# Recurrent Neural Networks

---

# Why sequence models?

deeplearning.ai
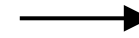
# Examples of sequence data

Speech recognition $\underset{x}{\longrightarrow}$ ⟿ (audio waveform) ⟶ "The quick brown fox jumped over the lazy dog." $y$

Music generation $\emptyset$ ⟶ (musical notation)
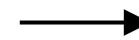
Sentiment classification "There is nothing to like in this movie." ⟶ ★☆☆☆☆

DNA sequence analysis → AGCCCCTGTGAGGAACTAG ⟶ AG<span style="color:red">CCCCTGTGAGGAACT</span>AG

Machine translation Voulez-vous chanter avec moi? ⟶ Do you want to sing with me?

Video activity recognition (images of runner) ⟶ Running

Name entity recognition → Yesterday, Harry Potter met Hermione Granger. ⟶ Yesterday, <span style="color:red">Harry Potter</span> met <span style="color:red">Hermione Granger</span>.

Andrew Ng

# Recurrent Neural Networks

---

# Notation

# Motivating example

NLP

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

$\rightarrow$   $x^{\langle 1 \rangle}$   $x^{\langle 2 \rangle}$   $x^{\langle 3 \rangle}$   - - - . .   $x^{\langle t \rangle}$   - - - .   $x^{\langle 9 \rangle}$

$T_x = 9$

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

$y^{\langle 1 \rangle}$   $y^{\langle 2 \rangle}$   $y^{\langle 3 \rangle}$   - - . . .   $y^{\langle 9 \rangle}$

$T_y = 9$

$X^{(i)\langle t \rangle}$     $T_x^{(i)} = 9$     15

$y^{(i)\langle t \rangle}$     $T_y^{(i)}$

Andrew Ng

# Representing words

$x^{<t>}$       $(x, y)$

$x \longrightarrow y$

x:        Harry Potter and Hermione Granger invented a new spell.
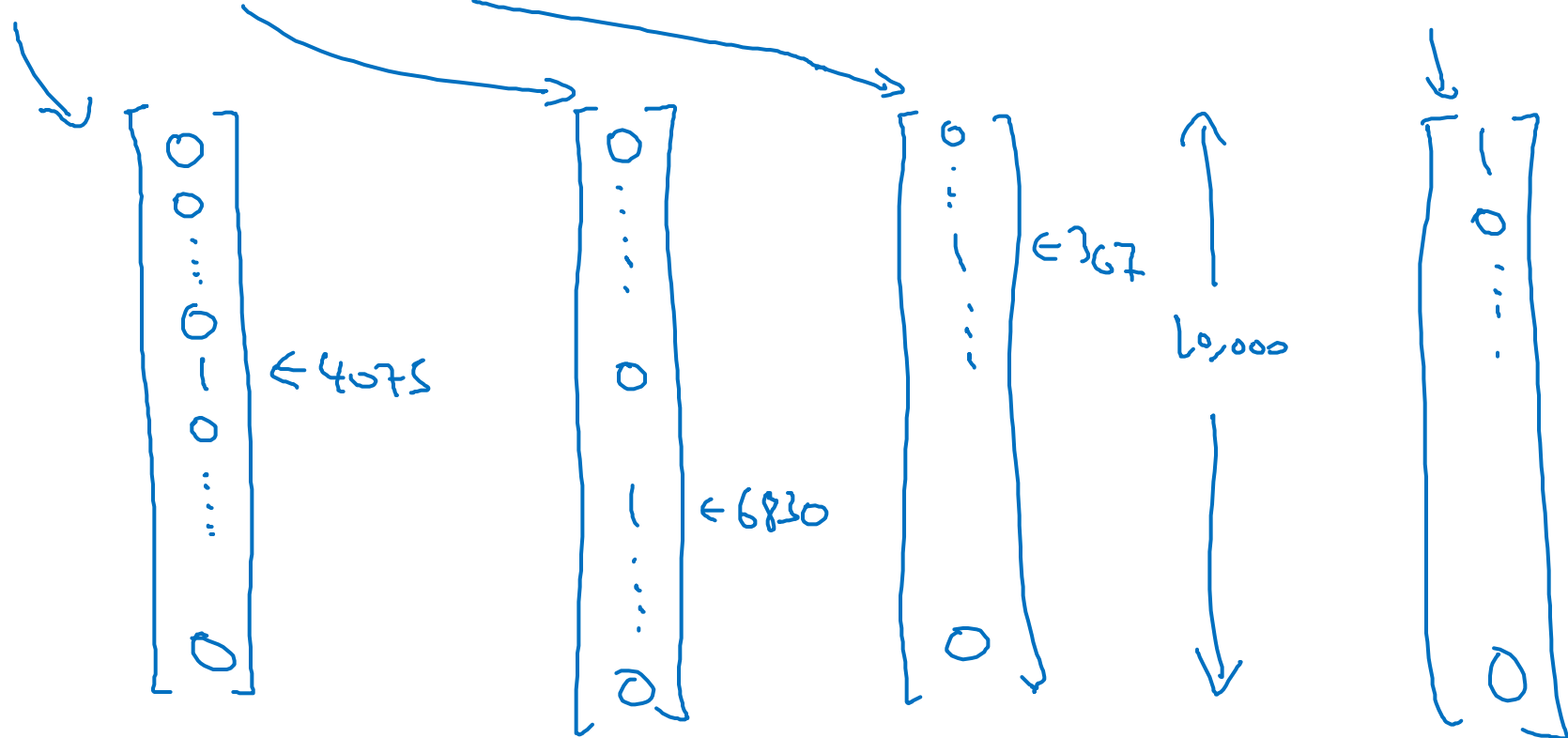
$x^{<1>}$   $x^{<2>}$   $x^{<3>}$        ...        $x^{<9>}$

$x^{<7>}$



Vocabulary

$\begin{bmatrix} a \\ aaron \\ \vdots \\ and \\ \vdots \\ harry \\ potter \\ \vdots \\ zulu \end{bmatrix}$ $\begin{matrix} 1 \leftarrow \\ 2 \\ \vdots \\ 367 \leftarrow \\ \vdots \\ 4075 \\ 6830 \\ \vdots \\ 10,000 \end{matrix}$

<UNK>    10,000

One-hot

Andrew Ng

# Representing words

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

$x^{<1>}$   $x^{<2>}$   $x^{<3>}$                          ...                          $x^{<9>}$

