

# Copyright Notice

These slides are distributed under the Creative Commons License.

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

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



deeplearning.ai

# Sequence to sequence models

---

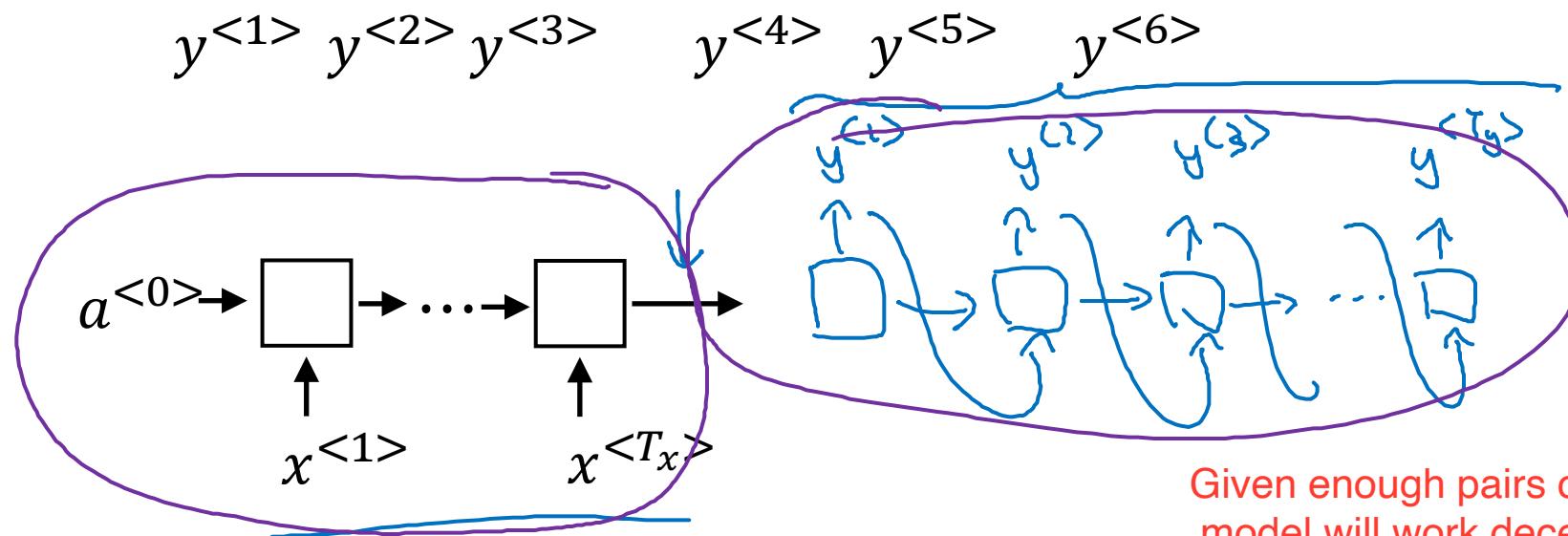
## Basic models

# Sequence to sequence model

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

Jane visite l'Afrique en septembre

→ Jane is visiting Africa in September.



Encoder network is built as a RNN and this could be a Gru or LSTM feeding the input, French words, one word at a time. After ingesting the input sequence the RNN then outputs a vector that represents the input sequence.

Decoded network takes as input the encoding output by the encoding network shown in blank on the left, and then can be trained to output the translation one word.

Given enough pairs of French and English sentences, this model will work decently well. This model simply uses an encoding network whose job is to find an encoding of the input French sentence, and then use a decoding network in order to generate the corresponding English translation.

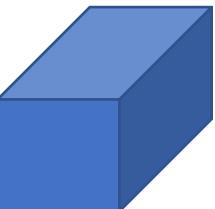
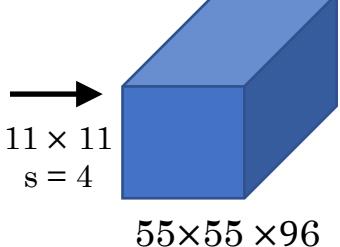
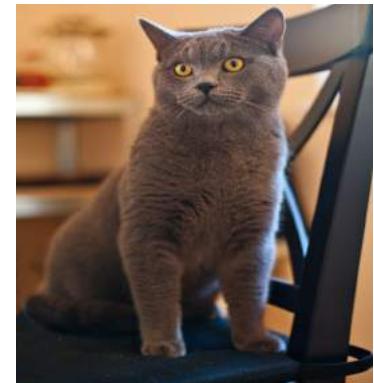
[Sutskever et al., 2014. Sequence to sequence learning with neural networks]

This architecture is very similar to image captioning.

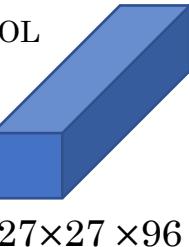
[Cho et al., 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation]

Andrew Ng

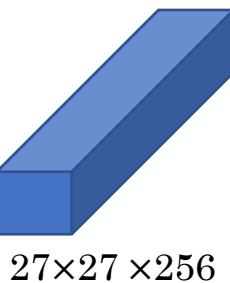
# Image captioning



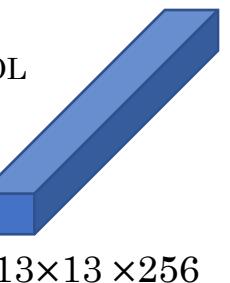
MAX-POOL  
 $3 \times 3$   
 $s = 2$



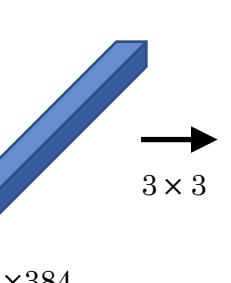
$5 \times 5$   
same



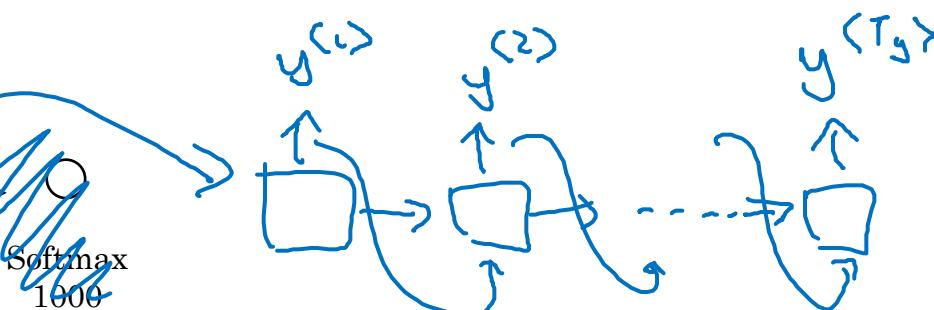
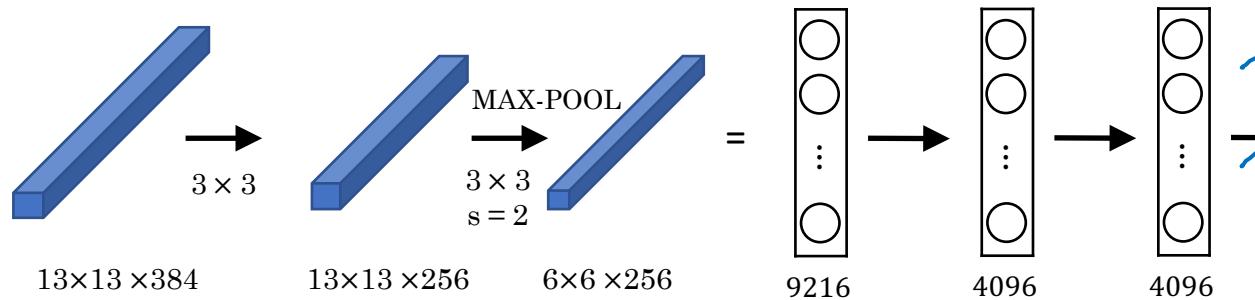
MAX-POOL  
 $3 \times 3$   
 $s = 2$



$3 \times 3$   
same

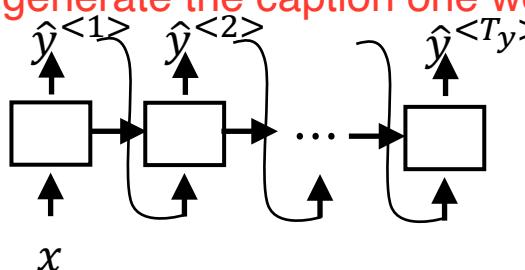


$y^{<1>} y^{<2>} y^{<3>} y^{<4>} y^{<5>} y^{<6>} \{$



Take this 4,096 vector and feed it to an RNN whose job is to generate the caption one word at a time.

You can input an image into a convolutional network, a pre-trained AlexNet, and have the network learn and encode a learner to the features of the input image—a 4,096 dimensional vector that represents the image in this example.



However, there are some differences between how to run/generate a sequence and how to synthesize novel text using a language model. Andrew Ng

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

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

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



deeplearning.ai

# Sequence to sequence models

---

There are some similarity between the sequence to sequence machine translation model and the language models.  
But there are some significant differences as well.

## Picking the most likely sentence

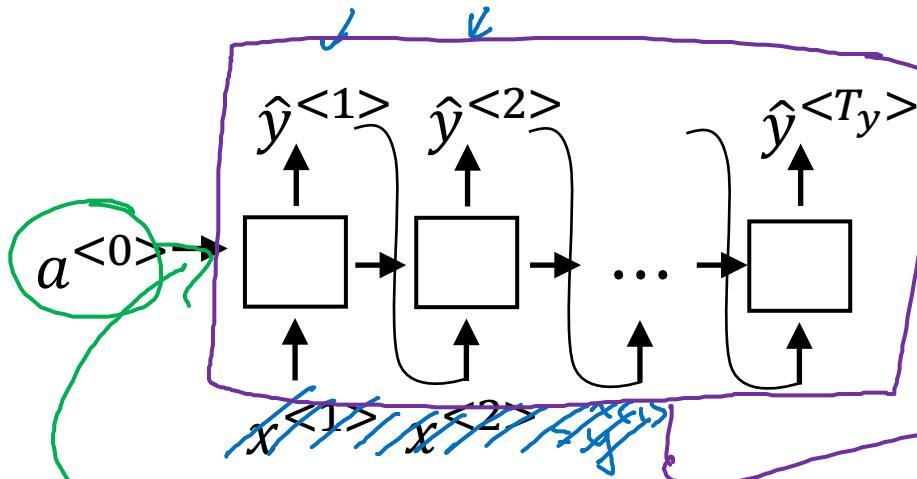
# Machine translation as building a conditional language model

Language model:

Instead of always starting along with the vector of all zeros, the machine translation model has an encoded network that figures out some representation for the input sentence.

