# Recurrent Neural Networks

## Why sequence models?

# Examples of sequence data

Speech recognition $\;x\;\longrightarrow\;$ "The quick brown fox jumped over the lazy dog." $\;y$

Music generation $\;\emptyset\;\longrightarrow$

Sentiment classification  "There is nothing to like in this movie." $\;\longrightarrow\;$ ★☆☆☆☆

DNA sequence analysis  AGCCCCTGTGAGGAACTAG $\;\longrightarrow\;$ AGCCCCTGTGAGGAACTAG

Machine translation  Voulez-vous chanter avec moi? $\;\longrightarrow\;$ Do you want to sing with me?

Video activity recognition  $\longrightarrow$ Running

Name entity recognition  Yesterday, Harry Potter met Hermione Granger. $\;\longrightarrow\;$ Yesterday, Harry Potter met Hermione Granger.

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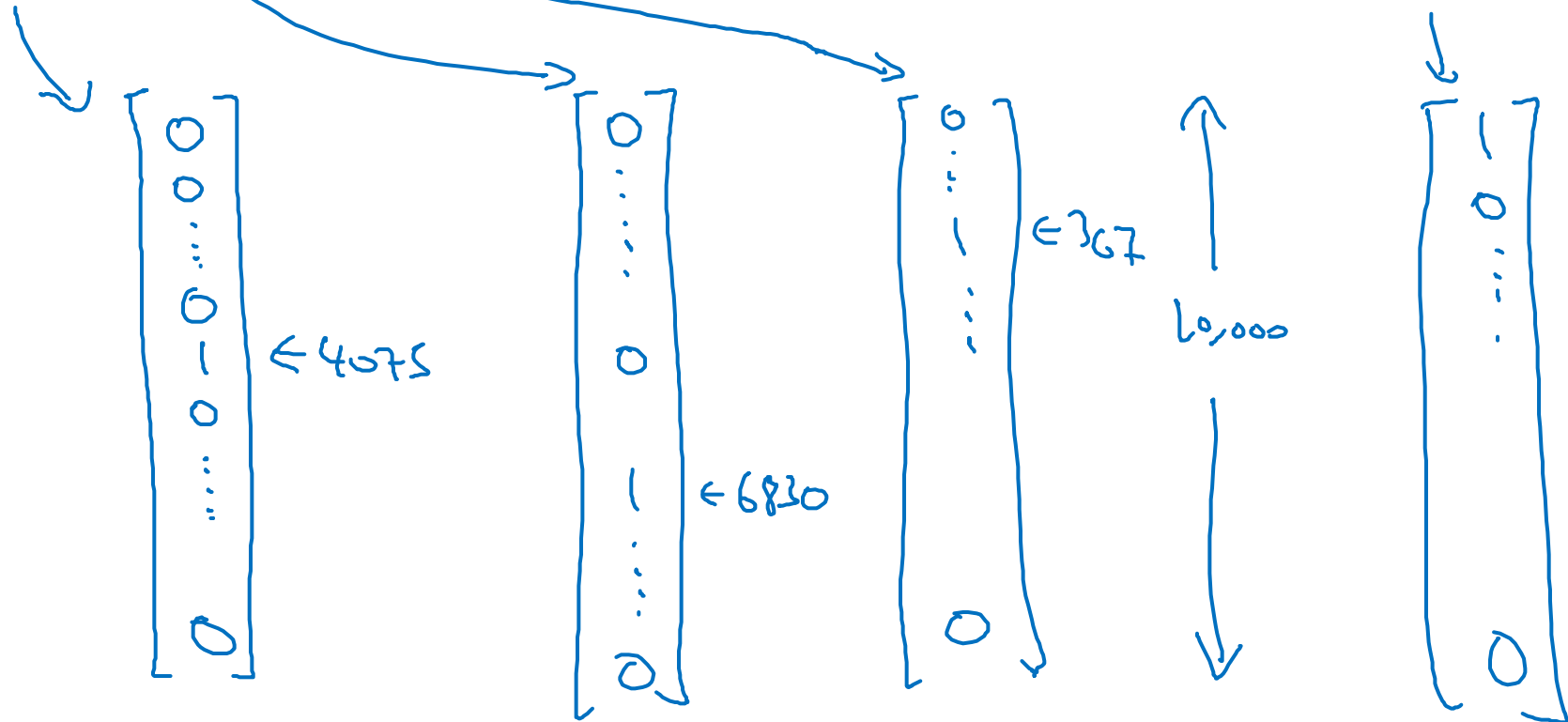
$x^{<1>}$   $x^{<2>}$   $x^{<3>}$      ...    $x^{<t>}$    $x^{<9>}$

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

$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$ ← 4075    $\begin{bmatrix} 0 \\ \vdots \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$ ← 6830    $\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$ ← 367   10,000    $\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$

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
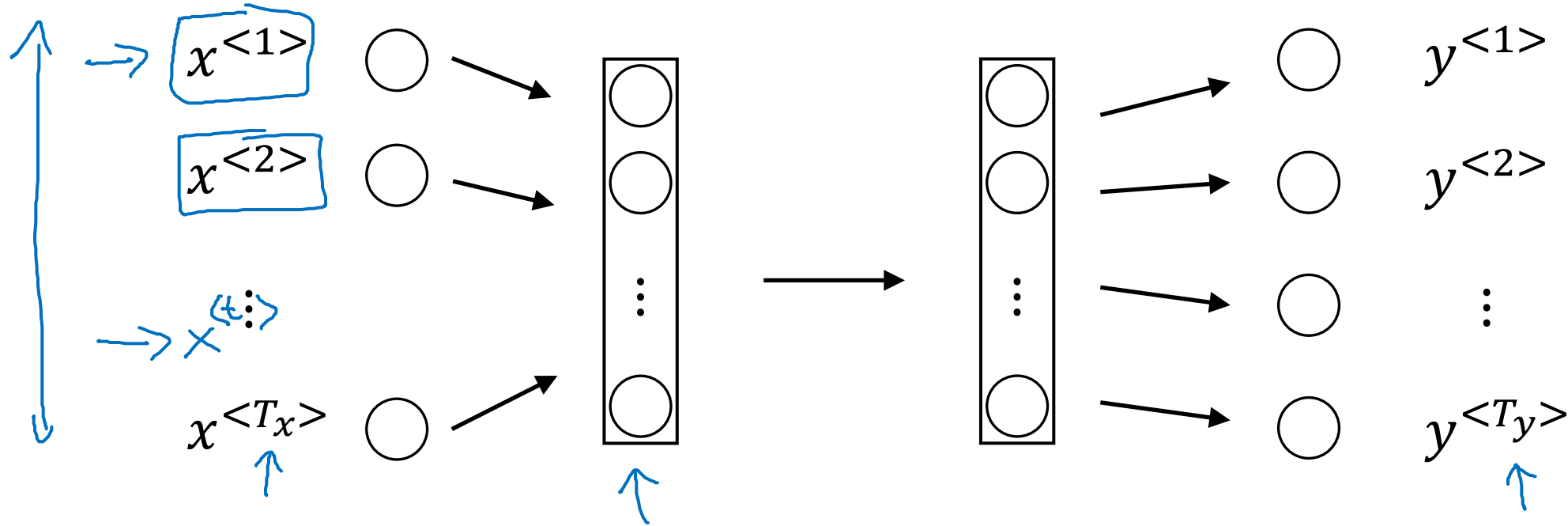Potter = 6830
Hermione = 4200
Gran... = 4000

