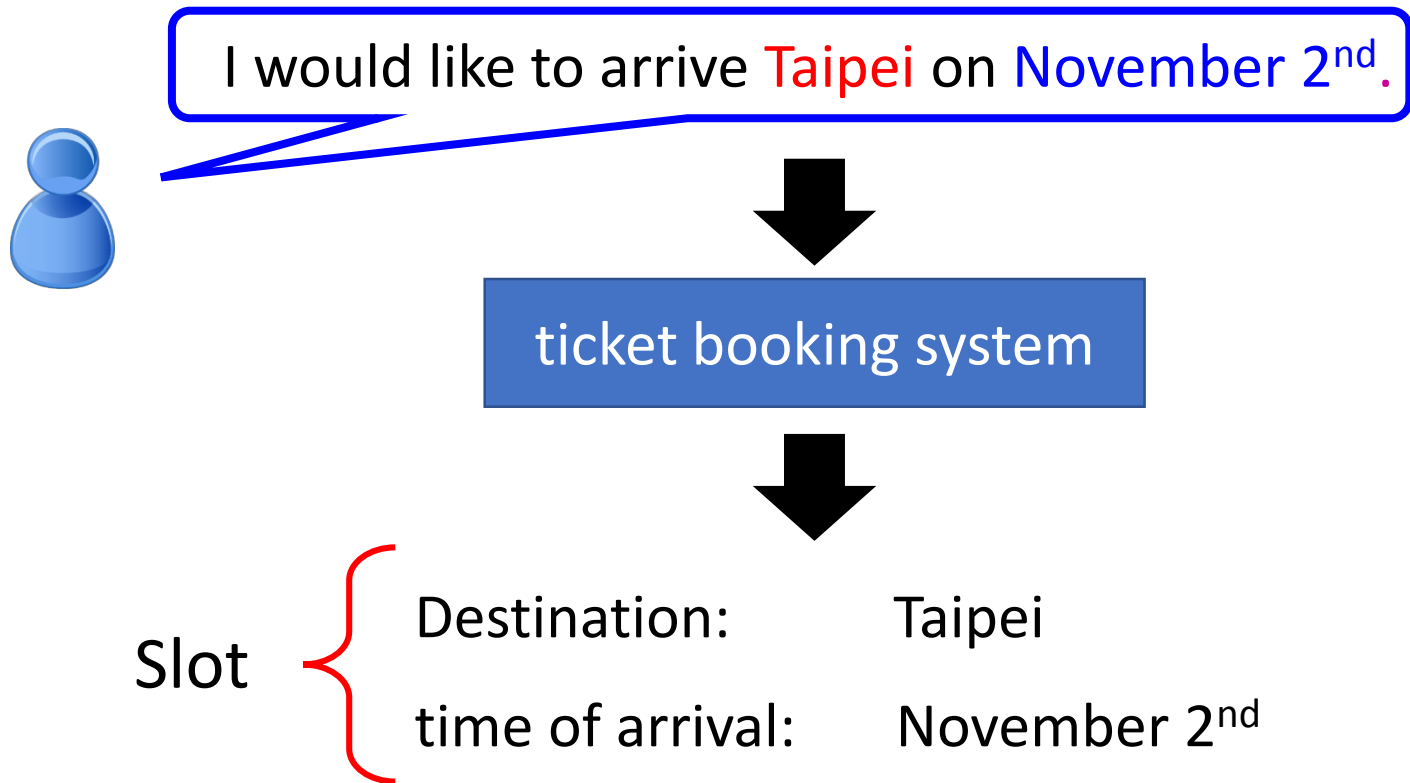


[http://speech.ee.ntu.edu.tw/~tlkagk/courses\\_ML20.html](http://speech.ee.ntu.edu.tw/~tlkagk/courses_ML20.html)

# Recurrent Neural Network (RNN)

# Example Application

- Slot Filling

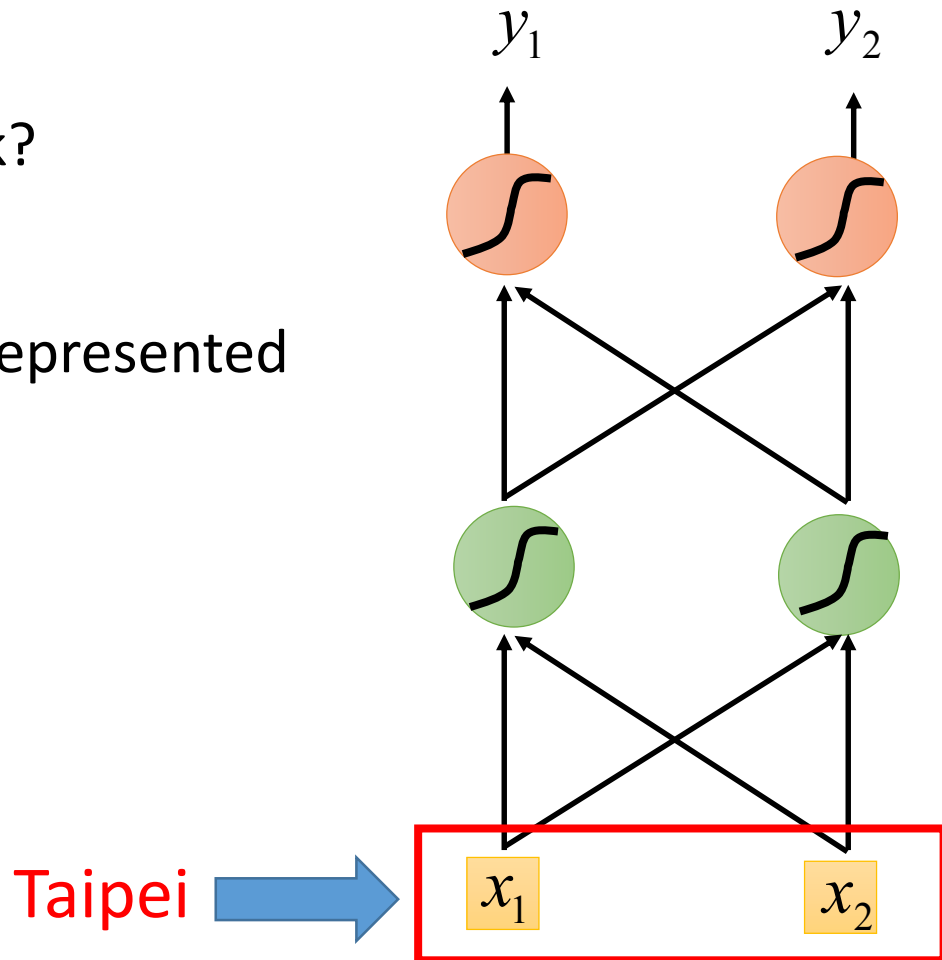


# Example Application

Solving slot filling by  
Feedforward network?

Input: a word

(Each word is represented  
as a vector)



# 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 ]

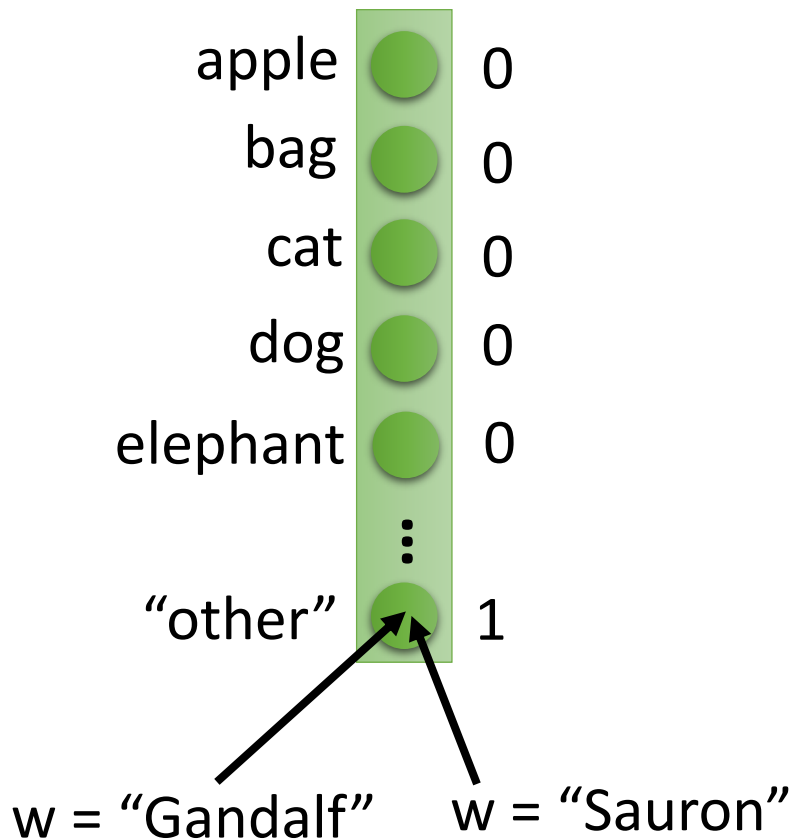
cat    = [ 0   0   1   0   0 ]

dog    = [ 0   0   0   1   0 ]

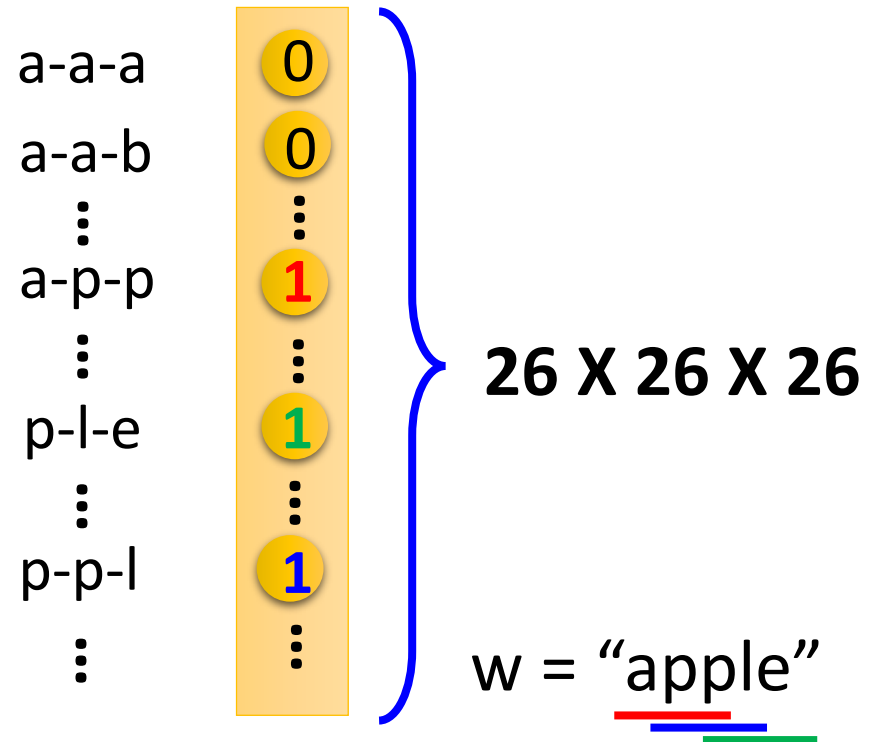
elephant = [ 0   0   0   0   1 ]

# Beyond 1-of-N encoding

## Dimension for “Other”



## Word hashing



# Example Application

Solving slot filling by  
Feedforward network?

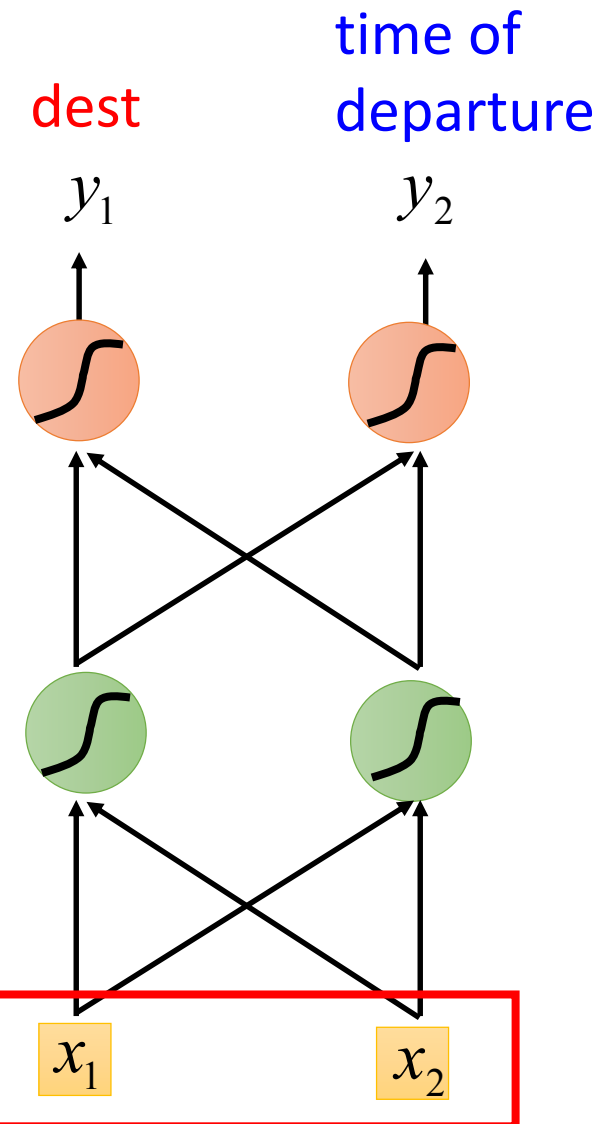
Input: a word

(Each word is represented  
as a vector)

Output:

Probability distribution that  
the input word belonging to  
the slots

Taipei →



# Example Application

arrive Taipei on November 2<sup>nd</sup>

other

dest

other

time

time

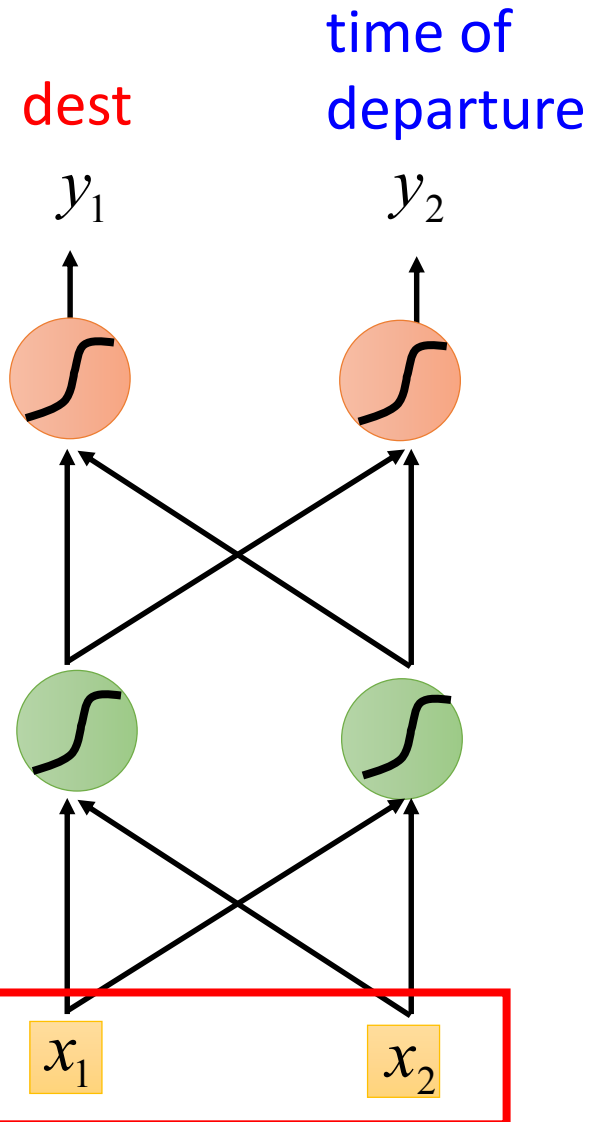
Problem?