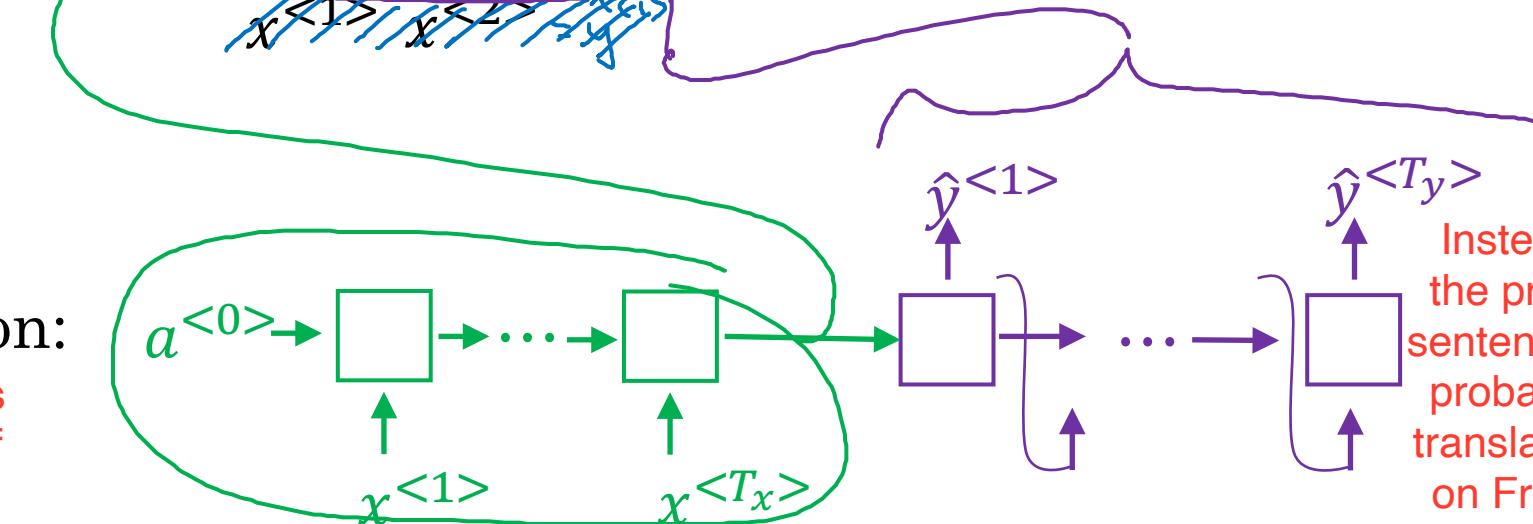
Machine translation:

Machine translation model takes the input sentence and starts off the decoded network with the representation of the input sentence rather than the representation of all zeros..



This language model allows you to estimate the probability of a sentence

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



“Conditional language model”

Andrew Ng

# Finding the most likely translation

Jane visite l'Afrique en septembre.

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

→ Jane is visiting Africa in September.

Given the input—French sentence, the model might tell you what the prob. of difference in corresponding English translation.

→ Jane is going to be visiting Africa in September.

→ In September, Jane will visit Africa.

You don't want to sample output at random. If you sampled words from the distribution ( $p$  of  $y$  given  $x$ ) at random, then you may get a different translation, e.g., the 2nd ~ 4th sentences

→ Her African friend welcomed Jane in September.

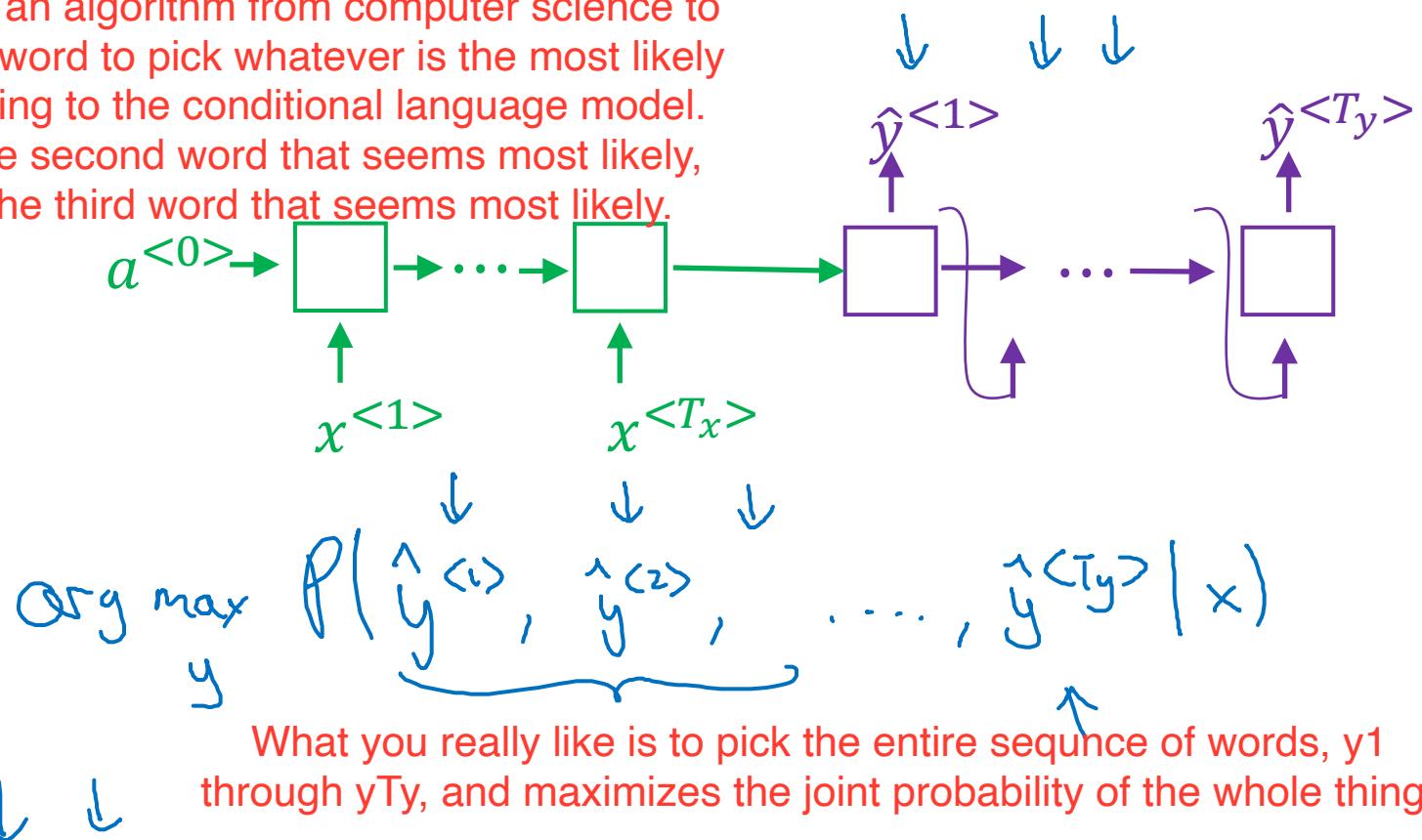
So, instead of trying to sample at random from the distribution, you would like to find the english sentence,  $y$ , that maximize the conditional probability.

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

The most common algorithm for this is “beam search”

# Why not a greedy search?

Greedy search is an algorithm from computer science to generate the first word to pick whatever is the most likely first word according to the conditional language model. And then pick the second word that seems most likely, and then pick the third word that seems most likely.



What you really like is to pick the entire sequence of words,  $y_1$  through  $y_{Ty}$ , and maximizes the joint probability of the whole thing.

→ Jane is visiting Africa in September.

→ Jane is going to be visiting Africa in September.

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

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

10,000 words in a dictionary and 10 word-long translation, then there are a huge space of possible sentences, which is impossible to rate them all. So, people commonly use an approximate search algorithm.

$$10,000$$

$$10$$

$$\frac{10,000}{10}$$

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

However, it is not always optimal to just pick one word at a time. For example, given “Jane is”, “going” can be more common than “visiting”.



