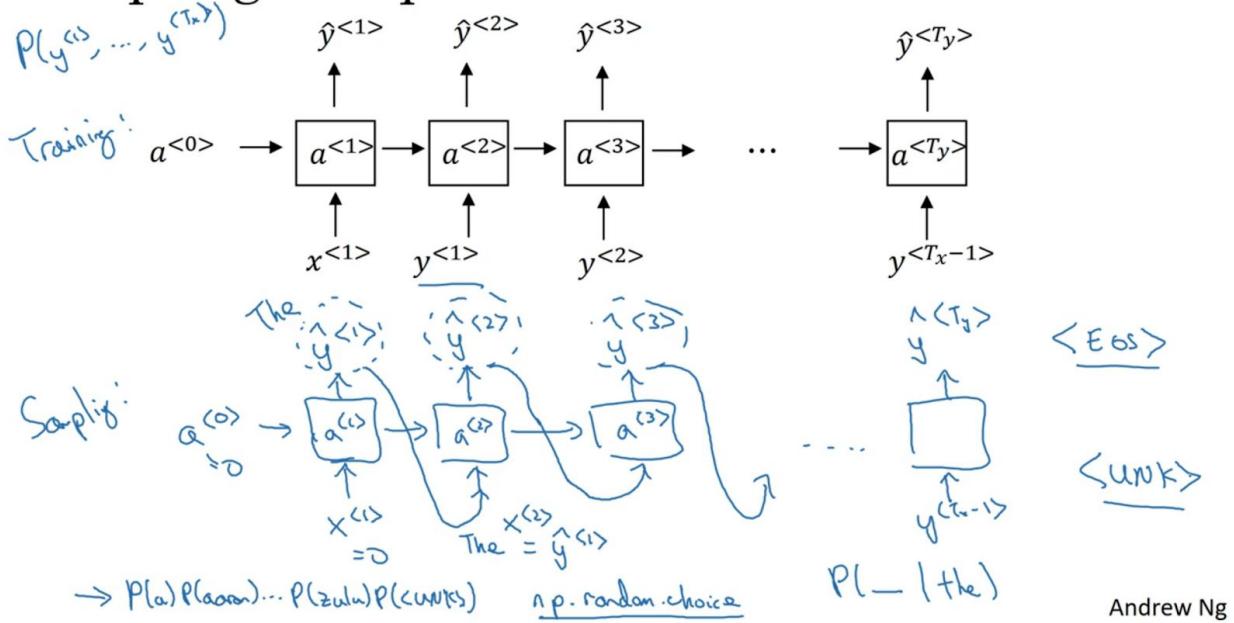
$\langle \text{UNK} \rangle$

RNN model



GENERATING A RANDOM SENTENCE FROM RNN:

Sampling a sequence from a trained RNN

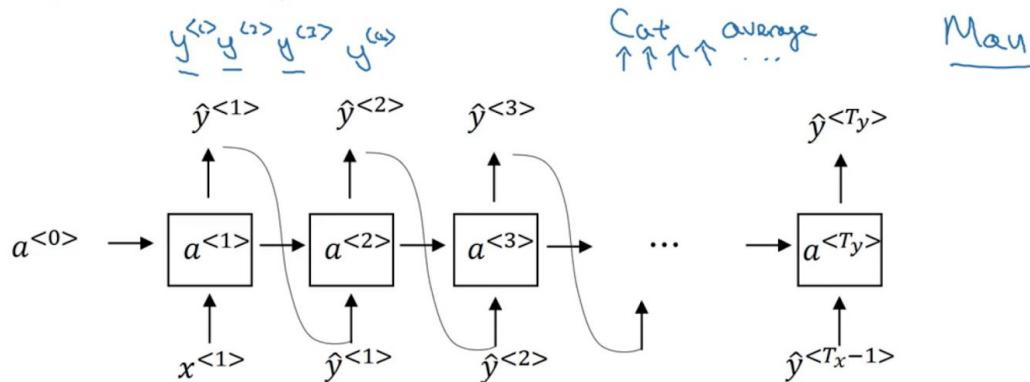


More expensive and challenging as 2 dozen or so characters are common for a few words. So not so common.

Character-level language model

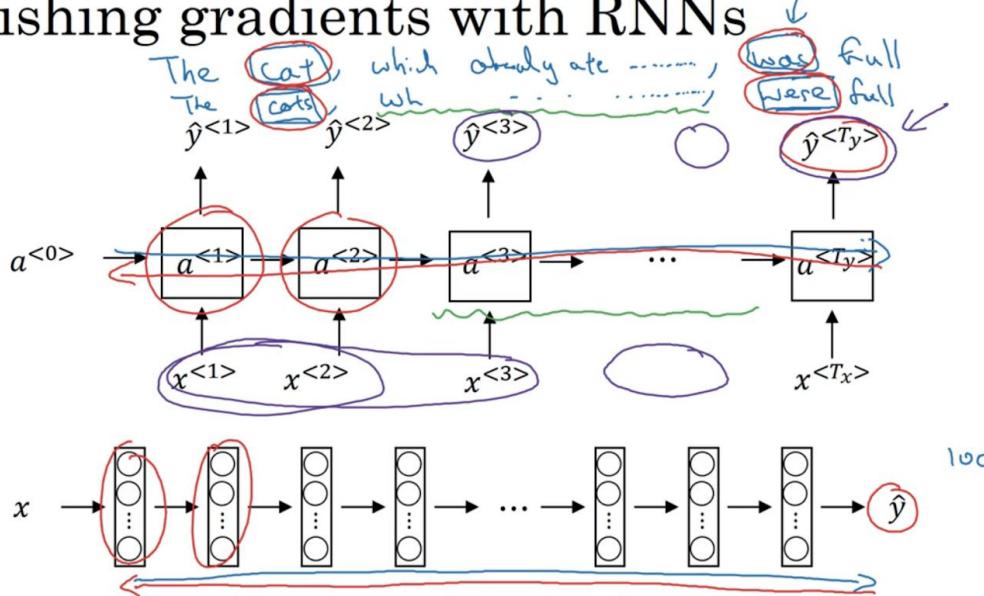
→ Vocabulary = [a, aaron, ..., zulu, <UNK>] ↵

$\rightarrow \text{Vocabulary} = [a, b, c, \dots, z, \cup, \circ, \rightarrow, ;, 0, \dots, 9, A, \dots, Z]$



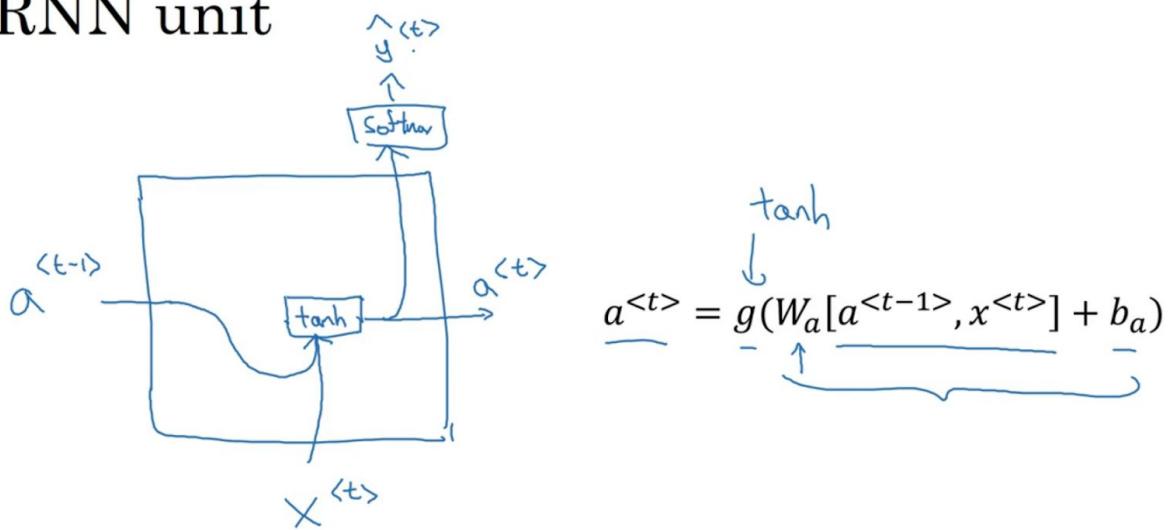
Gradient clipping can solve for gradient exploding. But vanishing gradient is harder to solve.

Vanishing gradients with RNNs



RNN Unit so far:

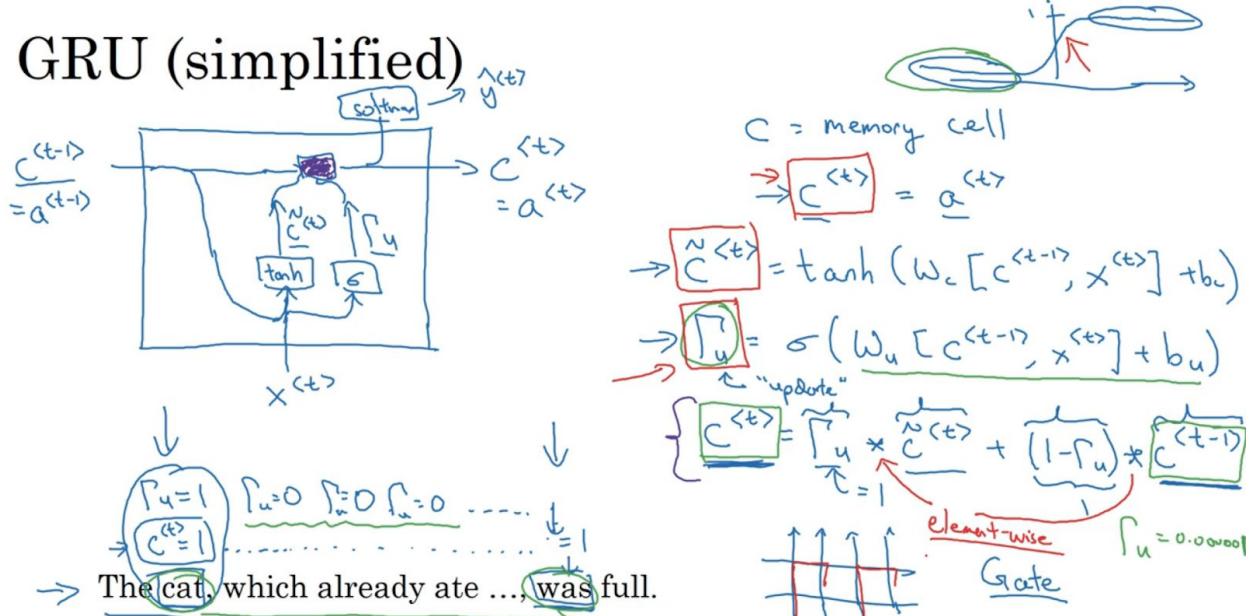
RNN unit



Basic difference is that you now create a candidate which gets either “chosen” or “ignored” depending on another parameter the Gate which also gets trained separately.

