



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 → 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

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

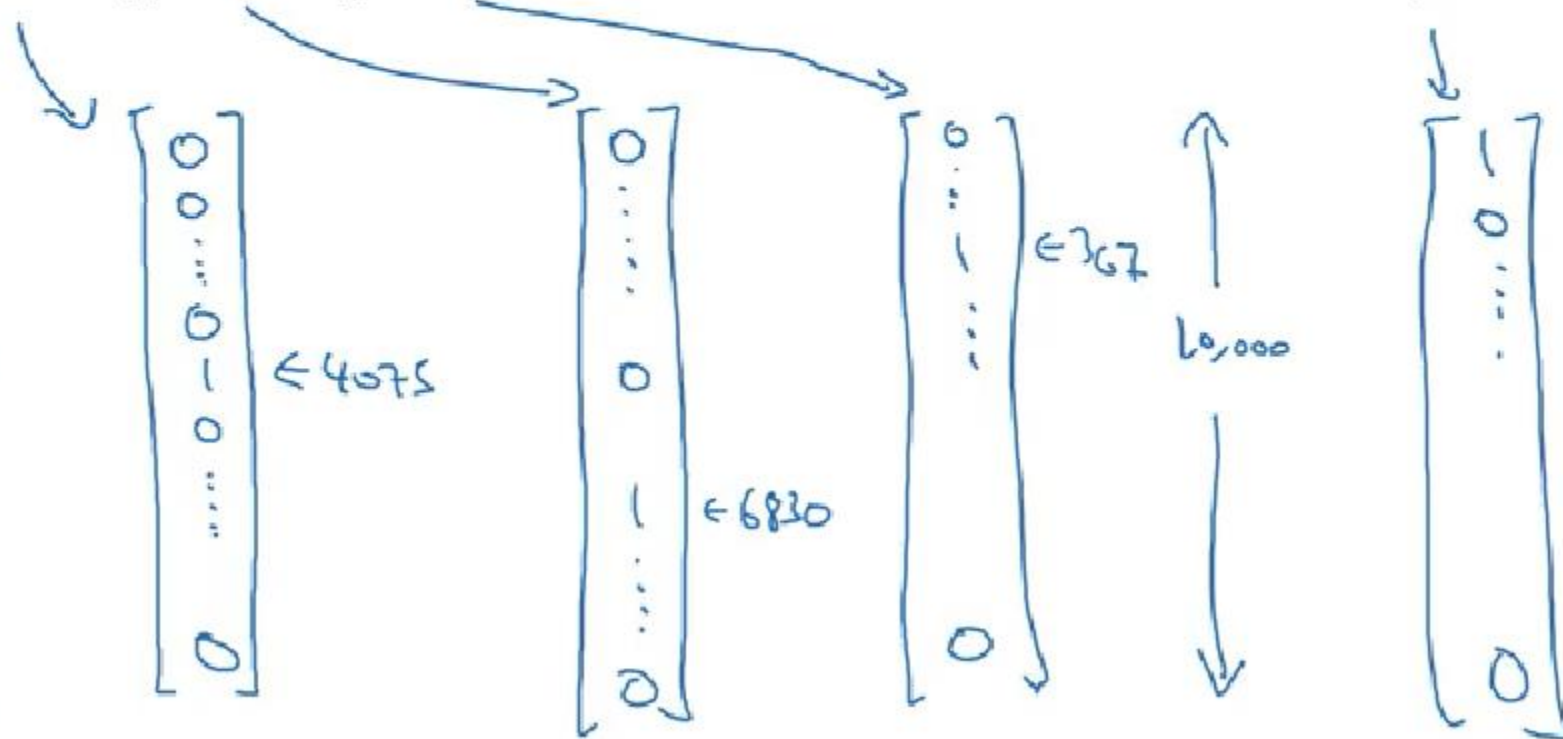
x: Harry Potter and Hermione Granger invented a new spell.

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

Vocabulary

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

<UNK> 10,000



One-hot

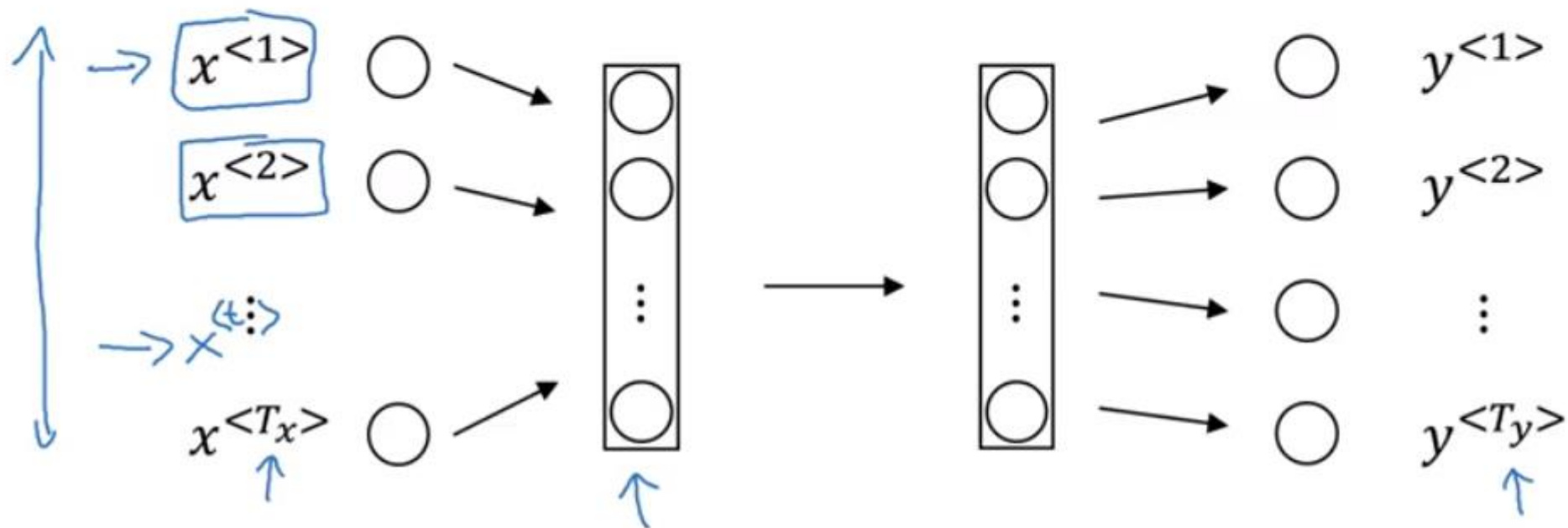


deeplearning.ai

Recurrent Neural Networks

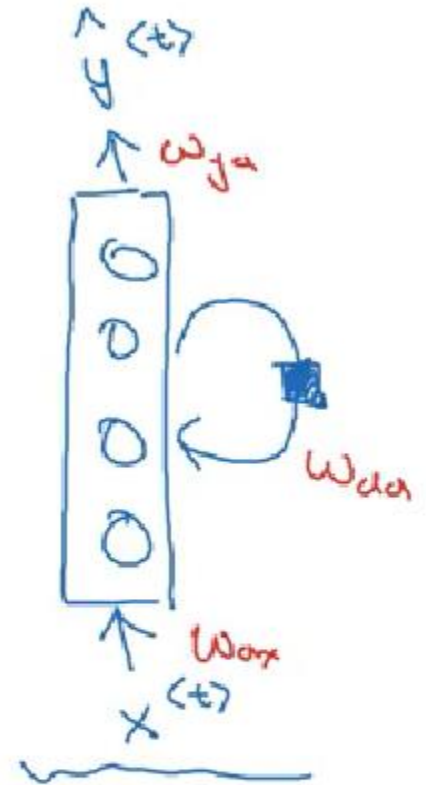
Recurrent Neural Network Model

Why not a standard network?



