

# Introduction to Neural Networks

Part I : Neural information processing

---

Web site of this course: <http://pattern-recognition.weebly.com>.



# Two topics

---

## Neural information processing

- Origins
- Perceptron
- Multilayer perceptron (MLP)
- Convolutional neural network (CNN)

## Training multilayer networks

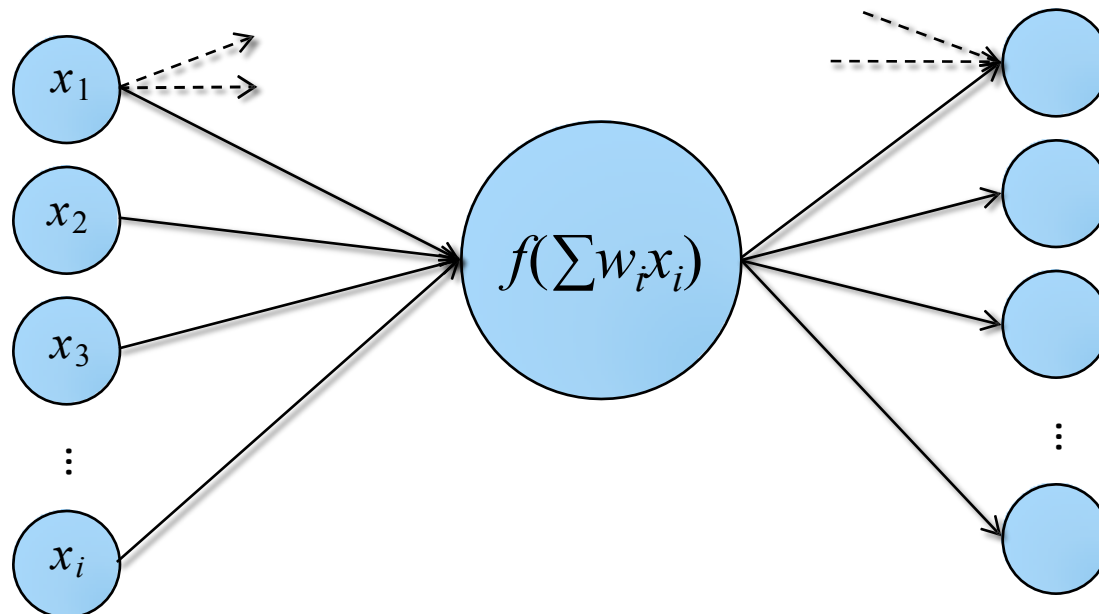
- Perceptron learning
- Optimization basics
- Stochastic gradient descent
- Backpropagation learning algorithms

# Neural Information Processing

- **Origins**
- Perceptron
- Multilayer perceptron
- Convolutional neural network

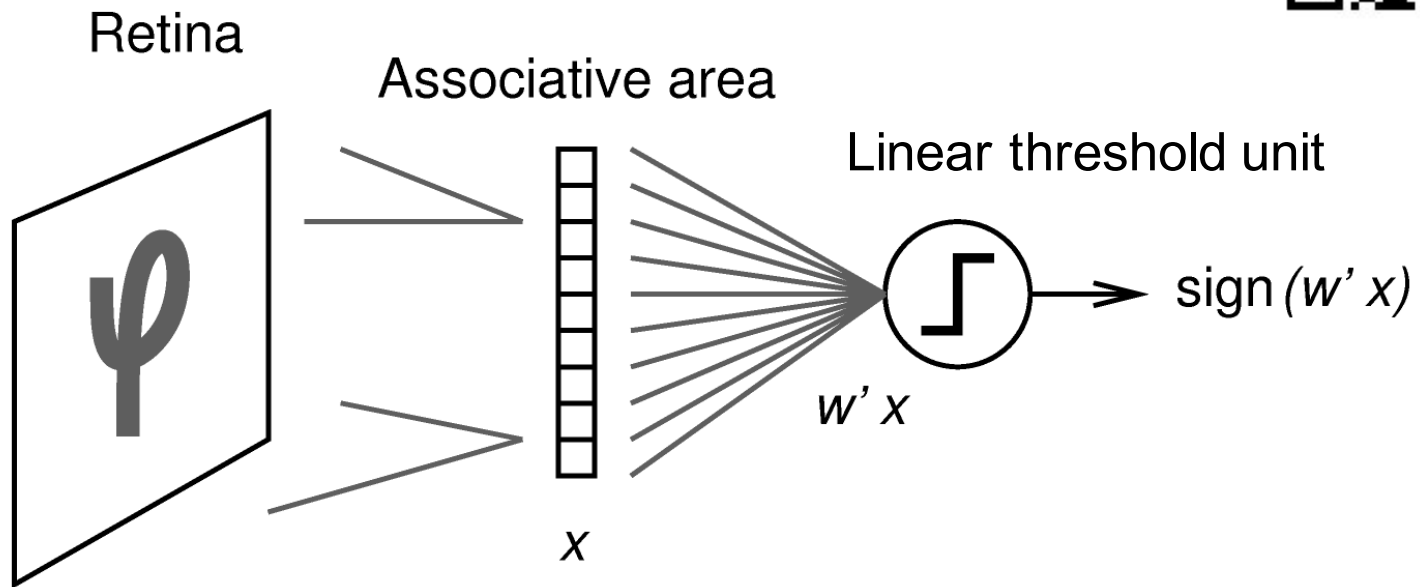
# McCulloch & Pitts (1943)

---



A simplified neuron model: the Linear Threshold Unit.

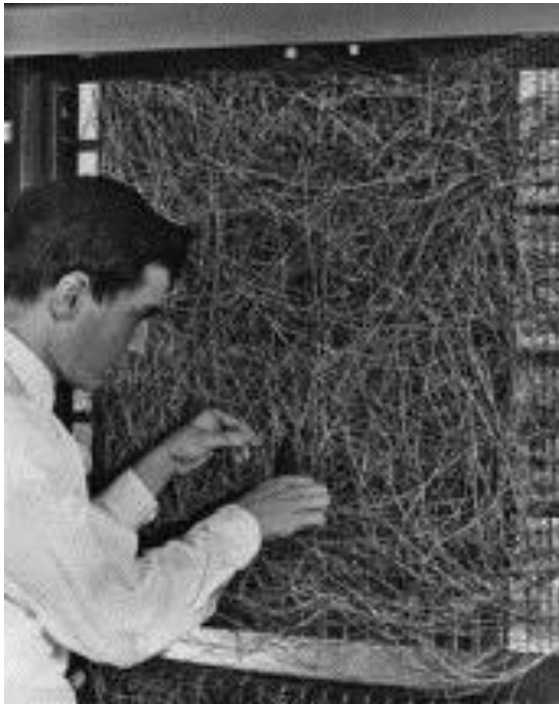
# The perceptron (1957)



Supervised learning of the weights  
using the Perceptron algorithm.

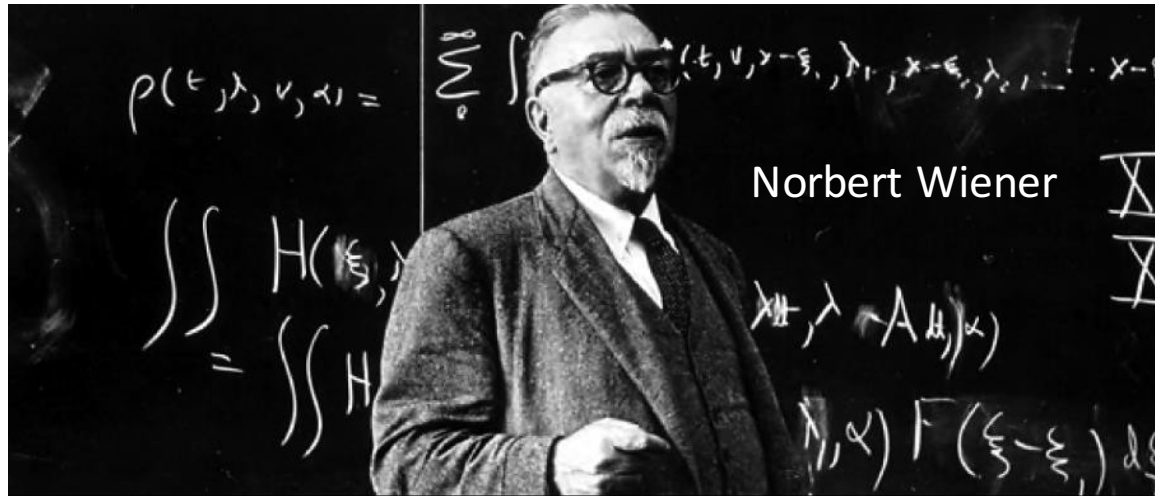
Rosenblatt 1957

# The perceptron is a machine



# Cybernetics (1948)

---



- Mature communication technologies, nascent computing technologies
- Redefining the man-machine boundary

# Two camps to design computers

---

Biological computer



Mathematical computer

$$\frac{\partial}{\partial \theta} \int_{\mathcal{R}_n} T(x) f(x, \theta) dx = \int_{\mathcal{R}_n} \frac{\partial}{\partial \theta} T(x) f(x, \theta) dx$$
$$\frac{\partial}{\partial a} \ln f_{a, \sigma^2}(\xi_1) = \frac{(\xi_1 - a)}{\sigma^2} f_{a, \sigma^2}(\xi_1) = \frac{1}{\sqrt{2\pi\sigma^2}} \left( \frac{\xi_1 - a}{\sigma} \right) e^{-\frac{(\xi_1 - a)^2}{2\sigma^2}}$$
$$\int_{\mathcal{R}_n} T(x) \cdot \frac{\partial}{\partial \theta} f(x, \theta) dx = M \left( T(\xi) \cdot \frac{\partial}{\partial \theta} \ln L(\xi, \theta) \right)$$
$$\int_{\mathcal{R}_n} T(x) \cdot \left( \frac{\partial}{\partial \theta} \ln L(x, \theta) \right) \cdot f(x, \theta) dx = \int_{\mathcal{R}_n} T(x) \cdot \left( \frac{\partial}{\partial \theta} \ln L(x, \theta) \right) f(x, \theta) dx$$
$$\frac{\partial}{\partial \theta} \int_{\mathcal{R}_n} T(x) f(x, \theta) dx = \int_{\mathcal{R}_n} \frac{\partial}{\partial \theta} T(x) f(x, \theta) dx$$

