



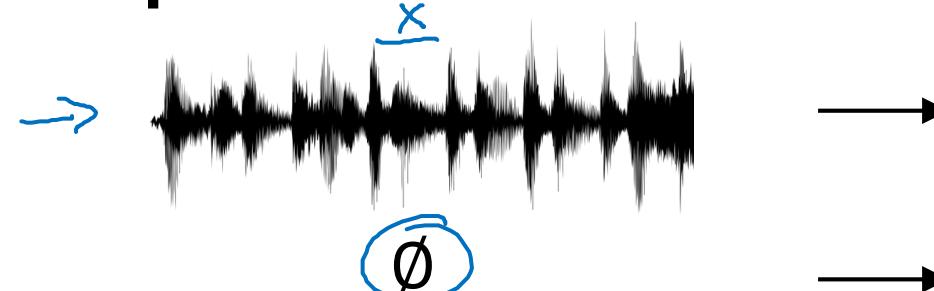
deeplearning.ai

Recurrent Neural Networks

Why sequence
models?

Examples of sequence data

Speech recognition



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

→ AGCCCCTGTGAGGAAC TAG



AGCCCCTGTGAGGAAC **TAG**

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.

$\rightarrow \underline{x^{<1>}}$ $x^{<2>}$ $x^{<3>}$ - - - - - $x^{<t>}$ - - - - - $x^{<9>}$

$T_x = 9$

$\rightarrow y:$

1 1 0 1 1 0 0 0 0
 $y^{<1>}$ $y^{<2>}$ $y^{<3>}$ - - - - - $y^{<9>}$

$T_y = 9$

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

$y^{(i)<t>}$

$T_x^{(i)} = 9$

$T_y^{(i)}$

15

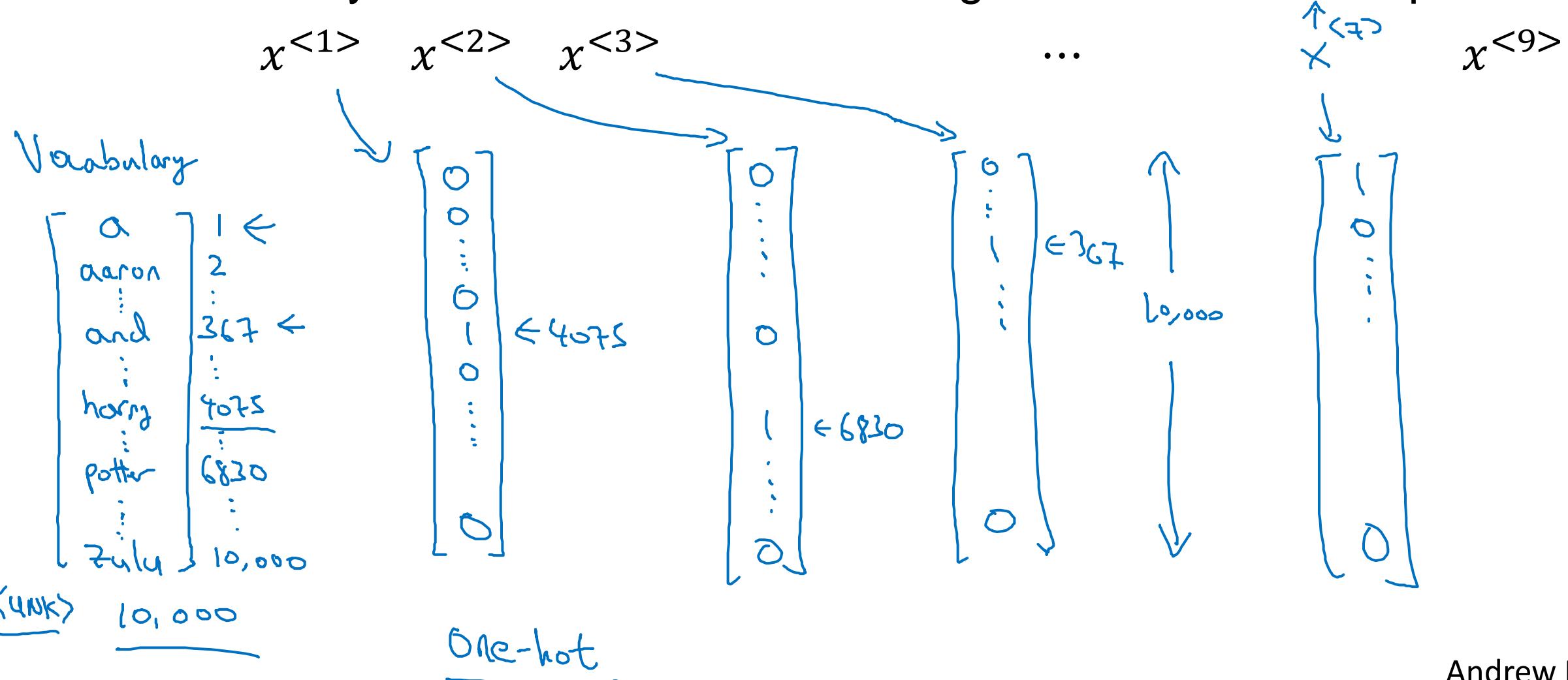
Representing words

$$x^{<\leftrightarrow>} \quad x \rightarrow y$$

(x, y)

$x:$

Harry Potter and Hermione Granger invented a new spell.



Representing words

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

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad \dots \quad x^{<9>}$

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

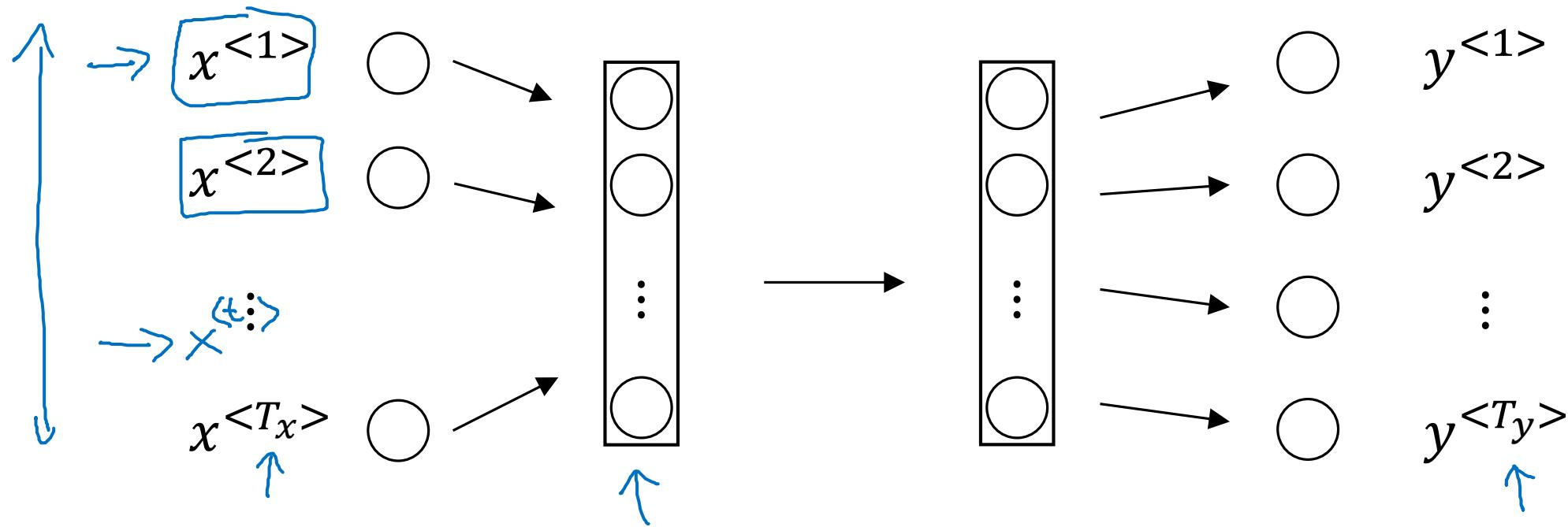


deeplearning.ai

Recurrent Neural Networks

Recurrent Neural Network Model

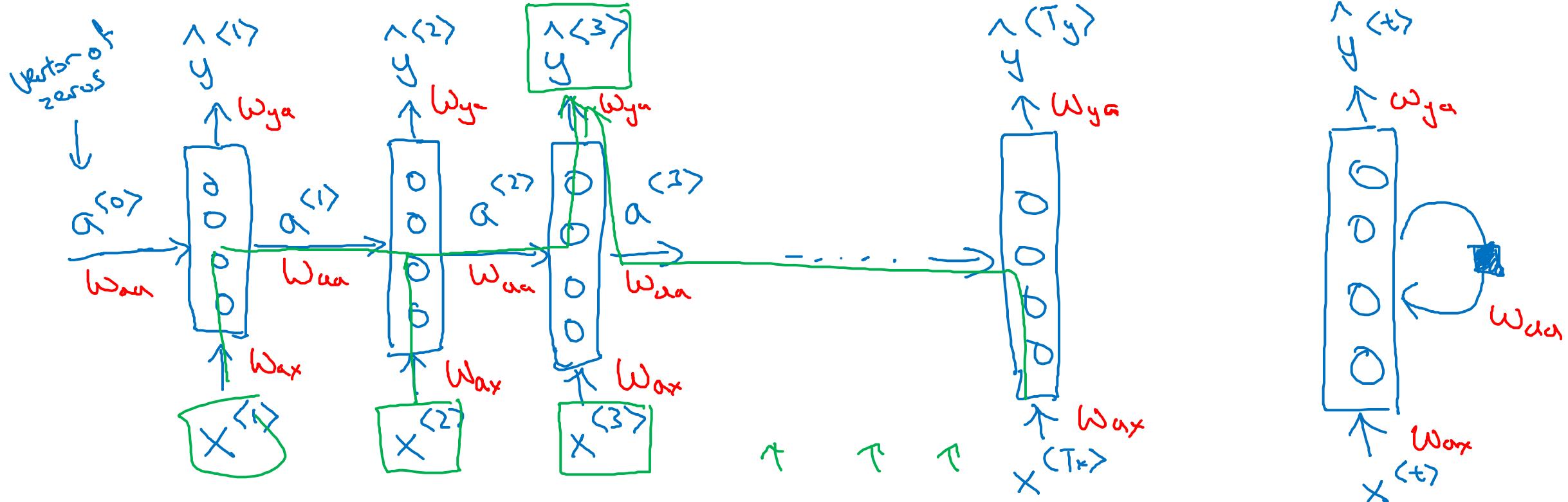
Why not a standard network?



