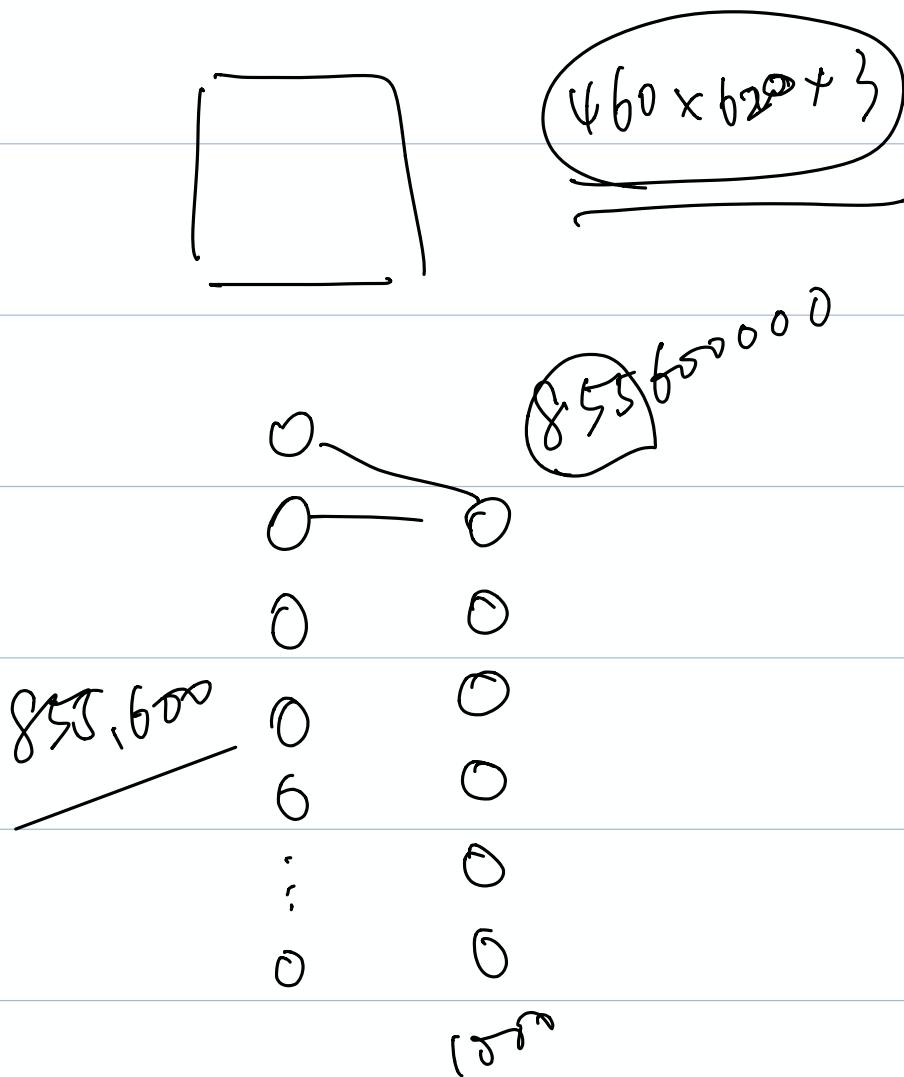


# why use convolutions?

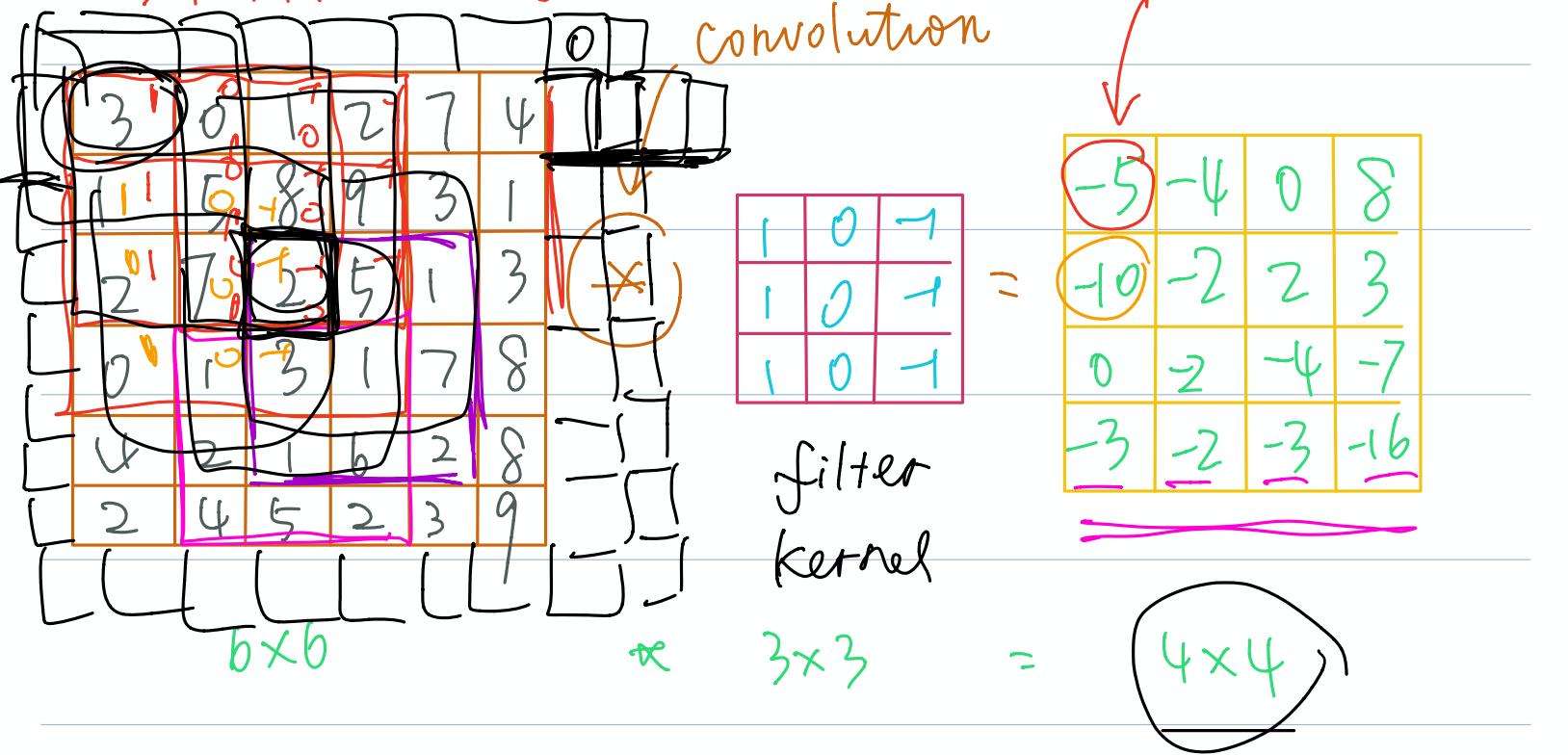
$64 \times 64 \times 3$  image      } what will happen if  
 $1000 \times 1000 \times 3$  image      } we flatten them and  
    only use fully connected  
    layers?



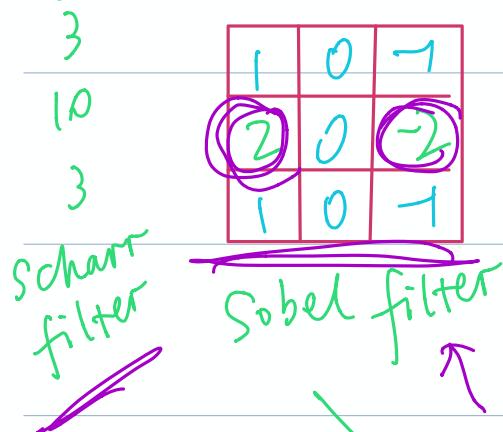
# Filter and kernel

## vertical edge detection

$$3 \times 1 + 1 \times 1 + 2 \times 1 + 0 \times 0 + 5 \times 0 + 7 \times 0 + (x-1) + 8 \times (-1) + 2 \times (-1) = -5$$



also



also can detect  
vertical edges

$3 \times 3$

$5 \times 5$

$9 \times 9$

$1 \times 1$

$45^\circ$   $30^\circ$

$45^\circ$  and  $30^\circ$

We can learn  
filters that  
better suit  
the task

# Padding

$$(6 \times 6) * (3 \times 3) \rightarrow (4 \times 4)$$

$$\underline{(n \times n)} * \underline{(f \times f)} \rightarrow \underline{(n-f+1 \times n-f+1)}$$

Reason: ① shrinks for large model

② the edge pixels are given less weight

$$p=0$$

$$\boxed{6 \times 6} \xrightarrow{\text{pad}} \boxed{(8 \times 8) * (3 \times 3)} \rightarrow \boxed{6 \times 6}$$

$\underbrace{n+2p \times 2p+n}_{(n+2p-f+1 \times n+2p-f+1)} \quad \underbrace{f \times f}$

valid same .

