

Convolutional Neural Net.

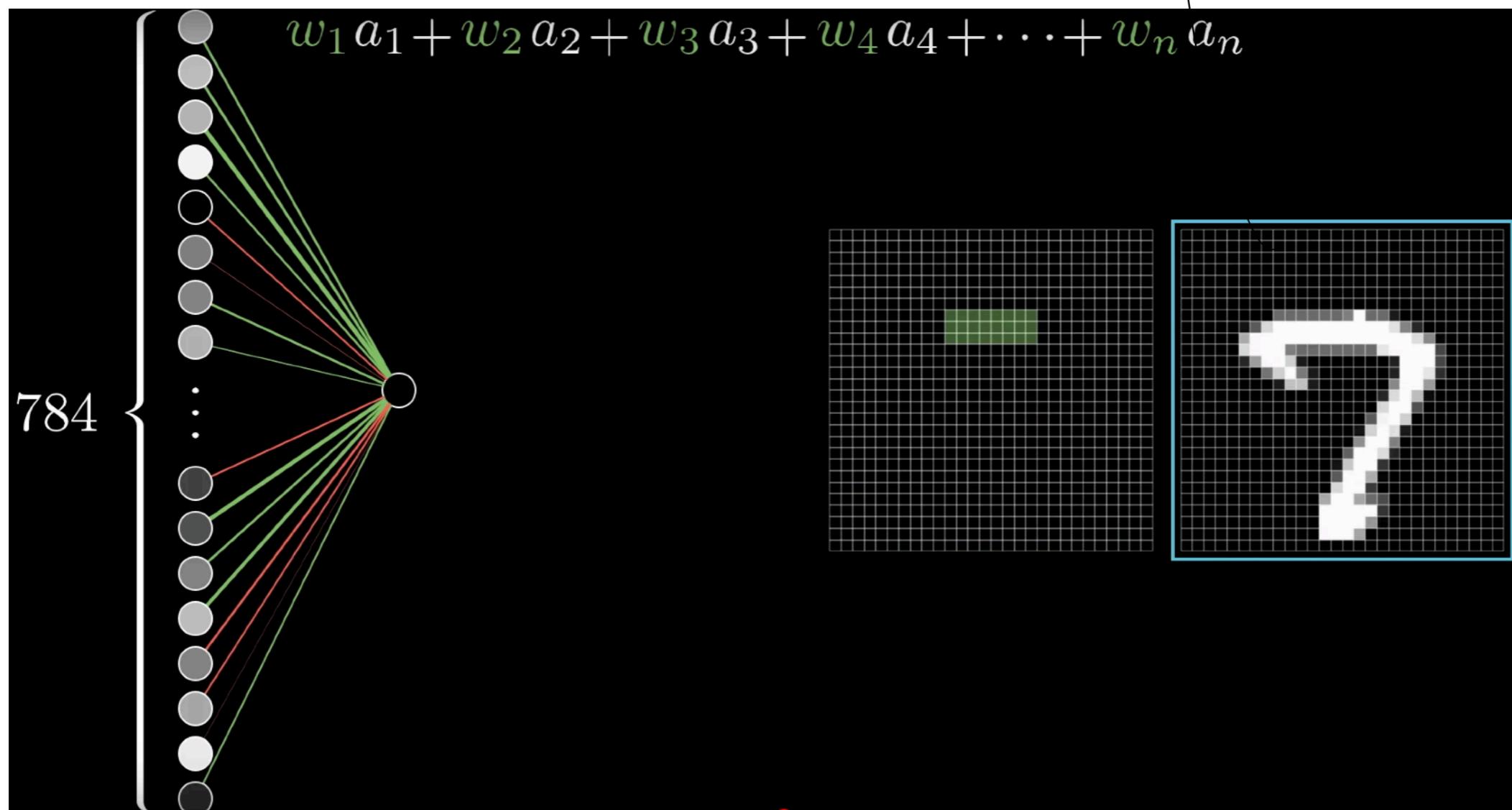
Seyoung Yun

CNN

- http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture5.pdf
- http://www.di.ens.fr/~lelarge/dldiy/slides/lecture_6/index.html#80

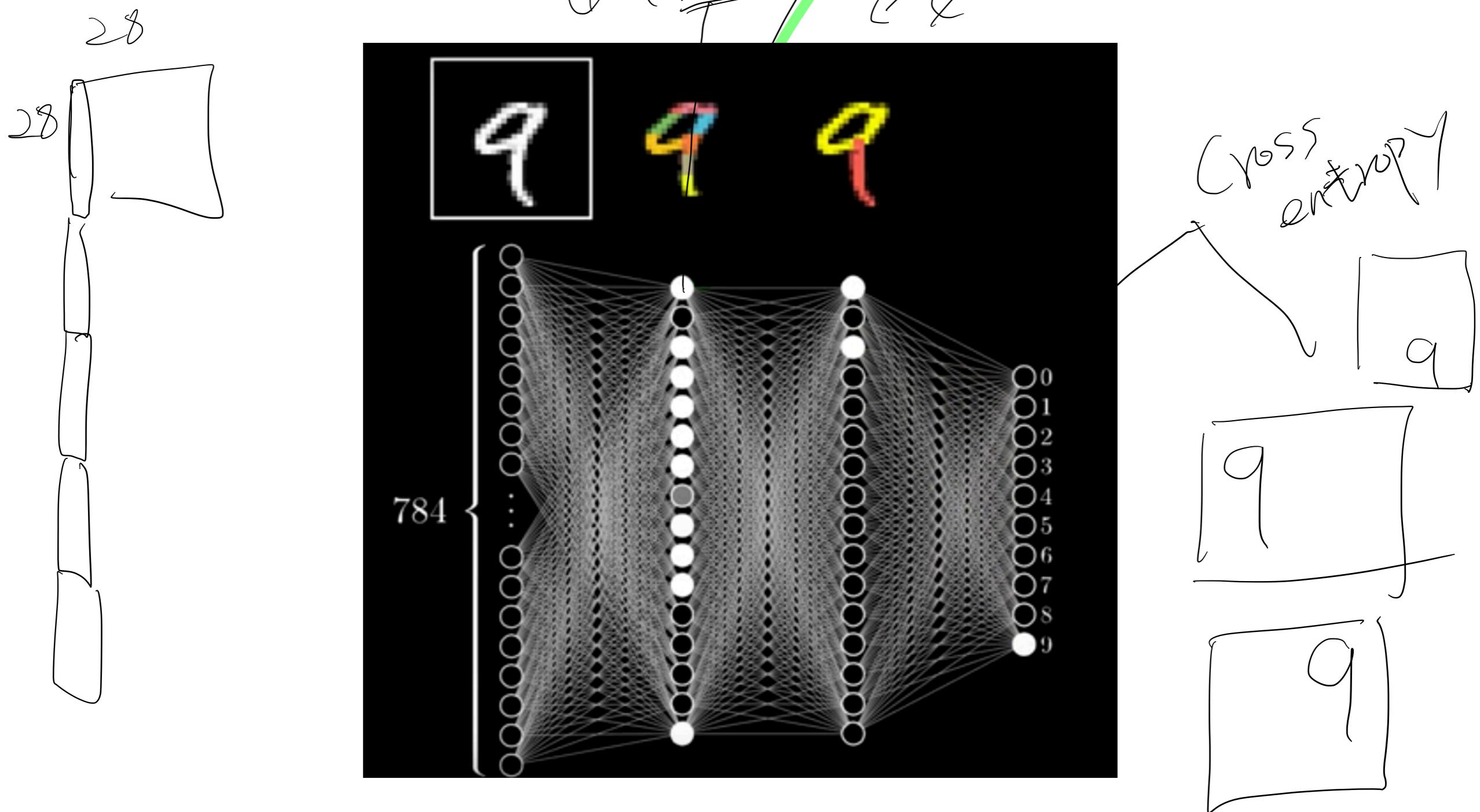
Why we need CNN?

