# Deep Learning Tutorial
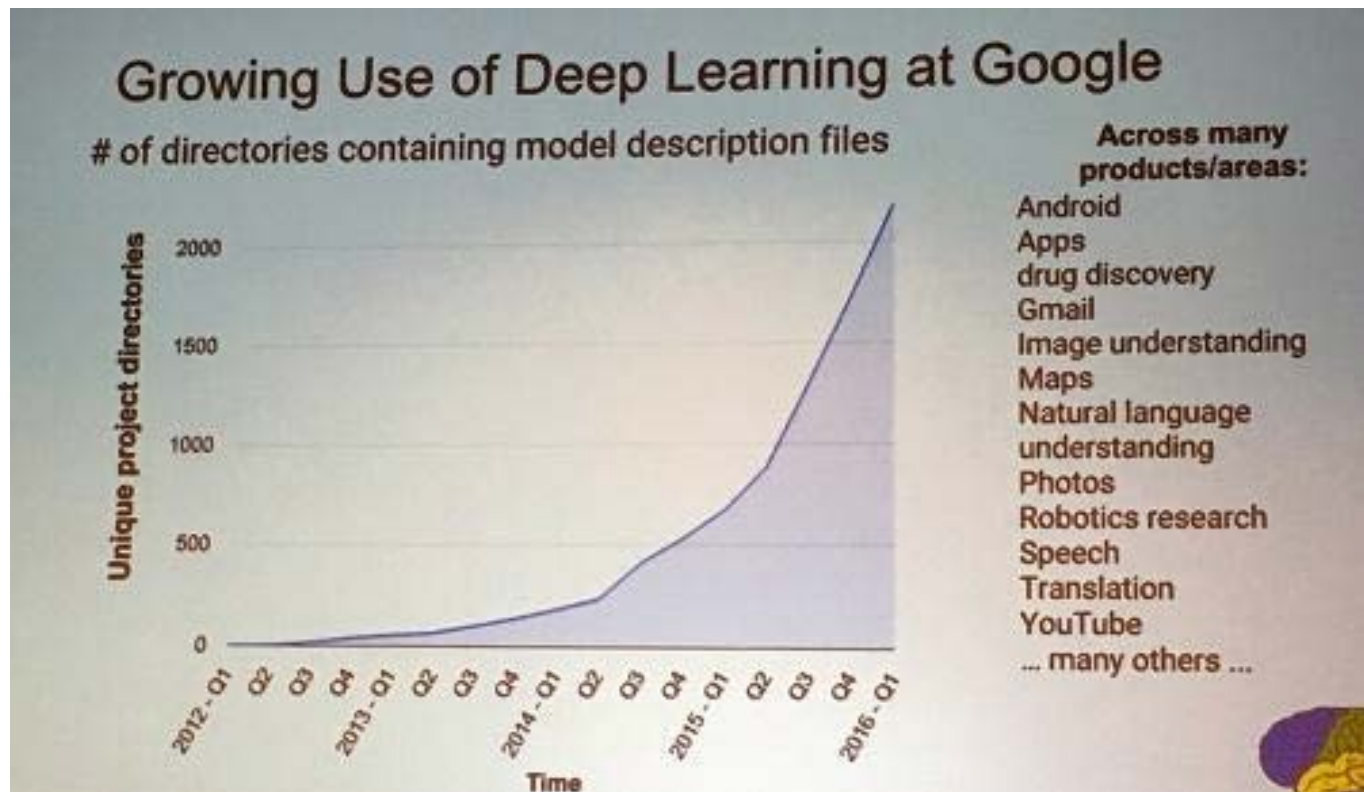
李宏毅

Hung-yi Lee

# Deep learning attracts lots of attention.

- I believe you have seen lots of exciting results before.

Growing Use of Deep Learning at Google

# of directories containing model description files

Across many products/areas:
Android
Apps
drug discovery
Gmail
Image understanding
Maps
Natural language understanding
Photos
Robotics research
Speech
Translation
YouTube
... many others ...

Deep learning trends at Google. Source: SIGMOD/Jeff Dean

This talk focuses on the basic techniques.

# Outline

Lecture I: Introduction of Deep Learning

Lecture II: Tips for Training Deep Neural Network

Lecture III: Variants of Neural Network

Lecture IV: Next Wave

# Lecture I:
# Introduction of Deep Learning

# Three Steps for Deep Learning

Step 1:
Network
Structure
→
Step 2:
Learning
Target
→
Step 3:
Learn!

天生的腦

# Three Steps for Deep Learning



Step 1: Network Structure → Step 2: Learning Target → Step 3: Learn!

based on training data

# Three Steps for Deep Learning

- Speech Recognition

$$f*(\quad\text{〜〜〜}\quad) = \text{"你好"}$$

- Handwritten Recognition

$$f*(\quad\text{2}\quad) = \text{"2"}$$

- Playing Go

$$f*(\quad\text{〜〜〜}\quad) = \text{"5-5"}$$
(step)

- Dialogue System

$$f*(\quad\text{"Hi"}\quad) = \text{"Hello"}$$
(what the user said)    (system response)

Step 3: Learn!

⬇

Pick the best function f*

# Three Steps for Deep Learning

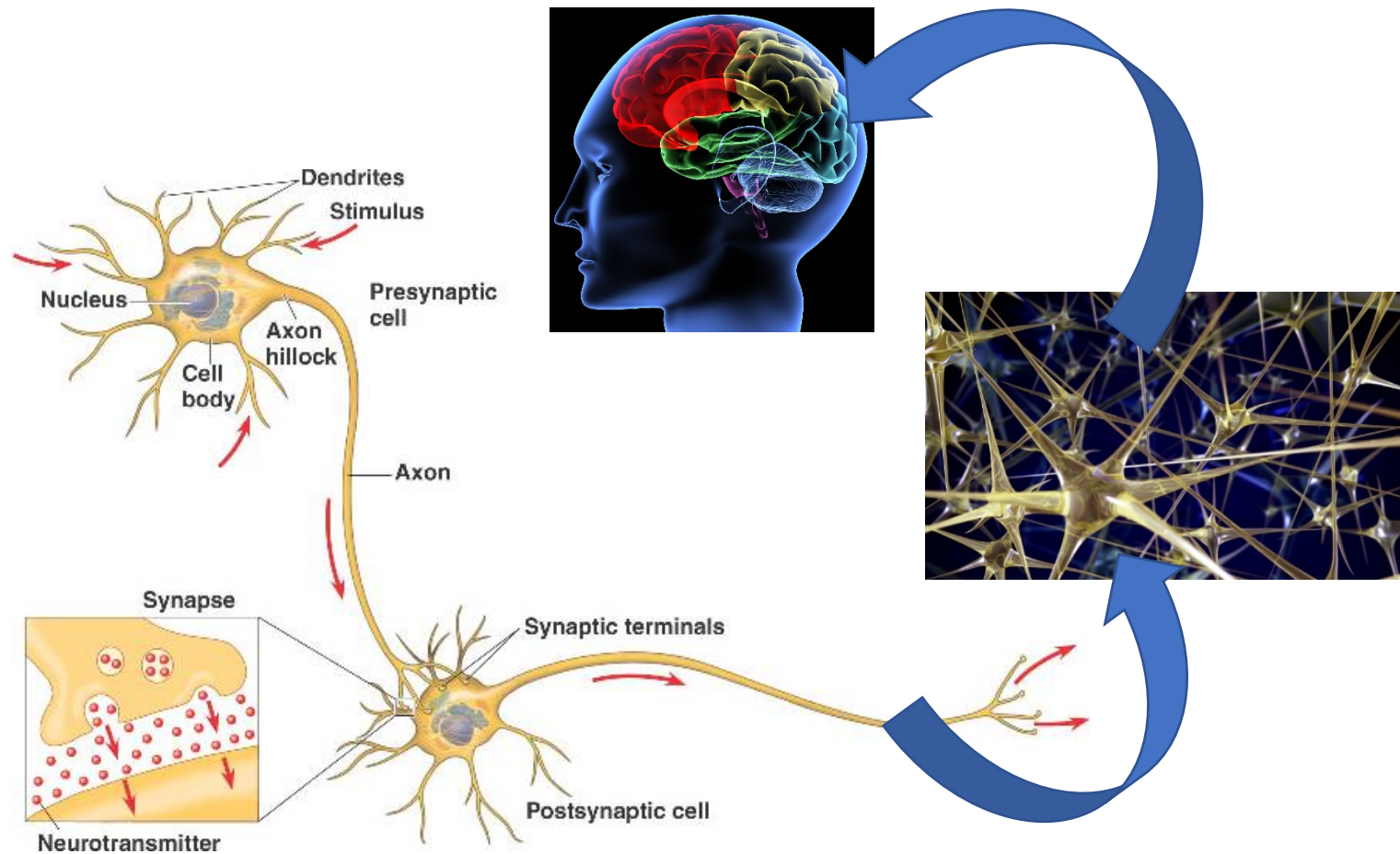| Step 1: Network Structure | → | Step 2: Learning Target | → | Step 3: Learn! |
|---|---|---|---|---|

Deep Learning is so simple ……



CDC.TENCENT.COM

# Three Steps for Deep Learning

# Human Brains

# Neural Network

## *Neuron*

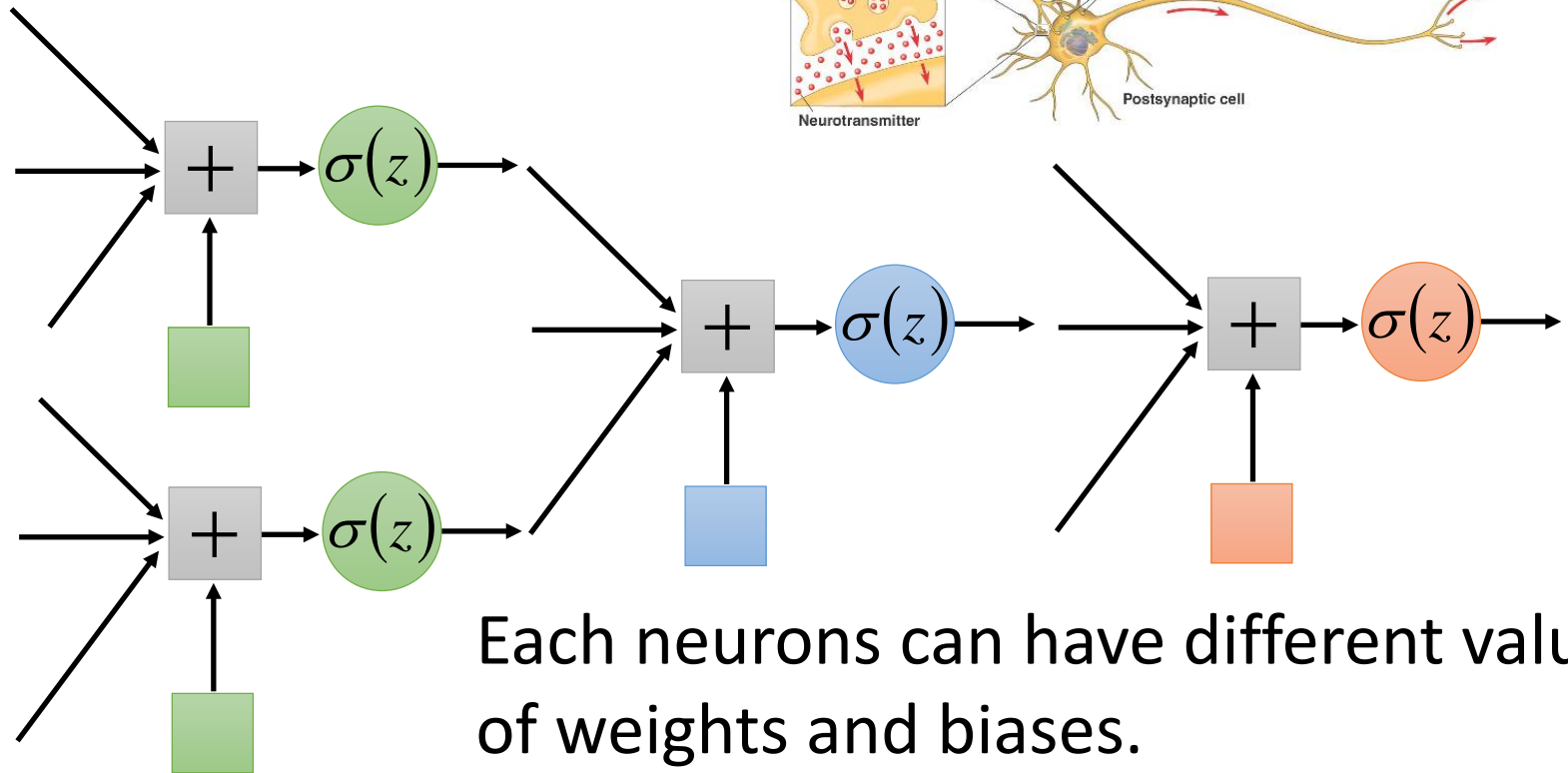$$z = a_1 w_1 + \cdots + a_k w_k + \cdots + a_K w_K + b$$



A simple function

$$\sigma(z)$$

Activation function

weights

$b$ bias

$z$  $a$

# Neural Network

**_Neuron_**

Sigmoid Function $\quad \sigma(z)$

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



2

1

-1

-2

1

-1

weights

+  4  $\sigma(z)$  →  0.98

1  bias

Activation function

# Neural Network

Different connections leads to different network structured



$+$   $\sigma(z)$

$+$   $\sigma(z)$

$+$   $\sigma(z)$

$+$   $\sigma(z)$

$+$   $\sigma(z)$

Each neurons can have different values of weights and biases.

Weights and biases are network parameters $\theta$
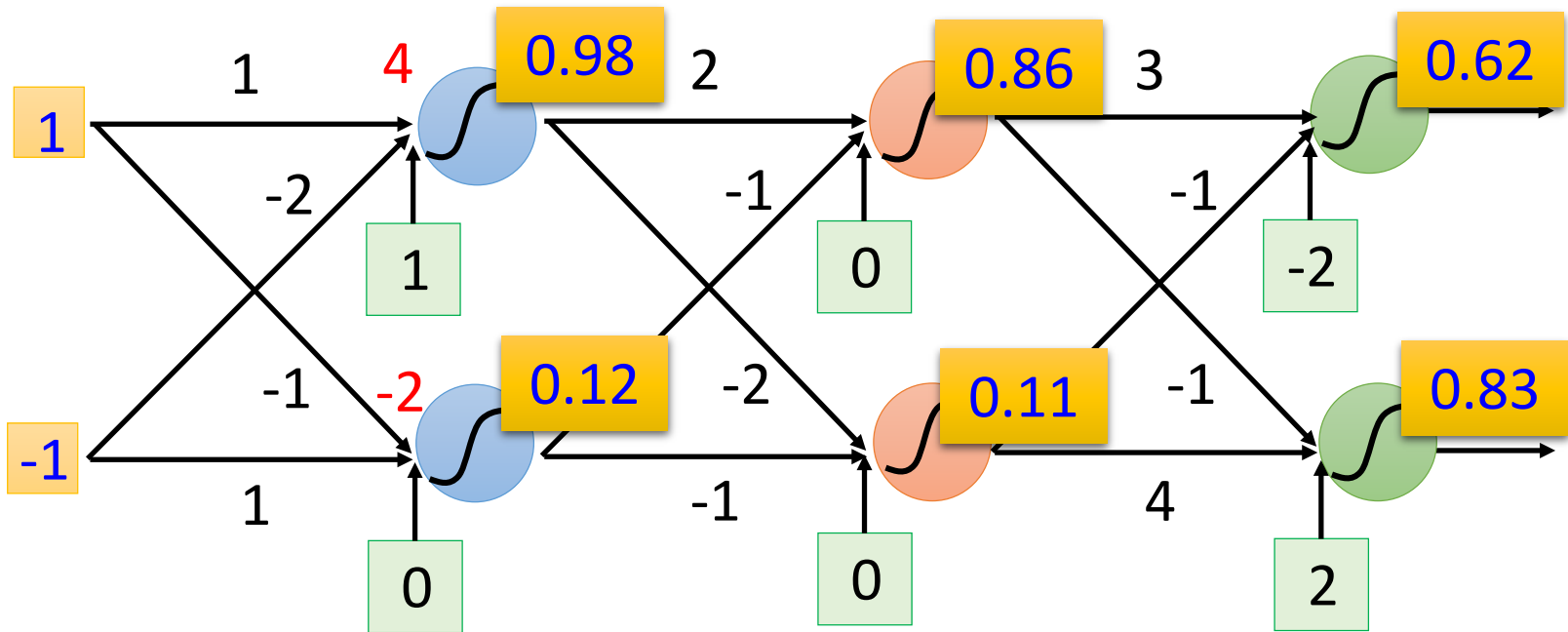
# Fully Connect Feedforward Network
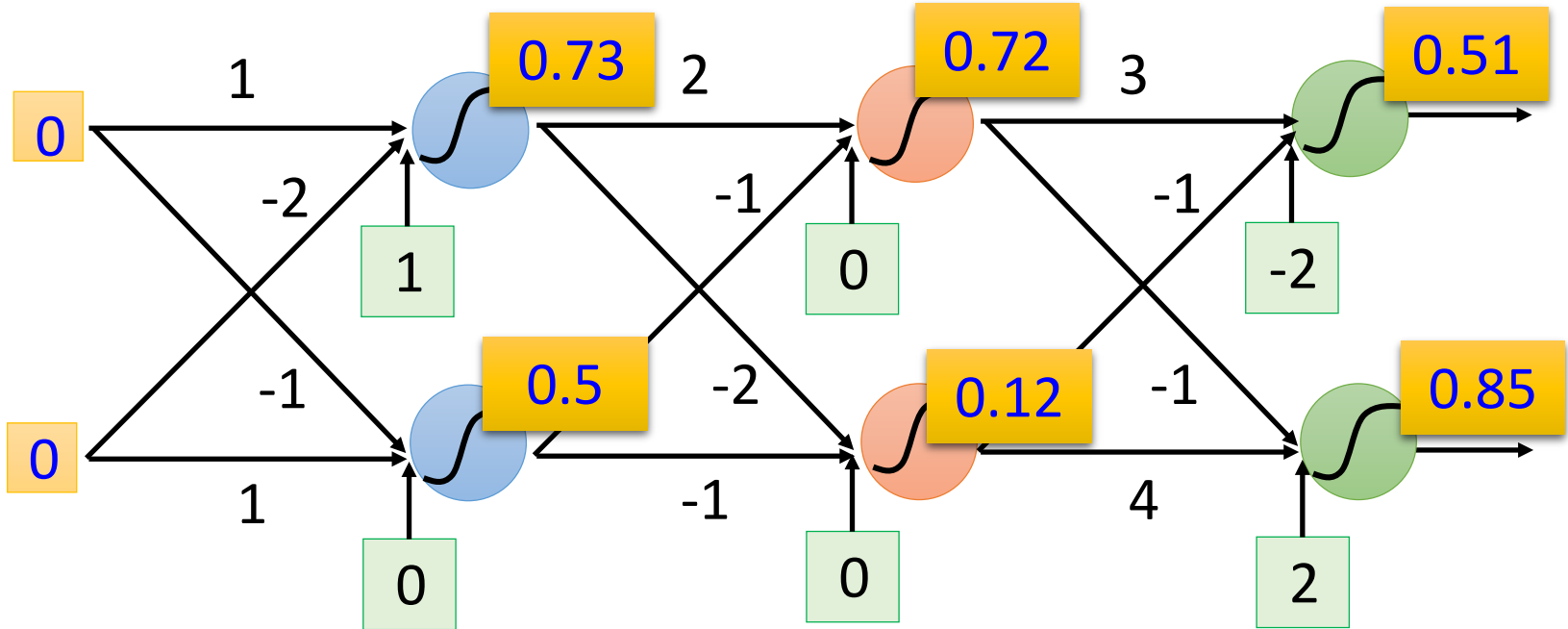


Sigmoid Function $\sigma(z)$

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

# Fully Connect Feedforward Network

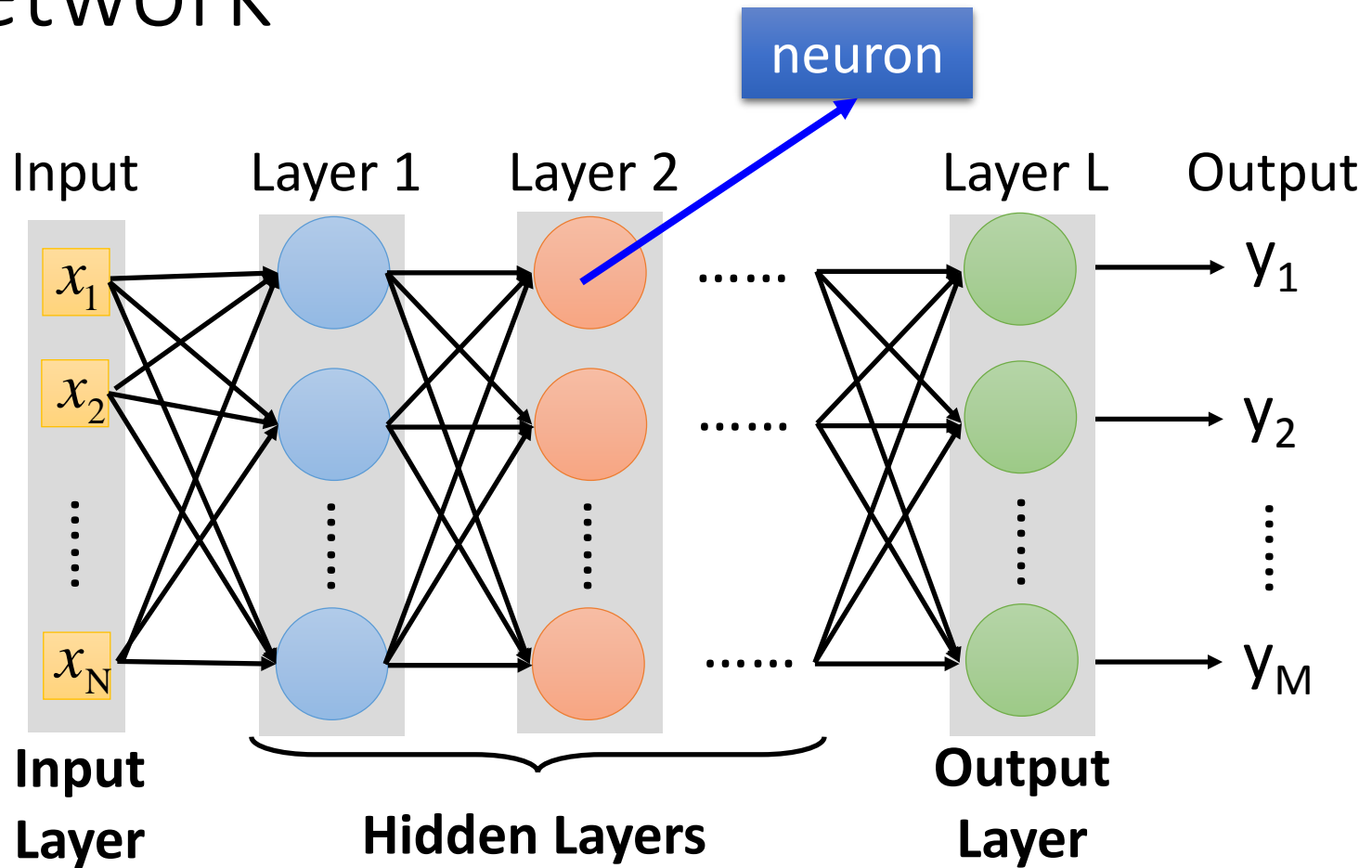# Fully Connect Feedforward Network



Network 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 $\theta$, define a function

Given network structure, define **_a function set_**
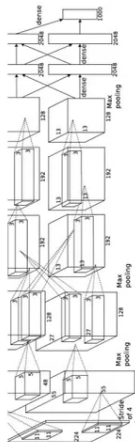
# Fully Connect Feedforward Network

neuron

Input    Layer 1    Layer 2    Layer L    Output

$x_1$    $x_2$    $x_N$

$y_1$
$y_2$
$y_M$

**Input Layer**    **Hidden Layers**    **Output Layer**

Deep means many hidden layers

# Ultra Deep Network

http://cs231n.stanford.edu/slides/winter1516_lecture8.pdf

22 layers

19 layers

8 layers

6.7%

7.3%

16.4%

AlexNet (2012)　　VGG (2014)　　GoogleNet (2014)

# Ultra Deep Network

**152 layers**

**101 layers**

This ultra deep network have special structure.

(Lecture IV)

3.57%

16.4%

7.3%

6.7%

AlexNet
(2012)

VGG
(2014)

GoogleNet
(2014)

Residual Net
(2015)

Taipei
101

# Fully Connect Feedforward Network - Matrix Operation



$$\sigma\left( W^1 \quad x \quad + \quad b^1 \right)$$

$$\sigma\left( W^2 \quad a^1 \quad + \quad b^2 \right)$$

$$\sigma\left( W^L \quad a^{L-1} \quad + \quad b^L \right)$$

# Fully Connect Feedforward Network - Matrix Operation



$$\boxed{y} = f(\boxed{x})$$

Using parallel computing techniques (e.g. GPU) to speed up matrix operation

$$= \sigma(\boxed{W^L} \cdots \sigma(\boxed{W^2} \sigma(\boxed{W^1} \boxed{x} + \boxed{b^1}) + \boxed{b^2}) \cdots + \boxed{b^L})$$

# Output Layer (Option)

- Softmax layer as the output layer

**_Ordinary Layer_**

$z_1 \longrightarrow \boxed{\sigma} \longrightarrow y_1 = \sigma(z_1)$

$z_2 \longrightarrow \boxed{\sigma} \longrightarrow y_2 = \sigma(z_2)$

$z_3 \longrightarrow \boxed{\sigma} \longrightarrow y_3 = \sigma(z_3)$

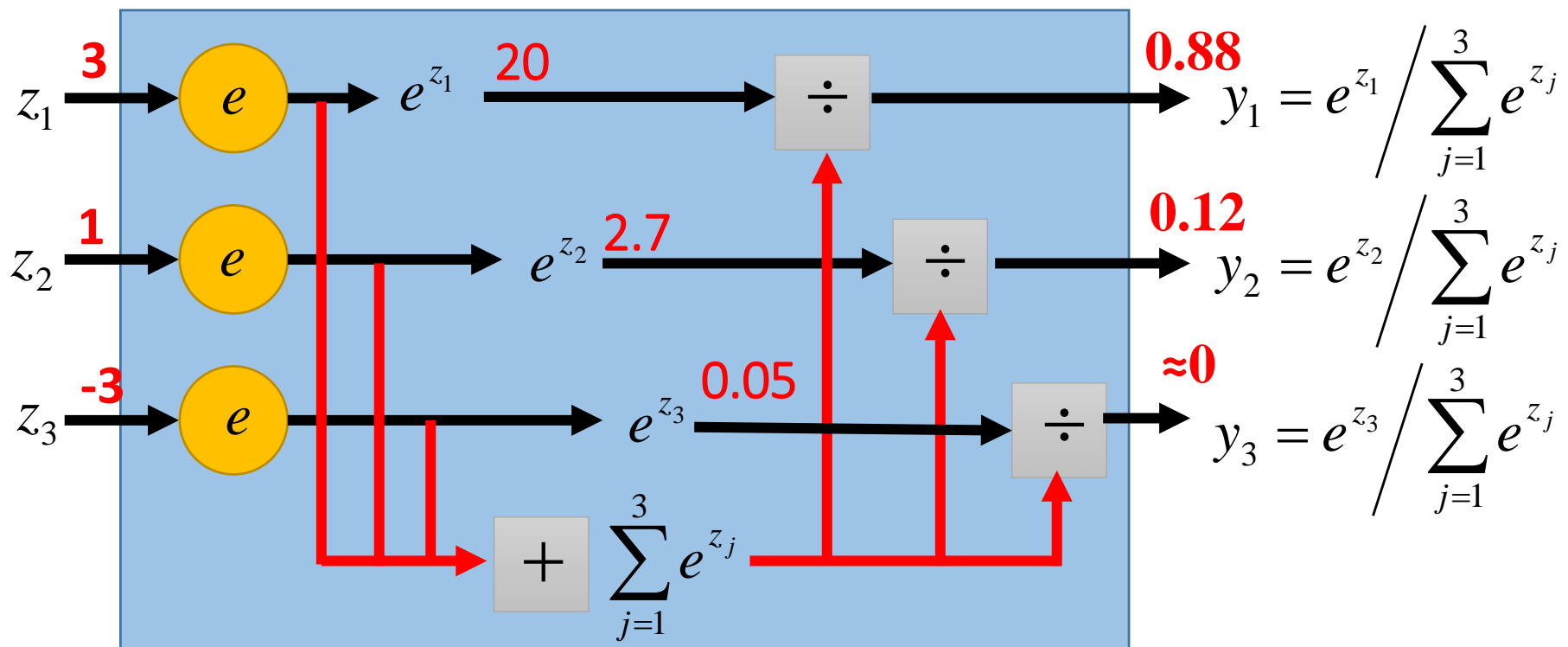In general, the output of network can be any value.

May not be easy to interpret

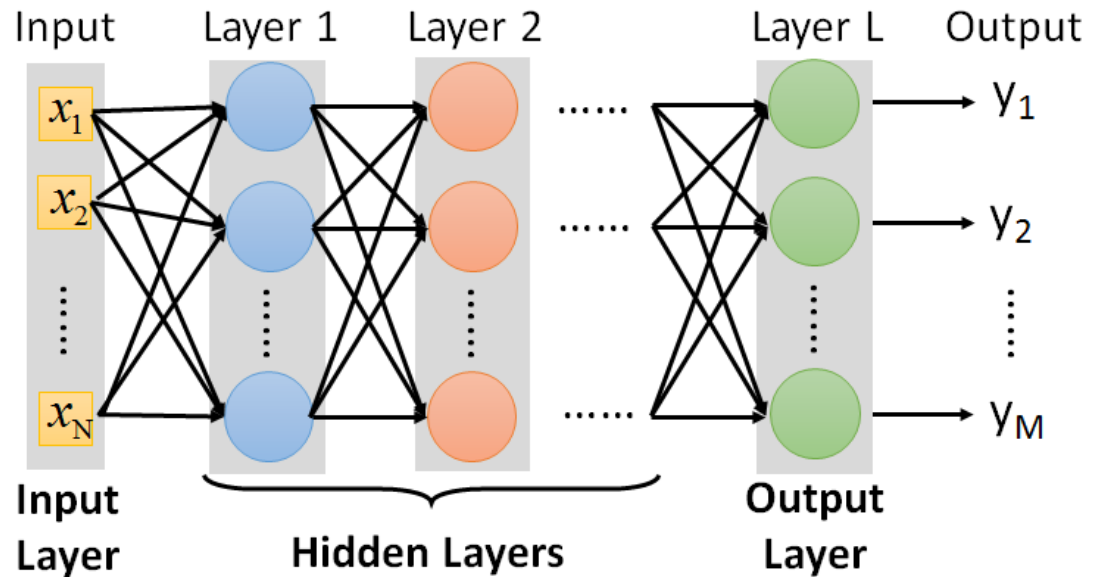# Output Layer (Option)

- Softmax layer as the output layer

**_Probability_**:
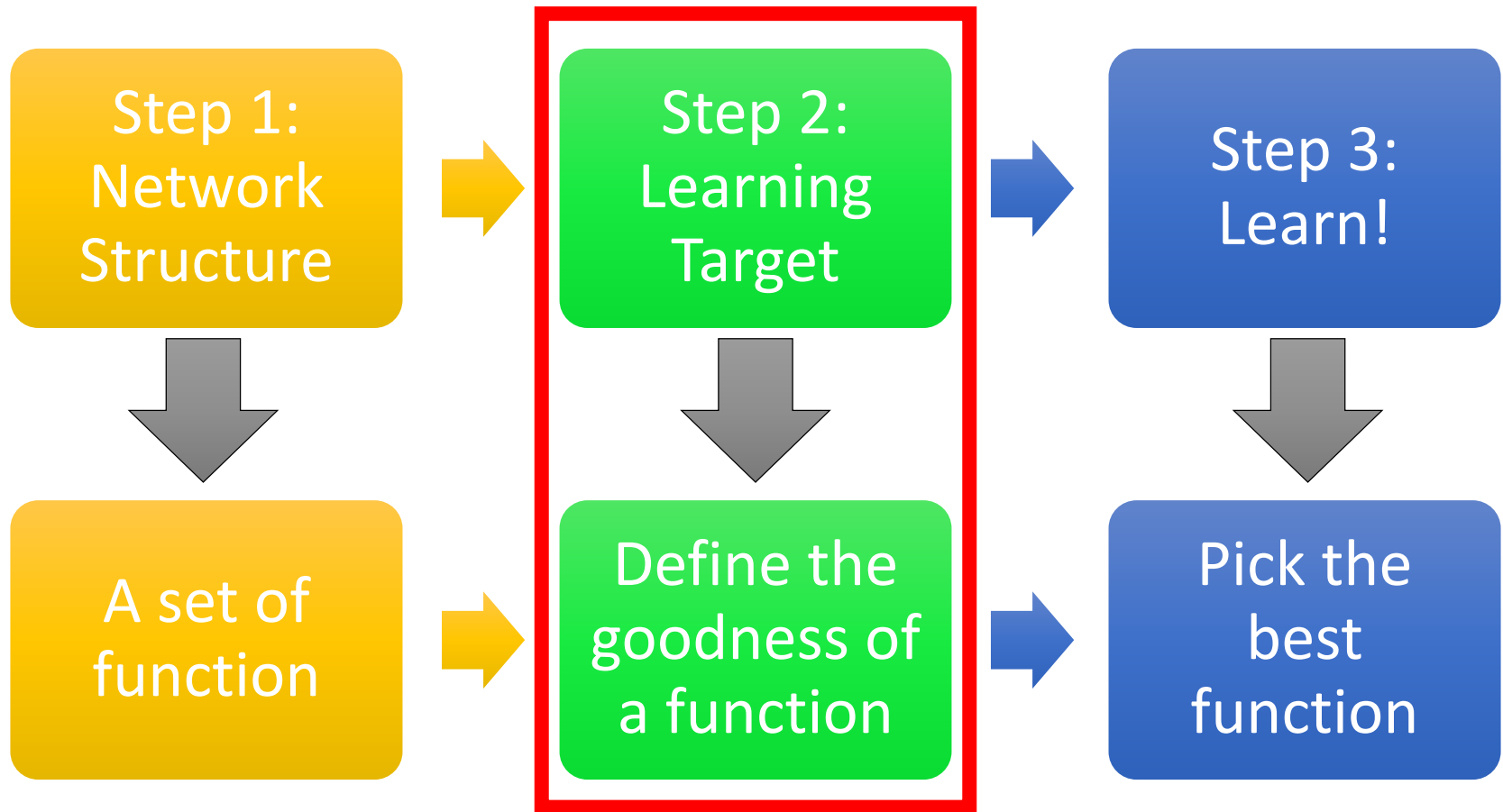- $1 > y_i > 0$
- $\sum_i y_i = 1$

**_Softmax Layer_**



**3** $z_1$   $e$   $e^{z_1}$ **20**   $\div$   **0.88** $y_1 = e^{z_1} \Big/ \sum_{j=1}^{3} e^{z_j}$

**1** $z_2$   $e$   $e^{z_2}$ **2.7**   $\div$   **0.12** $y_2 = e^{z_2} \Big/ \sum_{j=1}^{3} e^{z_j}$

**-3** $z_3$   $e$   $e^{z_3}$ **0.05**   $\div$   **≈0** $y_3 = e^{z_3} \Big/ \sum_{j=1}^{3} e^{z_j}$

$+$   $\sum_{j=1}^{3} e^{z_j}$

# FAQ



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

- Q: Can the structure be automatically determined?

# Three Steps for Deep Learning

# Example Application
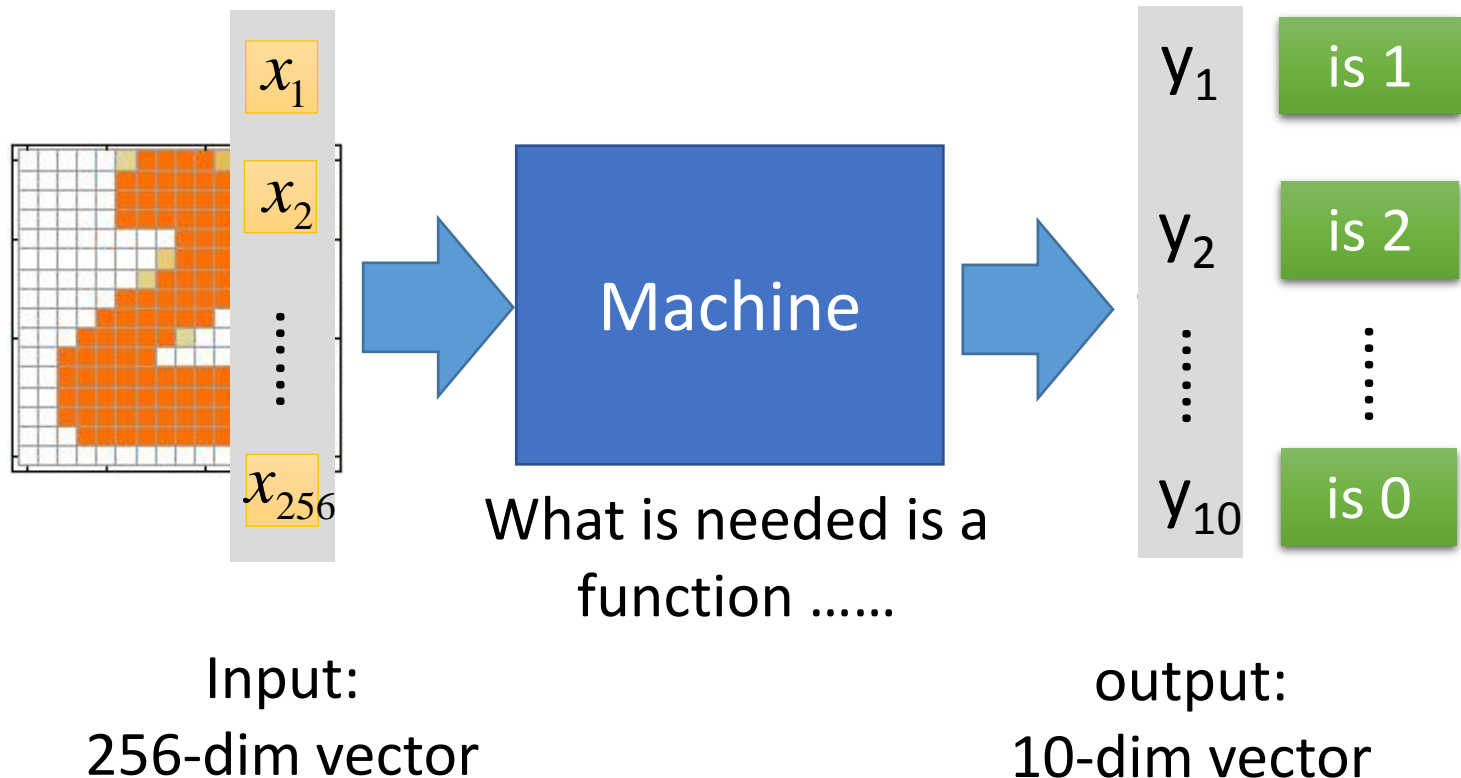


## Input



16 x 16 = 256

Ink → 1
No ink → 0

## Output

| | |
|---|---|
| 0.1 | is 1 |
| 0.7 | is 2 |
| 0.2 | is 0 |

The image is "2"

Each dimension represents the confidence of a digit.

# Example Application

- Handwriting Digit Recognition



Input:
256-dim vector

What is needed is a function ......

Machine

output:
10-dim vector

# Fully Connect Feedforward Network



neuron

Input    Layer 1    Layer 2    Layer L    Output

$x_1$

$x_2$

$x_{256}$

A function set containing the candidates for Handwriting Digit Recognition

$y_1$    is 1

$y_2$    is 2

$y_{10}$    is 0

**Input Layer**

**Hidden Layers**

**Output Layer**

# Fully Connect Feedforward Network

# Training Data

- Preparing training data: images and their labels



The learning target is defined on the training data.

