

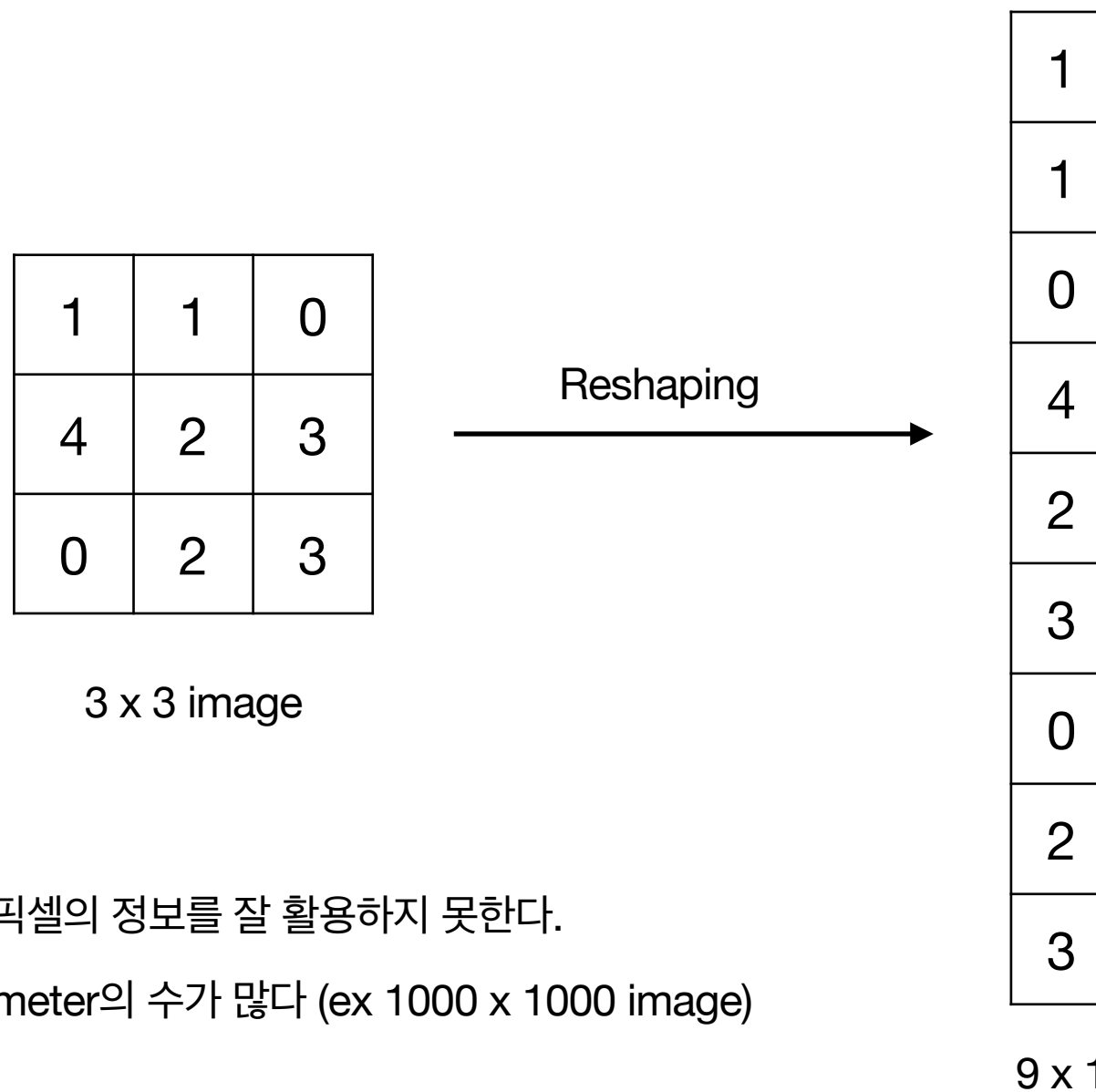
# **Lecture 2**

## **Convolution Neural Network**

**A-Task**

**Yeolkyo Choi**

# Image Treatment using FNN



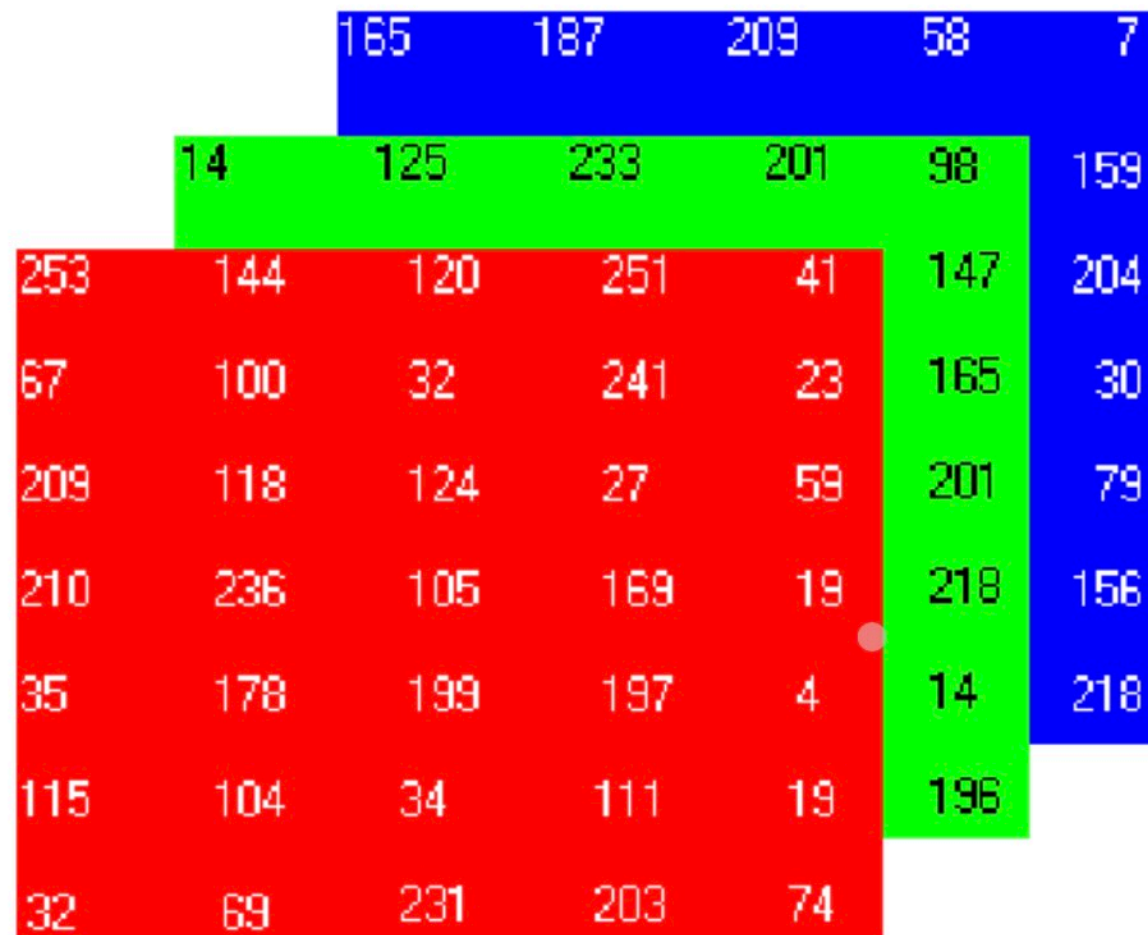
- 이웃 픽셀의 정보를 잘 활용하지 못한다.
- Parameter의 수가 많다 (ex 1000 x 1000 image)