- With a single neuron?

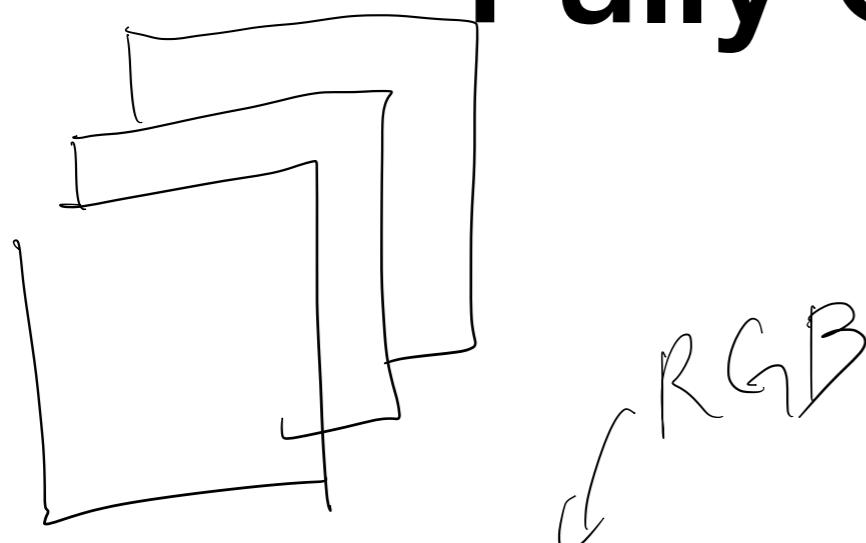


Why we need CNN?

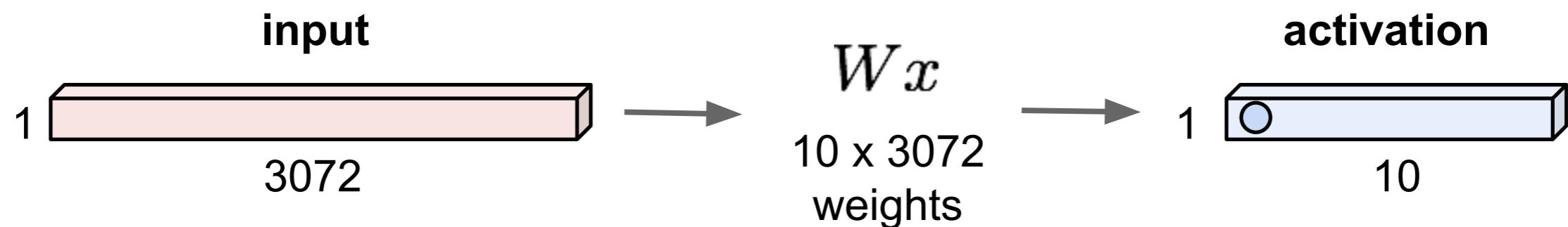
- With FNN?



