

# Recurrent Neural Network

Tianchu.Zhao@uts.edu.au

# Why Recurrent

# Sequence data

- Sequence data means that there exists a relation between data points
- The relation can be time, spatial, semantic ...
- Each sequence consists of multiple data points can have vary length

# Neural network on sequence data

- For vary length input, we need vary length neuron (parameters) in each layer
- If we capture a certain feature in one area, this feature can not be broadcasted to another area
- Can not accommodate the large amount of parameters
- To address the above problems, we use Recurrent Neural Network(RNN)

# Applications of RNN

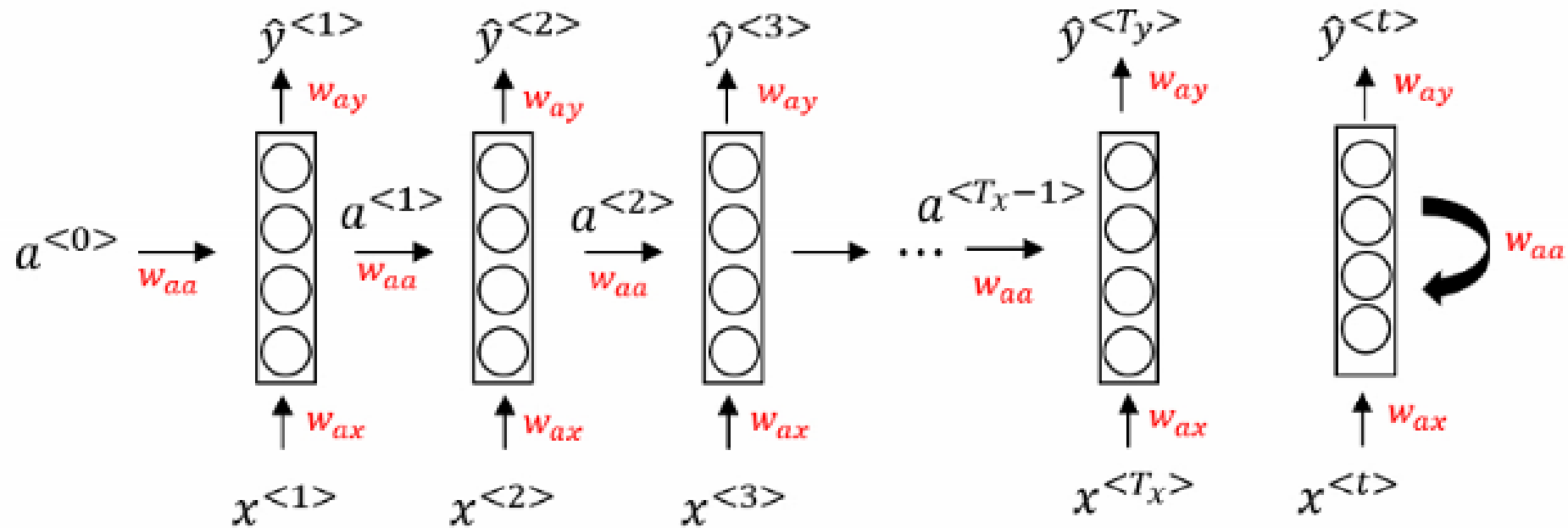
- name entity recognition: sentence to location of noun, action
  - eg, what is the noun(machine learning) action(training)
- sentiment analysis: sentence to score
- speech recognition: audio to text
- machine translation: French to English
- video activity recognition: video to "running"
- music generation: notes to music sheet/audio
- DNA/Protein analysis: sequence to property score

How Recurrent

# Mathematic notation

- $x$ , data point, each data point is a sequence
- $x^{(t)}$ , the  $t$ -th element in a sequence
- $y^{(t)}$ , the  $t$ -th label in a sequence
- $T_x$ , the input length
- $T_y$ , the output length
- put together we have
- $x^{(i)(t)}$ ,  $i$ -th sequence  $t$ -th element
- $y^{(i)(t)}$ ,  $i$ -th sequence  $t$ -th label
- $T_x^{(i)}$ , input length of  $i$ -th sequence

