# 人工智慧視覺運算方法

謝東佑

可測及可靠系統實驗室
(Testable And Reliable Systems Lab., TARS)
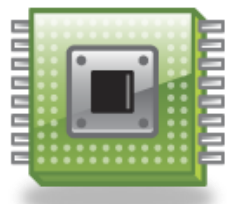
國立中山大學電機系

Office: 工EC-7038
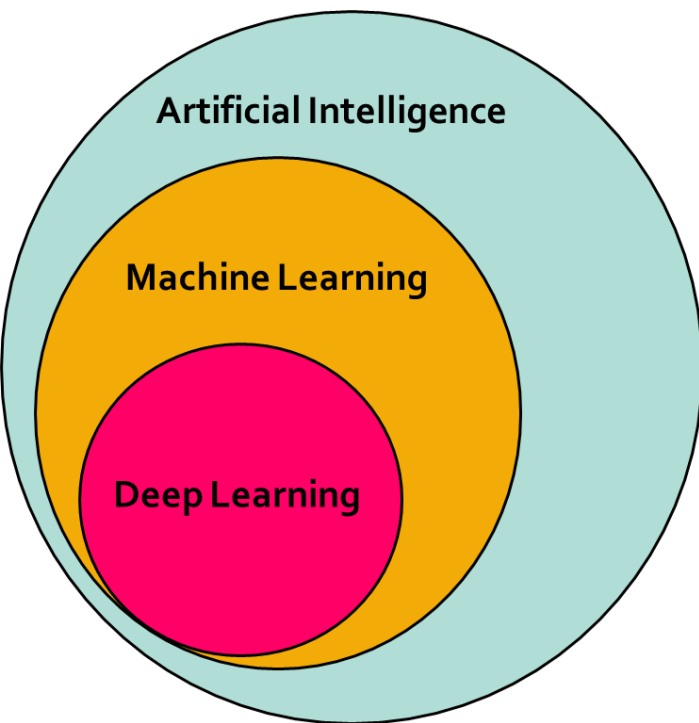
07-5252000 Ext. 4114

tyhsieh@mail.ee.nsysu.edu.tw
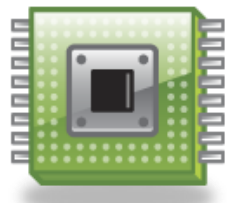
Keep feet on the ground

# AI vs Machine Learning vs Deep Learning

- AI: 模擬人類智慧
  - 結果有智慧就算
  - 一個擁有非常詳盡的 rule-based 系統也可以是 AI
- Machine learning是達成 AI 的一種方法
  - 從資料當中學習出 rules
  - 找到一個夠好的 function 能解決特定的問題
- Deep learning 是machine learning的一種
  - 從feature engineering 走向architecture engineering
  - 不再人工萃取特徵
  - 深層網路萃取更抽象特徵

# Deep Learning v.s. Feature Engineering



Raw data: pixel grid

Better features: clock hands' coordinates
{x1: 0.7, y1: 0.7}
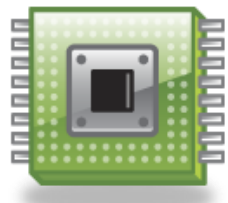{x2: 0.5, y2: 0.0}
{x1: 0.0, y2: 1.0}
{x2: -0.38, 2: 0.32}

Even better features: angles of clock hands
theta1: 45
theta2: 0
theta1: 90
theta2: 140

- **讓機器看時鐘報時**
- **直接看圖**
  - 要用CNN才行
  - 需要大量資料
- **放點工人智慧**
  - 用指針座標
  - 簡單的ML就可以
  - 少量資料就可以
- **更多工人智慧**
  - 用指針角度(像人看時鐘一樣)
  - 連ML都不用,查表就可以
  - 資料最少

對DL來講,好的特徵可以幫助你用較少資源與資料,
反過來,若你的資料資源很少,你會需要比較好的特徵(aka.更多工人智慧)

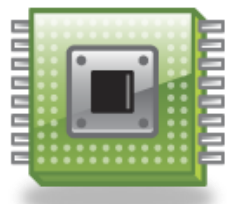# 影像處理 (Image Processing)

- **改變影像內容/本質，以方便**
  - 人眼辨識
  - 機器辨識

加強影線的邊緣線條，呈現更銳利的影像。見圖1.1



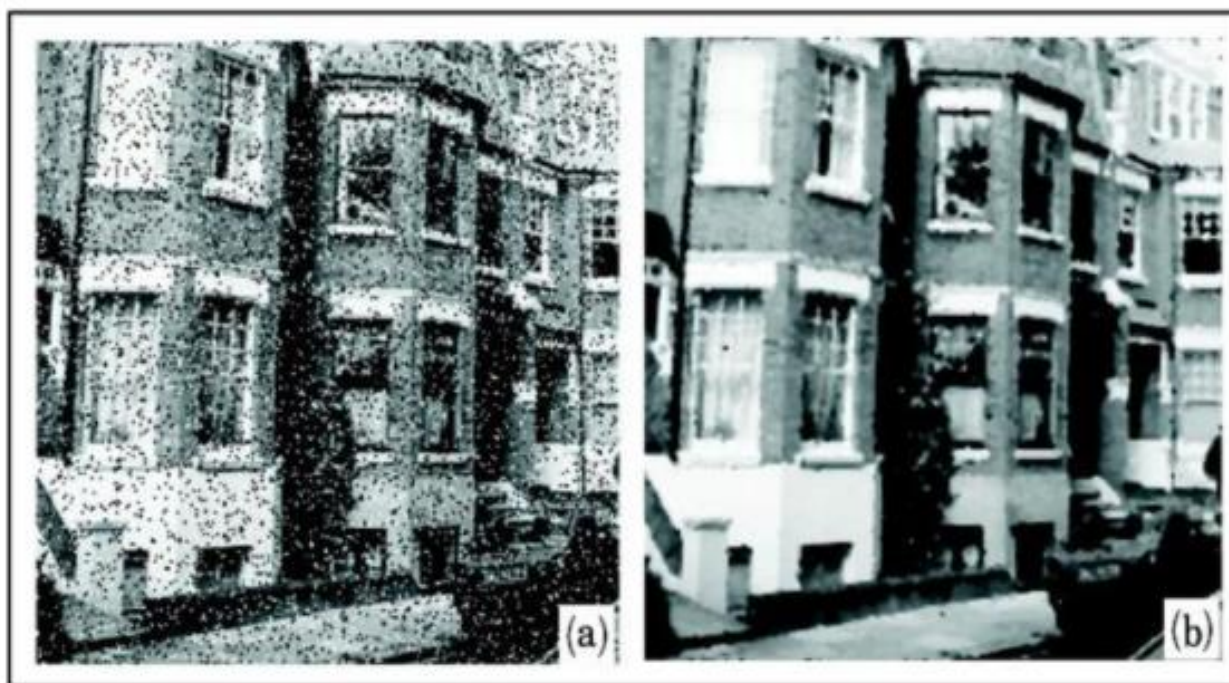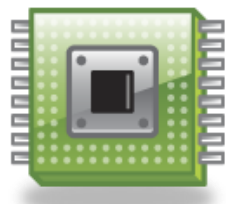圖1.1 影像銳利化 (a) 原始影像 (b) 銳利化結果

**NSYSUEE-TYHSIEH**

# 讓人看得更清晰

去除影像的雜訊。見圖1.2



圖1.2 去除影像雜訊 (a) 原始影像 (b) 去除雜訊結果
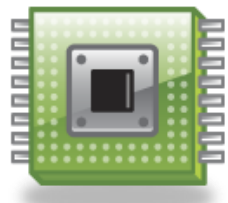
# 讓人看得更清晰

