



# 人工智能技术及应用

Artificial Intelligence and Application

---

# Recurrent Neural Network (RNN)





# Example Application

- Slot Filling



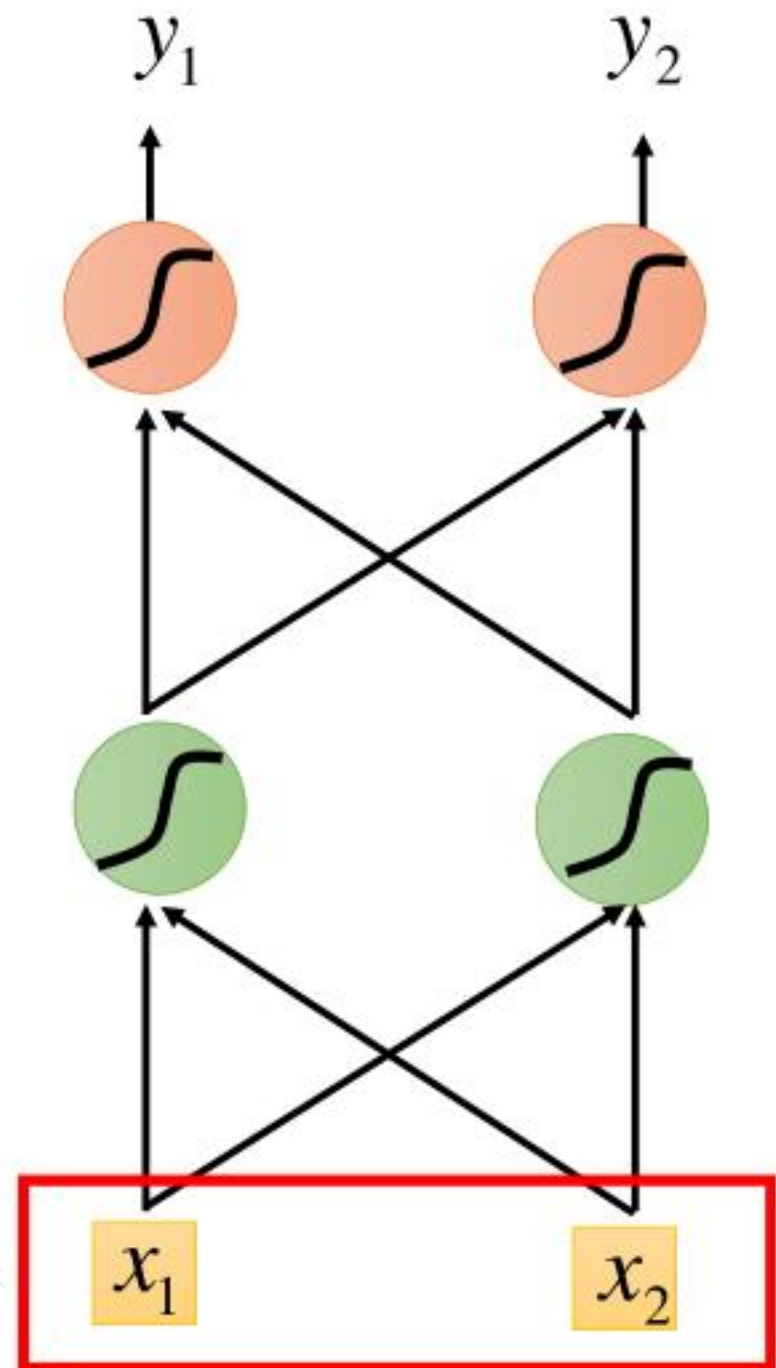
# Example Application

Solving slot filling by  
Feedforward network?

Input: a word

(Each word is represented  
as a vector)

Shanghai



# 1-of-N encoding

How to represent each word as a vector?

**1-of-N Encoding**    lexicon = {apple, bag, cat, dog, elephant}

The vector is lexicon size.

Each dimension corresponds  
to a word in the lexicon

The dimension for the word  
is 1, and others are 0

apple = [ 1   0   0   0   0 ]

bag    = [ 0   1   0   0   0 ]

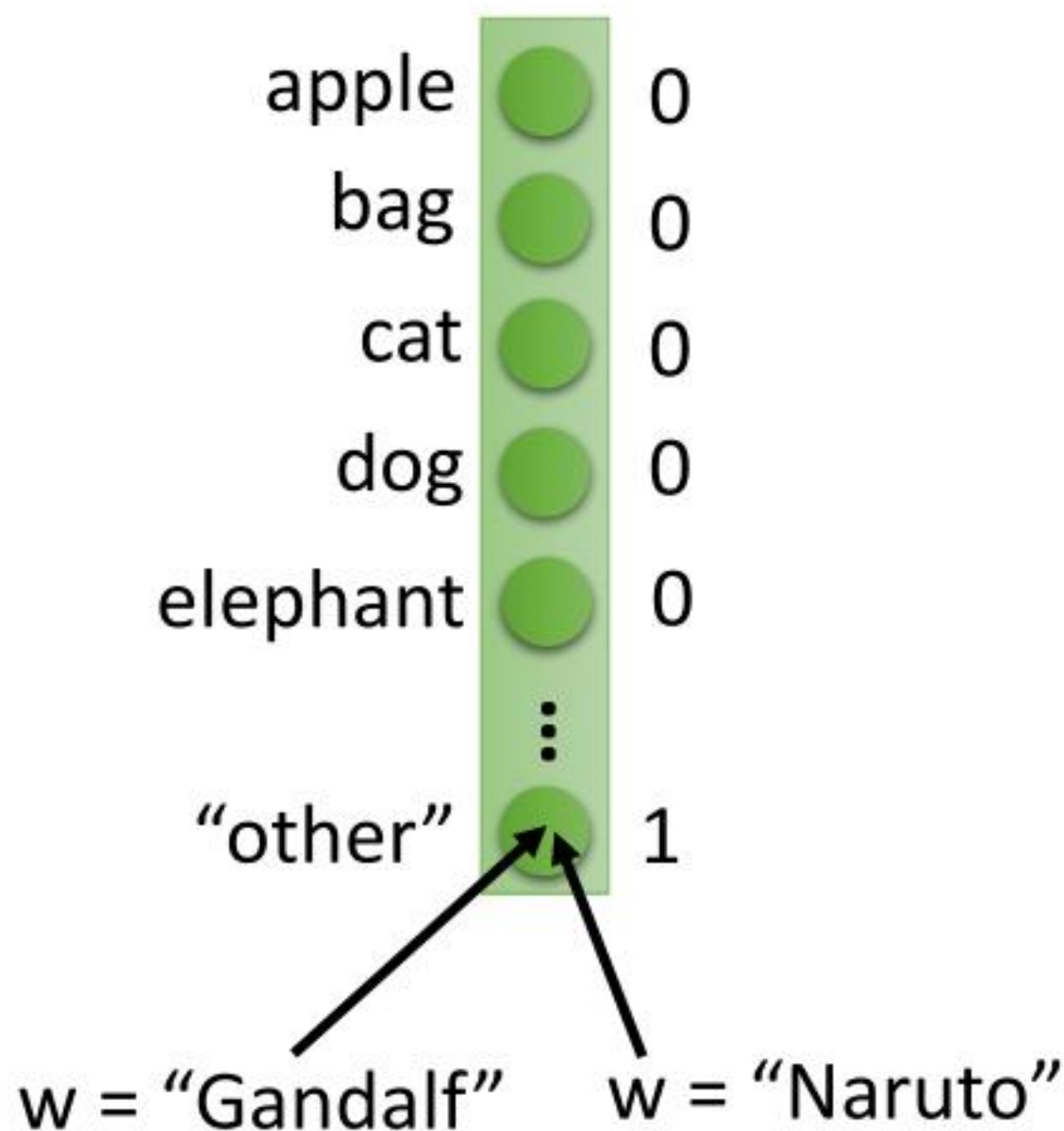
cat    = [ 0   0   1   0   0 ]

dog    = [ 0   0   0   1   0 ]

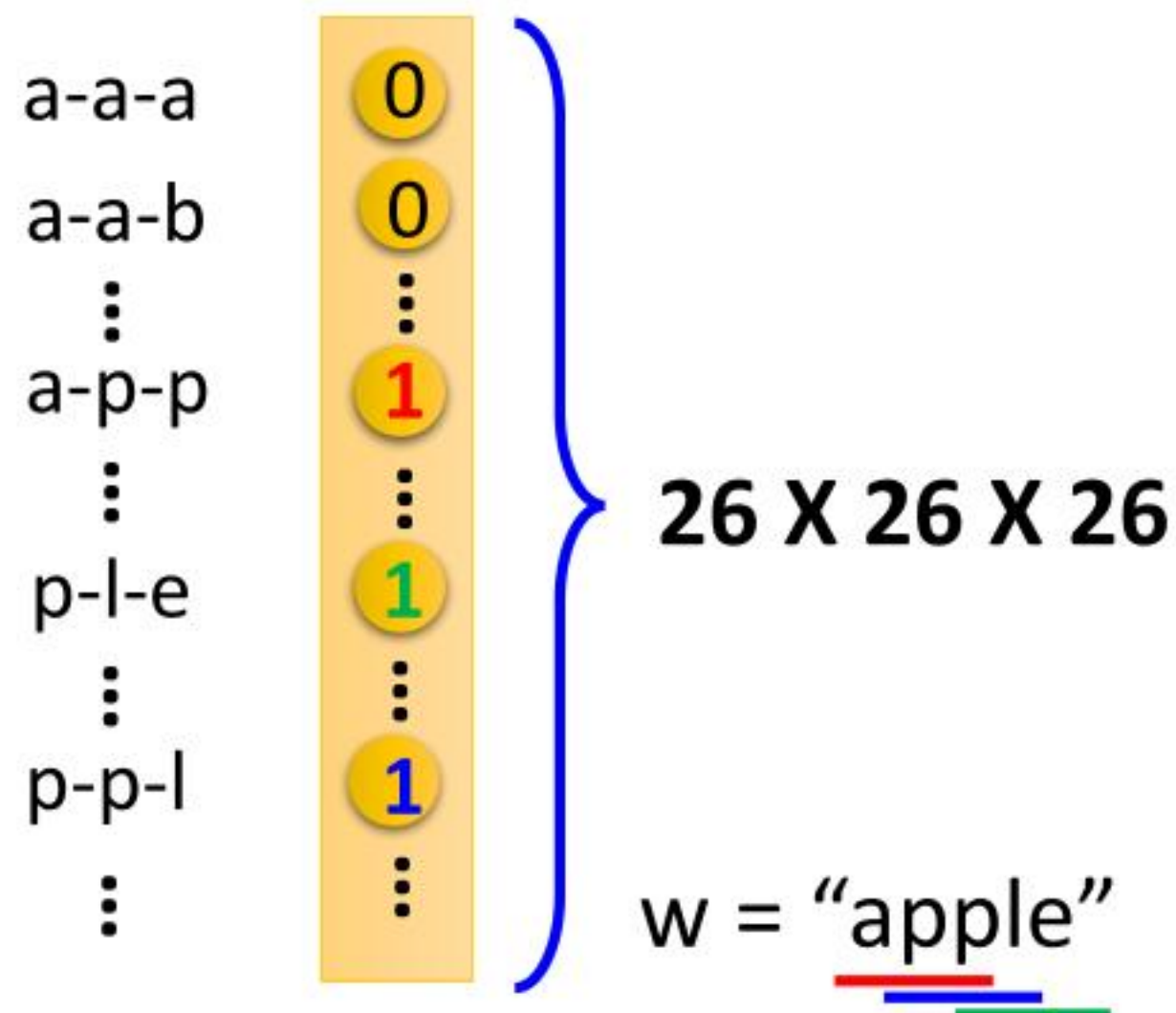
elephant = [ 0   0   0   0   1 ]

# Beyond 1-of-N encoding

## Dimension for "Other"



## Word hashing





# Example Application

Solving slot filling by  
Feedforward network?

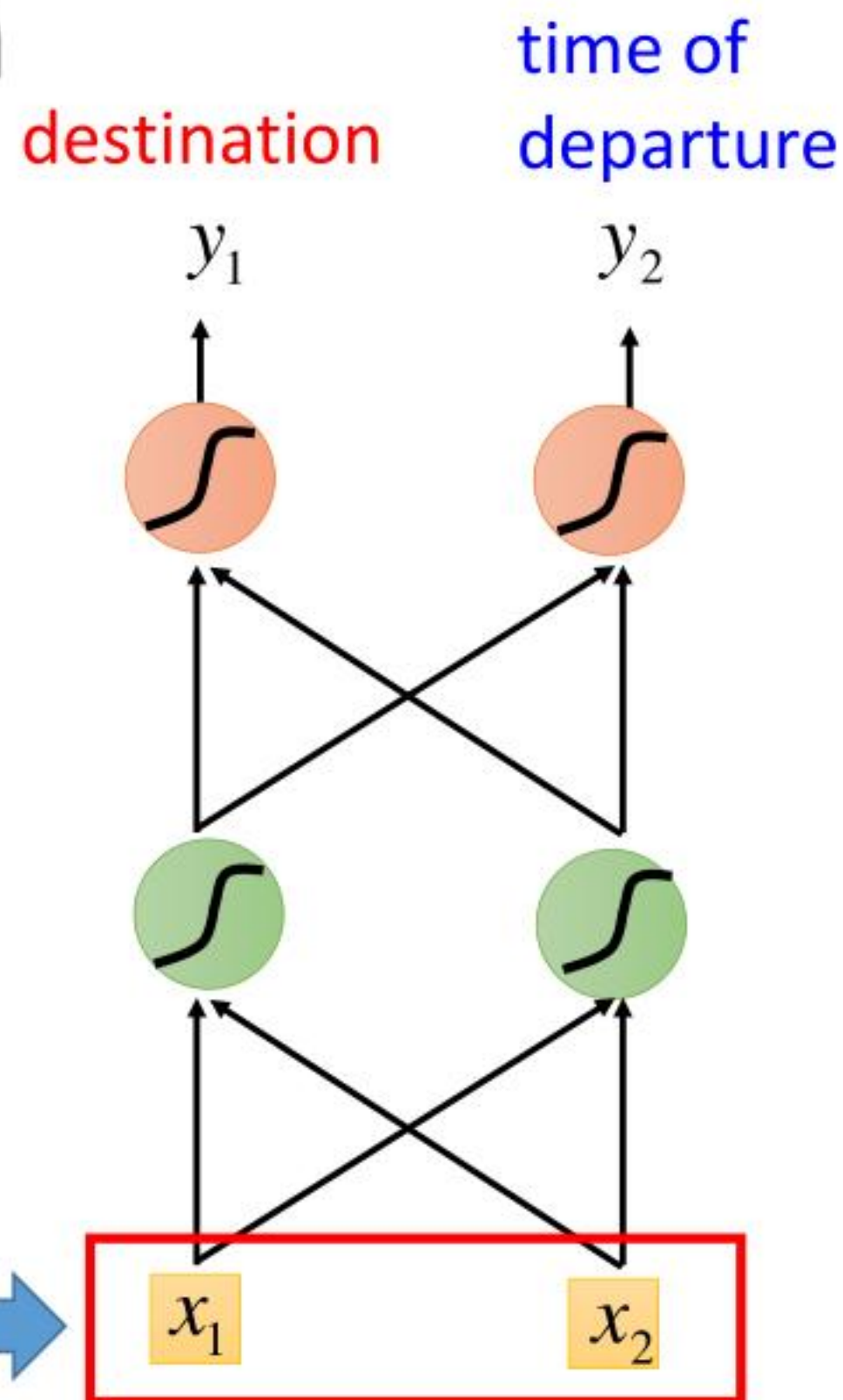
Input: a word

(Each word is represented  
as a vector)

Output:

Probability distribution that  
the input word belonging to  
the slots

Shanghai



# Example Application

arrive Shanghai on June 2<sup>nd</sup>

other destination other time time

Problem?

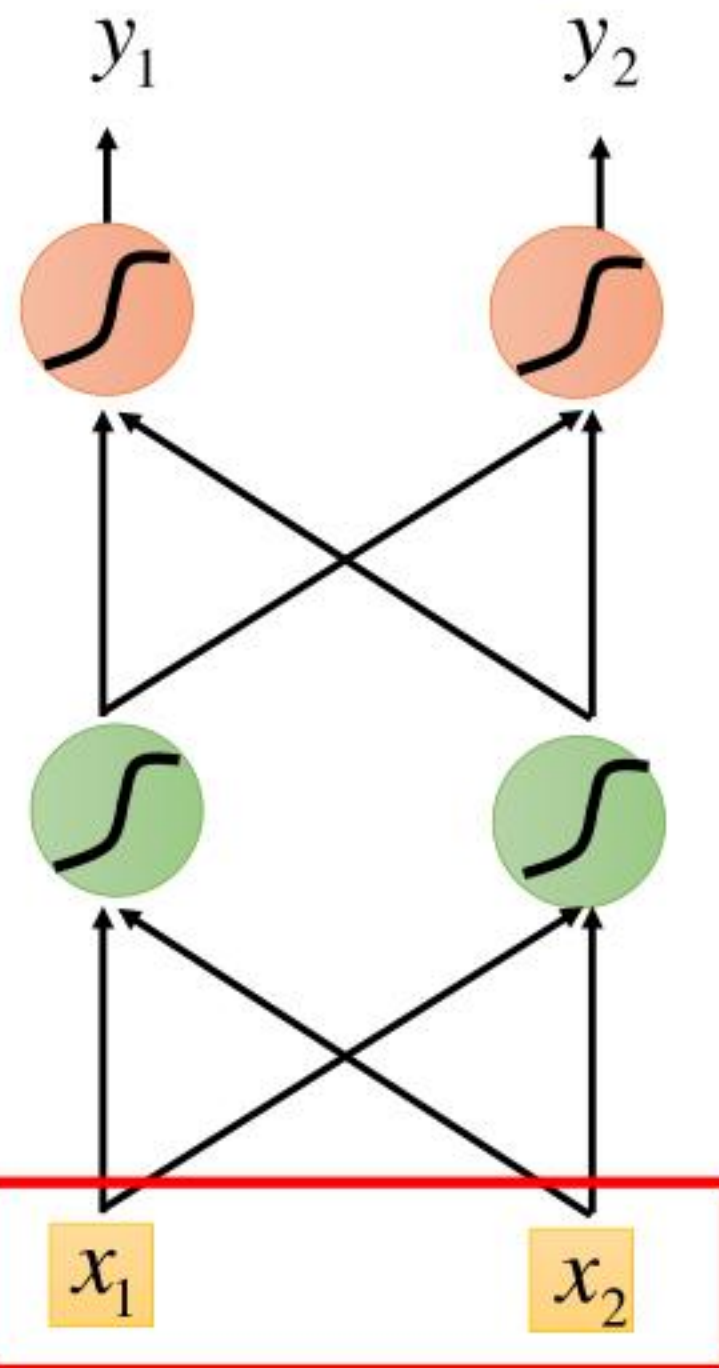
leave Shanghai on June 2<sup>nd</sup>

place of departure

Neural network  
needs memory!