Problems:

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

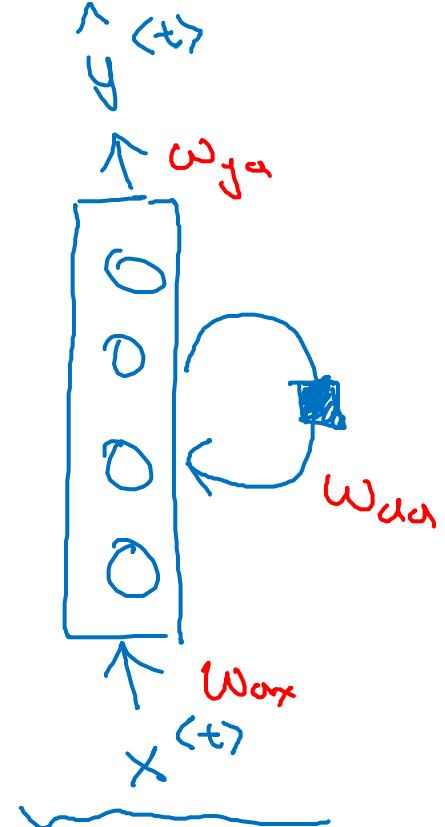
Recurrent Neural Networks



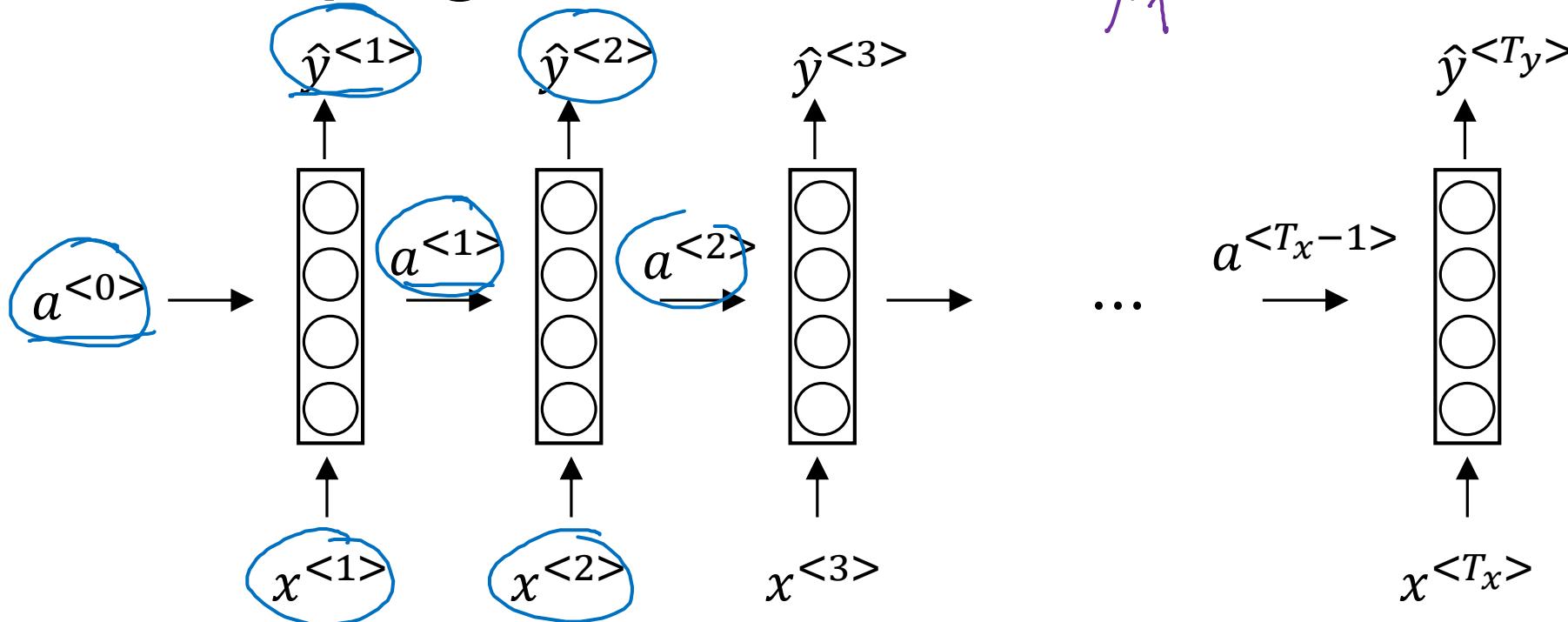
Bidirectional RNN (BRNN)

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

He said, "Teddy bears are on sale!"



Forward Propagation



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

$$a^{(i)} = g_i(W_a a^{(i-1)} + W_x x^{(i)} + b_a) \leftarrow \tanh \text{ or } \text{ReLU}$$

$$\hat{y}^{(i)} = g_i(W_y a^{(i)} + b_y) \leftarrow \text{Sigmoid}$$

$$a^{(t)} = g(W_a a^{(t-1)} + W_x x^{(t)} + b_a)$$

$$\hat{y}^{(t)} = g(W_y a^{(t)} + b_y)$$

Simplified RNN notation

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

W_{aa} (100, 100) W_{ax} (100, 10,000) $x^{(t)}$ (10,000)

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

$$\hat{y}^{(t)} = g(W_y a^{(t)} + b_y)$$

W_y (10,000) $a^{(t)}$ (100) b_y (10,000)

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

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

W_{aa} (100, 100) W_{ax} (100, 10,000)

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

$a^{(t-1)}$ (100) $x^{(t)}$ (10,000) $[a^{(t-1)}, x^{(t)}]$ (10,000)

