# Introduction to Deep Learning and Tensorflow

Veerasak Kritsanapraphan
Software Park Thailand

Deep learning
attracts lots of attention.

◎ Google Trends

Deep learning obtains many exciting results.

➡ The talks in this afternoon

This talk will focus on the technical part.

| | | | | |
|---|---|---|---|---|
| 2007 | 2009 | 2011 | 2013 | 2015 |

# Outline

- Part I: Introduction of Deep Learning

- Part II: Why Deep?

- Part III: Tips for Training Deep Neural Network

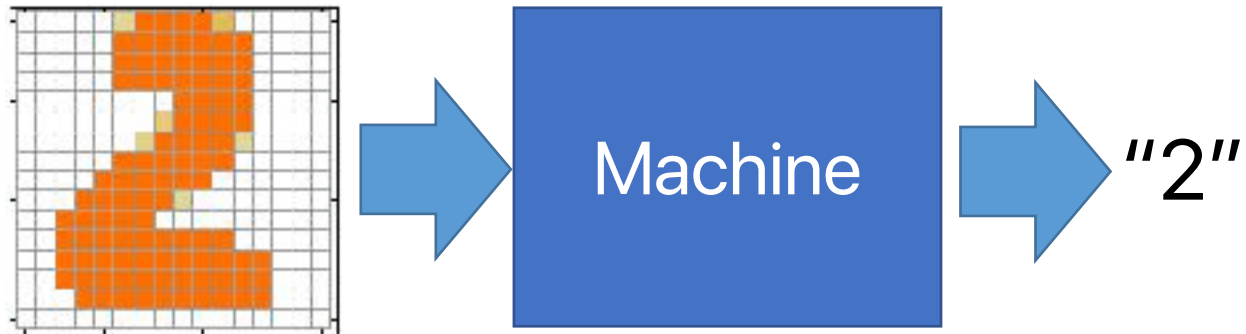- Part IV: Neural Network with Memory

# Part I:
# Introduction of
# Deep Learning

# Example Application

◎ Handwriting Digit Recognition

# Handwriting Digit Recognition

◎ Input

$x_1$

$x_2$

$x_{256}$

16 x 16 = 256

Ink → 1
No ink → 0

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

$x_1$

$x_2$

⋮

$x_{256}$

Machine

$$f: R^{256} \rightarrow R^{10}$$

$y_1$

$y_2$ ,

⋮

$y_{10}$

In deep learning, the function is represented by neural network

# Element of Neural Network

**_Neuron_**   $f: R^K \to R$



$$z = a_1 w_1 + a_2 w_2 + \boxed{?} + a_K w_K + b$$
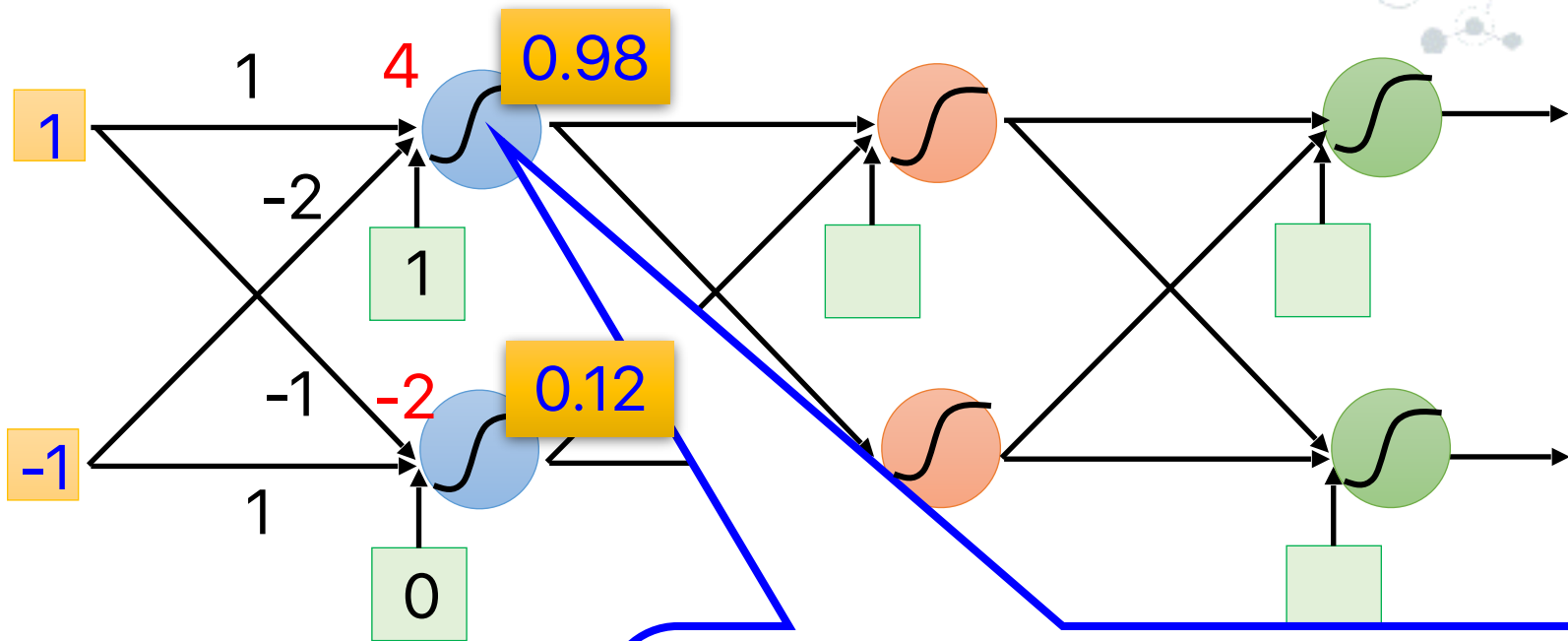
weights

bias

Activation function
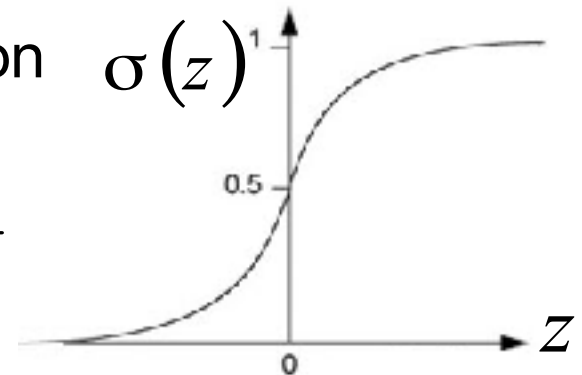
# Neural Network



Deep means many hidden layers

# Example of Neural Network



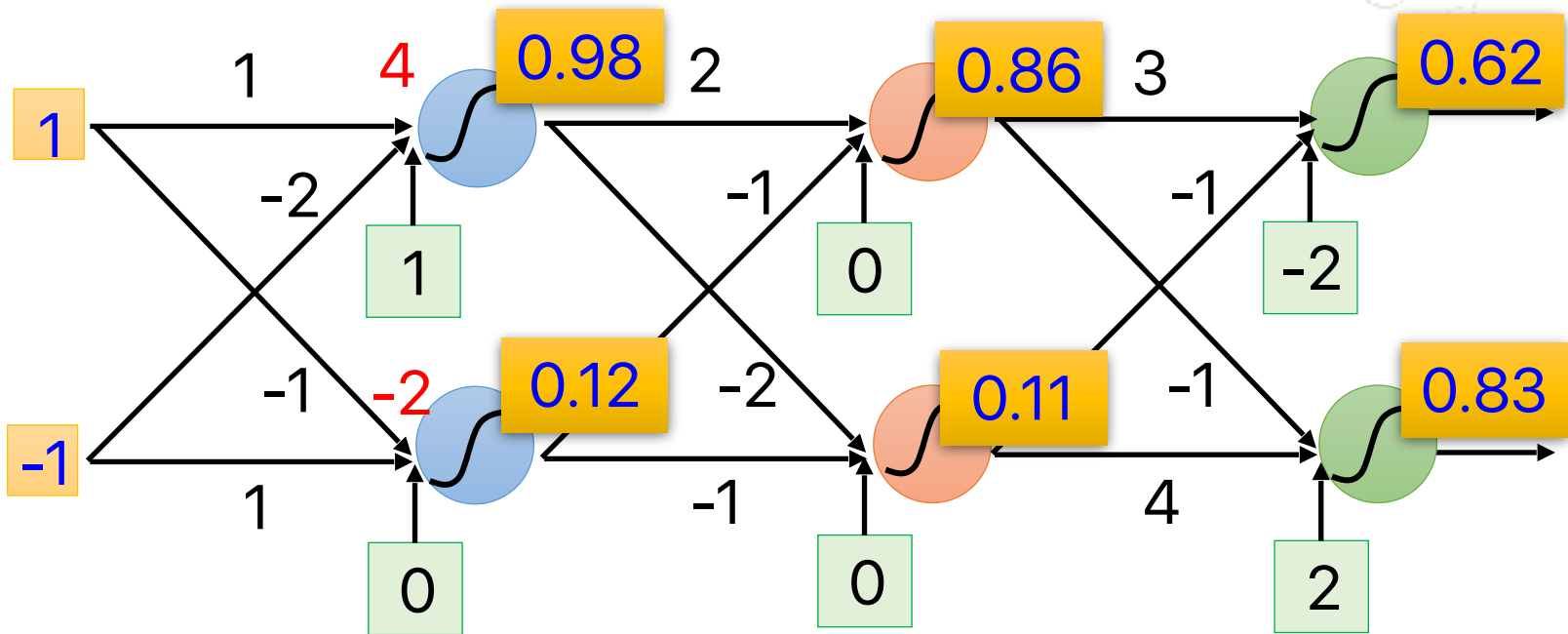Sigmoid Function $\sigma(z)$
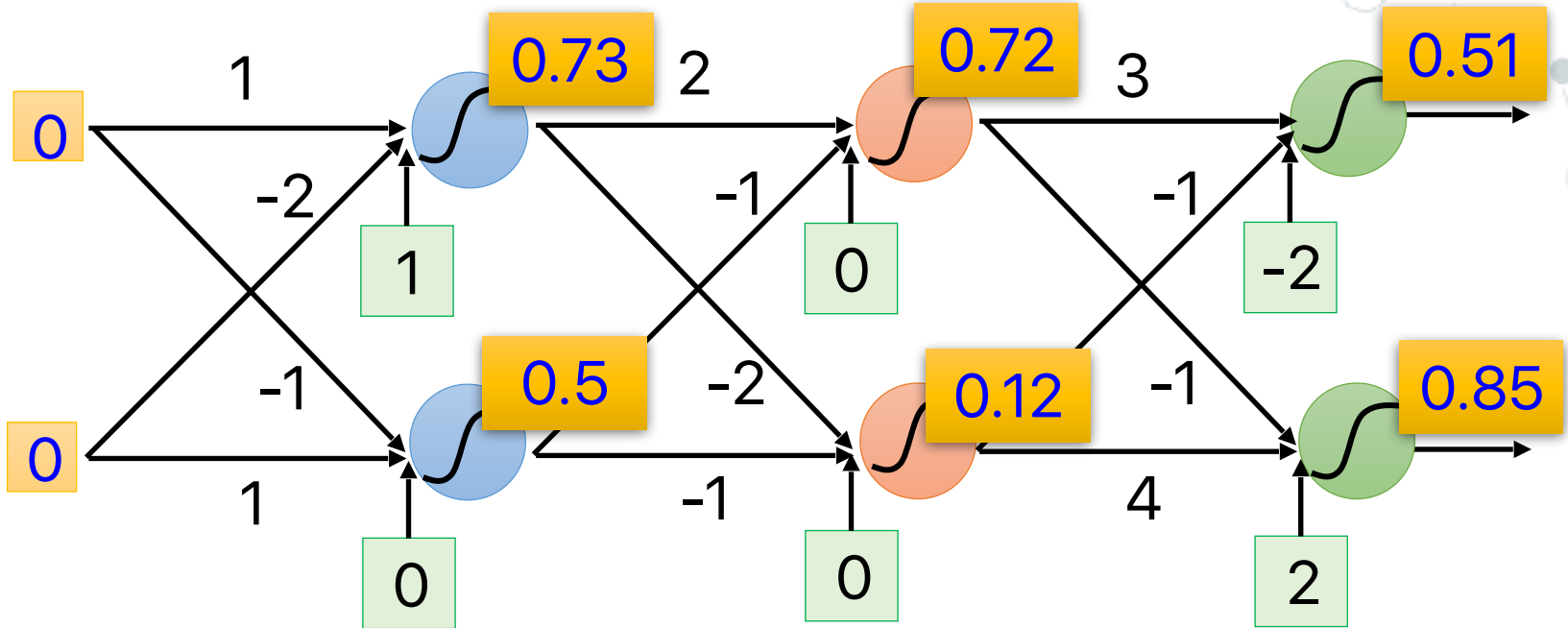
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

# Example of Neural Network

# Example of Neural Network



$$f: R^2 \rightarrow R^2$$

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

**Different parameters define different function**

# Matrix Operation



$$\sigma\left(\begin{bmatrix} 1 & -2 \\ -1 & 1 \end{bmatrix}\begin{bmatrix} 1 \\ -1 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix}\right) = \begin{bmatrix} 0.98 \\ 0.12 \end{bmatrix}$$

$$\begin{bmatrix} 4 \\ -2 \end{bmatrix}$$

# Neural Network



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

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

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

14

# Neural Network



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

Using parallel computing techniques
to speed up matrix operation

$$= \sigma(\,W^L \ldots \sigma(\,W^2\,\sigma(\,W^1\,x\, +\, b^1\,)\, +\, b^2\, \ldots +\, b^L\,)$$

# Softmax

◎ Softmax layer as the output layer

**_Ordinary Layer_**

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

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

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

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

May not be easy to interpret

# Softmax

◎ Softmax layer as the output layer

**_Softmax Layer_**



$$y_1 = e^{z_1} \bigg/ \sum_{j=1}^{3} e^{z_j}$$

$$y_2 = e^{z_2} \bigg/ \sum_{j=1}^{3} e^{z_j}$$

$$y_3 = e^{z_3} \bigg/ \sum_{j=1}^{3} e^{z_j}$$

$z_1$ — **3** — $e$ — $e^{z_1}$ — **20** — ÷ — **0.88**

$z_2$ — **1** — $e$ — $e^{z_2}$ — **2.7** — ÷ — **0.12**

$z_3$ — **-3** — $e$ — $e^{z_3}$ — **0.05** — ÷ — **≈0**

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

# How to set network parameters

$$\theta = \left\{ W^1, b^1, W^2, b^2, \cdots W^L, b^L \right\}$$



16 x 16 = 256

Ink → 1
No ink → 0

Set the network parameters such that ......

Input ......... m value

How to let the neural network achieve this

Input: $y_2$ has the maximum value

18

# Training Data

◎ Preparing training data: images and their labels

"5" "0" "4" "1"

"9" "2" "1" "3"

Using the training data to find the network parameters.

# Cost

Given a set of network parameters , each example has a cost value.



"1"

$x_1$
$x_2$
$x_{256}$

0.2
0.3
0.5

Cost

$L(\theta)$

1
0
0

target

Cost can be Euclidean distance or cross entropy of the network output and target

# Total Cost

For all training data ...



Total Cost:

$$C(\theta) = \sum_{r=1}^{R} L^r(\theta)$$

How bad the network parameters is on this task

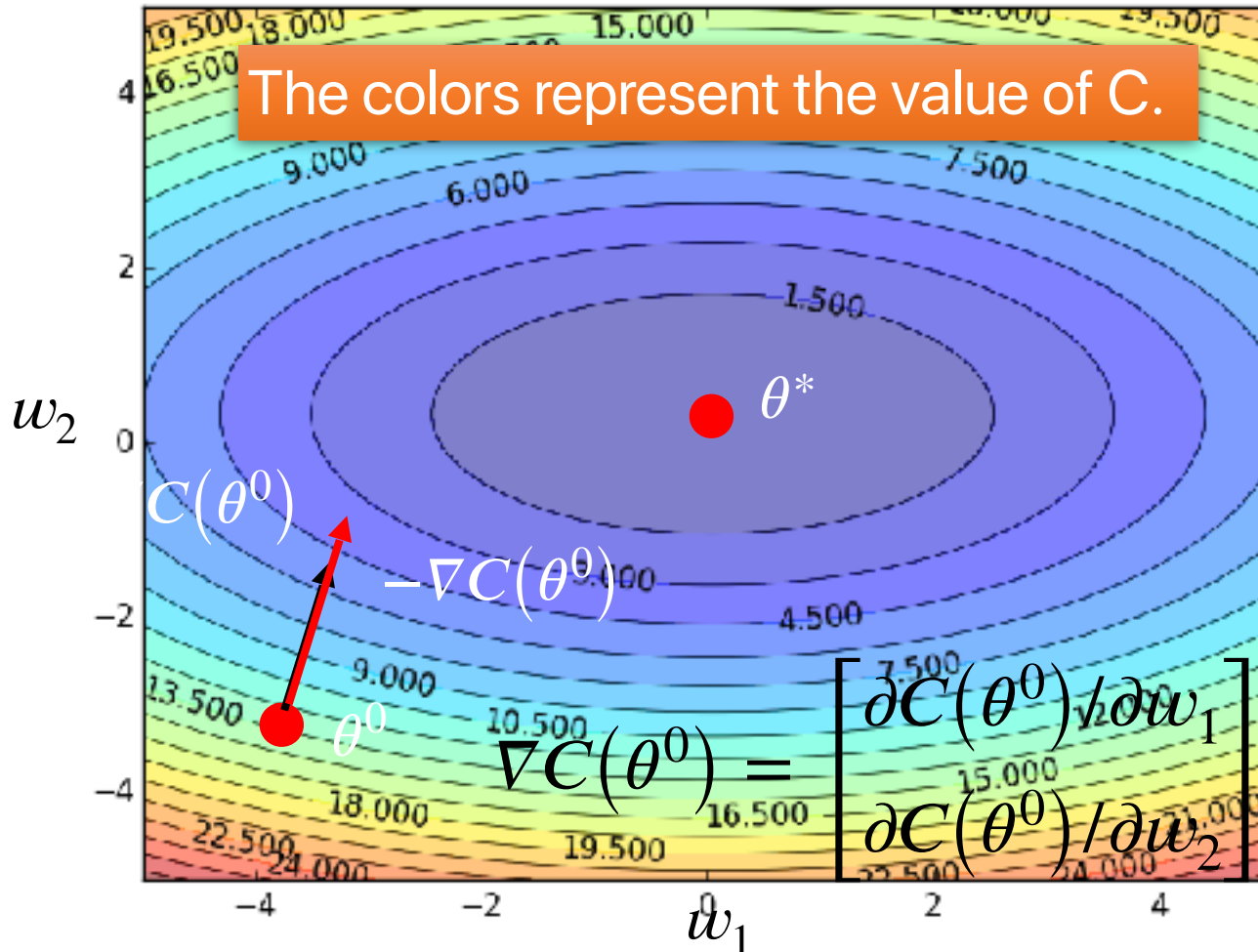Find the network parameters that minimize this value

# Gradient Descent

Assume there are only two parameters $w_1$ and $w_2$ in a network.

$$\theta = \{w_1, w_2\}$$

## Error Surface

The colors represent the value of C.

$$\theta^*$$

$$C(\theta^0)$$

$$-\nabla C(\theta^0)$$

$$\theta^0$$

$$\nabla C(\theta^0) = \begin{bmatrix} \partial C(\theta^0)/\partial w_1 \\ \partial C(\theta^0)/\partial w_2 \end{bmatrix}$$

$$w_1$$

Randomly pick a starting point

Compute the negative gradient at

$$\Rightarrow \quad -\nabla C(\theta^0)$$

Times the learning rate

$$\Rightarrow \quad -\eta \nabla C(\theta^0)$$

# Gradient Descent



Eventually, we would reach a minima .....

$-\eta\nabla C(\theta^1)$
$\theta^2 - \eta\nabla C(\theta^2)$
$-\nabla C(\theta^1)$
$-\nabla C(\theta^2)$
$\theta^1$
$\theta^0$

$w_2$

$w_1$

Randomly pick a starting point

Compute the negative gradient at

$$\implies -\nabla C(\theta^0)$$

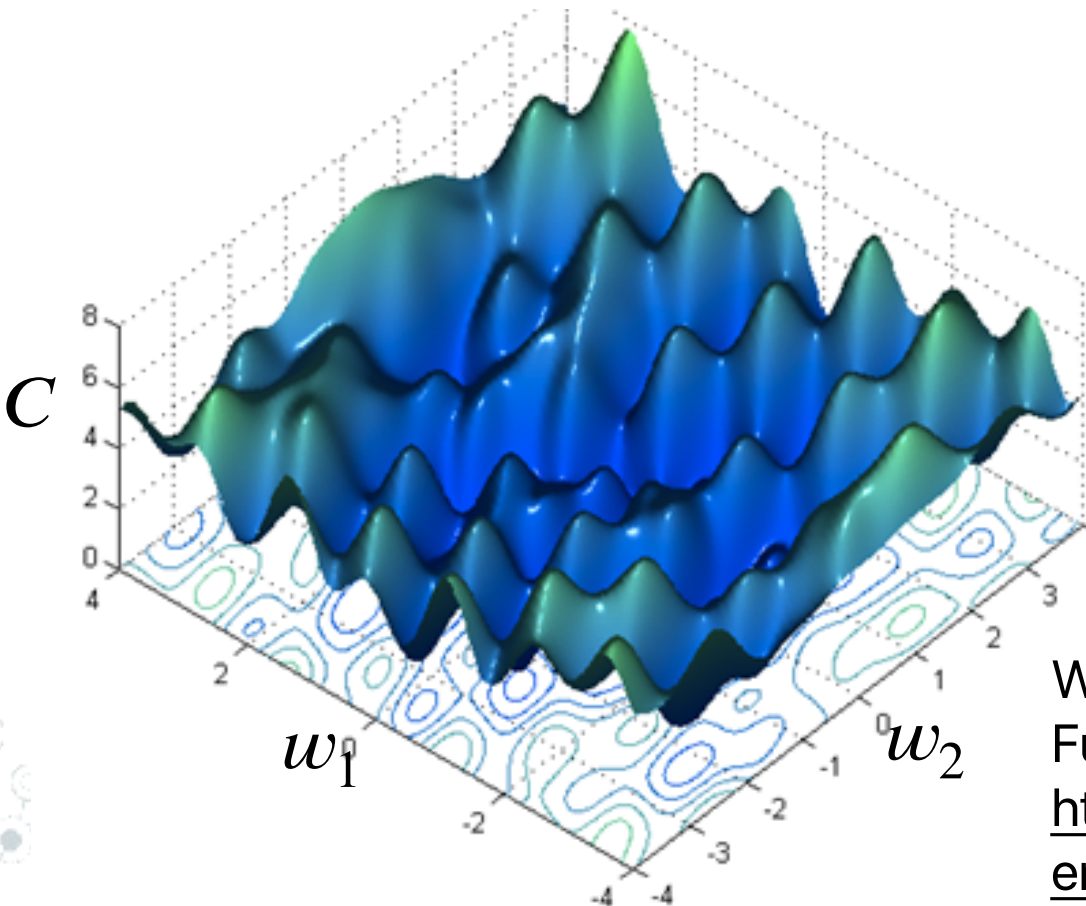Times the learning rate

$$\implies -\eta\nabla C(\theta^0)$$

# Local Minima

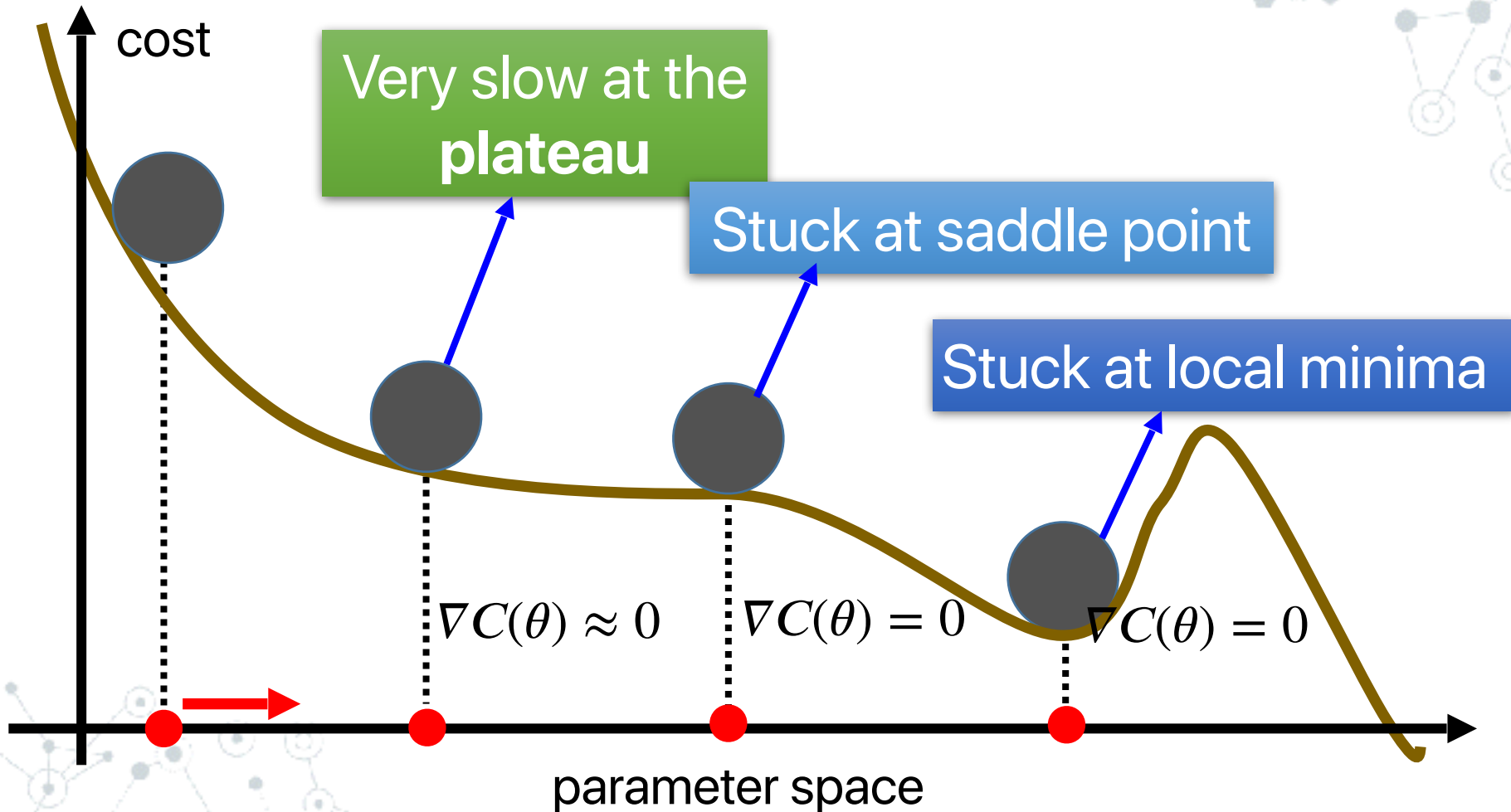◎ Gradient descent never guarantee global



