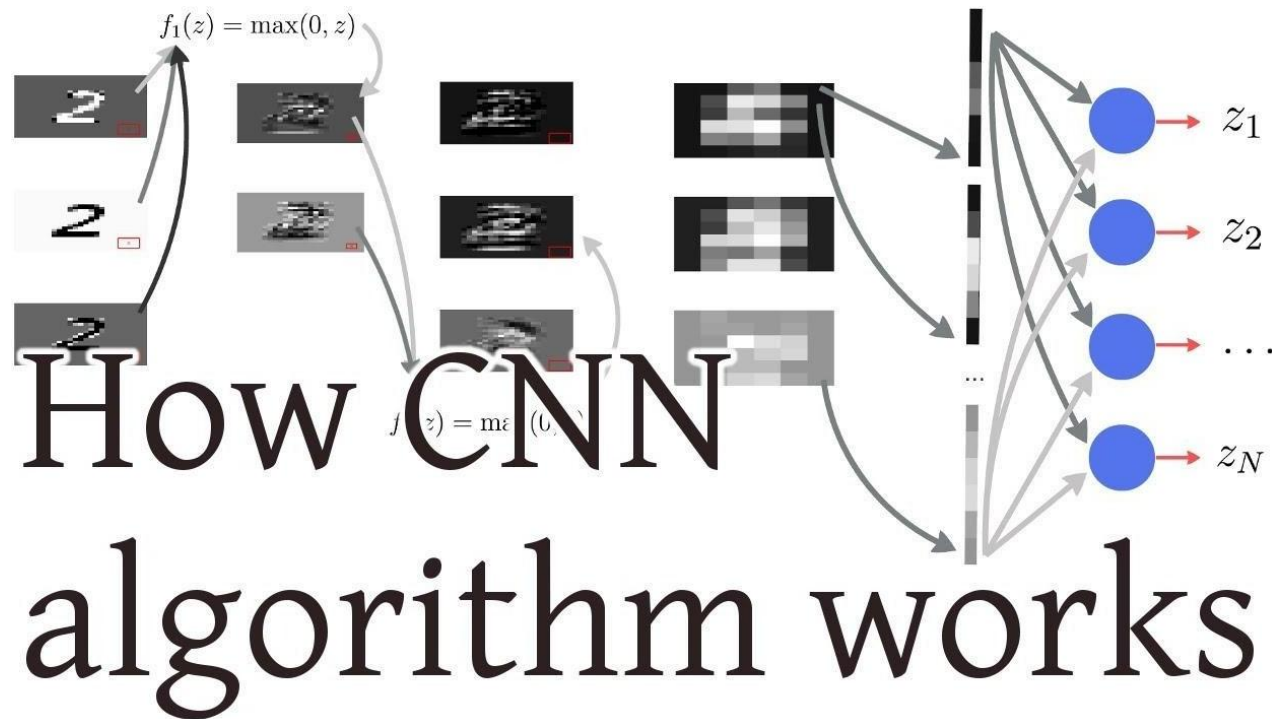


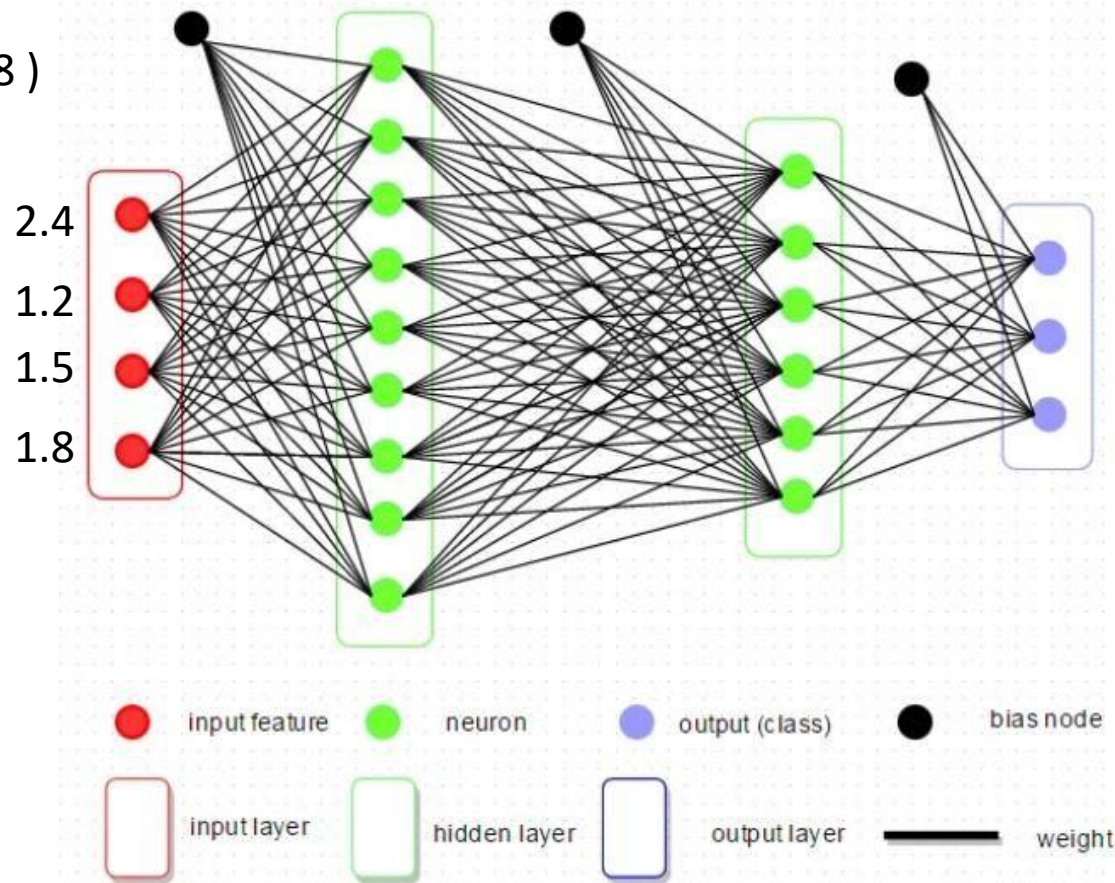
# Convolutional Neural Network (CNN)



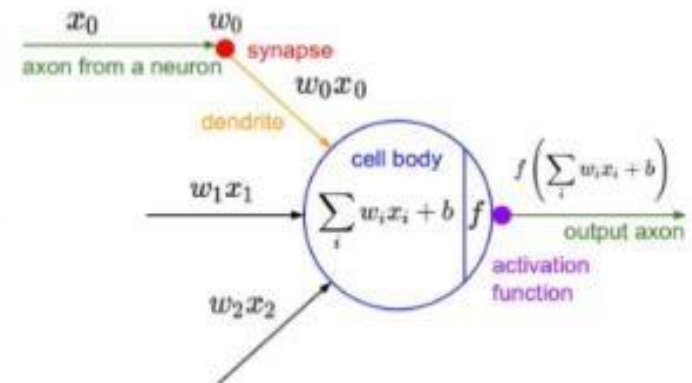
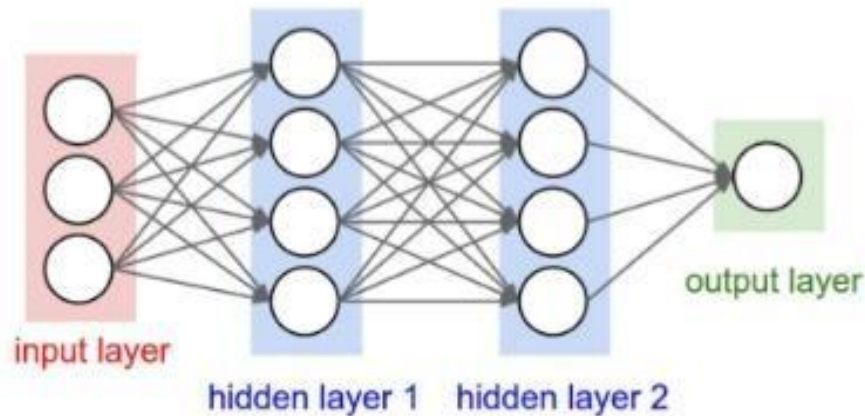
# [ 복습 ] Fully Connected Layer (전 결 합 계 층)

A 3-layers fully connected neural network (DNN)

Ex) Input vector  
( 2.4, 1.2, 1.5, 1.8 )



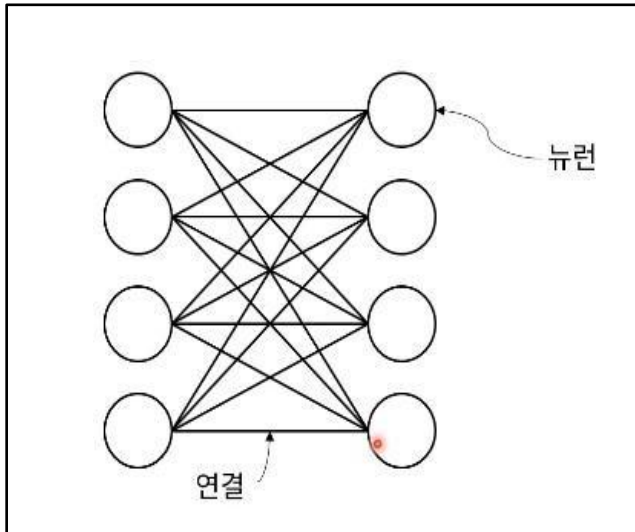
## Fully Connected Layers



# 1. What is Convolutional Layer?

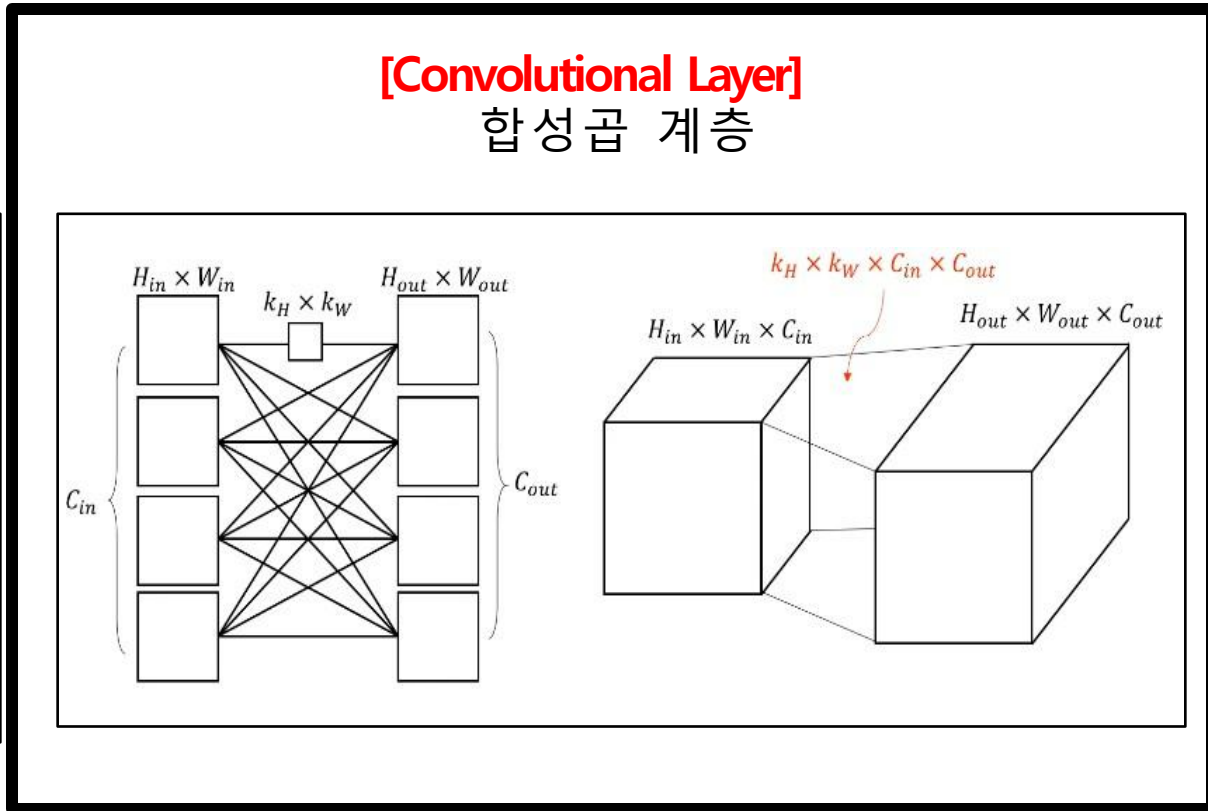
## [Fully Connected Layer]

전결합 계층



## [Convolutional Layer]

합성곱 계층

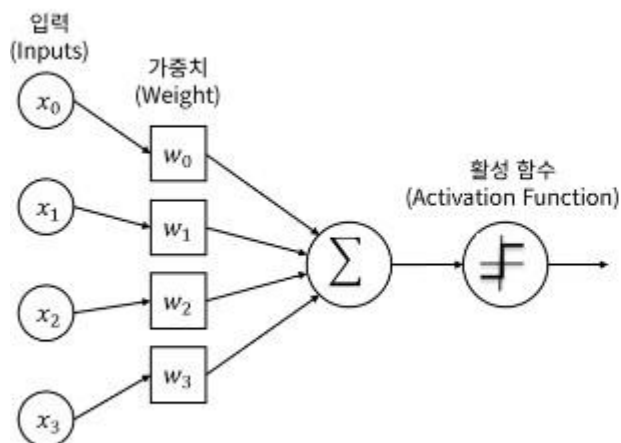


What's the difference?

# 1. What is Convolutional Layer?

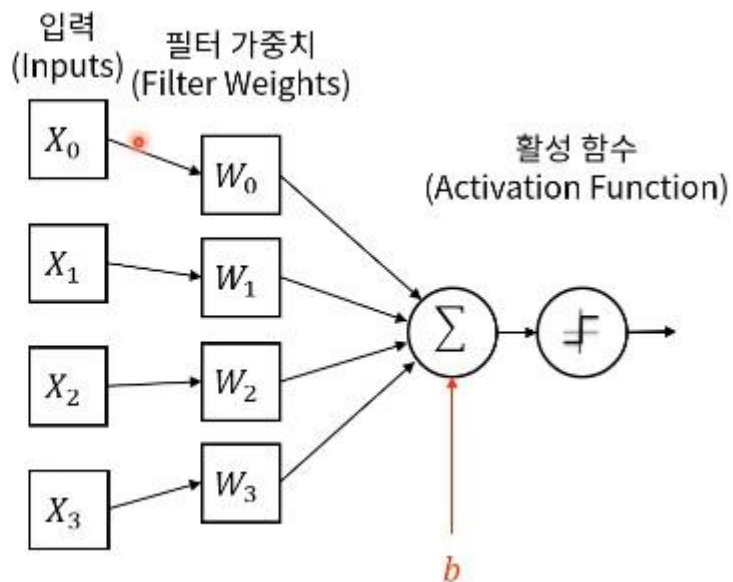
[Fully Connected Layer]

적결합 계층



[Convolutional Layer]