去除影像的動態模糊現象。見圖1.3



圖1.3 去除影像模糊現象 (a) 原始影像 (b) 去除模糊現象結果

取得影線邊緣線條，這個動作是為了測量影像中的物體。見圖1.4 (a與b)



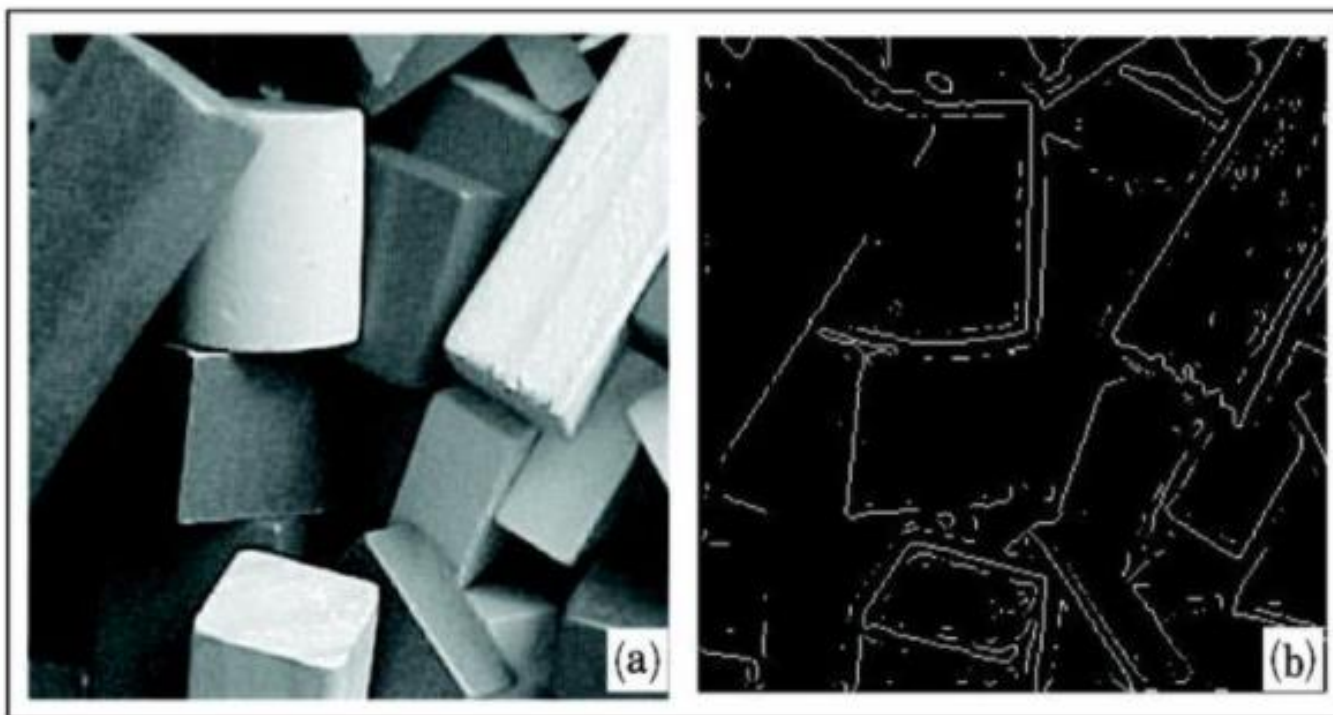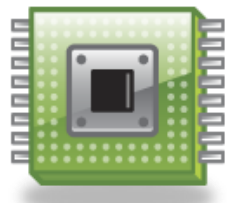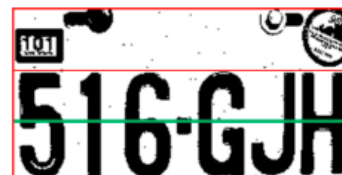圖1.4 取得影像邊緣線條 (a) 原始影像 (b) 物體邊緣線條

# 車牌辨識

相機或手機的影像 → 將影像縮圖至 320x240 → 將影像轉為灰階影像 → 邊緣偵測(Sobel) → 灰階、Sobel 做二值化 → Sobel 影像做中值濾波 → 車牌定位 → 字元切割 → 字元辨識 → 顯示車牌號碼

Gartner Hype Cycle for Artificial Intelligence, 2019

# AI VISUAL ALGORITHMS

# How Does A Computer Classify Pictures?

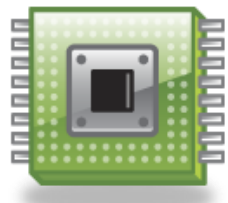- **A picture is only a group of pixels for a computer.**

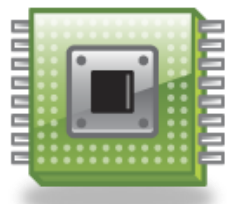- **Modern AI nets learn features of objects.**
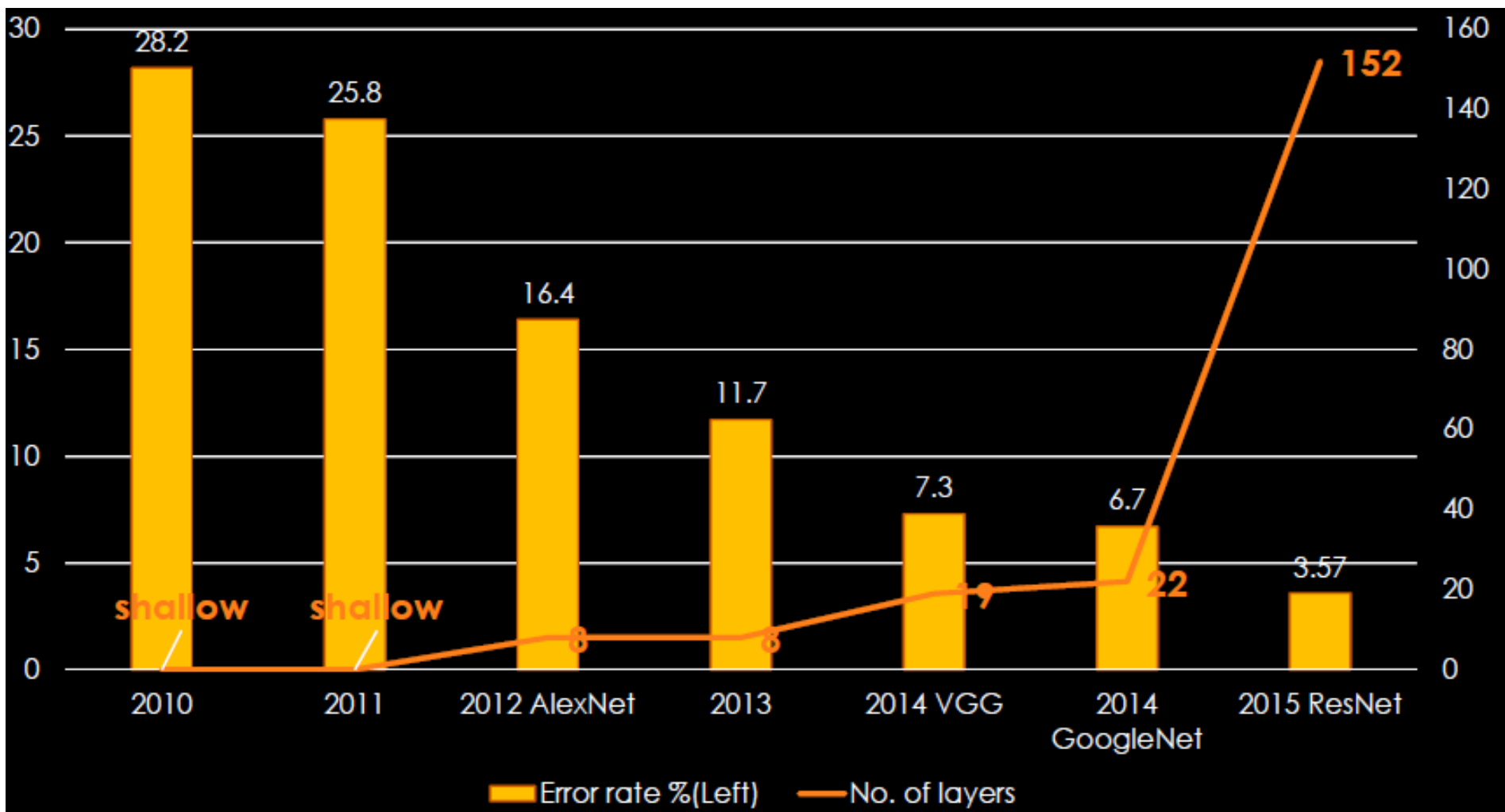


