



温州大學  
WENZHOU UNIVERSITY

# 深度学习-第十一章-序列模型

黄海广 副教授

2021年05月

- 01** 序列模型概述
- 02** 循环神经网络(RNN)
- 03** 长短期记忆(LSTM)
- 04** 双向循环神经网络
- 05** 深层循环神经网络

# 1.序列模型概述

3

## 01 序列模型概述

**02** 循环神经网络(RNN)

**03** 长短期记忆(LSTM)

**04** 双向循环神经网络

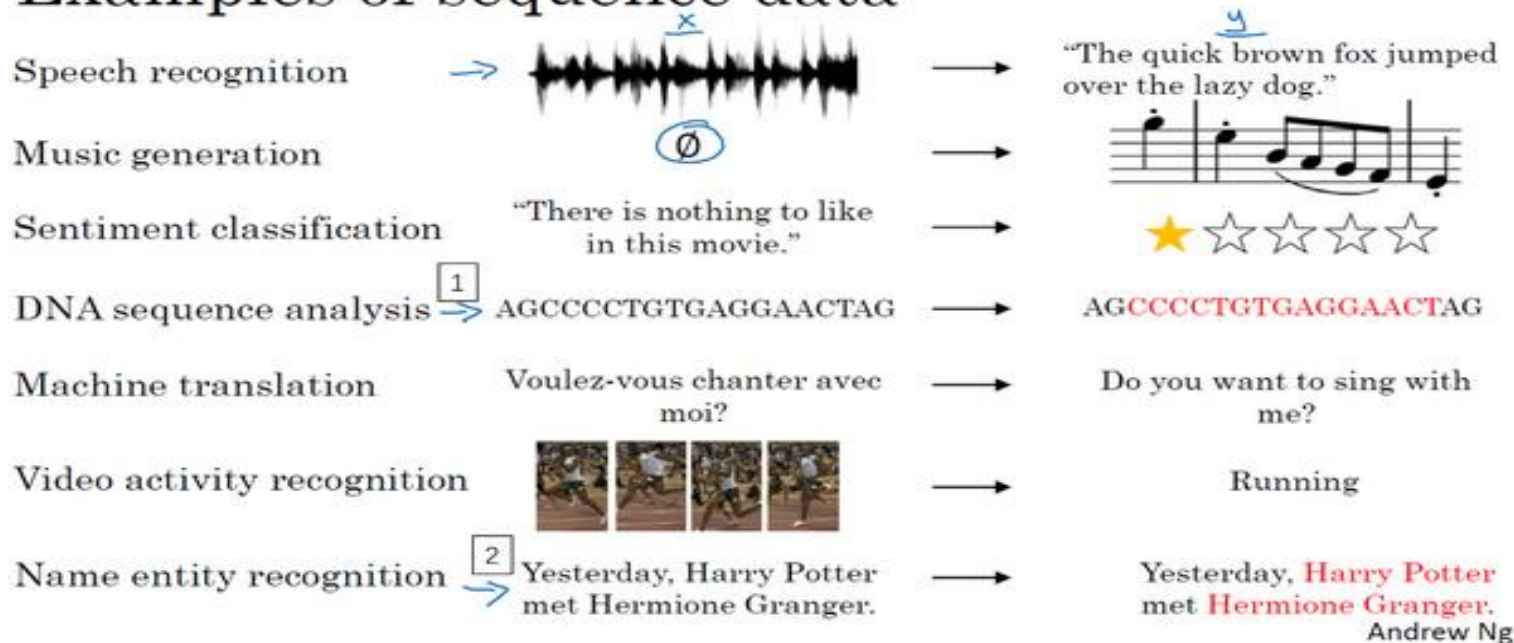
**05** 深层循环神经网络

# 1.序列模型概述

4

循环神经网络（**RNN**）之类的模型在语音识别、自然语言处理和其他领域中引起变革。

## Examples of sequence data



# 数学符号

5

在这里 $x^{<1>}$ 表示**Harry**这个单词，它就是一个第4075行是1，其余值都是0的向量（上图编号1所示），因为那是**Harry**在这个词典里的位置。

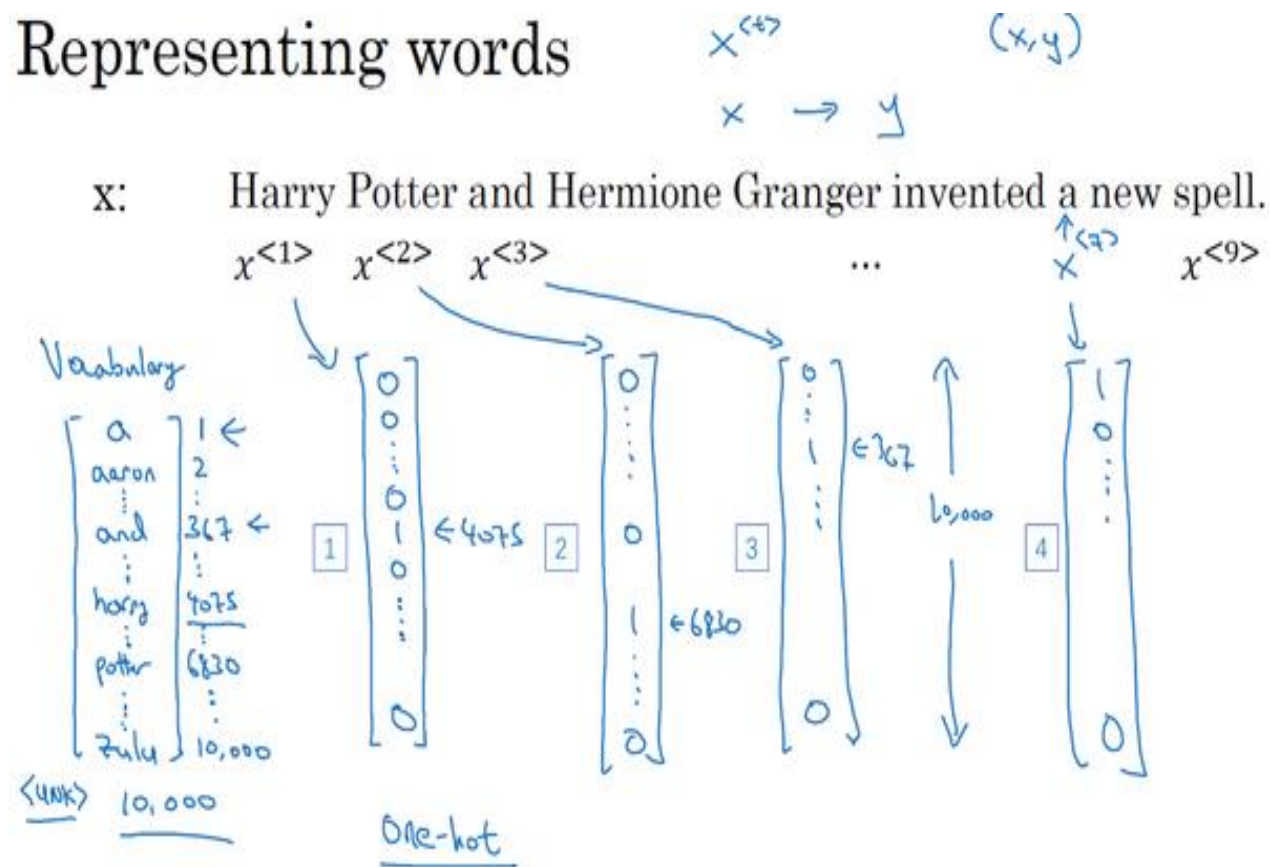
$x^{<2>}$ 是第6830行是1，其余位置都是0的向量（上图编号2所示）。

