

TSMC
Jan 23rd, 2016

Machine Learning Tutorial

YUN-NUNG (VIVIAN) CHEN

[HTTP://VIVIANCHEN.IDV.TW](http://vivianchen.idv.tw)



臺灣大學

National Taiwan University

Slide credit from Hung-Yi Lee and Mark Chang

Talk Outline

2

Part I: Introduction to
Machine Learning & Deep Learning



Part II: Variants of Neural Nets



Part III: Beyond Supervised Learning
& Recent Trends

Talk Outline

3

Part I: Introduction to
Machine Learning & Deep Learning



Part II: Variants of Neural Nets



Part III: Beyond Supervised Learning
& Recent Trends



4

PART I

Introduction to Machine Learning & Deep Learning

Part I: Introduction to ML & DL

5

- Basic Machine Learning
- Basic Deep Learning
- Toolkits and Learning Recipe

Part I: Introduction to ML & DL

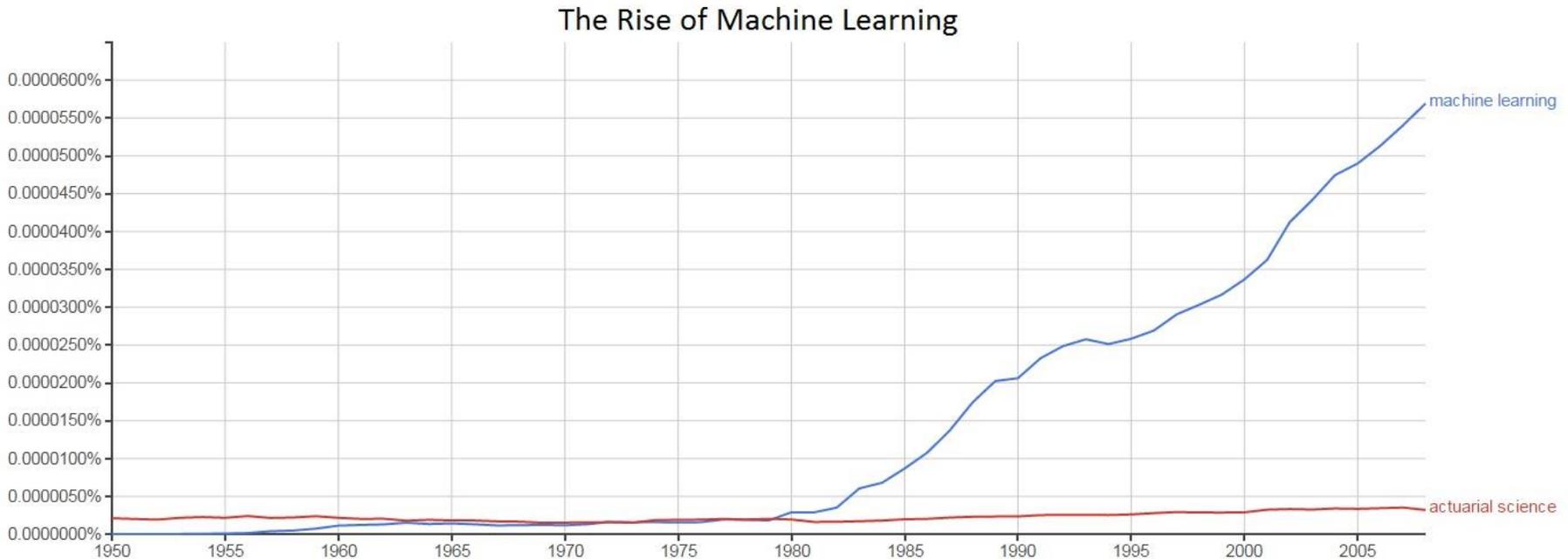
6

- **Basic Machine Learning**
- Basic Deep Learning
- Toolkits and Learning Recipe

Machine Learning

7

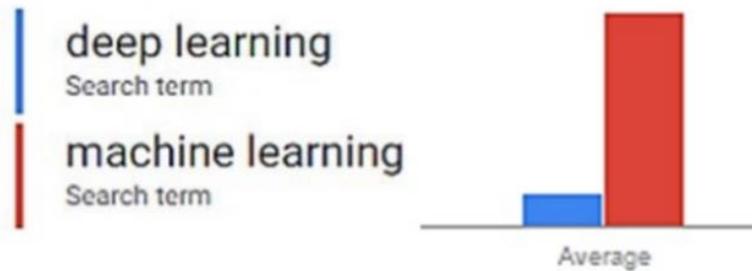
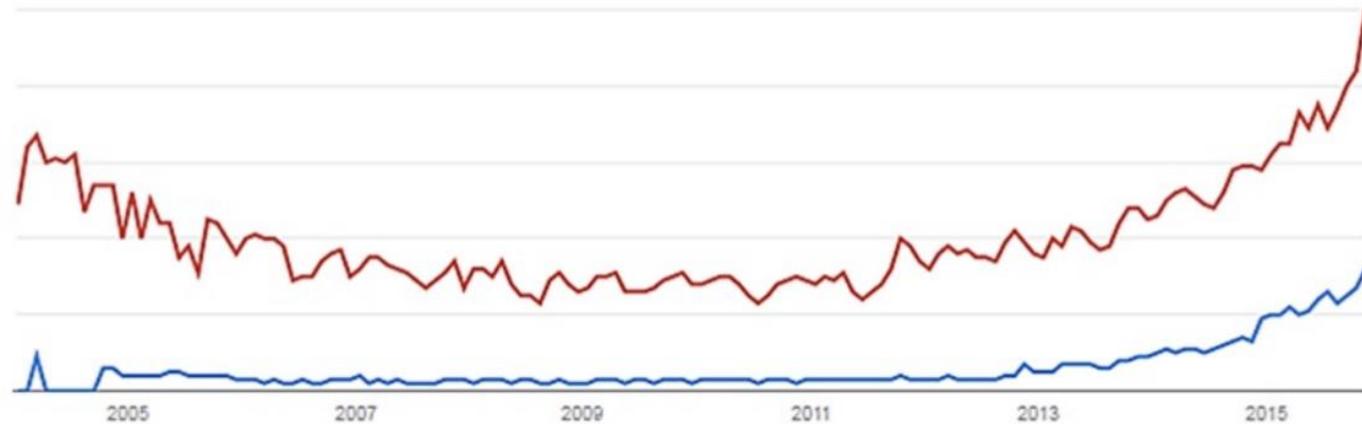
- Machine learning is rising rapidly in recent days



What the y-axis shows is this: of all the bigrams (two word letter combinations) contained in Google's sample of books written in English, what percentage of them are "machine learning" or "actuarial science"?

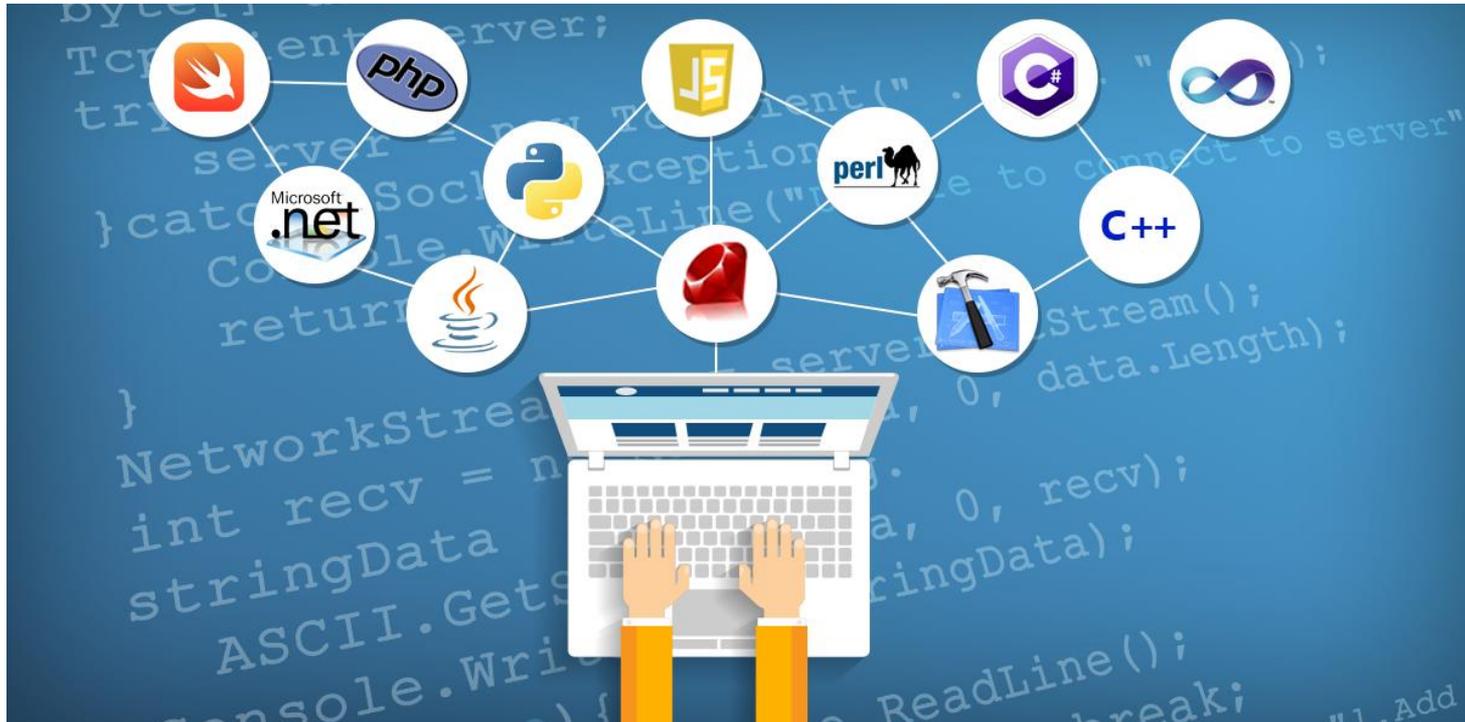
Recent Trend

8



What Computers Can Do?

9

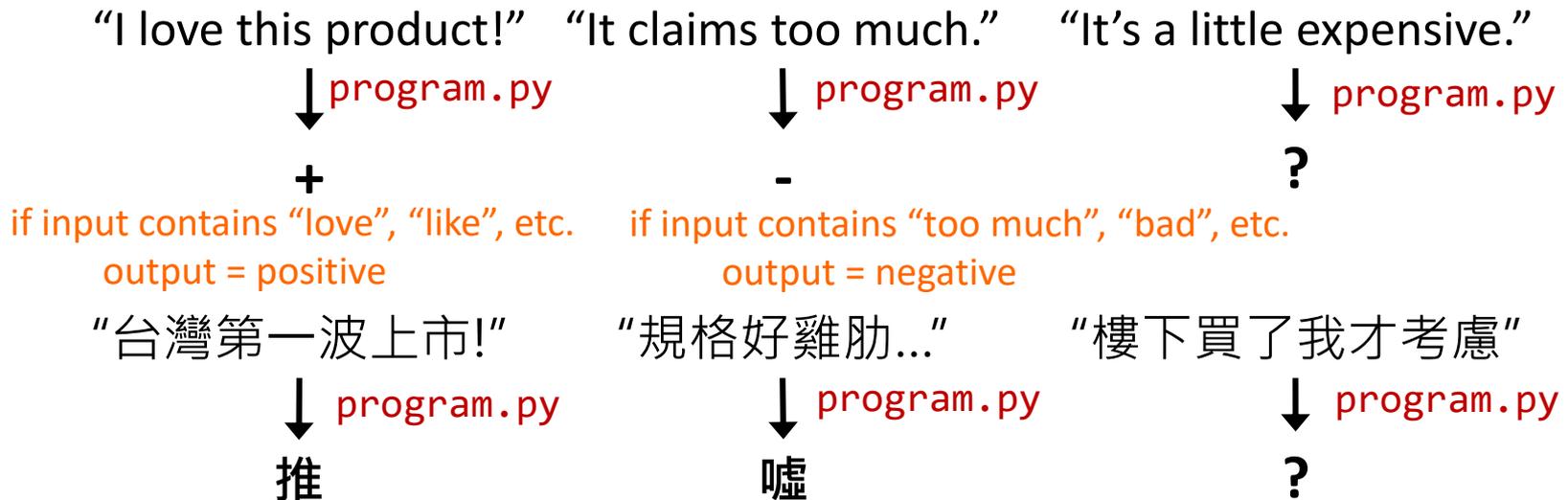


Programs can do the things you ask them to do

Program for Solving Tasks

10

- Task: predicting positive or negative given a product review



Some tasks are complex, and we don't know how to write a program to solve them.

Learning \approx Looking for a Function

11

- Task: predicting positive or negative given a product review

“I love this product!” “It claims too much.” “It’s a little expensive.”

↓ f
+

↓ f
-

↓ f
?

“台灣第一波上市!”

↓ f
推

“規格好雞肋...”

↓ f
噓

“樓下買了我才考慮”

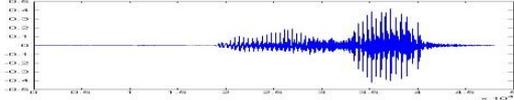
↓ f
?

Given a large amount of data, the machine learns what the function f should be.

Learning \approx Looking for a Function

12

□ Speech Recognition

$f(\text{  }) = \text{“你好”}$

□ Image Recognition

$f(\text{  }) = \text{cat}$

□ Go Playing

$f(\text{  }) = \text{5-5 (next move)}$

□ Dialogue System

$f(\text{ “台積電怎麼去” }) = \text{“地址為...
現在建議搭乘計程車”}$

Framework

$$f\left(\text{Image of a cat}\right) = \text{"cat"}$$

13

A set of
function

Model

$f_1, f_2 \dots$

$$f_1\left(\text{Image of a cat}\right) = \text{"cat"}$$

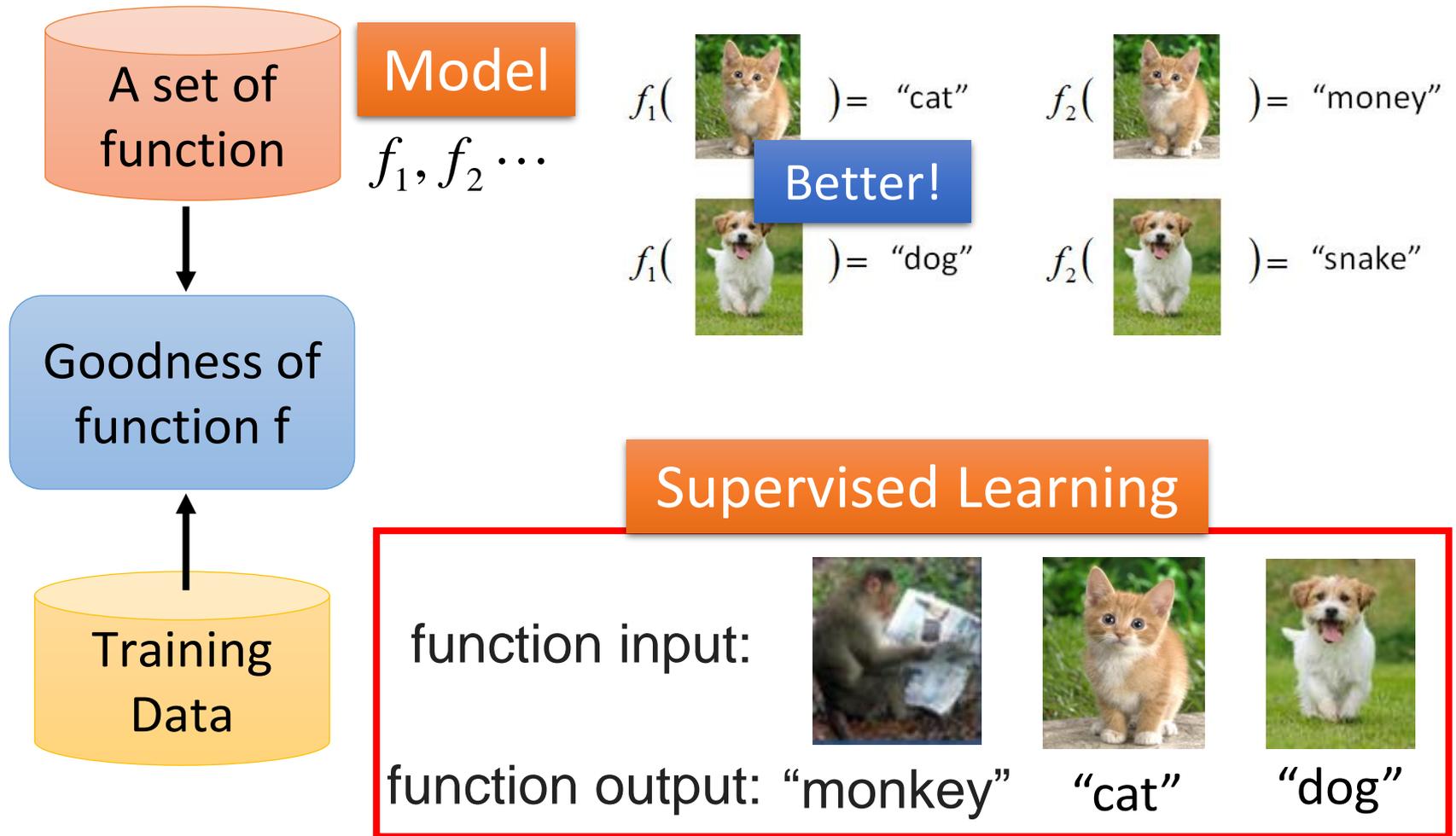
$$f_2\left(\text{Image of a cat}\right) = \text{"monkey"}$$

$$f_1\left(\text{Image of a dog}\right) = \text{"dog"}$$

$$f_2\left(\text{Image of a dog}\right) = \text{"snake"}$$

Framework

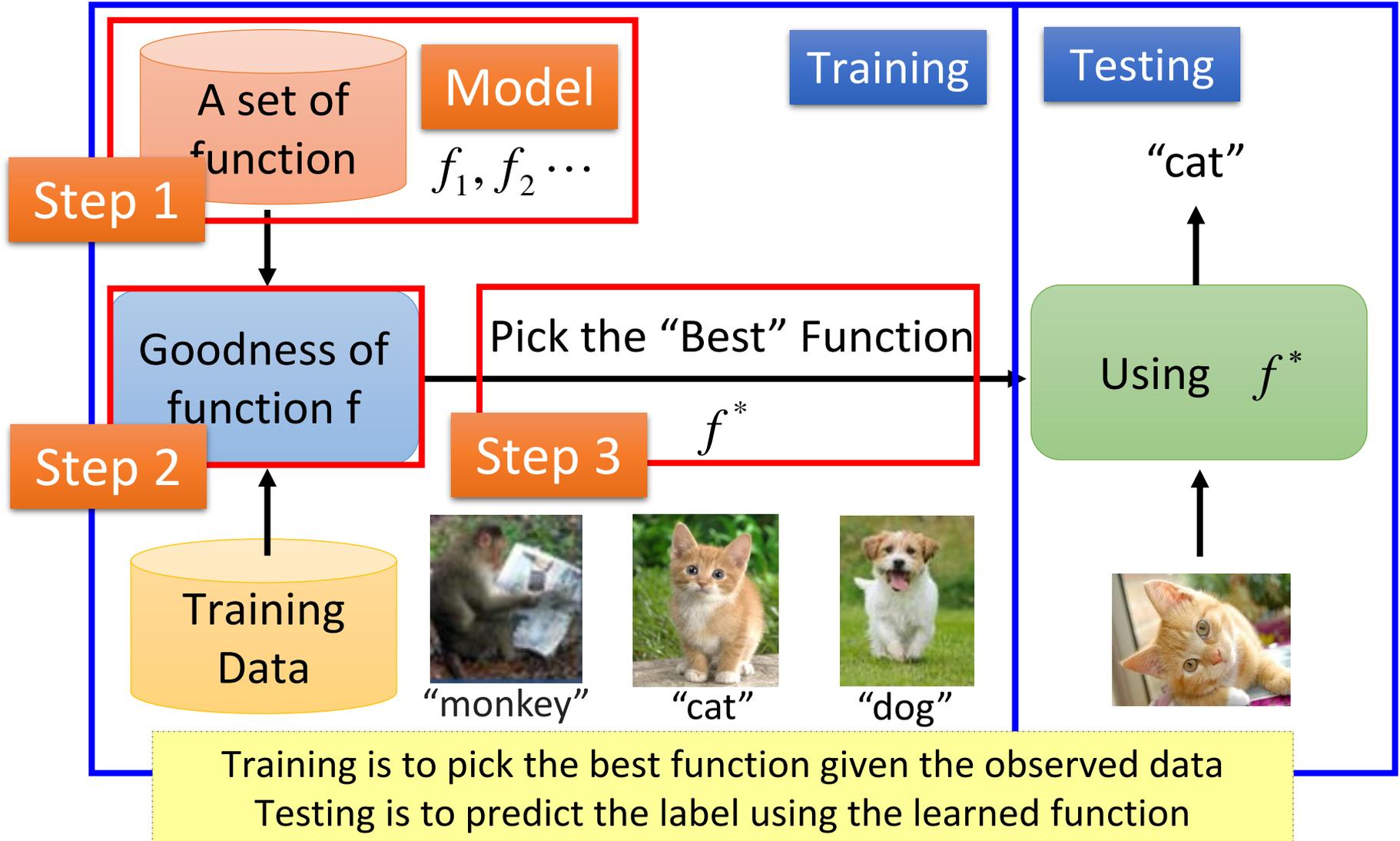
$$f\left(\text{img_cat}\right) = \text{"cat"}$$



Framework

$$f(\text{Image of a cat}) = \text{"cat"}$$

15

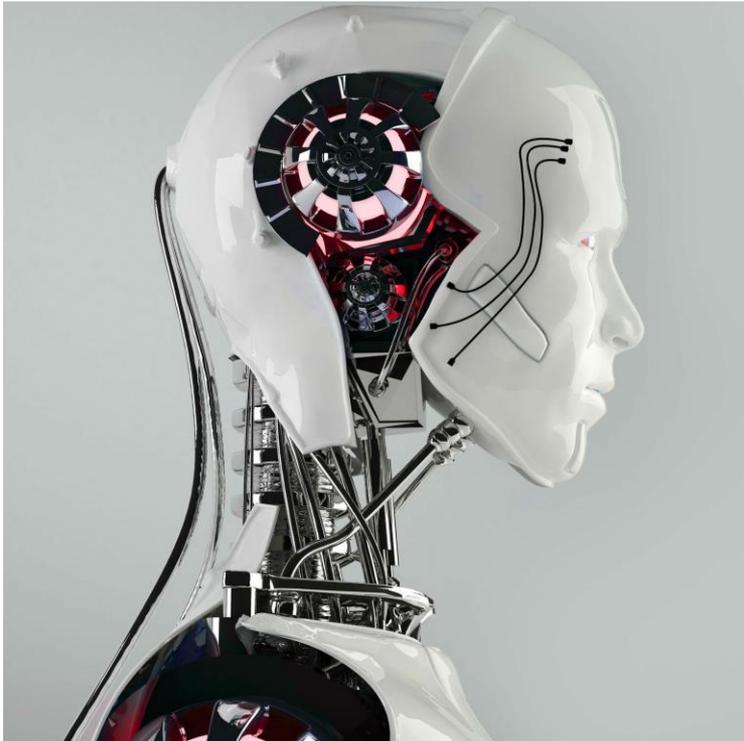


Why to Learn Machine Learning?

16

- AI Age
 - ▣ AI can work for most of labor work?
 - ▣ New job market → AI 訓練師

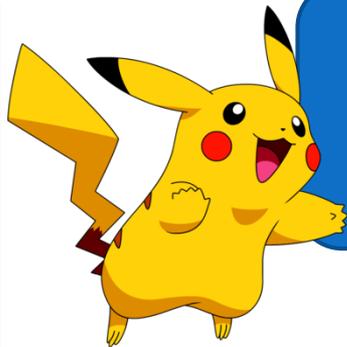
(Machine Learning Expert 機器學習專家、
Data Scientist 資料科學家)



AI 訓練師

17

機器不是自己會學嗎？
為什麼需要 AI 訓練師



戰鬥是寶可夢在打，
為什麼需要寶可夢訓練師？



AI 訓練師

Step 1:
define a set
of function



Step 2:
goodness of
function



Step 3: pick
the best
function

18

□ 寶可夢訓練師

- ▣ 挑選適合的寶可夢來戰鬥
 - 寶可夢有不同的屬性
- ▣ 召喚出來的寶可夢不一定聽話
 - E.g. 小智的噴火龍

→ 需要足夠的經驗

□ AI 訓練師

- ▣ 在 step 1，AI訓練師要挑選合適的模型
 - 不同模型適合處理不同的問題
- ▣ 不一定能在 step 3 找出 best function
 - E.g. Deep Learning

→ 需要足夠的經驗

AI 訓練師

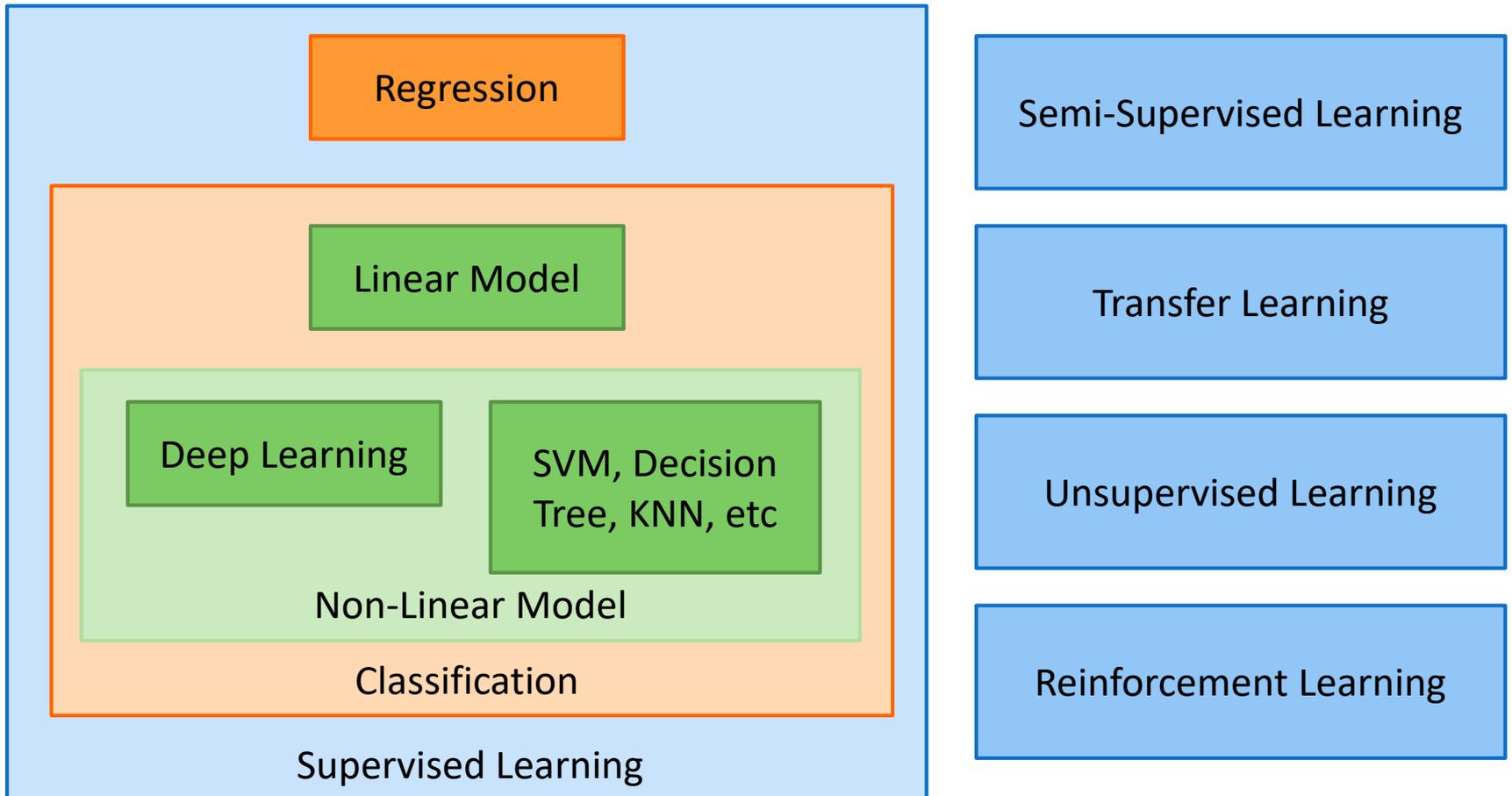
19

- 厲害的 AI ， AI 訓練師功不可沒
- 讓我們一起朝 AI 訓練師之路邁進



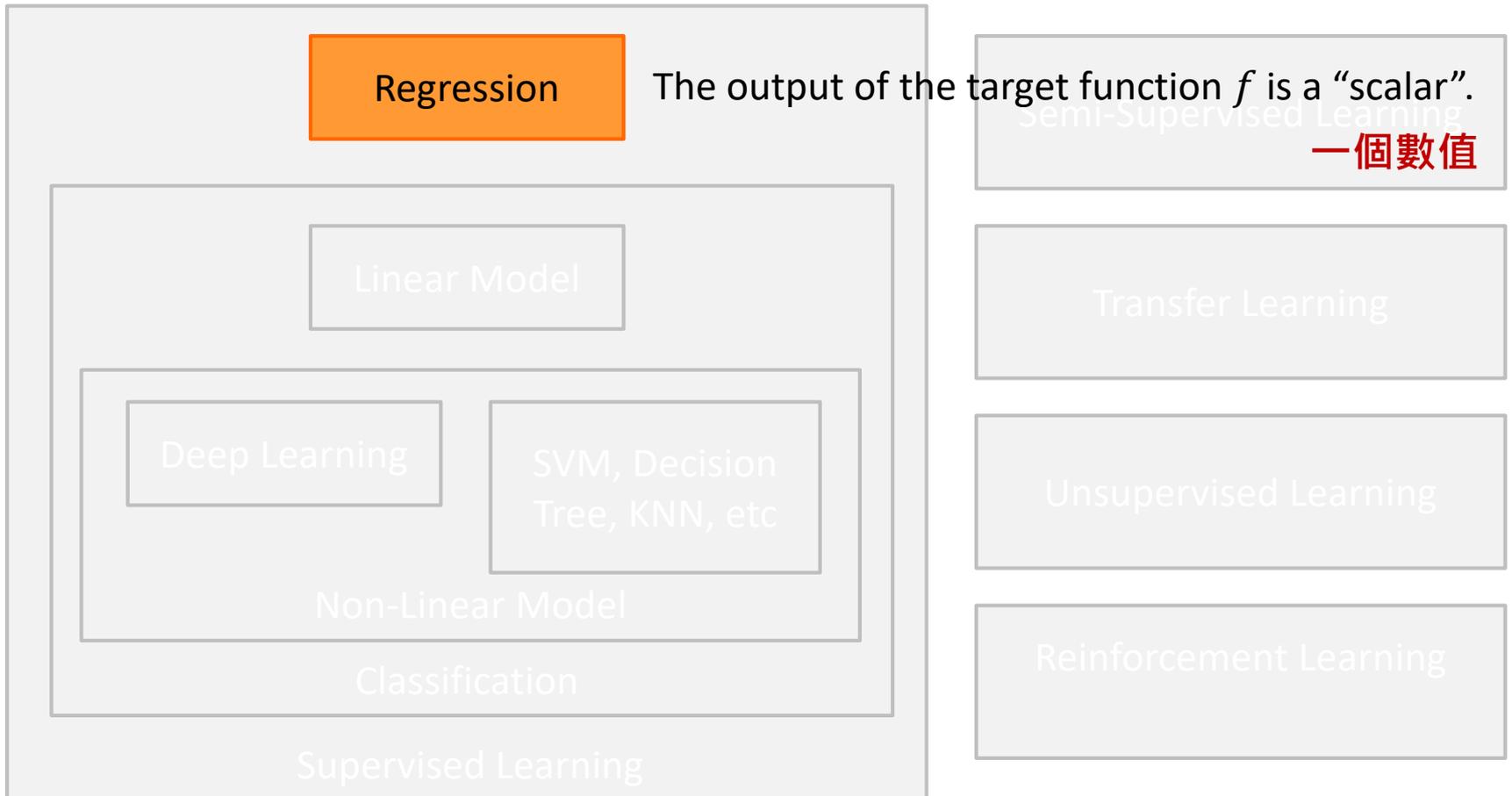
Machine Learning Map

Scenario Task Method



Machine Learning Map

Scenario Task Method



Regression

22

- Stock Market Forecast

$$f(\text{Dow Jones Industrial Average at today}) = \text{Dow Jones Industrial Average at tomorrow}$$


- Self-driving Car

$$f(\text{Current steering angle}) = \text{Directional angle}$$

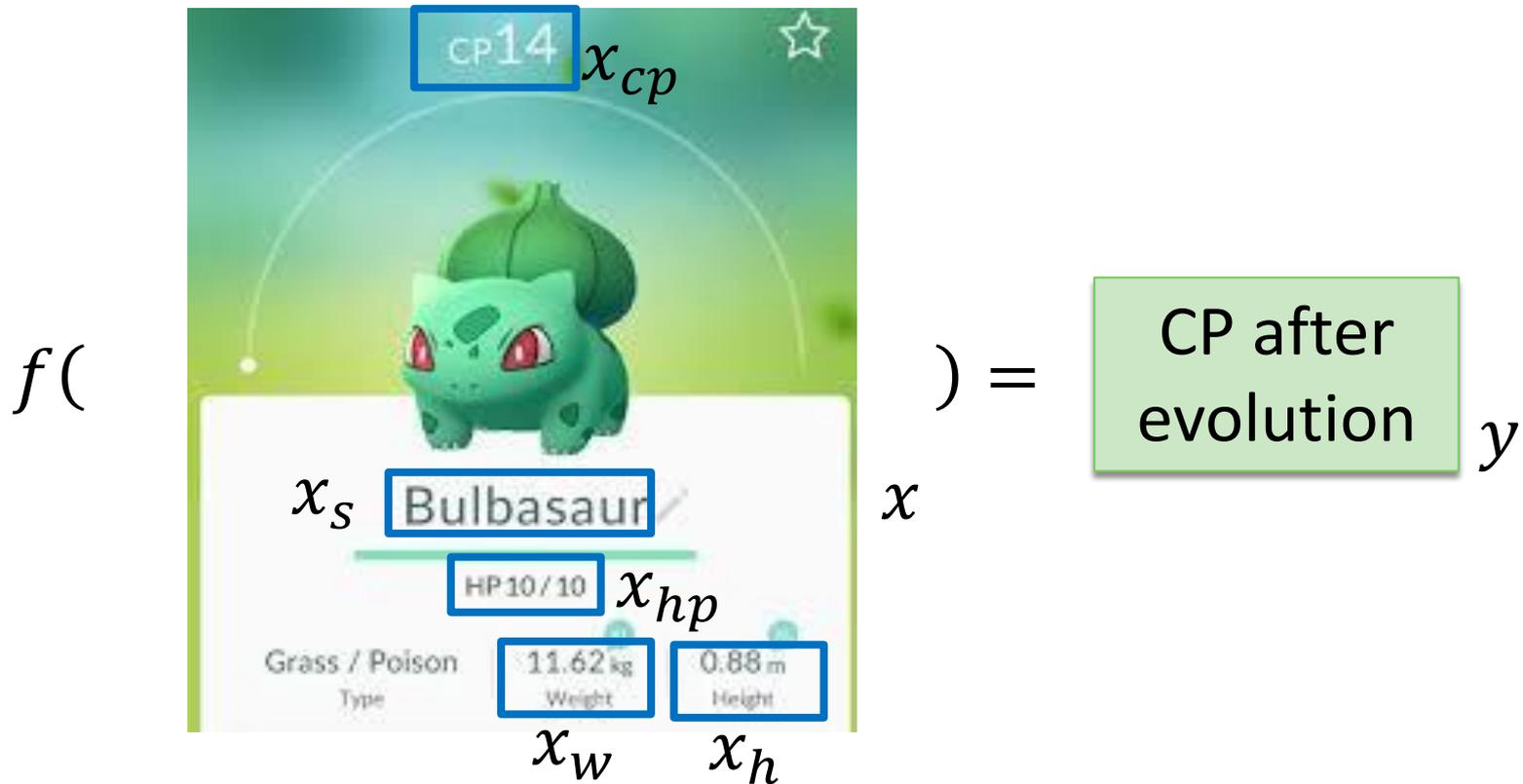

- Recommendation

$$f(\text{User A}, \text{Item B}) = \text{Purchase probability}$$

Example Application

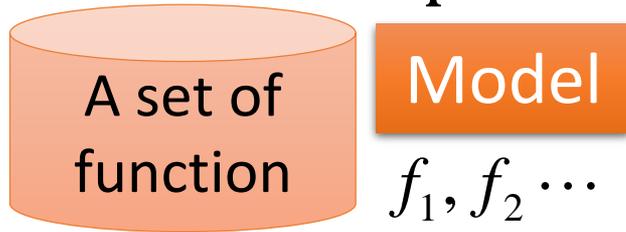
23

- Estimating the Combat Power (CP) of a pokemon after evolution



Step 1: Model

$$y = b + w \cdot x_{cp}$$



w and b are parameters
(can be any value)

$$f_1: y = 10.0 + 9.0 \cdot x_{cp}$$

$$f_2: y = 9.8 + 9.2 \cdot x_{cp}$$

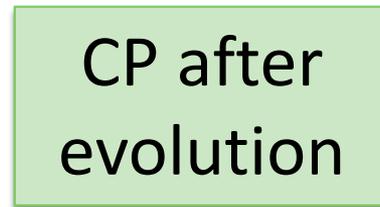
$$f_3: y = -0.8 - 1.2 \cdot x_{cp}$$

..... infinite

$f($



$x) =$



y

Linear model:

$$y = b + \sum w_i x_i$$

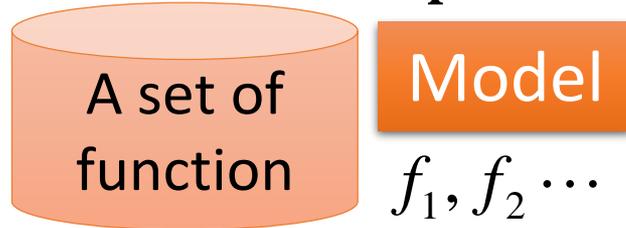
x_i : an attribute of input x feature

w_i : weight, b: bias

Step 2: Goodness of Function

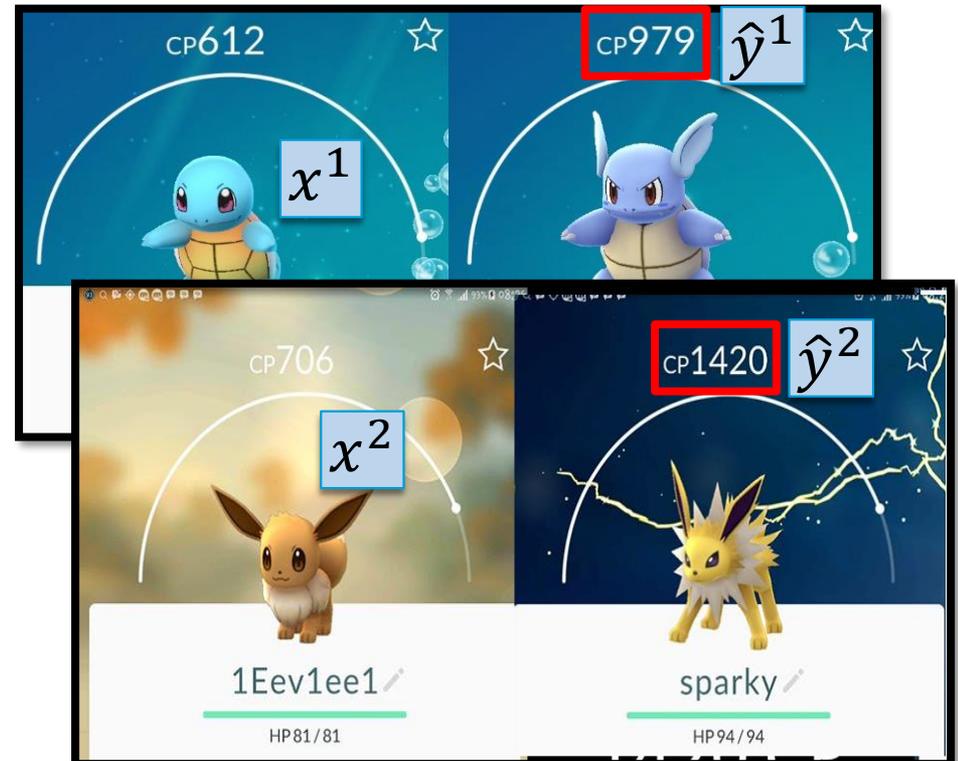
25

$$y = b + w \cdot x_{cp}$$



function
input:

function
output (scalar):



Step 2: Goodness of Function

26

□ Training data

- 1st pokemon:

$$(x^1, \hat{y}^1)$$

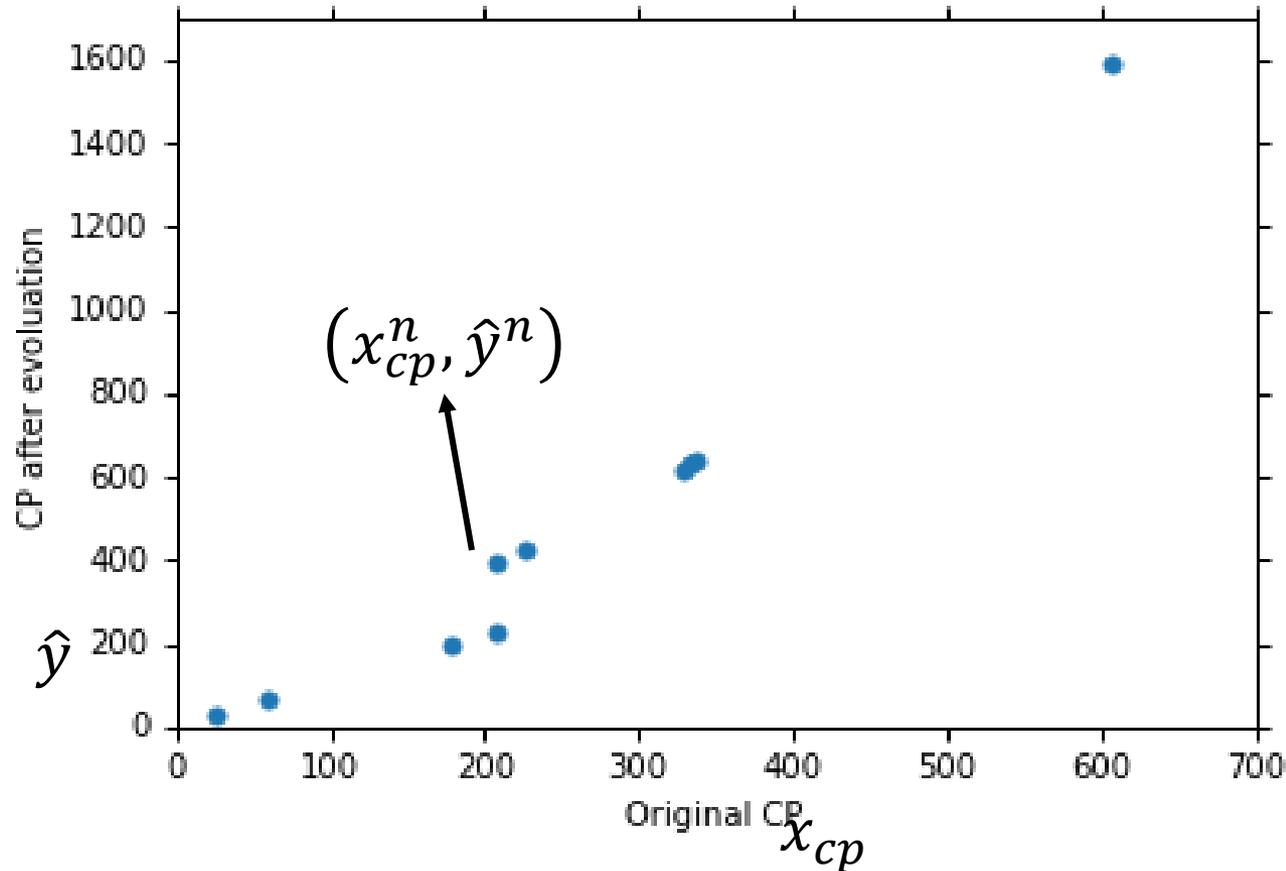
- 2nd pokemon:

$$(x^2, \hat{y}^2)$$

⋮

- 10th pokemon:

$$(x^{10}, \hat{y}^{10})$$



This is real data.

Step 2: Goodness of Function

27

$$y = b + w \cdot x_{cp}$$

A set of
function

Model

$f_1, f_2 \dots$

Goodness of
function f

Loss function L :

Input: a function, output: how bad it is

$$L(f) = L(w, b)$$

Estimated y based
on input function

$$= \sum_{n=1}^{10} \left(\hat{y}^n - \left(\underline{b + w \cdot x_{cp}^n} \right) \right)^2$$

Estimation error

Training
Data

Sum over examples

Step 2: Goodness of Function

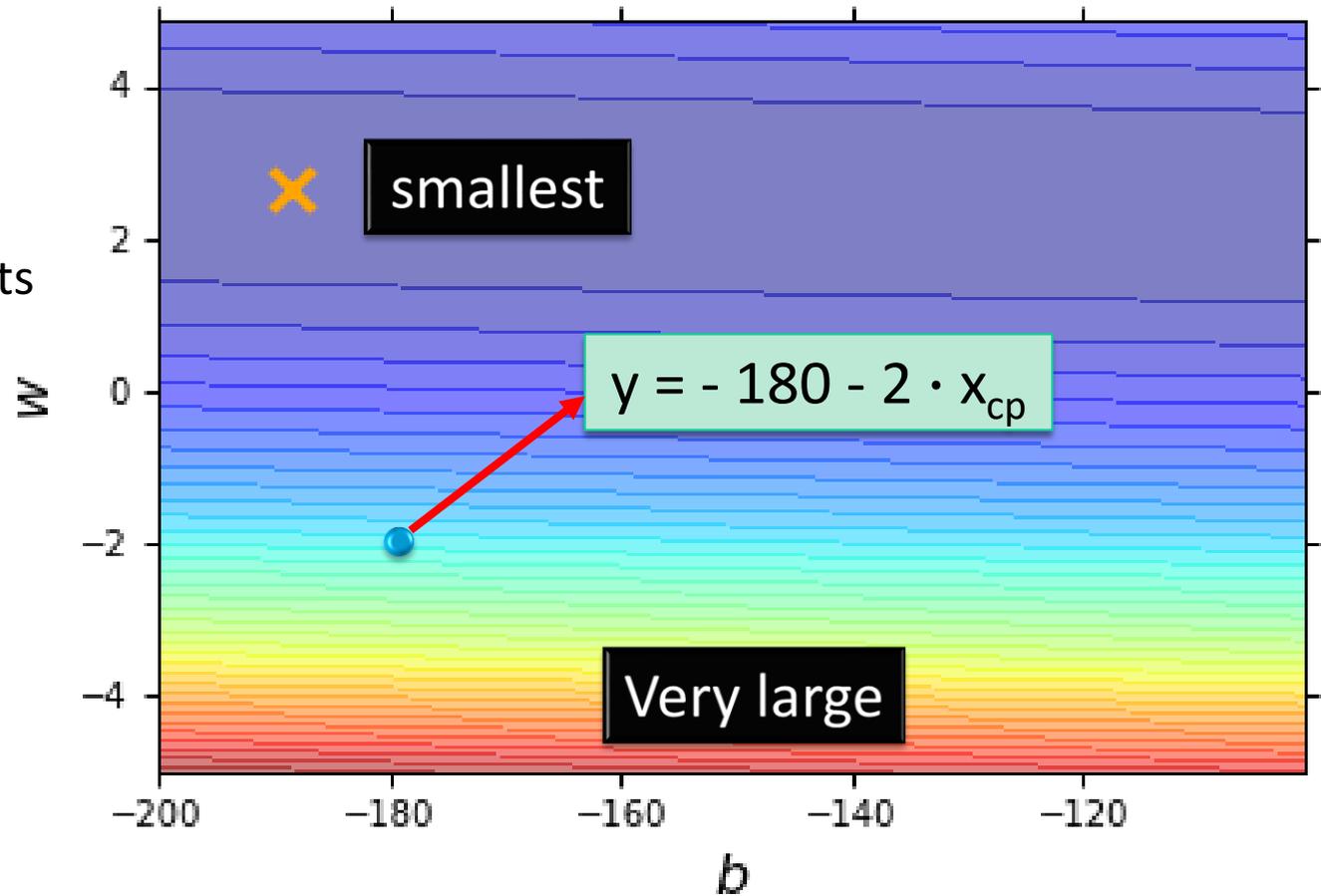
28

□ Loss Function

$$L(w, b) = \sum_{n=1}^{10} \left(\hat{y}^n - (b + w \cdot x_{cp}^n) \right)^2$$

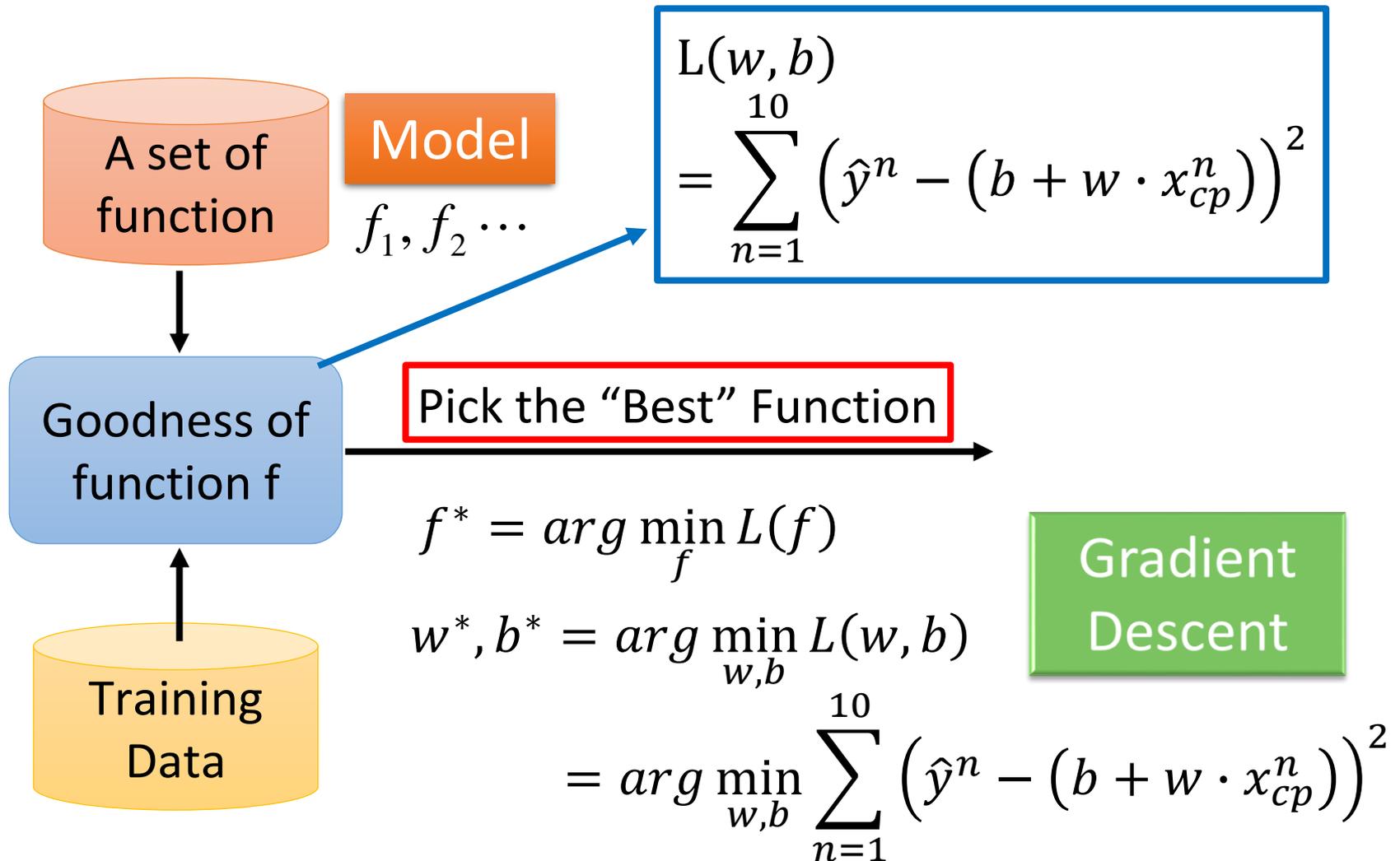
Each point in the figure is a function

The color represents $L(w, b)$



Step 3: Best Function

29



Step 3: Gradient Descent

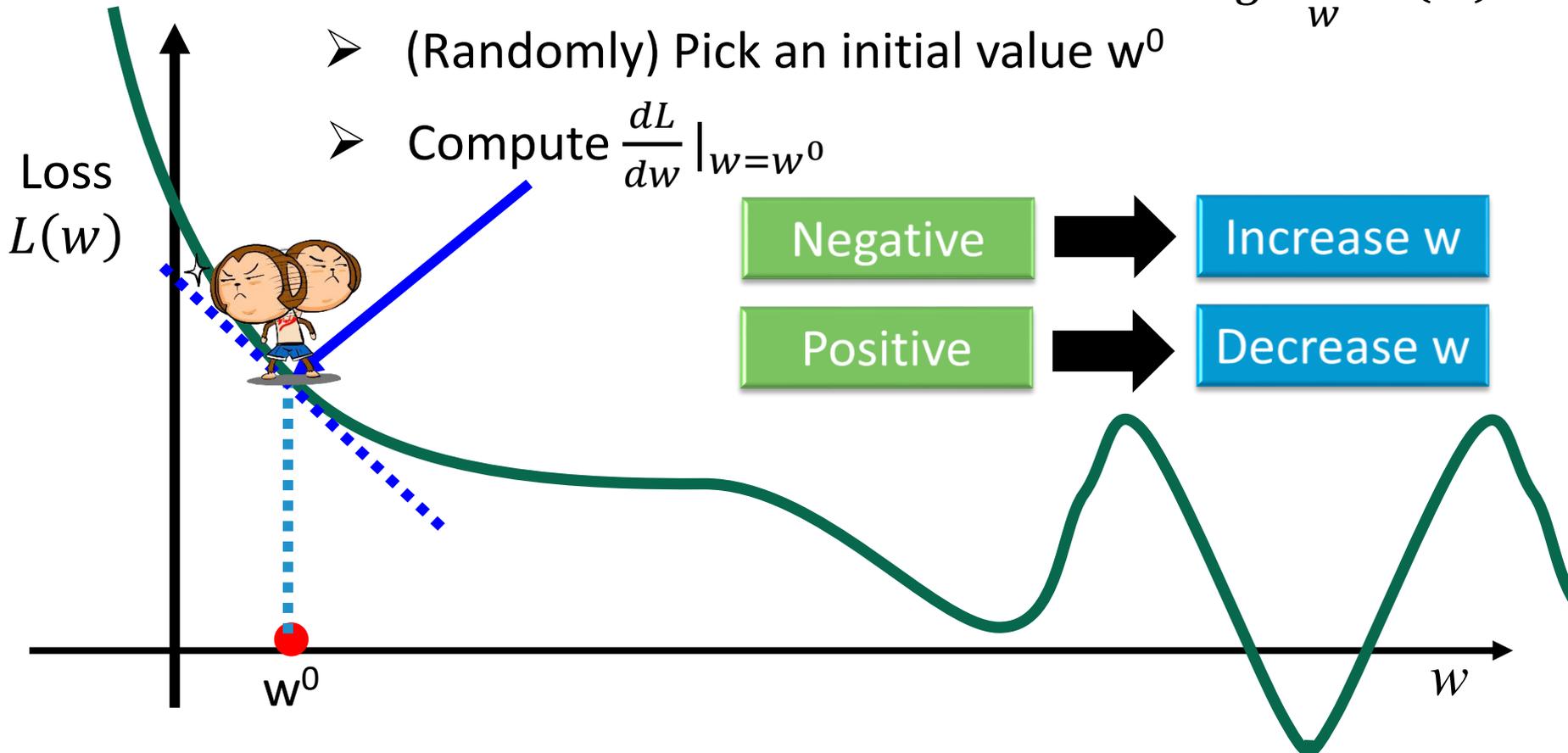
30

- Consider loss function $L(w)$ with one parameter w :

$$w^* = \underset{w}{\operatorname{arg\,min}} L(w)$$

➤ (Randomly) Pick an initial value w^0

➤ Compute $\frac{dL}{dw} \Big|_{w=w^0}$



Step 3: Gradient Descent

31

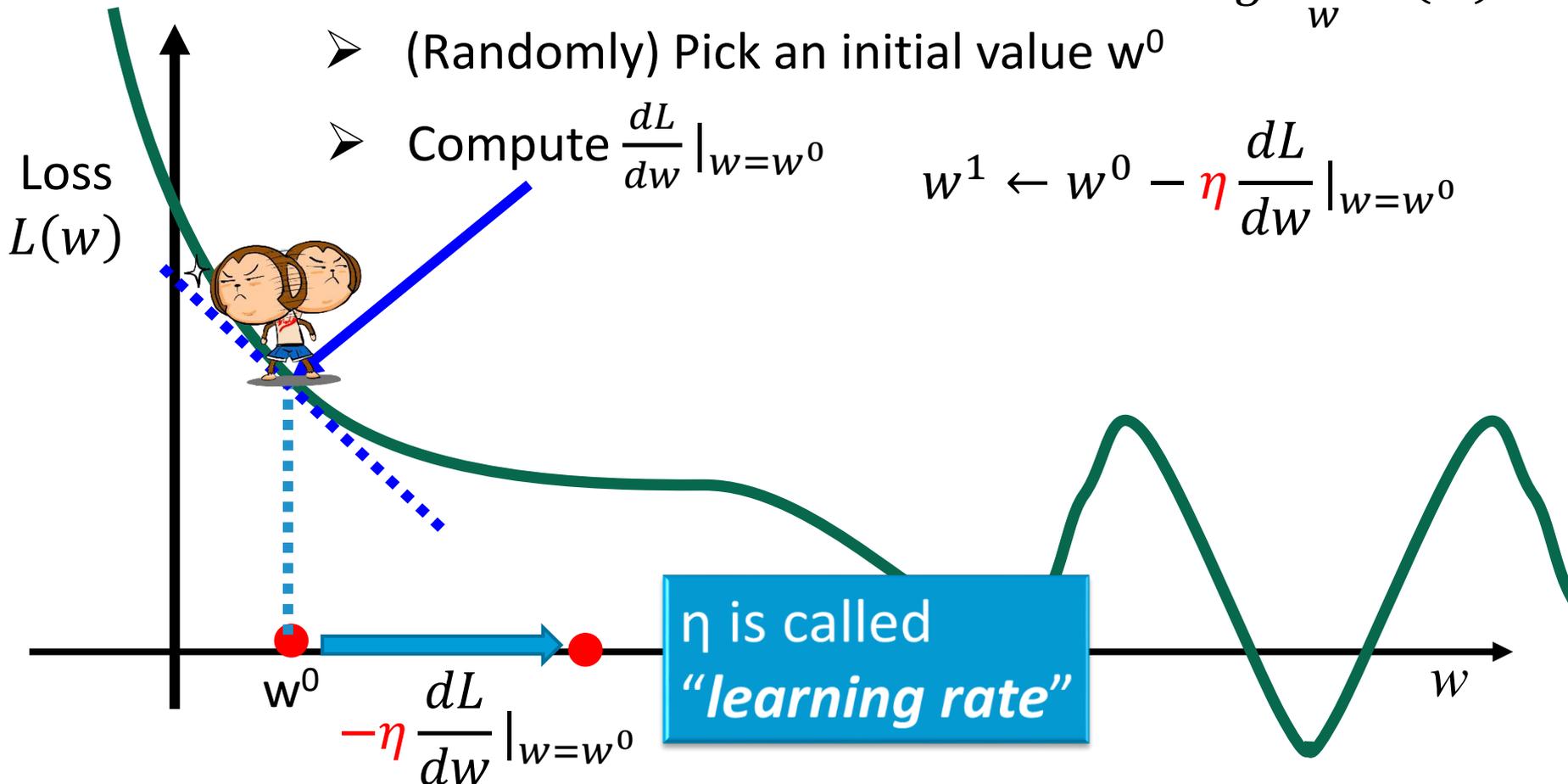
- Consider loss function $L(w)$ with one parameter w :

$$w^* = \underset{w}{\operatorname{arg\,min}} L(w)$$

➤ (Randomly) Pick an initial value w^0

➤ Compute $\frac{dL}{dw} \Big|_{w=w^0}$

$$w^1 \leftarrow w^0 - \eta \frac{dL}{dw} \Big|_{w=w^0}$$



Step 3: Gradient Descent

32

- Consider loss function $L(w)$ with one parameter w :

$$w^* = \underset{w}{\operatorname{arg\,min}} L(w)$$

➤ (Randomly) Pick an initial value w^0

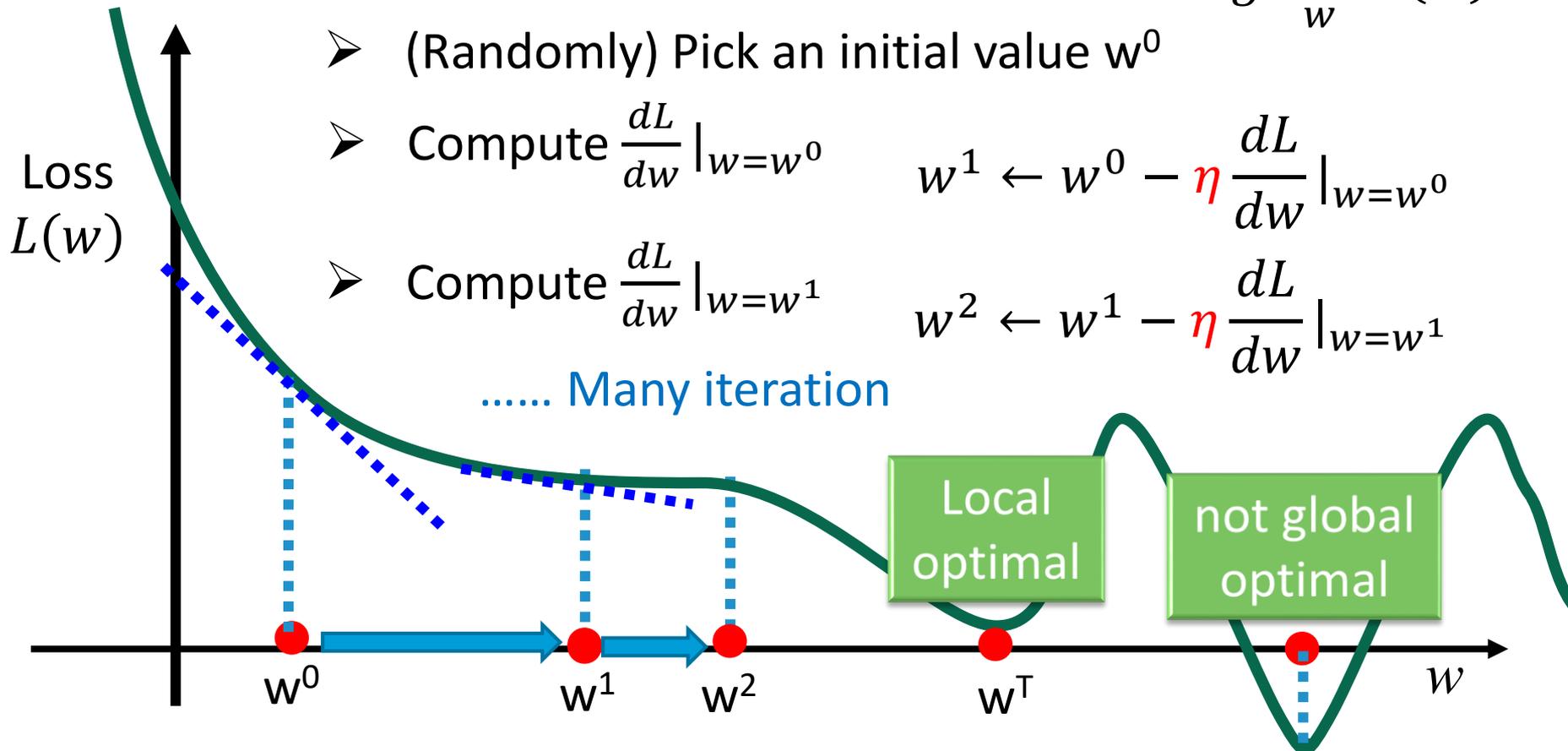
➤ Compute $\frac{dL}{dw} \Big|_{w=w^0}$

$$w^1 \leftarrow w^0 - \eta \frac{dL}{dw} \Big|_{w=w^0}$$

➤ Compute $\frac{dL}{dw} \Big|_{w=w^1}$

$$w^2 \leftarrow w^1 - \eta \frac{dL}{dw} \Big|_{w=w^1}$$

..... Many iteration



Step 3: Gradient Descent

33

- How about two parameters?

$$w^*, b^* = \arg \min_{w, b} L(w, b)$$

- (Randomly) Pick an initial value w^0, b^0

- Compute $\frac{\partial L}{\partial w} \Big|_{w=w^0, b=b^0}, \frac{\partial L}{\partial b} \Big|_{w=w^0, b=b^0}$

$$w^1 \leftarrow w^0 - \eta \frac{\partial L}{\partial w} \Big|_{w=w^0, b=b^0} \quad b^1 \leftarrow b^0 - \eta \frac{\partial L}{\partial b} \Big|_{w=w^0, b=b^0}$$

- Compute $\frac{\partial L}{\partial w} \Big|_{w=w^1, b=b^1}, \frac{\partial L}{\partial b} \Big|_{w=w^1, b=b^1}$

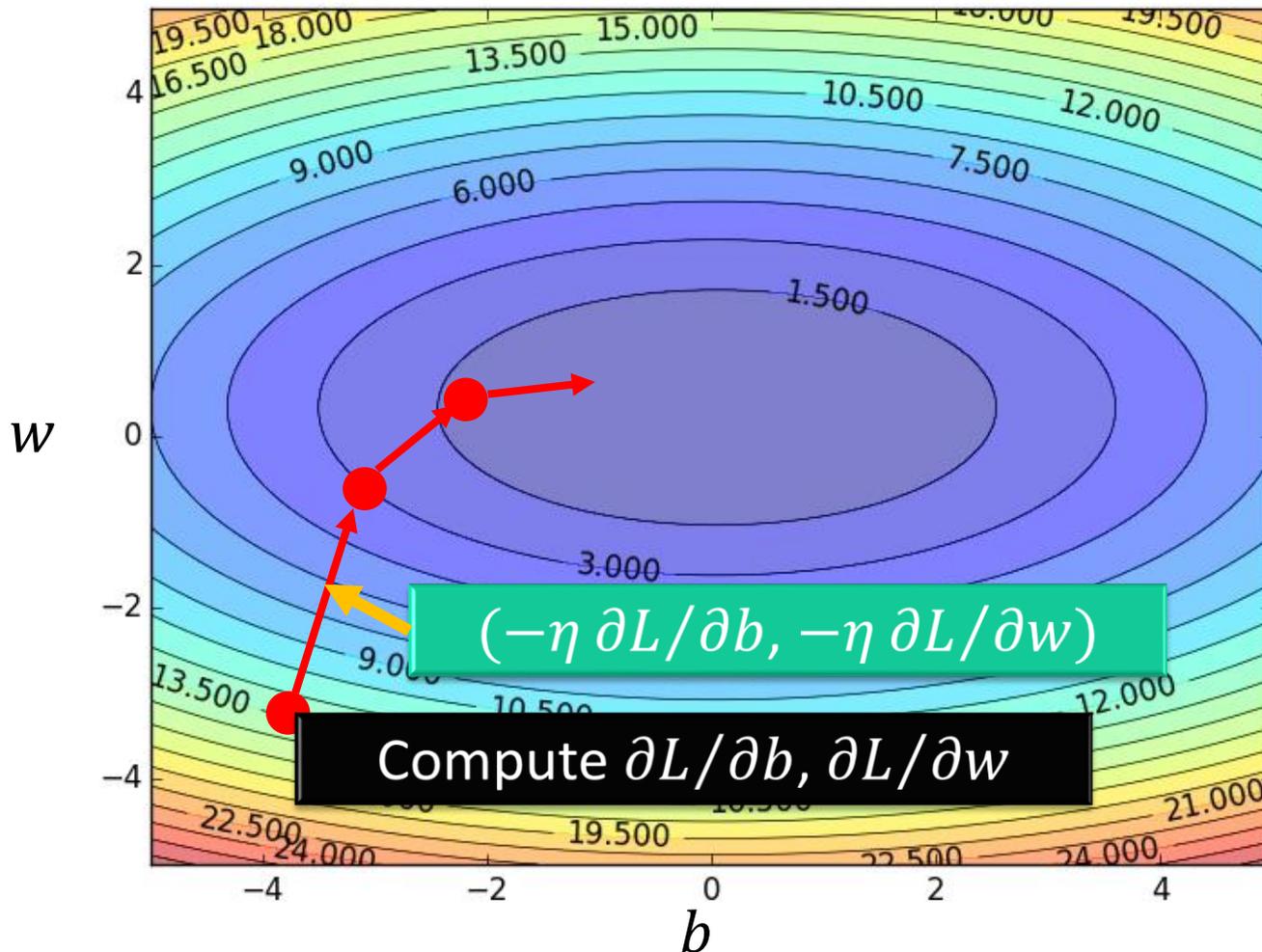
$$w^2 \leftarrow w^1 - \eta \frac{\partial L}{\partial w} \Big|_{w=w^1, b=b^1} \quad b^2 \leftarrow b^1 - \eta \frac{\partial L}{\partial b} \Big|_{w=w^1, b=b^1}$$

$$\nabla L = \begin{bmatrix} \frac{\partial L}{\partial w} \\ \frac{\partial L}{\partial b} \end{bmatrix} \text{gradient}$$

Step 3: Gradient Descent

34

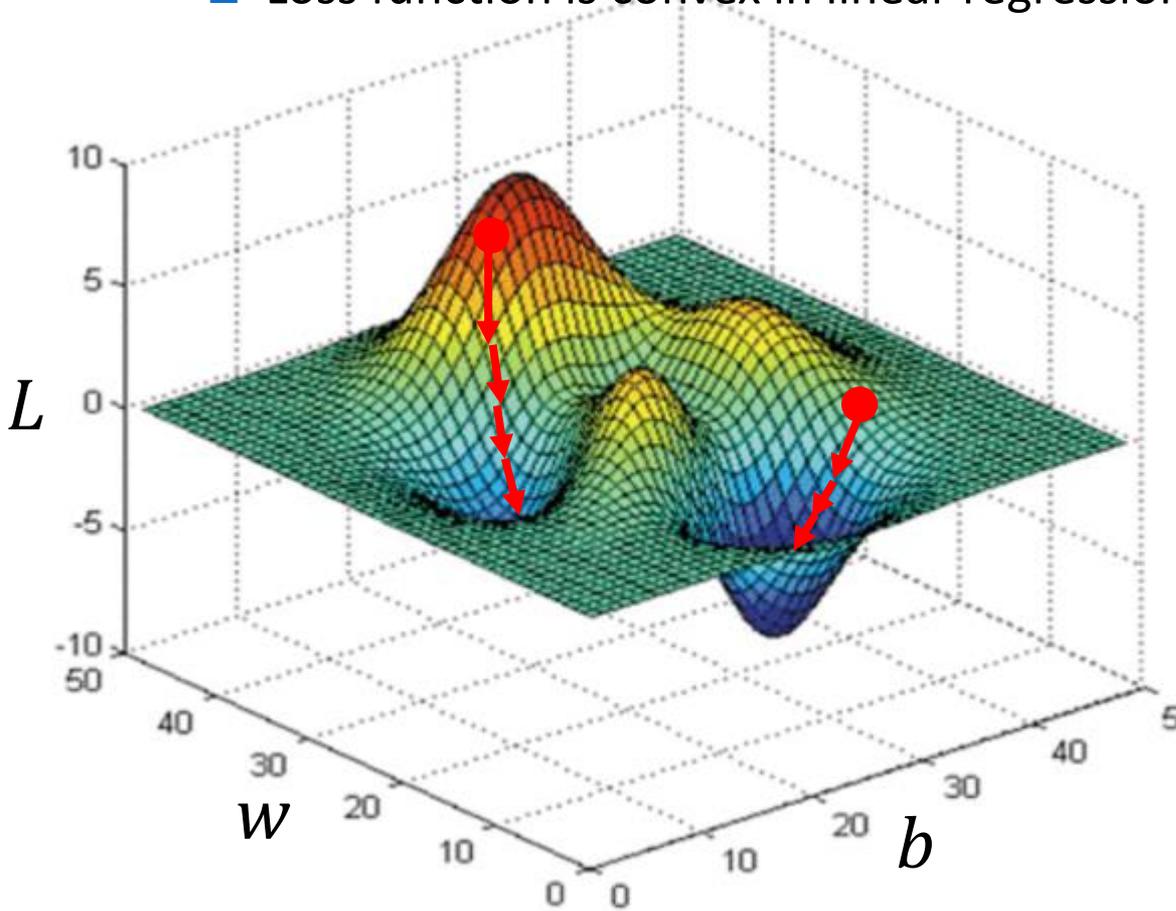
Color: Value of loss $L(w, b)$



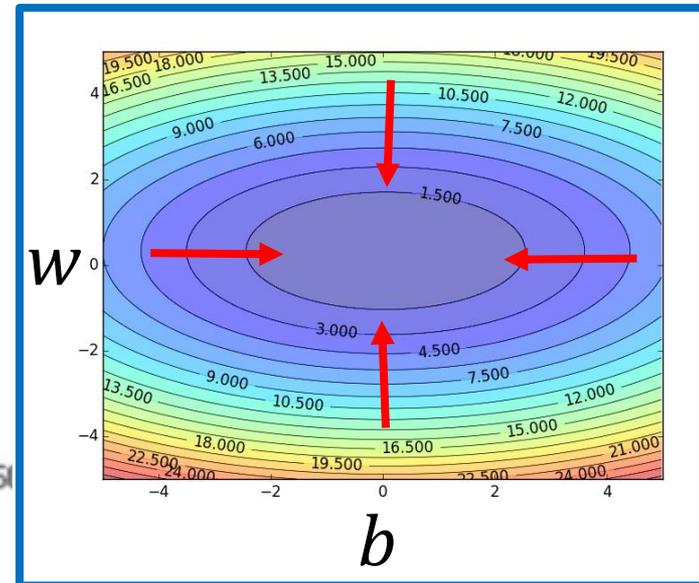
Step 3: Gradient Descent

35

- Local optimal
 - ▣ Loss function is convex in linear regression



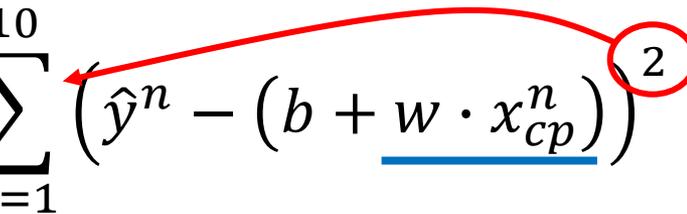
Linear regression →
No local optimal



Step 3: Gradient Descent

36

- Formulation of $\partial L / \partial w$ and $\partial L / \partial b$

$$L(w, b) = \sum_{n=1}^{10} \left(\hat{y}^n - (b + \underline{w \cdot x_{cp}^n}) \right)^2$$


$$\frac{\partial L}{\partial w} =? \sum_{n=1}^{10} 2 \left(\hat{y}^n - (b + w \cdot x_{cp}^n) \right)$$

$$\frac{\partial L}{\partial b} =?$$

Step 3: Gradient Descent

37

- Formulation of $\partial L / \partial w$ and $\partial L / \partial b$

$$L(w, b) = \sum_{n=1}^{10} \left(\hat{y}^n - \underline{(b + w \cdot x_{cp}^n)} \right)^2$$

$$\frac{\partial L}{\partial w} =? \sum_{n=1}^{10} 2 \left(\hat{y}^n - (b + w \cdot x_{cp}^n) \right) (-x_{cp}^n)$$

$$\frac{\partial L}{\partial b} =? \sum_{n=1}^{10} 2 \left(\hat{y}^n - (b + w \cdot x_{cp}^n) \right)$$

Learned Model

38

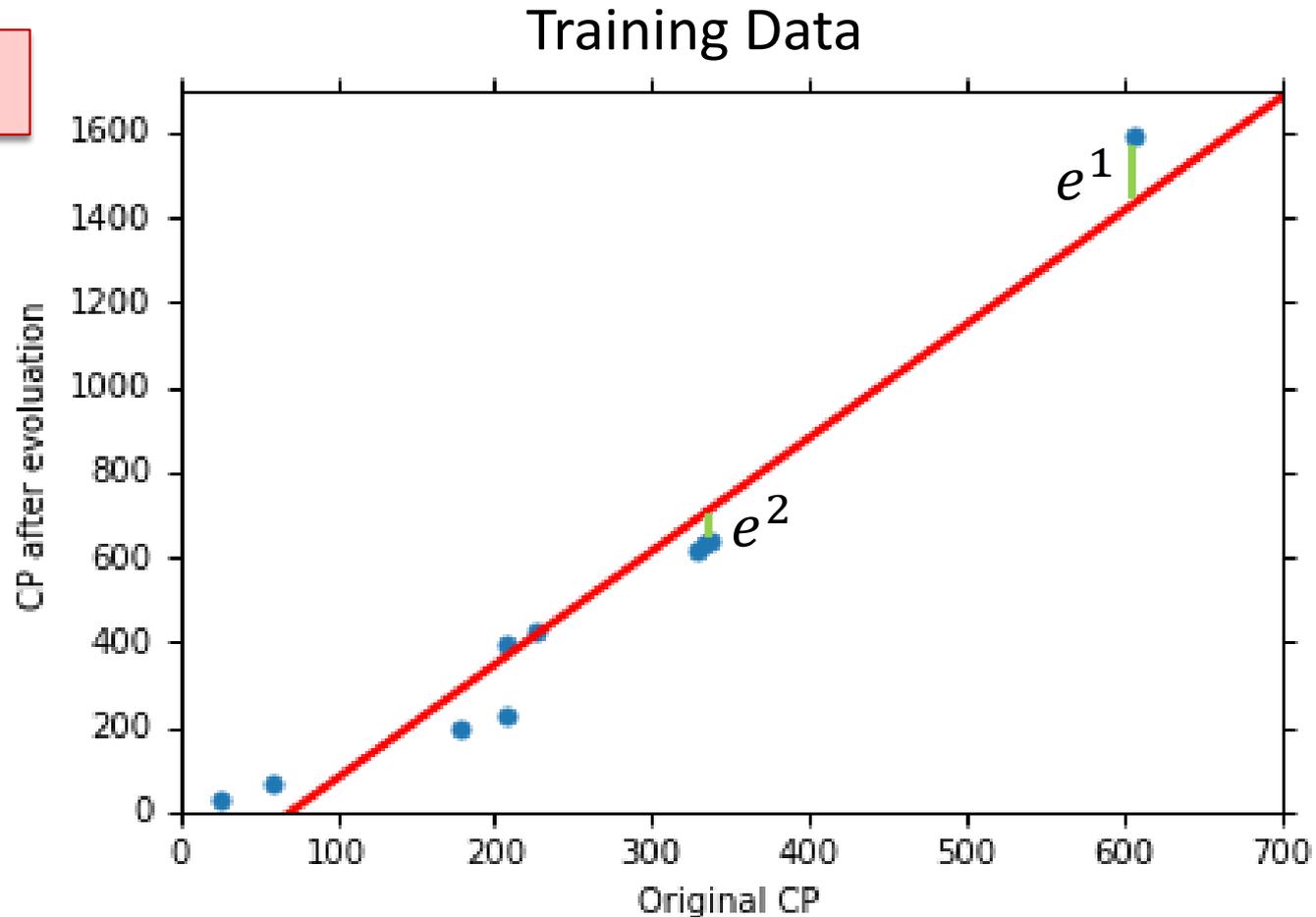
$$y = b + w \cdot x_{cp}$$

$$b = -188.4$$

$$w = 2.7$$

Average Error on
Training Data

$$= \sum_{n=1}^{10} e^n = 31.9$$



Model Generalization

What we really care about is the error on new data (testing data)

39

$$y = b + w \cdot x_{cp}$$

$$b = -188.4$$

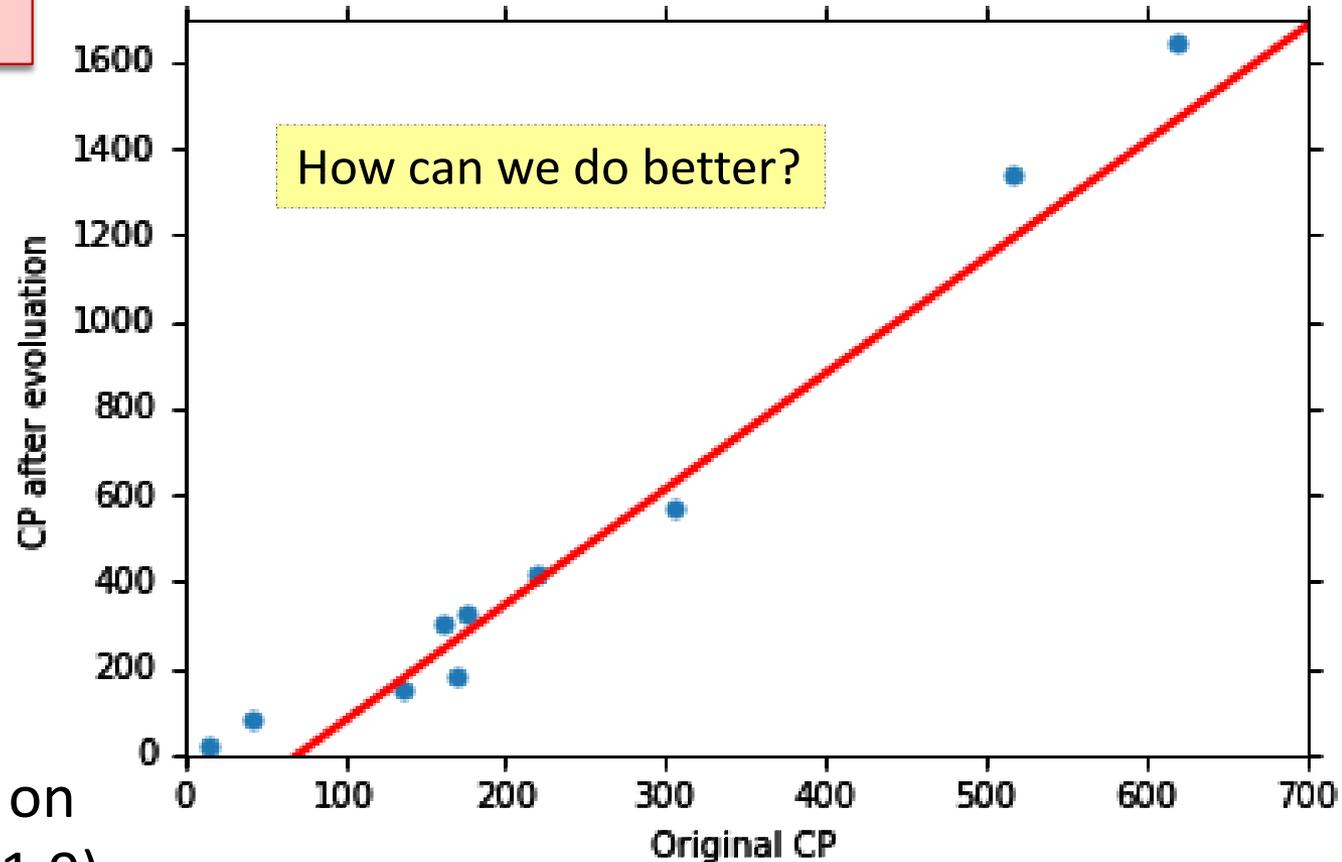
$$w = 2.7$$

Average Error on Testing Data

$$= \sum_{n=1}^{10} e^n = 35.0$$

> Average Error on Training Data (31.9)

Another 10 pokemons as testing data



Model Generalization

40

- Select another model

$$y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2$$

- Best function

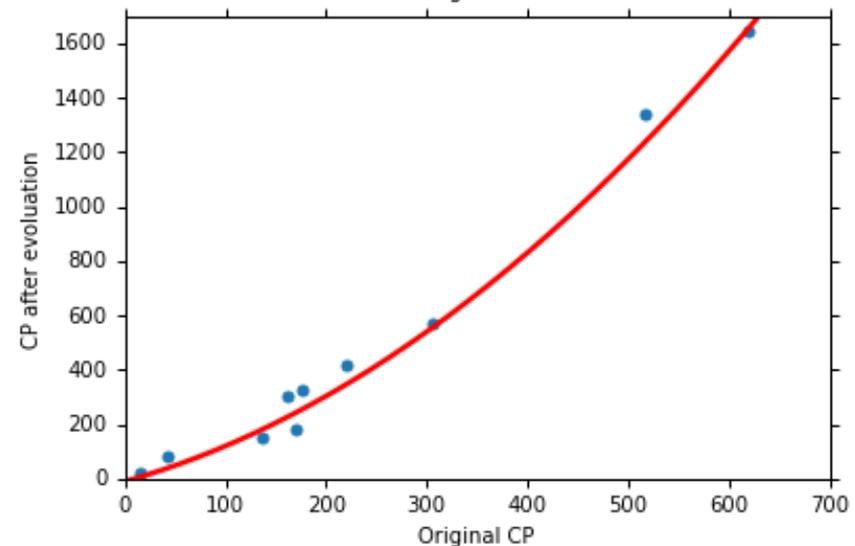
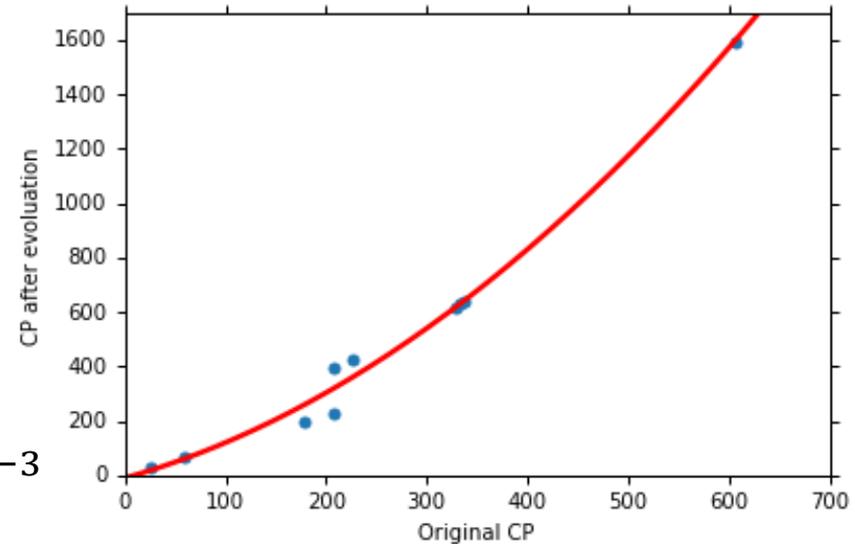
$$b = -10.3, w_1 = 1.0, w_2 = 2.7 \times 10^{-3}$$

Average Error = 15.4

- Testing

Average Error = 18.4

Better! Could it be even better?