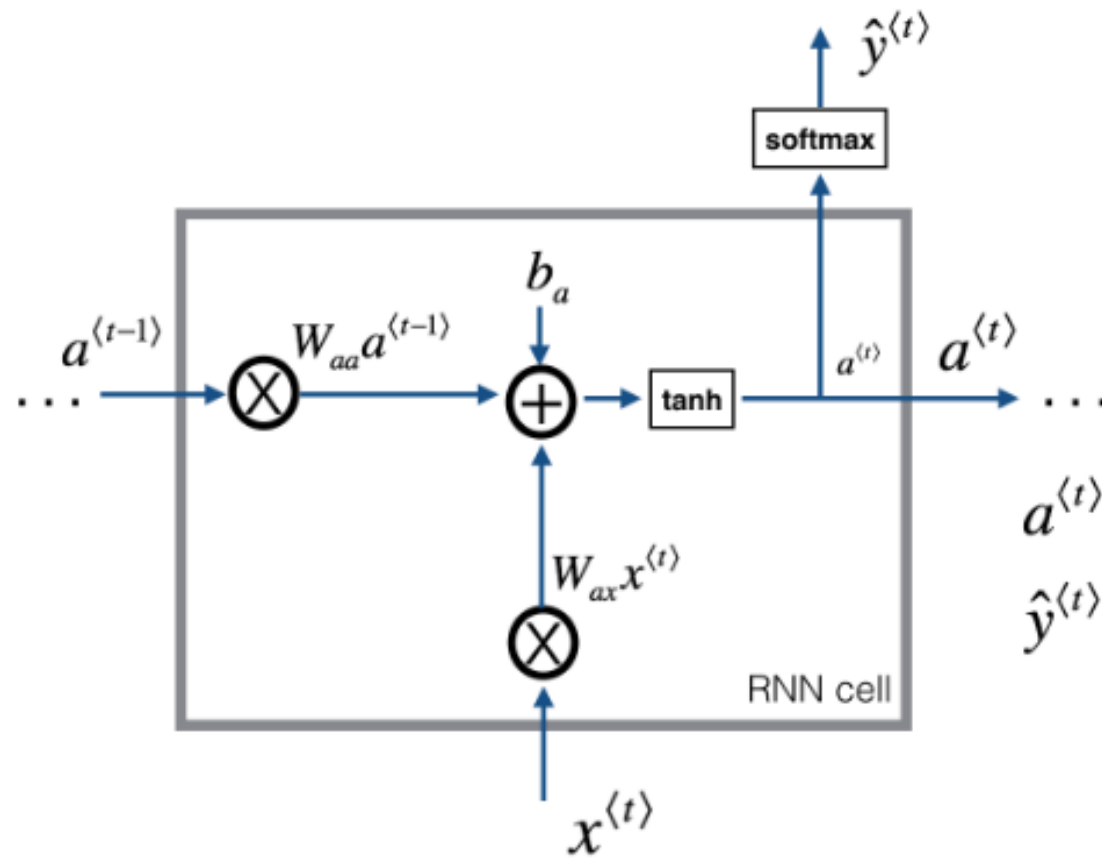
# RNN Structure



- The hidden layer of each corresponding time step will receive the activation value  $a^{(t-1)}$ , as a input,  $a^{(0)}$  is usually 0
- The parameters are input  $w_{ax}$ ,  $w_{aa}$ ,  $w_{ay}$  where we share parameters at every timestamp



# Zoom in each neuron



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

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

# Forward pass

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

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

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

- To simplify the calculation we can combine (append)  $W_{ax}$ ,  $W_{aa}$  to  $W_a$
- combine  $a^{(t-1)}$  and  $x^{(t)}$

