# Convolutional Neural Network (CNN)
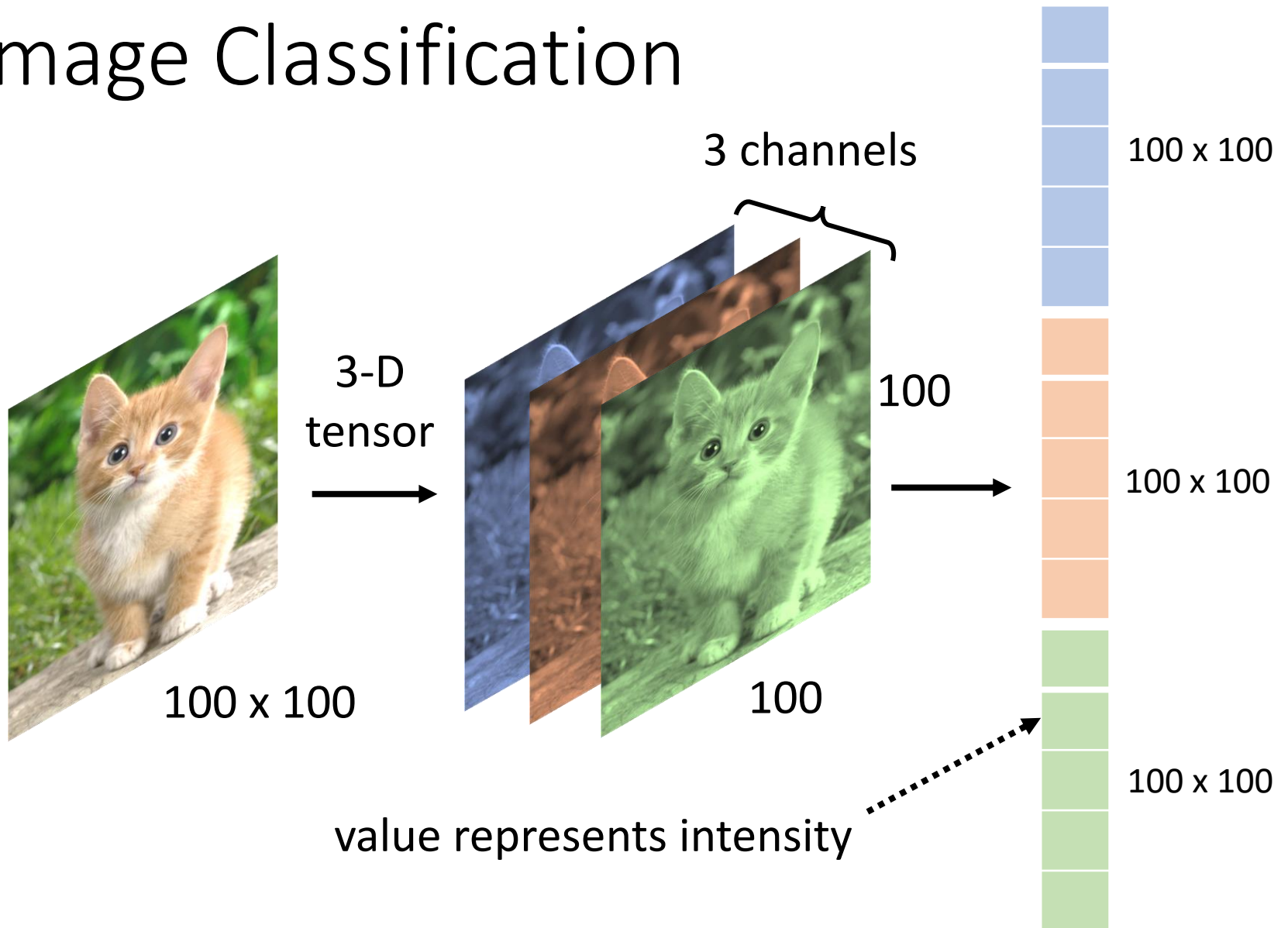
Network Architecture designed for Image

# Image Classification
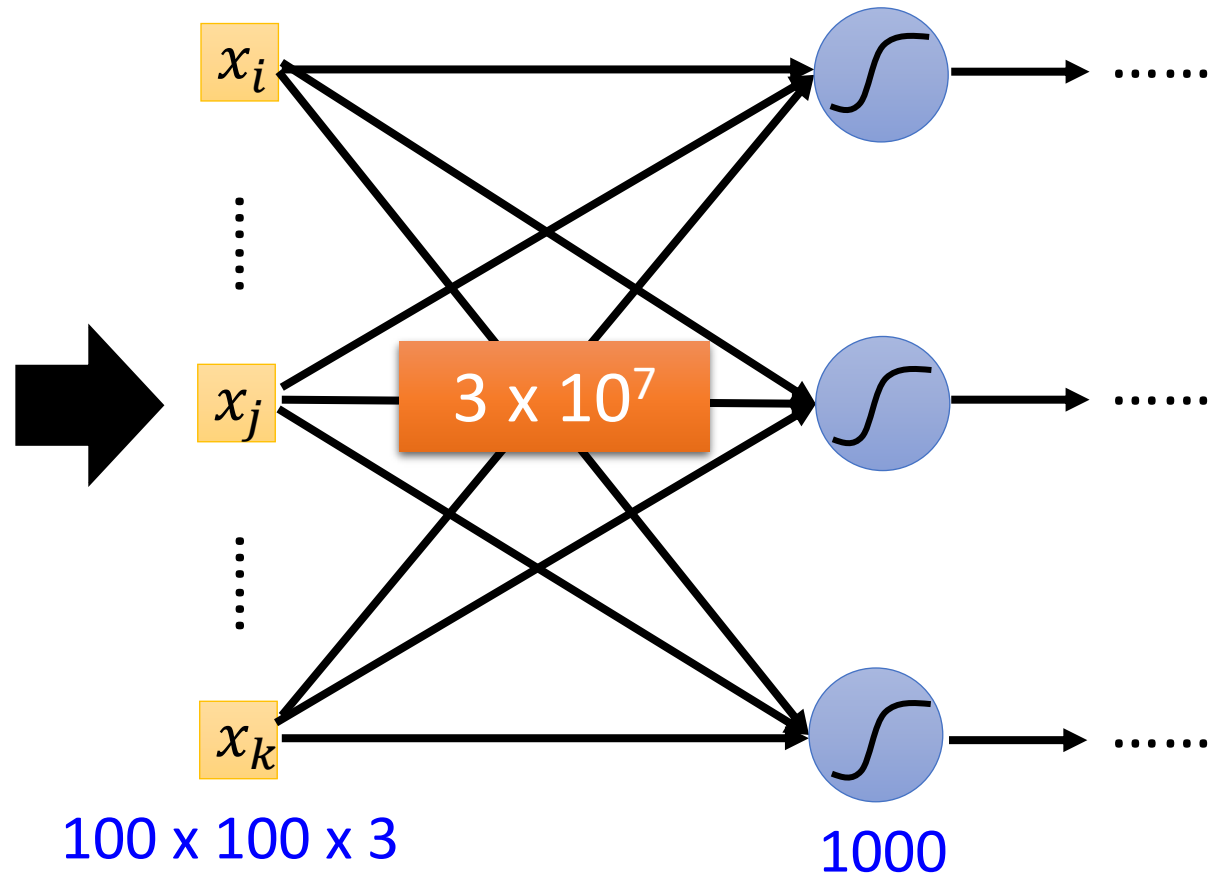


100 x 100

$$\begin{bmatrix} \vdots \\ 0.2 \\ 0.7 \\ 0.1 \\ \vdots \end{bmatrix} \quad \begin{matrix} \text{dog} \\ \text{cat} \\ \text{tree} \end{matrix} \quad \begin{bmatrix} \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \end{bmatrix}$$

Model

$y'$ ⟷ $\hat{y}$

Cross entropy

(All the images to be classified have the same size.)