## Convolution Neural Network의 장점

- 데이터가 가지고 있는 spatial 정보를 추출하는데 유리
- 적은 parameter로 이미지 데이터 분석에 용이

# Color Image

- Color image의 경우, 하나의 픽셀이 보통 3개의 색 정보를 지님 (Red, Green, Blue, a.k.a., RGB; called channel or depth = 3)
- 3차원의 array로 저장됨 (m x n 칼라 이미지의 경우, m x n x 3의 3차원 array로 데이터 저장)
- 각 cell은 하나의 색상 정보를 지니며 각각이 하나의 feature (독립변수)의 값이 됨.



# Kernel (Filter) in CNN

- 이미지 정보를 추출하기 위해 filter (kernel)을 적용
- 하나의 filter를 2차원으로 이동시키면서 정보를 추출
- filter의 크기는  $F \times F \times D$

D : 입력 이미지의 depth (channel 수)

F : filter의 가로 또는 세로의 길이

(흑백 이미지는  $D = 1$ , 즉  $F \times F \times 1$  filter 사용)

$W_{11}$	$W_{12}$	$W_{13}$
$W_{21}$	$W_{22}$	$W_{23}$
$W_{31}$	$W_{32}$	$W_{33}$

- Filter의 각 셀은 고유의 가중치를 가짐
- Filter 적용 시, filter가 적용되는 이미지 부분의 색 정보와 filter의 각 가중치 간에 내적 연산(합성곱) 수행  
(<https://setosa.io/ev/image-kernels/>)

1	10	122	143	155
1	20	124	142	156
2	22	110	130	150
3	24	108	122	140
4	18	11	80	100

$$w_{11} \times 1 + w_{12} \times 10 \dots + w_{33} \times 110$$

1	10	122	143	155
1	20	124	142	156
2	22	110	130	150
3	24	108	122	140
4	18	11	80	100

$$w_{11} \times 10 + w_{12} \times 122 \dots + w_{33} \times 130$$

■ ■ ■

1	10	122	143	155
1	20	124	142	156
2	22	110	130	150
3	24	108	122	140
4	18	11	80	100

$$w_{11} \times 110 + w_{12} \times 130 \dots + w_{33} \times 100$$

# Activation Map

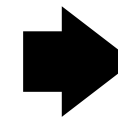
1	10	122	143	155
1	20	124	142	156
2	22	110	130	150
3	24	108	122	140
4	18	11	80	100

 $\otimes$ 

$W_{11}$	$W_{12}$	$W_{13}$
$W_{21}$	$W_{22}$	$W_{23}$
$W_{31}$	$W_{32}$	$W_{33}$

 $=$ 

$Z_{11}$	$Z_{12}$	$Z_{13}$
$Z_{21}$	$Z_{22}$	$Z_{23}$
$Z_{31}$	$Z_{32}$	$Z_{33}$



$f(z_{11})$	$f(z_{12})$	$f(z_{13})$
$f(z_{21})$	$f(z_{22})$	$f(z_{23})$
$f(z_{31})$	$f(z_{32})$	$f(z_{33})$

filter는 하나의 bias 값을 갖는다.

Activation map

- activation function을 적용하여 출력되는 값들을 모아놓은 것
- Feature map 또는 response map 이라고도 함

# Stride in CNN

stride는 filter를 한번에 옆으로 (혹은 아래로) 얼마나 움직이는지를 결정 (이는 hyperparameter)