# Learning Target



$x_1$
$x_2$
$\vdots$
$x_{256}$

16 x 16 = 256

Ink → 1
No ink → 0

Softmax

$y_1$ — is 1
$y_2$ — is 2
$\vdots$
$y_{10}$ — is 0

The learning target is ......

Input: [1] ⟹ $y_1$ has the maximum value

Input: [2] ⟹ $y_2$ has the maximum value

# Total Loss

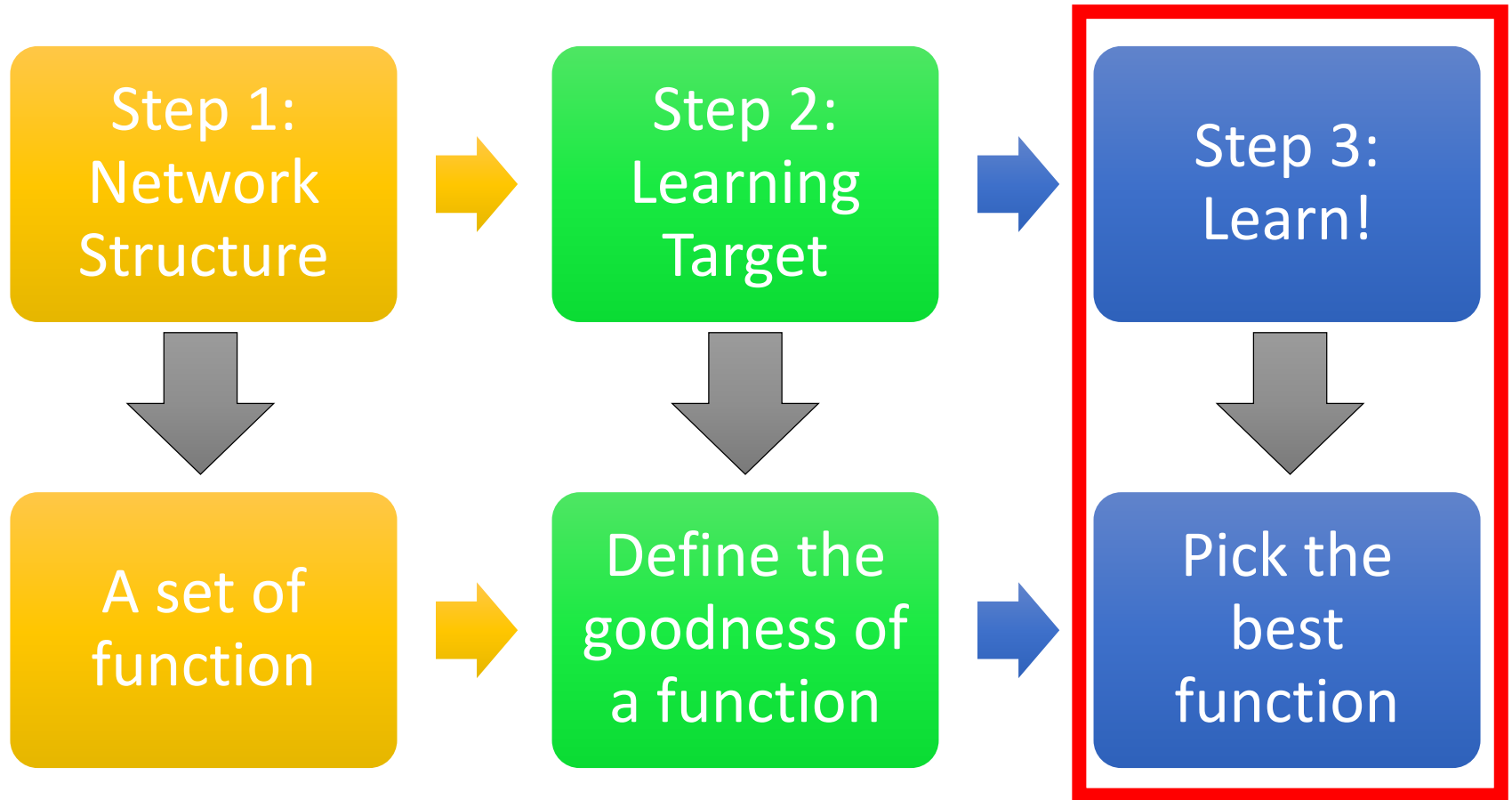## For all training data …



Total Loss:

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

As small as possible

Find *a function in function set* that minimize total loss L

Find *the network parameters* $\theta^*$ that minimize total loss L

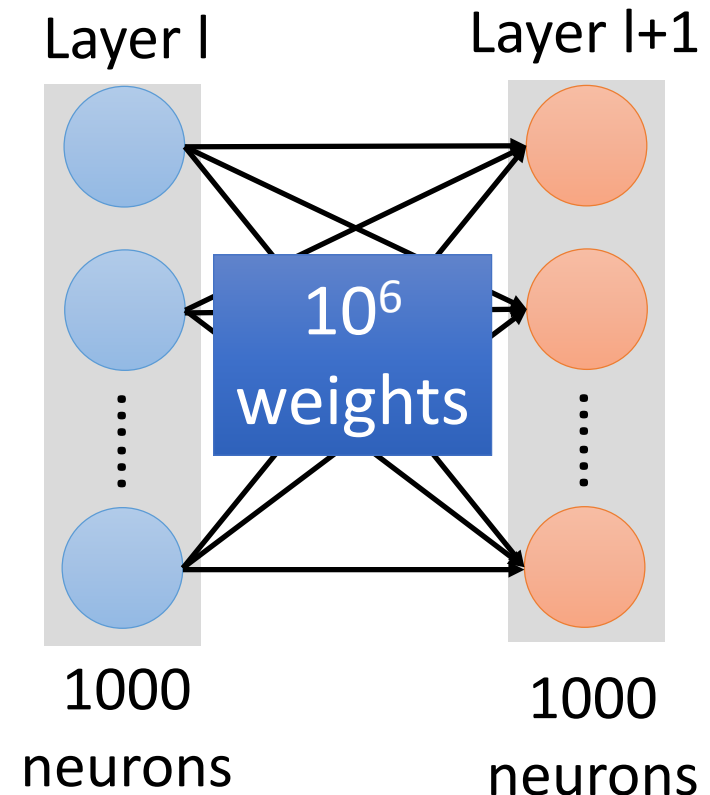# Three Steps for Deep Learning

# How to pick the best function

Find **_network parameters $\theta^*$_** that minimize total loss L

Enumerate all possible values

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

Millions of parameters

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

Layer l        Layer l+1

$10^6$ weights

1000 neurons

1000 neurons

# Gradient Descent

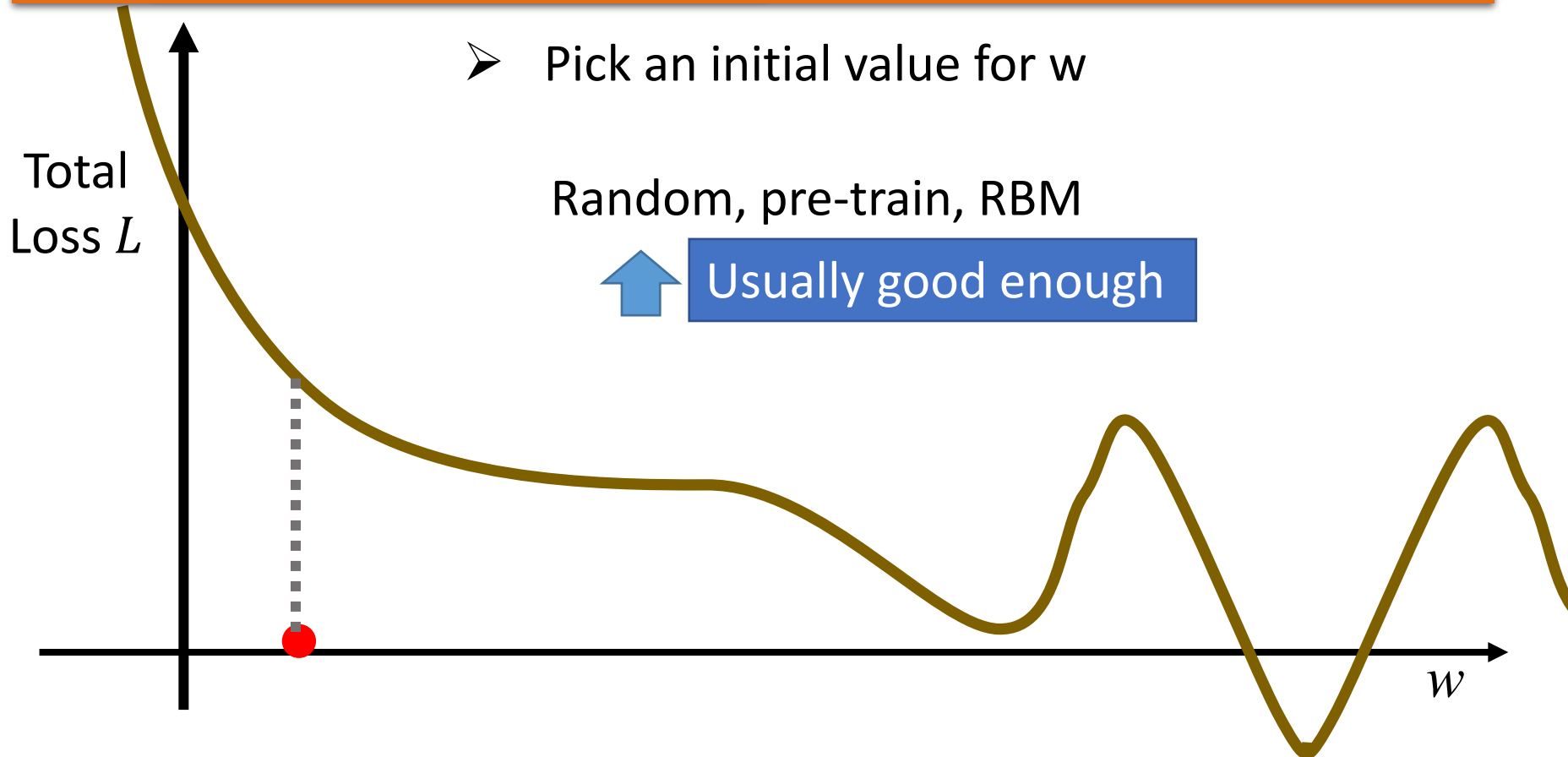Network parameters $\theta = \{w_1, w_2, \cdots, b_1, b_2, \cdots\}$

Find ***network parameters*** $\boldsymbol{\theta^*}$ that minimize total loss L

➢ Pick an initial value for w

Random, pre-train, RBM
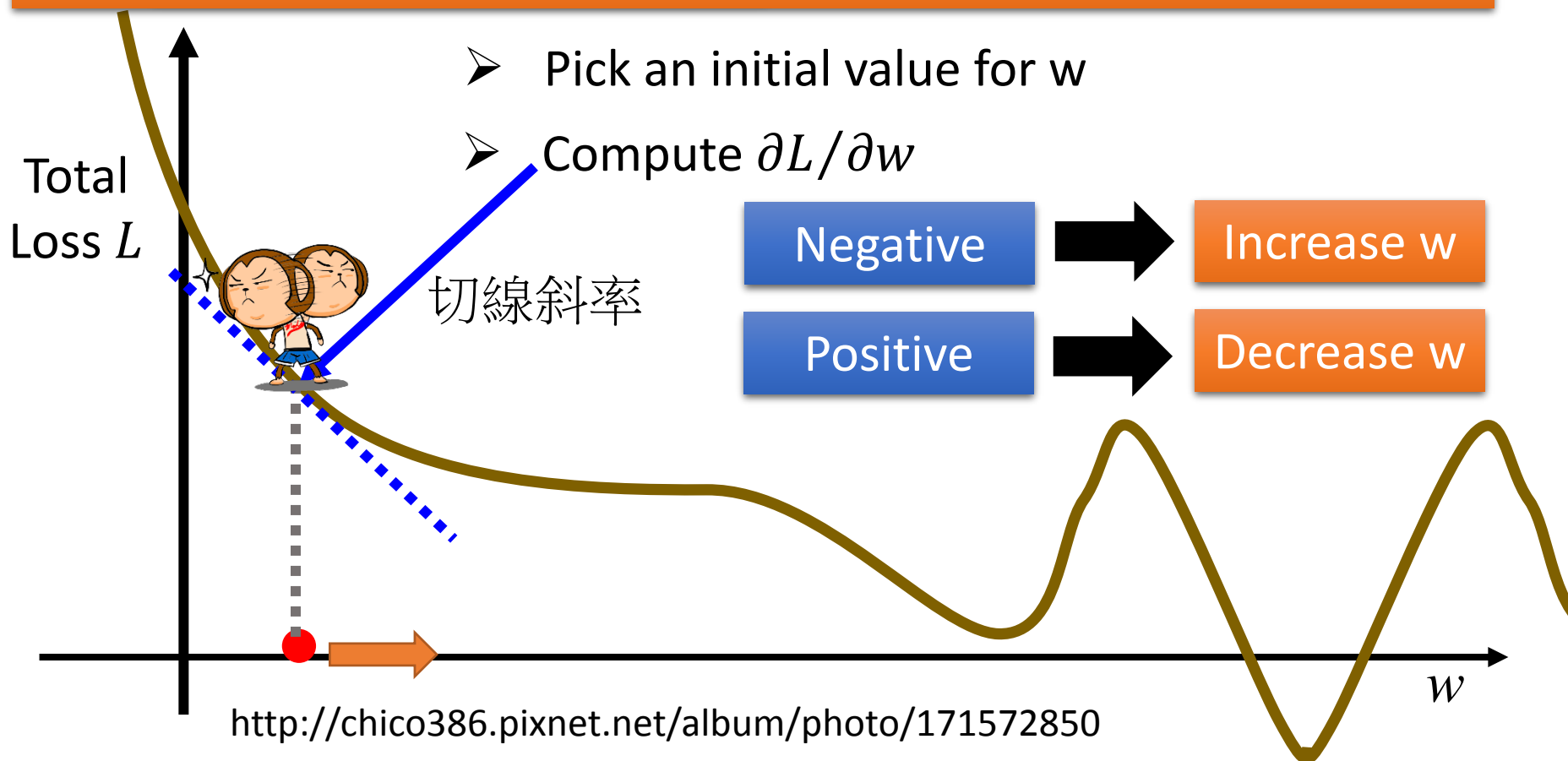
Usually good enough

Total Loss $L$

$w$

# Gradient Descent
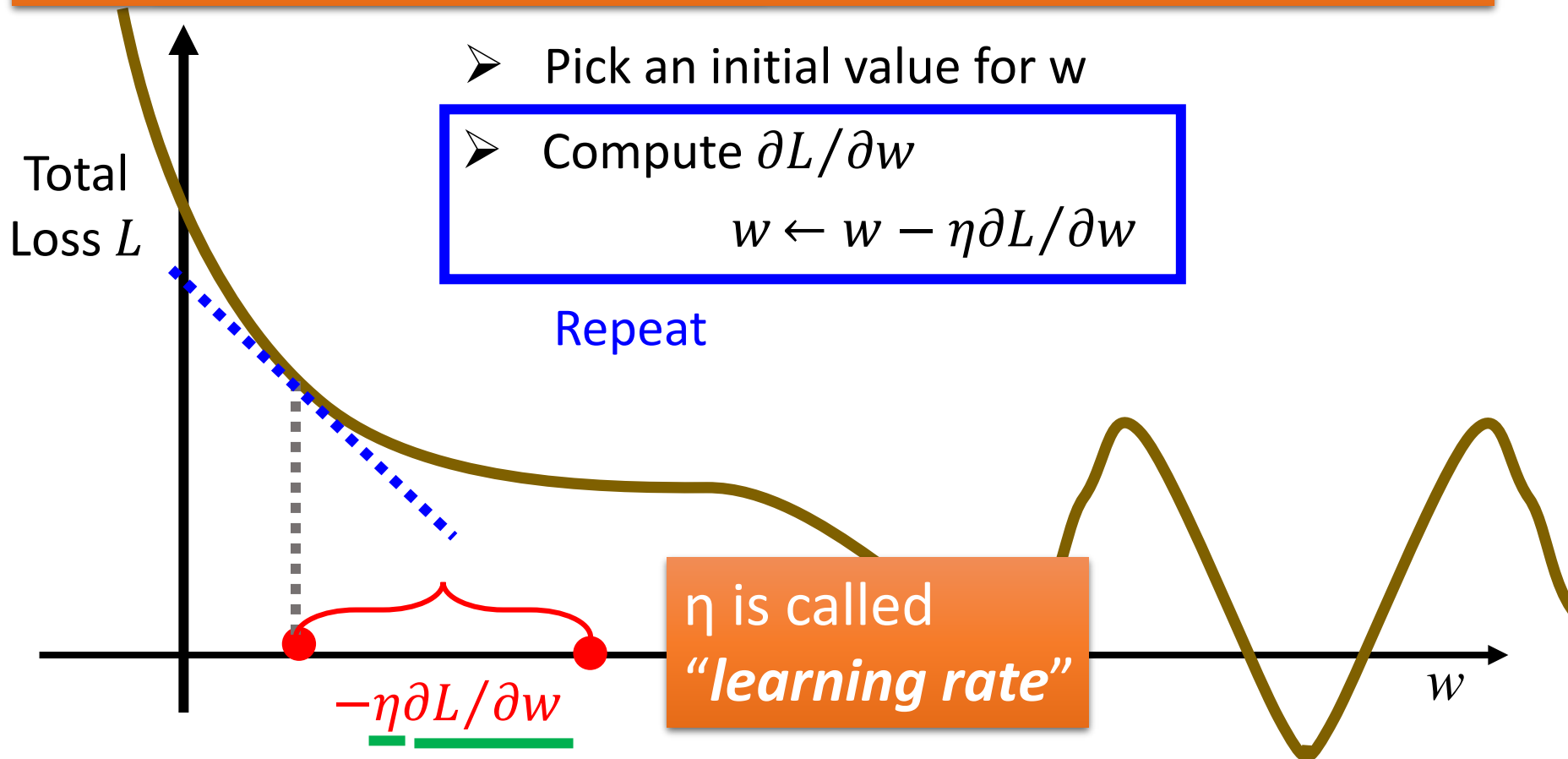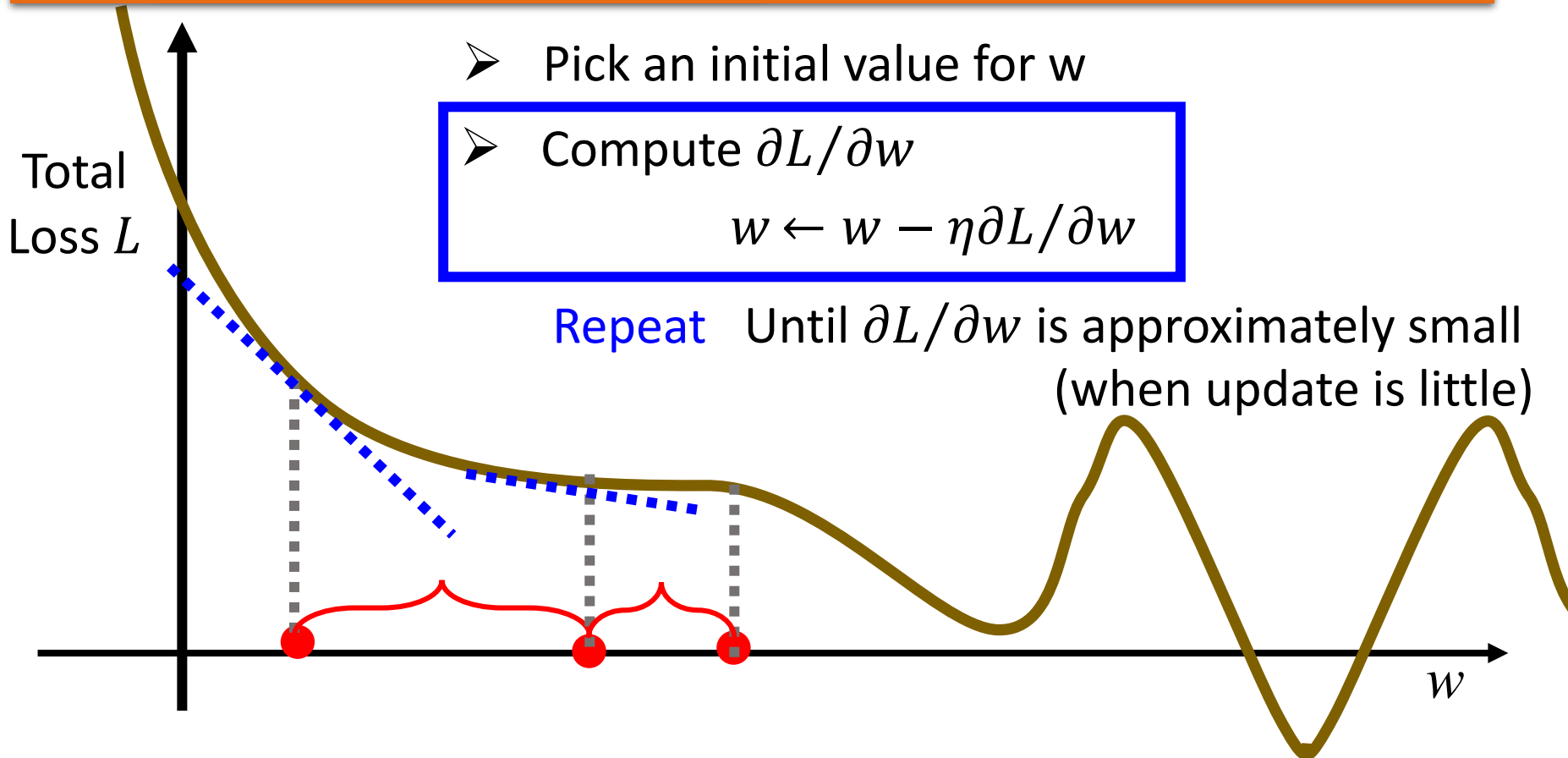
Network parameters $\theta = \{w_1, w_2, \cdots, b_1, b_2, \cdots\}$

Find **_network parameters $\theta^*$_** that minimize total loss L
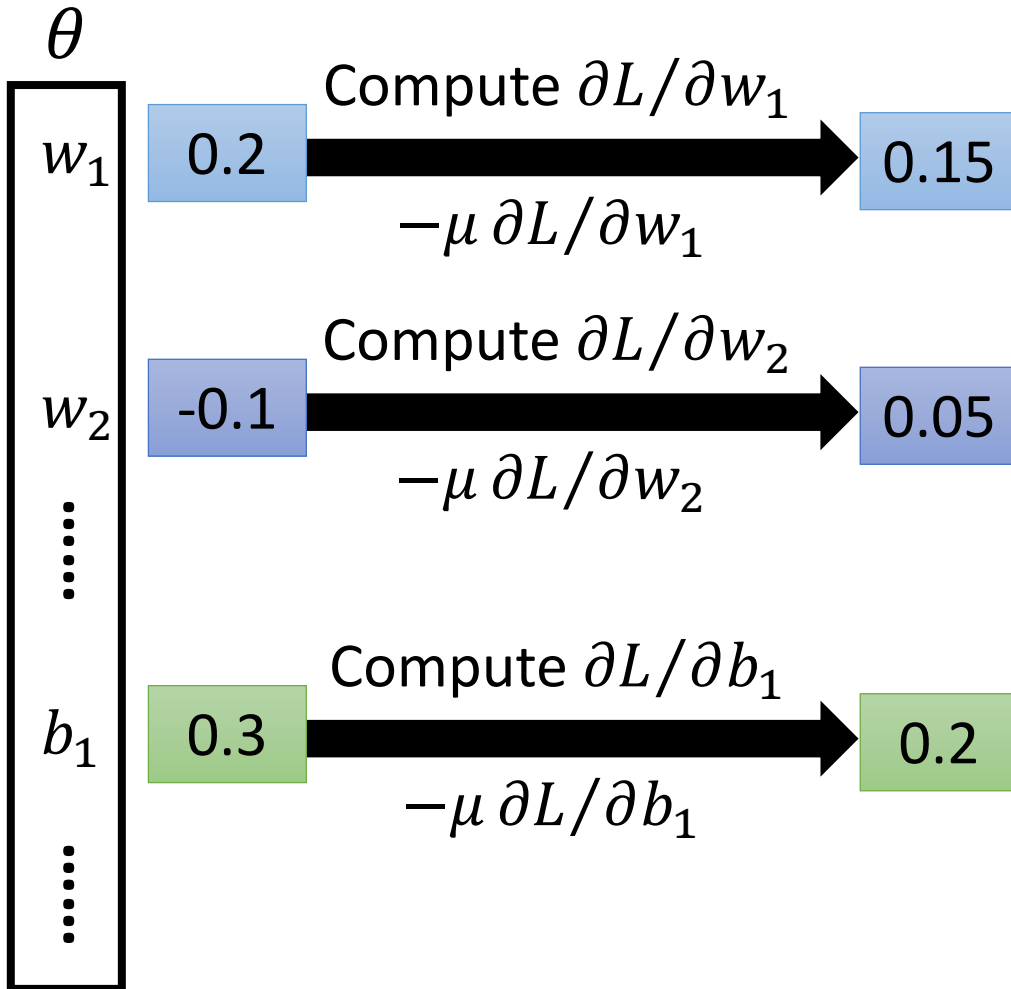
➢ Pick an initial value for w

➢ Compute $\partial L/\partial w$

Total Loss $L$

切線斜率

| Negative | → | Increase w |
| Positive | → | Decrease w |

$w$

http://chico386.pixnet.net/album/photo/171572850

# Gradient Descent

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

Find **_network parameters $\theta^*$_** that minimize total loss L

Total Loss $L$

➤ Pick an initial value for w

➤ Compute $\partial L/\partial w$

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

Repeat

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

η is called **_"learning rate"_**

$w$

# Gradient Descent

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

Find ***network parameters $\theta^*$*** that minimize total loss L

➢ Pick an initial value for w

➢ Compute $\partial L/\partial w$

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

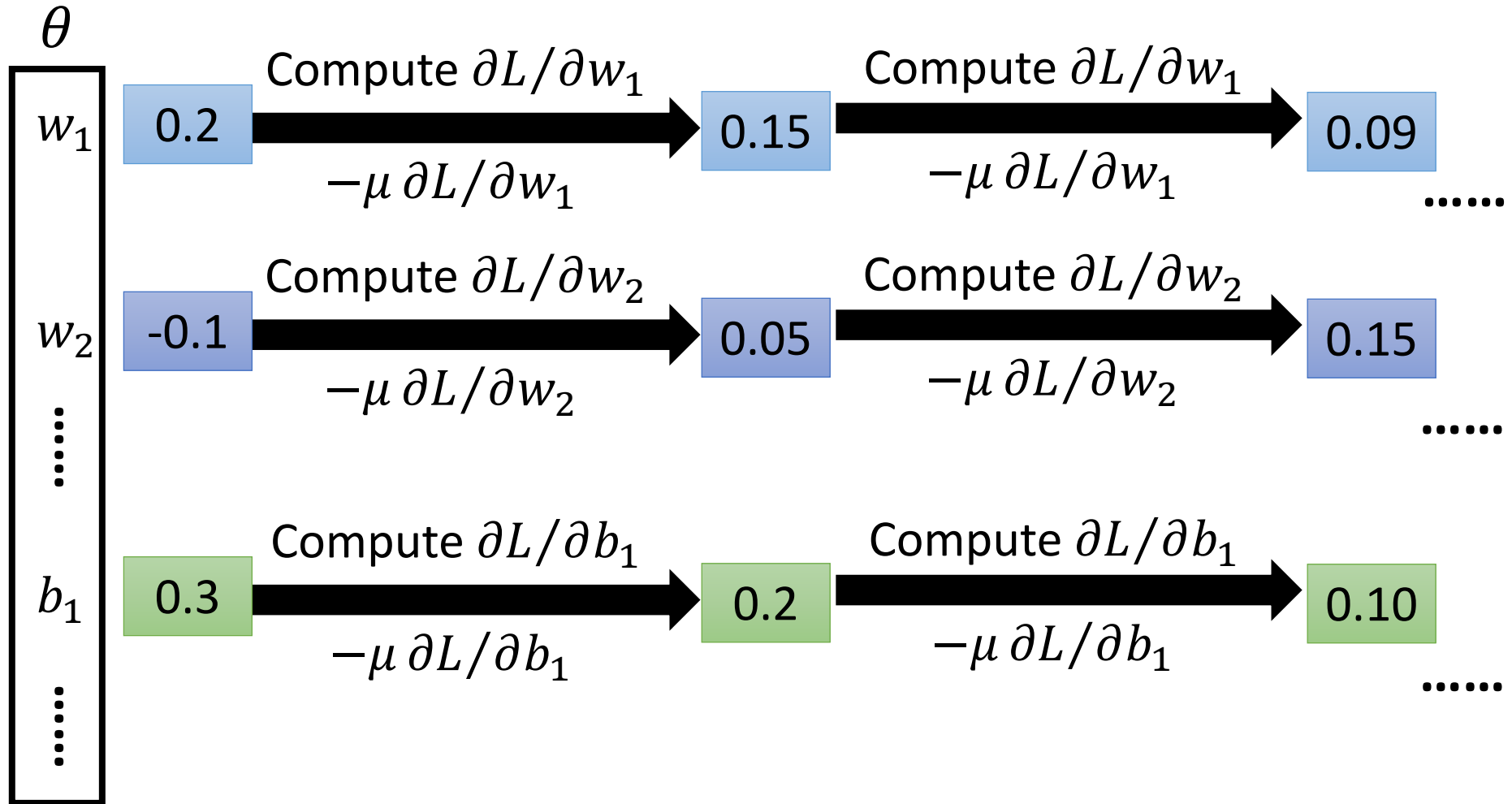Repeat   Until $\partial L/\partial w$ is approximately small (when update is little)

Total Loss $L$

$w$
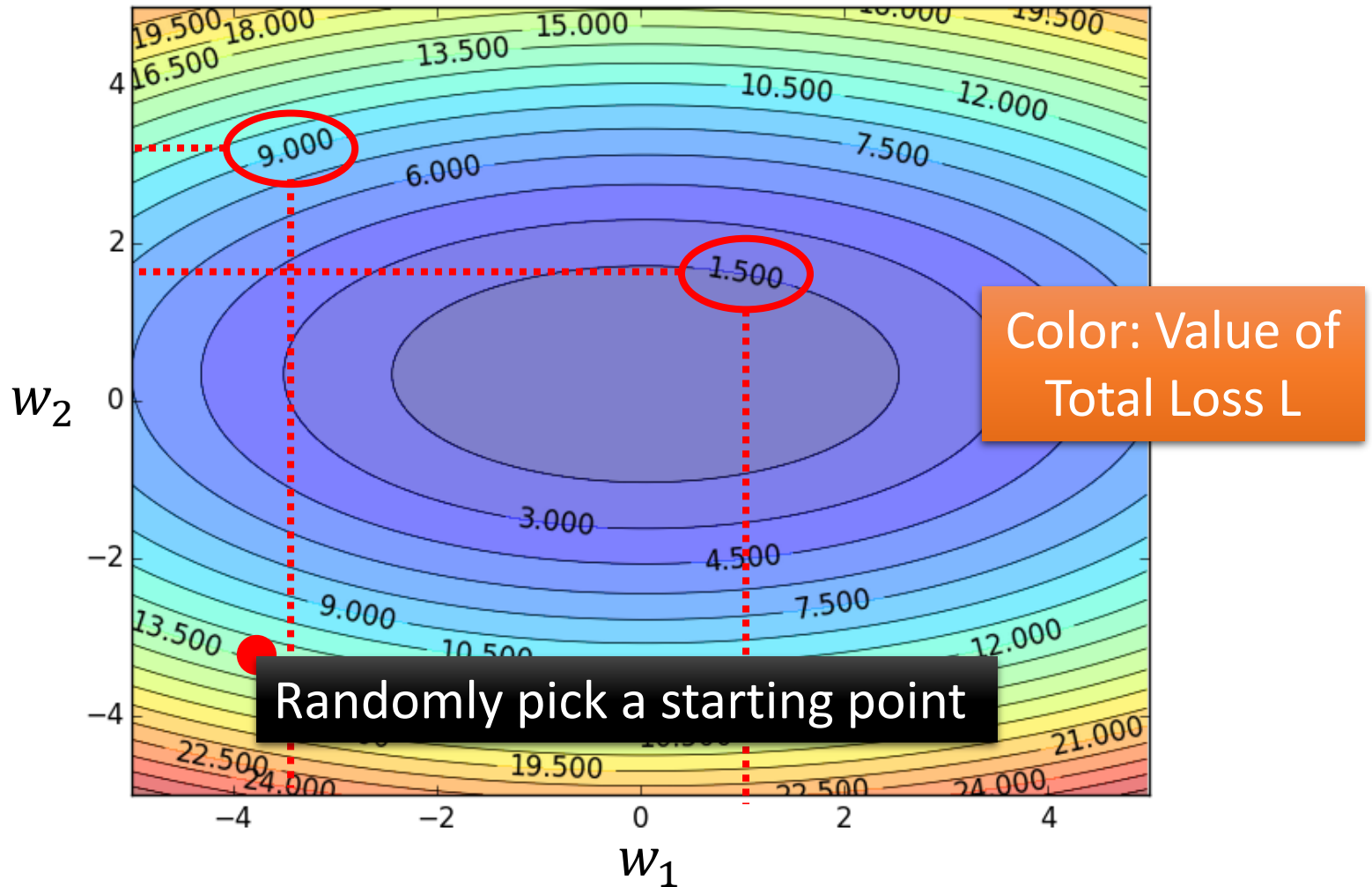
# Gradient Descent

$\theta$

$$\begin{bmatrix} w_1 \\ \\ w_2 \\ \vdots \\ \\ b_1 \\ \vdots \end{bmatrix}$$

Compute $\partial L / \partial w_1$

0.2 $\longrightarrow$ 0.15

$-\mu \, \partial L / \partial w_1$

Compute $\partial L / \partial w_2$

-0.1 $\longrightarrow$ 0.05

$-\mu \, \partial L / \partial w_2$

Compute $\partial L / \partial b_1$

0.3 $\longrightarrow$ 0.2

$-\mu \, \partial L / \partial b_1$

$$\nabla L = \begin{bmatrix} \dfrac{\partial L}{\partial w_1} \\[2ex] \dfrac{\partial L}{\partial w_2} \\[1ex] \vdots \\[1ex] \dfrac{\partial L}{\partial b_1} \\[1ex] \vdots \end{bmatrix}$$

gradient

# Gradient Descent

# Gradient Descent



Color: Value of Total Loss L

Randomly pick a starting point

# Gradient Descent

Hopfully, we would reach a minima …..



Color: Value of Total Loss L

$(-\eta \, \partial L / \partial w_1, \, -\eta \, \partial L / \partial w_2)$

Compute $\partial L / \partial w_1, \partial L / \partial w_2$

# Gradient Descent

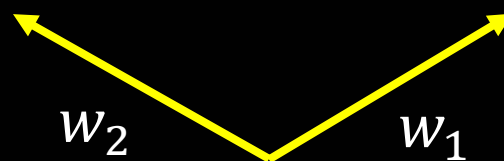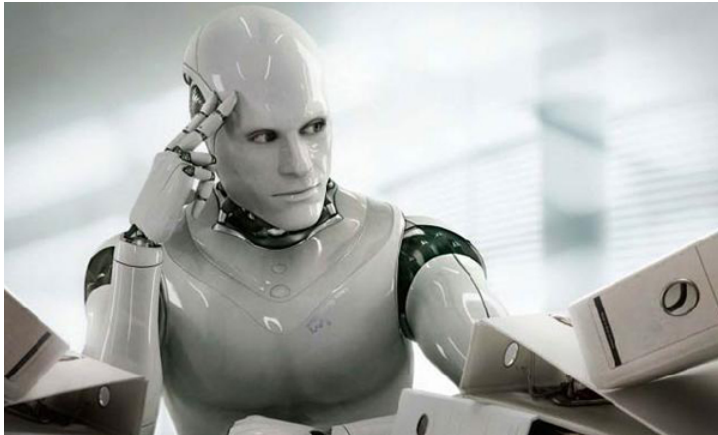- When considering multiple parameters together, do you see any problem?

# Local Minima

Who is Afraid of Non-Convex Loss Functions?
http://videolectures.net/eml07_lecun_wia/

- Gradient descent never guarantee global minima



Different initial point

Reach different minima, so different results

想像你在玩世紀帝國......
　　　　　沒有探索過的地方被戰霧覆蓋

$(-\eta \, \partial L / \partial w_1, -\eta \, \partial L / \partial w_2)$

Compute $\partial L / \partial w_1, \partial L / \partial w_2$

$w_2$　　　　　$w_1$

# Gradient Descent

This is the "learning" of machines in deep learning ......

➡ Even alpha go using this approach.

大家以為 Learning 是 ......          其實 Learning 只是 ......





I hope you are not too disappointed :p

# Backpropagation

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

  - Ref: http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS_2015_2/Lecture/DNN%20backprop.ecm.mp4/index.html



Don't worry about $\partial L / \partial w$, the toolkits will handle it.

# You can do lots of different things

```
┌─────────────┐      ┌─────────────┐      ┌─────────────┐
│   Step 1:   │  →   │   Step 2:   │  →   │   Step 3:   │
│   Network   │      │  Learning   │      │   Learn!    │
│  Structure  │      │   Target    │      │             │
└─────────────┘      └─────────────┘      └─────────────┘
       ↓                    ↓                    ↓
┌─────────────┐      ┌─────────────┐      ┌─────────────┐
│   A set of  │  →   │ Define the  │  →   │  Pick the   │
│  function   │      │ goodness of │      │    best     │
│             │      │ a function  │      │  function   │
└─────────────┘      └─────────────┘      └─────────────┘
```

# For example, you can do .......

- Image Recognition

# For example, you can do .......

**_Spam filtering_**

"Talk" in e-mail

EMAIL

Network

1/0

(Yes/No)

"free" in e-mail

1 (Yes)

SPAM

SPAM FILTER

0 (No)

(http://spam-filter-review.toptenreviews.com/)

# Document Classification



"stock" in document

"president" in document

Network

政治

經濟

體育

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

體育　　　　政治　　　　財經

# Playing Go



**19 x 19 image (matrix)**

Black: 1

white: -1

none: 0

→ **Network** →

Next move (19 x 19 positions)

19 x 19 vector

Fully-connected feedword network can be used

But CNN performs much better.

(Lecture III)

# Playing Go



蒐集一堆棋譜

Training:



Target:
天元 = 1
其他都是 0



Target:
五之 5 = 1
其他都是 0

# Concluding Remarks

- Deep Learning is simple & powerful!

但 Deep Learning 就像 雷神之槌

無法輕易被舉起來 ……

http://ent.ltn.com.tw/news/breaking news/1144545

# Lecture II:
# Tips for Training DNN

# Outline of Lecture II

"Hello World" for Deep Learning

Recipe of Deep Learning

# Keras

TensorFlow   or   theano

Very flexible

Need some effort to learn

Interface of TensorFlow or Theano

K

keras

Easy to learn and use
(still have some flexibility)

You can modify it if you can write TensorFlow or Theano

# Keras

- 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

# 使用 Keras 心得

# Example Application

- Handwriting Digit Recognition



28 x 28

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

"Hello world" for deep learning

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

# Keras

Step 1: Network Structure → Step 2: Learning Target → Step 3: Learn!



28x28

500

500

Softmax

$y_1$    $y_2$......    $y_{10}$

```
model = Sequential()
```

```
model.add( Dense( input_dim=28*28,
                  output_dim=500 ))
model.add( Activation('sigmoid') )
```

```
model.add( Dense( output_dim=500 ) )
model.add( Activation('sigmoid') )
```

```
model.add( Dense(output_dim=10 ) )
model.add( Activation('softmax') )
```

# Keras



Step 1: Network Structure → Step 2: Learning Target → **Step 3: Learn!**

Step 3.1: Configuration

```
model.compile(loss='mse',
              optimizer=SGD(lr=0.1),
              metrics=['accuracy'])
```

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

0.1

Step 3.2: Find the optimal network parameters

```
model.fit(x_train, y_train, batch_size=100, nb_epoch=20)
```
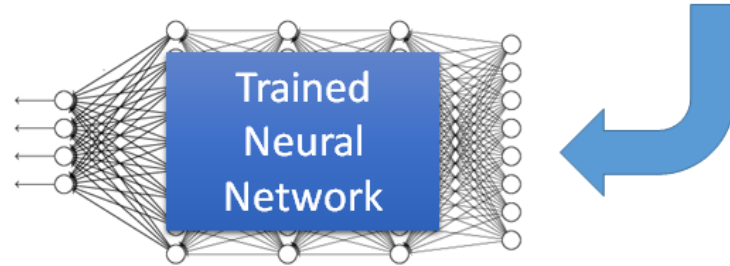
Training data
(Images)

Labels
(digits)

等一下再講

# Keras



Step 3.2: Find the optimal network parameters

```
model.fit(x_train, y_train, batch_size=100, nb_epoch=20)
```

numpy array

numpy array

28 x 28
=784

10

Number of training examples

Number of training examples

# Keras



Step 1: Network Structure → Step 2: Learning Target → Step 3: Learn!