$$W_{aa} [a^{(t-1)}, x^{(t)}] = W_{aa} a^{(t-1)} + W_{ax} x^{(t)}$$

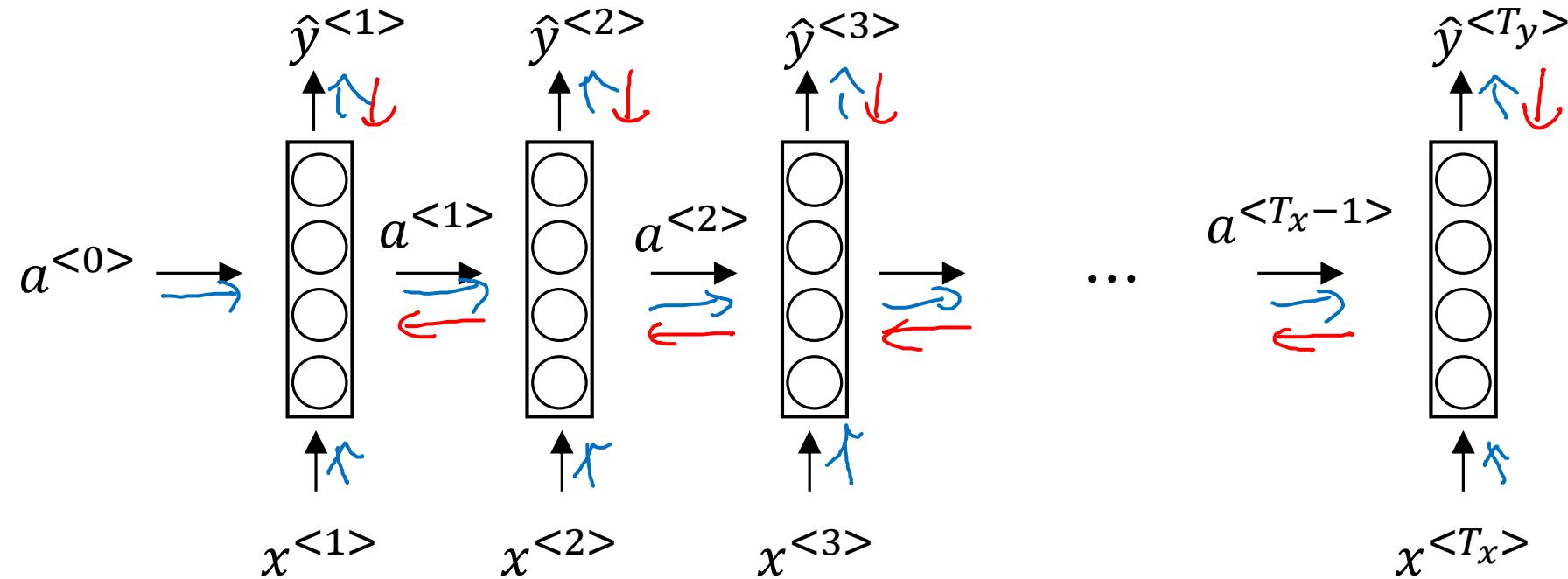


deeplearning.ai

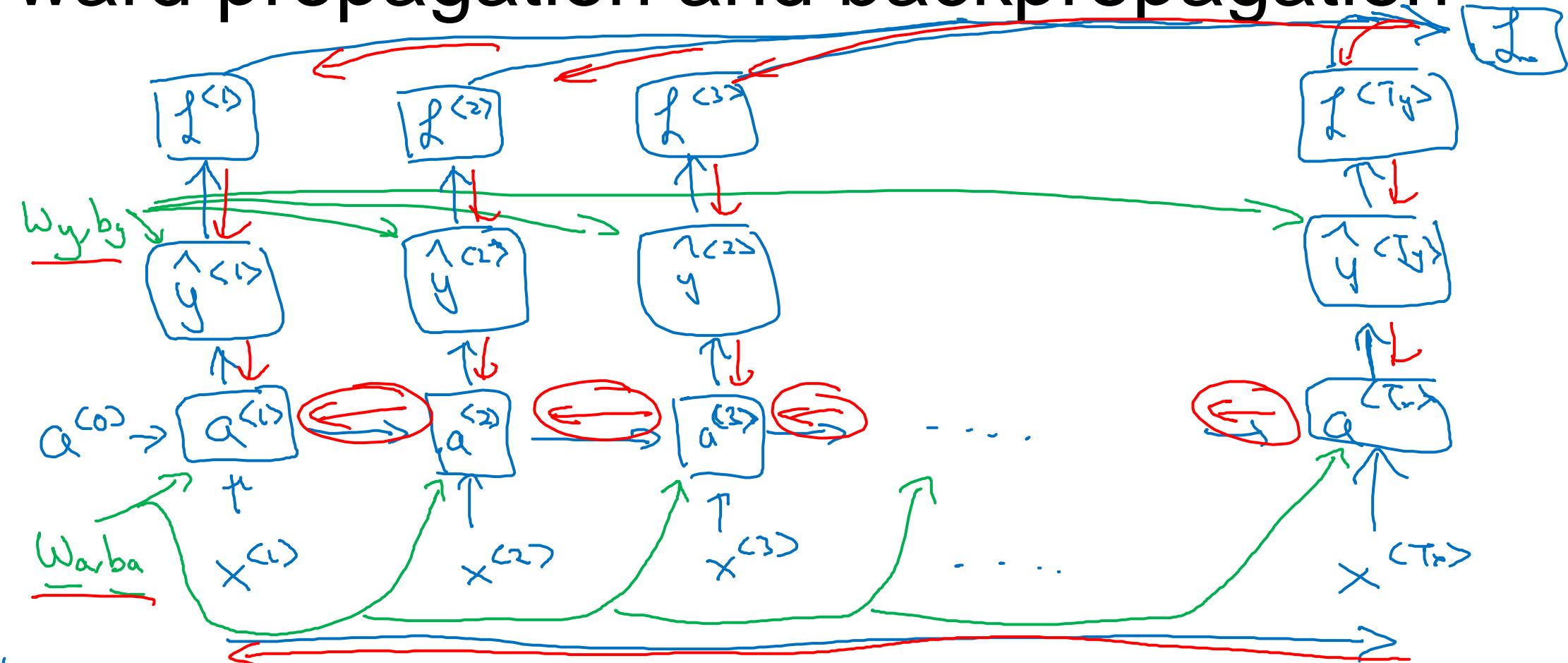
Recurrent Neural Networks

Backpropagation through time

Forward propagation and backpropagation



Forward propagation and backpropagation



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

$$\mathcal{L}(\hat{y}, y) = \sum_{t=1}^{T_{\text{out}}} \mathcal{L}^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

Backpropagation through time



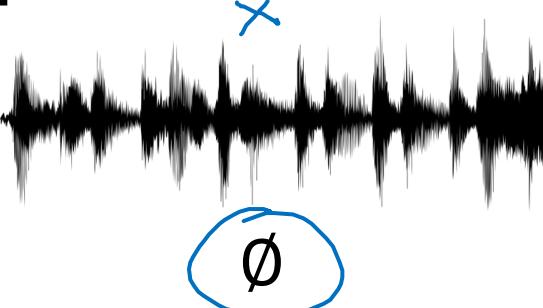
deeplearning.ai

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



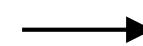
Sentiment classification

“There is nothing to like in
this movie.”



DNA sequence analysis

AGCCCCTGTGAGGAAC TAG



AGCCCCTGTGAGGAAC **TAG**

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.

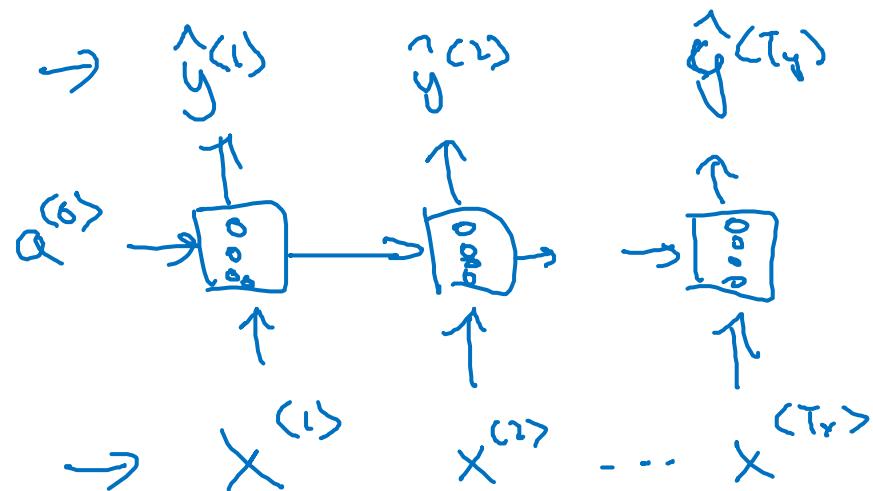


Yesterday, **Harry Potter** met
Hermione Granger.

