# Deep Learning
## Neural Network with Memory (1)

Hung-yi Lee

# Memory is important
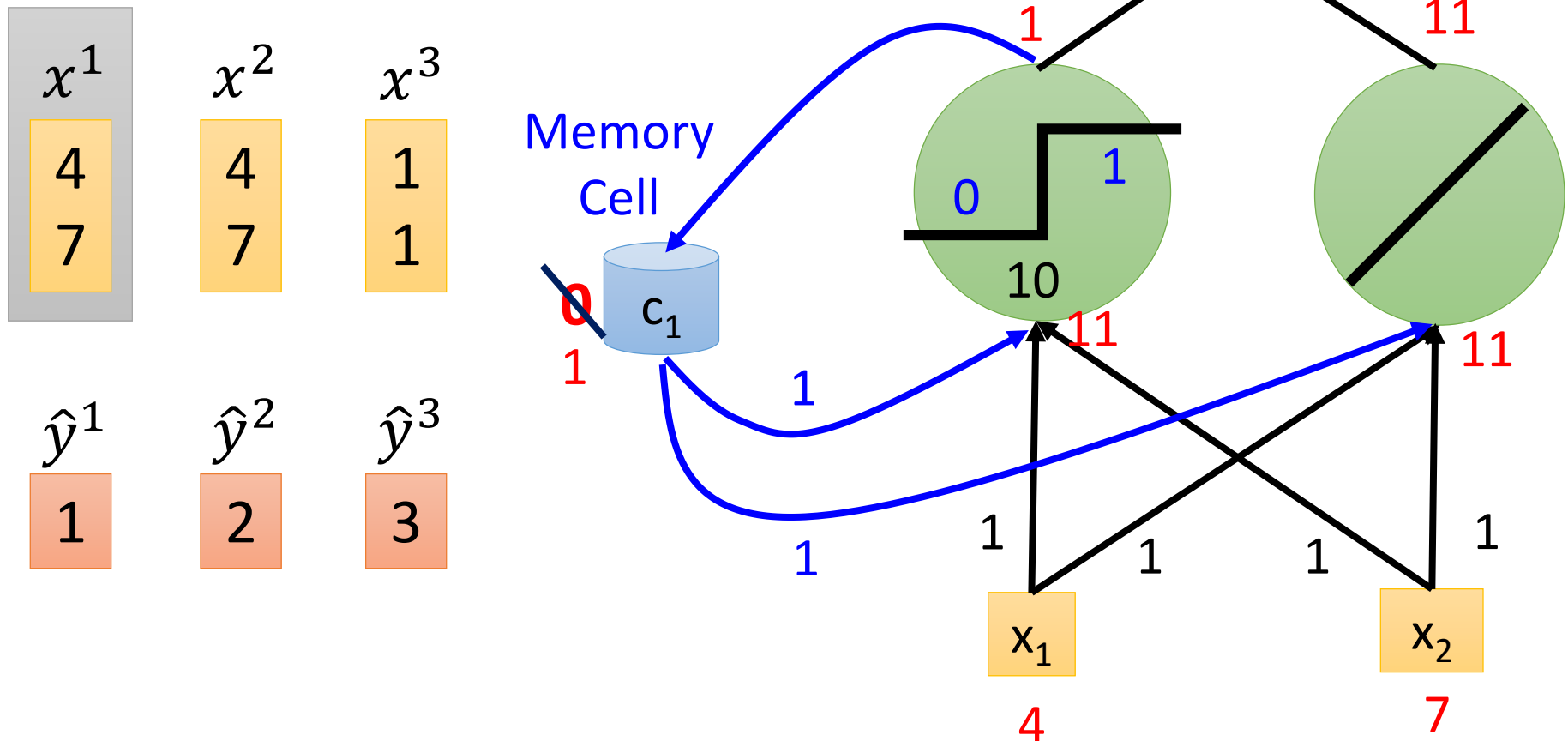
Input:
2 dimensions

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

| 4 | 4 | 1 |
| 7 | 7 | 1 |

Output:
1 dimension

$\hat{y}^1$ $\hat{y}^2$ $\hat{y}^3$

| 1 | 2 | 3 |

```
      1   1
    1   4   4
+   1   7   7
_____
    3   2   1
```