Images source: CC dataset
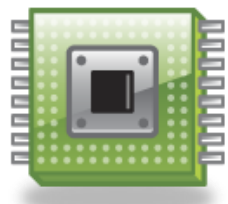
# Object Classification

- **Modern AI algorithms for object classification**
  - **AlexNet, 5 CNN layers and 3 FC layers, 2012**
  - **VGG, 16 CNN layers and 3 FC layers, 2014**
  - **GoogLenet, 21 CNN layers and 1 FC layer, 2014**
  - **ResNet, 151 CNN layers and 1 FC layer, 2015**
- **Foundation of object detection**
- **Limitation**
  - **One object in one picture, no localization**

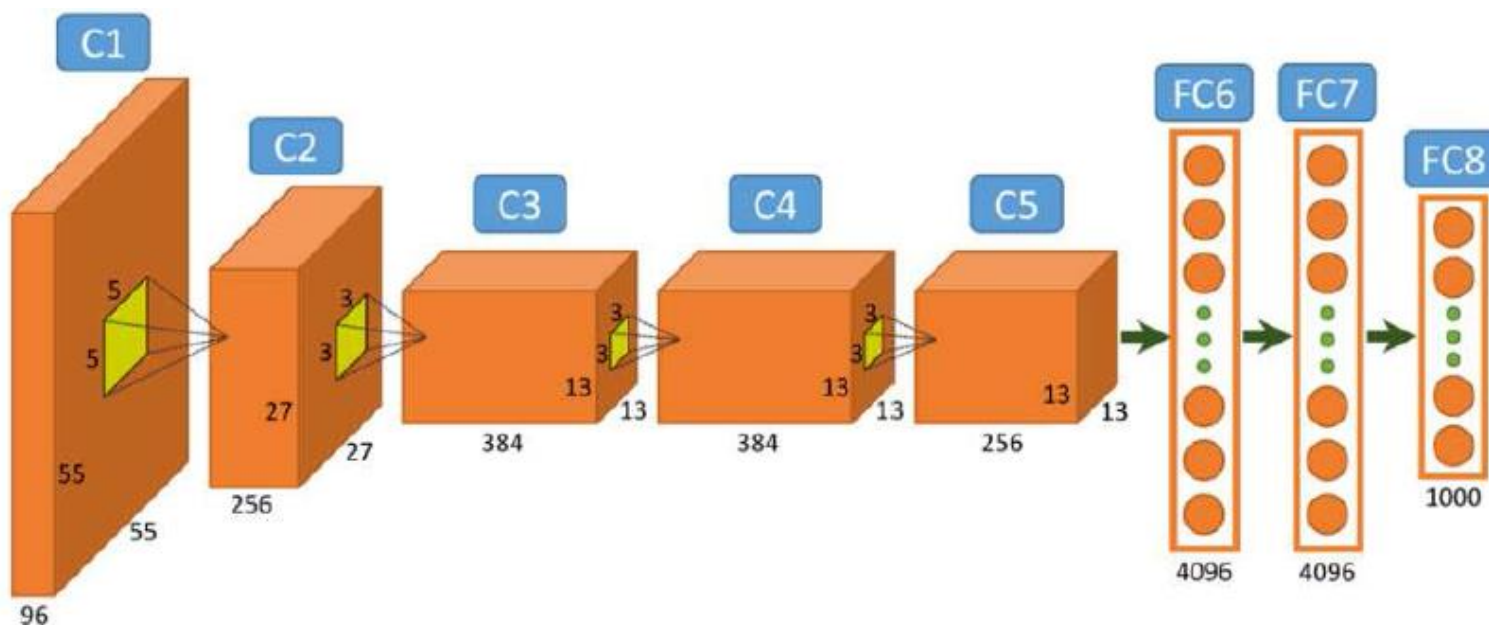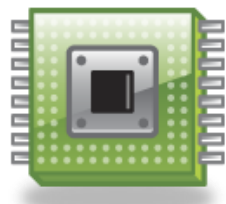# ILSVRC(IMAGENET Large Scale Visual Recognition Competition)
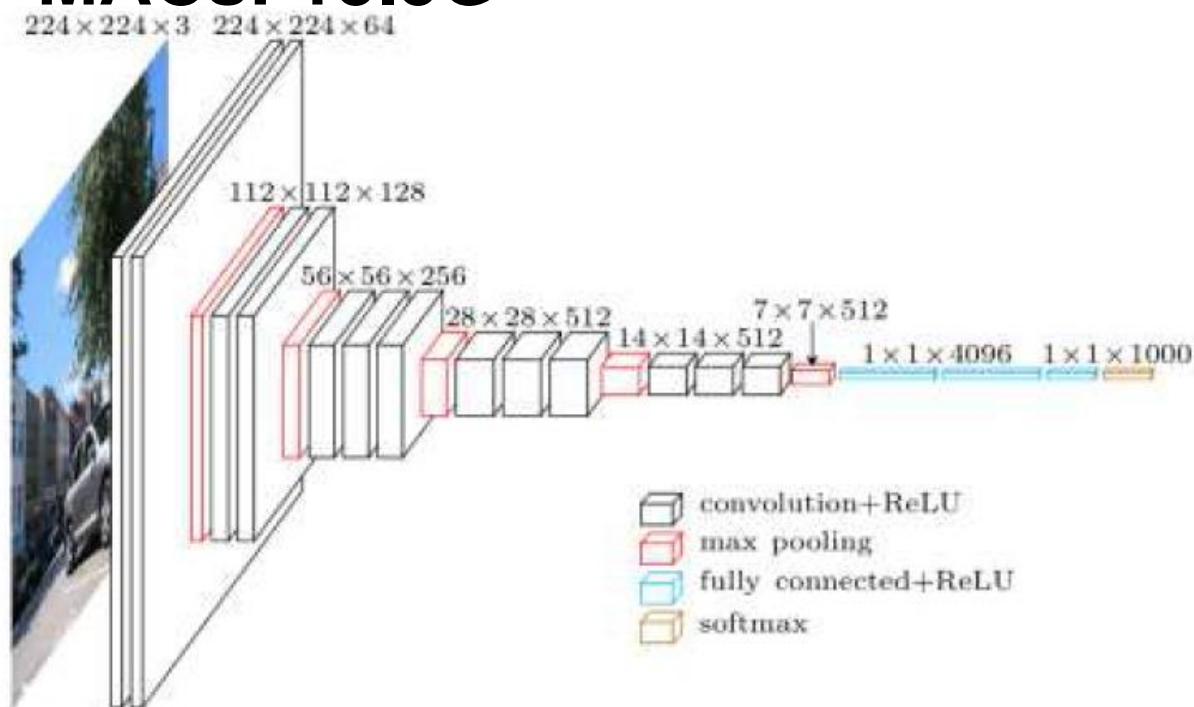
# AlexNet

