# CSC4007 Advanced Machine Learning
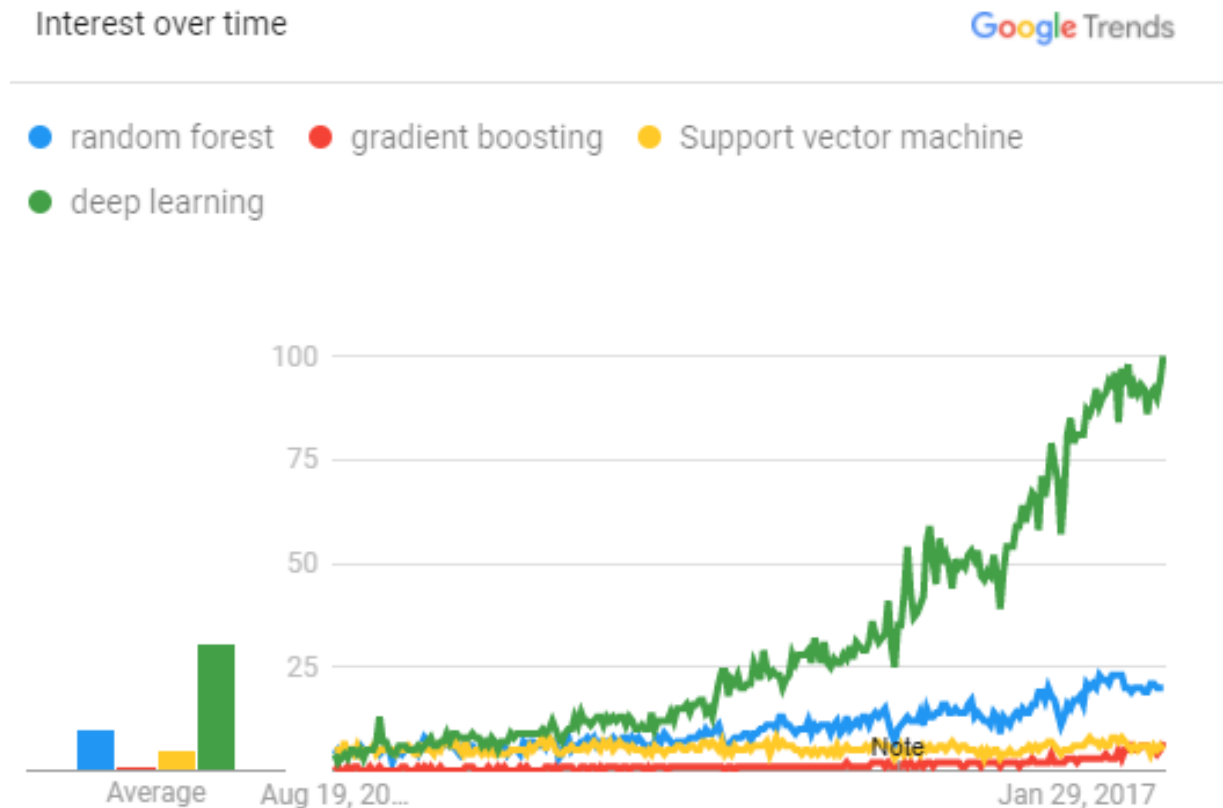
## **Lesson 08**: Deep Learning

by Vien Ngo
EEECS / ECIT / DSSC

# Outline

- Neural network basics and representation

- Perceptron learning, multi-layer perceptron

- Neural network training: Backpropagation

- **Modern neural network architecture (a.k.a Deep learning):**

  - **Convolutional neural network (CNN)**

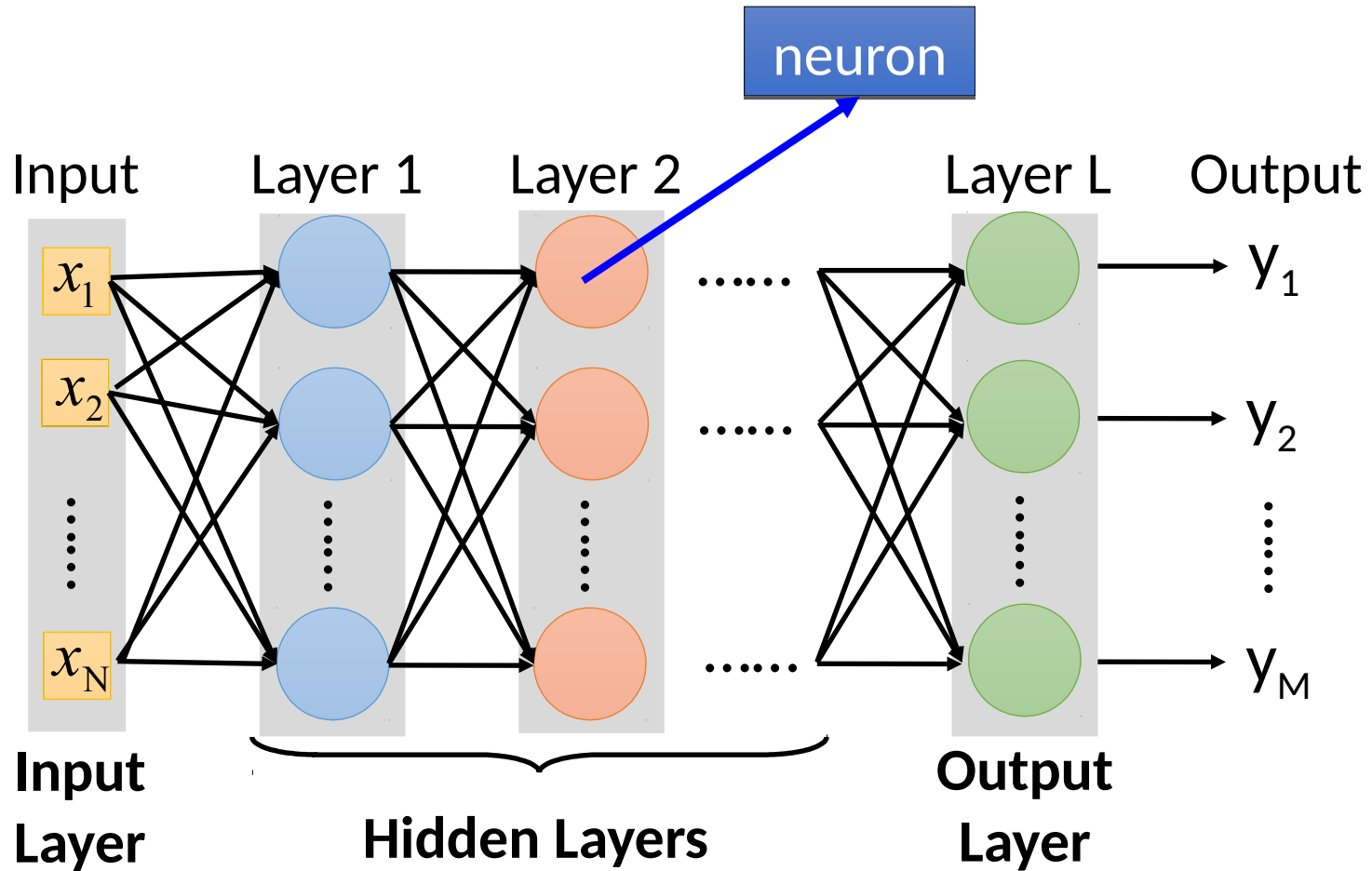  - **Recurrent neural network (RNN), long-short term memory network (LSTM)**

# Deep learning

- Google Trends: attracts lots of attention.
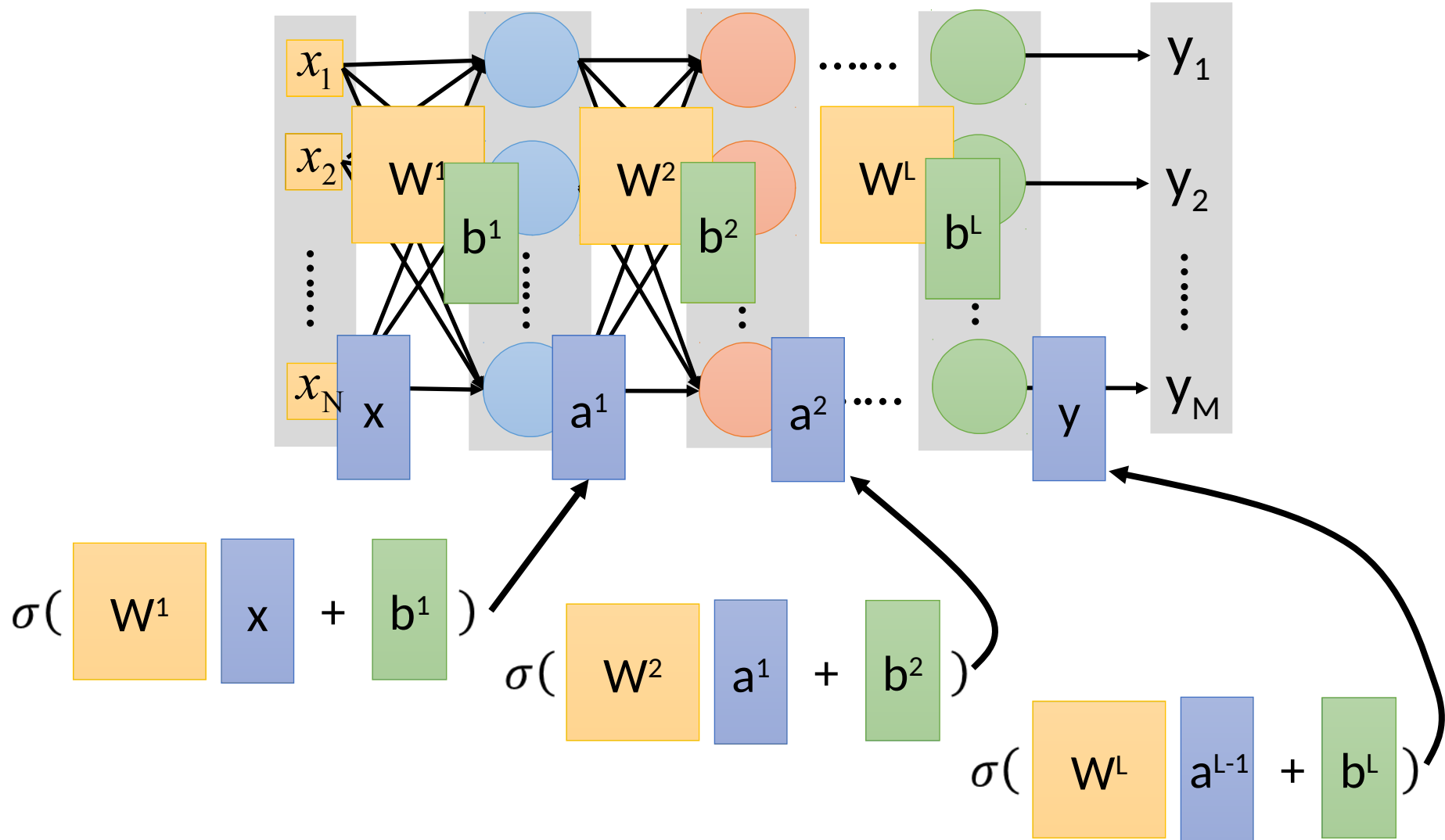  - Deep learning obtains many exciting results.

