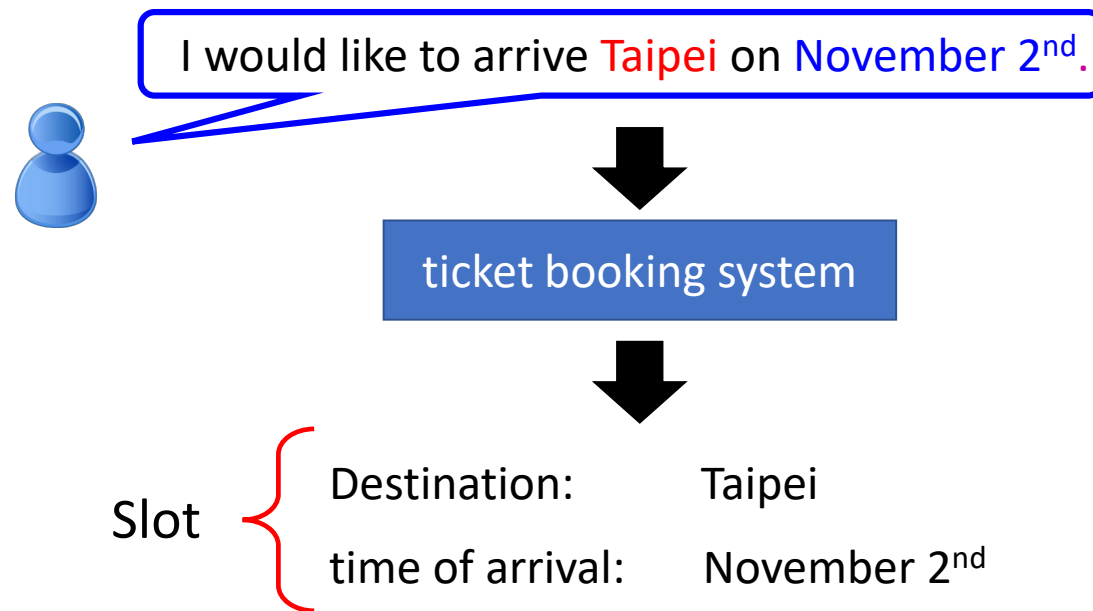


# Recurrent Neural Network (RNN)

# Example Application

- Slot Filling

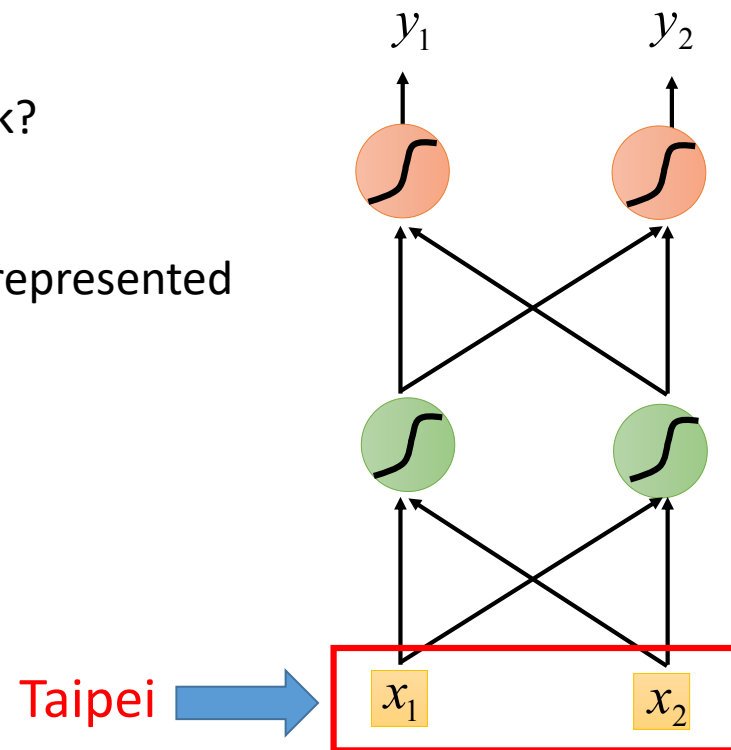


# Example Application

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.

apple = [ 1 0 0 0 0 ]

Each dimension corresponds  
to a word in the lexicon

bag = [ 0 1 0 0 0 ]

cat = [ 0 0 1 0 0 ]

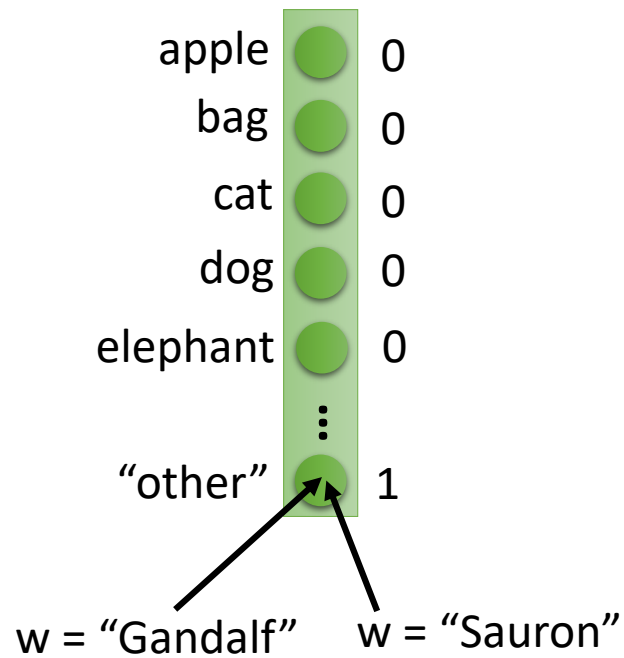
The dimension for the word  
is 1, and others are 0

dog = [ 0 0 0 1 0 ]

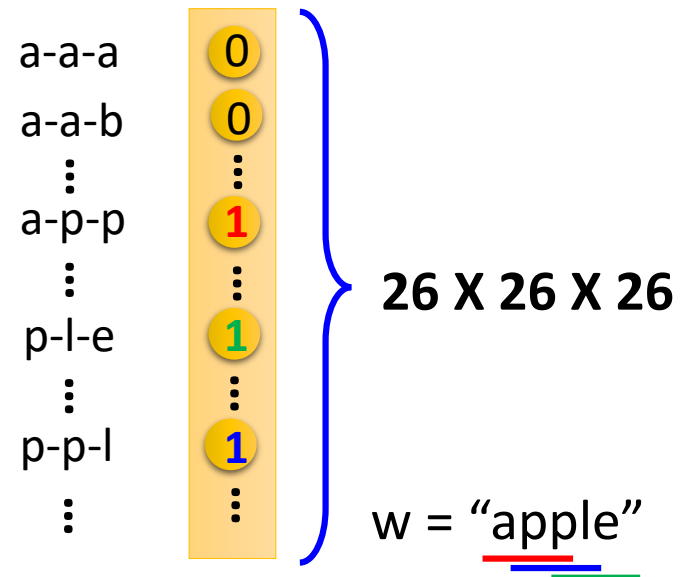
elephant = [ 0 0 0 0 1 ]

# Beyond 1-of-N encoding

## Dimension for "Other"



## Word hashing



# Example Application

Solving slot filling by  
Feedforward network?

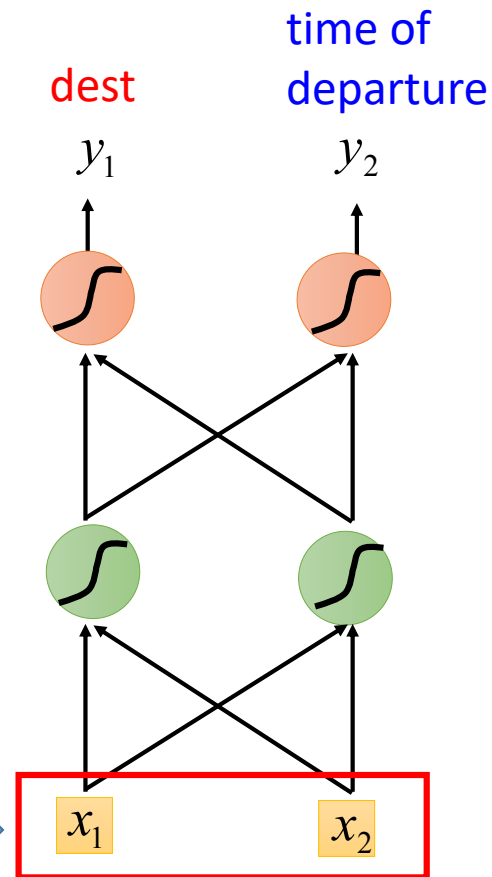
Input: a word

(Each word is represented  
as a vector)

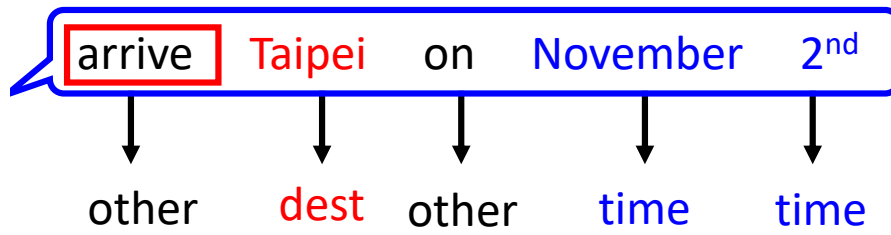
Output:

Probability distribution that  
the input word belonging to  
the slots

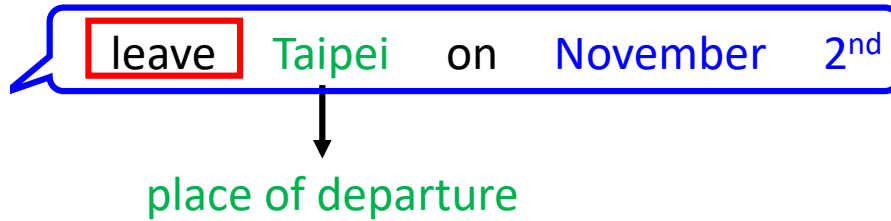
Taipei →



# Example Application

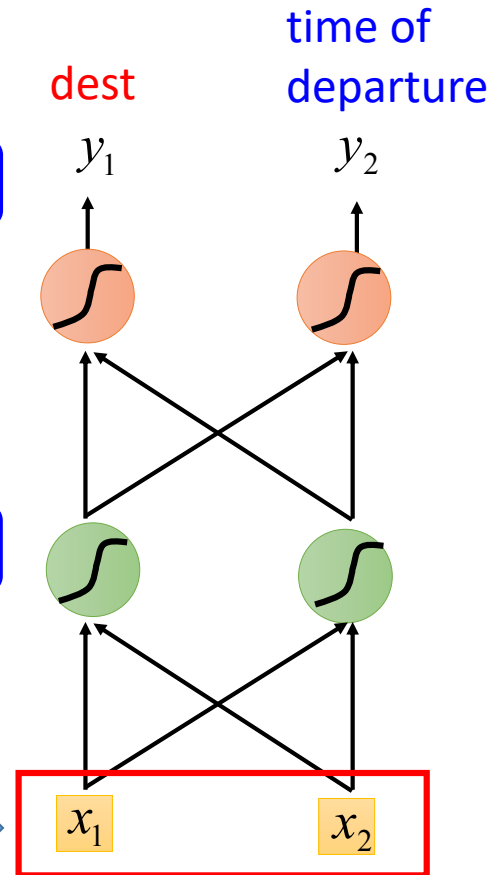


Problem?



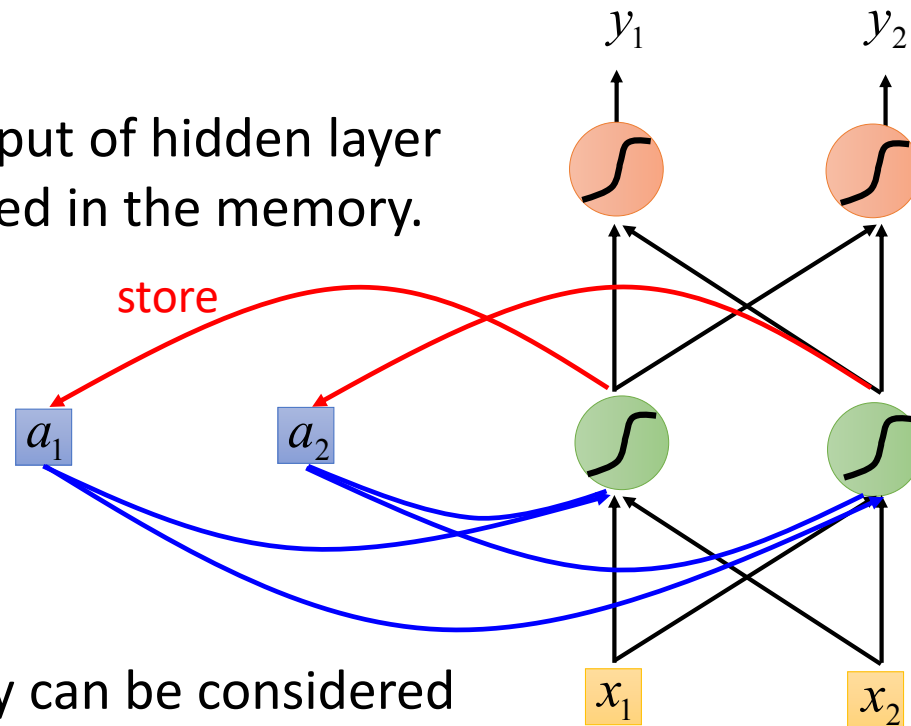
Neural network needs memory!

Taipei →



# Recurrent Neural Network (RNN)

The output of hidden layer are stored in the memory.

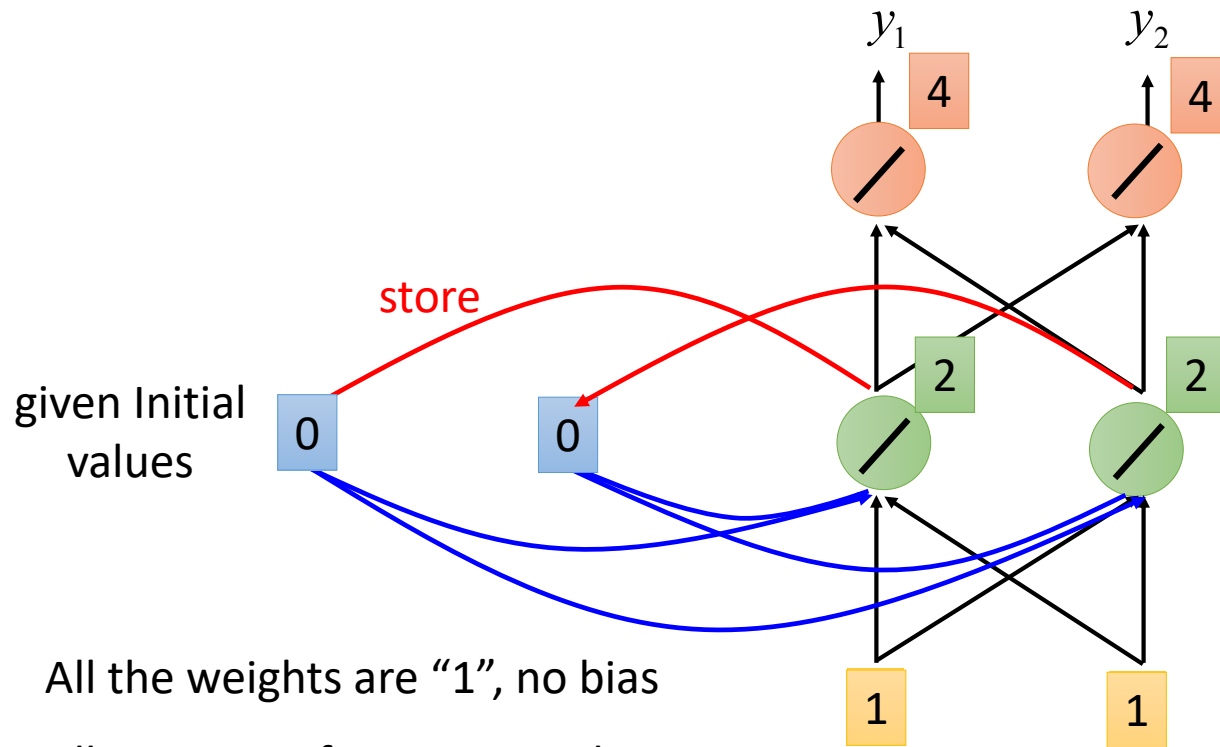


Memory can be considered as another input.



# Example

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



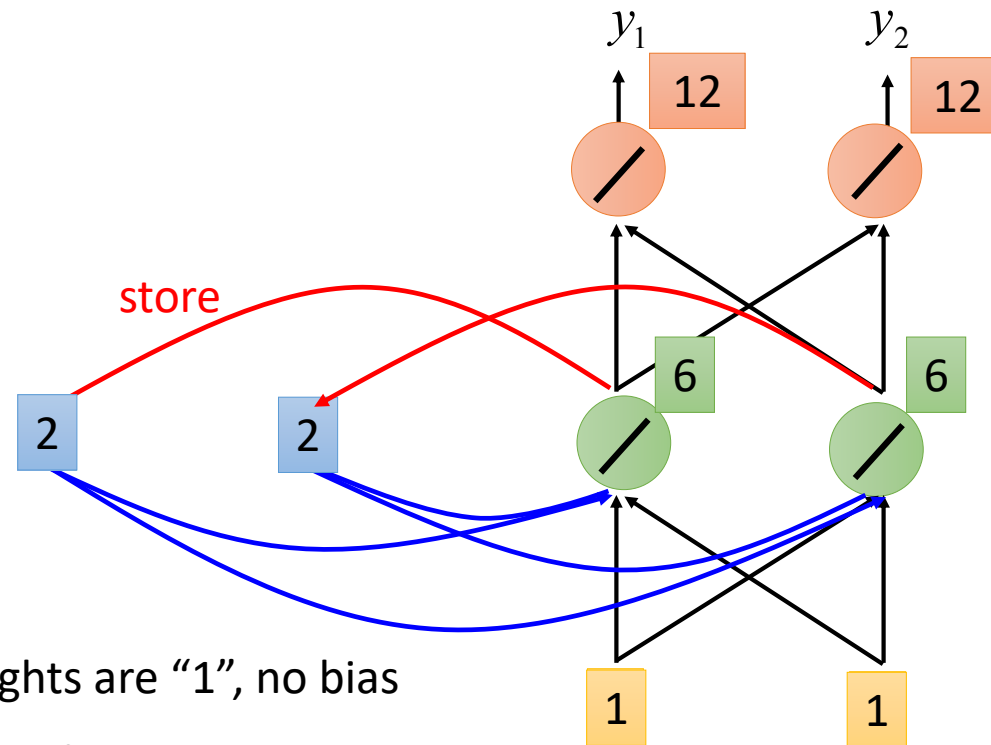
All the weights are "1", no bias

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

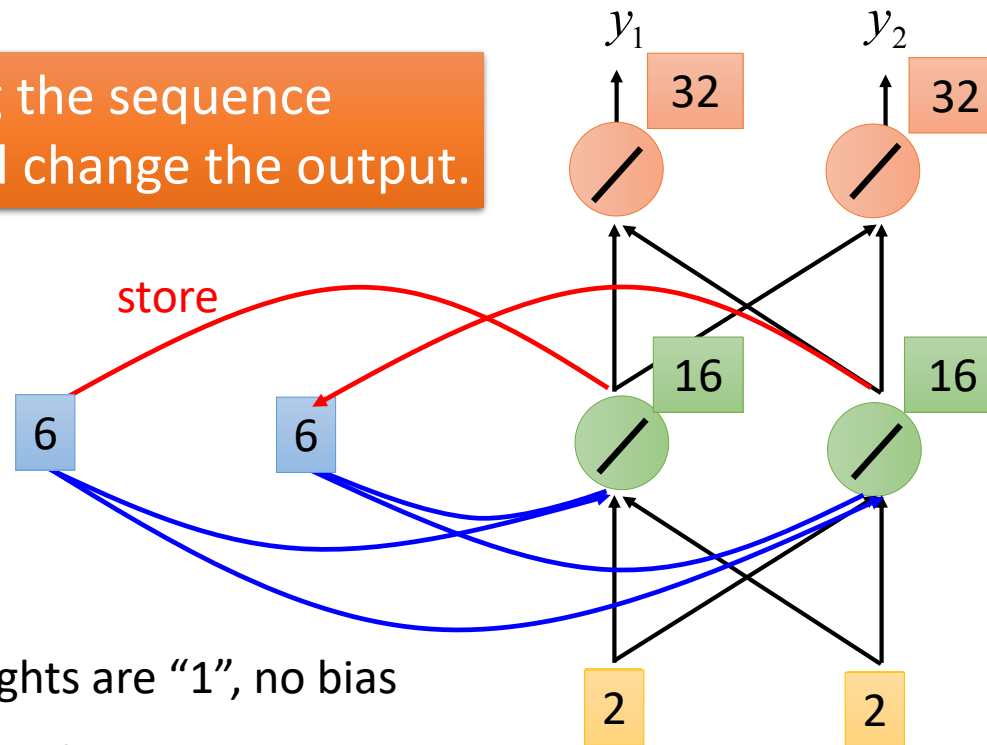
All activation functions are linear

# Example

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

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

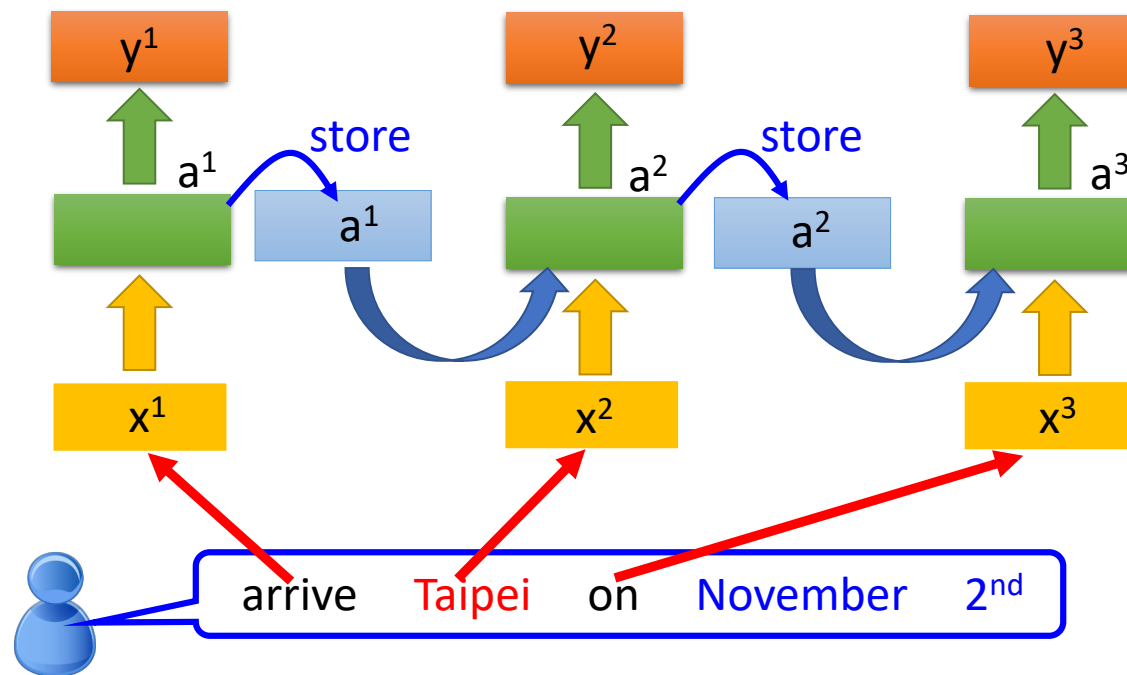
# RNN

The same network is used again and again.

