

# Recurrent Neural Network

# Examples of sequence data

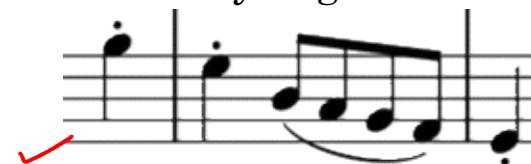
Speech recognition



T<sub>y</sub>  
“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 ✓

AGCCCCTGTGAGGAACTAG

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**.

$\underline{1000}$   $X^{(1)}$  : Harry Dotters and Hermione Granger invented a new spell.  
 $X^{(1)}$   $X^{(2)}$   $X^{(3)}$   $X^{(4)}$   $X^{(5)}$   $X^{(6)}$   $X^{(7)}$   $X^{(8)}$   $X^{(9)}$

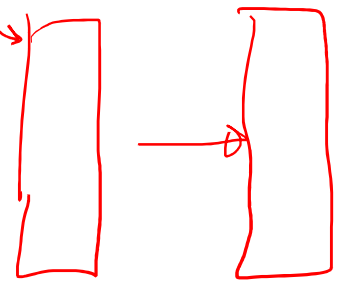
$\underline{10,000}$   $Y^{(1)}$  :  
 $Y^{(1)}$   $Y^{(2)}$   $Y^{(3)}$   $Y^{(4)}$   $Y^{(5)}$   $Y^{(6)}$   $Y^{(7)}$   $Y^{(8)}$   $Y^{(9)}$

$\begin{bmatrix} a \\ \vdots \\ \text{and} \\ \vdots \\ \text{Harry} \\ \vdots \\ \text{Zwiz} \\ \text{unk} \end{bmatrix}$   
 1  
 367  
 4075  
 10,000  
 10001

$X^{(i)} \langle t \rangle$   $T_X^{(i)}$   
 $Y^{(i)} \langle t \rangle$   $T_Y^{(i)}$

$X^{(1)}$   
 $X^{(2)}$   
 $X^{(3)}$   
 $X^{(4)}$   
 $X^{(5)}$   
 $X^{(6)}$   
 $X^{(7)}$   
 $X^{(8)}$   
 $X^{(9)}$

$X^{(1)}$   
 $X^{(2)}$   
 $X^{(3)}$   
 $X^{(4)}$   
 $X^{(5)}$   
 $X^{(6)}$   
 $X^{(7)}$   
 $X^{(8)}$   
 $X^{(9)}$



$Y^{(1)}$   
 $Y^{(2)}$   
 $Y^{(3)}$   
 $Y^{(4)}$   
 $Y^{(5)}$   
 $Y^{(6)}$   
 $Y^{(7)}$   
 $Y^{(8)}$   
 $Y^{(9)}$

# Motivating example

NLP

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

→  $x^{(1)}$   $x^{(2)}$   $x^{(3)}$  ...  $x^{(t)}$  ...  $x^{(9)}$   
 $T_x = 9$

→ y: 1 1 0 1 1 0 0 0 0  
 $y^{(1)}$   $y^{(2)}$   $y^{(3)}$  ...  $y^{(9)}$   
 $T_y = 9$