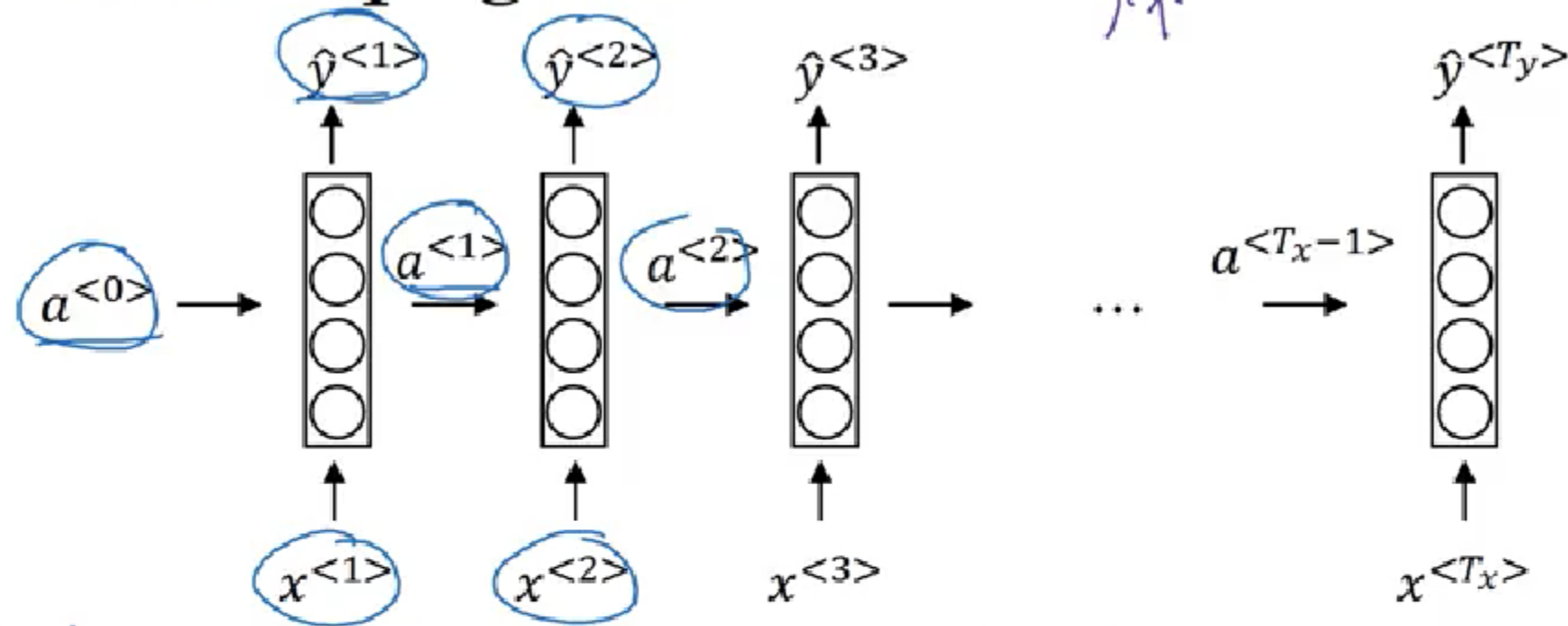
Problems:

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

$$T_x = T_y$$


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 100×100 matrix (indicated by a green box and arrows).
- $a^{<t-1>}$ is a 100 vector (indicated by a blue box and arrow).
- W_{ax} is a $100 \times 10,000$ matrix (indicated by a blue box and arrows).
- $x^{<t>}$ is a $10,000$ vector (indicated by a blue box and arrow).

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

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

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

- W_y is a 1×100 vector (indicated by a blue box and arrow).
- $a^{<t>}$ is a 100 vector (indicated by a blue box and arrow).
- b_y is a 1 scalar (indicated by a blue box and arrow).

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

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

- W_a is a 100×10100 matrix (indicated by a green box and arrows).
- $[a^{<t-1>}, x^{<t>}]$ is a 10100 vector (indicated by a blue box and arrow).
- b_a is a 100 vector (indicated by a blue box and arrow).

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

Diagram illustrating the dimensions of the weights in the fourth equation:

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

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

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

- $\begin{bmatrix} W_{aa} & W_{ax} \end{bmatrix}$ is a 100×10100 matrix (indicated by a green box and arrows).
- $\begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix}$ is a 10100 vector (indicated by a green box and arrows).
- $W_{aa}a^{<t-1>} + W_{ax}x^{<t>}$ is a 100 vector (indicated by a green box and arrow).

