Applications: machine translation, image captioning (sequence generation from image).

# Sequence to sequence models

---

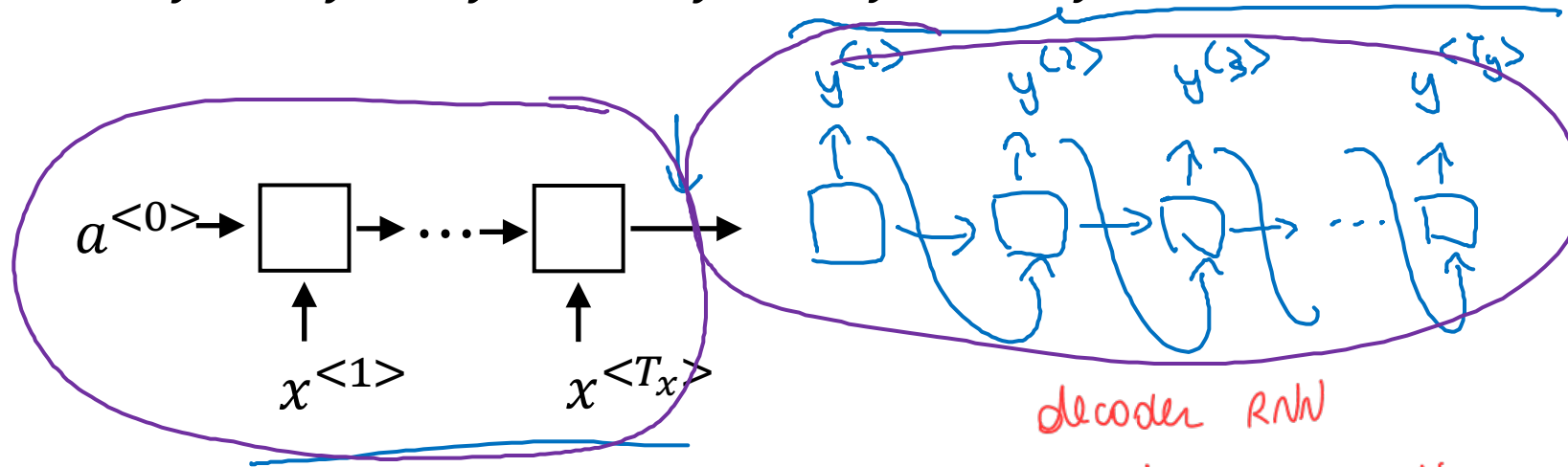# Basic models

deeplearning.ai

# Sequence to sequence model

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

Jane visite l'Afrique en septembre

⟶   Jane is visiting Africa in September.

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



$a^{<0>}$

$x^{<1>}$     $x^{<T_x>}$
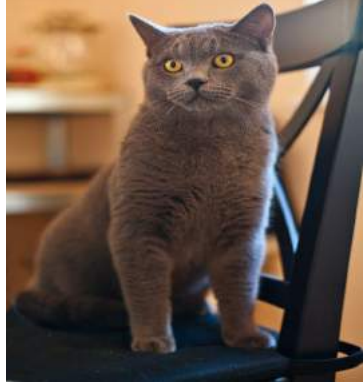
$y^{<1>}$   $y^{<2>}$   $y^{<3>}$   $y^{<T_y>}$

decoder RNN
produces decoding into different language
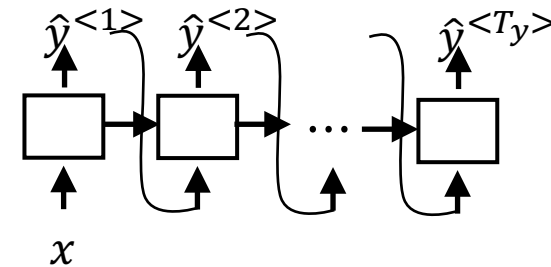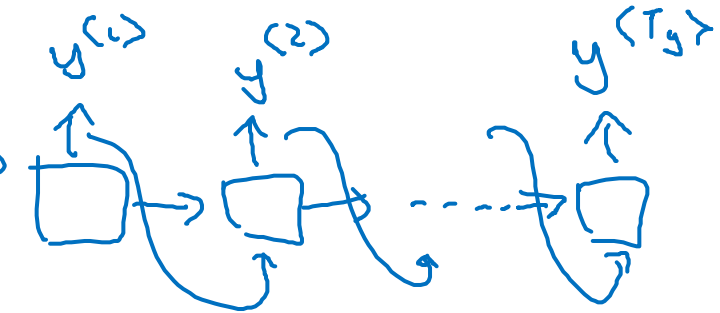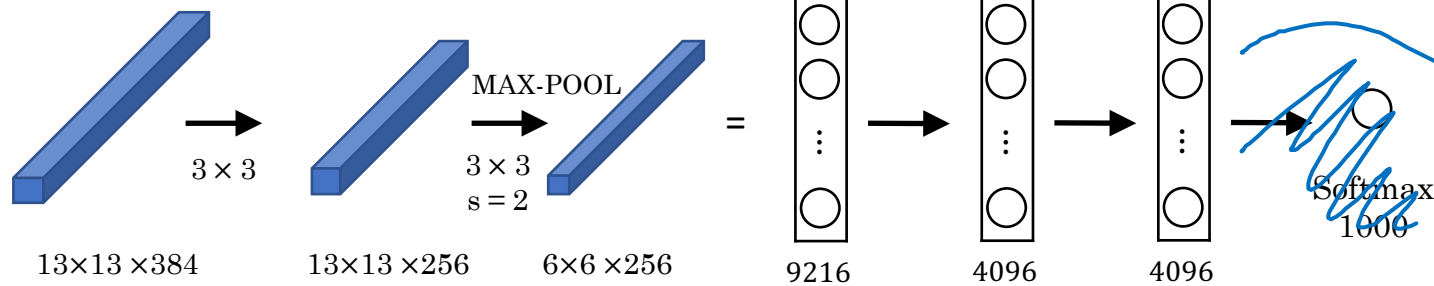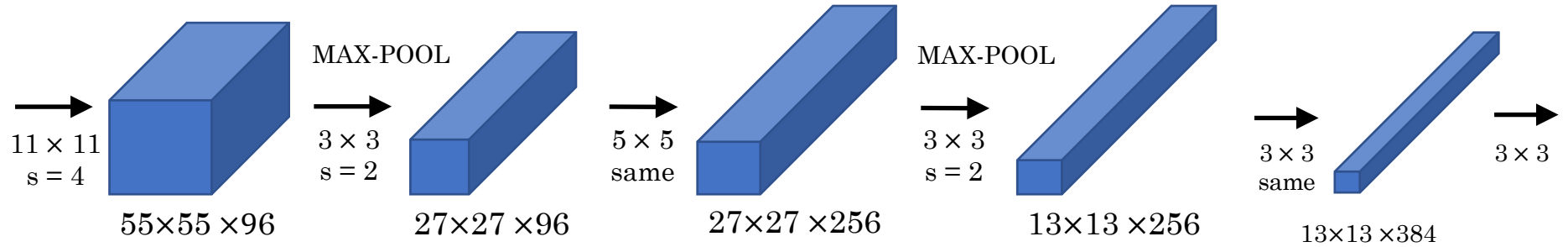
encoder RNN
encodes original message