image = 7 x 7 / filter = 3 x 3 / stride = 1

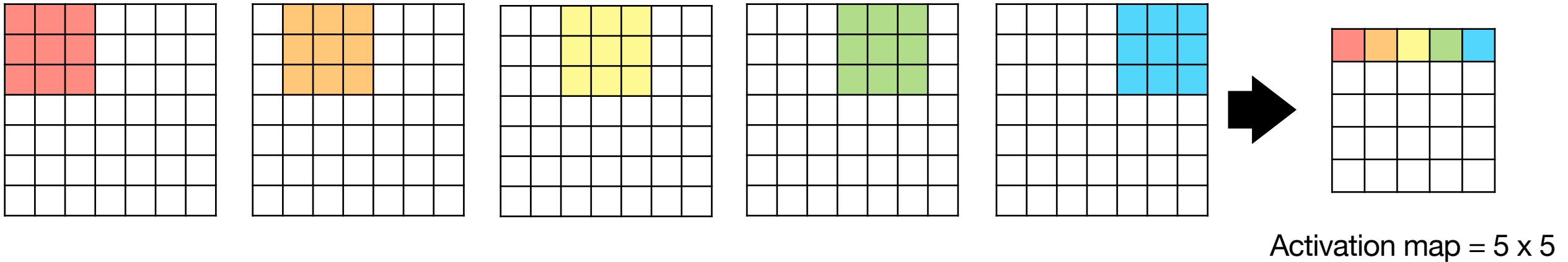
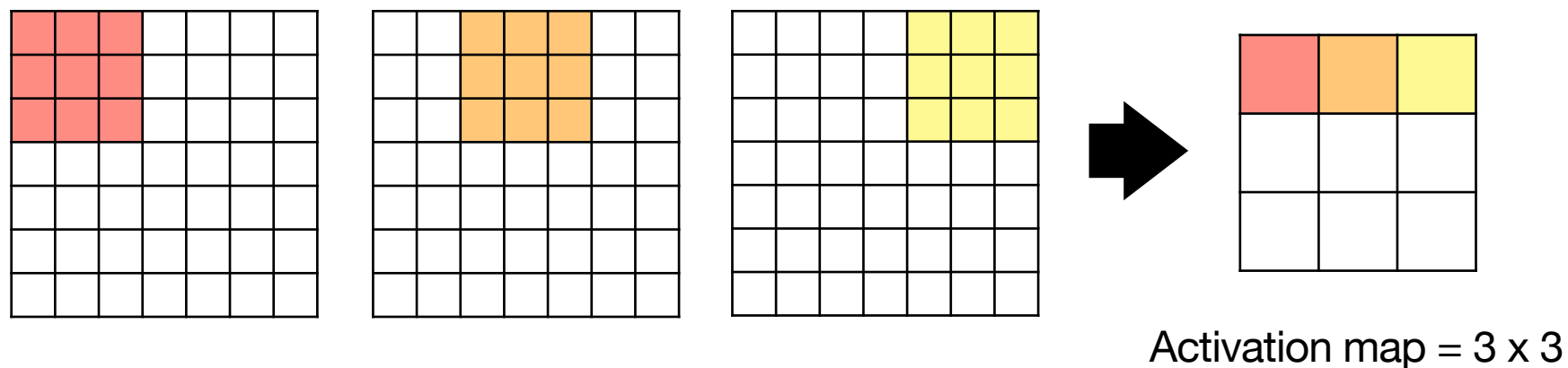


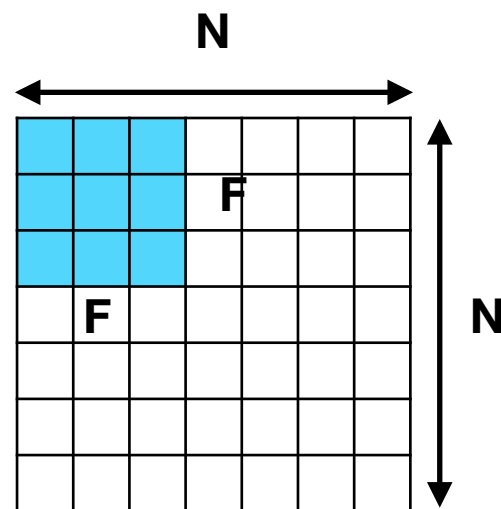
image = 7 x 7 / filter = 3 x 3 / stride = 2



Mainly for downsampling

- parameter 수 감소
- 학습속도 개선
- 정보 손실, 성능 저하

# Padding in CNN



**Output size:**

$$(N-F) / S + 1$$

e.g.  $N = 7, F = 3$

$$\text{stride } 1 \rightarrow (7-3)/1 + 1 = 5$$

$$\text{stride } 2 \rightarrow (7-3)/2 + 1 = 3$$

$$\text{stride } 3 \rightarrow (7-3)/3 + 1 = 2.33$$

In practice: Common to zero pad the border

0	0	0	0	0	0	0	0
0							0
0							0
0							0
0							0
0							0
0							0
0	0	0	0	0	0	0	0

e.g. input = 7 x 7

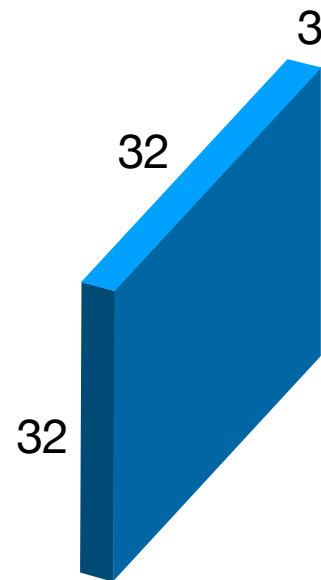
3 x 3 filter, applied with stride 1

pad with 1 pixel border → what is the output?

$$(N + 2P - F)/S + 1 \Rightarrow 7 \times 7 \text{ output!}$$

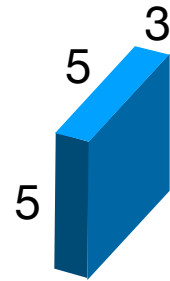
In general, common to see CONV layers with stride 1, filters of size  $F \times F$ , and zero-padding with  $(F-1)/2$  (will preserve size spatially)

# Filter with Multichannel



Color Image (RGB)