$x^{(i)(t)}$   $T_x^{(i)} = 9$  15  
 $y^{(i)(t)}$   $T_y^{(i)}$



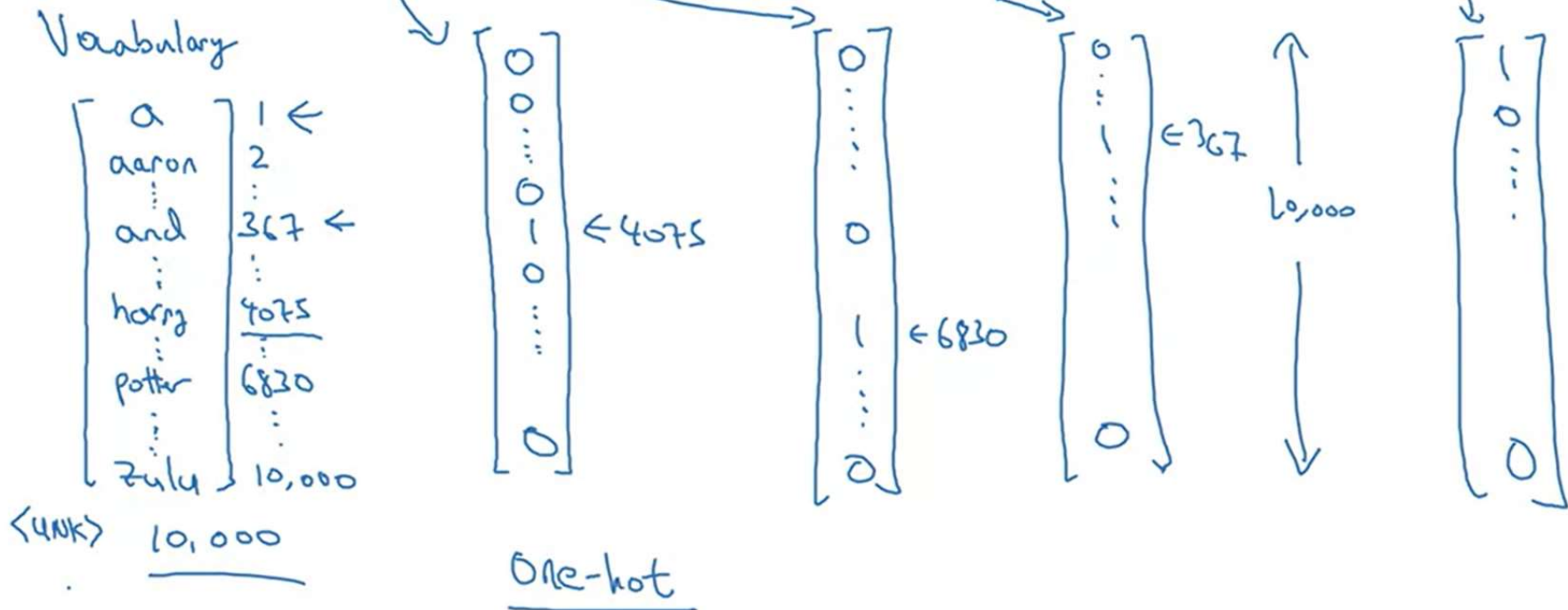
# Representing words

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

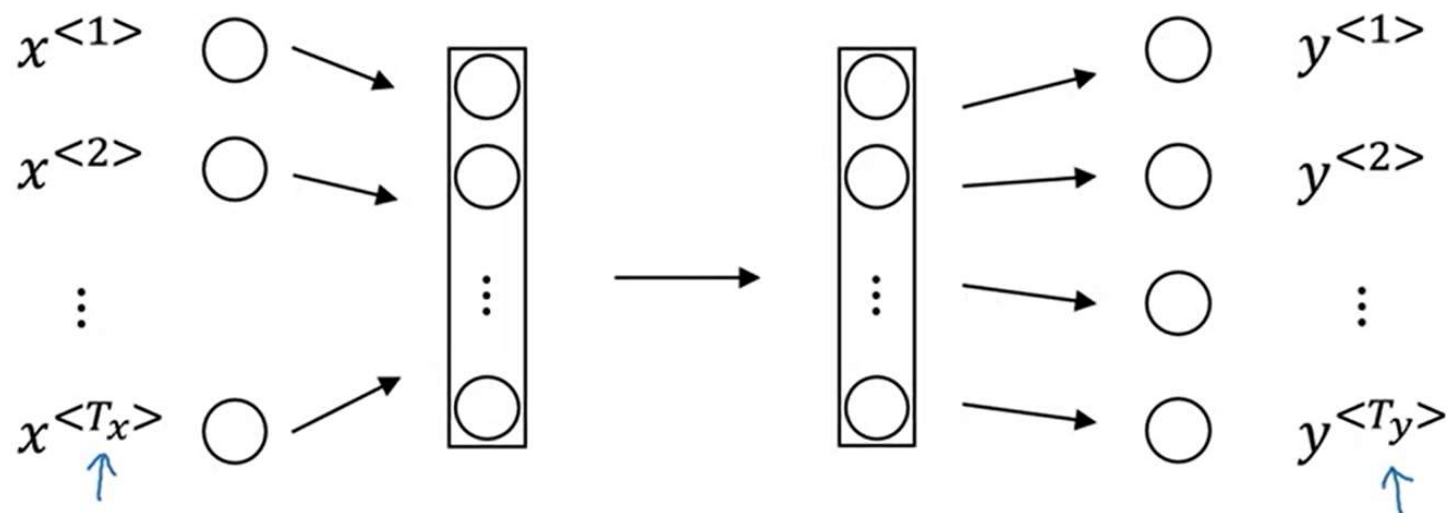
$$(x, y)$$