Network needs memory
to achieve this

# Memory is important

## Network with Memory

# Memory is important

## Network with Memory

$x^1$

| 4 |
| 7 |

$x^2$

| 4 |
| 7 |

$x^3$

| 1 |
| 1 |

$\hat{y}^1$

| 1 |

$\hat{y}^2$

| 2 |

$\hat{y}^3$

| 3 |

Memory Cell

**0**
1
1
0

$c_1$

$y$

3

-10

0

1

3

1

0
1
10

3

3

1

1

1

1

1

$x_1$

$x_2$
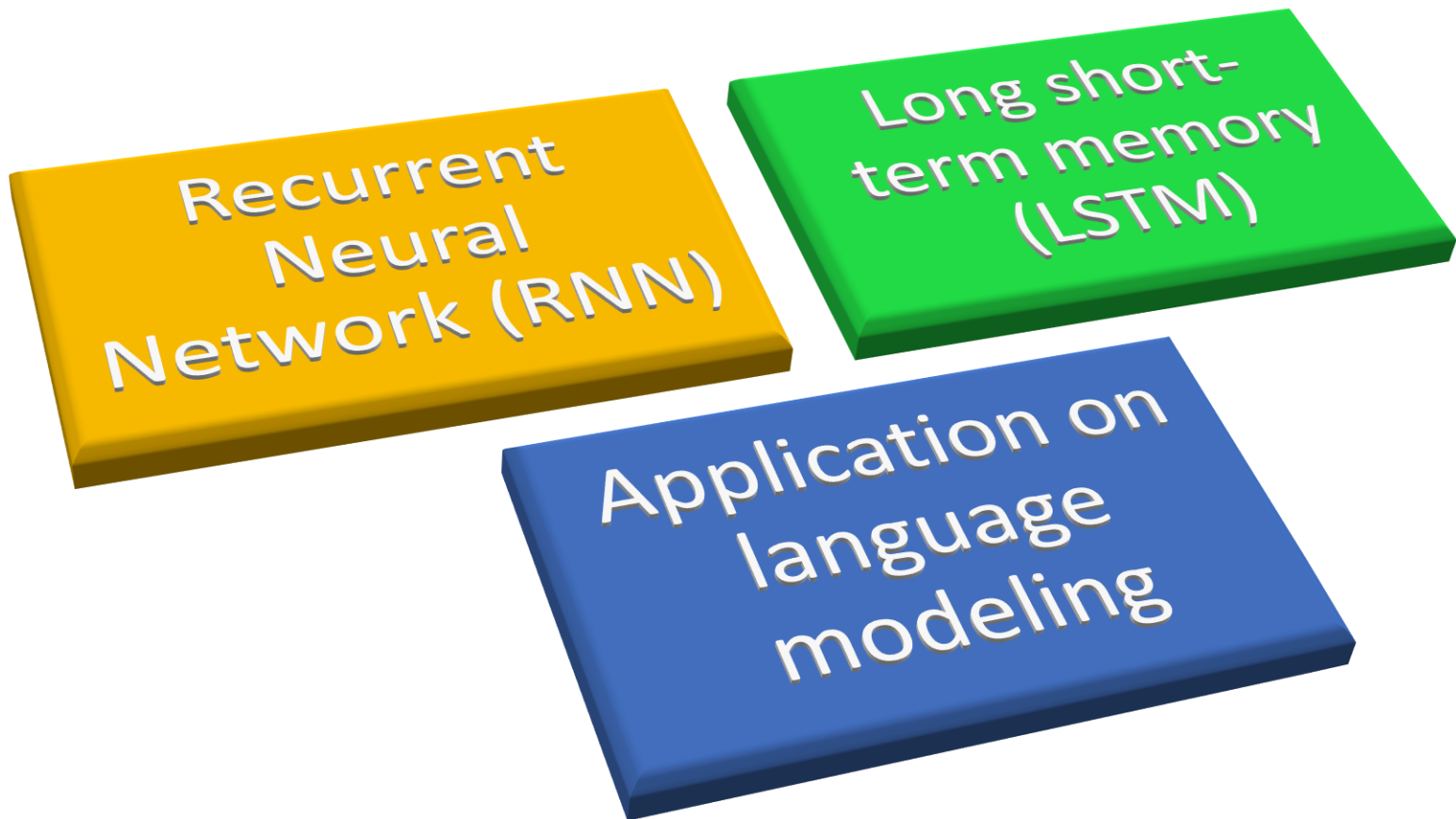
1

1

# Outline

Today we will simply assume we already know the parameters in the network.
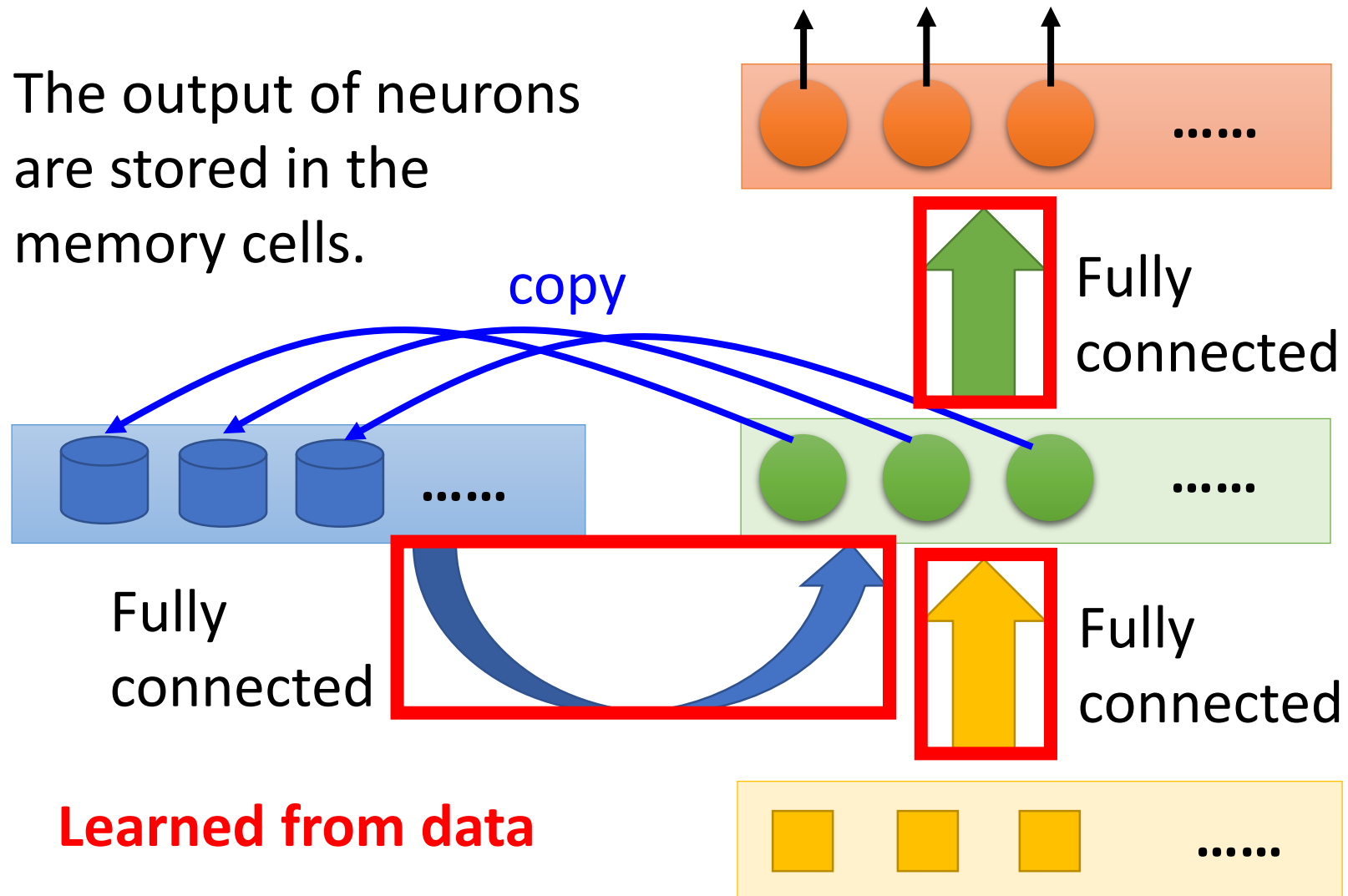Training will be discussed in the next week.

Recurrent Neural Network (RNN)

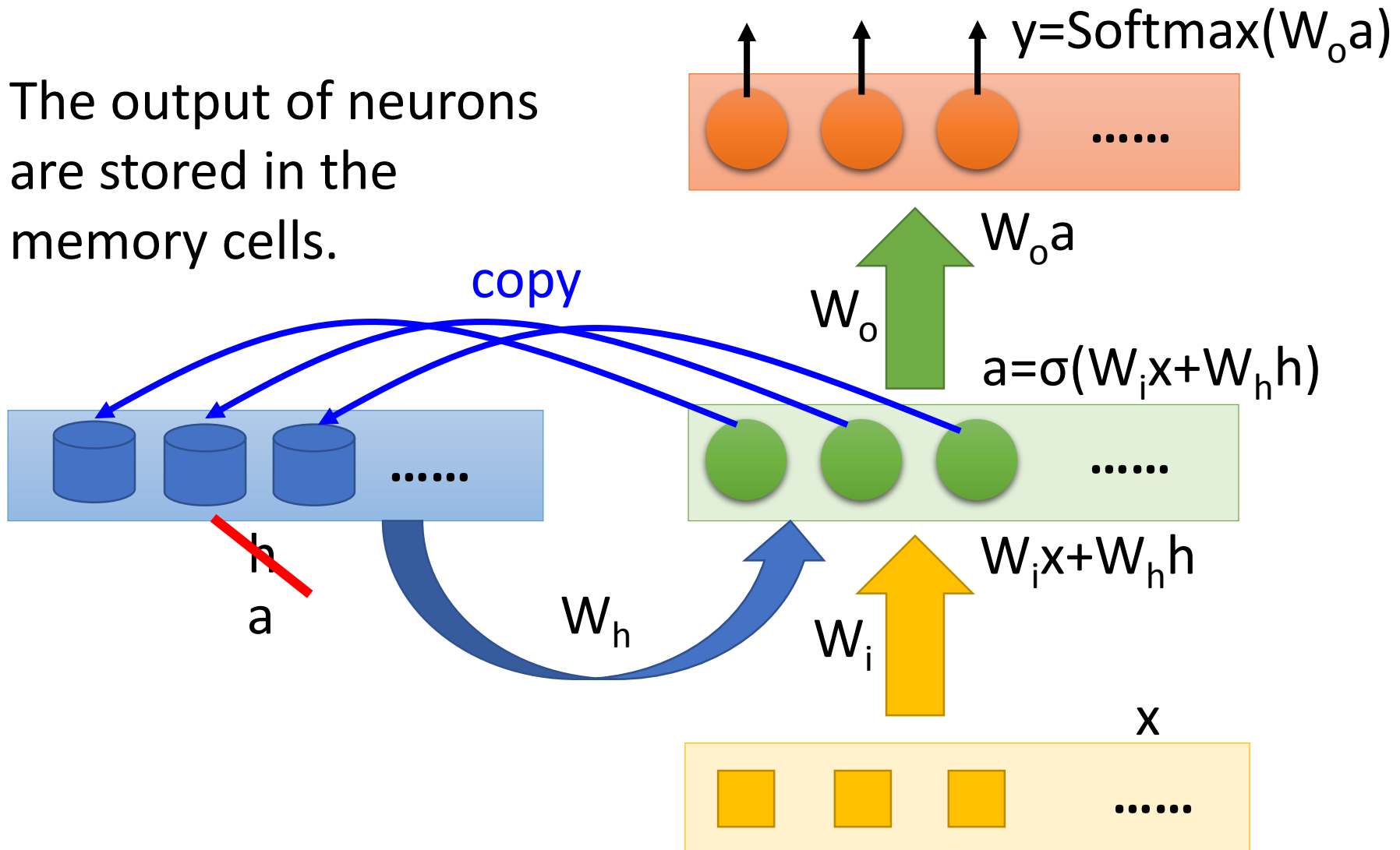Long short-term memory (LSTM)

