

# Understanding Deep Neural Networks

## Chapter Seven

# Sequence Learning

---

深度学习引论2020

Zhang Yi, *IEEE Fellow*  
Autumn 2020

# Outline

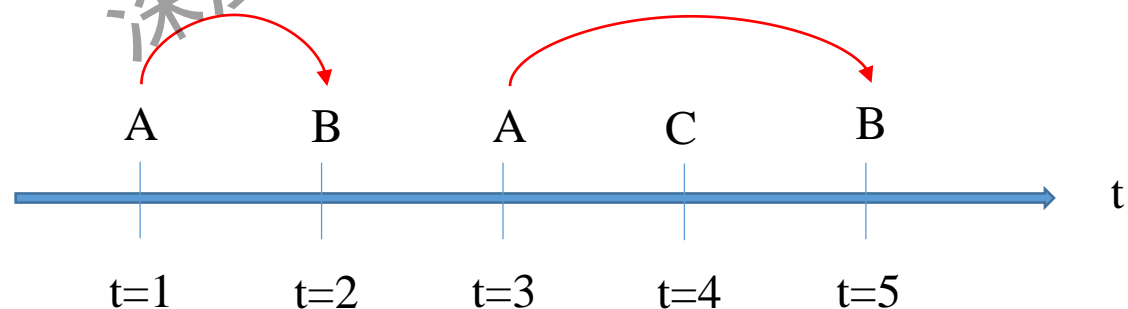
- A Sequence Recognizing Example
- Review of BP for Mono-target Output NNs
- BP Method for Multi-target Outputs NNs
- BP Algorithm for Multi-target Outputs NNs
- Illustrative Examples
- Assignment

# A Sequence Recognizing Example

## Generated Sequences

1. ABACB
2. CCBBA
3. CACCB
4. ACCCB
5. CACBC
6. AAACB
7. BAACB
8. CCBAB
9. BCCAB
10. CABAC
- .....

*Problem: Recognize A followed by B.*



# A Sequence Recognizing Example

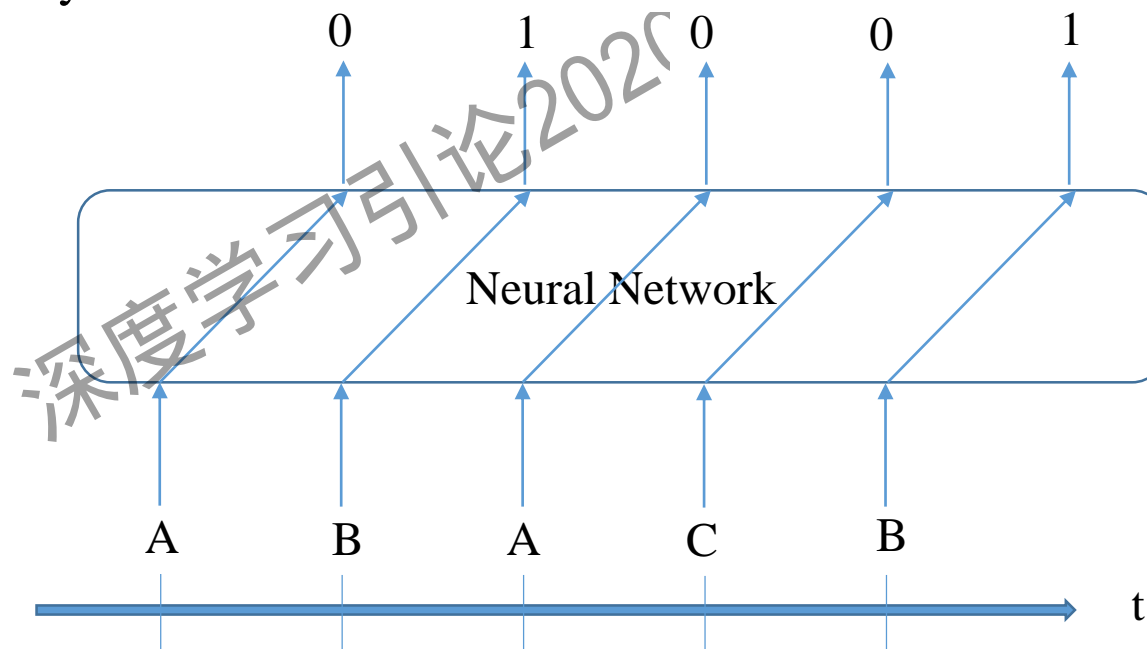
## *Recognize A followed by B Problem*

The task is to recognize A followed by B.

Generated Sequences

1. ABCAB
2. CCBBA
3. CACCB
4. ACCCB
5. CACBC
6. AAACB
7. BAACB
8. CCBAB
9. BCCAB
10. CABAC

.....



**Problem:** Can we use neural network to solve this problem?

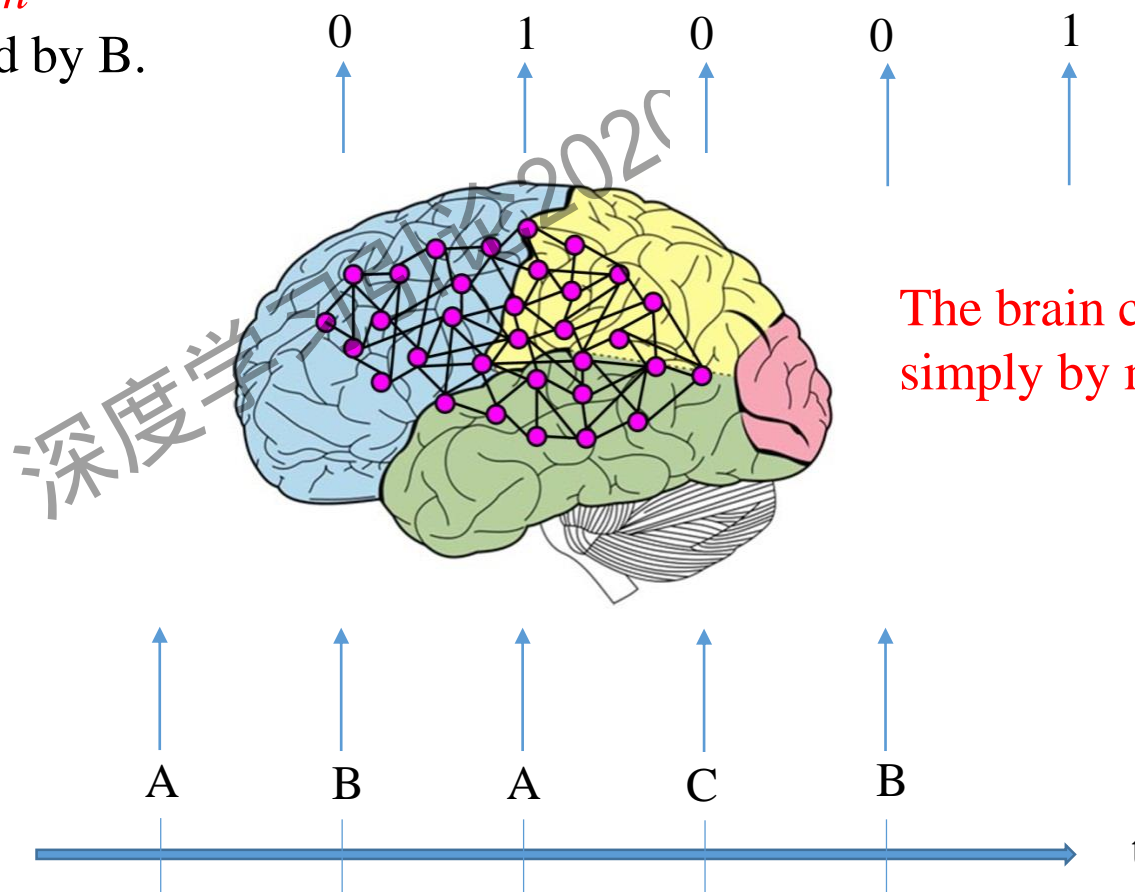
# A Sequence Recognizing Example

## *Recognize A followed by B Problem*

The task is to recognize A followed by B.

Generated Sequences

1. ABCAB
2. CCBBA
3. CACCB
4. ACCCB
5. CACBC
6. AAACB
7. BAACB
8. CCBAB
9. BCCAB
10. CABAC
- .....



The brain can solve this problem simply by memorizing the last A.

# A Sequence Recognizing Example

## *Recognize A followed by B Problem*

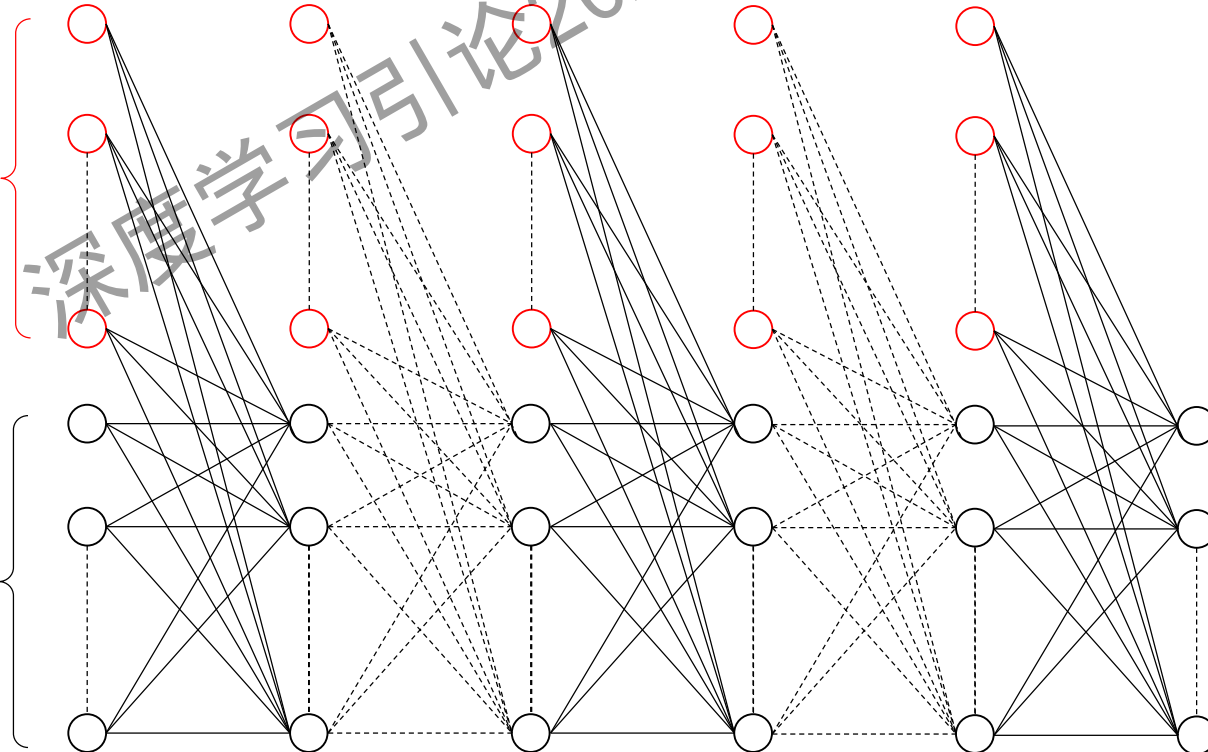
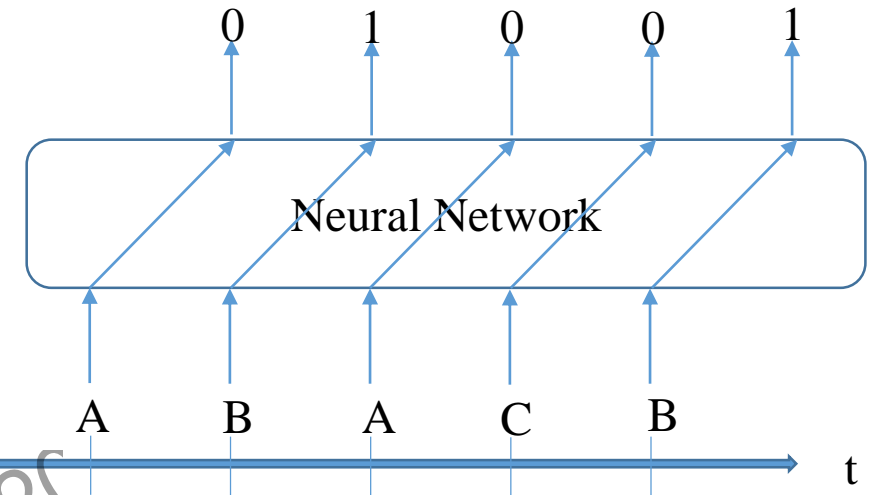
The task is to recognize A followed by B.

### Generated Sequences

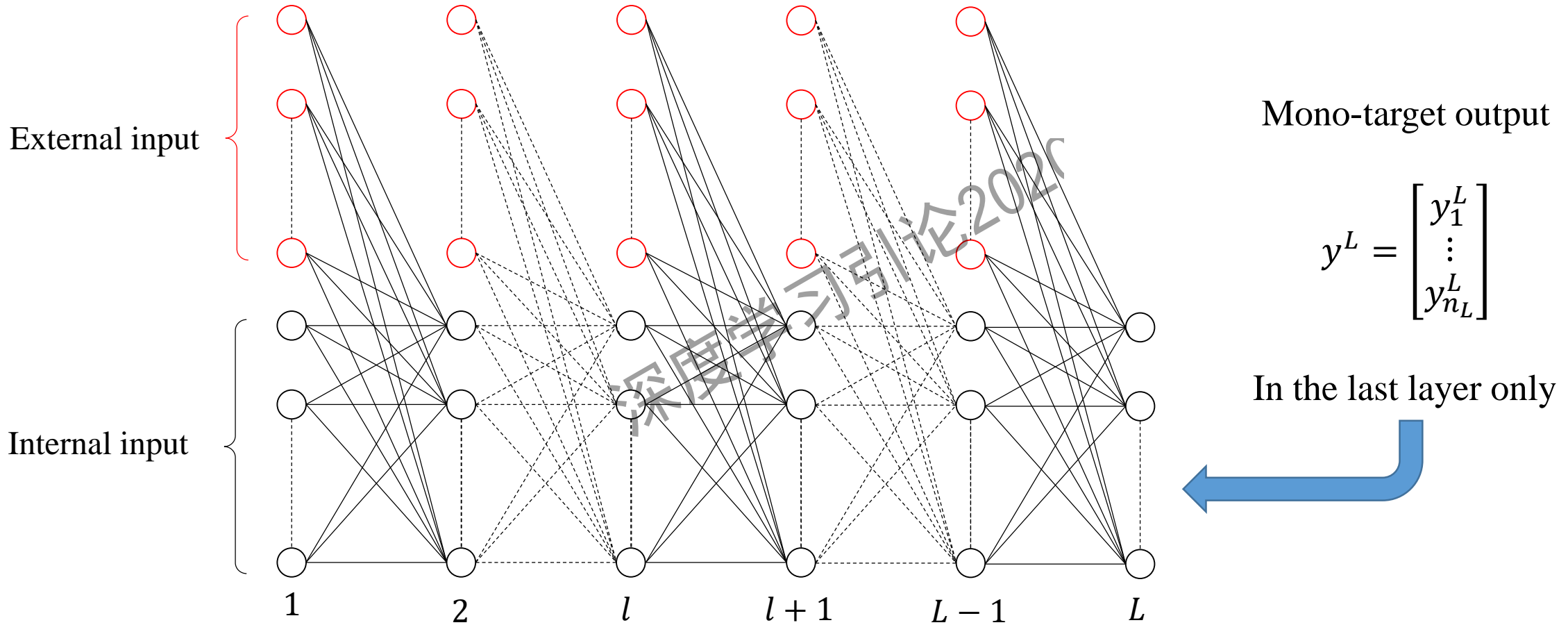
1. ABCAB
2. CCBBA
3. CACCB
4. ACCCB
5. CACBC
6. AAACB
7. BAACB
8. CCBAB
9. BCCAB
10. CABAC
- .....

External input

Internal input



# A Sequence Recognizing Example



Mono-target output network cannot solve the sequence recognizing problem.