Andrew Ng

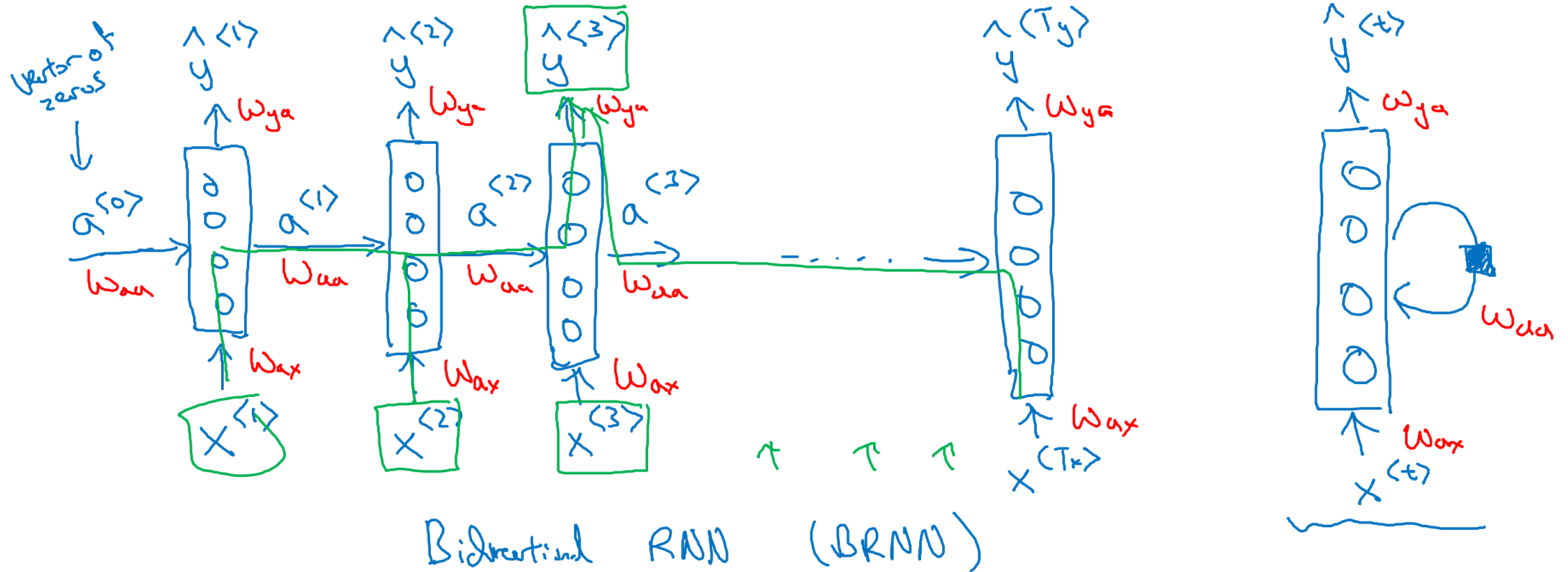deeplearning.ai

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.