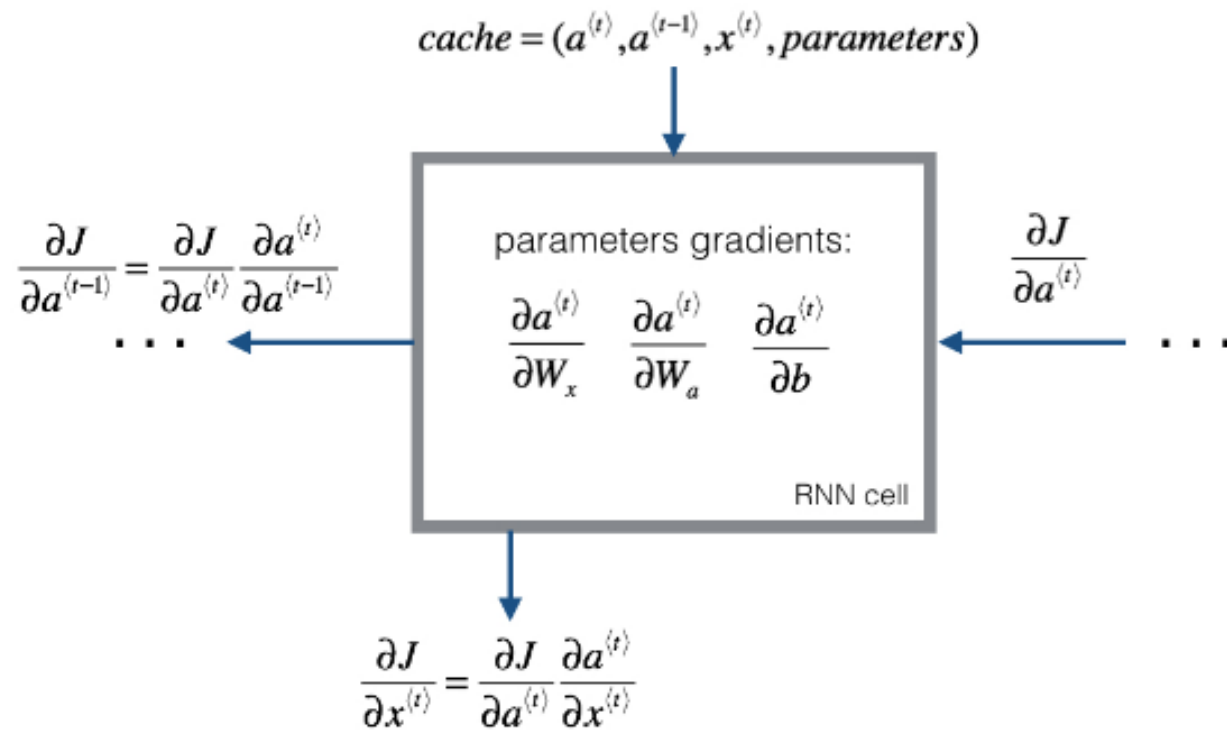
$$W_a = [W_{ax}, W_{aa}]$$

$$a^{(t)} = g_1(W_{aa}a^{(t-1)} + W_{ax}x^{(t)} + b_a) \quad \Longrightarrow \quad a^{(t)} = g_1(W_a[a^{(t-1)}, x^{(t)}] + b_a)$$

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

$$\hat{y}^{(t)} = g_2(W_y a^{(t)} + b_y)$$

# Backpropagation



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

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

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

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

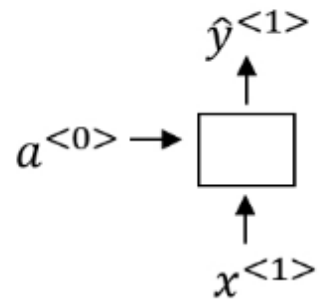
$$\frac{\partial a^{(t)}}{\partial b} = \sum_{batch} (1 - \tanh(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b)^2)$$

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

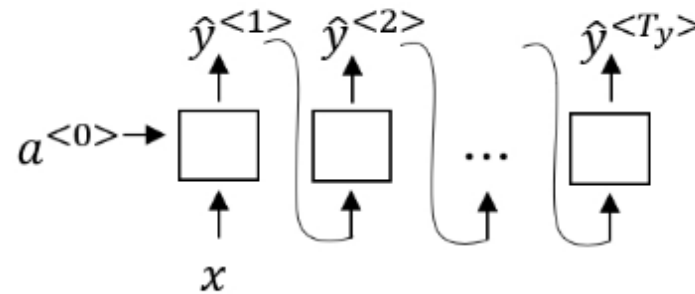
$$\frac{\partial a^{(t)}}{\partial a^{(t-1)}} = W_{aa}^T \cdot (1 - \tanh(W_{ax}x^{(t-1)} + W_{aa}a^{(t-1)} + b)^2)$$