[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]

Andrew Ng

# Image captioning

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

A   cat   sitting   on   a   chair



$11 \times 11$
$s = 4$

$55 \times 55 \times 96$

MAX-POOL

$3 \times 3$
$s = 2$

$27 \times 27 \times 96$

$5 \times 5$
same

$27 \times 27 \times 256$

MAX-POOL

$3 \times 3$
$s = 2$

$13 \times 13 \times 256$

$3 \times 3$
same

$13 \times 13 \times 384$

$3 \times 3$

$13 \times 13 \times 384$

$3 \times 3$

$13 \times 13 \times 256$

MAX-POOL

$3 \times 3$
$s = 2$

$6 \times 6 \times 256$

$= \quad$ 9216 $\quad$ 4096 $\quad$ 4096

Softmax
1000

$y^{(1)} \quad y^{(2)} \quad y^{<T_y>}$

$\hat{y}^{<1>} \quad \hat{y}^{<2>} \quad \hat{y}^{<T_y>}$

$x$

[Mao et. al., 2014. Deep captioning with multimodal recurrent neural networks]
[Vinyals et. al., 2014. Show and tell: Neural image caption generator]
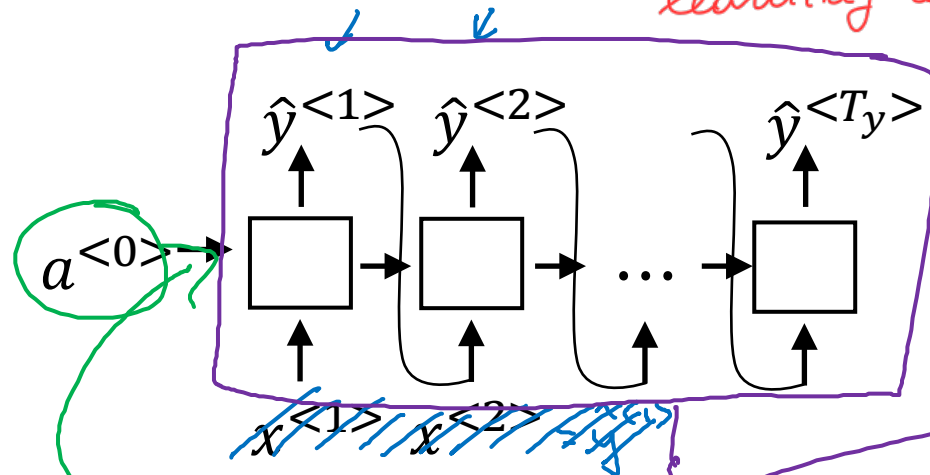[Karpathy and Li, 2015. Deep visual-semantic alignments for generating image descriptions]

Andrew Ng

# Machine translation as building a conditional language model



Language model:

Machine translation:
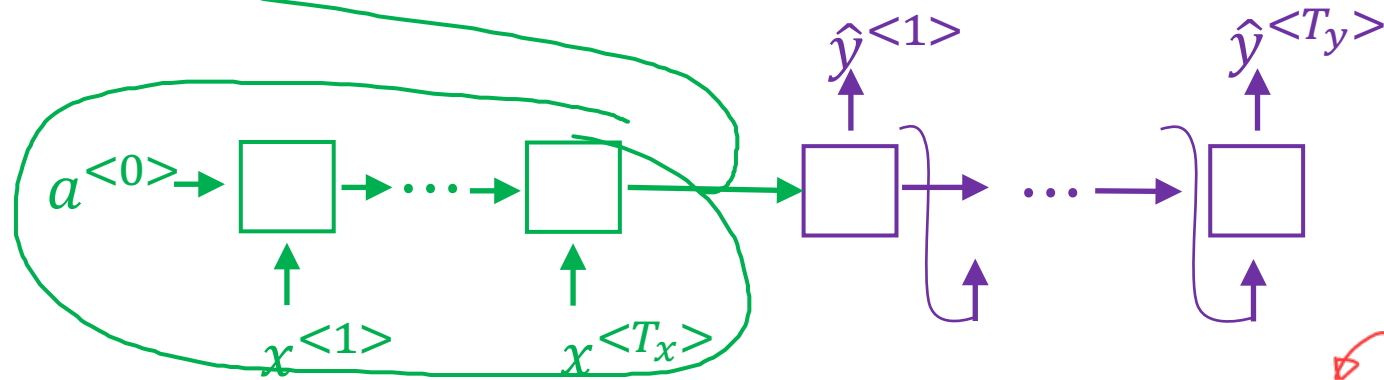
Ng sees machine translation as the task of learning a CONDITIONAL LANGUAGE MODEL. While in language modeling you learn

$$P(y^{(1)}, \ldots, y^{(T_y)})$$

Here you learn (approximate) the probability of a series of words in english conditioned by a series of words in french.

"Conditional language model"

$$P(y^{(1)}, \ldots, y^{(T_y)} \mid x^{(1)}, \ldots, x^{(T_x)})$$

Andrew Ng

# Finding the most likely translation

Jane visite l'Afrique en septembre.

English sentence | French sentence

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

→ 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.

While decoding you basically sample words from the learned probability distribution, one word at a time. But this could represent a suboptimal translation and the words sequence $y^{<1>}, y^{<2>}, \ldots, y^{<T_y>}$ could change every time.

=> Use a search algorithm to find the sequence $y^{<1>} \ldots y^{<T_y>}$ that maximizes the Prob.

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

Usually BEAM search is used instead of greedy approaches.

Andrew Ng

# Why not a greedy search?

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



$a^{<0>}$ → □ → .... → □ → □ → .... → □

$x^{<1>}$  $x^{<T_x>}$

$\hat{y}^{<1>}$  $\hat{y}^{<T_y>}$

$\text{arg max}_y \ P\left(\hat{y}^{<1>}, \hat{y}^{<2>}, \ldots, \hat{y}^{<T_y>} | x\right)$

$10,000$

$10$

$\dfrac{10,000^{10}}{}$

$P(y|x)$

→ Jane is visiting Africa in September.

→ Jane is going to be visiting Africa in September.

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

Andrew Ng

Beam search is an extension of greedy search where the B greediest options are selected at each step. • At step 1 the B=3 most likely words are selected (st. highest $P(y^{<1>}|x) = \hat{y}^{<1>1}$
2nd highest $P(y^{<1>}|x) = \hat{y}^{<1>2}$
3rd highest $P(y^{<1>}|x) = \hat{y}^{<1>3}$

• At step 2, for each of those 3 words, the set of possible following words is computed and the 3 most likely sequences $y^{<1>}, y^{<2>}$ are determined

• At step 3, repeat step 2 using the three sequences of words. Each of the B "branches" is closed when the sequence generates EOL.

This algorithm reduces the greediness of pure greedy search.

# Sequence to sequence models

# Beam search

# Beam search algorithm

$B = 3$

## Step 1

$$\begin{bmatrix} a \\ \vdots \\ in \\ \vdots \\ jane \\ \vdots \\ september \\ \vdots \\ zulu \end{bmatrix}$$

10000

$\rightarrow P(y^{<1>} \mid x)$



$\hat{y}^{<1>}$  softmax

$a^{<0>} \rightarrow \cdots \rightarrow$

$x^{<1>} \qquad x^{<T_x>}$

Andrew Ng

# Beam search algorithm

$(B = 3)$

Step 1          Step 2

$$\begin{bmatrix} a \\ \vdots \\ in \\ \vdots \\ jane \\ \vdots \\ september \\ \vdots \\ zulu \end{bmatrix}$$