Trained Neural Network

Save and load models

http://keras.io/getting-started/faq/#how-can-i-save-a-keras-model

How to use the neural network (testing):

case 1:
```
score = model.evaluate(x_test,y_test)
print('Total loss on Testing Set:', score[0])
print('Accuracy of Testing Set:', score[1])
```

case 2:
```
result = model.predict(x_test)
```

# Keras

- Using GPU to speed training
  - Way 1
    - THEANO_FLAGS=device=gpu0 python YourCode.py
  - Way 2 (in your code)
    - import os
    - os.environ["THEANO_FLAGS"] = "device=cpu"

# Outline of Lecture II

"Hello World" for Deep Learning

Recipe of Deep Learning

# *Recipe of Deep Learning*

# Do not always blame Overfitting



Not well trained

Training Data

Testing Data

Overfitting?

**Deep Residual Learning for Image Recognition**
http://arxiv.org/abs/1512.03385

# Recipe of Deep Learning



- Choosing proper loss
- Mini-batch
- New activation function
- Adaptive Learning Rate
- Momentum

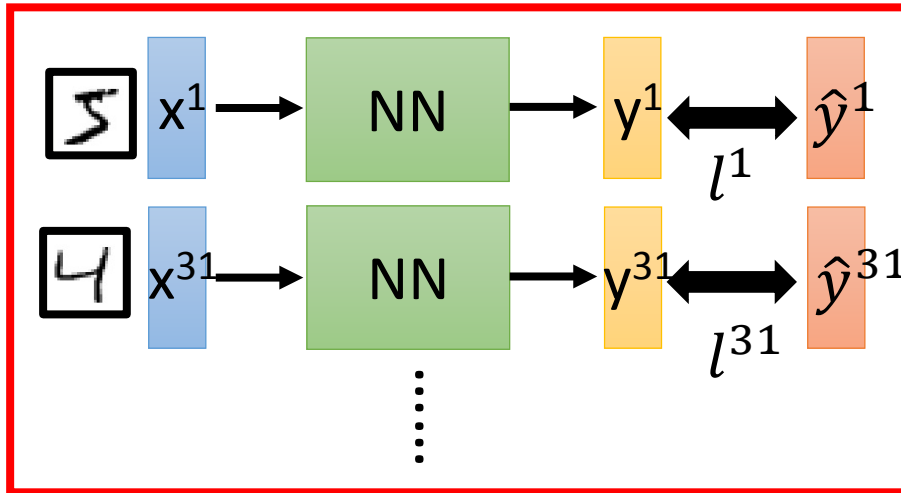Good Results on Testing Data?

YES

Good Results on Training Data?

YES

# Choosing Proper Loss

# Choosing Proper Loss

When using softmax output layer, choose cross entropy



Total Loss

Cross Entropy

Square Error

$w_1$

$w_2$

# Let's try it

Square Error

```
model.compile(loss='mse',
              optimizer=SGD(lr=0.1),
              metrics=['accuracy'])
```

Cross Entropy

```
model.compile(loss='categorical crossentropy',
              optimizer=SGD(lr=0.1),
              metrics=['accuracy'])
```

# Let's try it

Testing:

| | Accuracy |
|---|---|
| Square Error | 0.11 |
| Cross Entropy | 0.84 |


Training

# *Recipe of Deep Learning*

Choosing proper loss

Mini-batch

New activation function

Adaptive Learning Rate

Momentum

Good Results on Testing Data?

Good Results on Training Data?

YES

YES

```
model.fit(x_train, y_train, batch_size=100, nb_epoch=20)
```

We do not really minimize total loss!

# Mini-batch

- Randomly initialize network parameters
- Pick the 1st batch
  $$L' = l^1 + l^{31} + \cdots$$
  Update parameters once
- Pick the 2nd batch
  $$L'' = l^2 + l^{16} + \cdots$$
  Update parameters once
  $\vdots$
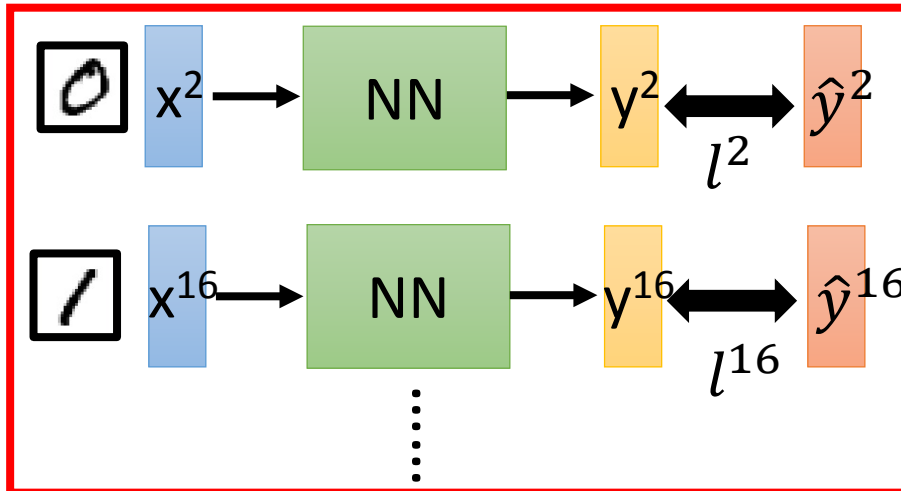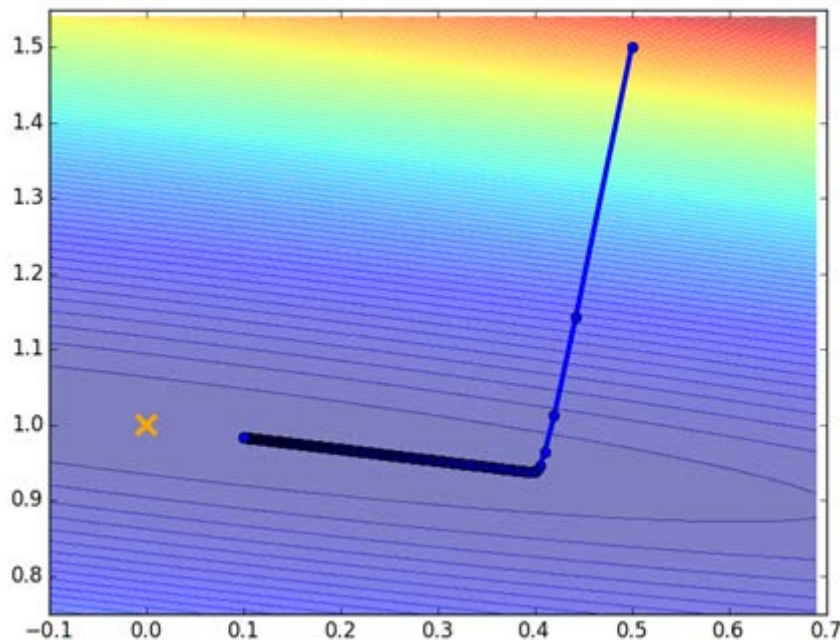- Until all mini-batches have been picked

one epoch

Repeat the above process

# Mini-batch

```
model.fit(x_train, y_train, batch_size=100, nb_epoch=20)
```

Mini-batch



100 examples in a mini-batch

Repeat 20 times

➢ Pick the 1st batch

$$L' = l^1 + l^{31} + \cdots$$

Update parameters once

➢ Pick the 2nd batch

$$L'' = l^2 + l^{16} + \cdots$$

Update parameters once

⋮

➢ Until all mini-batches have been picked

one epoch

# Mini-batch

**_Original Gradient Descent_**

**_With Mini-batch_**



Unstable!!!

The colors represent the total loss.

# Mini-batch is Better!

| | Accuracy |
|---|---|
| Mini-batch | 0.84 |
| No batch | 0.12 |



Training

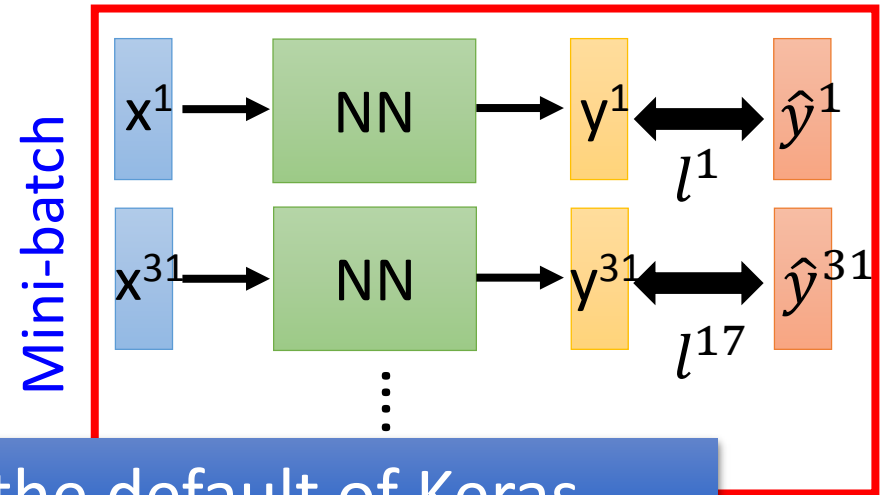Mini-batch

No batch

Accuracy

Epoch

# *Shuffle the training examples for each epoch*
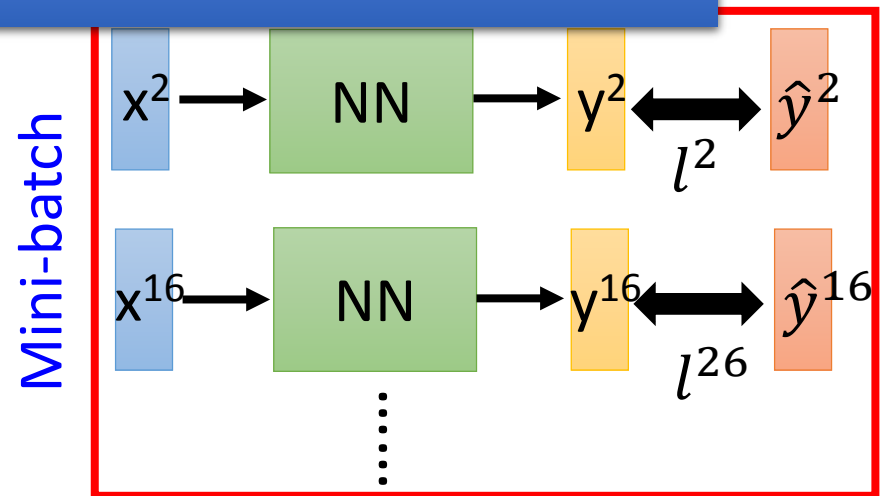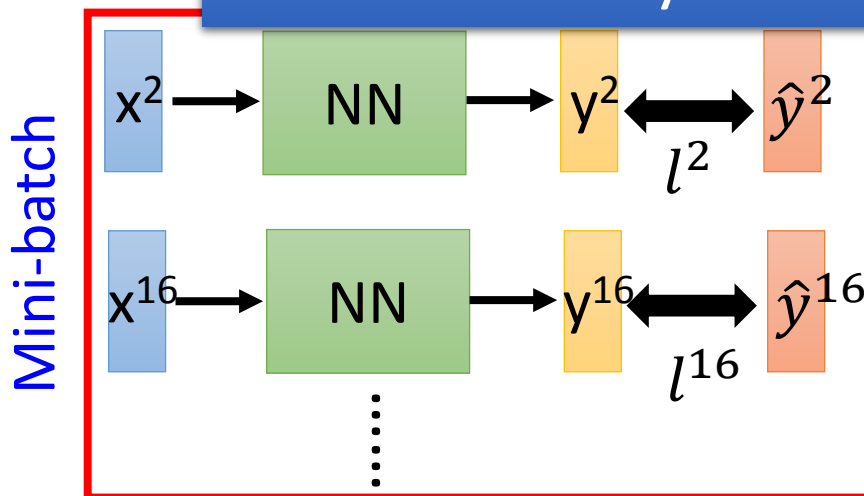
## *Epoch 1*



## *Epoch 2*

Don't worry. This is the default of Keras.

# Recipe of Deep Learning



- Choosing proper loss
- Mini-batch
- New activation function
- Adaptive Learning Rate
- Momentum

Good Results on Testing Data?

Good Results on Training Data?

YES

YES

# Hard to get the power of Deep …



Handwritting Digit Classification

Deeper usually does not imply better.

# Let's try it

| Testing: | Accuracy |
|---|---|
| 3 layers | 0.84 |
| 9 layers | 0.11 |



Training

3 layers

9 layers

# Vanishing Gradient Problem



Smaller gradients

Learn very slow

Almost random

Larger gradients

Learn very fast

Already converge

based on random!?

# Vanishing Gradient Problem

Smaller gradients

$x_1$

$x_2$

$\vdots$

$x_N$

$+\Delta w$

Small output

Large input

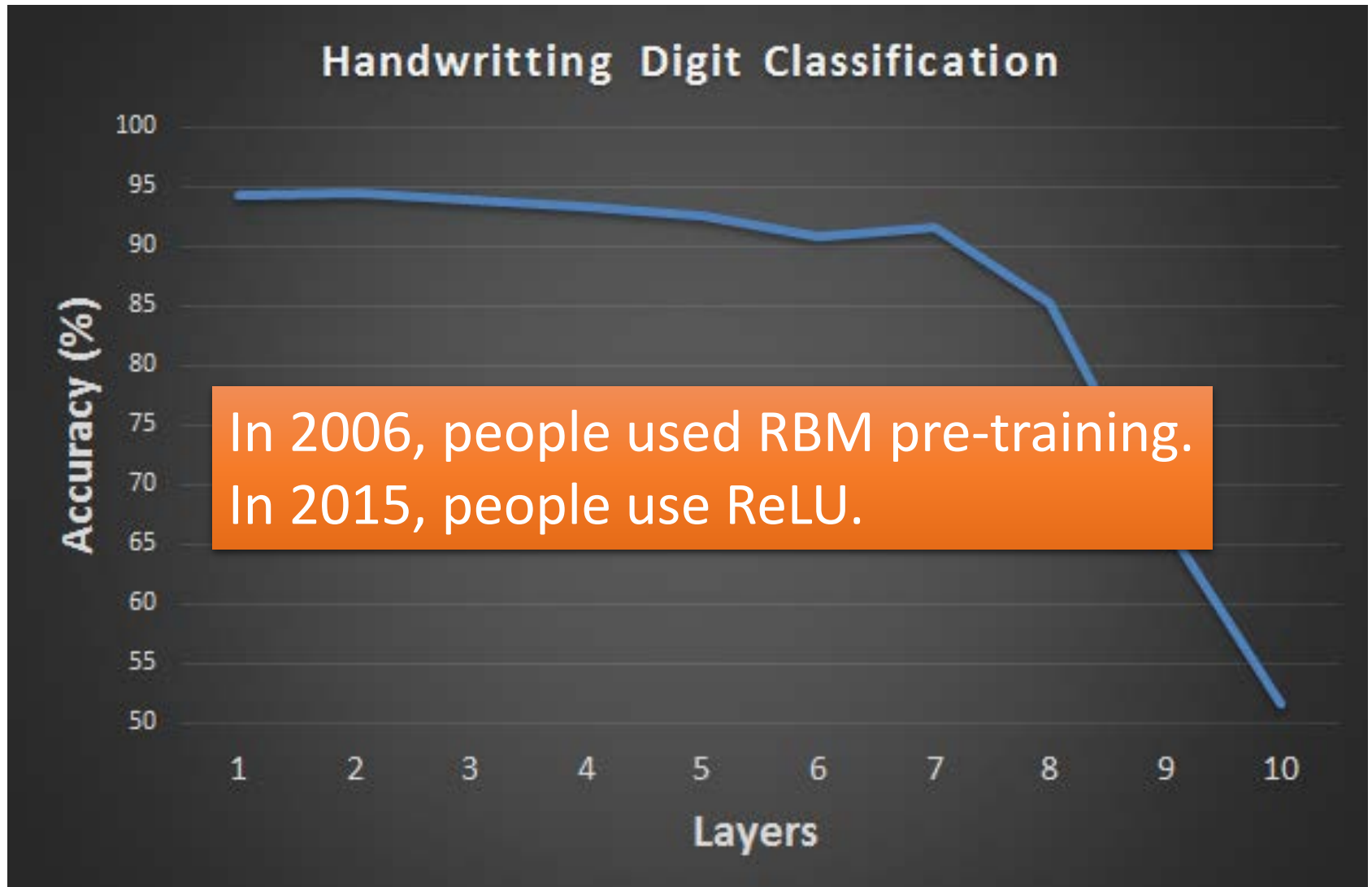Intuitive way to compute the derivatives …

$$\frac{\partial l}{\partial w} = ? \quad \frac{\Delta l}{\Delta w}$$

# Hard to get the power of Deep …



Handwritting Digit Classification

In 2006, people used RBM pre-training.
In 2015, people use ReLU.

# ReLU

- Rectified Linear Unit (ReLU)

$\sigma(z)$
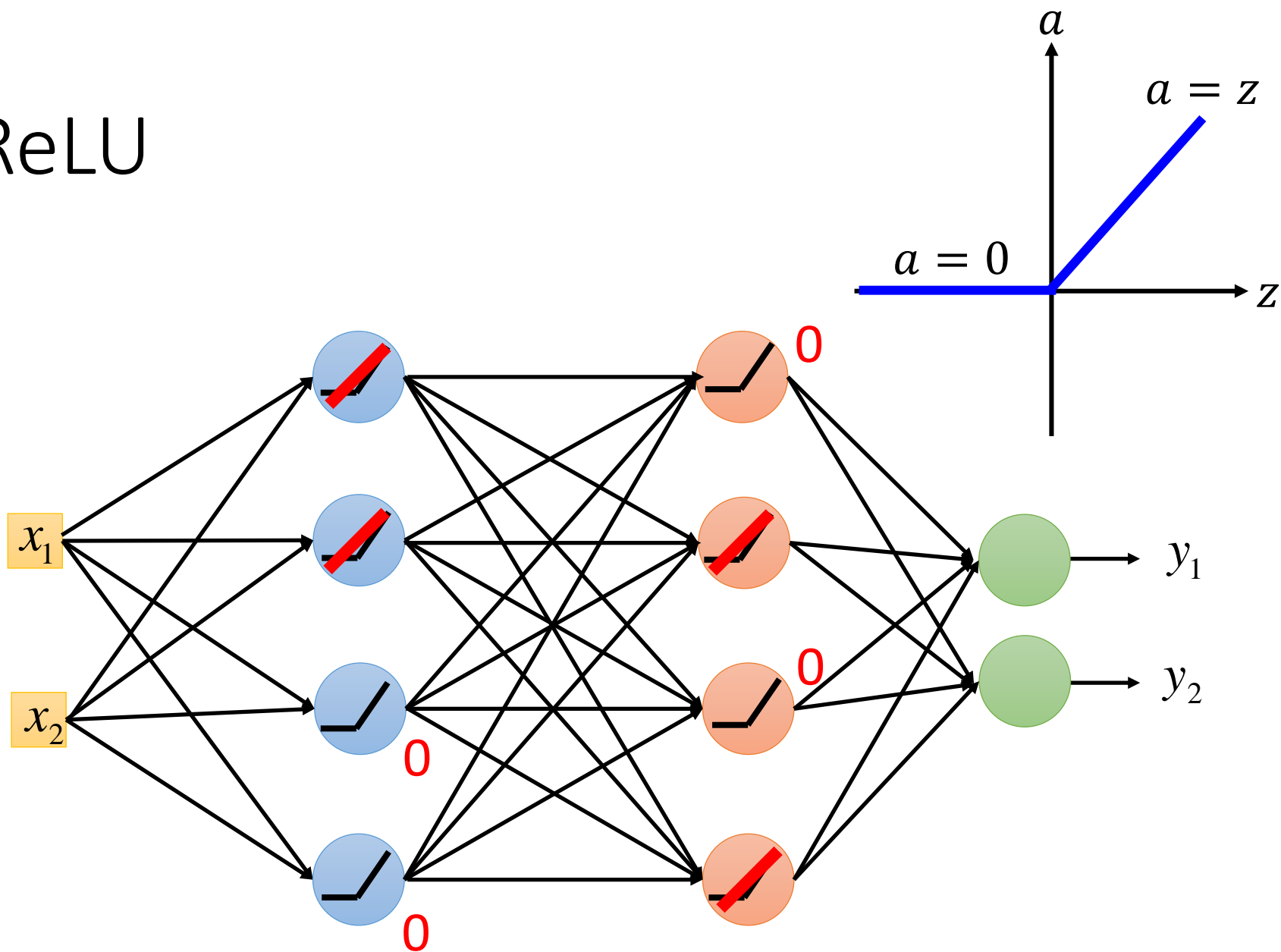


$a$

$a = z$

$a = 0$

$z$

[Xavier Glorot, AISTATS'11]
[Andrew L. Maas, ICML'13]
[Kaiming He, arXiv'15]

**_Reason:_**

1. Fast to compute

2. Biological reason

3. Infinite sigmoid with different biases
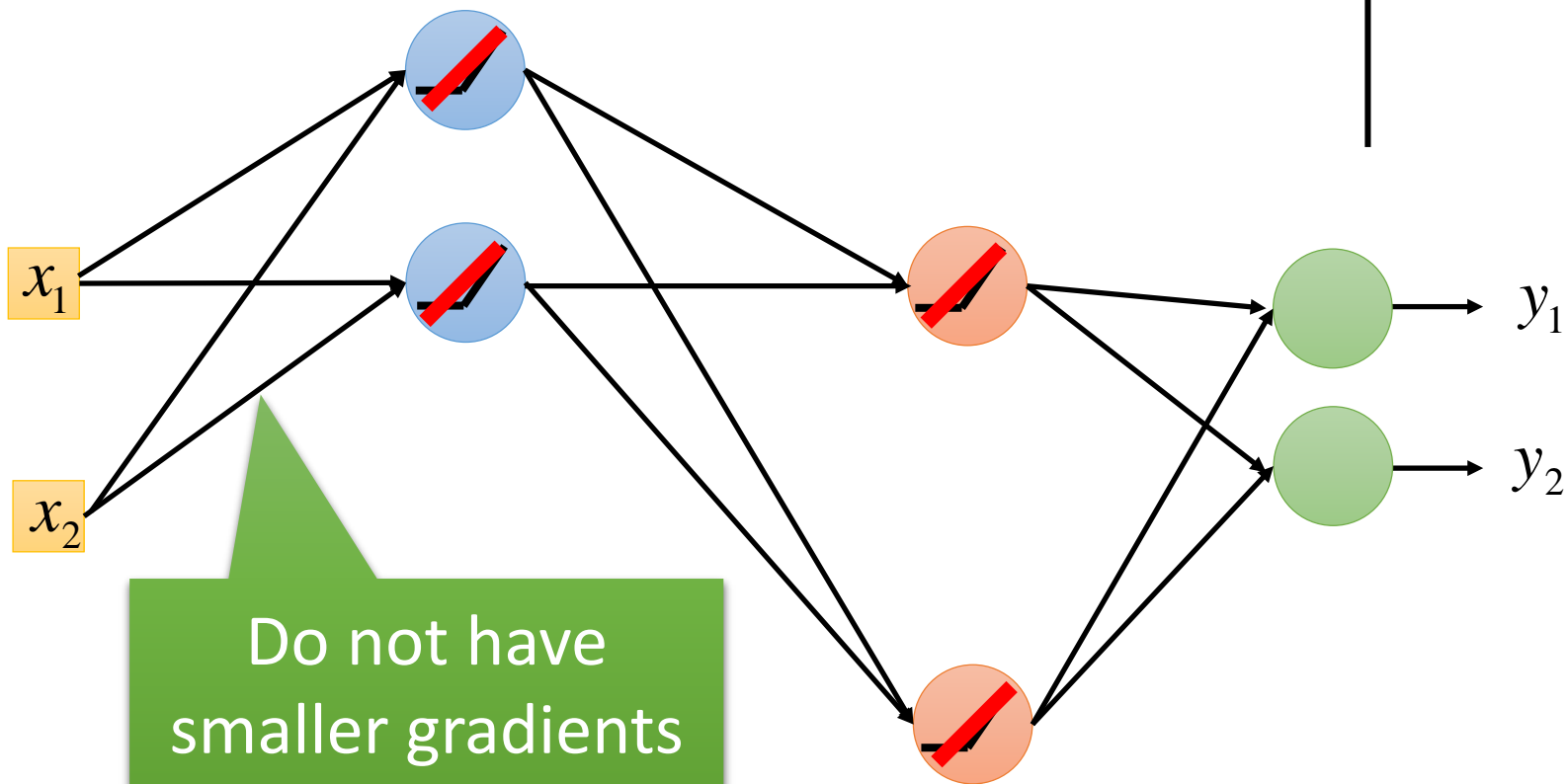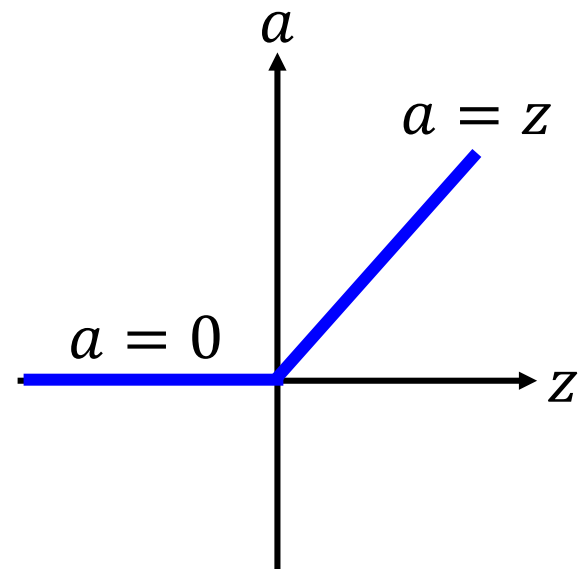
4. Vanishing gradient problem

ReLU

$a = z$

$a = 0$

$x_1$

$x_2$

$y_1$

$y_2$

0

0

0

0

ReLU

A Thinner linear network

Do not have smaller gradients

$x_1$

$x_2$

$y_1$

$y_2$

$a$

$a = z$

$a = 0$

$z$

# Let's try it

```
model.add( Activation('sigmoid') )
```

⬇

```
model.add( Activation('relu') )
```

# Let's try it

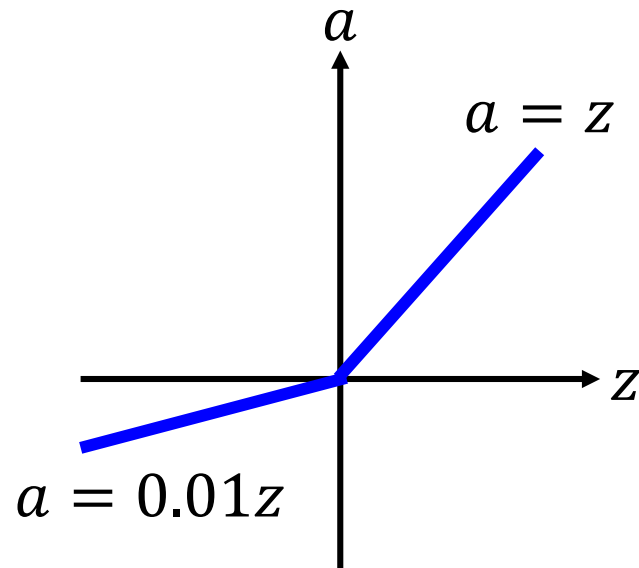| 9 layers | Accuracy |
|----------|----------|
| Sigmoid | 0.11 |
| ReLU | 0.96 |

- 9 layers

# ReLU - variant

## Leaky ReLU

$a$

$a = z$

$z$

$a = 0.01z$
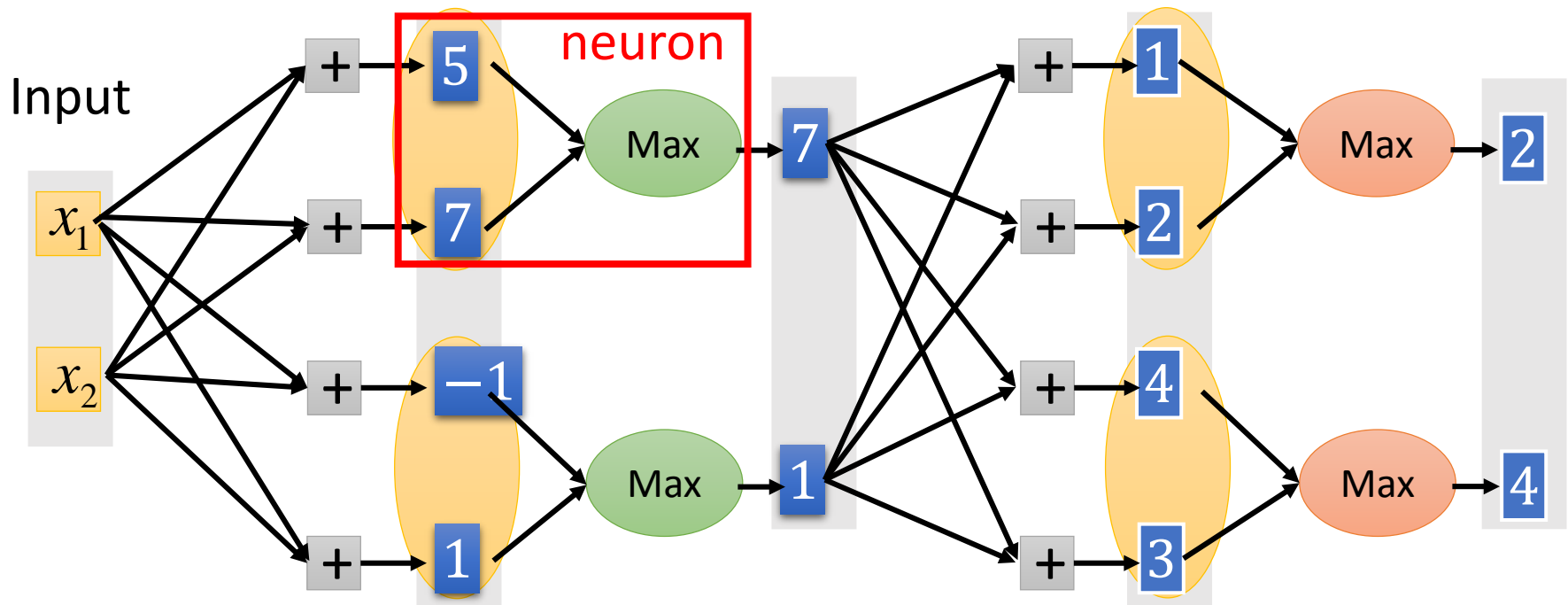
## Parametric ReLU

$a$

$a = z$

$z$

$a = \alpha z$

α also learned by gradient descent

# Maxout

ReLU is a special cases of Maxout

• Learnable activation function [Ian J. Goodfellow, ICML'13]



You can have more than 2 elements in a group.

# Maxout

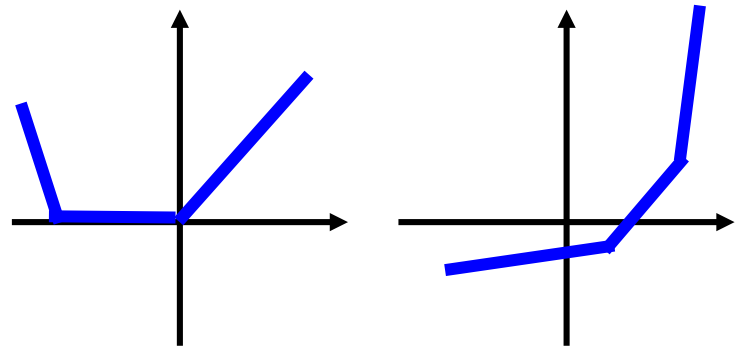- Learnable activation function [Ian J. Goodfellow, ICML'13]
  - Activation function in maxout network can be any piecewise linear convex function
  - How many pieces depending on how many elements in a group

**2 elements in a group**

**3 elements in a group**

# Recipe of Deep Learning



- Choosing proper loss
- Mini-batch
- New activation function
- Adaptive Learning Rate
- Momentum

Good Results on Testing Data? — YES

Good Results on Training Data? — YES

# Learning Rates

Set the learning rate η carefully



If learning rate is too large

Total loss may not decrease after each update

# Learning Rates

Set the learning rate η carefully

If learning rate is too large

Total loss may not decrease after each update

If learning rate is too small

Training would be too slow

# Learning Rates

- Popular & Simple Idea: Reduce the learning rate by some factor every few epochs.
  - At the beginning, we are far from the destination, so we use larger learning rate
  - After several epochs, we are close to the destination, so we reduce the learning rate
  - E.g. 1/t decay: $\eta^t = \eta / \sqrt{t+1}$

- Learning rate cannot be one-size-fits-all
  - Giving different parameters different learning rates

# Adagrad

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

Adagrad: $w \leftarrow w - \boxed{\eta_w} \partial L / \partial w$

Parameter dependent learning rate

$$\eta_w = \frac{\eta}{\sqrt{\sum_{i=0}^{t}(g^i)^2}}$$

constant

$g^i$ is $\partial L / \partial w$ obtained at the i-th update

Summation of the square of the previous derivatives

# Adagrad

$$\eta_w = \frac{\eta}{\sqrt{\sum_{i=0}^{t}(g^i)^2}}$$

$w_1$ | **g⁰** |
| 0.1 |

$w_2$ | **g⁰** |
| 20.0 |

Learning rate:

$$\frac{\eta}{\sqrt{0.1^2}} = \frac{\eta}{0.1}$$

$$\frac{\eta}{\sqrt{0.1^2 + 0.2^2}} = \frac{\eta}{0.22}$$

Learning rate:

$$\frac{\eta}{\sqrt{20^2}} = \frac{\eta}{20}$$

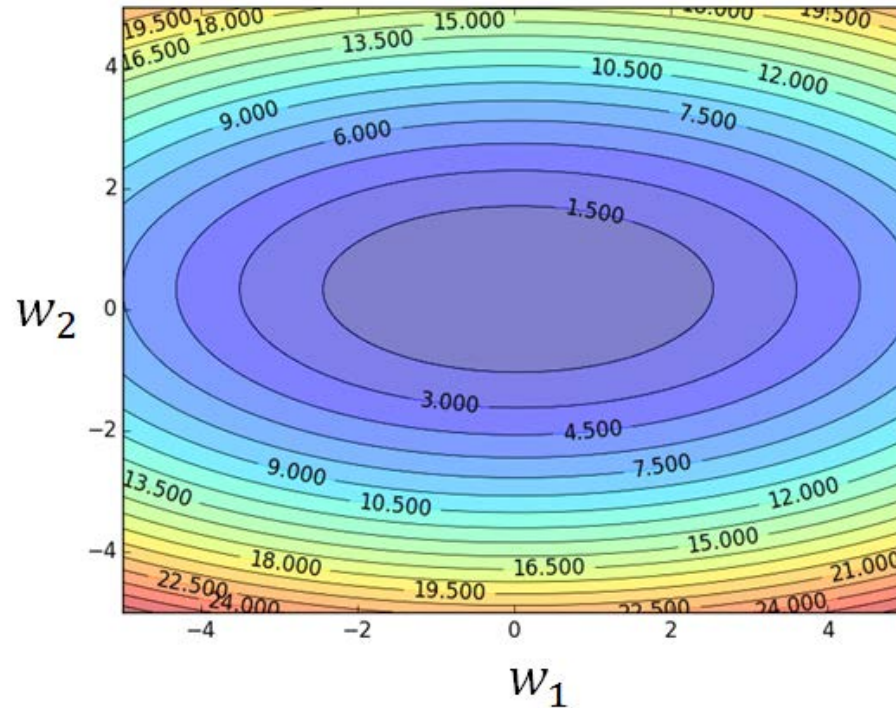$$\frac{\eta}{\sqrt{20^2 + 10^2}} = \frac{\eta}{22}$$

***Observation:*** 1. Learning rate is smaller and smaller for all parameters

2. Smaller derivatives, larger learning rate, and vice versa

Why?