And = 367
Invented = 4700
A = 1
New = 5976
Spell = 8376
Harry = 4075
Potter = 6830
Hermione = 4200
Gran… = 4000
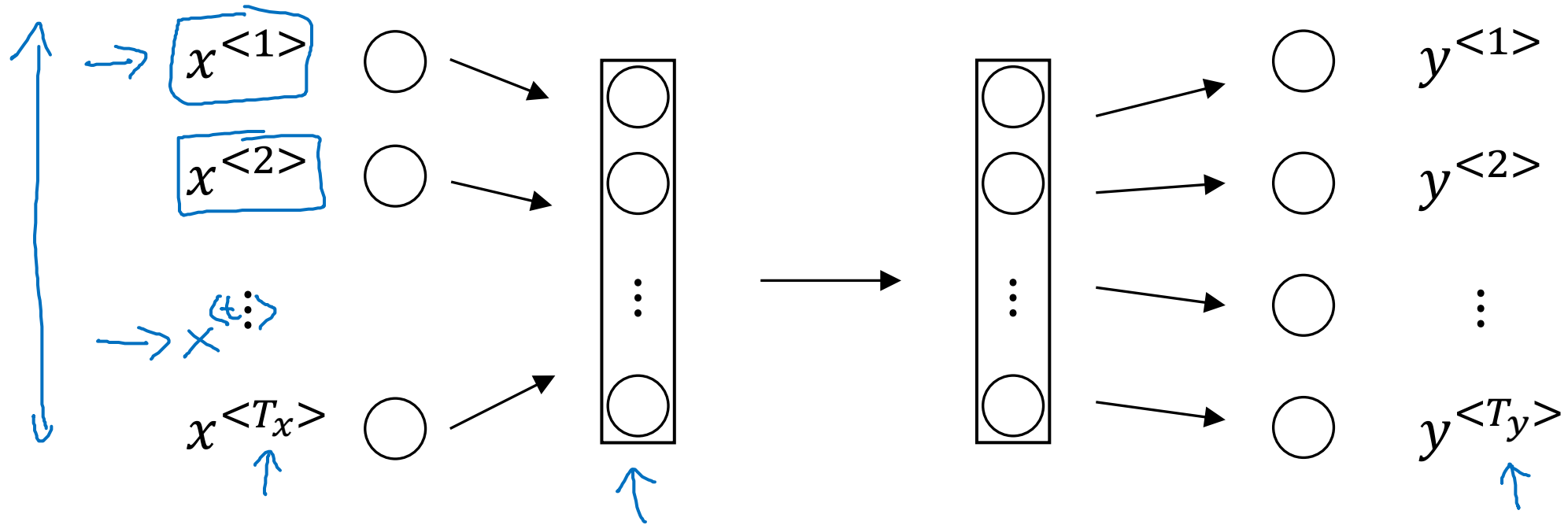
deeplearning.ai

# Recurrent Neural Networks
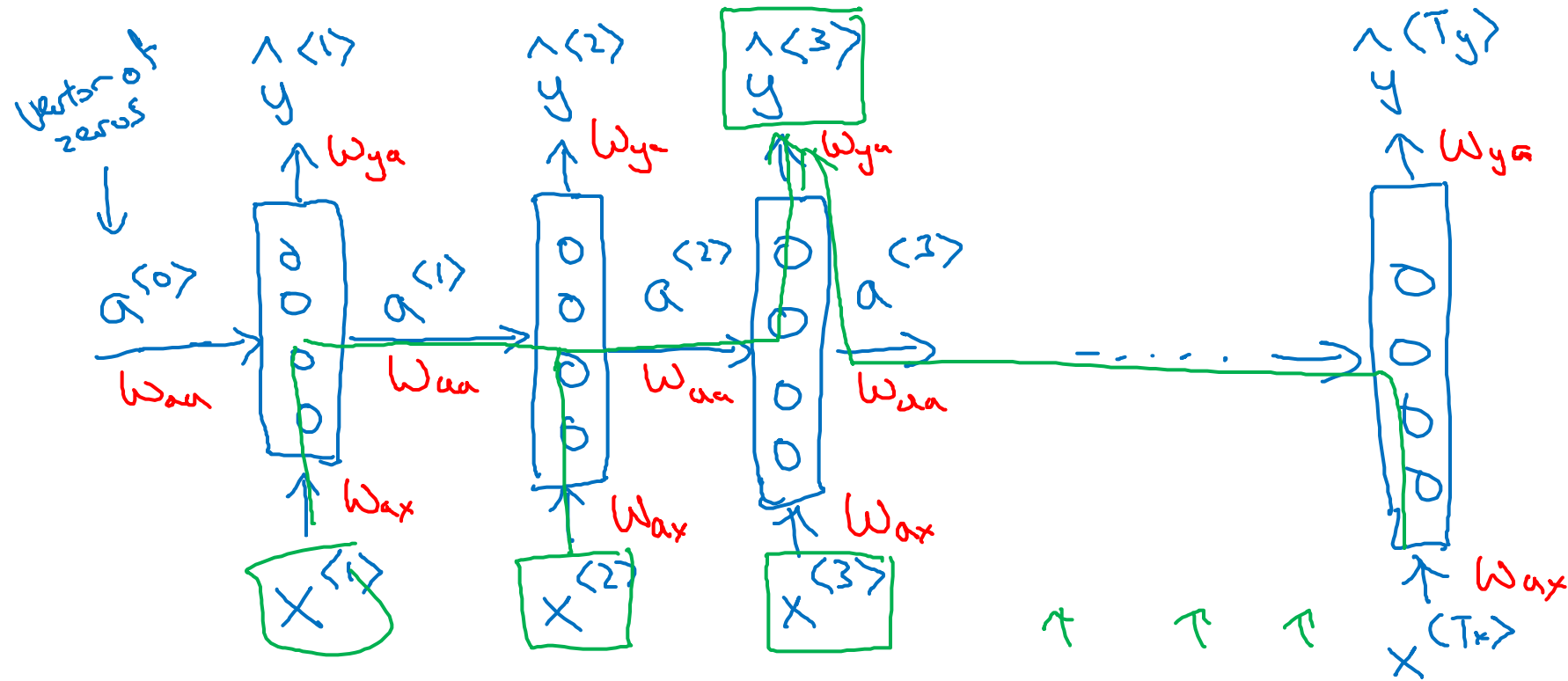
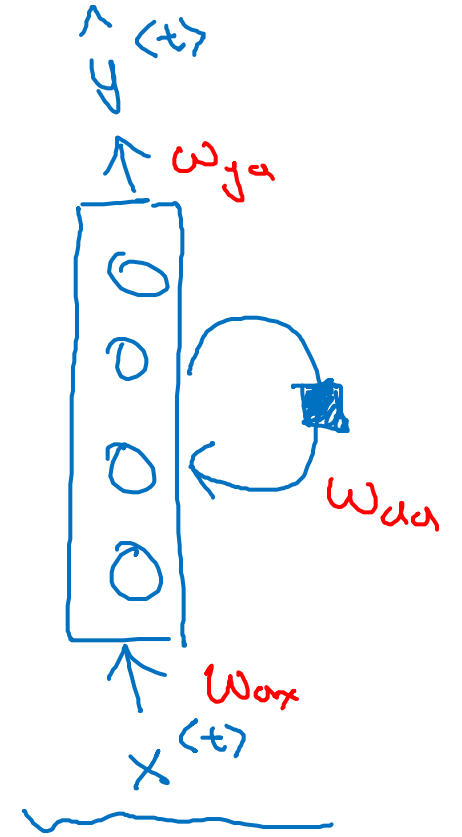Recurrent Neural Network Model