# Image Classification



3 channels

3-D tensor

100

100 x 100

100

100 x 100

100 x 100

100 x 100

value represents intensity

若一層有1000個neuron，
就需要3e7個參數
很容易overfitting

100 x 100

## *Fully Connected Network*

$x_i$

$x_j$

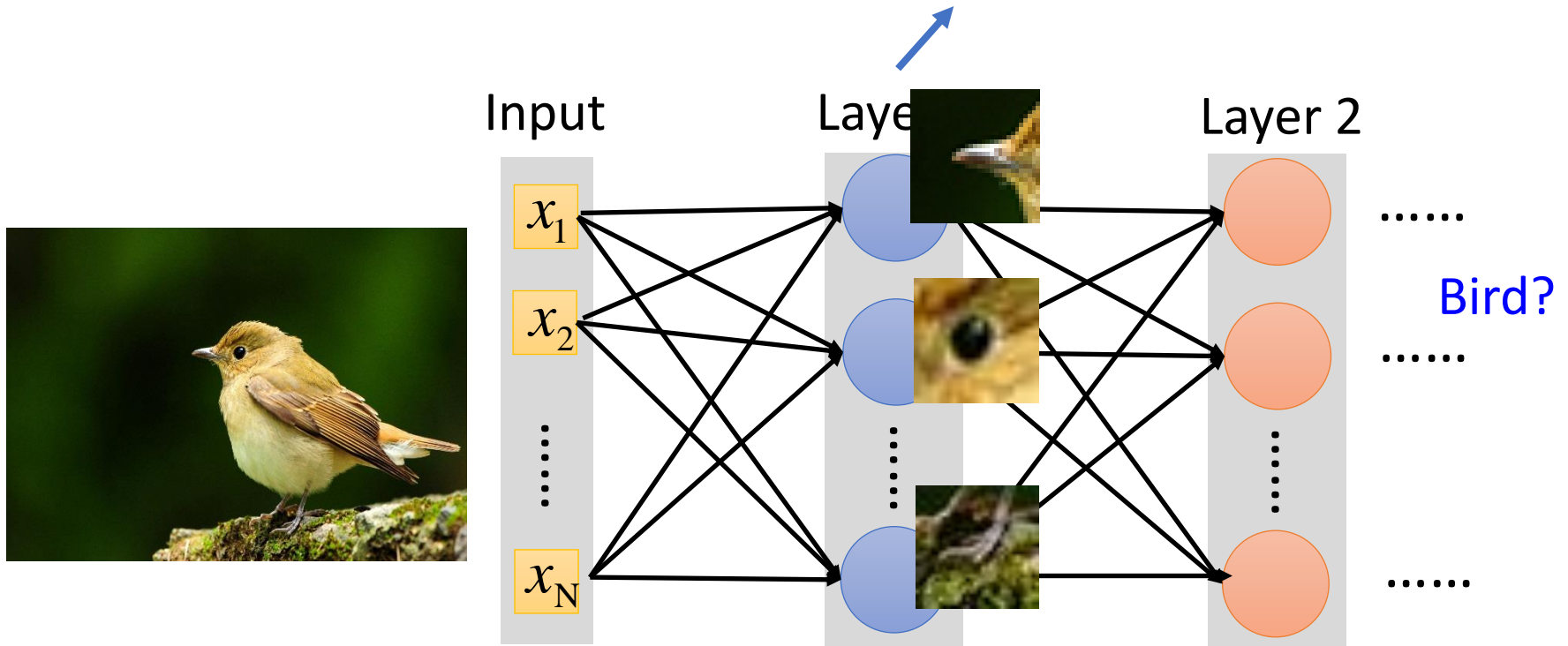$3 \times 10^7$

$x_k$

100 x 100

100 x 100 x 3

1000

······

······

······

Do we really need *"fully connected"*
in image processing?

4

# Observation 1

Identifying some critical patterns
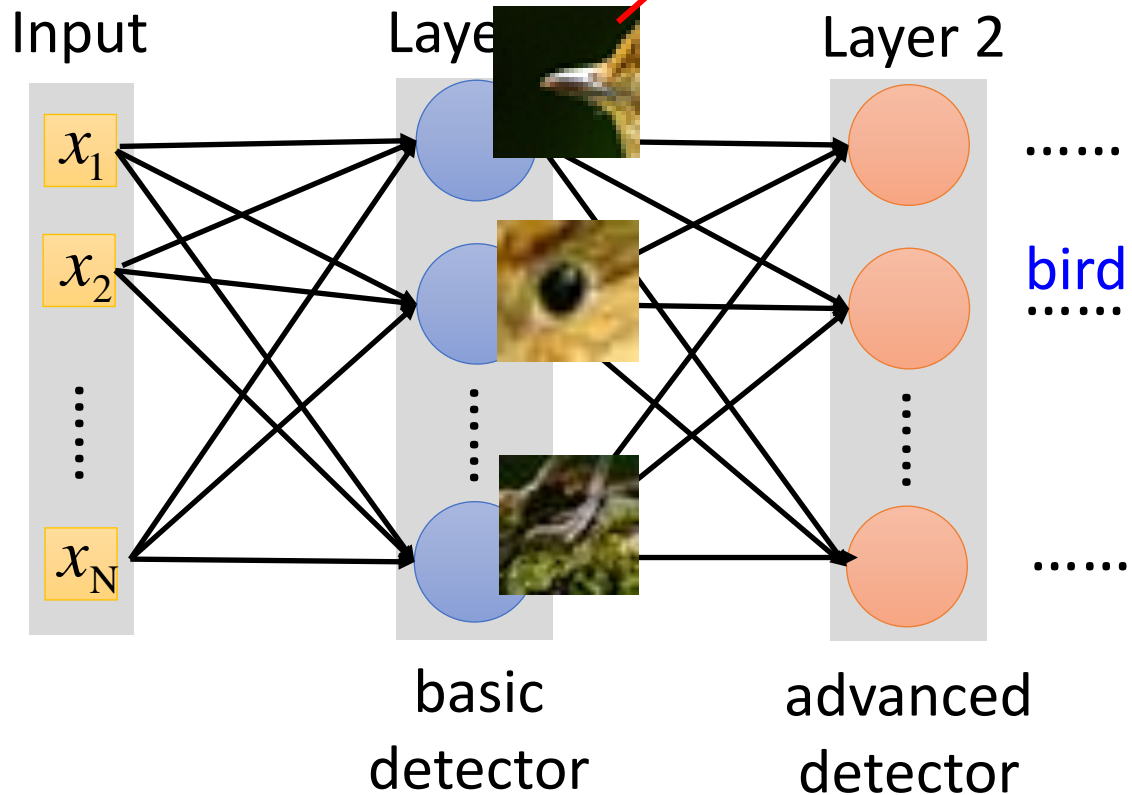


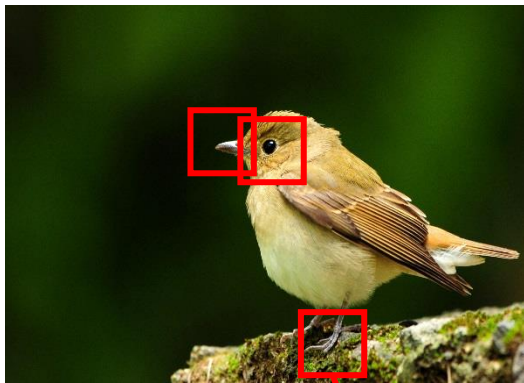Perhaps human also identify birds in a similar way ... ☺

https://www.dcard.tw/f/funny/p/233833012

# Observation 1

A neuron does not have to see the whole image.