Different initial point

Reach different minima, so different results
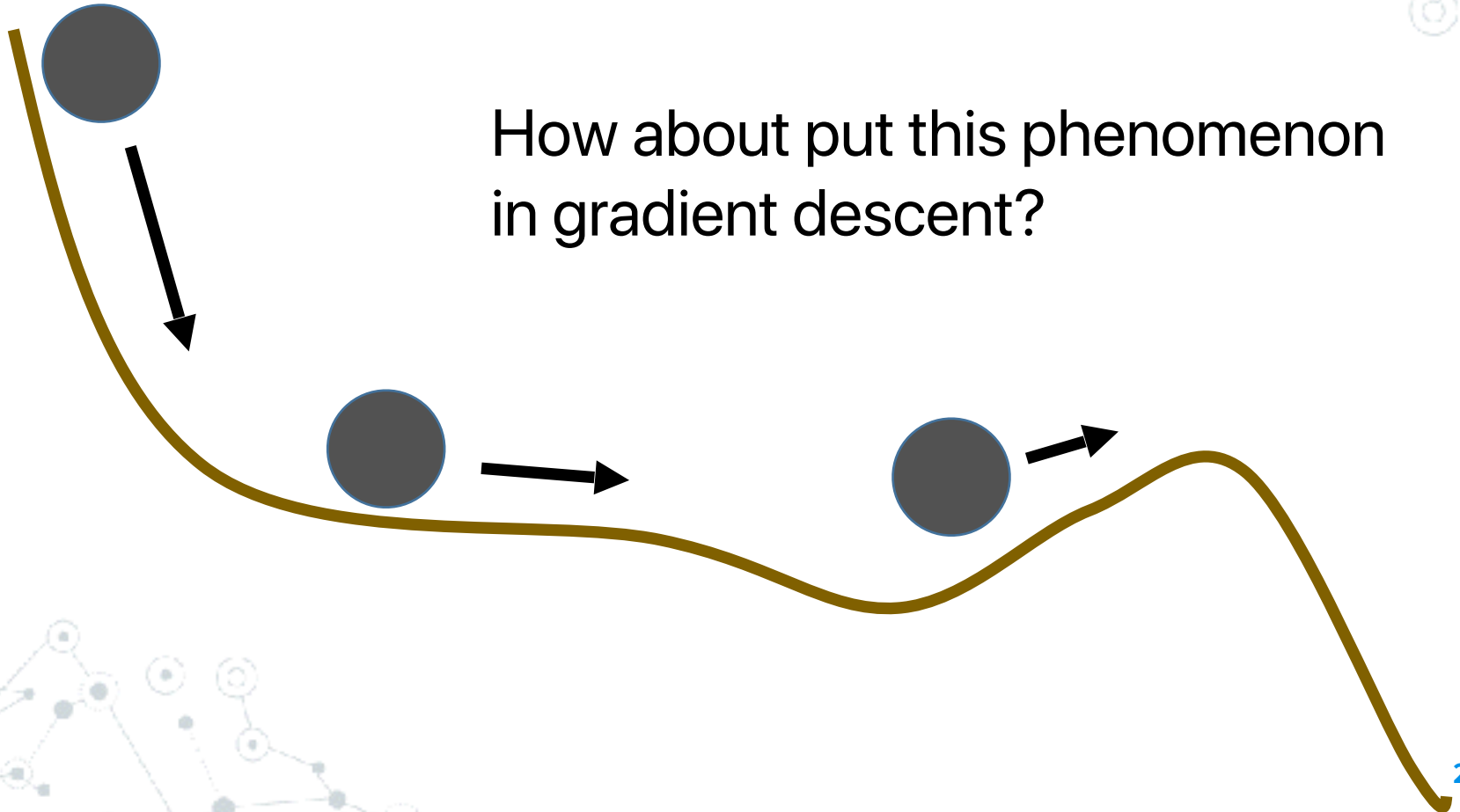
Who is Afraid of Non-Convex Loss Functions?
http://videolectures.net/eml07_lecun_wia/

# Besides local minima ......



cost

Very slow at the **plateau**
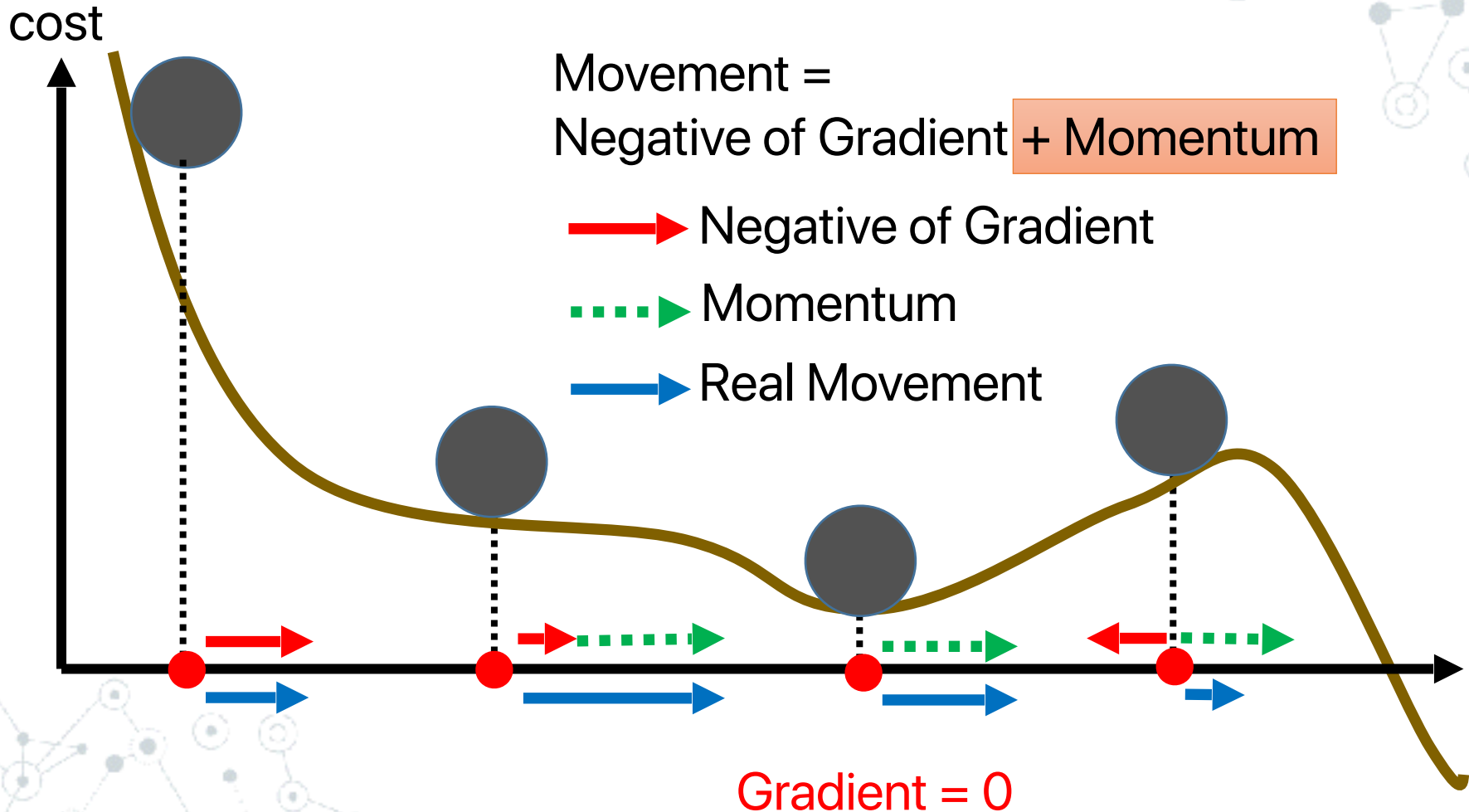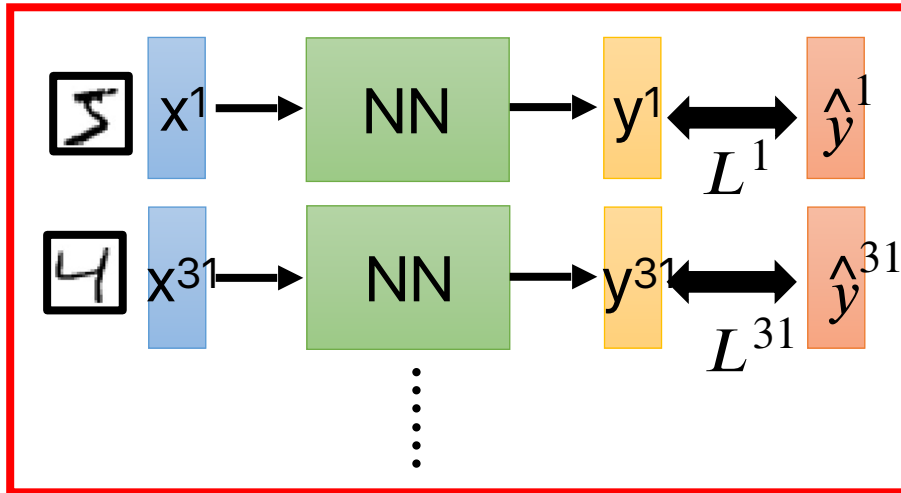
Stuck at saddle point

Stuck at local minima

$\nabla C(\theta) \approx 0$   $\nabla C(\theta) = 0$   $\nabla C(\theta) = 0$

parameter space

# In physical world ......

◎ Momentum

How about put this phenomenon in gradient descent?

# Momentum

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

cost

Movement =
Negative of Gradient + Momentum

→ Negative of Gradient

⋯▶ Momentum

→ Real Movement

Gradient = 0
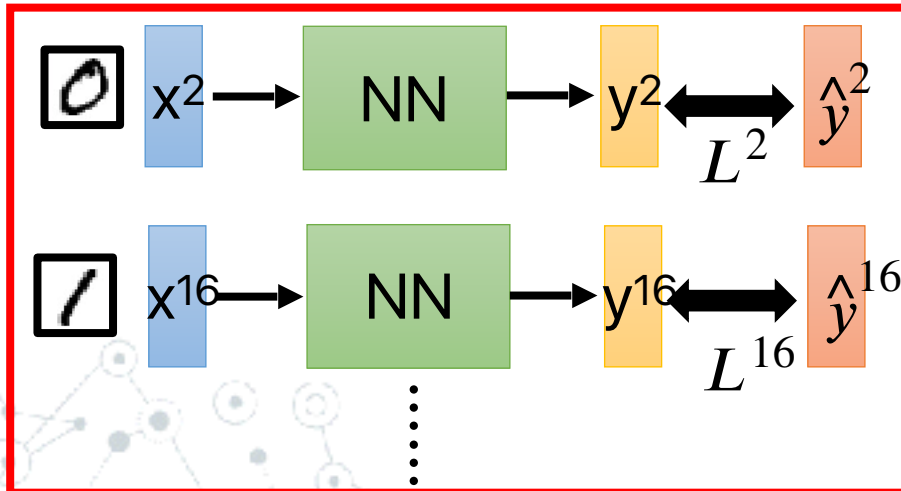
# Mini-batch



➢ Randomly initialize

➢ Pick the 1st batch

$$C = L^1 + L^{31} + \cdots$$

$$\theta^1 \leftarrow \theta^0 - \eta \nabla C(\theta^0)$$

➢ Pick the 2nd batch

$$C = L^2 + L^{16} + \cdots$$

$$\theta^2 \leftarrow \theta^1 - \eta \nabla C(\theta^1)$$

$$\vdots$$

C is different each time when we update parameters!

# Mini-batch

## Original Gradient Descent



## With Mini-batch



unstable

The colors represent the total C on all training data.
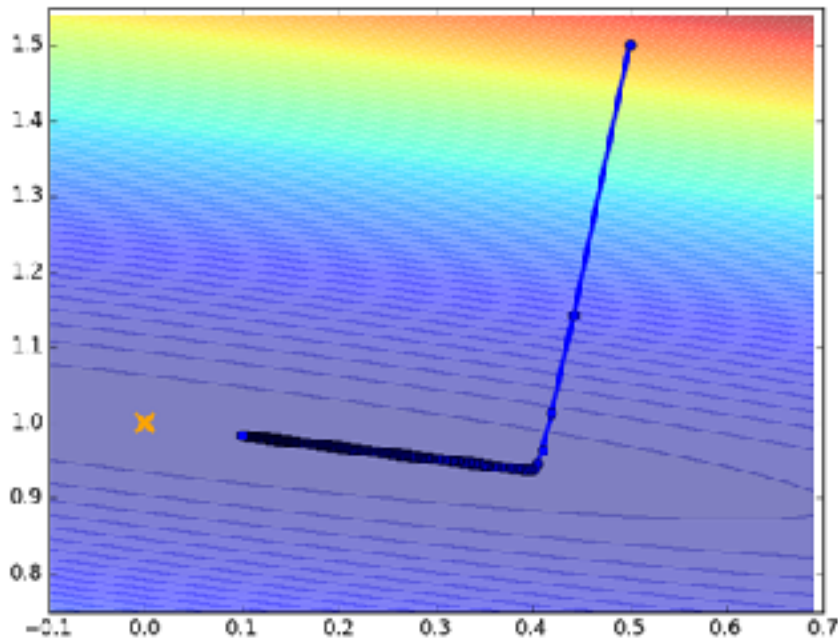
# Mini-batch

➤ Randomly initialize

➤ Pick the 1st batch

$$C = C^1 + C^{31} + \cdots$$

$$\theta^1 \leftarrow \theta^0 - \eta \nabla C(\theta^0)$$

➤ Pick the 2nd batch

$$C = C^2 + C^{16} + \cdots$$

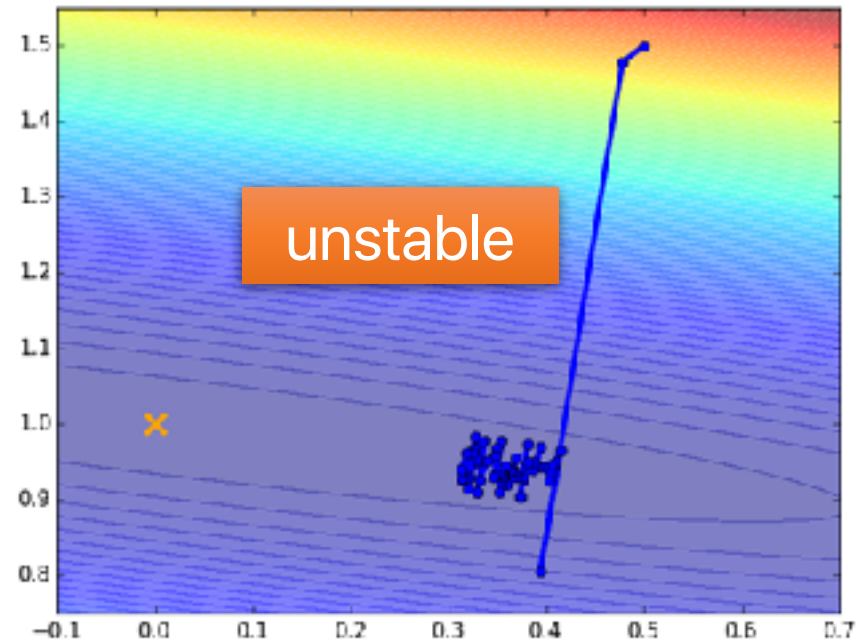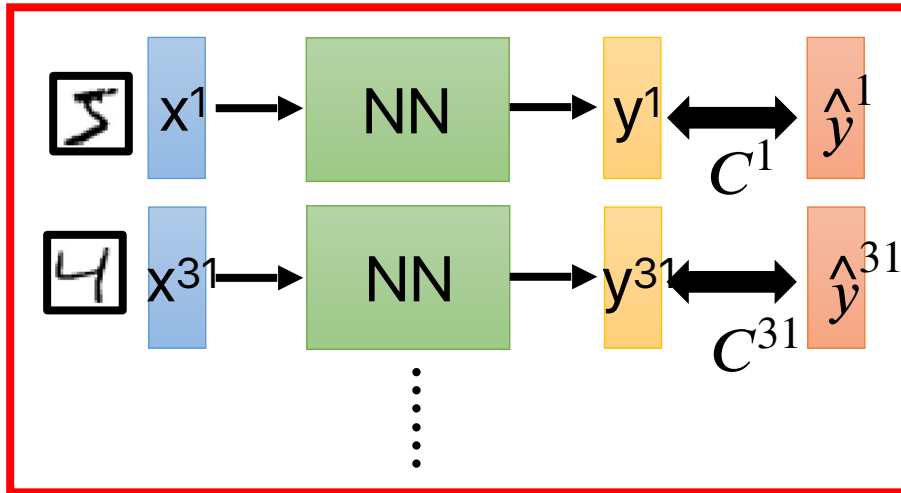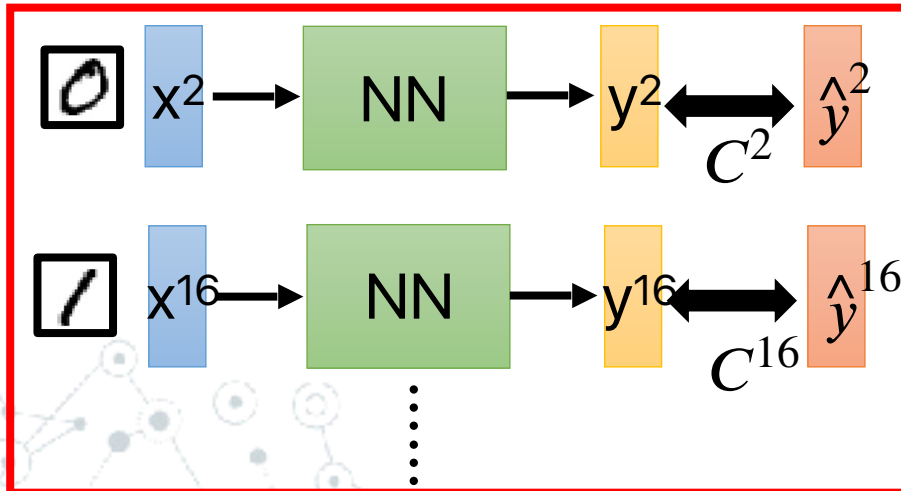$$\theta^2 \leftarrow \theta^1 - \eta \nabla C(\theta^1)$$
$$\vdots$$

➤ Until all mini-batches have been picked

one epoch

Repeat the above process

30

# Backpropagation

◎ A network can have millions of parameters.

- Backpropagation is the way to compute the gradients efficiently (not today)
- Ref: http://speech.ee.ntu.edu.tw/~tlkagk/courses/ MLDS_2015_2/Lecture/DNN%20backprop.ecm.mp4/ index.html

◎ Many toolkits can compute the gradients automatically



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

# Part II:
# 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.
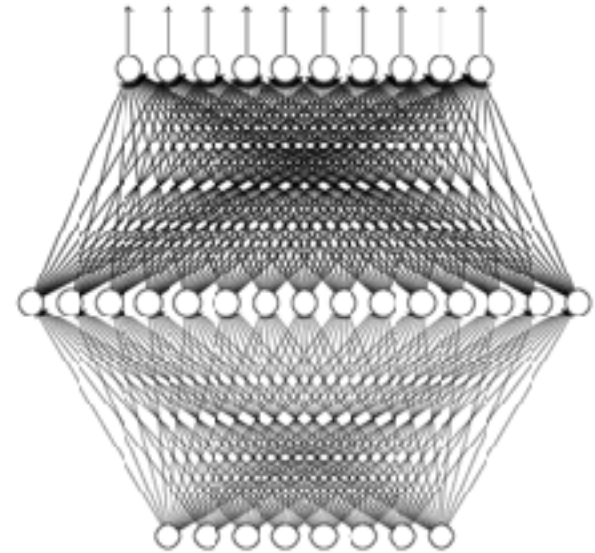
# Universality Theorem

Any continuous function f

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

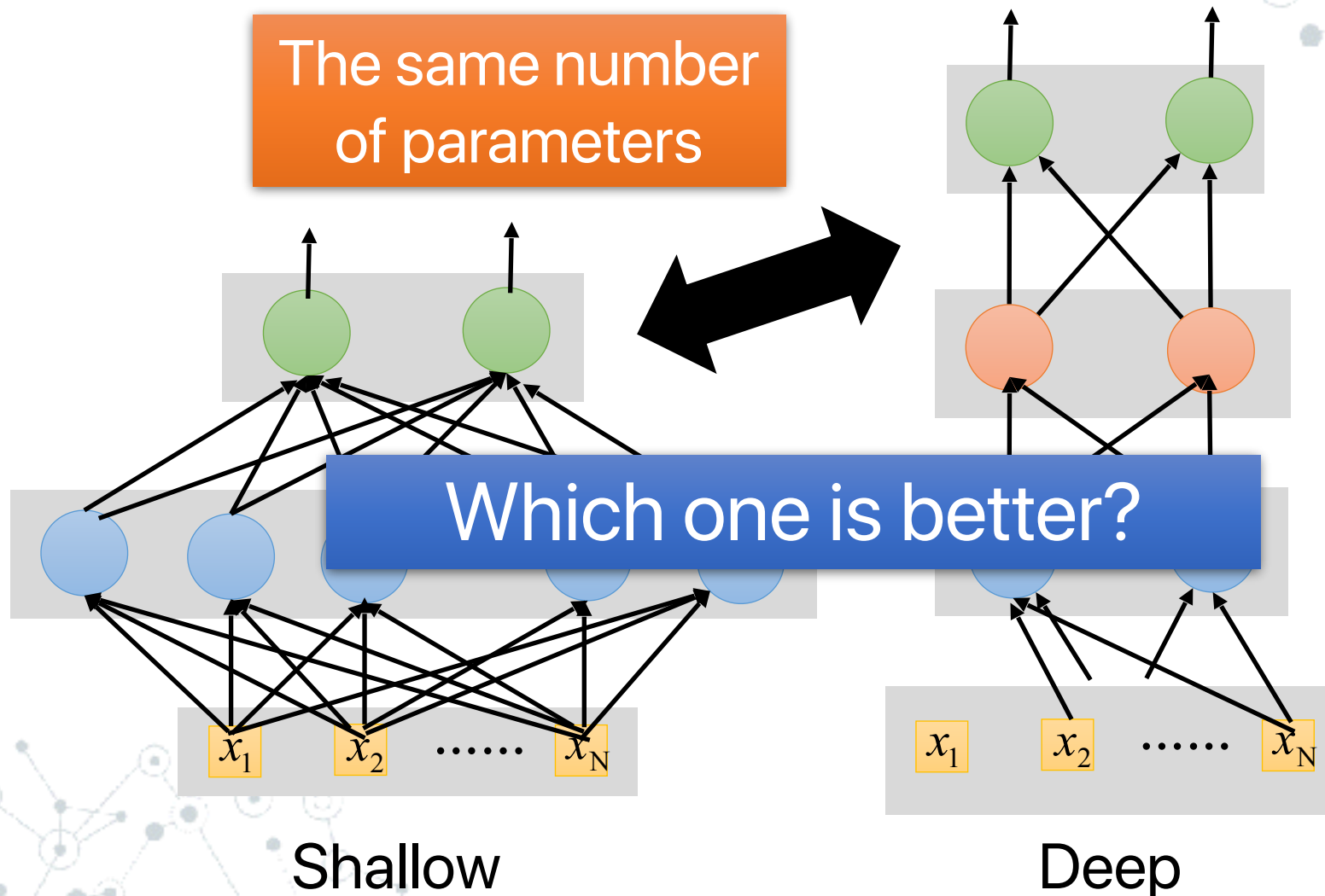Can be realized by a network with one hidden layer

(given **enough** hidden neurons)

Reference for the reason: http://neuralnetworksanddeeplearning.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?