**and**在词典里排第367，所以 $x^{<3>}$ 就是第367行是1，其余值都是0的向量（上图编号3所示）。

因为**a**是字典第一个单词， $x^{<7>}$ 对应**a**，那么这个向量的第一个位置为1，其余位置都是0的向量（上图编号4所示）。

**Unknow Word**的伪单词，用**<UNK>**作为标记。

## Representing words



如果你的词典大小是10,000的话，那么这里的每个向量都是10,000维的。



## 2.循环神经网络(RNN)

6

**01** 序列模型概述

**02** 循环神经网络(RNN)

**03** 长短期记忆(LSTM)

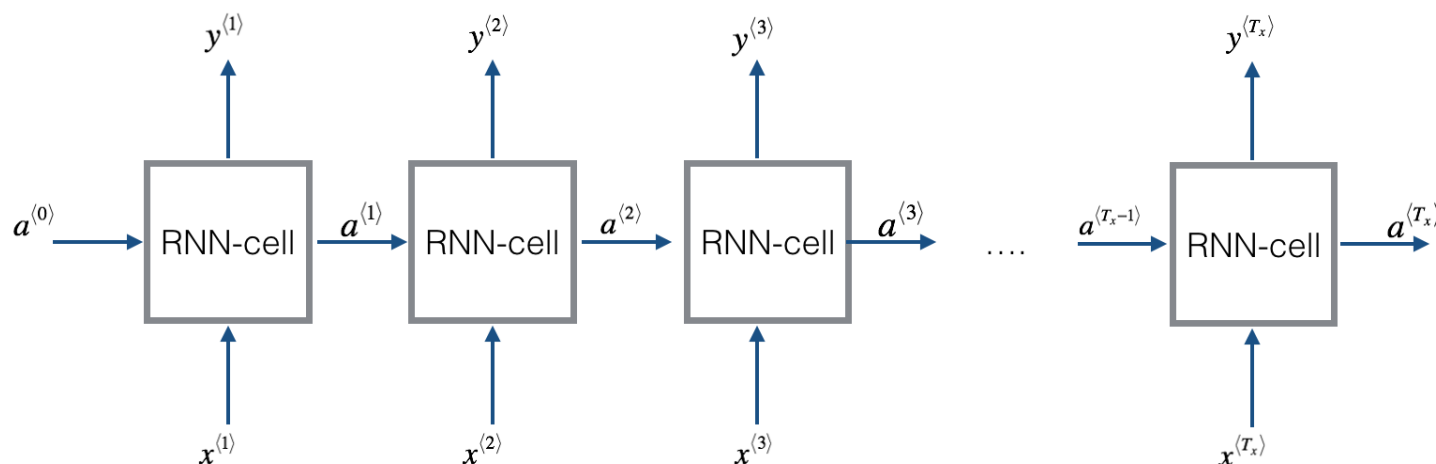
**04** 双向循环神经网络

**05** 深层循环神经网络

## 2.循环神经网络(RNN)

7

### RNN的前向传播



$$a^{<0>} = 0$$

$$a^{<1>} = g_1(W_{aa}a^{<0>} + W_{ax}x^{<1>} + b_a)$$

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

$$a^{<t>} = g_1(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g_2(W_{ya}a^{<t>} + b_y)$$