Need to see the whole image?

Input

Layer

Layer 2

$x_1$

$x_2$

$x_N$

...... bird ......

basic detector

advanced detector

Some patterns are much smaller than the whole image.

# Simplification 1

Receptive field

3 x 3

3 x 3

3 x 3

3 x 3 x 3 weights

bias

1

1 0 0 0 0 1
0 1 0 0 1 0
0 0 1 1 0 0
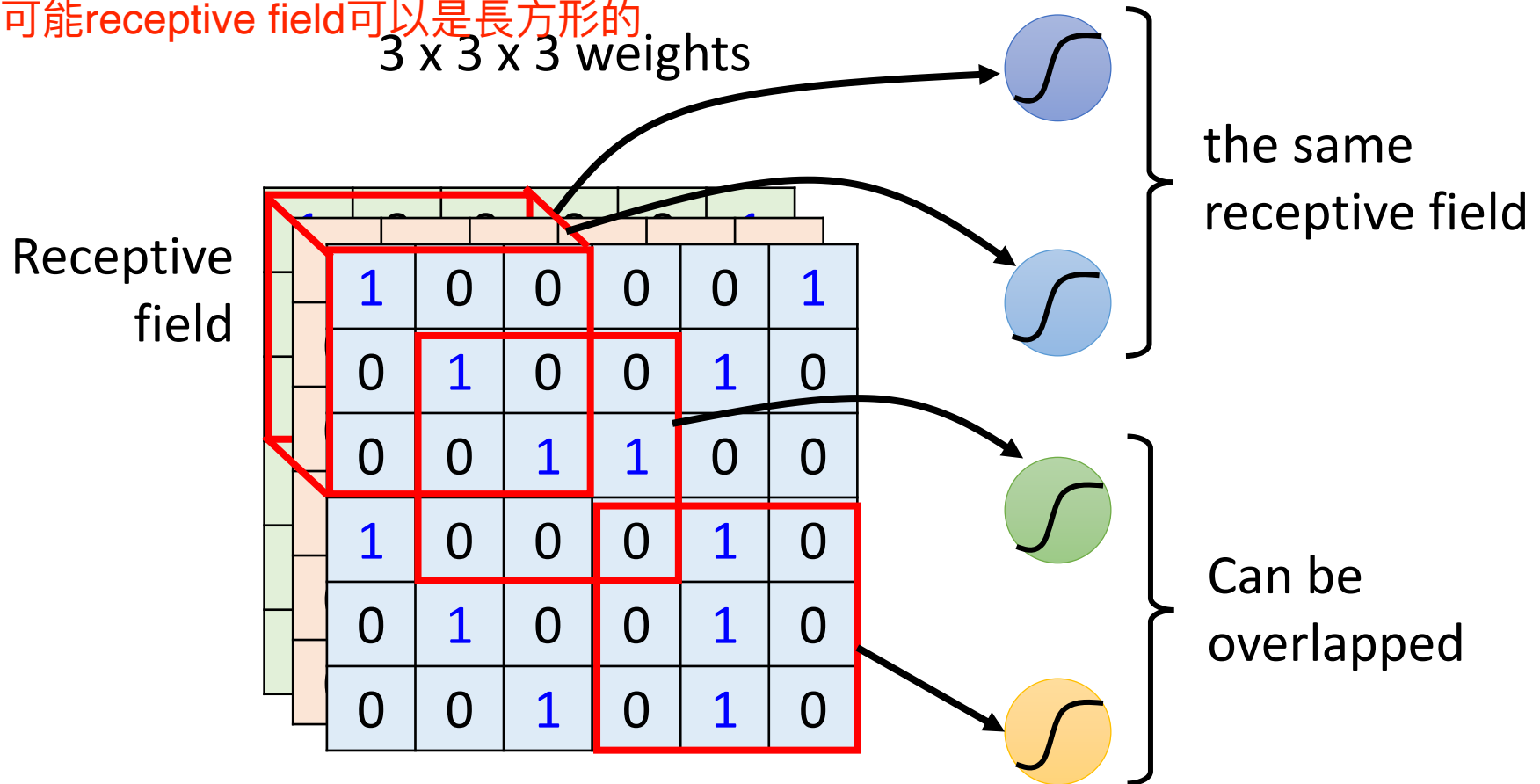1 0 0 0 1 0
0 1 0 0 1 0
0 0 1 0 1 0

# Simplification 1

架構的設計跟問題本身有關
有可能不一定要把rbg都計算
也有可能同一個receptive field可以接兩個以上的neuron
也有可能receptive field可以是長方形的
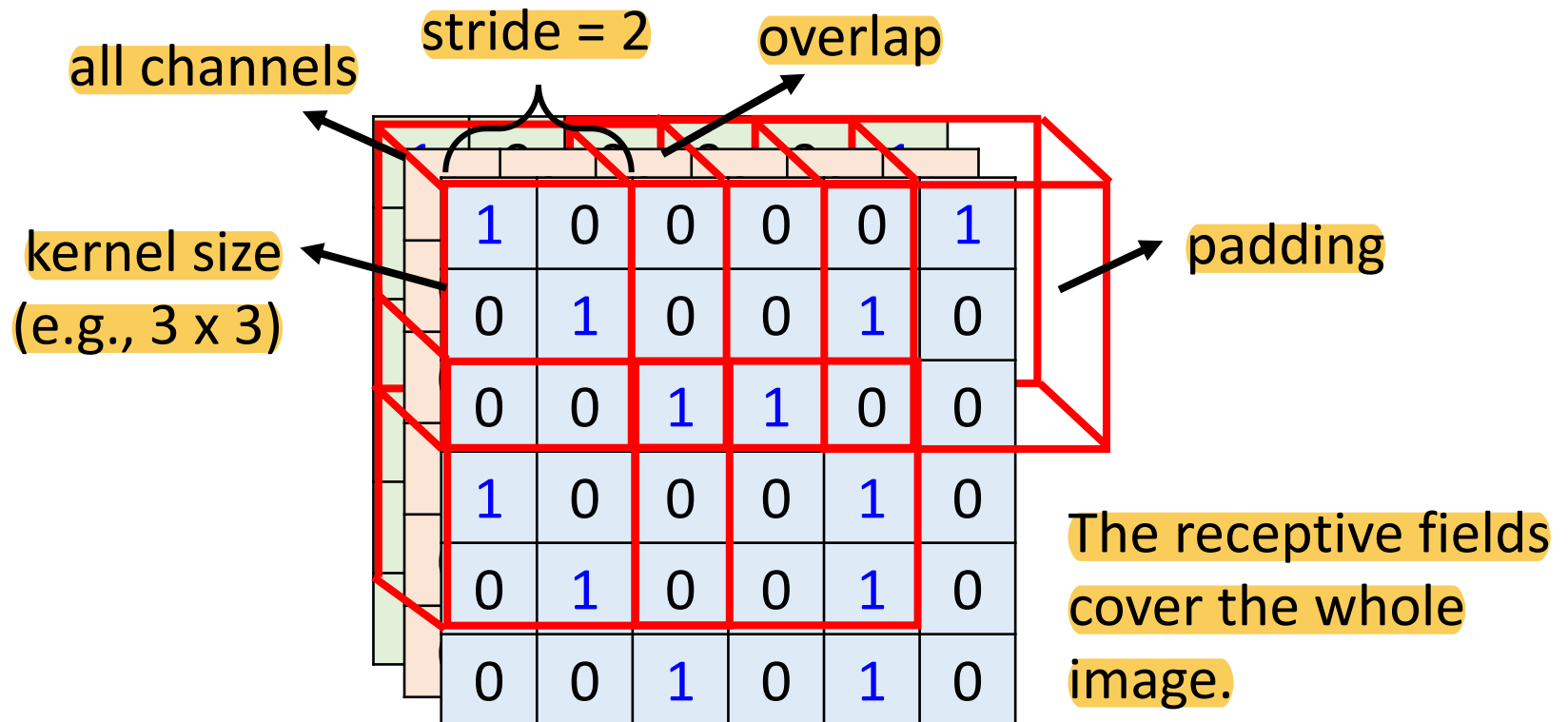
- Can different neurons have different sizes of receptive field?
- Cover only some channels?
- Not square receptive field?

