

深度學習實作與應用

Deep learning and its applications

Convolutional Networks

IM5062, Spring 2024

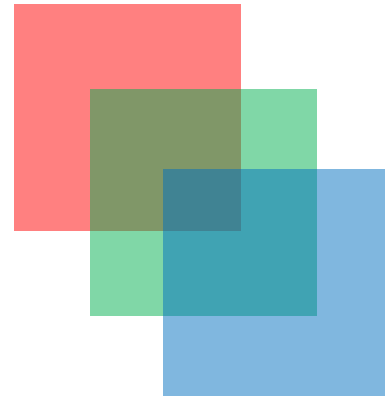
Some slides adopted from Alex Smola, Chien-Yao Wang

黃意婷

# Outline

- What is Convolution
- Use Convolution
- Convolution (d2l.ch7)
  - Kernel/Filter
  - Padding
  - Stride
- Pooling
  - Max/Average Pooling
- Multiple Channels
- Summary
- LeNet/AlexNet (d2l.ch8.1)
- CNN for Go & Text Classification
- Batch normalization (d2l.ch8.5)
- Network in Network (d2l.ch8.3)
- RasNet (d2l.ch8.6)

# Convolution



# Convolution

-1 是弱化  
1 是強化



1		

取平均



模糊化

0	-1	
	1	

強化中間



銳利化

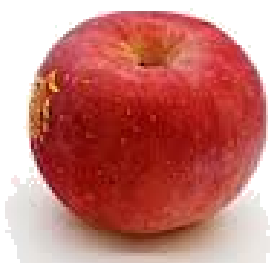
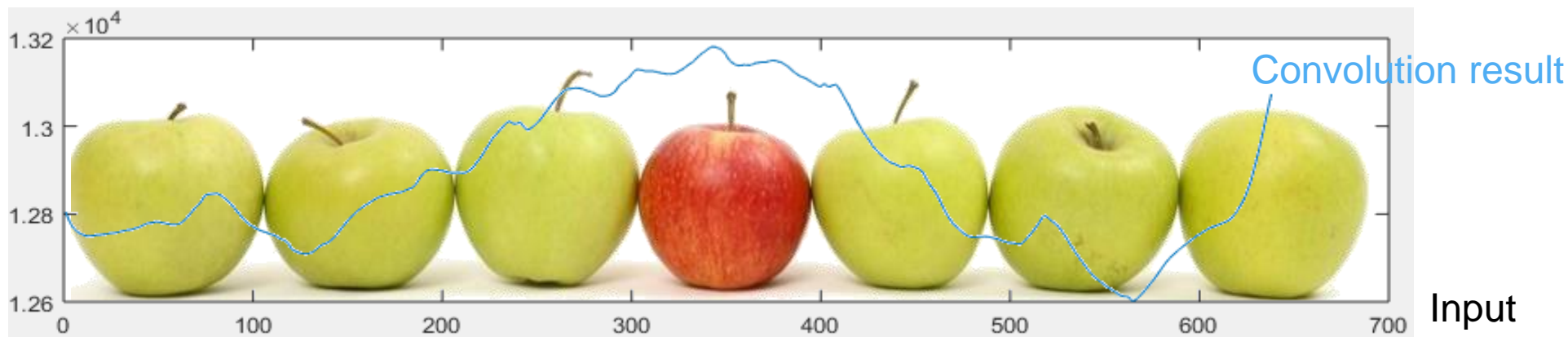
-1		
	0.5	1

邊緣化



浮雕感

# Why use Convolution



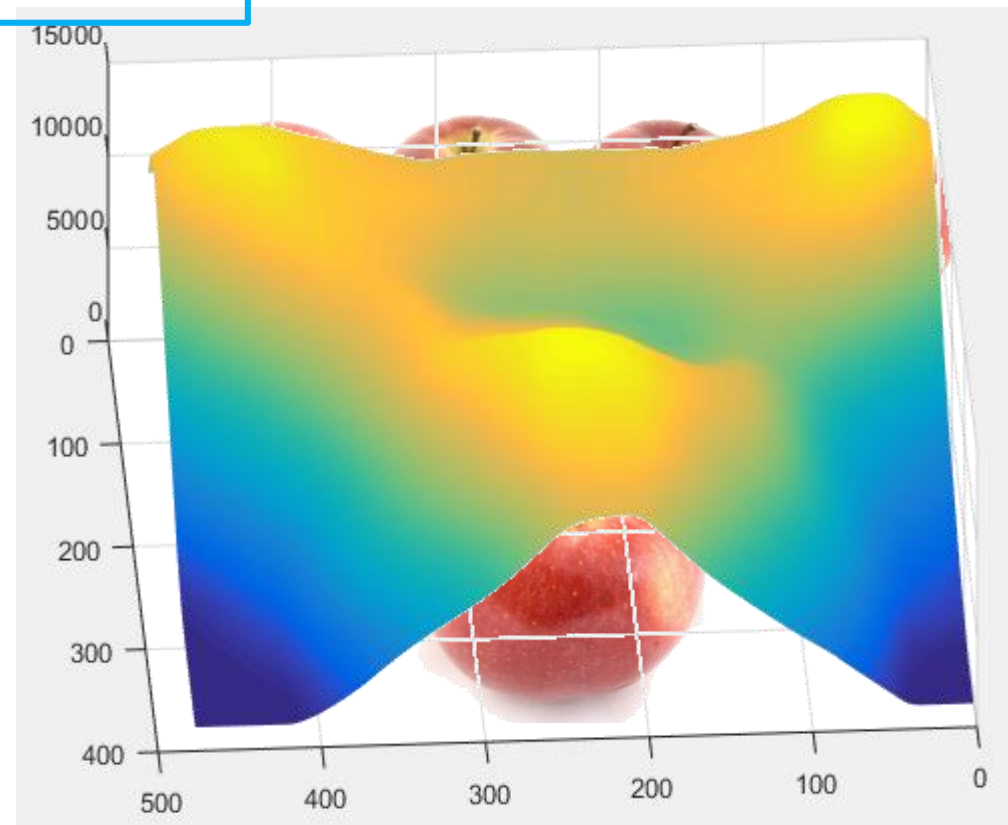
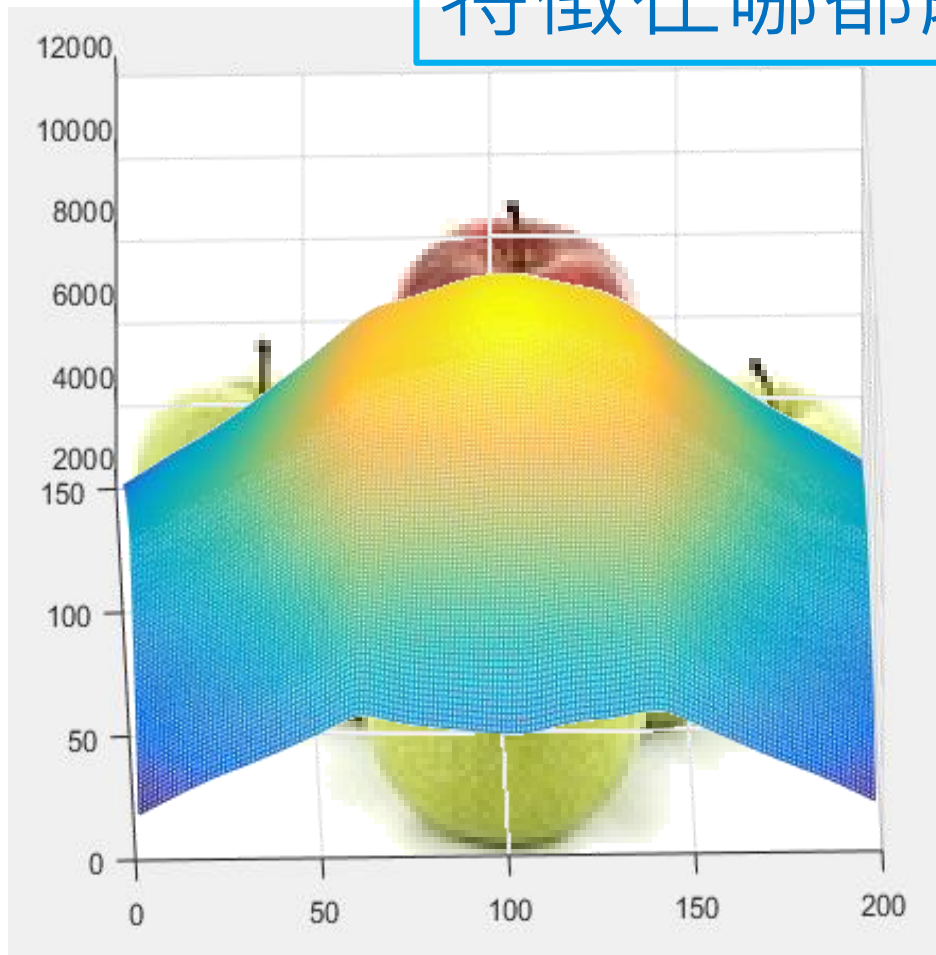
Filter

可以找到特定的特徵

# Why use Convolution

Translation Invariance

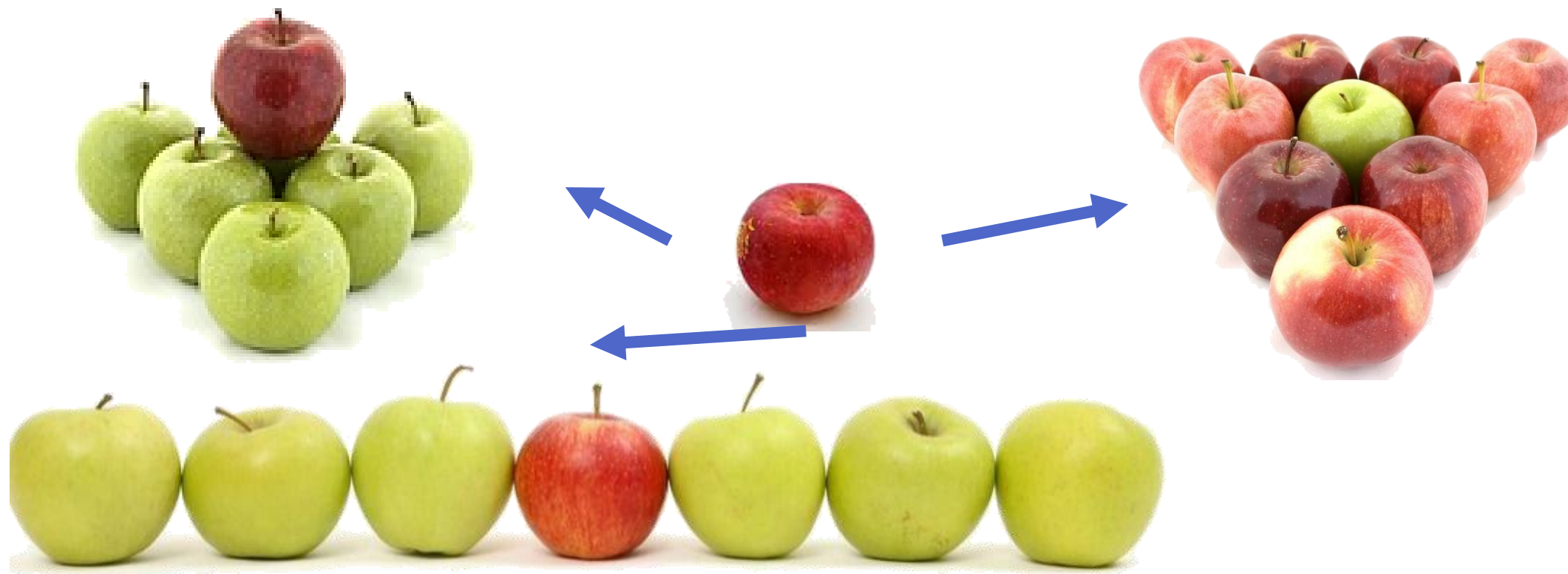
特徵在哪都能被找到





# Why use Convolution

任何輸入大小都能做



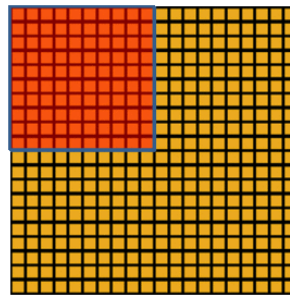
# Convolutional Neural Network

- Convolution
- Pooling
- Full Connection

1 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	0	0
0 <sub>x0</sub>	1 <sub>x1</sub>	1 <sub>x0</sub>	1	0
0 <sub>x1</sub>	0 <sub>x0</sub>	1 <sub>x1</sub>	1	1
0	0	1	1	0
0	1	1	0	0