- **CONV Layers: 5**

- **Fully Connected Layers: 3**

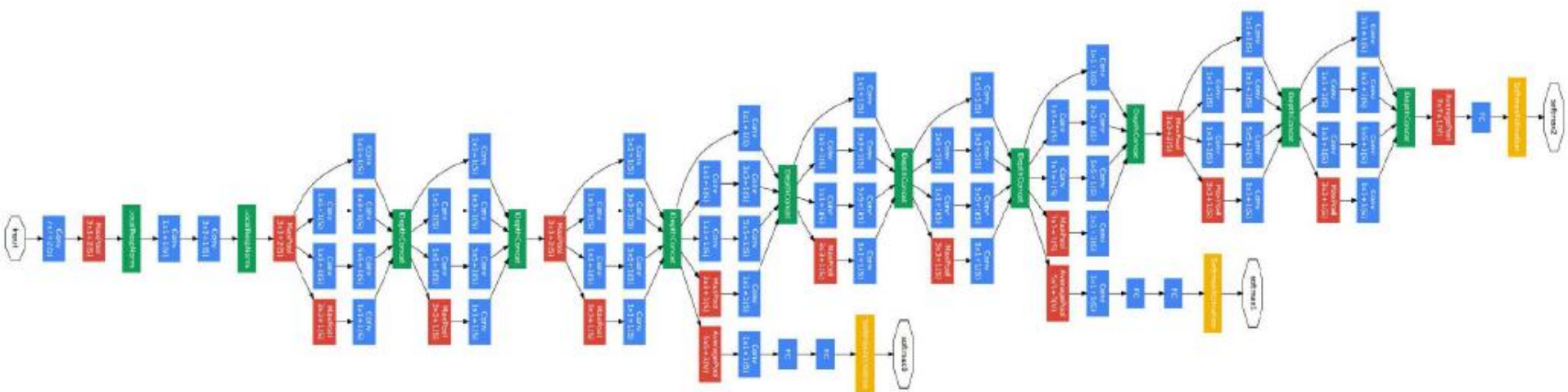- **Weights: 61M**

- **MACs: 724M**

# VGG

- **CONV Layers: 16**
- **Fully Connected Layers: 3**
- **Weights: 138M**
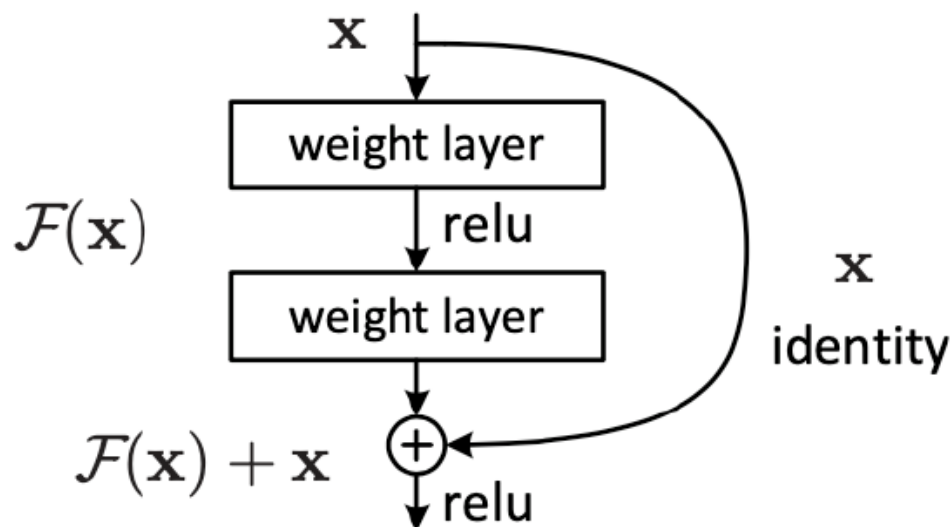- **MACs: 15.5G**

# GoogLenet

- **CONV Layers: 21**

- **Fully Connected Layers: 1**

- **Weights: 7.0M**

- **MACs: 1.43G**

# ResNet

- **Main idea**
  - **Residual layer**
- **CONV Layers: 151**
- **Fully Connected Layers: 1**
- **Weights: 25.5M**
- **MACs: 3.9G**



$\mathbf{x}$

weight layer

$\mathcal{F}(\mathbf{x})$ → relu

weight layer

$\mathbf{x}$ identity

$\mathcal{F}(\mathbf{x}) + \mathbf{x}$ → ⊕ → relu

# Idea of CNN (Convolutional Neural Network)

# **Convolution Layer**



Input Image          Feature
                     Detector

卷積運算          Feature Map

# **Convolution Layer**



Input Image ⊗ Feature Detector = Feature Map

卷積運算

# **Convolution Layer**



We create many feature maps to obtain our first convolution layer

Feature Maps

Input Image

Convolutional Layer

16種不同的Feature Detector

# **Convolution Layer**



利用Feature Detector萃取出物體的邊界

# 使用Relu函數去掉負值，更能淬煉出物體的形狀

# 使用Relu函數去掉負值，更能淬煉出物體的形狀



Black = negative; white = positive values

Only non-negative values

# 使用**Relu**函數去掉負值，更能淬煉出物體的形狀



Black = negative; white = positive values

Only non-negative values

# 其他函數



Sigmoid          ReLU         Leaky ReLU

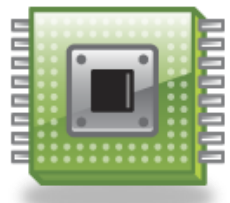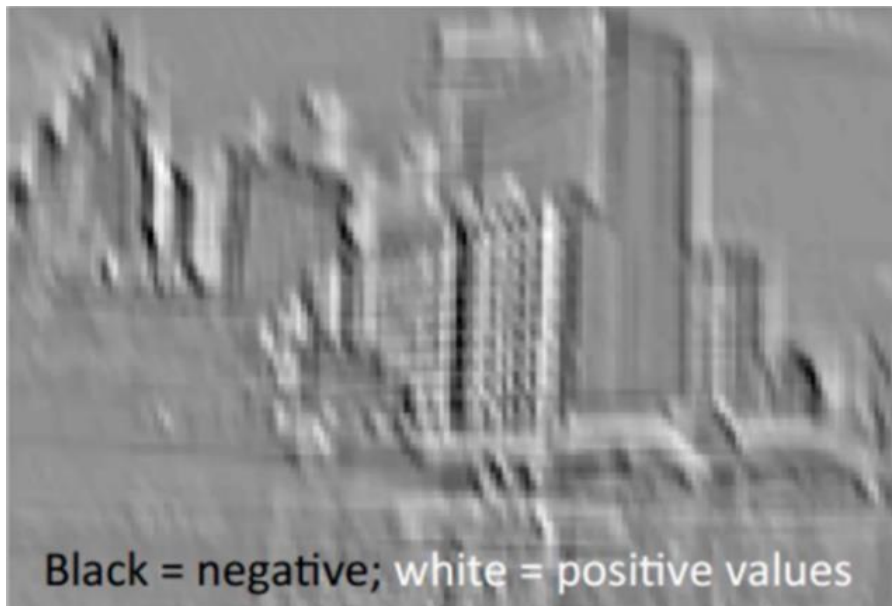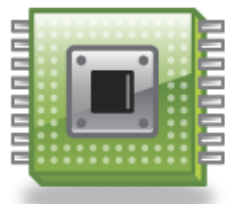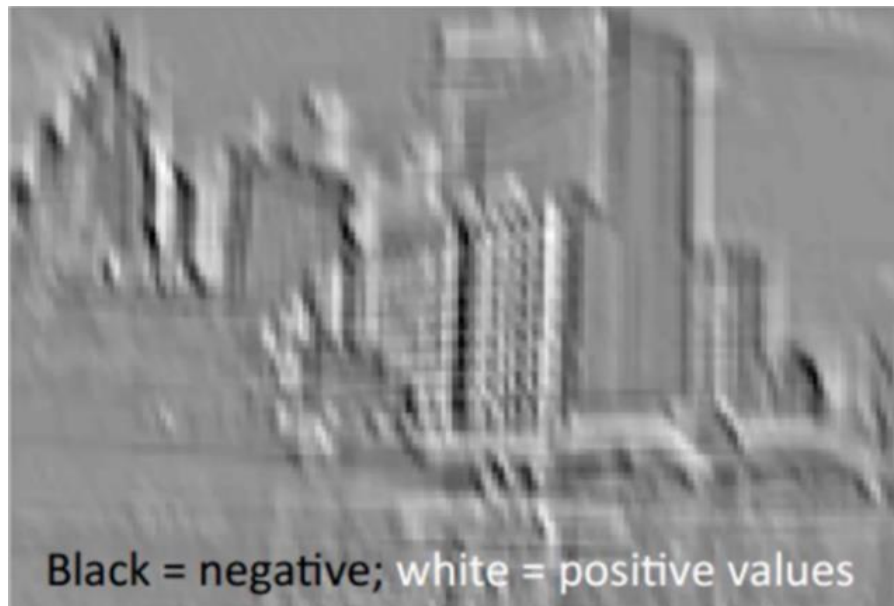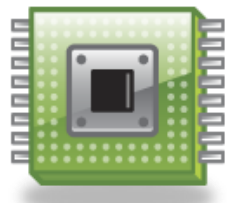# **Pooling Layer** 池化層