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

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



# Why not a standard network?



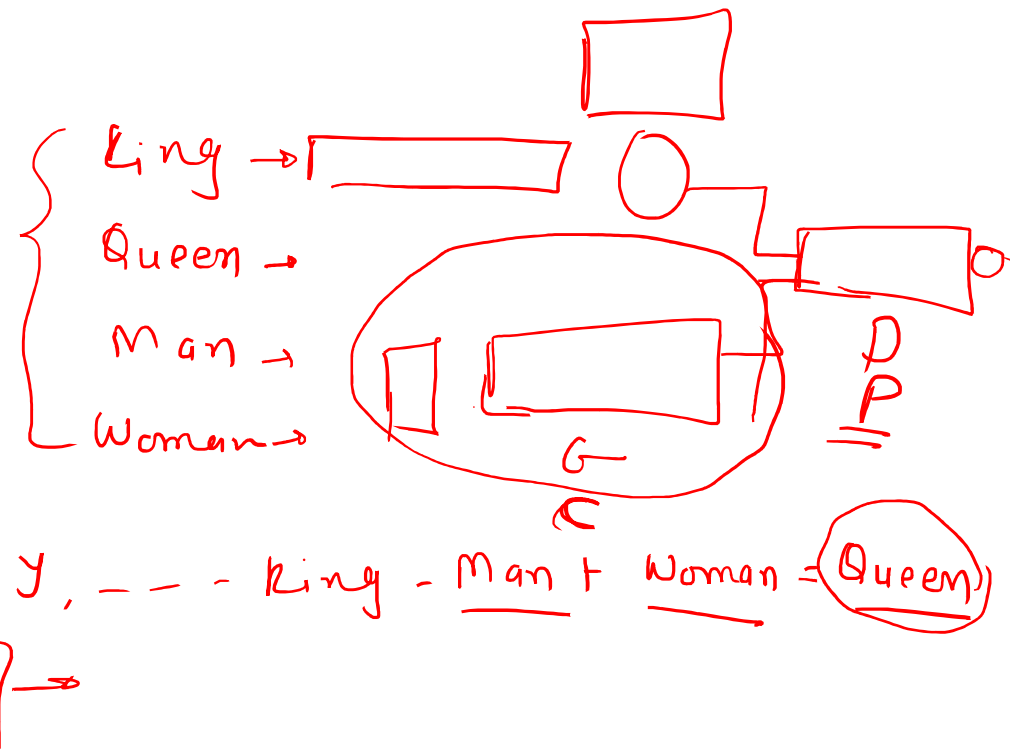
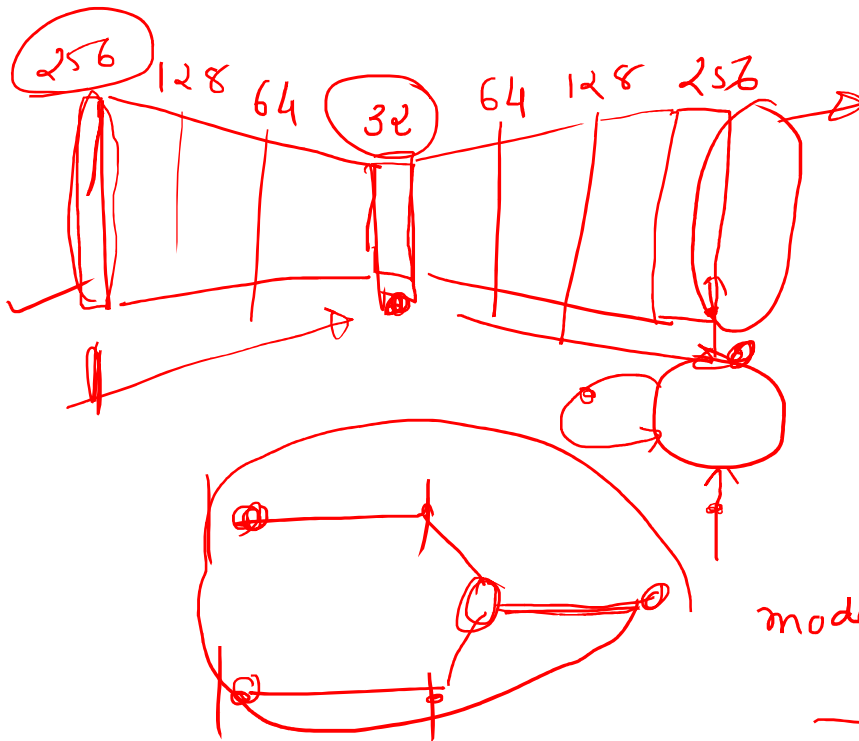
## Problems:

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

# Why not standard network?

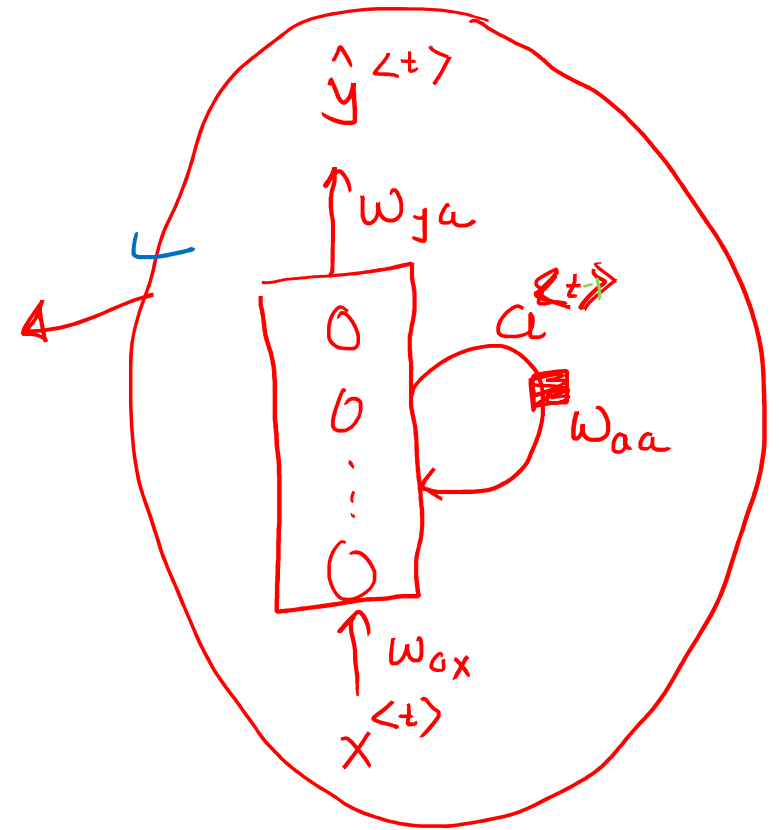
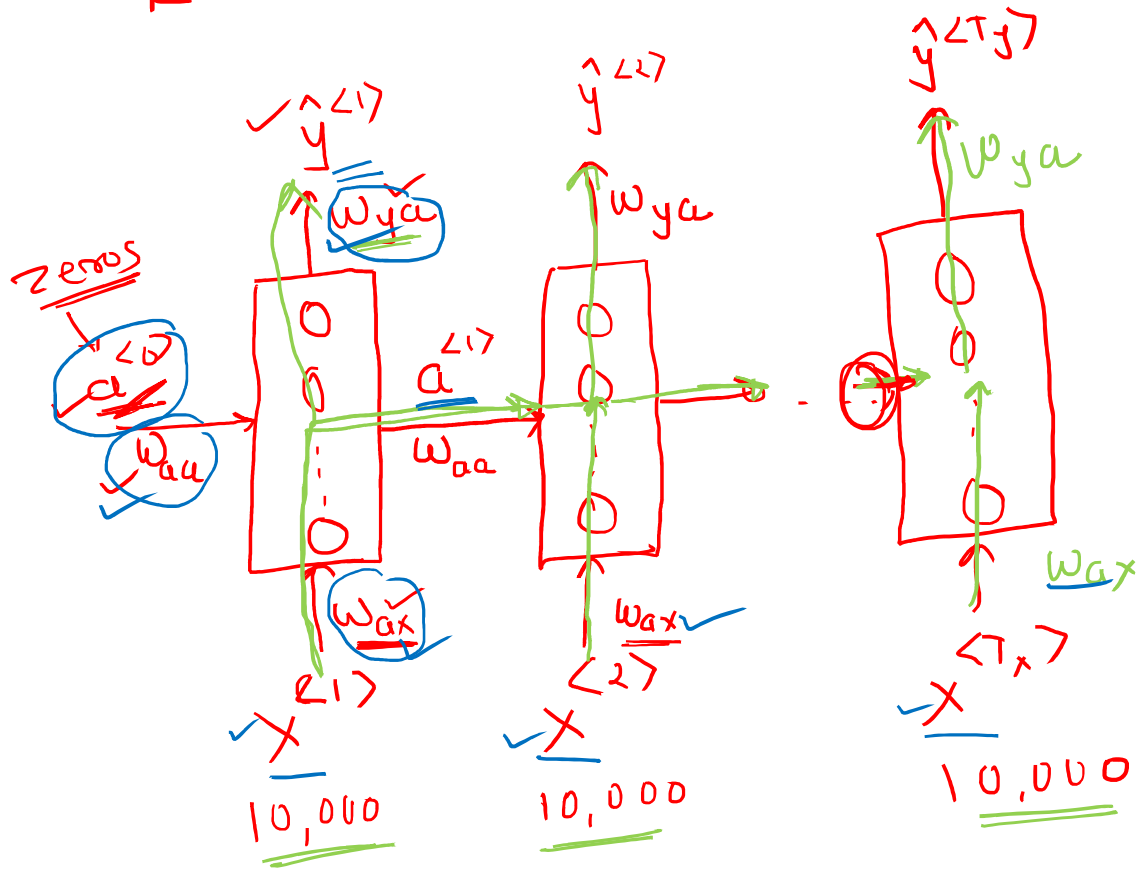
Image generator