32 x 32 pixel with depth (channel) = 3



5 x 5 x 3 filter

**How many weights?**

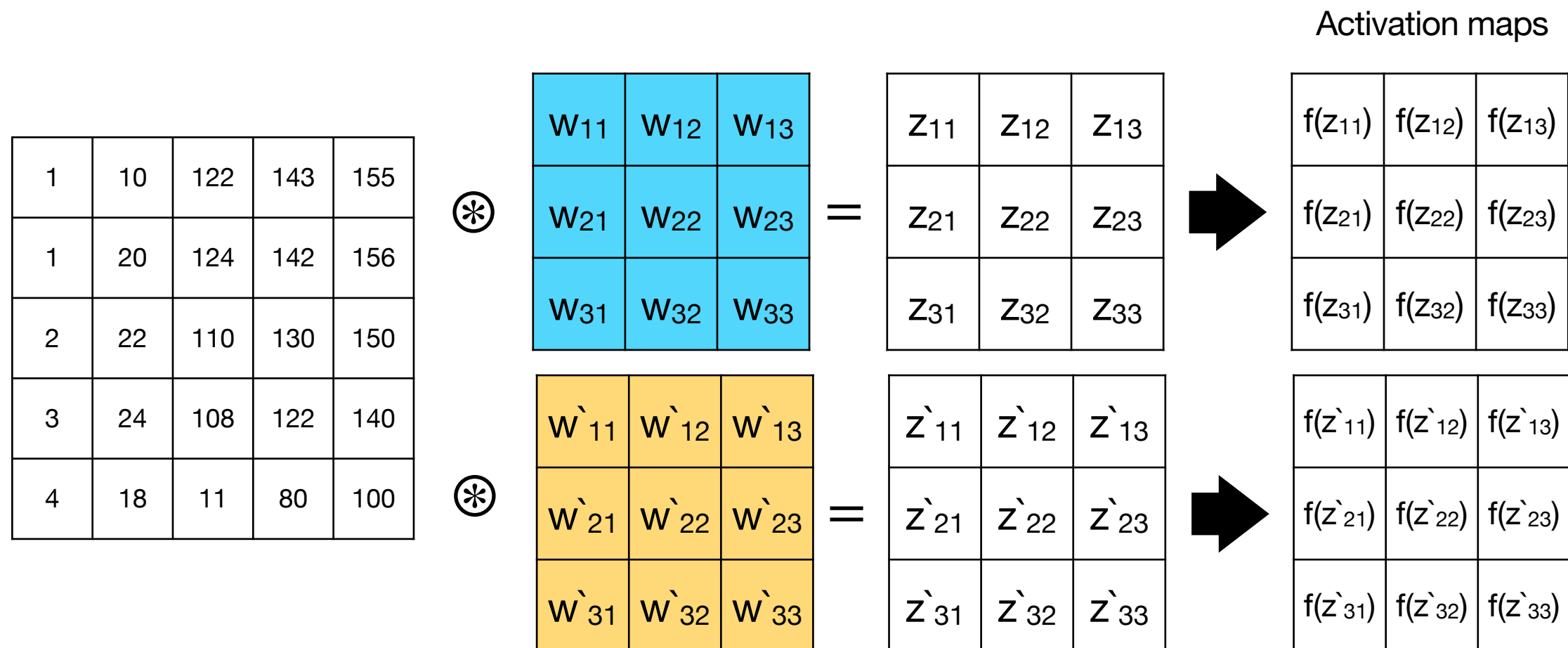
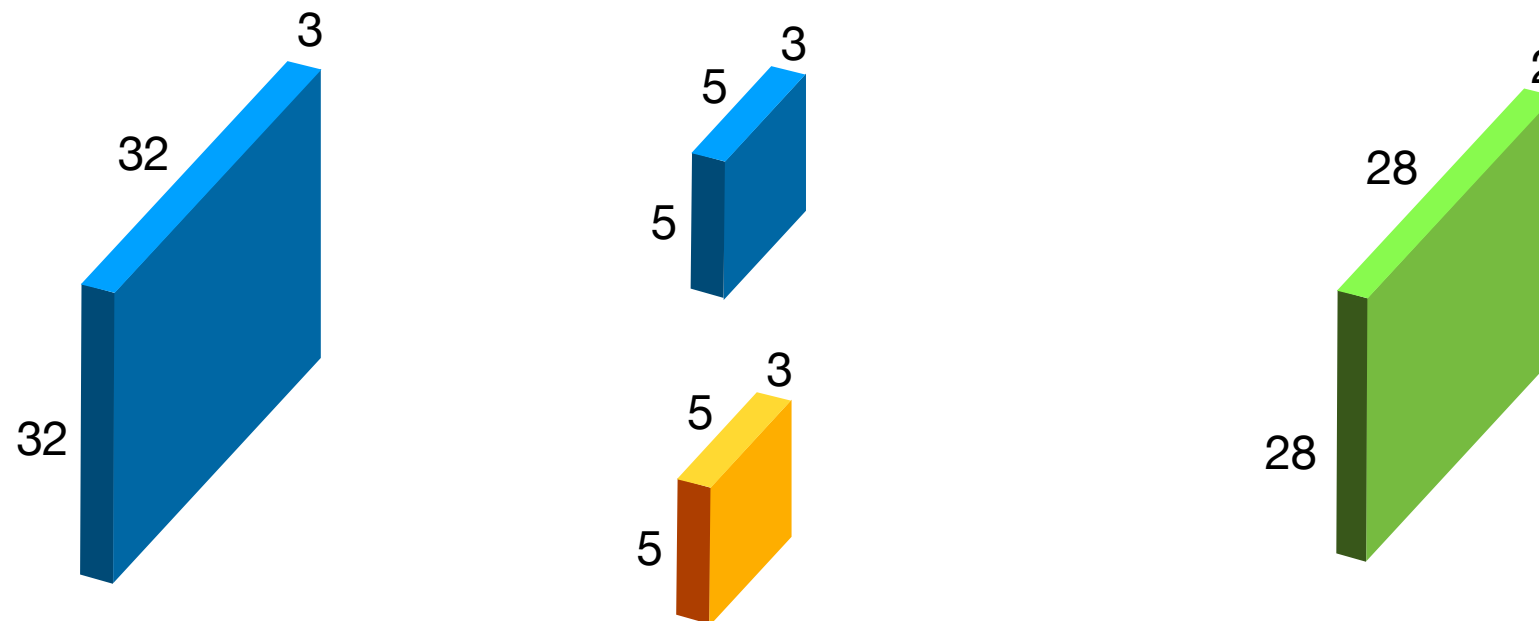
**How many bias?**