Gate. candidate + (1-Gate)* Previous candidate is used to decide the final cell state. This basically means that how much you should ignore this value compared to memorizing and continuing previous value.

GRU (simplified)



[Cho et al., 2014. On the properties of neural machine translation: Encoder-decoder approaches]
 [Chung et al., 2014. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling]

Andrew Ng

Because the gate and cells can be multiple dimensions, you can keep the cat in one dimension, food in another dimension etc.

In GRU the vanishing gradient problem does not exist because the gate can be very very small like 0.000001 basically just storing the previous value of cell state -- so no gradient vanishing!

Full GRU

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

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

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

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

Gamma R (r stands for relevance) is how much important is the current set of inputs and previous inputs at all. Its a very iterative approach to come up with these architectures - no fixed set. Its more empirical!

LSTM (more general form of GRU):

GRU and LSTM

GRU

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

$$\Gamma_u = \sigma(W_u[c^{$$

$$\Gamma_r = \sigma(W_r[c^{$$

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

\downarrow

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

$\uparrow \Gamma_e$

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

$$(\text{updat.}) \quad \Gamma_u = \sigma \left(W_u [a^{<t-1>} \times ^{<t>}] + b_u \right)$$

$$(\text{Forget}) \quad P_f = \phi(\omega_f [a^{<\leftrightarrow>} \times x^{<\leftrightarrow>}] + b_f)$$

$$P_0 = - (W_0 [a^{(t-1)}, x^{(t)}, \tilde{s}^t b_0])$$

$$a^{<\leftrightarrow} = \int_{\mathbb{R}} * c^{<\leftrightarrow}$$

c_{t-1} can also be used as peephole connection.

LSTM in pictures

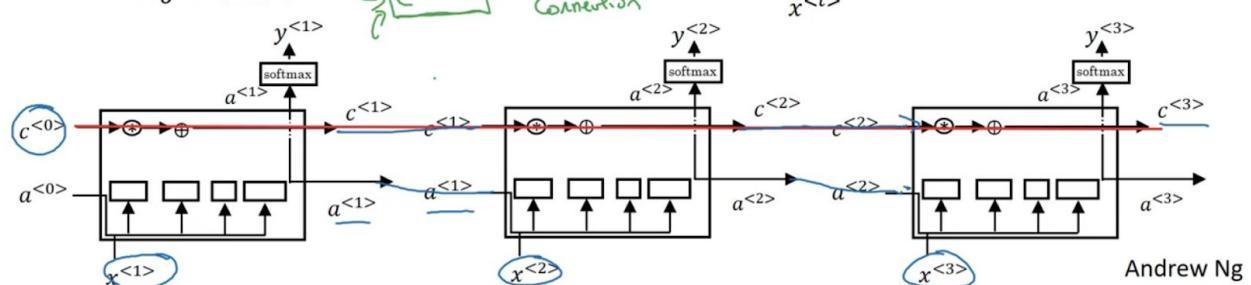
$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

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

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

$$\Gamma_o = \sigma(W_o [a^{$$

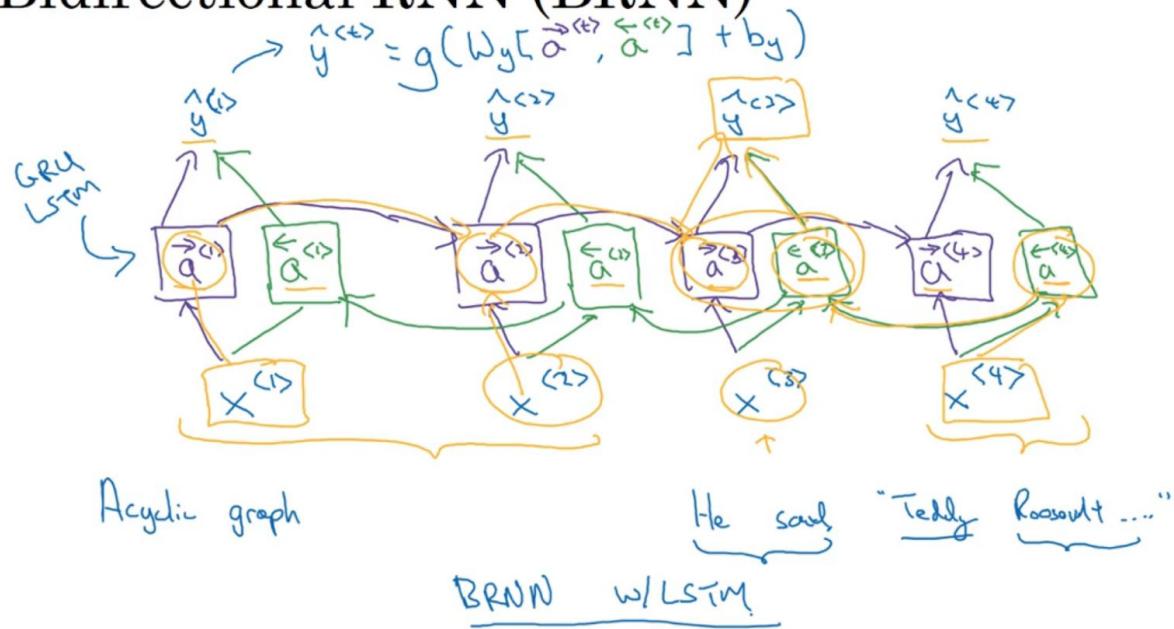
$$a^{<t>} \equiv \Gamma * \tanh c^{<t>}$$



Cell state flows through -- so there is no vanishing gradient!

LSTM is usually the first choice, but GRUs more and more often these days as there are only 2 gates compared to 3 gates so computationally more efficient.

Bidirectional RNN (BRNN)

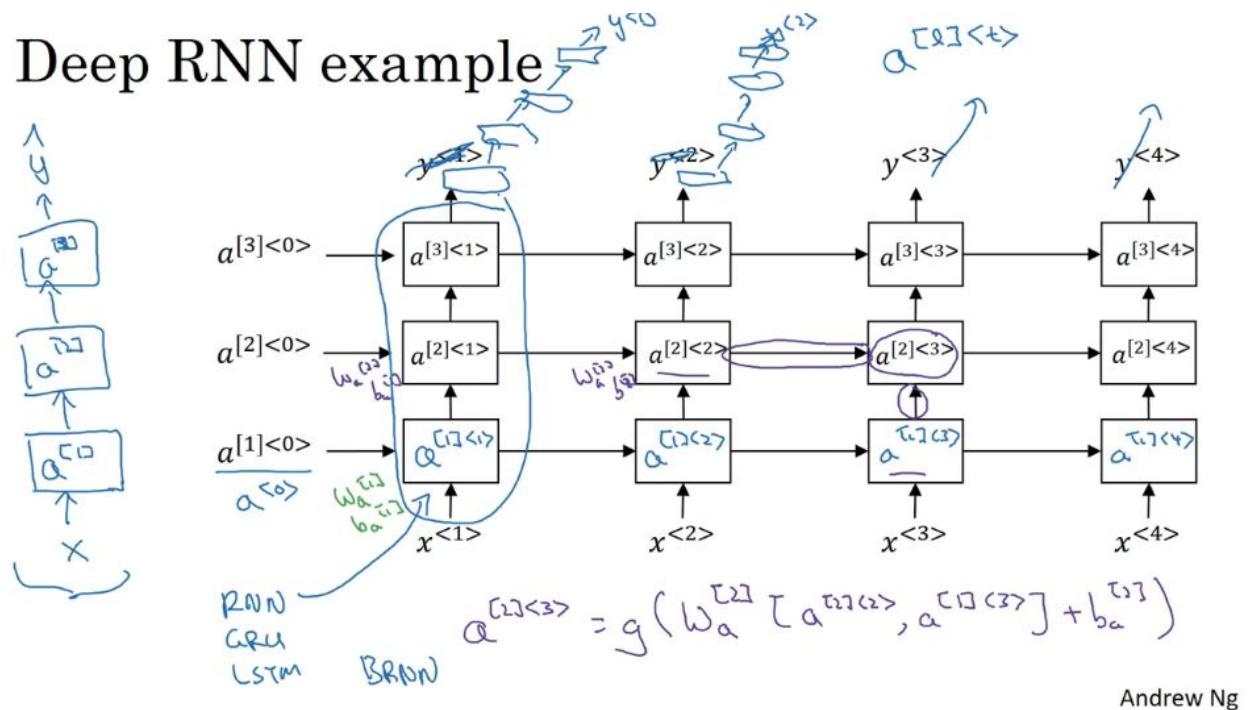


NLP tasks BRNN with LSTM is very common! Disadvantage is that you need entire sequence, you need to wait for person to stop talking before making a prediction. For complex speech recognition systems, you need more complex models.

For tasks where you can get entire sentence or sequence, BRNN is quite powerful and good!

DEEP RNN you can stack LSTM or GRU or RNN or BRNN and then hook them up in the end with a generic long deep network which does not have any temporal info.

Deep RNN example



WORD EMBEDDINGS (WEEK 2)

Word representation

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

$|V| = 10,000$

1-hot representation

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

Annotations show the vectors for 'Man' and 'Woman' being converted into their respective 1-hot representations o_{5391} and o_{9853} .

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

Andrew Ng

Can't transfer learning from orange juice to also apple juice with One Hot Encoding. Every word is different. Another disadvantage is that the word vectors are really very sparse.

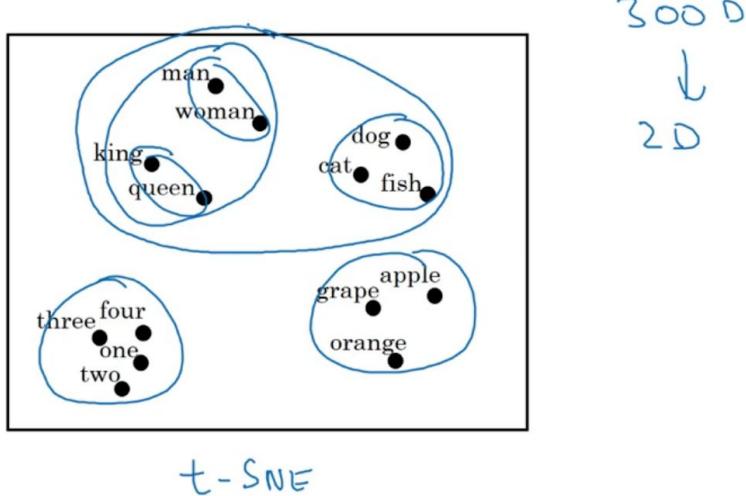
Featurized representation: word embedding

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	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 verb	⋮	⋮	⋮	⋮	I want a glass of orange juice.	
	e_{5391}	e_{9853}			I want a glass of apple juice.	Andrew Ng

Above instead of one hot -- you are using features. So 100 dimensional vector is basically 100 features like gender, age, food, size, color etc. and a dot product of the two (OR some other analogy function like subtraction _ more later_) can give "nearness".