- Which model to emulate : brain or mathematical logic ?
- **Mathematical logic won.**



# Computing with symbols

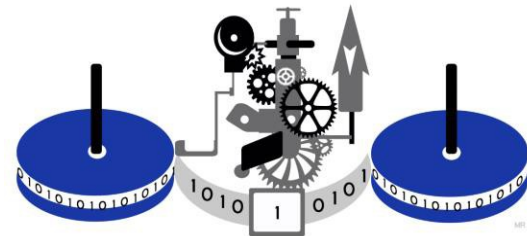
---

## General computing machines

- Turing machine
- von Neumann machine

## Engineering

- Programming  
(reducing a complex task into a collection of simple tasks.)
- Computer language
- Debugging
- Operating systems
- Libraries



```
public int[] extractMessage(int[] res) {
    for (int i = 0; i < MAX_RES_LEN; i++) buf[i] = 0;
    int loc = 0, i = 0;
    while (i < res.length) {
        int j = 0;
        while (j < loc; j++) res[j] = buf[j];
        return res;
    }
    int loc = 0, i = 0;
    while (i < res.length) {
        int j = 0;
        while (j < loc; j++) res[j] = buf[j];
        return res;
    }
}

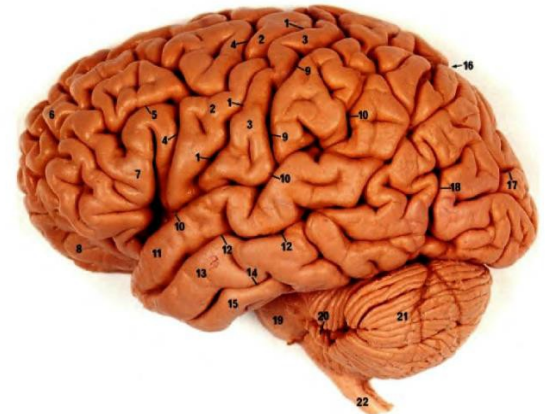
public int[] decodeMessage(int[] res) {
    for (int i = 0; i < MAX_RES_LEN; i++) buf[i] = 0;
    int loc = 0, i = 0;
    while (i < res.length) {
        int j = 0;
        while (j < loc; j++) res[j] = buf[j];
        return res;
    }
}
```

# Computing with the brain

---

## An engineering perspective of brain

- Compact
- Energy efficient (20 watts)
- $10^{12}$  Glial cells (power, cooling, support)
- $10^{11}$  Neurons (soma + wires)
- $10^{14}$  Connections (synapses)
- Volume = 50% glial cells + 50% wires.



## Could brain be a general computing machine?

- Basically, brain is
  - Slow for mathematical logic, arithmetic, etc.
  - Very fast for vision, speech, language, social interactions, etc.
- Because of brain evolution : vision -> language -> logic.

# Success stories

---

## Record performance

- MNIST (1988, 2003, 2012)
- ImageNet (2012)
- ...



MNIST



ImageNet

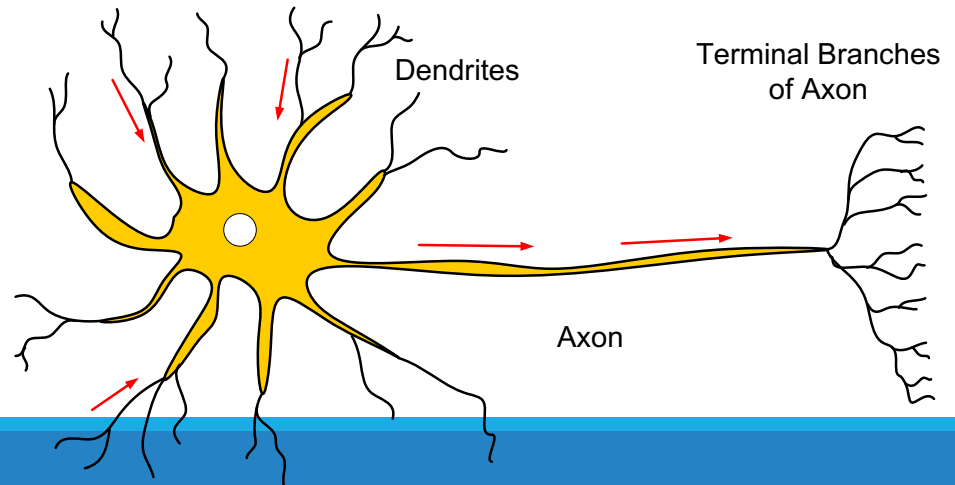
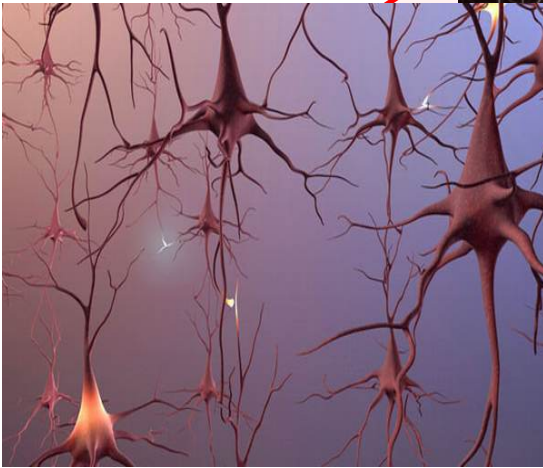
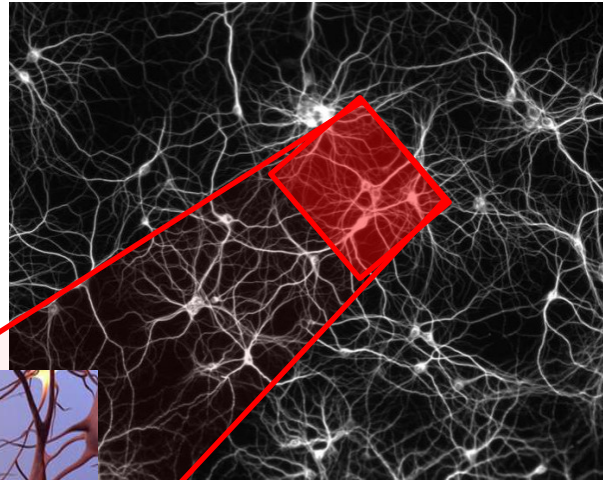
## Real applications

- Check reading (AT&T Bell Labs, 1995 – 2005)
- Optical character recognition (Microsoft OCR, 2000)
- Cancer detection from medical images (NEC, 2010)
- Object recognition (Google and Baidu's photo taggers, 2013)
- Speech recognition (Microsoft, Google, IBM switched in 2012)
- Natural Language Processing (NEC 2010)
- ...

# Neural Information Processing

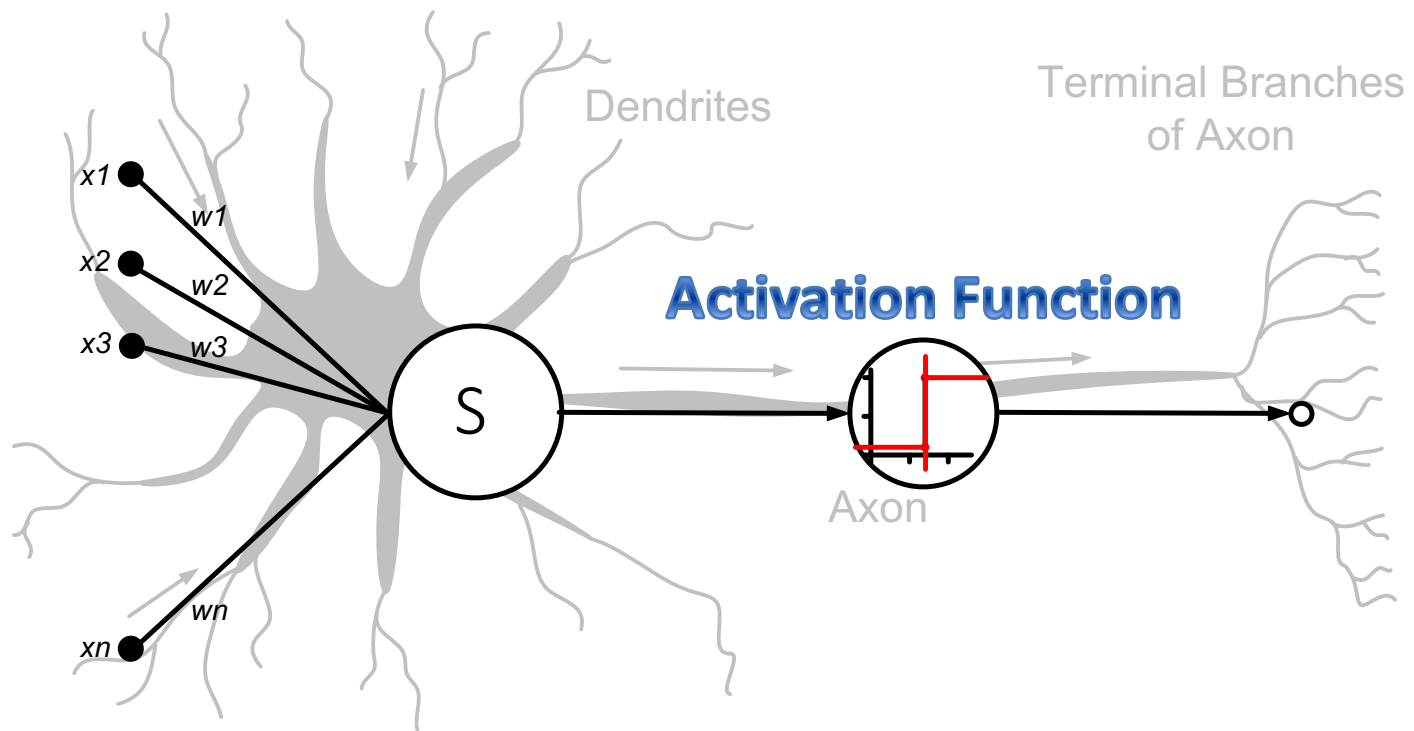
- Origins
- **Perceptron**
- Multilayer perceptron
- Convolutional neural network