3 x 3 x 3 weights

Receptive field

| 1 | 0 | 0 | 0 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

the same receptive field

Can be overlapped

# Simplification 1 – Typical Setting

經典的CNN架構

Each receptive field has a set of neurons (e.g., 64 neurons).



stride = 2

overlap

all channels

kernel size
(e.g., 3 x 3)

padding

The receptive fields cover the whole image.

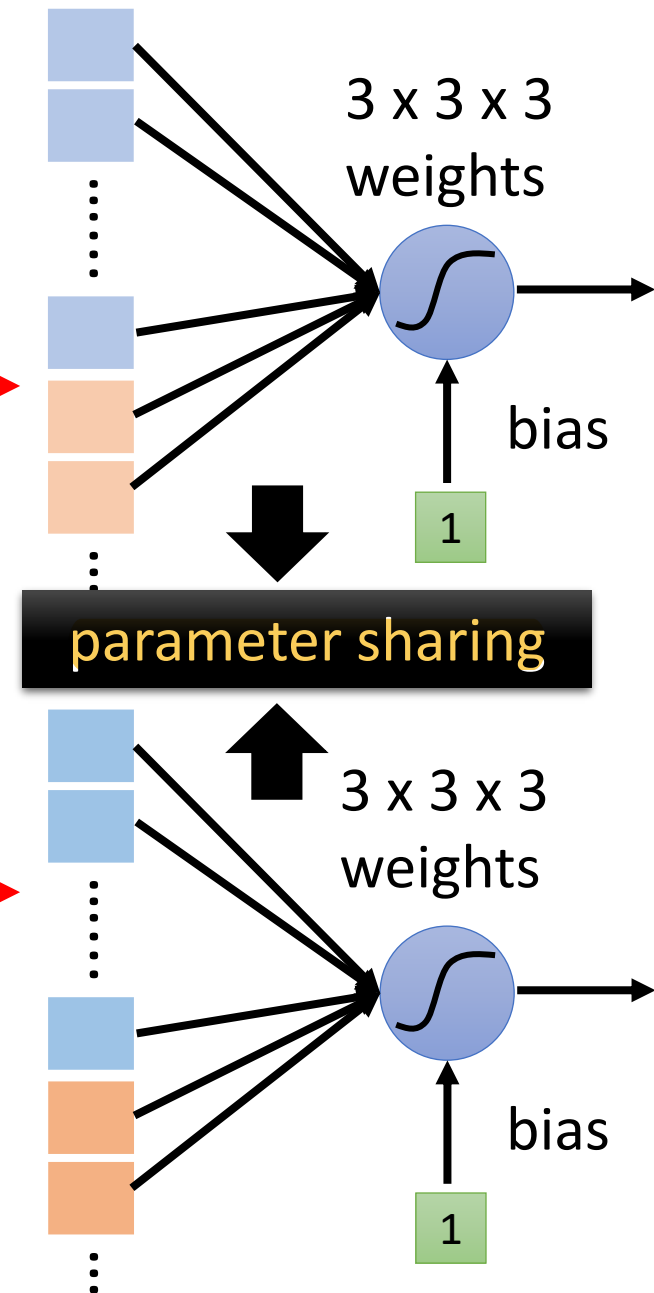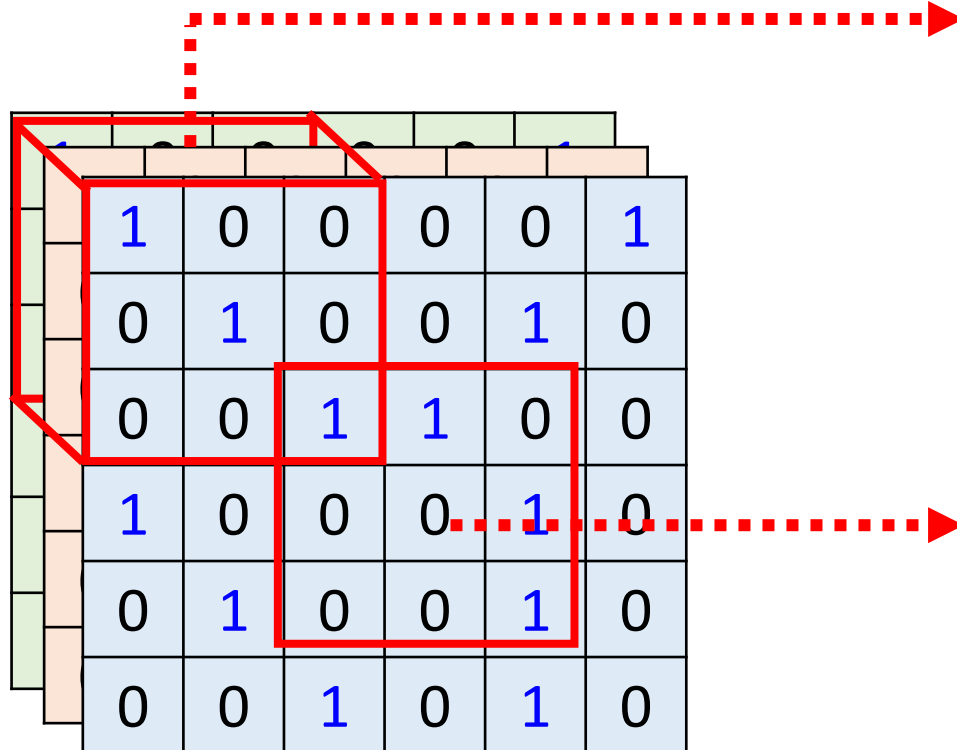# Observation 2  <span style="color:red">每一個receptive field在做的事情都很相近</span>

- The same patterns appear in different regions.



I detect "beak" in my receptive field.

Each receptive field needs a "beak" detector?

I detect "beak" in my receptive field.

