



deeplearning.ai

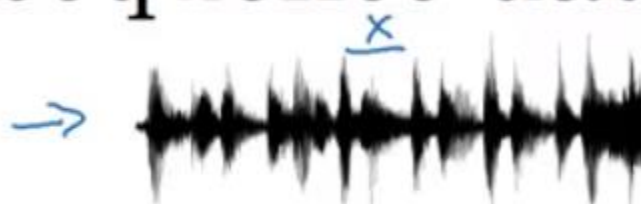
# Recurrent Neural Networks

---

## Why sequence models?

# Examples of sequence data

Speech recognition



“The quick <sup>y</sup> brown fox jumped  
over the lazy dog.”

Music generation



Sentiment classification

“There is nothing to like  
in this movie.”



DNA sequence analysis → AGCCCCTGTGAGGAACTAG



AG**CCCCTGTGAGGAACT**AG

Machine translation

Voulez-vous chanter avec  
moi?



Do you want to sing with  
me?

Video activity recognition



Running

Name entity recognition → Yesterday, Harry Potter  
met Hermione Granger.



Yesterday, **Harry Potter**  
met **Hermione Granger**.

Andrew Ng



deeplearning.ai

# Recurrent Neural Networks

---

## Notation

# Motivating example

NLP

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

$\rightarrow x^{(1)} \quad x^{(2)} \quad x^{(3)} \quad \dots \quad x^{(t)} \quad \dots \quad x^{(9)}$

$$T_x = 9$$

$\rightarrow y:$

$y^{(1)} \quad y^{(2)} \quad y^{(3)} \quad \dots \quad y^{(9)}$

$$T_y = 9$$

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

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

15

$y^{(i)(t)}$   
 $\uparrow$

$$T_y^{(i)}$$

# Representing words

$$x^{(t)} \rightarrow y^{(t)}$$
$$(x, y)$$

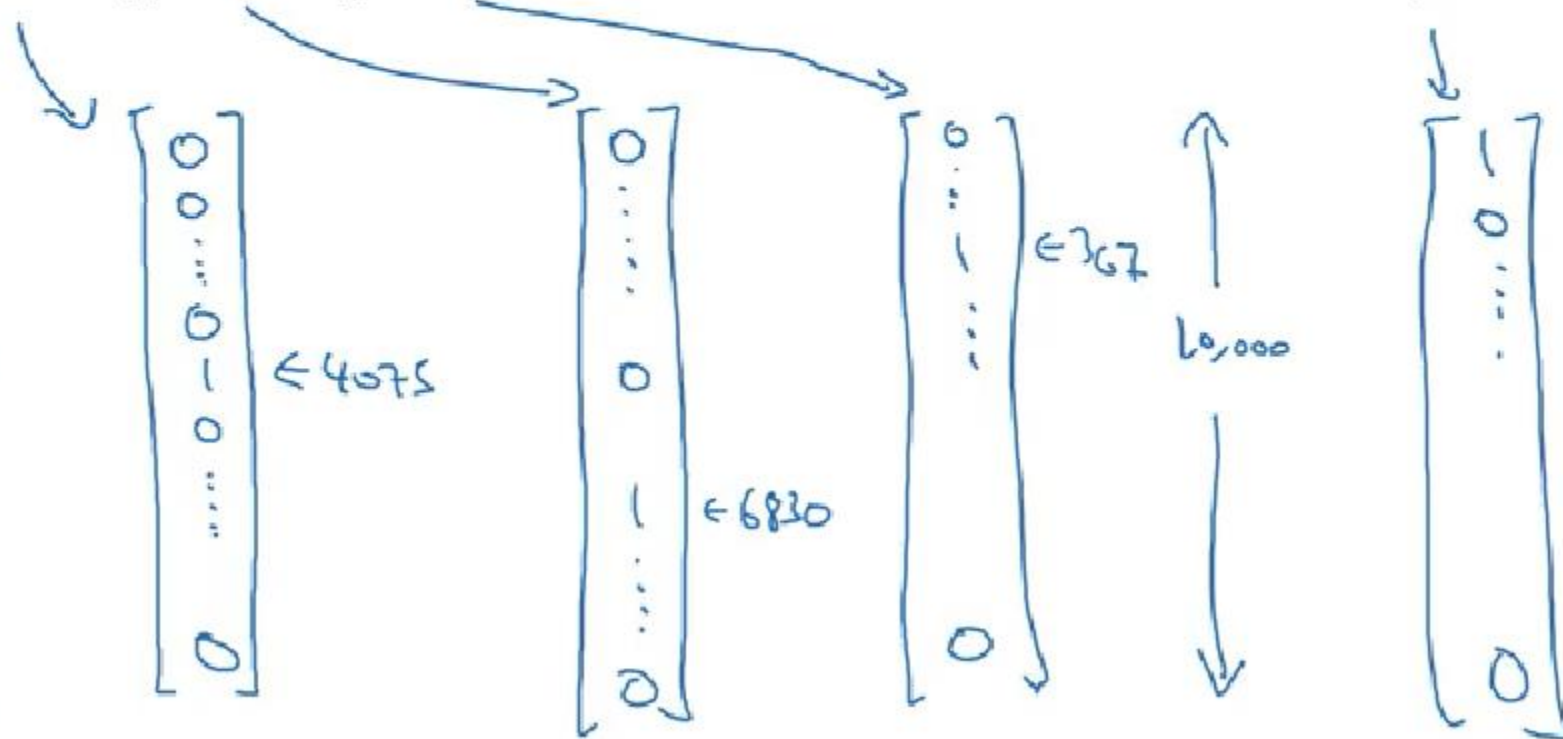
x: Harry Potter and Hermione Granger invented a new spell.