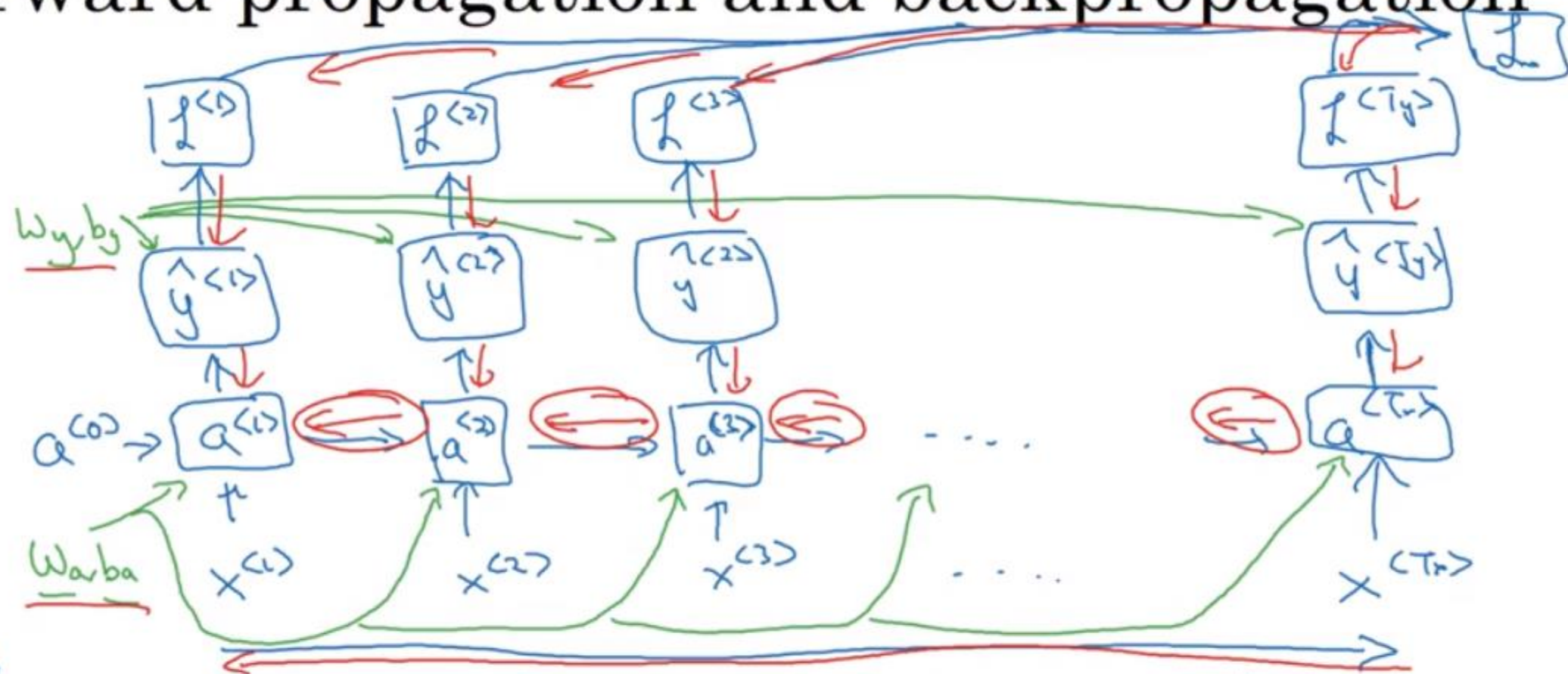


deeplearning.ai

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



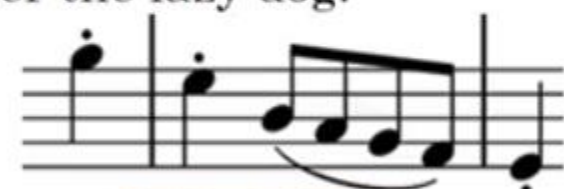
T_x

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

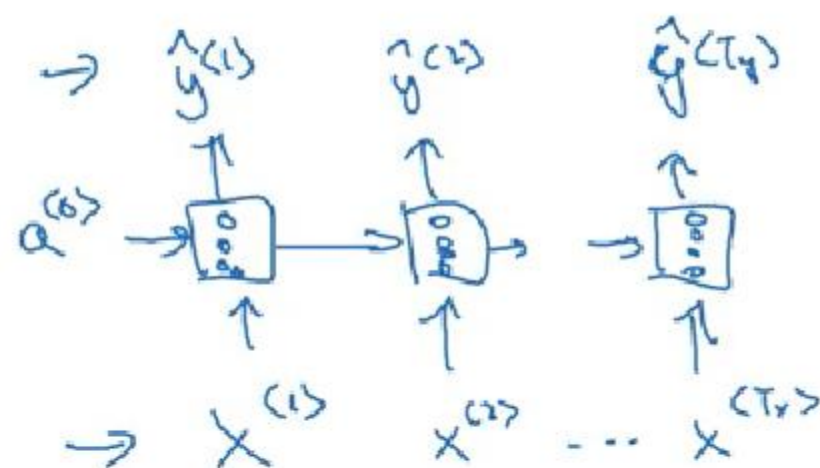
Name entity recognition

Yesterday, Harry Potter
met Hermione Granger.