**Activation map ?**

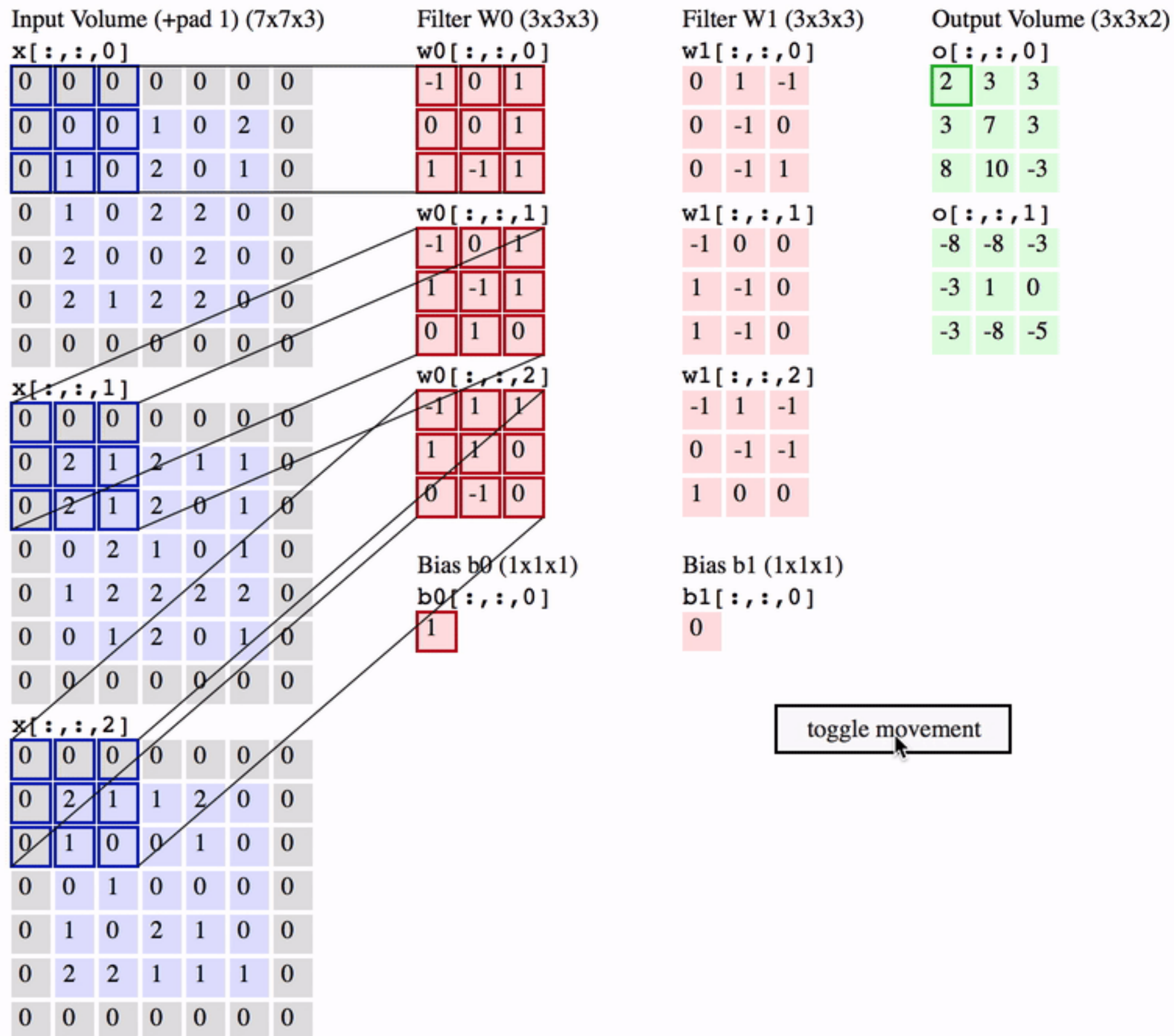


# Filter with Multichannel

하나의 이미지에 여러개의 필터를 적용하는 것이 가능

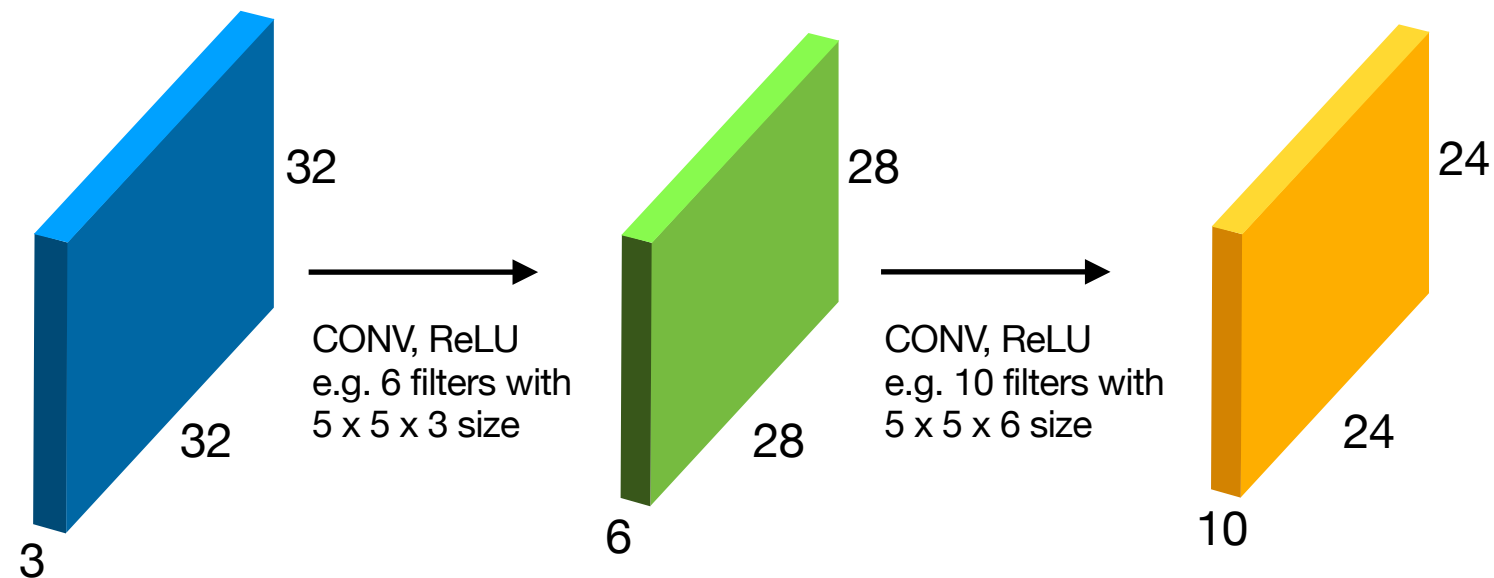


# How CNN Works



# How CNN Works

- 합성곱의 반복적용



위 그림은 32x32x3 color image에 5x5x3 크기의 filter 6개를 적용해서 나온 activation map (28x28x6)에 10개의 5x5x6 filter를 다시 한번 적용하여 24x24x10 크기의 activation map을 얻는 경우를 나타낸다.



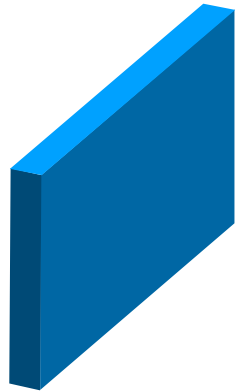
