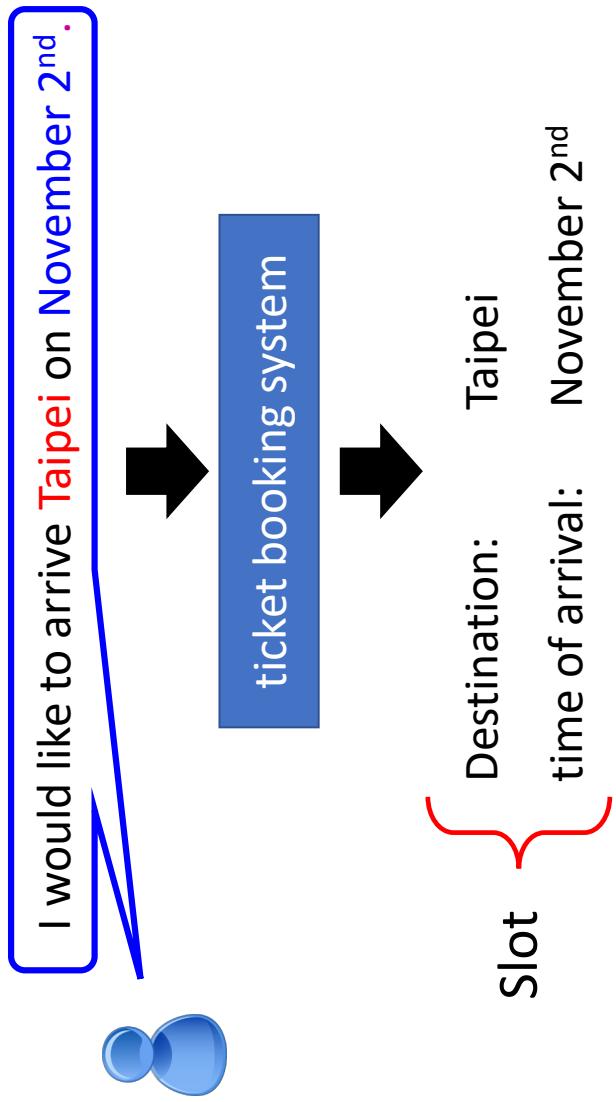


Recurrent Neural Network (RNN)

Example Application

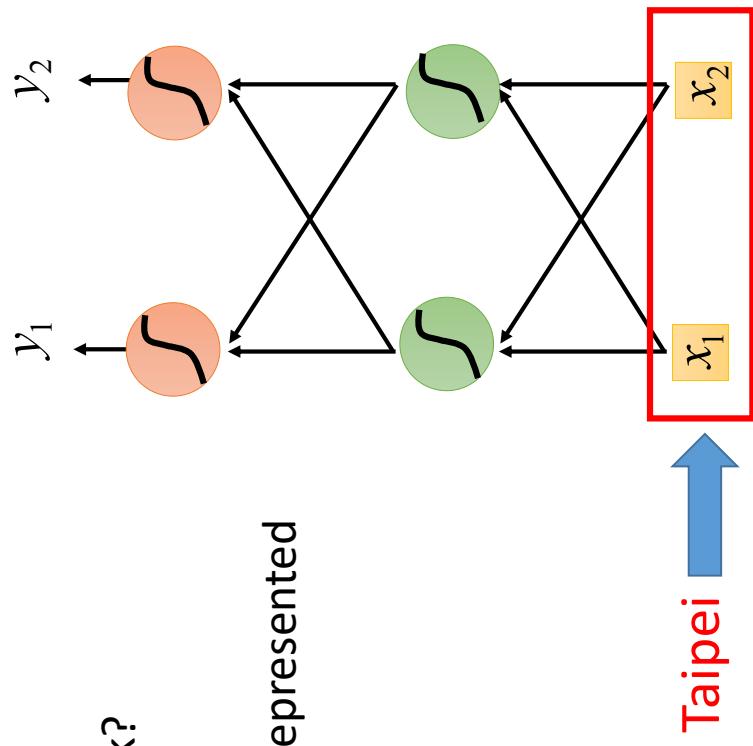
- Slot Filling



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}

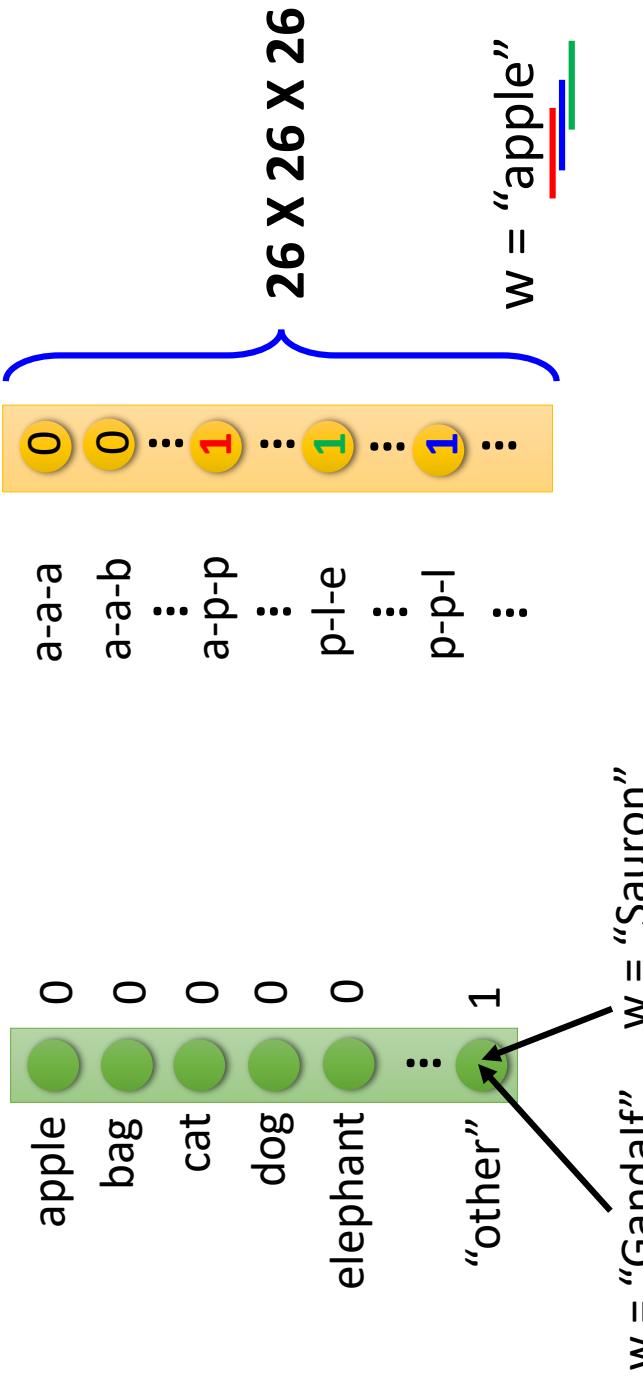
The vector is lexicon size.

Each dimension corresponds
to a word in the lexicon

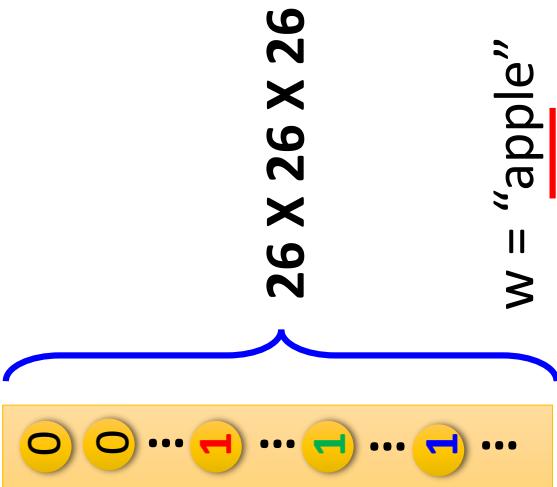
The dimension for the word
is 1, and others are 0
apple = [1 0 0 0 0]
bag = [0 1 0 0 0]
cat = [0 0 1 0 0]
dog = [0 0 0 1 0]
elephant = [0 0 0 0 1]

Beyond 1-of-N encoding

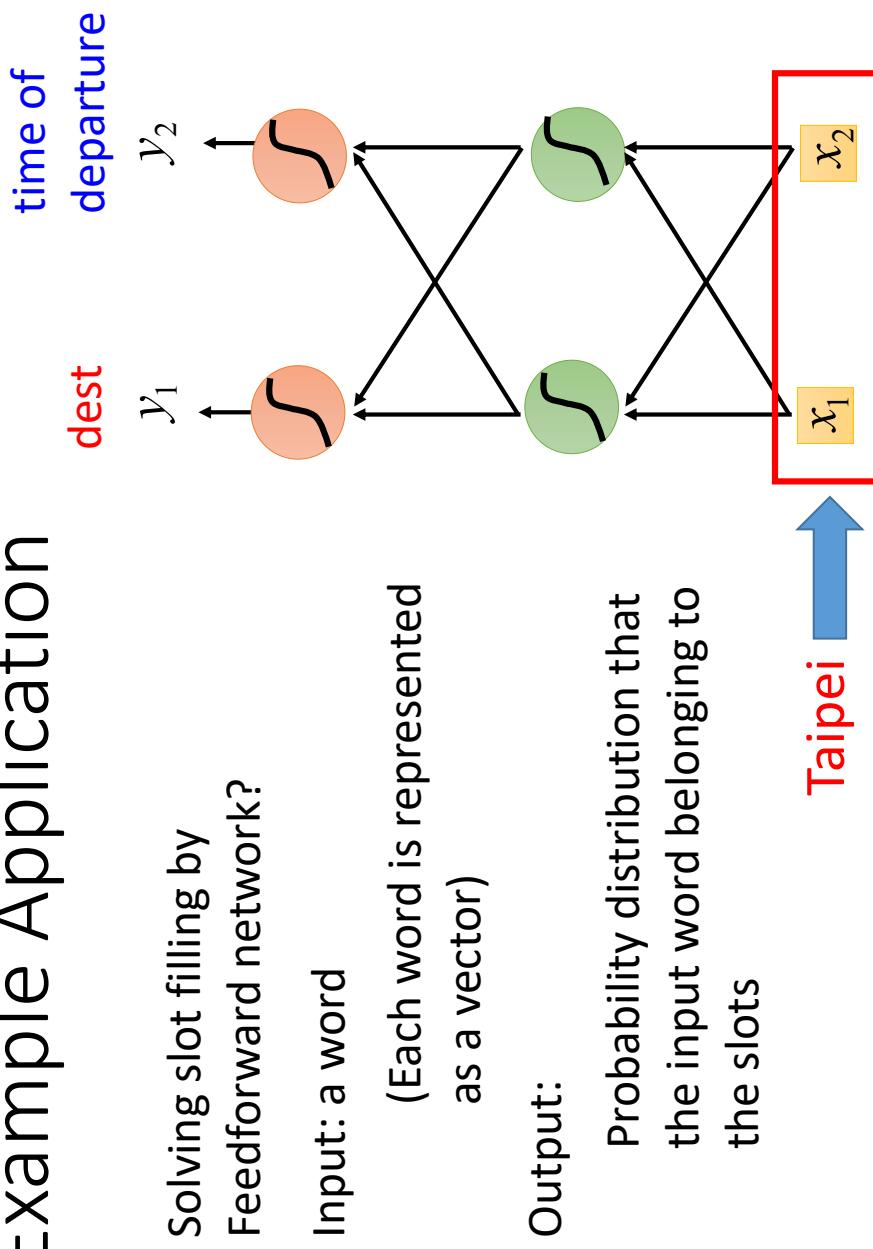
Dimension for “Other”



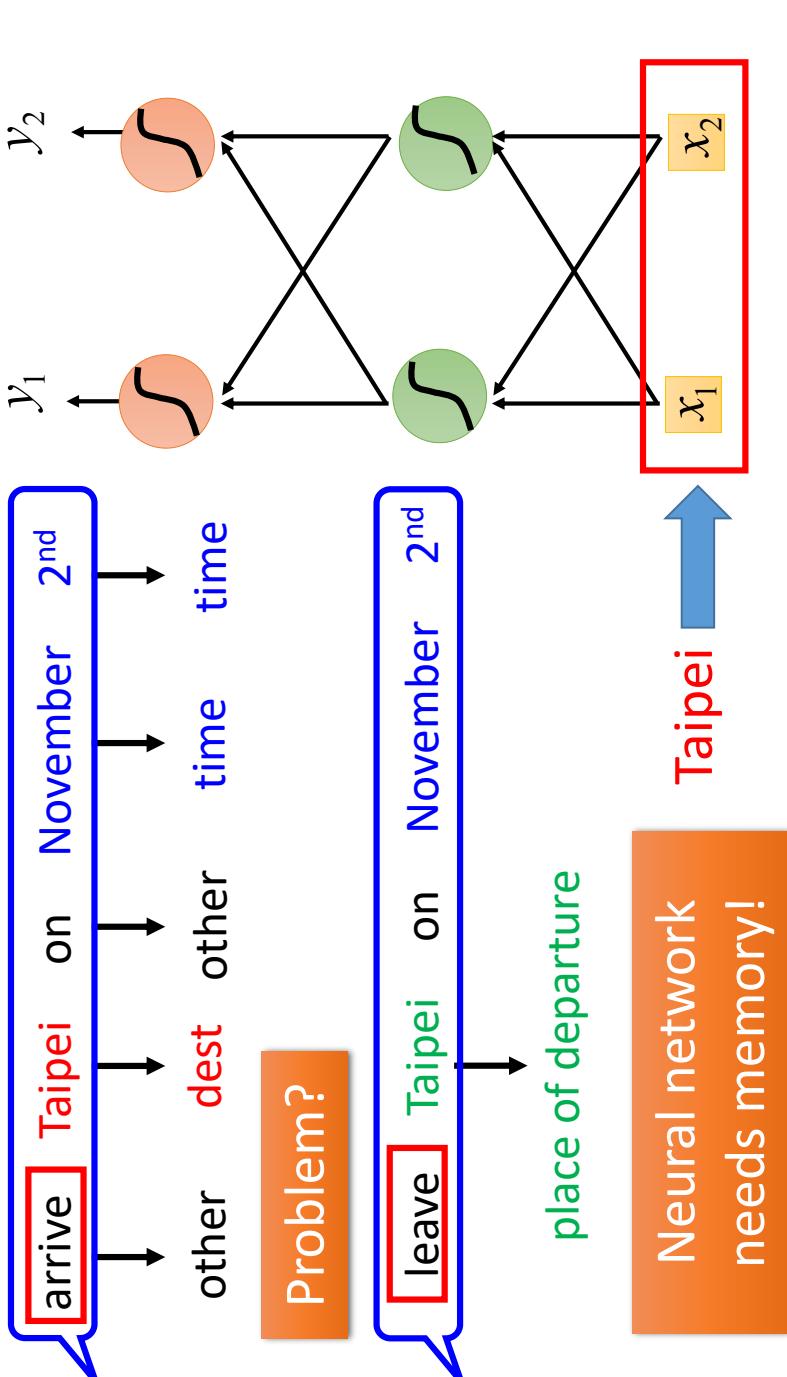
Word hashing



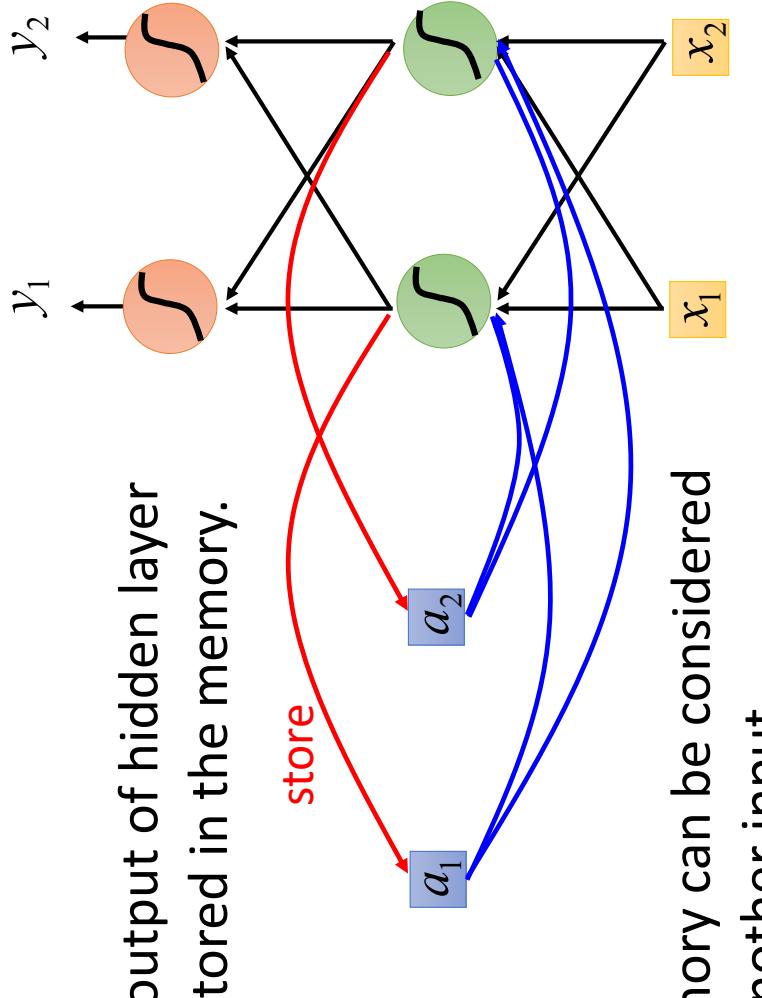
Example Application



Example Application



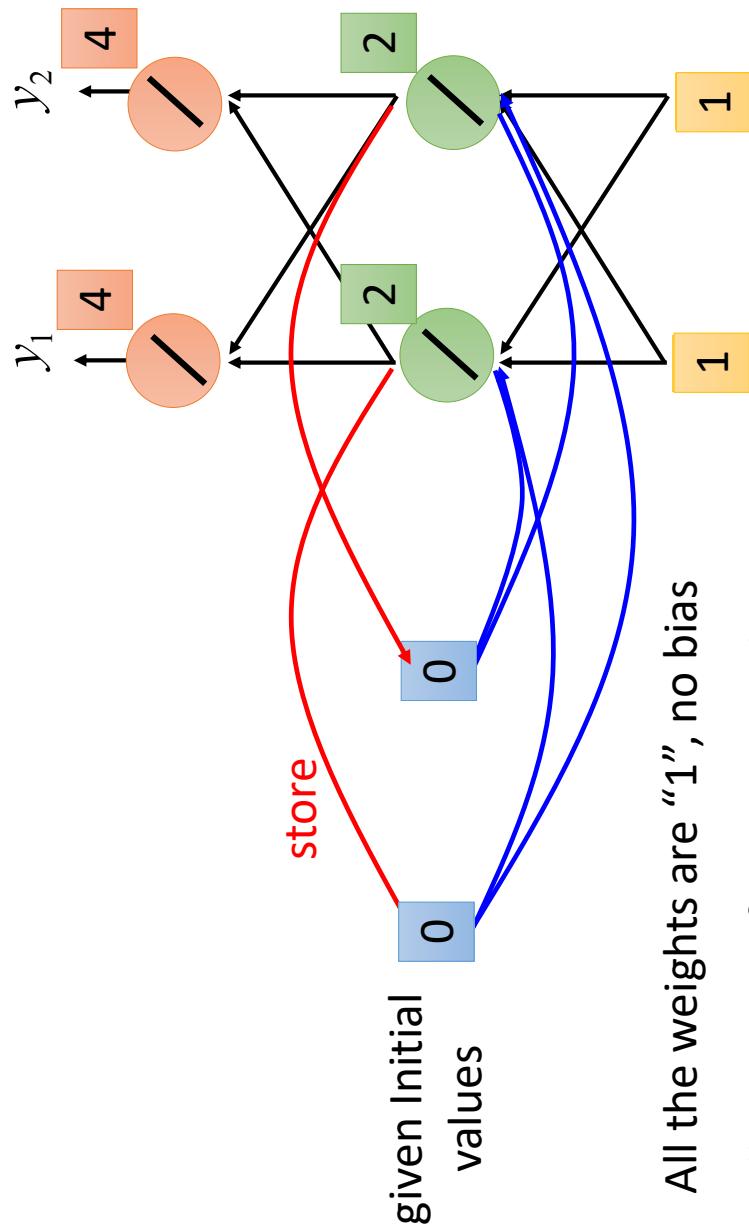
Recurrent Neural Network (RNN)