destination

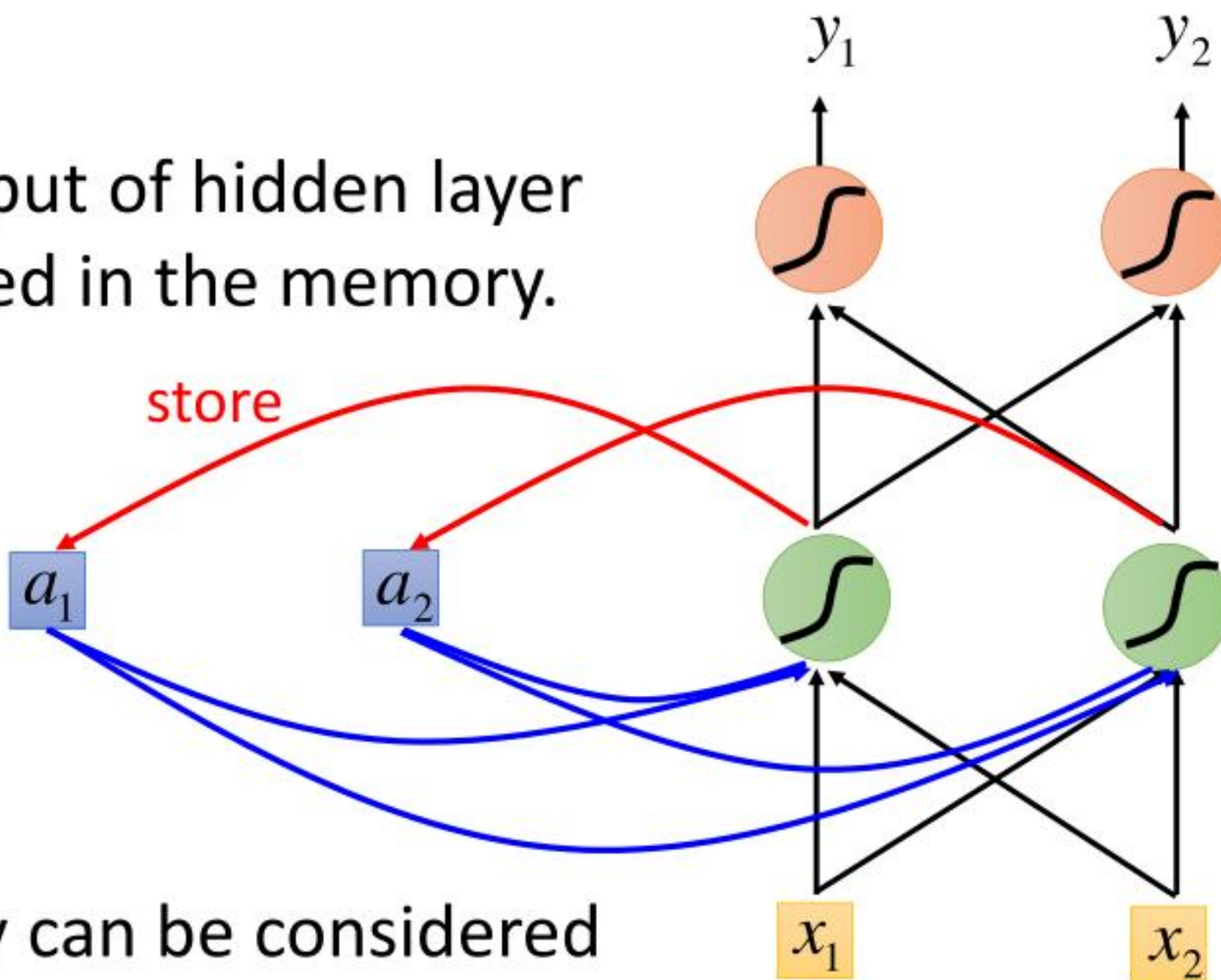
time of  
departure





# Recurrent Neural Network (RNN)

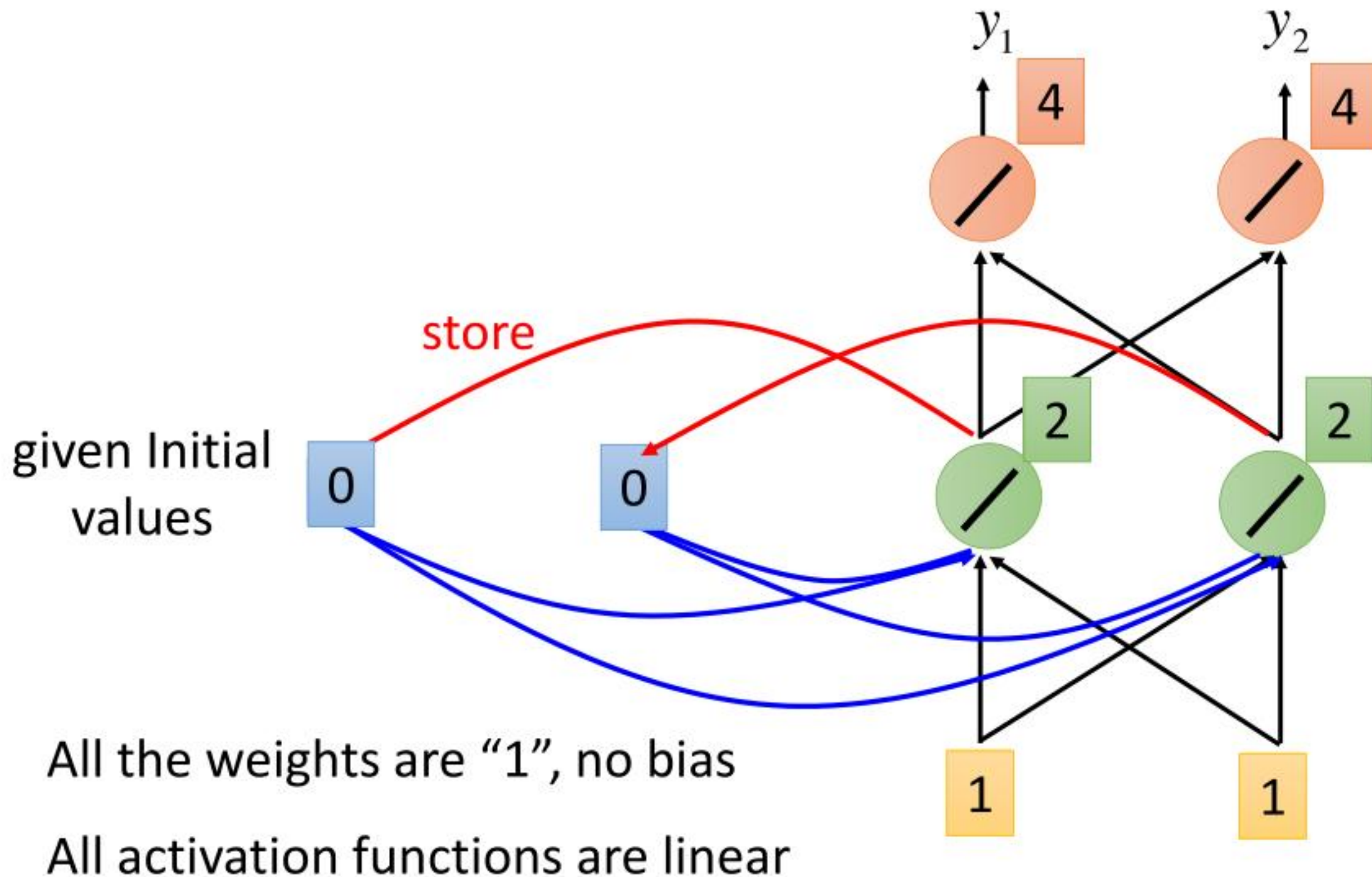
The output of hidden layer are stored in the memory.



Memory can be considered as another input.

# Example

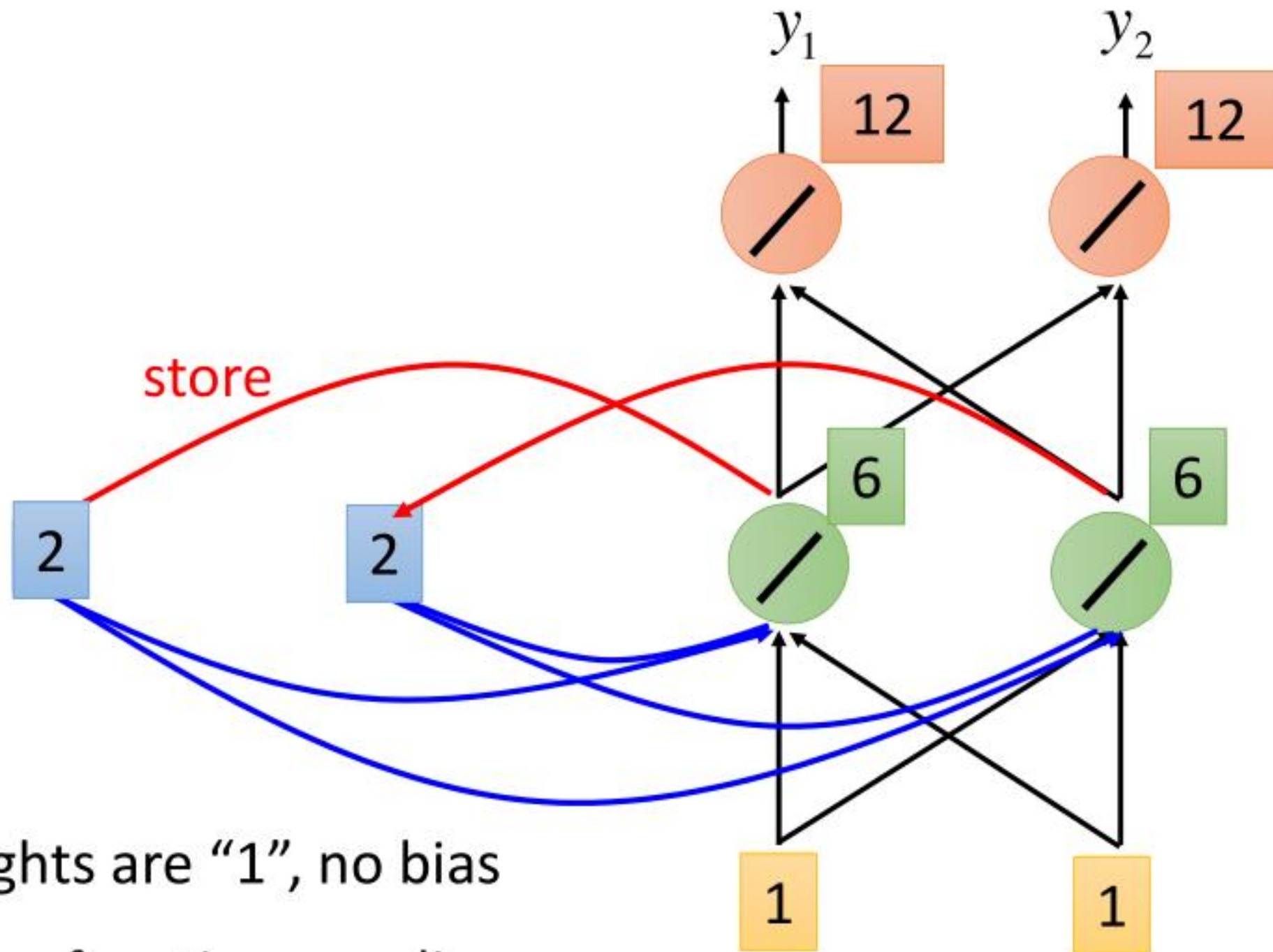
Input sequence:  $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} \dots$   
output sequence:  $\begin{bmatrix} 4 \\ 4 \end{bmatrix}$



# Example

Input sequence:  $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} \dots$

output sequence:  $\begin{bmatrix} 4 \\ 4 \end{bmatrix} \begin{bmatrix} 12 \\ 12 \end{bmatrix}$



All the weights are "1", no bias

All activation functions are linear

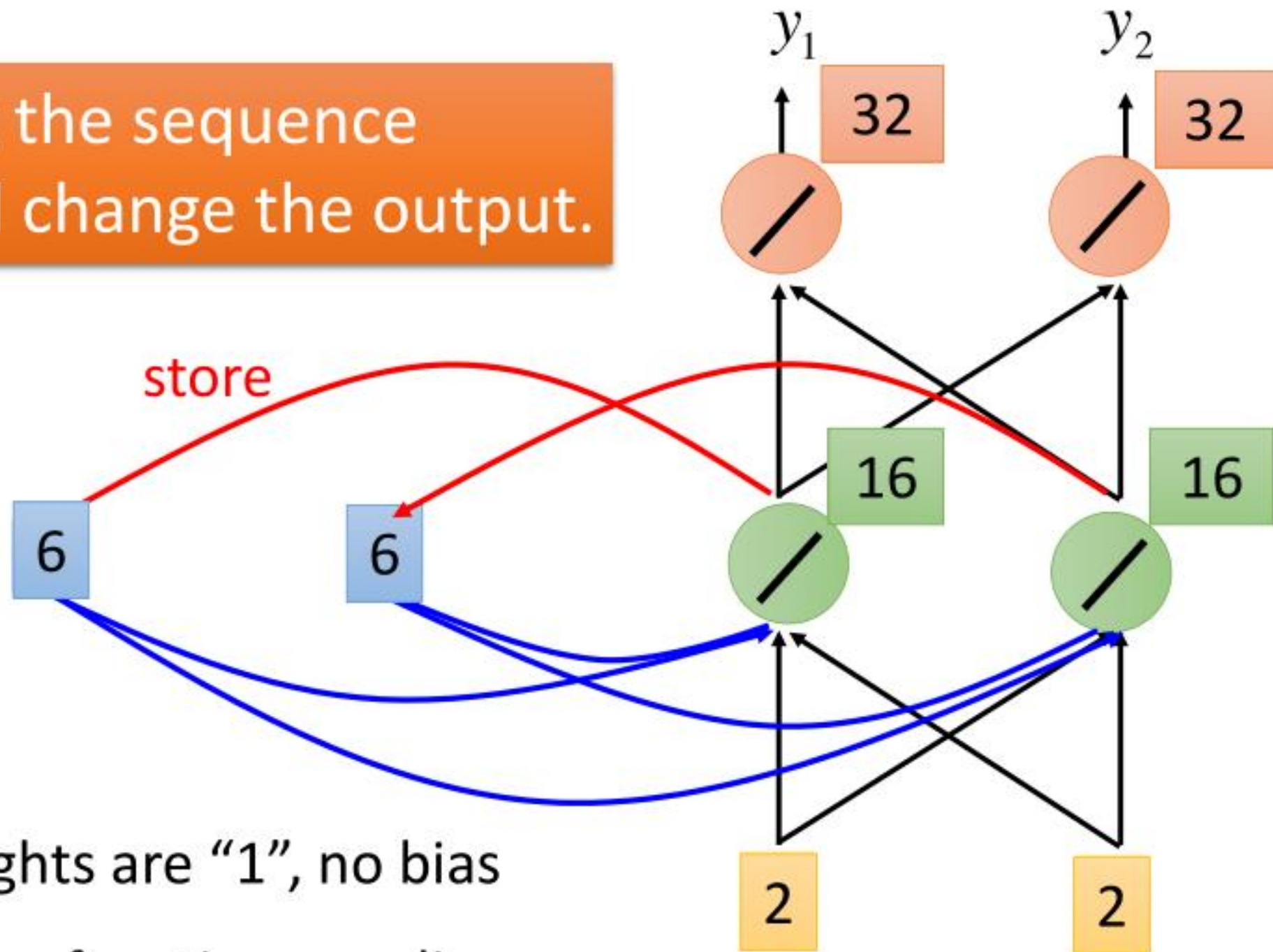


# Example

Input sequence:  $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} \dots$

output sequence:  $\begin{bmatrix} 4 \\ 4 \end{bmatrix} \begin{bmatrix} 12 \\ 12 \end{bmatrix} \begin{bmatrix} 32 \\ 32 \end{bmatrix}$

Changing the sequence order will change the output.



All the weights are "1", no bias

All activation functions are linear

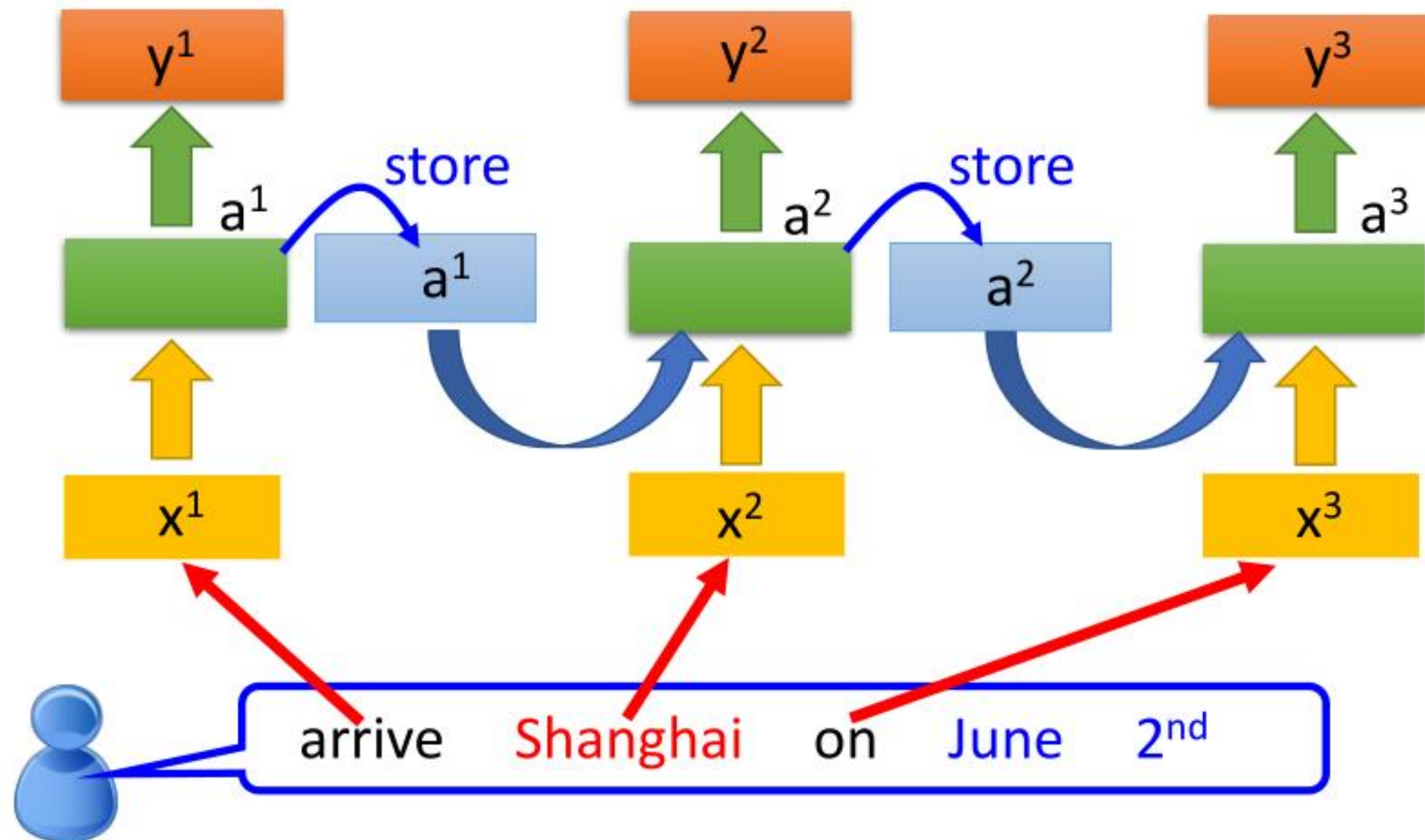
# RNN

The same network is used again and again.

