# Recurrent Neural Networks
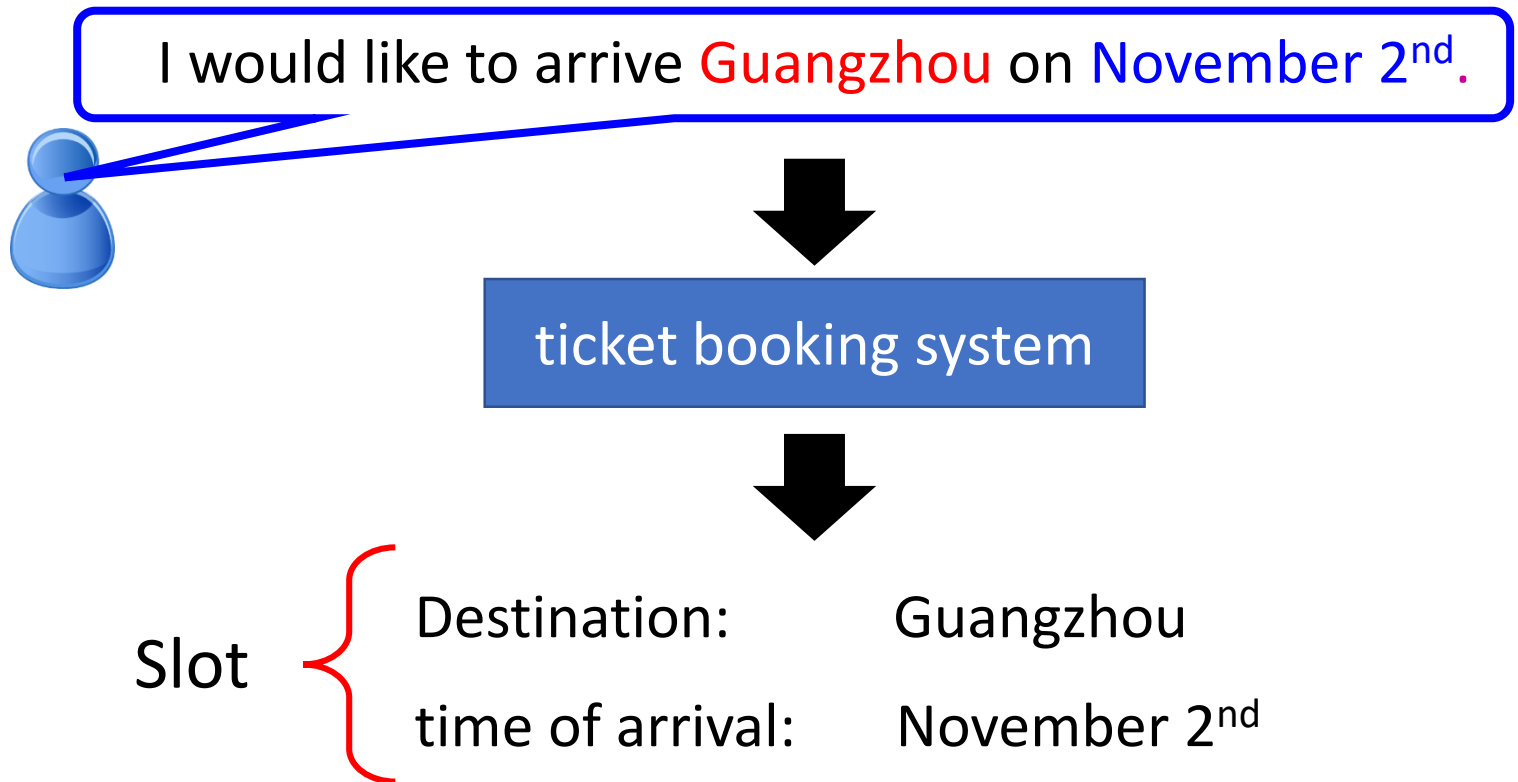
Hankui Zhuo

May 22, 2019

http://xplan-lab.org/IAA-Course

# Example Application

- Slot Filling

I would like to arrive Guangzhou on November 2$^{nd}$.

ticket booking system
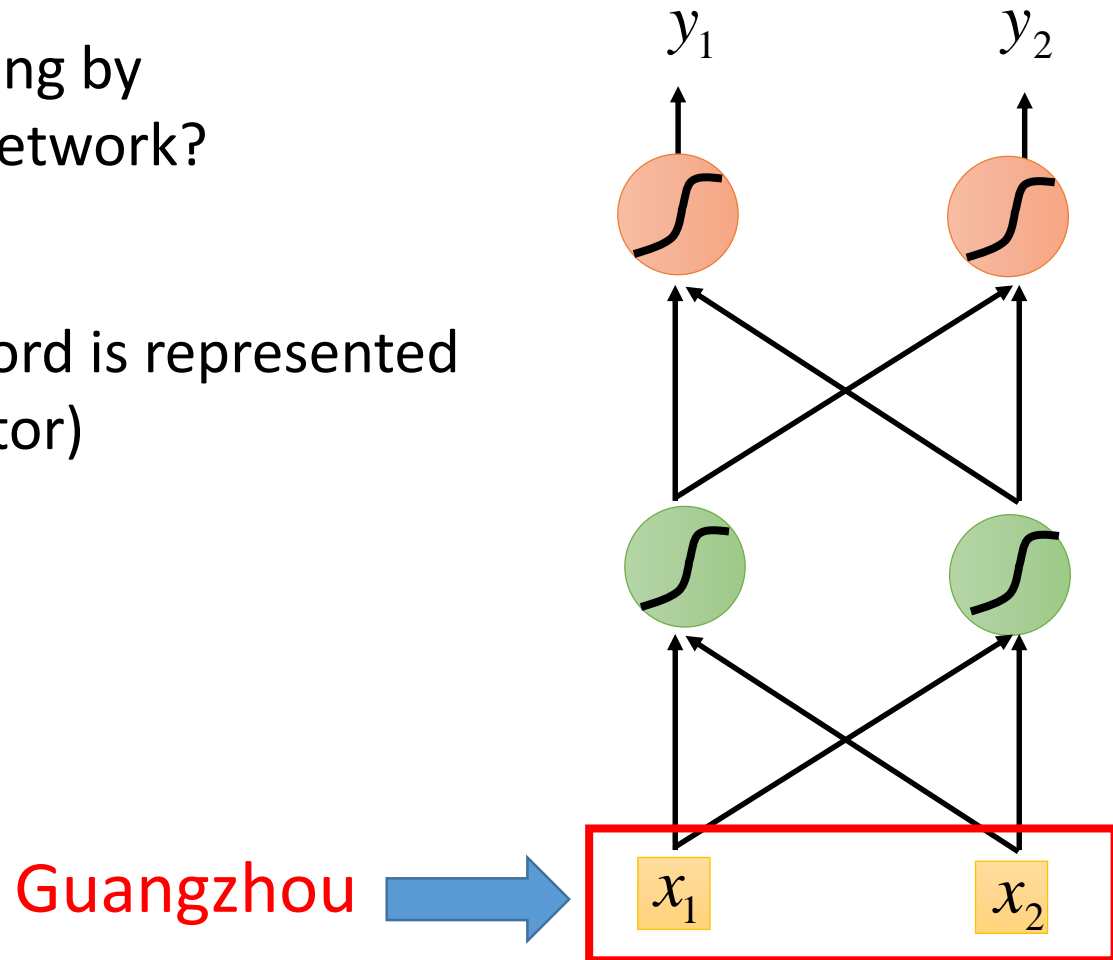
Slot

Destination:        Guangzhou

time of arrival:        November 2$^{nd}$

# Example Application

Solving slot filling by Feedforward network?