valid conv: no padding

Same conv: output size = input size.

also because  
there is  
a central  
filter

$$p = \frac{f-1}{2}$$

f is usually odd

Strided Convs:

$$(n \times n) \times (f \times f)$$

padding:  $p$

stride:  $s$

floor

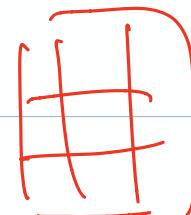
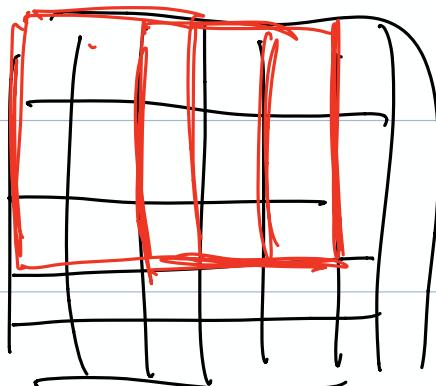
operation

$$\left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor$$

filter need to be

entirely inside

[2.5]

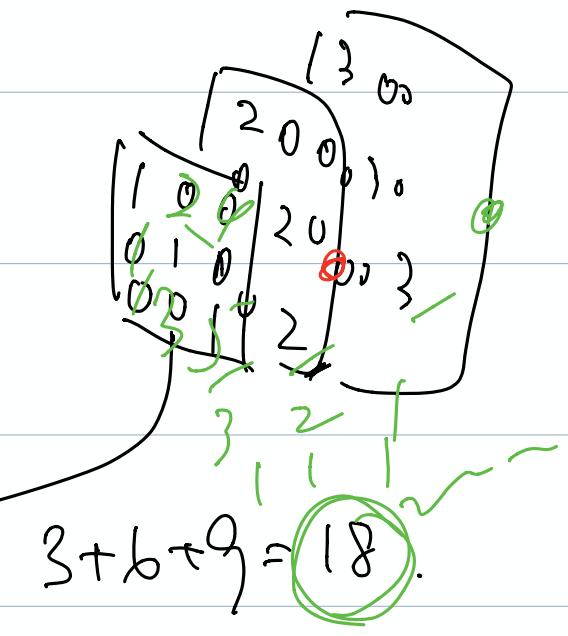
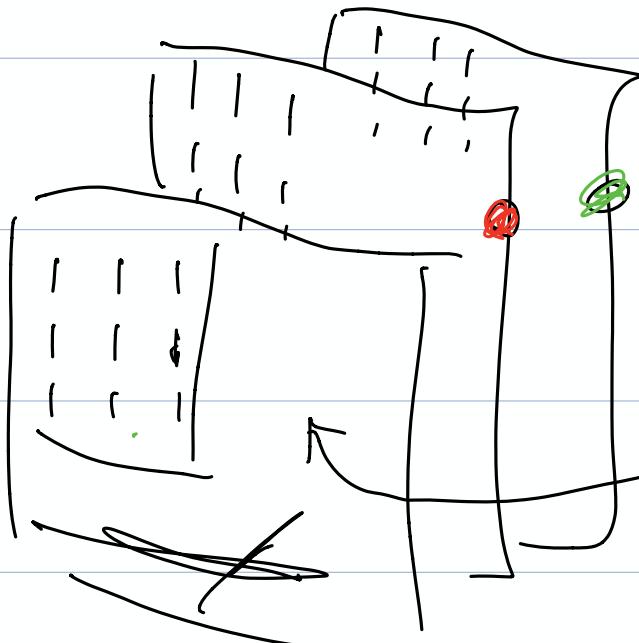