Andrew Ng

Examples of RNN architectures

$$T_x = T_y$$

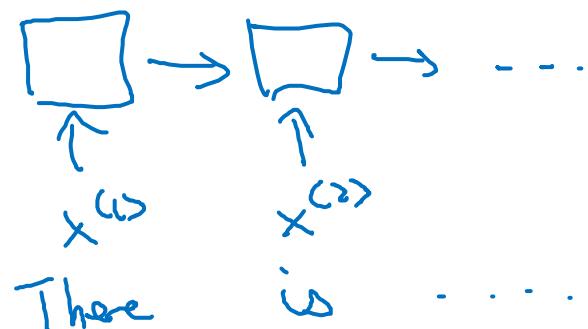


Many-to-many

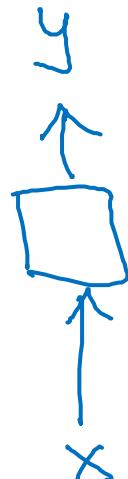
Sentiment classification

$x = \text{text}$

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

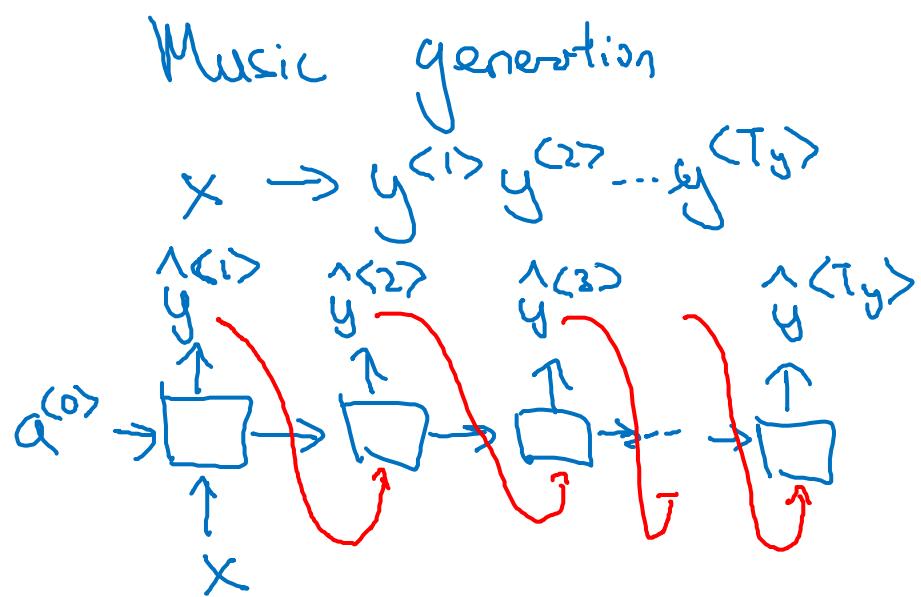


Many-to-one



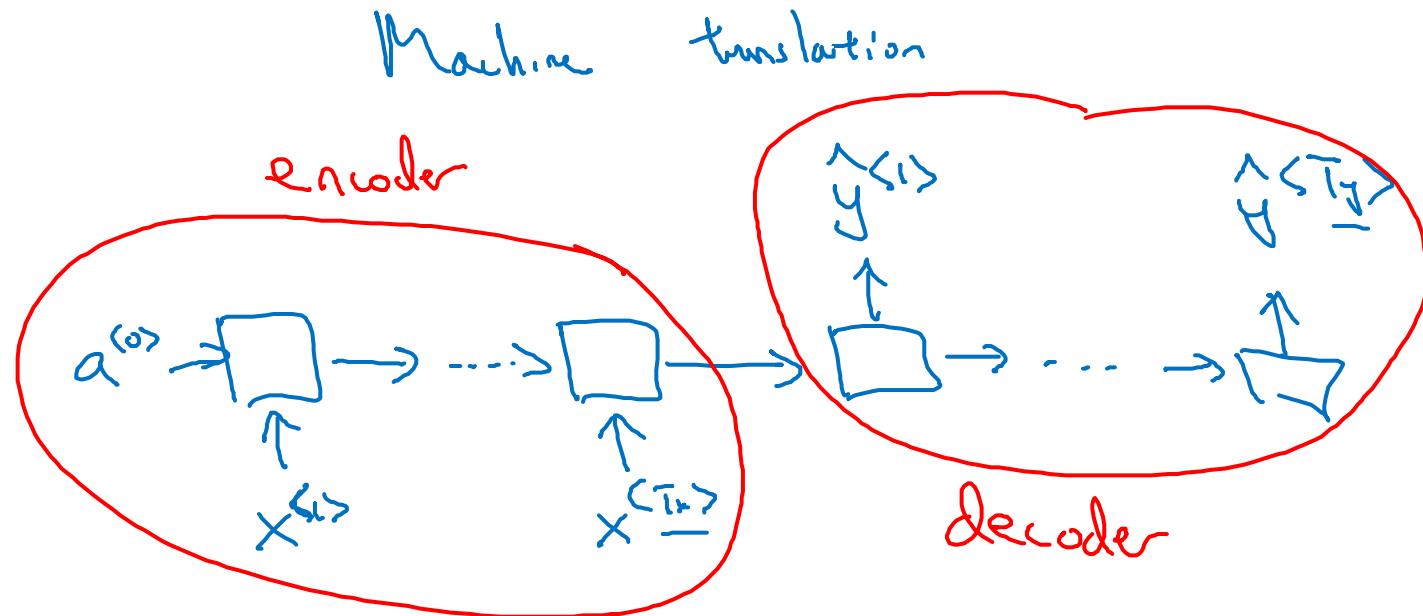
One-to-one

Examples of RNN architectures



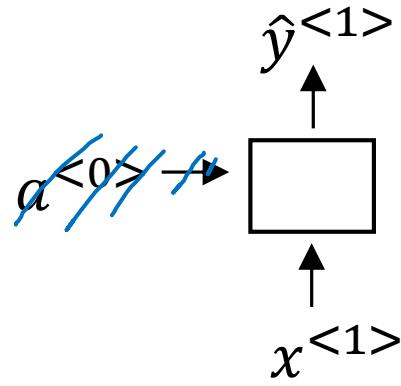
One-to-many

$$x = \phi$$

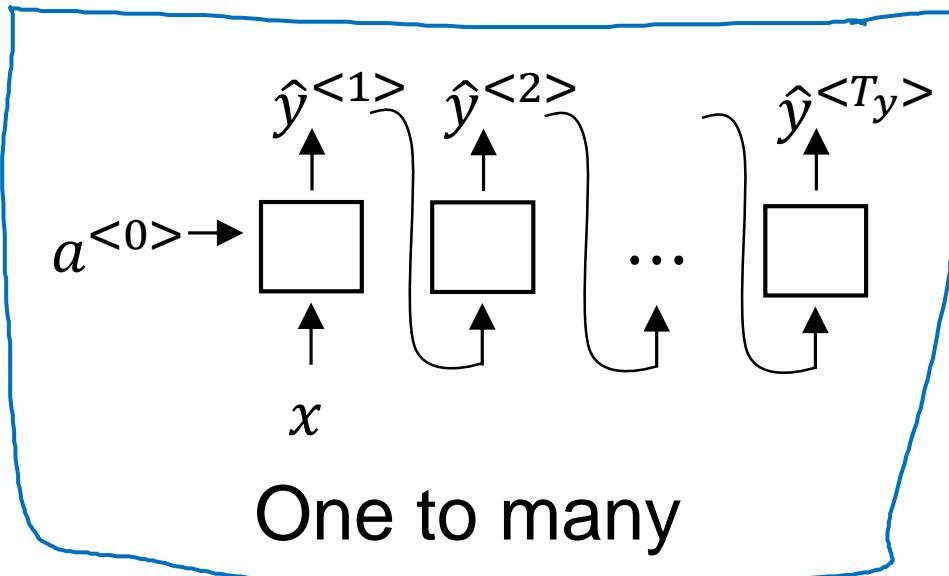


Many-to-many

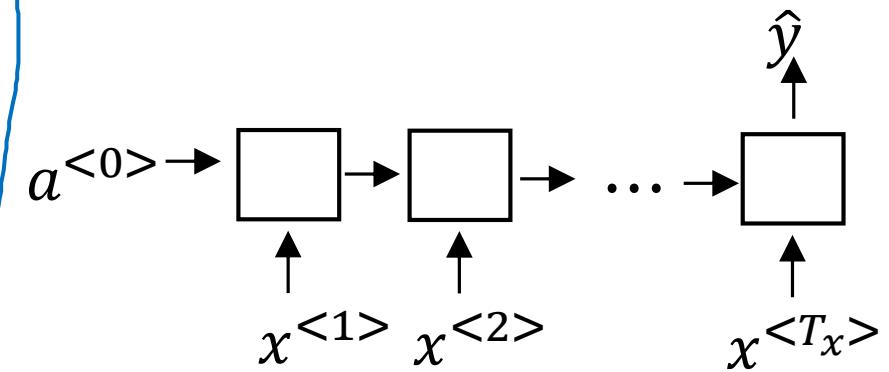
Summary of RNN types



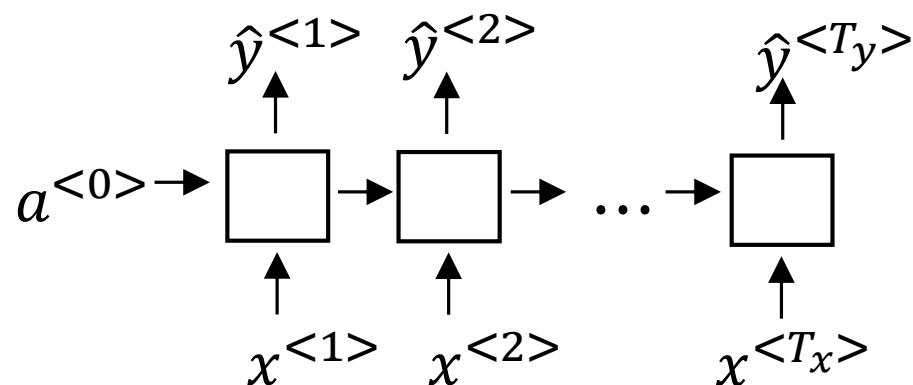
One to one



One to many

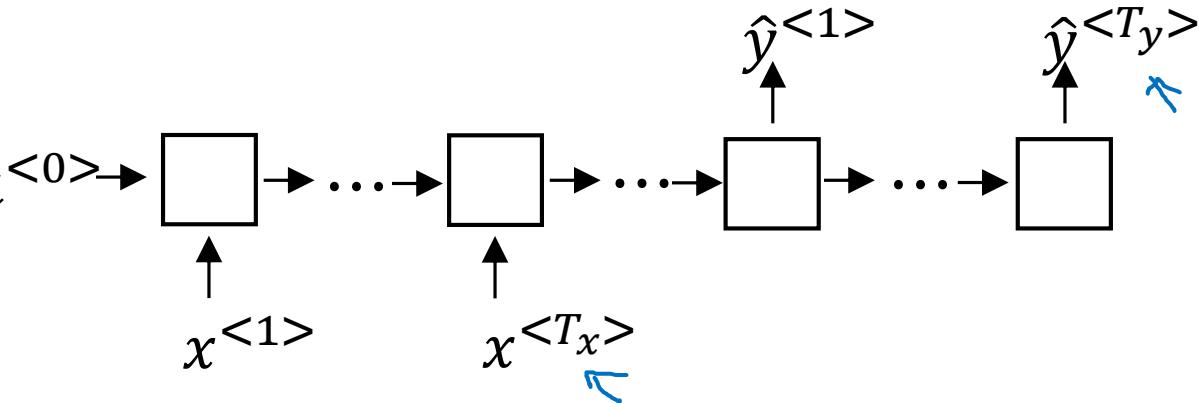


Many to one



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.

$$P(\text{The apple and pair salad}) = 3.2 \times 10^{-3}$$

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

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

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

Language modelling with an RNN

Training set: large corpus of english text.

Tokenize

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

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

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

...