# RNN Structuressssss

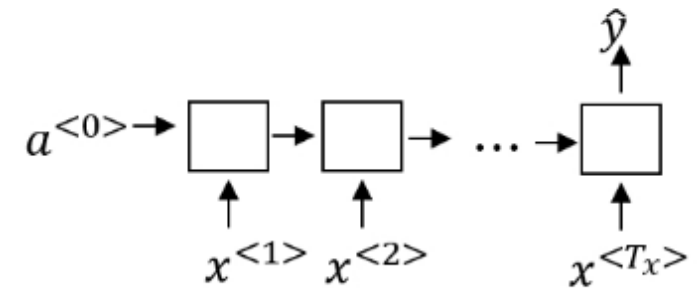
- we can have vary input output structure



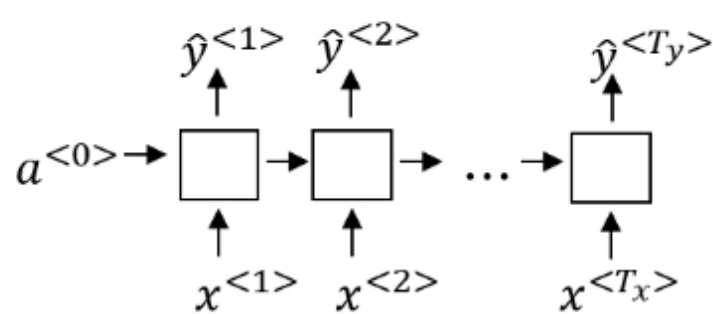
One to one



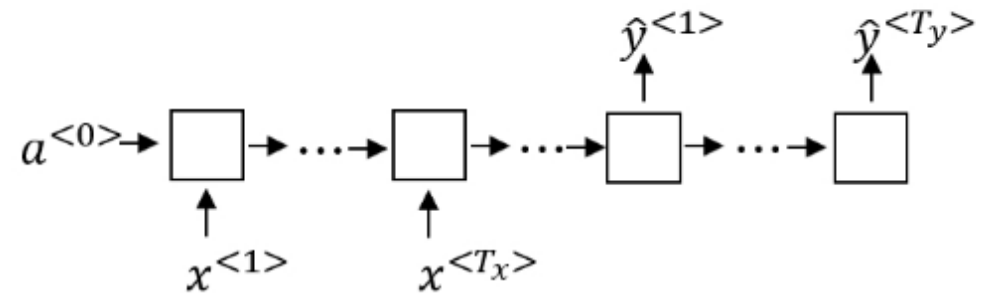
One to many



Many to one



Many to many



Many to many

# RNN Issues

- The cat, which already ate a bunch of food, was full
- The cats, which already ate a bunch of food, were full
- The verb at the end is dependent on the noun at the beginning (subject-verb agreement)
- but RNN is not good at learning this long term dependent relationship
- As we go deeper, vanishing gradient
- To combat these two issues we have GRU and LSTM

# Gated Recurrent Units (GRU)

- introduce additional variable call  $c$  in the recurrent unit which stands for memory cell
- We have the following new variables
- $c^{(t)}$ , output activation value  $a^{(t)}$
- $\tilde{c}^{(t)}$ , the next activation value  $a^{(t)}$  candidate
- $\Gamma_u$ , the Update Gate, it tells when to update memory cell  $c$  value

$$a^{(t)} = g_1(W_a[a^{(t-1)}, x^{(t)}] + b_a)$$



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

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

$$c^{(t)} = \Gamma_u \times \tilde{c}^{(t)} + (1 - \Gamma_u) \times c^{(t-1)}$$

$$a^{(t)} = c^{(t)}$$

- Complete version also includes a Relevance gate  $\Gamma_r$  to represent the relevance between  $c^{(t)}$  and  $c^{(t)}$