- **Max Pooling**

- 當圖片整個平移幾個**Pixel**的話對判斷上完全不會造成影響，以及有很好的抗雜訊功能

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

Max Pooling →

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

Feature Map

Pooled Feature Map

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

Max Pooling →

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

Feature Map

Pooled Feature Map

# **Fully Connected Layer** 全連接層

- 將之前的結果平坦化之後接到最基本的神經網絡



Pooled Feature Map

Flattening

# Fully Connected Layer 全連接層

■ 將之前的結果平坦化之後接到最基本的神經網絡



Input Image

Convolution → Convolutional Layer

Pooling → Pooling Layer

Flattening → Input layer of a future ANN
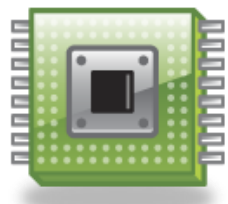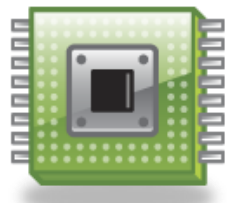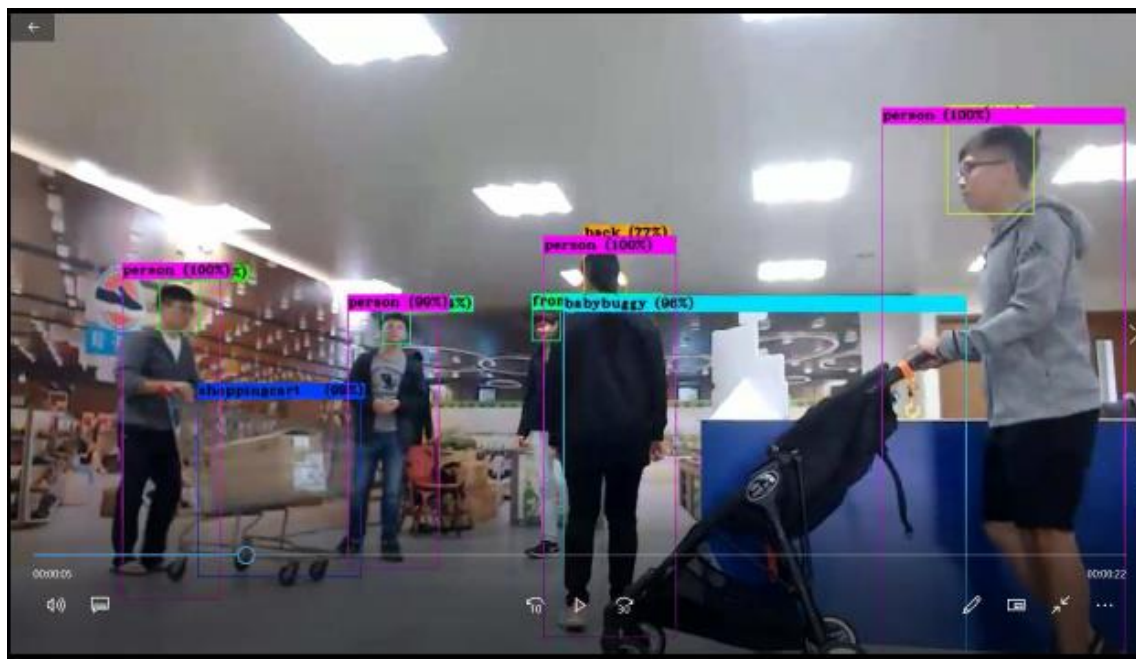
# Fully Connected Layer 全連接層
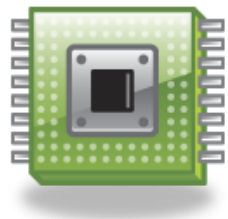
■ 將之前的結果平坦化之後接到最基本的神經網絡

# Objection Localization

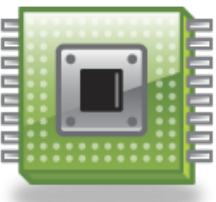- **Besides class, the computer needs to know the location of each object.**



Images source: FCU SoC Lab

# Modern AI Algorithms for Object Detection

- **RCNN (Region-based CNN), fast RCNN, faster RCNN**

- **YOLO (You Only Look Once)**

- **SSD (Single Shot Detection)**

# Object Detection

- **Two-stage object detection**
  - **Good detection accuracy but slow operation**
  - **Ex: Faster R-CNN**

| Input Image | → | Feature Extractor | → | Object localization |
|---|---|---|---|---|
| | | | | ↓ |
| | | | | Classification |

- **One-stage object detection**
  - **Fast operation and acceptable detection accuracy**
  - **Ex: SSD, YOLO**

| Input Image | → | Feature Extractor | → | Object localization and Classification |
|---|---|---|---|---|

# RCNN (Region-Based CNN), Fast RCNN, Faster RCNN

- **Two-stage ways**

| Region proposal (SS) | |
|---|---|
| Feature extraction (deep net) | |
| Classification (SVM) | (regression) |

RCNN
Slow in both training and testing

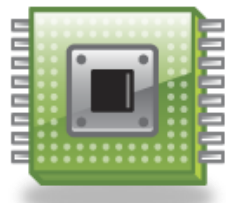| Region proposal (SS) |
|---|
| Feature extraction, Classification, Rect. refine   (deep net) |