Model Generalization

41

- Select another model

$$y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3$$

- Best function

$$b = 6.4, w_1 = 0.66, w_2 = 4.3 \times 10^{-3}, w_3 = 1.8 \times 10^{-6}$$

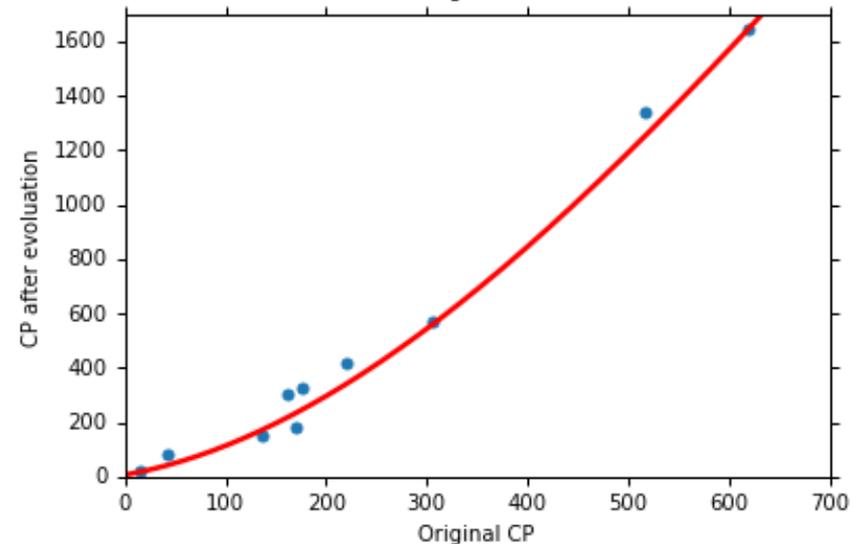
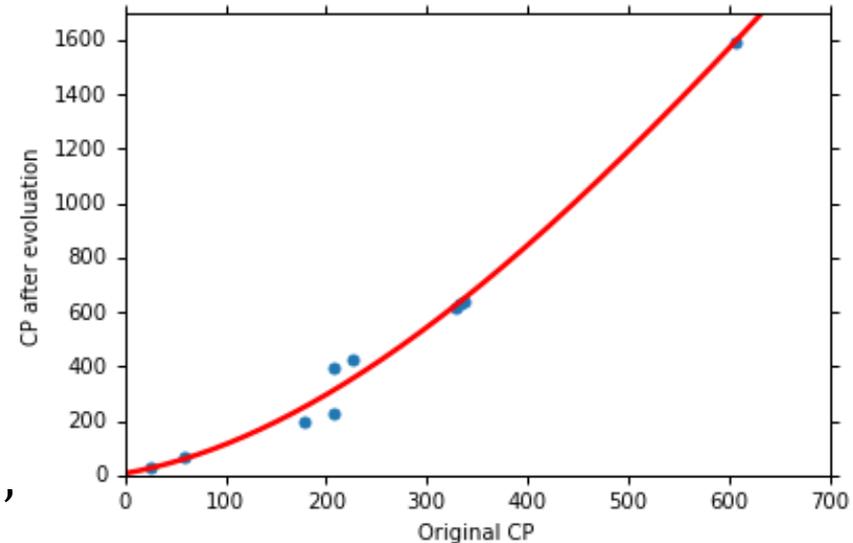
Average Error = 15.3

- Testing

Average Error = 18.1

Slightly better.

How about more complex model?



Model Generalization

42

- Select another model

$$y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3 + w_4 \cdot (x_{cp})^4$$

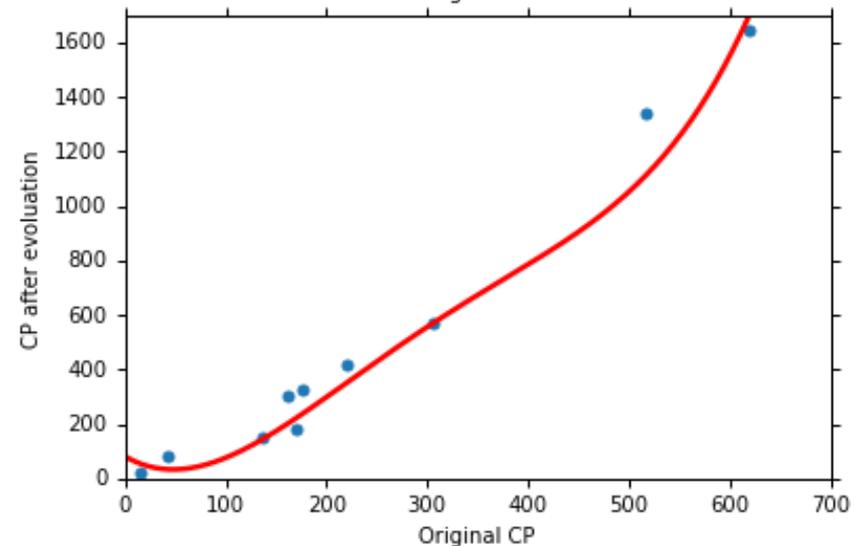
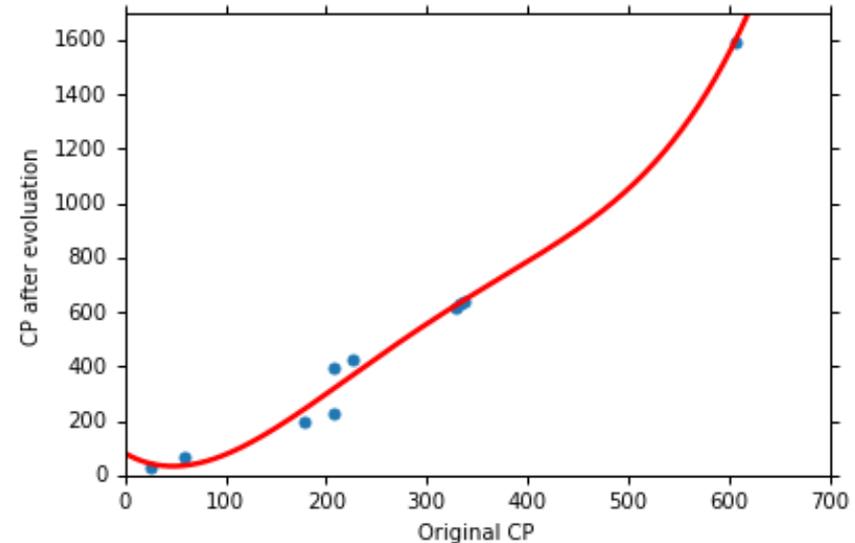
- Best function

Average Error = 14.9

- Testing

Average Error = 28.8

The results become worse



Model Generalization

43

- Select another model

$$y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3 + w_4 \cdot (x_{cp})^4 + w_5 \cdot (x_{cp})^5$$

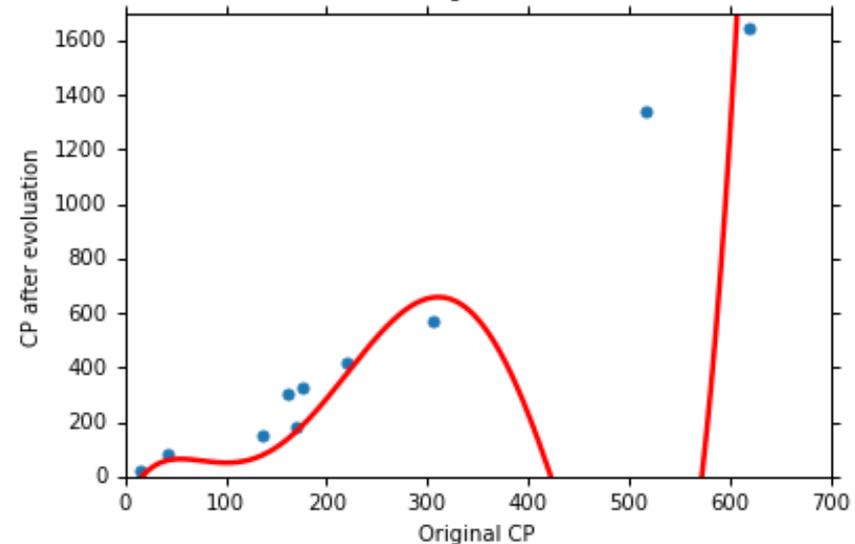
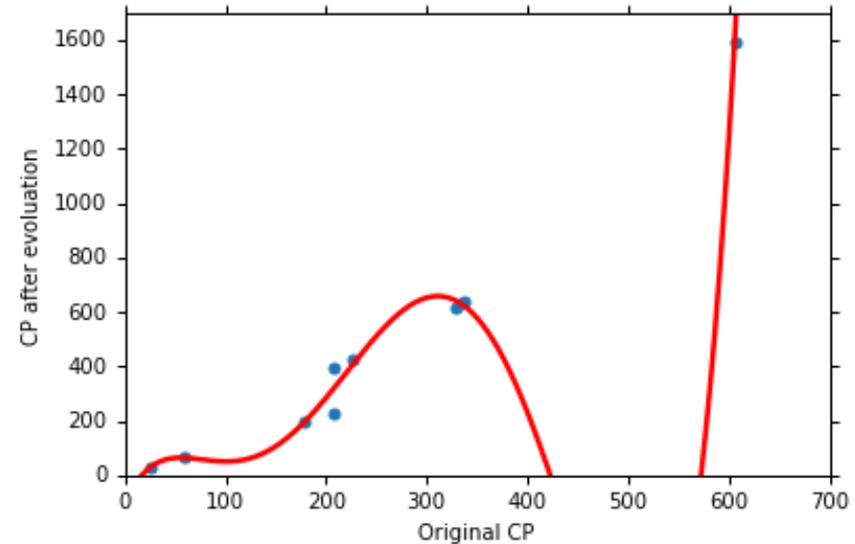
- Best function

Average Error = 12.8

- Testing

Average Error = 232.1

The results are so bad



Model Selection

44

1. $y = b + w \cdot x_{cp}$

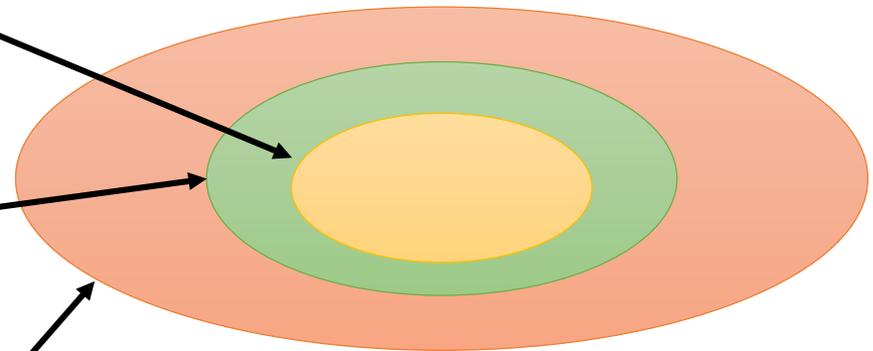
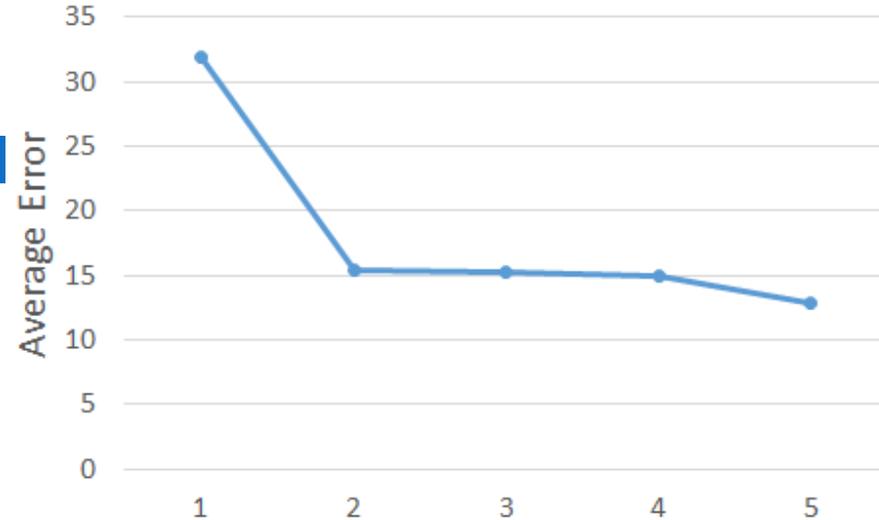
2. $y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2$

3. $y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3$

4. $y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3 + w_4 \cdot (x_{cp})^4$

5. $y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3 + w_4 \cdot (x_{cp})^4 + w_5 \cdot (x_{cp})^5$

Training Data



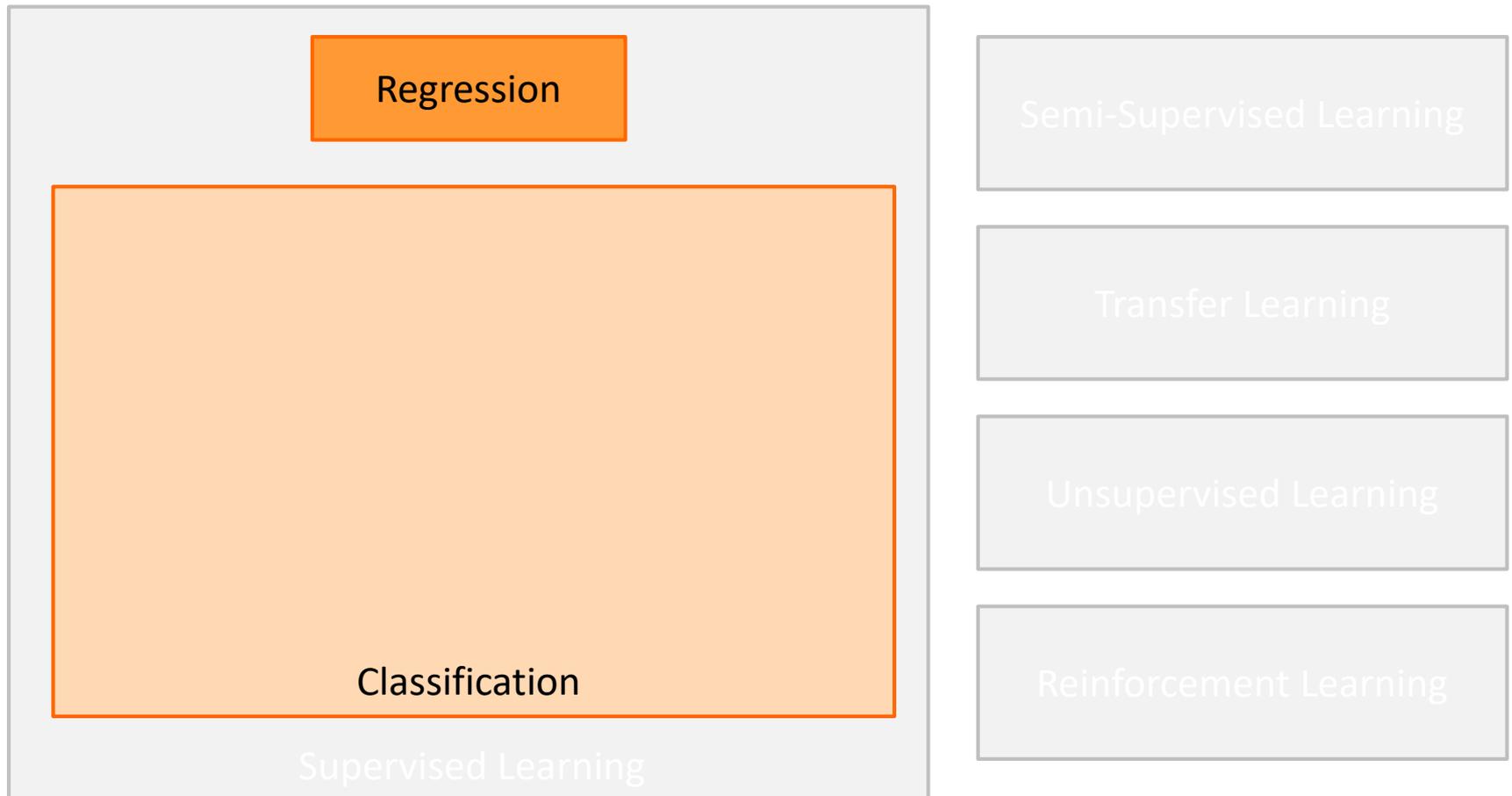
A more complex model yields lower error on training data.

If we can truly find the best function

Machine Learning Map

45

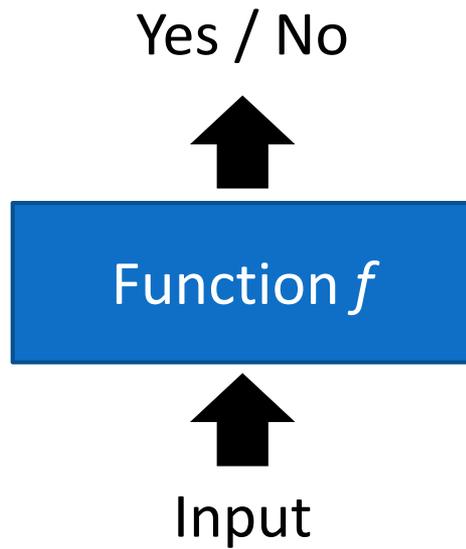
Scenario Task Method



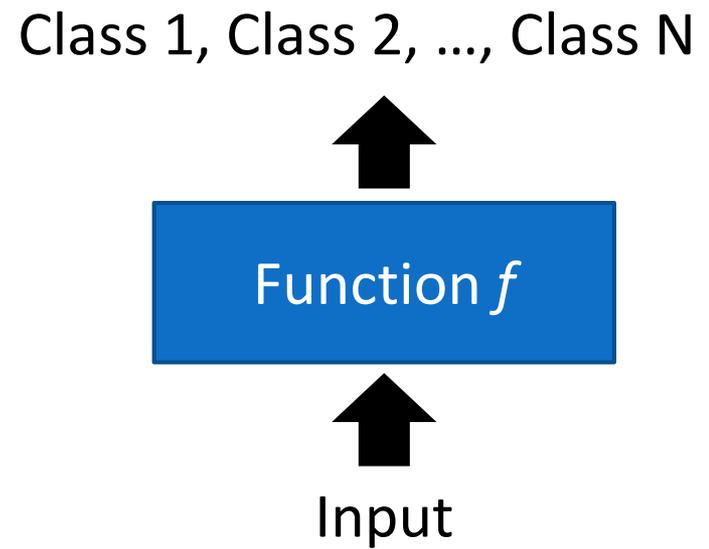
Classification

46

□ Binary Classification

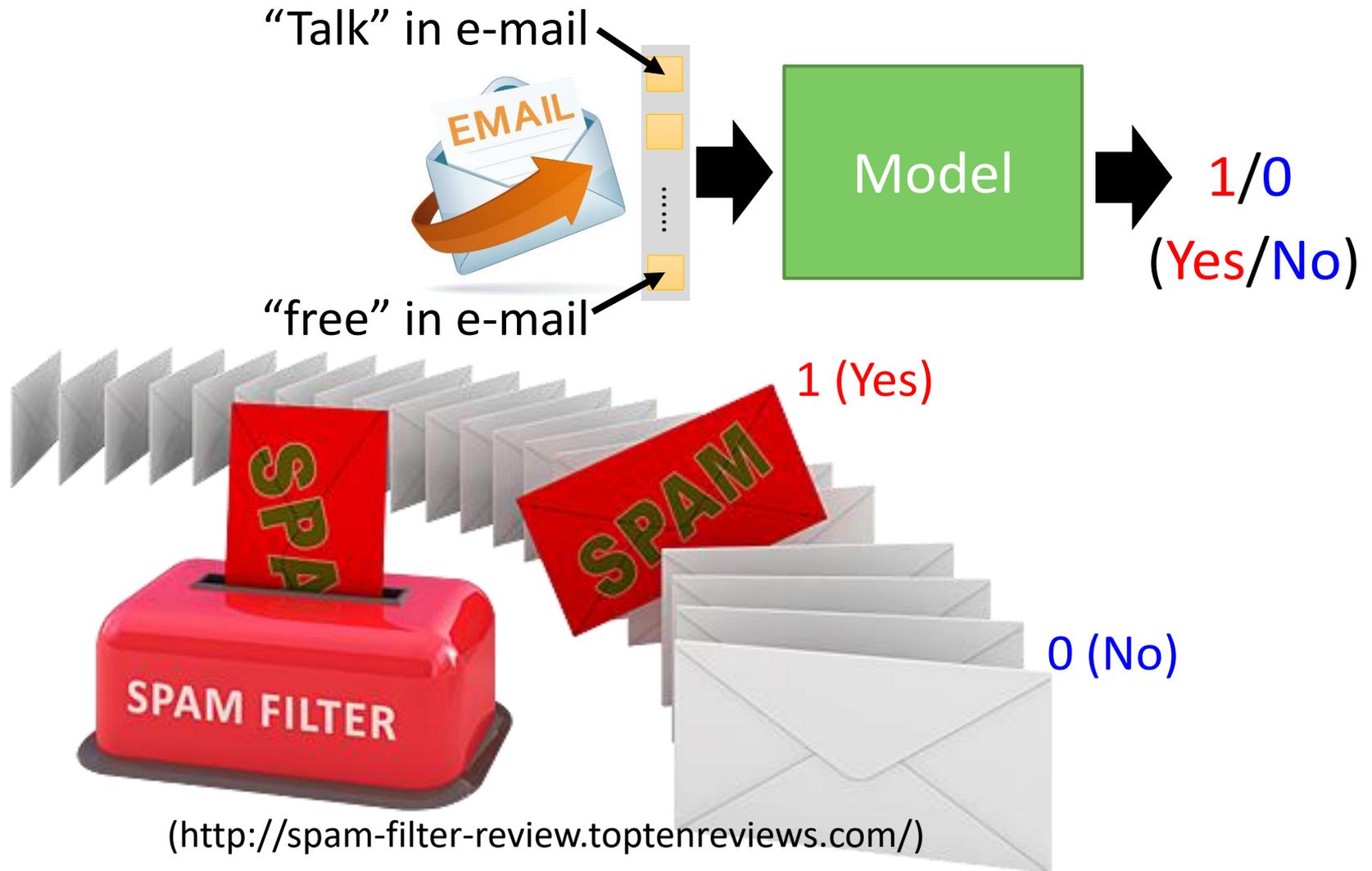


□ Multi-Class Classification



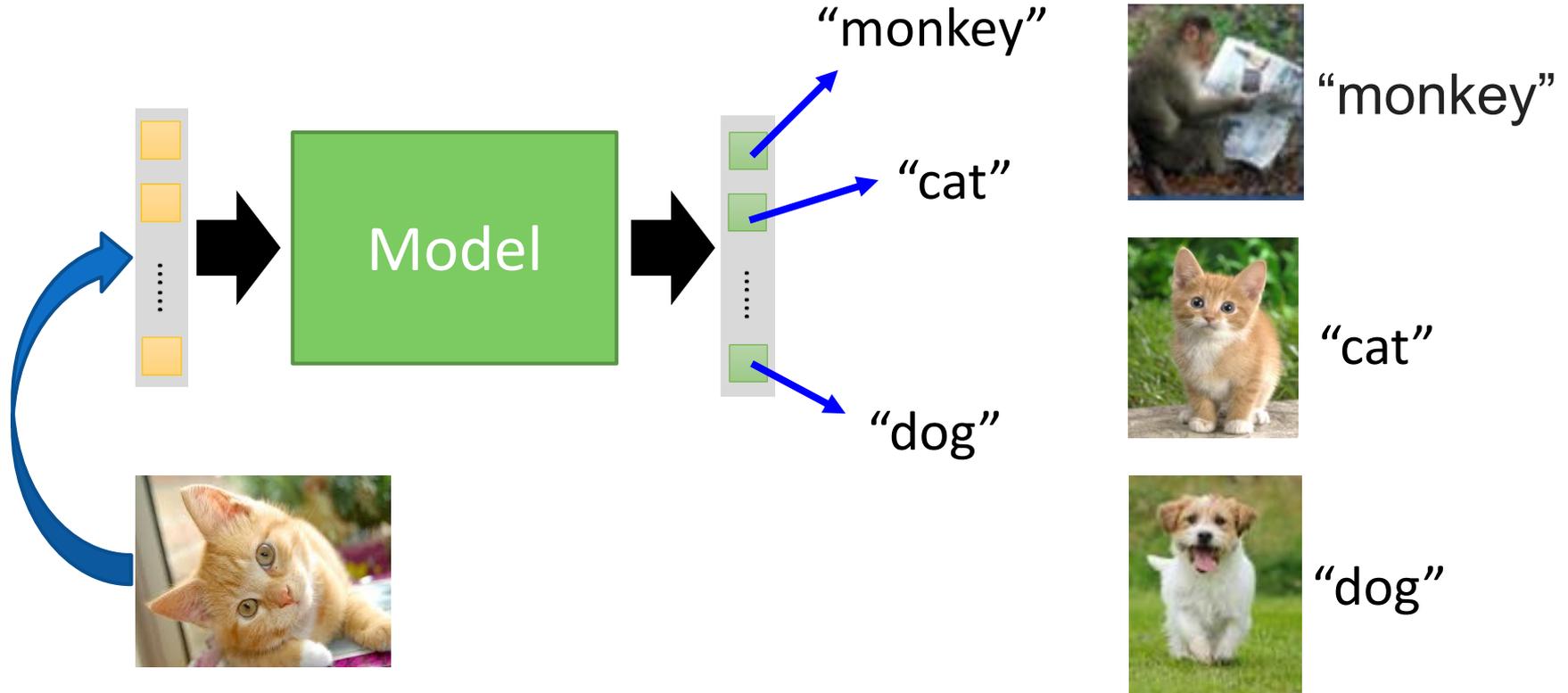
Binary Classification – Spam Filtering

47



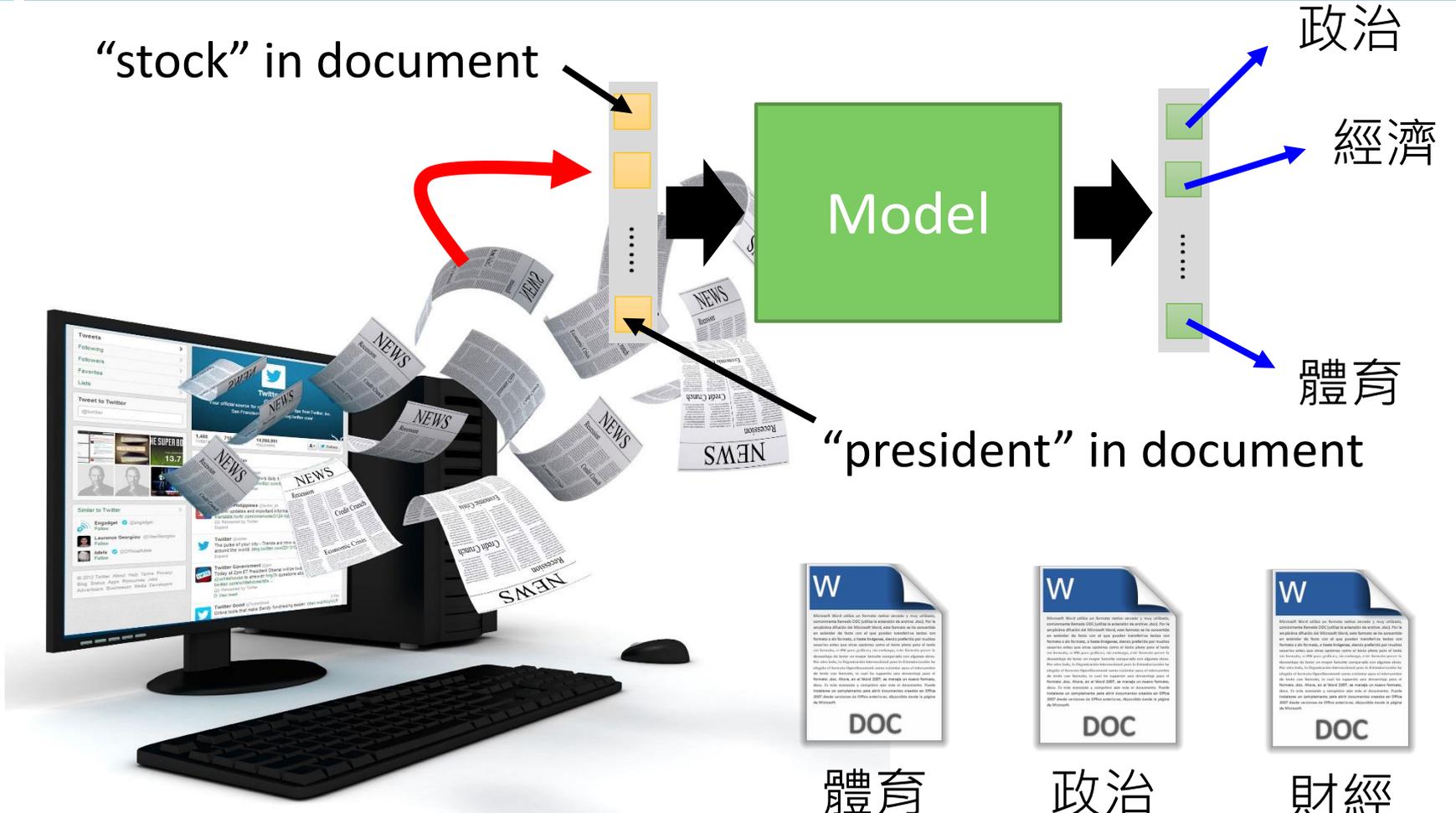
Multi-Class – Image Recognition

48



Multi-Class – Topic Classification

“stock” in document



政治

經濟

體育

“president” in document



體育



政治



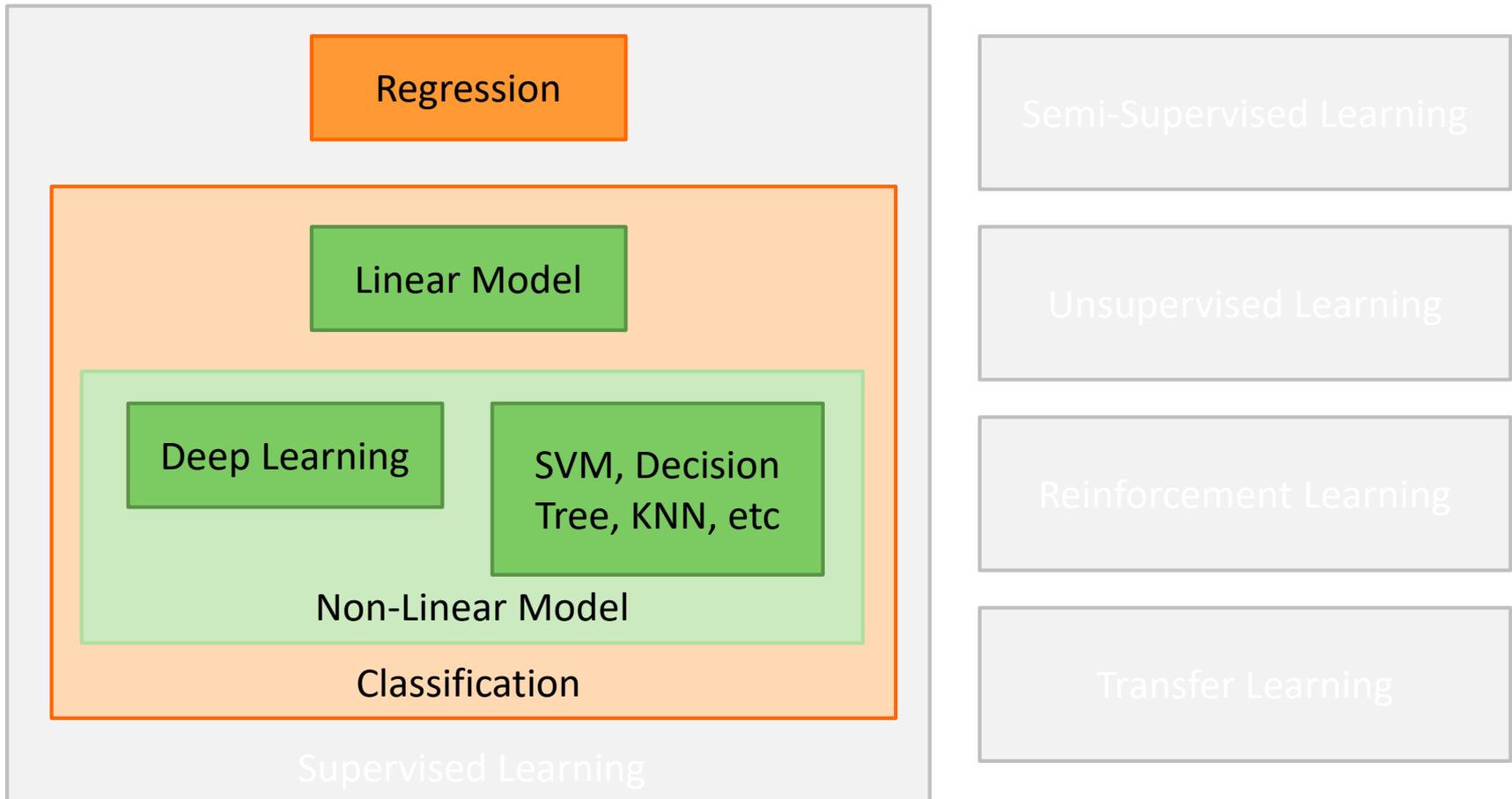
財經

<http://top-breaking-news.com/>

Machine Learning Map

50

Scenario Task Method



Part I: Introduction to ML & DL

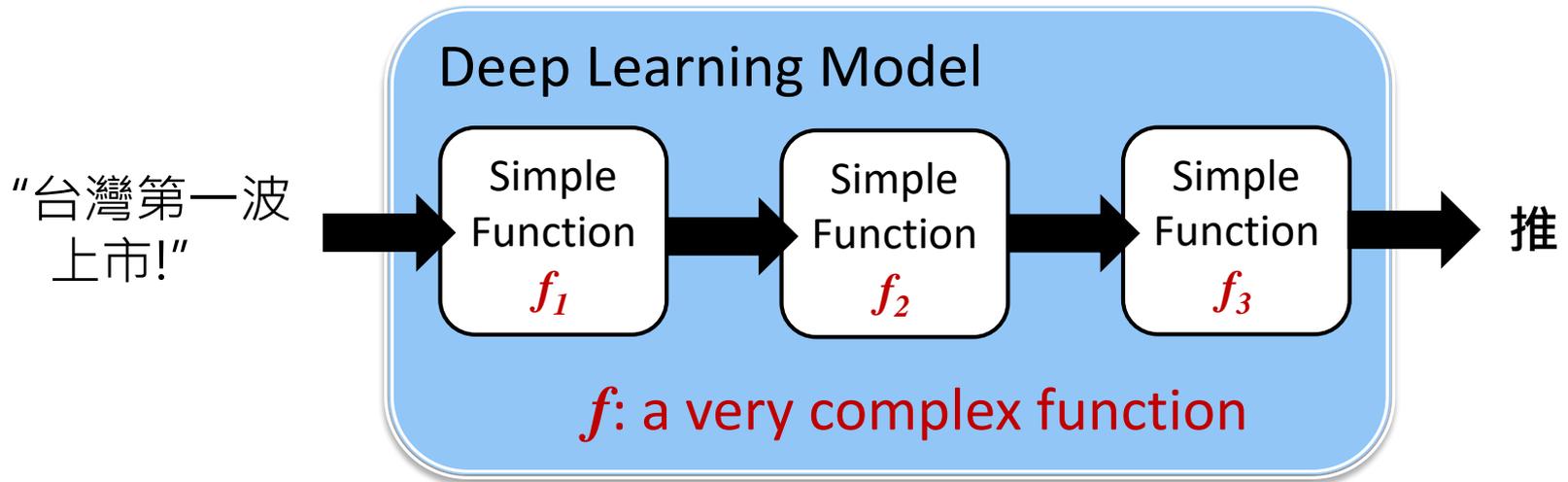
51

- Basic Machine Learning
- **Basic Deep Learning**
- Toolkits and Learning Recipe

Stacked Functions Learned by Machine

52

□ Production line (生產線)



End-to-end training: what each function should do is learned automatically

Deep learning usually refers to *neural network* based model

Three Steps for Deep Learning

53

Step 1: define a set of function



Step 2: goodness of function



Step 3: pick the best function

Three Steps for Deep Learning

54

Step 1: define a set of function

Neural Network

Step 2: goodness of function

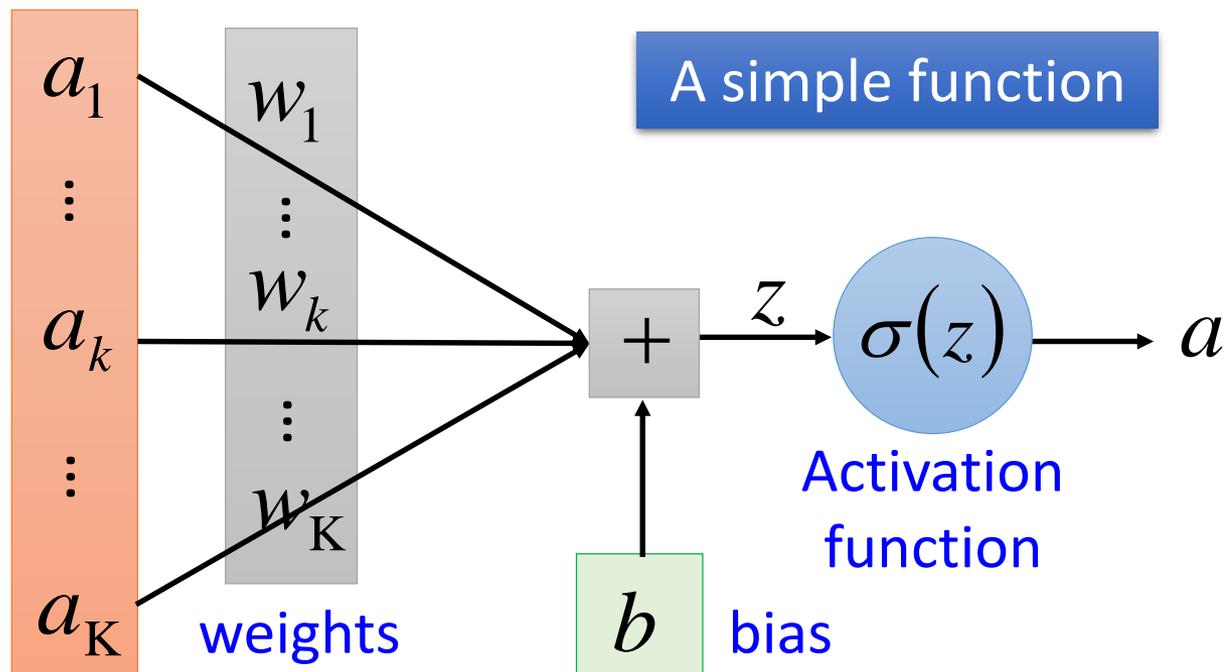
Step 3: pick the best function

Neural Network

55

Neuron

$$z = a_1 w_1 + \dots + a_k w_k + \dots + a_K w_K + b$$

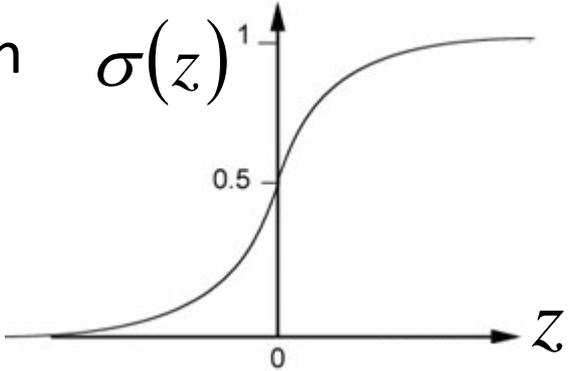


Neural Network

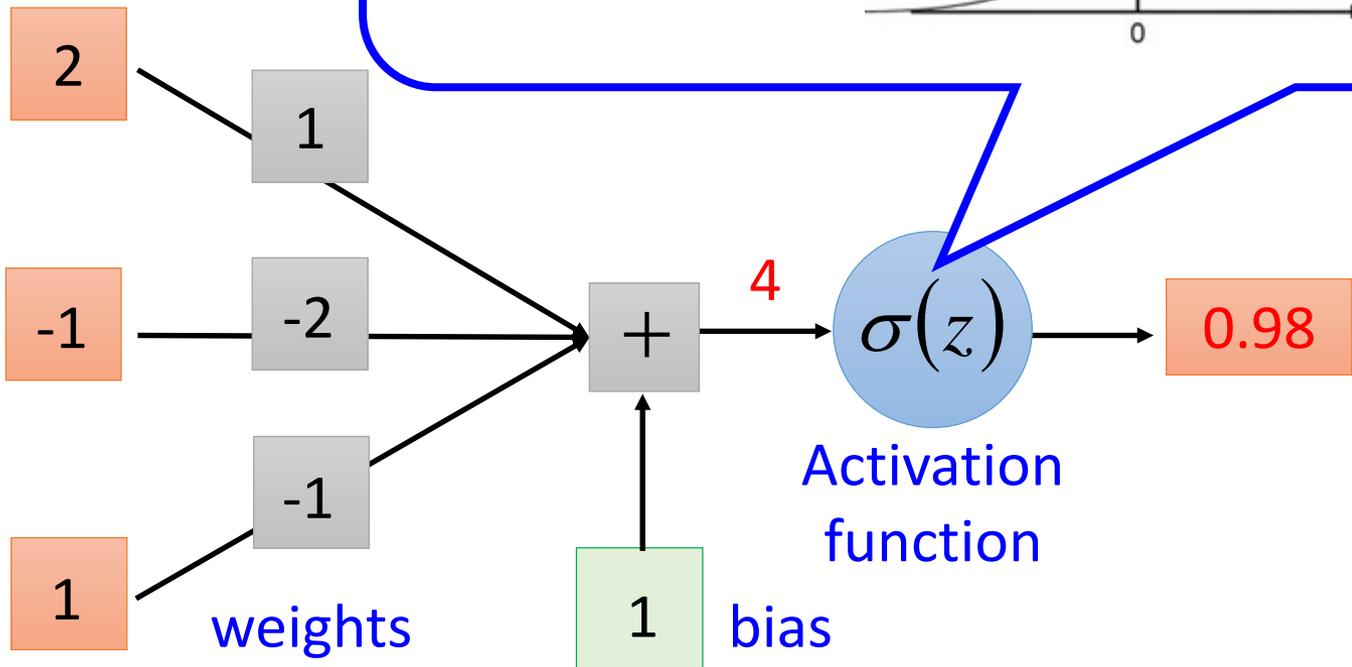
56

Neuron

Sigmoid Function $\sigma(z)$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$


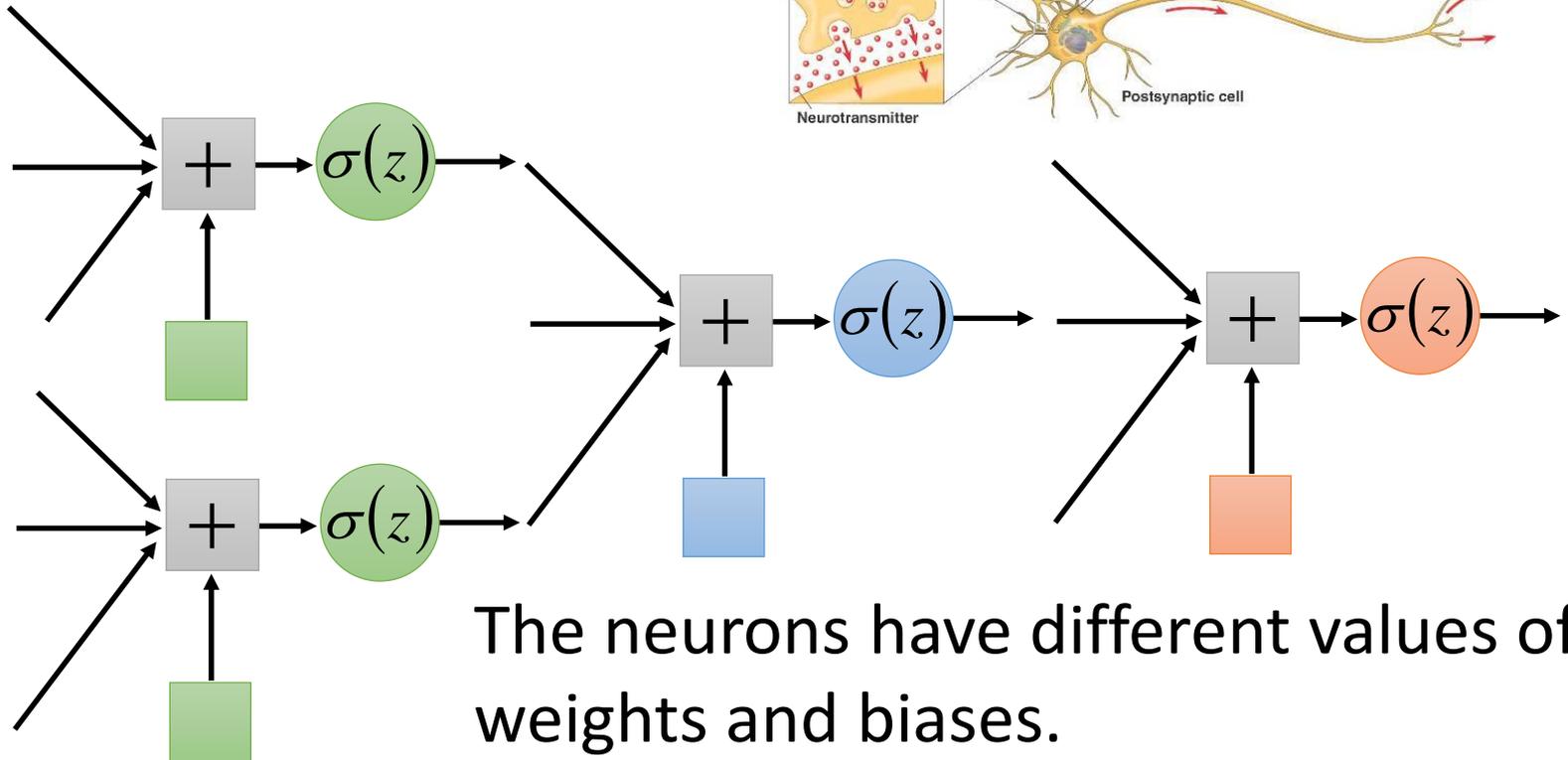
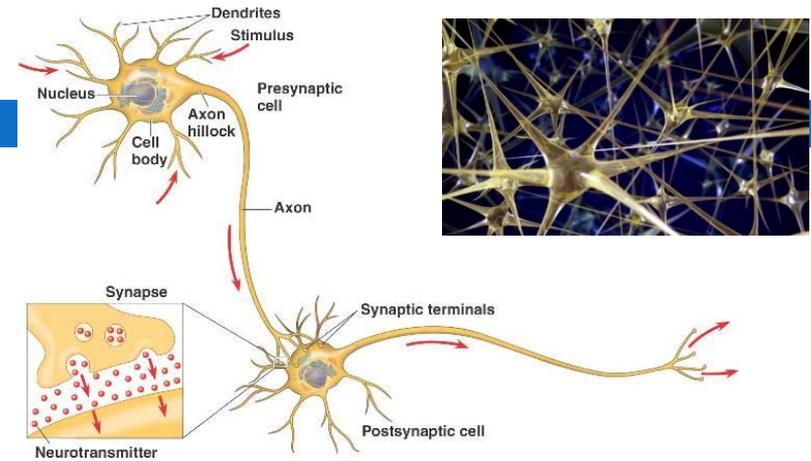
The graph shows the Sigmoid Function $\sigma(z)$ plotted against z . The curve is S-shaped, passing through the point (0, 0.5) and approaching 1 as z increases. The y-axis is labeled with 0.5 and 1, and the x-axis is labeled with 0 and z .



Neural Network

57

Different connections lead to different network structures

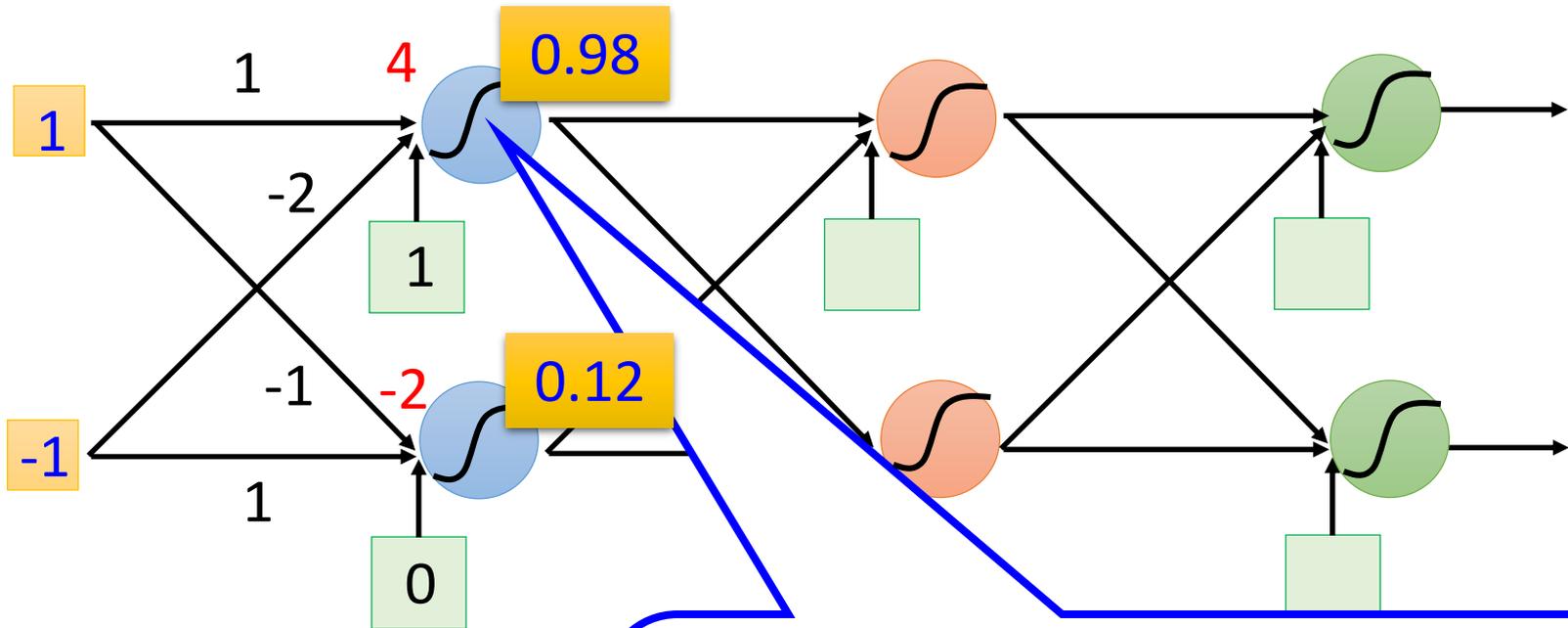


The neurons have different values of weights and biases.

Weights and biases are network parameters θ

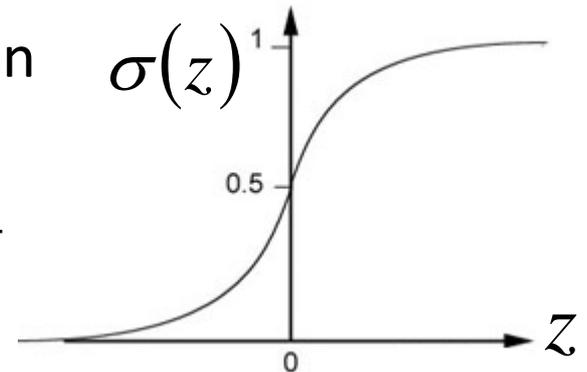
Fully Connected Feedforward Network

58



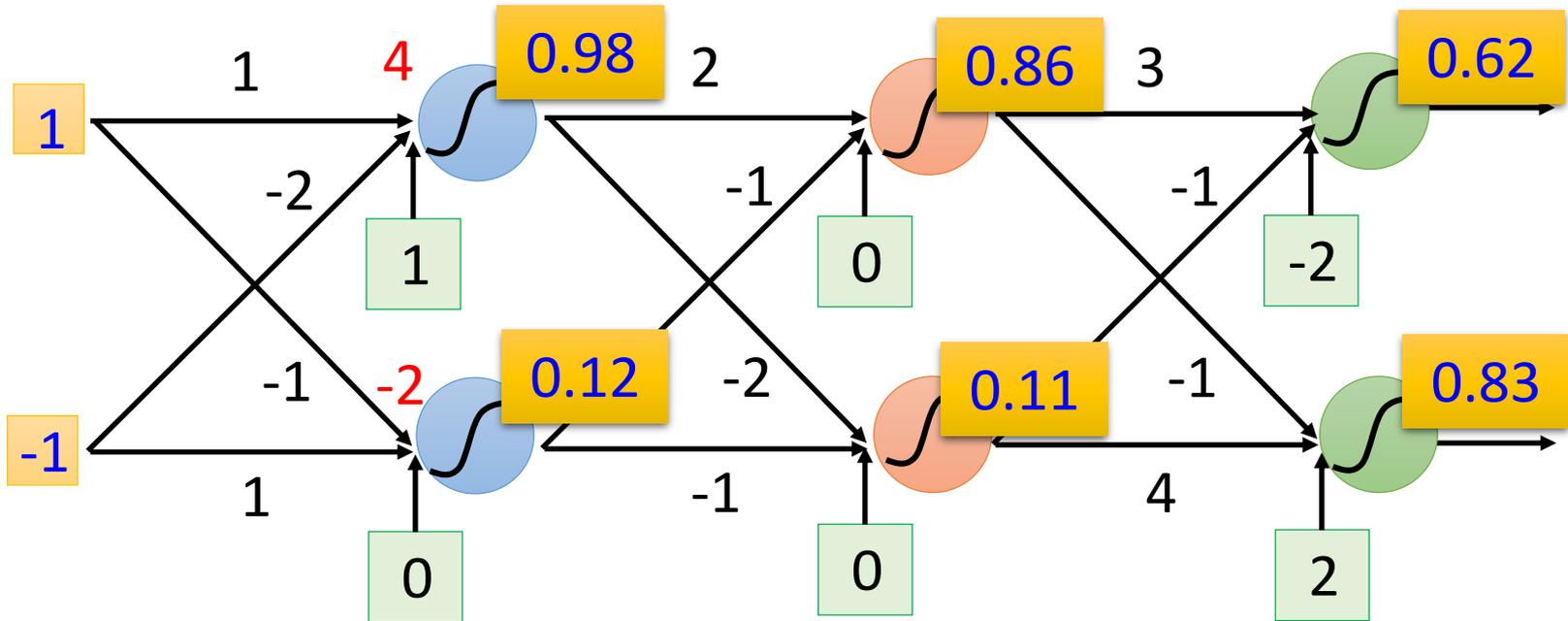
Sigmoid Function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



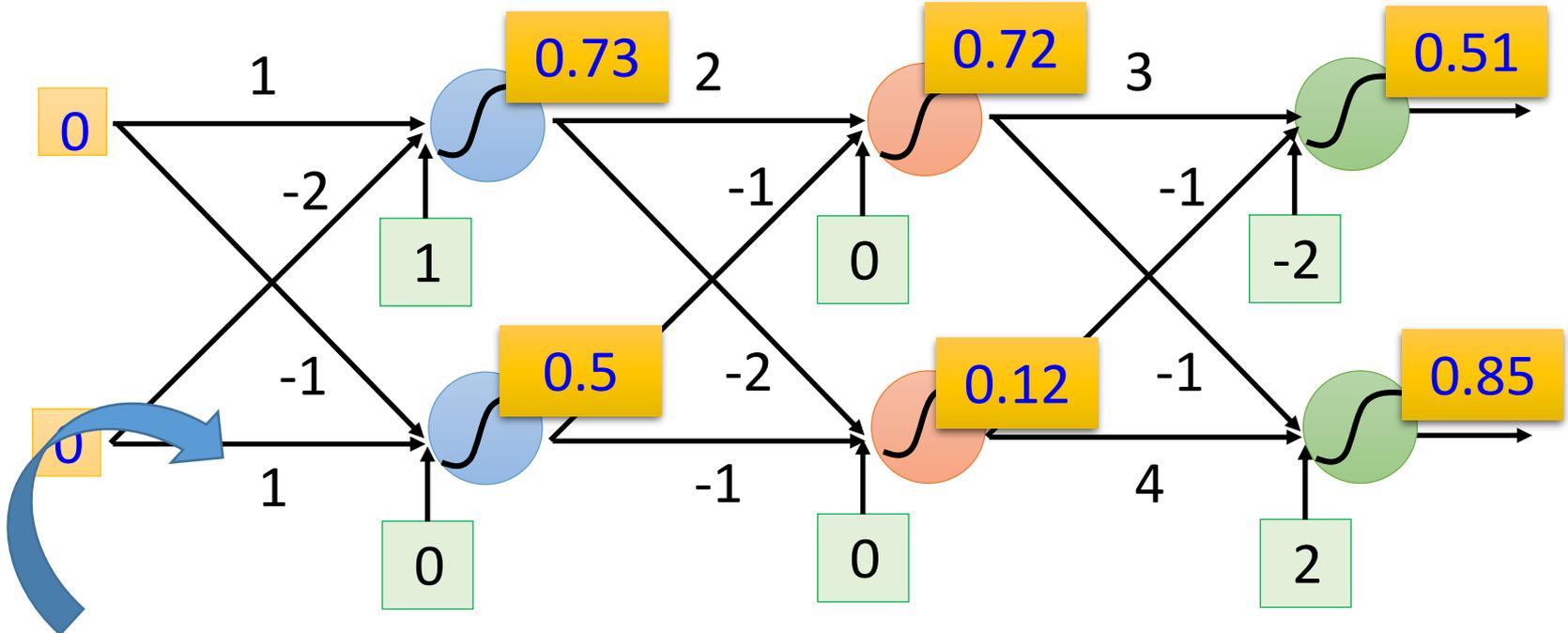
Fully Connected Feedforward Network

59



Fully Connected Feedforward Network

60



This is a function.

Input vector, output vector

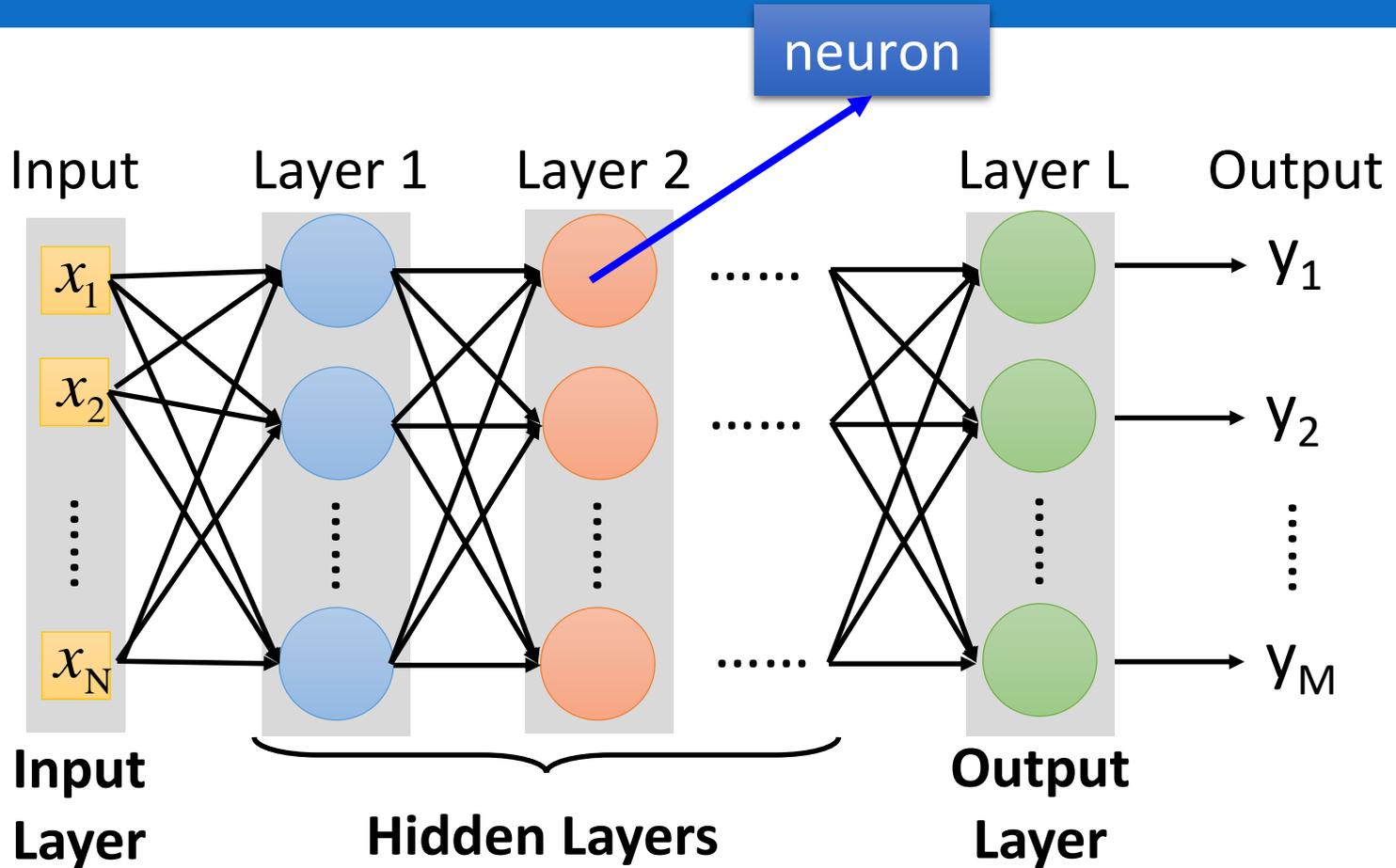
$$f\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix}\right) = \begin{bmatrix} 0.62 \\ 0.83 \end{bmatrix} \quad f\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}\right) = \begin{bmatrix} 0.51 \\ 0.85 \end{bmatrix}$$

Given parameters θ , define a function

Given network structure, define a function set

Fully Connect Feedforward Network

61



Deep means many hidden layers

Why Deep? Universality Theorem

62

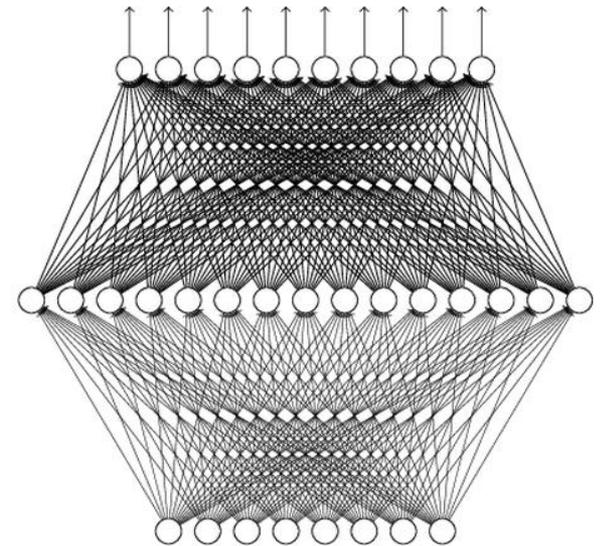
- Any continuous function f

$$f : \mathbb{R}^N \rightarrow \mathbb{R}^M$$

can be realized by a network with only hidden layer

- (given **enough** hidden neurons)

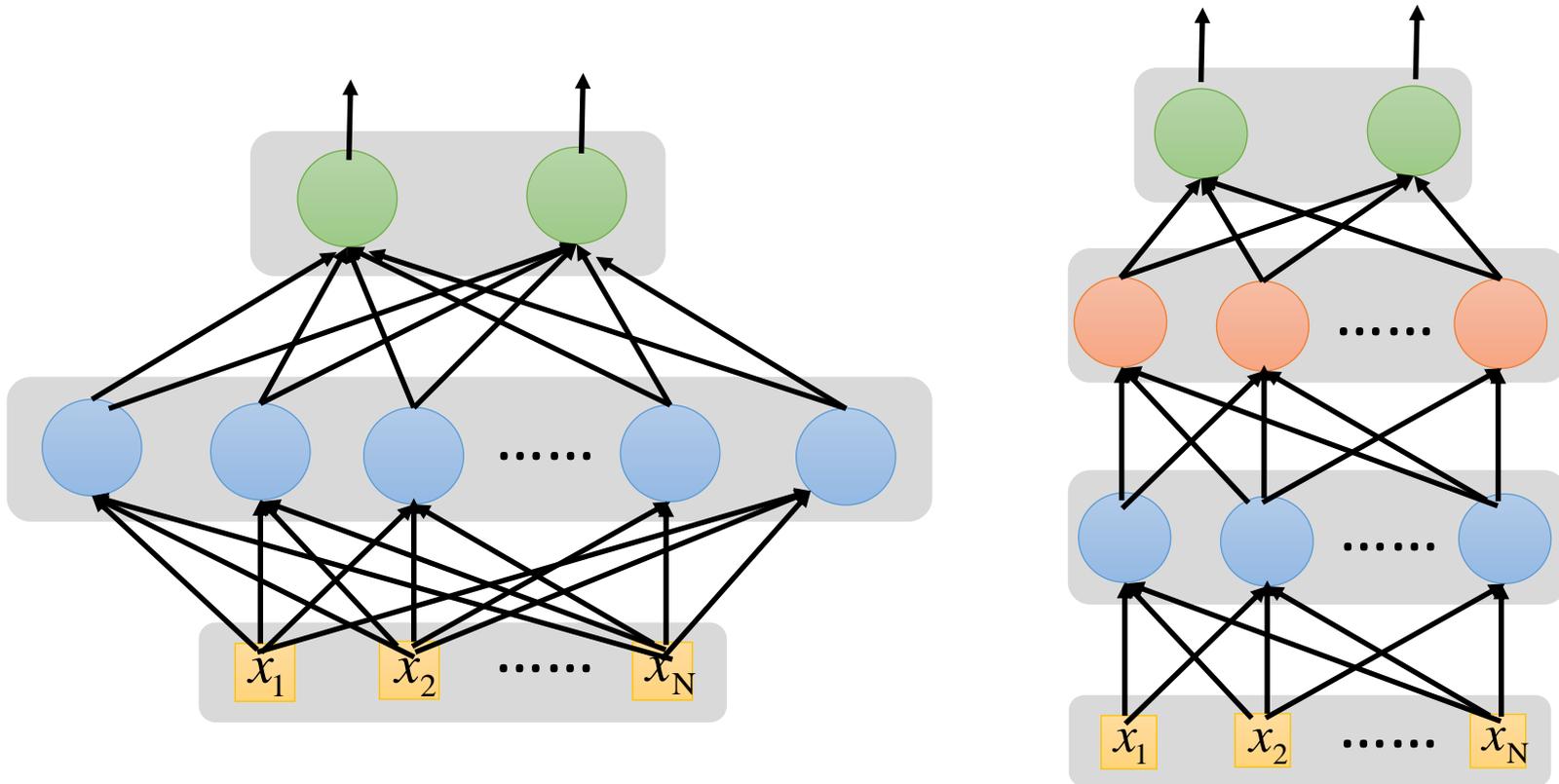
Why “deep” not “fat”?



Fat + Shallow v.s. Thin + Deep

63

- Two networks with the same number of parameters



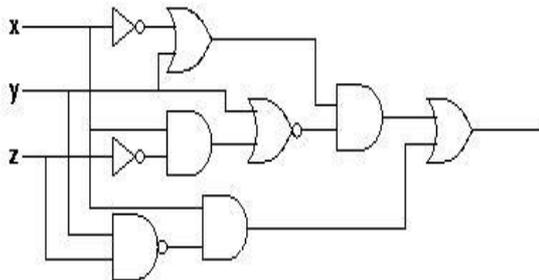
Why Deep

64

- Logic circuits
 - ▣ Consists of **gates**
 - ▣ **A two layers of logic gates** can represent **any Boolean function**.
 - ▣ Using multiple layers of logic gates to build some functions are much simpler



less gates needed



- Neural network
 - ▣ consists of **neurons**
 - ▣ **A hidden layer network** can represent **any continuous function**.
 - ▣ Using multiple layers of neurons to represent some functions are much simpler



less parameters



less data?

Deep = Many Hidden Layers

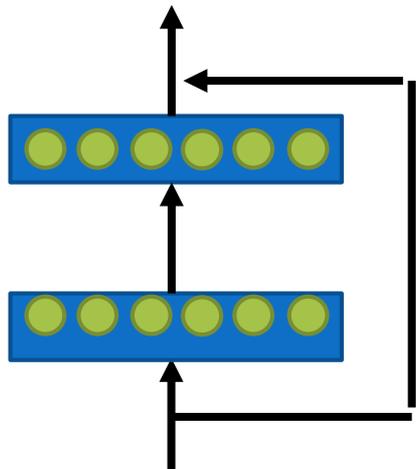
66

101 layers

152 layers

3.57%

Special structure



16.4%

AlexNet
(2012)

7.3%

VGG
(2014)

6.7%

GoogleNet
(2014)

Residual Net
(2015)

Taipei
101



Output Layer

67

- Softmax layer as the output layer

Ordinary Layer

$$z_1 \longrightarrow \sigma \longrightarrow y_1 = \sigma(z_1)$$

$$z_2 \longrightarrow \sigma \longrightarrow y_2 = \sigma(z_2)$$

$$z_3 \longrightarrow \sigma \longrightarrow y_3 = \sigma(z_3)$$

In general, the output of network can be any value.

May not be easy to interpret

Output Layer

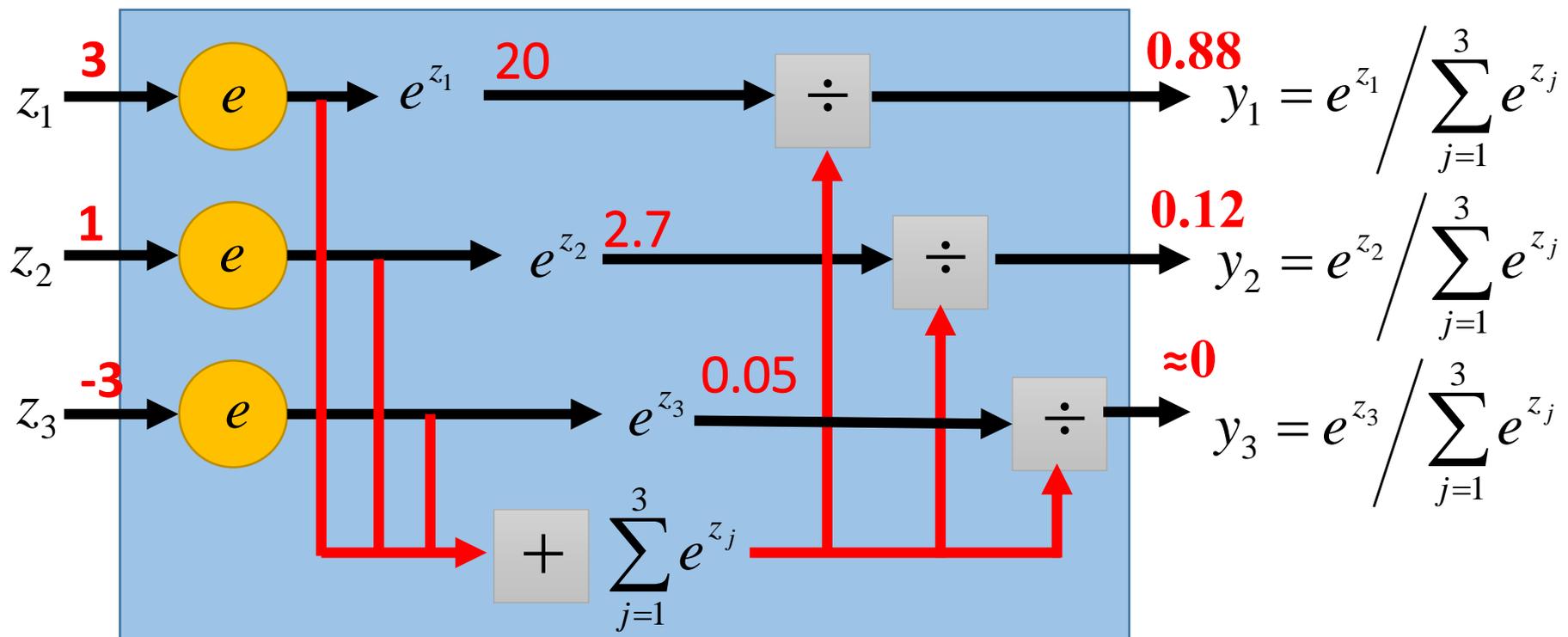
68

- Softmax layer as the output layer

Probability:

- $1 > y_i > 0$
- $\sum_i y_i = 1$

Softmax Layer

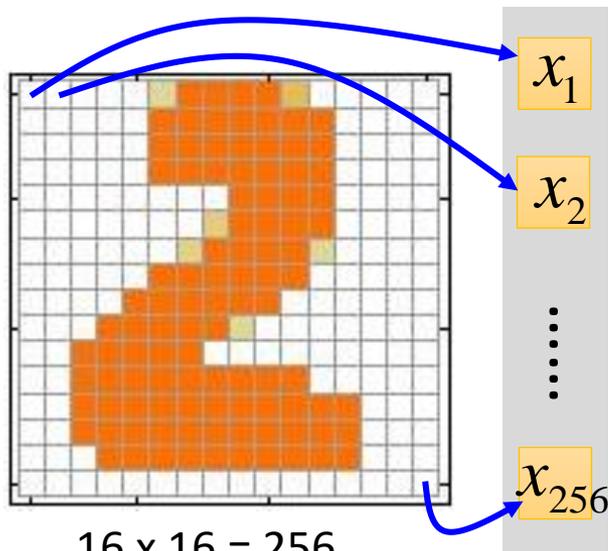


Example Application



69

□ Input

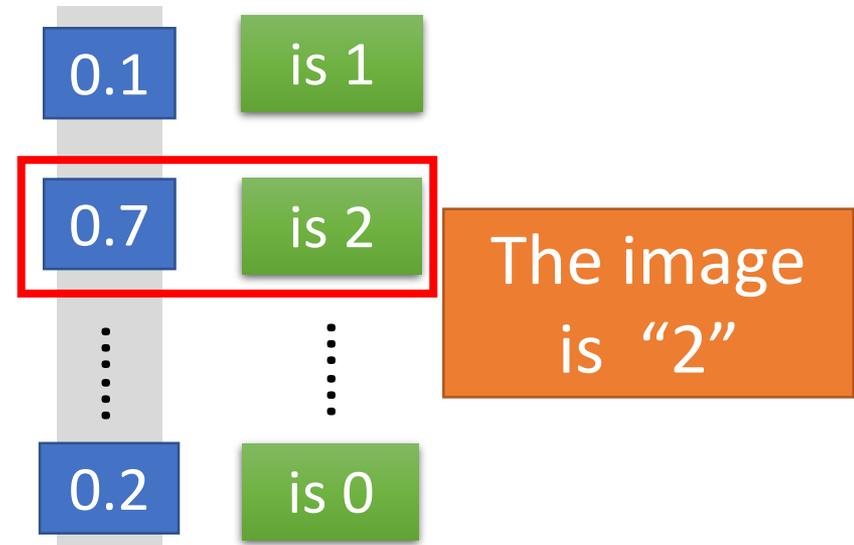


16 x 16 = 256

Ink \rightarrow 1

No ink \rightarrow 0

□ Output

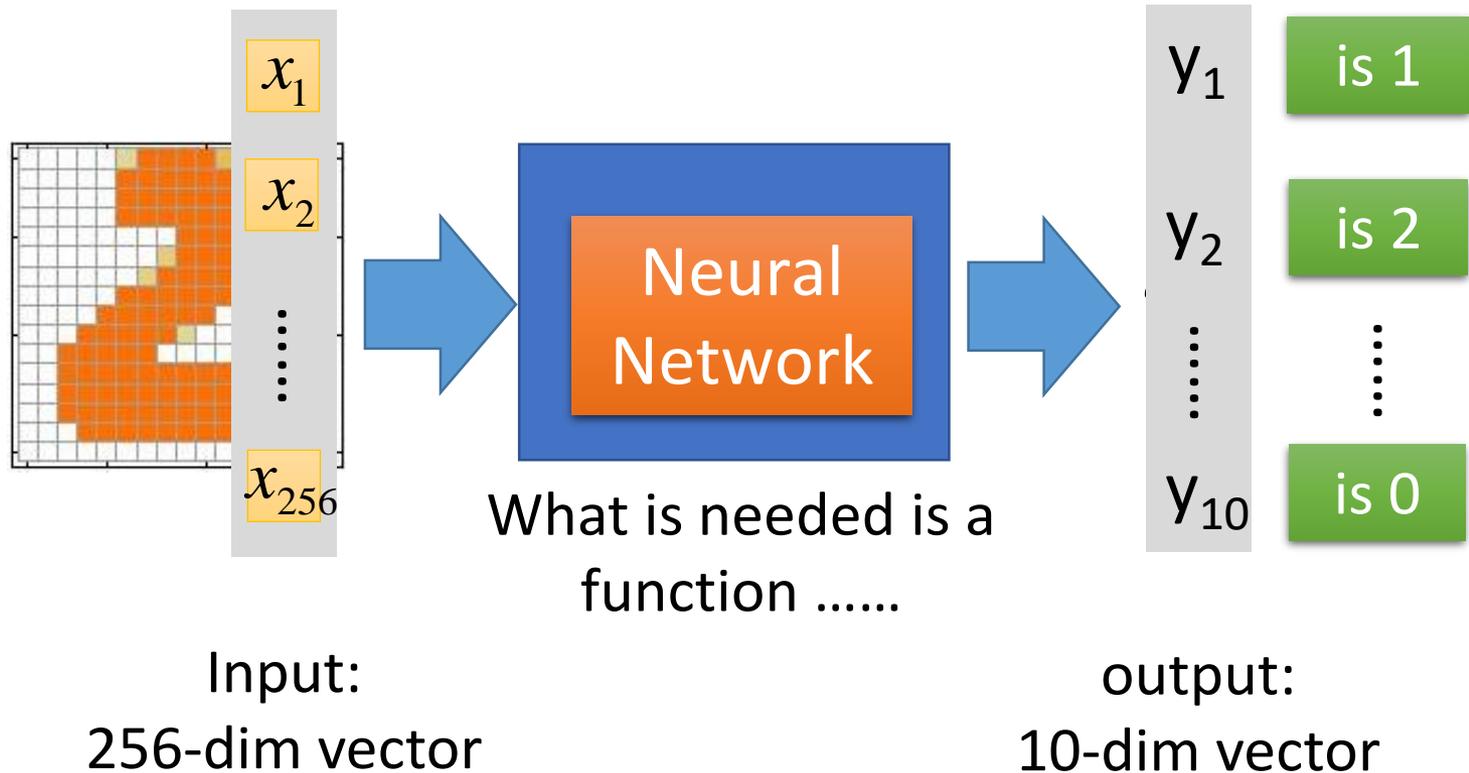


Each dimension represents the confidence of a digit.

Example Application

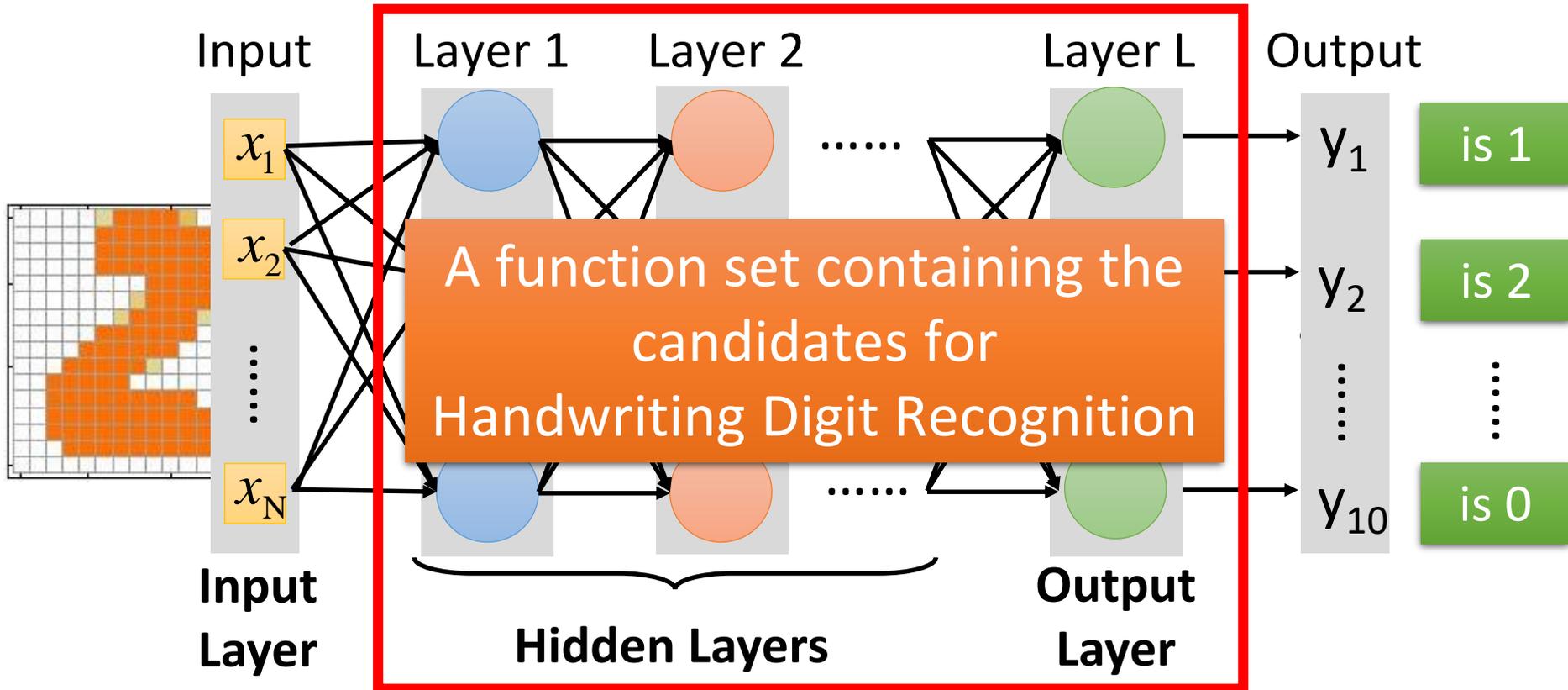
70

□ Handwriting Digit Recognition



Example Application

71



You need to decide the network structure to let a good function in your function set.

FAQ

72

- Q: How many layers? How many neurons for each layer?

Trial and Error

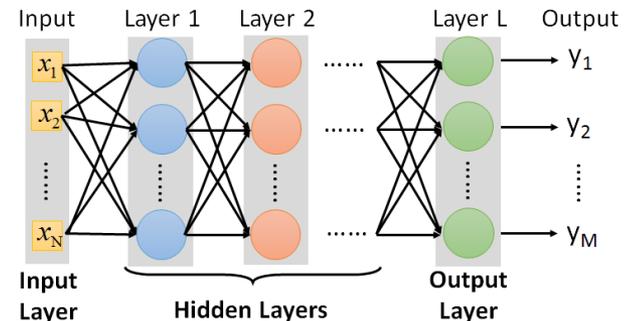
+

Intuition

- Q: Can we design the network structure?

Variants of Neural Networks
(next lecture)

- Q: Can the structure be automatically determined?
 - ▣ Yes, but not widely studied yet.



Three Steps for Deep Learning

73

Step 1: define a set of function



Step 2: goodness of function

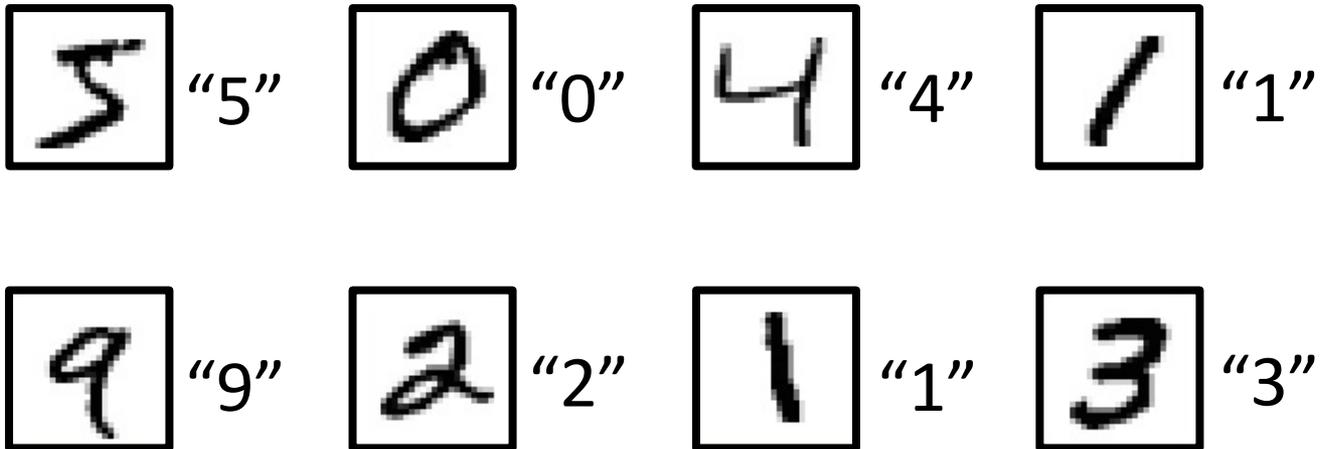


Step 3: pick the best function

Training Data

74

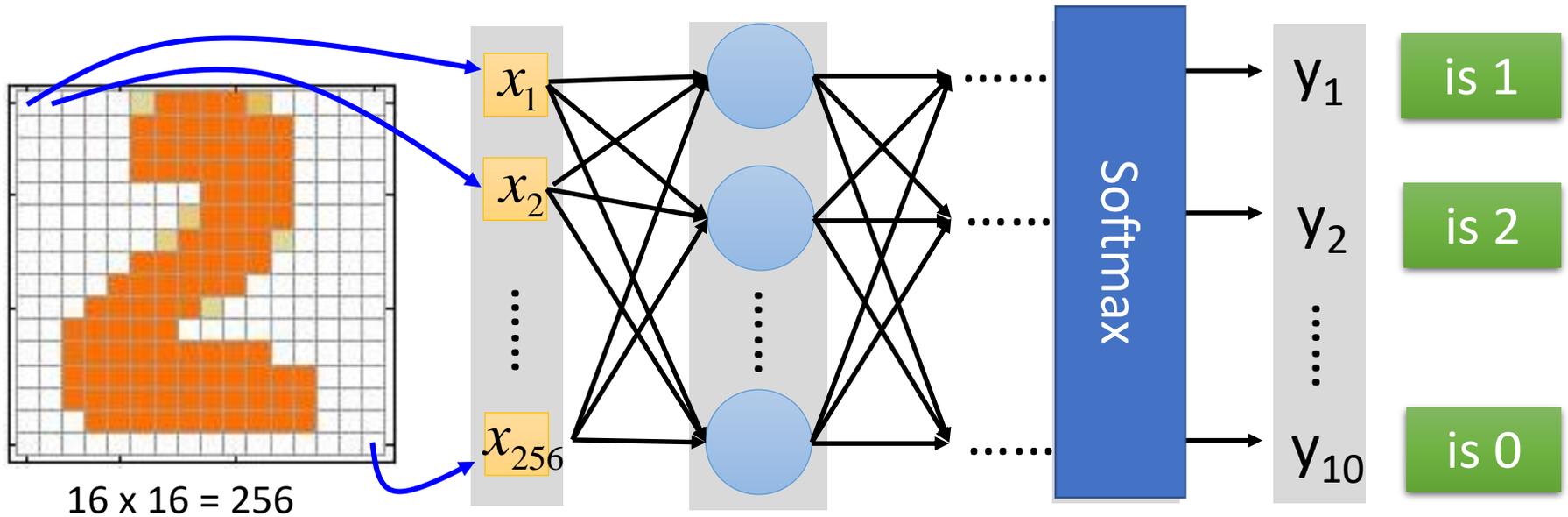
- Preparing training data: images and their labels



The learning target is defined on the training data.