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

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

$$c^{(t)} = \Gamma_u \times \tilde{c}^{(t)} + (1 - \Gamma_u) \times c^{(t-1)}$$

$$a^{(t)} = c^{(t)}$$



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

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

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

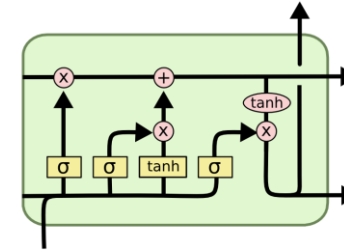
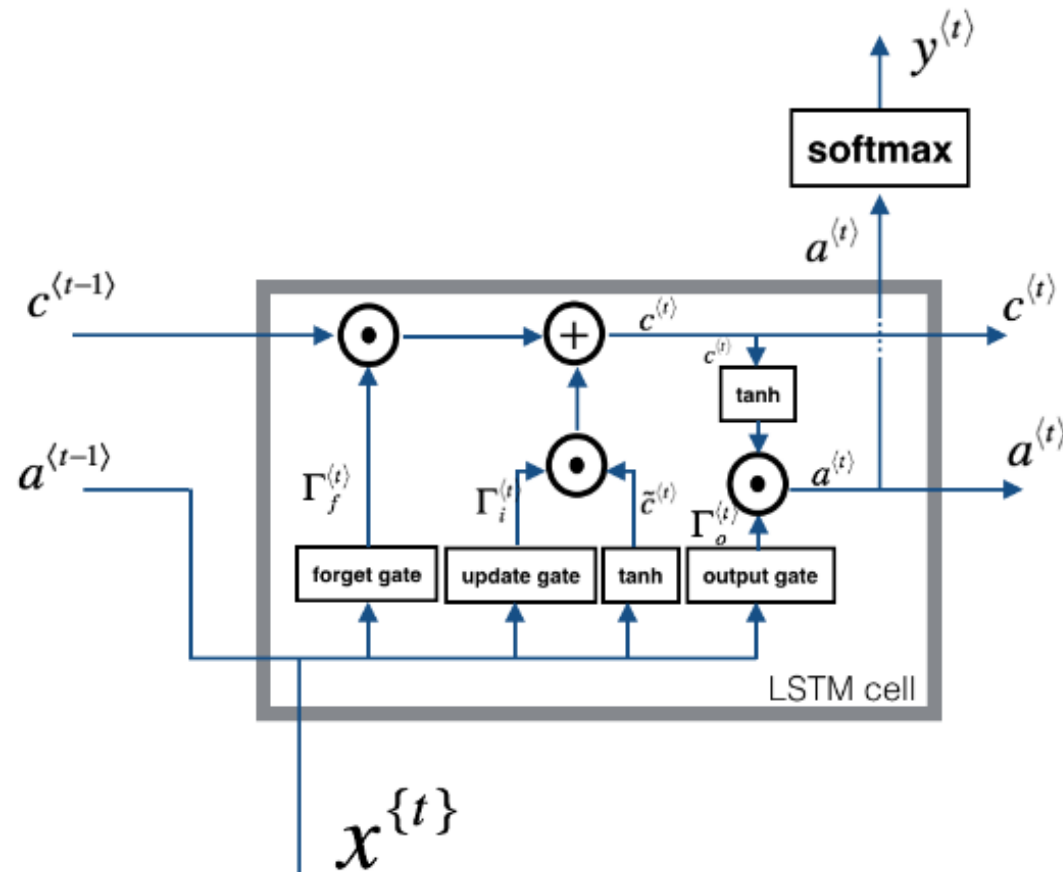
$$c^{(t)} = \Gamma_u \times \tilde{c}^{(t)} + (1 - \Gamma_u) \times c^{(t-1)}$$