deeplearning.ai

# Sequence to sequence models

---

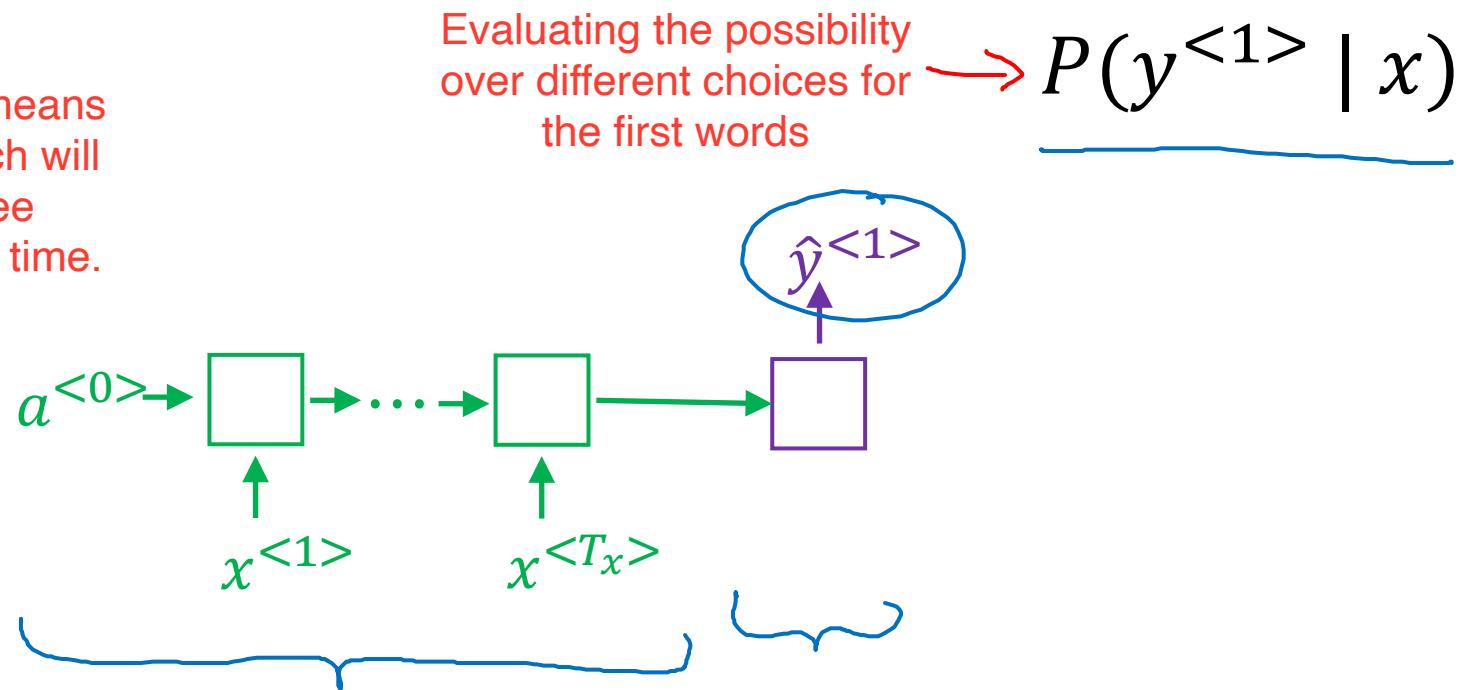
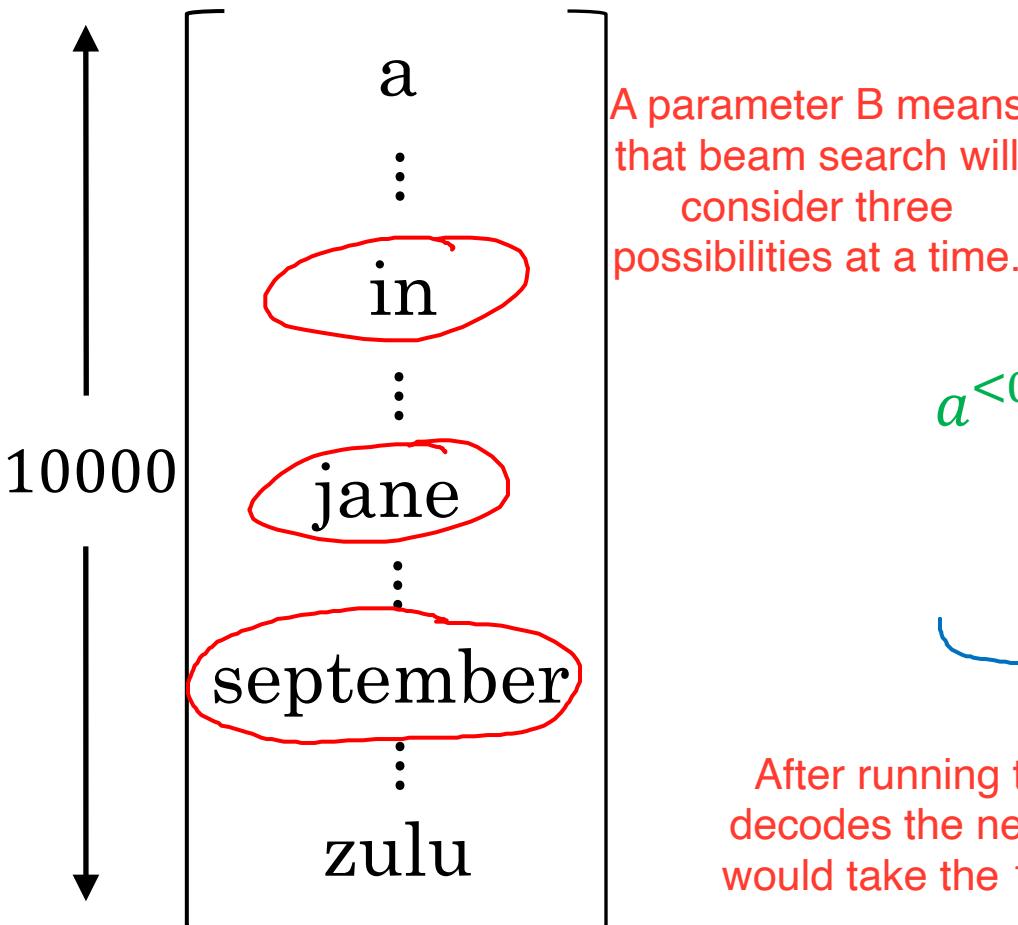
You want to output the best and the most likely English  
translation rather than a random English translation

## Beam search

# Beam search algorithm

B = 3 (beam width)

Step 1: to pick the first words of the English translation by considering multiple alternatives.



After running the input French sentence through the encoder network, the first step decodes the network (e.g., a softmax output — overall 10,000 possibilities). Then you would take the 10,000 possible outputs and keep in memory which were the top three.

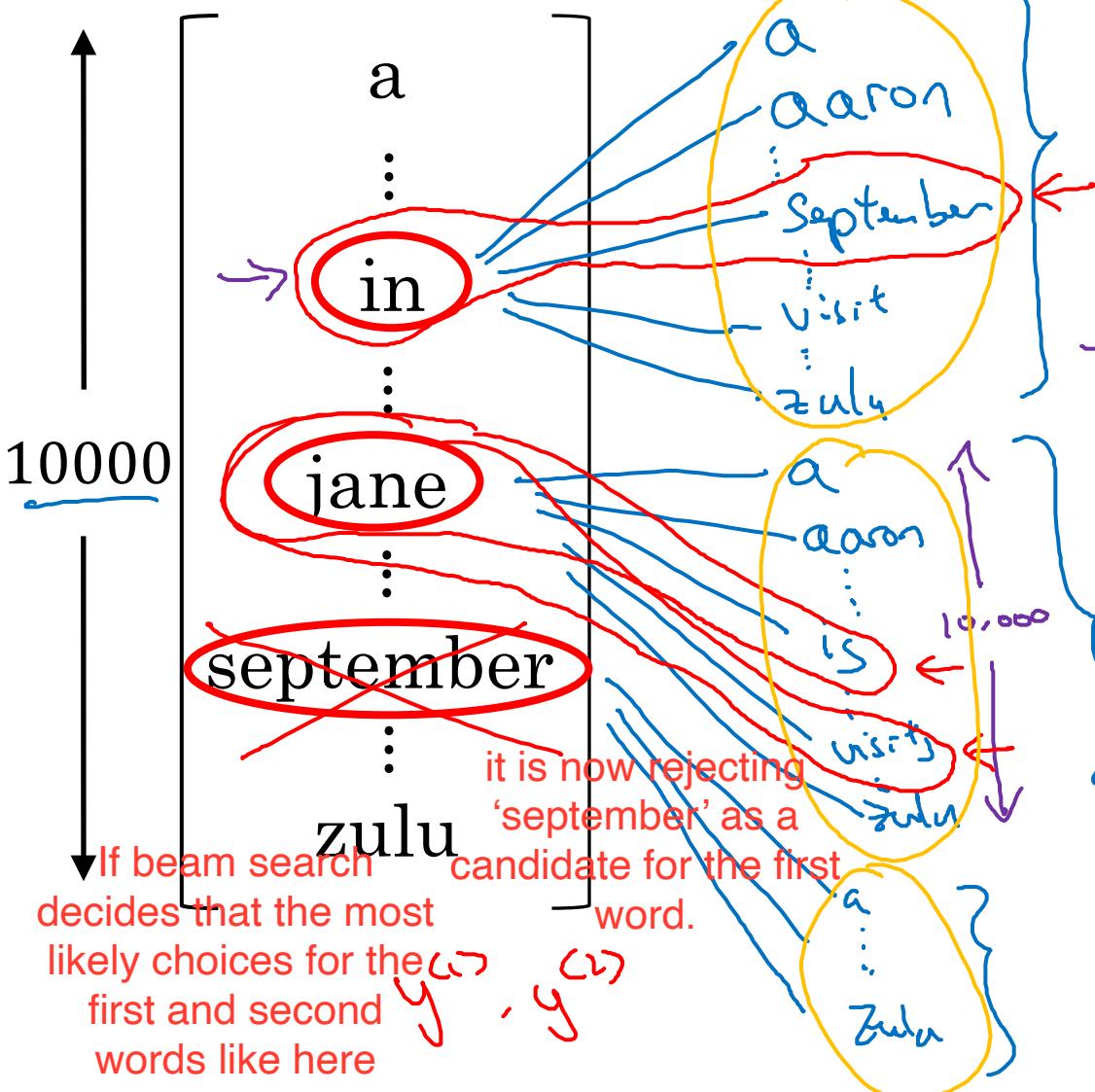
When  $B=3$ , you instantiate 3 copies of the network.