Larger derivatives

Smaller Learning Rate

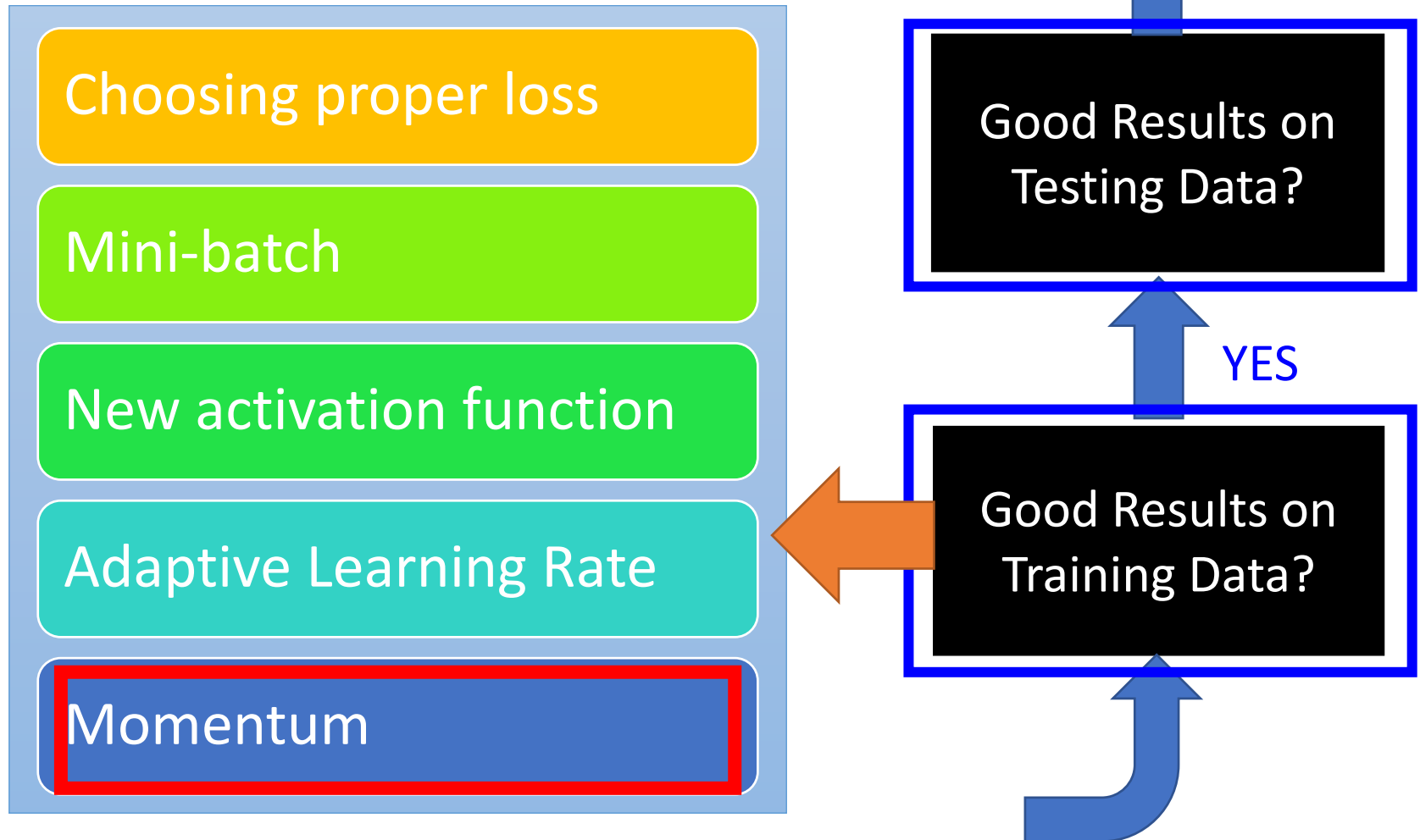Smaller Derivatives

Larger Learning Rate

$w_2$

$w_1$

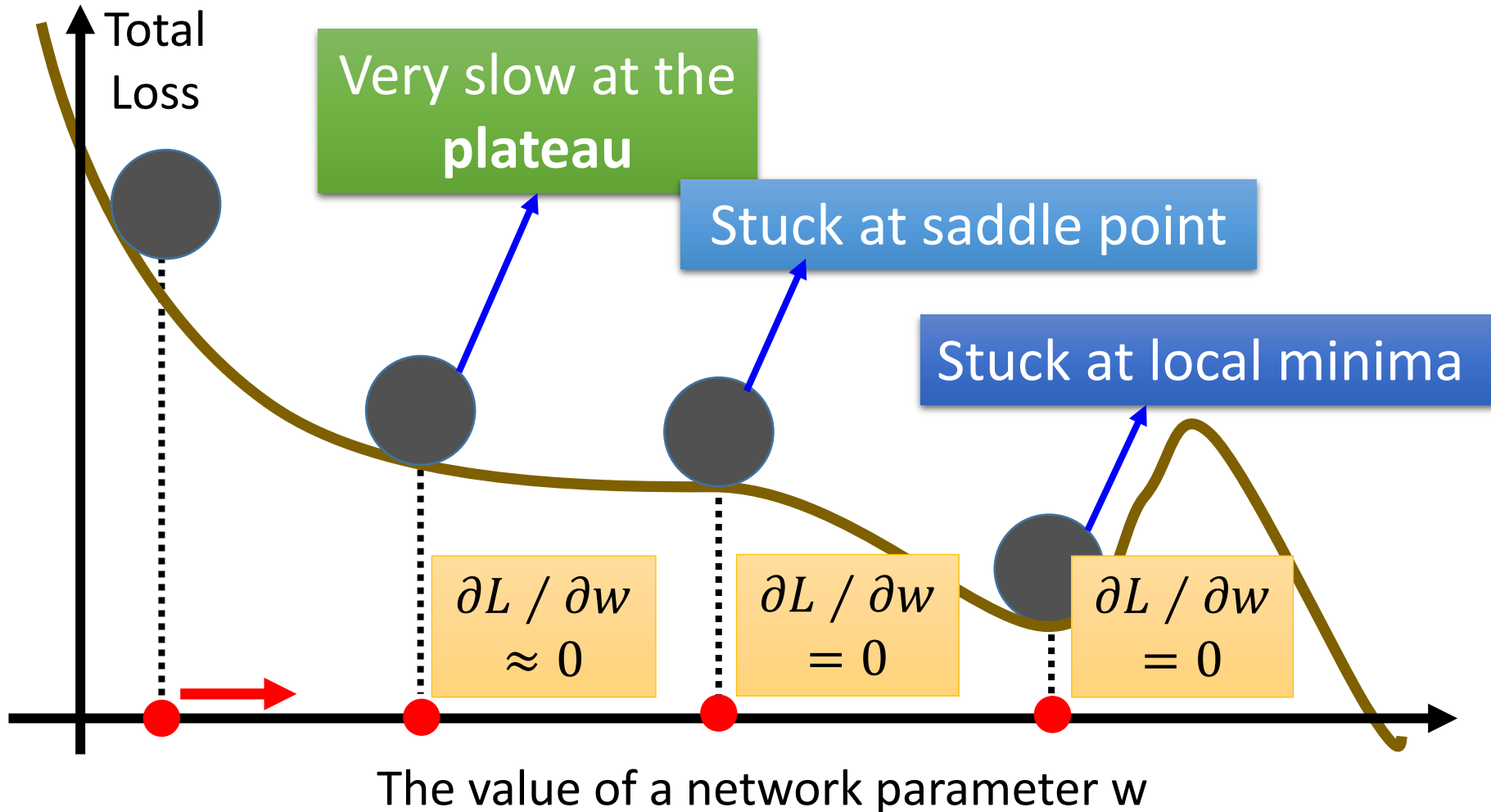2. Smaller derivatives, larger learning rate, and vice versa

Why?

# Not the whole story ……

- Adagrad [John Duchi, JMLR'11]
- RMSprop
  - https://www.youtube.com/watch?v=O3sxAc4hxZU
- Adadelta [Matthew D. Zeiler, arXiv'12]
- "No more pesky learning rates" [Tom Schaul, arXiv'12]
- AdaSecant [Caglar Gulcehre, arXiv'14]
- Adam [Diederik P. Kingma, ICLR'15]
- Nadam
  - http://cs229.stanford.edu/proj2015/054_report.pdf

# Hard to find optimal network parameters

# In physical world ……

- Momentum

How about put this phenomenon in gradient descent?

# Momentum

Still not guarantee reaching global minima, but give some hope ......

Movement =
Negative of $\partial L / \partial w$ + Momentum

→ Negative of $\partial L / \partial w$

⇢ Momentum

→ Real Movement

$\partial L / \partial w = 0$

# Adam

```
model.compile(loss='categorical_crossentropy',
              optimizer=SGD(lr=0.1),
              metrics=['accuracy'])
```

```
model.compile(loss='categorical_crossentropy',
              optimizer=Adam(),
              metrics=['accuracy'])
```

**Algorithm 1:** *Adam*, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation. $g_t^2$ indicates the elementwise square $g_t \odot g_t$. Good default settings for the tested machine learning problems are $\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. All operations on vectors are element-wise. With $\beta_1^t$ and $\beta_2^t$ we denote $\beta_1$ and $\beta_2$ to the power $t$.

**Require:** $\alpha$: Stepsize
**Require:** $\beta_1, \beta_2 \in [0, 1)$: Exponential decay rates for the moment estimates
**Require:** $f(\theta)$: Stochastic objective function with parameters $\theta$
**Require:** $\theta_0$: Initial parameter vector
  $m_0 \leftarrow 0$ (Initialize 1$^{\text{st}}$ moment vector)
  $v_0 \leftarrow 0$ (Initialize 2$^{\text{nd}}$ moment vector)
  $t \leftarrow 0$ (Initialize timestep)
  **while** $\theta_t$ not converged **do**
    $t \leftarrow t + 1$
    $g_t \leftarrow \nabla_\theta f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep $t$)
    $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate)
    $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate)
    $\widehat{m}_t \leftarrow m_t/(1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)
    $\widehat{v}_t \leftarrow v_t/(1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)
    $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \widehat{m}_t/(\sqrt{\widehat{v}_t} + \epsilon)$ (Update parameters)
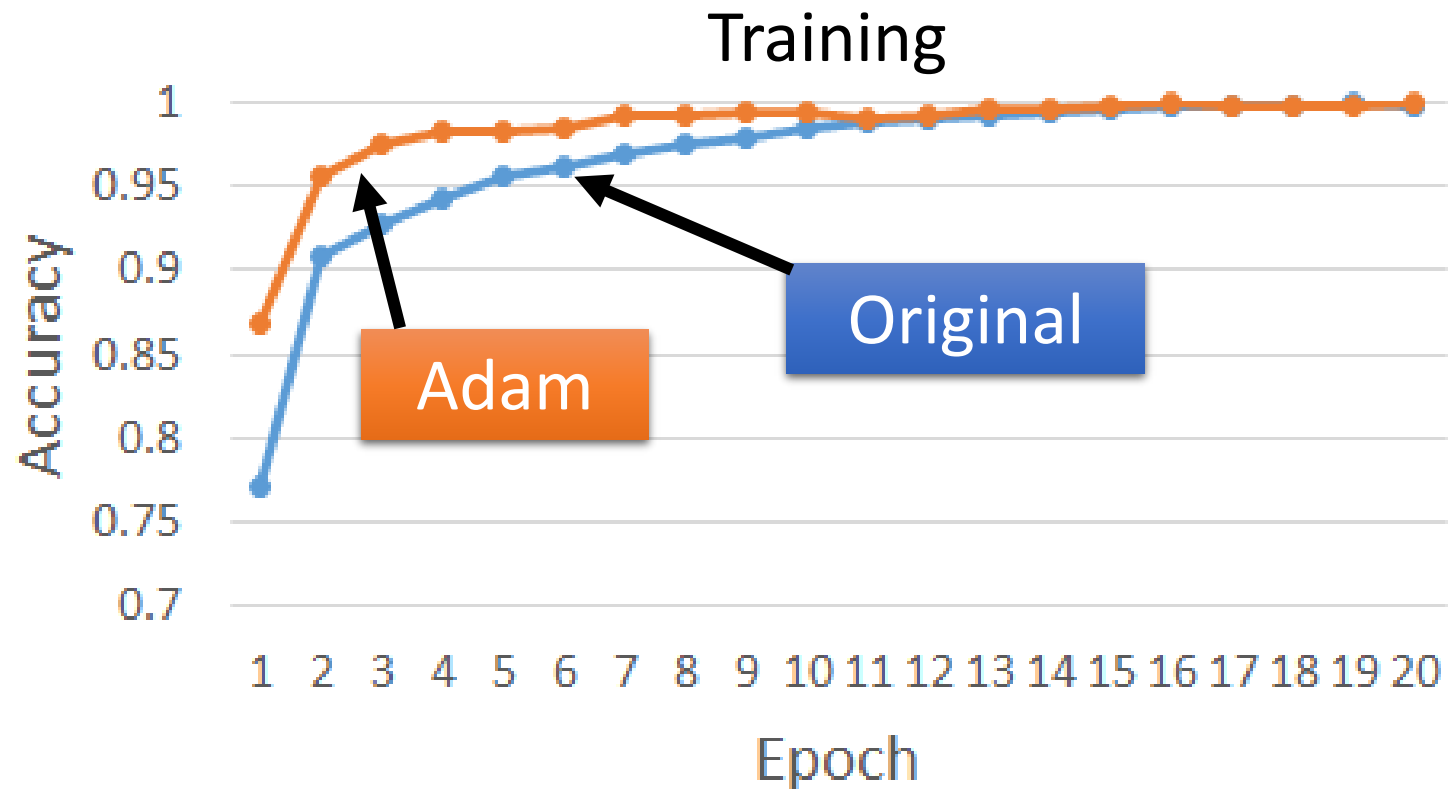  **end while**
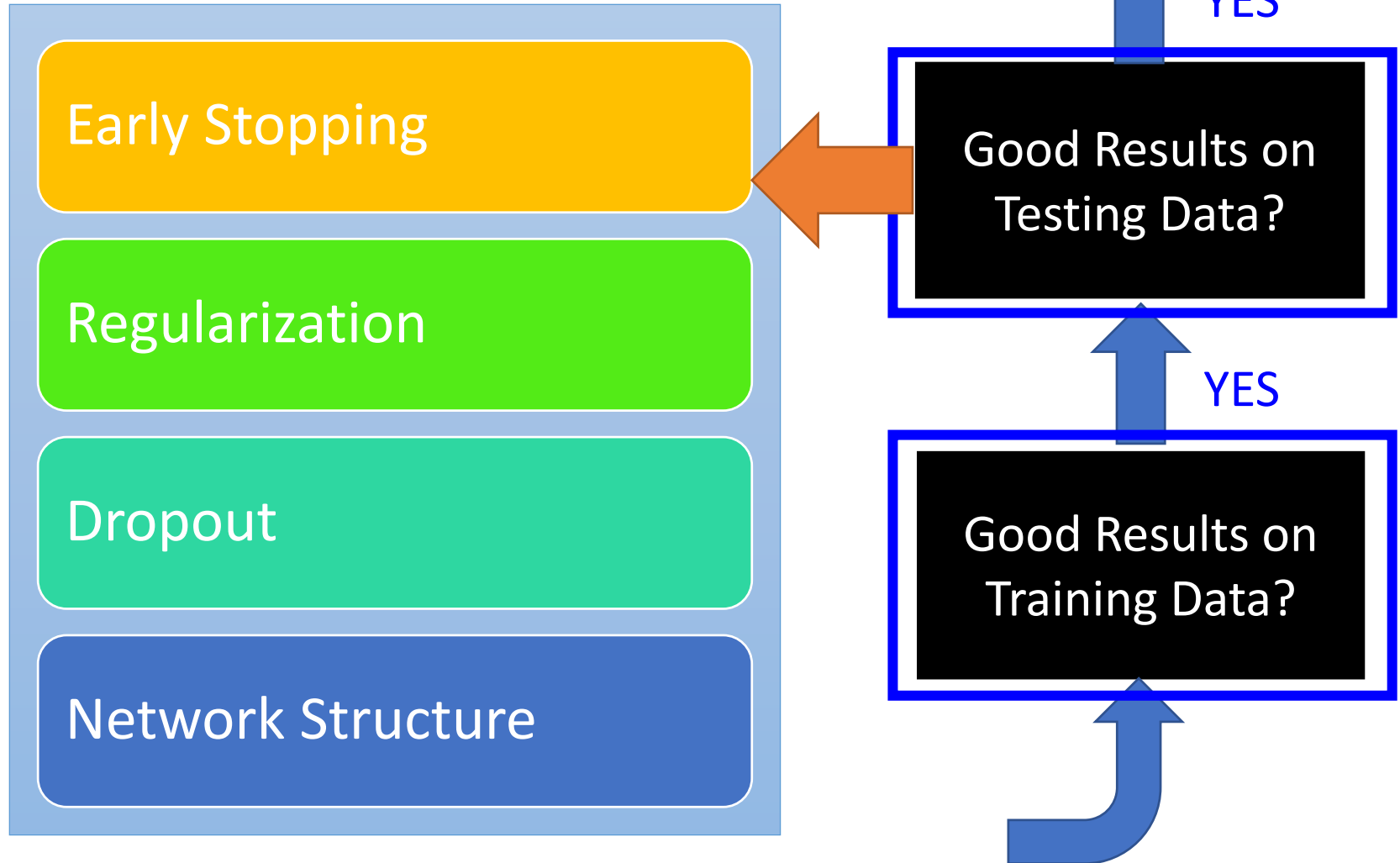  **return** $\theta_t$ (Resulting parameters)

# Let's try it

Testing:

| | Accuracy |
|---|---|
| Original | 0.96 |
| Adam | 0.97 |

- ReLU, 3 layer

Training

# Recipe of Deep Learning



Early Stopping

Regularization

Dropout

Network Structure
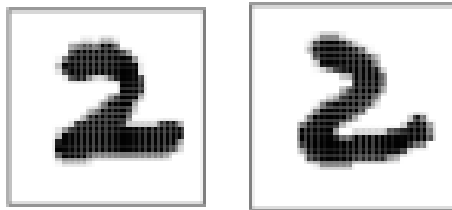
Good Results on Testing Data?
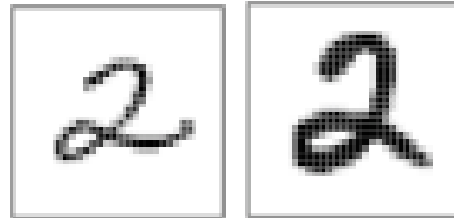
Good Results on Training Data?

YES

YES

# Why Overfitting?

- Training data and testing data can be different.

Training Data:



Testing Data:
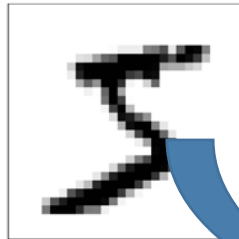


Learning target is defined by the training data.

The parameters achieving the learning target do not necessary have good results on the testing data.

# Panacea for Overfitting

- Have more training data
- **Create** more training data (?)

Handwriting recognition:

Original Training Data:

Created Training Data:

Shift 15°

# Why Overfitting?

- For experiments, we added some noises to the testing data

```
-1.36230370e-01,    1.03749340e-01,    1.15432226e-01,
 2.58670464e-01,    1.48774333e+00,    1.92885328e+00,
 1.70038673e+00,    2.46242981e+00,    1.21244572e+00,
-9.28660713e-01,    3.63209761e-01,   -1.81327713e+00,
-1.97910760e-01,    4.32874592e-01,   -5.40565788e-01,
 2.95630655e-01,    2.07984424e+00,   -1.84243292e+00,
-5.11166017e-01,   -5.80935128e-01,    1.06273647e+00,
 1.80551097e-02,    2.27983997e-02,   -1.67979148e+00,
 8.12423001e-01,   -6.25888706e-01,   -1.25027082e+00,
 6.15135458e-01,   -1.21394611e-01,   -1.28089527e+00,
 3.24609806e-01,    6.70569391e-01,    1.49161323e-01,
 8.01573609e-01,    6.43116741e-01,   -9.37629233e-02,
 1.74677366e+00,    6.80996008e-01,   -7.03717611e-01,
 1.02079749e-01,    1.19505614e+00,   -2.77959386e-01,
-5.21652916e-02,    3.53683601e-01,   -4.08310762e-01,
-1.81042967e+00,   -9.03308062e-01,    1.05404509e+00,
-9.80876877e-01,    3.52078891e-01,    6.65981840e-01,
 1.06550150e+00,   -2.28433613e-01,    3.64483904e-01,
-1.51484666e+00,   -7.52612872e-02,   -2.97058082e-01,
-7.27414382e-01,   -2.45875340e-01,   -1.27948942e-01,
-3.69310620e-01,   -2.62300428e+00,    2.11585073e+00,
 6.85561585e-01,   -1.57443985e-01,    1.38128777e+00,
 6.84265587e-02,    3.12536292e-01,    4.54253185e-01,
-7.88471875e-01,   -6.58403343e-02,   -1.41847985e+00,
-1.39753340e-01,   -5.55354856e-01,   -5.01917779e-01,
 6.93118522e-01,   -2.45360497e-01,   -1.26943186e+00,
-2.62323855e-01])

n [3]: x test[0]
```
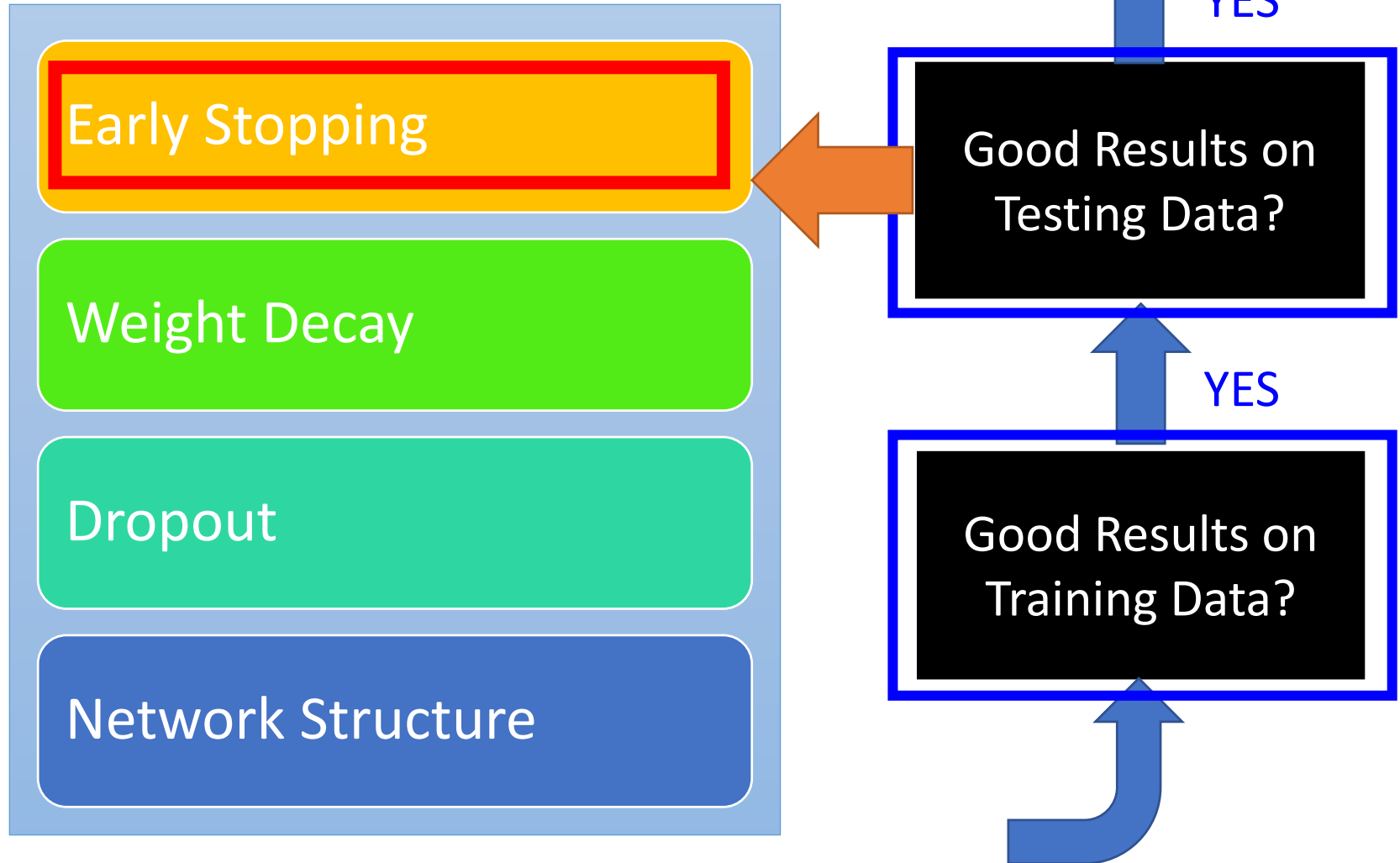
# Why Overfitting?

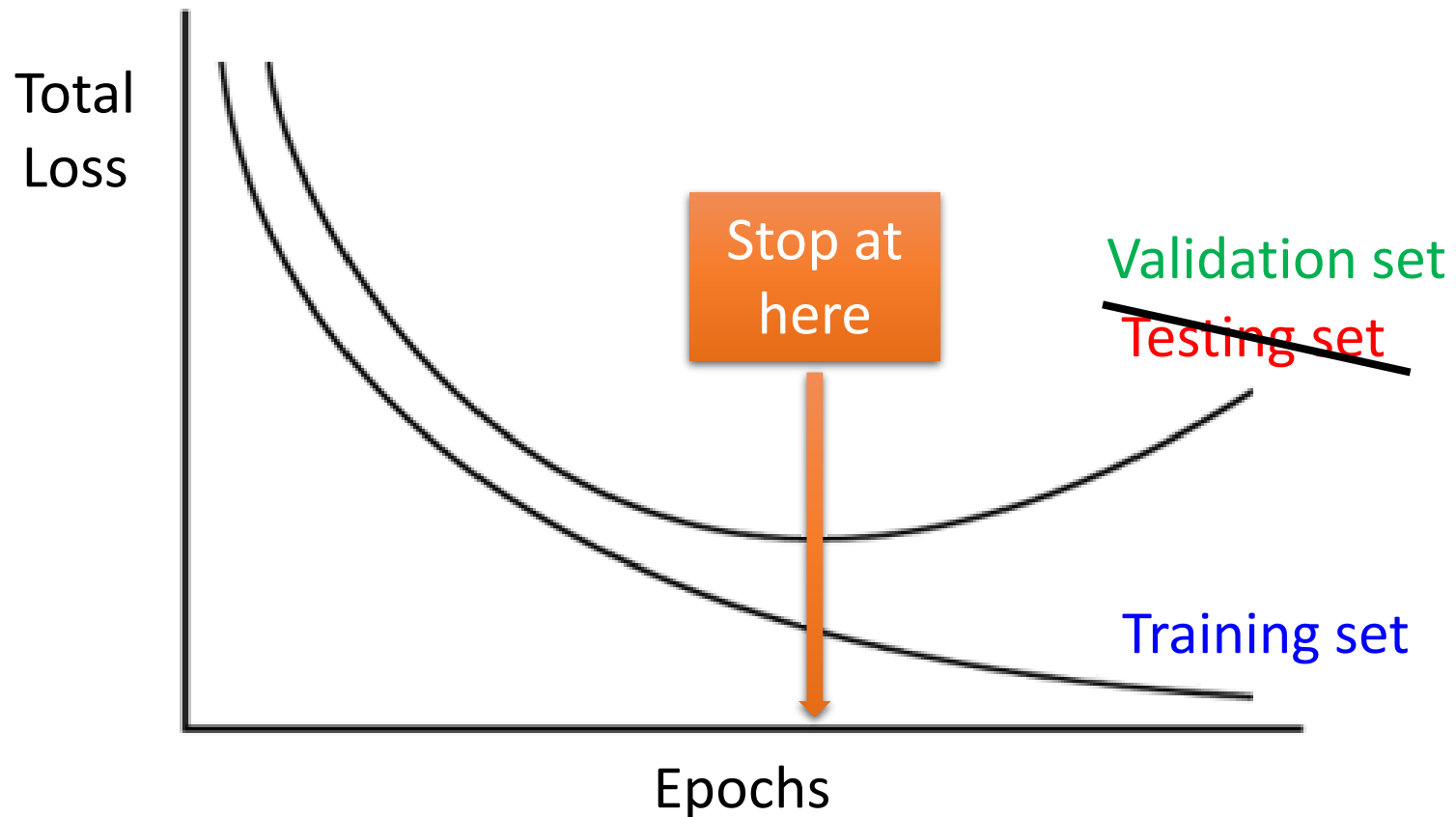- For experiments, we added some noises to the testing data

Testing:

| | Accuracy |
|---|---|
| Clean | 0.97 |
| Noisy | 0.50 |

Training is not influenced.

# Recipe of Deep Learning

**Early Stopping**

**Weight Decay**

**Dropout**

**Network Structure**

Good Results on Testing Data?

Good Results on Training Data?

YES

YES

# Early Stopping



Total Loss

Stop at here

Validation set

Testing set

Training set

Epochs

Keras: http://keras.io/getting-started/faq/#how-can-i-interrupt-training-when-the-validation-loss-isnt-decreasing-anymore

# Recipe of Deep Learning



Early Stopping

Weight Decay

Dropout

Network Structure

Good Results on Testing Data?

YES

Good Results on Training Data?

YES

# Weight Decay

- Our brain prunes out the useless link between neurons.

**Synaptic Density**

At birth    6 years old    14 years old

Source: Rethinking the Brain, Families and Work Institute, Rima Shore, 1997; Founders Network slide

Doing the same thing to machine's brain improves the performance.

# Weight Decay

Layer 1    Layer 2



$x_1$

$x_2$

$x_{256}$

Weight decay is one kind of regularization

Useless

Close to zero (萎縮了)

# Weight Decay

- Implementation

Original: $w \leftarrow w - \eta \dfrac{\partial L}{\partial w}$

$$\lambda = 0.01$$

Weight Decay: $w \leftarrow \boxed{0.99}\ w - \eta \dfrac{\partial L}{\partial w}$

Smaller and smaller

Keras: http://keras.io/regularizers/

# Recipe of Deep Learning

- Early Stopping
- Weight Decay
- **Dropout**
- Network Structure

Good Results on Testing Data?

YES

Good Results on Training Data?

YES

# Dropout

➢ **Each time before updating the parameters**
  ● Each neuron has p% to dropout

# Dropout

Thinner!

> **Each time before updating the parameters**
>
> - Each neuron has p% to dropout
>
>   ➡ **The structure of the network is changed.**
>
> - Using the new network for training

For each mini-batch, we resample the dropout neurons

# Dropout

➢ **No dropout**

● If the dropout rate at training is p%,
   all the weights times (1-p)%

● Assume that the dropout rate is 50%.
   If a weight $w = 1$ by training, set $w = 0.5$ for testing.

# Dropout - Intuitive Reason



➢ When teams up, if everyone expect the partner will do the work, nothing will be done finally.

➢ However, if you know your partner will dropout, you will do better.

➢ When testing, no one dropout actually, so obtaining good results eventually.

# Dropout - Intuitive Reason

- Why the weights should multiply (1-p)% (dropout rate) when testing?

**_Training of Dropout_**

Assume dropout rate is 50%

**_Testing of Dropout_**

No dropout



Weights from training
$$z' \approx 2z$$

Weights multiply (1-p)%
$$z' \approx z$$

# Dropout is a kind of ensemble.

**_Ensemble_**



Train a bunch of networks with different structures

# Dropout is a kind of ensemble.

***Ensemble***

# Dropout is a kind of ensemble.



**Training of Dropout**

M neurons

$2^M$ possible networks

➢Using one mini-batch to train one network

➢Some parameters in the network are shared

# Dropout is a kind of ensemble.

**_Testing of Dropout_**

testing data x



All the weights multiply (1-p)%

$y_1$        $y_2$        $y_3$

average

神秘

≈

y

# More about dropout

- More reference for dropout [Nitish Srivastava, JMLR'14] [Pierre Baldi, NIPS'13][Geoffrey E. Hinton, arXiv'12]

- Dropout works better with Maxout [Ian J. Goodfellow, ICML'13]

- Dropconnect [Li Wan, *ICML'13*]
  - Dropout delete neurons
  - Dropconnect deletes the connection between neurons

- Annealed dropout [S.J. Rennie, SLT'14]
  - Dropout rate decreases by epochs

- Standout [J. Ba, NISP'13]
  - Each neural has different dropout rate

# Let's try it

# Let's try it



No Dropout

Dropout

Step 1: Network Structure

Step 2: Learning Target

Step 3: Learn!

Good Results on Testing Data?

Good Results on Training Data?

Neural Network

YES

YES

NO

NO

Accuracy

1
0.9
0.8
0.7
0.6
0.5
0.4
0.3
0.2
0.1
0

1  2  3  4  5  6  7  8  9  10 11 12 13 14 15

Epoch

Training

Testing:

| | Accuracy |
|---|---|
| Noisy | 0.50 |
| + dropout | 0.63 |

# Concluding Remarks of Lecture II

# Recipe of Deep Learning

# Lecture III:
## Variants of Neural Networks

# Variants of Neural Networks

Convolutional Neural Network (CNN)

Consdiering the property of images

Recurrent Neural Network (RNN)

# Why CNN for Image?

- When processing image, the first layer of fully connected network would be very large



100

100

100 x 100 x 3

$3 \times 10^7$

1000

Softmax

Can the fully connected network be simplified by considering the properties of image recognition?

# Why CNN for Image

- 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



負責偵測
有沒有喙

這是喙

# Why CNN for Image

- The same patterns appear in different regions.

# Why CNN for Image

- The same patterns appear in different regions.



我負責看
圖中有沒有喙

一個 neuron 就夠了 ......

# Why CNN for Image

- Subsampling the pixels will not change the object

除非太過頭

bird

bird



subsampling

We can subsample the pixels to make image smaller

Less parameters for the network to process the image

# Convolutional Neural Network



Step 1:
Neural
Network

Step 2:
Learning
Target

Step 3:
Learn!

天生的腦

# The whole CNN

# The whole CNN



**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

Can repeat many times

Flatten

# The whole CNN



cat dog ......

# CNN – Convolution

**Those are the network parameters to be learned.**

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

Matrix

| -1 | 1 | -1 |
|----|---|----|
| -1 | 1 | -1 |
| -1 | 1 | -1 |

Filter 2

Matrix

⋮

**Property 1** Each filter detects a small pattern (3 x 3).

# CNN – Convolution

Filter 1

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

stride=1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

3    -1

# CNN – Convolution

Filter 1

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

If stride=2

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

3   -3

We set stride=1 below

# CNN – Convolution

Filter 1

stride=1

6 x 6 image

Property 2

# CNN – Convolution

Filter 2

| -1 | 1 | -1 |
|----|---|----|
| -1 | 1 | -1 |
| -1 | 1 | -1 |

stride=1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

Do the same process for every filter

| -1 | -1 | -1 | -1 |
|----|----|----|----|
| -1 | | | 1 |
| -1 | -1 | -2 | 1 |
| -1 | 0 | -4 | 3 |

Feature Map

4 x 4 image

# CNN – Zero Padding

| 1 | -1 | -1 |
|----|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

| 0 | 0 | 0 |   |   |   |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 0 | 0 | 1 | 0 |
|   | 0 | 0 | 1 | 1 | 0 | 0 |
|   | 1 | 0 | 0 | 0 | 1 | 0 |
|   | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
|   | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
|   |   |   |   | 0 | 0 | 0 |

6 x 6 image

You will get another 6 x 6 images in this way

➡ Zero padding

# CNN – Colorful image

Filter 1

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 2

| -1 | 1 | -1 |
|----|---|----|
| -1 | 1 | -1 |
| -1 | 1 | -1 |

Colorful image

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

# The whole CNN



cat dog ......

Fully Connected Feedforward network

Flatten

Convolution

Max Pooling

Convolution

Max Pooling

Can repeat many times

# CNN – Max Pooling

| 1 | -1 | -1 |
|---|---|---|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

| -1 | 1 | -1 |
|---|---|---|
| -1 | 1 | -1 |
| -1 | 1 | -1 |

Filter 2

| 3 | -1 | -3 | -1 |
|---|---|---|---|
| -3 | 1 | 0 | -3 |
| -3 | -3 | 0 | 1 |
| 3 | -2 | -2 | -1 |

| -1 | -1 | -1 | -1 |
|---|---|---|---|
| -1 | -1 | -2 | 1 |
| -1 | -1 | -2 | 1 |
| -1 | 0 | -4 | 3 |

# CNN – Max Pooling

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

Conv

Max Pooling

New image but smaller

| -1 | 1 |
|----|---|
| 0  | 3 |

2 x 2 image

Each filter is a channel

# The whole CNN



-1    1

0    3

A new image

Smaller than the original image

The number of the channel is the number of filters

Convolution

Max Pooling

Convolution

Max Pooling

Can repeat many times

# Flatten



Flatten

Fully Connected Feedforward network

# The whole CNN

(Ignoring the non-linear activation function after the convolution.)

Filter 1

6 x 6 image

Less parameters!

1: 1
2: 0
3: 0
4: 0
⋮
7: 0
8: 1
9: 0
10: 0
⋮
13: 0
14: 0
15: 1
16: 1
⋮

3

Only connect to 9 input, not fully connected

Filter 1

6 x 6 image

Less parameters!

Even less parameters!

1: 1
2: 0
3: 0
4: 0
⋮
7: 0
8: 1
9: 0
10: 0
⋮
13: 0
14: 0
15: 1
16: 1
⋮

3

-1

Shared weights

(Ignoring the non-linear activation function after the convolution.)

Input

Dim = 6 x 6 = 36

parameters =
36 x 32 = 1152

Dim = 4 x 4 x 2
= 32

image

convolution

Only 9 x 2 = 18
parameters

Max
pooling

# Convolutional Neural Network



Step 1: Neural Network → Step 2: Learning Target → Step 3: Learn!

Convolution, Max Pooling, fully connected

"monkey" ↔ 0
"cat" ↔ 1
"dog" ↔ 0

target

Learning: Nothing special, just gradient descent ......

# *CNN in Keras*

```
model.add(Convolution2D(32, 3, 3,
        border_mode='same',
        input_shape=(3,32,32)))
model.add(Activation('relu'))
```

```
model.add(Convolution2D(32, 3, 3))
model.add(Activation('relu'))
```

```
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.25))
```

```
model.add(Convolution2D(64, 3, 3,
        border_mode='same'))
model.add(Activation('relu'))
```

```
model.add(Convolution2D(64, 3, 3))
model.add(Activation('relu'))
```

```
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.25))
```

Code:
https://github.com/fchollet/keras/blob/master/examples/cifar10_cnn.py

```
model.add(Dense(10))
model.add(Activation('softmax'))
```

```
model.add(Dense(512))
model.add(Activation('relu'))
model.add(Dropout(0.5))
```
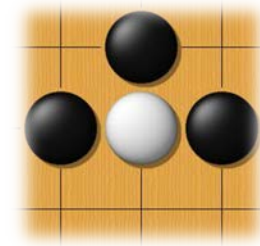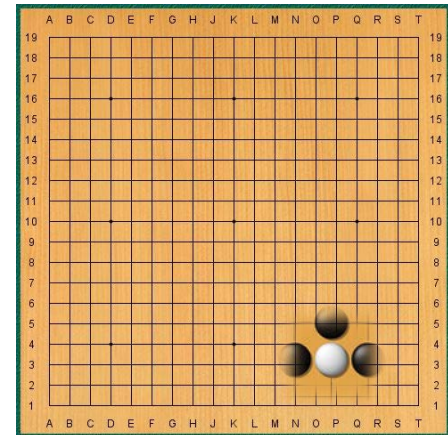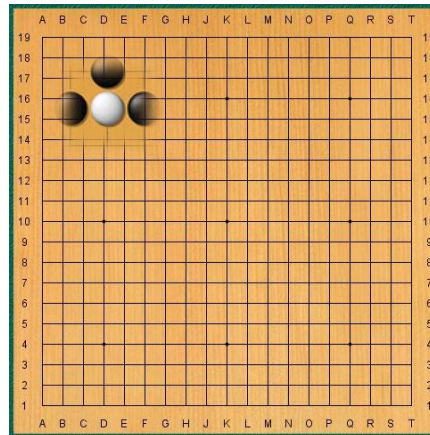
```
model.add(Flatten())
```

# Why CNN for playing Go?

- 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?

- 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 $23 \times 23$ image, then convolves $k$ filters of kernel size $5 \times 5$ 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 $21 \times 21$ image, then convolves $k$ filters of kernel size $3 \times 3$ with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size $1 \times 1$ 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.
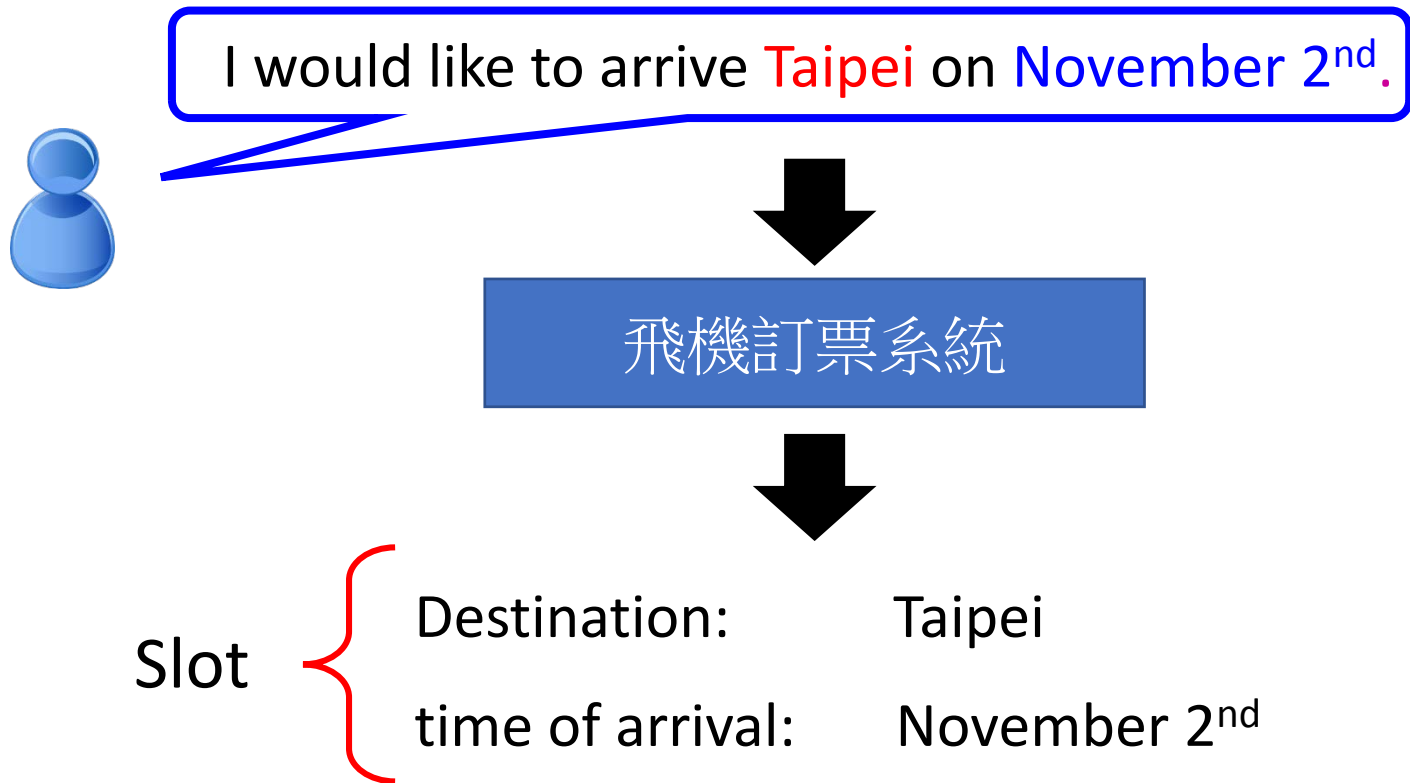
# Variants of Neural Networks

Convolutional Neural Network (CNN)

Recurrent Neural Network (RNN)

Neural Network with Memory

# Example Application

- Slot Filling

I would like to arrive Taipei on November 2nd.