Application on language modeling

# Outline

# Recurrent Neural Network (RNN)
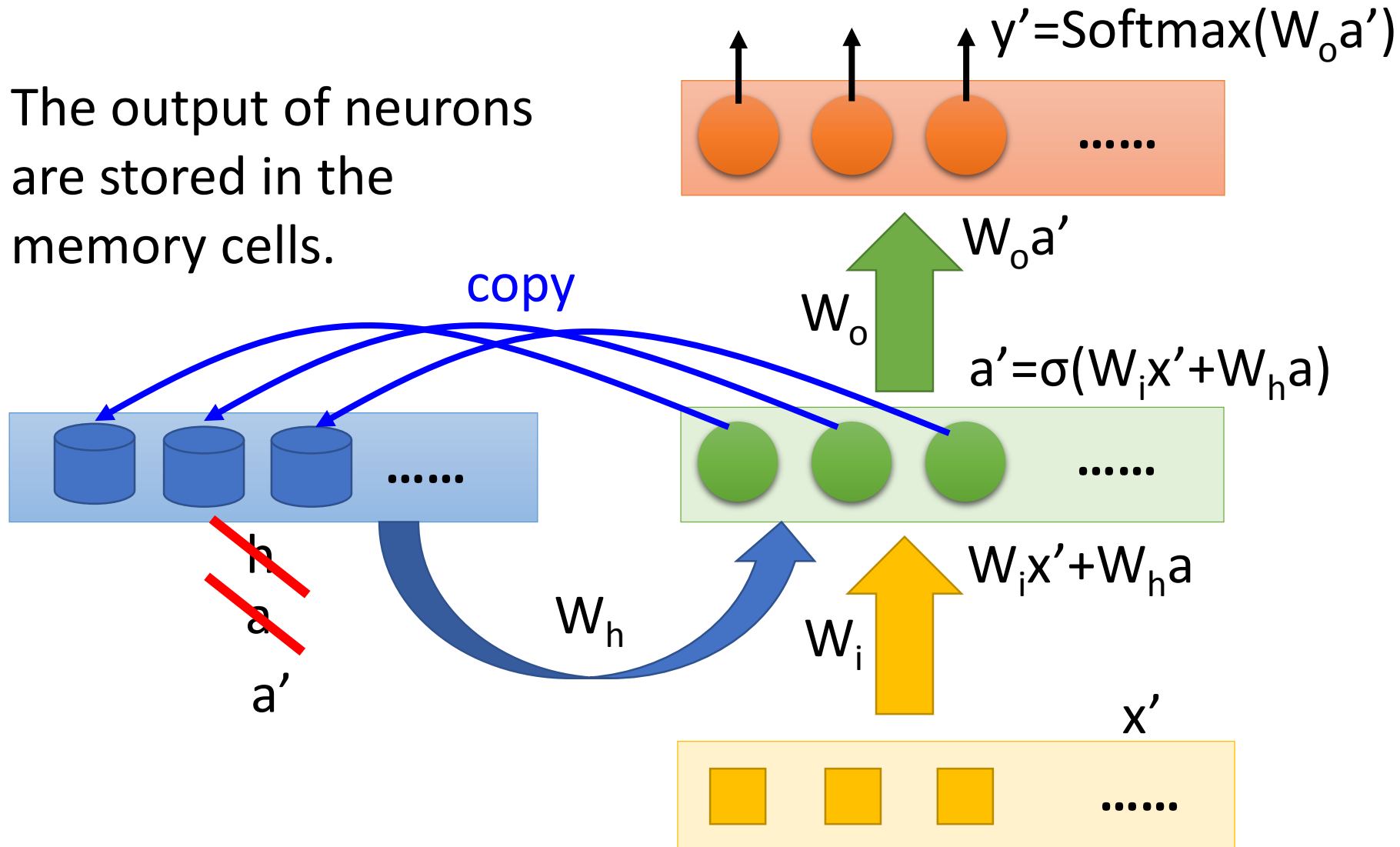


The output of neurons are stored in the memory cells.

copy

Fully connected

Fully connected

Fully connected

**Learned from data**

# Recurrent Neural Network (RNN)

The output of neurons are stored in the memory cells.

$y=\text{Softmax}(W_o a)$

copy

$W_o a$

$W_o$

$a=\sigma(W_i x + W_h h)$

$W_i x + W_h h$

h

a

$W_h$

$W_i$

x

# Recurrent Neural Network (RNN)

The output of neurons are stored in the memory cells.

$y'=Softmax(W_o a')$

copy

$a'=\sigma(W_i x'+W_h a)$

$W_o a'$

$W_o$

$W_i x'+W_h a$

$W_h$

$W_i$

h

a

a'

x'

# Outline

Recurrent Neural Network (RNN)

Long short-term memory (LSTM)

(simplified version)

Application on language modeling

# Long Short-term Memory (LSTM)

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

multiply

$z_o$

Output Gate

$f(z_o)$

$h(c')$

Activation function f is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

Forget Gate

$c$  $f(z_f)$

$c'$

$z_f$

$cf(z_f)$

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

$f(z_i)$  $g(z)f(z_i)$

$z_i$

Input Gate

multiply

$g(z)$

$g$

Block

$z$

Original Network:

> Simply replace the neurons with LSTM

$a_1$

$a_2$

Output Gate

$\int_f$

$\int_h$

Forget Gate

Cell

$\int_f$

Input Gate

$\int_f$

$\int_g$

Block

**4 times of parameters**

$x_1$

$x_2$

Input

# 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 |

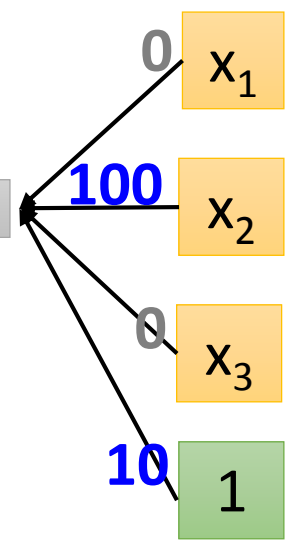| $y$ | 0 | 0 | 0 | 0 | 0 | 7 | 0 | 0 | 6 |
|---|---|---|---|---|---|---|---|---|---|

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.