$x^{(1)}$   $x^{(2)}$   $x^{(3)}$  ...  $x^{(9)}$

Vocabulary

a	1
aaron	2
...	...
and	367
...	...
harry	4075
potter	6830
...	...
zulu	10,000

<UNK> 10,000



One-hot



deeplearning.ai

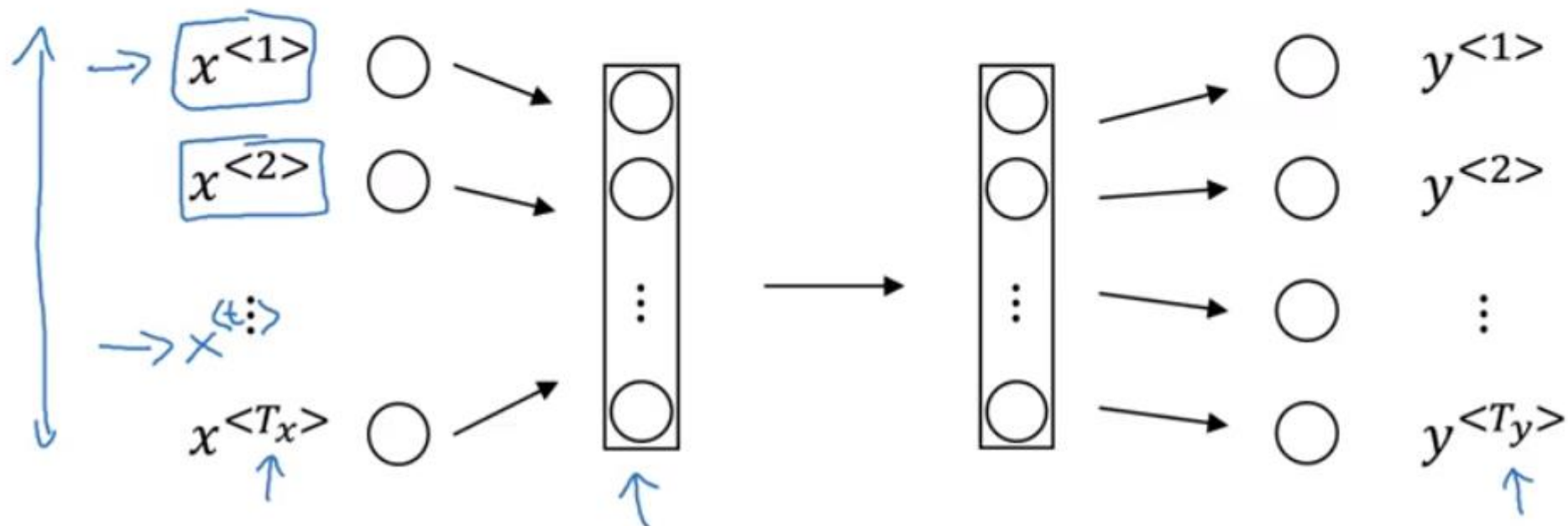
# Recurrent Neural Networks

---

## Recurrent Neural Network Model



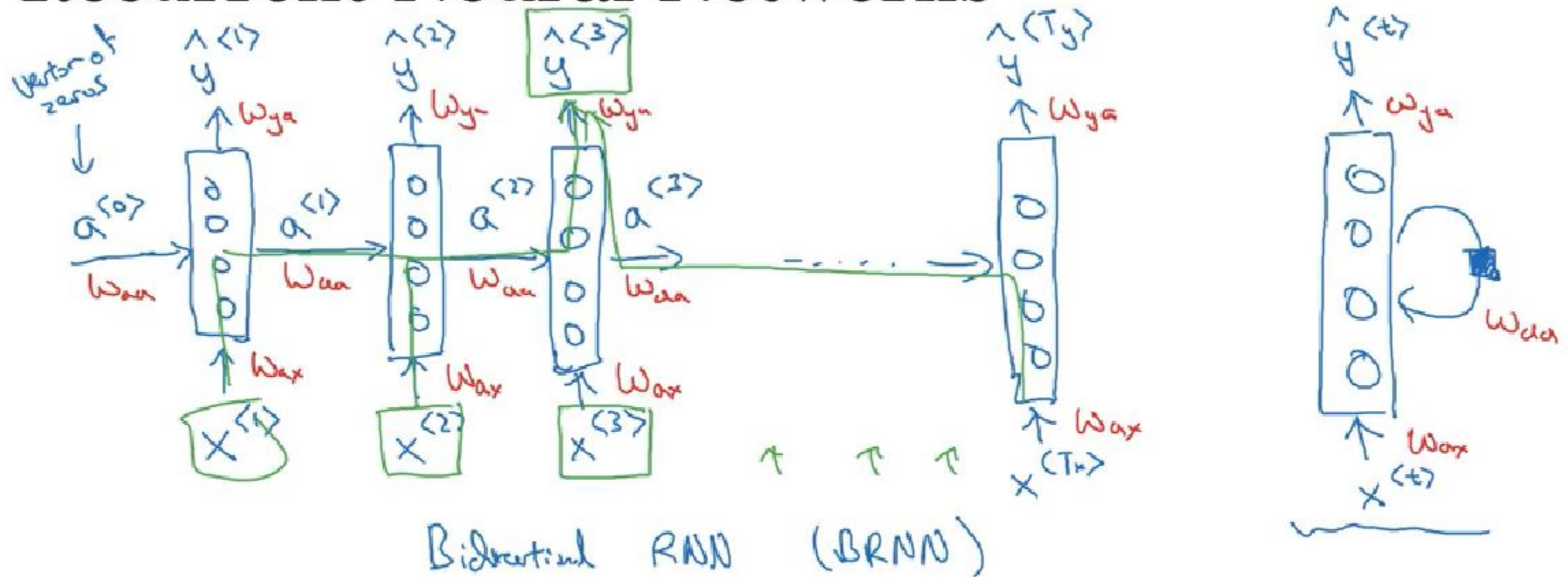
# Why not a standard network?



## Problems:

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

# Recurrent Neural Networks

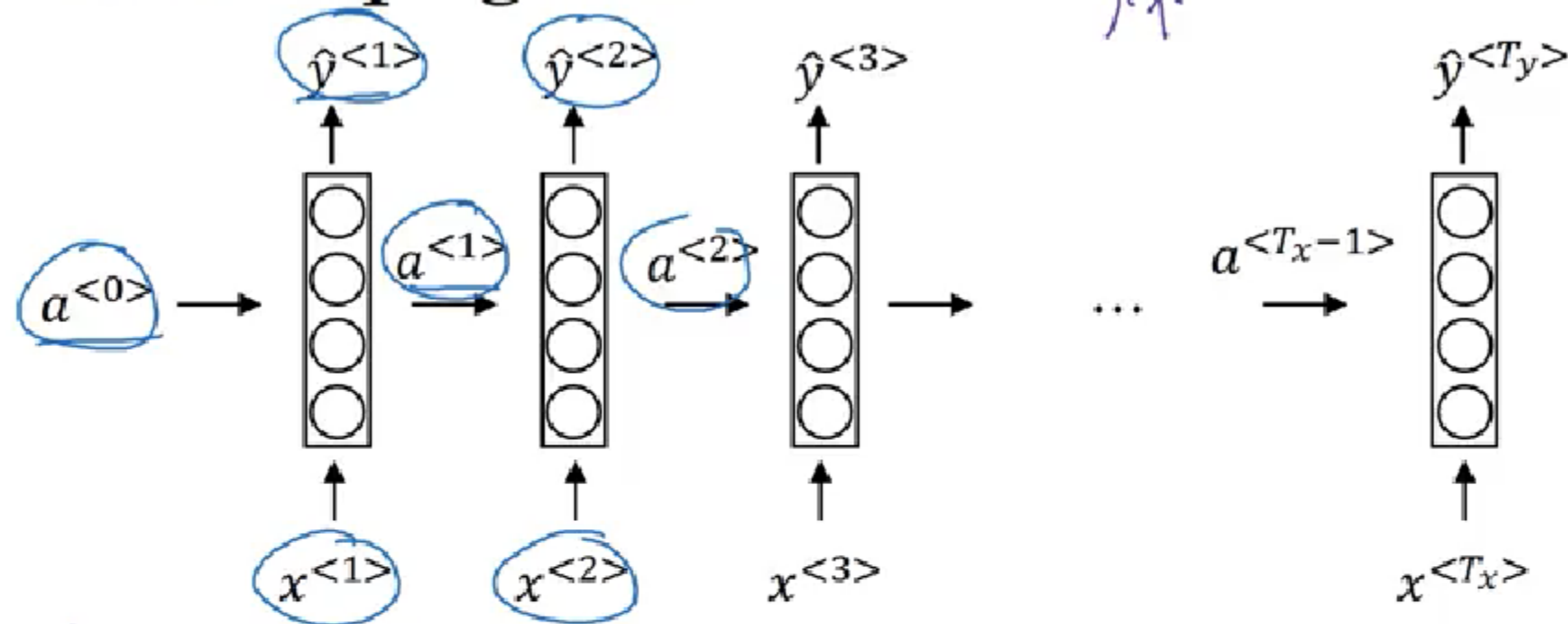


He said, "Teddy Roosevelt was a great President."

He said, "Teddy bears are on sale!"



# Forward Propagation



$$a^{(0)} = \vec{0}$$

$$\underline{a}^{(t)} = g_1(W_{aa} \underline{a}^{(t-1)} + \underline{W_{ax}} x^{(t)} + b_a) \leftarrow \tanh / \text{Relu}$$

$$\underline{\hat{y}}^{(t)} = g_2(\underline{W_{ya}} \underline{a}^{(t)} + b_y) \leftarrow \text{sigmoid}$$

$$\boxed{\begin{aligned} a^{(t)} &= g(W_{aa} a^{(t-1)} + W_{ax} x^{(t)} + b_a) \\ \hat{y}^{(t)} &= g(W_{ya} a^{(t)} + b_y) \end{aligned}}$$

# Simplified RNN notation

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

Diagram illustrating the dimensions of the weights and inputs in the first equation:

- $W_{aa}$  is a  $100 \times 100$  matrix (indicated by  $(100, 100)$  below it).
- $a^{<t-1>}$  is a  $100$ -dimensional vector (indicated by  $100$  below it).
- $W_{ax}$  is a  $100 \times 10,000$  matrix (indicated by  $(100, 10,000)$  below it).
- $x^{<t>}$  is a  $10,000$ -dimensional vector (indicated by  $10,000$  below it).

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

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

Diagram illustrating the dimensions of the weights and inputs in the second equation:

- $W_y$  is a  $1 \times 100$  weight vector (indicated by  $1$  and  $100$  below it).
- $a^{<t>}$  is a  $100$ -dimensional vector (indicated by  $100$  below it).
- $b_y$  is a  $1$ -dimensional bias vector (indicated by  $1$  below it).

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

Diagram illustrating the dimensions of the weights and inputs in the simplified notation:

- $W_a$  is a  $100 \times 10100$  matrix (indicated by  $(100, 10100)$  below it).
- $[a^{<t-1>}, x^{<t>}]$  is a  $10100$ -dimensional vector (indicated by  $10100$  below it).

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

Diagram illustrating the dimensions of the weights in the simplified notation:

- $W_{aa}$  is a  $100 \times 100$  matrix (indicated by  $100$  and  $100$  below it).
- $W_{ax}$  is a  $100 \times 10,000$  matrix (indicated by  $100$  and  $10,000$  below it).
- $W_a$  is a  $100 \times 10100$  matrix (indicated by  $(100, 10100)$  below it).

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

Diagram illustrating the dimensions of the inputs in the simplified notation:

- $a^{<t-1>}$  is a  $100$ -dimensional vector (indicated by  $100$  below it).
- $x^{<t>}$  is a  $10,000$ -dimensional vector (indicated by  $10,000$  below it).
- The concatenated vector  $[a^{<t-1>}, x^{<t>}]$  is a  $10100$ -dimensional vector (indicated by  $10100$  below it).