# Simplification 2

3 x 3 x 3 weights

bias

parameter sharing

3 x 3 x 3 weights

bias

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

1

1

12

# Simplification 2

$$\sigma(w_1 x_1 + w_2 x_2 + \cdots)$$

$x_1$

$x_2$

$w_1$

$w_2$

bias

1

兩個neuron的weight完全一樣

$$\sigma(w_1 x_1' + w_2 x_2' + \cdots)$$

$x_1'$

$x_2'$

$w_1$

$w_2$

bias

1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

Two neurons with the same receptive field would not share parameters.
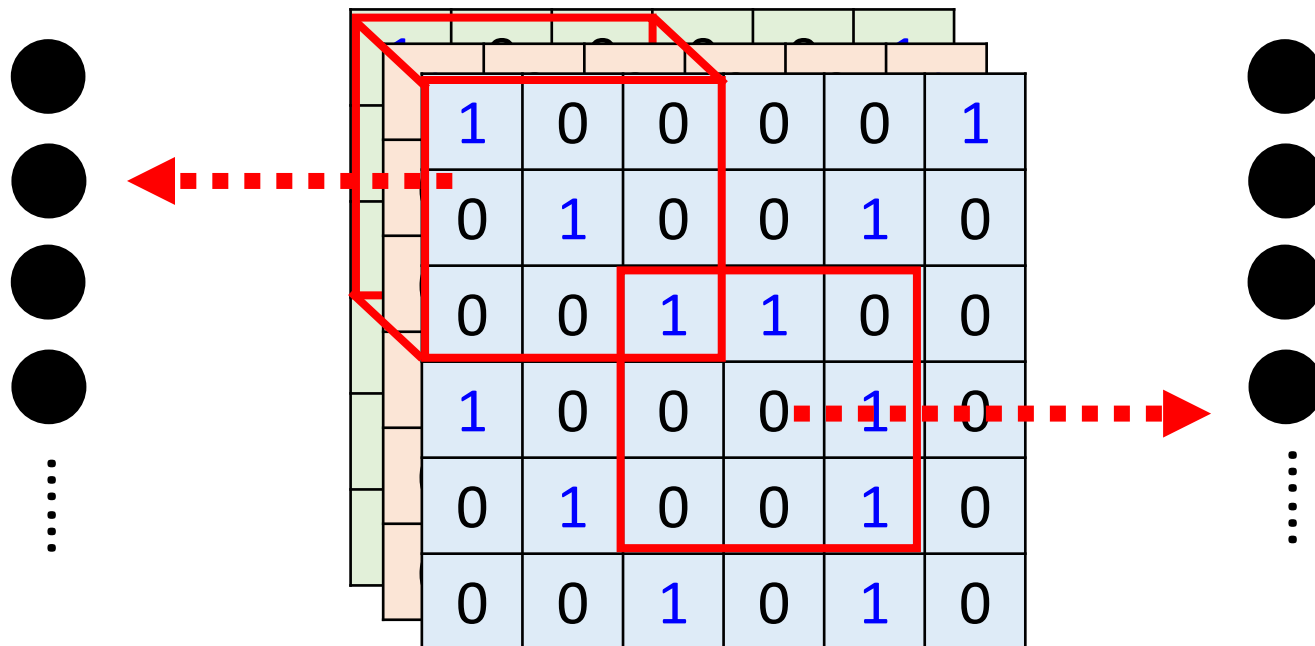
因為這樣輸出就都會是一樣的，沒有必要