Learning Target

75



Ink \rightarrow 1

No ink \rightarrow 0

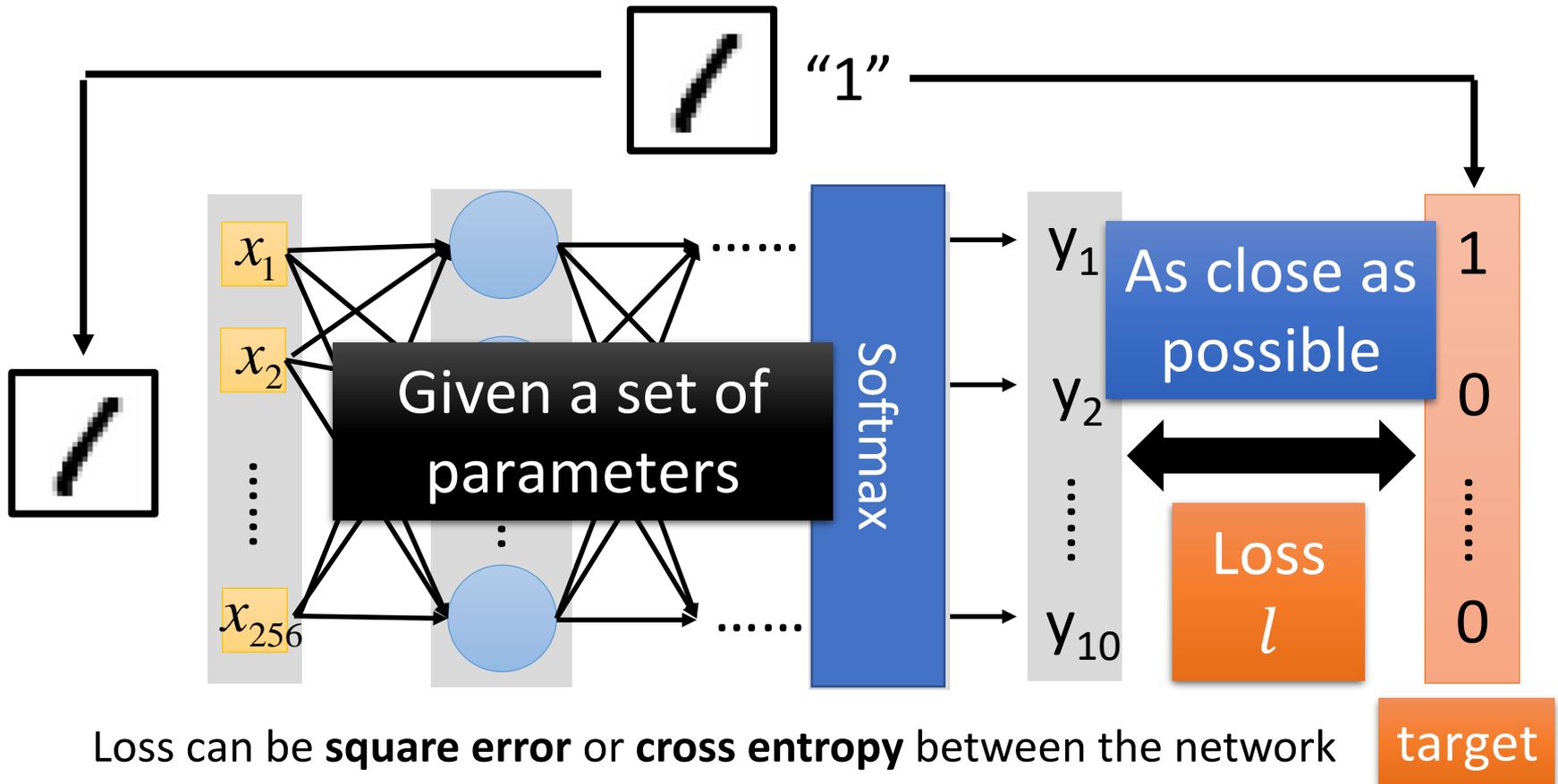
The learning target is

Input:  \rightarrow y_1 has the maximum value

Input:  \rightarrow y_2 has the maximum value

Loss

76



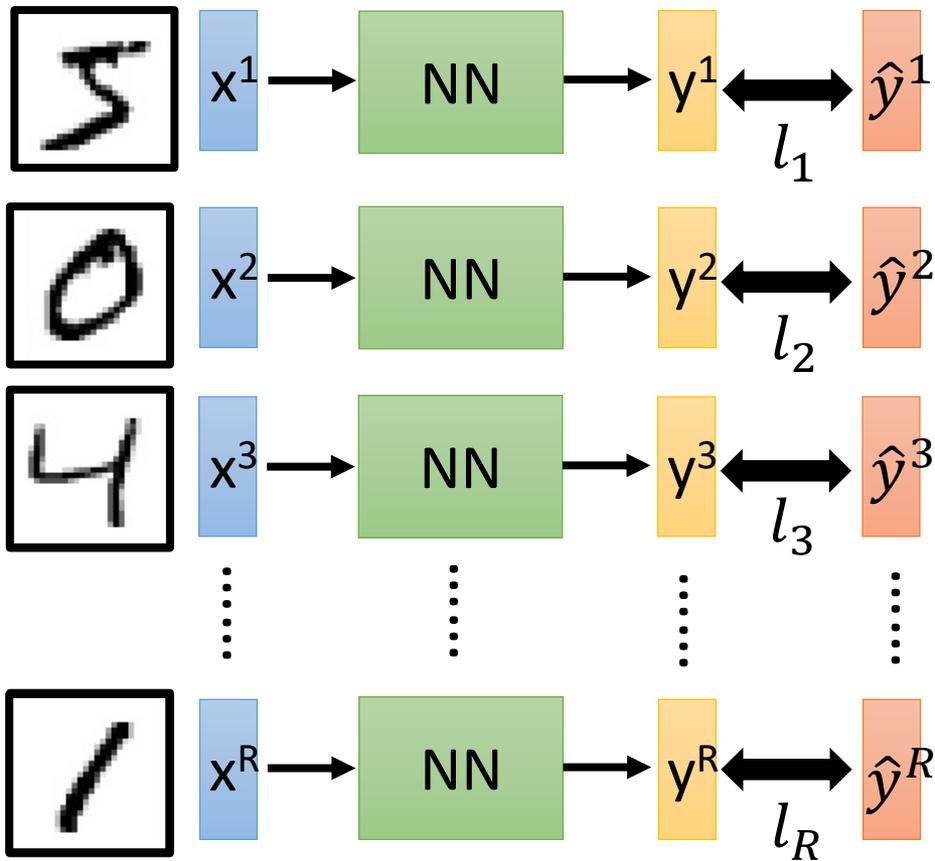
Loss can be **square error** or **cross entropy** between the network output and target

A good function should make the loss of all examples as small as possible.

Total Loss

77

- For all training data ...



Total Loss:

$$L = \sum_{r=1}^R l_r$$

As small as possible

Find a function in function set that minimizes total loss L

Find the network parameters θ^* that minimize total loss L

Three Steps for Deep Learning

78

Step 1: define a set of function



Step 2: goodness of function



Step 3: pick the best function

How to pick the best function

79

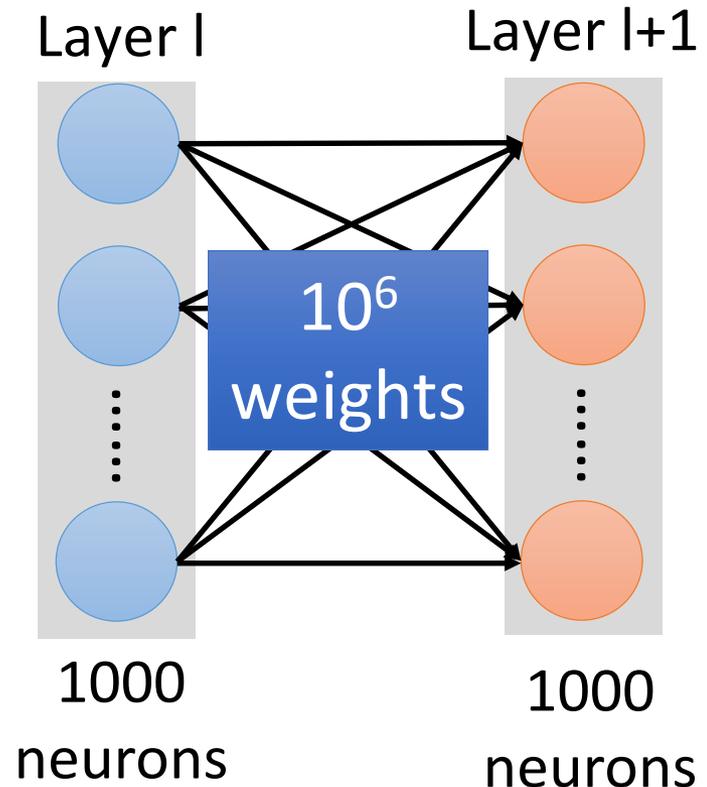
Find network parameters θ^* that minimize total loss L

Enumerate all possible values

Network parameters $\theta =$
 $\{w_1, w_2, w_3, \dots, b_1, b_2, b_3, \dots\}$

Millions of parameters

E.g. speech recognition: 8 layers and
1000 neurons each layer



Gradient Descent

Network parameters $\theta = \{w_1, w_2, \dots, b_1, b_2, \dots\}$

80

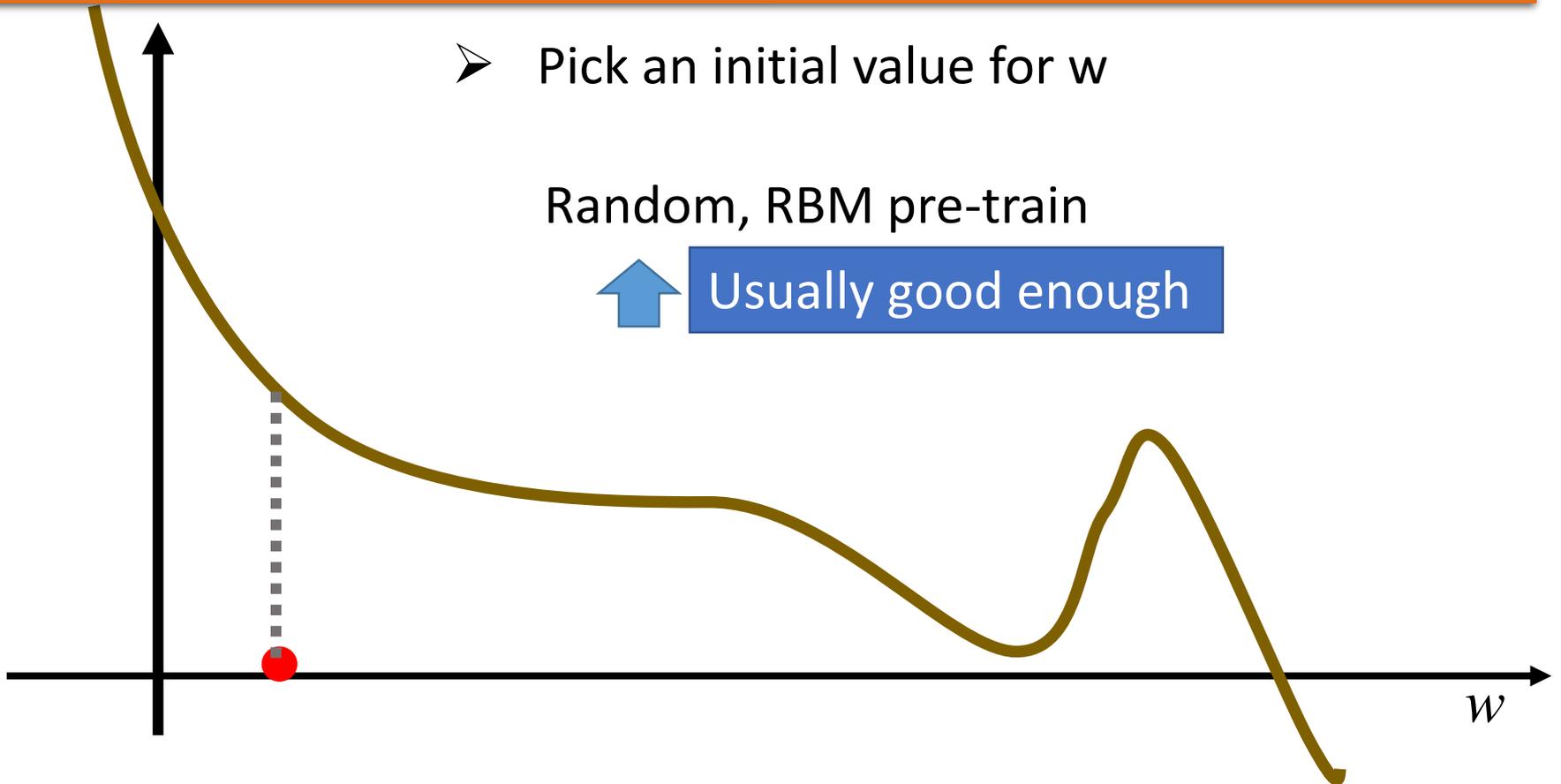
Find network parameters θ^* that minimize total loss L

- Pick an initial value for w

Random, RBM pre-train



Usually good enough



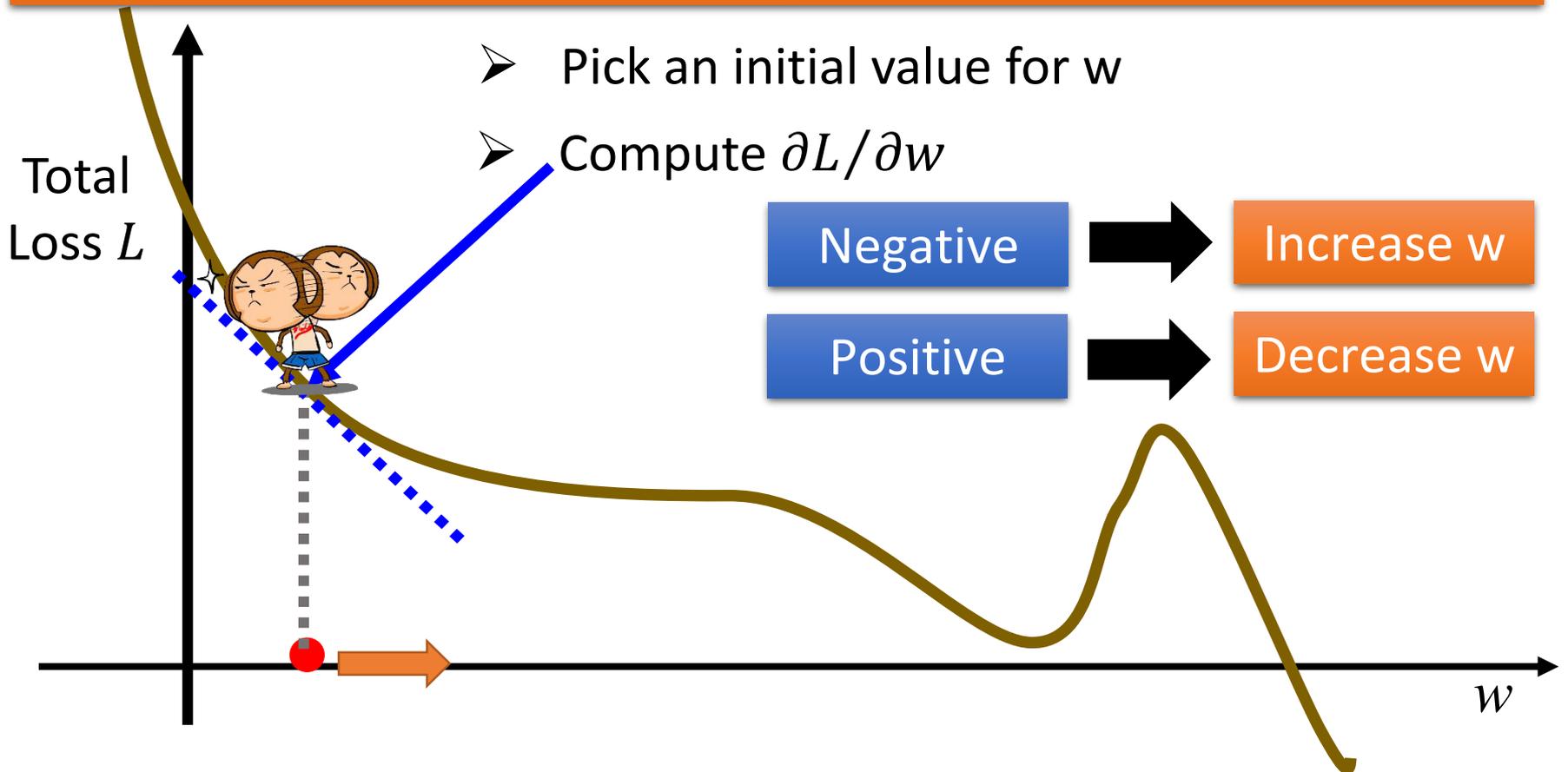
Gradient Descent

Network parameters $\theta = \{w_1, w_2, \dots, b_1, b_2, \dots\}$

81

Find network parameters θ^* that minimize total loss L

- Pick an initial value for w
- Compute $\partial L / \partial w$



Gradient Descent

Network parameters $\theta = \{w_1, w_2, \dots, b_1, b_2, \dots\}$

82

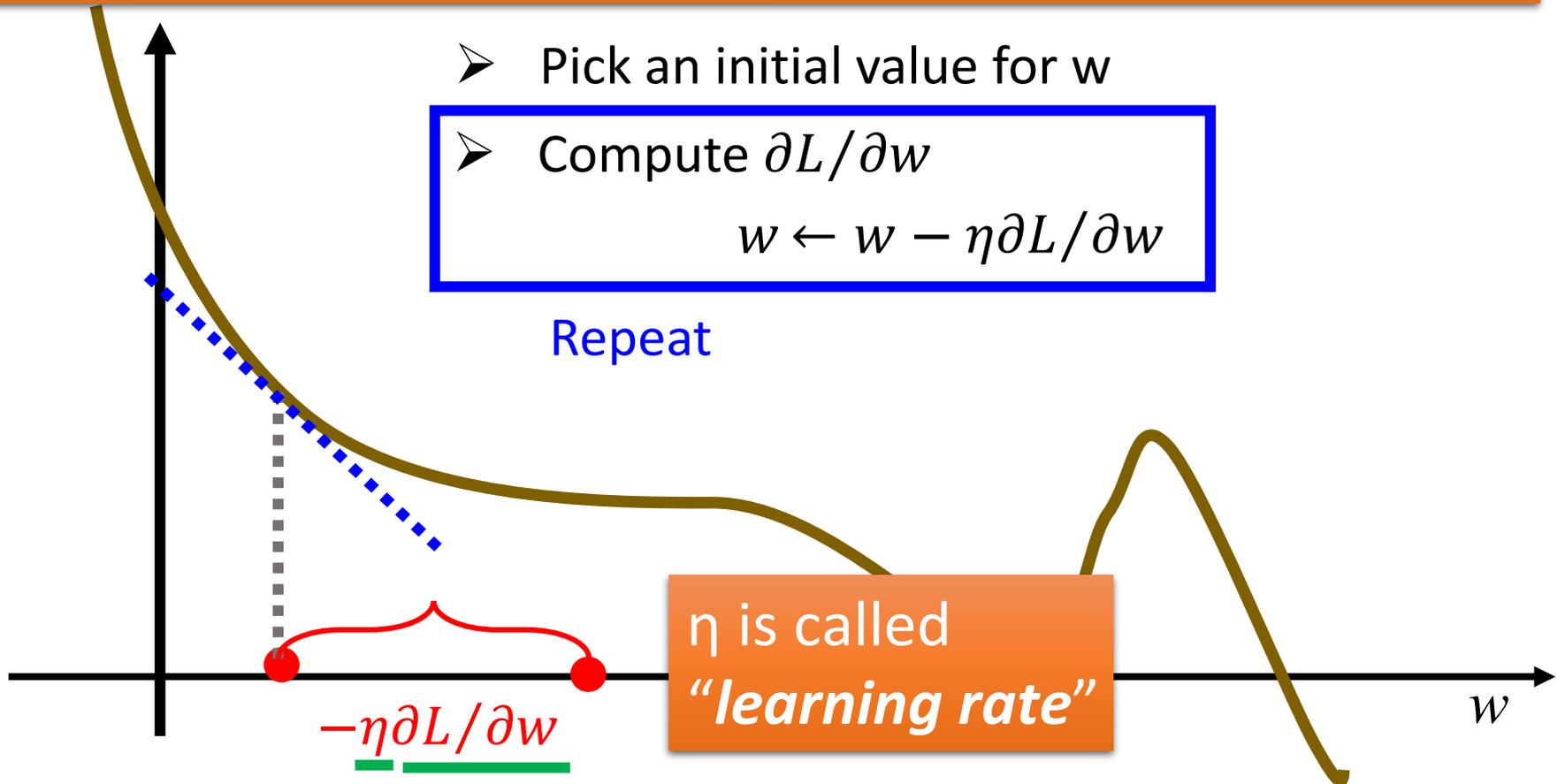
Find network parameters θ^* that minimize total loss L

➤ Pick an initial value for w

➤ Compute $\partial L / \partial w$

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

Repeat



Gradient Descent

Network parameters $\theta = \{w_1, w_2, \dots, b_1, b_2, \dots\}$

83

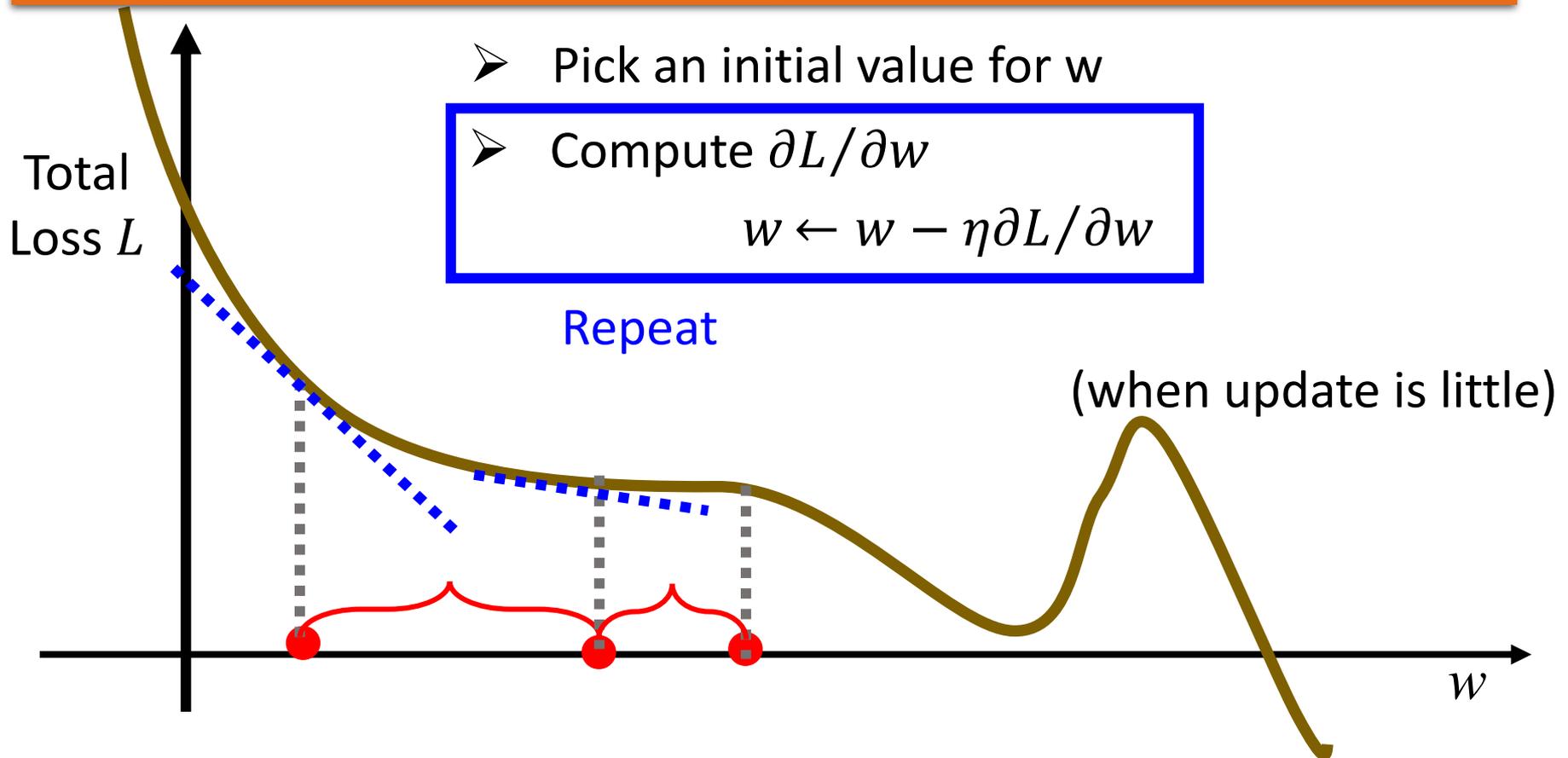
Find network parameters θ^* that minimize total loss L

➤ Pick an initial value for w

➤ Compute $\partial L / \partial w$

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

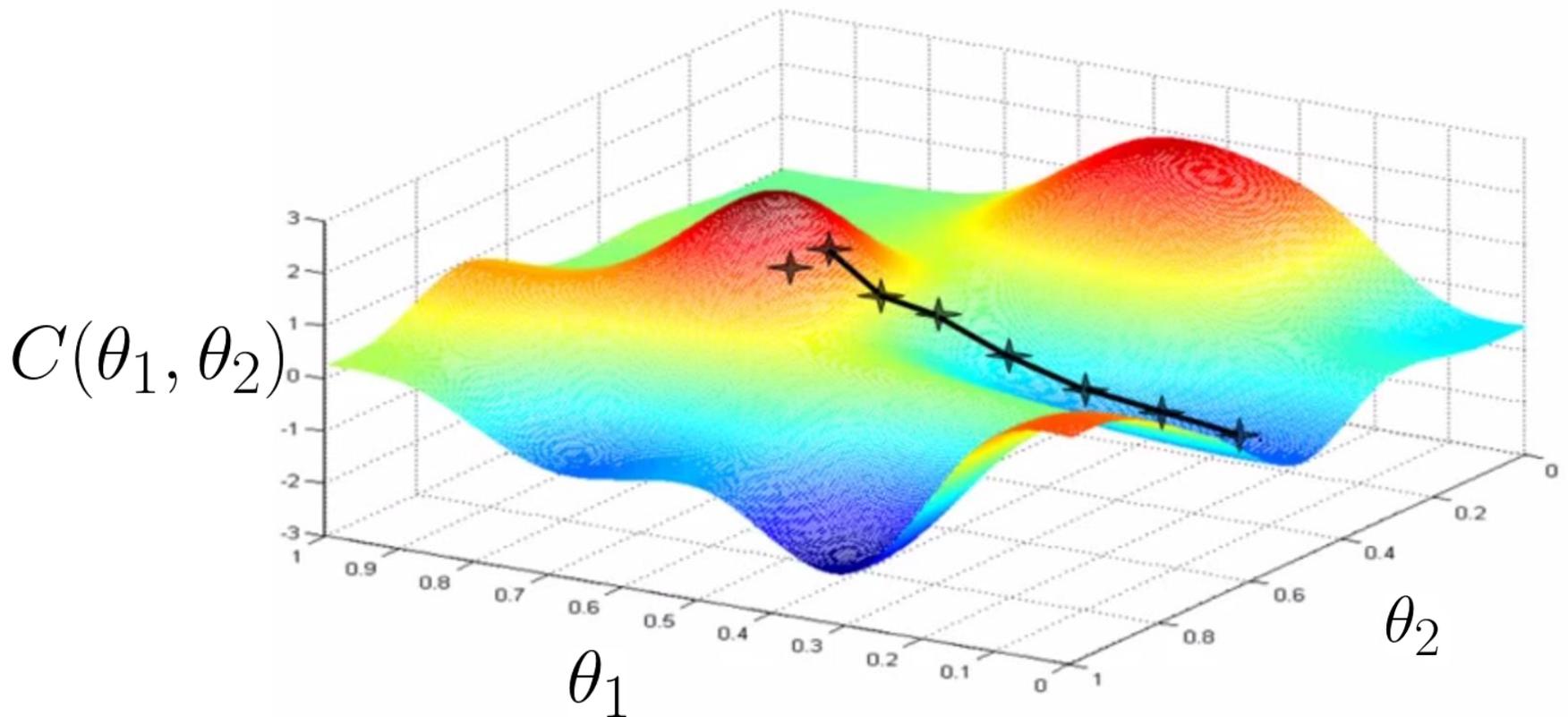
Repeat



Gradient Descent

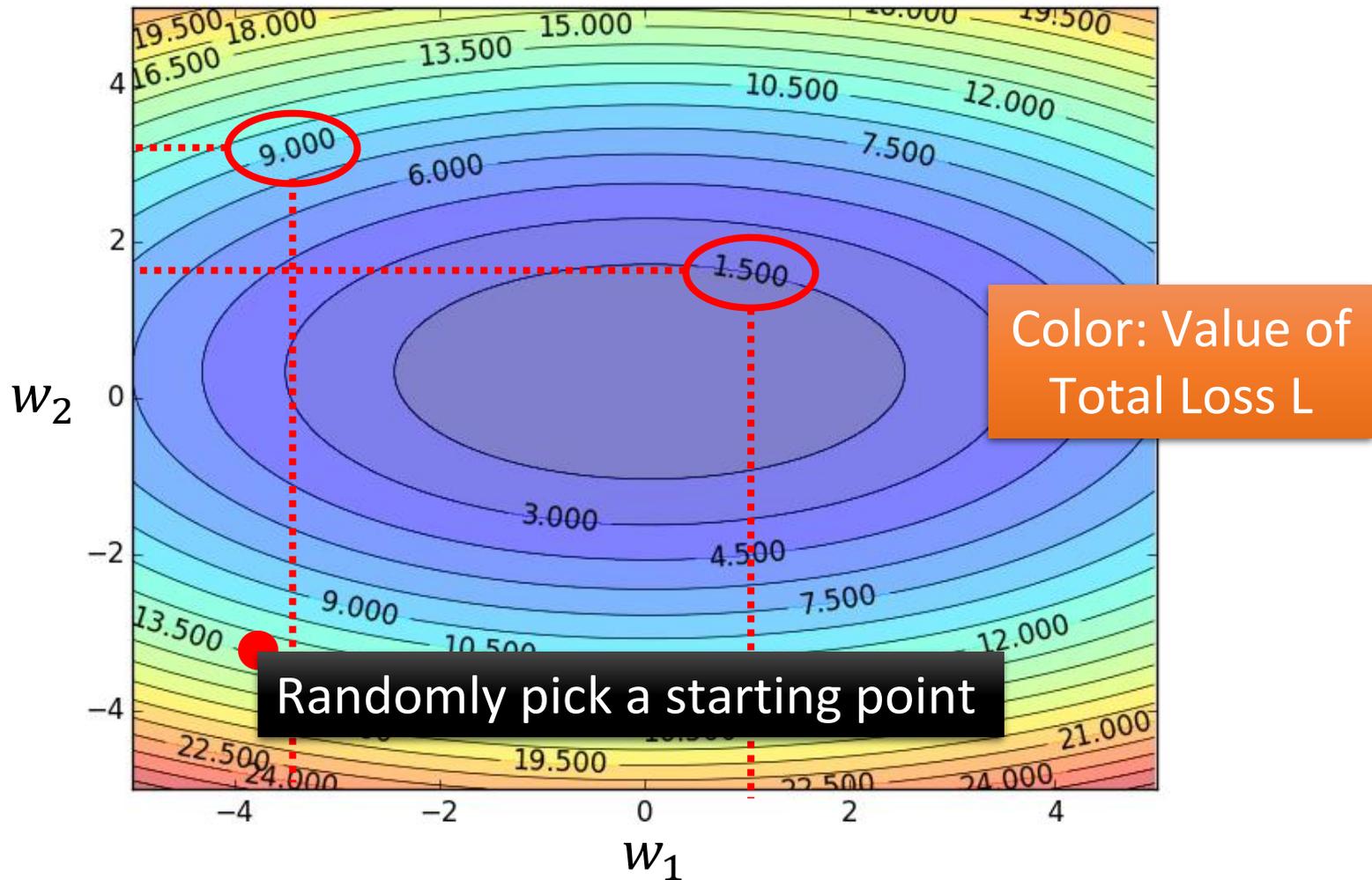
84

- Assume that θ has two variables $\{\theta_1, \theta_2\}$



Gradient Descent

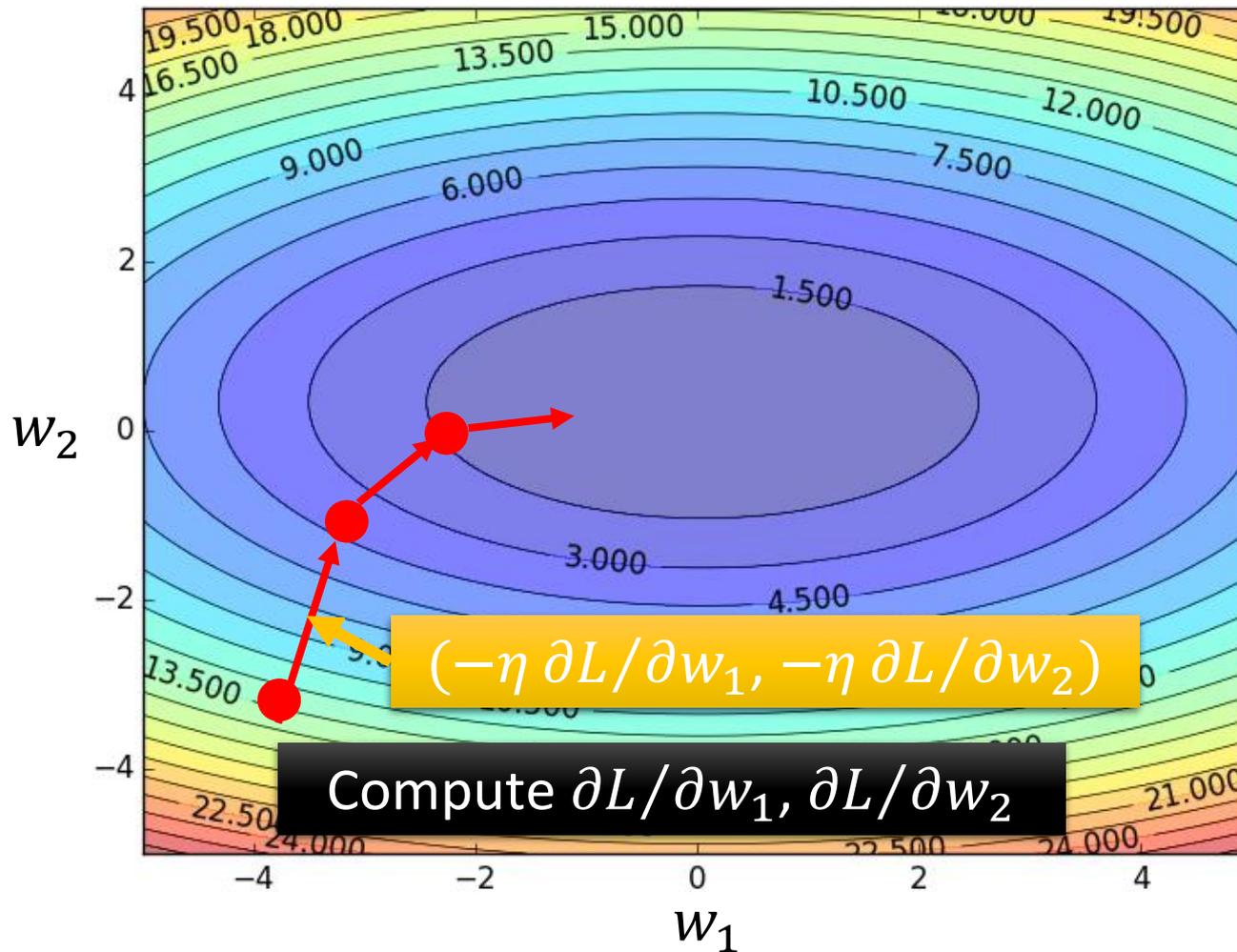
85



Gradient Descent

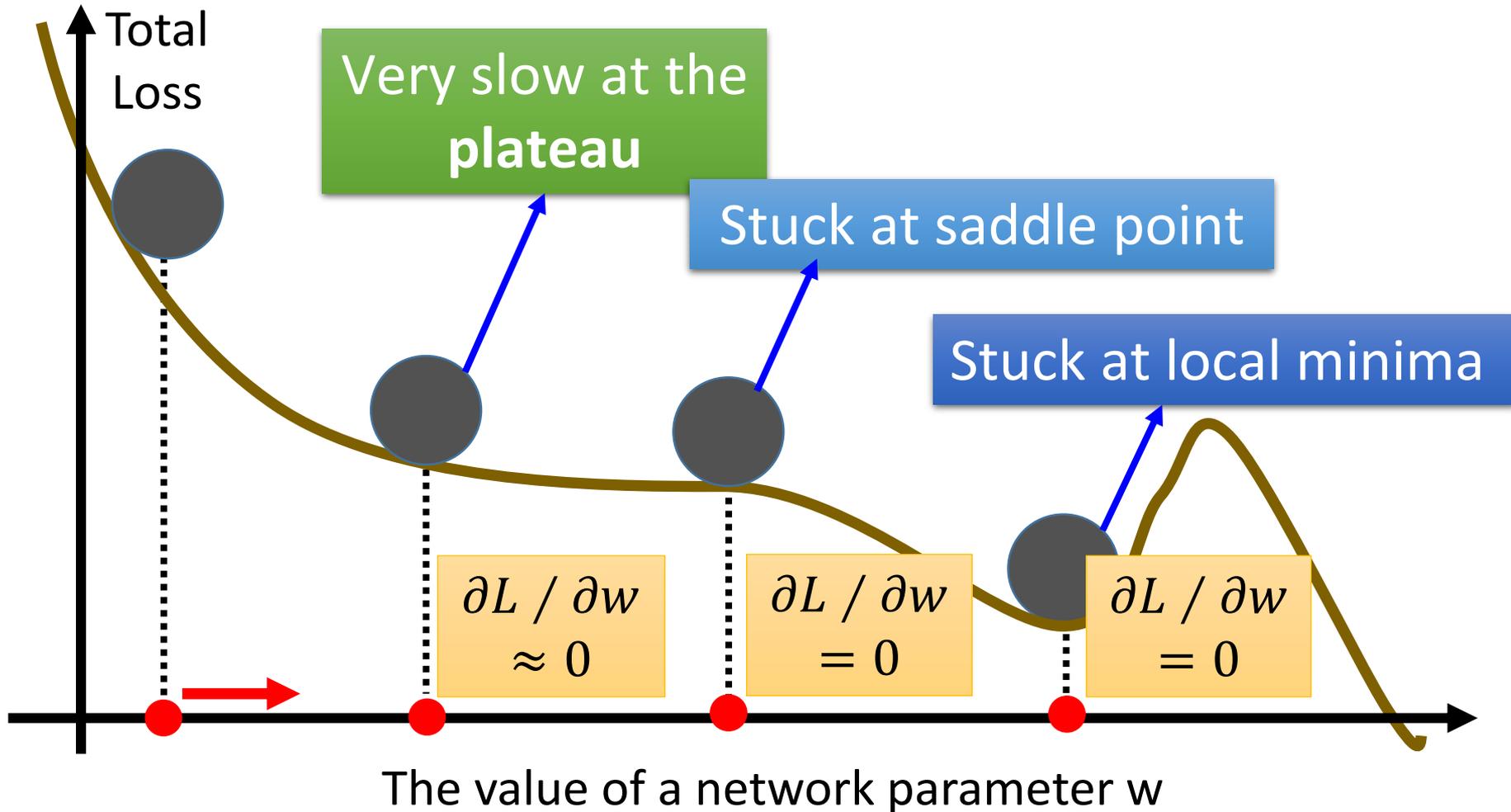
Hopfully, we would reach
a minima

86



Local Minima

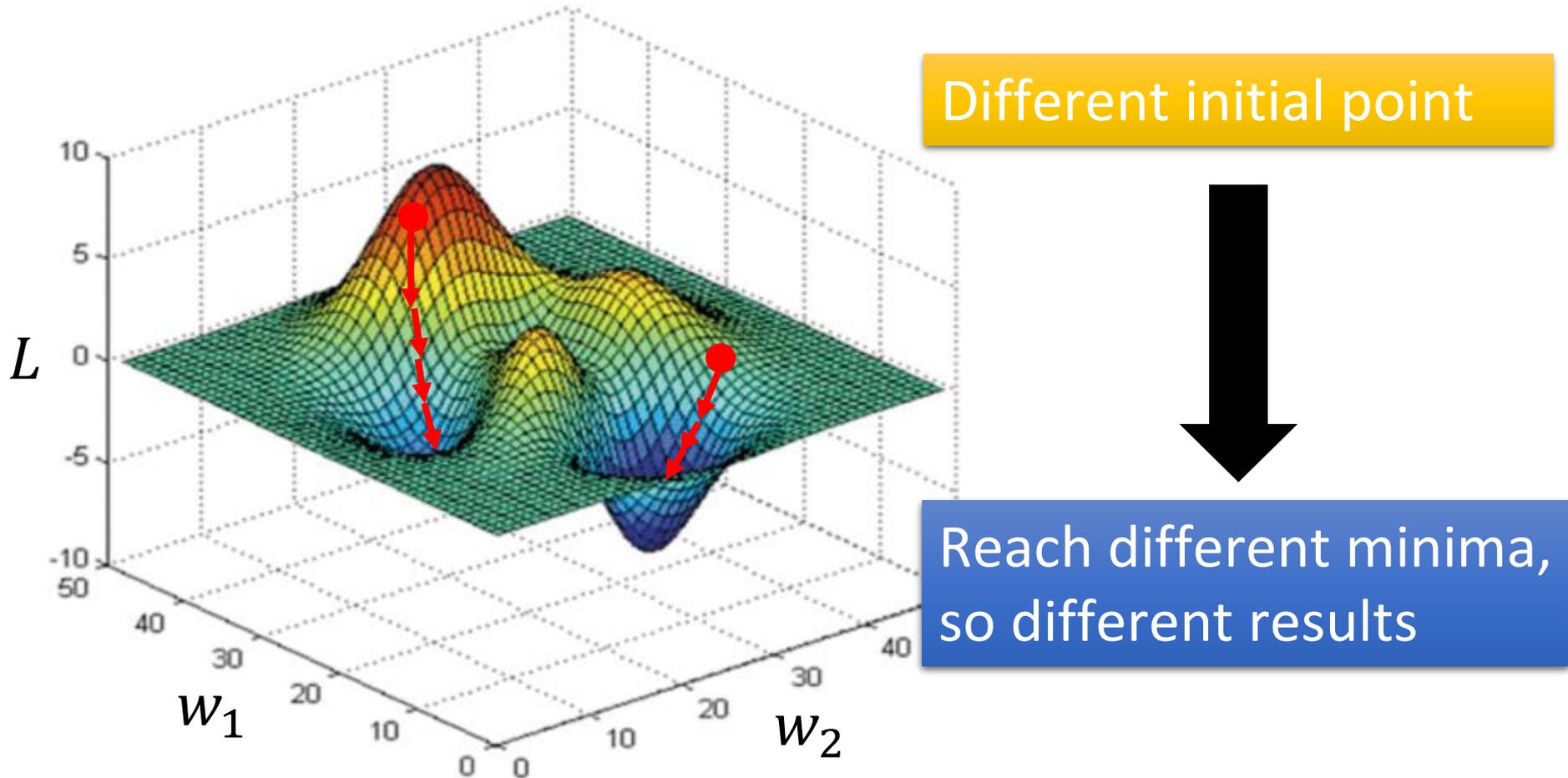
87



Local Minima

88

- Gradient descent never guarantee global



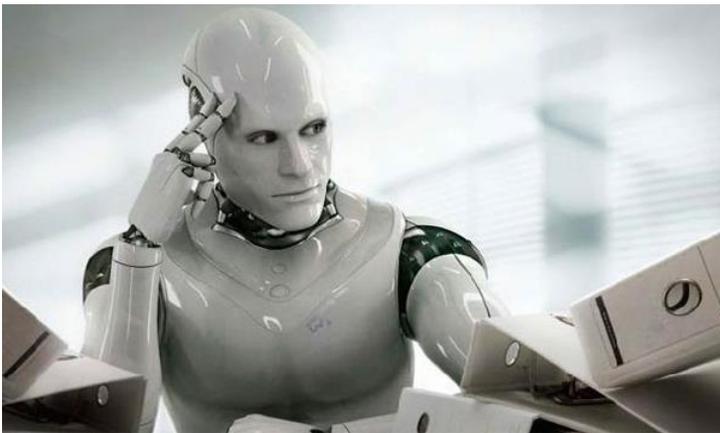
Gradient Descent

89

This is the “learning” of machines in deep learning

➔ Even AlphaGo using this approach.

People image



Actually



I hope you are not too disappointed :p

Part I: Introduction to ML & DL

90

- Basic Machine Learning
- Basic Deep Learning
- **Toolkits and Learning Recipe**

Deep Learning Toolkit

91

- Backpropagation: an efficient way to compute $\partial L / \partial w$ in neural network



Caffe



theano



Deep Learning library produced by Amazon

DSSTNE

Three Steps for Deep Learning

92



Deep Learning is so simple

Now If you want to find a function

If you have lots of function input/output (?) as training data

 You can use deep learning

Keras

93



or **theano**

Very flexible
Need some
effort to learn



Interface of
TensorFlow or
Theano



keras

Easy to learn and use
(still have some flexibility)
You can modify it if you can write
TensorFlow or Theano

Keras

94

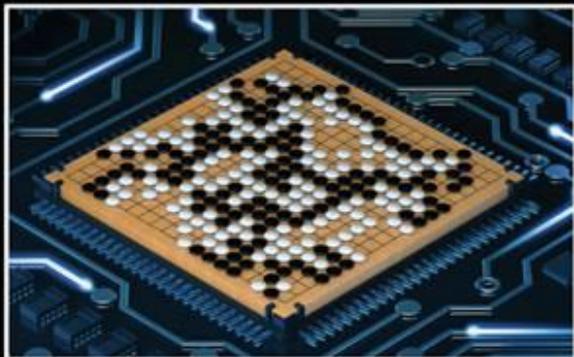
- François Chollet is the author of Keras.
 - ▣ He currently works for Google as a deep learning engineer and researcher.
- Keras means *horn* in Greek
- Documentation: <http://keras.io/>
- Example
 - ▣ <https://github.com/fchollet/keras/tree/master/examples>
- Step-by-step lecture by Prof. Hung-Yi Lee
 - ▣ Slide
http://speech.ee.ntu.edu.tw/~tlkagk/courses/ML_2016/Lecture/Keras.pdf
 - ▣ Lecture recording:
<https://www.youtube.com/watch?v=qetE6uUoLQA>

使用 Keras 心得

95

Deep Learning 研究生

感謝 沈昇勳 同學提供圖檔



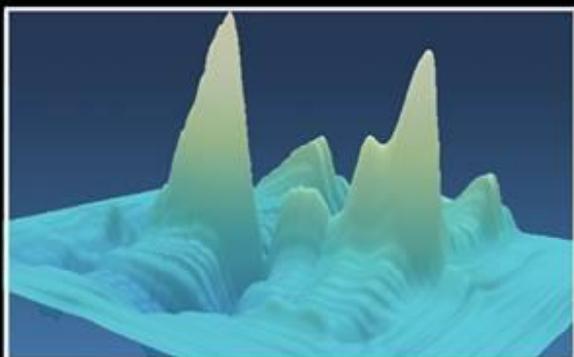
朋友覺得我在



我媽覺得我在



大眾覺得我在



指導教授覺得我在



我以為我在

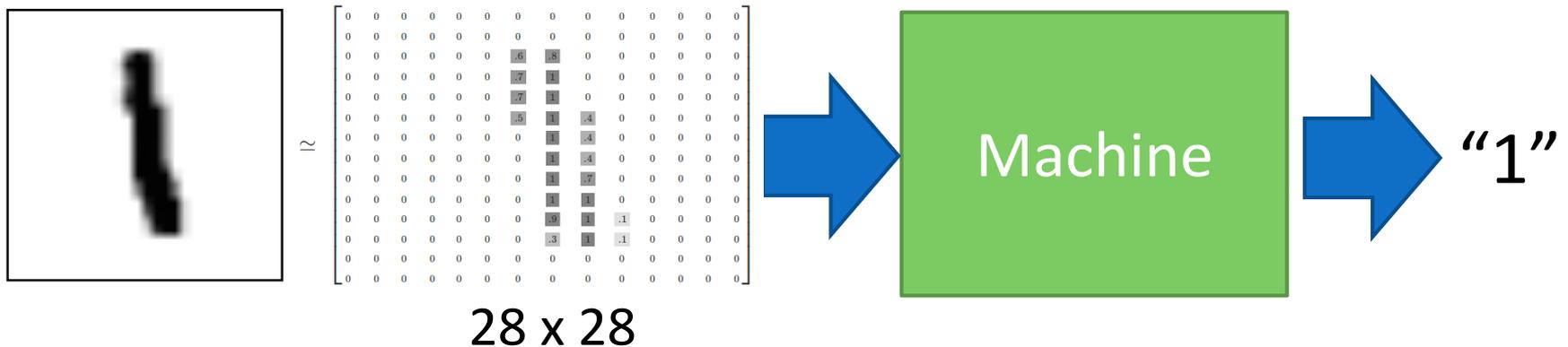


事實上我在

Example Application

96

□ Handwriting Digit Recognition



MNIST Data: <http://yann.lecun.com/exdb/mnist/>

“Hello world” for deep learning

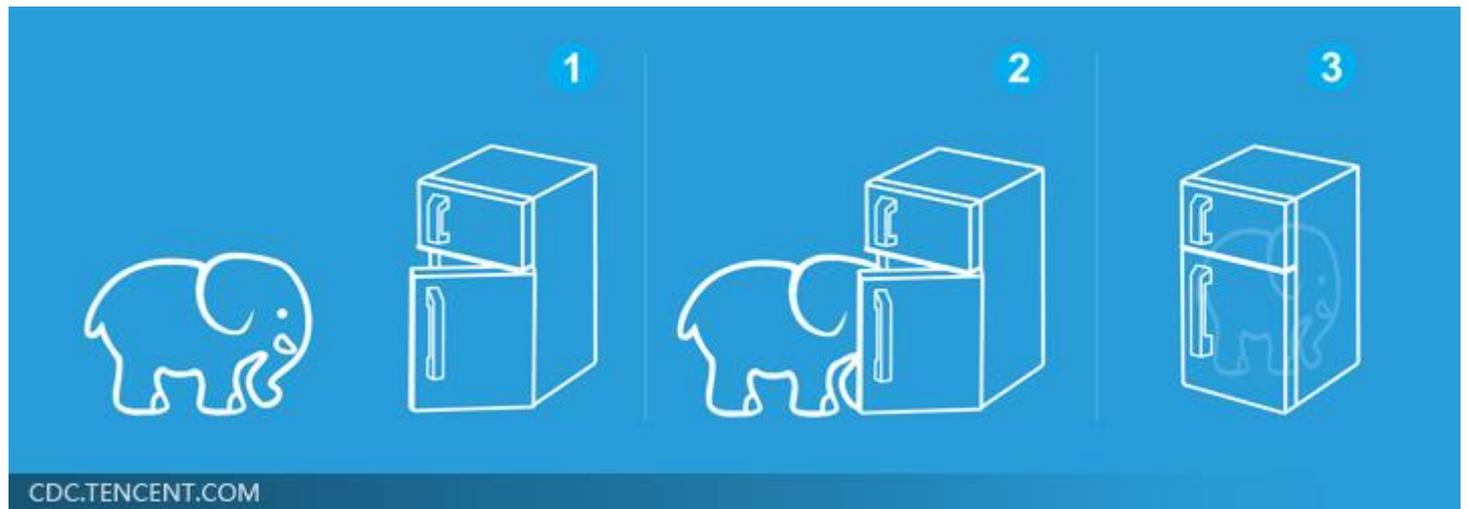
Keras provides data sets loading function: <http://keras.io/datasets/>

Three Steps for Deep Learning

97



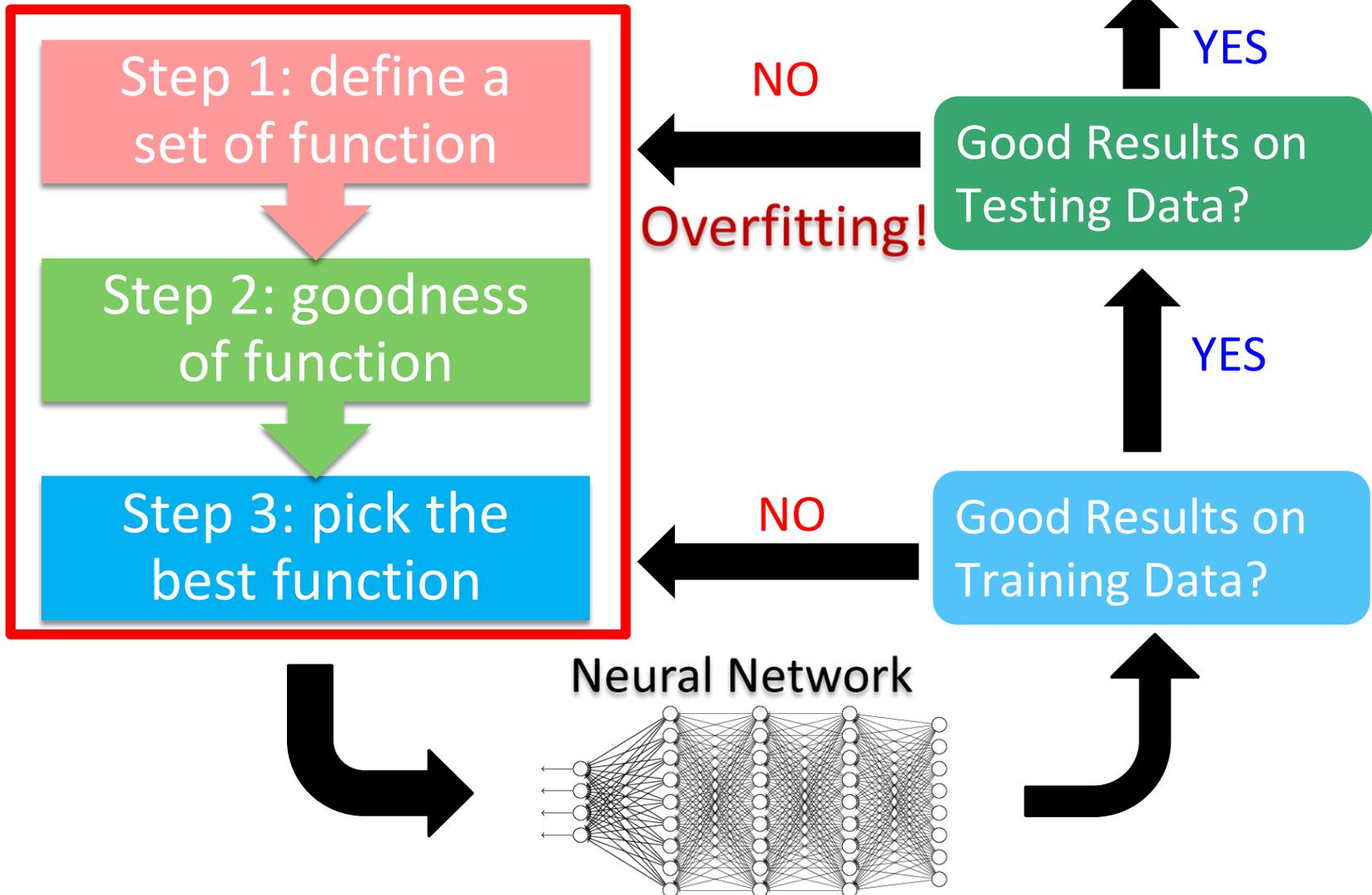
Deep Learning is so simple



Learning Recipe



98

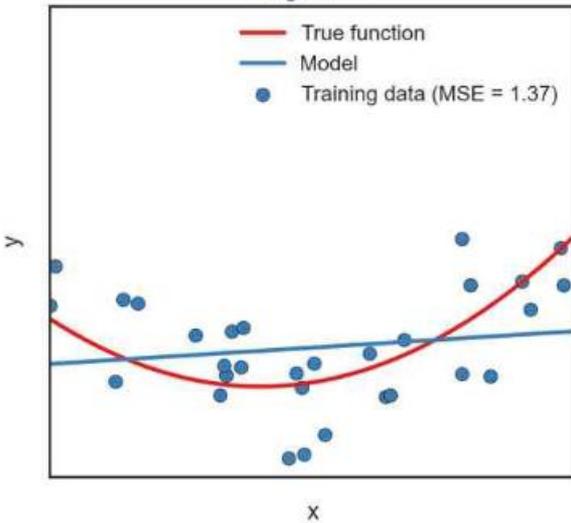


Overfitting

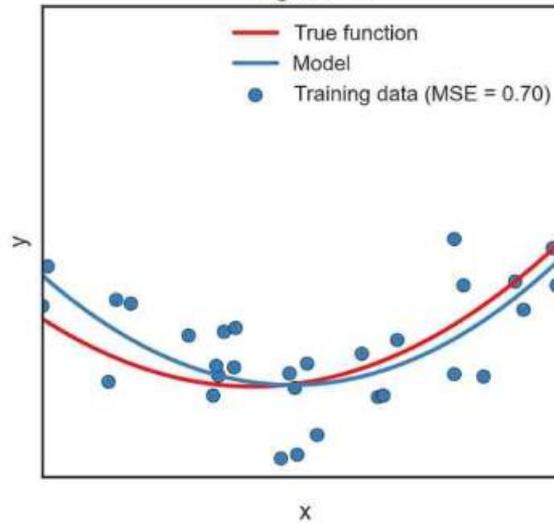
99

Fitting training data

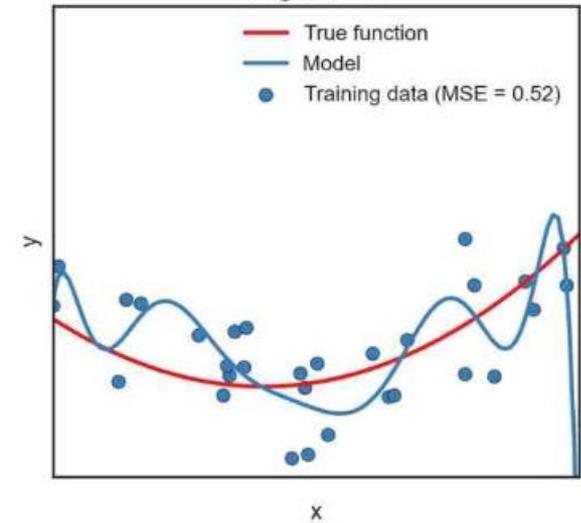
Degree = 1



Degree = 2



Degree = 10

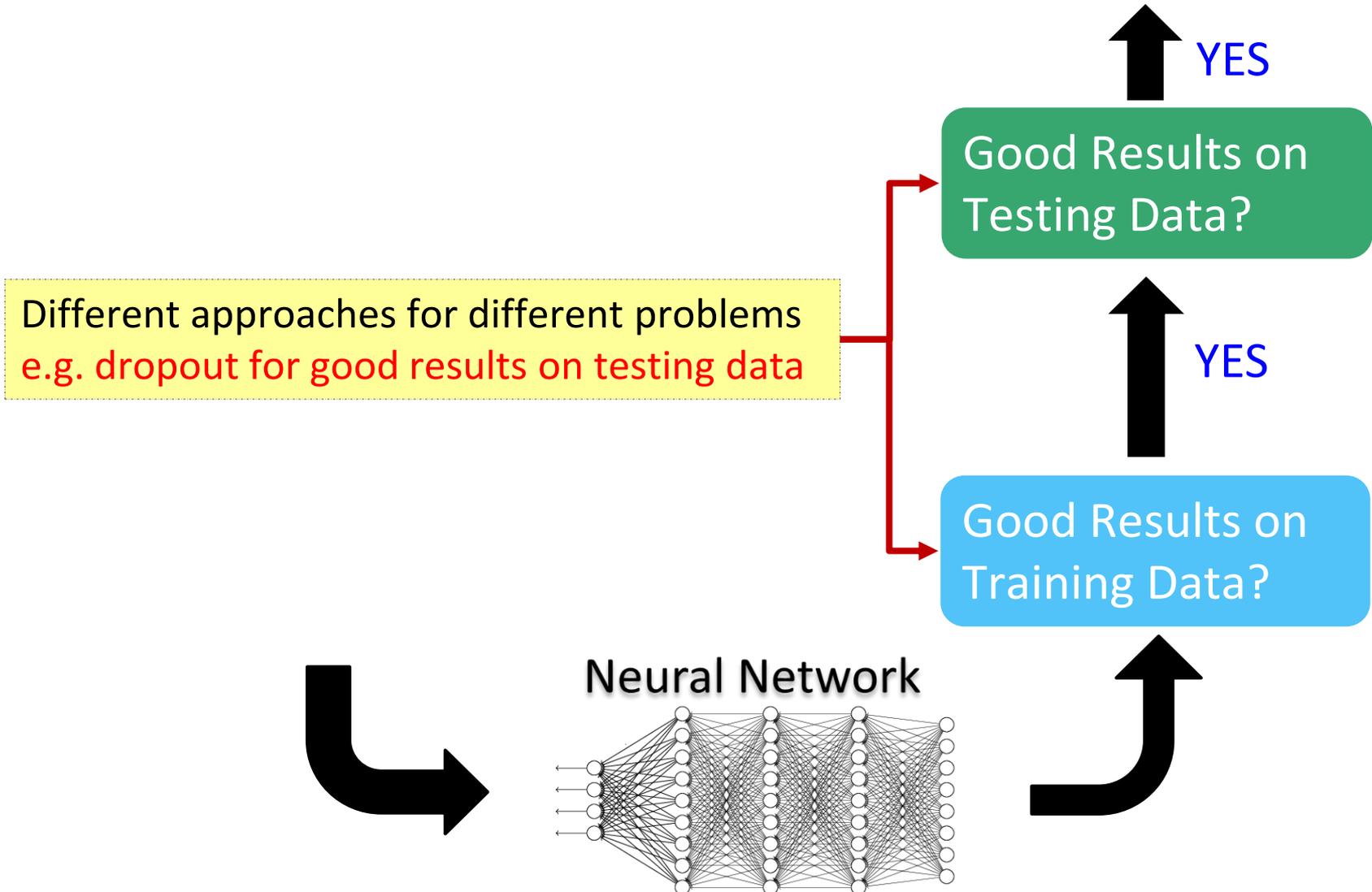


- Possible solutions
 - ▣ more training samples
 - ▣ some tips: dropout, etc.

Learning Recipe



100



Learning Recipe

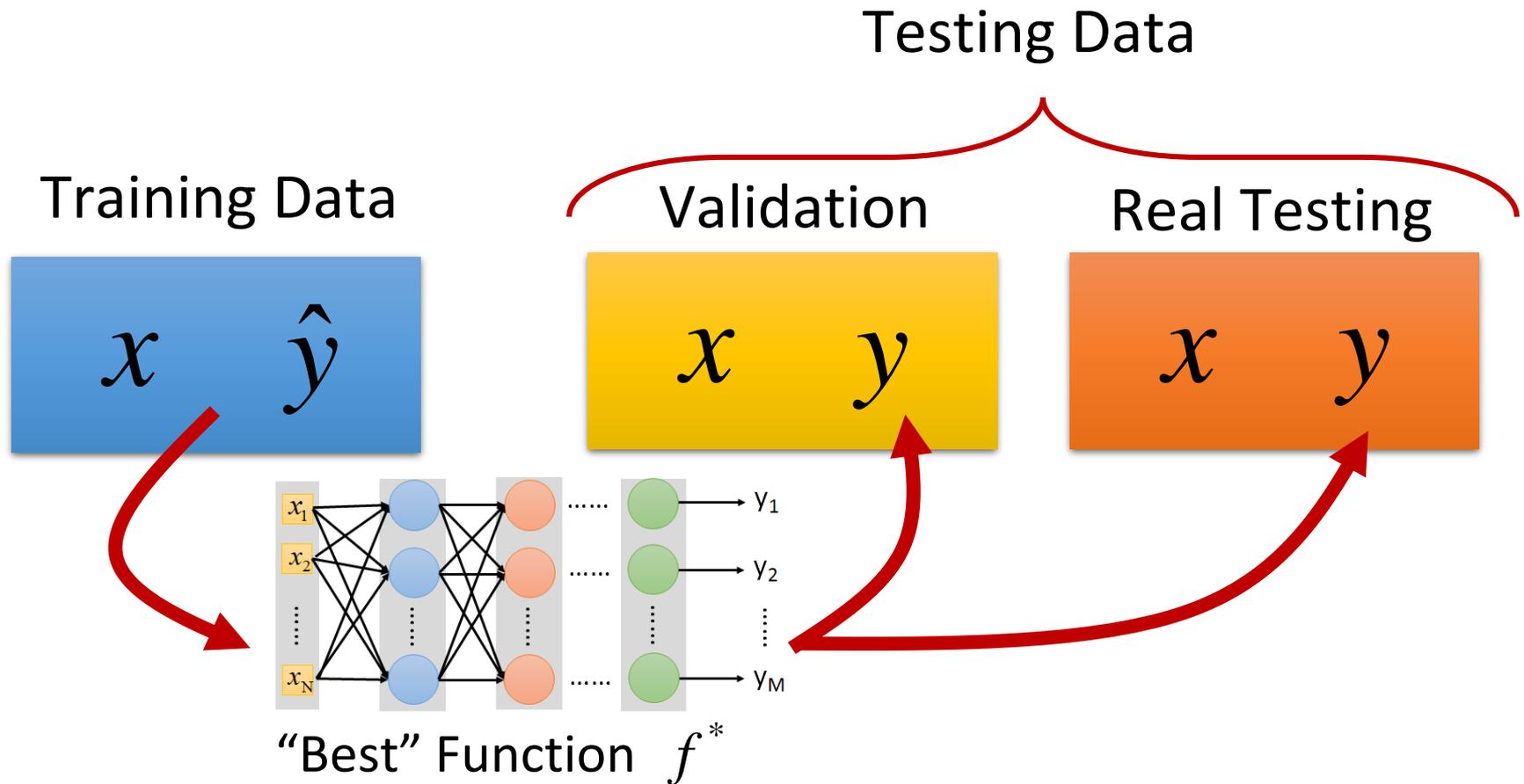


101



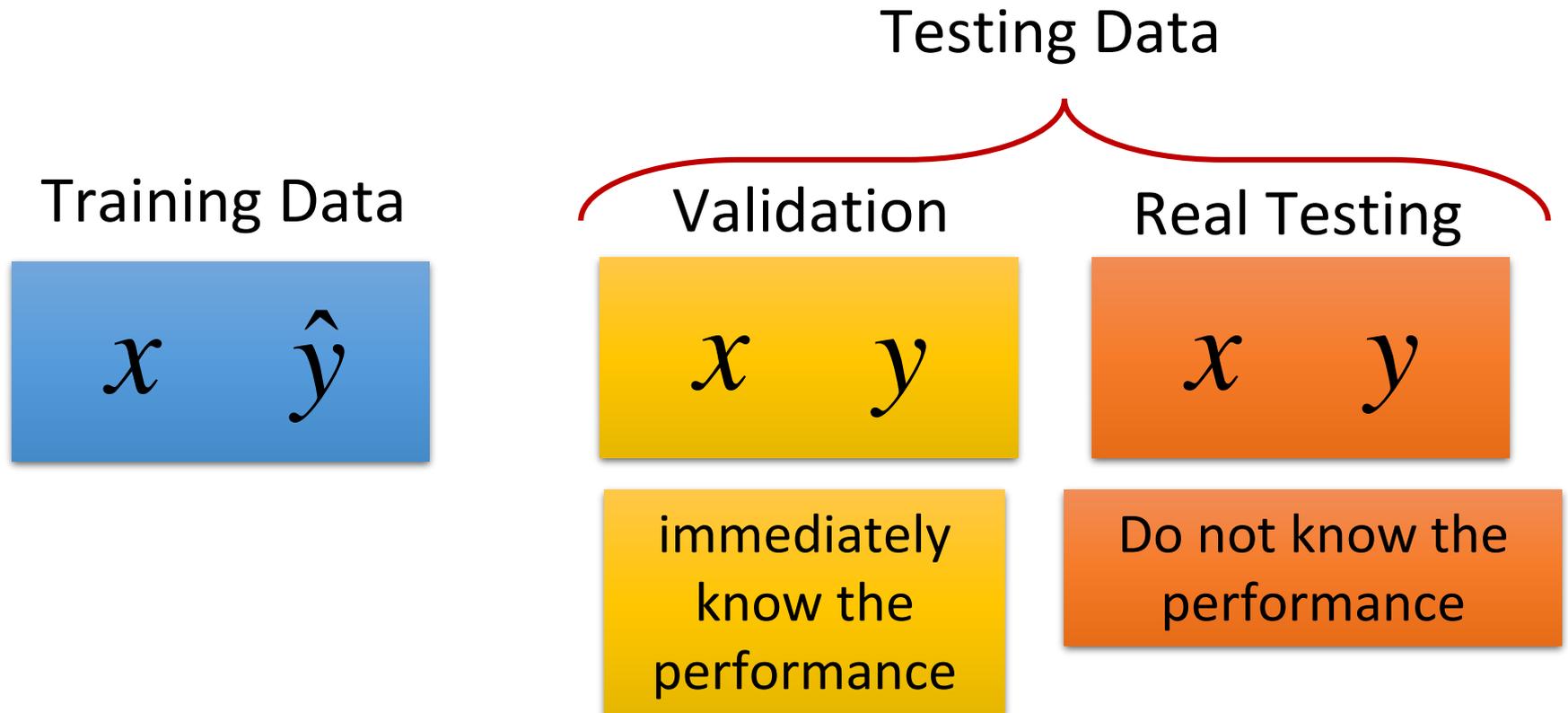
Learning Recipe

102



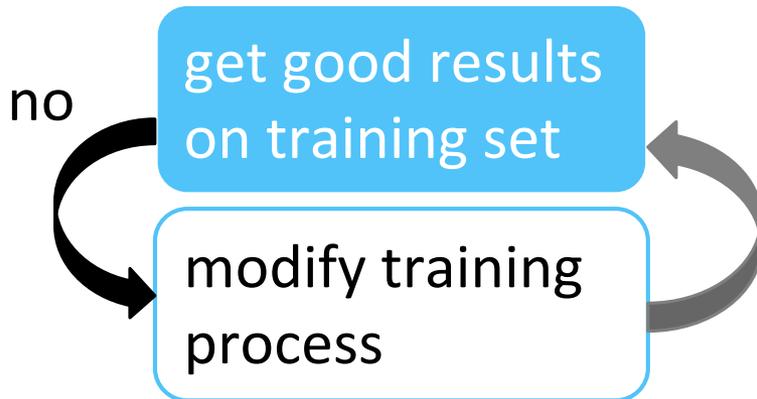
Learning Recipe

103



Learning Recipe

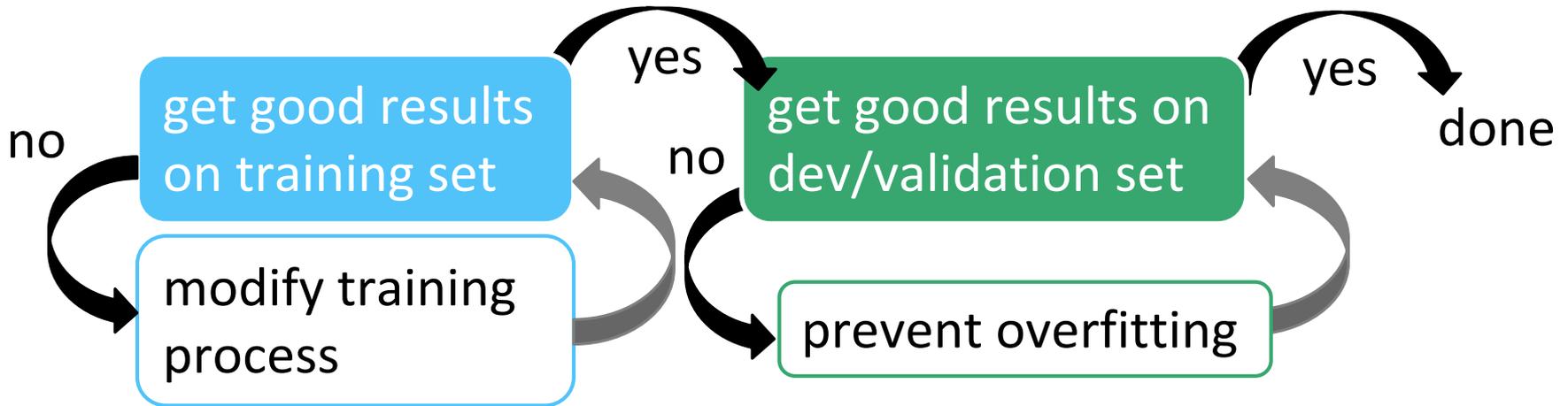
104



- Possible reasons
 - ▣ no good function exists: bad hypothesis function set
→ reconstruct the model architecture
 - ▣ cannot find a good function: local optima
→ change the training strategy

Learning Recipe

105

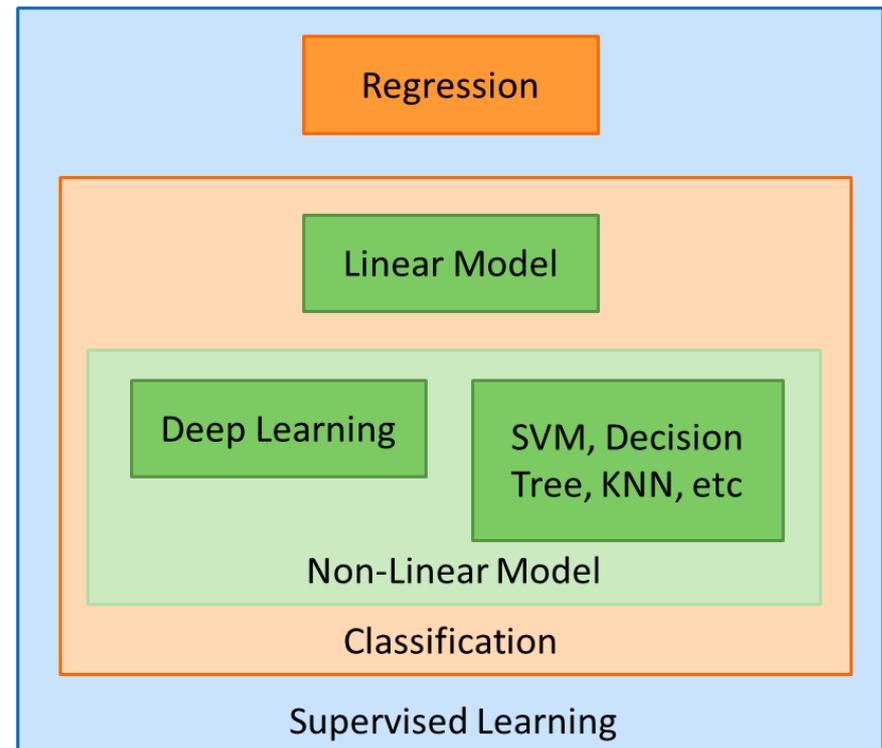


Better performance on training but worse performance on dev → overfitting

Concluding Remarks

106

- Basic Machine Learning
 1. Define a set of functions
 2. Measure goodness of functions
 3. Pick the best function
- Basic Deep Learning
 - ▣ Stacked functions



Talk Outline

107

Part I: Introduction to
Machine Learning & Deep Learning



Part II: Variants of Neural Nets



Part III: Beyond Supervised Learning
& Recent Trends

PART II

Variants of Neural Networks

PART II: Variants of Neural Networks

109

- Convolutional Neural Network (CNN)
- Recurrent Neural Network (RNN)

PART II: Variants of Neural Networks

110

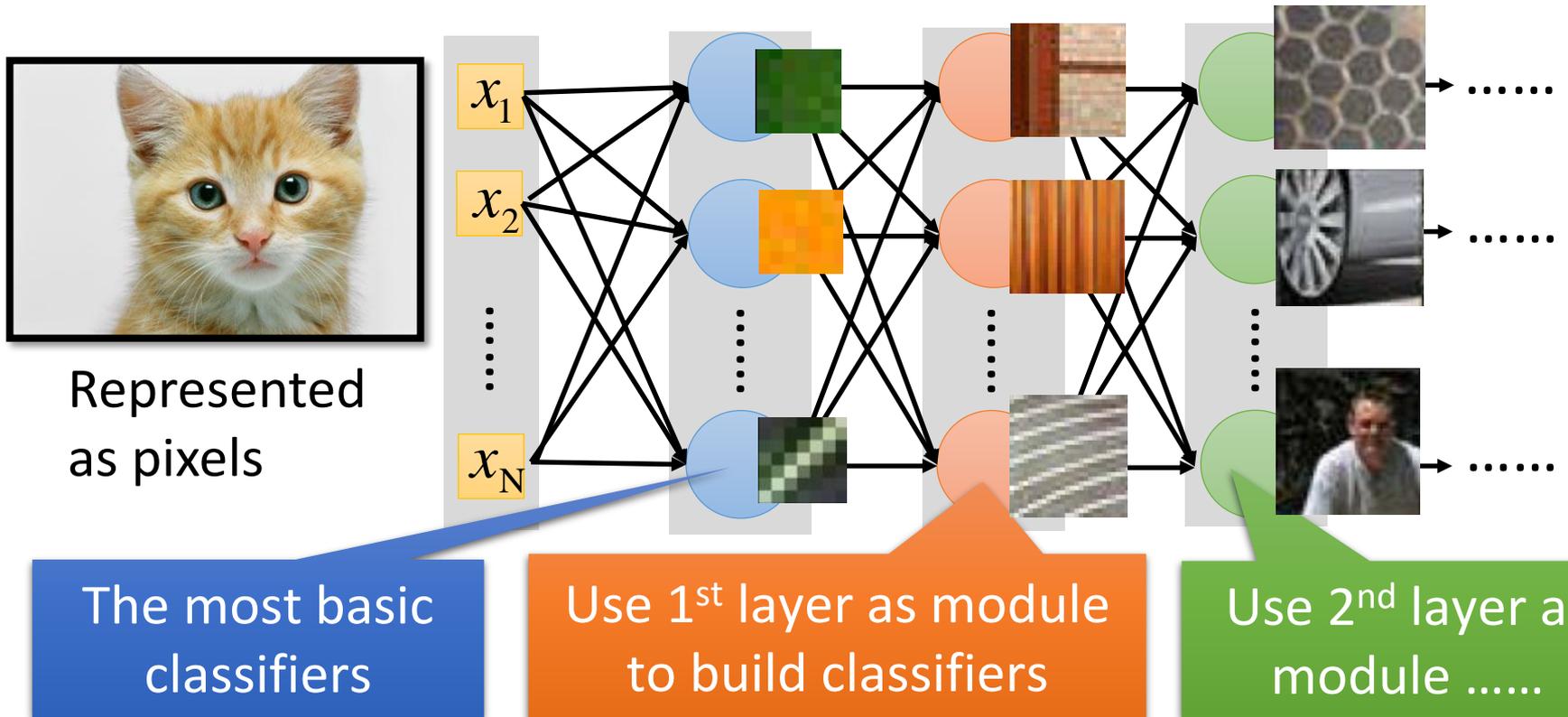
- **Convolutional Neural Network (CNN)**
- Recurrent Neural Network (RNN)

Widely used in image processing

Why CNN for Image

(Zeiler, M. D., *ECCV 2014*)

111



Can the network be simplified by considering the properties of images?

Why CNN for Image

112

- Some patterns are much smaller than the whole image

A neuron does not have to see the whole image to discover the pattern.

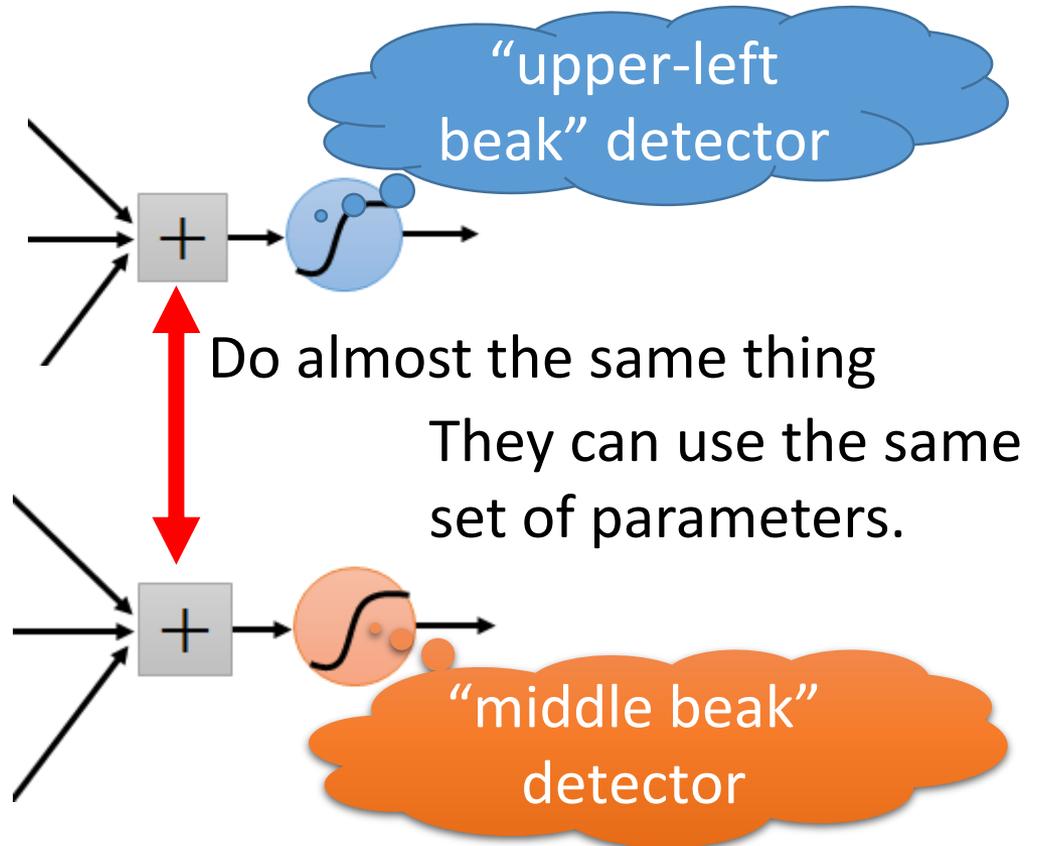
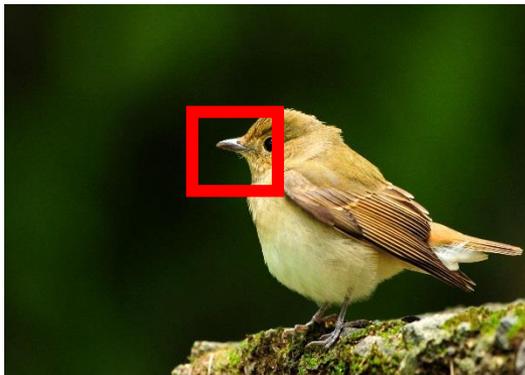
Connecting to small region with less parameters



Why CNN for Image

113

- The same patterns appear in different regions.