뒤로 갈수록 object의 윤곽정보를 포함함

# Pooling Layer

- makes the representations smaller and more manageable
- operates over each activation map independently

MAX POOLING

224x224x64



112x112x64



downsampling

Single depth slice

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

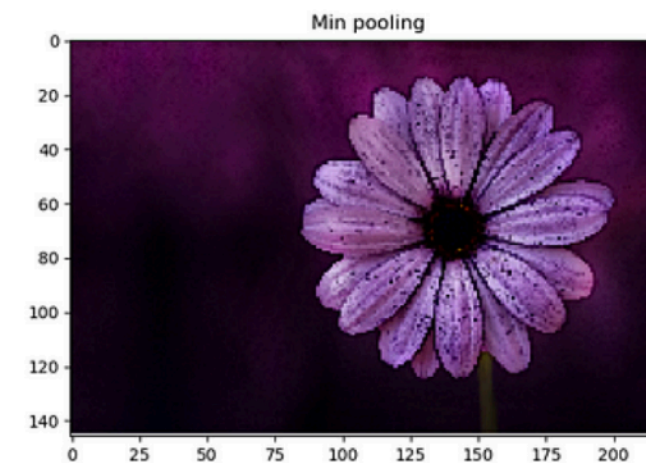
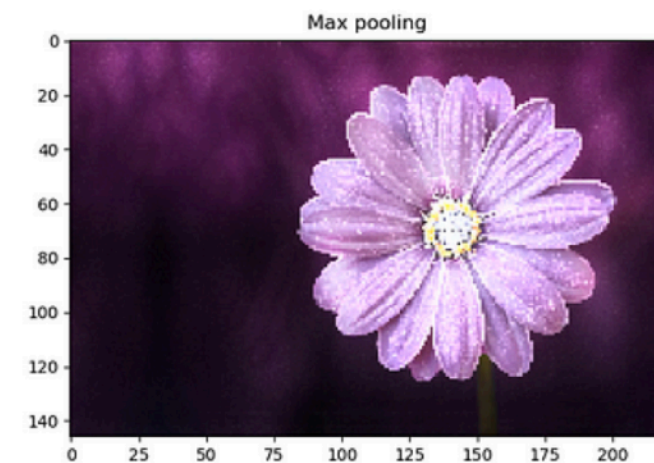
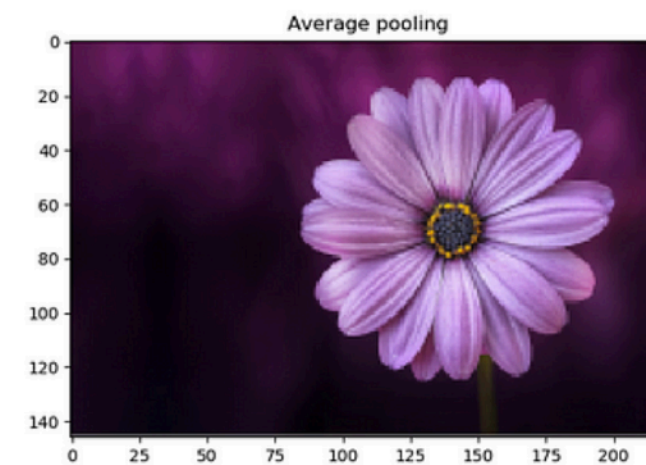
Common setting