# A Sequence Recognizing Example

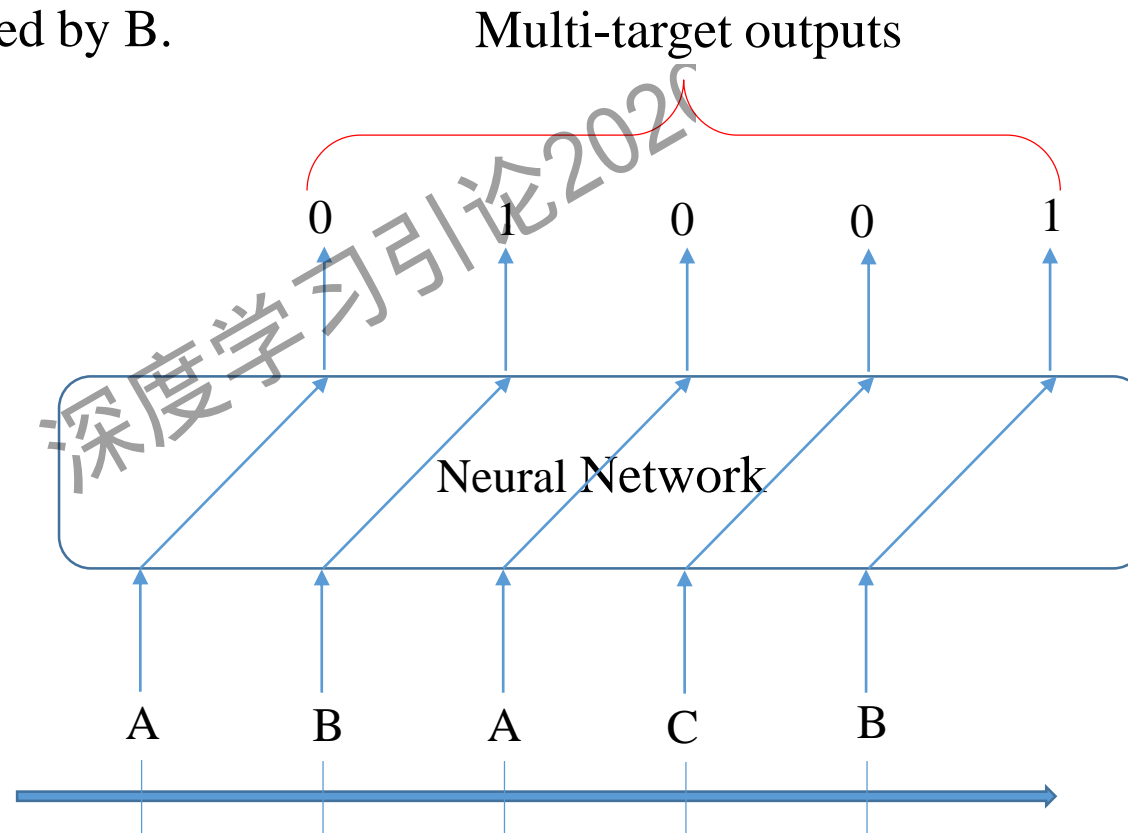
## *Recognize A followed by B Problem*

The task is to recognize A followed by B.

### Generated Sequences

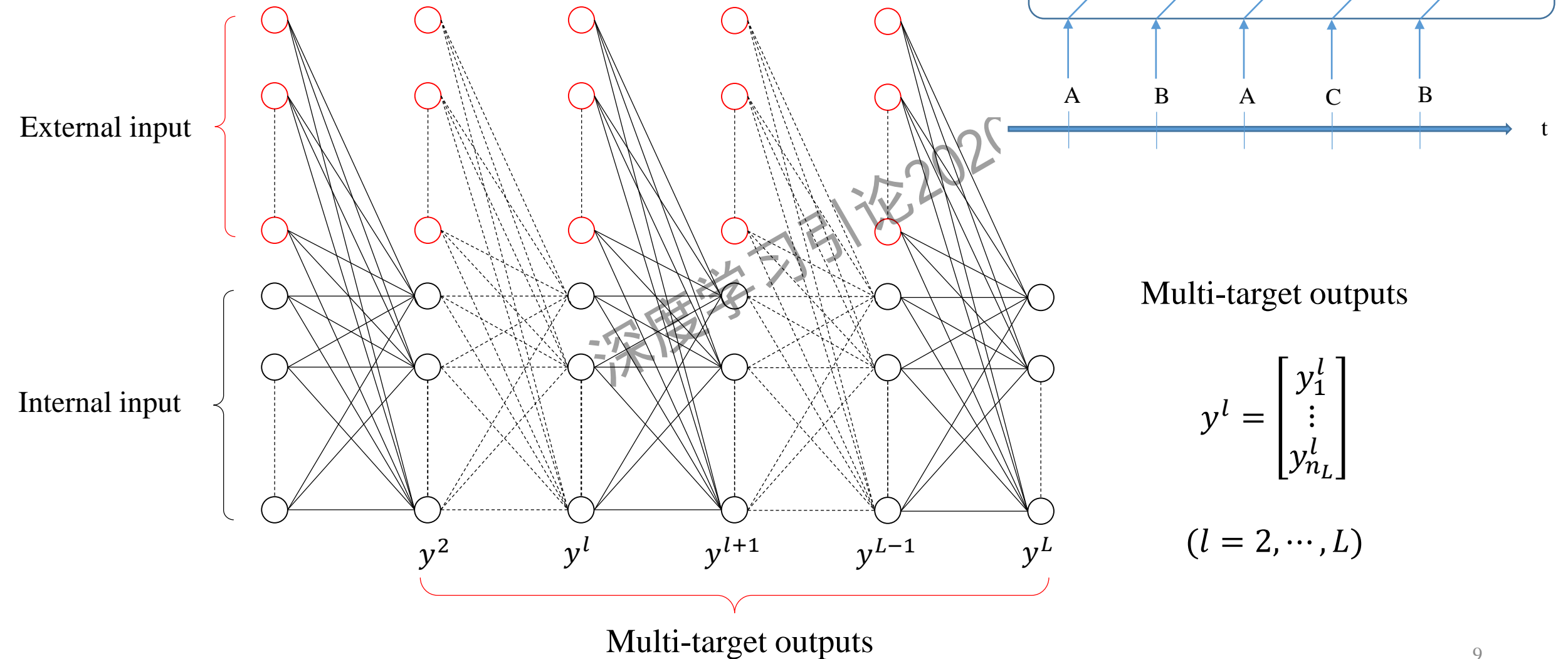
1. ABCAB
2. CCBBA
3. CACCB
4. ACCCB
5. CACBC
6. AAACB
7. BAACB
8. CCBAB
9. BCCAB
10. CABAC

.....

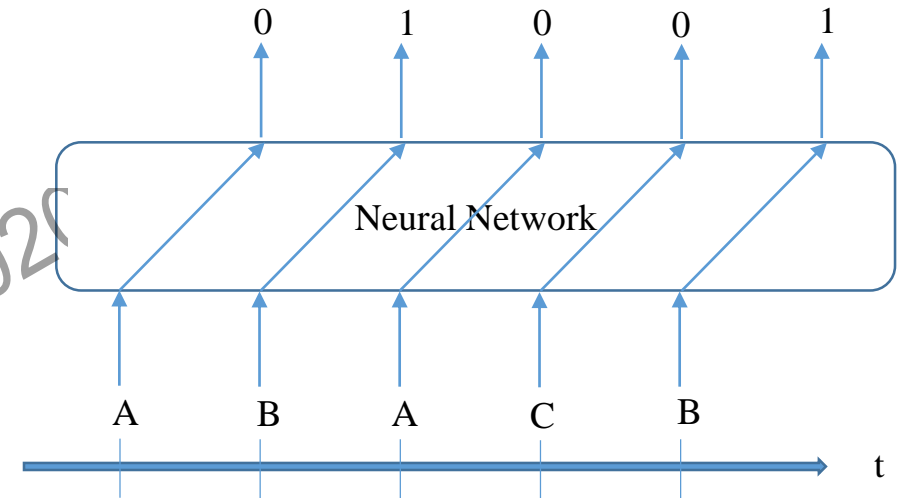
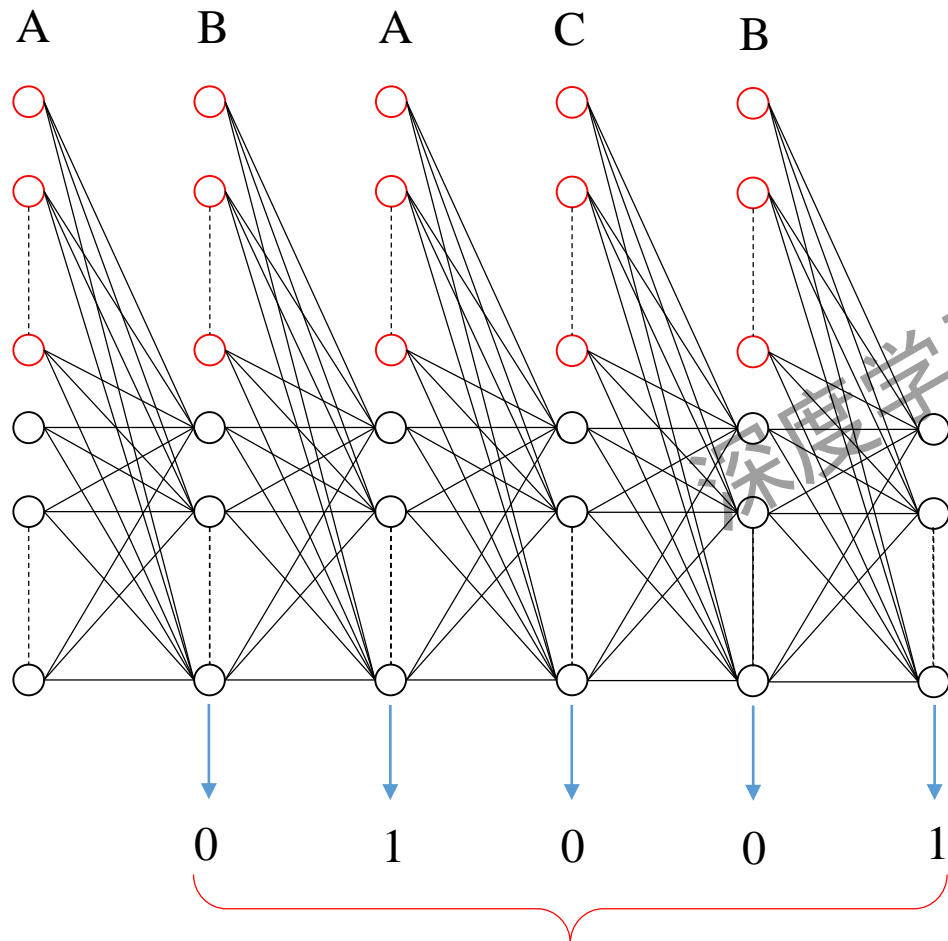




# A Sequence Recognizing Example

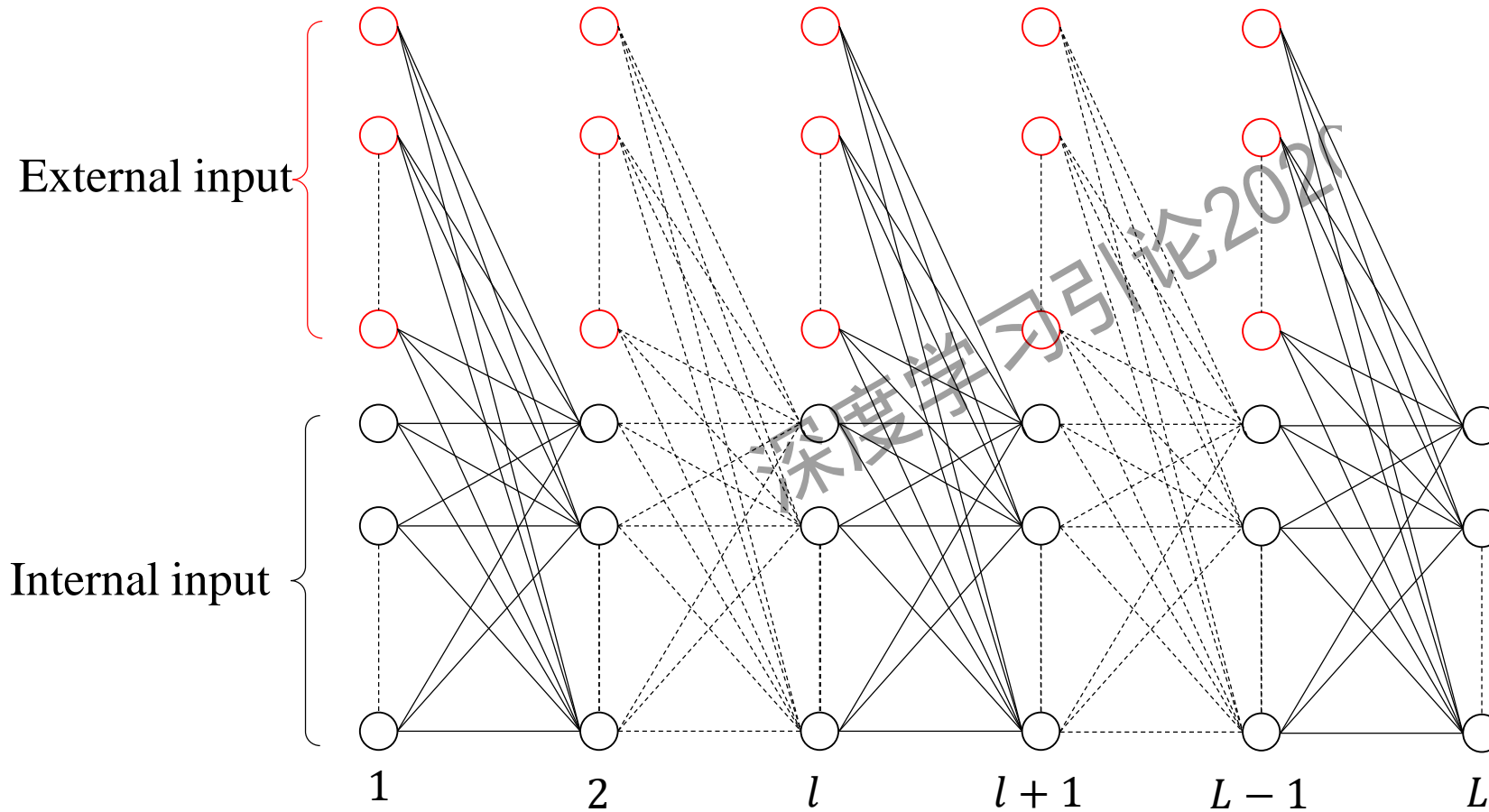


# A Sequence Recognizing Example



**Problem:**  
How to develop algorithm to train the network?

# A Sequence Recognizing Example



Mono-target output

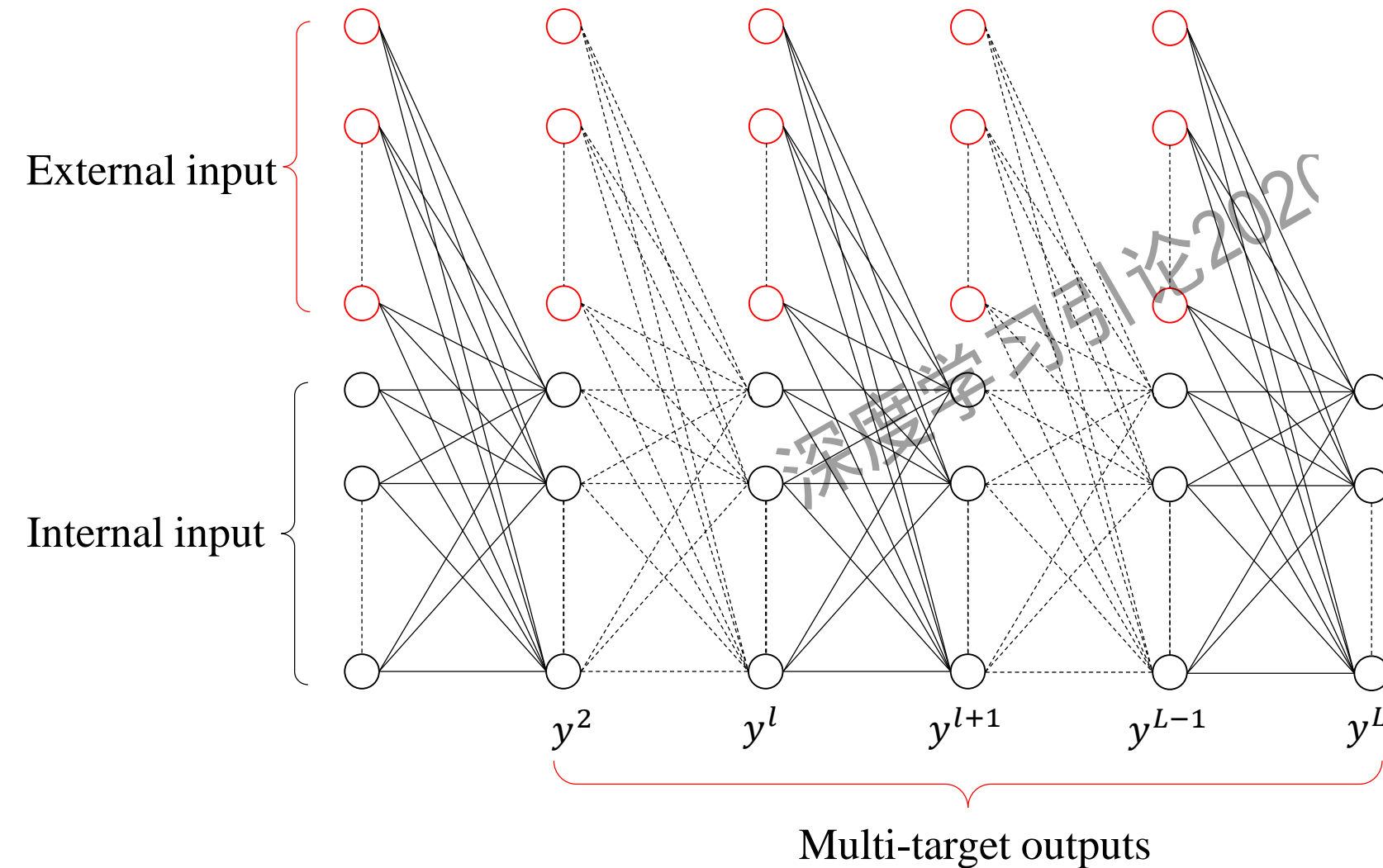
$$y^L = \begin{bmatrix} y_1^L \\ \vdots \\ y_{n_L}^L \end{bmatrix}$$

In the last layer only



We already developed BP Algorithm for mono-target output.