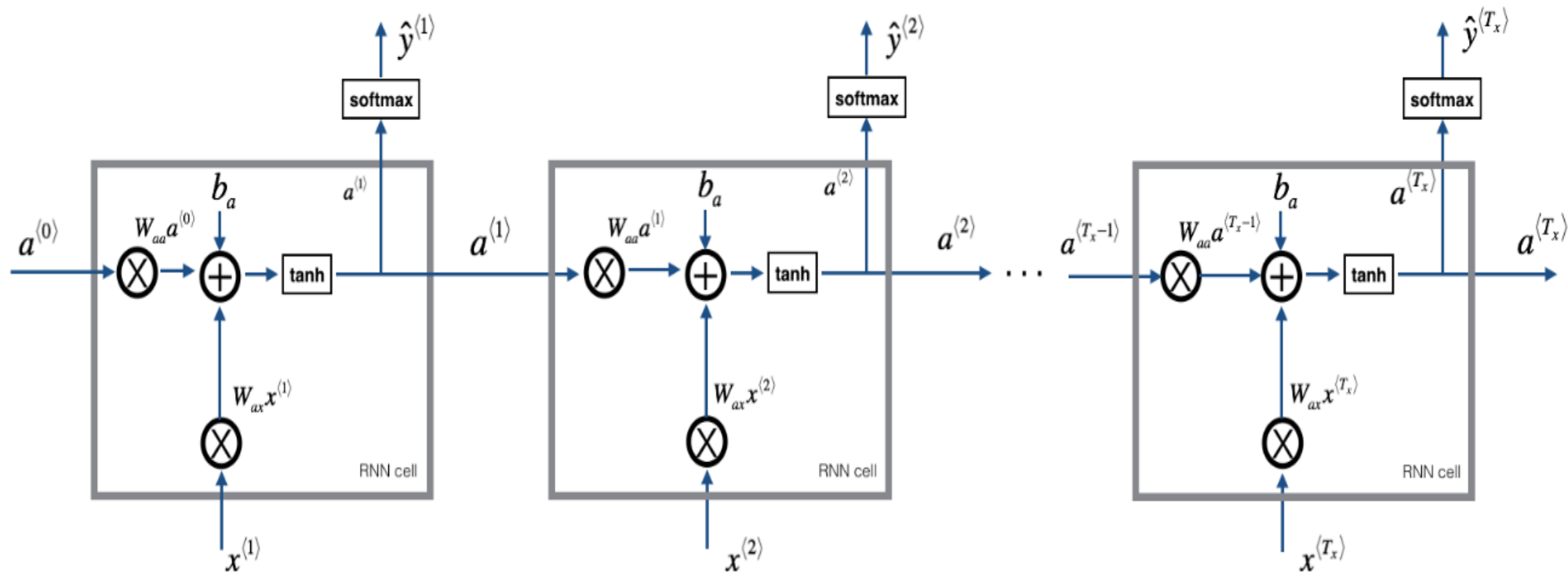
## 2.循环神经网络(RNN)

8

RNN的前向传播

$$a^{<t>} = g_1(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g_2(W_{ya}a^{<t>} + b_y)$$

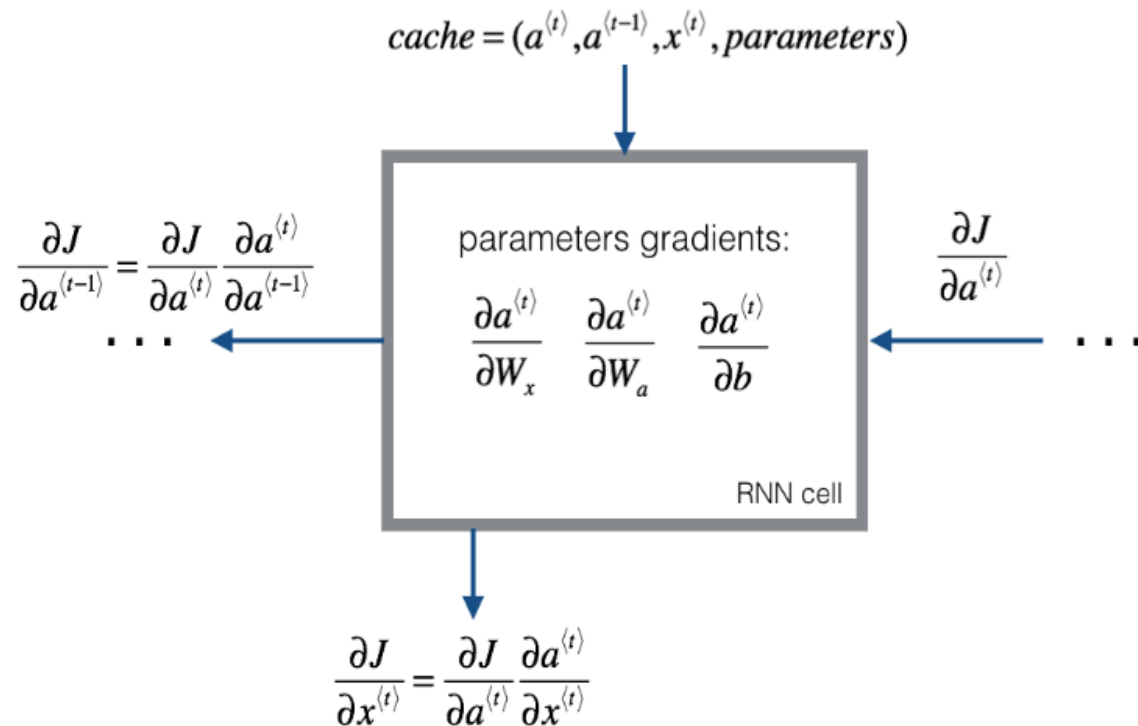




## 2.循环神经网络(RNN)

9

### RNN的反向传播



$$a^{(t)} = \tanh(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b)$$

$$\frac{\partial \tanh(x)}{\partial x} = 1 - \tanh(x)^2$$

$$\frac{\partial a^{(t)}}{\partial W_{ax}} = (1 - \tanh(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b)^2) x^{(t)T}$$

$$\frac{\partial a^{(t)}}{\partial W_{aa}} = (1 - \tanh(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b)^2) a^{(t-1)T}$$

$$\frac{\partial a^{(t)}}{\partial b} = \sum_{batch} (1 - \tanh(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b)^2)$$