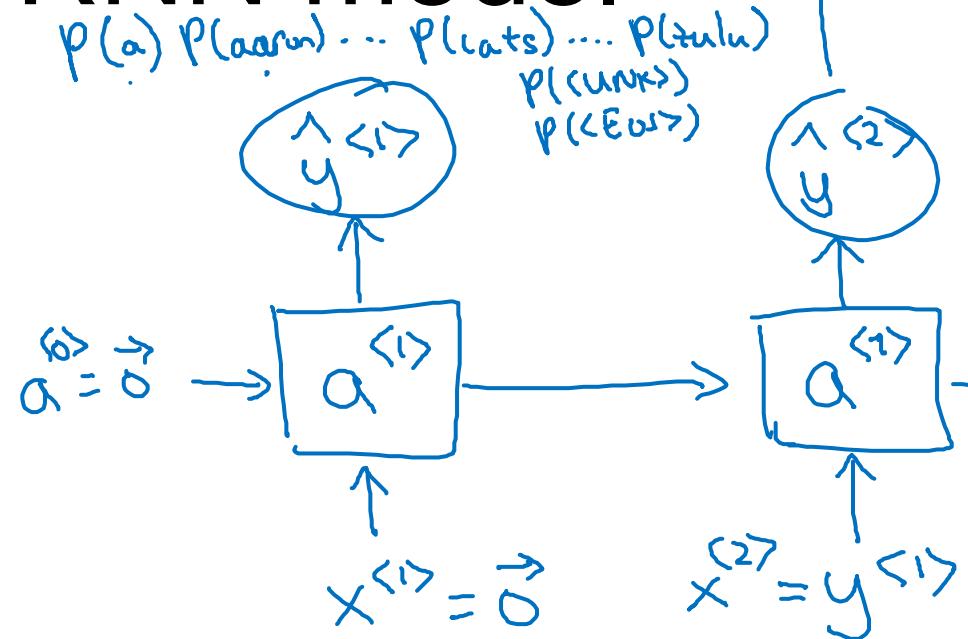
$y^{(8)}$ $y^{(9)}$

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

$\langle \text{UNK} \rangle$

10,000

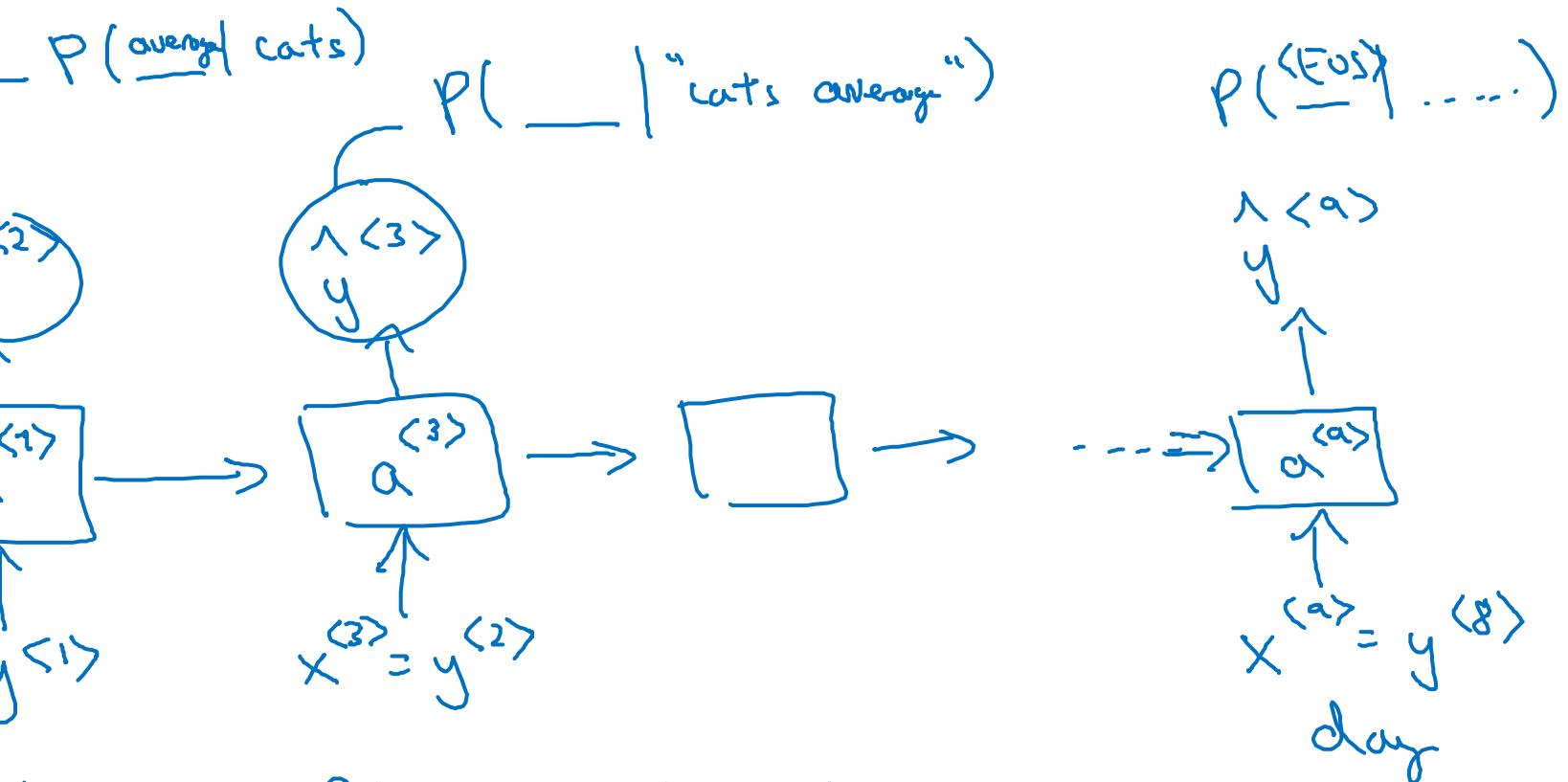
RNN model



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

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

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



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

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

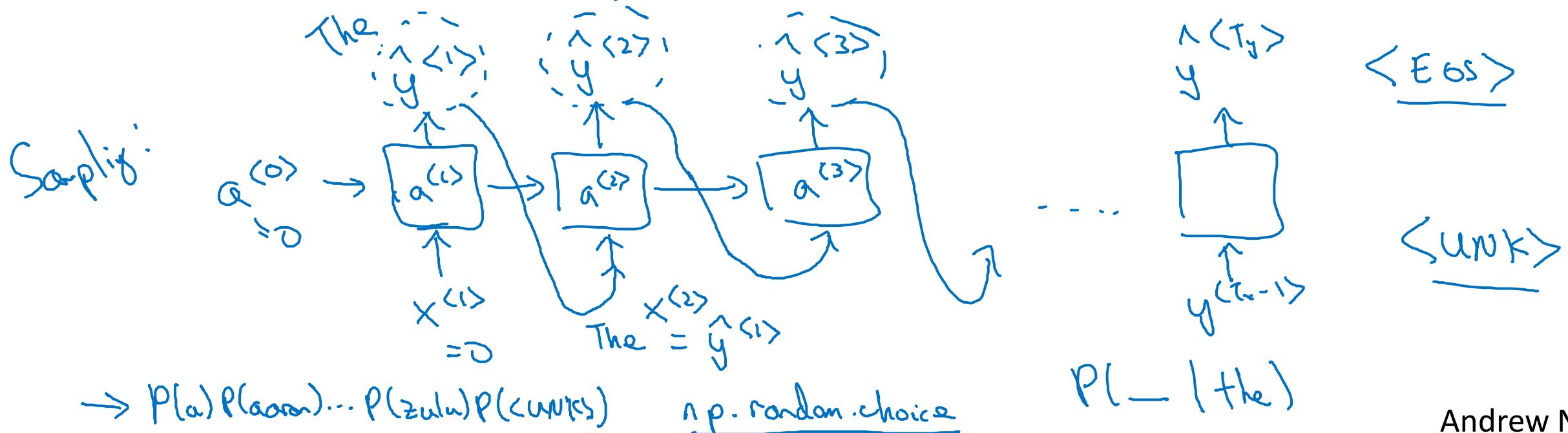
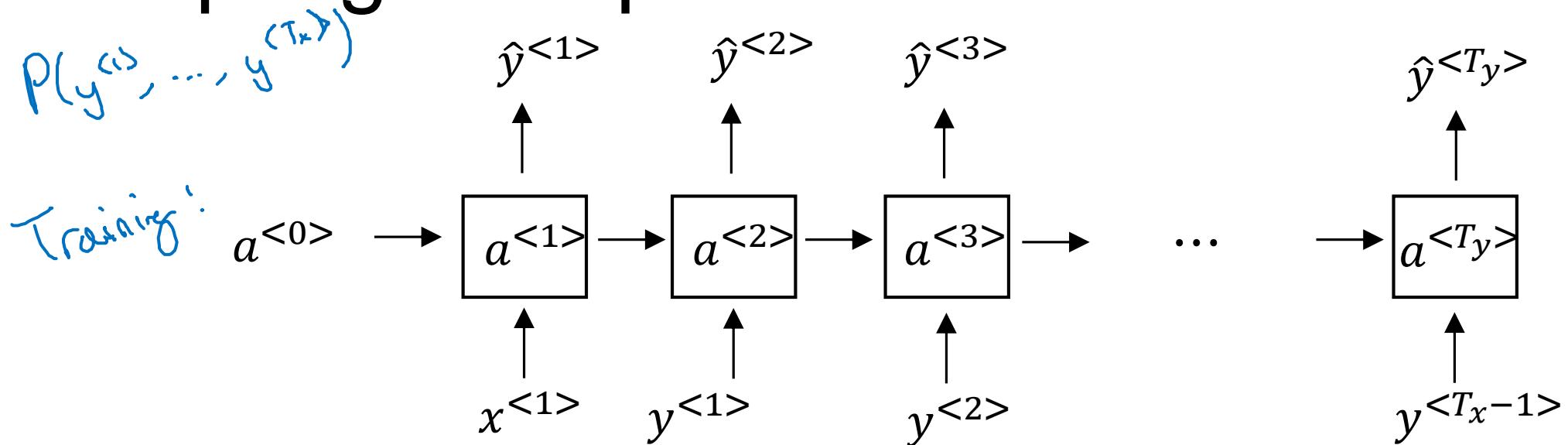


deeplearning.ai

Recurrent Neural Networks

Sampling novel
sequences

Sampling a sequence from a trained RNN



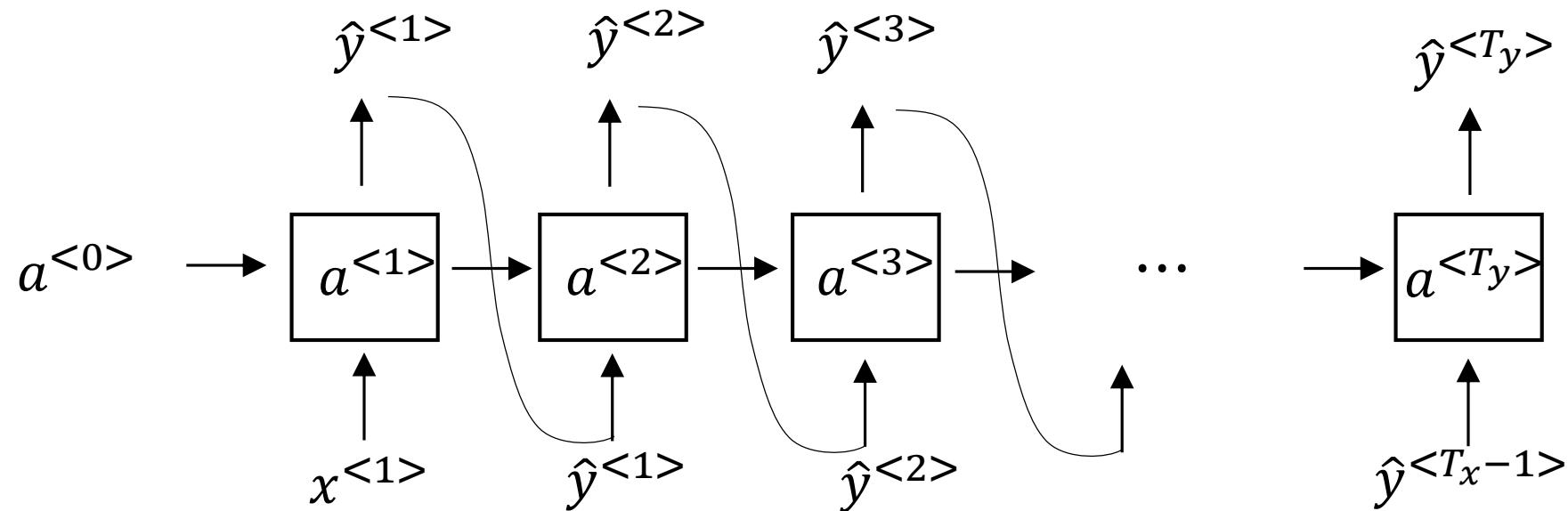
Character-level language model

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

$$y^{(0)} - y^{(1)} = y^{(2)} - y^{(3)}$$

Cat average
↑ ↑ ↑ ↑ . . .