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



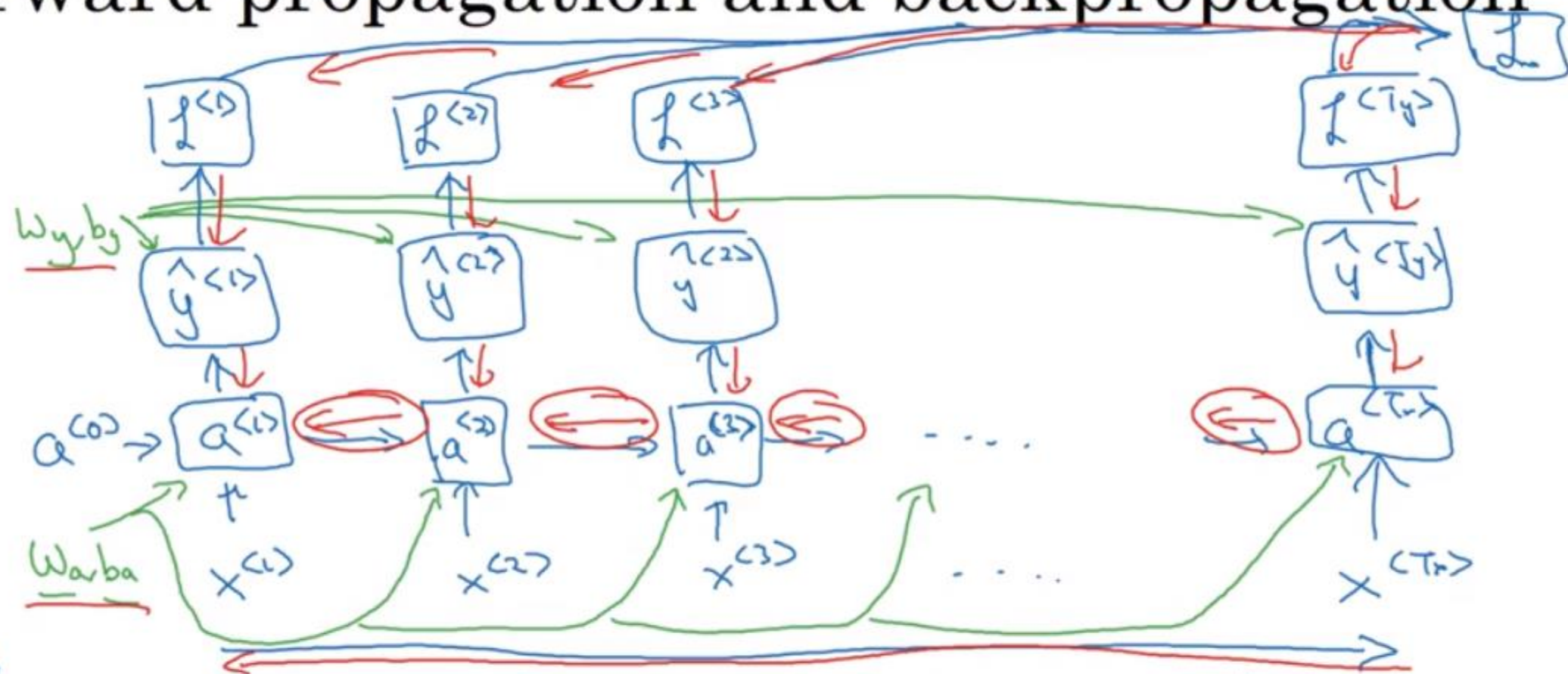
deeplearning.ai

# Recurrent Neural Networks

---

## Backpropagation through time

# 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



deeplearning.ai

# Recurrent Neural Networks

---

## Different types of RNNs



# Examples of sequence data

Speech recognition



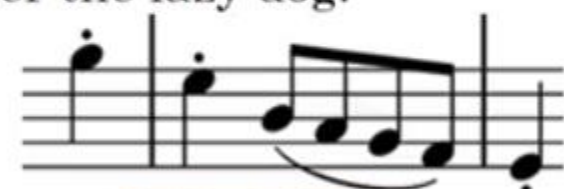
$T_x$

$T_y$

$y$

“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

AG**CCCCTGTGAGGAACTAG**

Machine translation

Voulez-vous chanter avec  
moi?

Do you want to sing with  
me?

Video activity recognition



Running

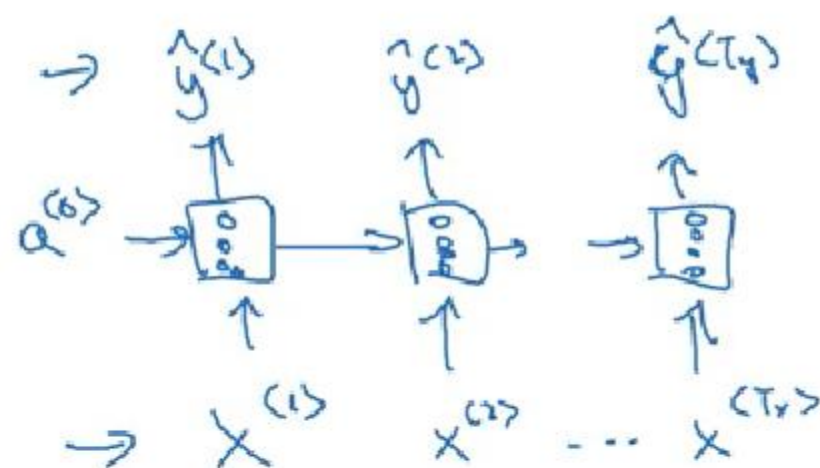
Name entity recognition

Yesterday, Harry Potter  
met Hermione Granger.

Yesterday, **Harry Potter**  
met **Hermione Granger**.