Probability of  
“arrive” in each slot

Probability of  
“**Shanghai**” in each slot

Probability of  
“on” in each slot





# RNN

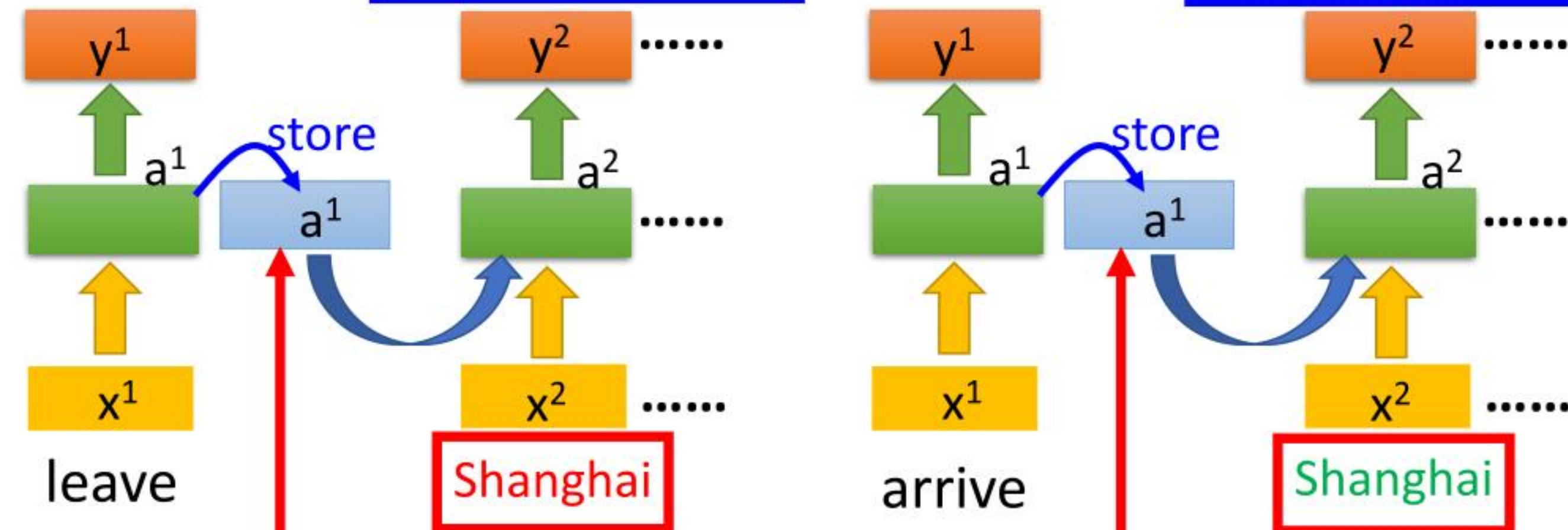
Different

Prob of "leave"  
in each slot

Prob of "**Shanghai**"  
in each slot

Prob of "arrive"  
in each slot

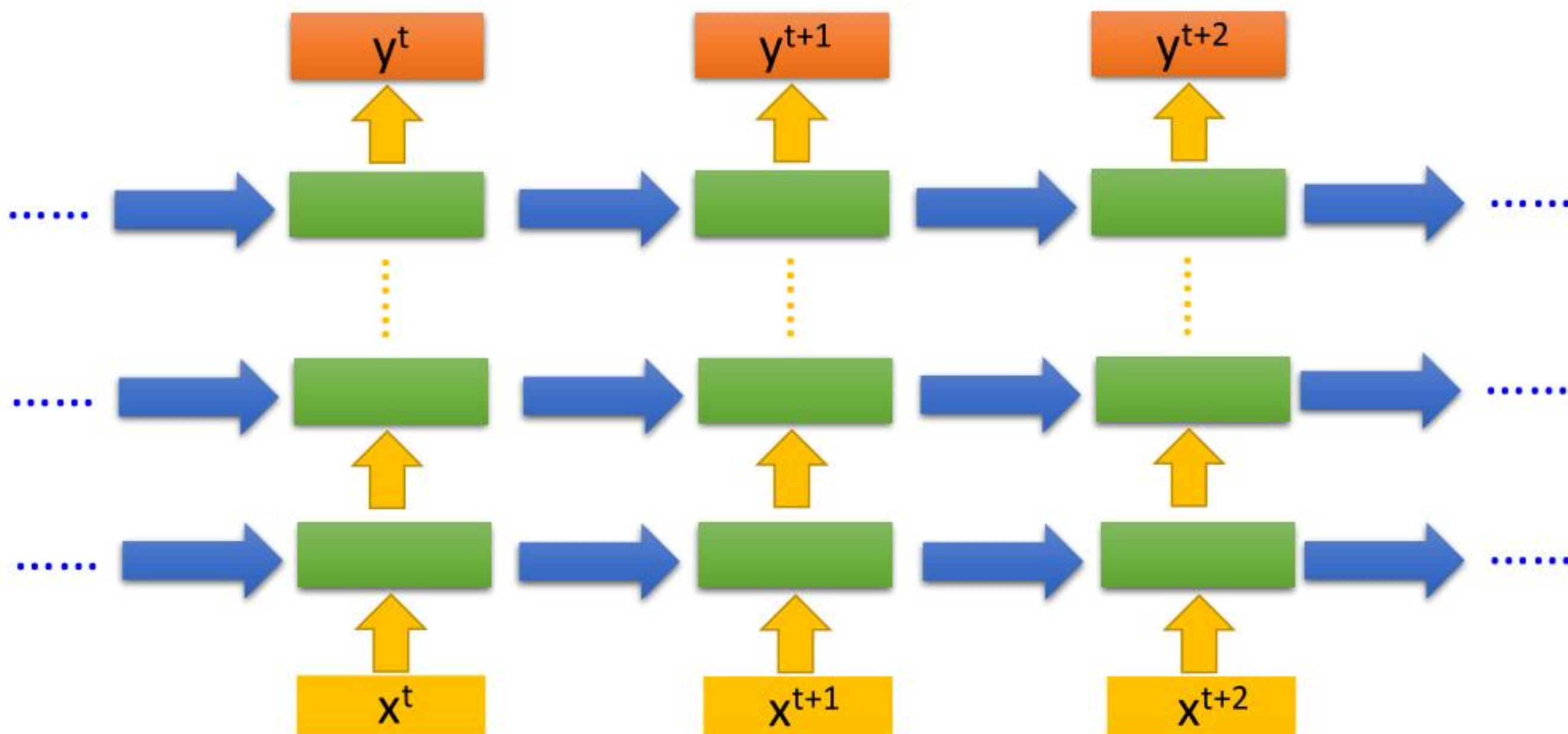
Prob of "**Shanghai**"  
in each slot



The values stored in the memory is different.

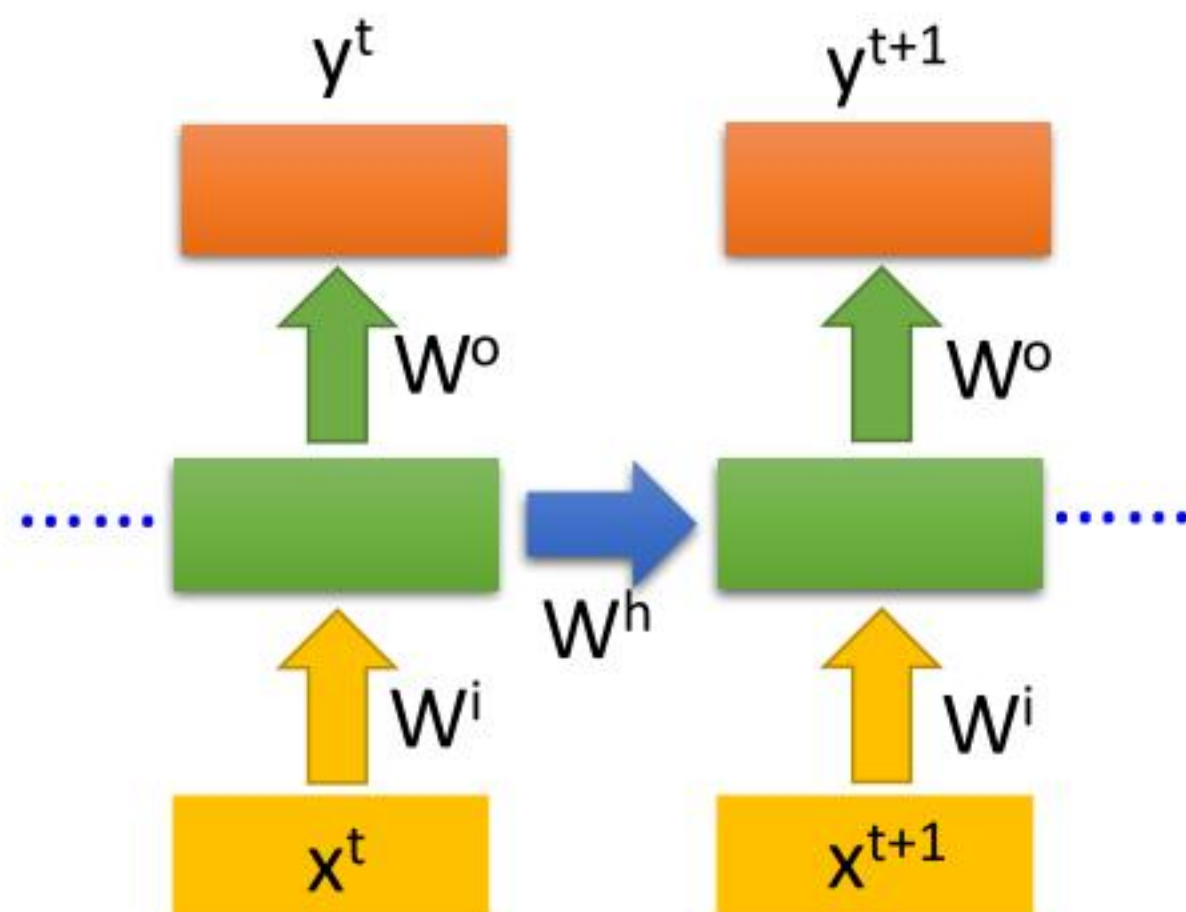


Of course it can be deep ...

