



Introduction to Deep Learning

Chapter 5: Convolutional Neural Networks (Part 1)

1

Pao-Ann Hsiung

National Chung Cheng University

Contents

- ▶ Introduction
- ▶ Classic Networks
- ▶ Object Detection
- ▶ YOLO Algorithm
- ▶ Face Recognition
- ▶ Neural Style Transfer

Introduction

- ▶ Deep Learning for Computer Vision
 - ▶ **Rapid advances** in computer vision due to deep learning
 - ▶ **Cross-fertilization** between computer vision and other fields for deep learning
 - ▶ For example, speech recognition vs. computer vision

Computer Vision Problems

Image Classification



→ Cat? (0/1)

Object detection



Neural Style Transfer



Deep Learning on Large Images

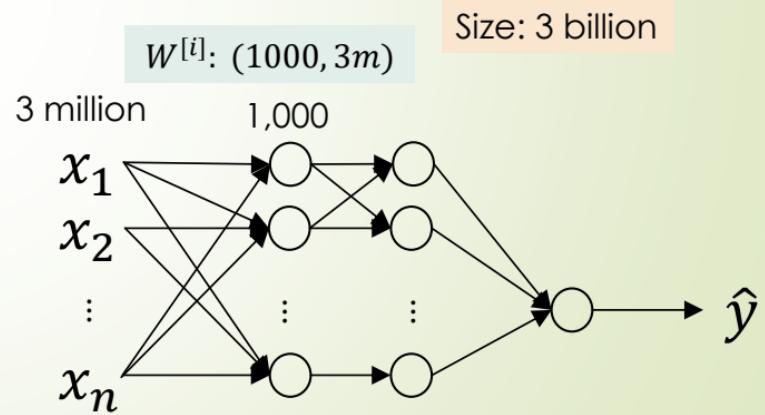


→ Cat? (0/1)

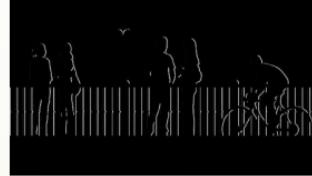
$64 \times 64 \times 3 = 12288$



$1000 \times 1000 \times 3 = 3 \text{ million}$

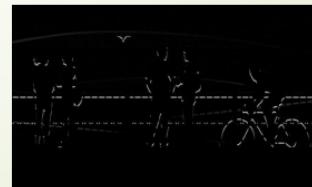


Computer Vision Problem



垂直邊界

vertical edges



水平邊界

horizontal edges

Vertical edge detection

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

convolution



為0所以vertical

1	0	-1
1	0	-1
1	0	-1

3 x 3 filter (kernel)

$$\begin{array}{c}
 \begin{array}{cccc}
 -5 & -4 & 0 & 9 \\
 -10 & -2 & 2 & 3 \\
 0 & -2 & -4 & -7 \\
 -3 & -2 & -3 & -16
 \end{array} \\
 =
 \end{array}$$

Vertical edge detection examples

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0



0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10



*

1	0	-1
1	0	-1
1	0	-1



=

0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0



*

1	0	-1
1	0	-1
1	0	-1



=

0	-30	-30	0
0	-30	-30	0
0	-30	-30	0
0	-30	-30	0



Vertical edge detection

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0

*

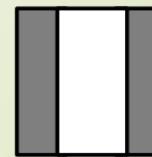
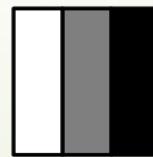
1	0	-1
1	0	-1
1	0	-1

=

0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0



*



Vertical and Horizontal Edge Detection

1	0	-1
1	0	-1
1	0	-1

Vertical

較混亂

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10

*

1	1	1
0	0	0
-1	-1	-1

1	1	1
0	0	0
-1	-1	-1

Horizontal

越邊數字越大

0	0	0	0
30	10	-10	-30
30	10	-10	-30
0	0	0	0

Learning to detect edges

透過filter偵測不同東西

1	0	-1
1	0	-1
1	0	-1

1	0	-1
2	0	-2
1	0	-1

3	0	-3
10	0	-10
3	0	-3

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

*

w_1	w_2	w_3
w_4	w_5	w_6
w_7	w_8	w_9

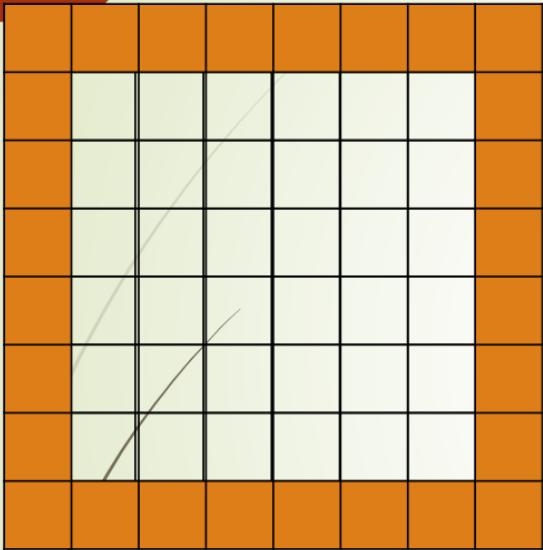
=

Sobel Filter

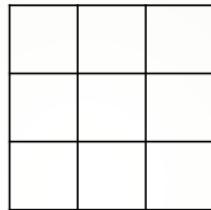
Scharr Filter

Padding

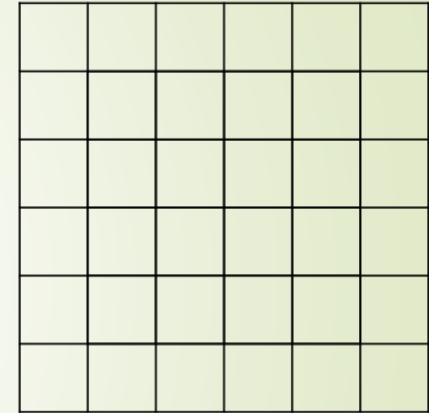
如果要保持原本的圖片大小



*



=



Size of input data: $n \times n$,

Size of filter: $f \times f$

Size of output: $(n - f + 1) \times (n - f + 1)$

Example: $6 - 3 + 1 = 4$, hence 4×4 , **size is reduced**

中間比較常使用，加上padding讓原本size跟圖片一樣大小

Use padding, extra border of 1 all around, ($p = 1$) gives output of

$$(n + 2p - f + 1) \times (n + 2p - f + 1) = 6 \times 6 \text{ (same as original data size)}$$

Valid and Same Convolutions

- ▶ “**Valid**” (no padding): $n \times n * f \times f \rightarrow (n - f + 1) \times (n - f + 1)$ 有縮小
 - ▶ $6 \times 6 * 3 \times 3 \rightarrow 4 \times 4$
- ▶ “**Same**”: Pad so that output size is the same as the input size Input 與 output 一樣
 - ▶ $n + 2p - f + 1 = n \Rightarrow p = \frac{f-1}{2}$ 透過 filter 決定 padding
 - ▶ $f = 3, p = \frac{3-1}{2} = 1, \text{ or } f = 5, p = 2$

Strided Convolutions

2	3	3	4	7	3	4	4	6	3	2	4	9	4
6	1	6	0	9	1	8	0	7	1	4	0	3	2
3	-3	4	4	8	3	3	4	8	3	9	4	7	4
7	1	8	0	3	1	6	0	6	1	3	0	4	2
4	-3	2	4	1	3	8	4	3	3	4	4	6	4
3	1	2	0	4	1	1	0	9	1	8	0	3	2
0	-1	1	0	3	-3	9	0	2	-3	1	0	4	3

*

3	4	4
1	0	2
-1	0	3

 3×3

圖片相較於走1次的小

91	100	83
69	91	127
44	71	74

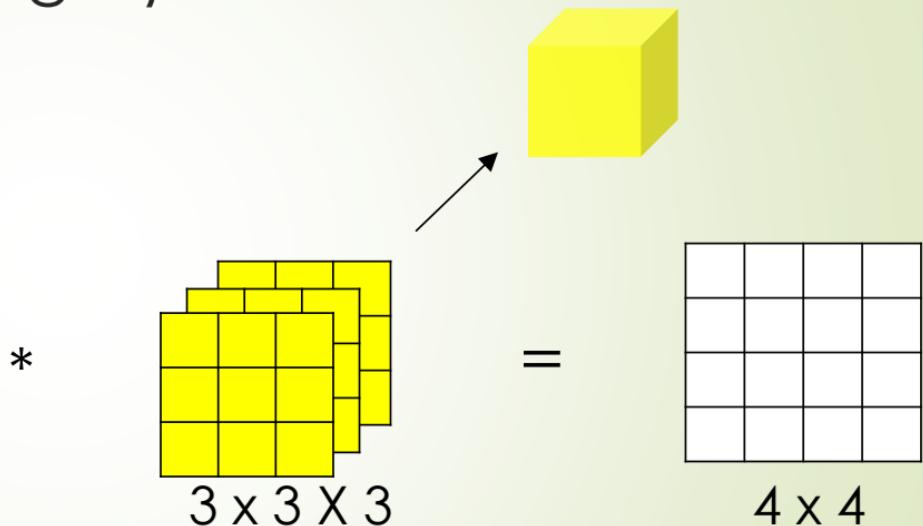
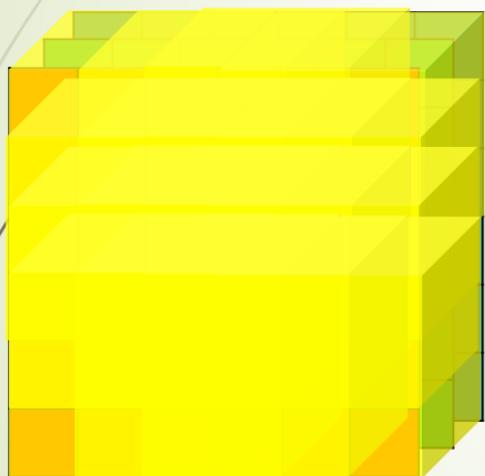
Stride = 2 每次走兩個pixel

 7×7

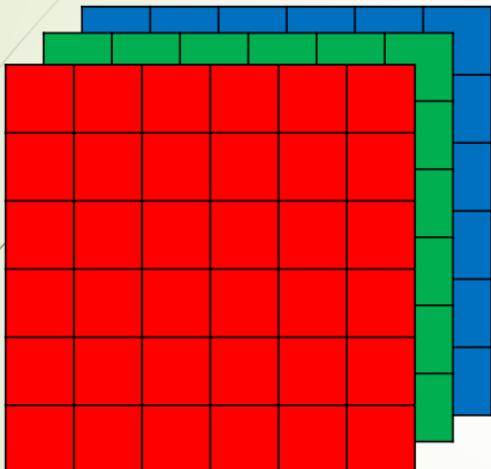
$$n \times n * f \times f \text{ (padding } p, \text{ stride } s) \Rightarrow \lfloor \frac{n + 2p - f}{s} + 1 \rfloor \times \lfloor \frac{n + 2p - f}{s} + 1 \rfloor$$

$$(7 + 0 - 3)/2 + 1 = 4/2 + 1 = 3$$

Convolutions over volumes (e.g., RGB images)



Multiple Filters



$$6 \times 6 \times 3$$

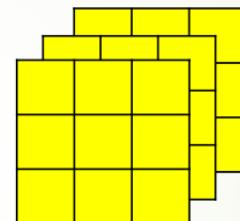
$$n \times n \times n_c * f \times f \times n_c$$

$$\Rightarrow (n - f + 1) \times (n - f + 1) \times n_f$$

$$6 \times 6 \times 3 * 3 \times 3 \times 3$$

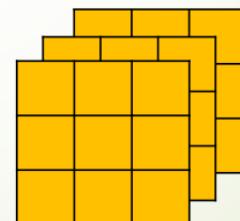
$$\Rightarrow 4 \times 4 \times 2 (\#filters)$$