飛機訂票系統

Slot

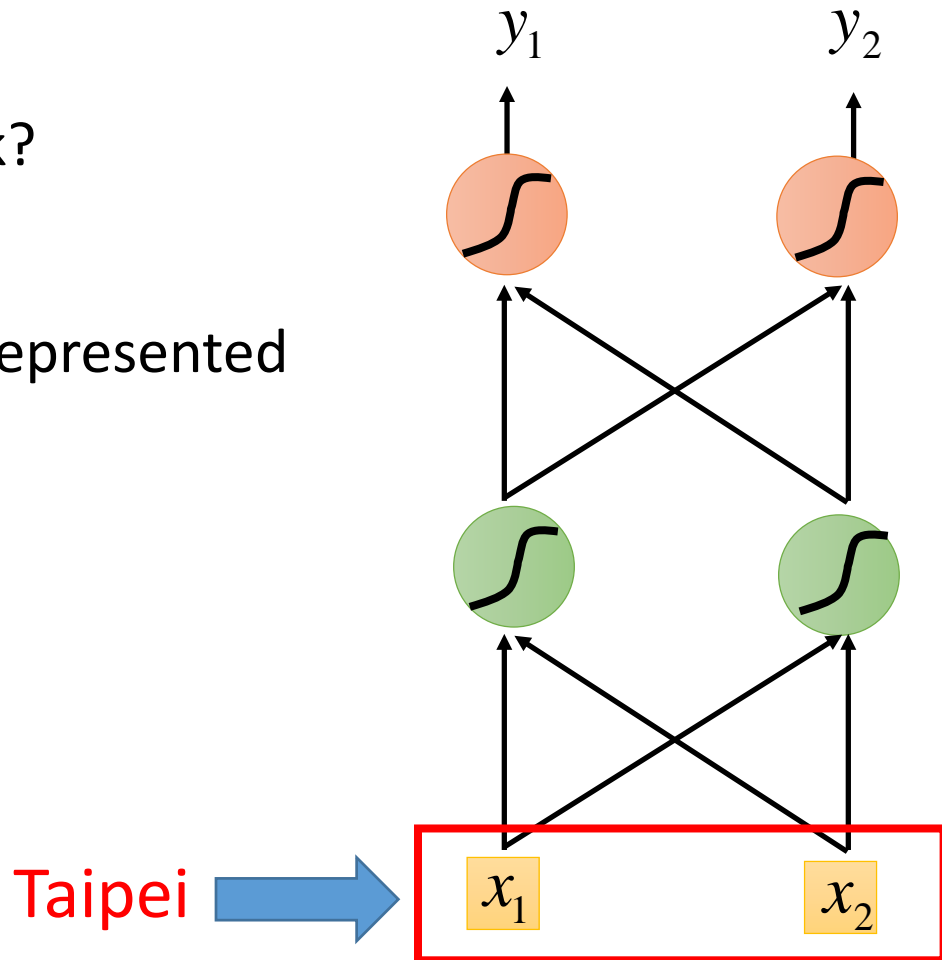Destination:     Taipei

time of arrival:     November 2nd

# Example Application

Solving slot filling by Feedforward network?

Input: a word

(Each word is represented as a vector)

# 1-of-N encoding

How to represent each word as a vector?

**_1-of-N Encoding_**    lexicon = {apple, bag, cat, dog, elephant}

The vector is lexicon size.

Each dimension corresponds to a word in the lexicon

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

apple = [ 1  0  0  0  0]

bag   = [ 0  1  0  0  0]

cat   = [ 0  0  1  0  0]

dog   = [ 0  0  0  1  0]

elephant  = [ 0  0  0  0  1]

# Beyond 1-of-N encoding

**_Dimension for "Other"_**

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

w = "Gandalf"    w = "Sauron"

**_Word hashing_**

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

**26 X 26 X 26**

w = "apple"

40

# Example Application



Solving slot filling by Feedforward network?

Input: a word

(Each word is represented as a vector)

Output:

Probability distribution that the input word belonging to the slots

# Example Application

# Recurrent Neural Network

Step 1:
Neural
Network

Step 2:
Learning
Target

Step 3:
Learn!

Recurrent Neural
Network (RNN)

天生的腦

http://onepiece1234567890.blogspot.tw/2013/12/blog-post_8.html

# Recurrent Neural Network (RNN)

The output of hidden layer are stored in the memory.

store

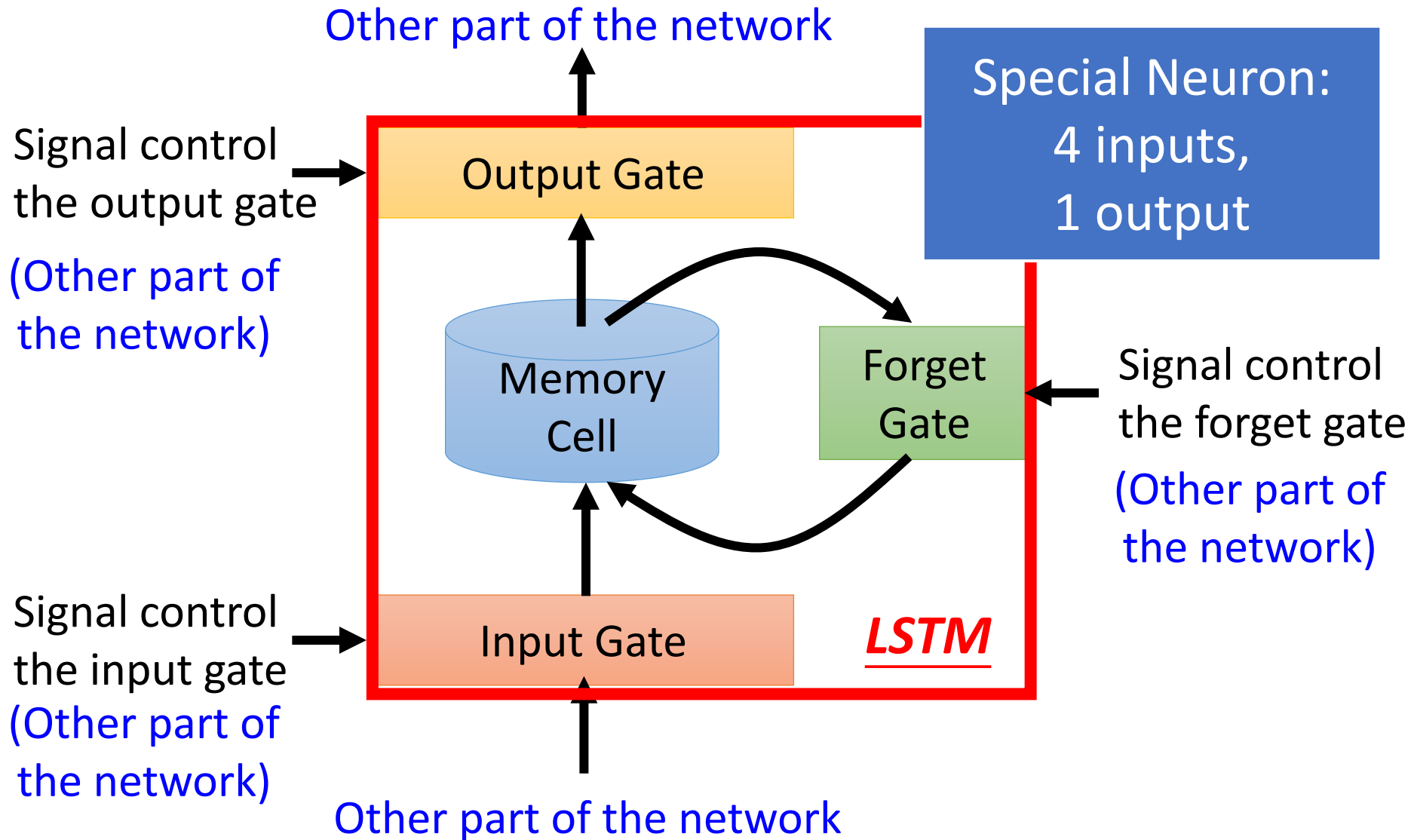Memory can be considered as another input.

$y_1$ $y_2$

$a_1$ $a_2$

$x_1$ $x_2$

# Of course it can be deep …

# Bidirectional RNN

# Long Short-term Memory (LSTM)

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

$z_o$ → **f** Output Gate

$f(z_o)$

multiply

$h(c')$

**h**

Forget Gate

$c$  $f(z_f)$

$c'$  **f** ← $z_f$

$cf(z_f)$

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

$z_i$ → **f** Input Gate

multiply

$g(z)$

**g**

Block

$z$

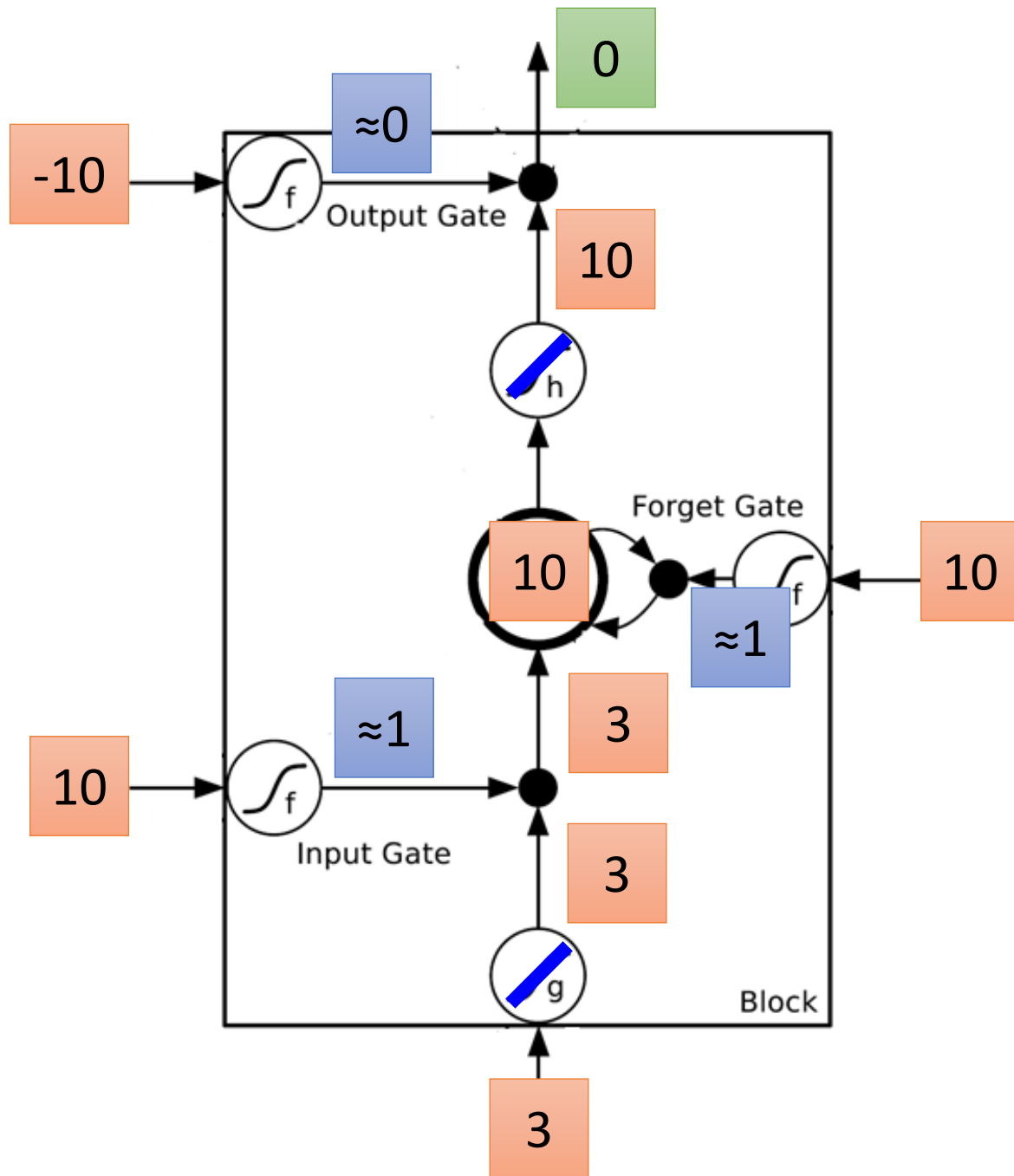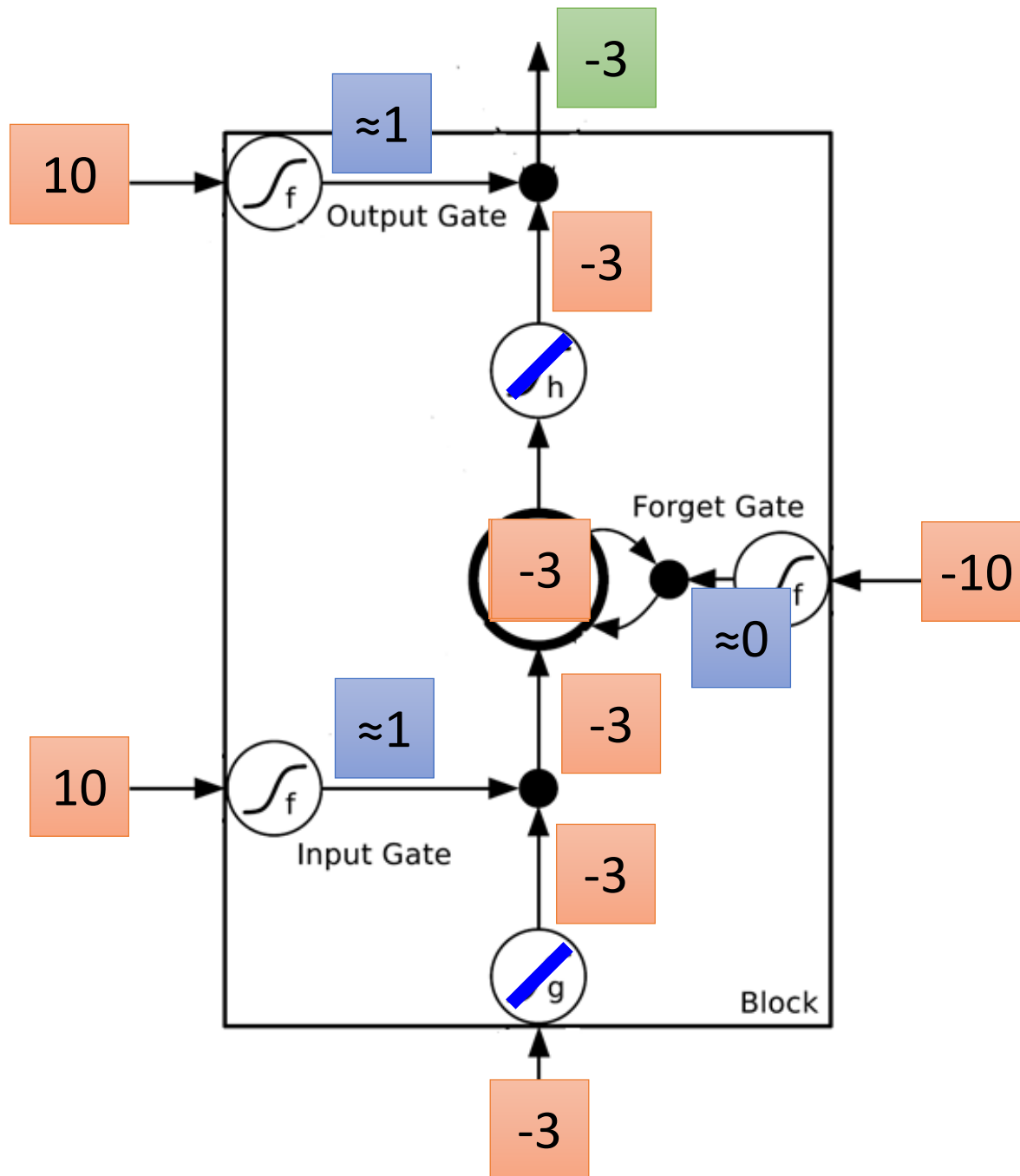Activation function f is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

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

# LSTM



$c^{t-1}$
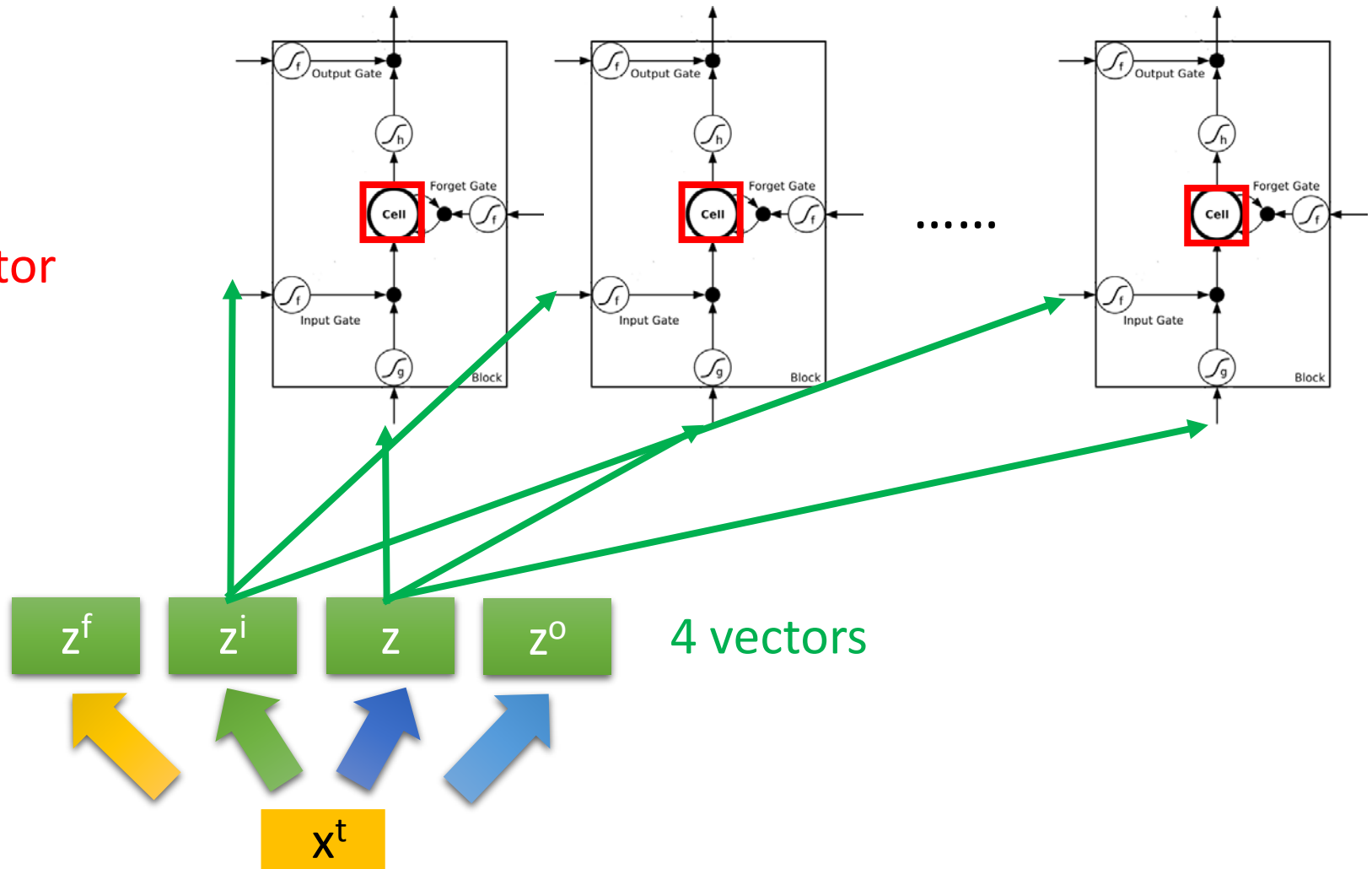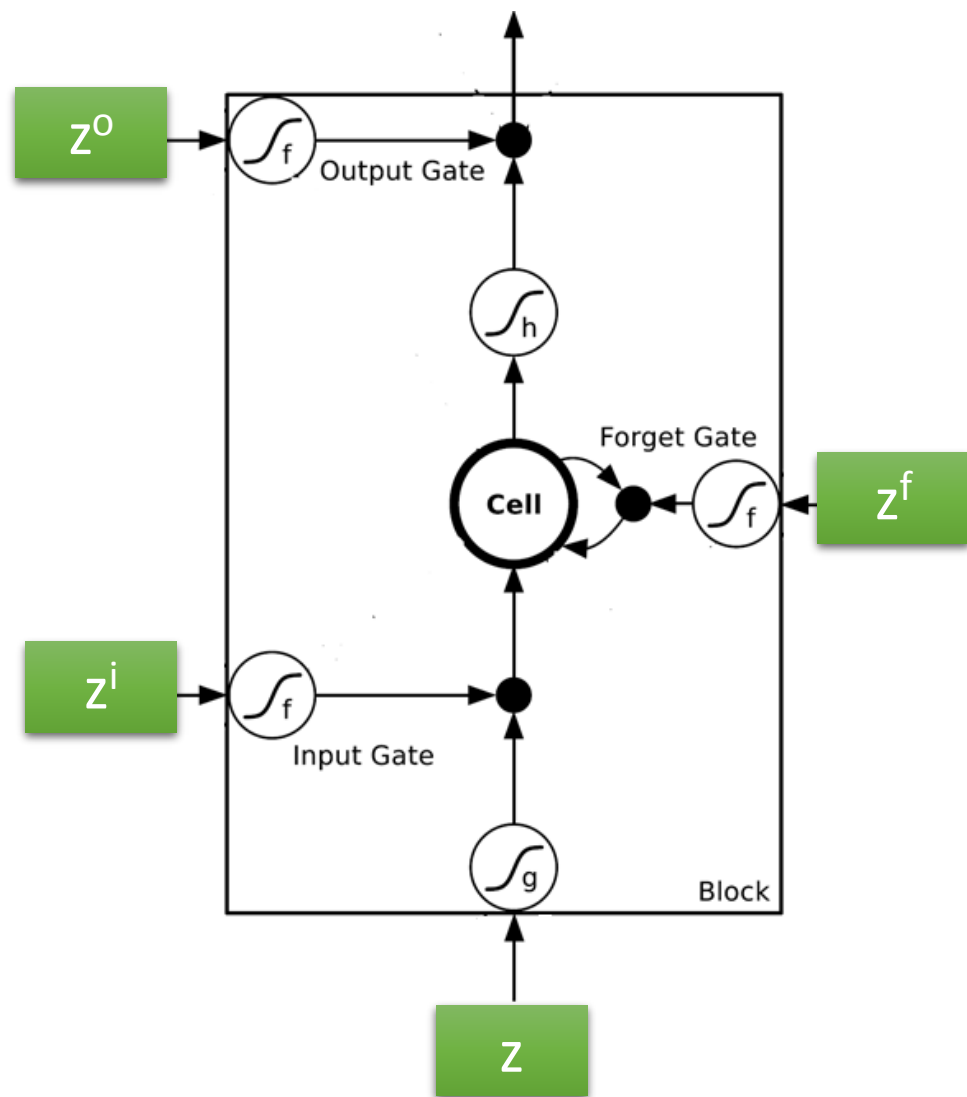
vector

$z^f$  $z^i$  $z$  $z^o$    4 vectors

$x^t$
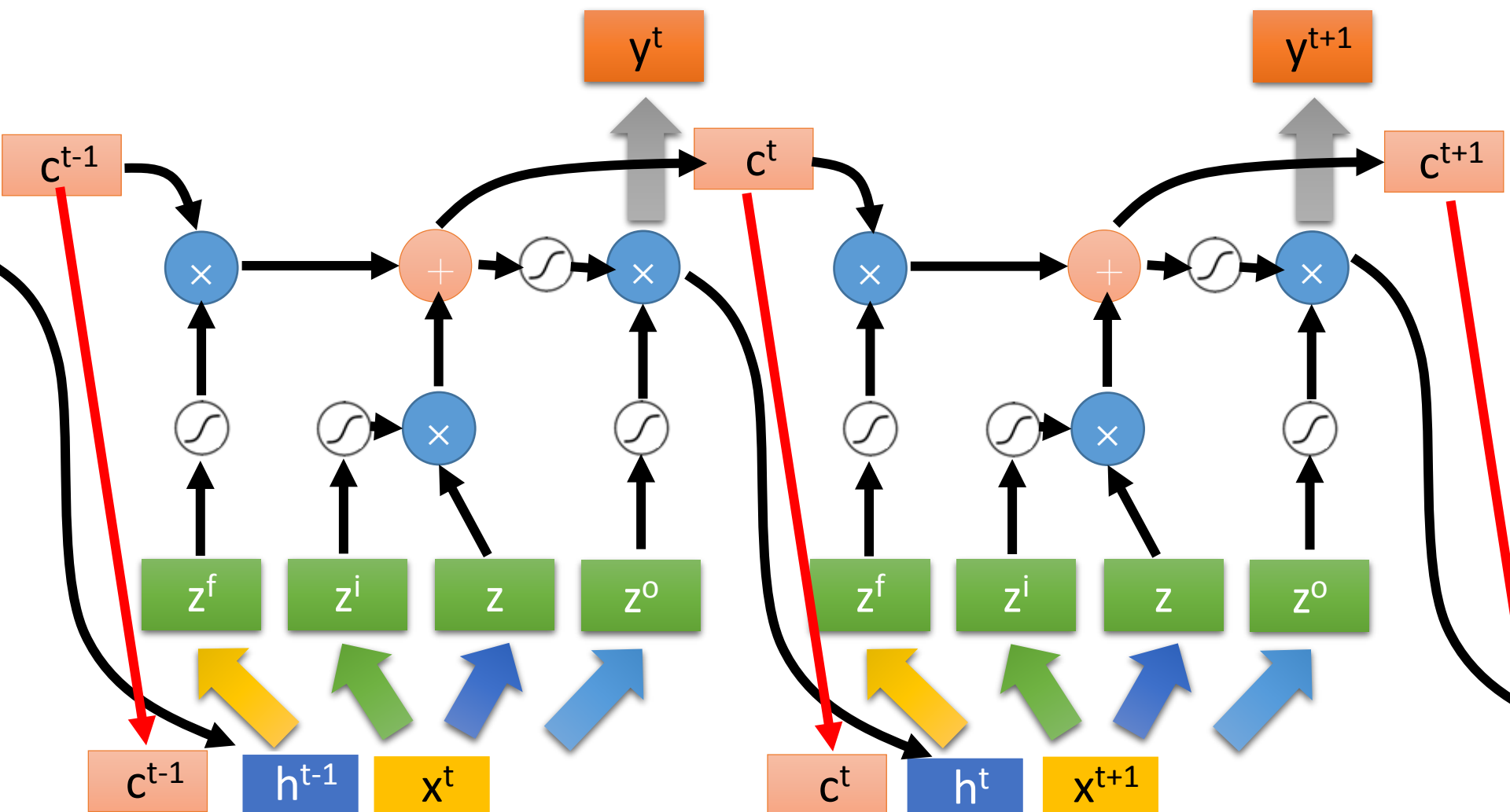
# LSTM

# LSTM

*Multiple-layer LSTM*

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

Keras supports "LSTM", "GRU", "SimpleRNN" layers

This is quite standard now.

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

# Recurrent Neural Network

Step 1:
Neural
Network

Step 2:
Learning
Target

Step 3:
Learn!

天生的腦

# *Learning Target*

# Recurrent Neural Network



Step 1: Neural Network → Step 2: Learning Target → Step 3: Learn!

天生的腦

# Learning



Backpropagation through time (BPTT)

copy

$a_1$

$a_2$

$w$

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

$y_1$

$y_2$

$x_1$

$x_2$

RNN Learning is very difficult in practice.

# Unfortunately ……

- RNN-based network is not always easy to learn

Real experiments on Language modeling



Total Loss
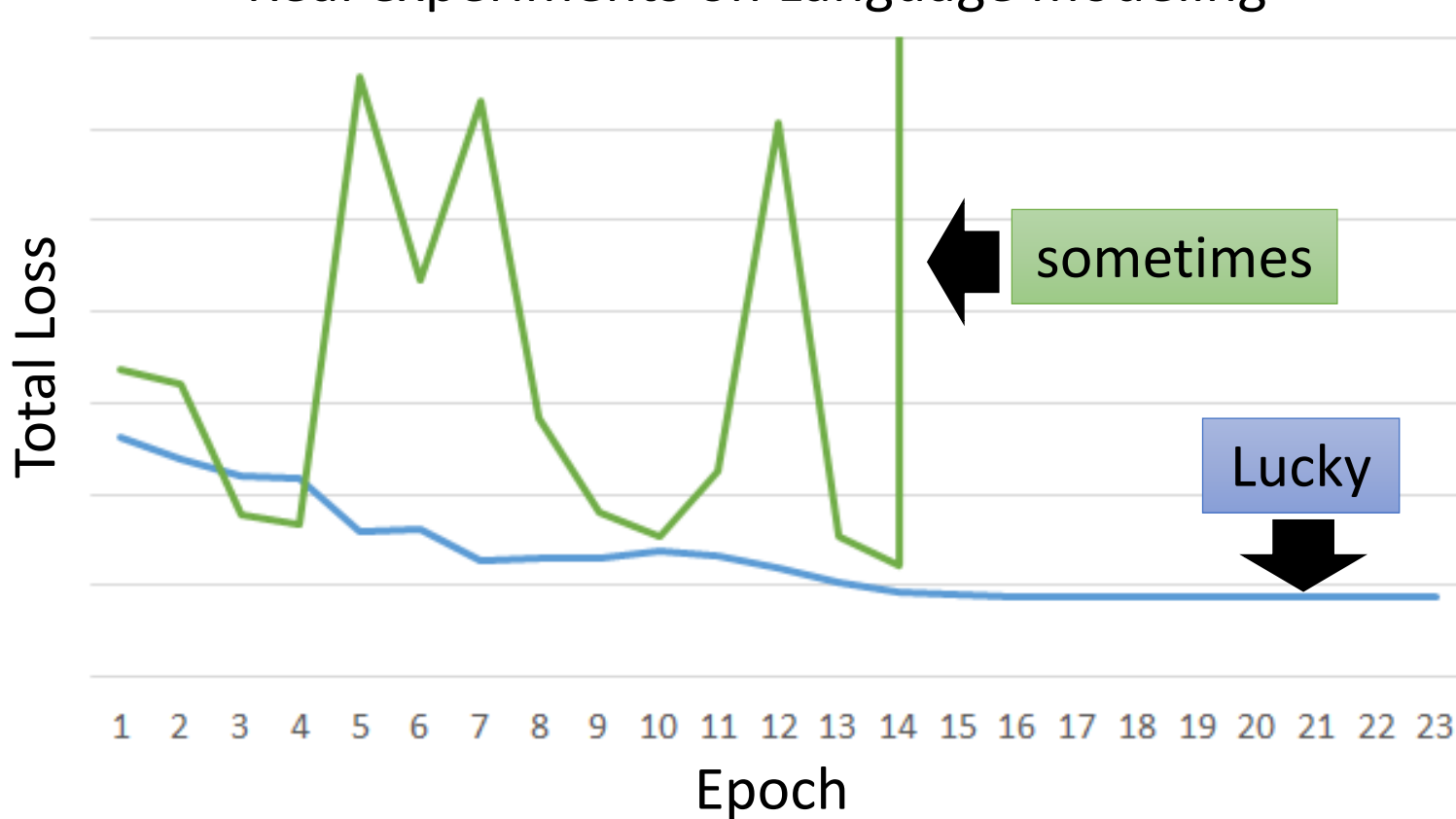
sometimes

Lucky

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23

Epoch

# The error surface is rough.



The error surface is either very flat or very steep.

Clipping

Total Loss

$w_2$

$w_1$

[Razvan Pascanu, ICML'13]

# Why?

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

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

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

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

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

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

$= w^{999}$

**_Toy Example_**

# Helpful Techniques

- Advance momentum method
  - Nesterov's Accelerated Gradient (NAG)



**Methods:**

— Gradient descent

— Momentum

— Nesterov's Accelerated Gradient (NAG)

Source:
http://www.cs.toronto.edu/~fritz/absps/momentum.pdf

# Helpful Techniques

- Long Short-term Memory (LSTM)
  - Can deal with gradient vanishing (not gradient explode)

  ➢ Memory and input are ***added***

  ➢ The influence never disappears unless forget gate is closed

  ➡ No Gradient vanishing
  (If forget gate is opened.)

# Helpful Techniques

• Gated Recurrent Unit (GRU)

Simplified LSTM

[Cho, EMNLP'14]

舊的不去、新的不來
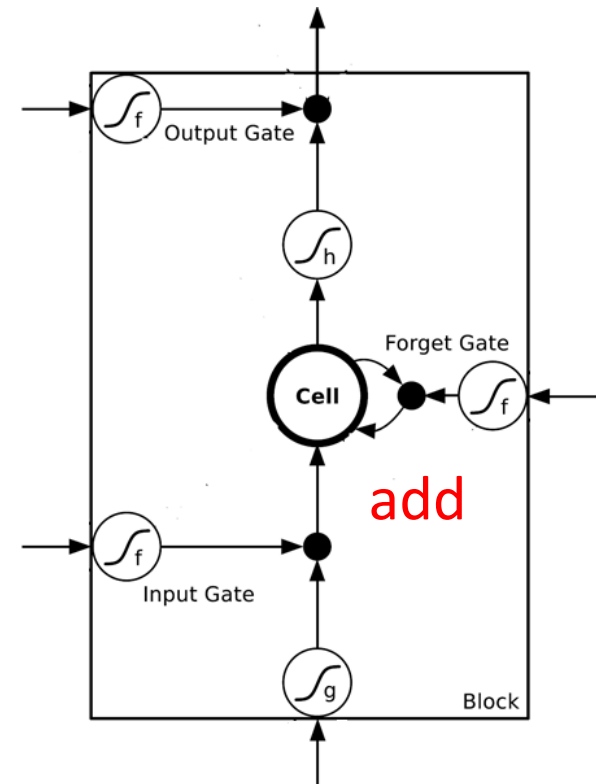
GRU has less parameters than LSTM

# Helpful Techniques

### Clockwise RNN



Output

$T_1$    $T_2$    $T_g$

Hidden

Input

[Jan Koutnik, JMLR'14]

### Structurally Constrained Recurrent Network (SCRN)



$y_t$

$U$    $V$

$R$    $h_t$    $P$    $s_t$    $\alpha$

$A$    $B$

$x_t$

[Tomas Mikolov, ICLR'15]

Vanilla RNN Initialized with Identity matrix + ReLU activation function [Quoc V. Le, arXiv'15]

➢ Outperform or be comparable with LSTM in 4 different tasks

# More Applications ......

Probability of "arrive" in each slot

Probability of "Taipei" in each slot

Probability of "on" in each slot



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

$a^3$

$a^2$

Input and output are both sequences with the same length

RNN can do more than that!

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

arrive    Taipei    on    November    2nd
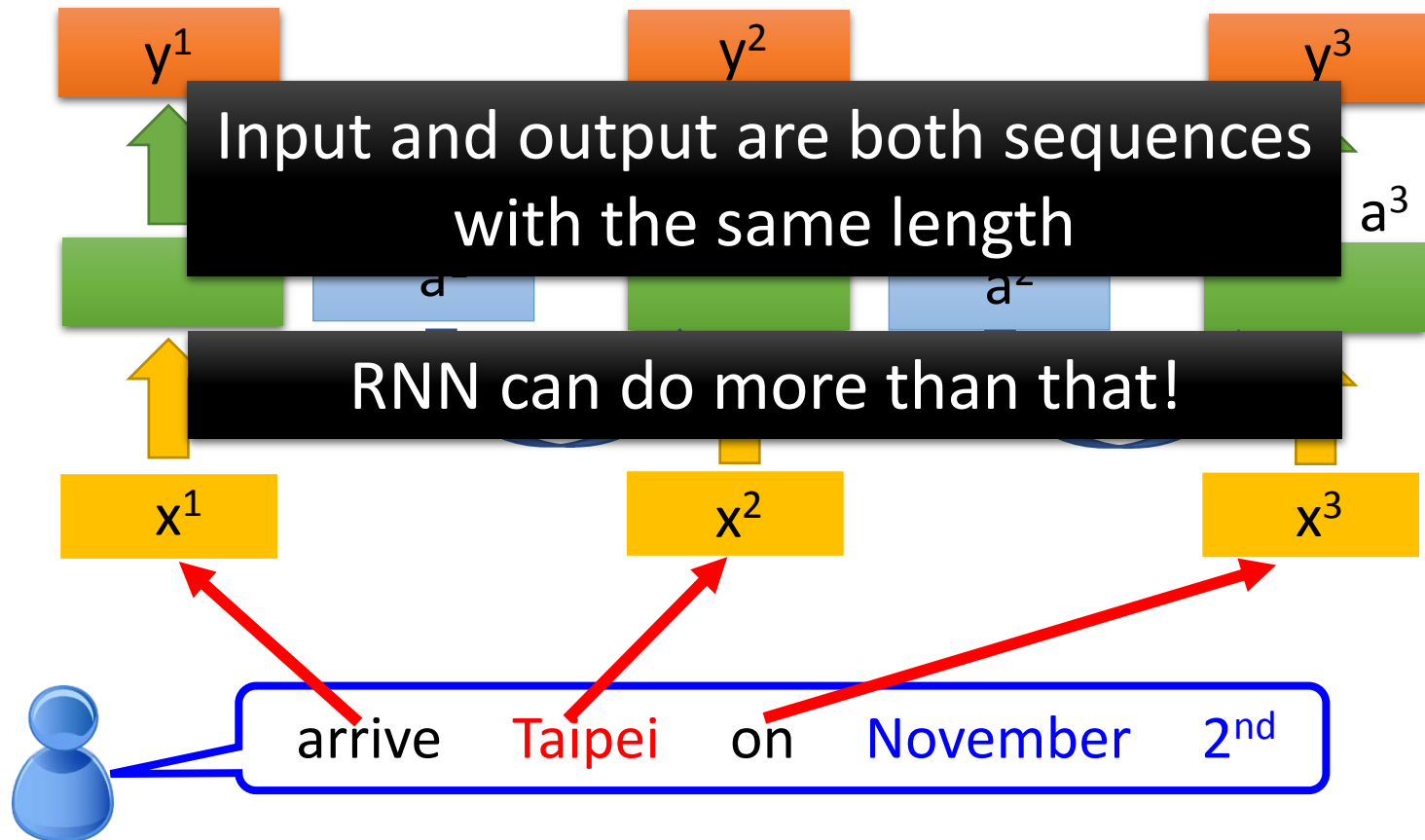
# Many to one

Keras Example:
https://github.com/fchollet/keras/blob/master/examples/imdb_lstm.py

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

### *Sentiment Analysis*

看了這部電影覺得很高興 …….

這部電影太糟了 …….

這部電影很棒 …….