Shallow

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 |

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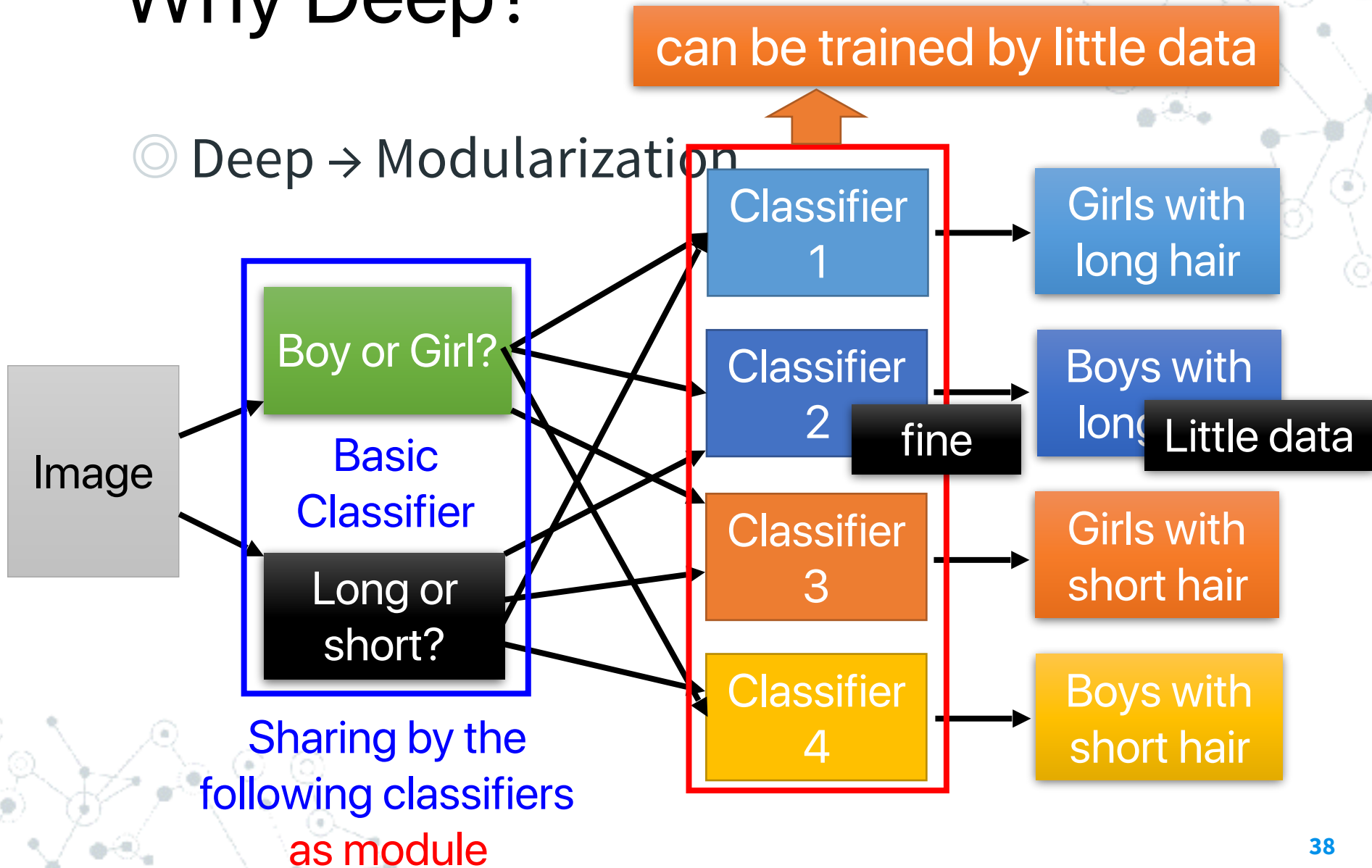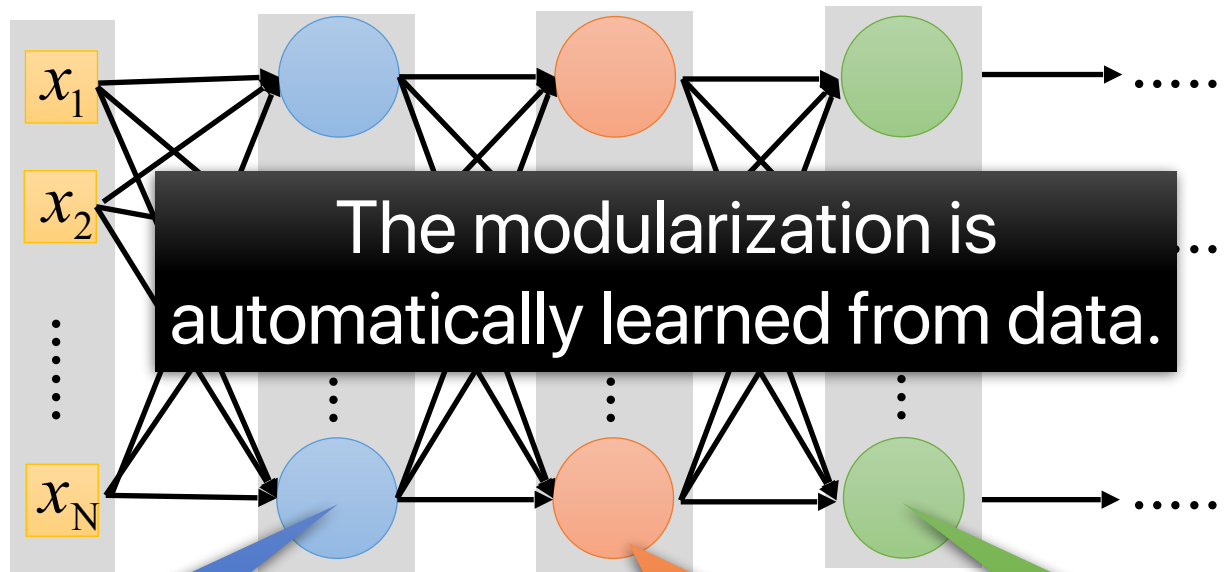
# Why Deep?

◎ Deep → Modularization



Classifier 1 → Girls with long hair

Classifier 2 → Boys with long hair — weak — Little examples

Classifier 3 → Girls with short hair

Classifier 4 → Boys with short hair

長髮女

長髮

短髮女

短髮男

Image

37

# Why Deep?

Deep Learning also works on small data set like TIMIT.

◎ Deep → Modularizat → Less training data?



The modularization is automatically learned from data.

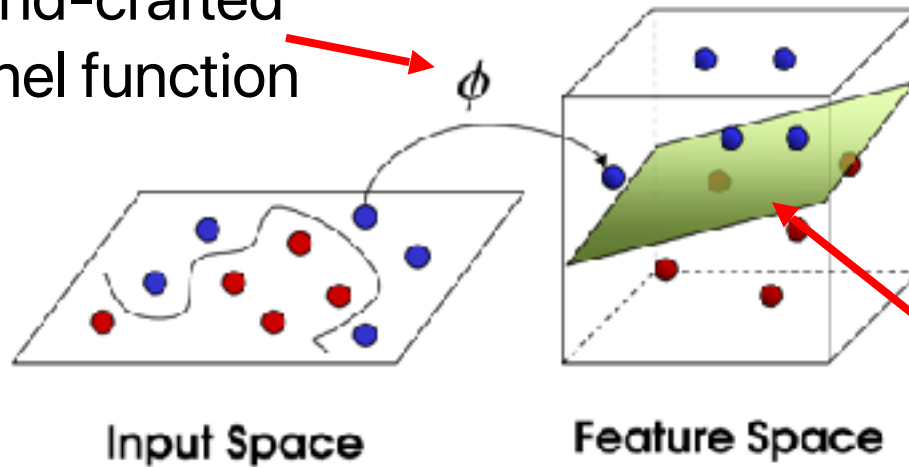The most basic classifiers

Use 1st layer as module to build classifiers

Use 2nd layer as module ......
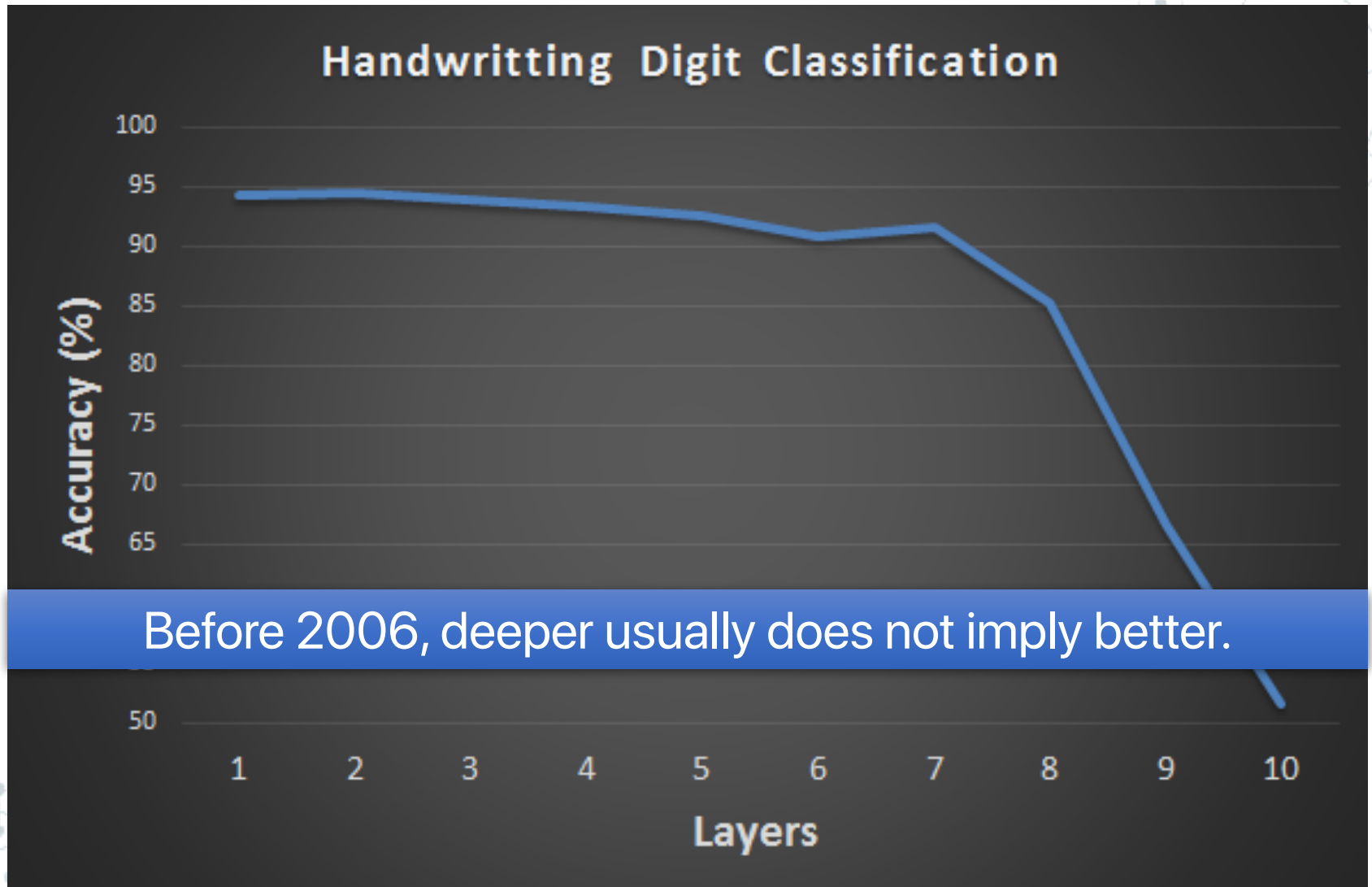
# SVM

Hand-crafted kernel function

$\phi$

Apply simple classifier

Input Space

Feature Space

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

# Deep Learning

Learnable kernel

$\phi(x)$
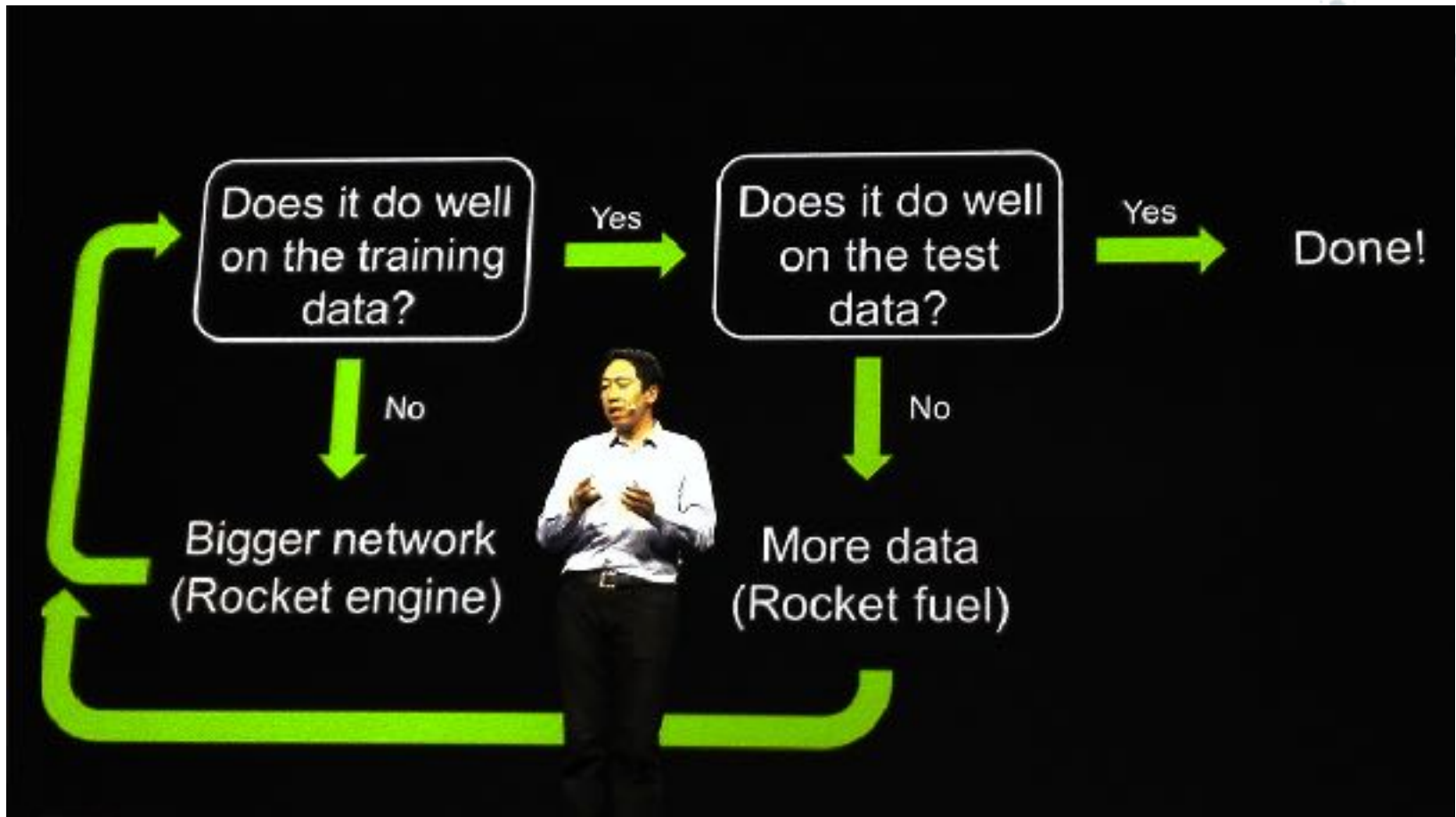
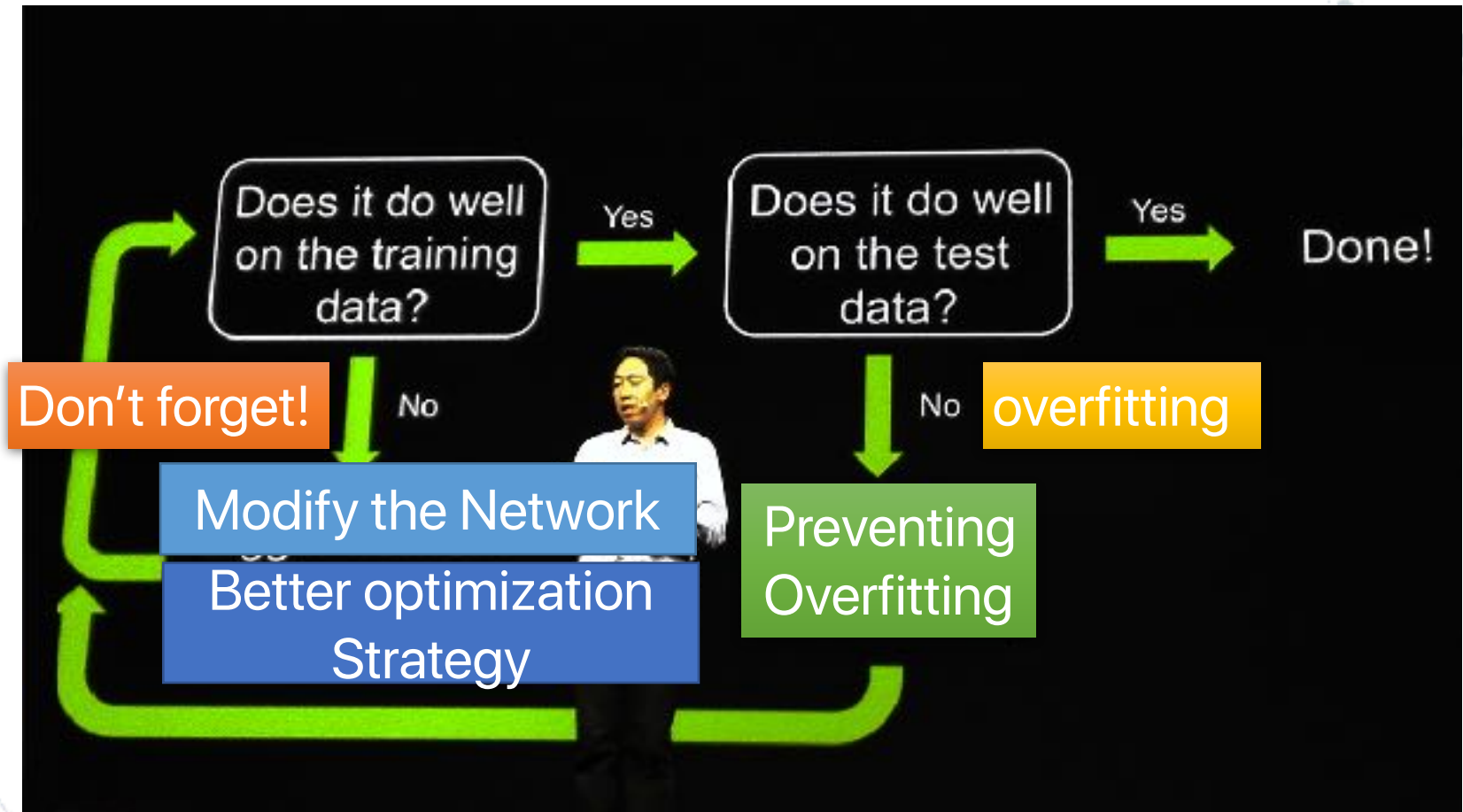simple classifier

$x$

$x_1$

$x_2$

$x_N$

...... ...... ......

$y_1$

$y_2$

$y_M$

40

# Hard to get the power of Deep ...



Before 2006, deeper usually does not imply better.

# Part III:
# Tips for Training DNN

# Recipe for Learning

# Recipe for Learning



Does it do well on the training data?

Yes → Does it do well on the test data? → Yes → Done!

No

Don't forget!

Modify the Network

Better optimization Strategy

No

overfitting

Preventing Overfitting

http://www.gizmodo.com.au/2015/04/the-basic-recipe-for-machine-learning-explained-in-a-single-powerpoint-slide/

# Recipe for Learning

**Modify the Network**
- New activation functions, for

**Better optimization Strategy**
- Adaptive learning rates

**Prevent Overfitting**
- Dropout

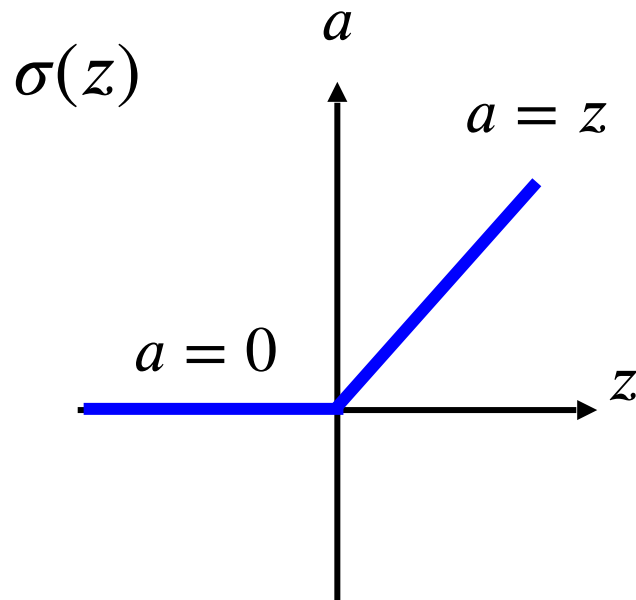Only use this approach when you already obtained good results on the training data.