# A Sequence Recognizing Example



*Problem:*

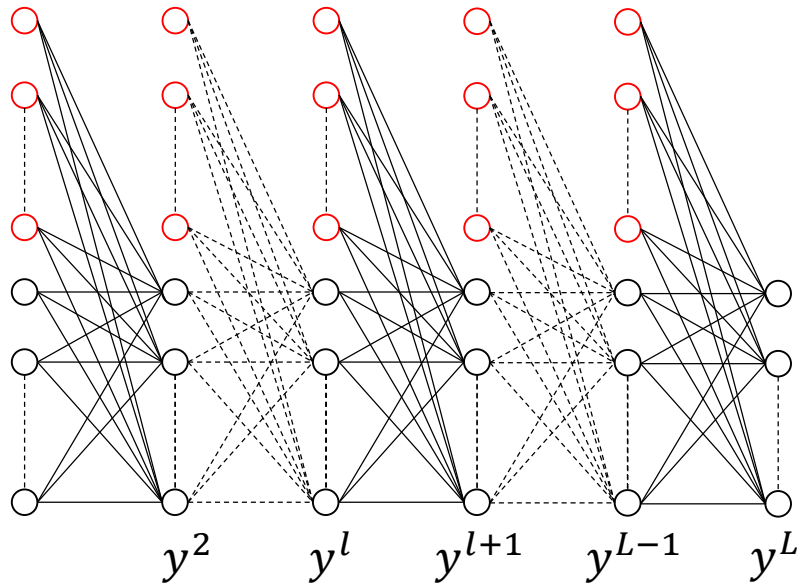
Can we develop learning algorithms for multi-target output similar to BP for mono-target output?

Multi-target outputs

$$y^l = \begin{bmatrix} y_1^l \\ \vdots \\ y_{n_L}^l \end{bmatrix}$$

$$(l = 2, \dots, L)$$

# A Sequence Recognizing Example



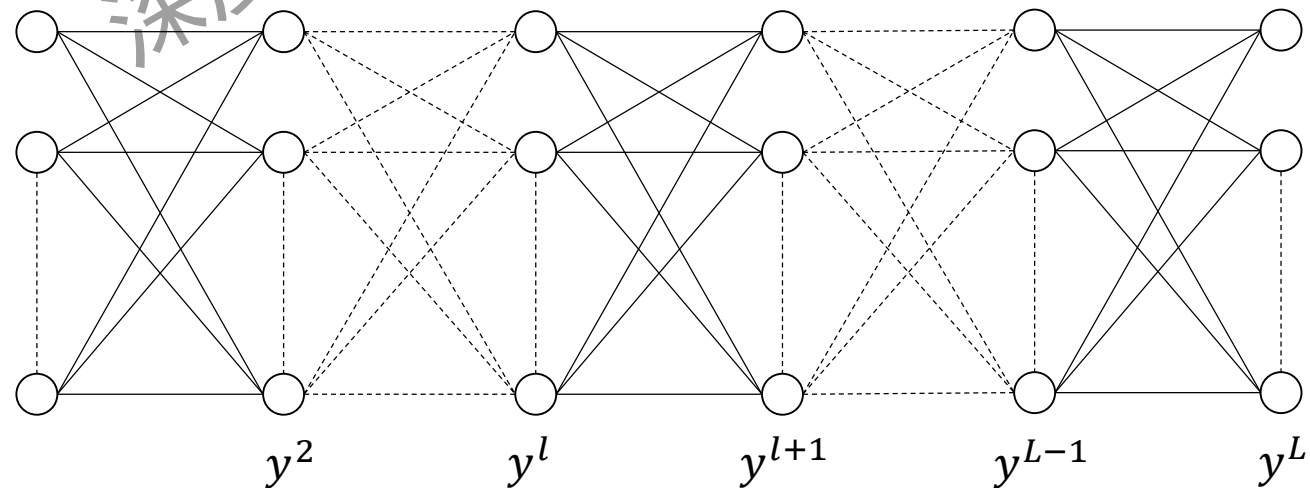
Combine the external and  
internal inputs in first layer

Equivalent

Multiple Target Outputs

$$y^l = \begin{bmatrix} y_1^l \\ \vdots \\ y_{n_L}^l \end{bmatrix}$$

$(l = 2, \dots, L)$



**Problem:**

Can we develop learning  
algorithms similar to BP?

# Outline

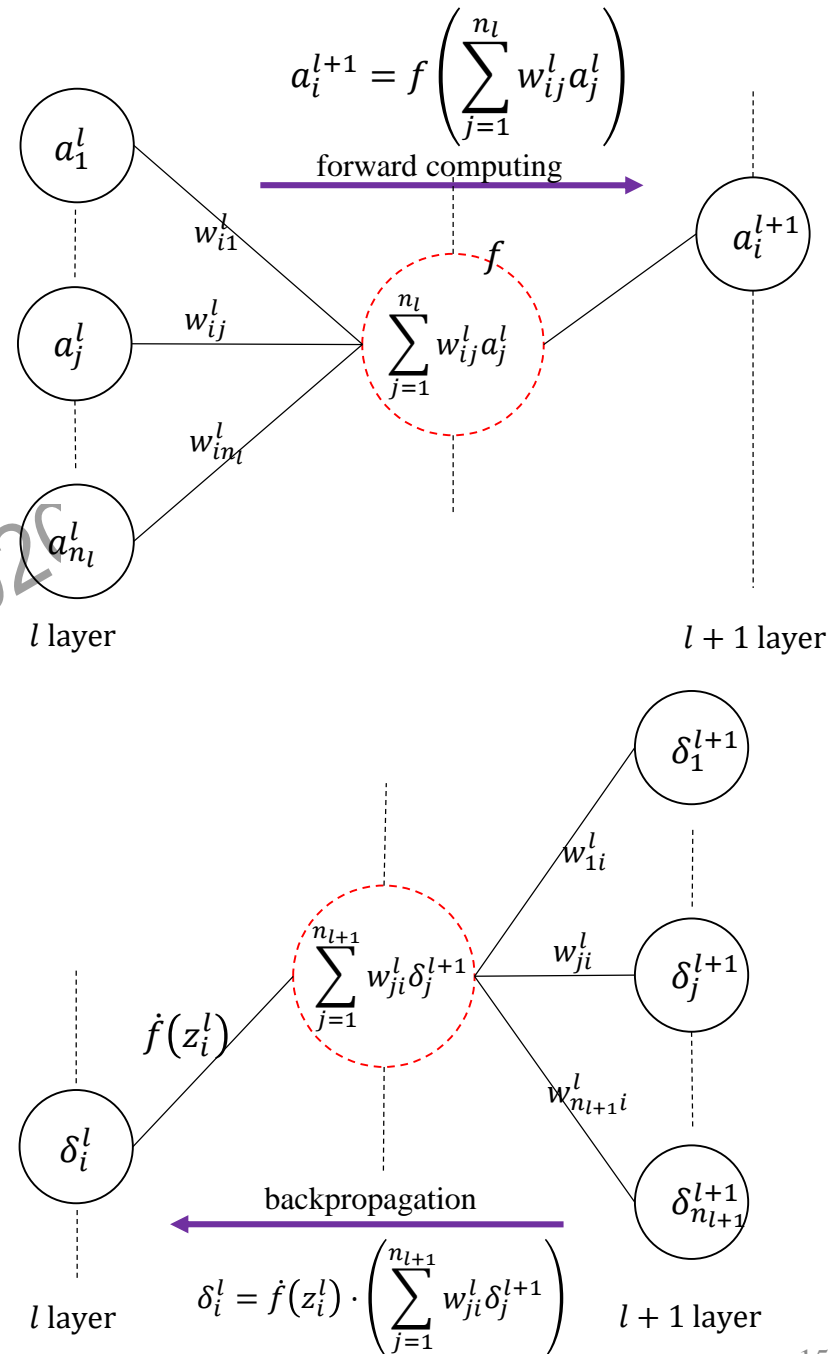
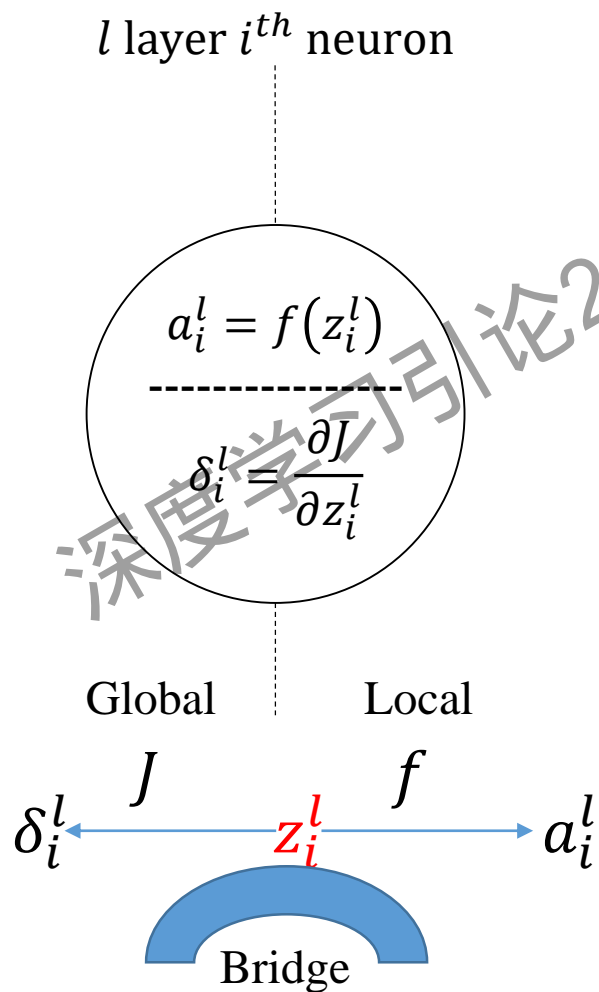
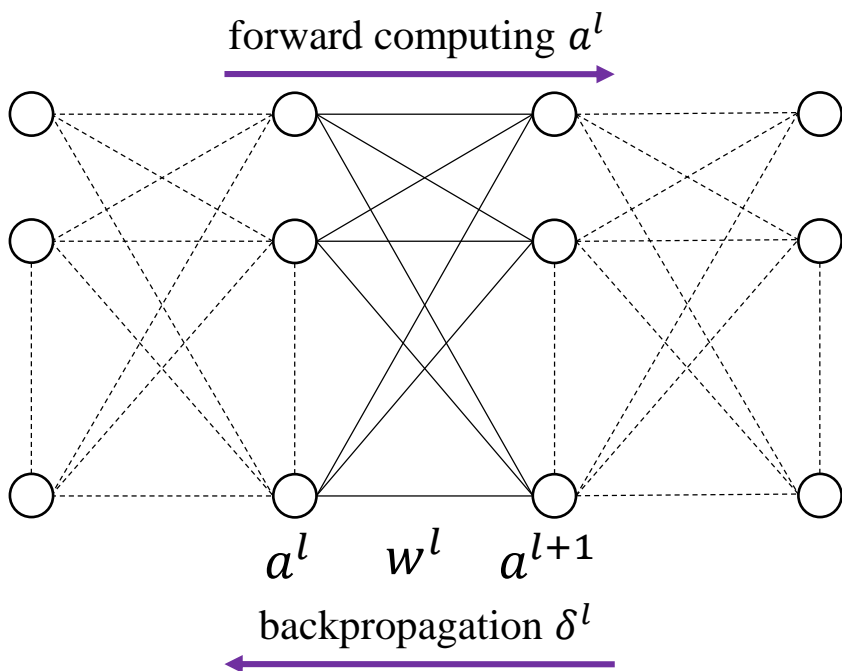
- A Sequence Recognizing Example
- Review of BP for Mono-target Output NNs
- BP Method for Multi-target Outputs NNs
- BP Algorithm for Multi-target Outputs NNs
- Illustrative Examples
- Assignment

# Only One Page to Understand BP

Cost function:  $J(w^1, \dots, w^{L-1})$

Updating rule:  $w_{ji}^l \leftarrow w_{ji}^l - \alpha \cdot \frac{\partial J}{\partial w_{ji}^l}$

Relationship:  $\frac{\partial J}{\partial w_{ji}^l} = \delta_j^{l+1} \cdot a_i^l$

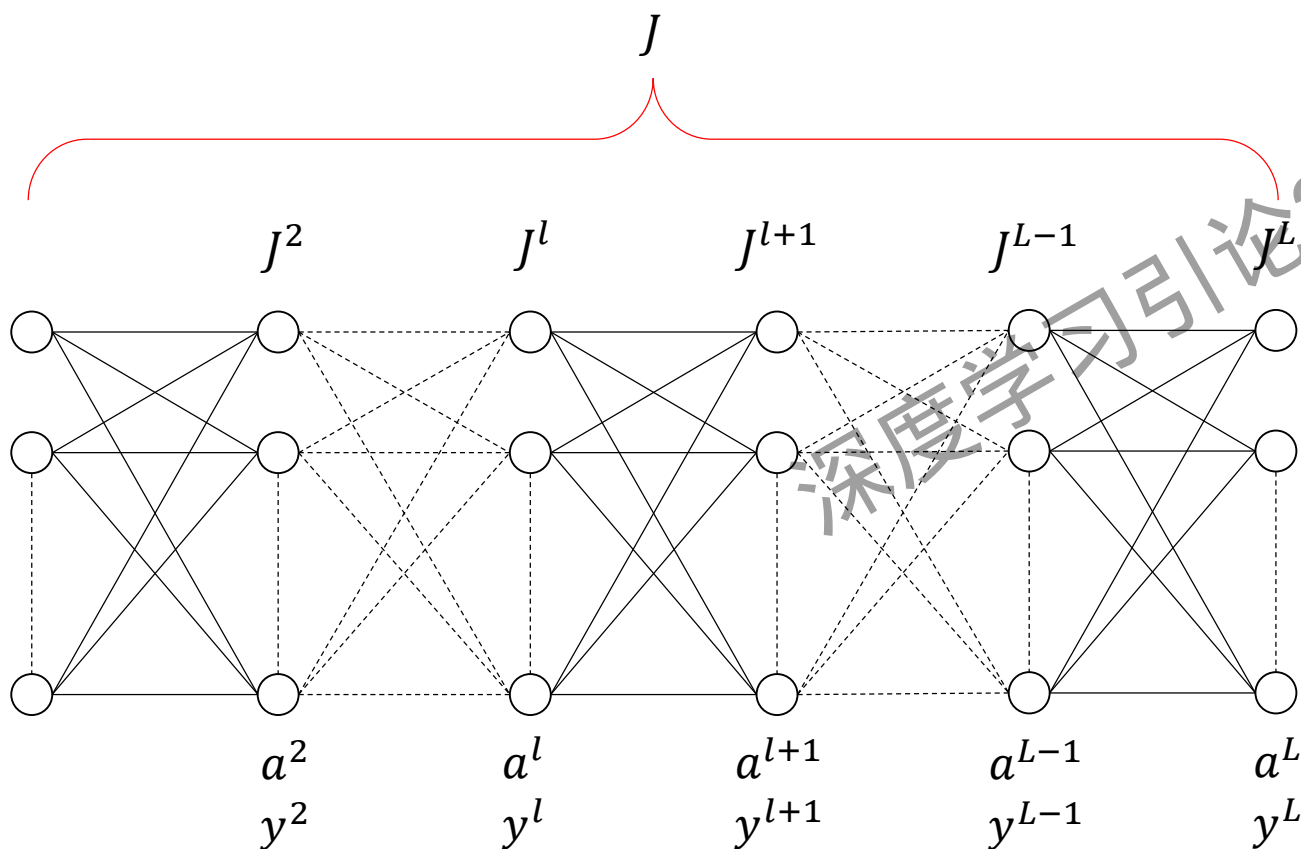


# Outline

- A Sequence Recognizing Example
- Review of BP for Mono-target Output NNs
- BP Method for Multi-target Outputs NNs
- BP Algorithm for Multi-target Outputs NNs
- Illustrative Examples
- Assignment