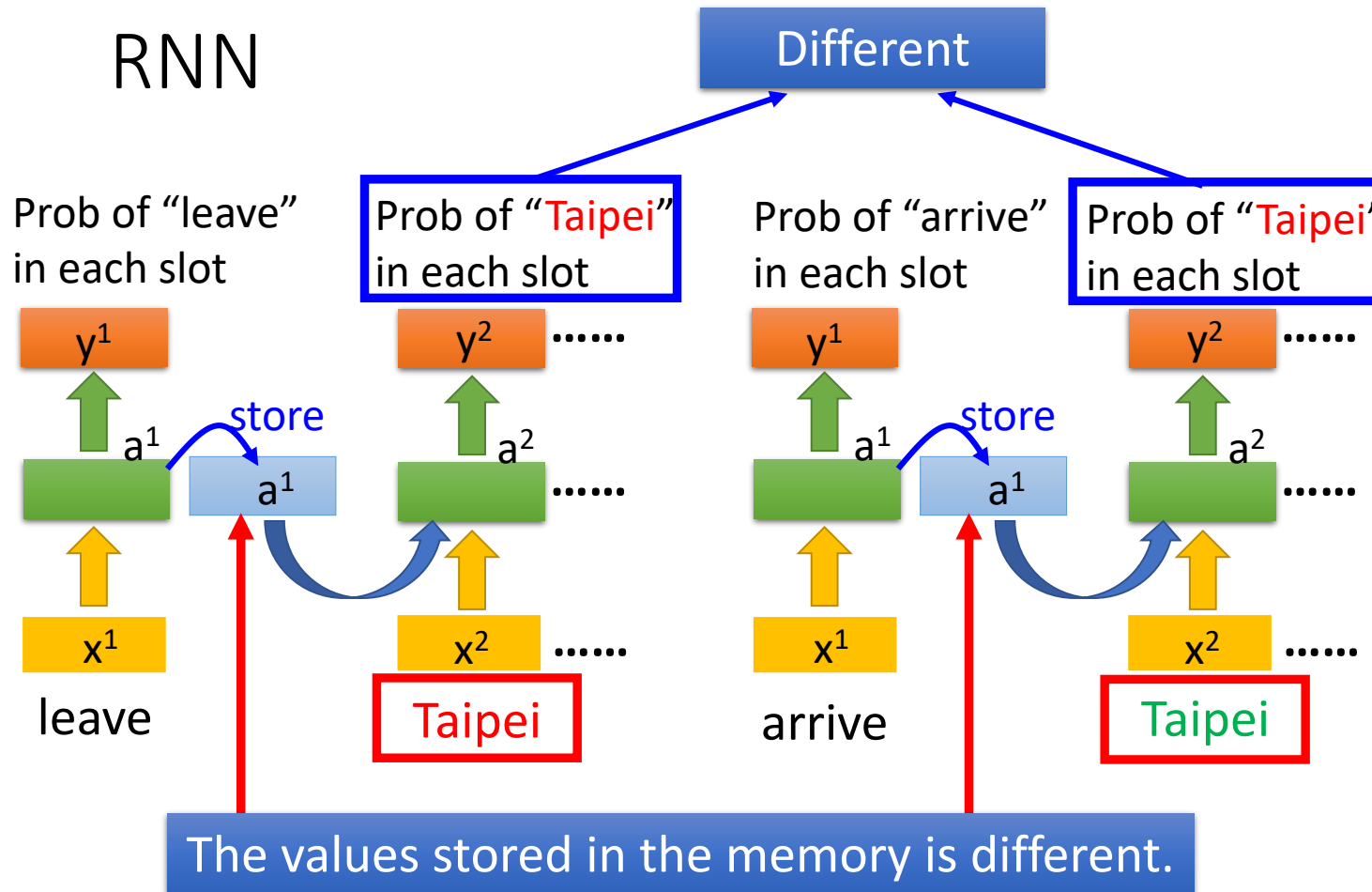
Probability of  
"arrive" in each slot

Probability of  
"Taipei" in each slot

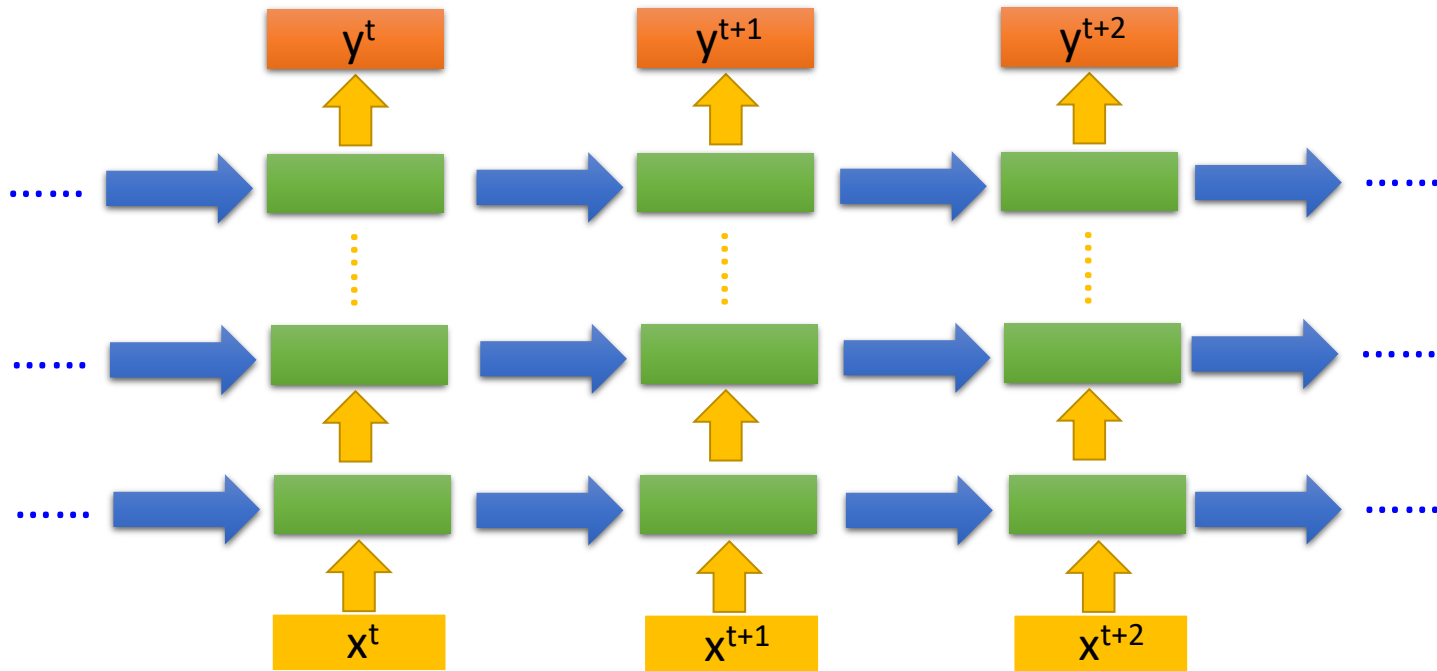
Probability of  
"on" in each slot



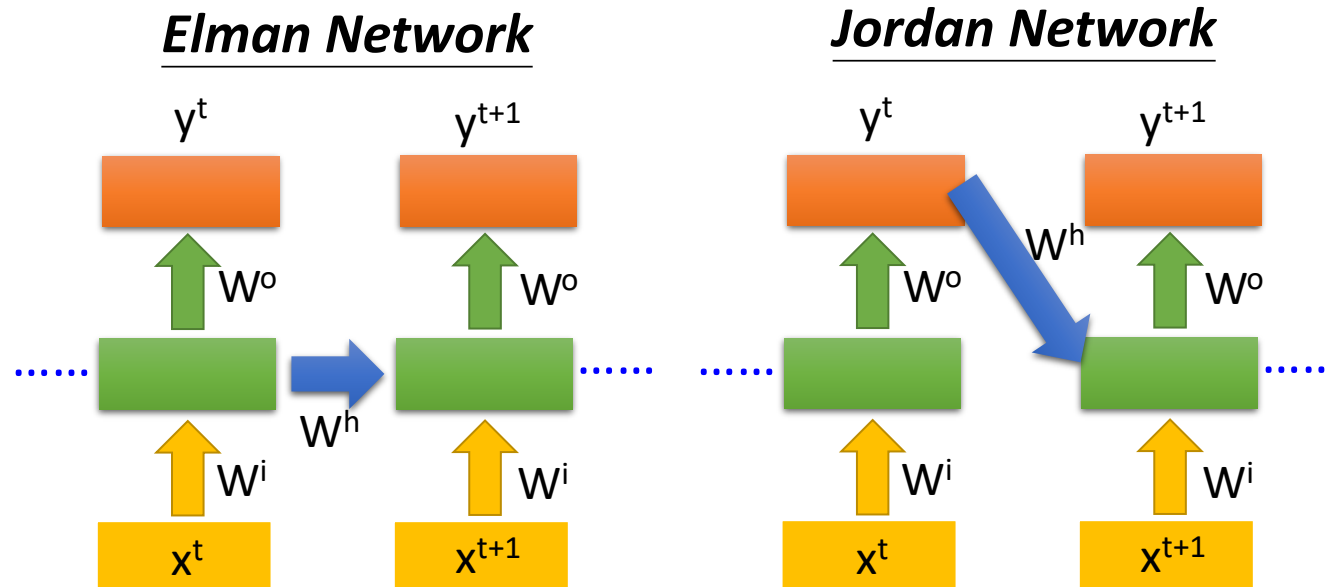
# RNN



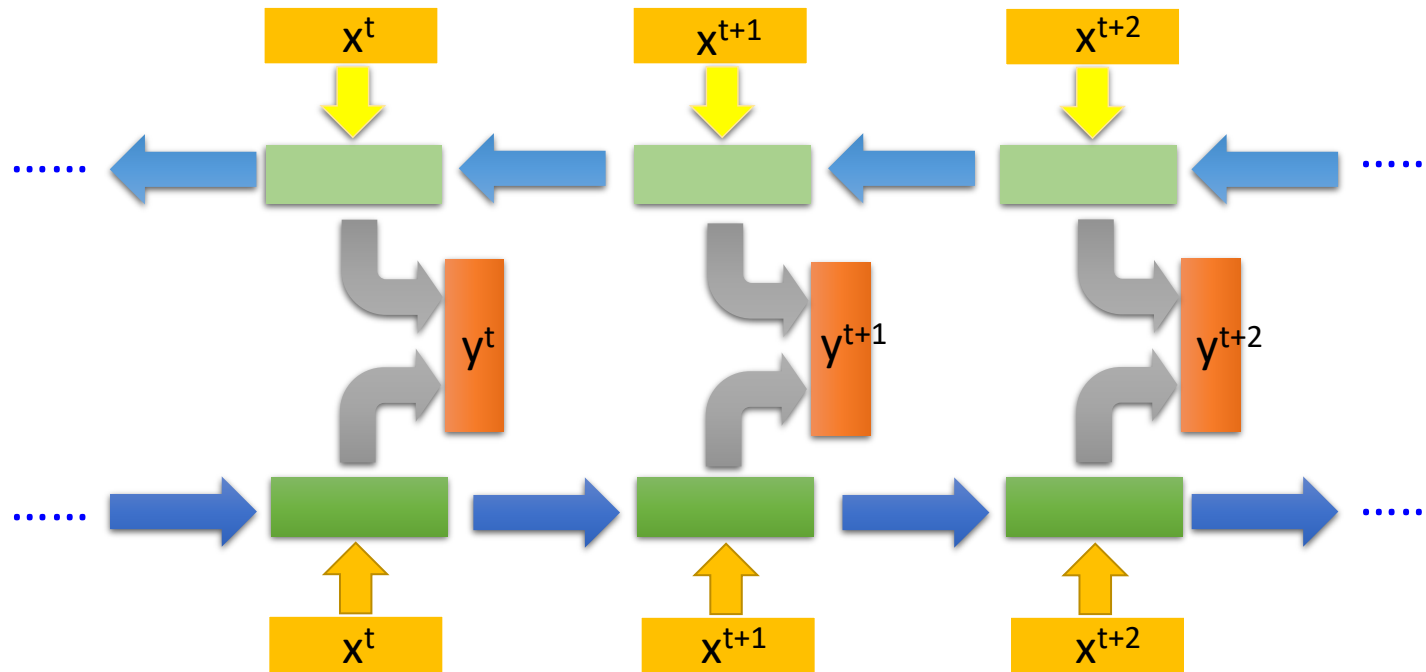
Of course it can be deep ...



# Elman Network & Jordan Network

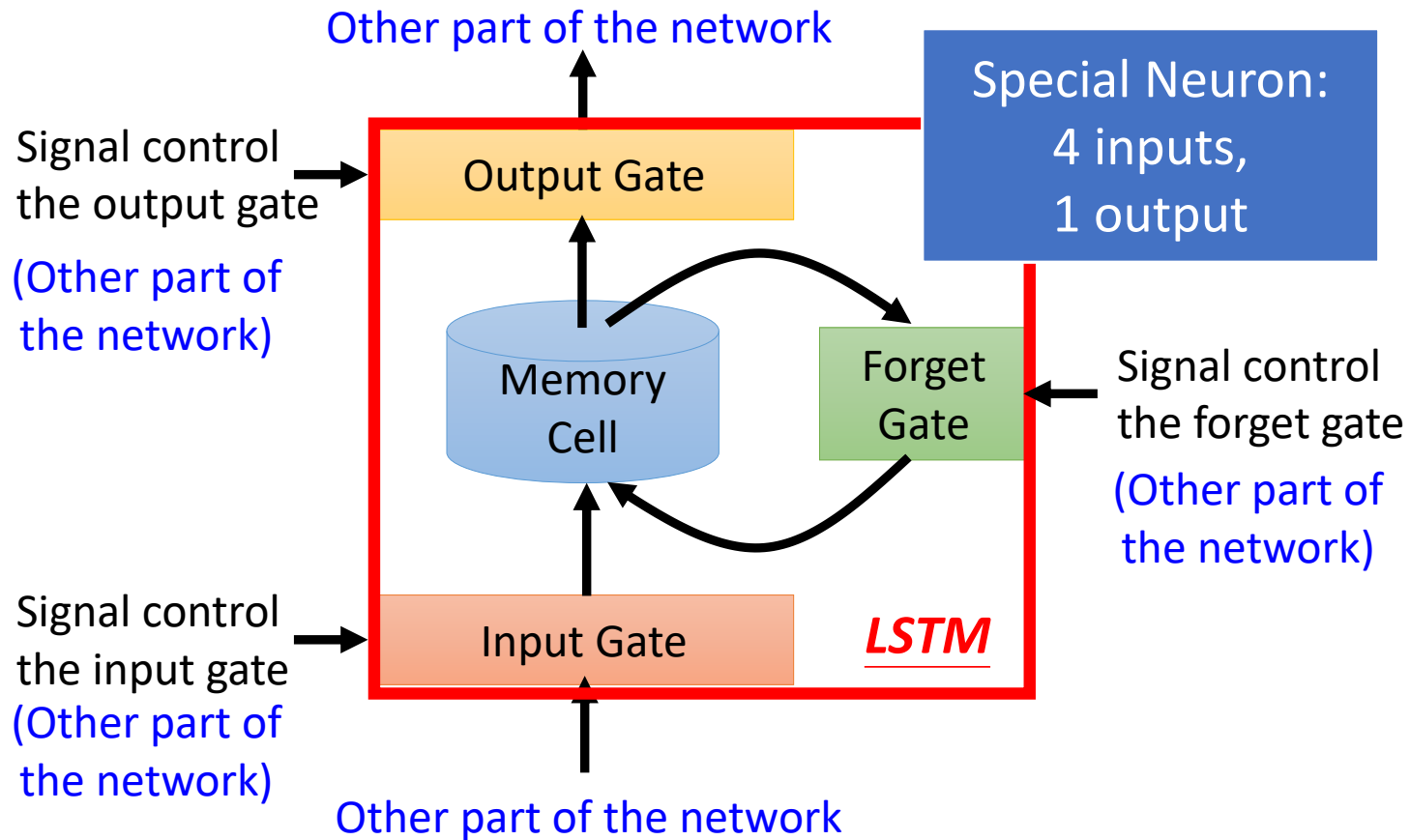


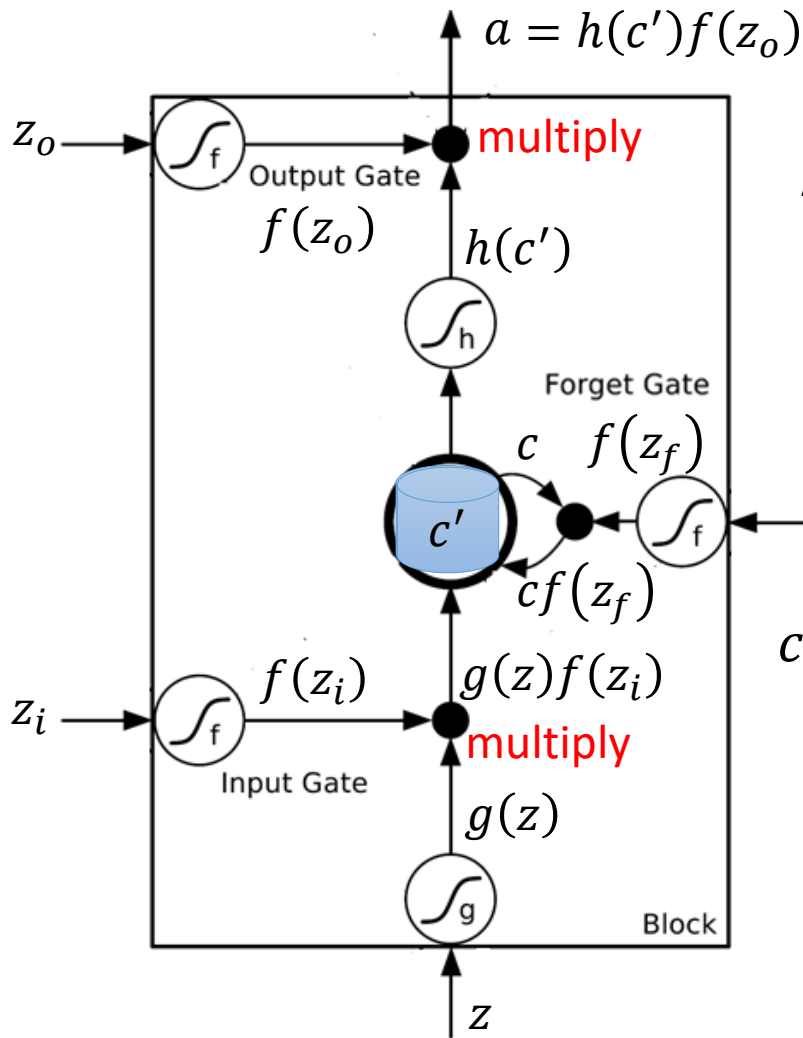
# Bidirectional RNN





# Long Short-term Memory (LSTM)





Activation function  $f$  is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