Example

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

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



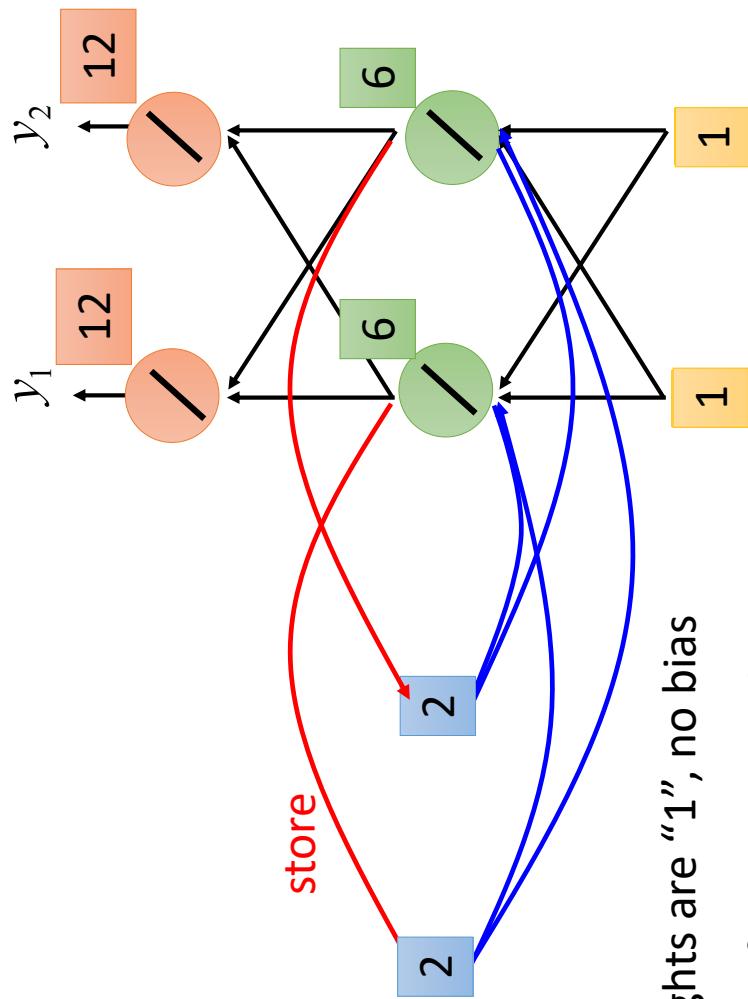
All the weights are “1”, no bias

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



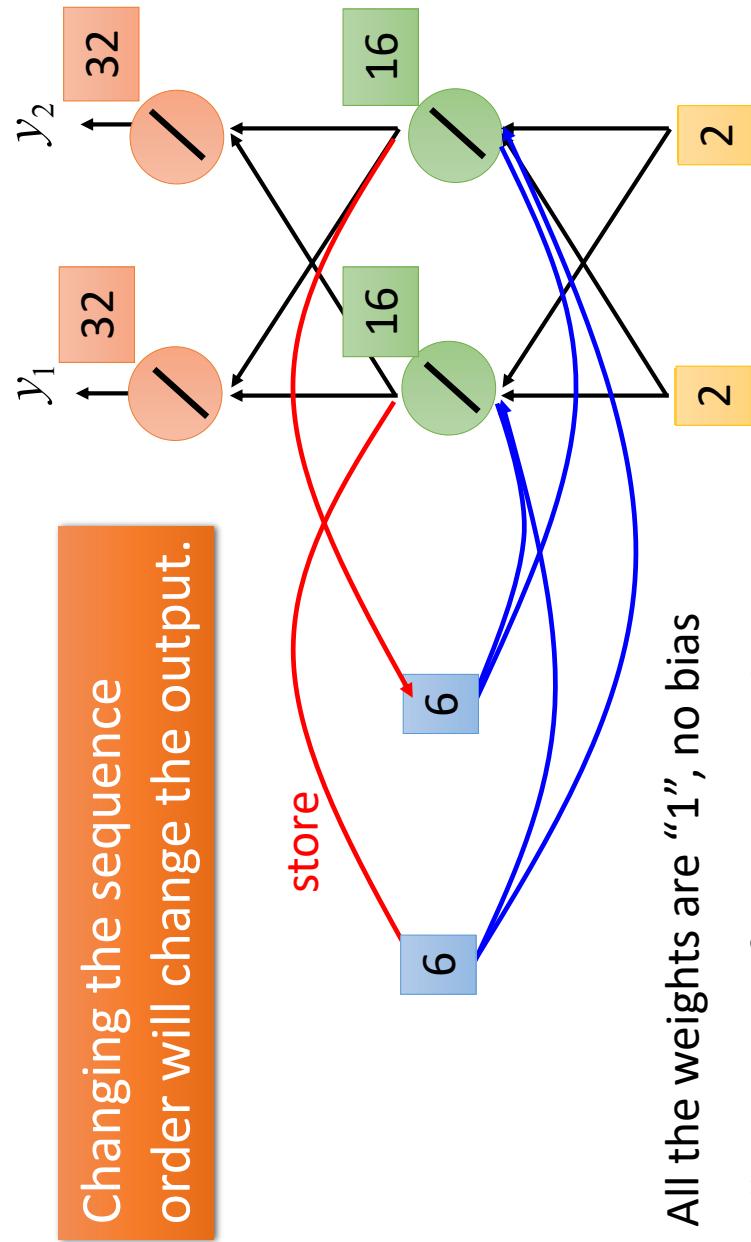
All the weights are “1”, no bias

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

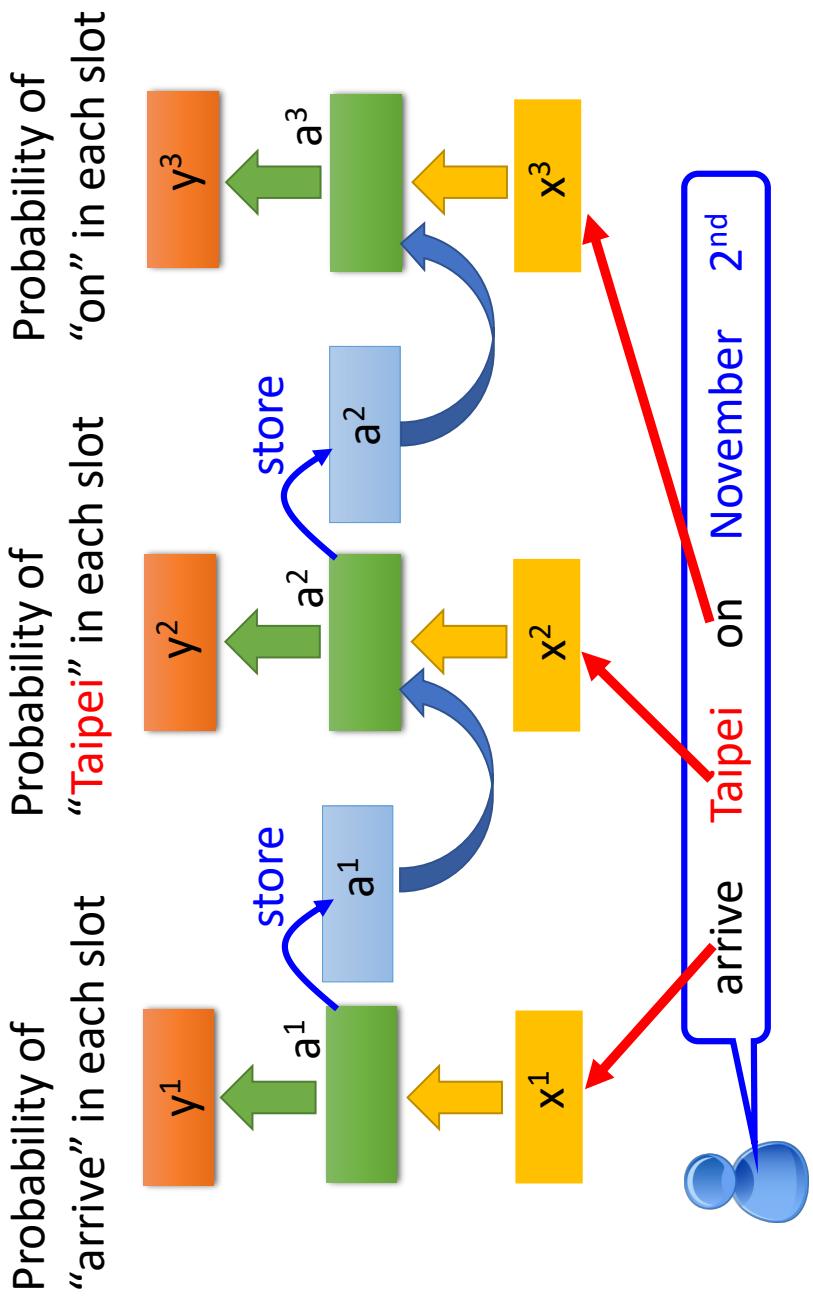


All the weights are “1”, no bias

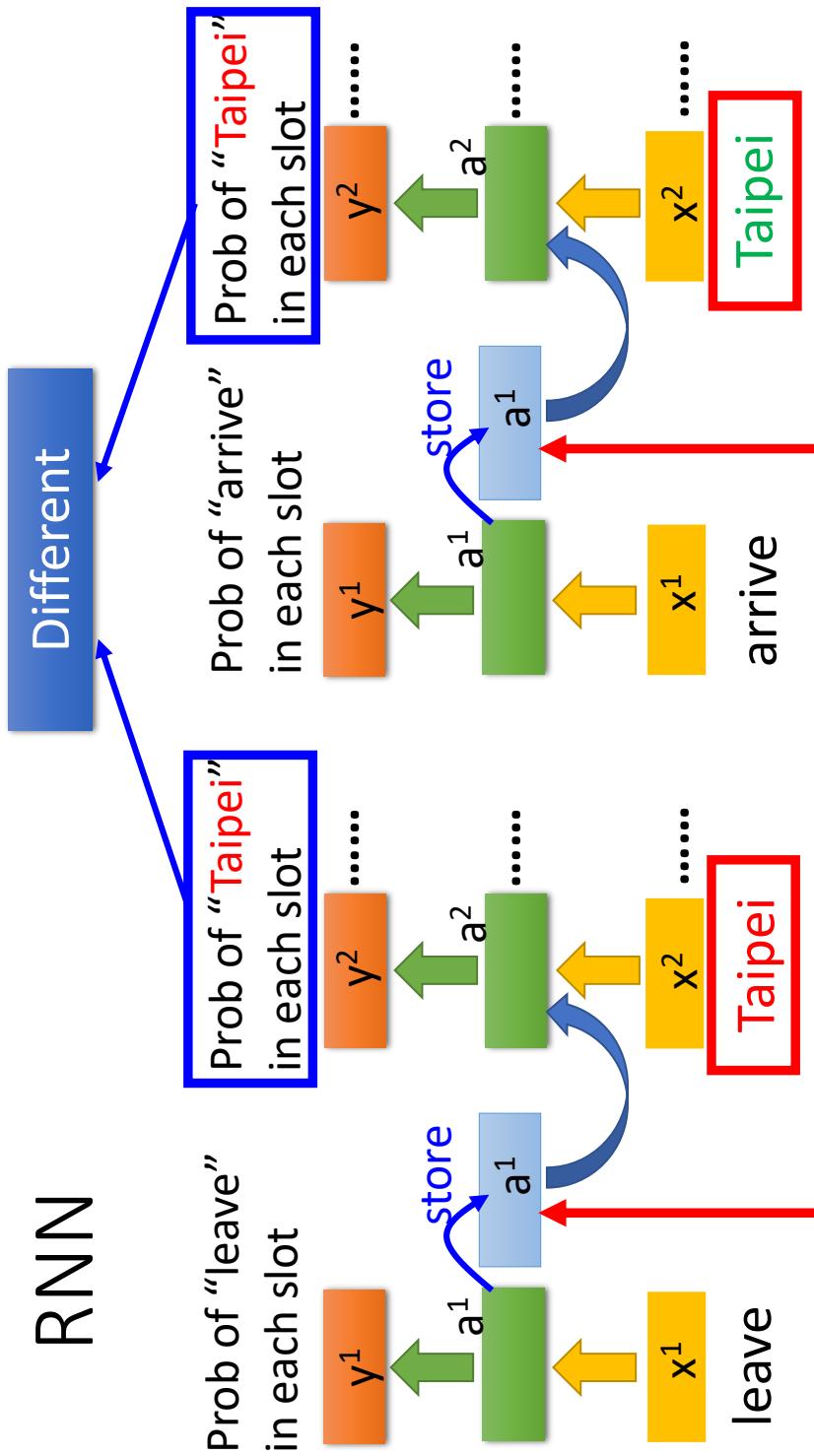
All activation functions are linear

RNN

The same network is used again and again.

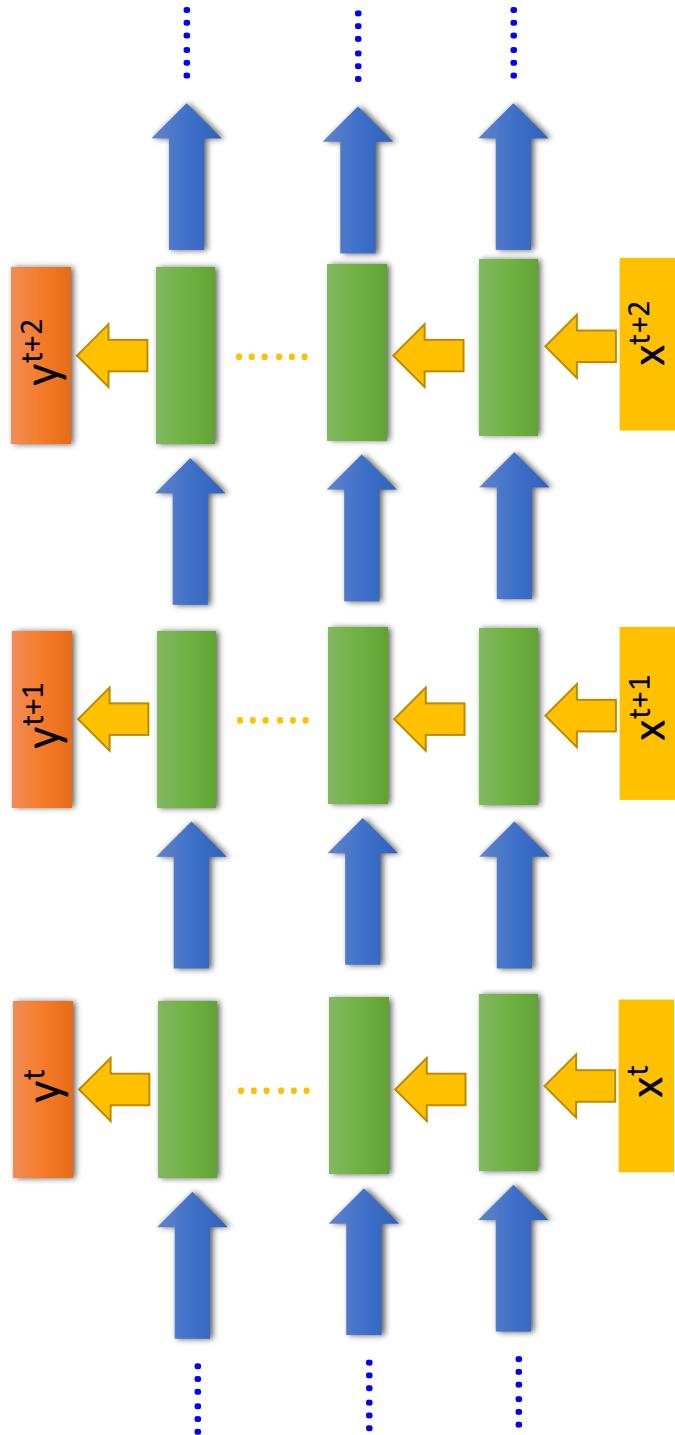


RNN

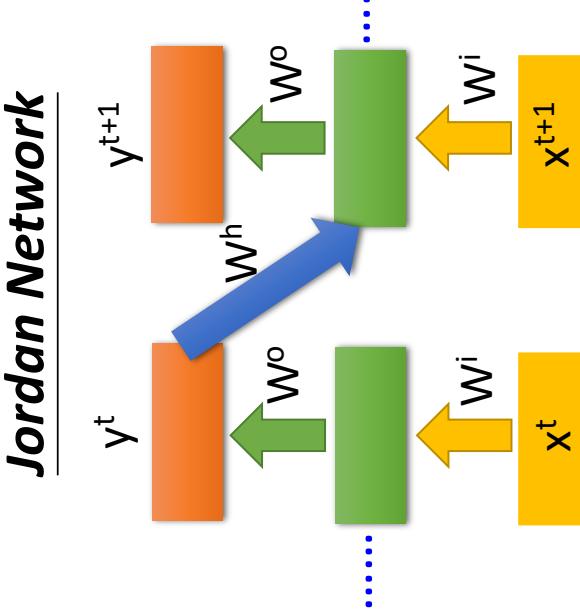
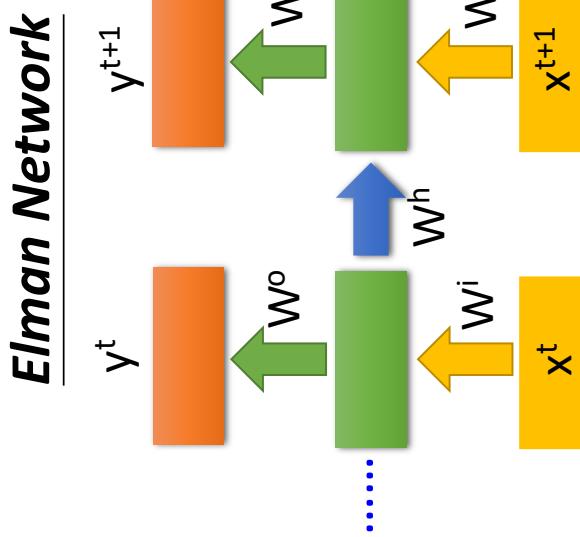


The values stored in the memory is different.

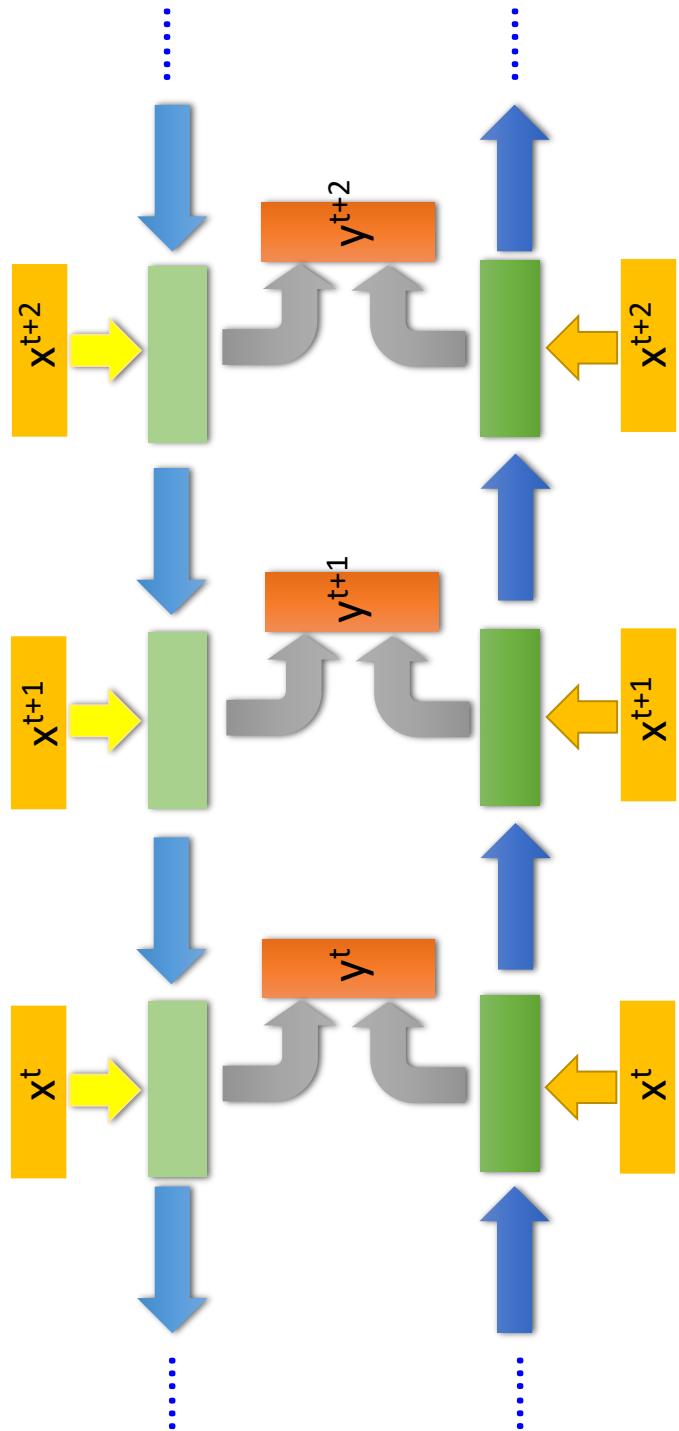
Of course it can be deep ...



