# Convolutional Neural Network (CNN)

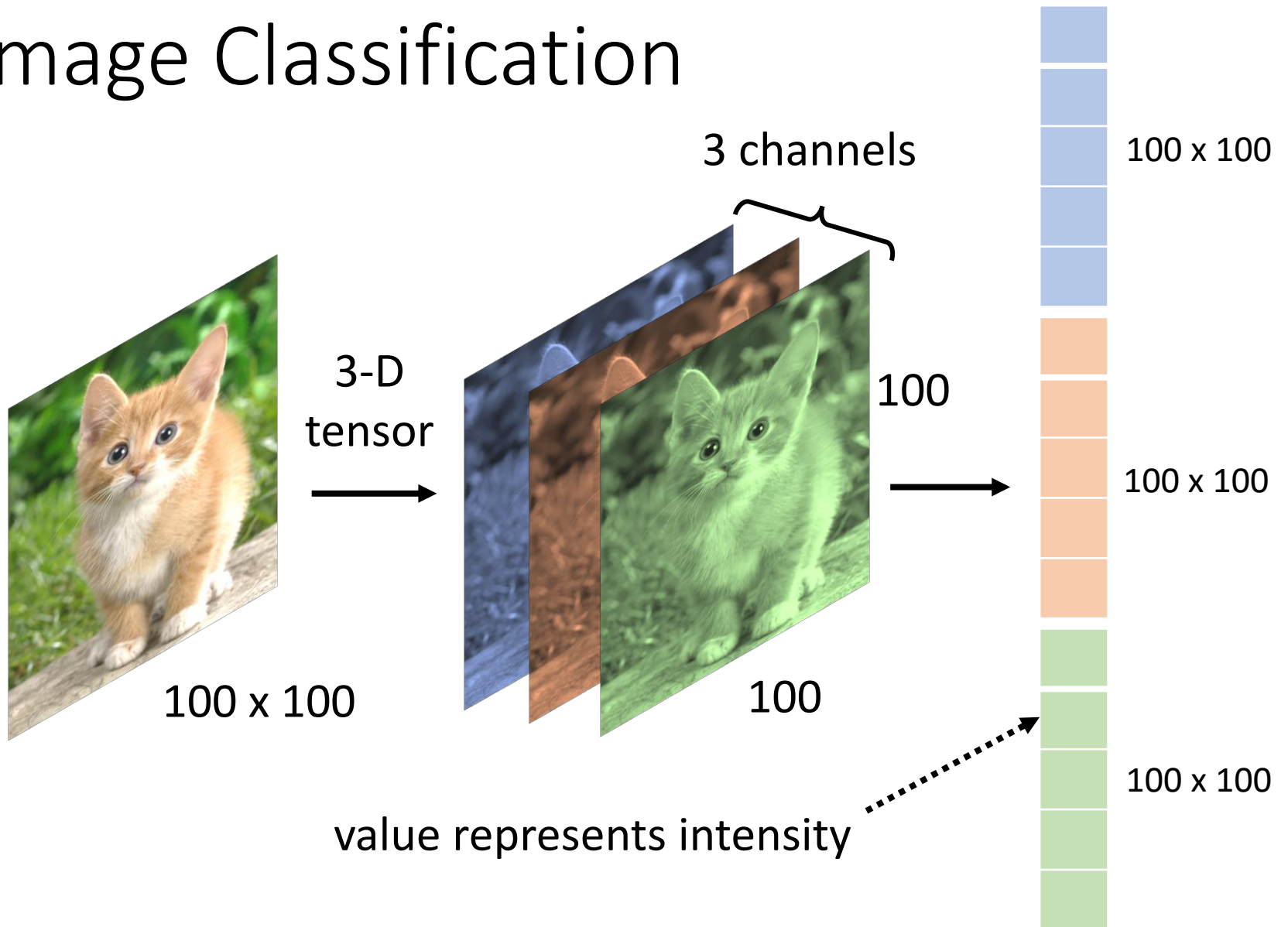Network Architecture designed for Image

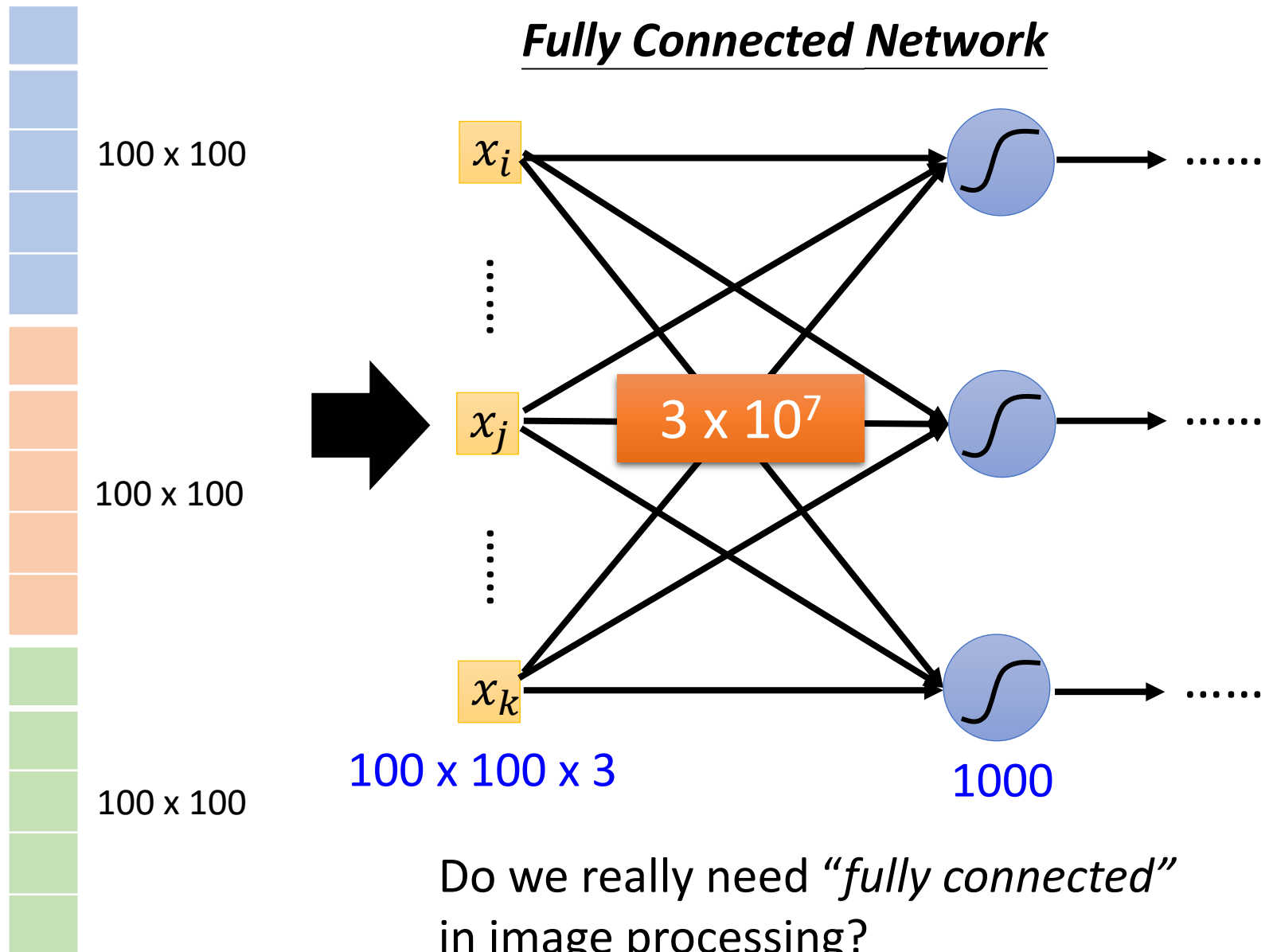# Image Classification



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

Model

$y'$ ⟷ $\hat{y}$

Cross entropy

100 x 100

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

# Image Classification



3-D tensor

3 channels

100

100 x 100

100 x 100

100

100 x 100

100 x 100

100 x 100

value represents intensity

# *Fully Connected Network*

100 x 100

100 x 100

100 x 100

$x_i$

$x_j$

$x_k$

3 x $10^7$

100 x 100 x 3

1000

......

......

......