# 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.

# Recurrent Neural Networks



Bidirectional RNN (BRNN)

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

He said, "Teddy bears are on sale!"

# Forward Propagation

$$a \leftarrow W_{ax} x^{(1)}$$



$$a^{<0>} = \vec{0}.$$

$$a^{(1)} = g_1(W_{aa} a^{<0>} + W_{ax} x^{(1)} + b_a) \leftarrow \tanh / ReLu$$

$$\hat{y}^{(1)} = g_2(W_{ya} a^{(1)} + b_y) \leftarrow Sigmoid$$

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

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

Andrew Ng

# Simplified RNN notation

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

$(100,100)$    100   $(100,10,000)$   10,000

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

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

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

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

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

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

Andrew Ng

# Recurrent Neural Networks

## Backpropagation through time

deeplearning.ai

# Forward propagation and backpropagation

# Forward propagation and backpropagation



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

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

Backpropagation through time

Andrew Ng

# Recurrent Neural Networks

## Different types of RNNs

deeplearning.ai

# Examples of sequence data

$T_x$    $T_y$

$x$    $y$

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   AGCCCCTGTGAGGAACTAG → AGCCCCTGTGAGGAACTAG

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.

Andrew Ng

# Examples of RNN architectures

$T_x = T_y$

$\hat{y}^{(1)}$     $\hat{y}^{(2)}$     $\hat{y}^{(T_y)}$

$a^{(0)}$

$x^{(1)}$    $x^{(2)}$   $\cdots$   $x^{(T_x)}$

Many-to-many

Sentiment classification
$x = text$
$y = 0/1$     $1 \cdots 5$

$y$

$x^{(1)}$     $x^{(2)}$     $x^{(T_x)}$

There    is    $\cdots$    movie

Many-to-one

$y$

$x$

one-to-one

# Examples of RNN architectures



Music generation

$x \rightarrow y^{<1>} y^{<2>} \ldots y^{<T_y>}$

One-to-many

$x = \phi$

Machine translation

encoder

decoder

Many - to - many

# Summary of RNN types



One to one

One to many

Many to one

Many to many    $T_x = T_y$

Many to many

Andrew Ng

# Recurrent Neural Networks

deeplearning.ai

# Language model and sequence generation

# What is language modelling?

Speech recognition

The apple and <u>pair</u> salad.

→ The apple and <u>pear</u> salad.

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

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

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

$P\left(y^{<1>}, y^{<2>}, \ldots, y^{<T_y>}\right)$

Andrew Ng

# Language modelling with an RNN

Training set: large corpus of english text.

Tokenize

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

$y^{<1>}$     $y^{<2>}$     $y^{<3>}$    . . .    $y^{<8>}$    $y^{<9>}$

$x^{<t>} = y^{<t-1>}$

The Egyptian Mau is a bread of cat. <EOS>

<UNK>

10,000

# RNN model

$P(a) \; P(aaron) \cdots P(cats) \cdots P(zulu)$
$P(<UNK>)$
$P(<EOS>)$

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

$P(\_\_ \mid \text{"cats average"})$

$P(<EOS> \mid \cdots)$



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

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

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

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

$= P(y^{<1>}) \; P(y^{<2>} \mid y^{<1>})$
$\quad P(y^{<3>} \mid y^{<1>}, y^{<2>})$

Andrew Ng

Recurrent Neural Networks

Sampling novel sequences

deeplearning.ai

# Sampling a sequence from a trained RNN



$P(y^{(1)}, ..., y^{(T_x)})$

Training:

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

$a^{<0>}$ → $a^{<1>}$ → $a^{<2>}$ → $a^{<3>}$ → ... → $a^{<T_y>}$

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

Sampling:

The $\hat{y}^{<1>}$   $\hat{y}^{<2>}$   $\hat{y}^{<3>}$   $\hat{y}^{<T_y>}$

$a^{(0)} = 0$ → $a^{(1)}$ → $a^{(2)}$ → $a^{(3)}$ → ... → 

$x^{(1)} = 0$   $x^{(2)} = \hat{y}^{(1)}$   $y^{(T_x-1)}$
The

$<EOS>$

$<UNK>$

→ $P(a) P(aaron) ... P(zulu) P(<UNK>)$   n.p.random.choice   $P(\_ | the)$

Andrew Ng

# Character-level language model

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

$\rightarrow$ Vocabulary = [ a, b, c, ..., z, ␣, ., , , ;, 0, ..., 9, A, ..., Z]

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

Cat ↑↑↑↑ average ...

May

# Sequence generation

## News

## Shakespeare

President enrique peña nieto, announced sench's sulk former coming football langston paring.

"I was not at all surprised," said hich langston.

"Concussion epidemic", to be examined.

The gray football the told some and this has on the uefa icon, should money as.

The mortal moon hath her eclipse in love.

And subject of this thou art another this fold.

When besser be my love to me see sabl's.

For whose are ruse of mine eyes heaves.