Positive (正雷)　Negative (負雷)　Positive (正雷)

超好雷
好雷
普雷
負雷
超負雷

我　覺　得　……　太　糟　了

# Many to Many (Output is shorter)

- Both input and output are both sequences, **_but the output is shorter._**
  - E.g. **_Speech Recognition_**

Output: "好棒" (character sequence)
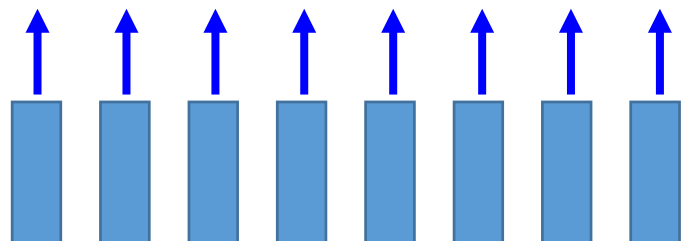
Trimming

Problem?
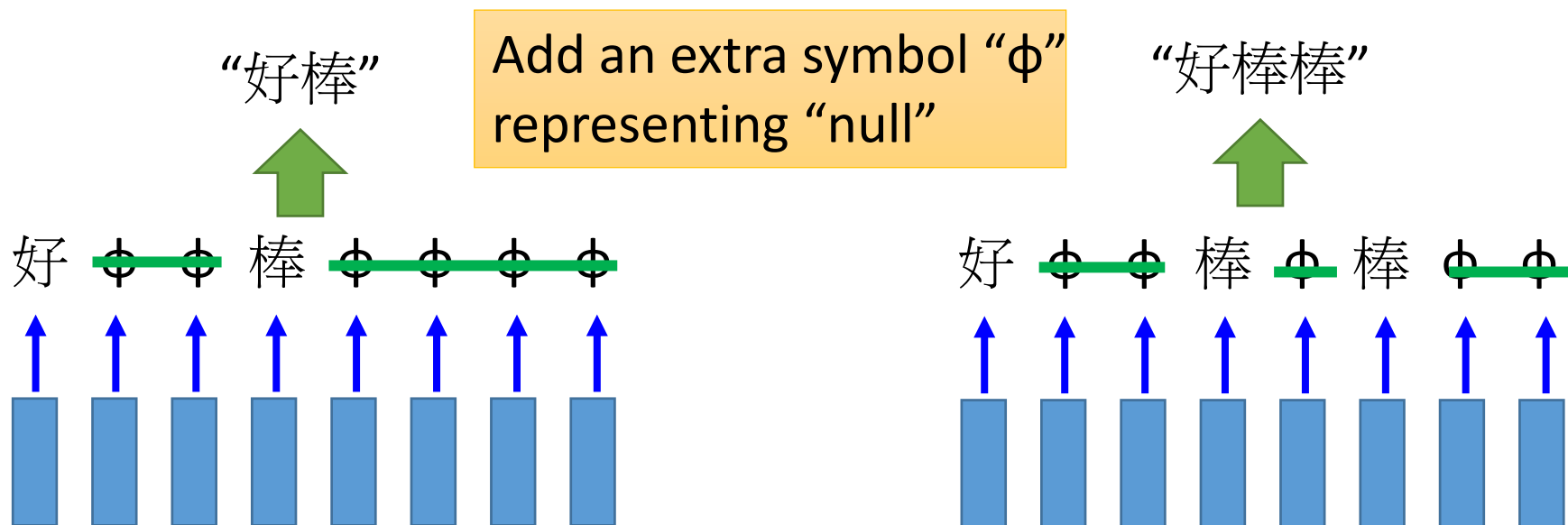
Why can't it be "好棒棒"

好 好 好 棒 棒 棒 棒 棒

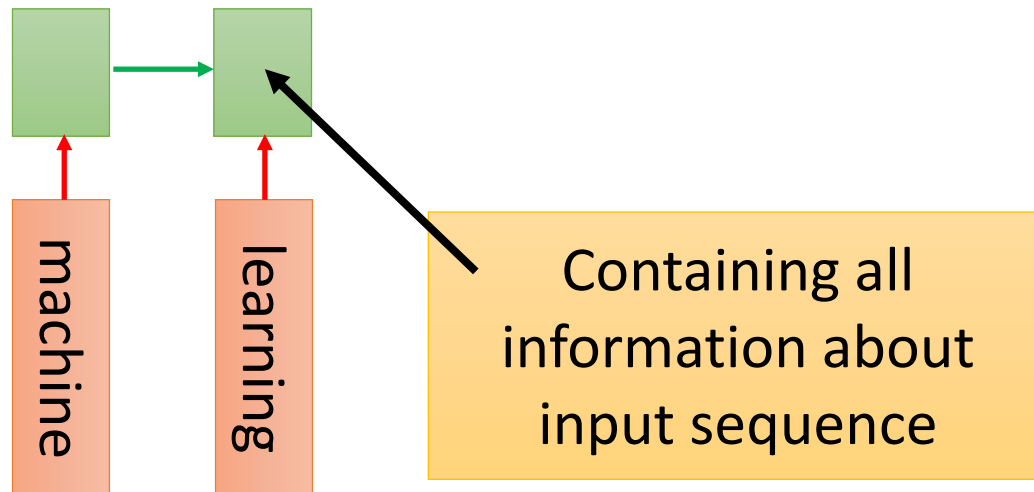Input: (vector sequence)

# Many to Many (Output is shorter)

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

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

"好棒"

Add an extra symbol "φ" representing "null"

"好棒棒"

好 φ φ 棒 φ φ φ φ

好 φ φ 棒 φ 棒 φ φ

# Many to Many (No Limitation)

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



machine  learning

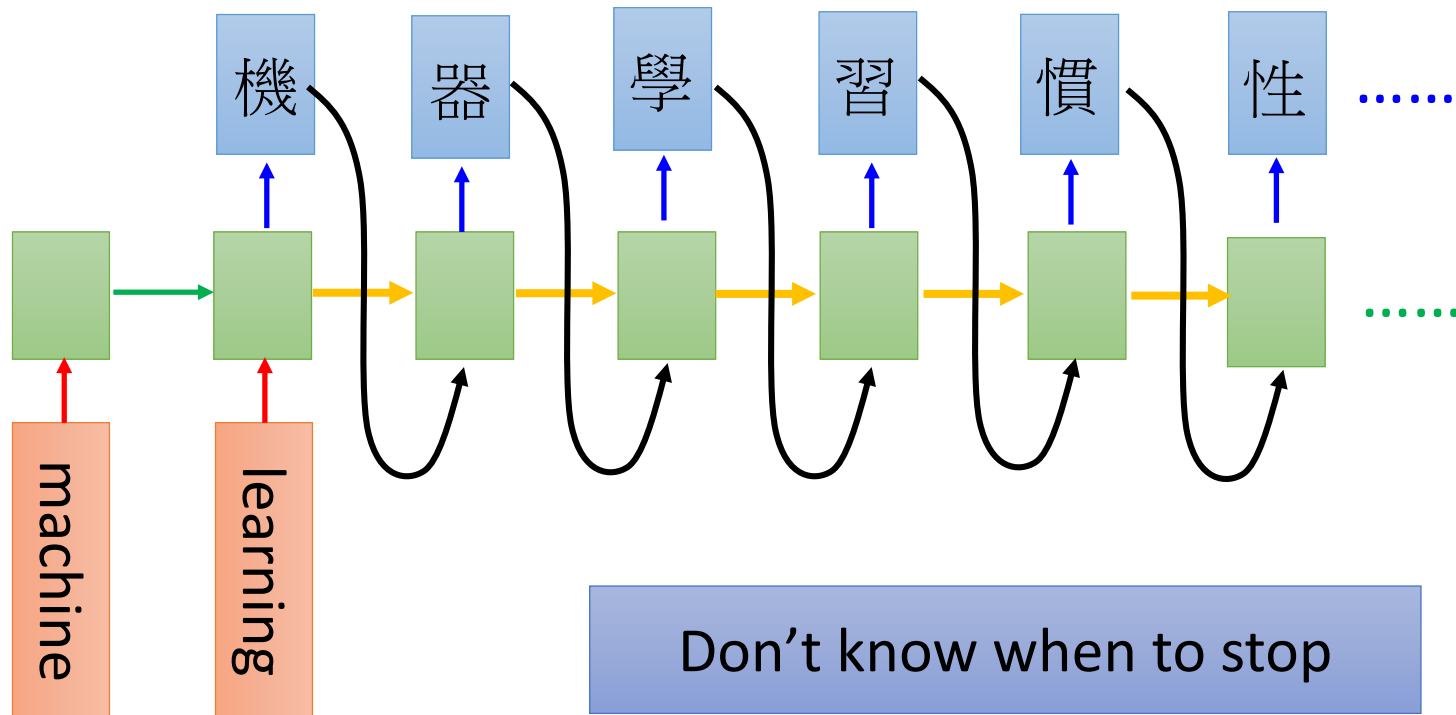Containing all information about input sequence

# Many to Many (No Limitation)

- Both input and output are both sequences **_with different lengths_**. → **_Sequence to sequence learning_**
  - E.g. **_Machine Translation_** (machine learning→機器學習)
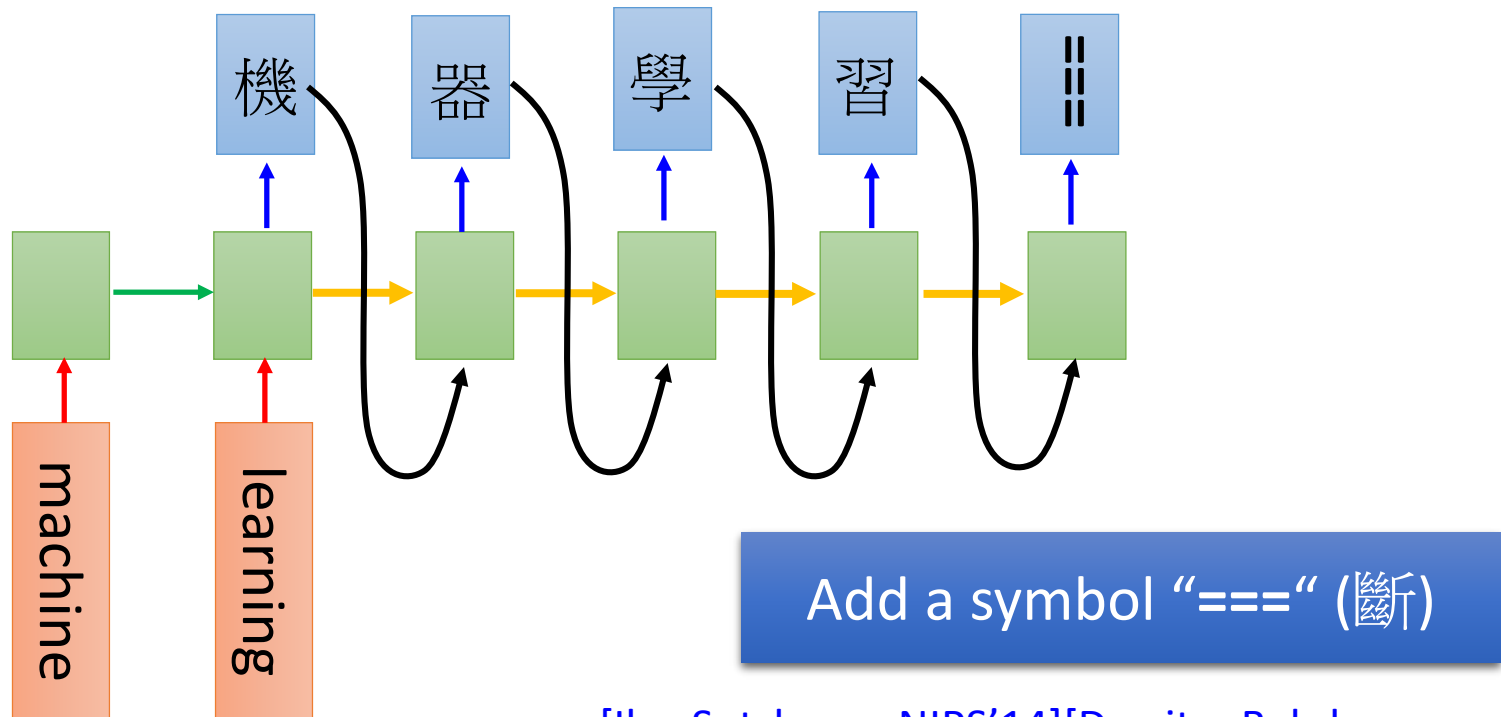
# Many to Many (No Limitation)

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

推 tlkagk:　　　=========斷==========

Ref:http://zh.pttpedia.wikia.com/wiki/%E6%8E%A5%E9%BE%8D%E6%8E%A8%E6%96%87 (鄉民百科)
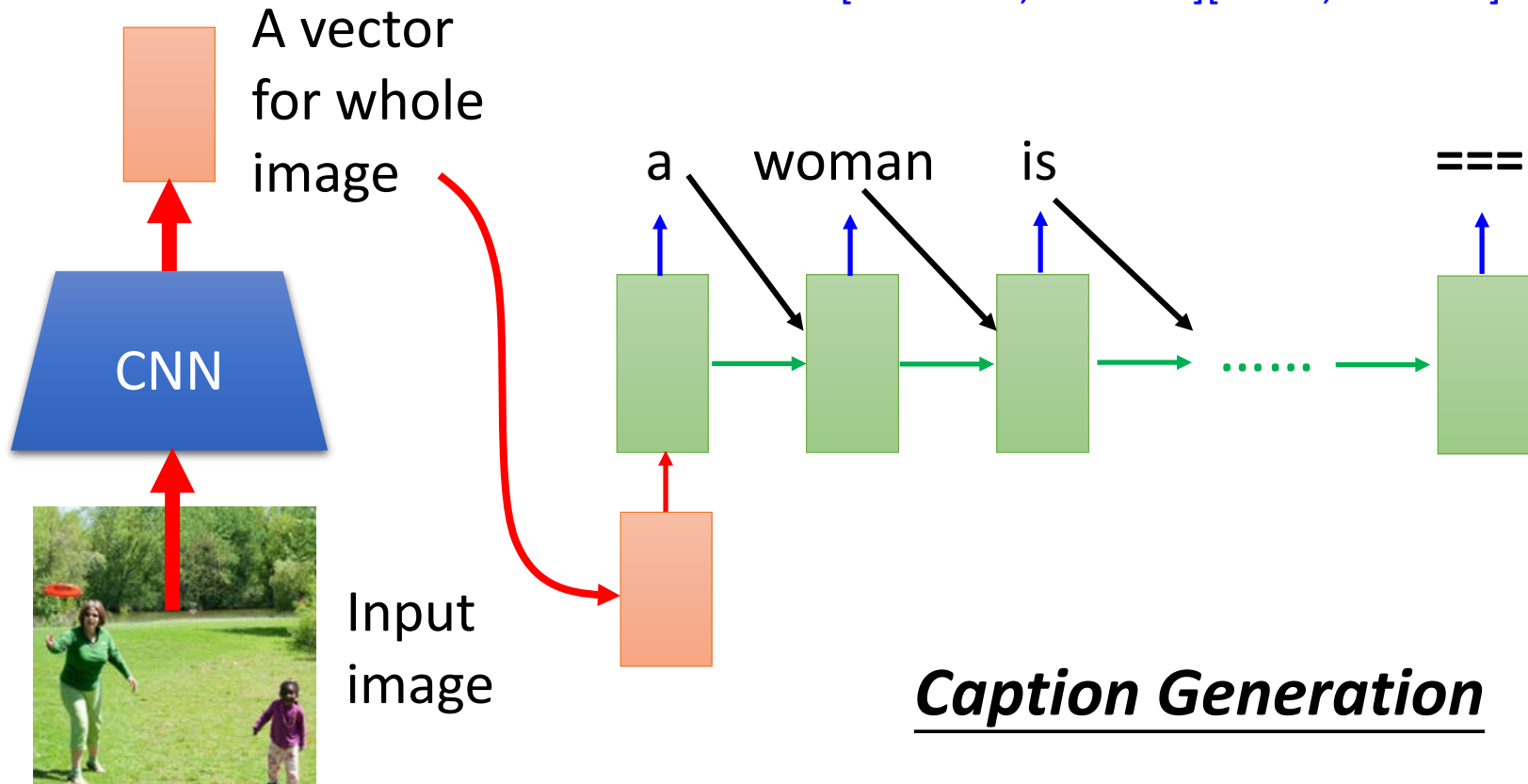
# Many to Many (No Limitation)

- Both input and output are both sequences **_with different lengths_**. → **_Sequence to sequence learning_**
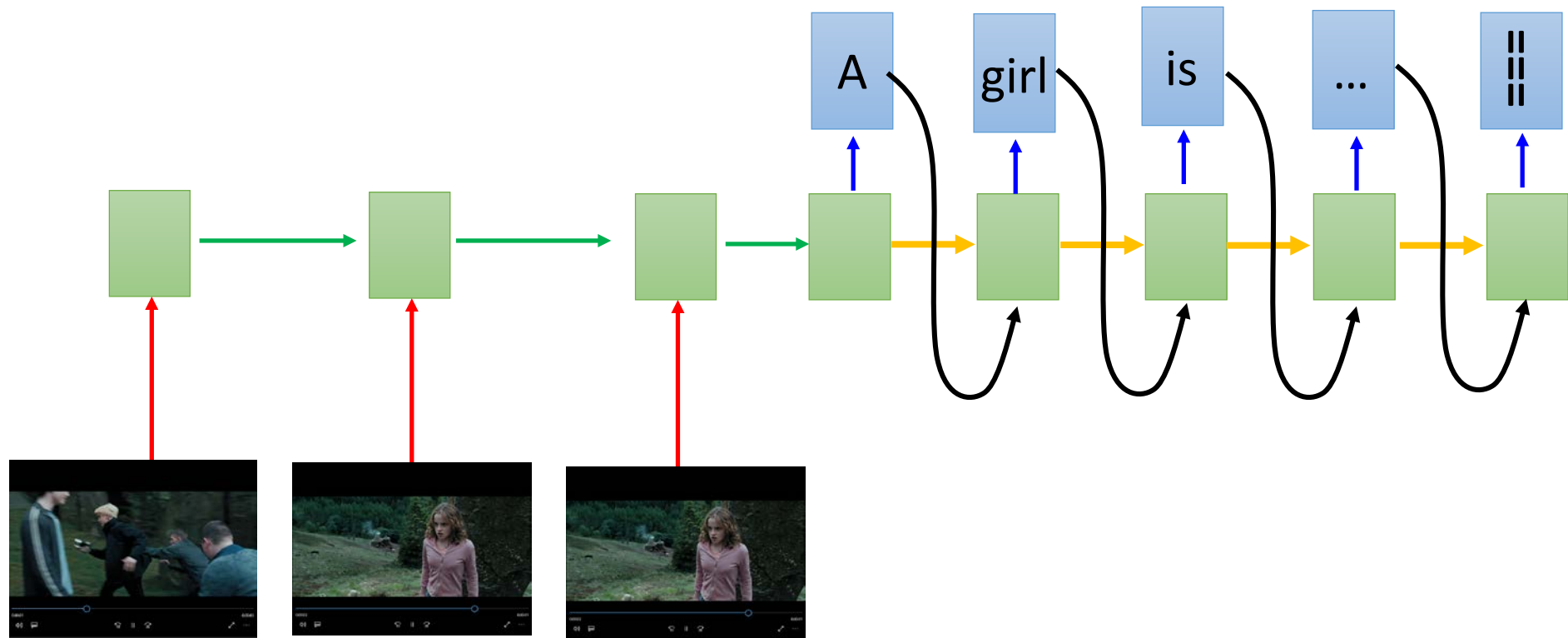    - E.g. **_Machine Translation_** (machine learning→機器學習)



Add a symbol "===" (斷)

[Ilya Sutskever, NIPS'14][Dzmitry Bahdanau, arXiv'15]

# One to Many

- Input an image, but output a sequence of words

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



A vector for whole image

CNN

Input image

a   woman   is   ===

***Caption Generation***

# Video Caption Generation



Video frames

# Concluding Remarks

Convolutional Neural Network (CNN)

Recurrent Neural Network (RNN)

# Lecture IV:
## Next Wave

# Outline

Ultra Deep Network

Attention Model

Reinforcement Learning

Realizing what the World Looks Like

Understanding the Meaning of Words

Why Deep?

# Ultra Deep Network

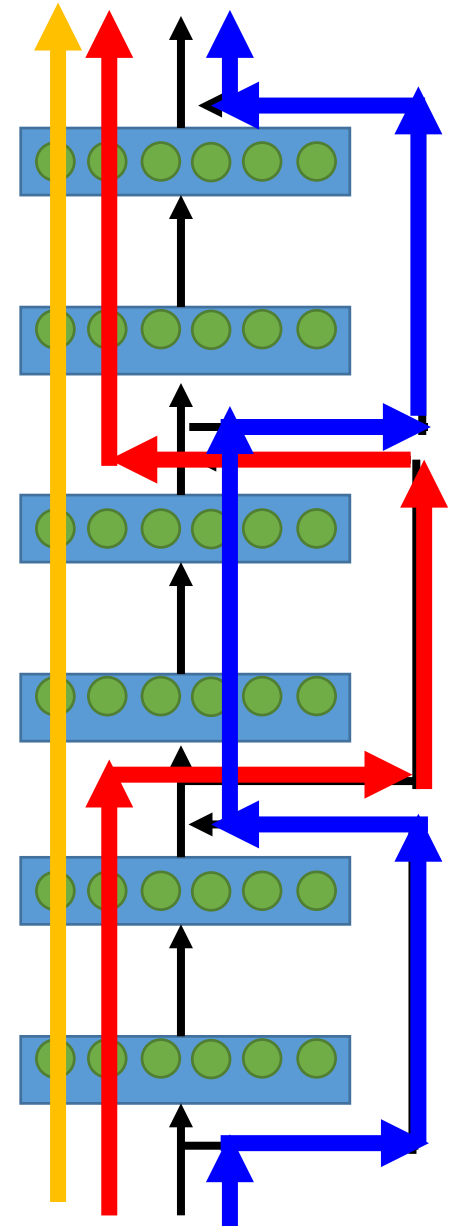- Ultra deep network is the ensemble of many networks with different depth.

**Ensemble** {
**6 layers**
**4 layers**
**2 layers**
}

Residual Networks are Exponential Ensembles of Relatively Shallow Networks
https://arxiv.org/abs/1605.06431
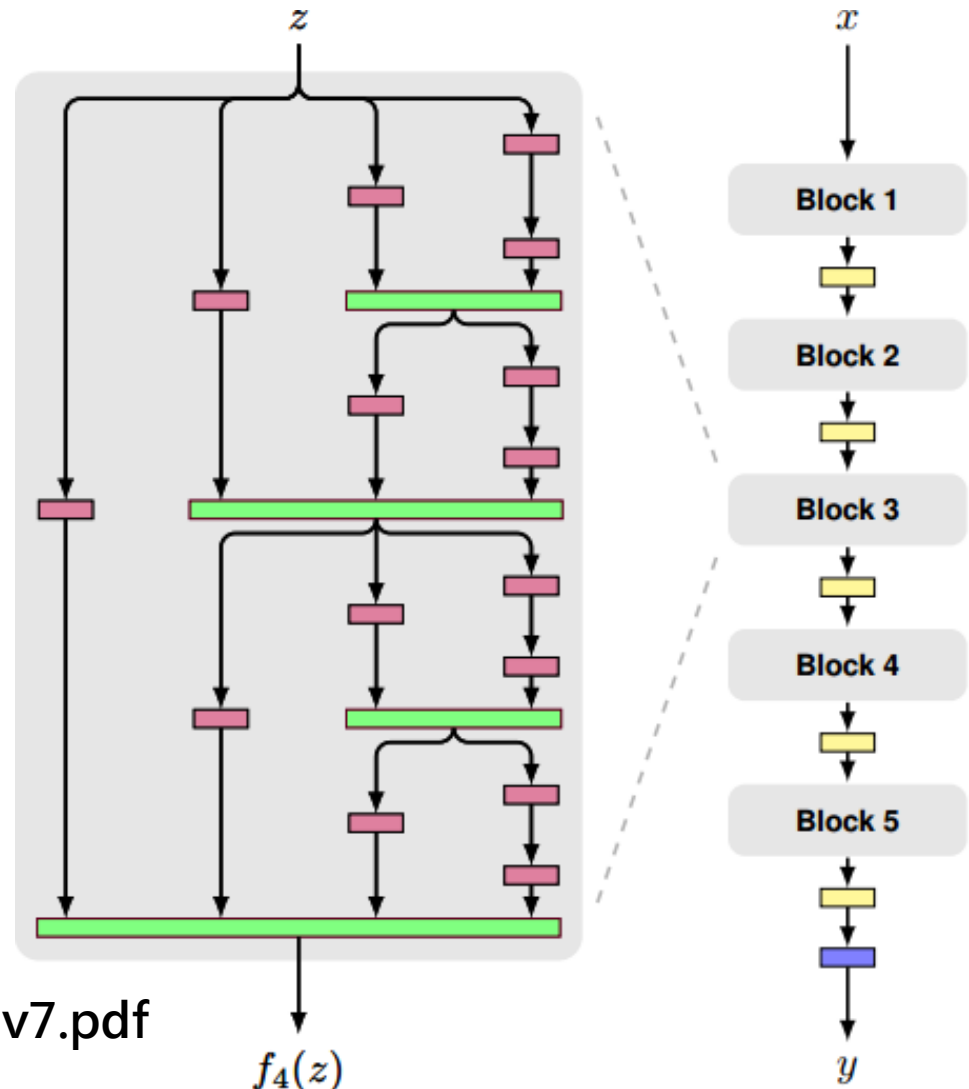
# Ultra Deep Network

- FractalNet

FractalNet: Ultra-Deep Neural Networks without Residuals
https://arxiv.org/abs/1605.07648

## Resnet in Resnet

Resnet in Resnet: Generalizing Residual Architectures
https://arxiv.org/abs/1603.08029

## Good Initialization?

All you need is a good init
http://arxiv.org/pdf/1511.06422v7.pdf

# Ultra Deep Network

- **Residual Network**

- **Highway Network**



copy

Gate
controller

copy

**Deep Residual Learning for Image
Recognition**
http://arxiv.org/abs/1512.03385

**Training Very Deep Networks**
https://arxiv.org/pdf/1507.0622
8v2.pdf

output layer

output layer

output layer

Highway Network automatically determines the layers needed!

Input layer

Input layer

Input layer

# Outline

Ultra Deep Network

Attention Model
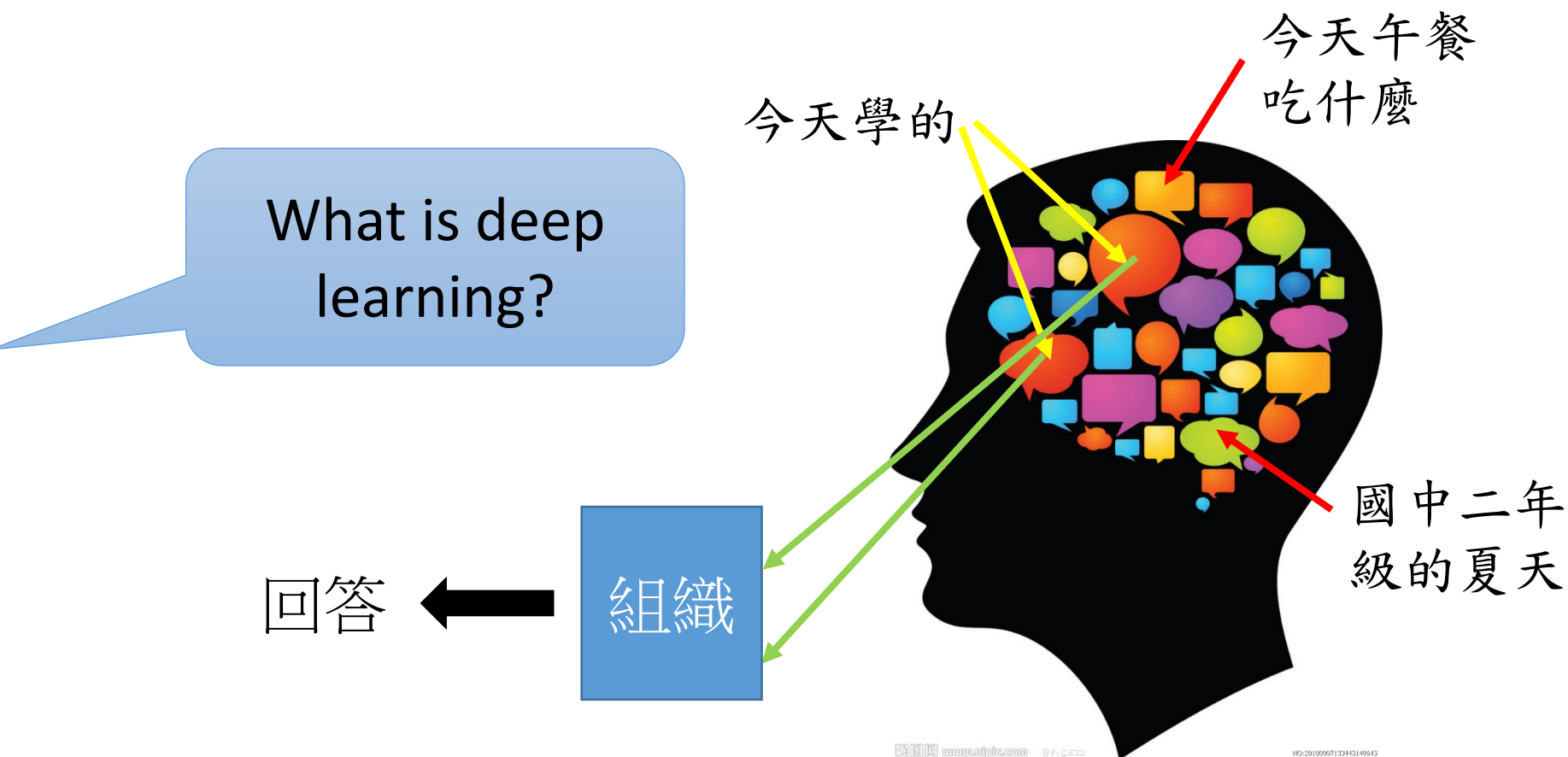
Reinforcement Learning

Realizing what the World Looks Like

Understanding the Meaning of Words
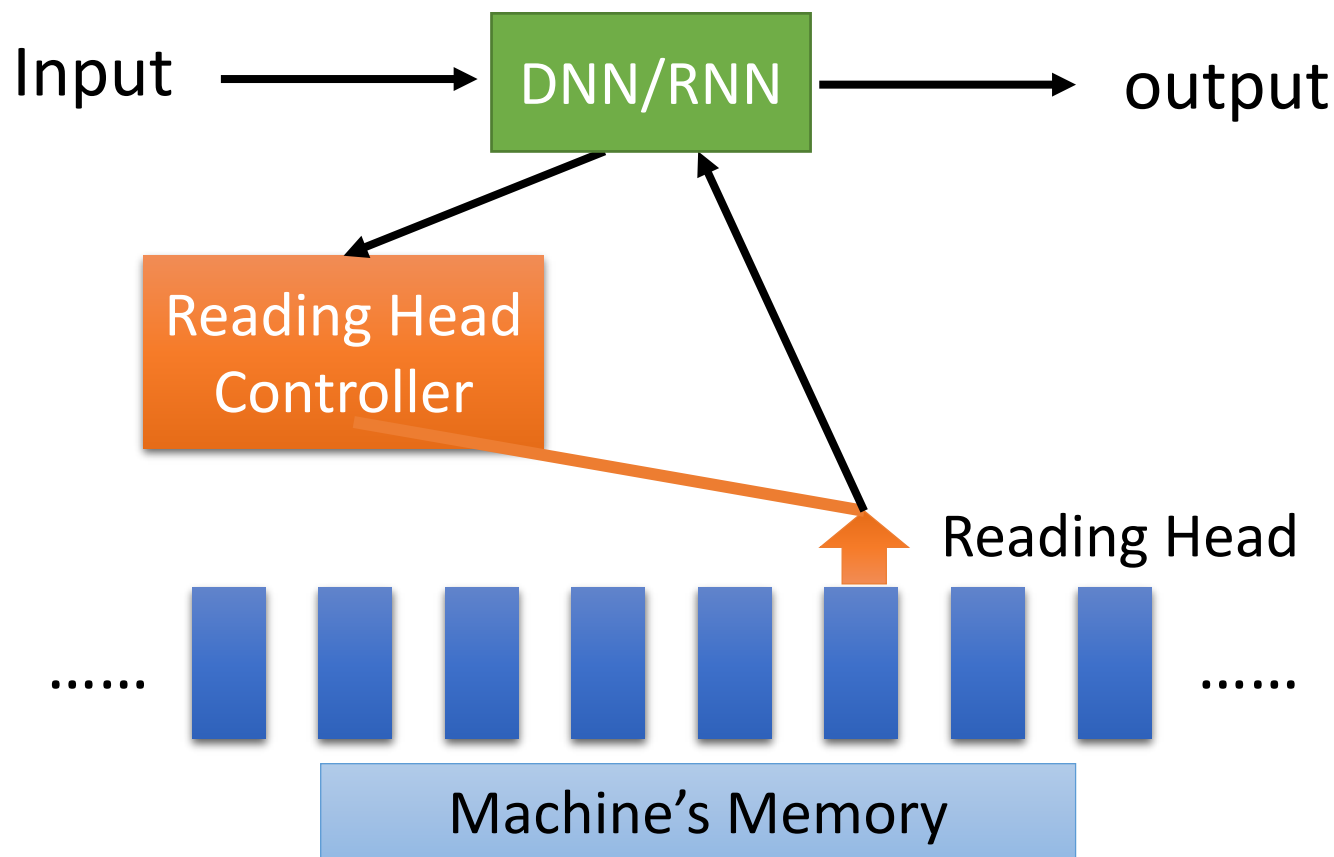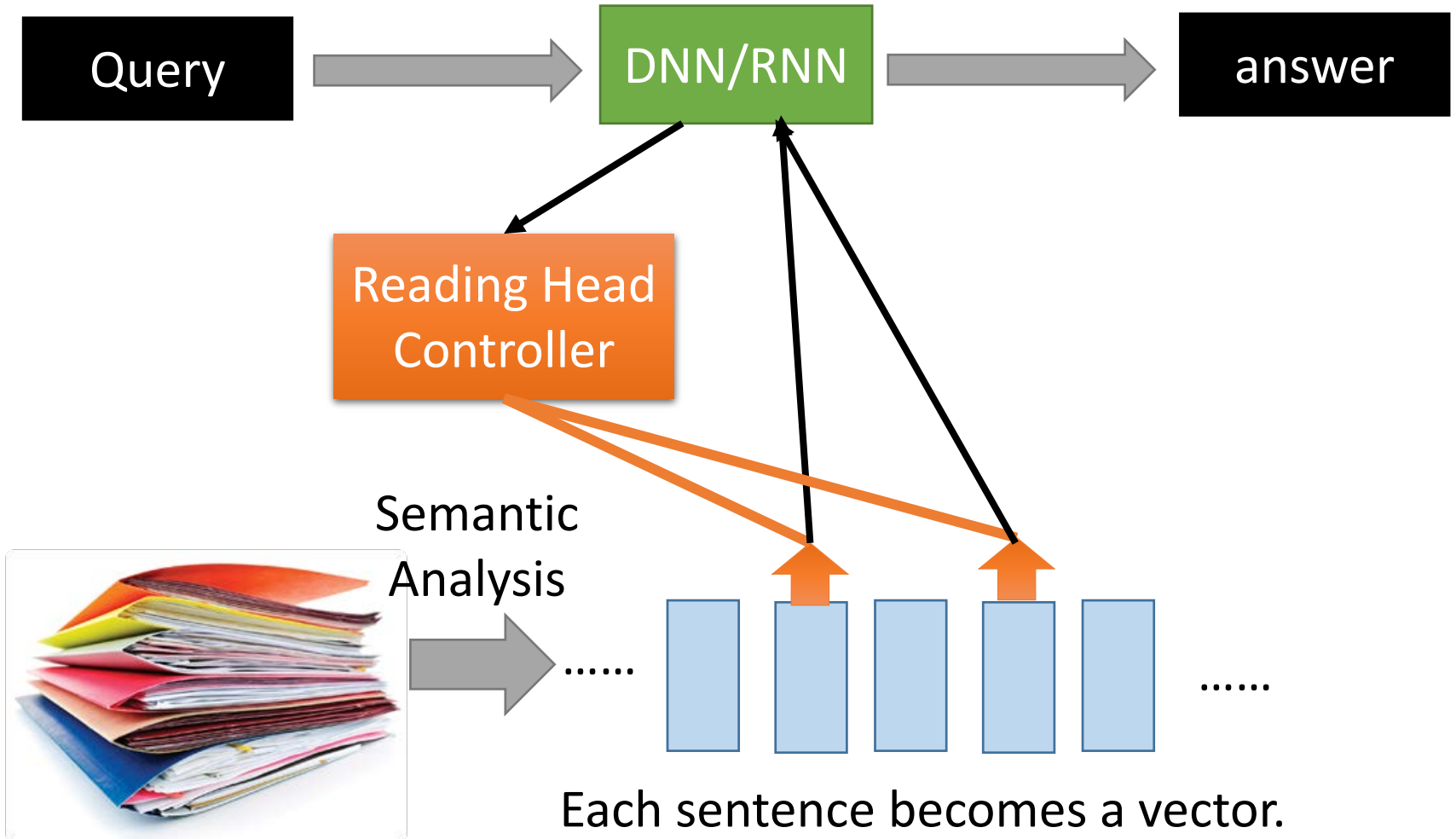
Why Deep?

# Attention-based Model



What is deep learning?

今天學的

今天午餐吃什麼

國中二年級的夏天

組織

回答

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

# Attention-based Model



Input → DNN/RNN → output

Reading Head Controller

Reading Head

...... Machine's Memory ......

# Reading Comprehension



Query → DNN/RNN → answer

Reading Head Controller

Semantic Analysis

Each sentence becomes a vector.

# Reading Comprehension

- End-To-End Memory Networks. S. Sukhbaatar, A. Szlam, J. Weston, R. Fergus. NIPS, 2015.

The position of reading head:

| 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** | | | | |

Keras has example:

https://github.com/fchollet/keras/blob/master/examples/babi_memnn.py

# Visual Question Answering
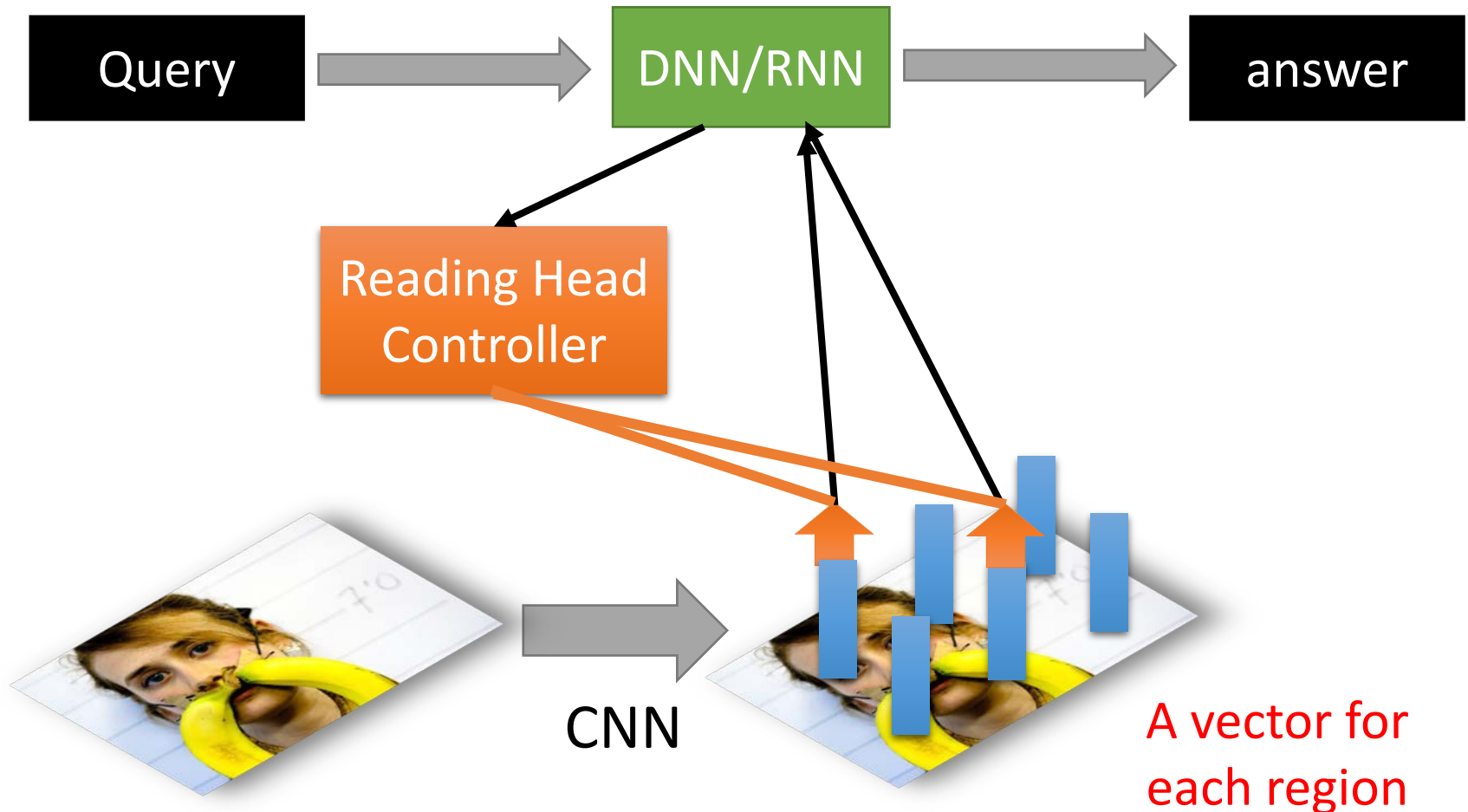


What is the mustache made of?