*



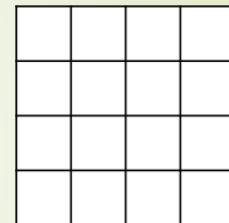
$$3 \times 3 \times 3$$

*



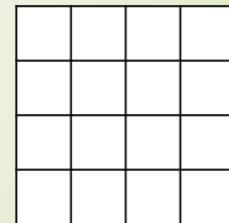
$$3 \times 3 \times 3$$

=



$$4 \times 4$$

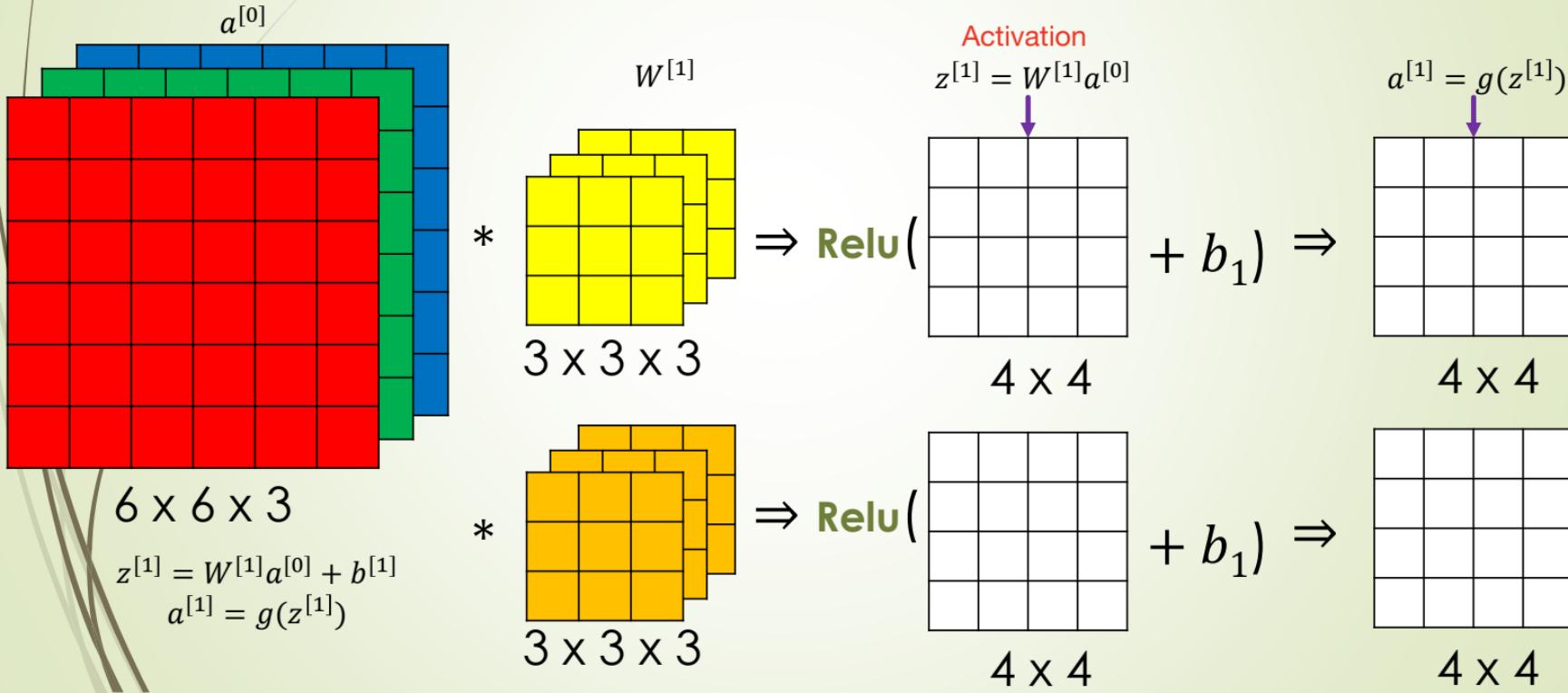
=



$$4 \times 4$$

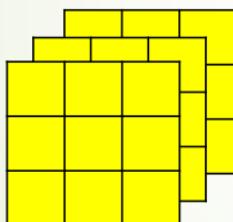
One Convolutional Layer

One Layer of Convolutional Neural Network

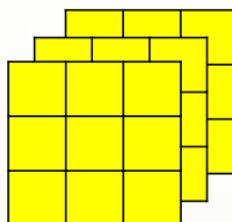


Number of Parameters in One Layer

- If you have 10 filters that are $3 \times 3 \times 3$ in one layer of a neural network, how many parameters does that layer have?

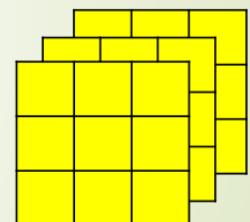


$3 \times 3 \times 3$ Bias



$3 \times 3 \times 3$

• • •



$3 \times 3 \times 3$

- $(3 \times 3 \times 3 + 1) \times 10 = 280$ parameters
- No matter how large the input image is, the number of parameters is fixed to 280 for 10 filters of size $3 \times 3 \times 3$.

Summary of notation

If layer 1 is a convolution layer:

- ▶ $f^{[l]}$ = filter size
- ▶ $p^{[l]}$ = padding
- ▶ $s^{[l]}$ = stride
- ▶ $n_c^{[l]}$ = number of filters
- ▶ Each filter is: $f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$
- ▶ Input Size: $n_H^{[l-1]} \times n_W^{[l-1]} \times n_c^{[l-1]}$
- ▶ Output Size: $n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$
- ▶ $n_H^{[l]} = \left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$
- ▶ $n_W^{[l]} = \left\lfloor \frac{n_W^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$

Activations: $a^{[l]} \Rightarrow n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$

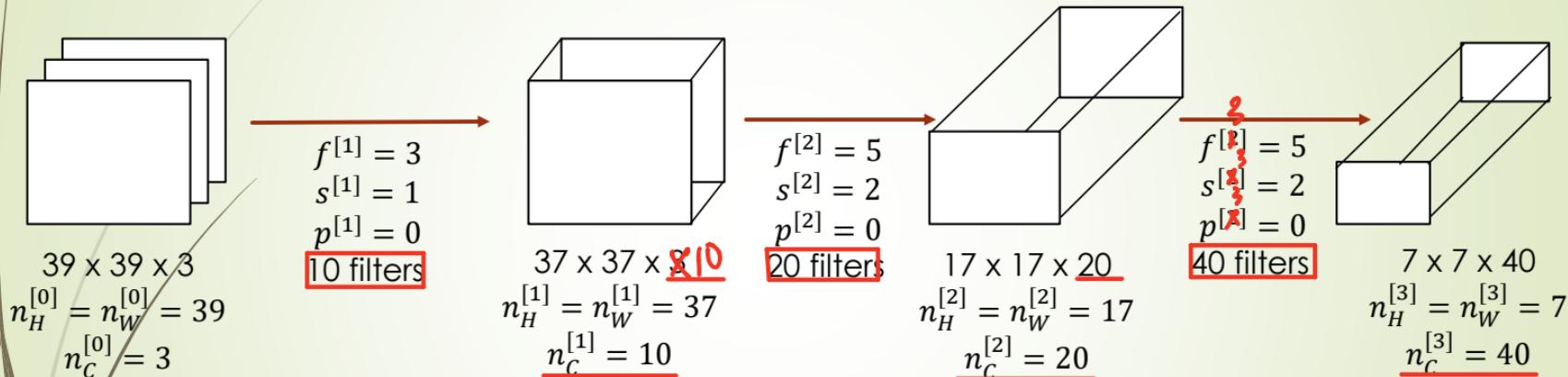
Weights: $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$

Bias: $n_c^{[l]}$

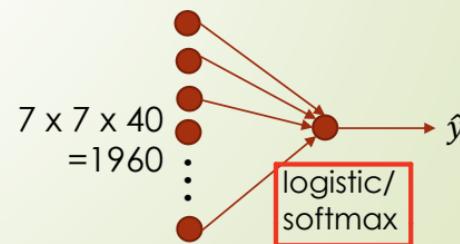
Example ConvNet

圖片變小，channel變多。

因為要辨識物件



$$\frac{n + 2p - f}{s} + 1 = \frac{39 + 0 - 3}{1} + 1 = 37$$



Types of layer in a convolutional network

- ▶ Convolutional (CONV)
- ▶ Pooling (POOL)
- ▶ Fully connected (FC)

Pooling Layer

Pooling Layer: Max Pooling

4個中最大

1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2

相素太大找最大的值縮小參數

Hyperparameters

$$f = 2$$

$$s = 2$$

Max Pooling



9	2
6	3

No parameters!