Input: a word

   (Each word is represented as a vector)

$y_1$    $y_2$

Guangzhou

$x_1$    $x_2$

# 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

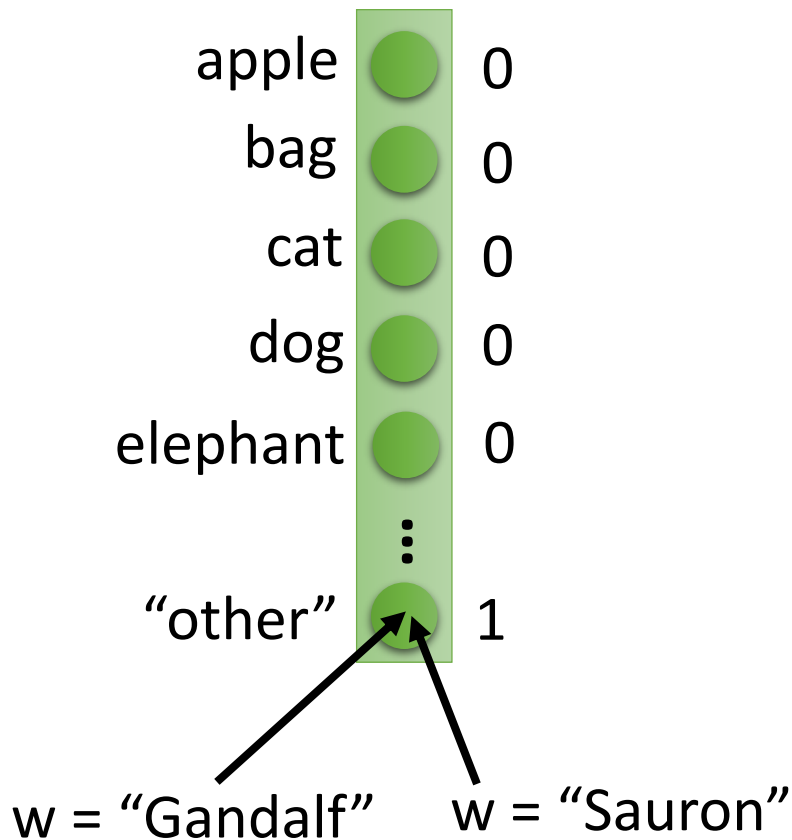apple = [ 1  0  0  0  0]

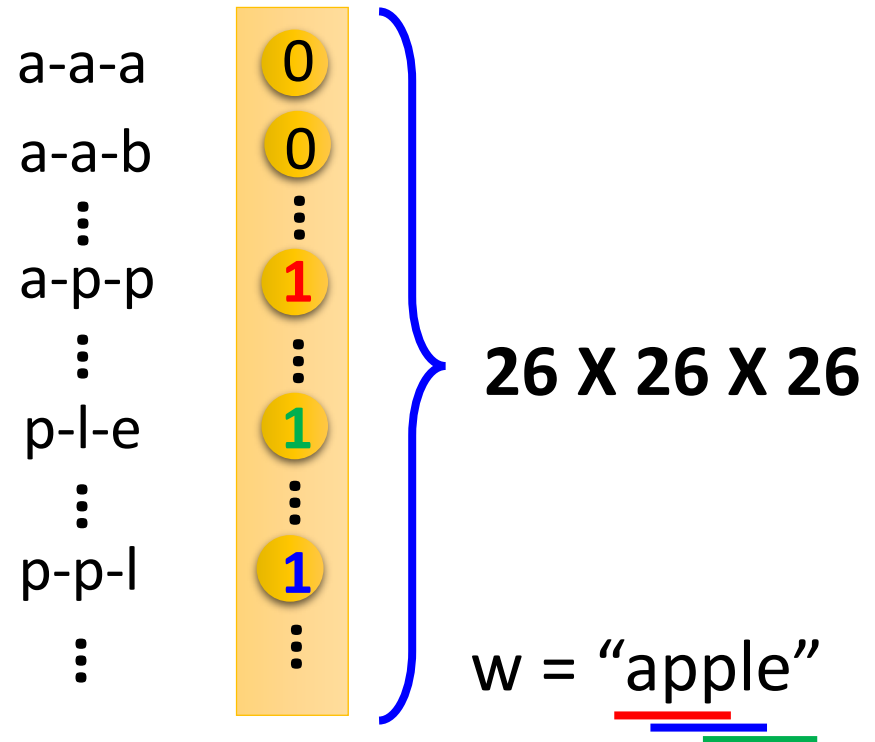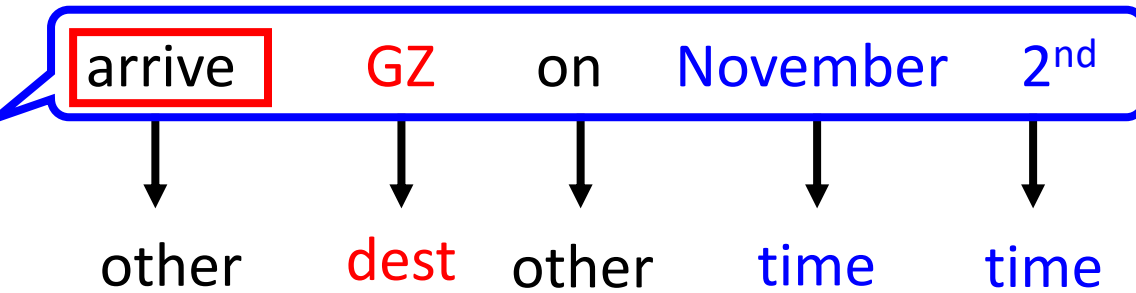bag   = [ 0  1  0  0  0]

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"*

| | |
|---|---|
| apple | 0 |
| bag | 0 |
| cat | 0 |
| dog | 0 |
| elephant | 0 |
| ⋮ | |
| "other" | 1 |

w = "Gandalf"     w = "Sauron"

## *Word hashing*

| | |
|---|---|
| a-a-a | 0 |
| a-a-b | 0 |
| ⋮ | ⋮ |
| a-p-p | 1 |
| ⋮ | ⋮ |
| p-l-e | 1 |
| ⋮ | ⋮ |
| p-p-l | 1 |
| ⋮ | ⋮ |

**26 X 26 X 26**

w = "apple"

# Example Application



dest    time of departure

Solving slot filling by Feedforward network?

Input: a word

(Each word is represented as a vector)

Output:

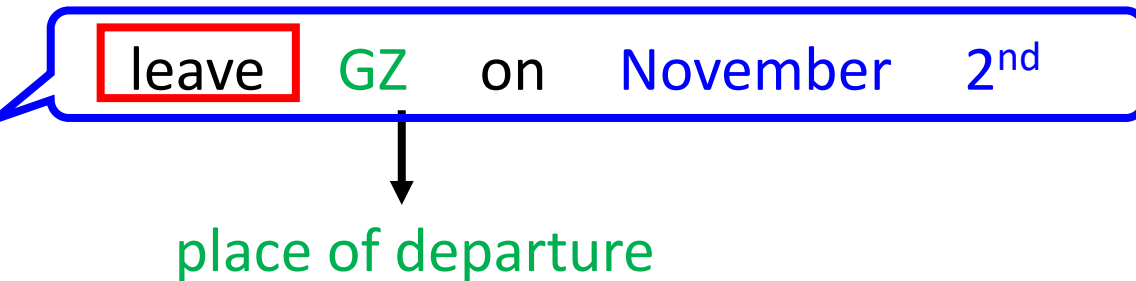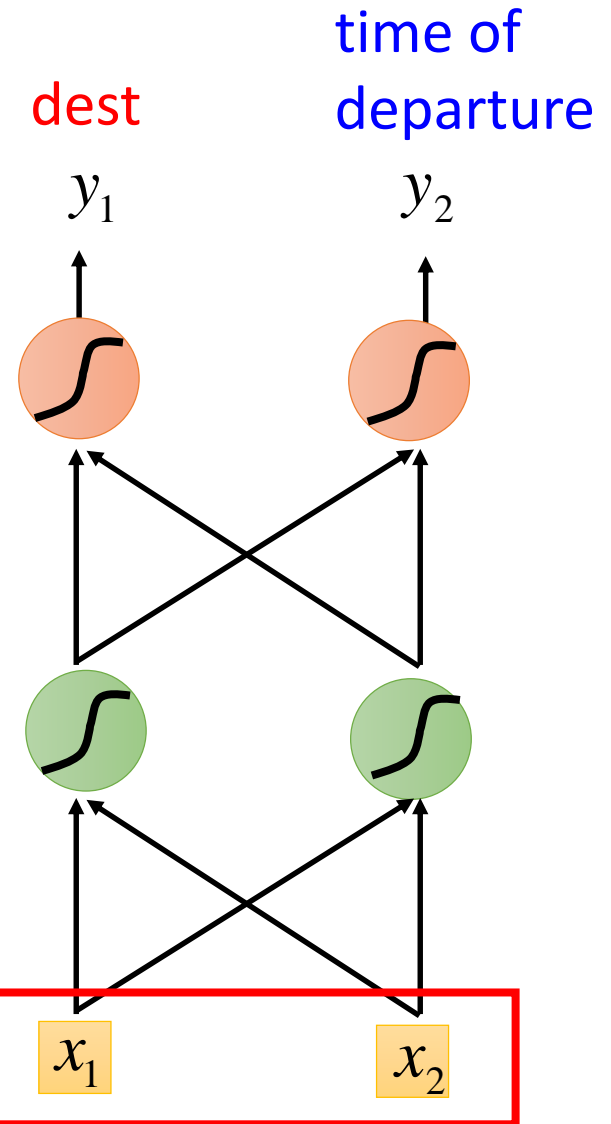Probability distribution that the input word belonging to the slots