Why CNN for Image

114

- Subsampling the pixels will not change the object

bird



subsampling

bird



We can subsample the pixels to make image smaller



Less parameters for the network to process the image

Three Steps for Deep Learning

115



Deep Learning is so simple

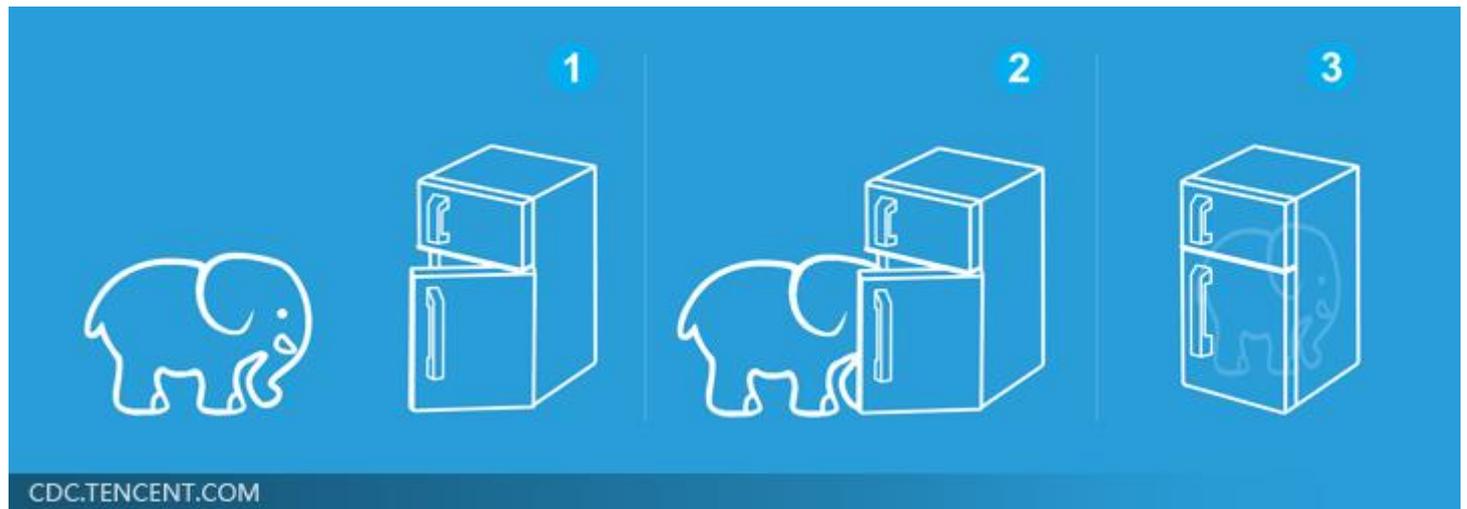


Image Recognition

116



mite

container ship

motor scooter

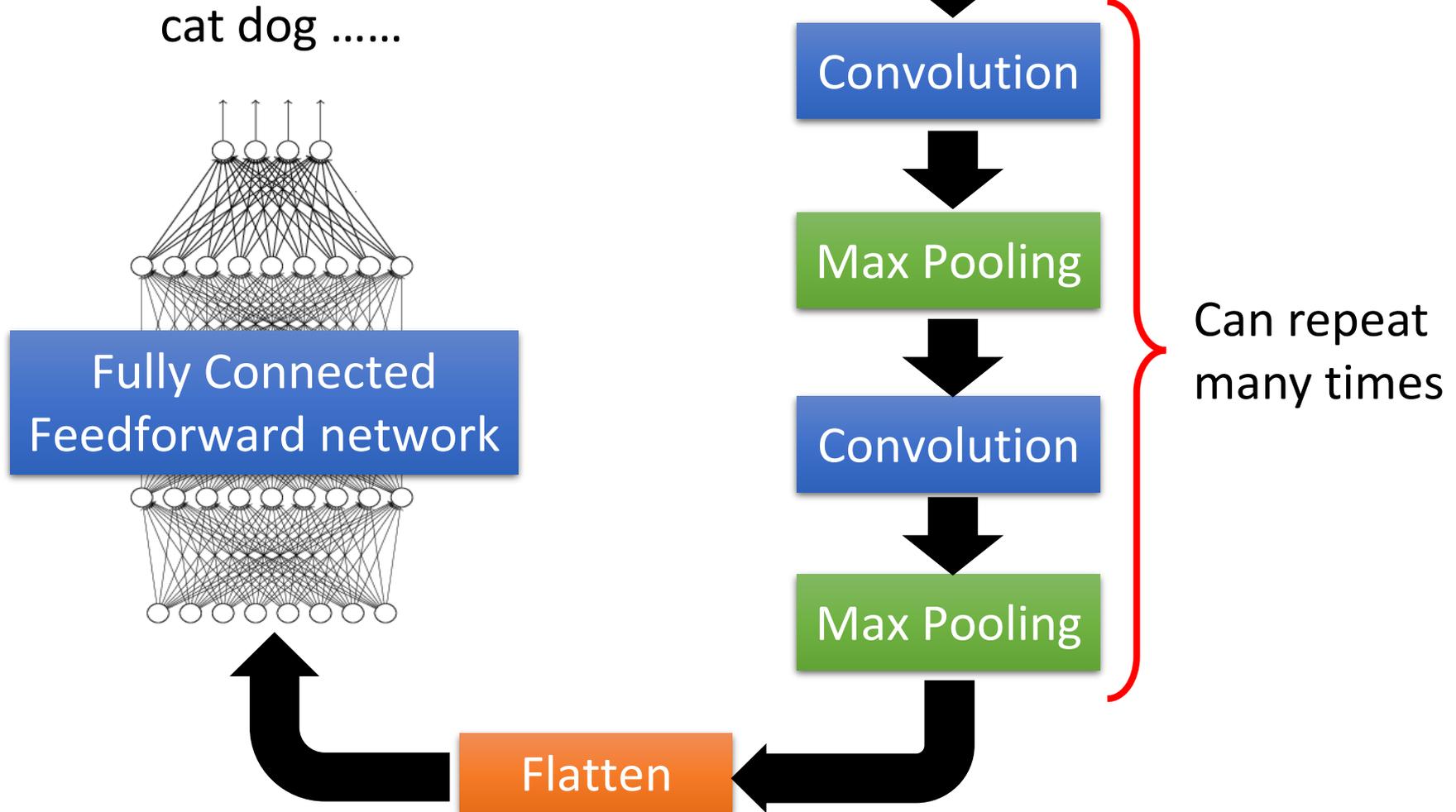
leopard

	mite		container ship		motor scooter		leopard
	black widow		lifeboat		go-kart		jaguar
	cockroach		amphibian		moped		cheetah
	tick		fireboat		bumper car		snow leopard
	starfish		drilling platform		golfcart		Egyptian cat

<http://www.cs.toronto.edu/~fritz/absps/imagenet.pdf>

The Whole CNN

117



The Whole CNN



118

Property 1

- Some patterns are much smaller than the whole image

Property 2

- The same patterns appear in different regions

Property 3

- Subsampling the pixels will not change the object

Convolution

Max Pooling

Convolution

Max Pooling

Flatten

Can repeat many times

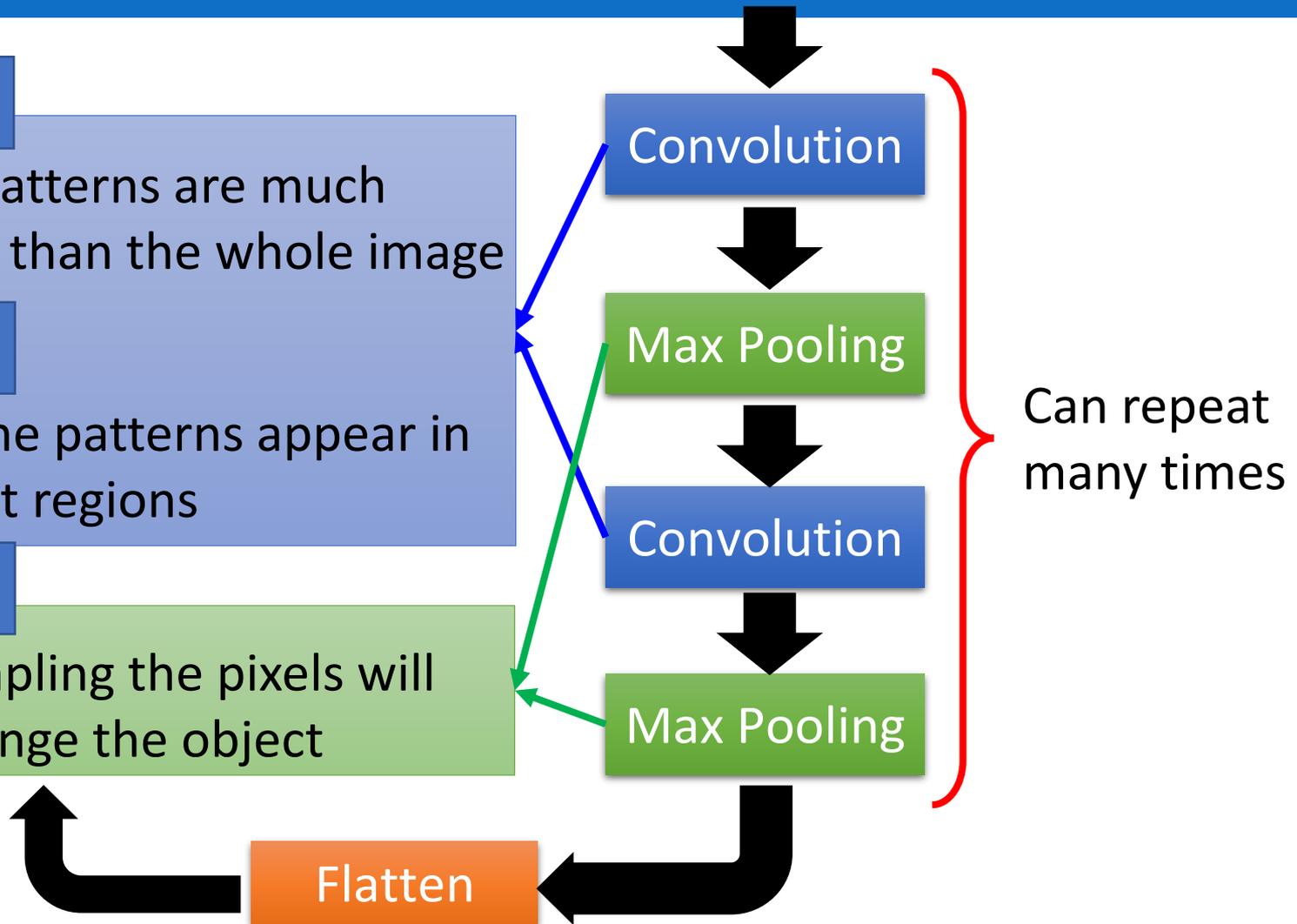
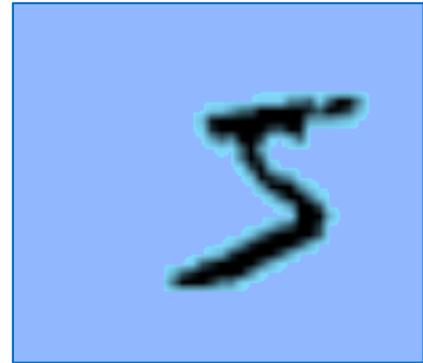


Image Recognition

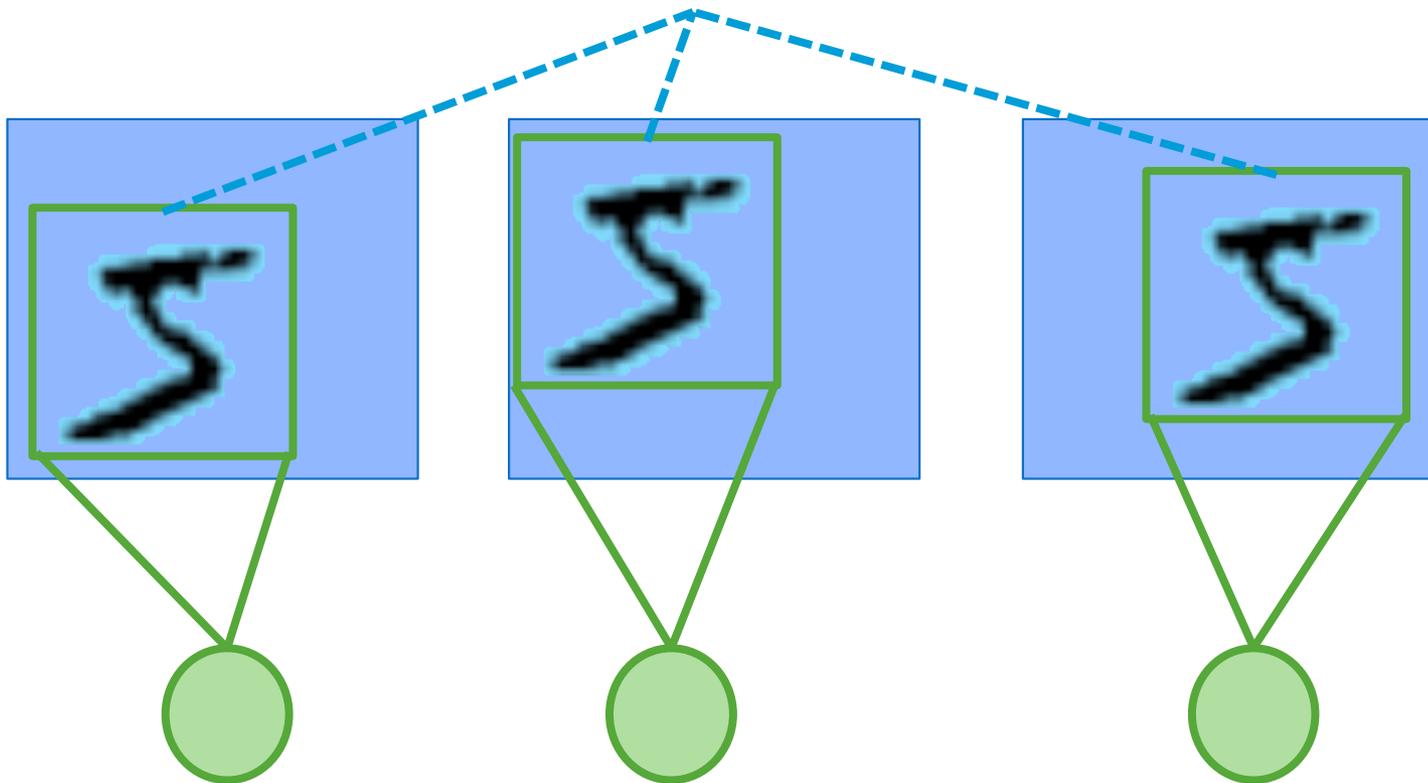
119



Local Connectivity

120

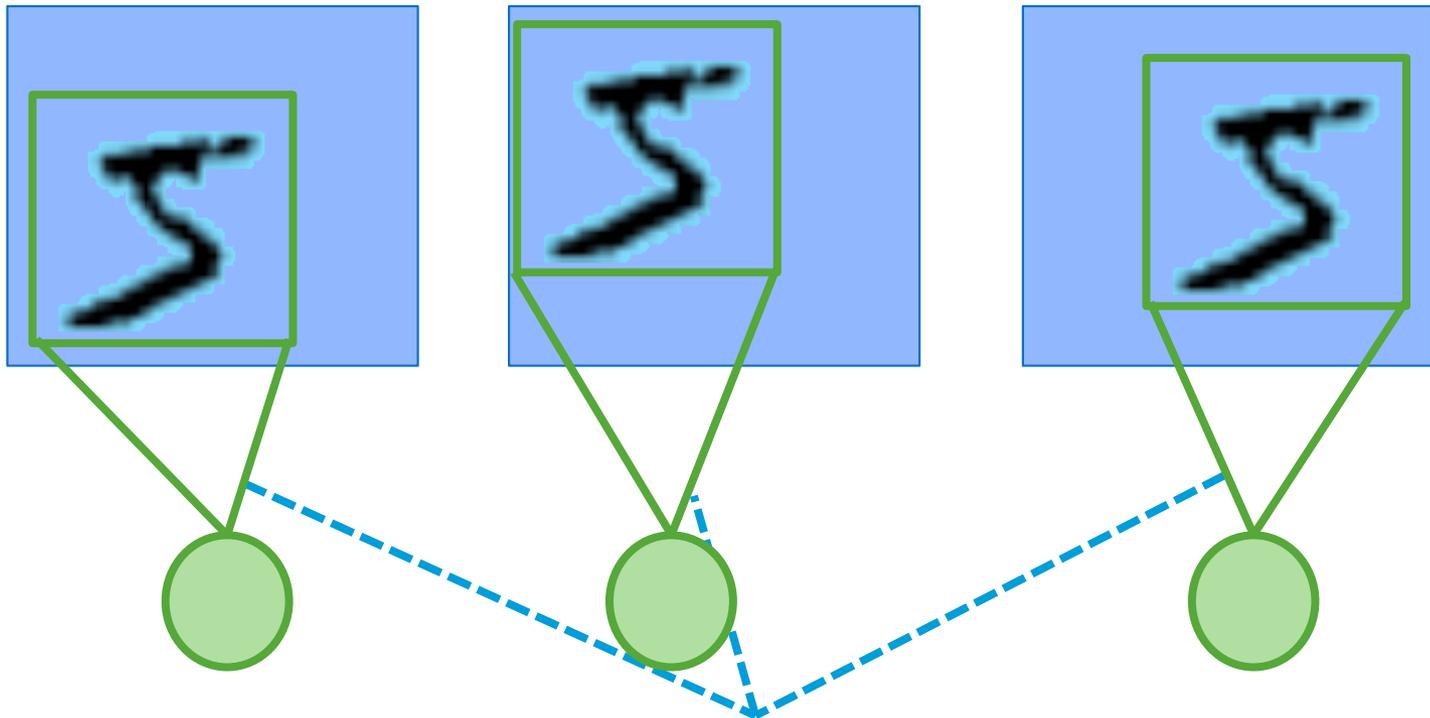
Neurons connect to a small region



Parameter Sharing

121

- The same feature in different positions

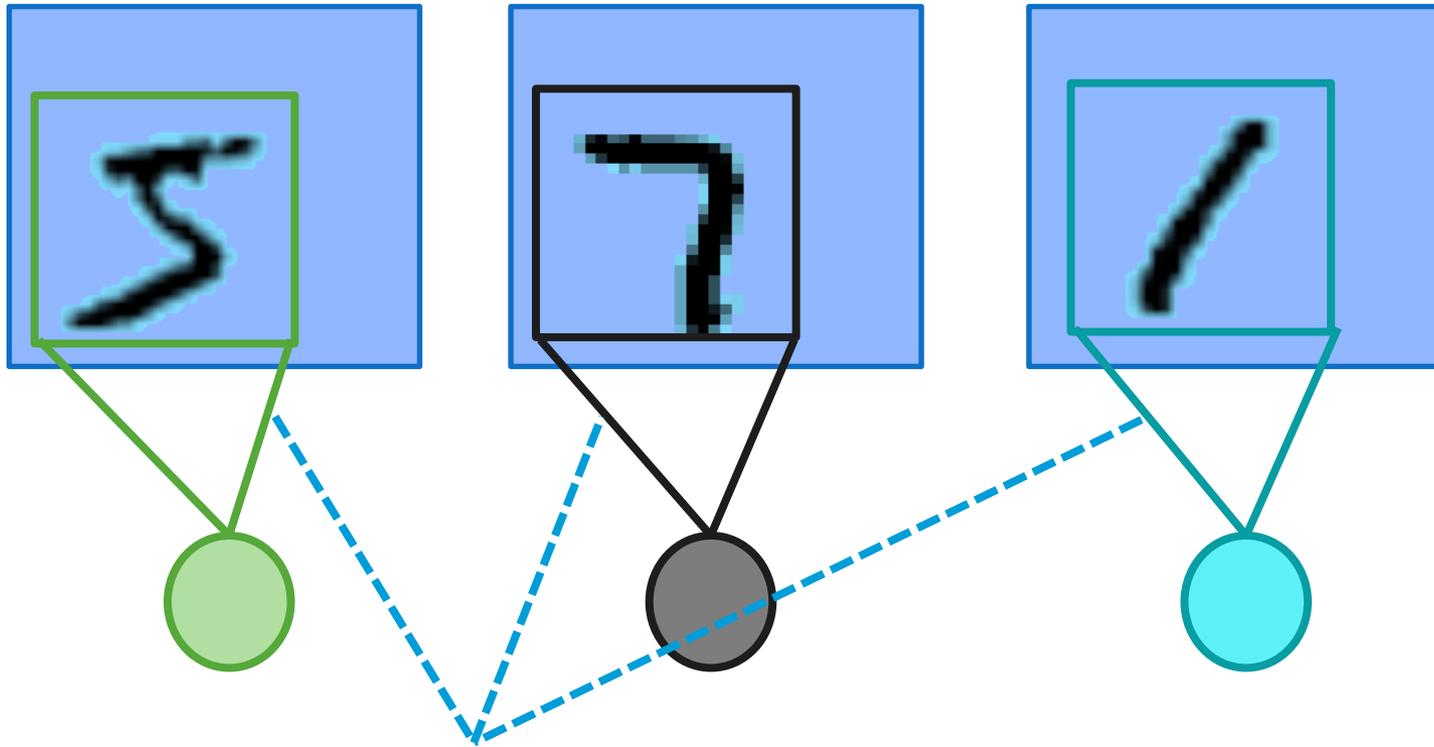


Neurons share the same weights

Parameter Sharing

122

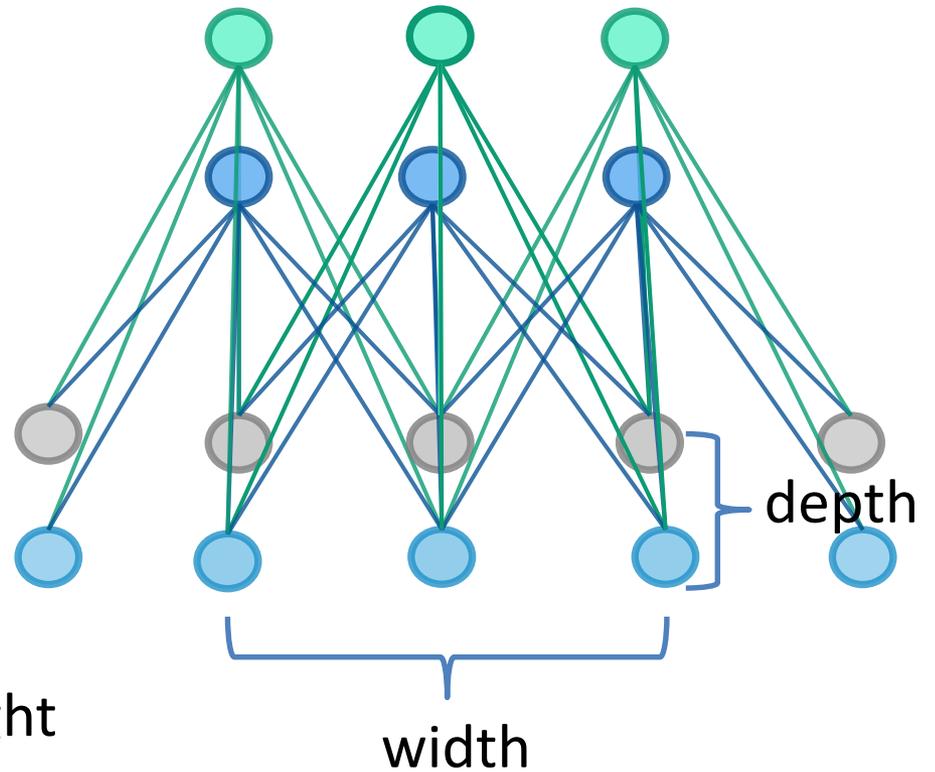
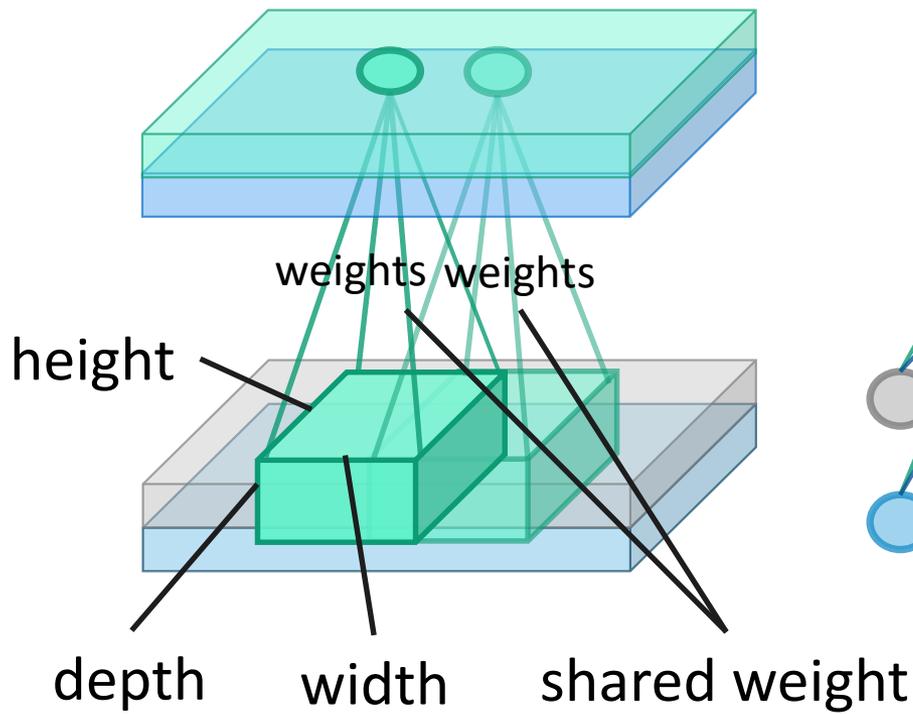
- Different features in the same position



Neurons have different weights

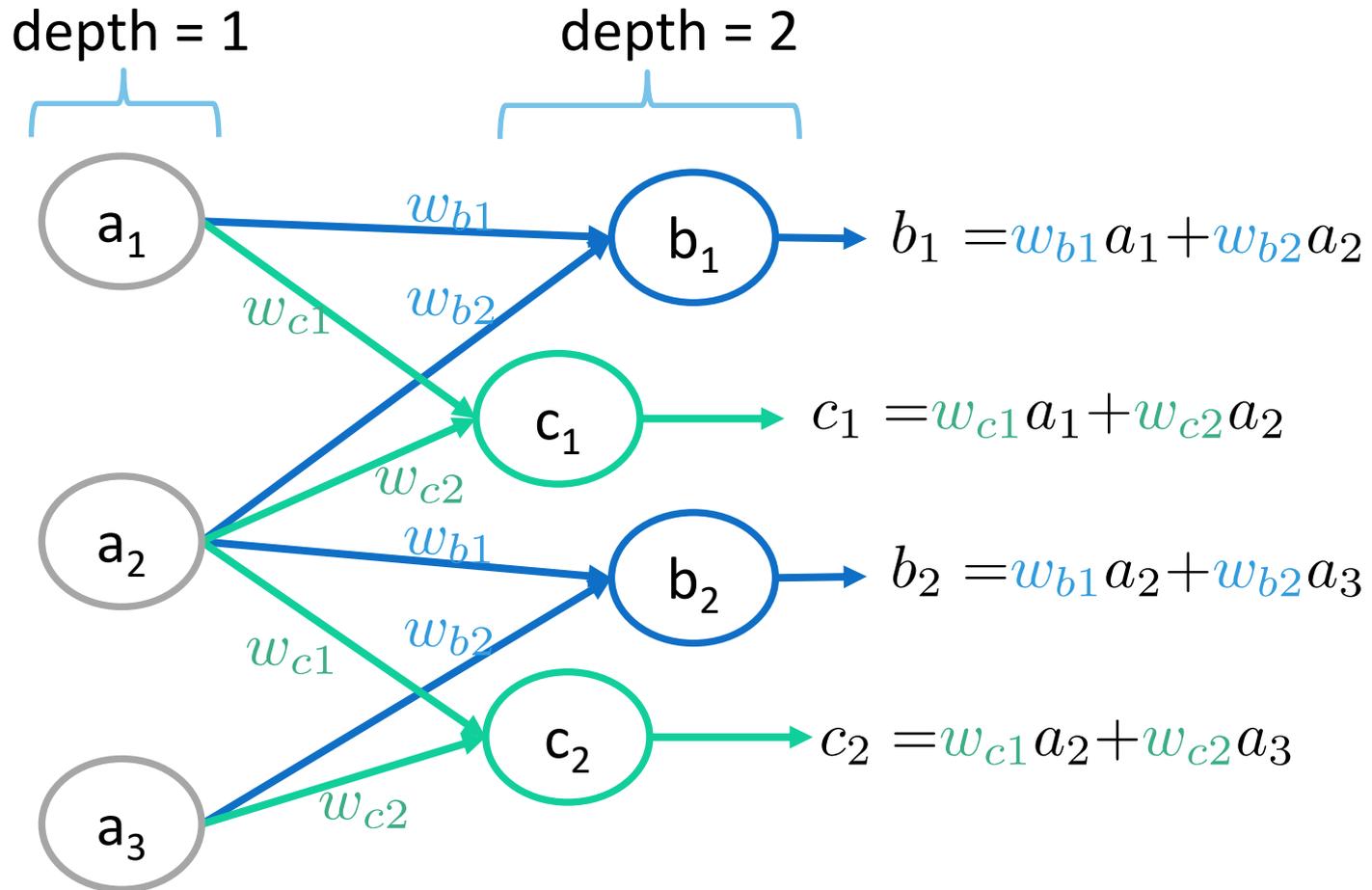
Convolutional Layers

123



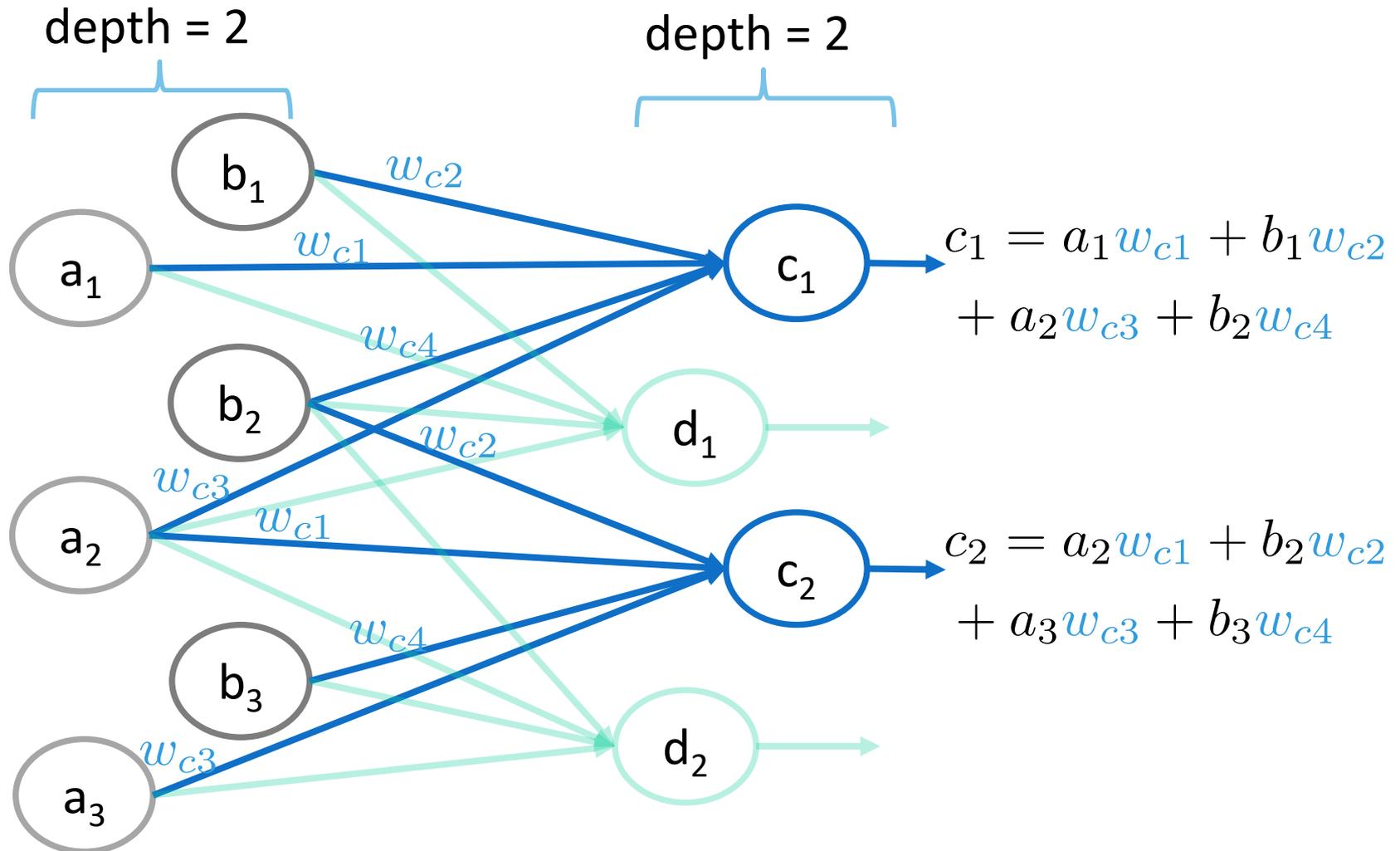
Convolutional Layers

124



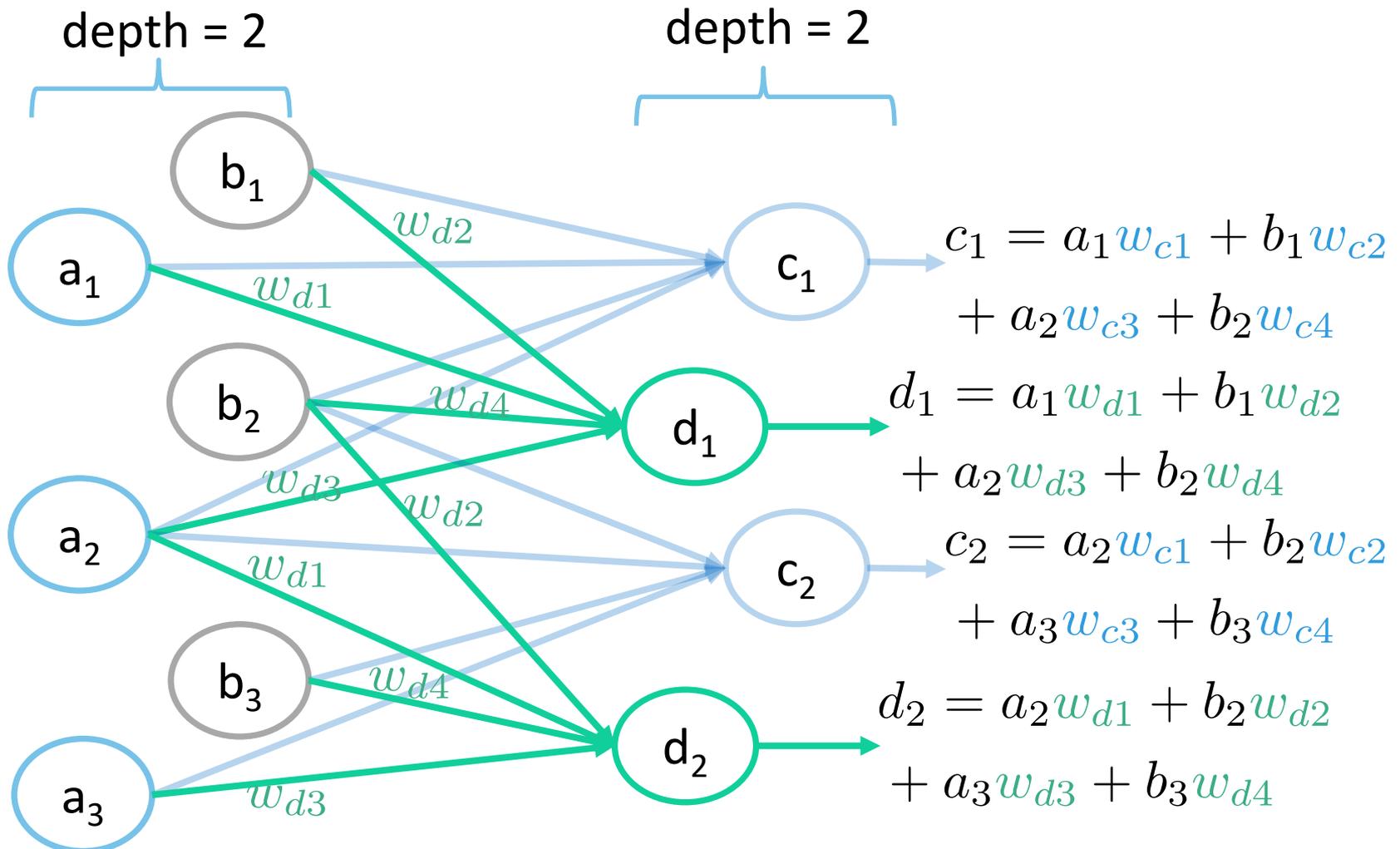
Convolutional Layers

125



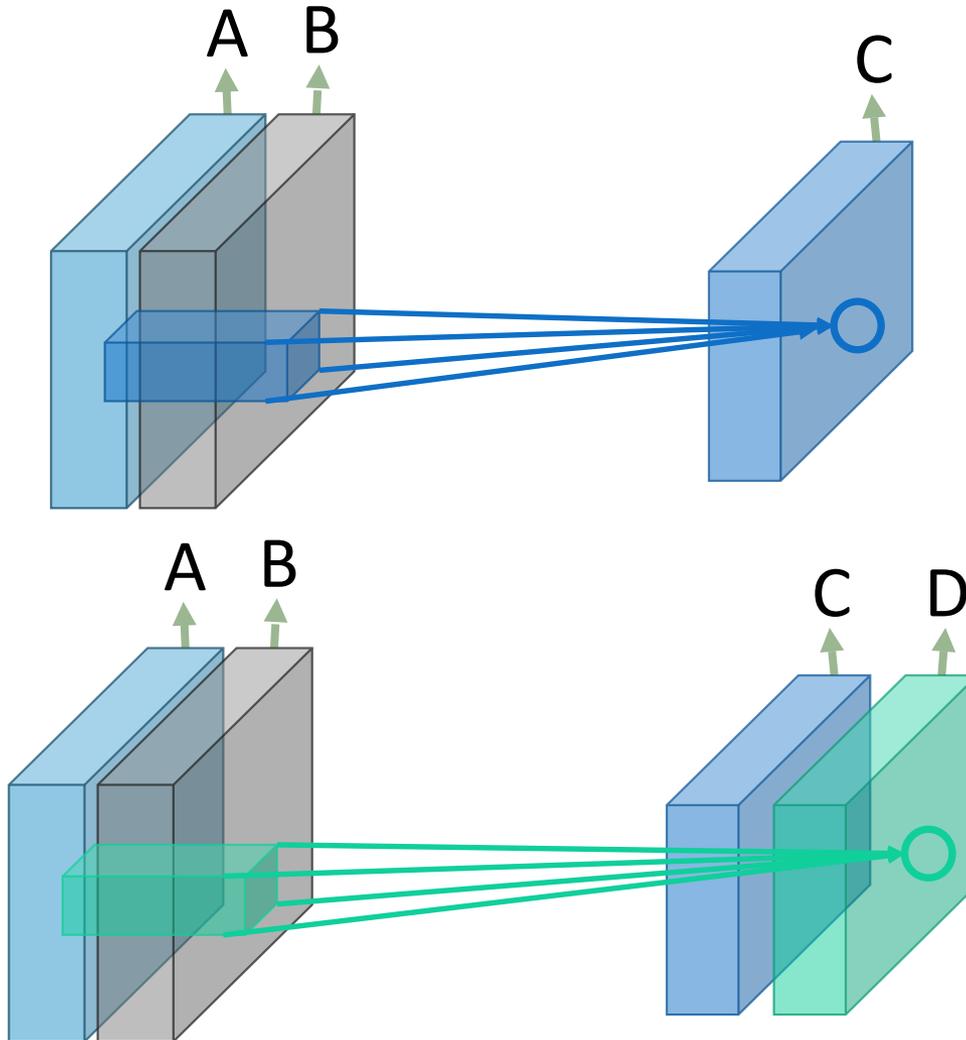
Convolutional Layers

126



Convolutional Layers

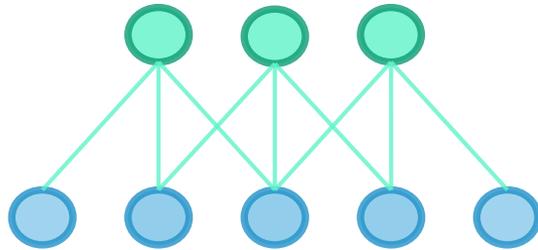
127



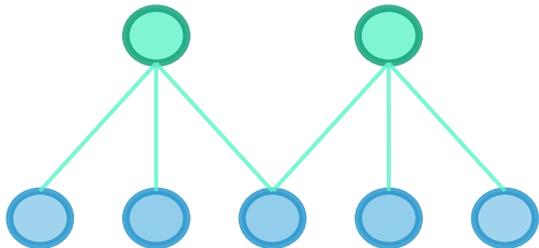
Hyper-parameters of CNN

128

□ Stride

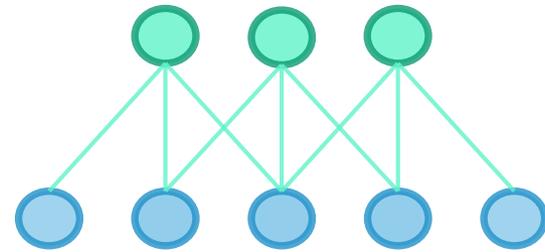


Stride = 1

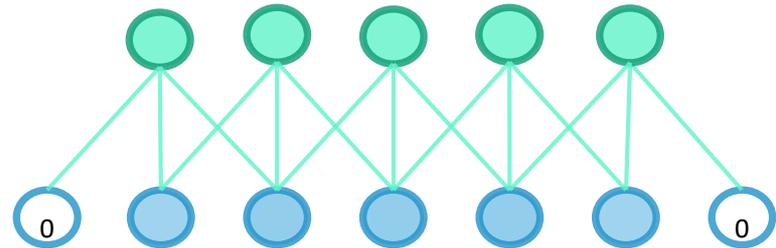


Stride = 2

□ Padding



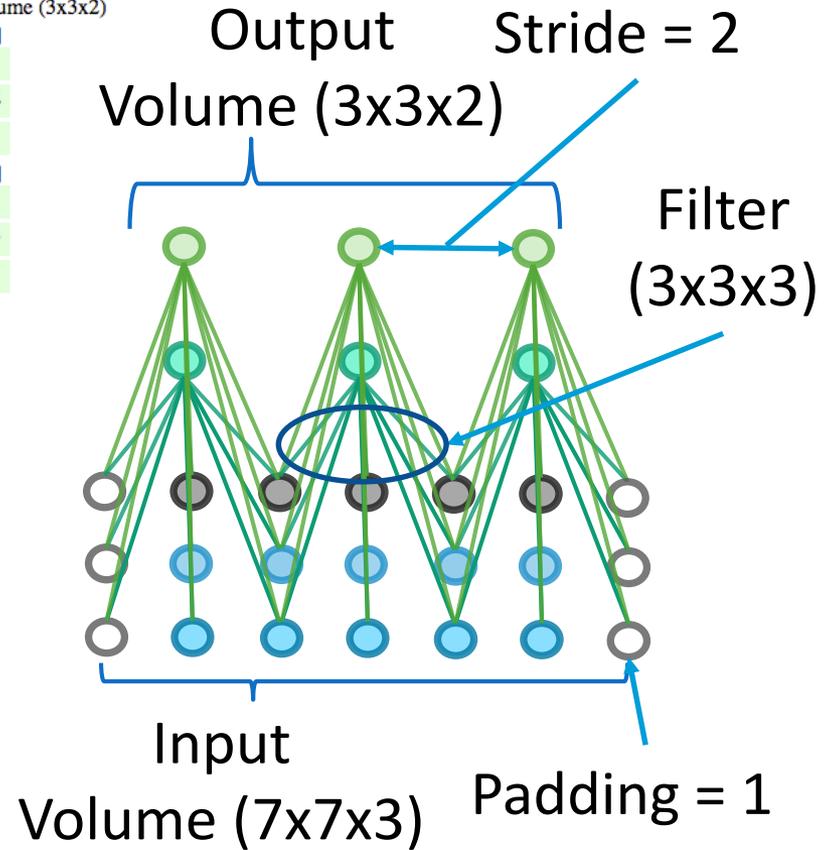
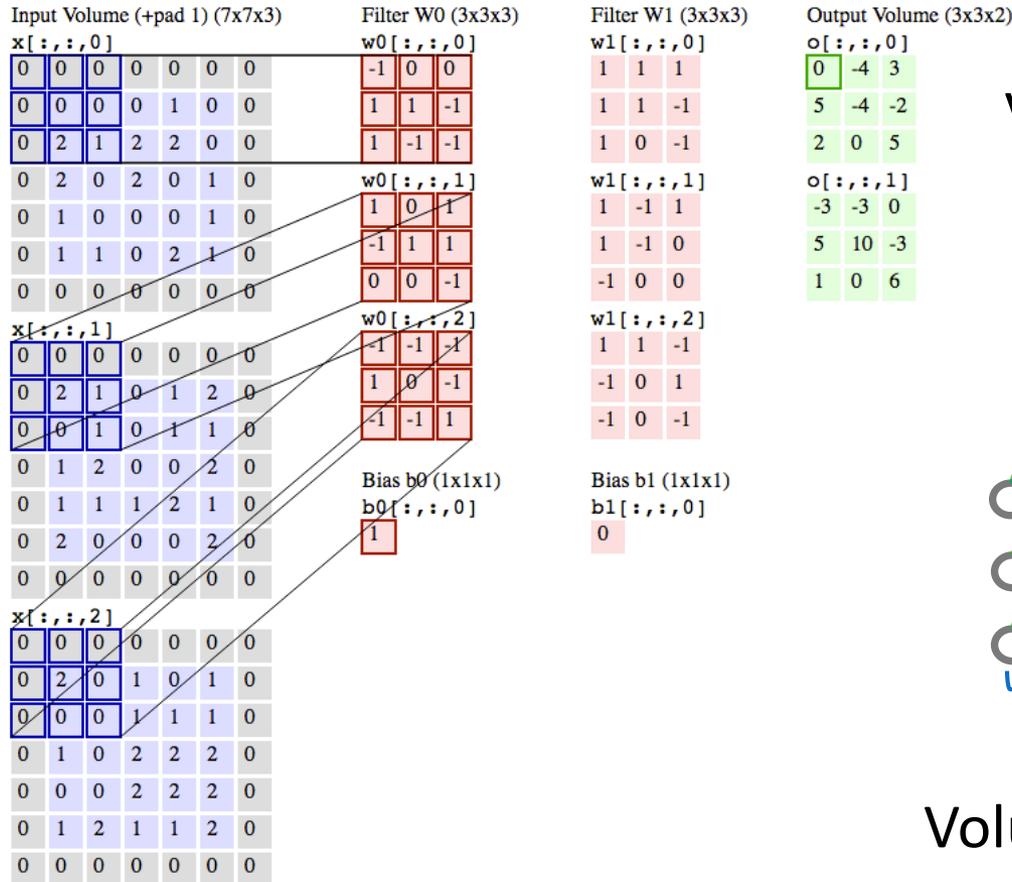
Padding = 0



Padding = 1

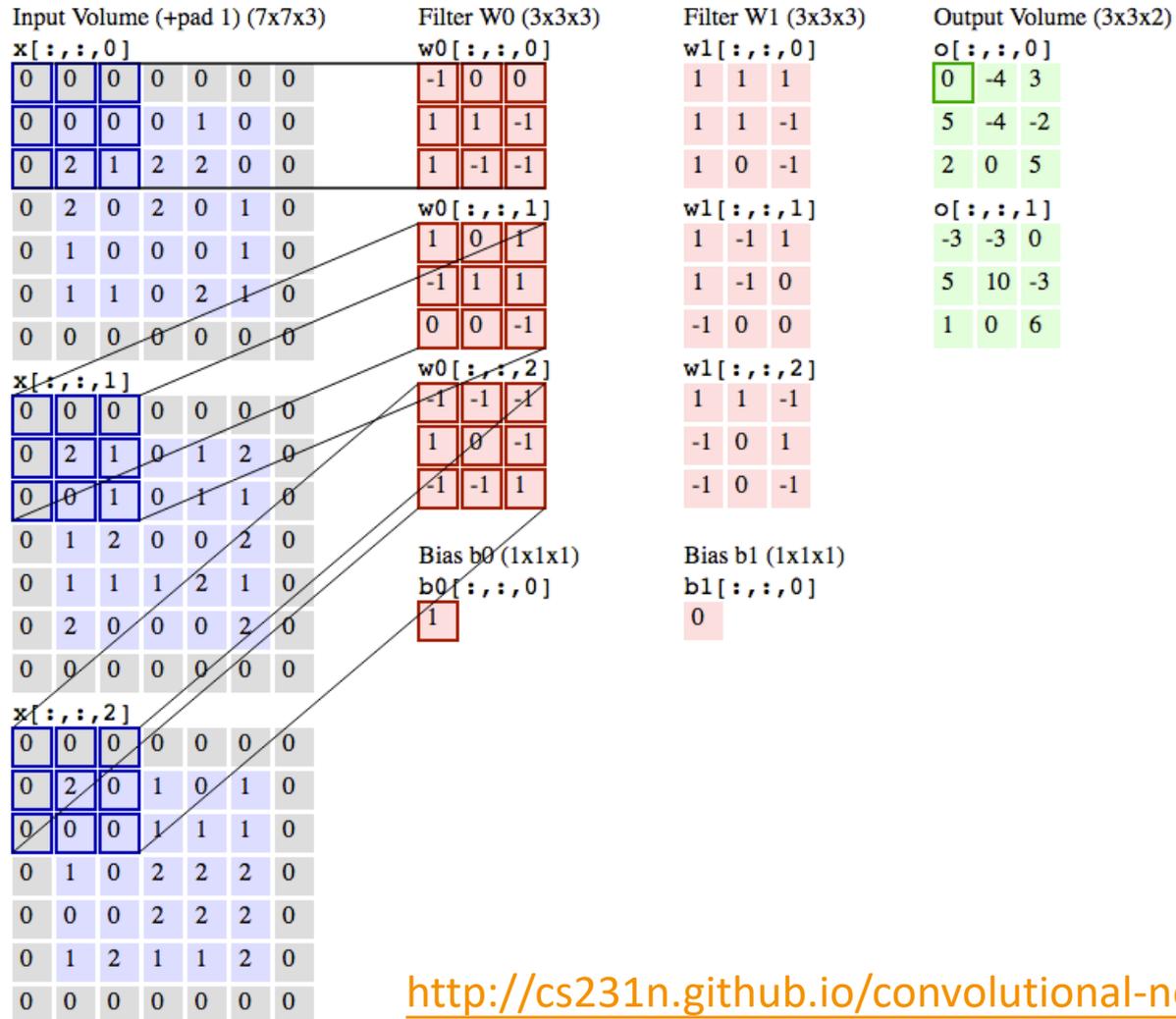
Example

129



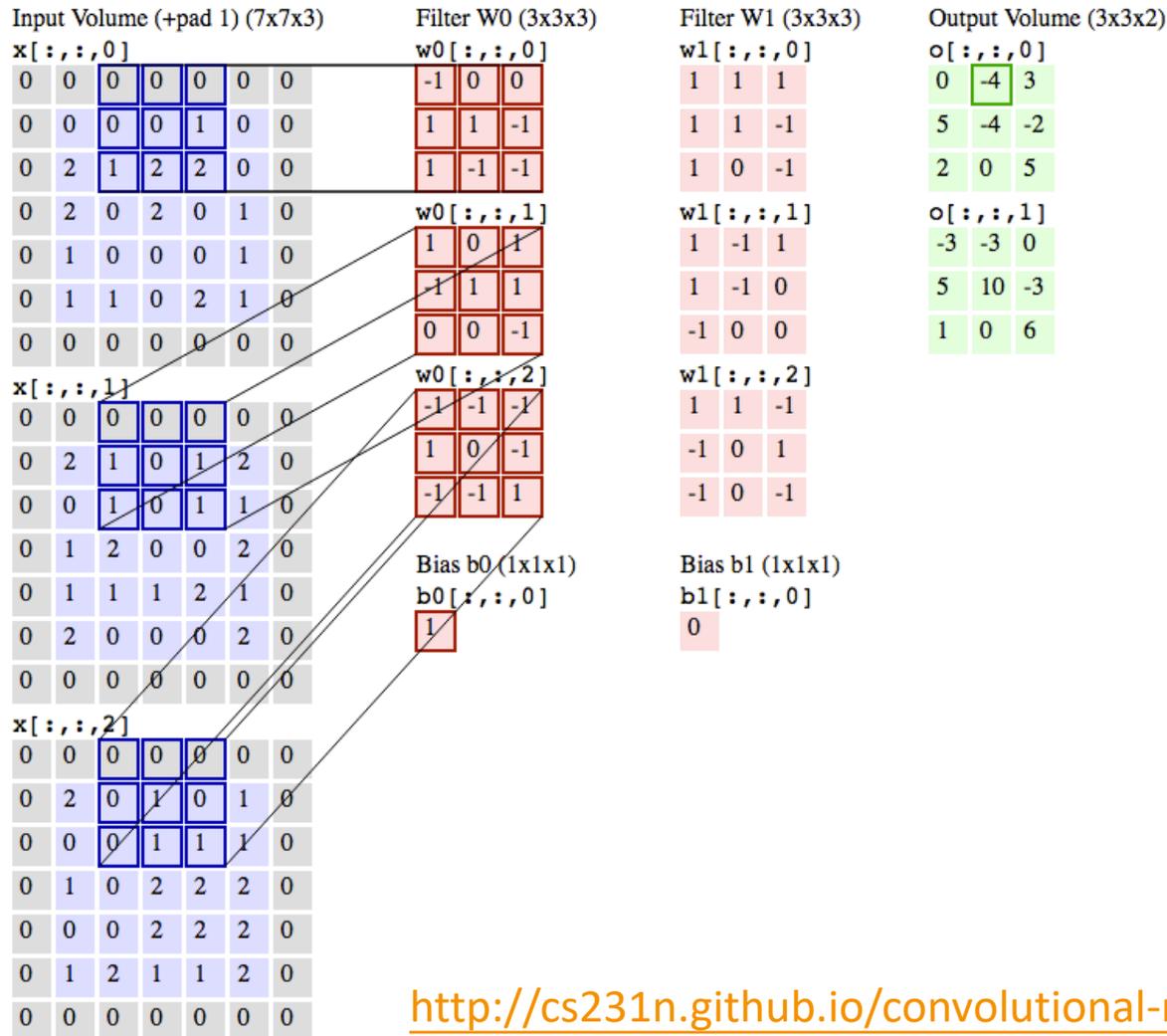
Convolutional Layers

130



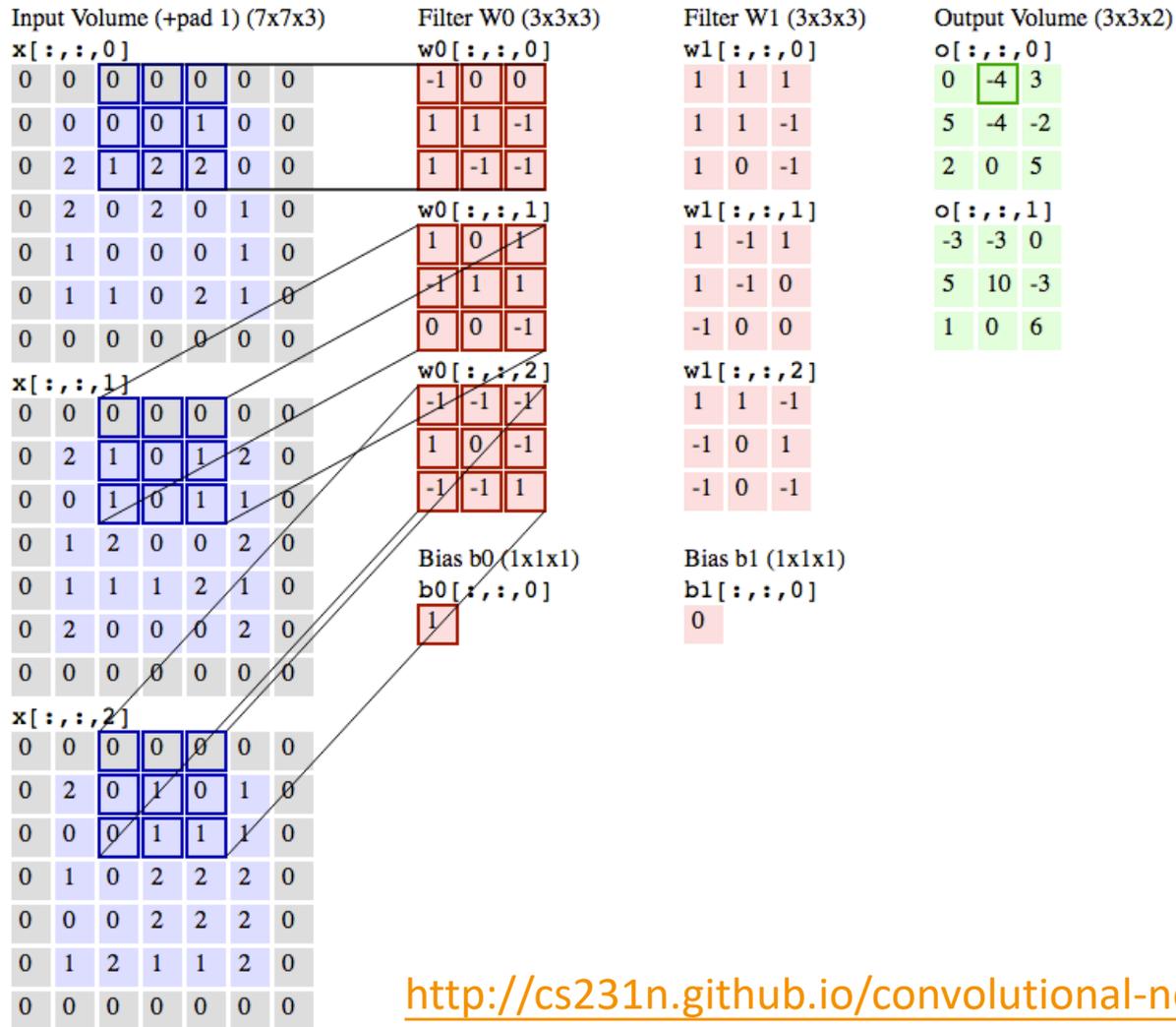
Convolutional Layers

131



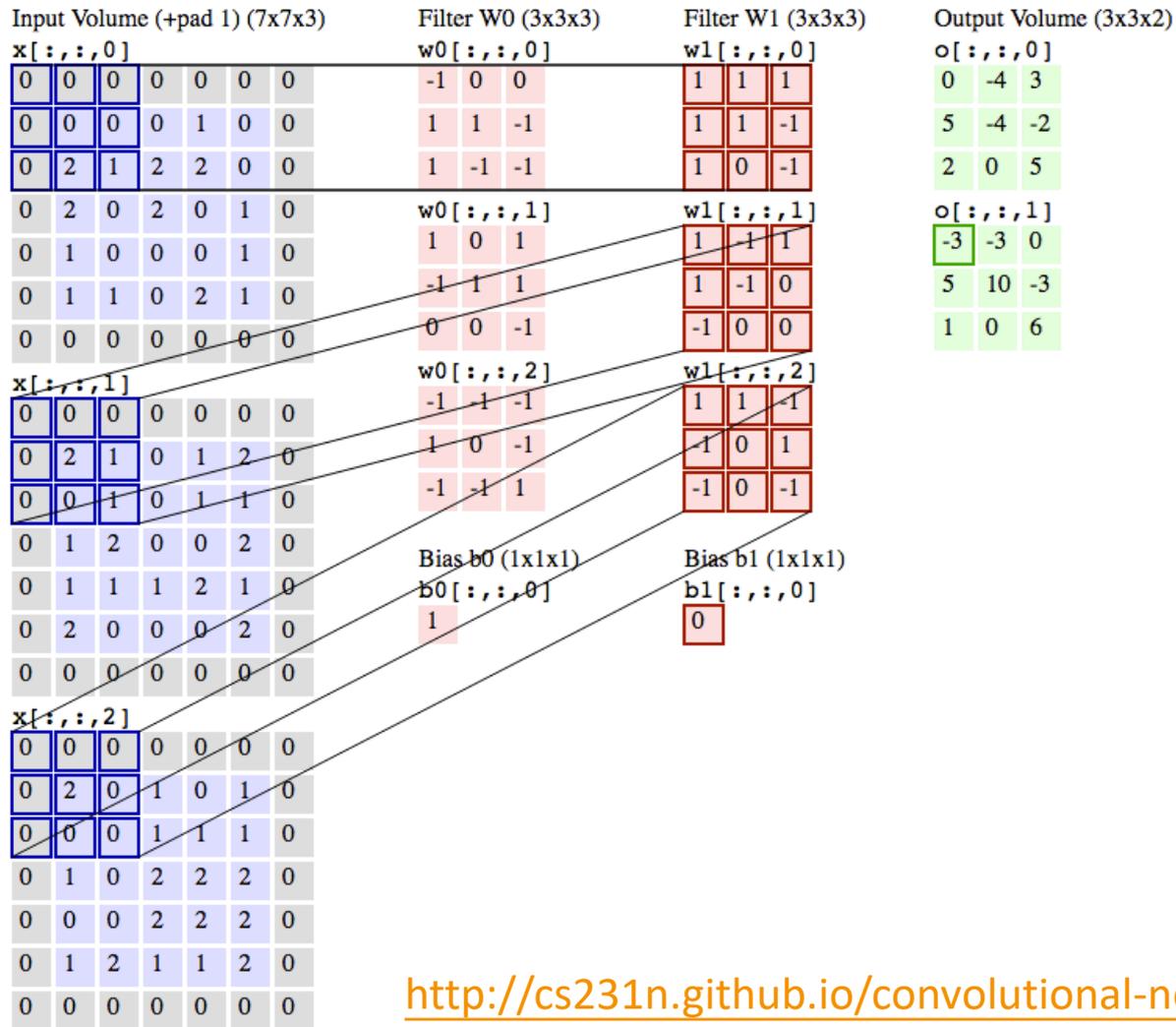
Convolutional Layers

132



<http://cs231n.github.io/convolutional-networks/>

Convolutional Layers



Pooling Layer

134

1	3	2	4
5	7	6	8
0	0	3	3
5	5	0	0

Maximum Pooling



7	8
5	3

$$\text{Max}(1, 3, 5, 7) = 7$$

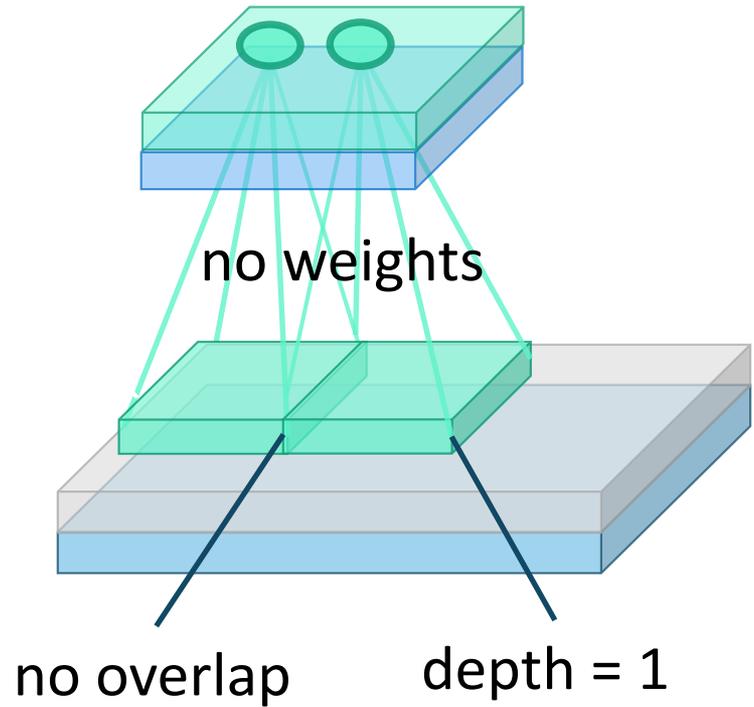
$$\text{Max}(0, 0, 5, 5) = 5$$

Average Pooling



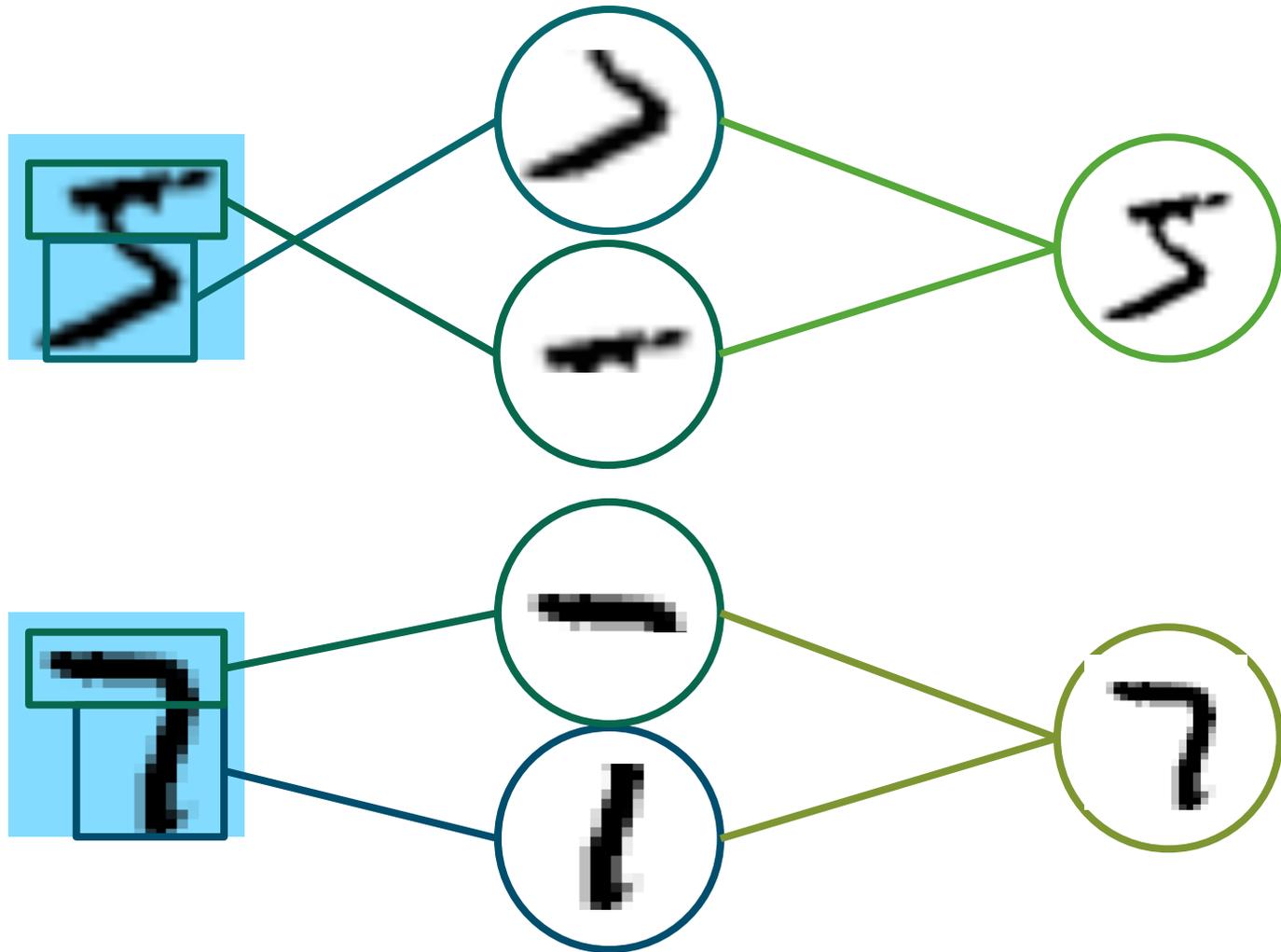
4	5
5	3

$$\text{Avg}(1, 3, 5, 7) = 4$$



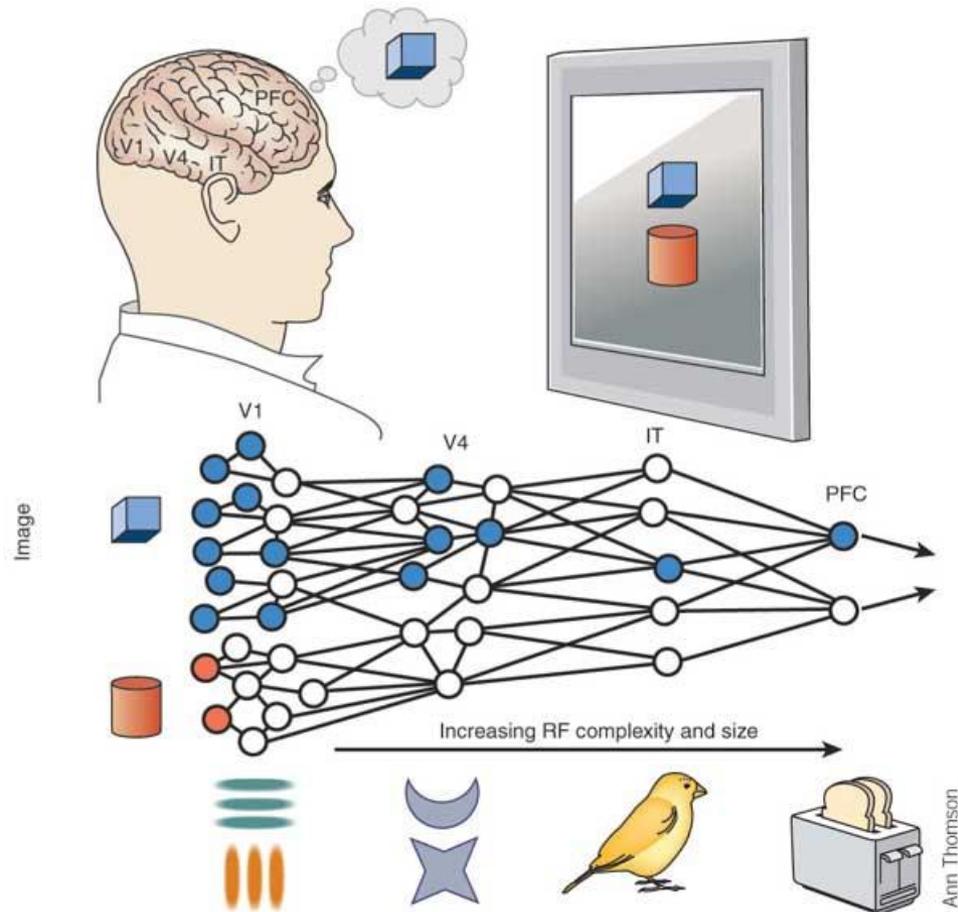
Why “Deep” Learning?

135



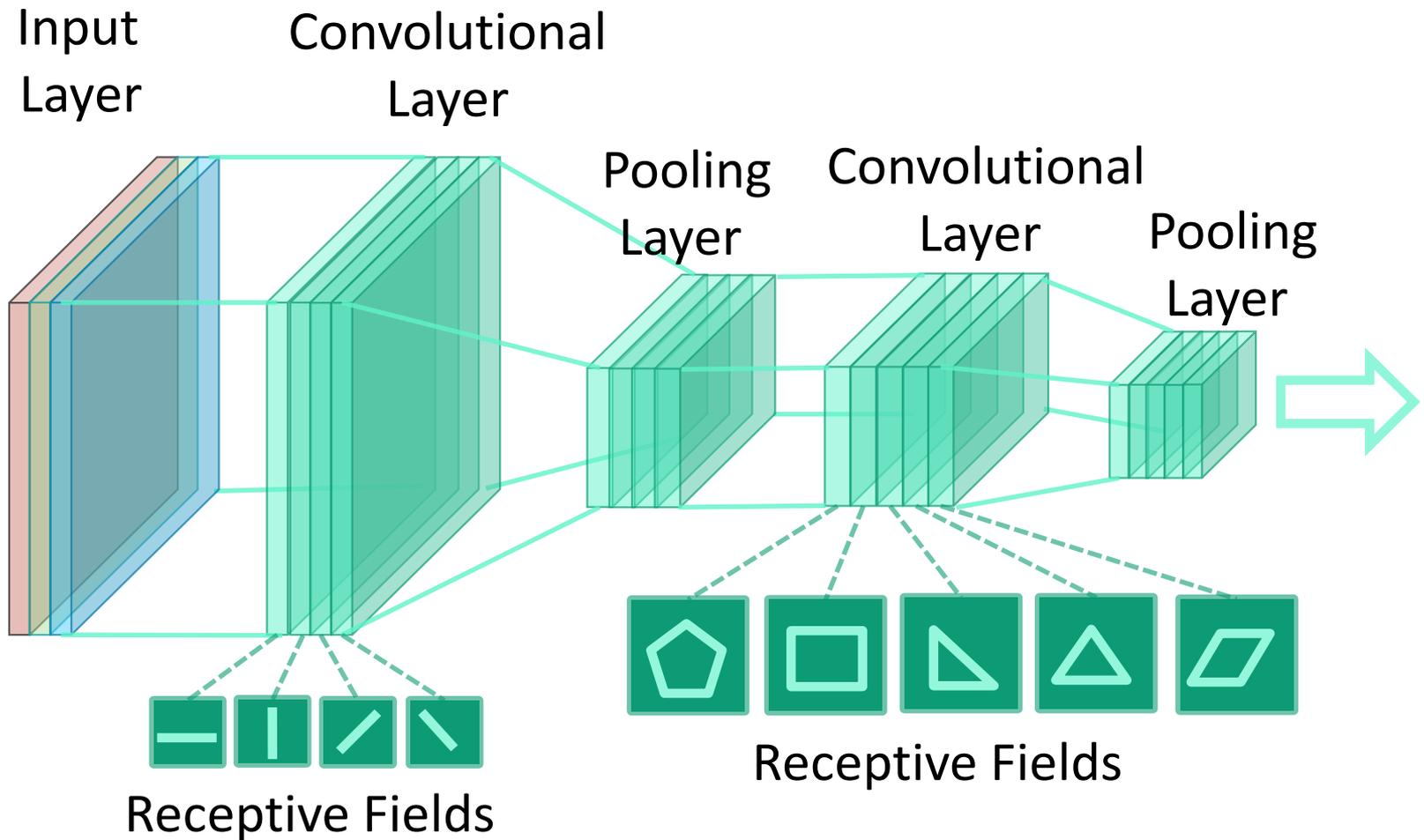
Visual Perception of Human

136



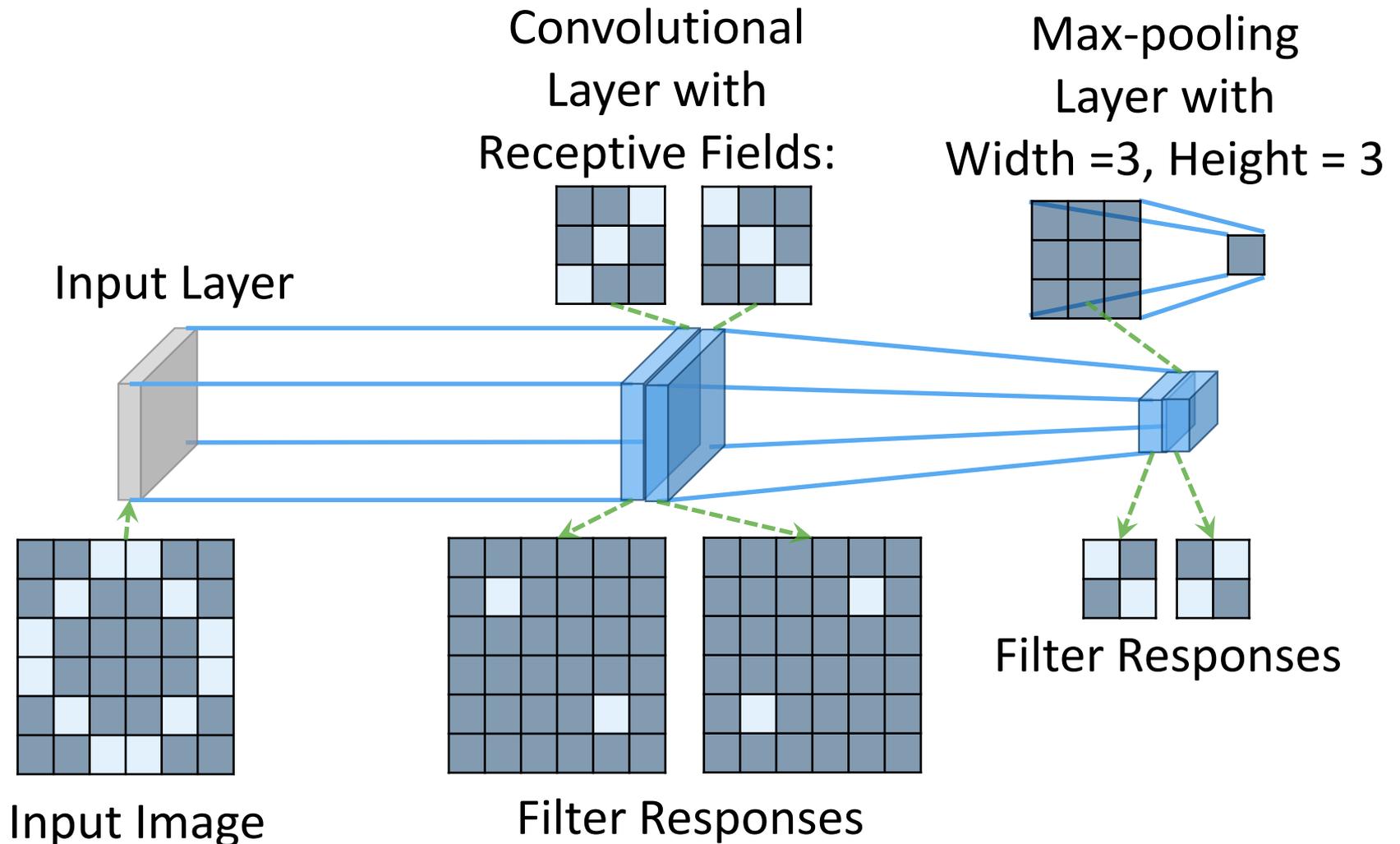
Visual Perception of Computer

137



Visual Perception of Computer

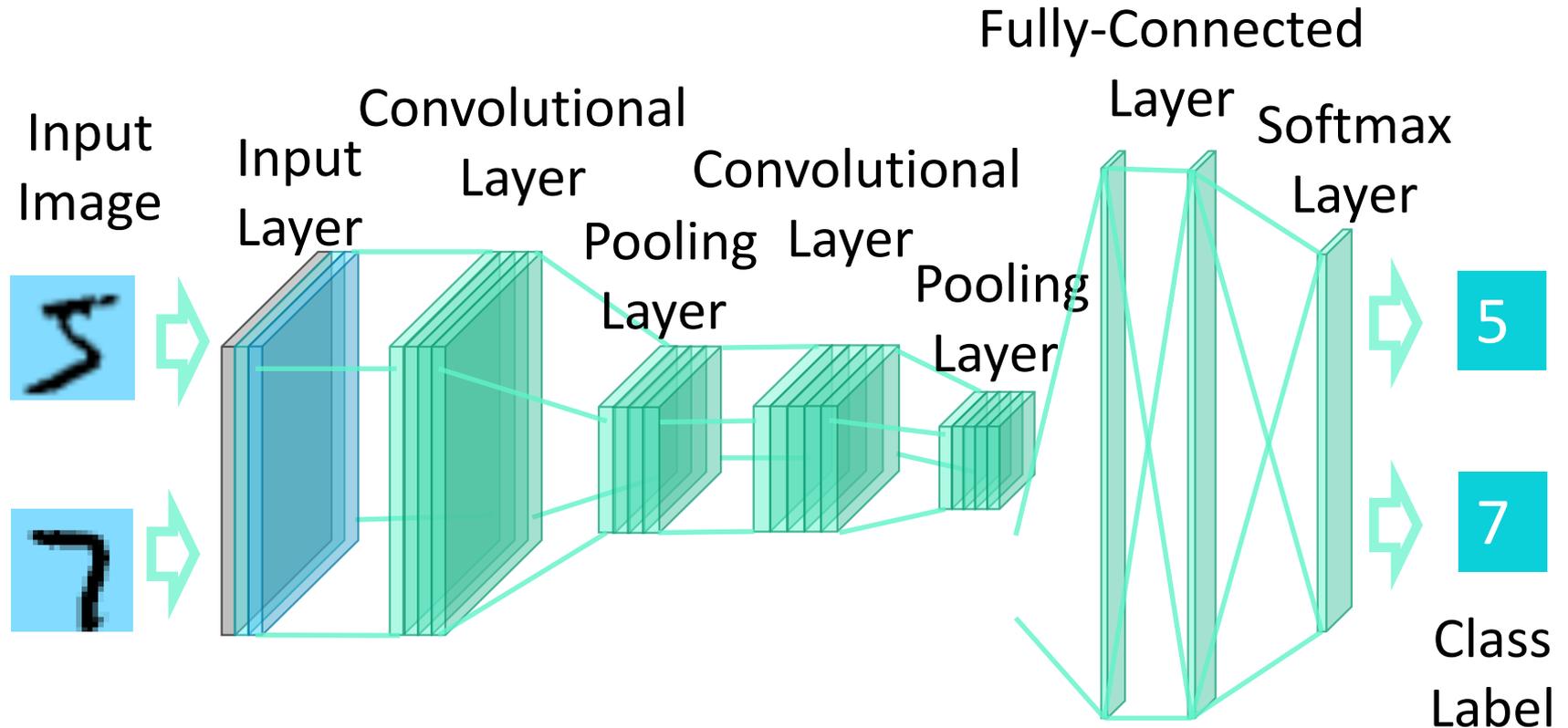
138



Fully-Connected Layer

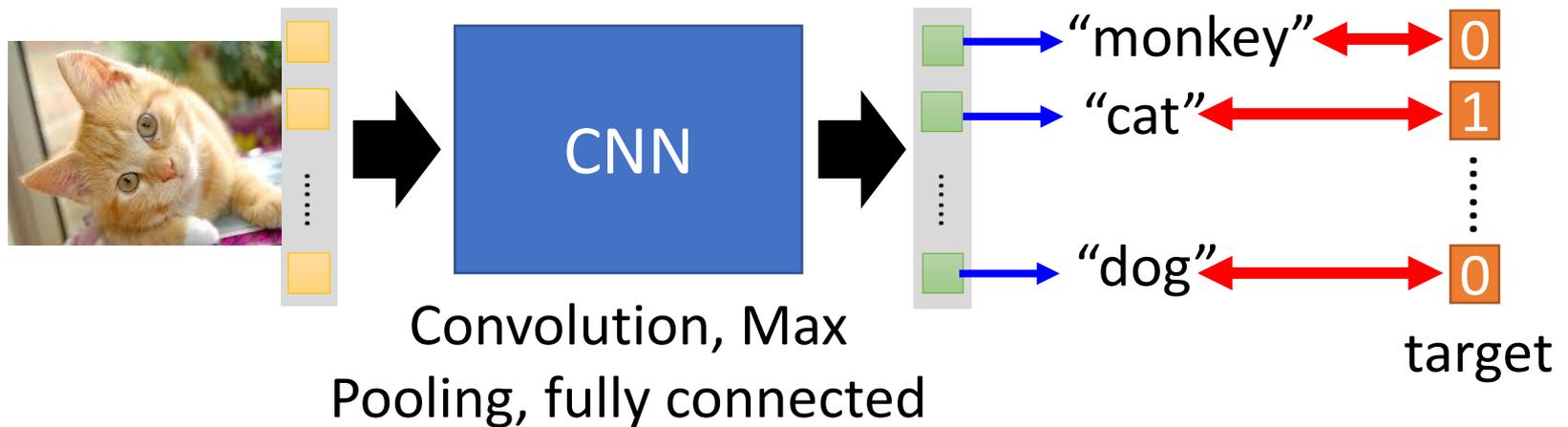
139

- Fully-Connected Layers : Global feature extraction
- Softmax Layer: Classifier



Convolutional Neural Network

140

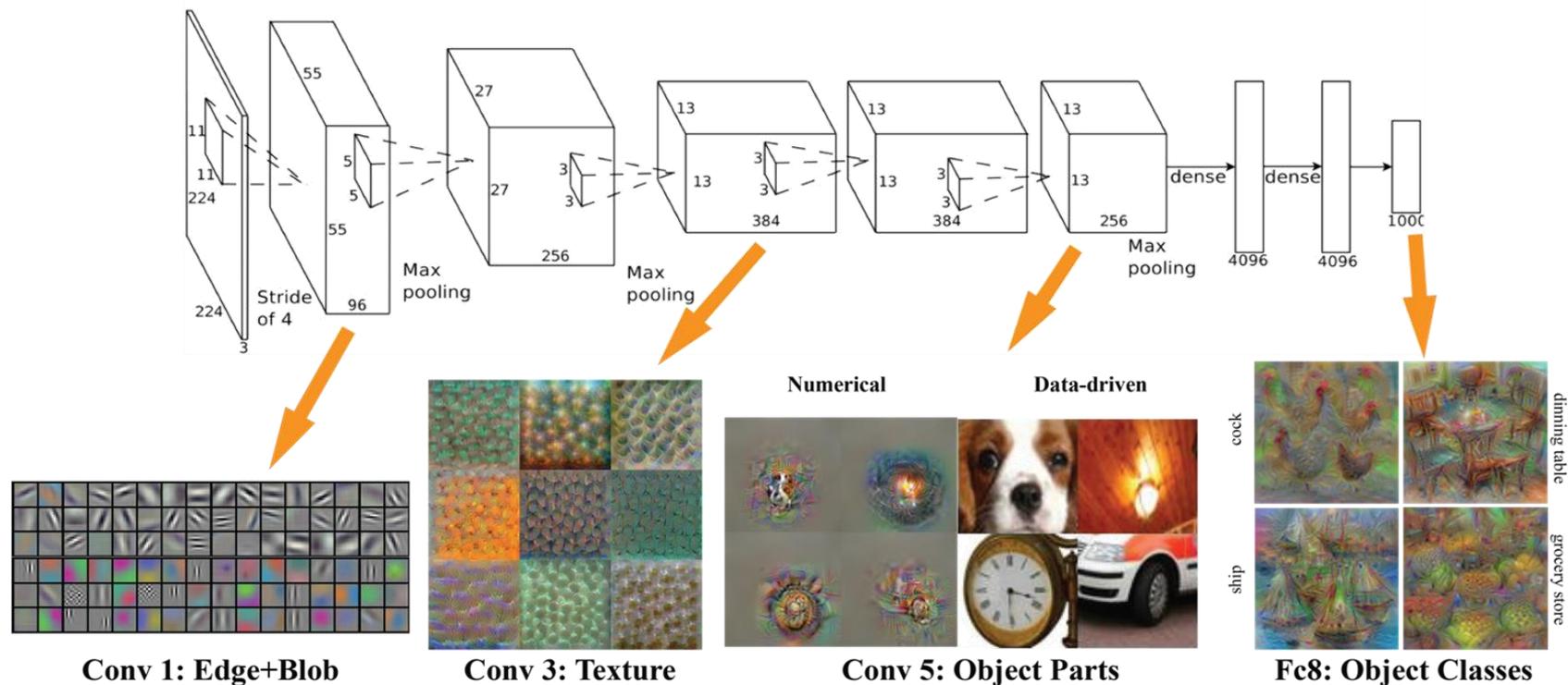


What CNN Learned

141

□ Alexnet

- ▣ <http://www.cs.toronto.edu/~fritz/absps/imagenet.pdf>



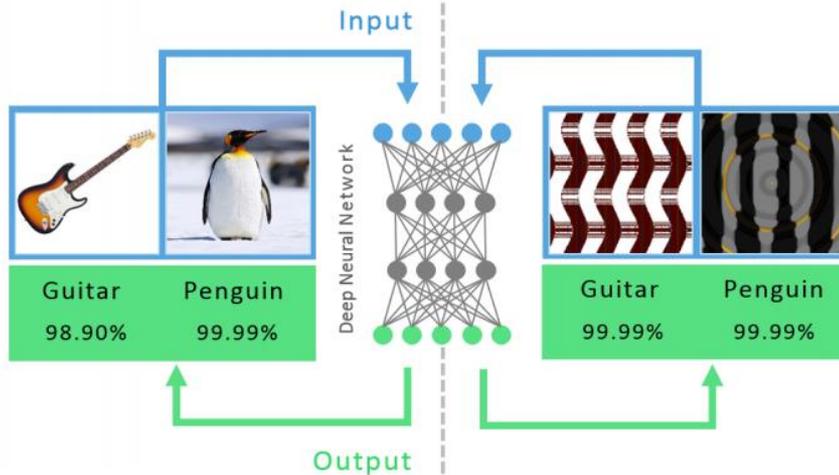
http://vision03.csail.mit.edu/cnn_art/data/single_layer.png

DNN are easily fooled

142

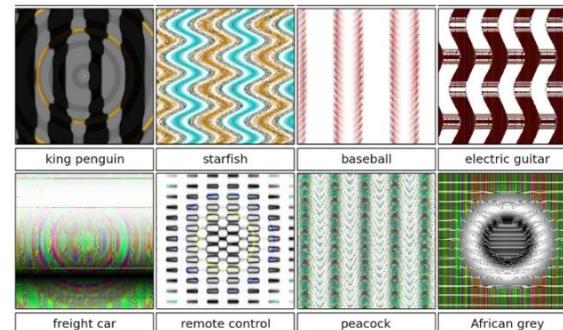
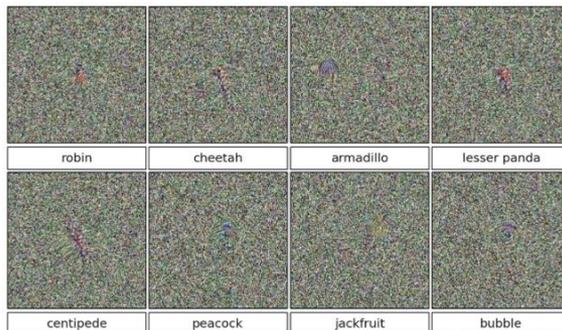
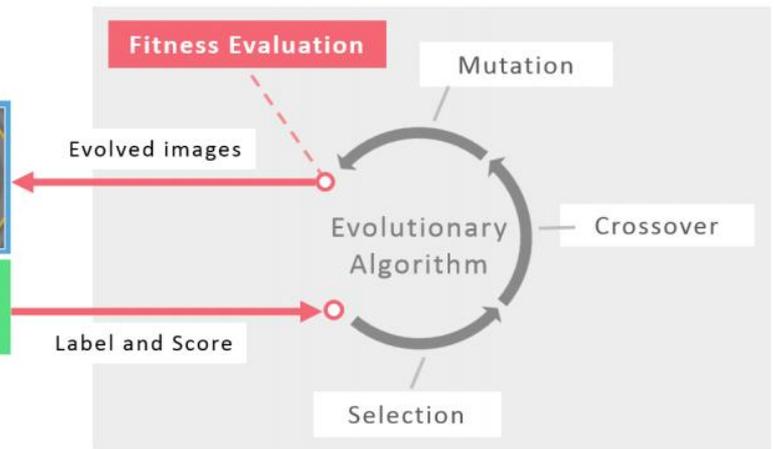
1

State-of-the-art DNNs can recognize real images with high confidence



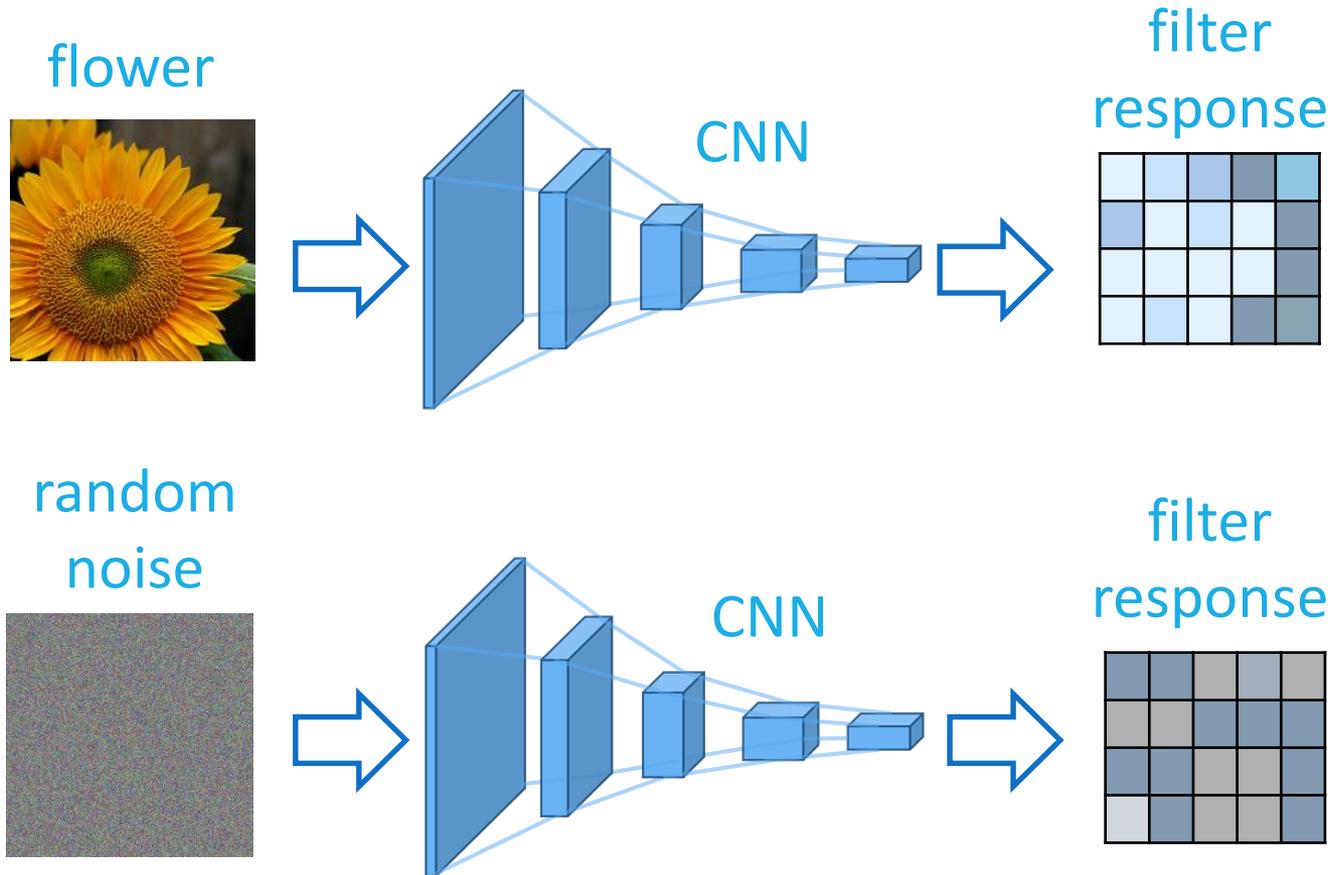
2

But DNNs are also easily fooled: images can be produced that are unrecognizable to humans, but DNNs believe with 99.99% certainty are natural objects