# BP Method for Multi-target Outputs NNs



Cost function

$$J^l = \frac{1}{2} \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2, (l = 2, \dots, L)$$

$$J = \sum_{l=2}^L J^l = \frac{1}{2} \sum_{l=2}^L \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2$$

Network Outputs      Multi-target outputs

$$a^l = \begin{bmatrix} a_1^l \\ \vdots \\ a_{n_l}^l \end{bmatrix}$$

$$(l = 2, \dots, L)$$

$$y^l = \begin{bmatrix} y_1^l \\ \vdots \\ y_{n_l}^l \end{bmatrix}$$

$$(l = 2, \dots, L)$$

# BP Method for Multi-target Outputs NNs

Steepest Descent Method

$$J^l = \frac{1}{2} \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2, (l = 2, \dots, L)$$

$$J = \sum_{l=2}^L J^l = \frac{1}{2} \sum_{l=2}^L \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2$$

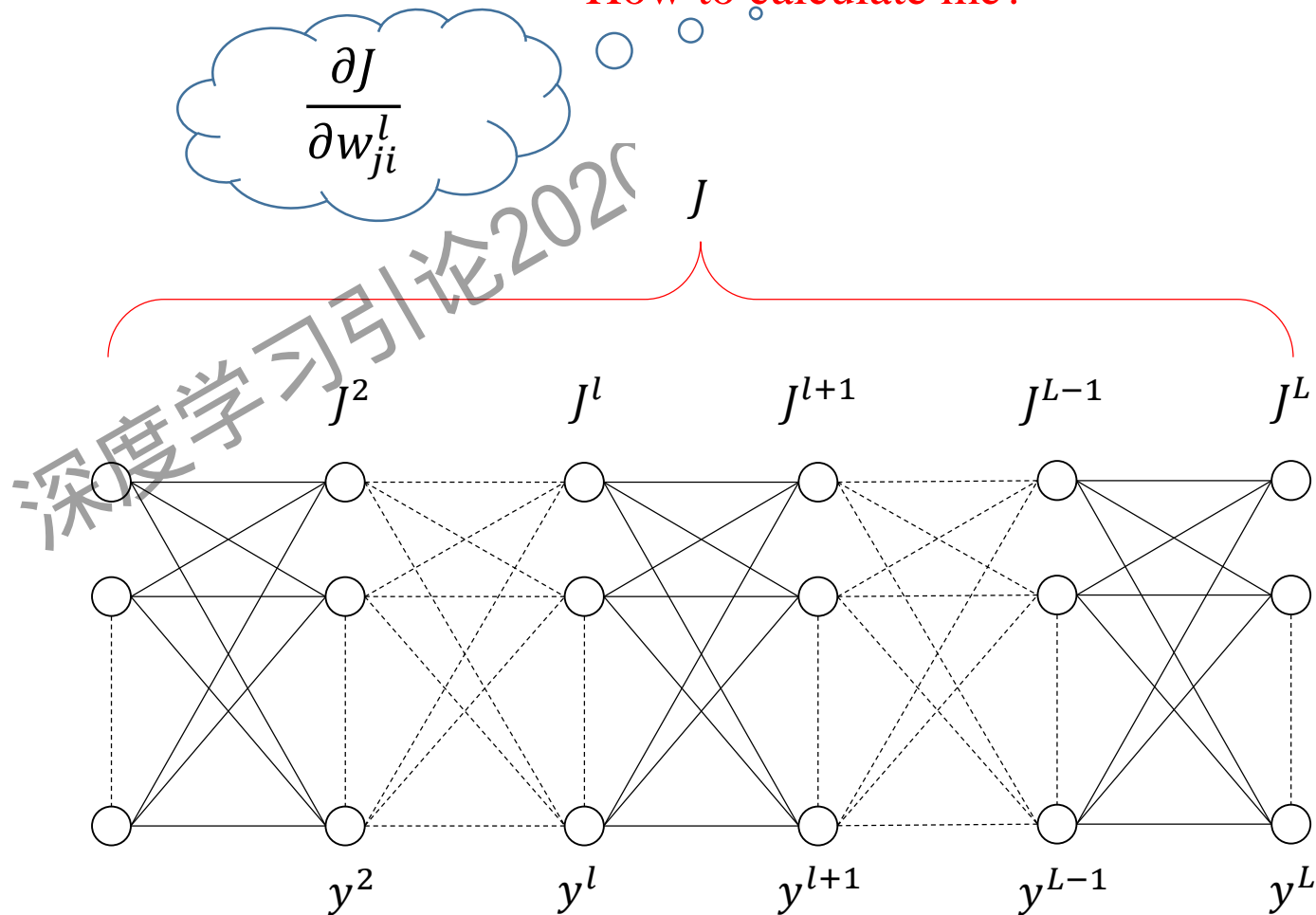
1. Computing

$$\frac{\partial J}{\partial w_{ji}^l}$$

2. Iterating

$$w_{ji}^l \leftarrow w_{ji}^l - \alpha \cdot \frac{\partial J}{\partial w_{ji}^l}$$

How to calculate me?



$l$  layer

$$a_i^l = f(z_i^l)$$

define  $\delta_i^l = \frac{\partial J}{\partial z_i^l}$

Relation between  $\delta_i^l$  and  $\frac{\partial J}{\partial w_{ji}^l}$

$$\frac{\partial J}{\partial w_{ji}^l} = \delta_j^{l+1} \cdot a_i^l$$

Why?

$$\frac{\partial J}{\partial w_{ji}^l} = \frac{\partial J}{\partial z_j^{l+1}} \cdot \frac{\partial z_j^{l+1}}{\partial w_{ji}^l} = \delta_j^{l+1} \cdot a_i^l$$

$$J^l = \frac{1}{2} \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2, (l = 2, \dots, L)$$

$$J = \sum_{l=2}^L J^l = \frac{1}{2} \sum_{l=2}^L \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2$$

$l + 1$  layer

$$\delta_j^{l+1} = \frac{\partial J}{\partial z_j^{l+1}}$$

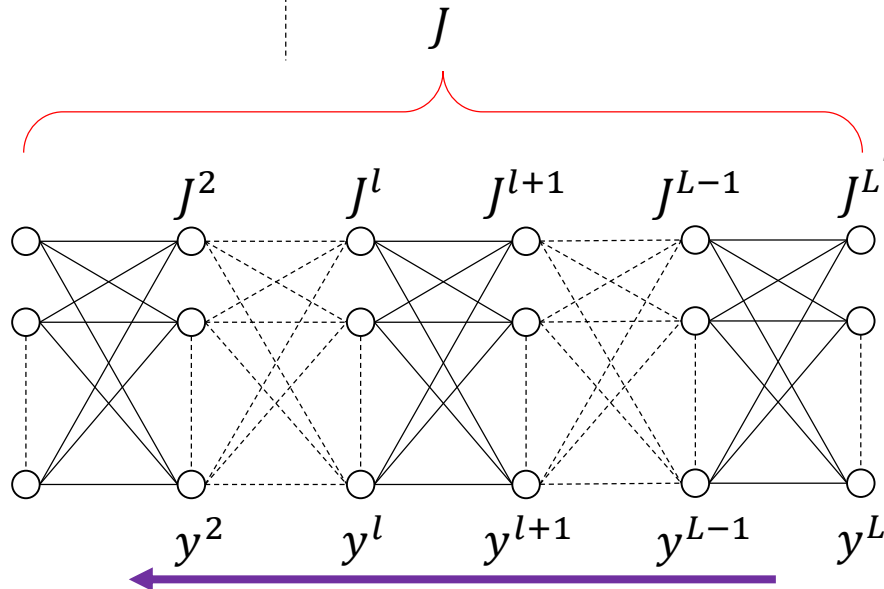
$$z_j^{l+1} = \sum_{i=1}^{n_l} w_{ji}^l a_i^l$$

$$\frac{\partial J}{\partial w_{ji}^l} = \delta_j^{l+1} \cdot a_i^l$$

$w_{ji}^l$

$$a_i^l = f(z_i^l)$$

$$\delta_i^l = \frac{\partial J}{\partial z_i^l}$$



Problem: Can we back propagate  $\delta^l$ ?

# Step 1: Calculating $\delta^L$ in Last Layer

$$J^l = \frac{1}{2} \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2, (l = 2, \dots, L)$$

$$J = \sum_{l=2}^L J^l = \frac{1}{2} \sum_{l=2}^L \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2$$

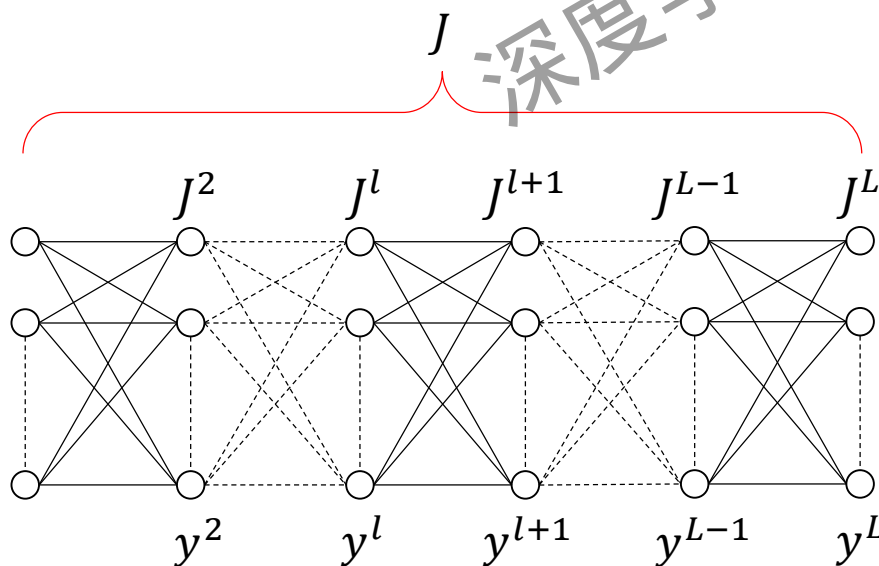
It holds that,

$$\delta_i^L = \frac{\partial J}{\partial z_i^L} = \frac{\partial J^L}{\partial z_i^L} = \frac{1}{2} \cdot \frac{\partial (a_i^L - y_i^L)^2}{\partial z_i^L} = (a_i^L - y_i^L) \cdot \frac{\partial a_i^L}{\partial z_i^L} = (a_i^L - y_i^L) \cdot f'(z_i^L)$$

$$a_i^L = f(z_i^L)$$

$$a_i^L = f(z_i^L)$$

$$\delta_i^L = \frac{\partial J}{\partial z_i^L}$$



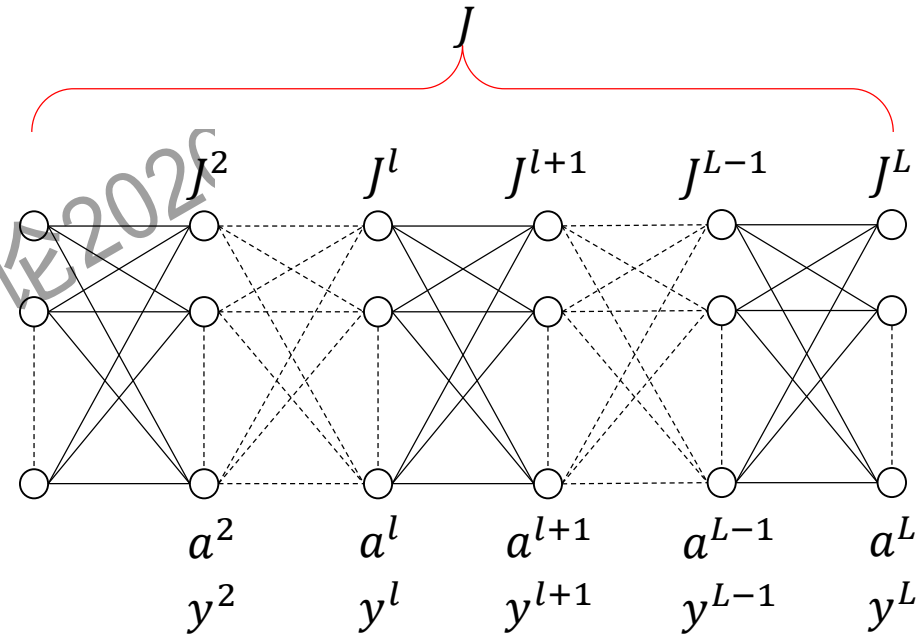
## Step 2: Relation Between $\delta^l$ and $\delta^{l+1}$

$$J^l = \frac{1}{2} \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2, (l = 2, \dots, L)$$

$$J = \sum_{l=2}^L J^l = \frac{1}{2} \sum_{l=2}^L \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2$$

$J$  may have an explicit dependence on  $z_i^l$ , it may also have an implicit dependence on  $z_i^l$  through later output values. To avoid ambiguity in interpreting partial derivatives, define  $z_i^l(*) = z_i^l$ .

$$\delta_i^l = \frac{\partial J}{\partial z_i^l} = \frac{\partial J}{\partial z_i^l(*)} \cdot \frac{\partial z_i^l(*)}{\partial z_i^l} + \sum_{j=1}^{n_{l+1}} \frac{\partial J}{\partial z_j^{l+1}} \cdot \frac{\partial z_j^{l+1}}{\partial z_i^l}$$



## Step 2: Relation Between $\delta^l$ and $\delta^{l+1}$

$$J^l = \frac{1}{2} \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2, (l = 2, \dots, L)$$

$$J = \sum_{l=2}^L J^l = \frac{1}{2} \sum_{l=2}^L \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2$$

$J$  may have an explicit dependence on  $z_i^l$ , it may also have an implicit dependence on  $z_i^l$  through later output values. To avoid ambiguity in interpreting partial derivatives, define  $z_i^l(*) = z_i^l$ .

$$\delta_i^l = \frac{\partial J}{\partial z_i^l} = \frac{\partial J}{\partial z_i^l(*)} \cdot \frac{\partial z_i^l(*)}{\partial z_i^l} + \sum_{j=1}^{n_{l+1}} \frac{\partial J}{\partial z_j^{l+1}} \cdot \frac{\partial z_j^{l+1}}{\partial z_i^l}$$

An Illustrate Example

$$J = x + y, y = \exp(x)$$

$$\frac{\partial J}{\partial x} = \frac{\partial J}{\partial x} + \frac{\partial J}{\partial y} \cdot \frac{\partial y}{\partial x}$$

$$x^* = x$$

$$\frac{\partial J}{\partial x} = \frac{\partial J}{\partial x^*} \cdot \frac{\partial x^*}{\partial x} + \frac{\partial J}{\partial y} \cdot \frac{\partial y}{\partial x}$$

## Step 2: Relation Between $\delta^l$ and $\delta^{l+1}$

$$\delta_i^l = \frac{\partial J}{\partial z_i^l} = \frac{\partial J}{\partial z_i^l(*)} \cdot \frac{\partial z_i^l(*)}{\partial z_i^l} + \sum_{j=1}^{n_{l+1}} \frac{\partial J}{\partial z_j^{l+1}} \cdot \frac{\partial z_j^{l+1}}{\partial z_i^l}$$

$$\frac{\partial J}{\partial z_i^l(*)} \cdot \frac{\partial z_i^l(*)}{\partial z_i^l} = \frac{\partial J^l}{\partial z_i^l} = \frac{1}{2} \cdot \frac{\partial (a_i^l - y_i^l)^2}{\partial z_i^l} = (a_i^l - y_i^l) \cdot \frac{\partial a_i^l}{\partial z_i^l} = (a_i^l - y_i^l) \cdot f'(z_i^l)$$