**Part III:**
**Tips for Training DNN**

New Activation Function

# ReLU

◎ Rectified Linear Unit (ReLU)

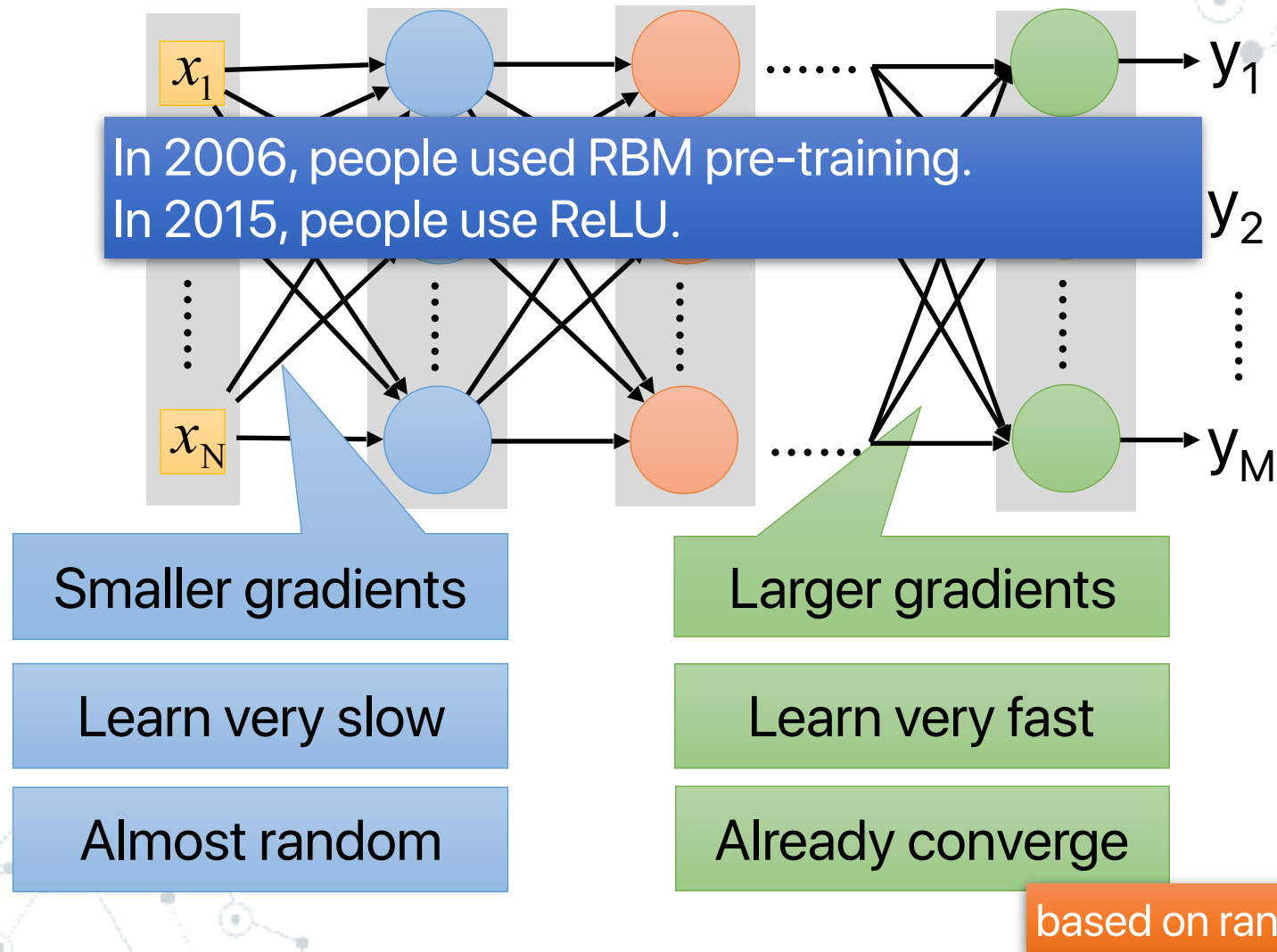$\sigma(z)$

$a$

$a = z$

$a = 0$

$z$

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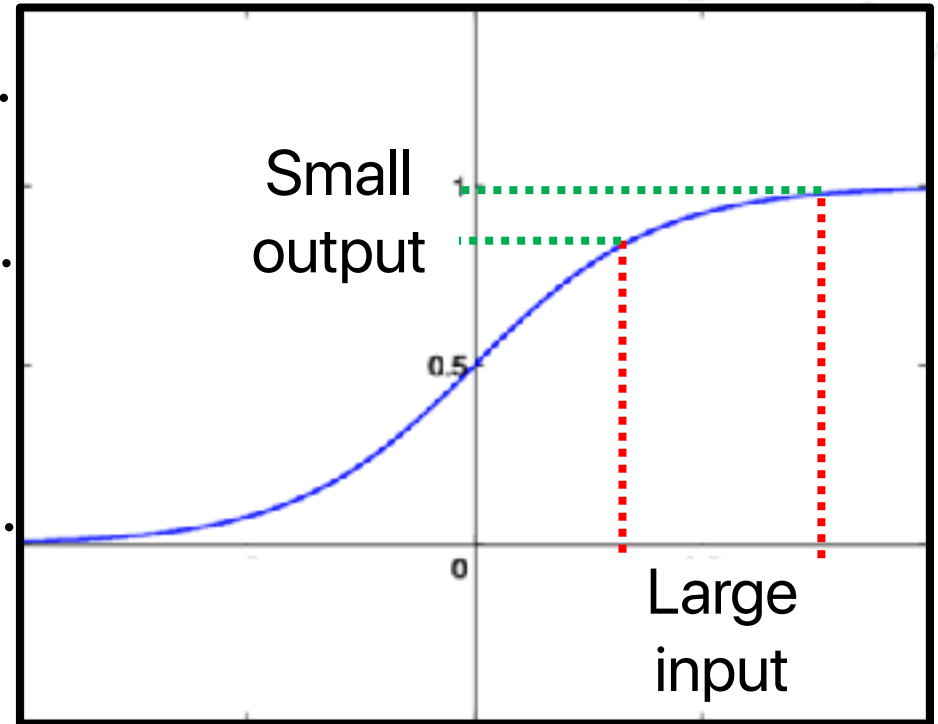[Andrew L. Maas, ICML'13]
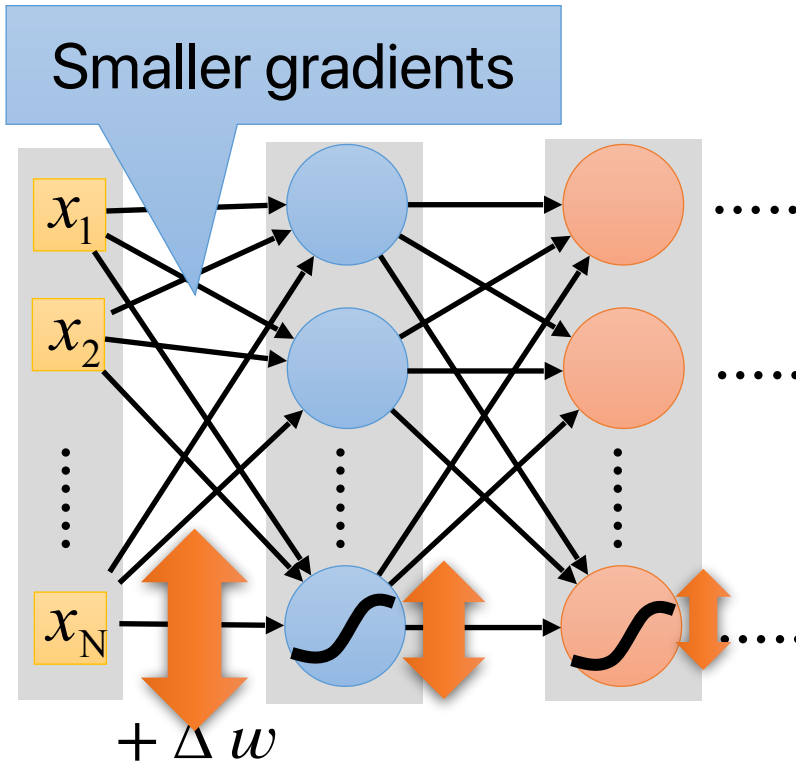[Kaiming He, arXiv'15]

## *Reason:*

1. Fast to compute

2. Biological reason

3. Infinite sigmoid with different biases

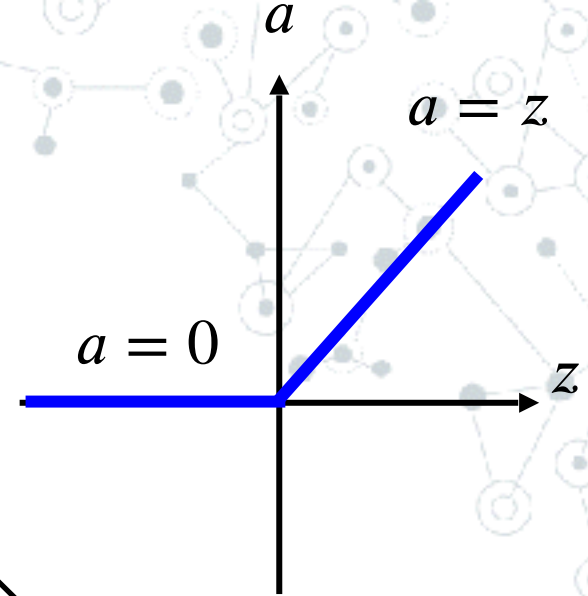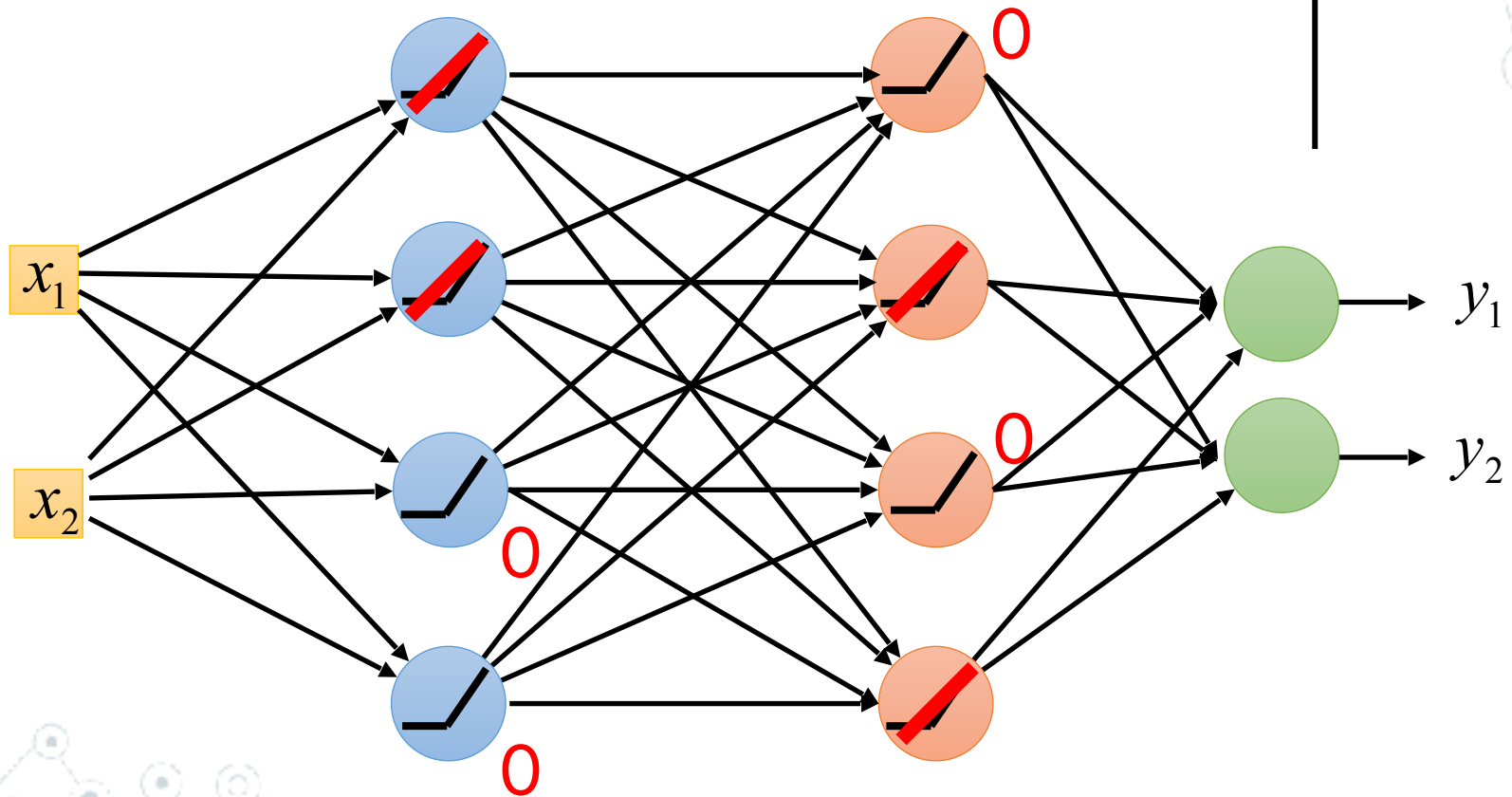4. Vanishing gradient problem

# Vanishing Gradient Problem



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

Smaller gradients

Learn very slow

Almost random

Larger gradients

Learn very fast

Already converge

based on random!?

# Vanishing Gradient Problem



Smaller gradients

$x_1$

$x_2$

$x_N$

$+ \, \Delta \, w$

Small output

Large input

Intuitive way to compute the gradient …

$$\frac{\partial C}{\partial w} = ? \, \frac{\Delta \, C}{\Delta \, w}$$

49

# ReLU

$a$

$a = z$

$a = 0$

$z$



$x_1$

$x_2$

0

0

0

0

$y_1$

$y_2$

50

# ReLU

A Thinner linear network

$a = z$

$a = 0$

$a$

$z$

$x_1$

$x_2$

Do not have smaller gradients

$y_1$

$y_2$

# Maxout

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

Input

$x_1$

$x_2$

neuron



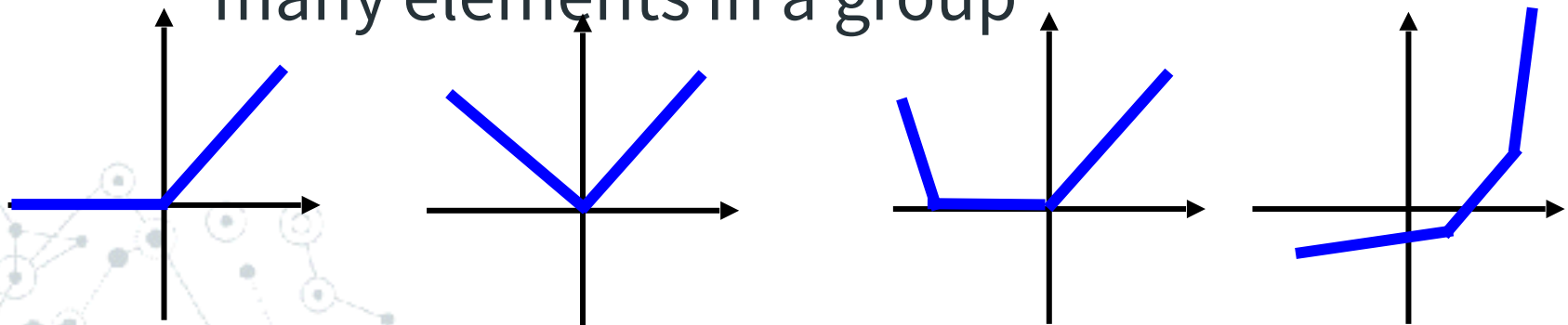You can have more than 2 elements in a group.

# Maxout

ReLU is a special cases of Maxout

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

○ Activation function in maxout network can be any piecewise linear convex function

2 elements in a group s dep 3 elements in a group many elements in a group
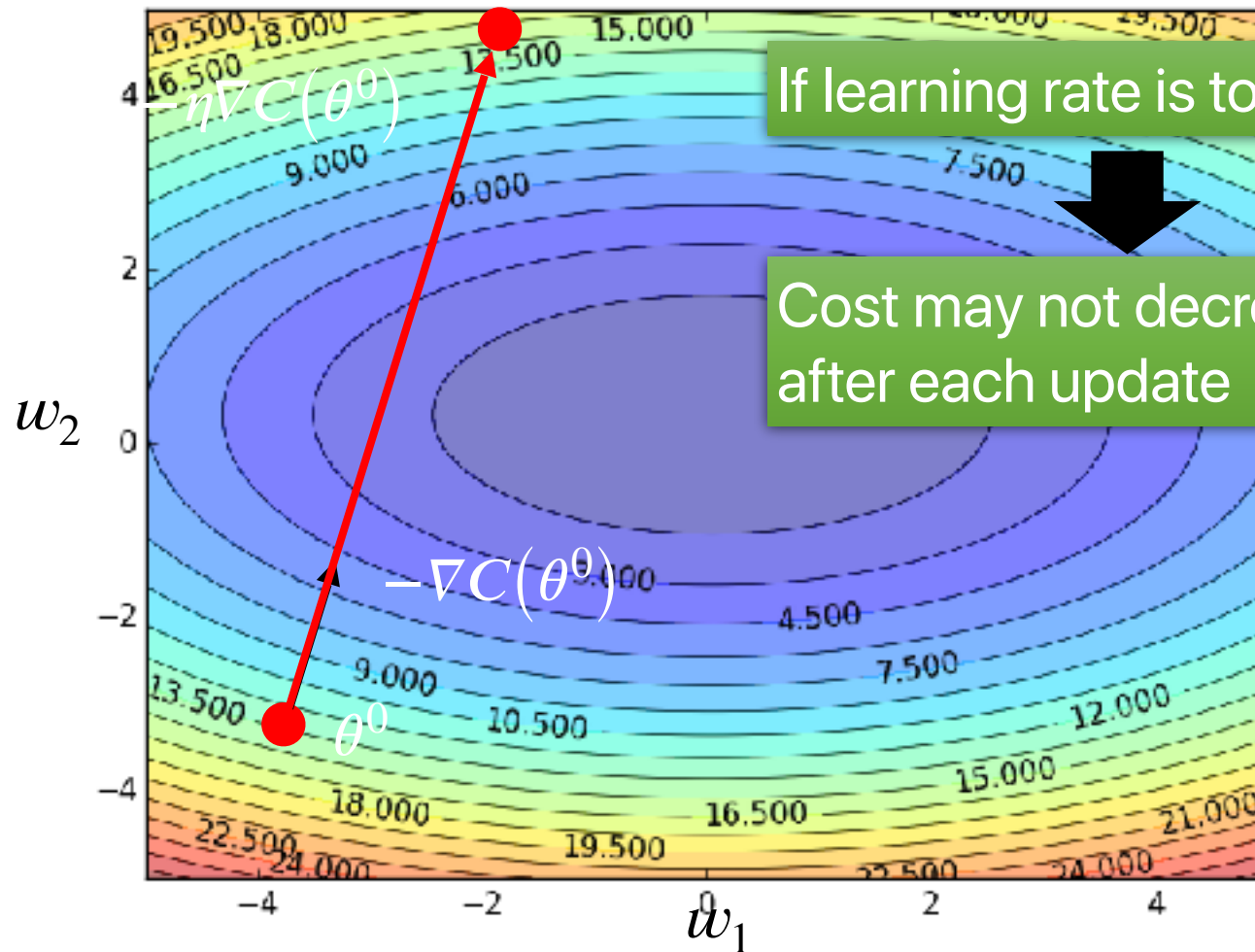
**Part III:**
**Tips for Training DNN**

Adaptive Learning Rate

# Learning Rate
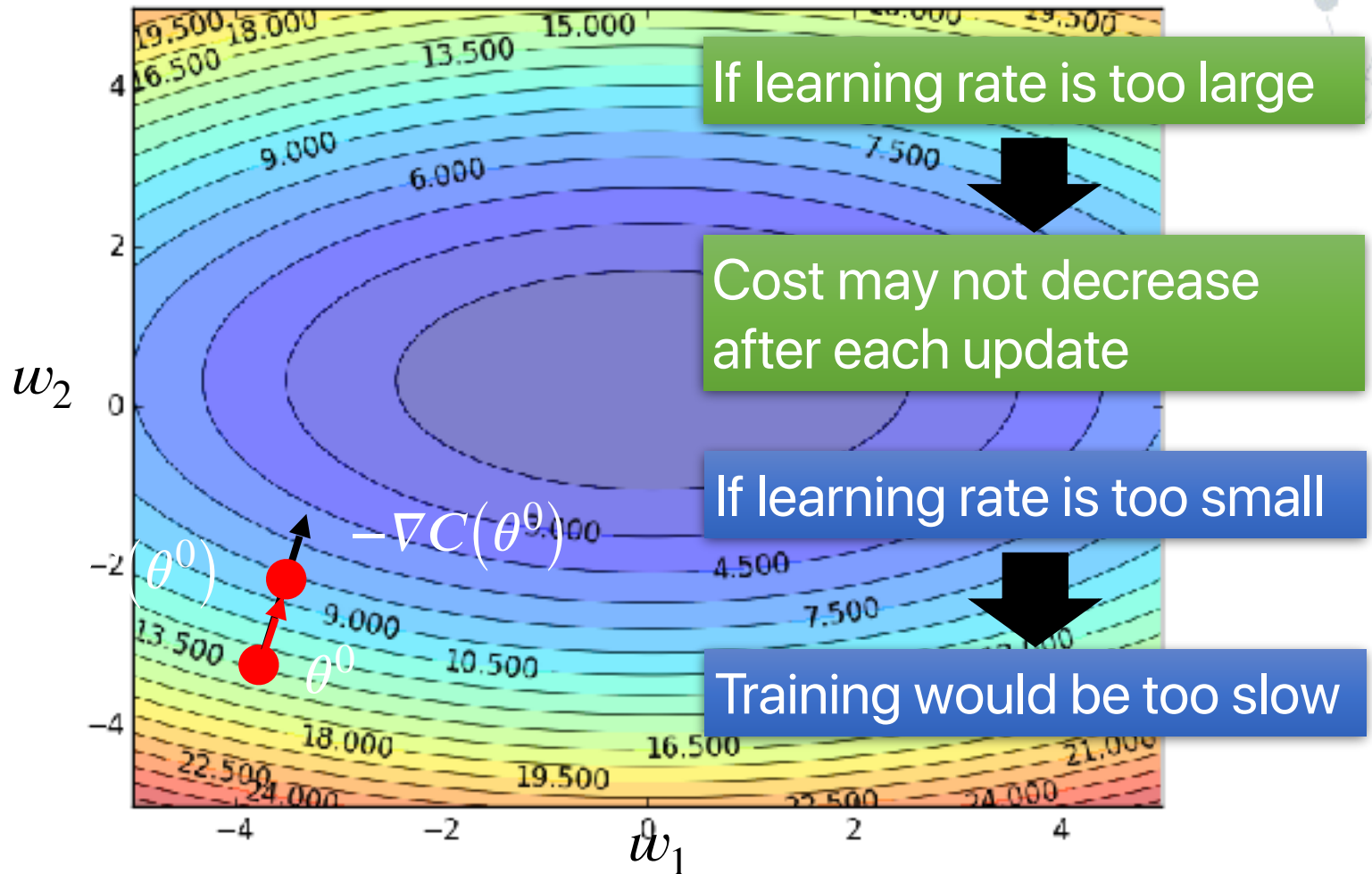
Set the learning rate η carefully

If learning rate is too large

Cost may not decrease after each update

$-\eta \nabla C(\theta^0)$

$-\nabla C(\theta^0)$

$\theta^0$

$w_2$

$w_1$

# Learning Rate

If learning rate is too large

Cost may not decrease after each update

If learning rate is too small

Training would be too slow

$-\nabla C(\theta^0)$

$(\theta^0)$

$\theta^0$

$w_2$

$w_1$

56

# Adagrad

Original Gradient Descent

$$\theta^t \leftarrow \theta^{t-1} - \eta \nabla C(\theta^{t-1})$$

Each parameter w are considered separately

$$w^{t+1} \leftarrow w^t - \eta_w g^t \qquad g^t = \frac{\partial C(\theta^t)}{\partial w}$$

Parameter dependent learning rate

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

constant

Summation of the square of the previous derivatives

# Adagrad

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

$w_1$

| 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:

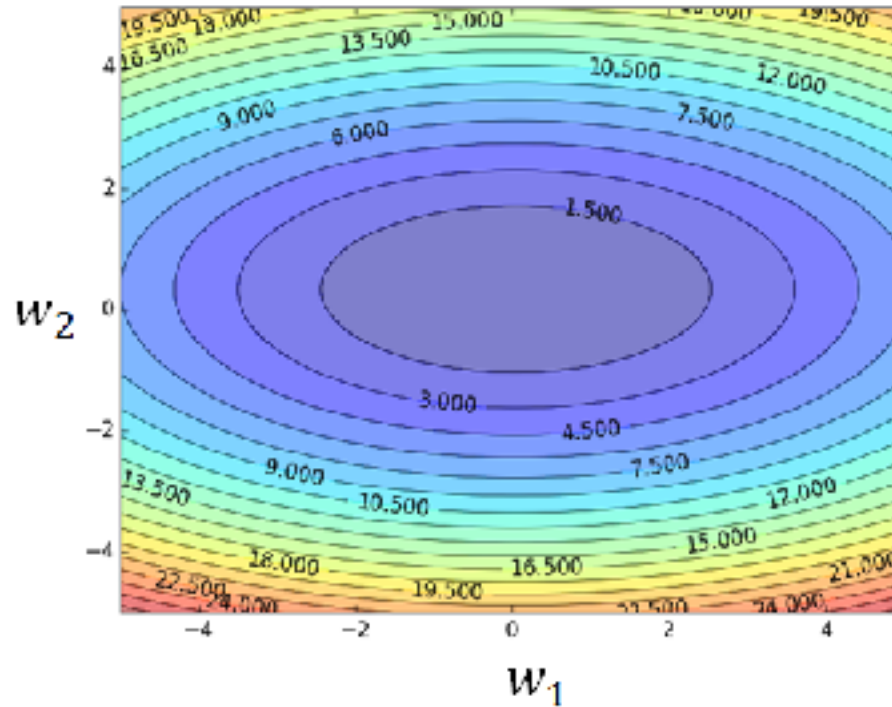$$\frac{\eta}{\sqrt{20^2}} = \frac{\eta}{20}$$