Image

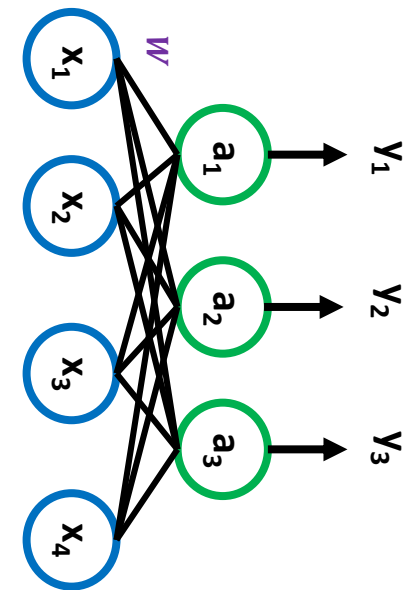
4		

Convolved  
FeatureConvolved  
feature

1	

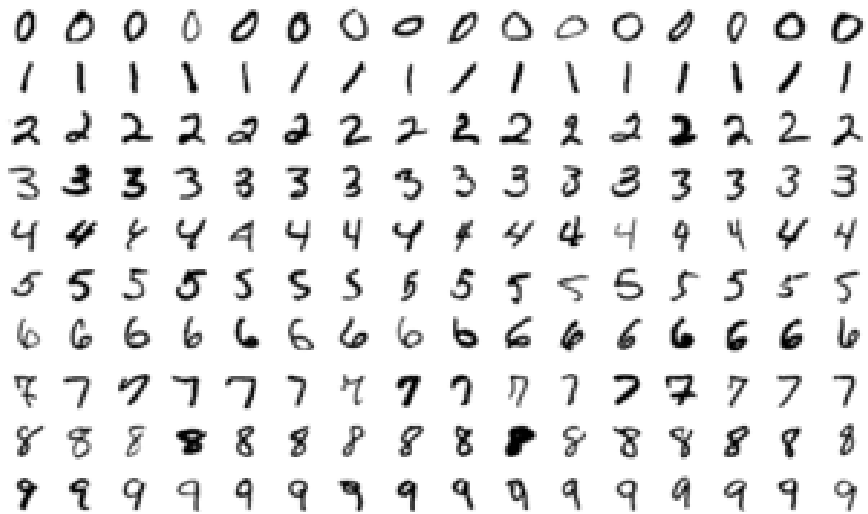
Pooled  
feature

Connection  
Full

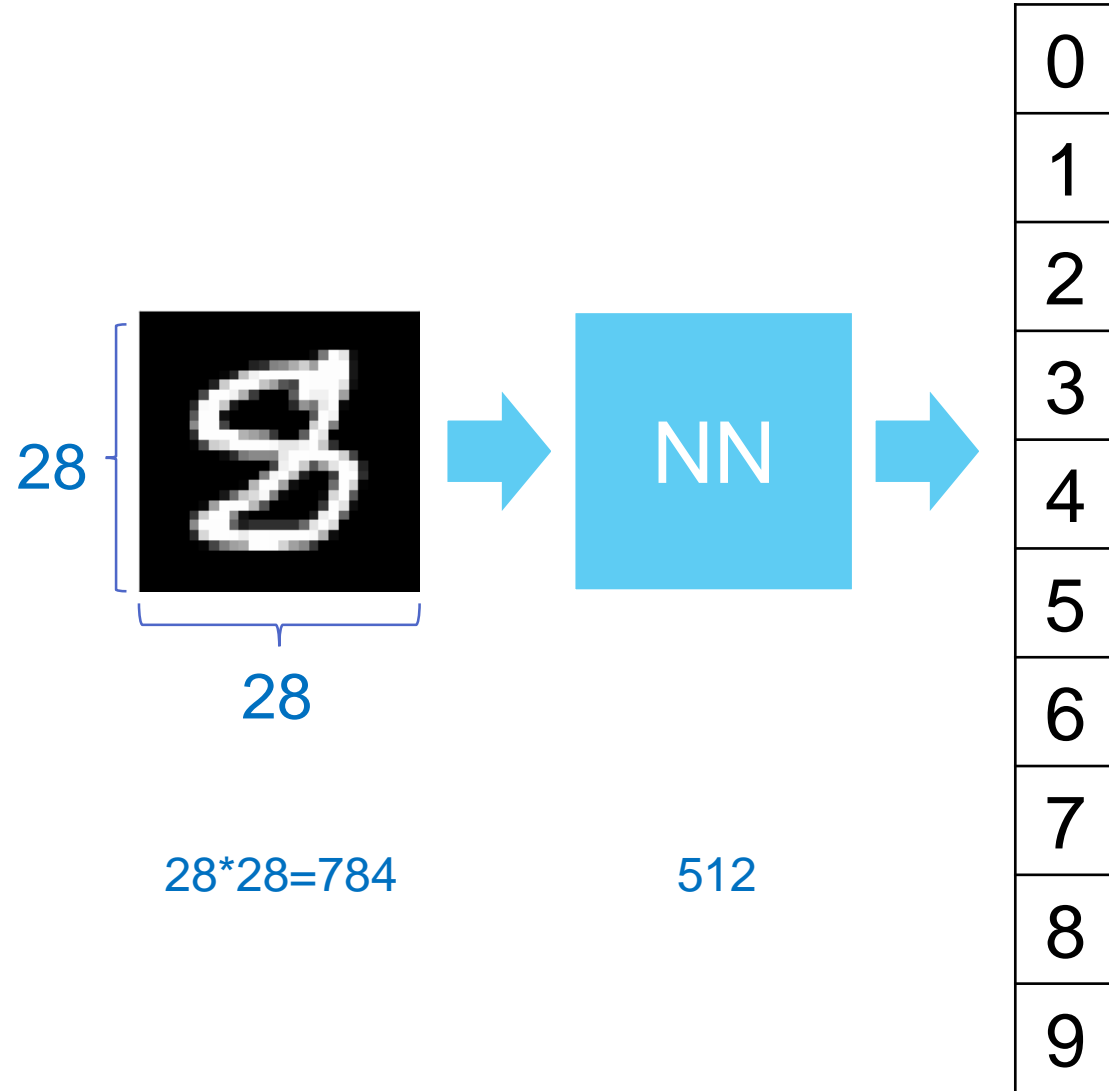




# Image classification

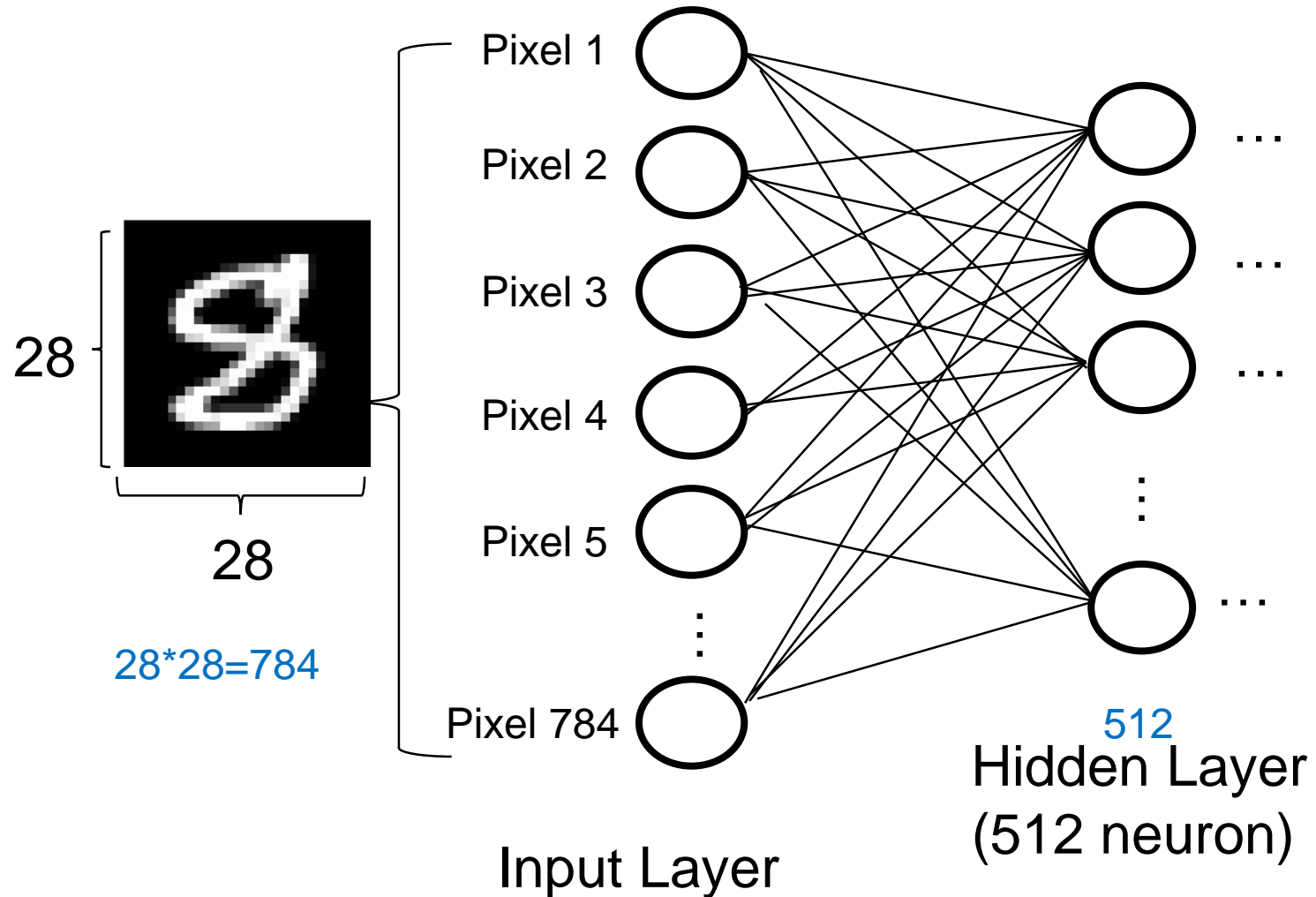


MNIST dataset



# Image classification

$$W=784*512=401,408$$



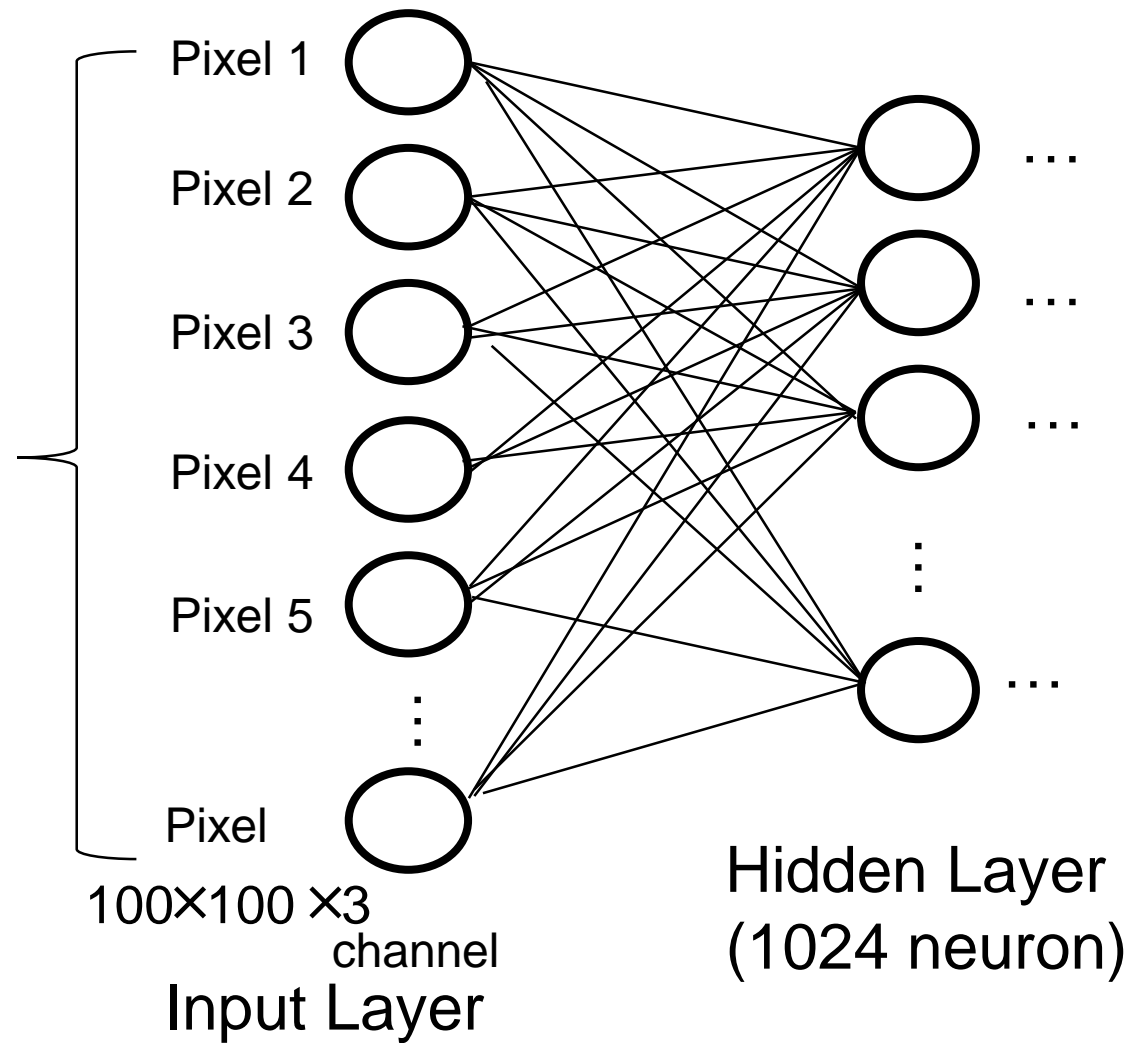
100



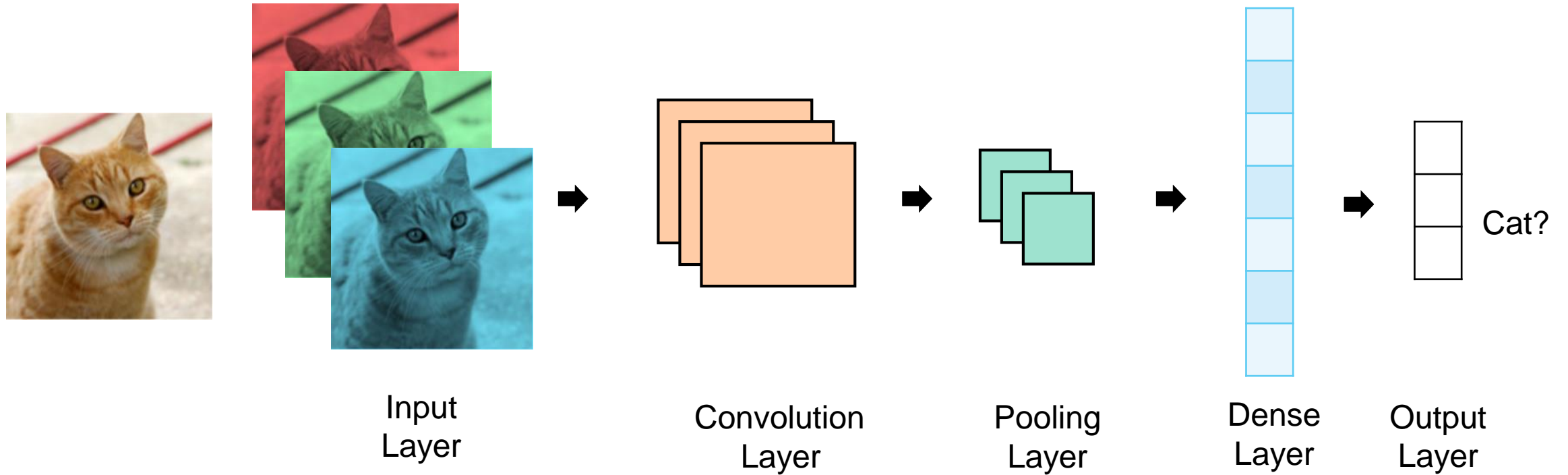
100



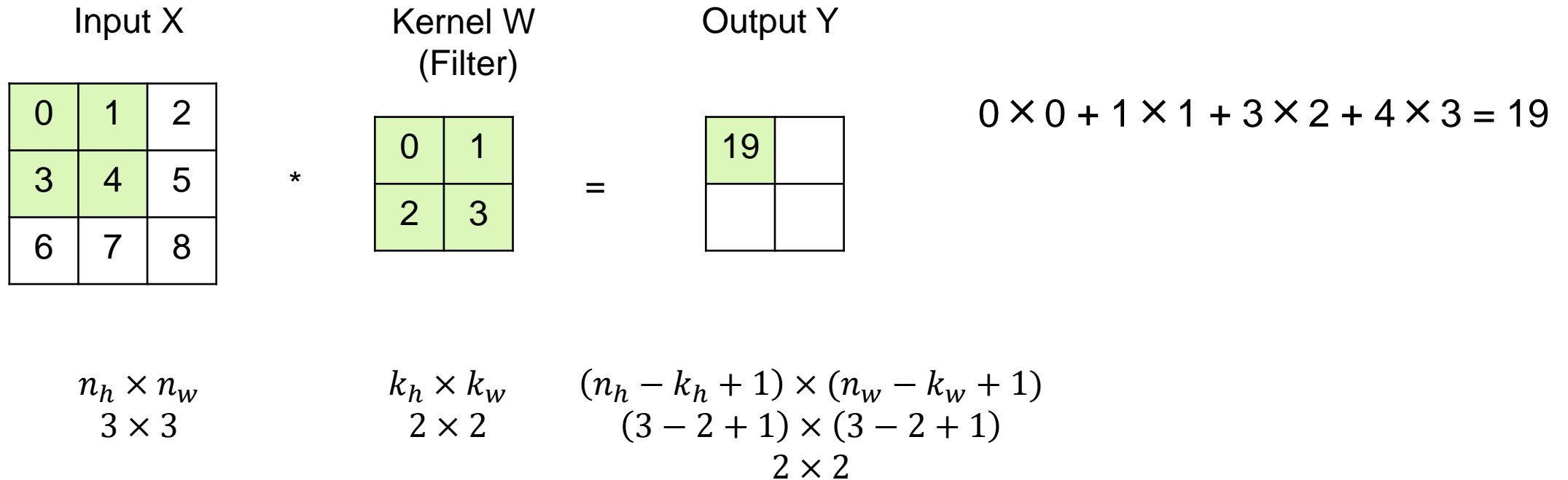
$$W=100*100*3*1024>3*10^7$$



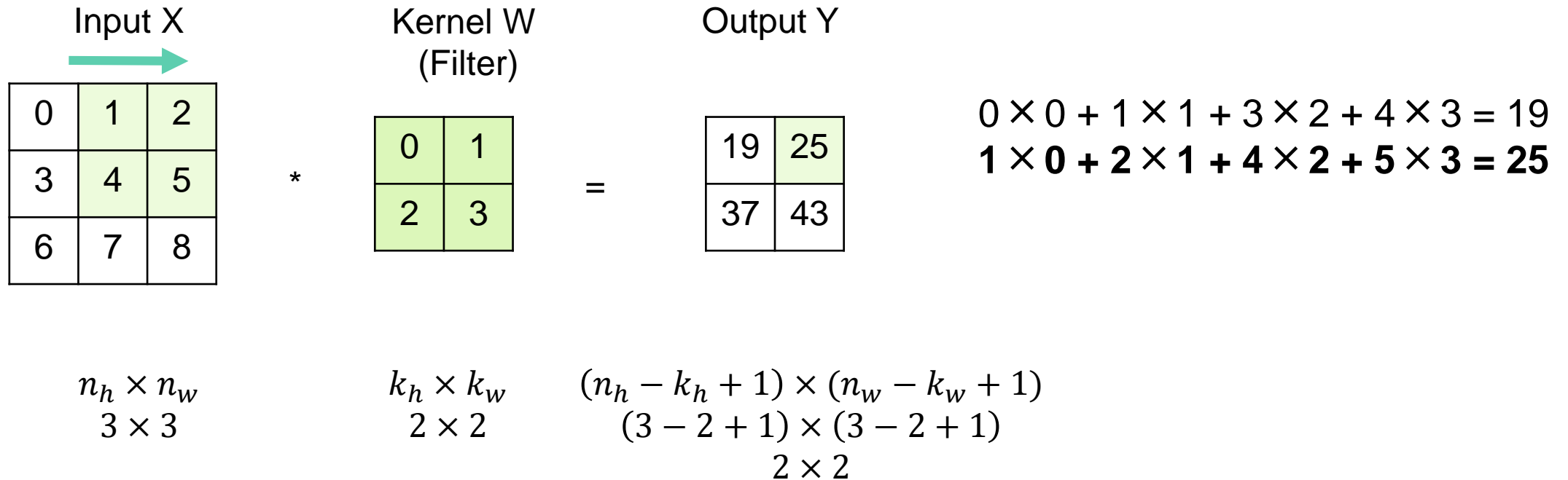
# Convolutional network



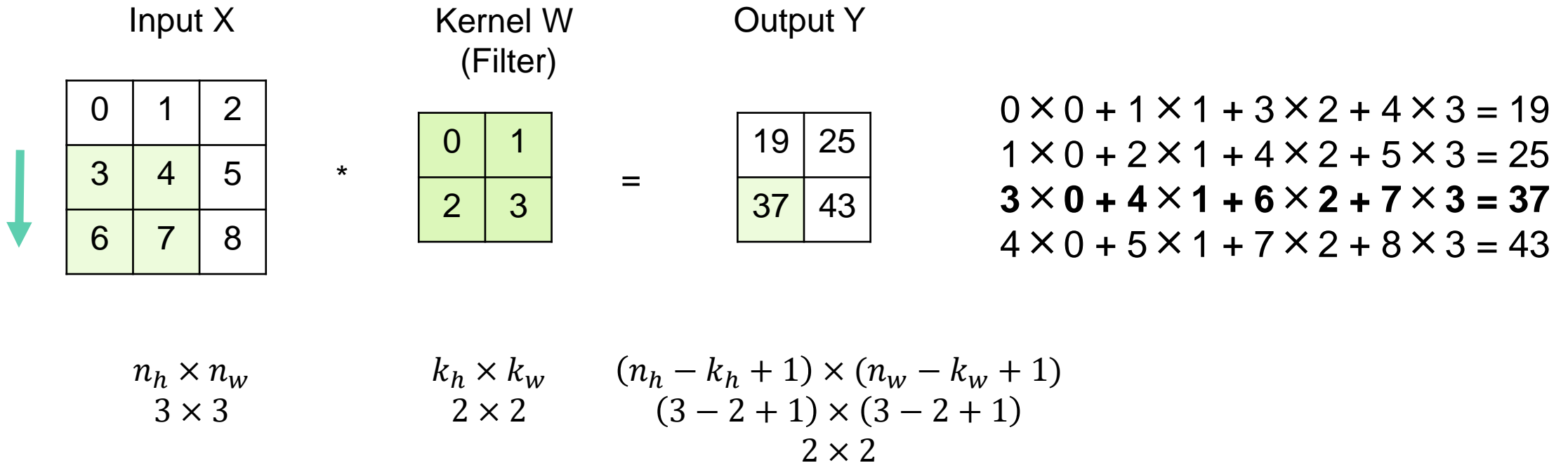
# 2D convolutional layer



# 2D convolutional layer



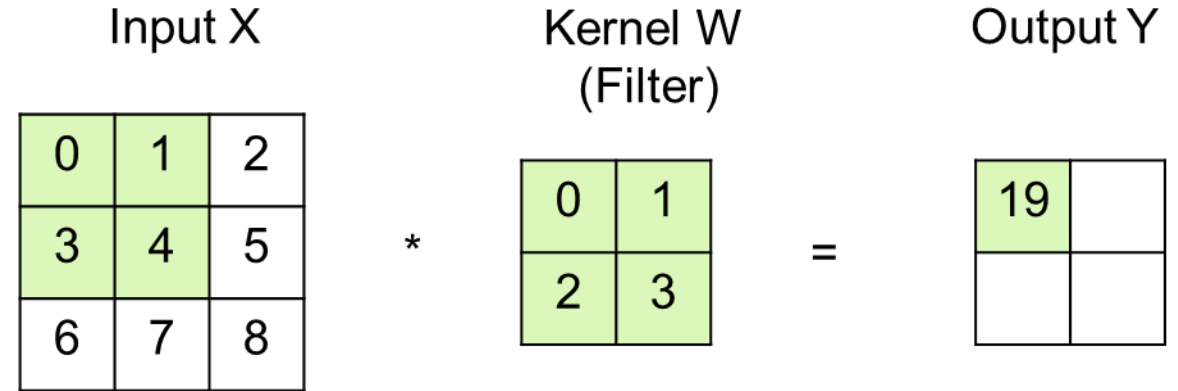
# 2D convolutional layer





# 2D convolutional layer

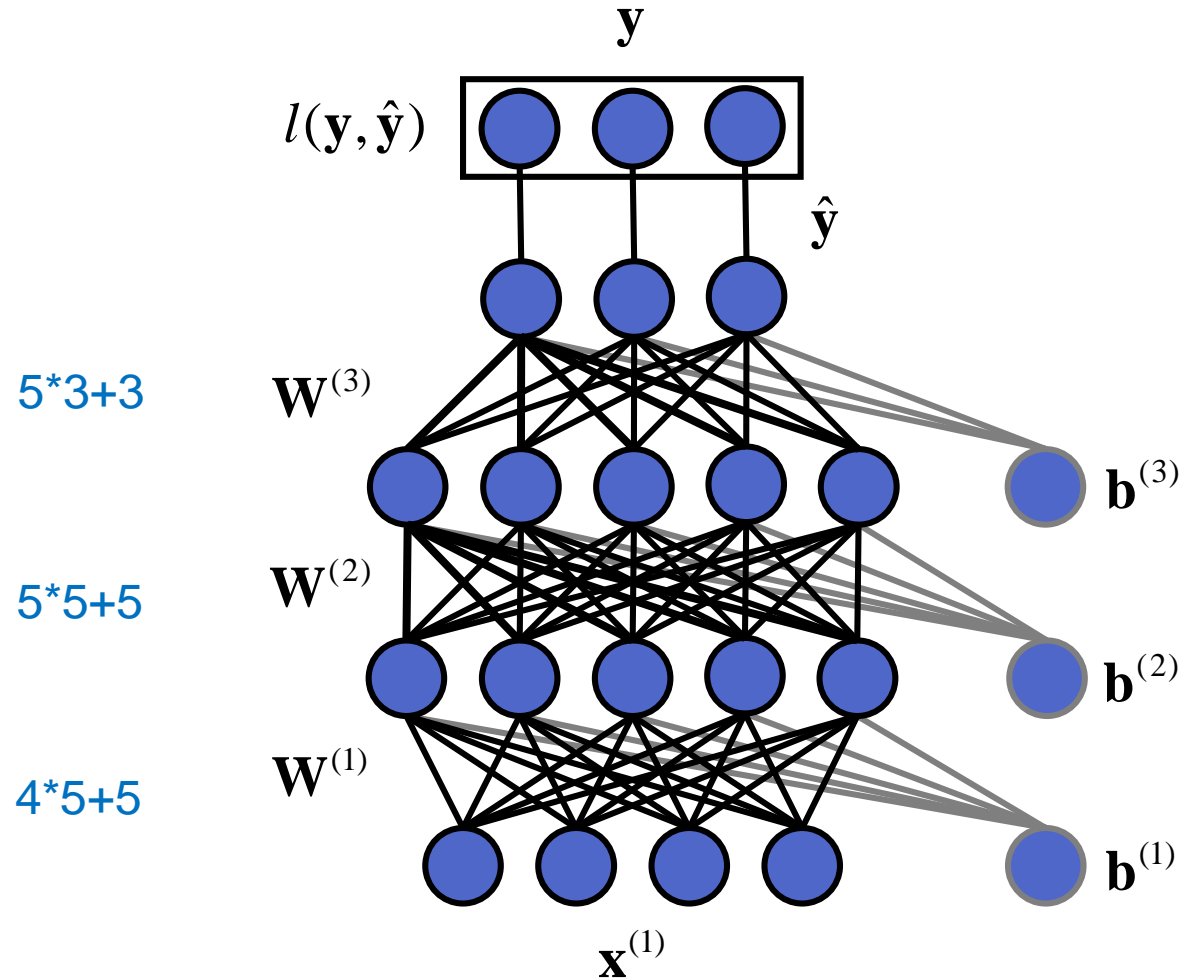
- $\mathbf{X}$ :  $n_h \times n_w$  input matrix
- $\mathbf{W}$ :  $k_h \times k_w$  kernel matrix
- $b$ : scalar bias
- $\mathbf{Y}$ :  $(n_h - k_h + 1) \times (n_w - k_w + 1)$  output matrix



$$\mathbf{Y} = \mathbf{X} \otimes \mathbf{W} + b$$

- $\mathbf{W}$  and  $b$  are learnable parameters.

# Fully-Connected Multilayer Perceptron

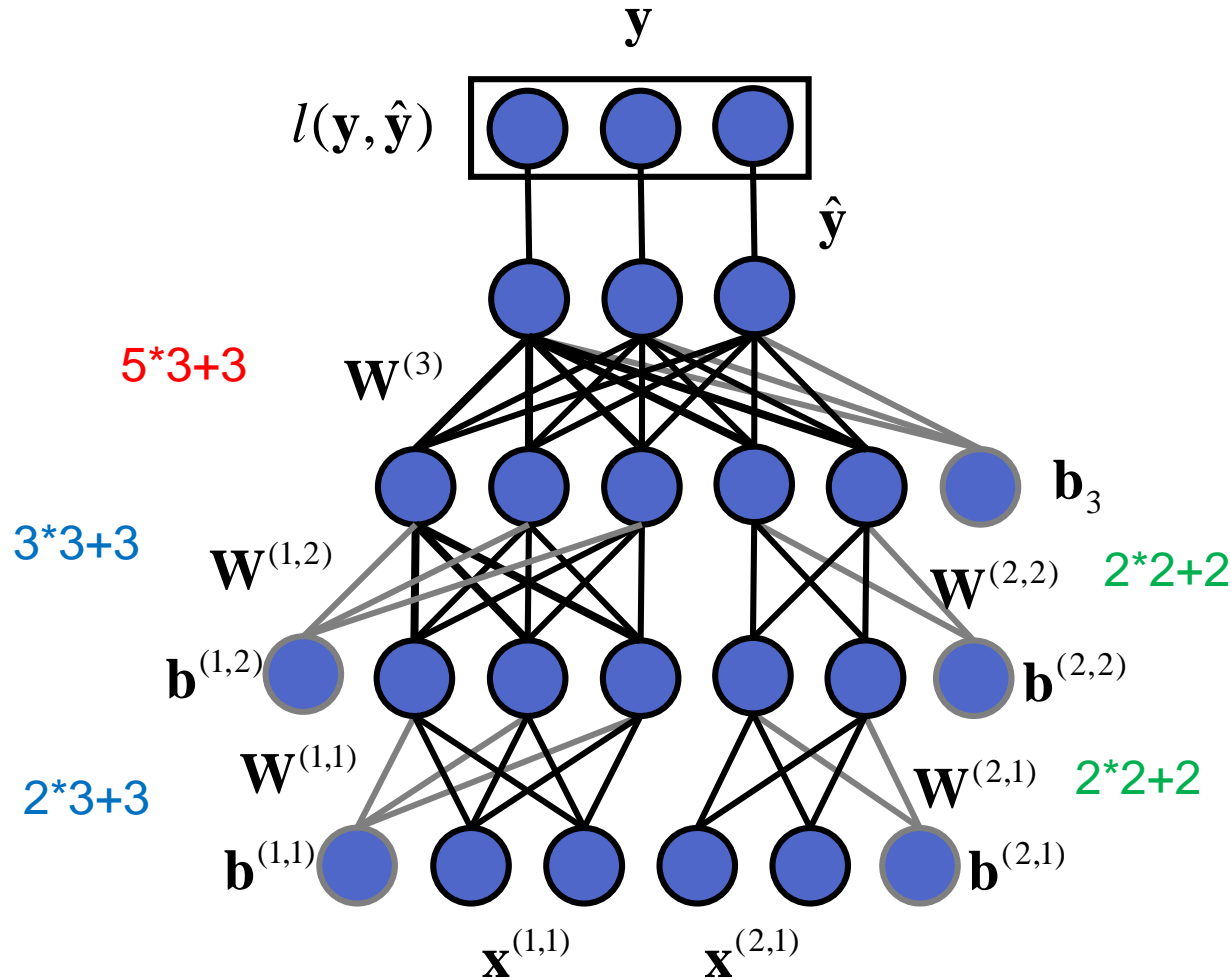


$$\mathbf{x}^{(2)} = \mathbf{W}^{(1)} \mathbf{x}^{(1)} + \mathbf{b}^{(1)}$$

$$\mathbf{x}^{(3)} = \mathbf{W}^{(2)} \mathbf{x}^{(2)} + \mathbf{b}^{(2)}$$

$$\hat{\mathbf{y}} = \mathbf{W}^{(3)} \mathbf{x}^{(3)} + \mathbf{b}^{(3)}$$

# Locally-Connected Multilayer Perceptron



$$\mathbf{x}^{(1,2)} = \mathbf{W}^{(1,1)} \mathbf{x}^{(1,1)} + \mathbf{b}^{(1,1)}$$