1

2

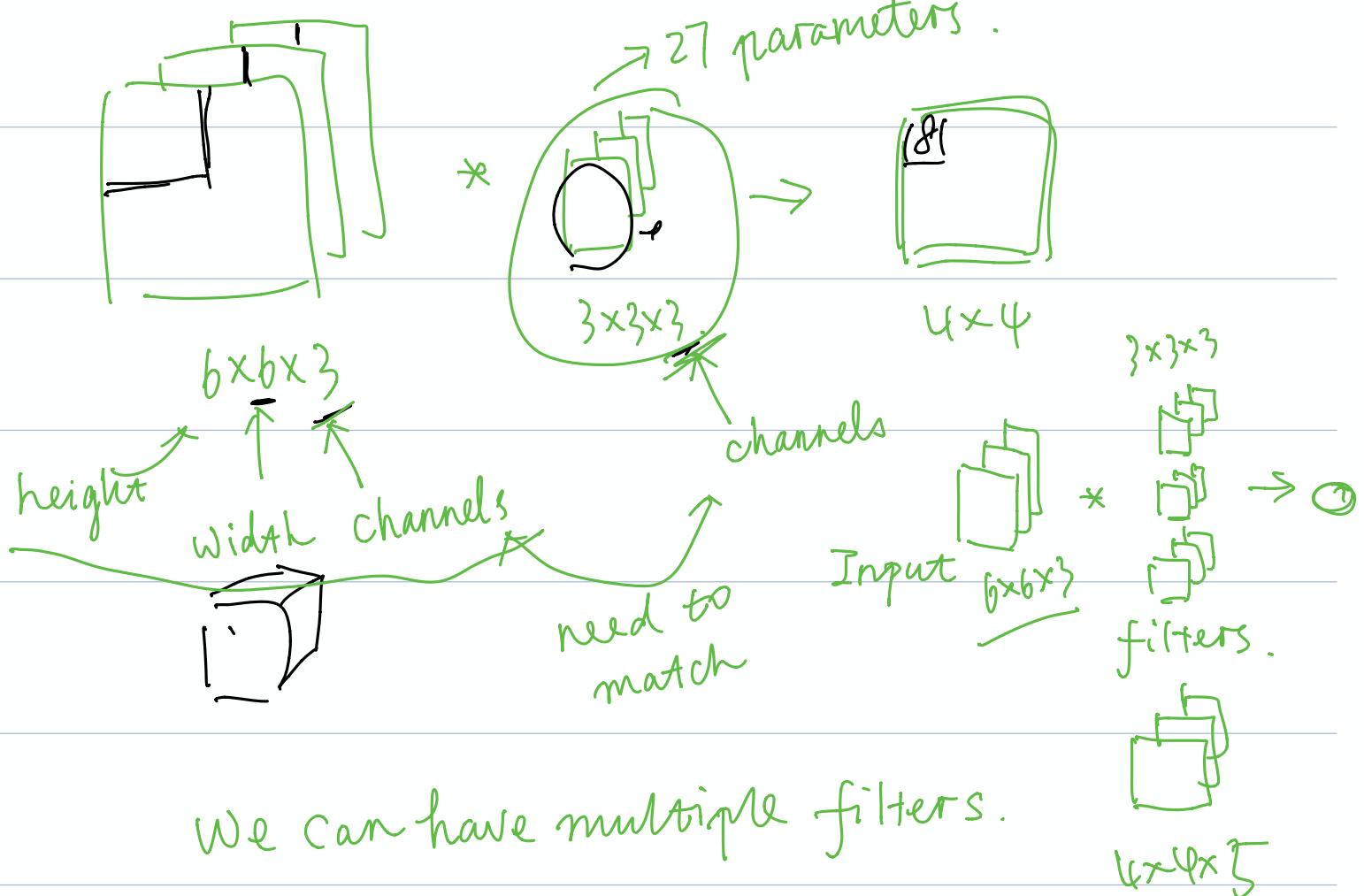
[2.9]  
2

RGB



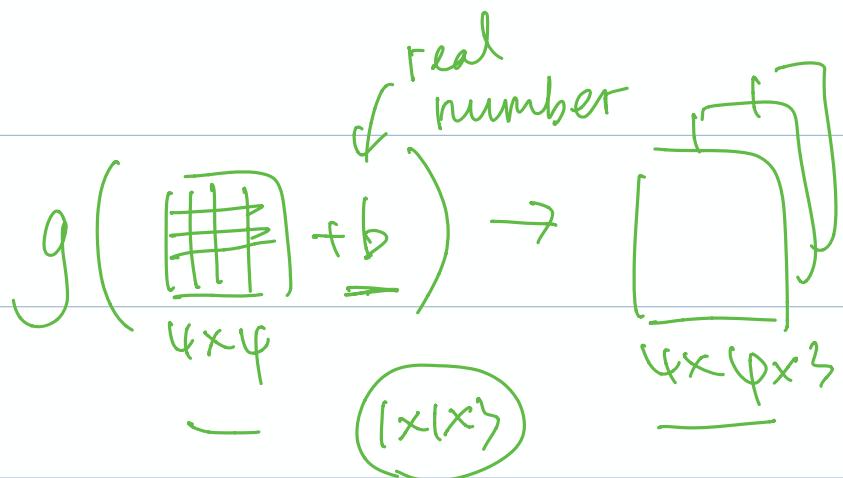
Conv over volumes:

w



and output will be

$$\begin{aligned}
 & \text{Input: } (f \times f \times n_c) \\
 & \text{Number of filters: } n_f \\
 & \text{Output: } (n_f \times f \times f \times n_c) \\
 & \text{Total number of parameters: } \\
 & \quad \text{Input size: } (f \times f \times n_c) \\
 & \quad \text{Number of filters: } n_f \\
 & \quad \text{Number of channels: } n_c
 \end{aligned}$$



$$g(a^{[l_0]} * w^{[l_1]} + b^{[l_1]}) = a^{[l_1]}$$

10 filters of  $3 \times 3 \times 3 \rightarrow 280$  parameters

$f^{[l]}$  = filter size

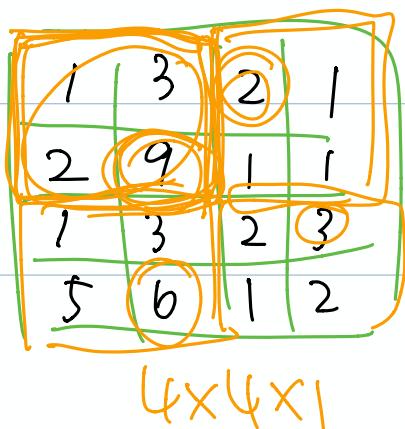
$p^{[l]}$  = padding

$s^{[l]}$  = stride

$n_c^{[l]}$  = # of filters

Pooling layer.

$$\begin{bmatrix} 1 & 3 \\ 2 & 9 \end{bmatrix} \rightarrow \begin{bmatrix} 15 \\ 4 \end{bmatrix}$$