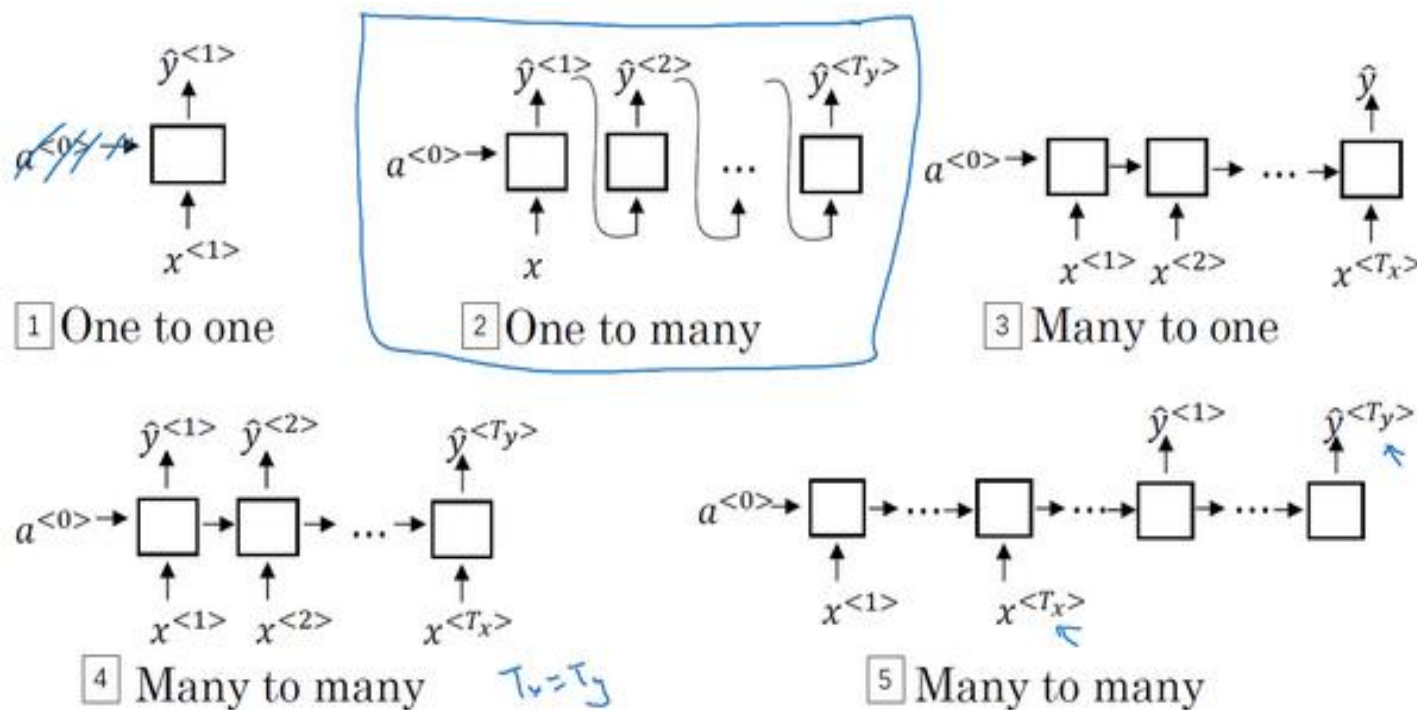
$$\frac{\partial a^{(t)}}{\partial x^{(t)}} = W_{ax}^T \cdot (1 - \tanh(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b)^2)$$

$$\frac{\partial a^{(t)}}{\partial a^{(t-1)}} = W_{aa}^T \cdot (1 - \tanh(W_{ax}x^{(t-1)} + W_{aa}a^{(t-1)} + b)^2)$$

## 2.循环神经网络(RNN)

10

### RNN的类型



## 2.循环神经网络(RNN)

11

语言模型和序列生成

**The apple and pear (pair) salad was delicious.**

第一句话的概率是:

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

而第二句话的概率是:

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

# 3.长短期记忆(LSTM)

12

**01** 序列模型概述

**02** 循环神经网络(RNN)

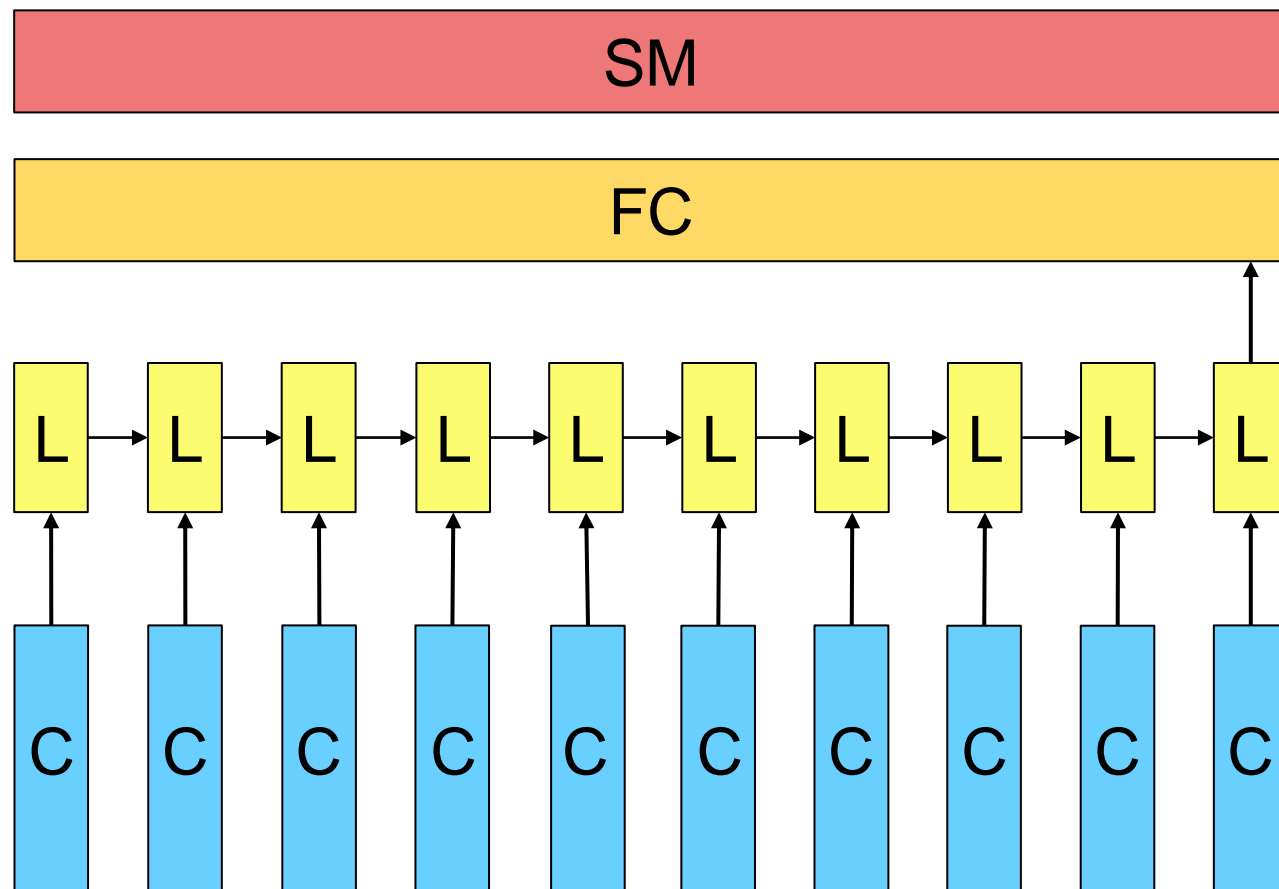
**03** 长短期记忆(LSTM)