leave Taipei on November 2<sup>nd</sup>

place of departure

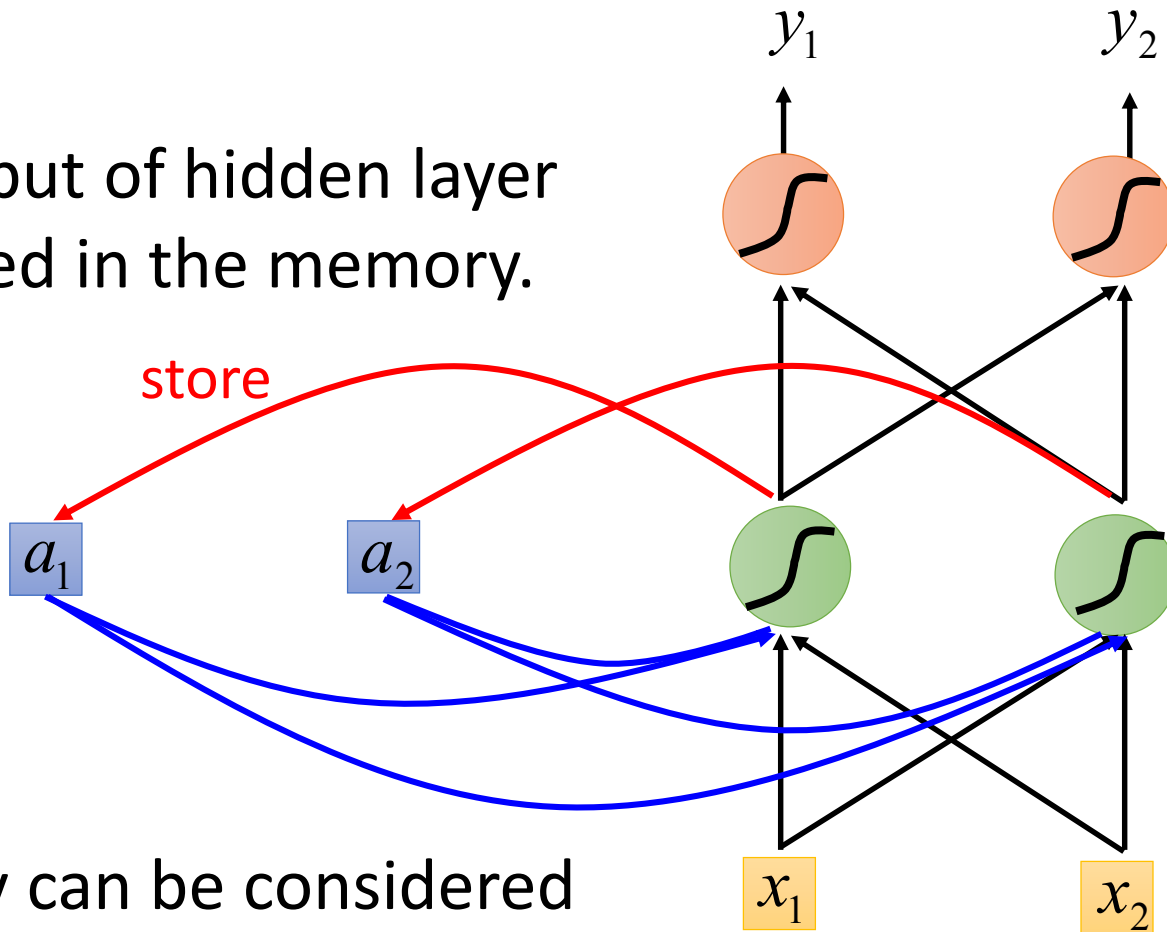
Neural network  
needs memory!

Taipei



# Recurrent Neural Network (RNN)

The output of hidden layer are stored in the memory.



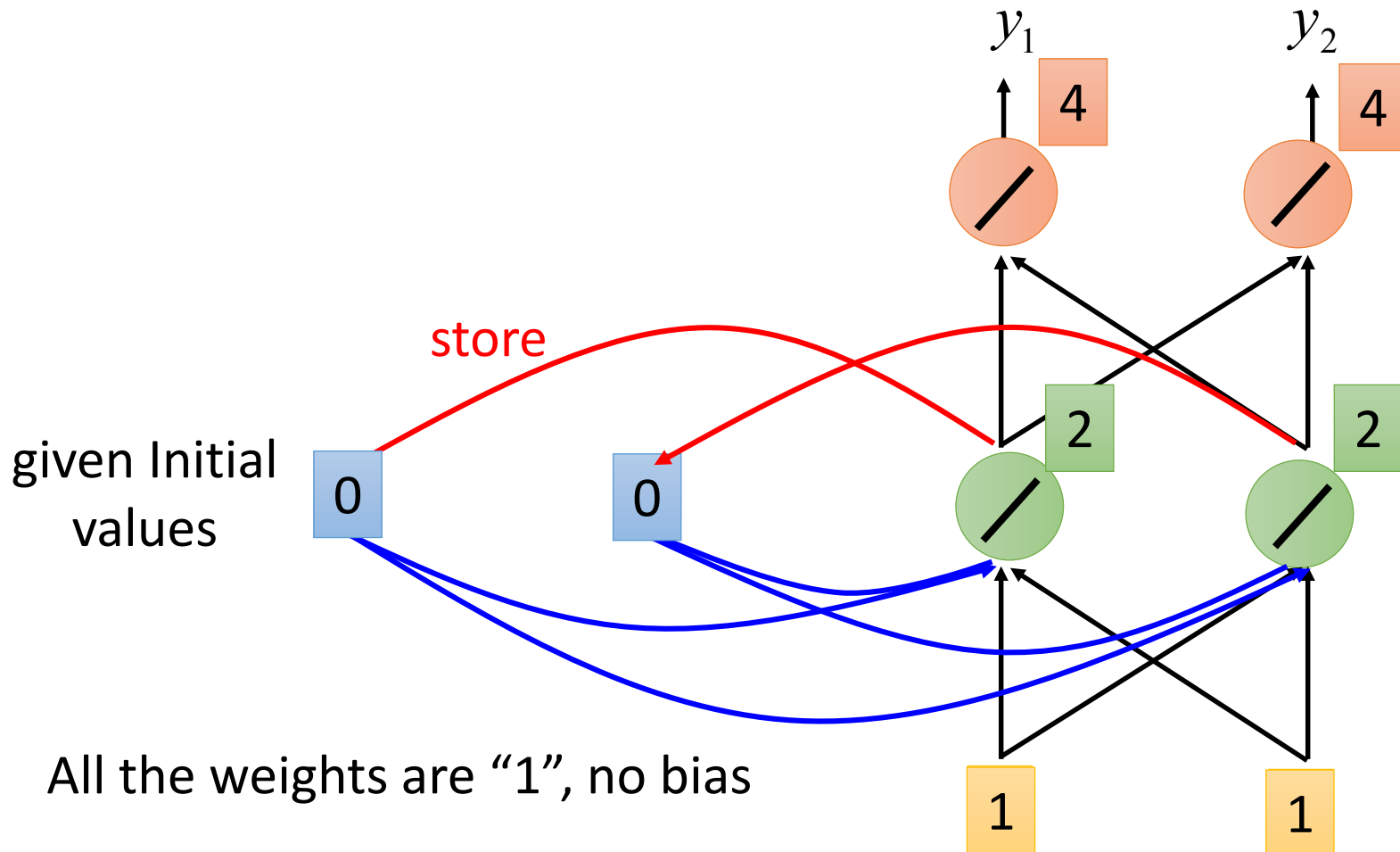
Memory can be considered as another input.