# Recurrent Neural Networks

deeplearning.ai

## Vanishing gradients with RNNs

# Vanishing gradients with RNNs

The (cat), which ately ate ........., was full

The cats, Wh      —— ...... ...... were full

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



$a^{<0>}$   $a^{<1>}$   $a^{<2>}$   $a^{<3>}$   $\cdots$   $a^{<T_y>}$

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

$x$   $\cdots$   $\hat{y}$

100

Exploding gradients.

NaN      Gradient clipping

# RNN unit



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

tanh

# GRU (simplified)

$C^{\langle t-1 \rangle}$
$= a^{\langle t-1 \rangle}$

softmax → $y^{\langle t \rangle}$

$\tilde{C}^{\langle t \rangle}$  $\Gamma_u$

tanh  $\sigma$

$x^{\langle t \rangle}$

$C^{\langle t \rangle}$
$= a^{\langle t \rangle}$

$\Gamma_u = 1$
$C^{\langle t \rangle} = 1$
$\Gamma_u = 0$  $\Gamma_u = 0$  $\Gamma_u = 0$ ......  $\downarrow = 1$

The cat, which already ate …, was full.

$C$ = memory cell

$C^{\langle t \rangle} = a^{\langle t \rangle}$

$$\tilde{C}^{\langle t \rangle} = \tanh\left(W_c\left[c^{\langle t-1 \rangle}, x^{\langle t \rangle}\right] + b_c\right)$$

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

"update"

$$C^{\langle t \rangle} = \Gamma_u * \tilde{C}^{\langle t \rangle} + (1-\Gamma_u) * C^{\langle t-1 \rangle}$$

$\Gamma_u = 1$

element-wise

$\Gamma_u = 0.000001$

Gate

[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

# Full GRU

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

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

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

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

LSTM

The cat, which ate already, was full.

# Recurrent Neural Networks

# LSTM (long short term memory) unit

# GRU and LSTM

## GRU

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

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

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

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

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

## LSTM

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

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

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

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

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

$$a^{<t>} = \Gamma_o * c^{<t>}$$

$\Gamma_f$

[Hochreiter & Schmidhuber 1997. Long short-term memory]          Andrew Ng

# LSTM units

### GRU

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

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

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

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

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

### LSTM

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

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

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

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

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

$$a^{<t>} = \Gamma_o * c^{<t>}$$

Andrew Ng

# LSTM in pictures

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

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

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

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

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

$$a^{<t>} = \Gamma_o * c^{<t>}$$

$c^{(t-1)}$ — peephole connection



Andrew Ng

Recurrent Neural Networks

Bidirectional RNN

deeplearning.ai

# Getting information from the future

He said, "Teddy bears are on sale!"

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



$\hat{y}^{<1>}$  $\hat{y}^{<2>}$  $\hat{y}^{<3>}$  $\hat{y}^{<4>}$  $\hat{y}^{<5>}$  $\hat{y}^{<6>}$  $\hat{y}^{<7>}$

$a^{<0>}$  $a^{<1>}$  $a^{<2>}$  $a^{<3>}$  $a^{<4>}$  $a^{<5>}$  $a^{<6>}$  $a^{<7>}$

RNN
GRU
LSTM

$x^{<1>}$  $x^{<2>}$  $x^{<3>}$  $x^{<4>}$  $x^{<5>}$  $x^{<6>}$  $x^{<7>}$

He  said,  "Teddy  bears  are  on  sale!"

# Bidirectional RNN (BRNN)

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



GRU
LSTM

$\hat{y}^{(1)}$    $\hat{y}^{(2)}$    $\hat{y}^{(3)}$    $\hat{y}^{(4)}$

$\overrightarrow{a}^{(1)}$    $\overleftarrow{a}^{(1)}$    $\overrightarrow{a}^{(2)}$    $\overleftarrow{a}^{(2)}$    $\overrightarrow{a}^{(3)}$    $\overleftarrow{a}^{(3)}$    $\overrightarrow{a}^{(4)}$    $\overleftarrow{a}^{(4)}$

$x^{(1)}$    $x^{(2)}$    $x^{(3)}$    $x^{(4)}$

Acyclic graph

He said, "Teddy Roosevelt ..."

BRNN w/ LSTM

Andrew Ng

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

$\longrightarrow$  Jane  is  visiting  Africa  in  September.

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



[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

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

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

$3 \times 3$
$s = 2$

$6 \times 6 \times 256$

$9216$

$4096$

$4096$

Softmax
1000

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

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

$\cdots$

$x$

[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

deeplearning.ai

Sequence to sequence models

Picking the most likely sentence

# Machine translation as building a conditional language model



Language model:

$a^{<0>}$

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

$x^{<1>}, x^{<2>}, x^{<T_y>}$

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

Machine translation:

$a^{<0>}$

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

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

"Conditional language model"

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

Andrew Ng

# Finding the most likely translation

Jane visite l'Afrique en septembre.

English

French

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

→ 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.

$$\underset{y^{<1>}, \ldots, y^{<T_y>}}{\arg\max} P(y^{<1>}, \ldots, y^{<T_y>} \mid x)$$

# Why not a greedy search?

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



$$\arg\max_{y} P(\hat{y}^{<1>}, \hat{y}^{<2>}, \ldots, \hat{y}^{<T_y>}|x)$$

$10,000$

$10$

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

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

⟶ Jane is visiting Africa in September.

⟶ Jane is going to be visiting Africa in September.

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

Sequence to sequence models

Beam search

deeplearning.ai

# Beam search algorithm

$B = 3$   (beam width)

## Step 1

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



$\begin{bmatrix} a \\ \vdots \\ in \\ \vdots \\ jane \\ \vdots \\ september \\ \vdots \\ zulu \end{bmatrix}$  10000

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow \square \rightarrow \hat{y}^{<1>}$

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

# Beam search algorithm

$(B = 3)$

**Step 1**   **Step 2**

10000

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

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

a, aaron, September, visit, zulu

a, aaron, is, visit, zulu

10,000

a, ..., zulu



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

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

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

Andrew Ng

# Beam search ($B = 3$)

$B = 1 \rightsquigarrow$ greedy search

in september

a
aaron
jane
zulu

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

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

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

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

jane is

a
visits
zulu

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

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

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

jane visits

a
africa
zulu

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

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

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

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

jane visits africa in september. <EOS>

Andrew Ng

Sequence to sequence models

Refinements to beam search

deeplearning.ai

# Length normalization

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

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

$\log$

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

$P(y|x) \leftarrow$

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

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

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

$\alpha = 0.7$

$\alpha = 1$

$\alpha = 0$

Andrew Ng

# Beam search discussion

Beam width B?

large B: better result, slower

small B: worse result, faster

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

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 $\arg \max\limits_{y} P(y|x)$.

deeplearning.ai

Sequence to sequence models

Error analysis on beam search

# Example

Jane visite l'Afrique en septembre.