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

# Observation 1

Identifying some critical patterns



Perhaps human also identify birds in a similar way … ☺

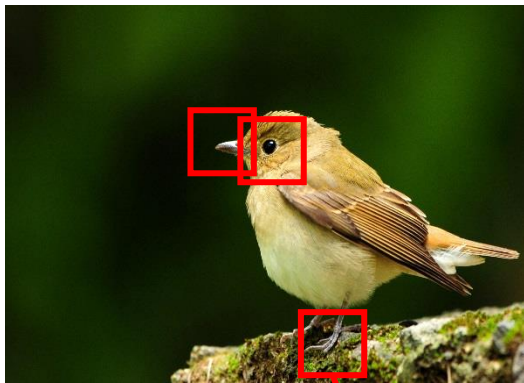https://www.dcard.tw/f/funny/p/233833012

# Observation 1

A neuron does not have to see the whole image.
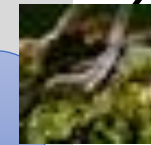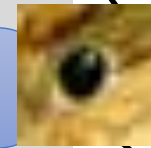
Need to see the whole image?



Input

Layer

Layer 2

$x_1$

$x_2$

$x_N$

...... 

bird

......

basic detector

advanced detector
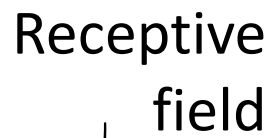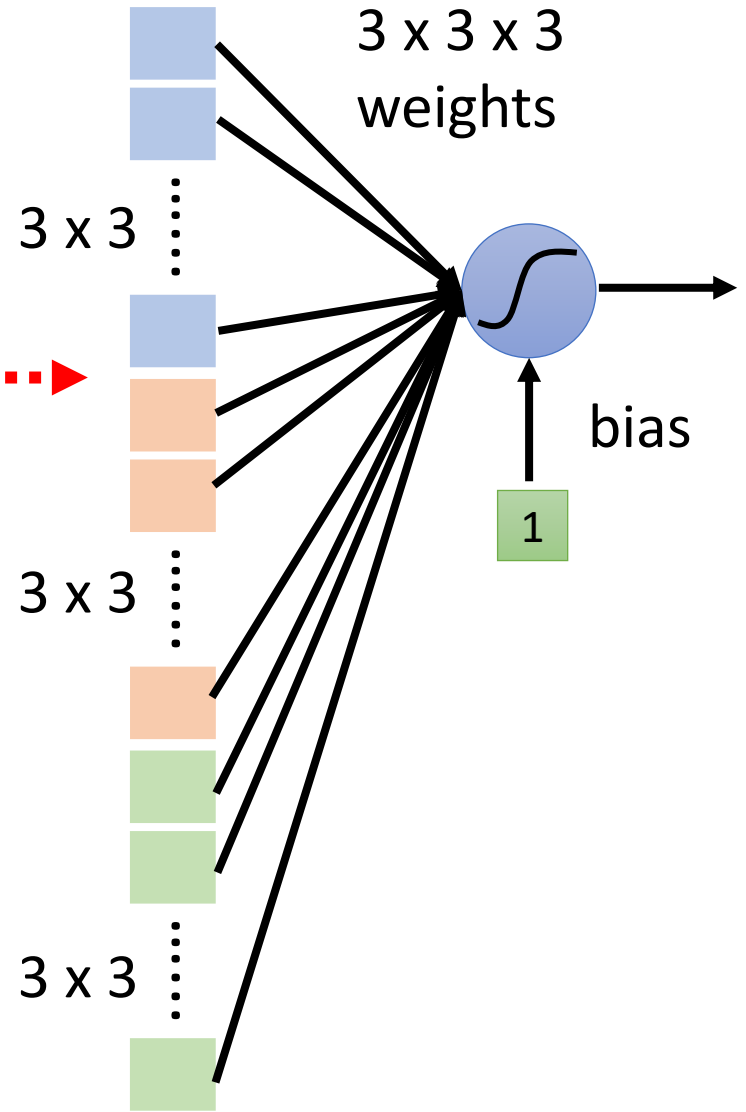
Some patterns are much smaller than the whole image.