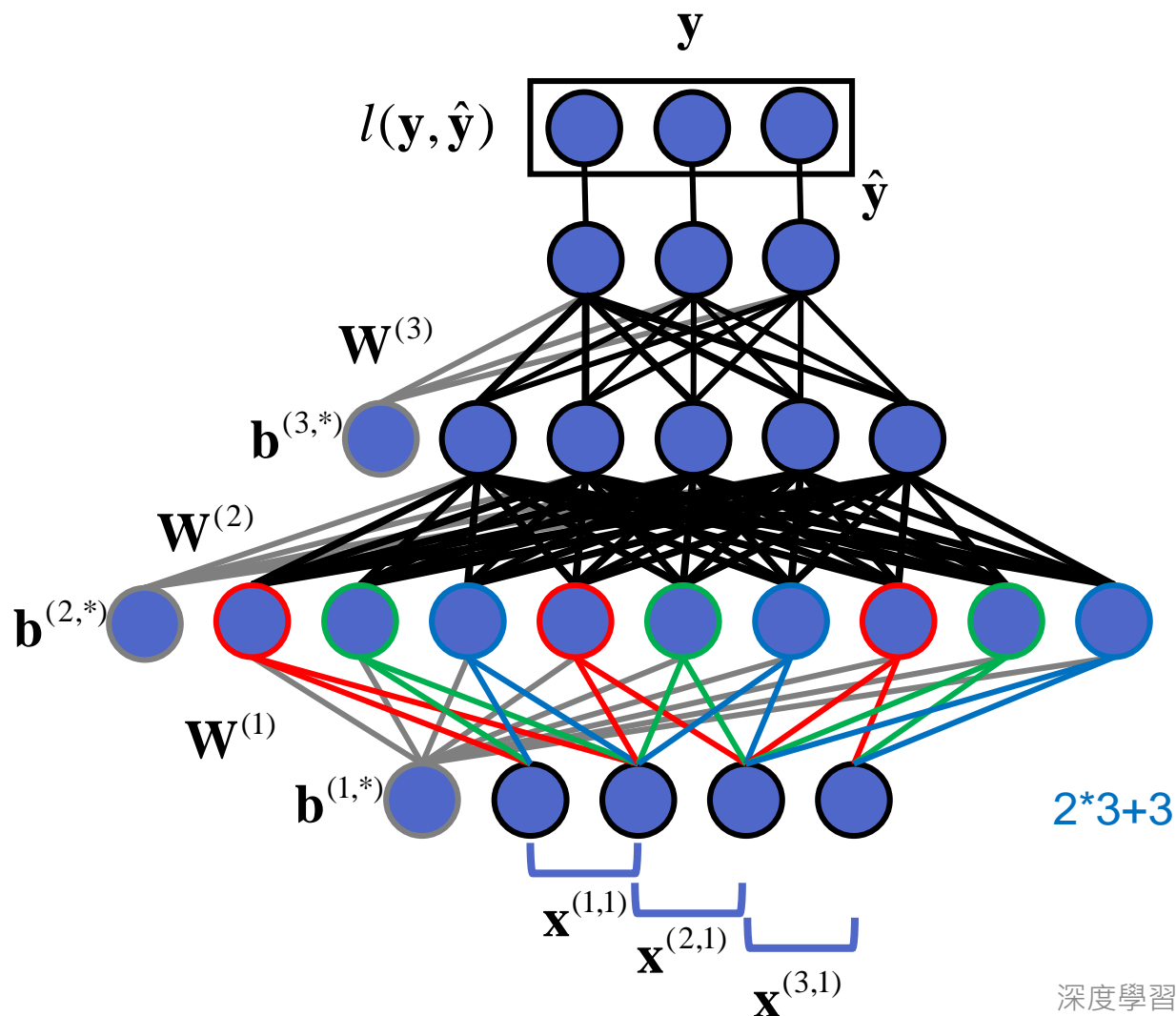
$$\mathbf{x}^{(1,3)} = \mathbf{W}^{(1,2)} \mathbf{x}^{(1,2)} + \mathbf{b}^{(1,2)}$$

$$\mathbf{x}^{(2,2)} = \mathbf{W}^{(2,1)} \mathbf{x}^{(2,1)} + \mathbf{b}^{(2,1)}$$

$$\mathbf{x}^{(2,3)} = \mathbf{W}^{(2,2)} \mathbf{x}^{(2,2)} + \mathbf{b}^{(2,2)}$$

$$\hat{\mathbf{y}} = \mathbf{W}^{(3)} [\mathbf{x}^{(1,3)}, \mathbf{x}^{(2,3)}] + \mathbf{b}^{(3)}$$

# Convolutional Neural Network

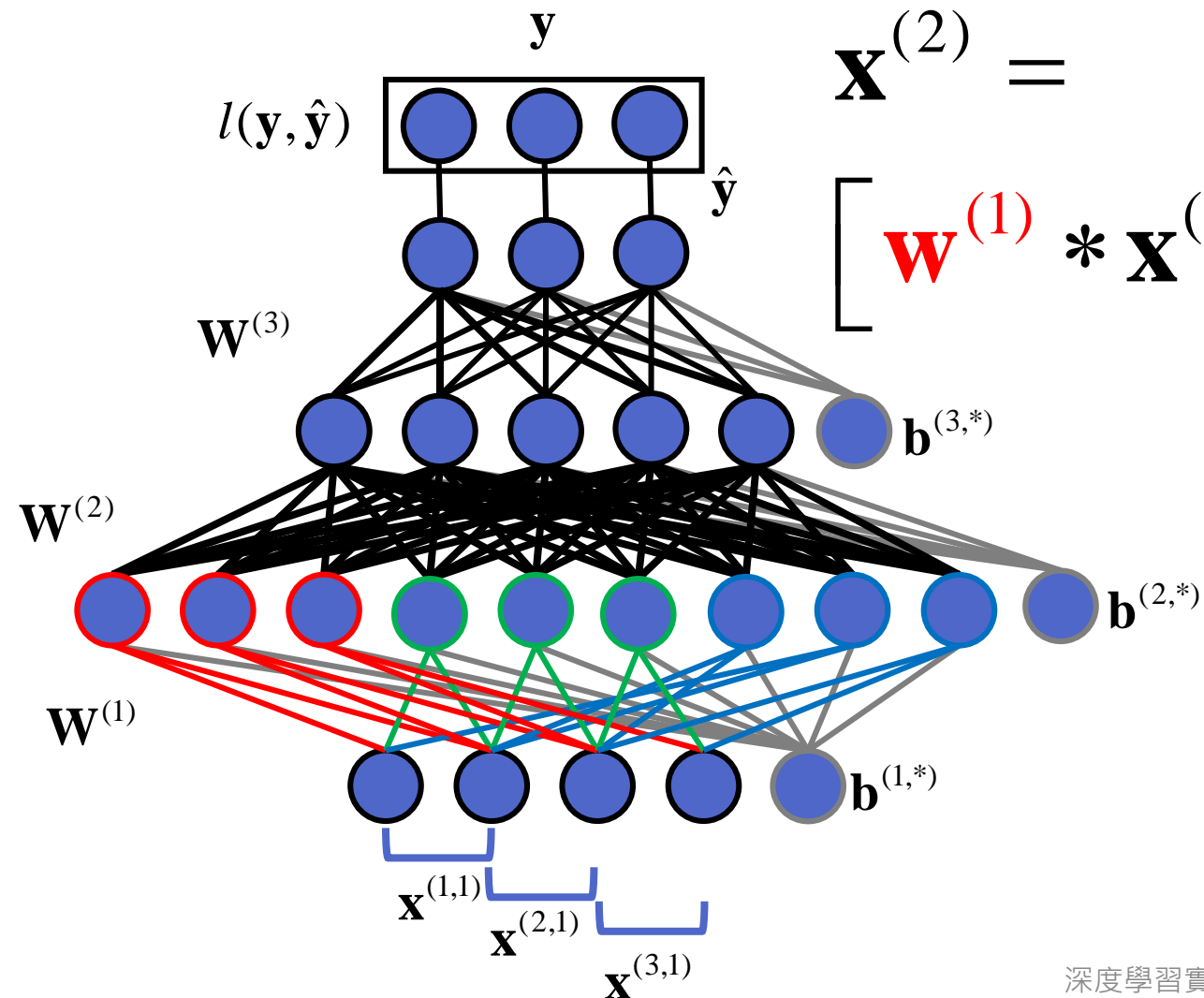


$$\mathbf{x}^{(2)} = \mathbf{W}^{(1)} * \mathbf{x}^{(1)} + \mathbf{b}^{(1)}$$

$$\mathbf{x}^{(3)} = \mathbf{W}^{(2)} \mathbf{x}^{(2)} + \mathbf{b}^{(2)}$$

$$\hat{y} = \mathbf{W}^{(3)} \mathbf{x}^{(3)} + \mathbf{b}^{(3)}$$

# Convolutional Neural Network

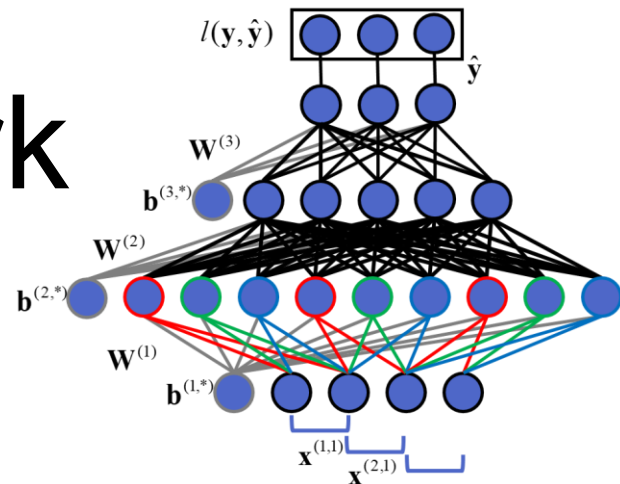


$$\mathbf{X}^{(2)} =$$

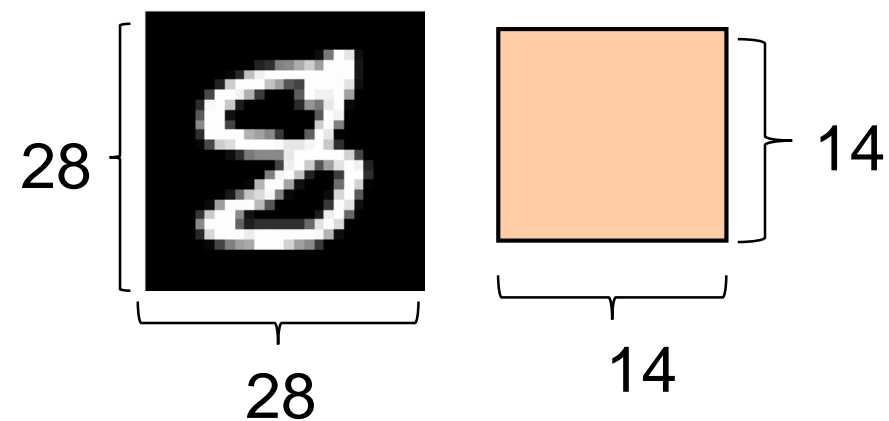
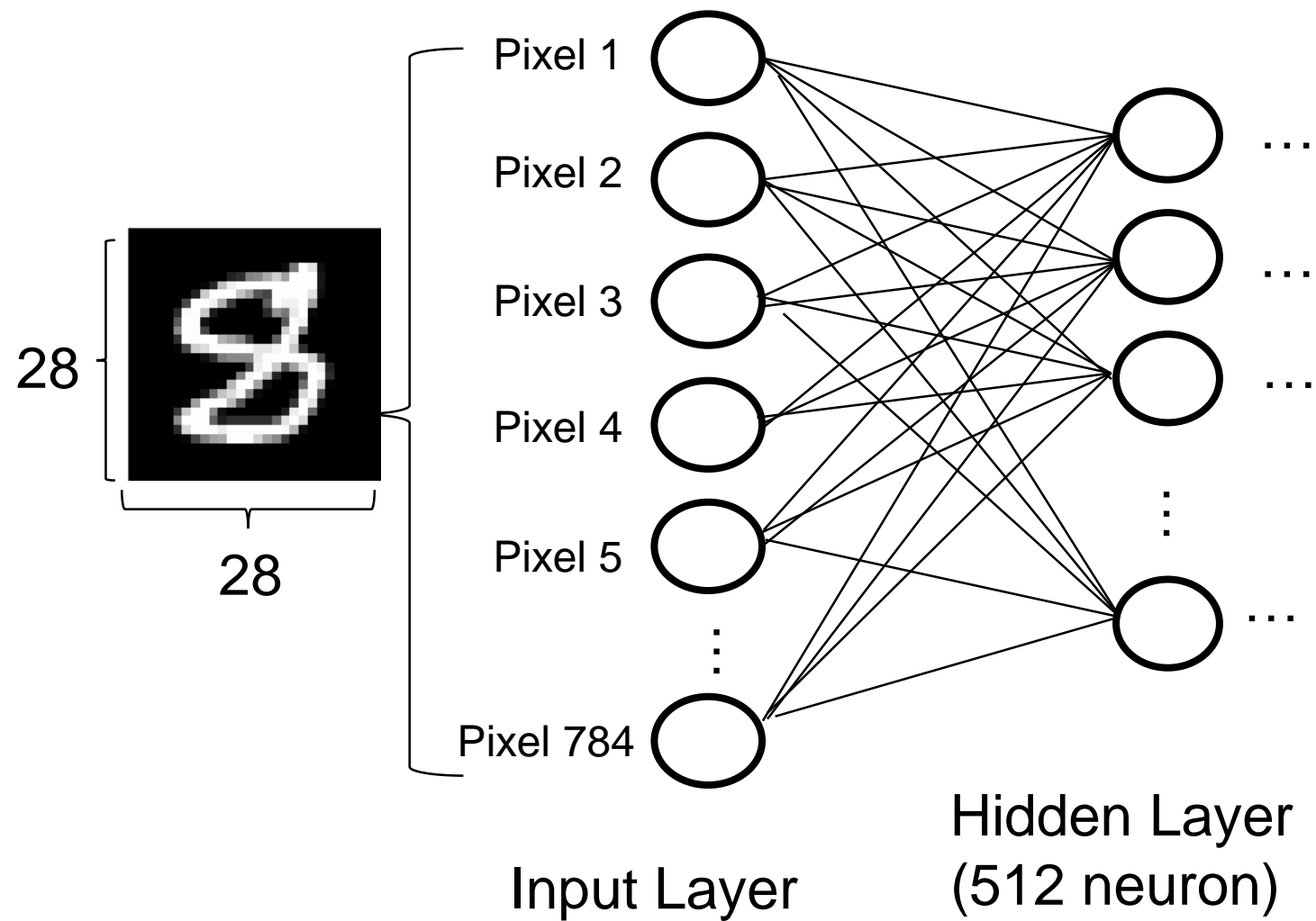
$$\left[ \mathbf{w}^{(1)} * \mathbf{x}^{(1)}, \mathbf{w}^{(1)} * \mathbf{x}^{(1)}, \mathbf{w}^{(1)} * \mathbf{x}^{(1)} \right] + \mathbf{b}^{(1)}$$

$$\mathbf{x}^{(3)} = \mathbf{W}^{(2)} \mathbf{x}^{(2)} + \mathbf{b}^{(2)}$$

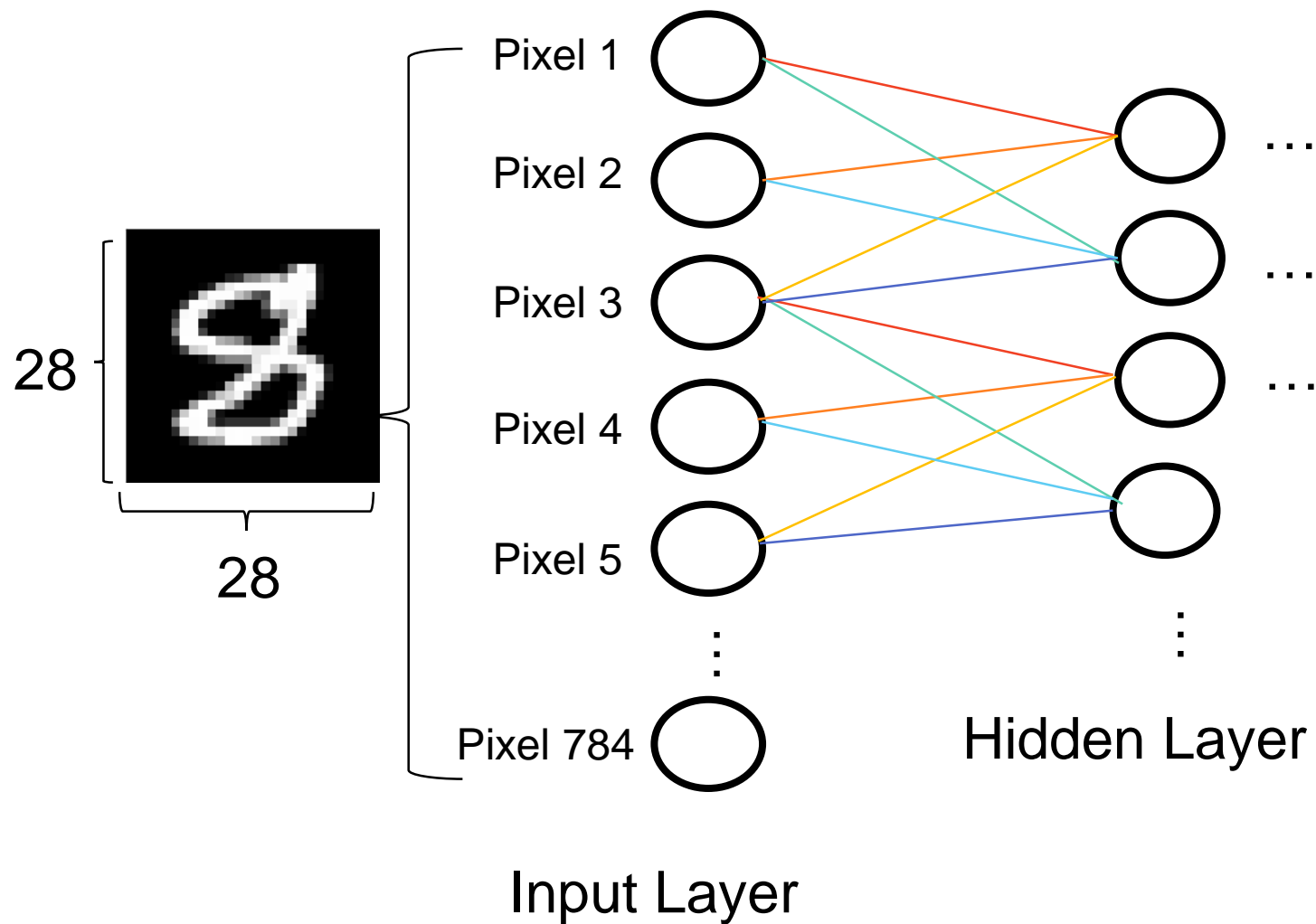
$$\hat{y} = \mathbf{W}^{(3)} \mathbf{x}^{(3)} + \mathbf{b}^{(3)}$$



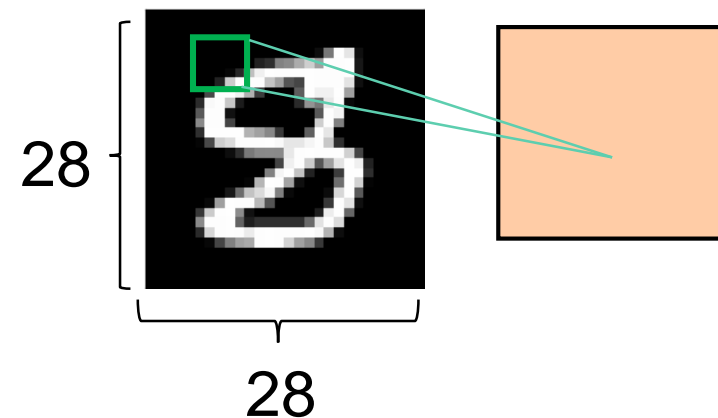
# DNN & CNN



# DNN & CNN



- Each neuron only connects to **part** of previous layer.  
→ parameter sharing  
→ Less parameters than fully connected layer.

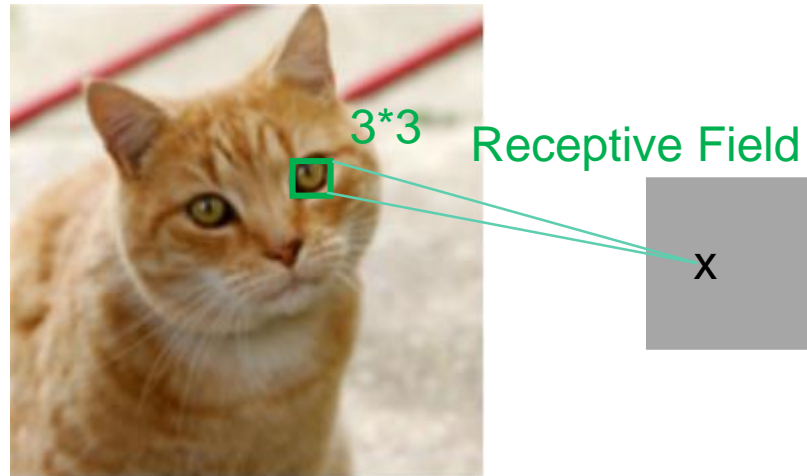


Due to parameter sharing  
Idea #1 – Translation Invariance  
Idea #2 – Locality



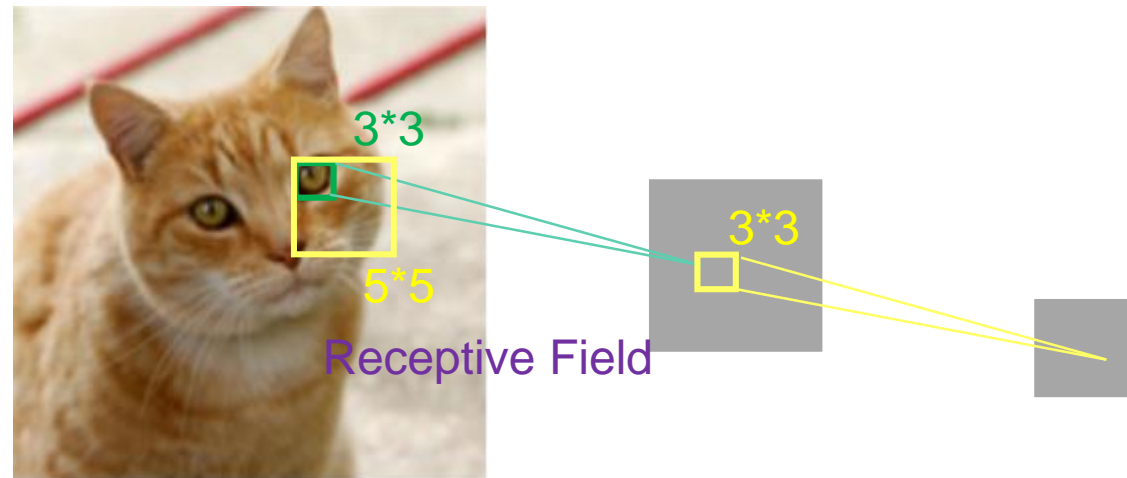
# Receptive Field (filter size)

- For any element  $x$  of some layer, its receptive field refers to all the elements (from all the previous layers) that may affect the calculation of  $x$  during the forward propagation.

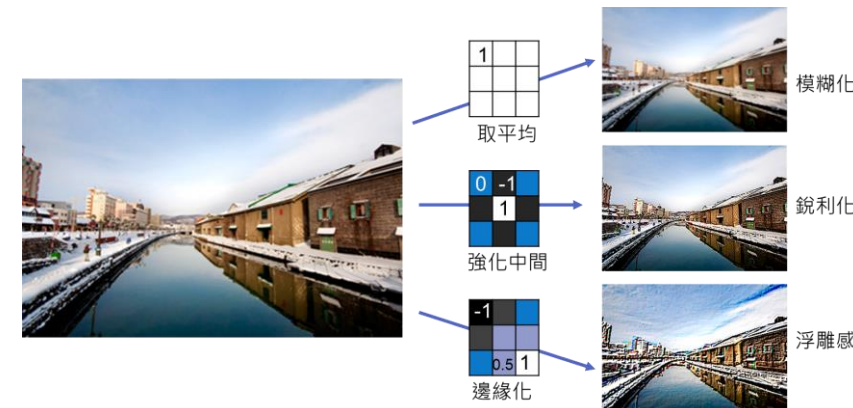


# Receptive Field (filter size)

- For any element  $x$  of some layer, its receptive field refers to all the elements (from all the previous layers) that may affect the calculation of  $x$  during the forward propagation.