Fully Connected Layer

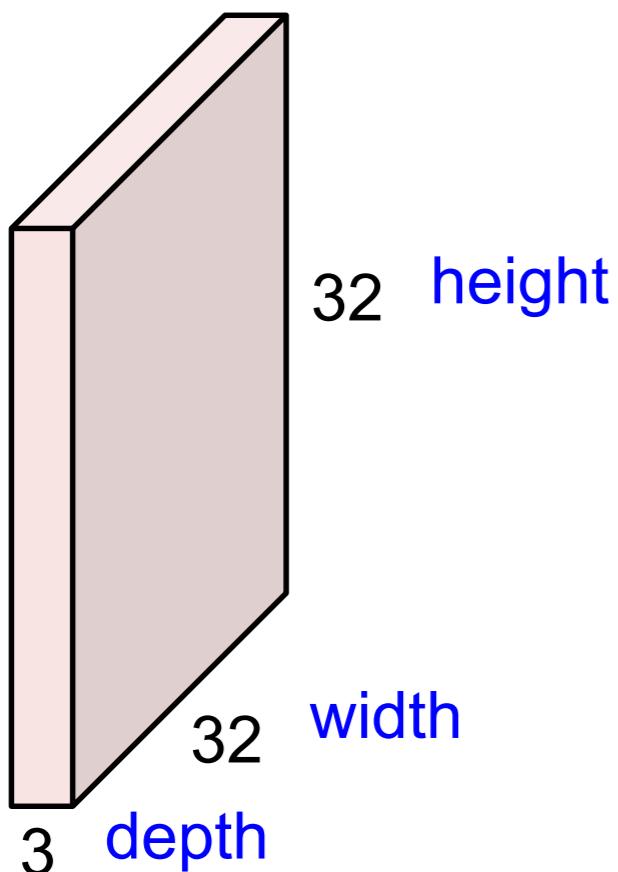


32x32x3 image -> stretch to 3072 x 1

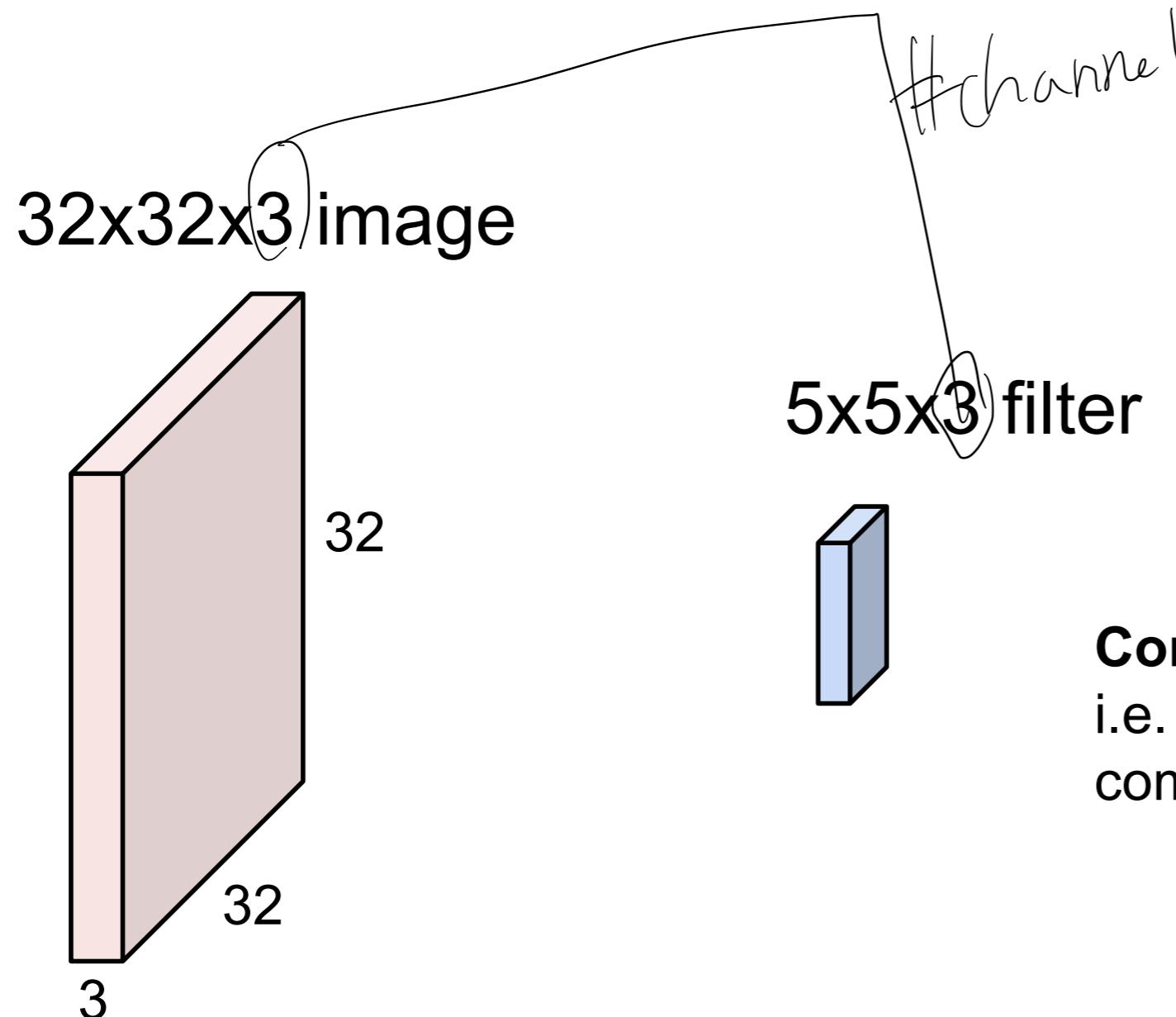


Convolution Layer

32x32x3 image -> preserve spatial structure

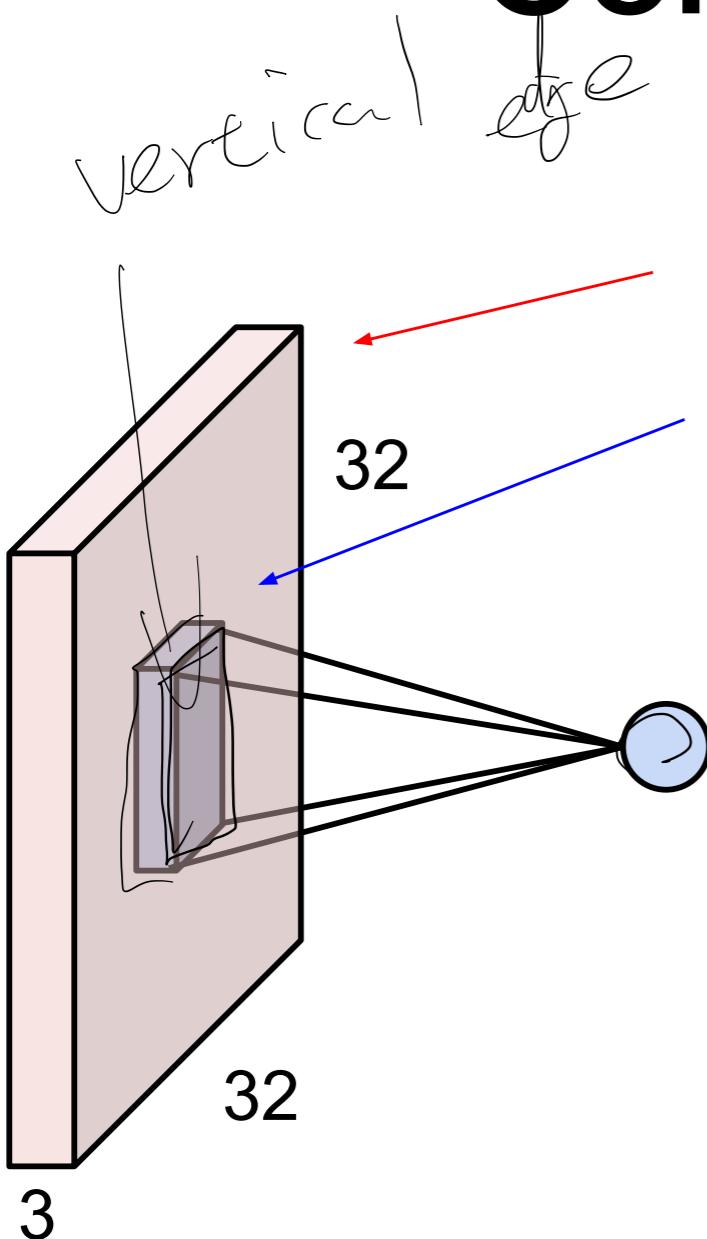


Convolution Layer



Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