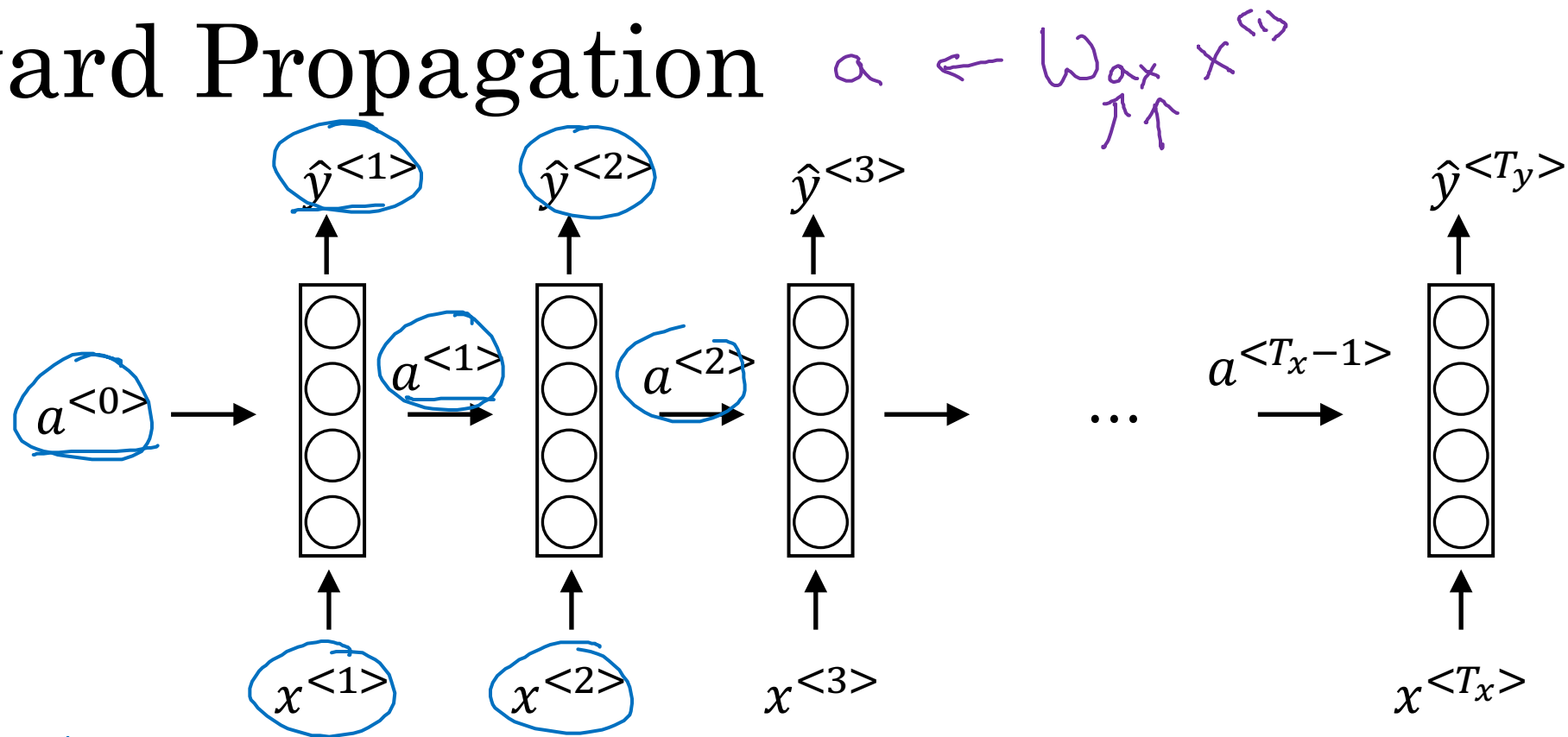# Recurrent Neural Networks

$T_x = T_y$



Bidirectional RNN (BRNN)

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

He said, "Teddy bears are on sale!"

Andrew Ng

# 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   10,000

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

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

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

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

$$\begin{bmatrix} W_{aa} & W_{ax} \end{bmatrix} = W_a$$

100   10000   (100, 10100)

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

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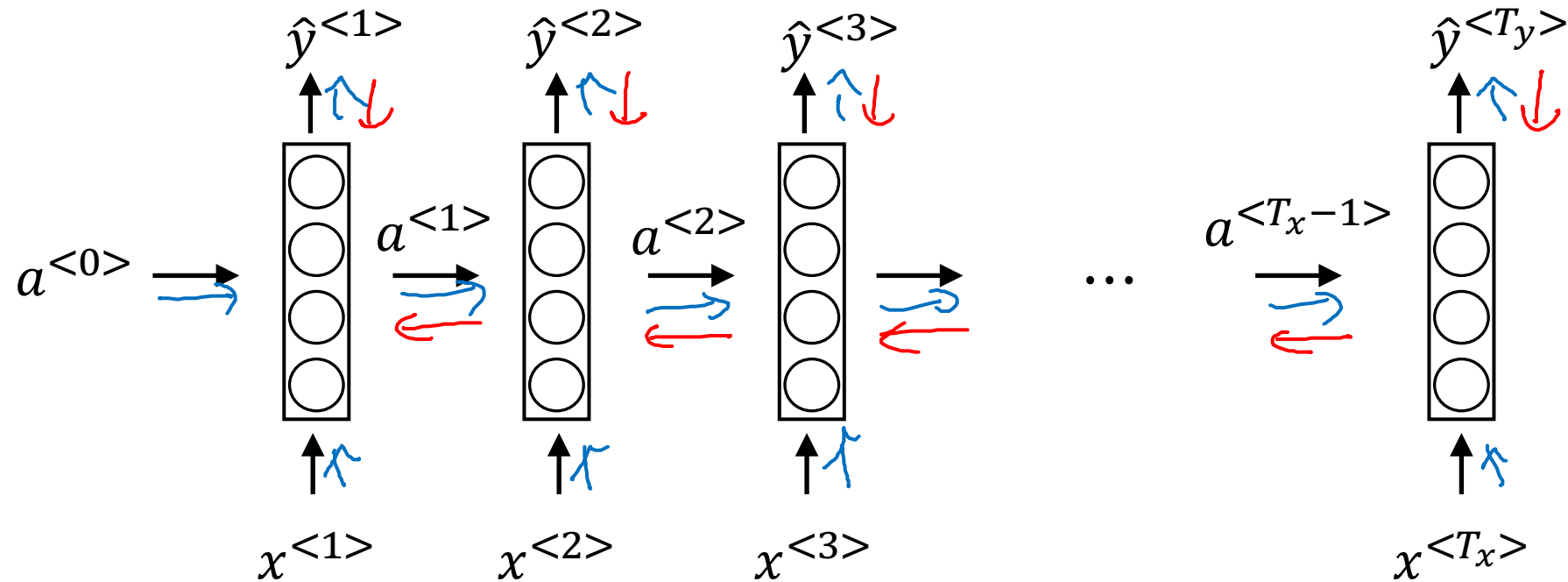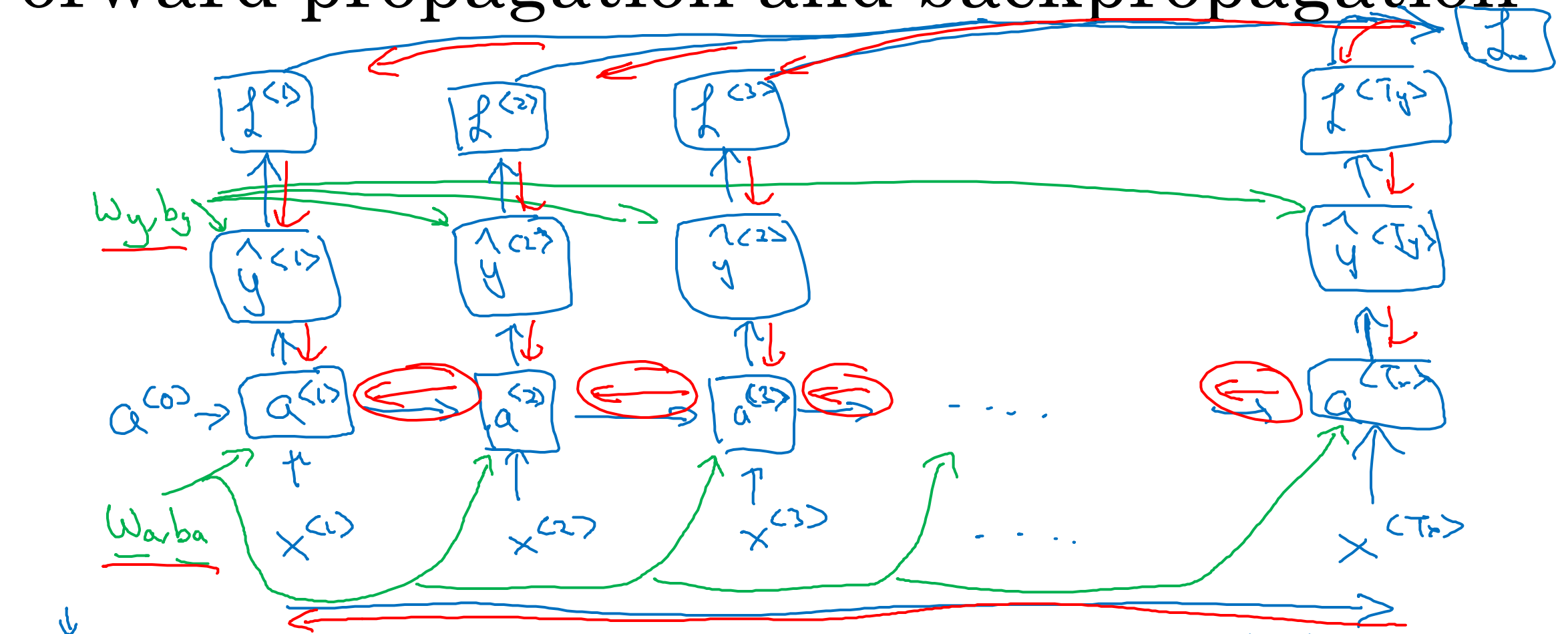
# Recurrent Neural Networks