# Object Edge Detection in Images



Input X

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

Kernel W  
(Filter)

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

=

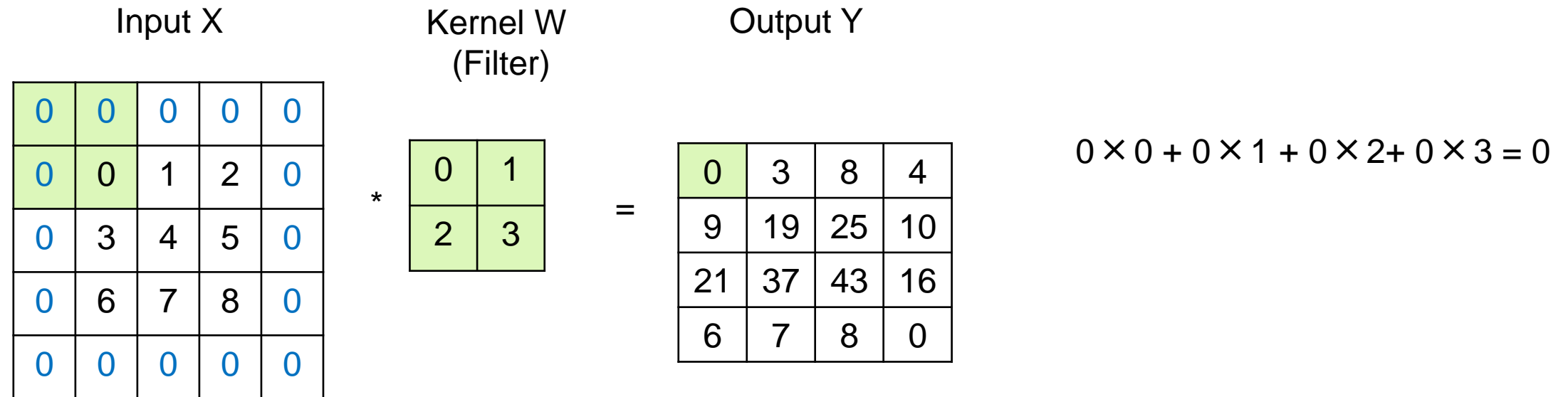
Output Y

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

Edge

# Padding

- Padding adds row/columns around input.



$$(n_h + p_h) \times (n_w + p_w)$$

$$(3 + 2) \times (3 + 2)$$

$$5 \times 5$$

$$k_h \times k_w$$

$$2 \times 2$$

$$(n_h - k_h + p_h + 1) \times (n_w - k_w + p_w + 1)$$

$$(3 - 2 + 2 + 1) \times (3 - 2 + 2 + 1)$$

$$4 \times 4$$

$$p_h \times p_w$$

$$2 \times 2$$

# Padding

- Padding adds row/columns around input.

Input X

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

Kernel W  
(Filter)

0	1
2	3

Output Y

0	3	8	4
9	19	25	10
21	37	43	16
6	7	8	0

$$(n_h + p_h) \times (n_w + p_w)$$

$$(3 + 2) \times (3 + 2)$$

$$5 \times 5$$
  

$$p_h \times p_w$$

$$2 \times 2$$

$$k_h \times k_w$$

$$2 \times 2$$

$$(n_h - k_h + p_h + 1) \times (n_w - k_w + p_w + 1)$$

$$(3 - 2 + 2 + 1) \times (3 - 2 + 2 + 1)$$

$$4 \times 4$$

$$0 \times 0 + 0 \times 1 + 0 \times 2 + 0 \times 3 = 0$$

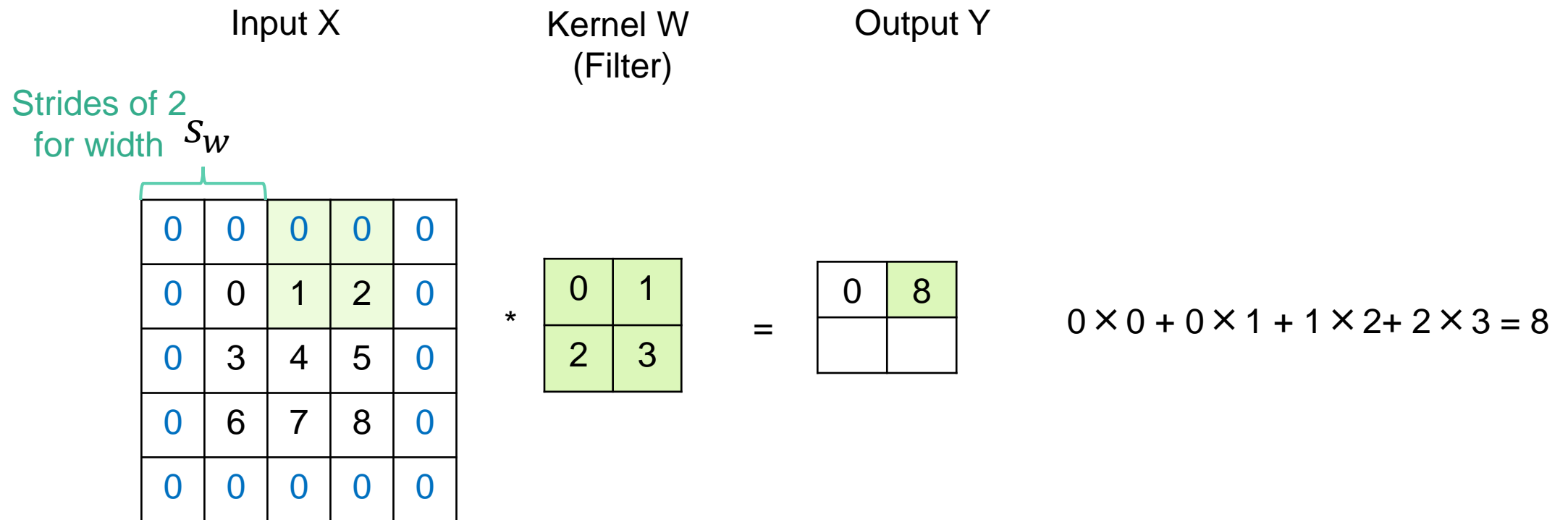
$$0 \times 0 + 0 \times 1 + 0 \times 2 + 1 \times 3 = 3$$

# Padding for consistent input/output shape

- Padding  $p_h$  rows and  $p_w$  columns, output shape will be  $(n_h - k_h + p_h + 1) \times (n_w - k_w + p_w + 1)$
- A common choice is  $p_h = k_h - 1$  and  $p_w = k_w - 1$ 
  - Odd  $k_h$  : pad  $\frac{p_h}{2}$  on both sides
    - e.g.,  $k_h = 3, \frac{p_h}{2}=1$ (上下各一);  $k_w = 3, \frac{p_w}{2}=1$ (左右各一)
  - Even  $k_h$ : pad  $\left\lceil \frac{p_h}{2} \right\rceil$  rows on the top,  $\left\lfloor \frac{p_h}{2} \right\rfloor$  rows on the bottom.
    - e.g.,  $k_h = 2, p_h = 1, \left\lceil \frac{p_h}{2} \right\rceil = \left\lceil \frac{1}{2} \right\rceil = 1, \left\lfloor \frac{p_h}{2} \right\rfloor = \left\lfloor \frac{1}{2} \right\rfloor = 0$

# Stride (How to move)

- Stride is the #rows/#columns per slide





# Stride

- Stride is the #rows/#columns per slide

Input X

Strides of 2  
for width  $s_w$

Strides of 3  
for height  $s_h$

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

Kernel W  
(Filter)

0	1
2	3

\*

=

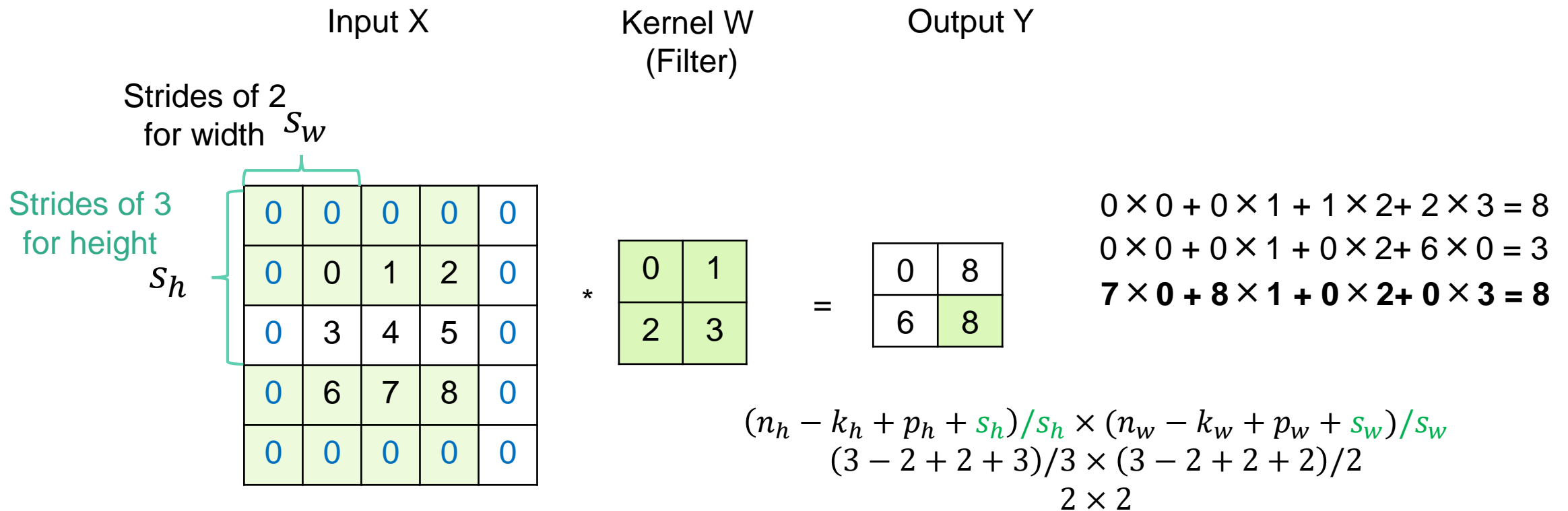
Output Y

0	8
6	

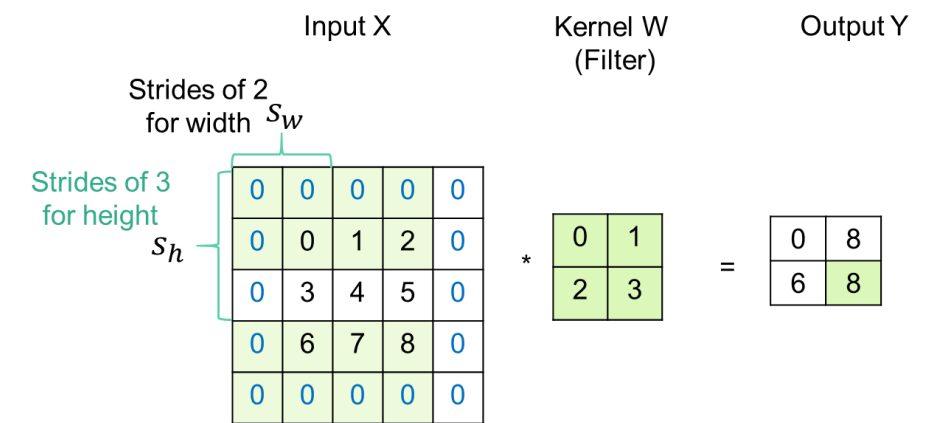
$$0 \times 0 + 0 \times 1 + 1 \times 2 + 2 \times 3 = 8$$
$$0 \times 0 + 0 \times 1 + 0 \times 2 + 6 \times 0 = 3$$

# Stride

- Stride is the #rows/#columns per slide



# Stride



- Given stride  $s_h$  for the height and stride  $s_w$  for the width, the output shape is

$$(n_h - k_h + p_h + s_h)/s_h \times (n_w - k_w + p_w + s_w)/s_w$$

- With  $p_h = k_h - 1$  and  $p_w = k_w - 1$  (超出範圍的不計算)  

$$\lfloor (n_h + s_h - 1)/s_h \rfloor \times \lfloor (n_w + s_w - 1)/s_w \rfloor$$

- If input height/width are divisible by strides  

$$(n_h/s_h) \times (n_w/s_w)$$

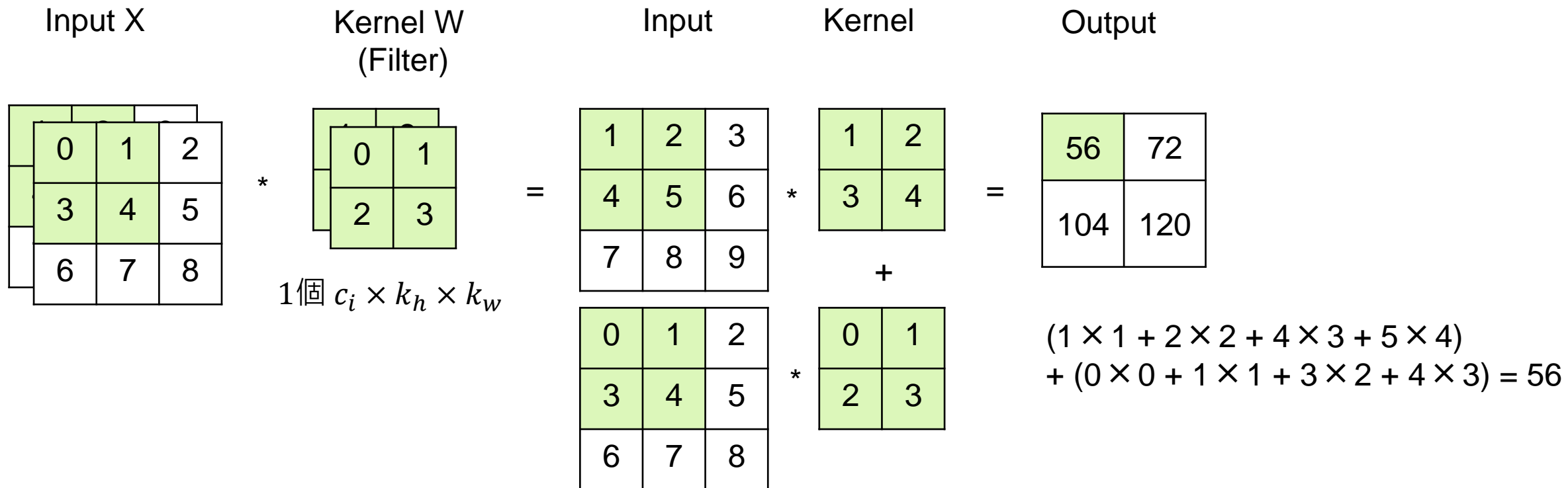
# Multiple Channels

- Color image may have three RGB channels
- Converting to grayscale loses information



# Multiple **Input** Channels

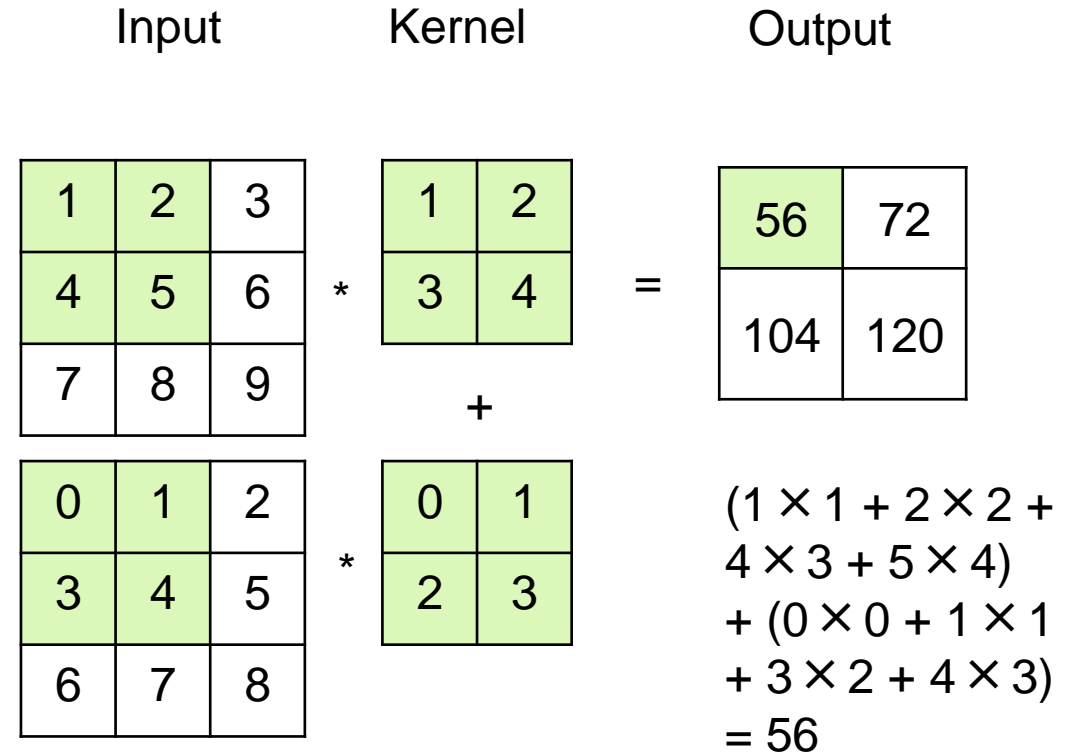
- Have a kernel for each channel, and then **sum** results over channels.



# Multiple **Input** Channels

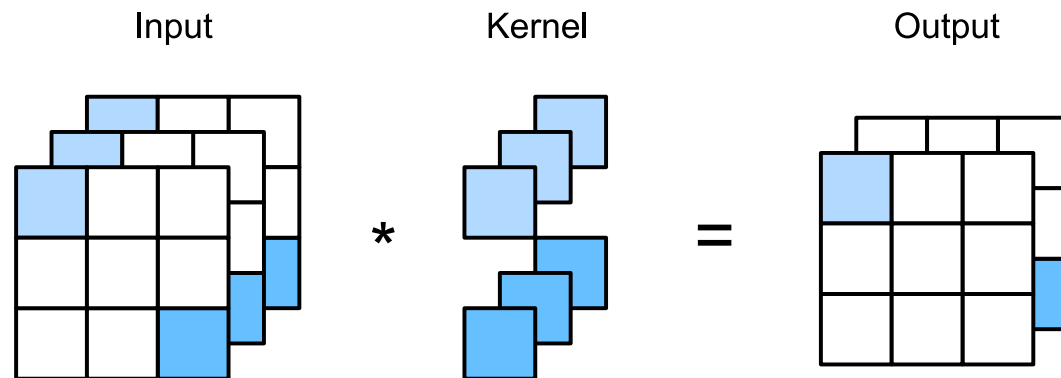
- $X: c_i \times n_h \times n_w$  input
- $W: c_i \times k_h \times k_w$  kernel
- $Y: m_h \times m_w$  output

$$Y = \sum_{i=0}^{c_i} X_{i,:,:} \otimes W_{i,:,:}$$

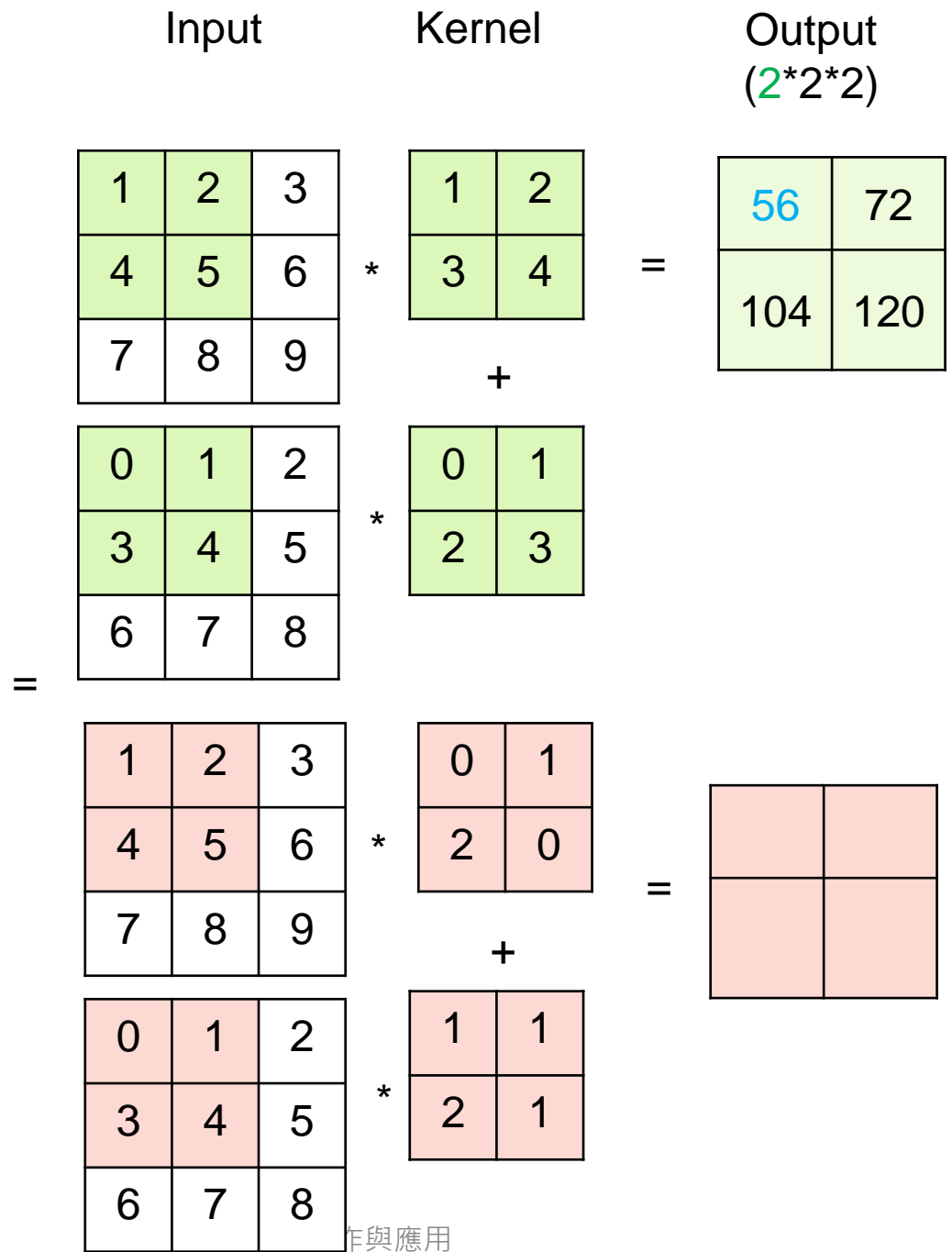
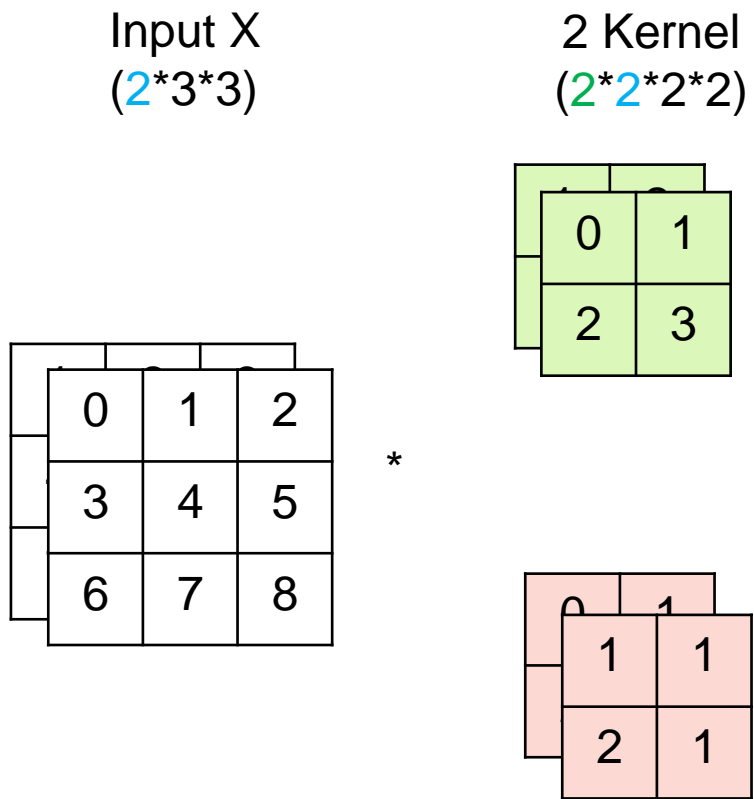