# Beam search algorithm

## Step 1

## Step 2

For each of these three choices, you consider what should be the second word



$$(B = 3)$$

To evaluate the probability of the second word,

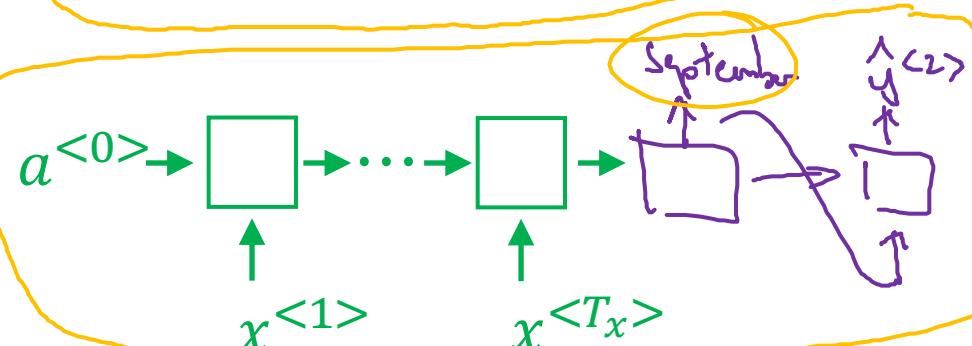
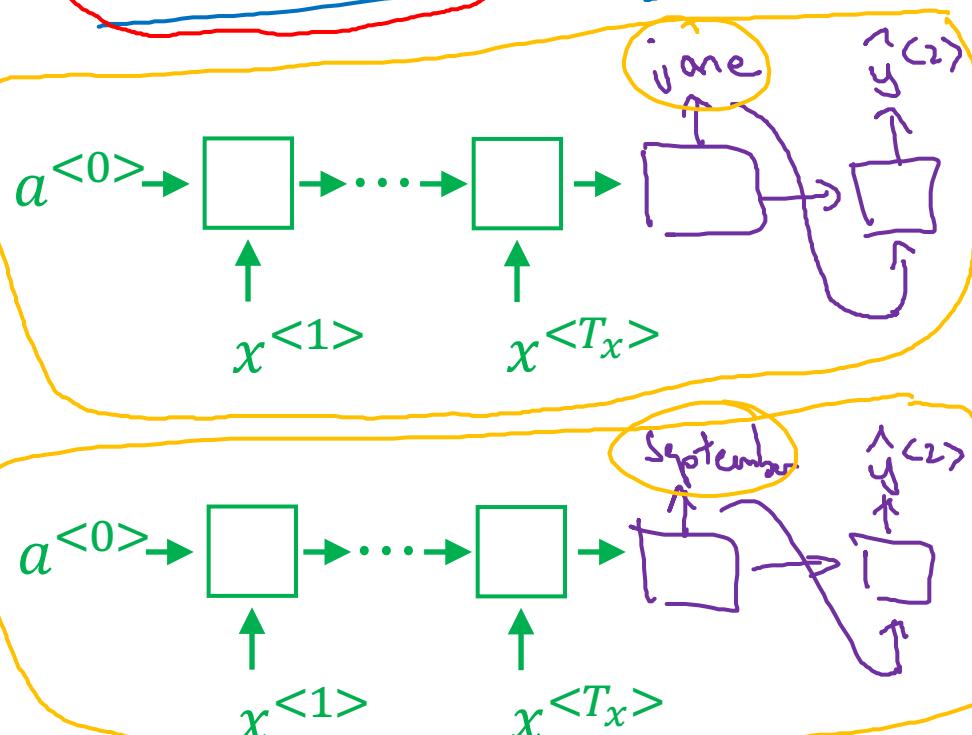
$P(y^{(2)} | x, "in")$

you use the new network fragments, remembering the decoder first output

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

$\text{Poly}^{(2)} | x; \text{one}^{\sim}$

Of 30,000 possibilities, 3 words would be memorized.



# Beam search ( $B = 3$ )

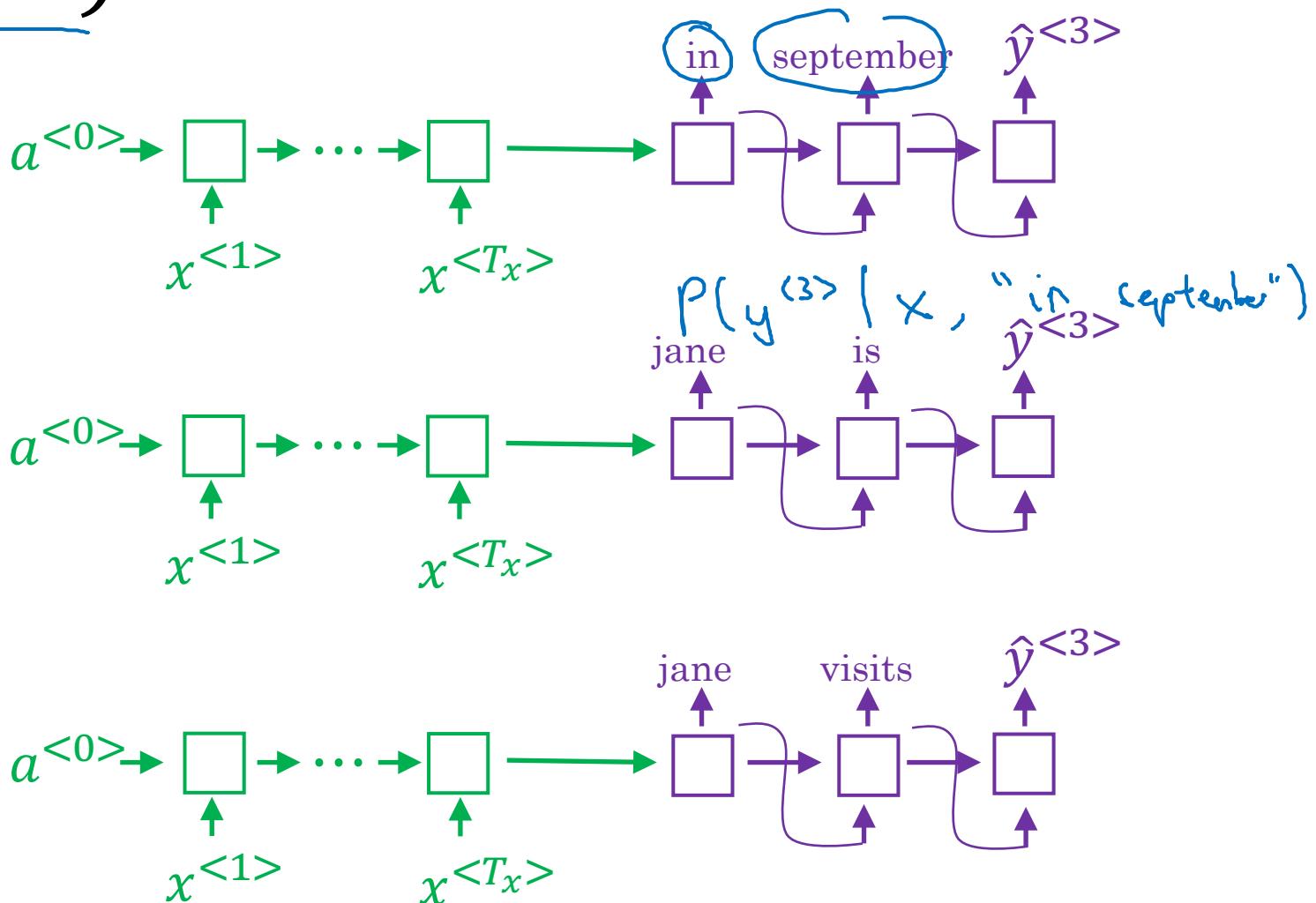
in september

jane is

jane visits

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

$B=1 \rightsquigarrow$  greedy search



jane visits africa in september. <EOS>



deeplearning.ai

# Sequence to sequence models

---

## Refinements to beam search

# Length normalization

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

product of probabilities

These probabilities are all numbers less than one, and multiplying such a tiny number can result in numerical underflow (i.e., too small for the floating point representation in a computer to store accurately).

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

If you have a very long sentence, the prob. of the sentence will be low because of multiplying many terms.

This obj. fn. may unnaturally tend to prefer very short translation

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

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

Normalize this by the number of words in the translation in order to reduce the penalty for putting longer translations.

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

Because logarithmic function is a strictly monotonically increasing function.  
 $\log P(y|x) \leftarrow$   
 their maximization will give the same results  
 $P(y|x) \leftarrow$

By taking logs, you end up with a more numerically stable algorithm that is less prone to numerical rounding errors or numerical underflow.

If alpha is one, it is full normalization.  
 If alpha is zero, it is no normalization.

$$\alpha = 0.7$$

$$\begin{cases} \alpha = 1 \\ \alpha = 0 \end{cases}$$

# Beam search discussion

Pros and cons of setting beam width to be very large vs. very small.

How to choose the beam width?

Beam width B?

$1 \rightarrow 3 \rightarrow 10, 100,$

$1000 \rightarrow 3000$

large B: better result, slower

small B: worse result, faster

If B is very large, you consider a lot of possibilities, and you tend to get a better results because you are consuming a lot of different options.

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

y



deeplearning.ai

# Sequence to sequence models

Beam search: Approximate search algorithm,  
also called a heuristic search algorithm

---

Error analysis can help you focus your time on  
doing the most useful work for your project.

## Error analysis on beam search

What if beam search makes a mistake? Learn how error analysis interacts with beam search and how you can figure out whether it is the beam search algorithm that is causing problems.

# Example

Jane visite l'Afrique en septembre.

→ RNN

→ Beam Search

B↑

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

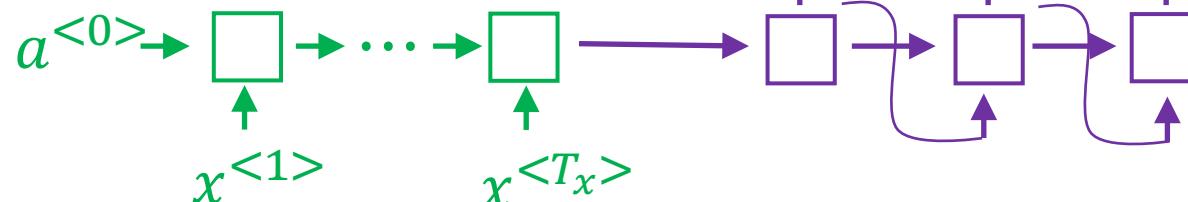
In your machine translation development set, the human provided the translation and call this  $y$  star.