# Example

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

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



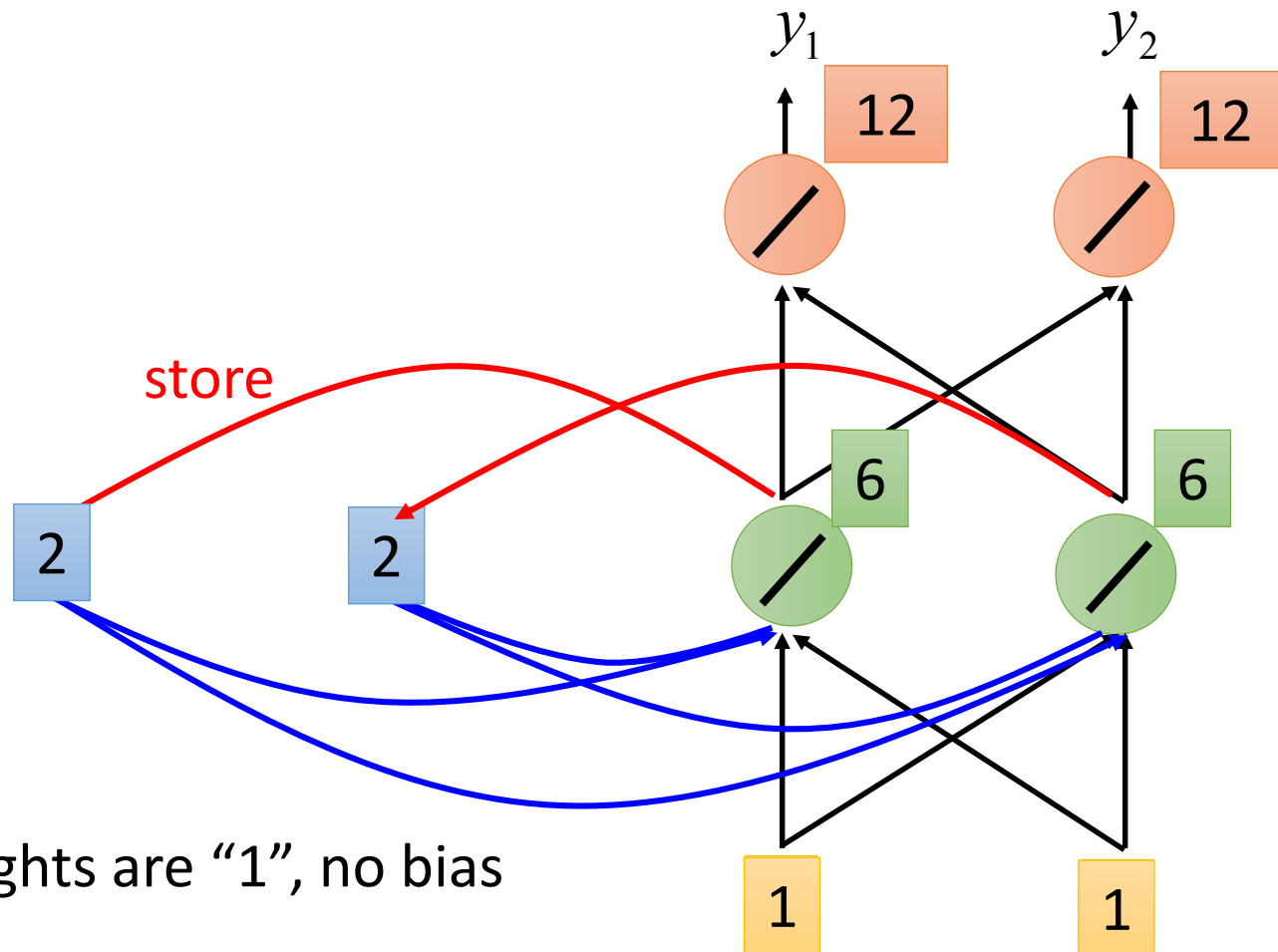
All the weights are "1", no bias

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



All the weights are "1", no bias

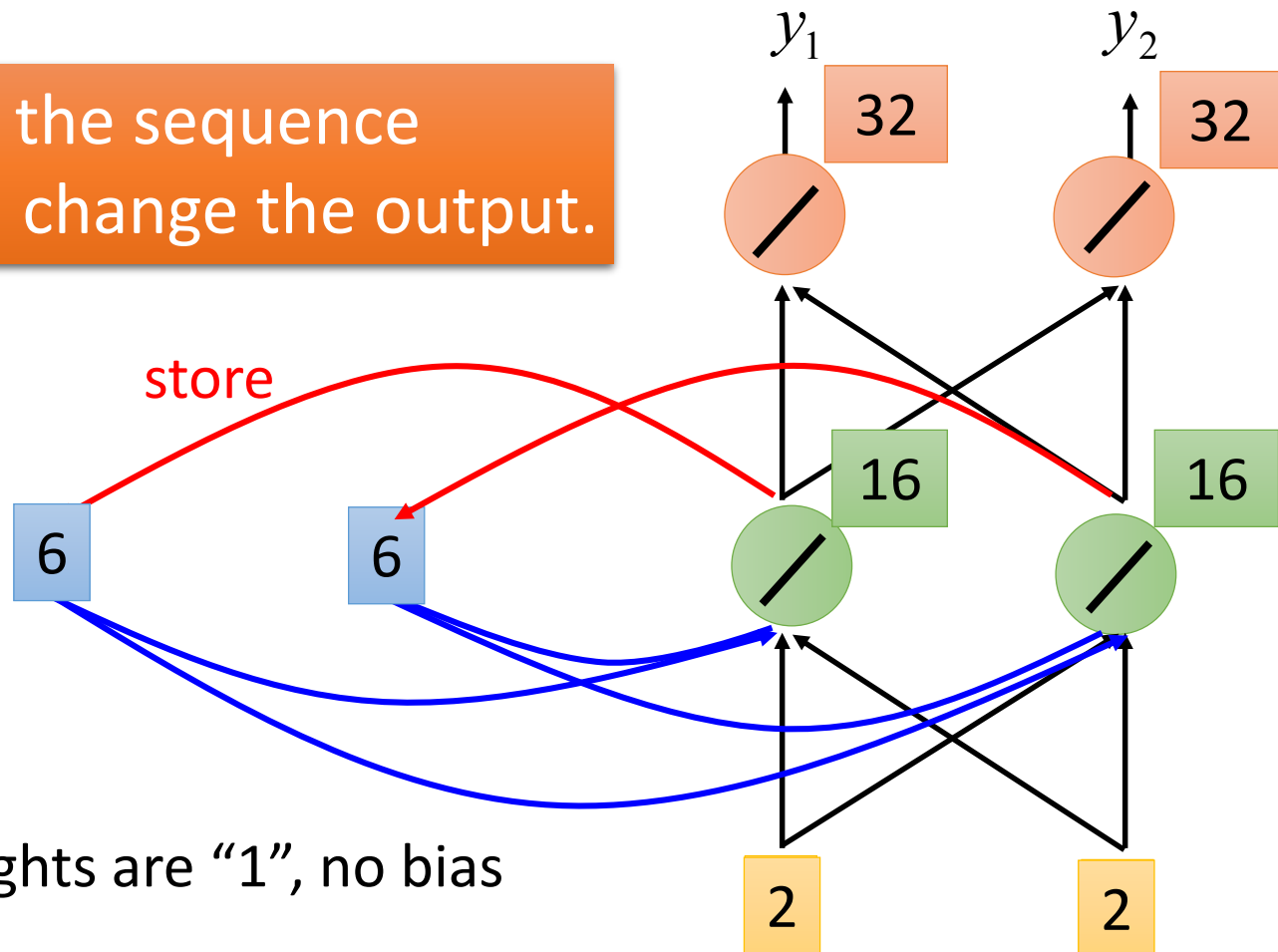
All activation functions are linear

# Example

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

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

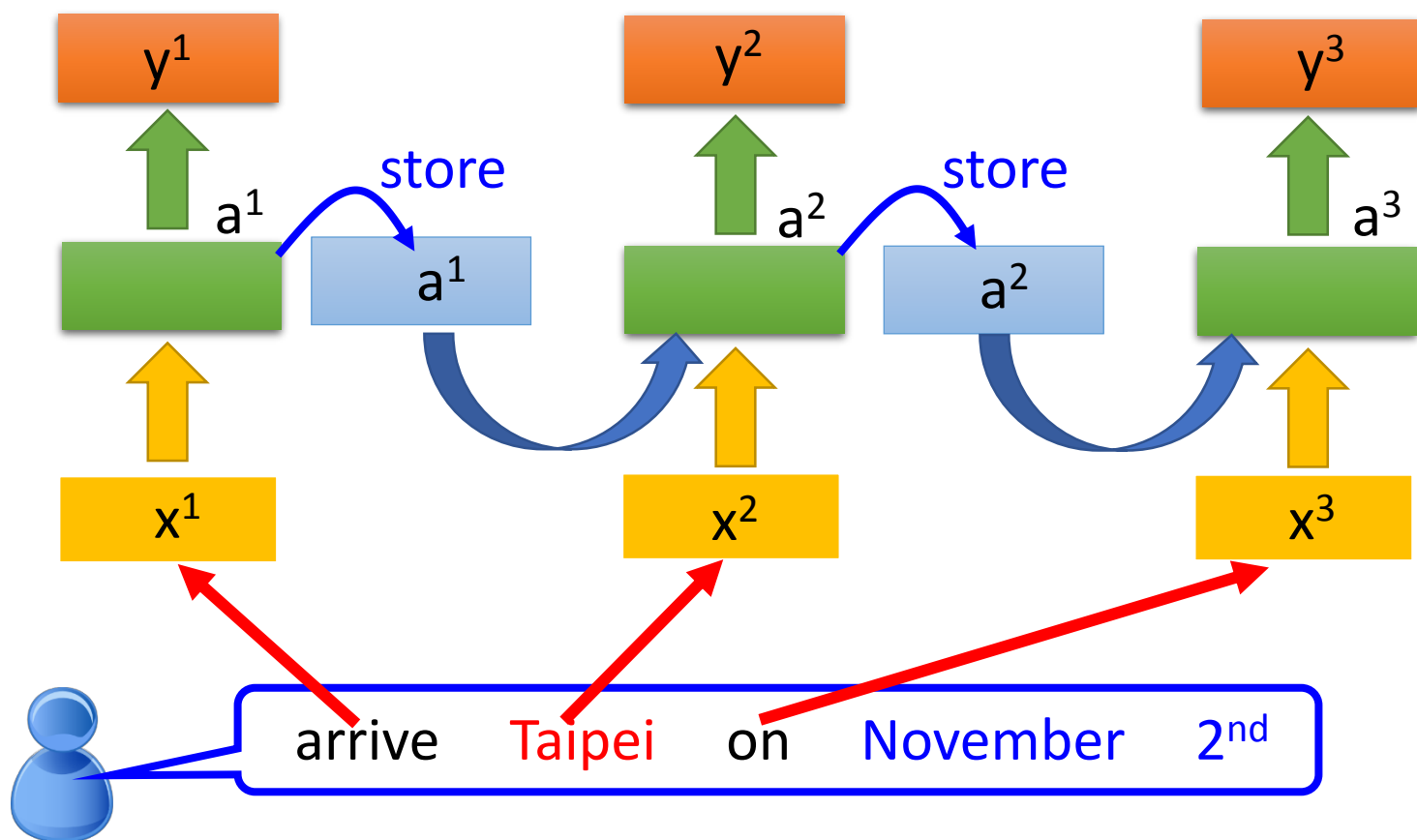
# RNN

The same network is used again and again.

Probability of  
“arrive” in each slot

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

Probability of  
“on” in each slot



# RNN

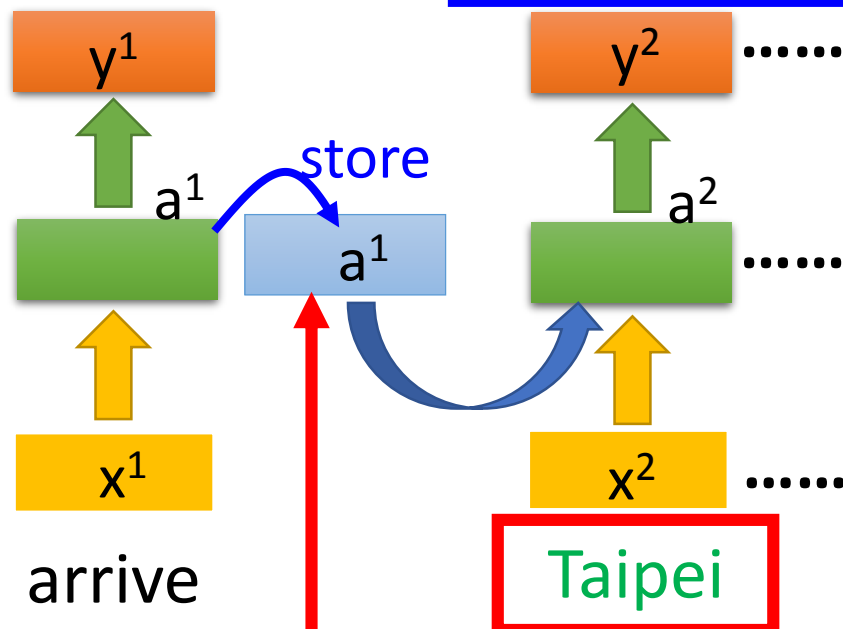
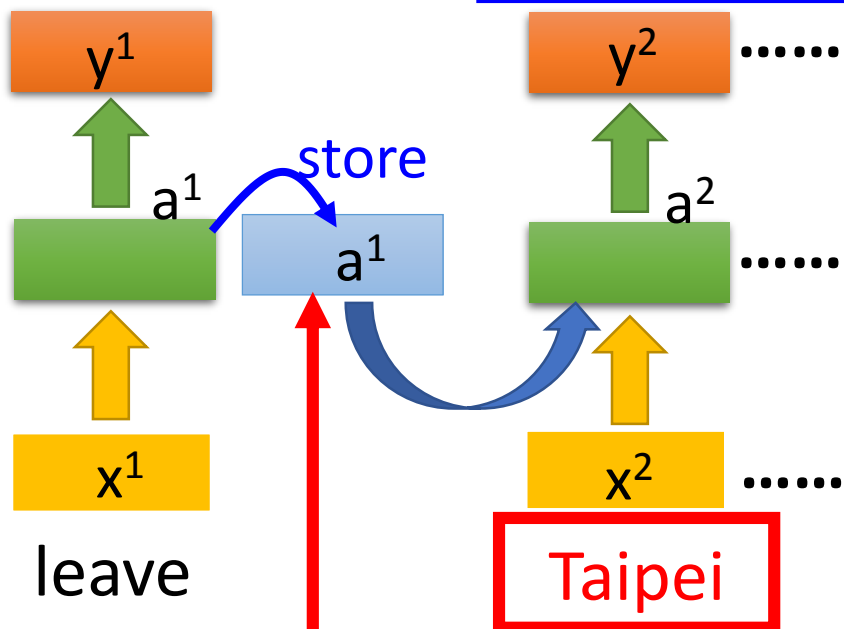
Different

Prob of “leave”  
in each slot

Prob of “**Taipei**”  
in each slot

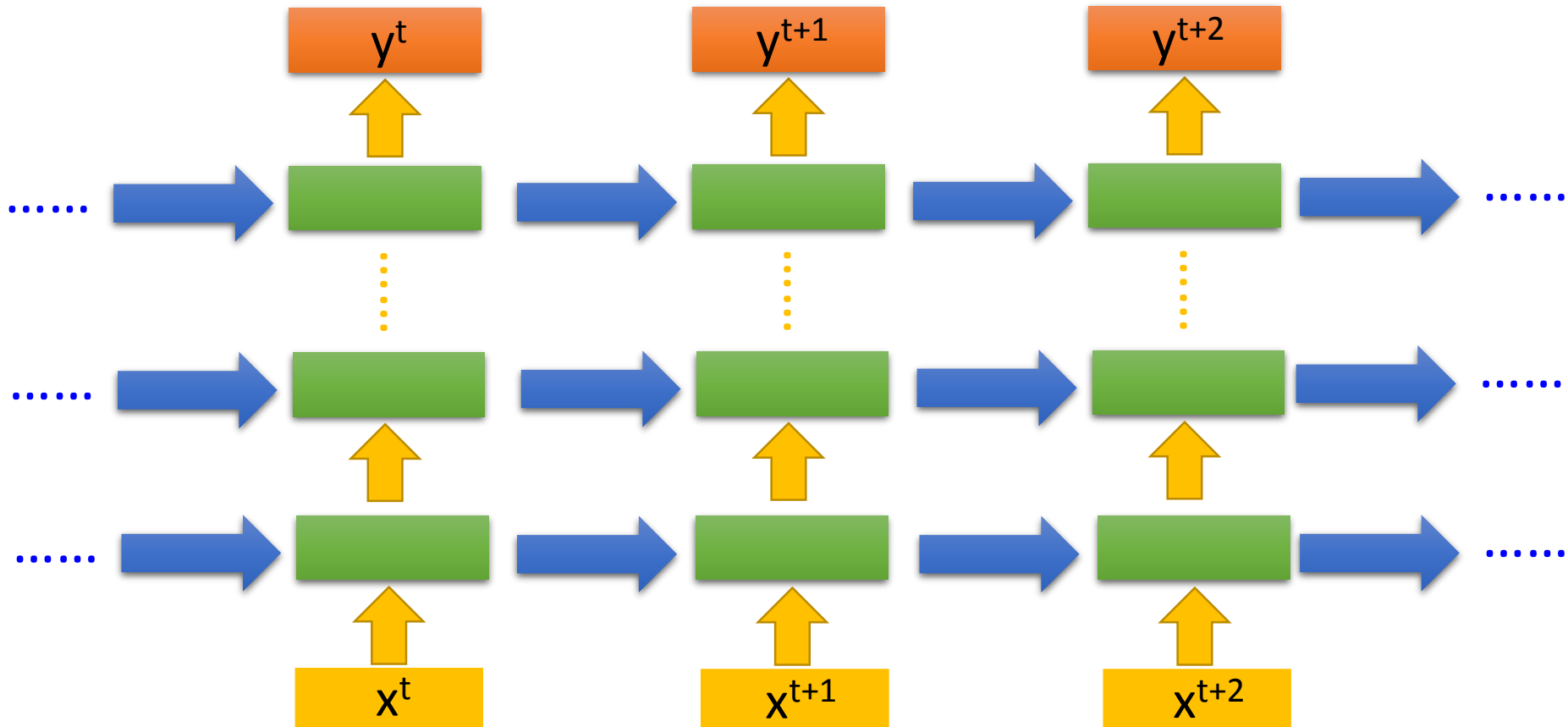
Prob of “arrive”  
in each slot

Prob of “**Taipei**”  
in each slot

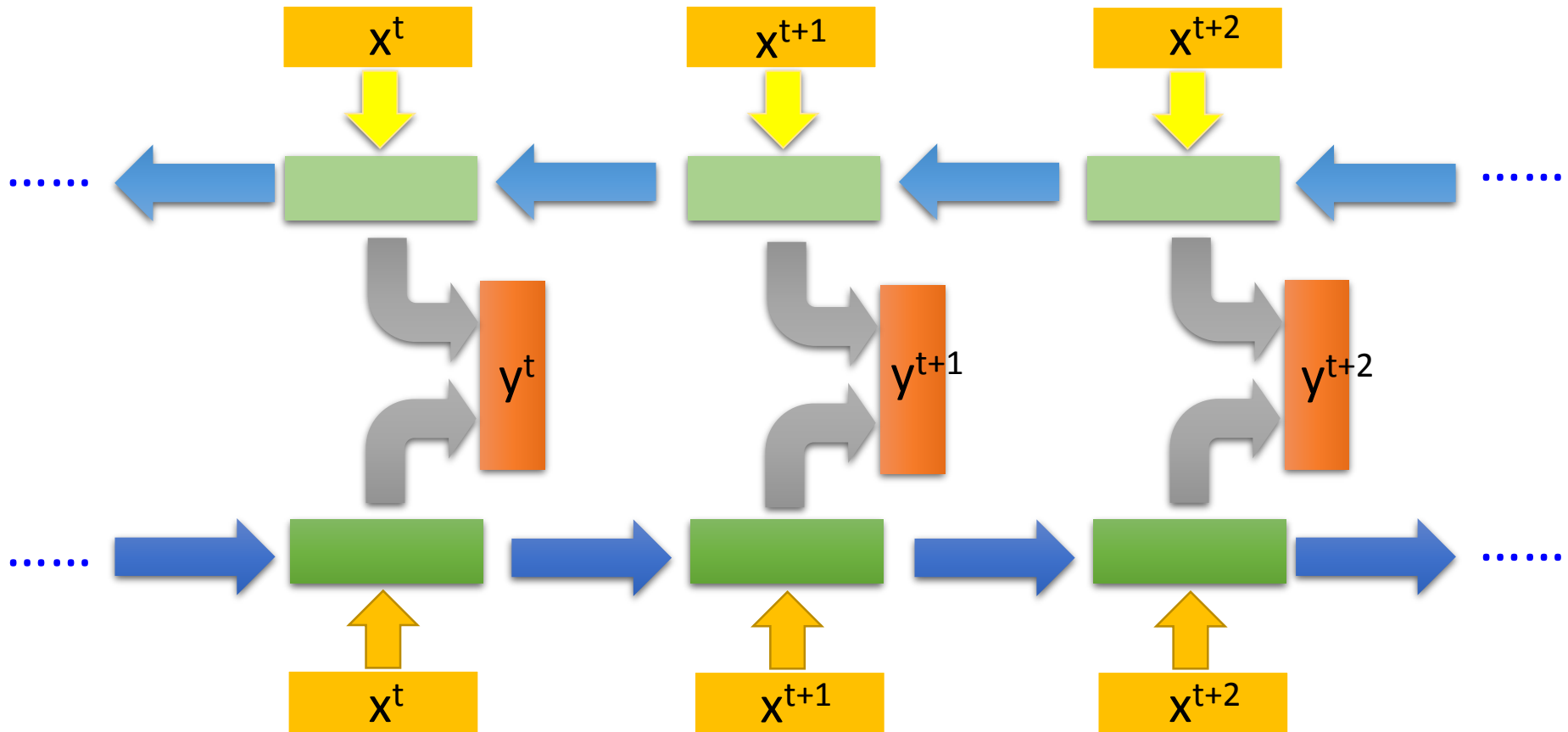


The values stored in the memory is different.