Pooling Layer: Max Pooling

1	3	2	1	3
2	9		1	5
1				2
8	3		1	0
5	6	1	2	9

5 x 5

Hyperparameters

$$\begin{aligned}f &= 3 \\s &= 1\end{aligned}$$

Max Pooling

9	9	5
9	9	5
8	6	9

3 x 3

Note: For multiple channels, the above max pooling is done for each channel

Pooling Layer: Average Pooling

1	3	2	1
2	9	1	1
1	4	2	3
5	6	1	2

Hyperparameters

$$\begin{array}{l} f = 2 \\ s = 2 \end{array}$$


375	125
4	2

Summary of Pooling

Pooling:減少size

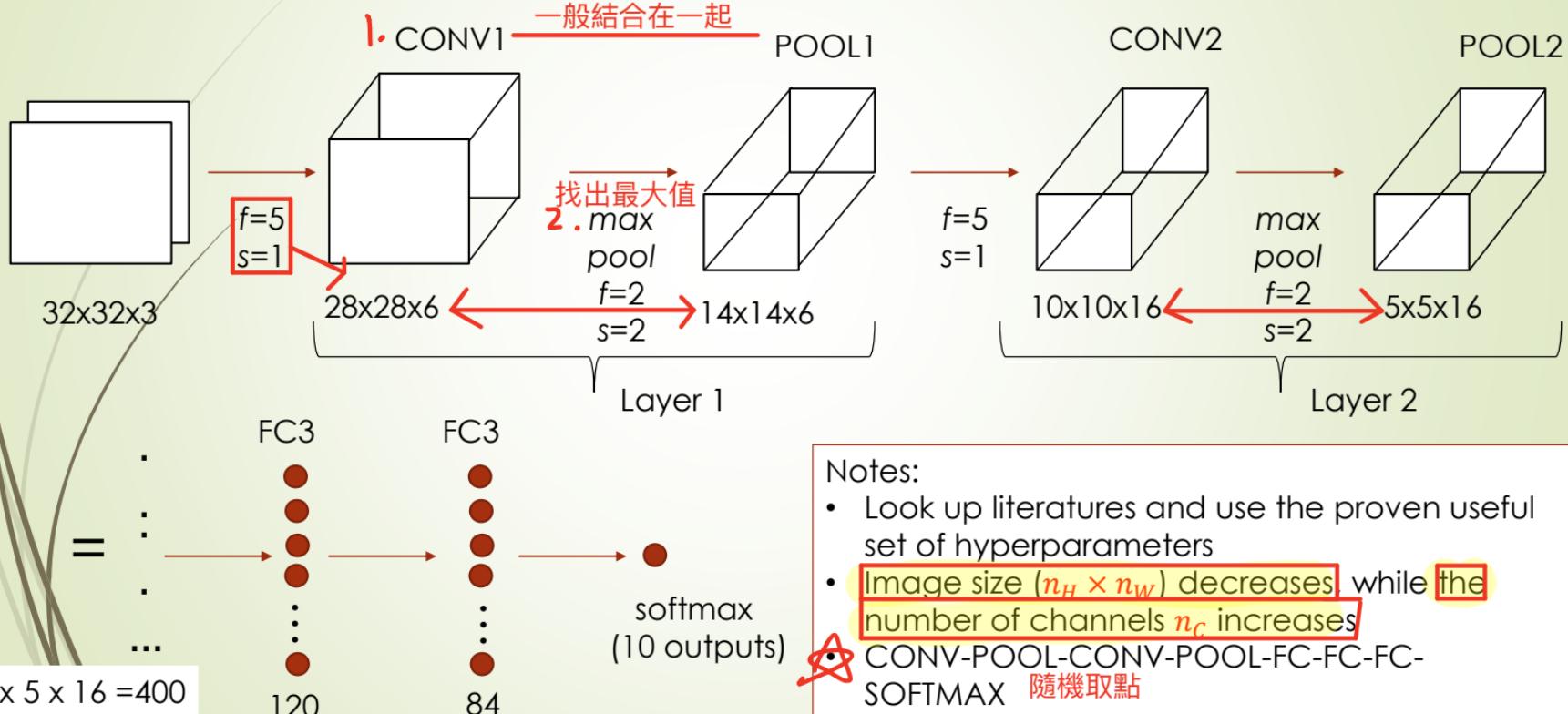
Padding:維持大小

Hyperparameters

- ▶ f : filter size
- ▶ s : stride
- ▶ Max or average pooling
- ▶ Usually, $p = 0$, no padding

$$n_H \times n_W \times n_C \rightarrow \left\lfloor \frac{n_H - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_H - f}{s} + 1 \right\rfloor \times n_C$$

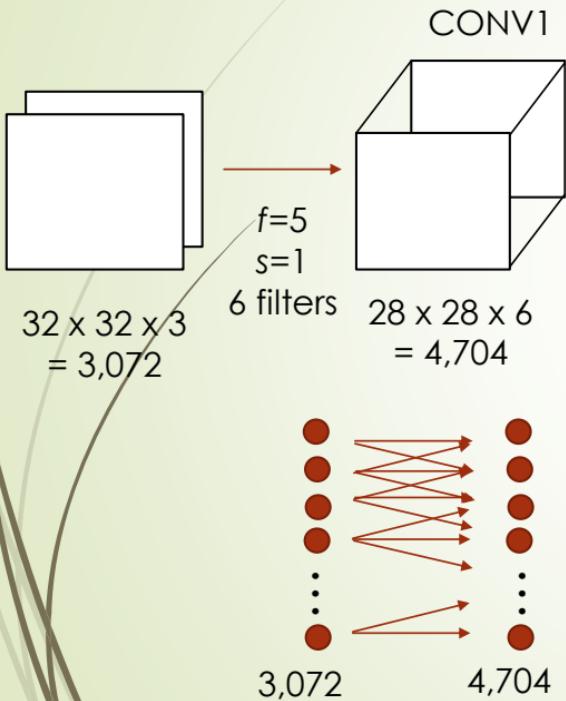
Neural Network Example (aka LeNet-5)



Neural Network Example

	Activation Shape	Activation Size	# Parameters
Input	(32, 32, 3)	3,072	0
CONV1 (f=5, s=1)	(28, 28, 8)	6,272	208
POOL1	(14, 14, 8)	1,568	0
CONV2 (f=5, s=1)	(10, 10, 16)	1,600	416
POOL2	(5, 5, 16)	400	0
FC3	(120, 1)	120	<u>48.001</u>
FC4	(84, 1)	84	<u>10.081</u> 隨機取點
Softmax	(10, 1)	10	841

Why Convolutions work?



$$6 \times (5 \times 5 + 1) = 156 \text{ parameters}$$

縮小參數

$$3,072 \times 4,704 \approx 14M$$

Large
Difference!!

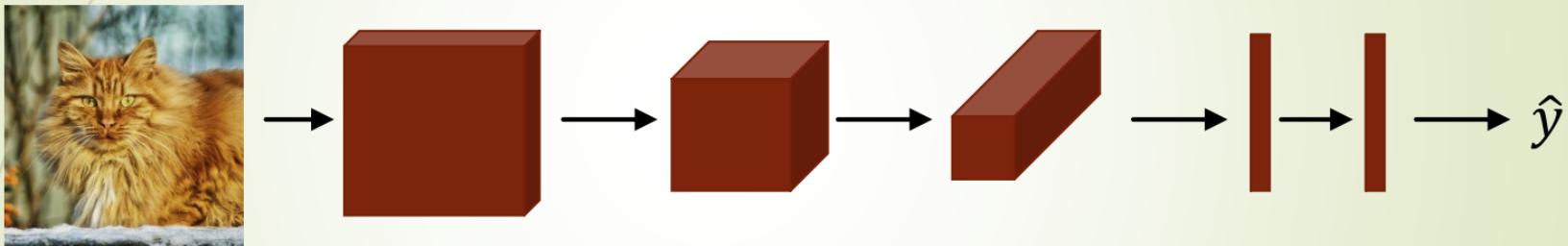
Why Convolutions work?

$$\begin{array}{|c|c|c|c|c|c|} \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 \end{array} * \begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array} = \begin{array}{|c|c|c|c|} \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 \end{array}$$

- ▶ **Parameter Sharing:** A feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image 使用同一個filter，可以移植到其他部分
- ▶ **Sparsity of connections:** In each layer, each output value depends only on a small number of inputs 一個像素點會與附近有相關

Putting it together

Training set $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$.



$$\text{Cost } J = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Use gradient descent to optimize parameters to reduce J

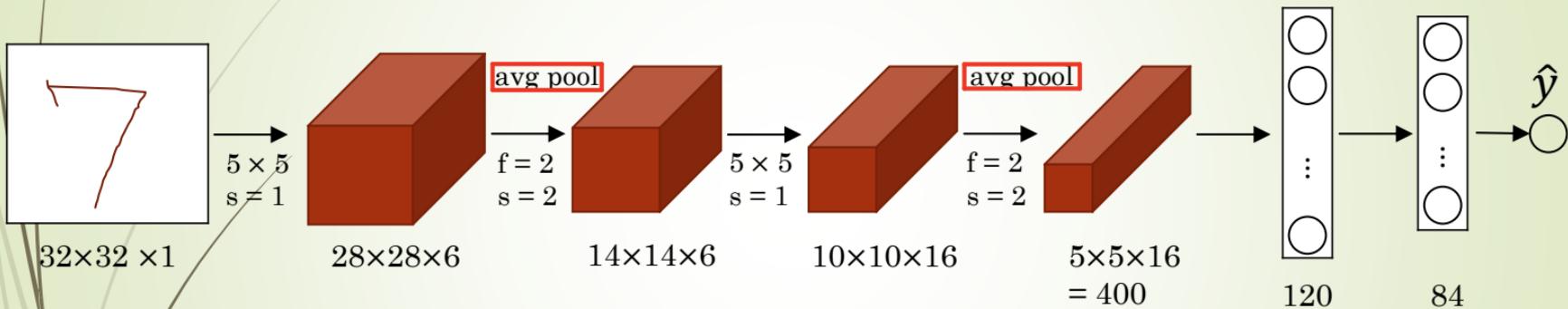
Case Studies

Le-Net5, AlexNet, VGG-16, ResNet

Why look at case studies?

- ▶ Often a trained NN can work also quite good in another application!
- ▶ Can use transfer learning
- ▶ Classic Networks
 - ▶ LeNet-5
 - ▶ AlexNet
 - ▶ VGG
- ▶ ResNet (Residual Network) with 152 layers (very deep)
- ▶ Inception

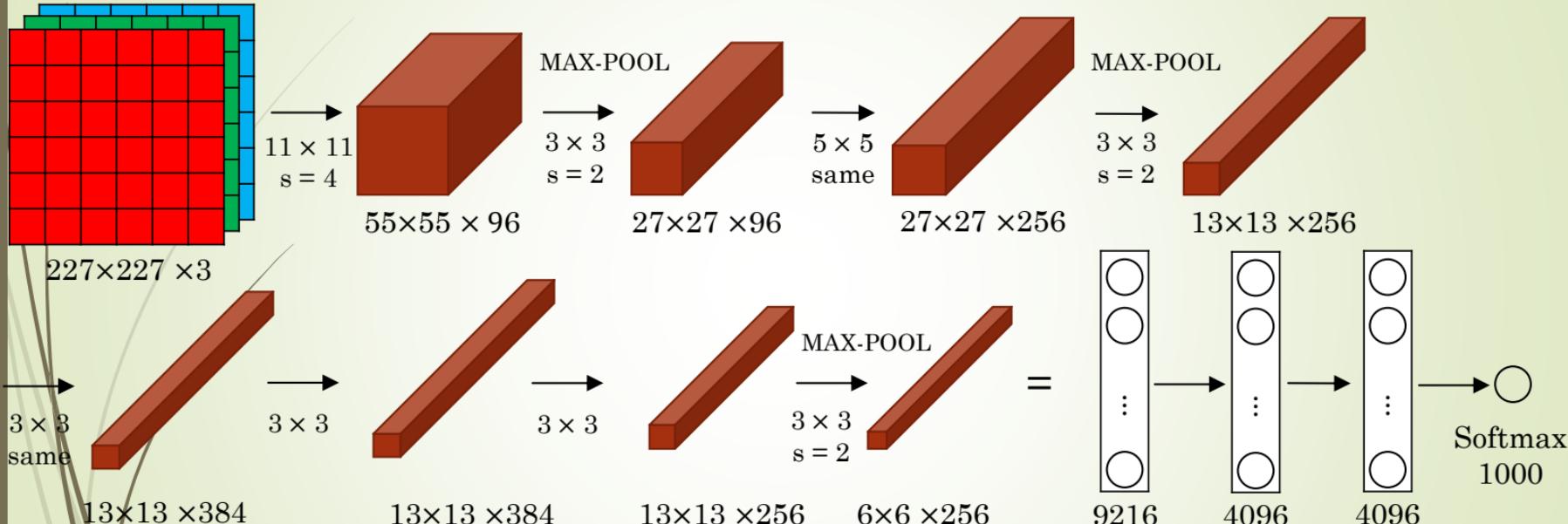
LeNet-5 像素辨識



- ▶ 60K parameters (small by modern standards)
 - ▶ n_H, n_W decreased with layers, n_C increased
 - ▶ CONV – POOL – CONV – POOL – FC – FC - Output
- ↑
Softmax