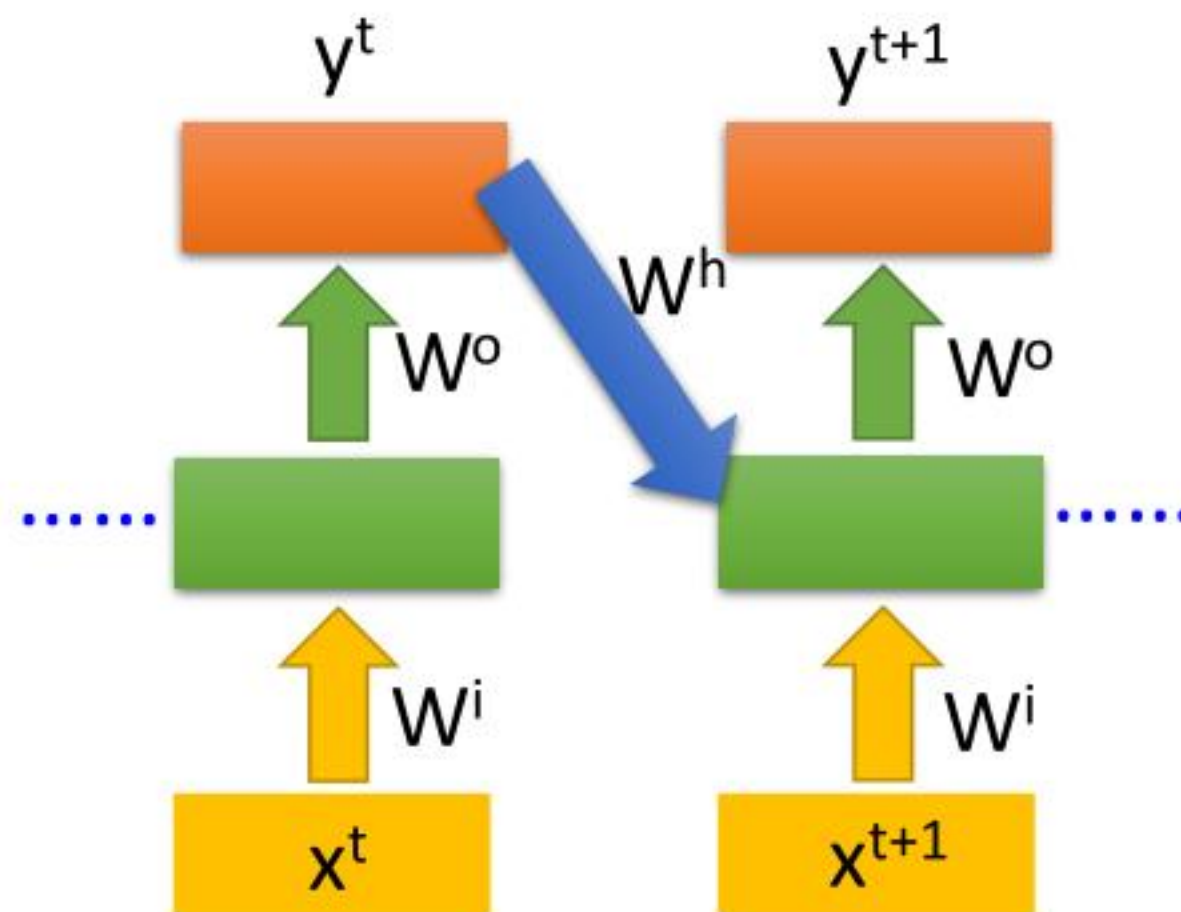


# Elman Network & Jordan Network

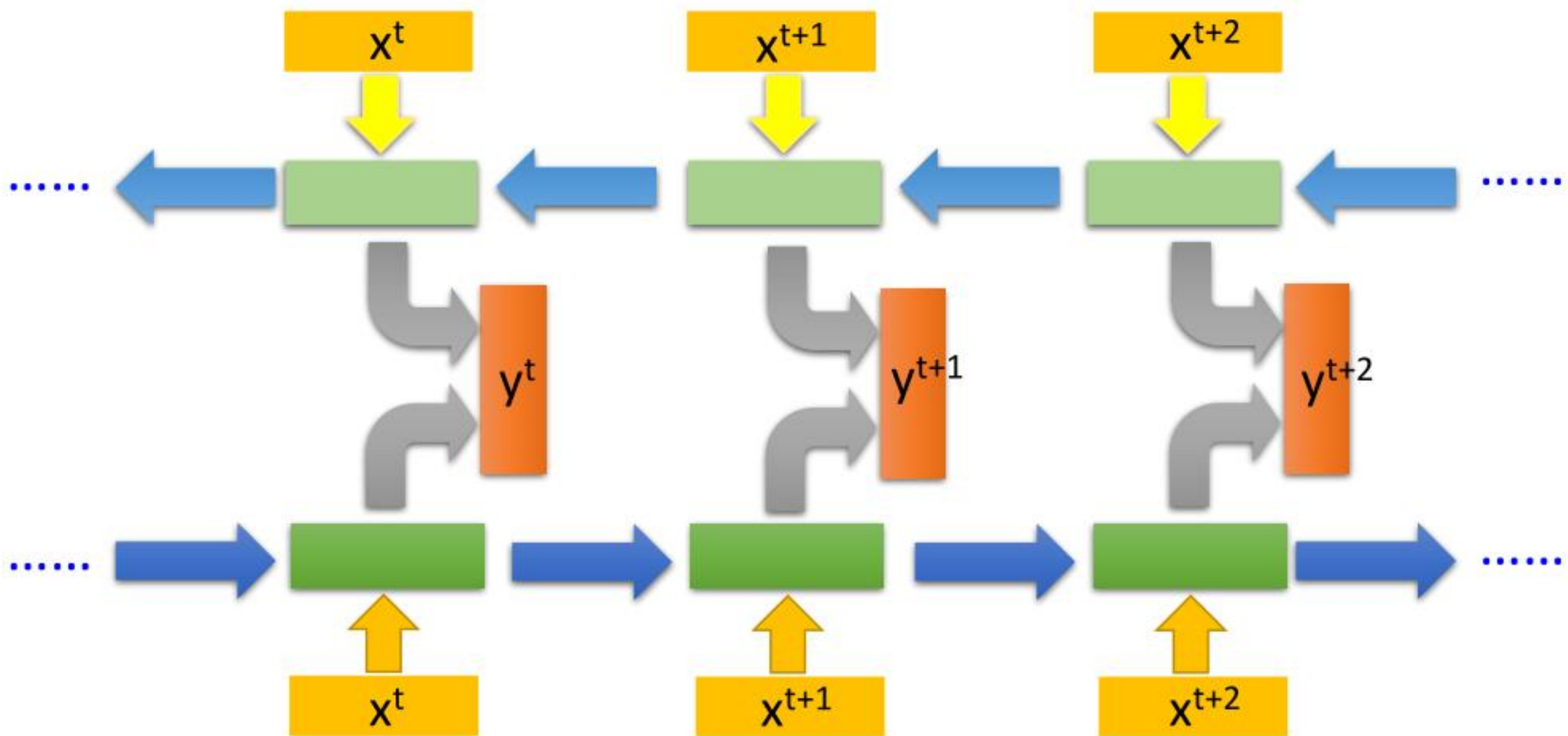
*Elman Network*



*Jordan Network*

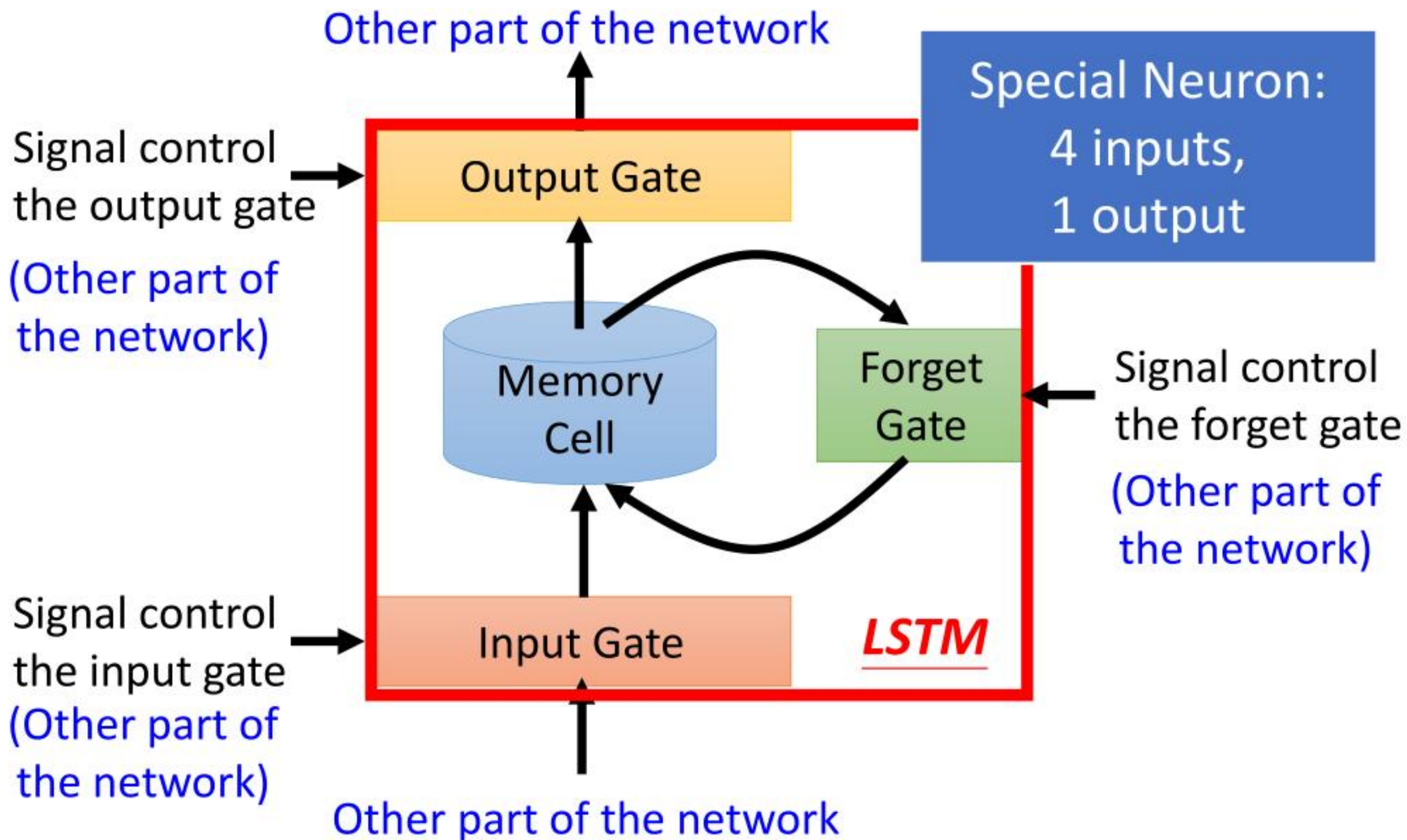


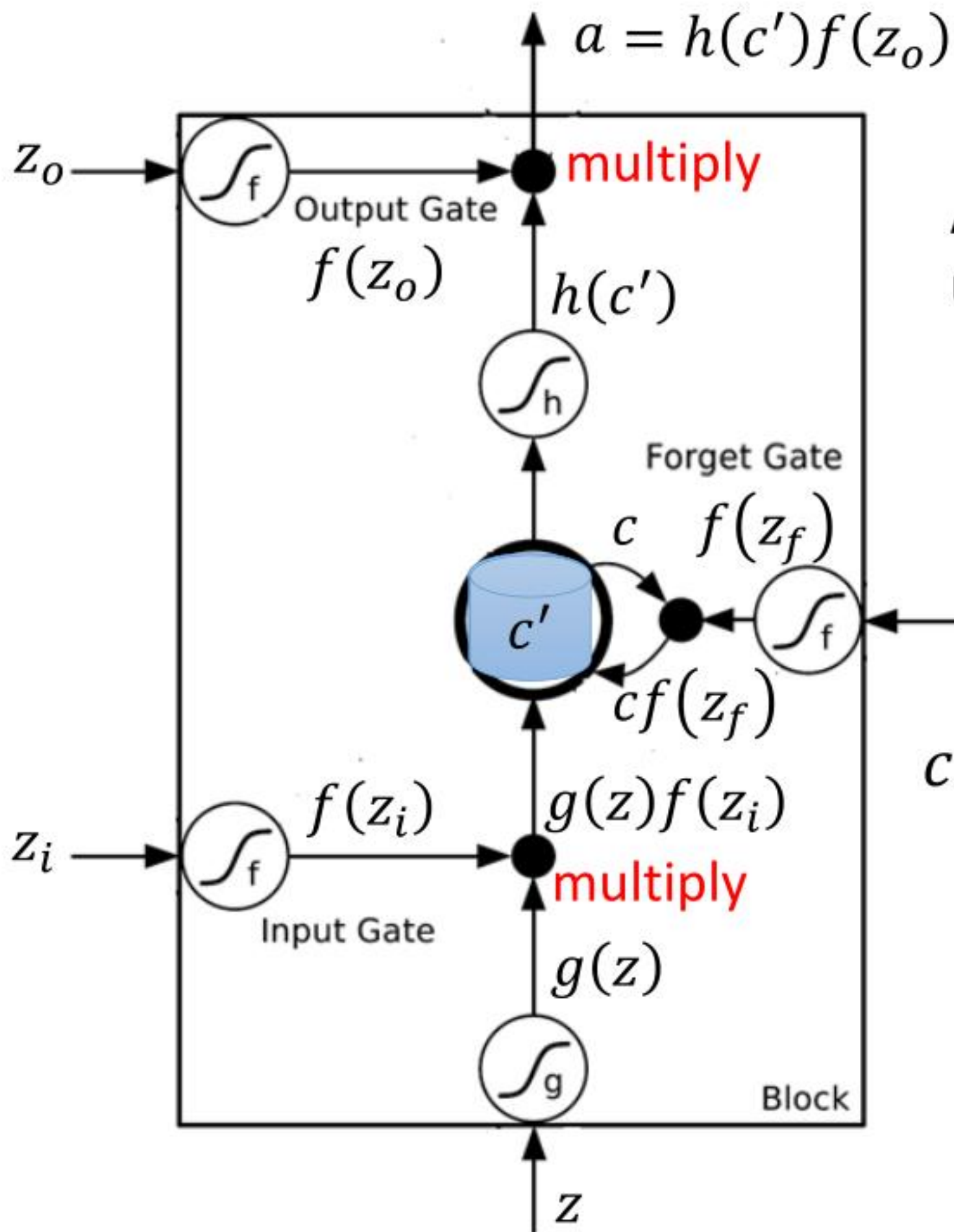
# Bidirectional RNN





# Long Short-term Memory (LSTM)



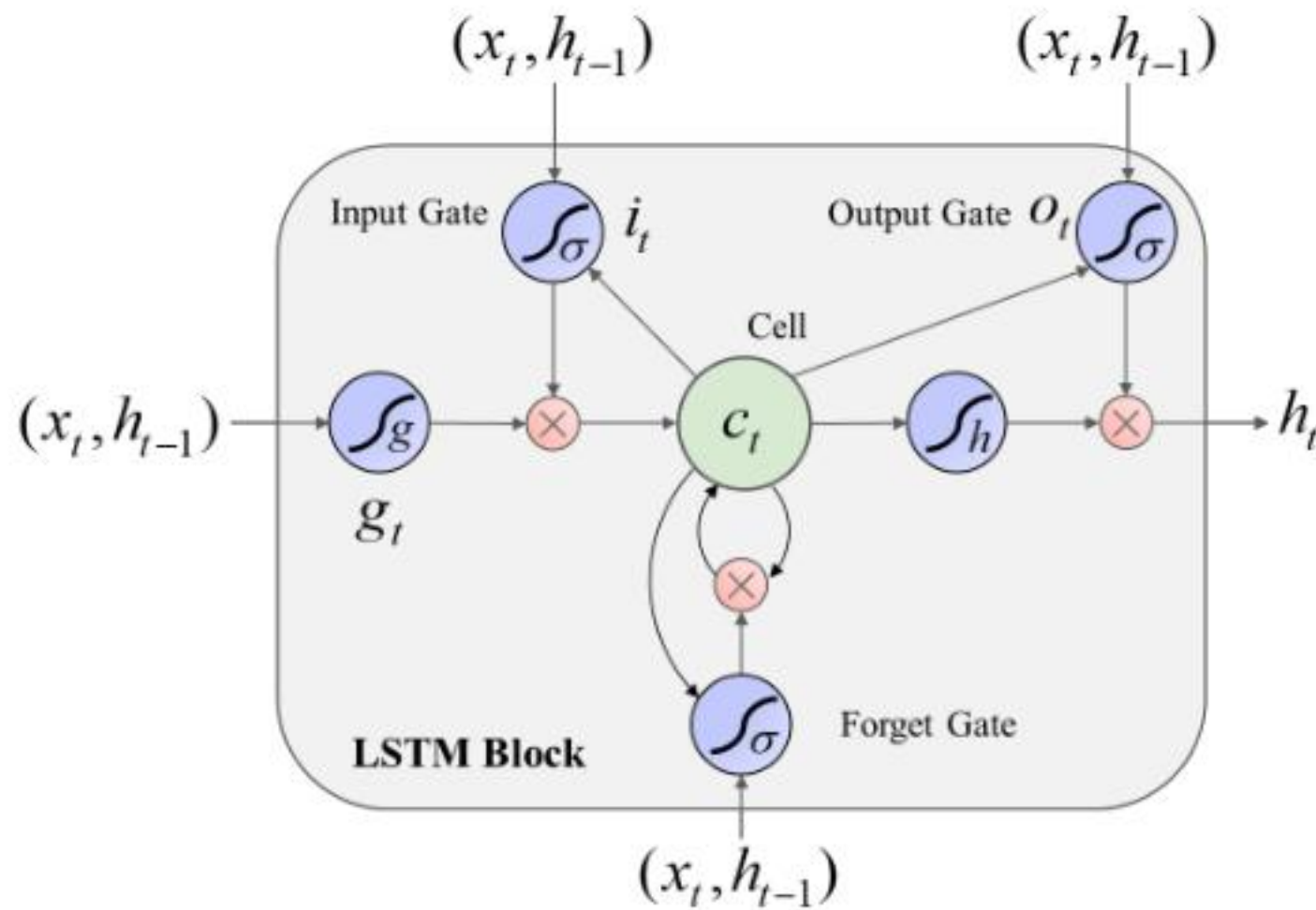


Activation function  $f$  is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

$$c' = g(z)f(z_i) + cf(z_f)$$



$$f_t = \text{sigmoid}(W_{f,x}x_t + W_{f,h}h_{t-1} + b_f)$$

$$i_t = \text{sigmoid}(W_{i,x}x_t + W_{i,h}h_{t-1} + b_i)$$

$$g_t = \tanh(W_{g,x}x_t + W_{g,h}h_{t-1} + b_g)$$

$$c_t = g_t \cdot i_t + c_{t-1} \cdot f_t$$

$$o_t = \text{sigmoid}(W_{o,x}x_t + W_{o,h}h_{t-1} + b_o)$$

$$h_t = \tanh(c_t) \cdot o_t$$



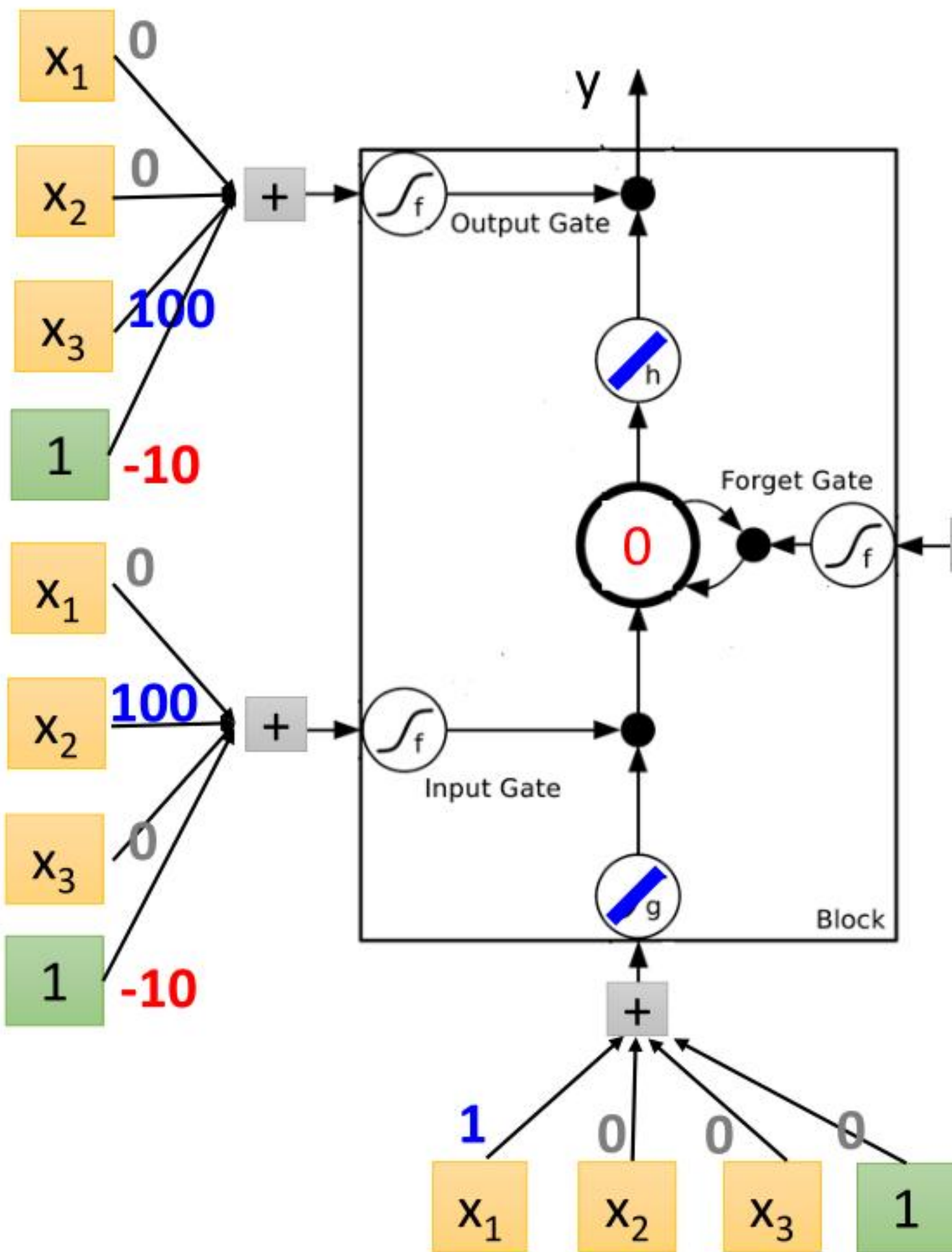
# LSTM - Example

	0	0	3	3	7	7	7	0	6
$x_1$	1	3	2	4	2	1	3	6	1
$x_2$	0	1	0	1	0	0	-1	1	0
$x_3$	0	0	0	0	0	1	0	0	1
$y$	0	0	0	0	0	7	0	0	6

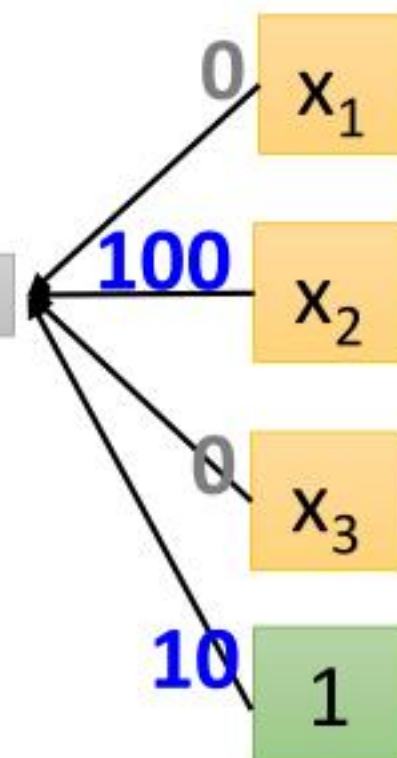
When  $x_2 = 1$ , add the numbers of  $x_1$  into the memory

When  $x_2 = -1$ , reset the memory

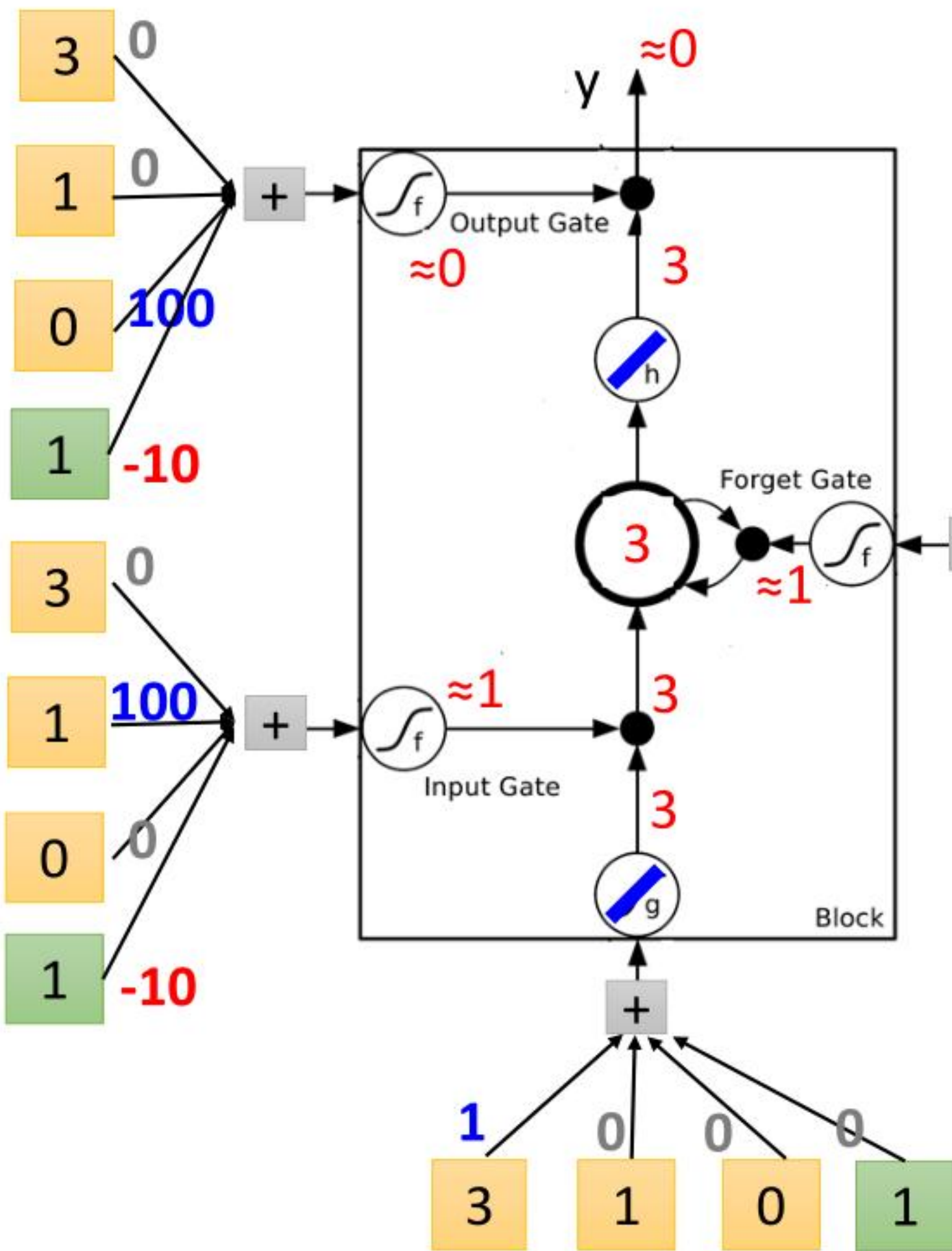
When  $x_3 = 1$ , output the number in the memory.