# Examples of RNN architectures

$$T_x = T_y$$

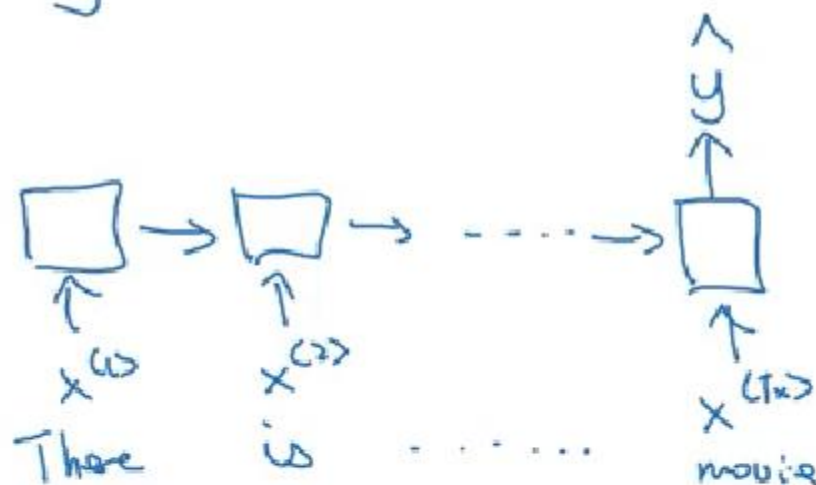


Many-to-many

Sentiment classification

$x = \text{text}$

$y = 0/1 \quad 1 \dots 5$

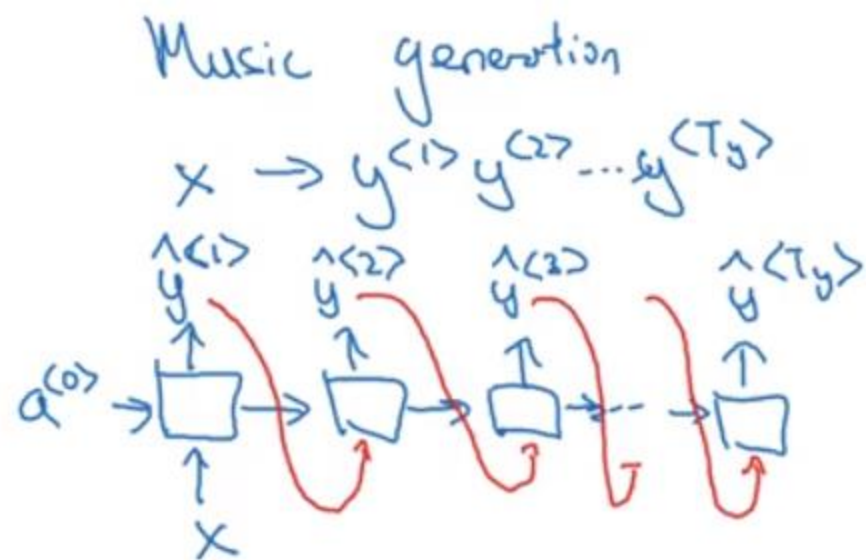


Many-to-one



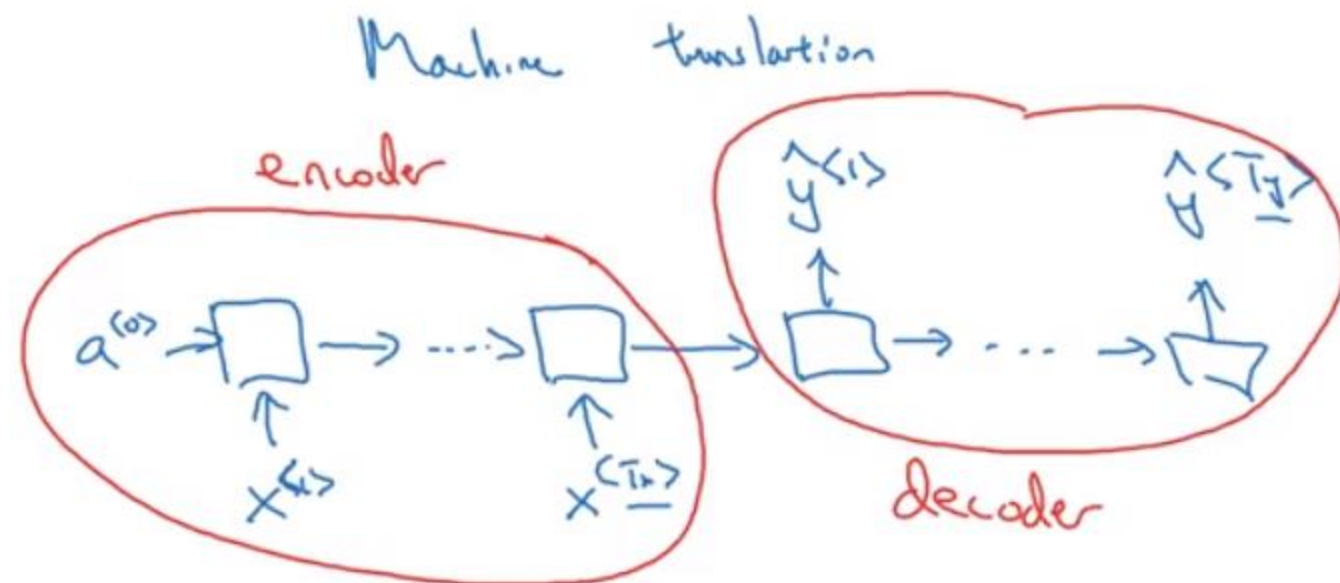
One-to-one

# Examples of RNN architectures



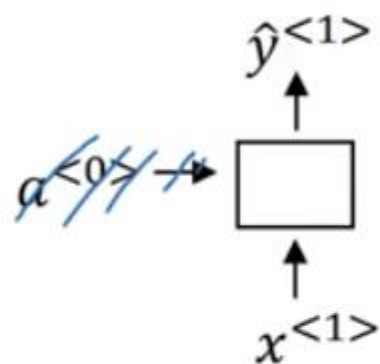
One-to-many

$$x = \phi$$

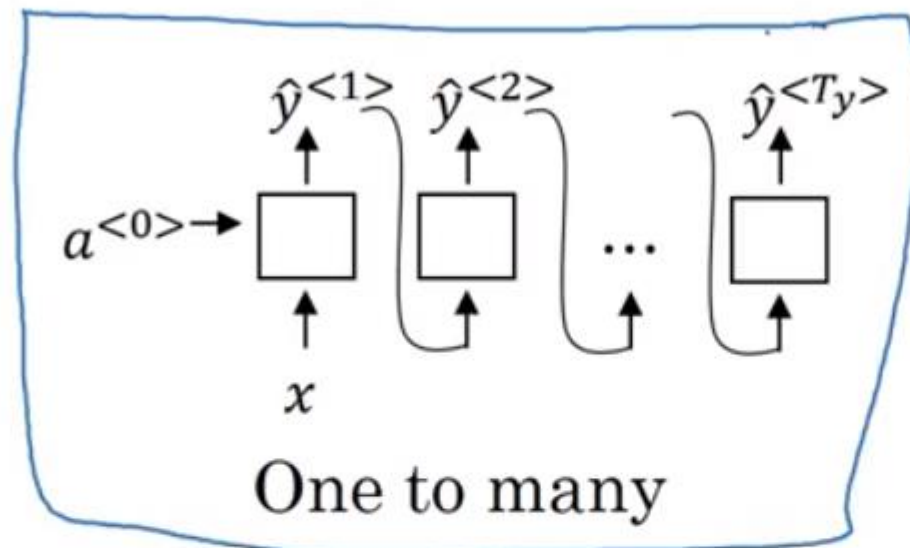


Many-to-many

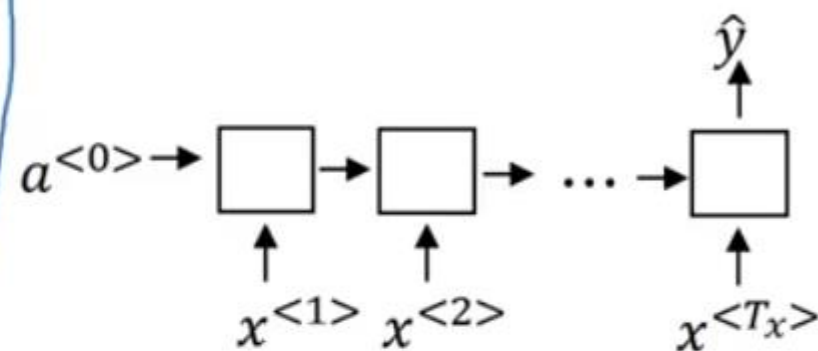
# Summary of RNN types



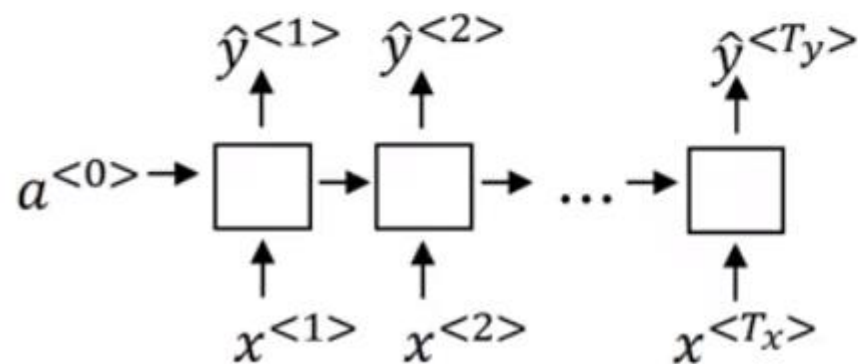
One to one



One to many

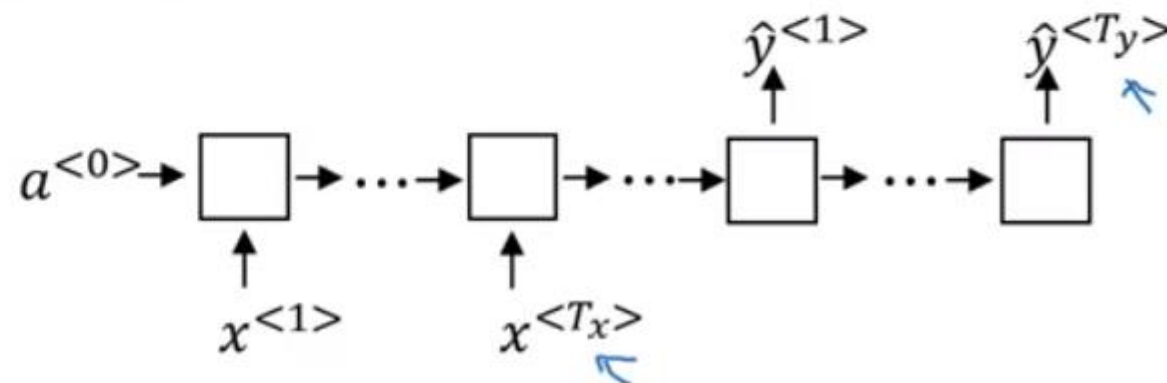


Many to one



Many to many

$T_x = T_y$



Many to many





deeplearning.ai

# Recurrent Neural Networks

---

Language model and  
sequence generation



# What is language modelling?

## Speech recognition

The apple and pair salad.

→ The apple and pear salad.

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

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

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

$$P(y^{(1)}, y^{(2)}, \dots, y^{(T)})$$

# Language modelling with an RNN

Training set: large corpus of english text.