Fast-RCNN
Few seconds per frame

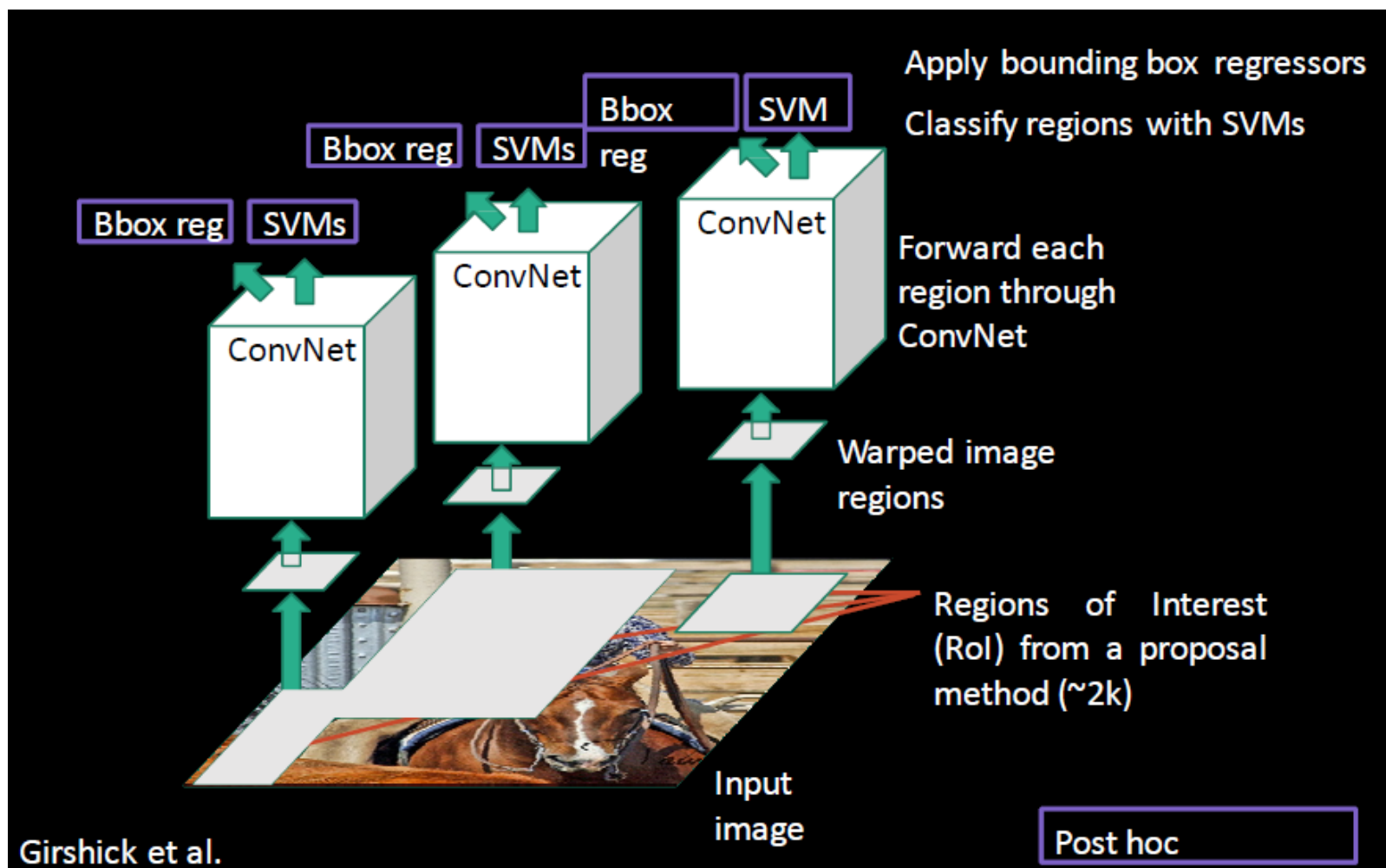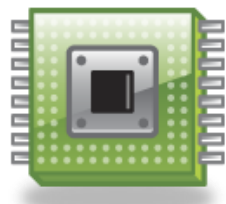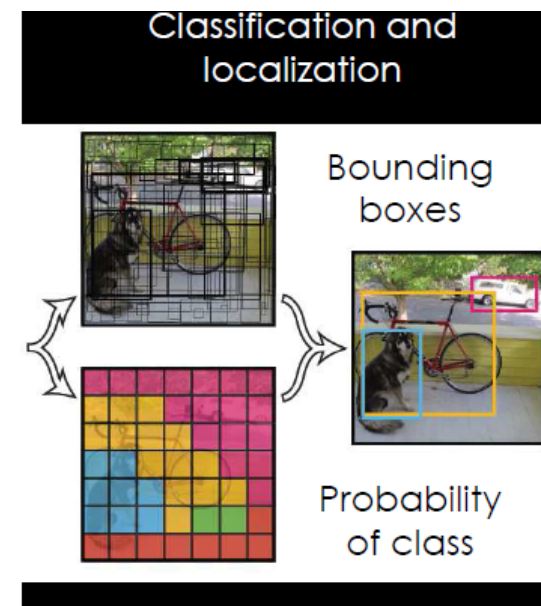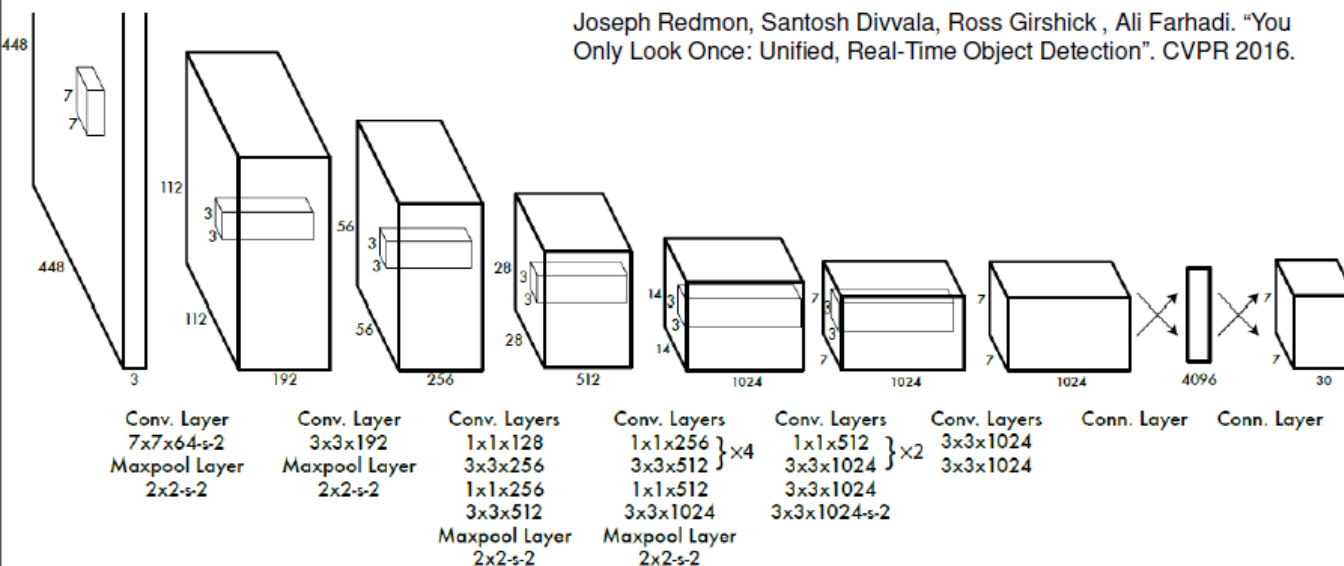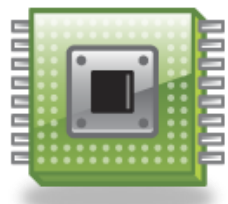| Region proposal, Feature extraction, Classification, Rect. refine   (deep net) |
|---|

Faster-RCNN
A dozen of fps on k40

# Two-Stage Ways

- ## One-stage way



Joseph Redmon, Santosh Divvala, Ross Girshick, Ali Farhadi. "You Only Look Once: Unified, Real-Time Object Detection". CVPR 2016.