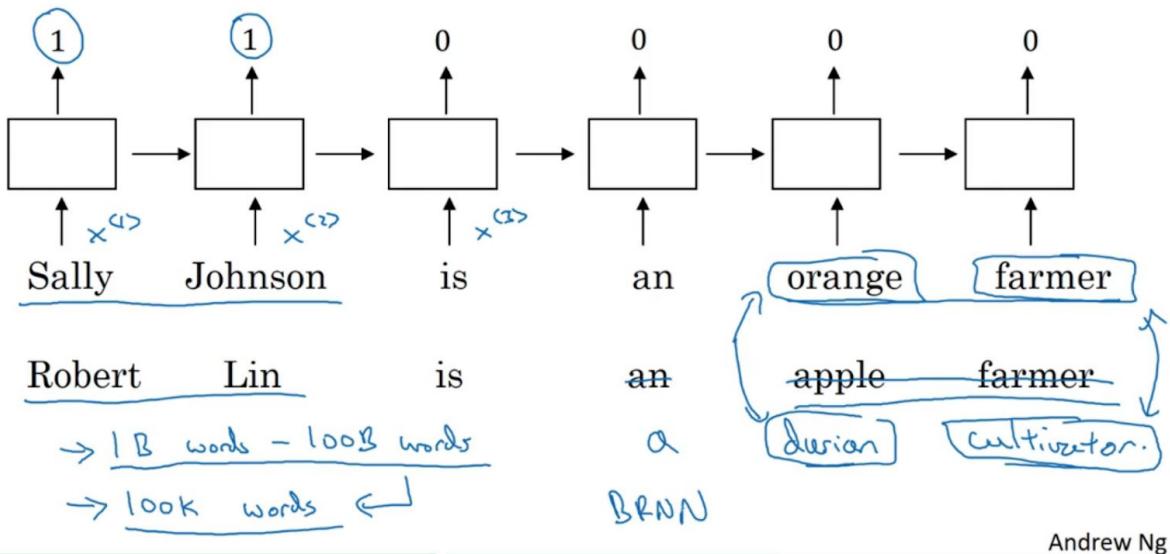
So apple and juice while they differ in colors, might have mostly similar features and hence can transfer learning about orange juice into apple juice.

Visualizing word embeddings



300 D to 2D so orthogonal dimension in same space

Named entity recognition example



TRANSFER LEARNING in RNNs is about using someone else's word embeddings from 1B or 100B words and use it for your scenario of 100k or so words. In real life you should use BRNN to get full context.

Transfer learning and word embeddings

1. Learn word embeddings from large text corpus. (1-100B words)

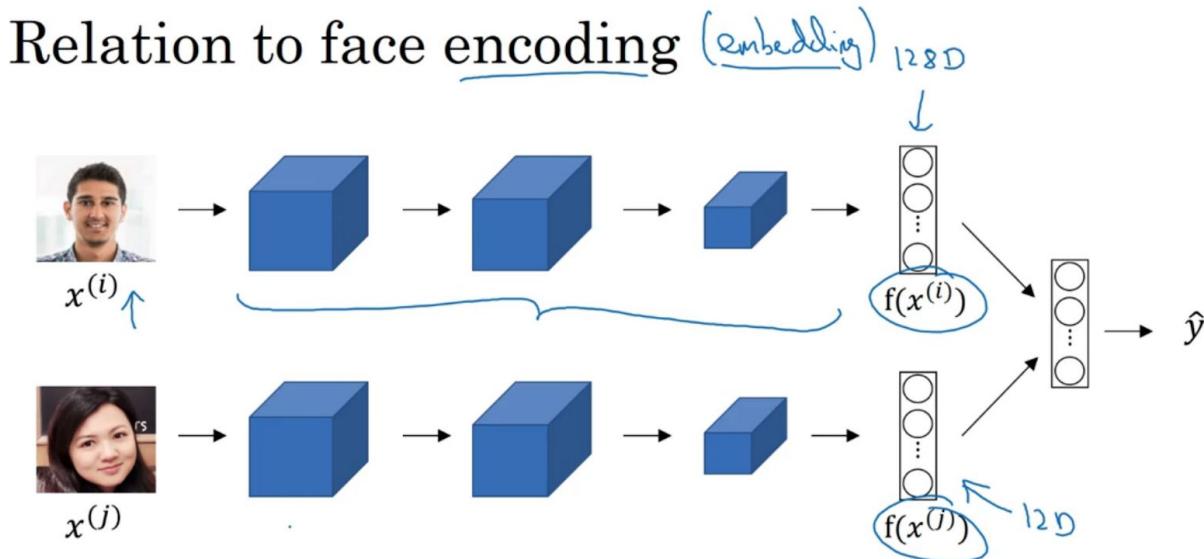
(Or download pre-trained embedding online.)

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.

If Step 2 has many many training data, then can do Step 3.

Word embeddings is useful for NLP tasks like Named Entity Recog, text summarization, coreference resolution, parsing etc. But not very useful for language modeling, machine translation (esp. if there is a lot of data)



IN CNN literature, you can compare if the two photos are exactly the same. In faces, you can continue to do this for unseen images. But for text its a fixed vocabulary, everything else is just UNKNOWN word.

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

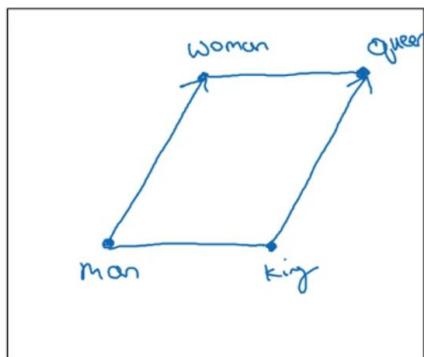
$e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\text{?}}$
 $e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\text{queen}}$

Mikolov et. al., 2013, Linguistic regularities in continuous space word representations]

Andrew Ng

Really influential paper!

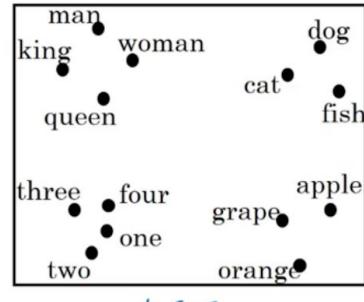
Analogy using word vectors



300D

Find word w: $\arg \max_w$

$300D \rightarrow 2D$



t-SNE

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

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

30 - 75%

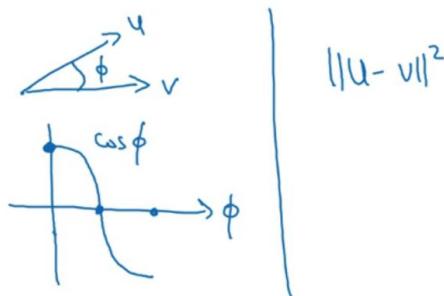
Andrew Ng

Please -- the parallelogram is more common in 300D Space! Not in TSNE!! Refer to original d dimensional space for cosine similarity

Cosine similarity

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

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



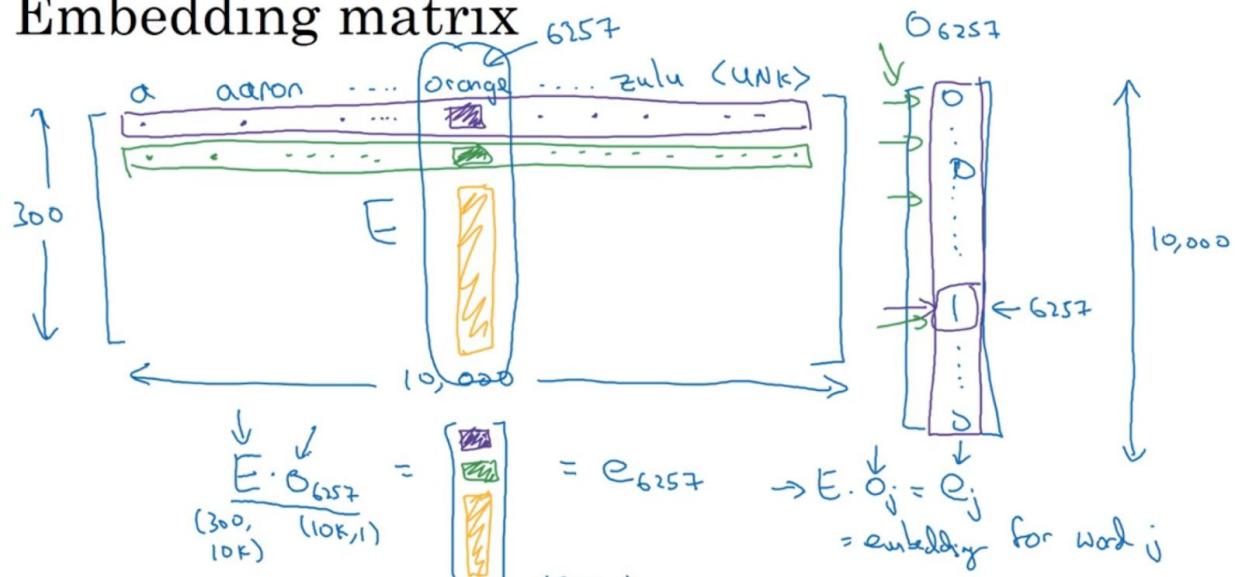
Man:Woman as Boy:Girl

Ottawa:Canada as Nairobi:Kenya

Big:Bigger as Tall:Taller

Yen:Japan as Ruble:Russia

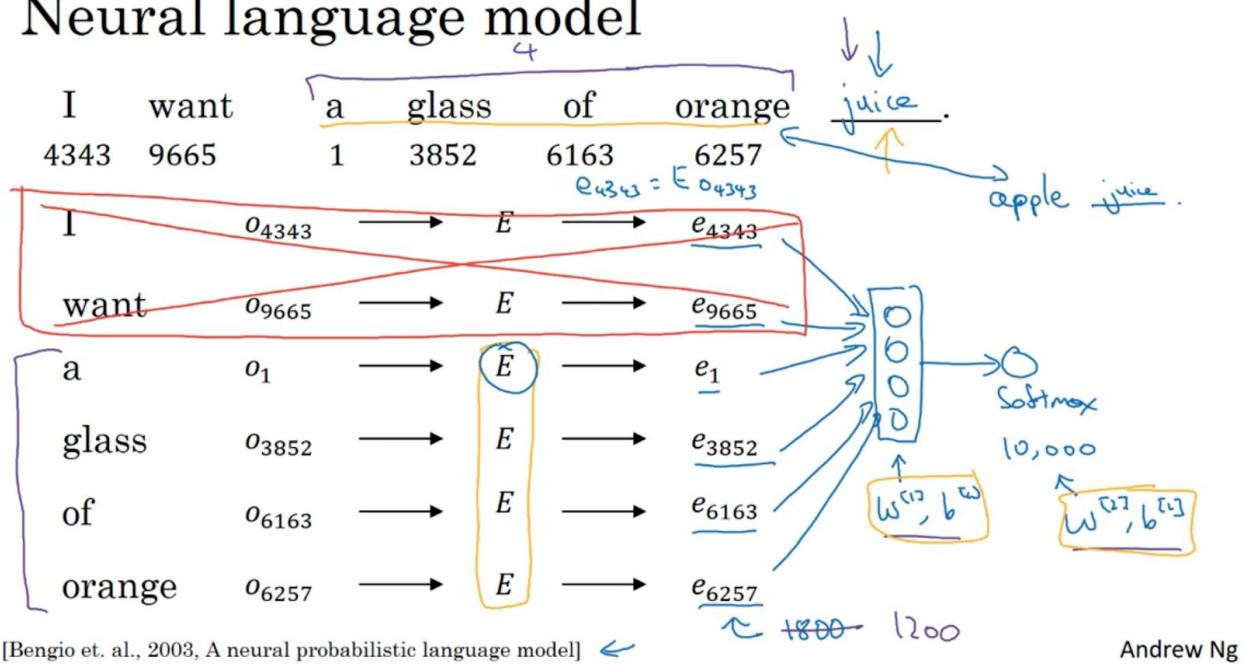
Embedding matrix



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

IN reality, you use Embedding layer in Keras which just calls the column you want instead of doing this so expensive calculation.