Tokenize

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

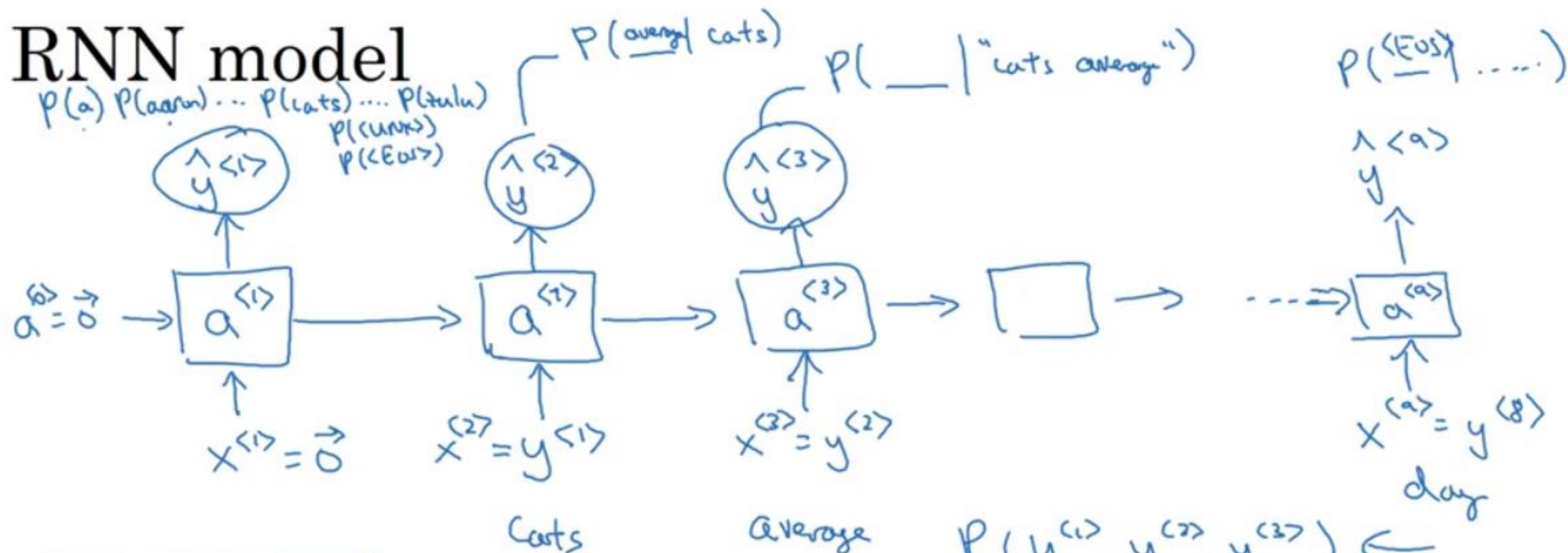
$y^{(1)}$     $y^{(2)}$     $y^{(3)}$    ...    $y^{(8)}$     $y^{(9)}$   
 $x^{(t)} = y^{(t-1)}$

The Egyptian ~~Mau~~ is a breed of cat.  $\langle \text{EOS} \rangle$

10,000

$\langle \text{UNK} \rangle$

# RNN model



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

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

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

$$P(y^{(1)}, y^{(2)}, y^{(3)}) \leftarrow$$

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

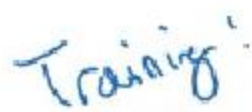


deeplearning.ai

# Recurrent Neural Networks

---

Sampling novel  
sequences

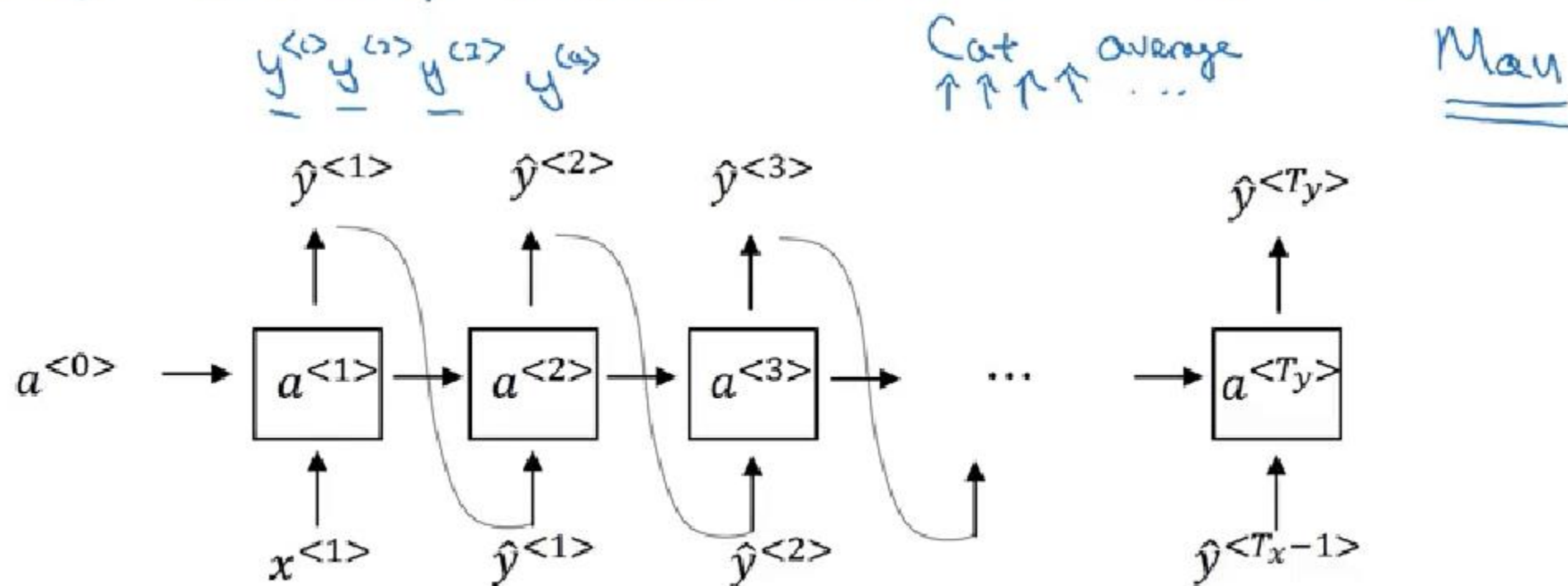
$$P(y^{(1)}, \dots, y^{(T_x)})$$




# Character-level language model

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

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



# Sequence generation

## News

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.

## Shakespeare

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.