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

***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

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]

◎ Adam [Diederik P. Kingma, ICLR'15]

◎ AdaSecant [Caglar Gulcehre, arXiv'14]

◎ "No more pesky learning rates" [Tom Schaul, arXiv'12]
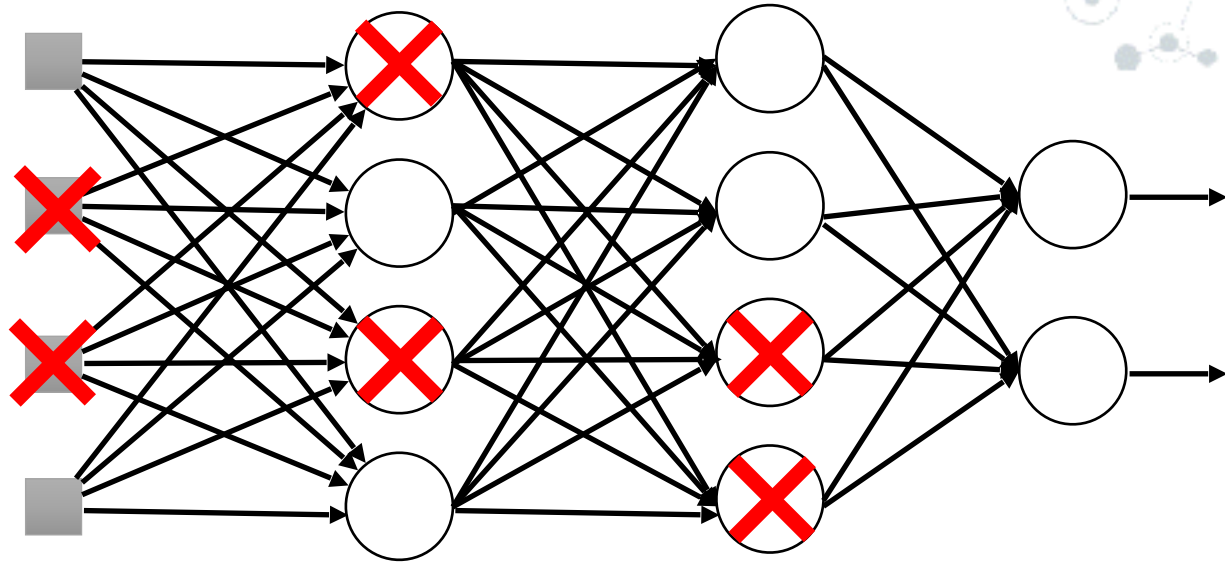
# Part III:
# Tips for Training DNN

## Dropout

# Dropout

Pick a mini-batch

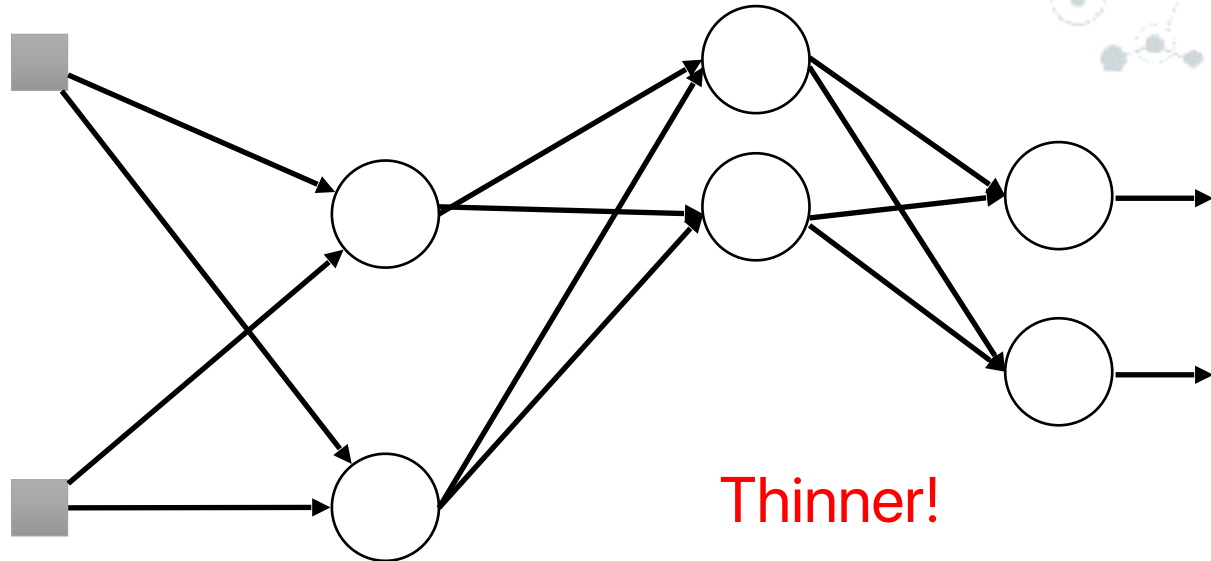$$\theta^t \leftarrow \theta^{t-1} - \eta \nabla C(\theta^{t-1})$$

**Training:**



➢ **Each time before computing the gradients**
- ● Each neuron has p% to dropout

# Dropout

Pick a mini-batch

$$\theta^t \leftarrow \theta^{t-1} - \eta \nabla C(\theta^{t-1})$$

**<u>Training:</u>**



Thinner!

➢ **Each time before computing the gradients**
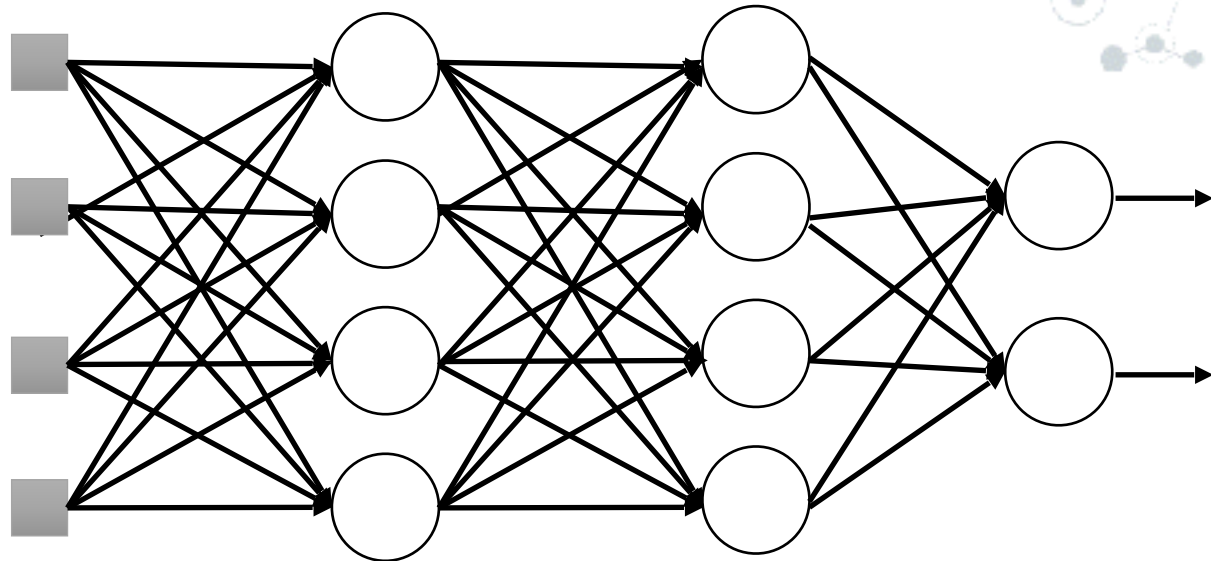   ● Each neuron has p% to dropout

   ➡ **The structure of the network is changed.**

   ● Using the new network for training

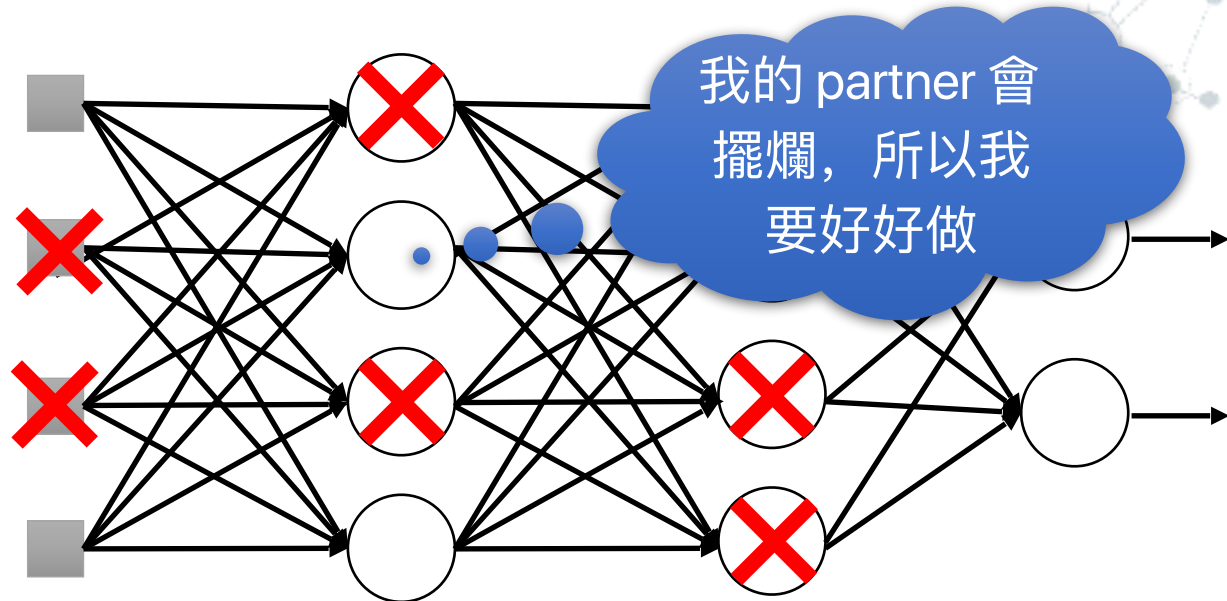For each mini-batch, we resample the dropout neurons

# 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   by training, set  for testing.
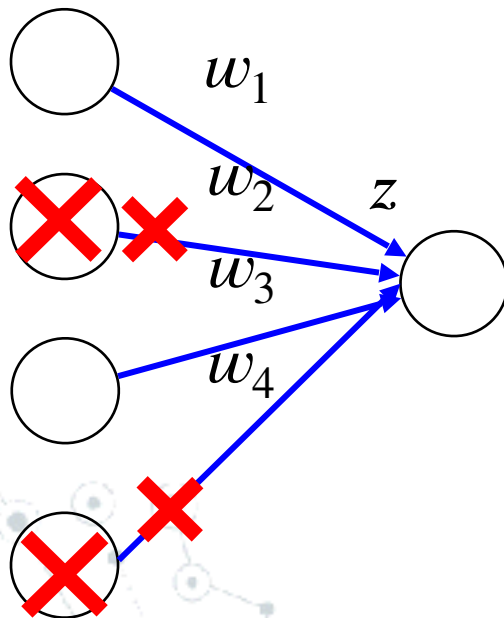
# Dropout - Intuitive Reason



➢ 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

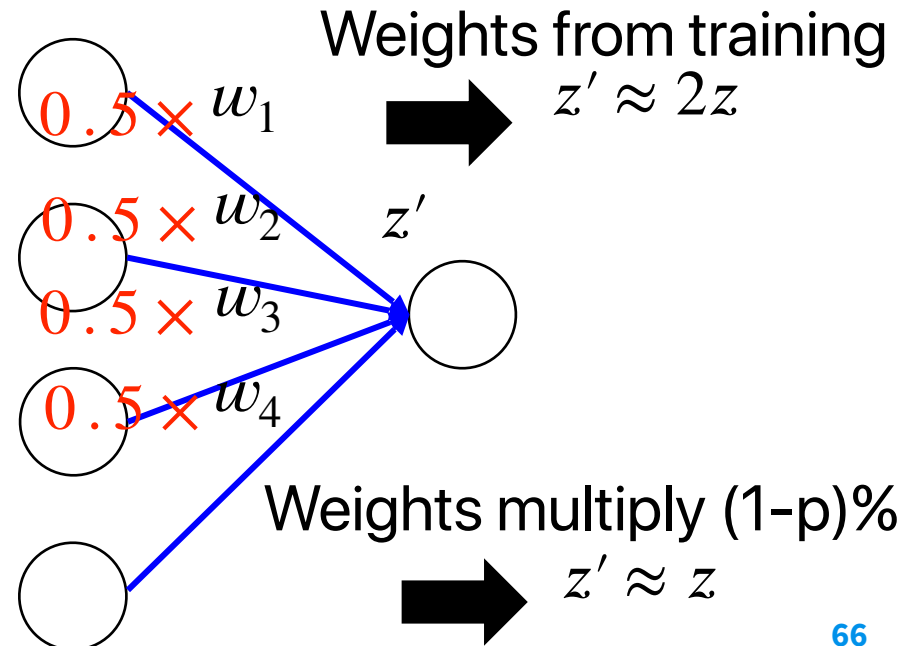◎ Why the weights should multiply (1-p)% (dropout rate) when testing?

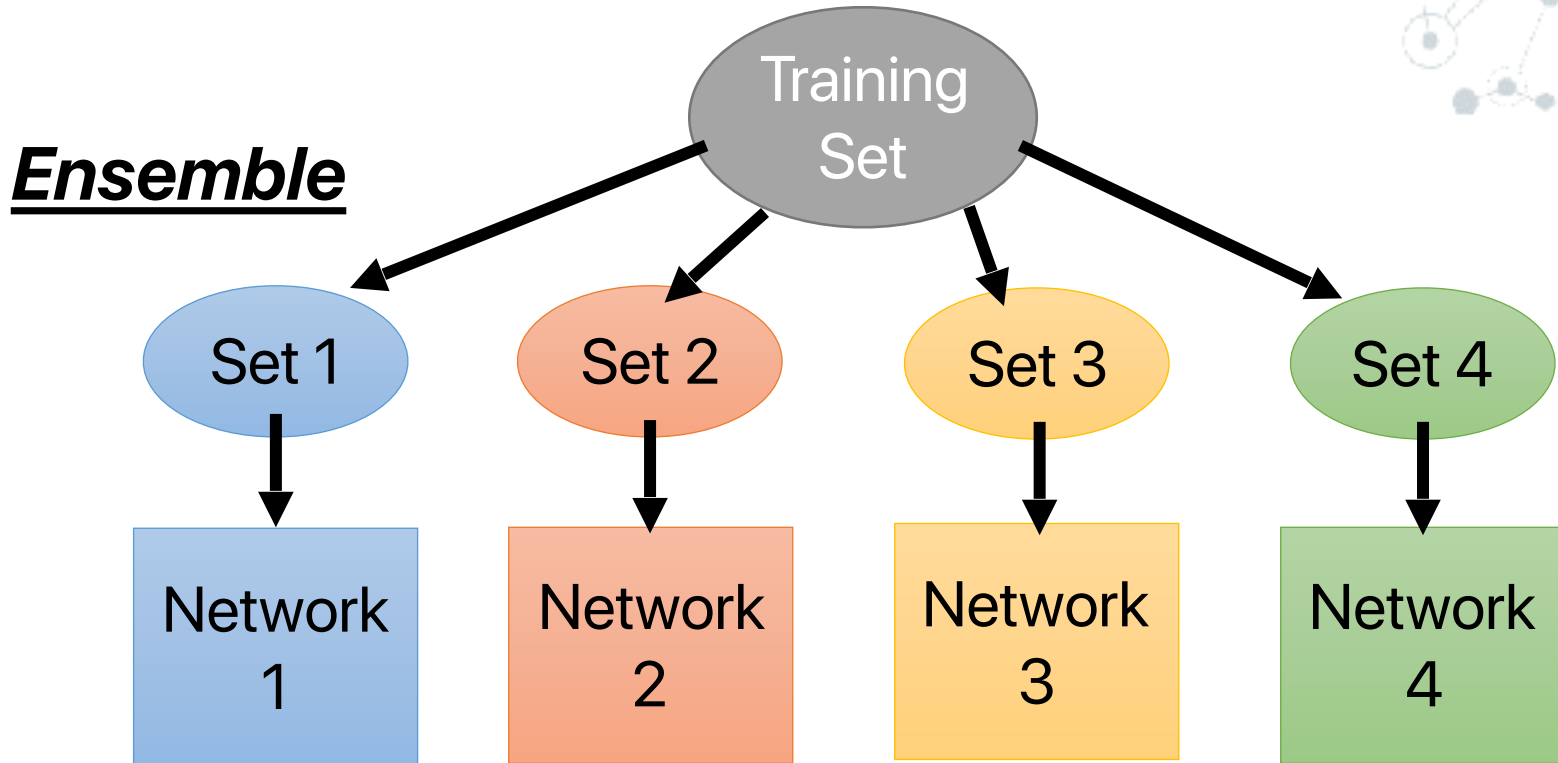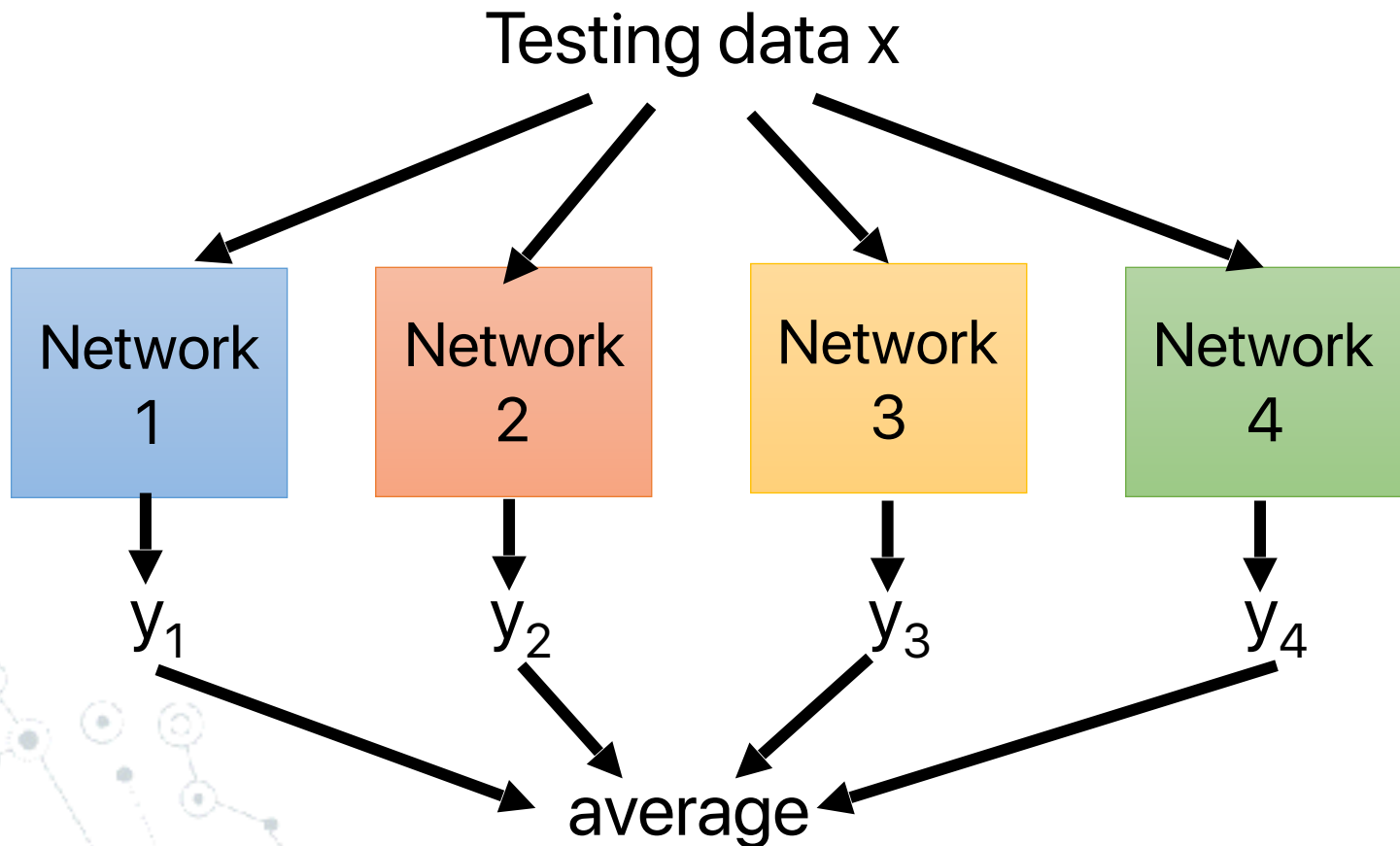**_Training of Dropout_**

Assume dropout rate is 50%



$w_1$

$w_2$  $z$

$w_3$

$w_4$

**_Testing of Dropout_**

No dropout

$0.5 \times w_1$

$0.5 \times w_2$  $z'$

$0.5 \times w_3$

$0.5 \times w_4$

Weights from training
$z' \approx 2z$

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

66

# Dropout is a kind of ensemble.



**_Ensemble_**

Training Set

Set 1 → Network 1

Set 2 → Network 2

Set 3 → Network 3

Set 4 → Network 4

Train a bunch of networks with different structures

# Dropout is a kind of ensemble.

***Ensemble***

Testing data x



| Network 1 | Network 2 | Network 3 | Network 4 |

$y_1$  $y_2$  $y_3$  $y_4$

average

# Dropout is a kind of ensemble.

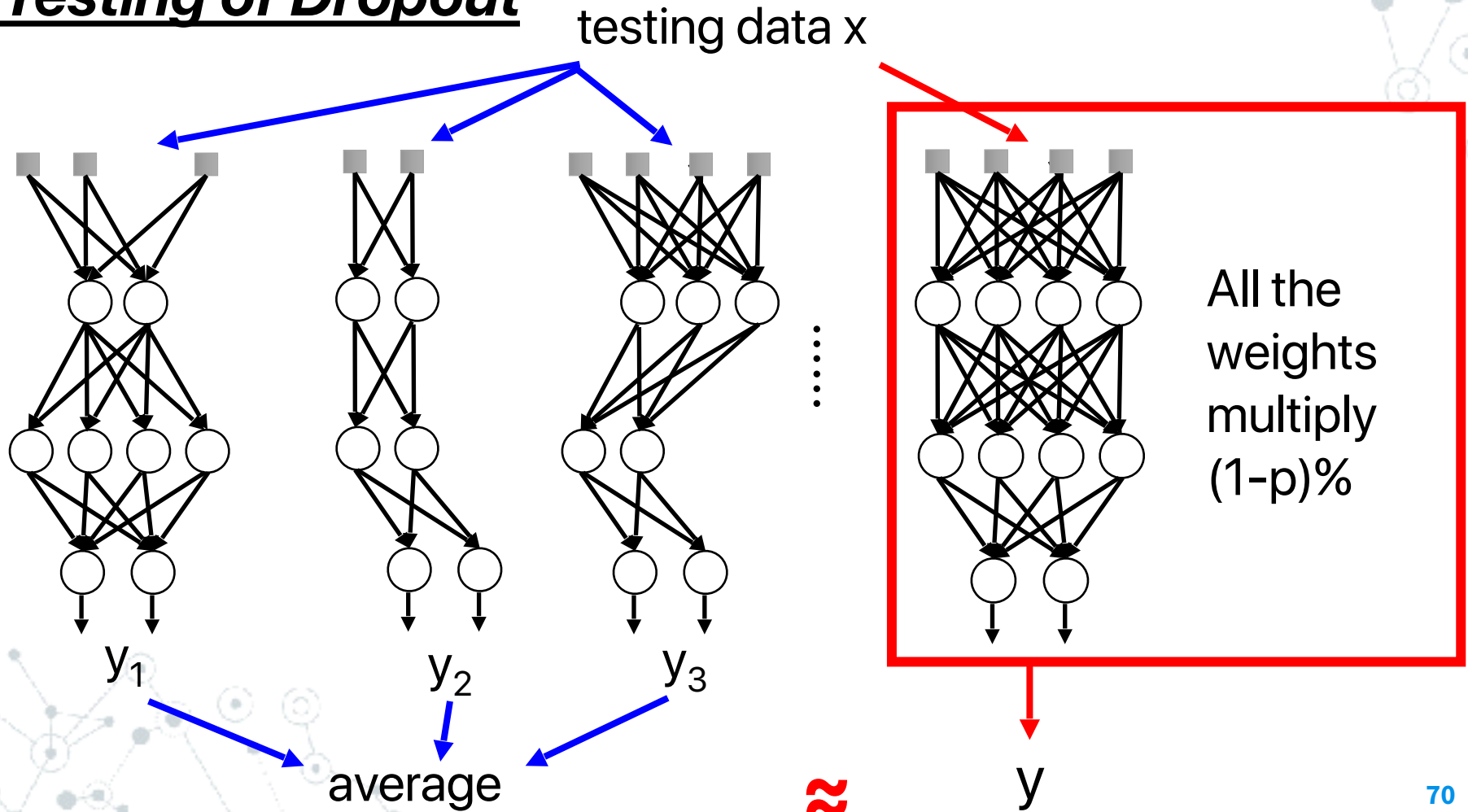| minibatch 1 | minibatch 2 | minibatch 3 | minibatch 4 |
|---|---|---|---|