# Biological Neuron



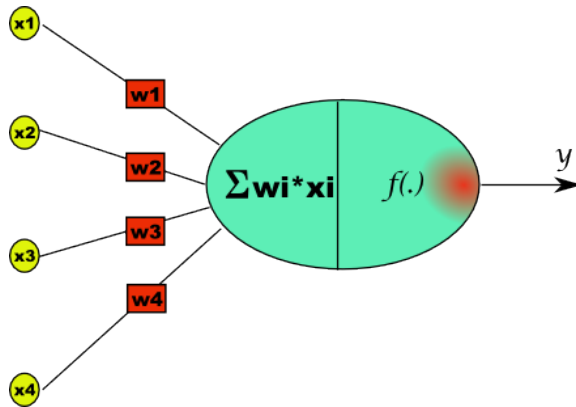
# Artificial Neuron

---

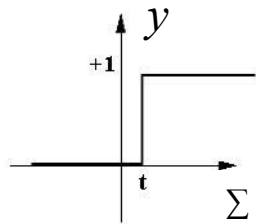


# Neuron's Mathematical Model

A neuron



$f$ : Transfer  
(Activation)  
Functions



Step Function

$$f(w^T x) = \begin{cases} 1, & \text{if } w^T x > t \\ 0, & \text{otherwise} \end{cases}$$

Without  $f$

$$y = \sum_{i=1}^4 w_i x_i = w^T x$$

A neuron is a "linear equation"

With  $f$

$$y = f(w^T x)$$

$$= \begin{cases} C1 \text{ (label 1)} & \text{if } w^T x > t \\ C2 \text{ (label 0)} & \text{if } w^T x \leq t \end{cases}$$

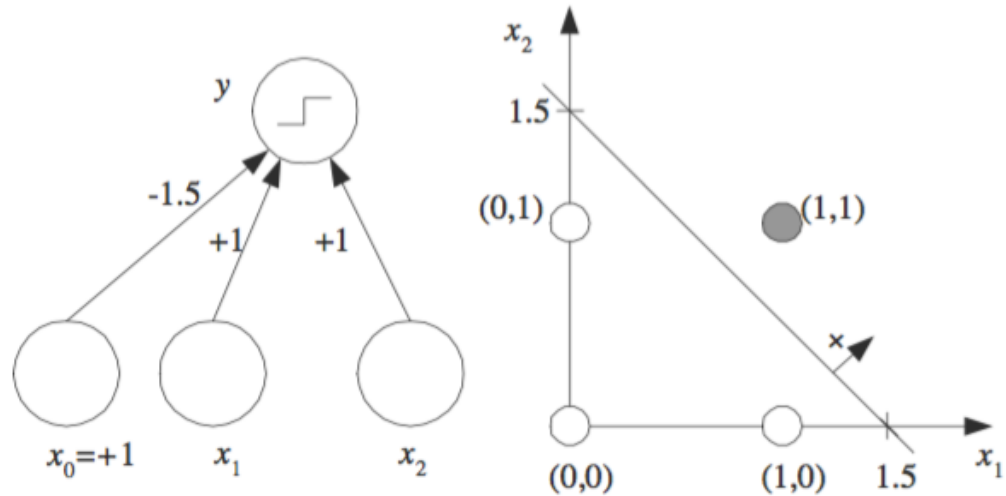
A neuron is a "linear classifier"

# "And" function Classification

---

$x_1$	$x_2$	$r$
0	0	0
0	1	0
1	0	0
1	1	1

Input and output for the AND function.



The perceptron that implements AND and its geometric interpretation.

- The discriminant is  $y = f(x_1 + x_2 - 1.5)$ .

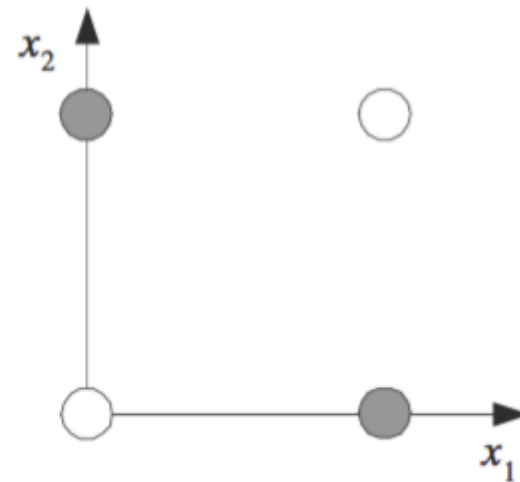


# "XOR" Function Classification

---

$x_1$	$x_2$	$r$
0	0	0
0	1	1
1	0	1
1	1	0

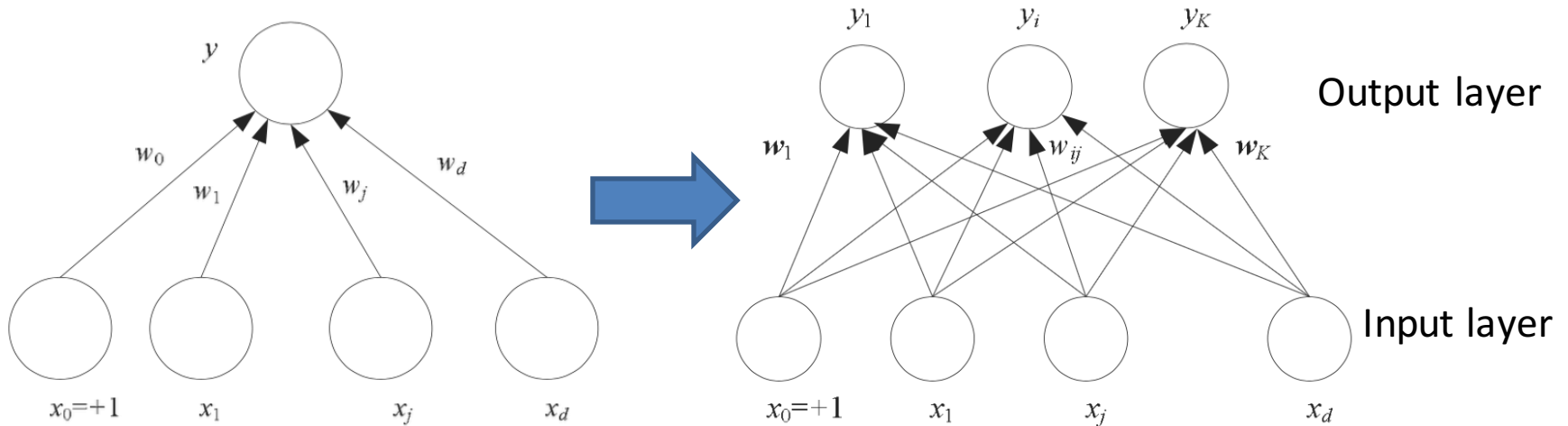
Input and output  
for the XOR function



XOR problem is not linearly separable.  
We cannot draw a line where  
the empty circles are on one side and  
the filled circles on the other side.

**Perceptron can not solve the XOR classification problem**

# The General Perceptron Model



$$o = \sum_{j=1}^d w_j x_j + w_0 = w^T x$$

$$w = [w_0, w_1, \dots, w_d]^T$$

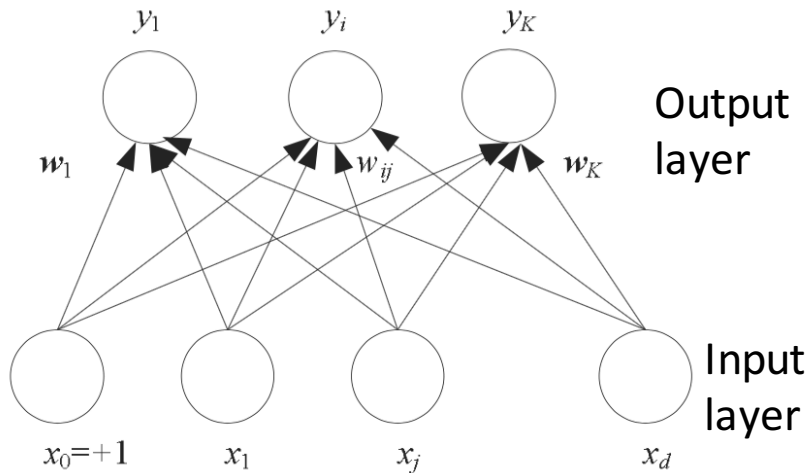
$$x = [1, x_1, \dots, x_d]^T$$

$$y = \begin{cases} 1, & \text{if } f(o) > t \\ 0, & \text{if } f(o) \leq t \end{cases} \quad (t \text{ is } 0 \text{ usually})$$

$$o_i = w_i^T x = \sum_{j=1}^d w_{ij} x_j + w_{i0}$$

$$y_i = \begin{cases} 1, & \text{if } f(o_i) > 0 \\ 0, & \text{if } f(o_i) \leq 0 \end{cases}$$

# Perceptron Model for Classification



## For a K-class classification problem

- For a given  $x$  with unknown class
- $x \in \text{class } k$ , if  $y_k = 1, y_j = 0$  for all  $j \neq k$
- Only one  $y_k$  can be 1

$$o_i = w_i^T x = \sum_{j=1}^d w_{ij} x_j + w_{i0}$$

$$y_i = \begin{cases} 1, & \text{if } f(o_i) > t_i \\ 0, & \text{if } f(o_i) \leq t_i \end{cases}$$

## A simplified K-class classification problem

- $x \in \text{class } k$ , if  $y_k = \max_i y_i$
- $y_i = w_i^T x = \sum_{j=1}^d w_{ij} x_j + w_{i0}$

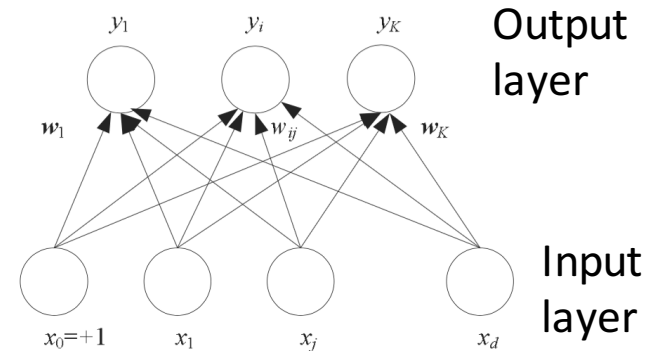
# Training the perceptron

## Testing for K-class classification problem

- For a given  $x$  with unknown class
- $x \in \text{class } k$ , if  $y_k = \max_i y_i$
- $y_i = w_i^T x = \sum_{j=1}^d w_{ij} x_j + w_{i0}$