May



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 lesser be my love to me see sabl's.
For whose are ruse of mine eyes heaves.

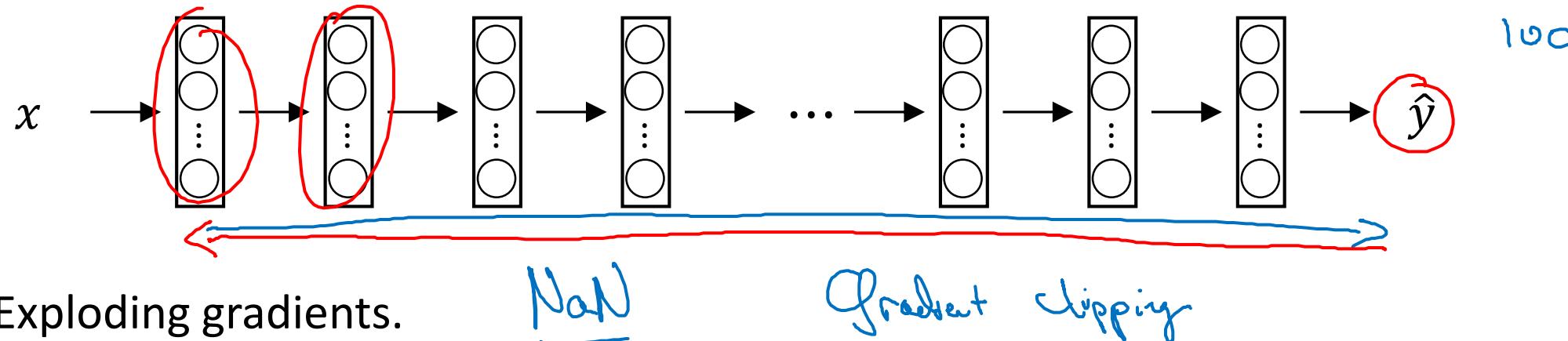
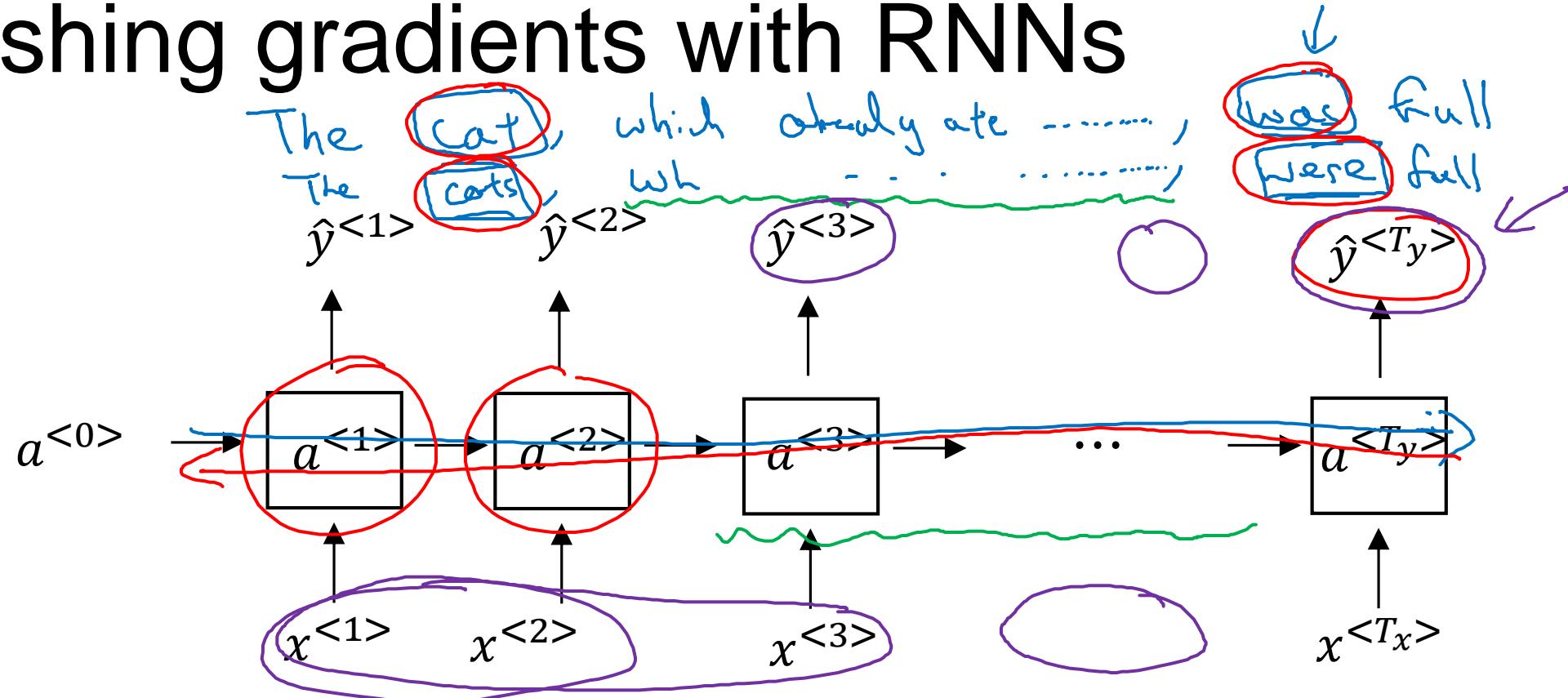