# Simplification 1



Receptive field

感受野

3 x 3 x 3 weights

3 x 3

3 x 3

3 x 3

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 |

8

# Simplification 1

- 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

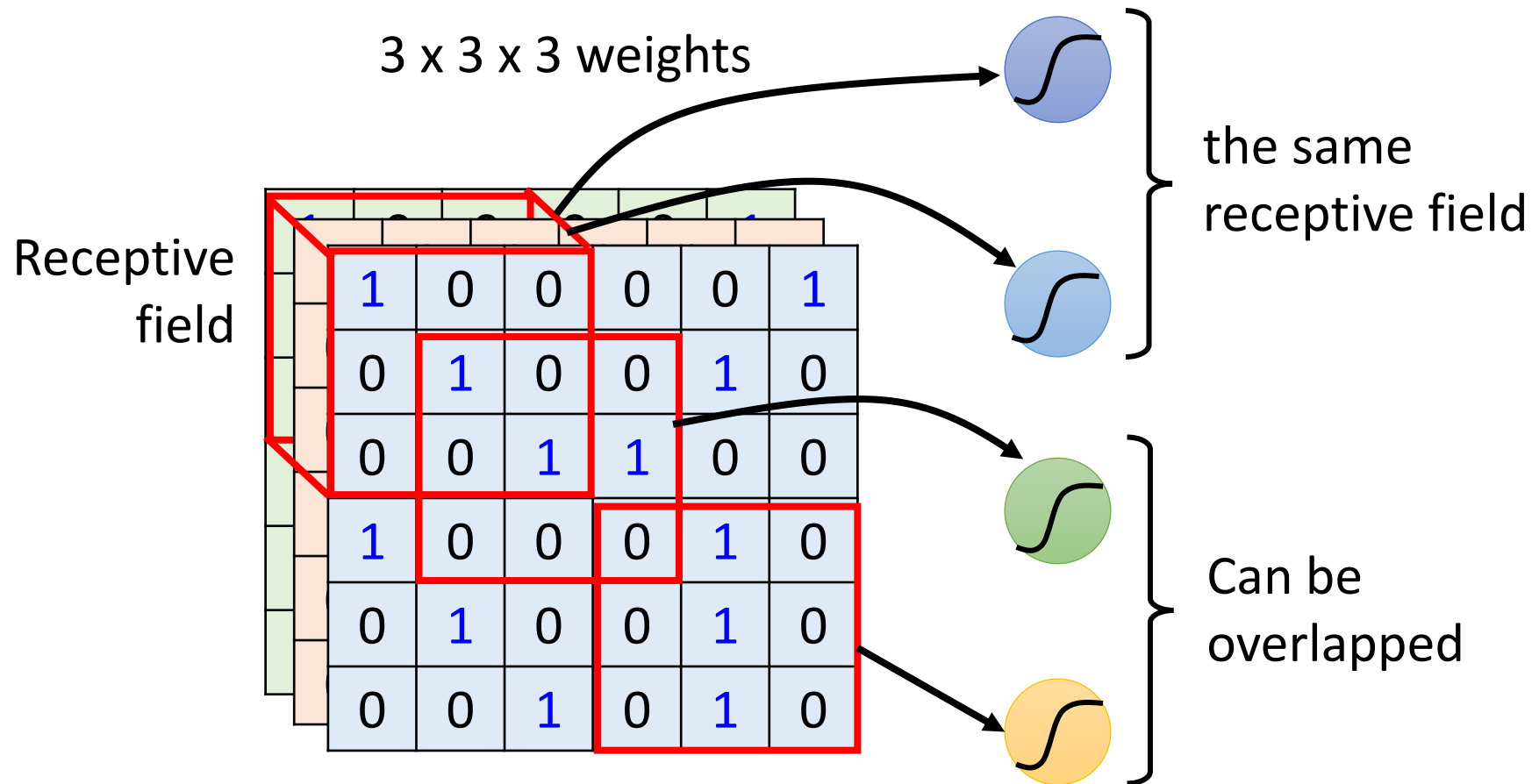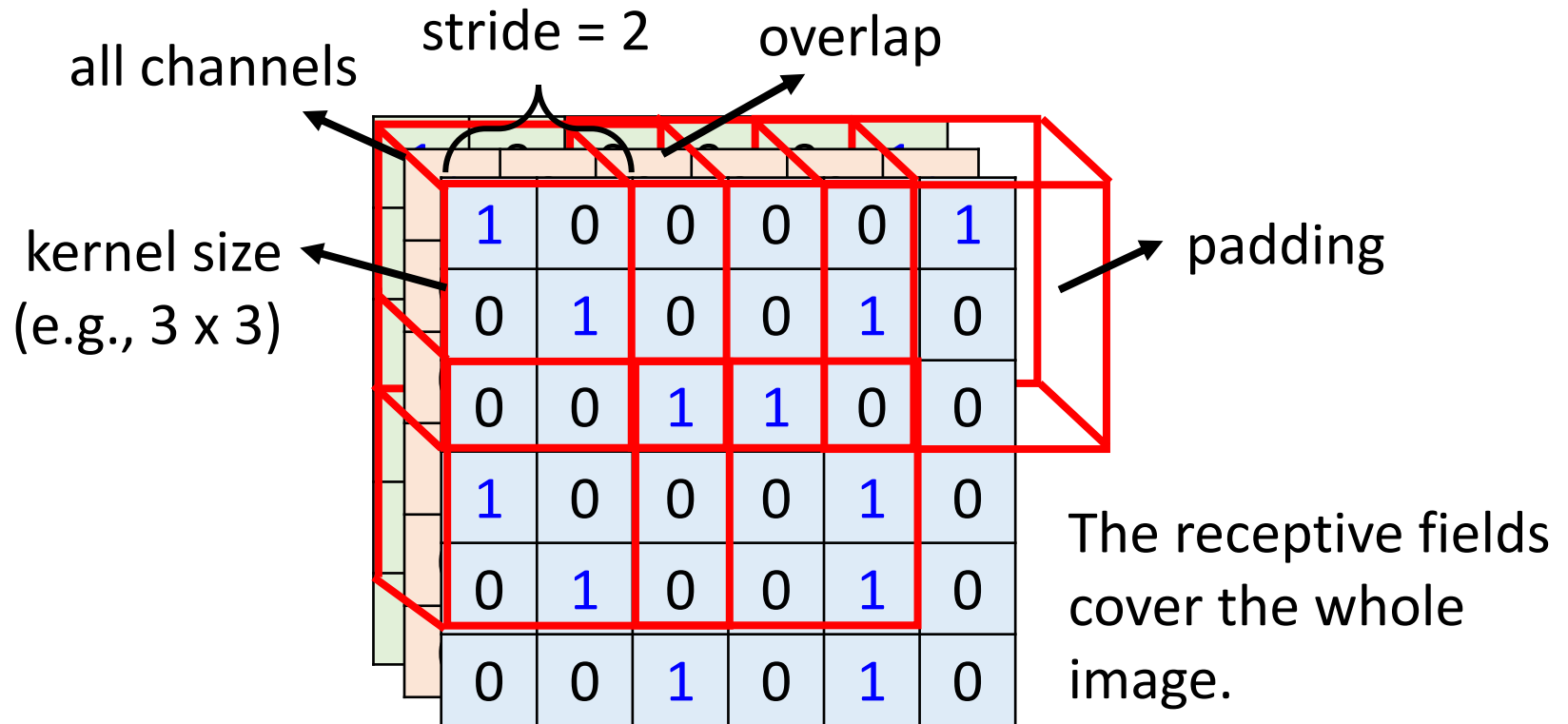# Simplification 1 – Typical Setting