That is

- A  $W$  represents a perceptron
- Given a  $W$ , then we can classify a pattern  $x$



$$W = [w_1, w_2, \dots, w_K]$$
$$= \begin{bmatrix} w_{11} & \cdots & w_{K1} \\ \vdots & \ddots & \vdots \\ w_{1d} & \cdots & w_{Kd} \end{bmatrix}$$

A **Machine Learning** problem:

how to obtain the  $W$  of a perceptron

- We need a set of training patterns  $(X, Y)$
  - We need a learning algorithm to learn  $W$  by  $(X, Y)$
- => **Perceptron learning algorithm  $A$ :  $W=A(X, Y)$**

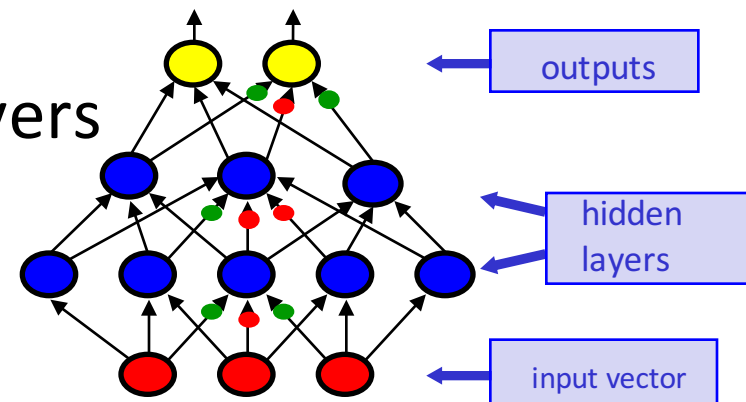
# Neural Information Processing

- Origins
- Perceptron
- **Multilayer perceptron (MLP)**
- Convolutional neural network

# Multi-Layer Perceptron (MLP)

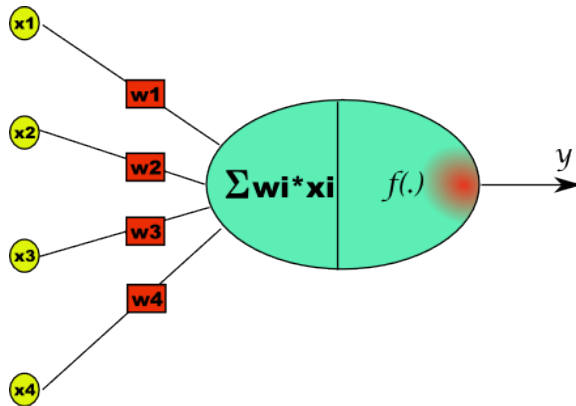
---

- Extension of perceptron from one layer into multiple layers
- Three differences compared with perceptron
  - Transfer(activation) function
    - From step function to **sigmoid function**
  - Layer
    - Add **hidden layer**
  - Training algorithm
    - **Backpropagation** learning algorithm to learn the weights



# Neuron's Mathematical Model

Neuron:



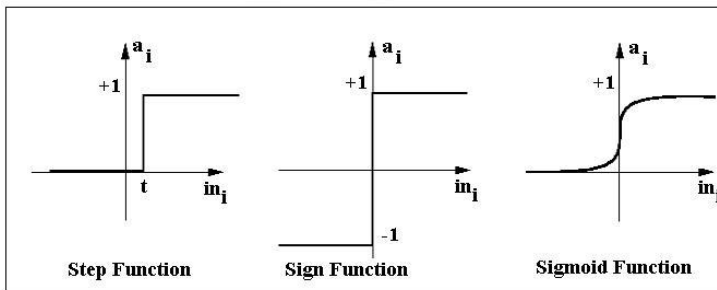
$f$  is a *step* function

$$y = f(w^T x)$$

$$= \begin{cases} C1 \text{ (label 1) if } w^T x > t \\ C2 \text{ (label 0) if } w^T x \leq t \end{cases}$$

A neuron is a "**linear classifier**"

$f$ : Transfer  
(Activation)  
Functions



$f$  is a *sigmoid* function

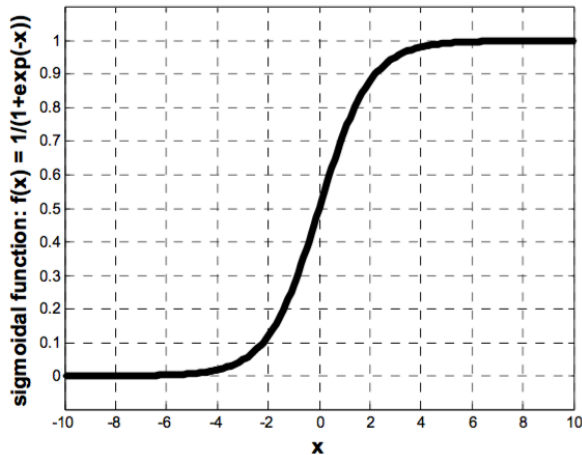
$$y = f(w^T x)$$

$$= \begin{cases} C1 \text{ if } f(w^T x) > 0.5 \\ C2 \text{ if } f(w^T x) \leq 0.5 \end{cases}$$

A neuron is still a "**linear classifier**"

# Transfer functions : Sigmoid

---

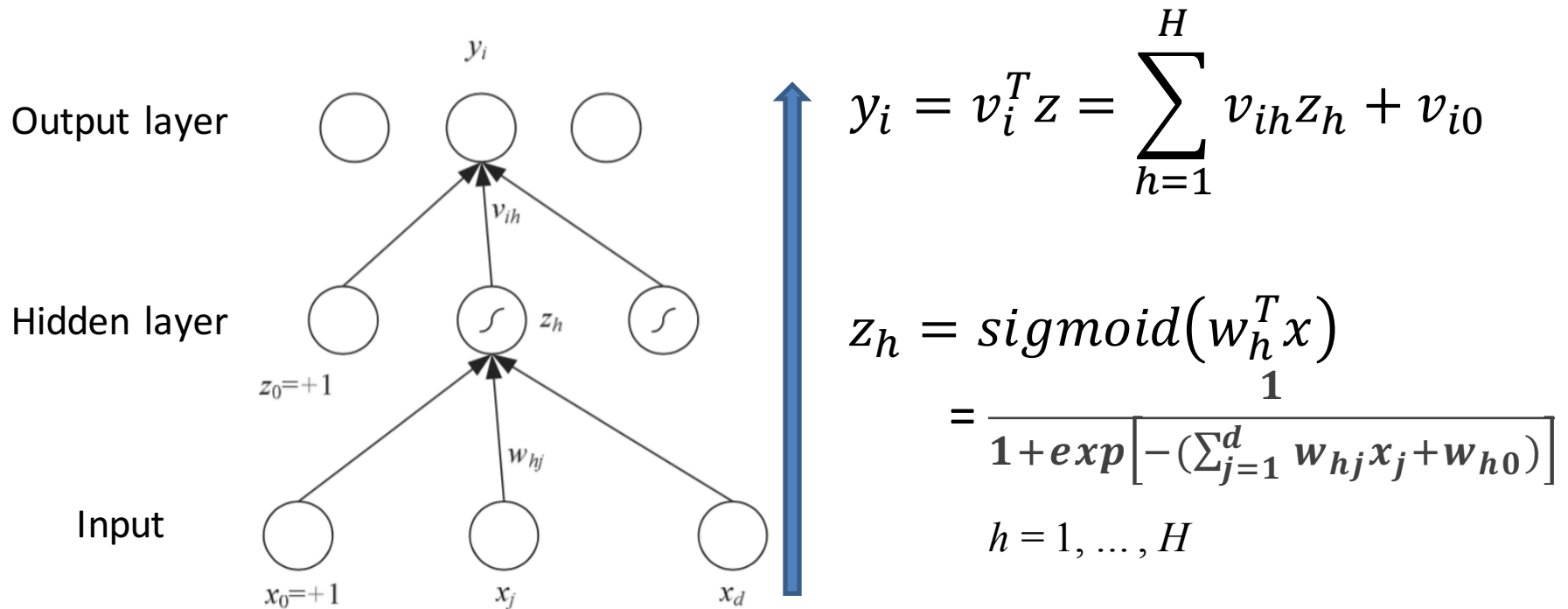


$$y = f(x) = \frac{1}{1 + e^{-x}}$$

$$y = f(w^T x) \\ = \frac{1}{1 + e^{-w^T x}}$$



# Propagation from input to output

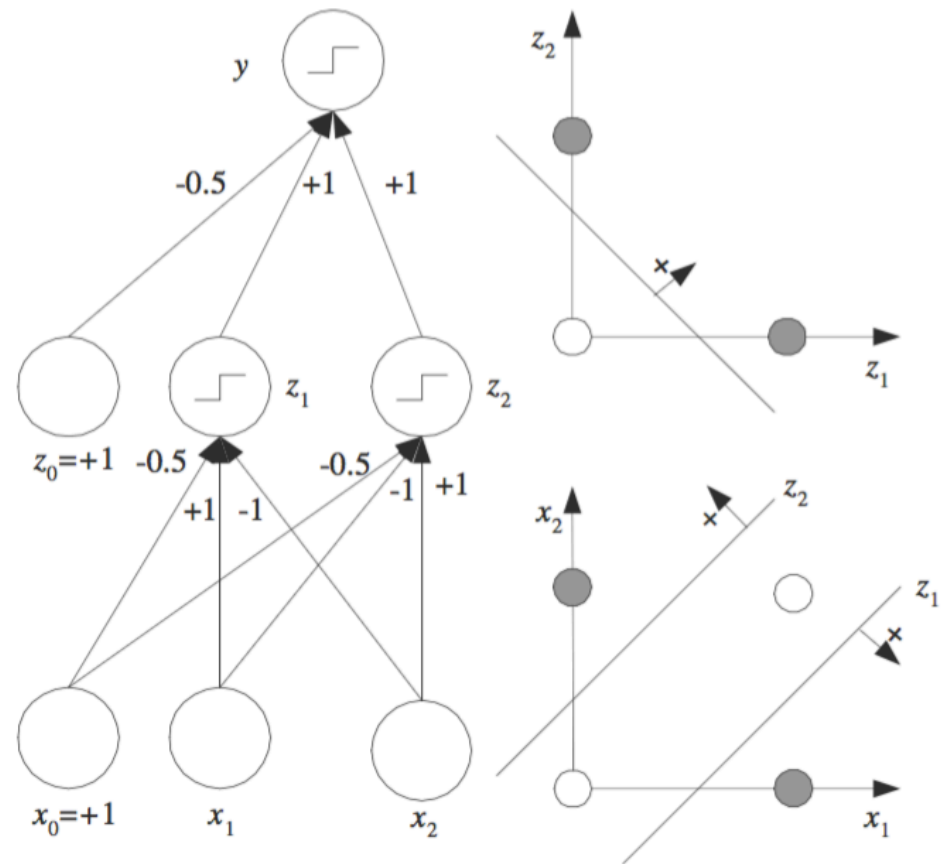


# MLP as a universal approximator

- The MLP that solves the XOR problem

$$\begin{aligned}
 & x_1 \text{ XOR } x_2 \\
 = & (x_1 \text{ AND } \sim x_2) \\
 & \text{OR } (\sim x_1 \text{ AND } x_2)
 \end{aligned}$$