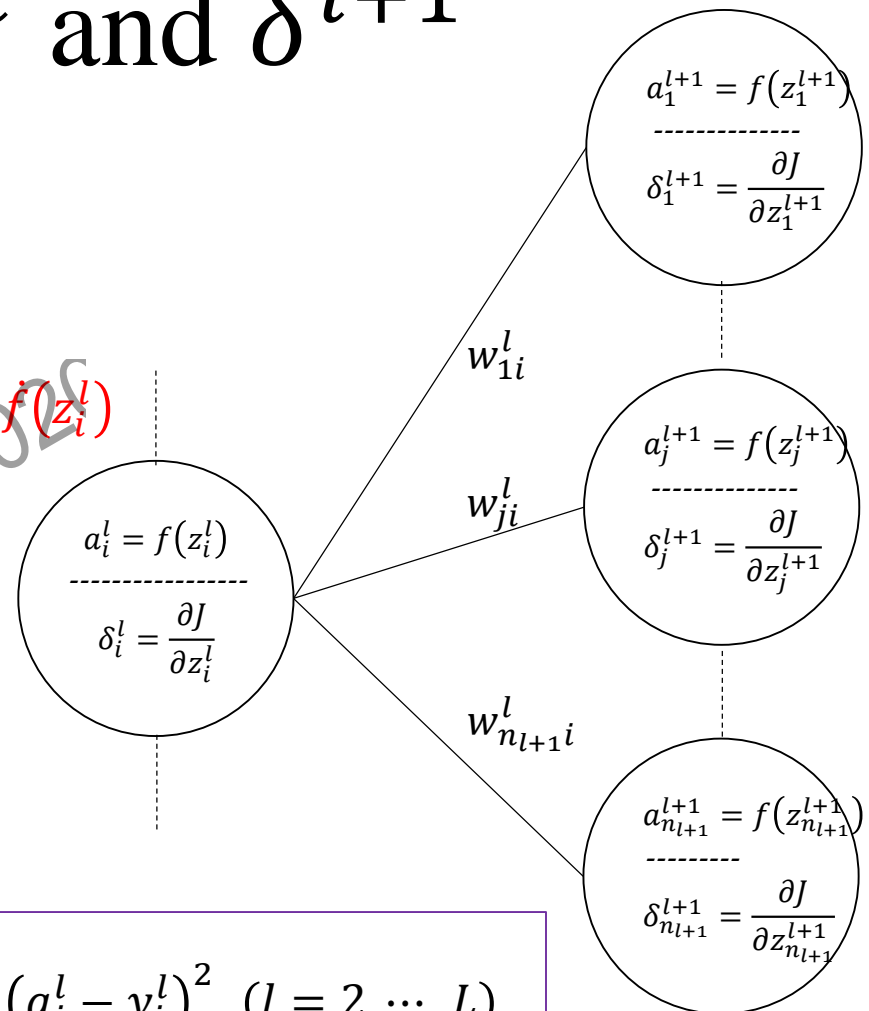
$$\sum_{j=1}^{n_{l+1}} \frac{\partial J}{\partial z_j^{l+1}} \cdot \frac{\partial z_j^{l+1}}{\partial z_i^l} = \sum_{j=1}^{n_{l+1}} \delta_j^{l+1} \cdot \frac{\partial z_j^{l+1}}{\partial z_i^l} = f'(z_i^l) \cdot \left( \sum_{j=1}^{n_{l+1}} \delta_j^{l+1} \cdot w_{ji}^l \right)$$

$$\begin{aligned} z_j^{l+1} &= \sum_{i=1}^{n_l} w_{ji}^l a_i^l \\ a_i^{l+1} &= f(z_i^{l+1}) \end{aligned}$$

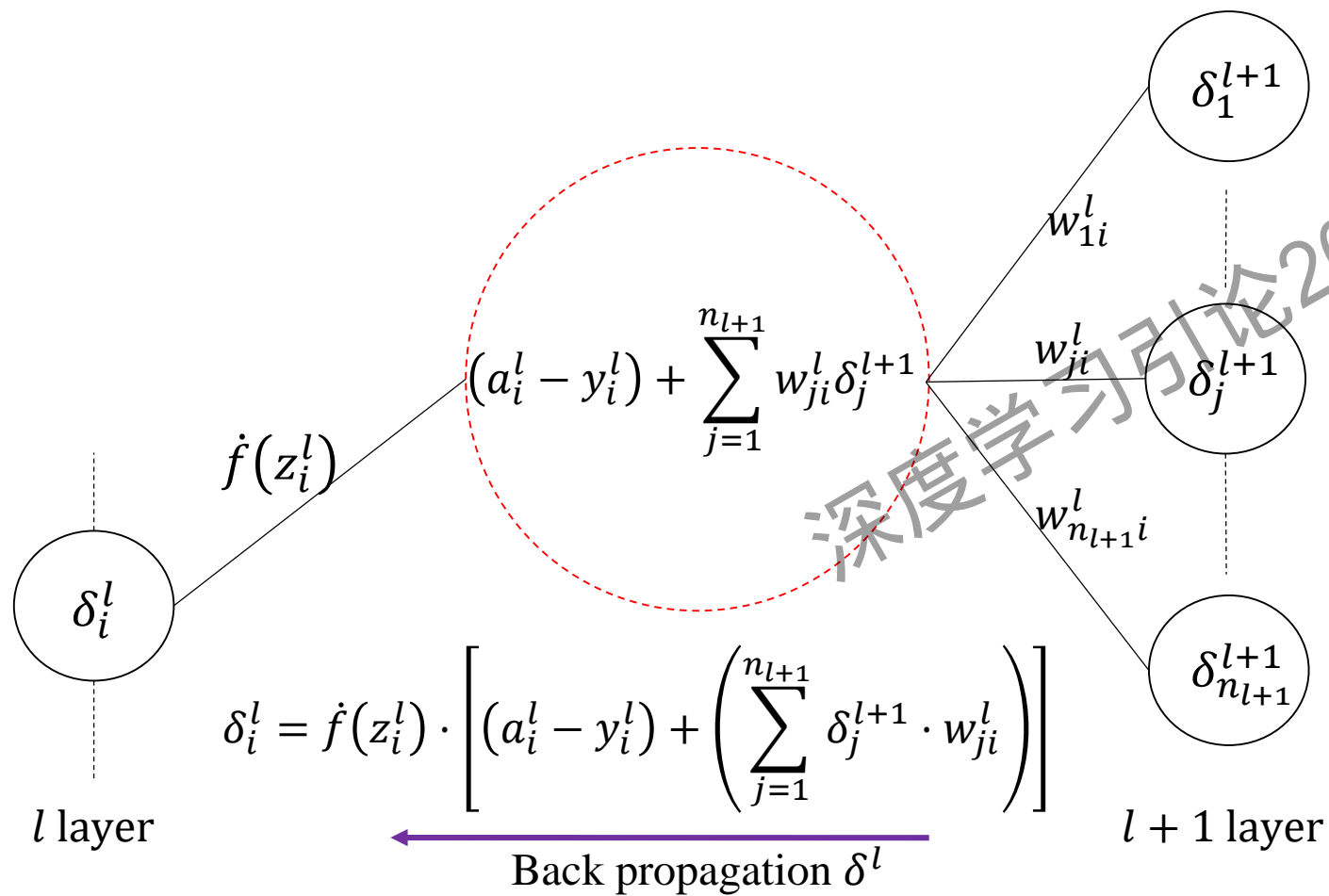
$$\frac{\partial z_j^{l+1}}{\partial z_i^l} = w_{ji}^l \cdot \frac{\partial a_j^{l+1}}{\partial z_i^l} = w_{ji}^l \cdot f'(z_i^l)$$

$$\delta_i^l = f'(z_i^l) \cdot \left[ (a_i^l - y_i^l) + \left( \sum_{j=1}^{n_{l+1}} \delta_j^{l+1} \cdot w_{ji}^l \right) \right]$$

$$\begin{aligned} J^l &= \frac{1}{2} \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2, (l = 2, \dots, L) \\ J &= \sum_{l=2}^L J^l = \frac{1}{2} \sum_{l=2}^L \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2 \end{aligned}$$

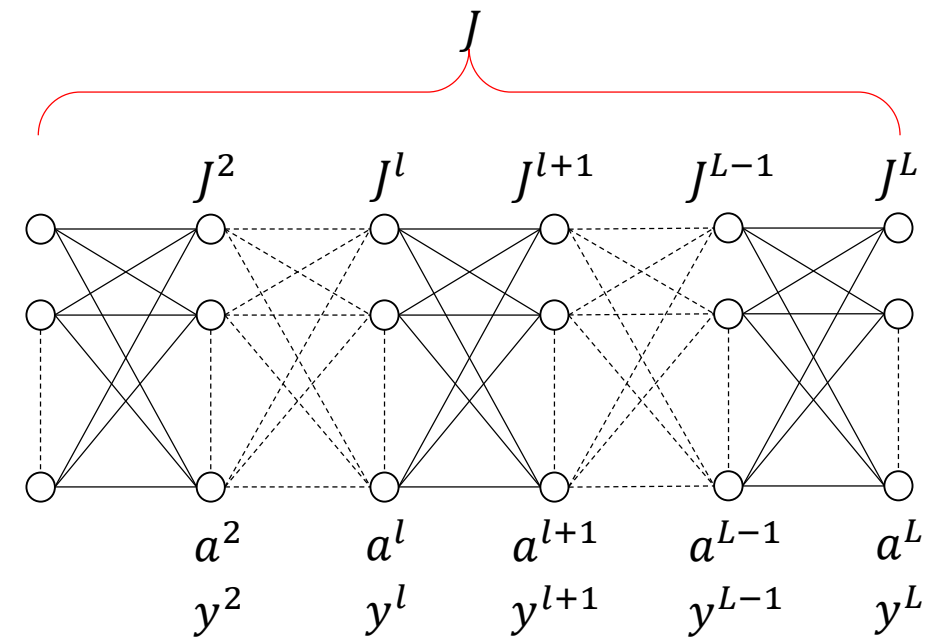


# Step 3: Backpropagation $\delta^l$



$$J^l = \frac{1}{2} \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2, (l = 2, \dots, L)$$

$$J = \sum_{l=2}^L J^l = \frac{1}{2} \sum_{l=2}^L \sum_{i=1}^{n_l} (a_i^l - y_i^l)^2$$



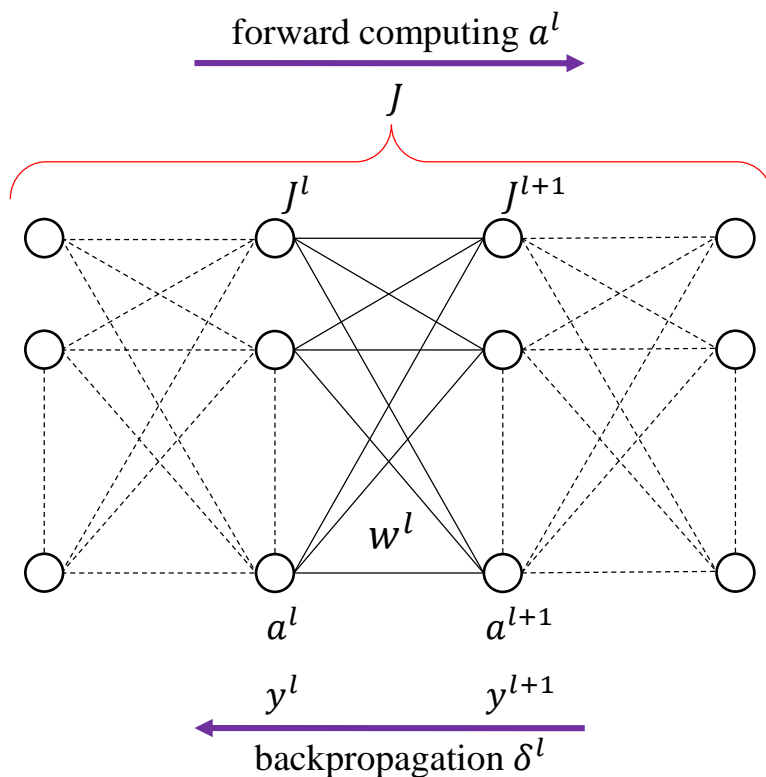


# Only One Page to Understand BP

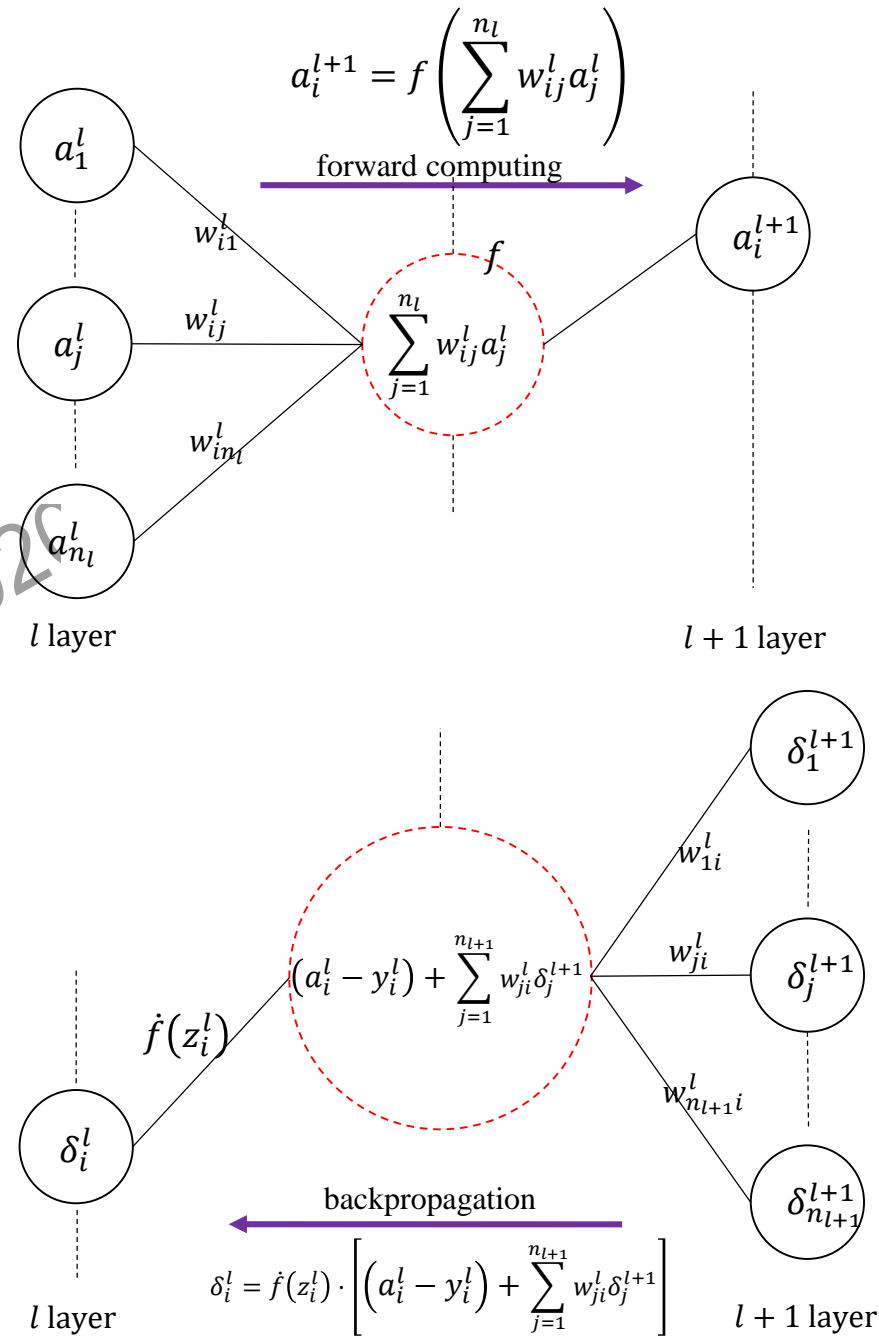
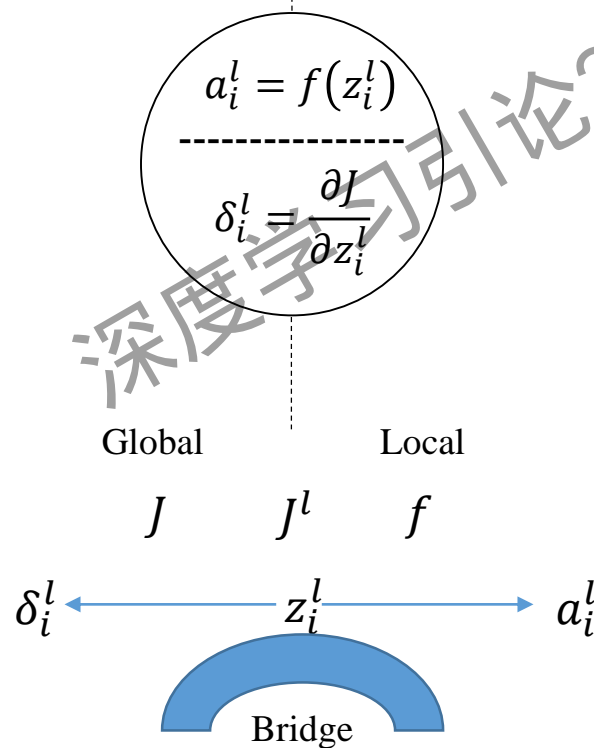
Cost function:  $J(w^1, \dots, w^{L-1})$

Updating rule:  $w_{ji}^l \leftarrow w_{ji}^l - \alpha \cdot \frac{\partial J}{\partial w_{ji}^l}$

Relationship:  $\frac{\partial J}{\partial w_{ji}^l} = \delta_j^{l+1} \cdot a_i^l$



$l$  layer  $i^{th}$  neuron



# Outline

- A Sequence Recognizing Example
- Review of BP for Mono-target Output NNs
- BP Method for Multi-target Outputs NNs
- BP Algorithm for Multi-target Outputs NNs
- Illustrative Examples
- Assignment

# The BP Algorithm

Step 1. Input the training data set  $D = \{(x, y)\}$

Step 2. Initial each  $w_{ij}^l$ , and choose a learning rate  $\alpha$ .