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

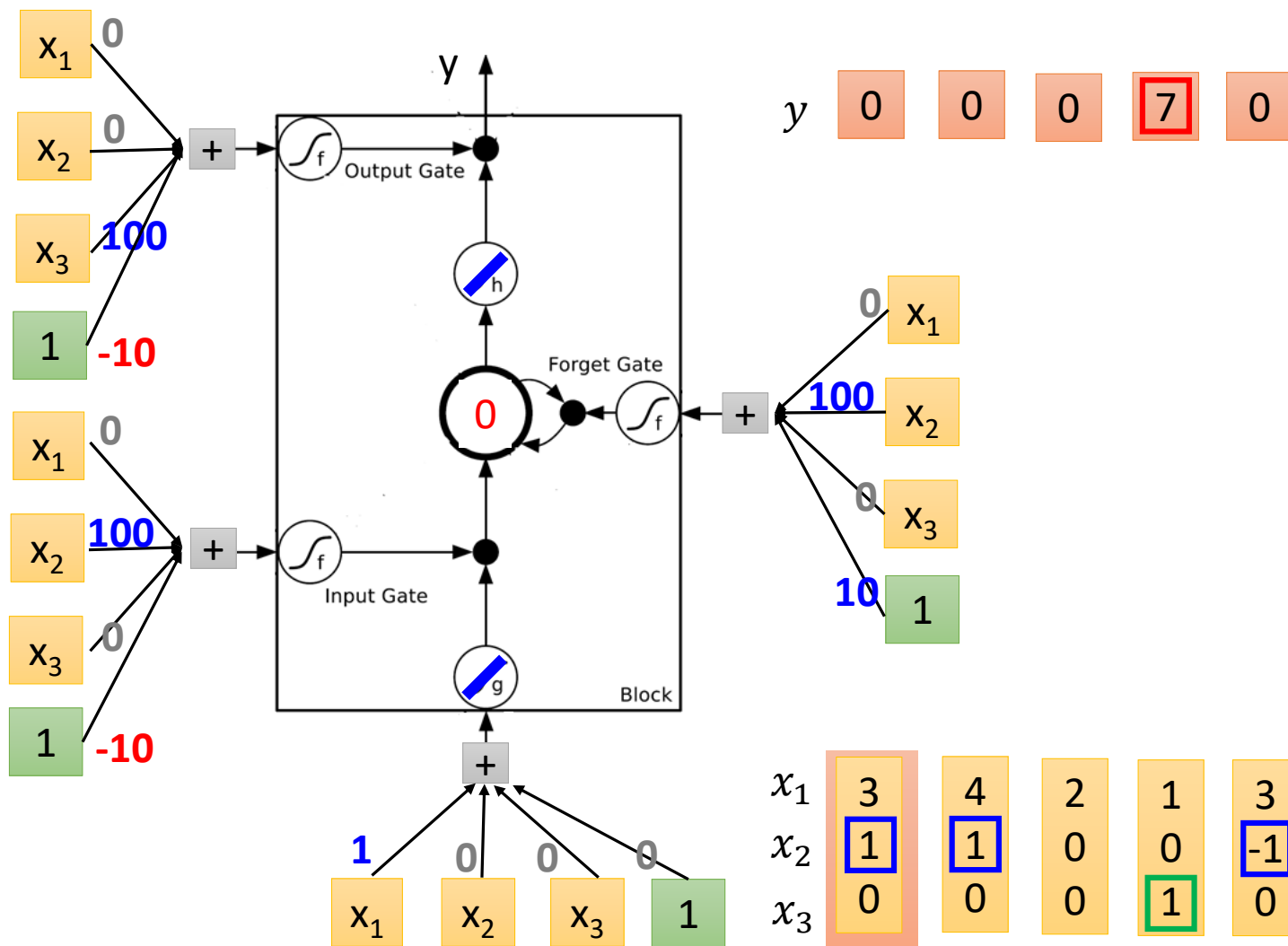
# LSTM - Example

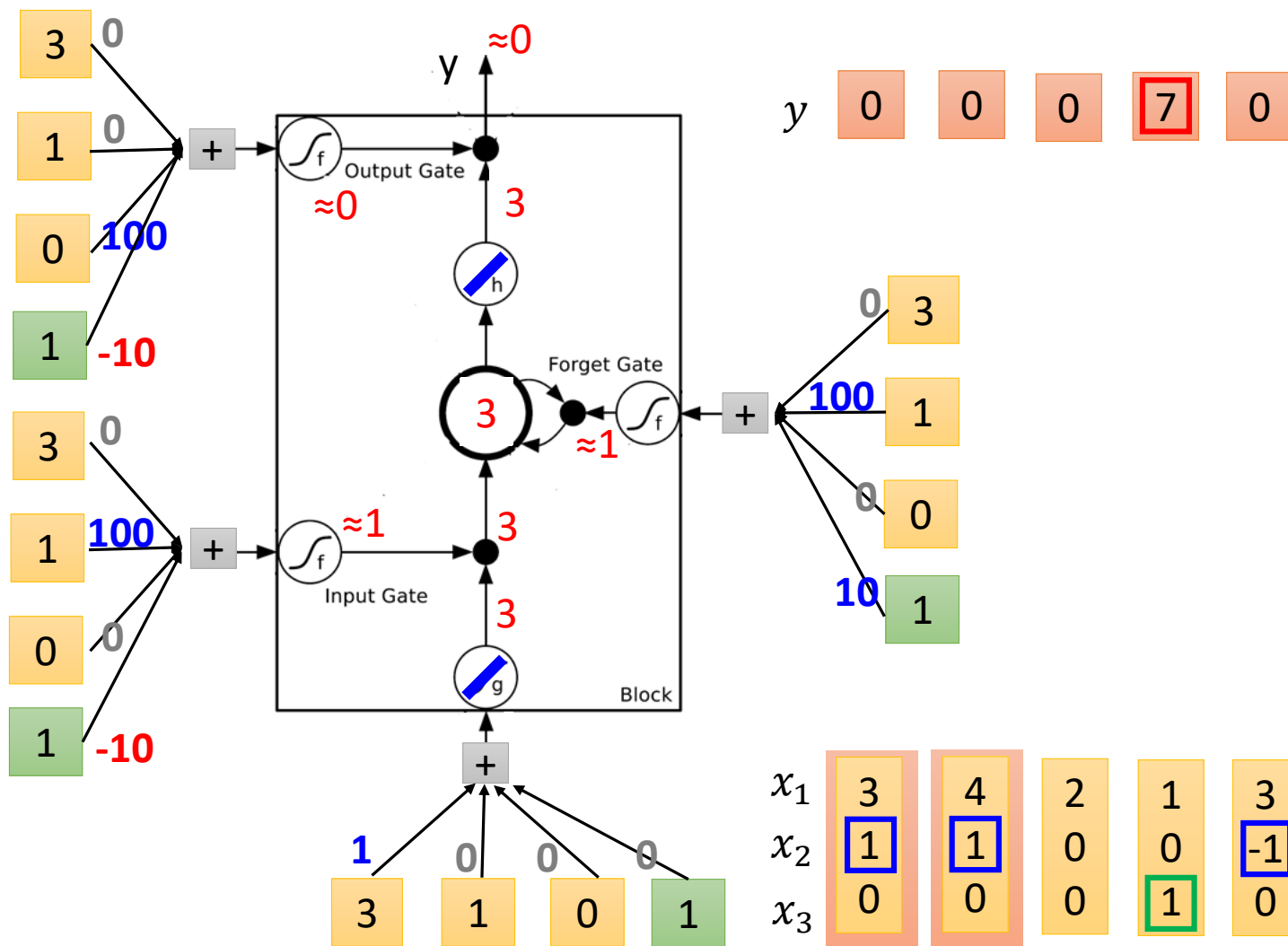
	0	0	3	3	7	7	7	0	6
$x_1$	1	3	2	4	2	1	3	6	1
$x_2$	0	1	0	1	0	0	-1	1	0
$x_3$	0	0	0	0	0	1	0	0	1
$y$	0	0	0	0	0	7	0	0	6

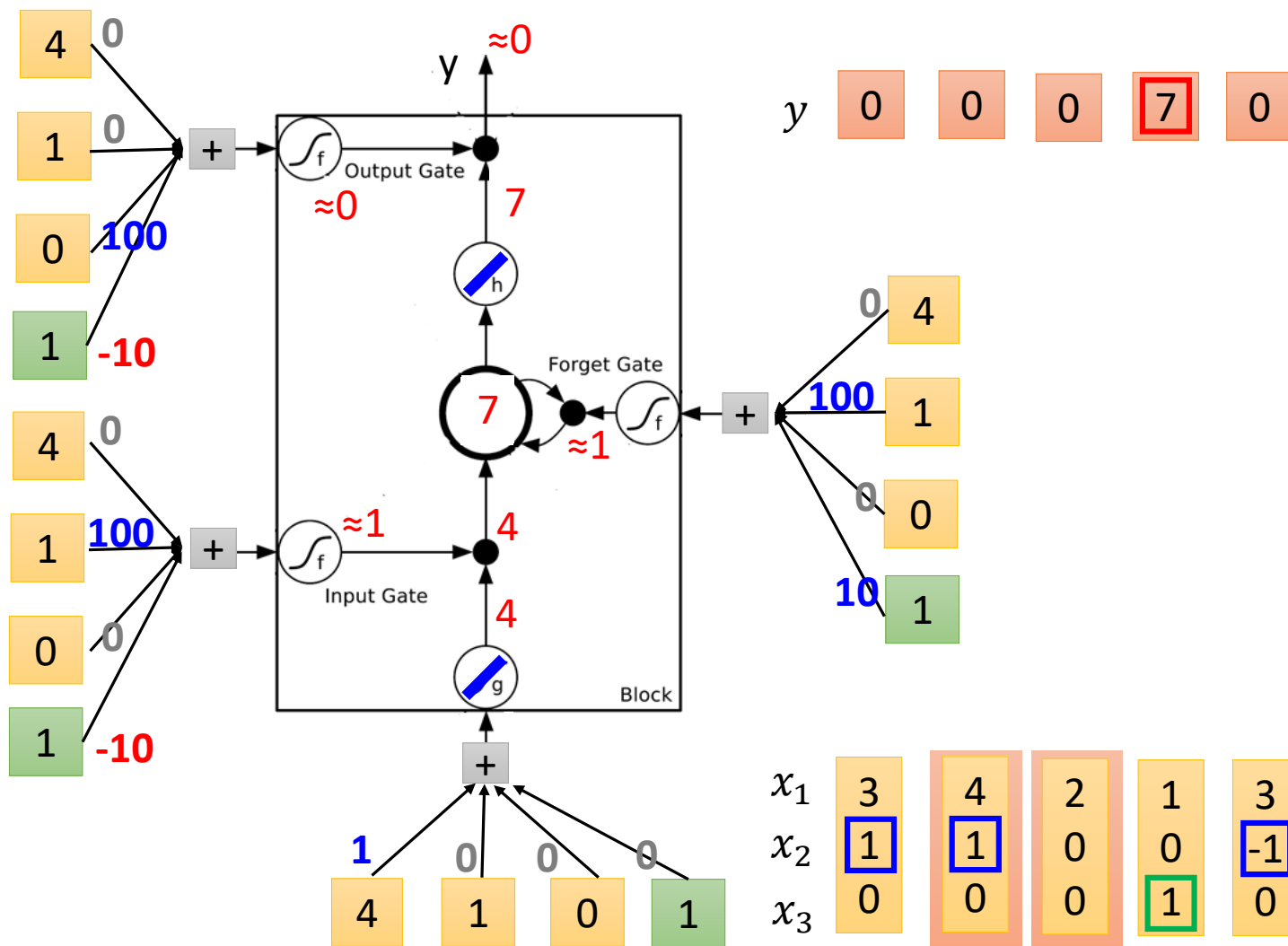
When  $x_2 = 1$ , add the numbers of  $x_1$  into the memory

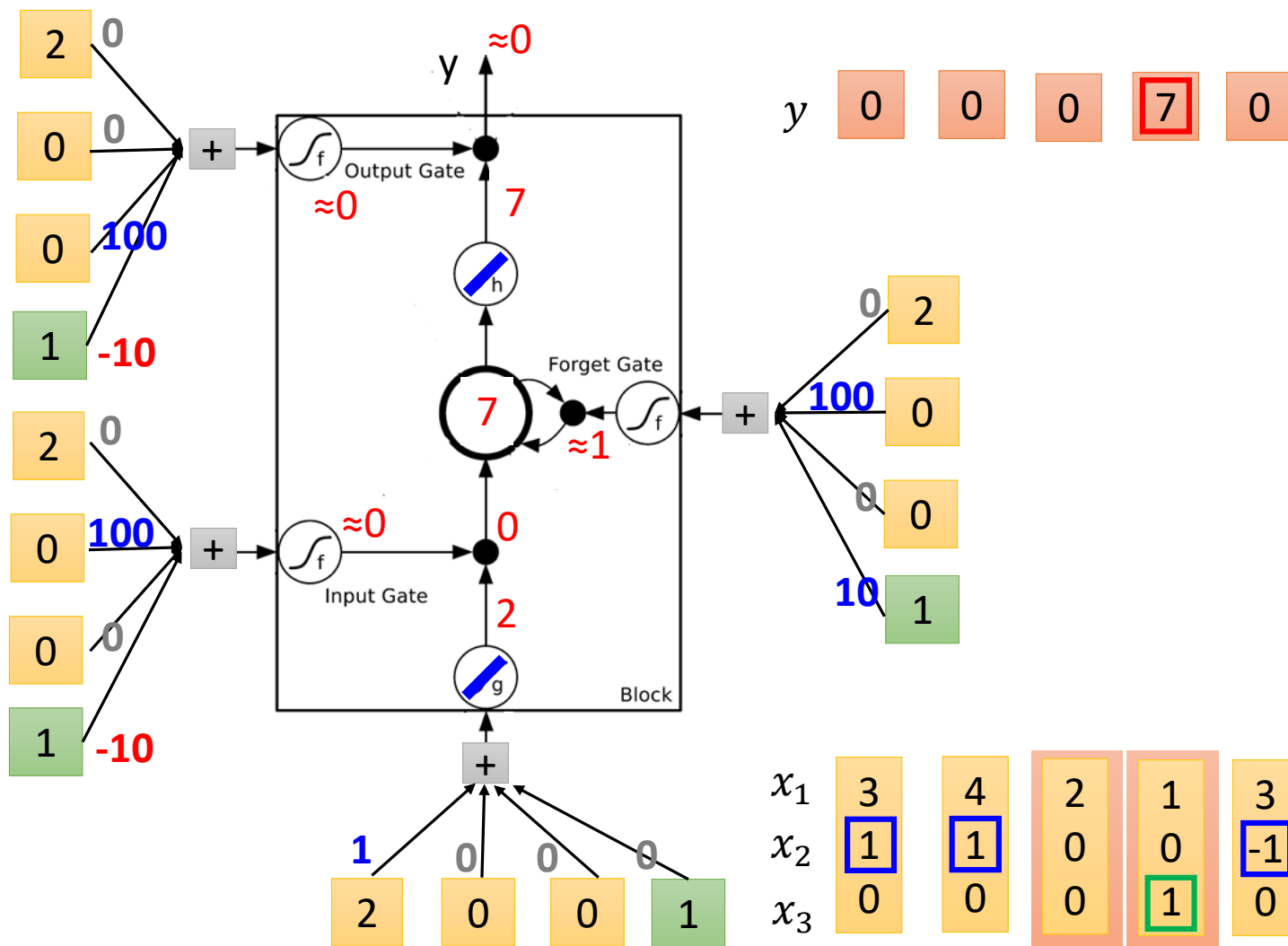
When  $x_2 = -1$ , reset the memory

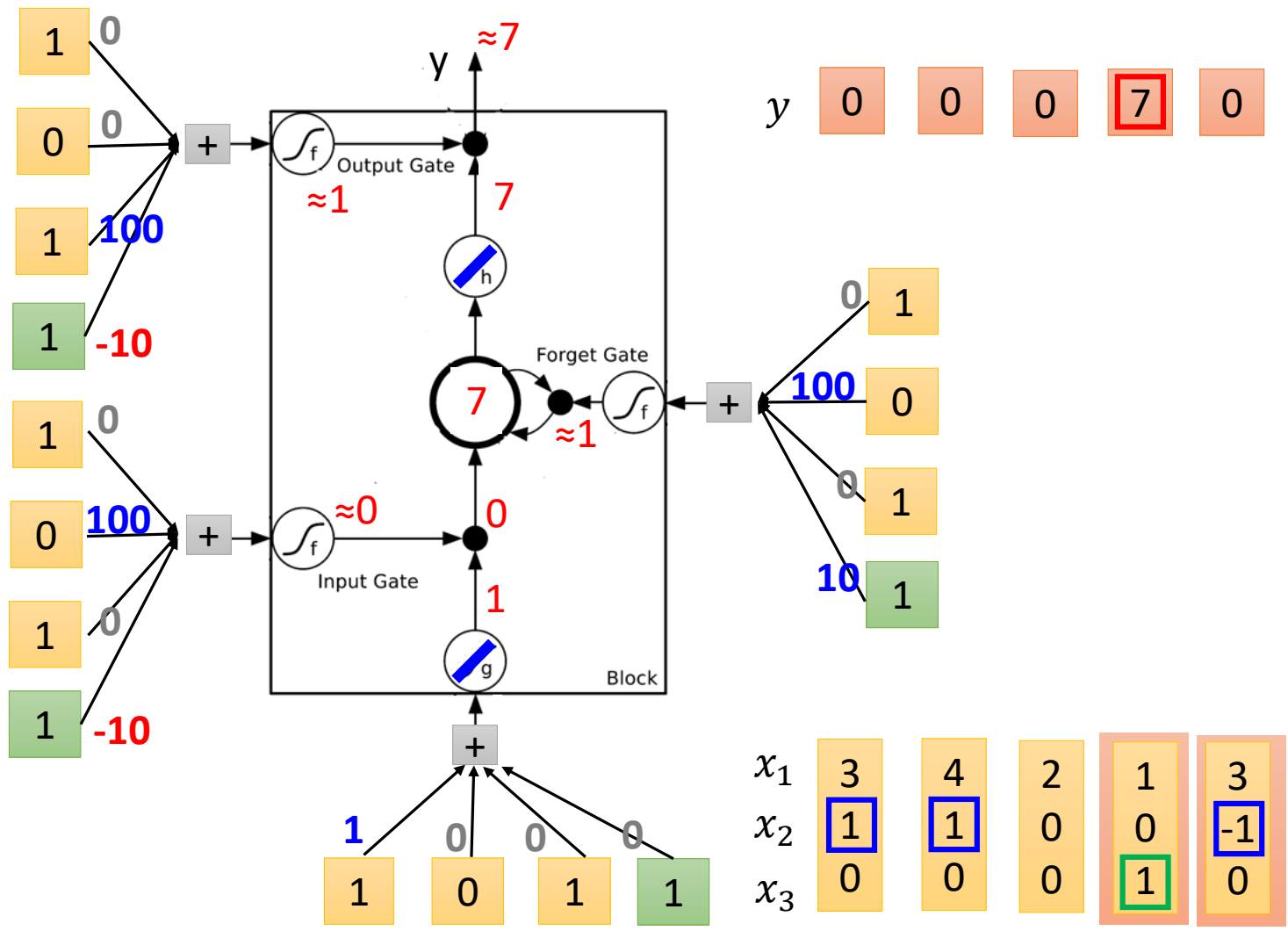
When  $x_3 = 1$ , output the number in the memory.