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

Examples of RNN architectures

$$T_x = T_y$$

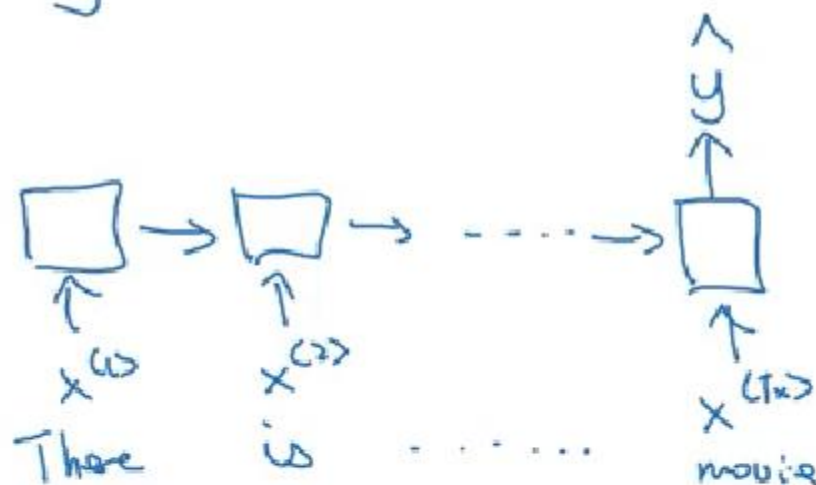


Many-to-many

Sentiment classification

$x = \text{text}$

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

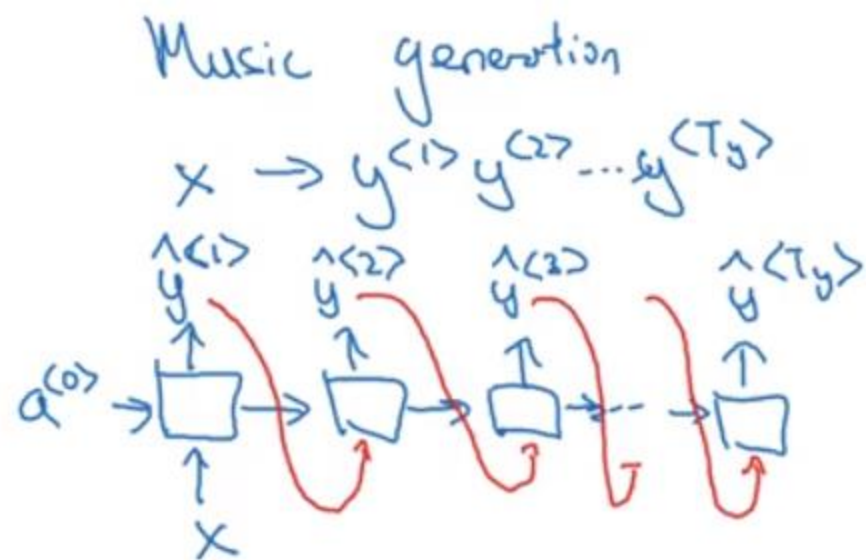


Many-to-one



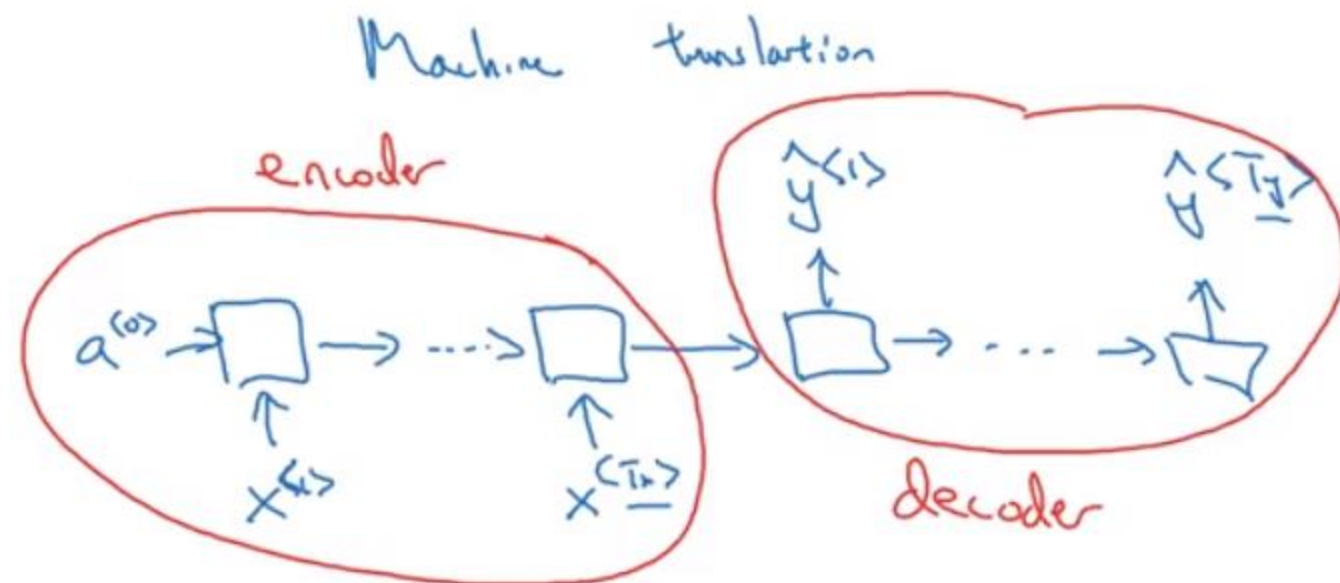
One-to-one

Examples of RNN architectures



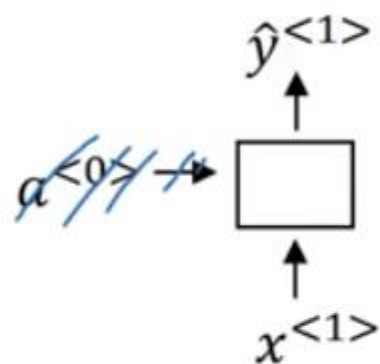
One-to-many

$$x = \phi$$

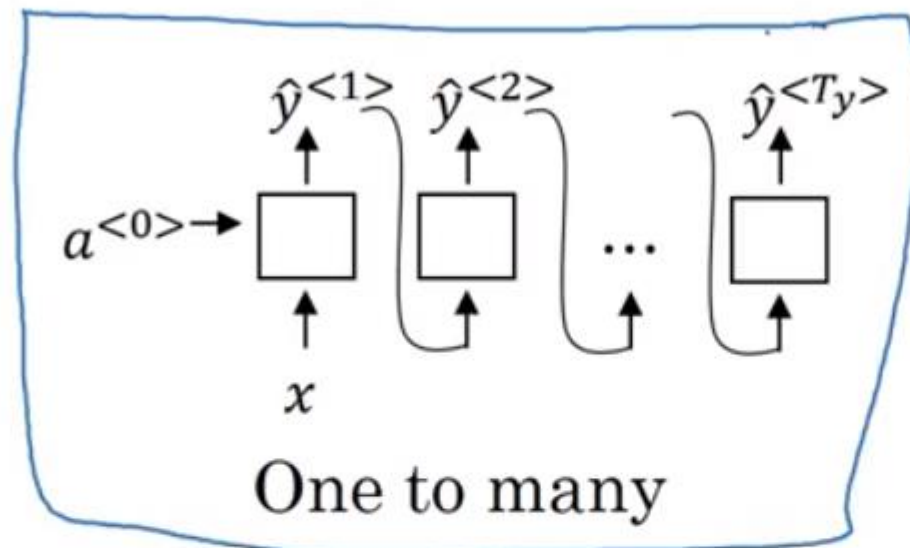


Many-to-many

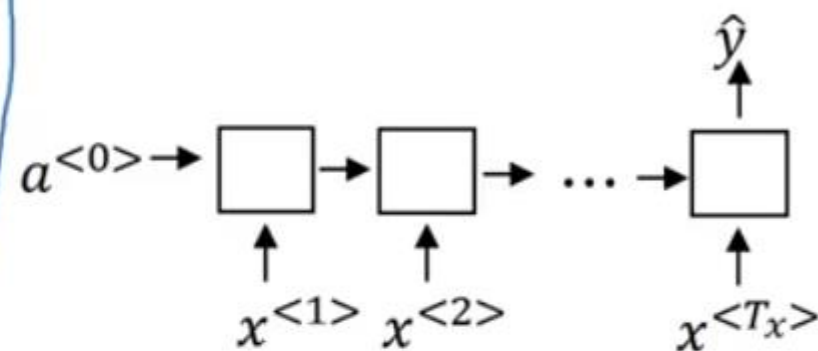
Summary of RNN types



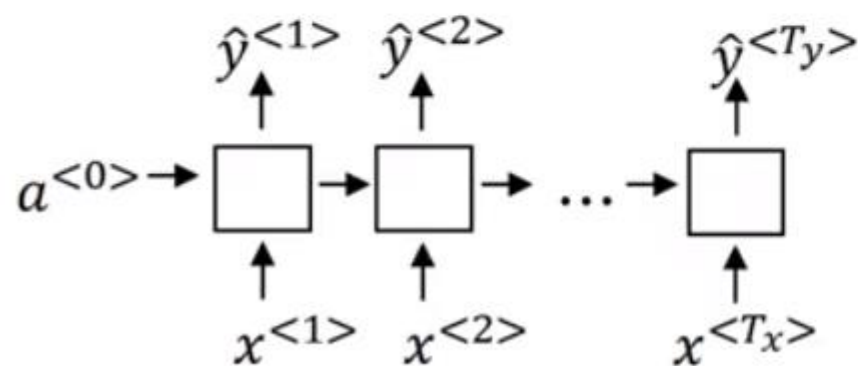
One to one



One to many

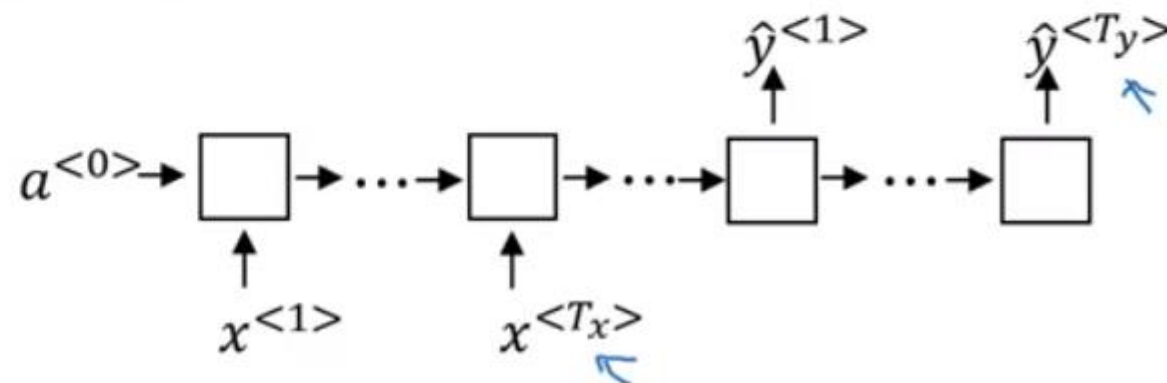


Many to one



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.

$P(\text{The apple and pair salad}) =$

$P(\text{The apple and pear salad}) =$

Language modelling with an RNN

Training set: large corpus of english text.

Cats average 15 hours of sleep a day.