deeplearning.ai

# Backpropagation through time

# Forward propagation and backpropagation

# Forward propagation and backpropagation



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

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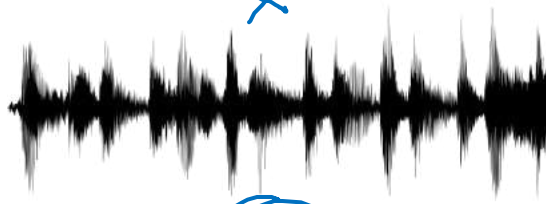
Backpropagation through time

Andrew Ng
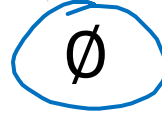
Recurrent Neural Networks

---

Different types
of RNNs

deeplearning.ai

# Examples of sequence data

$T_x$  $T_y$

Speech recognition
$x$

[audio waveform] $\longrightarrow$ "The quick brown fox jumped over the lazy dog."

$y$

Music generation
$\emptyset$ $\longrightarrow$ [musical notes]

Sentient classification
"There is nothing to like in this movie." $\longrightarrow$ ★☆☆☆☆

DNA sequence analysis
AGCCCCTGTGAGGAACTAG $\longrightarrow$ AGCCCCTGTGAGGAACTAG

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

Video activity recognition
[images of runner] $\longrightarrow$ Running

Name entity recognition
Yesterday, Harry Potter met Hermione Granger. $\longrightarrow$ Yesterday, Harry Potter met Hermione Granger.

Andrew Ng

# Examples of RNN architectures

$T_x = T_y$

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



Many-to-many

There is $\cdots\cdots$ movie

Many-to-one

one-to-one

Andrew Ng

# 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



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

$a^{<0>} \to$

$x^{<1>}$

One to one

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

$a^{<0>} \to$

$x$

One to many

$\hat{y}$

$a^{<0>} \to$

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

Many to one

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

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

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

Many to many $\quad T_x = T_y$

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

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

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

Many to many

Andrew Ng

Recurrent Neural Networks

Language model and sequence generation

deeplearning.ai

# What is language modelling?

Speech recognition

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

$\rightarrow$ The apple and <u>pear</u> salad.

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

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

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

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

# Language modelling with an RNN

Training set: large corpus of english text.

Tokenize

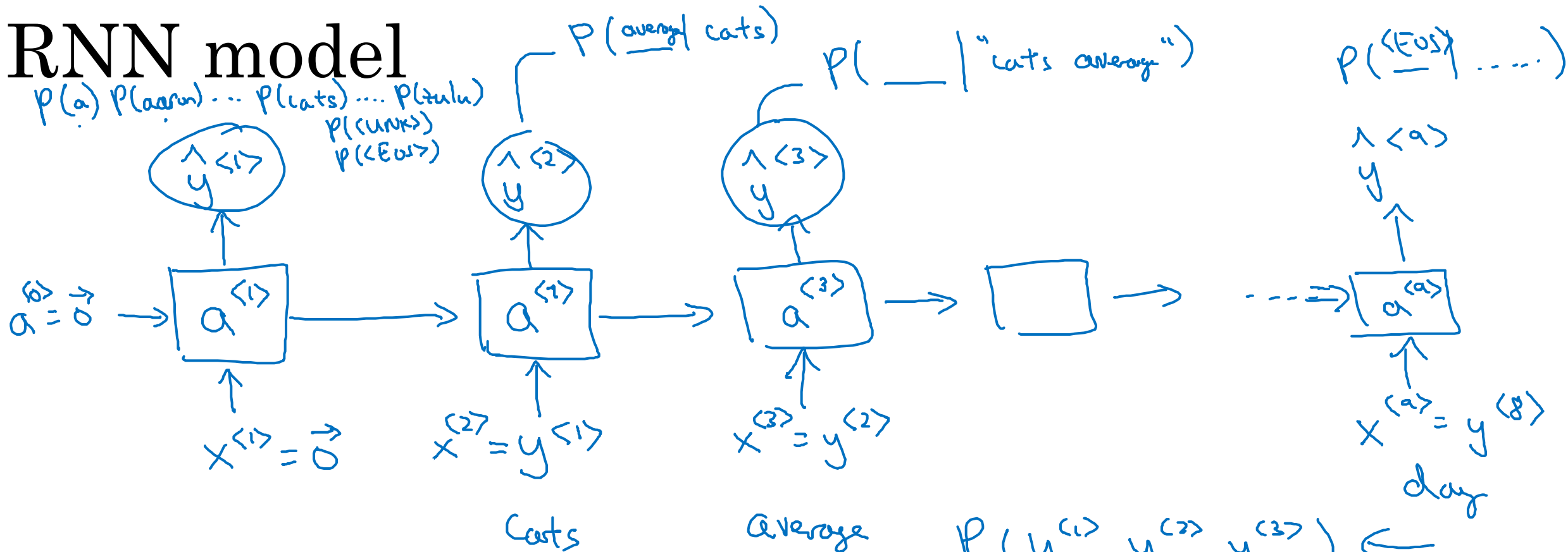Cats average 15 hours of sleep a day. $\downarrow$ <EOS>