# Deep Neural Network
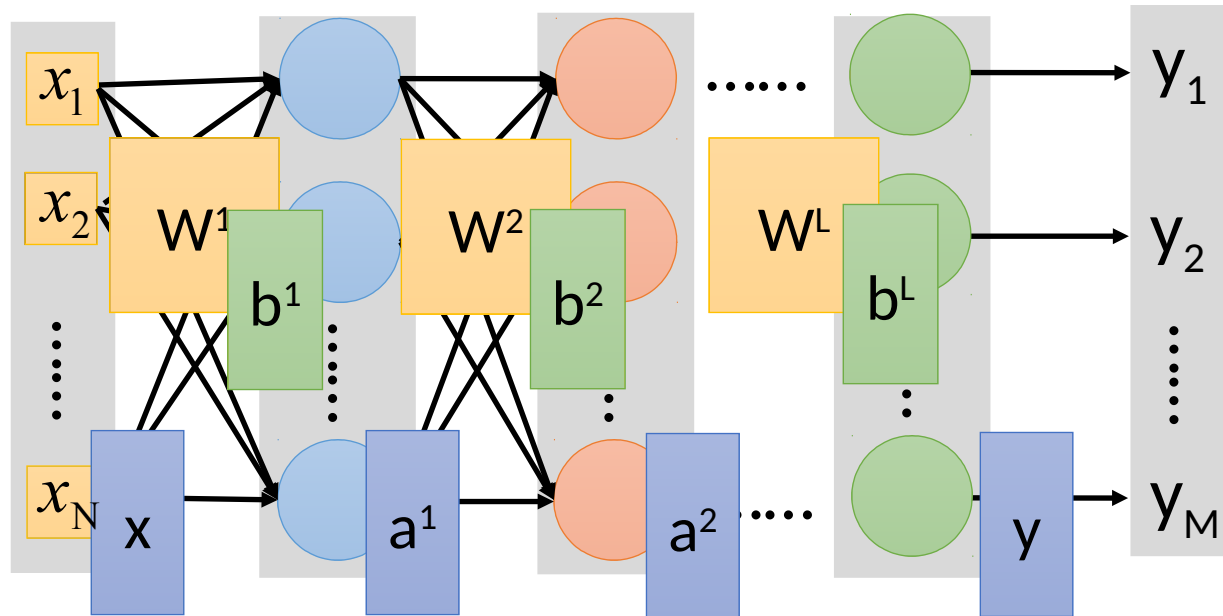
Example: a fully connected deep NN



neuron

Input    Layer 1    Layer 2    Layer L    Output

**Input Layer**

**Hidden Layers**

**Output Layer**

Deep means many hidden layers

# Deep Neural Network



$$\sigma(\; W^1 \; x \; + \; b^1 \;)$$

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

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

# Deep Neural Network: Forward propagation



Using parallel computing techniques to speed up matrix operation

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

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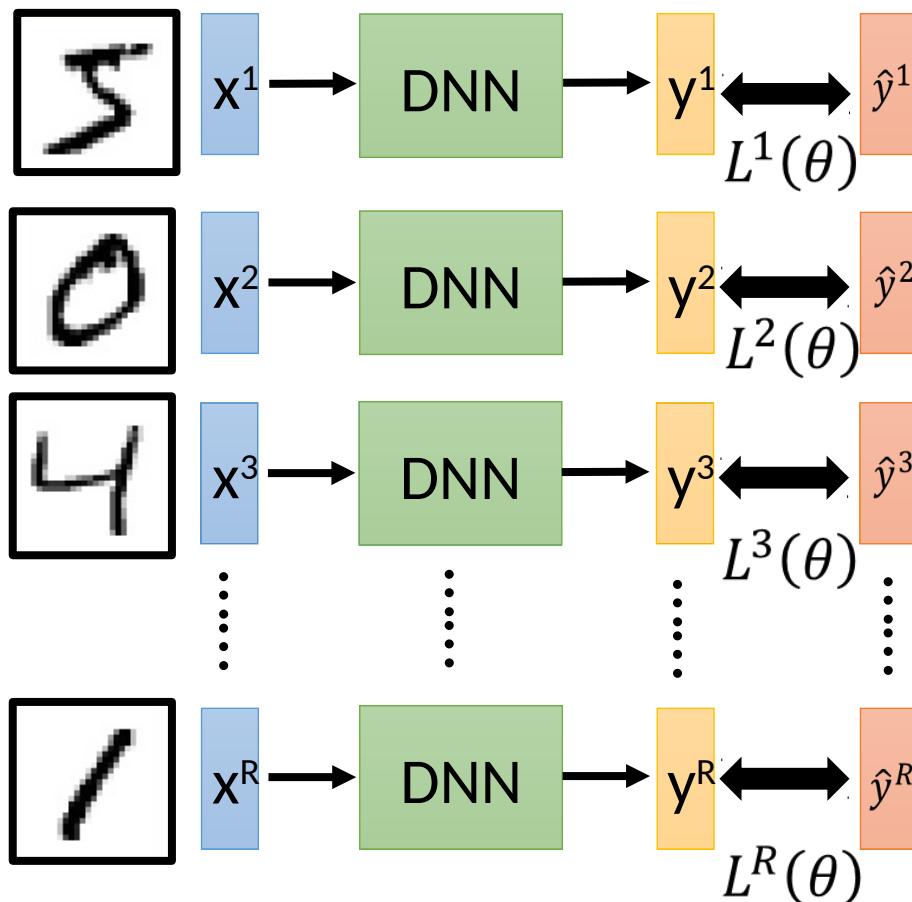
# Cost

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



$$L(\theta)$$

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 $\theta$ is on this task

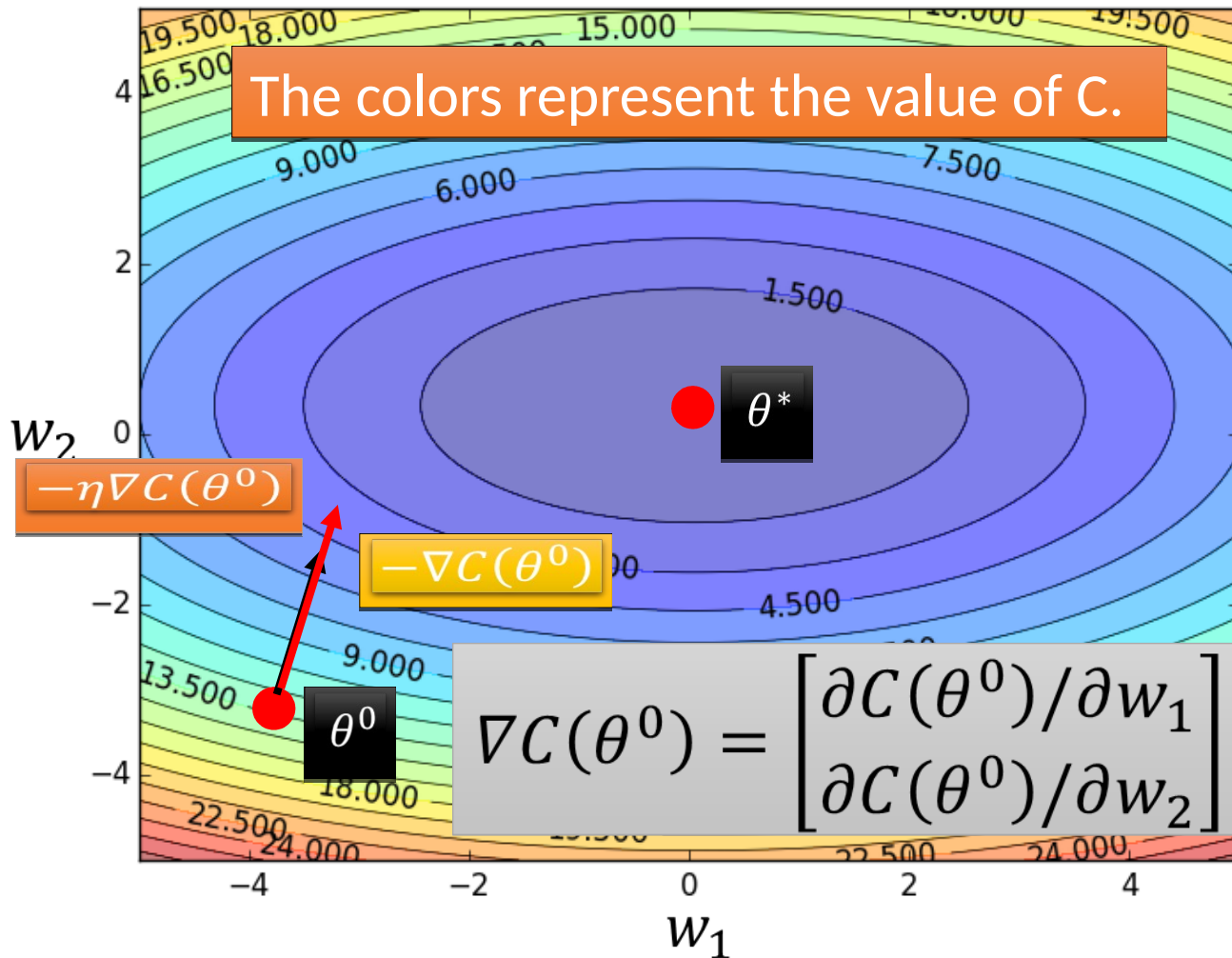Find the network parameters $\theta^*$ that minimize this value

# 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.

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

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

$\theta^*$

$\theta^0$

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

Randomly pick a starting point $\theta^0$

Compute the negative gradient at $\theta^0$
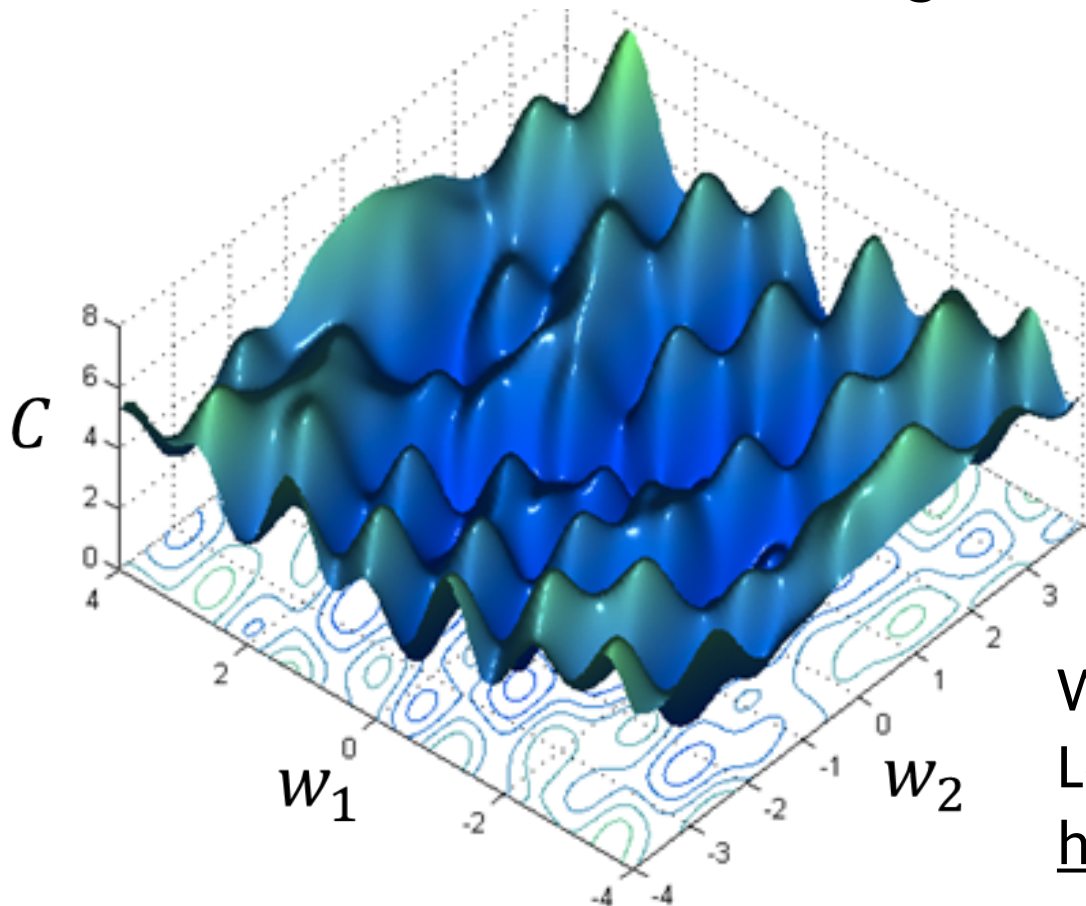
$\Rightarrow -\nabla C(\theta^0)$