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

RNN model

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

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

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

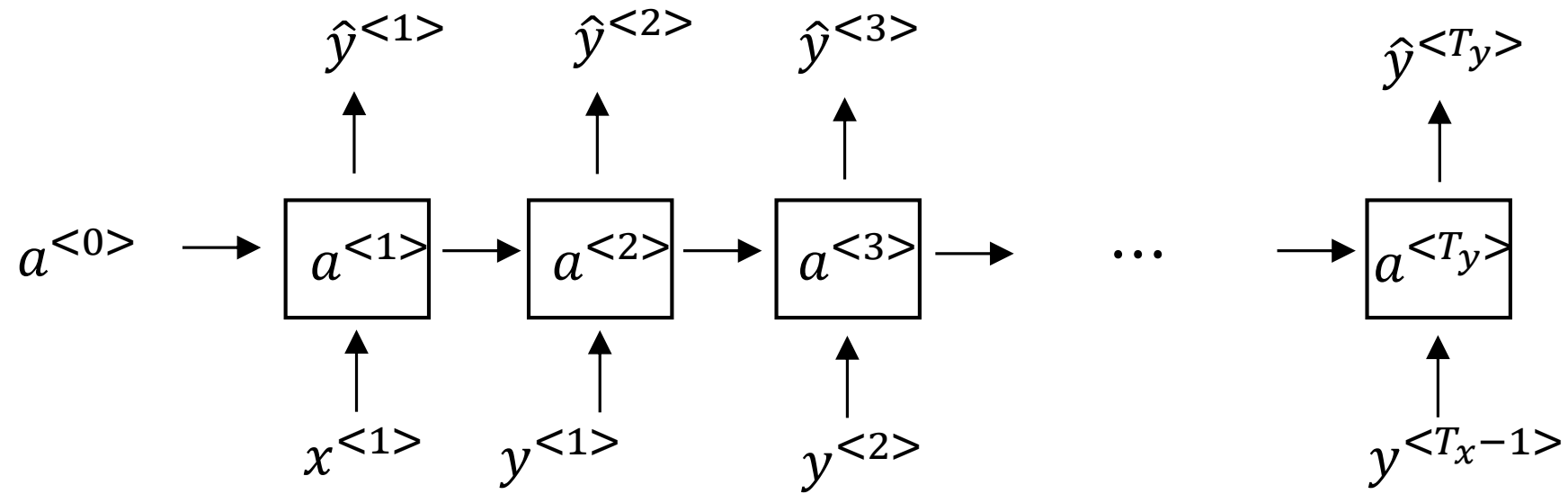


deeplearning.ai

Recurrent Neural Networks

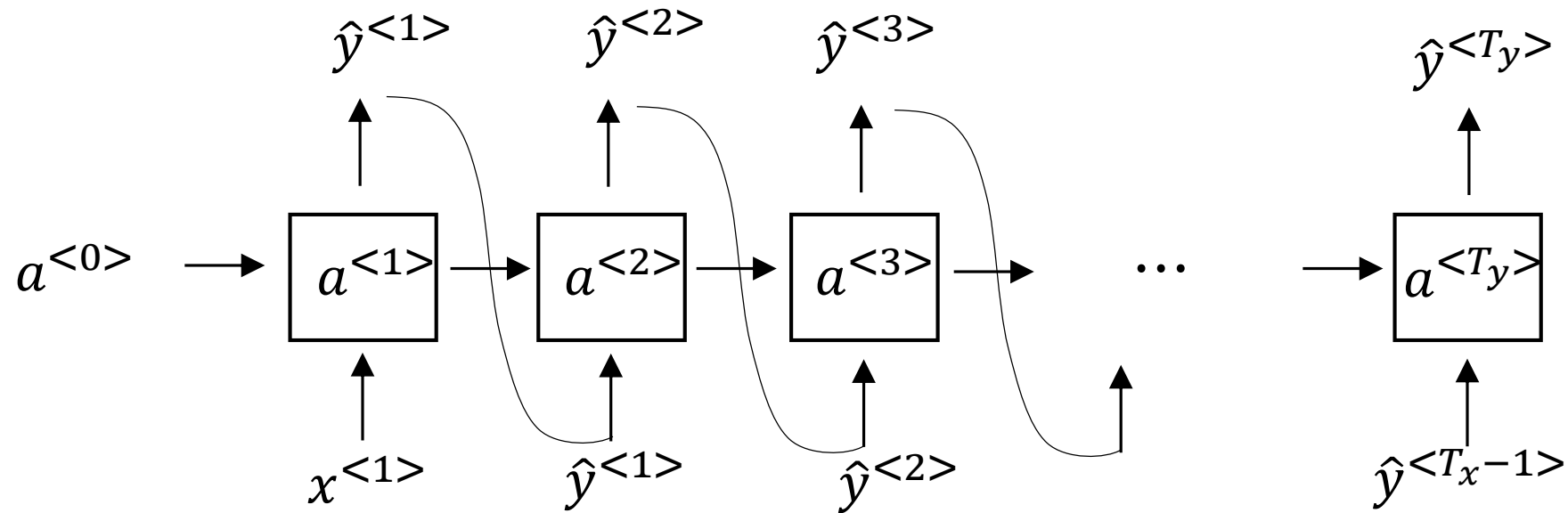
Sampling novel
sequences

Sampling a sequence from a trained RNN



Character-level language model

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



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.

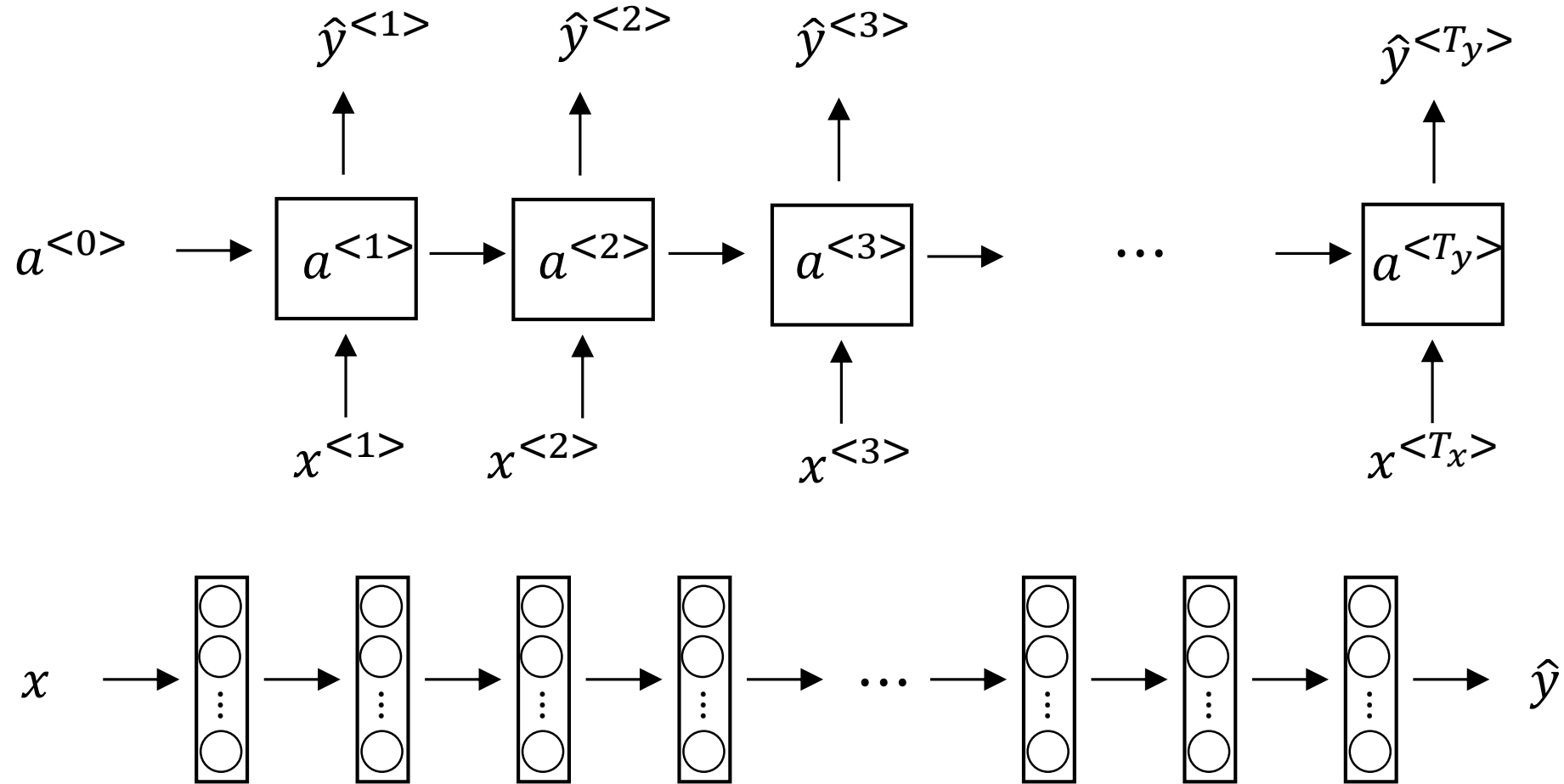


deeplearning.ai

Recurrent Neural Networks

Vanishing gradients with RNNs

Vanishing gradients with RNNs



Exploding gradients.