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

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

→ RNN

→ Beam Search

$\boxed{BT}$

RNN computes $P(y^* | x) \overset{>}{\underset{\leq}{}} P(\hat{y} | x)$

# Error analysis on beam search

$P(y^* \mid x)$

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

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

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

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

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

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

Conclusion: Beam search is at fault.

Case 2: $P(y^* \mid x) \lesssim P(\hat{y} \mid x) \leftarrow$

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

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. | $2 \times 10^{-10}$ | $1 \times 10^{-10}$ | 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. $(\hat{y})$

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

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

[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

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

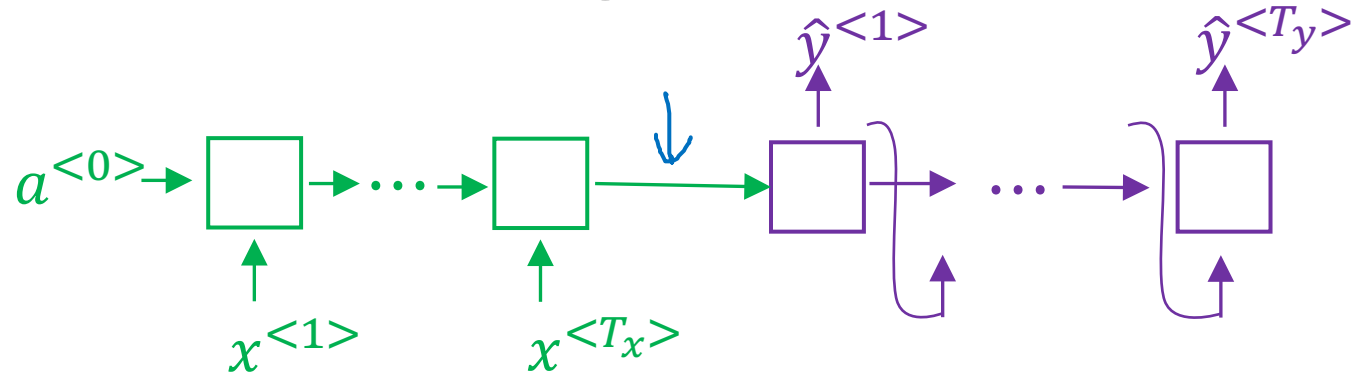[Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]     Andrew Ng

Sequence to sequence models

Attention model intuition

deeplearning.ai

# The problem of long sequences



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

$$x^{<1>} \qquad x^{<T_x>} \qquad \hat{y}^{<1>} \qquad \hat{y}^{<T_y>}$$
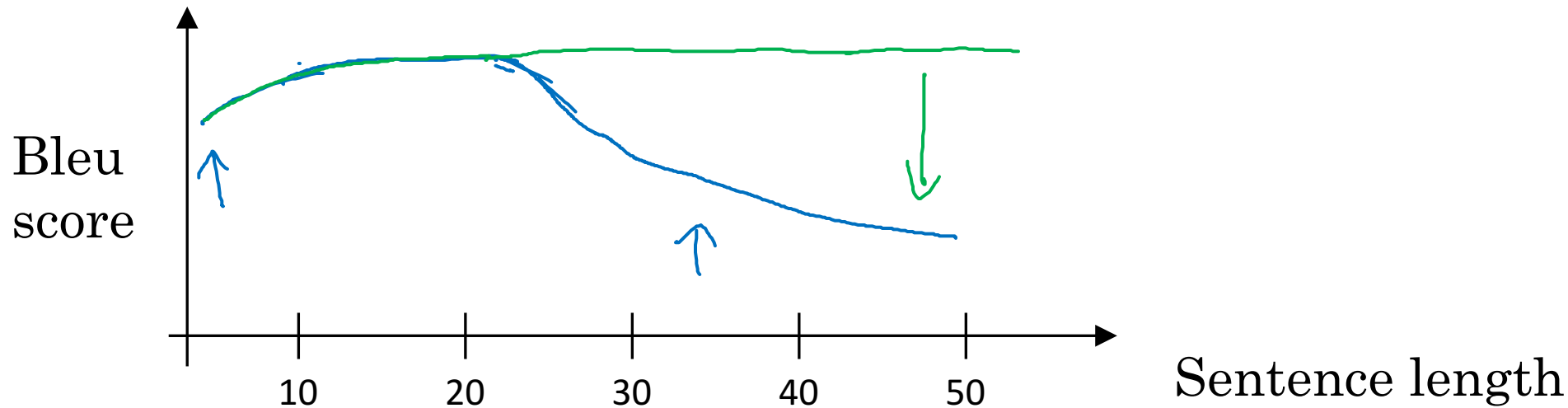
Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente d'y aller aussi.

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
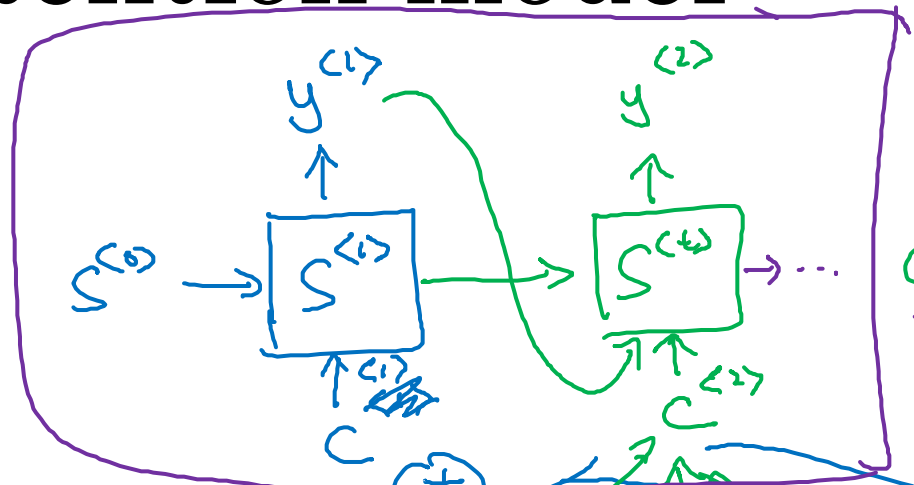
Andrew Ng