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

RNN computes

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

You run beam search on your learned RNN model, which ends up with the translation  $\hat{y}$



Your model has two components: 1) RNN model: a neural network model, the sequence to sequence model, which means an encoder and a decoder; 2) Beam search algorithm.

It is always tempting to collect more training data that never hurts. In similar way, it is always tempting to increase the beam width that never hurts. But just as getting more training data set might not get you to the level of performance you want. Likewise, increasing the beam width might not get you where you want to go.

Andrew Ng

# Error analysis on beam search

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

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

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

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

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

P(y star given x) as output by the RNN model is greater than P(y hat given x)

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

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

Since beam search algorithm chose  $\hat{y}$ , the way you got  $\hat{y}$  was you had an RNN that computes  $P(y|x)$

Conclusion: Beam search is at fault. Beam search fails to actually give you the value of  $y$  that maximizes  $P(y|x)$ . It chose  $\hat{y}$  that has a smaller value.  
Then, you may need to increase the beam width.

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

$y^*$  is a better translation than  $\hat{y}$ . But RNN predicted  $P(y^*|x) < P(\hat{y}|x)$ .  
Even though  $y^*$  is better translation in this example, the RNN says  $P(y^*)$  is less than  $P(\hat{y})$ , saying that  $y^*$  is less likely the output than  $\hat{y}$ .

Conclusion: RNN model is at fault.

Then, you could do a deeper layer of analysis to try to figure out if you want to add regularization, or get more training data, or try a different network architecture, or etc.

# Error analysis process

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

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



deeplearning.ai

BLEU stands for Bilingual Evaluation Understudy

# Sequence to sequence models

---

Bleu score  
(optional)

# Evaluating machine translation

French: Le chat est sur le tapis.

Reference 1: The cat is on the mat. 

Human 1 & 2 translation

Reference 2: There is a cat on the mat. 

MT output: the the the the the the.

Precision: 7 over 7

One way to measure how good the MT output is to look at each the words in the output and see if it appears in the reference: it is called a precision of the MT output.

Modified precision: 2 over 7

A modified precision measure gives each word credit only up to the maximum number of times it appears in the reference sentence. In reference 1, 'the' appears two, In reference 2, 'the' appears one. So, 'the' gets credit up to twice.

However, in the BLEU score, you don't want to just look at isolated words but at pairs of words as well.

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

Bleu  
bilingual evaluation understudy

BLEU score: given Machine generated translation, it allows you to automatically compute a score that measures how good is that machine translation. The intuition is so long as the Machine generated translation is pretty close to any of references by humans, then it will get high BLEU score

BLEU score could be a substitute for having humans evaluate every output of a machine translation system.

Andrew Ng

# Bleu score on bigrams

Bigrams: 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. ←

Possible bigrams

the cat

Count  
2 ←

Count clip  
1 ←

Maximum from  
references

cat the

1 ←

0

cat on

1 ←

1 ←

on the

1 ←

1 ←

the mat

1 ←

1 ←

Modified precision on bigrams.

$$\frac{4}{6}$$

# Bleu score on unigrams

If the MT output is exactly the same as either Reference 1 or 2, then all of these values P1 and P2 and so on, they will all be equal to 1.

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.

$$P_1, P_2, \dots = 1.0$$

To get a modified precision of 1, you have to have the output exactly be equal to one of the references.

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

↑  
Unigram

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

↑  
n-gram

# Bleu details

$p_n$  = Bleu score on n-grams only

$P_1, P_2, P_3, P_4$

Combined Bleu score:

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

$BP$  = brevity penalty

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

BLEU score is a useful single real number evaluation metric to use whenever you want your algorithm to generate a piece of text and you want to see whether it has similar meaning as a reference piece of text generated by humans.

But, this is not used for speech recognition, because there is usually one ground truth in speech recognition.

For image captioning, and multiple captions for a picture, it could be good.





deeplearning.ai

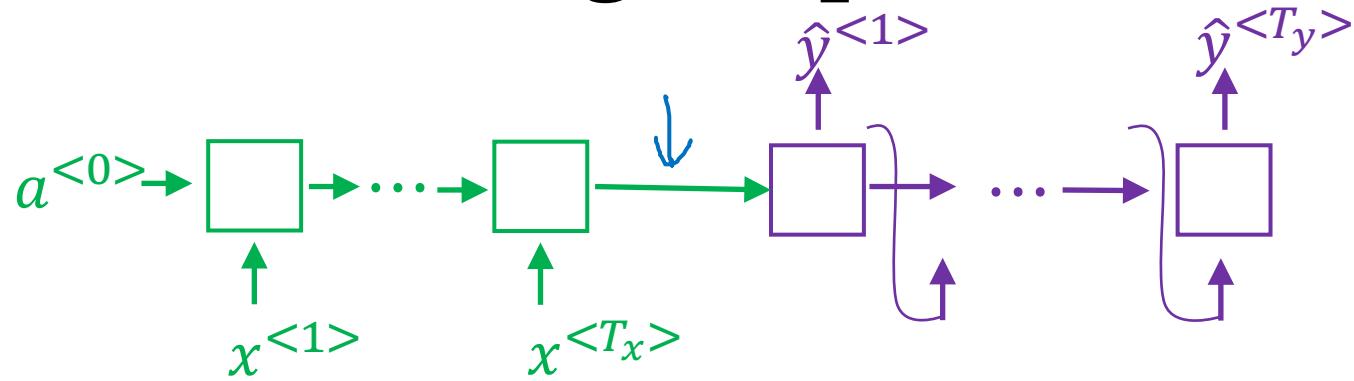
# Sequence to sequence models

---

Attention model: Modification to an Encoder-Decoder  
architecture for machine translation

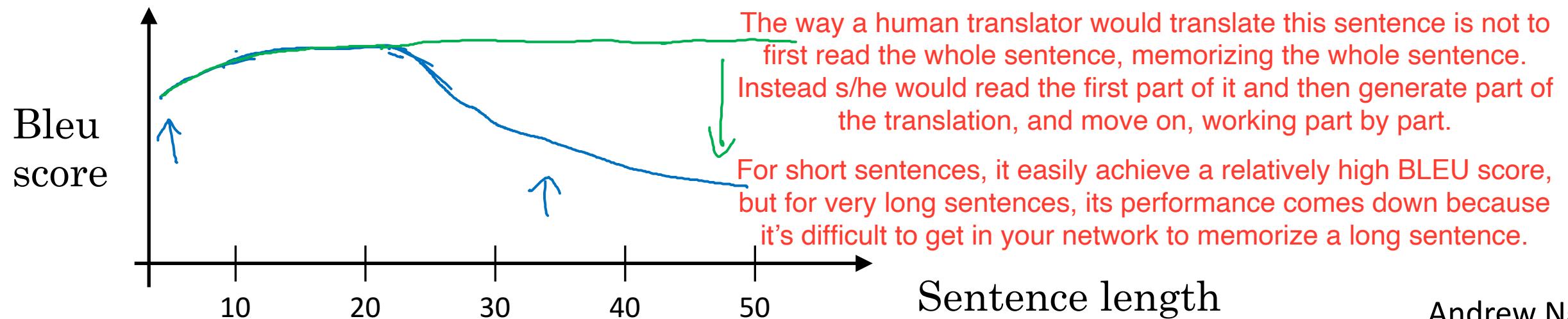
## Attention model intuition

# 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.



Like humans, Attention model looks at part of the sentence at a time.