Visualizing CNN

143



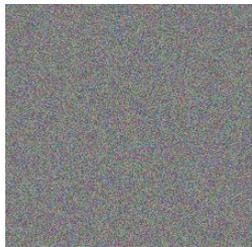
Gradient Ascent

144

- Magnify the filter response

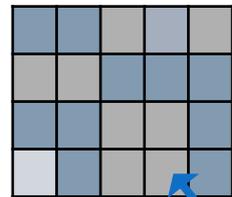
random

noise: \mathbf{x}



filter

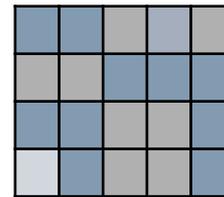
response: \mathbf{f}



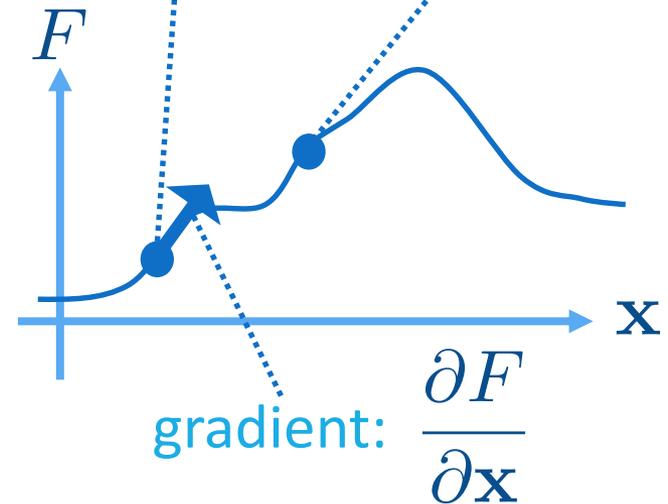
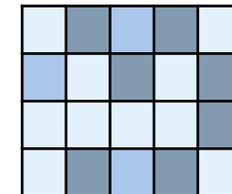
$f_{i,j}$

score: $F = \sum_{i,j} f_{i,j}$

lower
score



higher
score



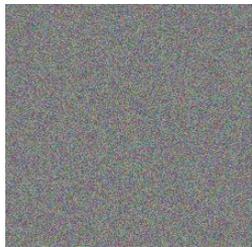
Gradient Ascent

145

- Magnify the filter response

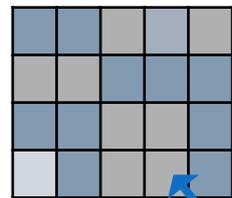
random

noise: \mathbf{x}



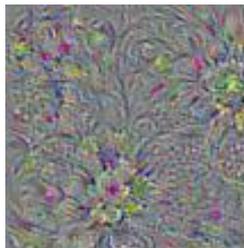
filter

response: \mathbf{f}



$f_{i,j}$

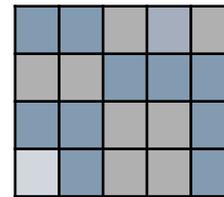
update \mathbf{x}



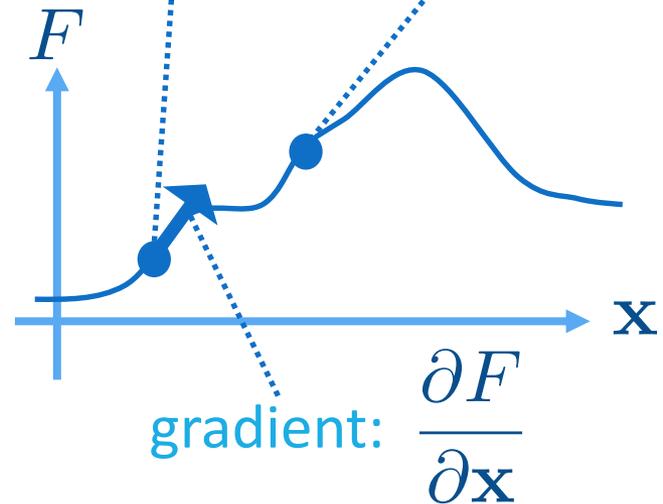
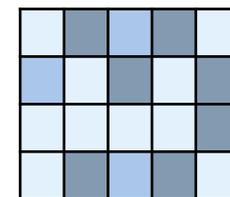
$$\mathbf{x} \leftarrow \mathbf{x} + \eta \frac{\partial F}{\partial \mathbf{x}}$$

learning rate

lower
score

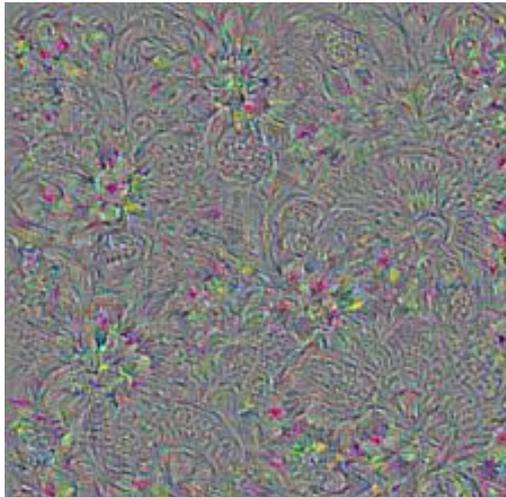
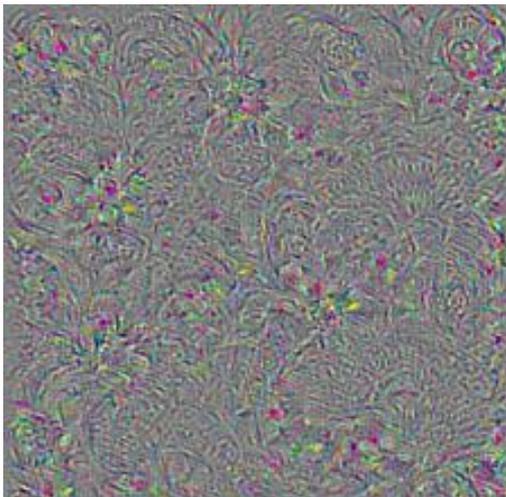
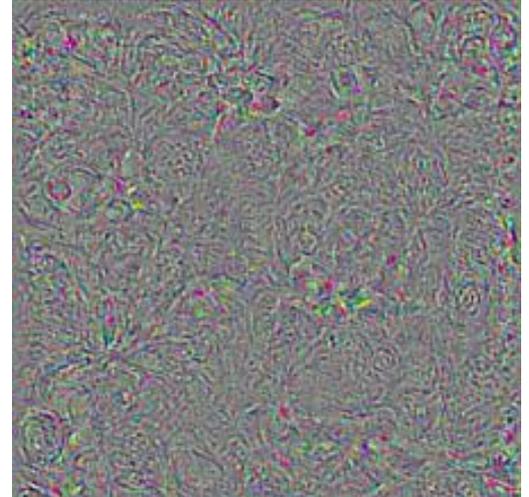
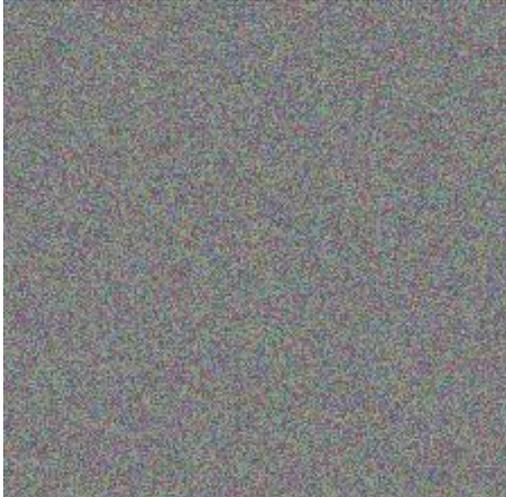


higher
score



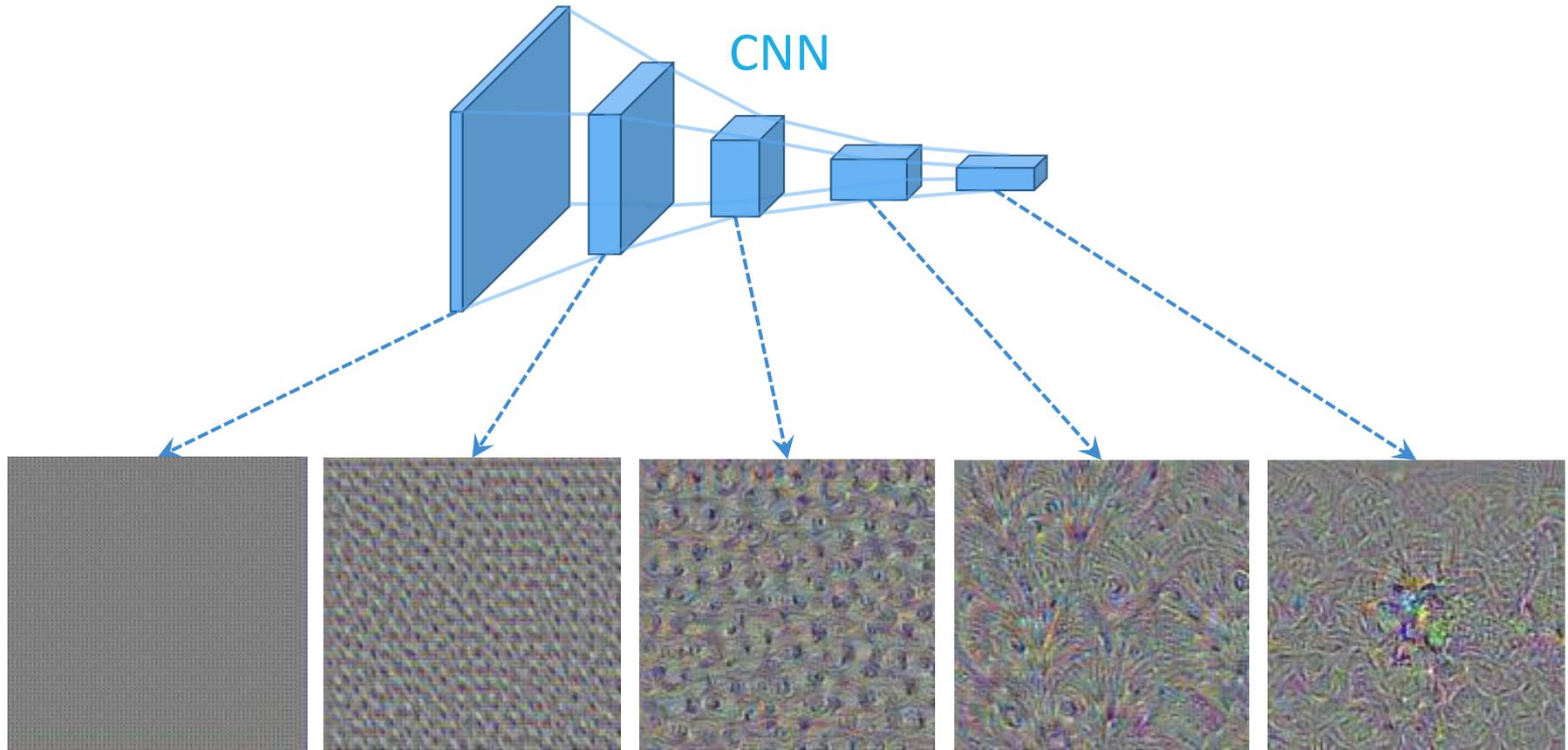
Gradient Ascent

146



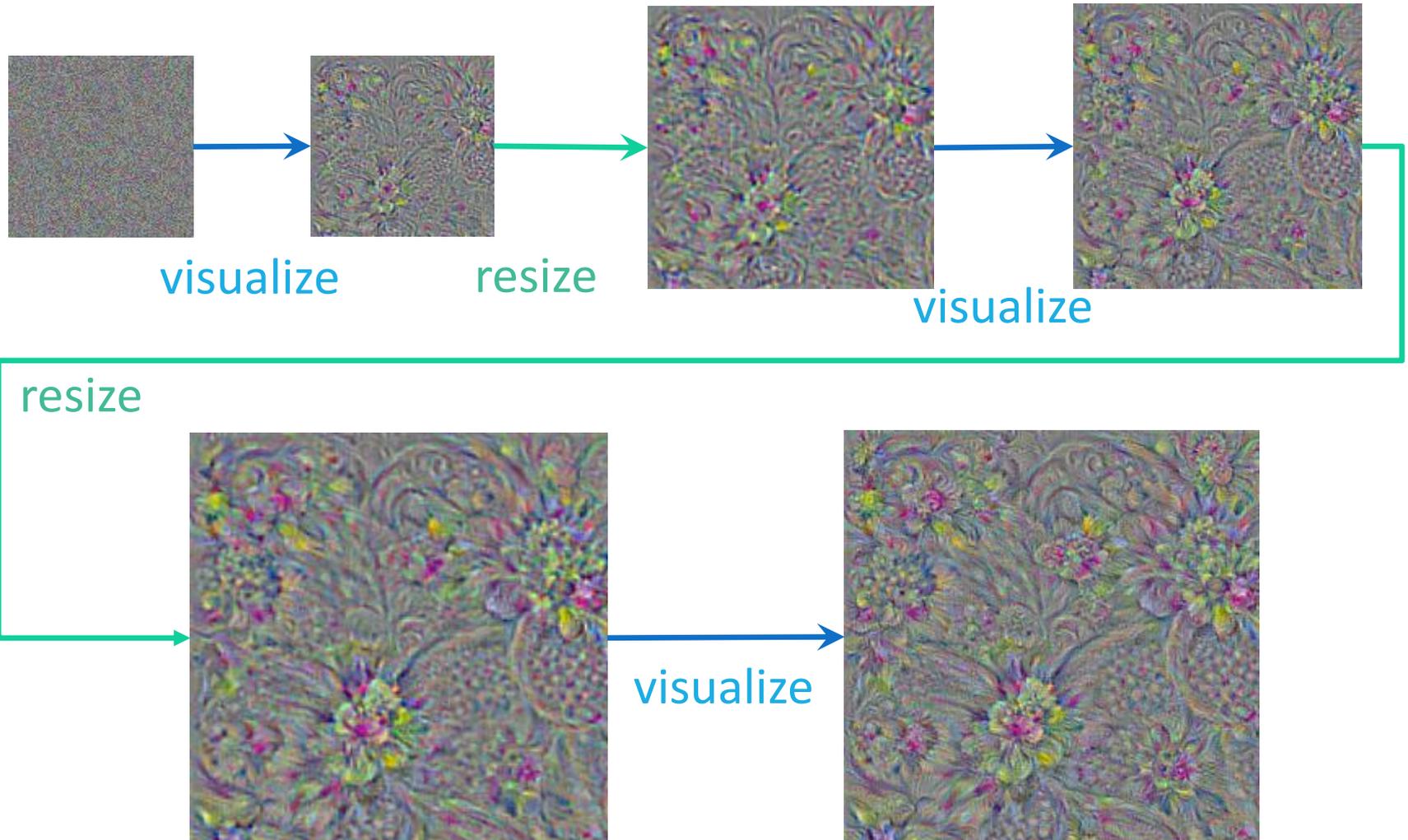
Different Layers of Visualization

147



Multiscale Image Generation

148

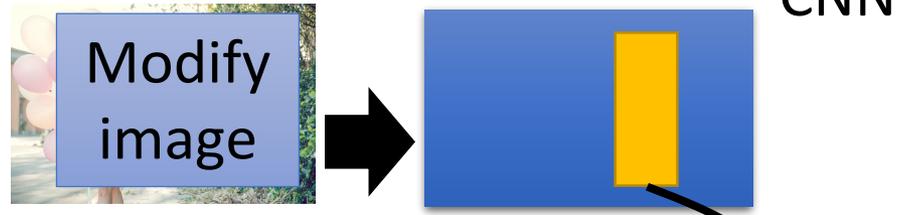


Multiscale Image Generation

149



Deep Dream

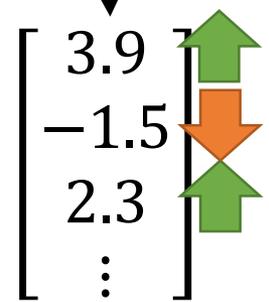


150

- Given a photo, machine adds what it sees



CNN exaggerates what it sees



<http://deepdreamgenerator.com/>

Deep Dream

<http://deepdreamgenerator.com/>

151

- Given a photo, machine adds what it sees



Deep Style

<http://deepdreamgenerator.com/>

152

- Given a photo, make its style like famous paintings



Deep Style

<http://deepdreamgenerator.com/>

153

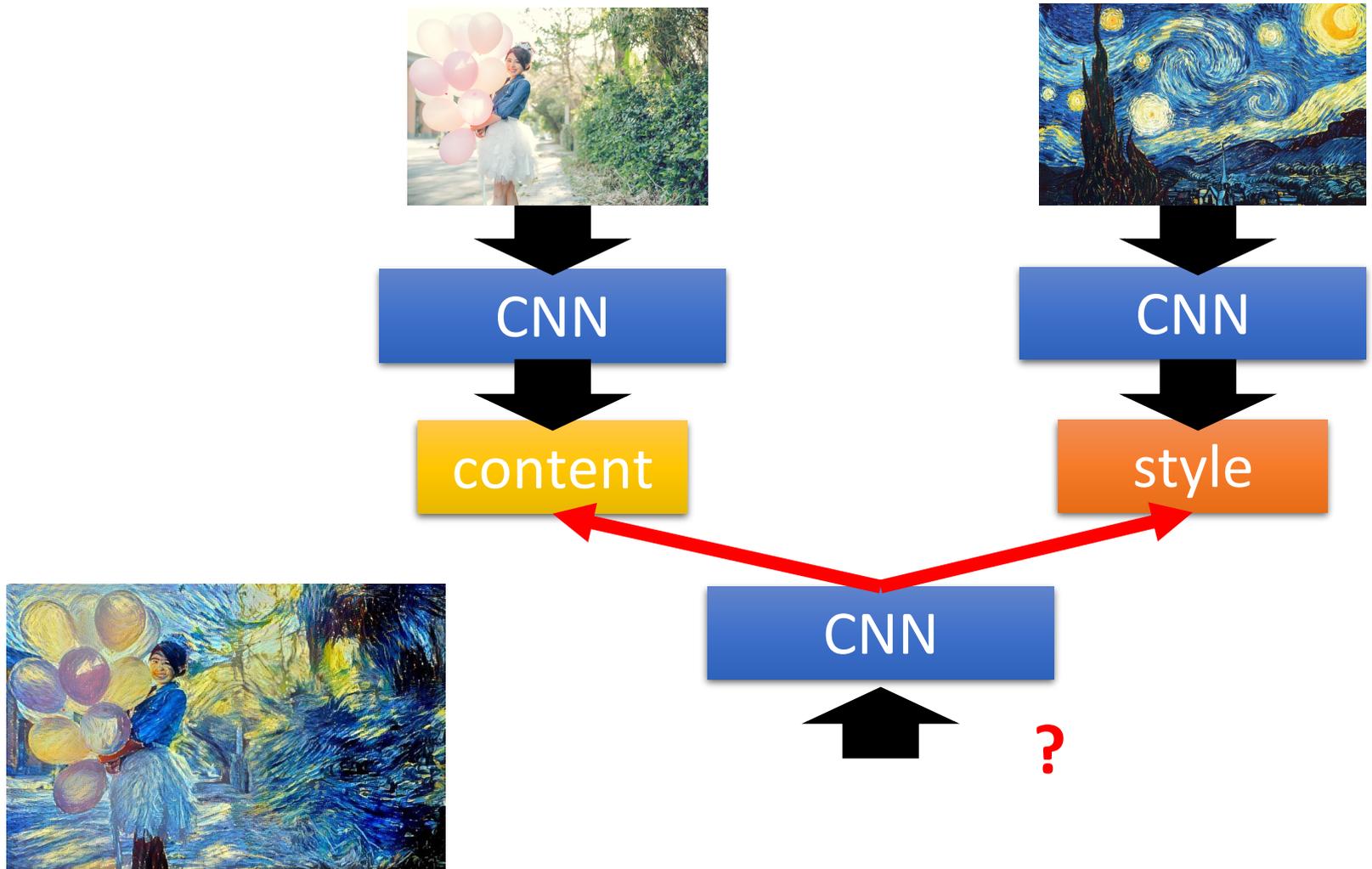
- Given a photo, make its style like famous paintings



Deep Style

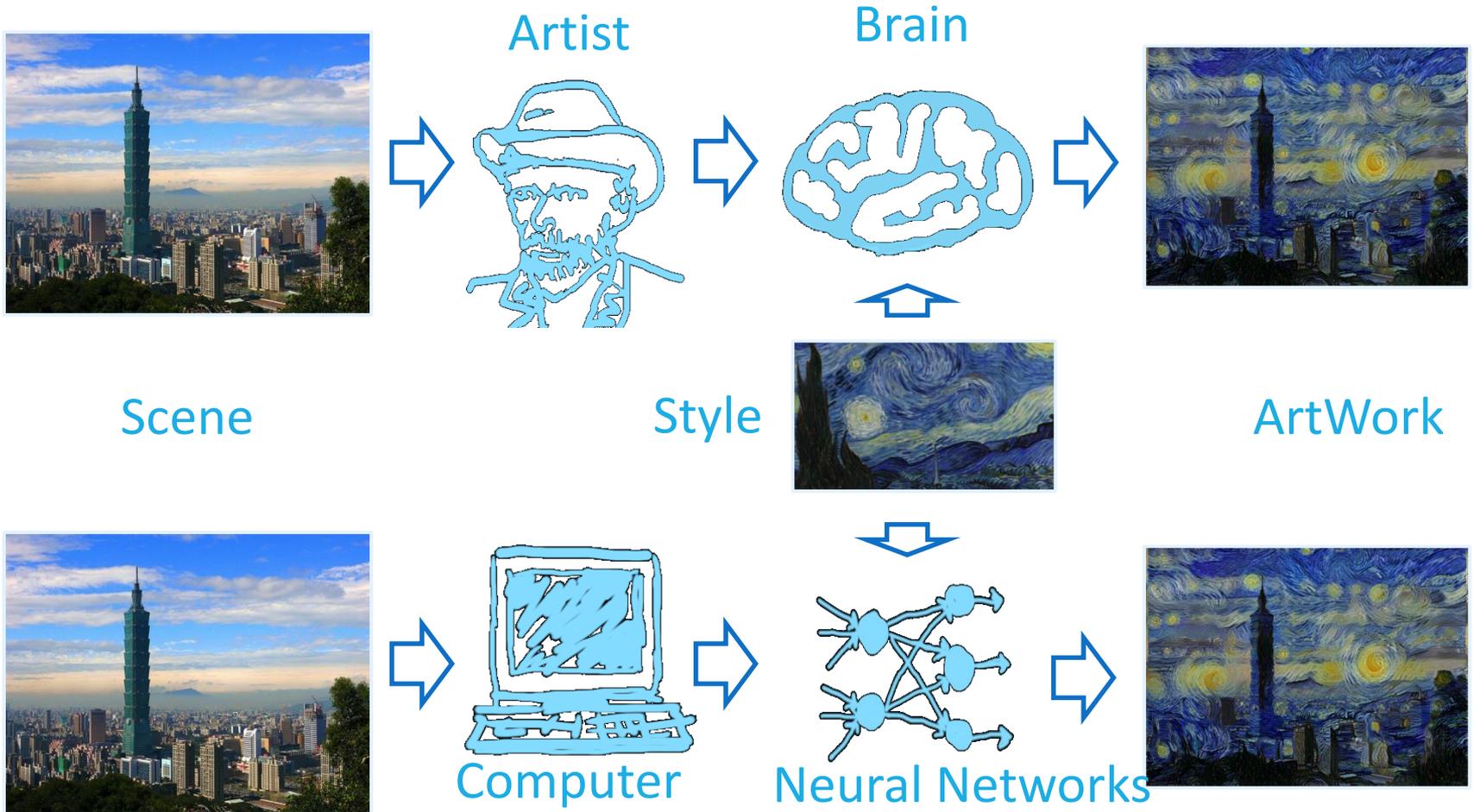
A Neural Algorithm of Artistic Style
<https://arxiv.org/abs/1508.06576>

154



Neural Art Mechanism

155



Go Playing

156



Black: 1
white: -1
none: 0



Network



Next move
(19 x 19
positions)

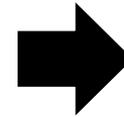
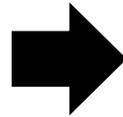
19 x 19 vector

Fully-connected feedforward
network can be used

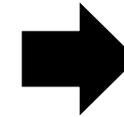
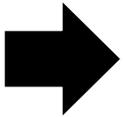
But CNN performs much better.

More Application: Playing Go

Training: record of previous plays 黑: 5之五 → 白: 天元 → 黑: 五之5 ...



Target:
“天元” = 1
else = 0



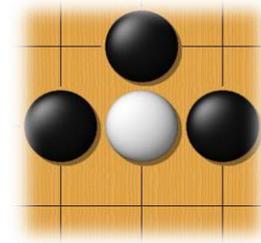
Target:
“五之5” = 1
else = 0

Why CNN for playing Go?

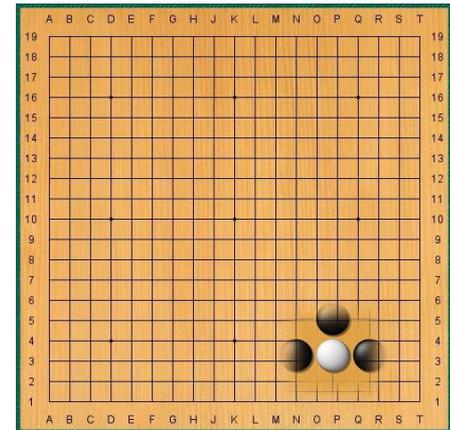
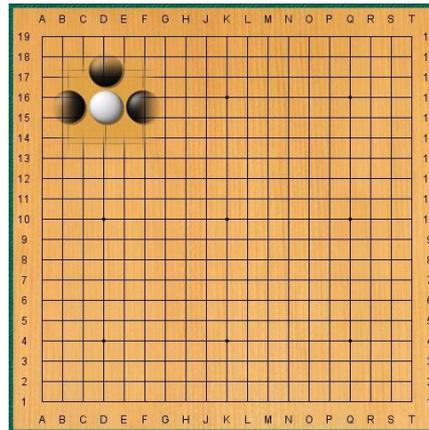
158

- Some patterns are much smaller than the whole image

Alpha Go uses 5 x 5 for first layer



- The same patterns appear in different regions



Why CNN for playing Go?

159

- Subsampling the pixels will not change the object



Max Pooling

How to explain this???

Neural network architecture. The input to the policy network is a $19 \times 19 \times 48$ image stack consisting of 48 feature planes. The first hidden layer zero pads the input into a 23×23 image, then convolves k filters of kernel size 5×5 with stride 1 with the input image and applies a rectifier nonlinearity. Each of the subsequent hidden layers 2 to 12 zero pads the respective previous hidden layer into a 21×21 image, then convolves k filters of kernel size 3×3 with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size 1×1 with stride 1, with a different bias for each position, and applies a softmax function. The **Alpha Go does not use Max Pooling** Extended Data Table 3 additionally show the results of training with $k = 128, 256$ and 384 filters.

PART II: Variants of Neural Networks

160

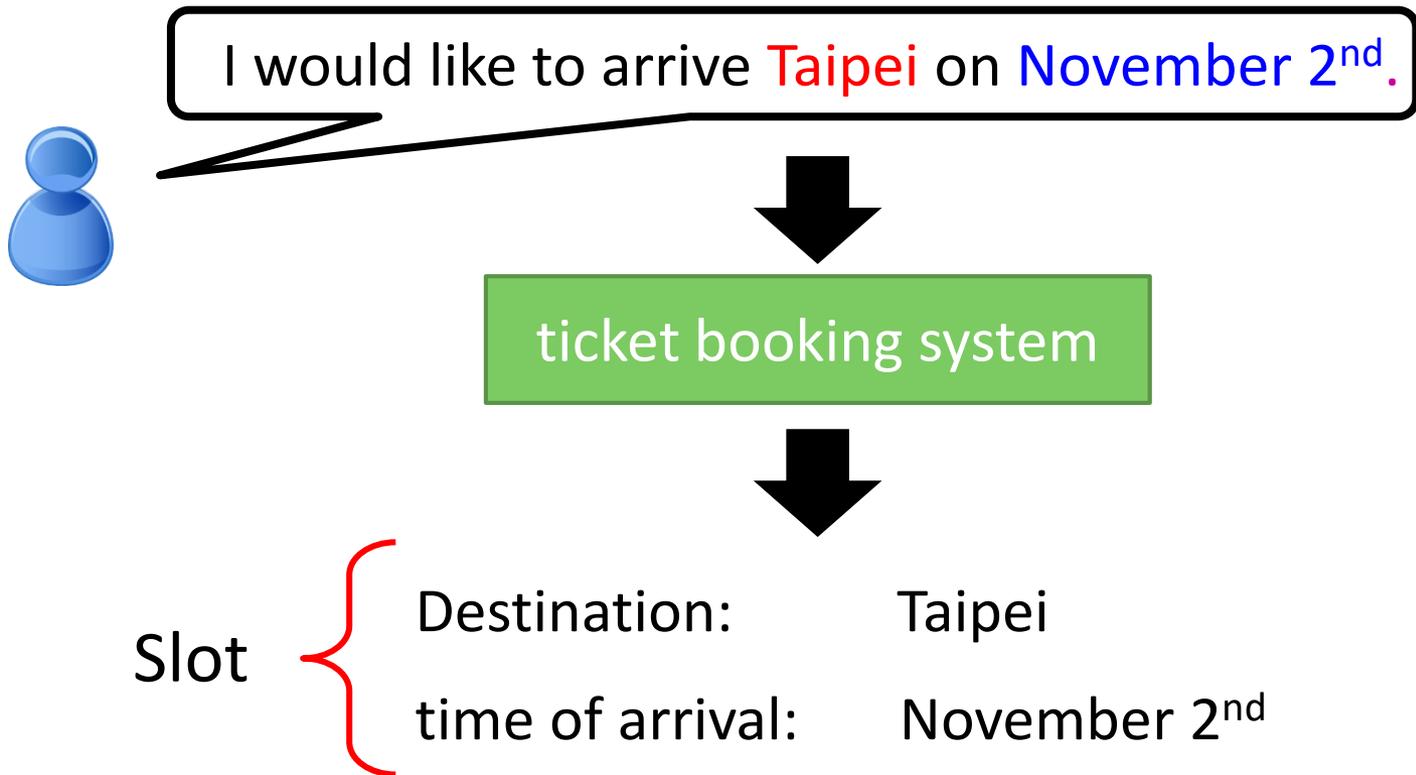
- Convolutional Neural Network (CNN)
- **Recurrent Neural Network (RNN)**

Neural Network with Memory

Example Application

161

□ Slot Filling



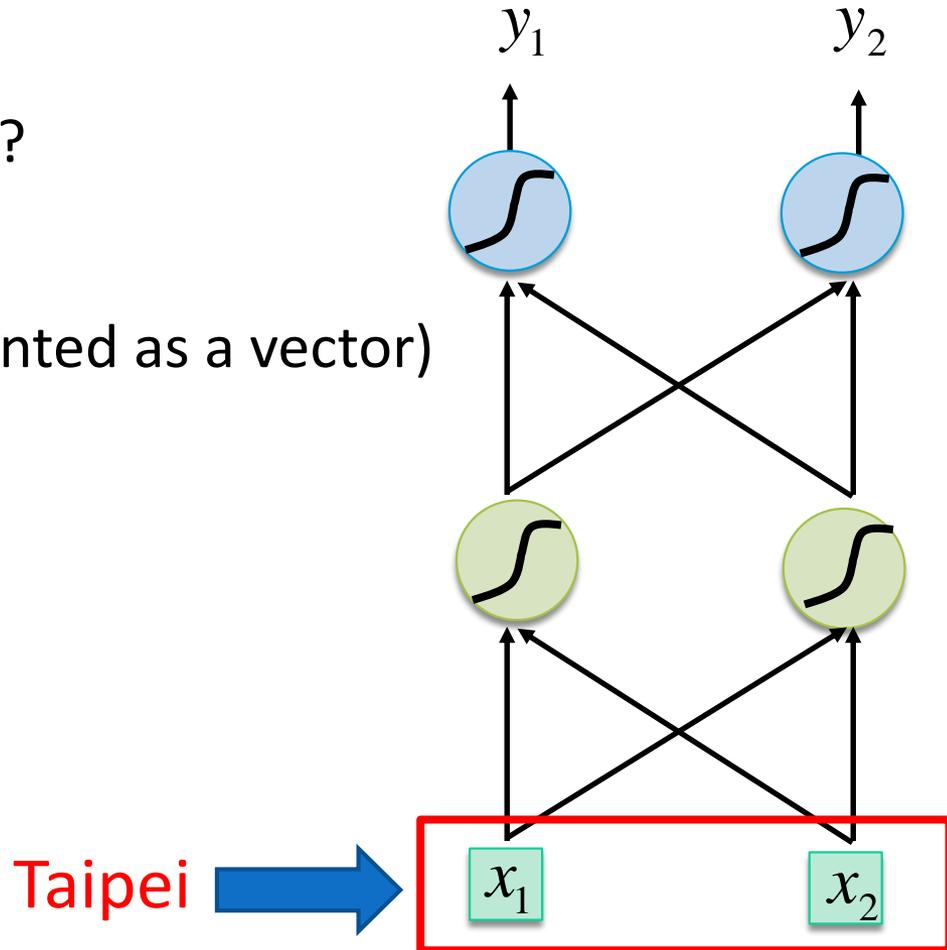
Example Application

162

Solving slot filling by
feedforward network?

Input: a word

(Each word is represented as a vector)



1-of-N encoding

163

How to represent each word as a vector?

1-of-N Encoding lexicon = {apple, bag, cat, dog, elephant}

The vector is lexicon size.

$$\text{apple} = [1 \ 0 \ 0 \ 0 \ 0]$$

Each dimension corresponds
to a word in the lexicon

$$\text{bag} = [0 \ 1 \ 0 \ 0 \ 0]$$

$$\text{cat} = [0 \ 0 \ 1 \ 0 \ 0]$$

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

$$\text{dog} = [0 \ 0 \ 0 \ 1 \ 0]$$

$$\text{elephant} = [0 \ 0 \ 0 \ 0 \ 1]$$

Example Application

164

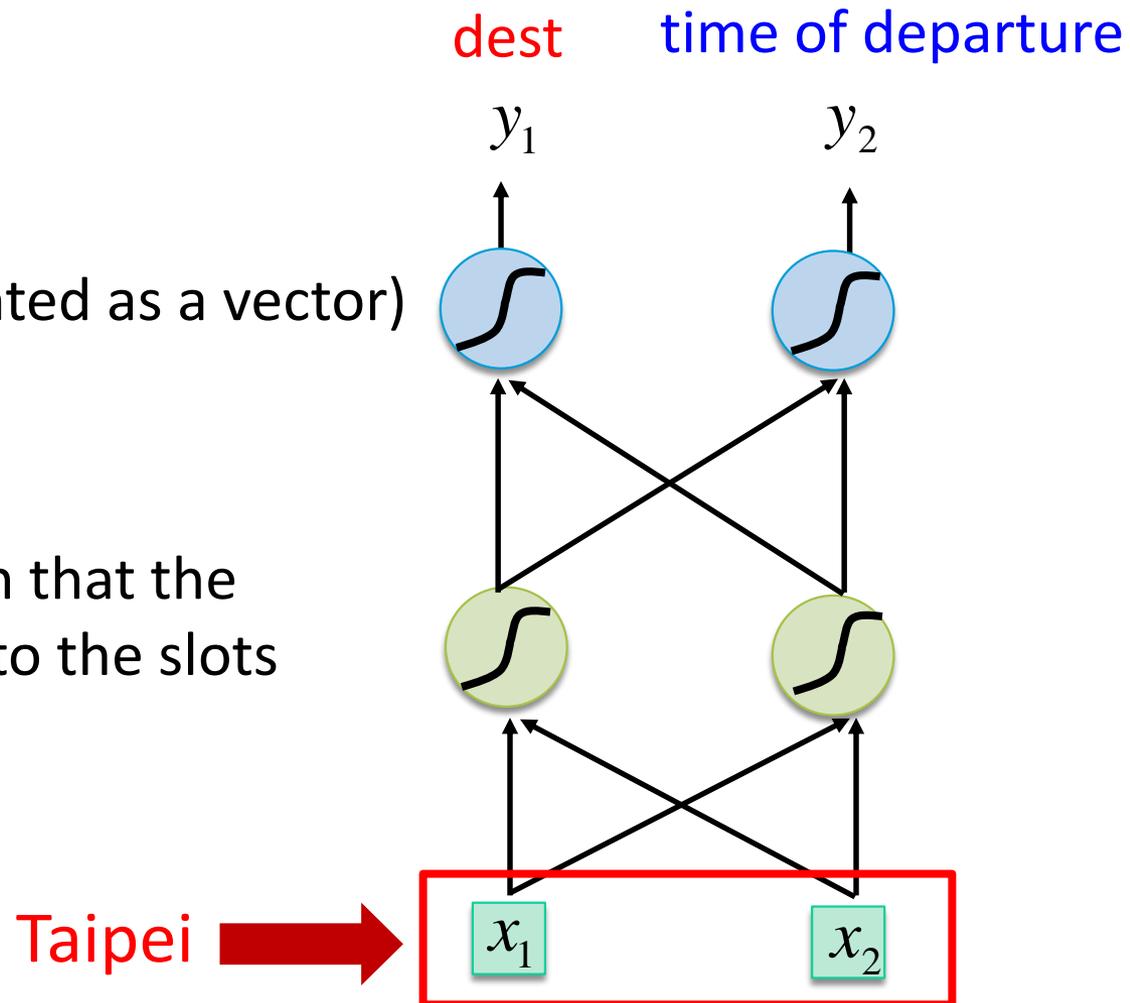
Solving slot filling by
feedforward network?

Input: a word

(Each word is represented as a vector)

Output:

probability distribution that the
input word belonging to the slots



Example Application

165

arrive Taipei on November 2nd

other dest other time time

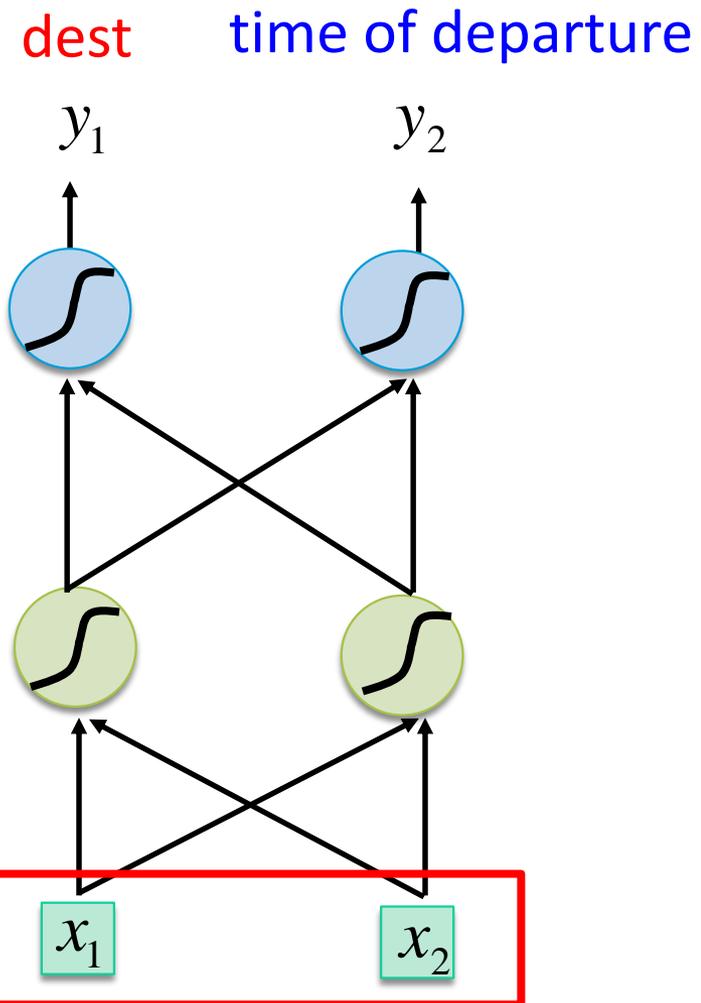
Problem?

leave Taipei on November 2nd

place of departure

Neural network needs memory!

Taipei →

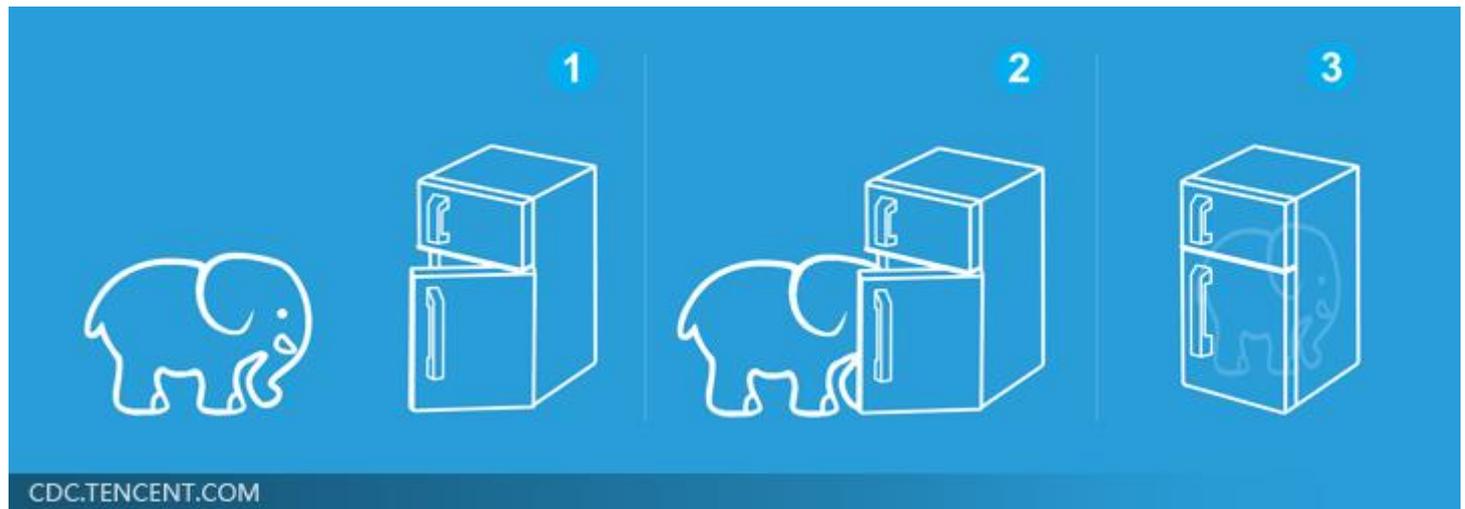


Three Steps for Deep Learning

166



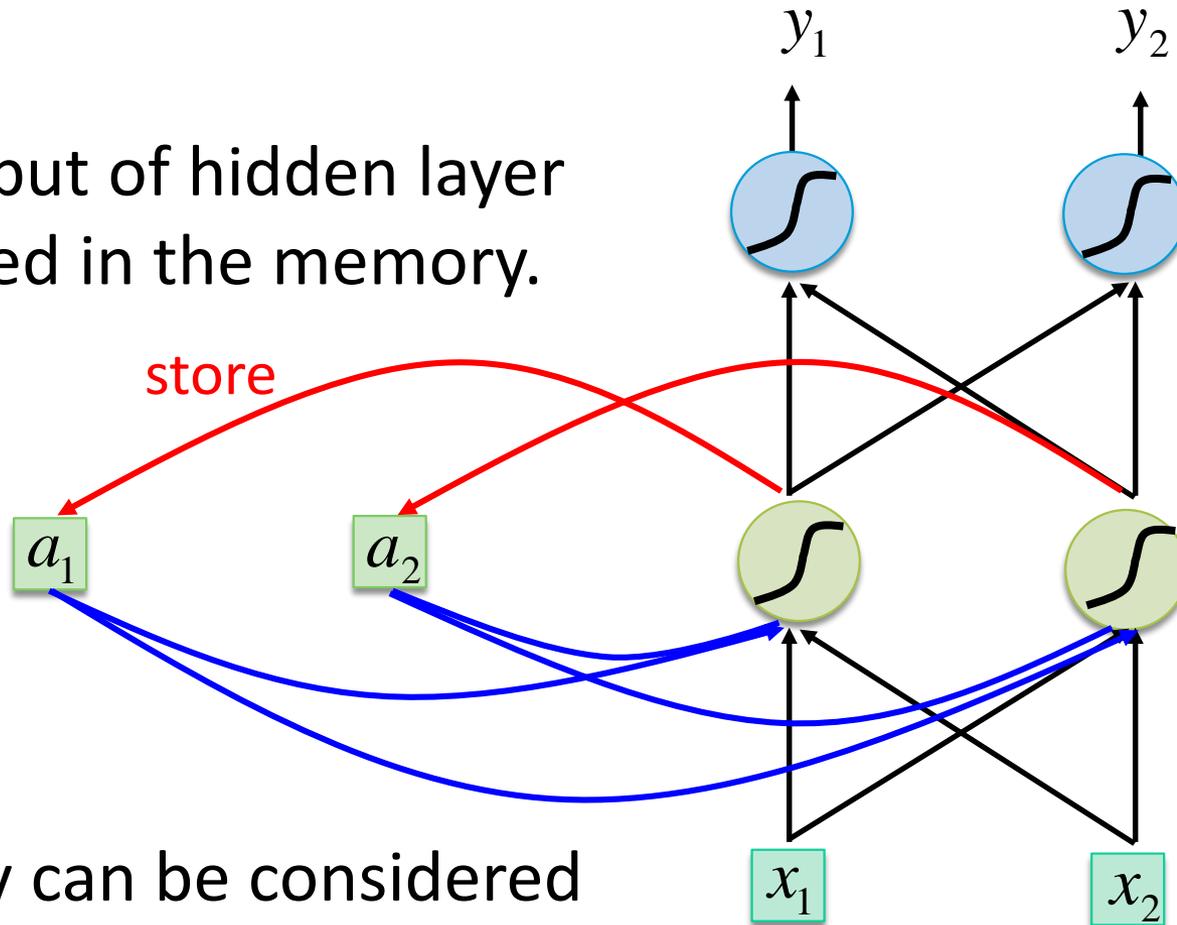
Deep Learning is so simple



Recurrent Neural Network (RNN)

167

The output of hidden layer are stored in the memory.



Memory can be considered as another input.

RNN

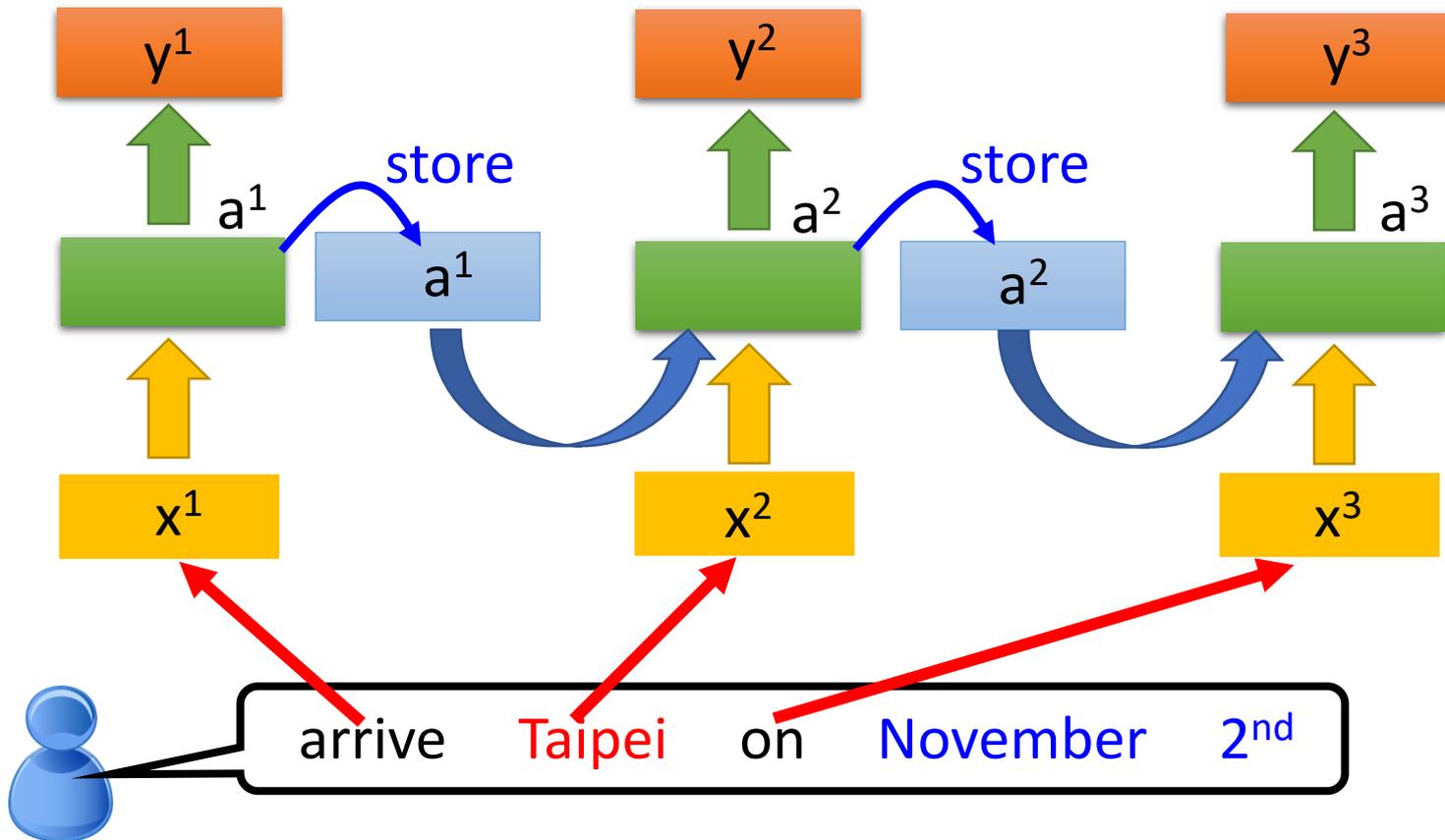
The same network is used again and again.

168

Probability of
“arrive” in each slot

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

Probability of
“on” in each slot



RNN

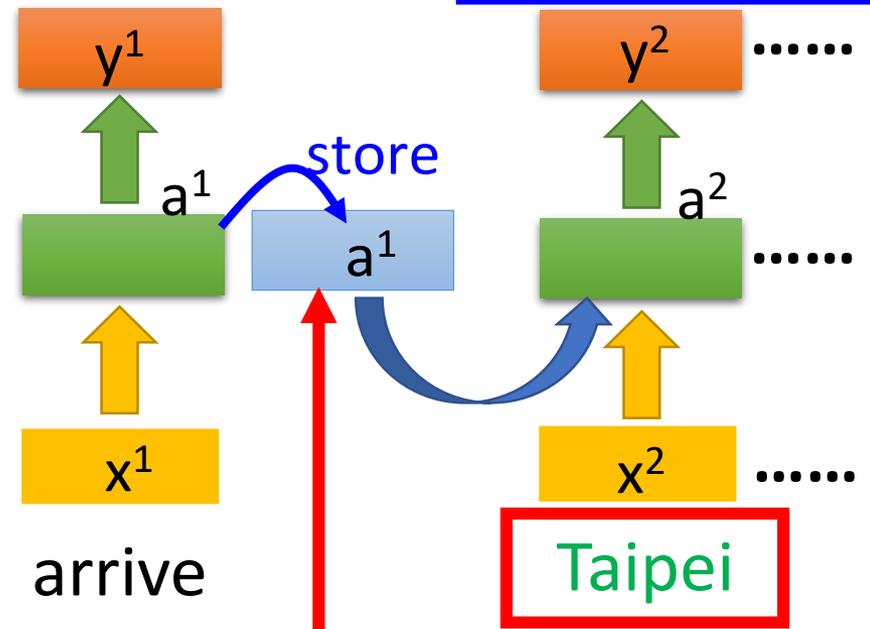
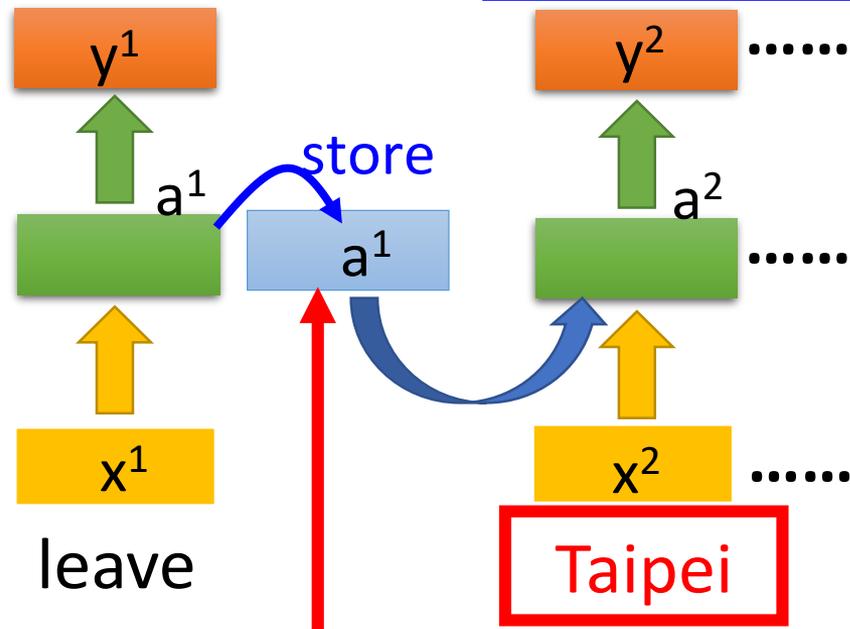
169

Different

Prob of "Taipei" in each slot

Prob of "arrive" in each slot

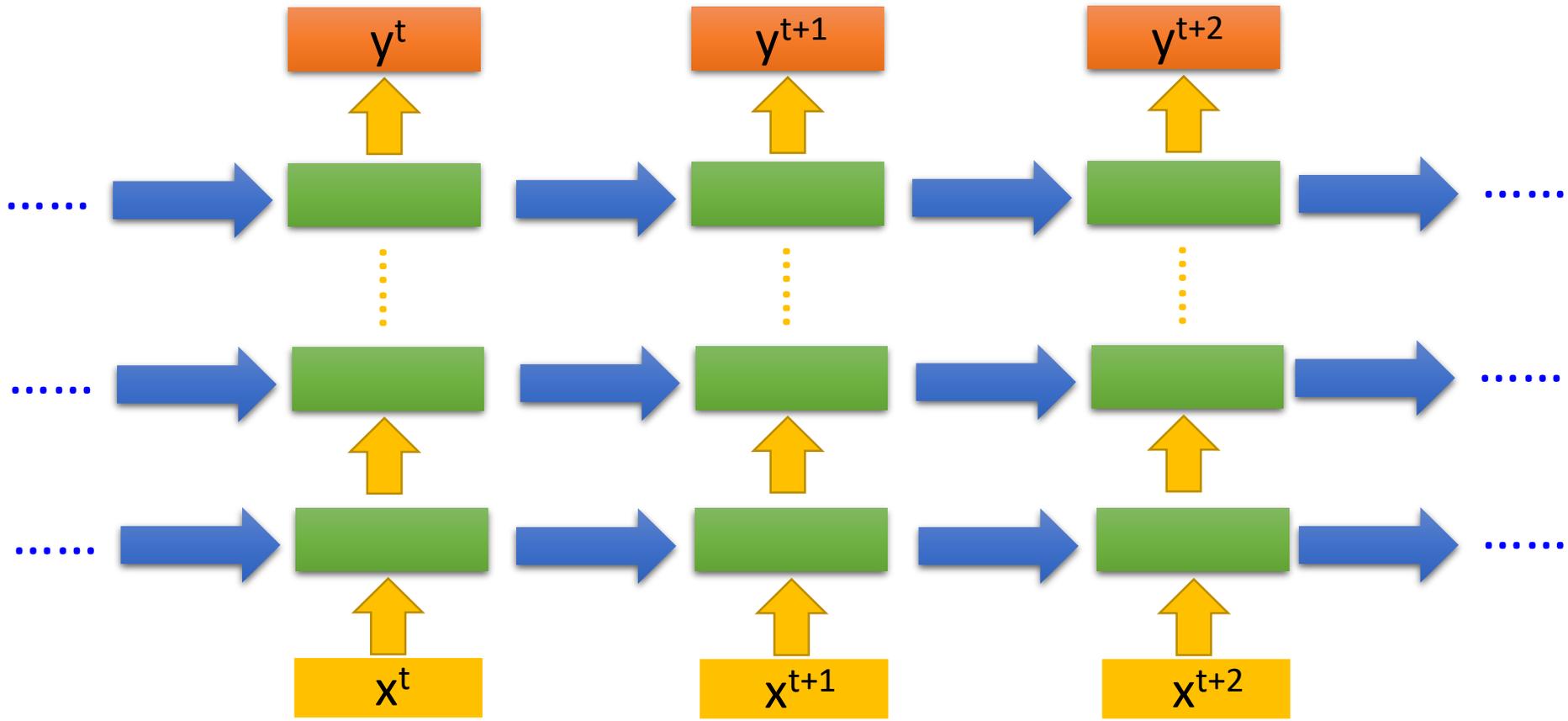
Prob of "Taipei" in each slot



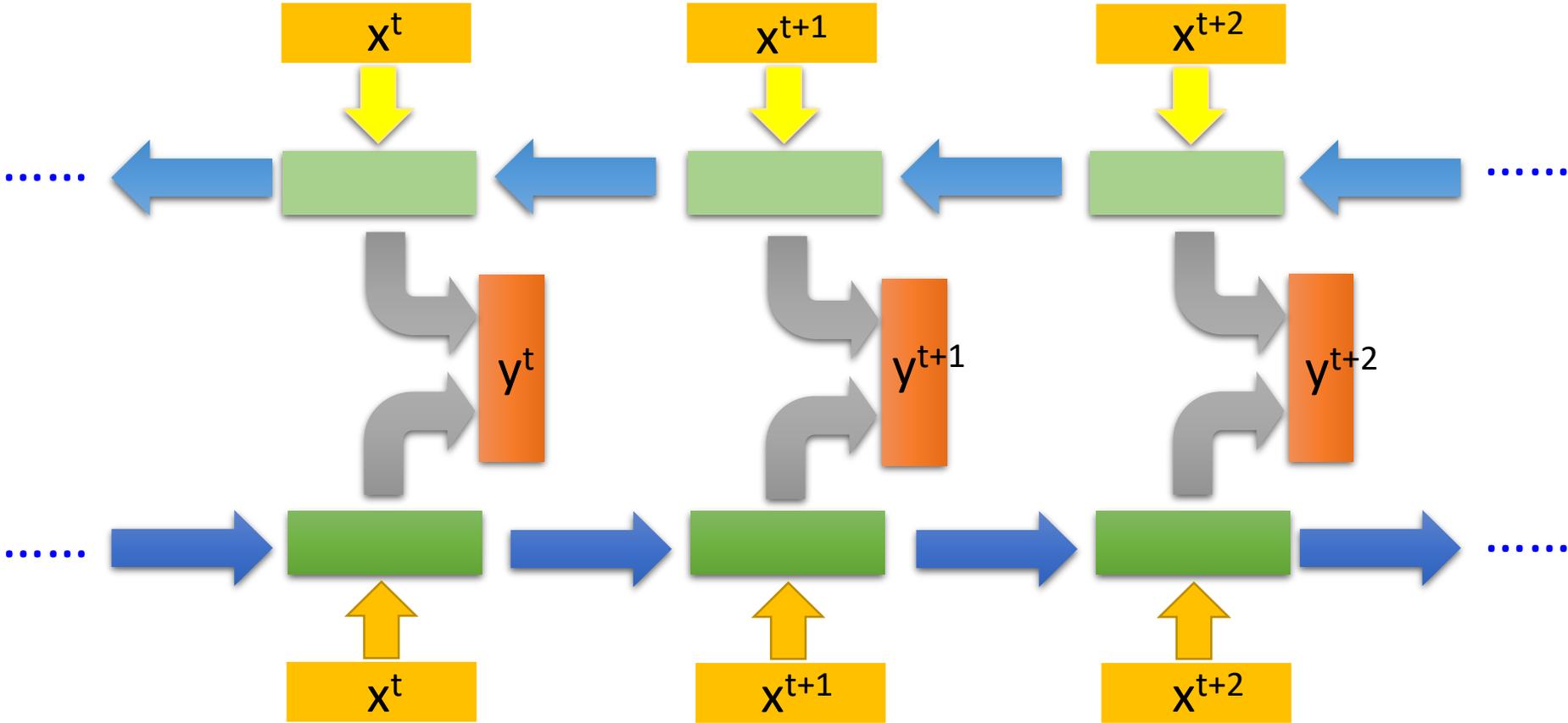
The values stored in the memory is different.

Deep RNN

170



Bidirectional RNN

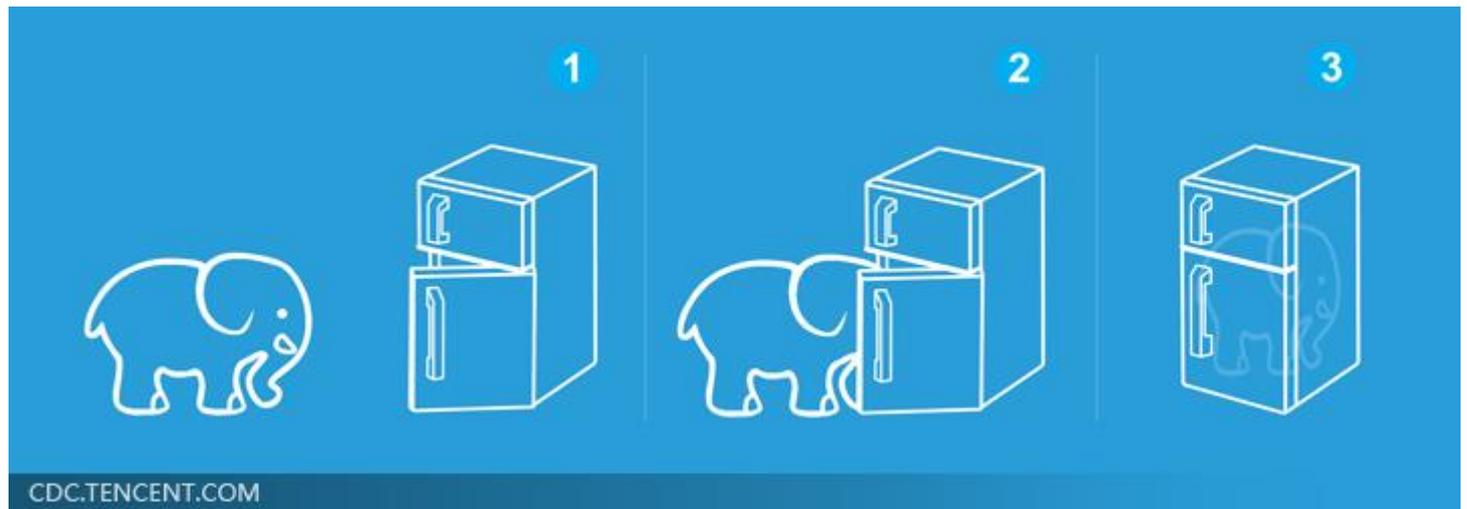


RNN

172

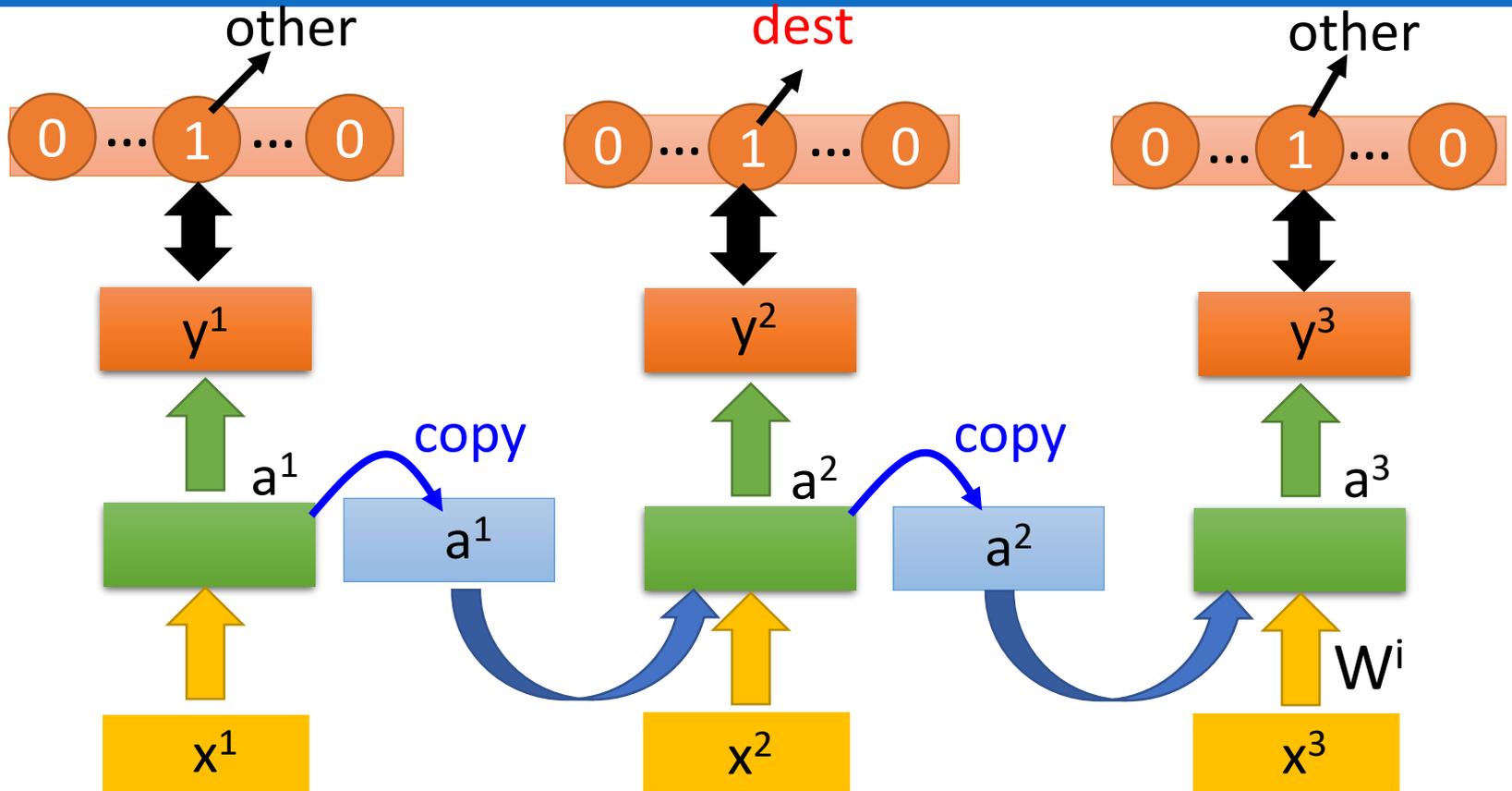


Deep Learning is so simple



Learning Target

173

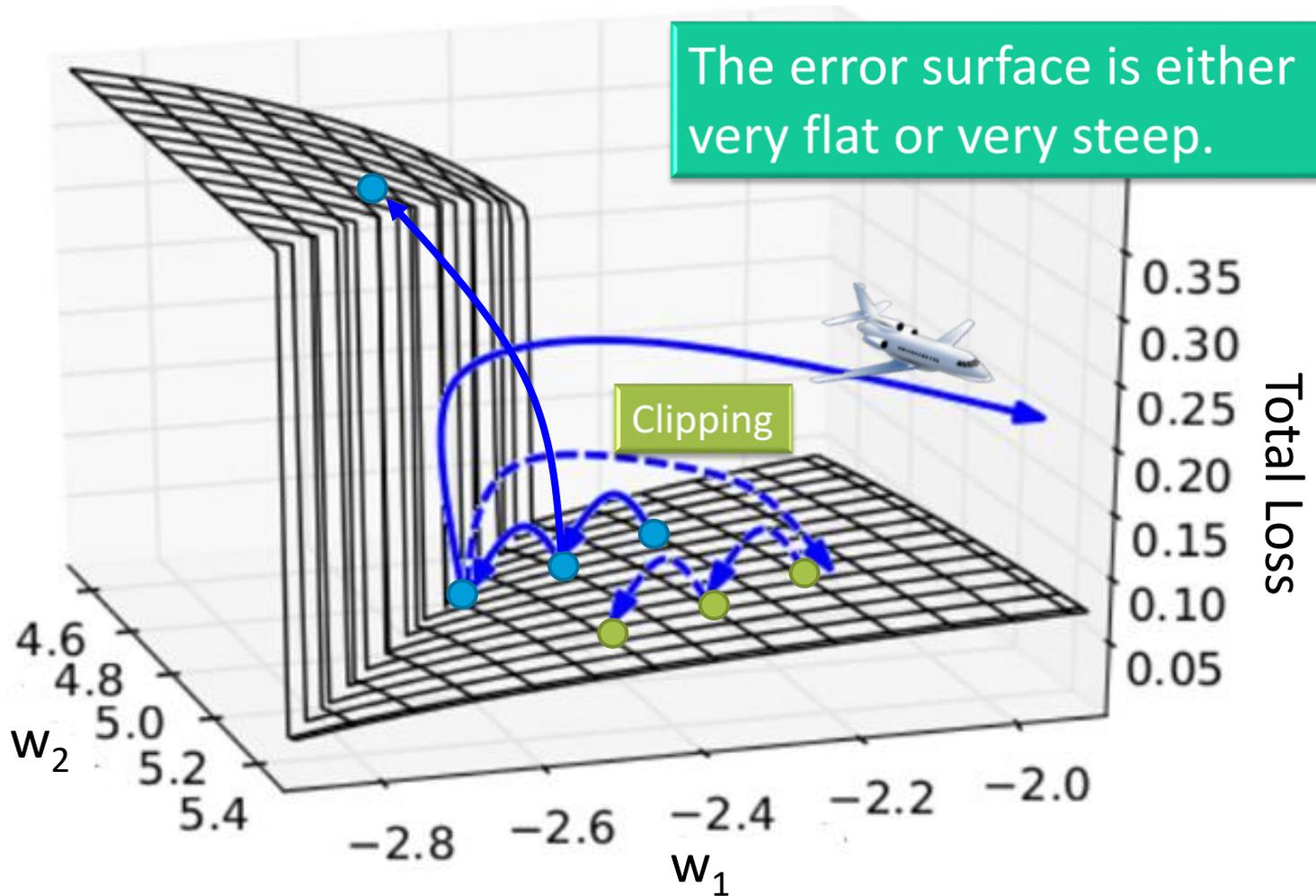


Training
Sentences:

arrive Taipei on November 2nd
other dest other time time

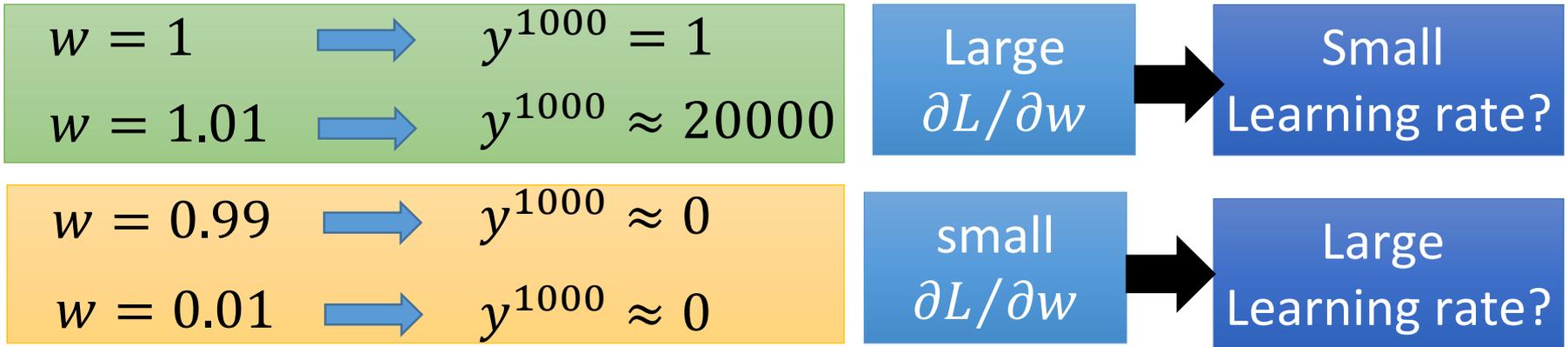
Rough Error Surface

174

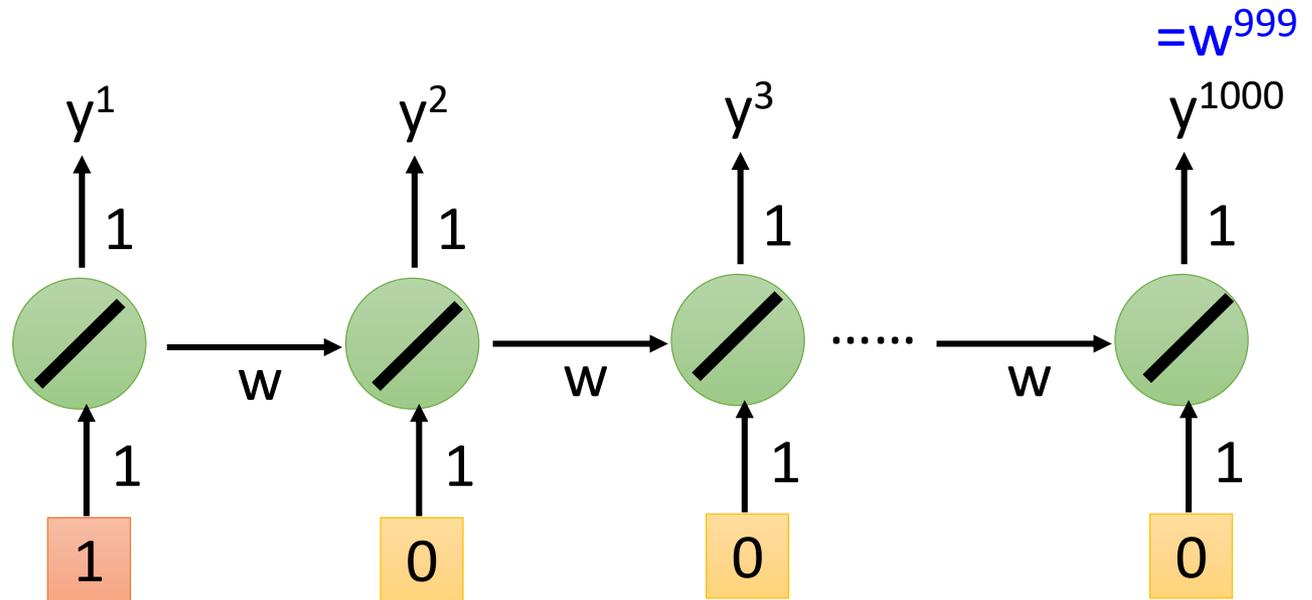


Rough Error Surface

175



Toy Example

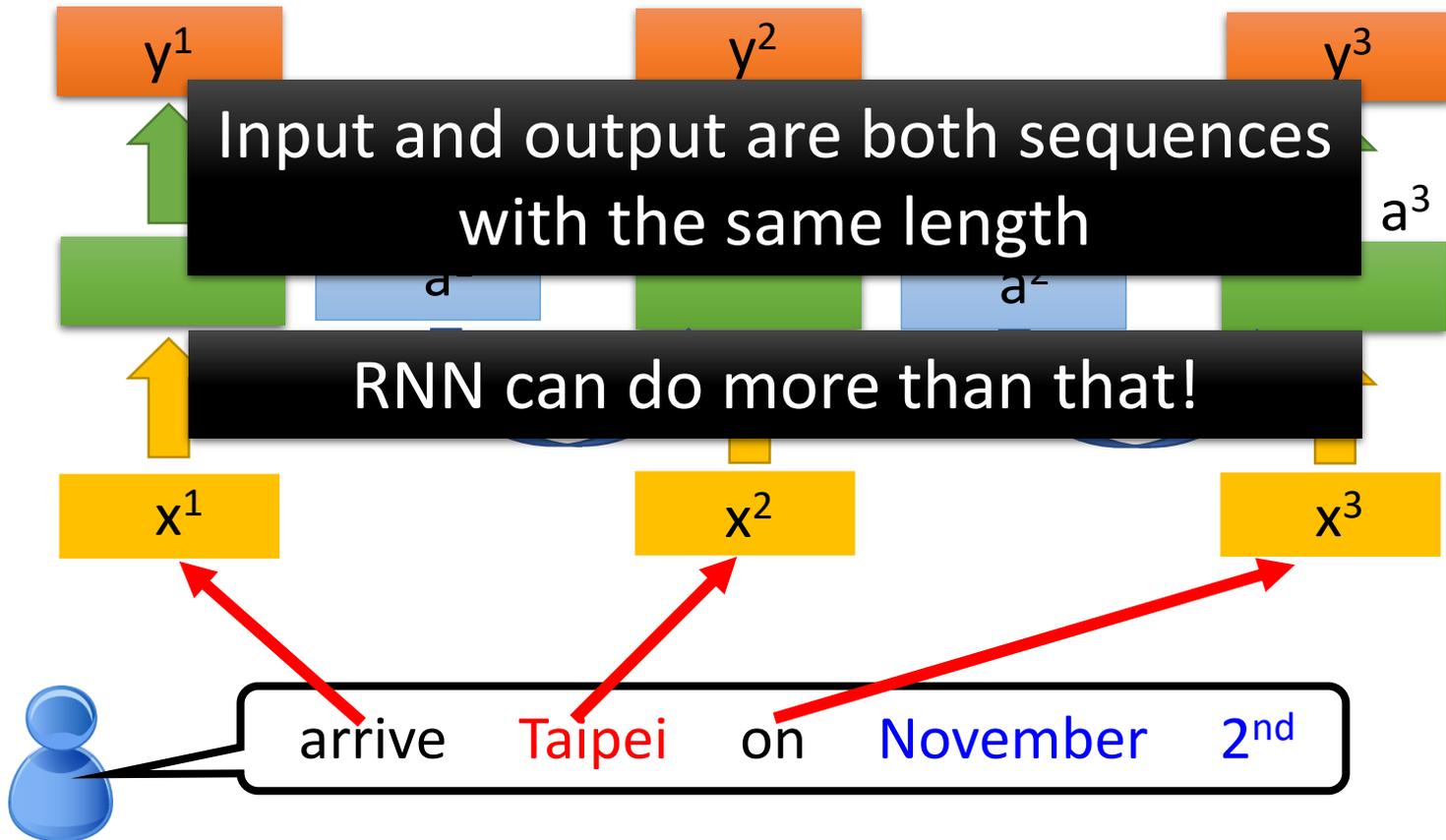


RNN Applications

Probability of
“arrive” in each slot

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

Probability of
“on” in each slot



Many-to-One

177

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

Sentiment Analysis

看了這部電影覺得很高興

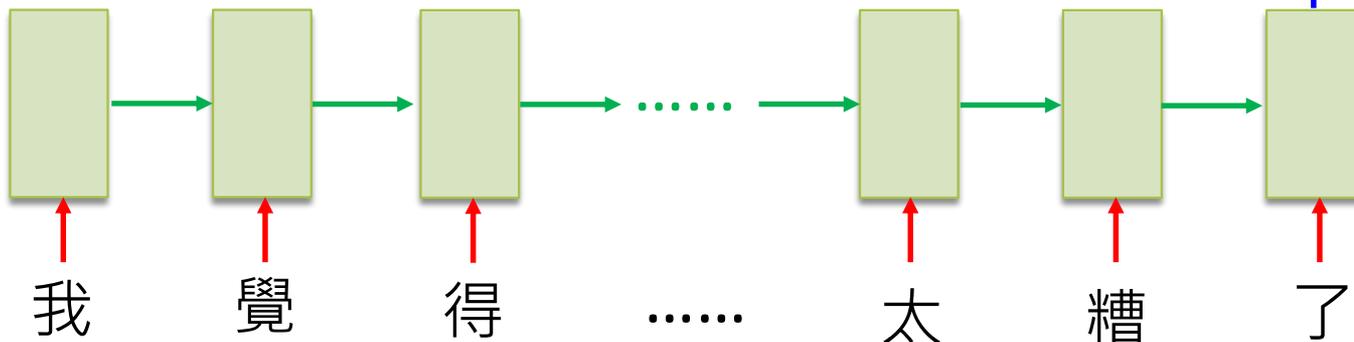
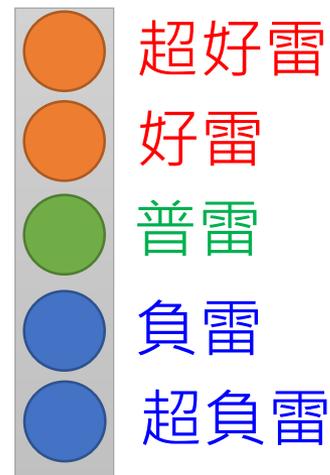
Positive (正雷)

這部電影太糟了

Negative (負雷)

這部電影很棒

Positive (正雷)



Many-to-Many (Output is shorter)

178

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

Problem?