AlexNet

辨識的東西越來越多

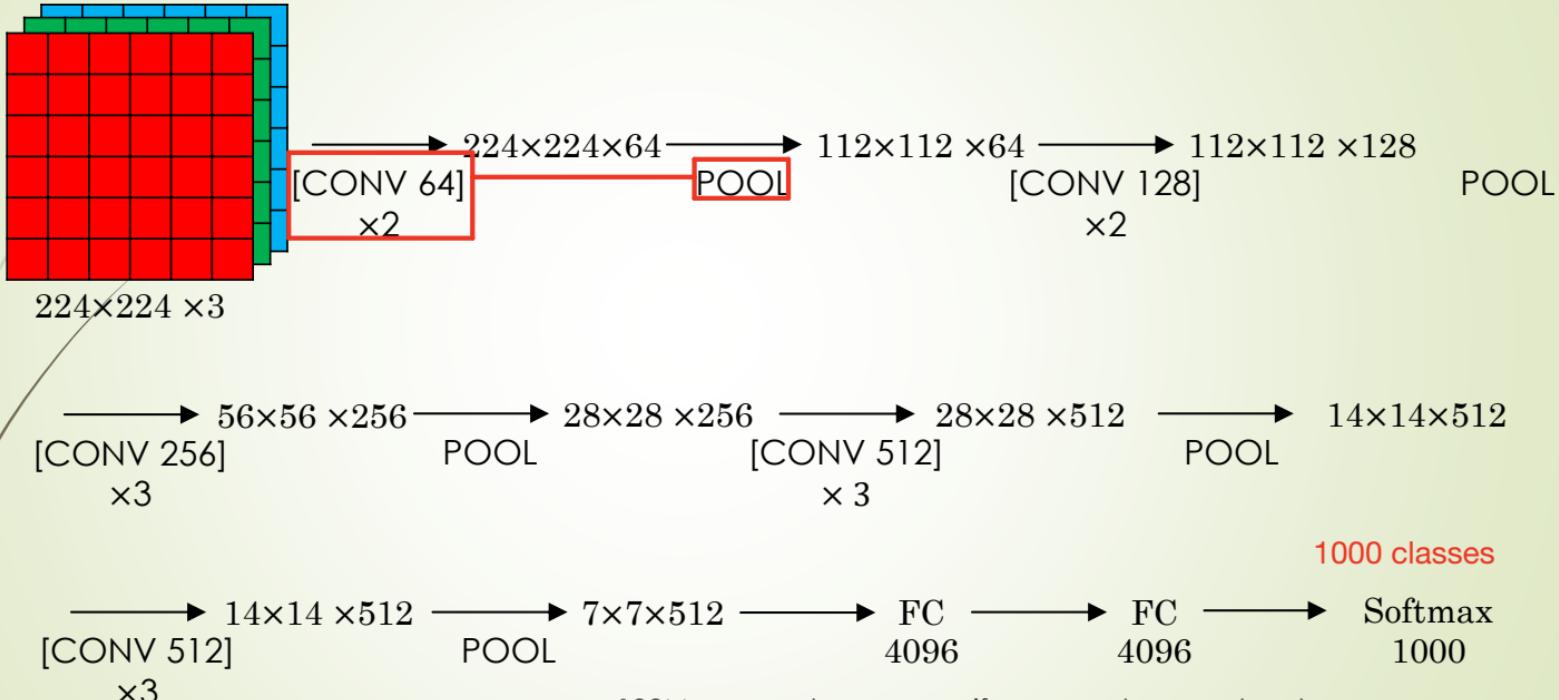


- Similar to LeNet-5, but much bigger (60M parameters)
- Instead of Sigmoid/Tanh, AlexNet used ReLU
- Multiple GPUs: in the olden days, it was difficult to train on multiple GPUs, it is fine now!
- Local Response Normalization (LRN): less useful now!

可以辨識1000個classes

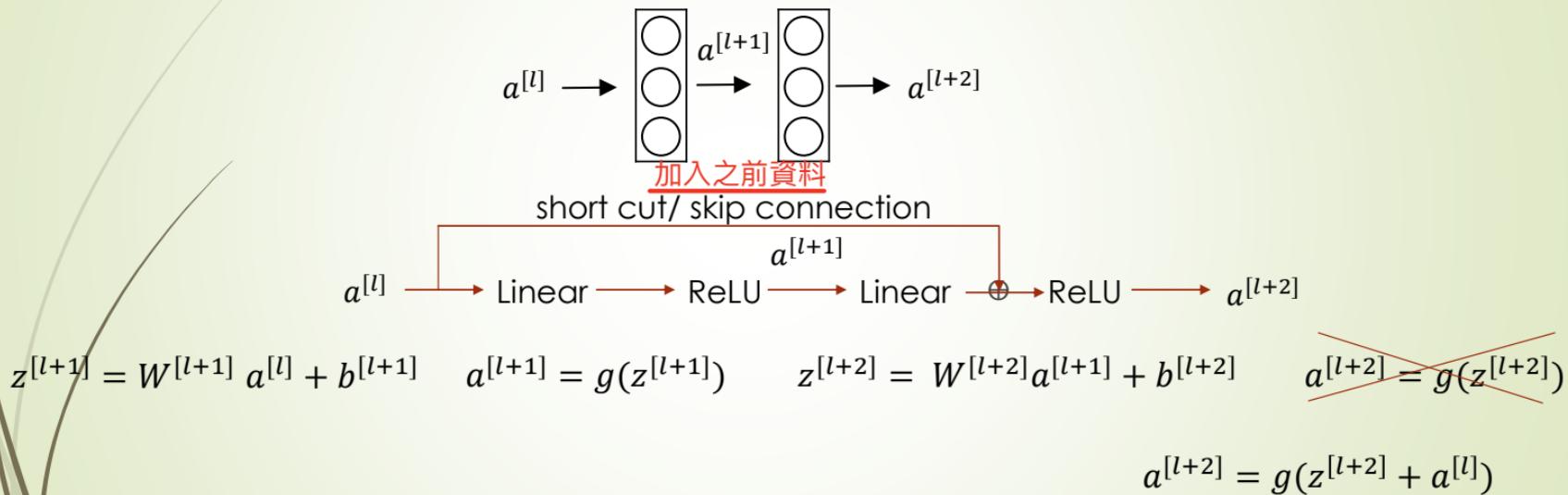
VGG-16

37

CONV = 3×3 filter, $s = 1$, sameMAX-POOL = 2×2 , $s = 2$ 

- ~138M parameters, very uniform, very large network
- VGG-19 a bigger network, almost similar to VGG-16 in performance

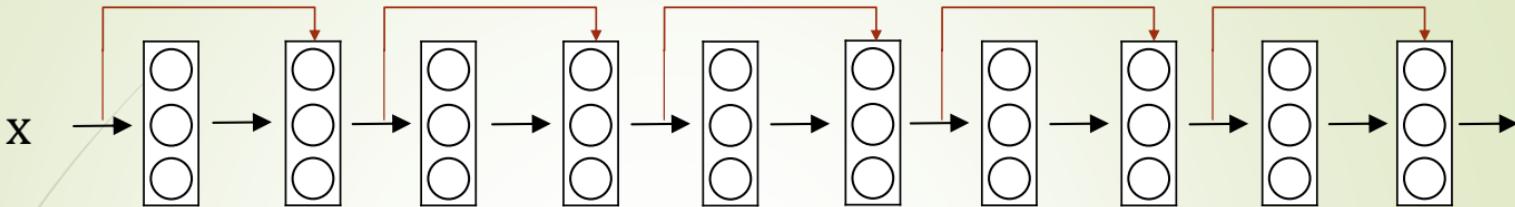
Residual Block (for Residual Networks or ResNets)



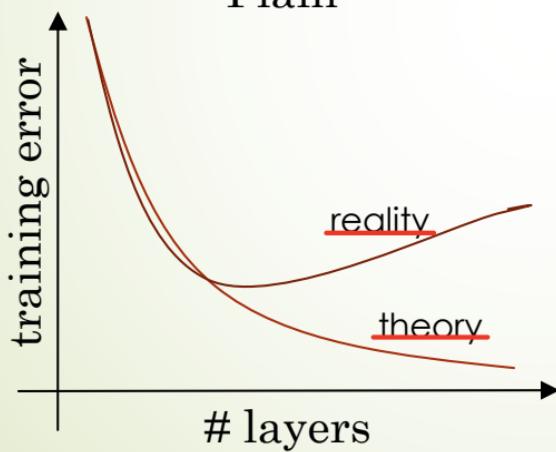
Residual Network

第一個超過100layer

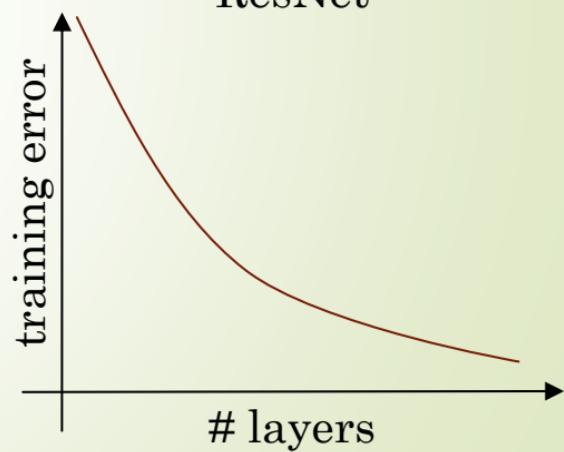
39



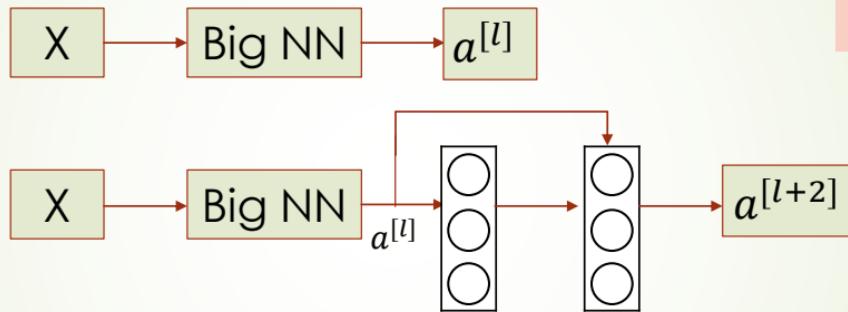
Plain



ResNet



Why do residual networks work?



Identity function is easy
for Residual Block to
learn!!!

之前學的有保留下來

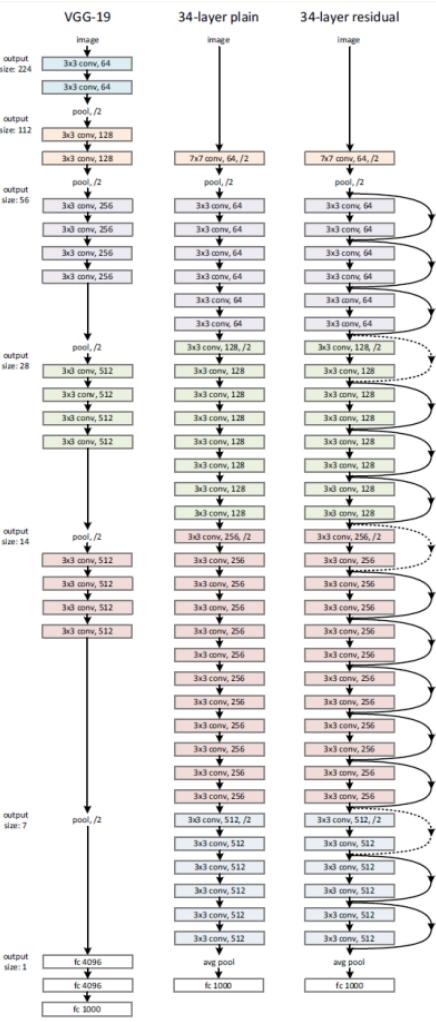
256 → ReLU, where $a \geq 0$

→ $a^{[l+2]} = g(z^{[l+2]} + a^{[l]}) = g(W^{[l+2]}a^{[l+1]} + W_s a^{[l]}) = g(a^{[l]}) = a^{[l]}, \text{ if } W^{[l+2]} = b^{[l+2]} = 0$

$\Re^{256 \times 128}$

128

ResNet



- ▶ Lots of 3×3 same CONV
 - ▶ Dimension preserved
 - ▶ Thus, $z^{[l+2]} + a^{[l]}$ makes sense

[He et al., 2015. Deep residual networks for image recognition]

Network in network and 1x1 convolution

What does a 1x1 convolution do?