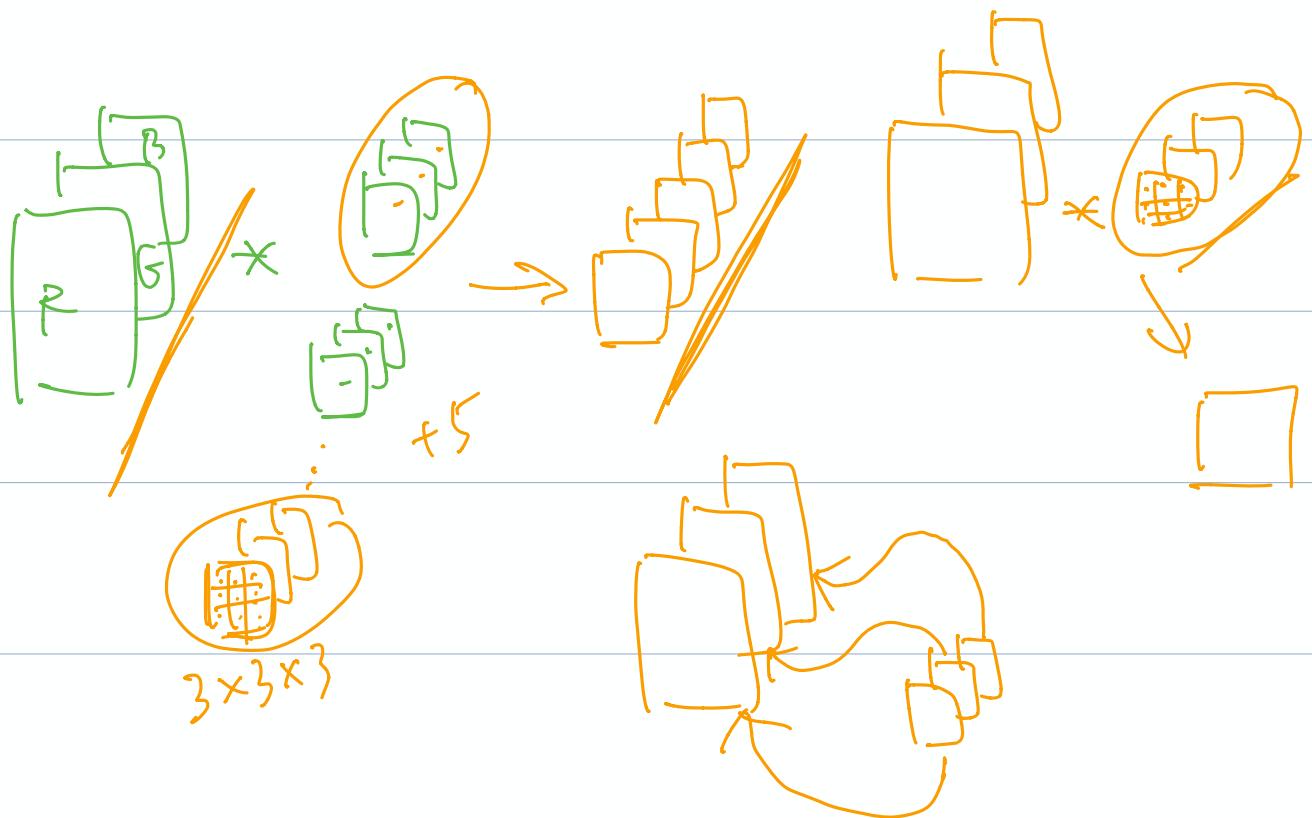


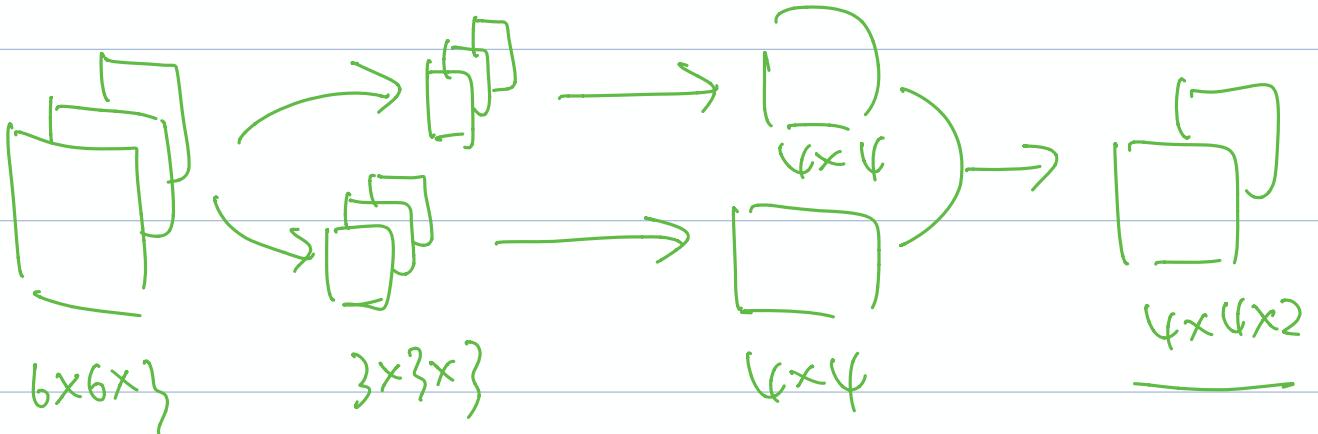
$$\begin{bmatrix} 9 & 2 \\ 6 & 3 \end{bmatrix}$$