**04** 双向循环神经网络

**05** 深层循环神经网络

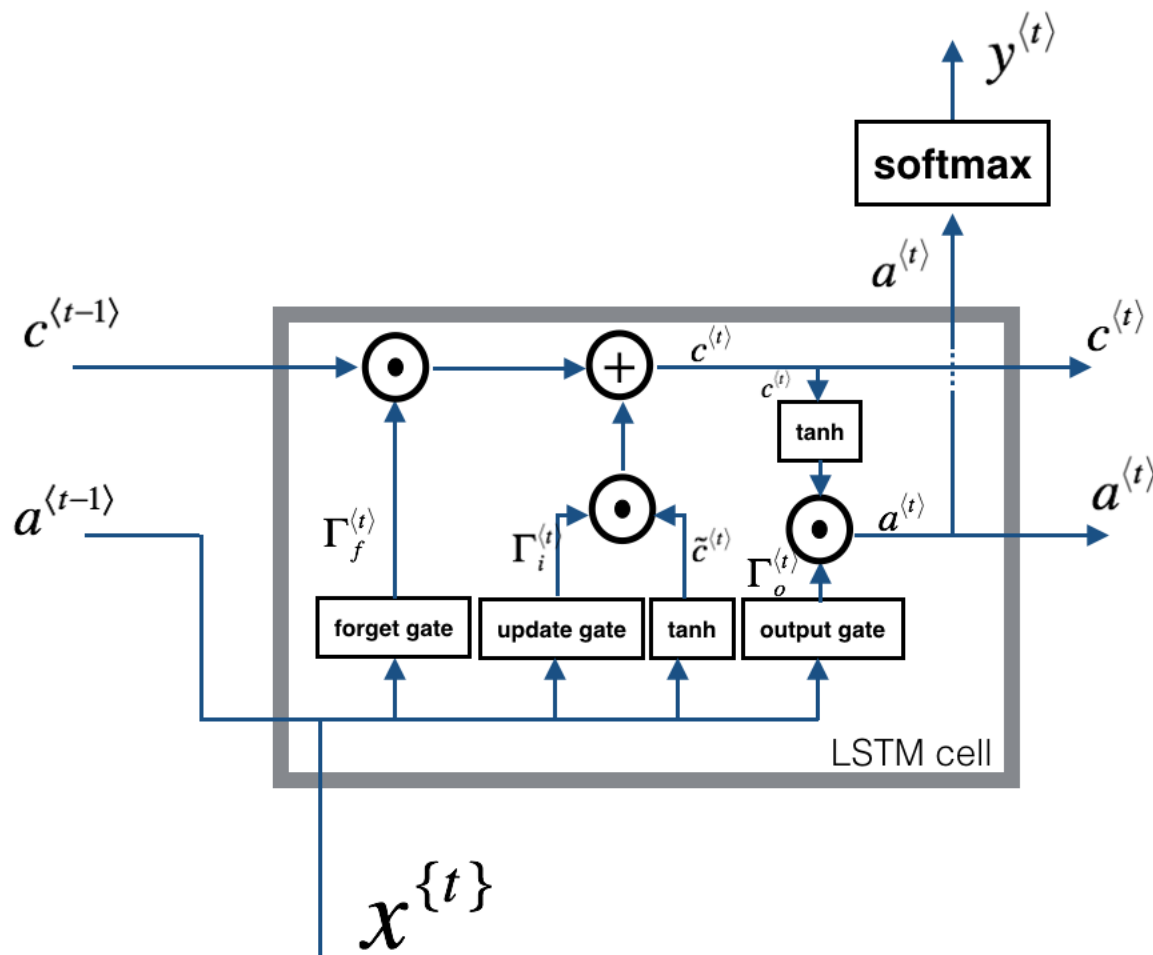
### 3.长短期记忆(LSTM)

13



### 3.长短期记忆(LSTM)

14



...

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

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

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

$$c^{\langle t \rangle} = \Gamma_f^{\langle t \rangle} \circ c^{\langle t-1 \rangle} + \Gamma_u^{\langle t \rangle} \circ \tilde{c}^{\langle t \rangle}$$

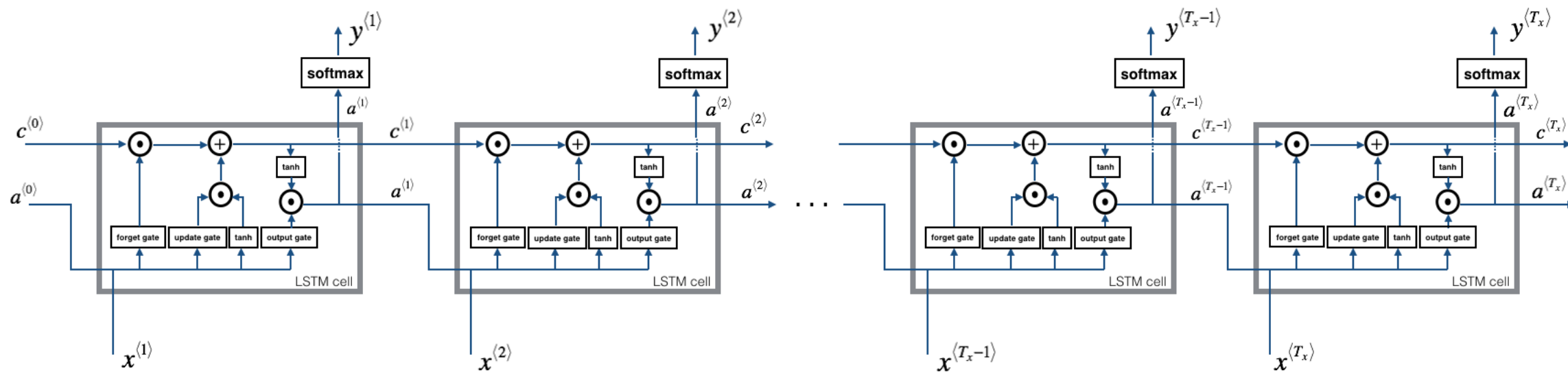
$$\Gamma_o^{\langle t \rangle} = \sigma(W_o[a^{\langle t-1 \rangle}, x^{\langle t \rangle}] + b_o)$$