F=2, S=2

F=3, S=2

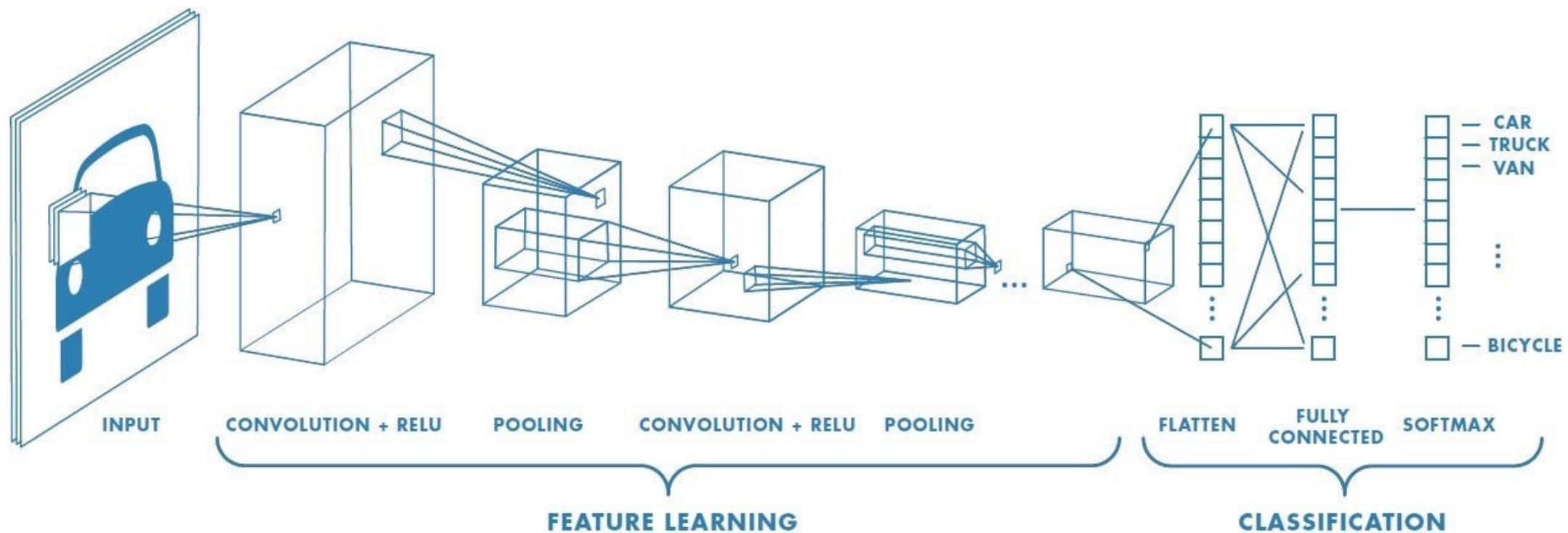
max pool with 2x2 filters  
and stride 2

6	8
3	4





# Fully Connected Layer



Fully Connected Layer

32x32x3 image > stretch to 3072 x 1

3072

$Wx$

10 x 3072



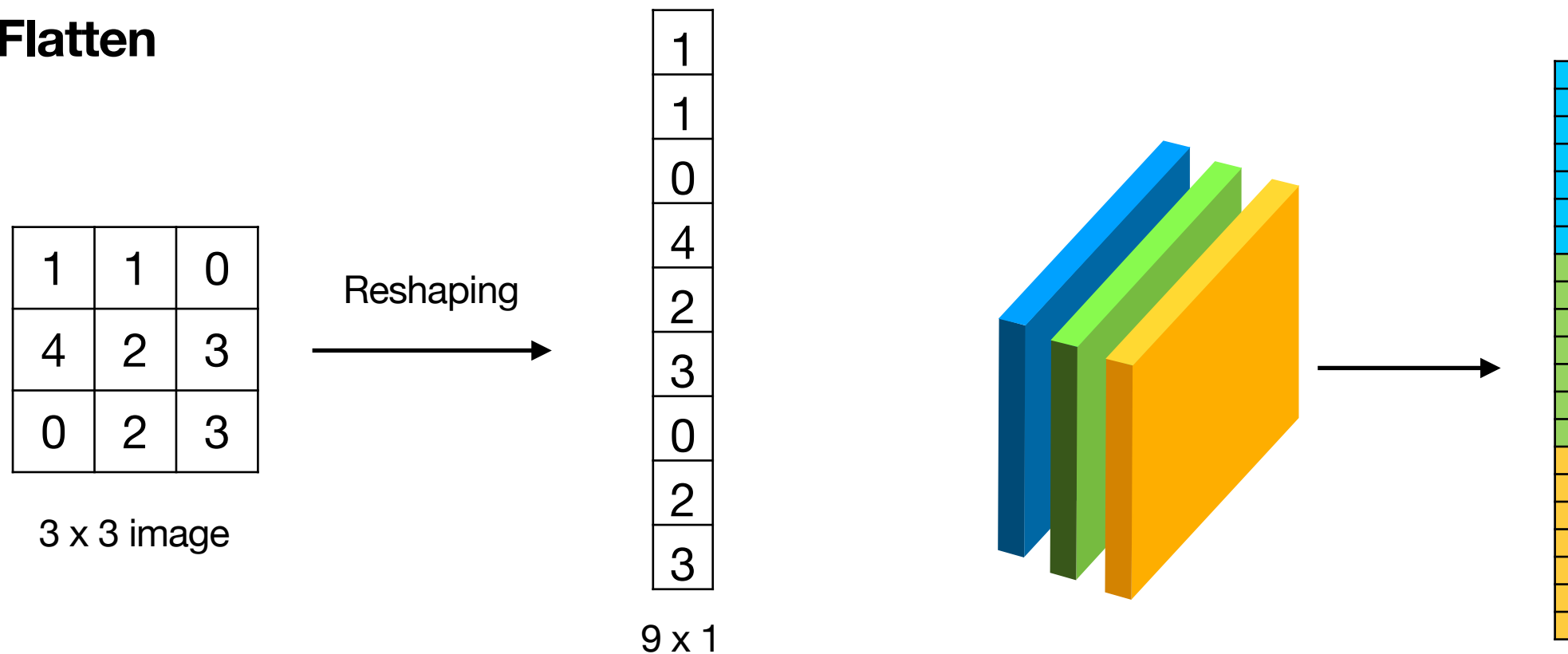
activation

1 number:

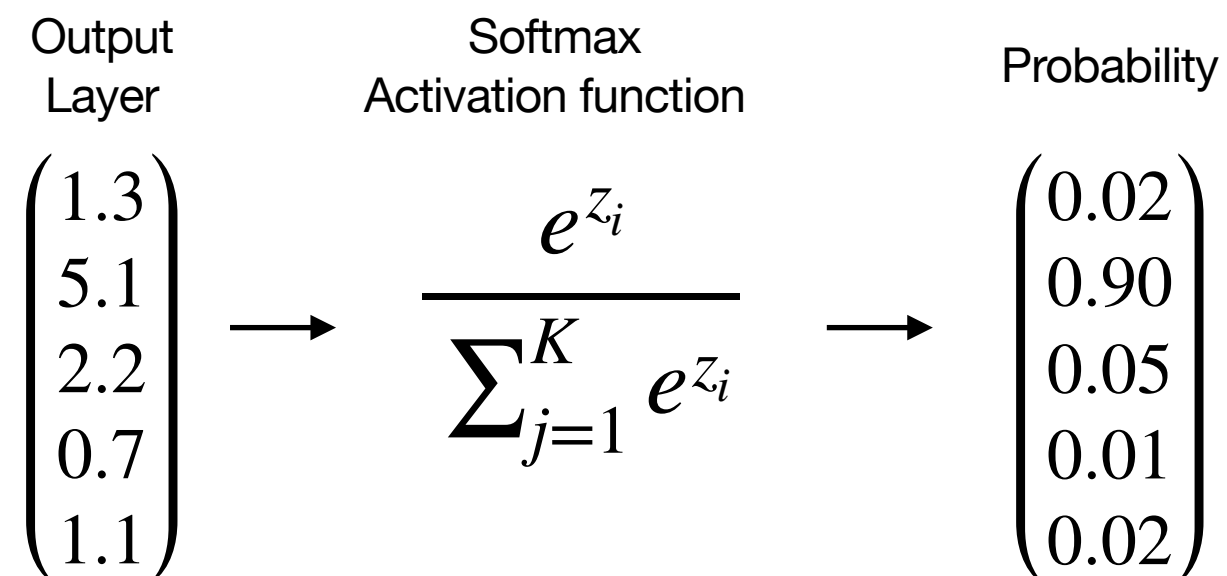
the result of taking a dot product between  $W$  and the input  
(a 3072 dimensional dot product)

# Flatten and Softmax

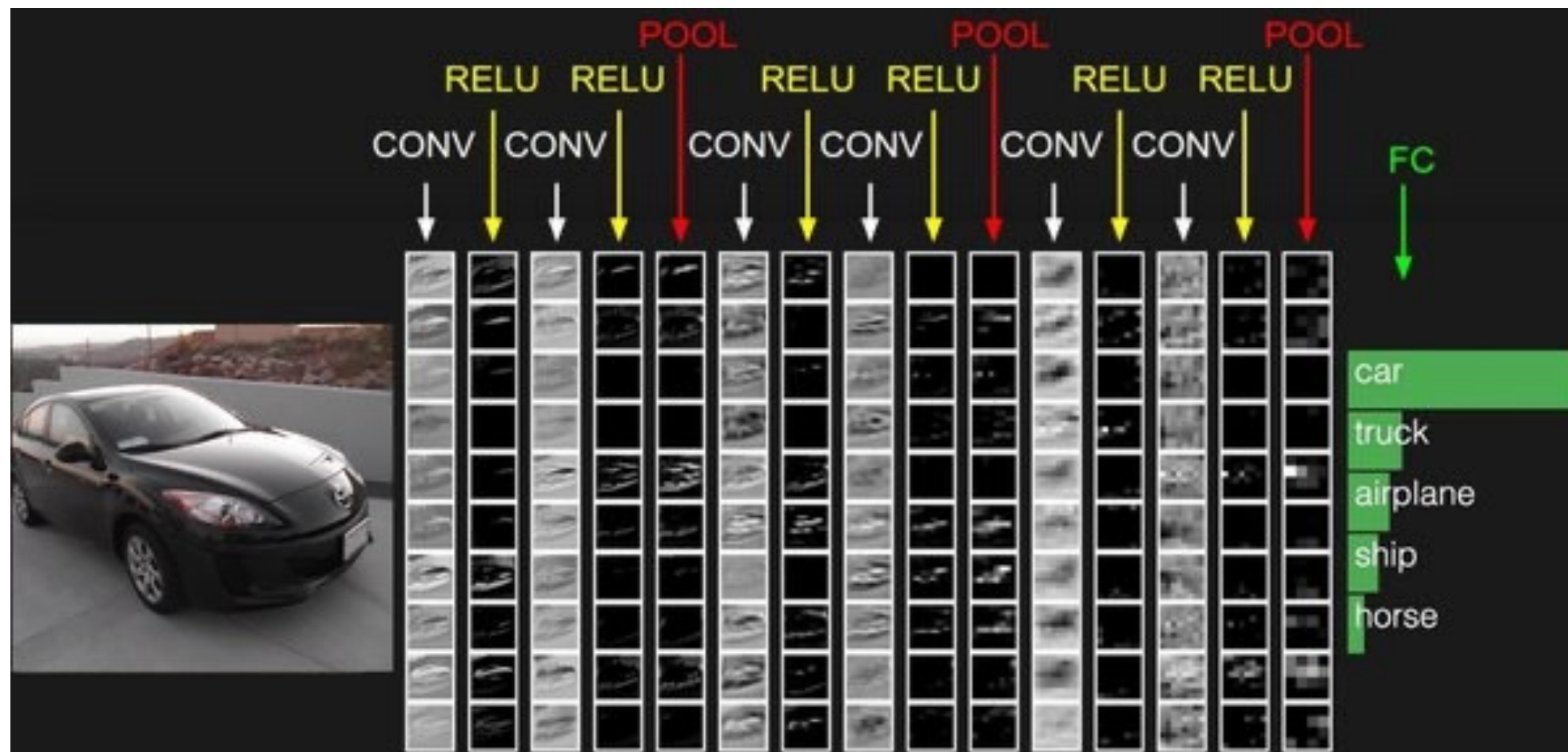
- Flatten**



- Softmax**



# CNN Summary



- CNNs stack of CONV, POOL, FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Typical architectures look like  $[(\text{CONV-RELU}) \times N - \text{POOL}] \times M - (\text{FC-RELU}) \times K - \text{SOFTMAX}$  where  $N$  is usually up to  $\sim 5$ ,  $M$  is large,  $0 \leq K \leq 2$
- Recent advances such as ResNet/GoogLeNet challenge this paradigm