10000

$y^{<1>}, y^{<2>}$

a
aaron
$\vdots$
September
visit
$\vdots$
zulu

a
aaron
$\vdots$
is
$\vdots$
visit
$\vdots$
zulu

a
$\vdots$
zulu

10,000

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow$

$x^{<1>}$        $x^{<T_x>}$

$\hat{y}^{<1>}$    $\hat{y}^{<2>}$

in

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

in

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

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow$

$x^{<1>}$        $x^{<T_x>}$

jane          $\hat{y}^{<2>}$

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

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow$

$x^{<1>}$        $x^{<T_x>}$

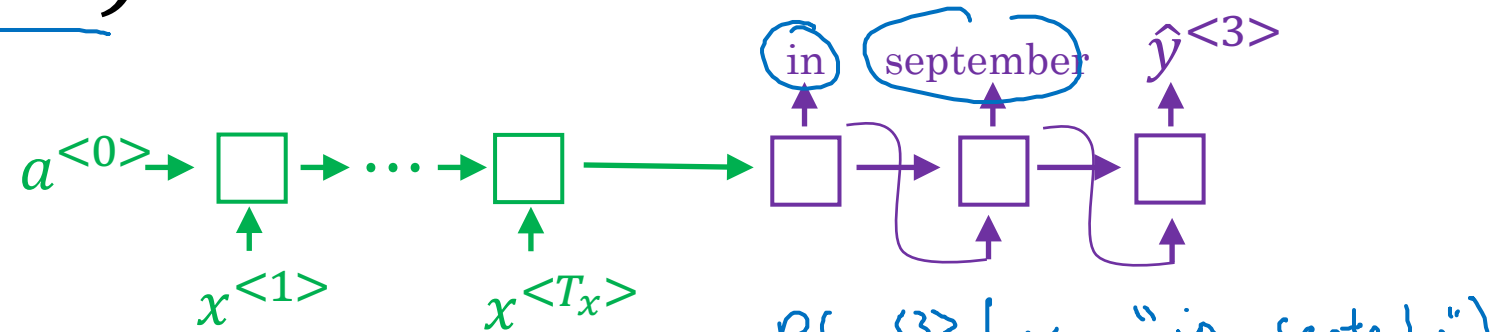September     $\hat{y}^{<2>}$

Andrew Ng

# Beam search ($B = 3$)

$B = 1 \rightsquigarrow$ greedy search

in september
- a
- aaron
- jane
- zulu

jane is
- a
- visits
- zulu

jane visits
- a
- africa
- zulu



in   september   $\hat{y}^{<3>}$

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow \square$

$x^{<1>}$   $x^{<T_x>}$

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

jane   is   $\hat{y}^{<3>}$

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow \square$

$x^{<1>}$   $x^{<T_x>}$

jane   visits   $\hat{y}^{<3>}$

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow \square$

$x^{<1>}$   $x^{<T_x>}$

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

jane visits africa in september. <EOS>

Andrew Ng

The B sequences returned by beam search may have different length because sequences are terminated when they are among the B most likely and terminate with a <EOL>. For example at the end of beam search you may have: "Jane likes apples<EOL>", "Jane likes eating apples<EOL>"

# Sequence to sequence models

"Jane eats the apple <EOL>".

Longer sentences are more penalized by the log loss.

Therefore one can normalize the loss by the length of the sequence.

# Refinements to beam search

deeplearning.ai

# Length normalization

What we do in sequence generation is $P(y^{<1>}, \dots, y^{<T_y>}|x) = P(y^{<1>}|x) P(y^{<2>}|x, y^{<1>}) \dots$
$$P(y^{<T_y>}|x, y^{<1>}, \dots, y^{<T_y-1>})$$

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

Which is essentially equivalent to

$\log$

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

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

(which is more numerically stable (avoid rounding errors due to very small values arising from many multiplications).

$$P(y|x) \leftarrow$$

This value is inversely proportional to the length of the sequence. (= Sentence becomes less likely the more long it is).

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

To remove this dependency we normalize this formula by the length of the sentence.

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

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

regulates the amount of normalization.

BUT HOW DOES THIS INFLUENCE SEQUENCE GENERATION IN PRACTICE?

Andrew Ng

# Beam search discussion

large B: better result, slower

small B: worse result, faster

Beam width B?

$1 \to 3 \to 10$  $100$,  $1000 \to 3000$

production system

research applications

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

Mispredictions in a sequence generation problem can be due to RNN errors or the beam search process. Error analysis can be performed on the mispredicted sentences (of the training or development set depending on whether to investigate sources for high bias or variance) to determine who's responsible for the error. For the mispredicted sample yhat collect the human ground truth y*. Compute the probabilities of yhat and t* by plugging them into the RNN. If $P(y^*|x) > P(\text{y-hat}|x)$, beam search is at fault. If $P(y^*|x) \leq P(\text{y-hat}|x)$, the RNN model is at fault. Beam search can be improved increasing the beam size B.

# Sequence to sequence models

---

# Error analysis on beam search

# Example

Jane visite l'Afrique en septembre.

$\rightarrow$ RNN

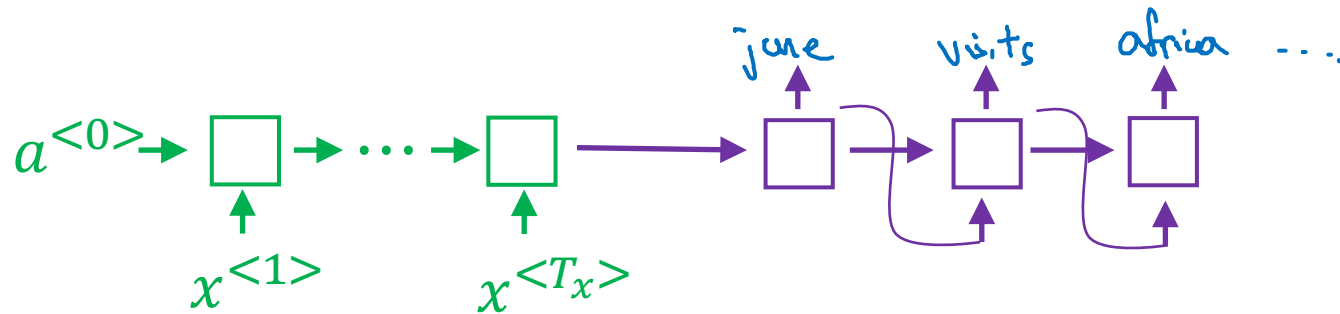$\rightarrow$ Beam Search

BT

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

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

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

jane    visits    africa ...