# YOLO V3

YOLO v1 (2015)

BN,

YOLO v2 (2017)

ResNet,

Multiscale

YOLO v3 (2018)

https://pjreddie.com/darknet/yolo/

Fast and accurate under mAP-50



| Method | mAP-50 | time |
|---|---|---|
| [B] SSD321 | 45.4 | 61 |
| [C] DSSD321 | 46.1 | 85 |
| [D] R-FCN | 51.9 | 85 |
| [E] SSD513 | 50.4 | 125 |
| [F] DSSD513 | 53.3 | 156 |
| [G] FPN FRCN | **59.1** | 172 |
| RetinaNet-50-500 | 50.9 | 73 |
| RetinaNet-101-500 | 53.1 | 90 |
| RetinaNet-101-800 | 57.5 | 198 |
| **YOLOv3-320** | 51.5 | **22** |
| **YOLOv3-416** | 55.3 | 29 |
| **YOLOv3-608** | 57.9 | 51 |

# SSD (Single Shot Detection)

- **+Multi-scale feature maps**
- **- FC layers**

# Performance Evaluation Indexes

- **TP, FP, TN, FN**
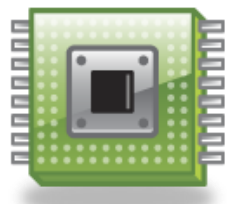- **Precision, Recall**
- **mAP (mean Average Precision)**

# TP, FP, TN, FN

- **TP: True Positive**
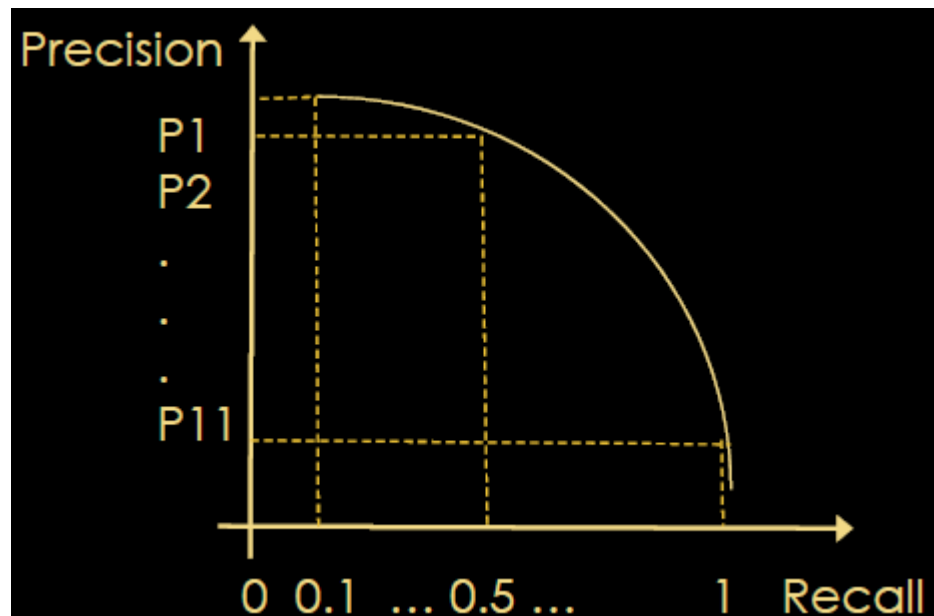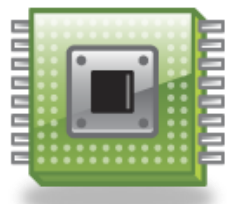- **FP: False Positive**
- **TN: True Negative**
- **FN: False Negative**

# Precision, Recall

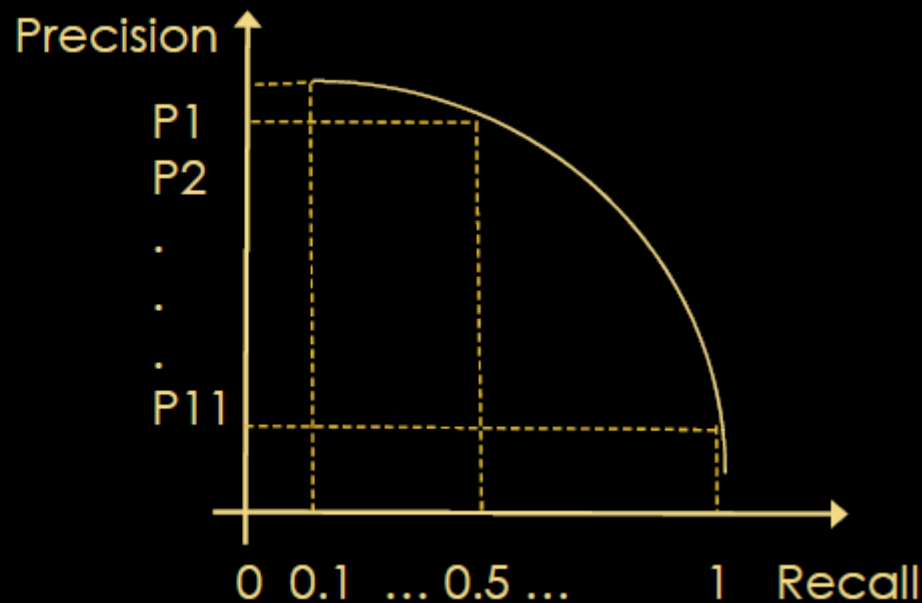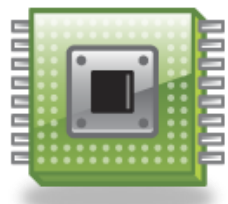$$Precision = \frac{TP}{TP+FP}$$

$$Recall = \frac{TP}{TP+FN}$$

# mAP (mean Average Precision)

- **AP: the average precision of precisions of different recalls**
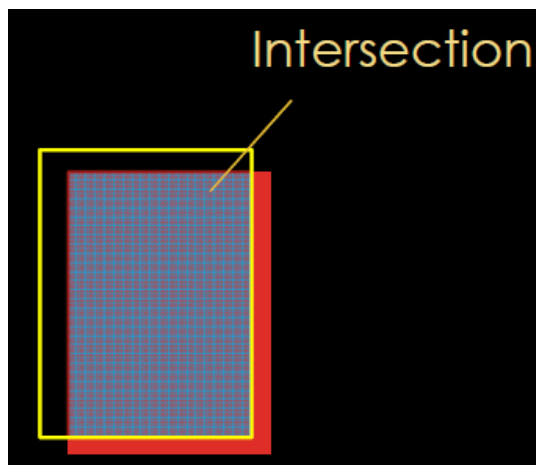
- **mAP: the mean of APs of different kinds of objects**
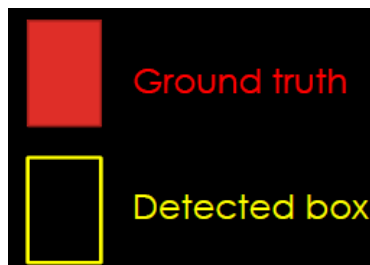
$$AP = \int_0^1 p(r)\, dr$$