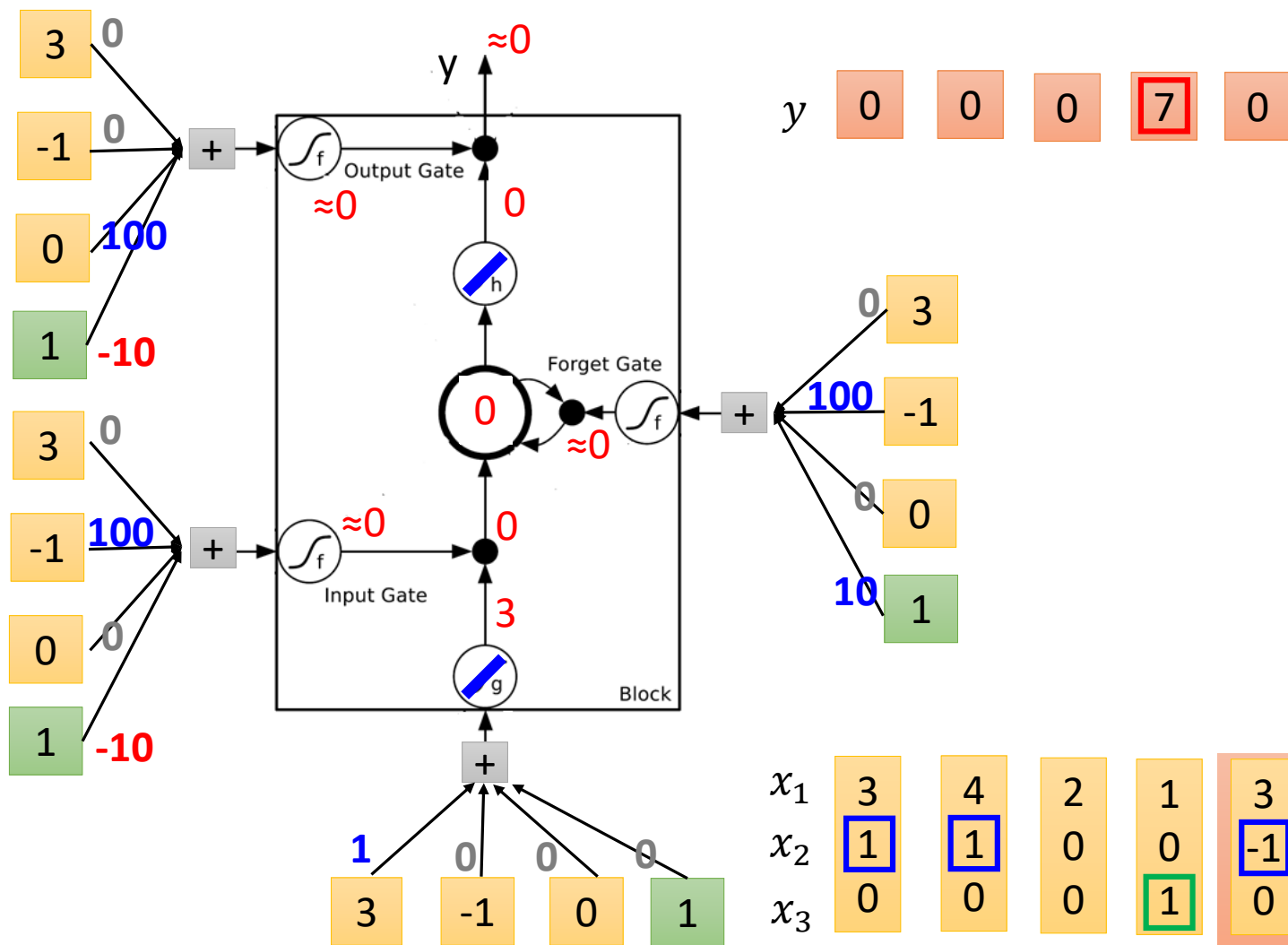






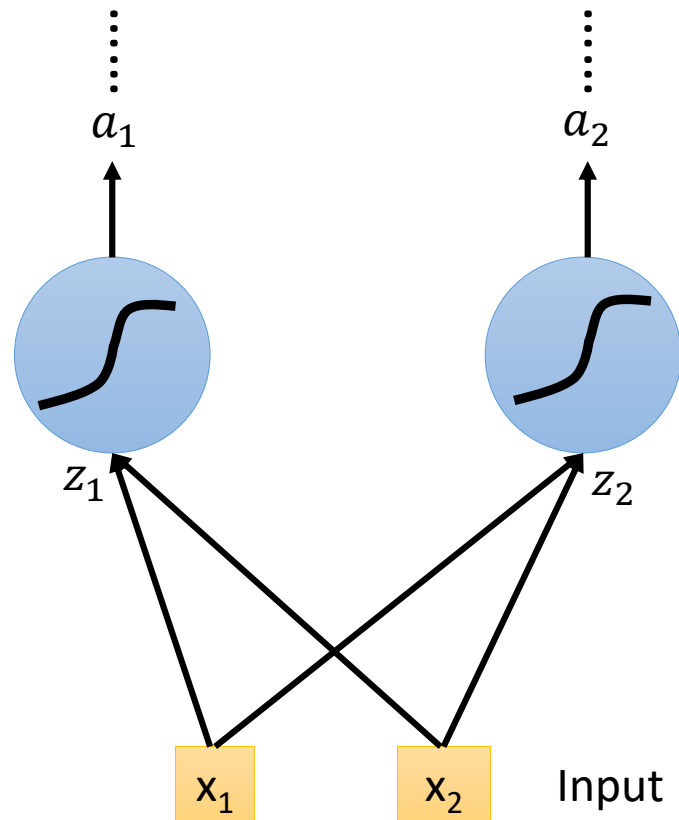


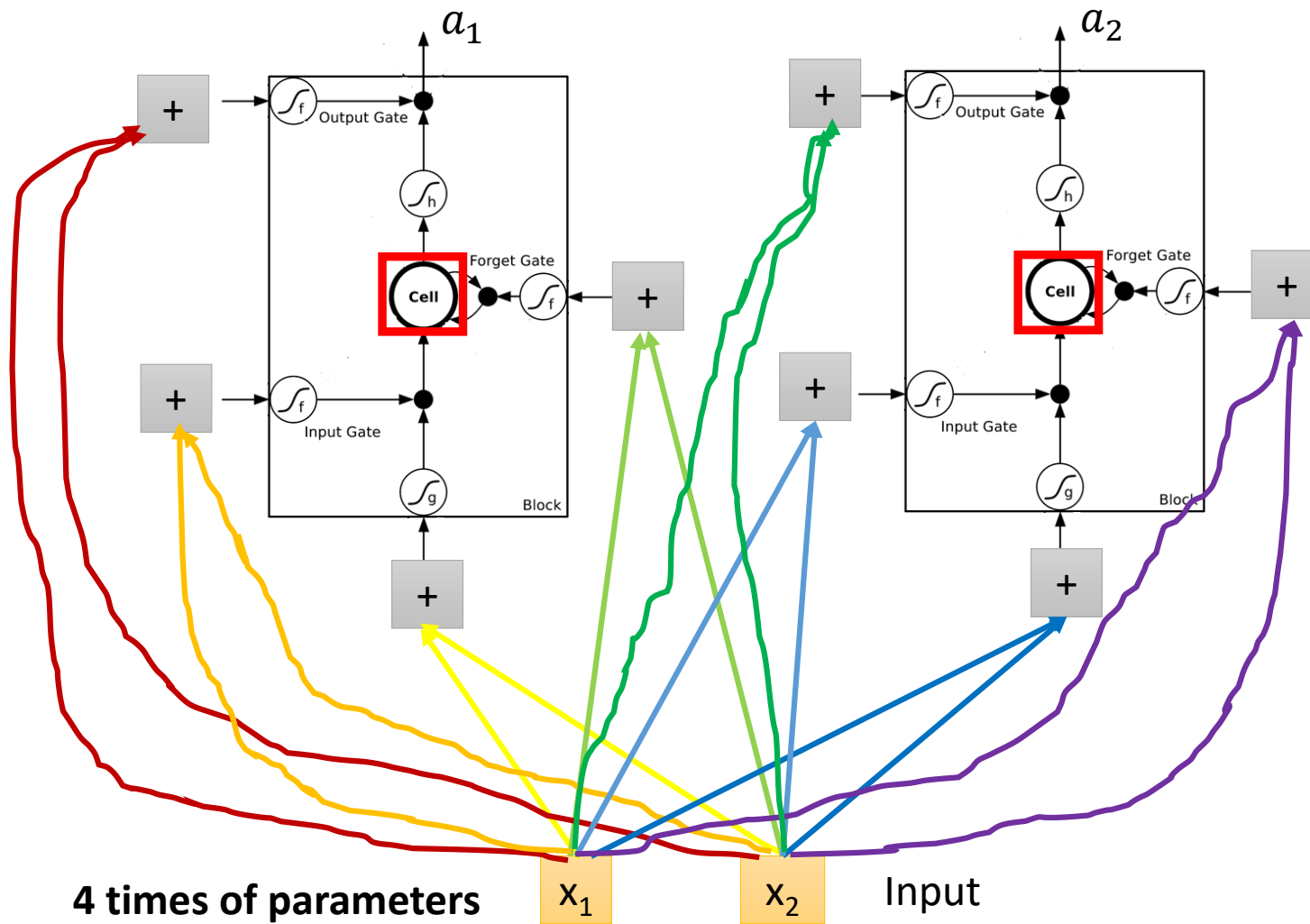




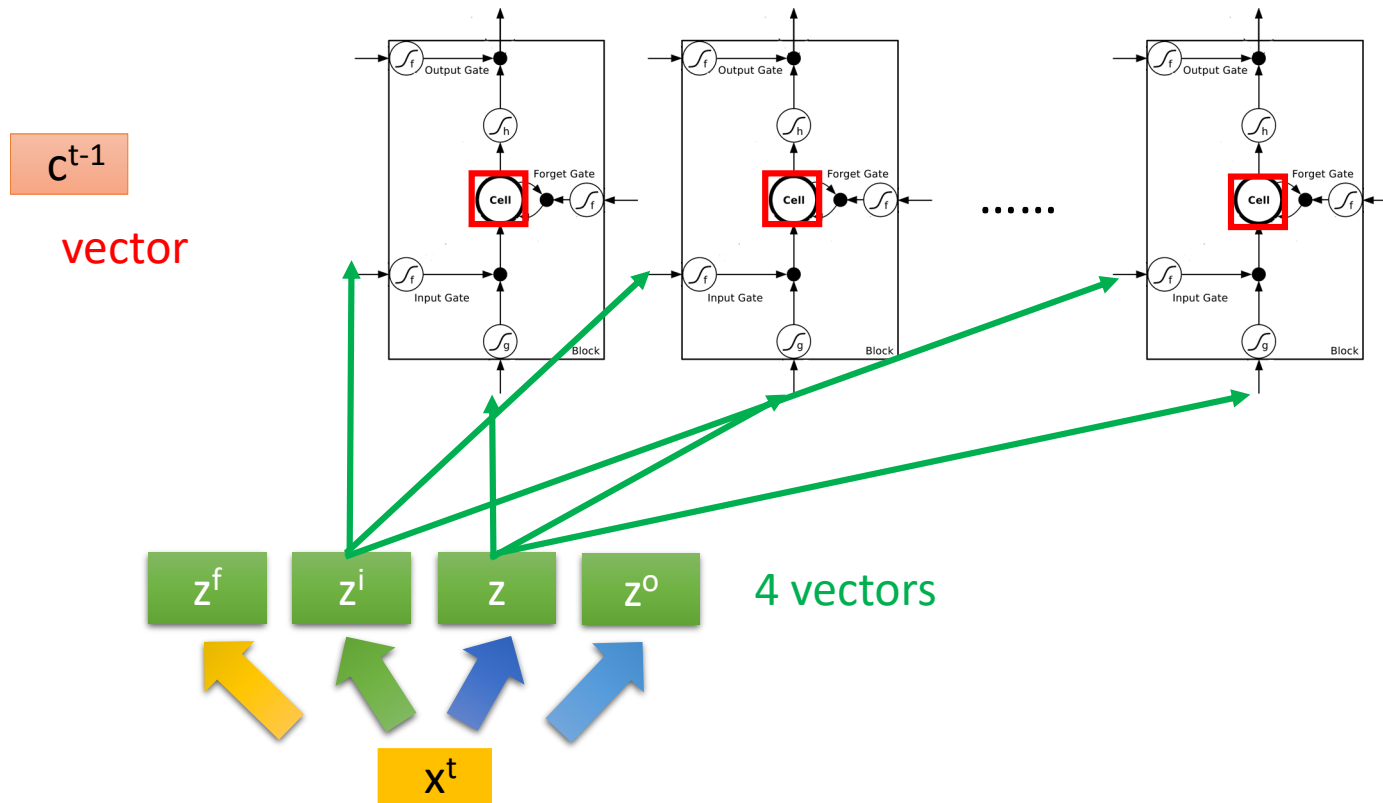
Original Network:

➤ Simply replace the neurons with LSTM

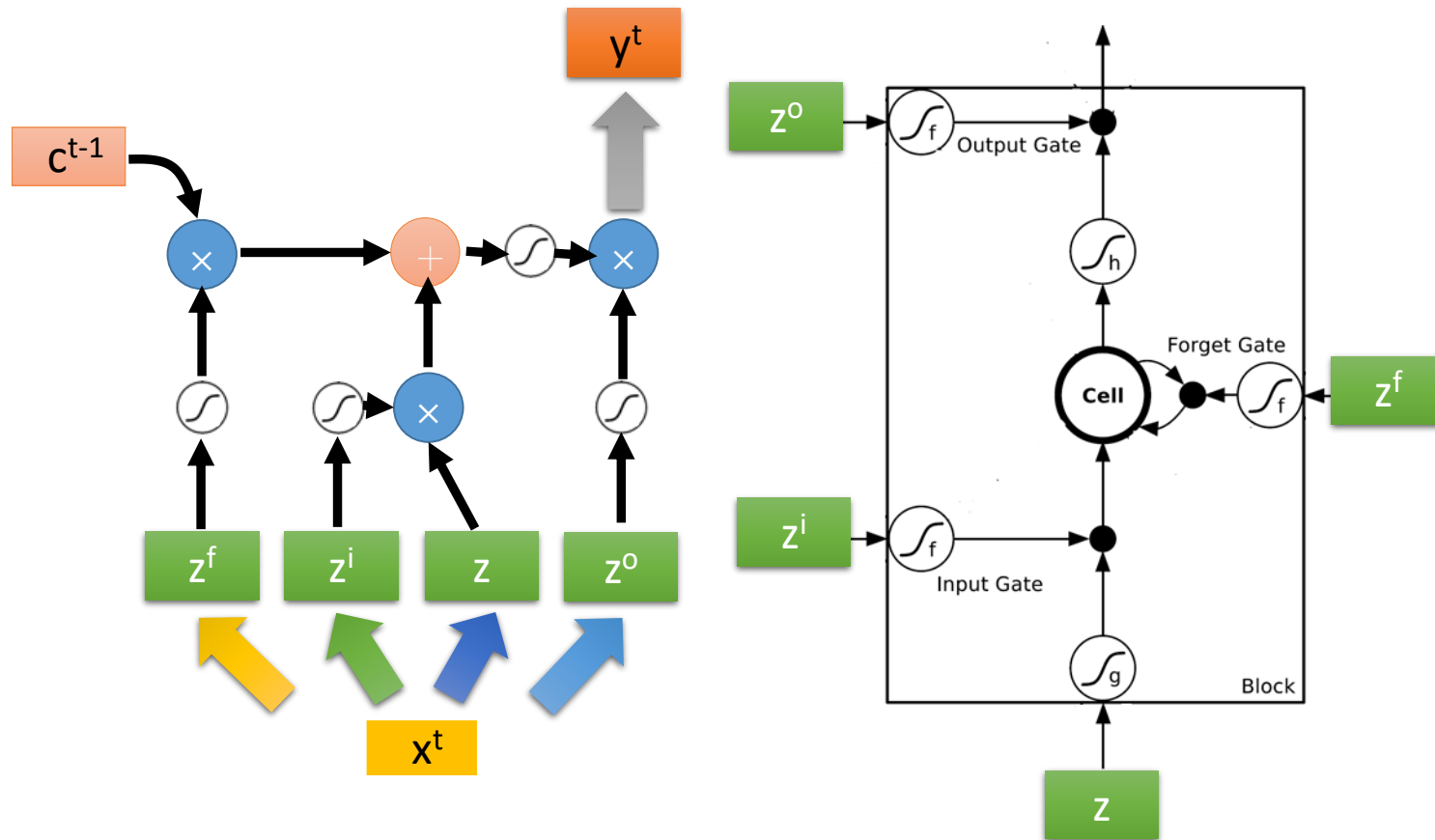




# LSTM

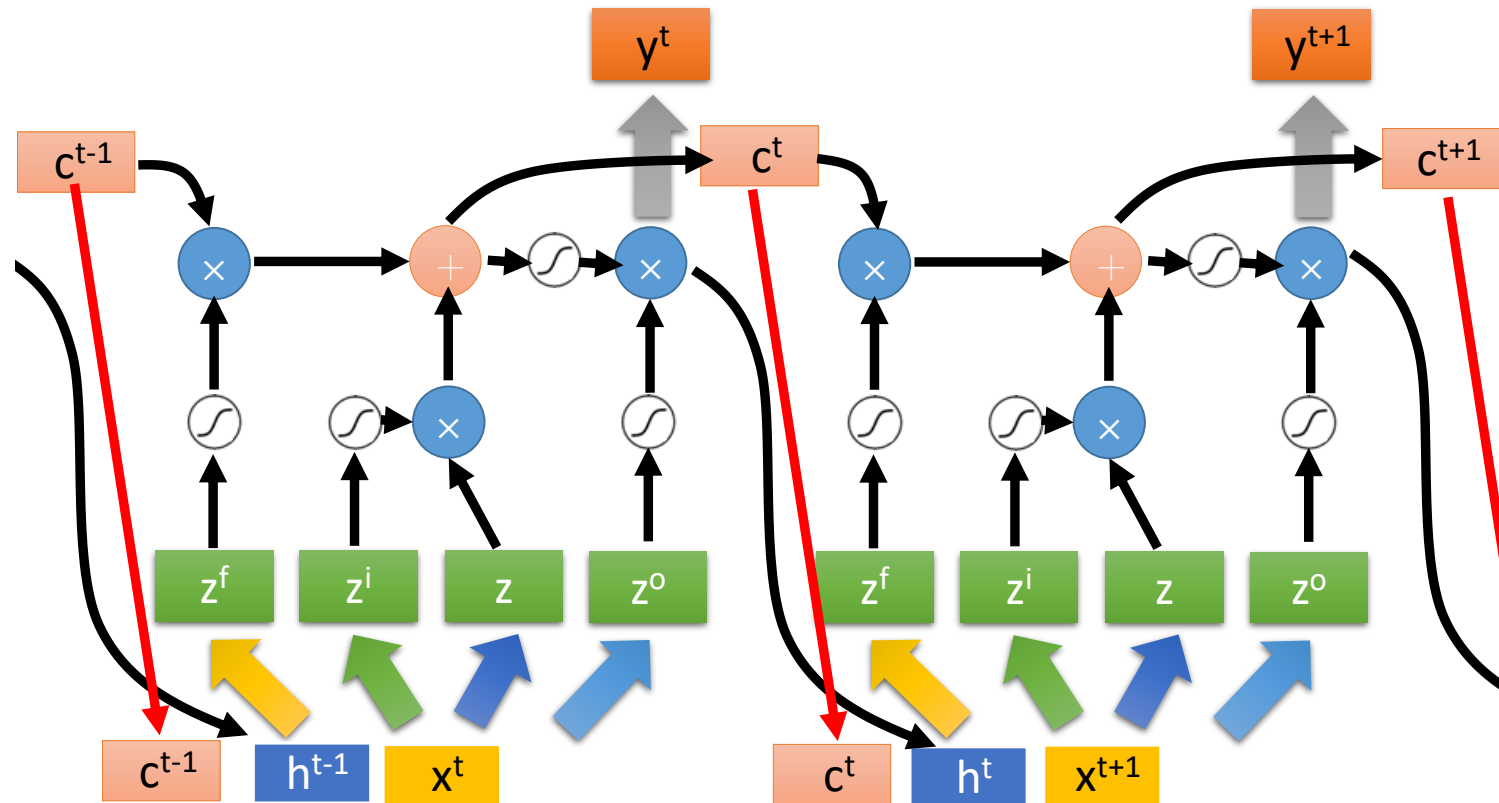


# LSTM

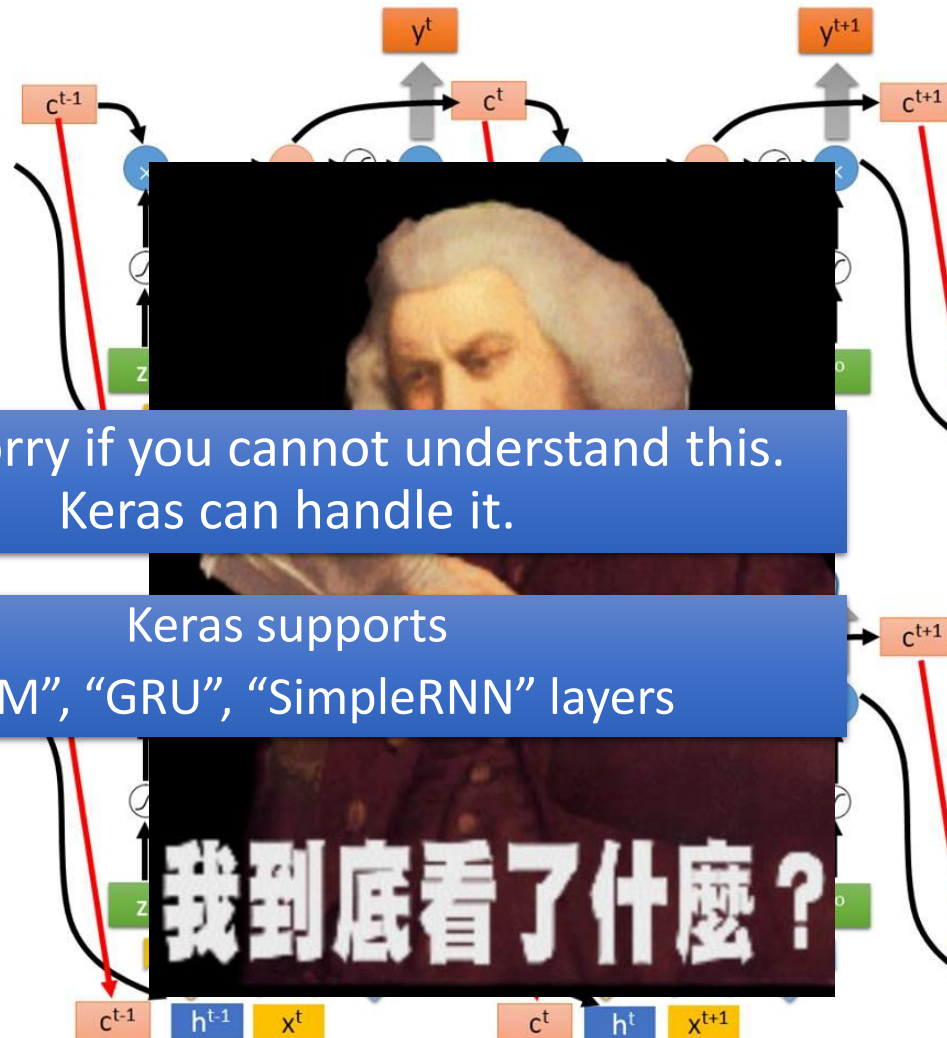


# LSTM

Extension: "peephole"