- The hidden units and the output have the threshold activation function with threshold at 0



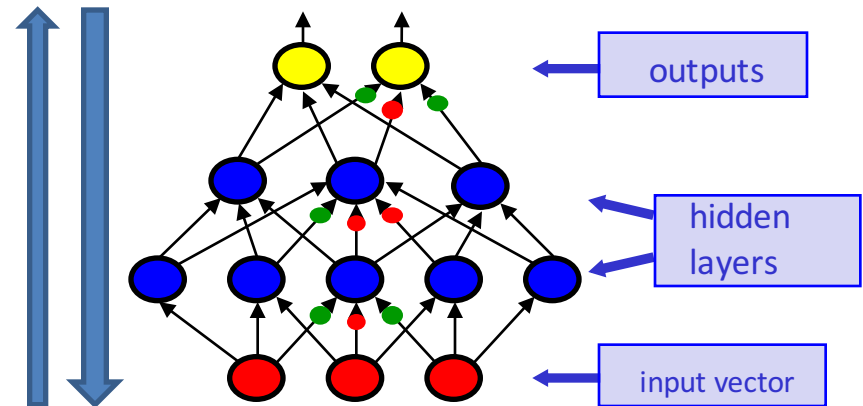
# Training the MLP: Backpropagation

## Testing for K-class classification problem

- For a given  $x$  with unknown class
- $x \in \text{class } k$ , if  $y_k = \max_i y_i$
- $y_i = v_i^T z = \sum_{h=1}^H v_{ih} z_h + v_{i0}$

That is

- A  $W$  represents a MLP
- Given a  $W$ , then we can classify a pattern  $x$



$$W = [w_1, \dots, w_K, v_1, \dots, v_H]$$

A **Machine Learning** problem:

how to obtain the  $W$  of a MLP

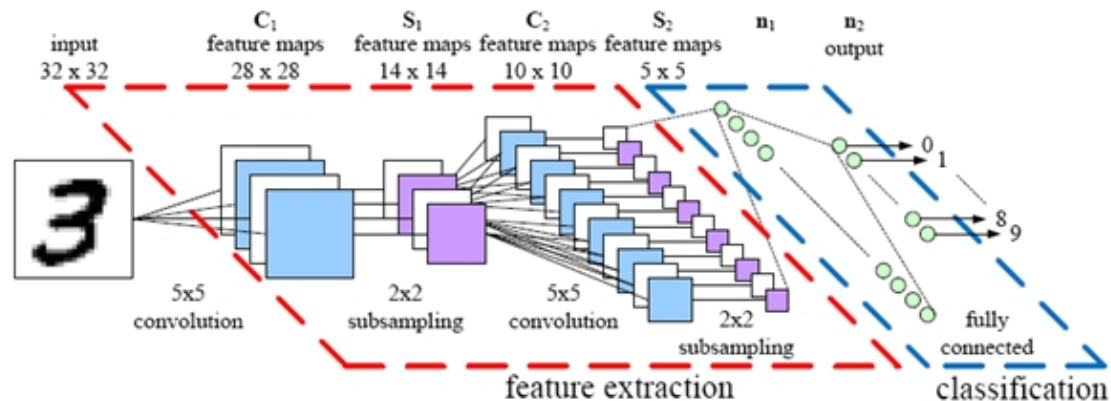
- We need a set of training patterns  $(X, Y)$
  - We need a learning algorithm to learn  $W$  by  $(X, Y)$
- => **Backpropagation learning algorithm  $B$ :  $W=B(X, Y)$**

# Neural Information Processing

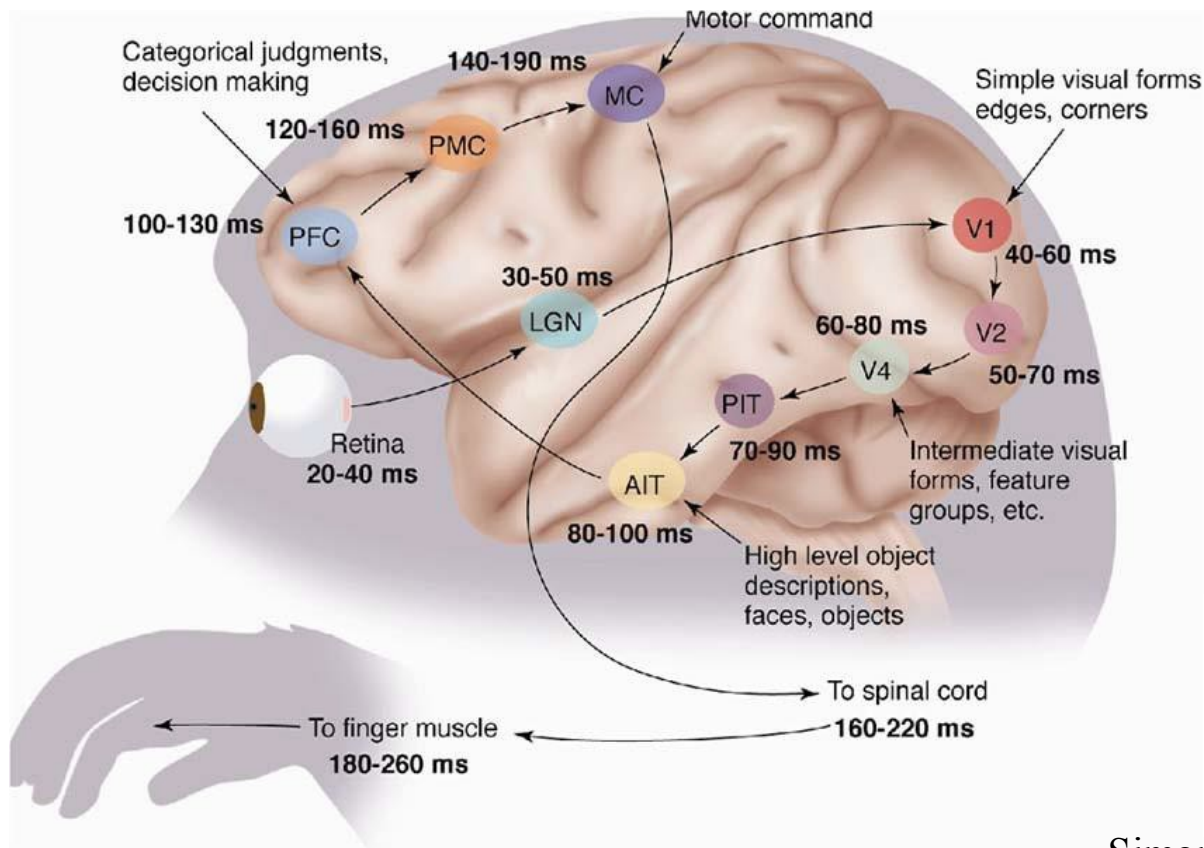
- Origins
- Perceptron
- Multilayer perceptron network
- **Convolutional neural networks (CNN)**

# What Is CNN

- An extension of MLP
  - *More hidden layers* to extract *features* of signal (image, speech)
  - Depth of network is "deep"
- CNN is one famous model of **deep neural networks (DNN)**



# Why CNN? Human vision is fast



Experiments show that humans can detect animals in a scene in less than 150 ms.

Simon Thorpe et al.,  
"Seeking Categories in the Brain,"  
Neuroscience 2001.

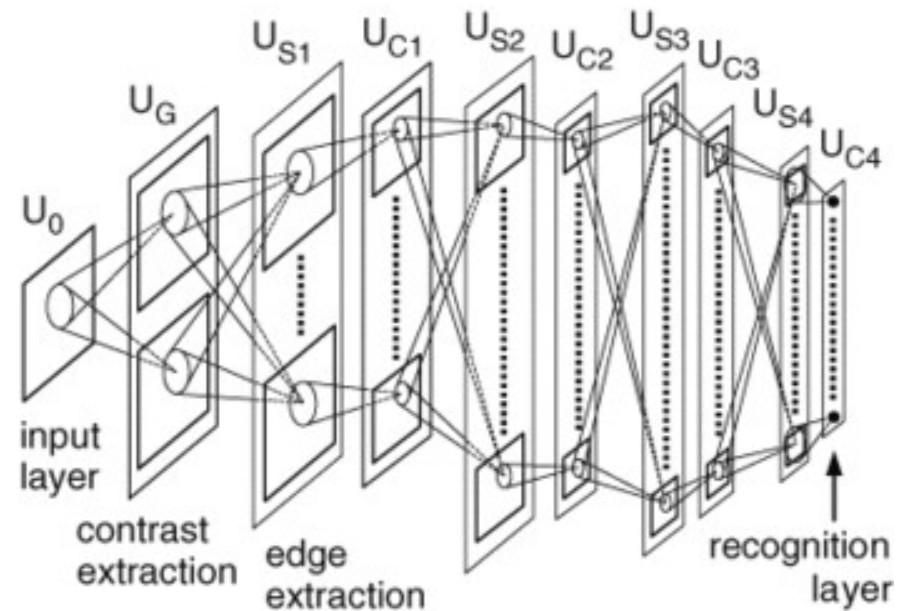
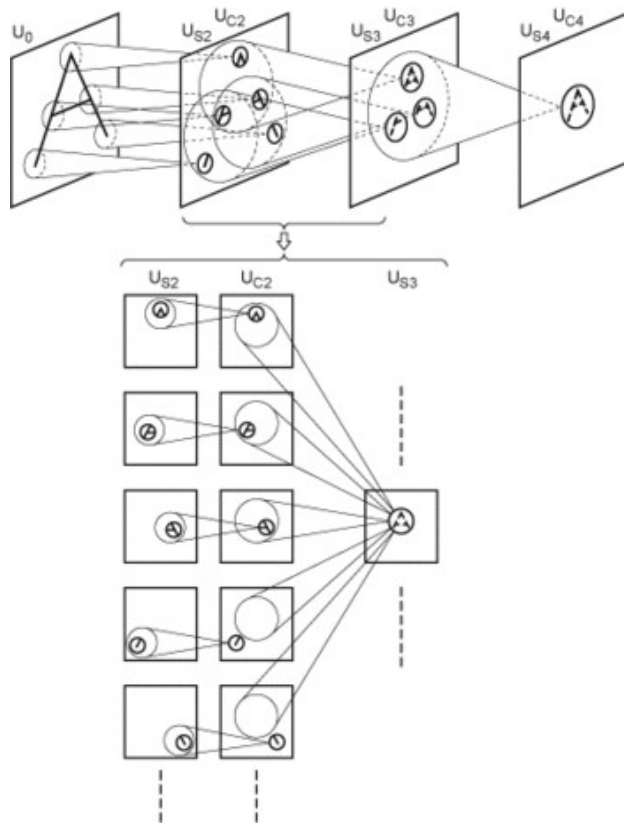
# Hubel & Wiesel (1962)