$$mAP = \frac{\sum_{i=1}^n AP_i}{n}$$

# Important Parameters
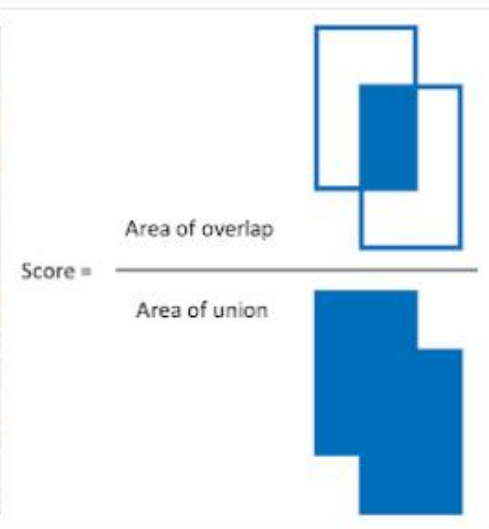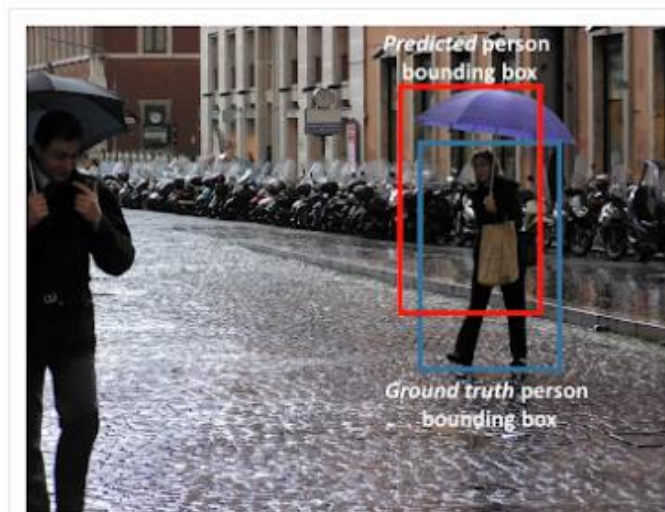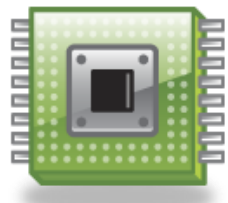
- **IoU (Intersection over Union)**
  - 一般**IoU>0.5**時為預測成功
- **Confidence threshold**

$$IoU(A, B) = \frac{A \cap B}{A \cup B}$$



Ground truth

Detected box



Intersection



Predicted person bounding box

Ground truth person bounding box

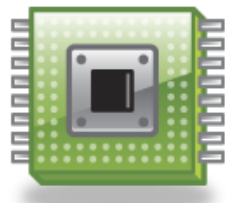$$Score = \frac{Area\ of\ overlap}{Area\ of\ union}$$

# Example 1

- Class: people
  - TP: 5
  - FP: 0
  - FN: 0
  - Precision: 5/5
  - Recall: 5/5

# Example 2



- Class: people
  - TP: 4
  - FP: 0
  - FN: 1
  - Precision: 4/4
  - Recall: 4/5

|  |  | Actual | |
|---|---|---|---|
|  |  | Positive | Negative |
| **Predictive** | Positive | TP | FP |
|  | Negative | FN | TN |

$$Precision = \frac{TP}{TP+FP}$$

$$Recall = \frac{TP}{TP+FN}$$

# Example 3



- Class: people
  - TP: 5
  - FP: 1
  - FN: 0
  - Precision: 5/6
  - Recall: 5/5

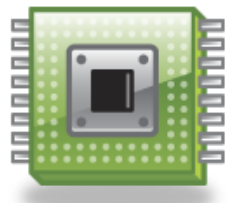| | | Actual | |
|---|---|---|---|
| | | Positive | Negative |
| **Predictive** | Positive | TP | FP |
| | Negative | FN | TN |

$$\text{Precision} = \frac{TP}{TP+FP}$$

$$\text{Recall} = \frac{TP}{TP+FN}$$

# Example 4

- **Assume**
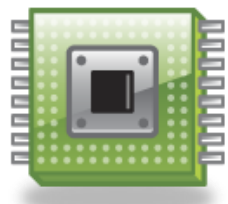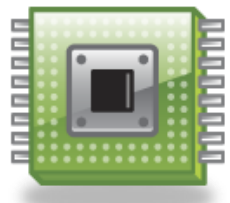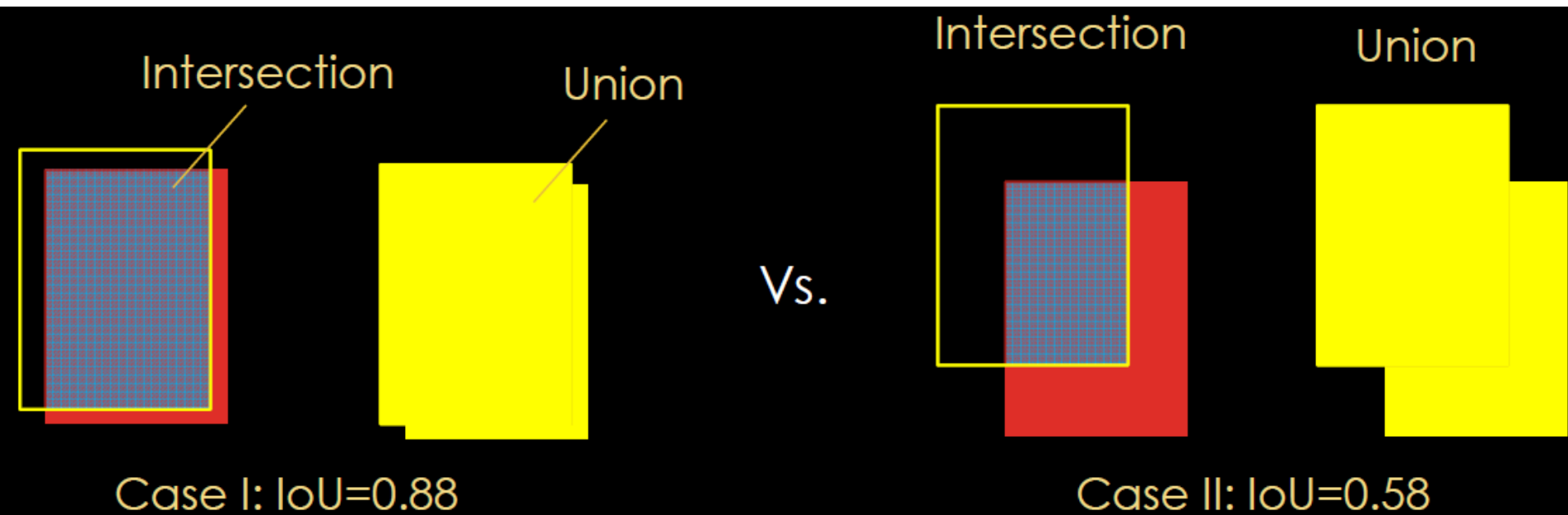  - **Recall** 1.0　0.9　0.8　0.7　0.6　0.5　0.4　0.3　0.2　0.1　0.0
  - **Precision** 0.70 0.74 0.78 0.82 0.85 0.89 0.93 0.96 0.98 0.99 1.00
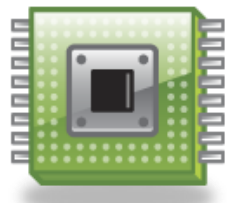  - **AP=(0.7+0.74+…+1)/11=0.88**

$$mAP = \frac{\sum_{i=1}^{n} AP_i}{n} \quad \text{for n classes}$$

# How Does IoU Affect AP?

- **Judging criteria of a nice shot**

# Commonly Used Indexes

- **AP-50: IoU=0.5 as the threshold**
  - **Both case I (IoU=0.88) and case II (IoU=0.58) get 1 TP**
- **AP-75: IoU=0.75 as the threshold**
  - **Case I (IoU=0.88) is TP, but case II (IoU=0.58) is not**
  - **Besides losing 1 TP, case II generates 1 FP and 1 FN simultaneously**
- **AP@[0.5 : 0.95]: from IoU=0.5 to IoU=0.95 with a step size of 0.05 (adopted in COCO dataset)**