Multiple-layer  
LSTM



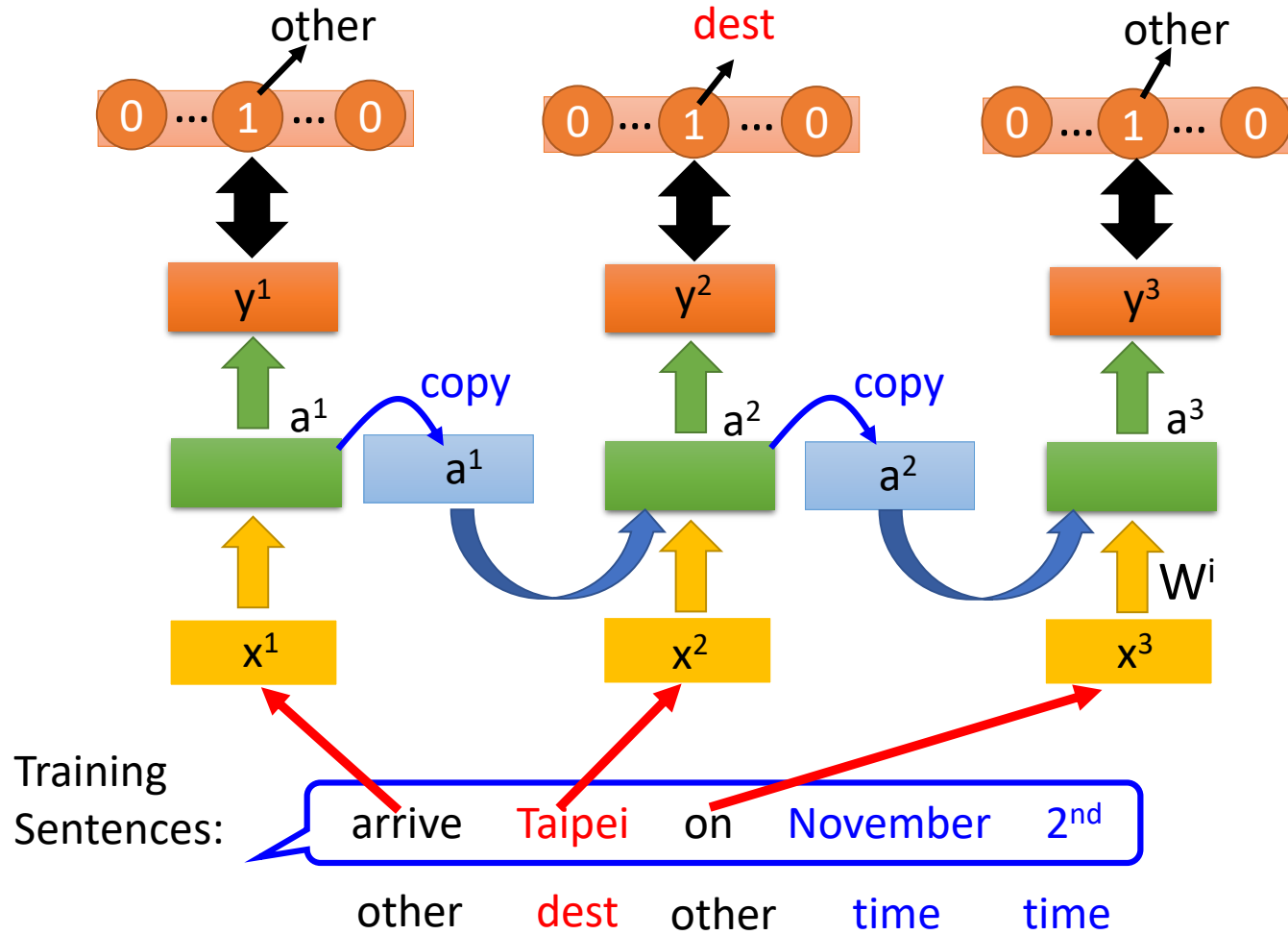
Don't worry if you cannot understand this.  
Keras can handle it.

Keras supports  
"LSTM", "GRU", "SimpleRNN" layers

This is quite  
standard now.

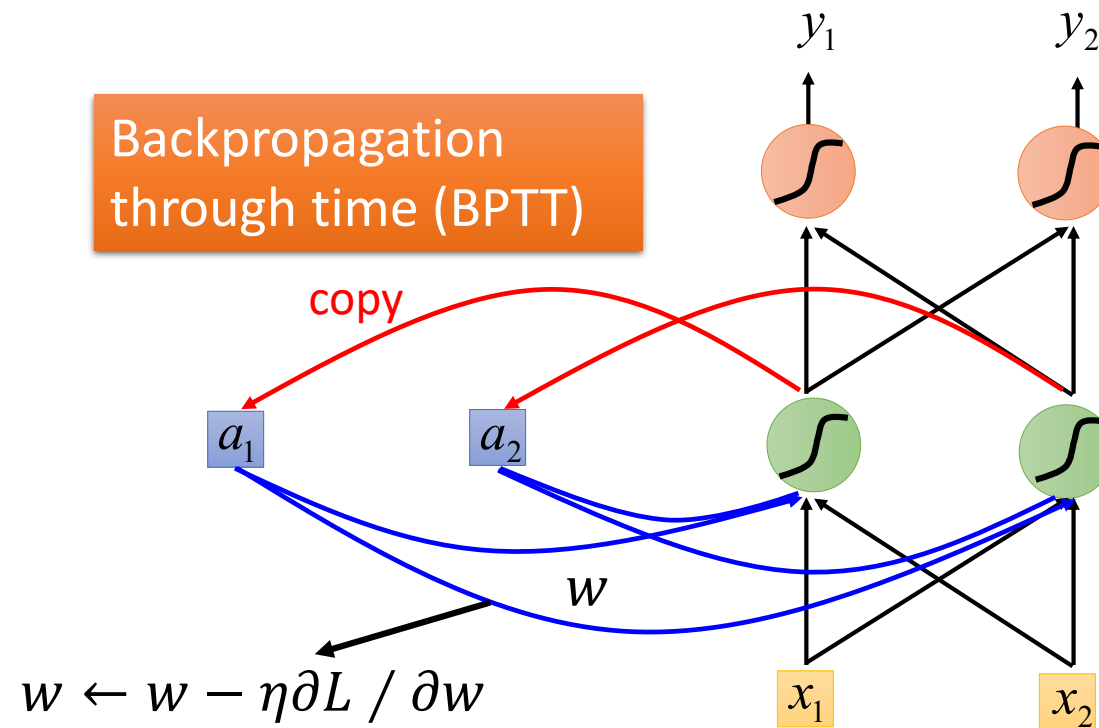
<https://img.komicolle.org/2015-09-20/src/14426967627131.gif>

## Learning Target





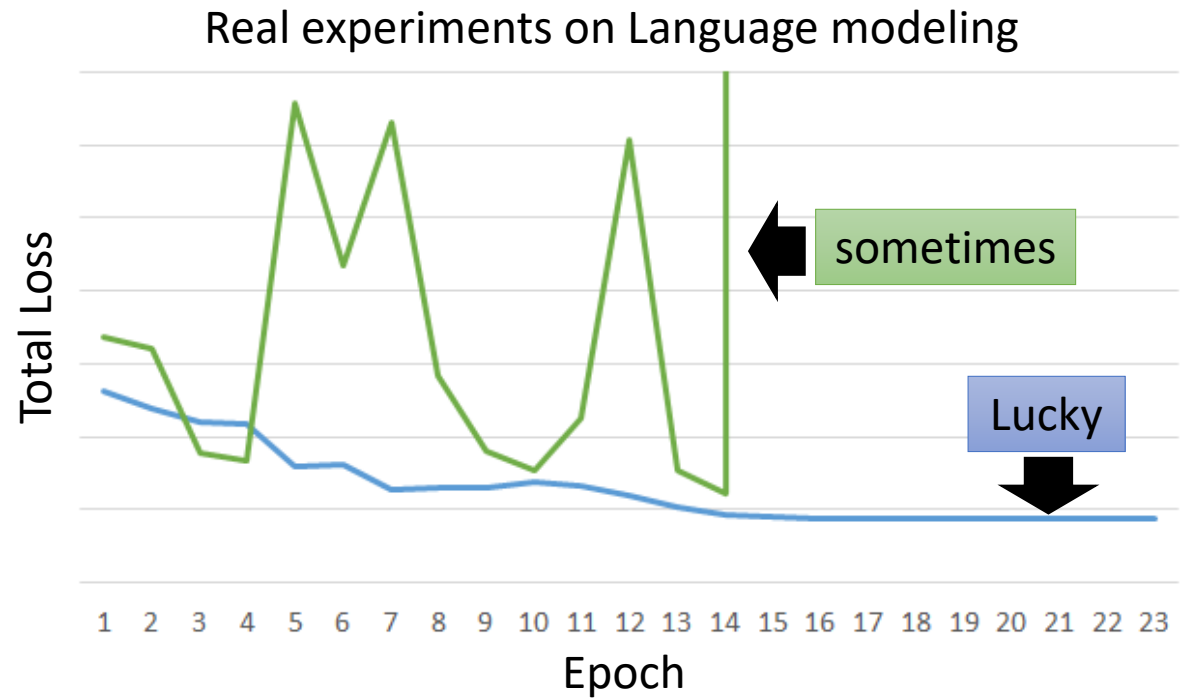
# Learning



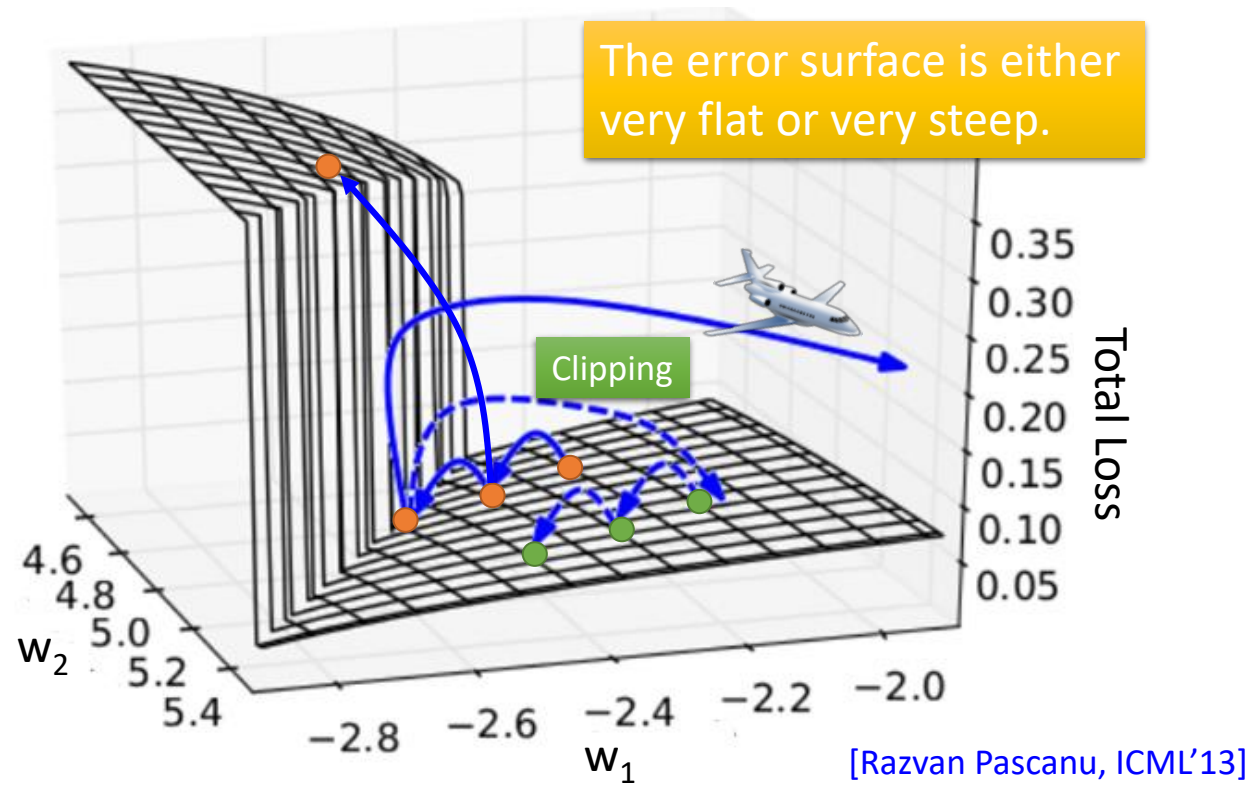
感謝 曾柏翔 同學  
提供實驗結果

## Unfortunately .....

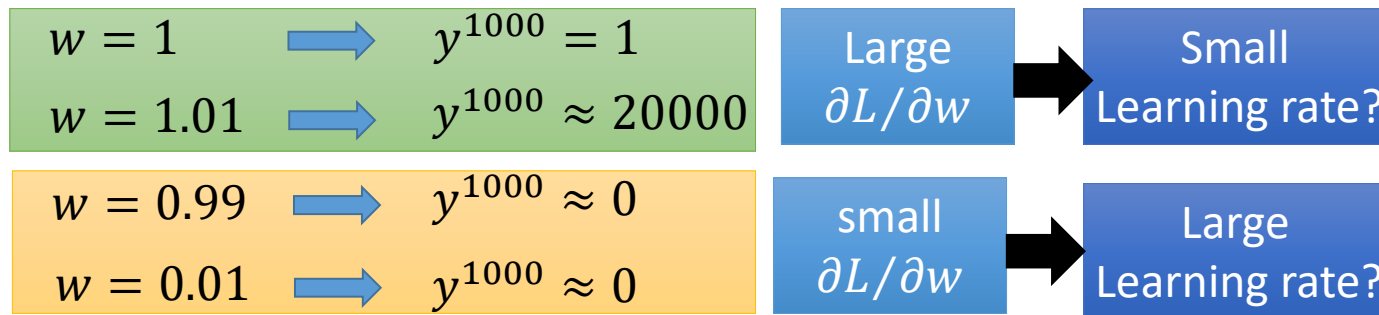
- RNN-based network is not always easy to learn



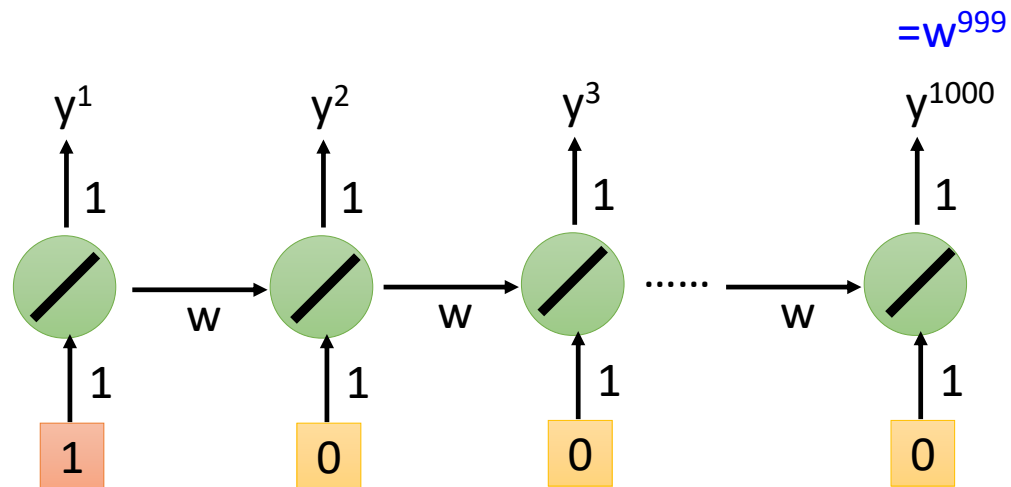
The error surface is rough.



# Why?



## Toy Example



# Helpful Techniques

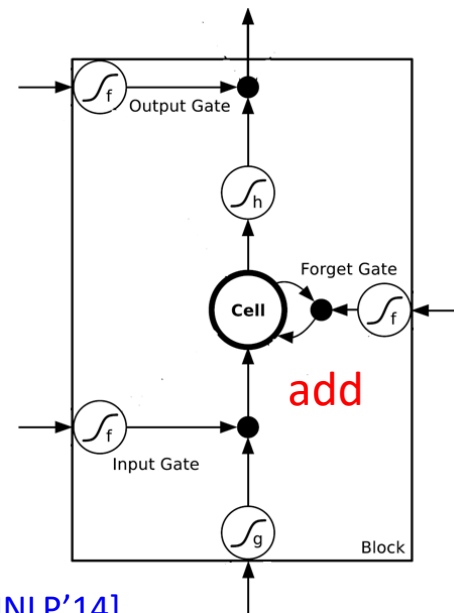
- Long Short-term Memory (LSTM)

- 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.)

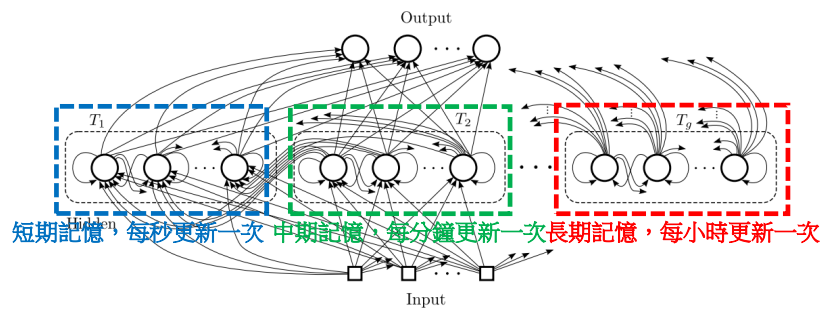
Gated Recurrent Unit (GRU):  
simpler than LSTM



[Cho, EMNLP'14]

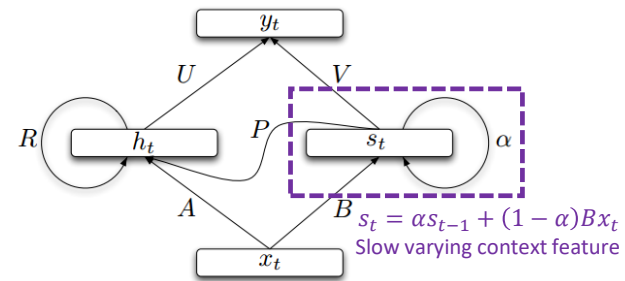
# Helpful Techniques

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]

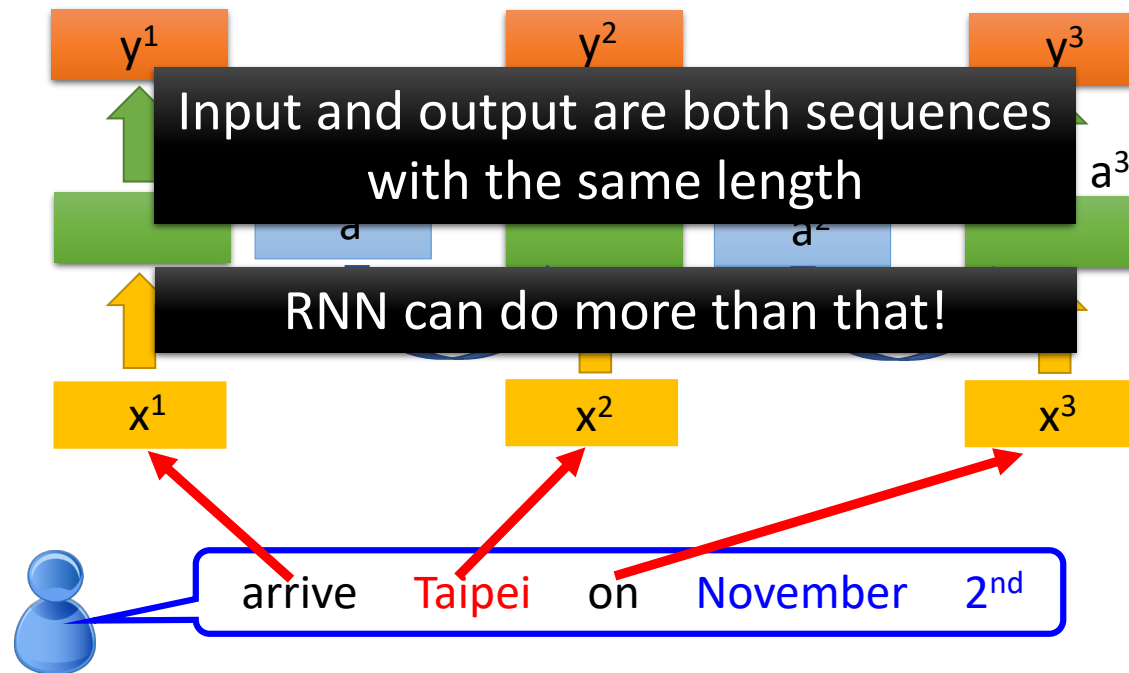
➤ Outperform or be comparable with LSTM in 4 different tasks

# More Applications .....

Probability of  
“arrive” in each slot

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

Probability of  
“on” in each slot



# Many to one

- Input is a vector sequence, but output is only one vector

## Sentiment Analysis

看了這部電影覺得很高興 .....