# Example Application
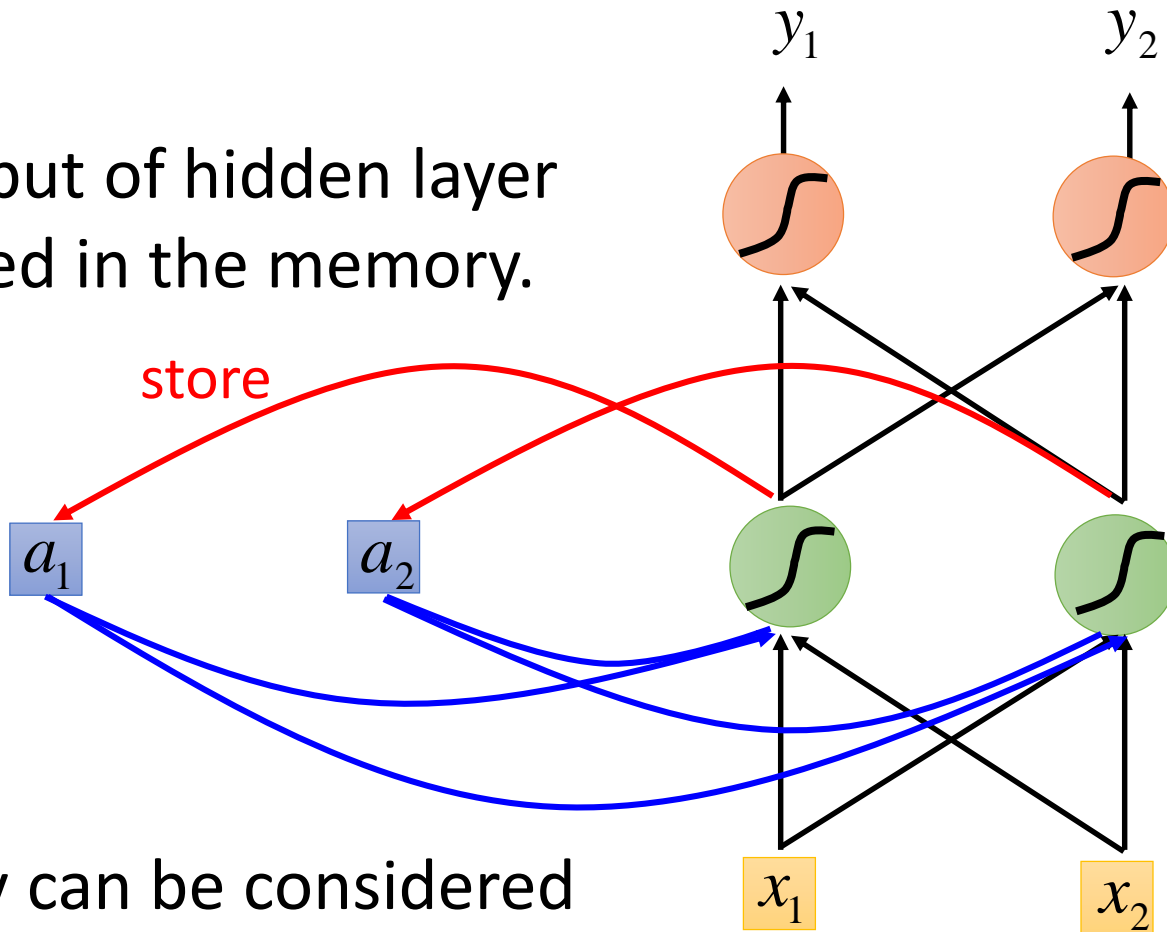
# Recurrent Neural Networks (RNNs)

The output of hidden layer are stored in the memory.
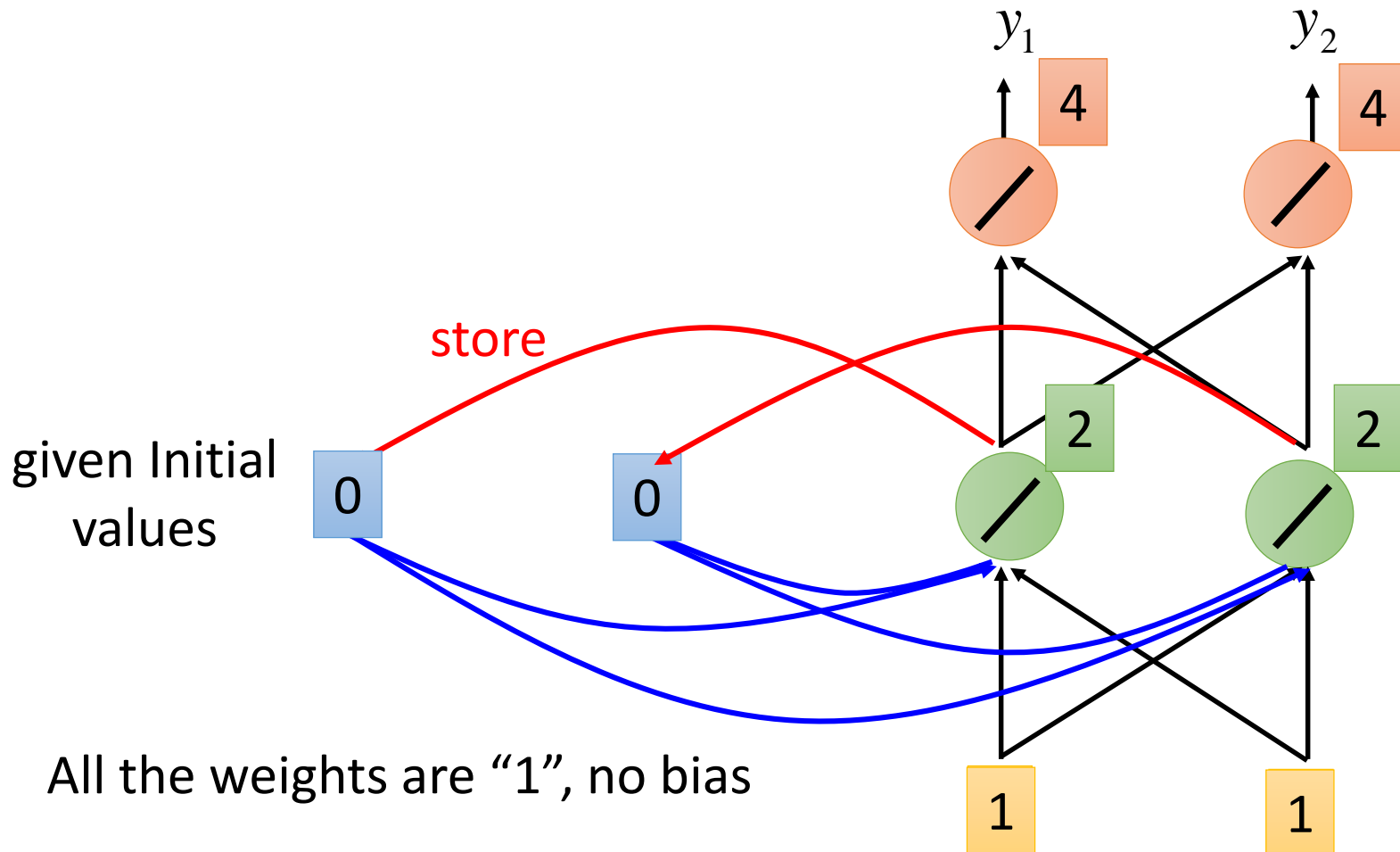
store

Memory can be considered as another input.