# Attention model intuition



[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

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

$y^{\langle 1\rangle}$   $y^{\langle 2\rangle}$

$S^{\langle 0\rangle} \rightarrow S^{\langle 1\rangle} \rightarrow S^{\langle t\rangle} \rightarrow \cdots$

$C^{\langle 2\rangle} = \sum_{t'} \alpha^{\langle 2, t'\rangle} a$

$a^{\langle t\rangle} = (\overrightarrow{a}^{\langle t\rangle}, \overleftarrow{a}^{\langle t\rangle})$

$\langle 1\rangle$

$C^{\langle 2\rangle}$

$C \oplus$

$\sum_{t'} \alpha^{\langle 1, t'\rangle} = 1$

$\alpha^{\langle 1,2\rangle}$   $\alpha^{\langle 1,3\rangle}$

$\alpha^{\langle 1,1\rangle}$

$C^{\langle 1\rangle} = \sum_{t'} \alpha^{\langle 1, t'\rangle} a^{\langle t'\rangle}$

$\overrightarrow{a}^{\langle 0\rangle} \rightarrow$    $\overrightarrow{a}^{\langle 1\rangle}$ $\overleftarrow{a}^{\langle 1\rangle}$    $\overrightarrow{a}^{\langle 2\rangle}$ $\overleftarrow{a}^{\langle 3\rangle}$    $\overrightarrow{a}^{\langle 3\rangle}$ $\overleftarrow{a}$    $\overrightarrow{a}^{\langle 5\rangle}$ $\overleftarrow{a}^{\langle 5\rangle}$ $\leftarrow \overleftarrow{a}^{\langle 6\rangle}$

$t'$

$x^{\langle 1\rangle}$         $x^{\langle 2\rangle}$         $x^{\langle 3\rangle}$         $x^{\langle 4\rangle}$         $x^{\langle 5\rangle}$

jane        visite        l'Afrique        en        septembre

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

Andrew Ng

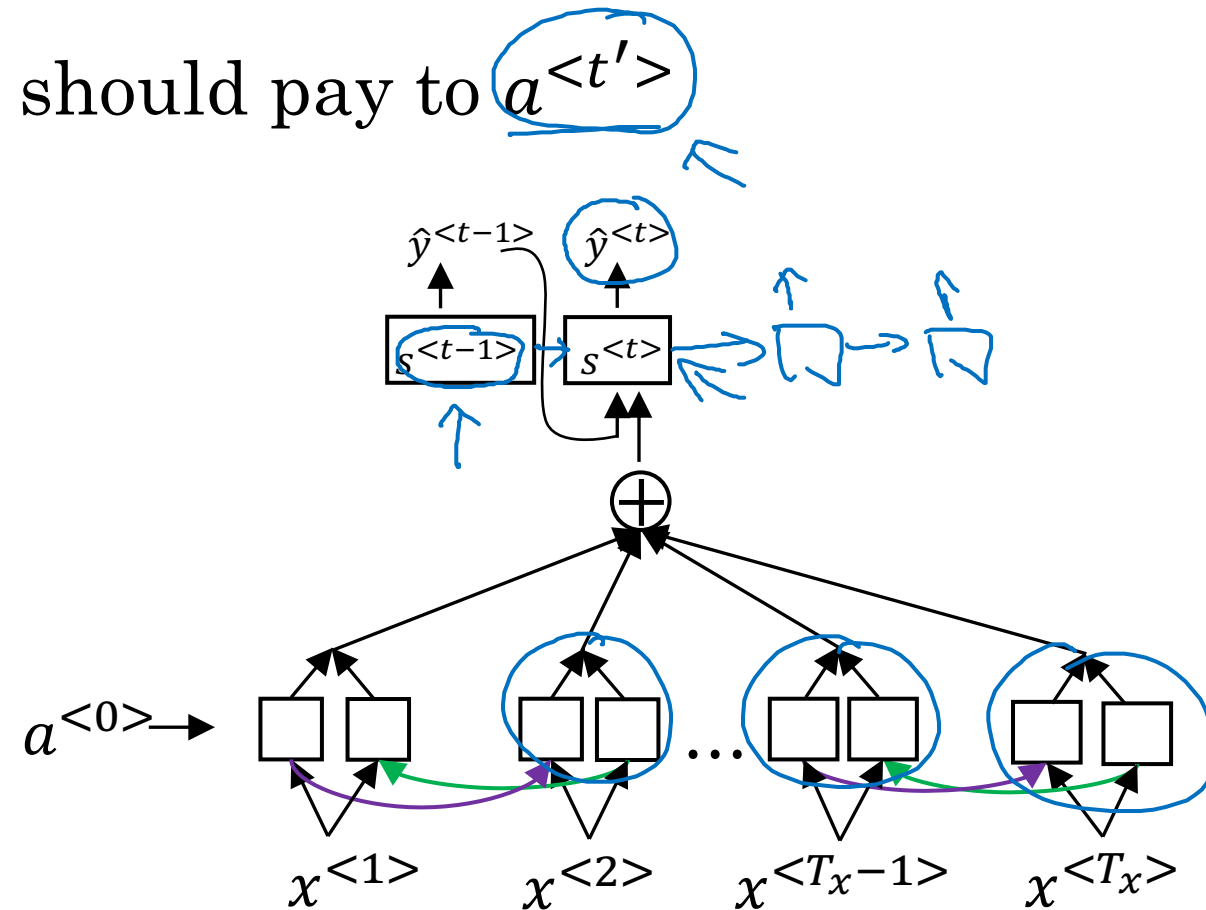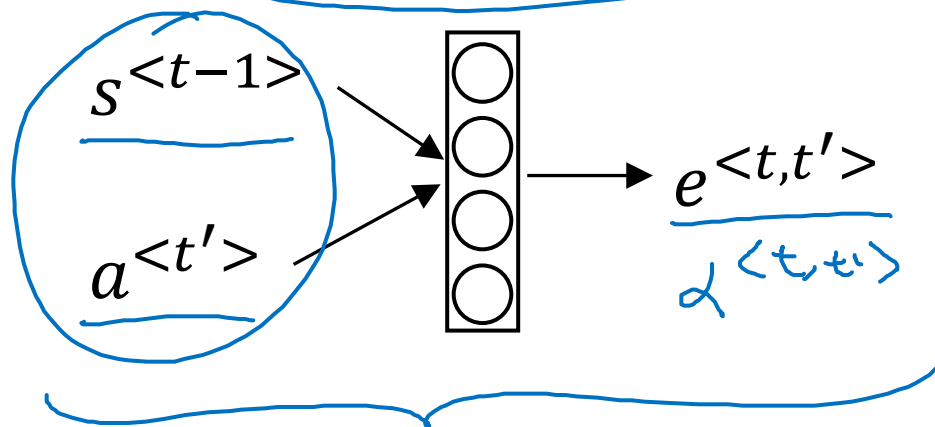# Computing attention $\alpha^{<t,t'>}$