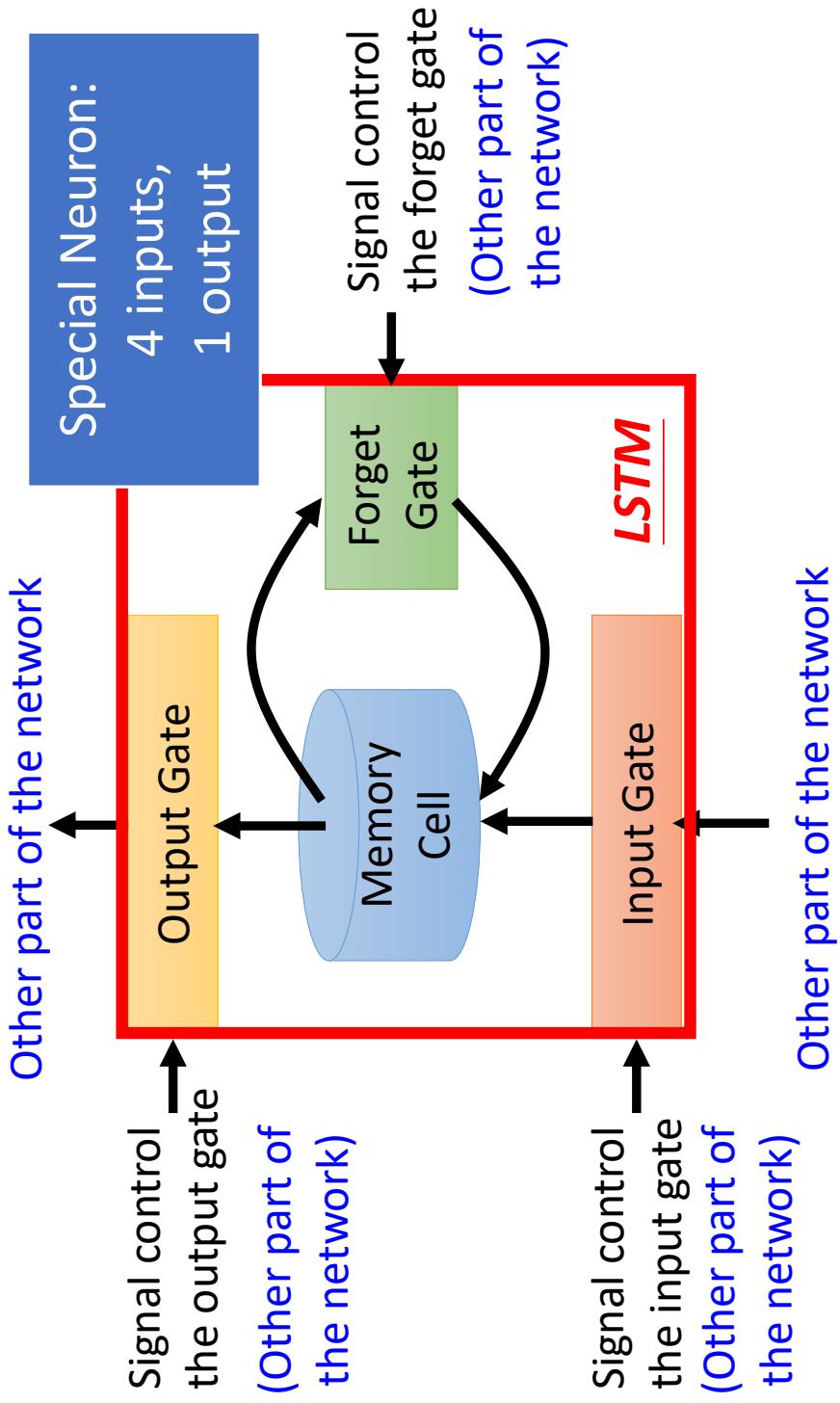
Elman Network & Jordan Network

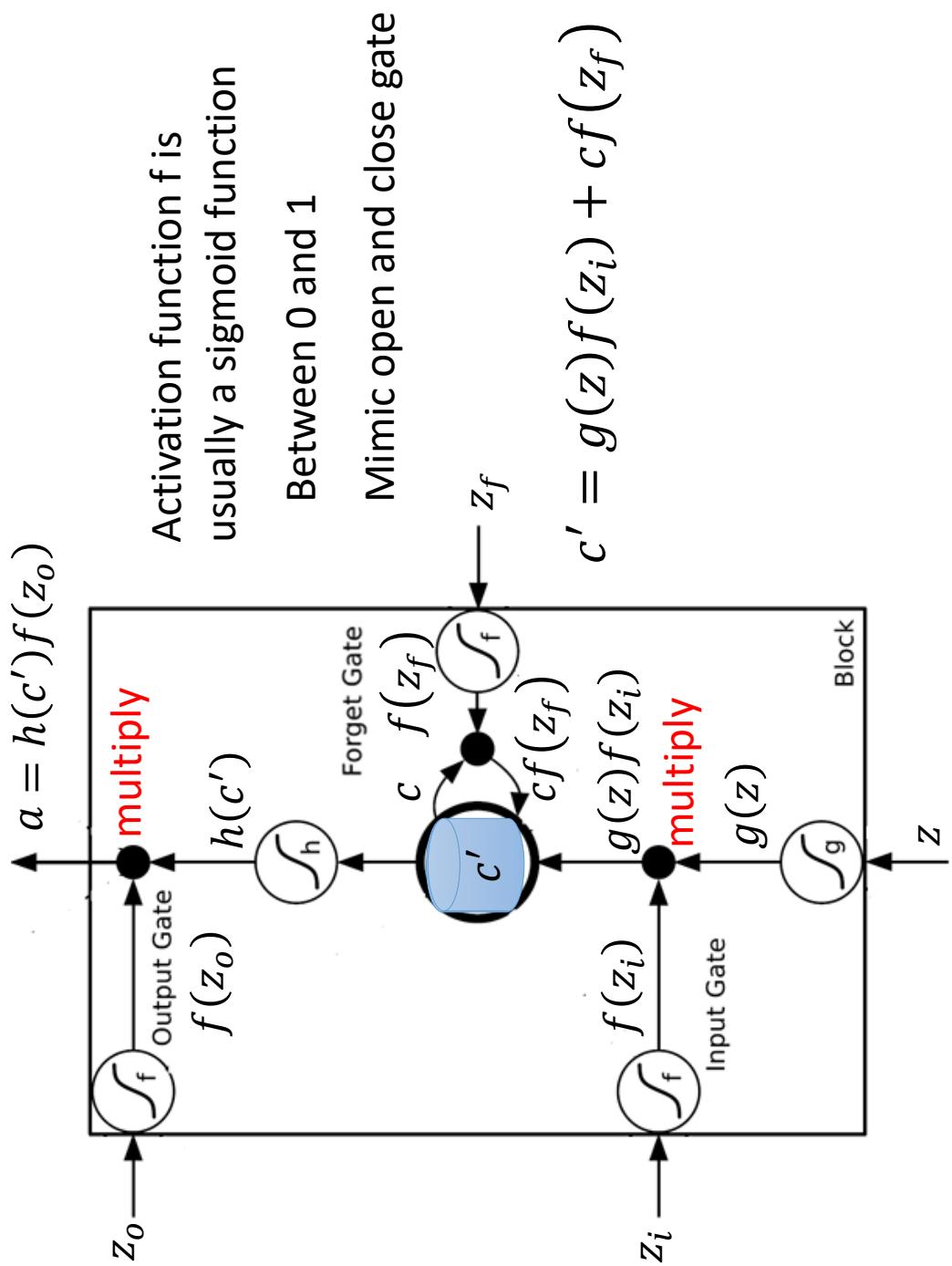


Bidirectional RNN

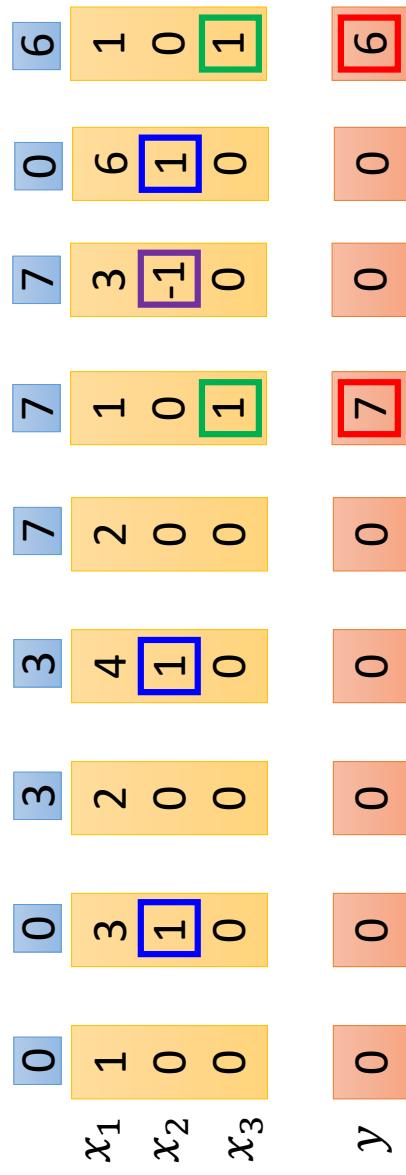


Long Short-term Memory (LSTM)





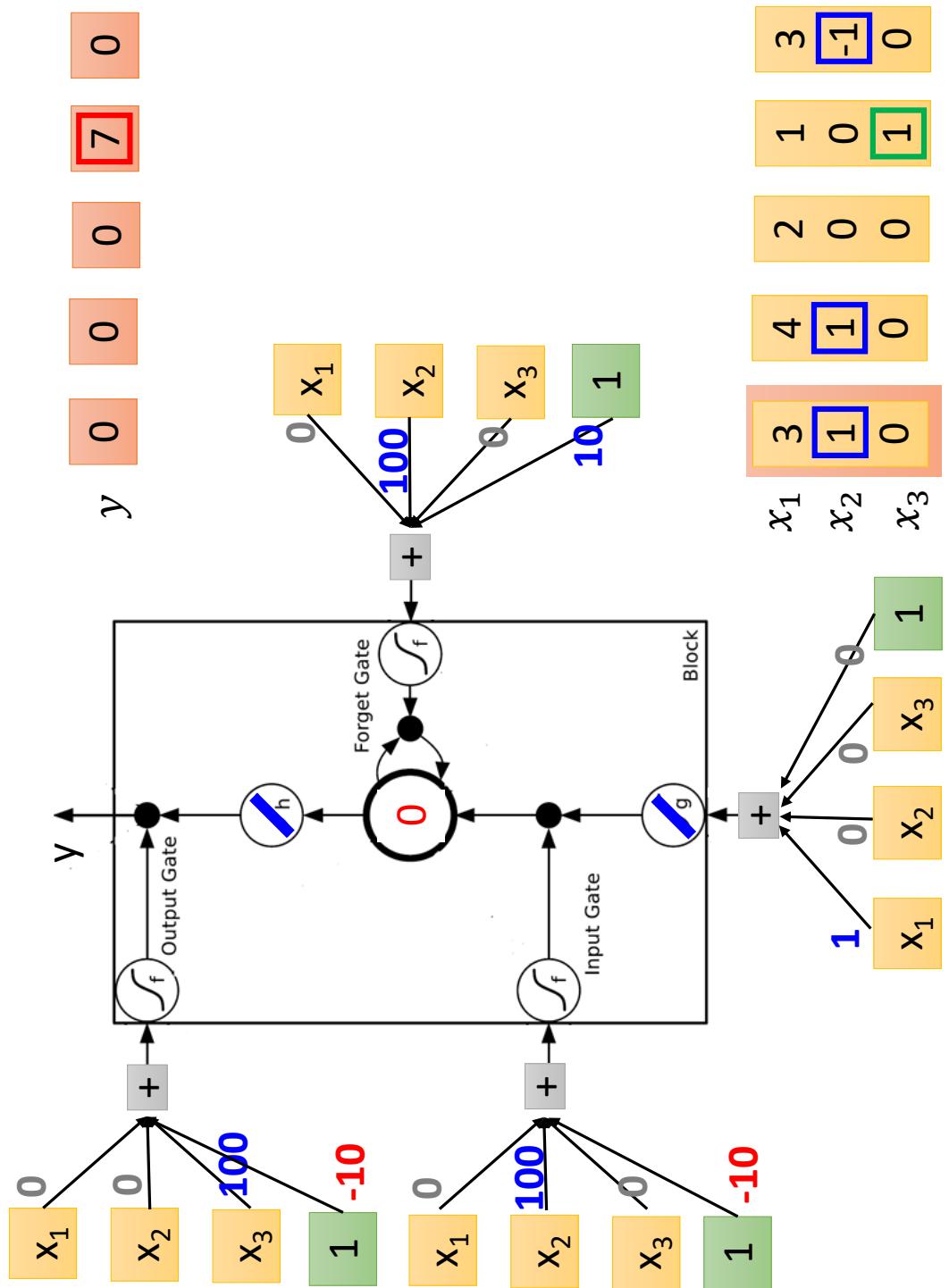
LSTM - Example

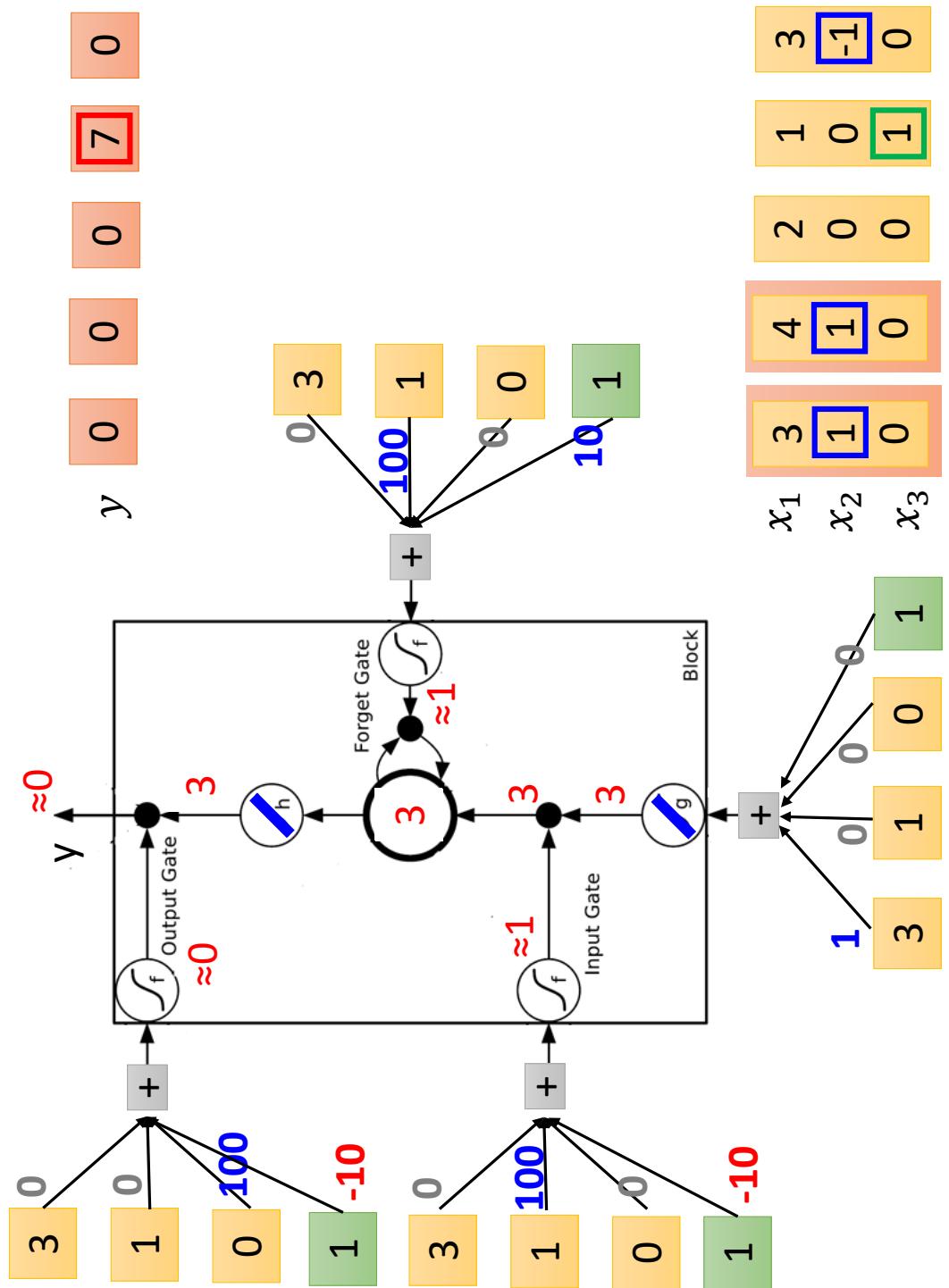


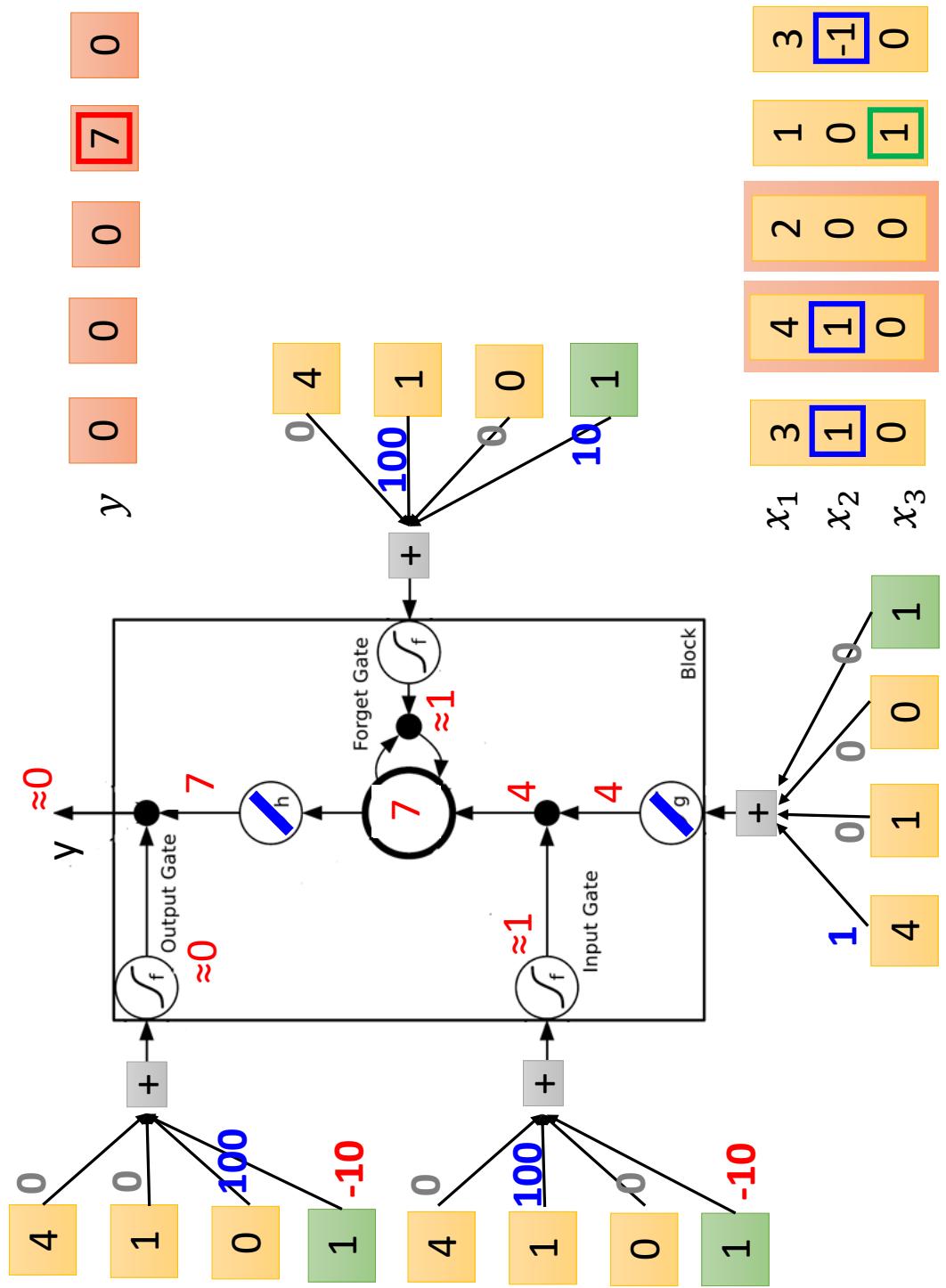
When $x_2 = 1$, add the numbers of x_1 into the memory