- Inputs, outputs can be different length in different example.
- Doesn't share features learned across different position of text.

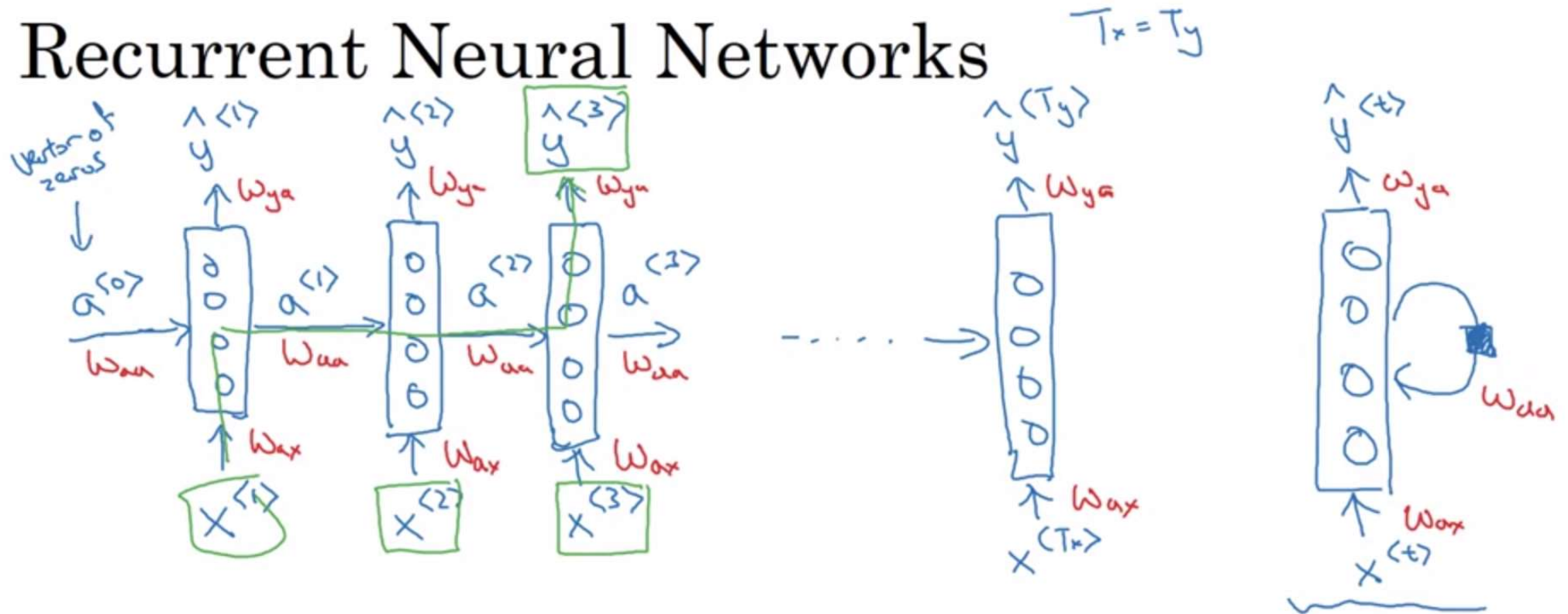




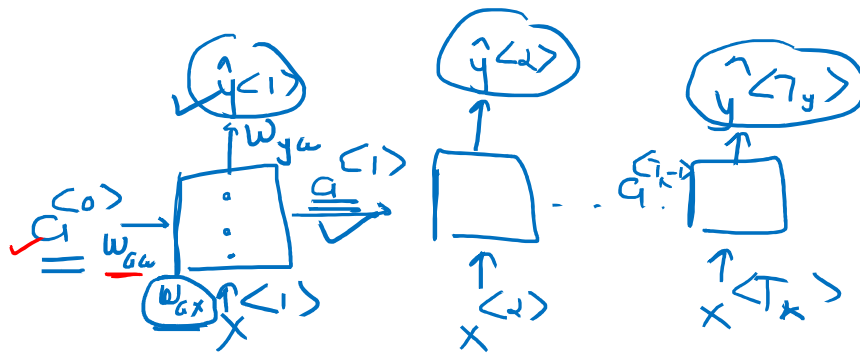
# RNN



# Recurrent Neural Networks



# RNN: forward propagation



$$\underline{a}^{<1>} = g_1 \left( \underbrace{w_{aa} a^{<0>}}_{\text{tanh/Relu}} + \underbrace{w_{ax} x^{<1>}}_{\text{0}} + b_a \right)$$