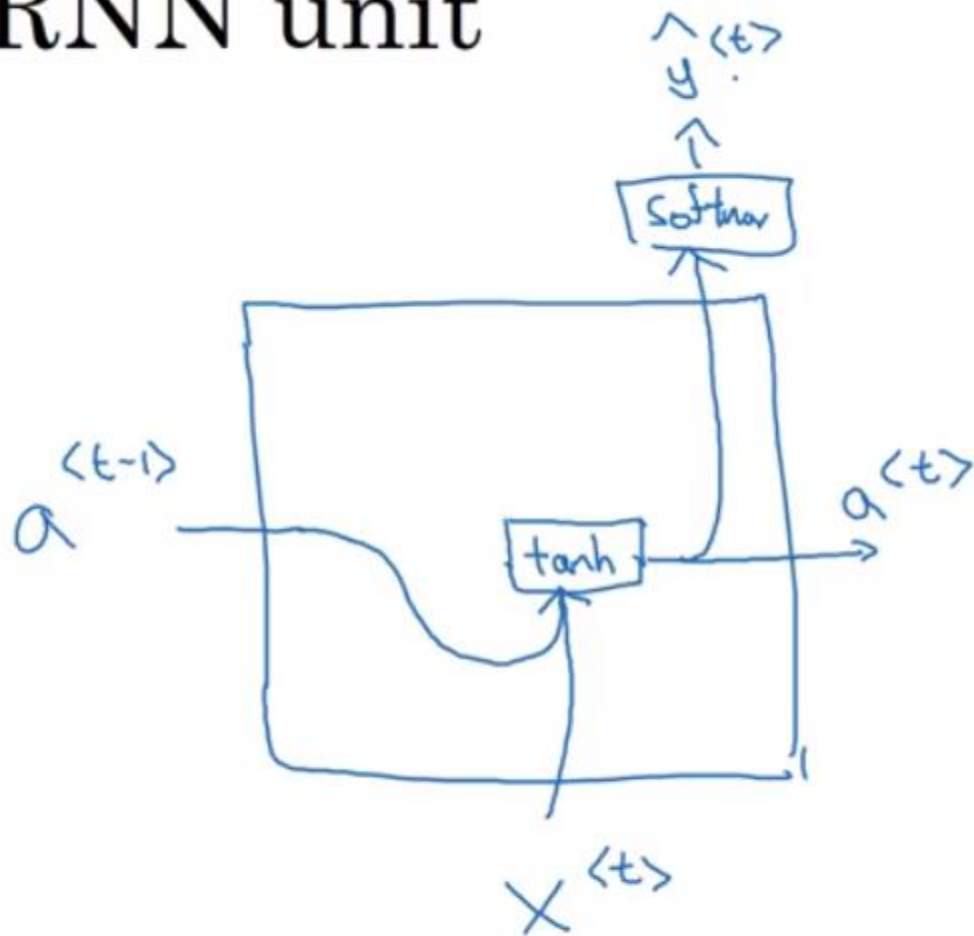


deeplearning.ai

Recurrent Neural Networks

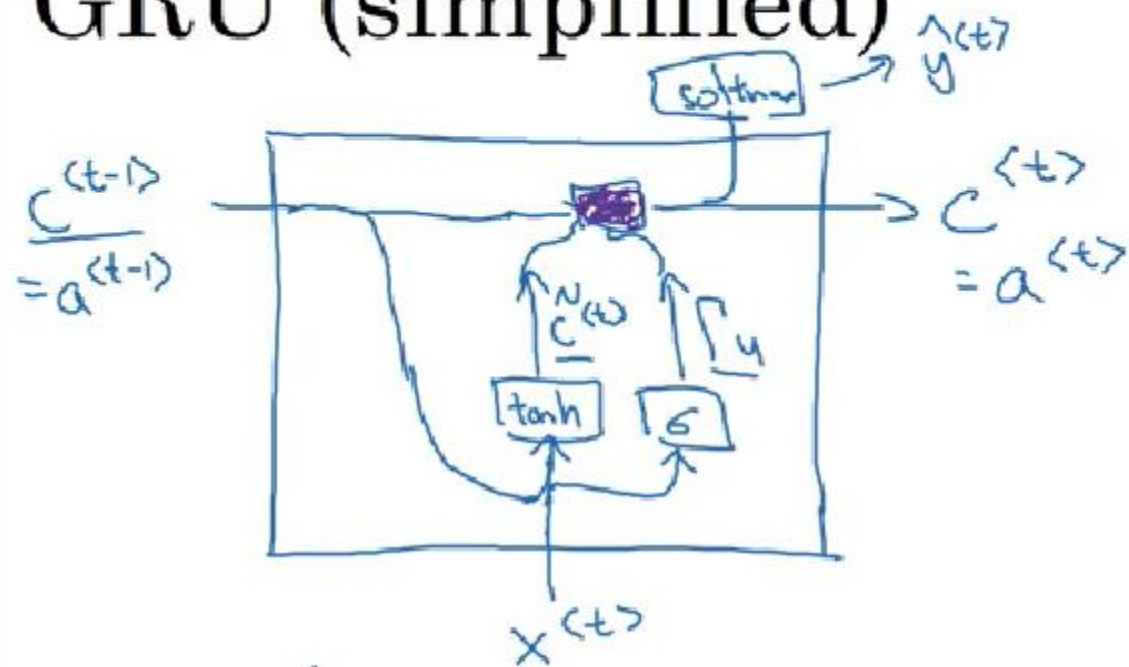
Gated Recurrent Unit (GRU)

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)



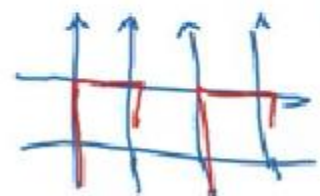
c = memory cell

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

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

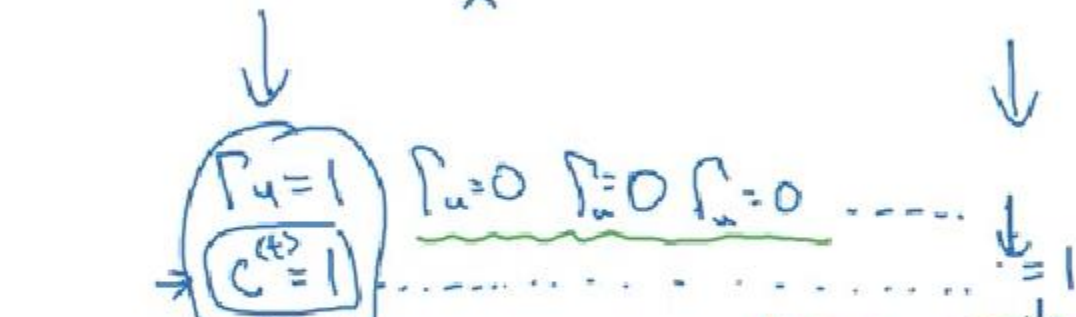
$$\rightarrow \boxed{\Gamma_u} = \sigma(w_u [c^{(t-1)}, x^{(t)}] + b_u)$$