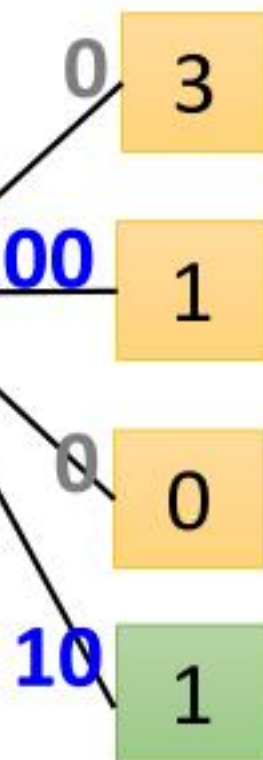
$y$  0 0 0 7 0



	3	4	2	1	3
$x_1$	3	4	2	1	3
$x_2$	1	1	0	0	-1
$x_3$	0	0	0	1	0

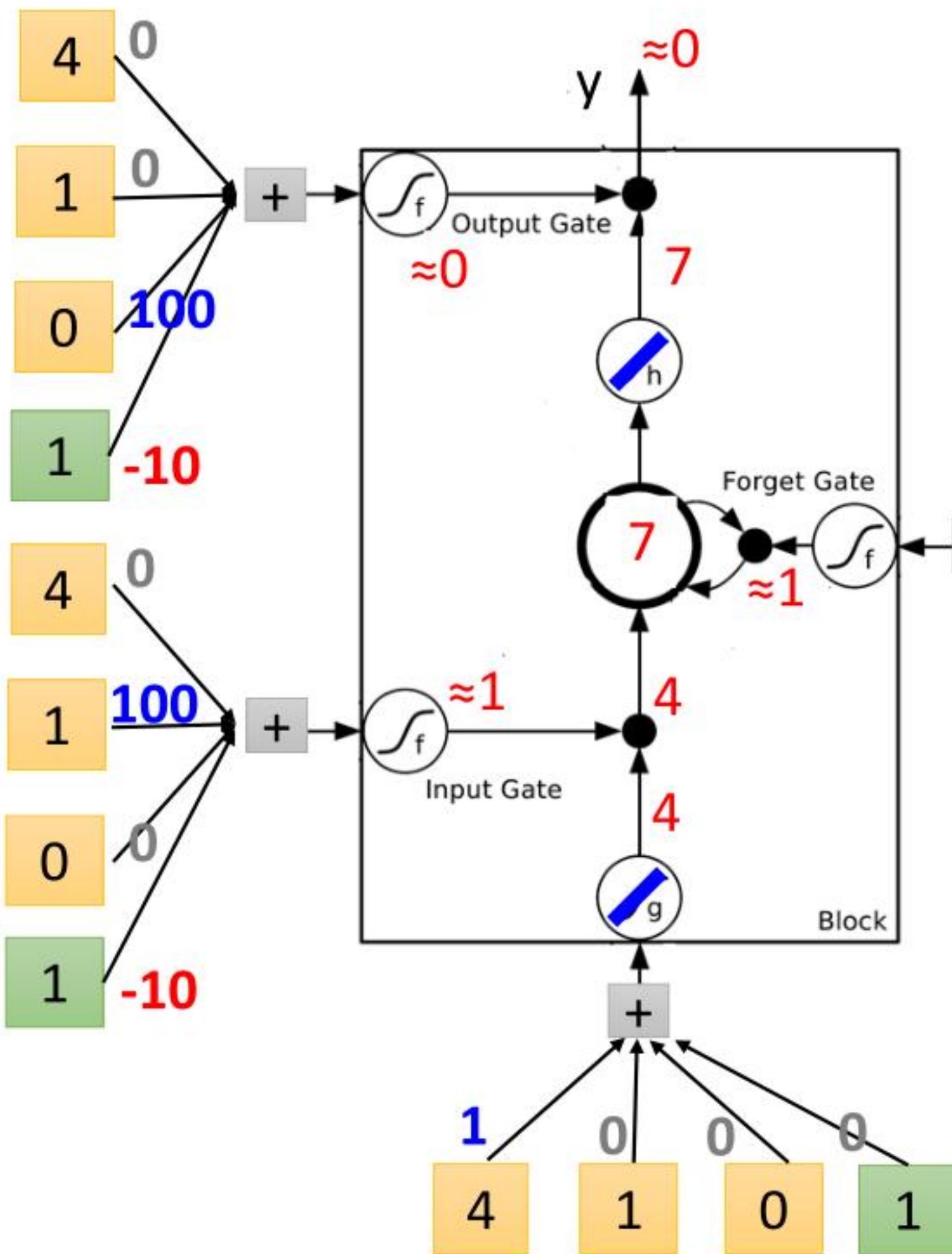


$y$  0 0 0 7 0

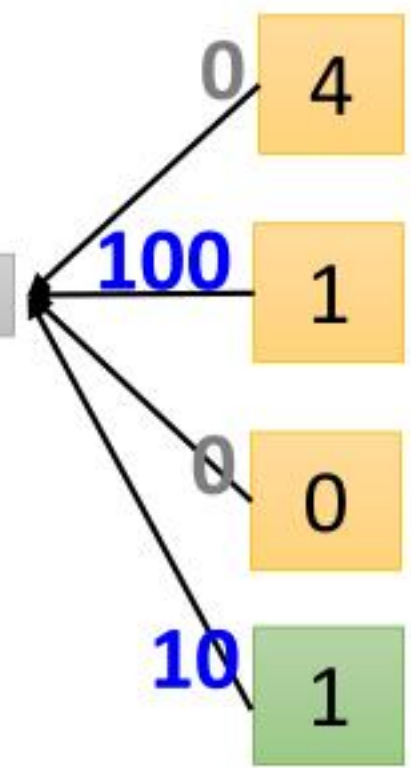


$x_1$  3 4 2 1 3  
 $x_2$  1 1 0 0 -1  
 $x_3$  0 0 0 1 0

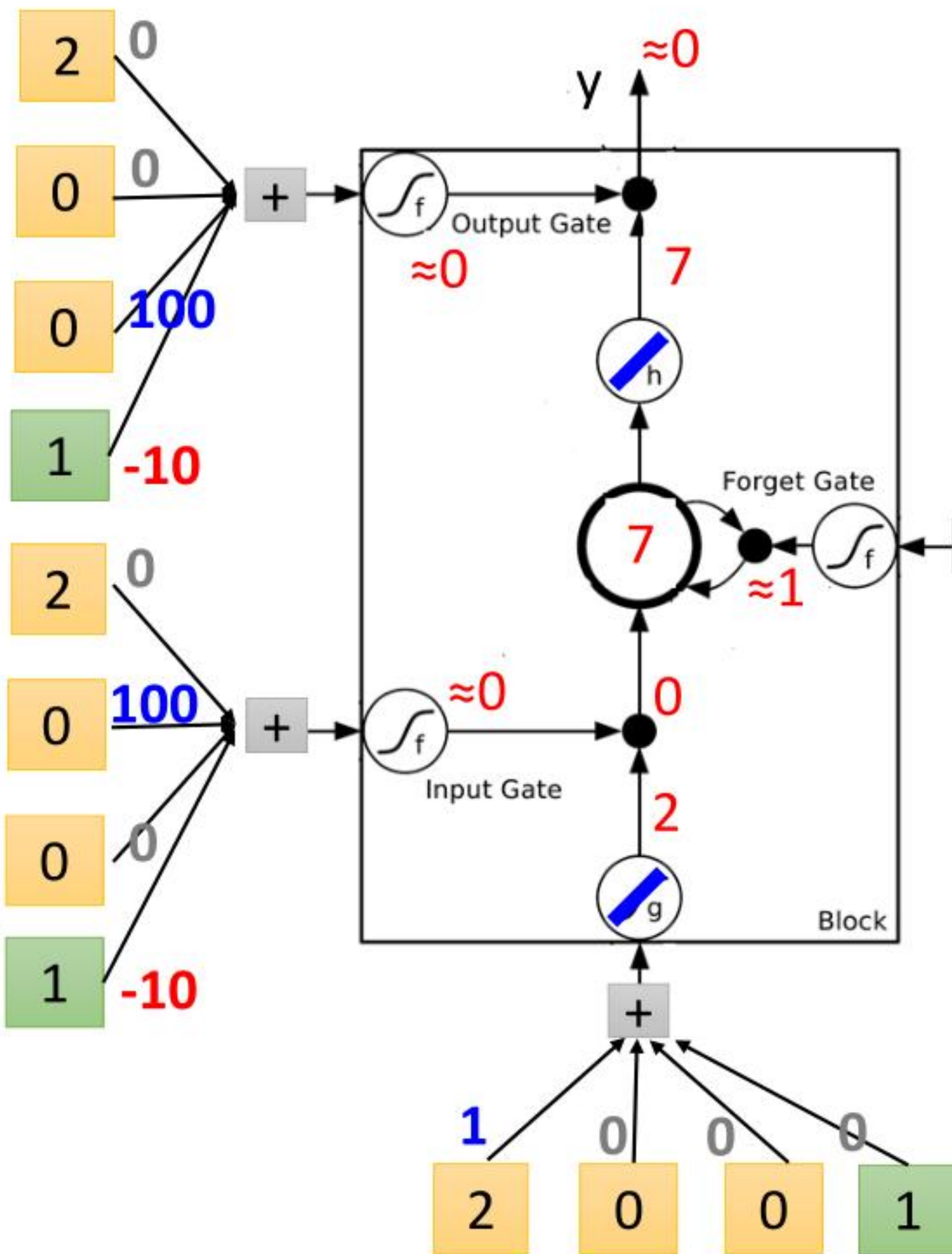




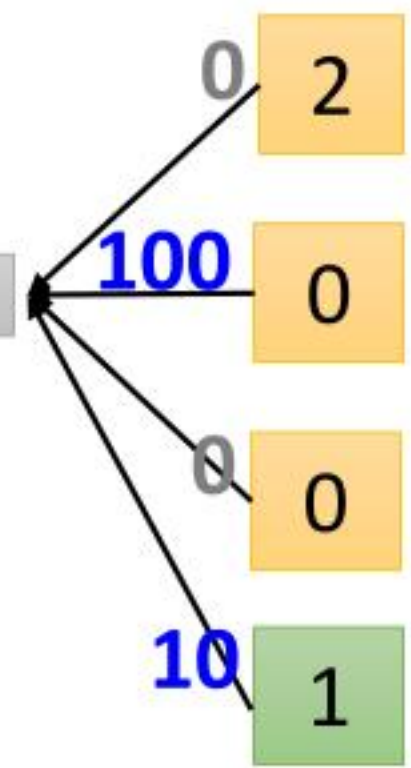
$y$  0 0 0 7 0



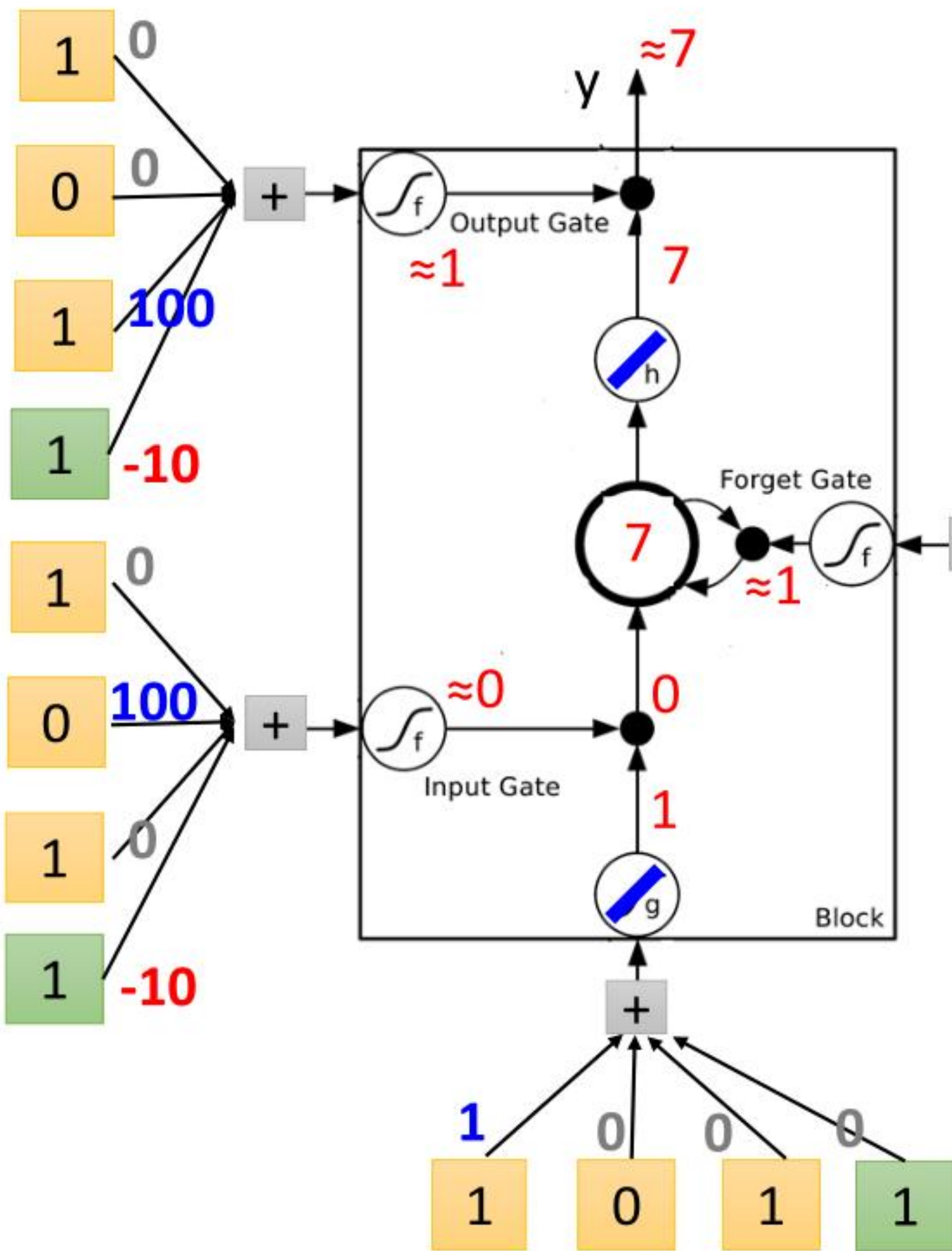
$x_1$  3 4 2 1 3  
 $x_2$  1 1 0 0 -1  
 $x_3$  0 0 0 1 0