Convolution Layer

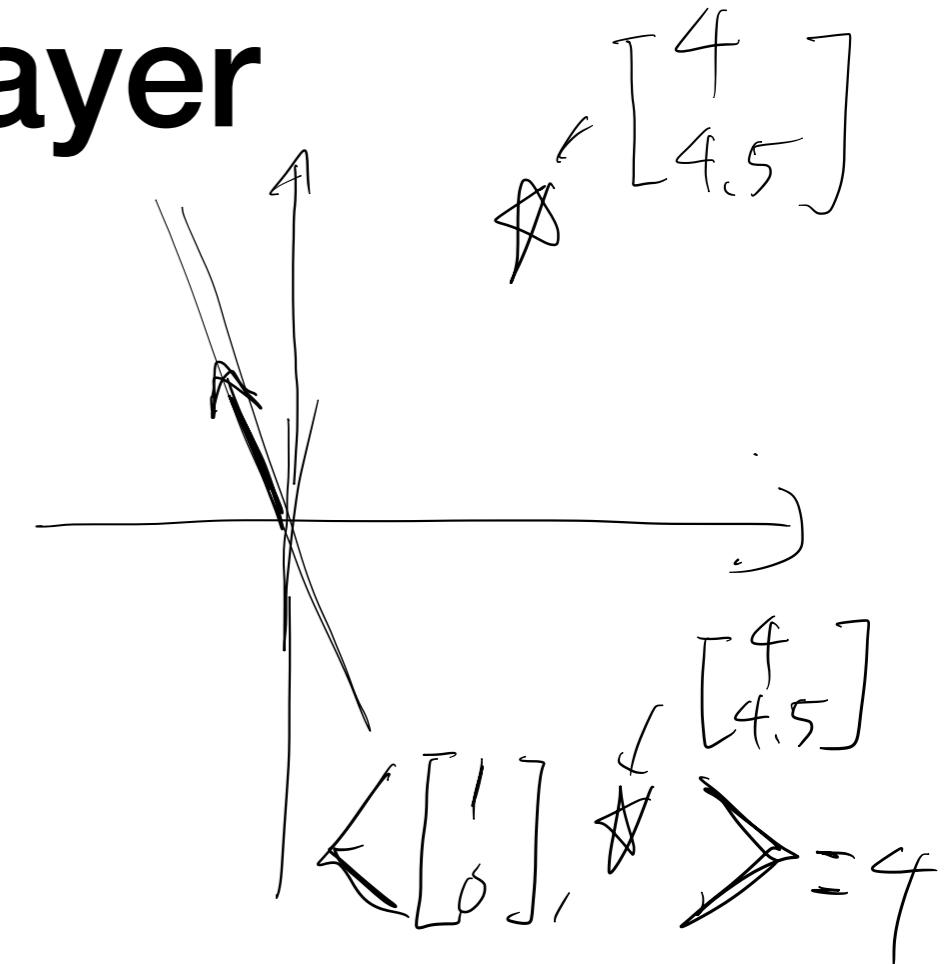


32x32x3 image
5x5x3 filter w

1 number:

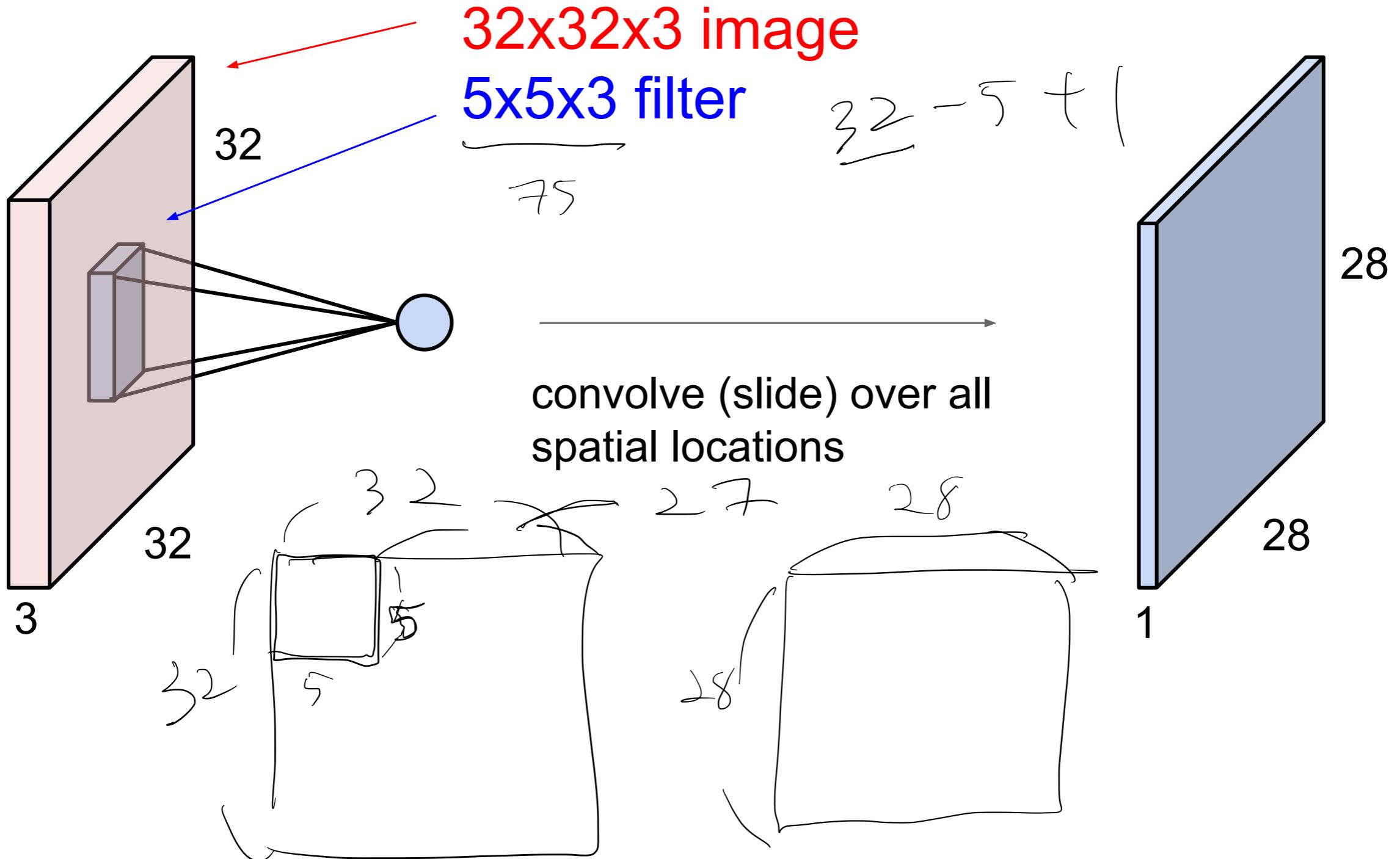
the result of taking a dot product between the filter and a small 5x5x3 chunk of the image
(i.e. $5 \times 5 \times 3 = 75$ -dimensional dot product + bias)