## Insights about early image processing in the brain

- Simple cells detect local features
- Complex cells pool local features in a retinotopic neighborhood



# The Neocognitron

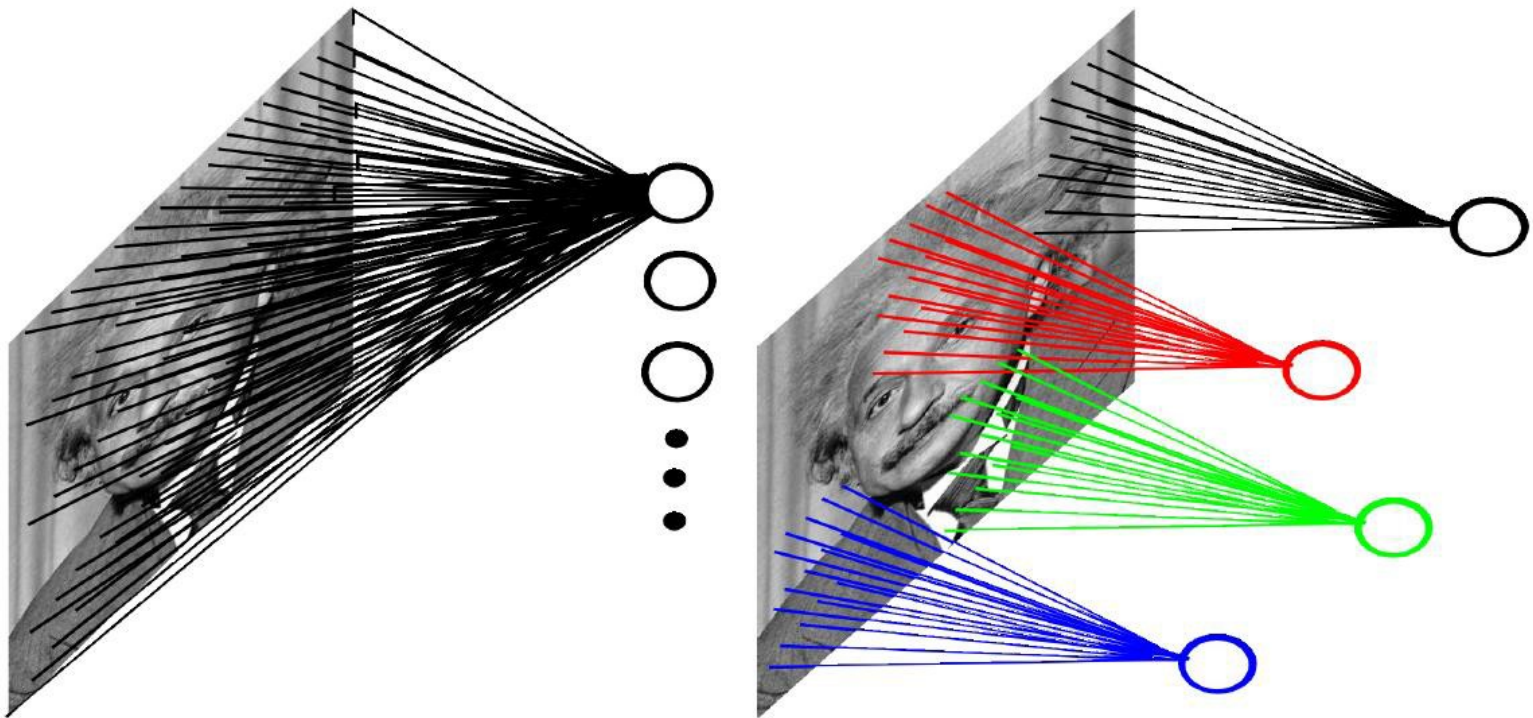


(Fukushima 1974-1982)