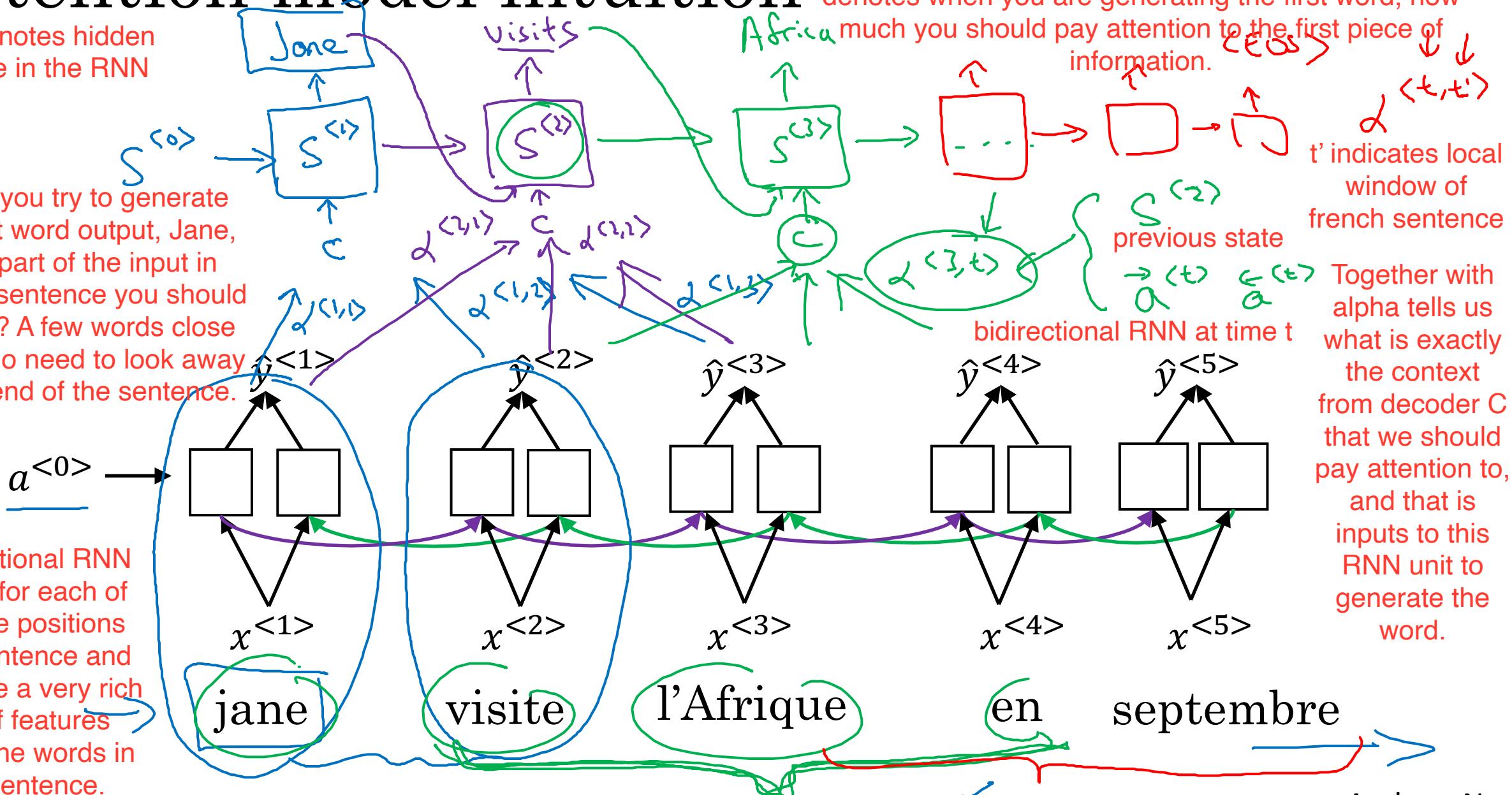
# Attention model intuition

$S$  denotes hidden state in the RNN

Attention model would compute a set of attention weights,  $(\alpha_{one, 1}), (\alpha_{one, 2}), (\alpha_{one, 3})$ , which denotes when you are generating the first word, how much you should pay attention to the first piece of information.

When you try to generate the first word output, Jane, what part of the input in French sentence you should look at? A few words close by but no need to look away at the end of the sentence.

Bidirectional RNN works for each of the five positions into sentence and compute a very rich set of features about the words in the sentence.