# Example

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

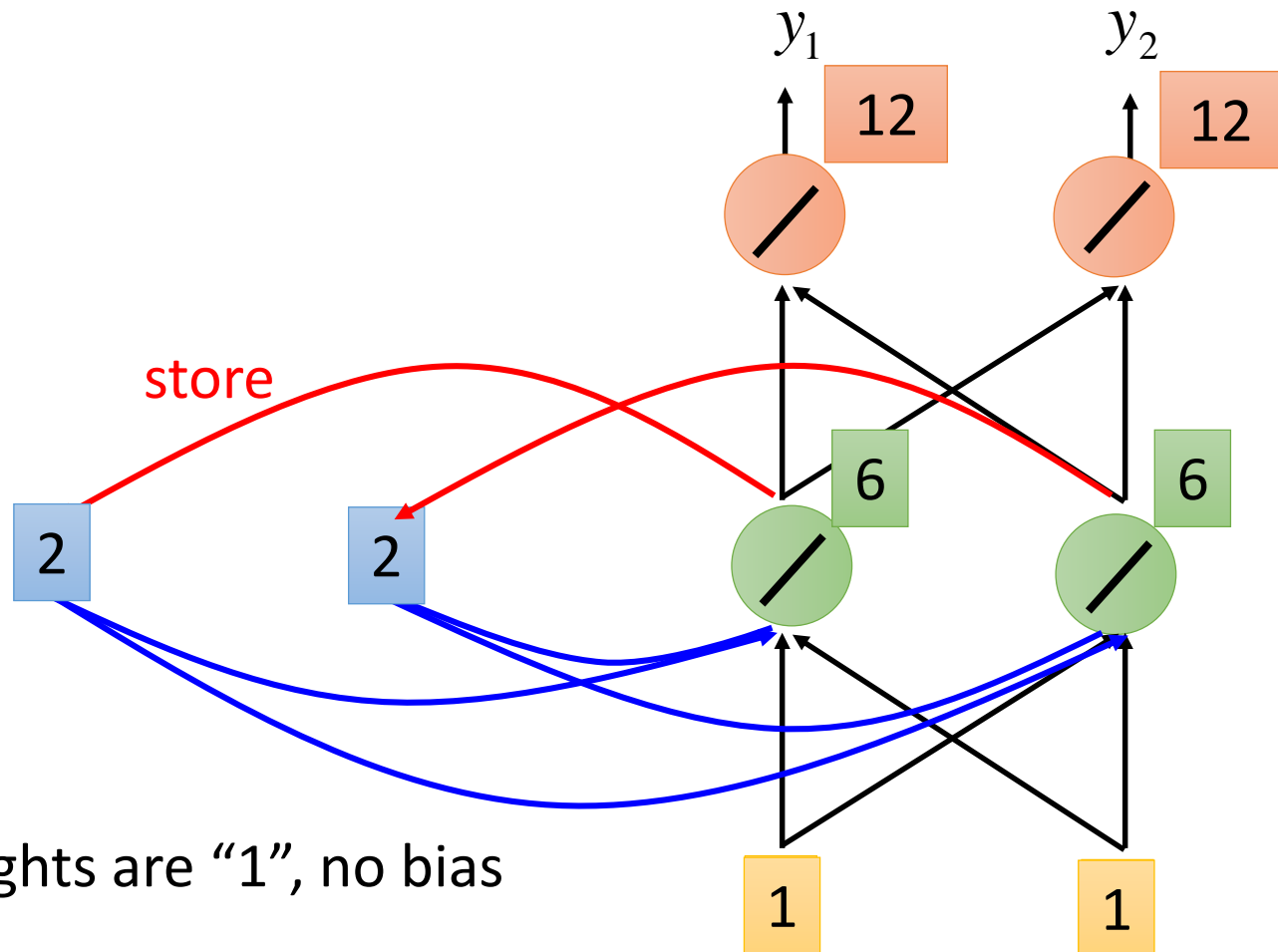output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix}$



$y_1$  4

$y_2$  4

store

given Initial values   0   0   2   2
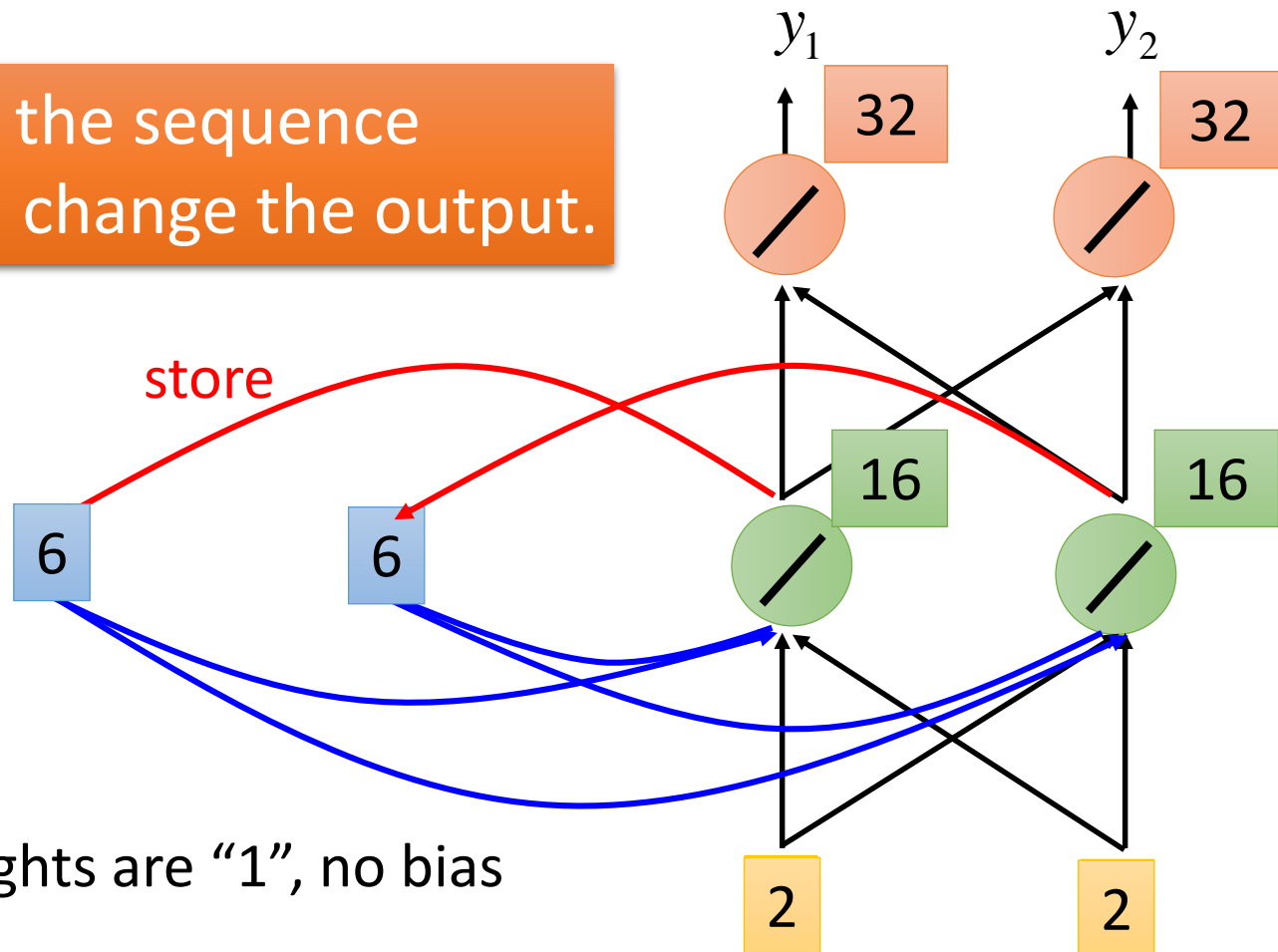
All the weights are "1", no bias