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

| 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 receptive fields cover the whole image.

# Observation 2

- The same patterns appear in different regions.

# Simplification 2

| 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 |

3 x 3 x 3 weights

bias

1

**parameter sharing**

3 x 3 x 3 weights

bias

1

# Simplification 2



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

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

Two neurons with the same receptive field would not share parameters.

13

# Simplification 2 – Typical Setting

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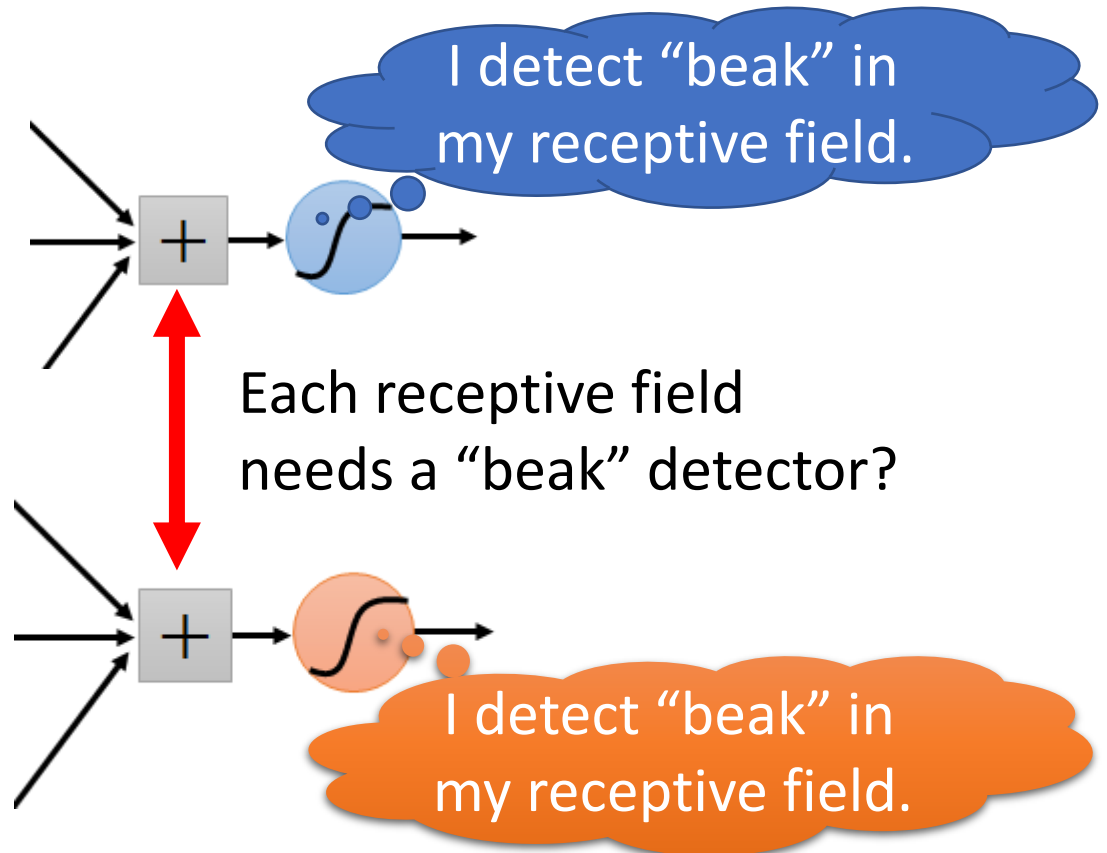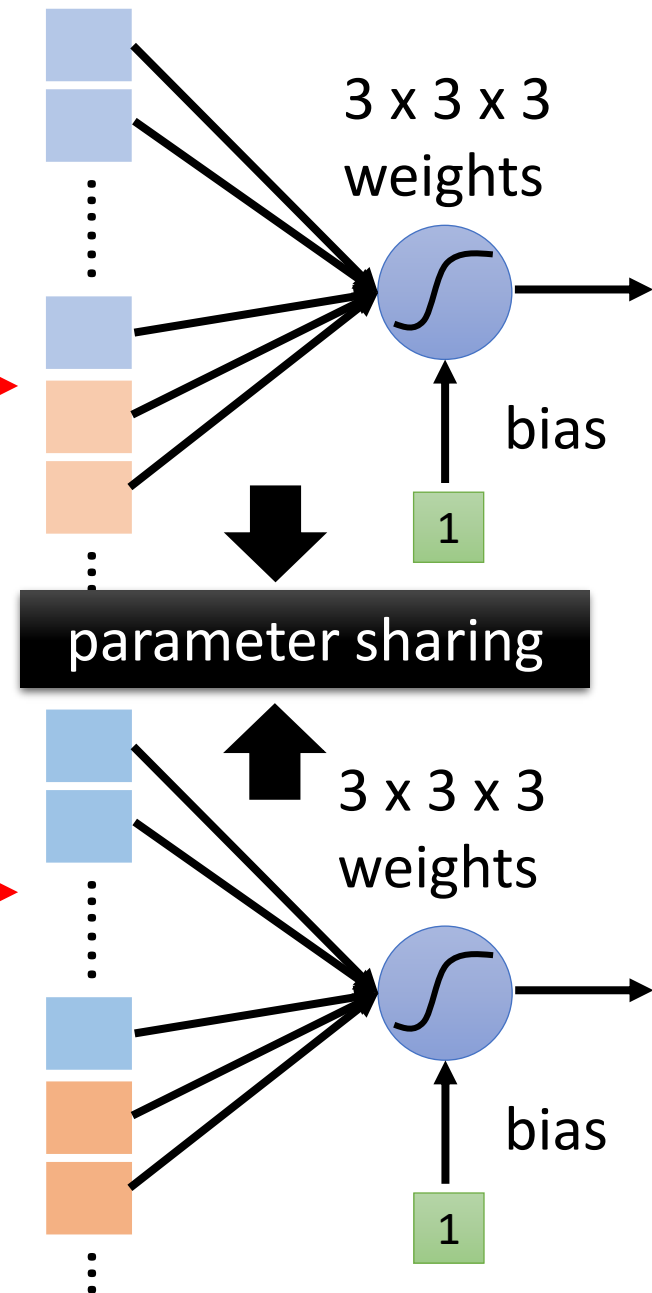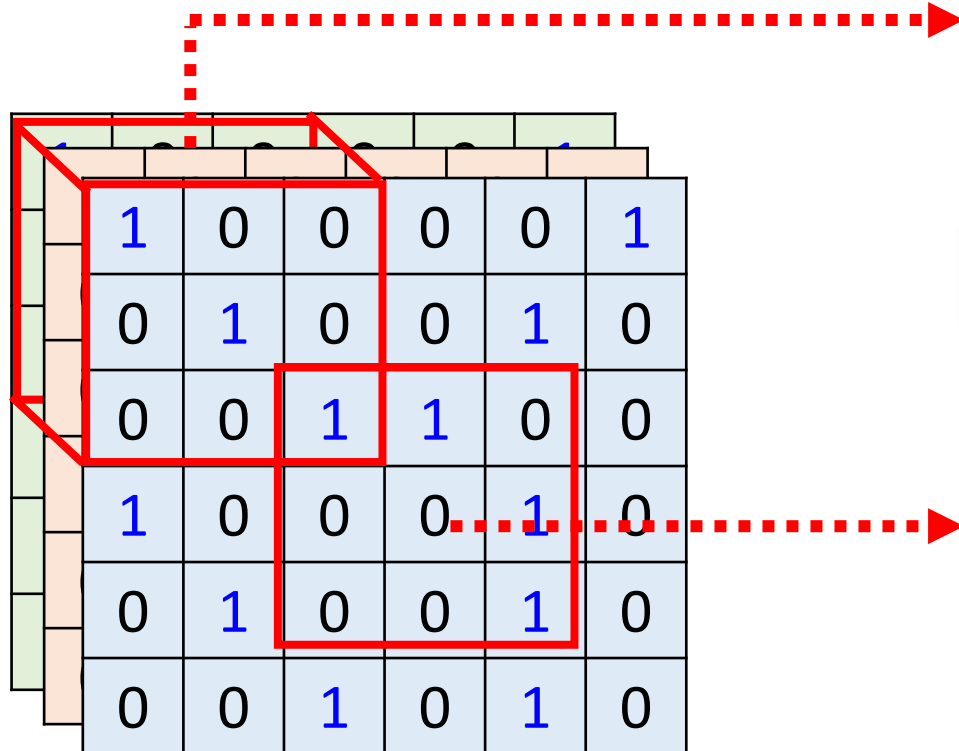
# Simplification 2 – Typical Setting