$$w^T x + b$$



$$(32 \times 32 \times 3) \times (28 \times 28) = ?$$

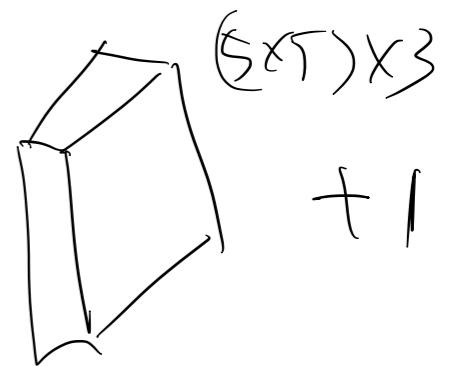
Convolution Layer



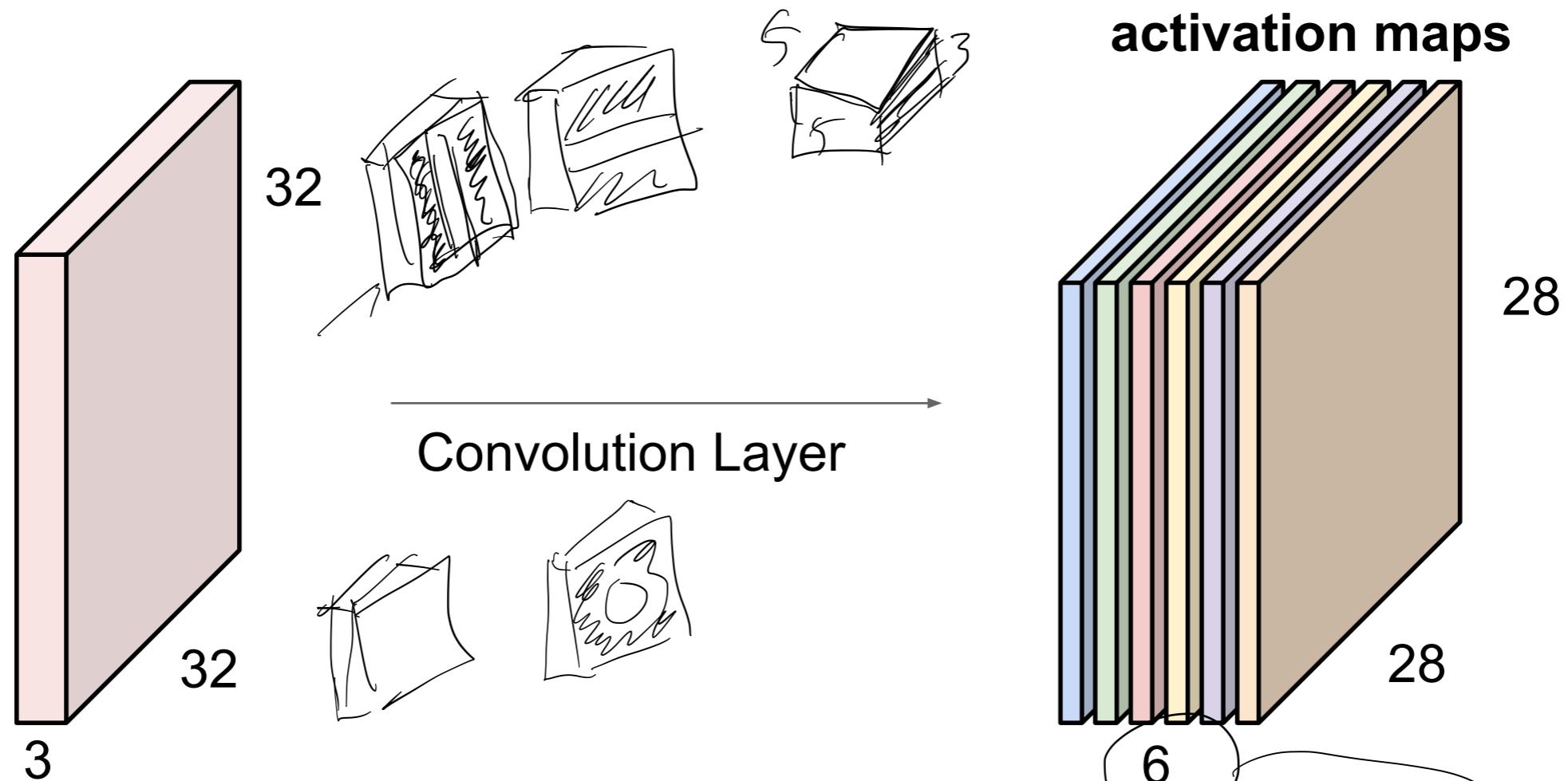
Convolution Layer



Convolution Layer



For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a “new image” of size $28 \times 28 \times 6$!

channel

What is Convolution?