***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



$y_1$     $y_2$     $y_3$

average

All the weights multiply $(1-p)\%$

$\approx$     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

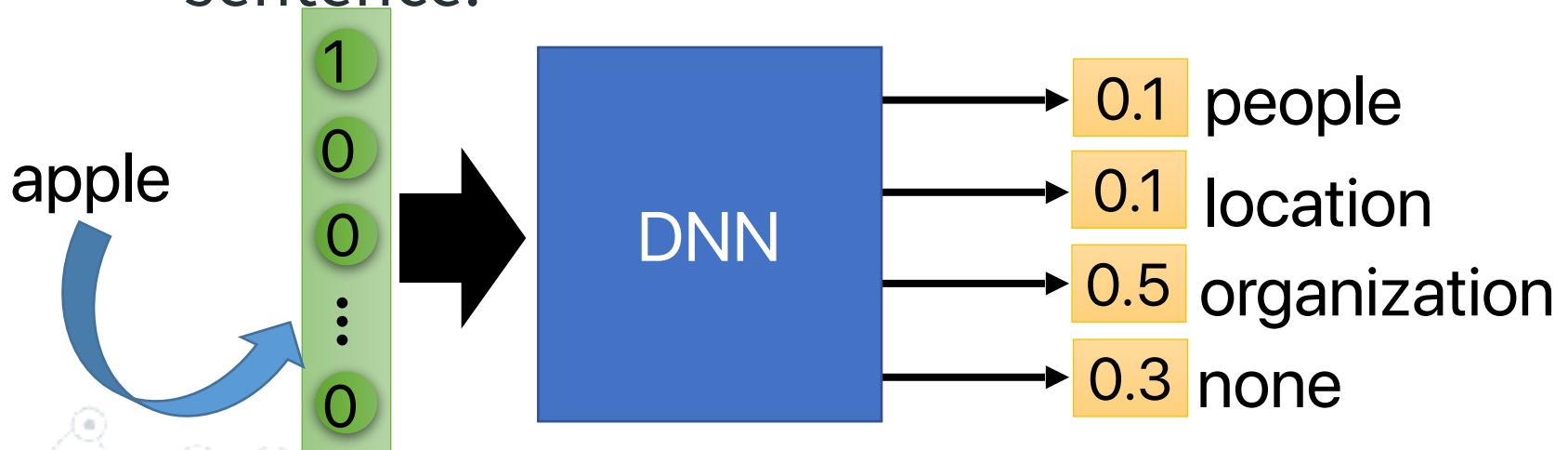- Standout [J. Ba, NISP'13]
  - Each neural has different dropout rate

# Part IV:
# Neural Network with Memory

# Neural Network needs Memory
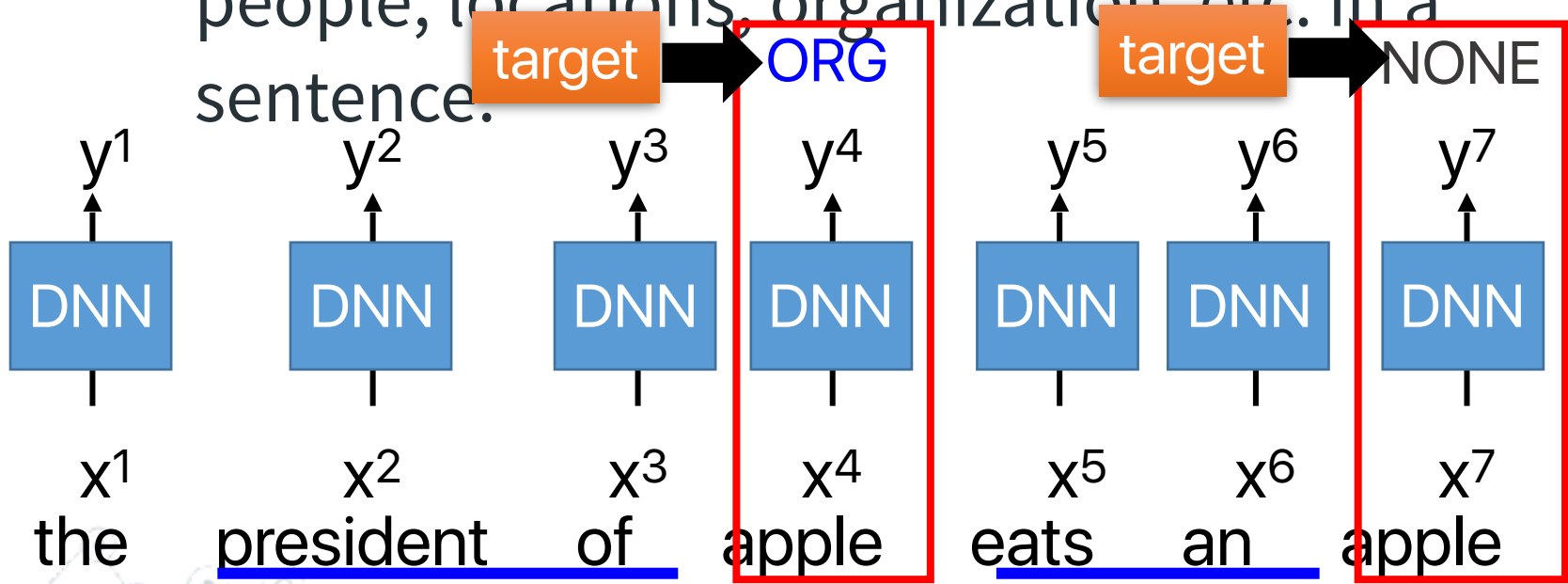
◎ Name Entity Recognition

○ Detecting named entities like name of people, locations, organization, etc. in a sentence.

apple

1
0
0
⋮
0

DNN

0.1 people

0.1 location

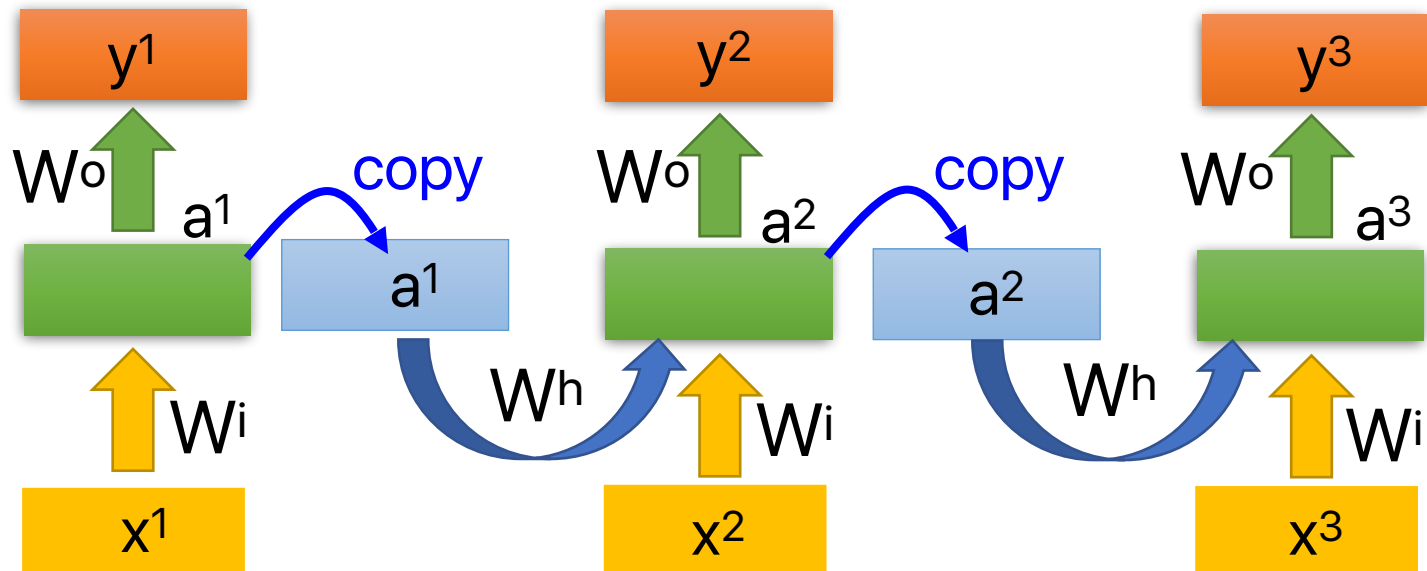0.5 organization

0.3 none

# Neural Network needs Memory

◎ Name Entity Recognition

○ Detecting named entities like name of people, locations, organization, etc. in a sentence.



$y^1$   $y^2$   $y^3$   $y^4$   $y^5$   $y^6$   $y^7$

target → ORG    target → NONE

DNN DNN DNN DNN DNN DNN DNN

$x^1$   $x^2$   $x^3$   $x^4$   $x^5$   $x^6$   $x^7$

the   president   of   apple   eats   an   apple

**DNN needs memory!**

# Recurrent Neural Network (RNN)

The output of hidden layer are stored in the memory.

copy

Memory can be considered as another input.

$y_1$

$y_2$

$a_1$

$a_2$

$x_1$

$x_2$

# RNN



The same network is used again and again.
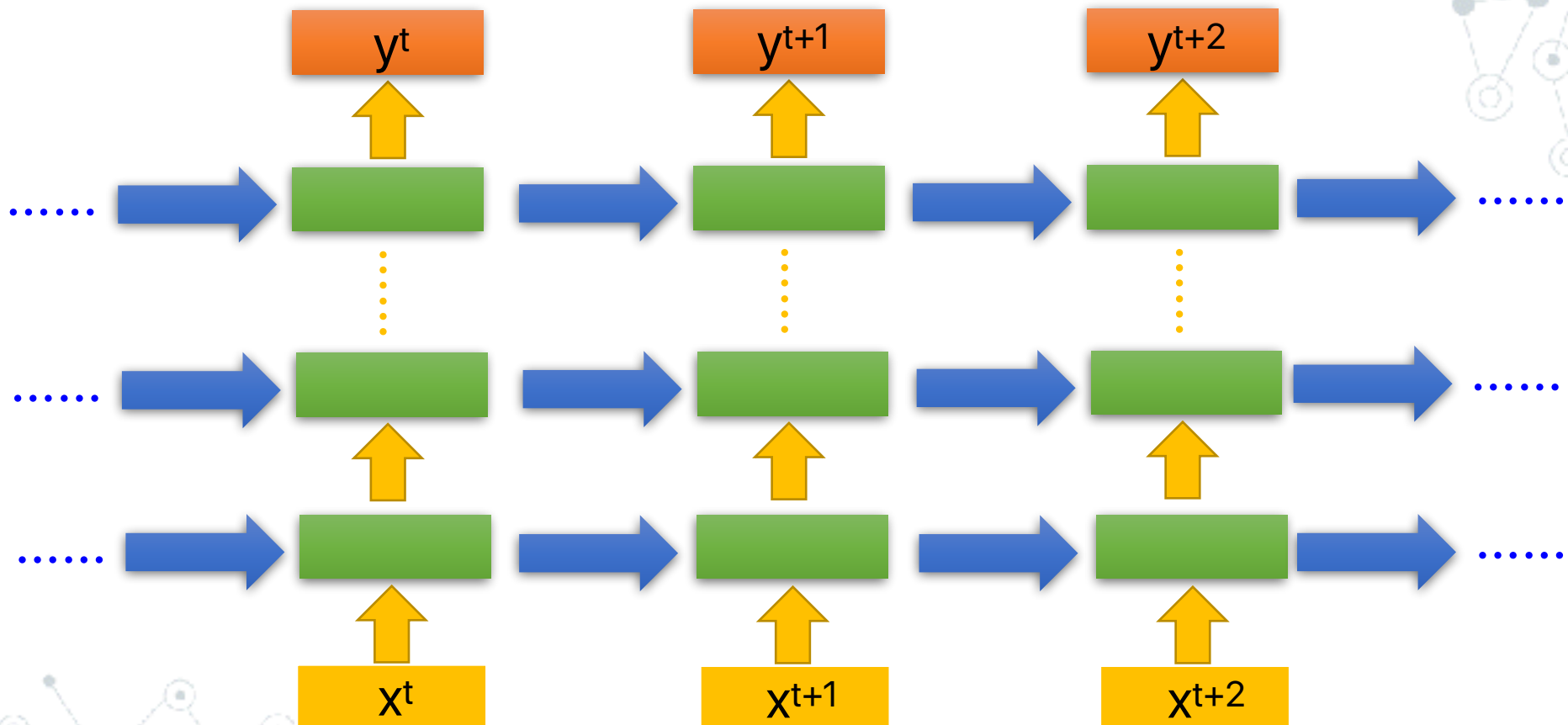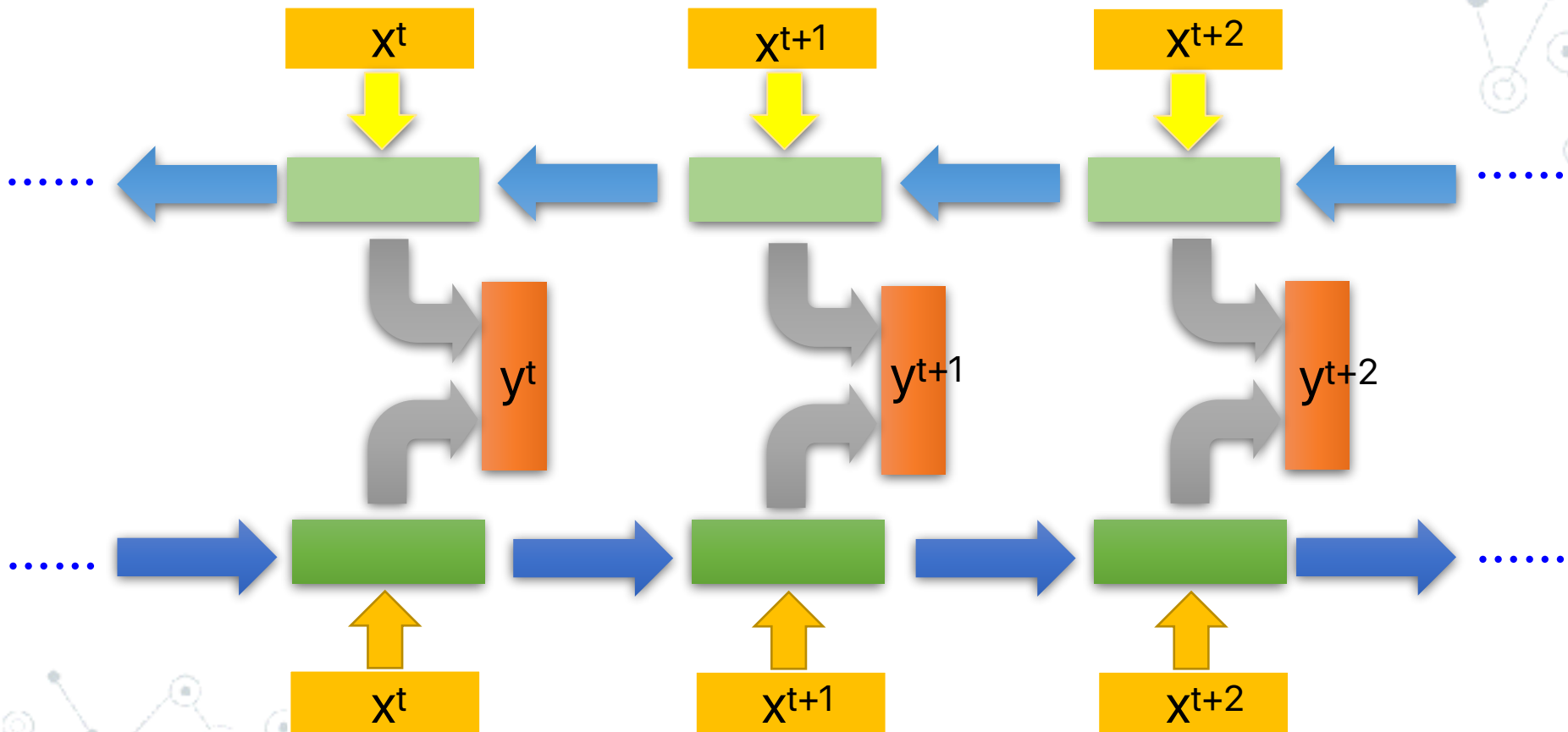Output $y^i$ depends on $x^1, x^2, \ldots x^i$

# RNN

How to train?

Find the network parameters to minimize the total cost:

Backpropagation through time (BPTT)

# Of course it can be deep ...

# Bidirectional RNN

# Many to Many (Output is shorter)

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

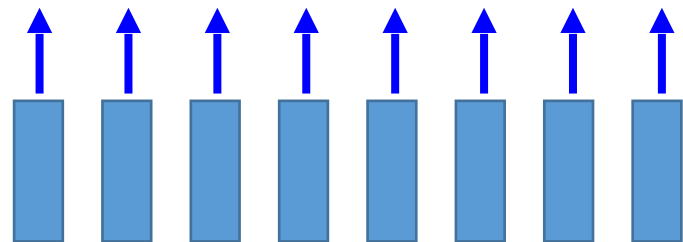　　○ E.g. ***Speech Recognition***

Output: "好棒" (character sequence)

Trimming

Problem?

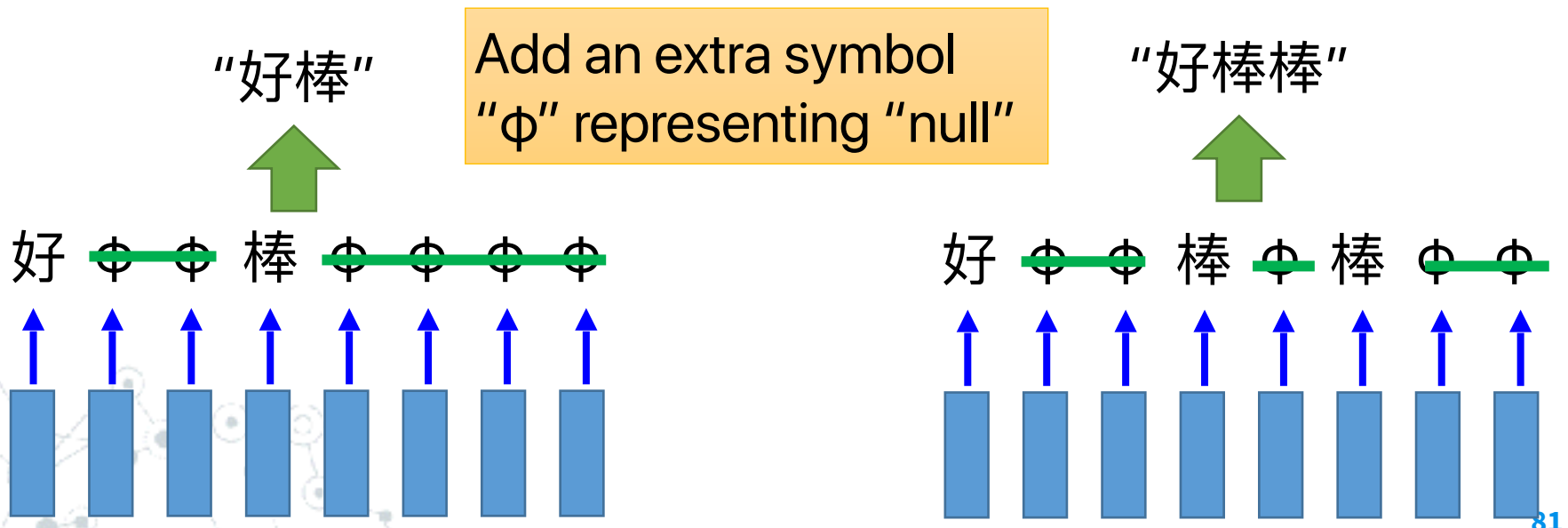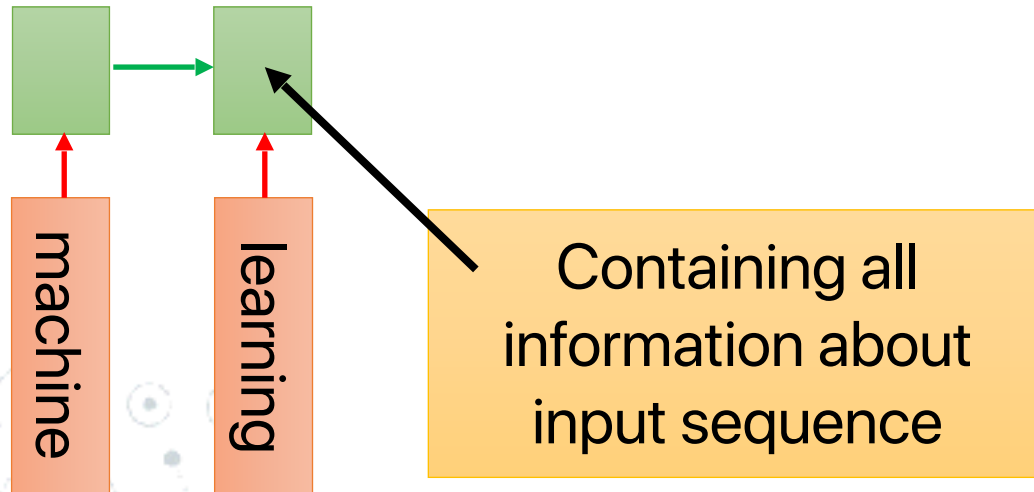Why can't it be "好棒棒"

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]
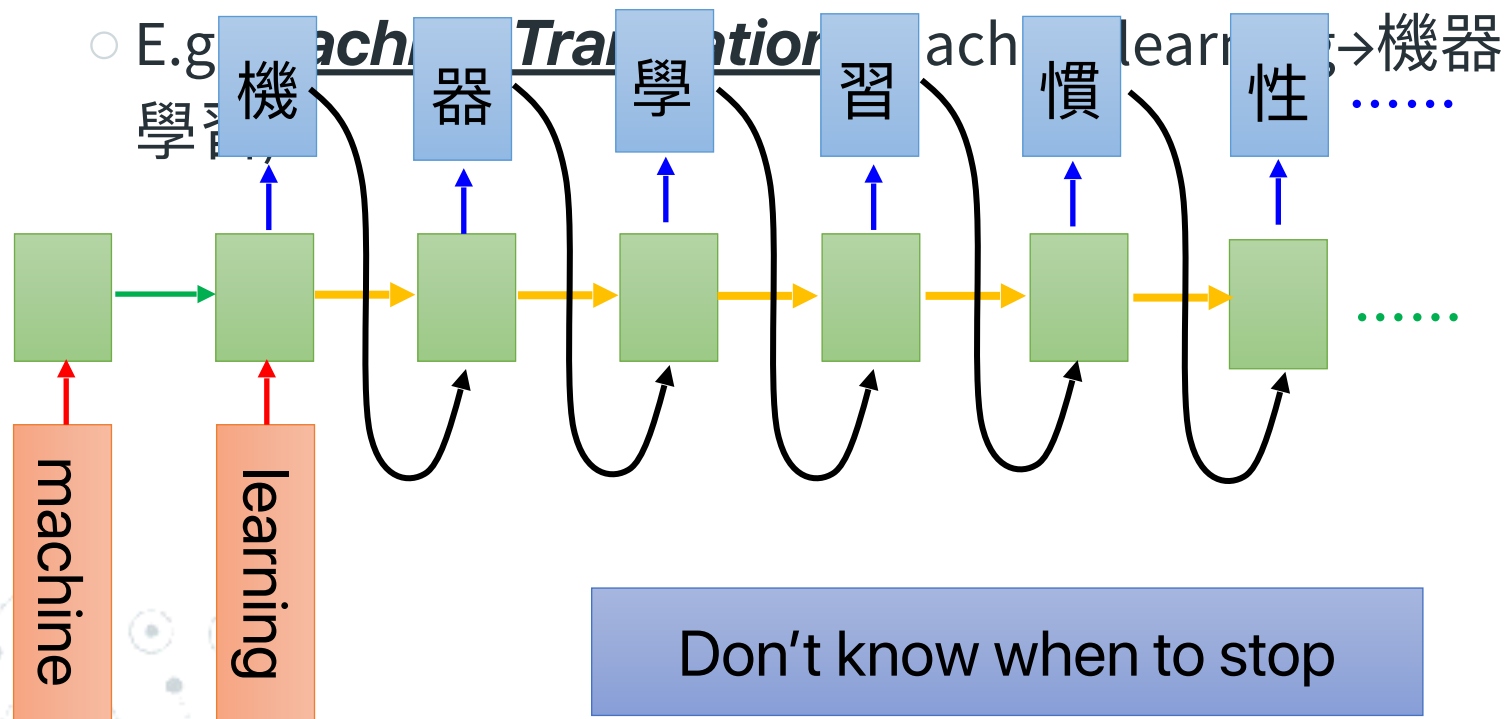
"好棒"