$y$  0 0 0 7 0



	$x_1$		$x_2$		$x_3$
	3		4		2
	1		1		0
	0		0		0
	0		0		0
	1		0		1
	0		0		0

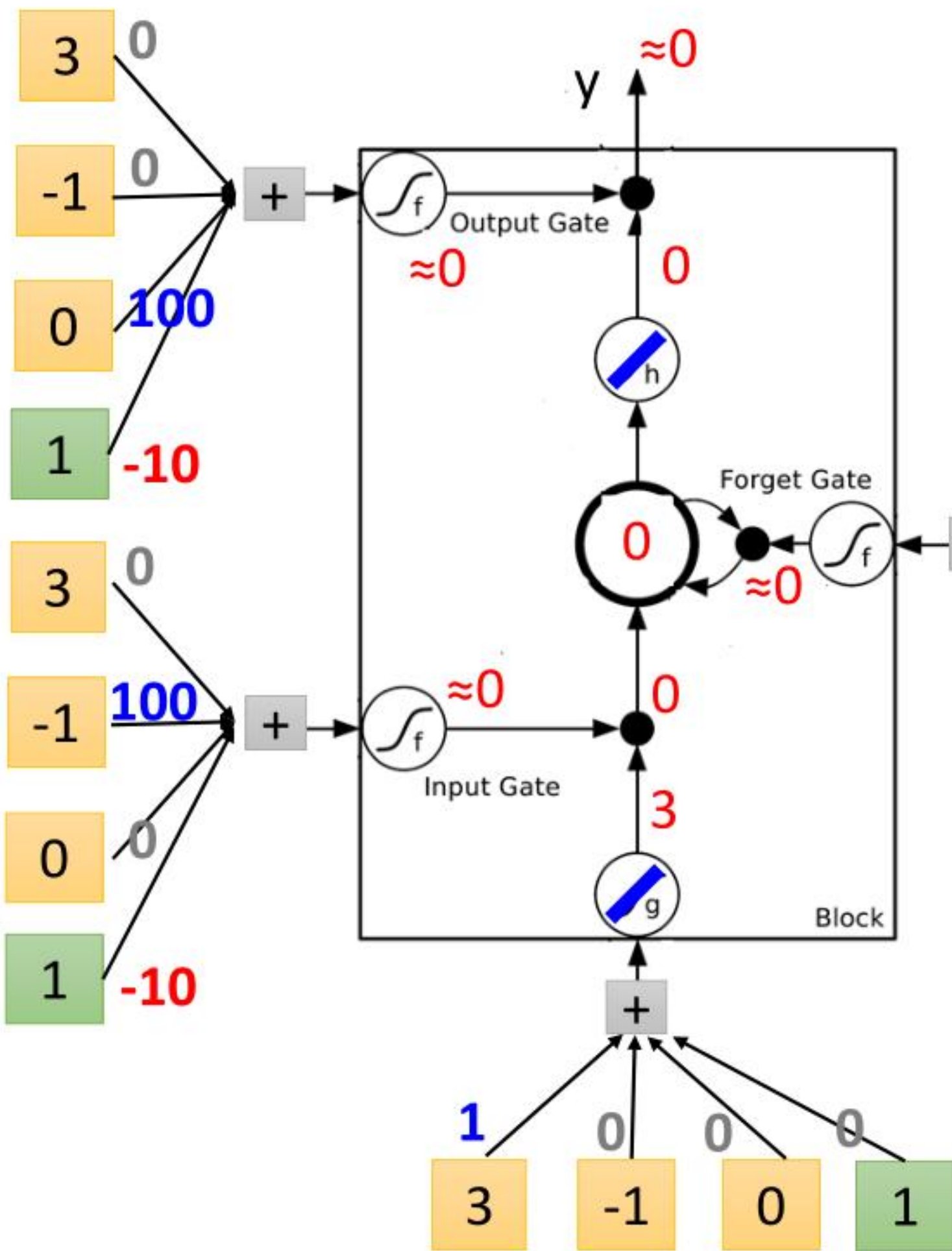


$y$  0 0 0 7 0

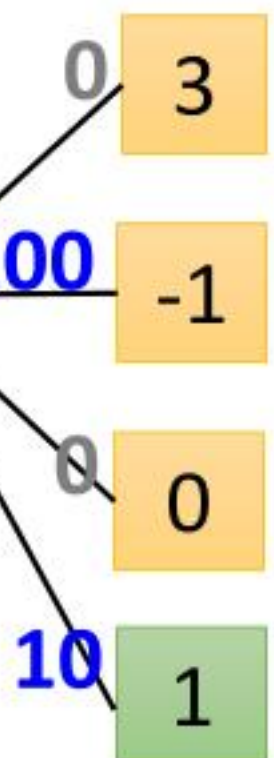
0 1  
100 0  
0 1  
10 1

$x_1$  3 4 2 1 3  
 $x_2$  1 1 0 0 -1  
 $x_3$  0 0 0 1 0





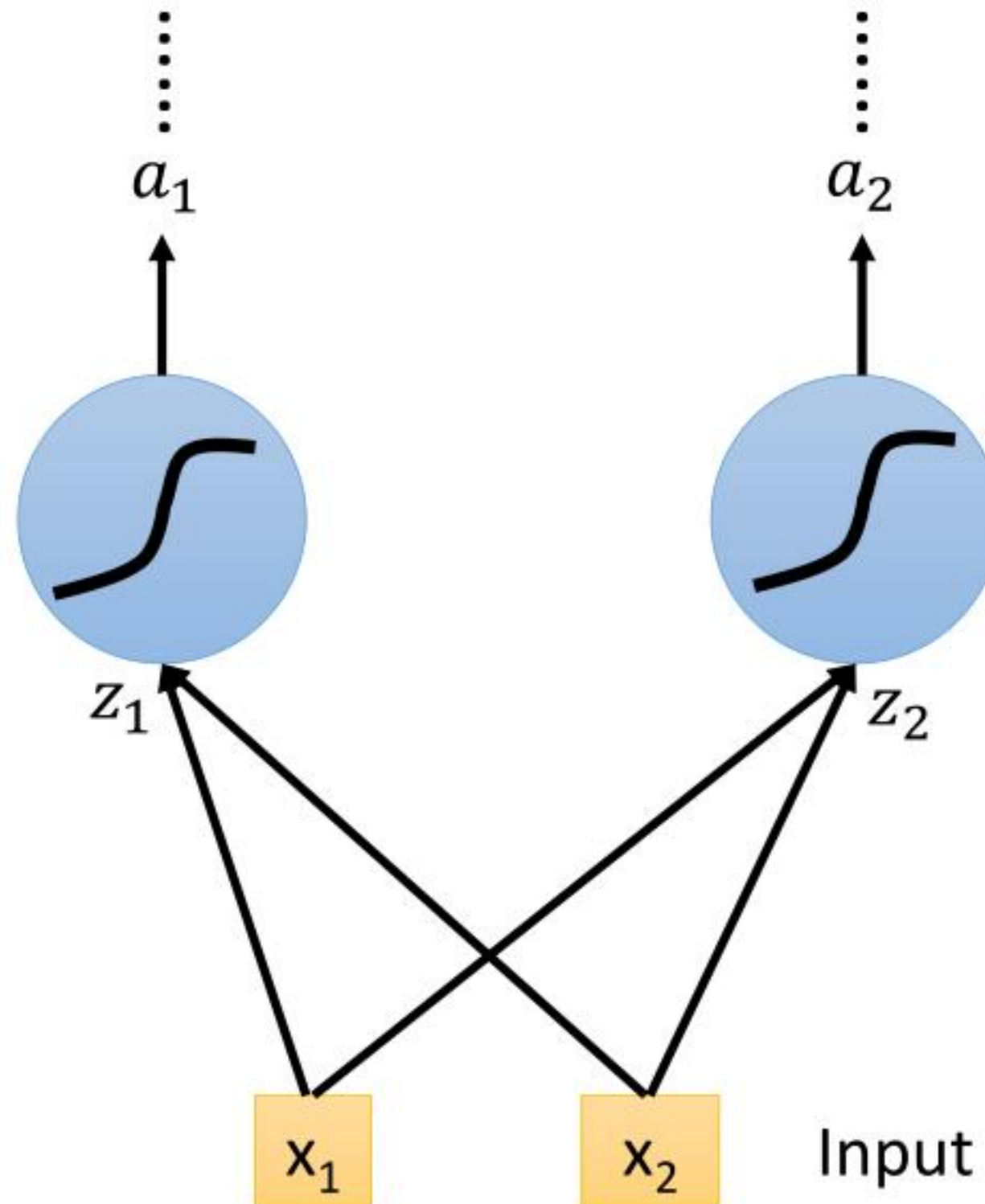
$y$  0 0 0 7 0

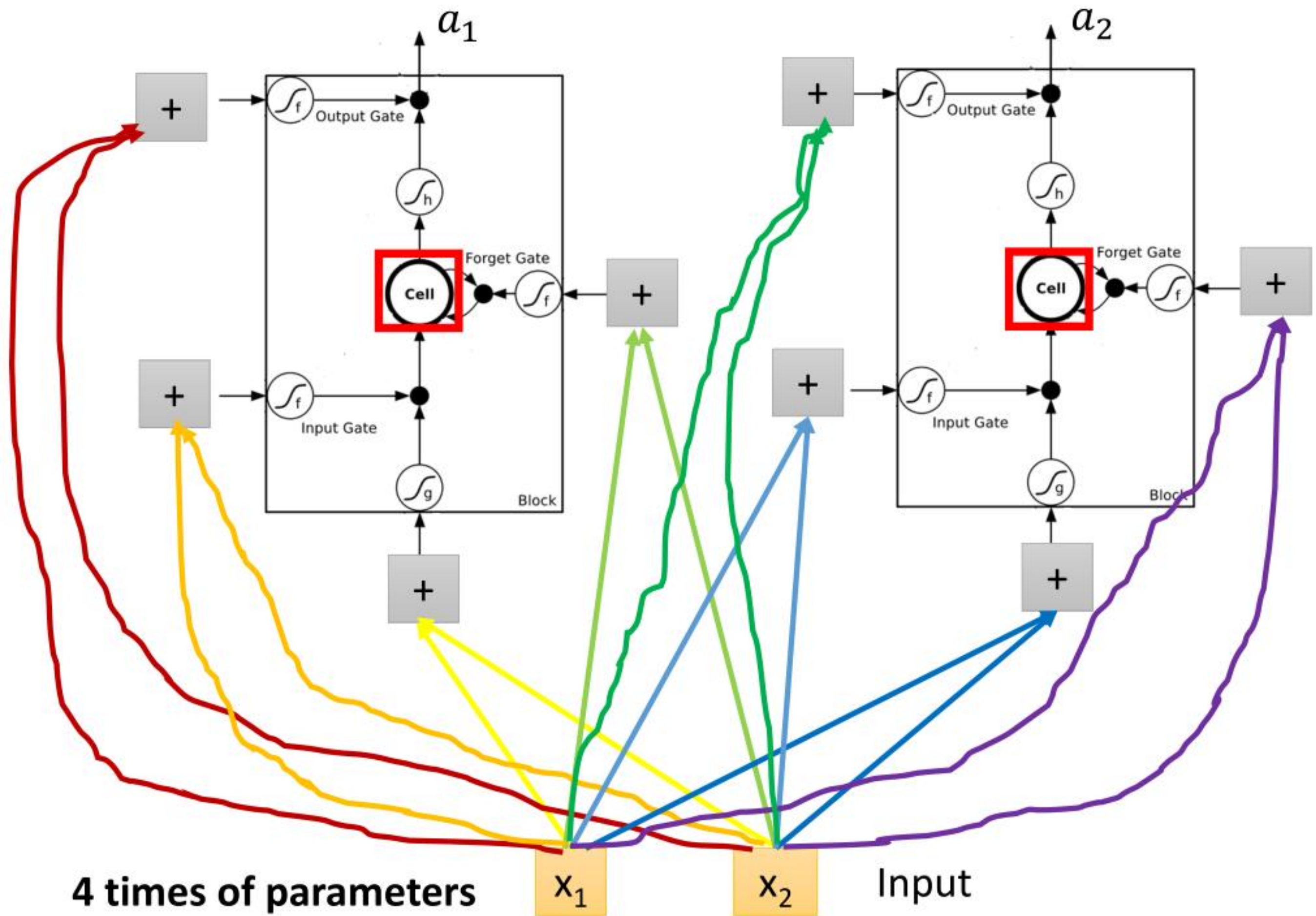


$x_1$  3 4 2 1 3  
 $x_2$  1 1 0 0 -1  
 $x_3$  0 0 0 1 0

## Original Network:

- Simply replace the neurons with LSTM



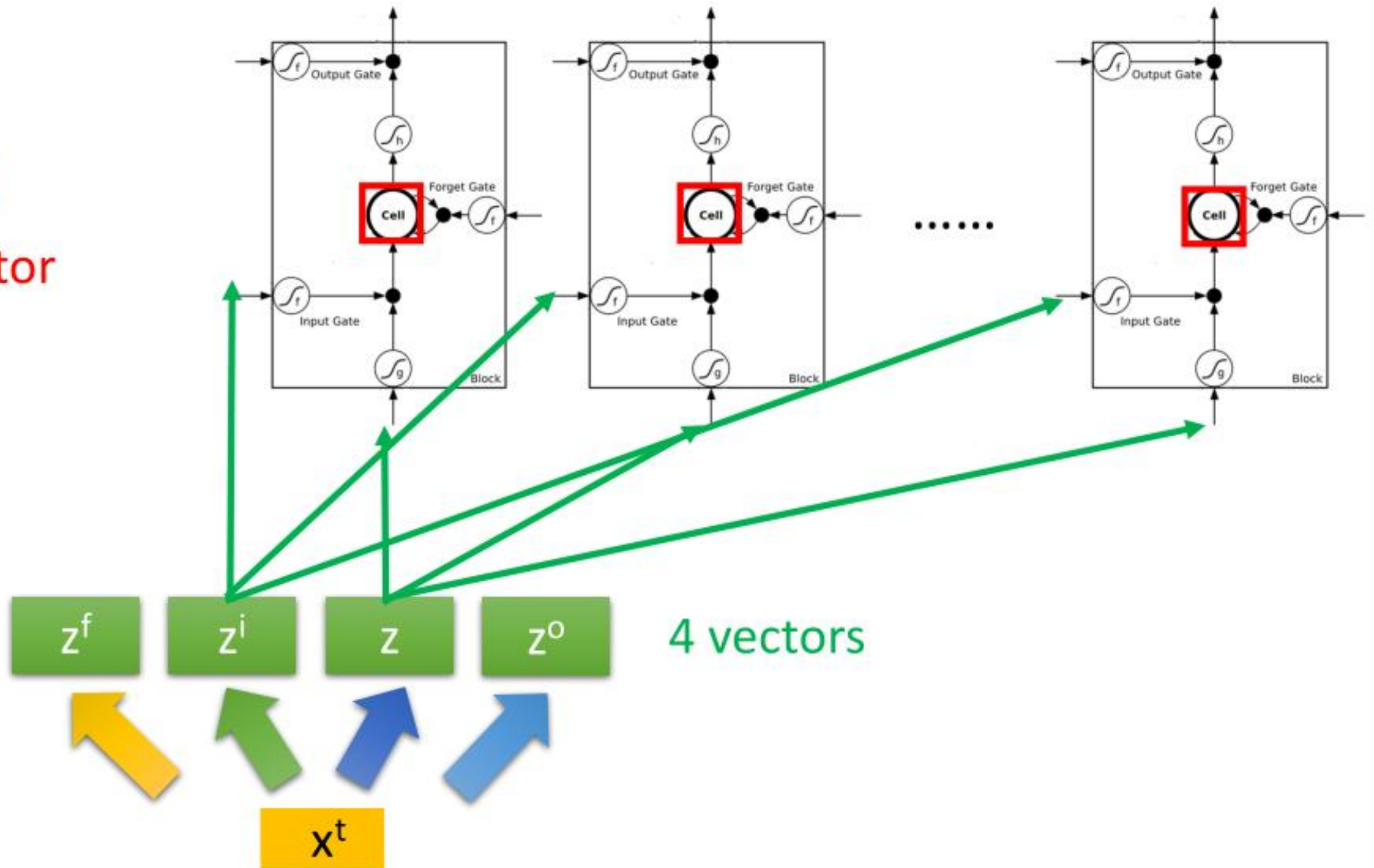




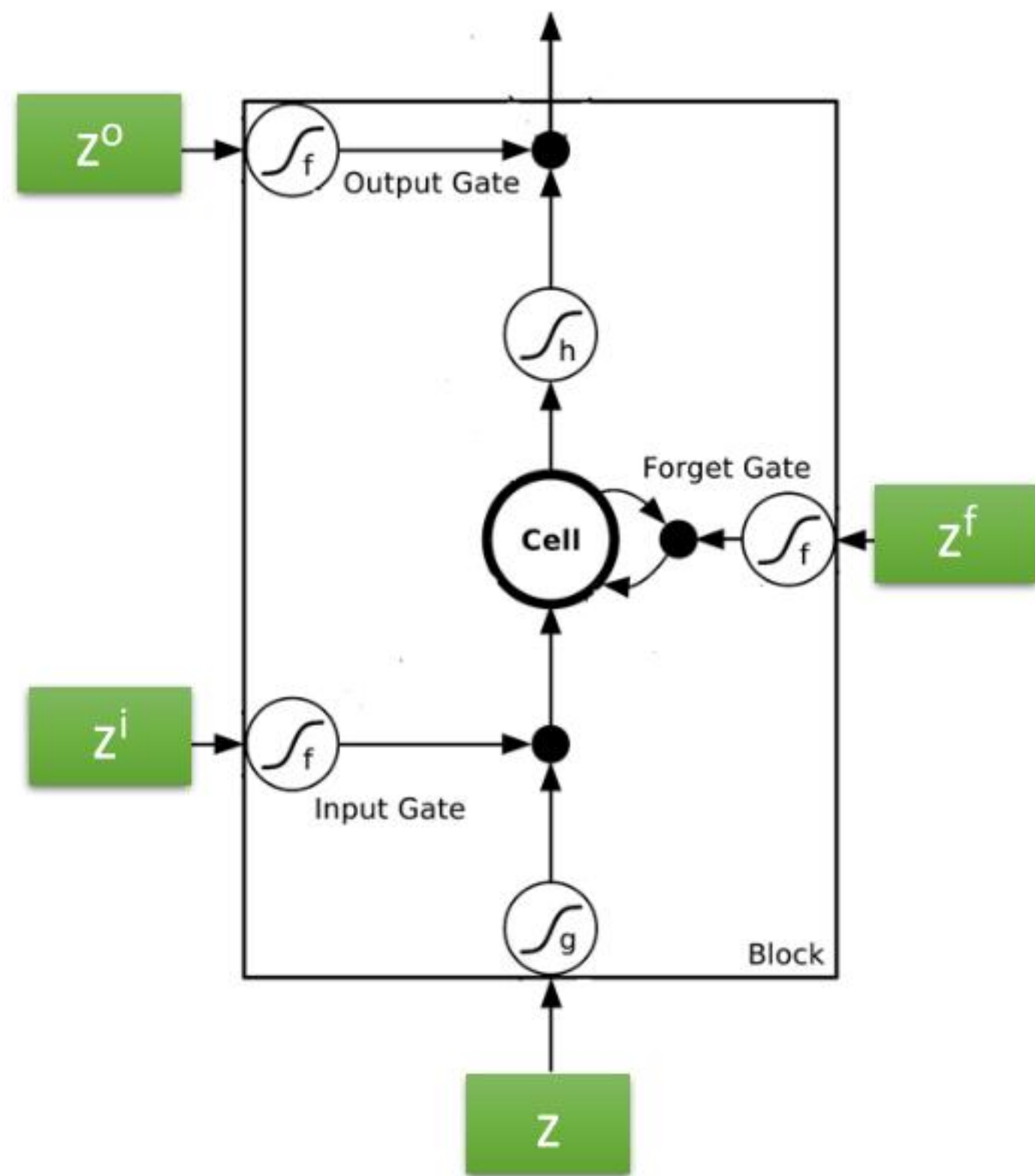
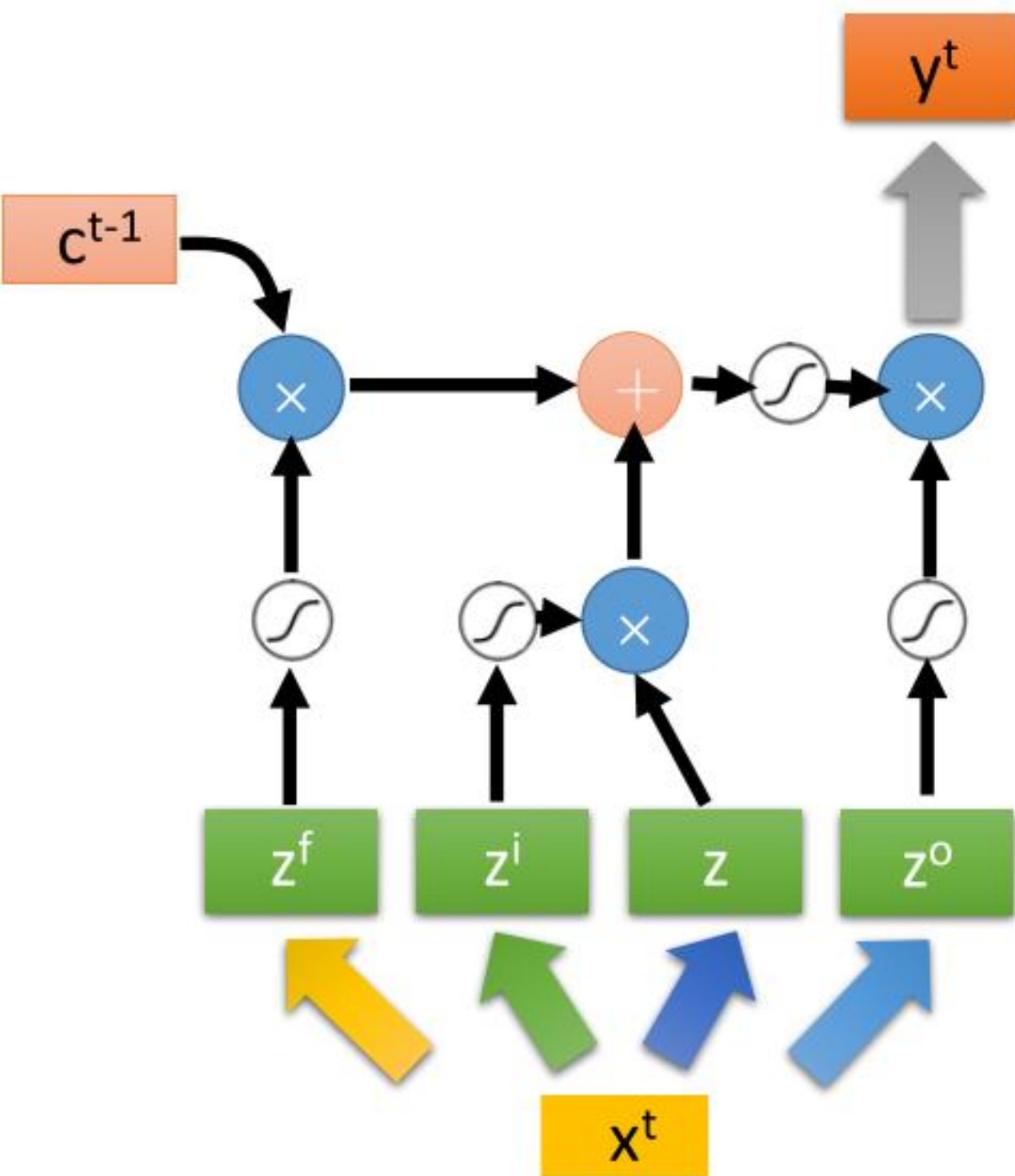
# LSTM