Step 3. For each mini-batch sample  $D_m \subseteq D$

for each  $x \in D_m$

$a^1 \leftarrow x$ ;

for  $l = 2:L$

$a^l \leftarrow fc(w^l, a^l)$ ;

end

$\delta^L = \frac{\partial J}{\partial z^L}$ ;

for  $l = L - 1:2$

$\delta^l \leftarrow bc(w^l, \delta^{l+1})$ ;

end

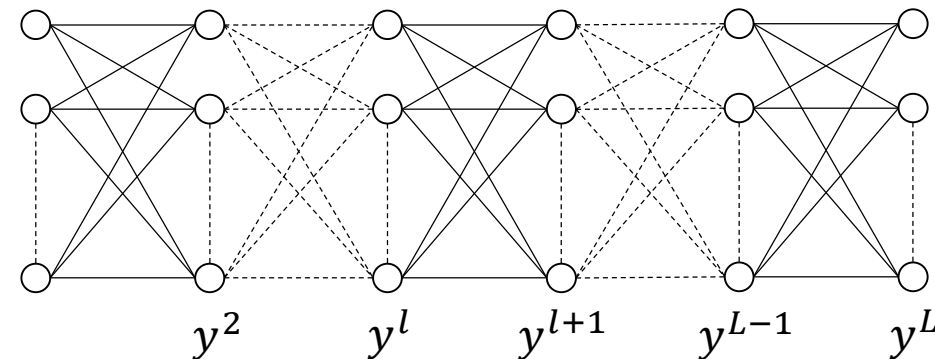
$\frac{\partial J}{\partial w_{ji}^l} \leftarrow \frac{\partial J}{\partial w_{ji}^l} + \delta_j^{l+1} \cdot a_i^l$ ;

end

Step 4. Updating

$w_{ji}^l \leftarrow w_{ji}^l - \alpha \cdot \frac{\partial J}{\partial w_{ji}^l}$ ;

Step 5. Return to Step 3 until each  $w^l$  converge.



function  $fc(w^l, a^l)$

for  $i = 1:n_{l+1}$

$$z_i^{l+1} = \sum_{j=1}^{n_l} w_{ij}^l a_j^l$$

$$a_i^{l+1} = f(z_i^{l+1})$$

end

Relationship:

$$\frac{\partial J}{\partial w_{ji}^l} = \delta_j^{l+1} \cdot a_i^l$$

function  $bc(w^l, \delta^{l+1})$

for  $i = 1:n_l$

$$\delta_i^l = f'(z_i^l) \cdot \left[ (a_i^l - y_i^l) + \left( \sum_{j=1}^{n_{l+1}} \delta_j^{l+1} \cdot w_{ji}^l \right) \right]$$

end

# Outline

- A Sequence Recognizing Example
- Review of BP for Mono-target Output NNs
- BP Method for Multi-target Outputs NNs
- BP Algorithm for Multi-target Outputs NNs
- Illustrative Examples
- Assignment

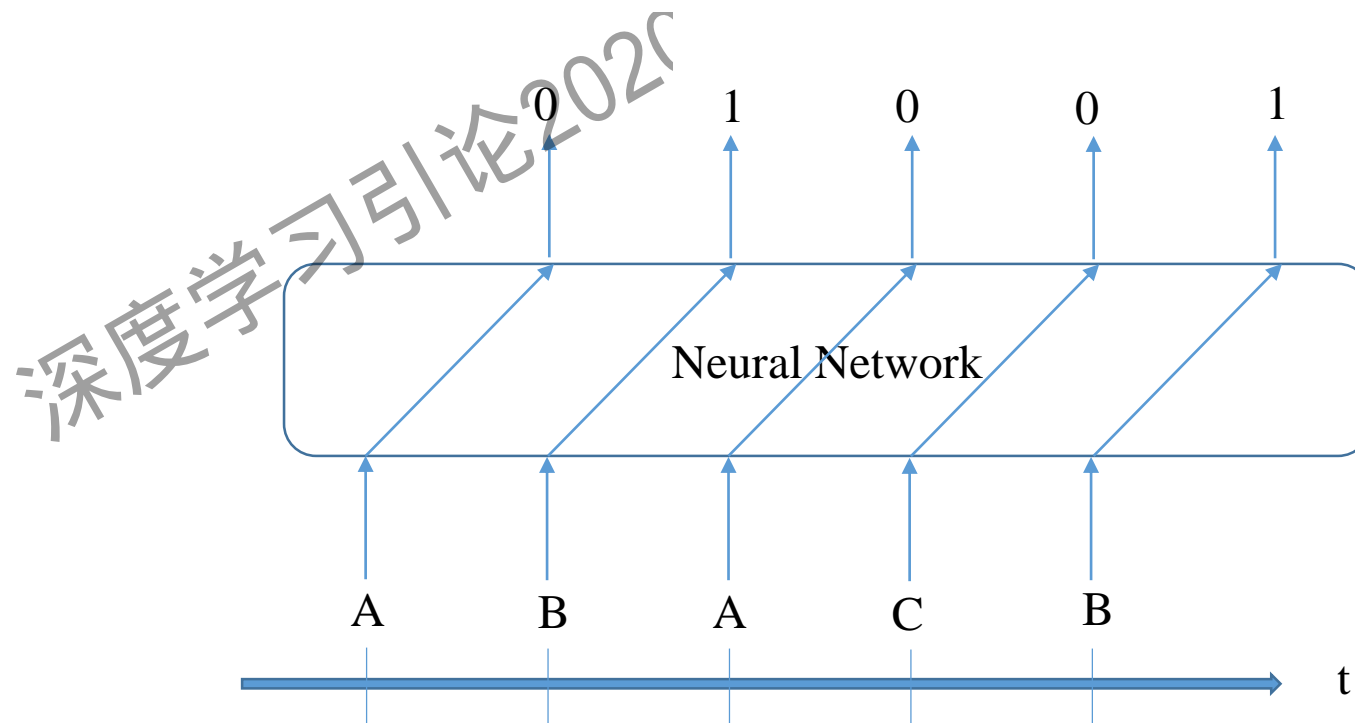
# Illustrative Example: Sequence Recognition

## *Recognize A followed by B Problem*

The task is to recognize A followed by B.

### Generated Sequences

1. ABACB
2. CCBBA
3. CACCB
4. ACCCB
5. CACBC
6. AAACB
7. BAACB
8. CCBAB
9. BCCAB
10. CABAC
- .....



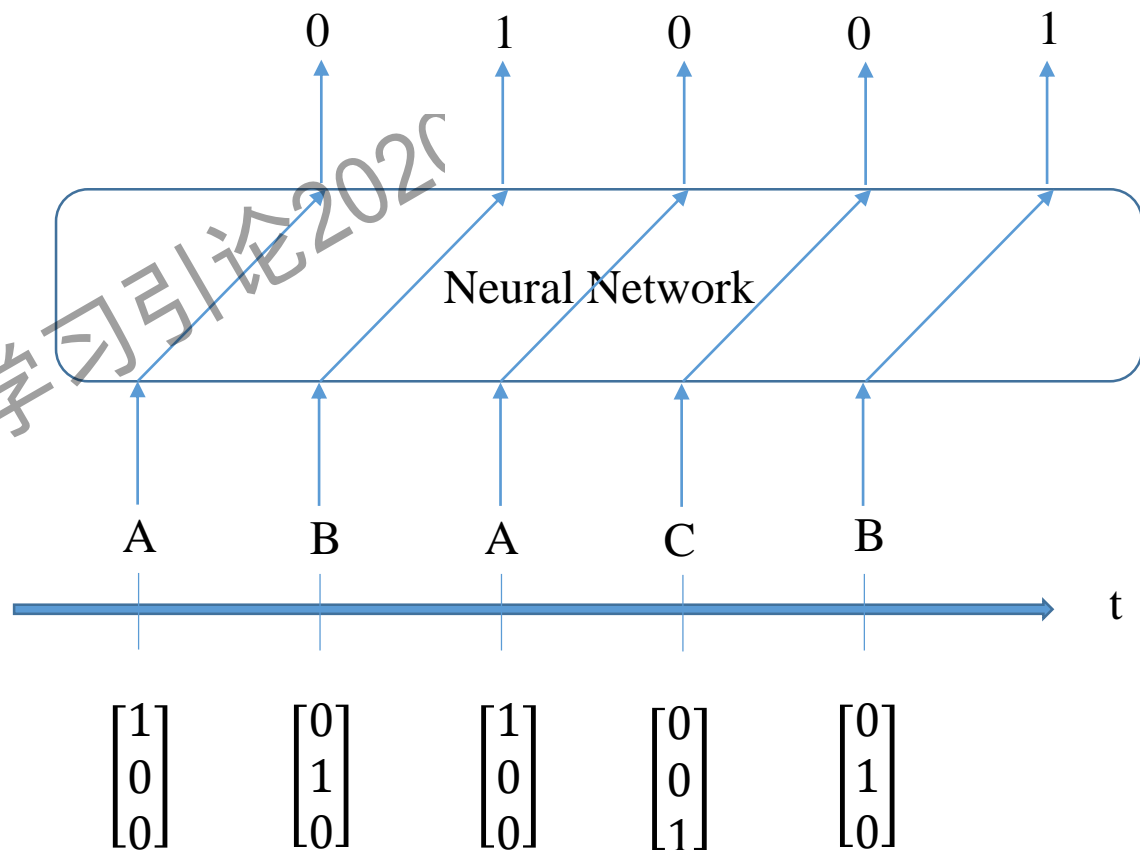
# Illustrative Example: Sequence Recognition

Coding the Inputs:

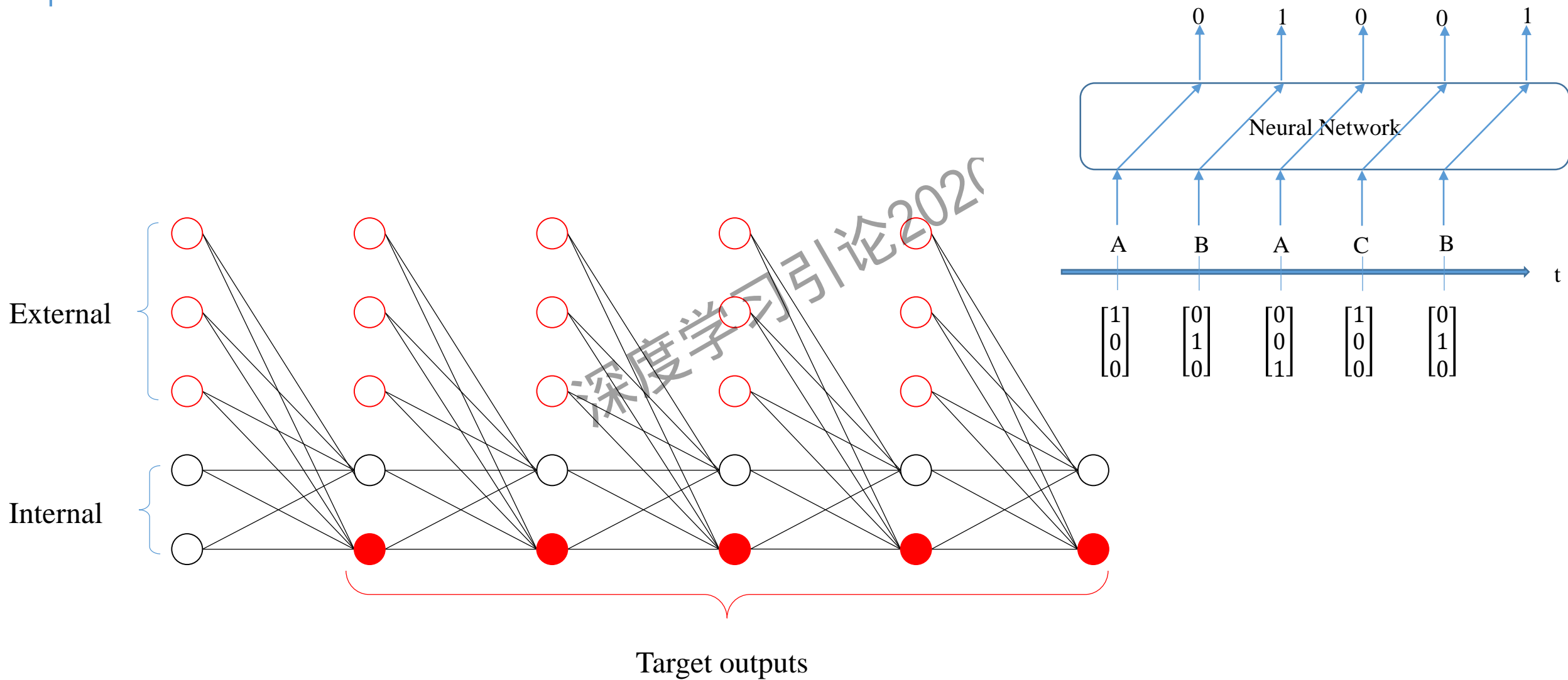
$$A = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad B = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad C = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

Generated Sequences

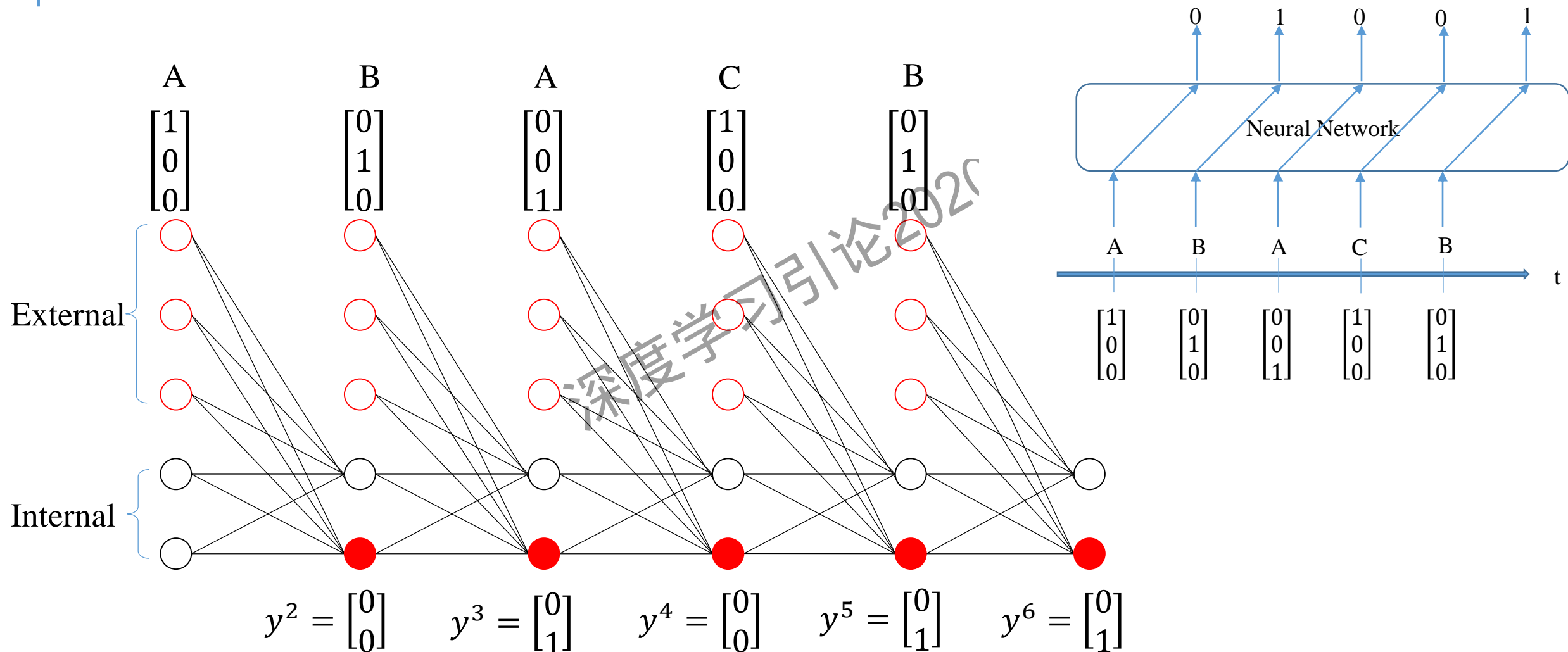
- |   |   |   |   |   |
|---|---|---|---|---|
| A   | B   | C   | A   | B   |
| $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ | $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$ | $\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$ | $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ | $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$ |
- |   |   |   |   |   |
|---|---|---|---|---|
| C   | A   | C   | C   | B   |
| $\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$ | $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ | $\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$ | $\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$ | $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$ |