Neural language model



[Bengio et. al., 2003, A neural probabilistic language model]

Andrew Ng

Use backprop to learn E and the weights. Stack all the e vectors - pass to deep layer and then use softmax to predict probabilities of next word.

Other context/target pairs

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

Context: Last 4 words.

4 words on left & right

Last 1 word

Nearby 1 word

skip gram

a glass of orange ? to go along with

orange ?

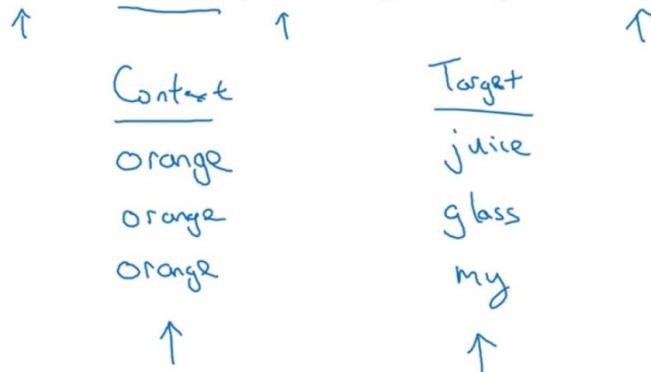
glass ?

Language Modelling - last 4 words is best! Like Probability of a sentence, pear or pair salad.

But just for word embeddings all the rest of the context vectors are also good enough. {NEED TO LEARN MORE HERE, SOME DOUBTS, why does softmax layer also need weights? Is this above method good enough to train E for first time or for refinement?}

Skip-grams

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



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

Andrew Ng

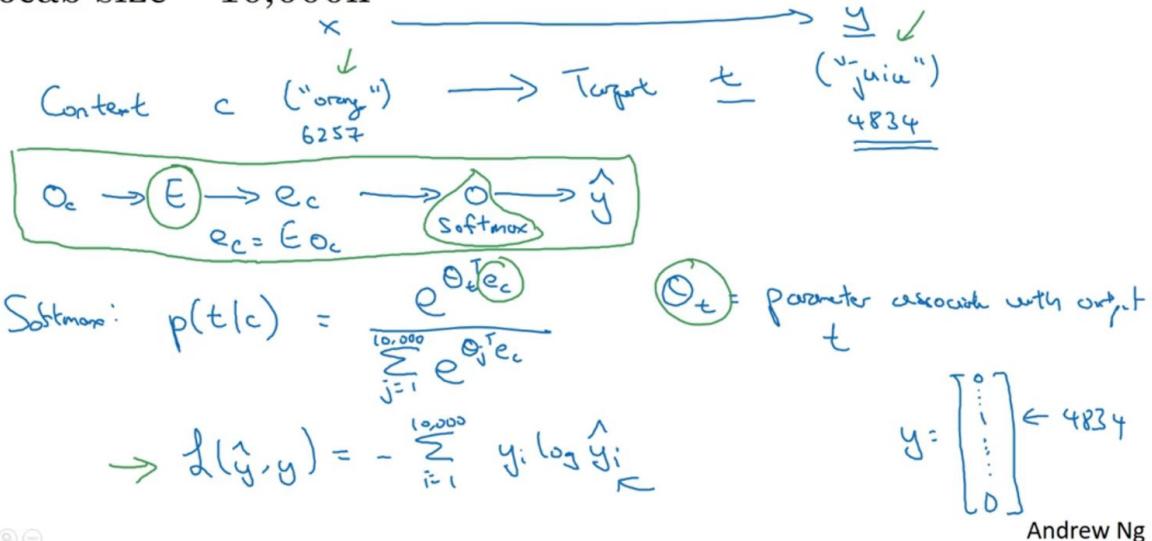
Randomly chosen context word in a window and a randomly chosen target word again in a window

SKIP GRAM MODEL:

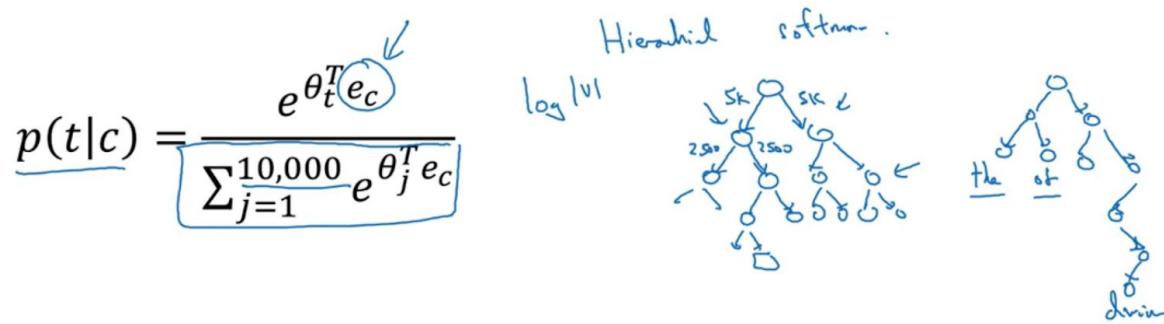
In stanford class CS224n, THETA there was v, and word emdeddings was u. Here its theta, and e. Every word has actually 2 vectors associated with them, one as context and one as target. When its the context, thaths the word embeddings.

Model

Vocab size = 10,000k



Problems with softmax classification



How to sample the context c ?

→ the, of, a, and, to, ...

→ orange, apple, durian

Q_{durian}

t

$c \rightarrow t$

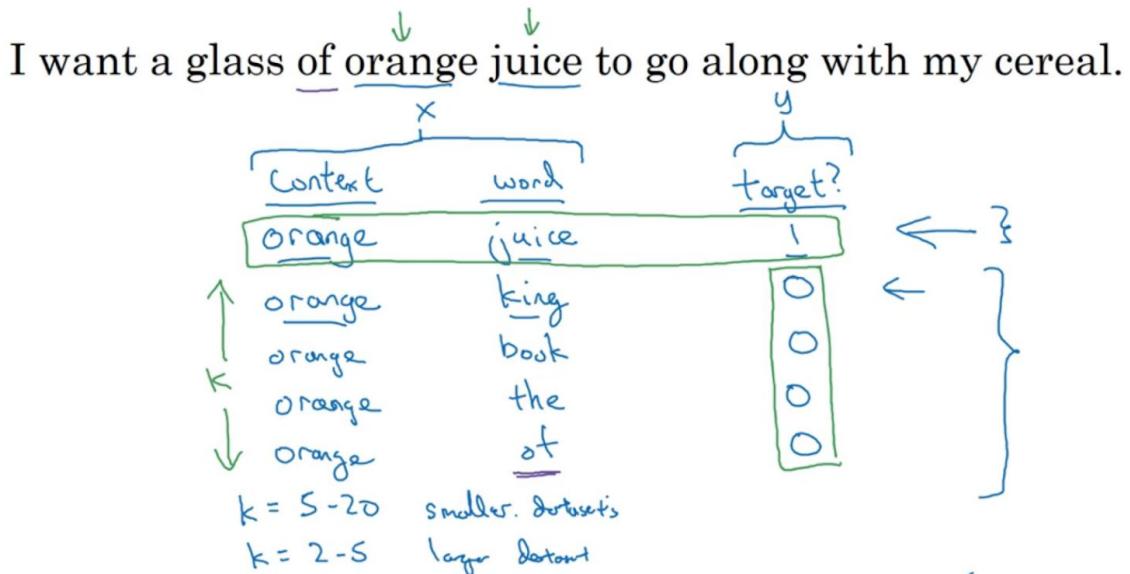
$P(c)$

Hierarchical softmax doesn't use softmax but just sigmoid logistic regression stuff to predict whether the target word is in left half or right half of vocab and so on. And hence $\log|V|$ complexity. In reality, more common words like the, of etc. are on the left and uncommon words like durian buried deep.

In terms of sampling context word c, dont want too many too, the, of, a etc. so there are ways to sample such that you spend some time on orange, apple etc. too.