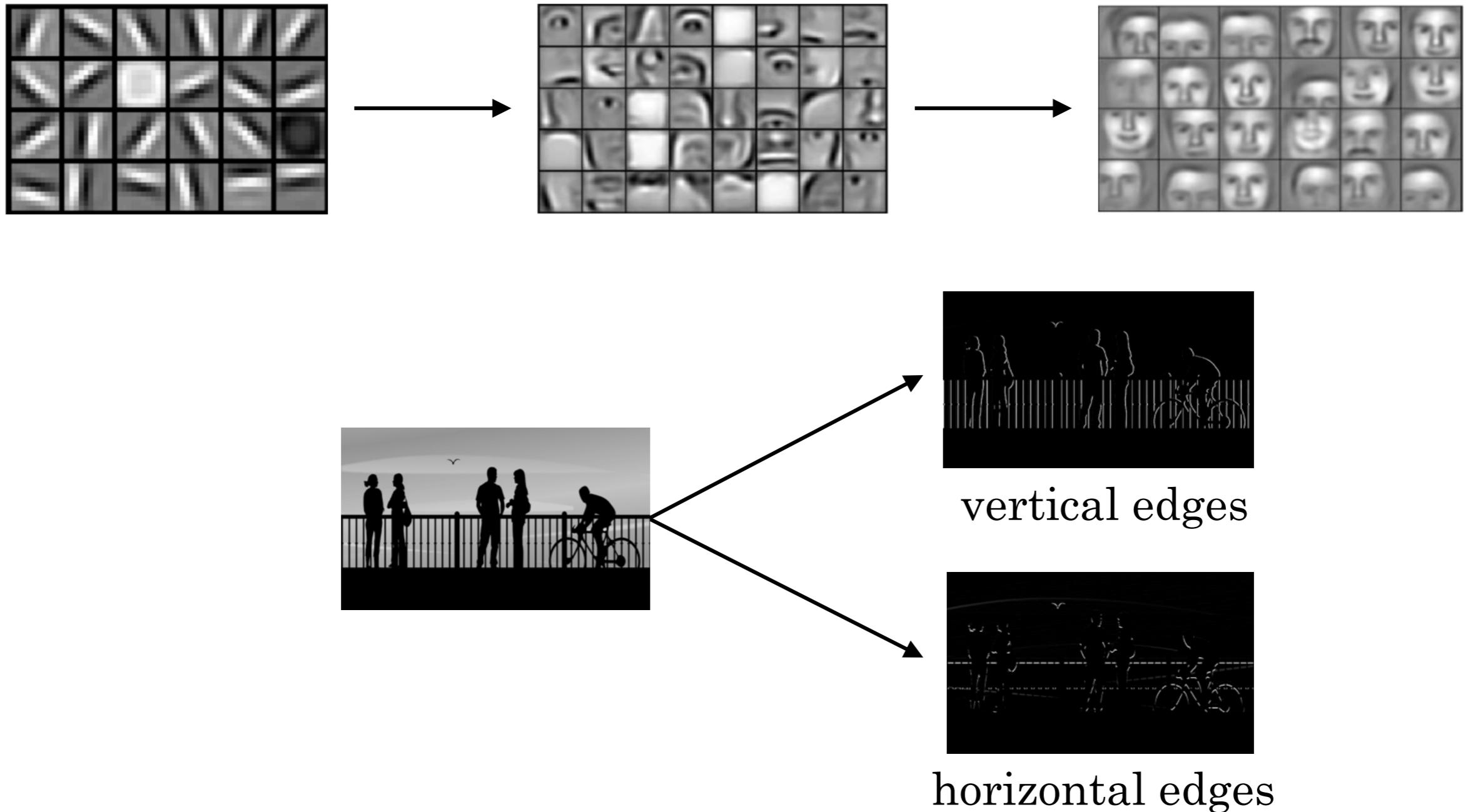
- Convolution

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(m, n)K(i - m, j - n).$$

- Cross-correlation

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(i + m, j + n)K(m, n).$$

Edge Detection



Edge Detection

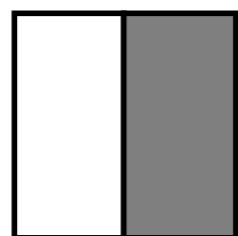
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0

*

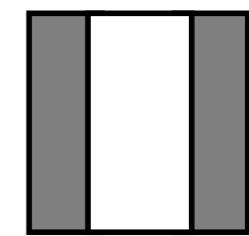
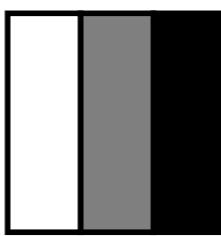
1	0	-1
1	0	-1
1	0	-1