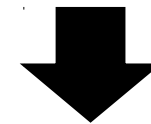Times the learning rate $\eta$

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

# Training DNN via Backpropagation

- Gradient descent never guarantee global minima
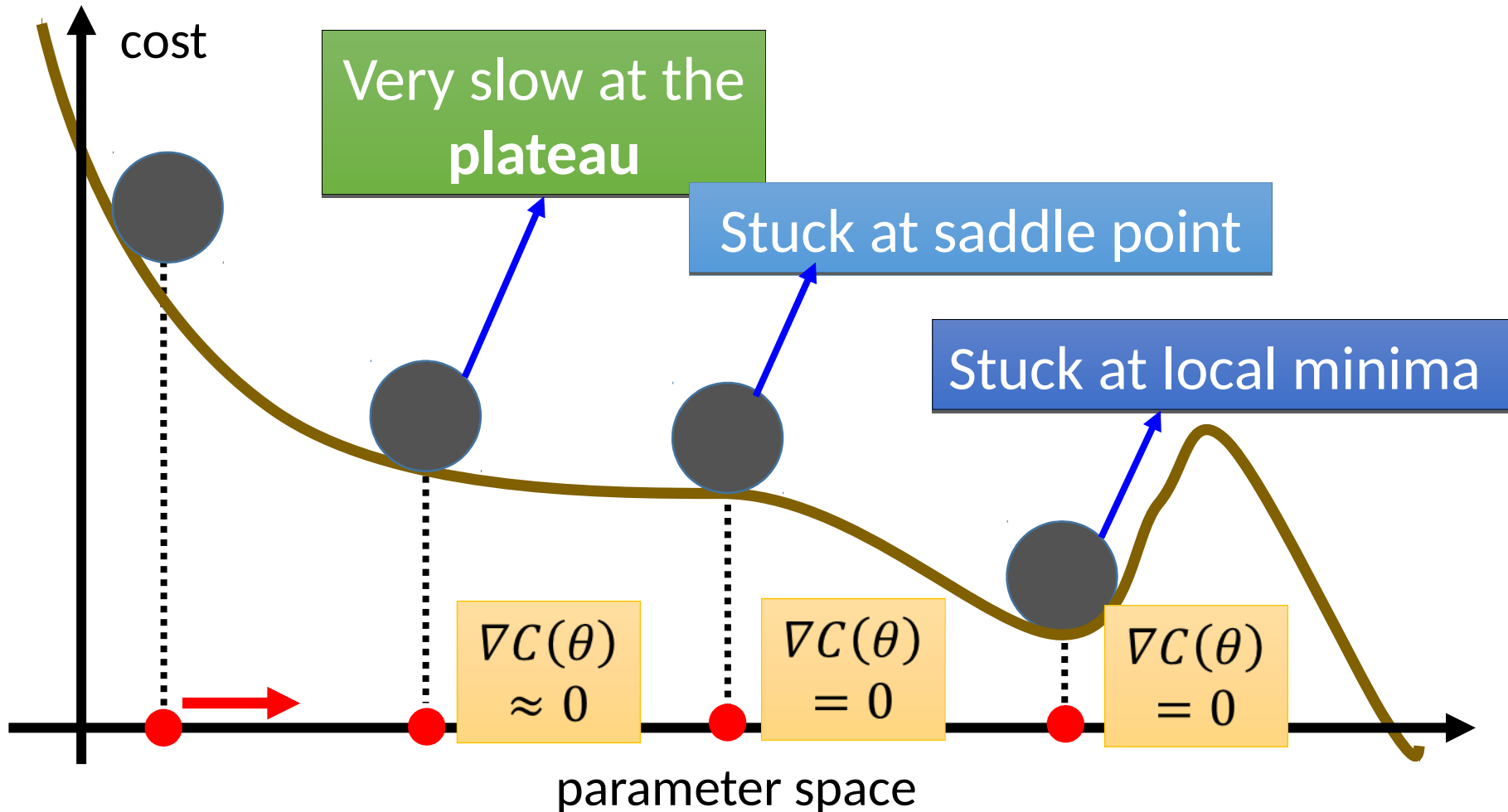


Different initial point $\theta^0$

Reach different minima, so different results

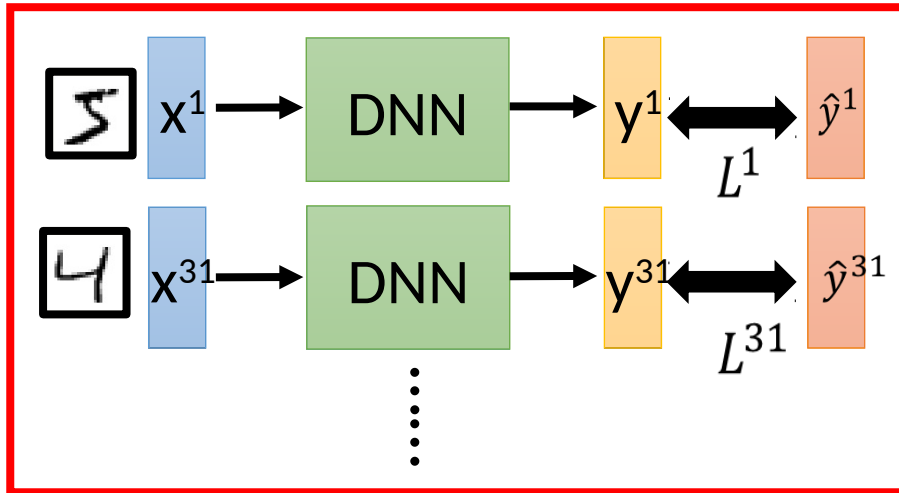Who is Afraid of Non-Convex Loss Functions?
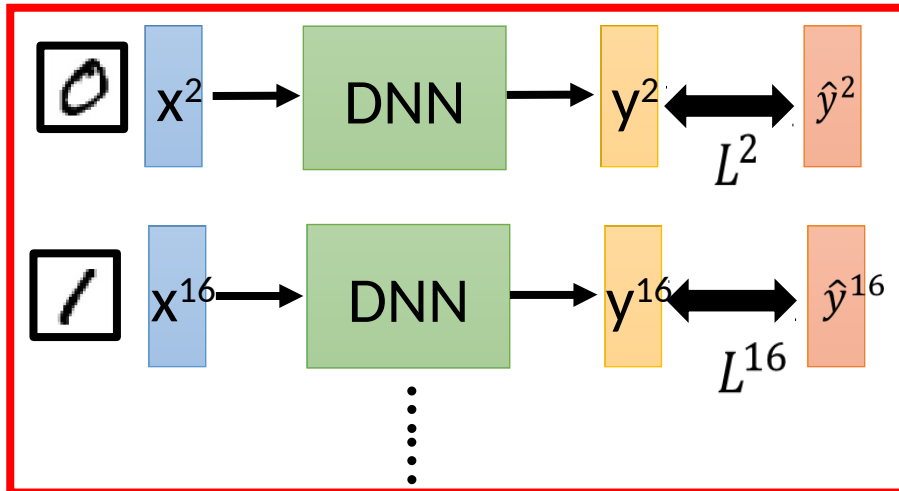http://videolectures.net/eml07_lecun_wia/

# 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
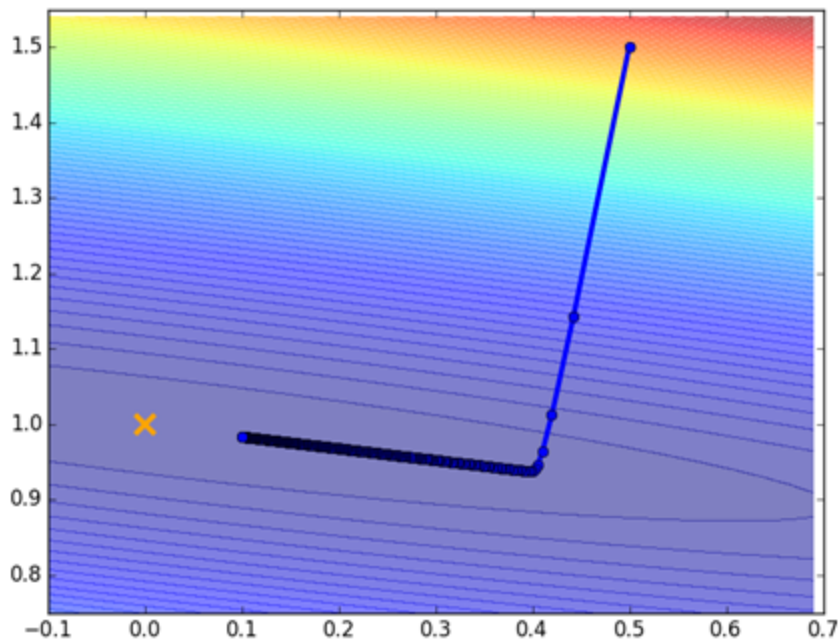
# Training with Mini-batch



➢ Randomly initialize $\theta^0$

➢ 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!

# Training with Mini-batch

## Original Gradient Descent

## With Mini-batch



unstable

The colors represent the total C on all training data.

# Training with Mini-batch

Faster    Better!

➤ Randomly initialize $\theta^0$

➤ Pick the 1st batch
$$C = C^1 + C^{31} + \cdots$$
$$\theta^1 \leftarrow \theta^0 - \eta \nabla C(\theta^0)$$

➤ Pick the 2nd batch
$$C = C^2 + C^{16} + \cdots$$
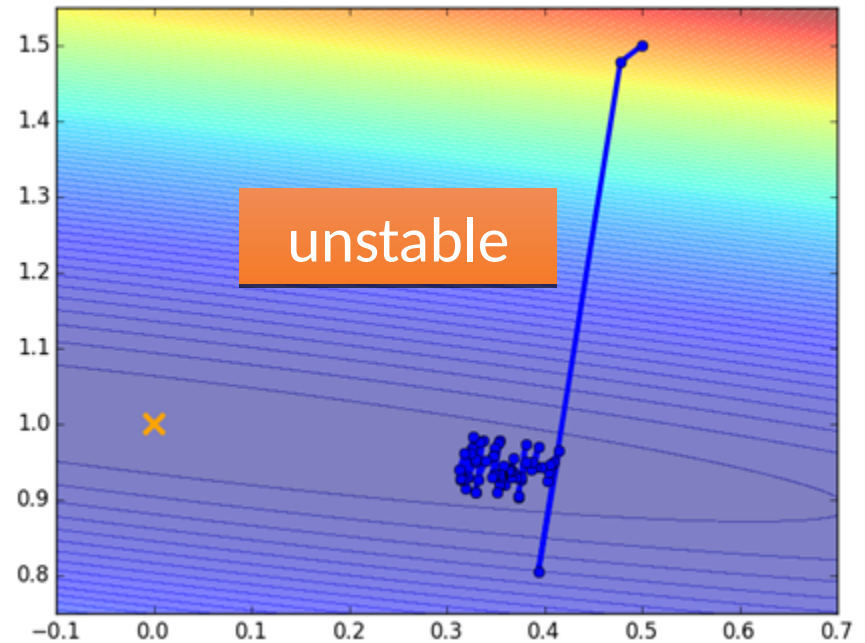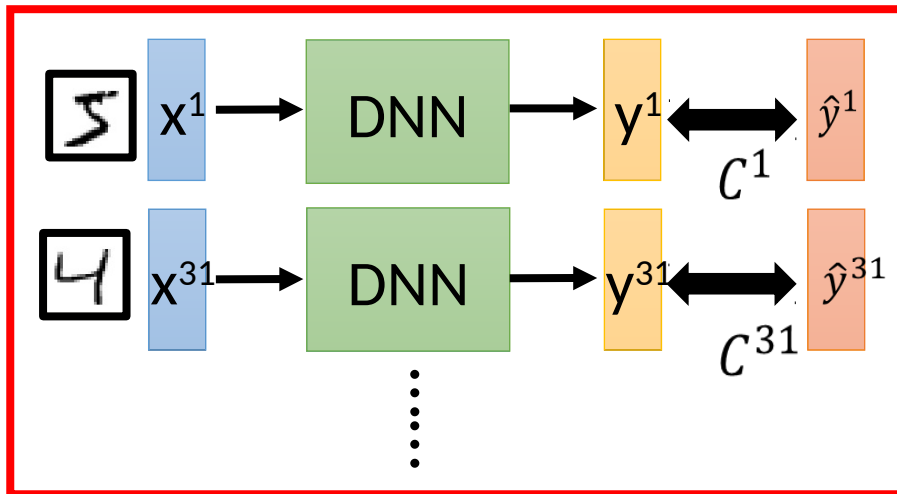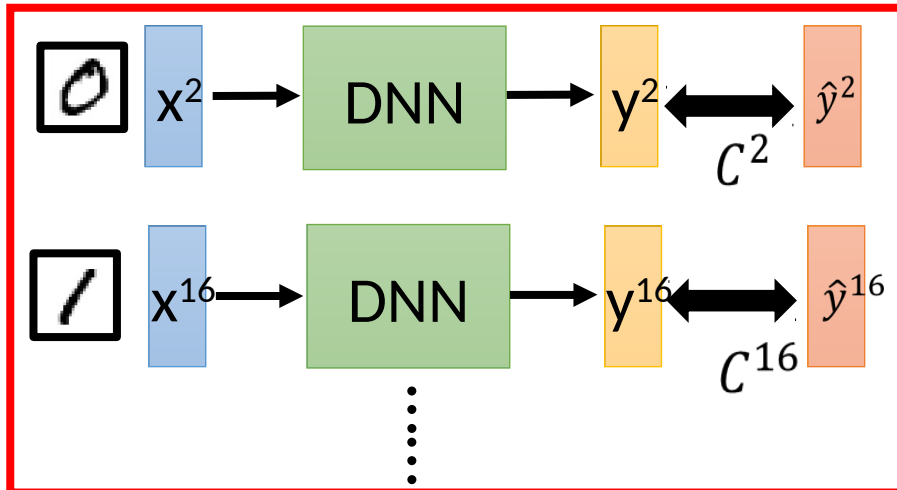$$\theta^2 \leftarrow \theta^1 - \eta \nabla C(\theta^1)$$
$$\vdots$$