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

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

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

<UNK>

10,000

# RNN model

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

$P(average\ |\ cats)$

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

$P(<EOS>\ |\ \dots\dots)$



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

$a^{<0>} = \vec{0}$    $a^{<1>}$    $a^{<2>}$    $a^{<3>}$    $a^{<a>}$

$x^{<1>} = \vec{0}$    $x^{<2>} = y^{<1>}$    $x^{<3>} = y^{<2>}$    $x^{<a>} = y^{<8>}$

Cats    average    day

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

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

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

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

Andrew Ng

# Recurrent Neural Networks

## Sampling novel sequences

deeplearning.ai

# Sampling a sequence from a trained RNN



$P(y^{<1>}, \ldots, y^{<T_x>})$

Training:

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

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

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

Sampling:

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

$a^{<0>} = 0 \rightarrow a^{<1>} \rightarrow a^{<2>} \rightarrow a^{<3>} \rightarrow \ldots$

$x^{<1>} = 0$    $x^{<2>} = \hat{y}^{<1>}$   The   $y^{<T_x-1>}$

$<EOS>$

$<UNK>$

$\rightarrow P(a)P(aaron)\ldots P(zulu)P(<UNK>)$    n.p. random.choice    $P(\_ \mid the)$

Andrew Ng

# Character-level language model

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

→ Vocabulary = [a, b, c, ..., z, ⌣, ., , , ;, 0, ..., 9, A, ..., Z]

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

Cat                    average
↑ ↑ ↑ ↑ ...

May



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

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

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

Andrew Ng

# 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.
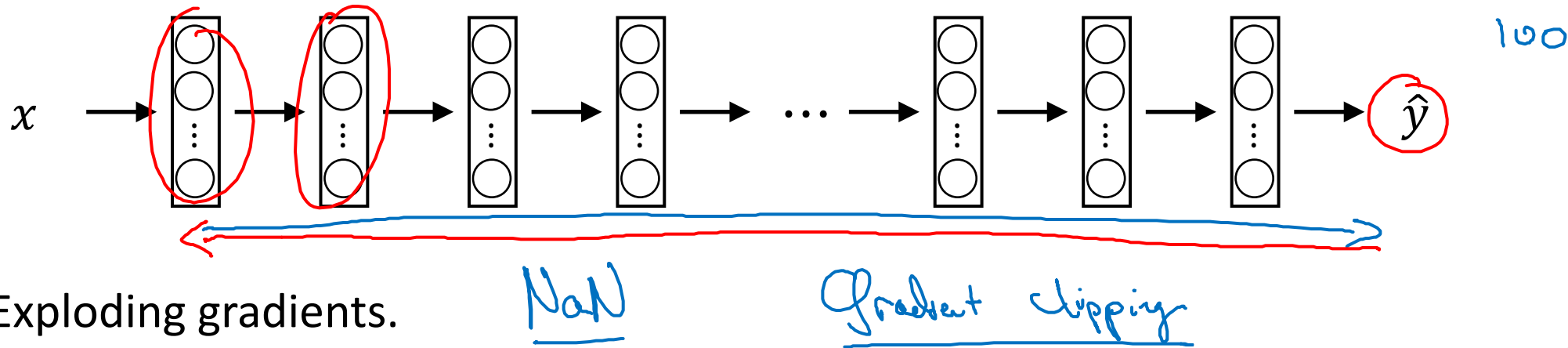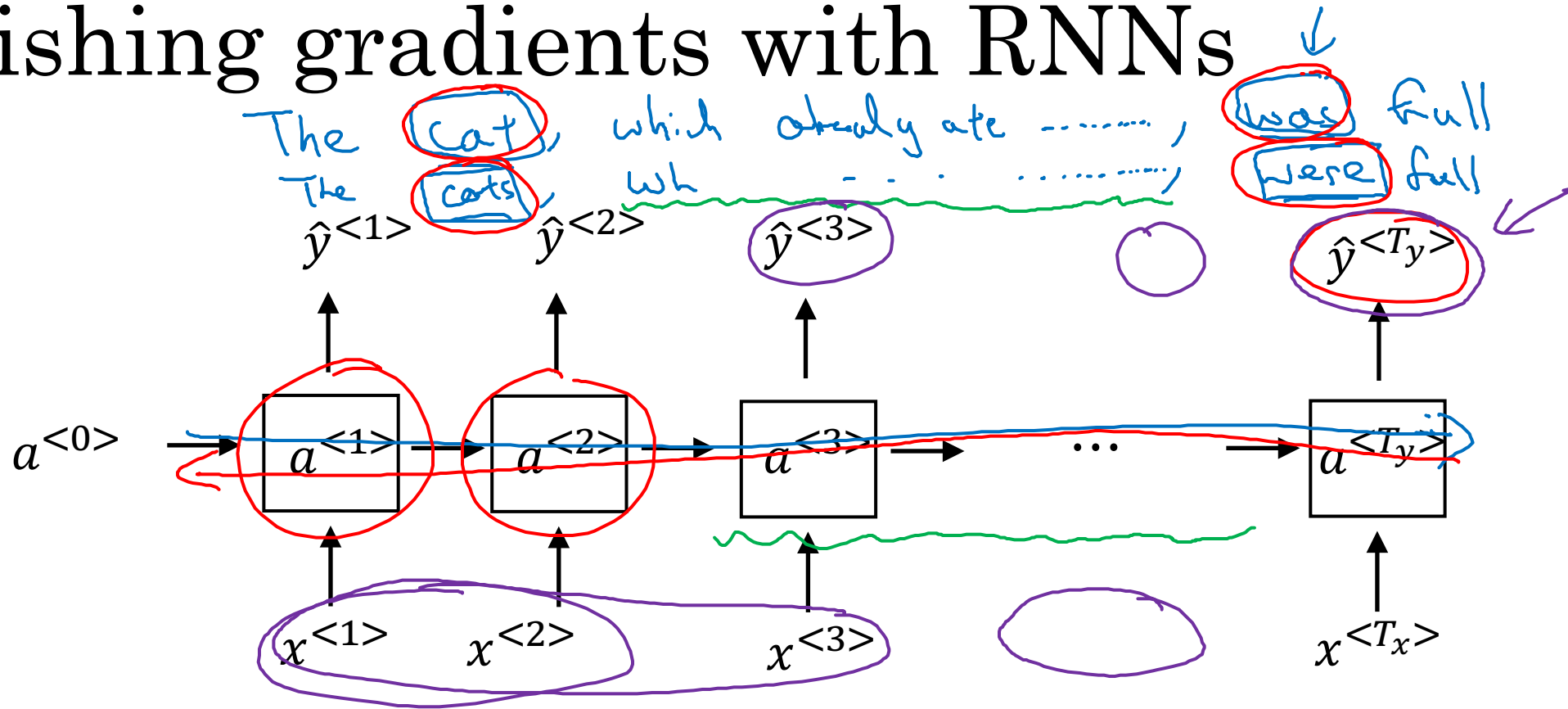
For whose are ruse of mine eyes heaves.