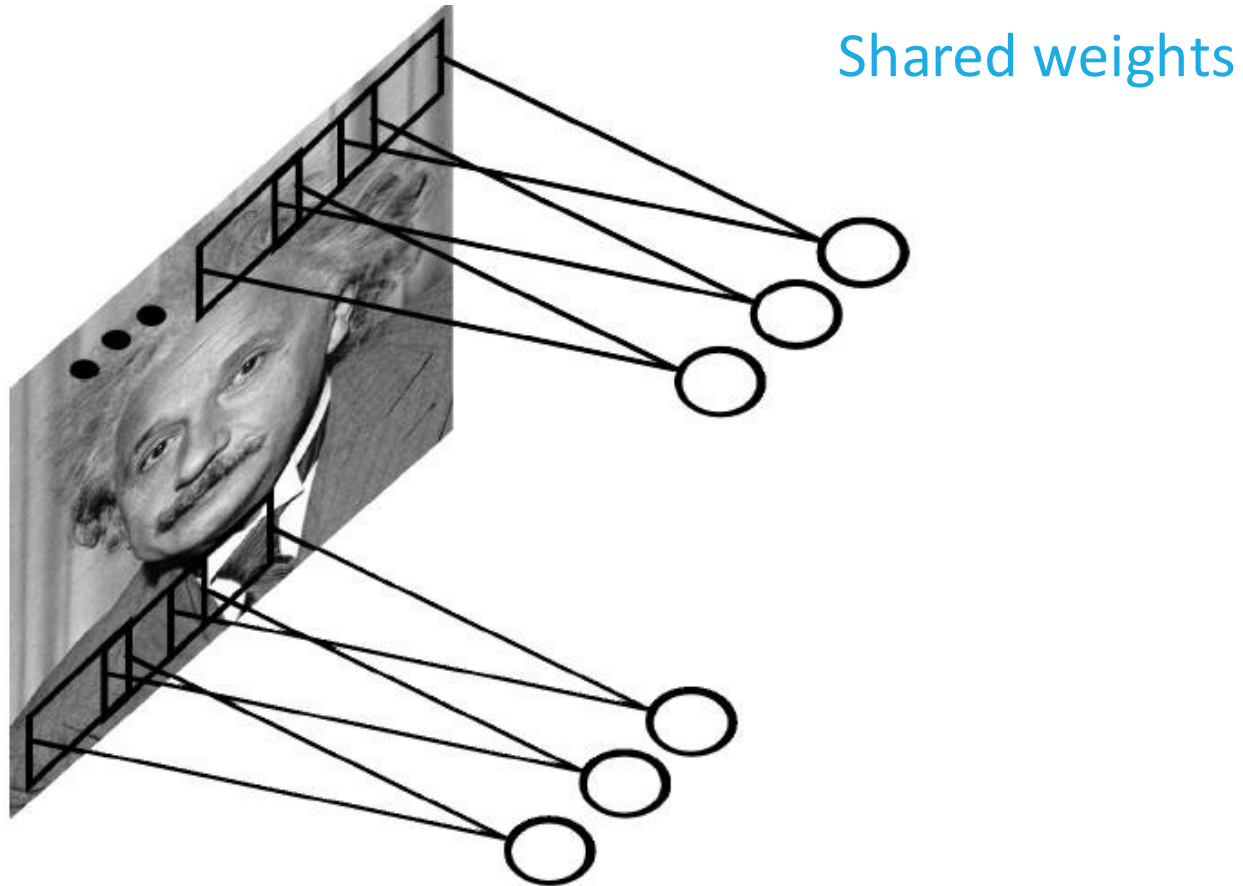
# Local connections

---



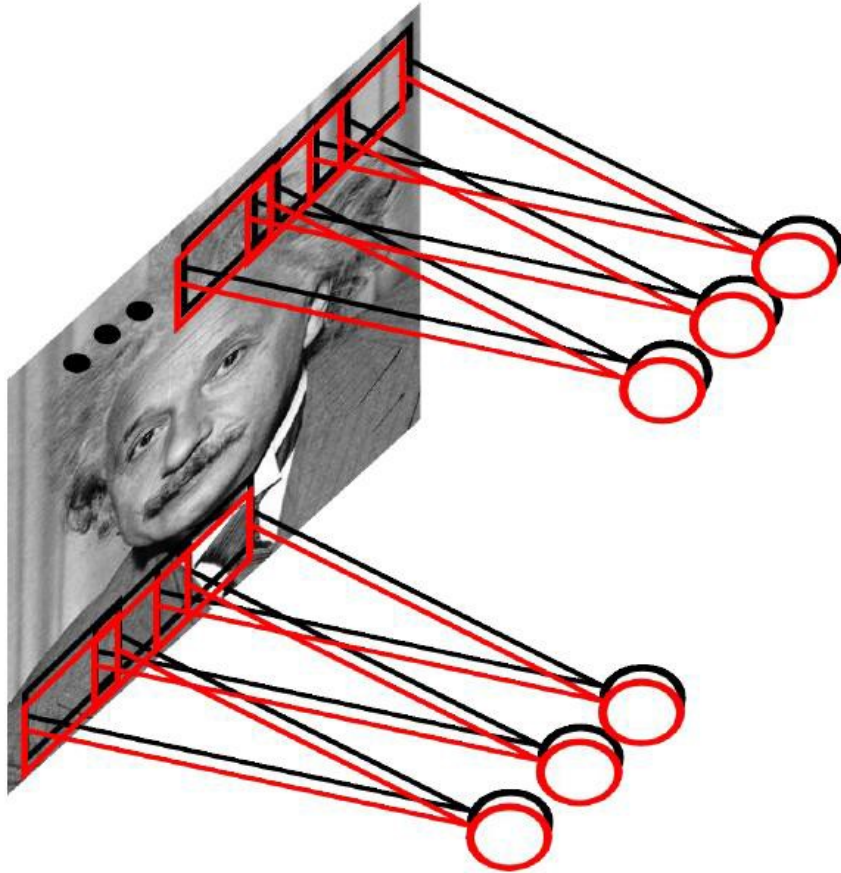
# Convolution

---



# Multiple convolutions

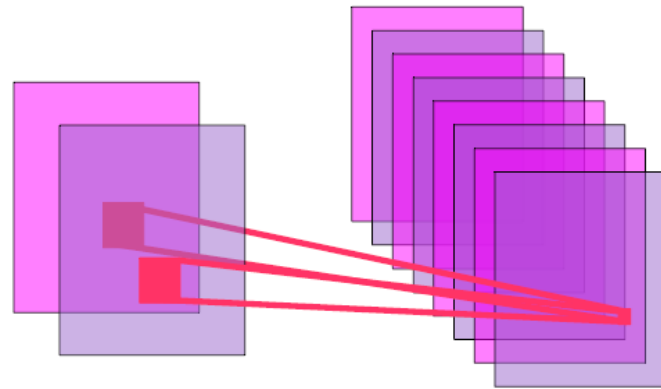
---



**CNN:**

Convolutional neural network

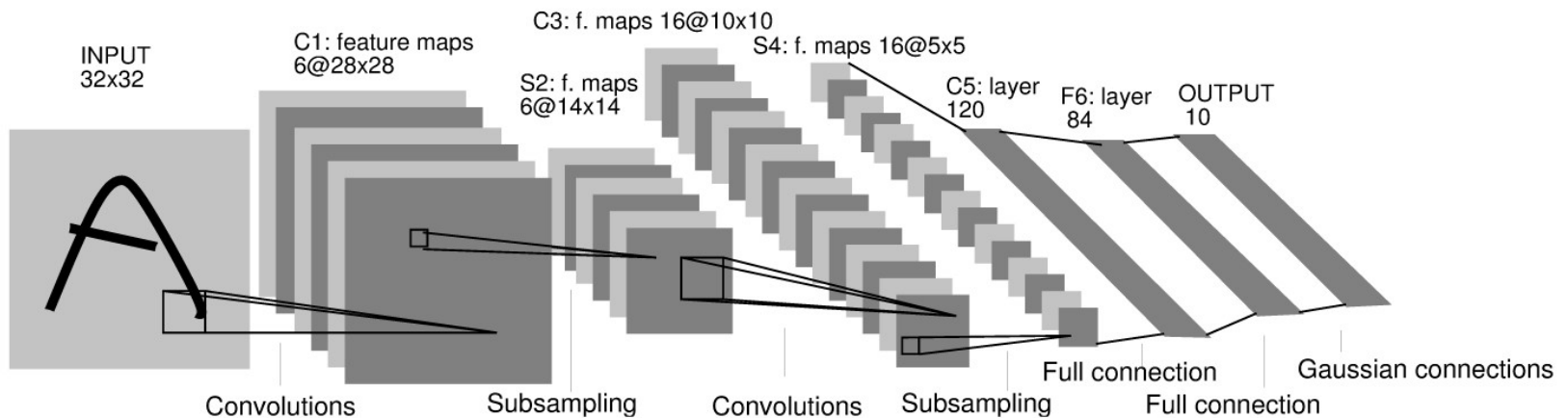
ConvNet



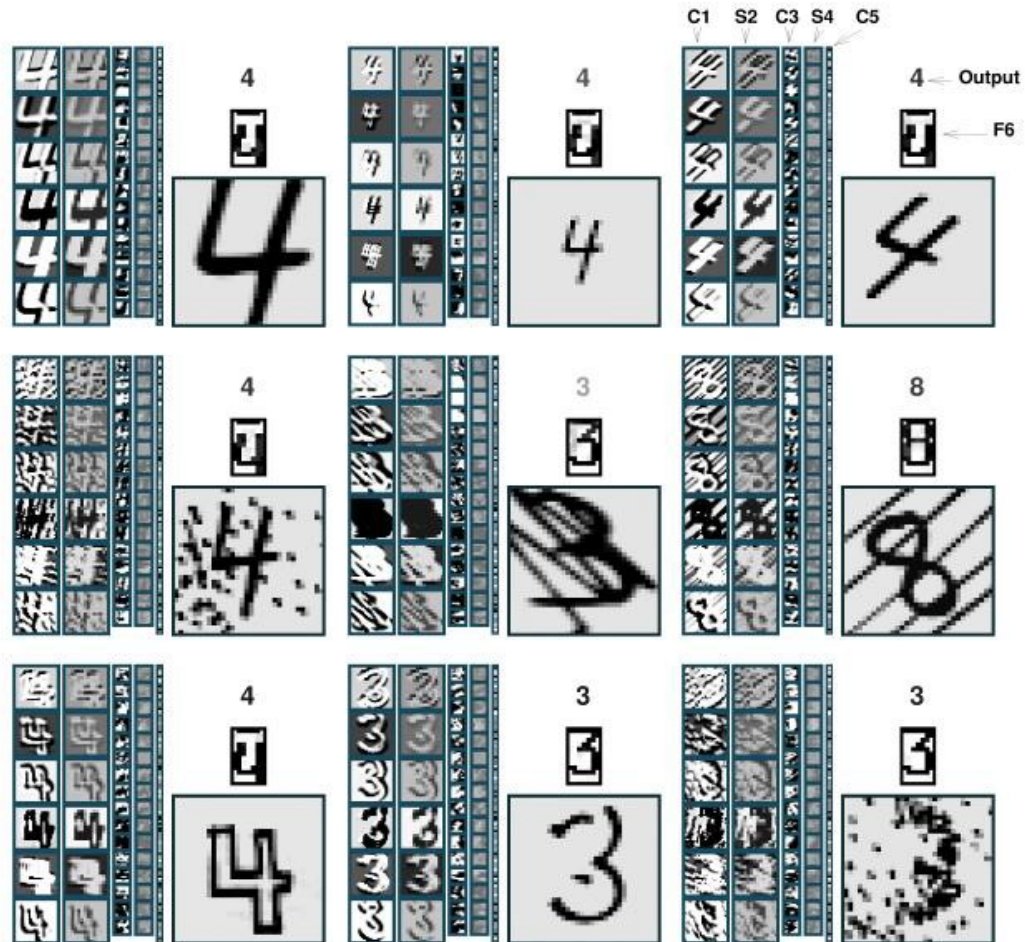
# CNNs in the 1990s

---

- 1989 Isolated handwritten character recognition (AT&T Bell Labs)
- 1991 Face recognition. Sonar image analysis. (Neuristique)
- 1993 Vehicle recognition. (Onera)
- 1994 Zip code recognition (AT&T Bell Labs)
- 1996 Check reading (AT&T Bell Labs)

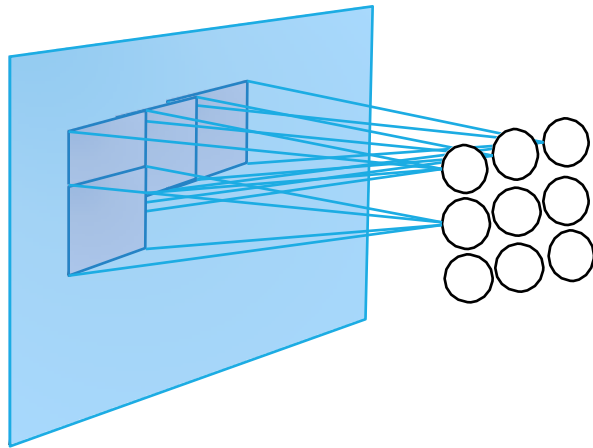


# Convnets in the 1990s



# Pooling

---



Name	Pooling formula
Average pool	$\frac{1}{s^2} \sum x_i$
Max pool	$\max\{x_i\}$
L2 pool	$\sqrt{\frac{1}{s^2} \sum x_i^2}$
L <sub>p</sub> pool	$\left(\frac{1}{s^2} \sum  x_i ^p\right)^{\frac{1}{p}}$

# Contrast Normalization

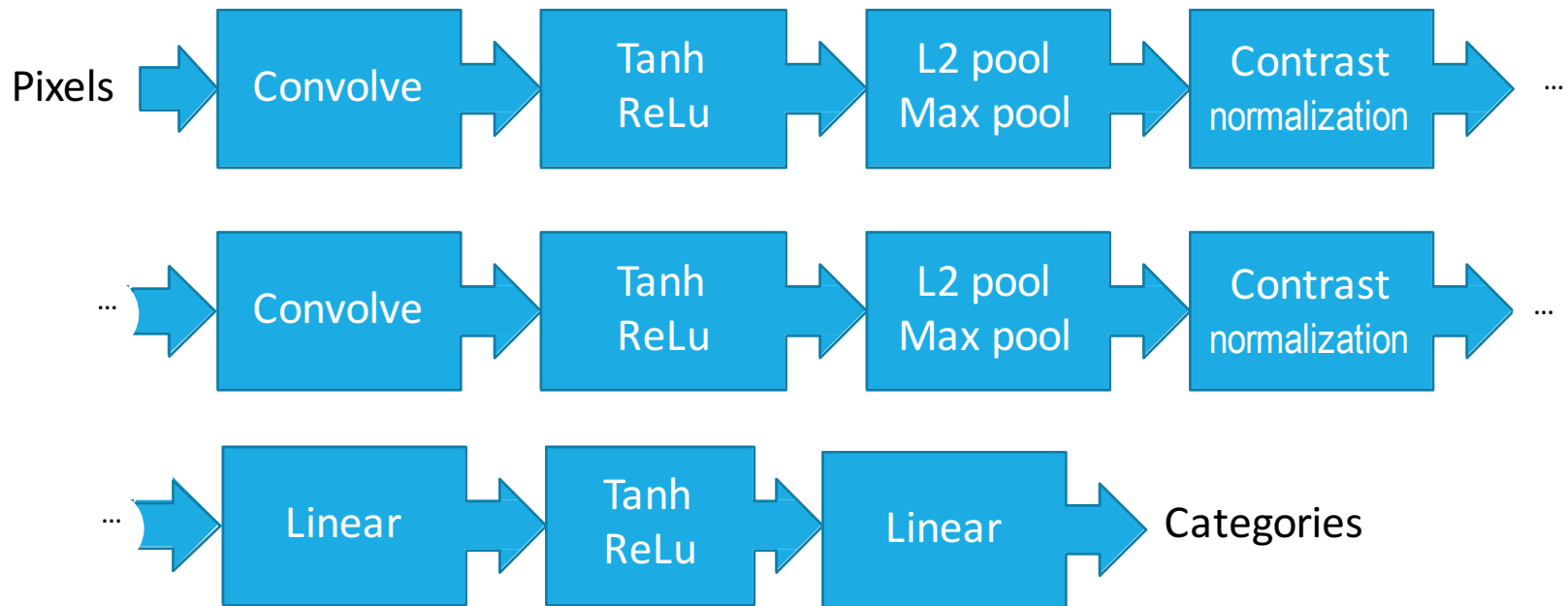
---

## Contrast normalization

- Subtracting a low-pass smoothed version of the layer
- Just another convolution in fact (with fixed coefficients)
- Lots of variants (per feature map, across feature maps, ...)
- Divisive normalization