➤ Until all mini-batches have been picked

one epoch

Repeat the above process

Mini-batch

$x^1$ → DNN → $y^1$ ⬌ $\hat{y}^1$
$C^1$

$x^{31}$ → DNN → $y^{31}$ ⬌ $\hat{y}^{31}$
$C^{31}$

Mini-batch

$x^2$ → DNN → $y^2$ ⬌ $\hat{y}^2$
$C^2$

$x^{16}$ → DNN → $y^{16}$ ⬌ $\hat{y}^{16}$
$C^{16}$

# 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

# 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.
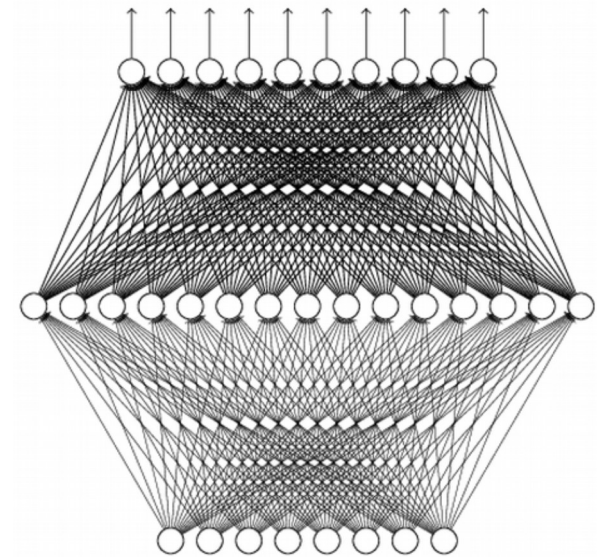
# Universality Theorem

Any continuous function f

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

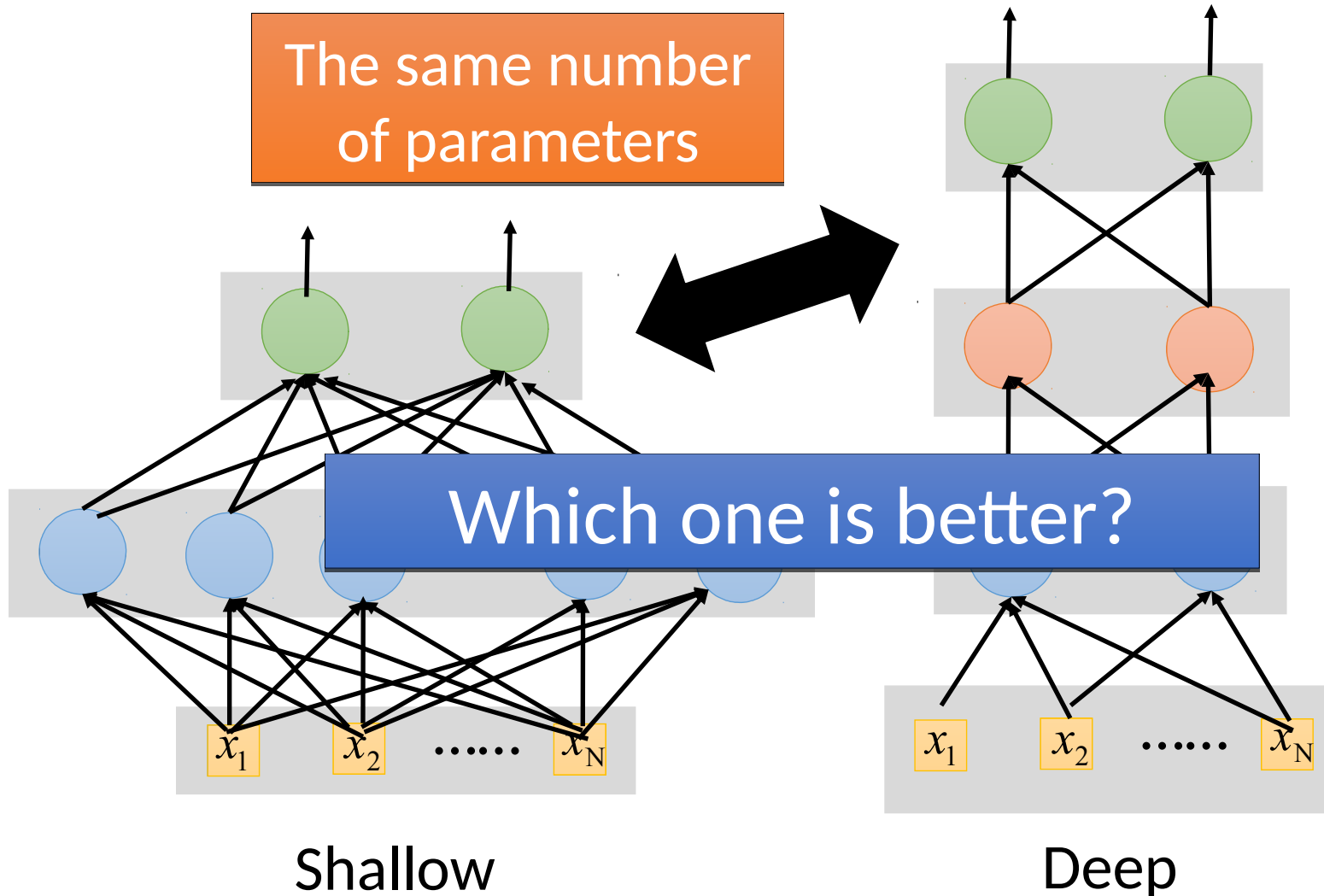Can be realized by a network with one hidden layer

(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?

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.
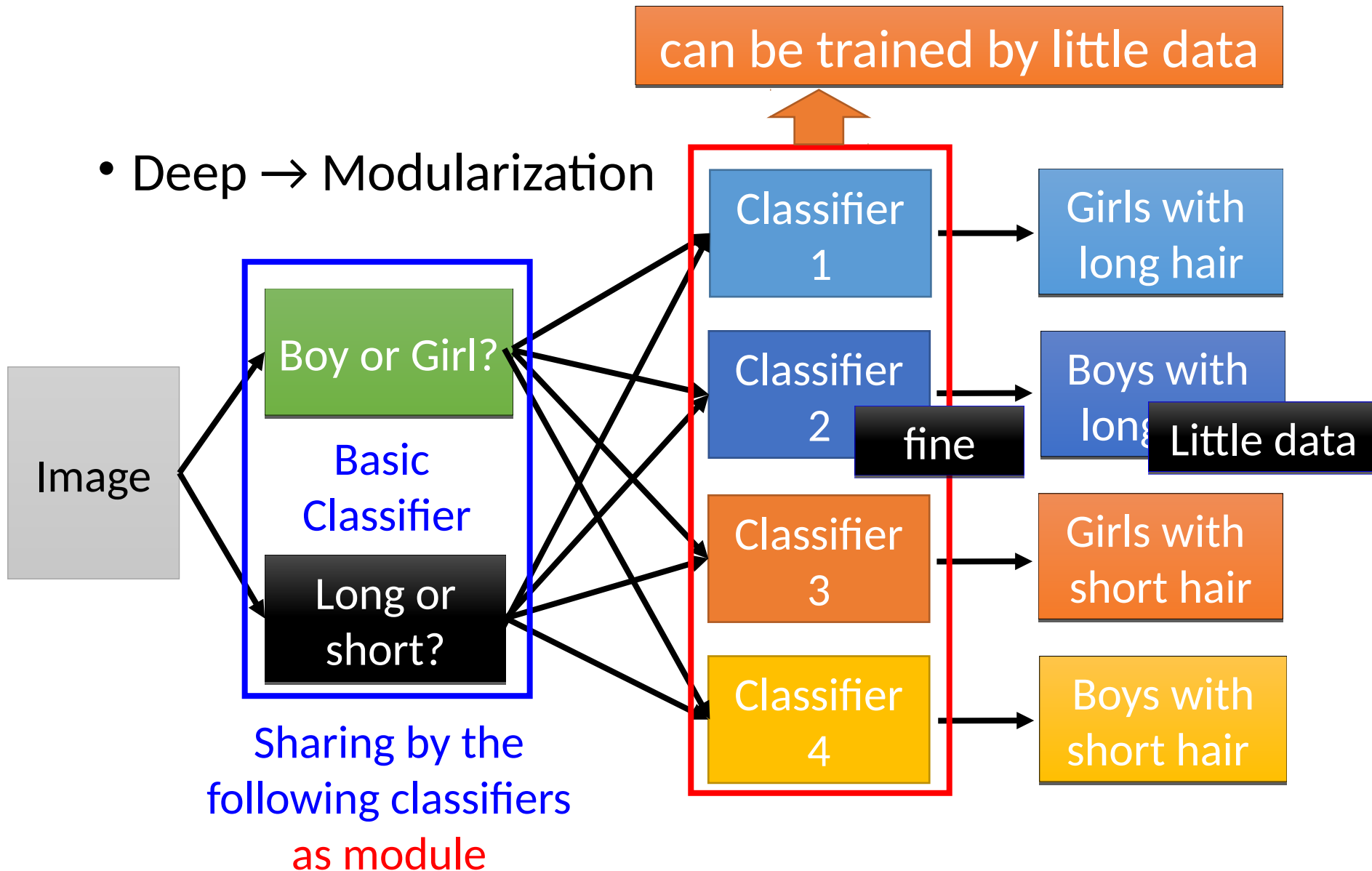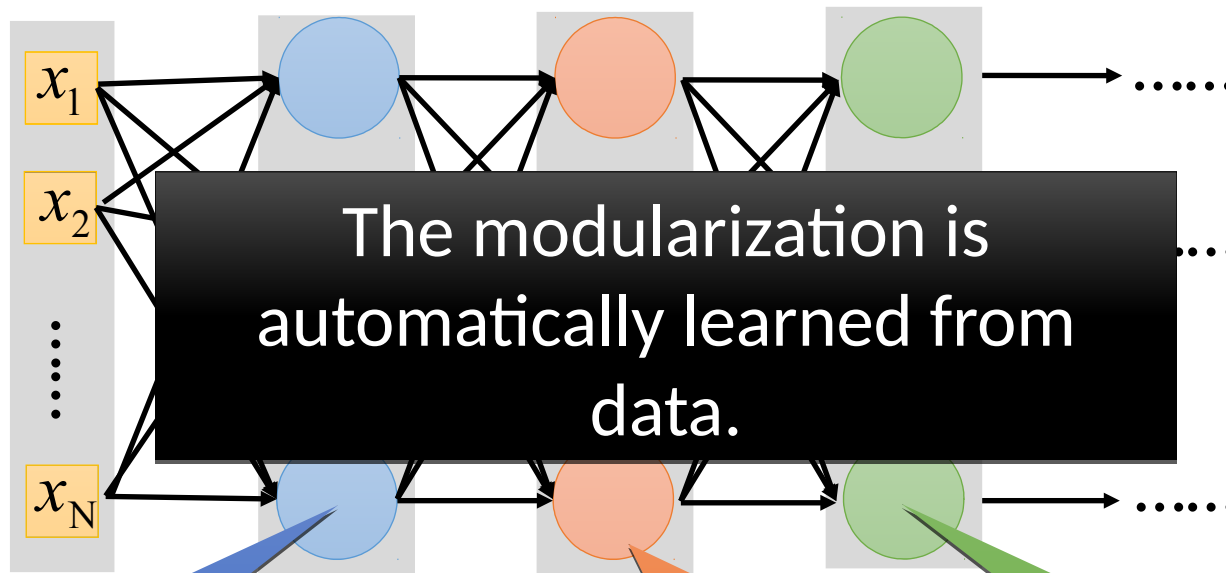
# Why Deep?