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

Fully Connected Layer ——→ Jack of all trades, master of none

Receptive Field

Parameter Sharing

Convolutional Layer ——→ Larger model bias (for image)

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

# Convolutional Layer
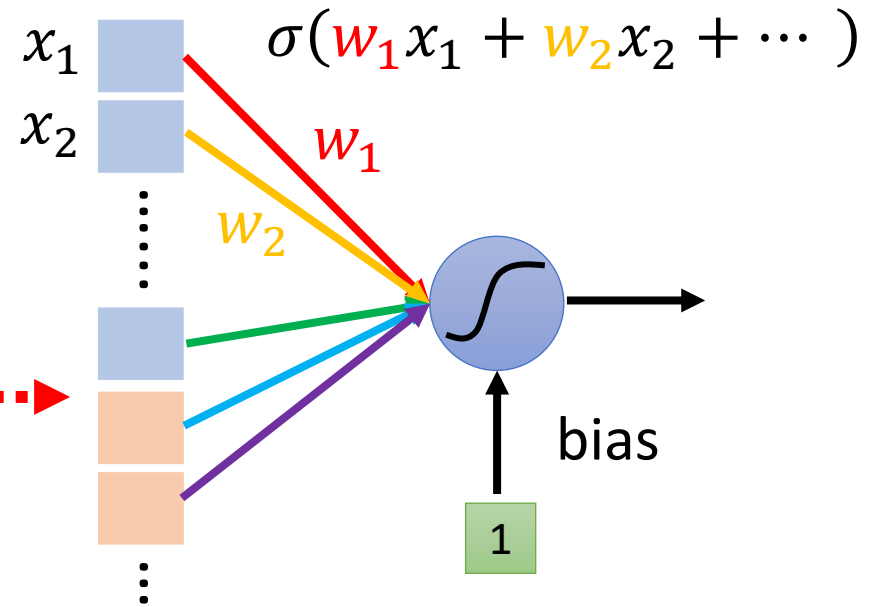


Convolution

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.)