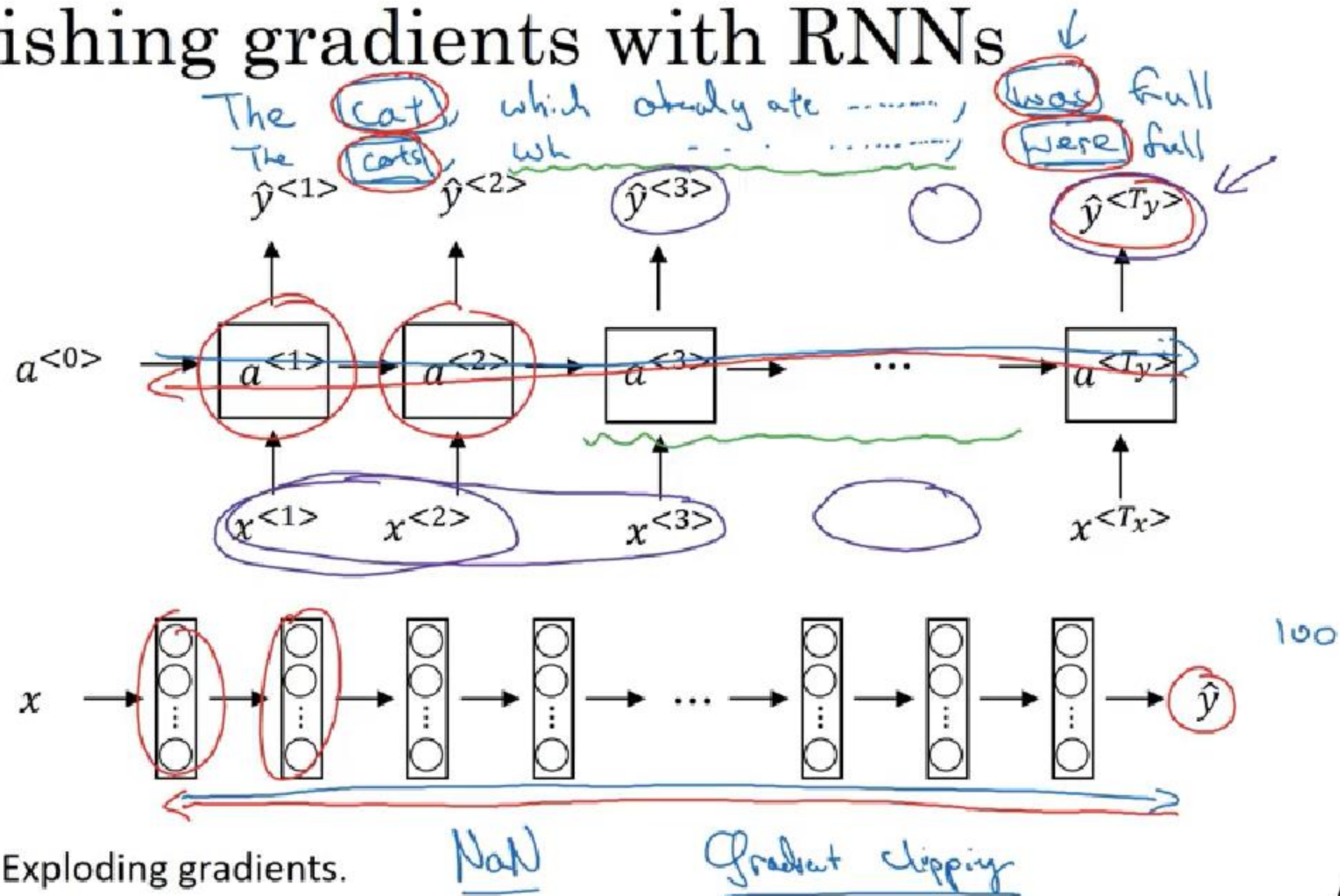
deeplearning.ai

# Recurrent Neural Networks

---

## Vanishing gradients with RNNs

# Vanishing gradients with RNNs





deeplearning.ai

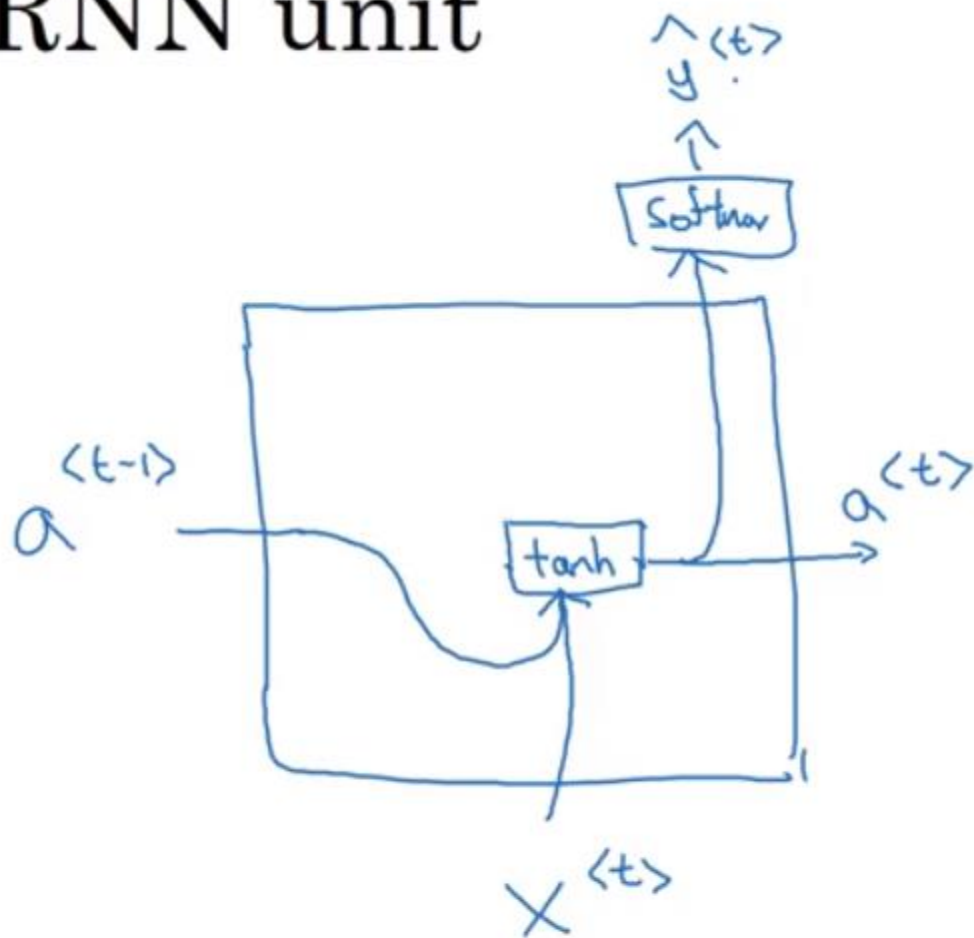
# Recurrent Neural Networks

---

## Gated Recurrent Unit (GRU)

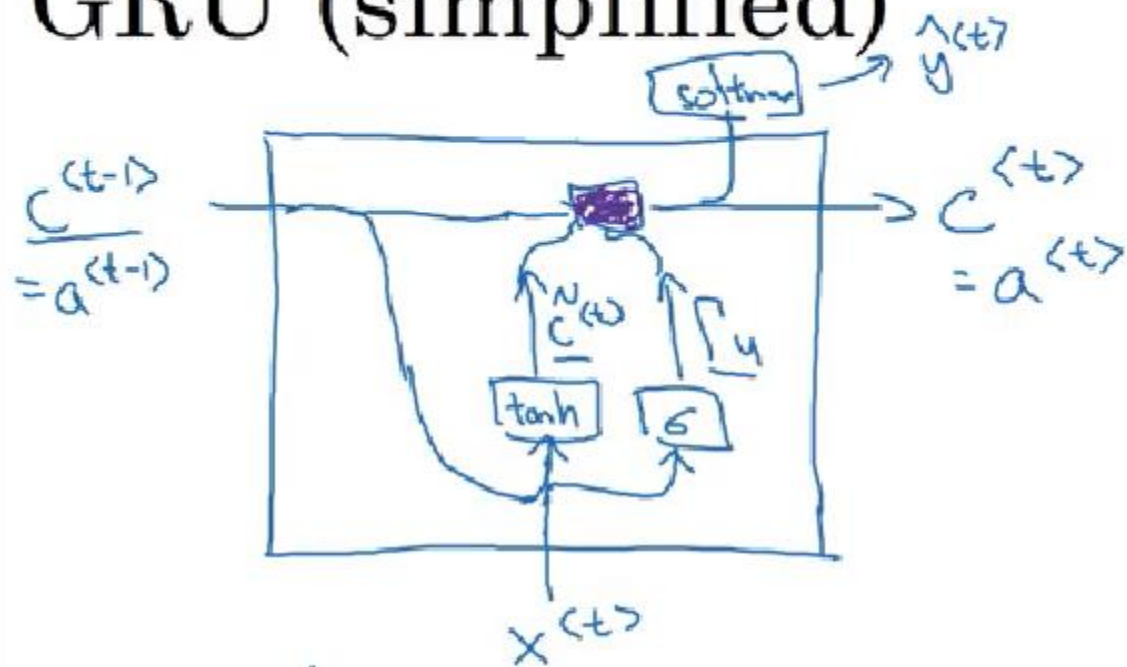


# RNN unit



$$\underline{a^{<t>}} = \overset{\text{tanh}}{\underset{\uparrow}{g}}(\underset{\uparrow}{W_a[a^{<t-1>}, x^{<t>}]} + \underline{b_a})$$