- Deep → Modularization

# Why Deep?

can be trained by little data

- Deep → Modularization



Image

Boy or Girl?

Basic Classifier

Long or short?

Sharing by the following classifiers as module

Classifier 1 → Girls with long hair

Classifier 2 → Boys with long

fine

Little data

Classifier 3 → Girls with short hair

Classifier 4 → Boys with short hair

# Why Deep?

- Deep → Modularization    → Less training data?



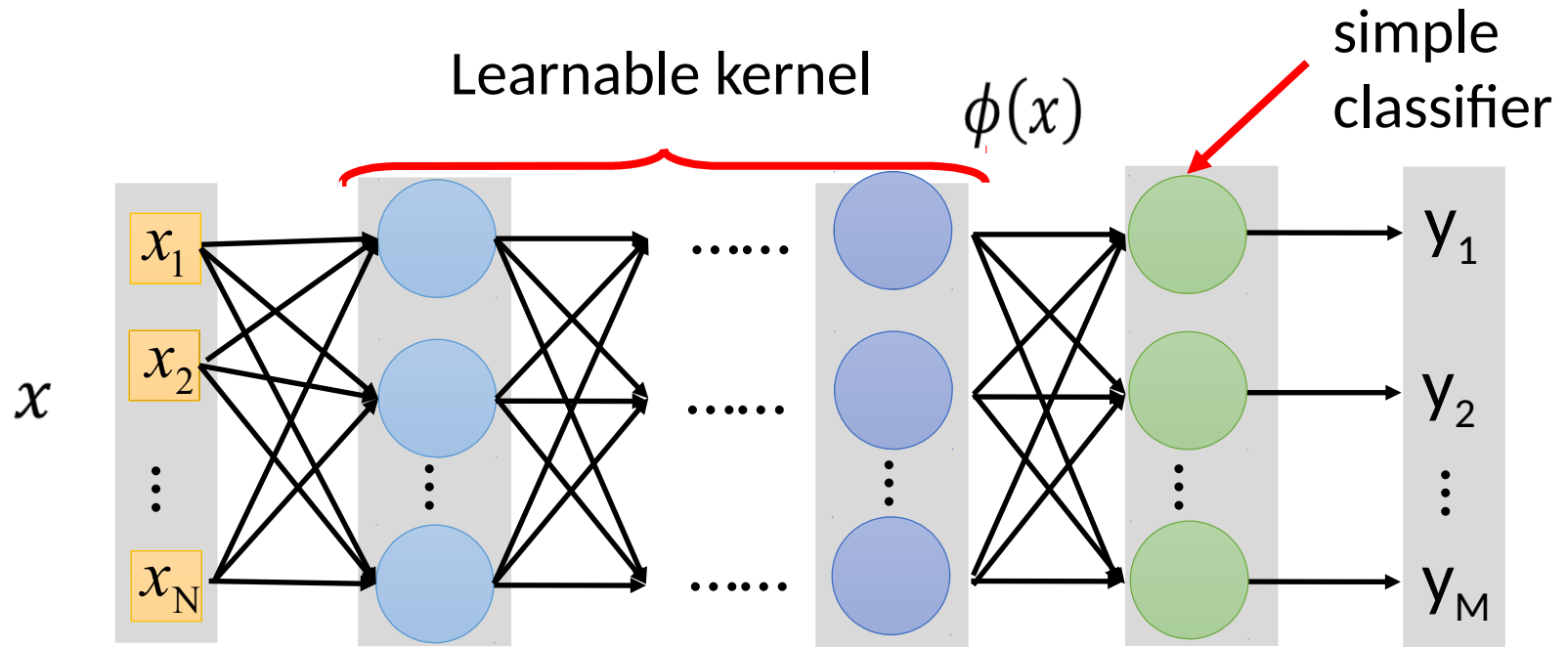The most basic classifiers

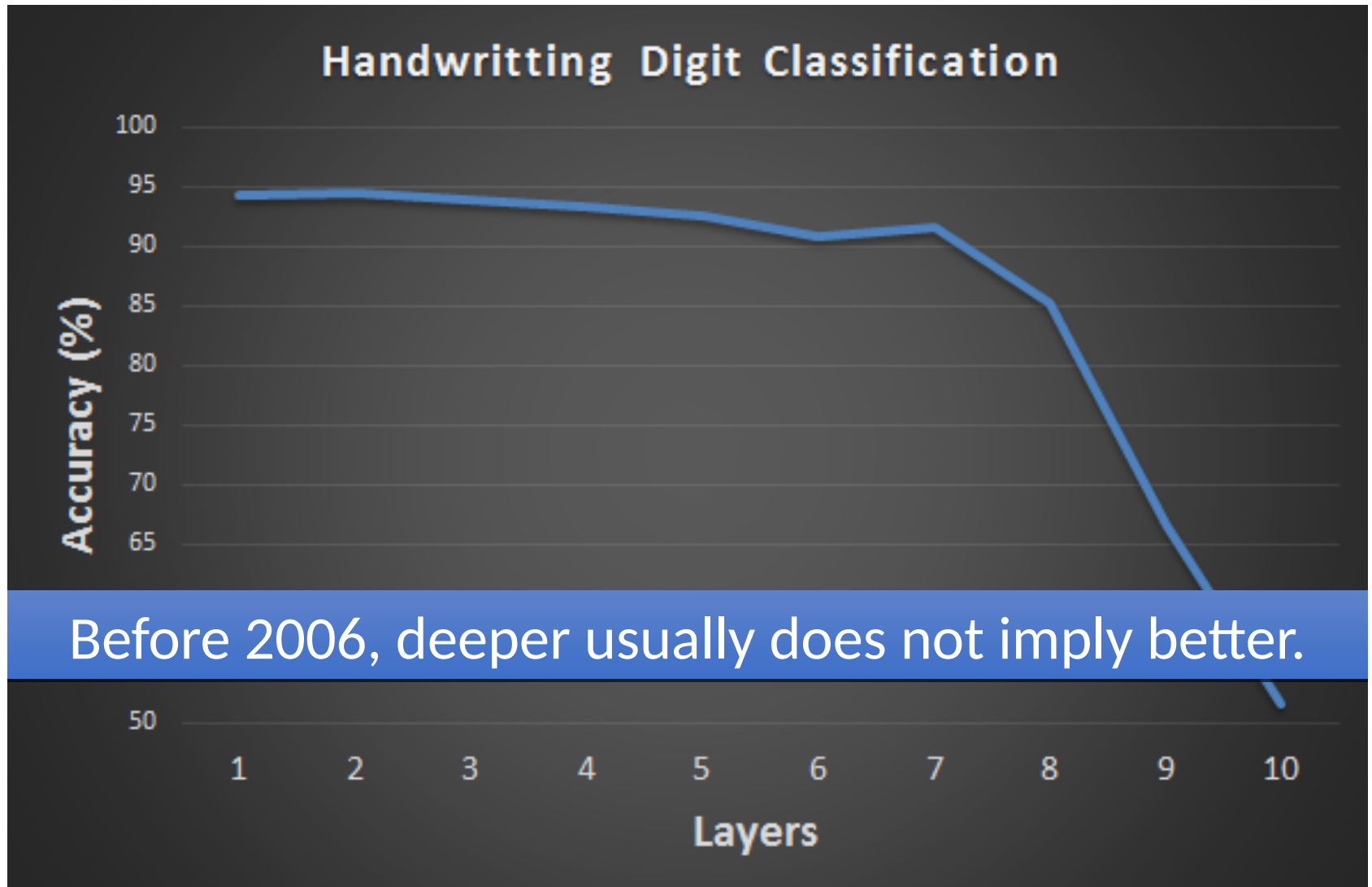Use 1st layer as module to build classifiers

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

The modularization is automatically learned from data.

*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

$x$  $x_1$  $x_2$  $\vdots$  $x_N$  ......  $y_1$  $y_2$  $\vdots$  $y_M$

# Hard to get the power of Deep



**Handwritting Digit Classification**

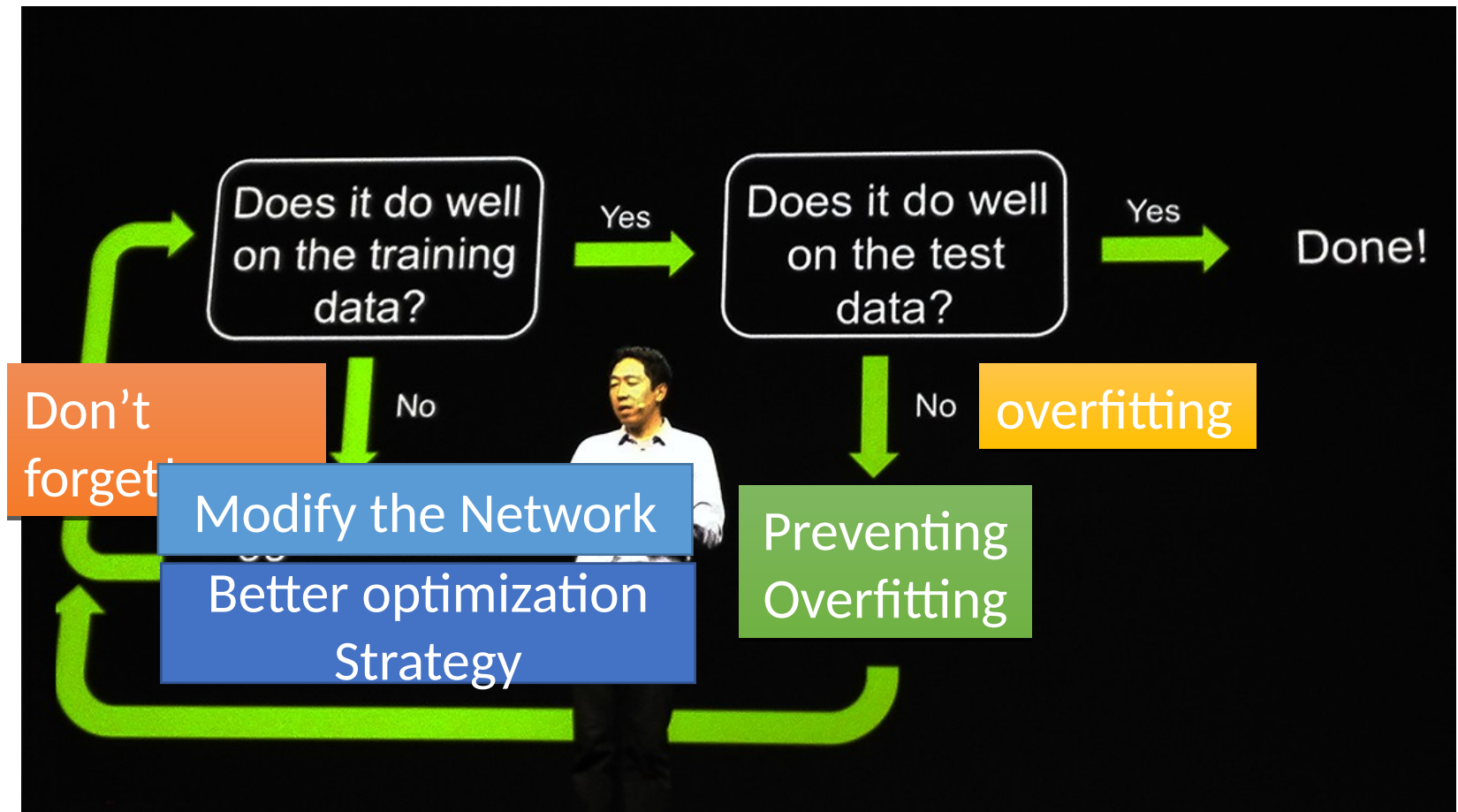Before 2006, deeper usually does not imply better.