NEGATIVE SAMPLING

Defining a new learning problem



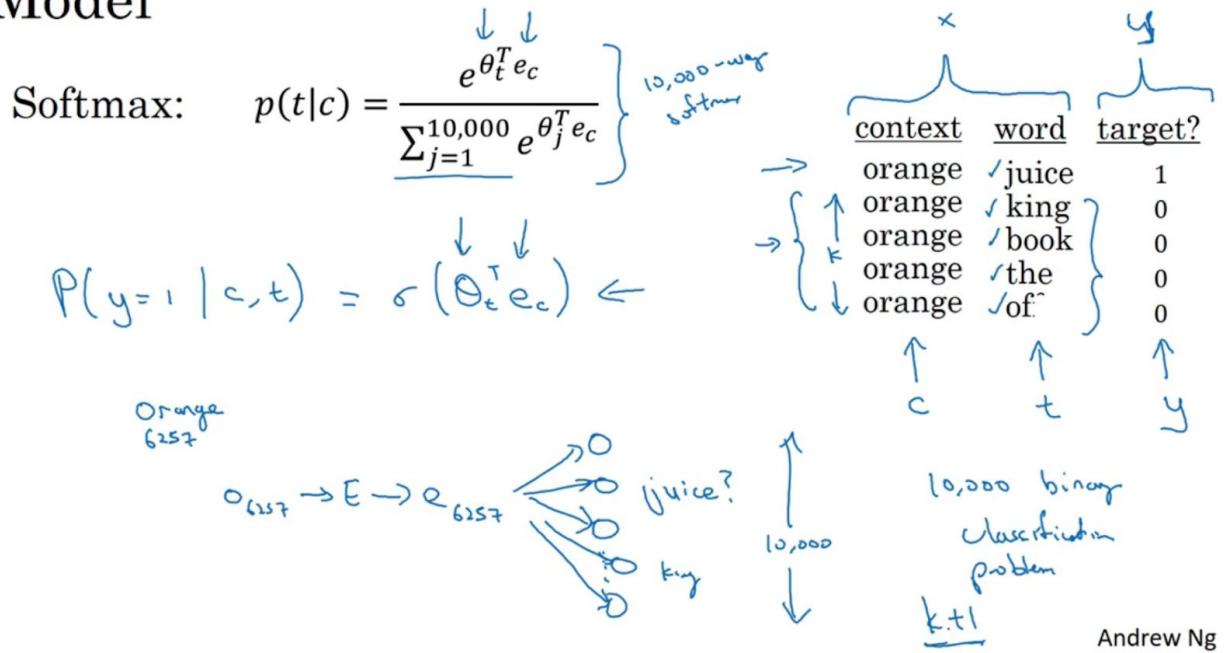
[Mikolov et. al., 2013. Distributed representation of words and phrases and their compositionality]

Andrew Ng

Orange from the window, and one word from window for target 1 , and some negative samples with one word from window and the rest from randomly chosen words in vocabulary.

Instead of 10000 softmax algorithm, you train $k + 1$ logistic regression models, so its much computationally easier.

Model



Selecting negative examples

the, of, and, ...

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

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

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

Neither most common negative words and nor least occurring, somewhere in between is better!
that's why 3/4th the freq of words

GloVe (global vectors for word representation)

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

c, t

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

\uparrow \uparrow \uparrow
 c t c

$$x_{ij} = x_{ji} \leftarrow$$

GLOVE uses counts of number of times a word appears in context of word j in vocab.

Model

$$\text{minimize} \quad \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(x_{ij}) (\Theta_i^T e_j + b_i + b_j - \log \underline{x_{ij}})^2$$

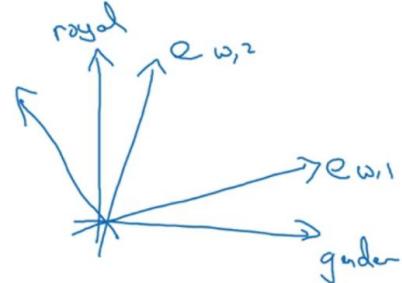
\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow
 f(x_{ij}) $\Theta_i^T e_j$ b_i b_j $-\log \underline{x_{ij}}$?

\uparrow \uparrow \uparrow \uparrow
 "weighty" term " $\Theta_i^T e_c$ " "0 log 0" = 0

$\Theta_i^T e_c$ Θ_i, e_i are symmetric
 this, is, at, a, ... $e_w^{(\text{final})} = \frac{e_w + \Theta_w}{2}$

A note on the featurization view of word embeddings

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



$$\text{minimize } \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) (\underbrace{\theta_i^T e_j + b_i - b'_j - \log X_{ij}}_0)^2$$

$\leftarrow (\mathbf{A}\boldsymbol{\theta}_i)^T (\mathbf{A}^T e_j) = \boldsymbol{\theta}_i^T \cancel{\mathbf{A}^T \mathbf{A}} e_j$

Andrew Ng

Not always human interpretable, but vector parallelogram math still works!

APPLICATIONS OF WORD EMBEDDINGS

Sentiment classification problem

$$x \longrightarrow y$$

The dessert is excellent.



Service was quite slow.



Good for a quick meal, but nothing special.



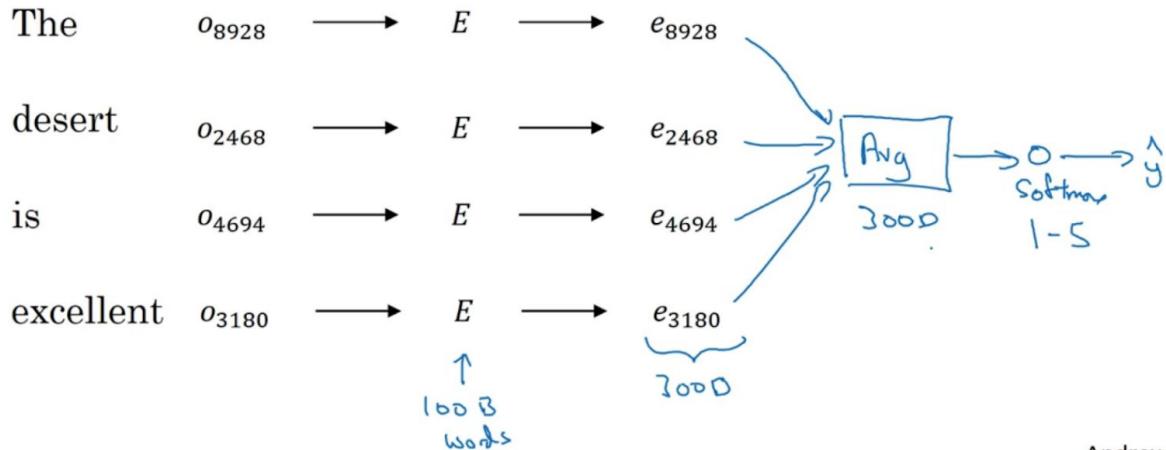
Completely lacking in good taste, good service, and good ambience.



You might not have a lot of training examples. Word embeddings can help if your training set is low.

Simple sentiment classification model

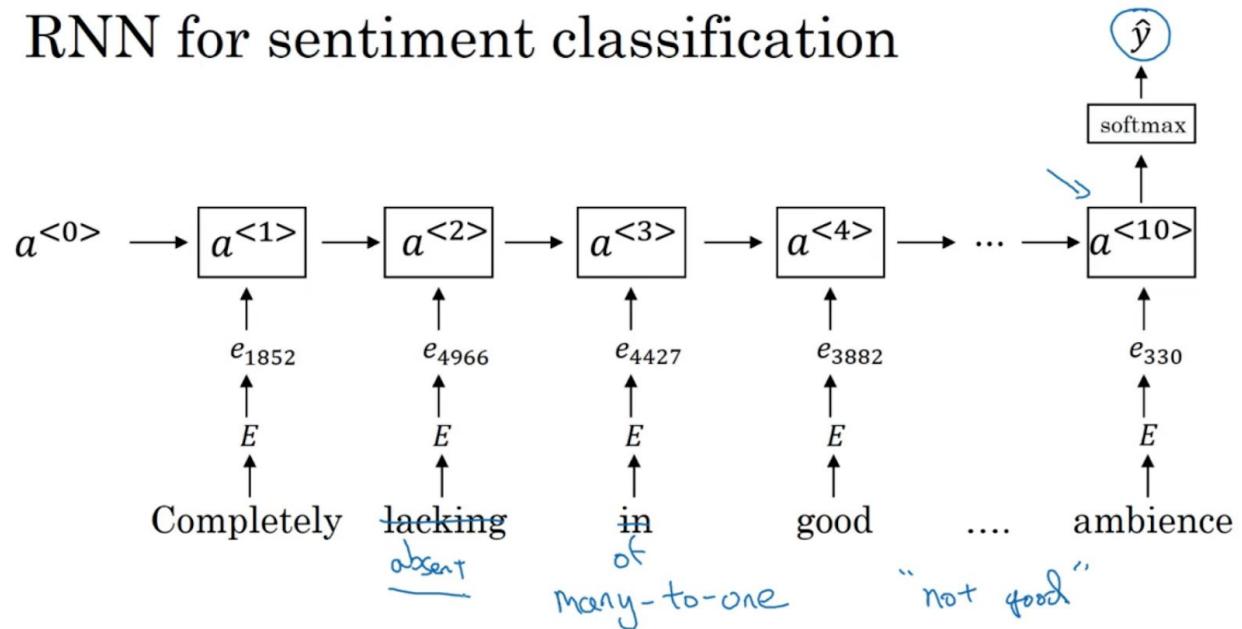
The dessert is excellent
8928 2468 4694 3180



Andrew Ng

Average is good because of 100 or 1000 word reviews, it just averages them and you can softmax the output. But there is no context!

RNN for sentiment classification



This is better for context!

The problem of bias in word embeddings

Man:Woman as King:Queen

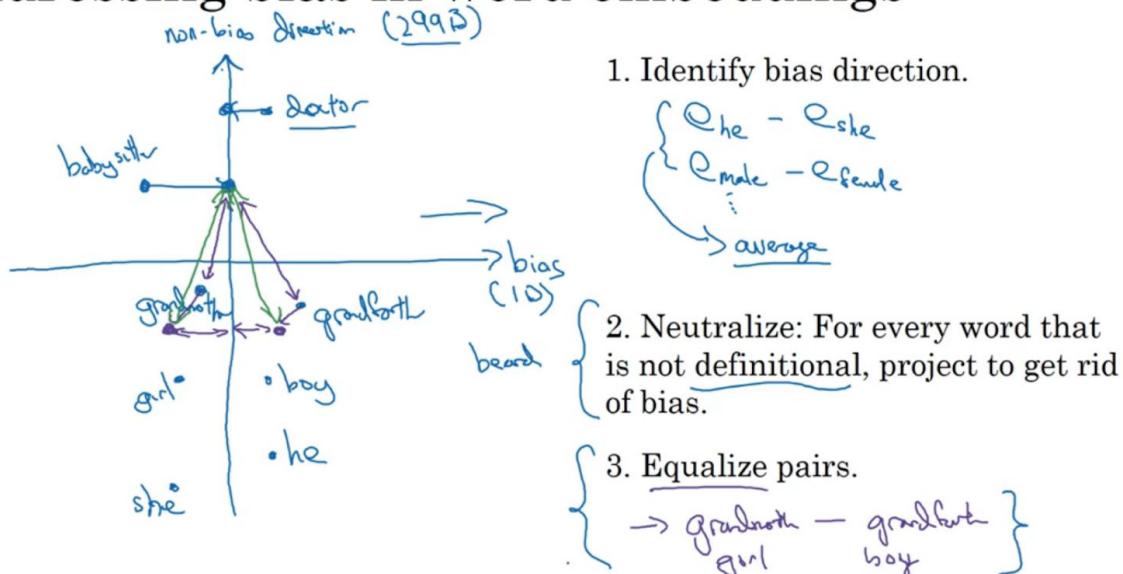
Man:Computer_Programmer as Woman:Homemaker

Father:Doctor as Mother:Nurse

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

[Bolukbasi et. al., 2016. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings] Andrew Ng

Addressing bias in word embeddings



[Bolukbasi et. al., 2016. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings] Andrew Ng

SEQUENCE to SEQUENCE MODELS (WEEK 3)

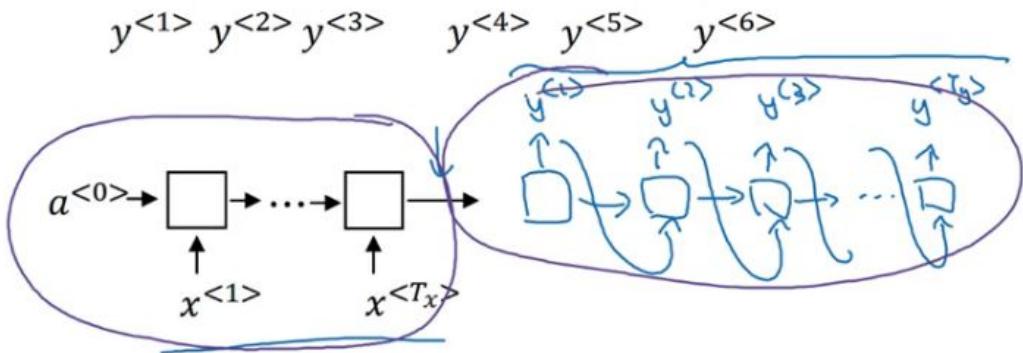
Encoder is a must for many to many (different lengths) RNN.