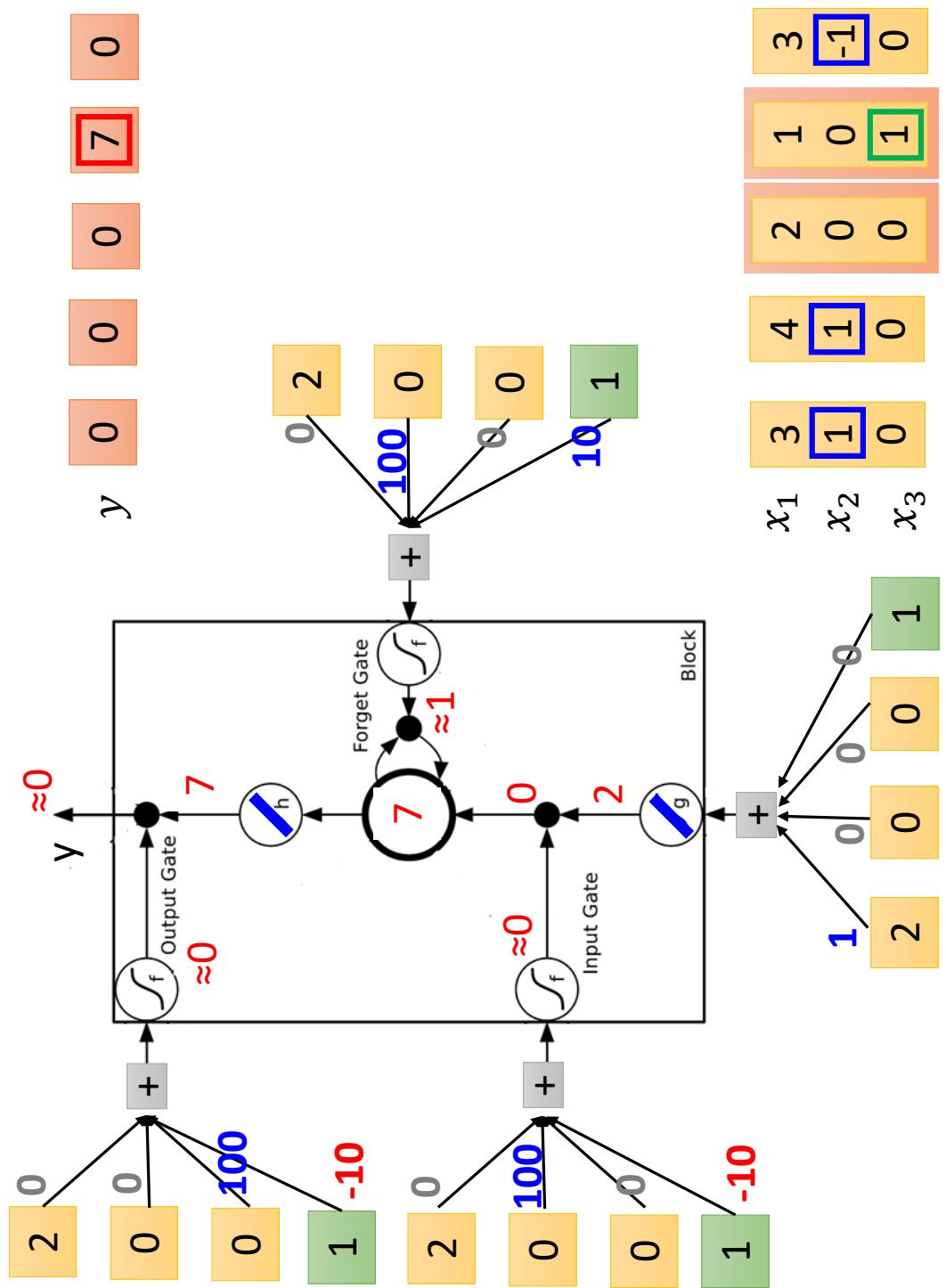
When $x_2 = -1$, reset the memory

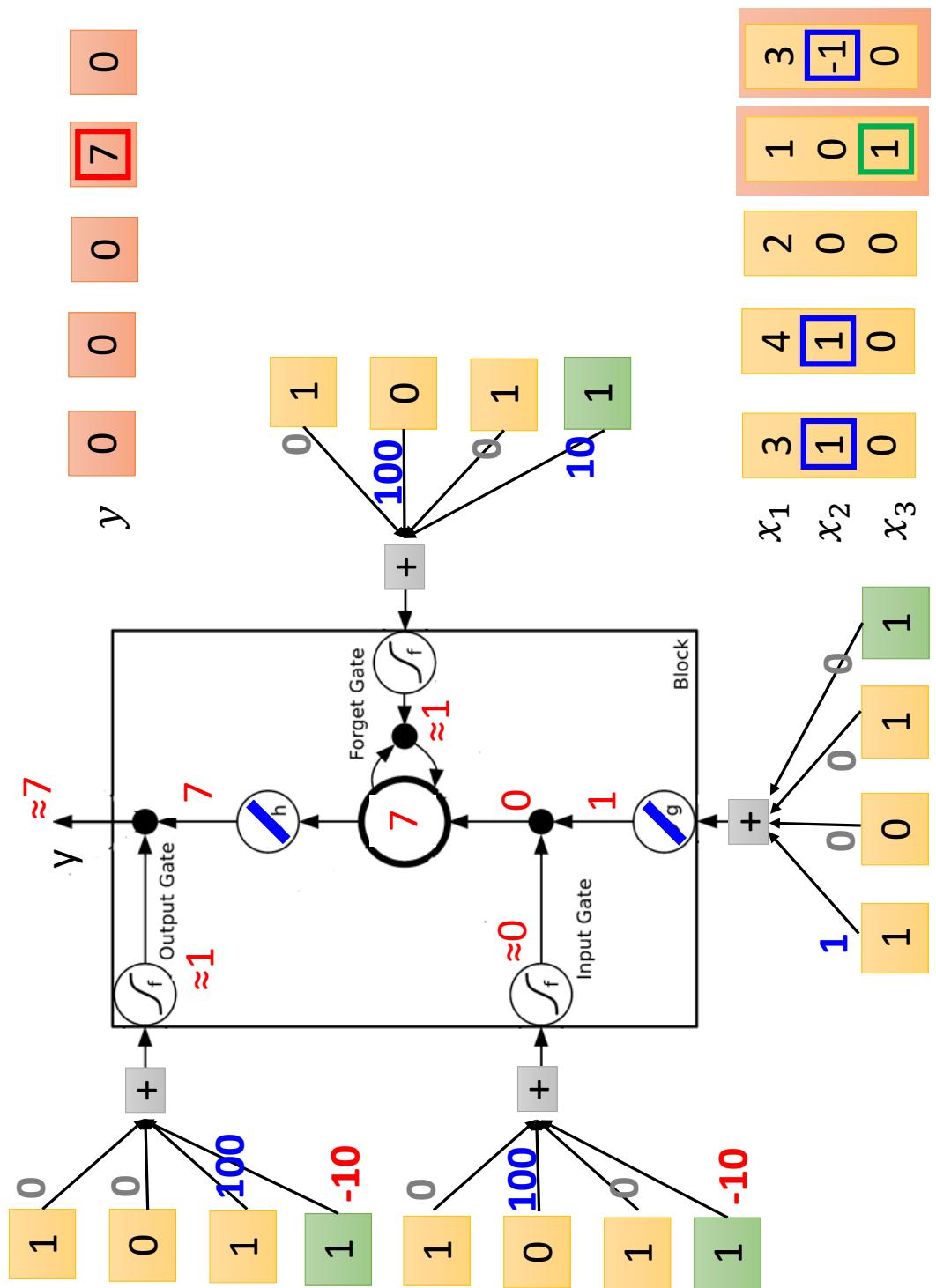
When $x_3 = 1$, output the number in the memory.

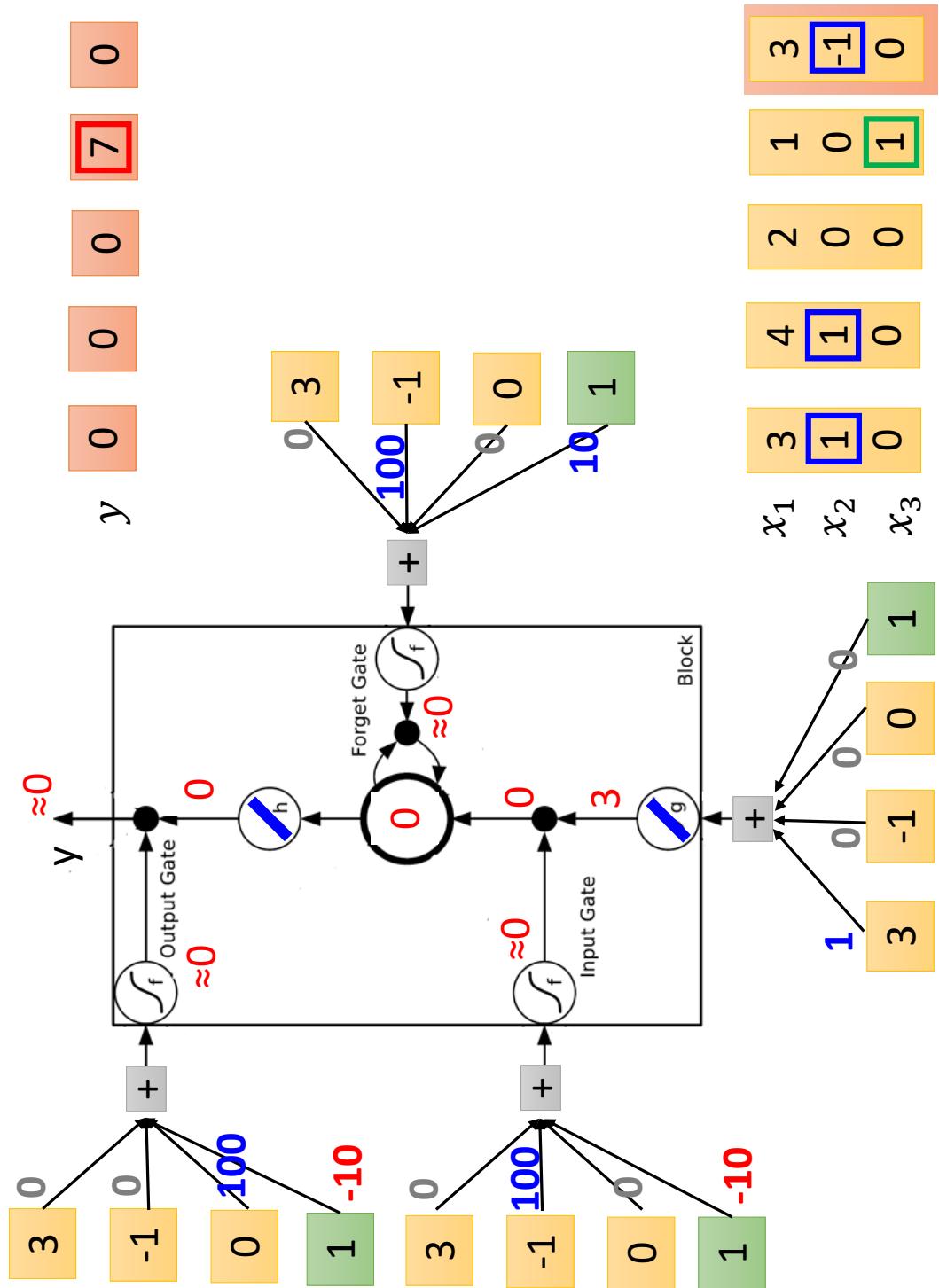






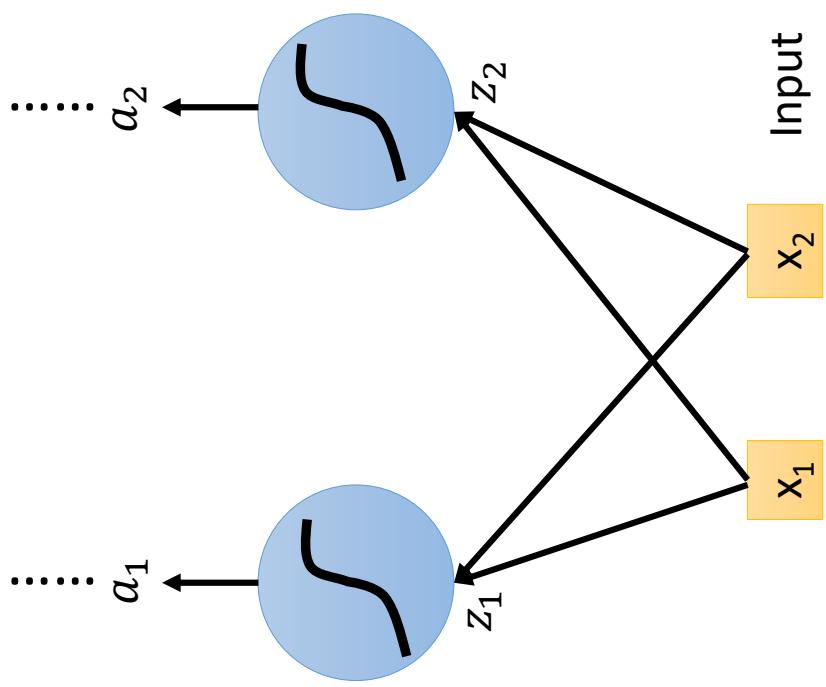


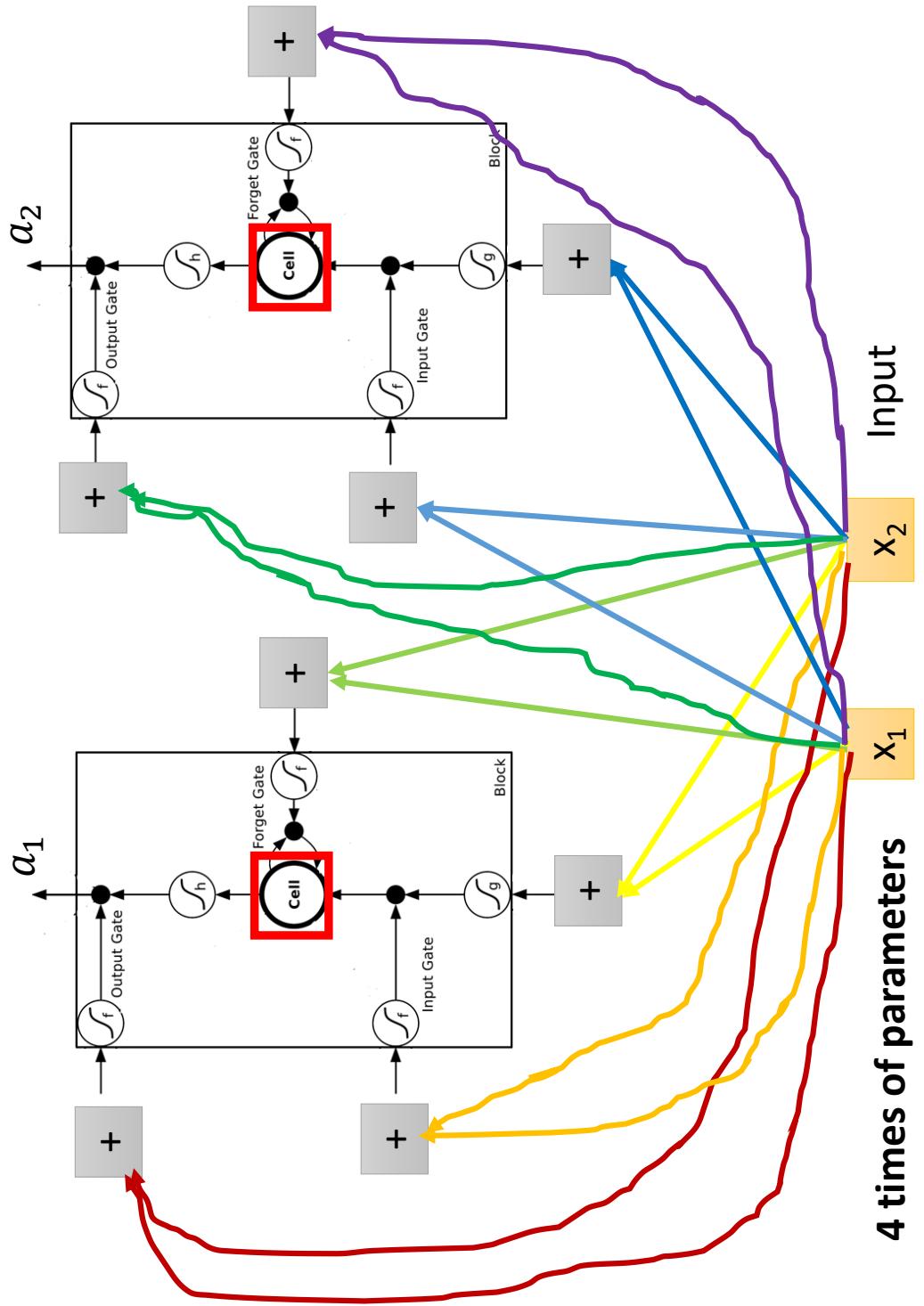




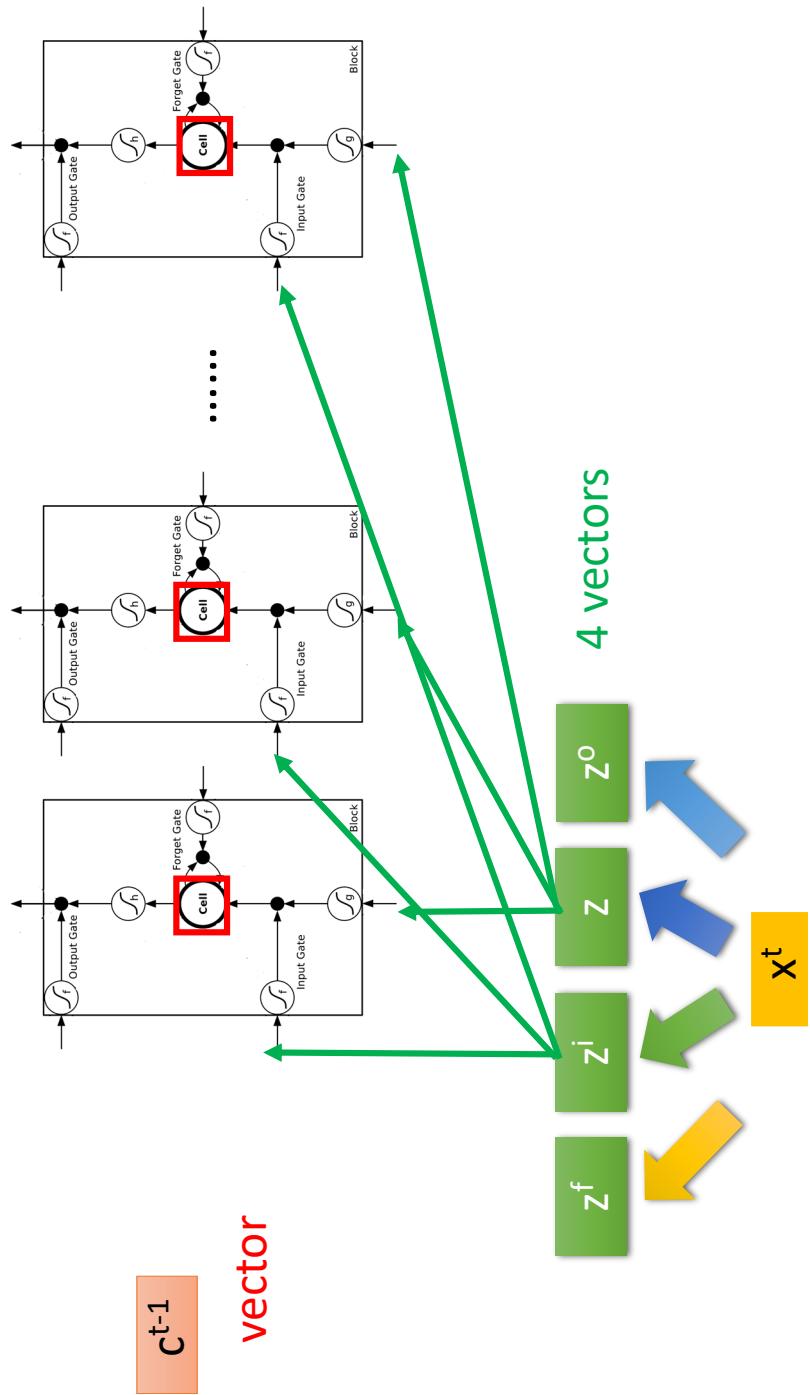
Original Network:

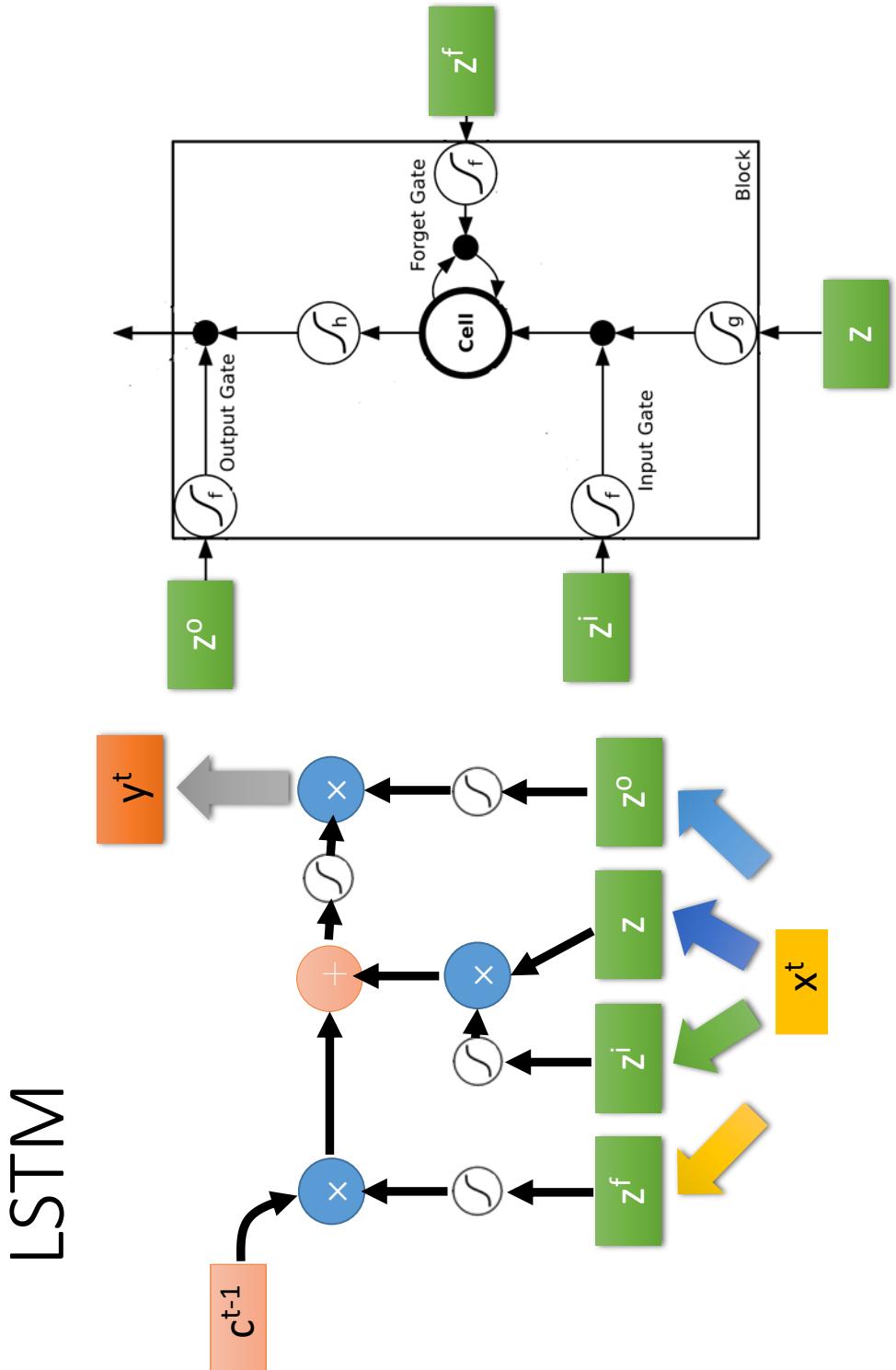
➤ Simply replace the neurons with LSTM





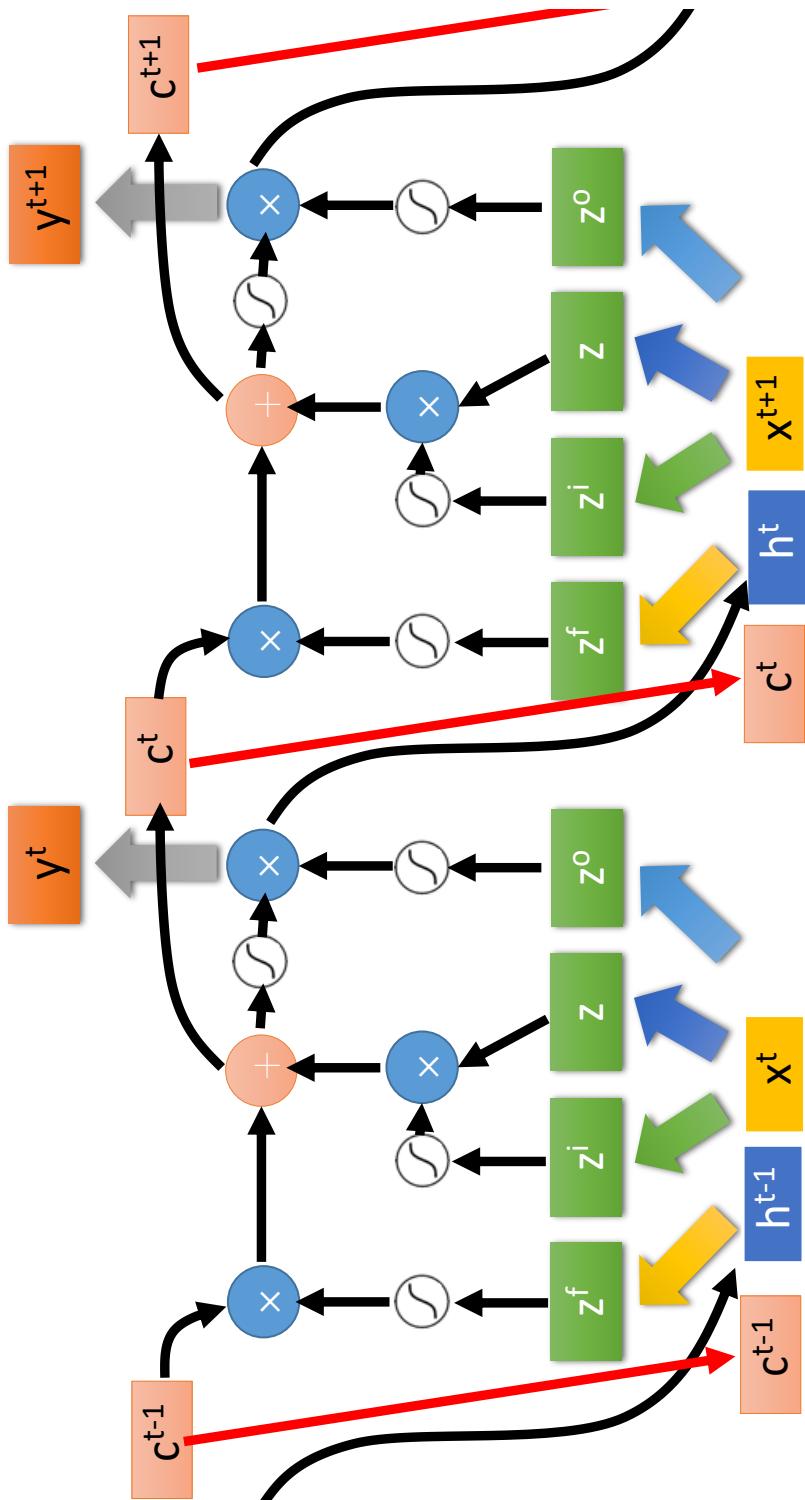
LSTM



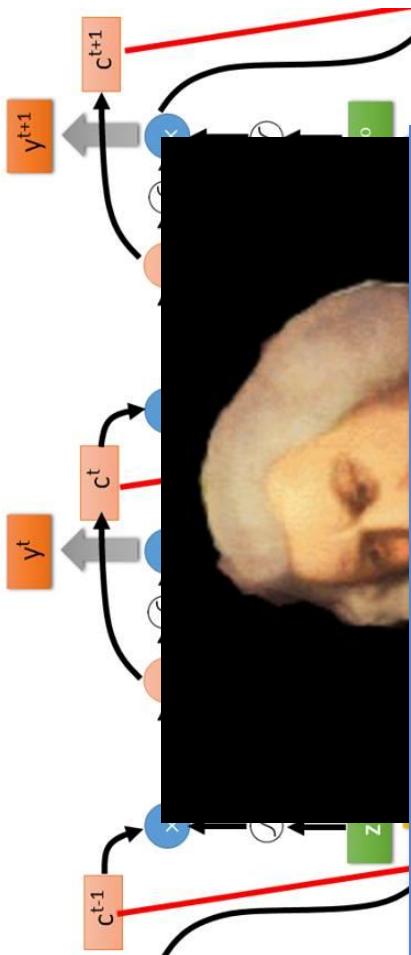


LSTM

Extension: "peephole"



Multiple-layer LSTM



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

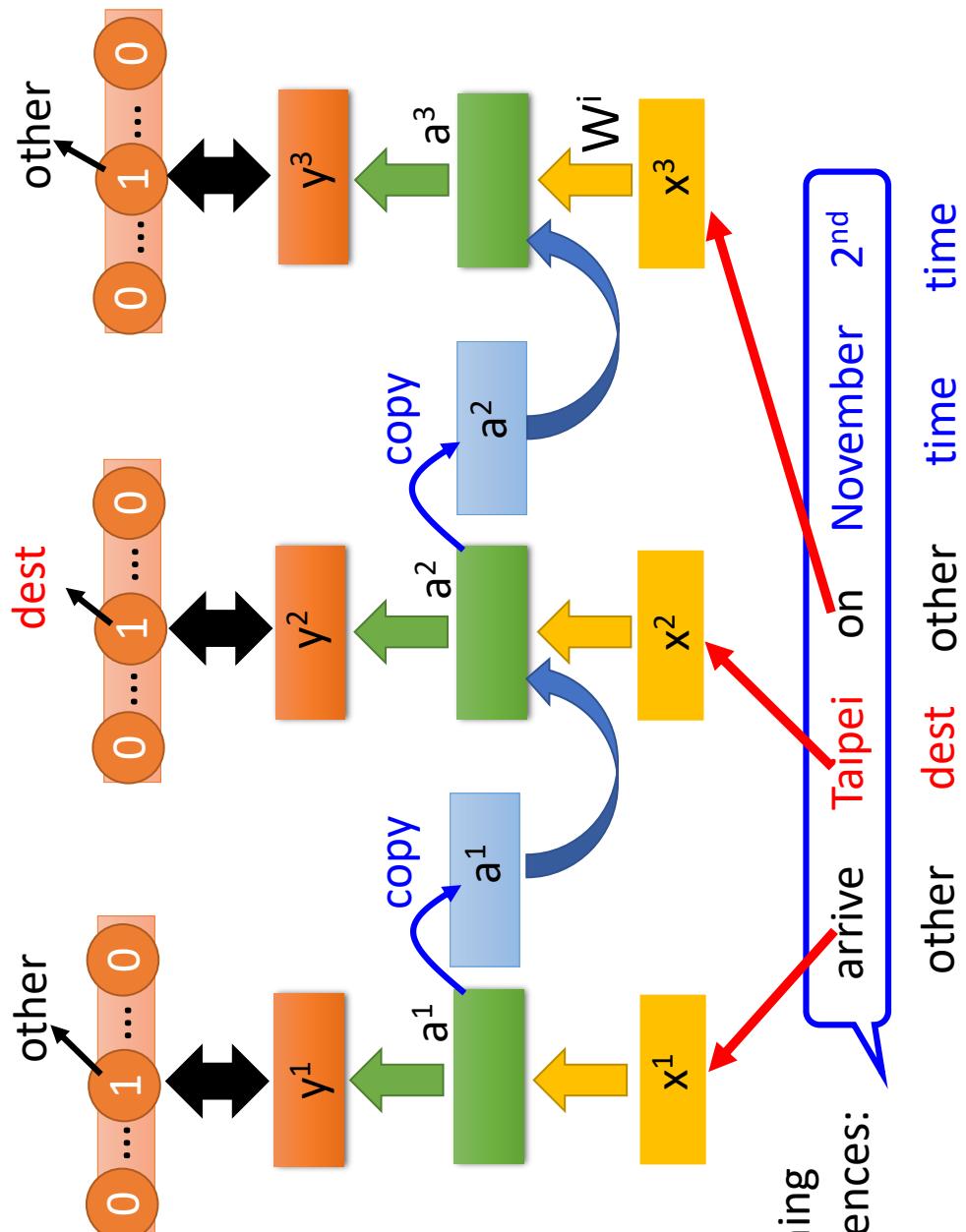
Keras supports
“LSTM”, “GRU”, “SimplerNN” layers

This is quite
standard now.

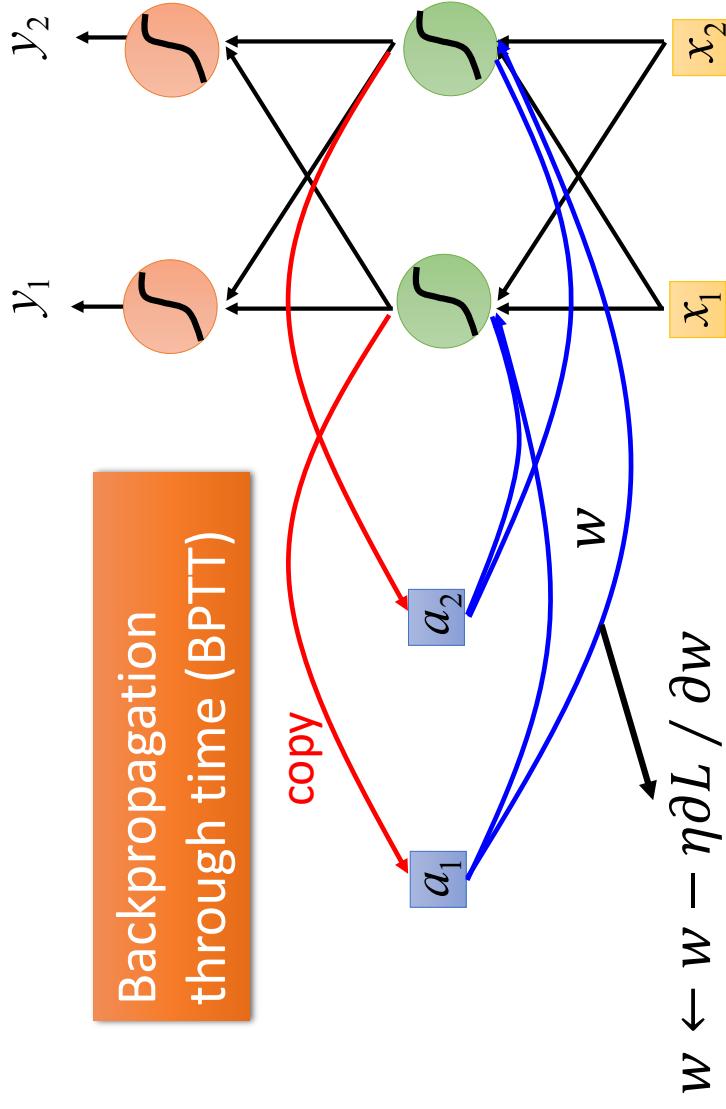
<https://img.komicolle.org/2015-09-20/src/14426967627131.gif>

c^{t-1} h^{t-1} x^t h^t c^t x^{t+1}

Learning Target



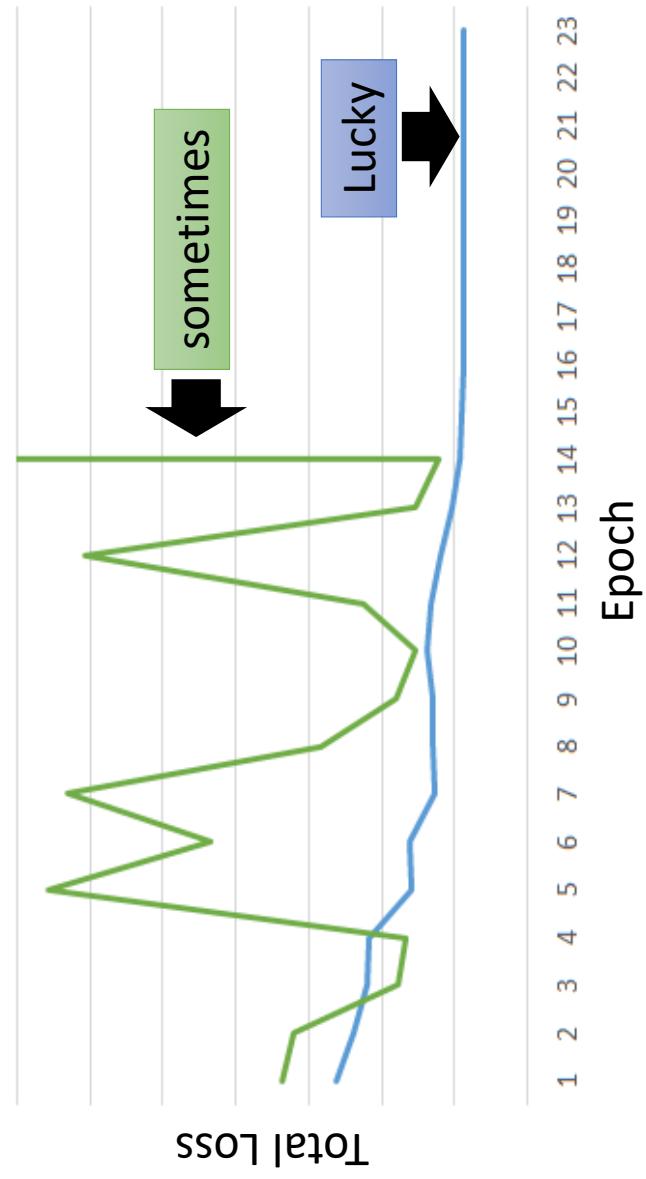
Learning



Unfortunately

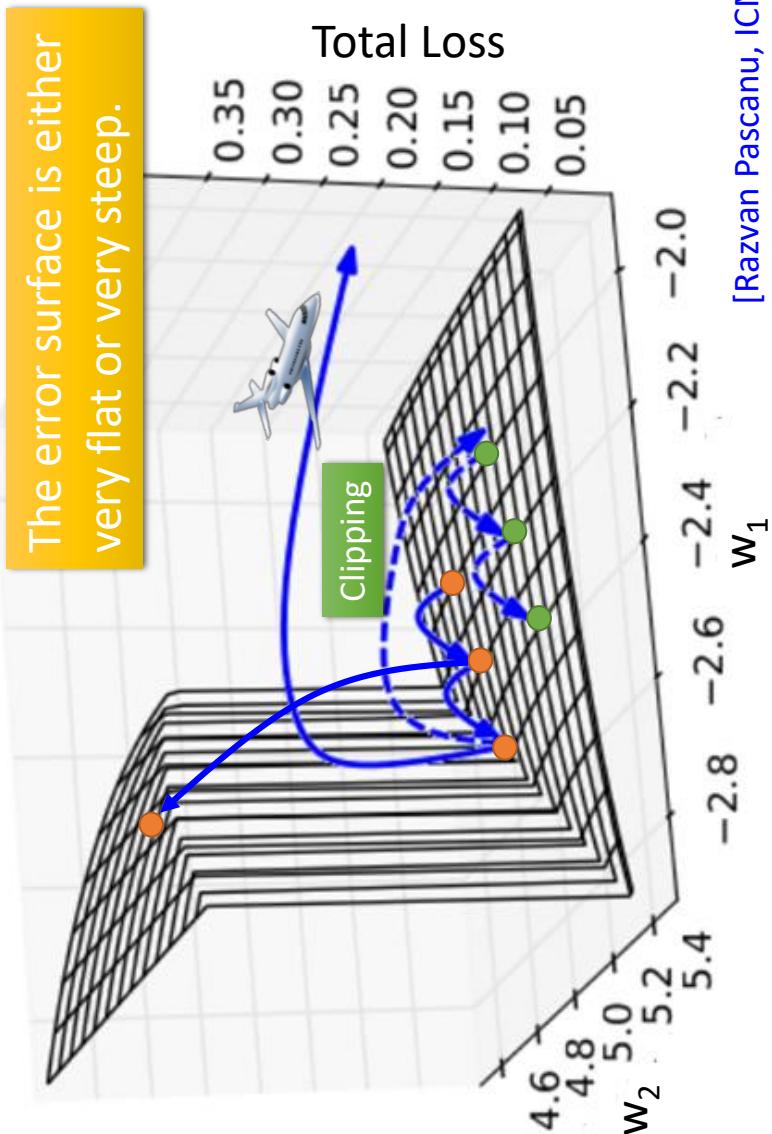
- RNN-based network is not always easy to learn

Real experiments on Language modeling



感謝曾柏翔同學
提供實驗結果

The error surface is rough.



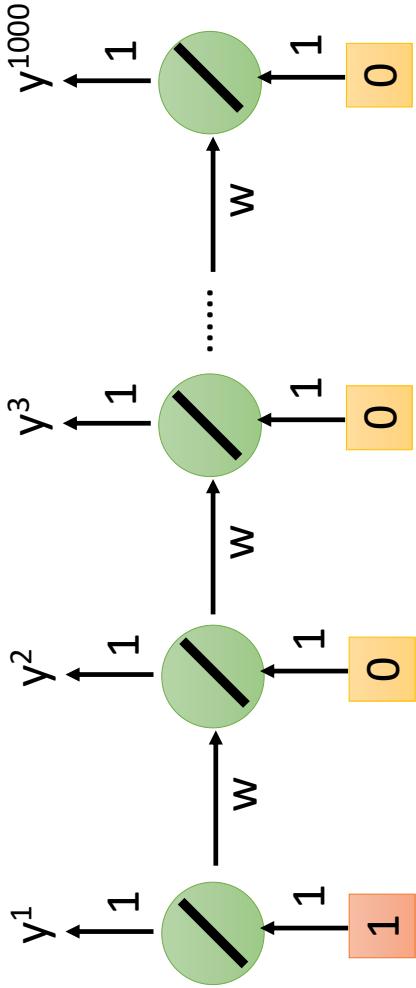
[Razvan Pascanu, ICML'13]

Why?

$$\begin{array}{ll} w = 1 & \rightarrow y^{1000} = 1 \\ w = 1.01 & \uparrow \\ w = 0.99 & \rightarrow y^{1000} \approx 0 \\ w = 0.01 & \uparrow \end{array}$$

$$\begin{array}{ll} \text{Large } \partial L / \partial w & \text{Small Learning rate?} \\ \text{small } \partial L / \partial w & \text{Large Learning rate?} \end{array}$$

$=w^{999}$

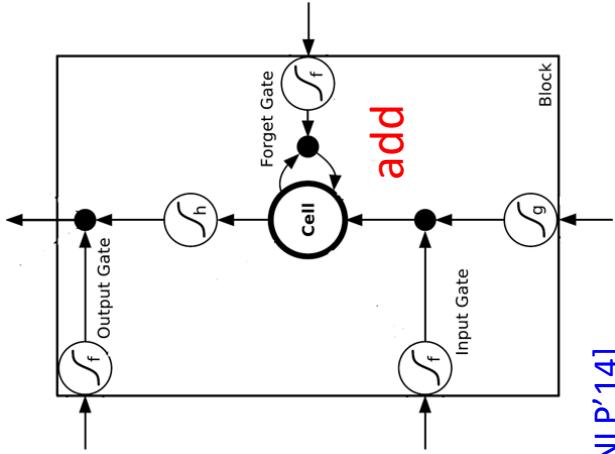


Toy Example

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.)