13

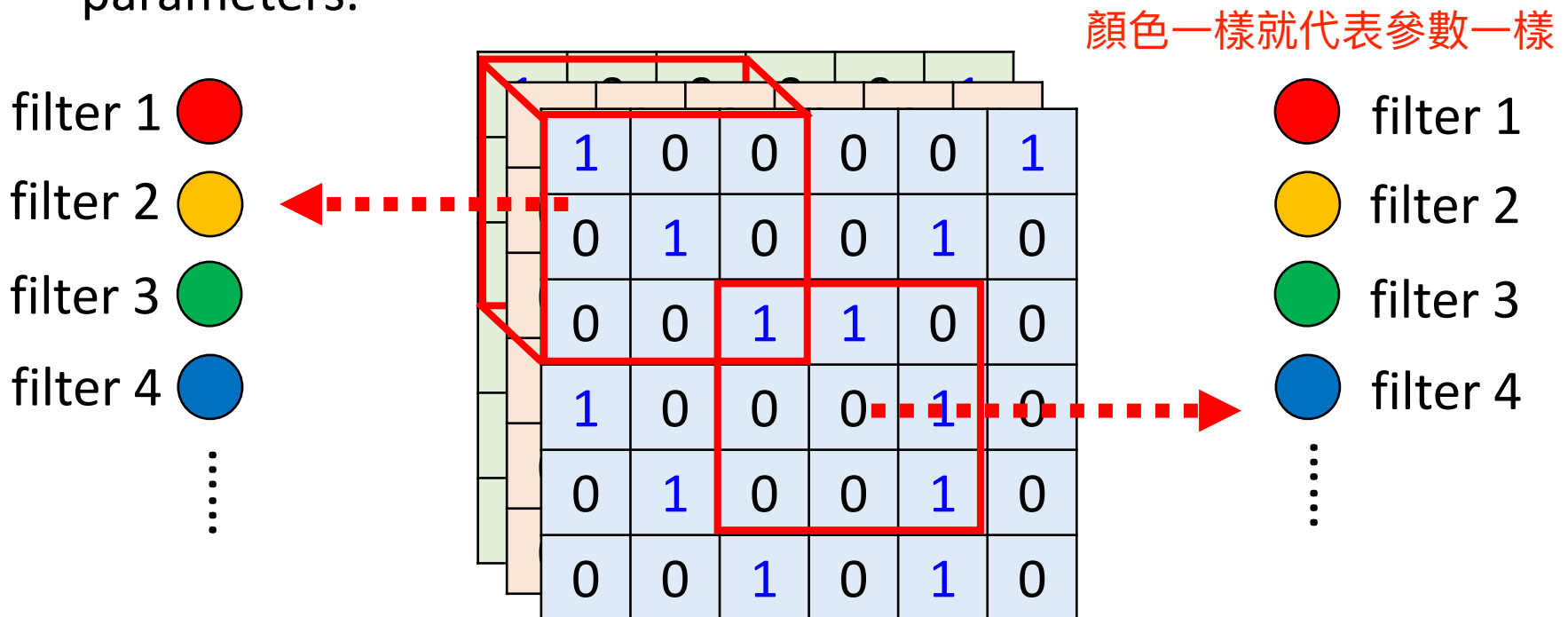# Simplification 2 – Typical Setting

Each receptive field has a set of neurons (e.g., 64 neurons).
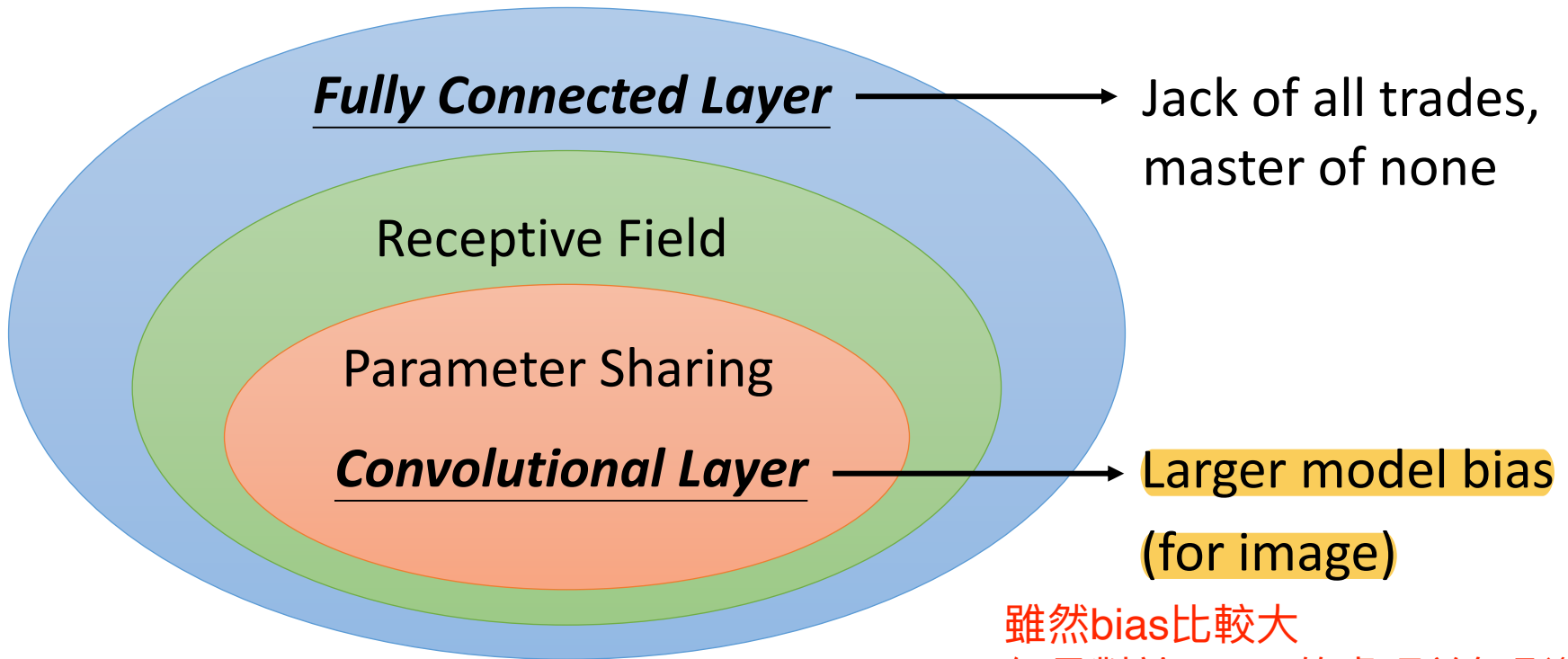
# Simplification 2 – Typical Setting

Each receptive field has a set of neurons (e.g., 64 neurons).

Each receptive field has the neurons with the same set of parameters.

# Benefit of Convolutional Layer

三種function set是包含關係

**_Fully Connected Layer_** ⟶ Jack of all trades, master of none

Receptive Field

Parameter Sharing

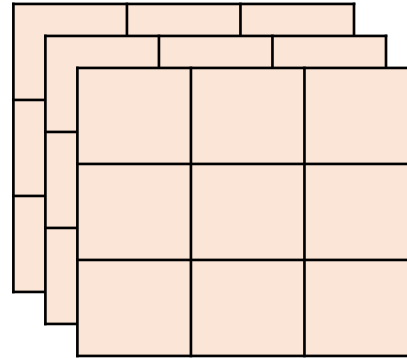**_Convolutional Layer_** ⟶ Larger model bias (for image)

雖然bias比較大
但是對於image的處理並無影響

- Some patterns are much smaller than the whole image.
- The same patterns appear in different regions.

16

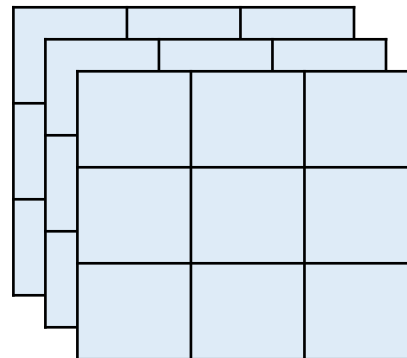# Convolutional Layer

Convolution

channel = 3  (colorful)

channel = 1  (black and white)

Filter 1
3 x 3 x channel
tensor

Filter 2
3 x 3 x channel
tensor

Each filter detects a small
pattern (3 x 3 x channel).

17

# Convolutional Layer

Consider channel = 1
(black and white image)

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

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

Filter 1

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

Filter 2

⋮

(The values in the filters
are unknown parameters.)

# Convolutional Layer

Filter 1

偵測左上到右下的斜線

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

stride=1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

| 3 | -1 | -3 | -1 |
|---|----|----|----|
| -3 | 1 | 0 | -3 |
| -3 | -3 | 0 | 1 |
| 3 | -2 | -2 | -1 |

19

# Convolutional Layer

**Filter 2**

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

偵測垂直線

stride=1

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

Do the same process for every filter

| -1 | -1 | -1 | -1 |
|----|----|----|----|
| -1 |    |    | 1 |
| -1 | -1 | -2 | 1 |
| -1 | 0 | -4 | 3 |