# Multiple Output Channels

- No matter how many inputs channels, so far we always get **single** output channel
- We can have **multiple** 3-D kernels, each one generates a output channel
- Input X:  $c_i \times n_h \times n_w$  ( $3 \times 3 \times 3$ )
- Output Y:  $c_o \times m_h \times m_w$  ( $2 \times 3 \times 3$ )
- Kernel W:  $c_o \times c_i \times k_h \times k_w$  ( $2 \times 1 \times 1 \times 1$ )







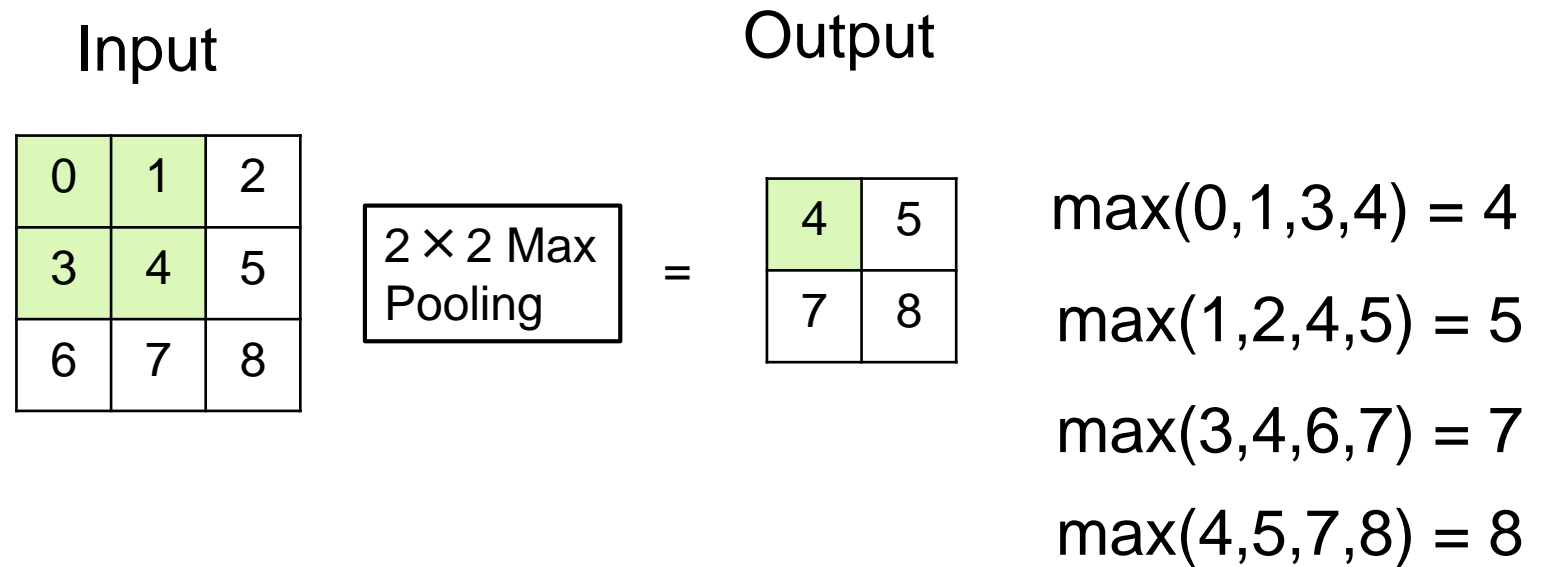
# Pooling

- Convolution is sensitive to position
- We need some degree of invariance to translation
  - Lighting, object positions, scales, appearance vary among images

# Pooling

- Max Pooling

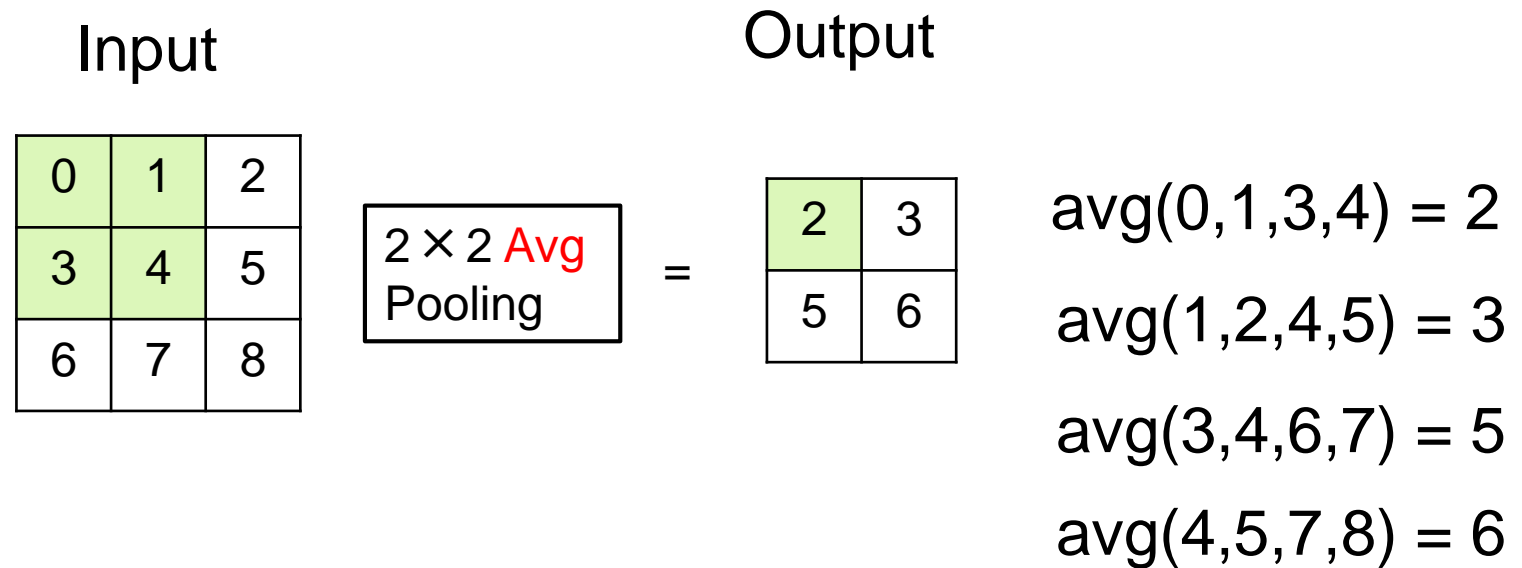
- Return the maximal value in the sliding windows
- The strongest pattern signal in a window



# Pooling

- Average Pooling

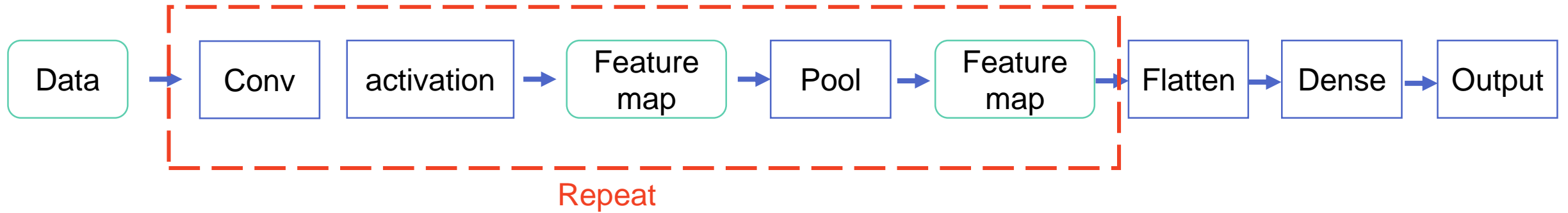
- Return the mean value in the sliding windows
- The average signal in a window



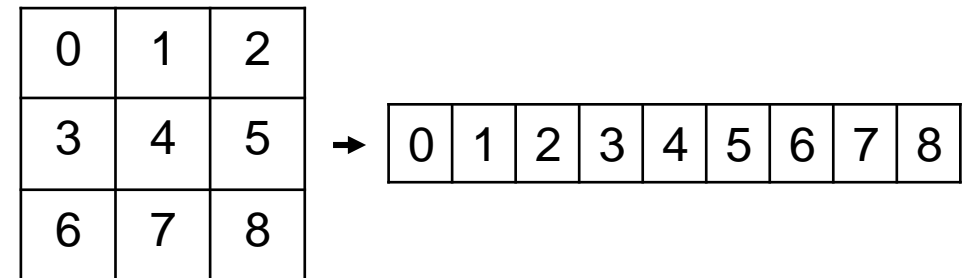
# Pooling

- Pooling layers have similar padding and stride as convolutional layers
- No learnable parameters
- Apply pooling for each input channel to obtain the corresponding output channel  
(每個channel獨立做pooling)
- # output channels = # input channels

# CNN loop and flatten





Flatten



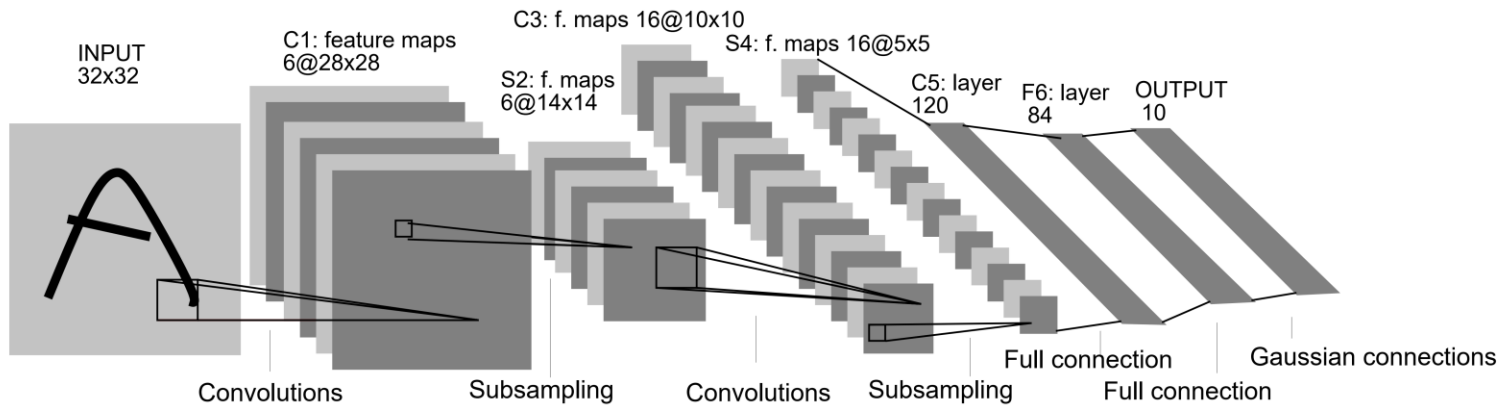
# Summary

- Convolutional layer
  - Reduced model capacity compared to dense layer
  - Efficient at detecting spatial patterns
  - Control output shape via padding, strides ( $o_h$  &  $o_w$ ) and (output) channels
- Max/Average Pooling layer
  - Provides some degree of invariance to translation

	MINST	ImageNet
		
Images	Gray image for hand-written digits	Color images with nature objects
Size	28 * 28	469 * 387
# examples (training)	60K	1.2M
# classes	10	1000

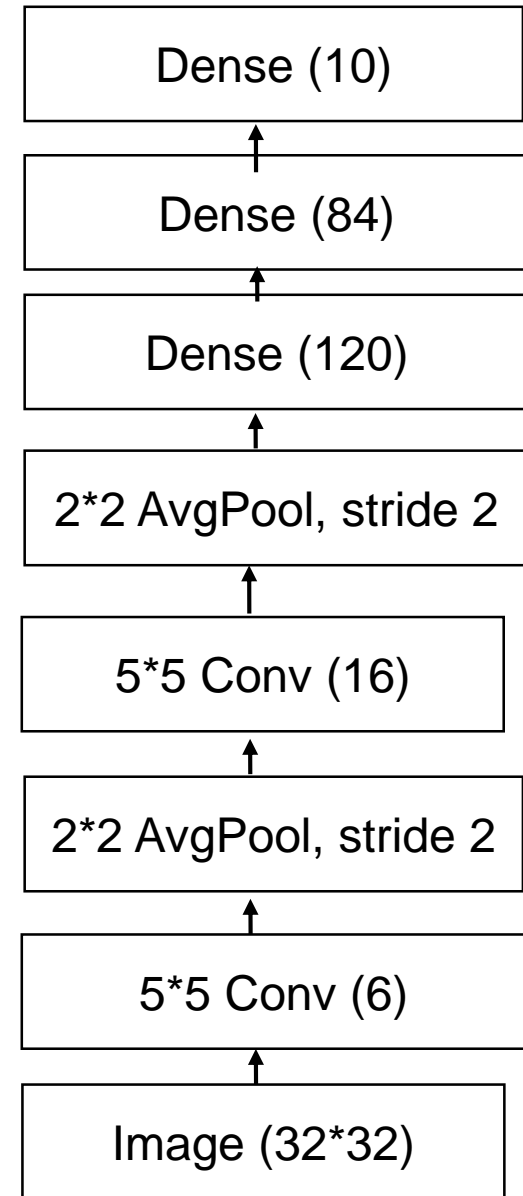


# LeNet



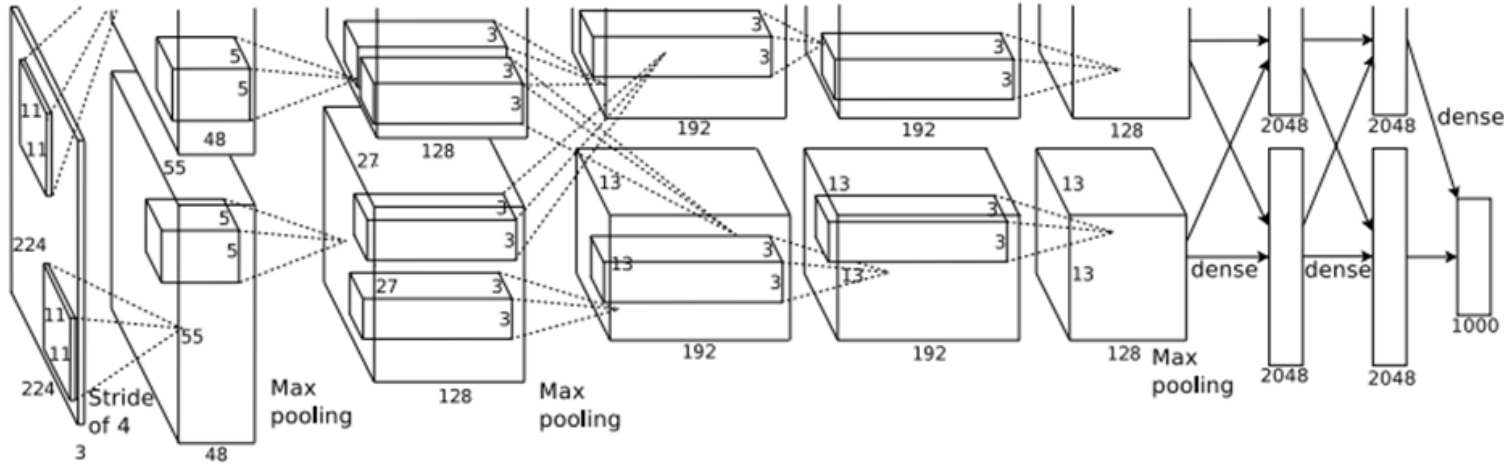
LeCun, Yann, Léon Bottou, Yoshua Bengio, and Patrick Haffner. "Gradient-based learning applied to document recognition." *Proceedings of the IEEE* 86, no. 11 (1998): 2278-2324.

Source: LeCun et al. (1998).

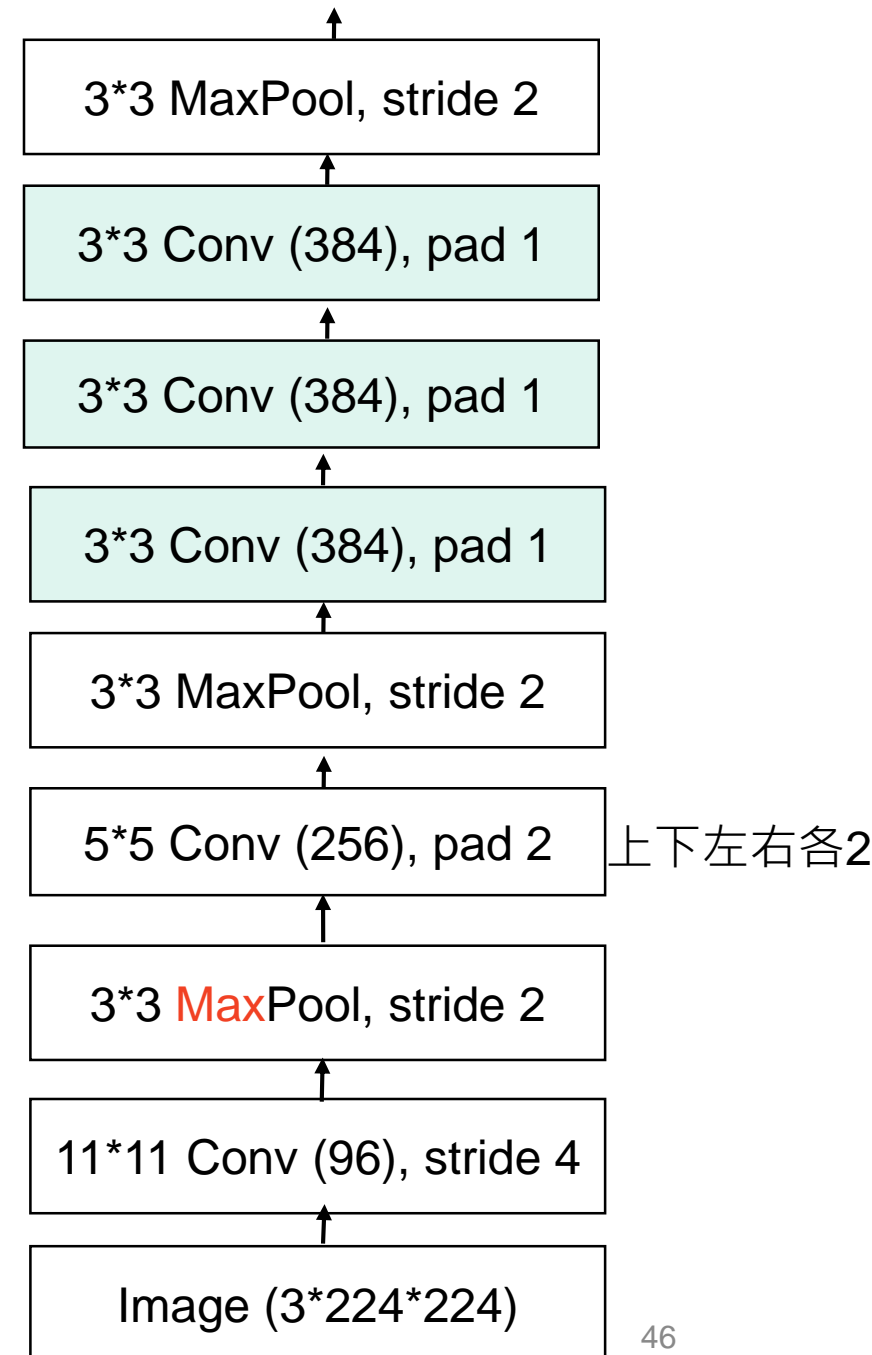


LeNet

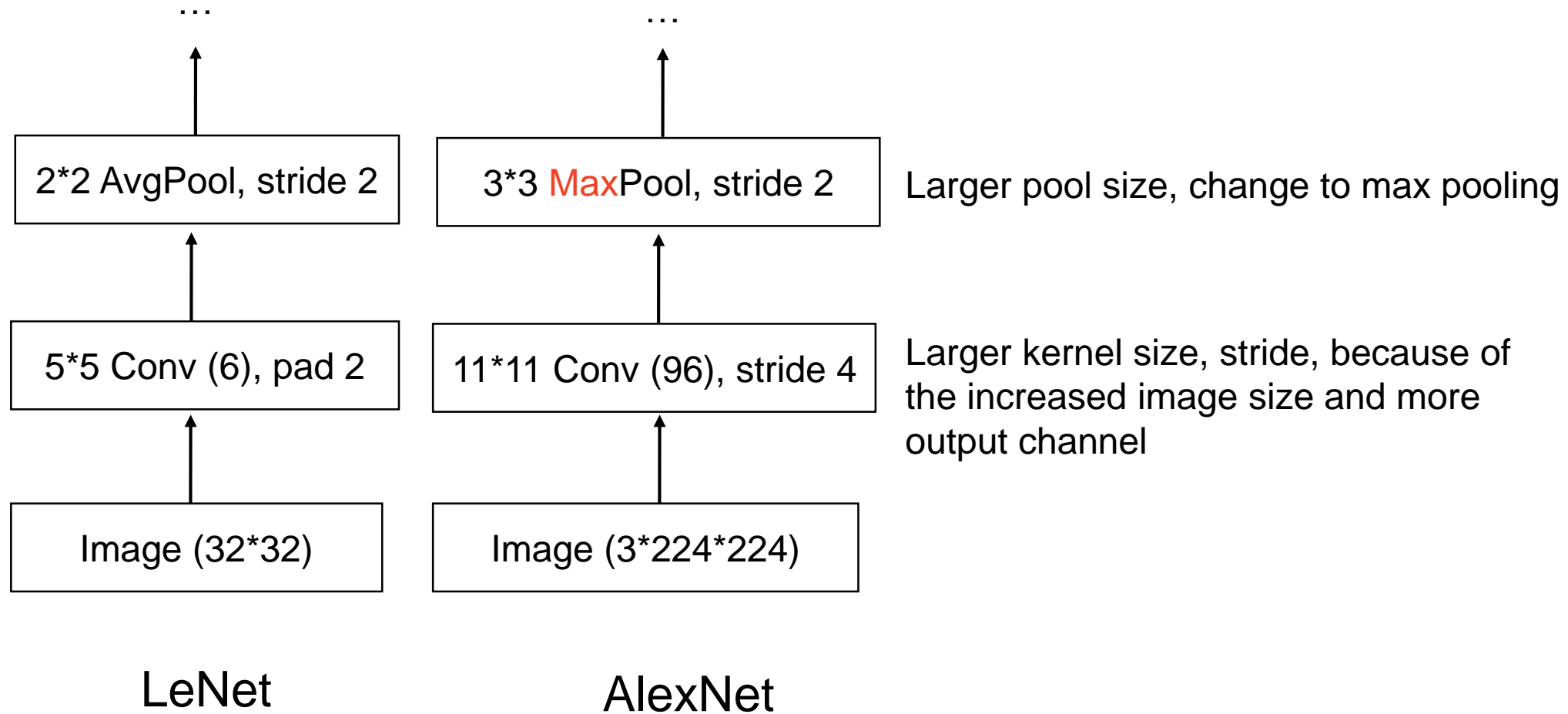
# AlexNet



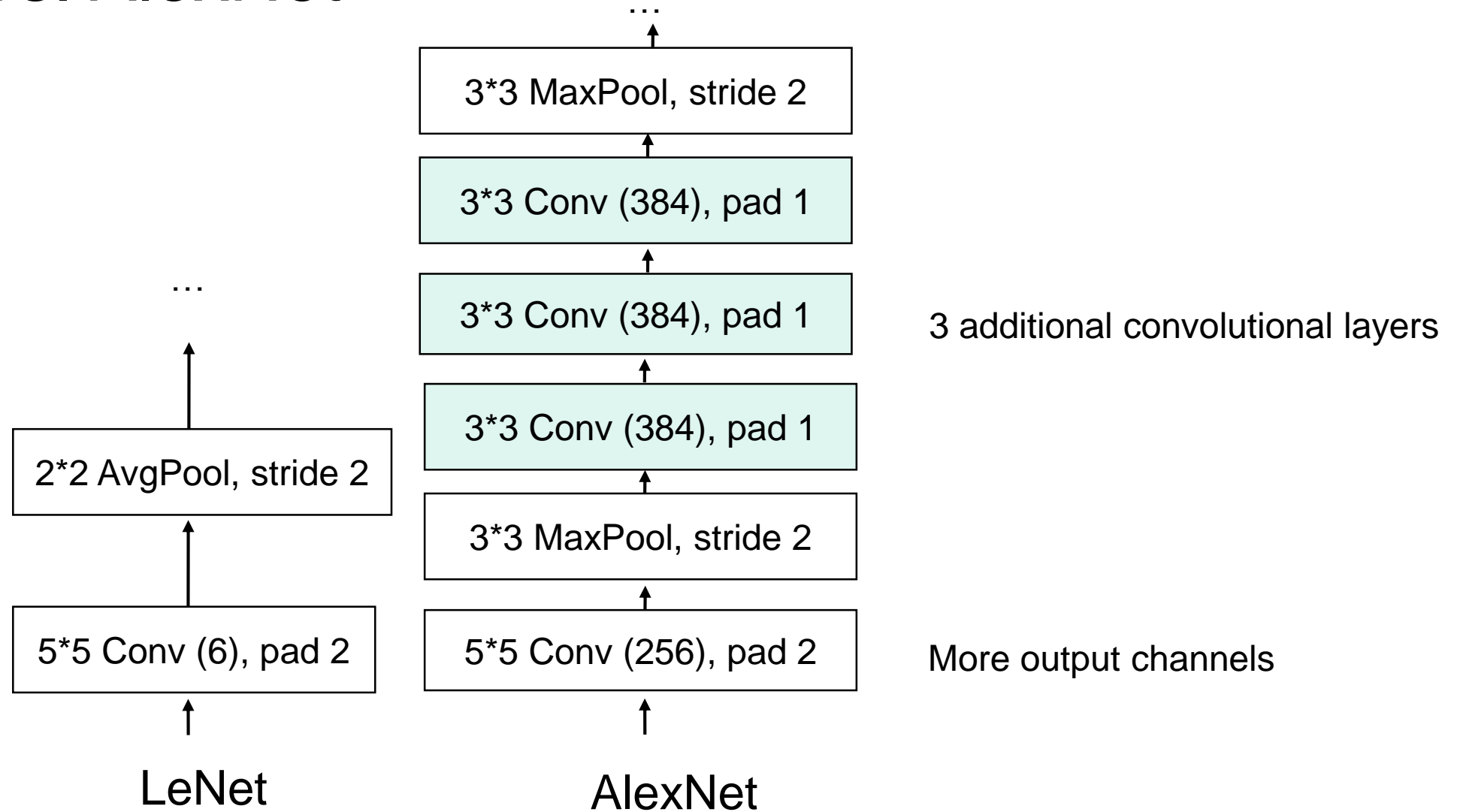
Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." *Advances in neural information processing systems* 25 (2012): 1097-1105.