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow \square$

$x^{<1>}$          $x^{<T_x>}$

Andrew Ng

# Error analysis on beam search

$P(y^* | x)$

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

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

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

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

$\text{arg max}_y P(y | x)$

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

Conclusion: Beam search is at fault.

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

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

Conclusion: RNN model is at fault.

Andrew Ng

# Error analysis process

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

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

Sequence to sequence models

Bleu score (optional)

deeplearning.ai

# 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:                     Modified precision:

Bleu
bilingual evaluation understudy

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

Andrew Ng

# 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 | Count$_{clip}$ |
|---|---|---|
| the cat | 2 ← | 1 ← |
| cat the | 1 ← | 0 |
| cat on | 1 ← | 1 ← |
| on the | 1 ← | 1 ← |
| the mat | 1 ← | 1 ← |

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

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

# Bleu score on unigrams

Example:    Reference 1: The cat is on the mat.

Reference 2: There is a cat on the mat.

$\rightarrow$ MT output: The cat the cat on the mat.   $(\hat{y})$

$$p_1 = \frac{\sum\limits_{unigram \in \hat{y}} count_{clip}(unigram)}{\sum\limits_{unigram \in \hat{y}} count(unigram)}$$

$$p_n = \frac{\sum\limits_{ngram \in \hat{y}} count_{clip}(ngram)}{\sum\limits_{ngram \in \hat{y}} count(ngram)}$$

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

# 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 = brevy penalty

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

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

Andrew Ng

The Attention Model is a modification to the Encoder-Decoder architecture for machine translation that allows it to perform better, especially on long sentences. Instead of memorizing the entire input sentence before translating, the Attention Model focuses on parts of the input while generating the output. Attention weights are used to determine how much focus should be given to each input word while generating a certain output.
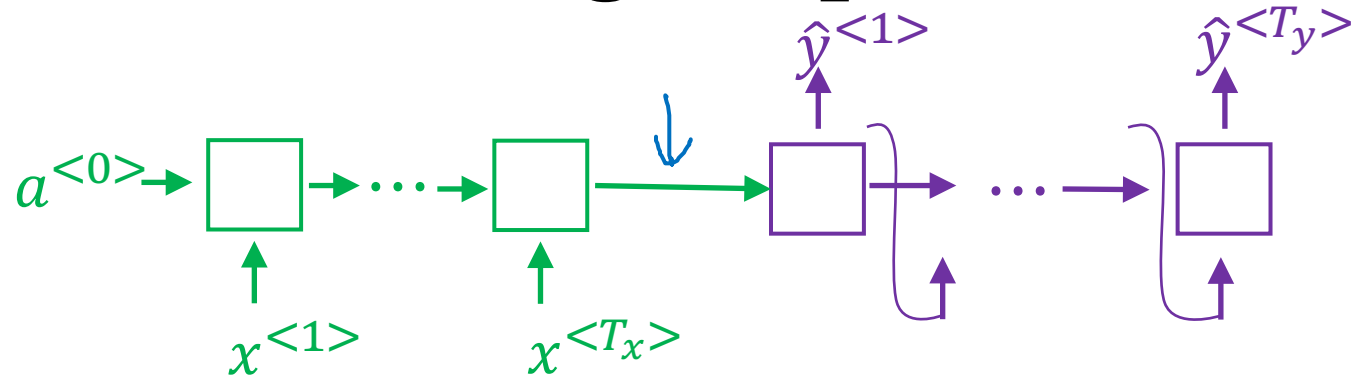
# Sequence to sequence models

# Attention model intuition
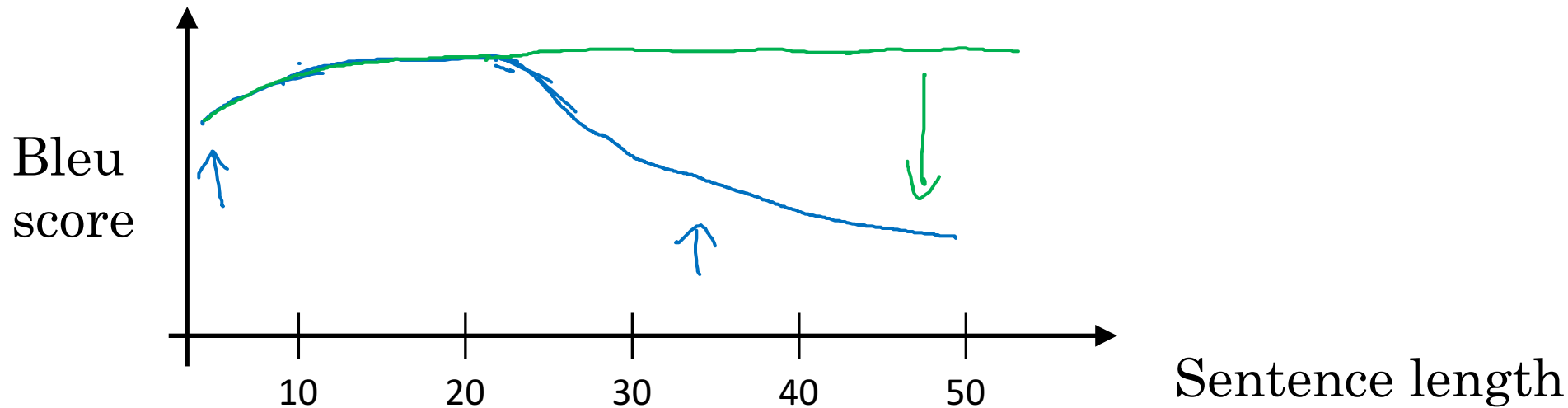
Specifically, the decoder at the step t receives yhat^<t-1> as input, along with a set of attention weights that tell which input words serve as context for the t-th prediction. The weights had been computed by the encoder prior to that. Note that if Tx and Ty are the lengths of the input and ouput sentences there will be Ty sets of attention weights and each set contains Tx weights.