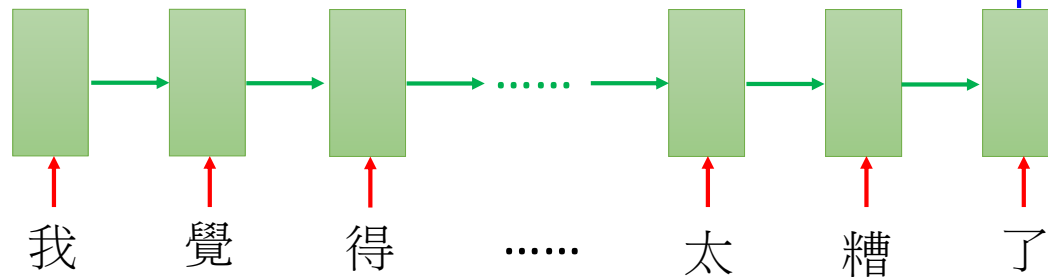
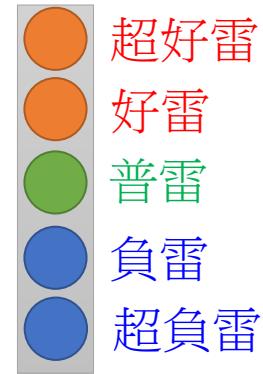
Positive (正雷)

這部電影太糟了 .....

Negative (負雷)

這部電影很棒 .....

Positive (正雷)

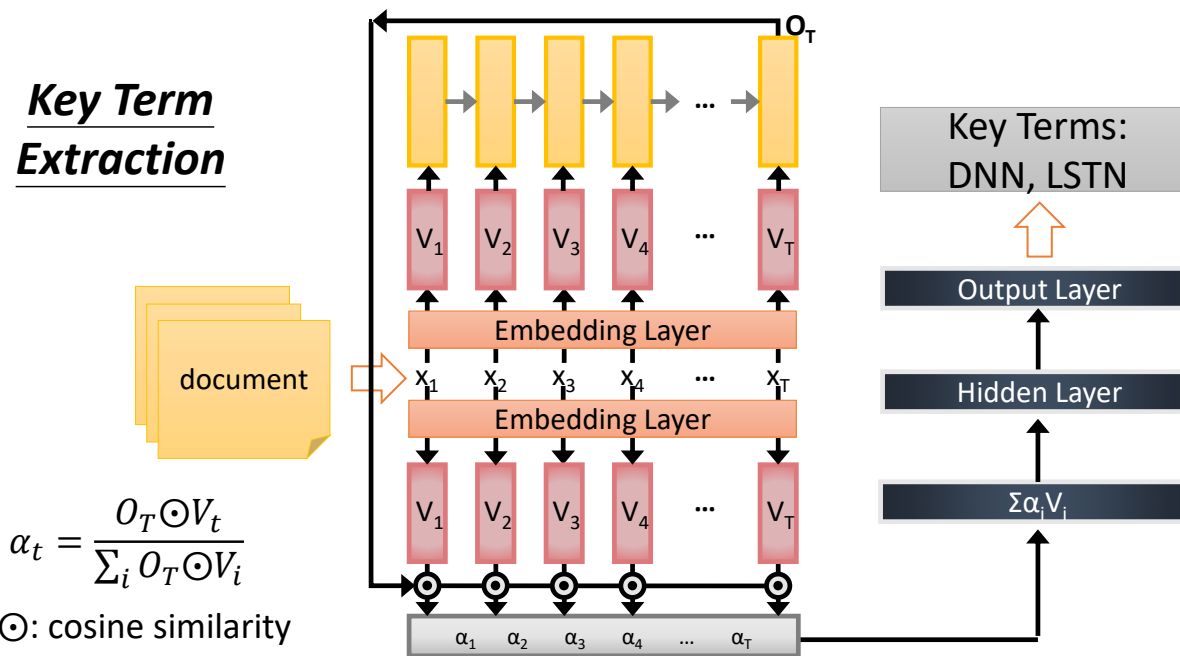




# Many to one

[Shen & Lee, *Interspeech 16*]

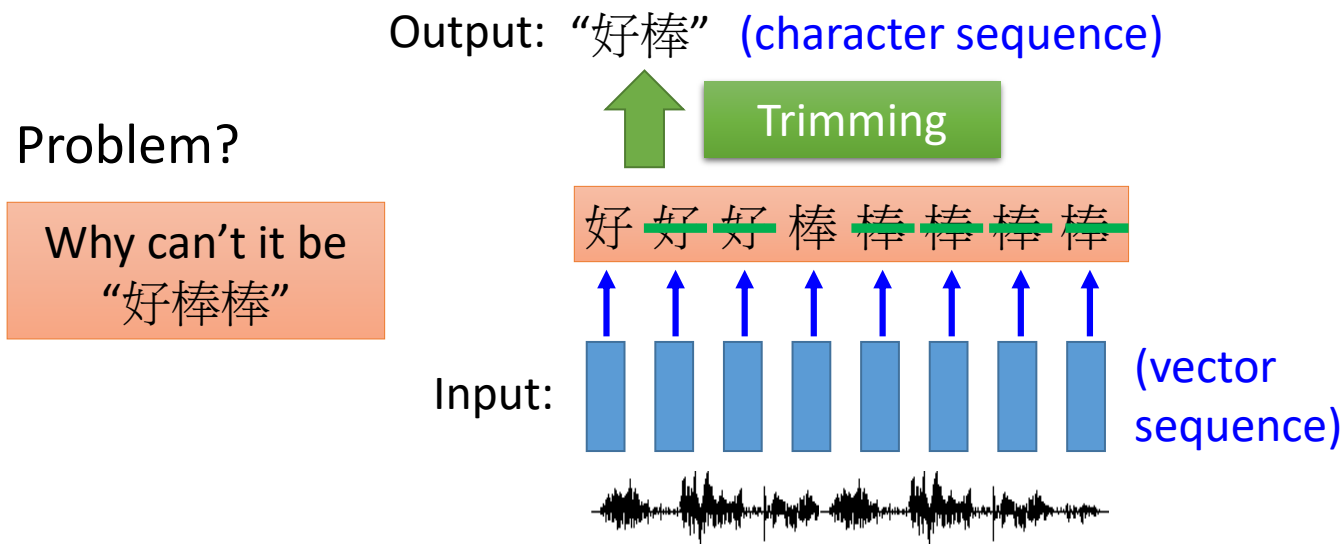
- Input is a vector sequence, but output is only one vector



Sheng-syun Shen, Hung-Yi Lee, "Neural Attention Models for Sequence Classification: Analysis and Application to Key Term Extraction and Dialogue Act Detection", the 17th Annual Conference of the International Speech Communication Association (INTERSPEECH'16), San Francisco, Sept. 2016

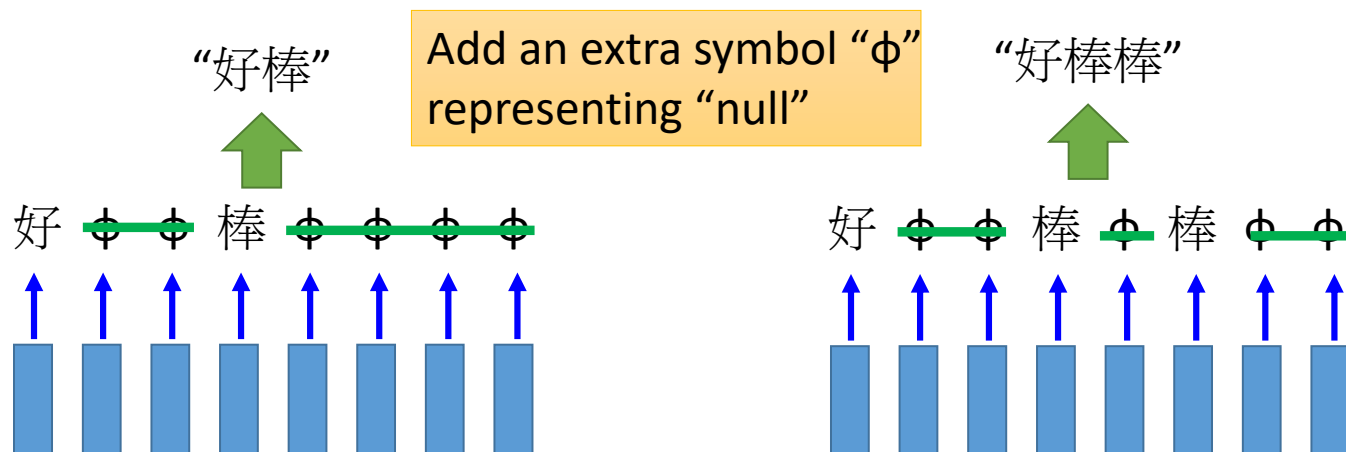
# Many to Many (Output is shorter)

- Both input and output are both sequences, but the output is shorter.
  - E.g. Speech Recognition



# Many to Many (Output is shorter)

- Both input and output are both sequences, **but the output is shorter.**
- Connectionist Temporal Classification (CTC) [Alex Graves, ICML'06][Alex Graves, ICML'14][Haşim Sak, Interspeech'15][Jie Li, Interspeech'15][Andrew Senior, ASRU'15]




# Many to Many (Output is shorter)


- CTC: Training


Acoustic Features: 

Label: 好 棒

All possible alignments are considered as correct.

  
好  $\phi$  棒  $\phi$   $\phi$   $\phi$

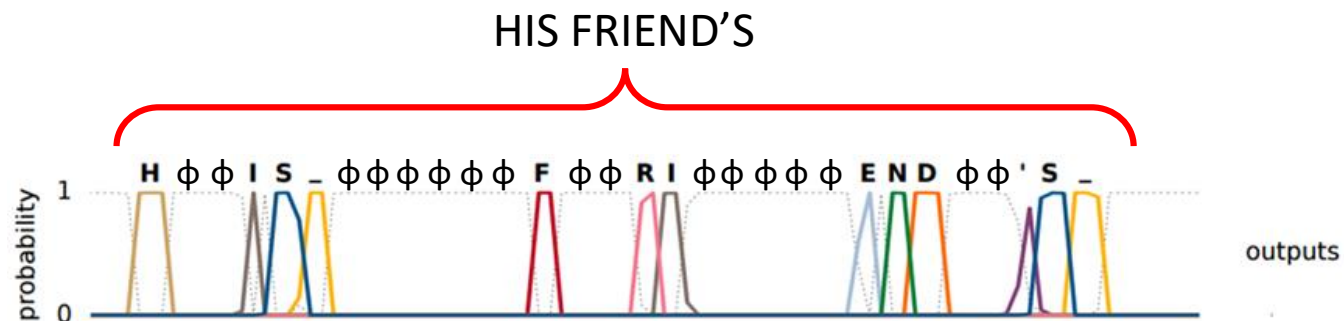
  
好  $\phi$   $\phi$  棒  $\phi$   $\phi$

  
好  $\phi$   $\phi$   $\phi$  棒  $\phi$

⋮

# Many to Many (Output is shorter)

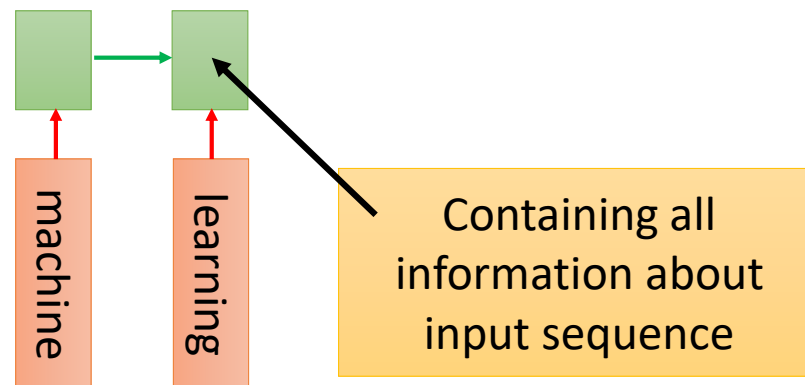
- CTC: example



Graves, Alex, and Navdeep Jaitly. "Towards end-to-end speech recognition with recurrent neural networks." *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*. 2014.

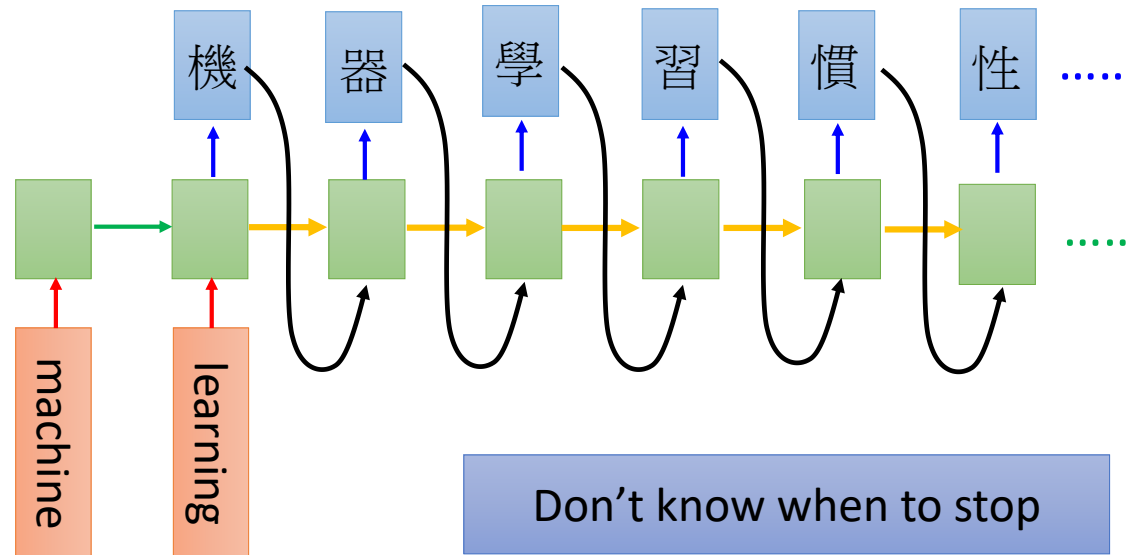
# Many to Many (No Limitation)

- Both input and output are both sequences *with different lengths.* → *Sequence to sequence learning*
  - E.g. *Machine Translation* (machine learning → 機器學習)



# Many to Many (No Limitation)

- Both input and output are both sequences with different lengths. → Sequence to sequence learning
  - E.g. Machine Translation (machine learning → 機器學習)



# Many to Many (No Limitation)

```
推  :      超      06/12 10:39
推  n:      人      06/12 10:40
推  tion:    正      06/12 10:41
→  host:    大      06/12 10:47
推  :      中      06/12 10:59
推  403:    天      06/12 11:11
推  :      外      06/12 11:13
推  527:    飛      06/12 11:17
→  990b:   仙      06/12 11:32
→  512:    草      06/12 12:15

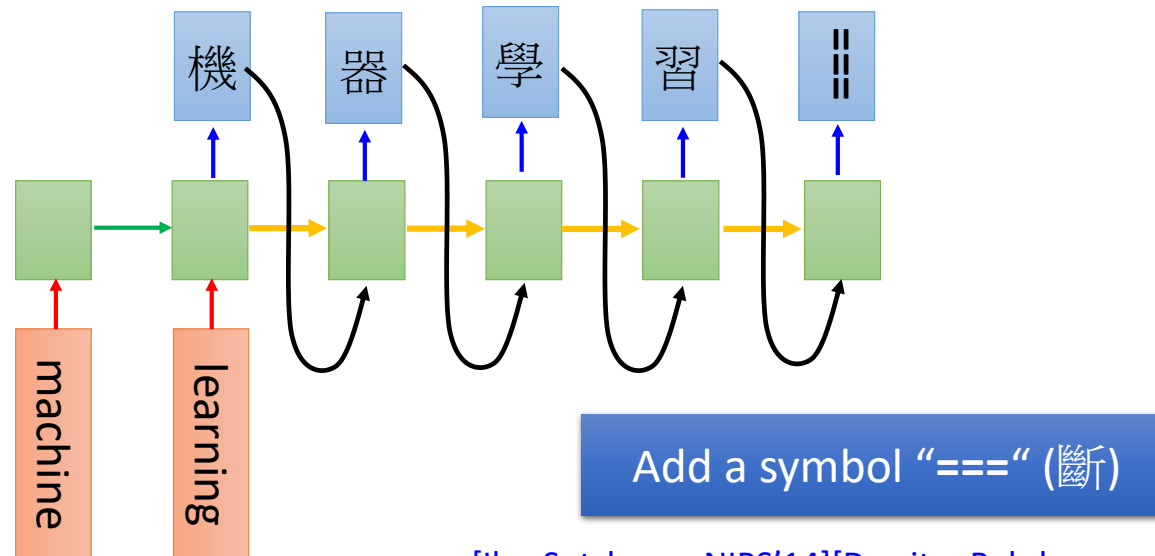
推 tlkagk:  =====斷=====
```

接龍推文是ptt在推文中的一種趣味玩法，與推齊有些類似但又有所不同，是指在推文中接續上一樓的字句，而推出連續的意思。該類玩法確切起源已不可知(鄉民百科)



# Many to Many (No Limitation)

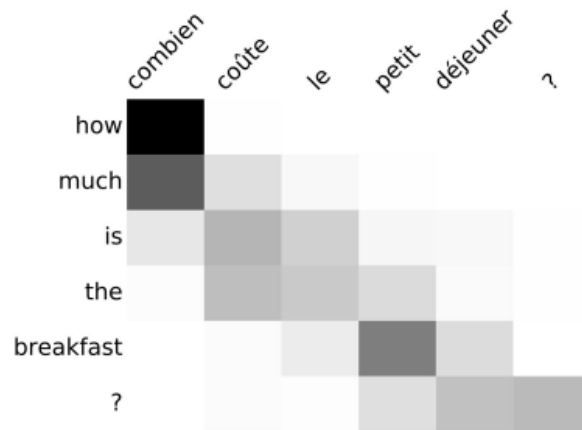
- Both input and output are both sequences with different lengths. → Sequence to sequence learning
  - E.g. Machine Translation (machine learning → 機器學習)



[Ilya Sutskever, NIPS'14][Dzmitry Bahdanau, arXiv'15]

# Many to Many (No Limitation)

- Both input and output are both sequences ***with different lengths***. → ***Sequence to sequence learning***
  - E.g. ***Machine Translation*** (machine learning → 機器學習)