18

# Convolutional Layer

Filter 1

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

stride=1

6 x 6 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 |

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

19

# Convolutional Layer

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

Filter 2

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

# *Convolutional Layer*



64 filters

Convolution

Convolution

"Image" with 64 channels

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

# *Multiple Convolutional Layers*



64 filters

Convolution

Convolution

"Image" with 64 channels

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

Filter:
3 x 3 x 64

64

# Multiple Convolutional Layers



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 |

# 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 |

bias

1

bias

1

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

# Convolutional Layer

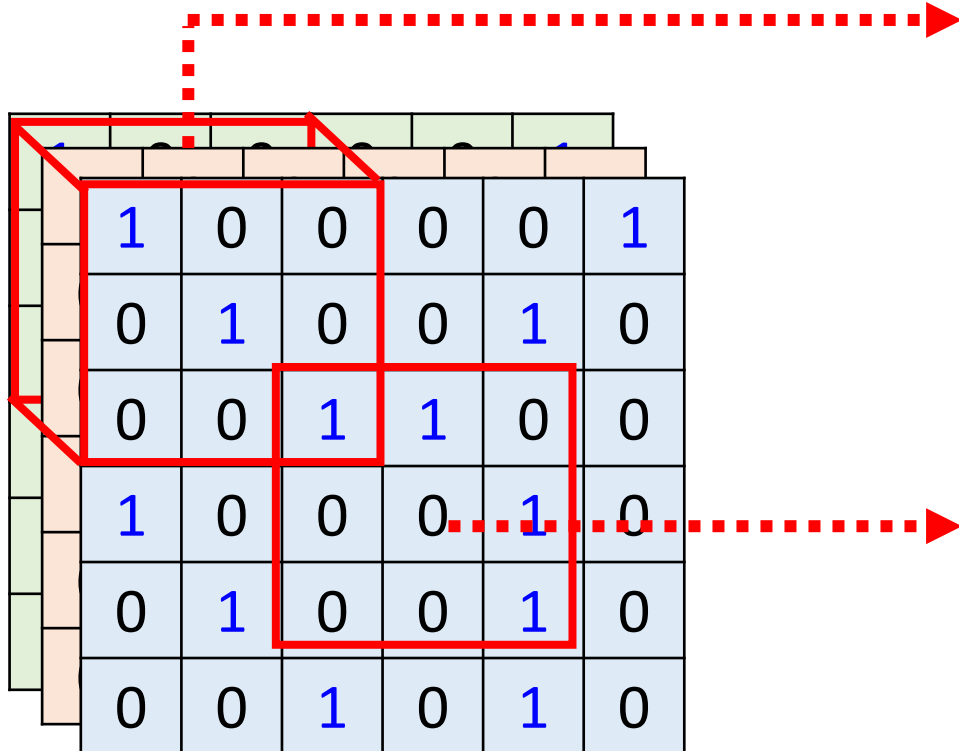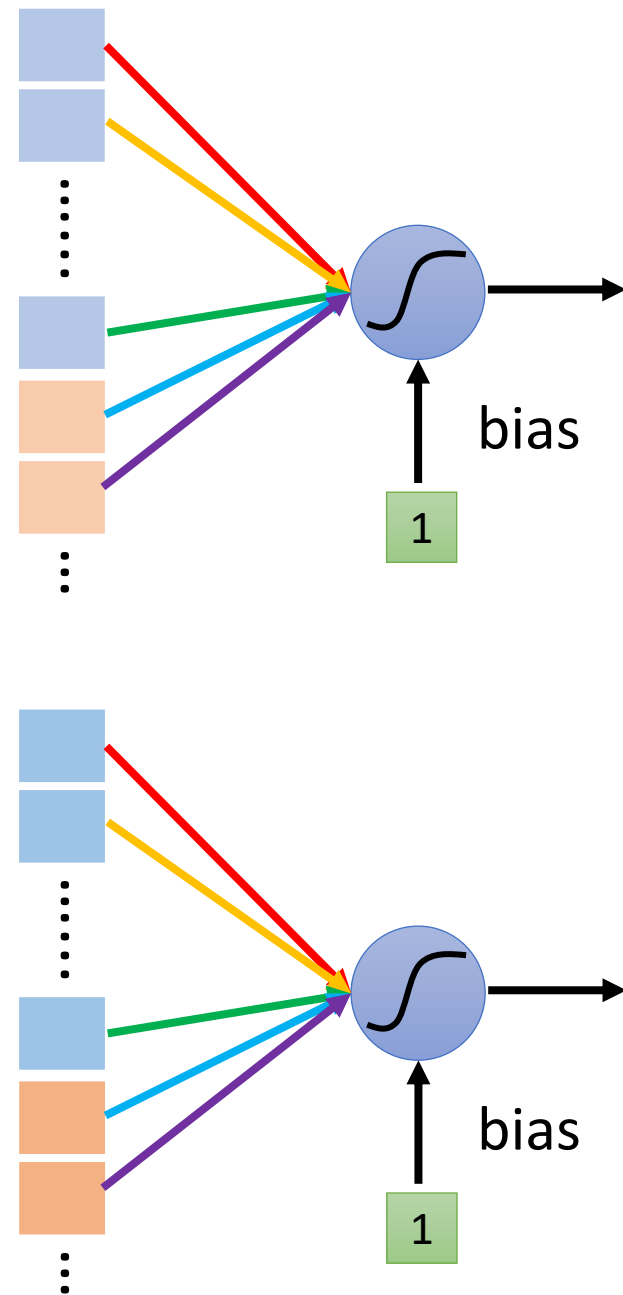| ***Neuron Version Story*** | ***Filter Version Story*** |
|---|---|
| Each neuron only considers a receptive field. | There are a set of filters detecting small patterns. |
| The neurons with different receptive fields share the parameters. | Each filter convolves over the input image. |

They are the same story.

# Observation 3

• Subsampling the pixels will not change the object

bird



subsampling

bird

# Pooling – Max Pooling

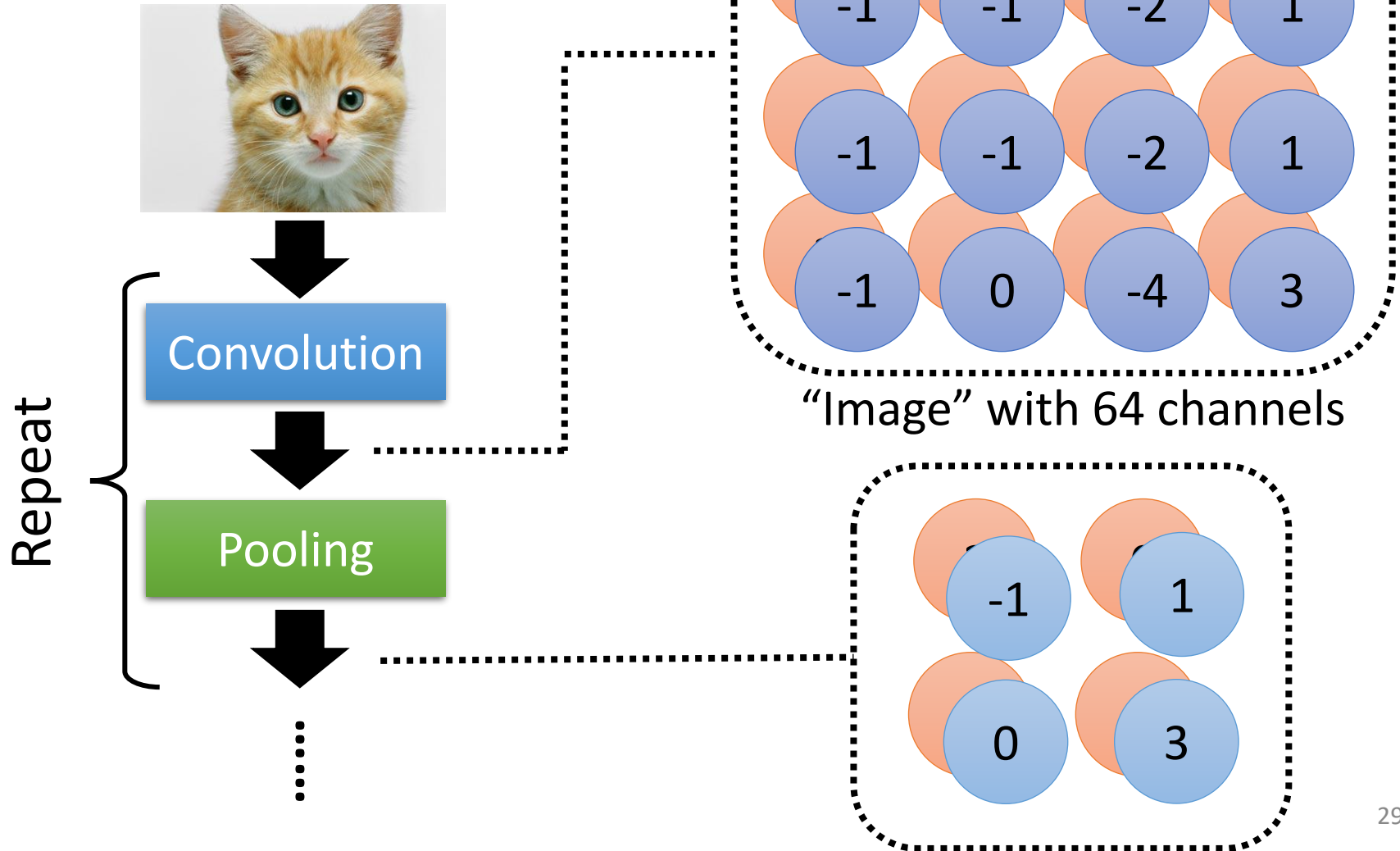| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