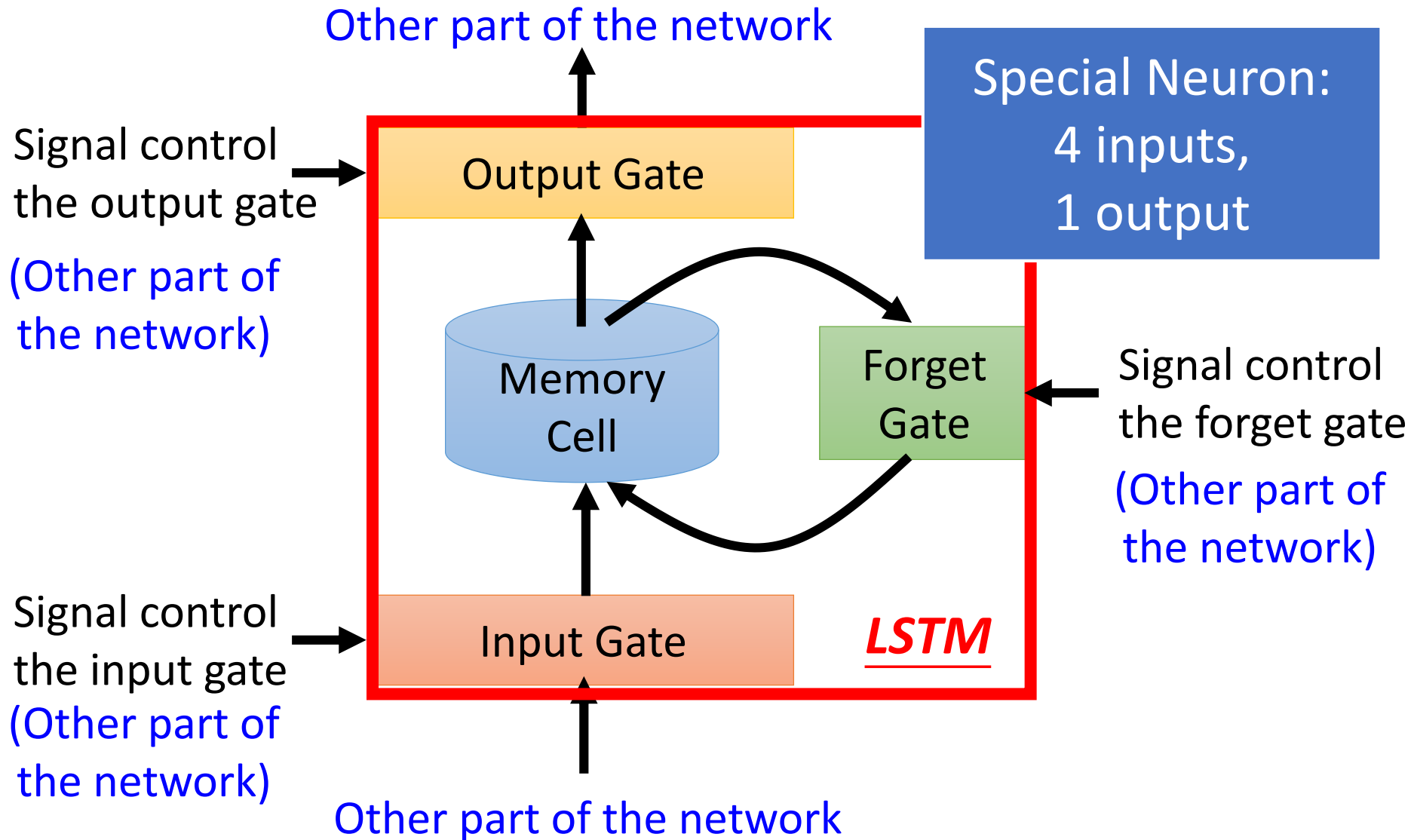
Of course it can be deep ...



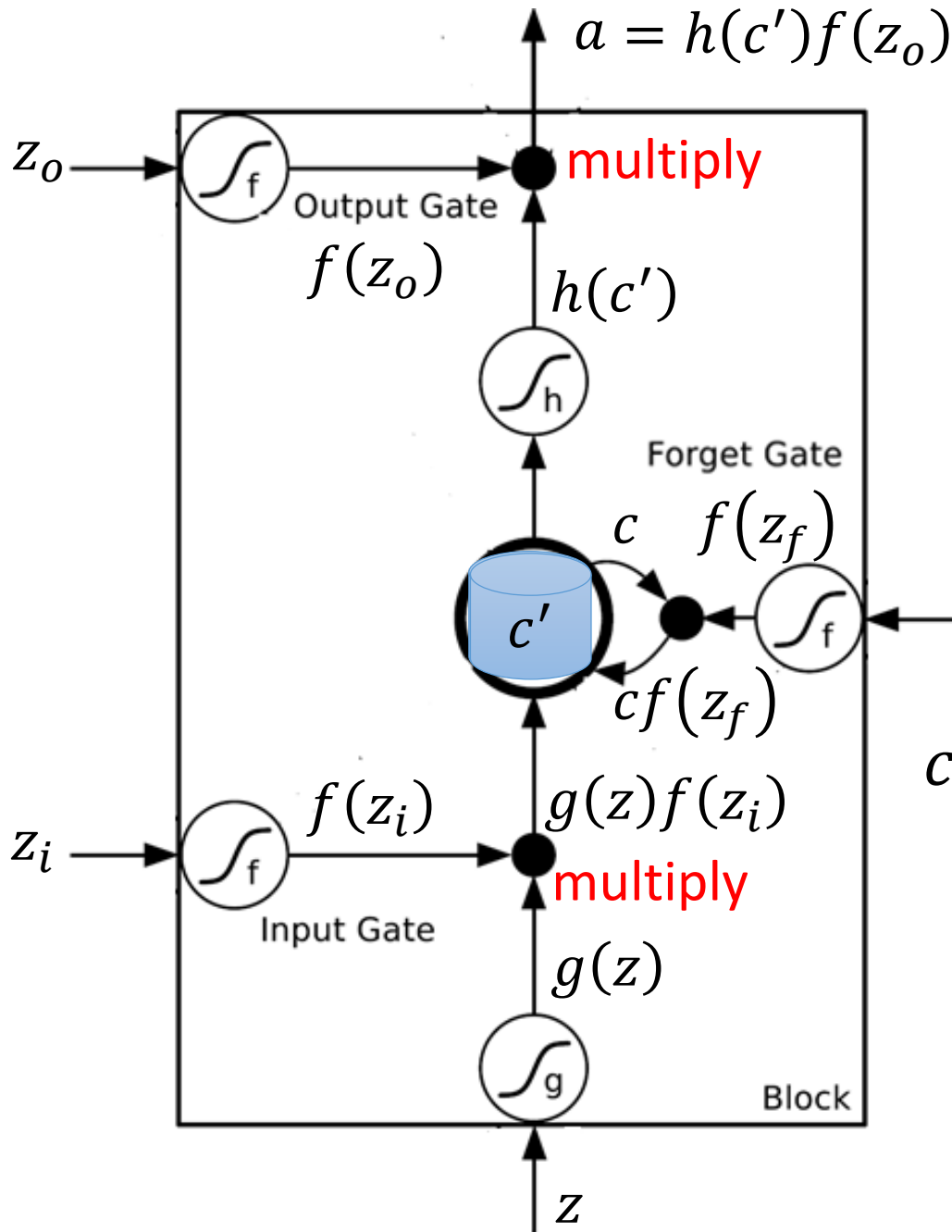
# Bidirectional RNN



# Long Short-term Memory (LSTM)







Activation function  $f$  is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

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

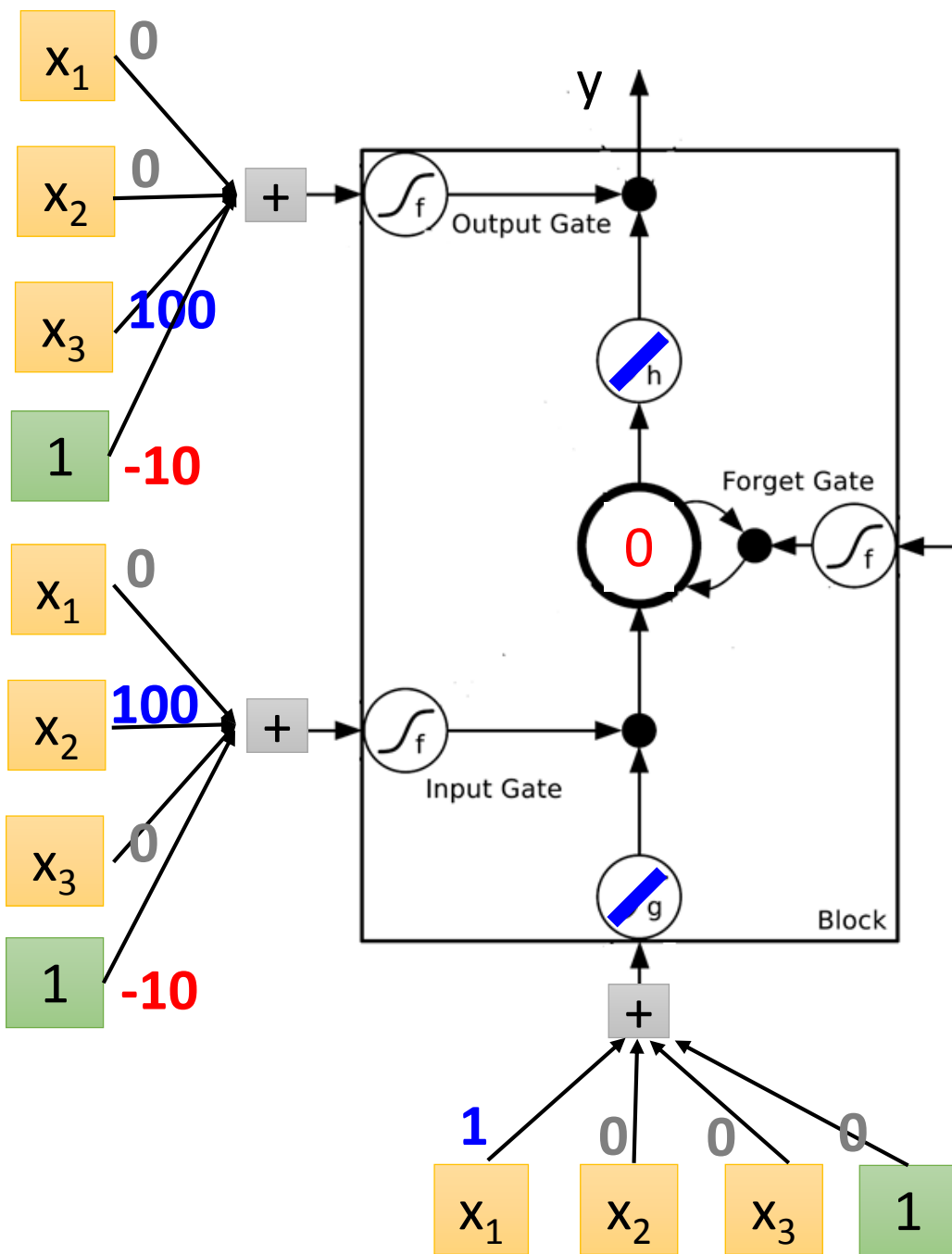
# LSTM - Example

|       |   |   |   |   |   |   |    |   |   |
|-------|---|---|---|---|---|---|----|---|---|
|       | 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 |

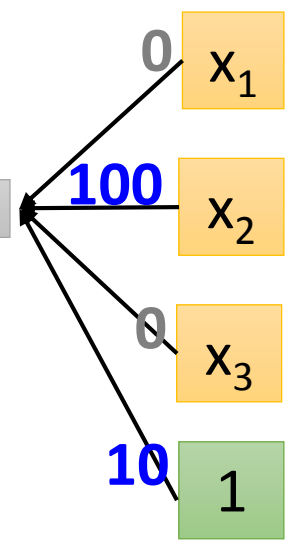
When  $x_2 = 1$ , add the numbers of  $x_1$  into the memory

When  $x_2 = -1$ , reset the memory

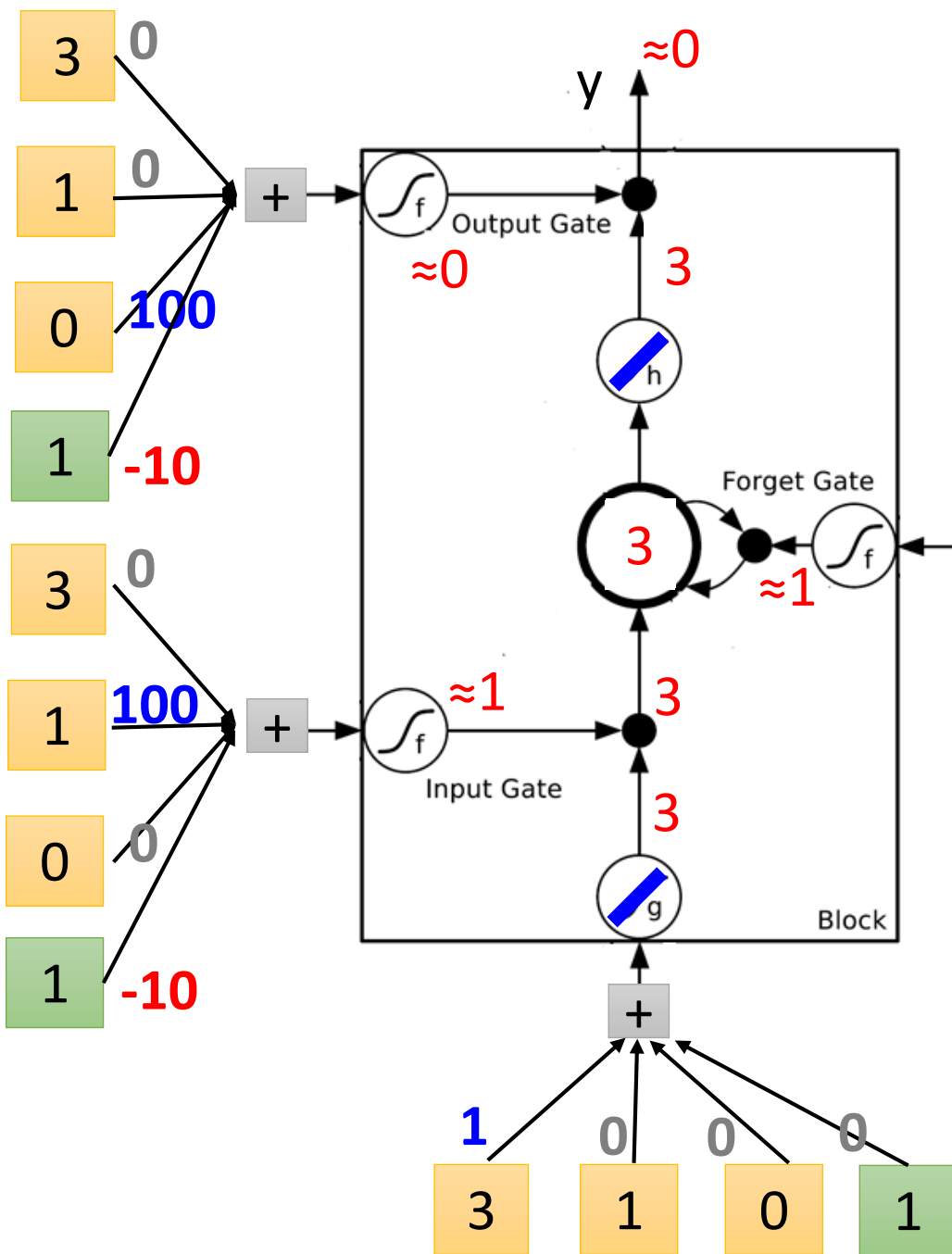
When  $x_3 = 1$ , output the number in the memory.