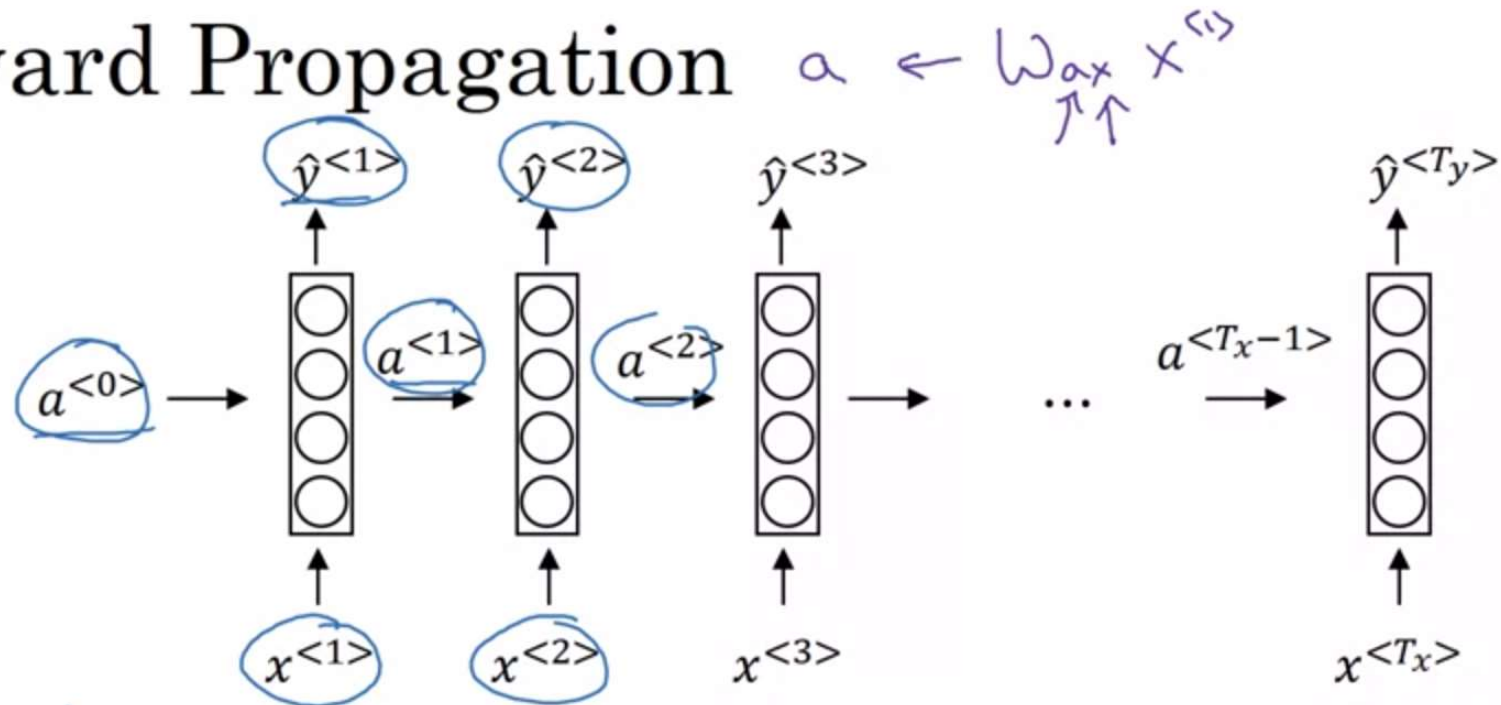
$$\hat{y}^{<1>} = g_2 \left( \underbrace{w_{ya} a^{<1>}}_{\text{Sigmoid}} + b_y \right)$$



$$w_a \begin{bmatrix} w_{aa} & w_{ax} \end{bmatrix} \begin{bmatrix} a^{<t-1>} \\ \vdots \\ x^{<t>} \end{bmatrix}$$

$$\left\{ \begin{aligned} a^{<t>} &= g_1^t (w_{aa} a^{<t-1>} + w_{ax} x^{<t>} + b_a) \\ \hat{y}^{<t>} &= g_2^t (w_{ya} a^{<t>} + b_y) \end{aligned} \right.$$