# GRU (simplified)



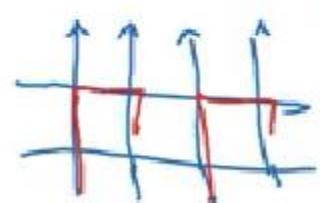
$C$  = memory cell

$$\rightarrow \boxed{C^{(t)}} = \underline{a}^{(t)}$$

$$\rightarrow \boxed{\tilde{C}^{(t)}} = \tanh(W_c [c^{(t-1)}, x^{(t)}] + b_c)$$

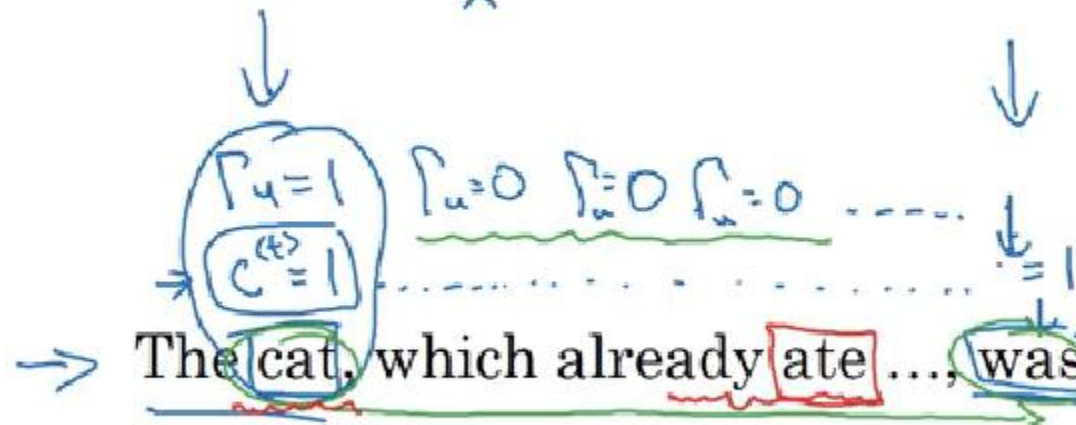
$$\rightarrow \boxed{\Gamma_u^{(t)}} = \sigma(W_u [c^{(t-1)}, x^{(t)}] + b_u)$$

$$\left\{ \boxed{C^{(t)}} = \underbrace{\Gamma_u^{(t)}}_{\leftarrow \text{"update"}} * \tilde{C}^{(t)} + (1 - \Gamma_u^{(t)}) * \boxed{C^{(t-1)}}$$



element-wise  
Gate

$$\Gamma_u = 0.000001$$



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

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

$$\begin{cases} \Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u) \\ \Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r) \end{cases}$$

LSTM

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

The cat, which ate already, was full.



deeplearning.ai

# Recurrent Neural Networks

---

LSTM (long short  
term memory) unit

# GRU and LSTM

## GRU

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

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

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

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

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

$\uparrow$   
 $\Gamma_f$

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

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

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



# LSTM in pictures

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

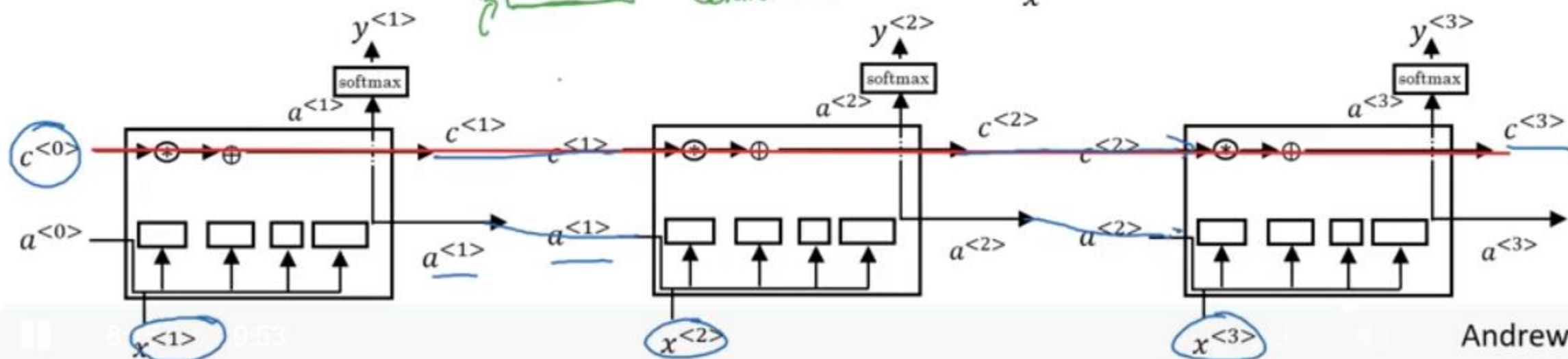
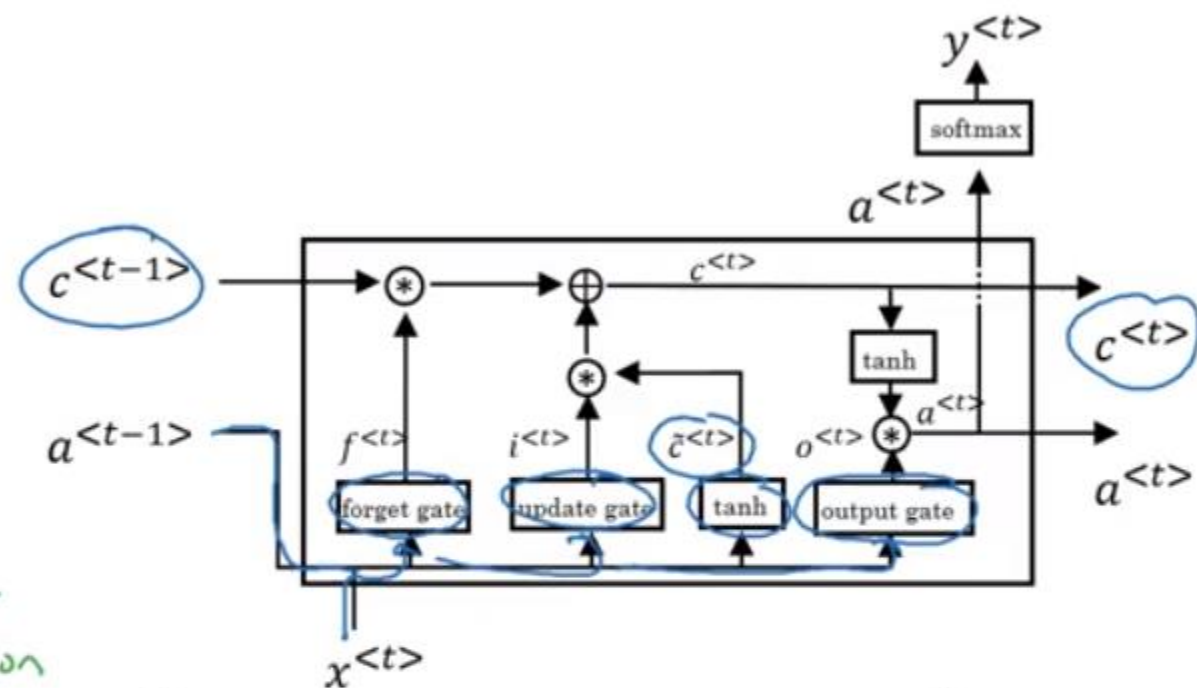
$$\Gamma_u = \sigma(W_u [a^{<t-1>}, x^{<t>}] + b_u)$$