# Illustrative Example: Sequence Recognition



# Illustrative Example: Sequence Recognition



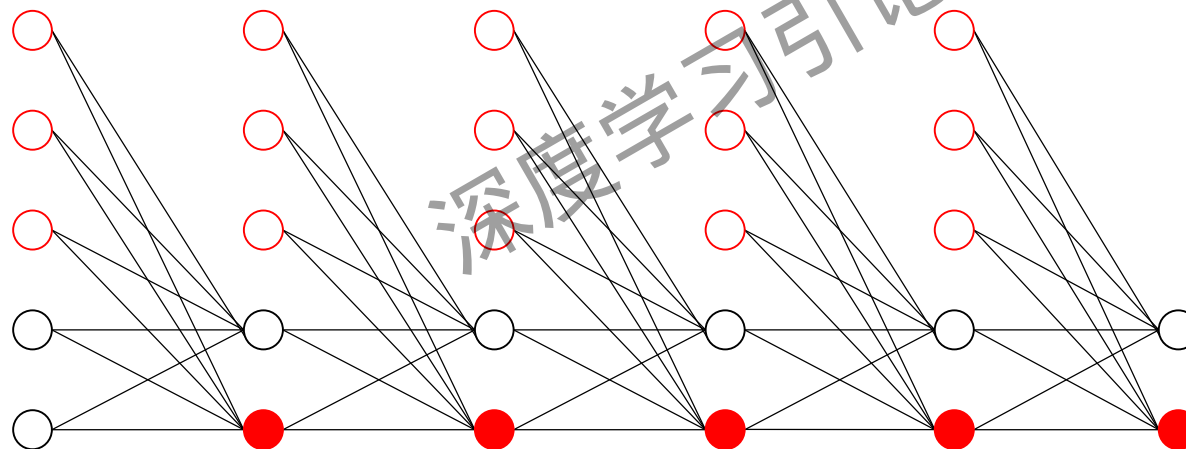


# Illustrative Example: Sequence Recognition

## Generated Training Sequences

1. ABCAB
2. CCBBA
3. CACCB
4. ACCCB
5. CACBC
6. AAACB
7. BAACB
8. CCBAB
9. BCCAB
10. CABAC

.....



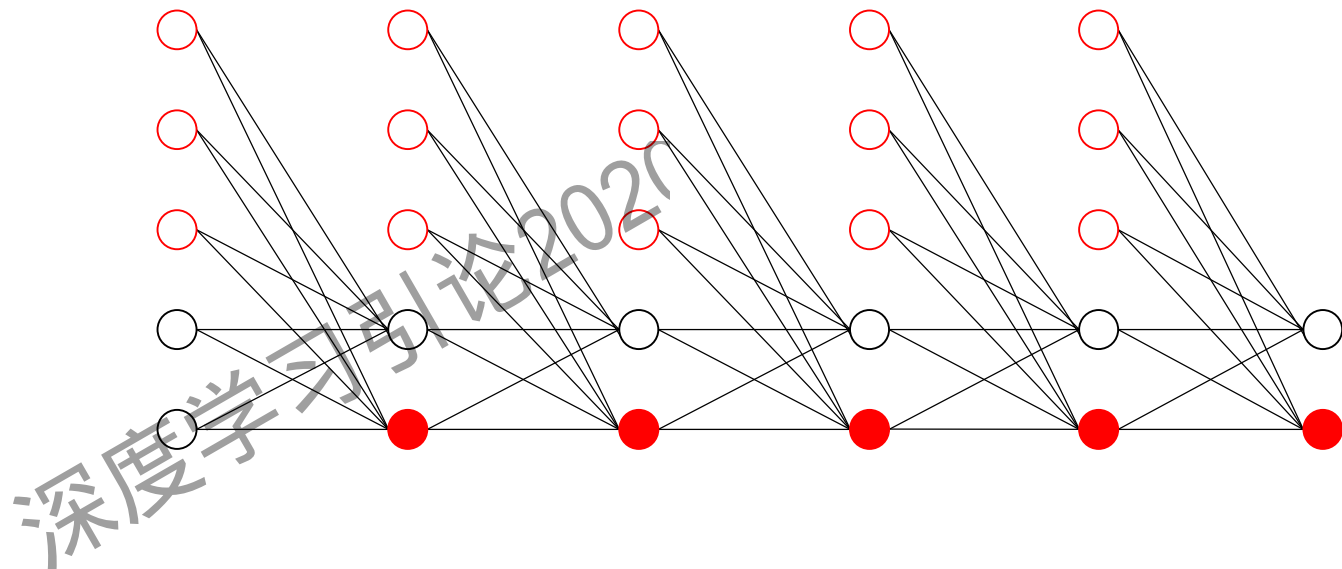
## Generated Testing Sequences

1. CBCAC
2. ACBBA
3. BACCB
4. ACBCB
5. AACBC
6. BAACB
7. AAACB
8. CCBAB
9. BBCAB
10. AABAC

.....

# Illustrative Example: Sequence Recognition

$$A = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad B = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad C = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$



Example outputs:

CAAC <b>B</b>	→	0.0184	0.0001	0.0211	0.0801	<b>0.9928</b>
A <b>B</b> BCA	→	0.0179	<b>0.9375</b>	0.0267	0.0012	0.0000
AAC <b>B</b> A	→	0.0179	0.0336	0.0286	<b>0.8722</b>	0.0000
CAC <b>B</b> B	→	0.0184	0.0001	0.0170	<b>0.8494</b>	0.0013
BCAAA	→	0.0182	0.0001	0.0001	0.0622	0.0018

# Illustrative Example: Sequence Recognition

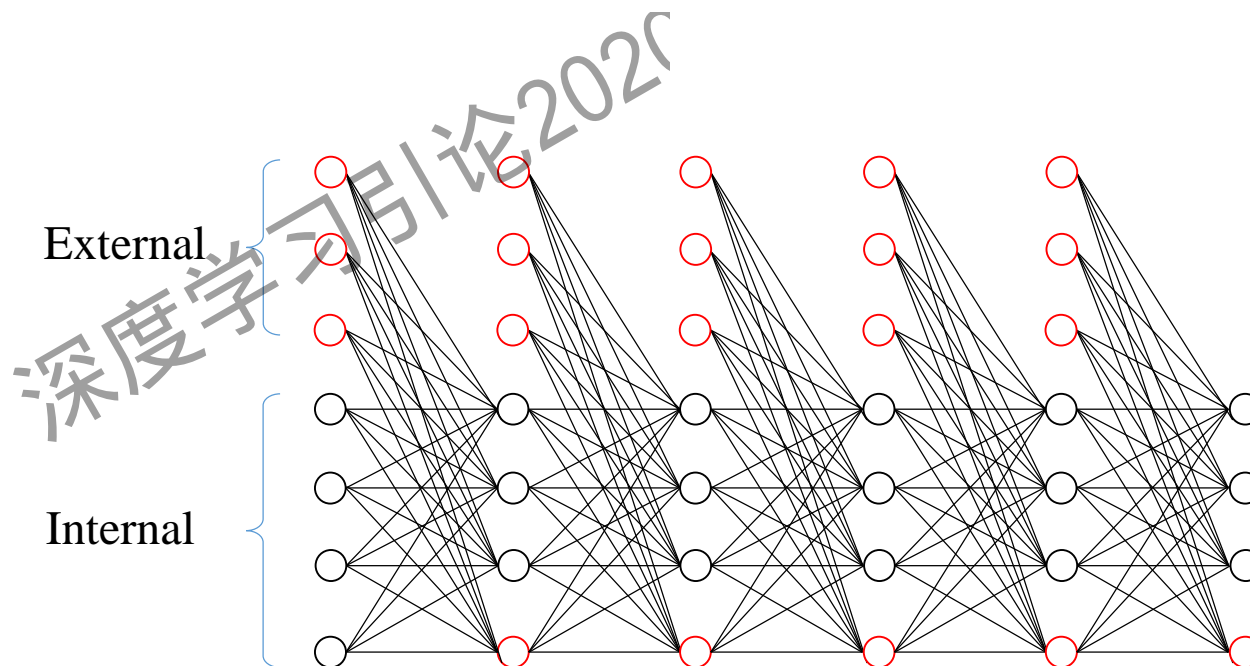
## *Recognize A followed by B Problem*

The task is to recognize A followed by B.

### Generated Sequences

1. ABCAB
2. CCBBA
3. CACCB
4. ACCCB
5. CACBC
6. AAACB
7. BAACB
8. CCBAB
9. BCCAB
10. CABAC

.....



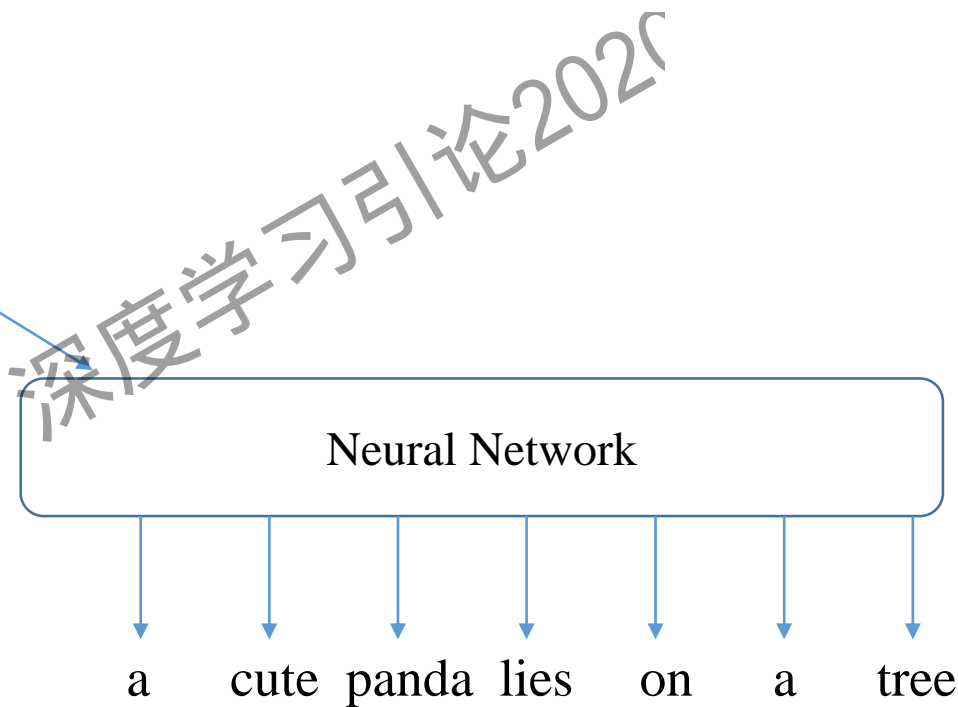
# Illustrative Example: Image Caption

## *Image Caption:*

The task is to describe the content of an image using properly formed English sentence.



image



# Illustrative Example: Image Caption

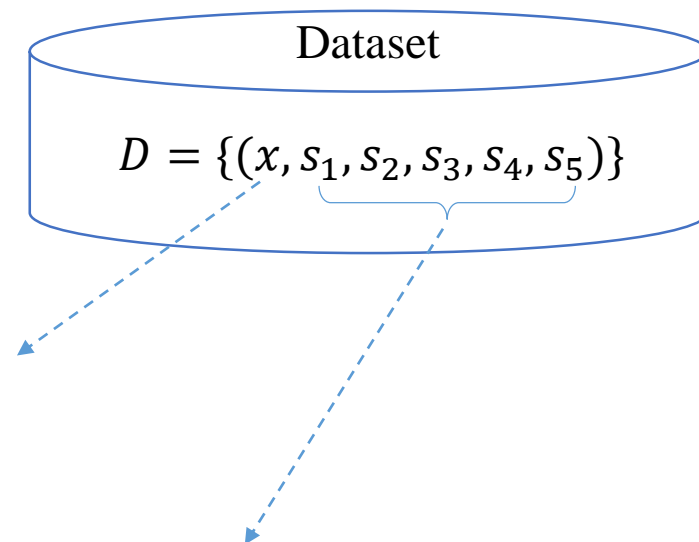
Dataset: COCO

COCO is a new image recognition, segmentation, and captioning dataset sponsored by Microsoft.

<http://mscoco.org/dataset/#download>

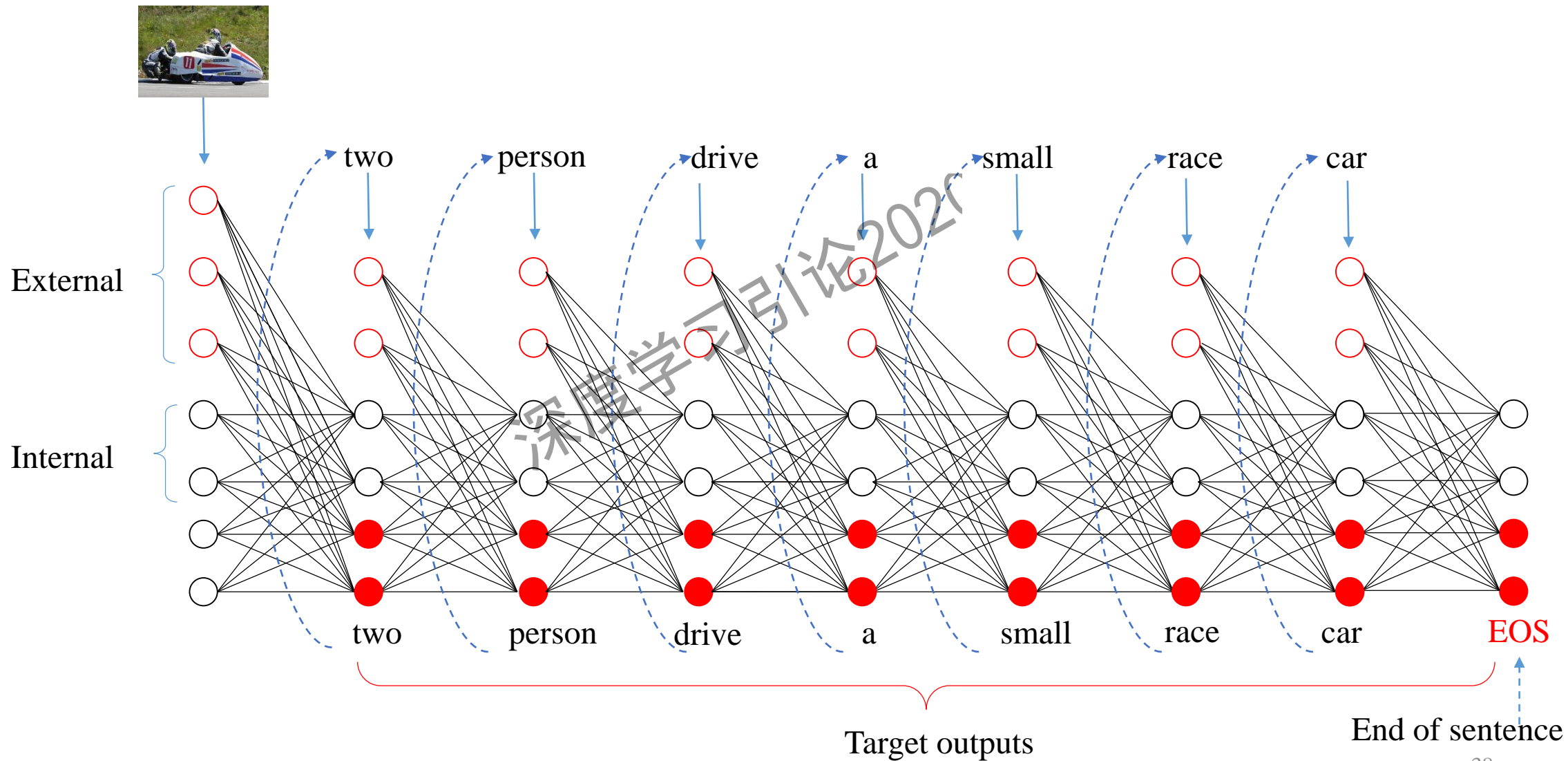
There are:

- 80,000 training samples
- 40,000 validation samples
- 40,000 test samples



1. Two person drive a small race car .
2. Two racer drive a white bike down a road .
3. Two motorist be ride along on their vehicle that be oddly design and color .
4. Two person be in a small race car drive by a green hill .
5. Two person in race uniform in a street car .