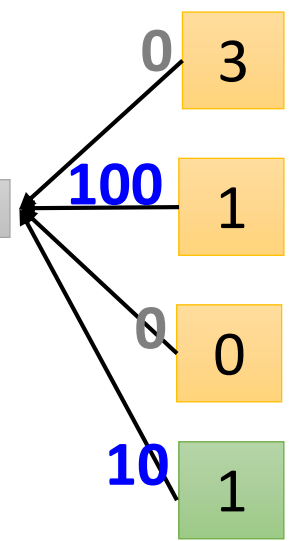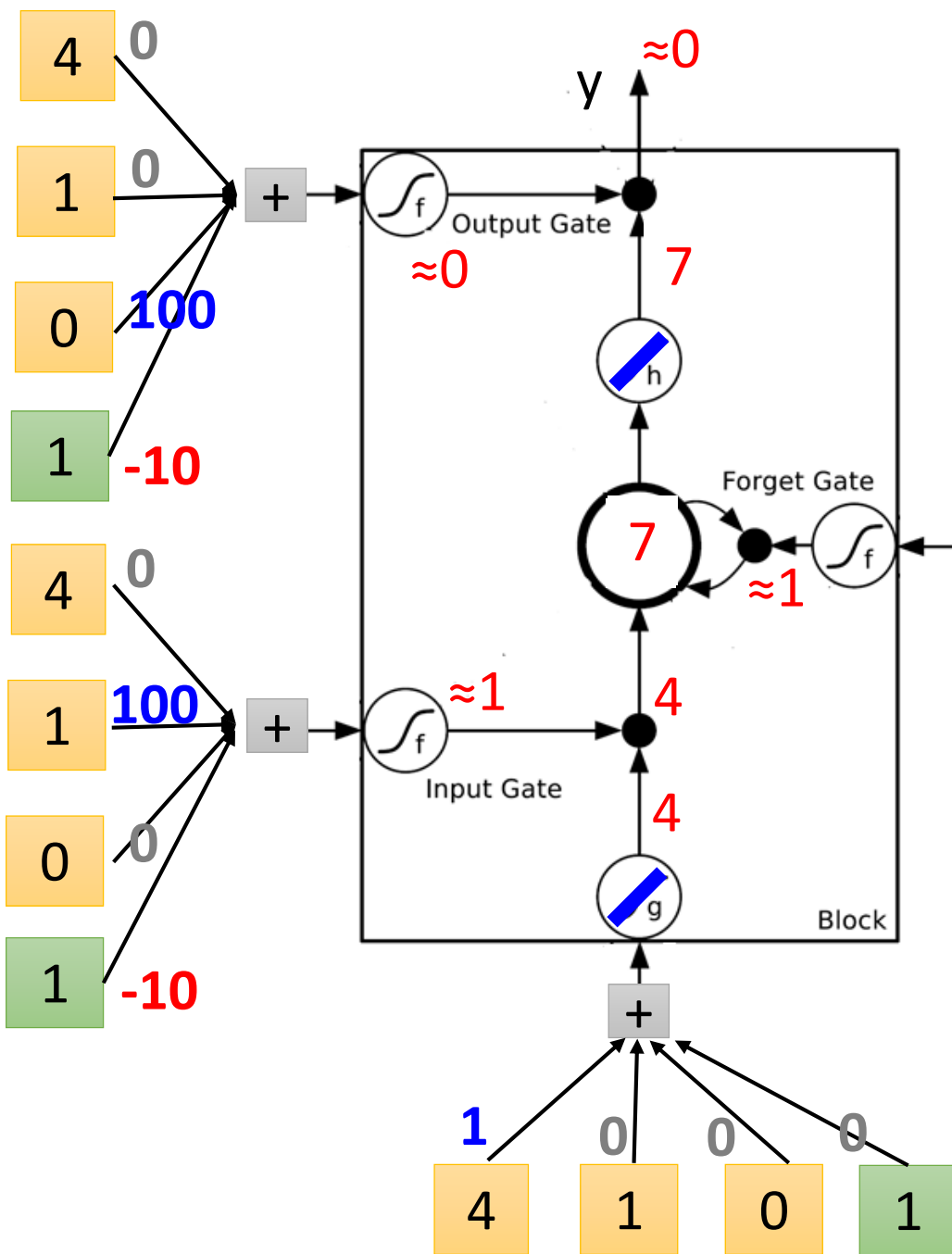| | | | | | |
|---|---|---|---|---|---|
| $y$ | 0 | 0 | 0 | 7 | 0 |

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

# Outline

Recurrent Neural Network (RNN)

Long short-term memory (LSTM)

Application on language modeling

# Language model (LM)

- Language model: Estimated the probability of word sequence
  - Word sequence: $w_1, w_2, w_3, ...., w_n$
  - $P(w_1, w_2, w_3, ...., w_n)$
- Useful in speech recognition
  - Different word sequence can have the same pronunciation



recognize speech
or
wreck a nice beach

If P(recognize speech)
>P(wreck a nice beach)

Output =
"recognize speech"

# Language model

P("wreck a nice beach")
=P(wreck|START)P(a|wreck)
P(nice|a)P(beach|nice)

- How to estimate $P(w_1, w_2, w_3, ...., w_n)$
- Collect a large amount of text data as training data
    - However, the word sequence $w_1, w_2, ...., w_n$ may not appear in the training data
- N-gram language model: $P(w_1, w_2, w_3, ...., w_n) =$ $P(w_1|START)P(w_2|w_1) ...... P(w_n|w_{n-1})$
- Estimate P(beach|nice) from training data

$$P(\text{beach}|\text{nice}) = \frac{C(nice\ beach)}{C(nice)}$$

Count of "nice beach" in the training data

Count of "nice" in the training data

# Language model - Smoothing

- Training data:
  - The dog ran ……
  - The cat jumped ……

This is called **language model smoothing**.