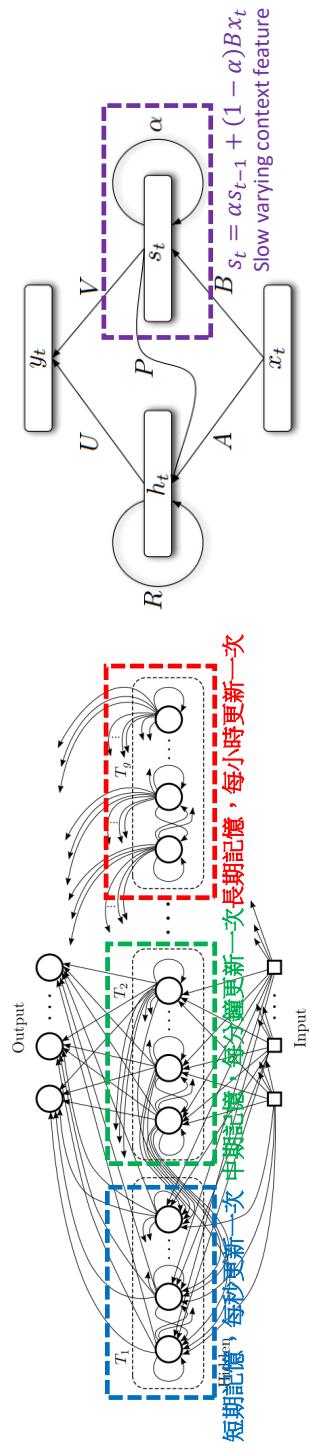
Gated Recurrent Unit (GRU):
simpler than LSTM

[Cho, EMNLP'14]

Helpful Techniques

Clockwise RNN

Structurally Constrained Recurrent Network (SCRN)



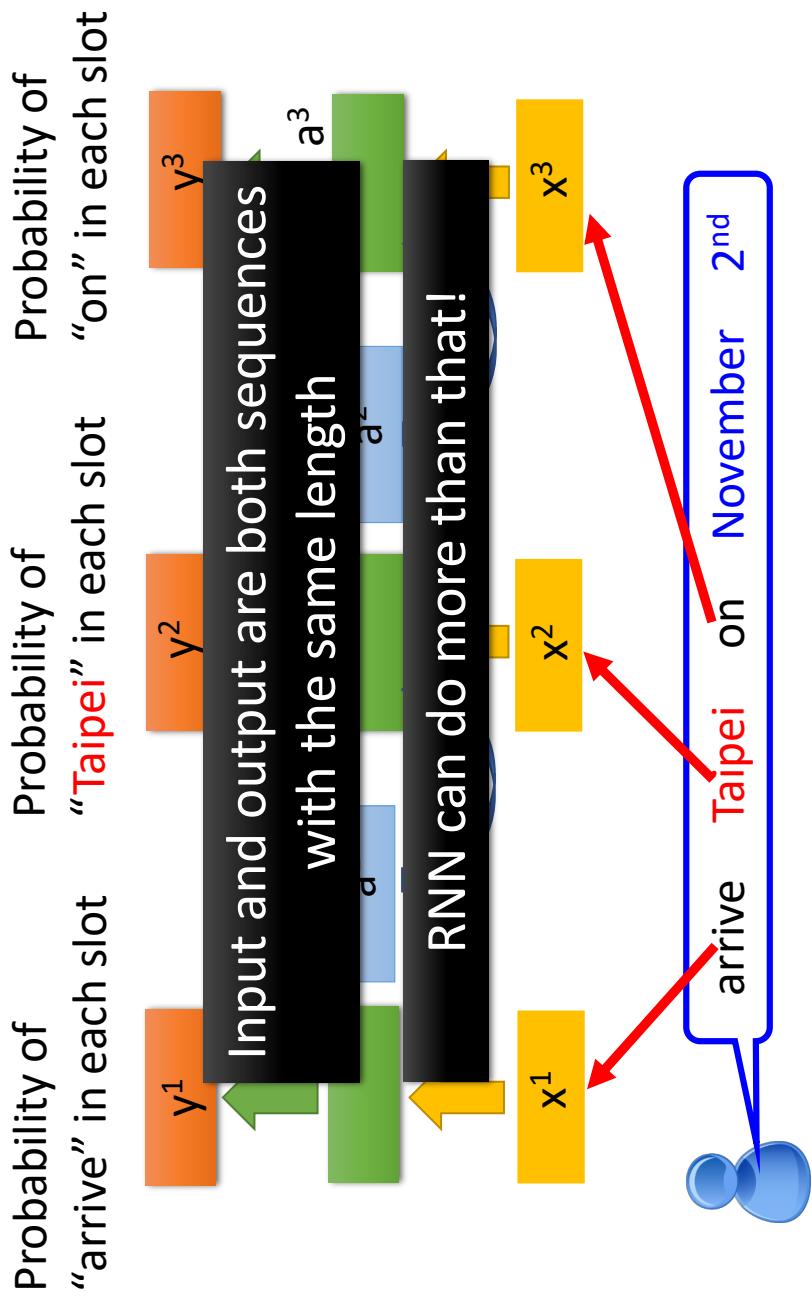
[Jan Koutník, JMLR'14]

[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



Many to one

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

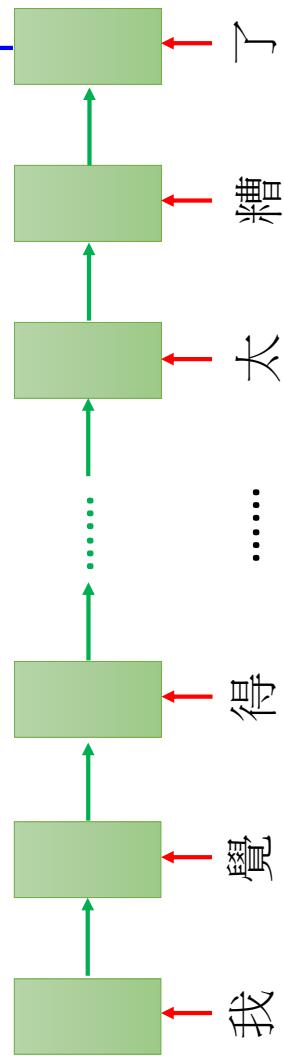
Sentiment Analysis

看了這部電影覺
得很高興.....

這部電影太糟了
.....

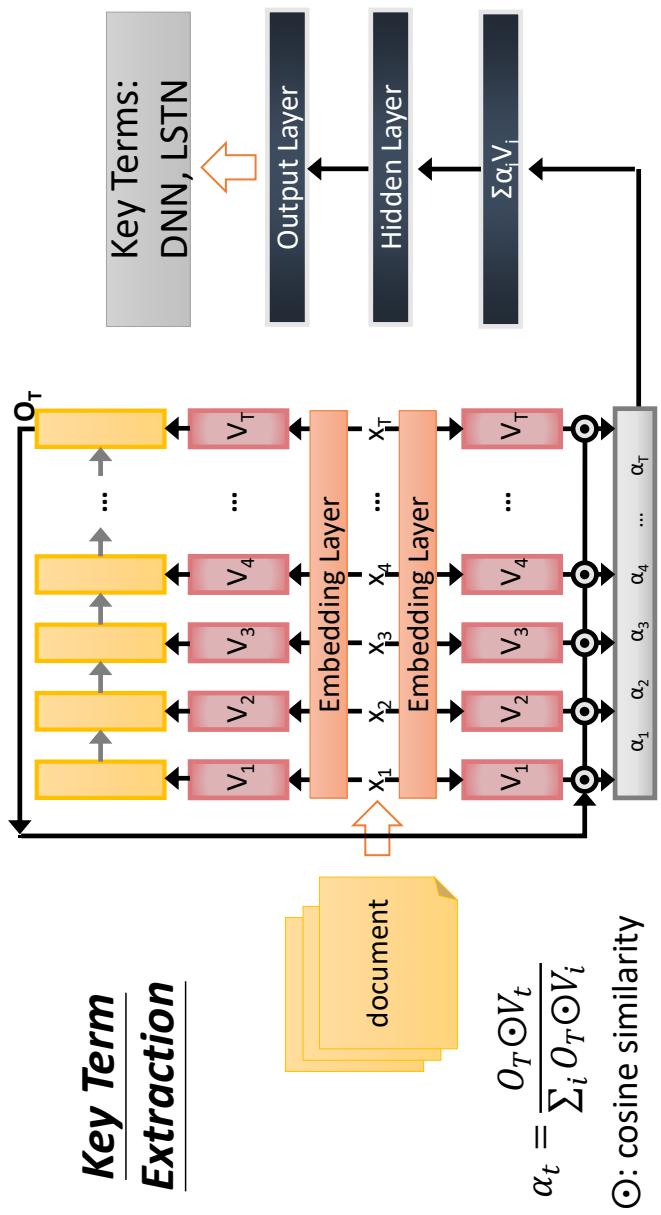
超好雷
好雷
普雷
負雷
超負雷

Positive (正雷) Negative (負雷)



Many to one

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



[Shen & Lee, Interspeech 16]

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 (INTERSPEECH16), San Francisco, Sept. 2016

Many to Many (Output is shorter)

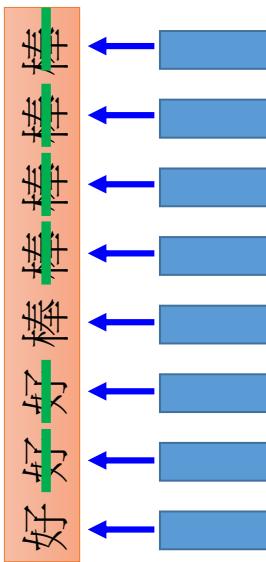
- Both input and output are both sequences, but the output is shorter.
- E.g. Speech Recognition

Problem?

Output: “好棒” (character sequence)



Why can't it be
“好棒棒”



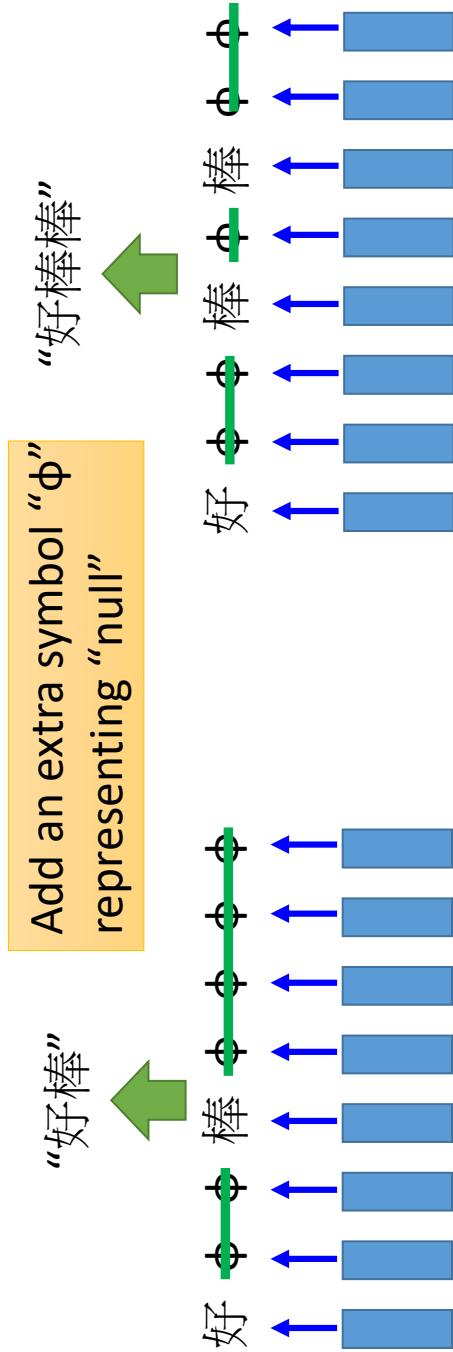
Input:



Many to Many (Output is shorter)

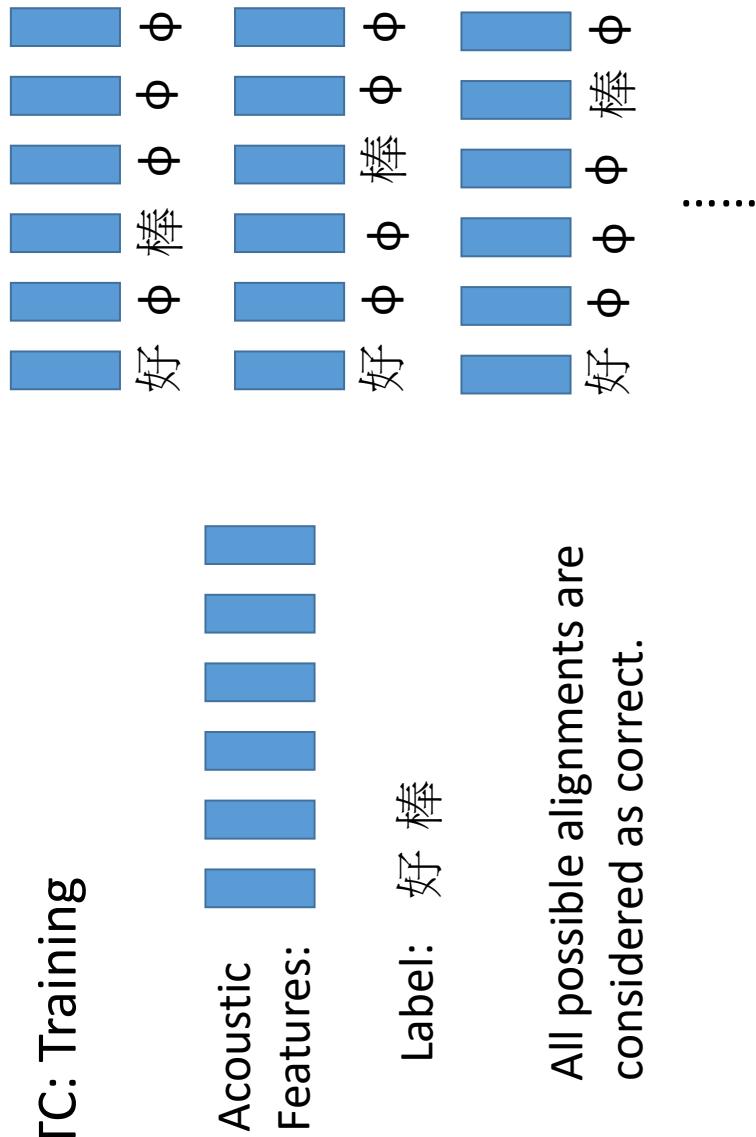
- Both input and output are both sequences, **but the output is shorter.**

- Connectionist Temporal Classification (CTC) [Alex Graves, ICM'06][Alex Graves, ICML'14][Hasim Sak, Interspeech'15][Jie Li, Interspeech'15][Andrew Senior, ASRU'15]



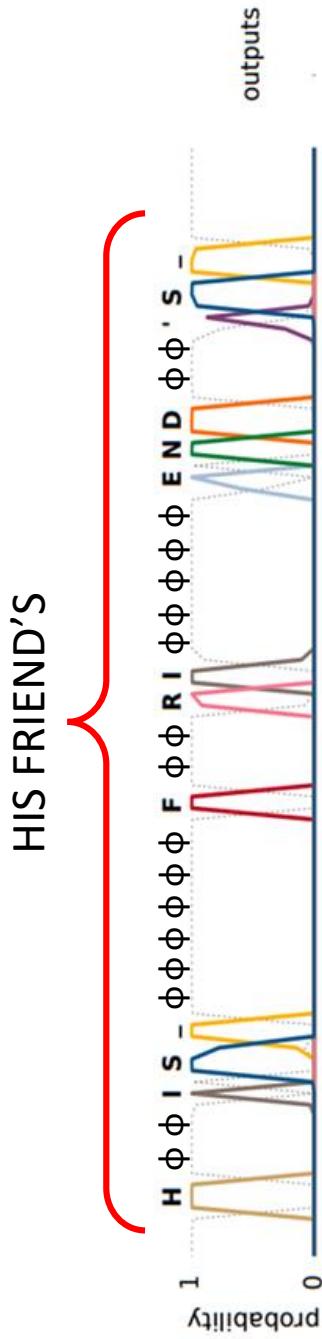
Many to Many (Output is shorter)

- CTC: Training



Many to Many (Output is shorter)

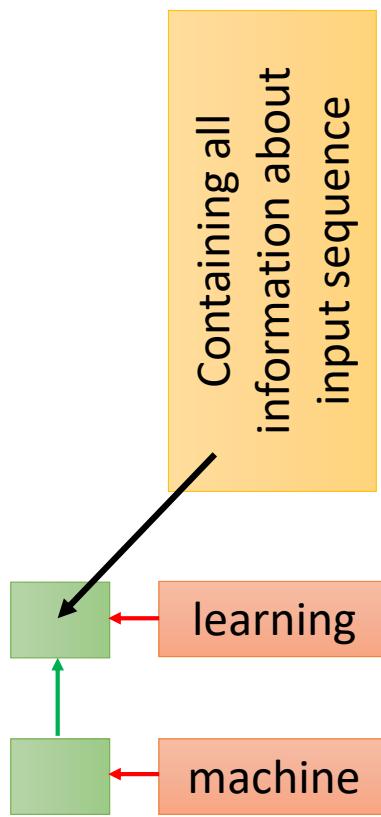
- CTC: example



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.

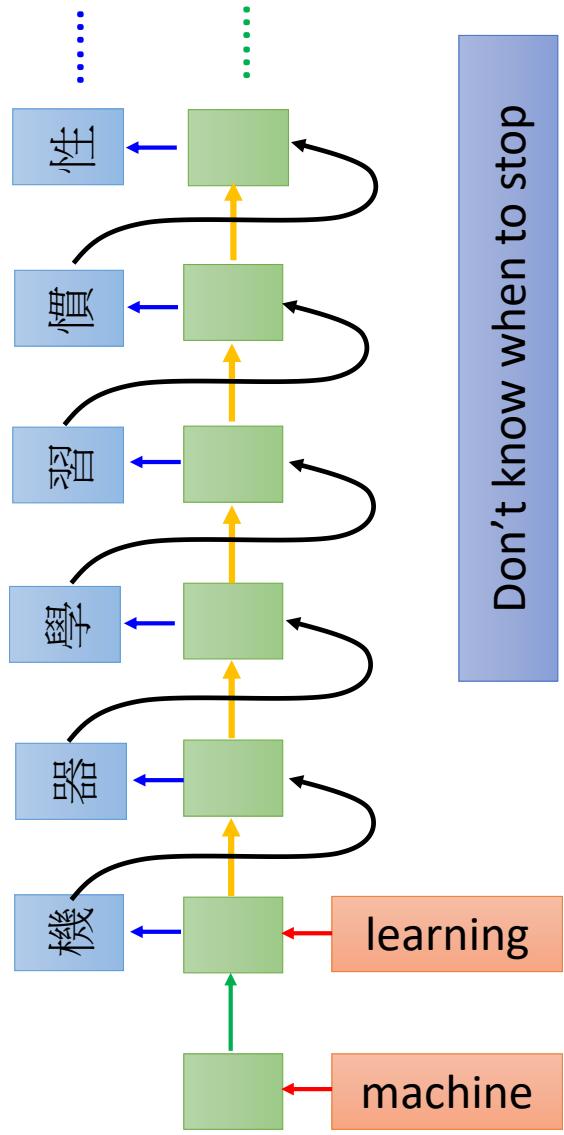
Many to Many (No Limitation)

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

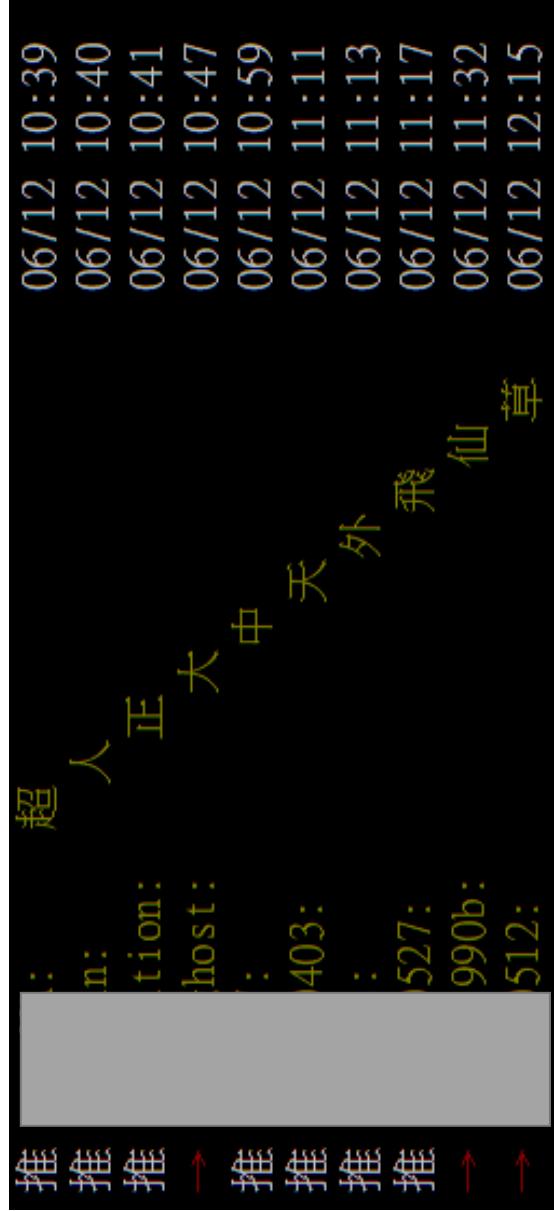


Many to Many (No Limitation)

- 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)

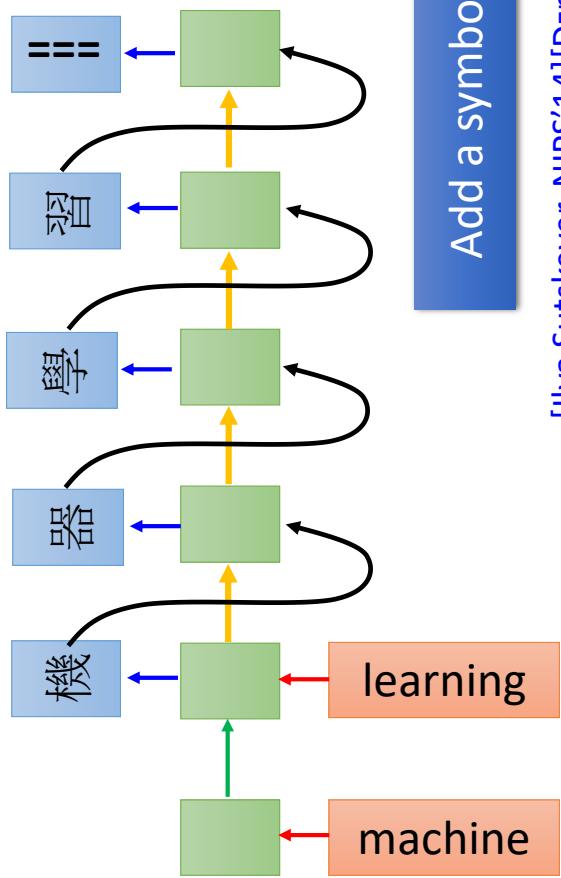


推 ttkagk: ======

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)