$T_x$   $T_y$

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

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



$s^{<t-1>}$

$a^{<t'>}$

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

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

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

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

$a^{<0>}$

$x^{<1>}$   $x^{<2>}$   $\cdots$   $x^{<T_x-1>}$   $x^{<T_x>}$

[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]
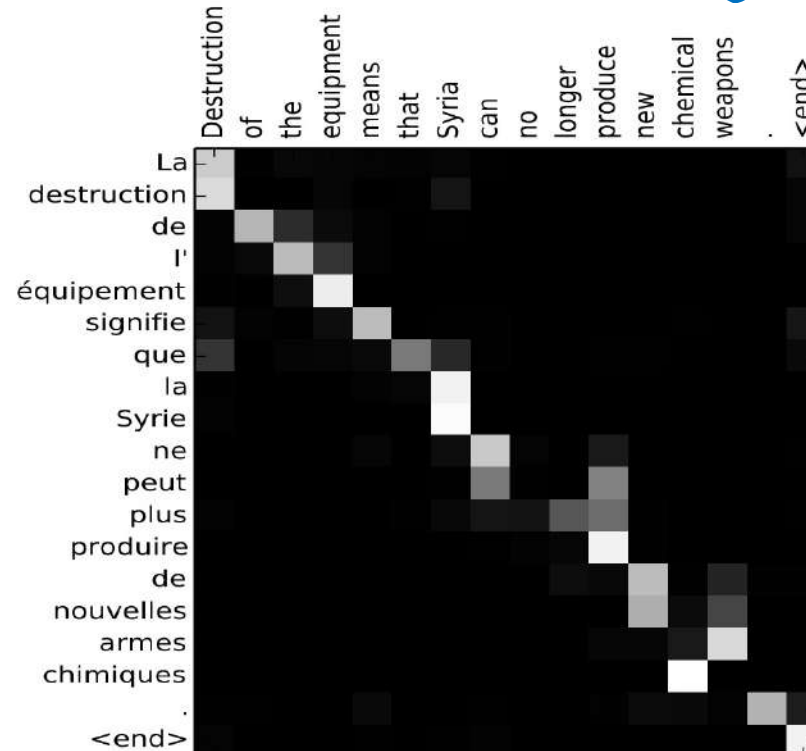
Andrew Ng