All activation functions are linear

1   1

# Example

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

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



$y_1$  12

$y_2$  12

store

2   2   6   6

All the weights are "1", no bias

1   1

All activation functions are linear

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

# Example

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

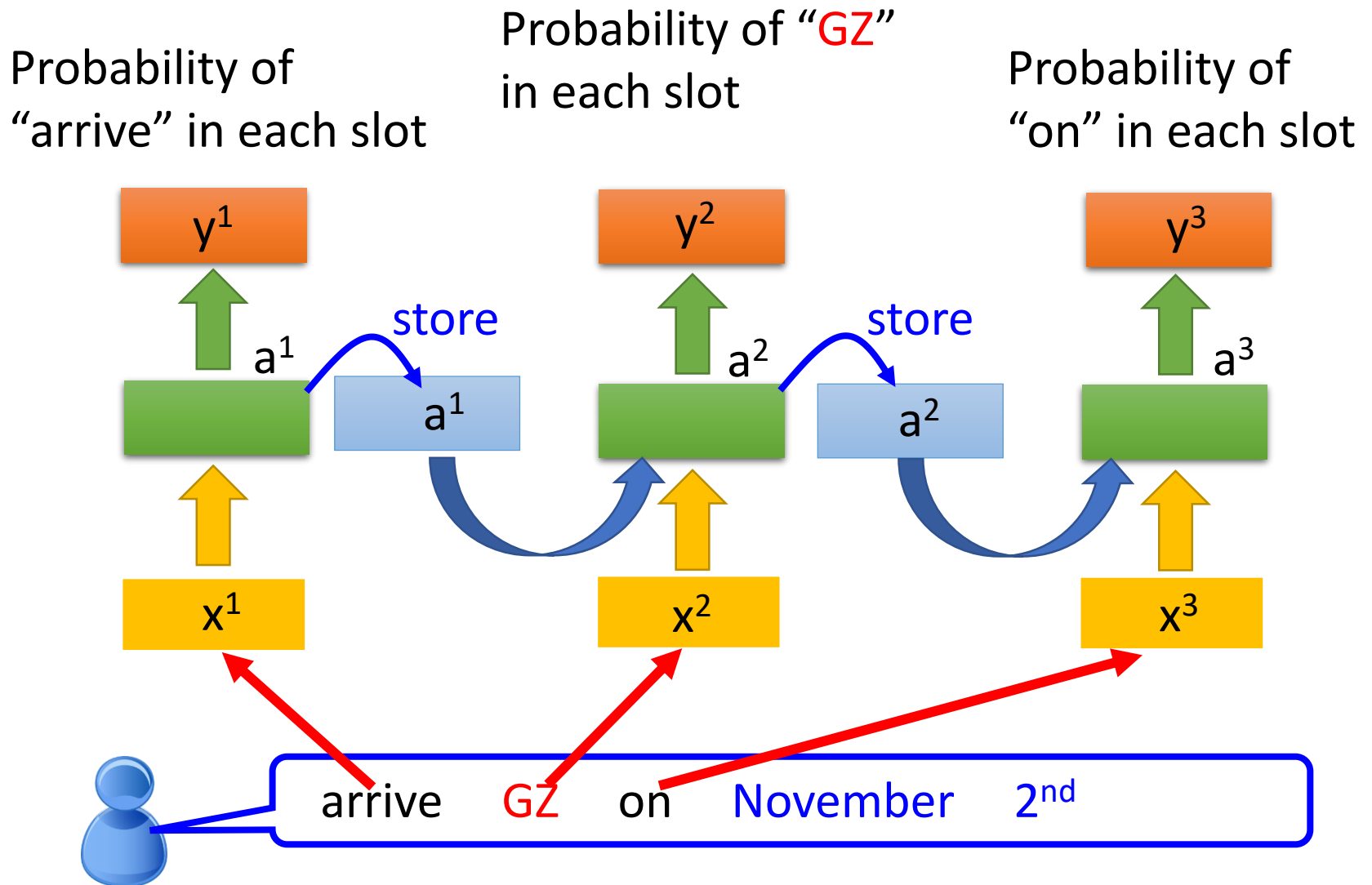$y_1$  $y_2$

Changing the sequence order will change the output.

store

All the weights are "1", no bias

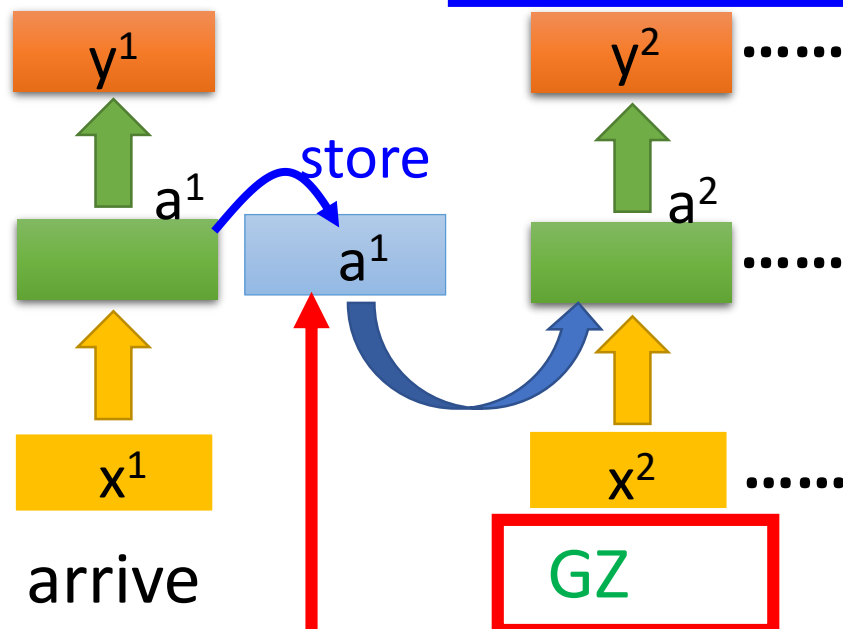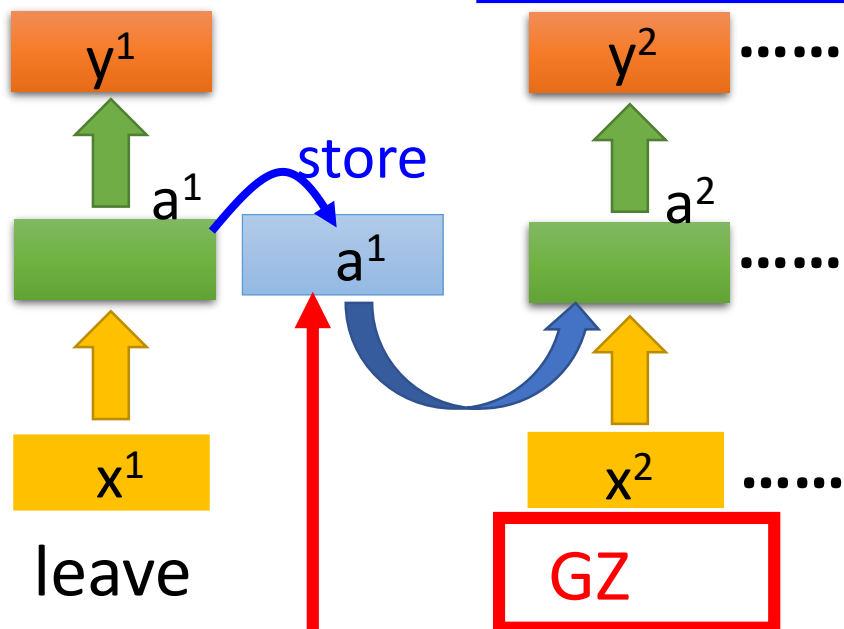All activation functions are linear

# RNN

The same network is used again and again.

Probability of "arrive" in each slot
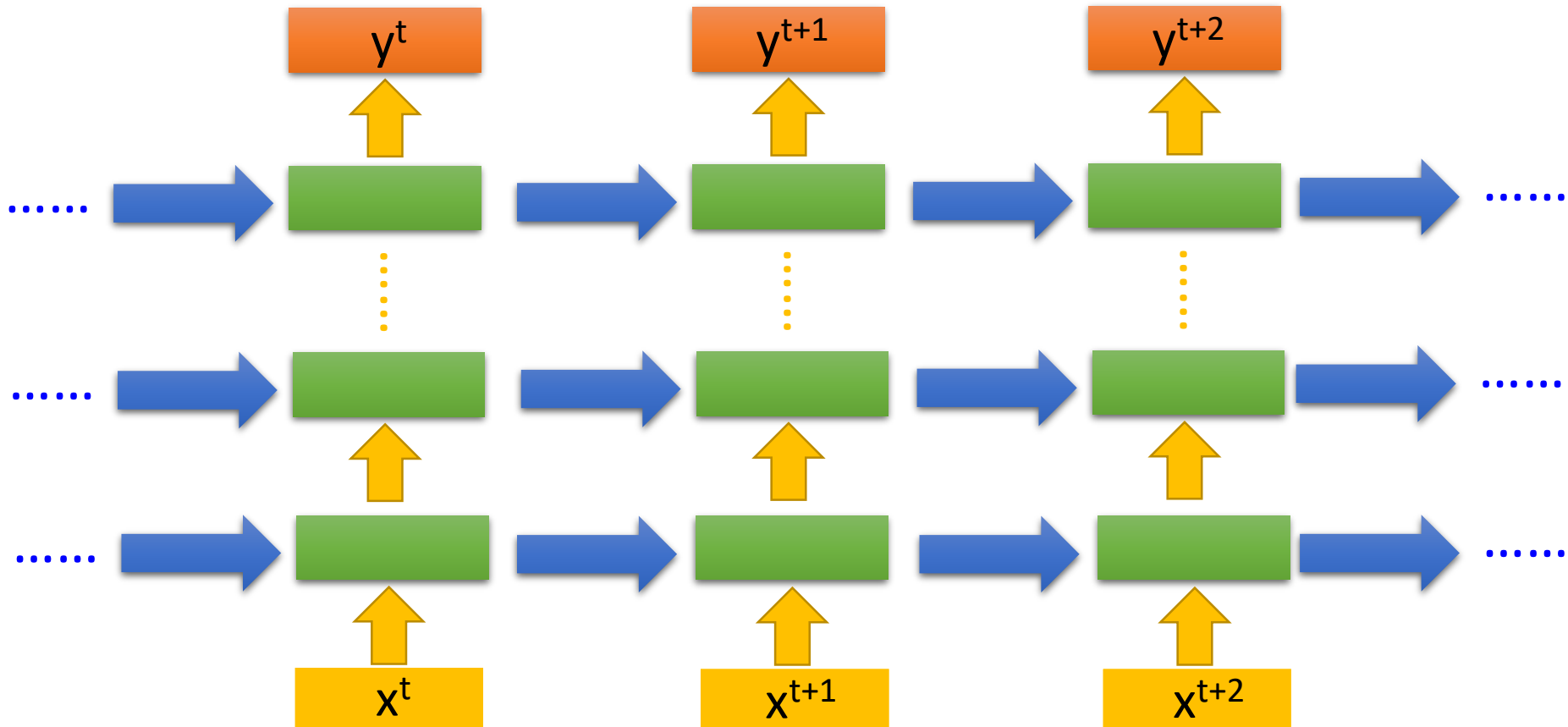
Probability of "GZ" in each slot