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)



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)



Don't know when to stop

# Many to Many (No Limitation)



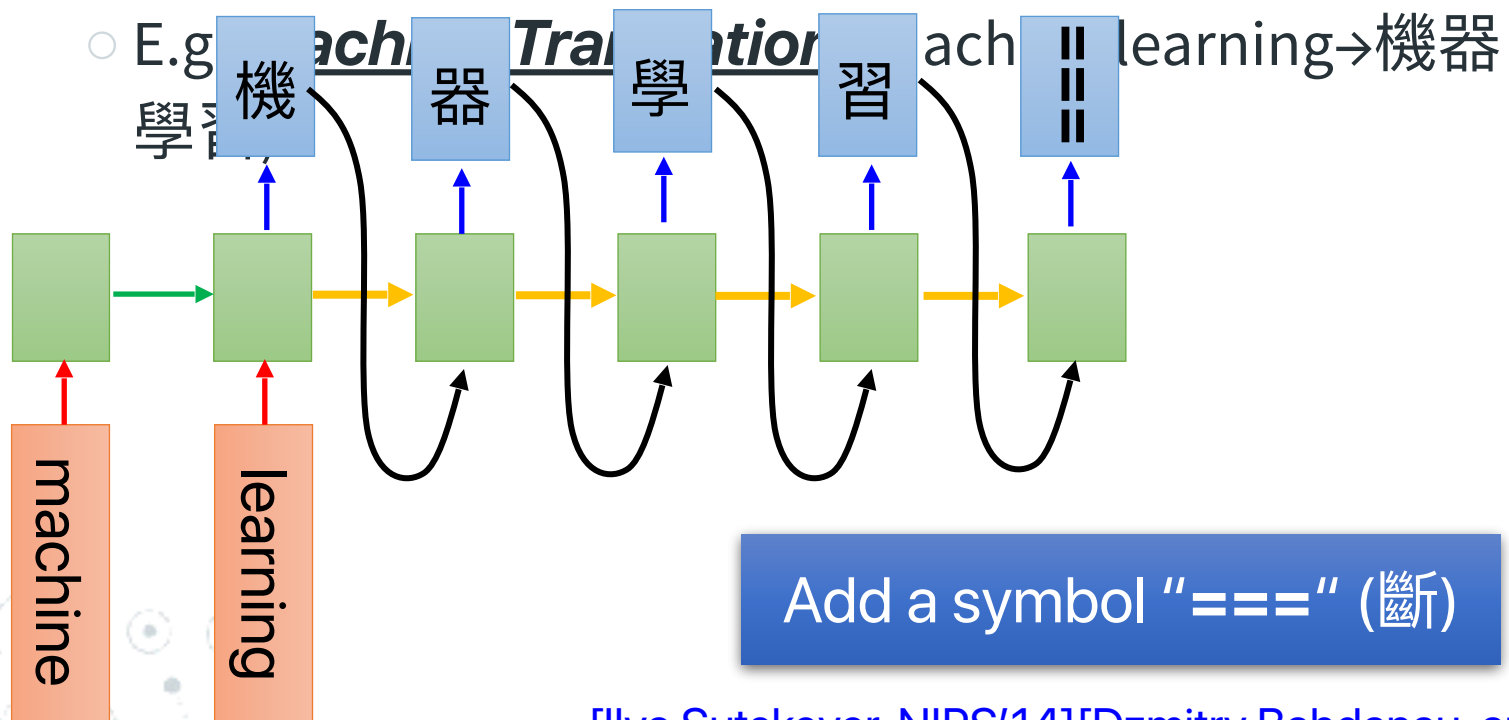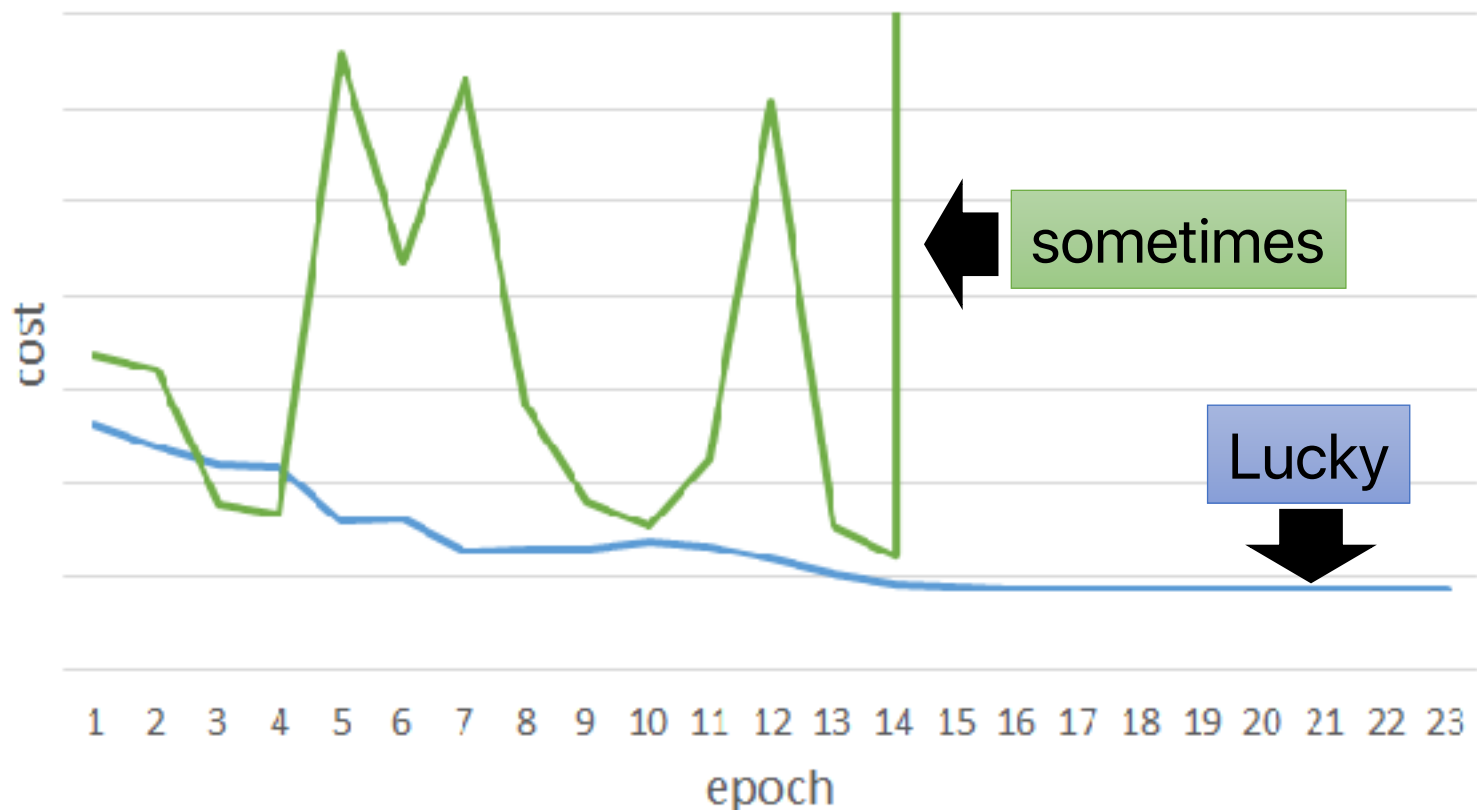推 tlkagk: =========斷=========

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***
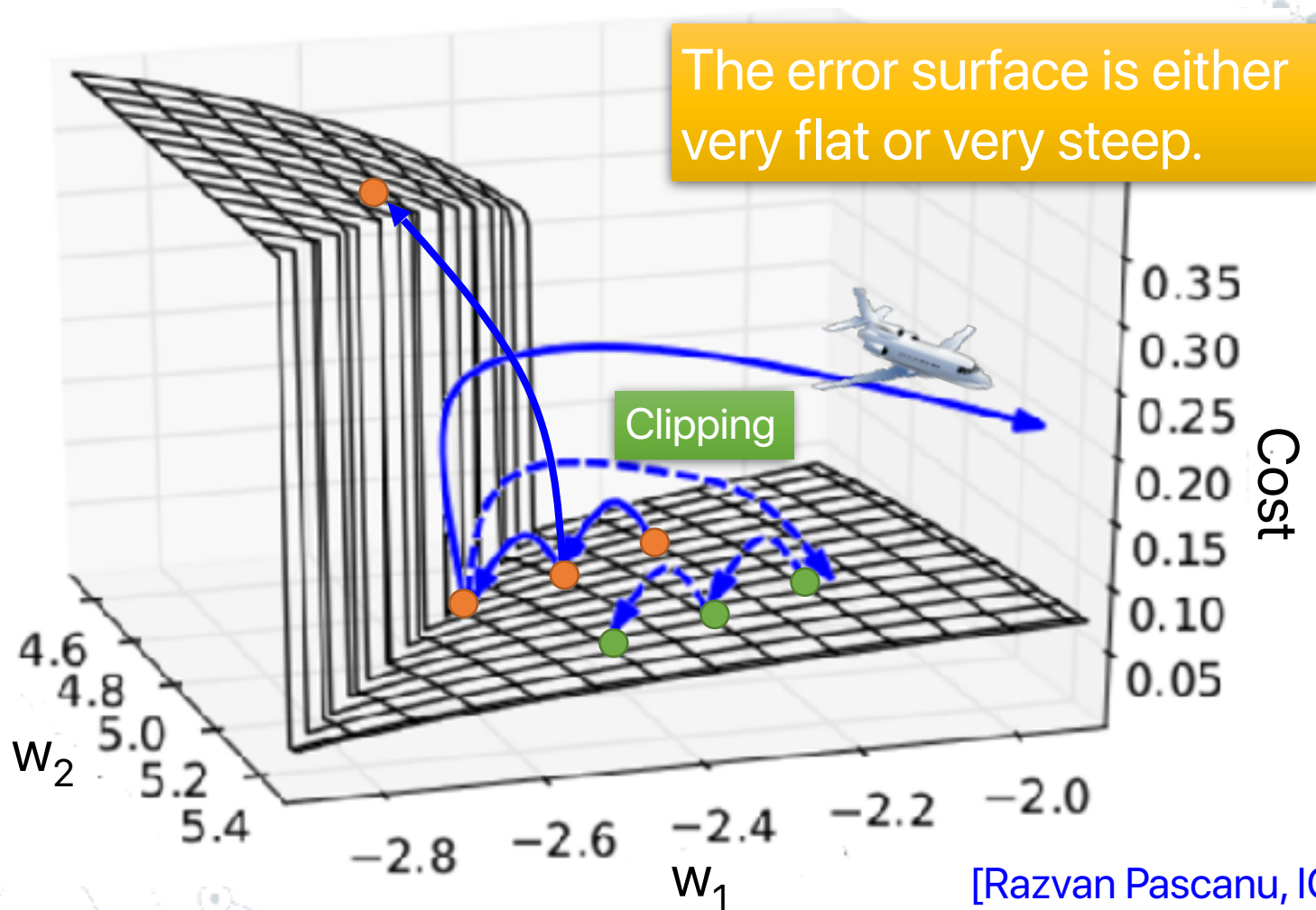
○ E.g. ***Machine Translation*** (machine learning→機器 學習)



Add a symbol "===" (斷)

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

# Unfortunately ......

◎ RNN-based network is not always easy to learn

Real experiments on Language modeling



cost

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

Cost

$w_2$

$w_1$

[Razvan Pascanu, ICML'13]

# Why?

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

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

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

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

| Large gradient | ⟹ | Small Learning rate? |
| small gradient | ⟹ | Large Learning rate? |

$= w^{999}$

## Toy Example



$y^1$    $y^2$    $y^3$    $y^{1000}$

1   1   1   1

w   w   ......   w

1   1   1   1

1   0   0   0

# Helpful Techniques

◎Nesterov's Accelerated Gradient (NAG):

　○Advance momentum method

◎RMS Prop

　○Advanced approach to give each parameter different learning rates

　○Considering the change of Second derivatives

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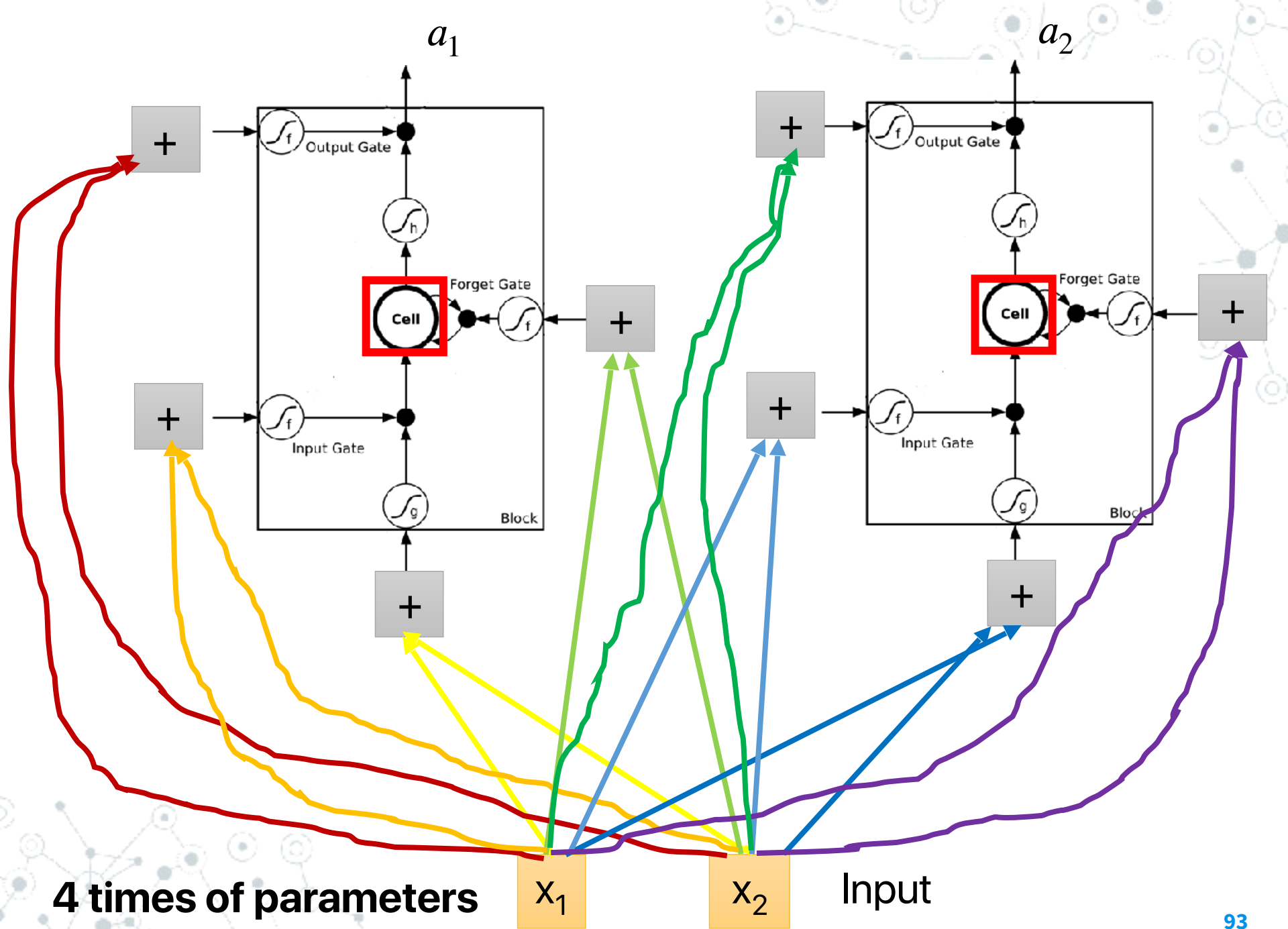
　○Can deal with gradient vanishing (not gradient explode)
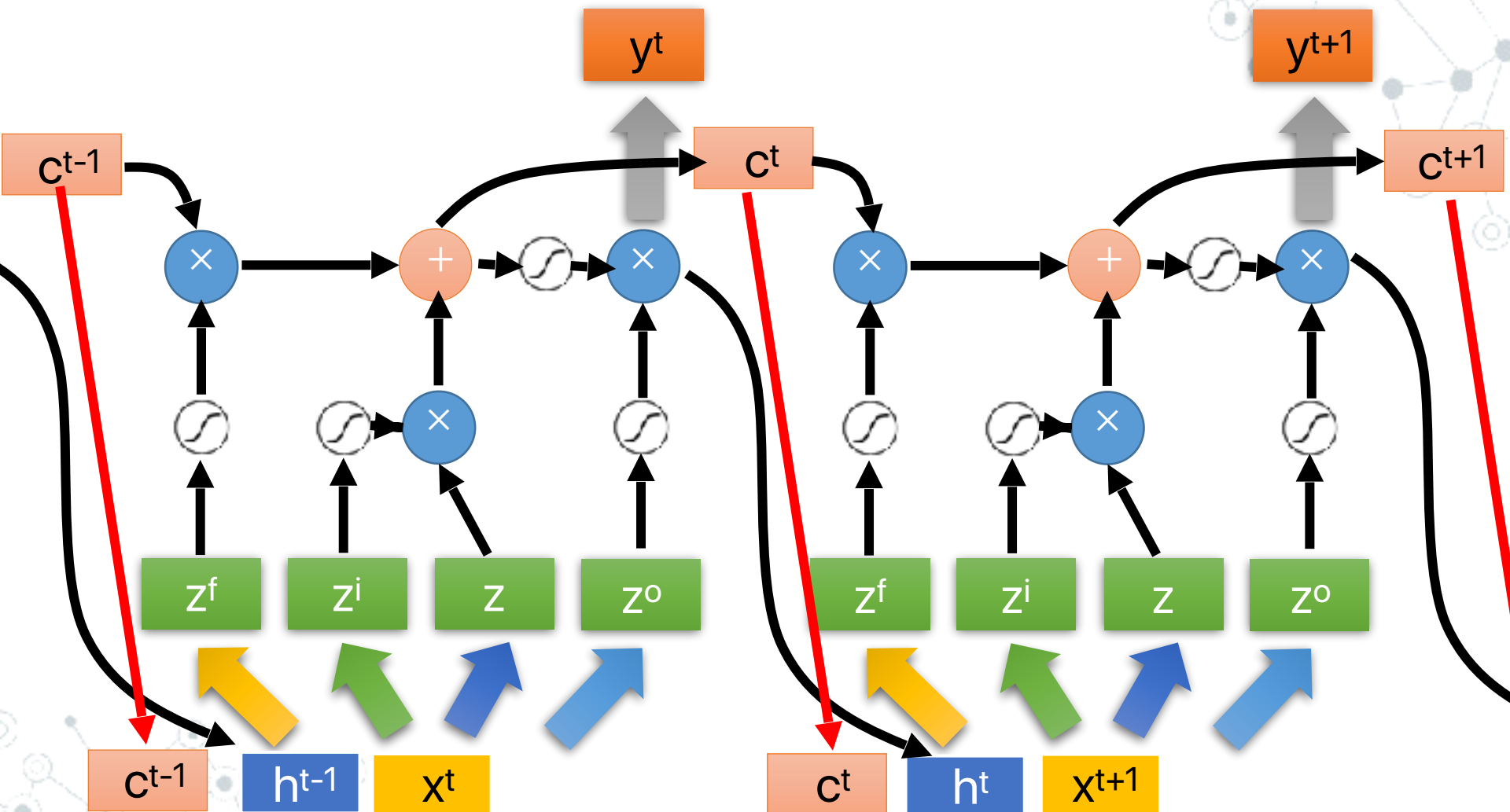
# Long Short-term Memory (LSTM)



Other part of the network

Signal control the output gate

(Other part of the network)

Output Gate

Special Neuron: 4 inputs, 1 output

Memory Cell

Forget Gate

Signal control the forget gate

(Other part of the network)

Signal control the input gate

(Other part of the network)

Input Gate

*LSTM*

Other part of the network

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

$z_o$

**Output Gate**

$f(z_o)$

<span style="color:red">multiply</span>

$h(c')$

h

**Forget Gate**

$c$

$f\left(z_f\right)$

$z_f$

$c'$

$c$

$cf\left(z_f\right)$

$f(z_i)$

$g(z)f(z_i)$

$z_i$

f

**Input Gate**

<span style="color:red">multiply</span>

$g(z)$

g

**Block**

$z$

Activation function f is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

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

Original Network:

➢Simply replace the neurons with LSTM



Input

$a_1$

$a_2$

Output Gate

Forget Gate

Cell

Input Gate

Block

**4 times of parameters**

$x_1$  $x_2$  Input

# LSTM

94

# Other Simpler Alternatives

Gated Recurrent Unit (GRU)

Structurally Constrained Recurrent Network (SCRN)



[Cho, EMNLP'14]



[Tomas Mikolov, ICLR'15]

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

➢ Outperform or be comparable with LSTM in 4 different tasks

# What is the next wave?

◎ Attention-based Model

Internal memory or information from output

...... Reading Head ...... Writing Head

Reading Head Controller

Writing Head Controller

Input x → DNN/ LSTM → output y

Already applied on speech recognition, caption generation, QA, visual QA

# What is the next wave?

◎ Attention-based Model

◎ End-To-End Memory Networks. S. Sukhbaatar, A. Szlam, J. Weston, R. Fergus. arXiv Pre-Print, 2015.

◎ Neural Turing Machines. Alex Graves, Greg Wayne, Ivo Danihelka. arXiv Pre-Print, 2014

◎ Ask Me Anything: Dynamic Memory Networks for Natural Language Processing. Kumar et al. arXiv Pre-Print, 2015

◎ Neural Machine Translation by Jointly Learning to Align and Translate. D. Bahdanau, K. Cho, Y. Bengio; International Conference on Representation Learning 2015.

◎ Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. Kelvin Xu et. al.. arXiv Pre-Print, 2015.

◎ Attention-Based Models for Speech Recognition. Jan Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, Yoshua Bengio. arXiv Pre-Print, 2015.

◎ Recurrent models of visual attention. V. Mnih, N. Hees, A. Graves and K. Kavukcuoglu. In NIPS, 2014.

◎ A Neural Attention Model for Abstractive Sentence Summarization. A. M. Rush, S. Chopra and J. Weston. EMNLP 2015.

# Concluding Remarks

# Concluding Remarks

◎ Introduction of deep learning

◎ Discussing some reasons using deep learning

◎ New techniques for deep learning

○ ReLU, Maxout

○ Giving all the parameters different learning rates

○ Dropout

◎ Network with memory

○ Recurrent neural network

○ Long short-term memory (LSTM)

# Reading Materials

◎ "Neural Networks and Deep Learning"
   ○ written by Michael Nielsen
   ○ http:// neuralnetworksanddeeplearning.com/

◎ "Deep Learning" (not finished yet)
   ○ Written by Yoshua Bengio, Ian J. Goodfellow and Aaron Courville
   ○ http://www.iro.umontreal.ca/~bengioy/ dlbook/

# Thank you!

# Appendix

# Matrix Operation



$$\sigma\left(\begin{bmatrix} W \end{bmatrix}\begin{bmatrix} x \end{bmatrix} + \begin{bmatrix} b \end{bmatrix}\right) = \begin{bmatrix} 0.98 \\ 0.12 \end{bmatrix} = a$$

# Why Deep? – Logic Circuits

◎ A two levels of basic logic gates can represent any Boolean function.