# Attention examples

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

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

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

Audio data

Speech recognition

deeplearning.ai

# Speech recognition problem

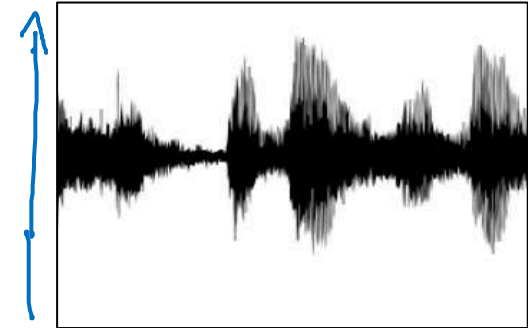$x$                $y$

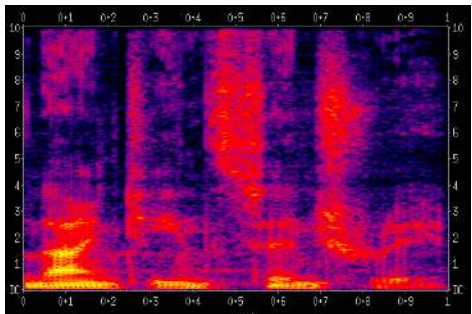audio clip        transcript
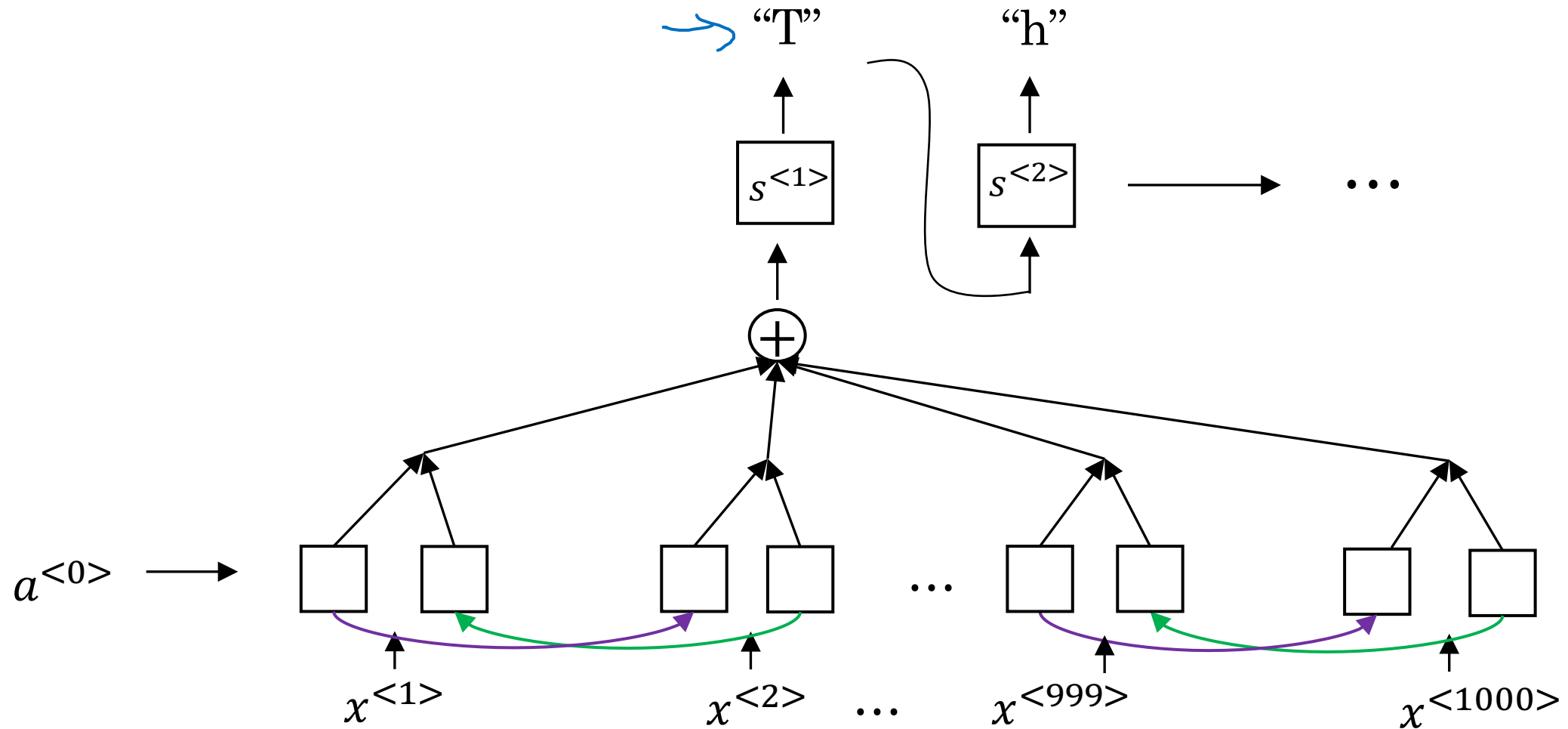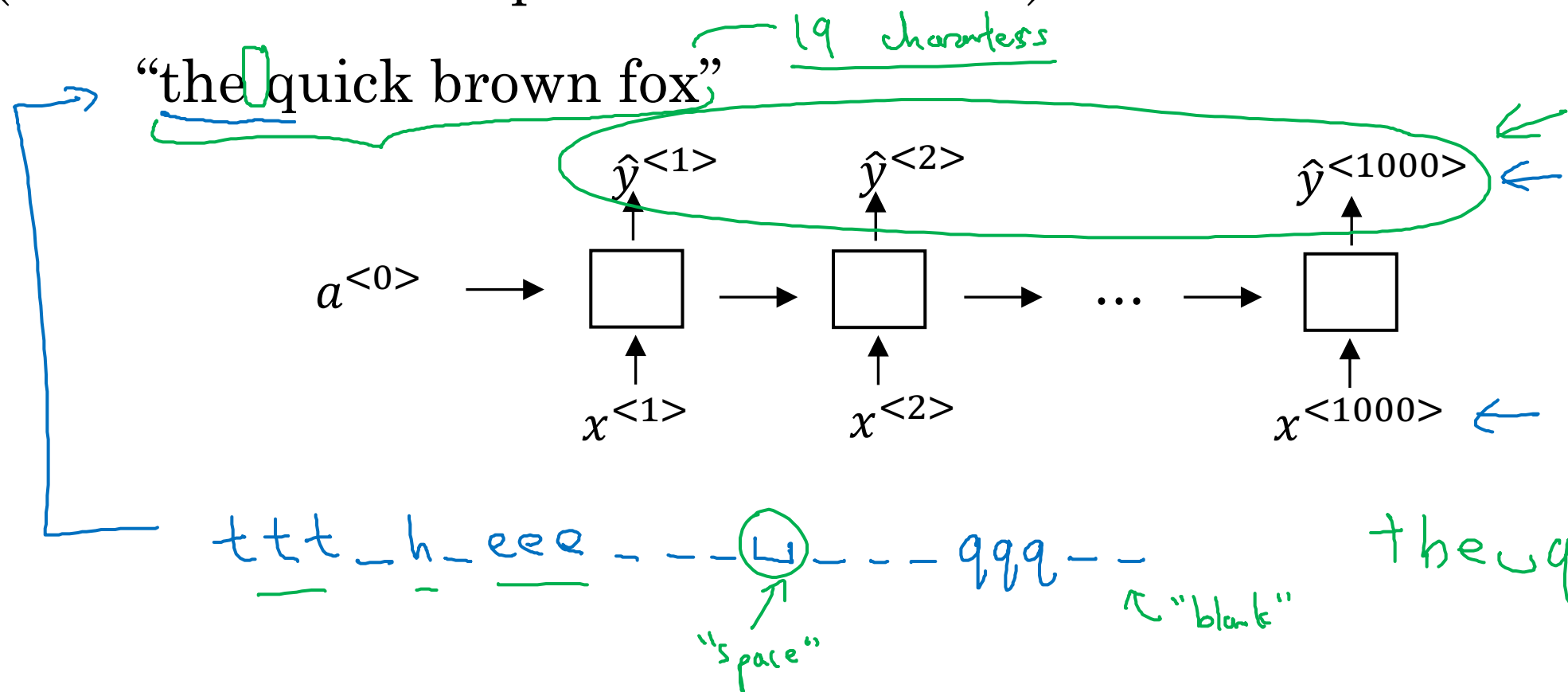


"the quick brown fox"

phonemes: de kwik braun

300h

3000h

100,000h

Andrew Ng

# Attention model for speech recognition

# CTC cost for speech recognition

(Connectionist temporal classification)

"the quick brown fox"    19 characters



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

$a^{<0>} \rightarrow$  $\boxed{\phantom{x}} \rightarrow \boxed{\phantom{x}} \rightarrow \cdots \rightarrow \boxed{\phantom{x}}$

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

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
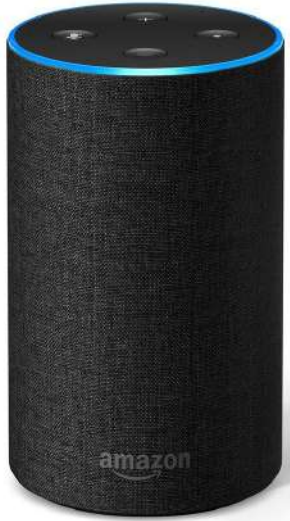
deeplearning.ai

Audio data

Trigger word
detection

# What is trigger word detection?



Amazon Echo
(Alexa)

Baidu DuerOS
(xiaodunihao)

Apple Siri
(Hey Siri)

Google Home
(Okay Google)

Andrew Ng

# Trigger word detection algorithm