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

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

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

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





deeplearning.ai

# Recurrent Neural Networks

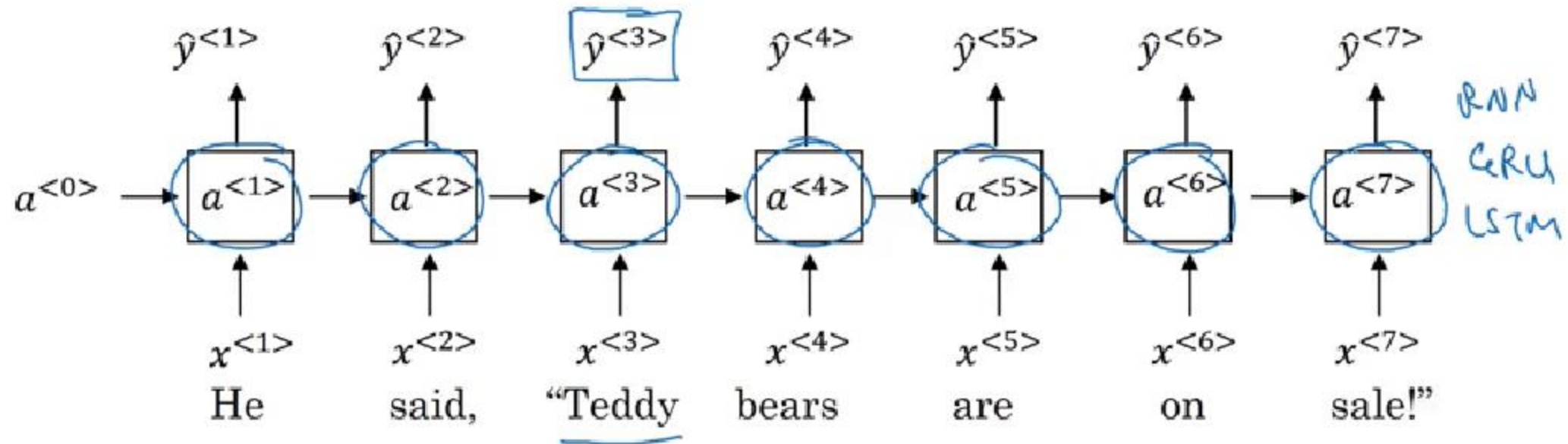
---

## Bidirectional RNN

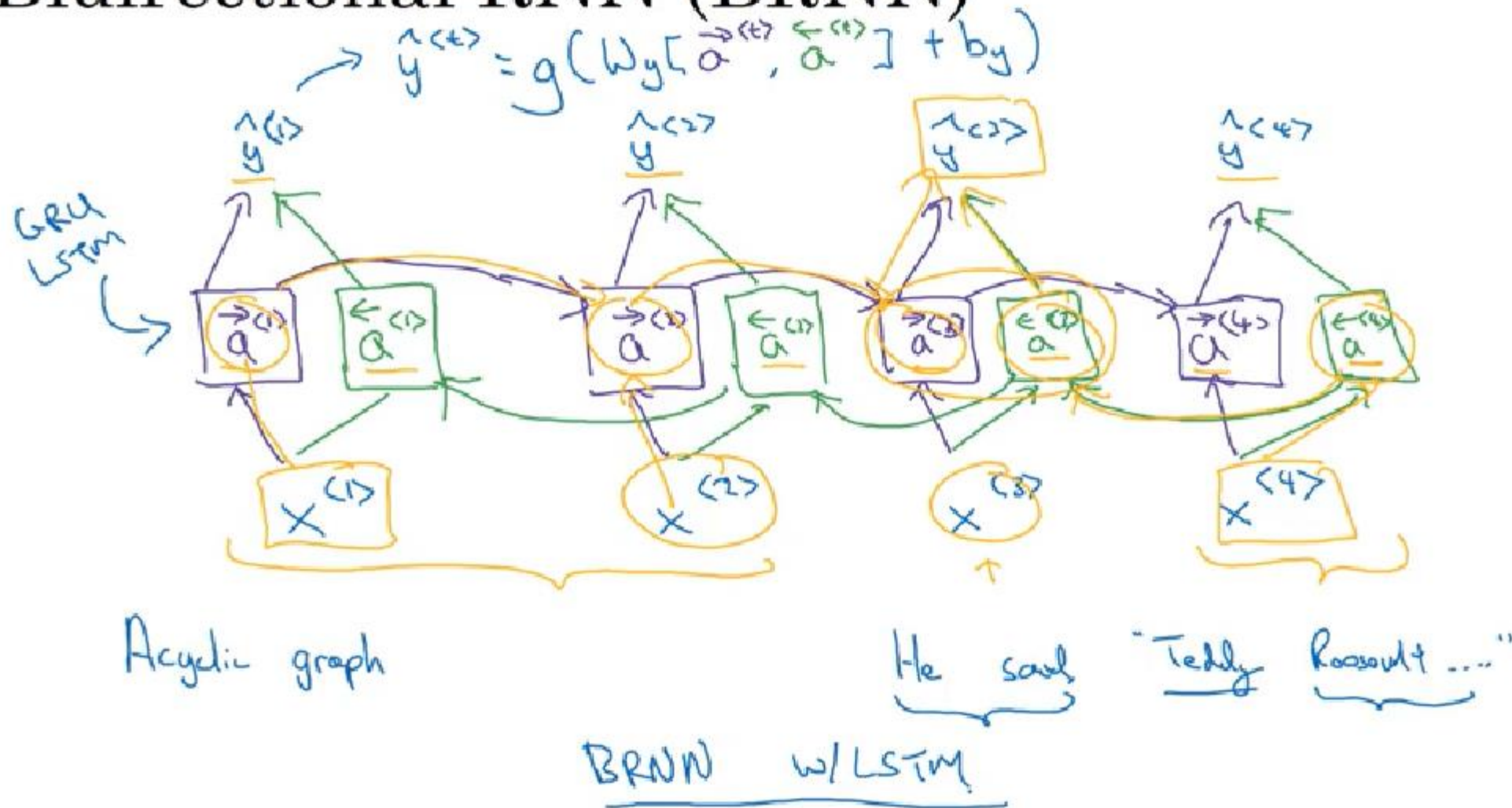
# Getting information from the future

He said, "Teddy bears are on sale!"

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



# Bidirectional RNN (BRNN)





deeplearning.ai

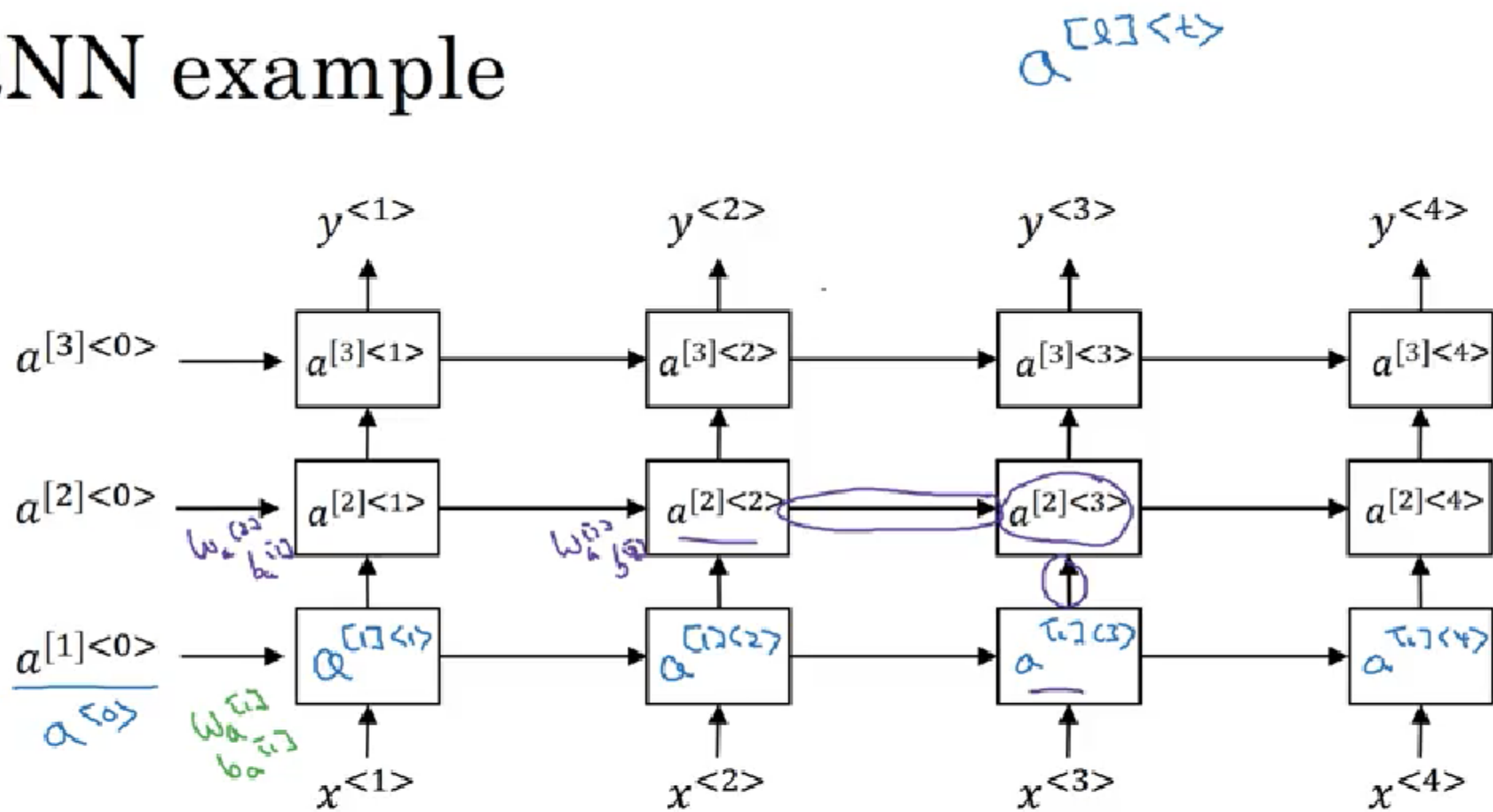
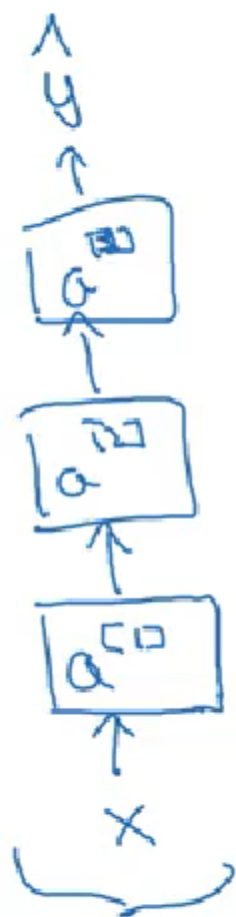
# Recurrent Neural Networks

---

## Deep RNNs



# Deep RNN example



$$a^{[2]<3>} = g(w_a^{[2]} [a^{[1]<2>}, a^{[1]<3>}] + b_a^{[2]})$$

# Deep RNN example