Probability of "on" in each slot



$y^1$   $y^2$   $y^3$

store   store

$a^1$   $a^2$   $a^3$

$a^1$   $a^2$

$x^1$   $x^2$   $x^3$

arrive   GZ   on   November   2nd

# Of course it can be deep ...

# Elman Network & Jordan Network

# Bidirectional RNN

# Long Short-term Memory (LSTM)

$$a = h(c')f(z_o)$$

$z_o$ → Output Gate, $f(z_o)$

multiply

$h(c')$

$c$ $f(z_f)$

Forget Gate

$c'$

$z_f$

$cf(z_f)$

Activation function f is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

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

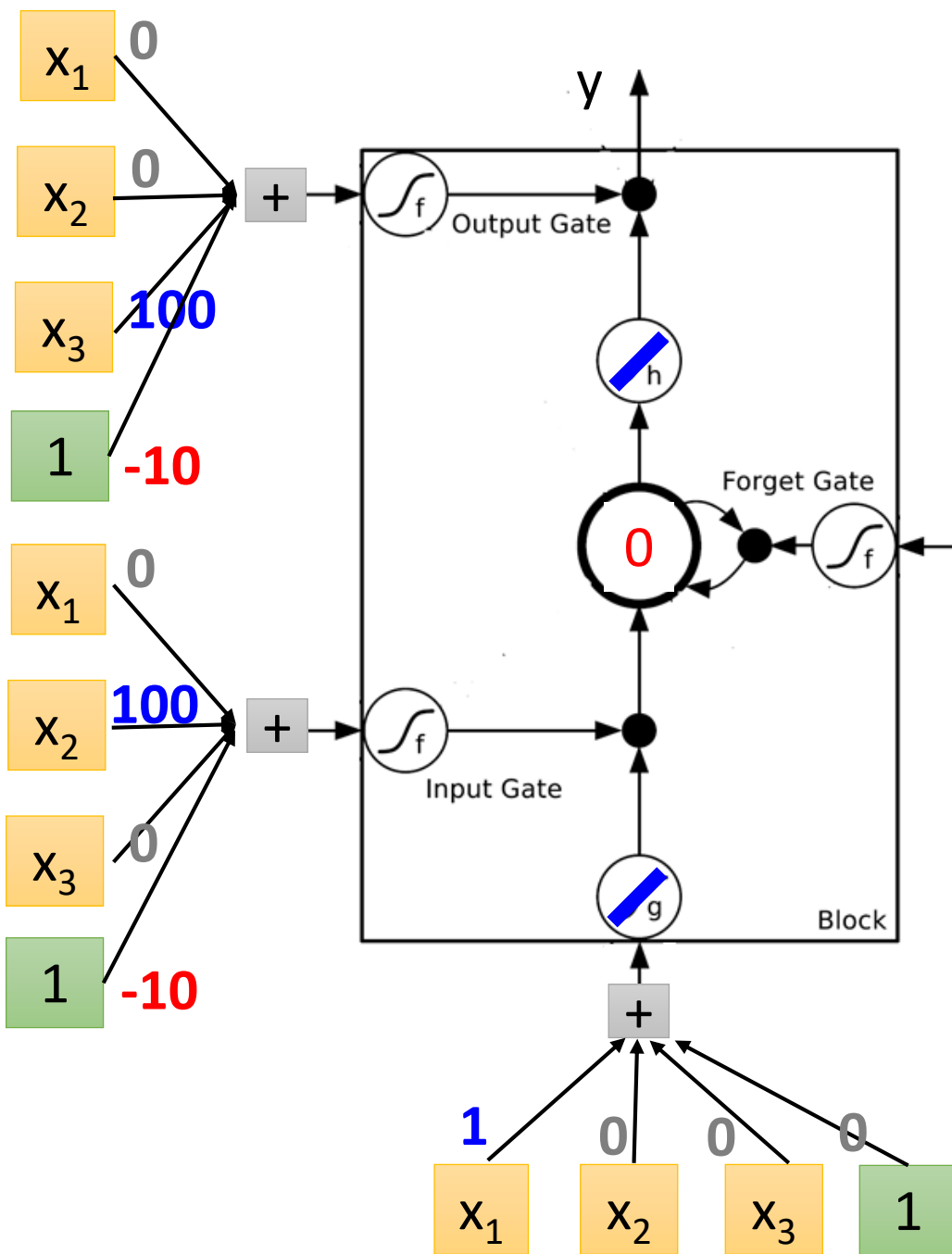$z_i$ → Input Gate, $f(z_i)$ $g(z)f(z_i)$

multiply

$g(z)$

$z$

Block

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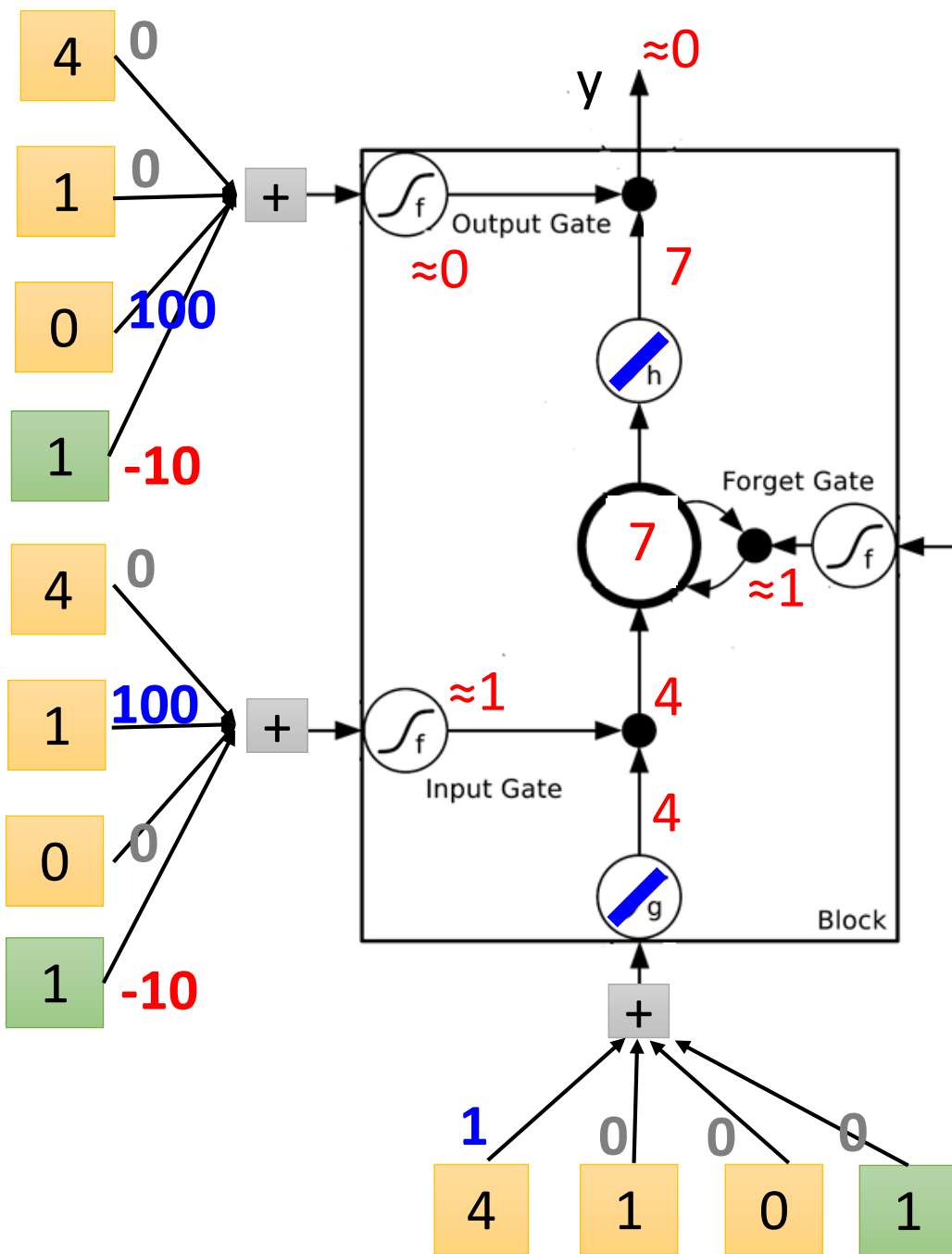
When $x_3 = 1$, output the number in the memory.