AI System → bananas
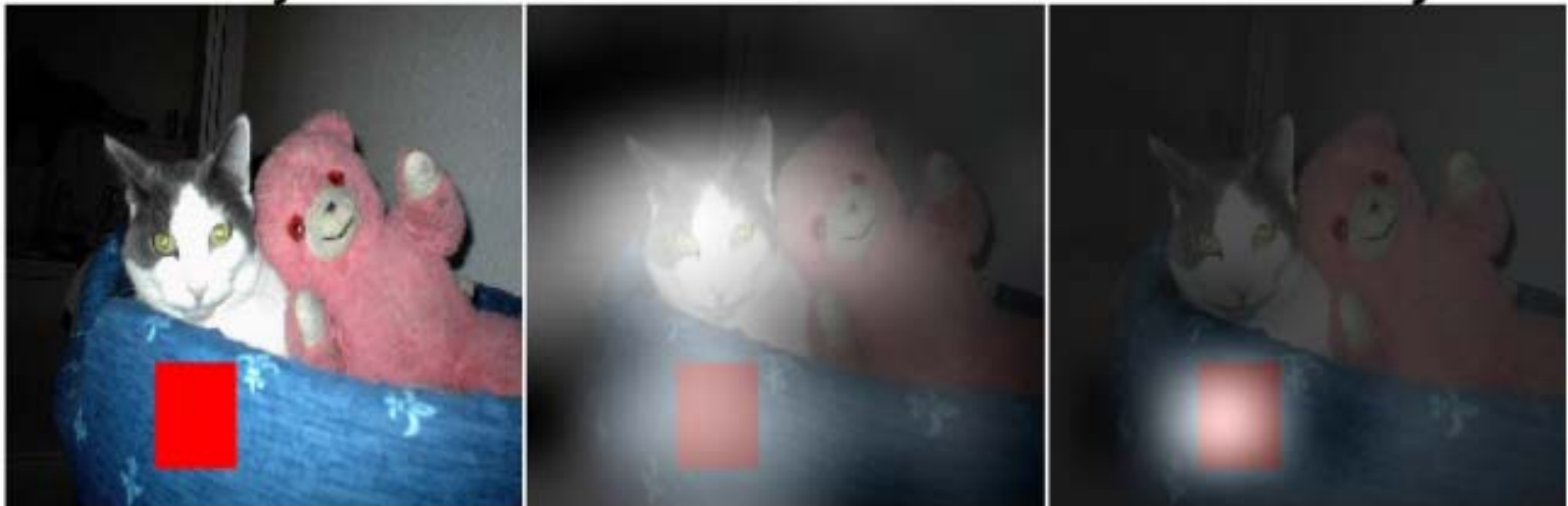
# Visual Question Answering

# Visual Question Answering

- Huijuan Xu, Kate Saenko. Ask, Attend and Answer: Exploring Question-Guided Spatial Attention for Visual Question Answering. arXiv Pre-Print, 2015



Is there a red square on the bottom of the cat?
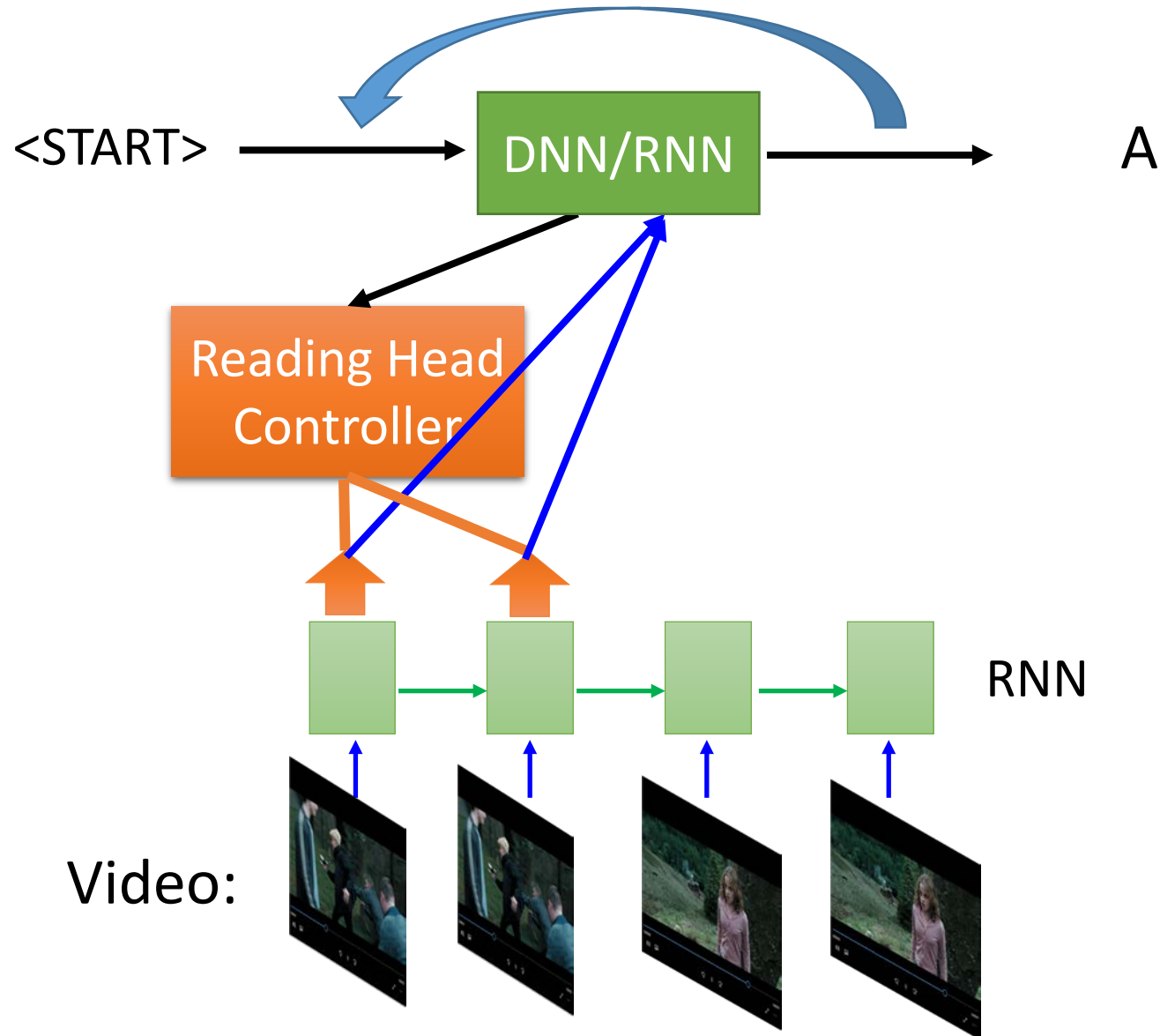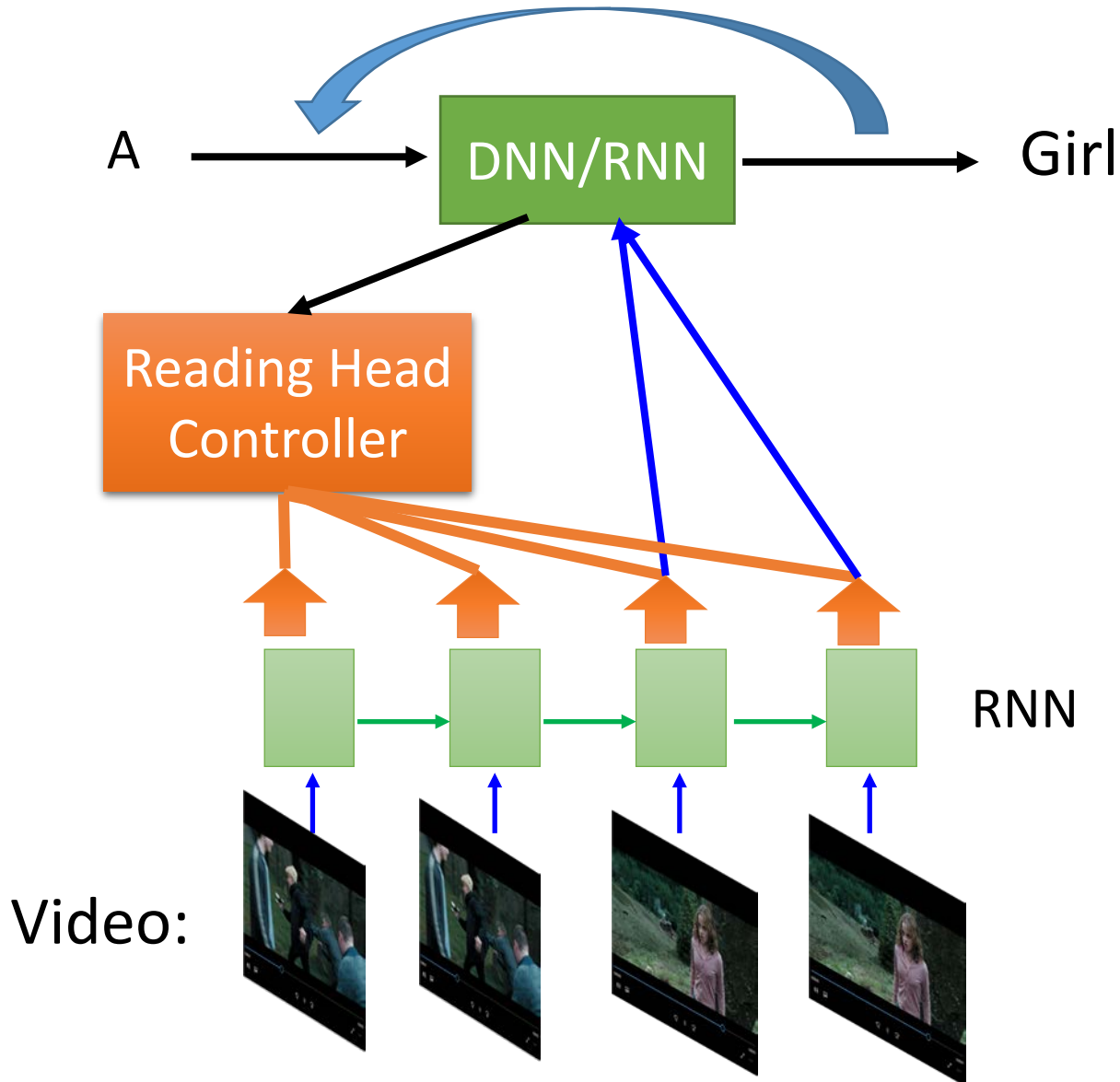GT: yes    Prediction: yes

# *Video Caption Generation*

Memory: video frames

Output: video description

<START> → DNN/RNN → A

Reading Head Controller

RNN

Video:

# Video Caption Generation

Memory: video frames

Output: video description

A → DNN/RNN → Girl

Reading Head Controller

RNN

Video:

# Video Caption Generation

- Demo: 曾柏翔、盧宏宗、吳柏瑜

# Outline

Ultra Deep Network

Attention Model

Reinforcement Learning
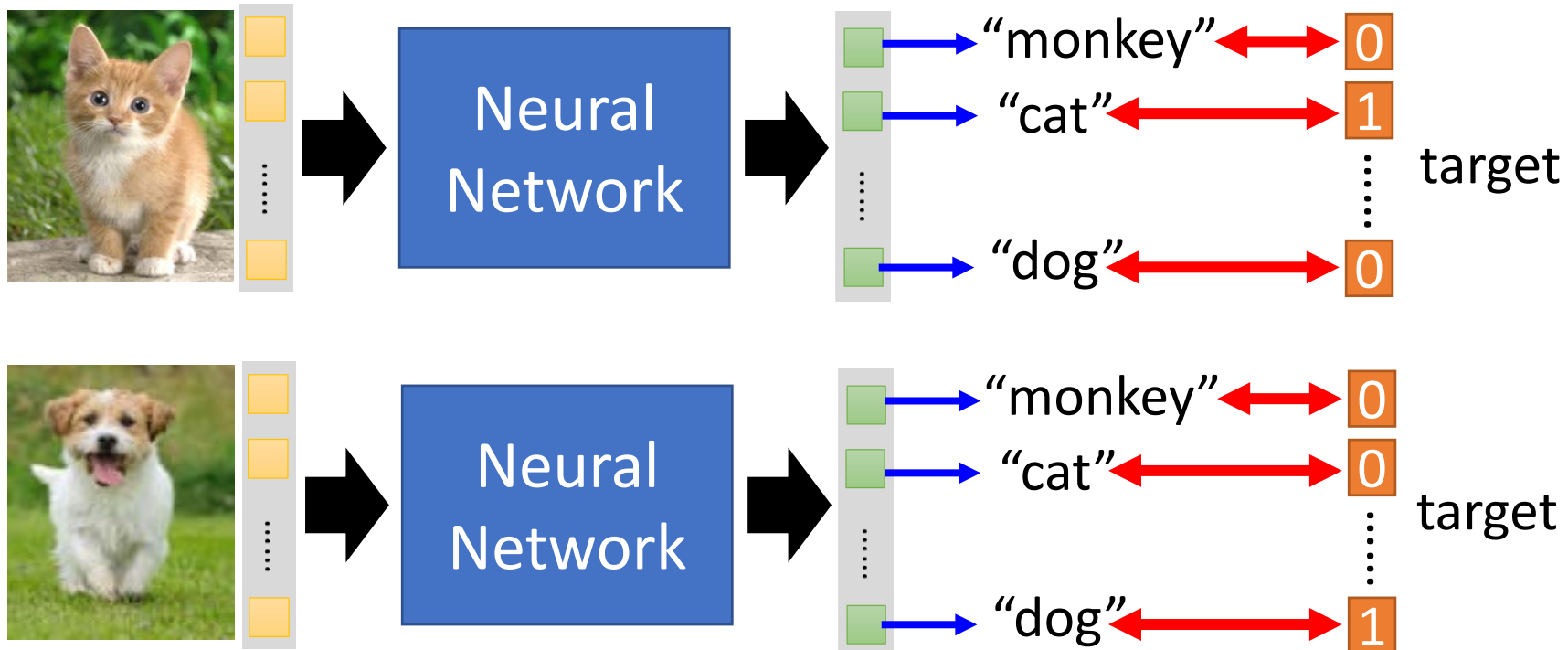
Realizing what the World Looks Like

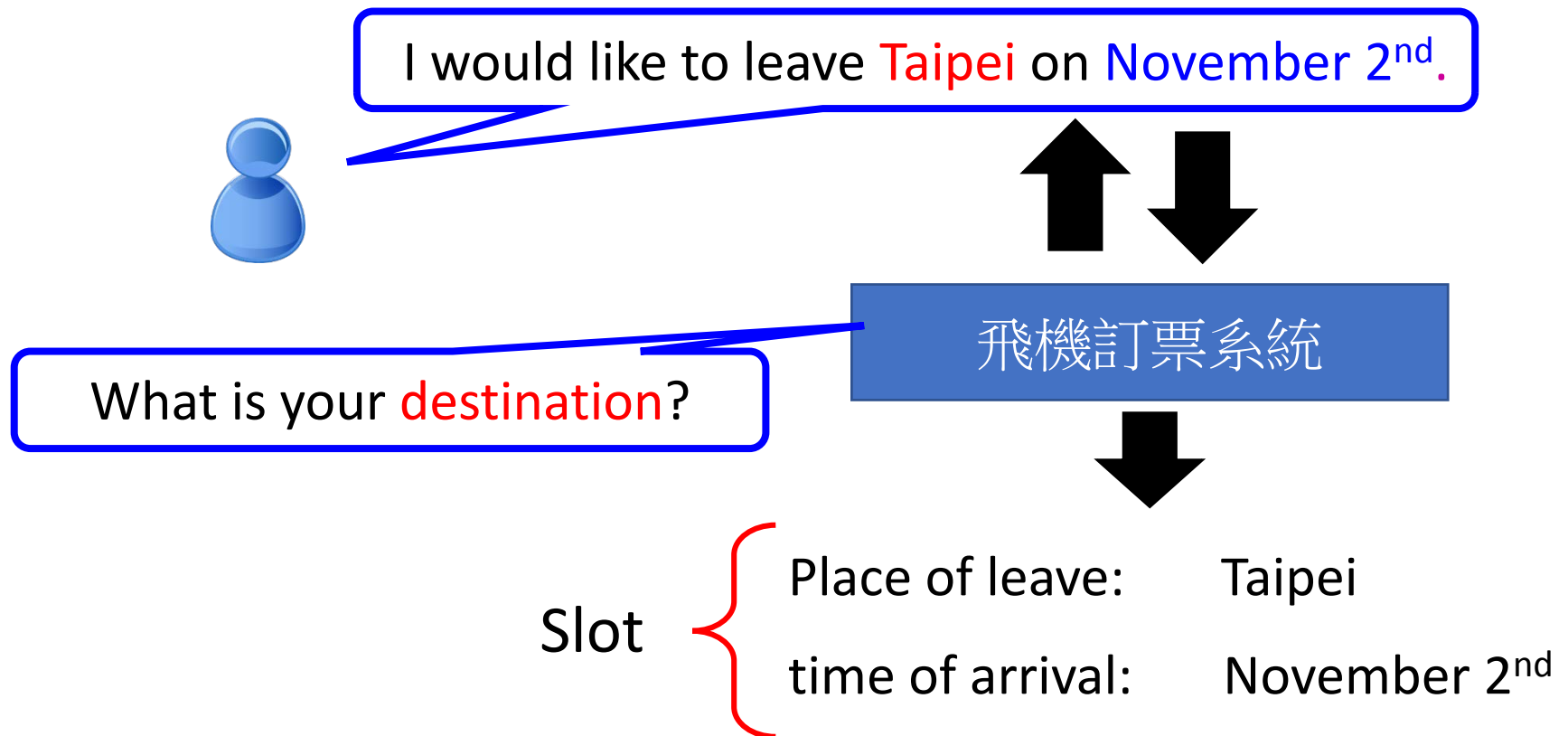Understanding the Meaning of Words

Why Deep?

# Only Supervised Learning until now ......

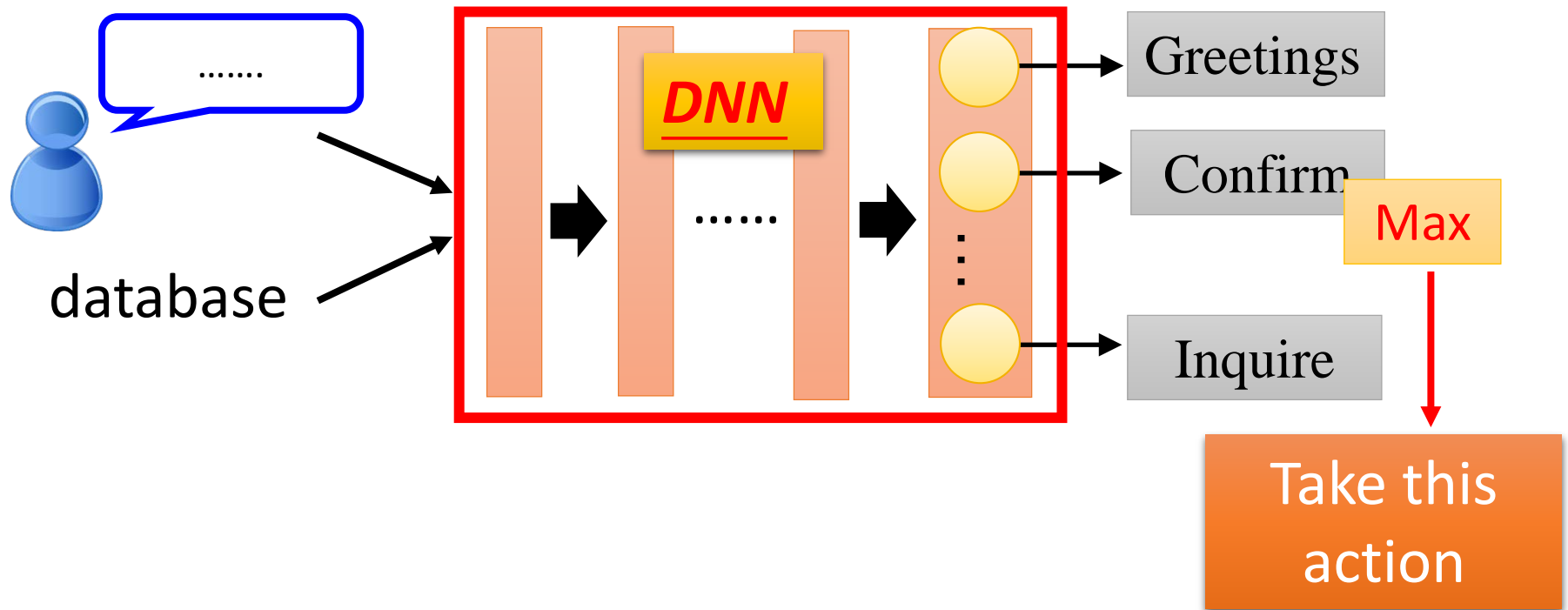- Network is a function. In supervised learning, the input-output pair is given in the training data

# Supervised v.s. Reinforcement

- Example: Dialogue Agent for 訂票系統

I would like to leave Taipei on November 2nd.

飛機訂票系統

What is your destination?

Slot {
Place of leave:    Taipei

time of arrival:    November 2nd

# Supervised v.s. Reinforcement
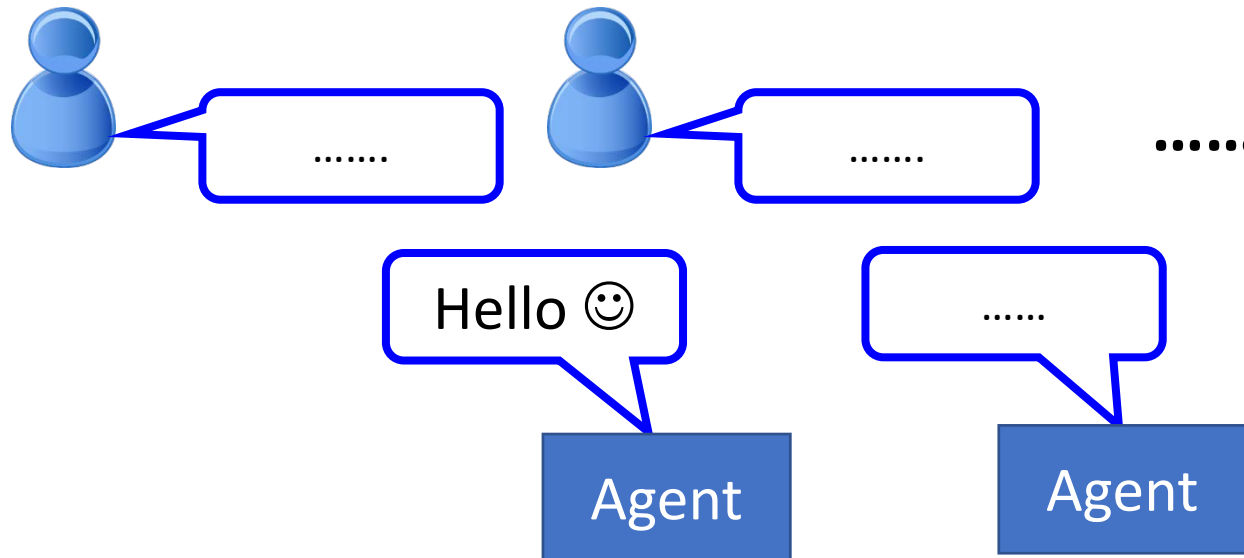
- Example: Dialogue Agent for 訂票系統

# Supervised v.s. Reinforcement

- Supervised

 "Hello"     You have to "greeting"
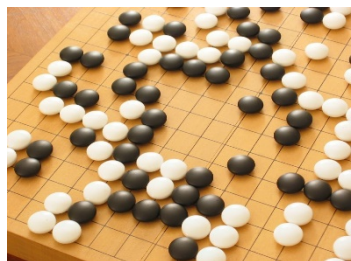
 %$#$%#     You have to "confirm"

- Reinforcement

# Supervised v.s. Reinforcement

- Playing GO
  - Supervised: 看著棋譜學

 下一步: "5-5"

 下一步: "3-3"

  - Reinforcement Learning

初手天元 ⮕ ……下了好幾百手 …… ⮕ Win!

Alpha Go is supervised learning + reinforcement learning.

# To learn deep reinforcement learning ……

- Lectures of David Silver
  - http://www0.cs.ucl.ac.uk/staff/D.Silver/web/Teaching.html
  - 10 堂課 (1:30 each)
- Deep Reinforcement Learning
  - http://videolectures.net/rldm2015_silver_reinforcement_learning/

# Outline

Ultra Deep Network
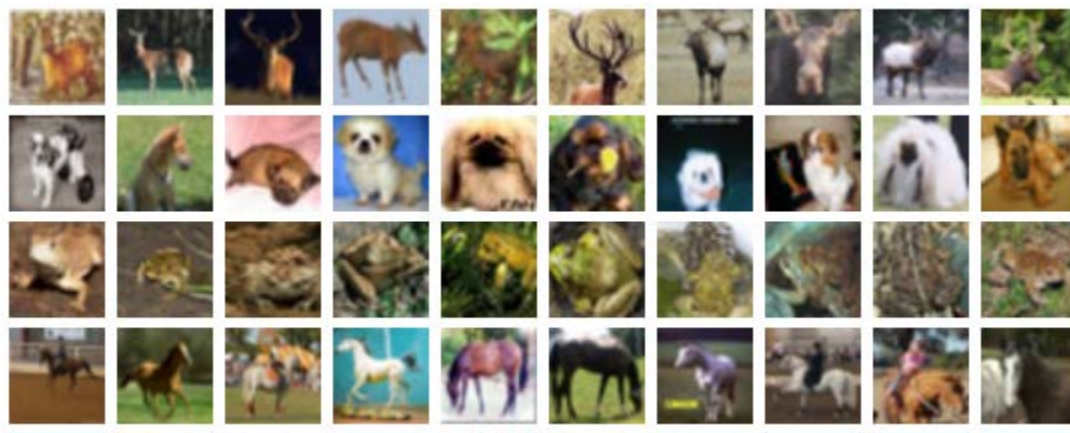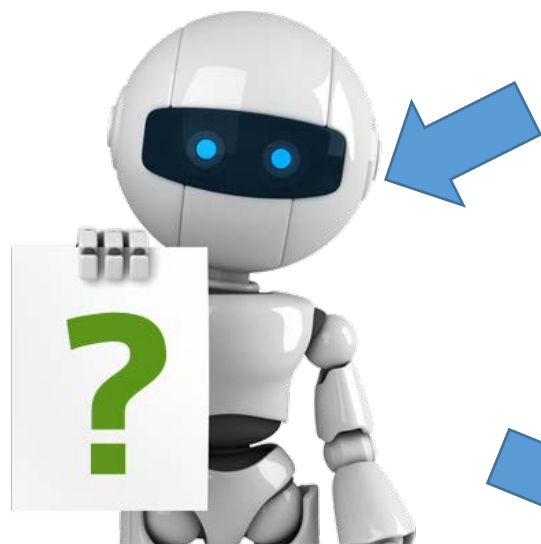
Attention Model

Reinforcement Learning

Realizing what the World Looks Like

Understanding the Meaning of Words

Why Deep?

# Does machine know what the world look like?

Ref: https://openai.com/blog/generative-models/



Draw something!

# Deep Dream

- Given a photo, machine adds what it sees ......



http://deepdreamgenerator.com/

# Deep Dream

- Given a photo, machine adds what it sees ......



http://deepdreamgenerator.com/

# Deep Style

- Given a photo, make its style like famous paintings



https://dreamscopeapp.com/

# Deep Style

- Given a photo, make its style like famous paintings



https://dreamscopeapp.com/

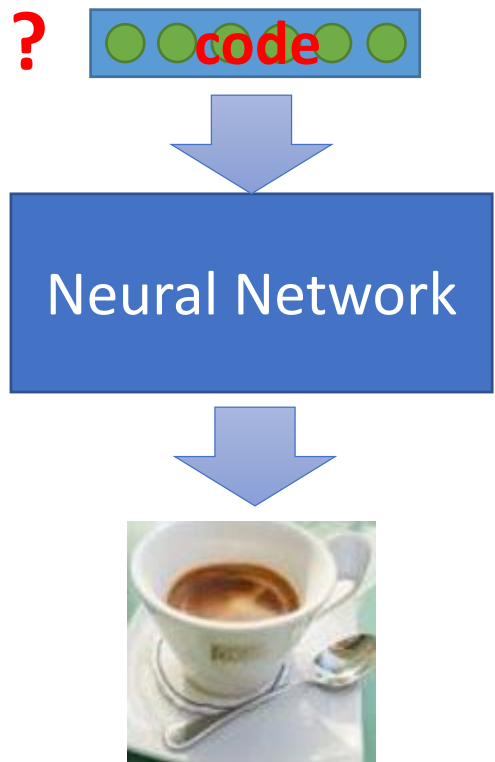*Deep Style*

# Generating Images (無中生有)

- Training a decoder to generate images is unsupervised

? [code]

Neural Network

Training data is a lot of images

# Auto-encoder

Not state-of-the-art approach



code → NN Decoder →

Learn together

→ NN Encoder → code

As close as possible

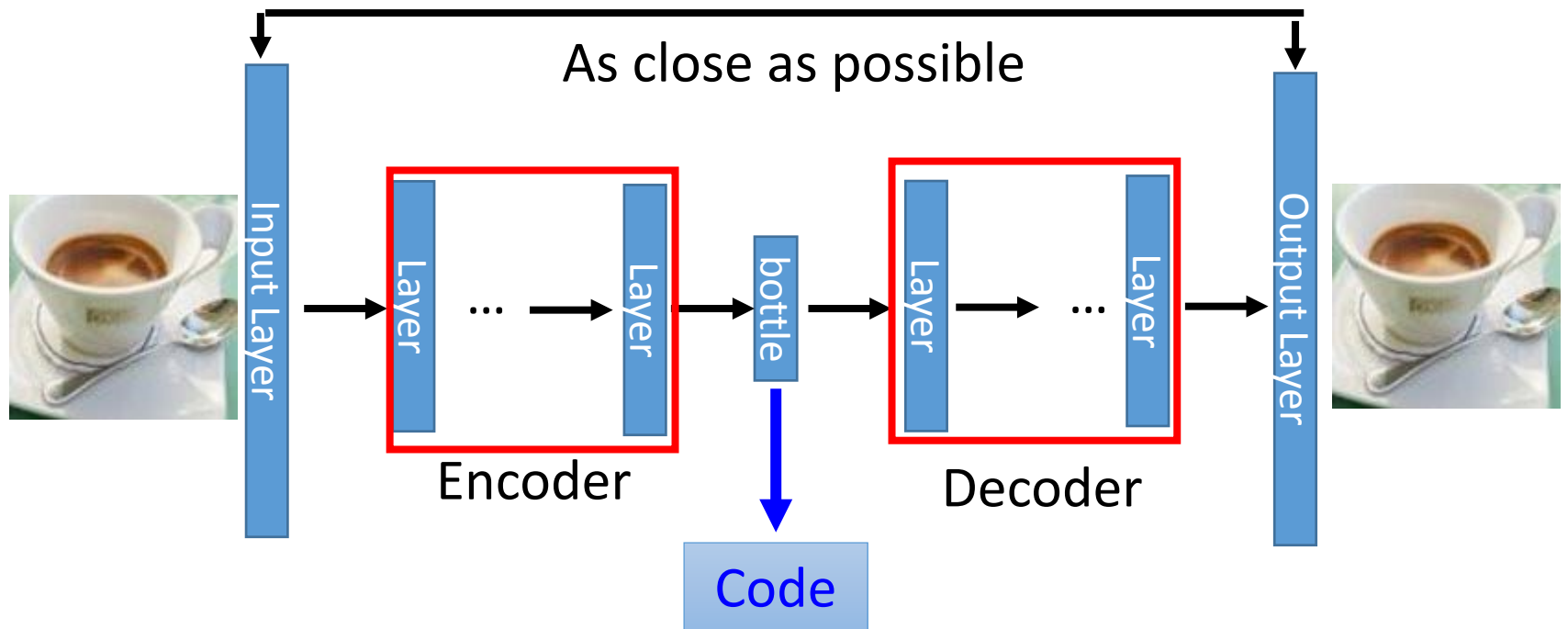Input Layer → Layer ... → Layer → bottle → Layer → ... Layer → Output Layer

Encoder

Decoder

Code

# Generating Images

- Training a decoder to generate images is unsupervised

- Variation Auto-encoder (VAE)
  - Ref: **Auto-Encoding Variational Bayes, https://arxiv.org/abs/1312.6114**

- Generative Adversarial Network (GAN)
  - Ref: **Generative Adversarial Networks, http://arxiv.org/abs/1406.2661**

# Which one is machine-generated?



Ref: https://openai.com/blog/generative-models/

# Generating Images by RNN

# Generating Images by RNN

- **Pixel Recurrent Neural Networks**
  - https://arxiv.org/abs/1601.06759



Real World

# Outline

Ultra Deep Network

Attention Model

Reinforcement Learning

Realizing what the World Looks Like

Understanding the Meaning of Words

Why Deep?

# Machine Reading

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



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

# Machine Reading

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

Word Vector / Embedding

# Machine Reading

- Generating Word Vector/Embedding is <span style="color:red">unsupervised</span>

Apple

Neural Network

**?**

Training data is a lot of text



https://garavato.files.wordpress.com/2011/11/stacksdocuments.jpg?w=490
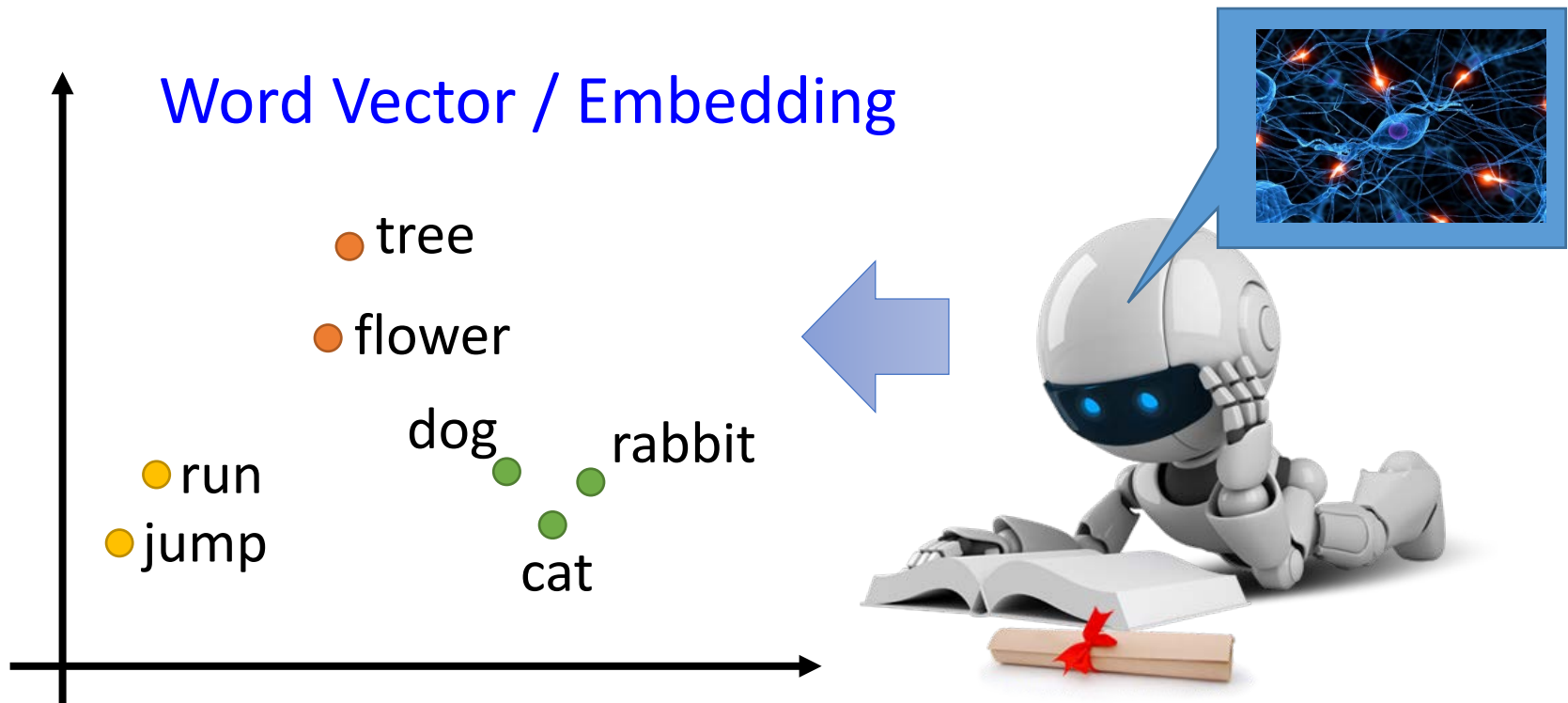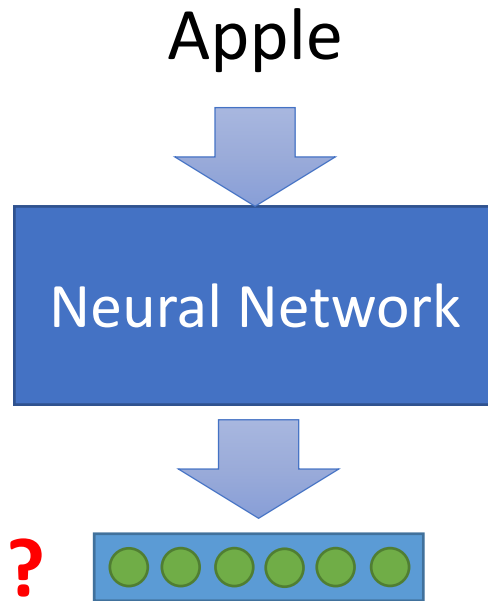
# Machine Reading

- 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
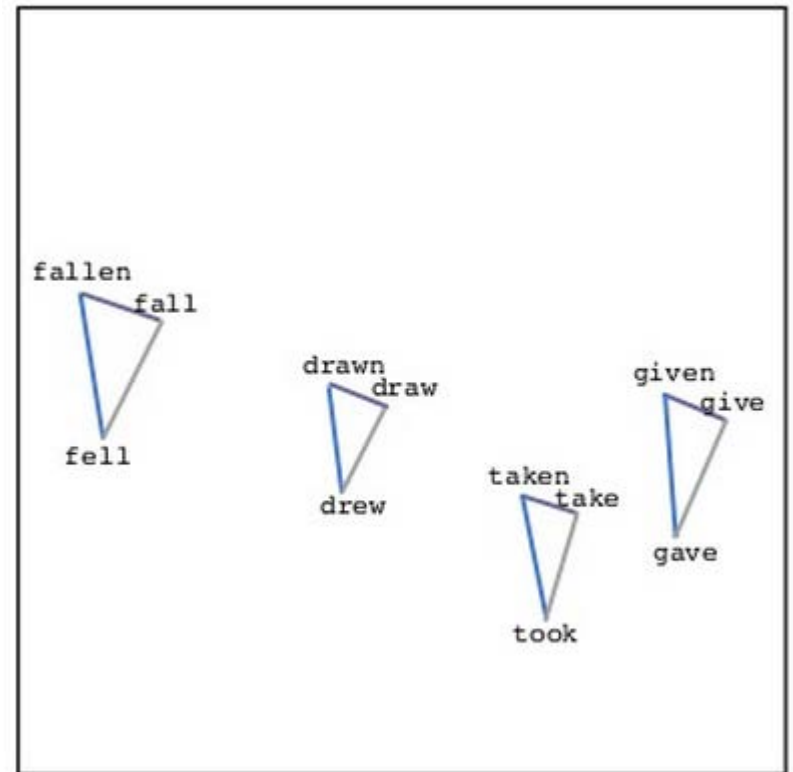
You shall know a word by the company it keeps

馬英九 520宣誓就職

蔡英文 520宣誓就職

# Word Vector



Source: http://www.slideshare.net/hustwj/cikm-keynotenov2014

# Word Vector

$$V(Germany)$$
$$\approx V(Berlin) - V(Rome) + V(Italy)$$

- Characteristics

$$V(hotter) - V(hot) \approx V(bigger) - V(big)$$
$$V(Rome) - V(Italy) \approx V(Berlin) - V(Germany)$$
$$V(king) - V(queen) \approx V(uncle) - V(aunt)$$

- Solving analogies

Rome : Italy = Berlin : ?