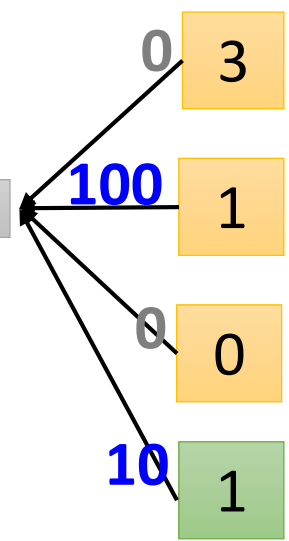
$y$  0 0 0 7 0



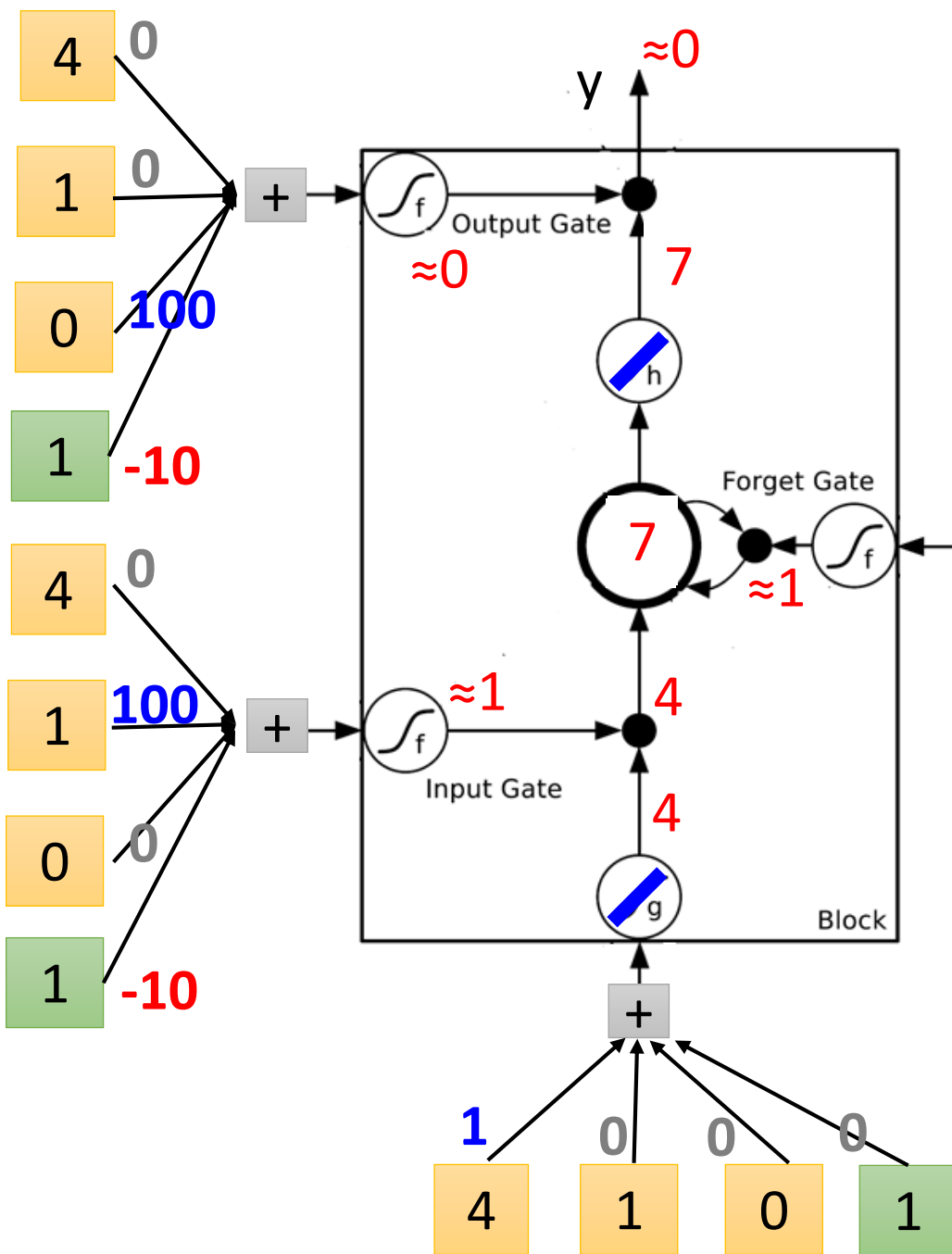
|       |   |   |   |   |    |
|-------|---|---|---|---|----|
| $x_1$ | 3 | 4 | 2 | 1 | 3  |
| $x_2$ | 1 | 1 | 0 | 0 | -1 |
| $x_3$ | 0 | 0 | 0 | 1 | 0  |



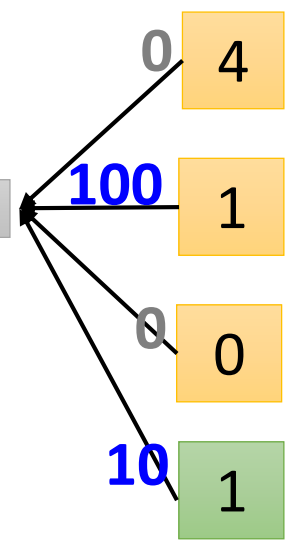
$y$  0 0 0 7 0



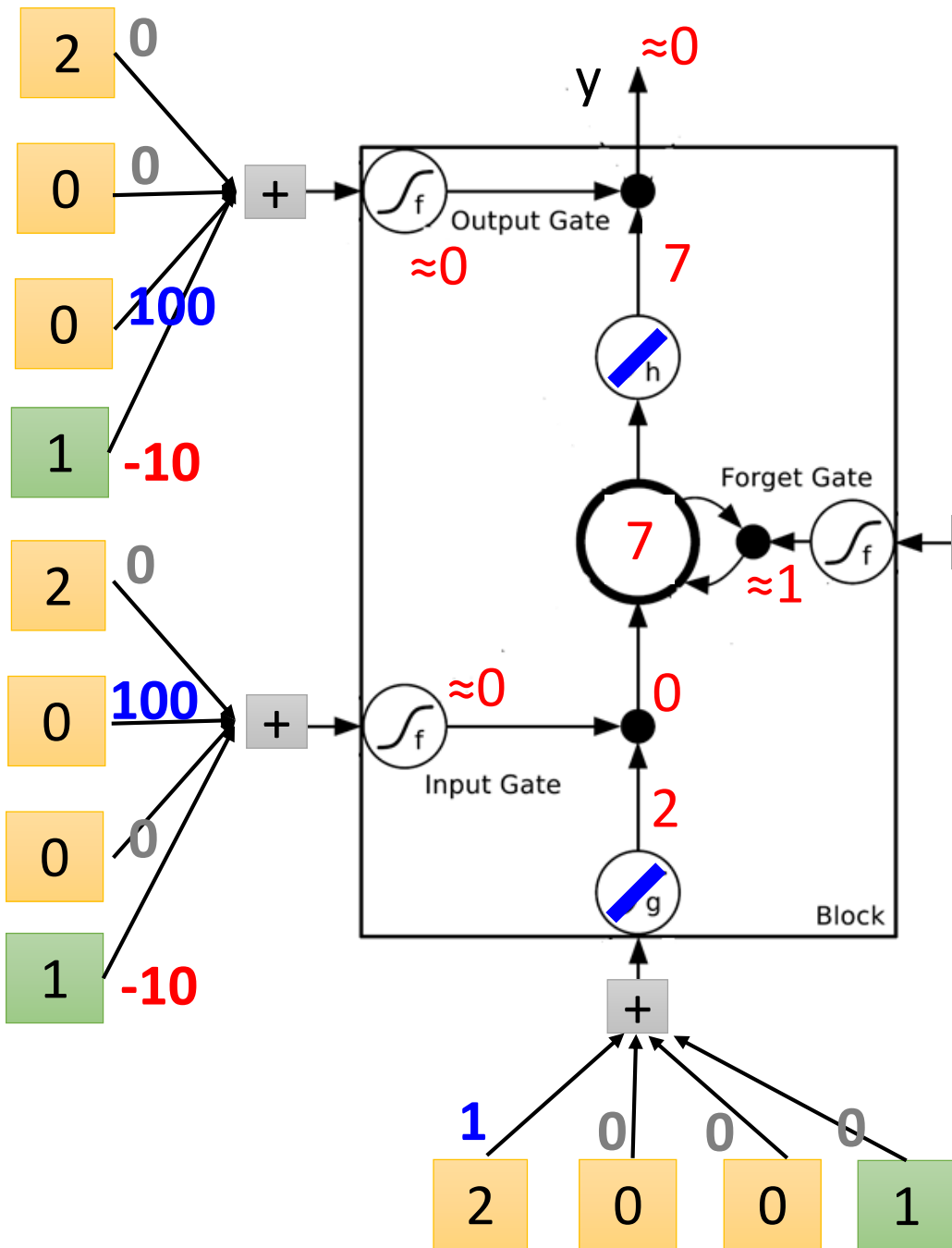
|  | $x_1$ | $x_2$ | $x_3$ |
|--|-------|-------|-------|
|  | 3     | 4     | 2     |
|  | 1     | 1     | 0     |
|  | 0     | 0     | 1     |
|  | 0     | 0     | 0     |
|  | -1    |       |       |



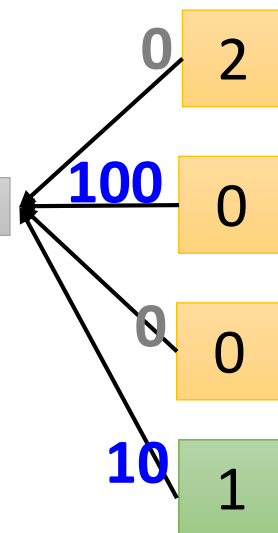
y 0 0 0 7 0



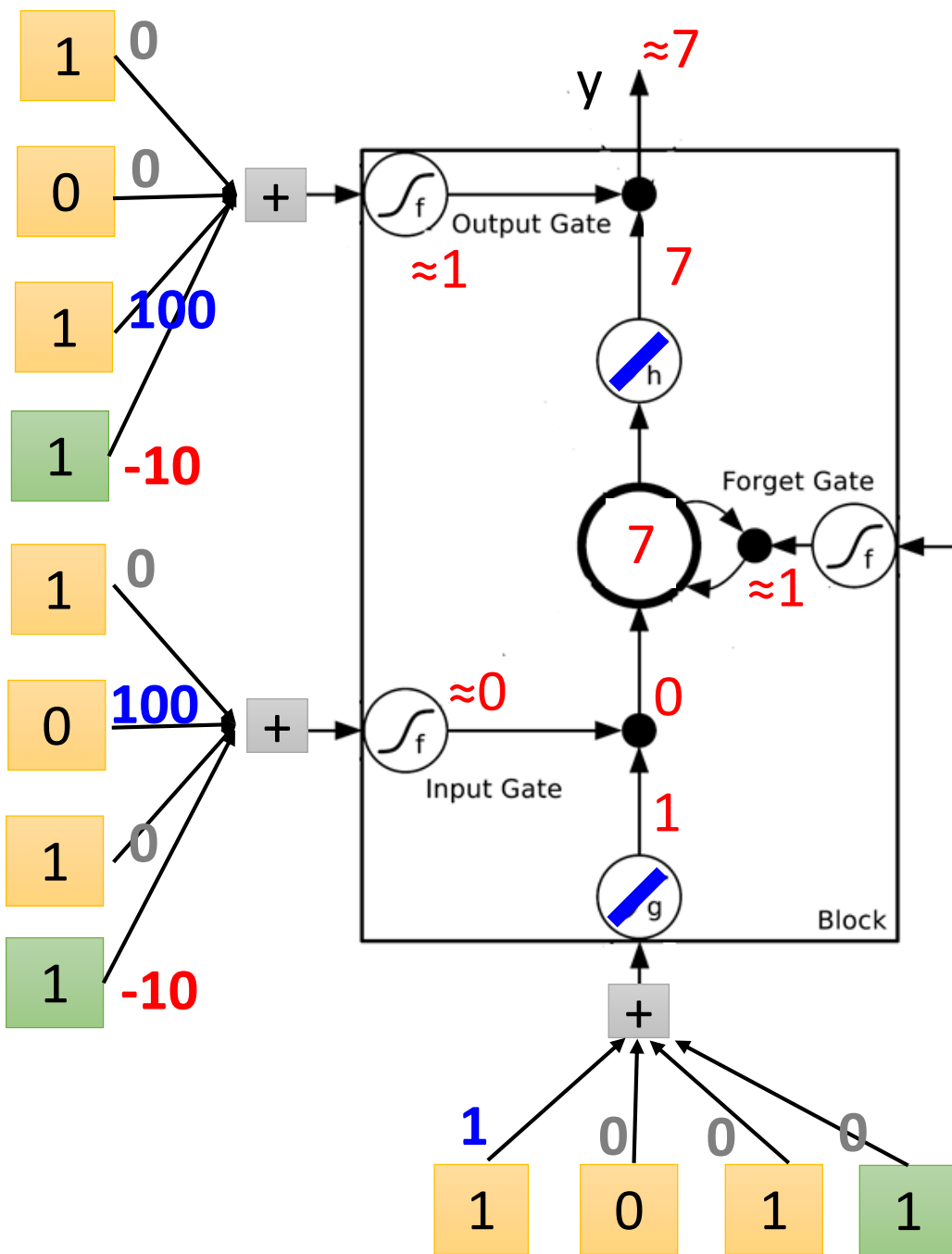
|       | $x_1$ | $x_2$ | $x_3$ |
|-------|-------|-------|-------|
| $x_1$ | 3     | 4     | 2     |
| $x_2$ | 1     | 1     | 0     |
| $x_3$ | 0     | 0     | 1     |