Original Network:

> Simply replace the neurons with LSTM

$a_1$

$a_2$

Output Gate

Forget Gate

Cell

Input Gate

Block

**4 times of parameters**

$x_1$

$x_2$

Input

# LSTM

# LSTM

# LSTM

Extension: "peephole"

# *Learning Target*

# Learning



Backpropagation through time (BPTT)

copy

$a_1$  $a_2$

$w$

$w \leftarrow w - \eta\, \partial L \,/\, \partial w$

$y_1$  $y_2$

$x_1$  $x_2$

# Unfortunately ……

- RNN-based network is not always easy to learn

Real experiments on Language modeling

# The error surface is rough.



The error surface is either very flat or very steep.

Clipping

Total Loss

0.35
0.30
0.25
0.20
0.15
0.10
0.05

$w_2$

4.6
4.8
5.0
5.2
5.4

$w_1$

−2.8  −2.6  −2.4  −2.2  −2.0

[Razvan Pascanu, ICML'13]

# Why?

$w = 1$ $\Longrightarrow$ $y^{1000} = 1$

$w = 1.01$ $\Longrightarrow$ $y^{1000} \approx 20000$

$w = 0.99$ $\Longrightarrow$ $y^{1000} \approx 0$

$w = 0.01$ $\Longrightarrow$ $y^{1000} \approx 0$

Large $\partial L / \partial w$ ➡ Small Learning rate?

small $\partial L / \partial w$ ➡ Large Learning rate?

$= w^{999}$

**_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



add

[Cho, EMNLP'14]

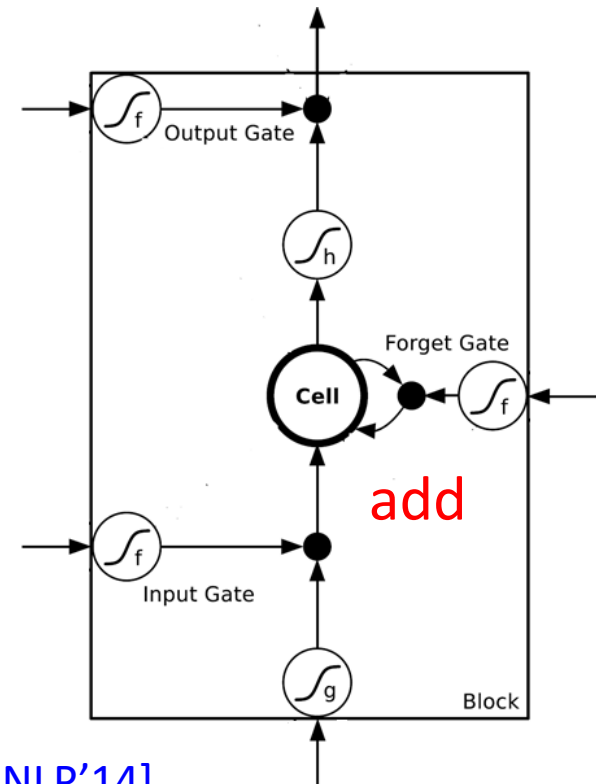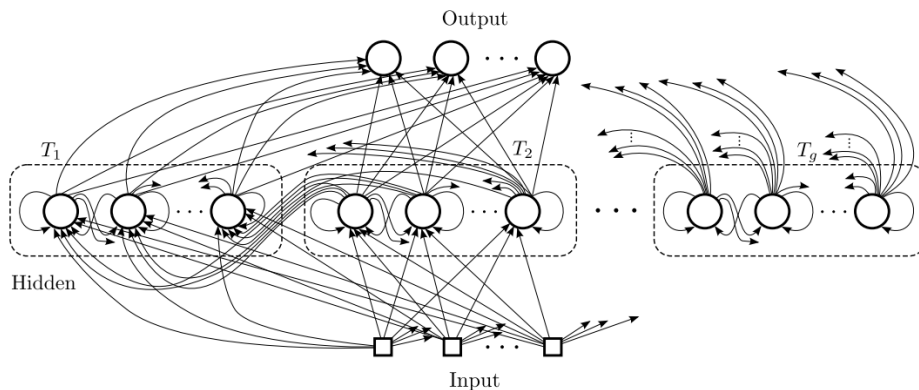# Helpful Techniques

Clockwise RNN

Structurally Constrained
Recurrent Network (SCRN)



[Jan Koutnik, JMLR'14]

[Tomas Mikolov, ICLR'15]

Vanilla RNN Initialized with Identity matrix + ReLU activation function [Quoc V. Le, arXiv'15]
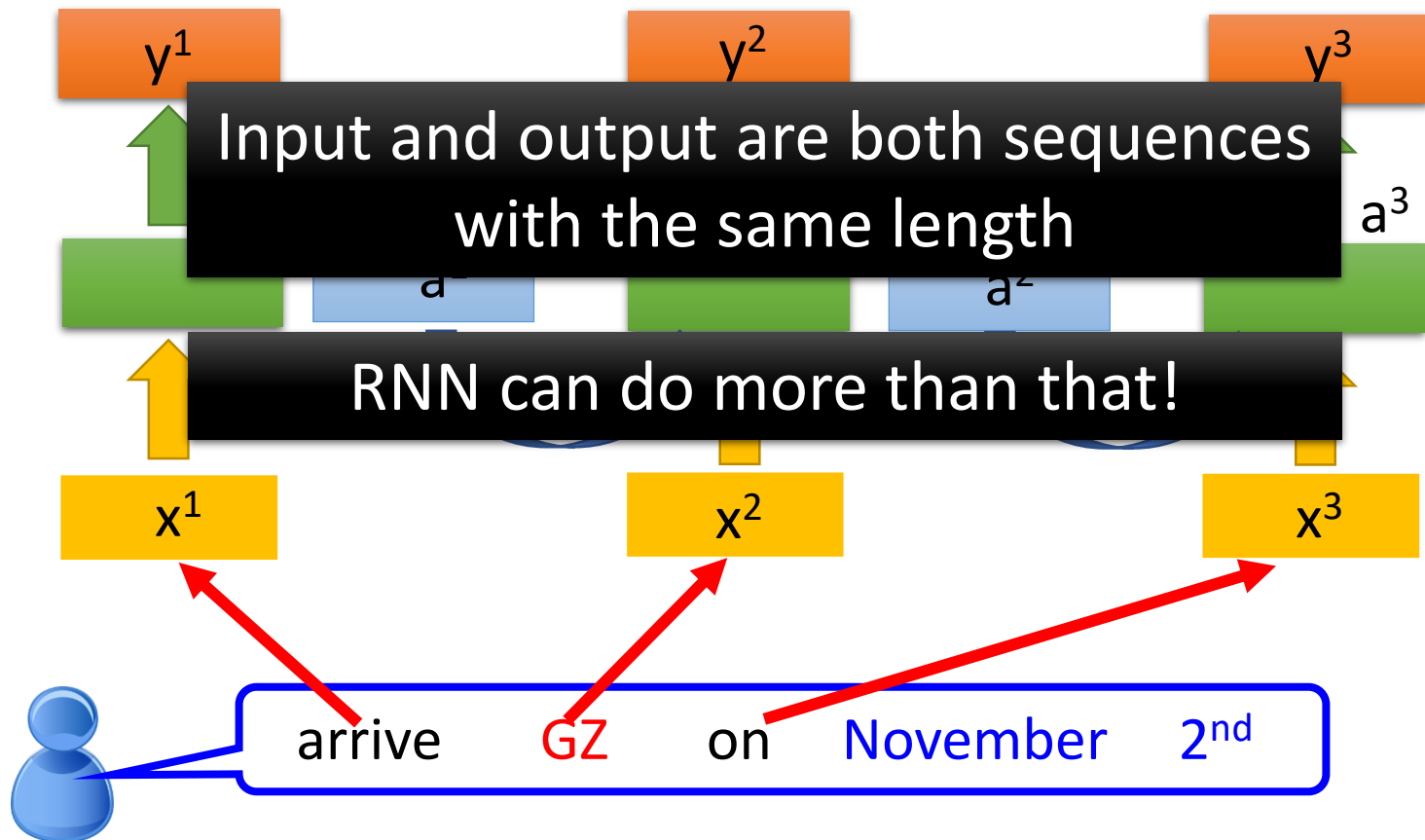
➢ Outperform or be comparable with LSTM in 4 different tasks

# More Applications ......

Probability of "arrive" in each slot

Probability of "GZ" in each slot

Probability of "on" in each slot



$y^1$    $y^2$    $y^3$

$a^3$

$a^2$

**Input and output are both sequences with the same length**

**RNN can do more than that!**

$x^1$    $x^2$    $x^3$

arrive    GZ    on    November    2nd

# Many to one

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

***Sentiment Analysis***

看了这部电影觉
得很高兴 …….