$$\begin{matrix} f=2 \\ s=2 \end{matrix} \} \text{ hyperparameters}$$

No parameters  
to learn

Average Pooling





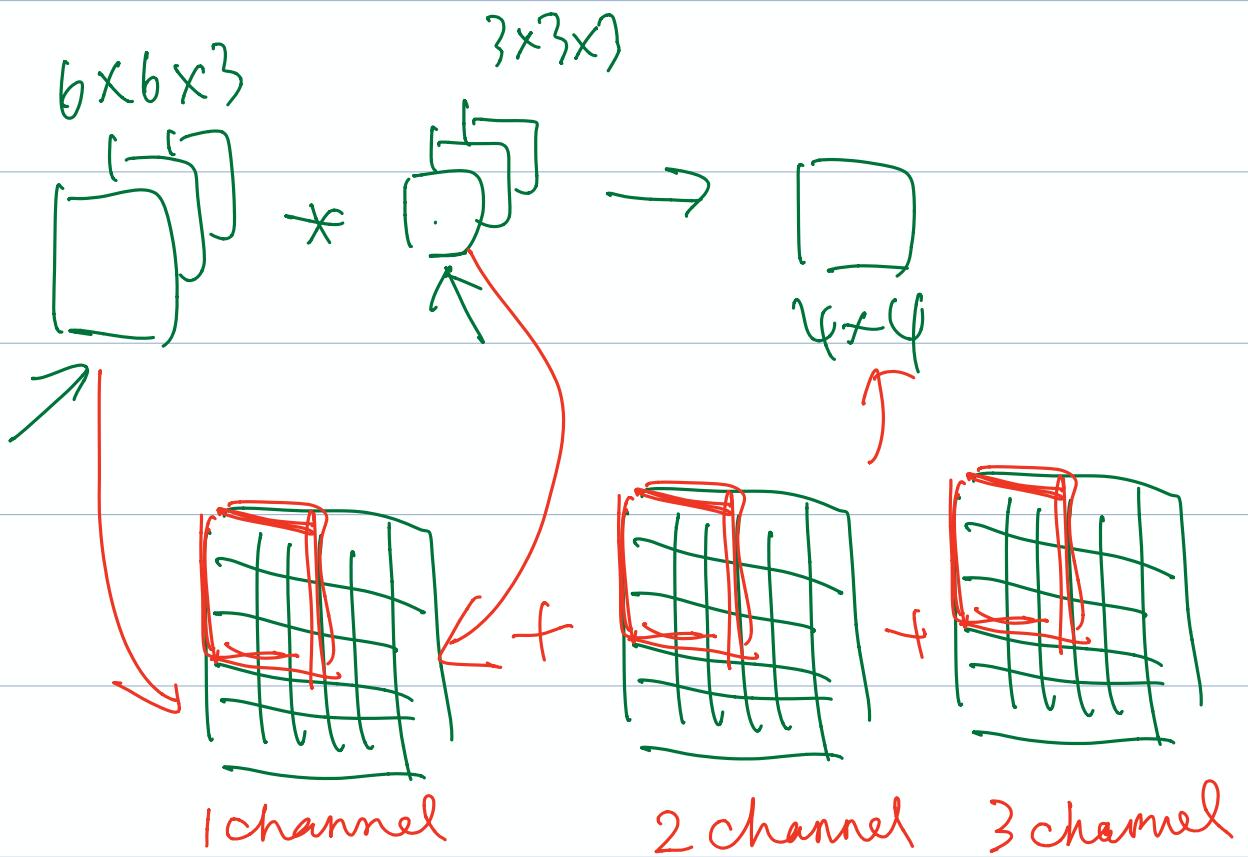
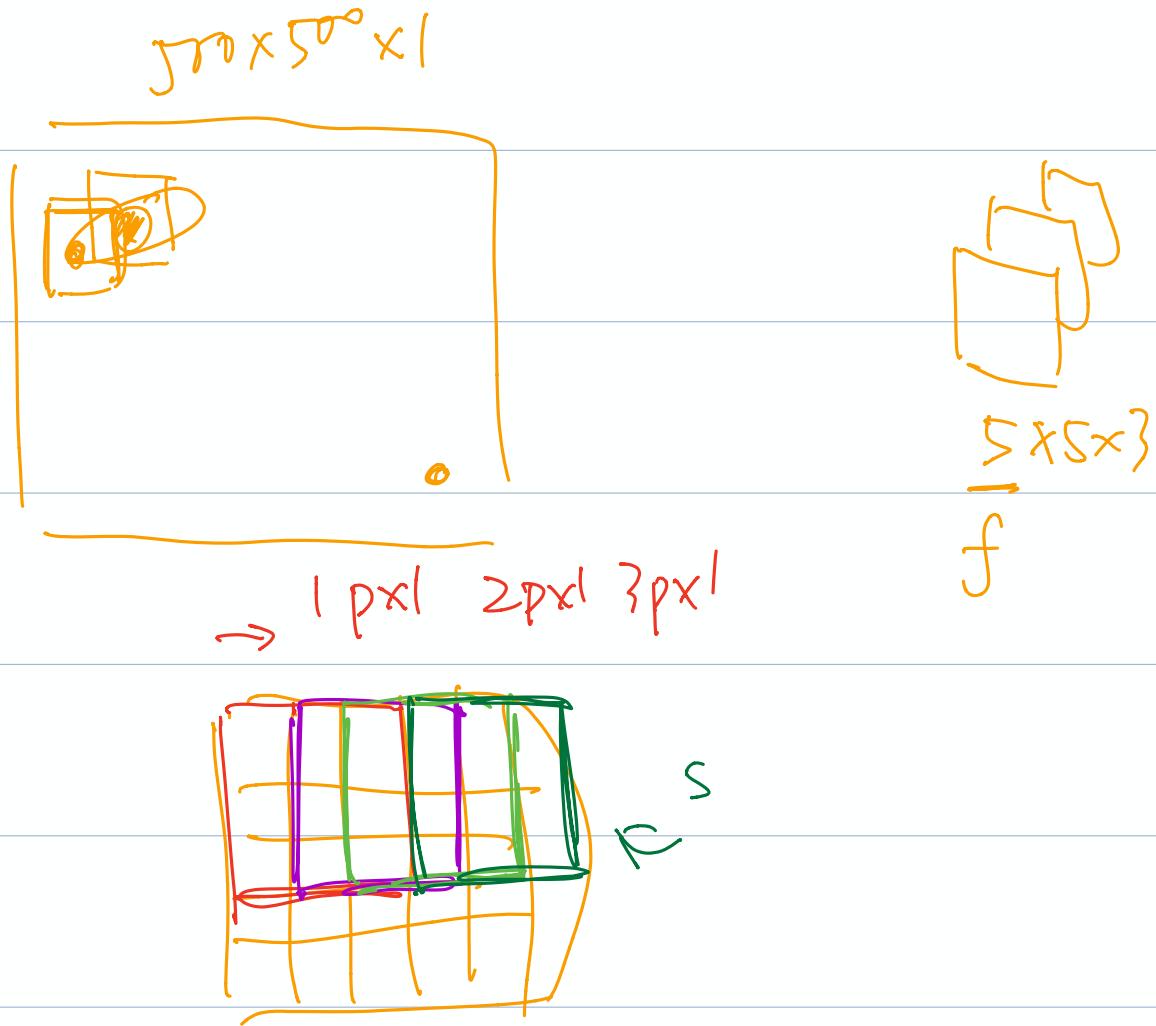
ReLU  $\left( \begin{array}{|c|c|c|} \hline 0 & -1 & 2 \\ \hline 3 & 4 & 5 \\ \hline -6 & 10 & 12 \\ \hline \end{array} \right) + 1 = \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 0 & 1 & 13 \\ \hline 0 & 1 & 13 \\ \hline \end{array}$