합성곱 계층



[ Input ] 입력으로 들어오는게 **사진/영상**

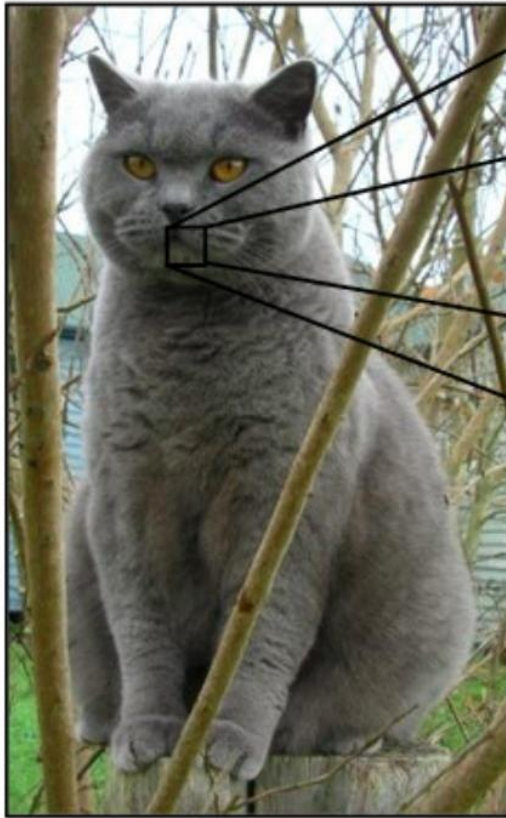
[ Weight ] 가중치가 **필터**라는 점

[ 연산 방식 ] 단순 곱셈이 아니라 **합성곱**

**Not so different!**

# 1. What is Convolutional Layer?

(1) input



08	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	87	08
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	45	04	56	62	00
81	49	31	73	55	79	14	29	93	71	40	67	58	88	30	03	49	13	36	65
52	70	95	23	04	60	11	42	60	24	68	56	01	32	56	71	37	02	36	91
22	31	16	71	51	67	83	89	41	92	36	54	22	40	40	28	66	33	13	80
24	47	35	80	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12	50
32	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64	70
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21
24	55	38	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63	72
21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31	33	95
78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57
86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04	89	55	40
04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33	27	98	66
55	56	68	87	57	62	20	72	03	46	33	67	46	55	12	32	63	93	53	69
04	42	16	73	35	39	11	24	94	72	18	08	46	29	32	40	62	76	36	
20	69	36	41	72	30	23	88	34	63	99	69	82	67	59	85	74	04	36	16
20	73	35	29	78	31	90	01	74	31	49	71	48	54	53	16	23	57	05	54
01	70	54	71	83	51	54	69	16	92	33	48	61	43	52	01	89	19	67	48

What the computer sees

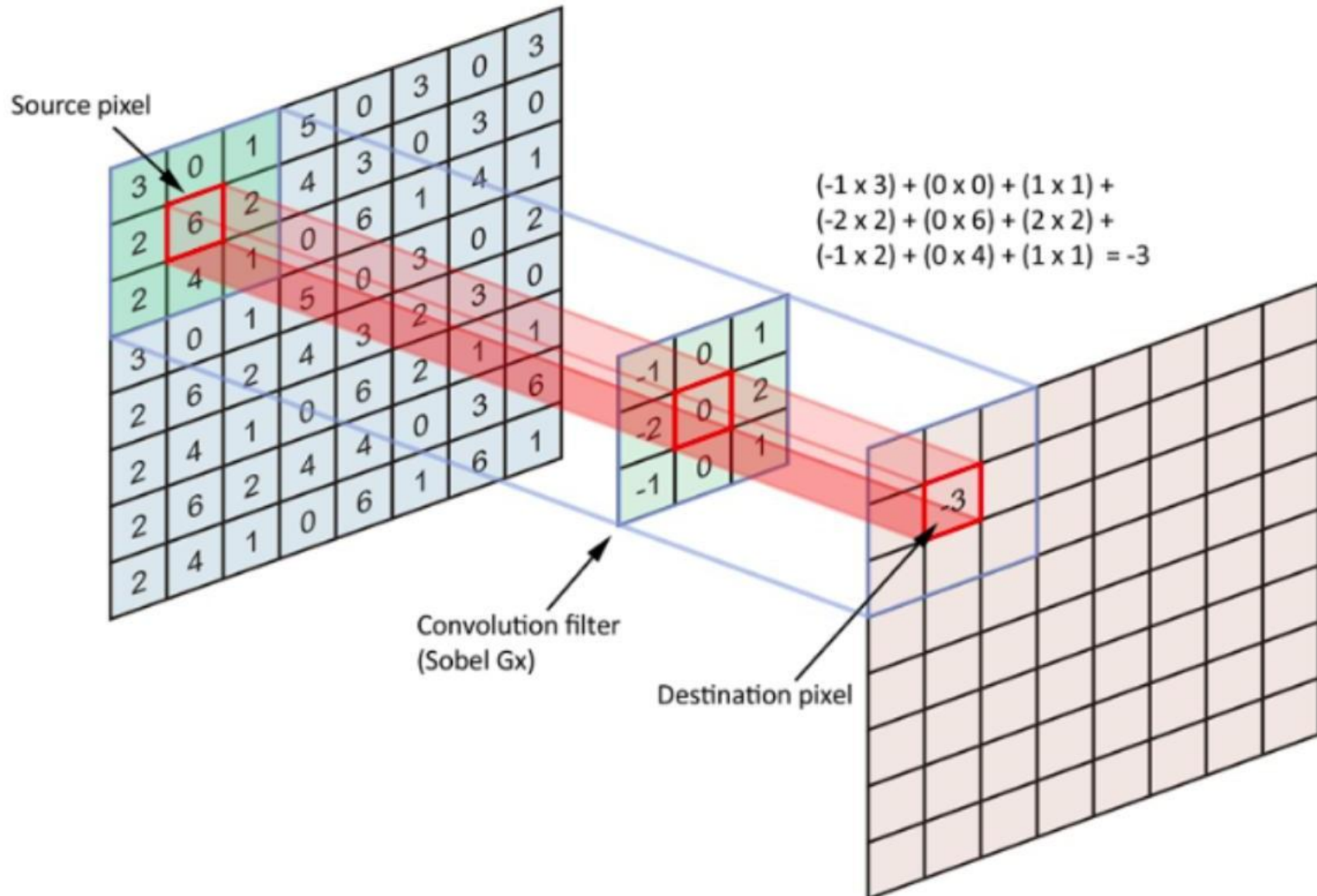
image classification →

- 82% cat
- 15% dog
- 2% hat
- 1% mug



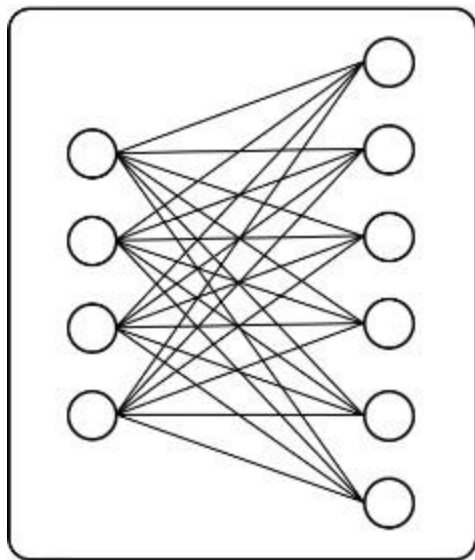
# 1. What is Convolutional Layer?

(2) **Filter** ( = weight matrix ) & (3) **합성곱**



## 2. Convolutional Layer – 수 식 적 표 현

### (1) FC의 수 식 적 표 현



$$W = [\mathbf{w}_0, \mathbf{w}_1, \dots, \mathbf{w}_{M-1}]^T$$

$$\mathbf{b} = [b_0, b_1, \dots, b_{M-1}]^T$$

$$y_0 = a(\mathbf{w}_0^T \mathbf{x} + b_0)$$

$$y_1 = a(\mathbf{w}_1^T \mathbf{x} + b_1)$$

$$\vdots$$

$$y_{M-1} = a(\mathbf{w}_{M-1}^T \mathbf{x} + b_{M-1})$$

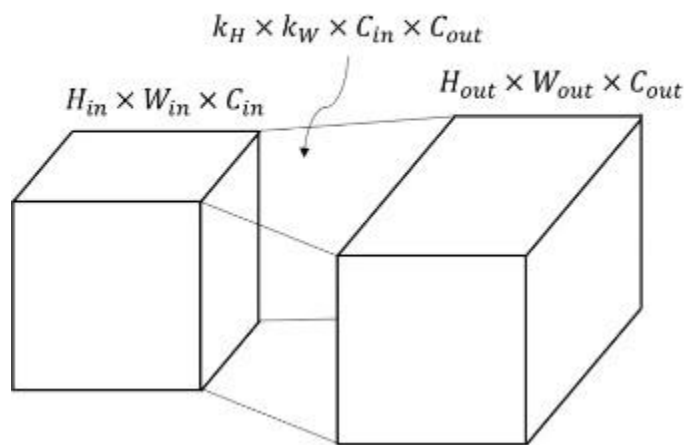
$$\mathbf{y} = a(W\mathbf{x} + \mathbf{b})$$

뉴런들이 곱해진 뒤, 모두 더해짐! ( Matrix 곱 연산으로 표현 )



## 2. Convolutional Layer – 수 식 적 표 현

### (2) Conv의 수 식 적 표 현



보통 3x3, 5x5, 7x7 등을 사용  $\mathbb{R}^{k_H \times k_W \times C_{in} \times C_{out}}$

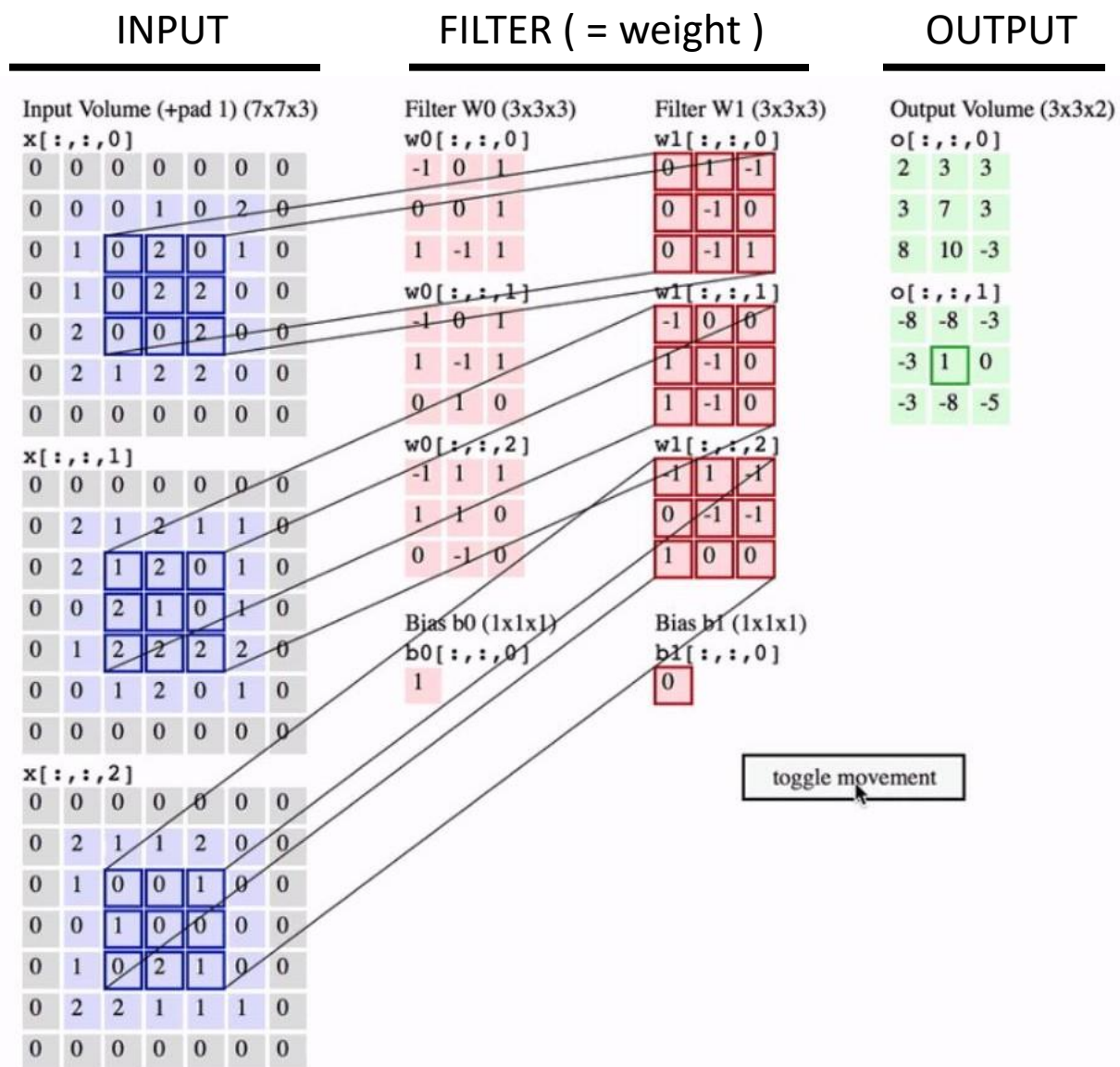
$$W = \begin{bmatrix} W_{0,0} & \cdots & W_{0,M-1} \\ \vdots & \ddots & \vdots \\ W_{N-1,0} & \cdots & W_{N-1,M-1} \end{bmatrix}$$

$$\mathbf{b} = [b_0, b_1, \dots, b_{M-1}]^T$$

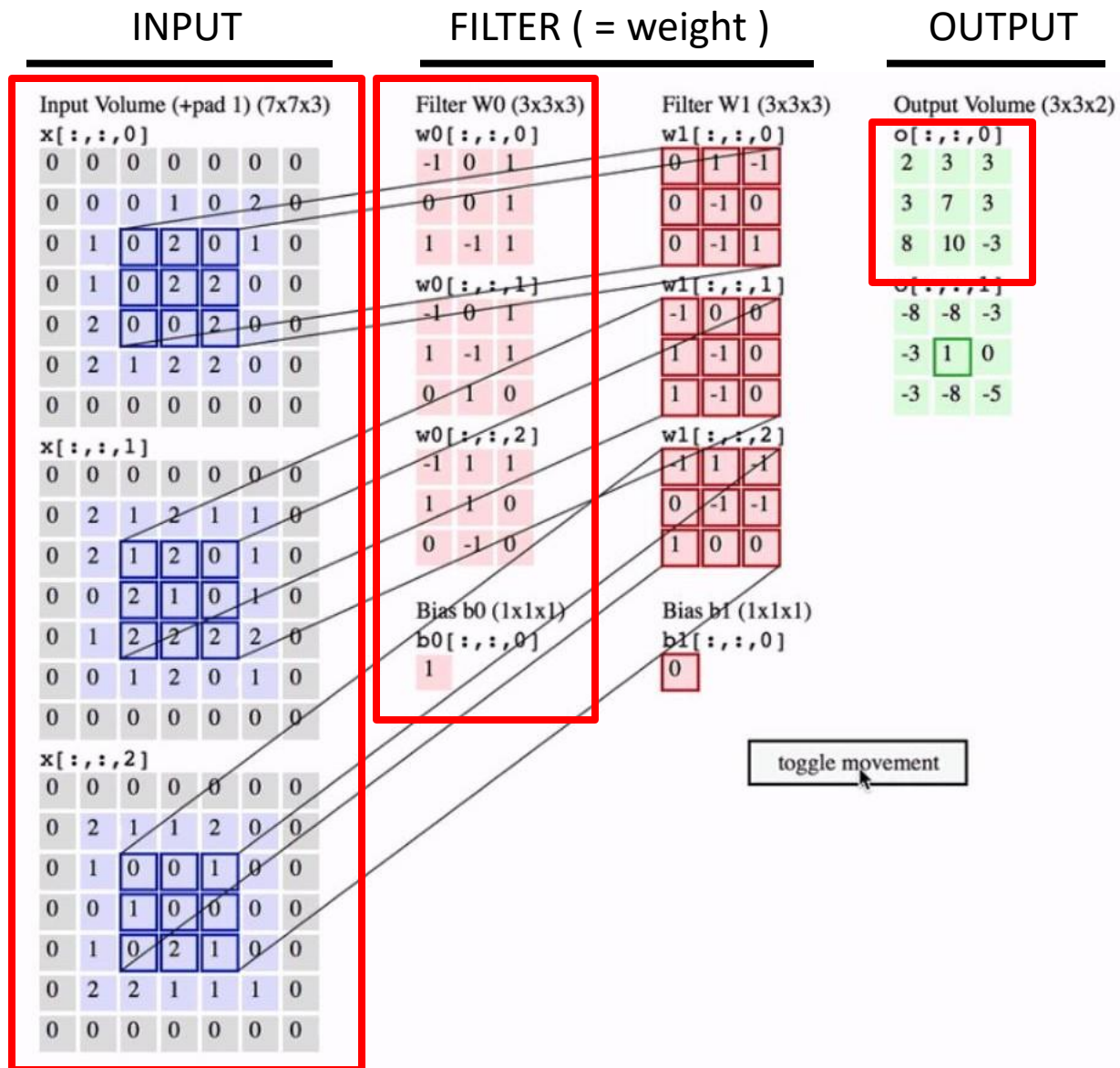
$$Y_{i,j} = a(W_{i,j} * X_i + b_j)$$

C(in) x C(out)번의 합성곱 연산이 이루어짐!