P( jumped | dog ) = ~~0~~  0.0001

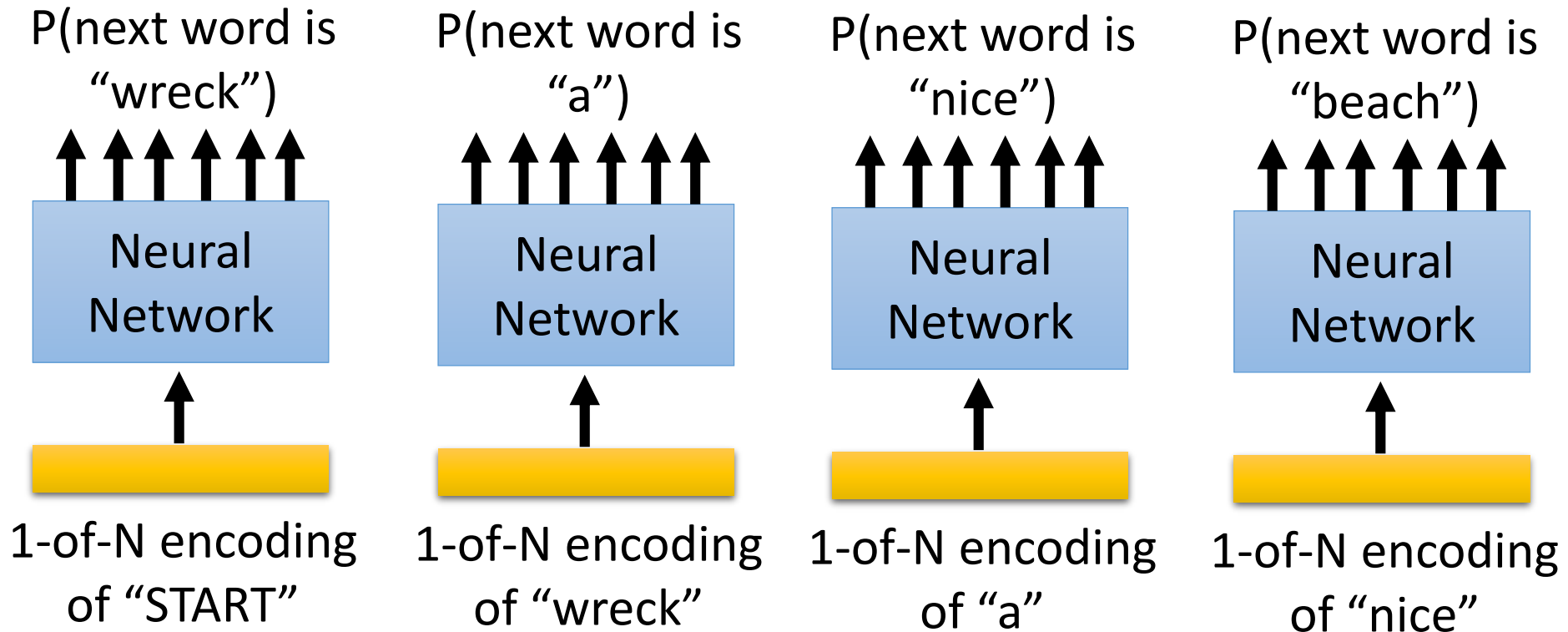P( ran | cat ) = ~~0~~  0.0001

Give some small probability

➢ The probability is not accurate.

➢ The phenomenon happens because we cannot collect all the possible text in the world as training data.
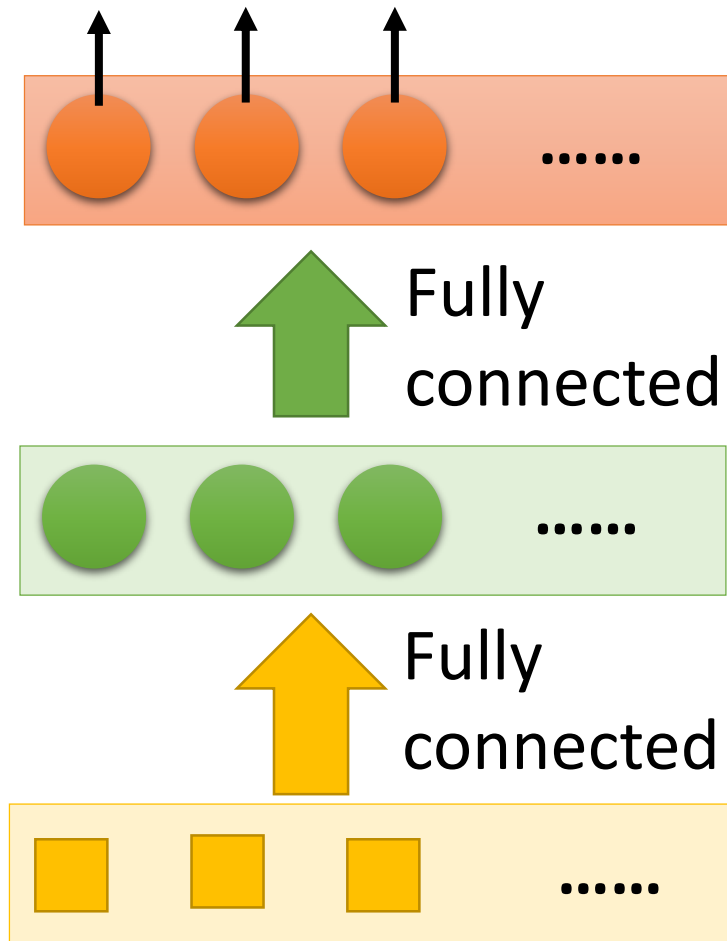
# Neural-network based LM

P("wreck a nice beach")

=P(wreck|START)P(a|wreck)P(nice|a)P(beach|nice)

P(b|a): not from count, but the NN that can predict the next word.

P(next word is "wreck")

Neural Network

1-of-N encoding of "START"

P(next word is "a")

Neural Network

1-of-N encoding of "wreck"

P(next word is "nice")

Neural Network

1-of-N encoding of "a"

P(next word is "beach")

Neural Network

1-of-N encoding of "nice"