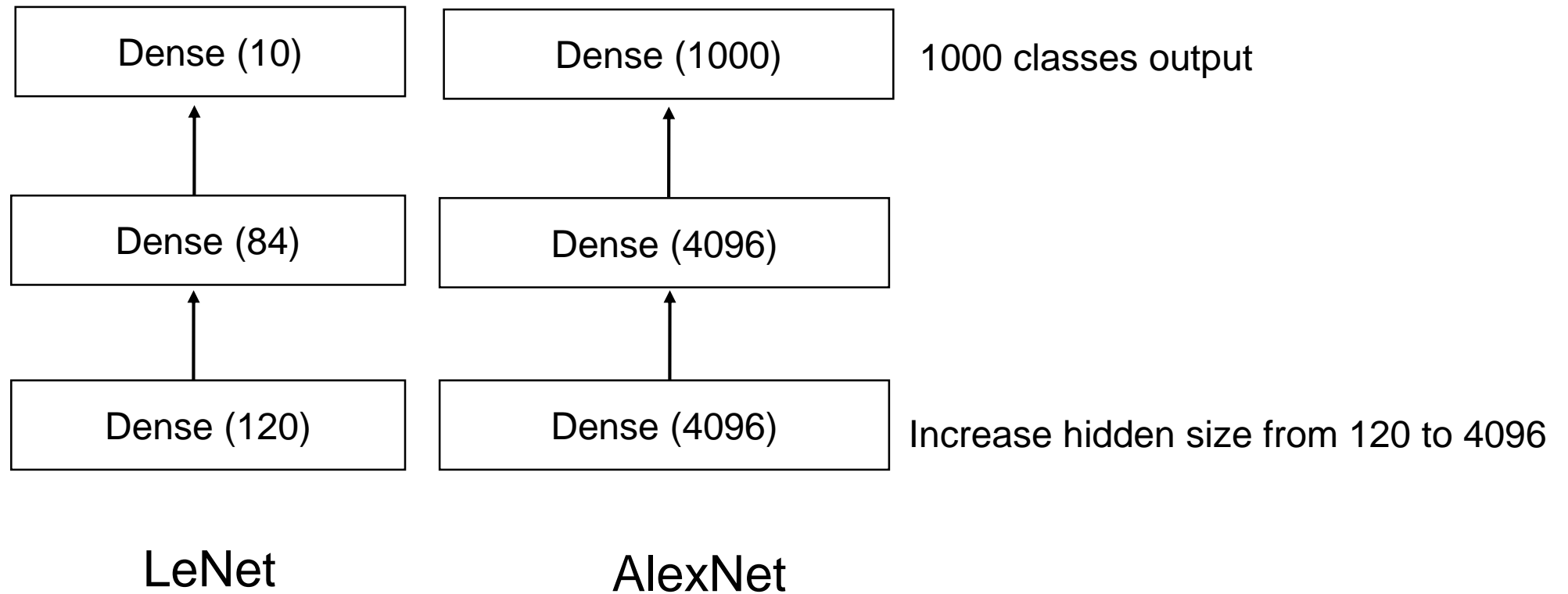
# LeNet vs. AlexNet



# LeNet vs. AlexNet



# LeNet vs. AlexNet



# AlexNet

- AlexNet won ImageNet competition in 2012
- Deeper and bigger LeNet
- Key modifications
  - Change activation from sigmoid to ReLU (no more vanish gradient)
  - Add a dropout layer after two hidden dense layers (better robustness / regularization)
  - MaxPooling

# 1D, 2D and 3D Convolution

- 1-D

$$y_{i,j} = \sum_{a=1}^h v_a x_{i+a}$$

- Text
- Time series
- Voice

- 2-D

$$y_{i,j} = \sum_{a=1}^h \sum_{b=1}^w v_{a,b} x_{i+a,j+b}$$

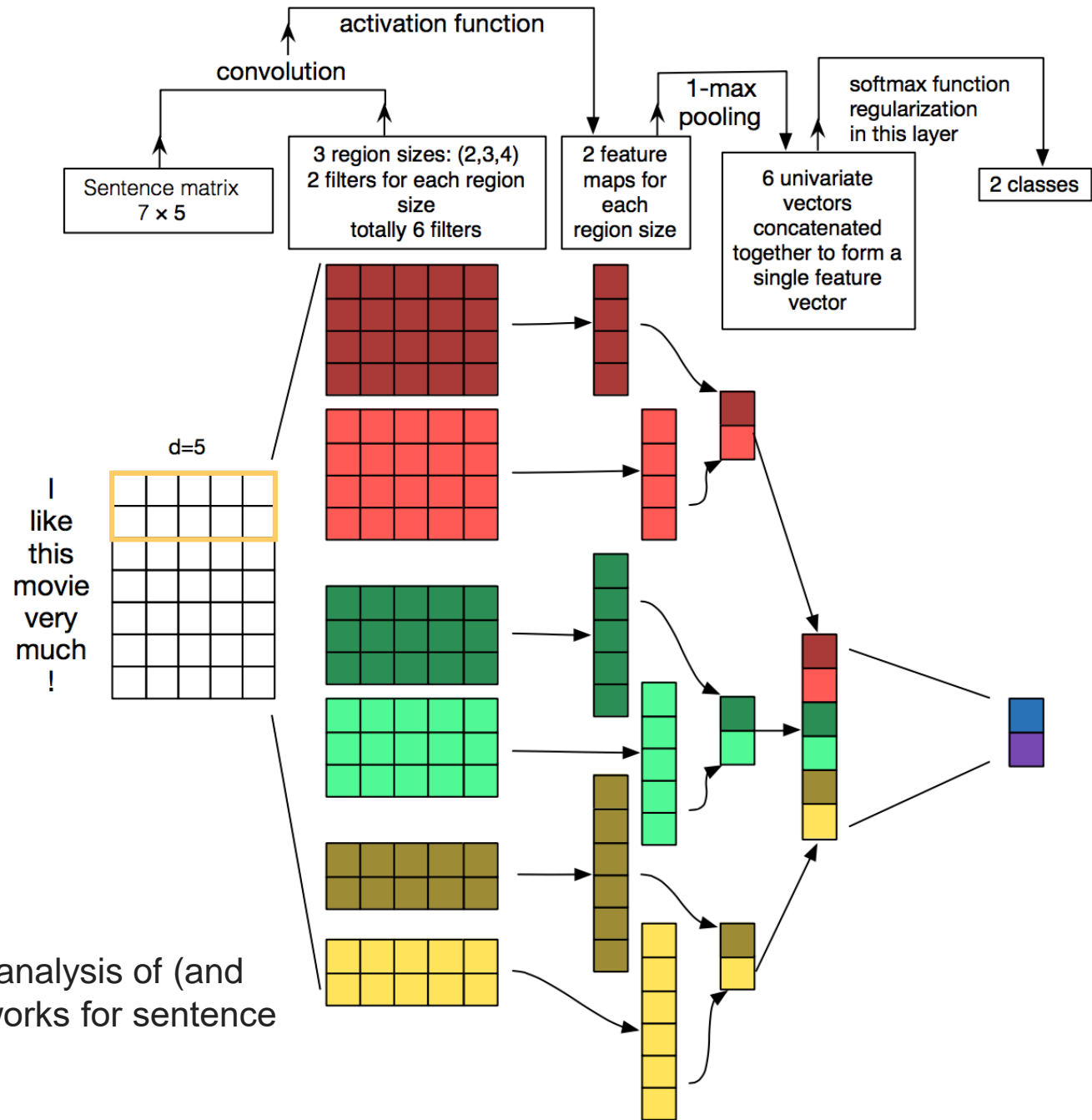
- Images
- Voice

- 3-D

$$y_{i,j,k} = \sum_{a=1}^h \sum_{b=1}^w \sum_{c=1}^d v_{a,b,c} x_{i+a,j+b,k+c}$$

- Video
- Medical images

# Text Classification



Zhang, Ye, and Byron C. Wallace. "A sensitivity analysis of (and practitioners' guide to) convolutional neural networks for sentence classification." IJCNLP2017.



# Problems (Internal Covariate Shift)

- Changes in model parameters during learning change the distributions of the outputs of each hidden layer.
  - Loss occurs at last layer
  - Data is at bottom layer
  - Bottom layers (weights) change – everything changes
- This means that later layers need to adapt to these (often noisy) changes during training.
- Can we avoid changing last layers while learning first layer?

# Recall input standardization

- Numerical inputs
  - Standardize: rescaling features to **zero mean** and **unit variance**

$$x = \frac{x - \mu}{\sigma}$$

- First, it proves convenient for optimization.
- Second, because we do not know *a priori* which features will be relevant, we do not want to penalize coefficients assigned to one feature more than on any other.

# Batch Normalization

$$x = \frac{x - \mu}{\sigma}$$

Input: Values of  $x$  over a mini-batch:  $B = \{x_1, \dots, x_m\}$

Learnable parameters:  $\gamma, \beta$

Output:  $\{y_i = BN_{\gamma, \beta}(x_i)\}$

// Fix mean and variance

$$\mu_B = \frac{1}{|B|} \sum_{i \in B} x_i, \quad \sigma_B^2 = \frac{1}{|B|} \sum_{i \in B} (x_i - \mu_B)^2$$

// normalize

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

// scale and shift

$$y_i = \gamma \hat{x}_i + \beta$$

# Batch normalization during testing

- Solution 1:
  - Mean and variance were calculated using the whole training data, rather than mini-batch.
- Solution 2:
  - The moving average of mean and variance of the batches were computed during training
  - $\text{running\_mean} = \text{momentum} * \text{running\_mean} + (1 - \text{momentum}) * \text{sample\_mean}$
  - $\text{running\_var} = \text{momentum} * \text{running\_var} + (1 - \text{momentum}) * \text{sample\_var}$

e.g. momentum = 0.1

# Batch Normalization 好處

- 解決 **Internal Covariate Shift** 的問題
- 有正則化的效果 (可以不使用Dropout)
- 減緩梯度消失
- 加速收斂

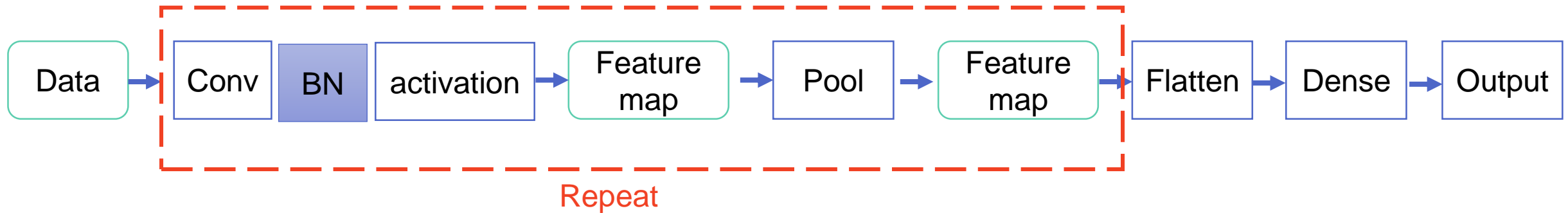
# Batch Normalization

- Fully connected layer (**before** activation function):

$$\mathbf{h} = \phi(\text{BN}(\mathbf{W}\mathbf{x} + \mathbf{b}))$$

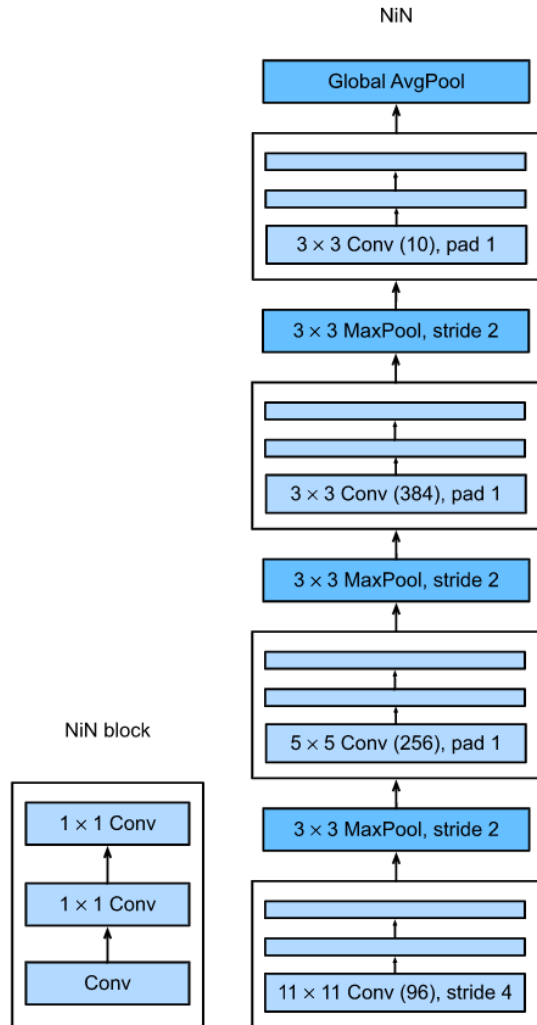
- CNN:
  - We can apply batch normalization **after** the convolution and **before** the nonlinear activation function.
  - When the convolution has multiple output channels, one normalization per channel, each channel has its own scale and shift parameters.
- Before: 速度快，因為BN可以merge到weight & bias
- After: 效果佳，真的讓output distribution被正規化

# Batch Normalization in CNN



- Assume that our minibatches contains  $m$  examples,
  - The output of convolution has height  $p$  and width  $q$
  - $m \cdot p \cdot q$  elements per channel are normalized.

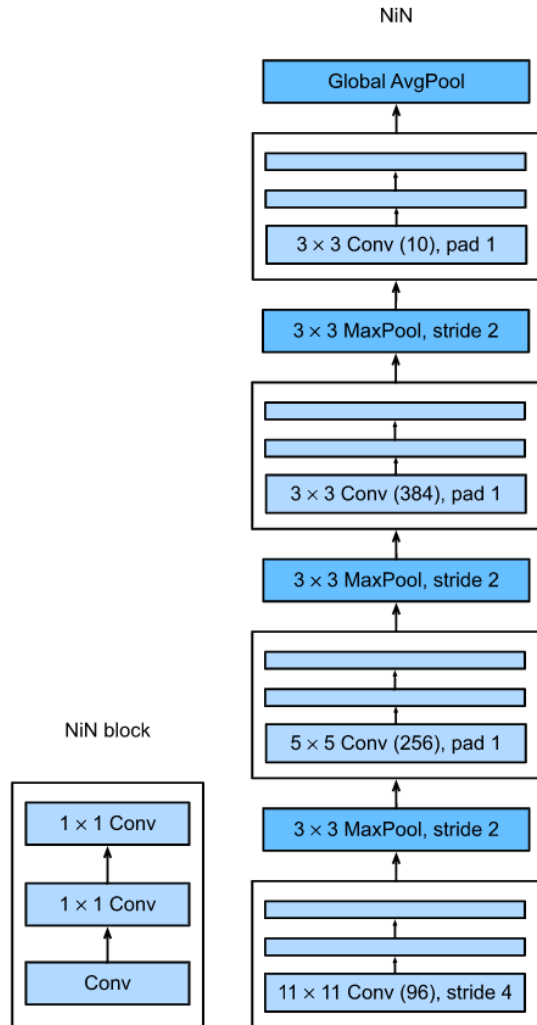
# Network in Network (NiN)



- The *network in network* (NiN) blocks ([Lin et al., 2013](#)):
  - (i) use  $1 \times 1$  convolutions to add local nonlinearities across the channel activations
  - (ii) use global average pooling to integrate across all locations in the last representation layer.
    - Note that global average pooling would not be effective, were it not for the added nonlinearities.

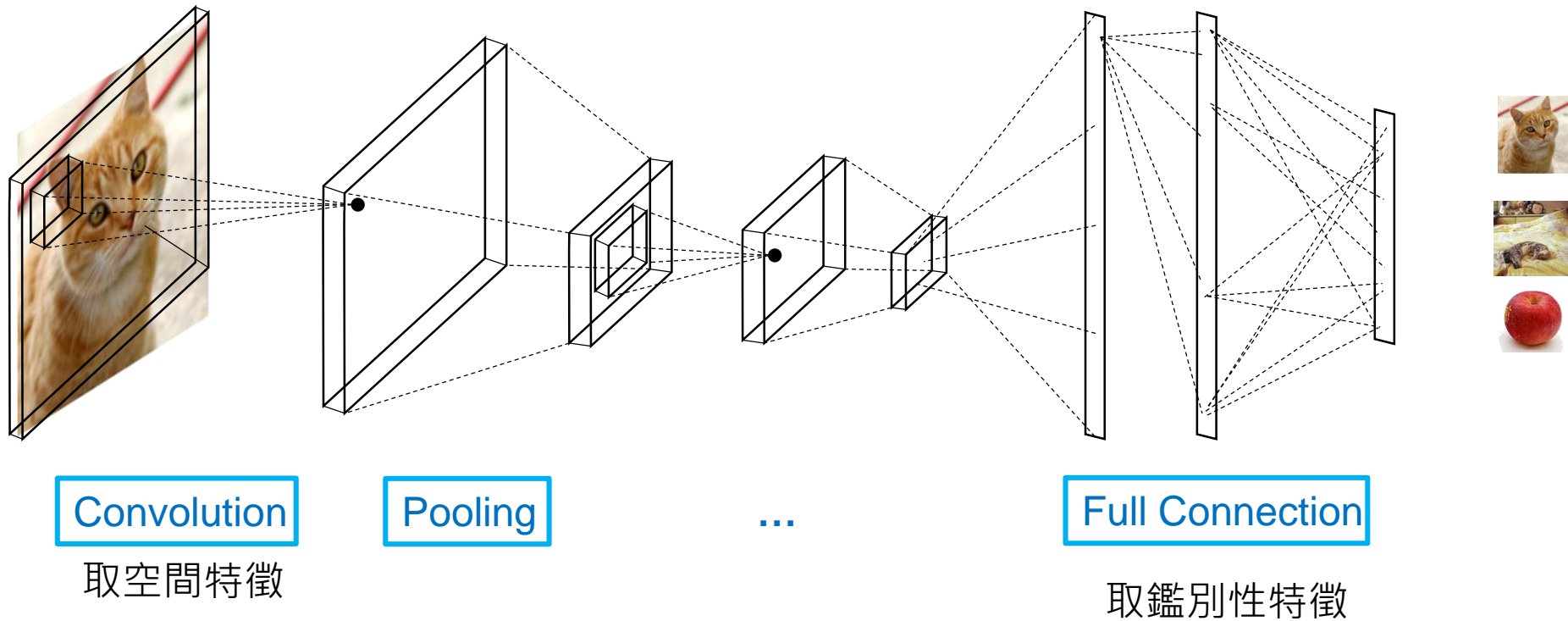


# Network in Network (NiN)

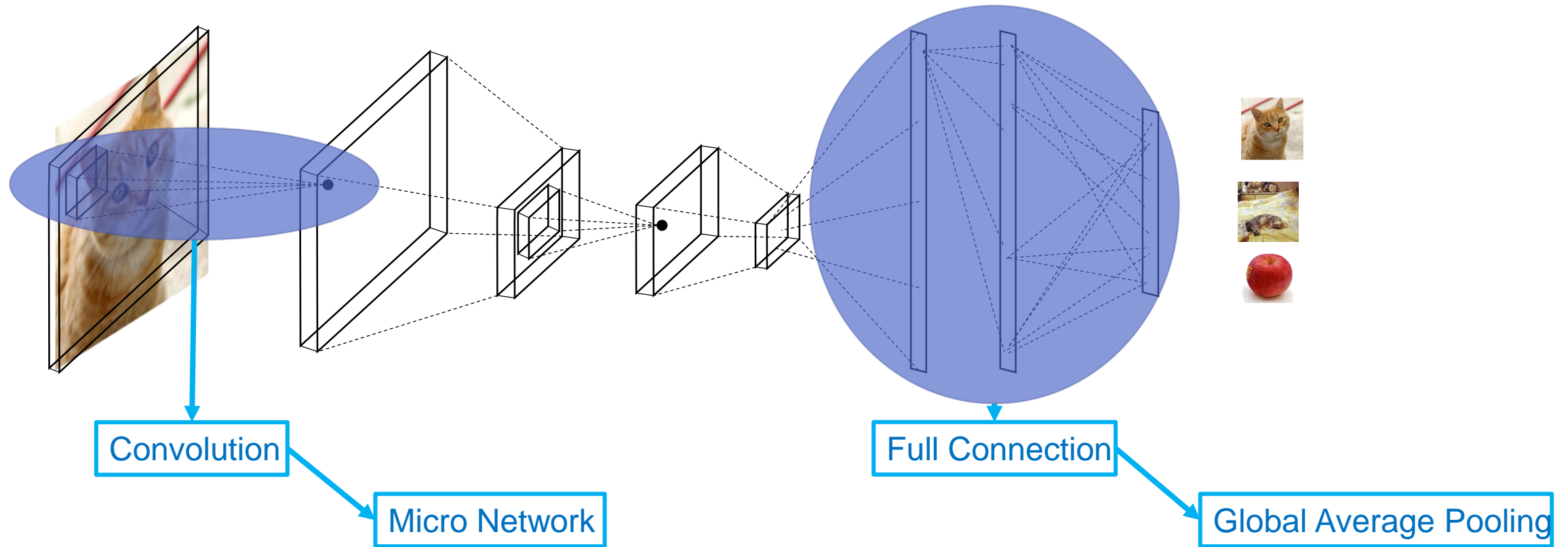


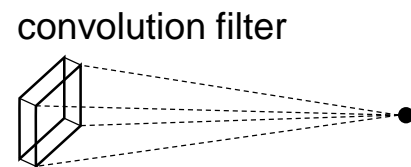
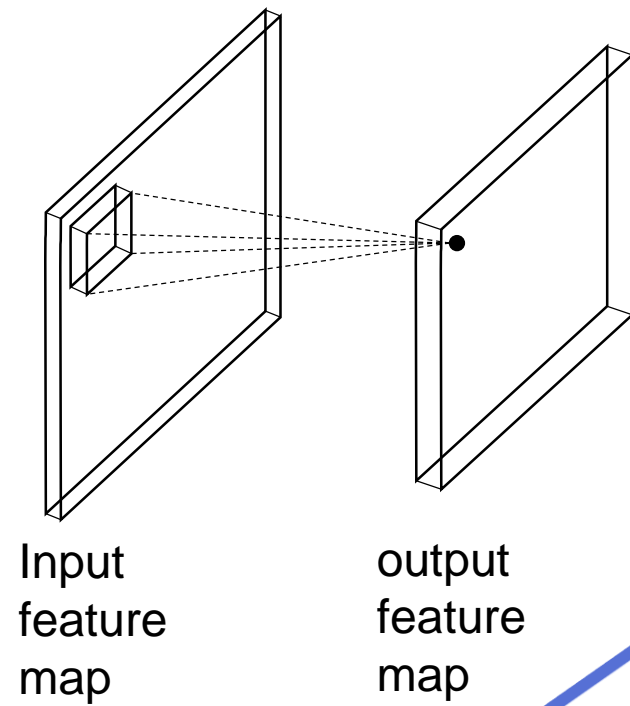
- NiN block:
  - The resulting  $1 \times 1$  convolution can be thought as a fully connected layer acting independently on each pixel location.