**Feature Map**

所有filter偵測出來的值
就稱為feature map

# *Convolutional Layer*



"Image" with 64 channels

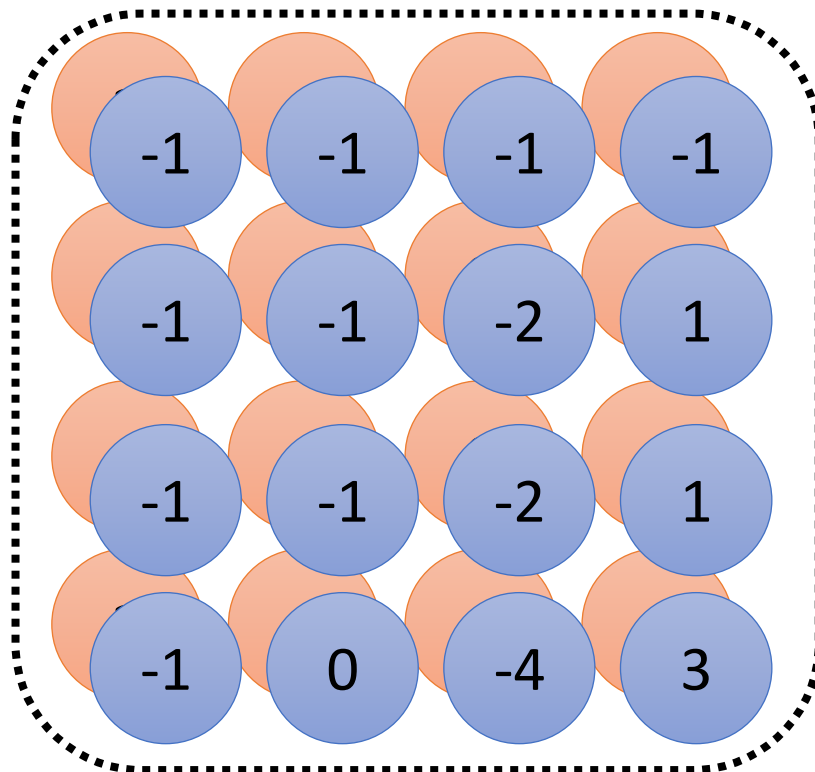64 filters

Convolution

Convolution
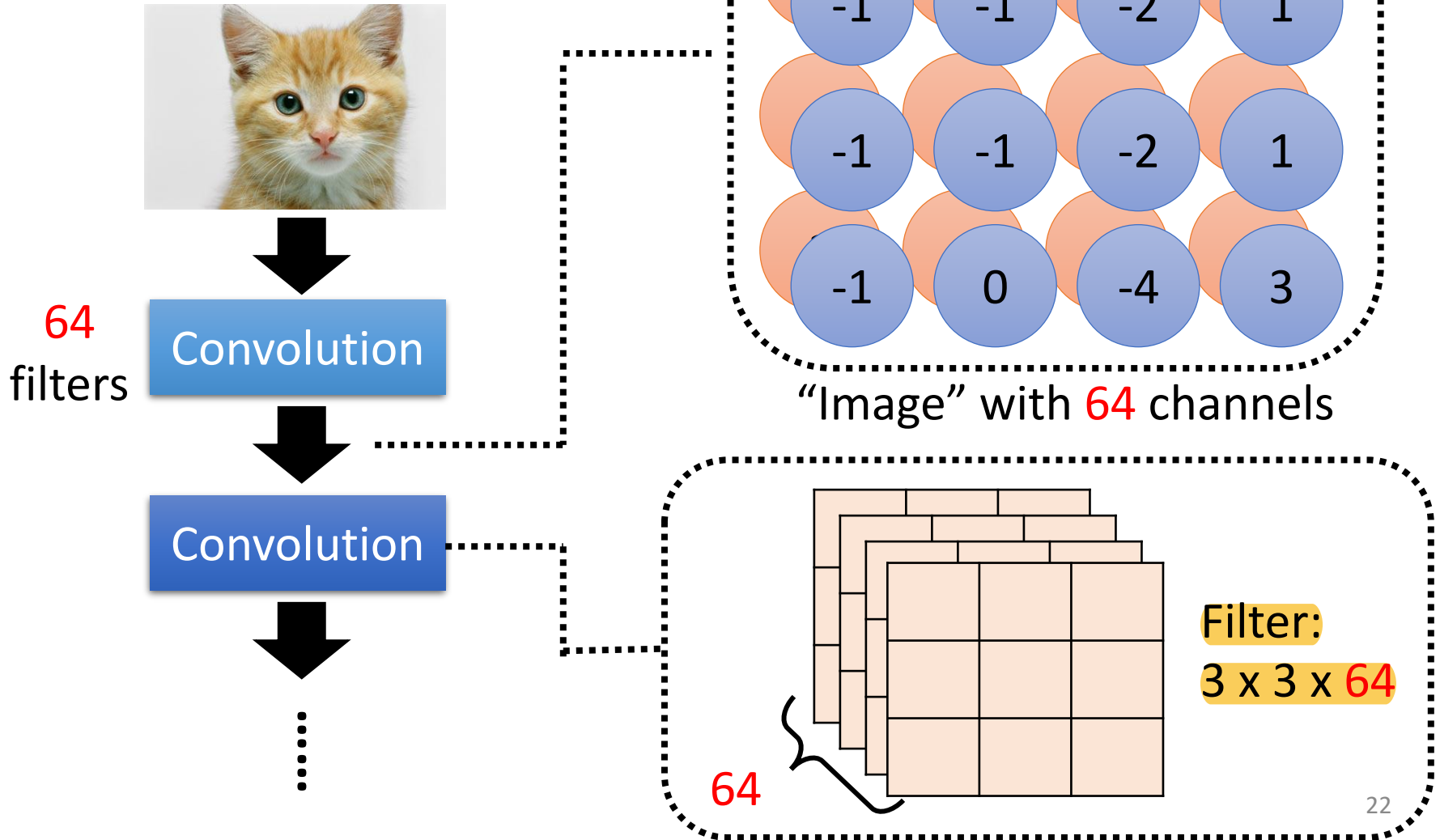
可以看做是一張新的圖片
而且有64個channels
 （如果有64個filters）
 （原本的圖片為3個channel：red, green, blue）

而且圖片會變小，接著又可以繼續再做convolution

# *Multiple Convolutional Layers*

64 filters

[ Convolution ]

[ Convolution ]



|  -1 | -1 | -1 | -1 |
| -1 | -1 | -2 |  1 |
| -1 | -1 | -2 |  1 |
| -1 |  0 | -4 |  3 |

"Image" with 64 channels

Filter:
3 x 3 x 64

64

22

# *Multiple Convolutional Layers*

若第二層filter僅有3x3
但是對應於原本的圖片
已經涵蓋5x5了

64 filters

Convolution

Convolution

| 1 | 0 | 0 | 0 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

| -1 | -1 | -1 | -1 |
| -1 | -1 | -2 | 1 |
| -1 | -1 | -2 | 1 |
| -1 | 0 | -4 | 3 |

23

# Comparison of Two Stories



Receptive field

Filter
3 x 3 x channel tensor

(ignore bias in this slide)

The neurons with different receptive fields **share the parameters**.

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

**Each filter convolves over the input image.**

bias

1

bias

1

# Convolutional Layer

| ***Neuron Version Story*** | ***Filter Version Story*** |
|---|---|
| Each neuron only considers a receptive field. | There are a set of filters detecting small patterns. |
| The neurons with different receptive fields share the parameters. | Each filter convolves over the input image. |

They are the same story.

# Observation 3

- Subsampling the pixels will not change the object

例如將原先圖片的偶數列與偶數行拿掉
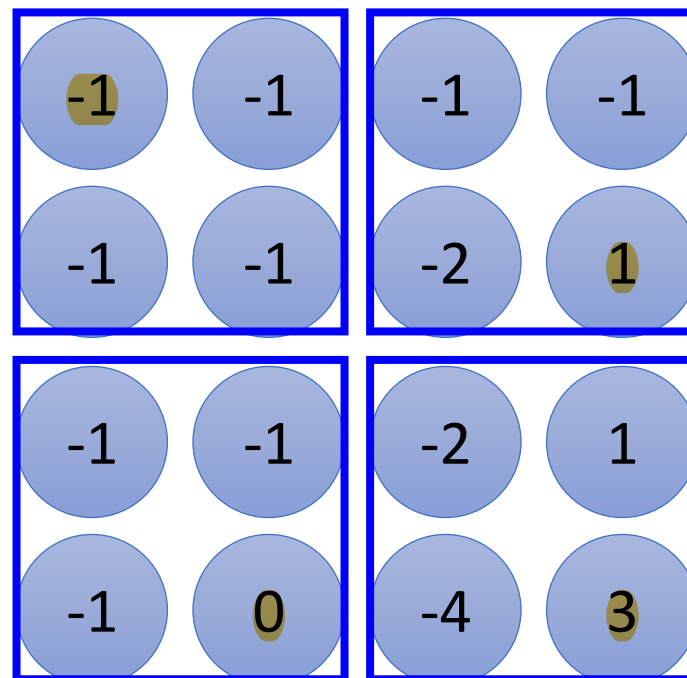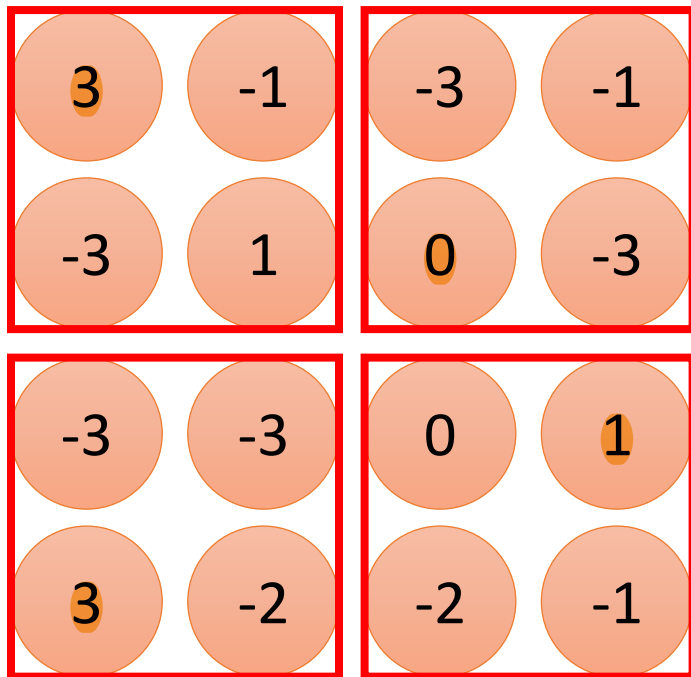看起來還是會很像