$$a^{\langle t \rangle} = \Gamma_o^{\langle t \rangle} \circ \tanh(c^{\langle t \rangle})$$

# 3.长短期记忆(LSTM)

15

## LSTM的前向传播





### 3.长短期记忆(LSTM)

16

#### LSTM的反向传播

$$d\Gamma_o^{\langle t \rangle} = da_{next} * \tanh(c_{next}) * \Gamma_o^{\langle t \rangle} * (1 - \Gamma_o^{\langle t \rangle})$$

$$d\tilde{c}^{\langle t \rangle} = dc_{next} * \Gamma_i^{\langle t \rangle} + \Gamma_o^{\langle t \rangle} (1 - \tanh(c_{next})^2) * i_t * da_{next} * \tilde{c}^{\langle t \rangle} * (1 - \tanh(\tilde{c})^2)$$

$$d\Gamma_u^{\langle t \rangle} = dc_{next} * \tilde{c}^{\langle t \rangle} + \Gamma_o^{\langle t \rangle} (1 - \tanh(c_{next})^2) * \tilde{c}^{\langle t \rangle} * da_{next} * \Gamma_u^{\langle t \rangle} * (1 - \Gamma_u^{\langle t \rangle})$$

$$d\Gamma_f^{\langle t \rangle} = dc_{next} * \tilde{c}_{prev} + \Gamma_o^{\langle t \rangle} (1 - \tanh(c_{next})^2) * c_{prev} * da_{next} * \Gamma_f^{\langle t \rangle} * (1 - \Gamma_f^{\langle t \rangle})$$

# 3.长短期记忆(LSTM)

17

## LSTM的反向传播

### parameter derivatives

$$dW_f = d\Gamma_f^{\langle t \rangle} * \begin{pmatrix} a_{prev} \\ x_t \end{pmatrix}^T$$

$$dW_c = d\tilde{c}^{\langle t \rangle} * \begin{pmatrix} a_{prev} \\ x_t \end{pmatrix}^T$$

$$dW_o = d\Gamma_o^{\langle t \rangle} * \begin{pmatrix} a_{prev} \\ x_t \end{pmatrix}^T$$

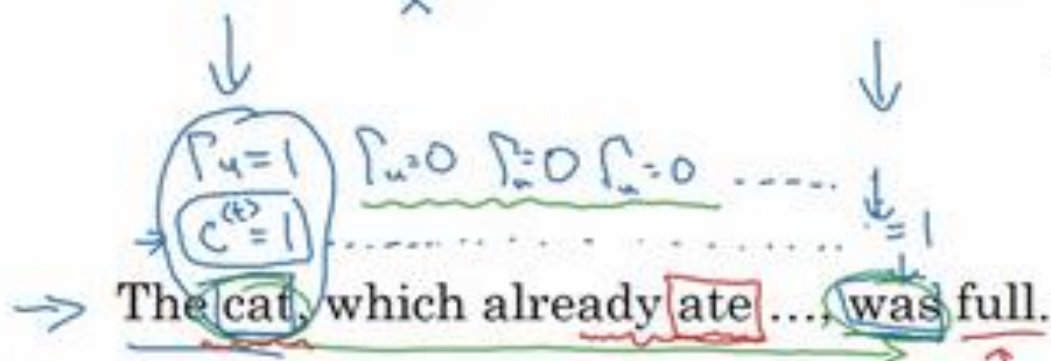
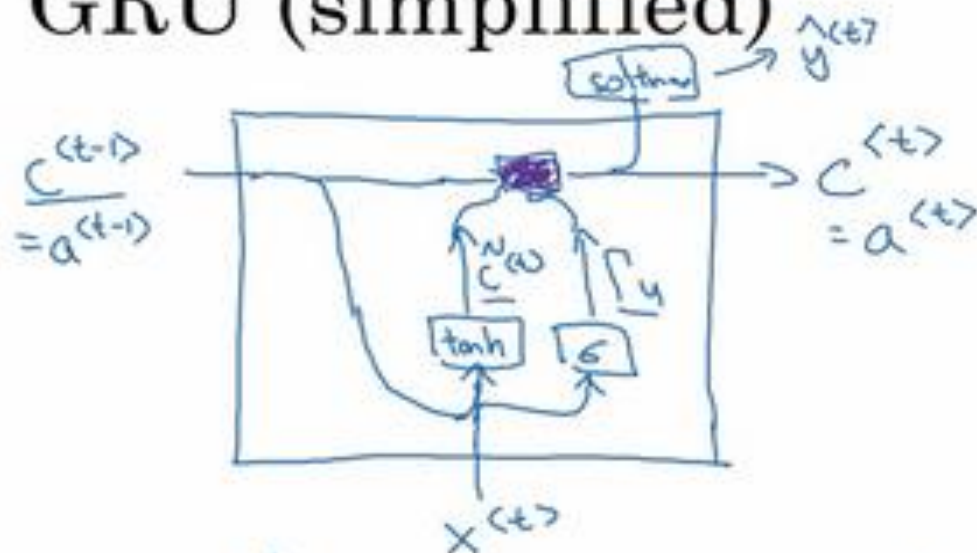
To calculate  $db_f, db_u, db_c, db_o$  you just need to sum across the horizontal (axis= 1) axis on  $d\Gamma_f^{\langle t \rangle}, d\Gamma_u^{\langle t \rangle}, d\tilde{c}^{\langle t \rangle}, d\Gamma_o^{\langle t \rangle}$  respectively. Note that you should have the `keep_dims = True` option.