deeplearning.ai

Recurrent Neural Networks

Vanishing gradients with RNNs

Vanishing gradients with RNNs



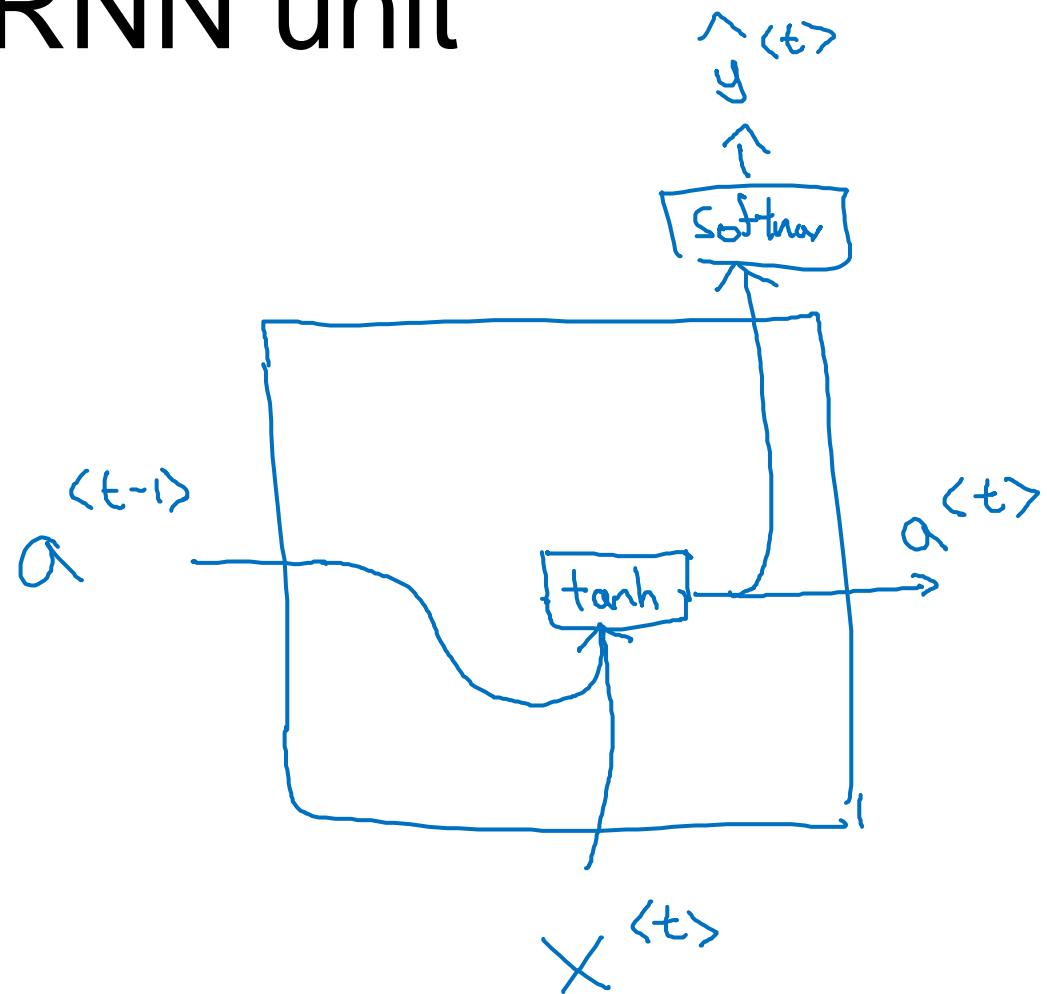


deeplearning.ai

Recurrent Neural Networks

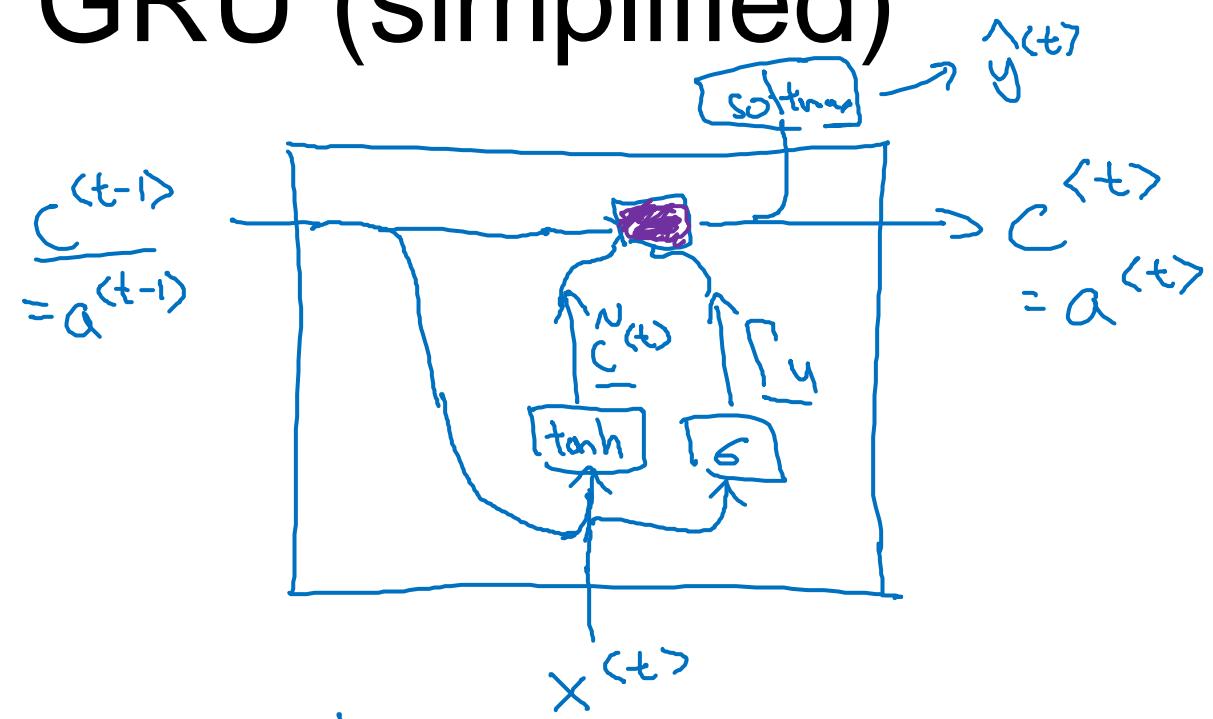
Gated Recurrent Unit (GRU)

RNN unit

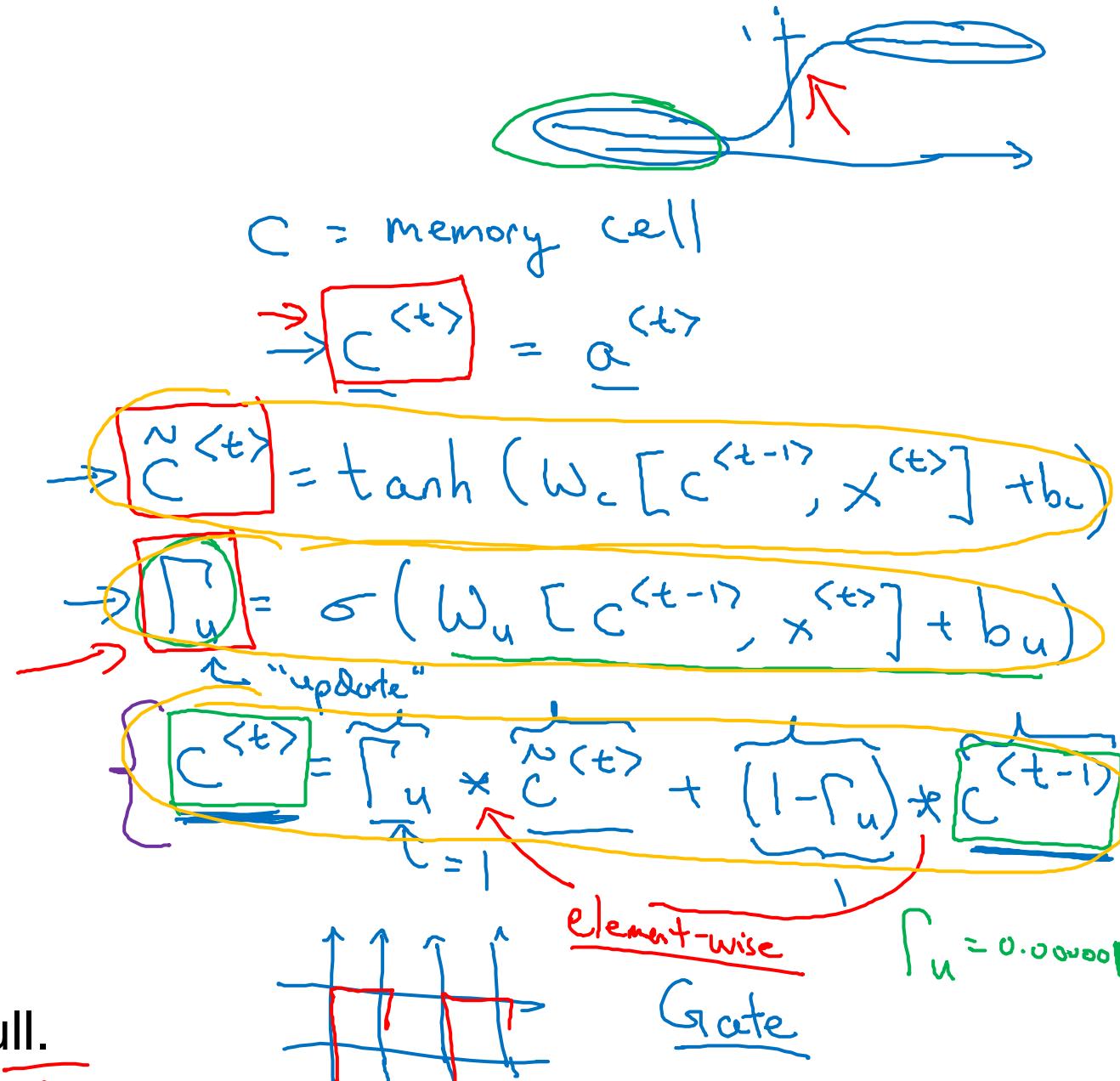


$$\underline{a^{(t)}} = \text{tanh}(W_a[\underline{a^{(t-1)}, x^{(t)}] + b_a})$$

GRU (simplified)



$f_u = 1$
 $f_u = 0$ $i_u = 0$ $N_c = 0$ $N_u = 0$ \dots
 \rightarrow 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

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

$$\begin{matrix} u \\ r \end{matrix} \quad \left\{ \begin{matrix} \Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u) \\ \Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r) \end{matrix} \right.$$

LSTM

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

The cat, which ate already, was full.



deeplearning.ai

Recurrent Neural Networks

LSTM (long short term memory) unit

GRU and LSTM

GRU

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

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

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

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

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

$$\Gamma_f$$

(output)