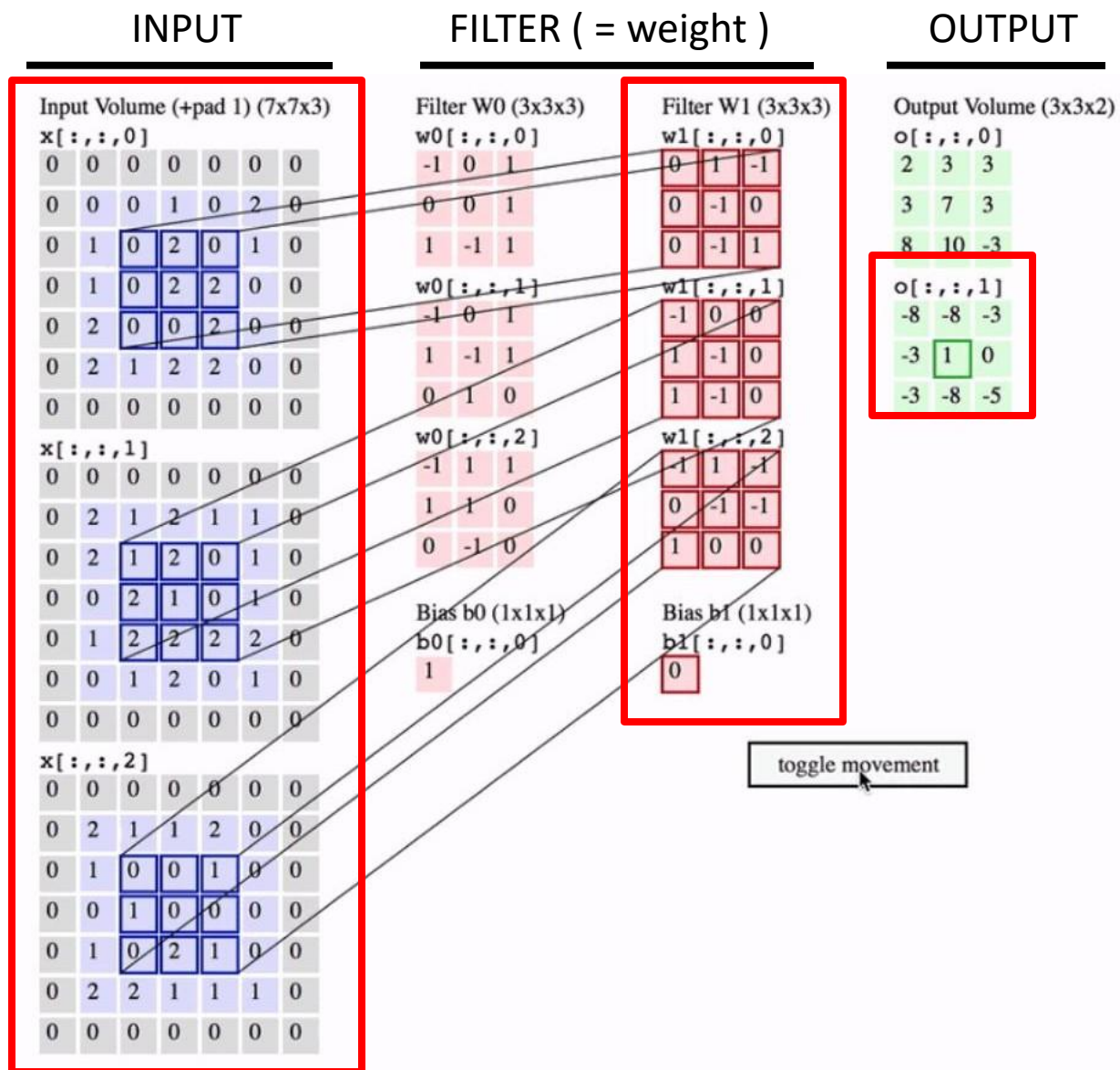
## 2. Convolutional Layer – 수 식 적 표 현



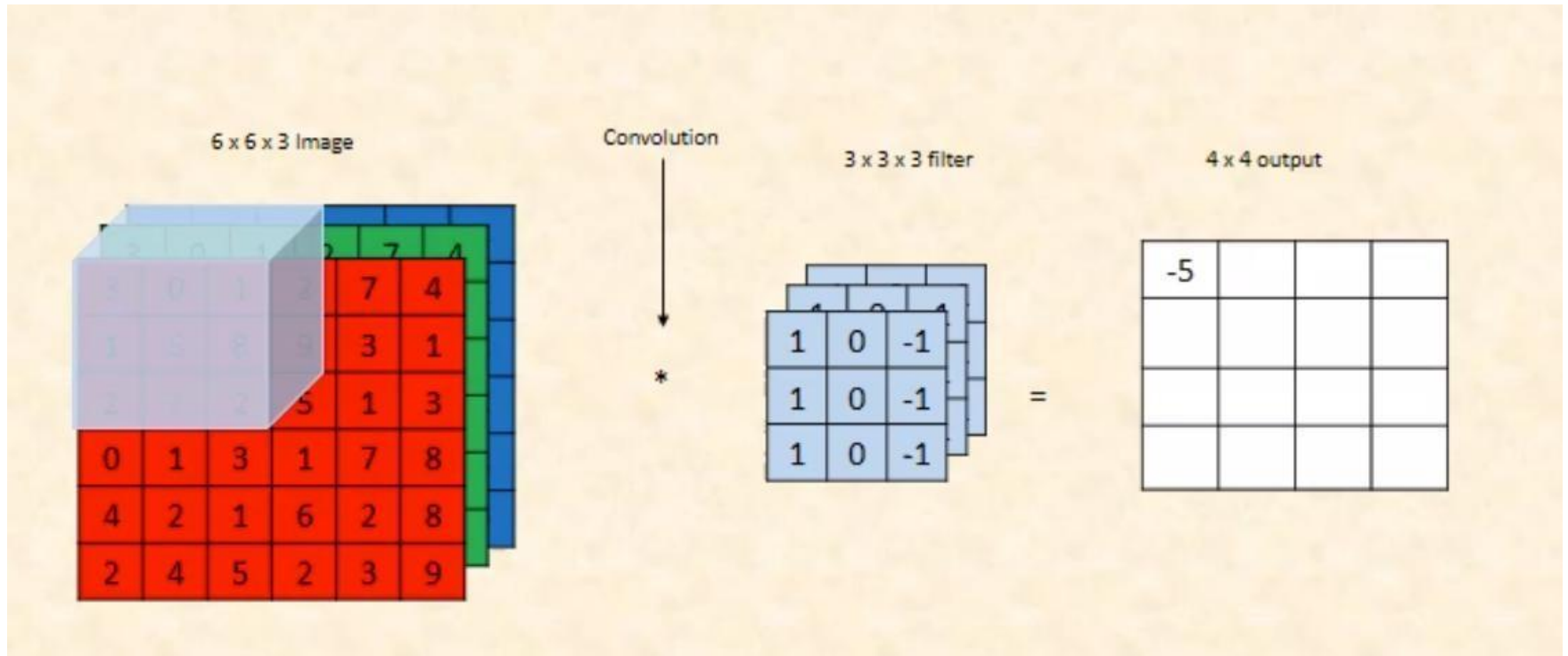
## 2. Convolutional Layer – 수 식 적 표 현



## 2. Convolutional Layer – 수 식 적 표 현

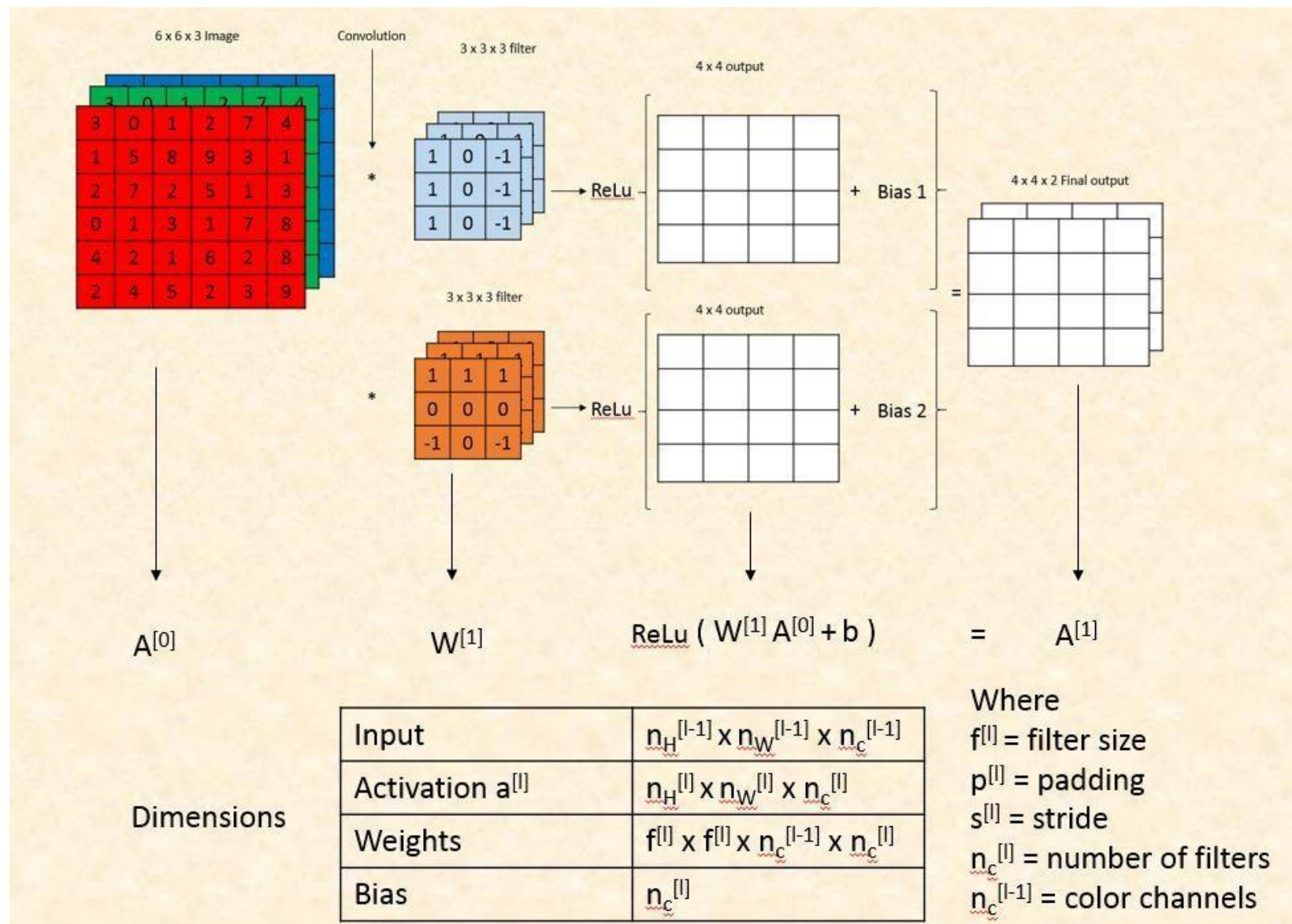


## 2. Convolutional Layer – 수 식 적 표 현





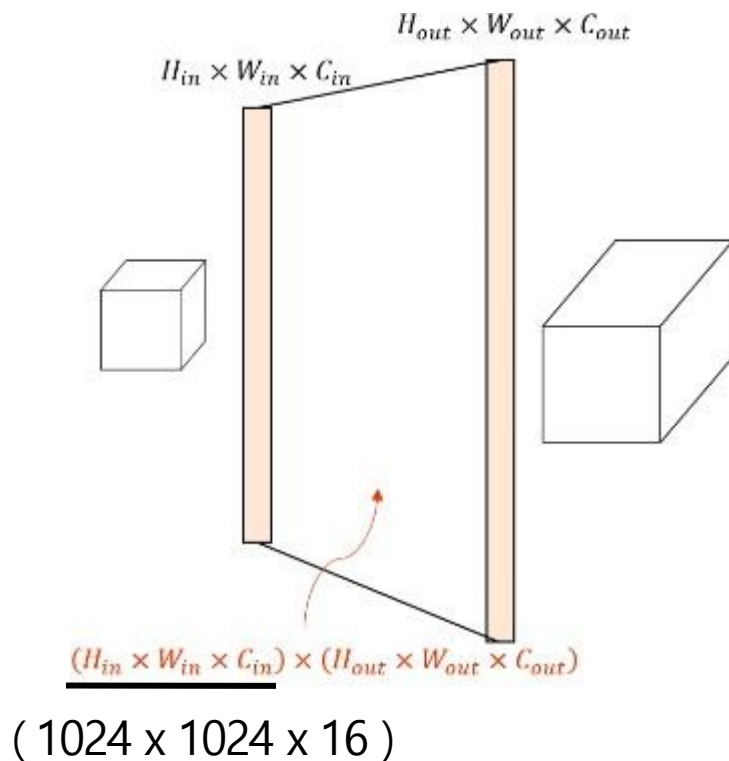
## 2. Convolutional Layer – 수 식 적 표 현



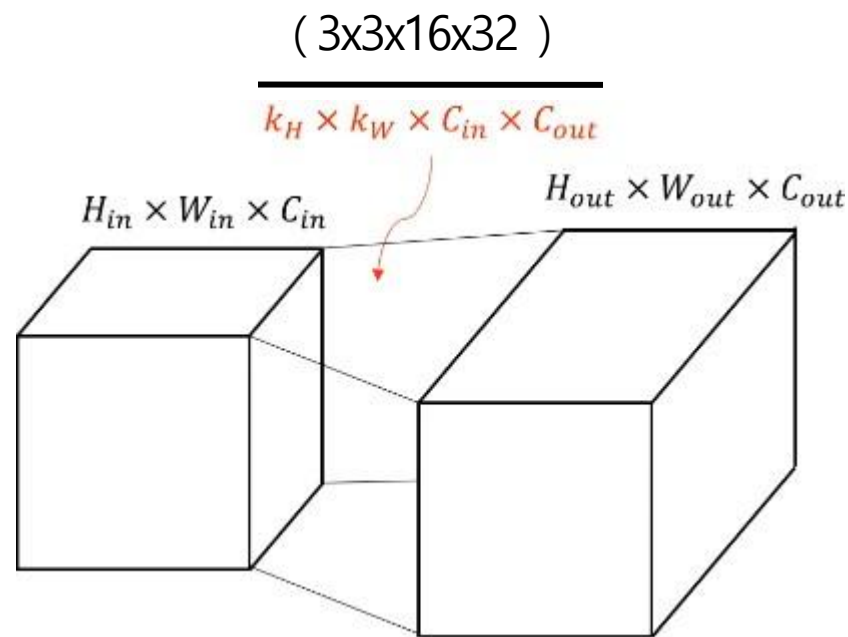


### 3. Why Convolutional Layer?

[ Fully Connected를 쓸 경우 ]



[ Convolutional Layer를 쓸 경우 ]



Ex) input 이미지가 1024x1024 크기일 때

FC 대신 Convolutional Layer 쓰는(써야하는) 이유?