# Recipe for Learning



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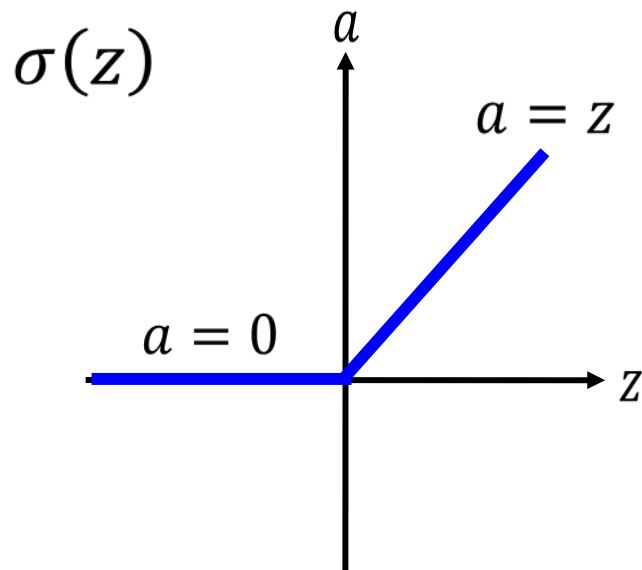
# Recipe for Learning



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

# Training DNN

- New Activation Function

- Adaptive Learning Rate

- Network Regularization
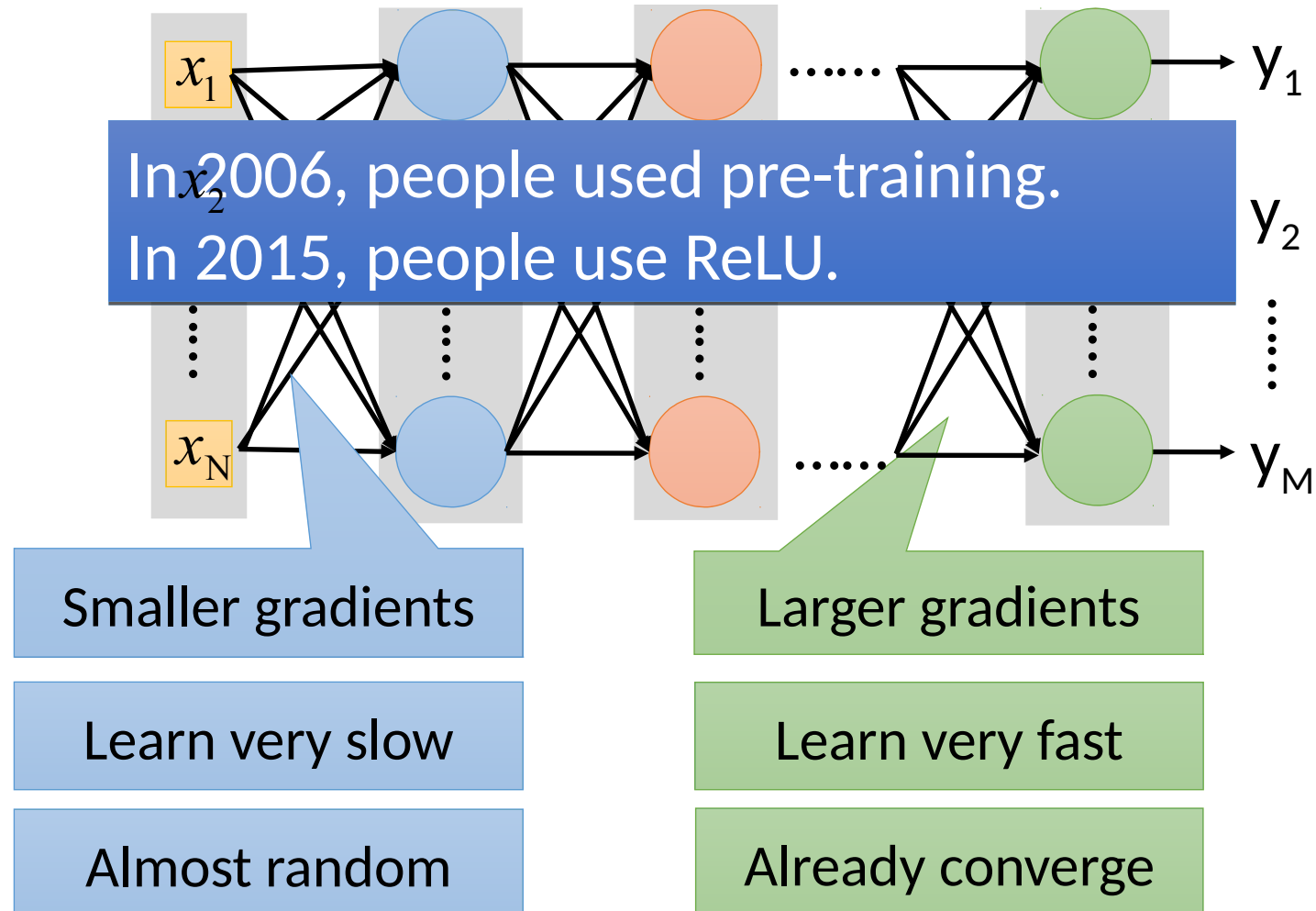
# ReLU

- Rectified Linear Unit (ReLU)

$\sigma(z)$



$a$

$a = z$

$a = 0$

$z$

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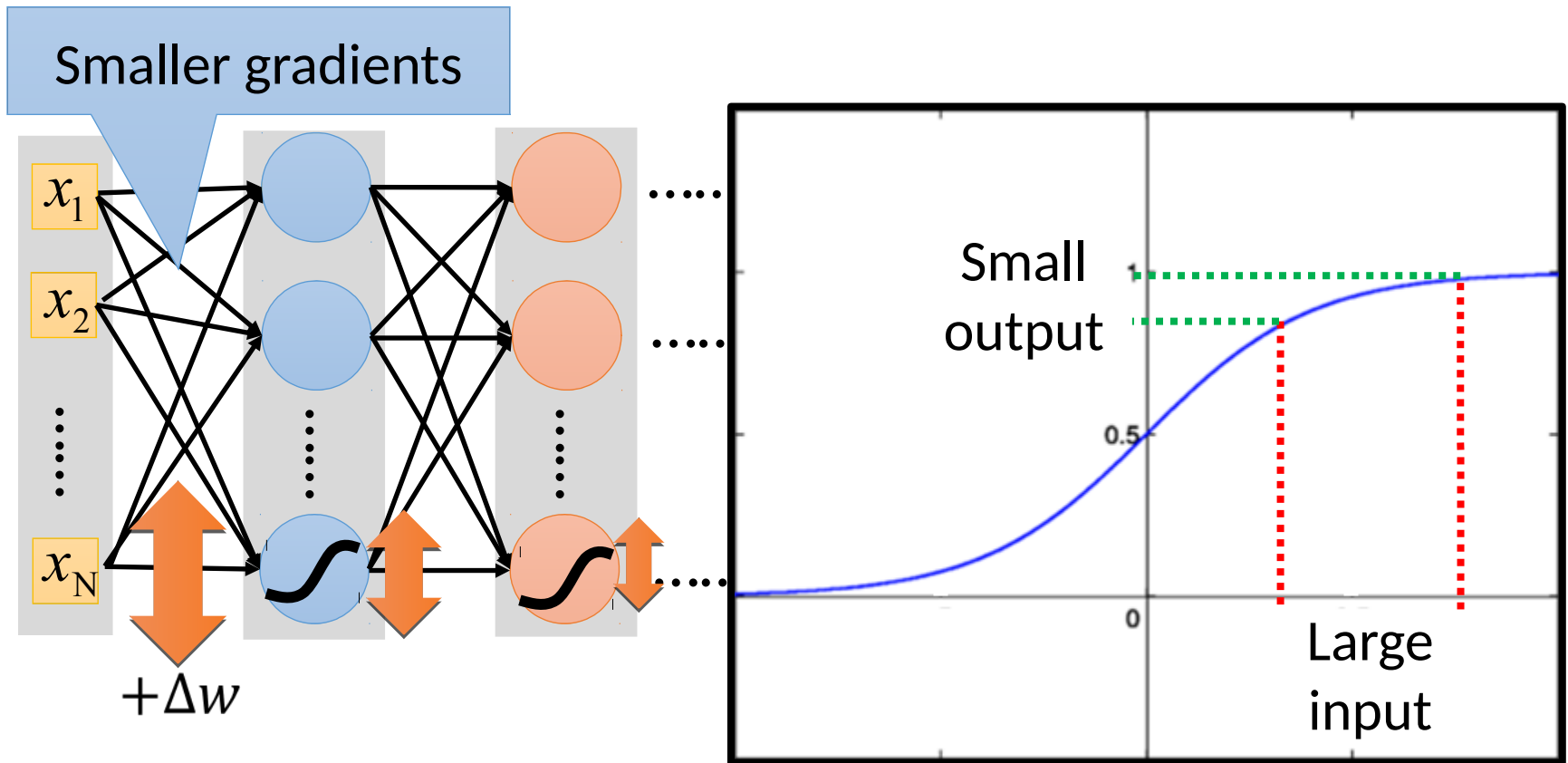[Kaiming He, arXiv'15]