LSTM

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

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

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

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

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

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

LSTM in pictures

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

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

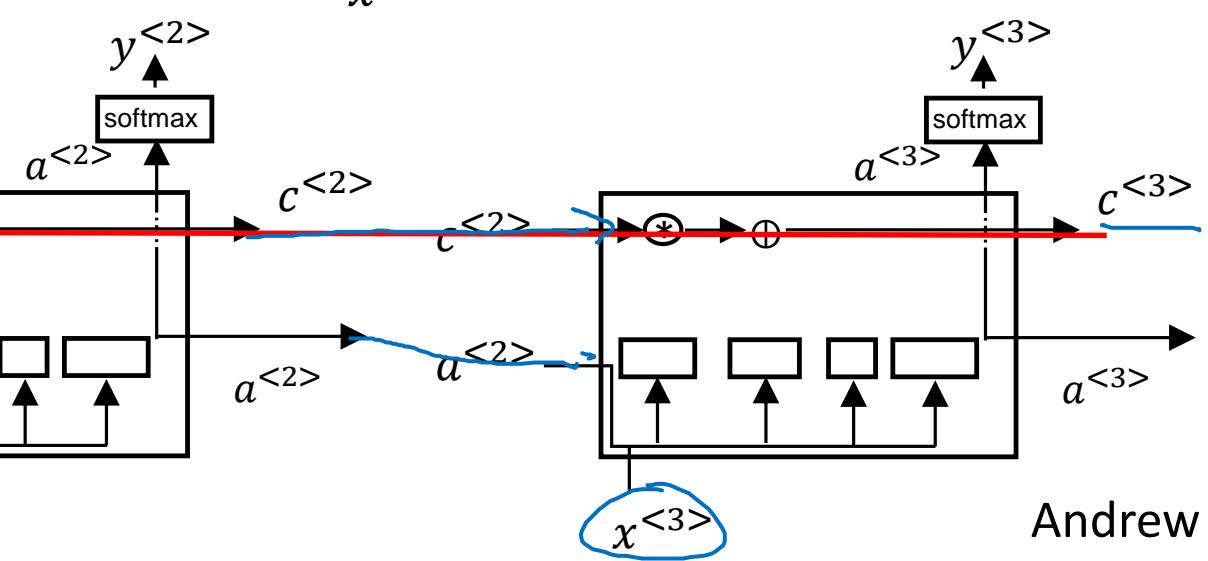
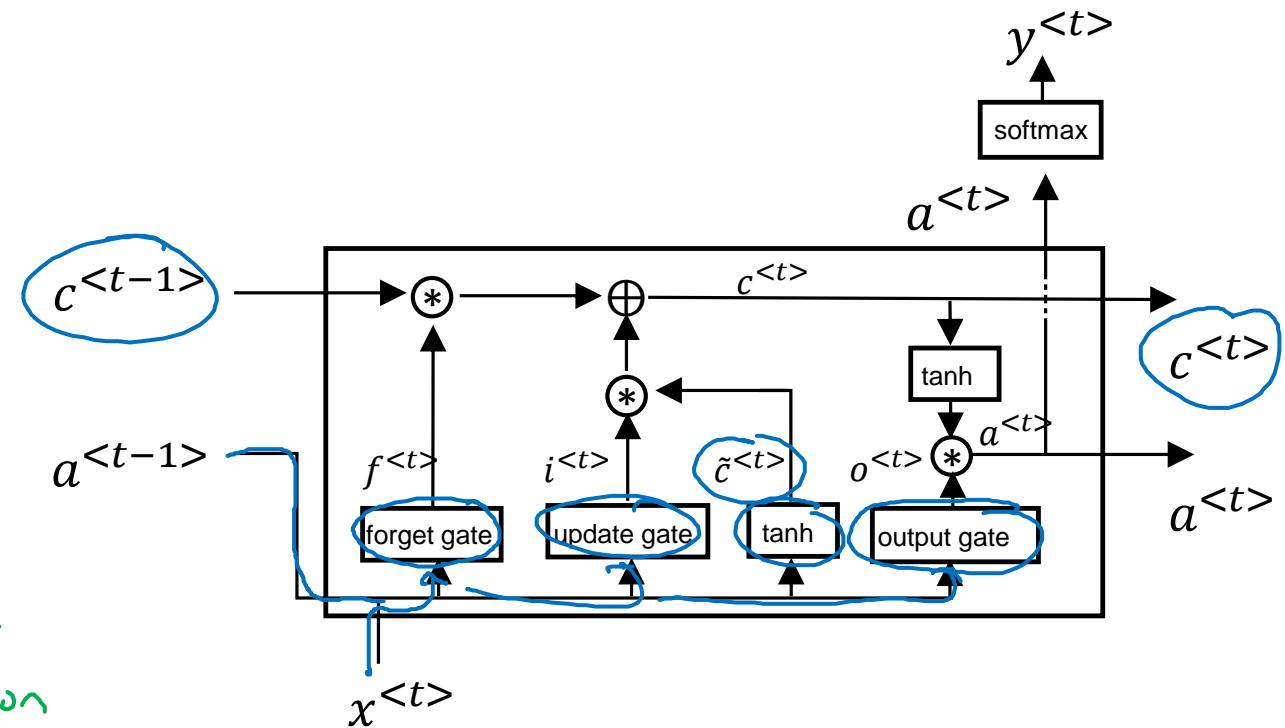
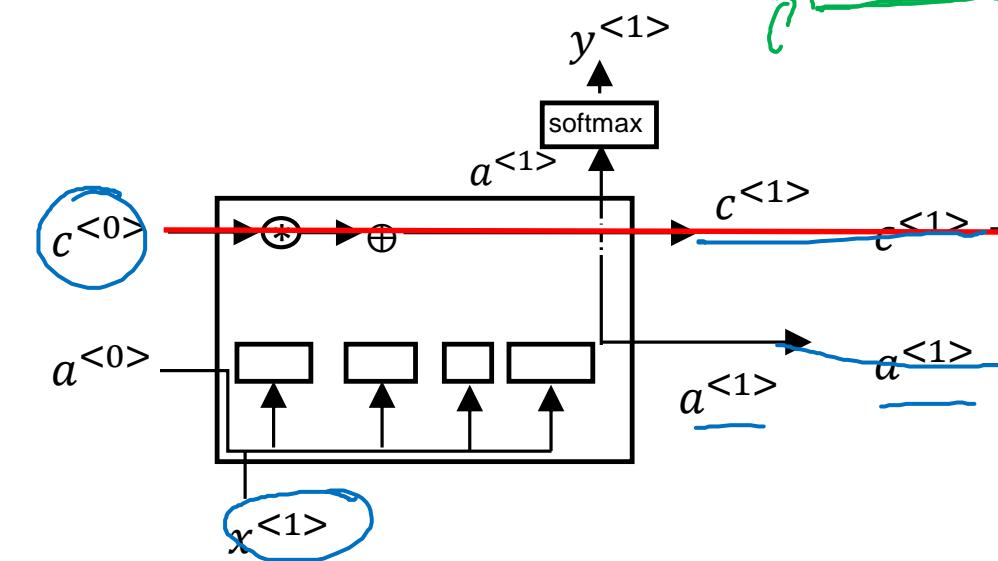
$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

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

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

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

peephole connection



Andrew Ng



deeplearning.ai

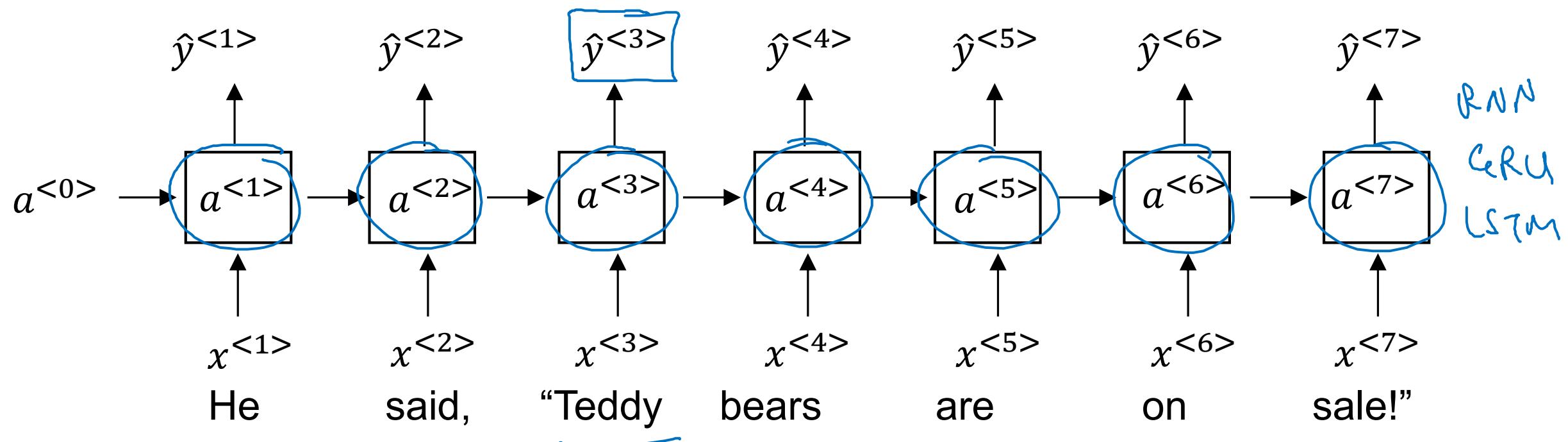
Recurrent Neural Networks

Bidirectional RNN

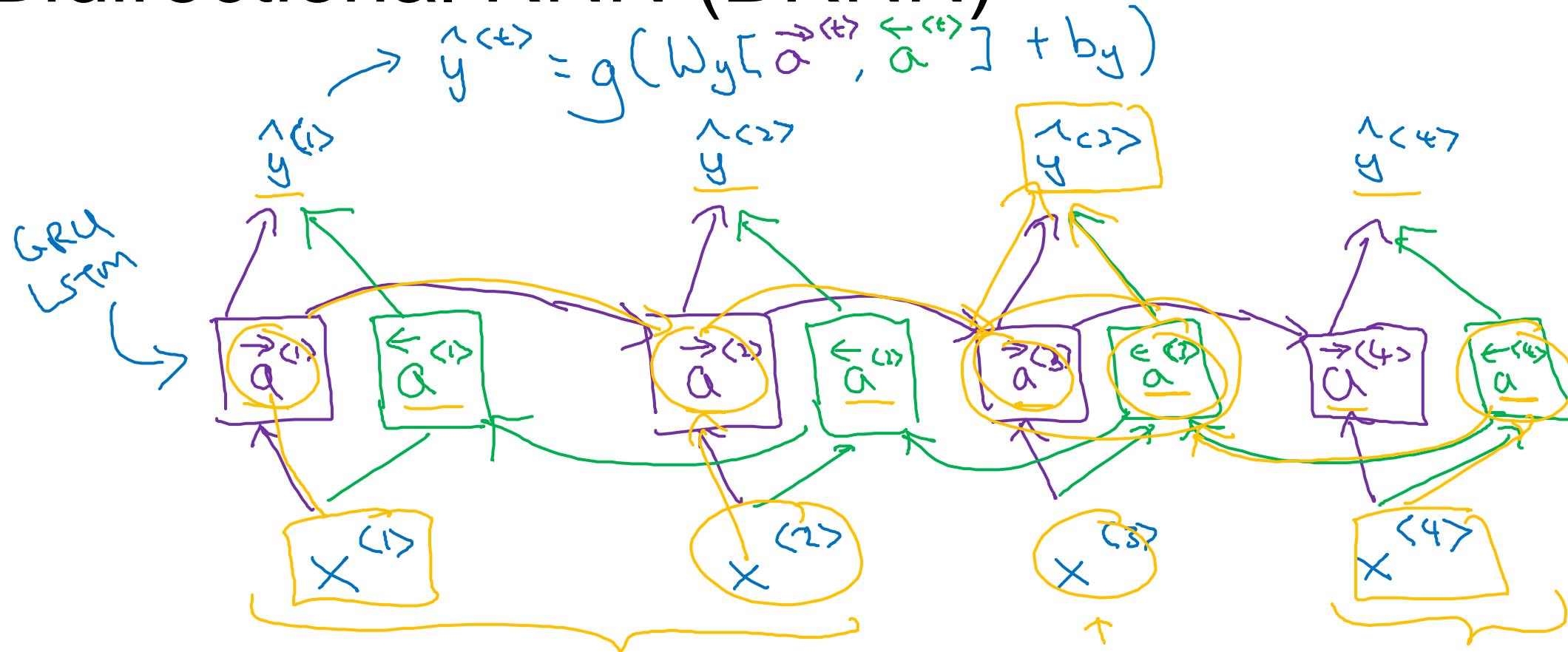
Getting information from the future

He said, “Teddy bears are on sale!”

He said, “Teddy Roosevelt was a great President!”



Bidirectional RNN (BRNN)



Acyclic graph

BRNN w/ LSTM

He said,

"Teddy Roosevelt ..."