# Neural-network based LM



The hidden layer of the related words are close.
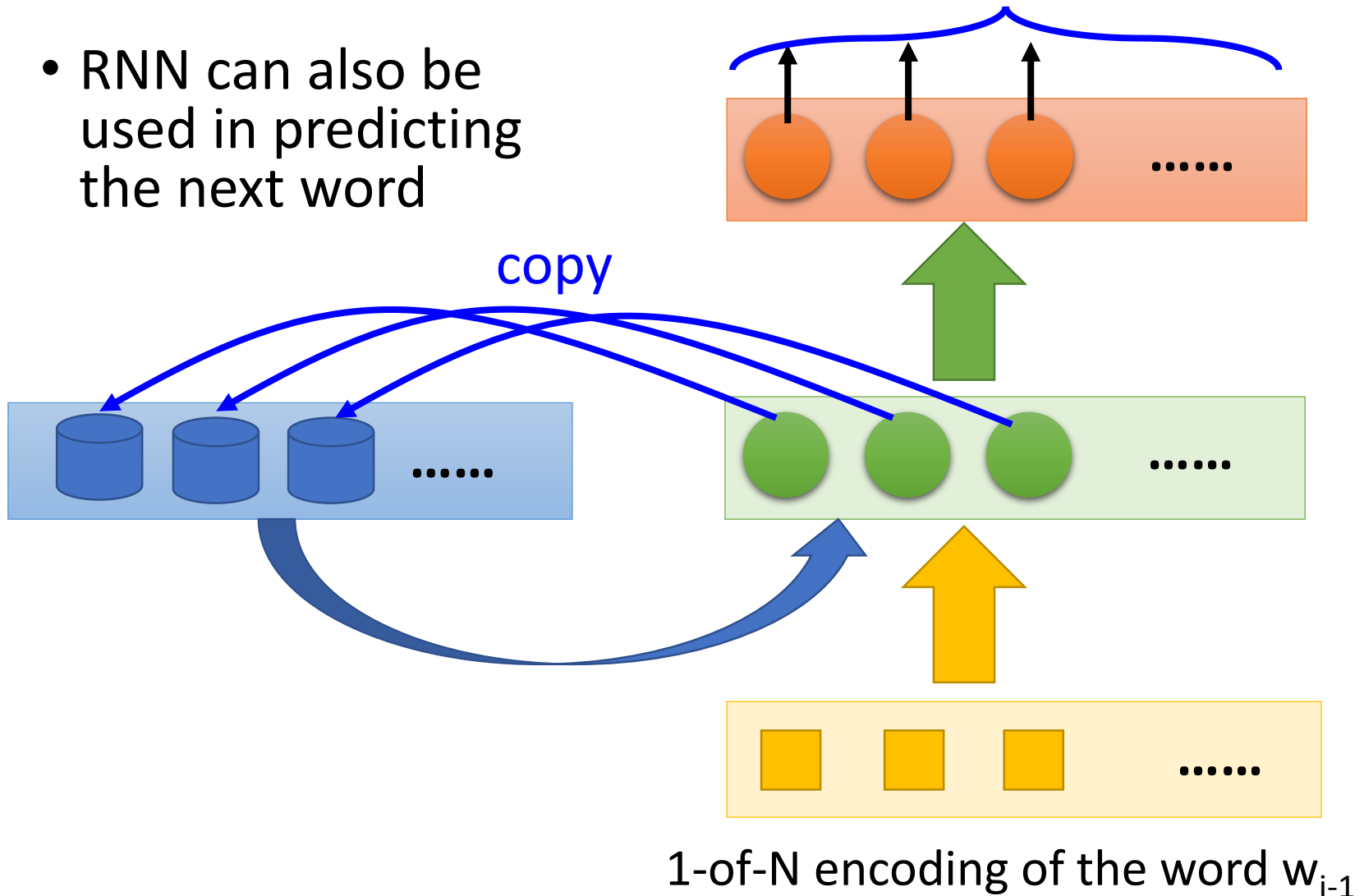
If P(jump|dog) is large, then P(jump|cat) increase accordingly.

(even there is not "... cat jump ..." in the data)

Smoothing is automatically done.

# RNN-based LM
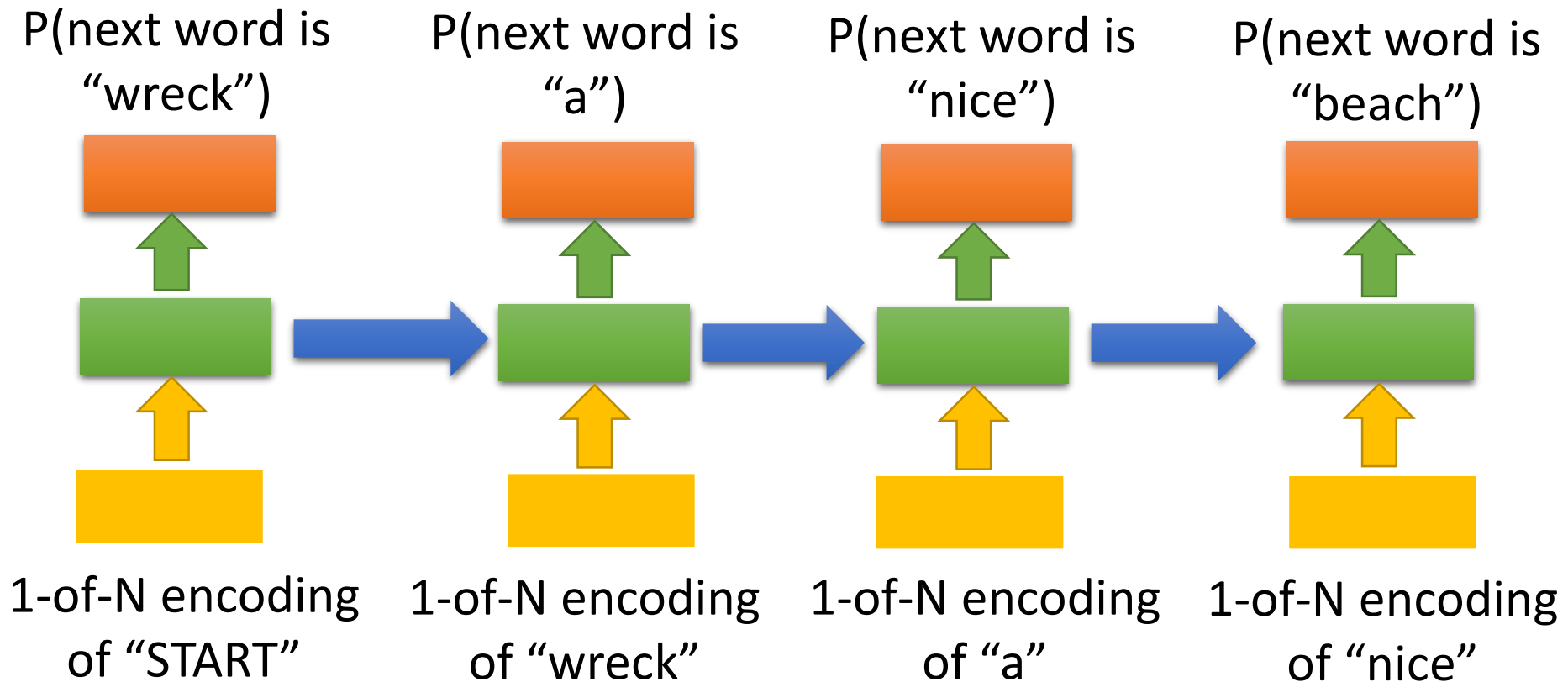
The probability for each word as the next word $w_i$

- RNN can also be used in predicting the next word

copy

1-of-N encoding of the word $w_{i-1}$

# RNN-based LM

P(next word is "wreck")    P(next word is "a")    P(next word is "nice")    P(next word is "beach")

1-of-N encoding of "START"    1-of-N encoding of "wreck"    1-of-N encoding of "a"    1-of-N encoding of "nice"

➢ Model long-term information

➢ Can also consider LSTM

# Another Applications

- Composer
    - http://people.idsia.ch/~juergen/blues/
- Sentence generation
    - http://www.cs.toronto.edu/~ilya/rnn.html

# Summary – Network with Memory

Recurrent Neural Network (RNN)

Long short-term memory (LSTM)

Application on language modeling

Thank You

# Appendix

# Long term memory

# RNN - Training

Training the parameters to let $y$ close to $\hat{y}$
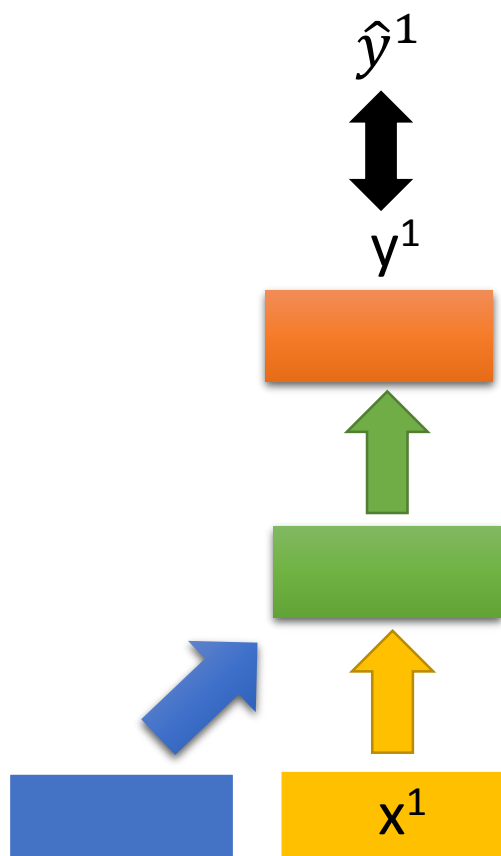
$\hat{y}^1$ $\hat{y}^2$ $\hat{y}^3$