(a) Machine translation alignment

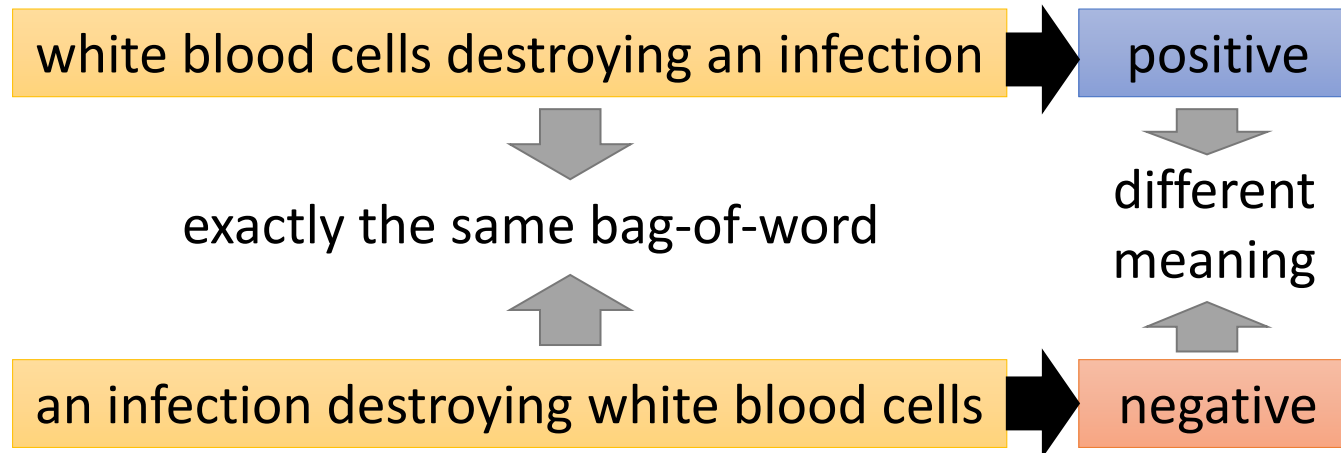


(b) Speech translation alignment

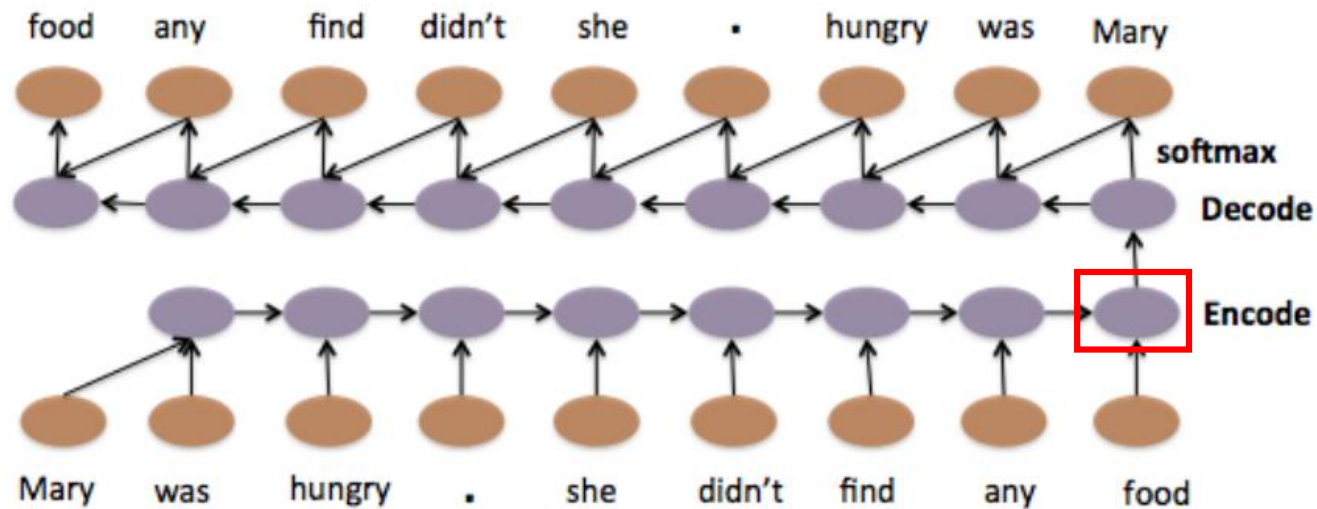
Figure 1: Alignments performed by the attention model during training

# Sequence-to-sequence Auto-encoder - Text

- To understand the meaning of a word sequence, the order of the words can not be ignored.

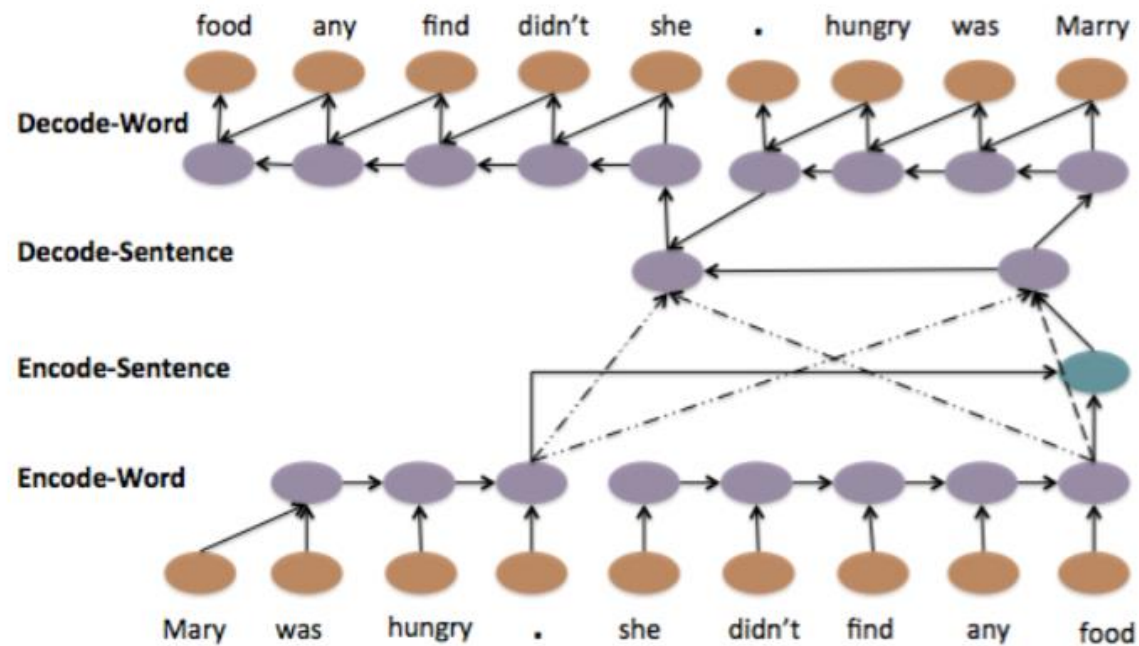


# Sequence-to-sequence Auto-encoder - Text



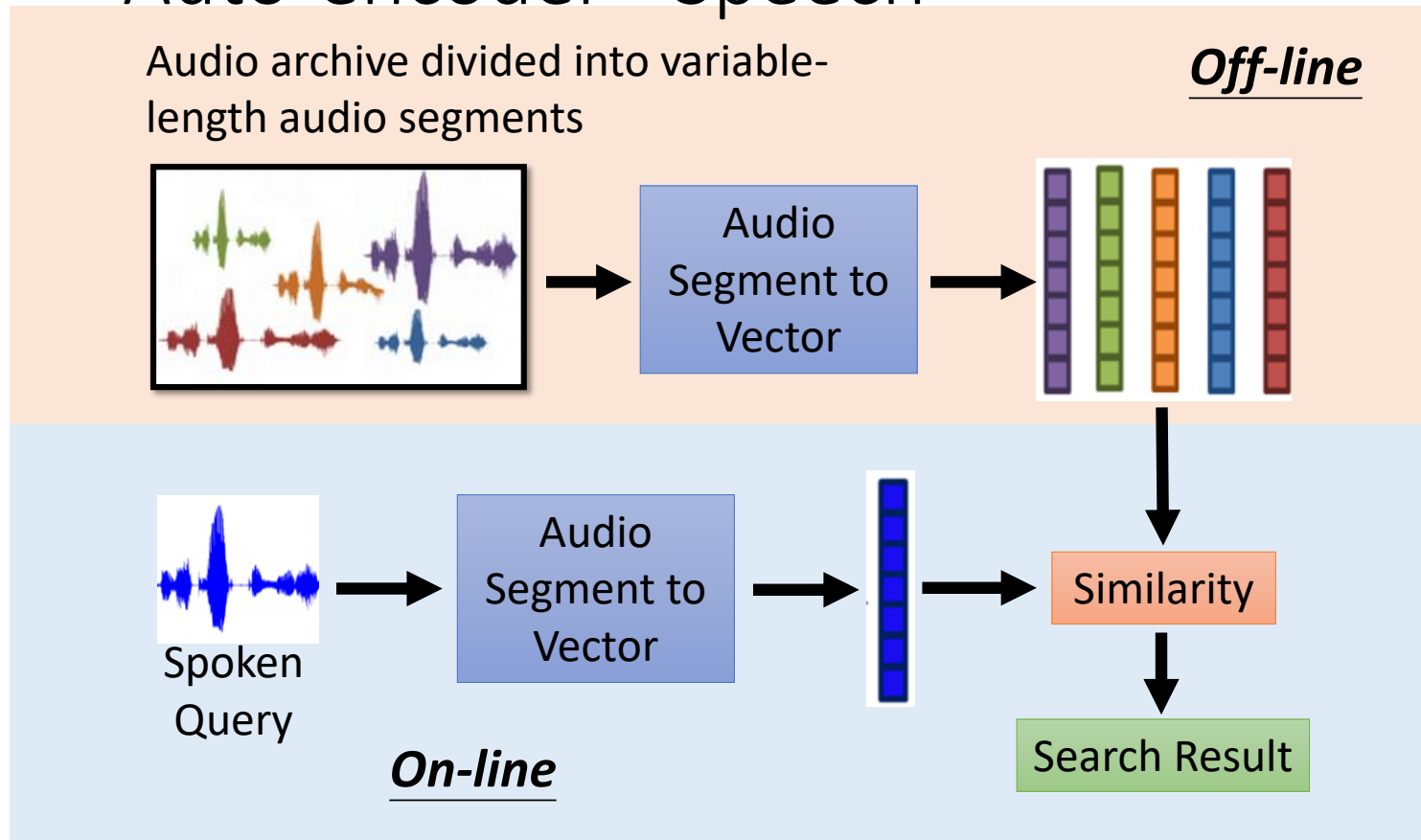
Li, Jiwei, Minh-Thang Luong, and Dan Jurafsky. "A hierarchical neural autoencoder for paragraphs and documents." *arXiv preprint arXiv:1506.01057*(2015).

# Sequence-to-sequence Auto-encoder - Text

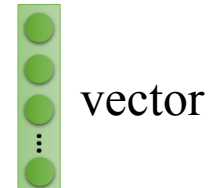
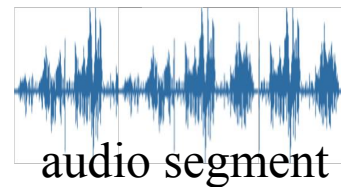


Li, Jiwei, Minh-Thang Luong, and Dan Jurafsky. "A hierarchical neural autoencoder for paragraphs and documents." *arXiv preprint arXiv:1506.01057*(2015).

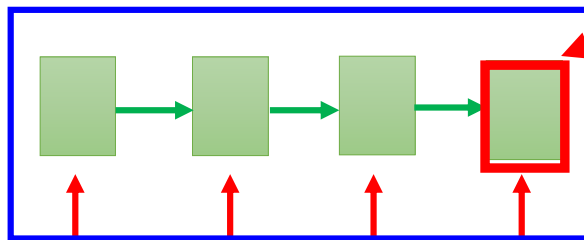
# Sequence-to-sequence Auto-encoder - Speech



# Sequence-to-sequence Auto-encoder - Speech



## RNN Encoder



The values in the memory represent the whole audio segment

The vector we want

How to train RNN Encoder?

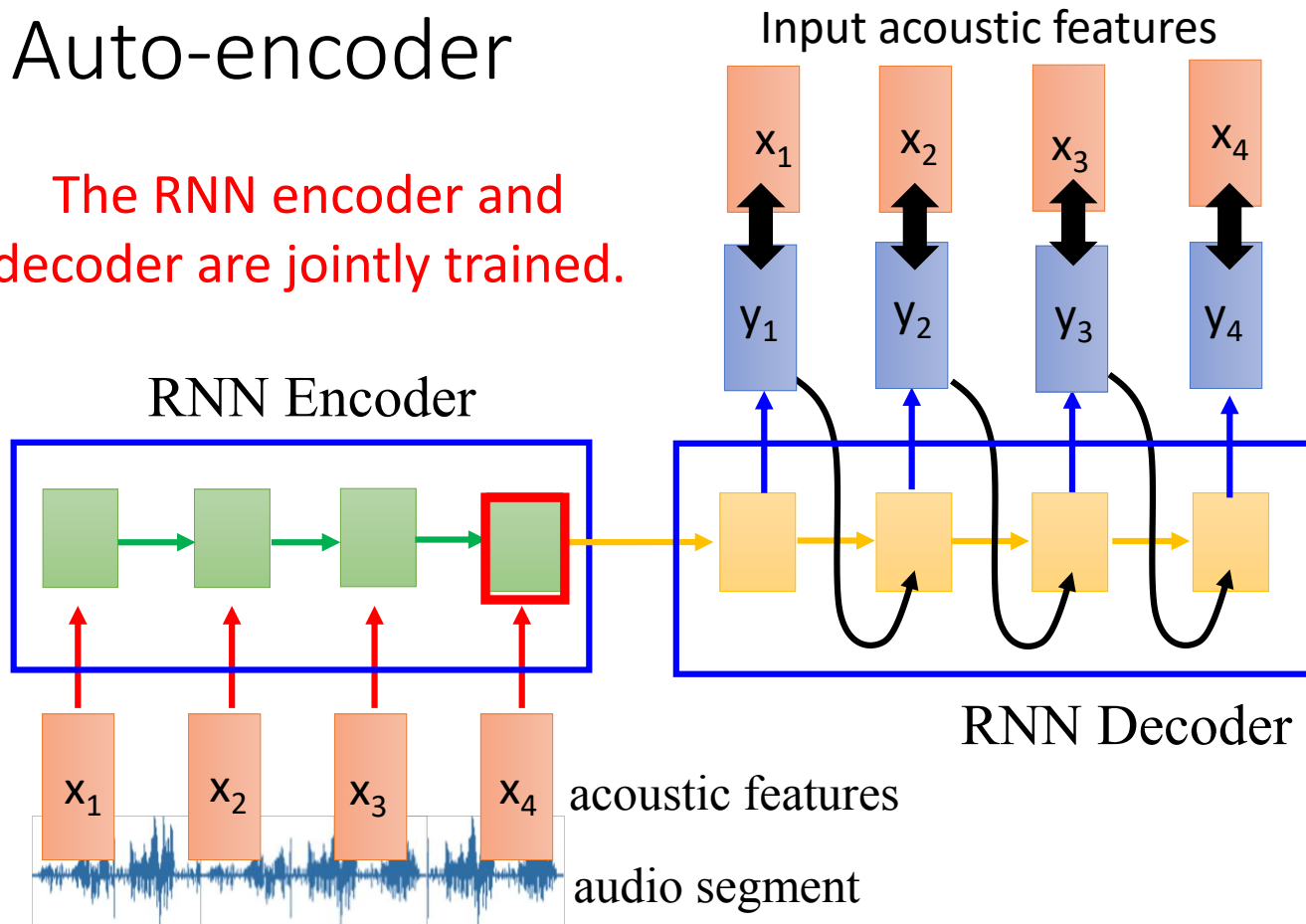


acoustic features

audio segment

# Sequence-to-sequence Auto-encoder

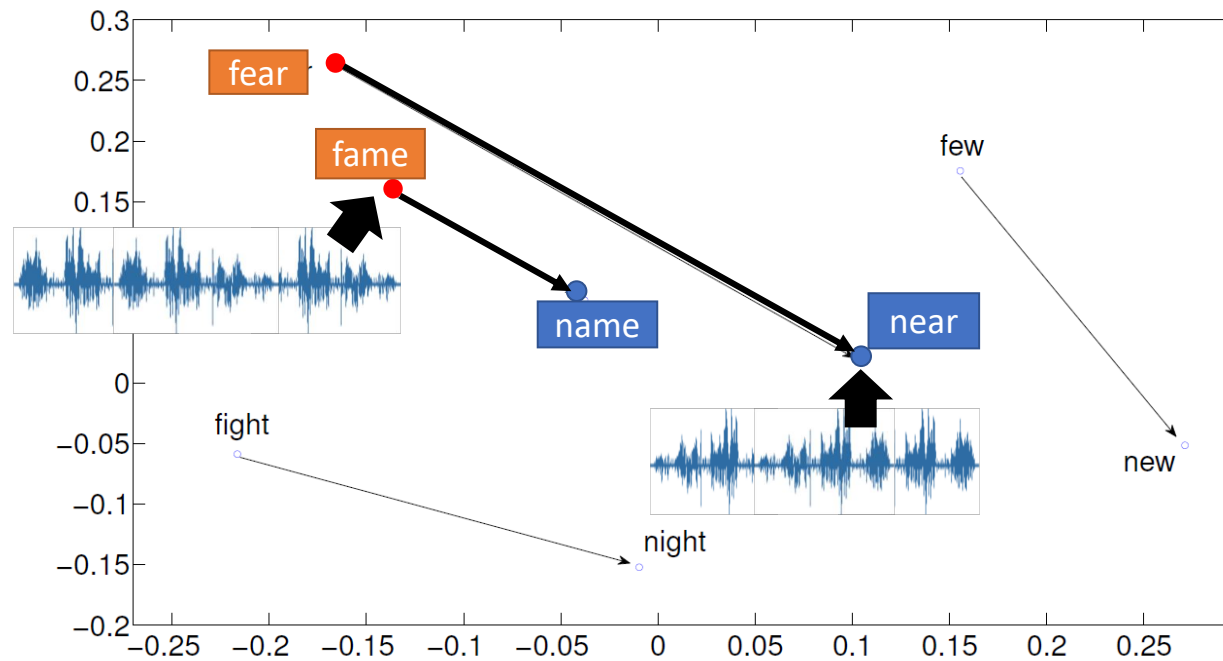
The RNN encoder and decoder are jointly trained.



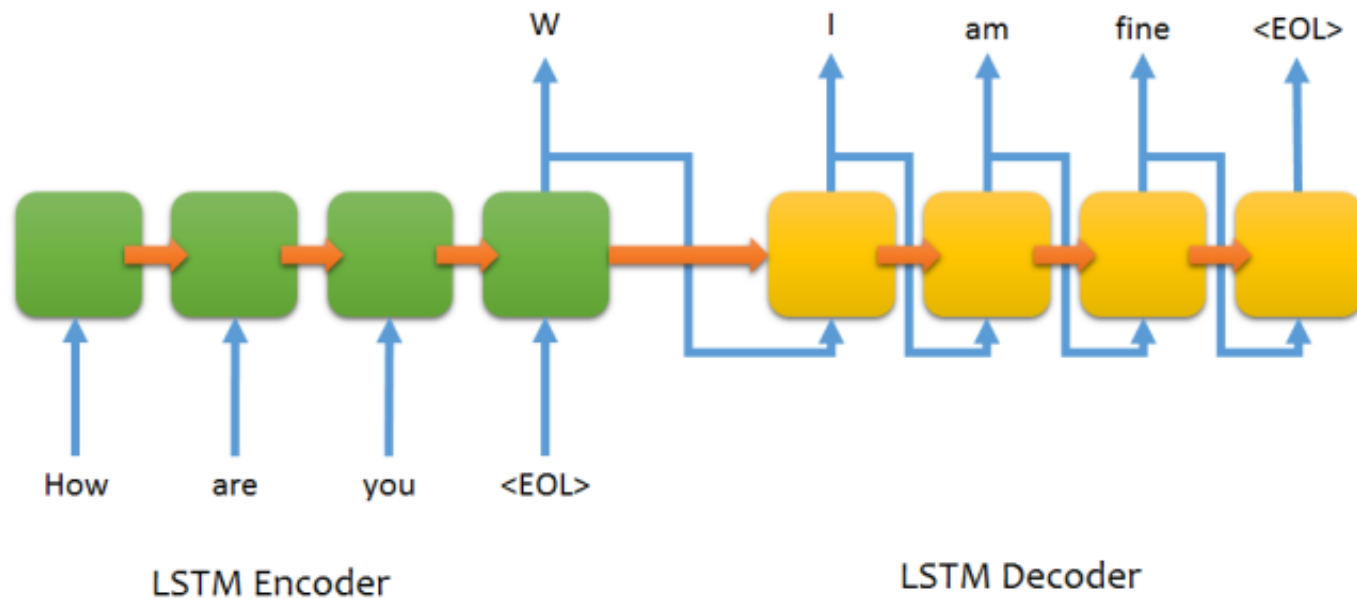


# Sequence-to-sequence Auto-encoder - Speech

- Visualizing embedding vectors of the words

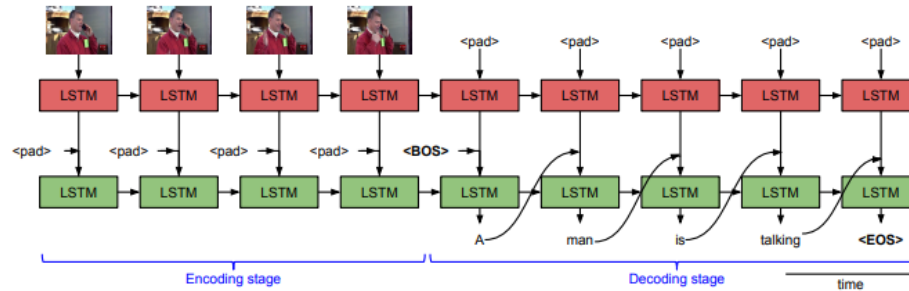


# Demo: Chat-bot



電視影集 (~40,000 sentences)、美國總統大選辯論

# Video Caption Generation



## Correct descriptions.



S2VT: A man is doing stunts on his bike.