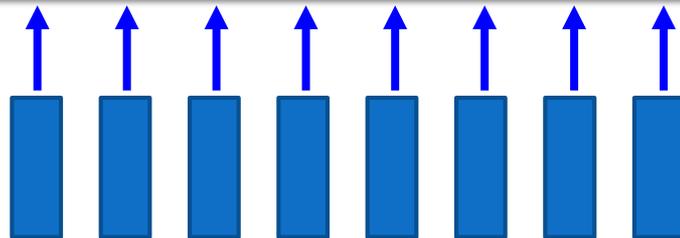
Why can't it be
“好棒棒”

Output: “好棒” (character sequence)



Trimming

好 好 好 棒 棒 棒 棒 棒



Input:

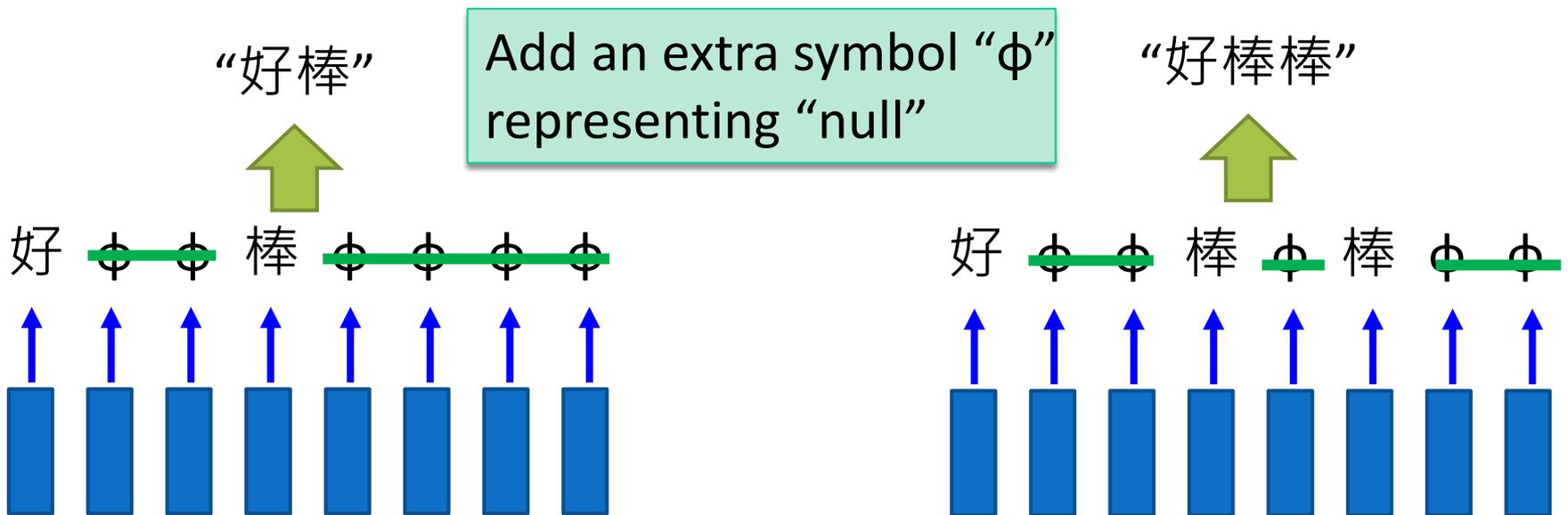
(vector
sequence)



Many-to-Many (Output is shorter)

179

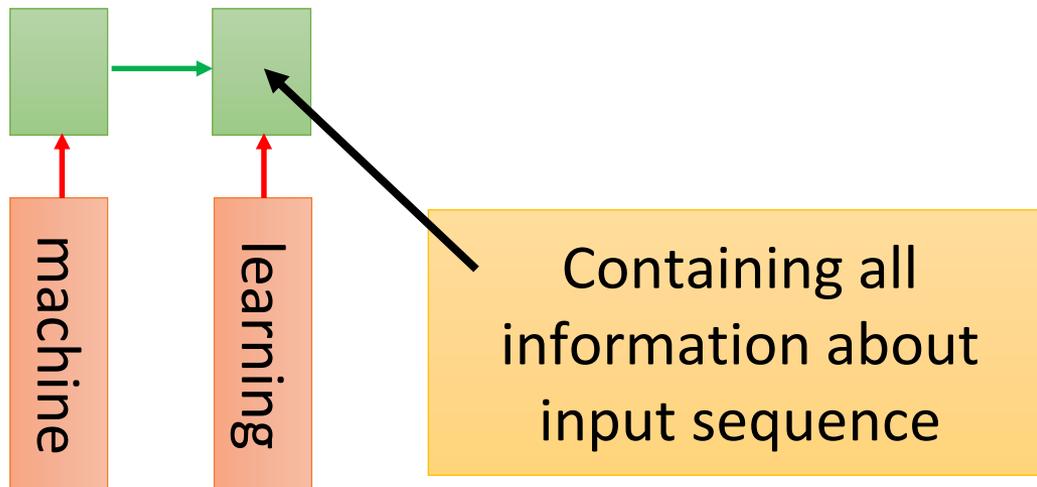
- Both input and output are both sequences, **but the output is shorter.**
- ▣ Connectionist Temporal Classification (CTC)



Many-to-Many (Output has no limitation)

180

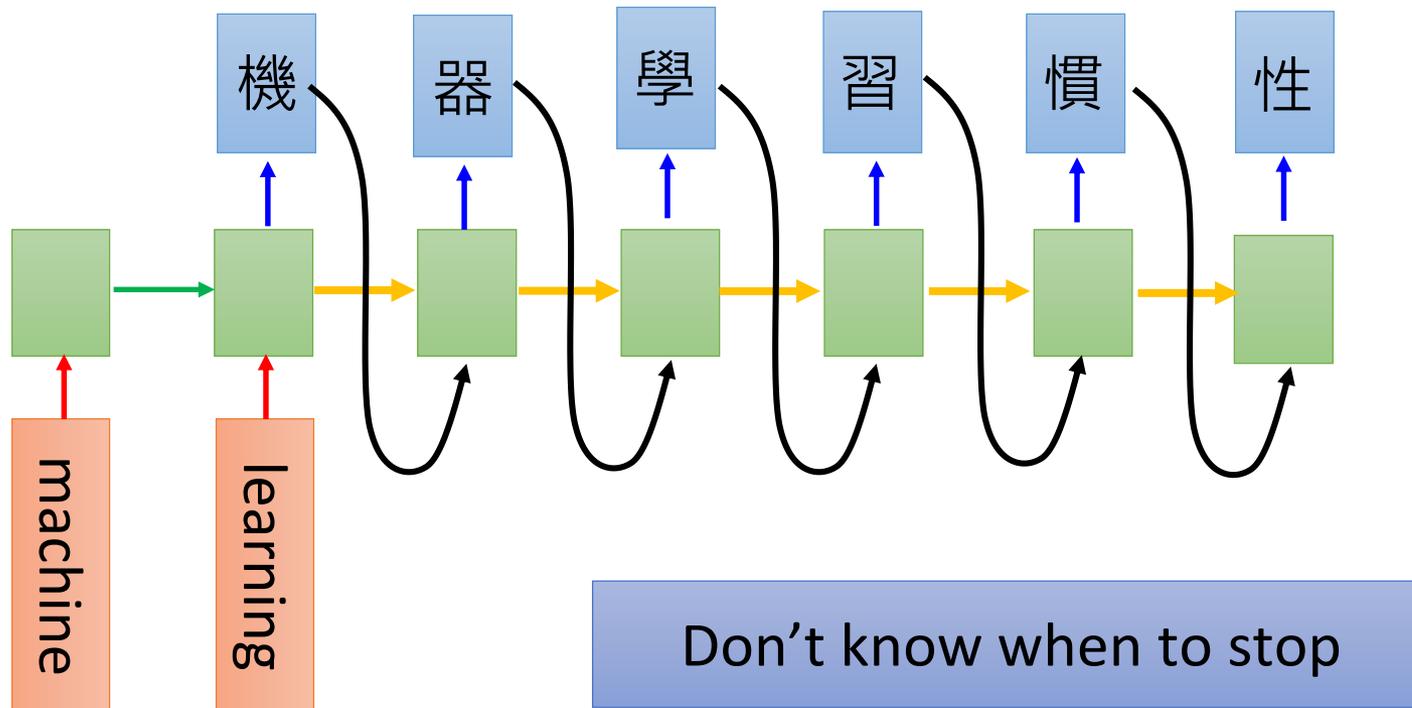
- Both input and output are both sequences *with different lengths.* → *Sequence to sequence learning*
 - E.g. *Machine Translation* (machine learning → 機器學習)



Many-to-Many (Output has no limitation)

181

- 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 (Output has no limitation)