***Reason:***

1. Fast to compute

2. Biological reason

3. Infinite sigmoid with different biases

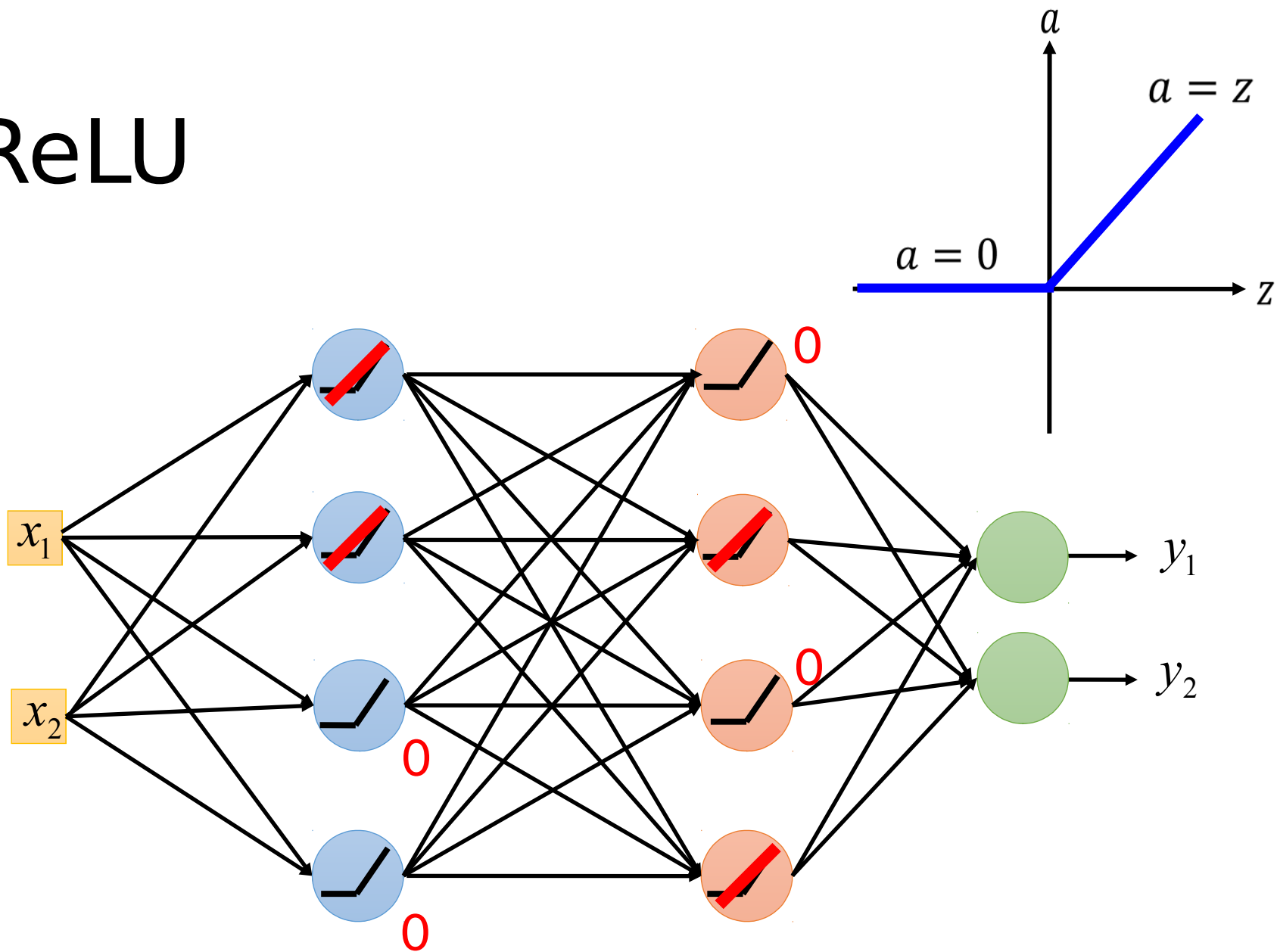4. Vanishing gradient problem

# Vanishing Gradient Problem



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

Smaller gradients

Learn very slow

Almost random

Larger gradients

Learn very fast

Already converge

# Vanishing Gradient Problem



Smaller gradients

$x_1$

$x_2$

$x_N$

$+\Delta w$

Small output

Large input

Intuitive way to compute the gradient …

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

# ReLU

# ReLU



A Thinner linear network

Do not have smaller gradients

# Learning Rate

Set the learning rate η carefully



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

If learning rate is too large

Cost may not decrease after each update

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

$\theta^0$

# Learning Rate

Can we give different parameters different learning rates?



If learning rate is too large

Cost may not decrease after each update

If learning rate is too small

Training would be too slow

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

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

$\theta^0$

# 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]

# Regularization: via Dropout

Pick a mini-batch

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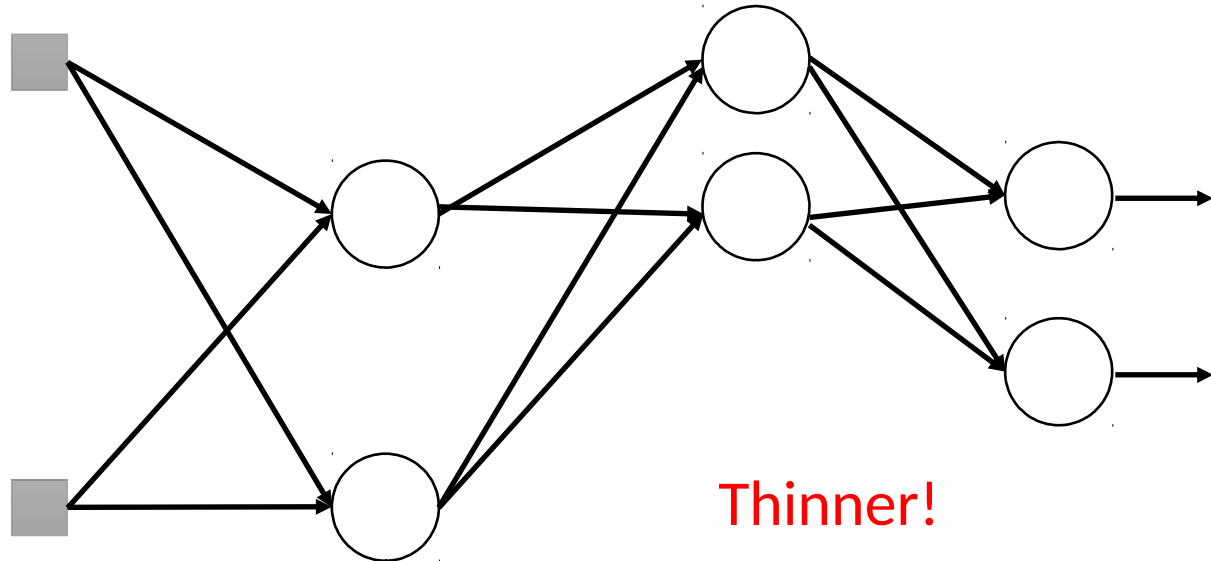
**Training:**



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

# Dropout

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

**Training:**



Thinner!

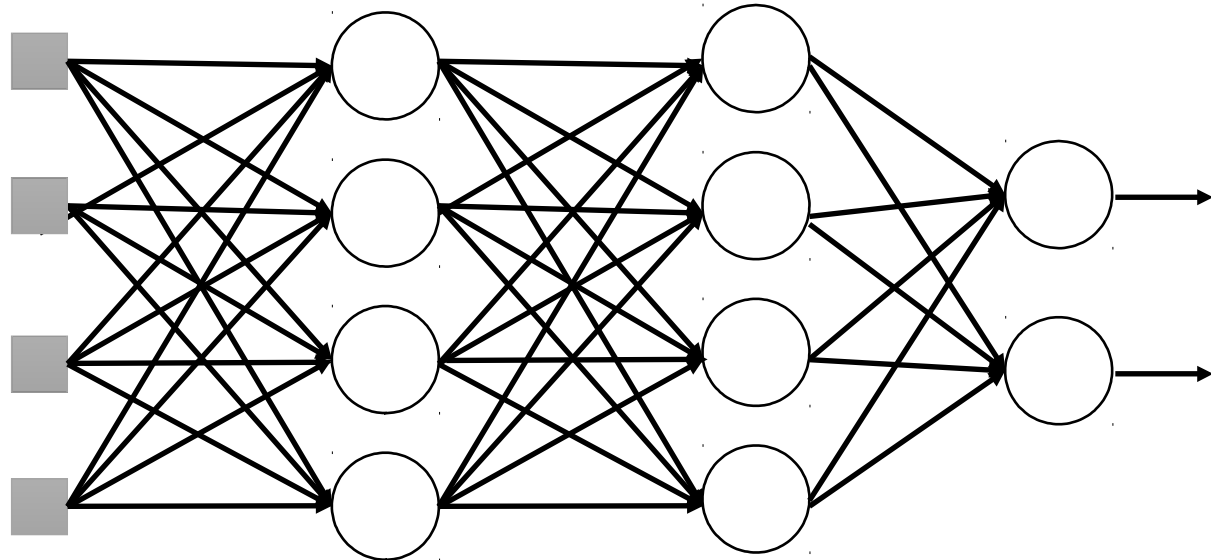➢ **Each time before computing the gradients**

- Each neuron has p% to dropout

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

- Using the new network for training

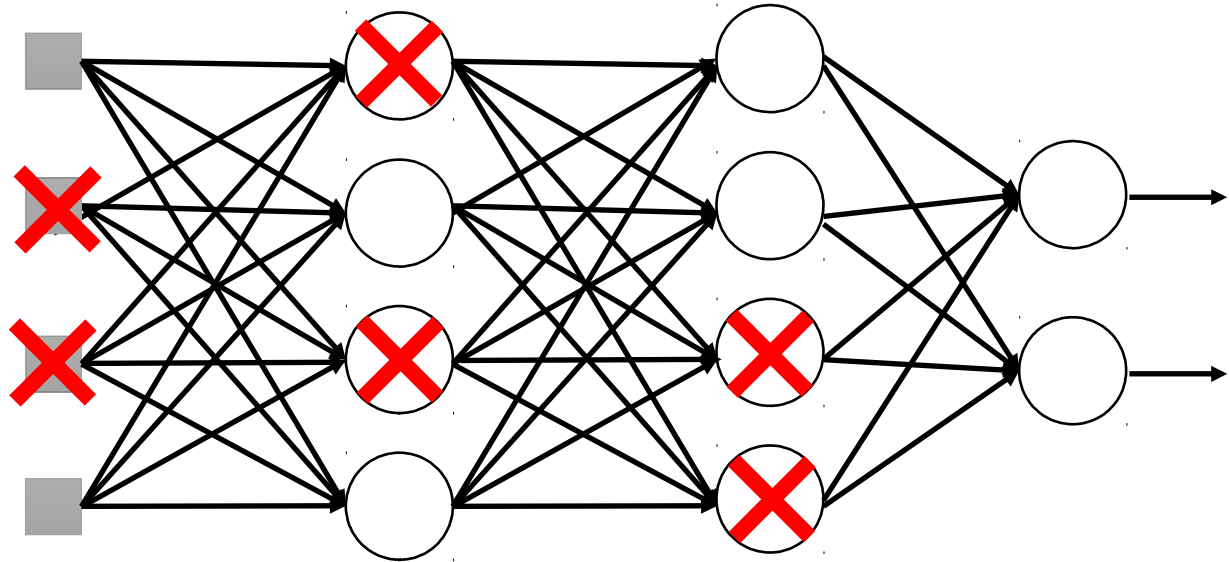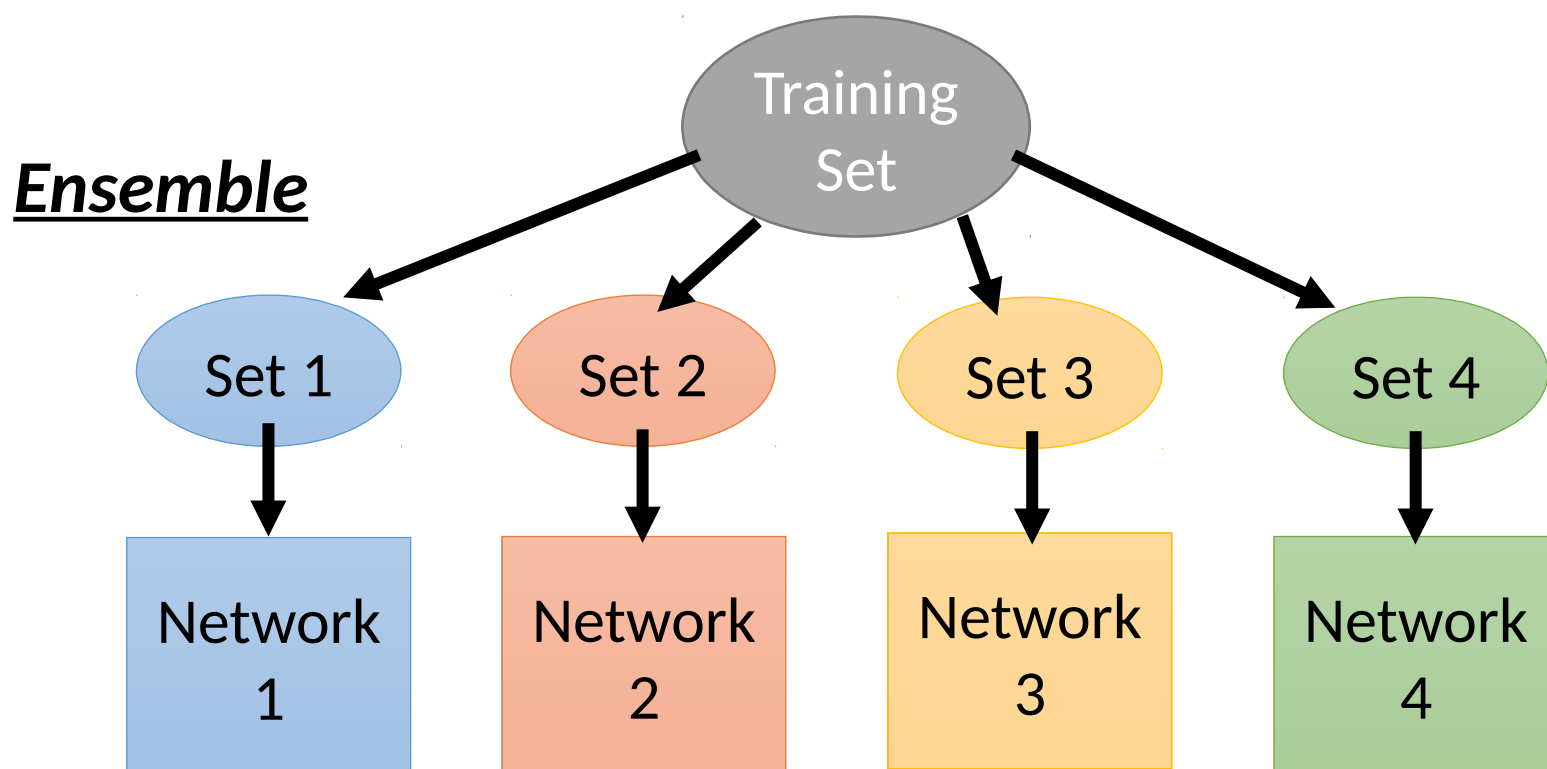For each mini-batch, we resample the dropout neurons