Encoder to Decoder

Sequence to sequence model

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

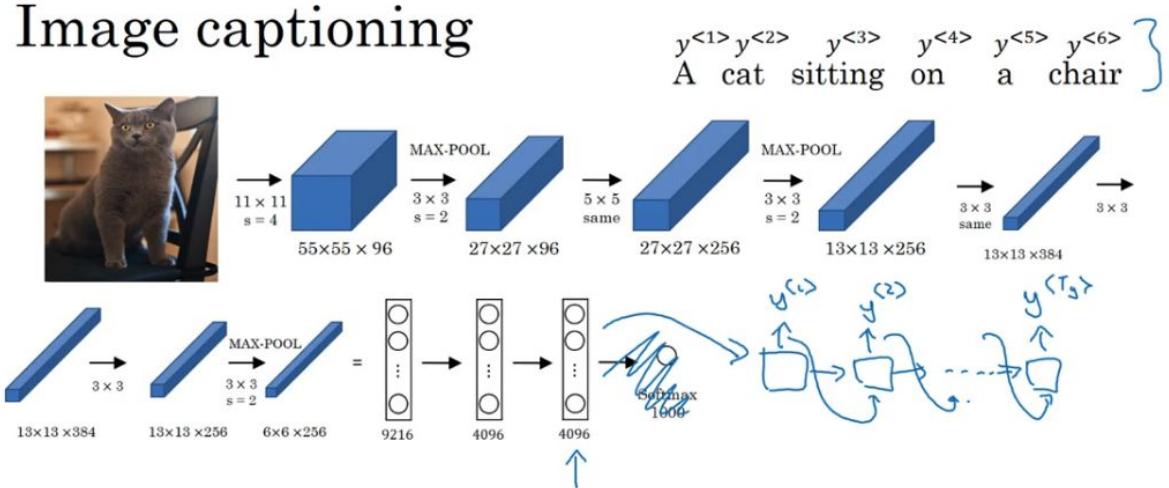
→ Jane is visiting Africa in September.



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

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

Image captioning



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

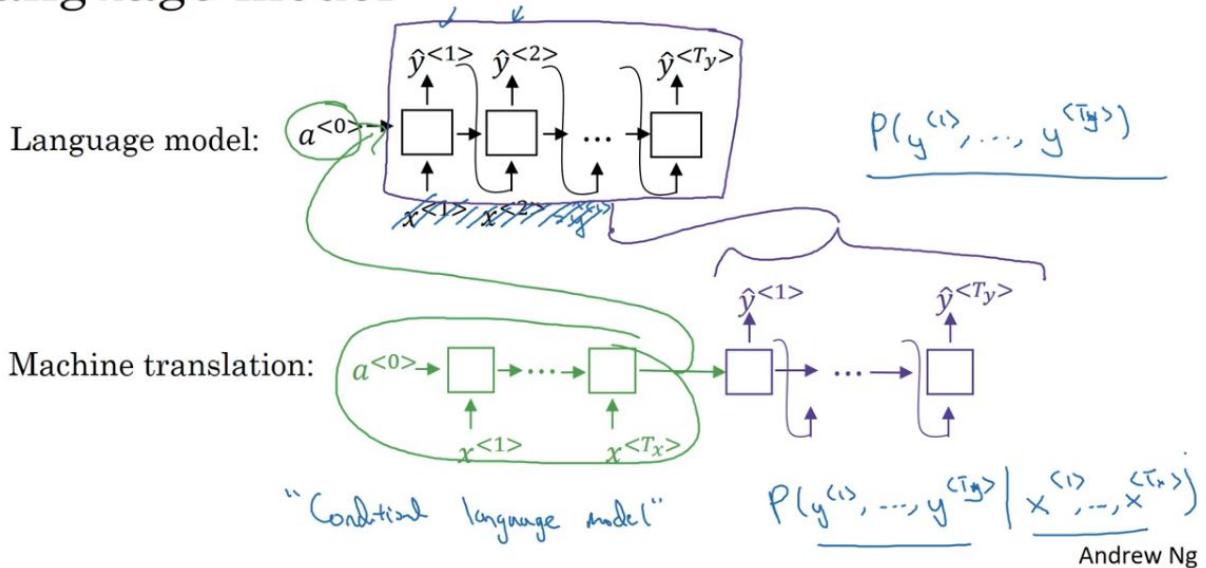
[Vinyals et. al., 2014. Show and tell: Neural image caption generator]

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

Andrew Ng

Probability of a general sentence vs conditional probability given a French sentence

Machine translation as building a conditional language model



Finding the most likely translation

Jane visite l'Afrique en septembre.

$$\underbrace{P(y^{<1>} , \dots , y^{<T_y>} | x)}_{\substack{\text{English} \\ \downarrow \\ \text{French}}}$$

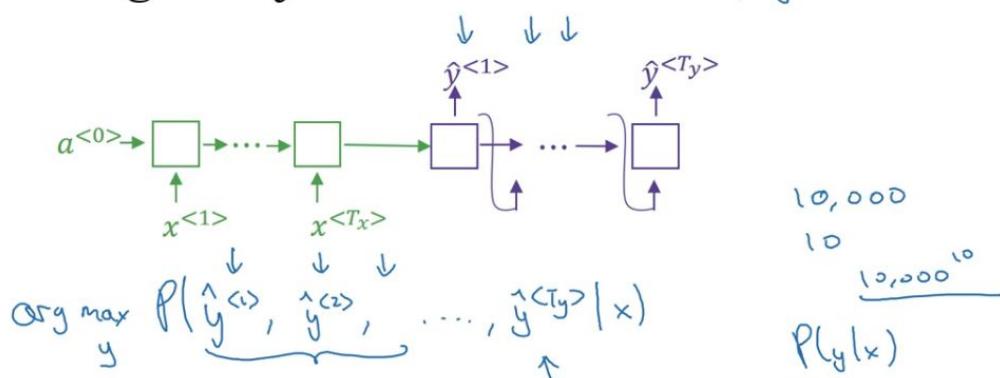
- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in September.
- In September, Jane will visit Africa.
- Her African friend welcomed Jane in September.

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

CANT SAMPLE AT RANDOM, you want arg max of the probability.

Greedy Search: Pick first word $P(y_1|x)$ and then pick the best 2nd word and then best 3rd word. Does Not work!

Why not a greedy search?



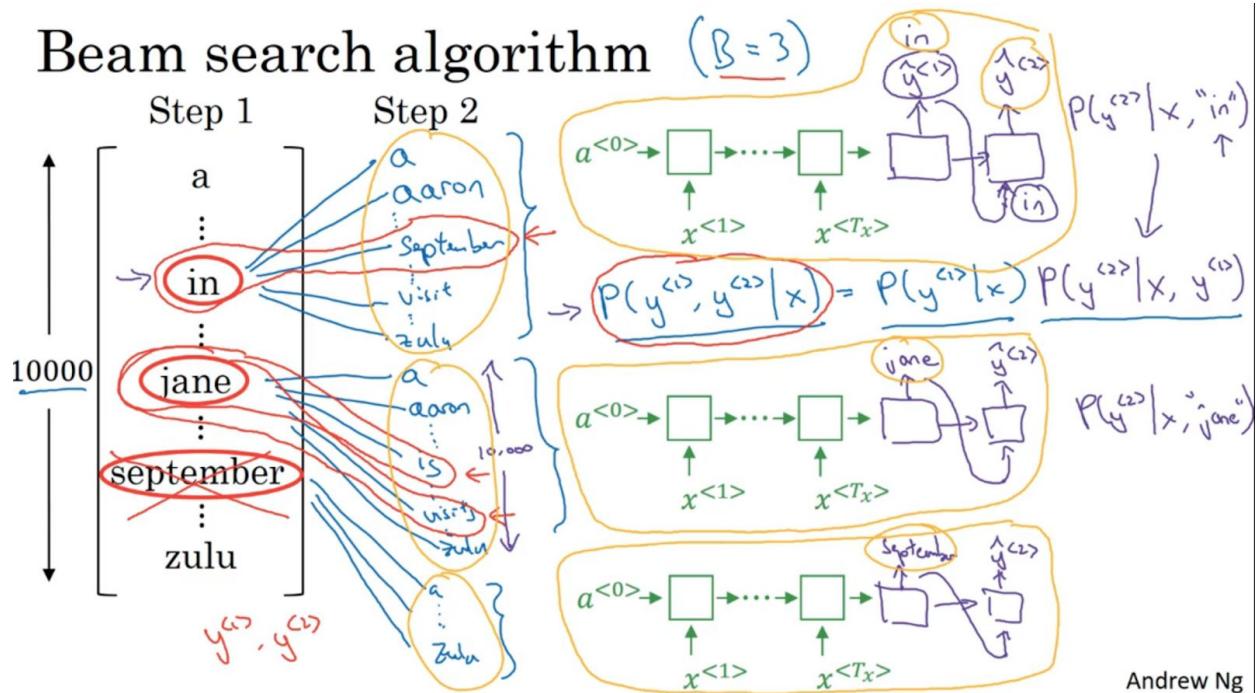
- Jane is visiting Africa in September.
 - Jane is going to be visiting Africa in September.
- $P(\text{Jane is going } | x) > P(\text{Jane is visit } | x)$

Andrew Ng

Therefore we need beam search that picks max probability of a sentence together!

BEAM SEARCH

So instead of exploring 10000^10 possibilities for 10 word sentence, beam search works as follows. It finds top 3 candidates for first word, and then for each word looks at 10k words each for 2nd word. Now it picks top 3 candidates for $P(y_2, y_1|x)$ by multiplying $P(y_1|x)$ from the top 3 list and then new probability. It then picks top 3 pairs again and so on. So the three pairs might be In September, Jane is, Jane goes.