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



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

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

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

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

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

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

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

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

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

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

LSTM

LSTM in pictures

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

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

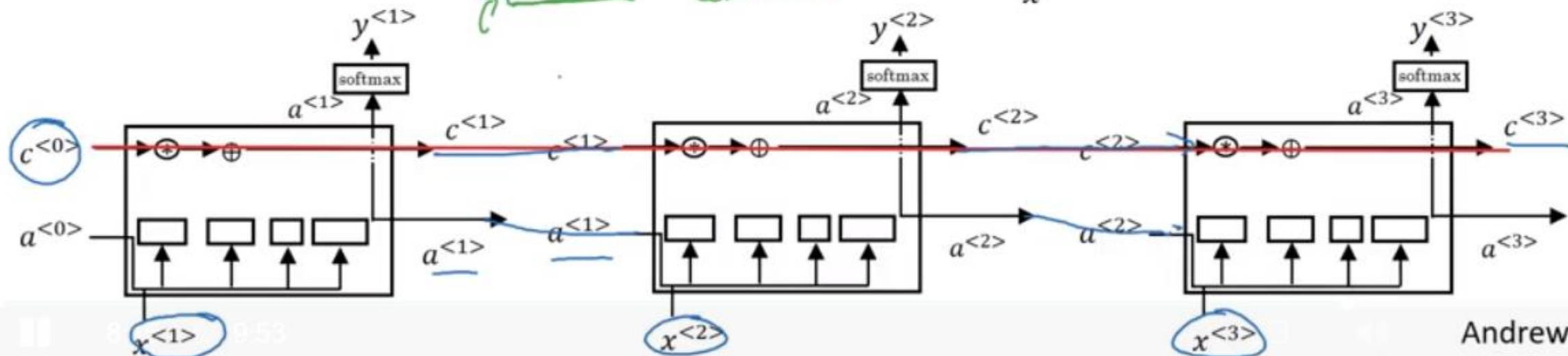
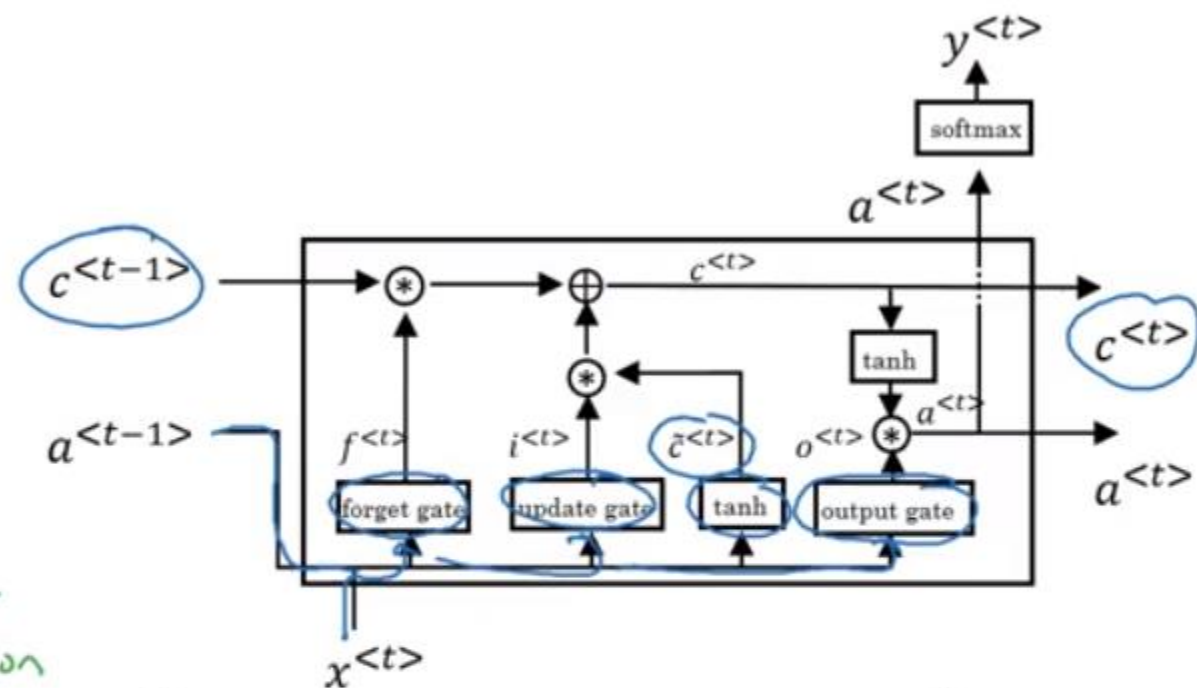
$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

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

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

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

peephole
connection





deeplearning.ai

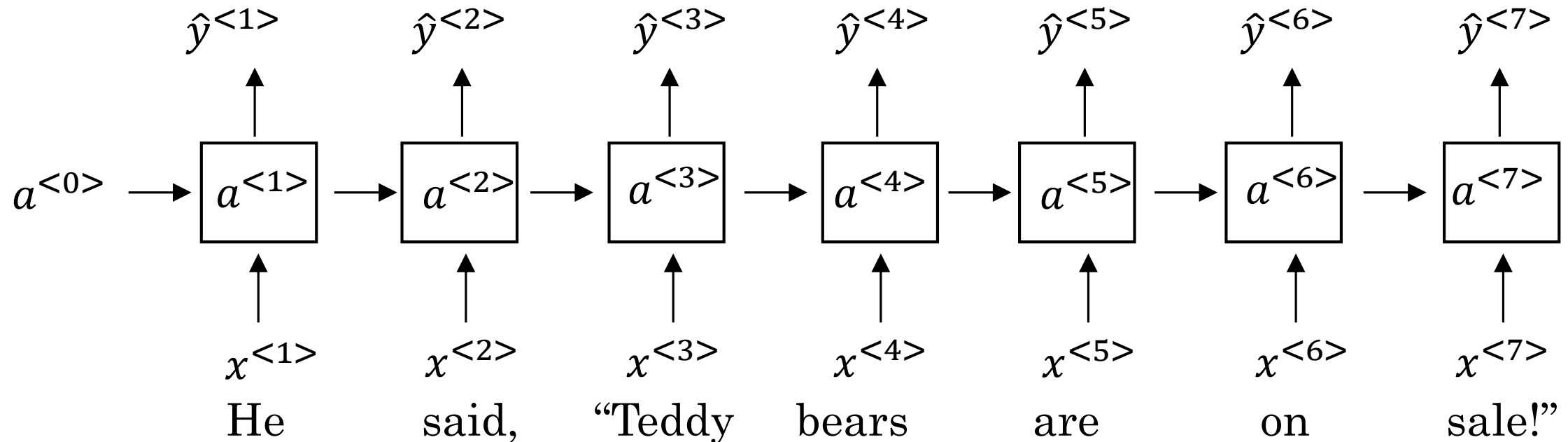
Recurrent Neural Networks

Bidirectional RNN

Getting information from the future

He said, “Teddy bears are on sale!”