# The problem of long sequences

$$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow \square \rightarrow \cdots \rightarrow \square$$

$$x^{<1>} \qquad x^{<T_x>} \qquad \hat{y}^{<1>} \qquad \hat{y}^{<T_y>}$$

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.

Bleu score

Sentence length
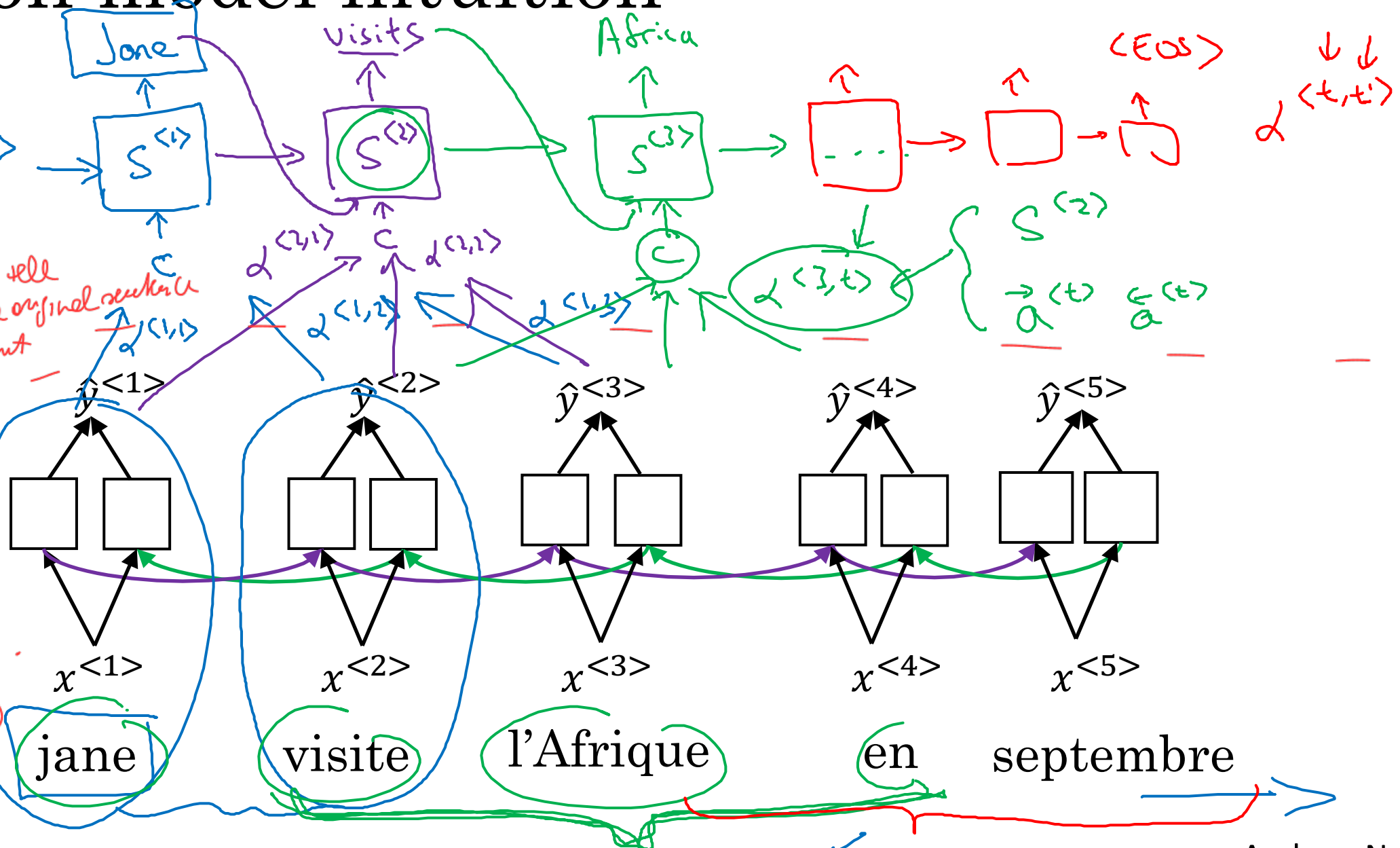
10    20    30    40    50

# Attention model intuition



② USUAL RNN (decoder) takes as input "attention weights" at each step that tell it where to look in the original sequence to produce the current output.

$S^{<0>}$

Jane

visits

Africa

<EOS>

$\alpha^{<t,t'>}$

$S^{<1>}$   $S^{<2>}$   $S^{<3>}$

$\alpha^{<2,1>}$  c  $\alpha^{<2,2>}$

$\alpha^{<3,t>}$

$\begin{cases} S^{<2>} \\ \rightarrow a^{(t)} \leftarrow a^{(t)} \end{cases}$

$\alpha^{<1,1>}$   $\alpha^{<1,2>}$   $\alpha^{<1,3>}$

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

① $a^{<0>}$

Bidirectional RNN used as encoder. (learns context = attention weights) somehow.