1	2	3	6	5	8
3	5	5	1	3	4
2	1	3	4	9	3
4	7	8	5	7	9
1	5	3	7	4	8
5	4	9	8	3	5

*

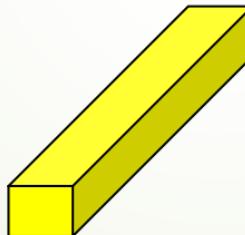
2

=

2	4	6		

 $6 \times 6 \times 32$

*

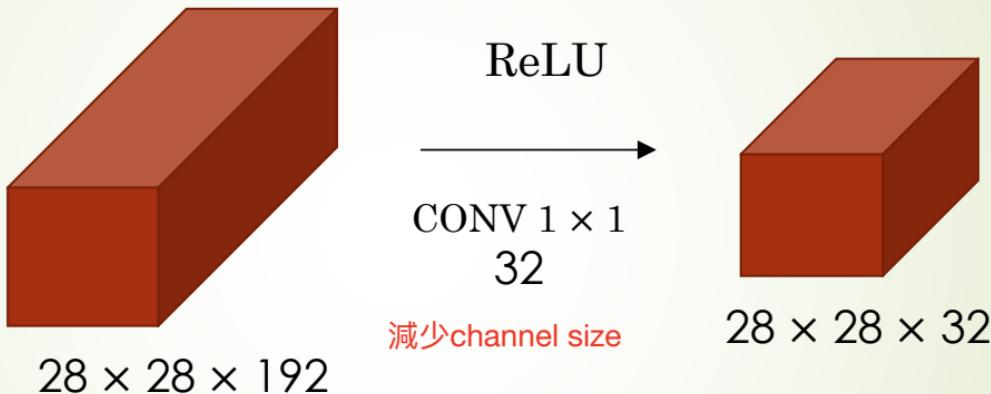
 $1 \times 1 \times 32$

=

 $6 \times 6 \times \# \text{ filters}$

Using 1x1 convolutions

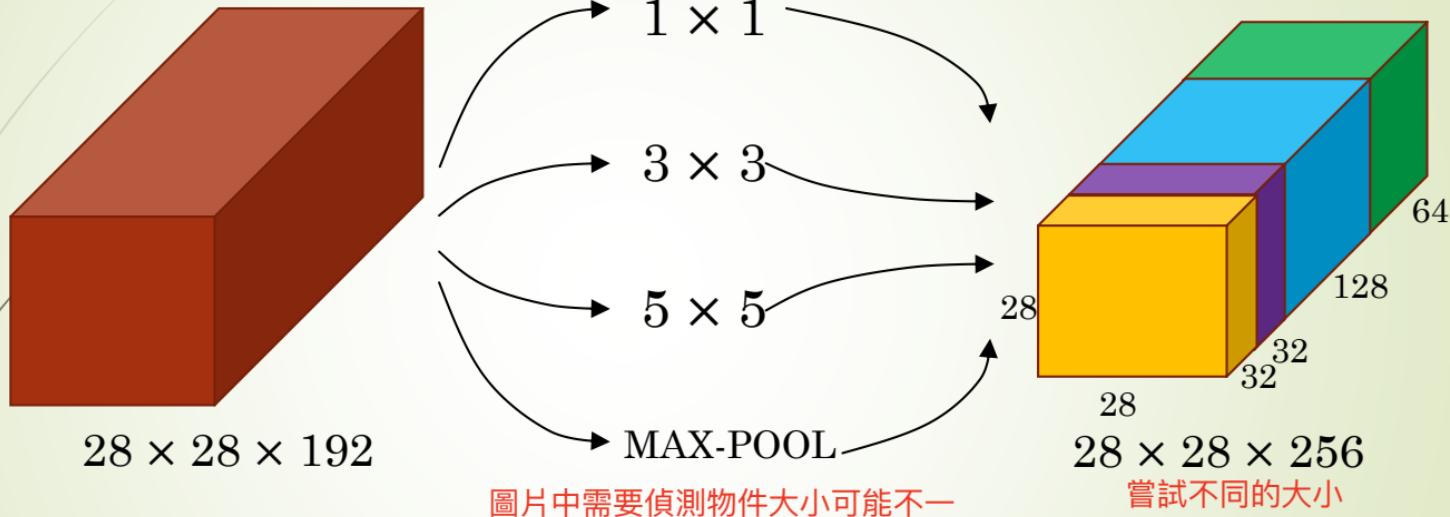
降低深度，但不降低原輸入二維的維度



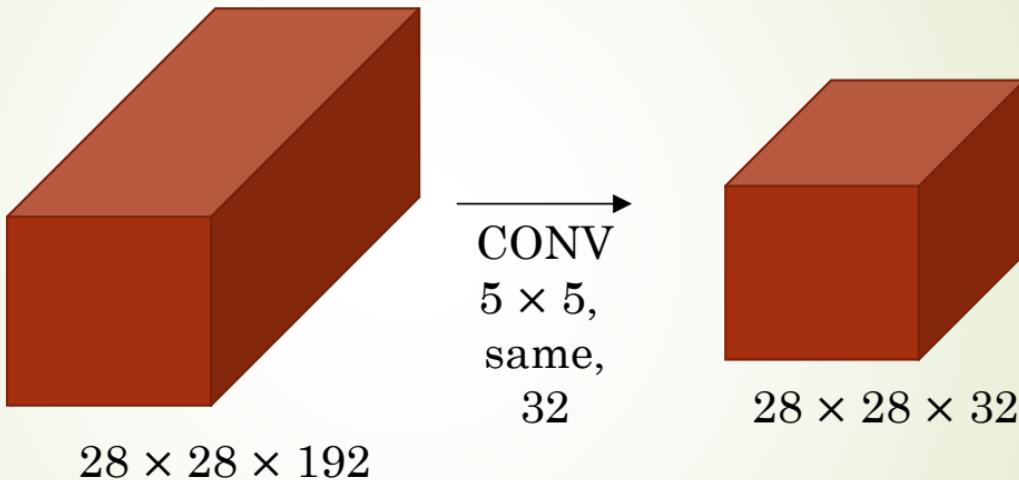
- ▶ Number of channels reduced from 192 to 32
非線性
- ▶ Adds non-linearity (when #channels not reduced)

Inception Network

Motivation for inception network

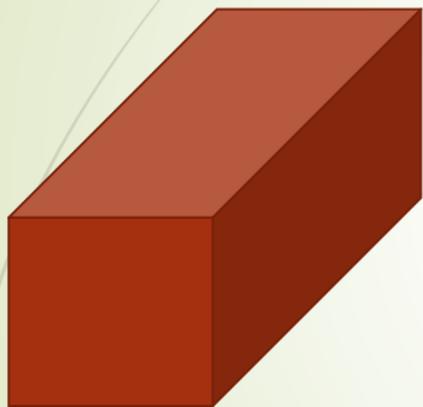


The problem of computational cost

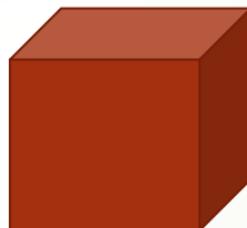


- ▶ 32 filters, each filter is of size $5 \times 5 \times 192$
- ▶ #multiplies = $(28 \times 28 \times 32) \times (5 \times 5 \times 192) = 120 \text{ M}$

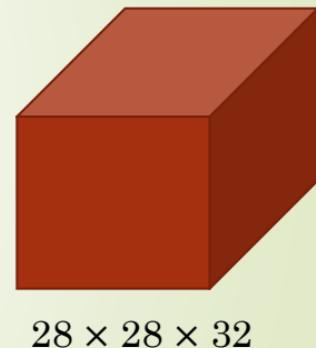
Using 1×1 convolution



$\xrightarrow{\text{CONV}}$
 $1 \times 1,$
16,
 $1 \times 1 \times 192$



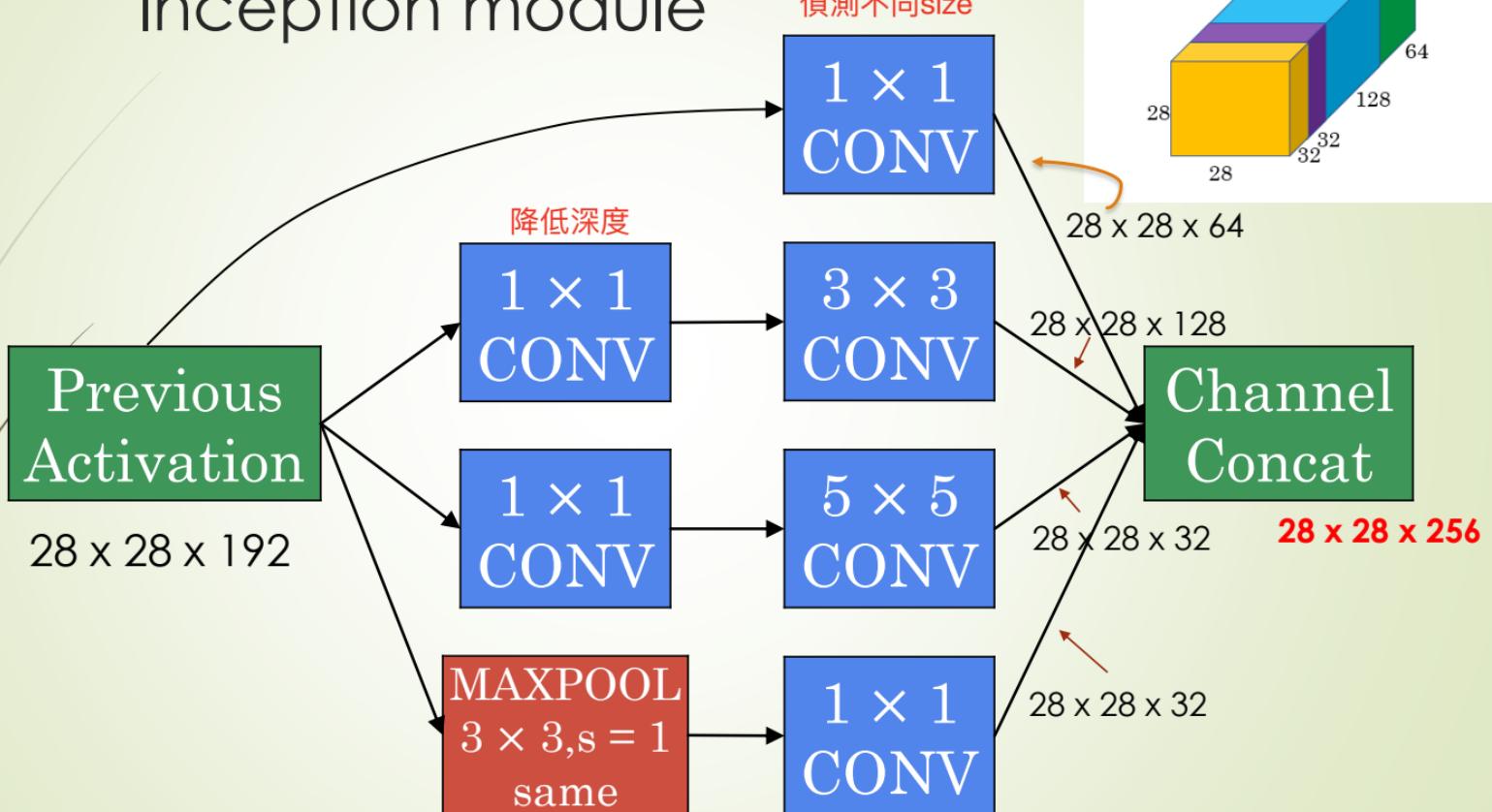
$\xrightarrow{\text{CONV}}$
 $5 \times 5,$
32,
 $5 \times 5 \times 16$