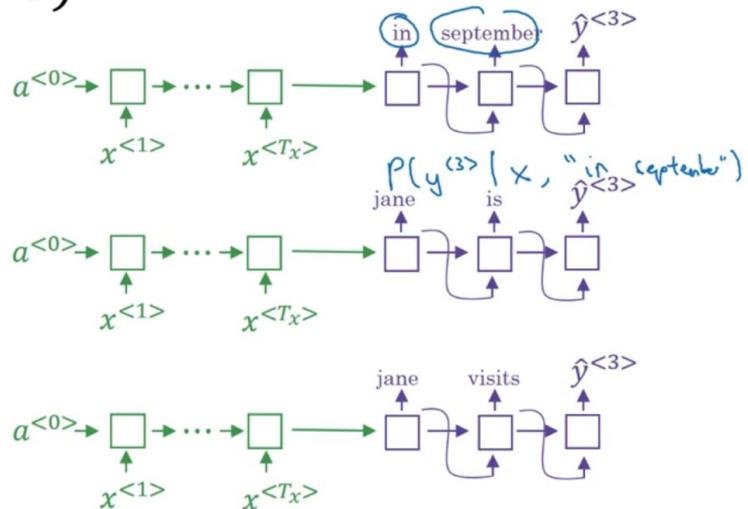


Beam search ($B = 3$)

in september
 aaron
 jane
 zulu

jane is
 visits
 zulu

jane visits
 africa
 zulu



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

jane visits africa in september. <EOS>

Andrew Ng

If $B=1$ then this becomes greedy search algorithm. If $B=3$ or so it can be better.

Instead of product of probabilities, its better to take sum of p's.

You should also take average of probas as you dont want to penalize longer sentences.

Length normalization

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

\log

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

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

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

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

Andrew Ng

Alpha of 0.7 works pretty well!

Beam search discussion

Beam width B?

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

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

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

B of 3 to 10 is common for commercial applications, for research from 100 to 3000.

How to know if RNN or Beam Search is at fault?

Error analysis on beam search

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

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

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

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

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

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

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

Conclusion: Beam search is at fault.

Case 2: $P(y^*|x) \leq P(\hat{y}|x)$ \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.	2×10^{-10}	1×10^{-10}	B
...	...	—	—	R
...	...	—	—	Q
				R R R :)

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

Andrew Ng

If RNN is at fault, then add training data, regularization or try other architectures.

Evaluating machine translation

French: Le chat est sur le tapis.

→ Reference 1: The cat is on the mat. ←

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

→ MT output: the the the the the the the

Precision: $\frac{7}{7}$ Modified precision: $\frac{2}{7}$ ← Count ("the")

Count clip ("the") ← Count ("the")

Bleu
bilingual evaluation under study

Only attribute the twice or the maximum number of times it occurs in any one reference.

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. ↪

	Count	CountClip	
the cat	2 ↪	1 ↪	
cat the	1 ↪	0	
cat on	1 ↪	1 ↪	
on the	1 ↪	1 ↪	
the mat	1 ↪	1 ↪	

$\frac{4}{6}$

Only credit bigram for the maximum number of times it occurs in the reference sentence.

Bleu score on unigrams

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. (↑)

$$P_1 = \frac{\sum_{\text{unigrams} \in \hat{y}} \text{Count}_{\text{clip}}(\text{unigram})}{\sum_{\text{unigrams} \in \hat{y}} \text{Count}(\text{unigram})} \quad P_n = \frac{\sum_{n\text{-grams} \in \hat{y}} \text{Count}_{\text{clip}}(n\text{-gram})}{\sum_{n\text{-grams} \in \hat{y}} \text{Count}(n\text{-gram})}$$

Bleu details

p_n = Bleu score on n-grams only

P_1, P_2, P_3, P_4

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

BP = brevity penalty

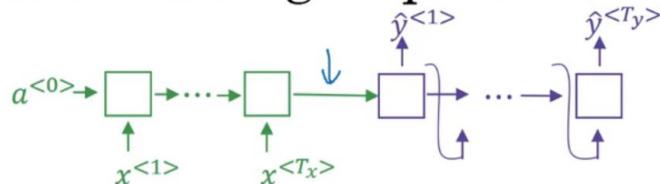
$$\text{BP} = \begin{cases} 1 & \text{if } \text{MT_output_length} > \text{reference_output_length} \\ \exp(1 - \text{MT_output_length}/\text{reference_output_length}) & \text{otherwise} \end{cases}$$

Sometimes RNN can make very small sentences and can get high precision ! SO penalise if short sentences! And not penalize if sentences are bigger than references!

ATTENTION MODEL

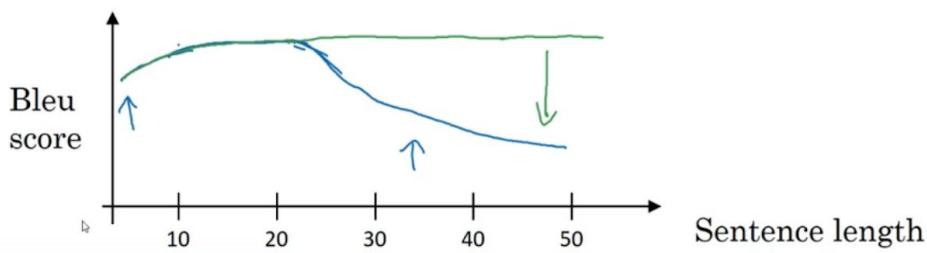
Hard to memorize everything in encoder -- so need something better

The problem of long sequences



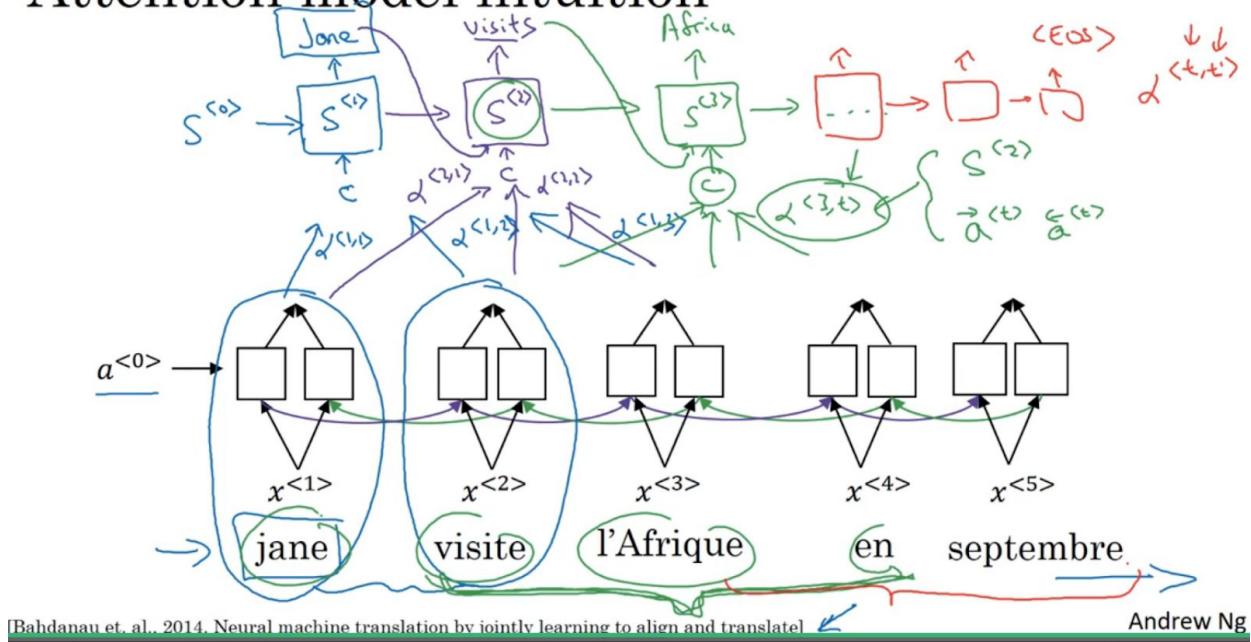
Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente d'y aller aussi.

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.