# Convolutional Neural Network



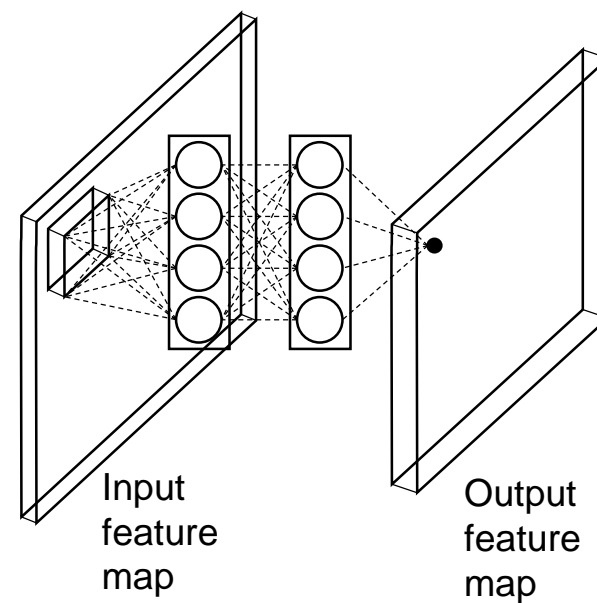
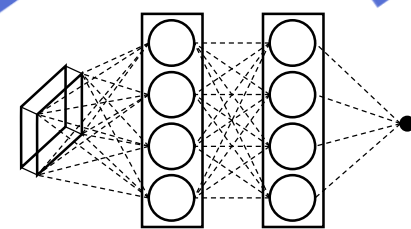
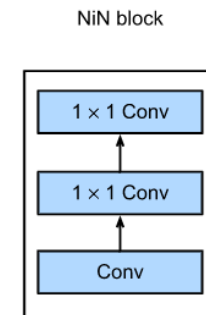
# Network In Network





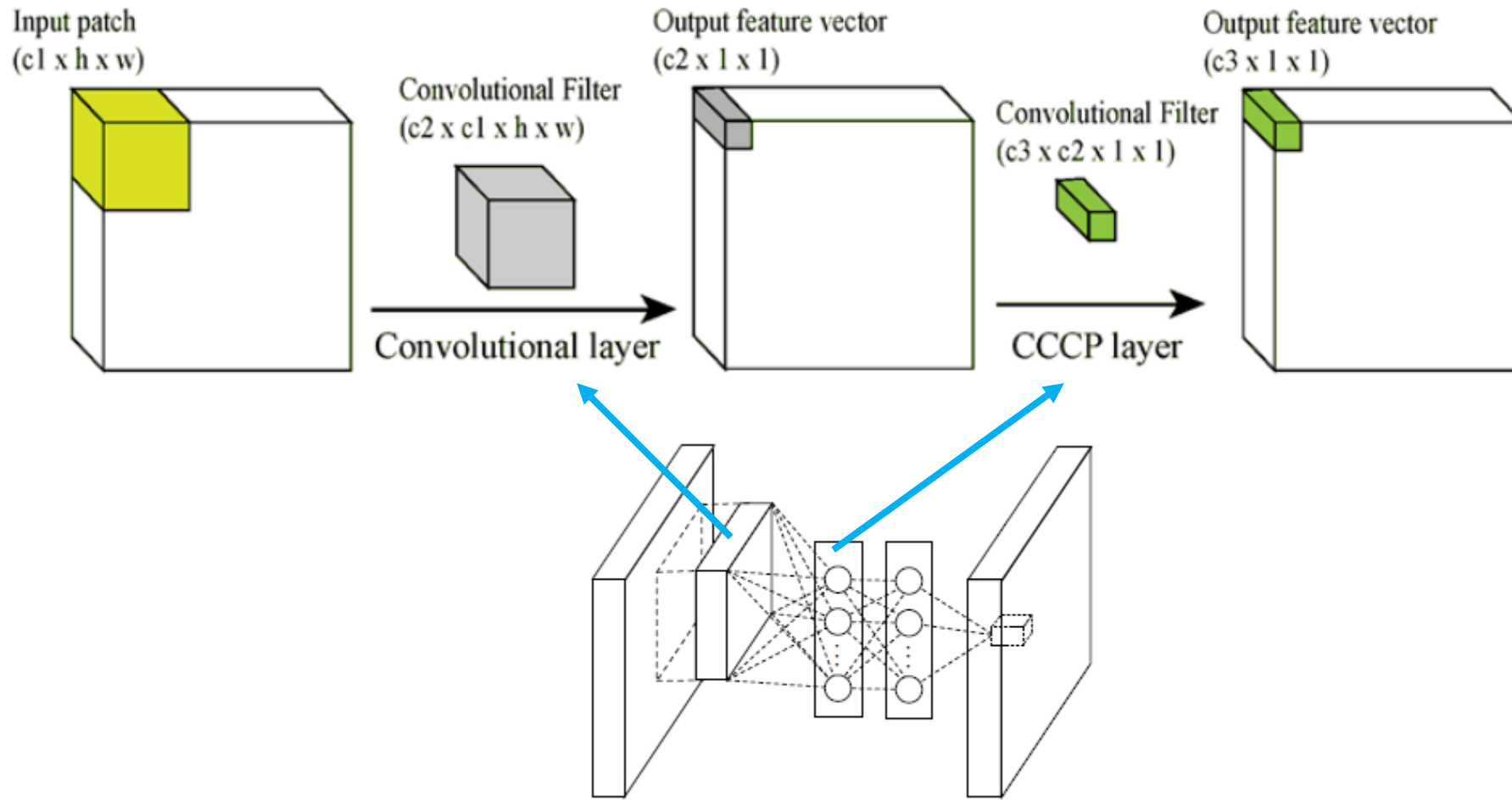
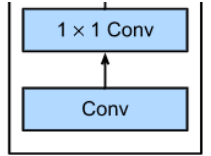
Convolution

Micro Network

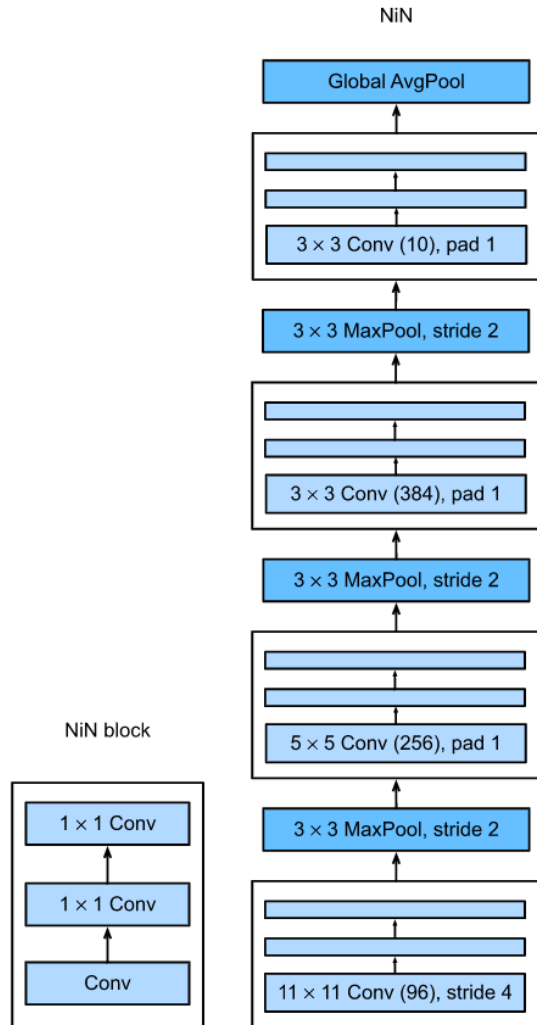


# Network In Network

- Cascade Cross Channel Parametric Pooling = 1 x 1 Convolution

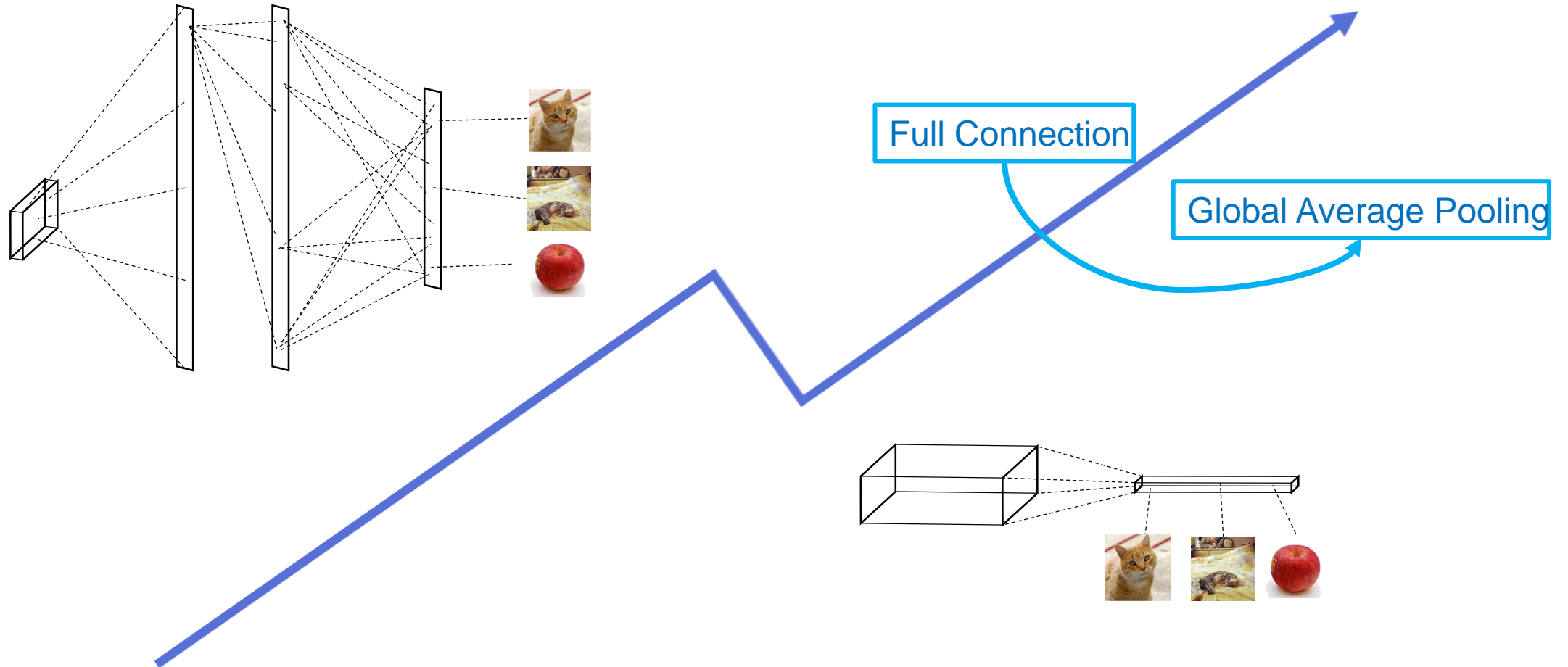


# Network in Network (NiN)

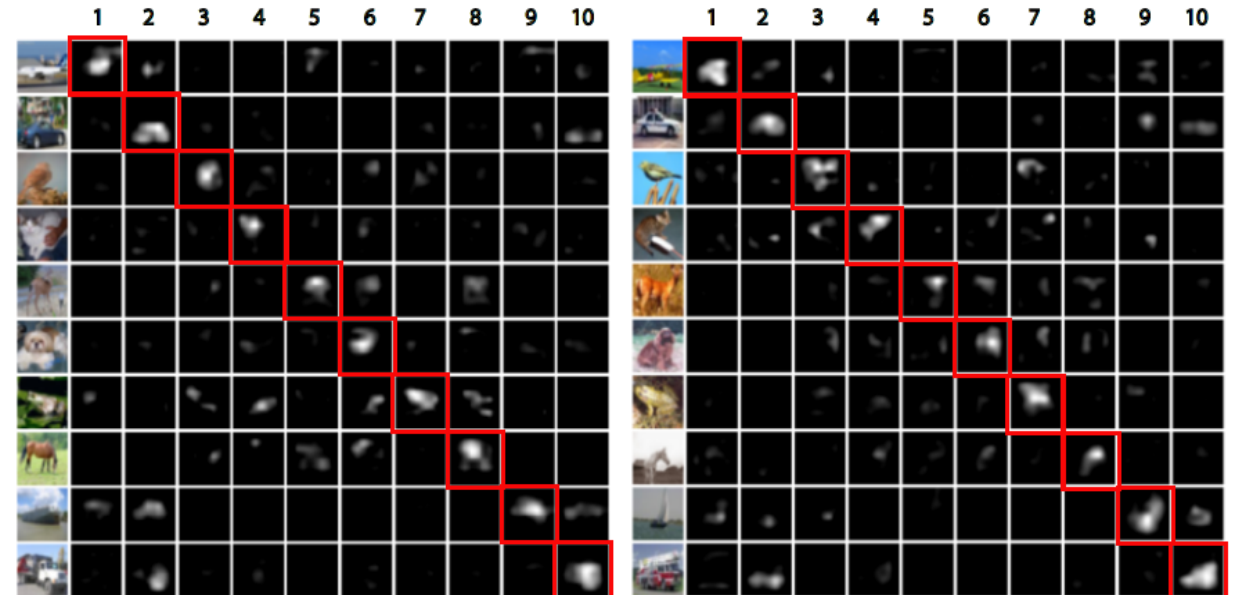
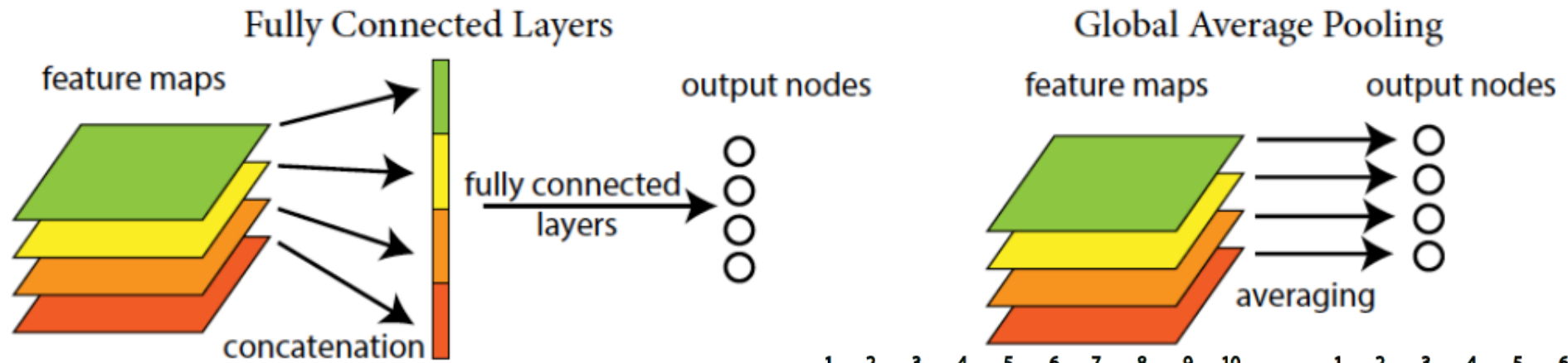


- Global AvgPool:
  - NiN avoids fully connected layers altogether.
  - Instead, NiN uses a NiN block with a number of output channels equal to the number of label classes, followed by a *global average pooling layer*, yielding a vector of logits.
  - This design significantly reduces the number of required model parameters, albeit at the expense of a potential increase in training time.

# Global Average Pooling



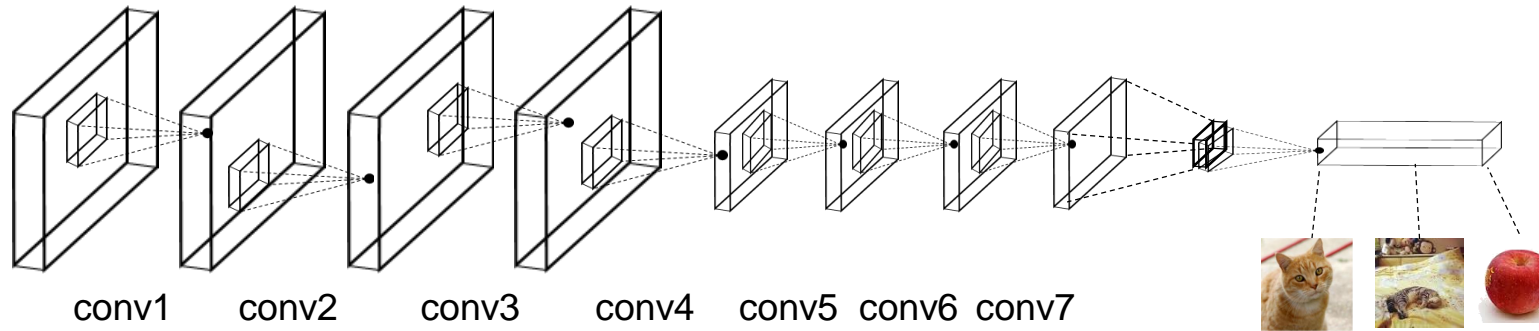
# Network In Network



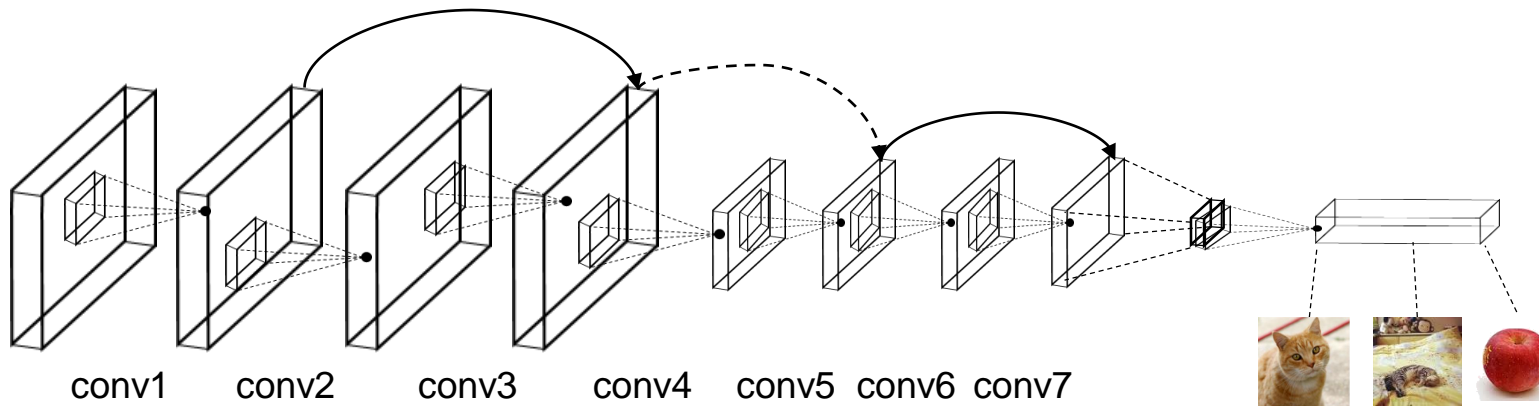


# What are Residual Networks

- Plain Networks: stacking of convolutional layers. (LeNet, AlexNet, VGG, ...)

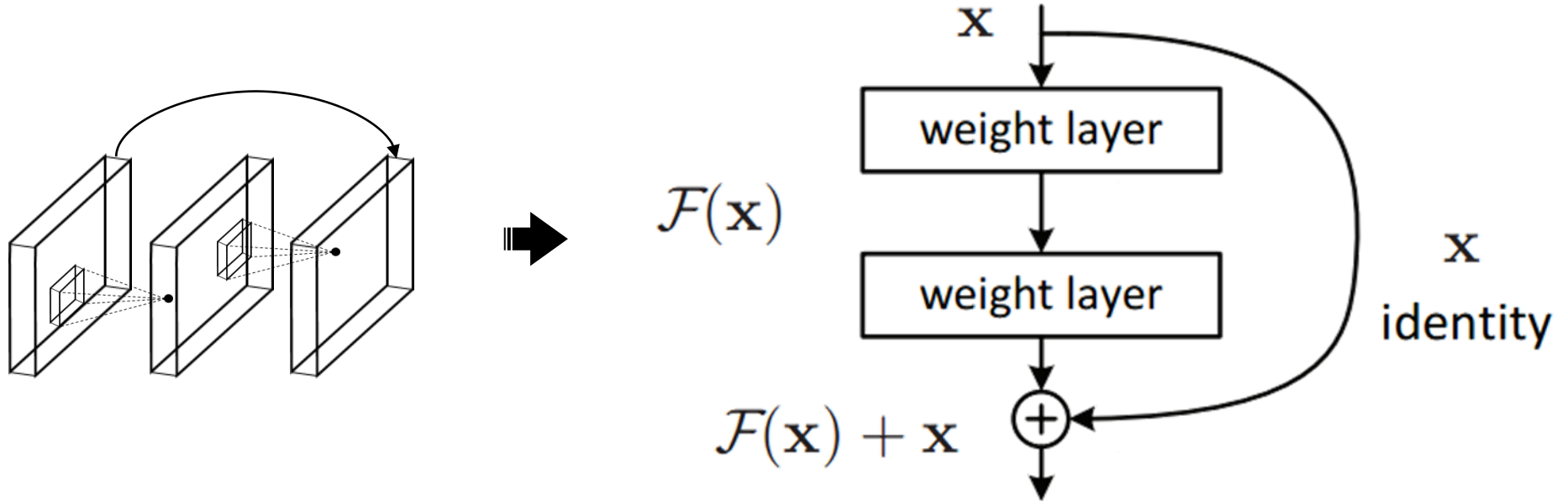


- Residual Networks: shortcut connections between feature maps.