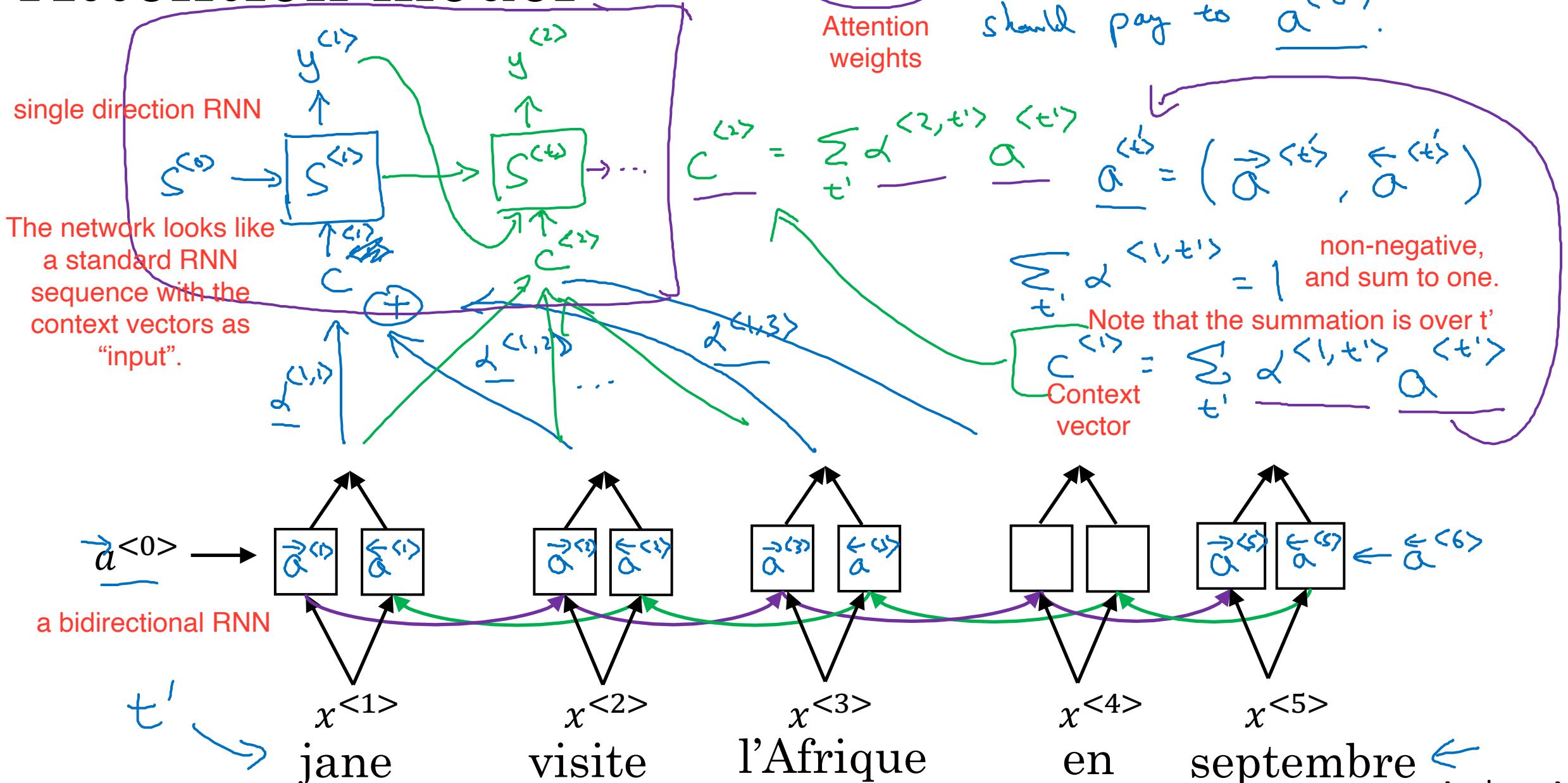
deeplearning.ai

# Sequence to sequence models

---

## Attention model

# Attention model

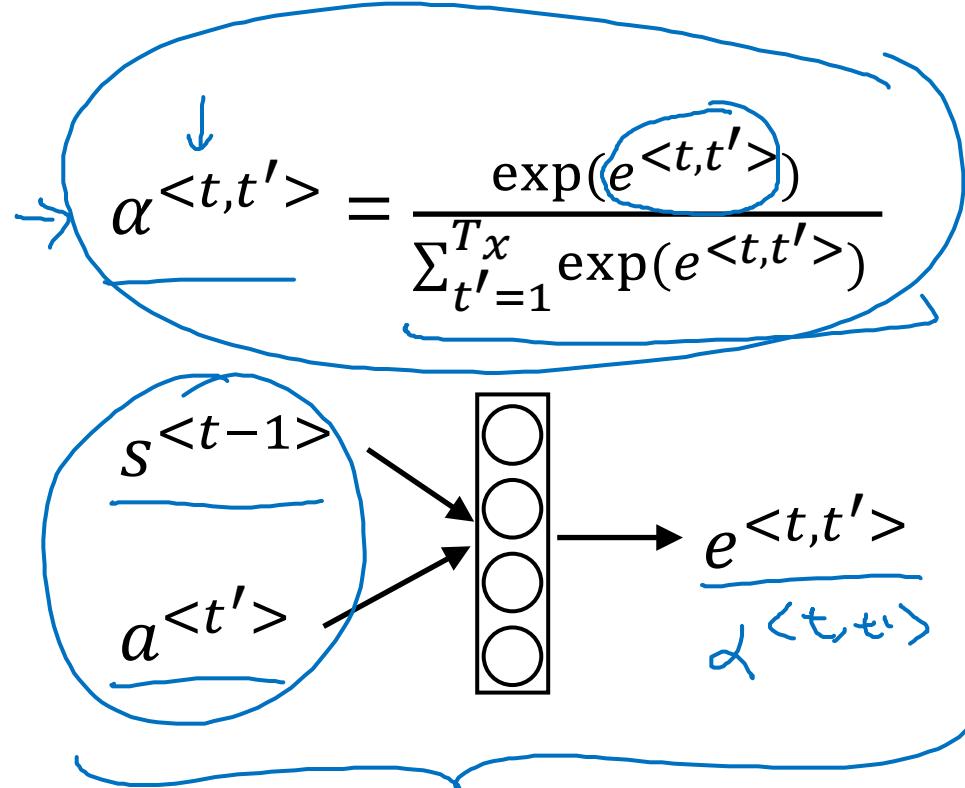


# Computing attention $\alpha^{}$

T<sub>x</sub>

T<sub>y</sub>

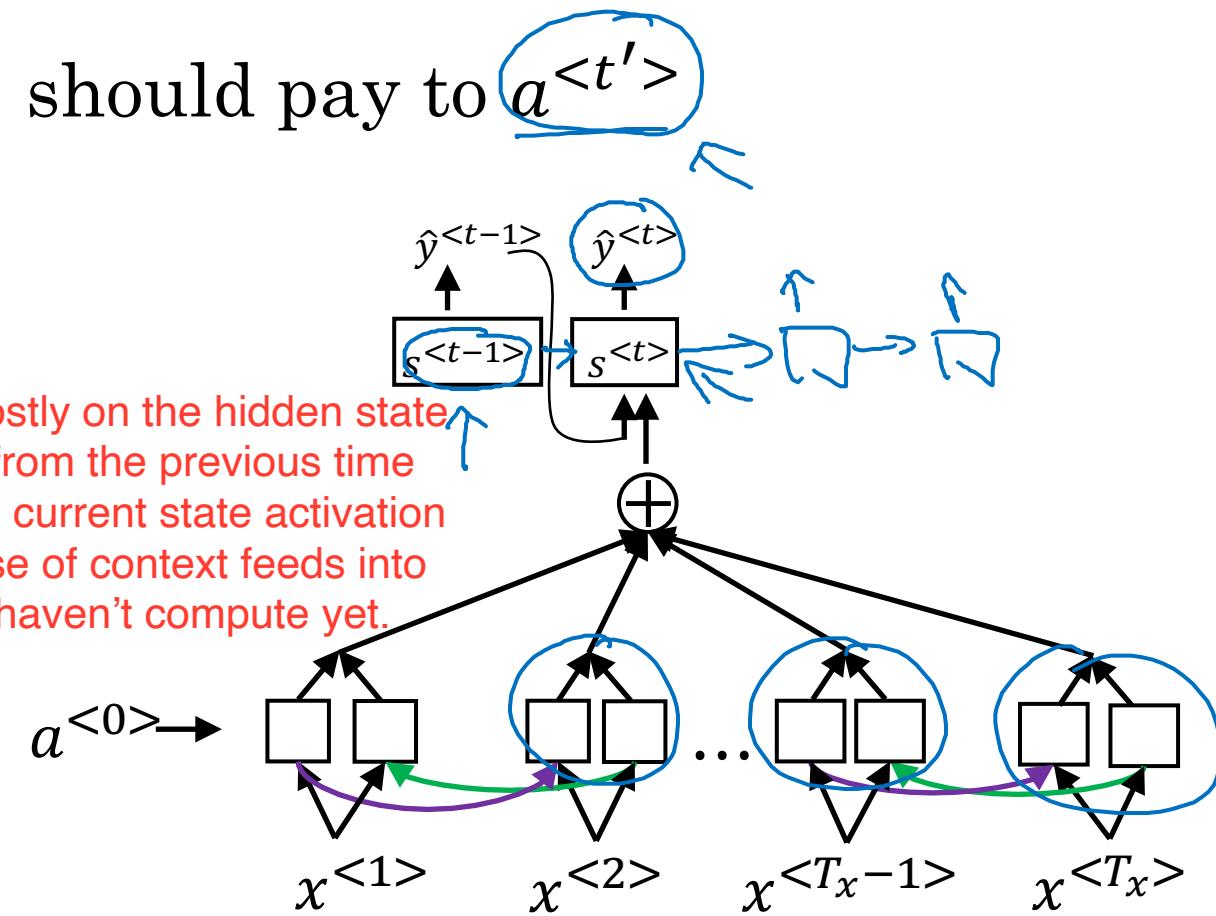
$\alpha^{}$  = amount of attention  $y^{}$  should pay to  $a^{}$



small neural network

Downside of the algorithm: it takes quadratic time/cost to run

Depends mostly on the hidden state activation from the previous time step, not the current state activation yet (because of context feeds into this) you haven't compute yet.



If you have  $T_x$  words as input and  $T_y$  words as output, the total number of the attention parameters will be  $T_x$  times  $T_y$

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

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

Andrew Ng

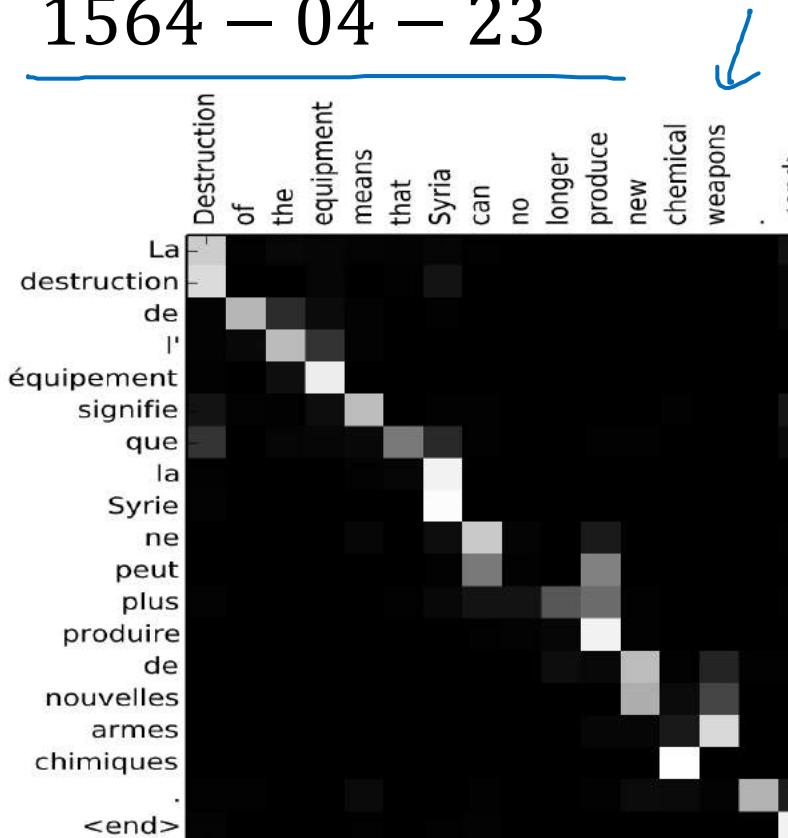
# Attention examples

normalized format

July 20th 1969 → 1969 – 07 – 20

23 April, 1564 → 1564 – 04 – 23

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



different colors are plotted for the magnitude of the different attention weights.