=

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



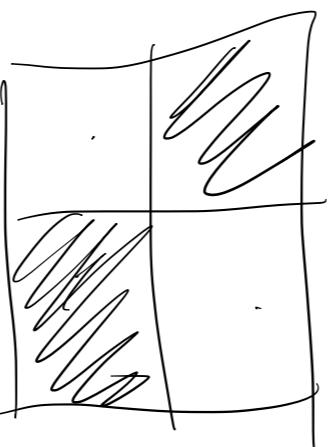
*



Edge Detection

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

Vertical



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

Horizontal

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

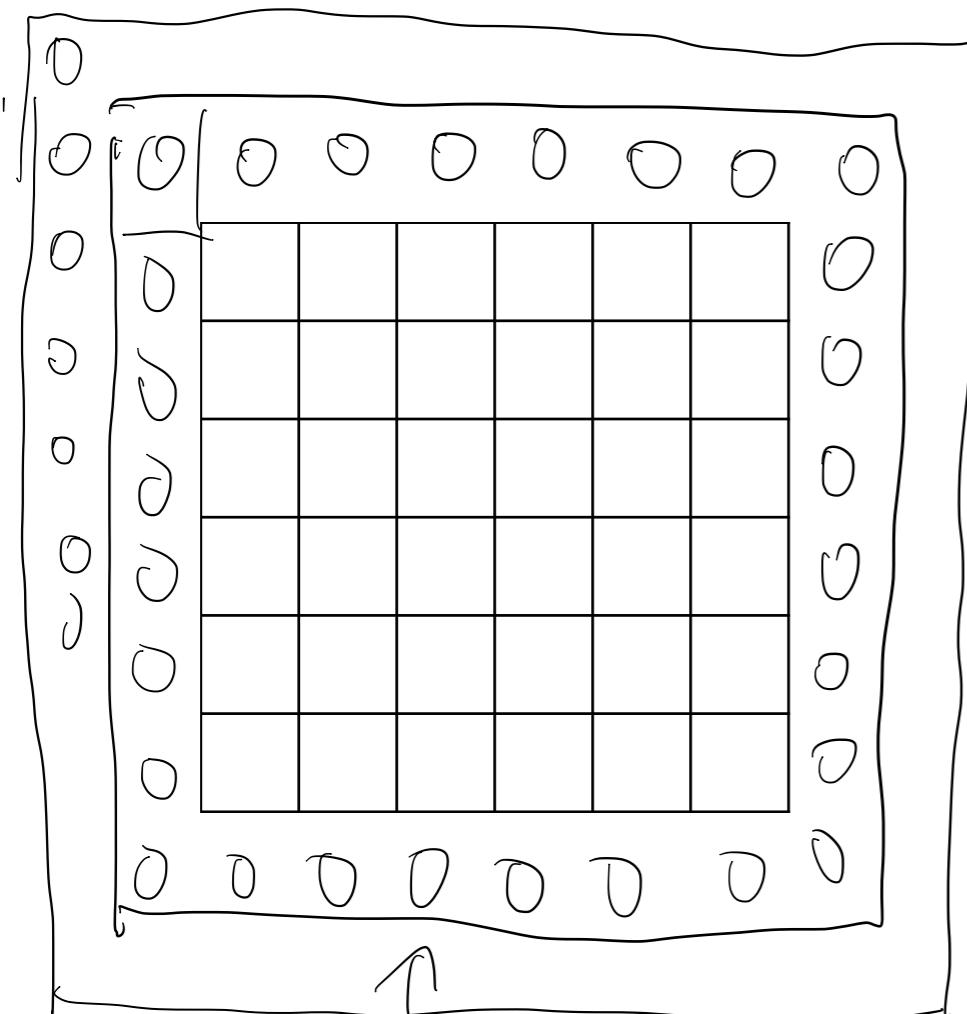
*

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

=

0	0	0	0
30	10	-10	-30
30	10	-10	-30
0	0	0	0

Padding



$$6 \times 6 \rightarrow 8 \times 8$$

$$\rightarrow 10 \times 10$$

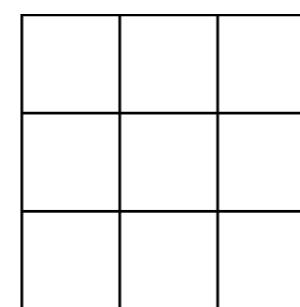
$d \times d$

P

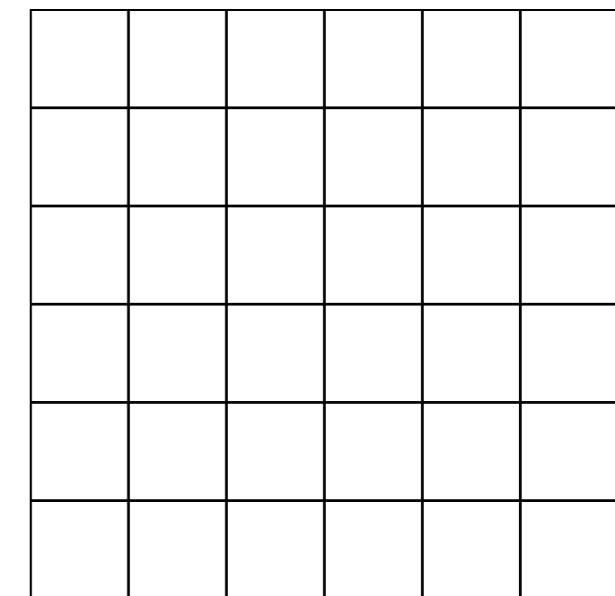
padding size

f
 P

filter size

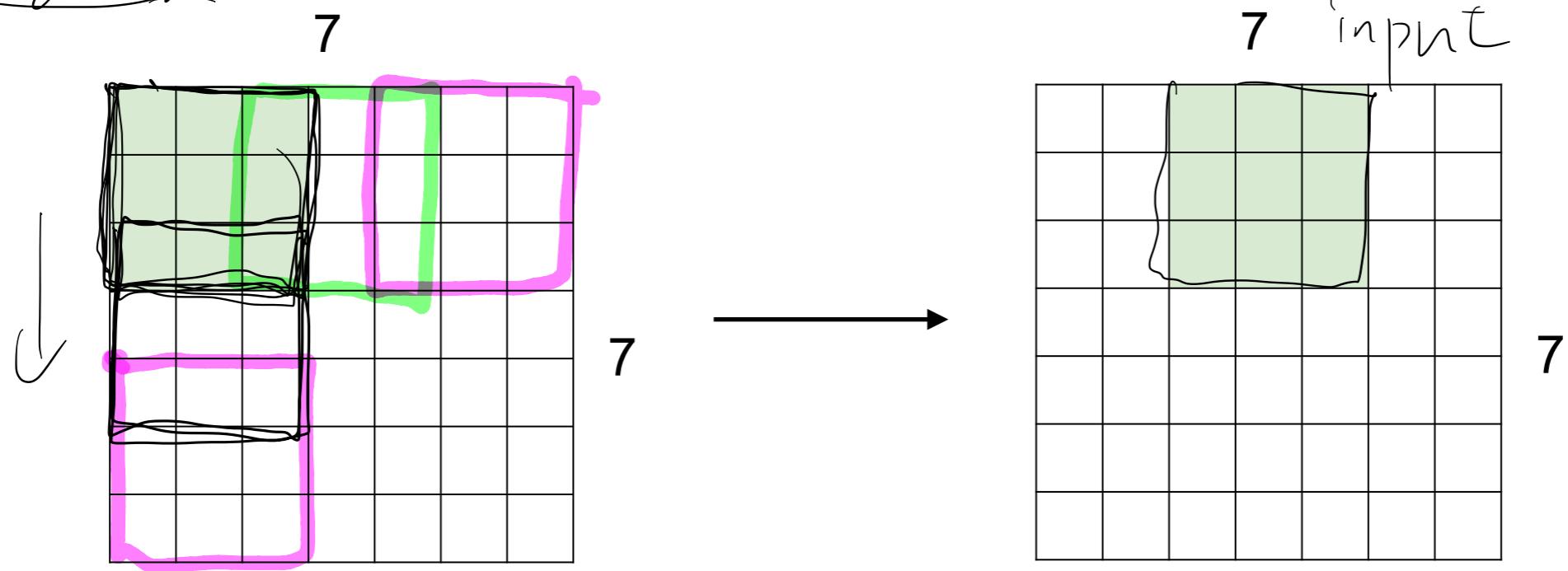
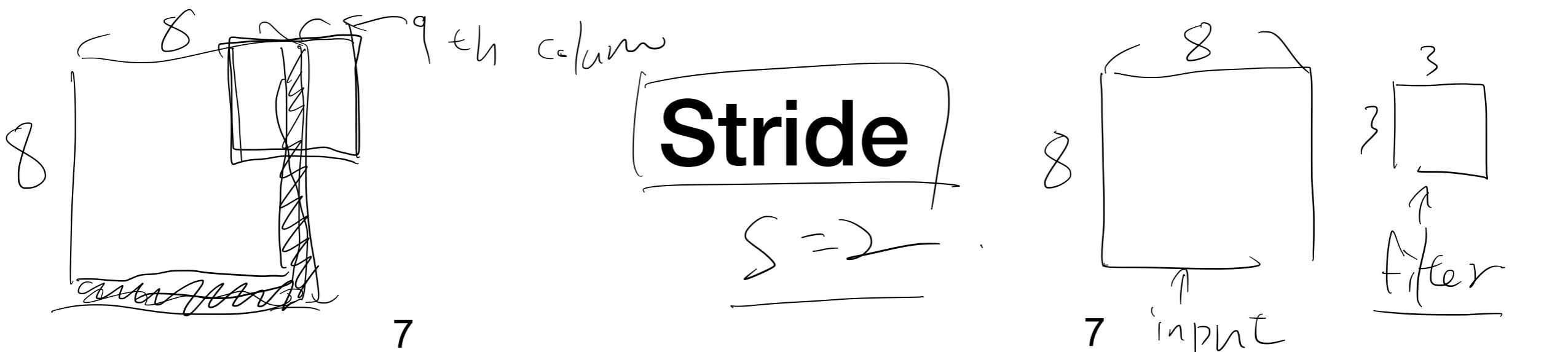


=



↑

$$(d+2P)-f + 1$$



$$\Rightarrow \begin{matrix} 3 \\ \parallel \\ 3 \end{matrix} \in \left[\frac{(d+2 \cdot p) - f}{s} + 1 \right]$$