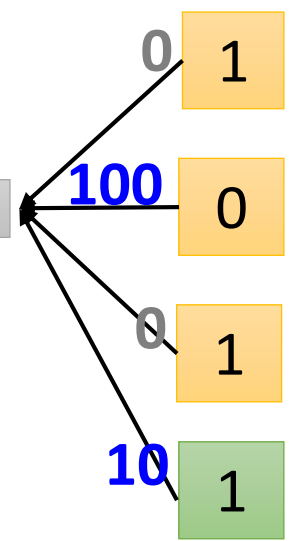
$y$  0 0 0 7 0



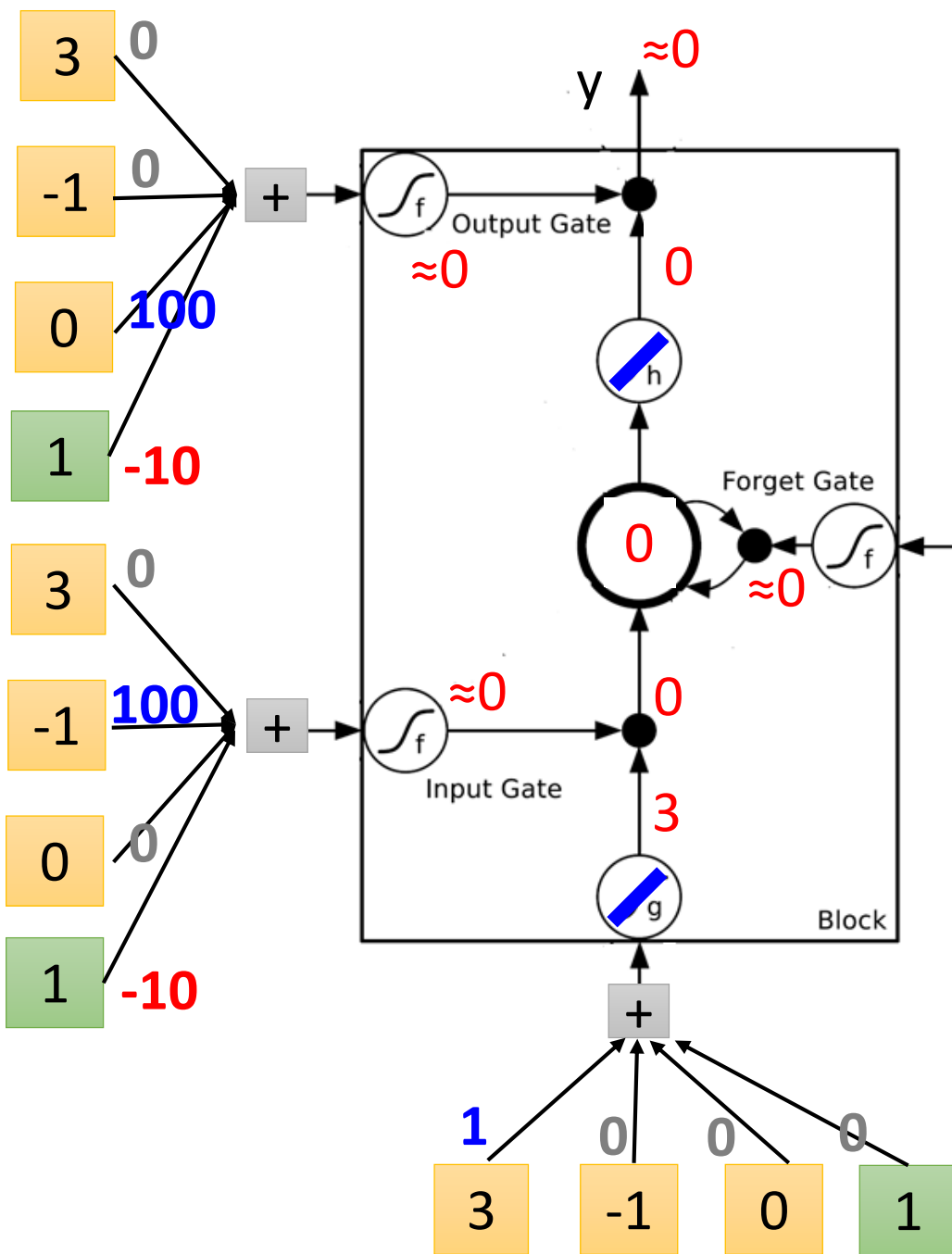
$x_1$  3 4 2 1 3  
 $x_2$  1 1 0 0 -1  
 $x_3$  0 0 0 1 0



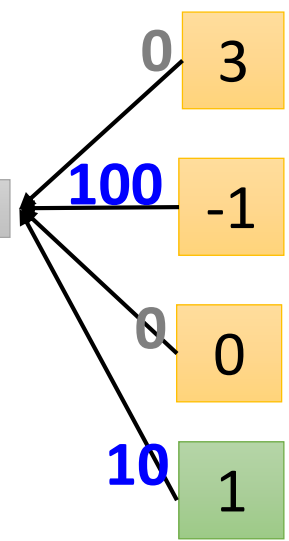
$y$  0 0 0 7 0



|       | $x_1$ |  | $x_2$ |  |   |   |    |
|-------|-------|--|-------|--|---|---|----|
|       | 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

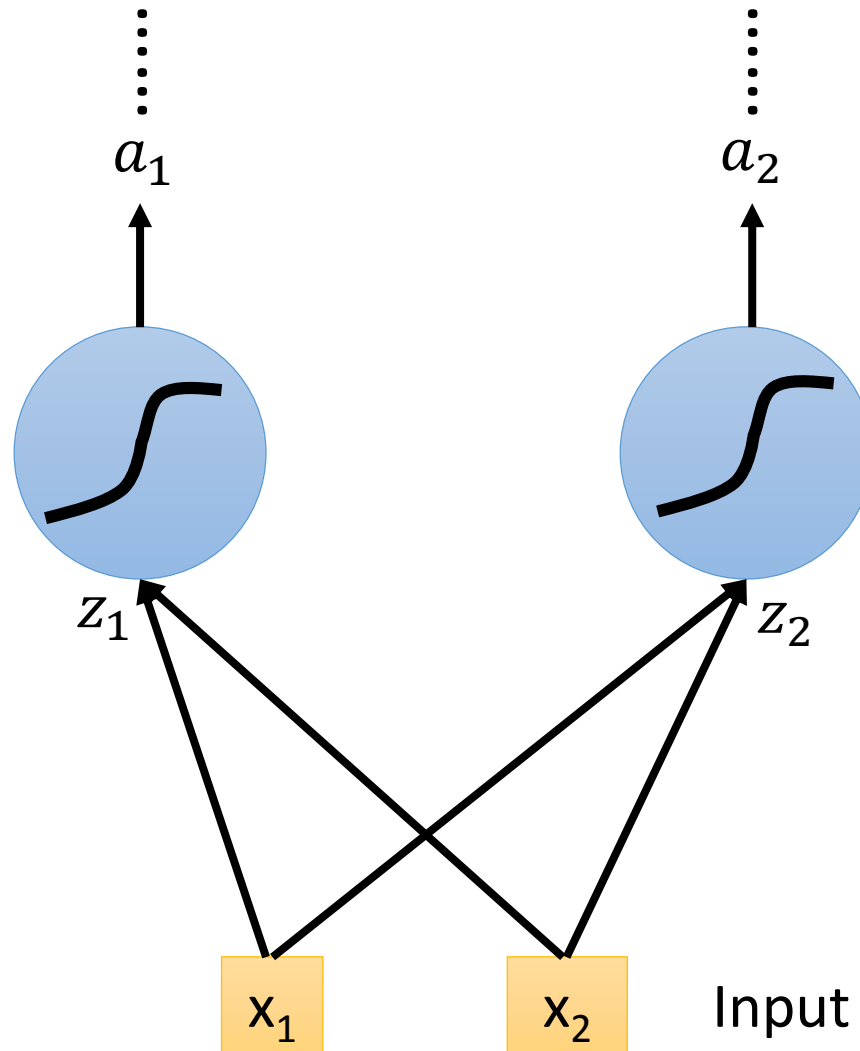


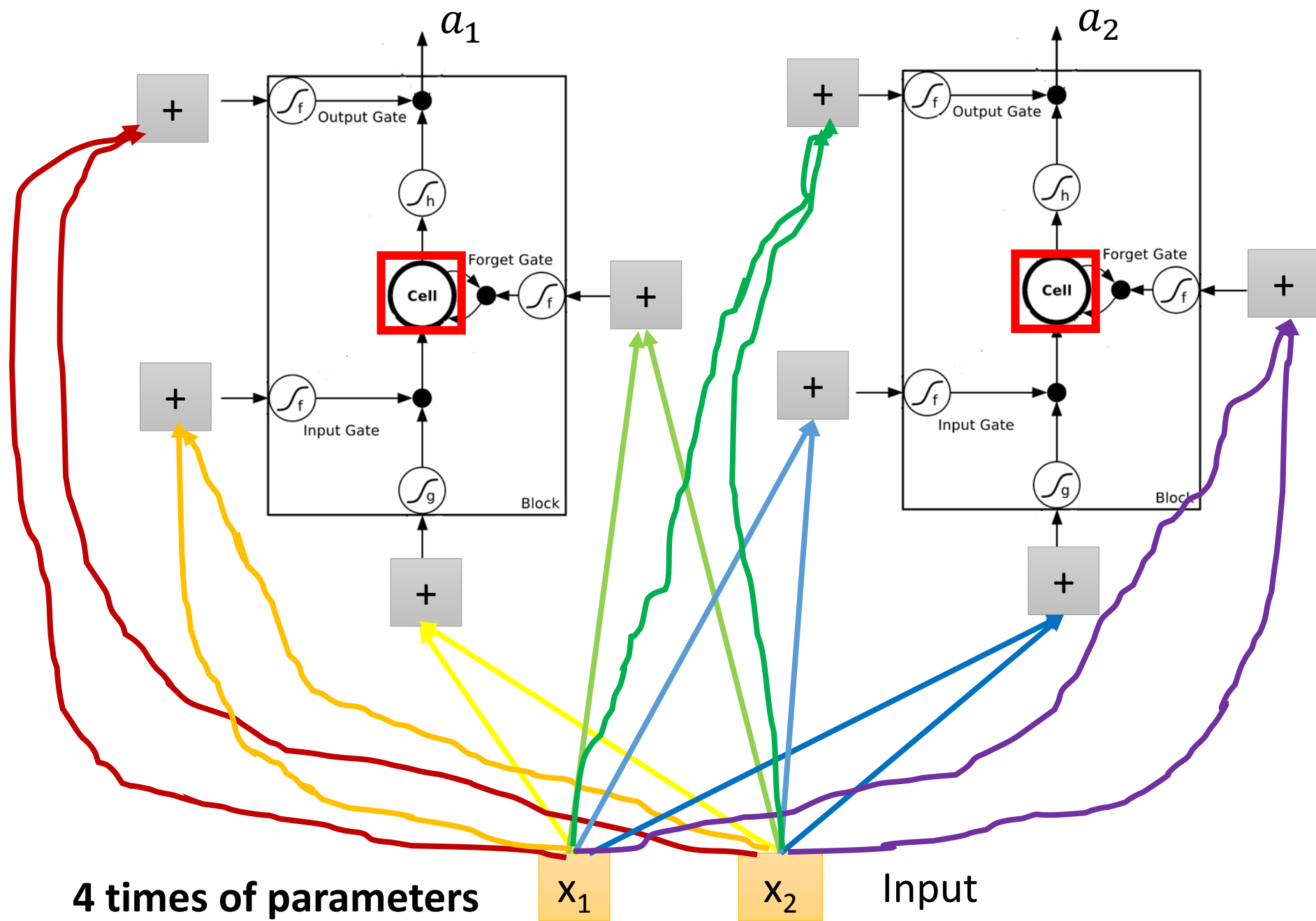
|       |   |   |   |   |    |
|-------|---|---|---|---|----|
| $x_1$ | 3 | 4 | 2 | 1 | 3  |
| $x_2$ | 1 | 1 | 0 | 0 | -1 |
| $x_3$ | 0 | 0 | 0 | 1 | 0  |