deeplearning.ai

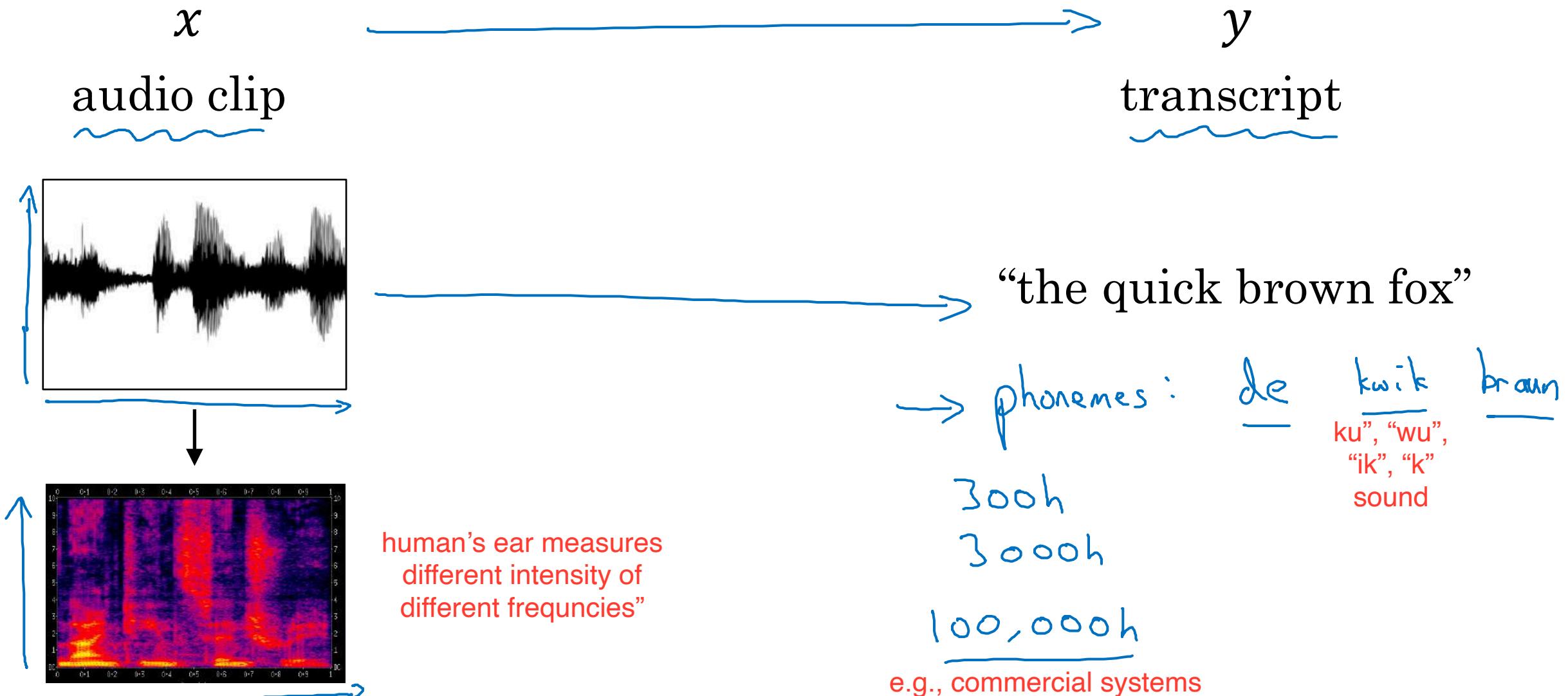
Audio data

---

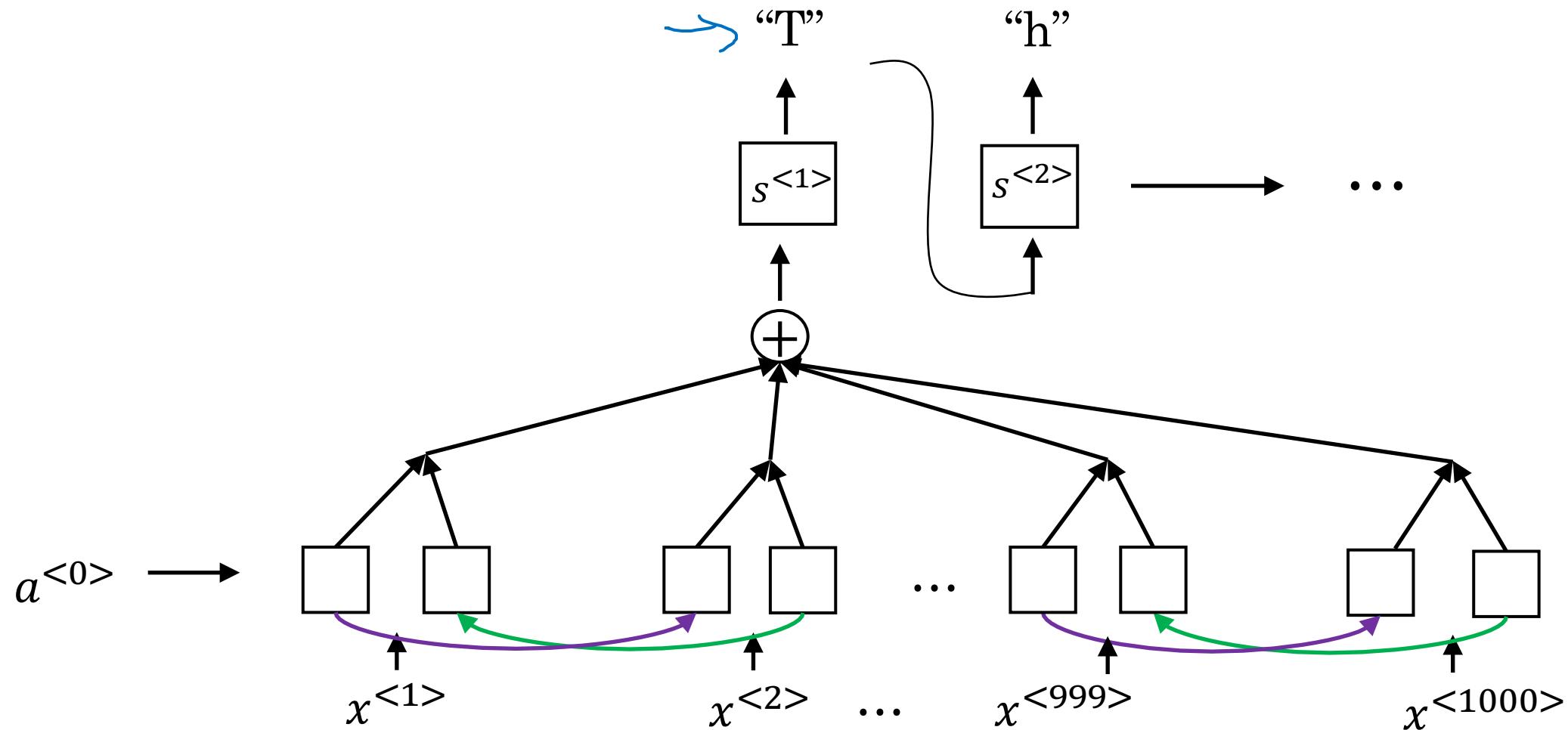
Speech recognition

Application of sequence to sequence model

# Speech recognition problem



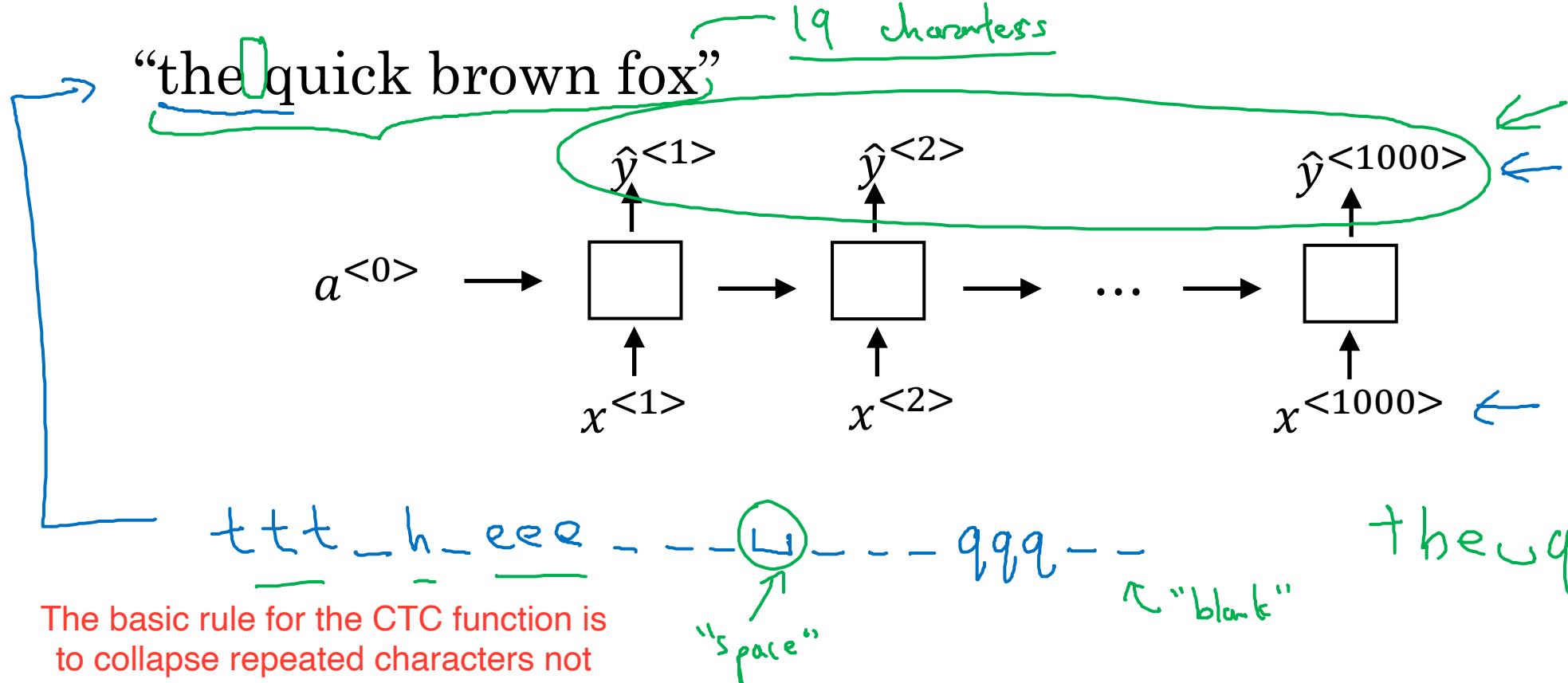
# Attention model for speech recognition



# CTC cost for speech recognition

(Connectionist temporal classification)

In practice, a bidirectional RNN would be used.



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



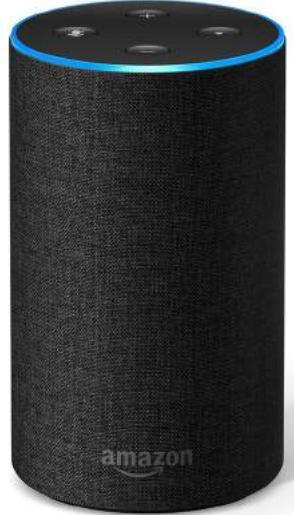
deeplearning.ai

Audio data

---

Trigger word  
detection

# What is trigger word detection?



Amazon Echo  
(Alexa)



Baidu DuerOS  
(xiaodunihao)

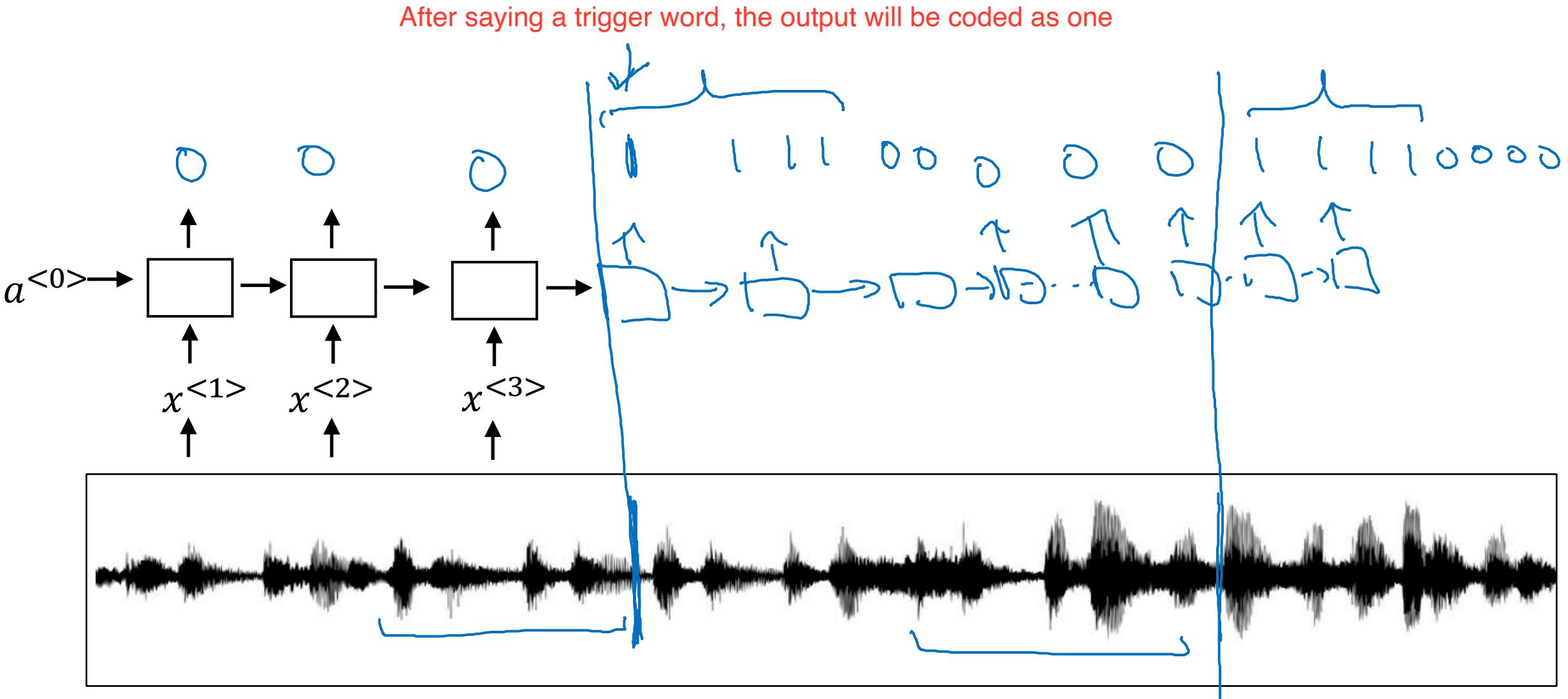


Apple Siri  
(Hey Siri)



Google Home  
(Okay Google)

# Trigger word detection algorithm





deeplearning.ai

# Conclusion

---

## Summary and thank you

# Specialization outline

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

# Deep learning is a super power

Please buy this  
from shutterstock  
and replace in  
final video.



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

Thank you.

- Andrew Ng