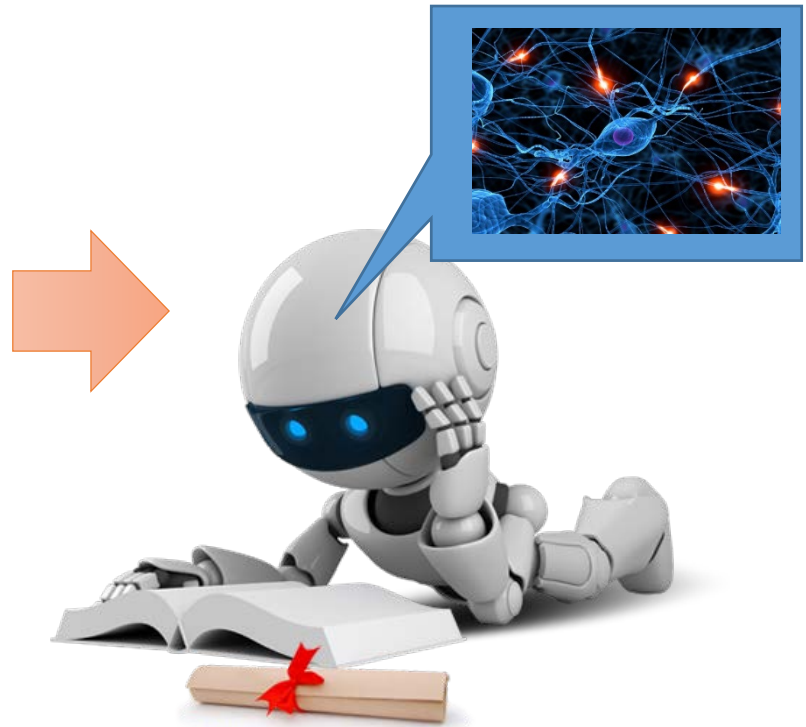Compute $V(Berlin) - V(Rome) + V(Italy)$

Find the word w with the closest V(w)

# Machine Reading

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



Machine learns to understand 鄉民用語 via reading the posts on PTT

# Demo

- Model used in demo is provided by 陳仰德
- Part of the project done by 陳仰德、林資偉
- TA: 劉元銘
- Training data is from PTT (collected by 葉青峰)

# Outline

Ultra Deep Network

Attention Model

Reinforcement Learning

Realizing what the World Looks Like

Understanding the Meaning of Words

Why Deep?

# Deeper is Better?

| Layer X Size | Word Error Rate (%) |
|---|---|
| 1 X 2k | 24.2 |
| 2 X 2k | 20.4 |
| 3 X 2k | 18.4 |
| 4 X 2k | 17.8 |
| 5 X 2k | 17.2 |
| 7 X 2k | 17.1 |
|  |  |

Not surprised, more parameters, better performance

Seide, Frank, Gang Li, and Dong Yu. "Conversational Speech Transcription Using Context-Dependent Deep Neural Networks." *Interspeech*. 2011.
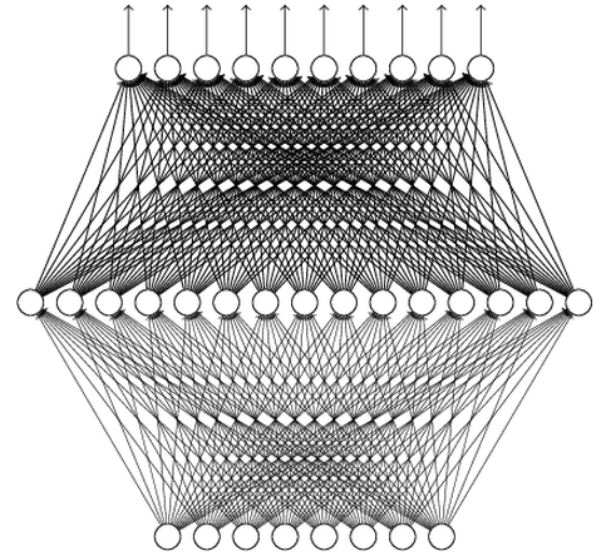
# Universality Theorem

Any continuous function f

$$f : R^N \rightarrow R^M$$

Can be realized by a network with one hidden layer
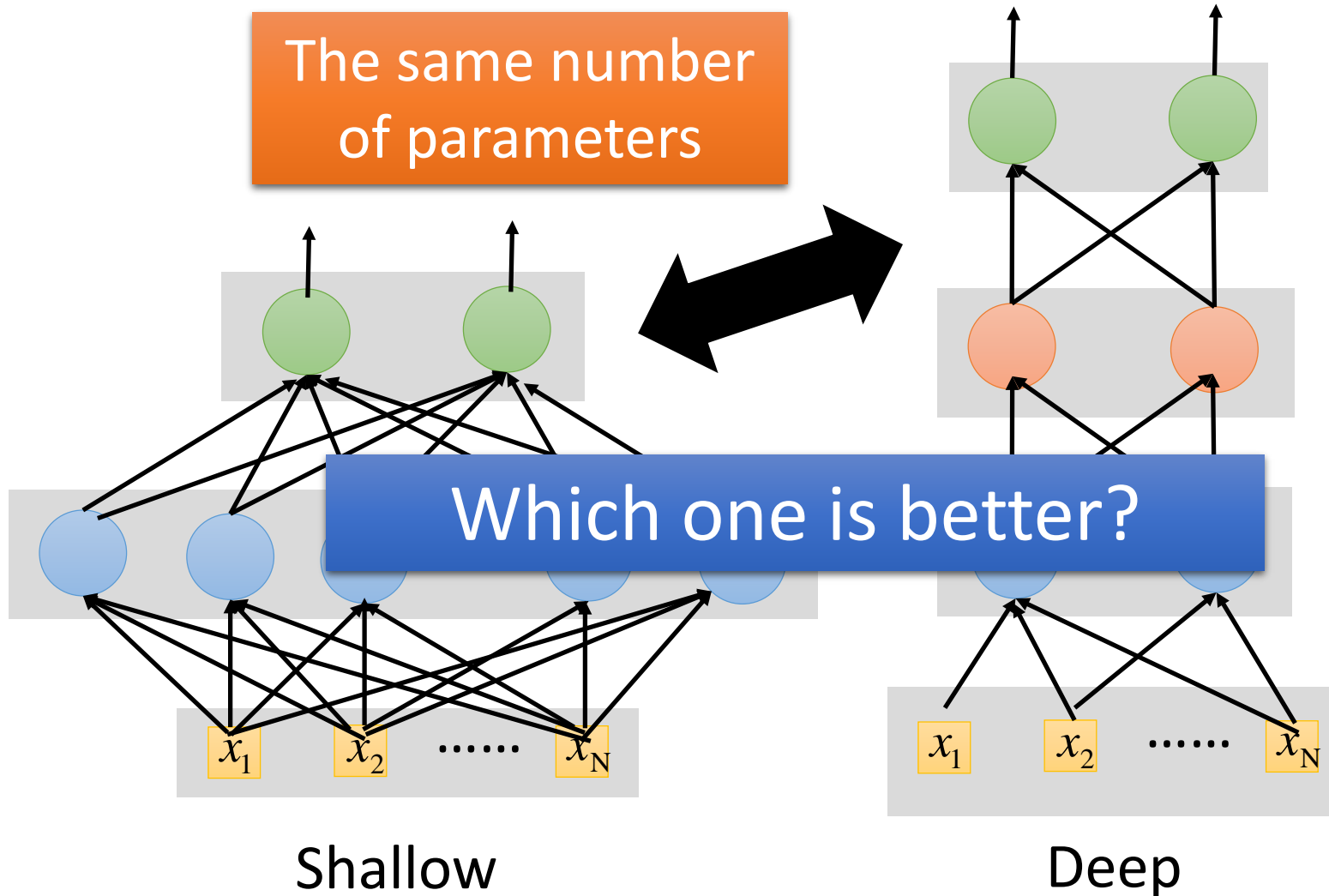
(given **enough** hidden neurons)



Reference for the reason:
http://neuralnetworksandde
eplearning.com/chap4.html

Why "Deep" neural network not "Fat" neural network?

# Fat + Short v.s. Thin + Tall



The same number of parameters

Which one is better?

$x_1$  $x_2$  ......  $x_N$

Shallow

$x_1$  $x_2$  ......  $x_N$

Deep

# Fat + Short v.s. Thin + Tall

| Layer X Size | Word Error Rate (%) | Layer X Size | Word Error Rate (%) |
|---|---|---|---|
| 1 X 2k | 24.2 | | |
| 2 X 2k | 20.4 | | |
| 3 X 2k | 18.4 | | |
| 4 X 2k | 17.8 | | |
| 5 X 2k | 17.2 | 1 X 3772 | 22.5 |
| 7 X 2k | 17.1 | 1 X 4634 | 22.6 |
| | | 1 X 16k | 22.1 |

Why?

Seide, Frank, Gang Li, and Dong Yu. "Conversational Speech Transcription Using Context-Dependent Deep Neural Networks." *Interspeech*. 2011.
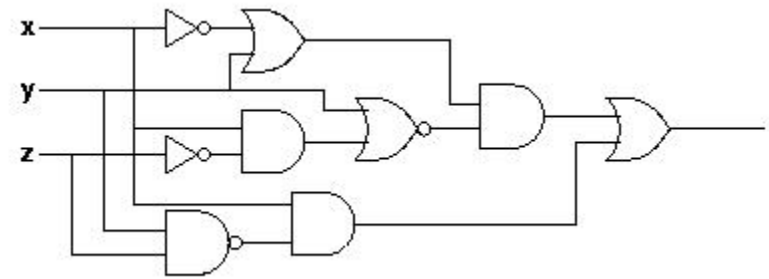
# Analogy



**Logic circuits**

- 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**

- 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
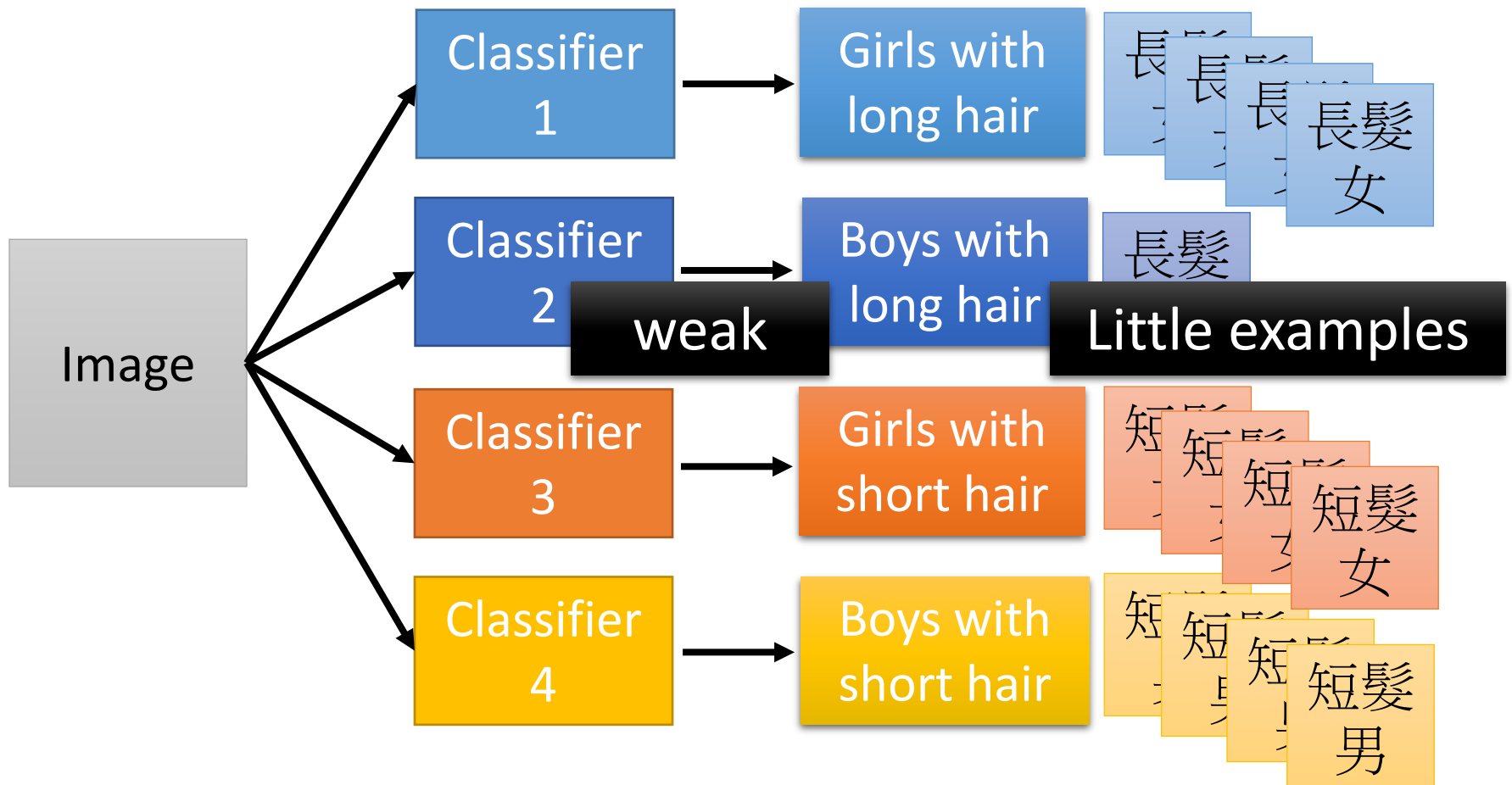
  → less parameters → less data?
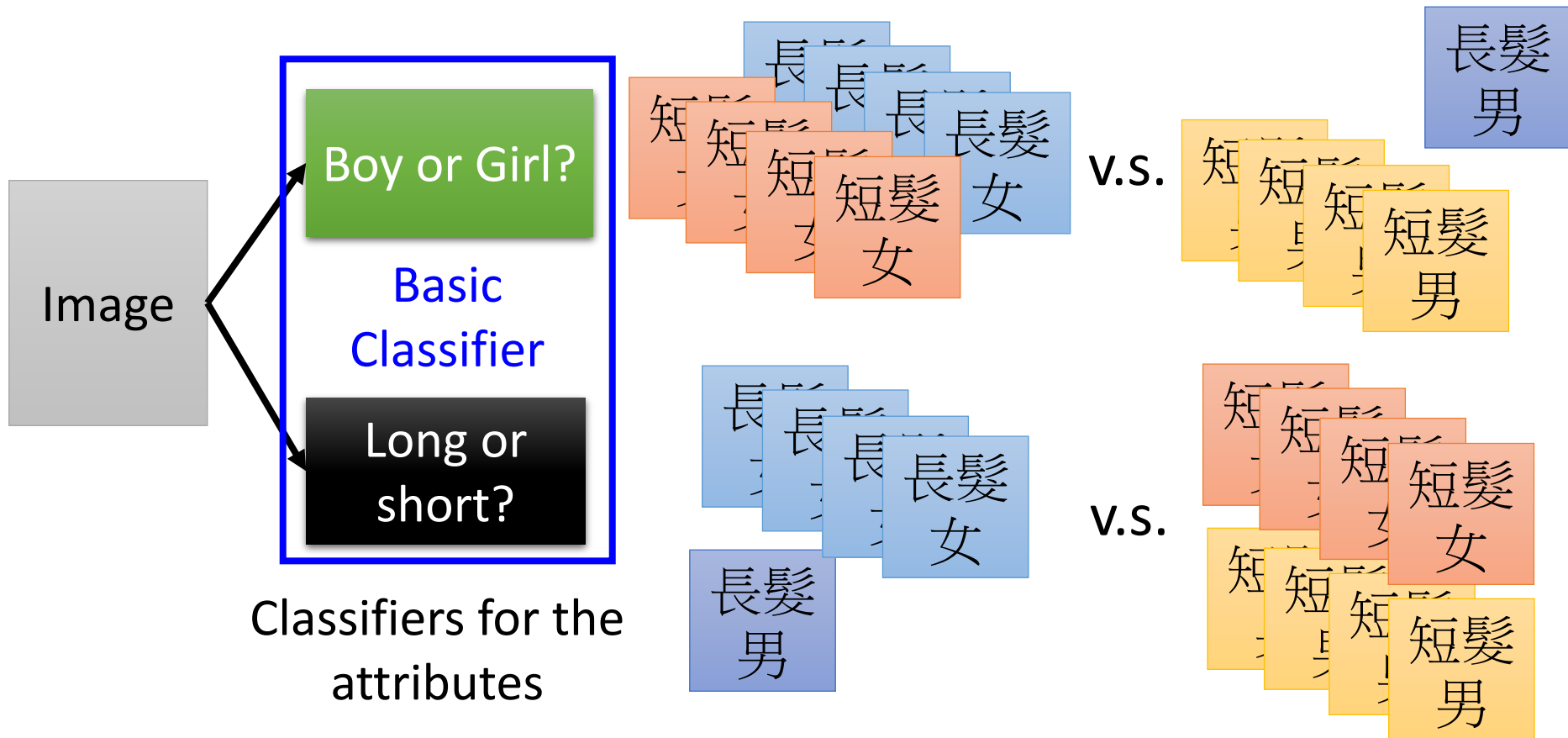
This page is for EE background.

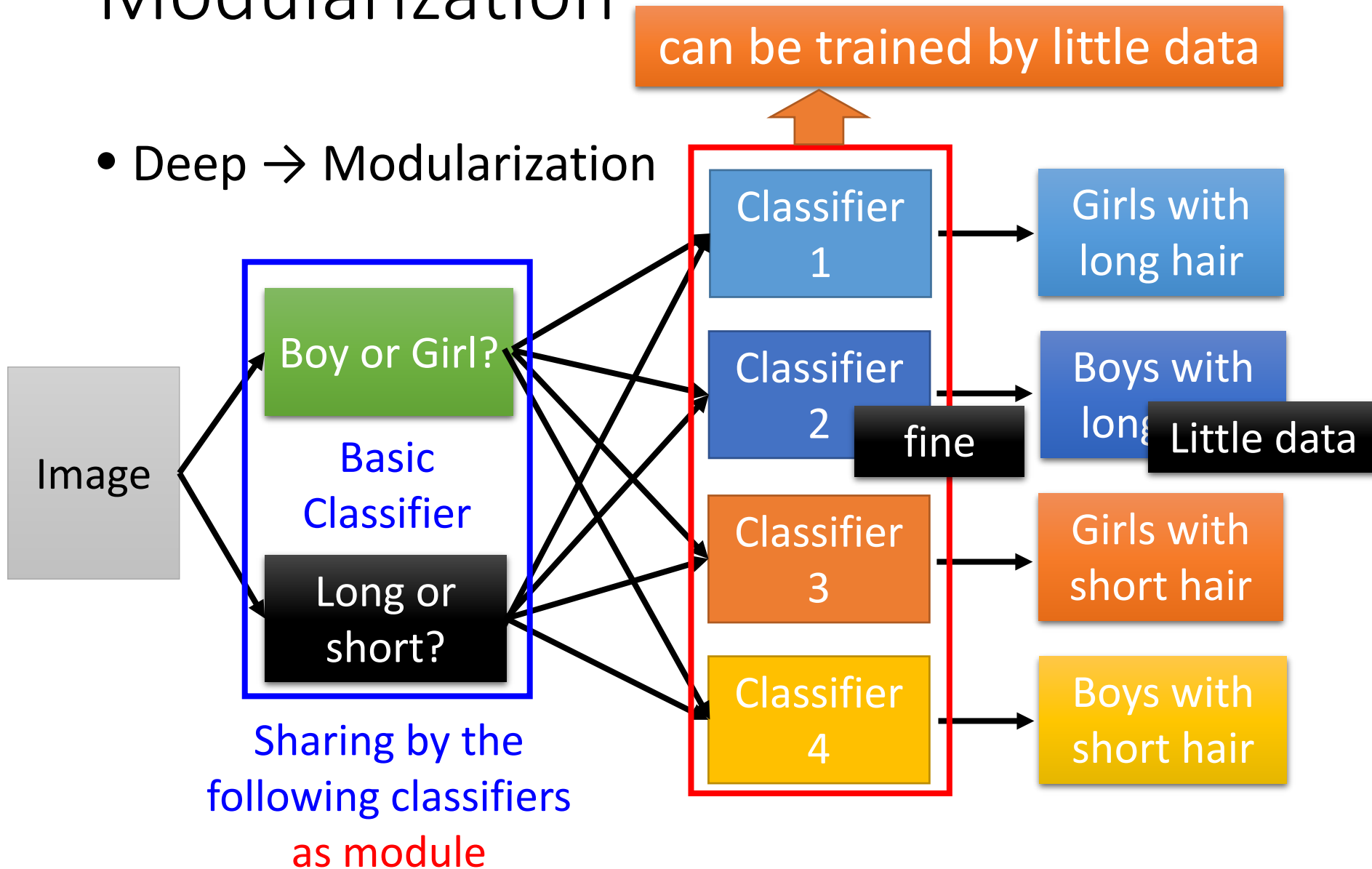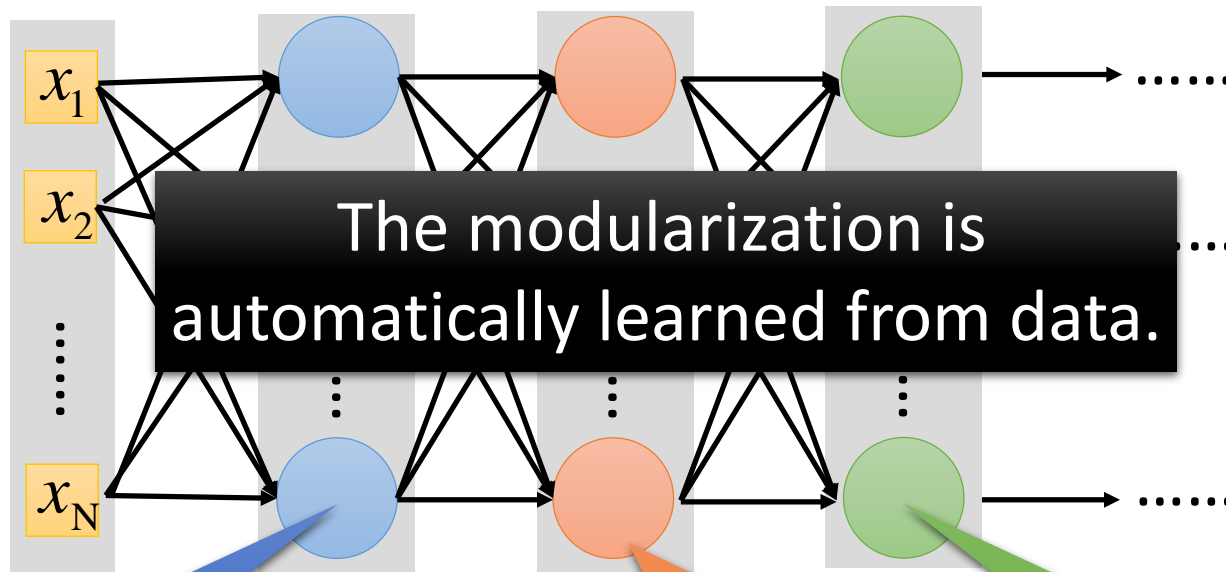# Modularization

- Deep → Modularization

# Modularization

- Deep → Modularization → Less training data?



The modularization is automatically learned from data.

The most basic classifiers

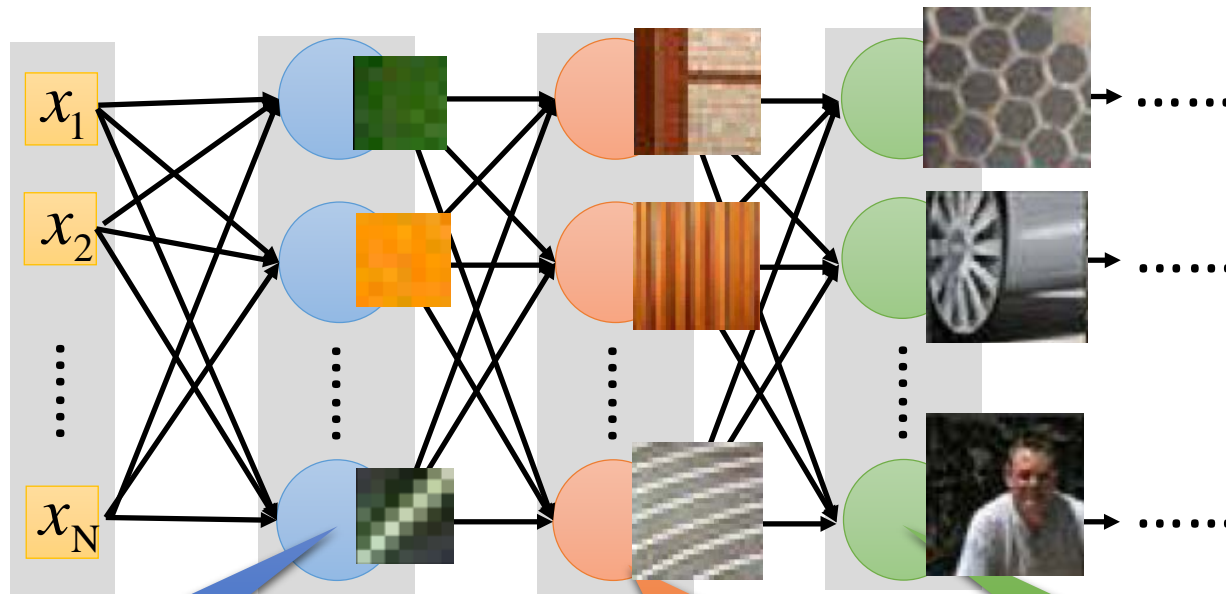Use 1$^{st}$ layer as module to build classifiers

Use 2$^{nd}$ layer as module ……

# Modularization

Reference: Zeiler, M. D., & Fergus, R. (2014). Visualizing and understanding convolutional networks. In *Computer Vision–ECCV 2014* (pp. 818-833)

- Deep → Modularization



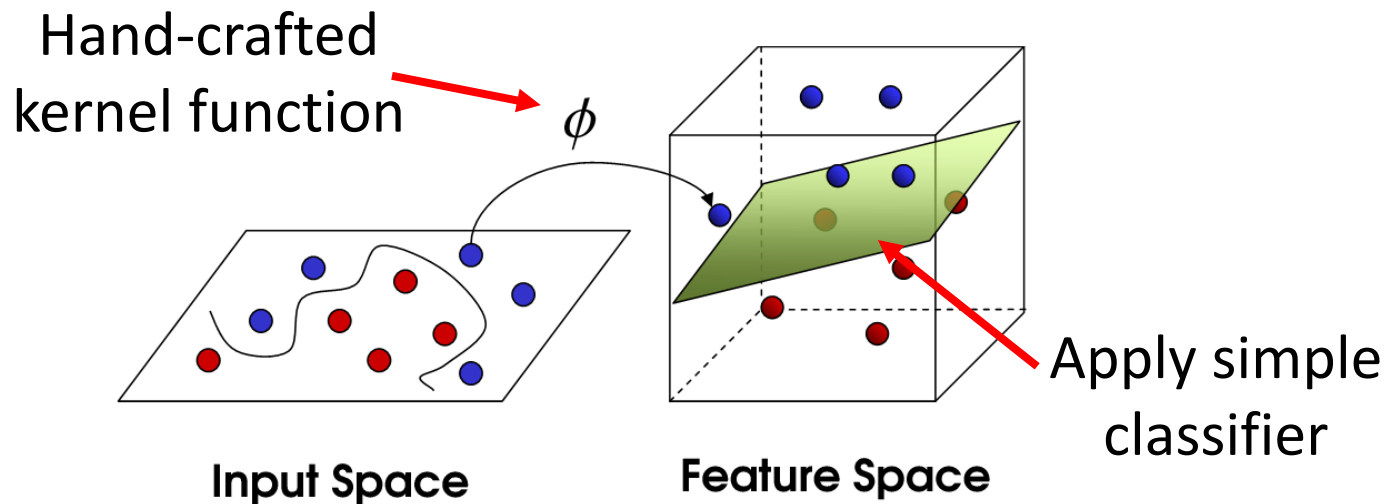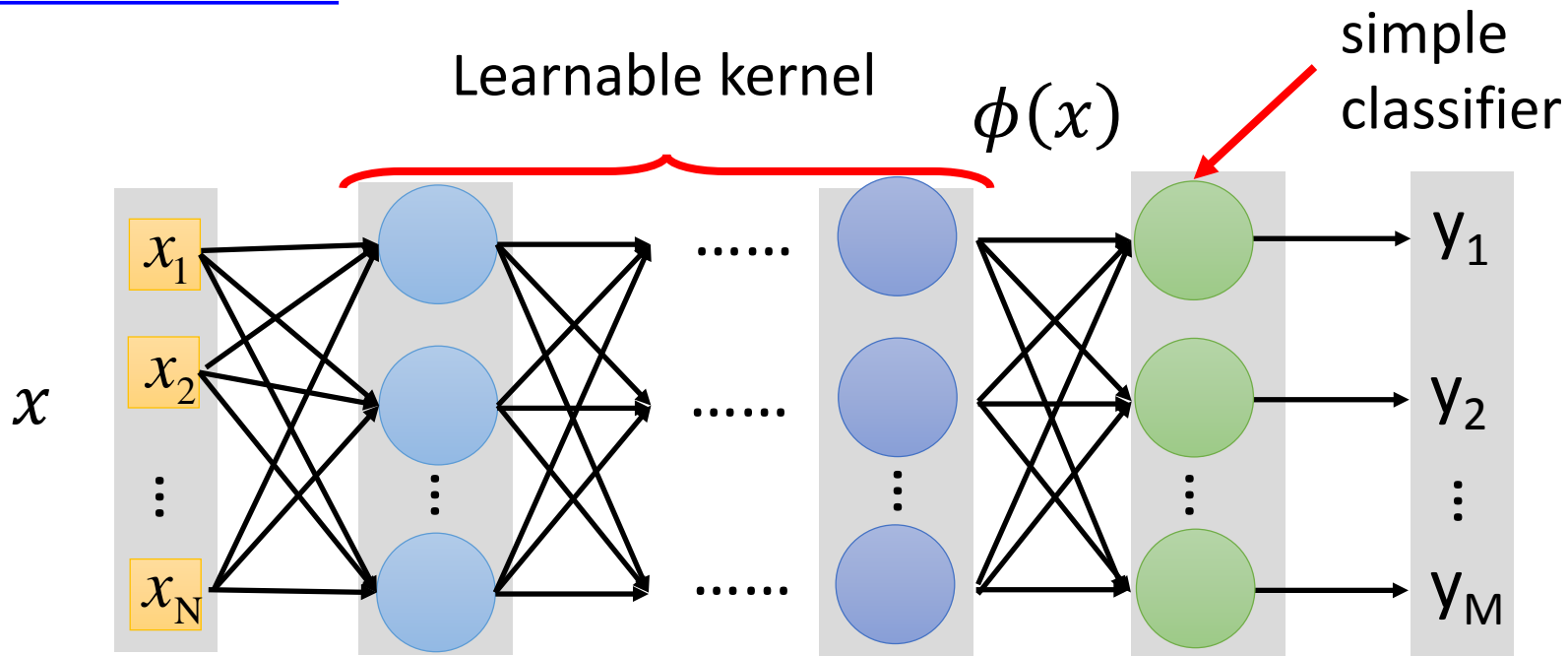The most basic classifiers

Use 1st layer as module to build classifiers

Use 2nd layer as module ......

**SVM**

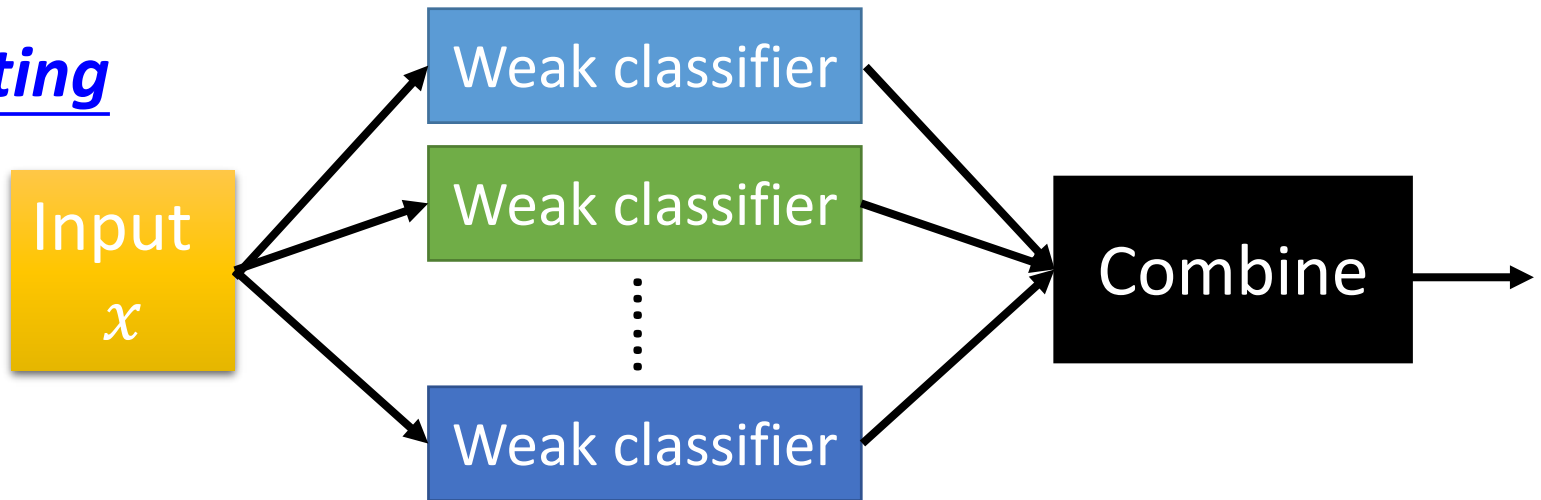Hand-crafted kernel function → $\phi$

Input Space

Feature Space

Apply simple classifier

Source of image: http://www.gipsa-lab.grenoble-inp.fr/transfert/seminaire/455_Kadri2013Gipsa-lab.pdf

**Deep Learning**

Learnable kernel

$\phi(x)$

simple classifier

$x$ : $x_1$, $x_2$, ..., $x_N$ → ...... → $y_1$, $y_2$, ..., $y_M$

*Boosting*

Input
$x$

Weak classifier

Weak classifier

Weak classifier

Combine

*Deep Learning*

Weak classifier

Boosted weak classifier

Boosted Boosted weak classifier

$x_1$

$x_2$

$x_N$

# More Reasons

- Do Deep Nets Really Need To Be Deep? (by Rich Caruana)
- http://research.microsoft.com/apps/video/default.aspx?id=232373&r=1



Do deep nets really need to be deep?

Rich Caruana
Microsoft Research

Lei Jimmy Ba
MSR Intern, University of Toronto

Thanks also to: Gregor Urban, Krzysztof Geras, Samira Kahou, Abdelrahman Mohamed, Jinyu Li, Rui Zhao, Jui-Ting Huang, and Yifan Gong



Yes!

Thank You

Any Questions?

keynote of Rich Caruana at ASRU 2015

# Concluding Remarks

# Today's Lecture

Lecture I: Introduction of Deep Learning

Lecture II: Tips for Training Deep Neural Network

Lecture III: Variants of Neural Network
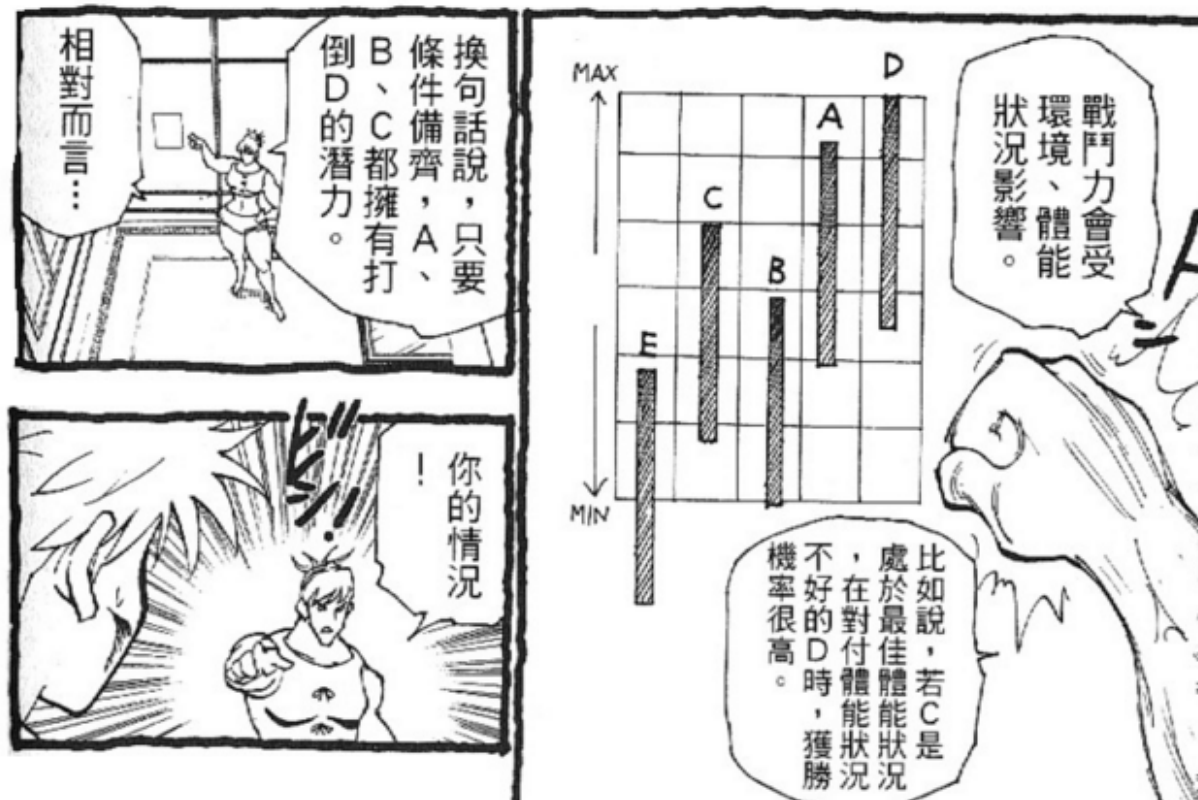
Lecture IV: Next Wave

# Some Opinions

• Also learn other machine learning methods

不想知道 Deep Learning 以外的方法

# Some Opinions

- In some situations, the simpler machine learning methods can be very powerful.

# Some Opinions

• 寒武纪大爆炸



已經有一些生物滅絕了

# Some Opinions

- Deep Learning is still at the phase of "神農嘗百草"

- Lots of questions still do not have answers

- However, probably also easy to enter

http://orchid.shu.edu.tw/upload/article/20110927181605_1_pic.png

# 如果你想 "深度學習 深度學習"

- "Neural Networks and Deep Learning"
  - written by Michael Nielsen
  - http://neuralnetworksanddeeplearning.com/
- "Deep Learning"
  - Written by Yoshua Bengio, Ian J. Goodfellow and Aaron Courville
  - http://www.iro.umontreal.ca/~bengioy/dlbook/
- Course: Machine learning and having it deep and structured
  - http://speech.ee.ntu.edu.tw/~tlkagk/courses_MLSD15_2.html

# 給資料科學愛好者

- 台大電機系於台大電信所成立「資料科學與智慧網路組」，開始招收碩、博士生
- 今年秋天開始報名