# Residual Block of Residual Networks

- Residual Block



# Why Residual Networks work (3/5)

- Plain Networks are just like...

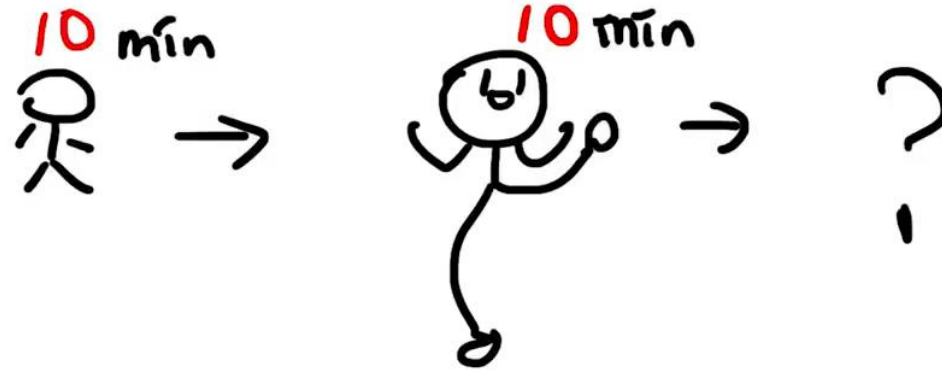


# Why Residual Networks work (4/5)

- Residual Networks know what following layers know.

<https://www.youtube.com/watch?v=mIHdYlwPZFA>

A 繪師 → B 繪師 → C 繪師



規則的部分呢跟上一次一樣喔

# Why Residual Networks work (5/5)

- Plain Networks vs Residual Networks



## Plain Networks

Each layer need to learn all things  
by information from next layer

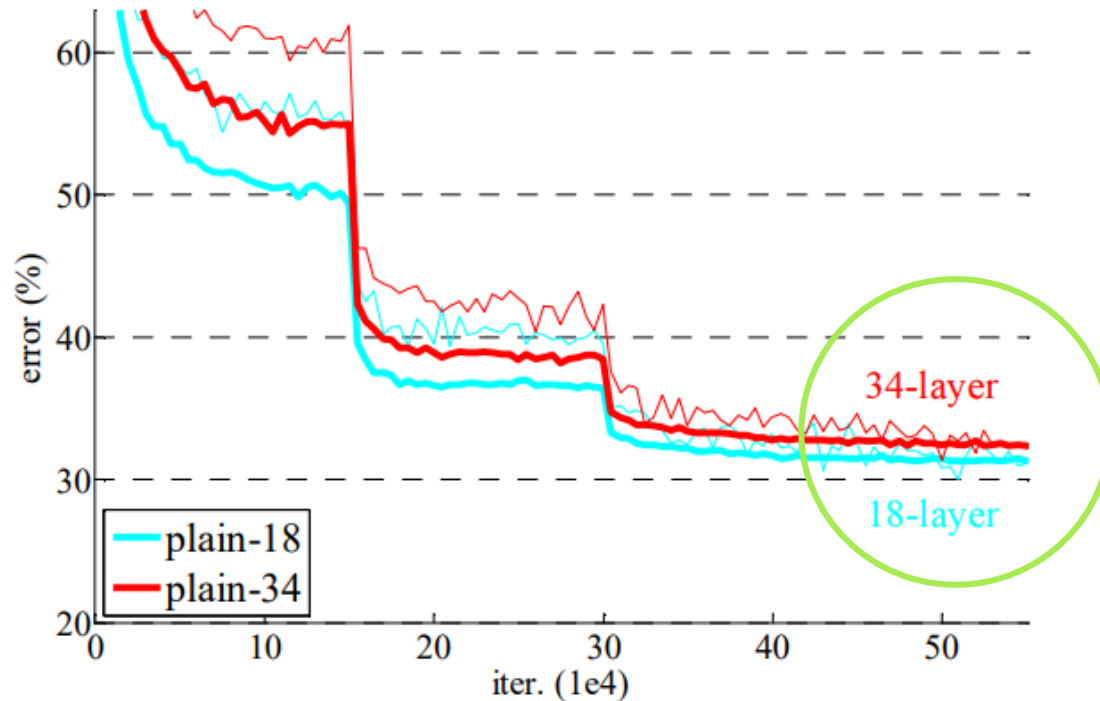


## Residual Networks

Each layer only need to learn things  
that other layers not yet learned

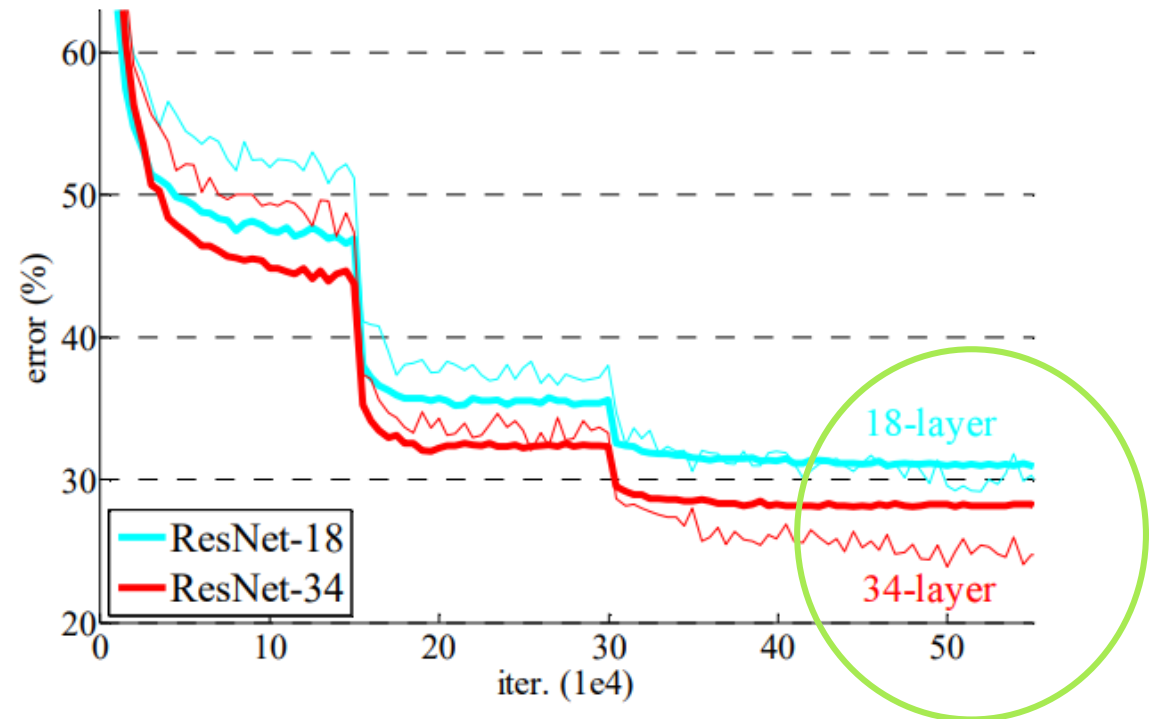
# Why Residual Networks are important

- Residual networks make **super deep networks learn well**.



**Plain Networks**

#layers ↑, Accuracy ↓

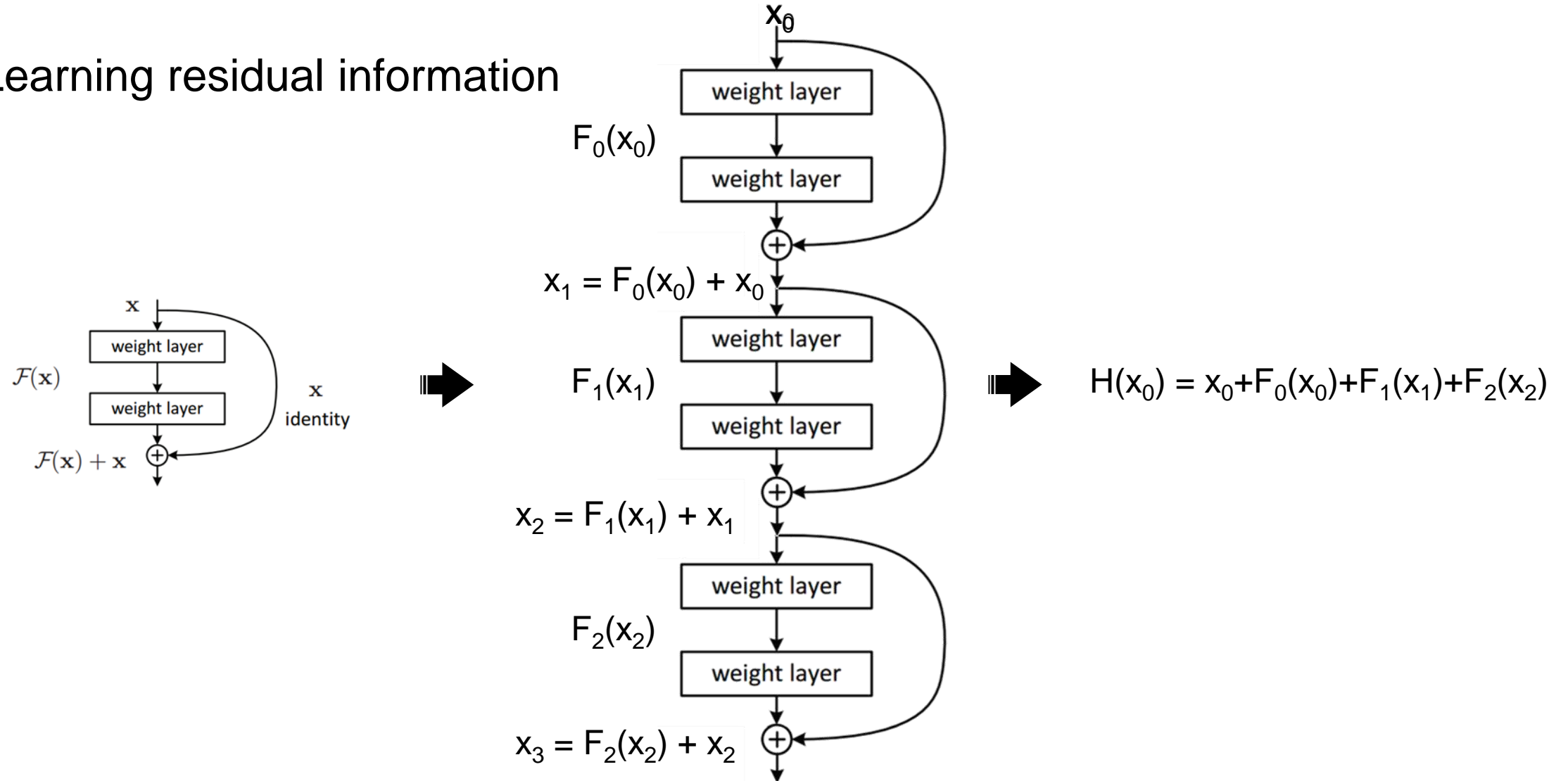


**Residual Networks**

#layers ↑, Accuracy ↑

# Why Residual Networks work (1/5)

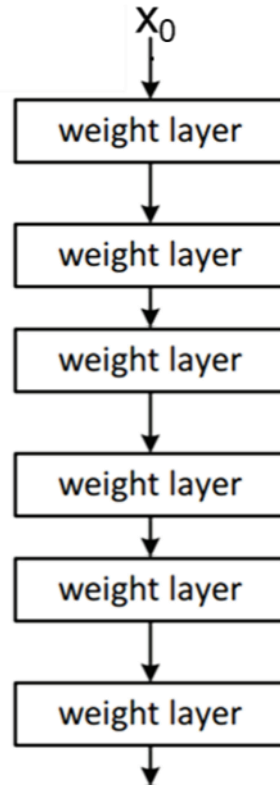
- Learning residual information





# Why Residual Networks work (2/5)

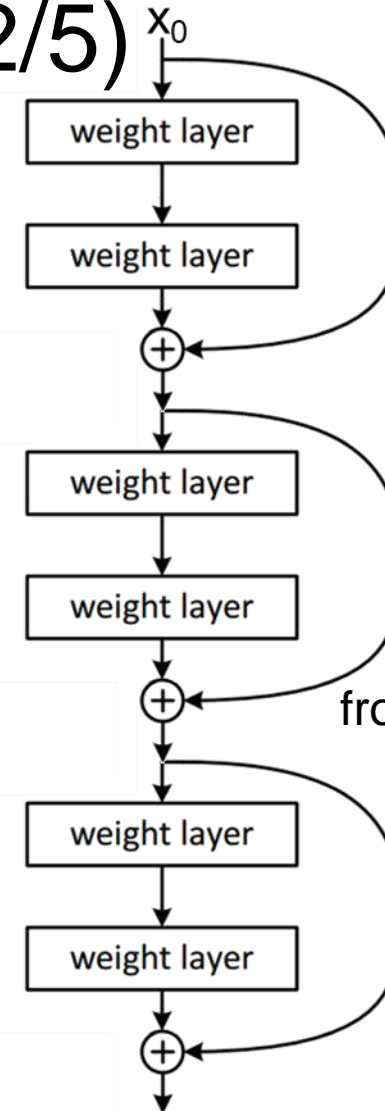
- Plain Networks vs Residual Networks



**Plain Networks**

$$H(x_0) = F_2(F_1(F_0(x_0)))$$

Can only learn information  
from next layer



**Residual Networks**

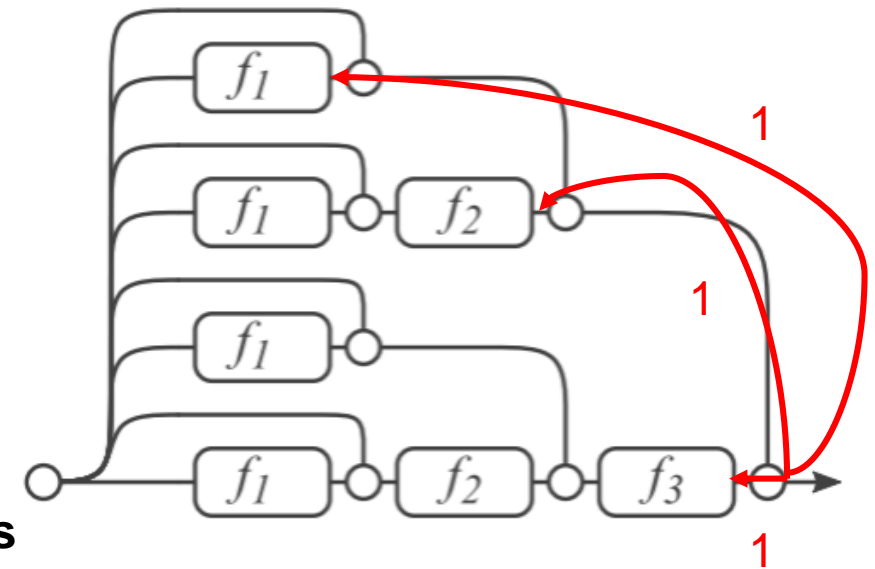
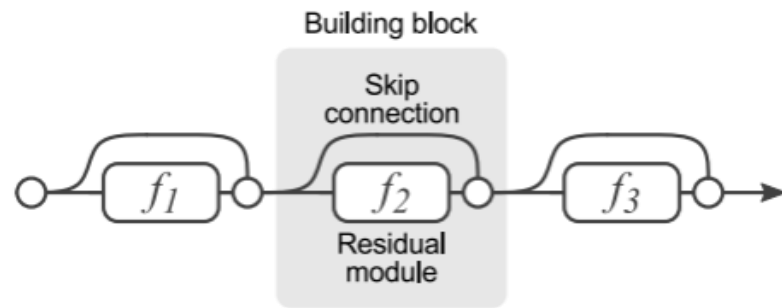
$$H(x_0) = x_0 + F_0(x_0) + F_1(x_1) + F_2(x_2)$$

Directly learn information  
from **label** and following layers

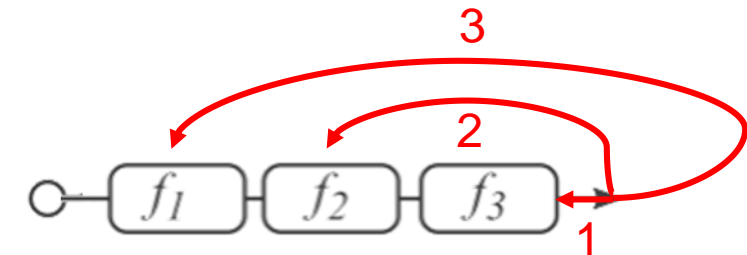
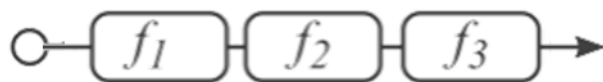


# Properties of Residual Networks

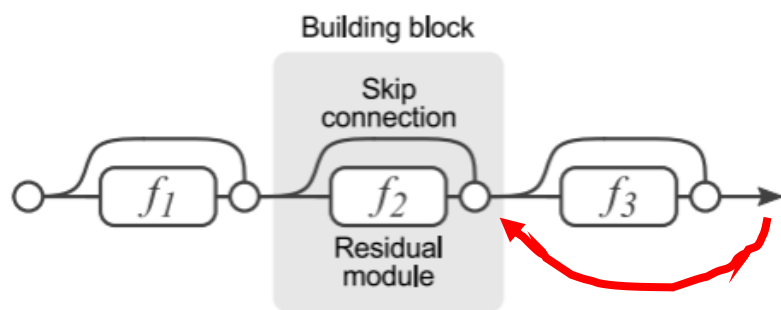
- Residual networks behave like ensembles of relatively shallow networks, it avoid the vanishing gradient problem (think about chain rule) by short paths.



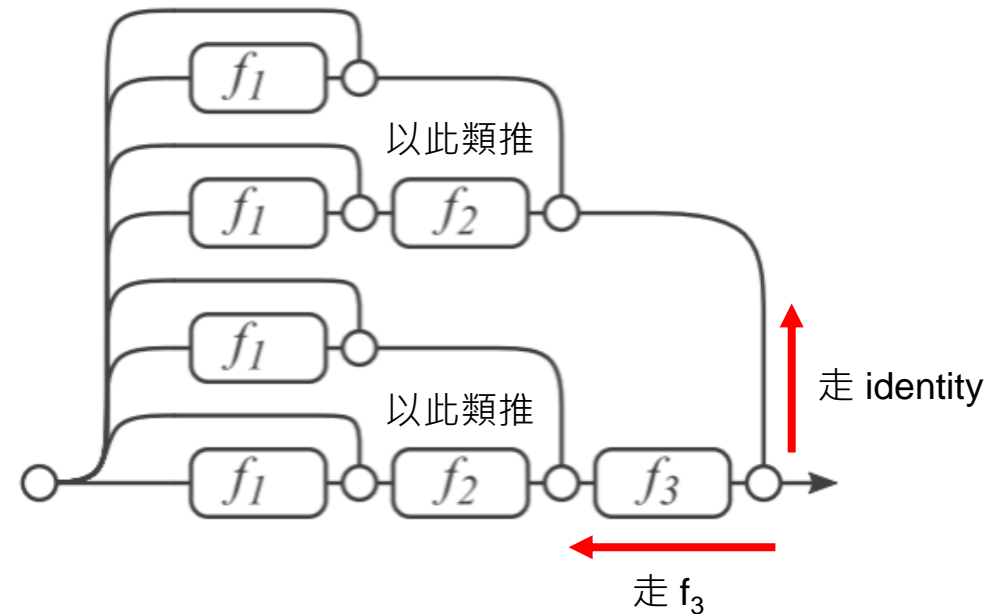
**Residual Networks**



**Plain Networks**



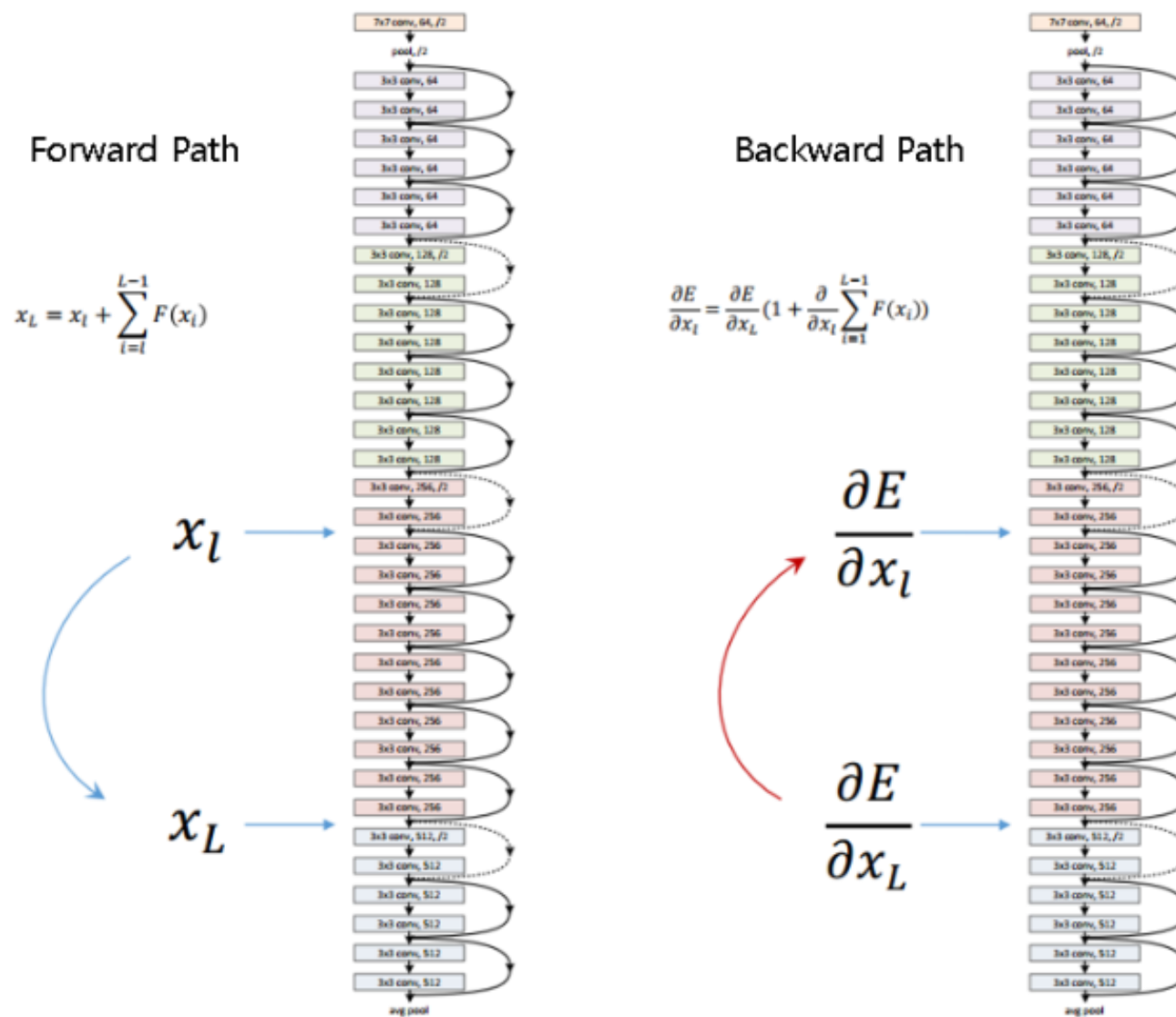
從一個  $\circ$  走到前一個  $\circ$  都有兩條路可以選擇 (1. 走identity; 2. 走  $f_i$ )。  
 所以三個residual layers就有  $2 \times 2 \times 2$  共 8 條路。



我們把 (1. 走identity; 2. 走  $f_i$ ) 展開，  
 就可以得到等效結構。

# Gradient computation of Residual Networks

- Simply add gradients from following layers.



# Assignment 2

- Will be announced on 3/29.
- Will be addressed next week.
- Due date will be 4/8.