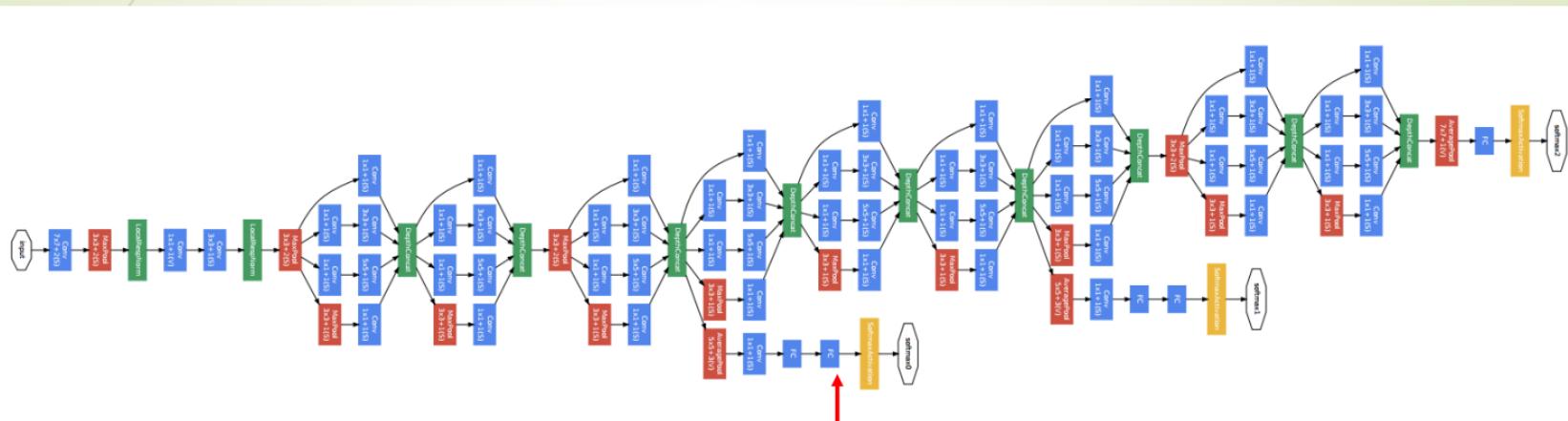


- ▶ $(28 \times 28 \times 16) \times 192 = 2.4 \text{ M}$
- ▶ $(28 \times 28 \times 32) \times (5 \times 5 \times 16) = 10 \text{ M}$
- ▶ $2.4 + 10 = 12.4 \text{ M}$ multiplications (1/10 that of original 5×5)
- ▶ #additions similar to #multiplications

Inception Network

Inception module





Side branches uses softmax to produce predictions

Initially, Google developed GoLeNet!

From the Inception Movie



Object Detection

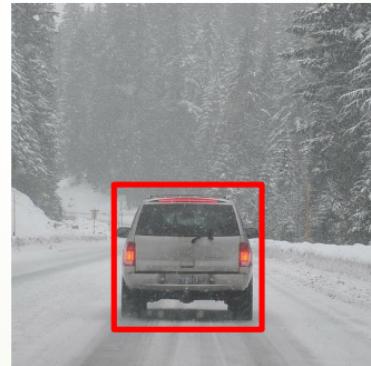
What are localization and detection?

Image classification



“Car”

Classification with
localization



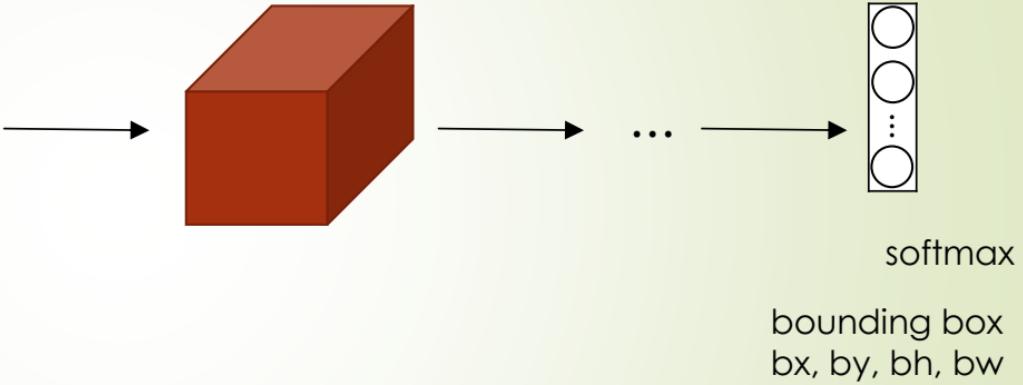
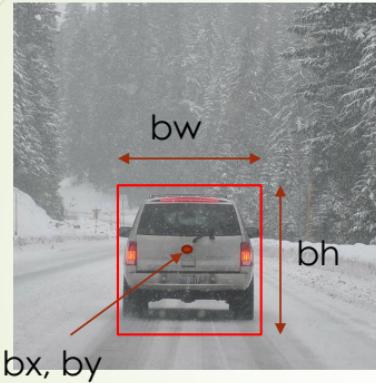
“Car”

Detection



“Multiple Objects”

Classification with Localization



- 1 - pedestrian 先找出中點
- 2 - car 再找出邊長
- 3 - motorcycle
- 4 - background

$bx = 0.5$
 $by = 0.7$
 $bh = 0.3$
 $bw = 0.4$

Defining the target label y

- ▶ 1. pedestrian
- ▶ 2. car
- ▶ 3. motorcycle
- ▶ 4. background

- ▶ Need to output bx , by , bh , bw ,
class label (1-4)

Is there an object?

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

Don't care.
沒有在object裡

$$\mathcal{L}(\hat{y}, y) = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2 \quad \text{if } y_1 = 1$$

$$(\hat{y}_1 - y_1)^2 \quad \text{if } y_0 = 1$$

Landmark Detection



b_x, b_y, b_h, b_w

- ▶ $l_{1x}, l_{1y}, l_{2x}, l_{2y}, l_{3x}, l_{3y}, \dots l_{64x}, l_{64y}$: key **landmarks**
- ▶ For **facial emotions**, **body postures**, etc.
- ▶ Landmarks should be **consistent** across different input images



偵測情緒



Body posture

Car detection example

一個圖片有很多物件



透過CNN偵測物件

Training set: (closely cropped images)

x

y



1 需要裁切圖片



1



1



0



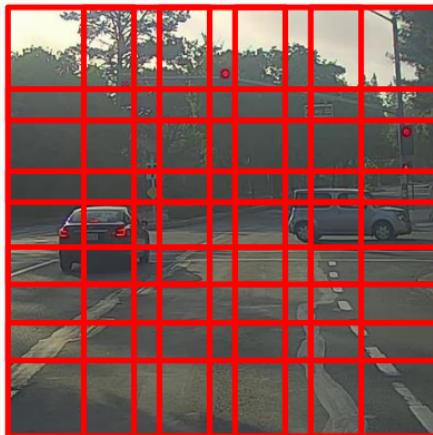
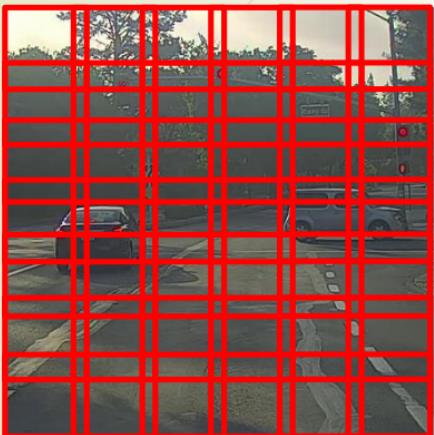
0



Conv
Net

y

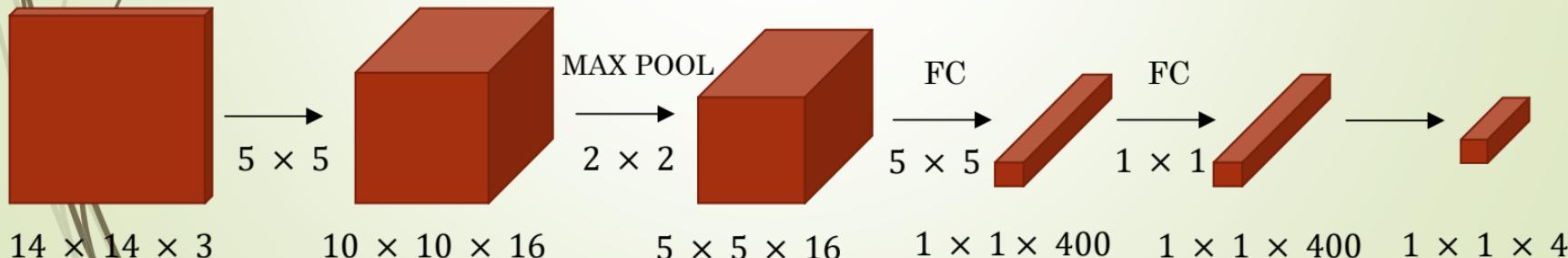
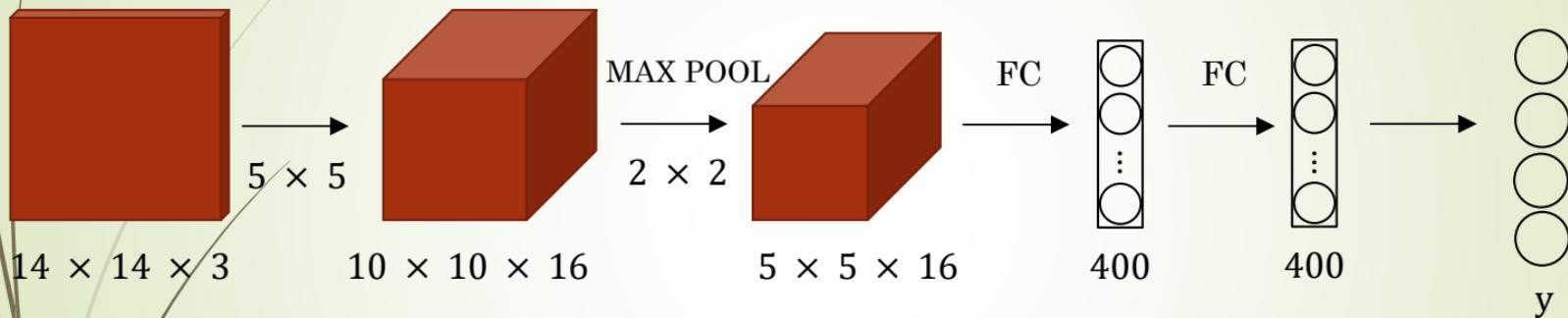
Sliding windows detection