Andrew Ng

Recurrent Neural Networks

deeplearning.ai

Vanishing gradients with RNNs

# Vanishing gradients with RNNs

The cat, which already 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

Recurrent Neural Networks

Gated Recurrent Unit (GRU)

deeplearning.ai

# RNN unit

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

Softmax

$a^{<t-1>}$

$a$

tanh

$a^{<t>}$

$x^{<t>}$

tanh

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

# GRU (simplified)



$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"

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

$\Gamma_u = 1$

element-wise

$\Gamma_u = 0.000001$

Gate

$\Gamma_u = 1$   $\Gamma_u = 0$  $\Gamma_u = 0$  $\Gamma_u = 0$  .......  $= 1$

$c^{\langle t \rangle} = 1$

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

[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

deeplearning.ai

---

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

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

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

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

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

$\tanh(c^{<t>})$

$\Gamma_f$

[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

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

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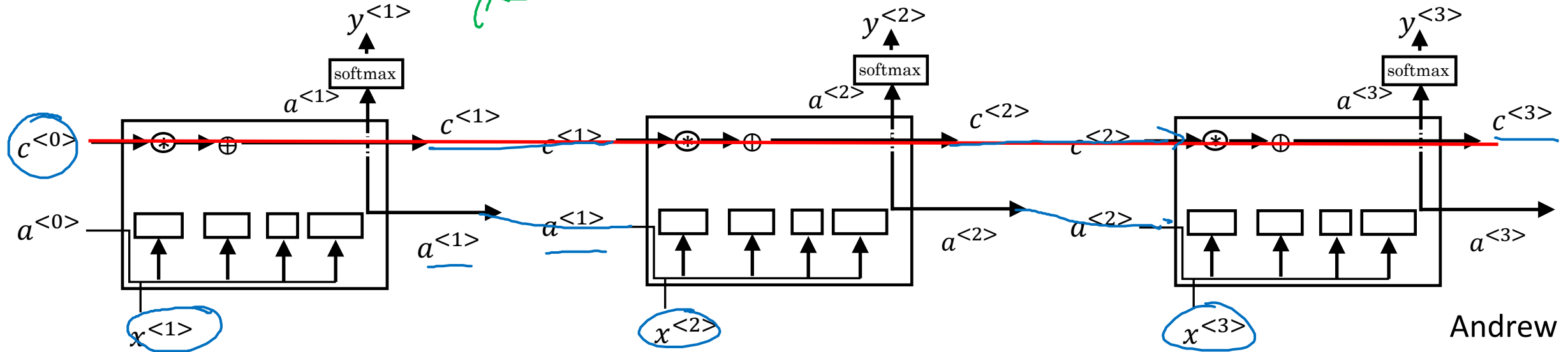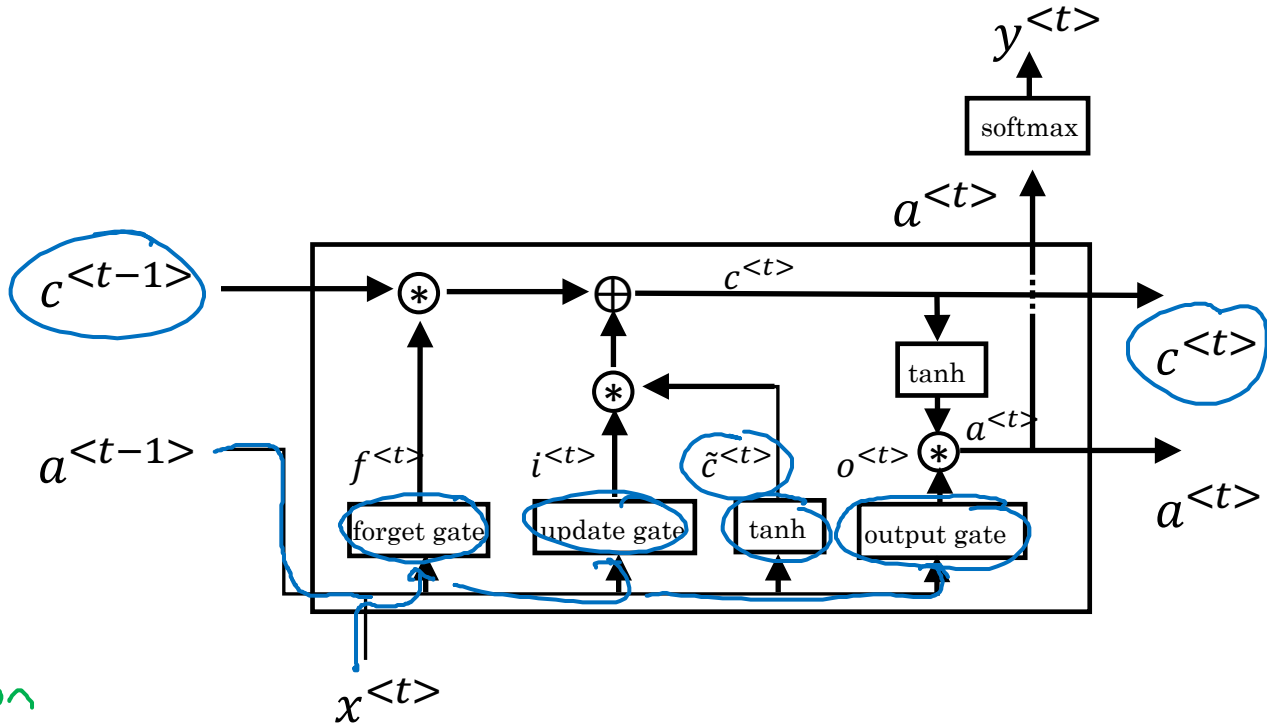
$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