Positive (正面)

这部电影太糟
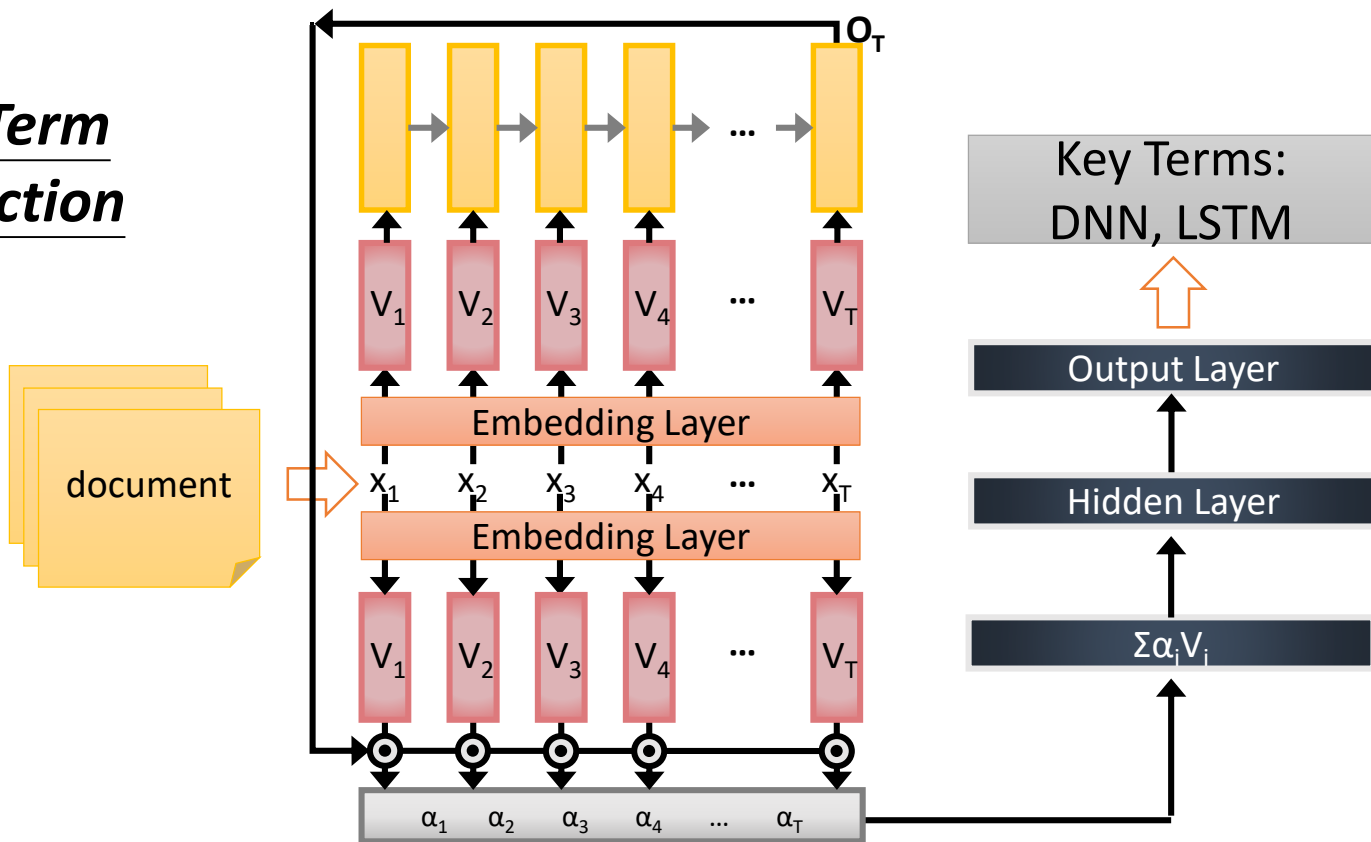了 …….

Negative (负面)

这部电影很
棒 …….

Positive (正面)

超好
好
一般
差
超差

我　　觉　　得　　……　　太　　糟　　了

# Many to one

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



*Key Term Extraction*

# Many to Many (Output is shorter)

- Both input and output are both sequences, ___but the output is shorter.___
    - E.g. ___Speech Recognition___

Output: "好棒" (character sequence)

Trimming

Problem?

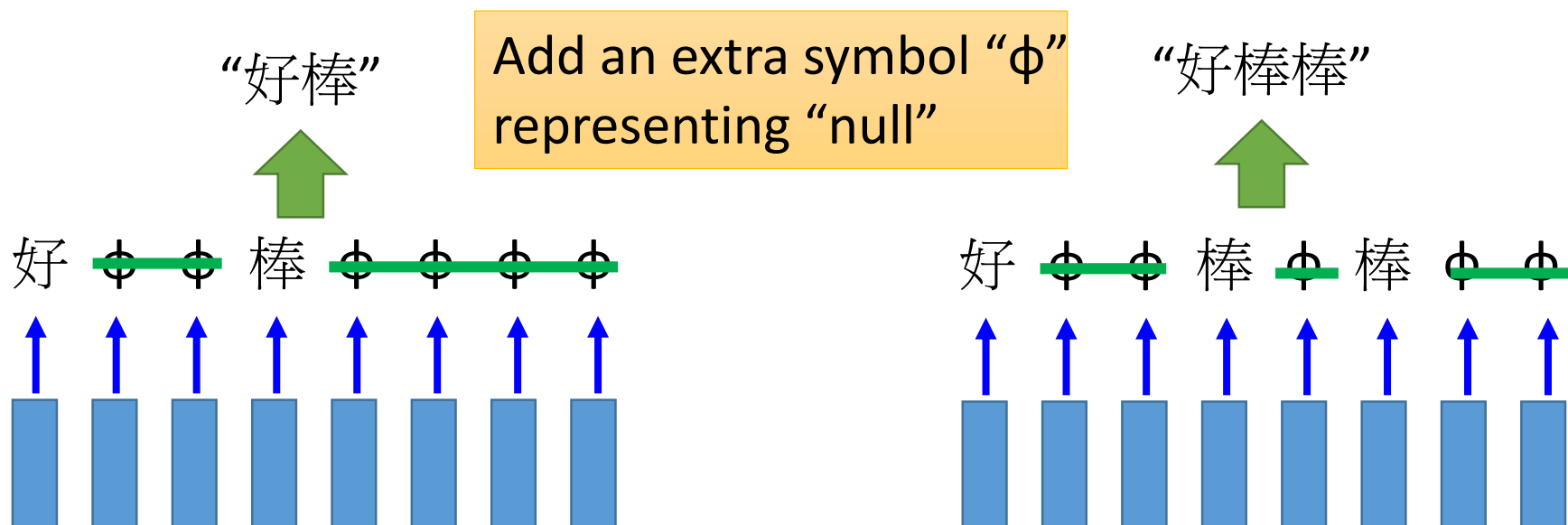Why can't it be "好棒棒"

好 好 好 棒 棒 棒 棒 棒

Input: (vector sequence)

# 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]
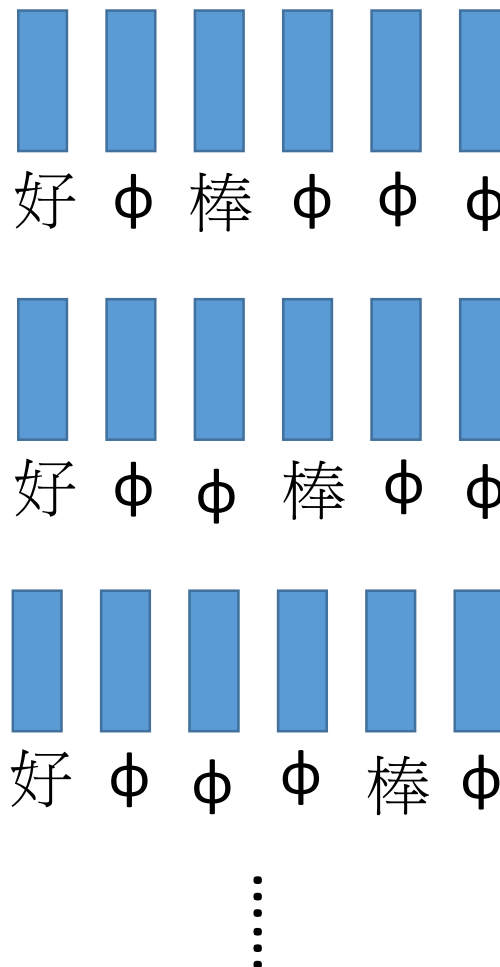
"好棒"

Add an extra symbol "φ" representing "null"

"好棒棒"

好 φ φ 棒 φ φ φ φ

好 φ φ 棒 φ 棒 φ φ

# Many to Many (Output is shorter)

- CTC: Training

Acoustic Features:

Label:　好 棒

All possible alignments are considered as correct.

好 ф 棒 ф ф ф

好 ф ф 棒 ф ф

好 ф ф ф 棒 ф

⋮

# 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.

# The End!