**"Parameter 수의 감소!"**

## 4. 용 어 설 명

### (1) Channel

RED Channel



Green Channel



Blue Channel

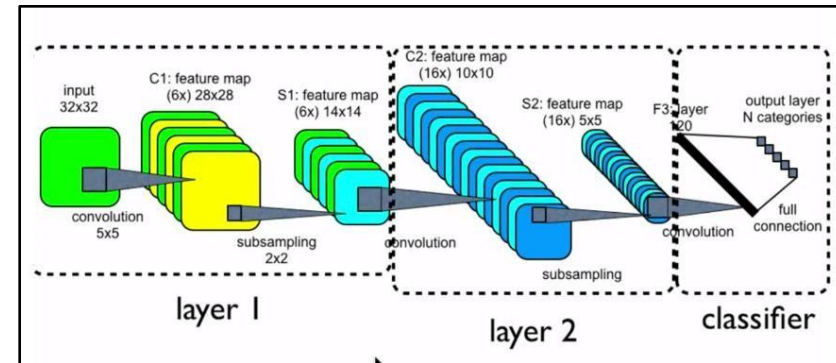
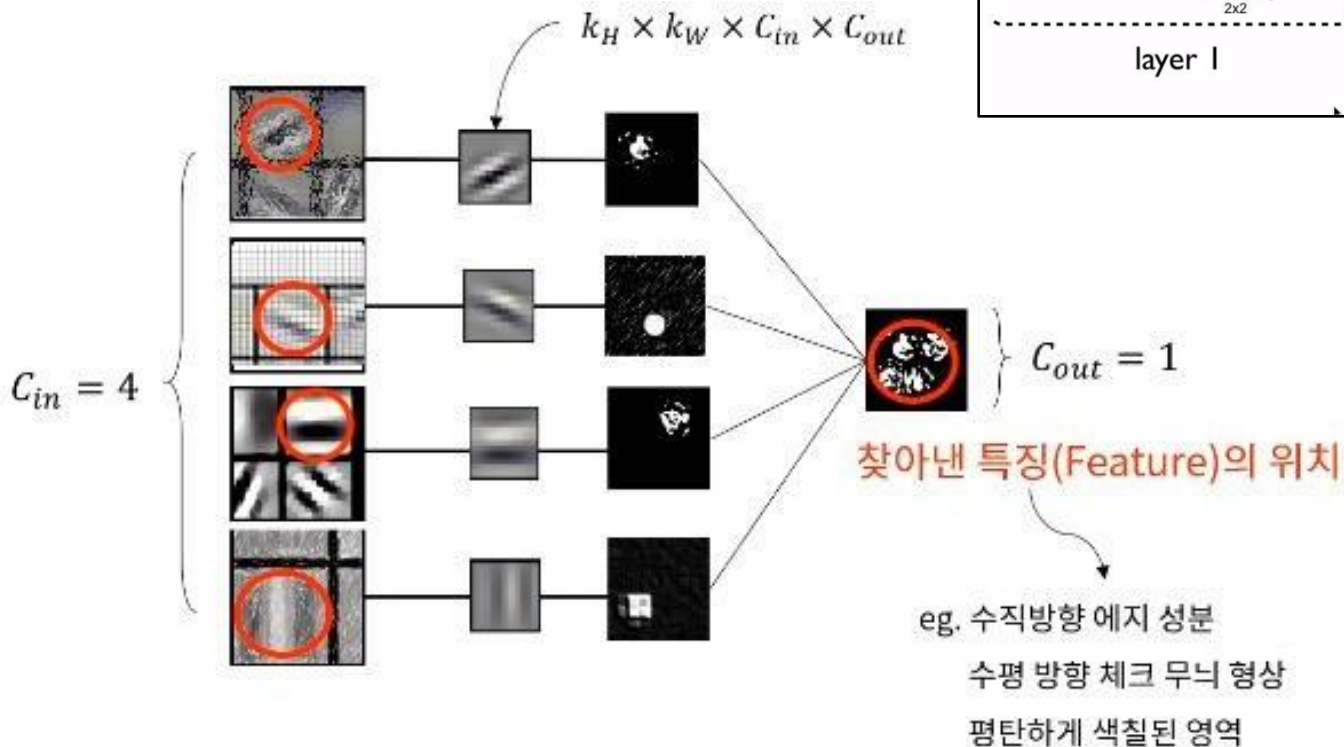


이미지 출처: [https://en.wikipedia.org/wiki/Channel\\_\(digital\\_image\)](https://en.wikipedia.org/wiki/Channel_(digital_image))

# 4. 용 어 설 명

## (2) Feature Map

합성곱을 통해서 만들어진 출력!

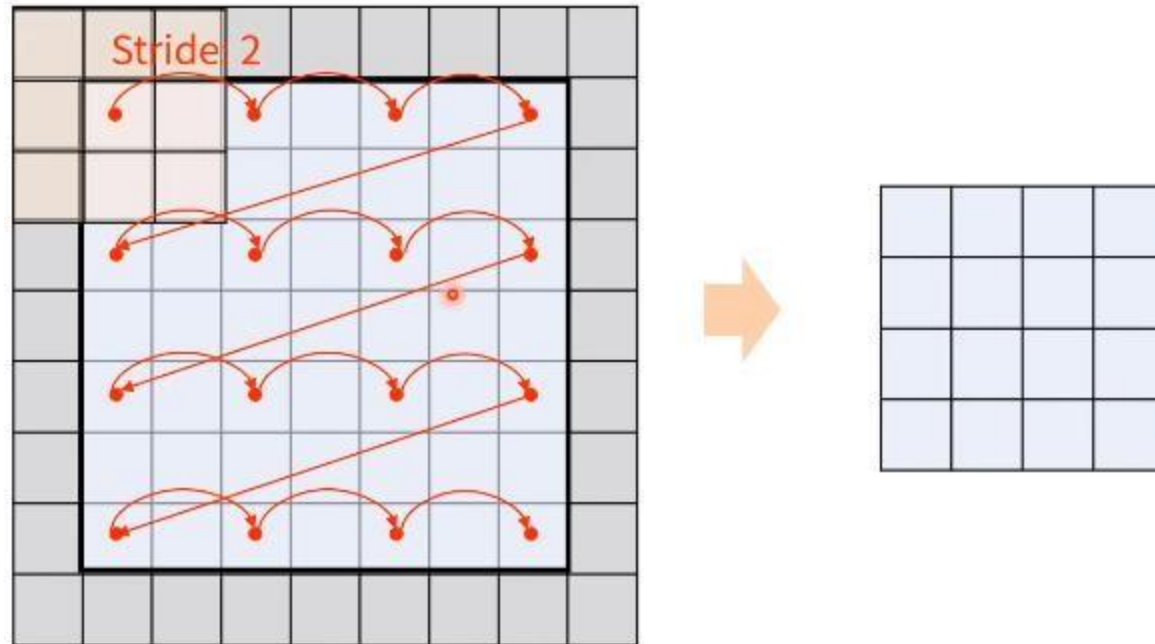


합성곱 계층의 의미?

여러 채널에서 특별한 "특징"이 나타내는 위치를 찾아내는 역할!

## 4. 용 어 설 명

### (3) Stride



합성곱 계산시, 커널을 이동시키는 거리!

Stride를 크게 할 수록, Feature Map의 크기는 줄어듦

## 4. 용 어 설 명

### (4) Padding

0	0	0	0	0	0
0	35	19	25	6	0
0	13	22	16	53	0
0	4	3	7	10	0
0	9	8	1	3	0
0	0	0	0	0	0

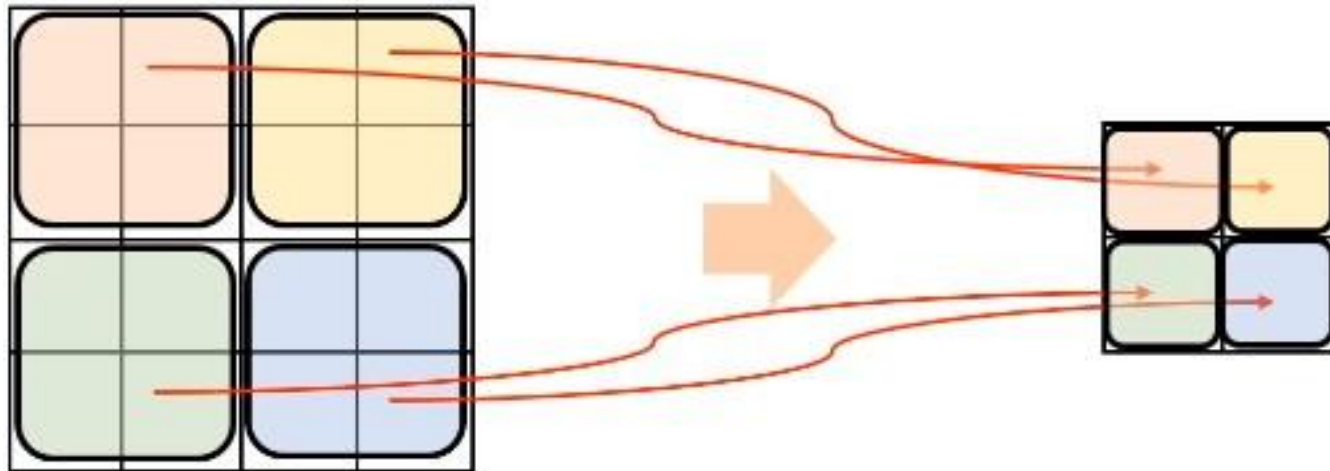
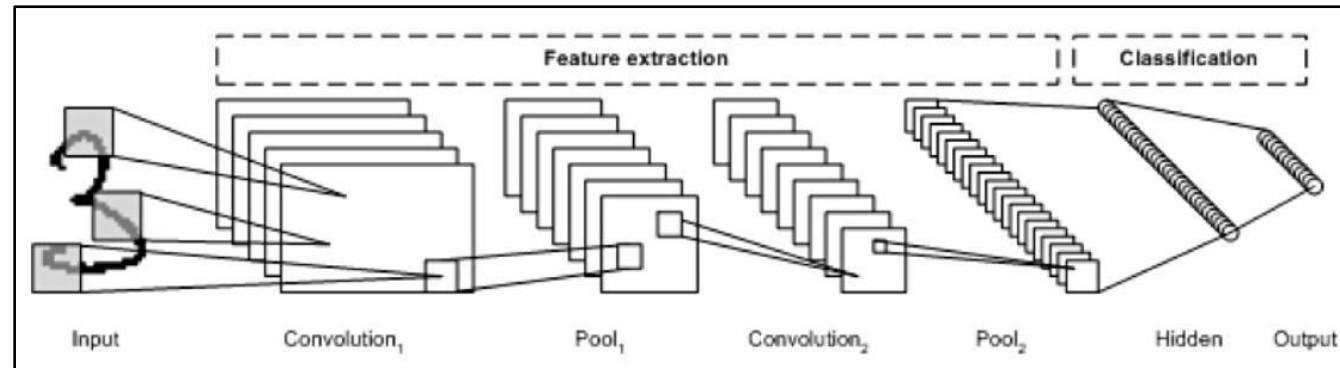
주위에 숫자(주로 0)을 둘러싸우는 것! ( For 크기 보존 )

합성 곱 연산시, 필터(커널)의 크기에 따라 영상의 크기가 줄어드는 문제 발생

ex) 크기가  $(2N+1)$ 인 커널 -> 상하좌우로  $N$ 개의 zero-padding!

# 4. 용 어 설 명

## (5) Pooling

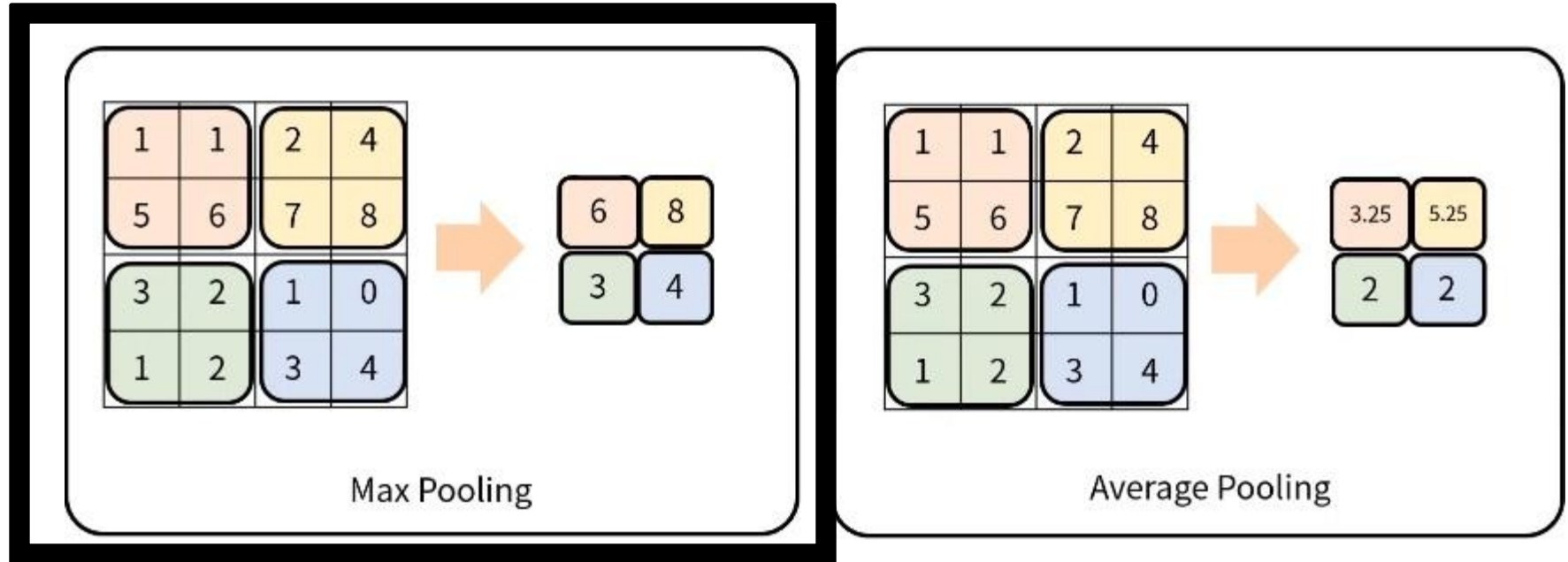


**Sampling**이라고 보면 됨! (여러 화소를 하나의 화소로 축약! )  
즉, "영상/사진의 크기가 줄어들고, 정보가 종합(요약)된다!"



## 4. 용 어 설 명

### (5) Pooling



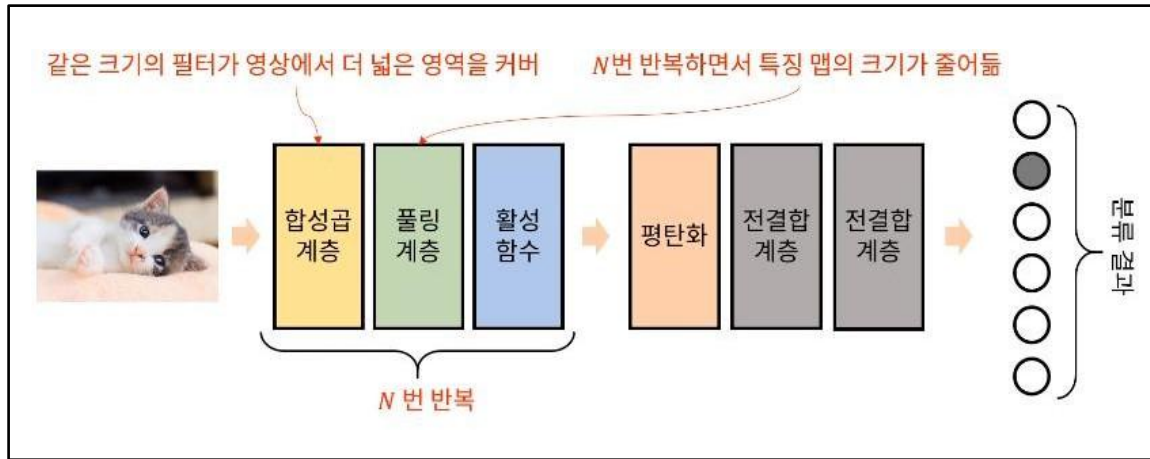
Pooling Layer ( vs Convolutional Layer )

- 1) 학습대상 parameter 없음
- 2) Pooling Layer를 통과하면 행렬의 크기 감소
- 3) Pooling Layer를 통해서 채널 수 변경 없음

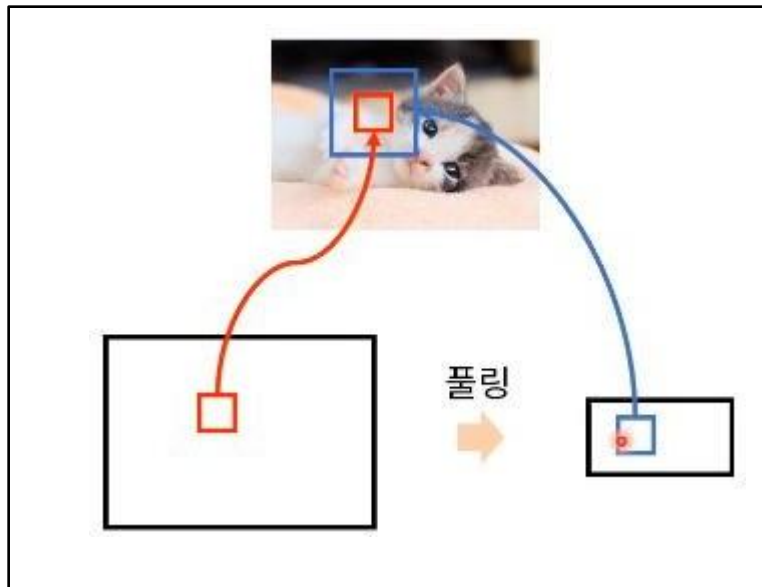
## 4. 용 어 설 명

### (5) Pooling

why pooling? Receptive Field?