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

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

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

peephole connection
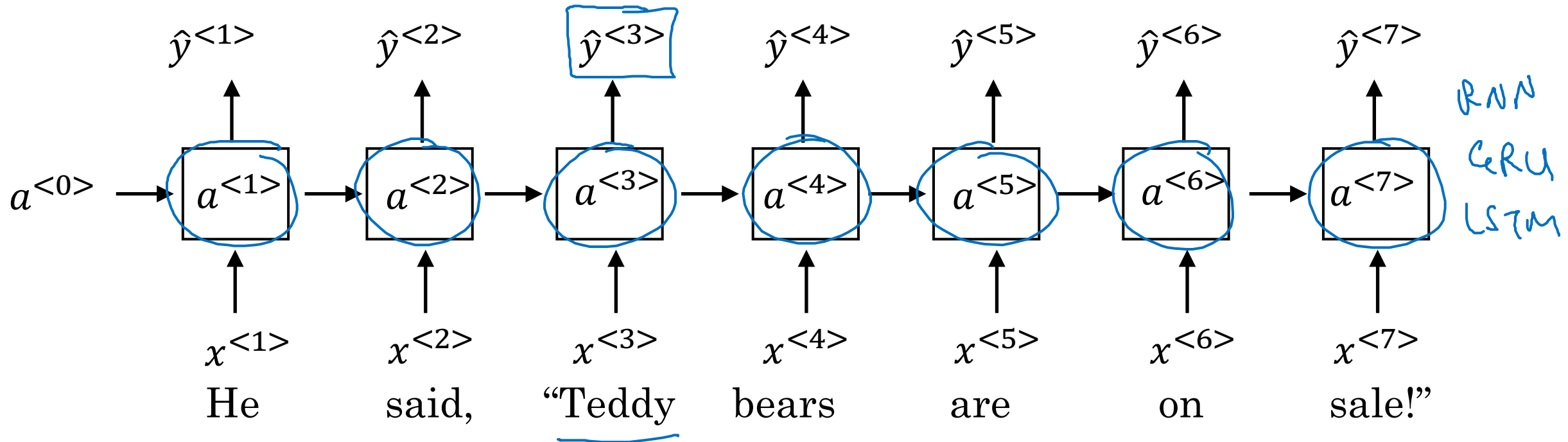
$C^{<t-1>}$



Andrew Ng

# Recurrent Neural Networks
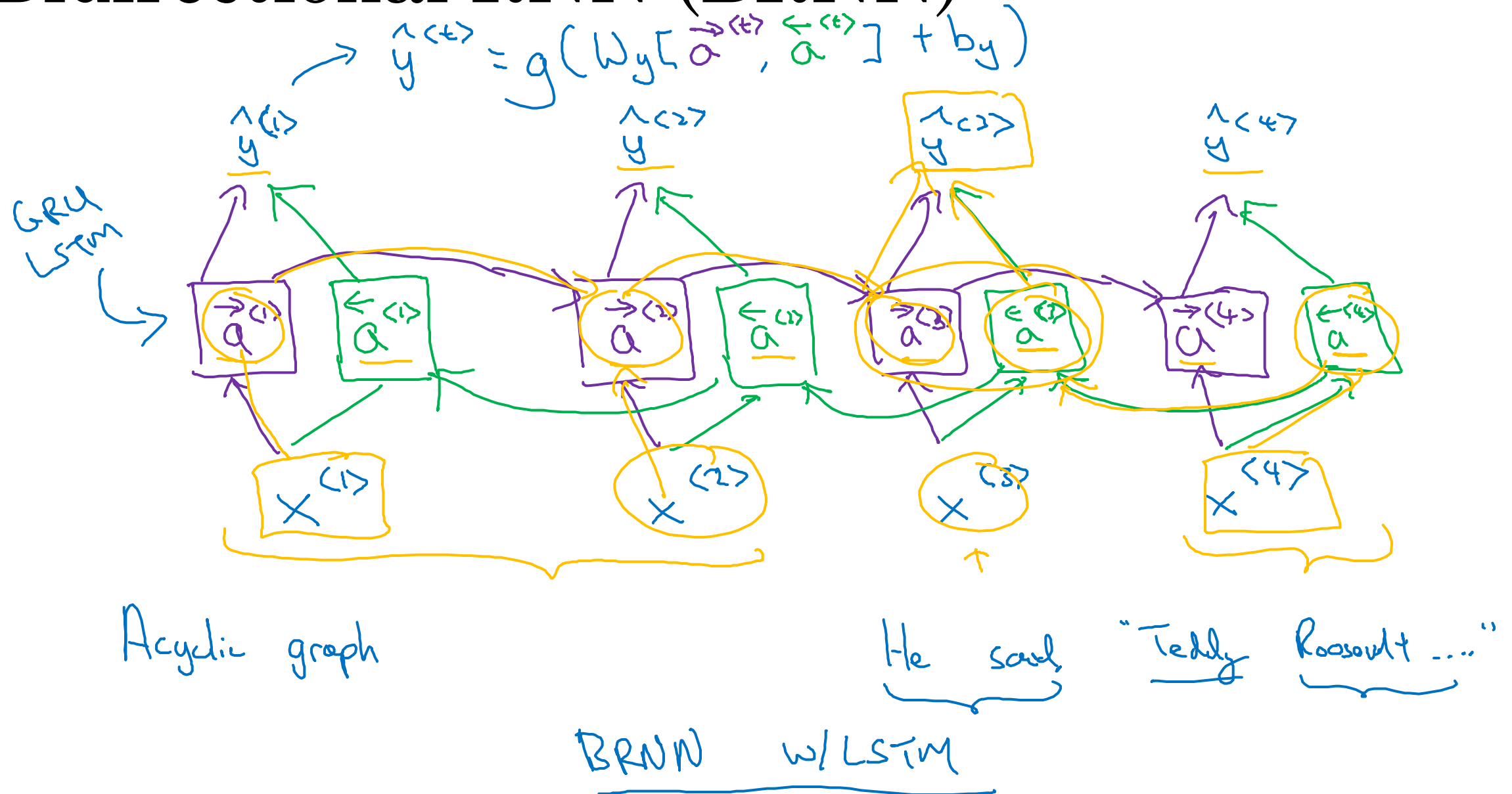
deeplearning.ai

Bidirectional RNN

# Getting information from the future

He said, "Teddy bears are on sale!"