Finally, you will compute the derivative with respect to the previous hidden state.

# 3.长短期记忆(LSTM)

18

## GRU (simplified)



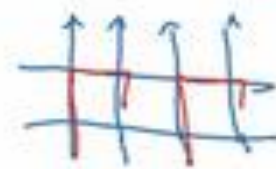
$C$  = memory cell

$$\textcircled{1} \rightarrow \underline{C}^{(t)} = \underline{a}^{(t)}$$

$$\textcircled{2} \rightarrow \tilde{C}^{(t)} = \tanh(W_c [C^{(t-1)}, x^{(t)}] + b_c)$$

$$\textcircled{3} \rightarrow \Gamma_u = \sigma(W_u [C^{(t-1)}, x^{(t)}] + b_u)$$

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



element-wise  $\textcircled{4}$   
Gate  
 $\Gamma_u = 0.00001$

[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

## 2.循环神经网络(RNN)

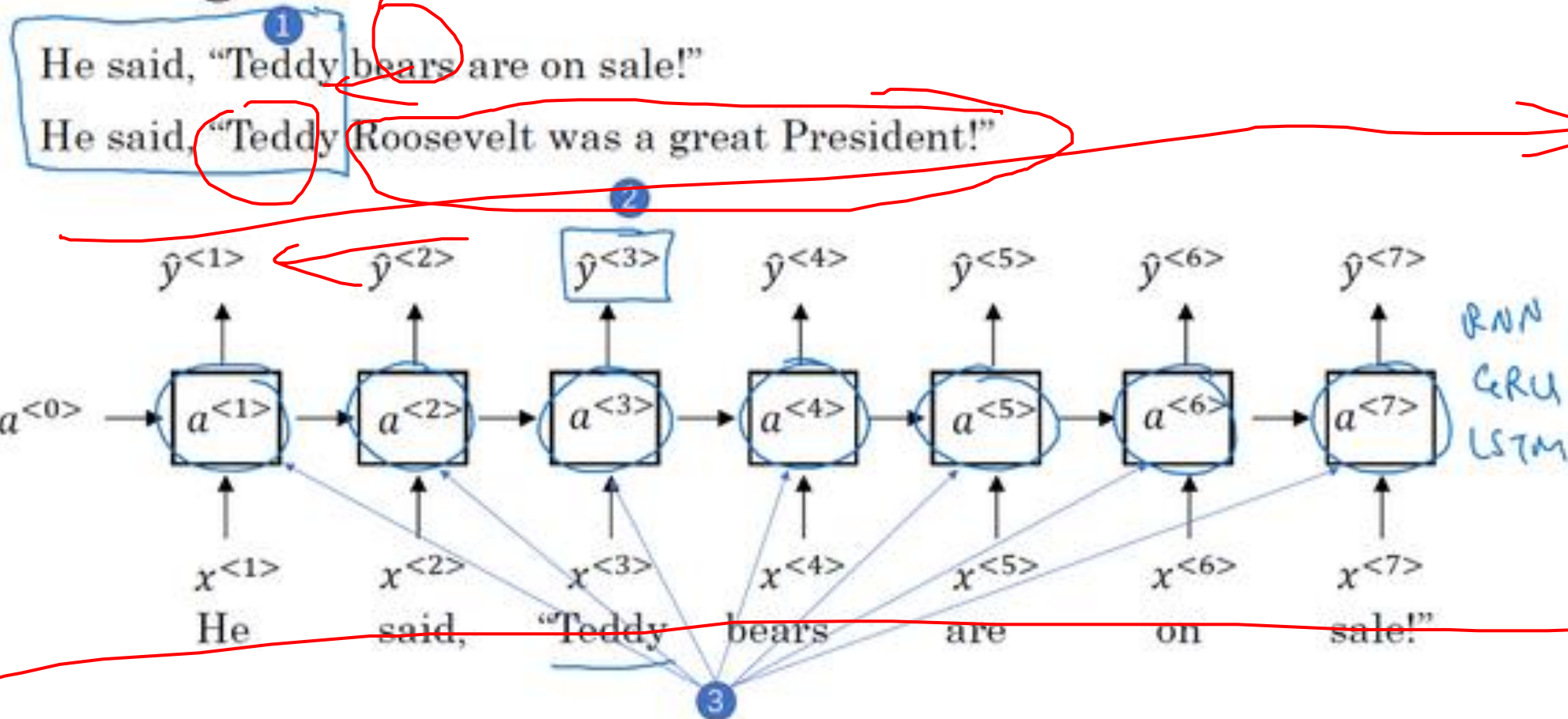
19

- 01** 序列模型概述
- 02** 循环神经网络(RNN)
- 03** 长短期记忆(LSTM)
- 04** 双向循环神经网络
- 05** 深层循环神经网络