Diagram showing the formula for calculating the output size. It features a 3x3 input matrix with a circled '3' at the top-left. A handwritten note '=>' is to its left. To the right is a 3x3 filter with a circled '3' at the top-left. Below the filter is a bracket with the formula $\left[\frac{(d+2 \cdot p) - f}{s} + 1 \right]$. The variable 's' is written below the bracket, and a plus sign '+' is followed by another bracket.

Example

input : 100×100

filter : 7×7

padding : 1

stride : 2. 95

Output ?

$$\left\lfloor \frac{100 + 2 - 7}{2} \right\rfloor + 1$$

$$4 \quad 7 \quad + 1$$

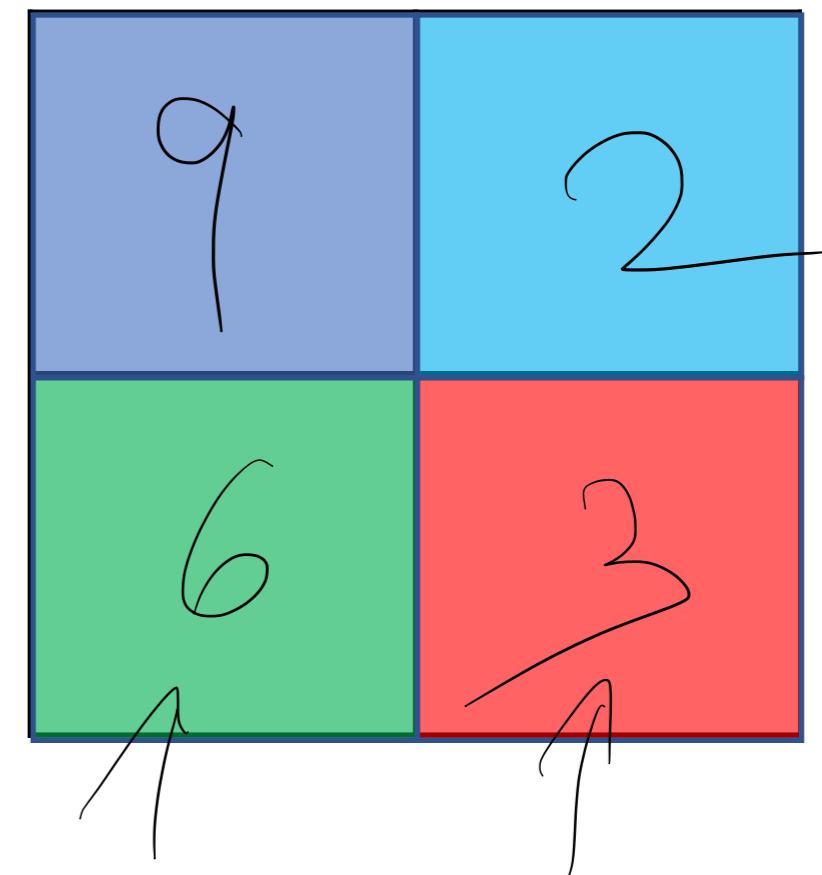
Types of layer in a convolutional network

- Convolution
- Pooling
- Fully connected

Max pooling

$f = 2 \times 2$
 $S = 2$

1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2



(d-f+1)

CNN