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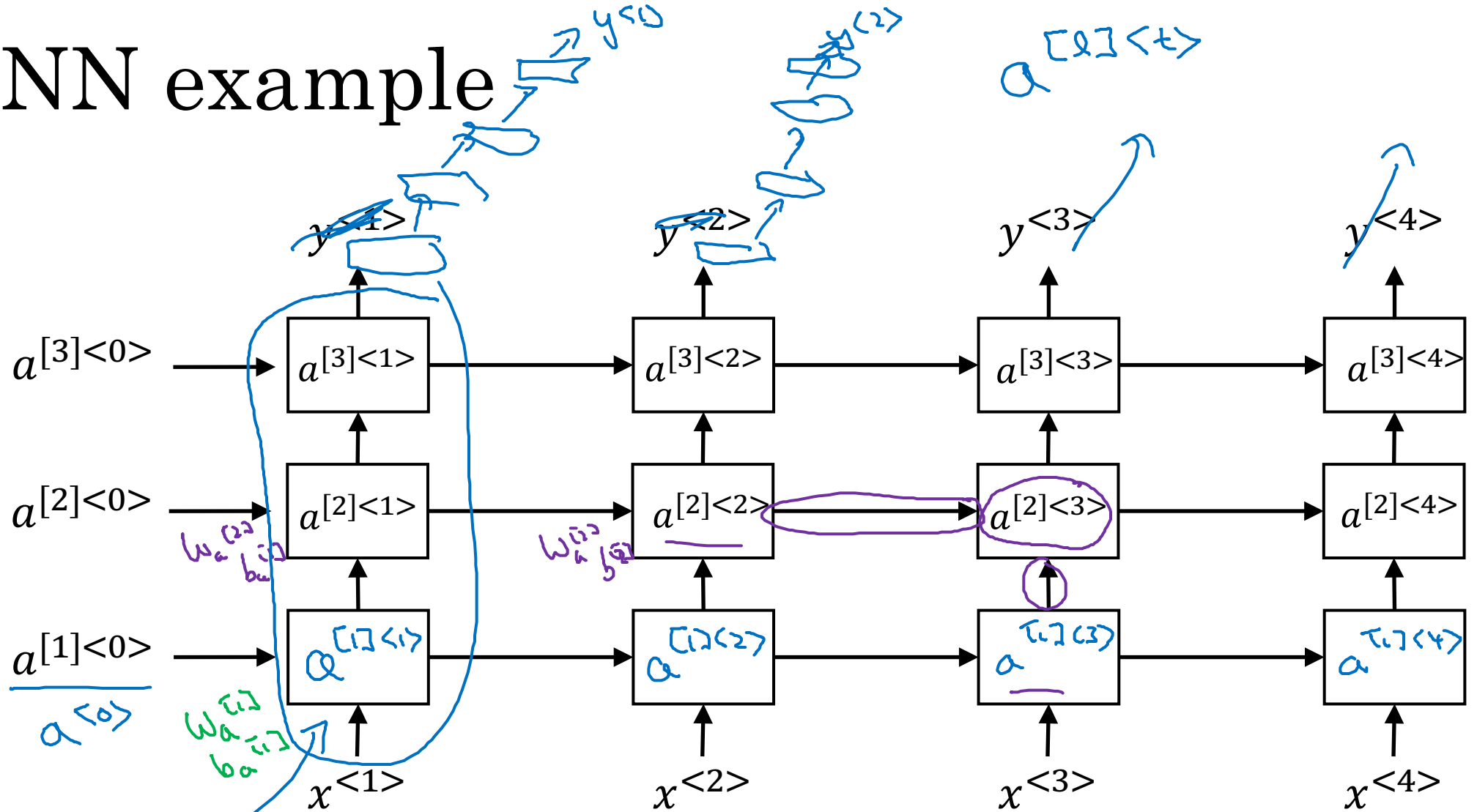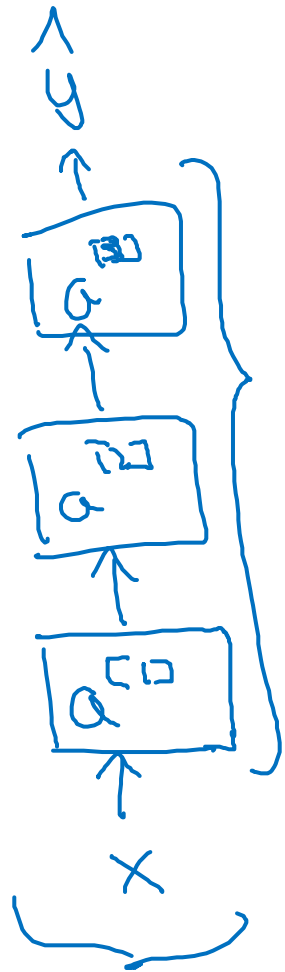


Andrew Ng

# Bidirectional RNN (BRNN)

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



GRU
LSTM

Acyclic graph

He said, "Teddy Roosevelt ..."

BRNN w/ LSTM

Andrew Ng

# Deep RNN example

$a^{[2]<t>}$

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

| | | | |
|---|---|---|---|
| $a^{[3]<0>} \rightarrow$ | $a^{[3]<1>}$ | $a^{[3]<2>}$ | $a^{[3]<3>}$ | $a^{[3]<4>}$ |

$a^{[3]<0>} \rightarrow \boxed{a^{[3]<1>}} \rightarrow \boxed{a^{[3]<2>}} \rightarrow \boxed{a^{[3]<3>}} \rightarrow \boxed{a^{[3]<4>}}$

$W_a^{[2]}, b_a^{[2]}$

$a^{[2]<0>} \rightarrow \boxed{a^{[2]<1>}} \rightarrow \boxed{a^{[2]<2>}} \rightarrow \boxed{a^{[2]<3>}} \rightarrow \boxed{a^{[2]<4>}}$

$W_a^{[1]}, b_a^{[1]}$

$a^{[1]<0>} \rightarrow \boxed{a^{[1]<1>}} \rightarrow \boxed{a^{[1]<2>}} \rightarrow \boxed{a^{[1]<3>}} \rightarrow \boxed{a^{[1]<4>}}$

$\underline{a^{[0]}}$

$W_a^{[1]}, b_a^{[1]}$

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

$\hat{y}$

$\boxed{a^{[3]}}$

$\boxed{a^{[2]}}$

$\boxed{a^{[1]}}$

$X$

RNN
GRU
LSTM     BRNN

$$a^{[2]<3>} = g\left(W_a^{[2]}\left[a^{[2]<2>}, a^{[1]<3>}\right] + b_a^{[2]}\right)$$

Andrew Ng