# Forward Propagation

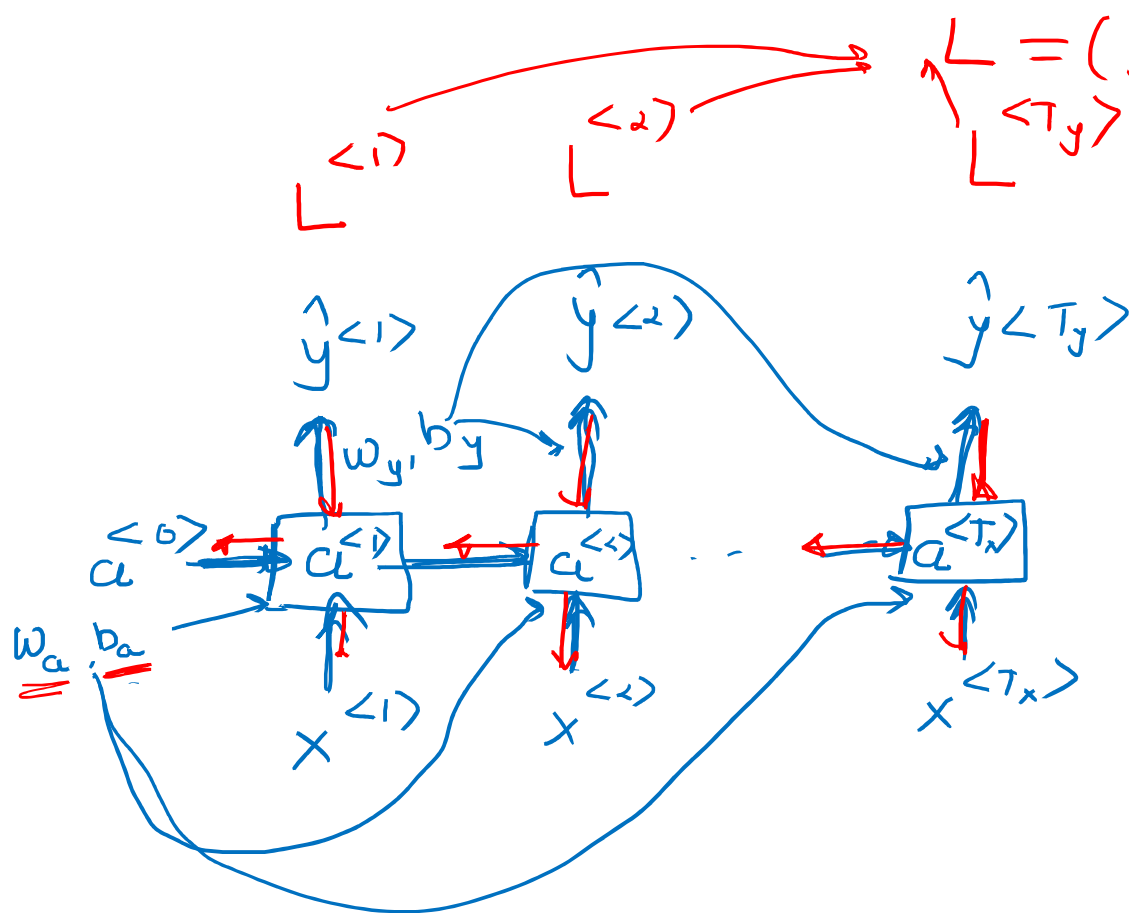


$$a^{<0>} = \vec{0}.$$

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

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

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



Back propagation through time

$$L = (y^*, y)$$

$$L^{<T_y>}$$

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

# Simplified RNN notation

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

Annotations:  $W_{aa}$  is  $(100, 100)$ ,  $W_{ax}$  is  $(100, 10,000)$ . A green arrow indicates the sequence  $a^{<t-1>}$  and  $x^{<t>}$  are inputs to the function  $g$ .

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

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

Annotations:  $W_y$  is  $(100, 10100)$ . A blue arrow indicates the sequence  $a^{<t>}$  is input to the function  $g$ .

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

Annotations:  $W_a$  is  $(100, 10100)$ . A green circle highlights  $W_a$  and a purple box highlights the input vector  $[a^{<t-1>}, x^{<t>}]$ . A blue arrow indicates the sequence  $a^{<t-1>}$  and  $x^{<t>}$  are inputs to the function  $g$ .

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

Annotations:  $W_{aa}$  is  $(100, 100)$ ,  $W_{ax}$  is  $(100, 10,000)$ . A green box highlights the matrix  $\begin{bmatrix} W_{aa} & W_{ax} \end{bmatrix}$  and a blue arrow indicates the sequence  $a^{<t-1>}$  and  $x^{<t>}$  are inputs to the function  $g$ .

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

Annotations:  $a^{<t-1>}$  is  $100$ ,  $x^{<t>}$  is  $10,000$ . A green circle highlights the vector  $\begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix}$  and a blue arrow indicates the sequence  $a^{<t-1>}$  and  $x^{<t>}$  are inputs to the function  $g$ .

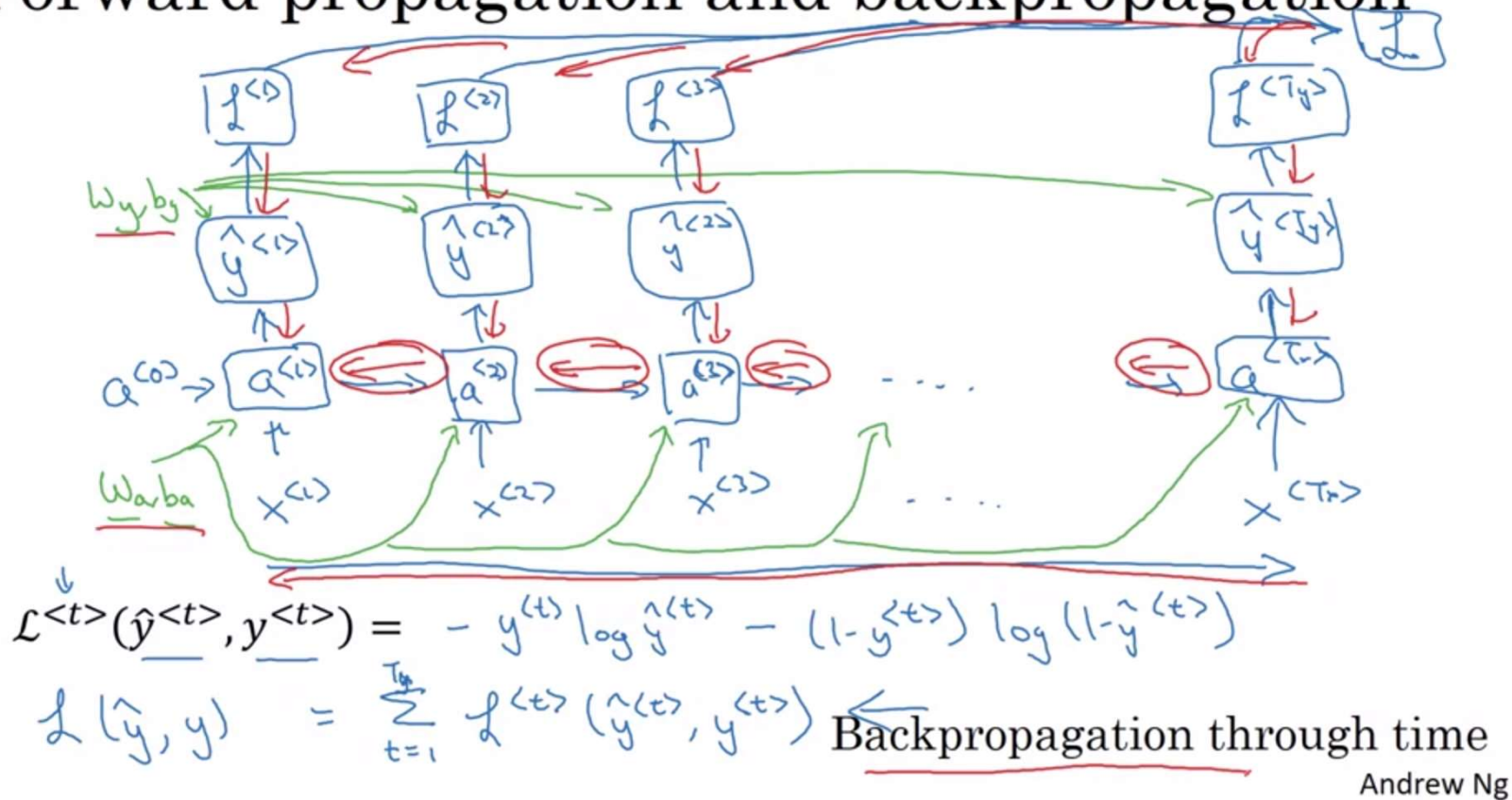
$$\begin{bmatrix} W_{aa} & W_{ax} \end{bmatrix} \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} = W_{aa}a^{<t-1>} + W_{ax}x^{<t>}$$