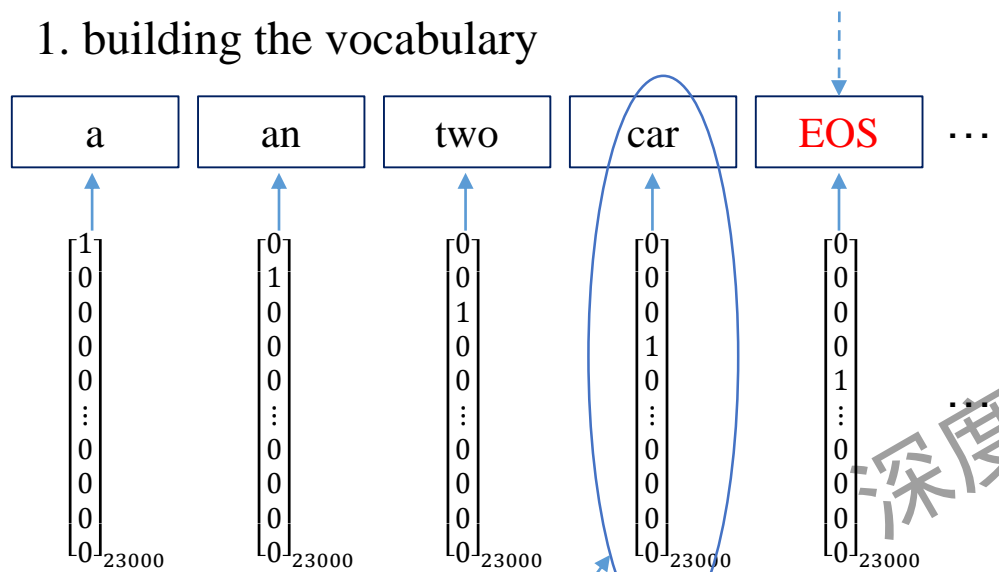
# Illustrative Example: Image Caption



# Illustrative Example: Image Caption

## Coding the Inputs:

### 1. building the vocabulary

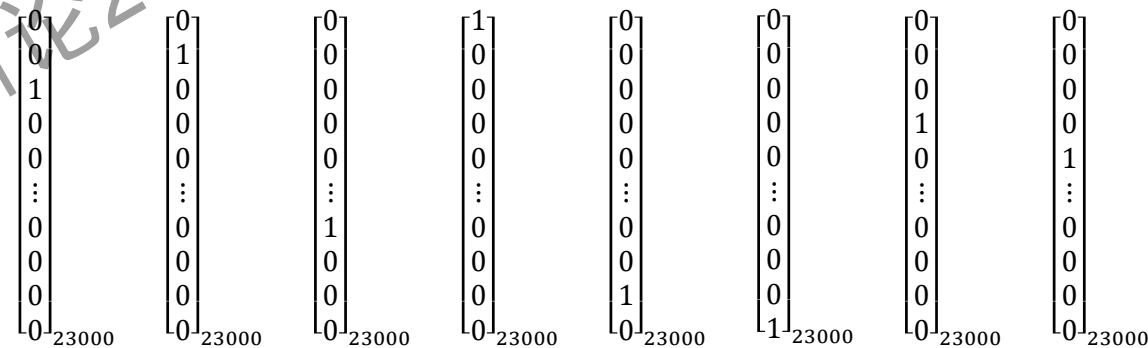


One-hot word vector

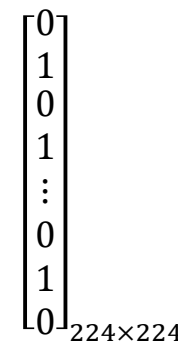
End of sentence

### 2. coding the sentence

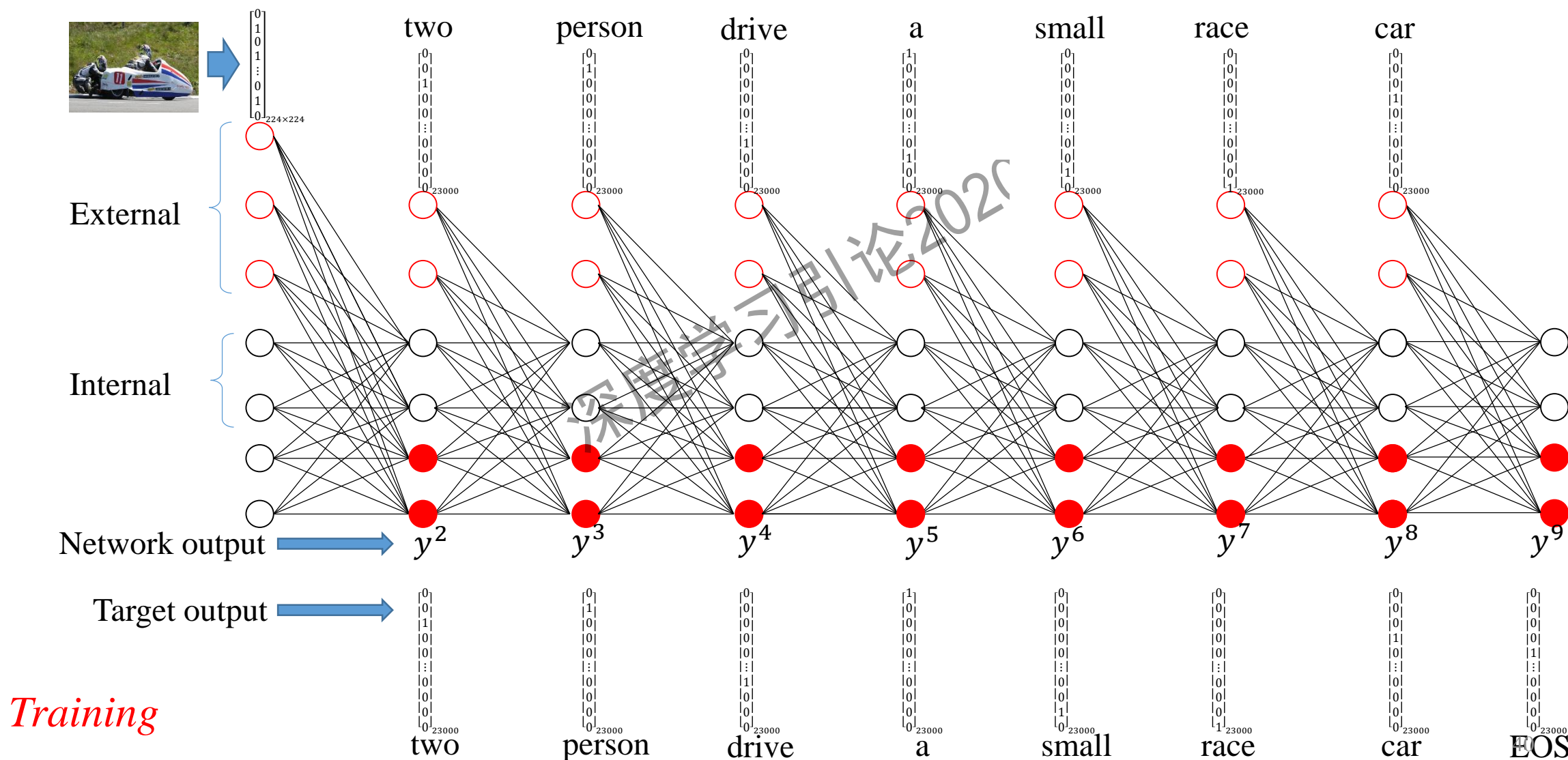
$s =$  two person drive a small race car EOS



### 3. digitizing the image

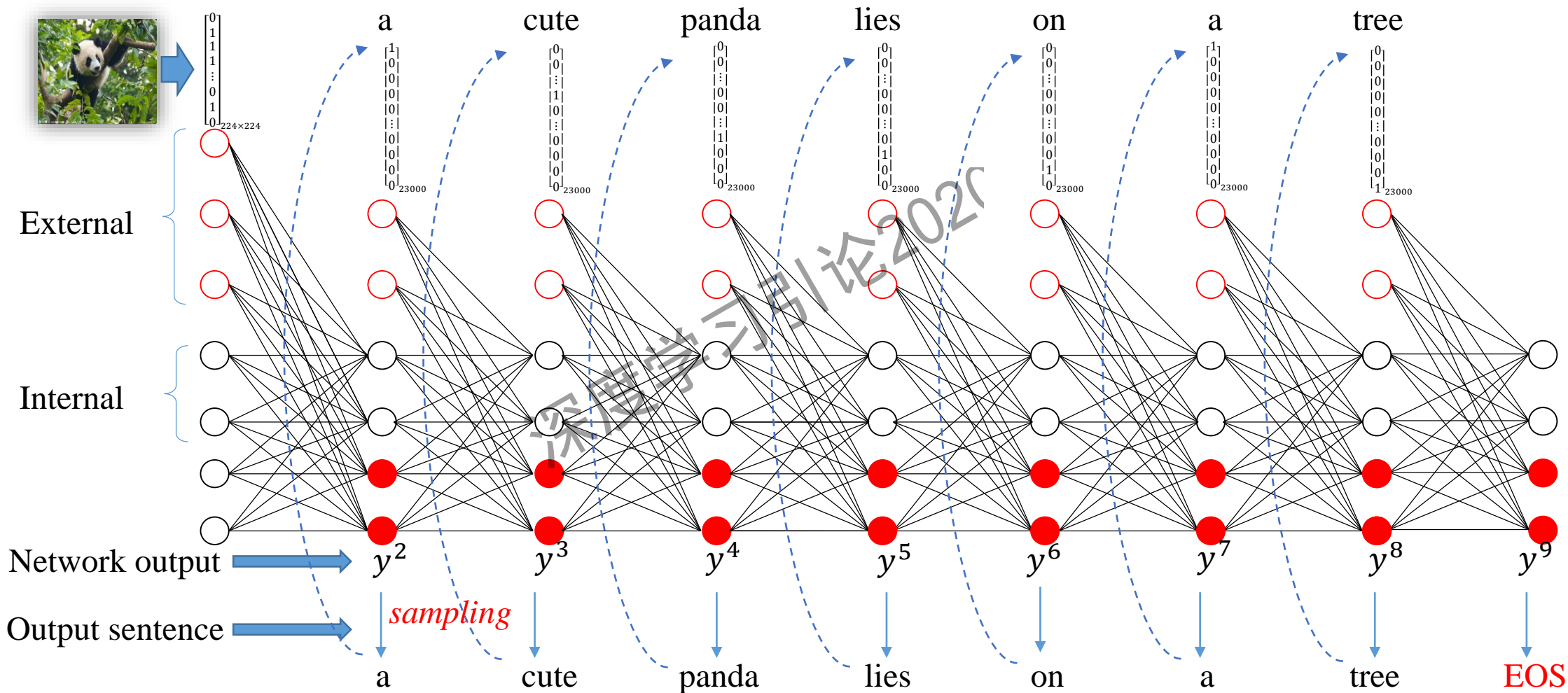


# Illustrative Example: Image Caption





# Illustrative Example: Image Caption



*Testing*

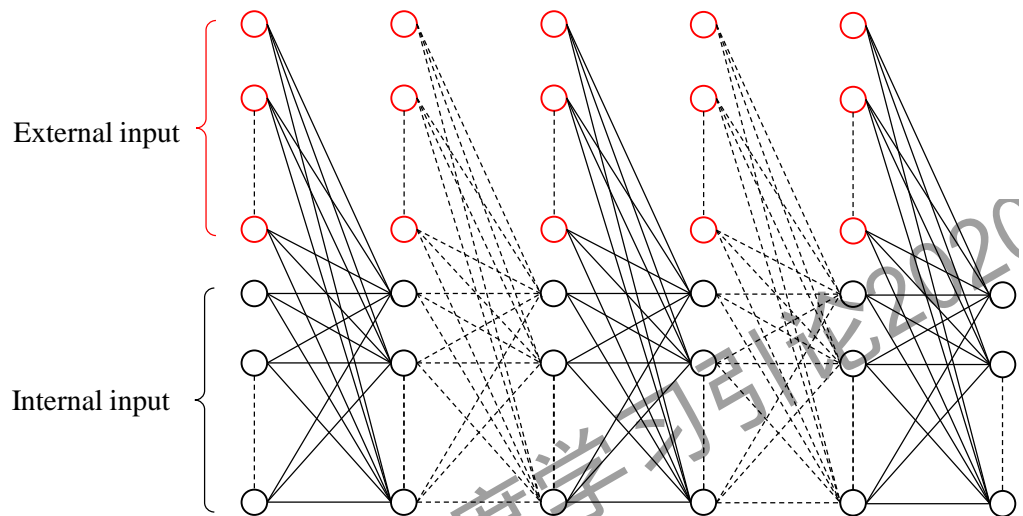
*Generate the word one by one, until get the word 'EOS'*

# Illustrative Example: Poem Creating



陆游  
卜算子·咏梅

驿外断桥边，  
寂寞开无主。  
已是黄昏独自愁，  
更著风和雨。  
无意苦争春，  
一任群芳妒。  
零落成泥碾作尘，  
只有香如故。



Can artificial neural networks create poem?



毛泽东  
卜算子·咏梅

风雨送春归，  
飞雪迎春到。  
已是悬崖百丈冰，  
犹有花枝俏。  
俏也不争春，  
只把春来报。  
待到山花烂漫时，  
她在丛中笑。



# Illustrative Example: Poem Creating



## 卜算子·咏梅

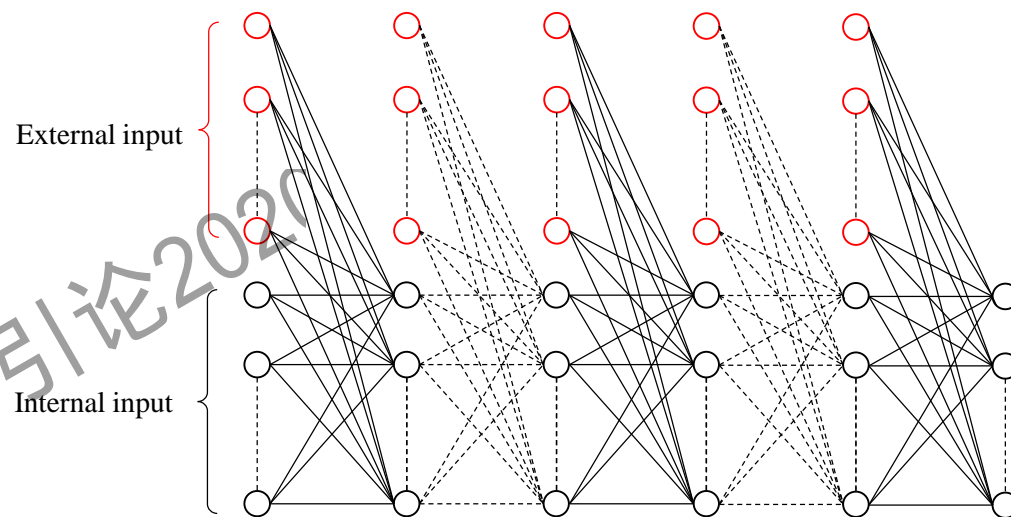
花谢早疏篱，  
几度陶潜里。  
永日梅花昔底寒，  
比向梅花妒。  
荣悴幻非凡，  
谓是娇芳伴。  
后著金陵几日时，  
中酒争先理。

## 卜算子·咏梅

朱阁见幽芳，  
露叶梅花里。  
玉屑琼台地屑琼，  
瑞鹊惊飞尾。  
便倚彩毫归，  
枝上簪盐谱。  
别作千秋一笑随，  
好趁伊家笑。

## 卜算子·咏梅

小试买梅花，  
并蒂栖香粉。  
肯向红蕖似竹姿，  
一叶清风许。  
心思寺炉高，  
深院松间曲。  
万纸参差故与黎，  
效我何知道。



# Outline

- A Sequence Recognizing Example
- Review of BP for Single Target Output NNs
- BP Method for Multiple Target Outputs NNs
- BP Algorithm for Multiple Target Outputs NNs
- Illustrative Examples
- Assignment

# Assignment

## *Assignment:*

Design a multi-target outputs neural network to learn to complete sequence.

The first two items of a sequence uniquely determine the remaining four.

Training Dataset

AA1212	AC1231	AD1221	AE1213
BA2312	BB2323	BC2331	BE2313
CB3123	CC3131	CD3121	CE3113
DA2112	DB2123	DC2131	DD2121
EA1312	EB1323	ED1321	EE1313

Testing Dataset

AB1223	BD2321	CA3112	DE2113	EC1331
--------	--------	--------	--------	--------

*The End*

深度学习引论2025

# A Sequence Recognizing Example