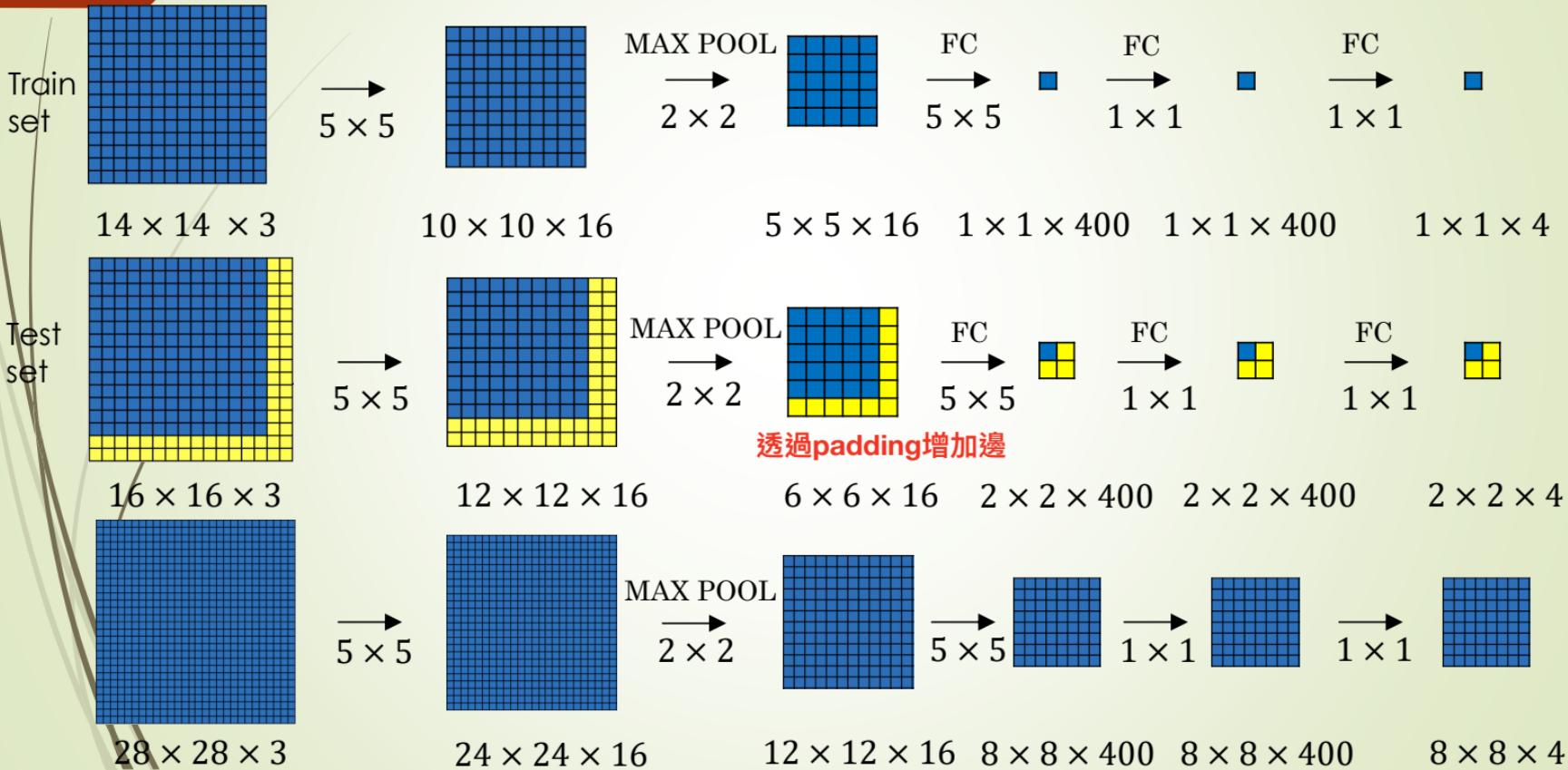
- Computational cost is high due to sliding windows (infeasibly slow)
- Solution: Implement convolutionally!

Object Detection: Convolutional Implementation of Sliding Windows

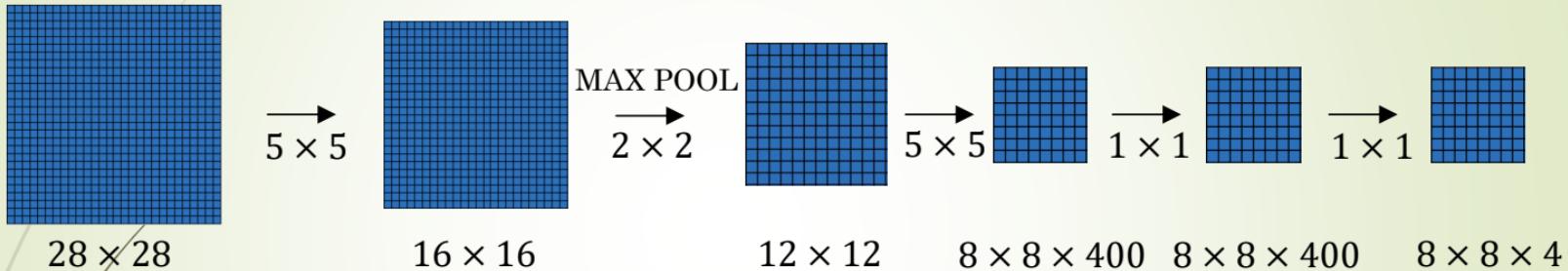
Turning FC layer into convolutional layers



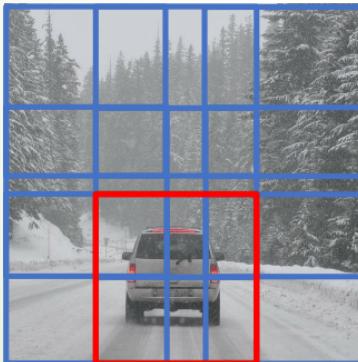
Convolutional Implementation of Sliding Windows



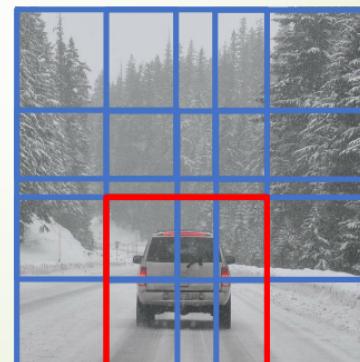
Convolutional Implementation of Sliding Windows



根據filter偵測車子



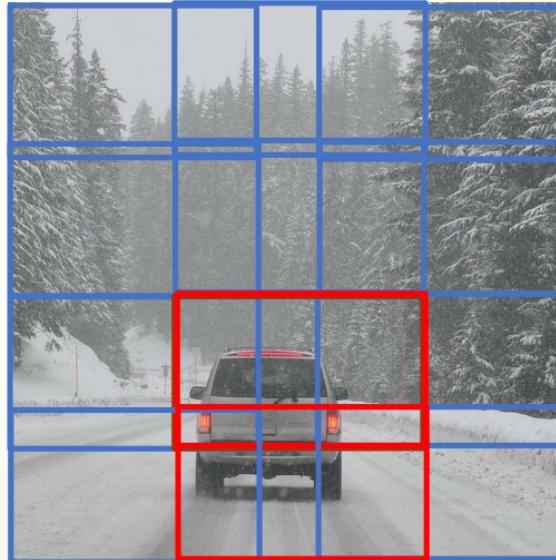
透過許多CNN偵測物件



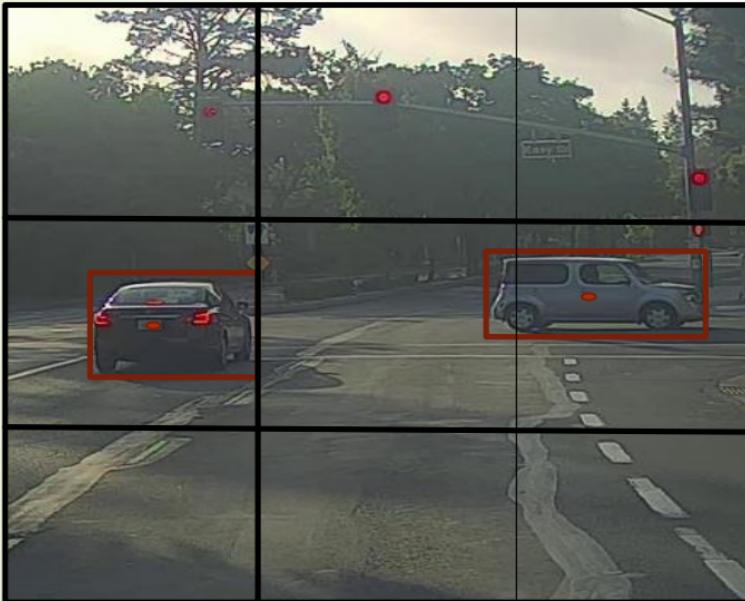
Bounding
Box not very
accurate!

Output accurate bounding boxes

- ▶ None of the boxes matches the car
- ▶ How to output more accurate bounding boxes?



YOLO Algorithm (You Only Look Once)



$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix} = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

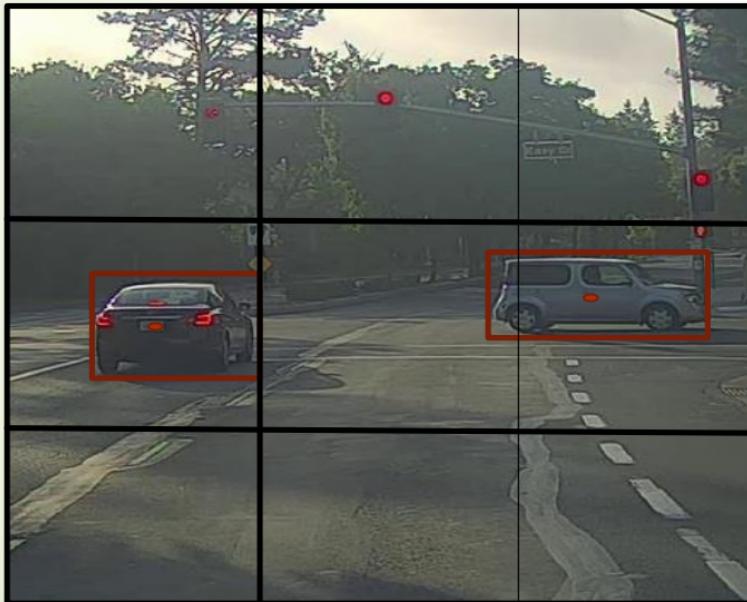
- ▶ Output size: $3 \times 3 \times 8$
 - ▶ Can be finer: $19 \times 19 \times 8$
- ▶ Object is assigned to only that grid cell that contains the **center** of the object
- ▶ Very efficient and popular!