# 4.双向循环神经网络

20

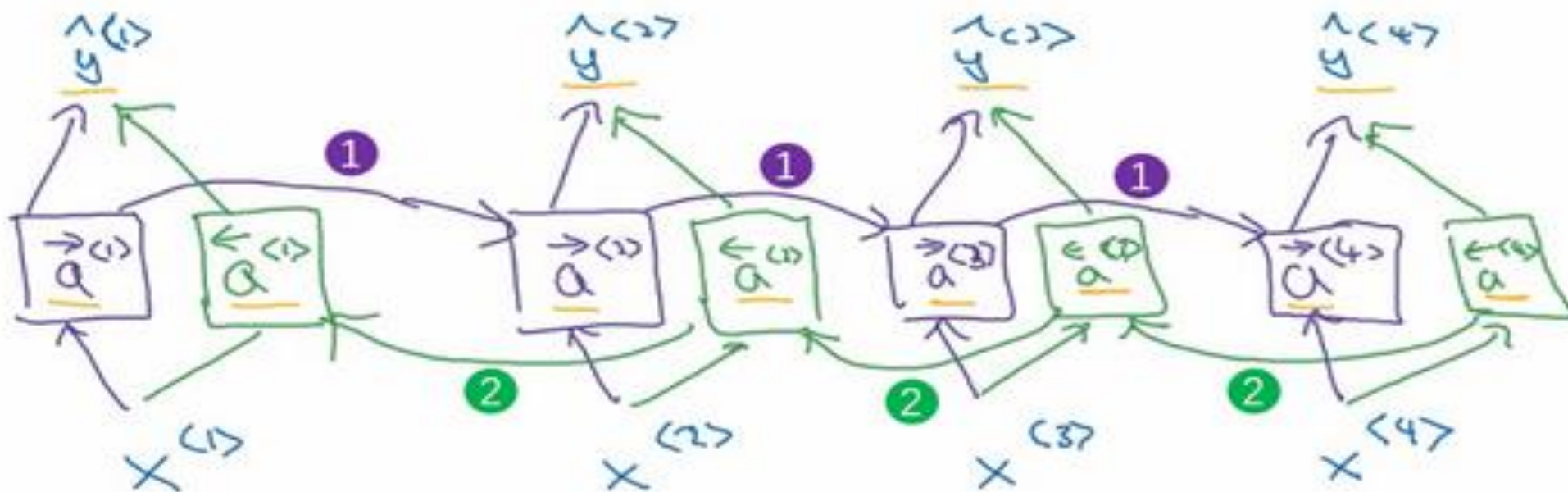
Getting information from the future



## 4.双向循环神经网络

21

### Bidirectional RNN (BRNN)



Acyclic graph

# 5. 深层循环神经网络

22

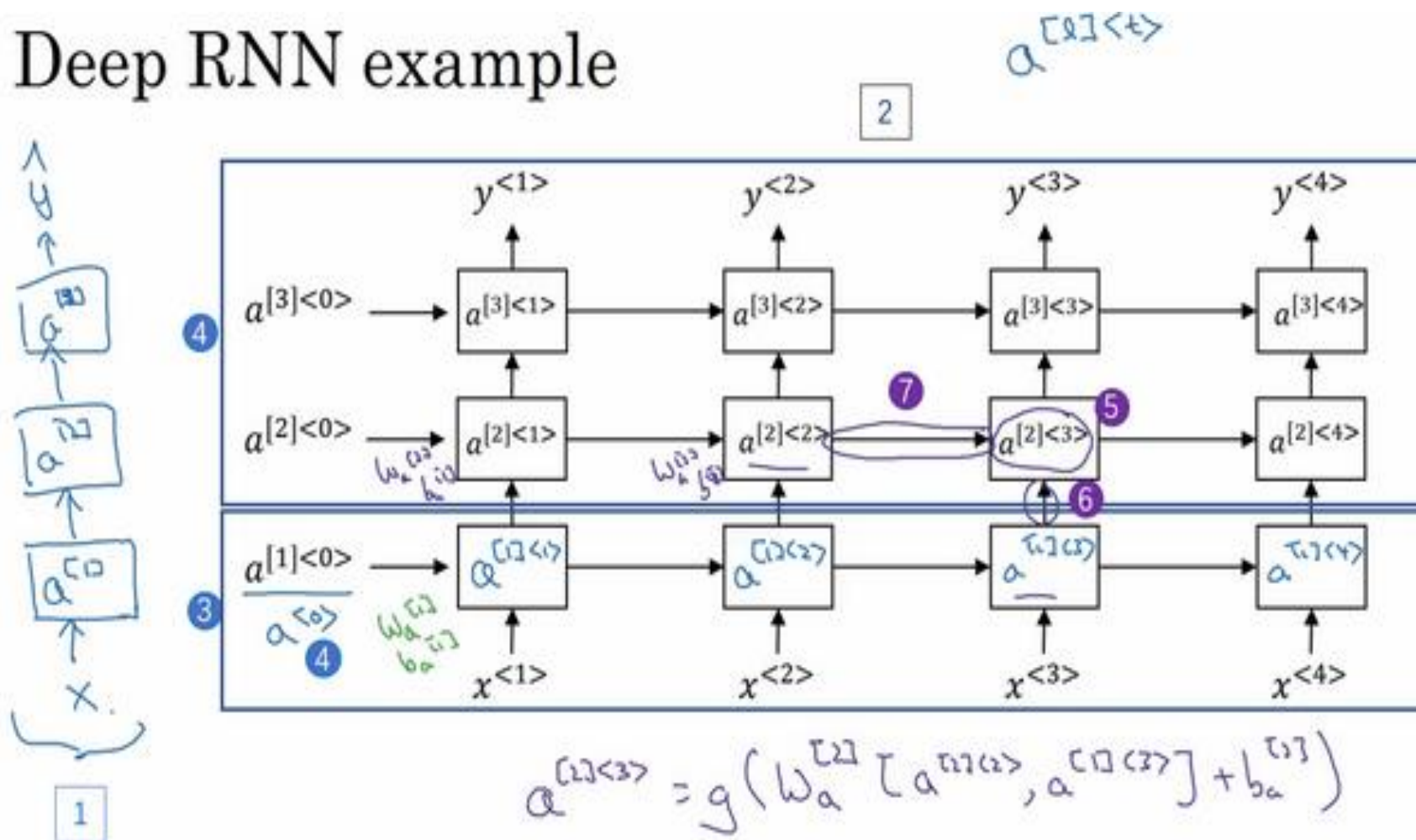
- 01** 序列模型概述
- 02** 循环神经网络(RNN)
- 03** 长短期记忆(LSTM)
- 04** 双向循环神经网络
- 05** 深层循环神经网络



## 5. 深层循环神经网络

23

## Deep RNN example



1. IAN GOODFELLOW等, 《深度学习》, 人民邮电出版社, 2017
2. Andrew Ng, <http://www.deeplearning.ai>

谢谢!