## Original Network:

- Simply replace the neurons with LSTM

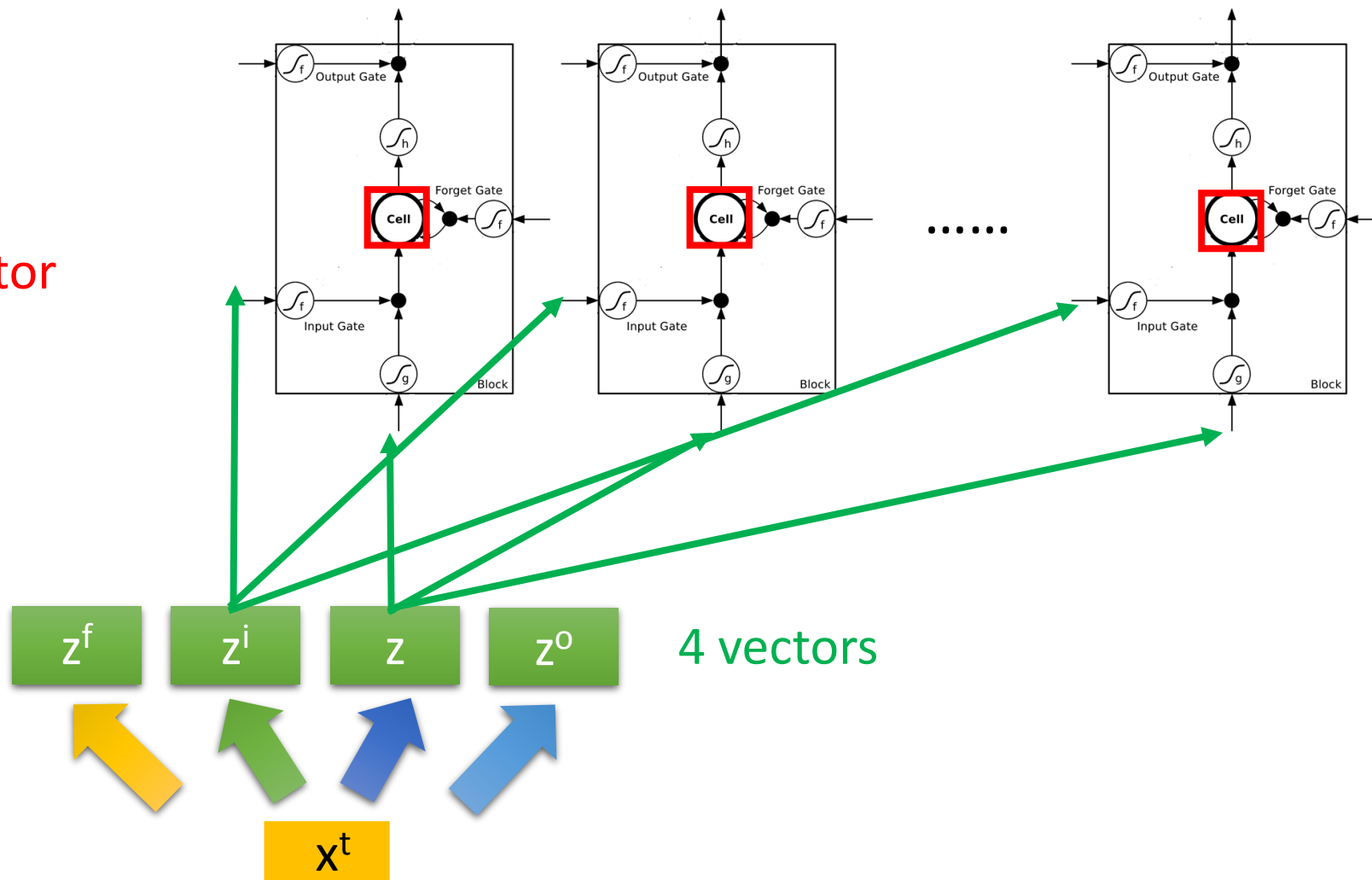




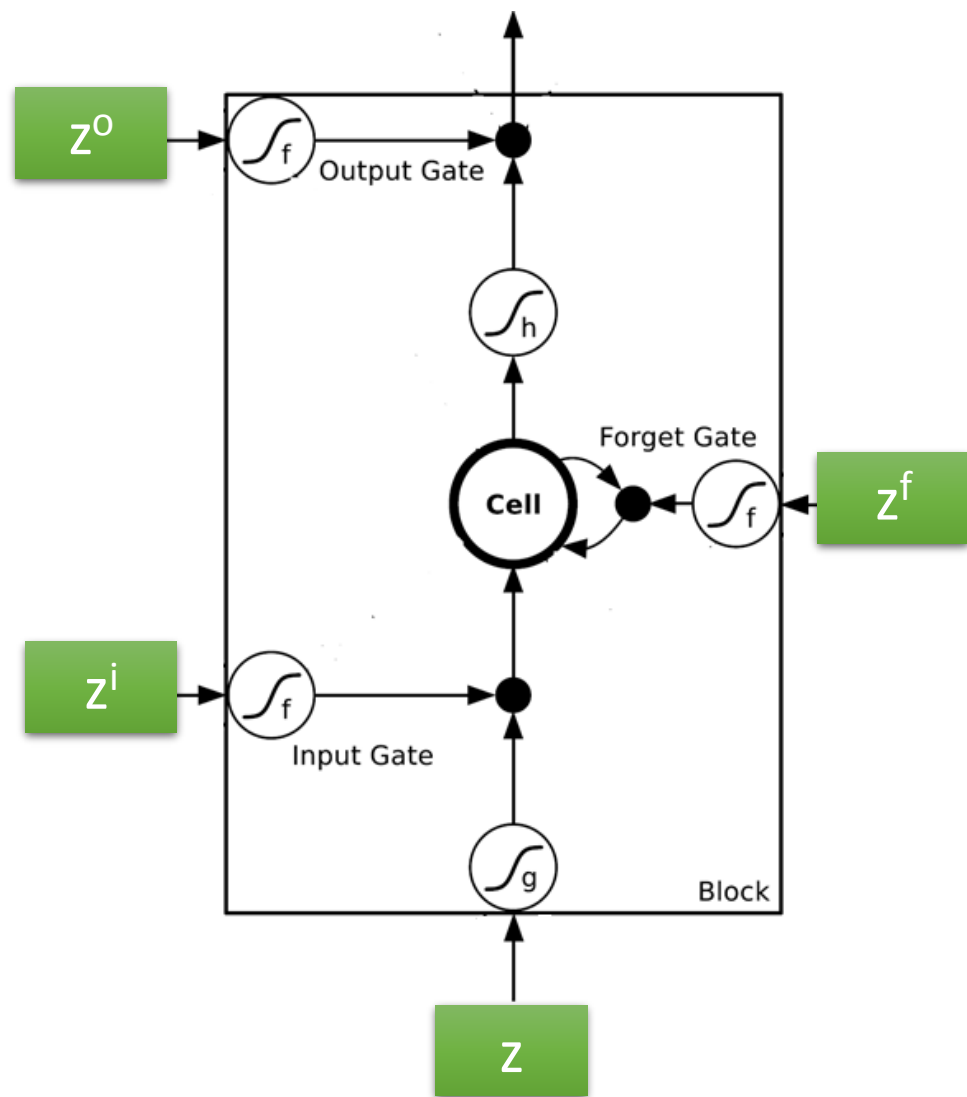
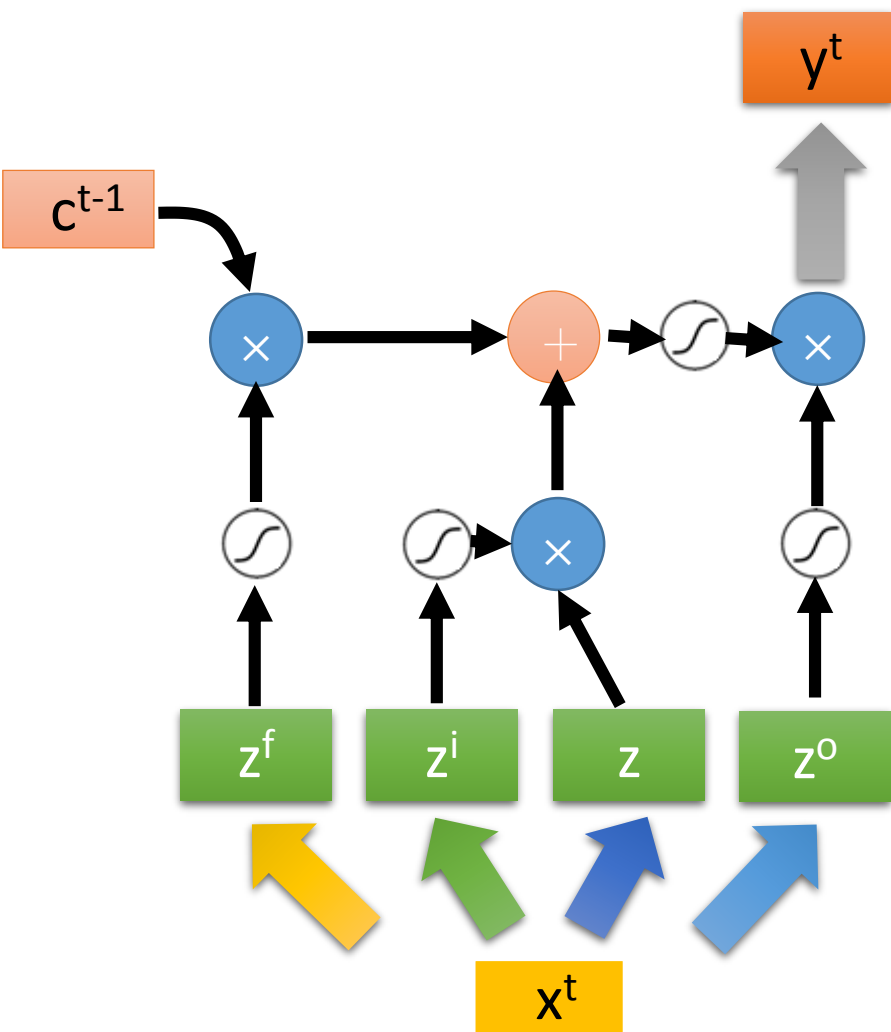
# LSTM

 $C^{t-1}$ 

vector

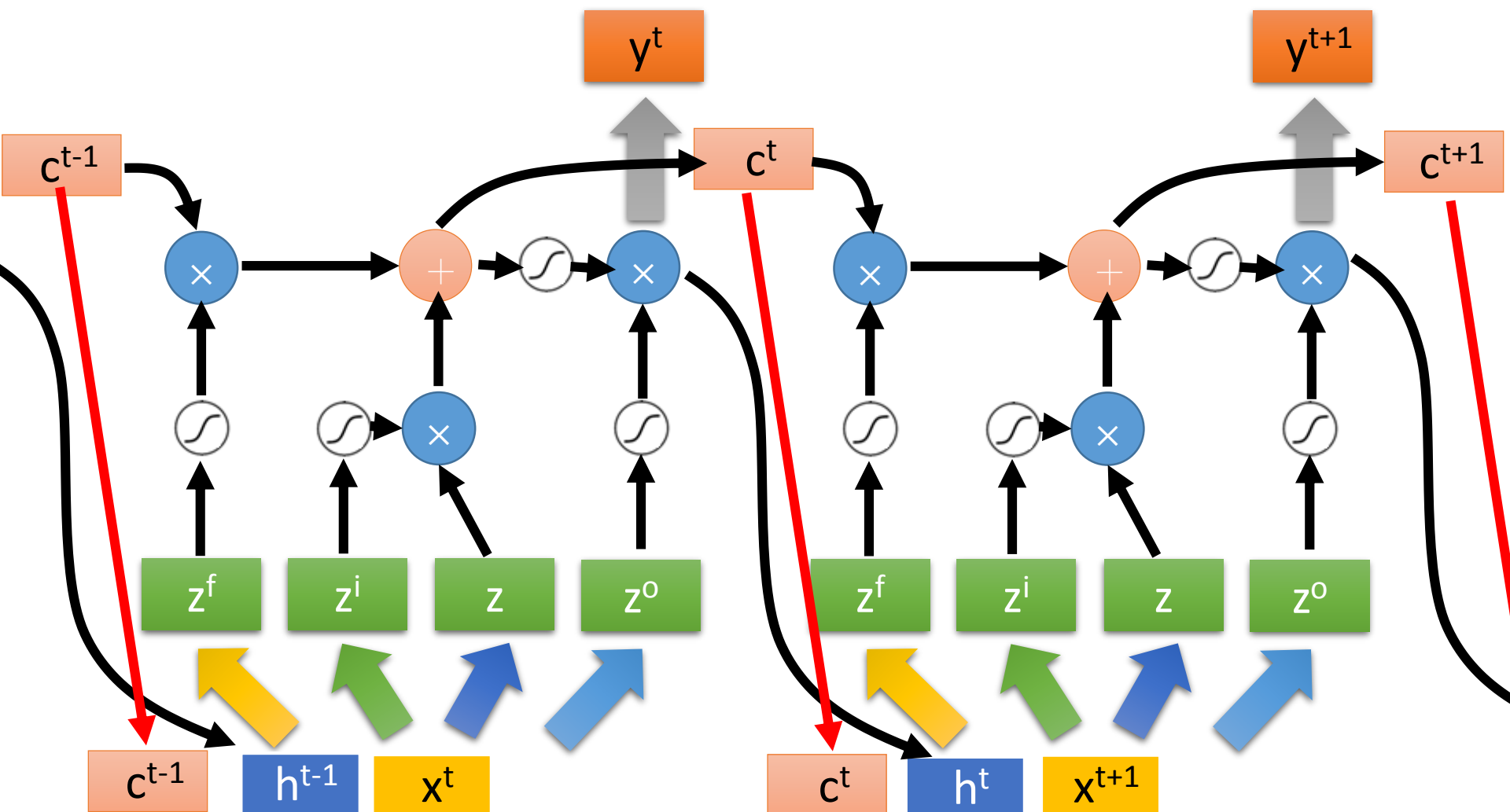


# LSTM



# LSTM

Extension: "peephole"



# Multiple-layer LSTM

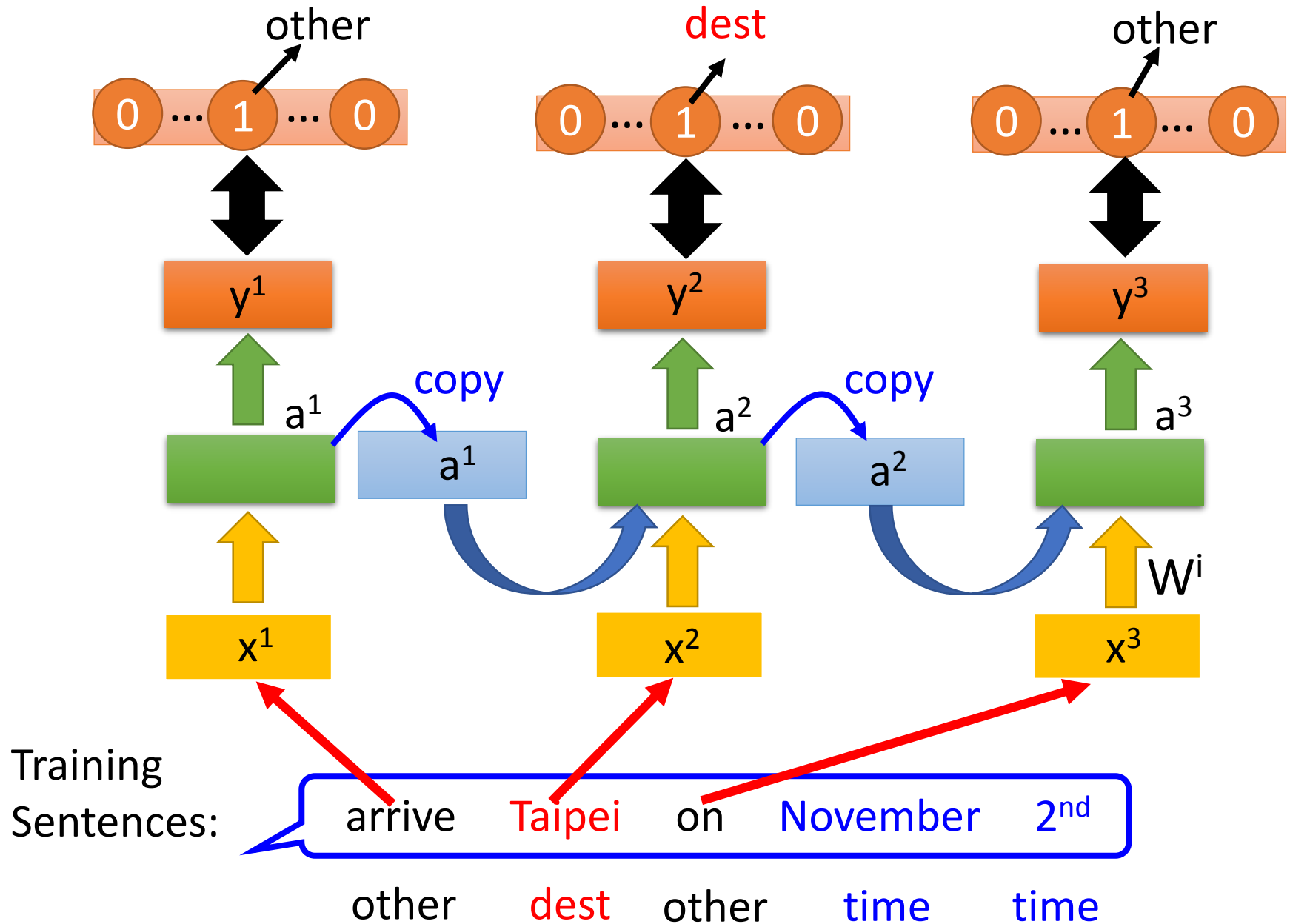
Don't worry if you cannot understand this.  
Keras can handle it.

Keras supports  
“LSTM”, “GRU”, “SimpleRNN” layers

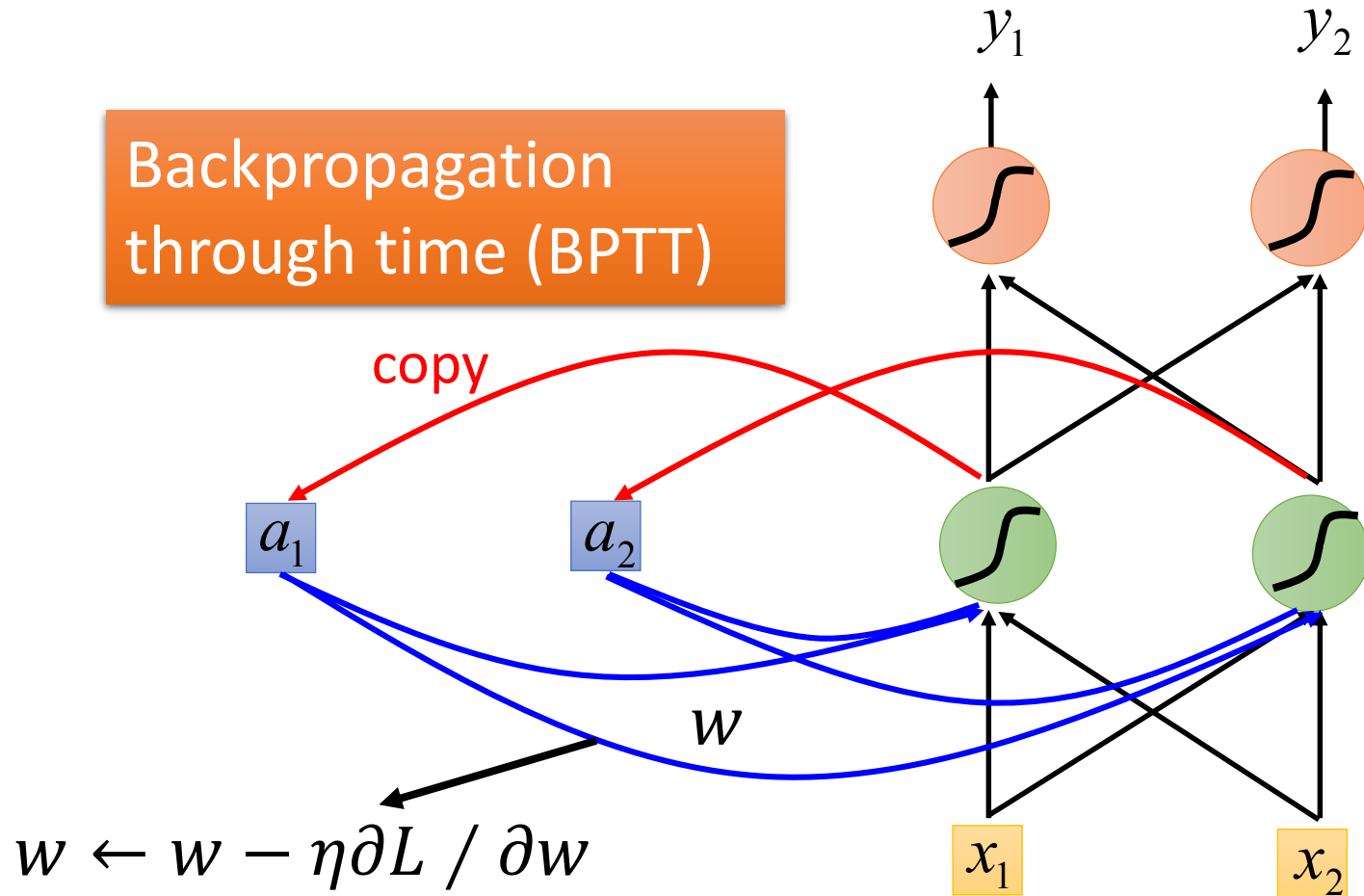
This is quite  
standard now.