(h,w)영상 크기 작아지고,  
(c)채널 수 늘어나게 될 것!

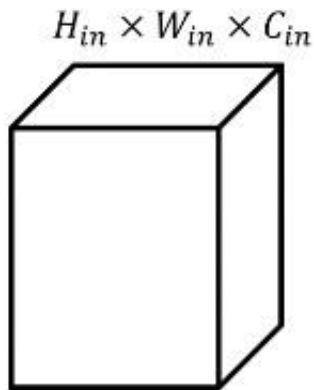


같은 크기의 filter여도, Pooling에 의해  
작아진 feature map에 적용되면, “원본 영상에서  
차지하는 범위(=receptive field)” 넓다

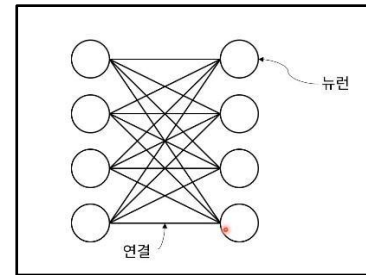
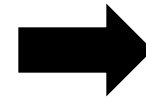
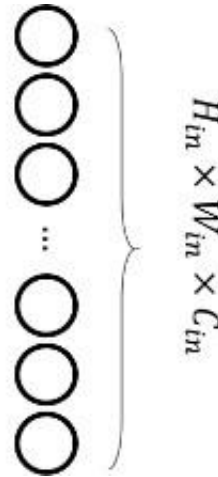
## 4. 용 어 설 명

### (6) Flatten

Ex)  $1024 \times 1024 \times 3$



Ex) 3,145,728



**평탄화!** -> Vector로 만들어줌

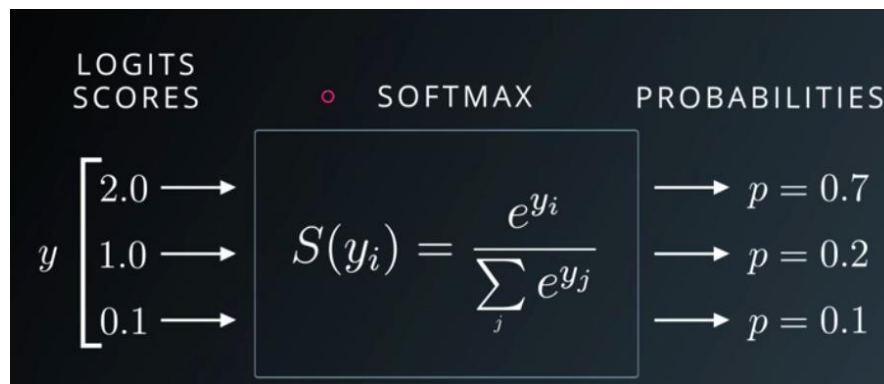
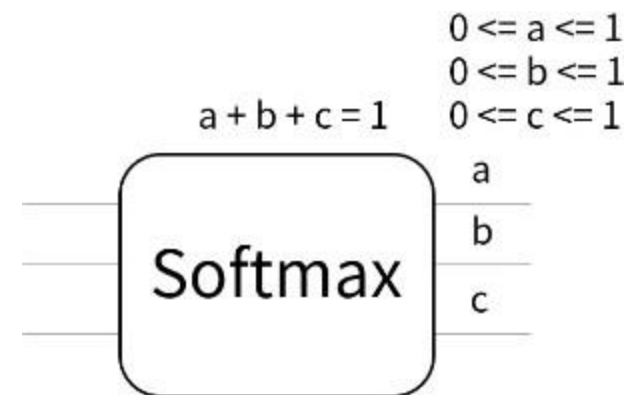
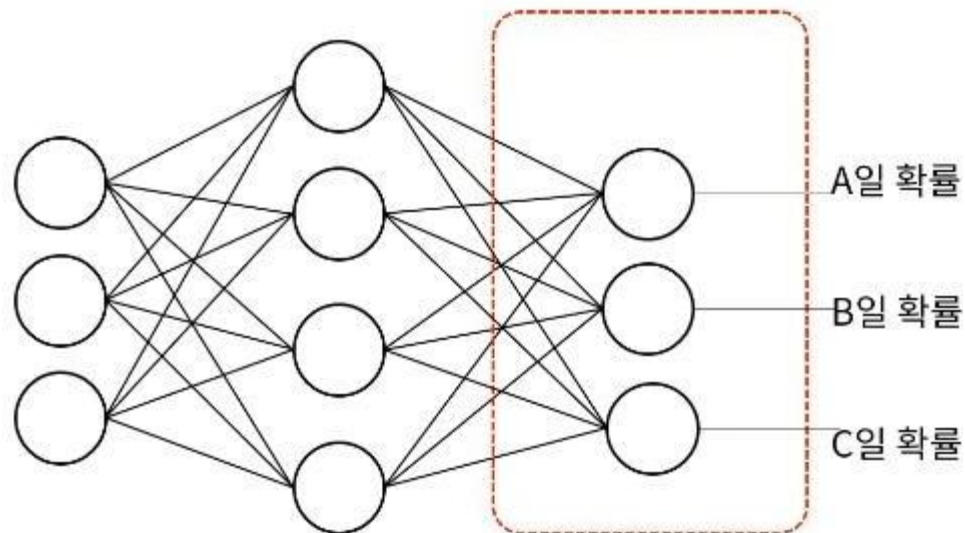
단지 쪽 피는 것을 의미! (연산 과정 X)

Convolutional Layer & Fully Connected 사이를 연결해주는 역할

# 4. 용 어 설 명

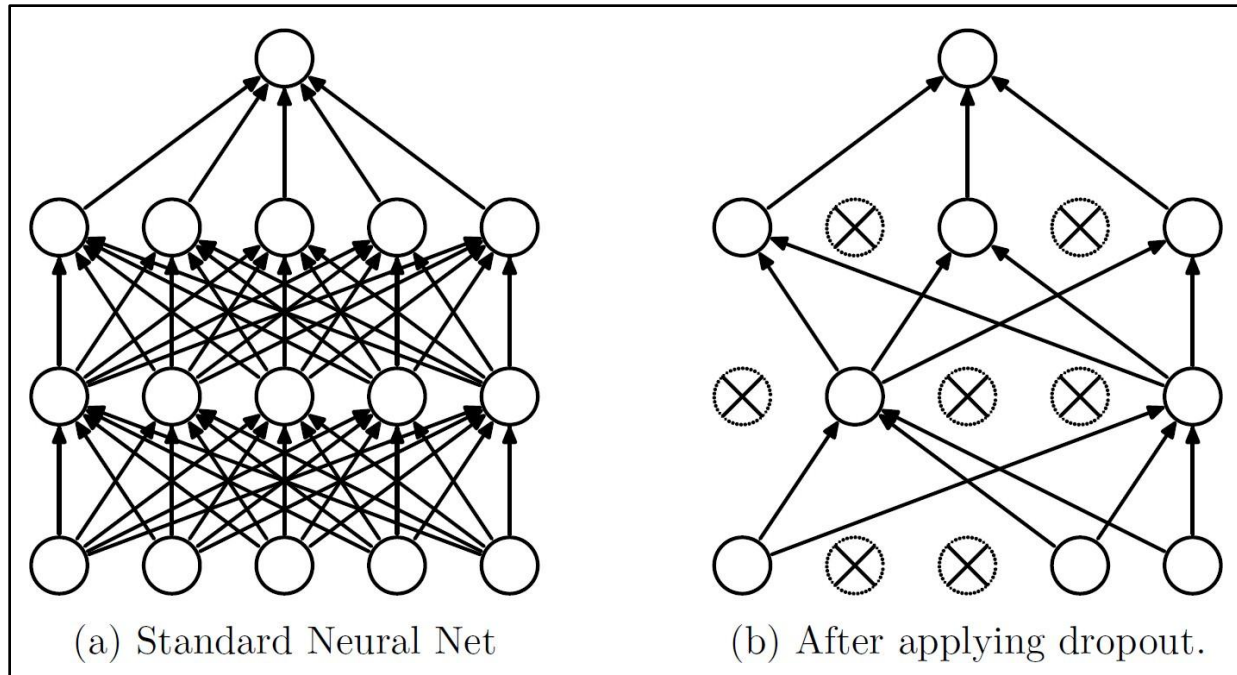
## (7) Softmax Function

for Multi-class Classification



# 5. Dropout

**To prevent overfitting!**



Training 할 때만!  
( testing할 때는 X )

Ex) Dropout rate 0.3 : 30%의 node는 사용 X ( less parameter )

**Overfitting 문제가 해결되는 이유를, 직관적으로 생각해보자!**

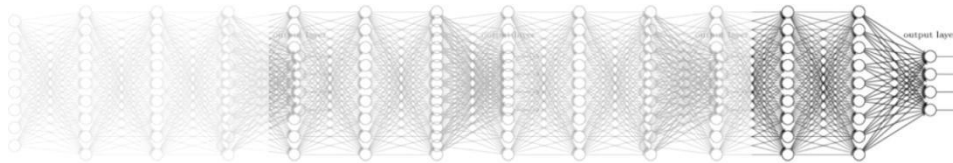
100가지의 일, 100명의 사람이 하는 경우? 20명의 사람이 하는 경우?

# 6. Gradient Vanishing Problem

## Gradient Vanishing?

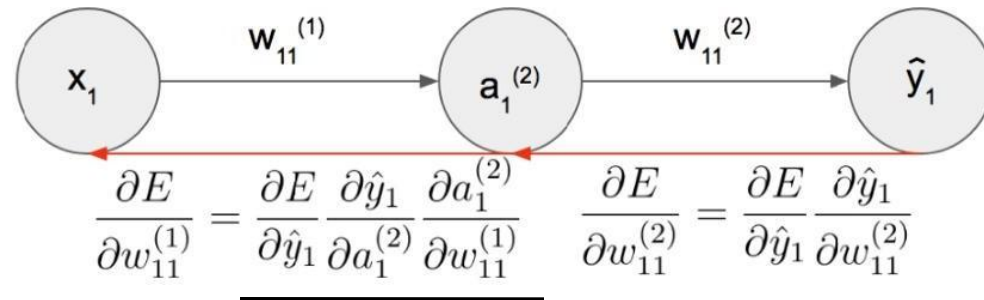
(기울기 소실 문제)

Vanishing gradient (NN winter2: 1986-2006)



DEEP 할수록 gradient vanishing! WHY?

( Backpropagation의 수식을 생각해 보자 )



( input layer 쪽으로 갈수록, 미분값의 곱이 많아짐! )



## 6. Gradient Vanishing Problem

### Derivative of Sigmoid Function

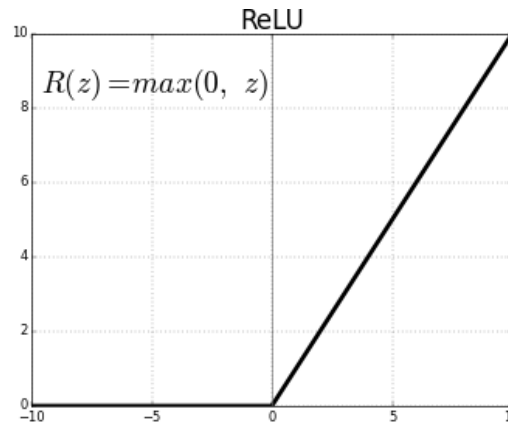
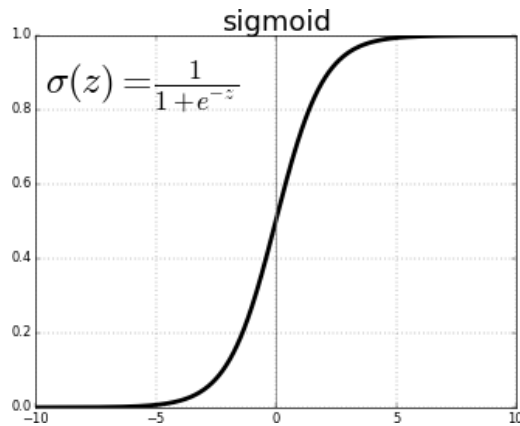
$$\begin{aligned}\frac{\partial \times_3}{\partial Z_3} &= \frac{\partial}{\partial Z_3} \left[ \frac{1}{1 + e^{-Z_3}} \right] \\&= \frac{\partial}{\partial Z_3} (1 + e^{-Z_3})^{-1} \\&= -(1 + e^{-Z_3})^{-2} (-e^{-Z_3}) \\&= \frac{e^{-Z_3}}{(1 + e^{-Z_3})^2} \\&= \frac{1}{1 + e^{-Z_3}} \cdot \frac{e^{-Z_3}}{1 + e^{-Z_3}} \\&= \frac{1}{1 + e^{-Z_3}} \cdot \frac{(1 + e^{-Z_3}) - 1}{1 + e^{-Z_3}} \\&= \frac{1}{1 + e^{-Z_3}} \cdot \left( \frac{1 + e^{-Z_3}}{1 + e^{-Z_3}} - \frac{1}{1 + e^{-Z_3}} \right) \\&= \frac{1}{1 + e^{-Z_3}} \cdot \left( 1 - \frac{1}{1 + e^{-Z_3}} \right) \\&= \times_3 \cdot (1 - \times_3) \longrightarrow \text{Maximum : 0.25}\end{aligned}$$

# 6. Gradient Vanishing Problem

## Gradient Vanishing?

How to Solve?

1) Activation Function으로 Sigmoid -> ReLU

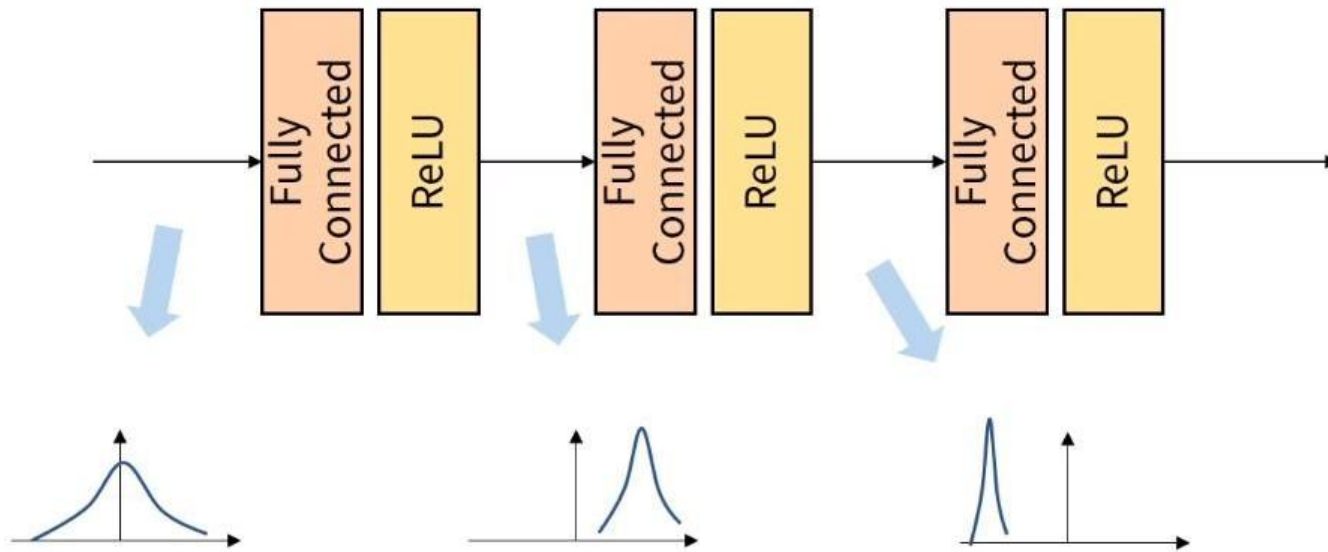


2) 적절한 weight 초기값 : Xavier Initialization, He Initialization

3) Batch Normalization

# 7. Batch Normalization

## Internal Covariance Shift



Training 과정에서, parameter의 변화로 인해

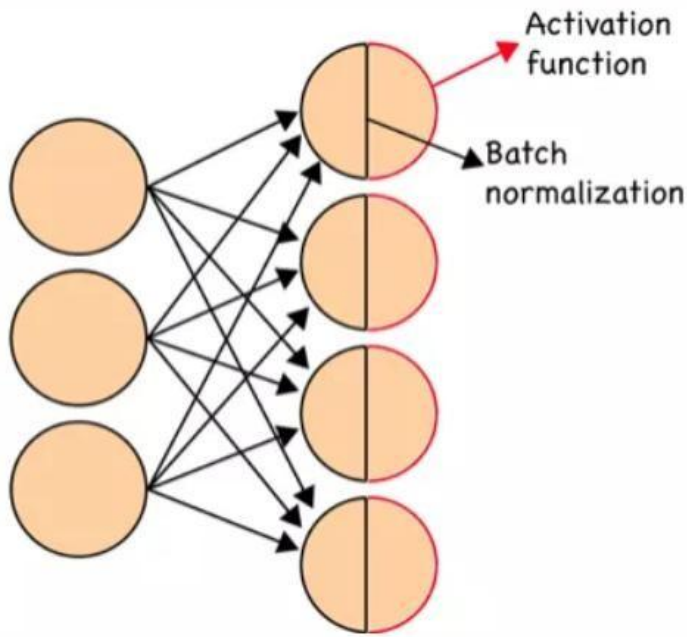
**Input data의 distribution이 달라지는** 현상! -> 불안정한 training

이 문제를 해결하는 방법이 **Batch Normalization!**

이로써 안정적인 training을 가능하게 하고, gradient vanishing problem 해결!