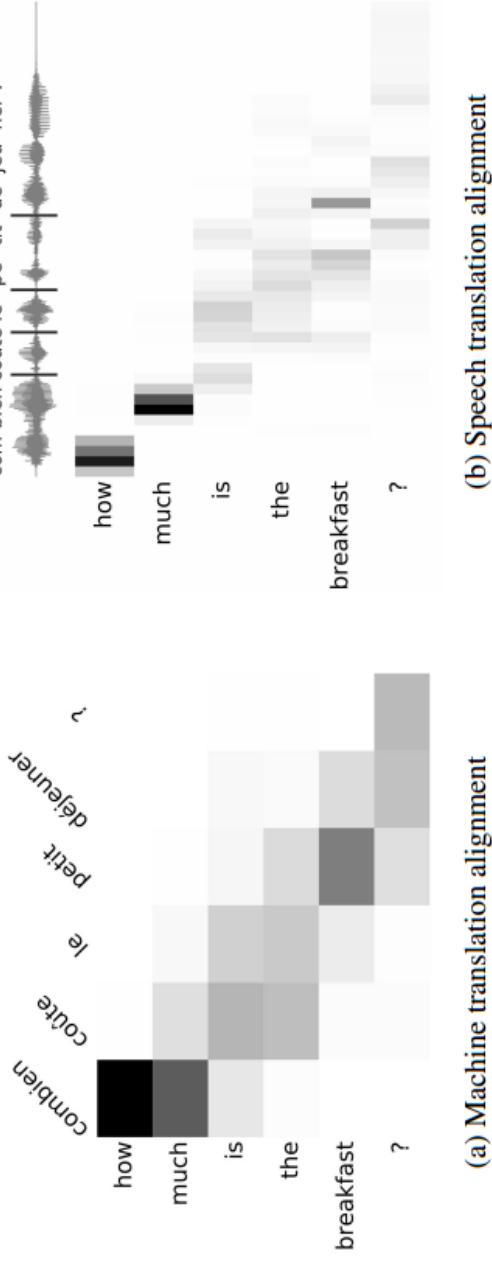
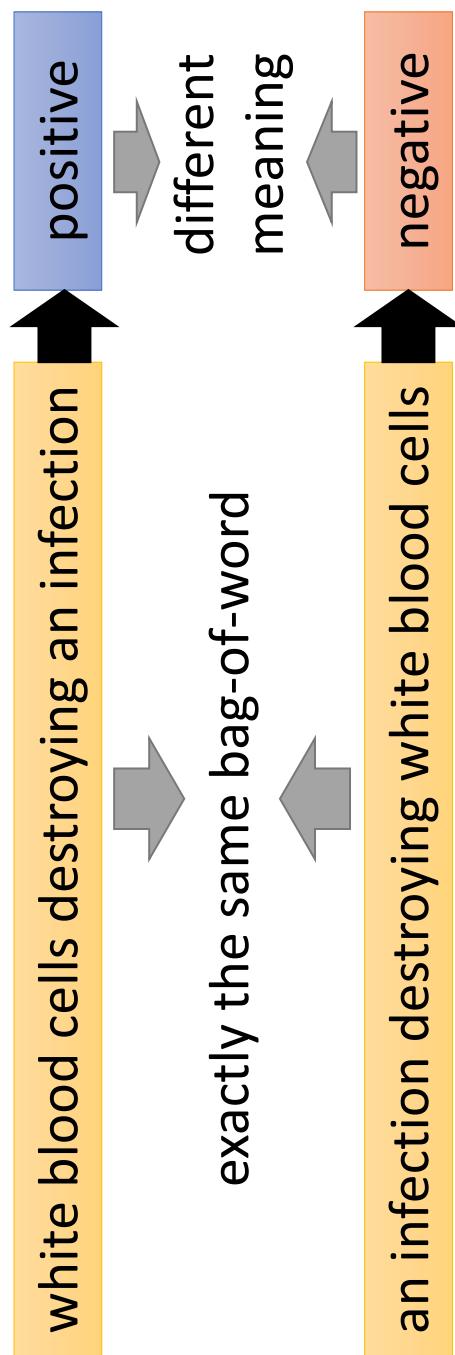


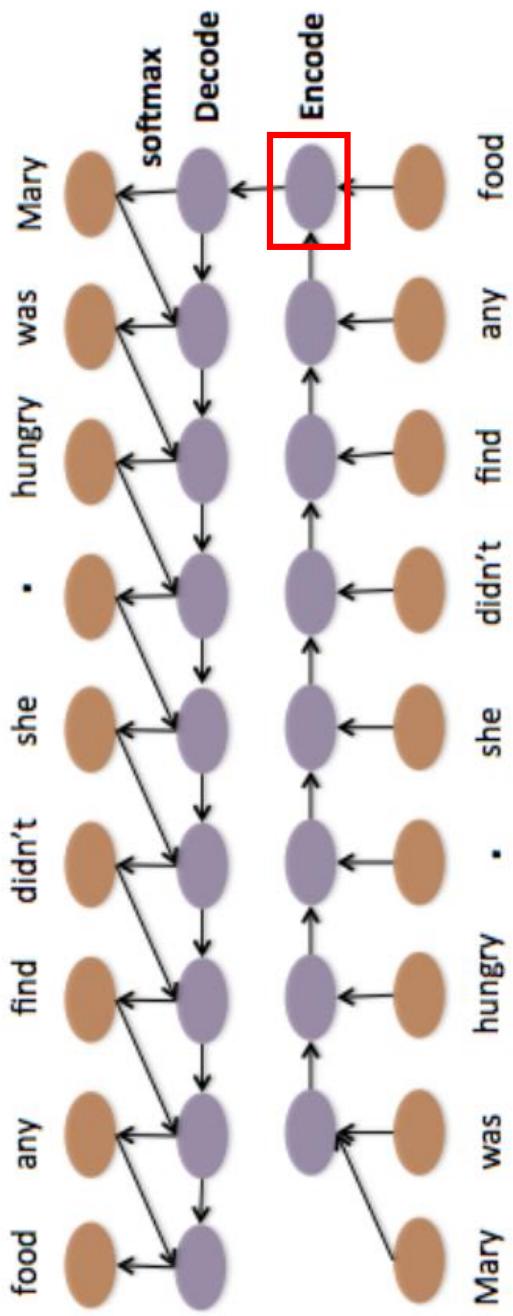
Figure 1: Alignments performed by the attention model during training

Sequence-to-sequence Auto-encoder - Text

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

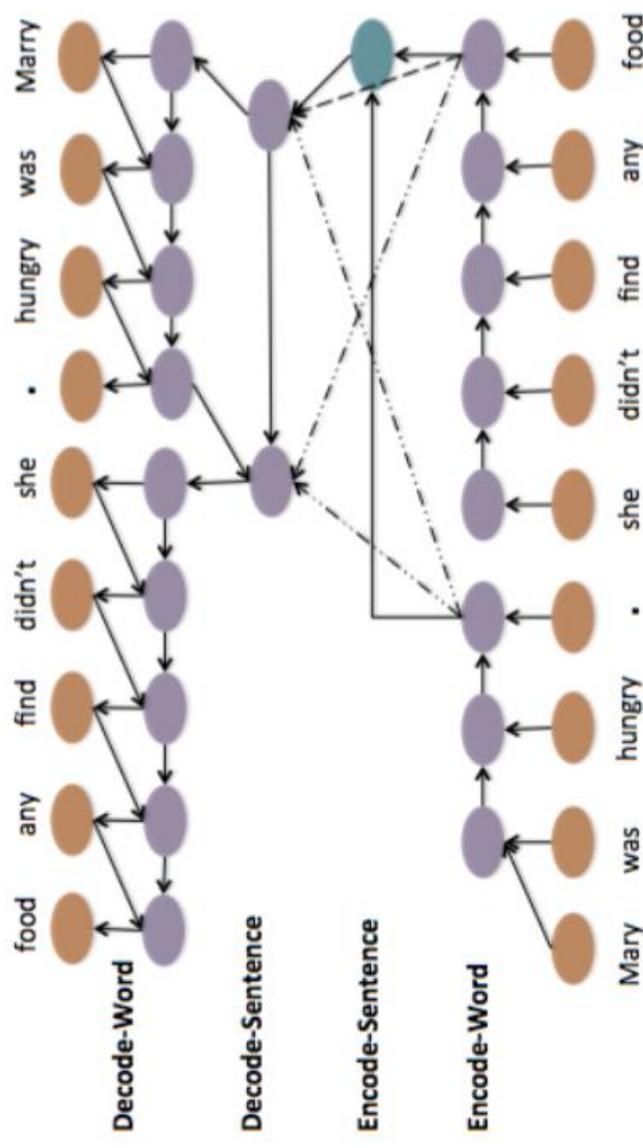


Sequence-to-sequence Auto-encoder - Text



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).

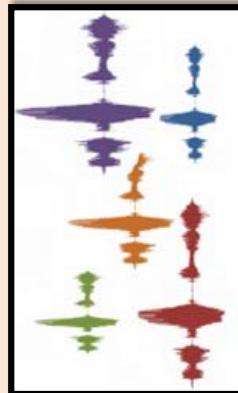
Sequence-to-sequence Auto-encoder - Speech

Audio archive divided into variable-length audio segments

Off-line



Audio
Segment to
Vector



Audio
Segment to
Vector



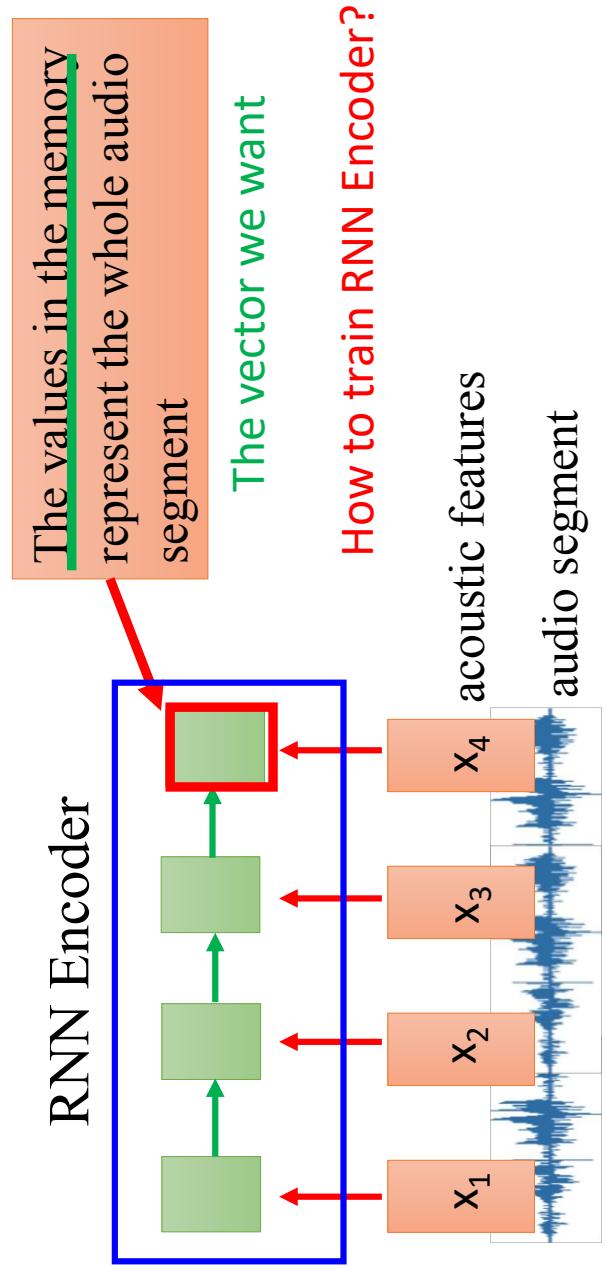
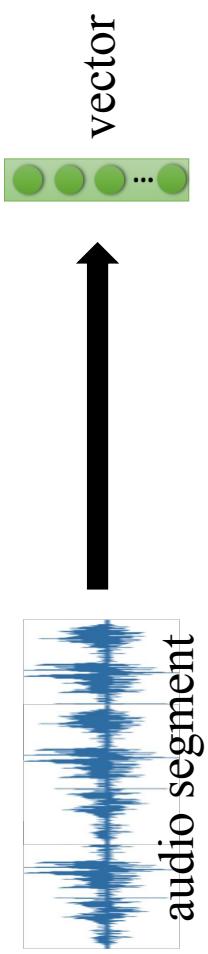
Similarity

Search Result

On-line

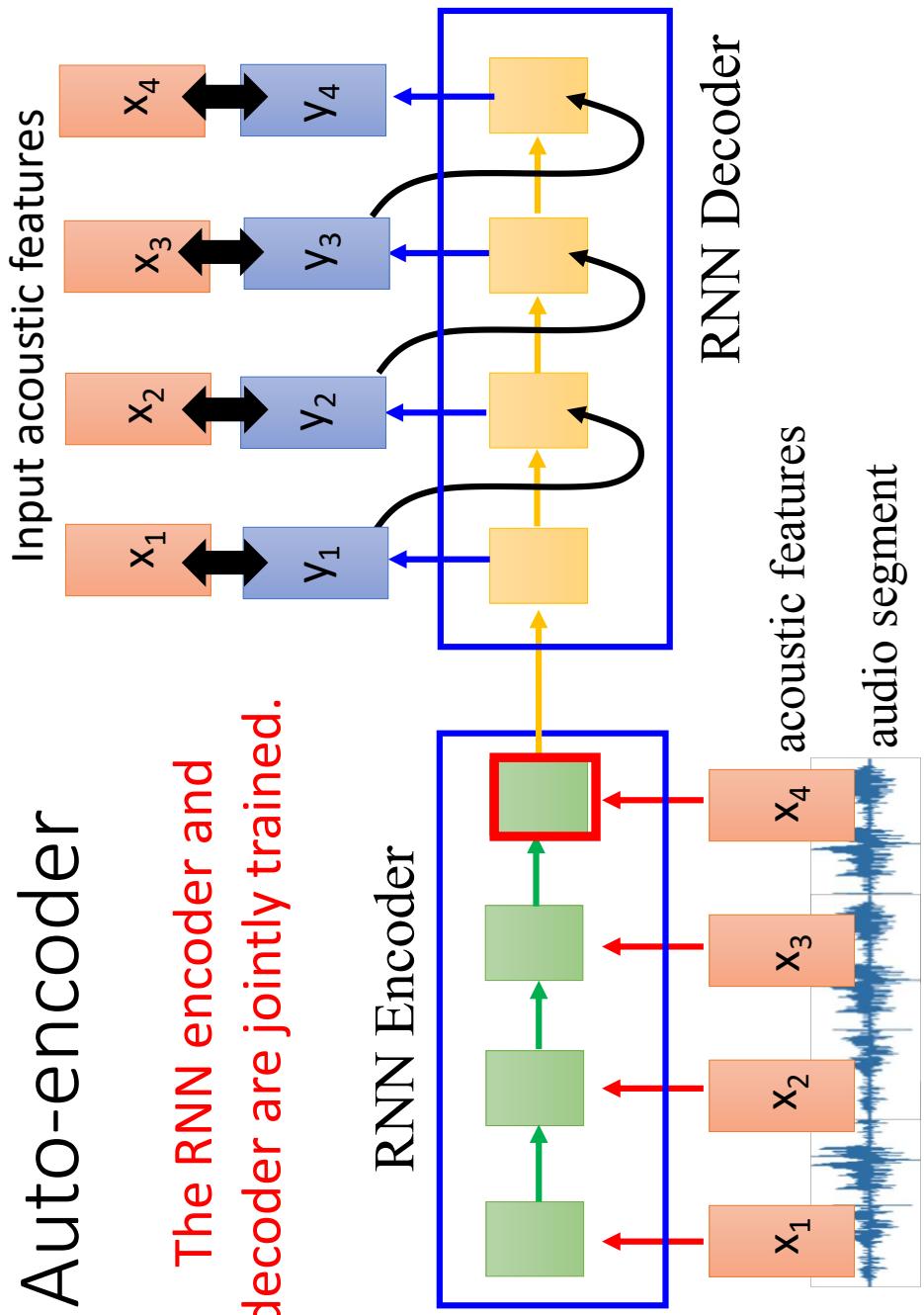
Spoken
Query

Sequence-to-sequence Auto-encoder - Speech



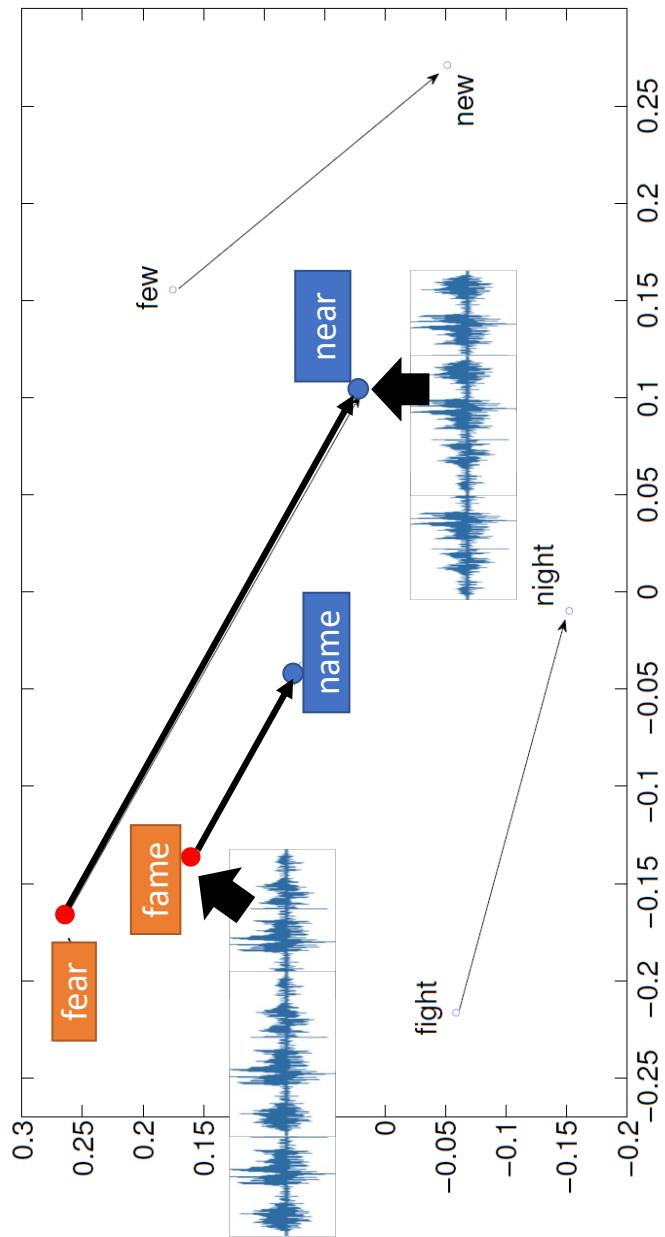
Sequence-to-sequence Auto-encoder

The RNN encoder and
decoder are jointly trained.