我到底看了什麼？

# Learning Target



# Learning

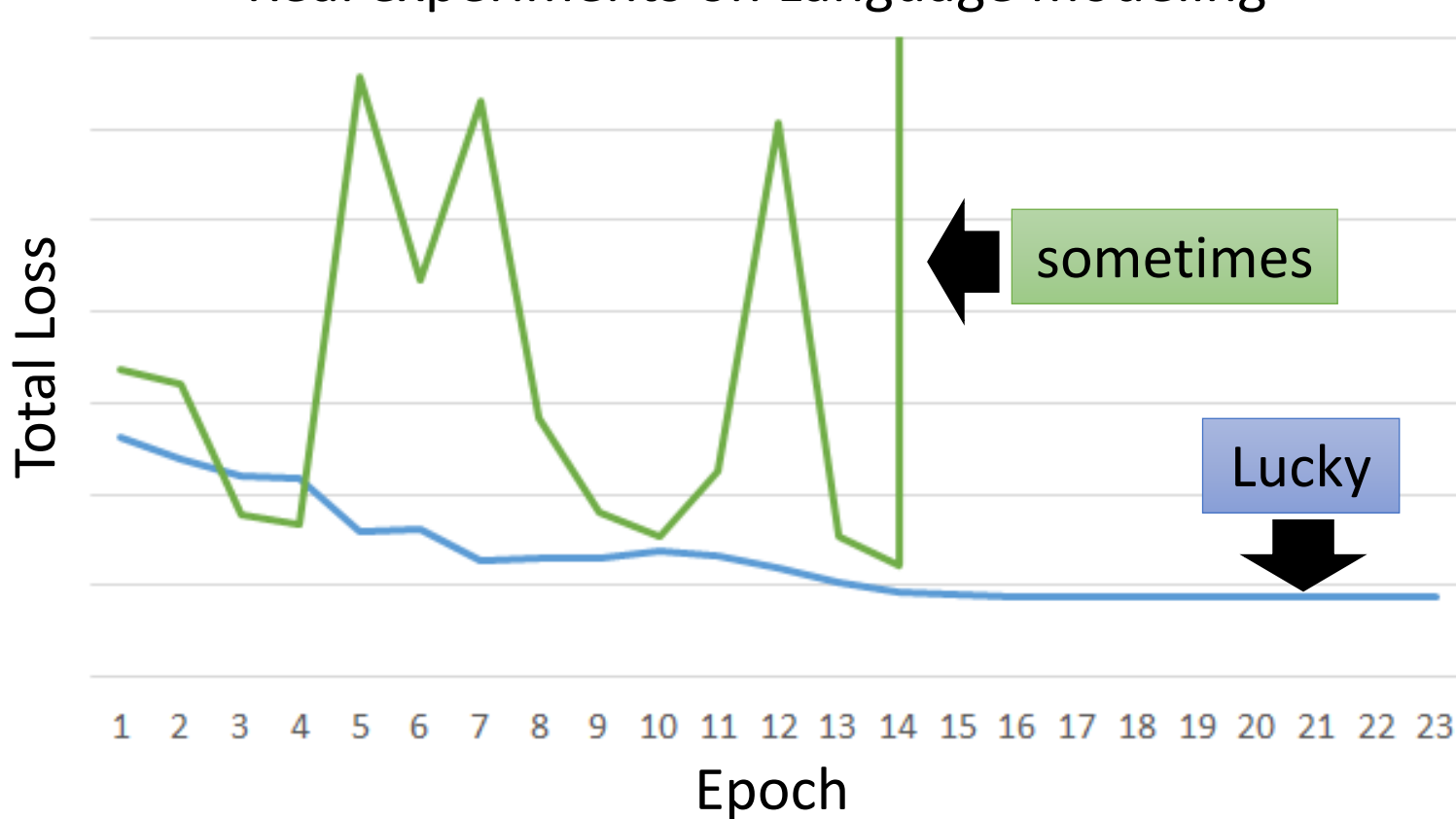




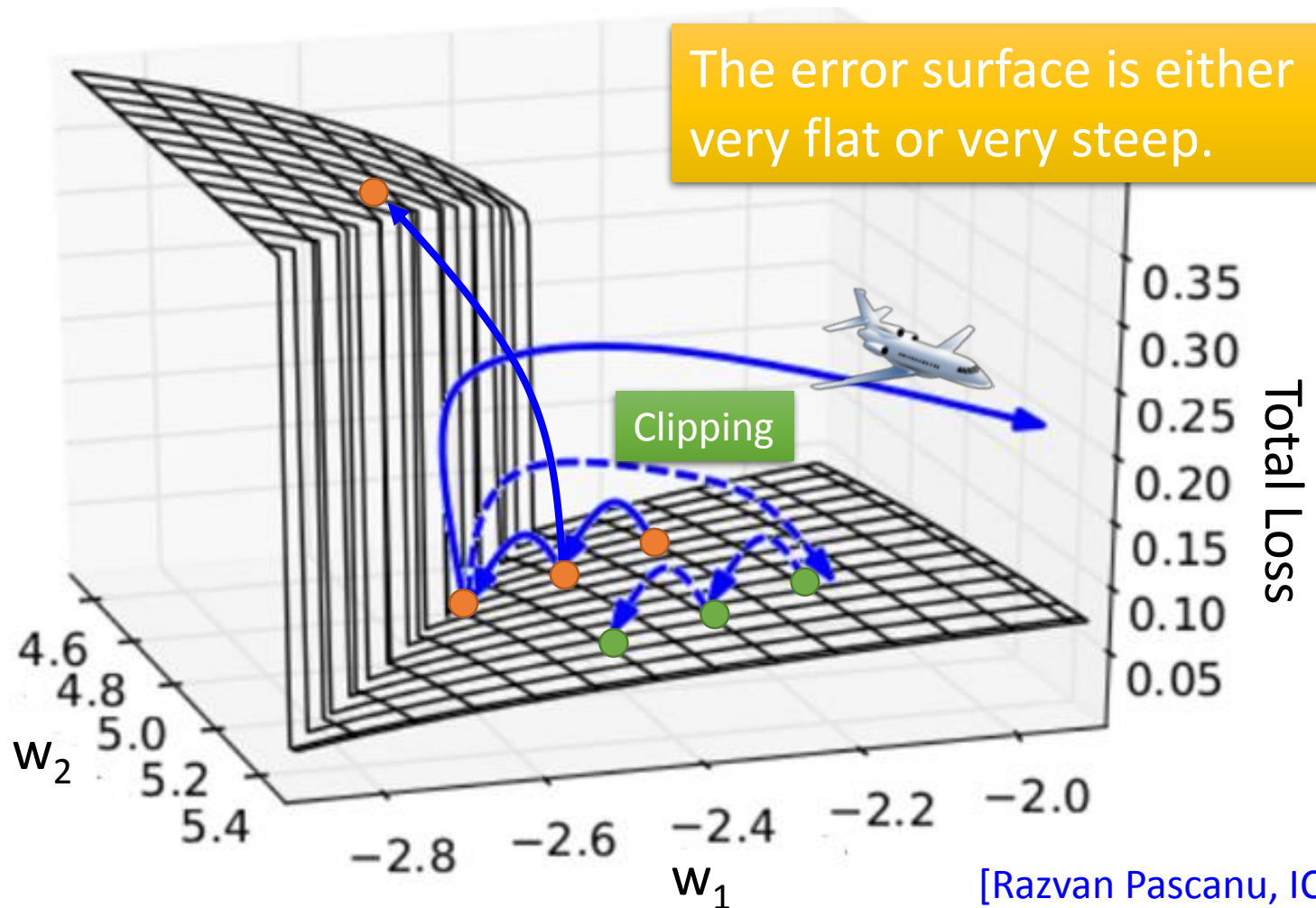
# Unfortunately .....

- RNN-based network is not always easy to learn

Real experiments on Language modeling



# The error surface is rough.



# Why?

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

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

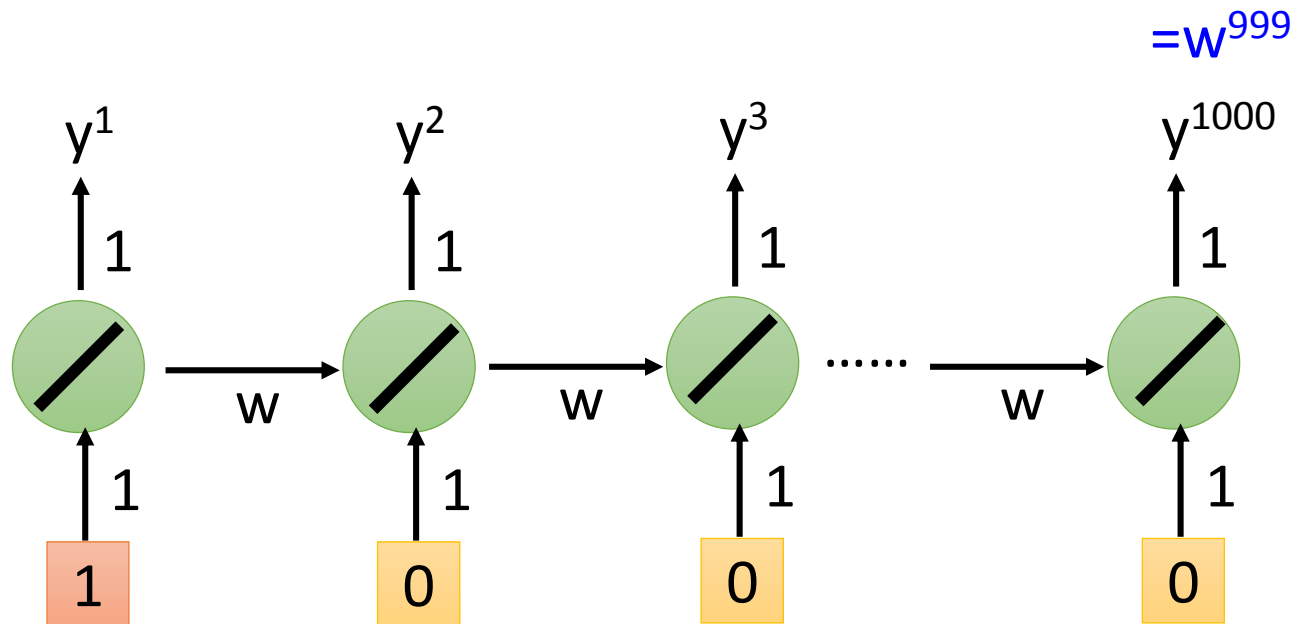
Large  
 $\partial L / \partial w$

Small  
Learning rate?

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)

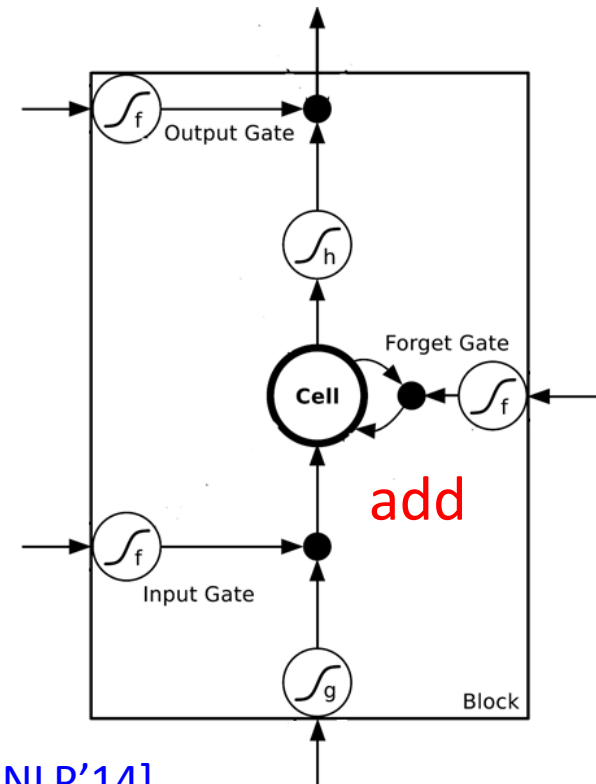
- Memory and input are

- added***

- The influence never disappears unless forget gate is closed

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

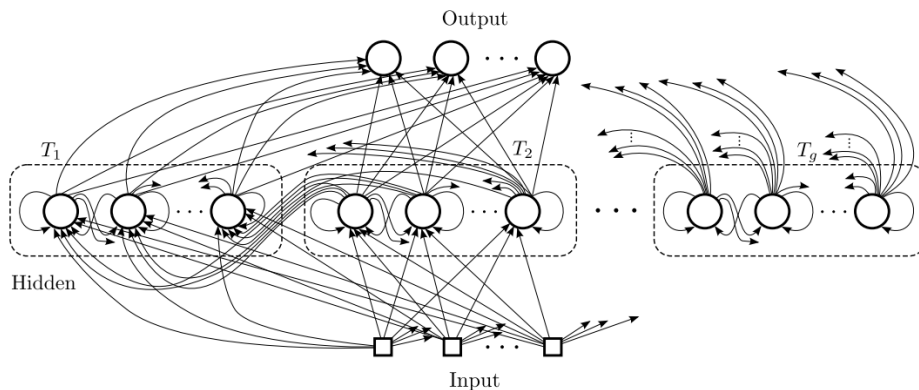
Gated Recurrent Unit (GRU):  
simpler than LSTM



[Cho, EMNLP'14]

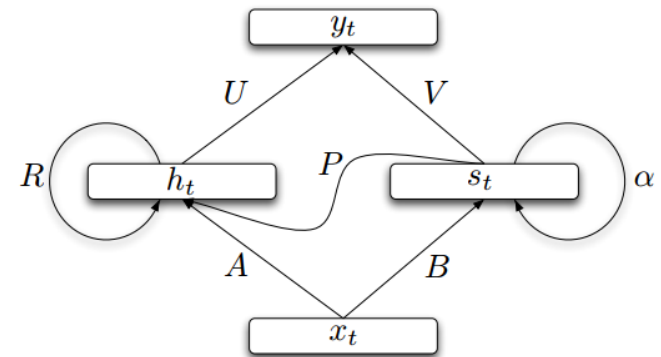
# Helpful Techniques

## Clockwise RNN



[Jan Koutnik, JMLR'14]

## Structurally Constrained Recurrent Network (SCRN)



[Tomas Mikolov, ICLR'15]

Vanilla RNN Initialized with Identity matrix + ReLU activation function [Quoc V. Le, arXiv'15]

➤ Outperform or be comparable with LSTM in 4 different tasks