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

Filter 2

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

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

# Convolutional Layers + Pooling

Repeat

Convolution

Pooling

| -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



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
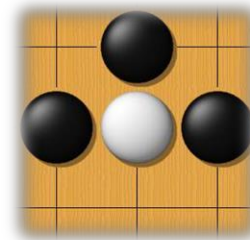
Black: 1

white: -1

none: 0

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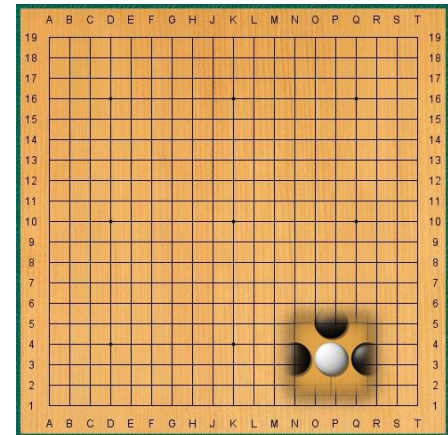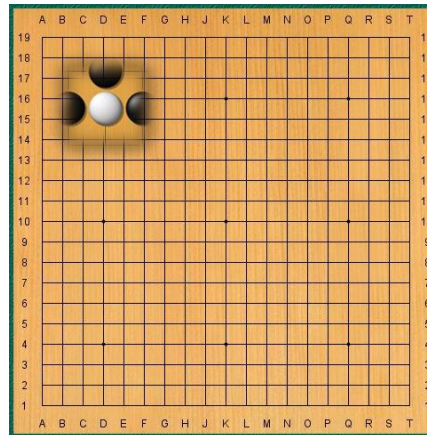
But CNN performs much better.

# 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

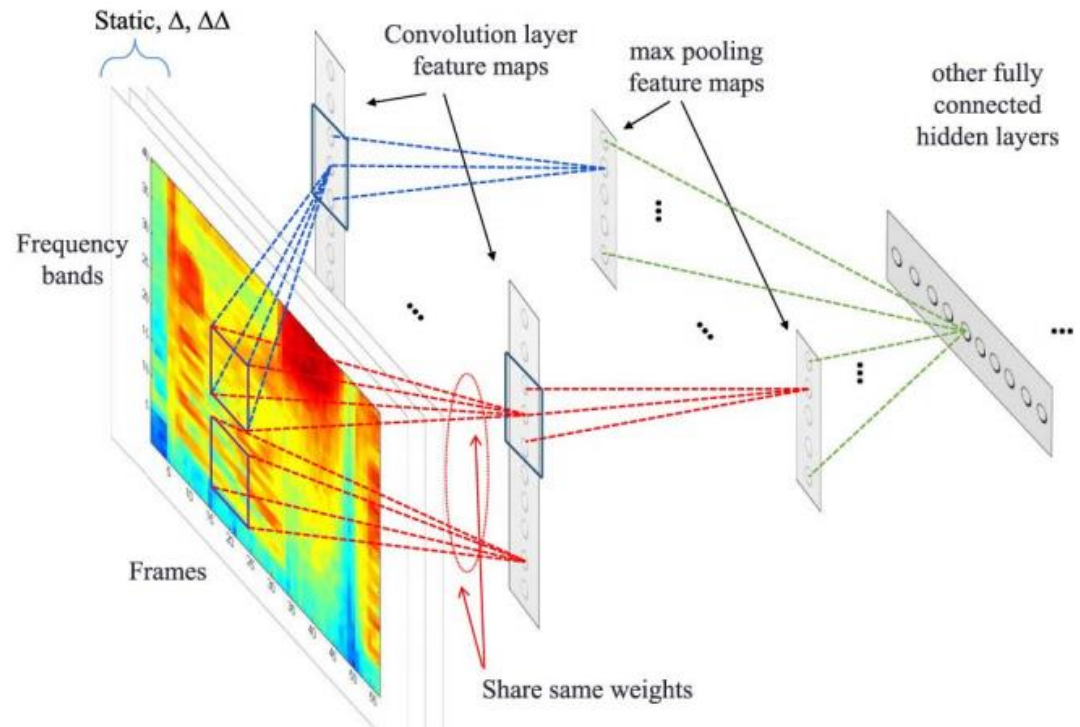➡️ **Pooling** **How to explain this???**

**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 Tabl̄̄ 256 and 384 filters

**Alpha Go does not use Pooling ……**

# *More Applications*



## Speech

https://dl.acm.org/doi/10.110
9/TASLP.2014.2339736

## Natural Language Processing

https://www.aclweb.org/anth
ology/S15-2079/

# To learn more …

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

*Spatial Transformer Layer*

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