 $C^{t-1}$ 

vector

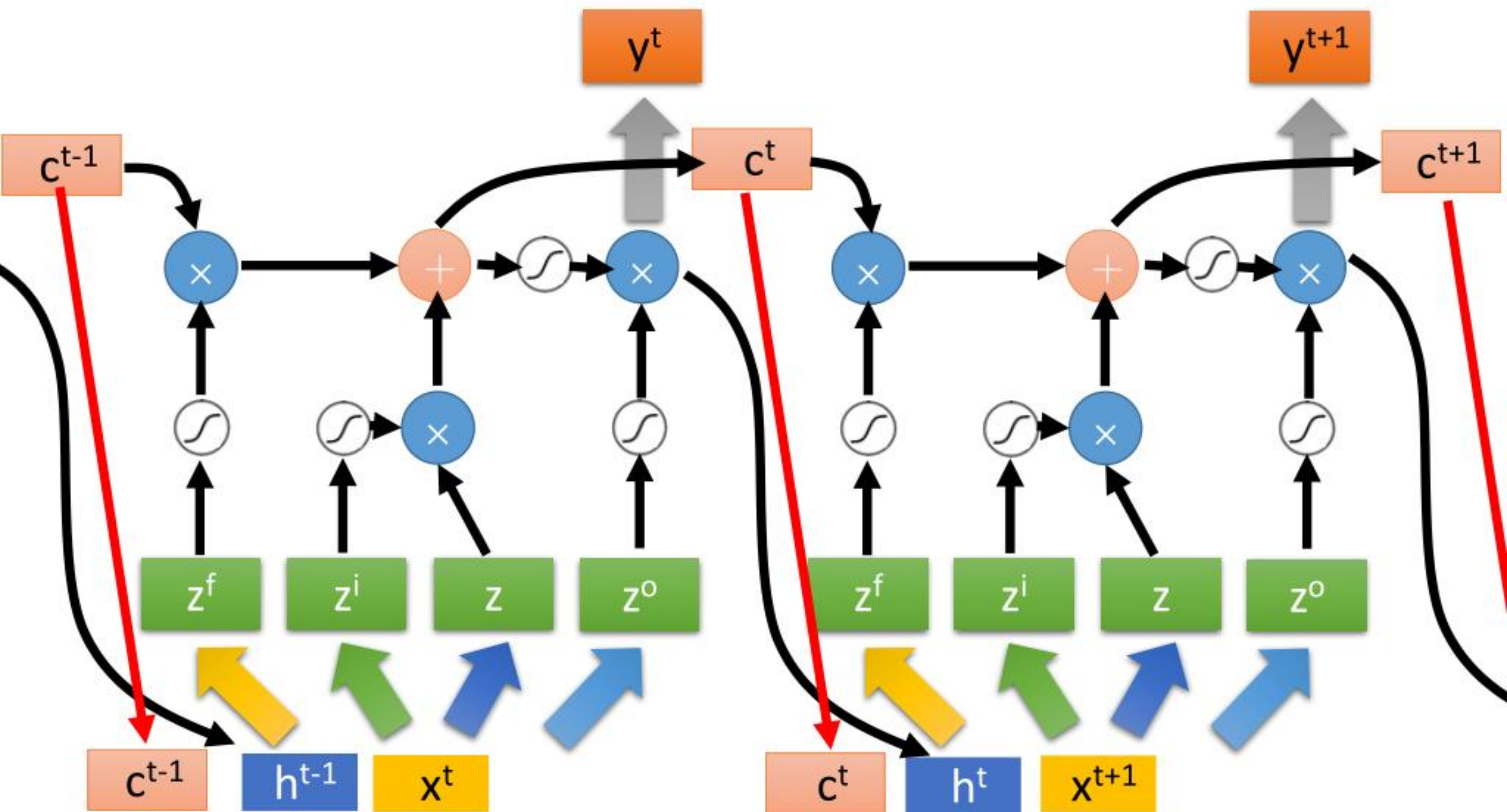


# LSTM



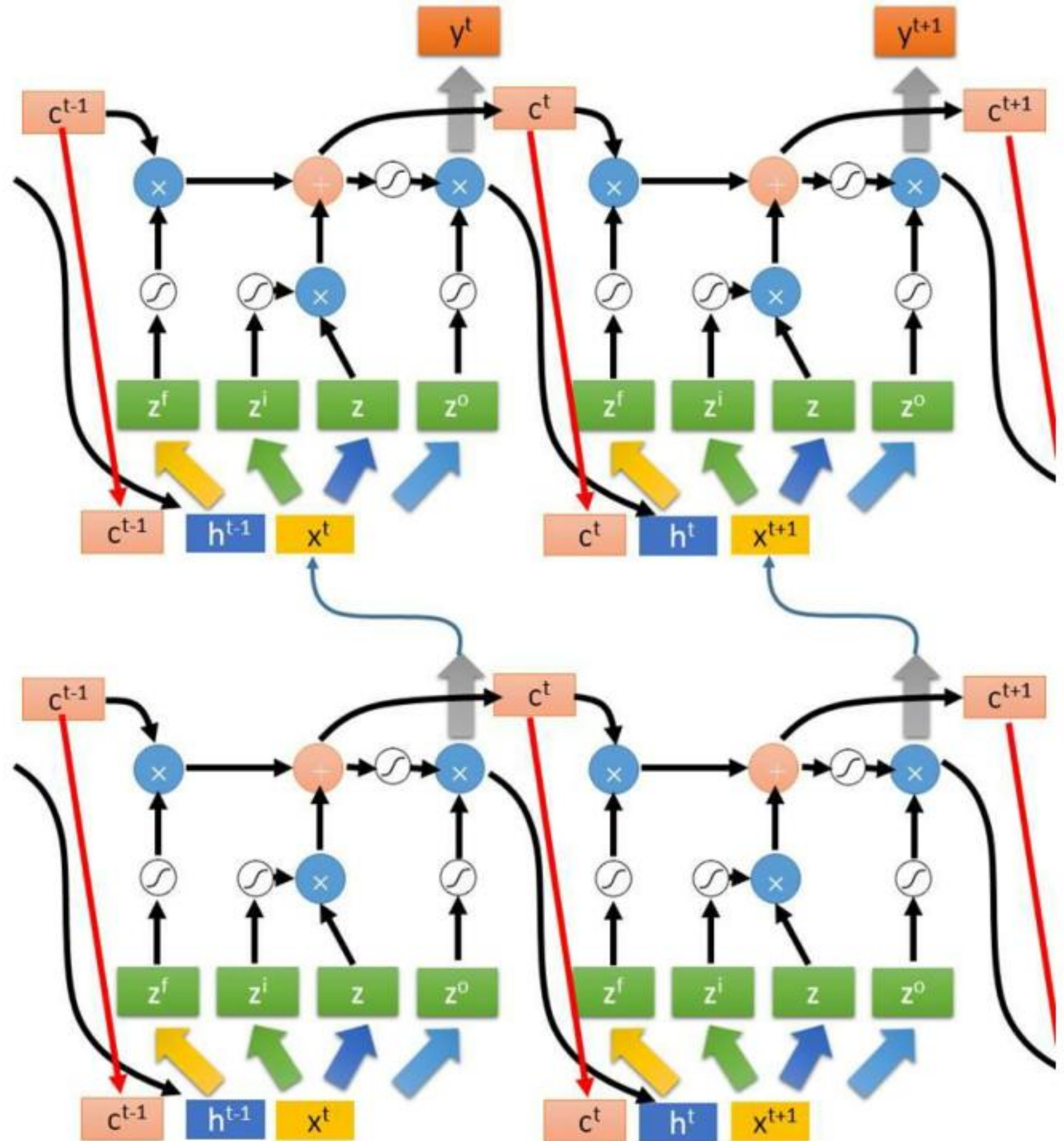
# LSTM

Extension: "peephole"



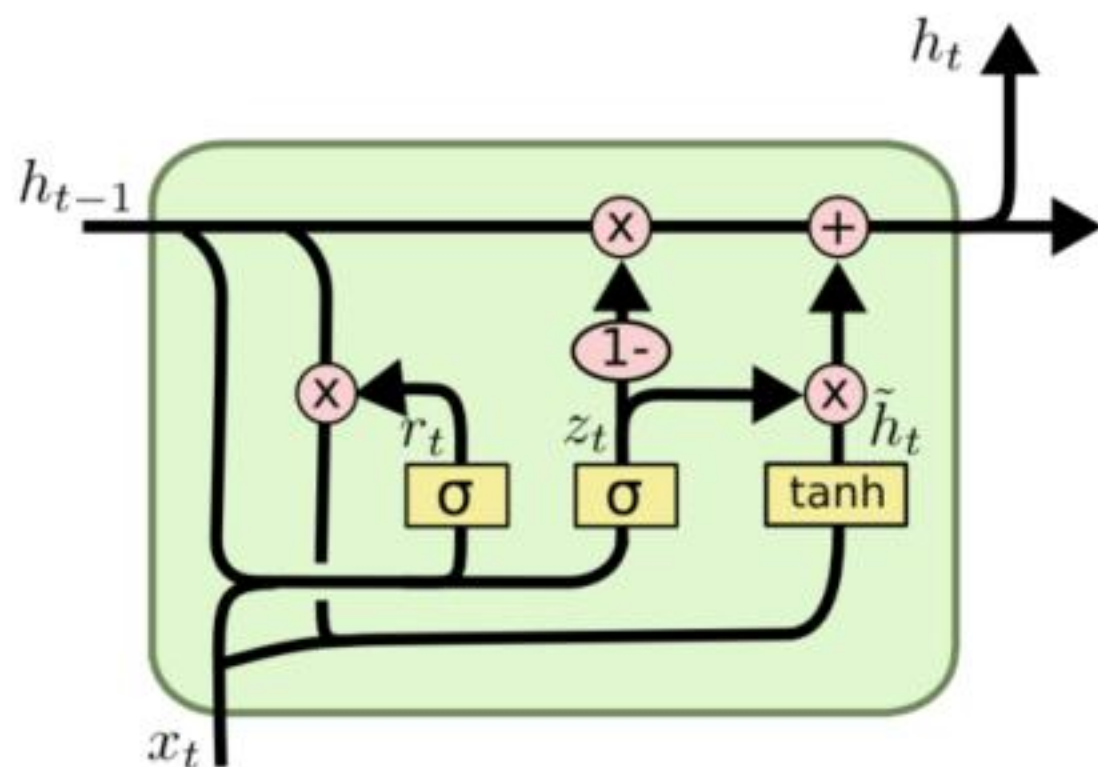


# Multiple-layer LSTM



This is quite  
standard now.

## Gated Recurrent Unit (GRU)



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

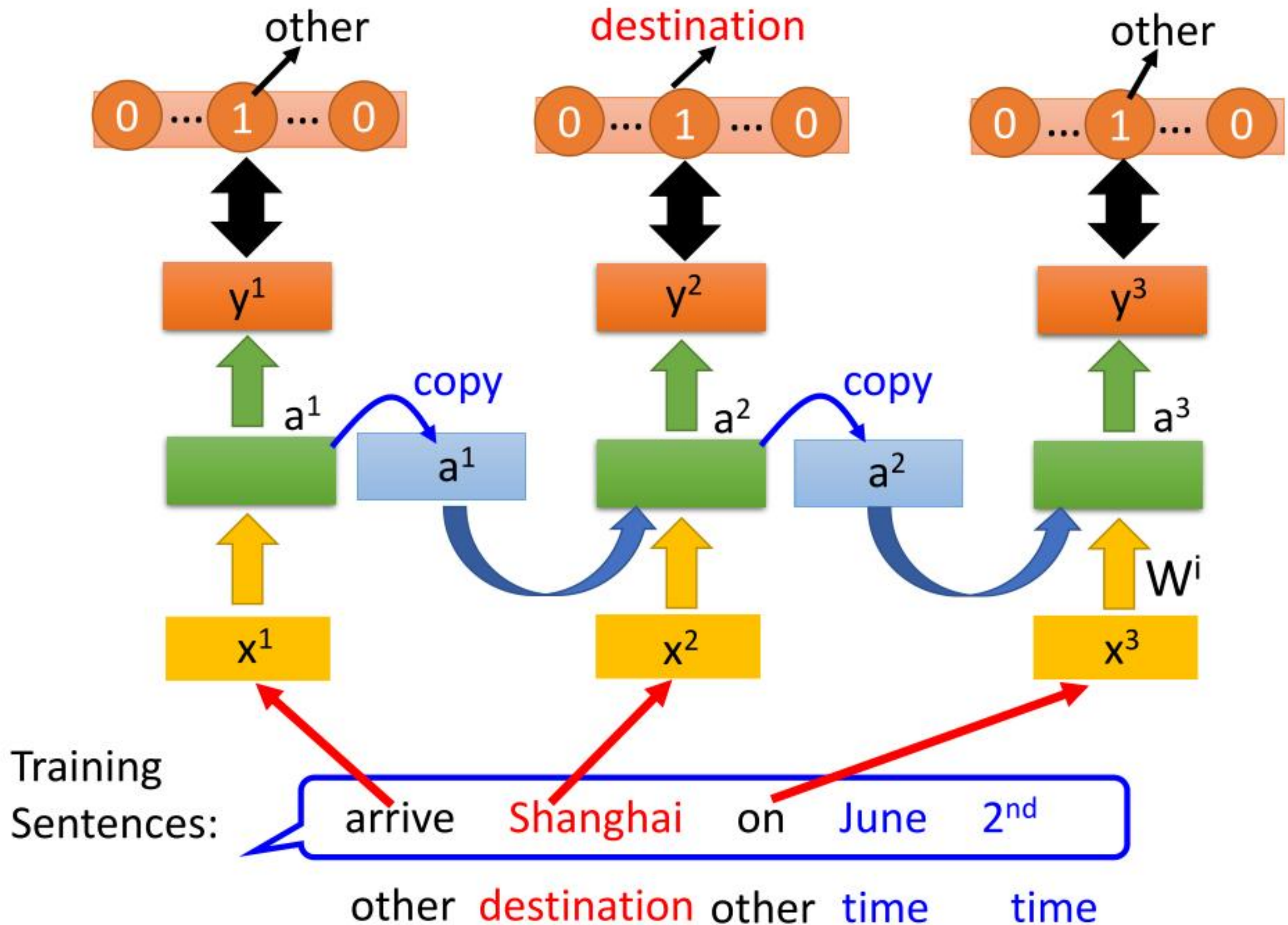
$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Only two gates: reset gate and update gate.