# 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 $w = 1$ 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.

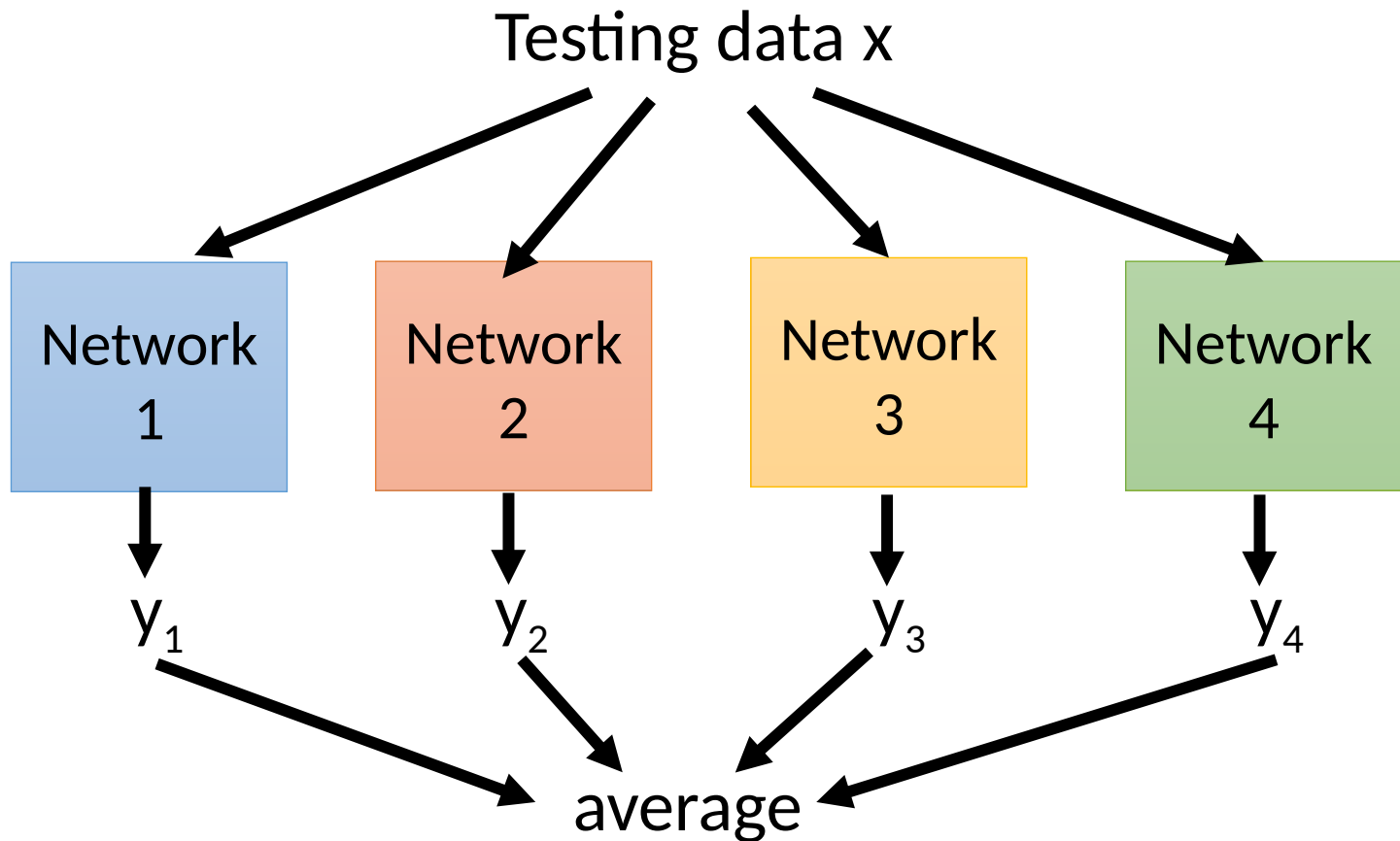➢ When testing, no one dropout actually, so obtaining good results eventually.
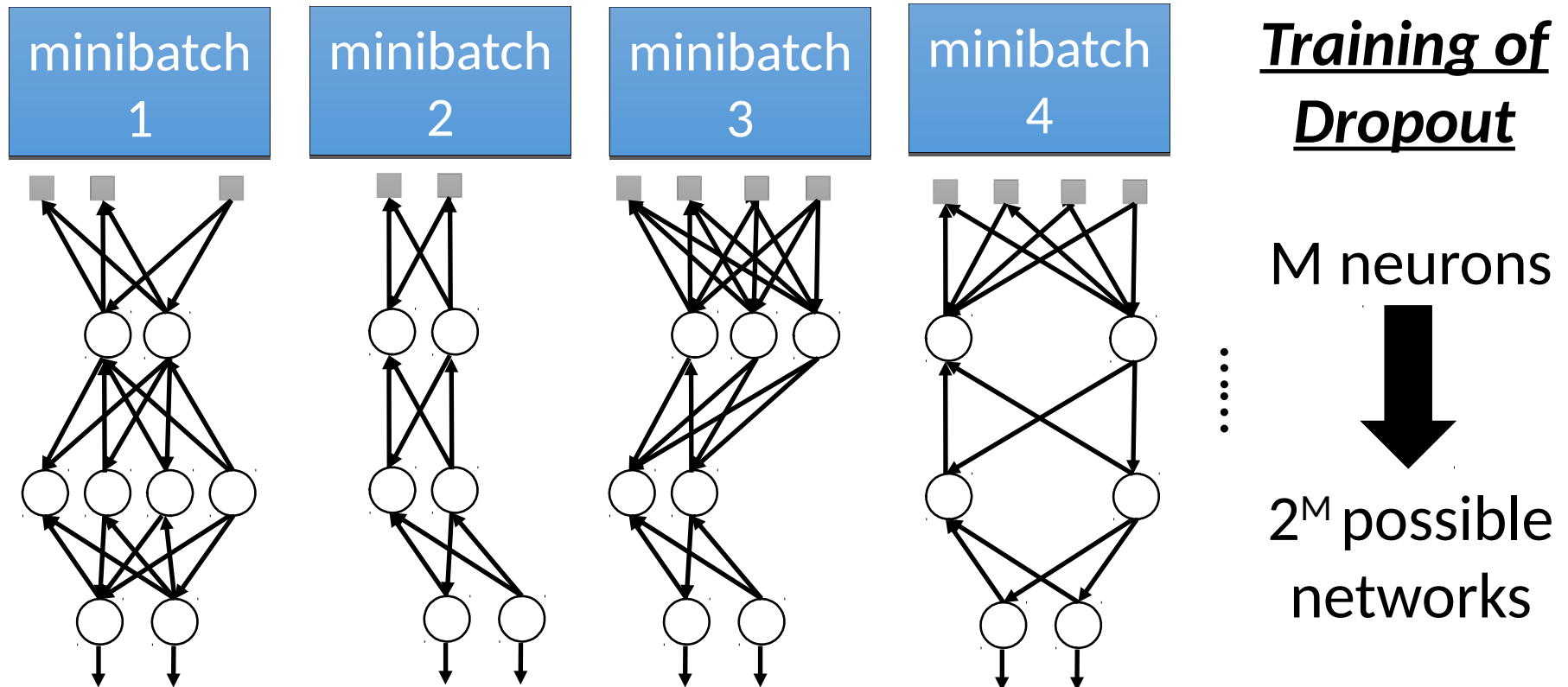
# Dropout is a kind of ensemble.

**_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
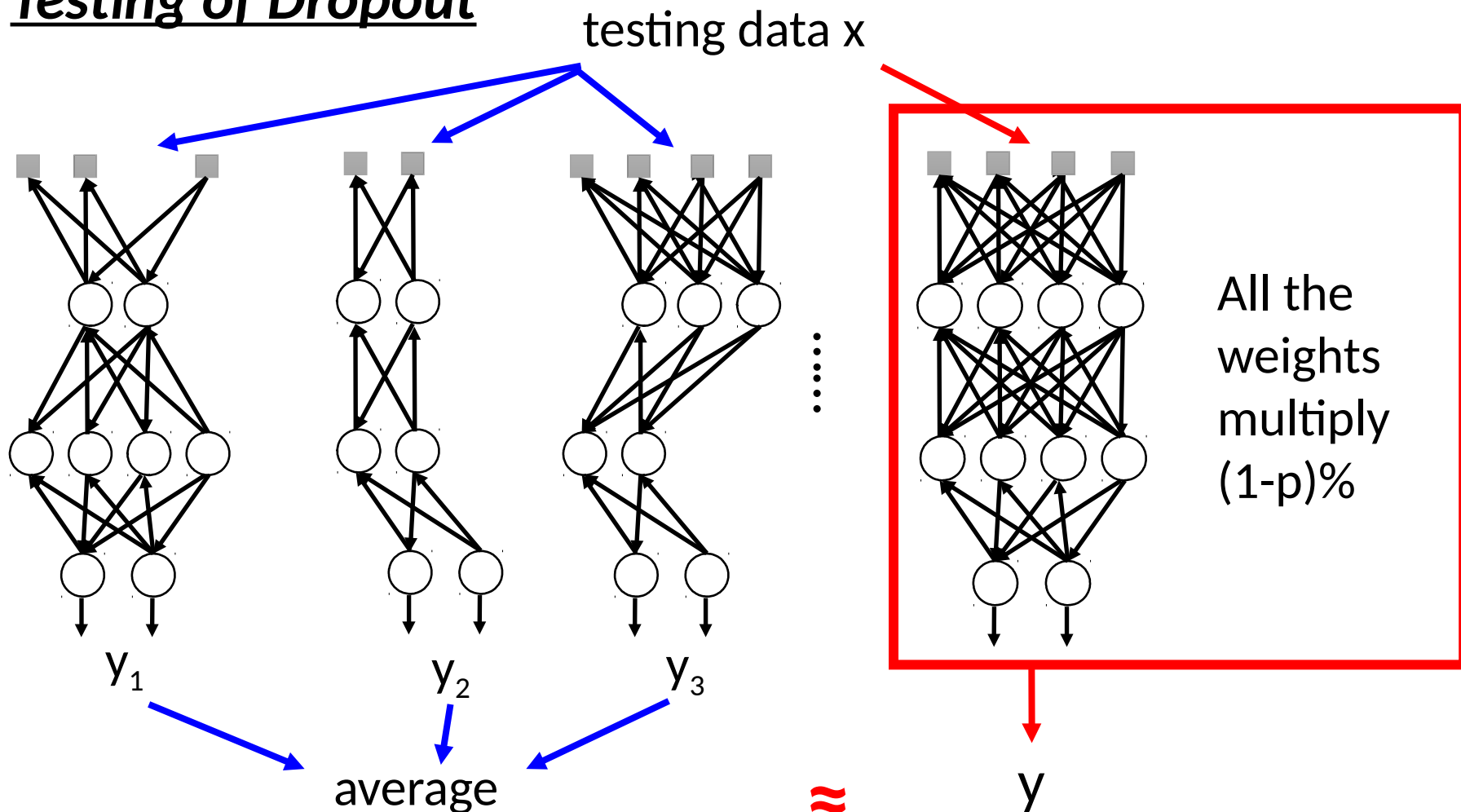
➤Using one mini-batch to train one network
➤Some parameters in the network are shared

# Dropout is a kind of ensemble

testing data x

All the weights multiply (1-p)%

$y_1$

$y_2$

$y_3$

average

≈

y

# Summary

- Introduction to deep learning

    - Fully connected neural networks

- Some training issues and solutions

    - Adaptive learning rate

    - New activation functions

    - Dropout

- Next lectures:

    - Convolution neural networks

    - Recurrent neural networks