# 7. Batch Normalization

## Batch Normalization



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

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

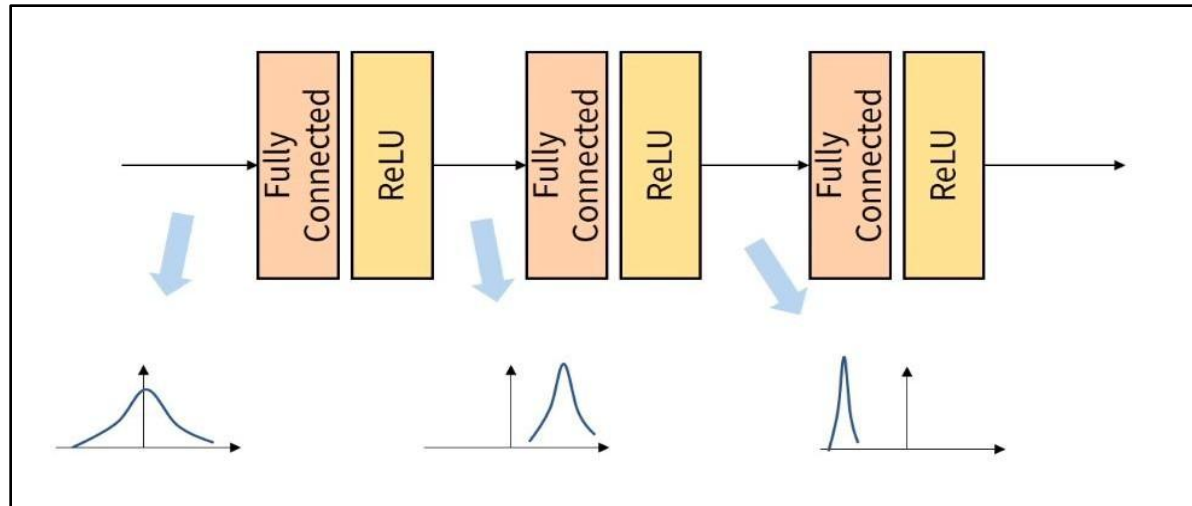
training data 전체에 대하여 normalize하는 것이 좋겠지만,  
Mini-batch Gradient Descent 방식을 사용하게 되면,  
parameter의 update Mini-batch단위로 이루어지기  
때문에, Mini-batch를 단위로 Batch Normalization(BN) 실시!

Internal Covariance Shift 문제를 해결하기 위해, " 층의 input의 distribution을  
평균 0, 표준편차 1인 input으로 normalize" 시키는 방법

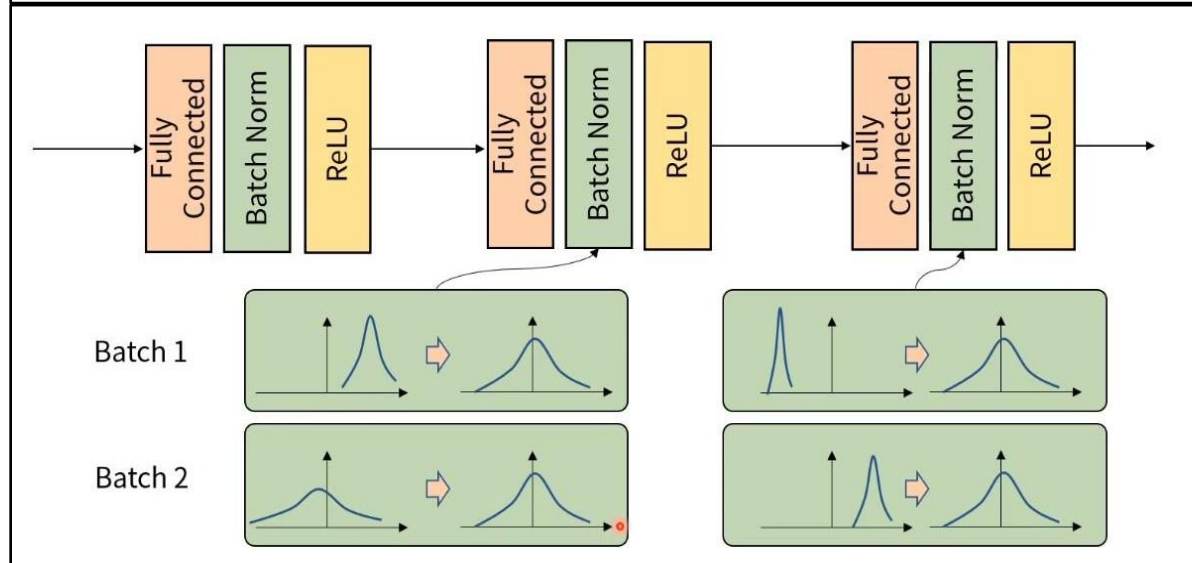
# 7. Batch Normalization

## Internal Covariance Shift

B N (X)



B N (O)



# 7. Batch Normalization

## Batch Normalization

< BN의 장점 >

1) **Gradient vanishing** 문제 해결!

2) **Regularization** 효과가 있음 ( overfitting 방지 )

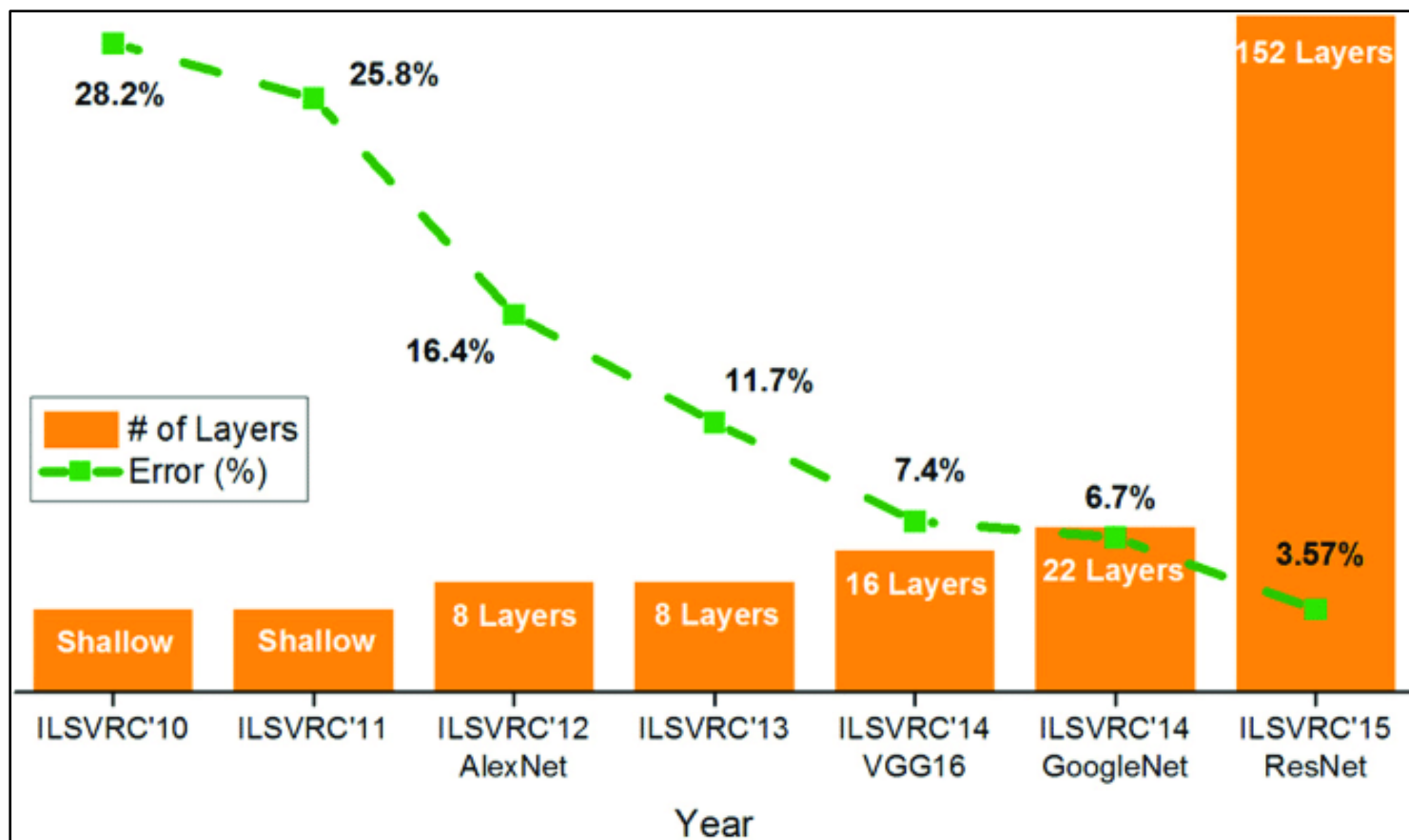
- (regularization 효과를 주는) Dropout을 제외할 수 있게 해줌

- Dropout 경우 효과는 좋지만 학습 속도가 느려진다는 단점!

3) Gradient의 **scale 영향을 덜 받음!**



## 7. 심화 CNN



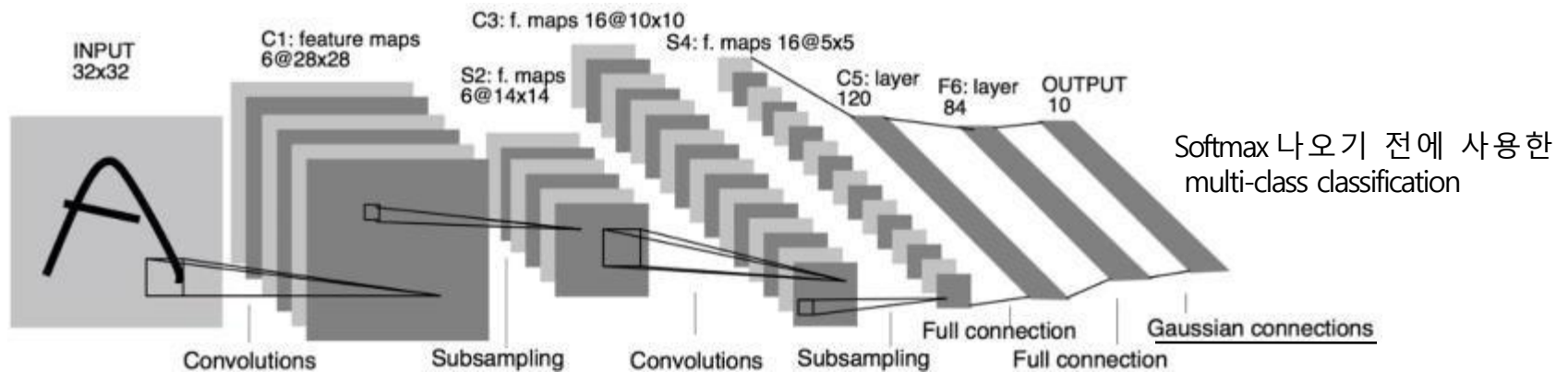
CNN의 대표 모델들에 대해 알아보자.

LeNet-5, AlexNet, VGG16, GoogLeNet, ResNet

# 7. 심화 CNN

## 1. LeNet-5

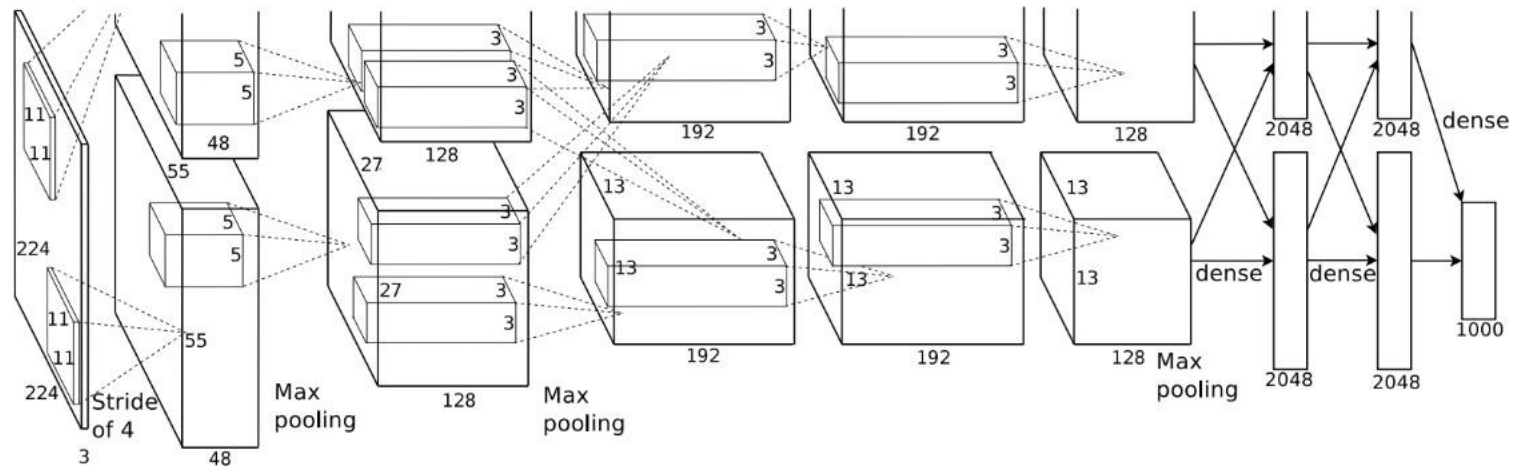
6채널의 28\*28 feature map



- 1998. CNN의 기본구조 창립
- MNIST data 사용

# 7. 심화 CNN

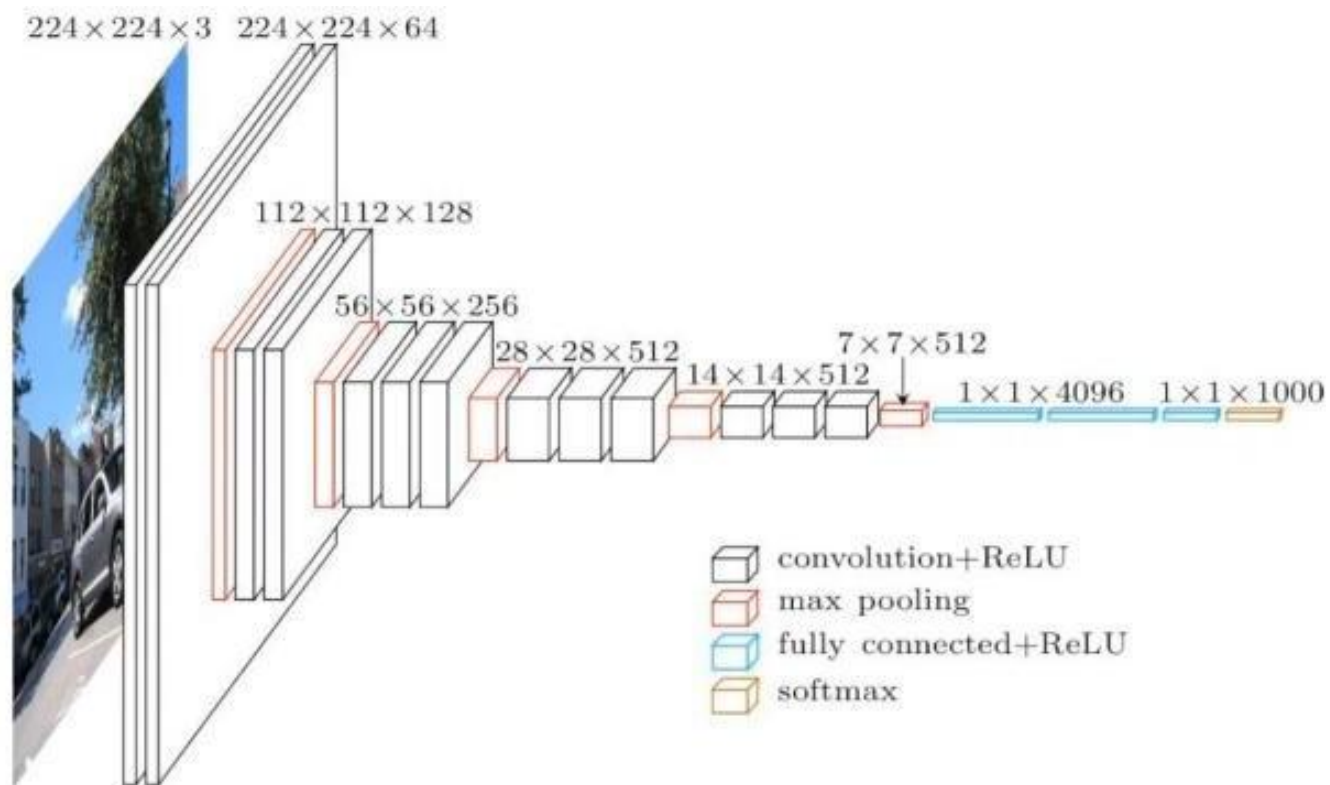
## 2. AlexNet



- Conv, Max Pooling, Dropout 5%
- Activation function : ReLU
- Batch Stochastic Gradient Descent