◎ However, no one uses two levels of logic gates to build computers

◎ Using multiple layers of logic gates to build some functions are much simpler (less gates needed).

**Boosting**

Input → Weak classifier / Weak classifier / ⋮ / Weak classifier → Combine →

**Deep Learning**

Weak classifier

Boosted weak classifier

Boosted Boosted weak classifier

$x_1$
$x_2$
⋮
$x_N$

......
......
......

# Maxout

Input

$z$

$w$

ReLU $\rightarrow$ $a$

$x$

$b$

$1$

Input

$w$

$+$ $z_1$

$b$

$0$

$+$ $z_2$

$0$

Max $\rightarrow$ $a$

$max\left\{ z_1, z_2 \right\}$

$x$

$1$

$a$

$z = wx + b$

$x$

$a$

$z_1 = wx + b$

$x$

$0$

# Maxout

Input

$w$

$z$

ReLU $\rightarrow$ $a$

$x$

$b$

1

Input

$w$

$+$ $z_1$

$b$

$x$

$w'$

$+$ $z_2$

Max $\rightarrow$ $a$

$b'$

1

$max\left\{ z_1, z_2 \right\}$

Learnable Activation Function

$a$

$z = wx + b$

$x$

$a$

$z_1 = wx + b$

$x$

$z_2 = w'x + b'$

107

# Getting Started

```
import
tensorflow as tf
```

# Graphs and Sessions

**Data Flow Graphs**

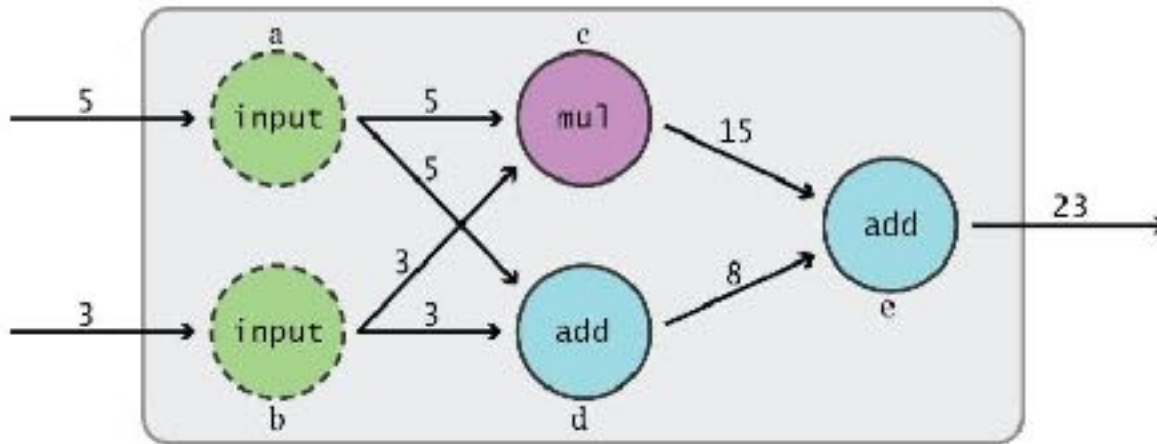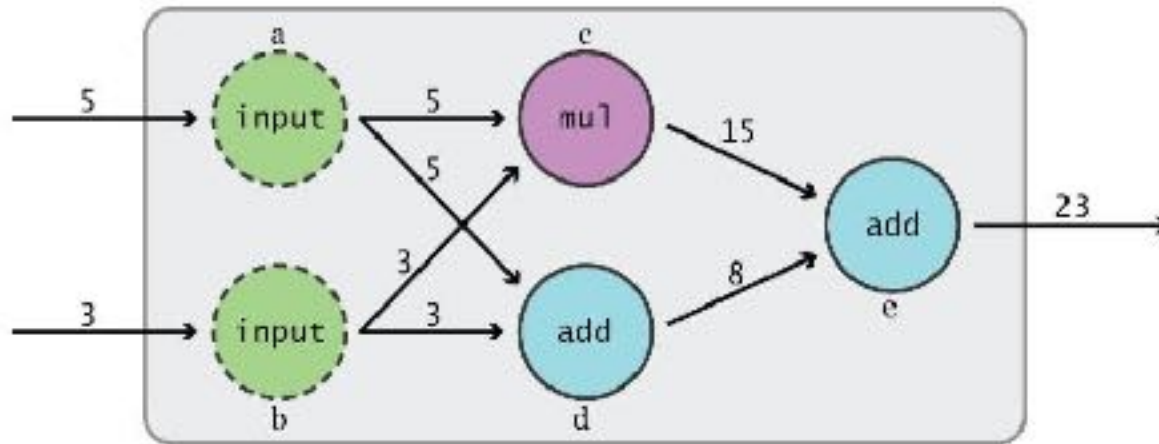# TensorFlow separates definition of computations from their execution



Graph from *TensorFlow for Machine Intelligence*

**Data Flow Graphs**

# Phase 1: assemble a graph

# Phase 2: use a session to execute operations in the graph



Graph from *TensorFlow for Machine Intelligence*

**Data Flow Graphs**

# Phase 1: assemble a graph

# Phase 2: use a session to execute operations in the graph

This might change in the future with eager mode!!



Graph from *TensorFlow for Machine Intelligence*

# What's a tensor?

## What's a tensor?

**An n-dimensional array**

0-d tensor: scalar (number)

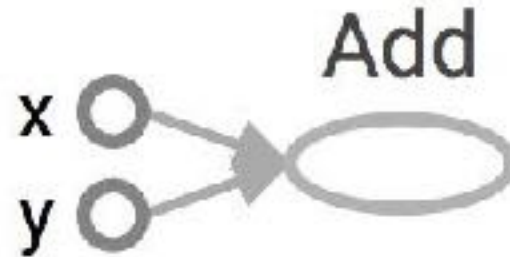1-d tensor: vector

2-d tensor: matrix

and so on

## Data Flow Graphs

```
import tensorflow as tf
a = tf.add(3, 5)
```
Visualized by TensorBoard

## Data Flow Graphs

```
import tensorflow as tf
a = tf.add(3, 5)
```
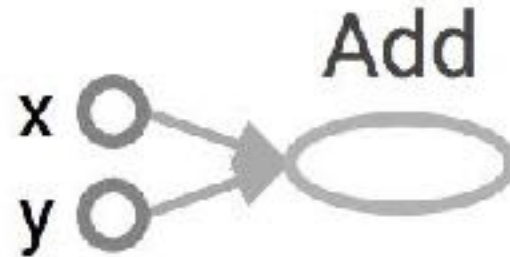Visualized by TensorBoard

Why x, y?

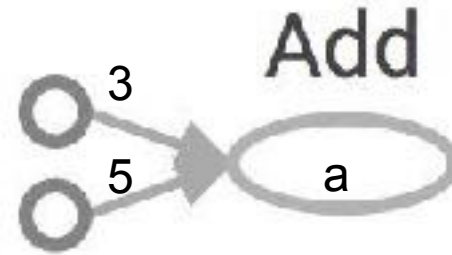TF automatically names the nodes when you don't explicitly name them.
x = 3
y = 5

## Data Flow Graphs

```
import tensorflow as tf
a = tf.add(3, 5)
```
Interpreted?

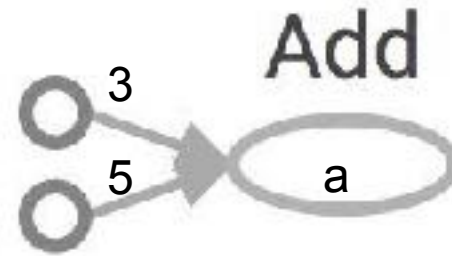Nodes: operators, variables, and constants
Edges: tensors

## Data Flow Graphs

```
import tensorflow as tf
a = tf.add(3, 5)
```

Interpreted?
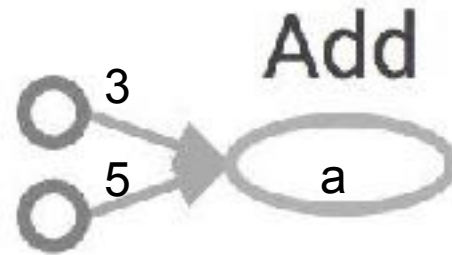
Nodes: operators, variables, and constants
Edges: tensors

Tensors are data.
TensorFlow = tensor + flow = data + flow

Add

3
5
a

## Data Flow Graphs

```
import tensorflow as tf
a = tf.add(3, 5)
print(a)
```



```
>> Tensor("Add:0", shape=(),
dtype=int32)
```
(Not 8)

**How to get the value of a?**

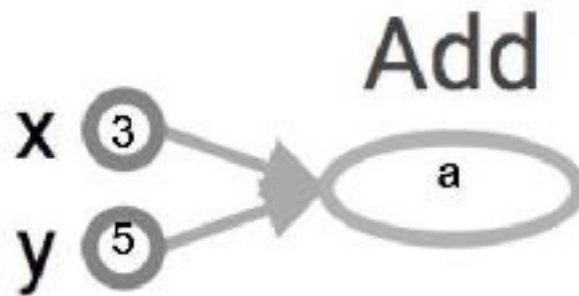Create a **session**, assign it to variable sess so we can call it later

Within the session, evaluate the graph to fetch the value of a

Create a **session**, assign it to variable sess so we can call it later

Within the session, evaluate the graph to fetch the value of a

```
import tensorflow as tf
a = tf.add(3, 5)
sess = tf.Session()
print(sess.run(a))
sess.close()
```
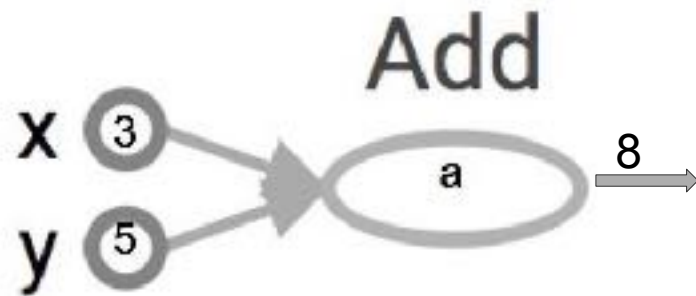
Create a **session**, assign it to variable sess so we can call it later

Within the session, evaluate the graph to fetch the value of a

```
import tensorflow as tf
a = tf.add(3, 5)
sess = tf.Session()
print(sess.run(a))
sess.close()
```
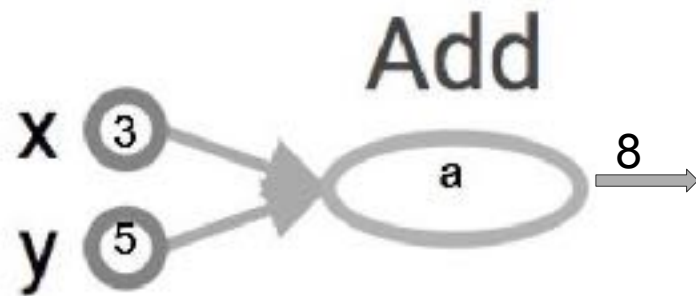
**How to get the value of a?**

Create a **session**, assign it to variable sess so we can call it later

Within the session, evaluate the graph to fetch the value of a

```
import tensorflow as tf
a = tf.add(3, 5)
sess = tf.Session()
with tf.Session() as sess:
        print(sess.run(a))
sess.close()
```

Add

x ③

y ⑤

a

8

**tf.Session()**

A Session object encapsulates the environment in which Operation objects are executed, and Tensor objects are evaluated.
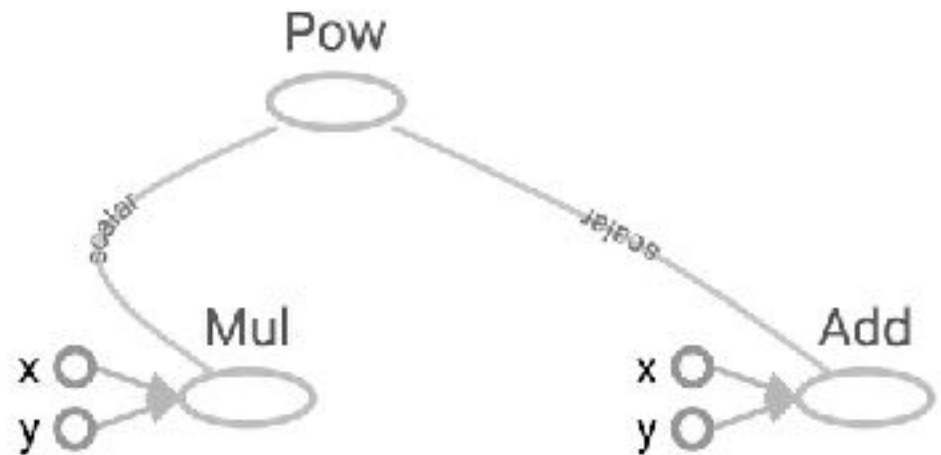
**tf.Session()**

A Session object encapsulates the environment in which Operation objects are executed, and Tensor objects are evaluated.

Session will also allocate memory to store the current values of variables.

## More graph
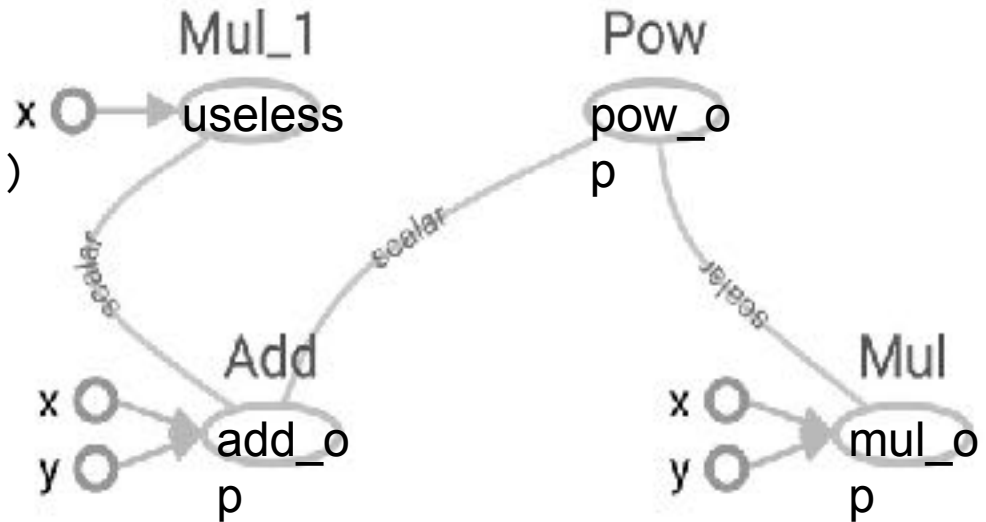
```
x = 2
y = 3
op1 = tf.add(x, y)
op2 = tf.multiply(x, y)
op3 = tf.pow(op2, op1)
with tf.Session() as sess
        op3 = sess.run(op
```
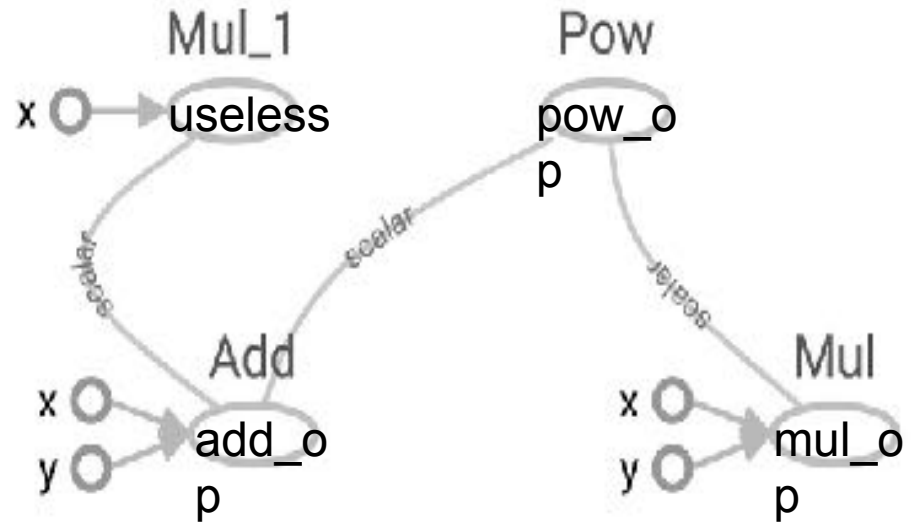
Visualized by TensorBoard

## Subgraphs

```
x = 2
y = 3
add_op = tf.add(x, y)
mul_op = tf.multiply(x, y)
useless = tf.multiply(x, add_op)
pow_op = tf.pow(add_op, mul_op)
with tf.Session() as sess:
        z = sess.run(pow_op)
```
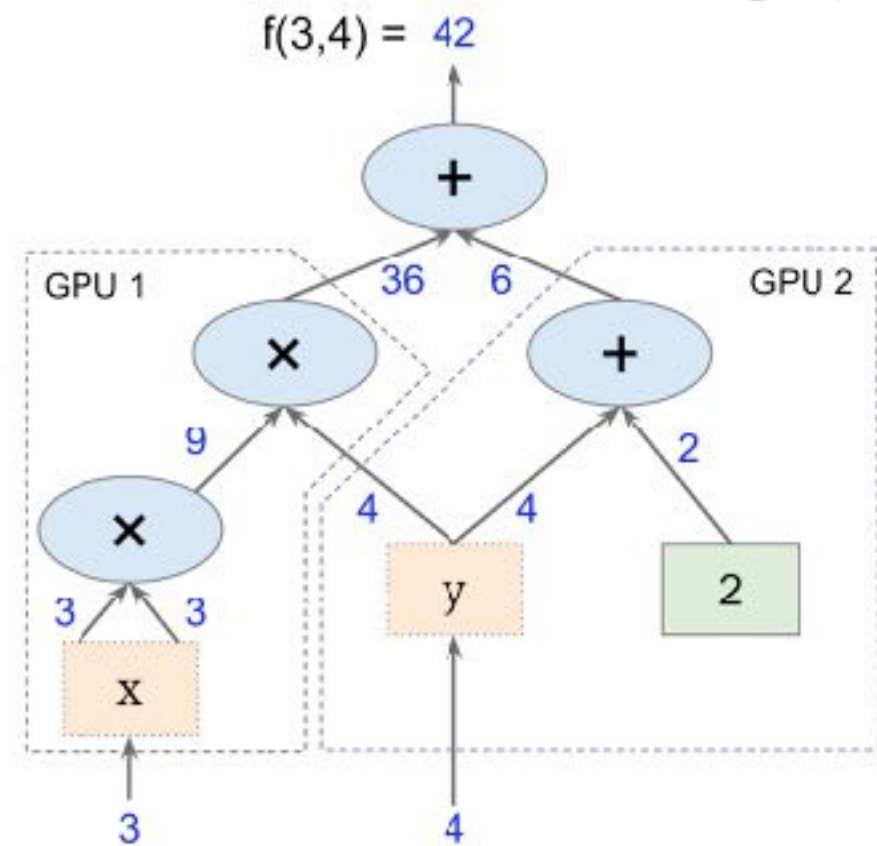
## Subgraphs

```
x = 2
y = 3
add_op = tf.add(x, y)
mul_op = tf.multiply(x, y)
useless = tf.multiply(x, add_op)
pow_op = tf.pow(add_op, mul_op)
with tf.Session() as sess:
        z, not_useless = sess.run
```

## Subgraphs

Possible to break graphs into several chunks and run them parallelly across multiple CPUs, GPUs, TPUs, or other devices

Example: AlexNet

Graph from *Hands-On Machine Learning with Scikit-Learn and TensorFlow*

# Distributed Computation

To put part of a graph on a specific CPU or GPU:

```python
# Creates a graph.
with tf.device('/gpu:2'):
  a = tf.constant([1.0, 2.0, 3.0, 4.0, 5.0, 6.0], name='a')
  b = tf.constant([1.0, 2.0, 3.0, 4.0, 5.0, 6.0], name='b')
  c = tf.multiply(a, b)

# Creates a session with log_device_placement set to True.
sess =
tf.Session(config=tf.ConfigProto(log_device_placement=True))

# Runs the op.
print(sess.run(c))
```

## tf.Graph()

create a graph:

```
g = tf.Graph()
```

## tf.Graph()

to add operators to a graph, set it as default:

```
g = tf.Graph()
with g.as_default():
        x = tf.add(3, 5)
sess = tf.Session(graph=g)
with tf.Session() as sess:
        sess.run(x)
```

## tf.Graph()

To handle the default graph:

```
g = tf.get_default_graph()
```

## tf.Graph()

Do not mix default graph and user created graphs

```
g = tf.Graph()

# add ops to the default graph
a = tf.constant(3)

# add ops to the user created graph
with g.as_default():
        b = tf.constant(5)
```

## tf.Graph()

Do not mix default graph and user created graphs

```
g1 = tf.get_default_graph()
g2 = tf.Graph()

# add ops to the default graph
with g1.as_default():
        a = tf.Constant(3)

# add ops to the user created graph
with g2.as_default():
        b = tf.Constant(5)
```

# Install Tensorflow

# Install Tensorflow

To install the library we will create an environment in Anaconda with **python 3.5** we name it **tensorflow**. However, you may choose your own desired name for it. Open command prompt (or terminal) and type:

```
conda create --name tensorflow python=3.5
```

Once the environment is created, we can activate the environment:

**(for Windows):**

```
activate tensorflow
```

**(for Linux & Mac):**

```
source activate tensorflow
```

**(for Windows):**

*(CPU version):*

```
pip install --upgrade tensorflow
```

*(GPU version):*

```
pip install --upgrade tensorflow-gpu
```

## For Mac - CPU

```
pip install --ignore-installed --upgrade https://storage.googleapis.com/
tensorflow/mac/cpu/tensorflow-1.10.0-py3-none-any.whl
```

## For Linux - CPU

```
pip install --ignore-installed --upgrade https://storage.googleapis.com/
tensorflow/linux/cpu/tensorflow-1.10.0-cp35-cp35m-linux_x86_64.whl
```
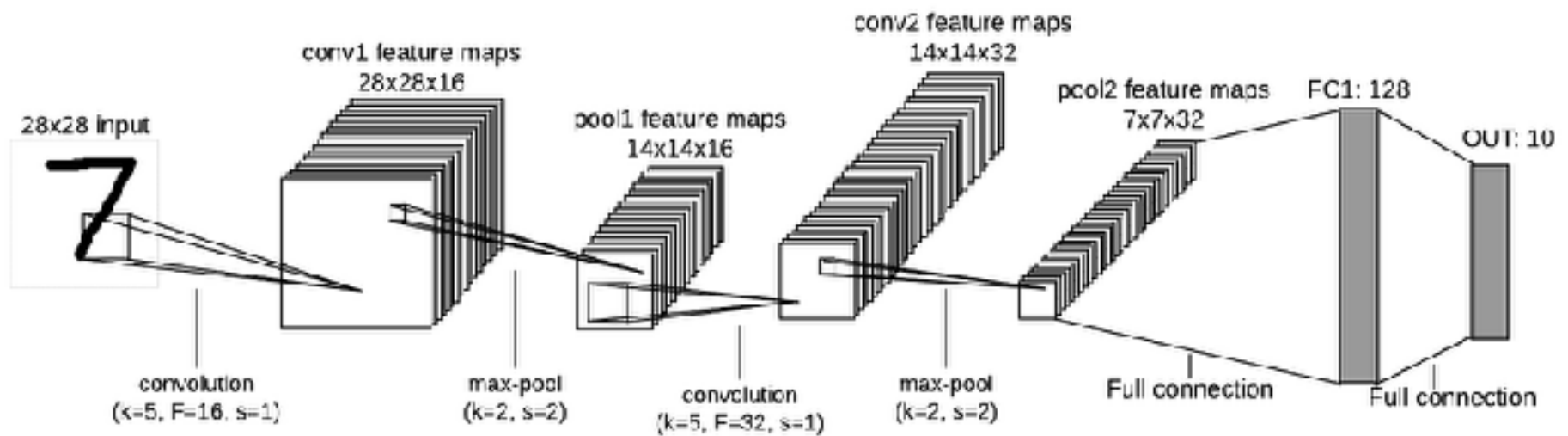
## For Linux - GPU

```
pip install --ignore-installed --upgrade https://storage.googleapis.com/
tensorflow/linux/gpu/tensorflow_gpu-1.10.0-cp35-cp35m-linux_x86_64.whl
```

## Install Jupiter

◎ conda install jupyter

# CNN (Convolution Neural Networks)

# See Jupiter Notebook

# Thank you