S2VT: A herd of zebras are walking in a field.



S2VT: A young woman is doing her hair.



S2VT: A man is shooting a gun at a target.

## Relevant but incorrect descriptions.



S2VT: A small bus is running into a building.



S2VT: A man is cutting a piece of a pair of a paper.



S2VT: A cat is trying to get a small board.



S2VT: A man is spreading butter on a tortilla.

## Irrelevant descriptions.



S2VT: A man is pouring liquid in a pan.



S2VT: A polar bear is walking on a hill.



S2VT: A man is doing a pencil.



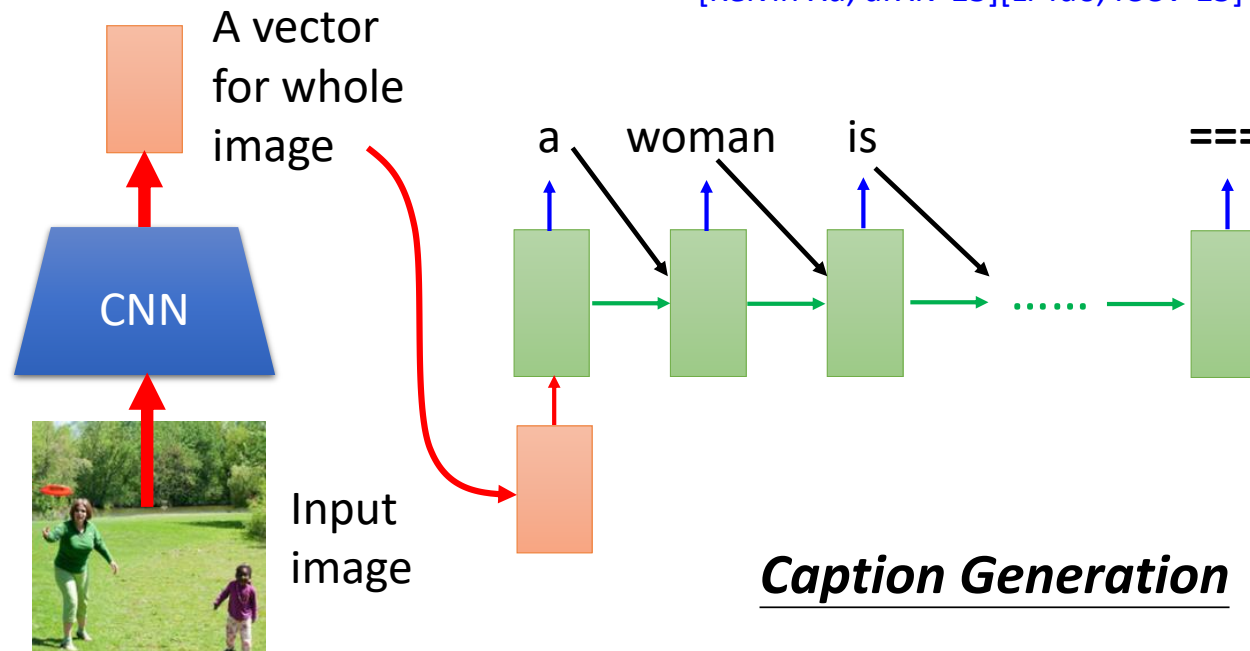
S2VT: A black clip to walking through a path.

Subhashini Venugopalan, Marcus Rohrbach, Jeffrey Donahue, Raymond Mooney, Trevor Darrell, and Kate Saenko. 2015. Sequence to Sequence -- Video to Text. In Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV) (ICCV '15). IEEE Computer Society, Washington, DC, USA, 4534-4542.

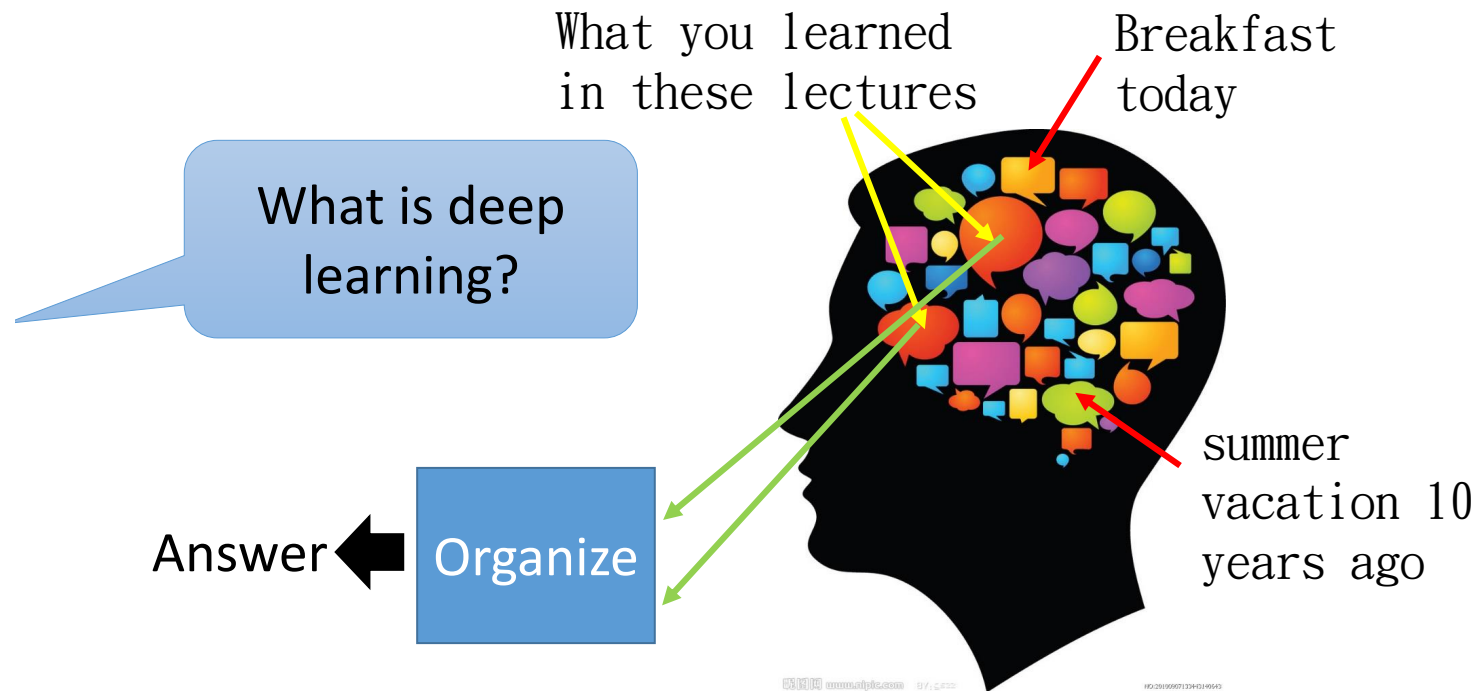
# Demo: Image Caption Generation

- Input an image, but output a sequence of words

[Kelvin Xu, arXiv'15][Li Yao, ICCV'15]

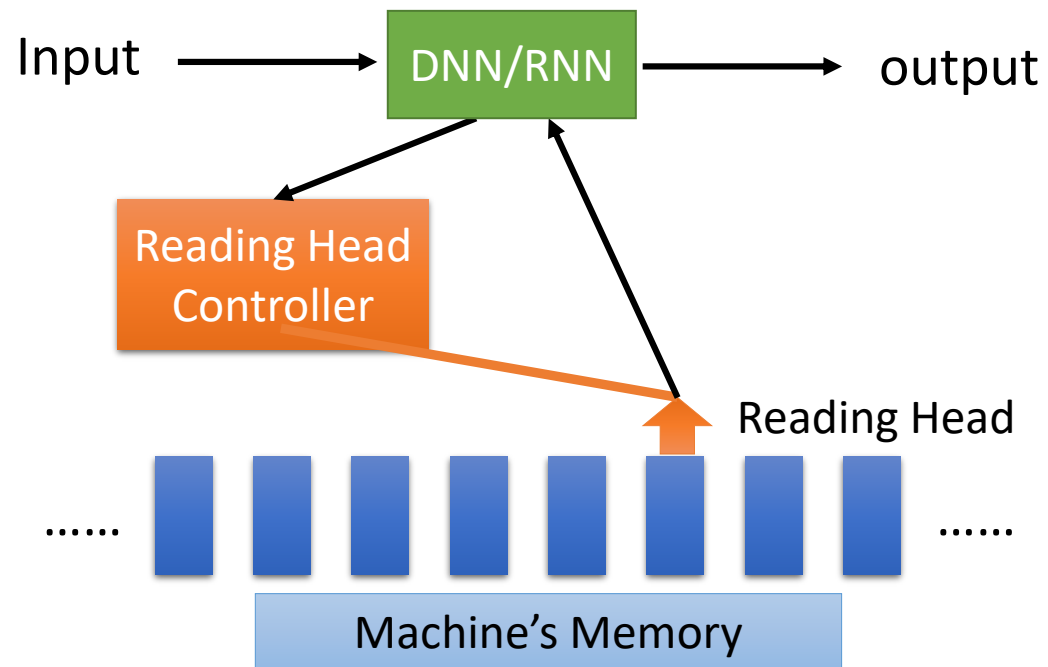


# Attention-based Model



[http://henrylo1605.blogspot.tw/2015/05/blog-post\\_56.html](http://henrylo1605.blogspot.tw/2015/05/blog-post_56.html)

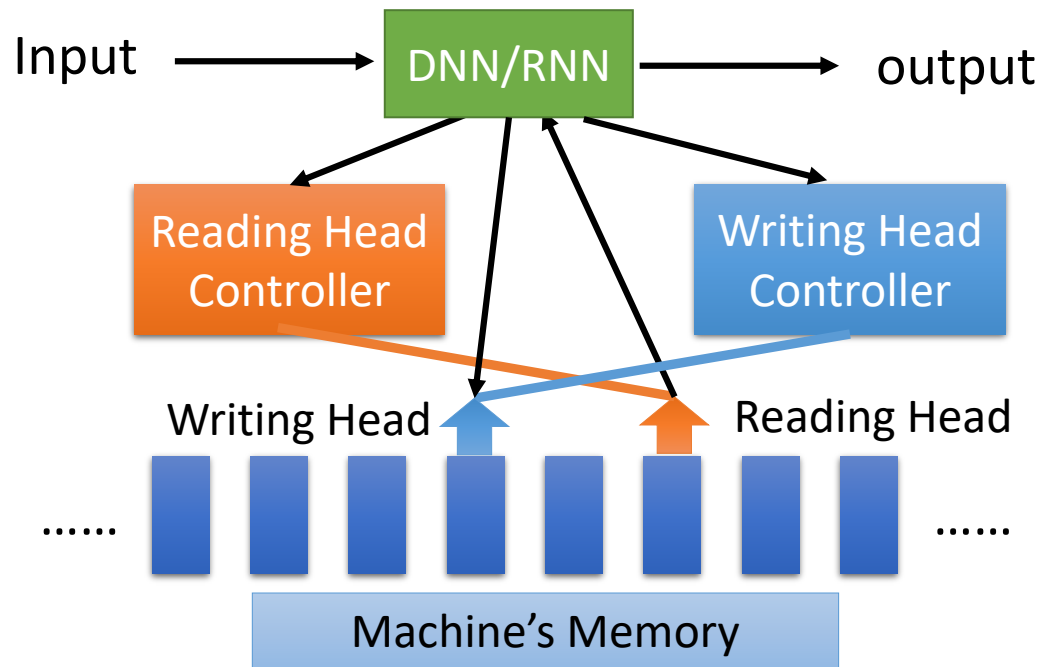
# Attention-based Model



Ref:

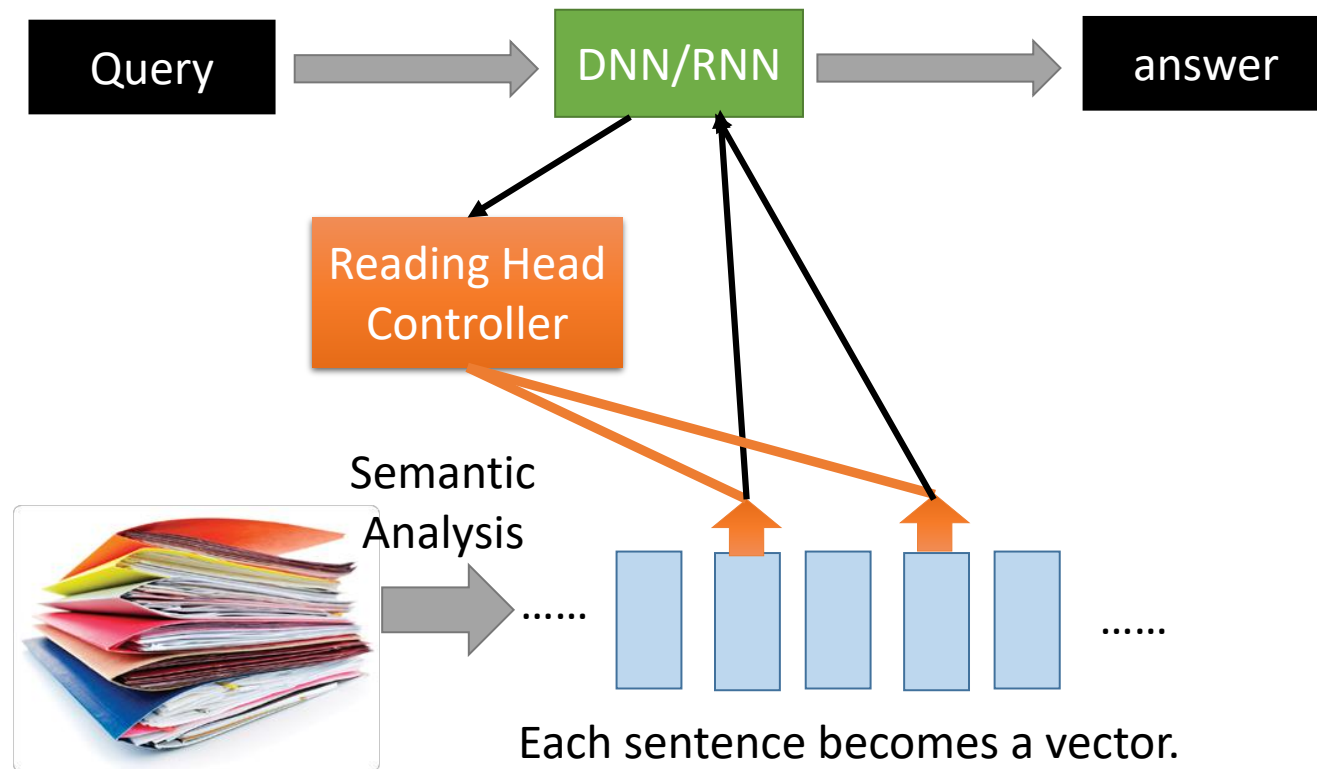
[http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS\\_2015\\_2/Lecture/Attain%20\(v3\).e cm.mp4/index.html](http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS_2015_2/Lecture/Attain%20(v3).e cm.mp4/index.html)

# Attention-based Model v2



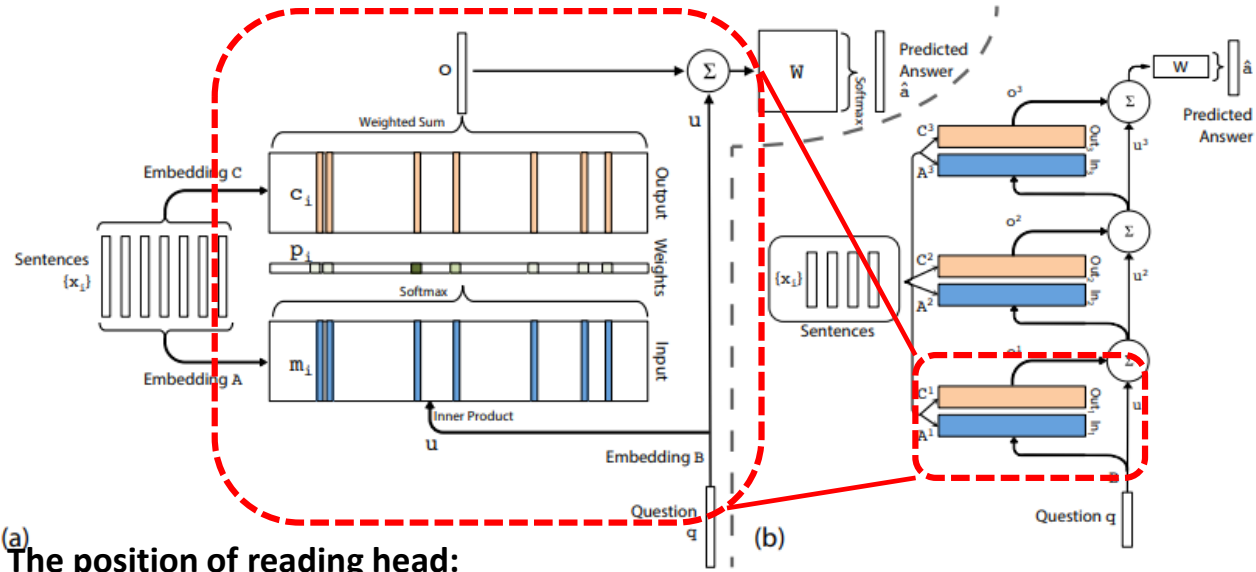
Neural Turing Machine

# Reading Comprehension





# Reading Comprehension

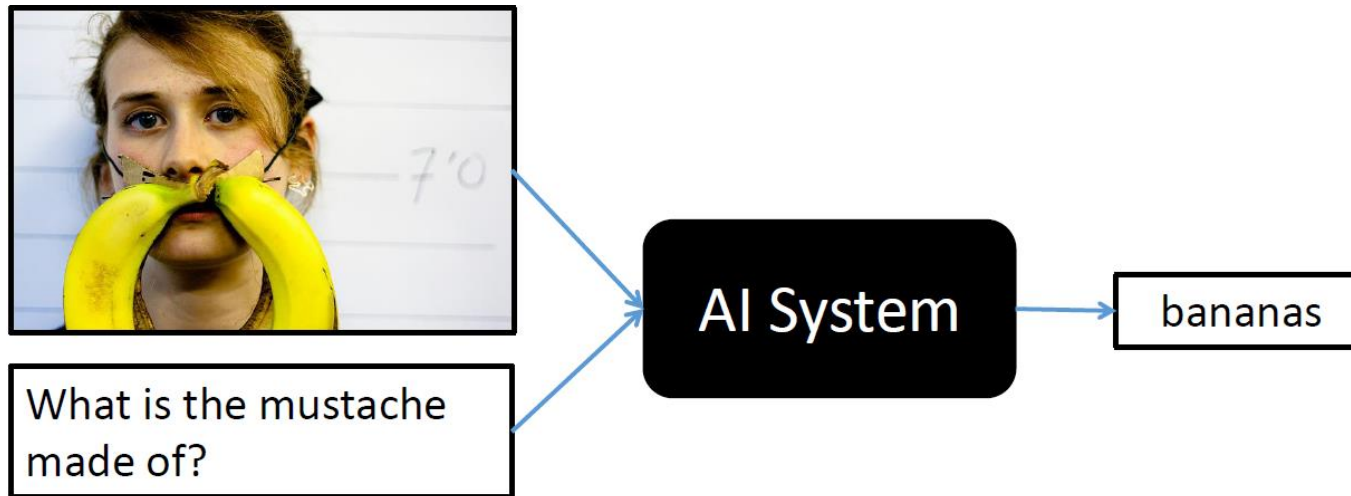


Story (16: basic induction)	Support	Hop 1	Hop 2	Hop 3
Brian is a frog.	yes	0.00	0.98	0.00
Lily is gray.		0.07	0.00	0.00
Brian is yellow.	yes	0.07	0.00	1.00
Julius is green.		0.06	0.00	0.00
Greg is a frog.	yes	0.76	0.02	0.00
<b>What color is Greg? Answer: yellow Prediction: yellow</b>				

End-To-End Memory Networks. S. Sukhbaatar, A. Szlam, J. Weston, R. Fergus. NIPS, 2015.

**Keras example:** [https://github.com/fchollet/keras/blob/master/examples/babi\\_memnn.py](https://github.com/fchollet/keras/blob/master/examples/babi_memnn.py)

# Visual Question Answering



source: <http://visualqa.org/>

# Visual Question Answering