# 7. 심화 CNN

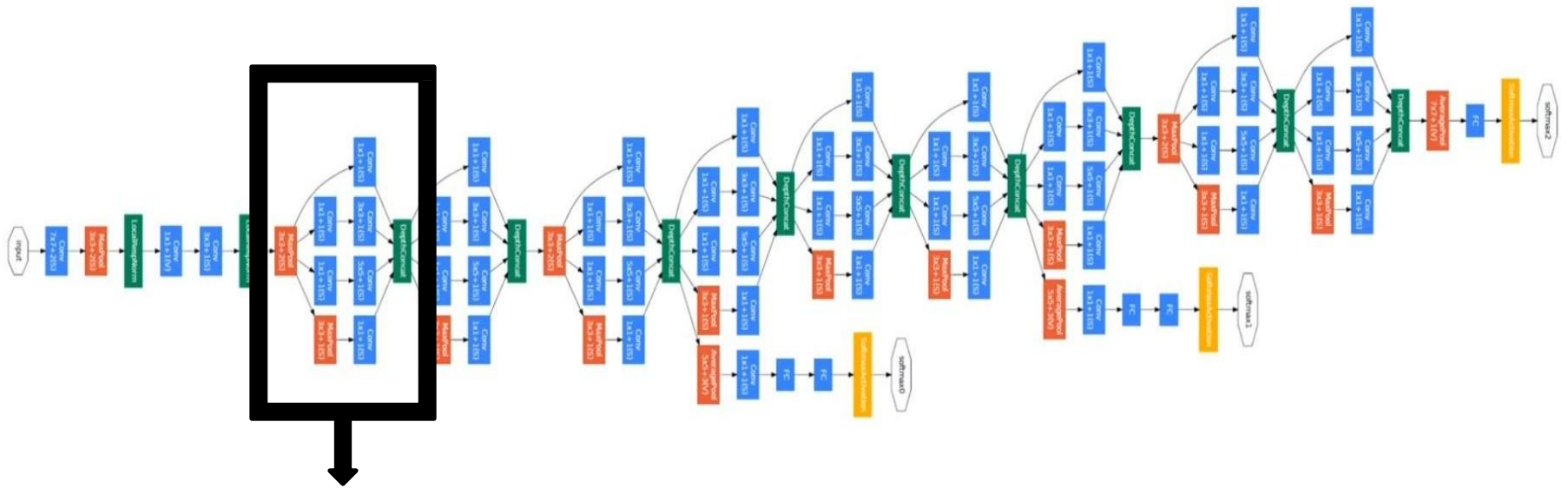
## 3. VGG-16



# 7. 심화 CNN

## 4. GoogLeNet (Inception)

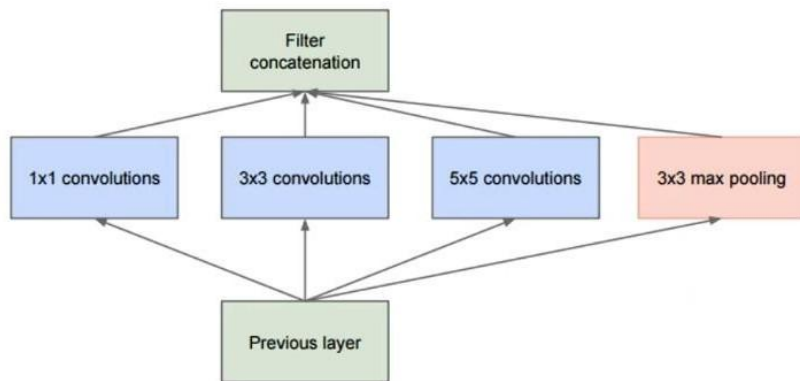
“Let's go Deeper and Deeper”



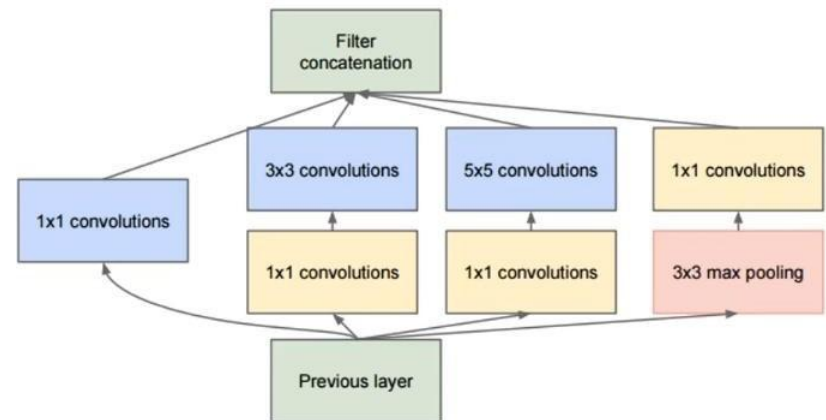
Inception 모듈에 대해 먼저 이해해야!

# 7. 심화 CNN

## 4. GoogLeNet (Inception)



(a) Inception module, naïve version



(b) Inception module with dimension reductions

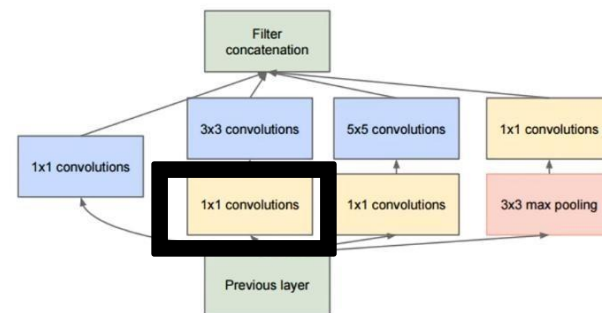
**Inception module** : 다양한 크기의 합성곱 계층을 한 번에 계산!

1x1 convolution : for 연산량 줄이기! (bottle neck 구조라고도 함)

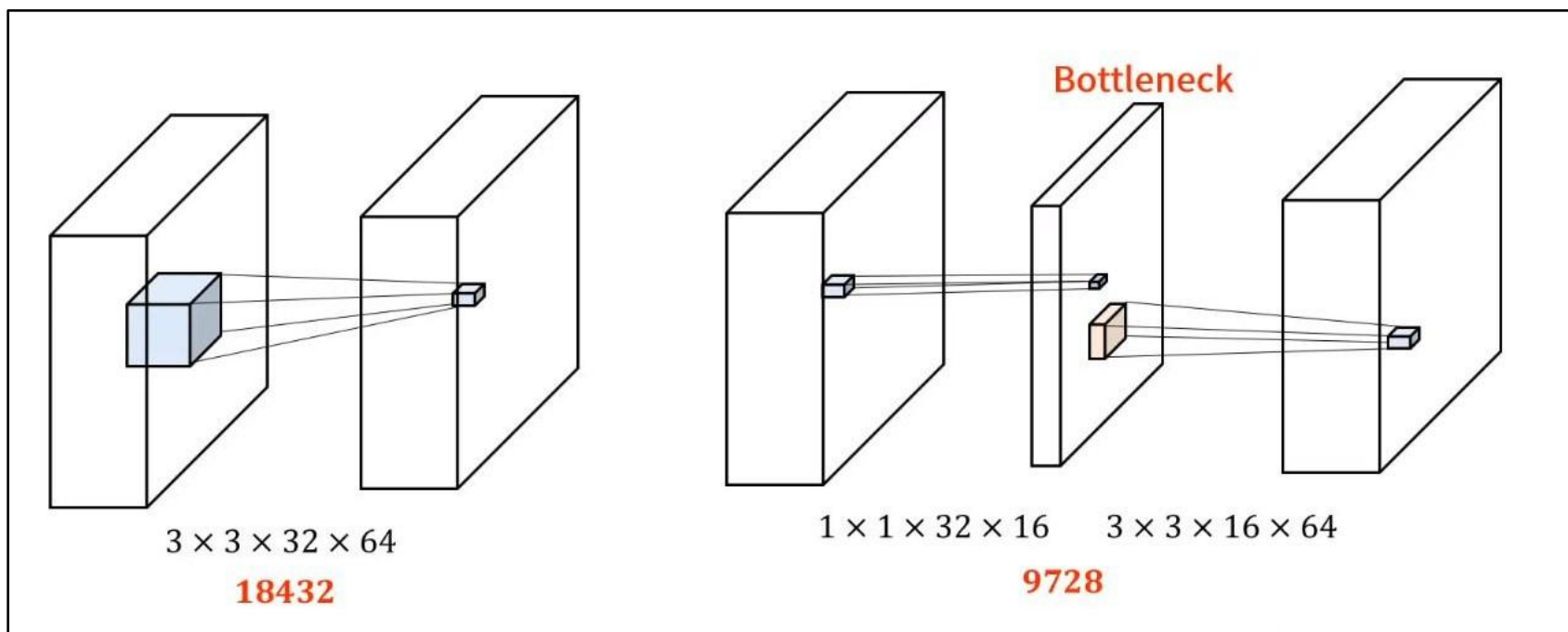


# 7. 심화 CNN

## 4. GoogLeNet (Inception)



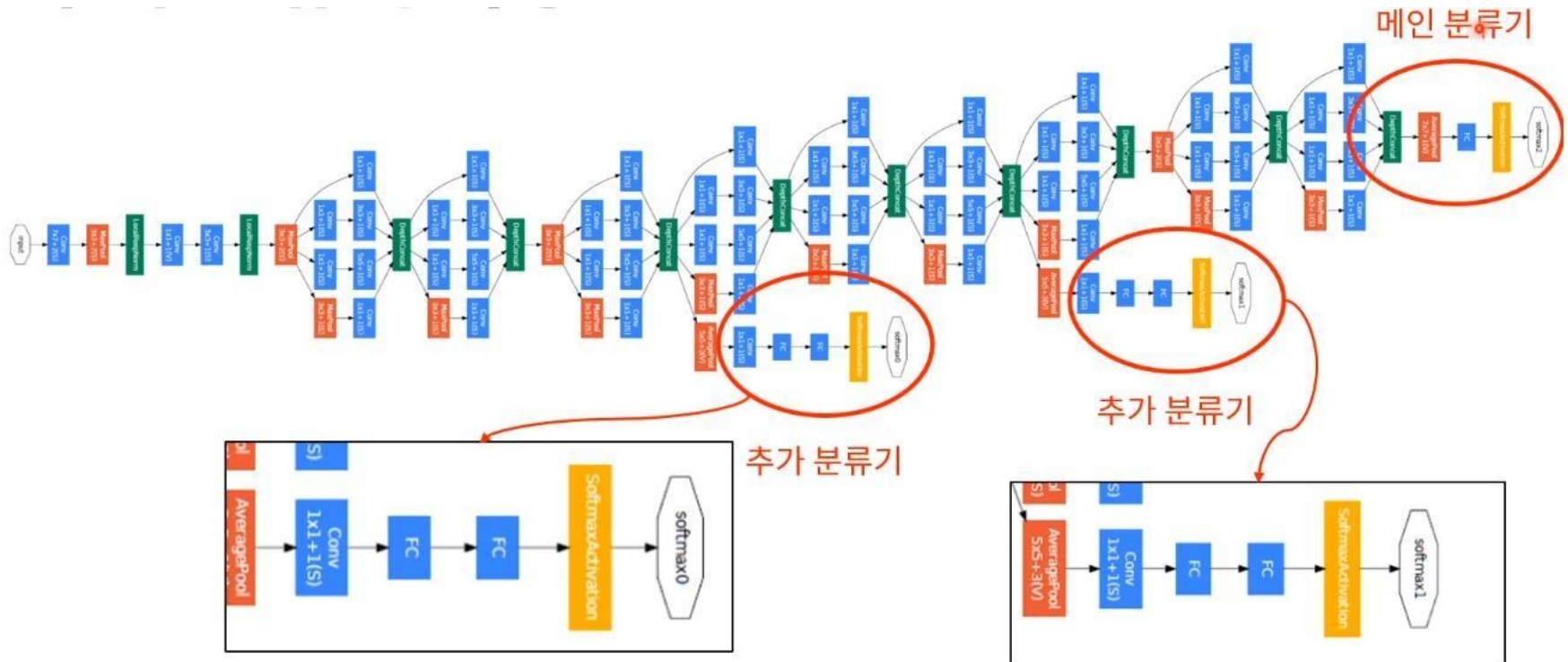
(b) Inception module with dimension reductions



1x1 Conv를 하면, 왜 연산량이 줄어들까?