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



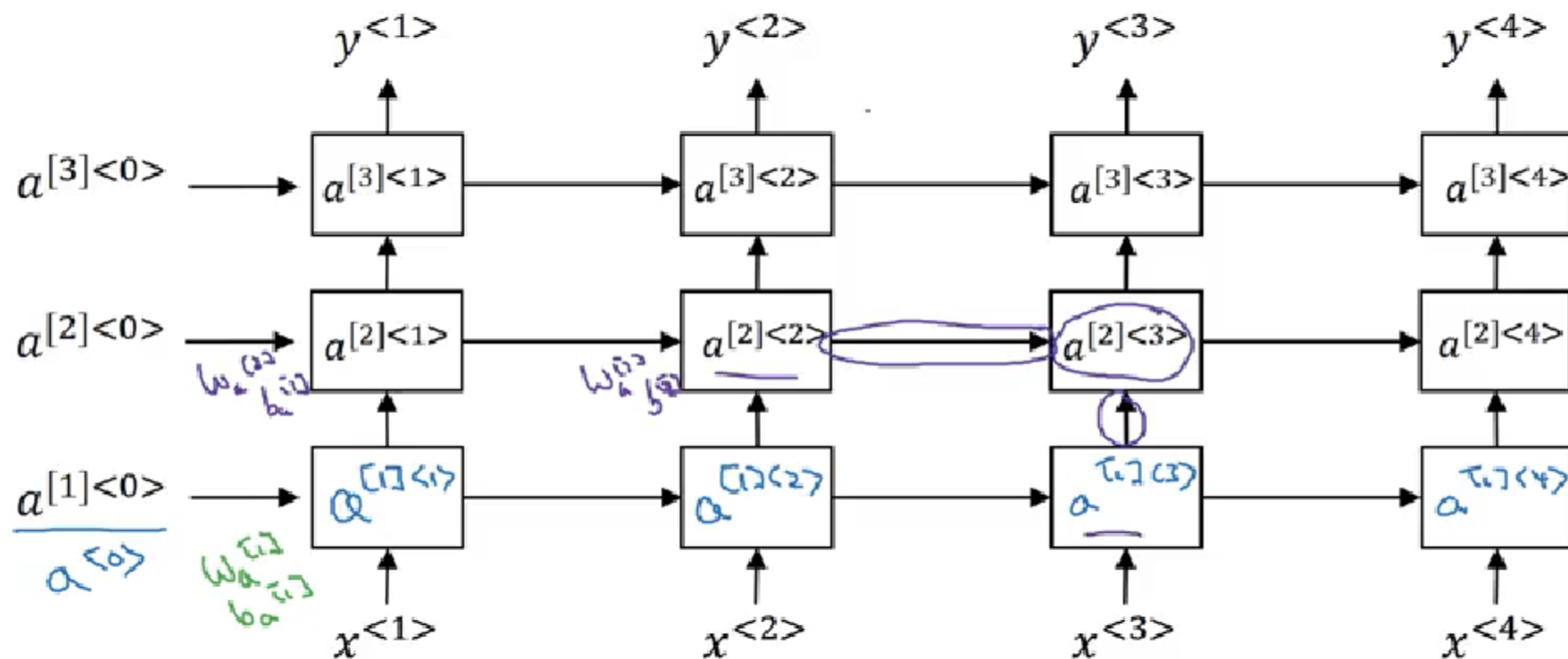
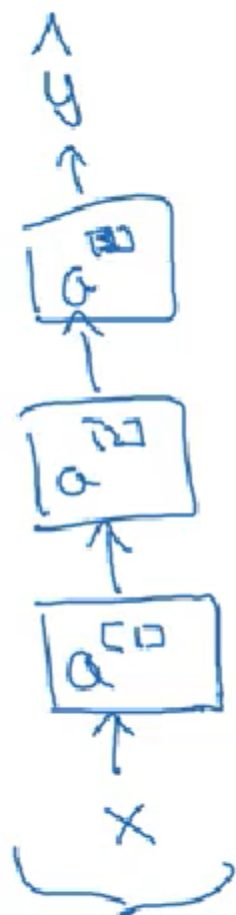
Bidirectional RNN (BRNN)



deeplearning.ai

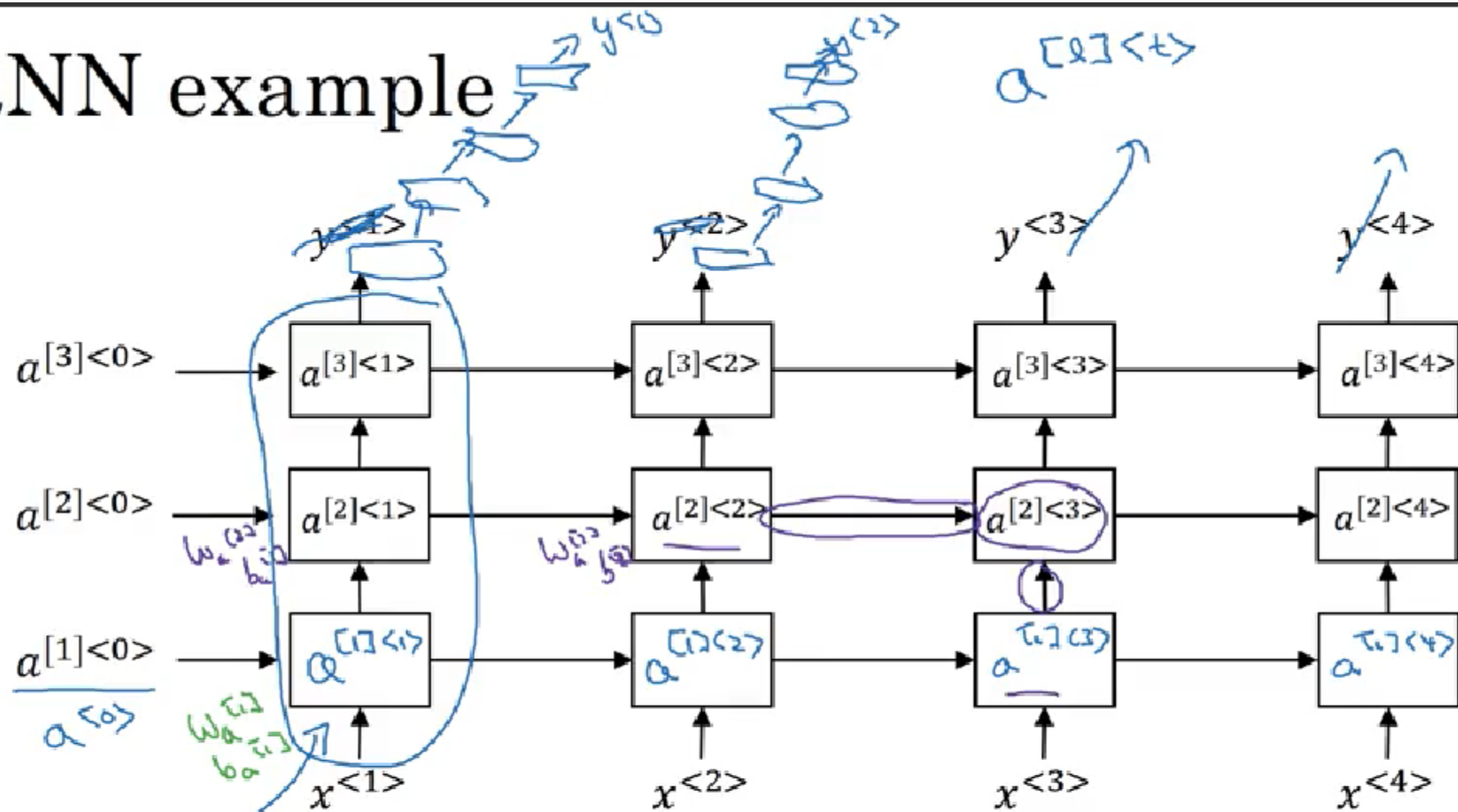
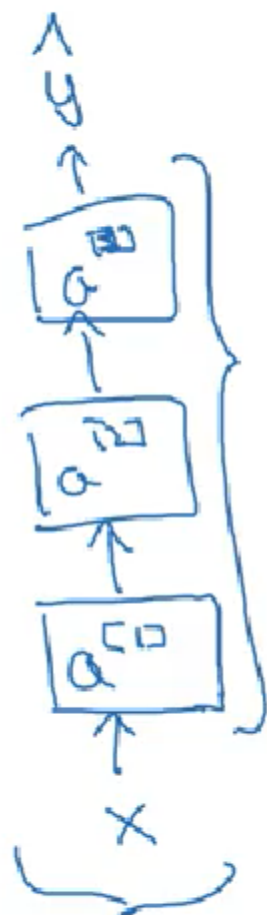
Recurrent Neural Networks

Deep RNNs

$$Q_{[2]} \langle t \rangle$$


$$a^{[12] \langle 3 \rangle} = g (w_a^{[12]} [a^{[12] \langle 2 \rangle}, a^{[1] \langle 3 \rangle}] + b_a^{[1]})$$

Deep RNN example



RNN
GRU
LSTM

BRNN

$$a^{[2]\langle 3 \rangle} = g(w_a^{[2]} [a^{[1]\langle 2 \rangle}, a^{[1]\langle 3 \rangle}] + b_a^{[2]})$$