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

jane    visite    l'Afrique    en    septembre

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

Andrew Ng

Sequence to
sequence models

Attention model

deeplearning.ai

# Attention model



$\alpha^{<t, t'>}$ — amount of "attention" $y^{<t>}$ should pay to $a^{<t'>}$.

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

$$a^{<t>} = (\overrightarrow{a}^{<t>}, \overleftarrow{a}^{<t'>})$$

$$\sum_{t'} \alpha^{<1, t'>} = 1$$

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

$y^{<1>}$   $y^{<2>}$

$S^{<0>} \rightarrow S^{<1>} \rightarrow S^{<2>} \rightarrow \dots$

$c$ $\oplus$   $c^{<2>}$

$\alpha^{<1,1>}$   $\alpha^{<1,2>}$   $\alpha^{<1,3>}$

$\overrightarrow{a}^{<0>} \rightarrow$

$\overrightarrow{a}^{<1>}$ $\overleftarrow{a}^{<1>}$   $\overrightarrow{a}^{<2>}$ $\overleftarrow{a}^{<3>}$   $\overrightarrow{a}^{<3>}$ $\overleftarrow{a}^{<3>}$   $\overrightarrow{a}^{<5>}$ $\overleftarrow{a}^{<5>}$ $\leftarrow \overleftarrow{a}^{<6>}$

$t'$

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

jane   visite   l'Afrique   en   septembre

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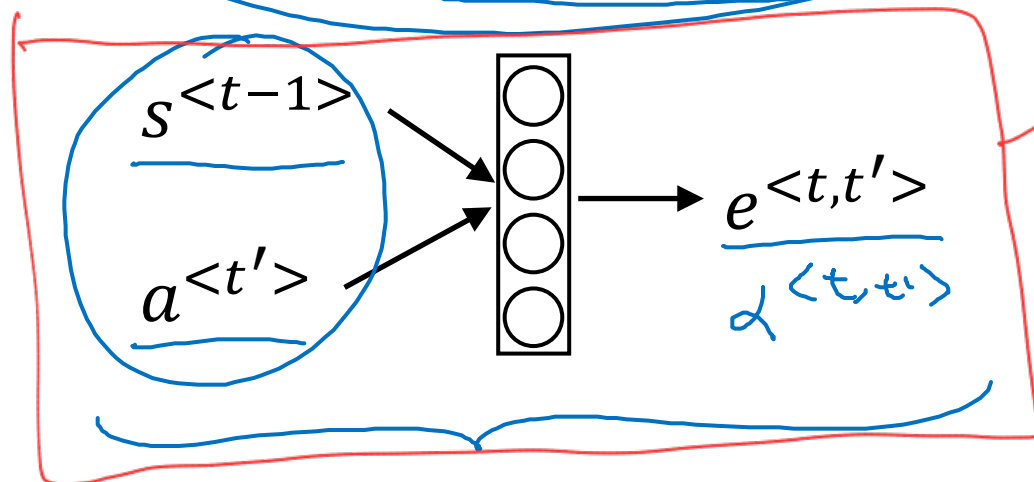
Andrew Ng

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

$T_x$    $T_y$

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

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

$s^{<t-1>}$

$a^{<t'>}$

$e^{<t,t'>}$

$\propto^{<t,t'>}$

this node is plugged into the whole ENC-DEC architecture

Attention Weights

$\hat{y}^{<t-1>}$    $\hat{y}^{<t>}$

$s^{<t-1>}$    $s^{<t>}$

$s^{<t-1>}$

AW  $e^{<t-1,1>}$  $\propto^{<t-1,1>}$  soft...