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

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

*UNROLL*

$\hat{y}^1$

$y^1$

$x^1$

$\hat{y}^2$

$y^2$

$x^1$ $x^2$

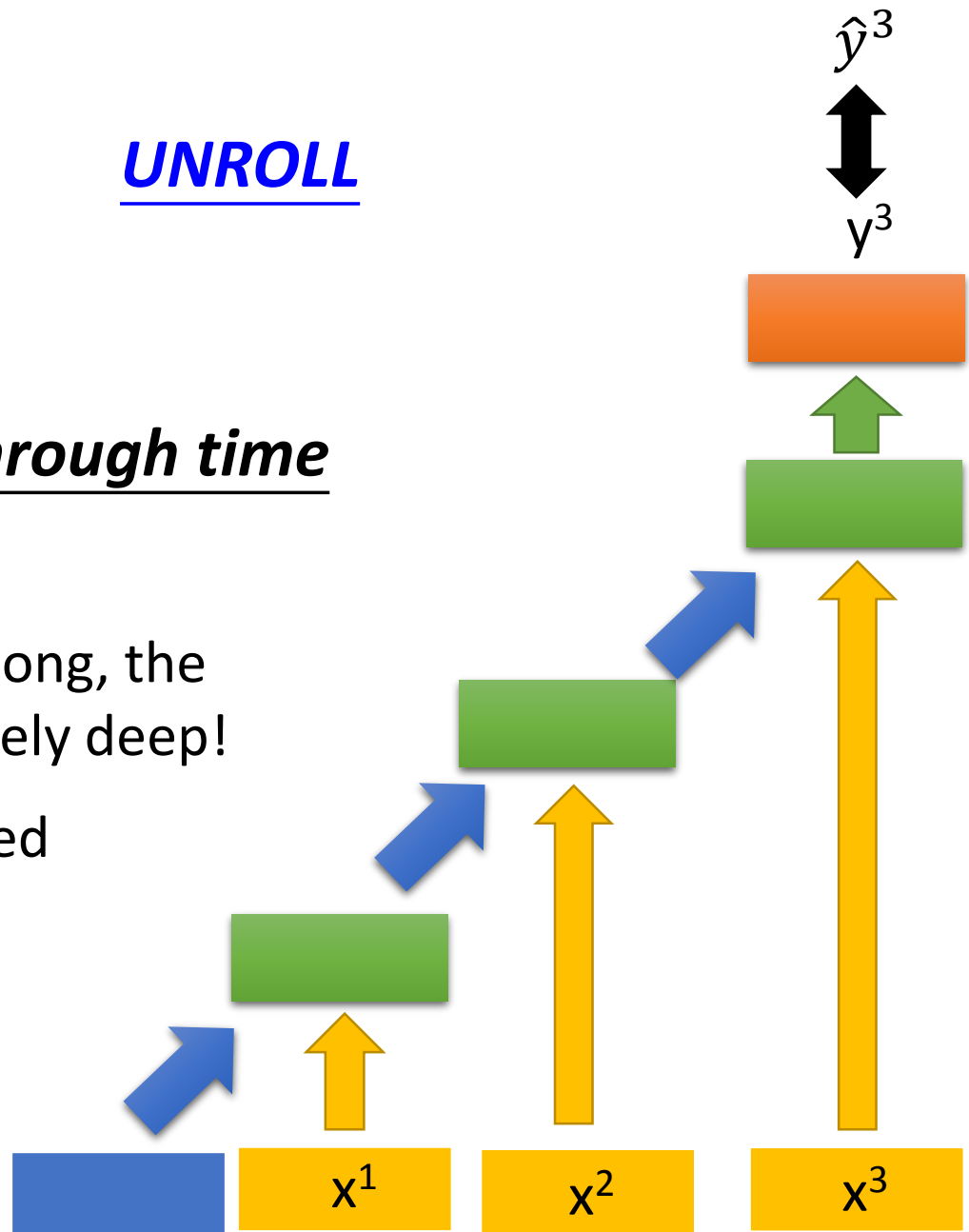*UNROLL*

## *Backpropagation through time (BPTT)*

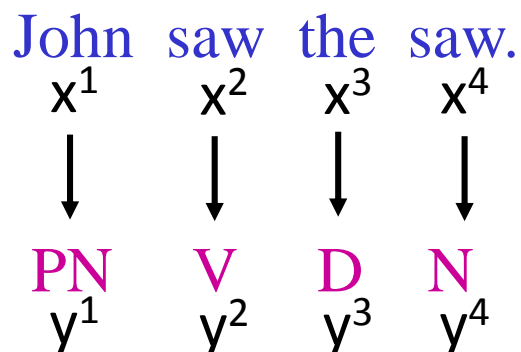When the sequence is long, the network can be extremely deep!

Some weights are shared

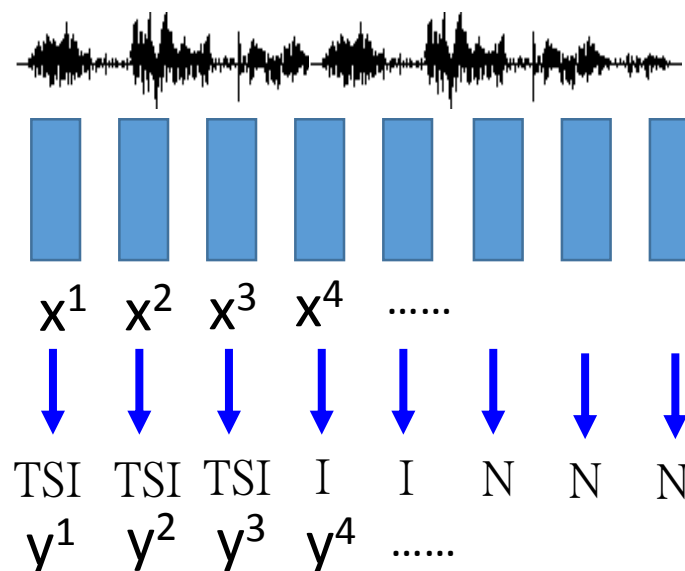# The task that needs memory
# Structured Leargning v.s. RNN