Annotations: A green box highlights the matrix  $\begin{bmatrix} W_{aa} & W_{ax} \end{bmatrix}$  and a blue arrow indicates the sequence  $a^{<t-1>}$  and  $x^{<t>}$  are inputs to the function  $g$ .



# Backpropagation through time

## Forward propagation and backpropagation





# Different RNN Architectures

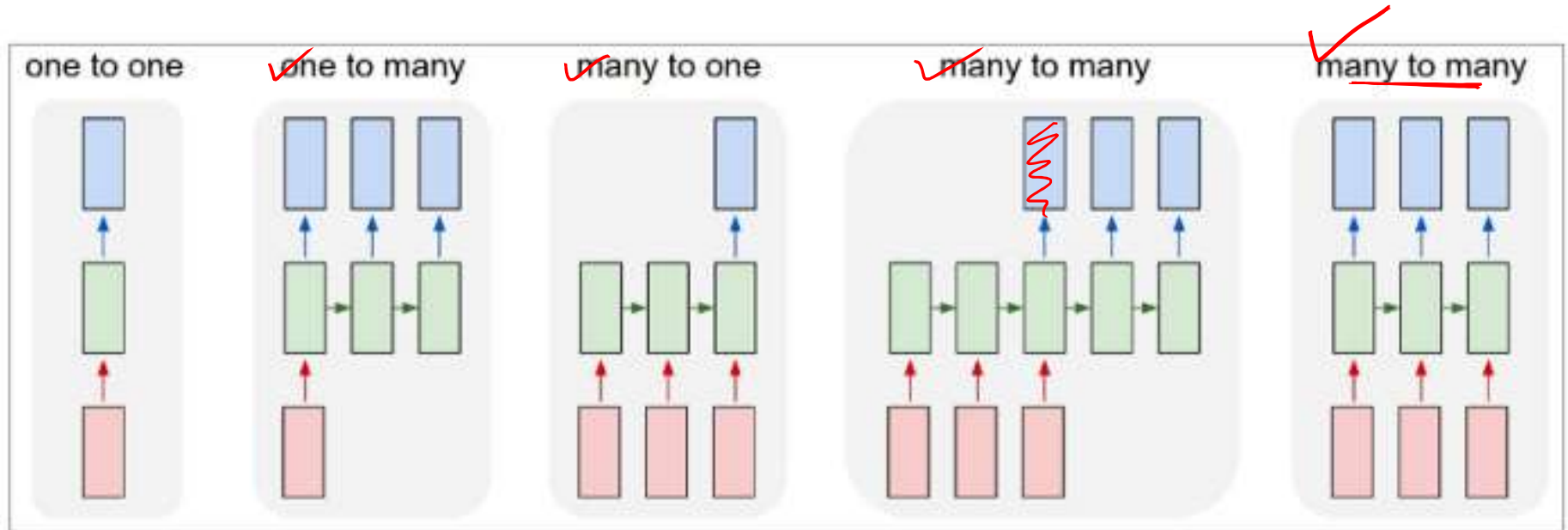


image classification

sentiment classification

Name entity recognition

image captioning,  
✓ Music Generation

Machine translation  
 $T_x \neq T_y$

$T_x = T_y$

# Advantages of Recurrent Neural Network

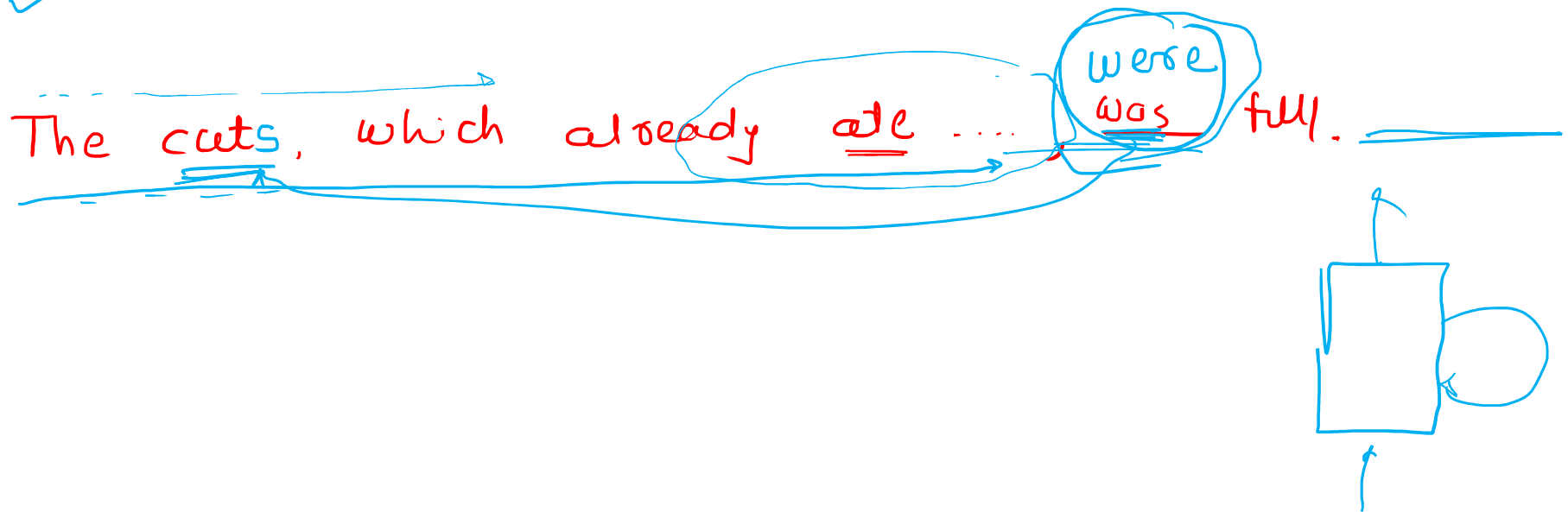
- It is useful in time series prediction only because of the feature to remember previous inputs as well.

# Disadvantages of Recurrent Neural Network

- Vanishing gradients or short term memory in RNN
- Exploding gradients. →
- Training an RNN is a very difficult task.

GRU, LSTM

bi directional



# Vanishing gradients

- As our RNN processes more steps, it has trouble retaining information from the previous steps. So, in some cases, this might not be a problem, where the word just depends on its previous neighboring word.

RP3

- For eg. the sentence: I can speak French very well.

RP2

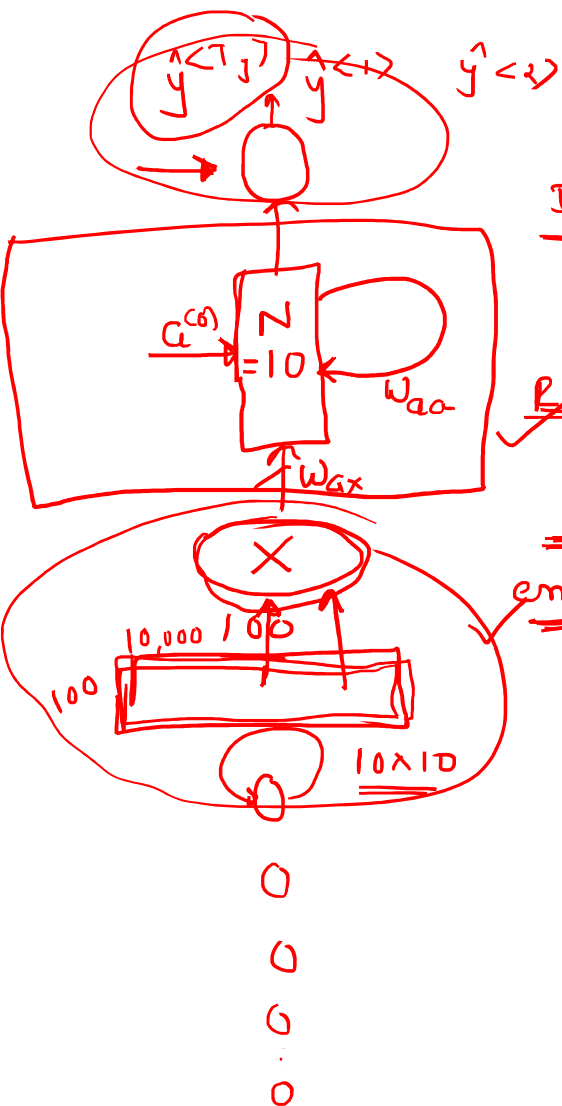
- But consider this sentence: I am going to France, the language spoken there is \_\_\_\_\_.

## Slide 20

---

- RP2** Now the answer “French” here has a dependency on the word France, which is far away from it. This type of dependency is known as long term dependency, and the normal structure of RNN fails to operate over these.  
RONAKKUMAR PATEL, 28-Sep-20
- RP3** Now the word ‘well’ in this sentence is very intuitive to come at this place and RNN can handle such sentences effectively.  
RONAKKUMAR PATEL, 28-Sep-20

- Short-Term Memory and the vanishing gradient is due to the nature of back-propagation; an algorithm used to train and optimize neural networks.
- When doing back propagation, each node in a layer calculates its gradient with respect to the effects of the gradients, in the layer before it. So if the adjustments to the layers before it is small, then adjustments to the current layer will be even smaller.



Dense - 1

$$w_{ay} + b_y = 11$$

RNN  $\Rightarrow$

$$w_{ax} + w_{aa} + b_a = 1110$$

embd. vocab size

$$10,000 \times 100 = 1,000,000$$

the which

100



Any  
Questions