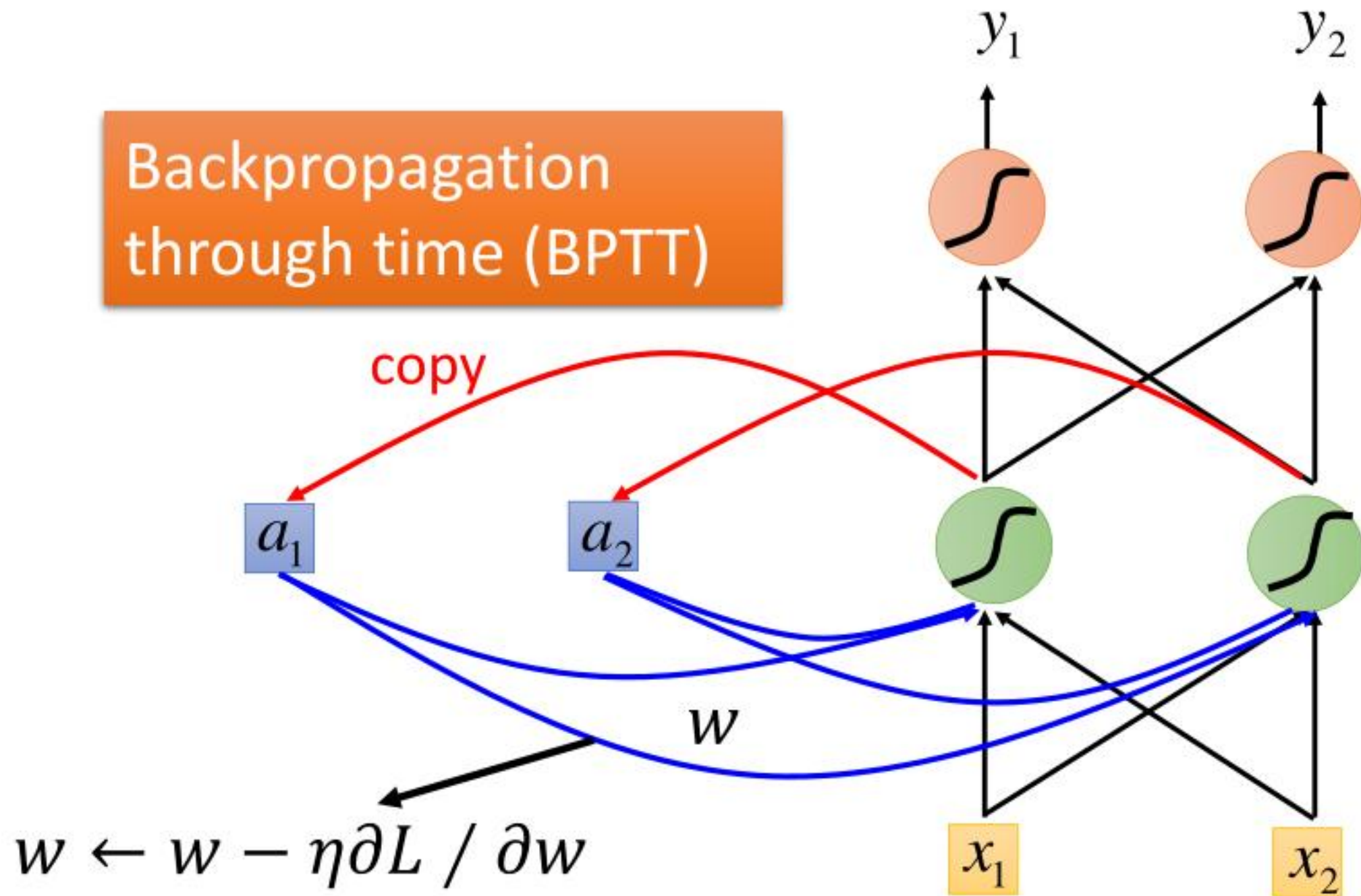


# Learning Target





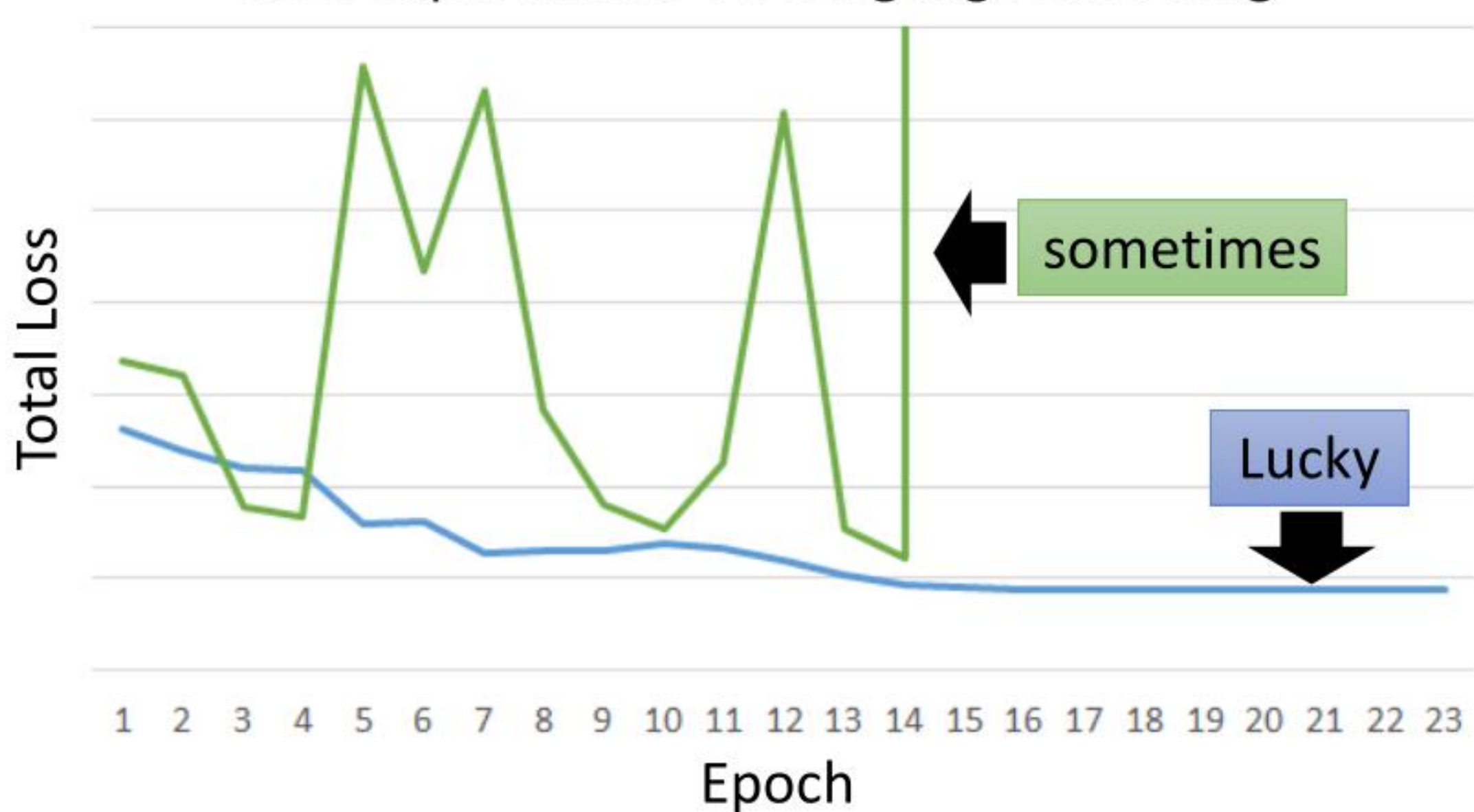
# Learning



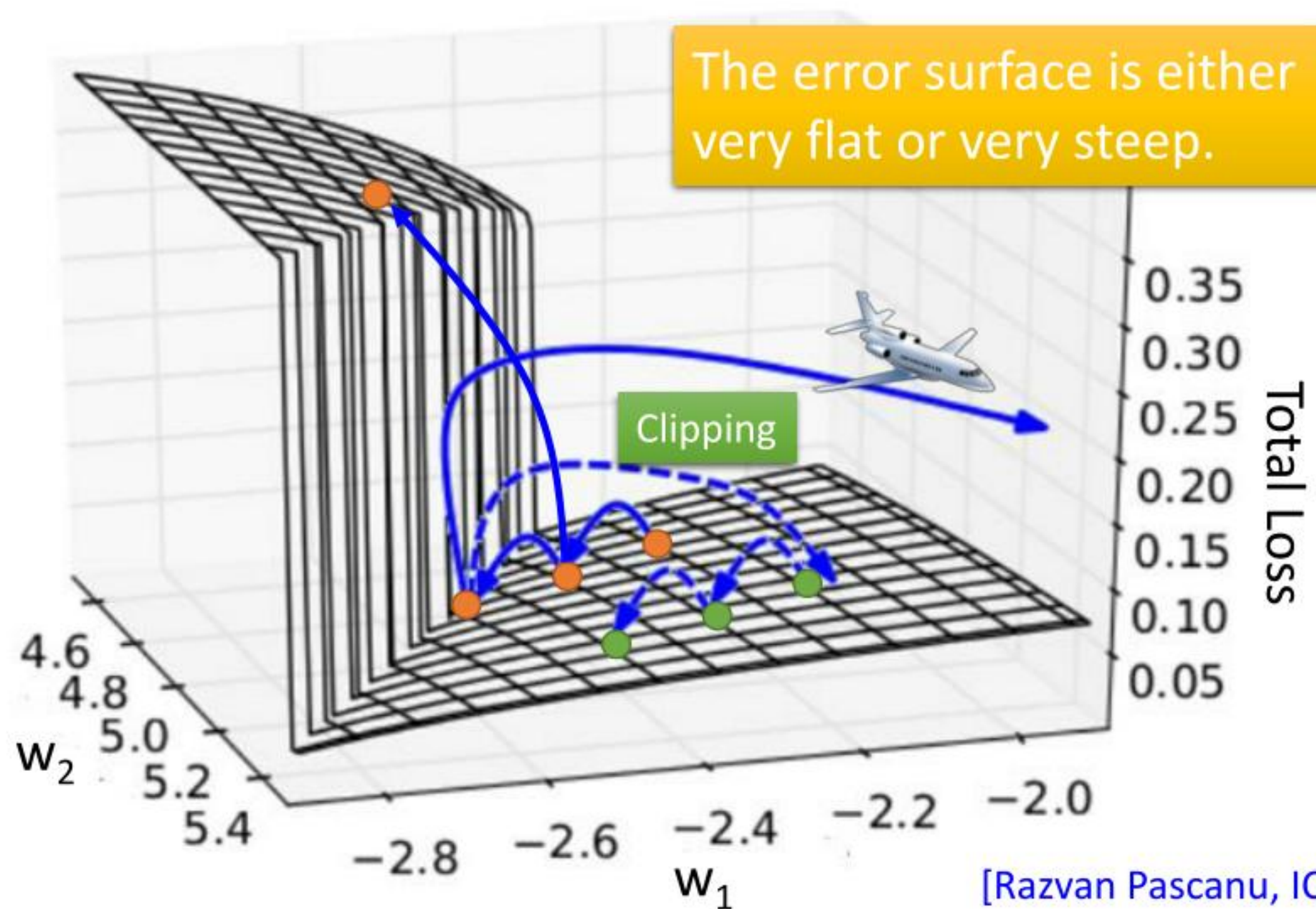
# Unfortunately .....

- RNN-based network is not always easy to learn

Real experiments on Language modeling



# The error surface is rough.





# Why?

$$\begin{array}{ll} w = 1 & \longrightarrow y^{1000} = 1 \\ w = 1.01 & \longrightarrow y^{1000} \approx 20000 \end{array}$$

Large  
 $\partial L / \partial w$

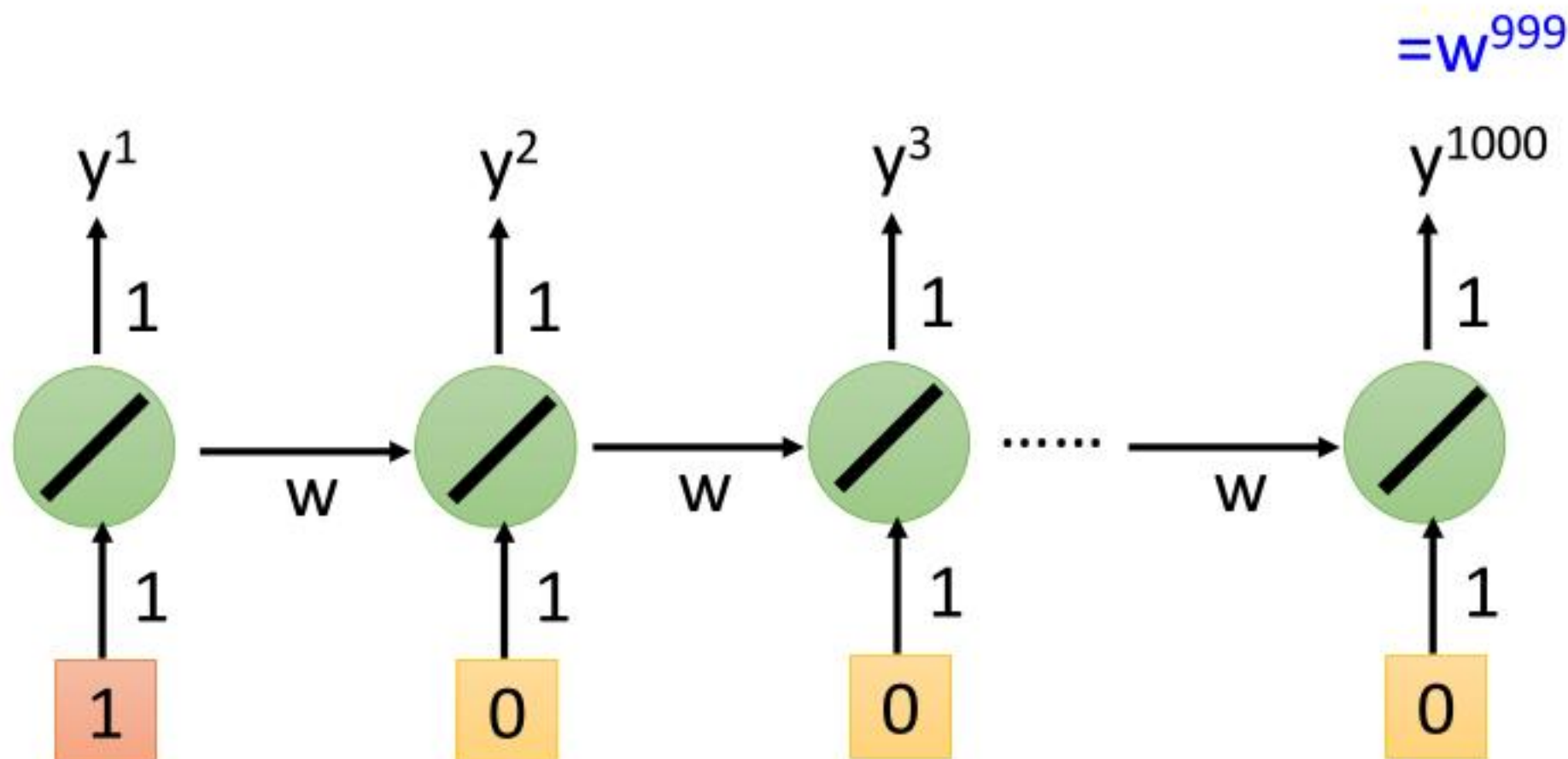
Small  
Learning rate?

$$\begin{array}{ll} w = 0.99 & \longrightarrow y^{1000} \approx 0 \\ w = 0.01 & \longrightarrow y^{1000} \approx 0 \end{array}$$

small  
 $\partial L / \partial w$

Large  
Learning rate?

## Toy Example



# Helpful Techniques

- Long Short-term Memory (LSTM)

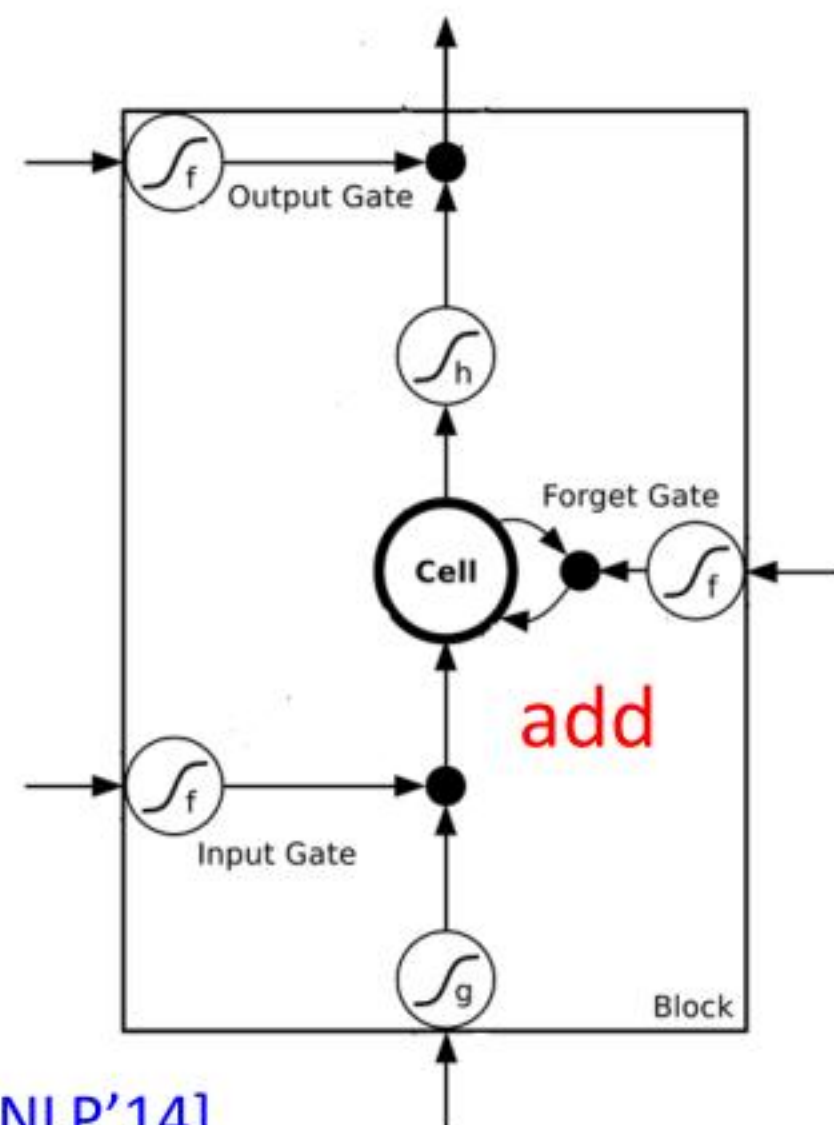
- Can deal with gradient vanishing (not gradient explode)

- Memory and input are **added**

- The influence never disappears unless forget gate is closed

➡ No Gradient vanishing  
(If forget gate is opened.)

Gated Recurrent Unit (GRU):  
simpler than LSTM



[Cho, EMNLP'14]