Specify the bounding boxes

Bound box

在偵測

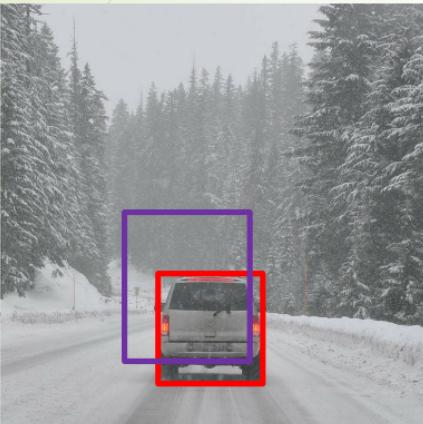
Car



$$\mathbf{y} = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0.4 \\ 0.3 \\ 0.9 \\ 0.5 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Object Detection: Intersection over union

Evaluating object localization



Intersection over Union (IoU)
= (size of $\square \cap \square$) / (size of $\square \cup \square$)

使用兩個box的交集

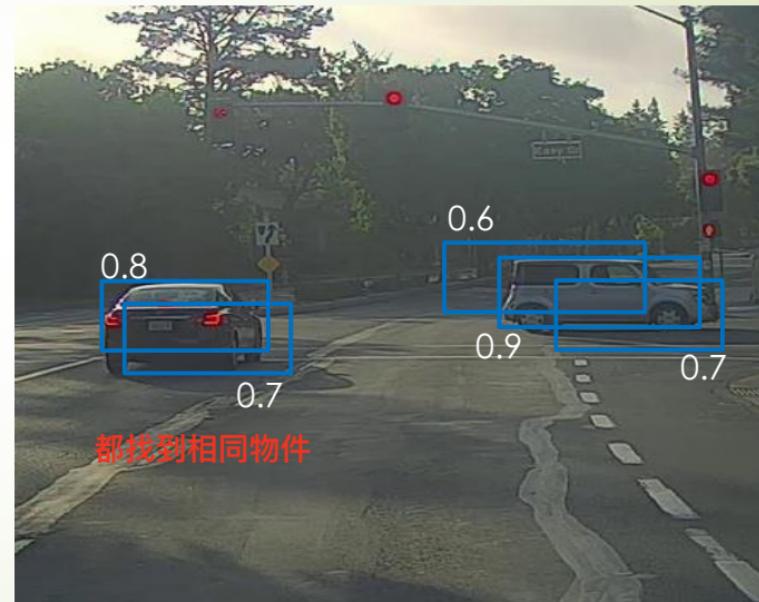
“Correct” if IoU ≥ 0.5 如果交集大於0.5

More generally, IoU is a measure of the overlap between two bounding boxes.

Object Detection Non-Max Suppression

Non-max suppression example

- More than one grid cells claim that to have one object (multiple detections) 找出機率大的
- Leave only the grid with the max probability p_c , and suppress the grids with non-max claims



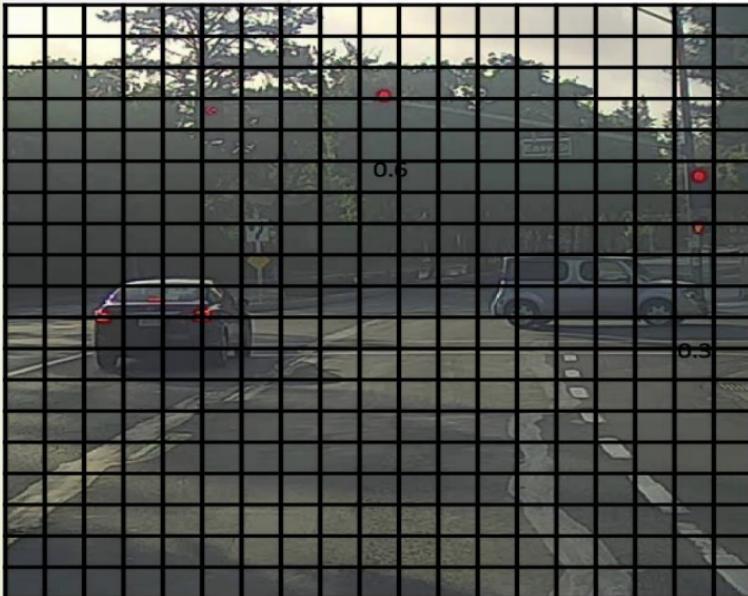
Non-max suppression algorithm

YOLO

找bounding box

使用IOU

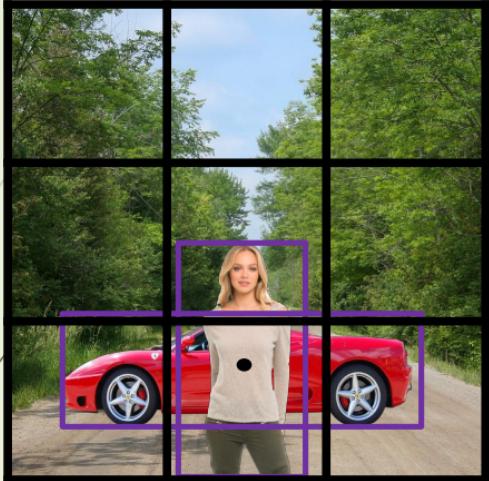
最後使用Non-max 方法



- ▶ Each output prediction is $\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$
- ▶ Discard all boxes with $p_c \leq 0.6$.
- ▶ While there are any remaining boxes:
 - ▶ Pick the box with the **largest p_c** . Output that as a prediction.
 - ▶ Discard any remaining box with **$\text{IoU} \geq 0.5$ with the box output in the previous step**.

Object Detection Anchor Boxes

Overlapping objects 兩個物件重疊



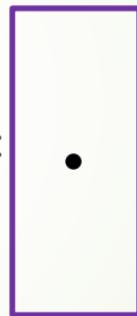
$$\mathbf{y} = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

- Problem: Centers overlap and cannot detect two objects using previous method

使用兩個 anchor box

- Method: Use 2 anchor boxes

Anchor box 1:



兩個物件不同大小
可以使用

Anchor box 2:



$$\mathbf{y} = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

For anchor
box 1

For anchor
box 2

Anchor box algorithm

Previously:

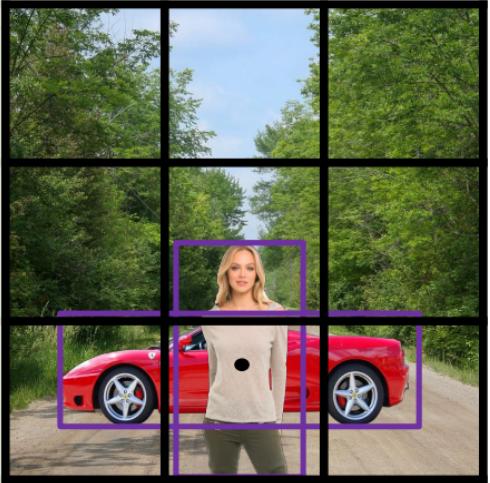
Each object in training image is assigned to grid cell that contains that object's midpoint.

- ▶ Previously, size of output y : $3 \times 3 \times 8$
- ▶ Now, size of output y : $3 \times 3 \times 16$

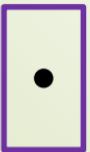
With two anchor boxes:

Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IoU.
i.e., assigned to (Grid cell, Anchor box)

Anchor box example



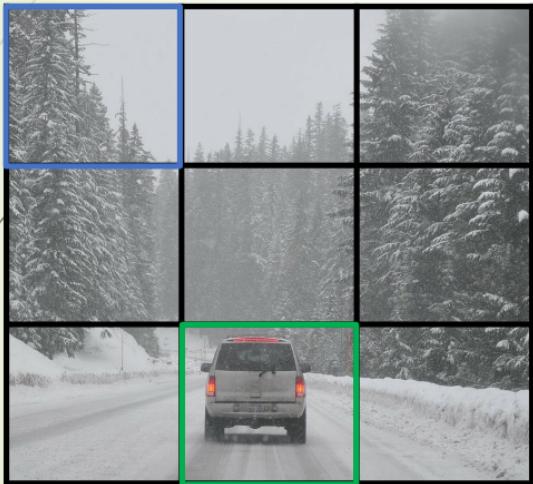
Anchor box 1: Anchor box 2:



$$\begin{aligned}
 \text{y} = & \begin{cases} \text{路人} \\ \text{車子} \end{cases} = \begin{cases} \begin{pmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{pmatrix} \\ \begin{pmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{pmatrix} \end{cases} \text{ or } \begin{cases} \begin{pmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 1 \\ 0 \\ 0 \end{pmatrix} \\ \begin{pmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 1 \end{pmatrix} \end{cases} \\
 & \begin{pmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{pmatrix} \quad \begin{pmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ ? \\ 0 \end{pmatrix}
 \end{aligned}$$

Object Detection YOLO Algorithm

Training



y is $3 \times 3 \times 2 \times 8$

↑
#anchors ↑
 5+#classes

- 1 - pedestrian
- 2 - car
- 3 - motorcycle

$y =$



Anchor
box 1



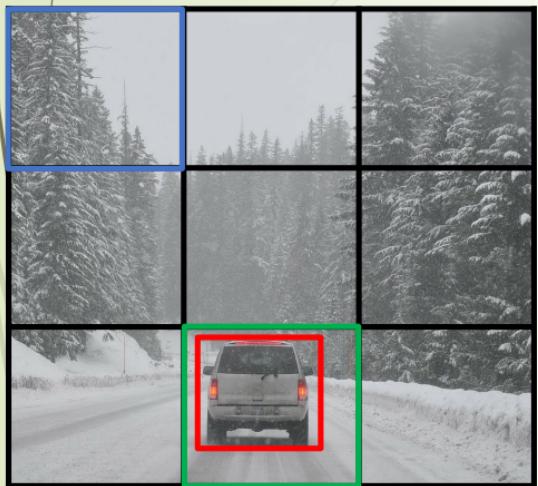
Anchor
box 2

$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

$$\begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

$$\begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Making Predictions



→

1

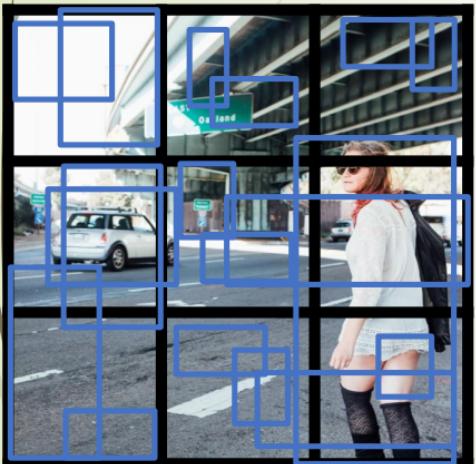
→

A small red rectangular prism is shown, representing a 3D volume element.

$$3 \times 3 \times 2 \times 8$$

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

Outputting the non-max suppressed outputs



- For each grid call, get 2 predicted bounding boxes.
- Get rid of low probability predictions.
- For each class (pedestrian, car, motorcycle) use non-max suppression to generate final predictions.

使用不同box

刪除機率低的IOU

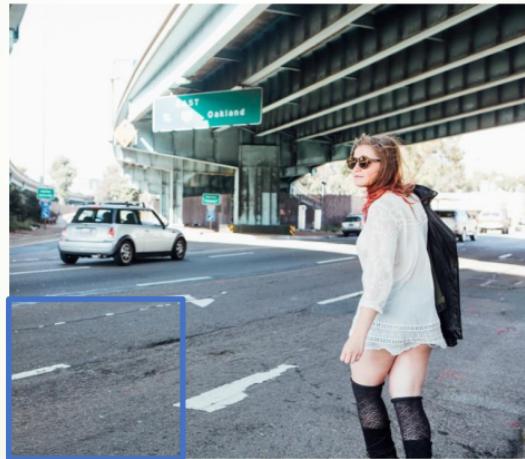
產出

The background features a dark vertical bar on the left and a light green gradient on the right. Overlaid are several thin, curved lines in various shades of brown and grey. A prominent red arrow points from the left towards the center, containing the white text '80'.

80

Region Proposal

Region proposal: R-CNN



只偵測藍色的部分



Faster algorithms

R-CNN:

Propose regions. Classify proposed regions one at a time. Output label + bounding box.

Fast R-CNN:

Propose regions. Use convolution implementation of sliding windows to classify all the proposed regions.

Faster R-CNN: Use convolutional network to propose regions.

[Girshik et. al, 2013. Rich feature hierarchies for accurate object detection and semantic segmentation]
[Girshik, 2015. Fast R-CNN]

[Ren et. al, 2016. Faster R-CNN: Towards real-time object detection with region proposal networks]