# 7. 심화 CNN

## 4. GoogLeNet (Inception)



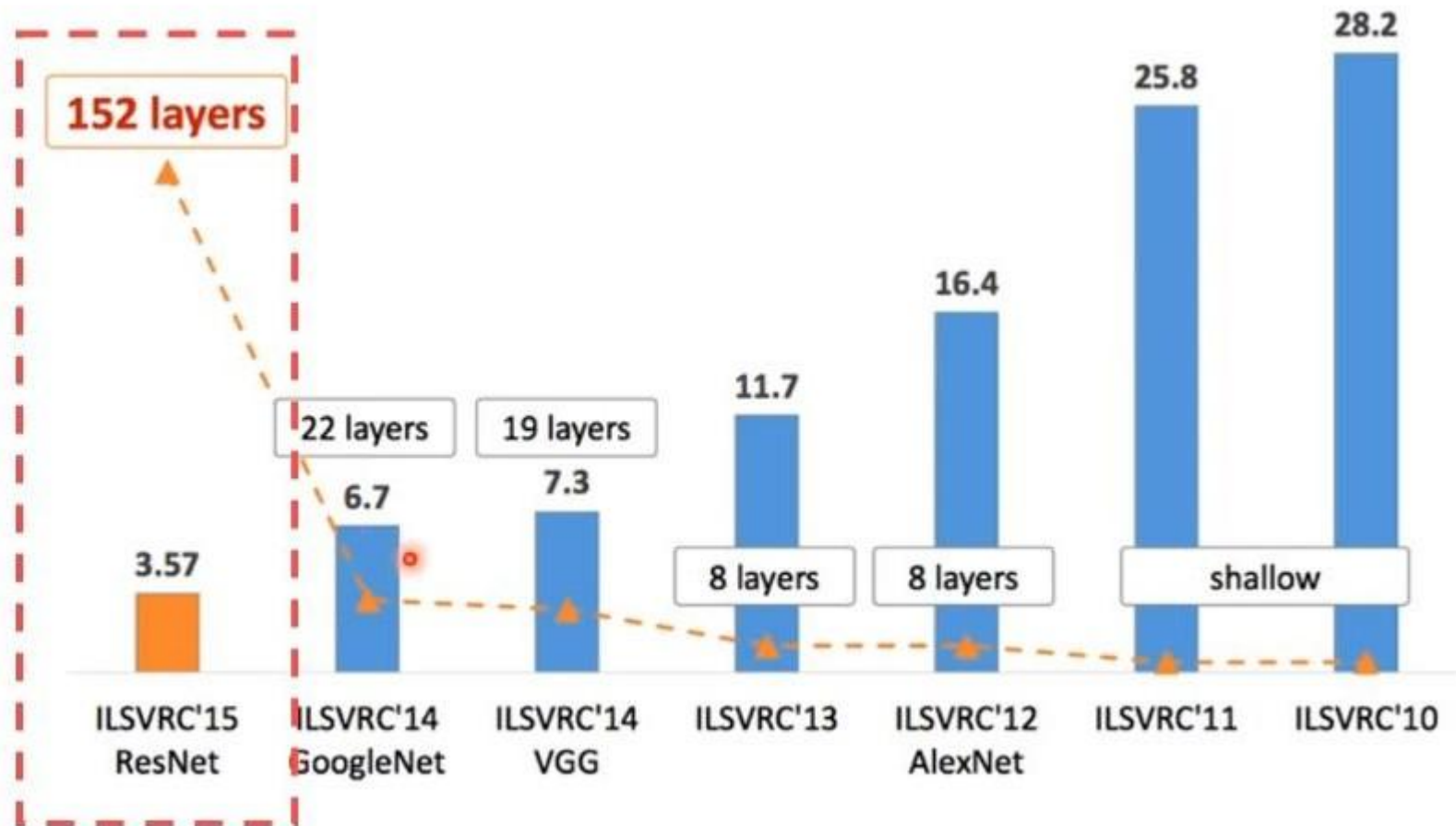
역전파에서 기울기 소실이 발생하는 것을 방지하기 위해 추가 분류기 존재

## 7. 심화 CNN

### 5. ResNet

152 Layers? How that deep?

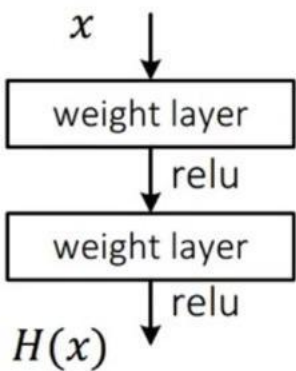
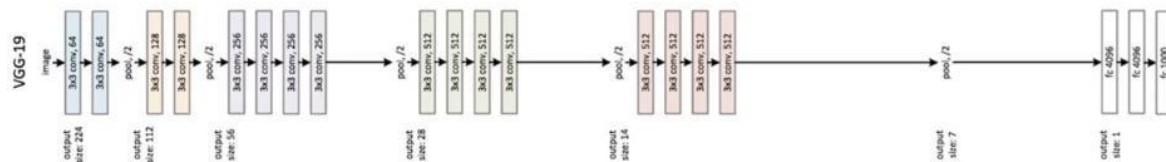
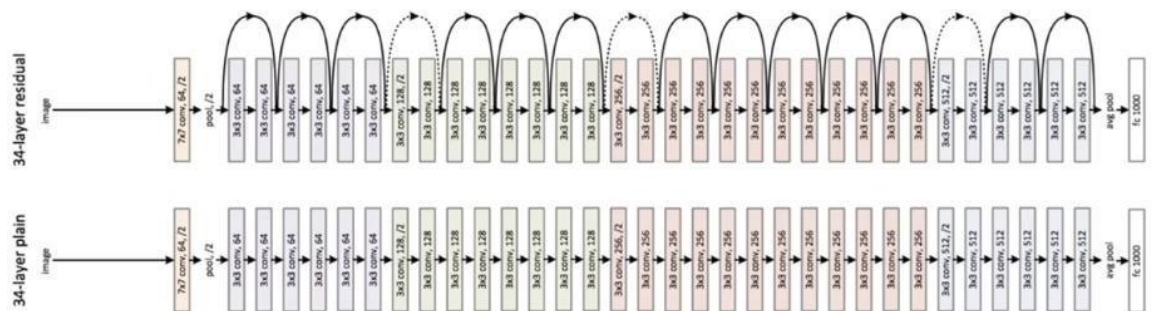
### Residual Block



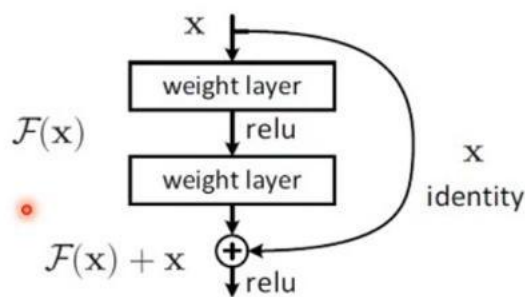
# 7. 심화 CNN

## 5. ResNet

사이에 사이에 있는 Skip-Connection이 한 몫함!



일반적인 구조



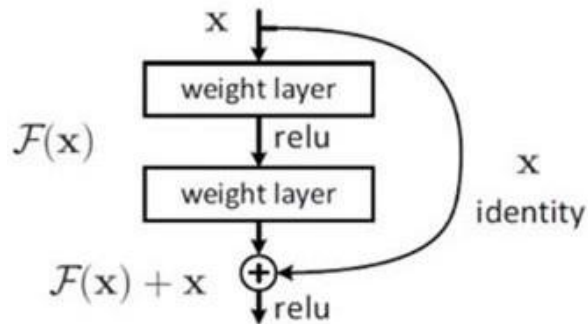
Residual 구조

Feature를 추출하기 전/후를 더함!

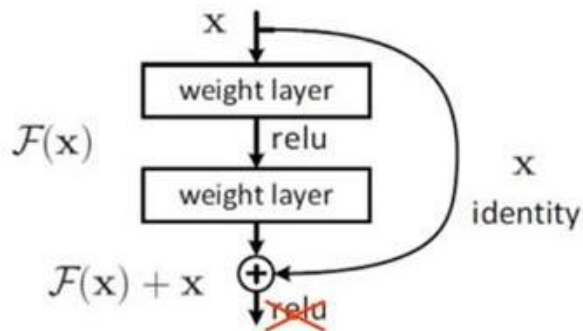
# 7. 심화 CNN

## 5. ResNet

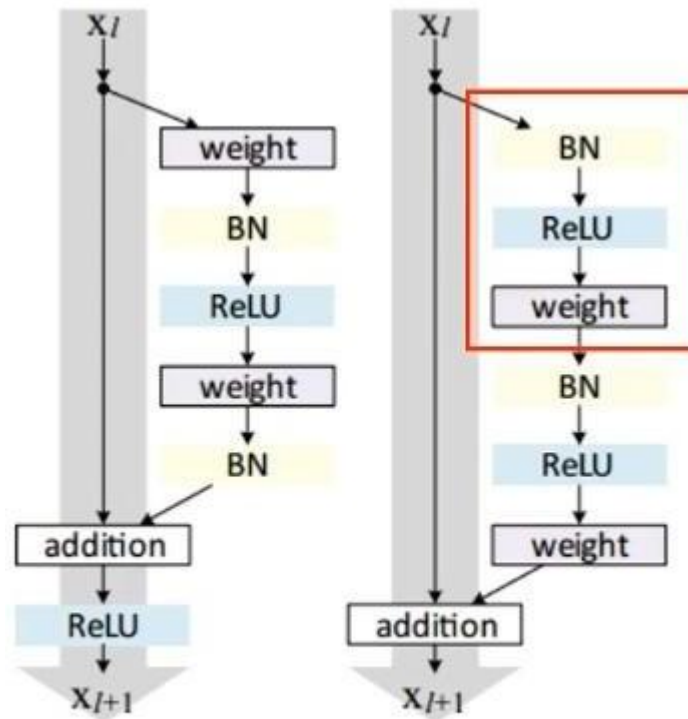
Feature map을 추출한 뒤, activation function을 하던 것이 일반적이었으나, 개선된 구조에서는 "Pre-Activation" (activation function이 먼저 들어감!)



Residual 구조



Identity Mapping

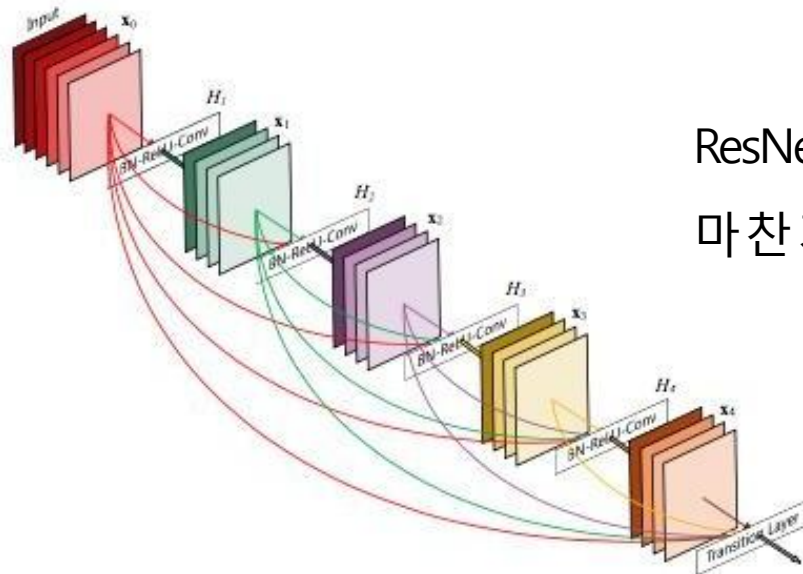


(a) original

(b) proposed

# 7. 심화 CNN

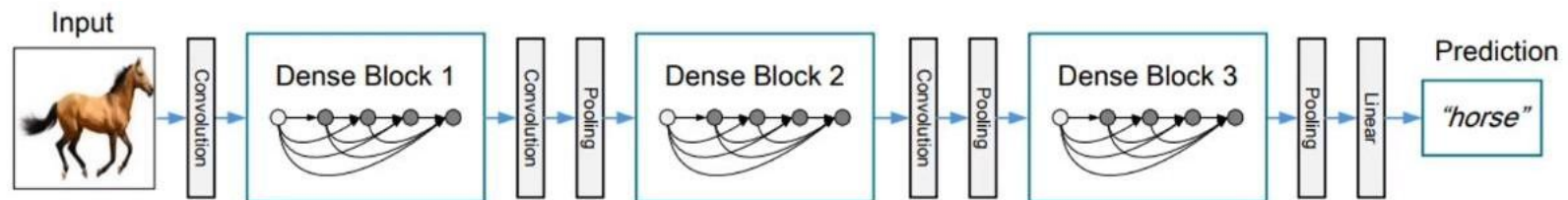
## 6. DenseNet



ResNet 아이디어의 연장선 상!

마찬가지로 Pre-Activation ( BN-ReLU-Conv )

**Figure 1:** A 5-layer dense block with a growth rate of  $k = 4$ . Each layer takes all preceding feature-maps as input.

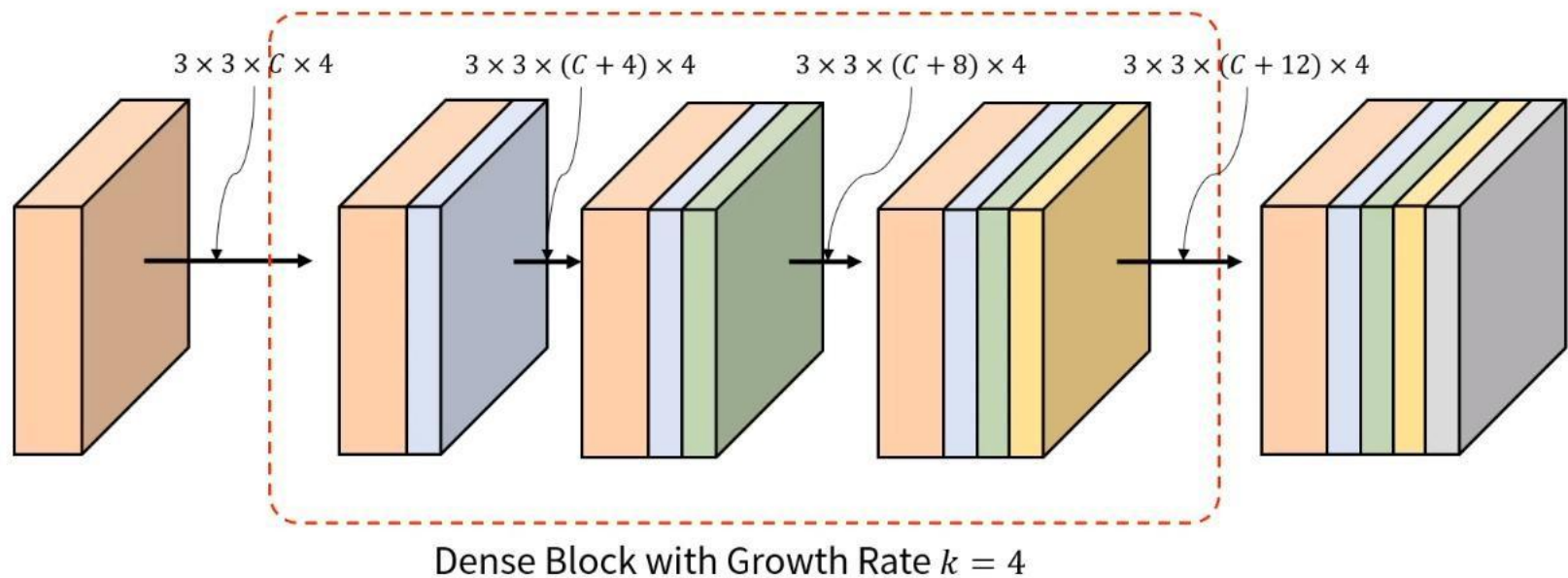


**Figure 2:** A deep DenseNet with three dense blocks. The layers between two adjacent blocks are referred to as transition layers and change feature-map sizes via convolution and pooling.



## 7. 심화 CNN

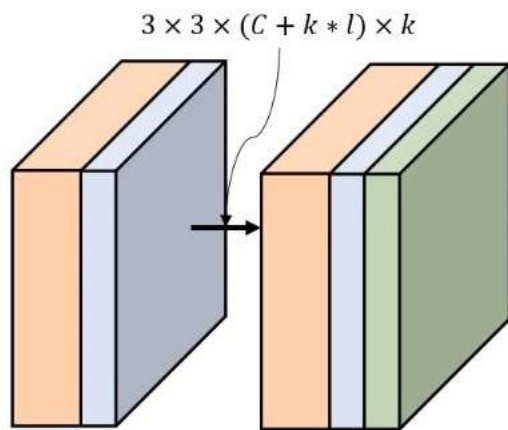
### 6. DenseNet



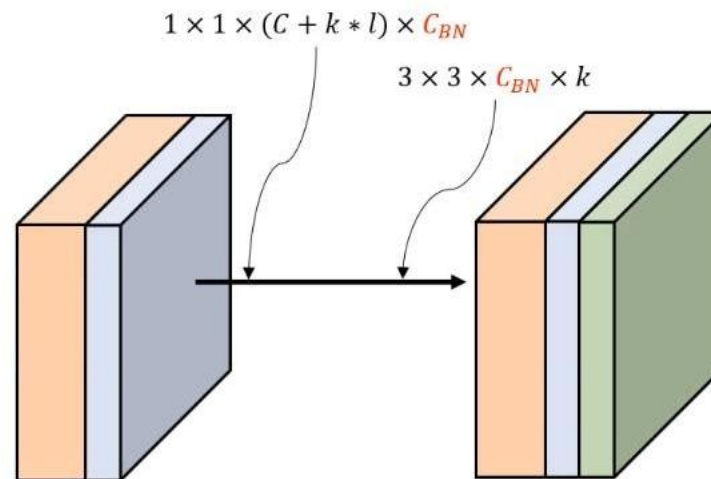
쉽게 말해, "이전 feature map에 누적해서 concatenate!"

# 7. 심화 CNN

## 6. DenseNet



일반 Dense block 계산



Bottleneck 구조

Deep -> 연산량 too much!

그래서 "1x1 Conv (= Bottleneck Layer)" 사용 함

## 8. 실습