$$a^{(t)} = c^{(t)}$$



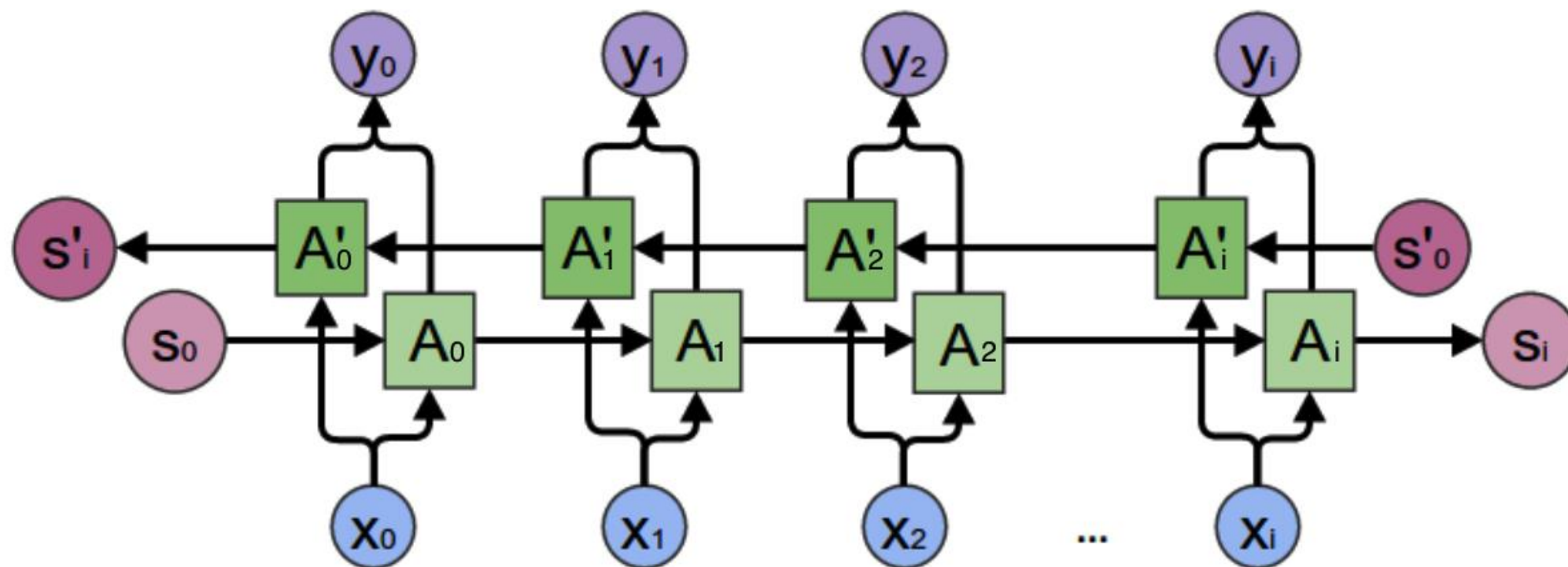
# Long Short Term Memory (LSTM)

- more flexible and powerful than GRU
- it involves Forget Gate  $\Gamma_f$  and output gate  $\Gamma_o$