bird



bird

subsampling

# Pooling – Max Pooling

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

Filter 1

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

Filter 2

max pooling，
就是在kernel size上取最大的當代表

| | |
|---|---|
| 3 | -1 |
| -3 | 1 |

| | |
|---|---|
| -3 | -1 |
| 0 | -3 |

| | |
|---|---|
| -3 | -3 |
| 3 | -2 |

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

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

| | |
|---|---|
| -1 | -1 |
| -2 | 1 |

| | |
|---|---|
| -1 | -1 |
| -1 | 0 |

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

# Convolutional Layers + Pooling

通常架構就是convolution和pooling在交替使用



Repeat

**Convolution**

**Pooling**

pooling存在的理由是為了降低運算
但可能對資料有所傷害
因次近年來，運算能力變強
越來越多人把pooling拿掉

| -1 | -1 | -1 | -1 |
|----|----|----|----|
| -1 | -1 | -2 | 1 |
| -1 | -1 | -2 | 1 |
| -1 | 0 | -4 | 3 |

"Image" with 64 channels

| -1 | 1 |
|----|----|
| 0 | 3 |

# The whole CNN

最後把圖片拉直再丟進去一個fully connected layers

cat dog ......

softmax

**Fully Connected Layers**



Convolution

Pooling

Convolution

Pooling

Flatten

# Application: Playing Go



**19 x 19 matrix (image)**

**Network**

Next move (19 x 19 positions)

19 x 19 classes

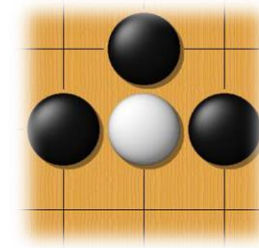48 channels in Alpha Go

Black: 1

white: -1

none: 0

Fully-connected network can be used
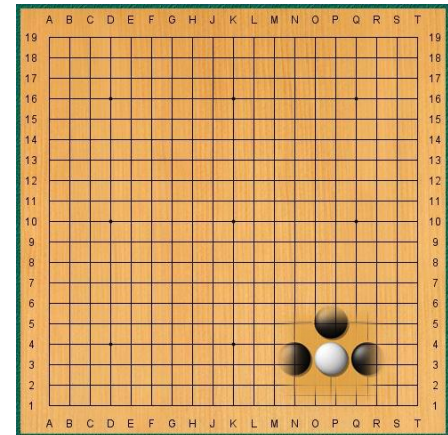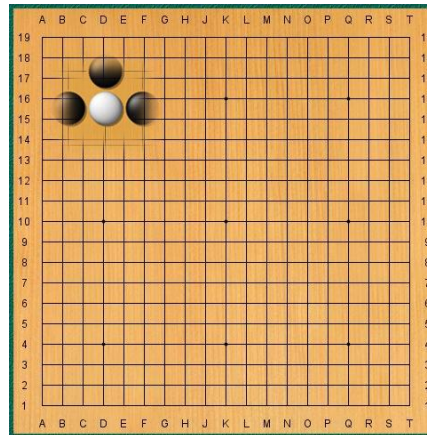
But CNN performs much better.

# Why CNN for Go playing?

- Some patterns are much smaller than the whole image

  Alpha Go uses 5 x 5 for first layer

- The same patterns appear in different regions.

# Why CNN for Go playing?

- Subsampling the pixels will not change the object

**Pooling**　　How to explain this???

**Neural network architecture.** The input to the policy network is a $19 \times 19 \times 48$ image stack consisting of 48 feature planes. The first hidden layer zero pads the input into a $23 \times 23$ image, then convolves $k$ filters of kernel size $5 \times 5$ with stride 1 with the input image and applies a rectifier nonlinearity. Each of the subsequent hidden layers 2 to 12 zero pads the respective previous hidden layer into a $21 \times 21$ image, then convolves $k$ filters of kernel size $3 \times 3$ with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size $1 \times 1$ with stride 1, with a different bias for each position, and applies a softmax function. The match version of AlphaGo used $k = 192$ filters; Fig. 2b and Extended Data Tabl_____256 and 384 filters

Alpha Go does not use Pooling ......

# *More Applications*

每一種cnn的應用
都會針對問題去設計架構
並不是圖片的cnn就能套用在任意的情境上

## Speech

https://dl.acm.org/doi/10.1109/TASLP.2014.2339736

## Natural Language Processing

https://www.aclweb.org/anthology/S15-2079/
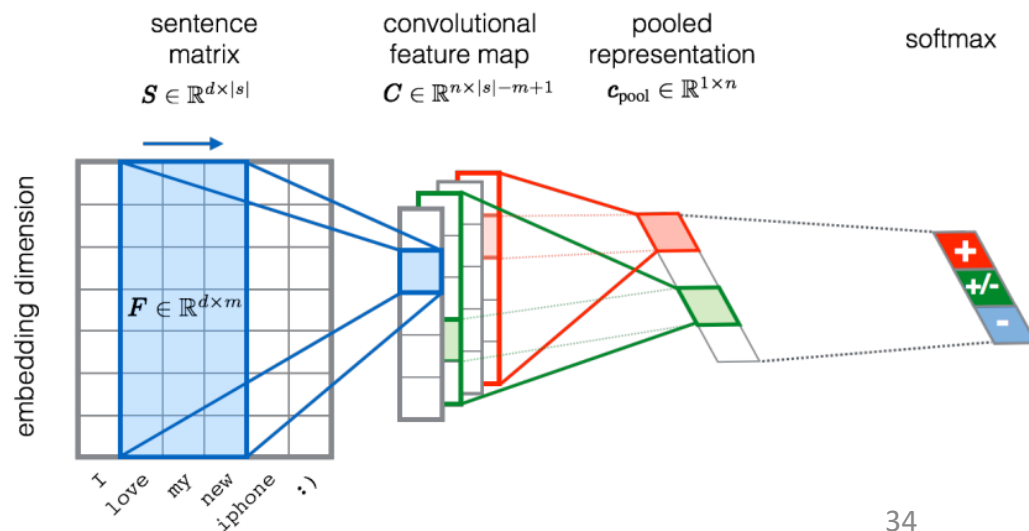
cnn辨識不出來這兩張圖片
雖然data augmentation可以處理這種問題

# To learn more ...

但是data augementation也不是所有放大或旋轉的角度都有包含
因此若能讓cnn自己學會scaling和rotation是最好的（就是spatial transformer layer）

- CNN is not invariant to scaling and rotation (we need data augmentation ☺).



*Spatial Transformer Layer*



https://youtu.be/SoCywZ1hZak
(in Mandarin)