## POS Tagging

John   saw   the   saw.
$x^1$    $x^2$    $x^3$    $x^4$

↓      ↓      ↓      ↓

PN    V    D    N
$y^1$    $y^2$    $y^3$    $y^4$

## Speech Recognition

$x^1$ $x^2$ $x^3$ $x^4$ ……

↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓

TSI TSI TSI I I N N N
$y^1$ $y^2$ $y^3$ $y^4$ ……
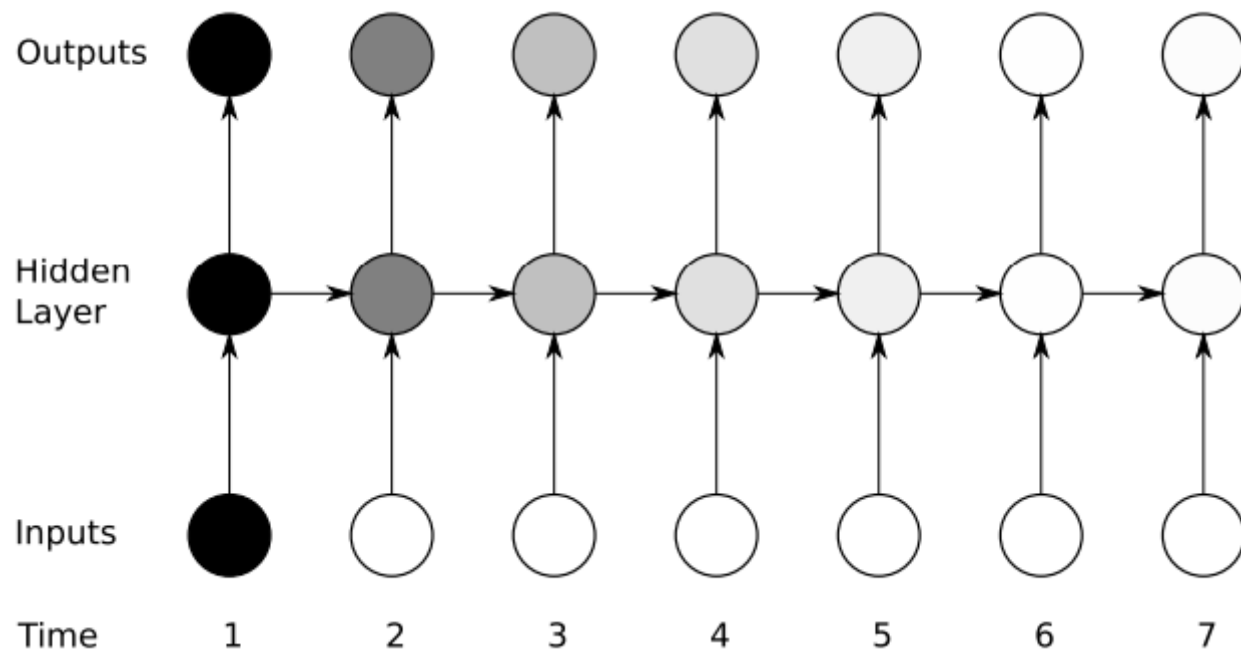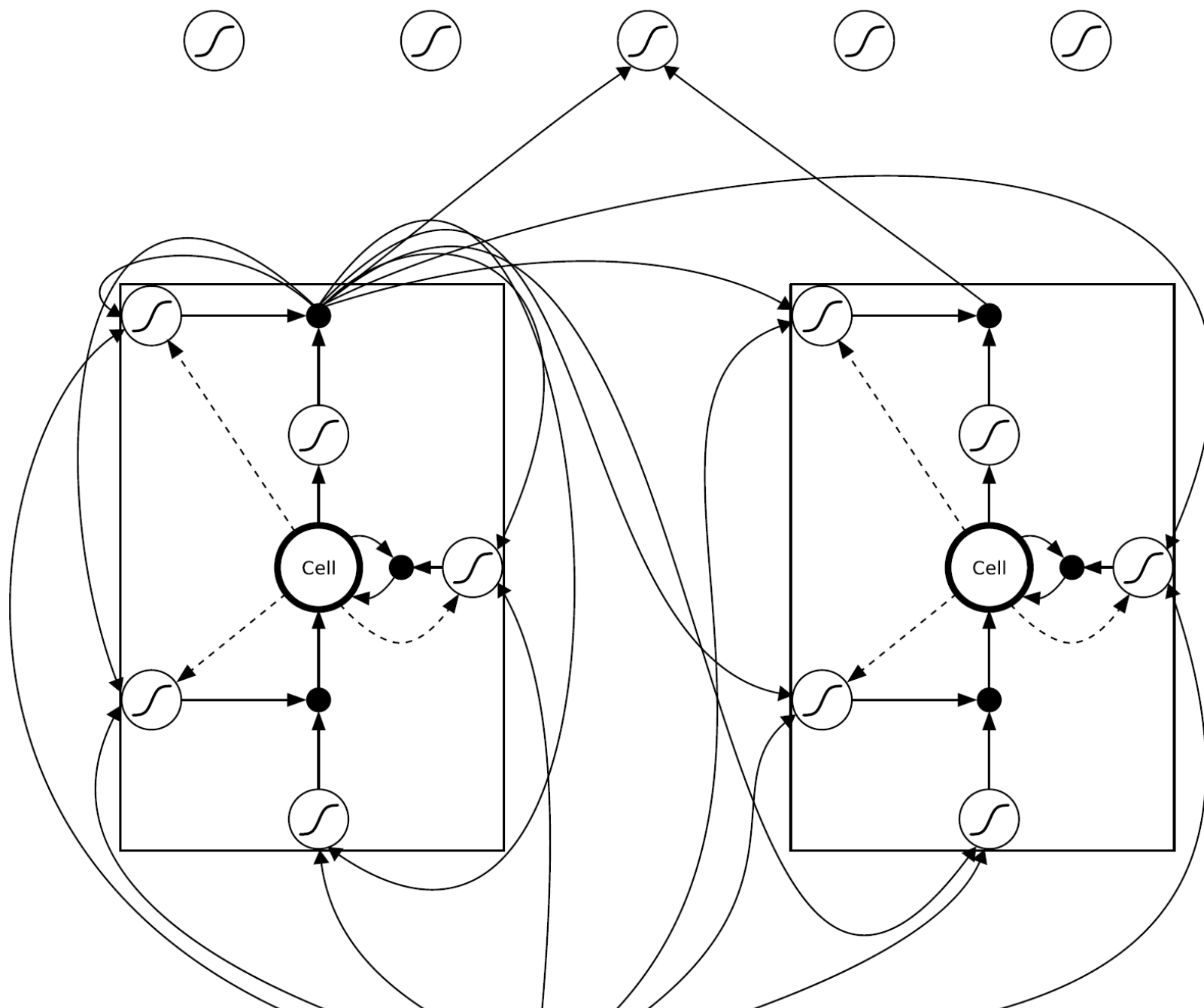
➢ Structured learning can also deal with these problems

➢ What are their difference?

# Outlook

- Speech recognition
  - http://www.cs.toronto.edu/~fritz/absps/RNN13.pdf


- Structured + Deep
  - http://research.microsoft.com/pubs/210167/rcrf_v9.pdf

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| Outputs | | | | | | | |
| Hidden Layer | | | | | | | |
| Inputs | | | | | | | |
| Time | | | | | | | |

# Neural-network based LM

- Training
  - Training data: "Hello how are you …… "

Input:                                    Target:

"Hello"  ➡️  **Neural Network**  ➡️  "how" have the largest probability
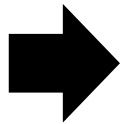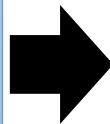
"how"                                     "are" have the largest probability

"are"                                     "you" have the largest probability