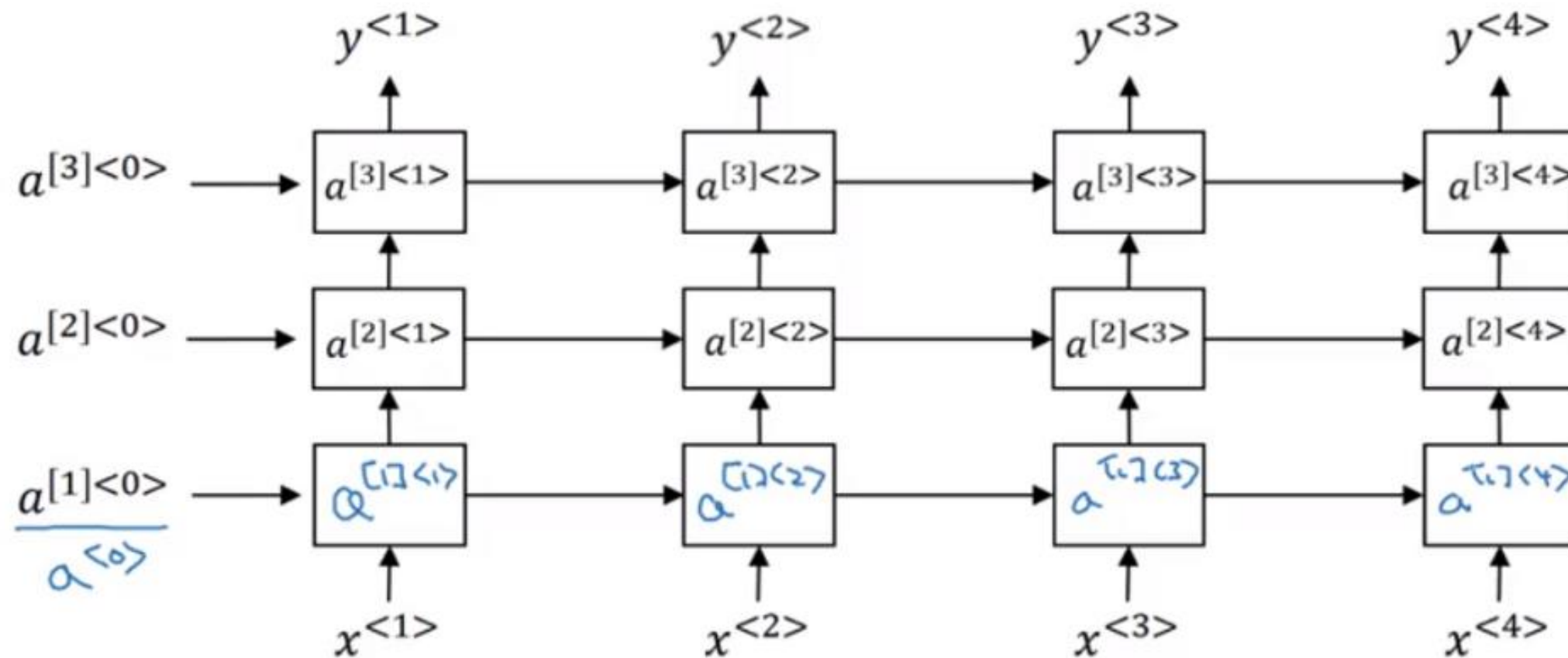


$$\begin{aligned}\Gamma_f^{(t)} &= \sigma(W_f[a^{(t-1)}, x^{(t)}] + b_f) \\ \Gamma_u^{(t)} &= \sigma(W_u[a^{(t-1)}, x^{(t)}] + b_u) \\ \tilde{c}^{(t)} &= \tanh(W_c[a^{(t-1)}, x^{(t)}] + b_c) \\ c^{(t)} &= \Gamma_f^{(t)} \circ c^{(t-1)} + \Gamma_u^{(t)} \circ \tilde{c}^{(t)} \\ \Gamma_o^{(t)} &= \sigma(W_o[a^{(t-1)}, x^{(t)}] + b_o) \\ a^{(t)} &= \Gamma_o^{(t)} \circ \tanh(c^{(t)})\end{aligned}$$