# CNNs in the 2010s

---





# Convnets in the 2000s

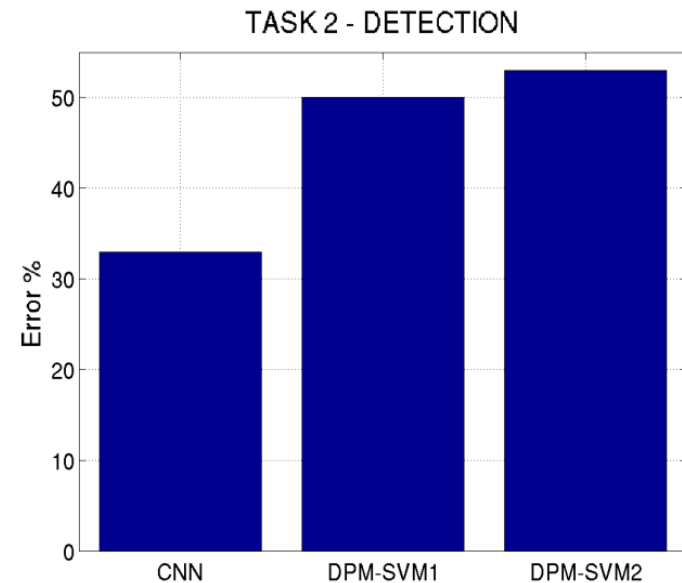
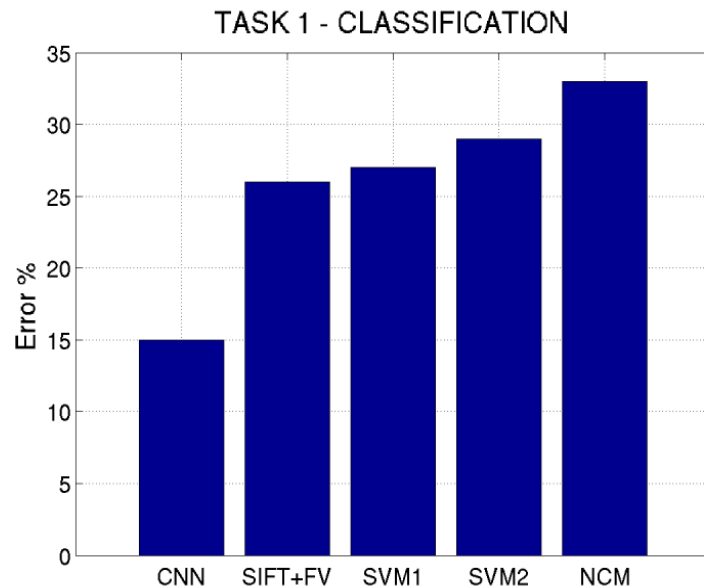
---

- OCR in natural images [2011]. Streetview house numbers (NYU)
- Traffic sign recognition [2011]. GTRSB competition (IDSIA, NYU)
- Pedestrian detection [2013]. INRIA datasets (NYU)
- Volumetric brain segmentation [2009]. Connectomics (MIT)
- Human action recognition [2002,2011]. Smartcatch (NEC), Hollywood II (SF)
- Object recognition [2004,2012]. Norb (NEC), ImageNet (UofT)
- Scene parsing [2010-2012]. Stanford bldg, Barcelona (NEC, NYU)
- Medical image analysis [2008]. Cancer detection (NEC)

# ImageNet 2012 competition

---

- Object recognition. 1000 categories. 1.2M examples





# ImageNet CNN



**mite**

**container ship**

**motor scooter**

**leopard**

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



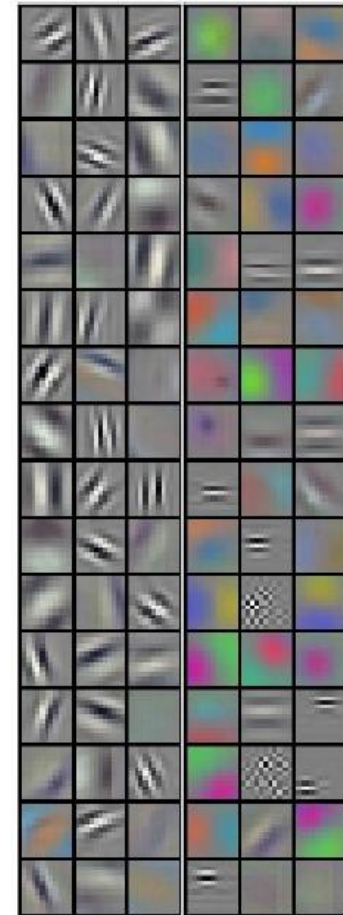
**grille**

**mushroom**

**cherry**

**Madagascar cat**

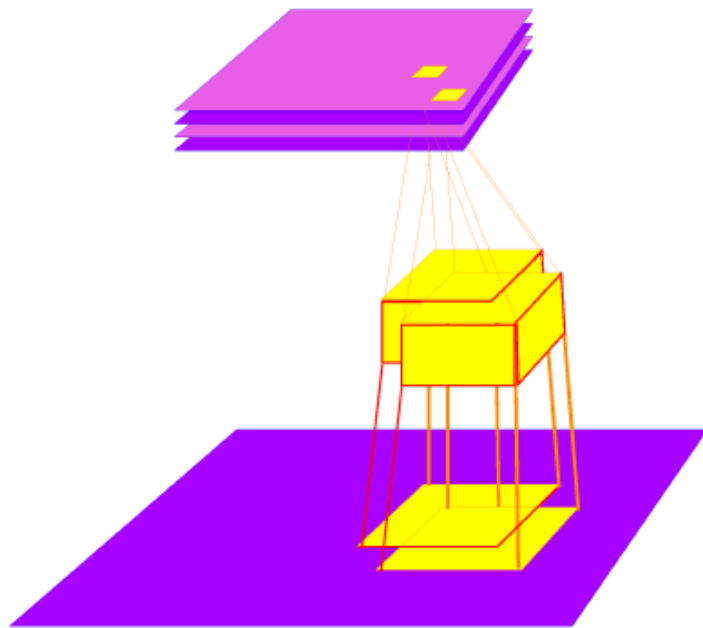
<p>convertible</p> <p>grille</p> <p>pickup</p> <p>beach wagon</p> <p>fire engine</p>	<p>agaric</p> <p>mushroom</p> <p>jelly fungus</p> <p>gill fungus</p> <p>dead-man's-fingers</p>	<p>dalmatian</p> <p>grape</p> <p>elderberry</p> <p>ffordshire bullterrier</p> <p>currant</p>	<p>squirrel monkey</p> <p>spider monkey</p> <p>titi</p> <p>indri</p> <p>howler monkey</p>
--	--	--	---



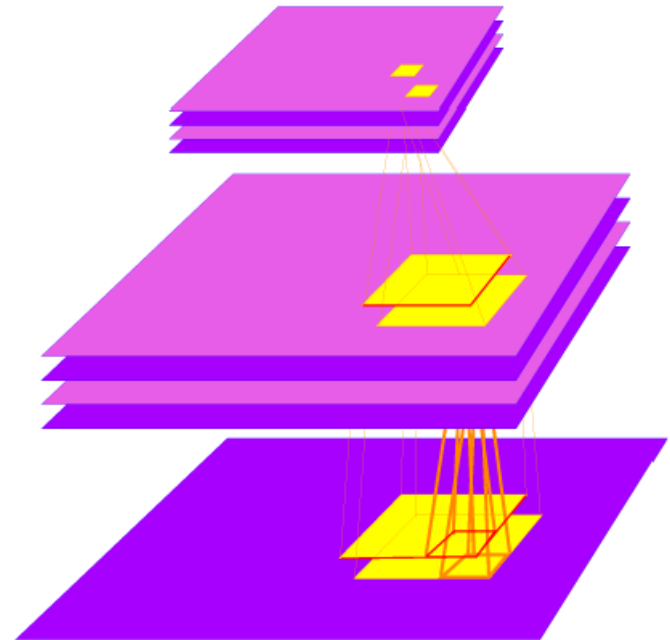
# Replicated CNNs

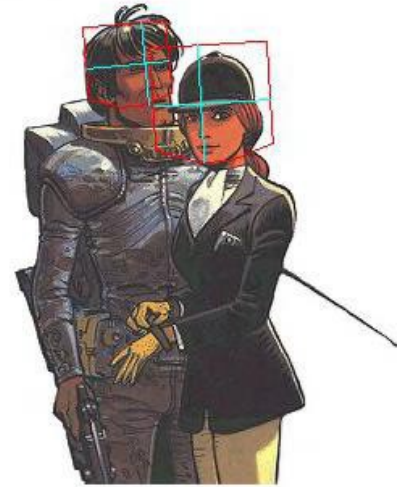
---

Wrong way



Right way





Today :

# Neural Information Processing

- Origins
- Perceptron
- Multilayer perceptron (MLP)
- Convolutional neural network (CNN)

Next :

# Training multilayer networks

- Perceptron learning
- Optimization basics
- Stochastic gradient descent
- Backpropagation learning algorithms