$a^{<1>}$

AW +in

$a^{<0>}$

$x^{<1>}$    $x^{<2>}$    ...    $x^{<T_x-1>}$    $x^{<T_x>}$

[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

July 20th 1969 $\longrightarrow$ $1969 - 07 - 20$

23 April, 1564 $\longrightarrow$ $1564 - 04 - 23$

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

deeplearning.ai

Audio data

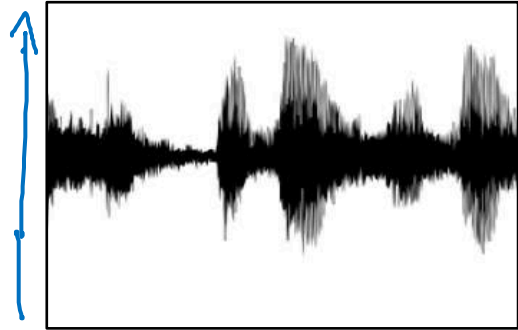Speech recognition

# Speech recognition problem
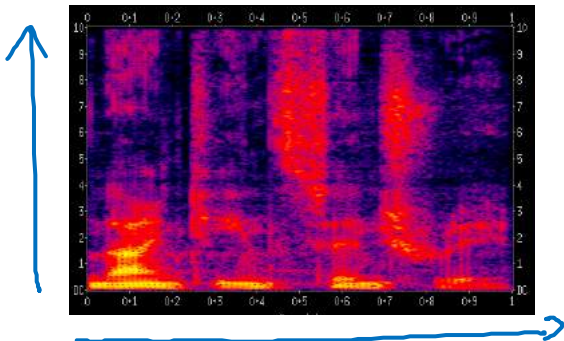
$x$

audio clip

$y$
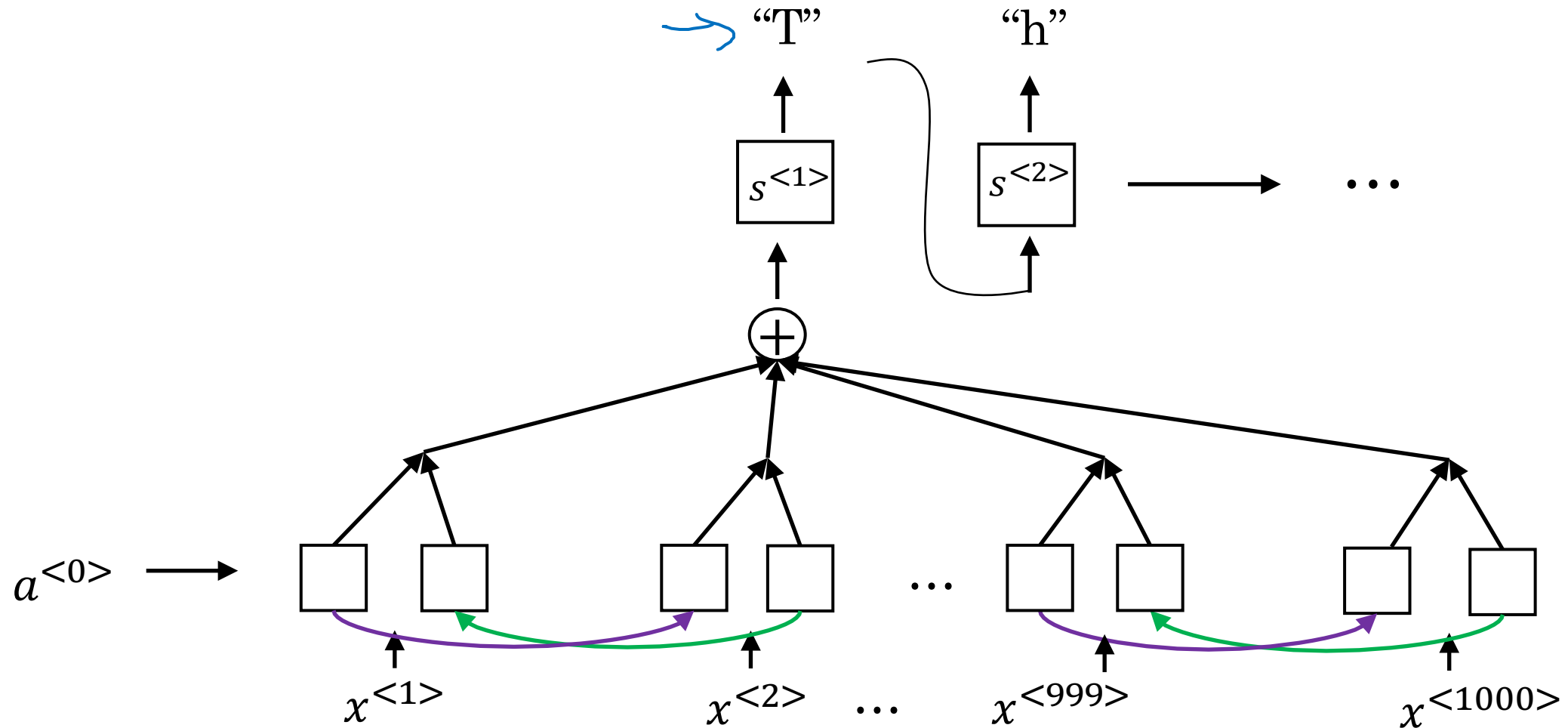
transcript



"the quick brown fox"

phonemes: de kwik braun

300h

3000h

100,000h

# Attention model for speech recognition

# CTC cost for speech recognition

(Connectionist temporal classification)

- Use some number of inputs and outputs
-

"the quick brown fox"    — 19 characters



$a^{<0>} \rightarrow$ [ ] $\rightarrow$ [ ] $\rightarrow \cdots \rightarrow$ [ ]

$\hat{y}^{<1>} \quad \hat{y}^{<2>} \quad \hat{y}^{<1000>}$

$x^{<1>} \quad x^{<2>} \quad x^{<1000>}$

ttt_h_eee - - - - ⊔ - - - - qqq - -        theuq

"space"      "blank"

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
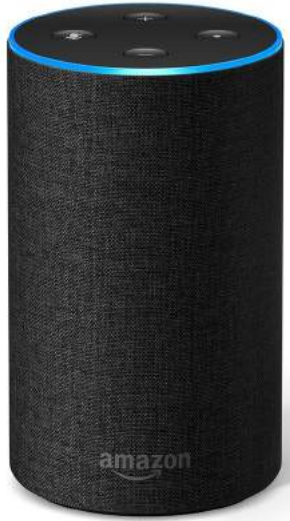
Audio data

Trigger word
detection

deeplearning.ai

# What is trigger word detection?

Amazon Echo
(Alexa)

Baidu DuerOS
(xiaodunihao)

Apple Siri
(Hey Siri)

Google Home
(Okay Google)

Andrew Ng

# Trigger word detection algorithm

# 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

Andrew Ng

# Deep learning is a super power

Please buy this from shutterstock and replace in final video.



Andrew Ng

Thank you.

- Andrew Ng