# Bidirectional RNN (BRNN)

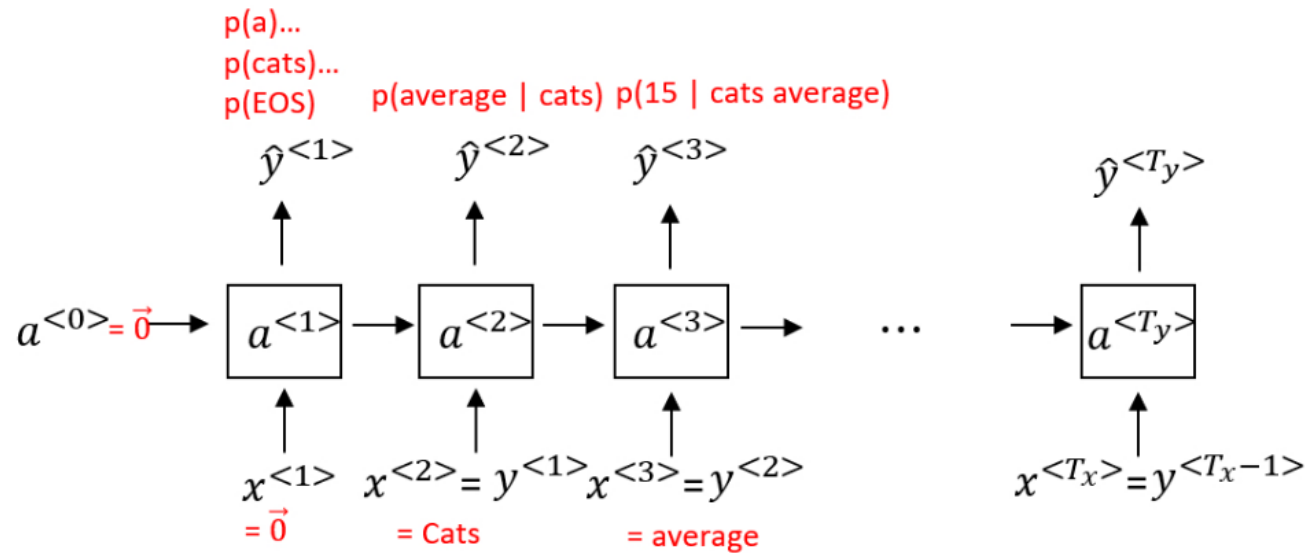


# Deep RNN (DRNN)



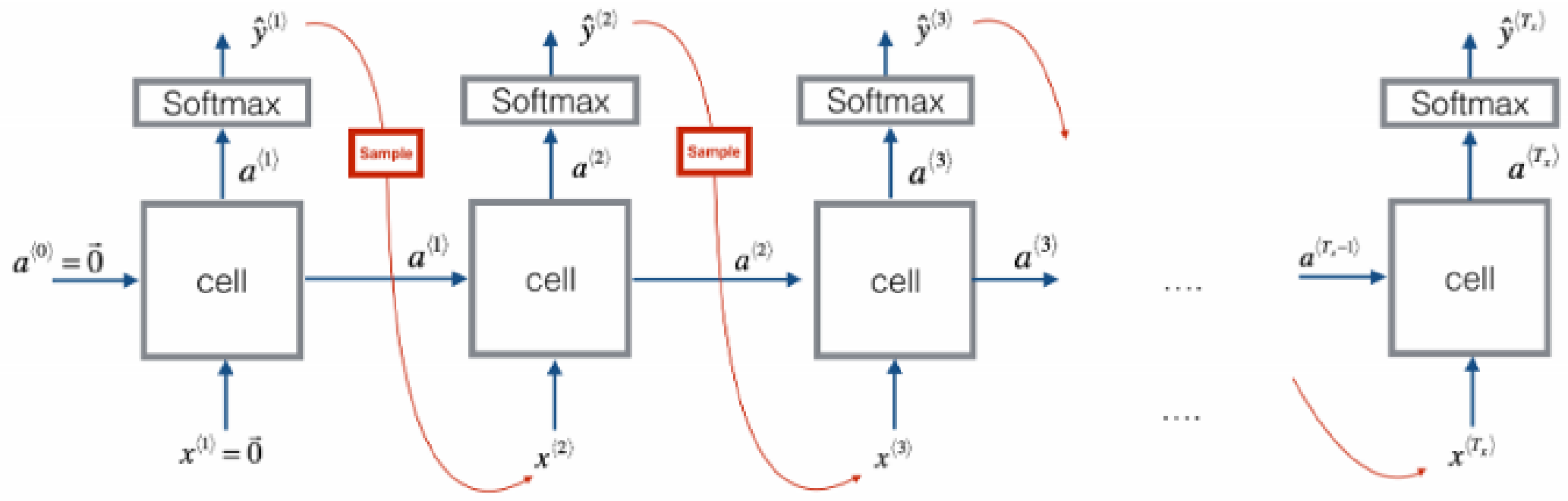
# Simple RNN example

- Language Model is to describe fact base on language to model language using mathematic, it is able to predict the probability of element occurring in a sequence
- we use corpus (text compose of many sentences)
- the first step to build the model is tokenize (dictionaries words)
- we then encode each word into word vector and add EOS (End of Sentence) to signify the end of sentence



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

- after we trained a language model, we can use sampling to learn what we've learnt
- first input  $a^{(0)}$   $x^{(1)}$  are 0, we get a distribution of softmax
- then we sample a random word
- we keep sampling till EOS
- then we get some sentence



# Grand Plan (besides the lecture material)

- (✓) Neural Network Foundation
- (✓) Ensemble/Optimise your neural network
- (✓) Convolution Neural Network
- (✓) Principal Component Analysis / Big data
- (✓) Reinforcement Learning
- (✓) Recurrent Neural Network
- (□) Natural Language Processing
- (□) Generative Adversarial Network