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

推 tlkagk:  =====斷=====
```

Many-to-Many (Output has no limitation)

183

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

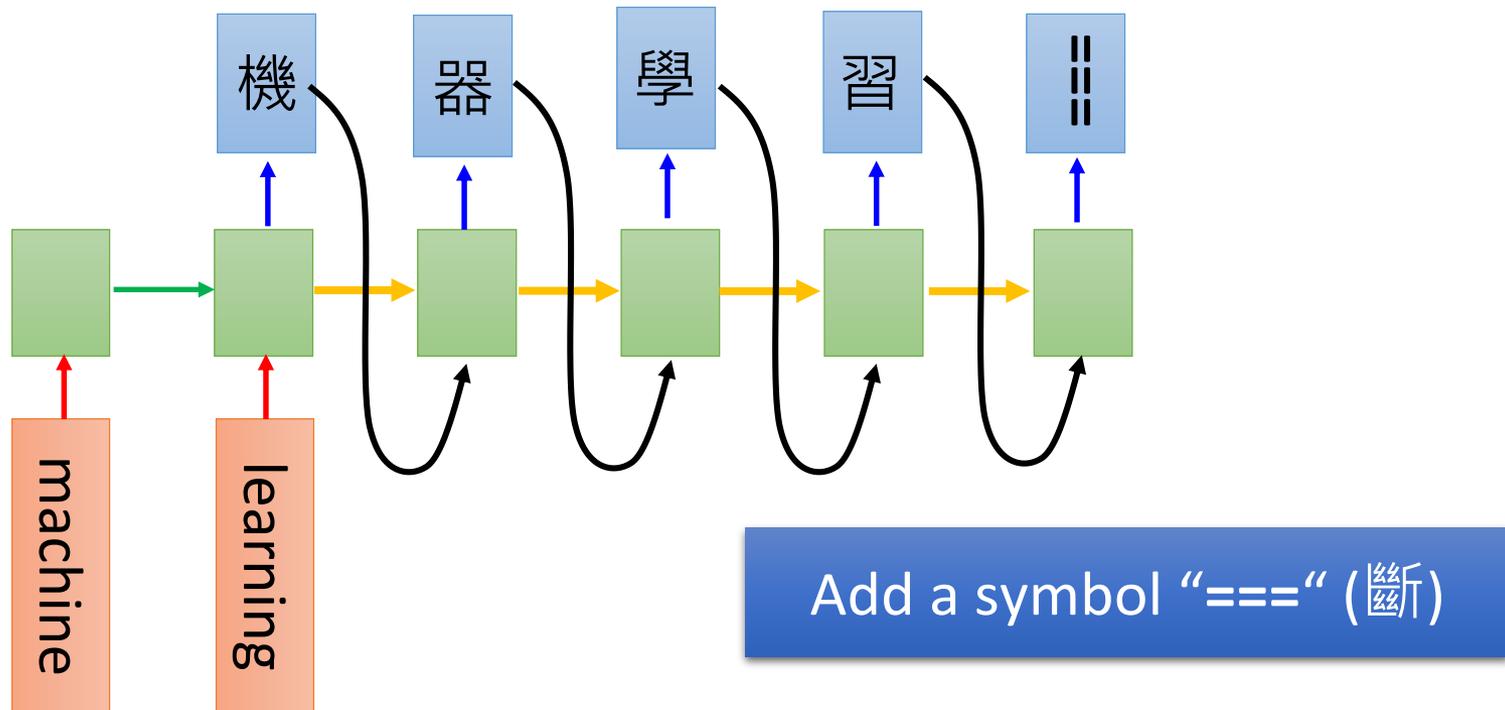


Image Caption Generation

- Input an image, but output a sequence of words

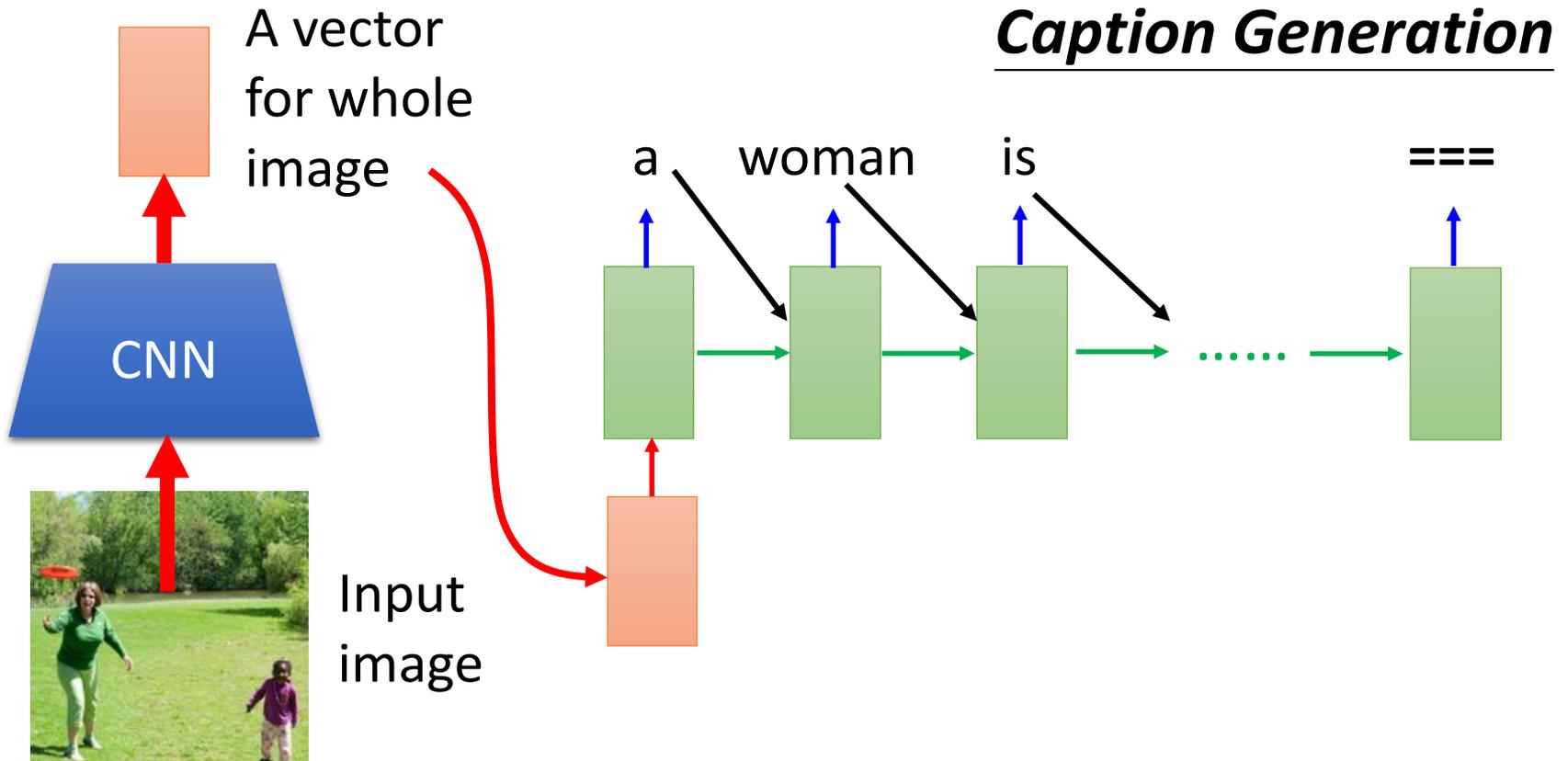


Image Caption Generation

185

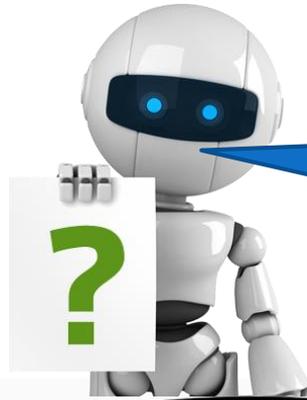


Video Caption Generation

186



Video



A girl is running.



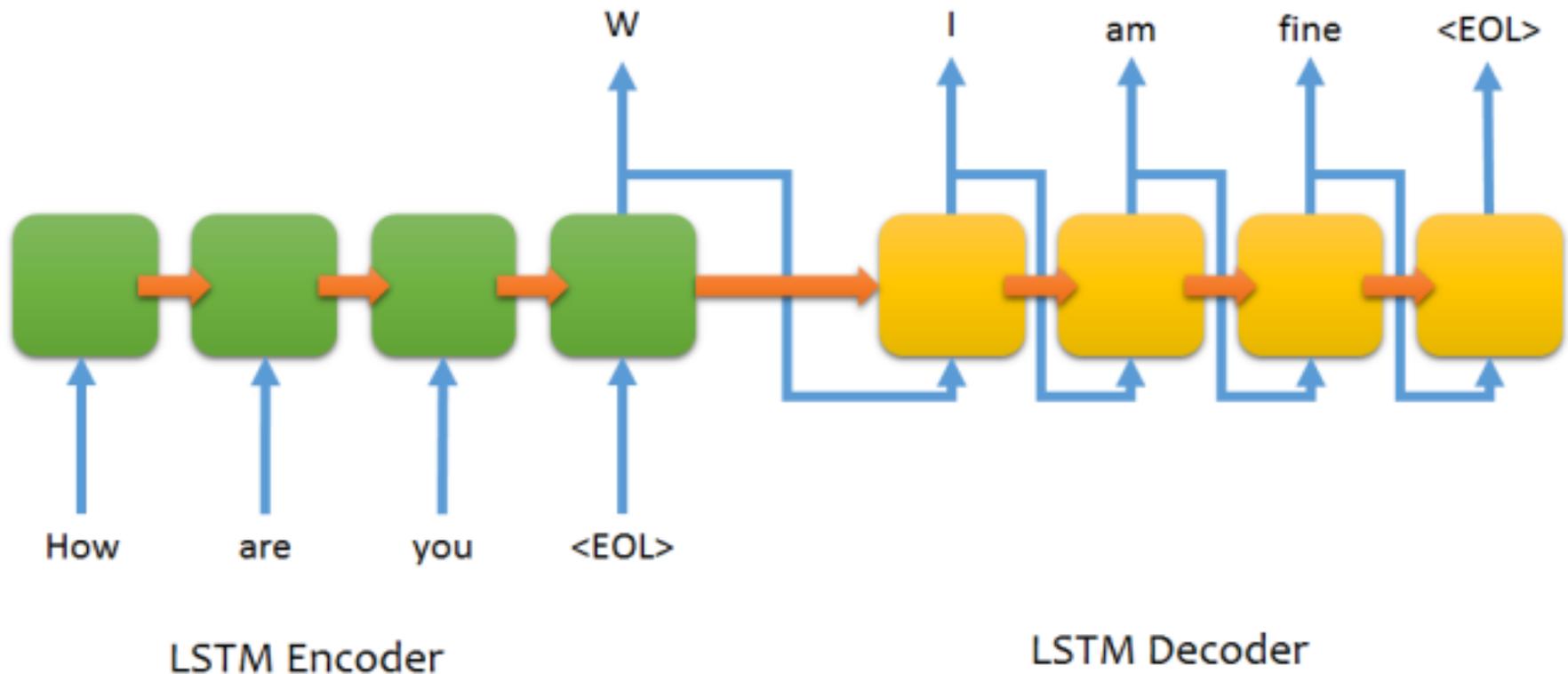
A group of people is knocked by a tree.



A group of people is walking in the forest.

Chit-Chat Bot

187



電視影集 (~40,000 sentences)、美國總統大選辯論

Sci-Fi Short Film - SUNSPRING

188

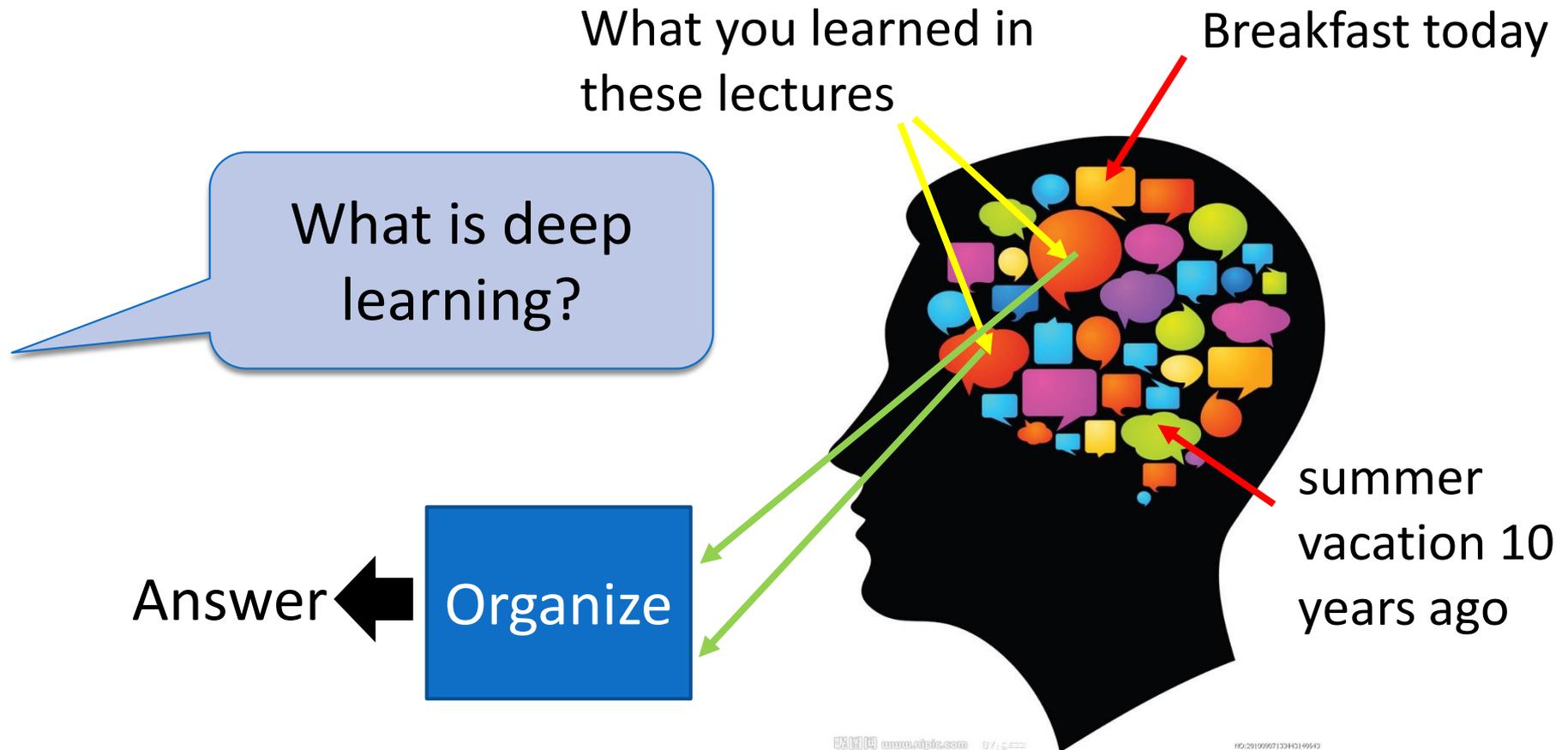
A close-up shot of a hand pulling a dark-colored drawer from a desk. The word "SUNSPRING" is printed in large, white, bold, sans-serif capital letters on the front of the drawer. On top of the desk, there are several items: a small brass vase, a blue book, and a white box. The background is a blurred office setting with a window.

SUNSPRING

<https://www.youtube.com/watch?v=LY7x2Ihqj>

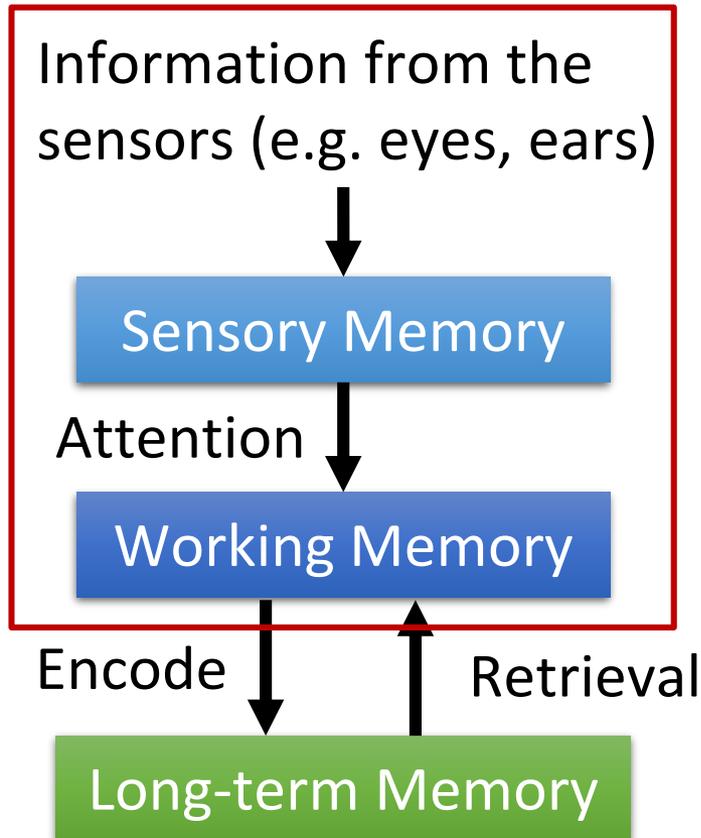
Attention and Memory

189



Attention on Sensory Info

190



When the input is a very long sequence or an image

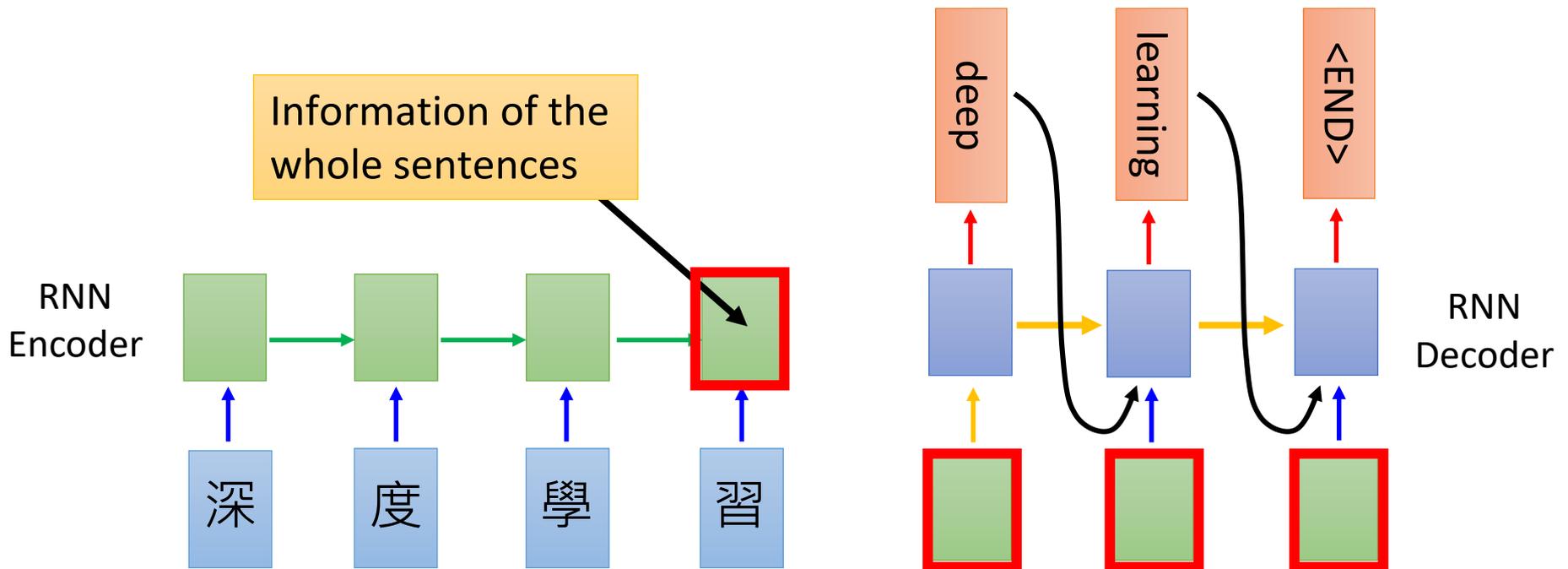


Pay attention on partial of the input object each time

Machine Translation

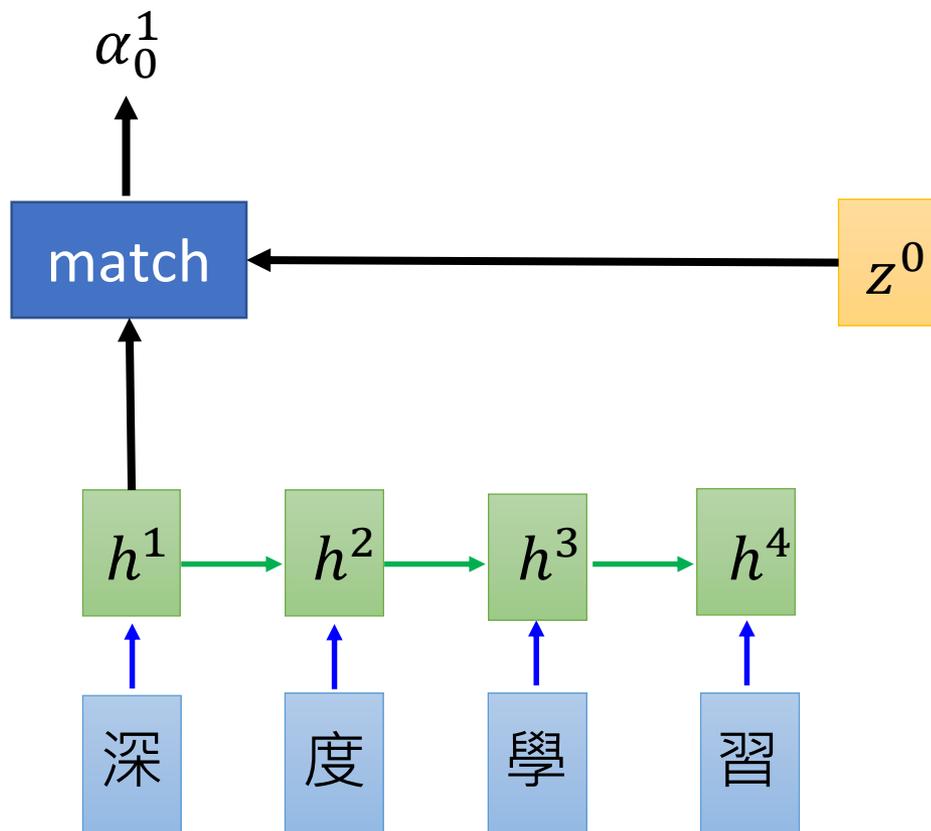
191

- Sequence-to-sequence learning: both input and output are both sequences *with different lengths*.
- E.g. 深度學習 → deep learning



Machine Translation with Attention

192



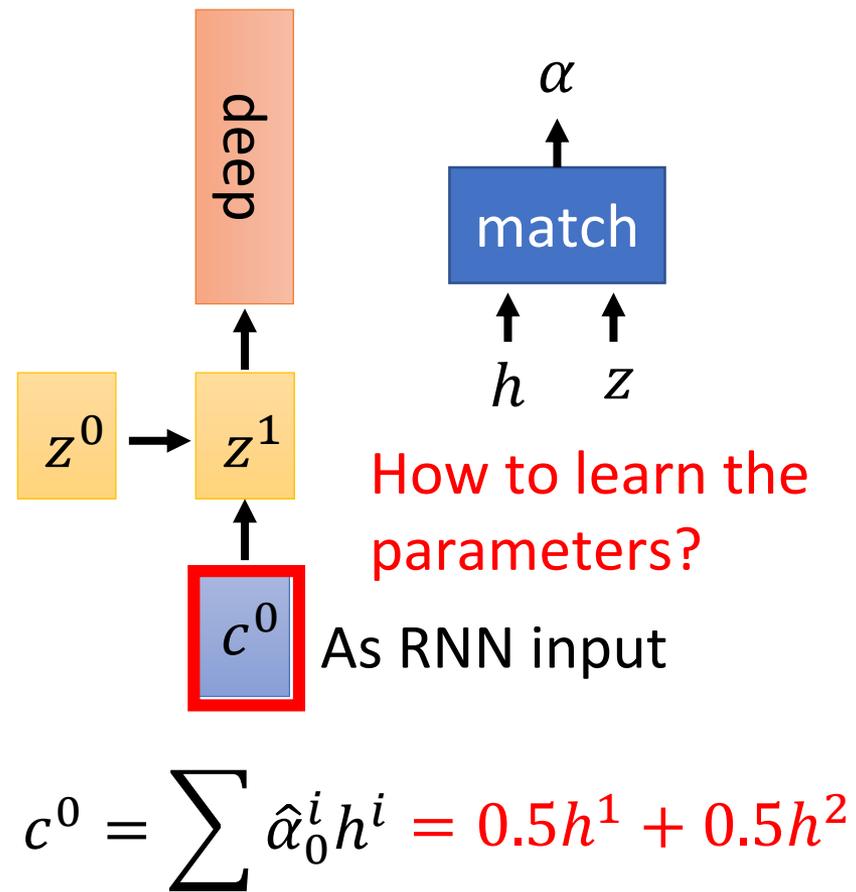
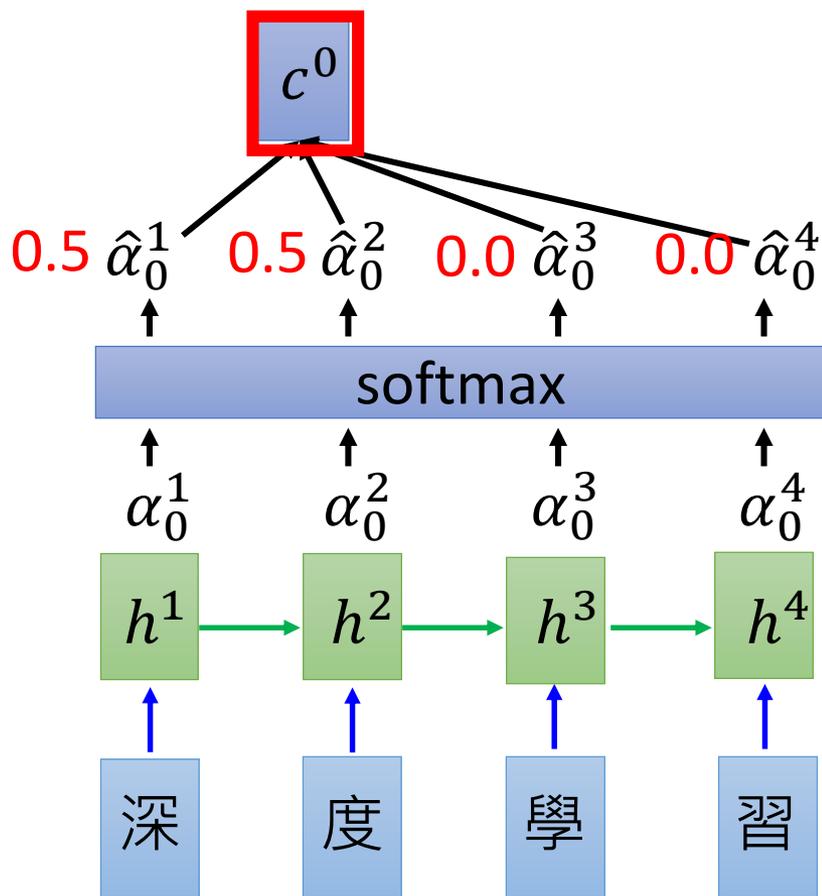
What is **match** ?

- Cosine similarity of z and h
- Small NN whose input is z and h , output a scalar
- $\alpha = h^T W z$

How to learn the parameters?

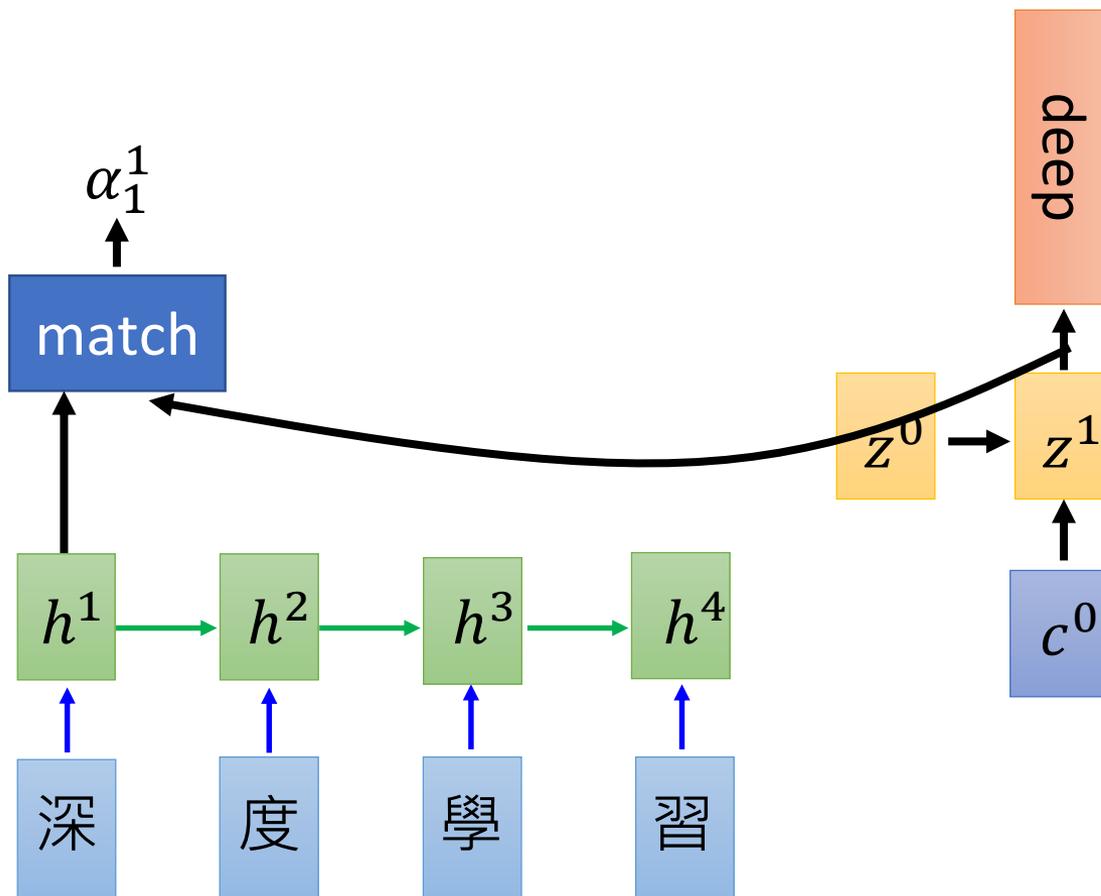
Machine Translation with Attention

193



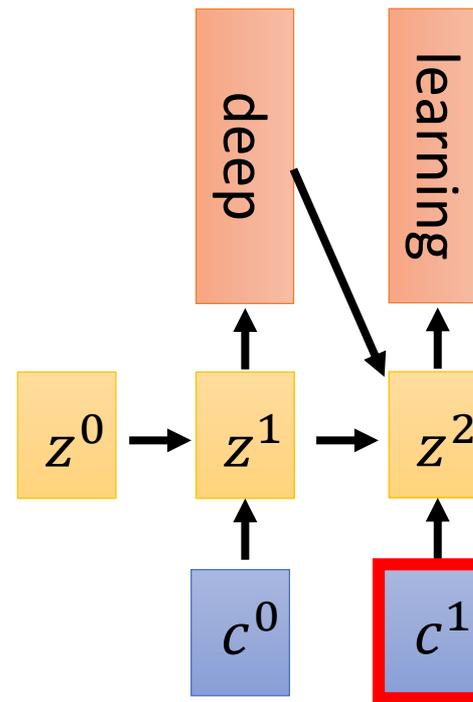
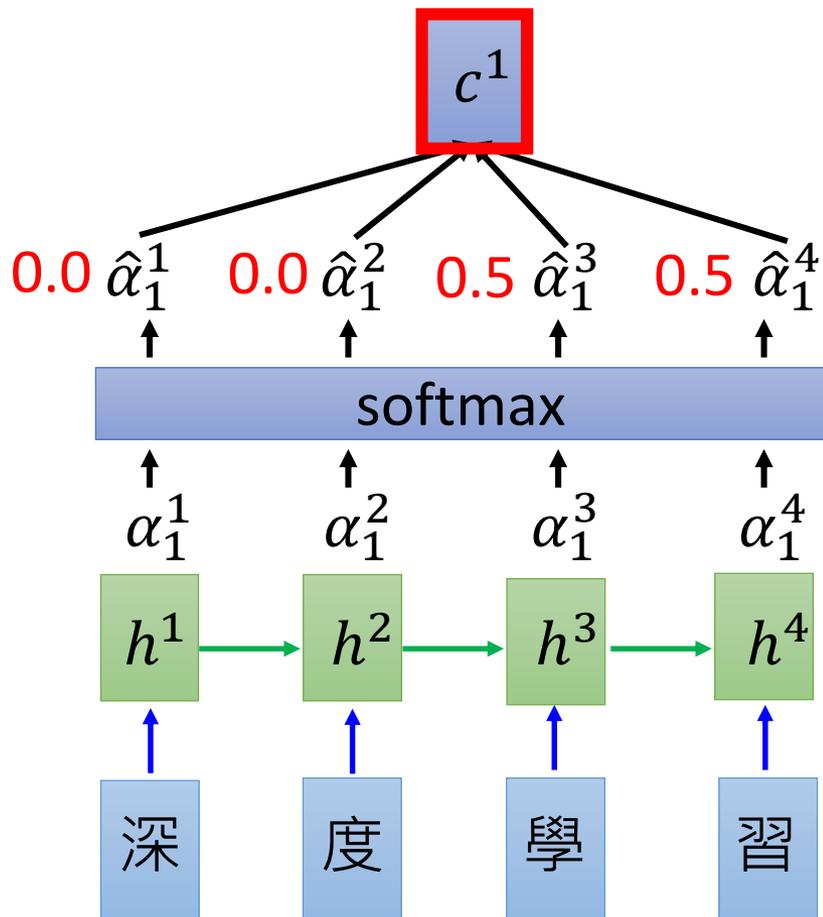
Machine Translation with Attention

194



Machine Translation with Attention

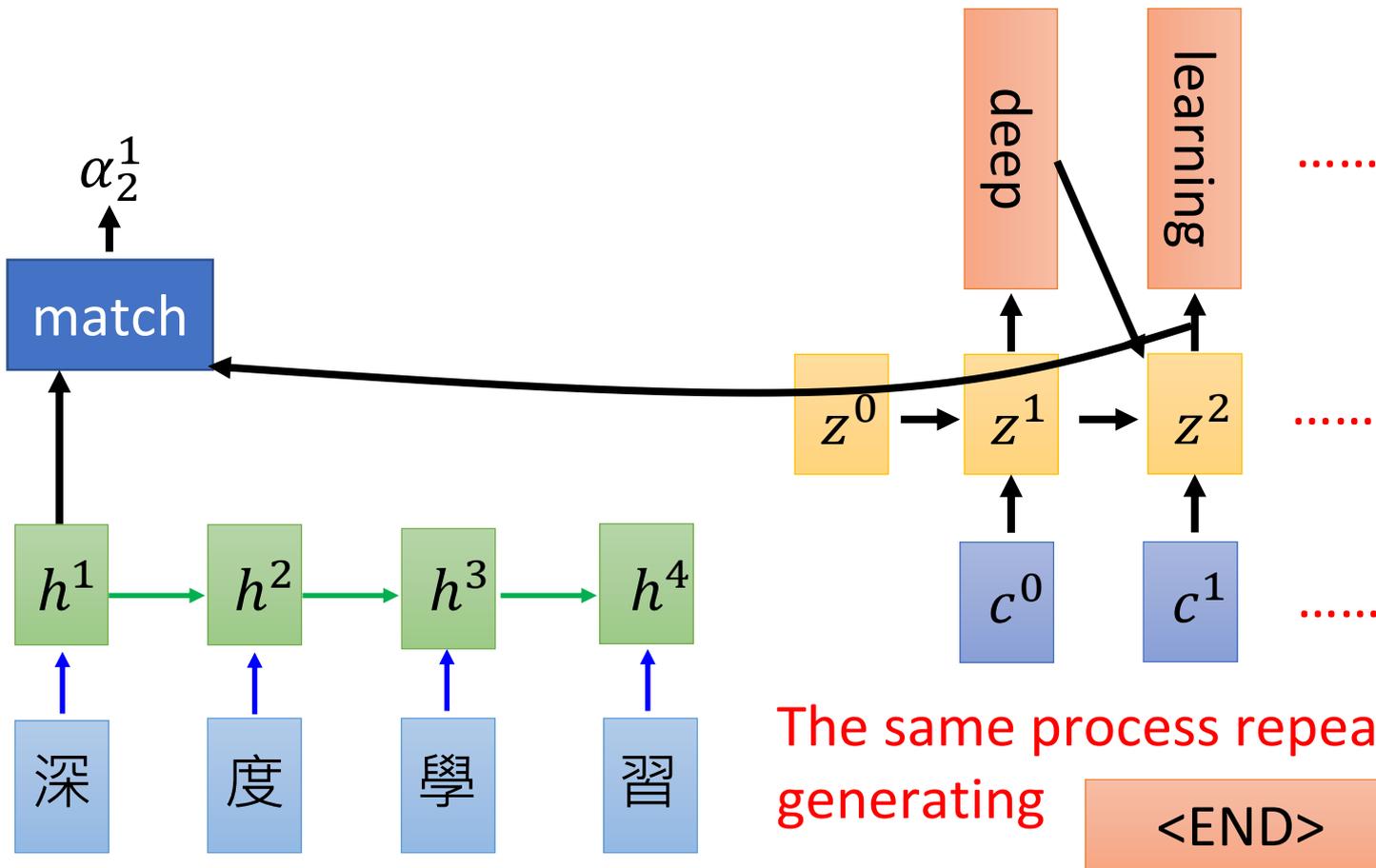
195



$$c^1 = \sum \hat{\alpha}_1^i h^i = 0.5h^3 + 0.5h^4$$

Machine Translation with Attention

196



Speech Recognition with Attention

197

Alignment between the Characters and Audio

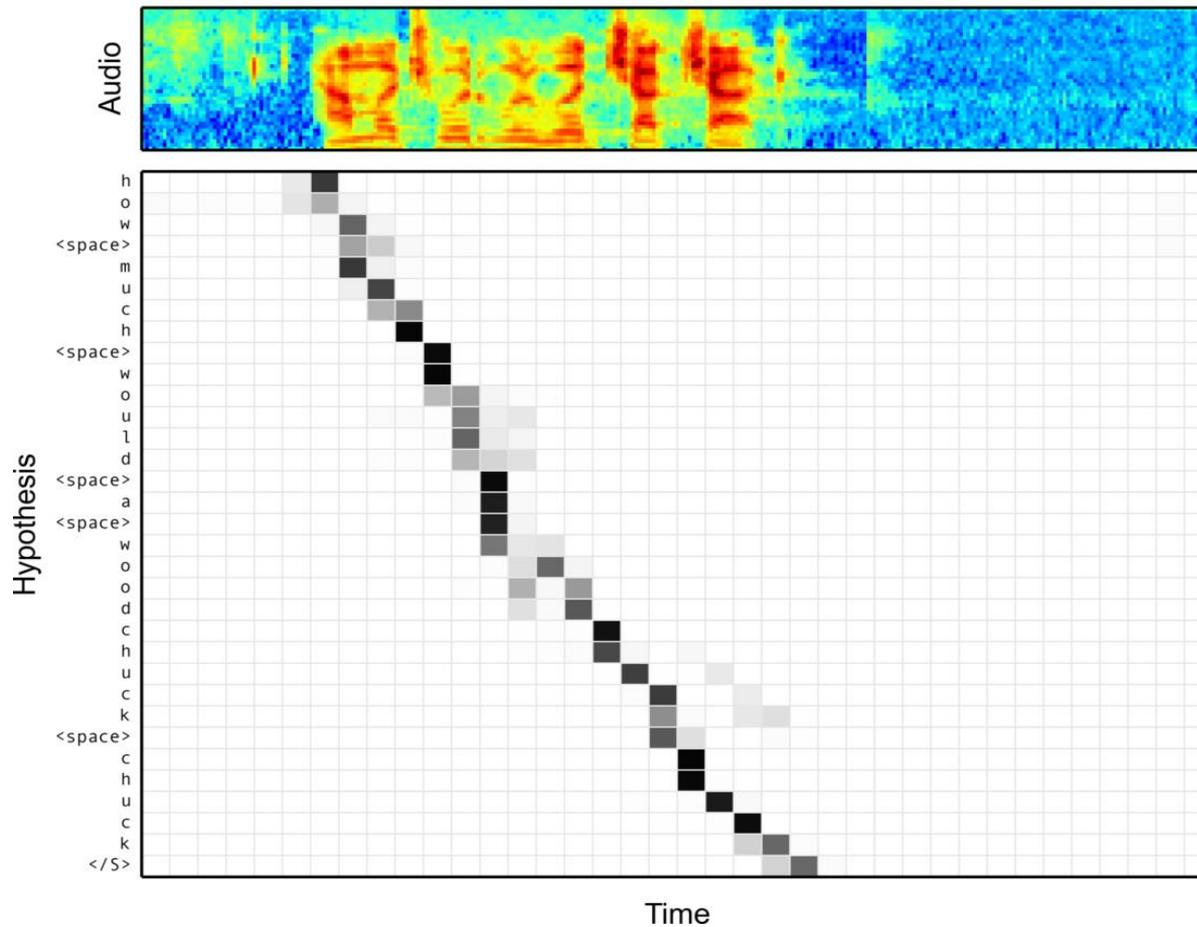


Image Captioning

198

- Input: image
- Output: word sequence

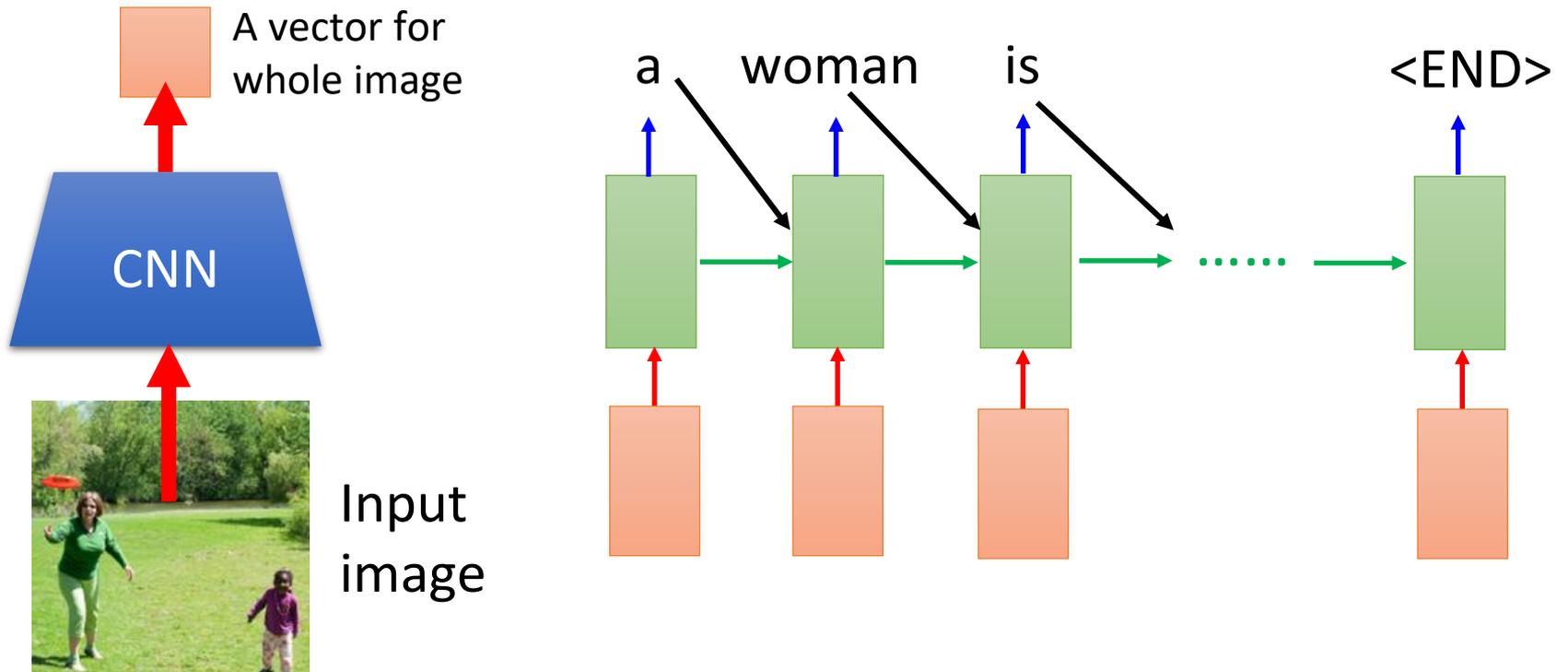


Image Captioning with Attention

199

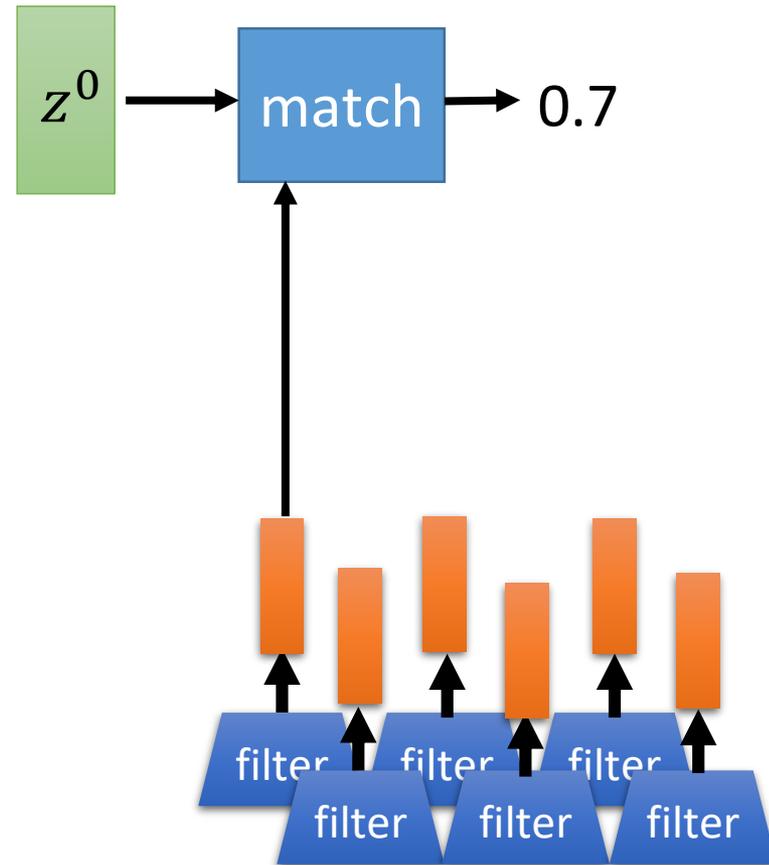
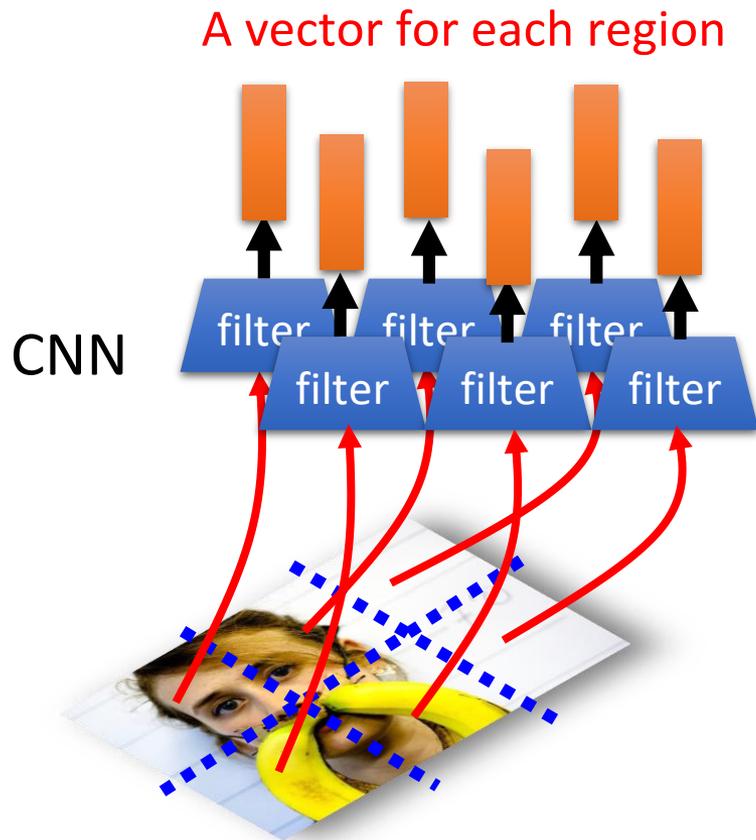


Image Captioning with Attention

200

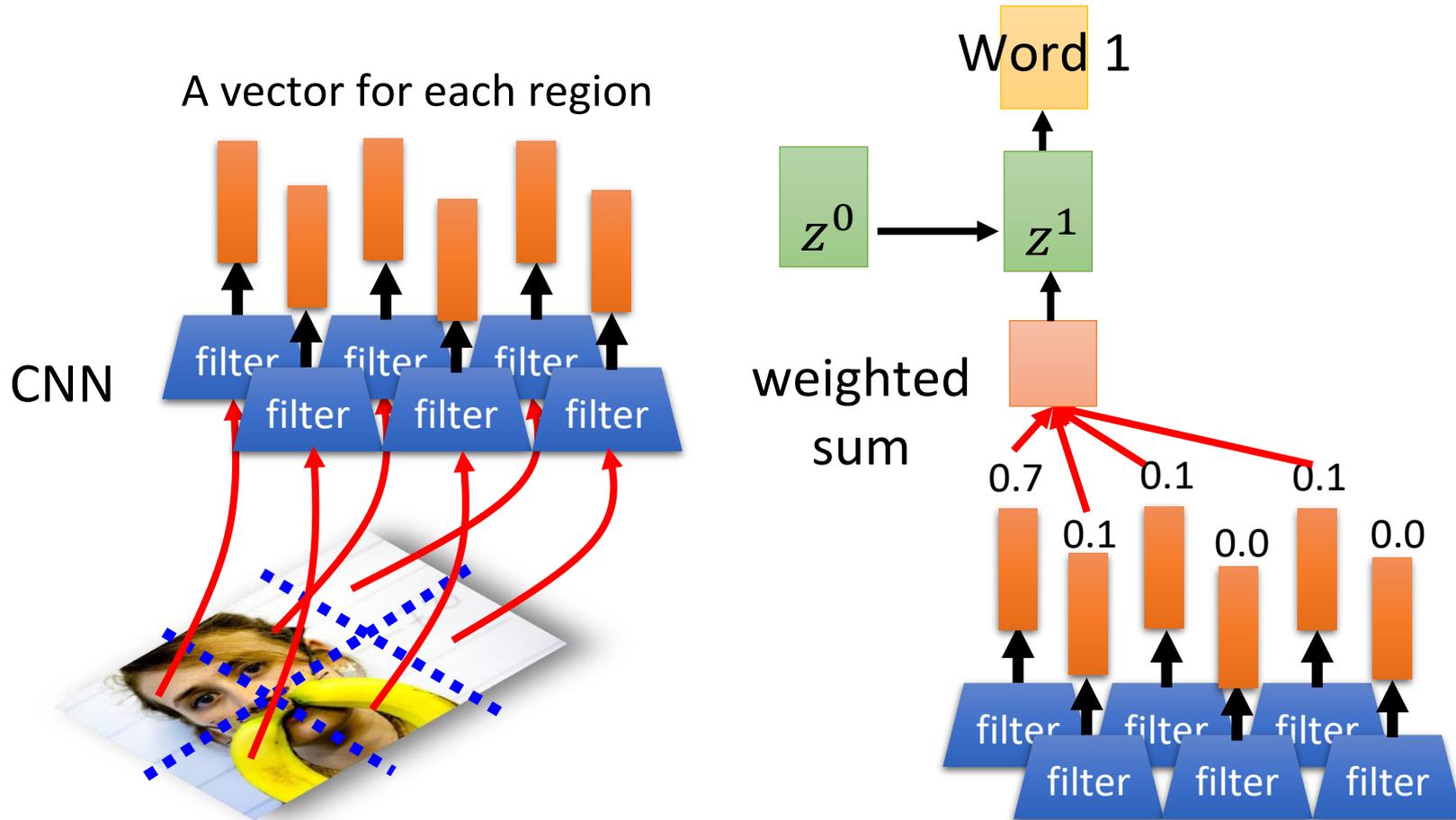


Image Captioning with Attention

201

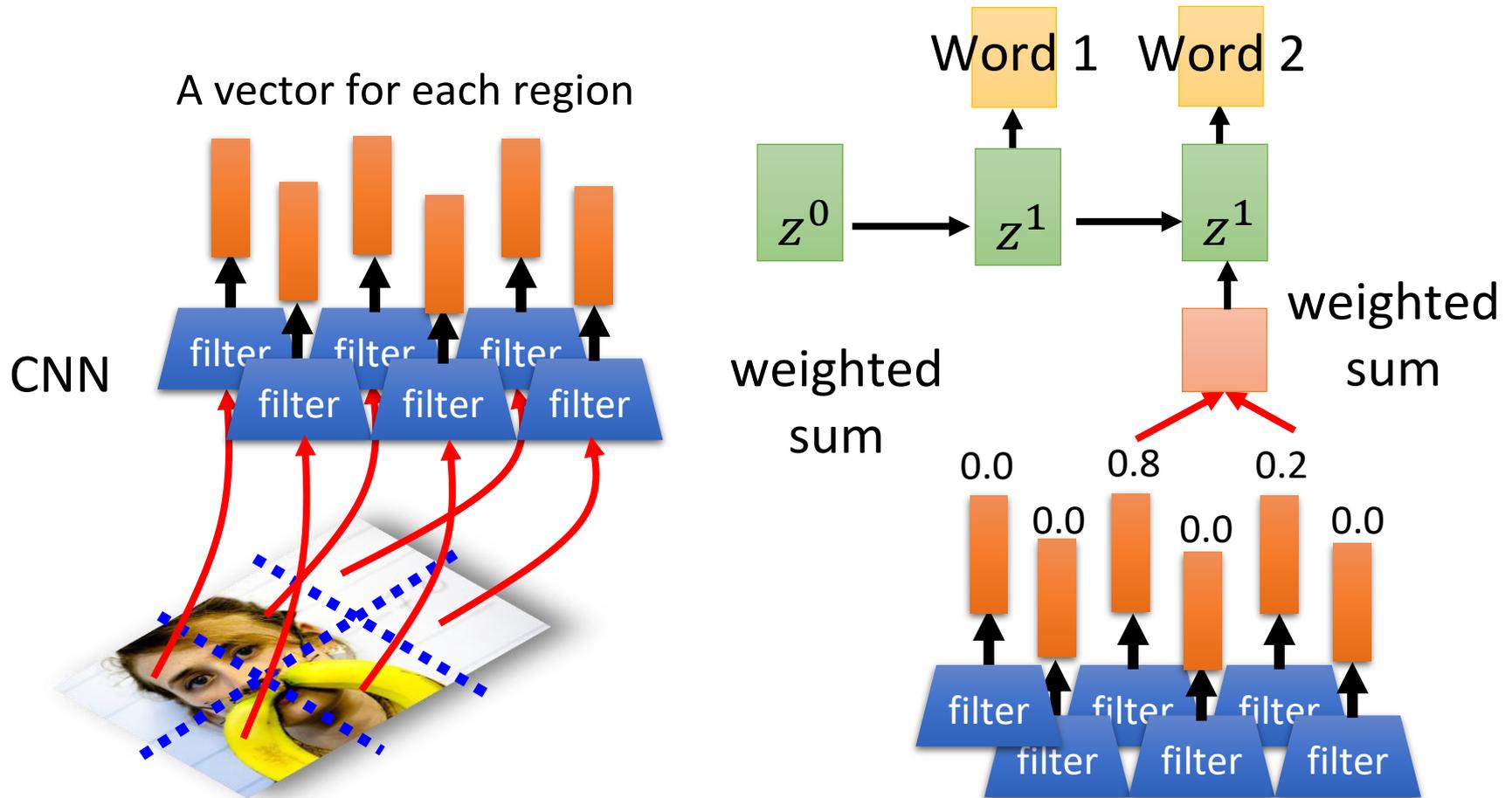


Image Captioning

202

□ Good examples



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

Image Captioning

203

□ Bad examples



A large white bird standing in a forest.



A woman holding a clock in her hand.



A man wearing a hat and a hat on a skateboard.



A person is standing on a beach with a surfboard.



A woman is sitting at a table with a large pizza.



A man is talking on his cell phone while another man watches.

Video Captioning

204



Ref: A man and a woman ride a motorcycle
A **man** and a **woman** are **talking** on the **road**

Video Captioning

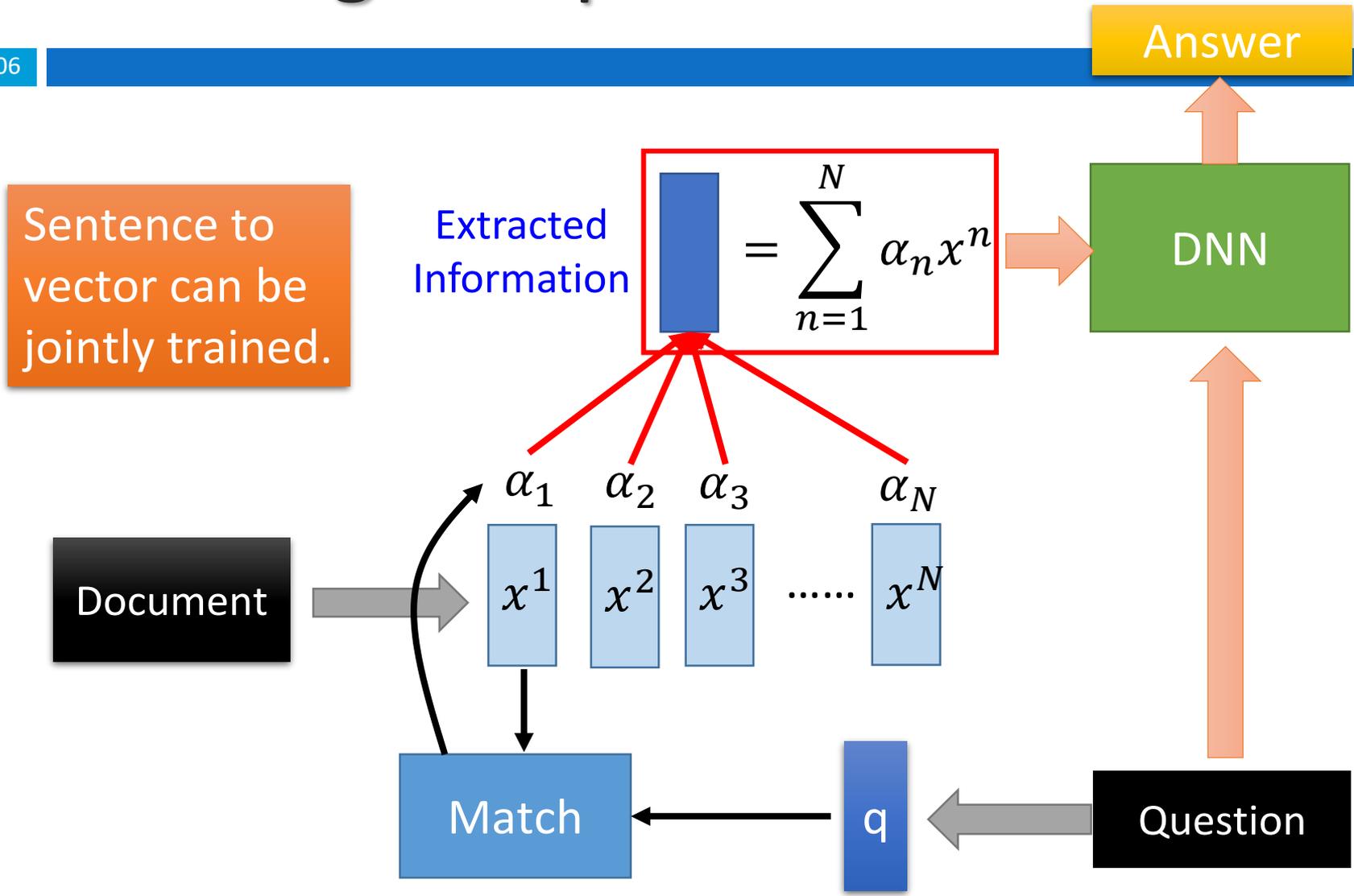
205



Ref: A woman is frying food
Someone is **frying** a **fish** in a **pot**

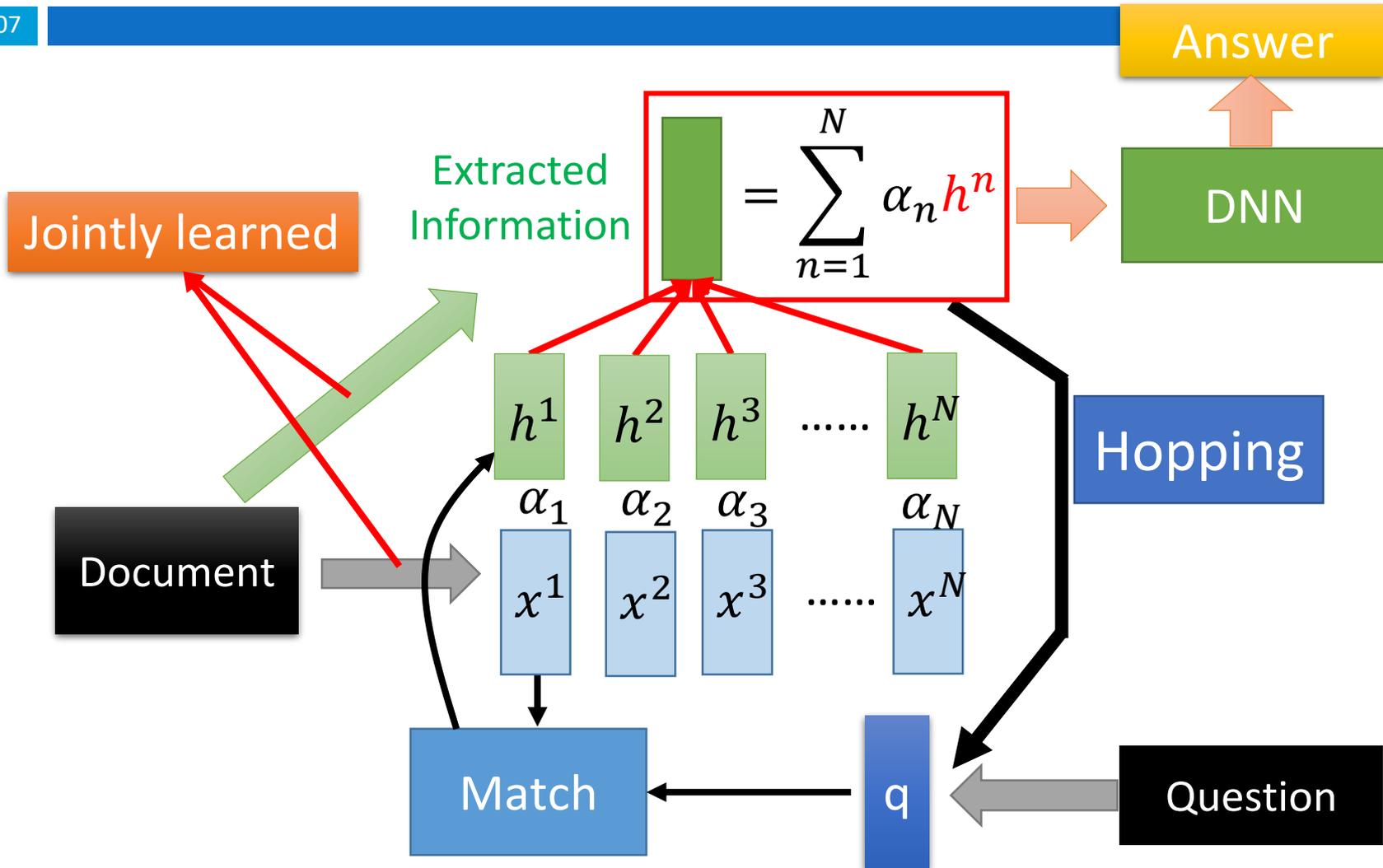
Reading Comprehension

206



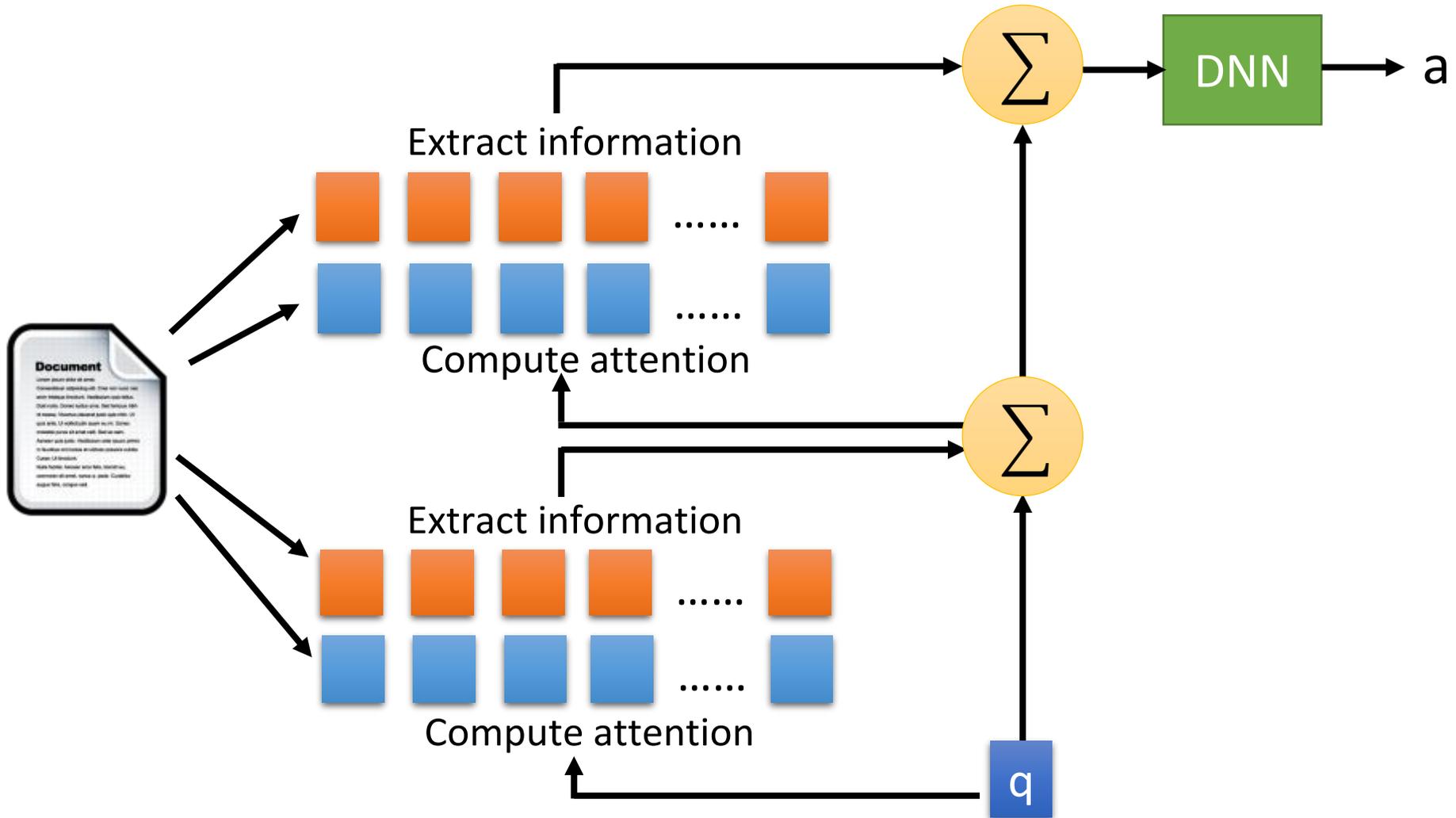
Reading Comprehension

207



Memory Network

208



Memory Network

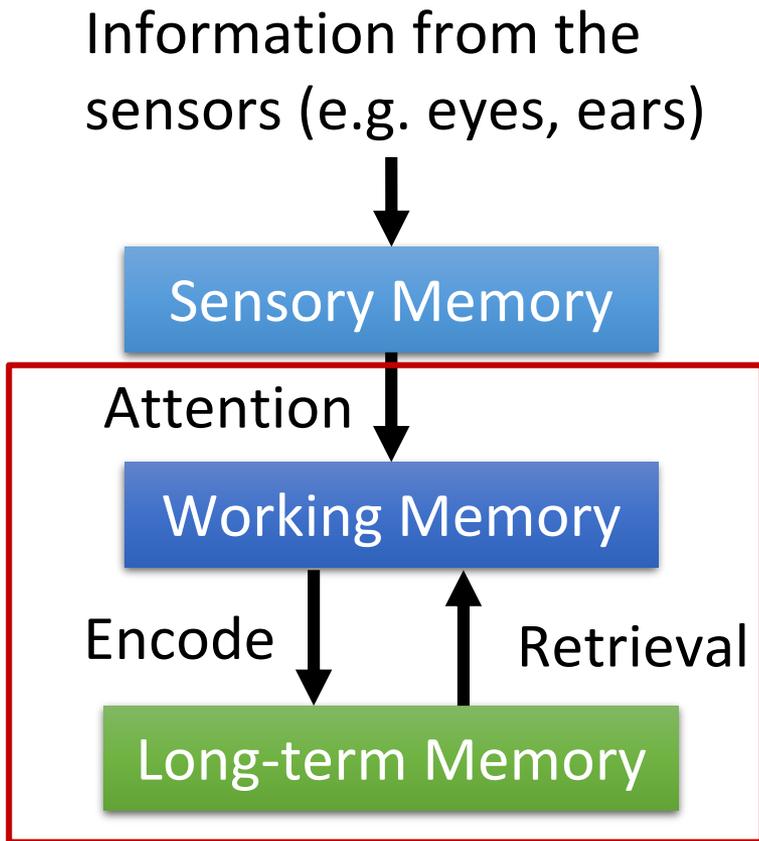
209

- Multi-hop performance analysis

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
What color is Greg? Answer: yellow Prediction: yellow				

Attention on Memory

210



When the input is a very long sequence or an image

➔ Pay attention on partial of the input object each time

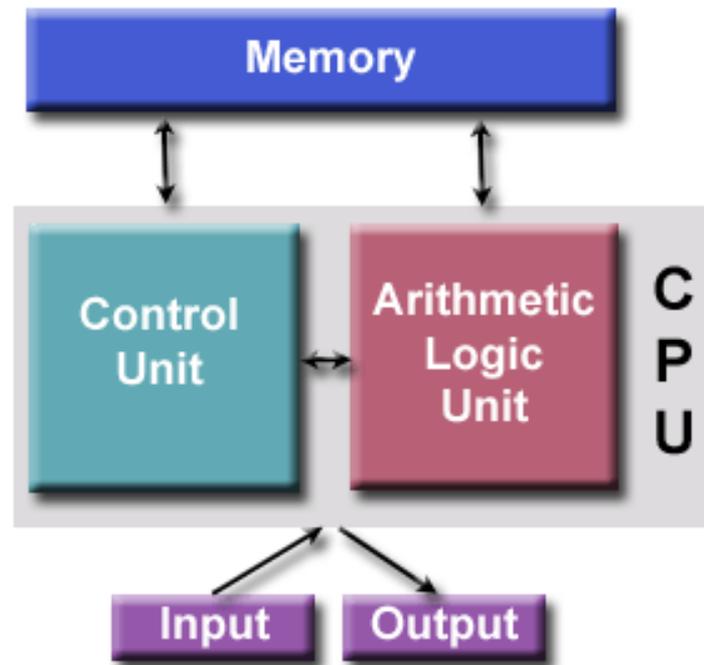
In RNN/LSTM, larger memory implies more parameters

➔ Increasing memory size will not increasing parameters

Neural Turing Machine

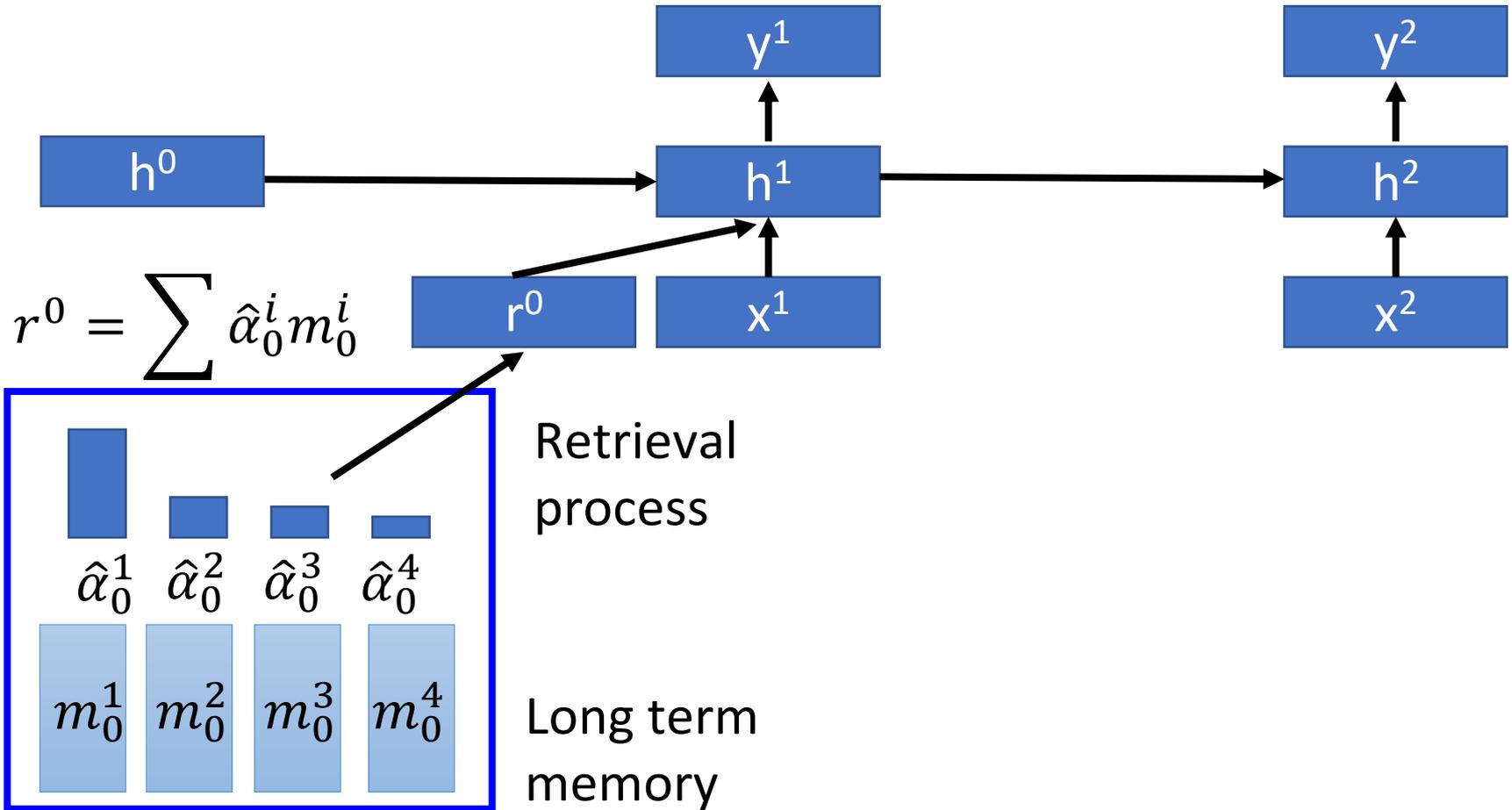
211

- Von Neumann architecture
- Neural Turing Machine is an advanced RNN/LSTM.



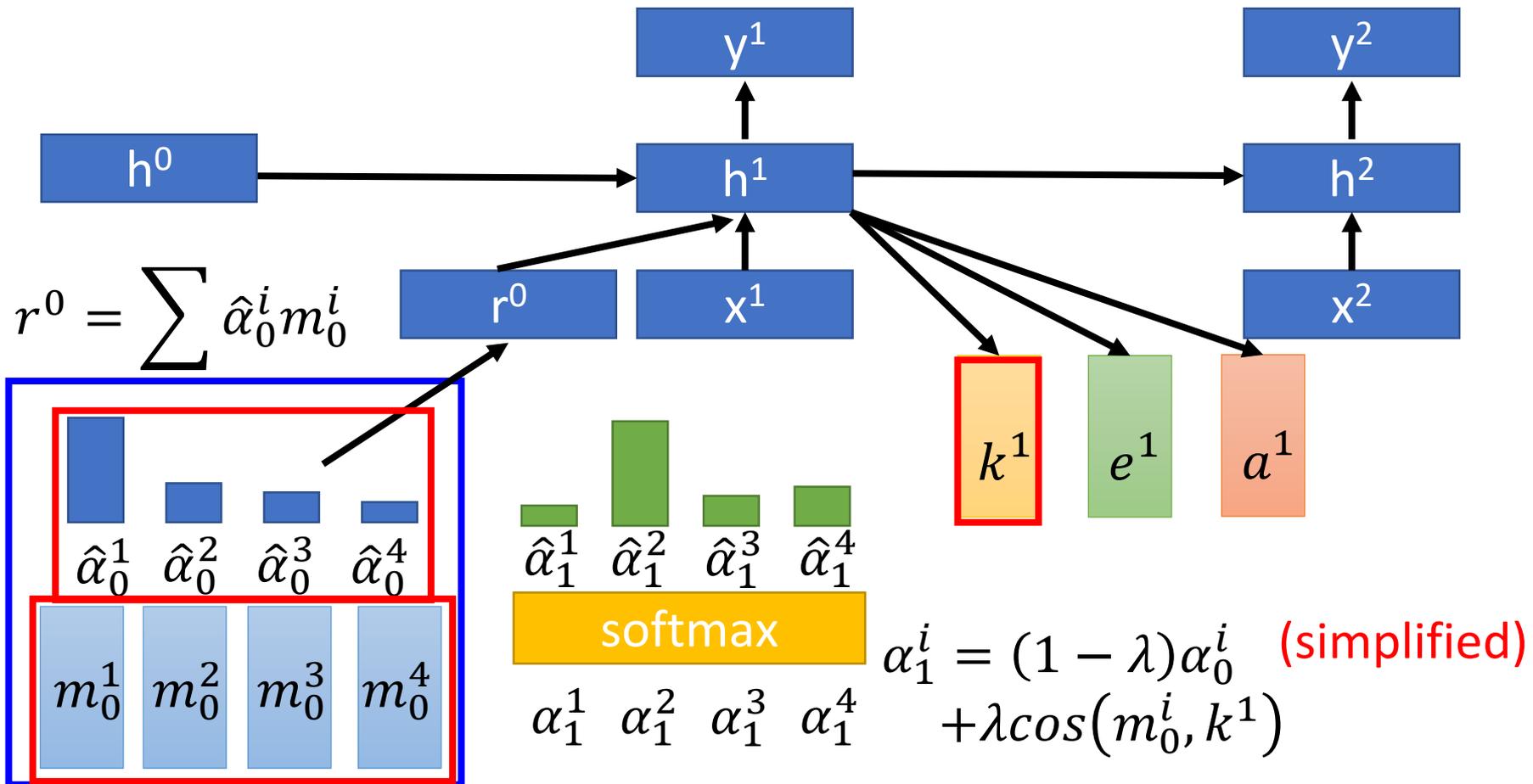
Neural Turing Machine

212



Neural Turing Machine

213

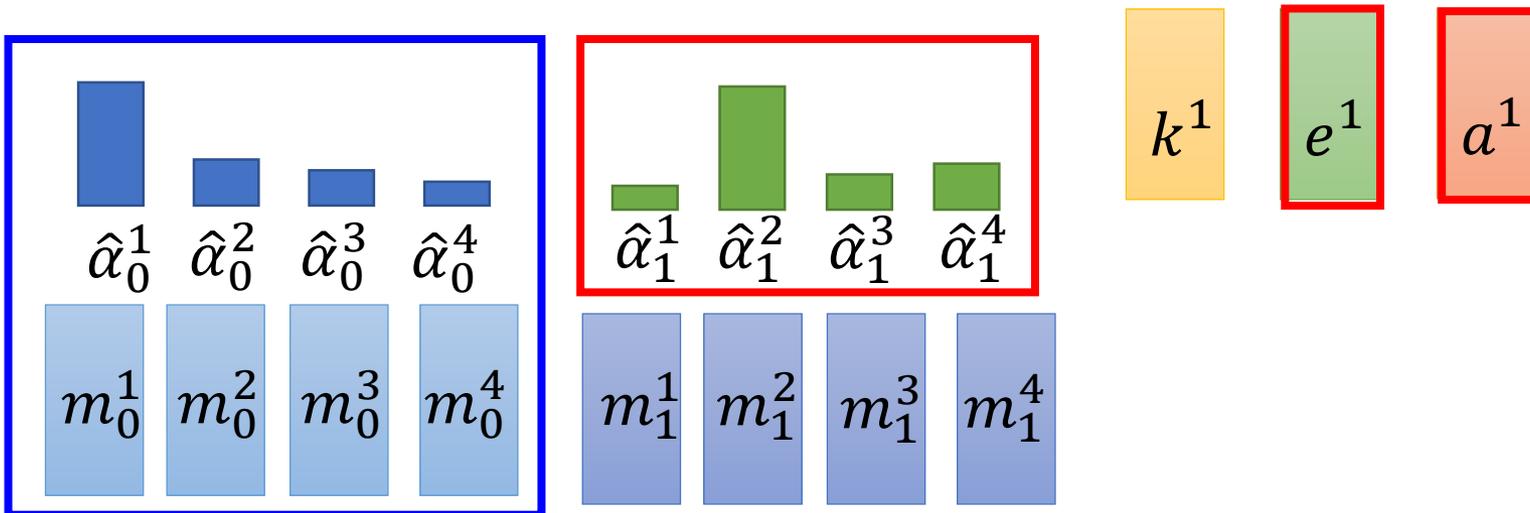


Neural Turing Machine

214

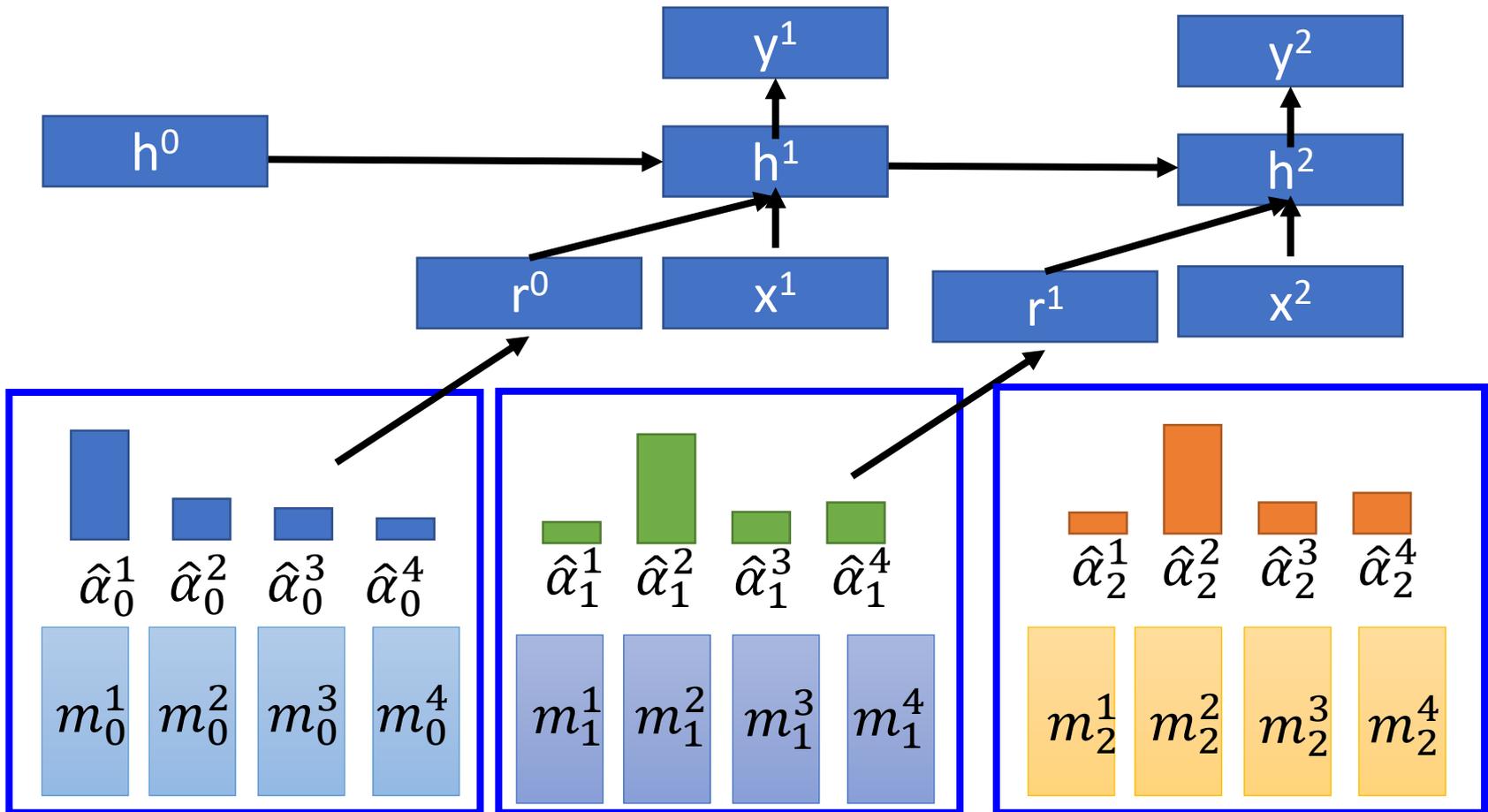
$$m_1^i = m_0^i * \begin{pmatrix} 1 & -\hat{\alpha}_1^i \\ & e^1 \end{pmatrix} + \hat{\alpha}_1^i a^1 \quad \rightarrow \text{Encode process}$$

(element-wise)



Neural Turing Machine

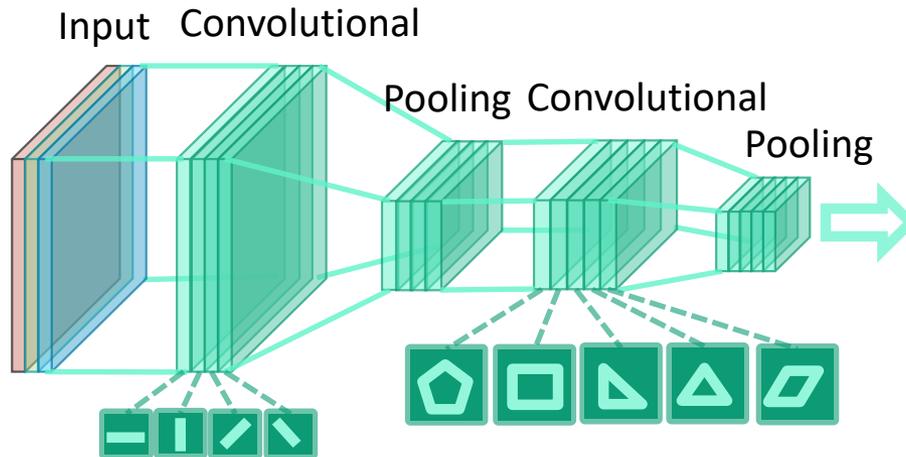
215



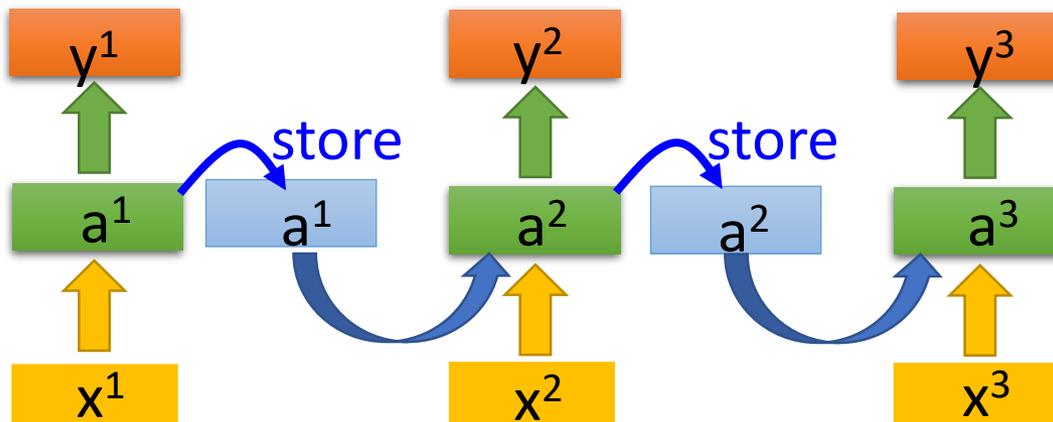
Concluding Remarks

216

□ Convolutional Neural Network (CNN)



□ Recurrent Neural Network (RNN)



Talk Outline

217

Part I: Introduction to
Machine Learning & Deep Learning



Part II: Variants of Neural Nets



Part III: Beyond Supervised Learning
& Recent Trends

218

PART III

Beyond Supervised Learning & Recent Trend

Introduction

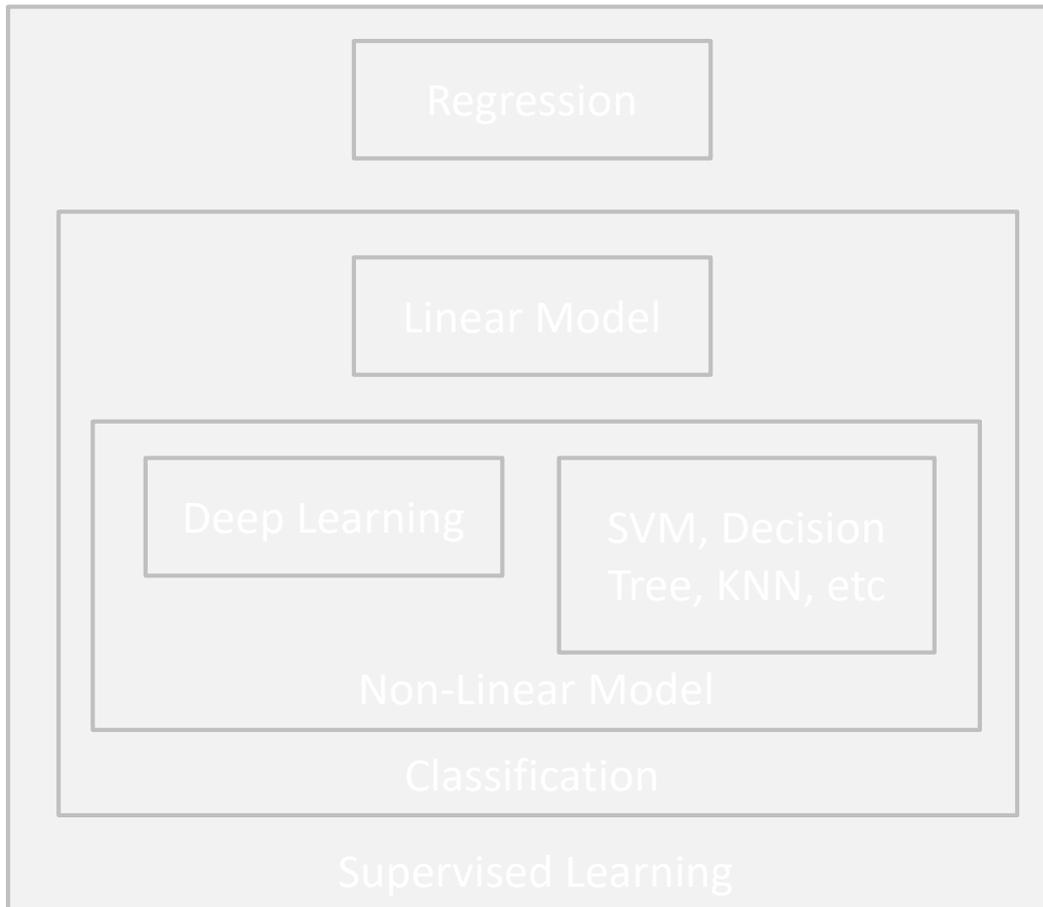
219

- Big data \neq Big annotated data
- Machine learning techniques include:
 - ▣ Supervised learning (if we have labelled data)
 - ▣ Reinforcement learning (if we have an environment for reward)
 - ▣ Unsupervised learning (if we do not have labelled data)

What can we do if there is no sufficient labelled training data?

Machine Learning Map

Scenario Task Method



Semi-Supervised Learning

Transfer Learning

Unsupervised Learning

Reinforcement Learning

Outline

221

- Semi-Supervised Learning
- Transfer Learning
- Unsupervised Learning
 - 化繁為簡 Representation Learning
 - 無中生有 Generative Model
- Reinforcement Learning

Outline

222

- **Semi-Supervised Learning**
- Transfer Learning
- Unsupervised Learning
 - ▣ 化繁為簡 Representation Learning
 - ▣ 無中生有 Generative Model
- Reinforcement Learning

Semi-Supervised Learning

223

Labelled
data



cat



dog

Unlabelled
data

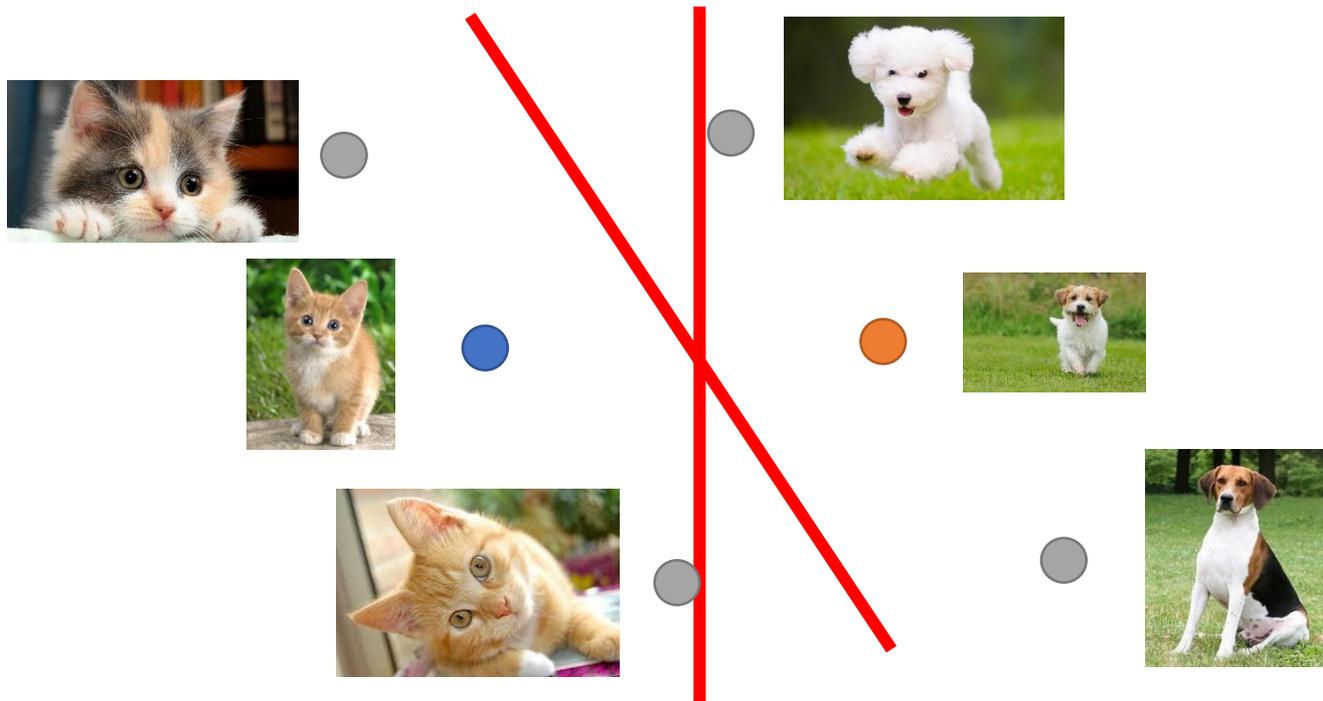


(Image of cats and dogs without labeling)

Semi-Supervised Learning

224

- Why semi-supervised learning helps?



The distribution of the unlabeled data provides some cues

Outline

225

- Semi-Supervised Learning
- **Transfer Learning**
- Unsupervised Learning
 - ▣ 化繁為簡 Representation Learning
 - ▣ 無中生有 Generative Model
- Reinforcement Learning

Transfer Learning

226

Labelled
data



cat



dog

Labeled
data



elephant



elephant



tiger



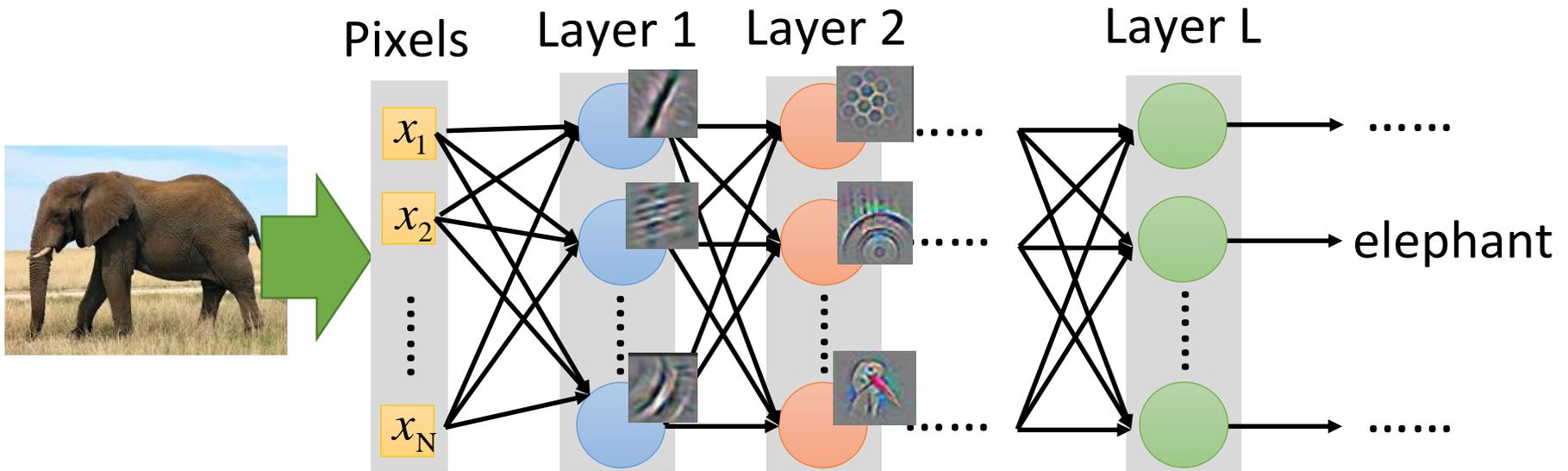
tiger

Not related to the task considered

Transfer Learning

227

- Widely used on image processing
 - ▣ Using sufficient labeled data to learn a CNN
 - ▣ Using this CNN as feature extractor



Transfer Learning Example

228

研究生 online

漫畫家 online

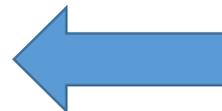
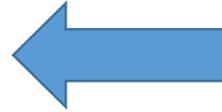
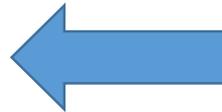
研究生
生存守則

研究生

指導教授

跑實驗

投稿期刊

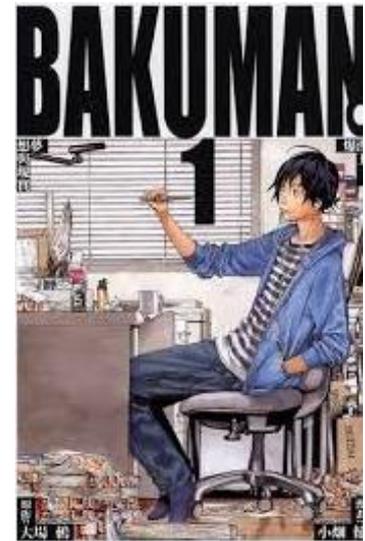


漫畫家

責編

畫分鏡

投稿 jump



爆漫王

Outline

229

- Semi-Supervised Learning
- Transfer Learning
- Unsupervised Learning
 - 化繁為簡 Representation Learning
 - 無中生有 Generative Model
- Reinforcement Learning

Outline

230

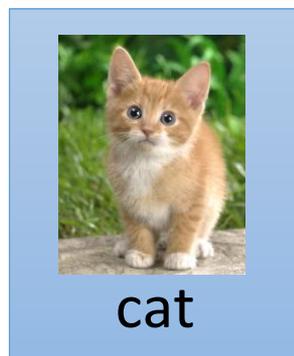
- Semi-Supervised Learning
- Transfer Learning
- **Unsupervised Learning**
 - ▣ 化繁為簡 Representation Learning
 - ▣ 無中生有 Generative Model
- Reinforcement Learning

Unsupervised Learning

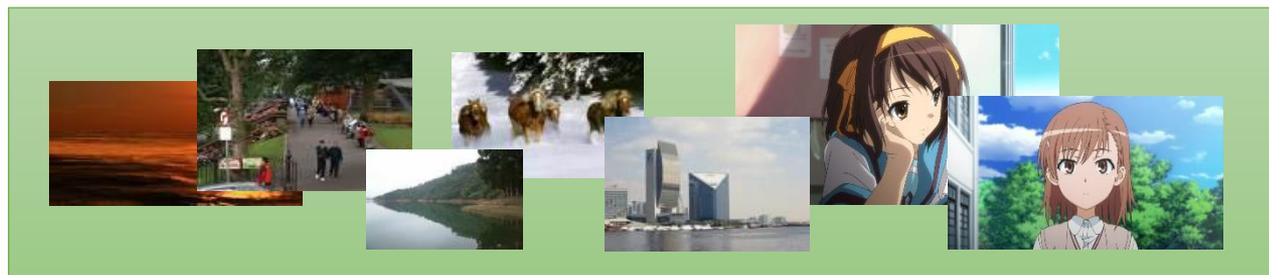
231

- The unlabeled data sometimes is not related to the task

Labelled
data



Unlabeled
data



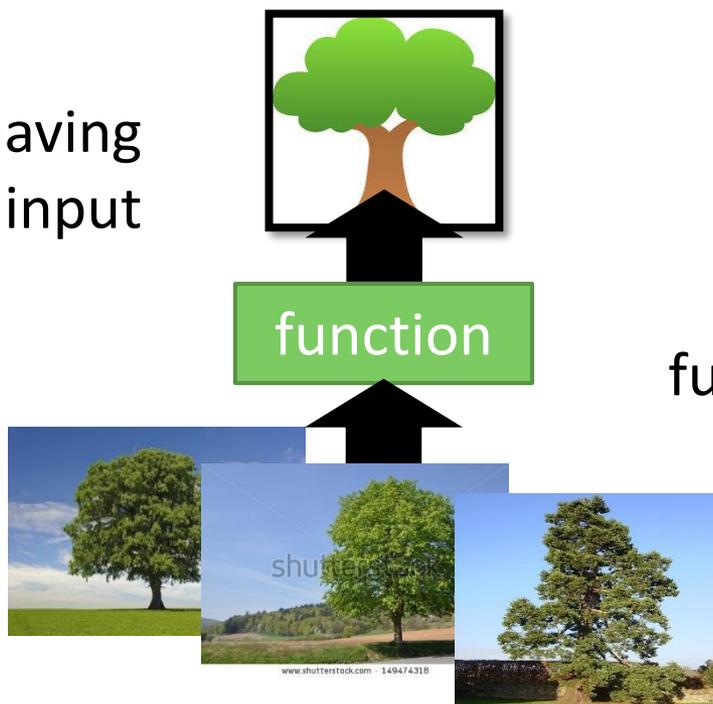
(Just crawl millions of images from the Internet)

Unsupervised Learning

232

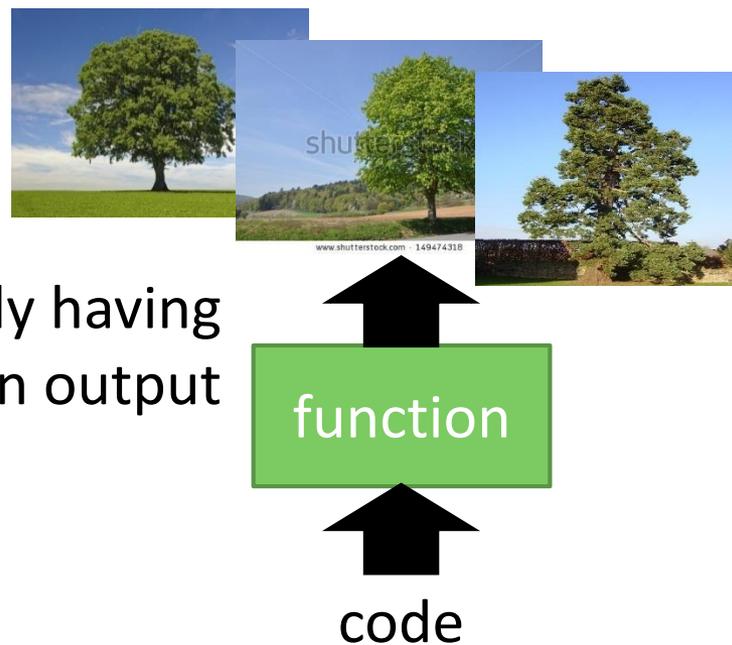
□ 化繁為簡

only having
function input



□ 無中生有

only having
function output



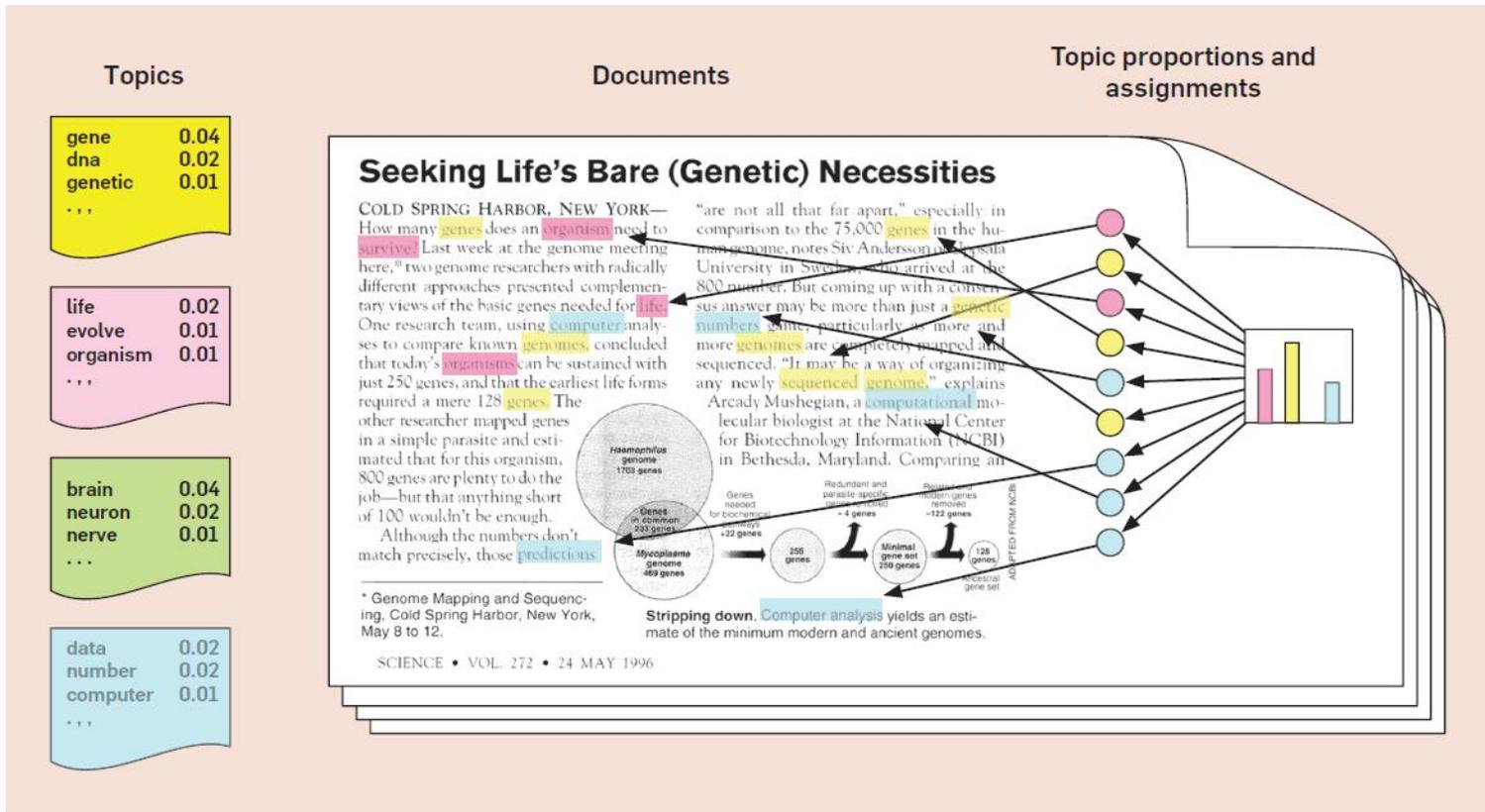
Unsupervised Learning

233

- How does self-taught learning work?
- Why does unlabeled and unrelated data help the tasks?

Finding latent factors that control the observations

Latent Factors for Documents



Latent Factors for Recommendation

236

單純呆

A



傲嬌

B

C



Latent Factor Exploitation

237

- Handwritten digits



The handwritten images are composed of **strokes**

Strokes (Latent Factors)



No. 1



No. 2



No. 3



No. 4

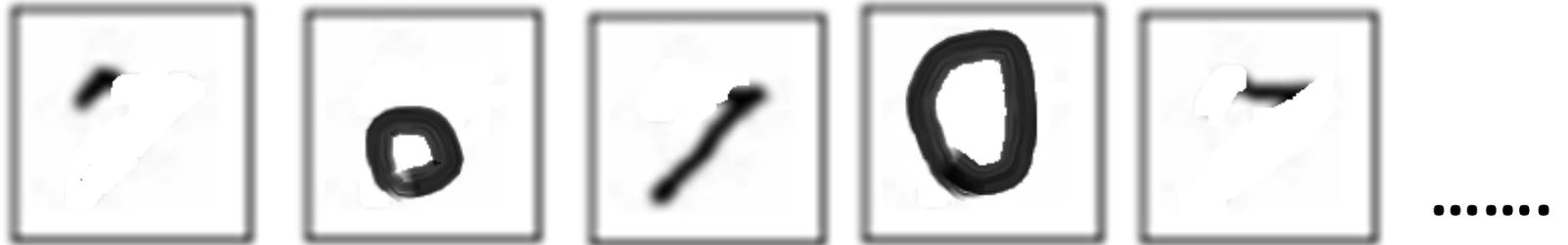


No. 5

.....

Latent Factor Exploitation

Strokes (Latent Factors)



No. 1

No. 2

No. 3

No. 4

No. 5

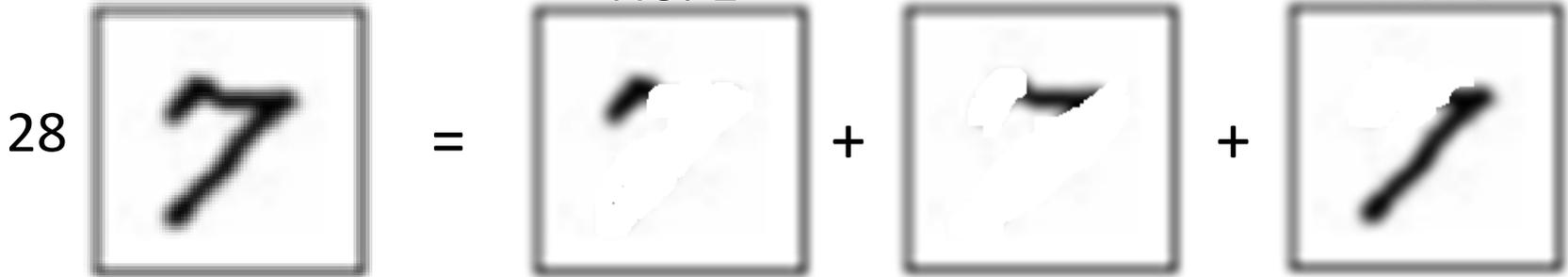
.....

28

No. 1

No. 3

No. 5



Represented by
28 X 28 = 784 pixels

[1 0 1 0 1 0]
(simpler representation)

Outline

239

- Semi-Supervised Learning
- Transfer Learning
- Unsupervised Learning
 - 化繁為簡 Representation Learning
 - 無中生有 Generative Model
- Reinforcement Learning

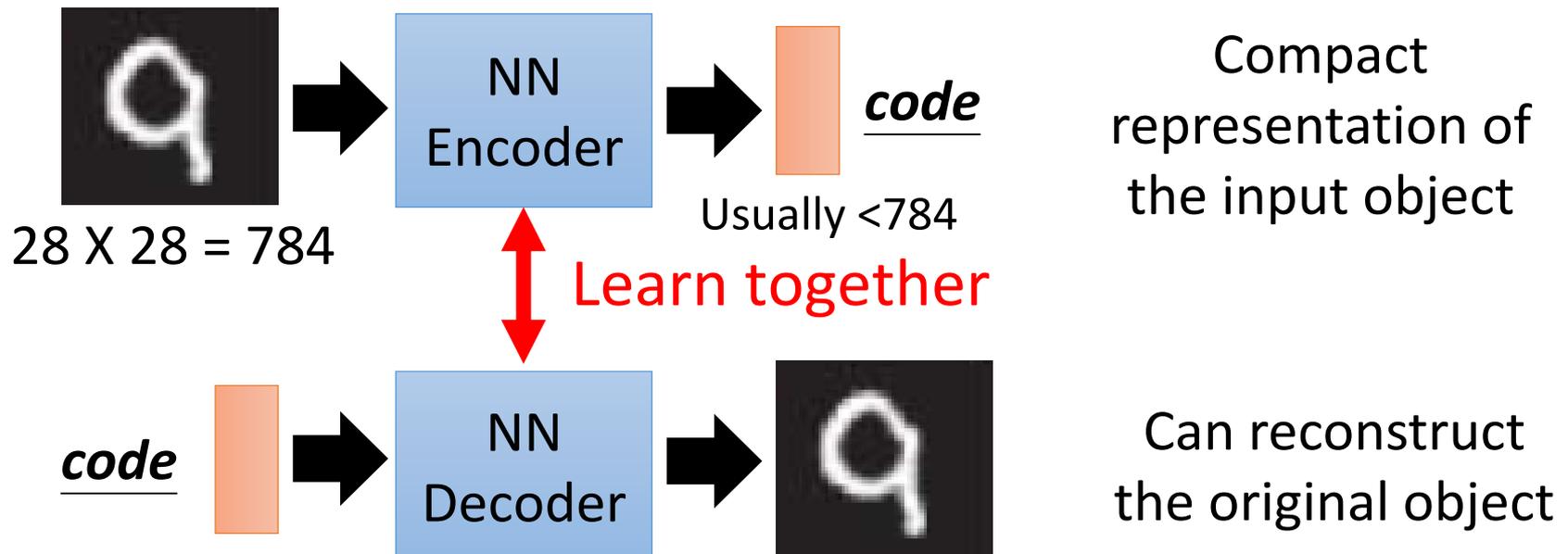
Autoencoder

240

- Represent a digit using 28 X 28 dimensions
- Not all 28 X 28 images are digits

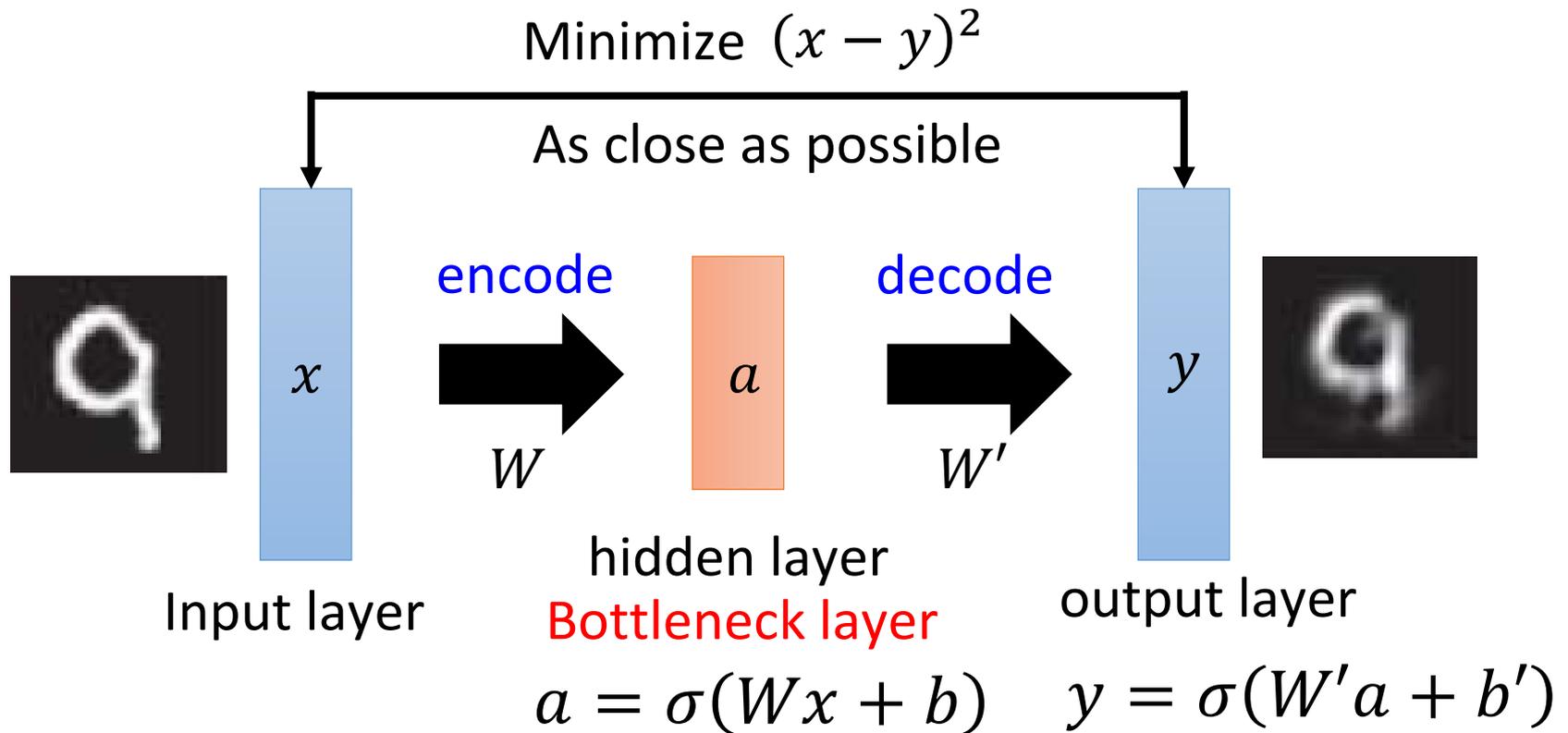


Idea: represent the images of digits in a more compact way



Autoencoder

241

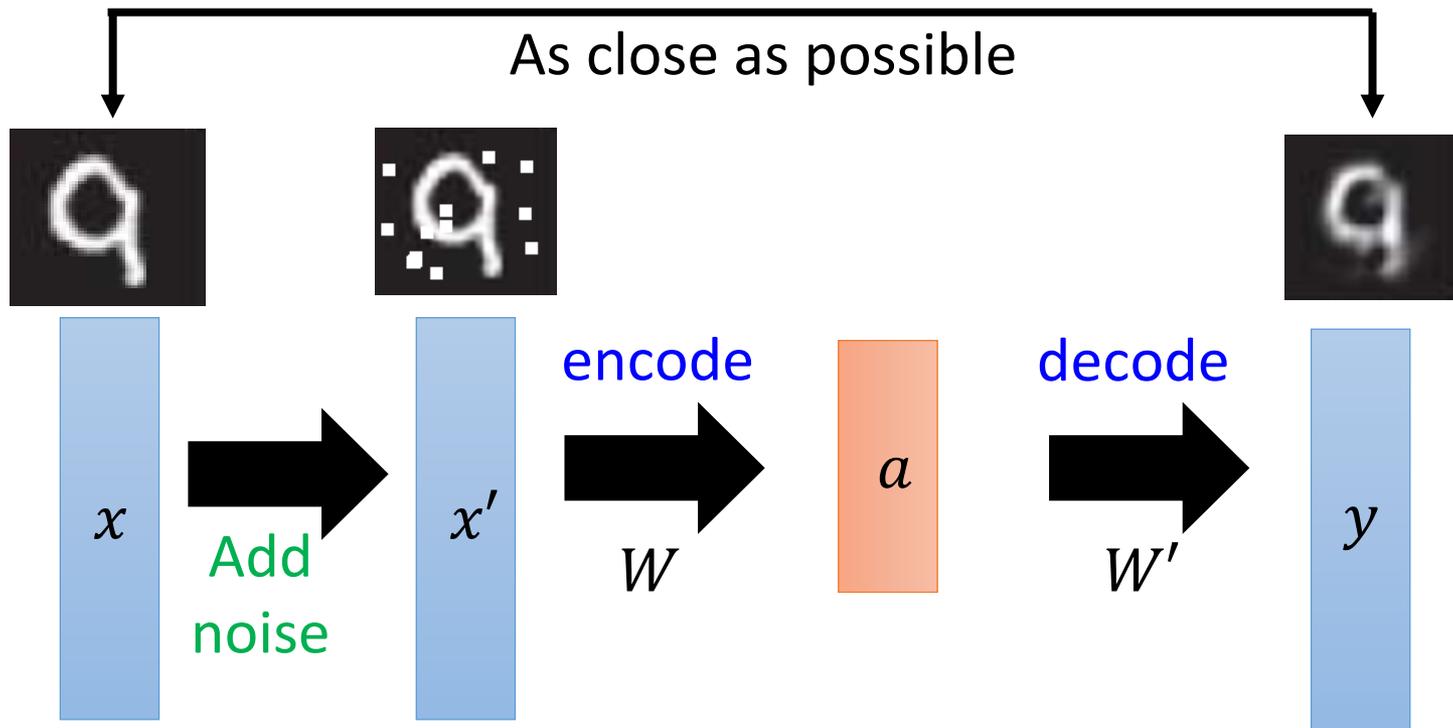


Output of the hidden layer is the code

Autoencoder

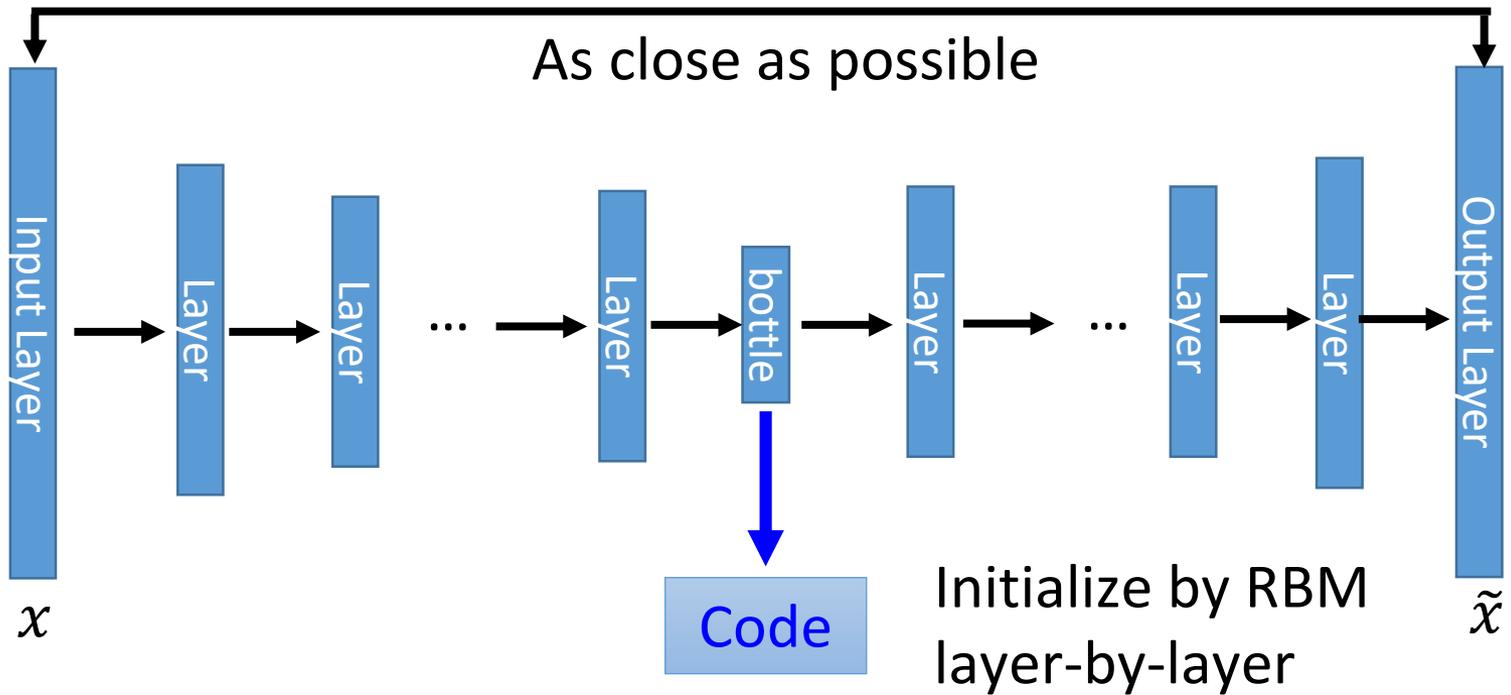
242

- De-noising auto-encoder



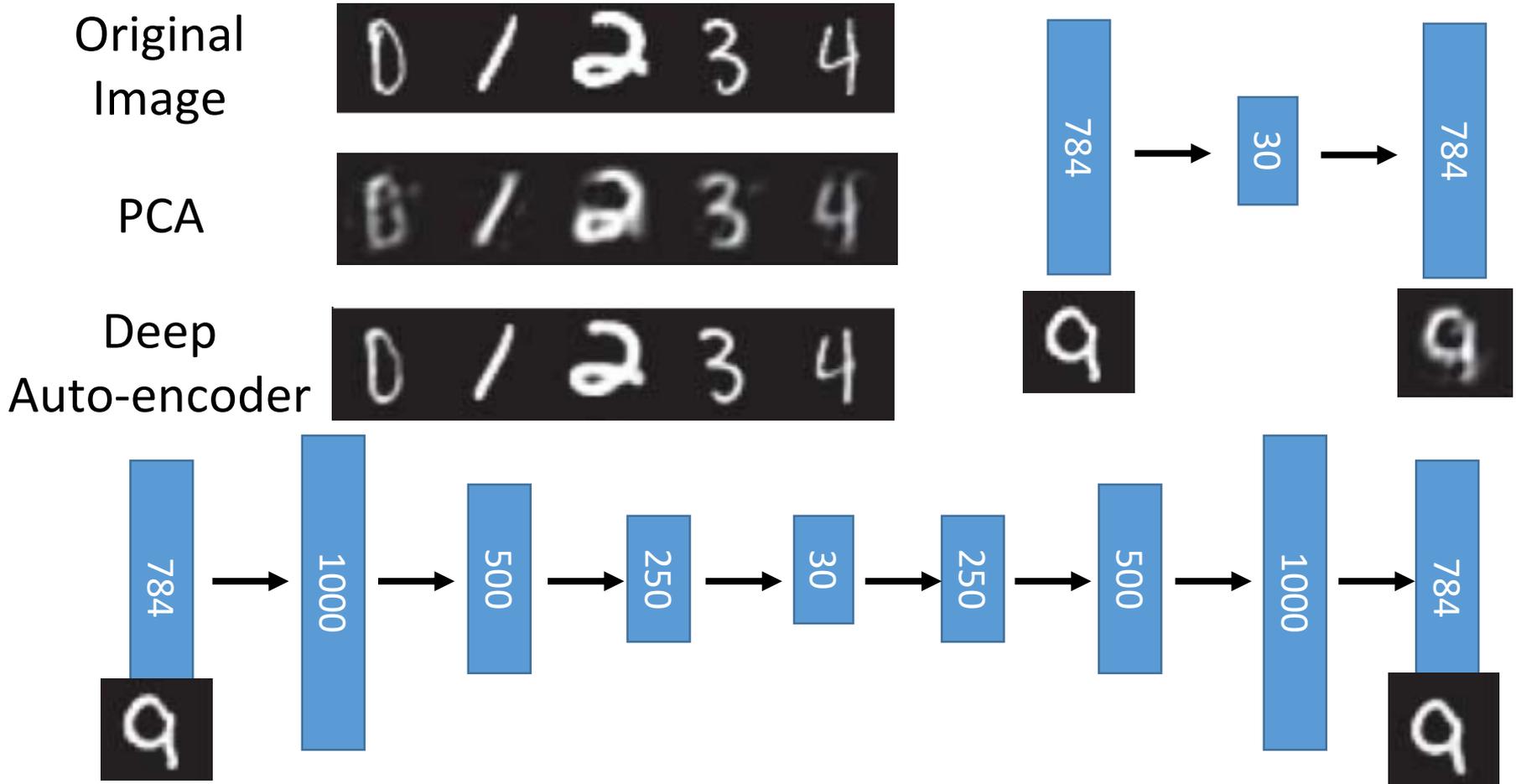
Deep Autoencoder

243

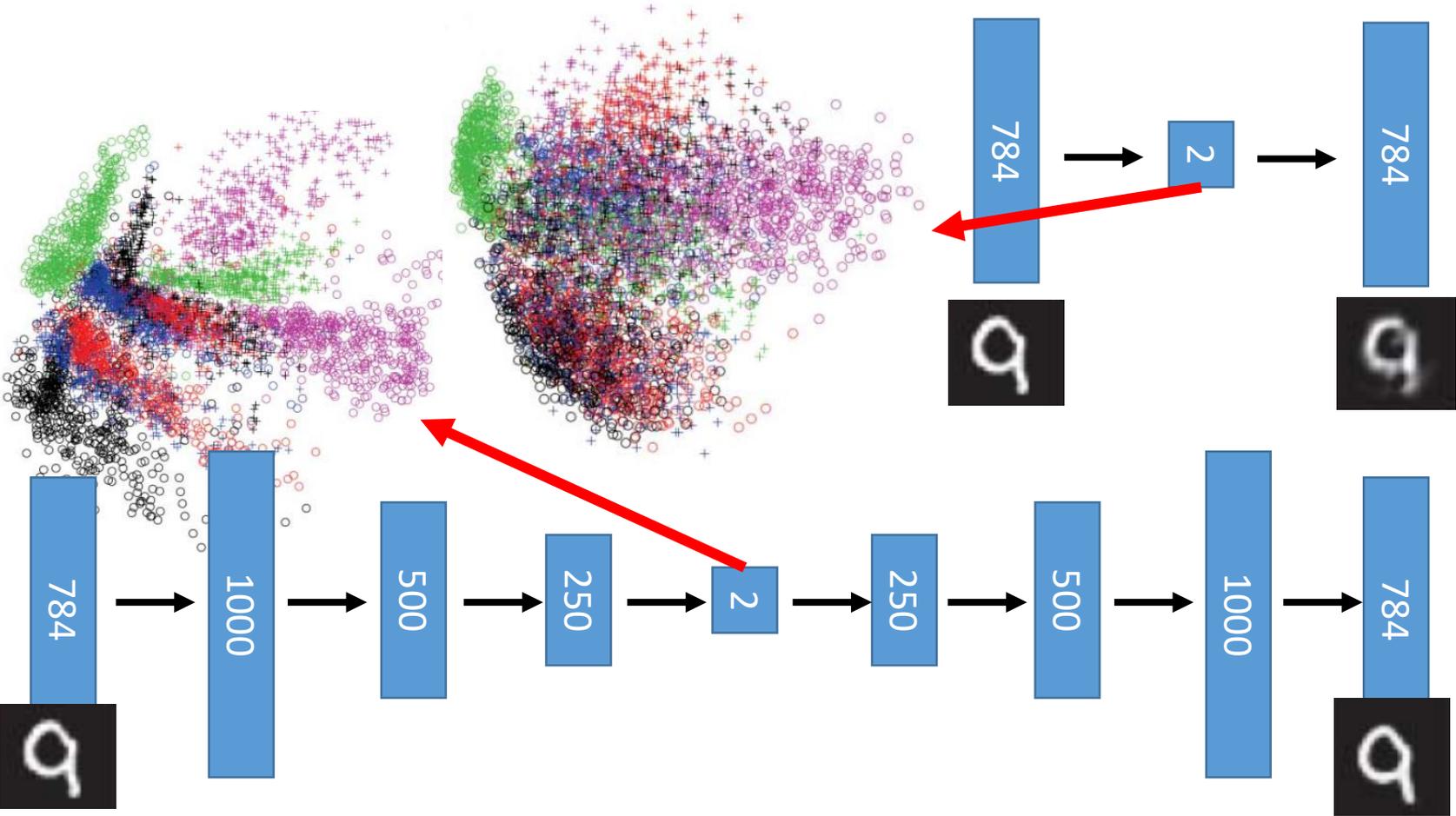


Deep Autoencoder

244



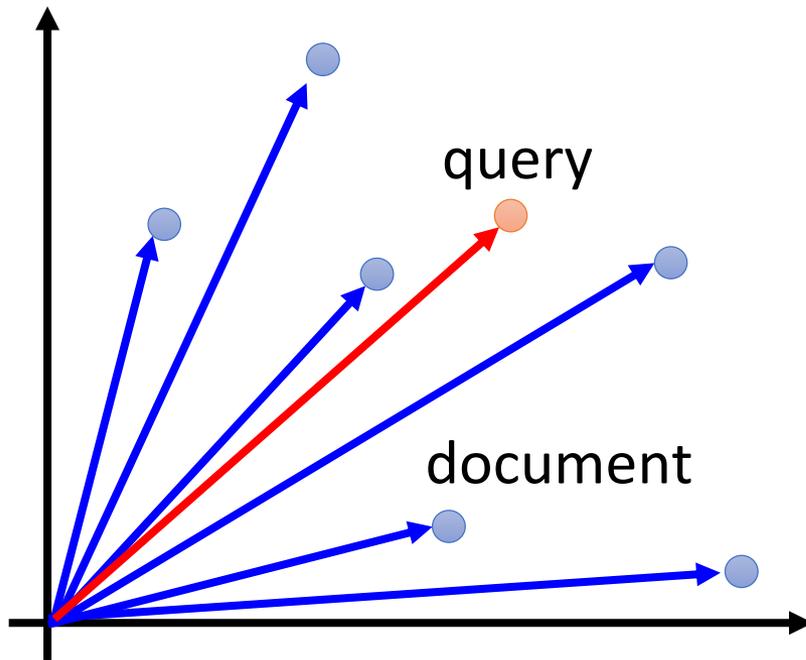
Feature Representation



Auto-encoder – Text Retrieval

246

Vector Space Model



Bag-of-words

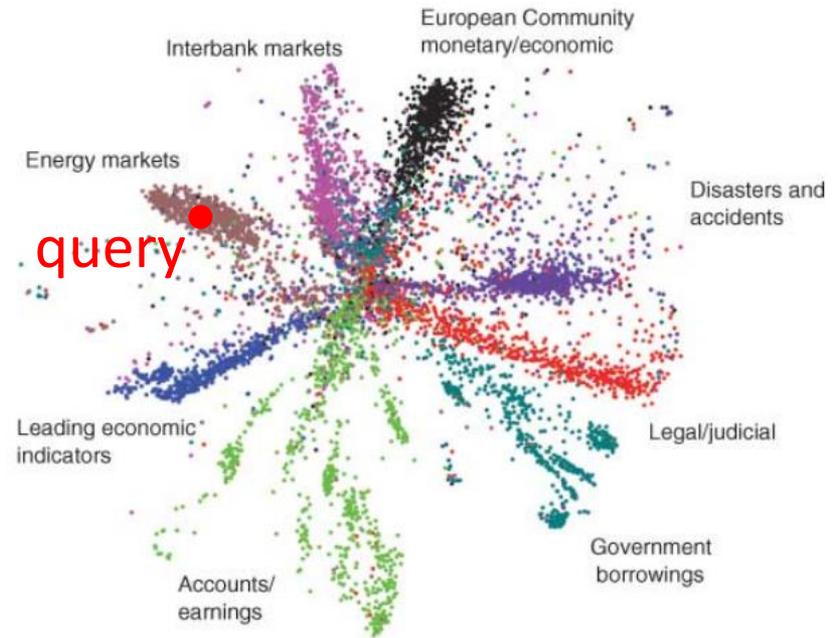
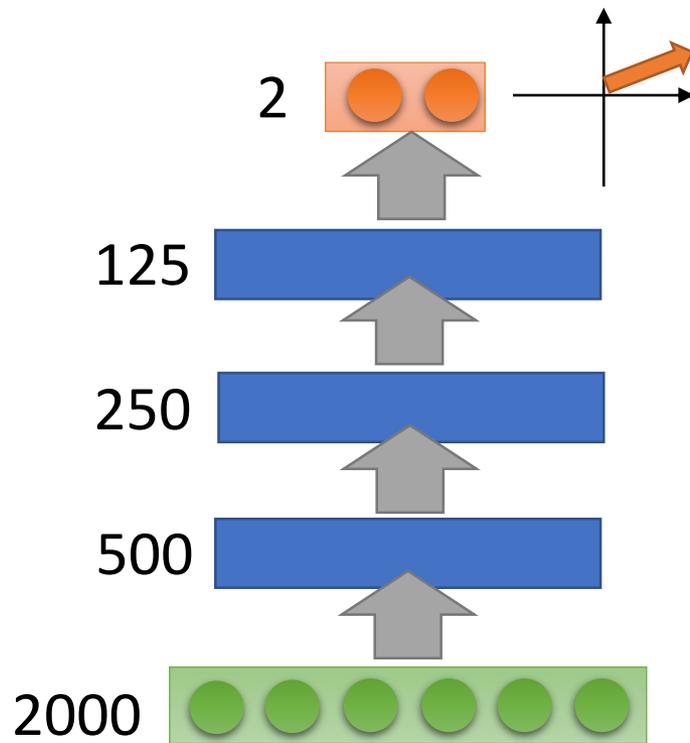
word string:
"This is an apple"

this	●	1
is	●	1
a	●	0
an	●	1
apple	●	1
pen	●	0
⋮		

Semantics are not considered

Autoencoder – Text Retrieval

247



The documents talking about the same thing will have close code

Bag-of-words (document or query)

Autoencoder – Similar Image Retrieval

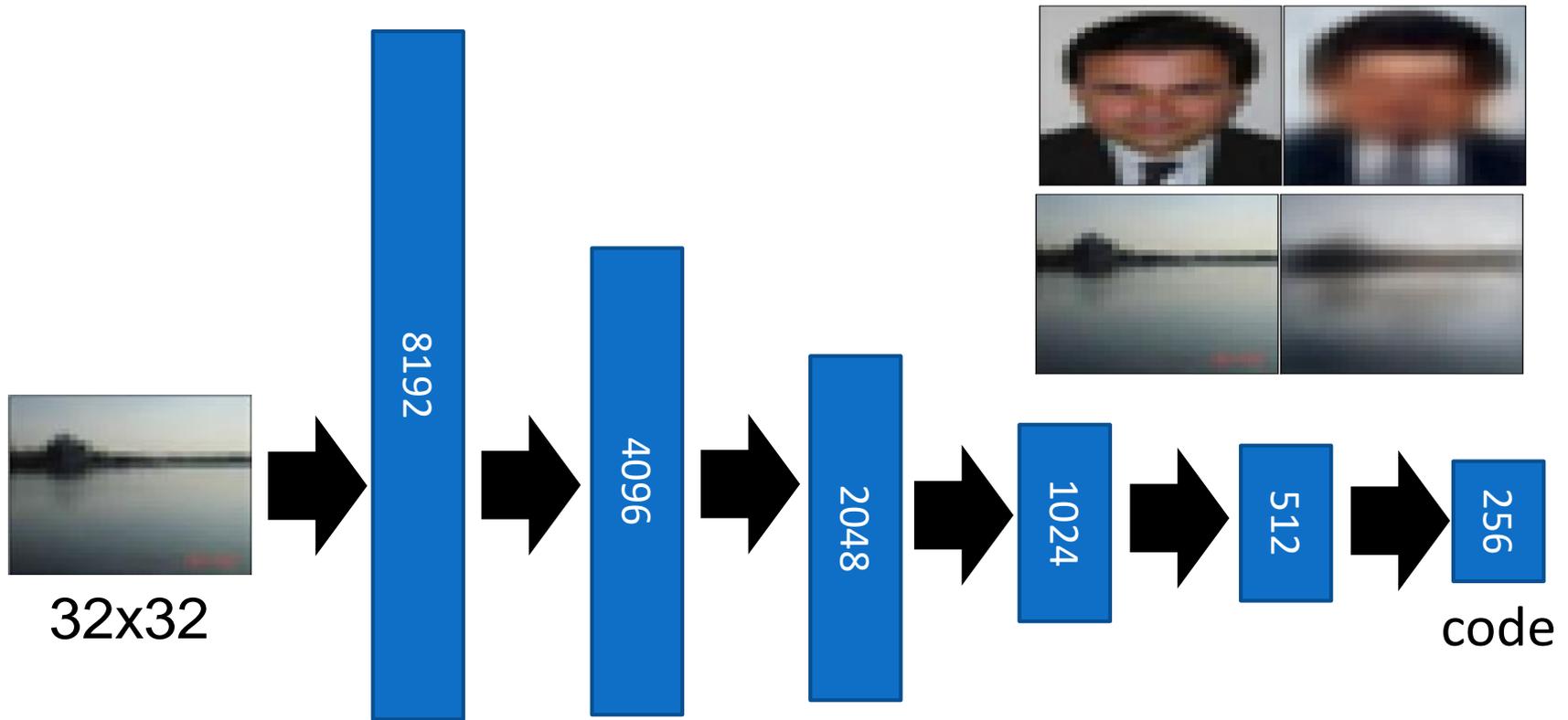
248

- Retrieved using Euclidean distance in pixel intensity space



Autoencoder – Similar Image Retrieval

249

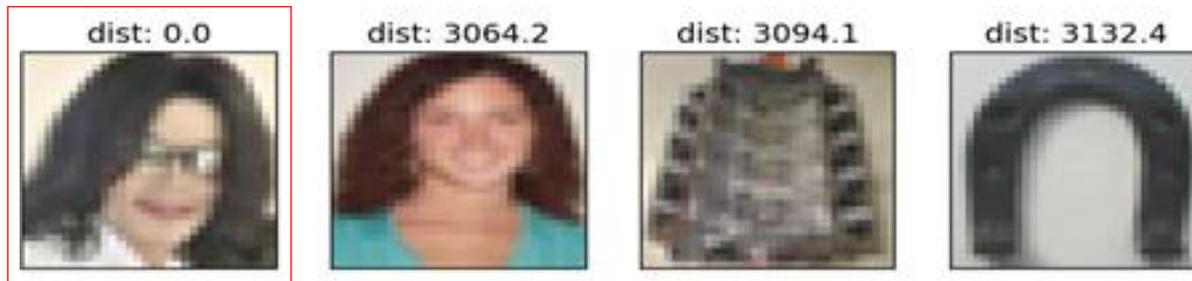


(crawl millions of images from the Internet)

Autoencoder – Similar Image Retrieval

250

- Images retrieved using Euclidean distance in pixel intensity space



- Images retrieved using 256 codes

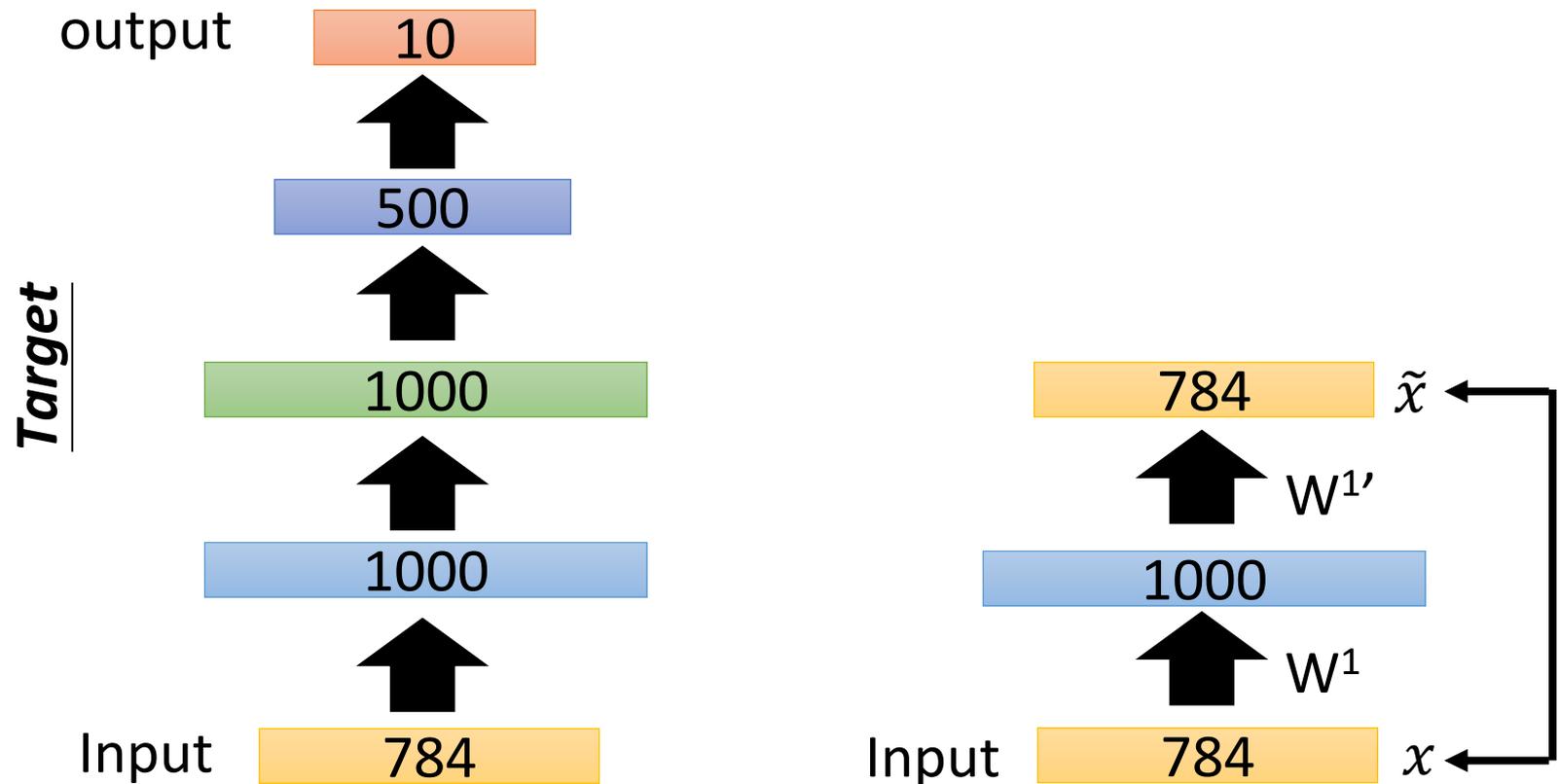


Learning the useful latent factors

Autoencoder for DNN Pre-Training

251

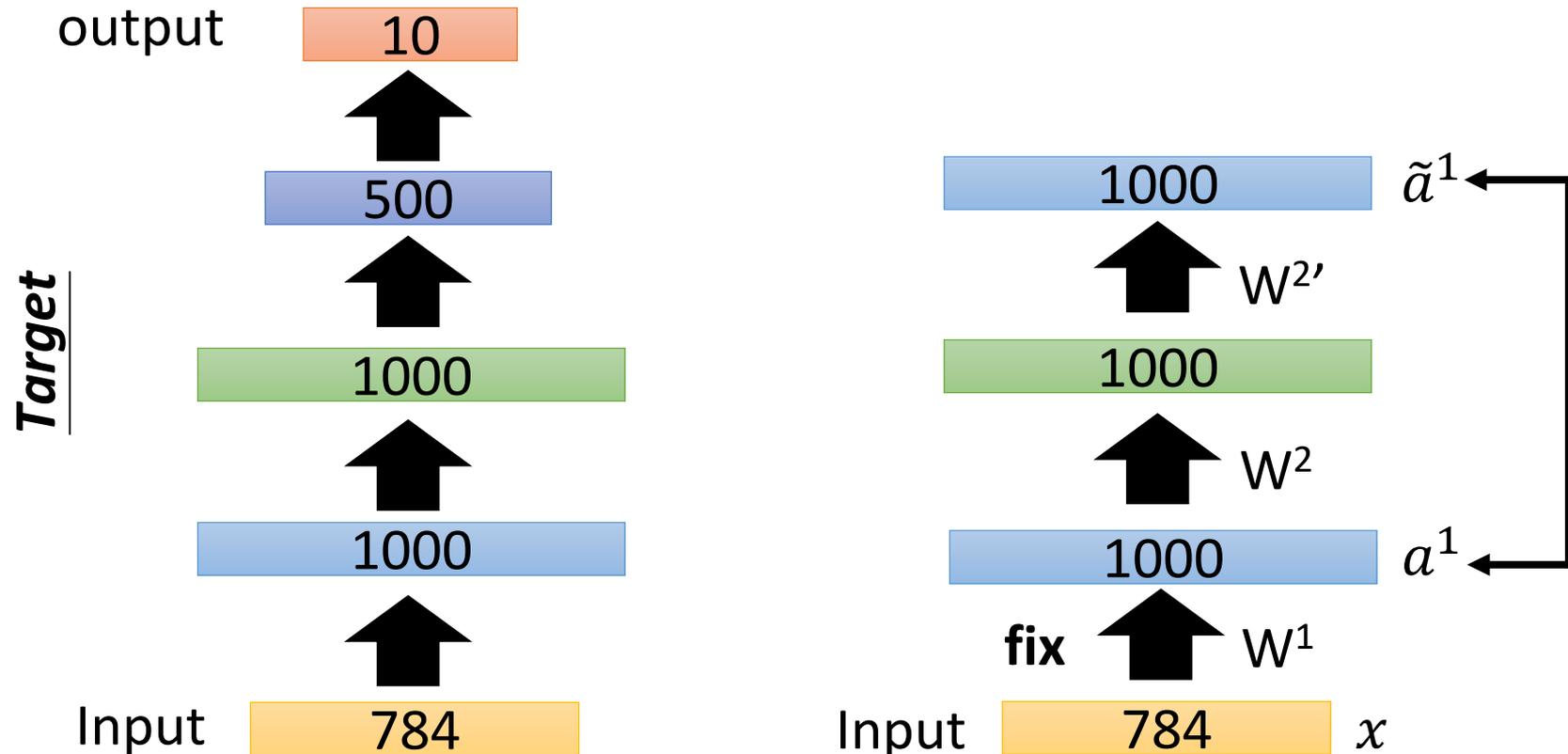
- Greedy layer-wise pre-training *again*



Autoencoder for DNN Pre-Training

252

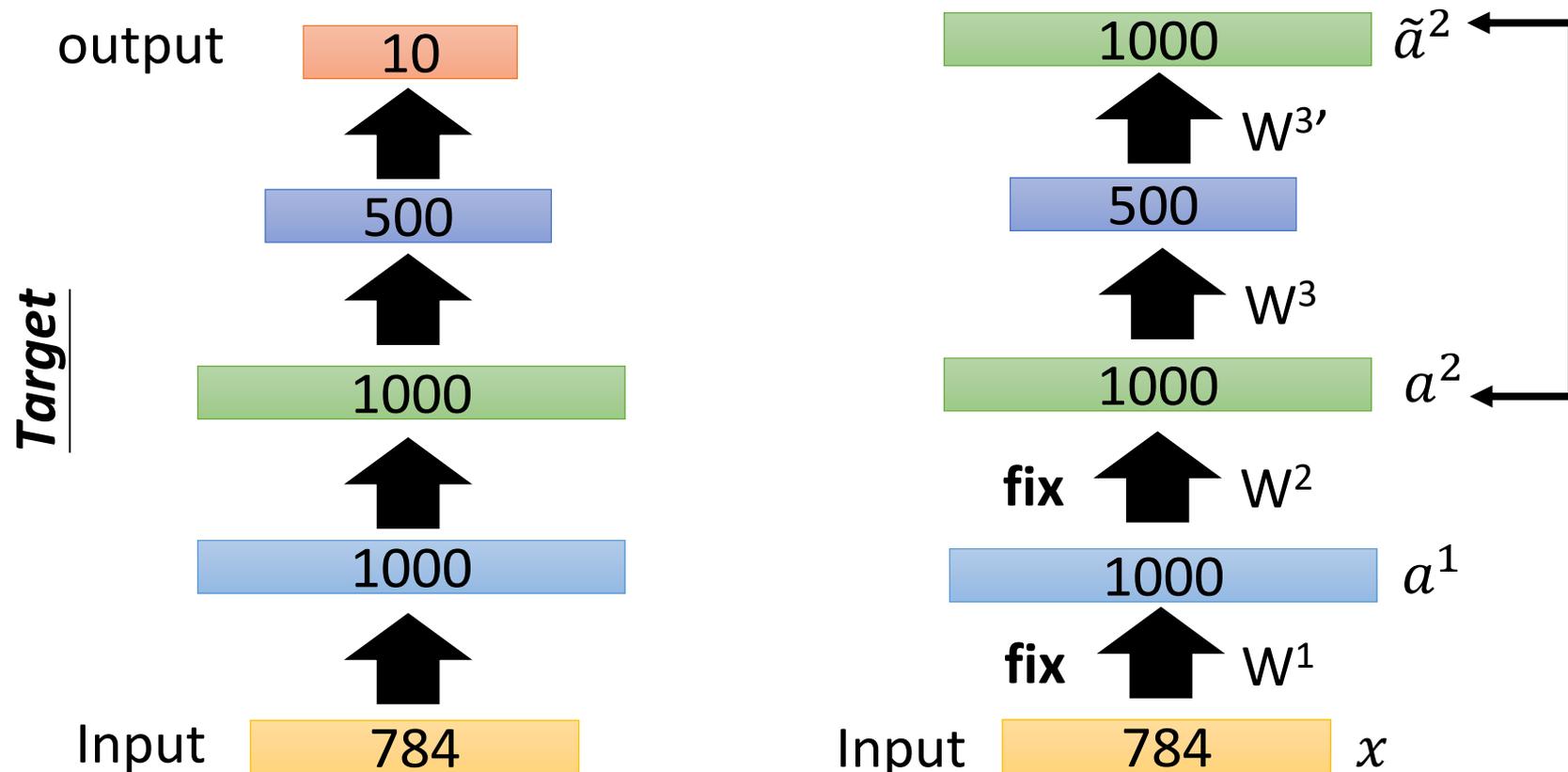
- Greedy layer-wise pre-training *again*



Autoencoder for DNN Pre-Training

253

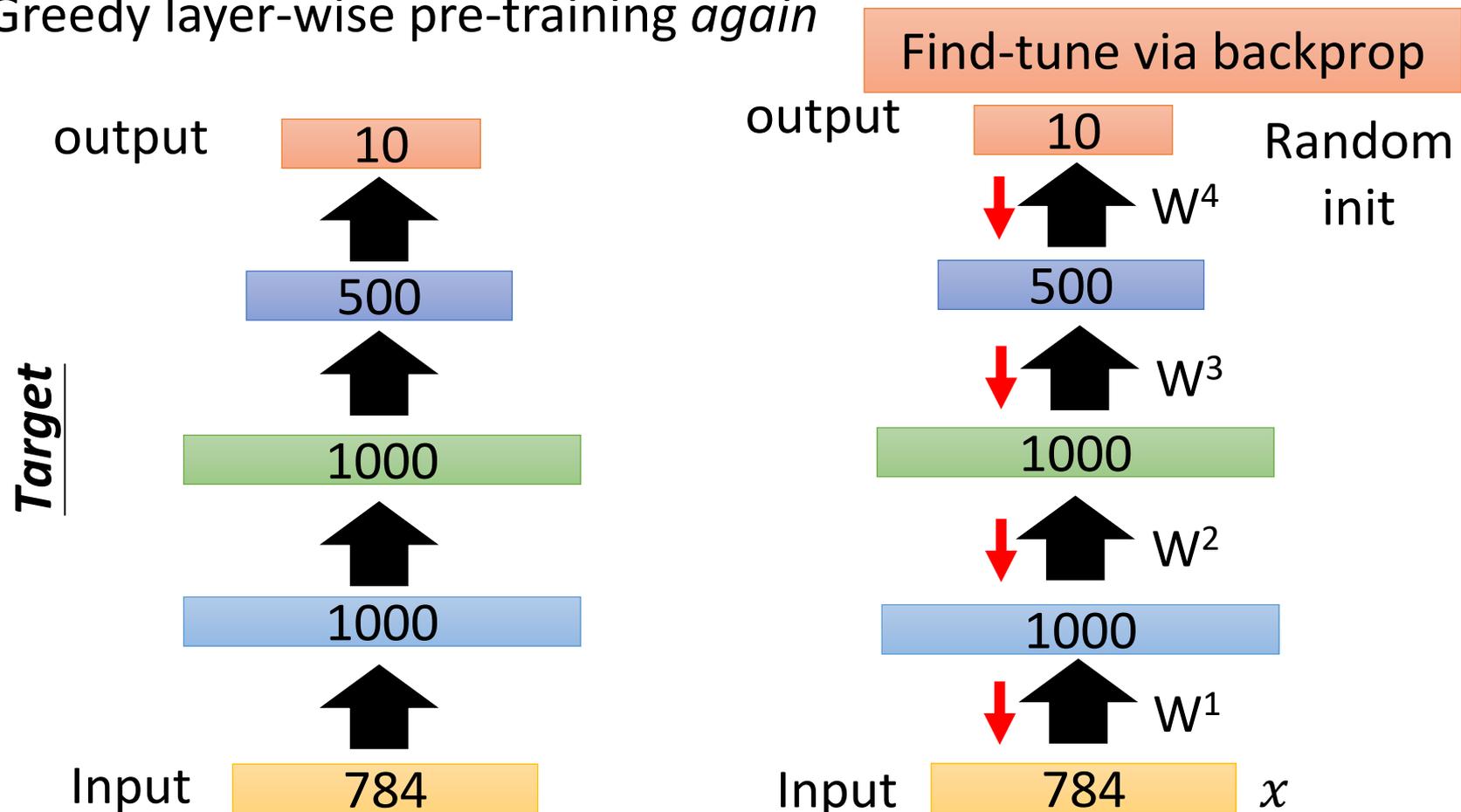
- Greedy layer-wise pre-training *again*



Autoencoder for DNN Pre-Training

254

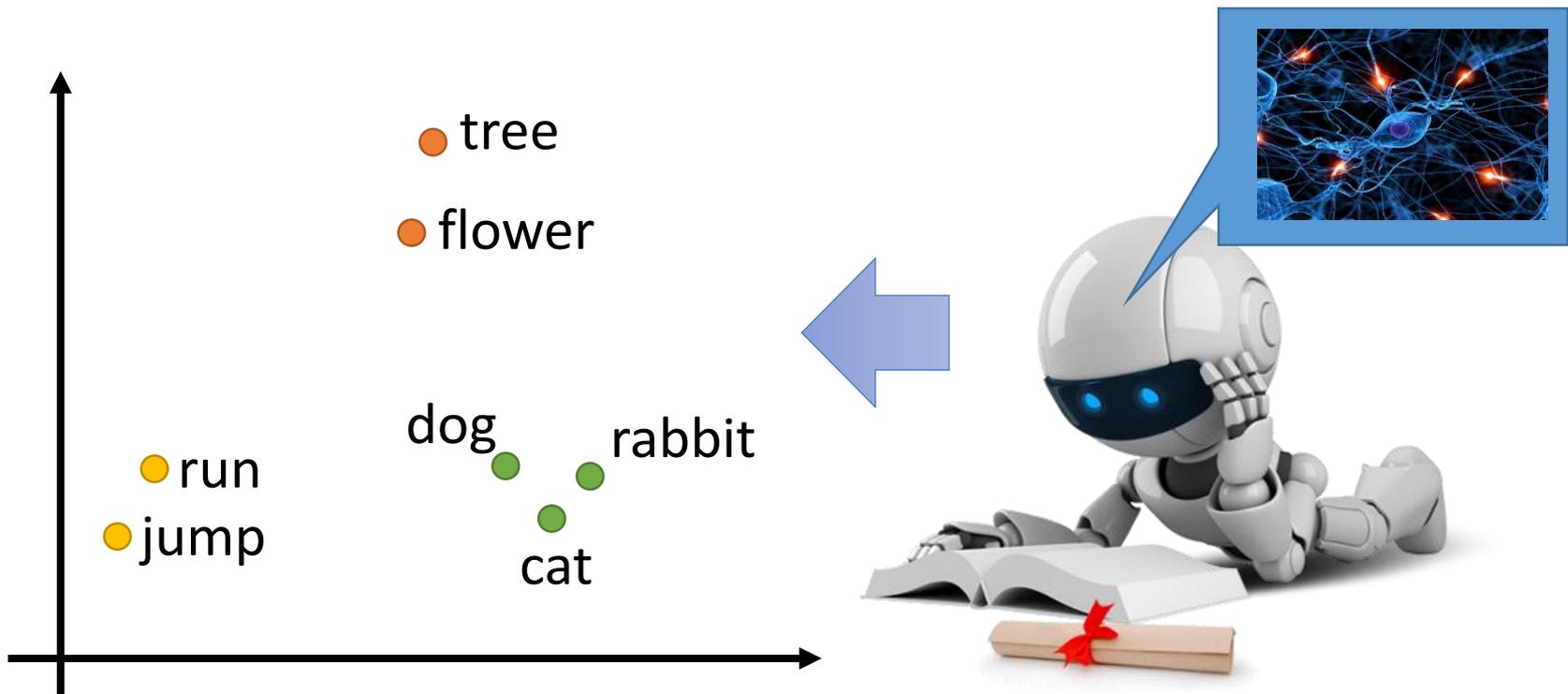
- Greedy layer-wise pre-training *again*



Word Vector/Embedding

255

- Machine learn the meaning of words from reading a lot of documents without supervision



Word Embedding

256

- Machine learn the meaning of words from reading a lot of documents without supervision
- A word can be understood by its context

蔡英文、馬英九 are something very similar

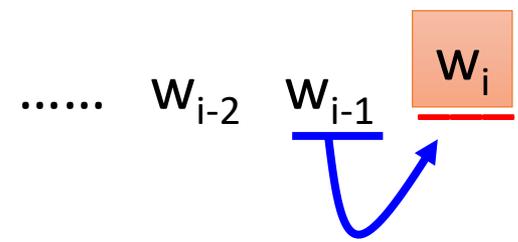
You shall know a word by the company it keeps

馬英九 520宣誓就職

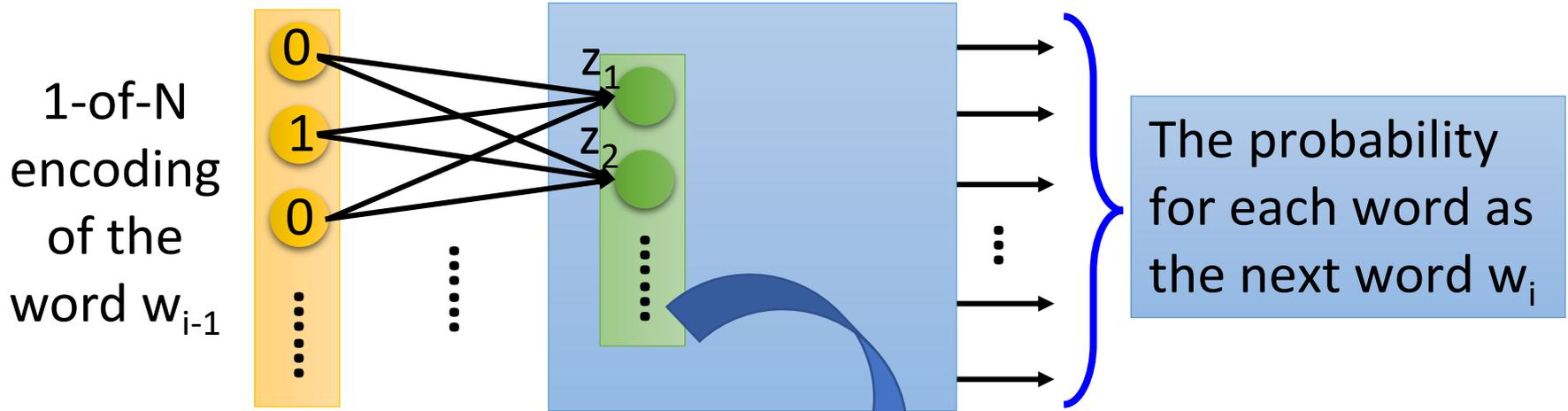
蔡英文 520宣誓就職



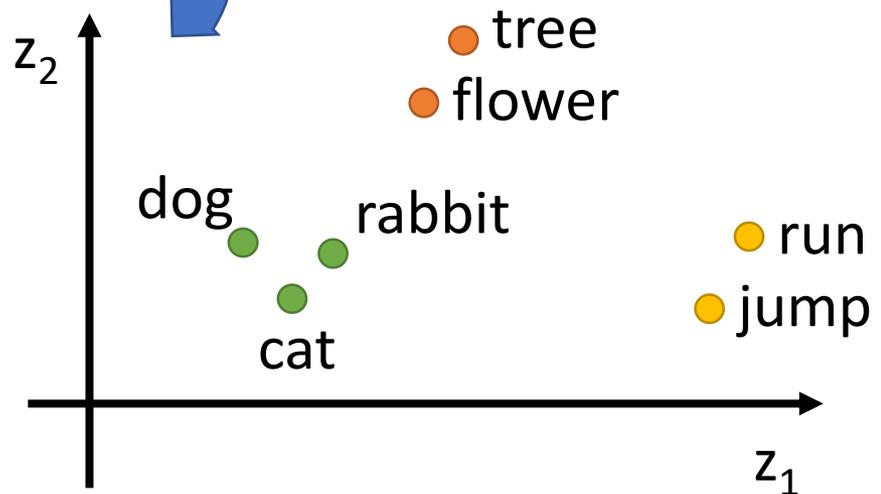
Prediction-Based



257



- Take out the input of the neurons in the first layer
- Use it to represent a word w
- Word vector, word embedding feature: $V(w)$

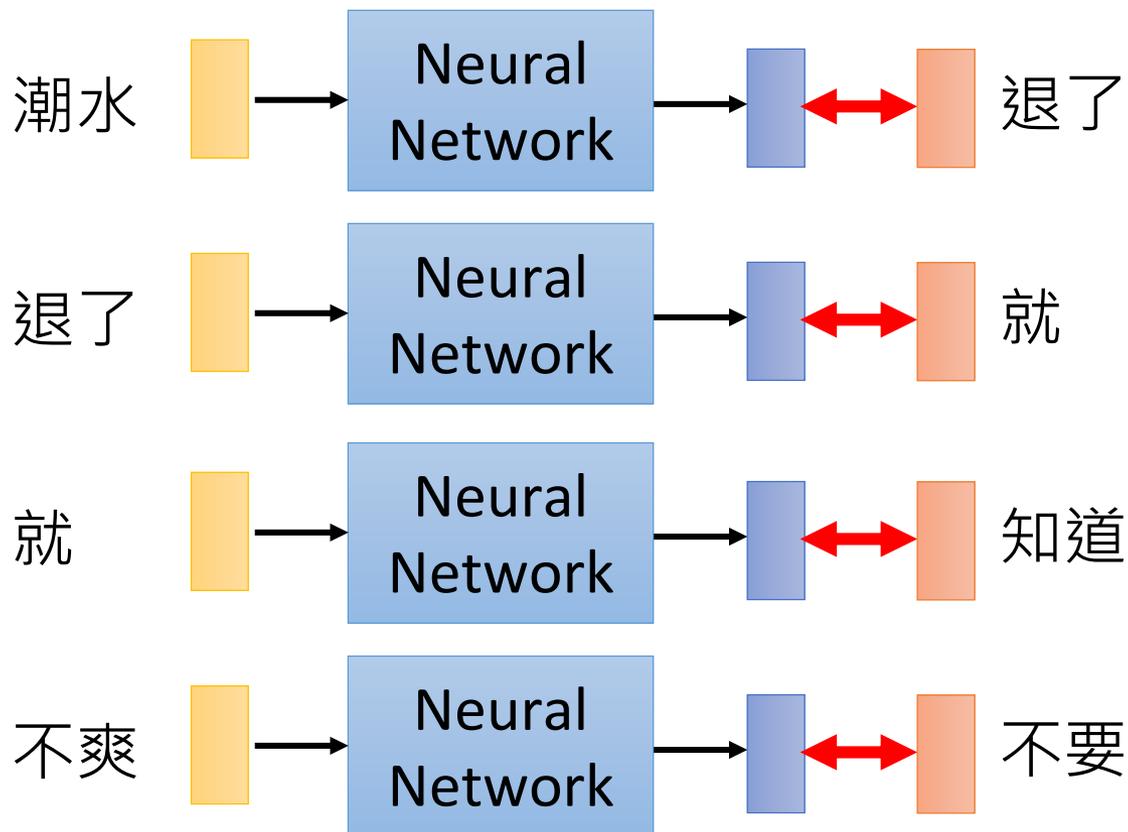


Prediction-Based

258

Collect data:

潮水 退了 就 知道 ...
不爽 不要 買 ...
公道價 八萬 一 ...
.....

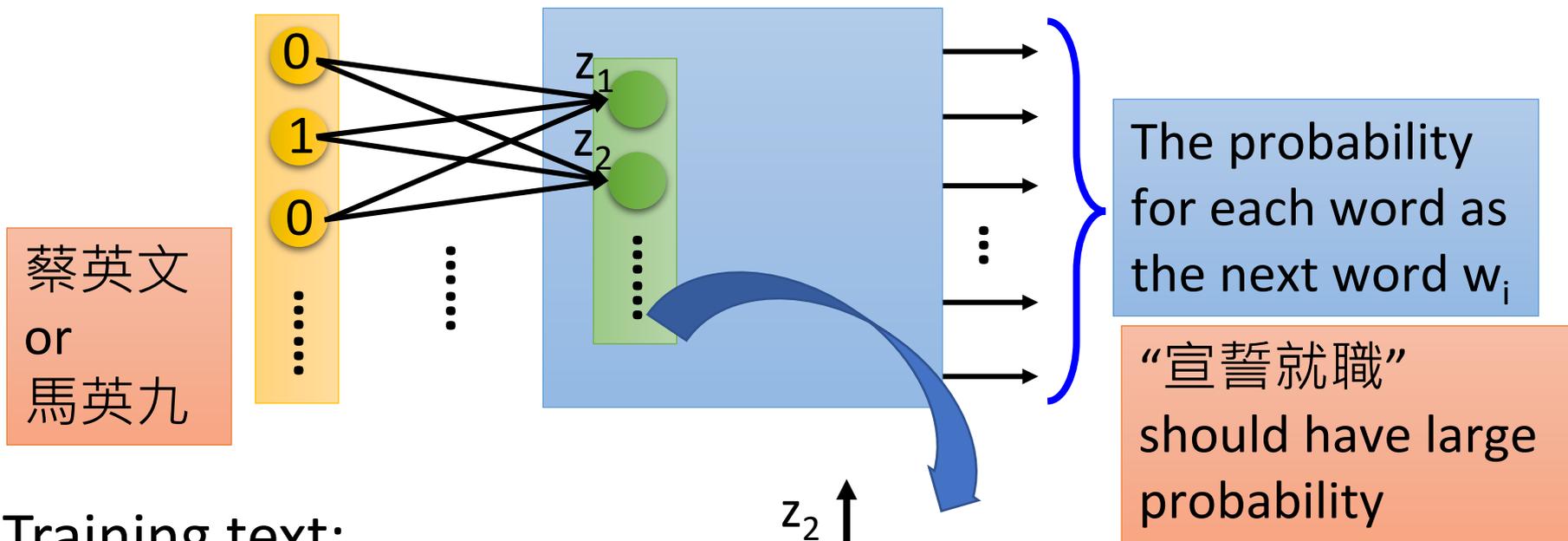


Minimizing cross entropy

Prediction-Based

You shall know a word by the company it keeps

259



Training text:

..... 蔡英文 宣誓就職

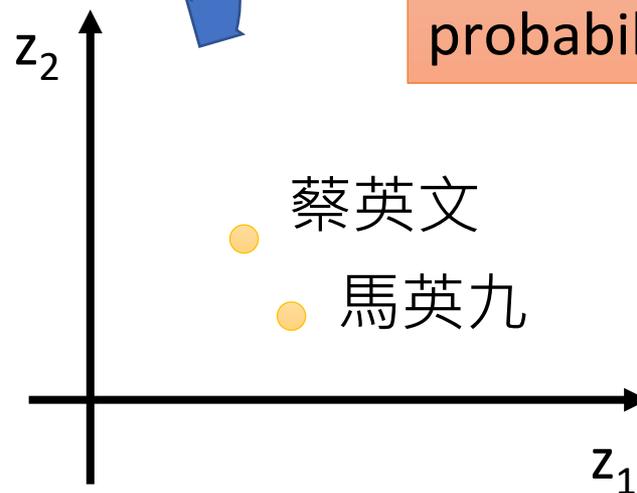
w_{i-1}

w_i

..... 馬英九 宣誓就職

w_{i-1}

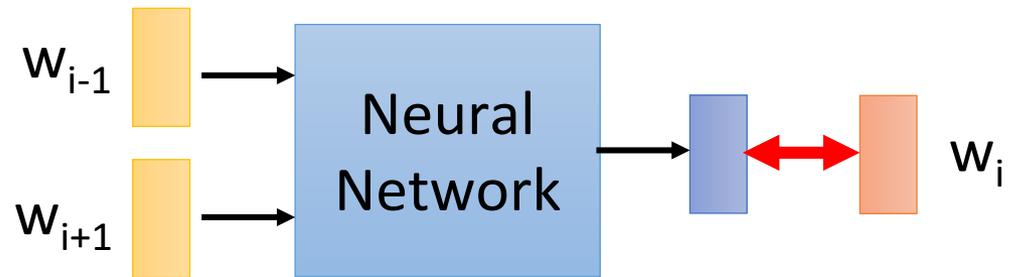
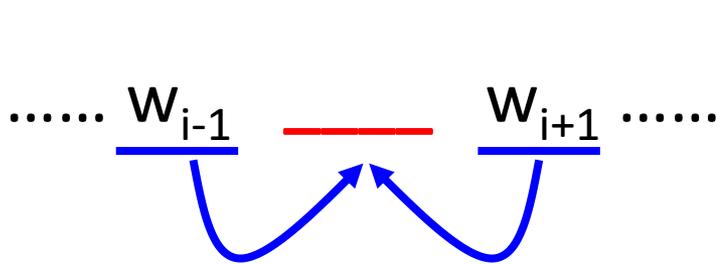
w_i



Various Architectures

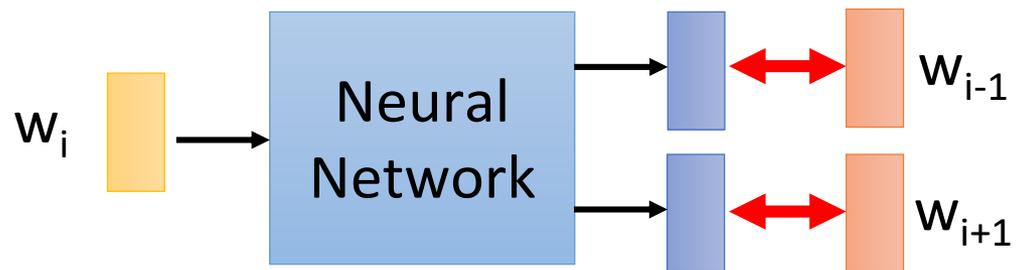
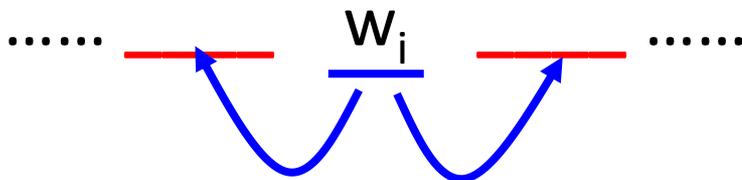
260

□ Continuous bag of word (CBOW) model



predicting the word given its context

□ Skip-gram

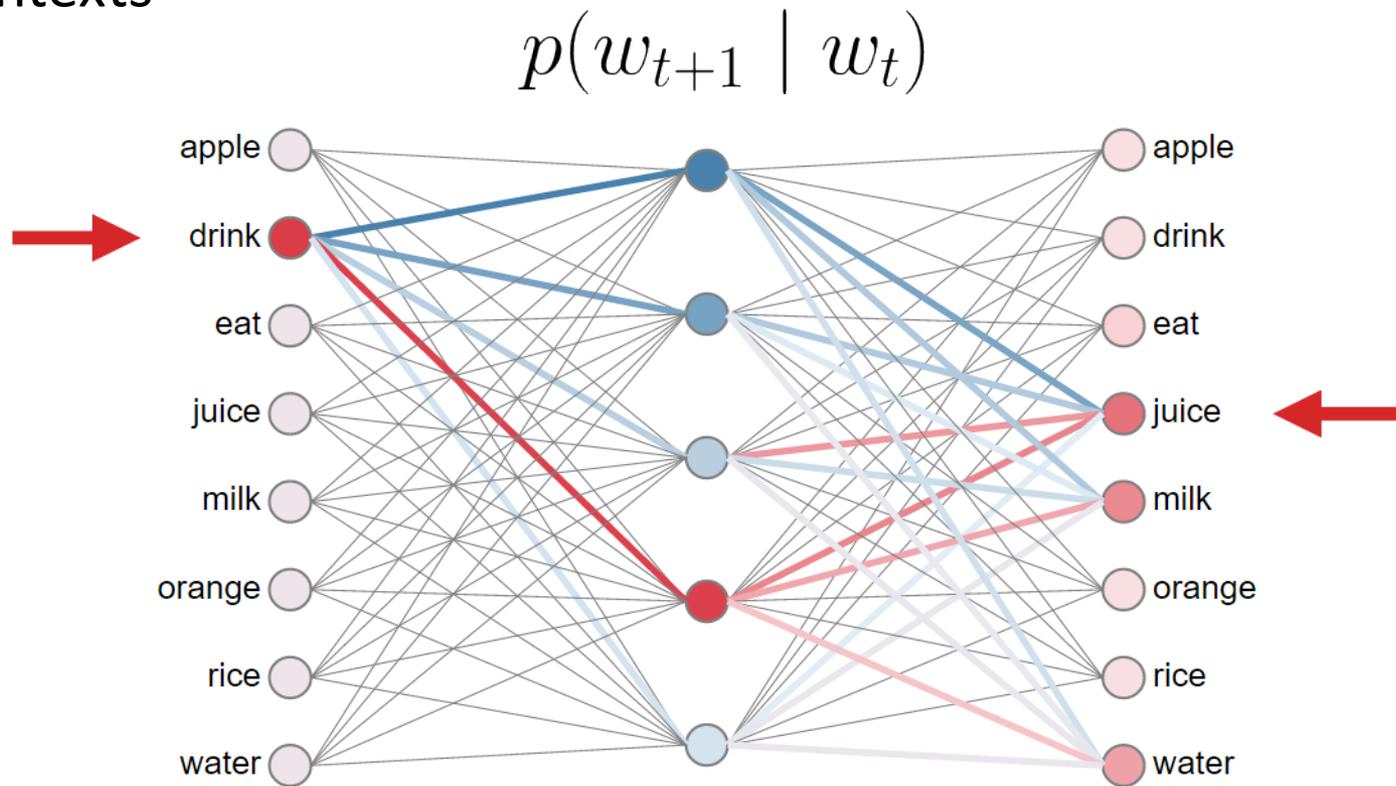


predicting the context given a word

Word2Vec LM

261

- Goal: predicting the next words given the proceeding contexts

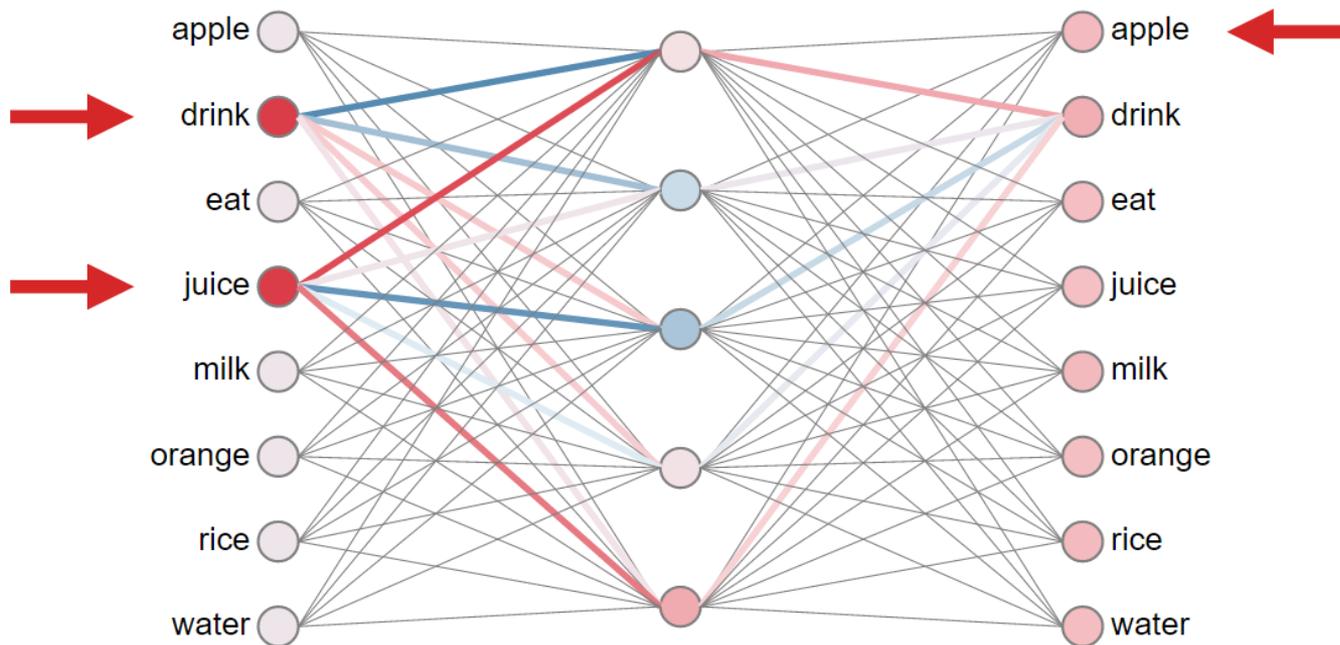


Word2Vec CBOW

262

- Goal: predicting the target word given the surrounding words

$$p(w_t \mid w_{t-m}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+m})$$

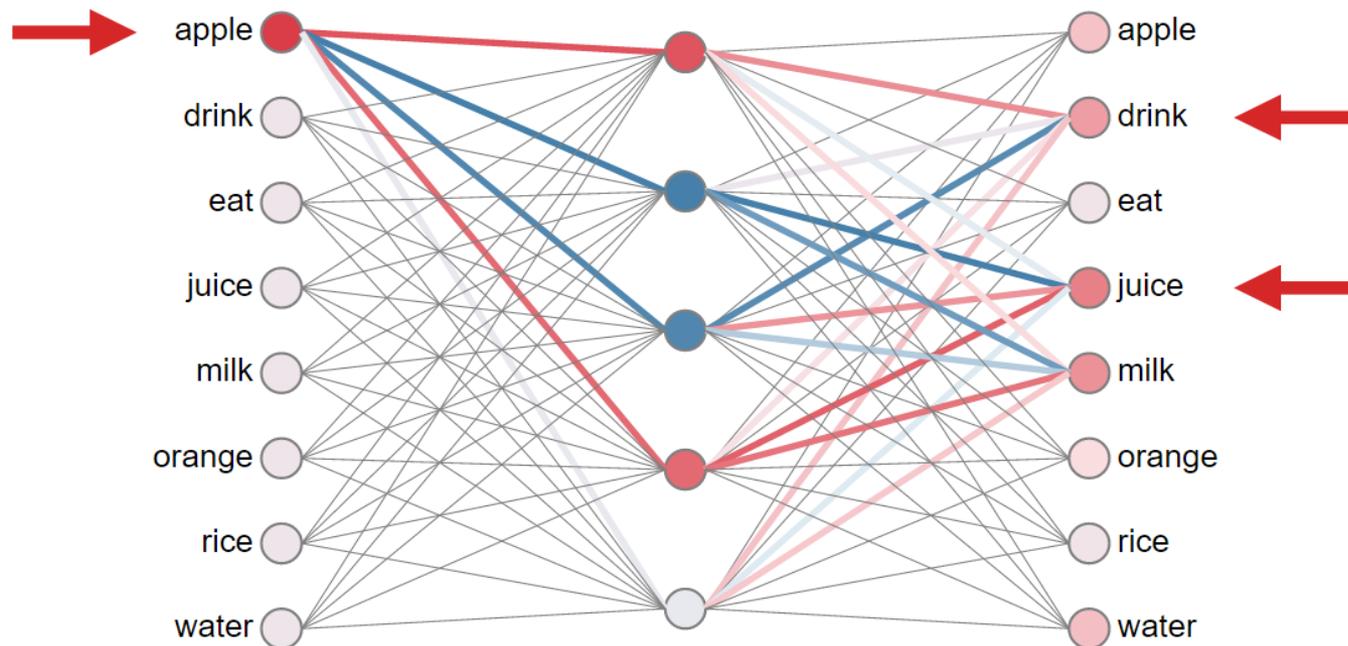


Word2Vec Skip-Gram

263

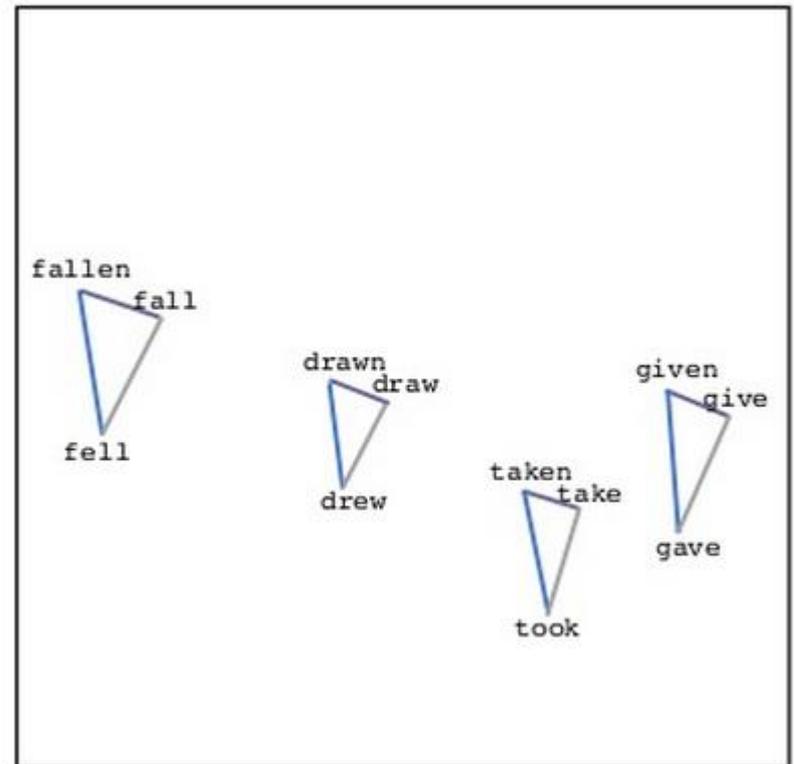
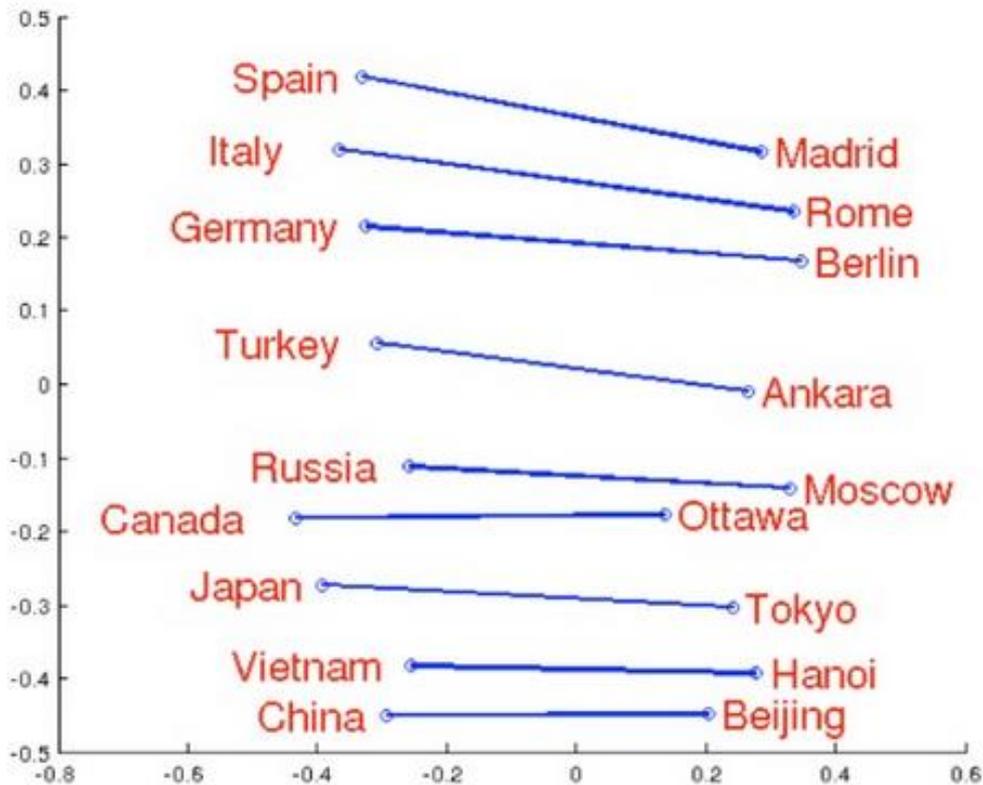
- Skip-gram training data:

apple | drink^juice,orange | eat^apple,rice | drink^juice,juice | drink^milk,milk
k | drink^rice,water | drink^milk,juice | orange^apple,juice | apple^drink,milk
| rice^drink,drink | milk^water,drink | water^juice,drink | juice^water



Word Embedding

264



Word Embedding

265

□ Characteristics

$$V(\textit{hotter}) - V(\textit{hot}) \approx V(\textit{bigger}) - V(\textit{big})$$

$$V(\textit{Rome}) - V(\textit{Italy}) \approx V(\textit{Berlin}) - V(\textit{Germany})$$

$$V(\textit{king}) - V(\textit{queen}) \approx V(\textit{uncle}) - V(\textit{aunt})$$

□ Solving analogies

$$\textit{Rome} : \textit{Italy} = \textit{Berlin} : ?$$

Compute $V(\textit{Berlin}) - V(\textit{Rome}) + V(\textit{Italy})$

Find the word w with the closest $V(w)$

$$\begin{aligned} & V(\textit{Germany}) \\ & \approx V(\textit{Berlin}) - V(\textit{Rome}) + V(\textit{Italy}) \end{aligned}$$

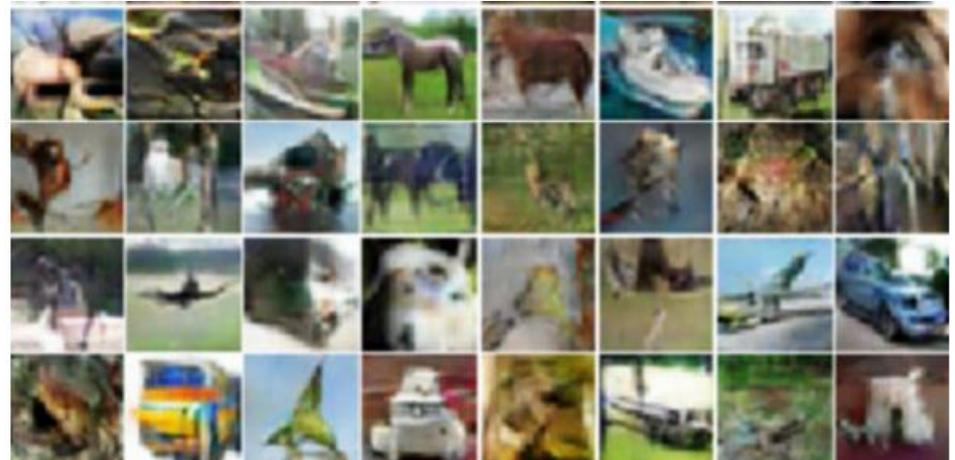
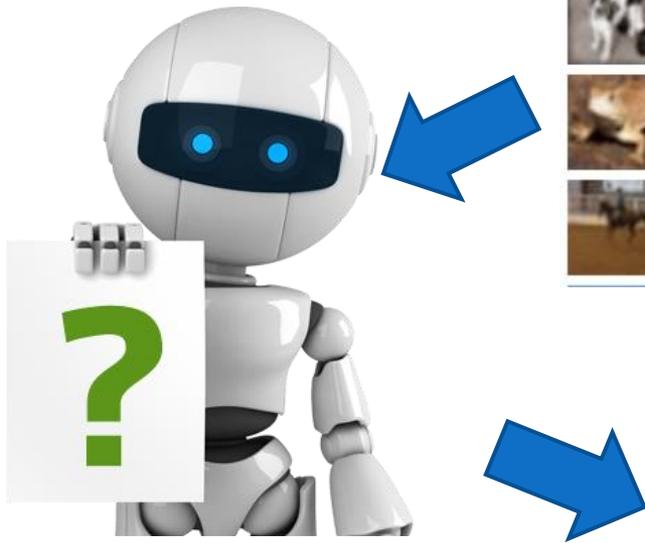
Outline

266

- Semi-Supervised Learning
- Transfer Learning
- **Unsupervised Learning**
 - 化繁為簡 Representation Learning
 - 無中生有 Generative Model
- Reinforcement Learning

Creation

267



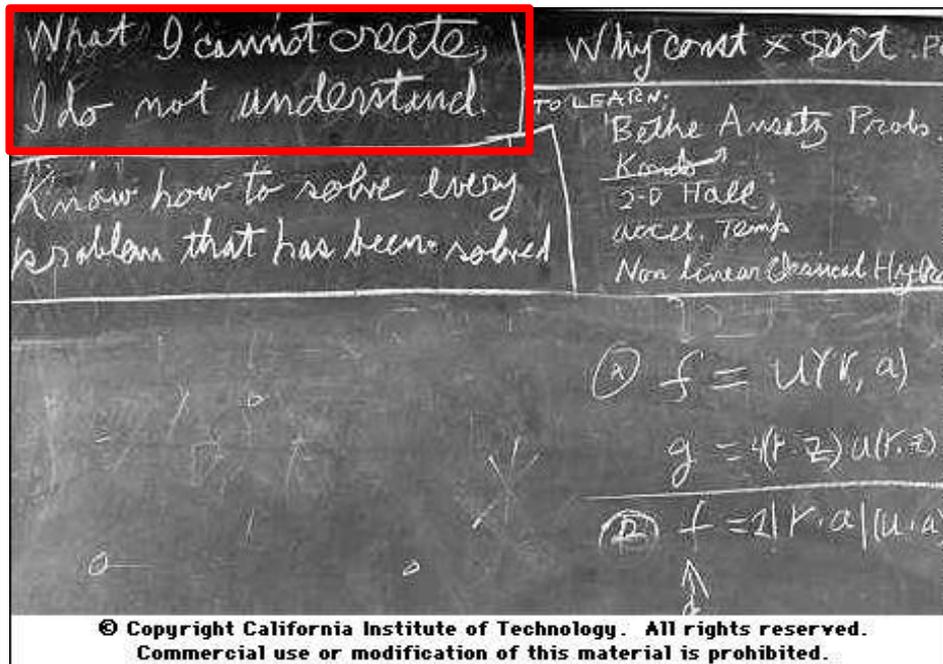
Draw something!

Creation

268

□ Generative Models

- <https://openai.com/blog/generative-models/>



What I cannot create,
I do not understand.

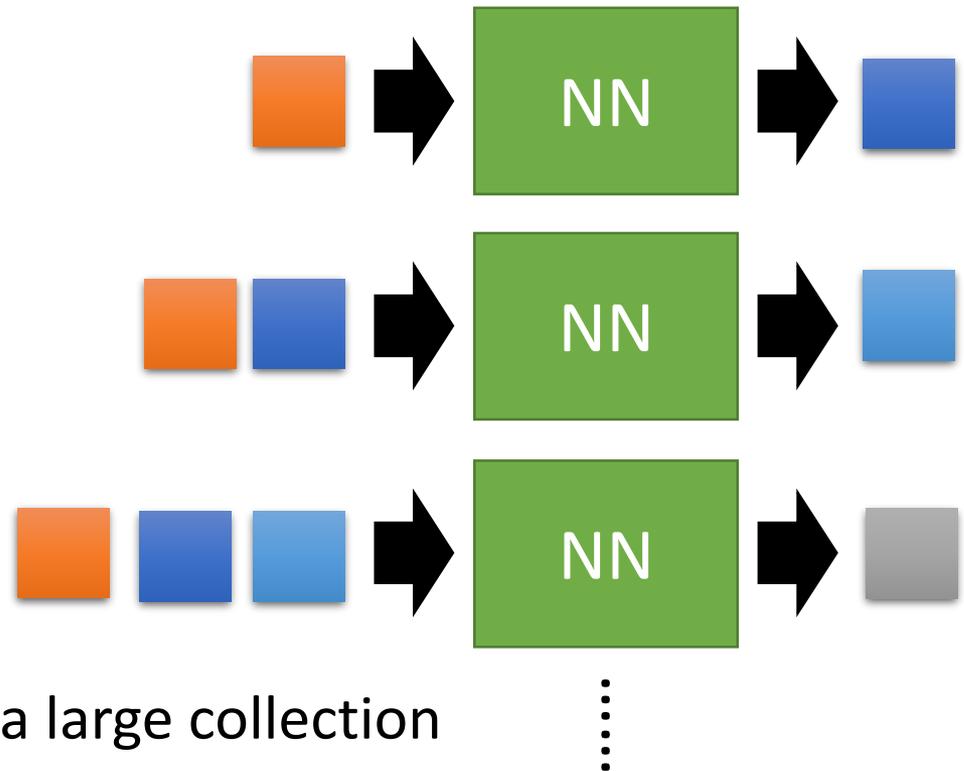
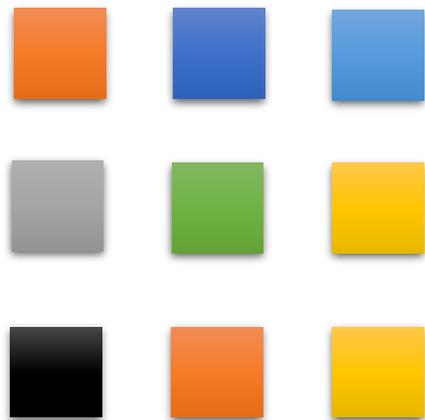
Richard Feynman

PixelRNN

269

- To create an image, generating a pixel each time

E.g. 3 x 3 images



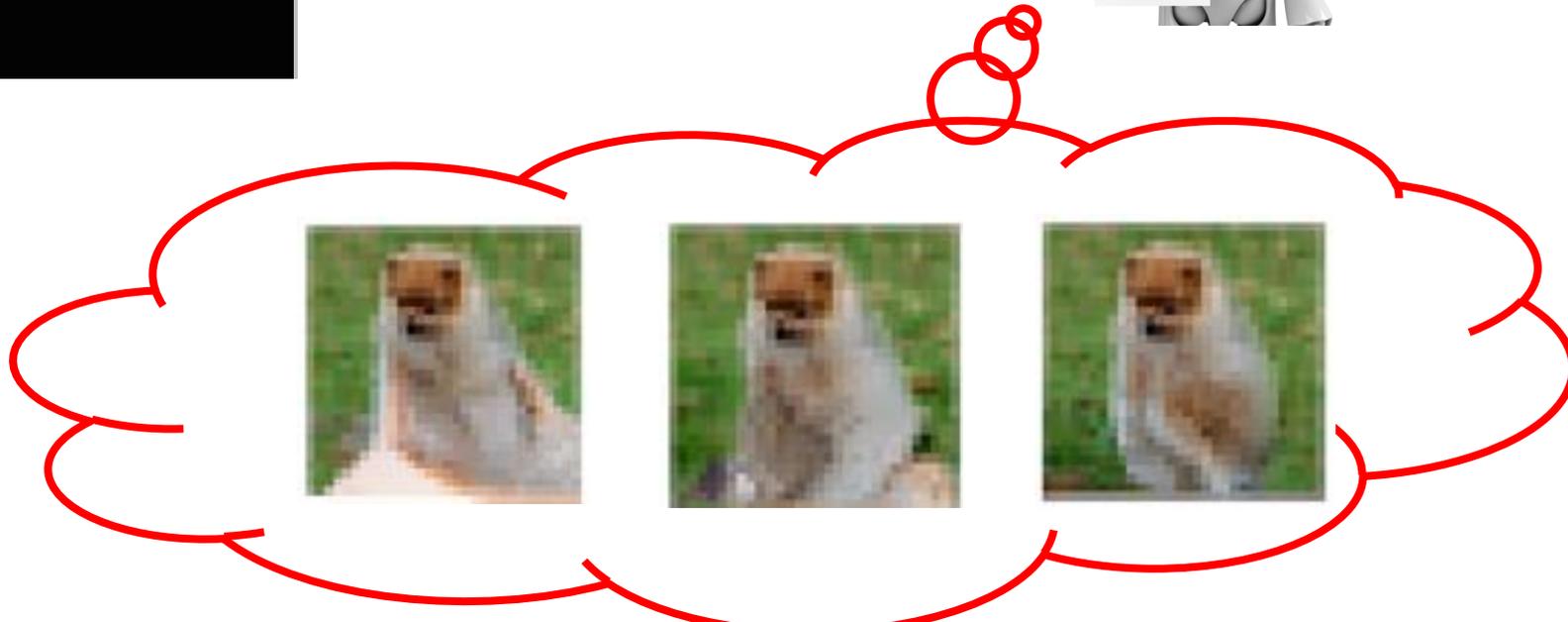
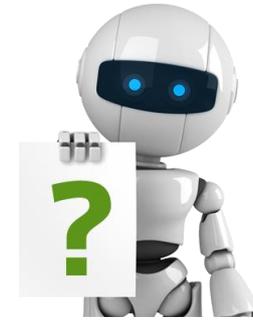
Can be trained just with a large collection of images without any annotation

PixelRNN

270

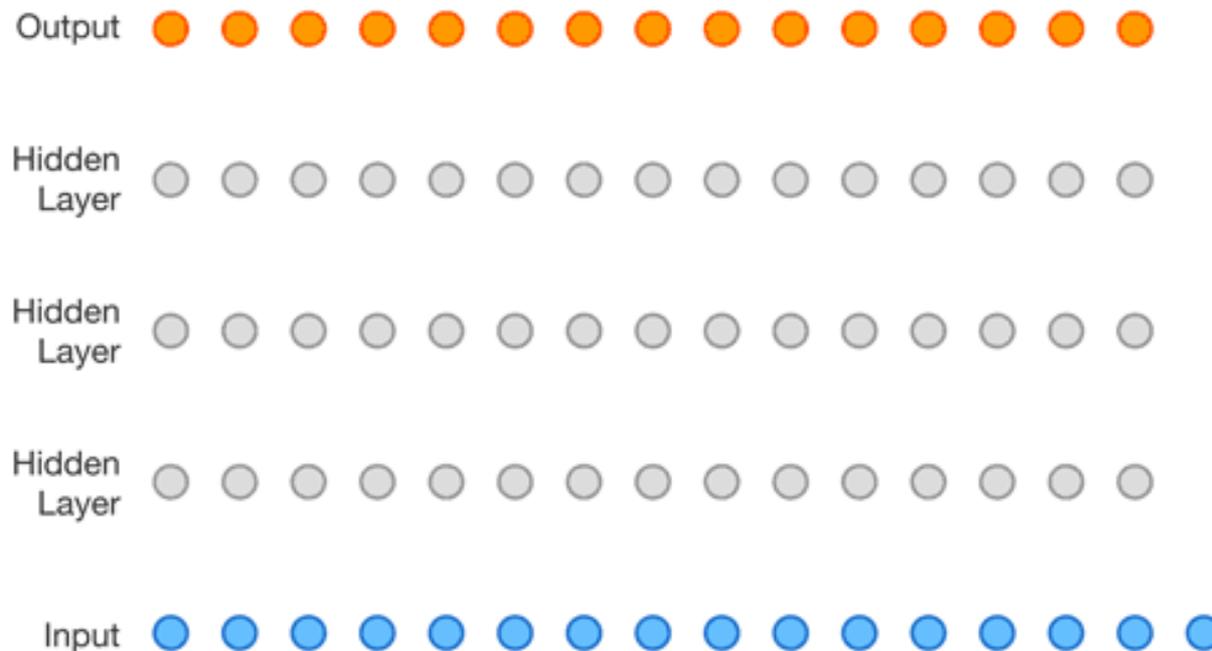


Real
World



PixelRNN – beyond Image

271



Audio: Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, Koray Kavukcuoglu, WaveNet: A Generative Model for Raw Audio, arXiv preprint, 2016

Video: Nal Kalchbrenner, Aaron van den Oord, Karen Simonyan, Ivo Danihelka, Oriol Vinyals, Alex Graves, Koray Kavukcuoglu, Video Pixel Networks , arXiv preprint, 2016

Generative Adversarial Network (GAN)

272

What are some recent and potentially upcoming breakthroughs in unsupervised learning?



Yann LeCun, Director of AI Research at Facebook and Professor at NYU

Written Jul 29 · Upvoted by Joaquin Quiñonero Candela, [Director Applied Machine Learning at Facebook](#) and Huang Xiao



Adversarial training is the coolest thing since sliced bread.

I've listed a bunch of relevant papers in a previous answer.

Expect more impressive results with this technique in the coming years.

What's missing at the moment is a good understanding of it so we can make it work reliably. It's very finicky. Sort of like ConvNet were in the 1990s, when I had the reputation of being the only person who could make them work (which wasn't true).

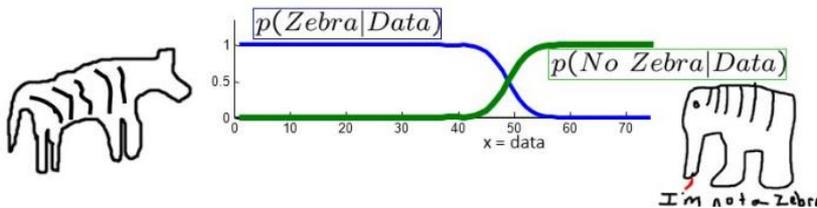
Ref: Generative Adversarial Networks, <http://arxiv.org/abs/1406.2661>

Discriminative v.s. Generative Models

273

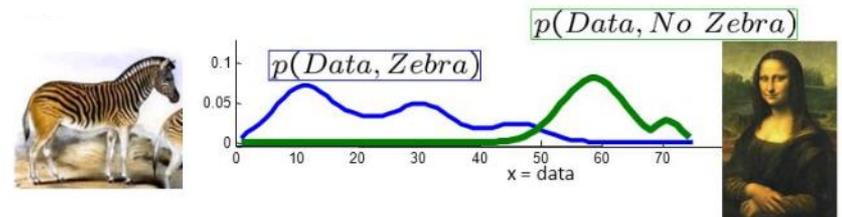
□ Discriminative

- learns a function that maps the input data (x) to some desired output class label (y)
 - directly learn the conditional distribution $P(y/x)$



□ Generative

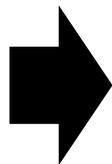
- tries to learn the joint probability of the input data and labels simultaneously, i.e. $P(x,y)$
 - can be converted to $P(y/x)$ for classification via Bayes rule



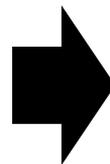
Advantage: generative models have the potential to understand and explain the underlying structure of the input data even when there are no labels

擬態的演化

274

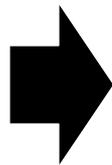


棕色



葉脈

蝴蝶不是棕色



蝴蝶沒有葉脈



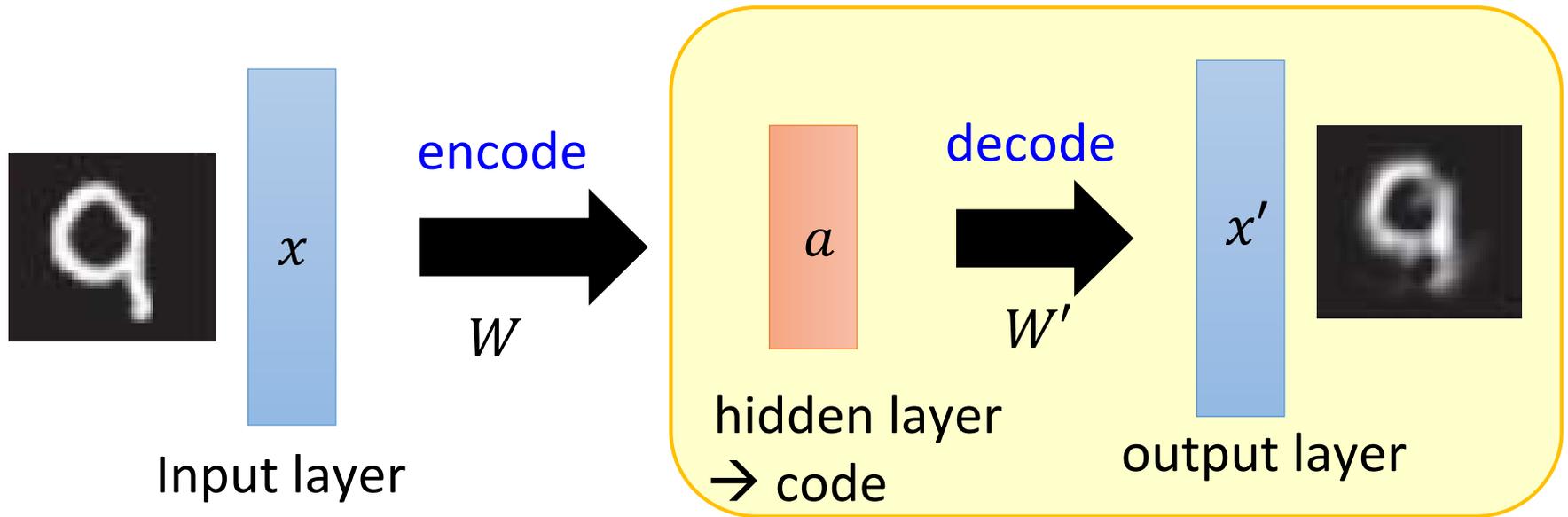
.....



Generator

275

- Decoder from autoencoder as generator

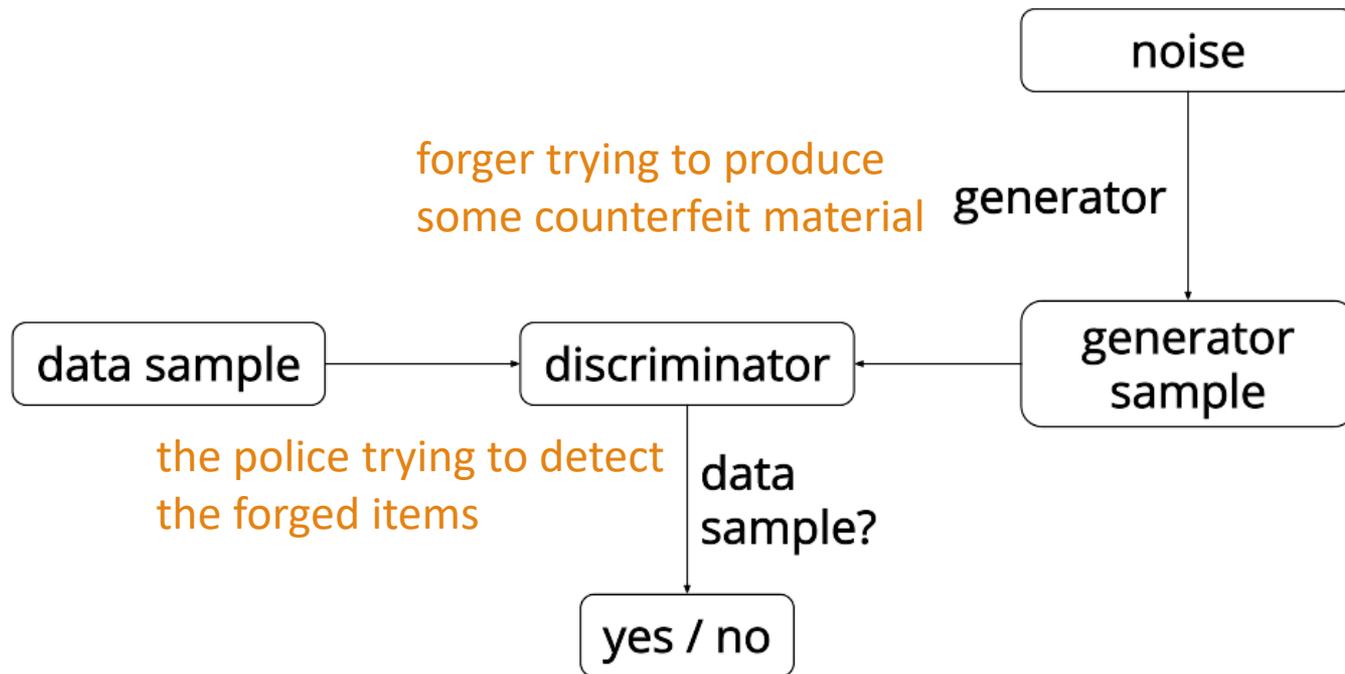


The generator is to generate the data from the code

Generative Adversarial Networks (GAN)

276

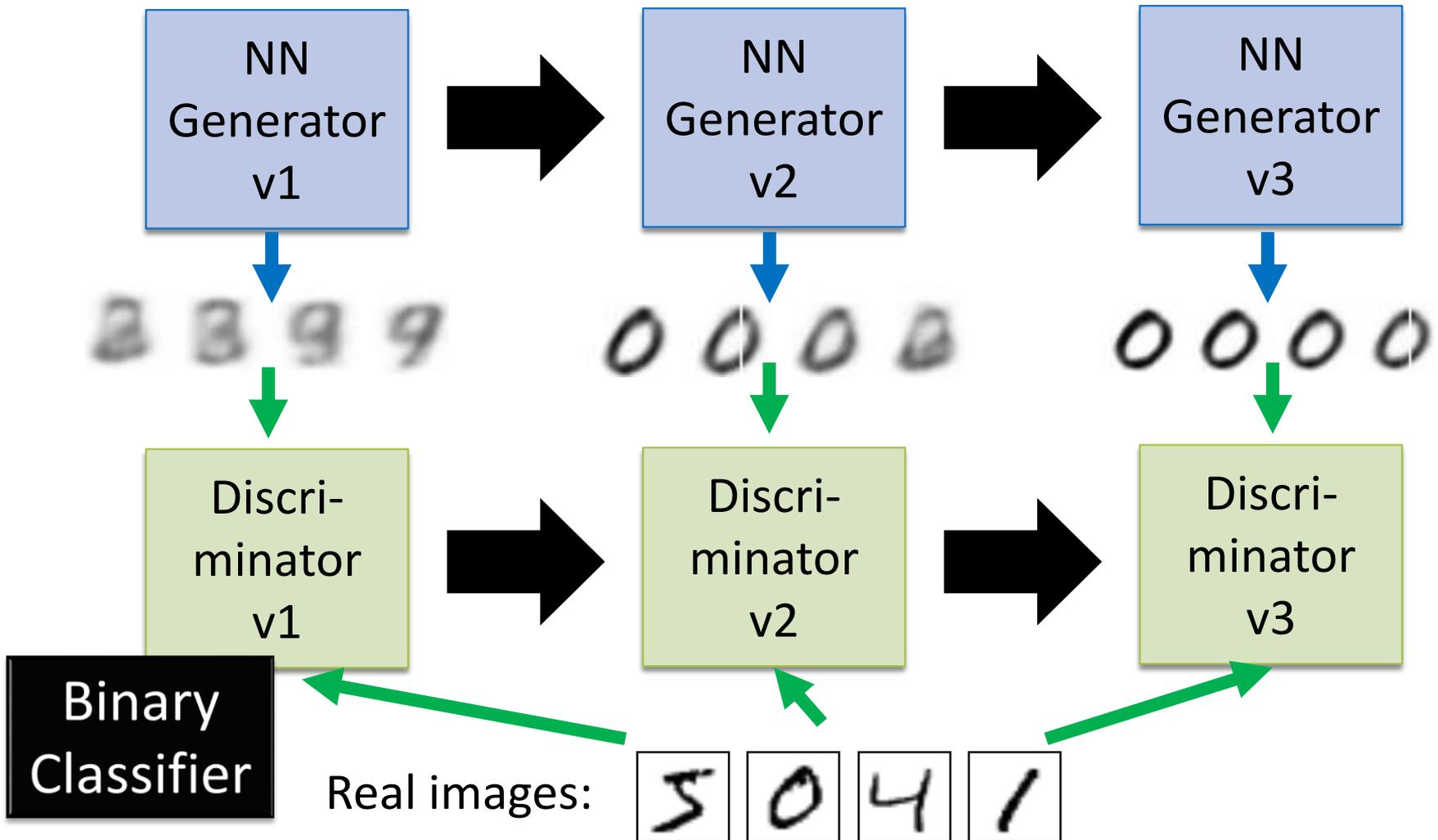
- Two competing neural networks: generator & discriminator



Training two networks jointly → the generator knows how to adapt its parameters in order to produce output data that can fool the discriminator

Generator Evolution

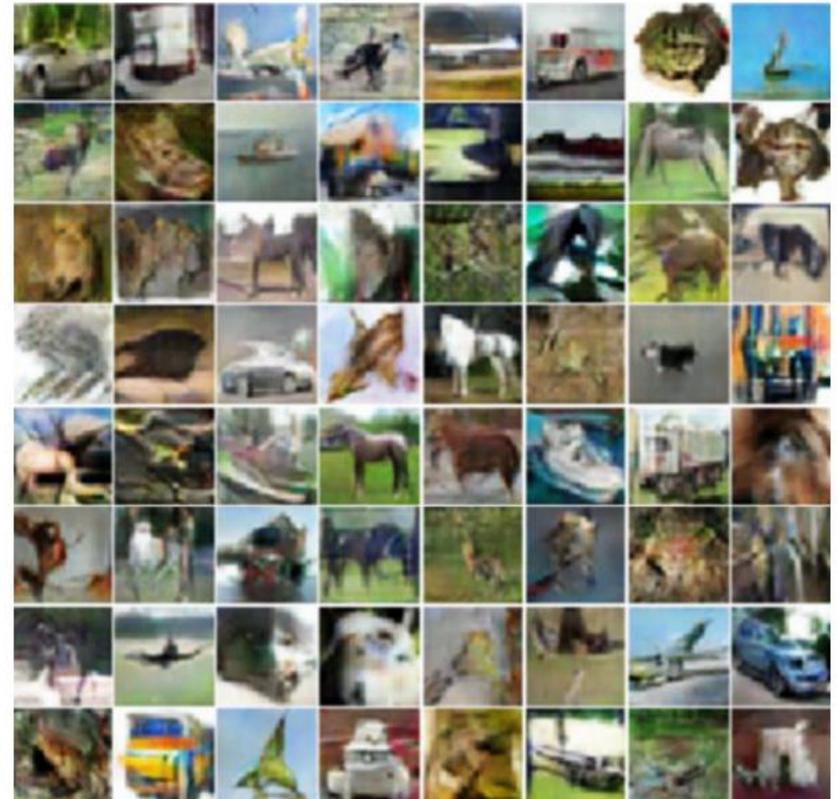
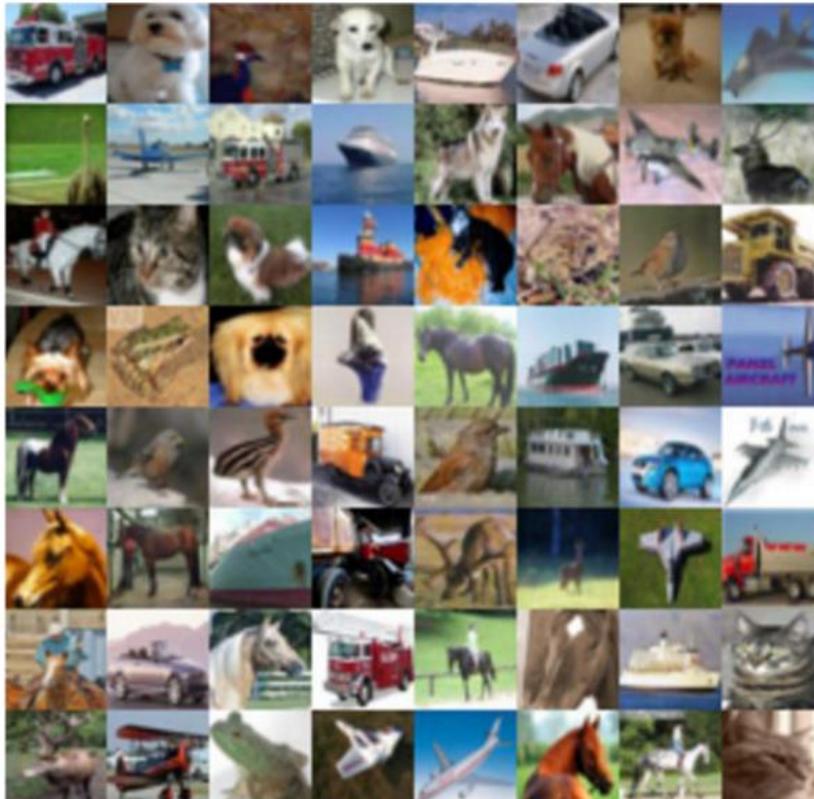
277



Cifar-10

278

- Which one is machine-generated?



<https://openai.com/blog/generative-models/>

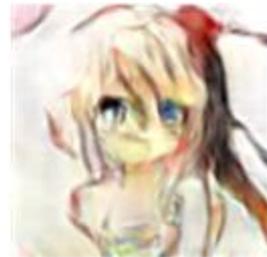
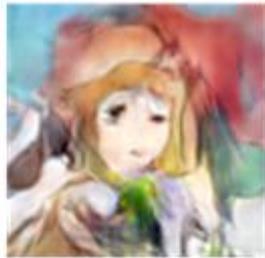
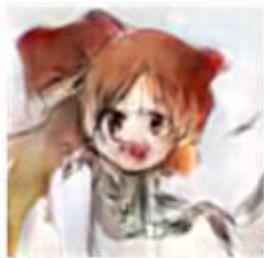
Generated Bedrooms

279



Comics Drawing

280



Comics Drawing

281



元画像



-赤髪+金髪



-赤目+青目



+制服+セーラー

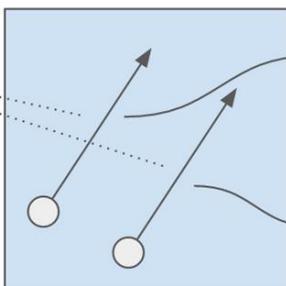


+笑顔+口開き



+青背景

長髪化ベクトル



一番左のキャラクターが元画像で、
右に行くほど長髪化ベクトルを強く足している

Pokémon Creation

282

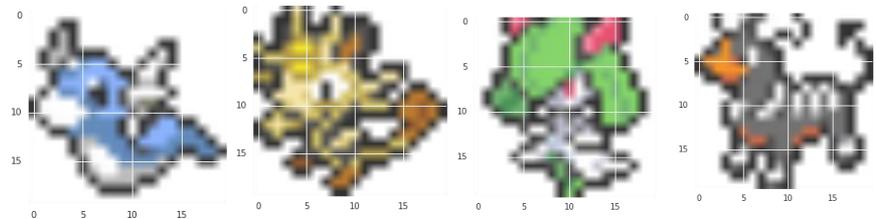
- Small images of 792 Pokémon's
 - ▣ Can machine learn to create new Pokémons?

Don't catch them! Create them!

- Source of image:
[http://bulbapedia.bulbagarden.net/wiki/List_of_Pok%C3%A9mon_by_base_stats_\(Generation_VI\)](http://bulbapedia.bulbagarden.net/wiki/List_of_Pok%C3%A9mon_by_base_stats_(Generation_VI))

Original image is 40 x 40

Making them into 20 x 20



Pokémon Creation

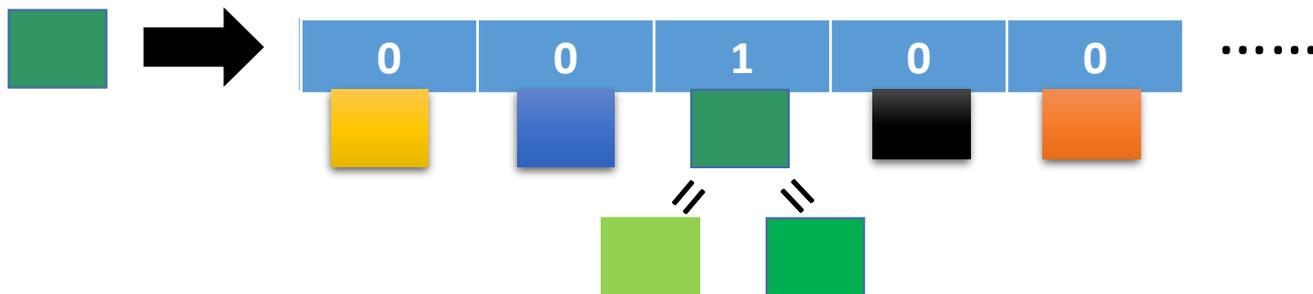
283

- Each pixel is represented by 3 numbers (corresponding to RGB)



R=50, G=150, B=100

- Each pixel is represented by a 1-of-N encoding feature



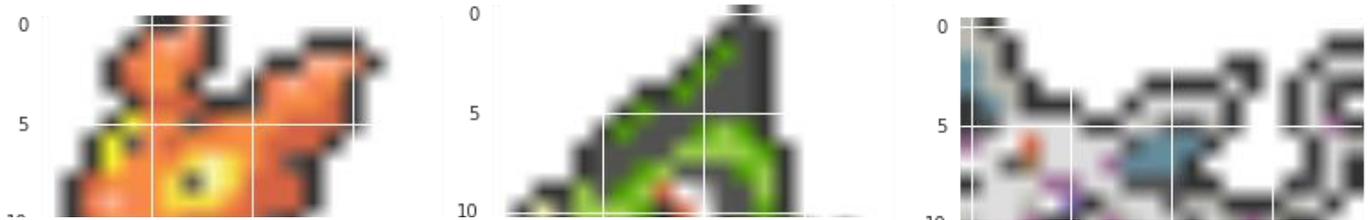
Clustering the similar color  167 colors in total

Pokémon Creation

284

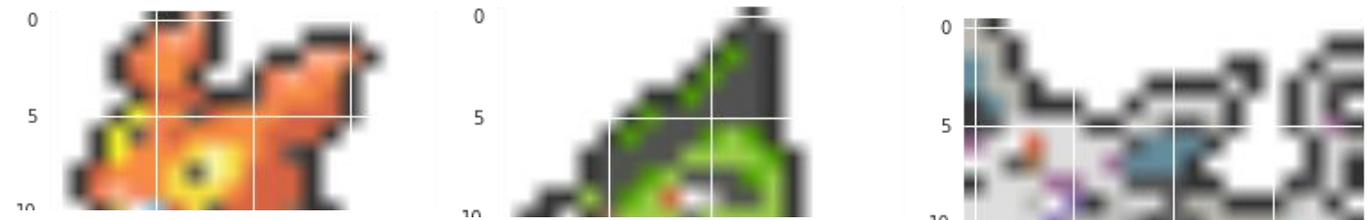
Real
Pokémon

Never seen
by machine!



It is difficult to evaluate generation.

Cover 50%



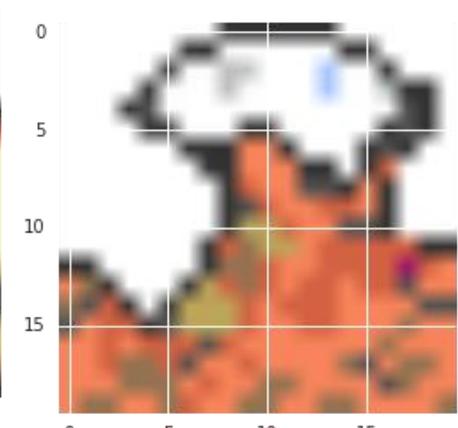
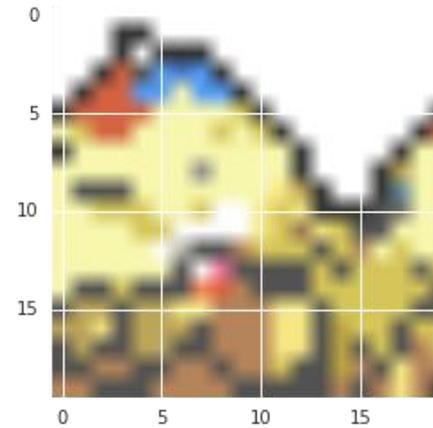
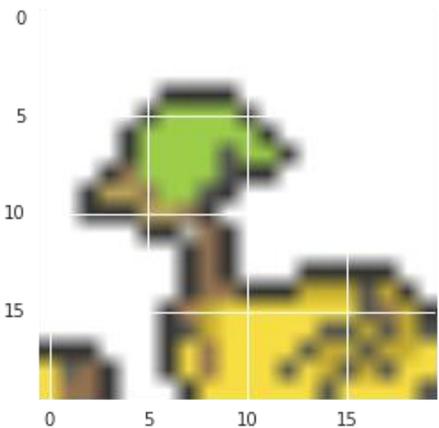
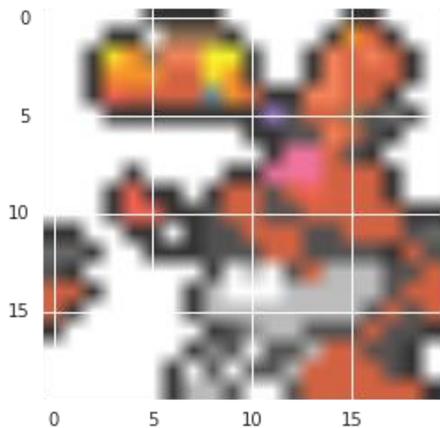
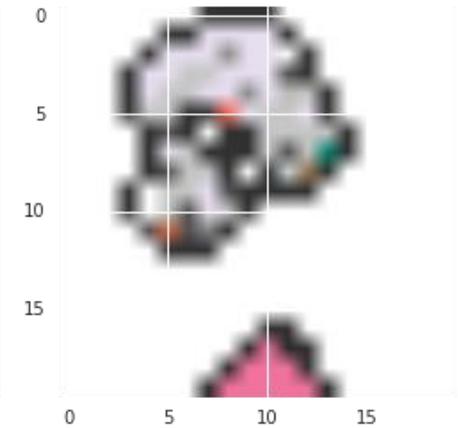
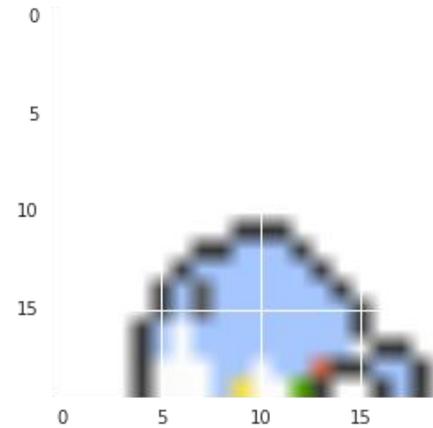
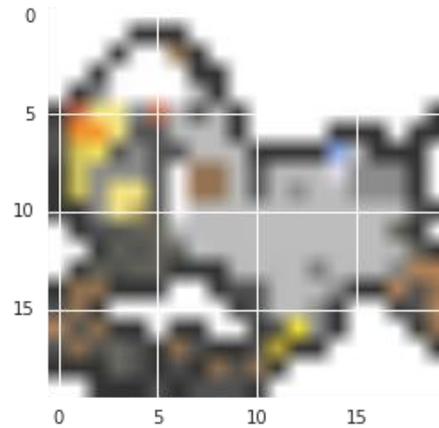
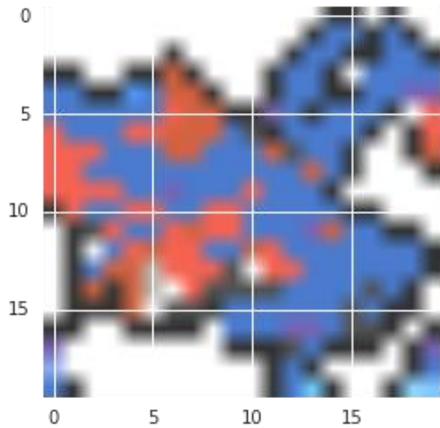
Cover 75%



Pokémon Creation

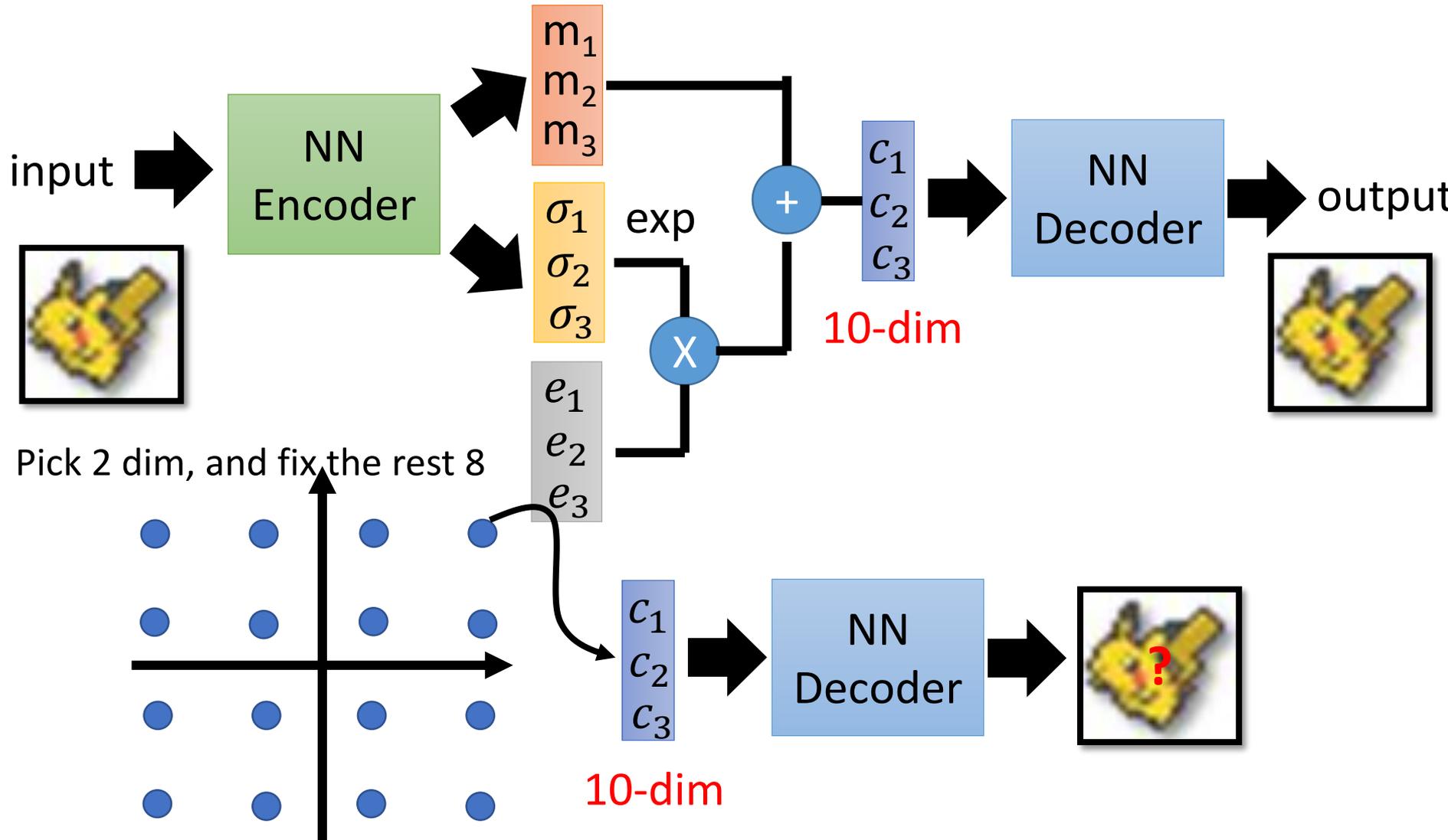
Drawing from scratch
Need some randomness

285

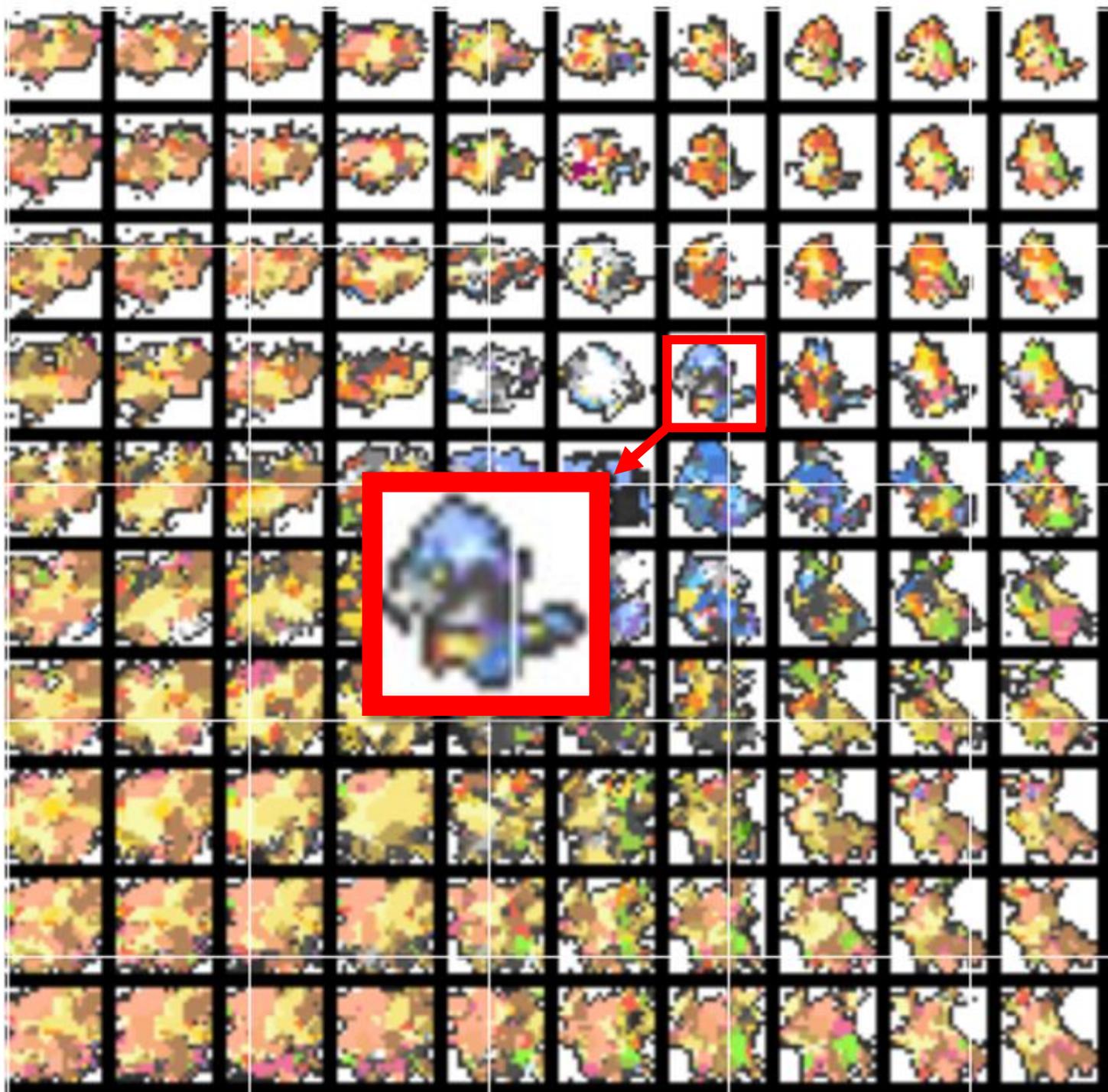


Pokémon Creation

286







Pokémon Creation - Data

289

- Original image (40 x 40):
http://speech.ee.ntu.edu.tw/~tlkagk/courses/ML_2016/Pokemon_creation/image.rar
- Pixels (20 x 20):
http://speech.ee.ntu.edu.tw/~tlkagk/courses/ML_2016/Pokemon_creation/pixel_color.txt
- Each line corresponds to an image, and each number corresponds to a pixel
 - http://speech.ee.ntu.edu.tw/~tlkagk/courses/ML_2016/Pokemon_creation/colormap.txt