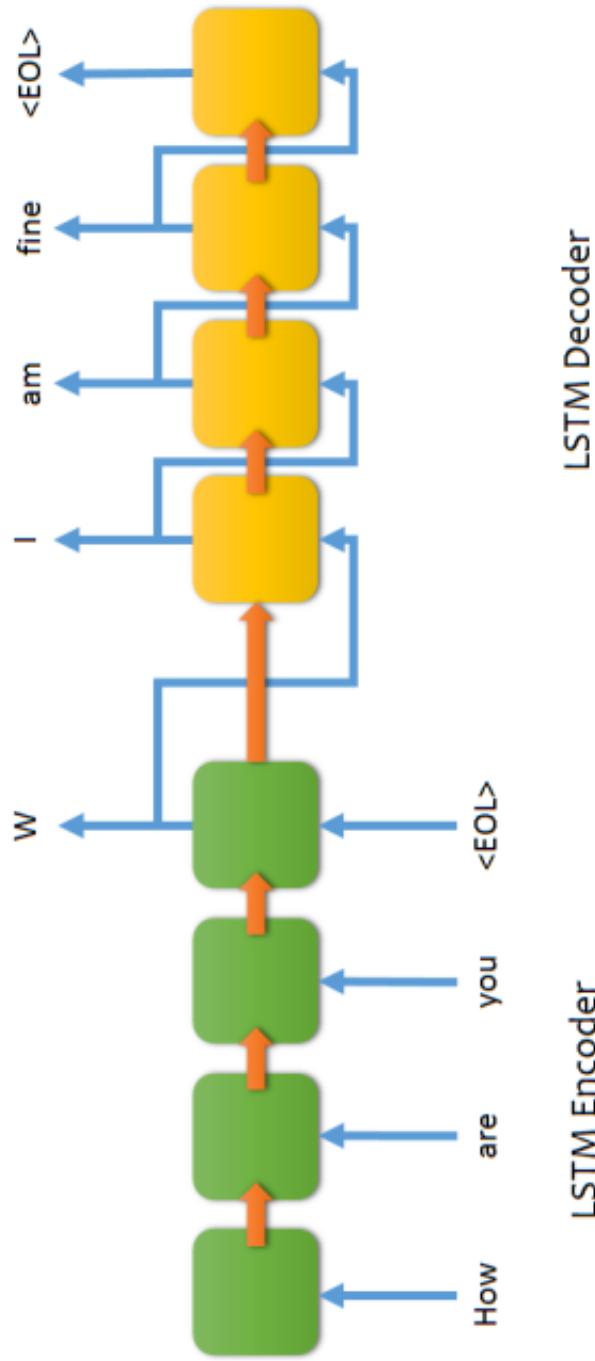


Sequence-to-sequence Auto-encoder - Speech

- Visualizing embedding vectors of the words

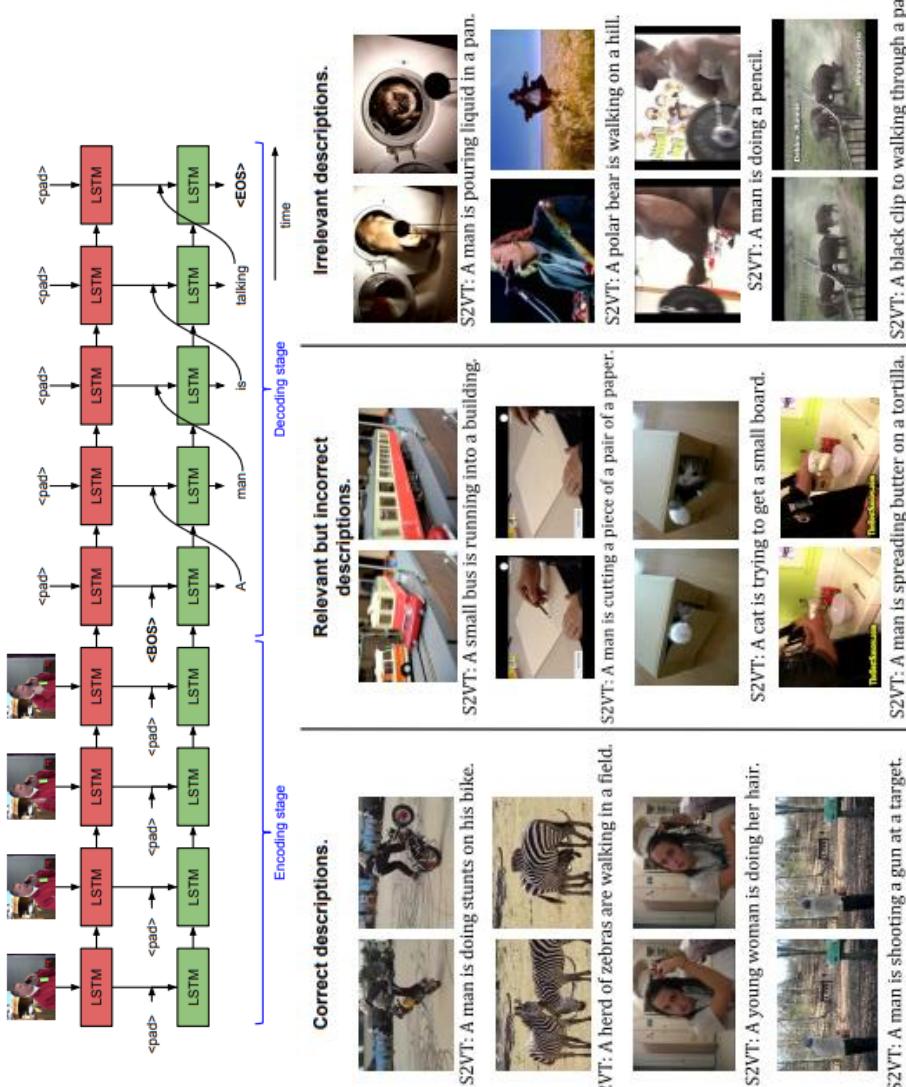


Demo: Chat-bot



電視影集 (~40,000 sentences)、美國總統大選辯論

Video Caption Generation

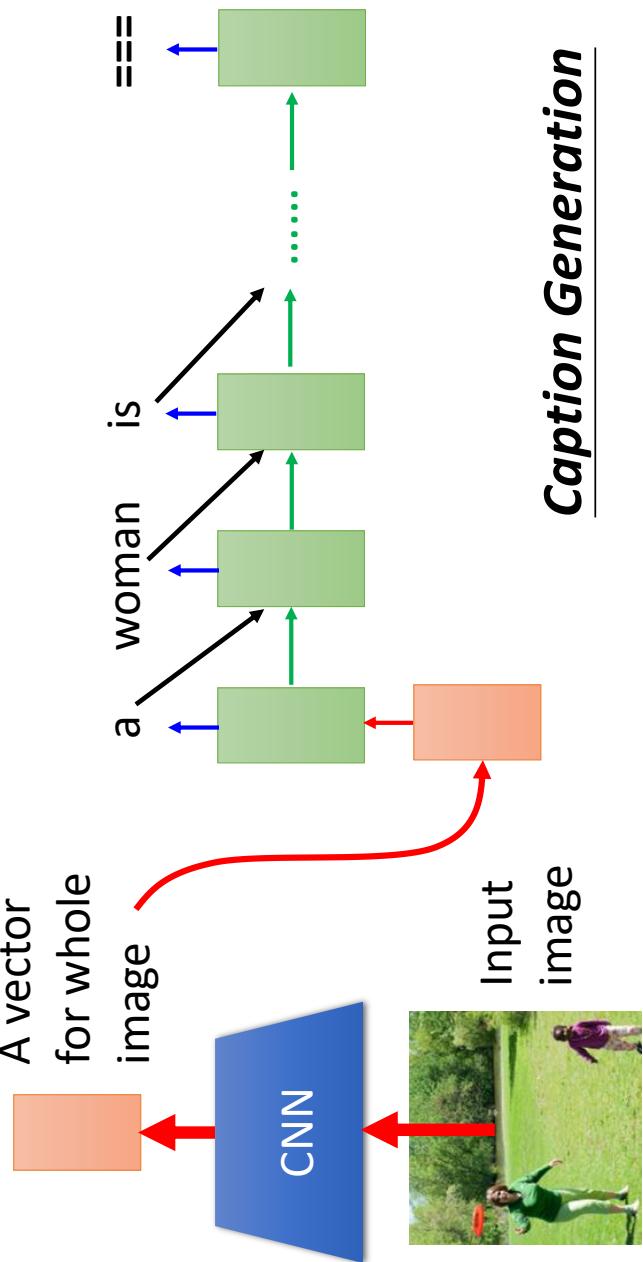


Subhashini Venugopalan, Marcus Rohrbach, Jeffrey Donahue, Raymond Mooney, Trevor Darrell, and Kate Saenko. 2015. Sequence to Sequence -- Video to Text. In Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV) (ICCV '15). IEEE Computer Society, Washington, DC, USA, 4534-4542.

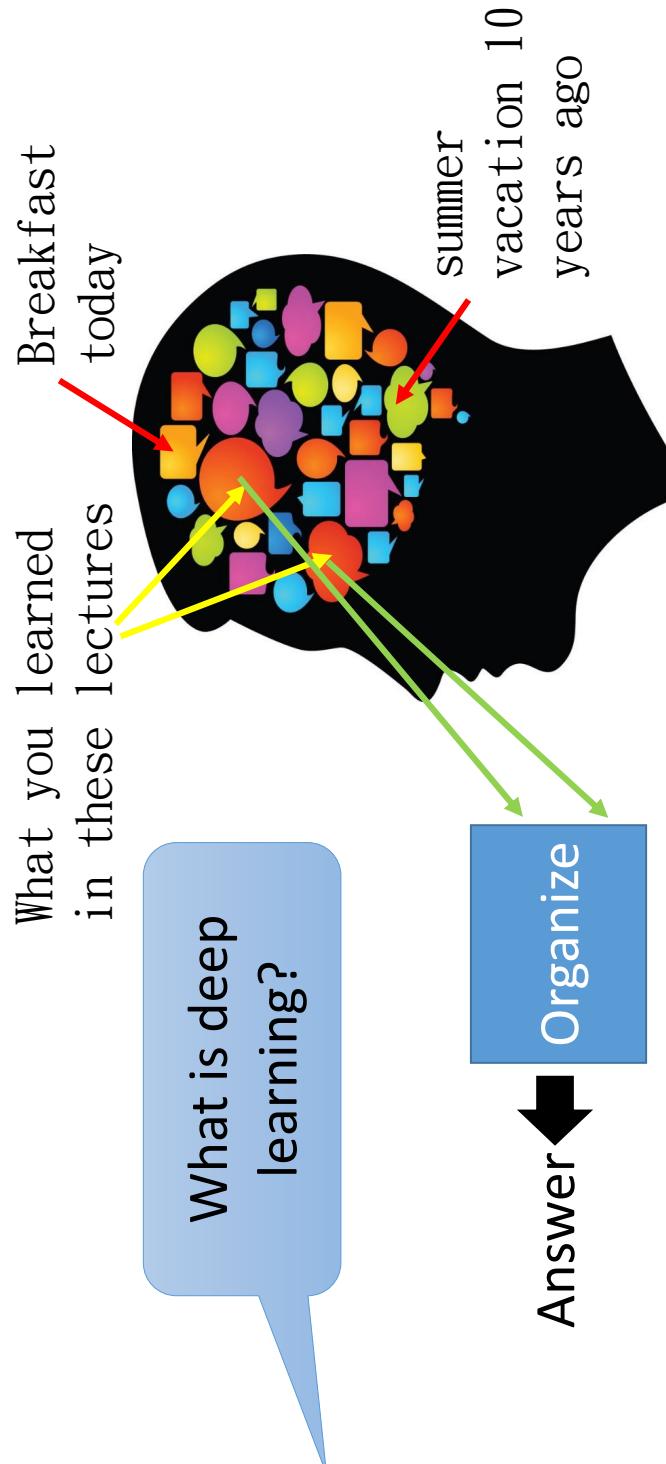
Demo: Image Caption Generation

- Input an image, but output a sequence of words

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

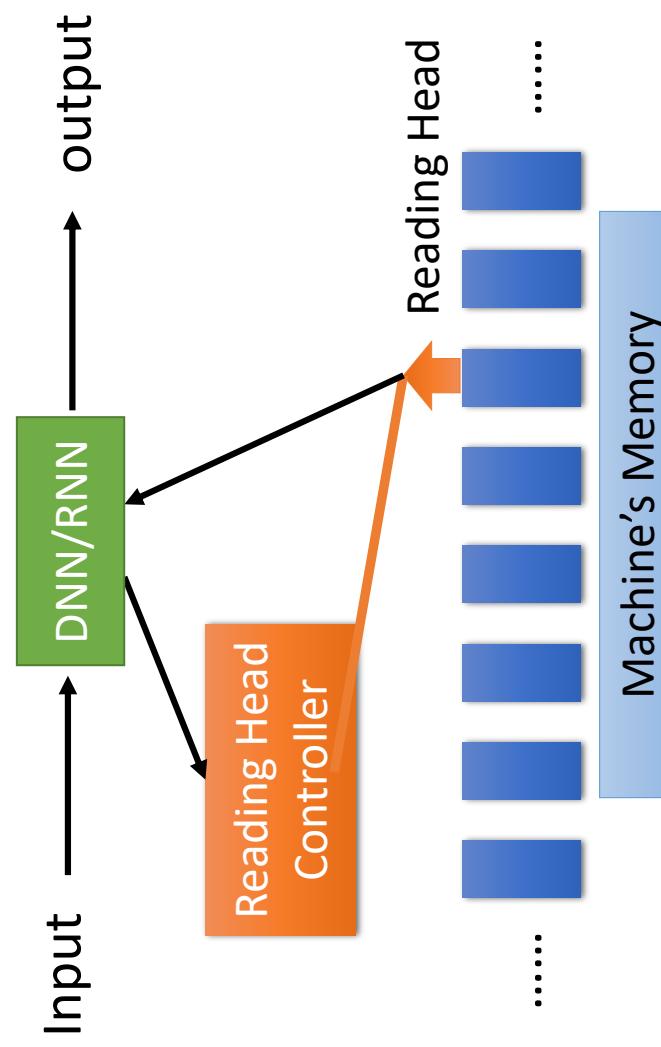


Attention-based Model



http://henrylo1605.blogspot.tw/2015/05/blog-post_56.html

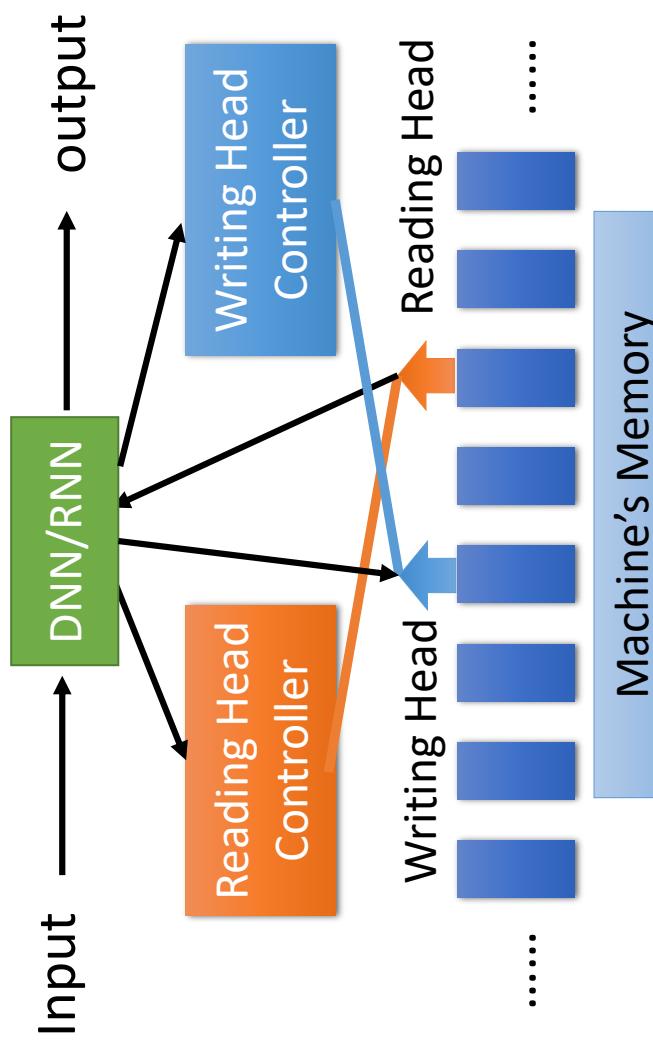
Attention-based Model



Ref:

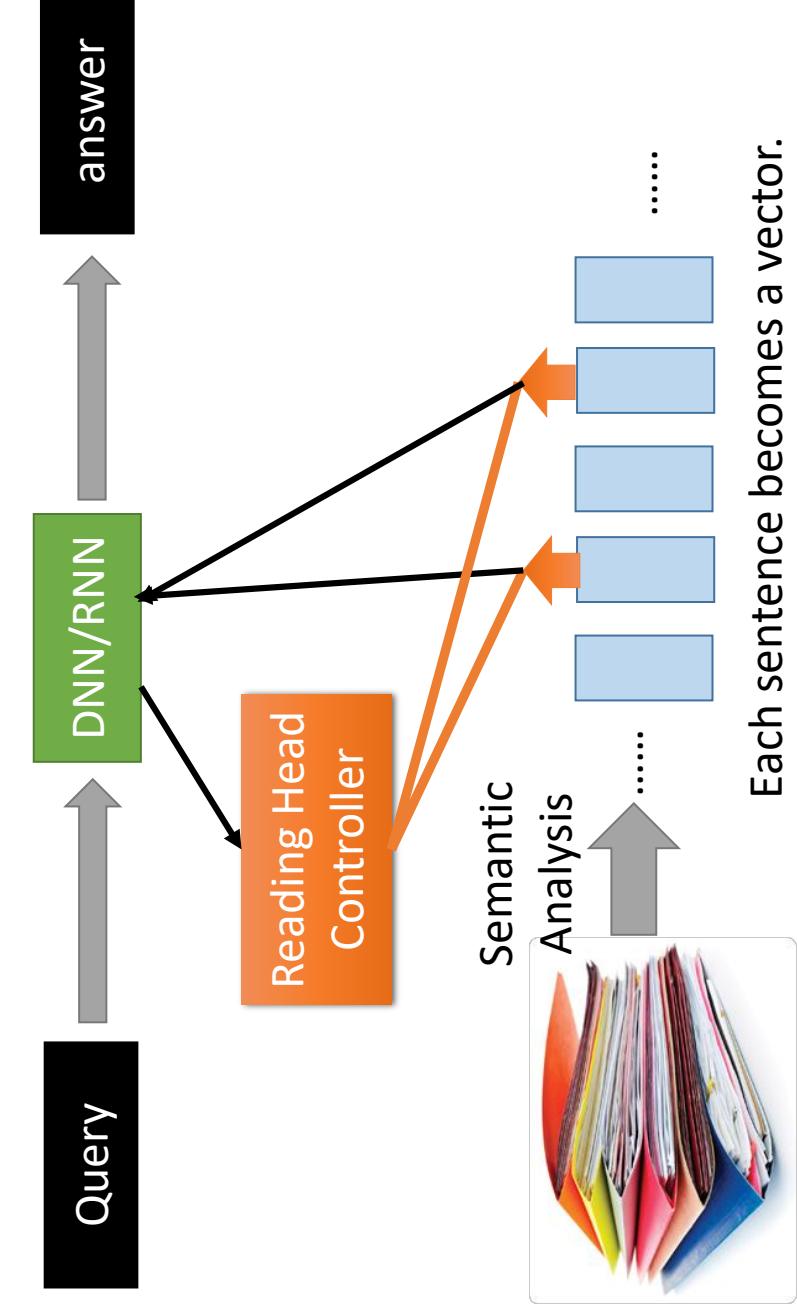
[http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS_2015_2/Lecture/Attain%20\(v3\).e cm.mp4/index.html](http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS_2015_2/Lecture/Attain%20(v3).e cm.mp4/index.html)

Attention-based Model v2

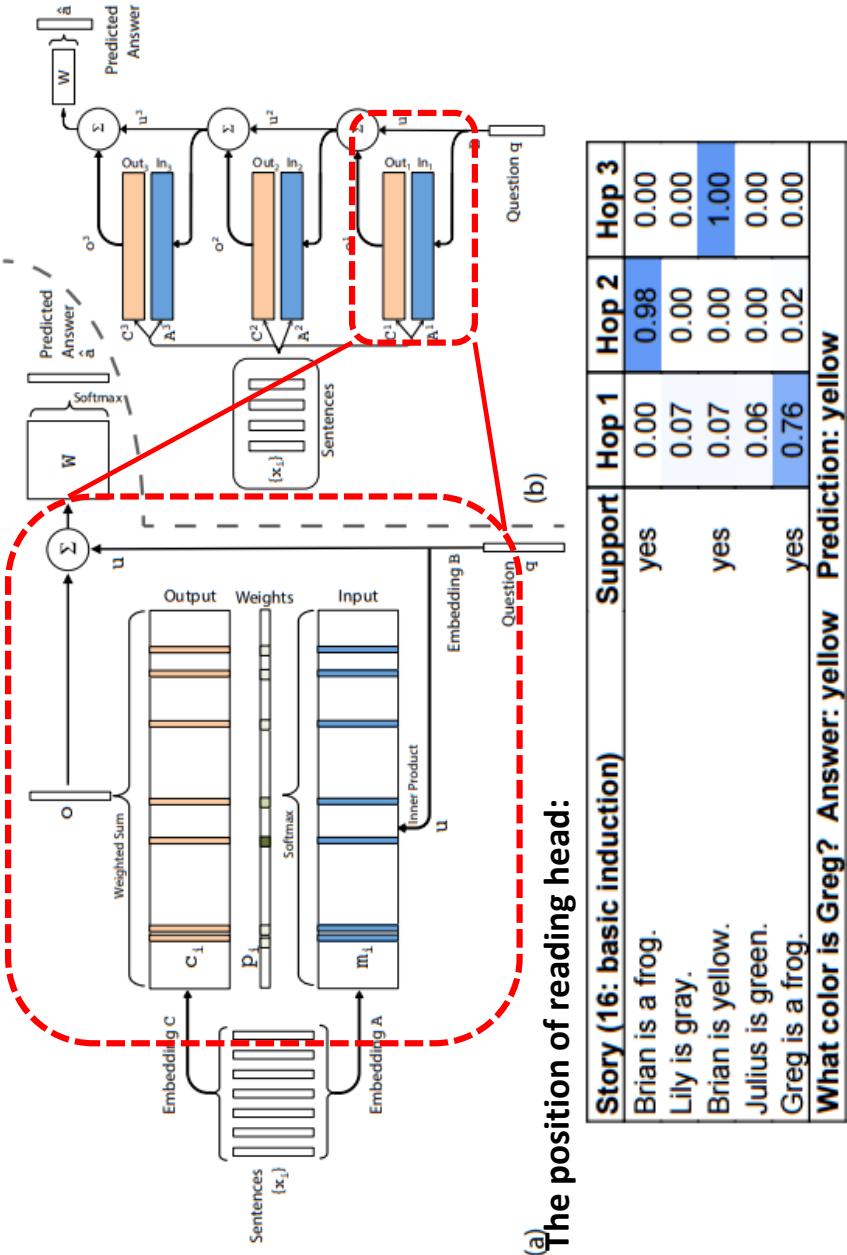


Neural Turing Machine

Reading Comprehension

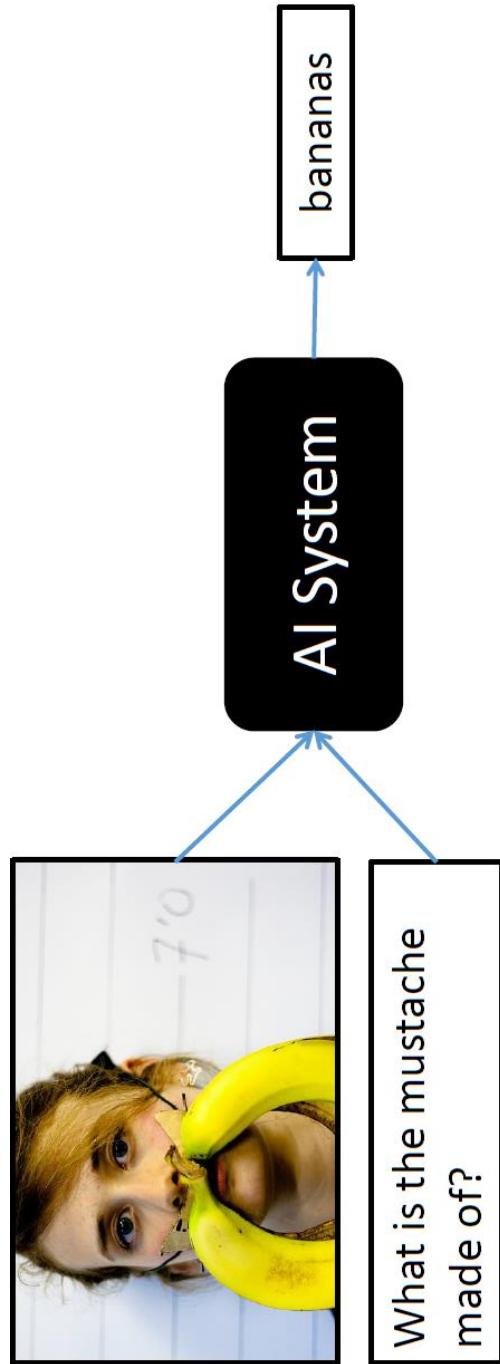


Reading Comprehension



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

Visual Question Answering



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

Visual Question Answering