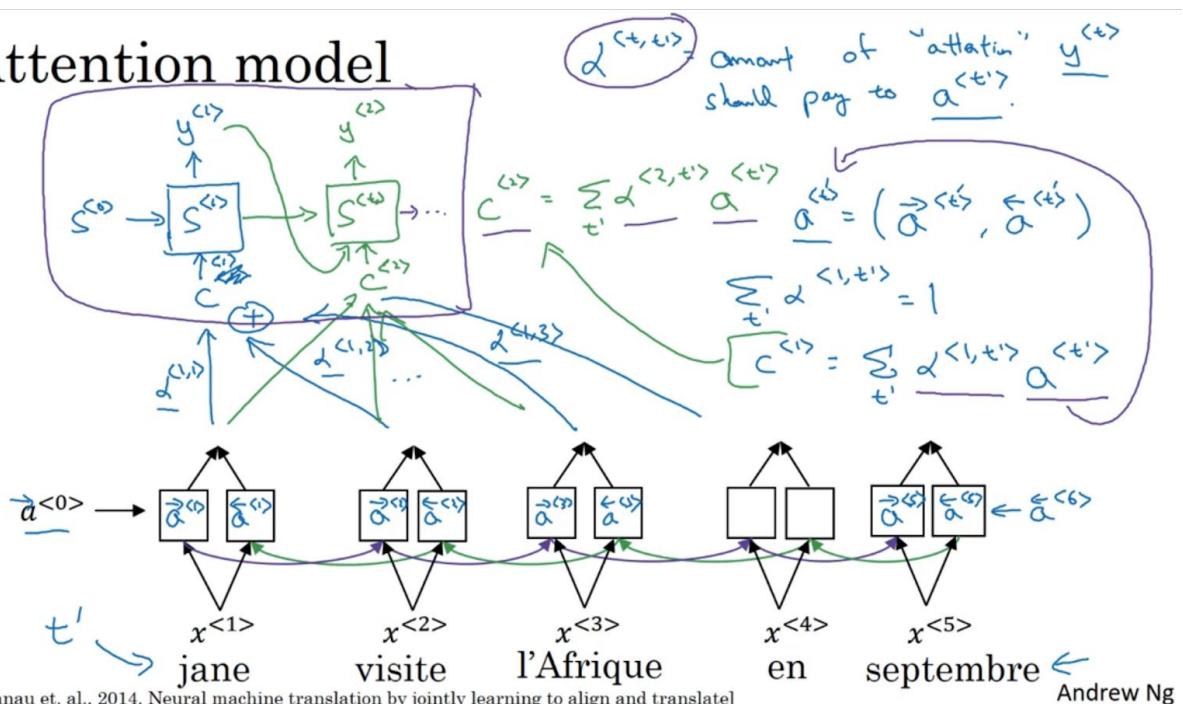
Andrew Ng

Attention model intuition

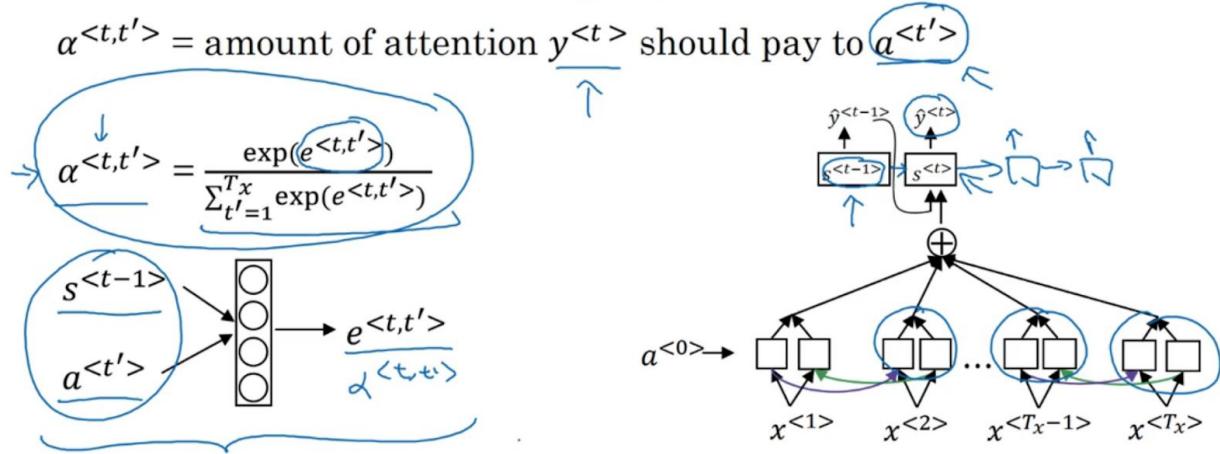


Alpha decides the context, which few words to consider while translating the sentence at the current step.

Attention model



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



[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate] \leftarrow
 [Xu et. al., 2015. Show, attend and tell: Neural image caption generation with visual attention] \leftarrow

Small neural network is used to train the weights of activations for each step.

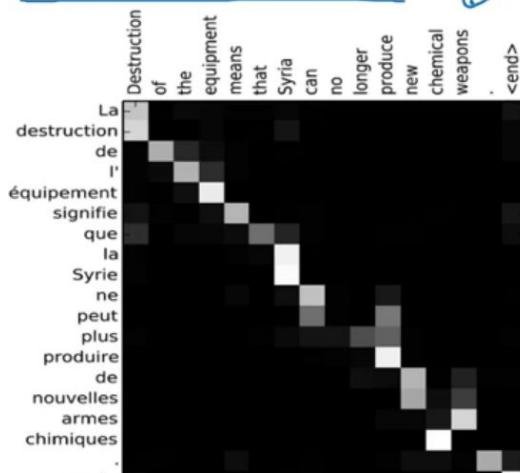
One of the papers above show how we can use attention while doing image captioning too -- paying attention to part of images while doing captioning work

Attention examples

July 20th 1969 → 1969 – 07 – 20

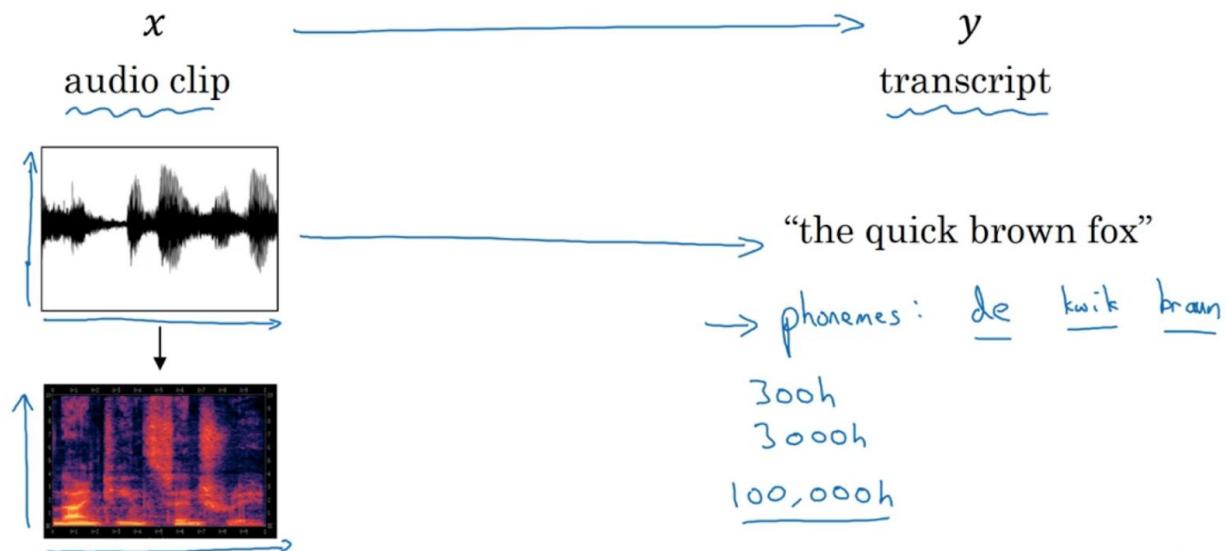
23 April, 1564 → 1564 – 04 – 23

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

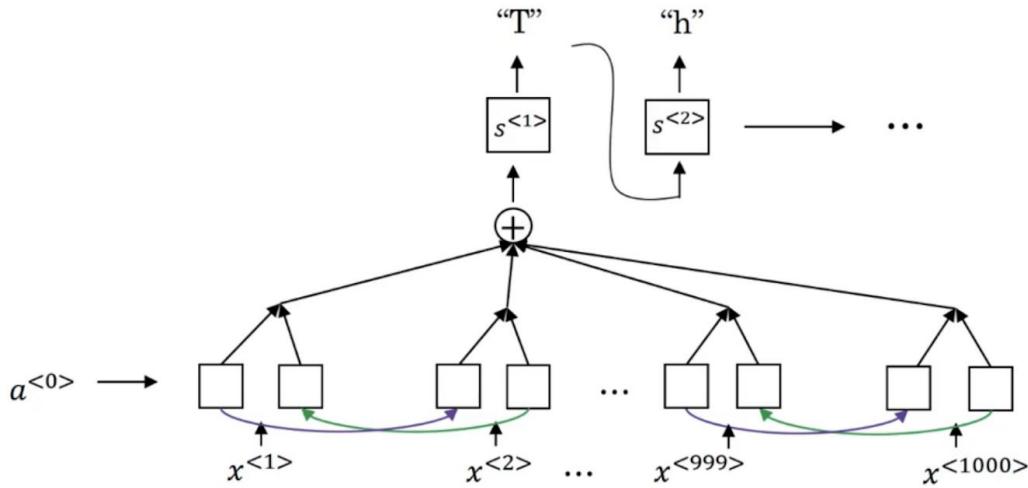


SPEECH RECOGNITION

Speech recognition problem

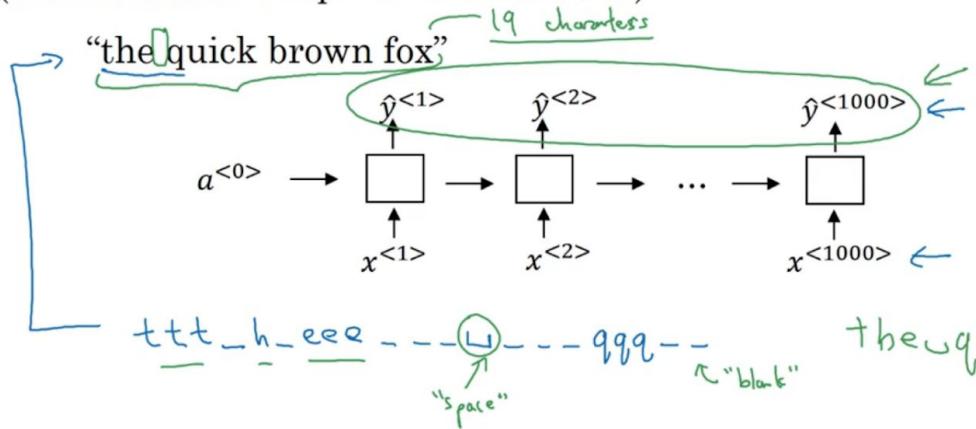


Attention model for speech recognition



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

What is trigger word detection?



Amazon Echo
(Alexa)



Baidu DuerOS
(xiaodunihao)

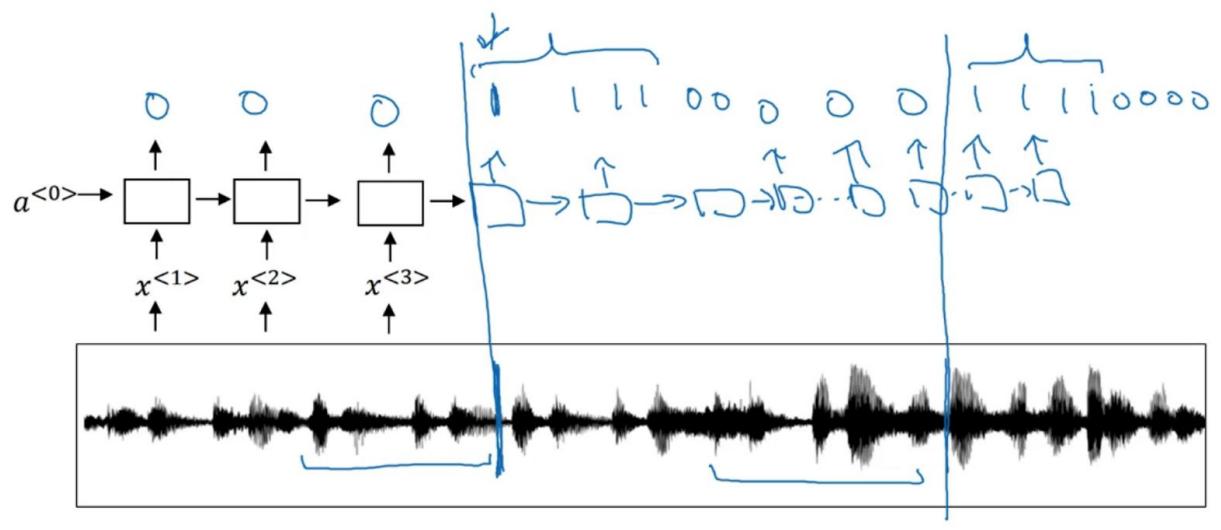


Apple Siri
(Hey Siri)



Google Home
(Okay Google)

Trigger word detection algorithm



Andrew Ng