```
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 19 41 34 0 0 19 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 1 2 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 1 44 74 44 51 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 1 21 80 80 81 0 0 0 0 0 0 0 0
0 0 0 0 0 1 2 3 18 35 22 0 5 2 0 0 0 0 0 0
93 94 93 93 85 95 38 96 97 98 99 99 67 99 9
0 0 0 0 0 0 1 106 106 106 106 106 61 107 0
```

.....

0	→	255 255 255
1	→	53 53 53
2	→	49 49 49
		186 186 186
		51 51 51
		54 54 54
		187 187 187
		83 83 83
		50 51 52
		251 251 251
		52 52 52

⋮

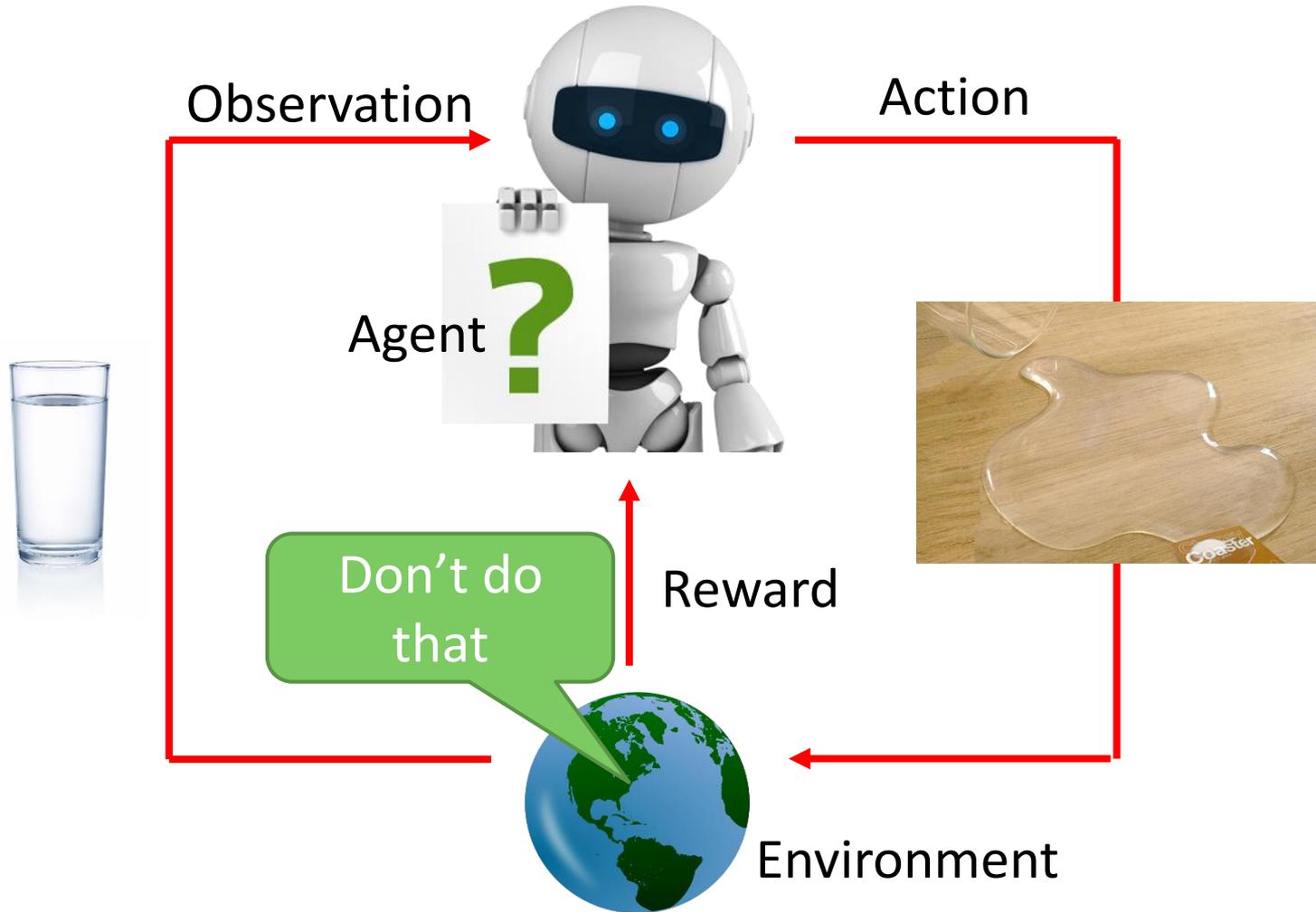
Outline

290

- Semi-Supervised Learning
- Transfer Learning
- Unsupervised Learning
 - ▣ 化繁為簡 Representation Learning
 - ▣ 無中生有 Generative Model
- Reinforcement Learning

Reinforcement Learning

291



Reinforcement Learning

292



Agent learns to take actions to maximize expected reward.

Supervised v.s. Reinforcement

293

□ Supervised

Learning from teacher



"Hello"

Say "Hi"



"Bye bye"

Say "Good bye"

□ Reinforcement



.....



.....

.....

Hello 😊

.....

Learning from critics

Agent

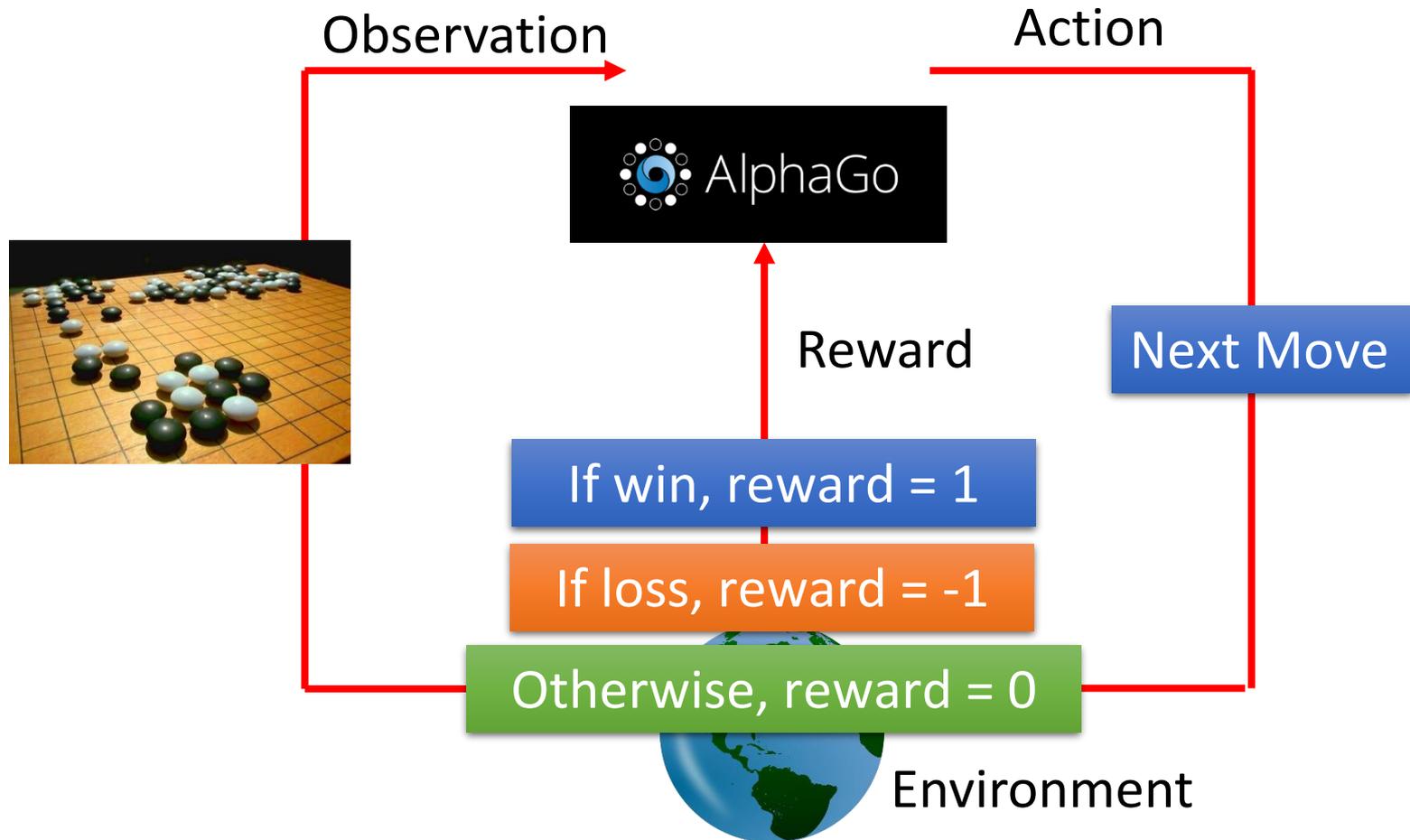
Agent



Bad

Scenario of Reinforcement Learning

294



Agent learns to take actions to maximize expected reward.

Supervised v.s. Reinforcement

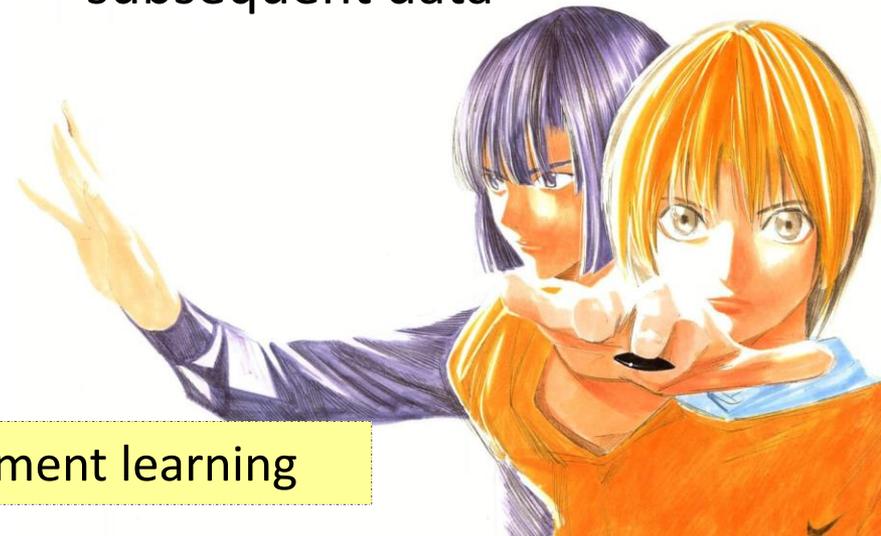
295

□ Supervised Learning

- Training based on supervisor/label/annotation
- Feedback is instantaneous
- Time does not matter

□ Reinforcement Learning

- Training only based on reward signal
- Feedback is delayed
- Time matters
- Agent actions affect subsequent data



AlphaGo uses supervised learning + reinforcement learning

Reinforcement Learning

296

- RL is a general purpose framework for **decision making**
 - ▣ RL is for an *agent* with the capacity to *act*
 - ▣ Each *action* influences the agent's future *state*
 - ▣ Success is measured by a scalar *reward* signal
 - ▣ Goal: *select actions to maximize future reward*



RL Difficulty

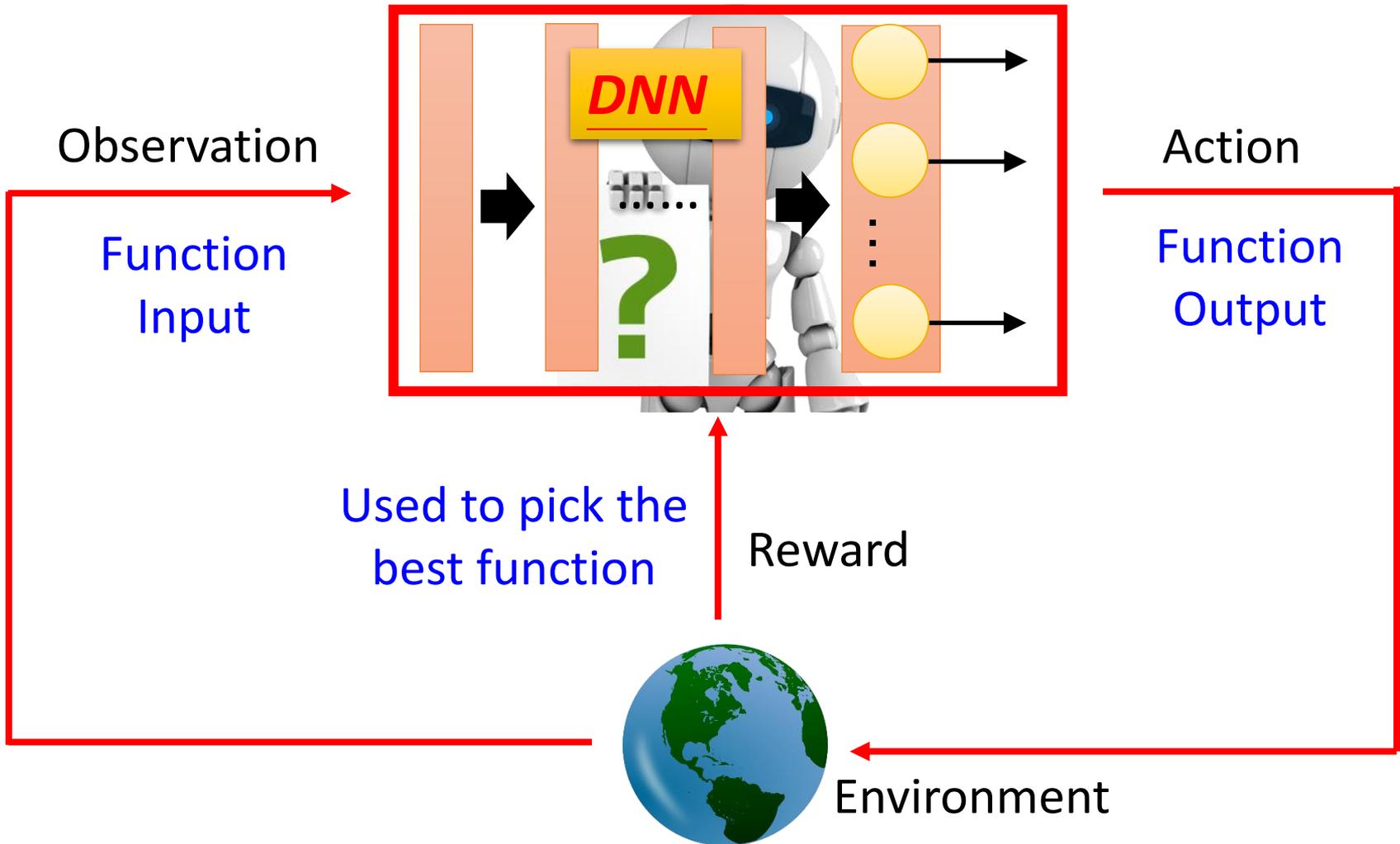
297

- Goal: select actions to maximize total future reward
 - ▣ Actions may have long-term consequences
 - ▣ Reward may be delayed
 - ▣ It may be better to sacrifice immediate reward to gain more long-term reward



Deep Reinforcement Learning

298



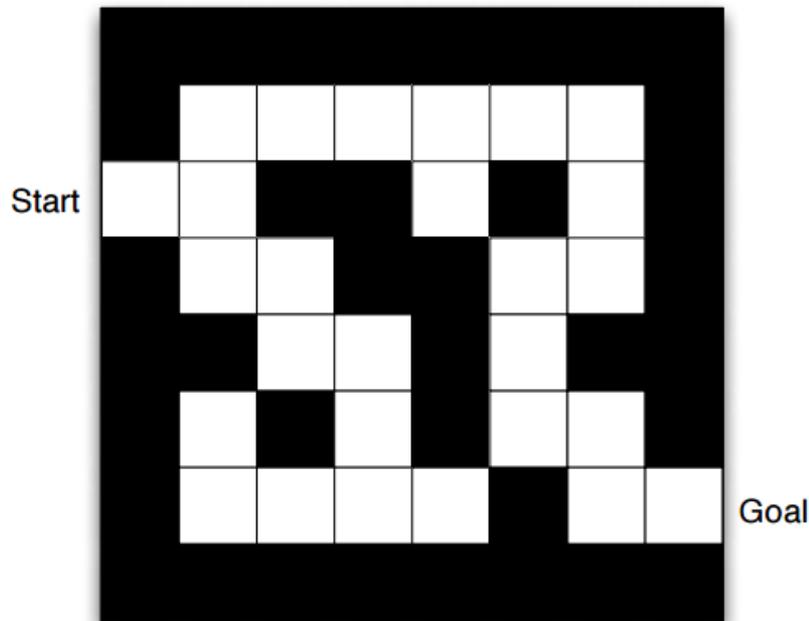
Major Components in an RL Agent

299

- An RL agent may include one or more of these components
 - ▣ **Policy**: agent's behavior function
 - ▣ **Value function**: how good is each state and/or action
 - ▣ **Model**: agent's representation of the environment

Maze Example

300

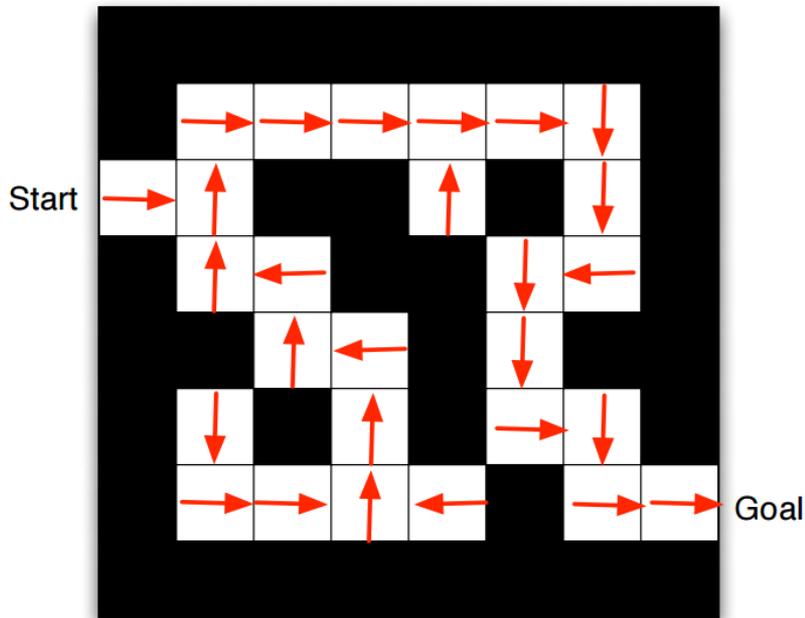


- Rewards: -1 per time-step
- Actions: N, E, S, W
- States: agent's location

Maze Example: Policy

301

- Rewards: -1 per time-step
- Actions: N, E, S, W
- States: agent's location

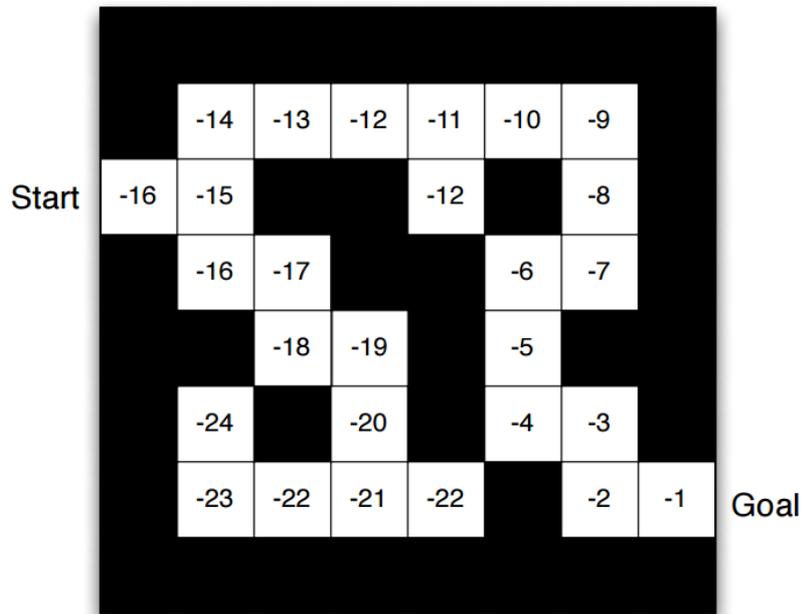


Arrows represent policy $\pi(s)$ for each state s

Maze Example: Value Function

302

- Rewards: -1 per time-step
- Actions: N, E, S, W
- States: agent's location



Numbers represent value $Q_{\pi}(s)$ of each state s

E2E RL Task-Completion Bot

303

□ Simulated User

- Generate interactions based on a predefined fake goal
- Automatically learn strategy by training on the simulated data

User goal: *Two tickets for “the witch” tomorrow 9:30 PM at regal meridian 16, Seattle.*

RULE BASED SYSTEM

User: Can I get 2 tickets for the witch?

Agent: Which theater do you want?

User: Regal meridian 16

Agent: What time would you like to see it?

User: 9:30 pm

Agent: What date would you like to see it?

User: Tomorrow!

Agent: How many people are you?

User: Just two adults.

Agent: What is your location?

User: Seattle!

Agent: What kind of video format do you like?

User: Oh, any kind is okay.

Agent: Do you have any price constraint?

User: No, I don't care.

Agent: Great - I was able to purchase 2 tickets for you to see the witch tomorrow at regal meridian 16 theater in Seattle at 9:30 pm.

User: Thank you.

Agent: Thank you!

REINFORCEMENT LEARNING SYSTEM

User: Can I get 2 tickets for the witch?

Agent: What date would you like to see it?

User: Tomorrow

Agent: Which theater do you prefer?

User: Regal meridian 16

Agent: What is your location?

User: Seattle

Agent: What time would you like to see it?

User: 9:30 pm.

Agent: Great - I was able to purchase 2 tickets for you to see the witch tomorrow at regal meridian 16 theater in Seattle at 9:30 pm.

User: Thanks.

Agent: Thanks!

More applications

304

- AlphaGo, Playing Video Games, Dialogue
- Flying Helicopter
 - ▣ <https://www.youtube.com/watch?v=0JL04JJjocc>
- Driving
 - ▣ <https://www.youtube.com/watch?v=0xo1Ldx3L5Q>
- Google Cuts Its Giant Electricity Bill With DeepMind-Powered AI
 - ▣ <http://www.bloomberg.com/news/articles/2016-07-19/google-cuts-its-giant-electricity-bill-with-deepmind-powered-ai>

Concluding Remarks

305

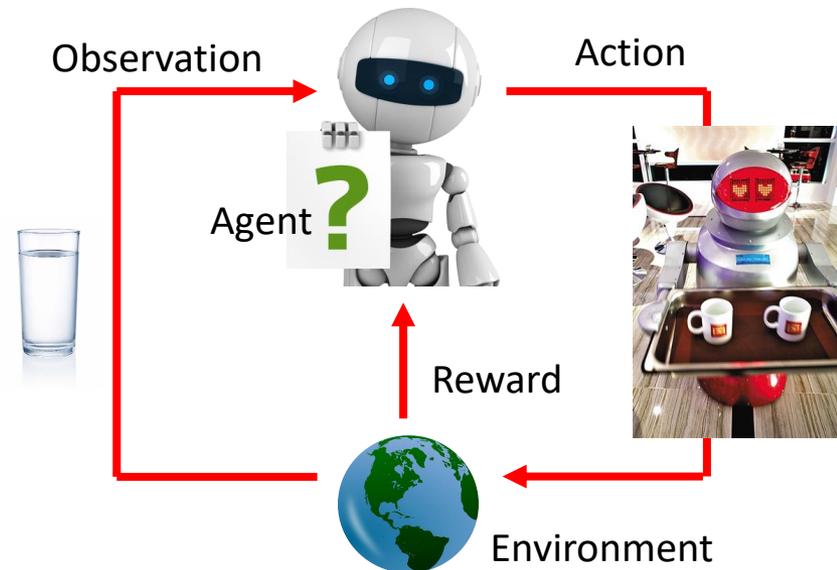
- Semi-Supervised Learning
- Transfer Learning
- Unsupervised Learning
 - ▣ 化繁為簡 Representation Learning
 - ▣ 無中生有 Generative Model
- Reinforcement Learning

Semi-Supervised Learning

Transfer Learning

Unsupervised Learning

Reinforcement Learning



如何成為武林高手

306

- 內外兼修
 - ▣ 內功充沛，恃強克弱
 - ▣ 招數精妙，以快打慢
- Machine Learning & Deep Learning 也需要內外兼修
 - ▣ 內力：運算資源
 - ▣ 招數：各種技巧
- 內力充沛,平常的招式也有可能發會巨大的威力