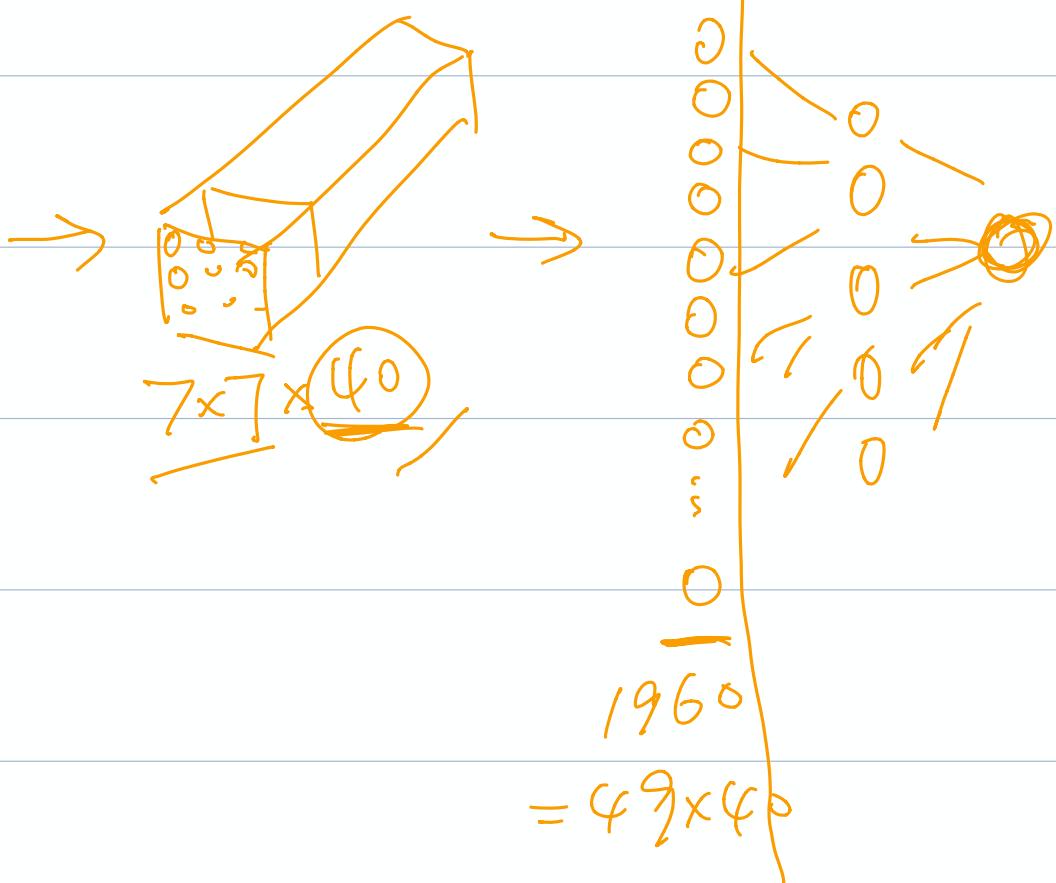
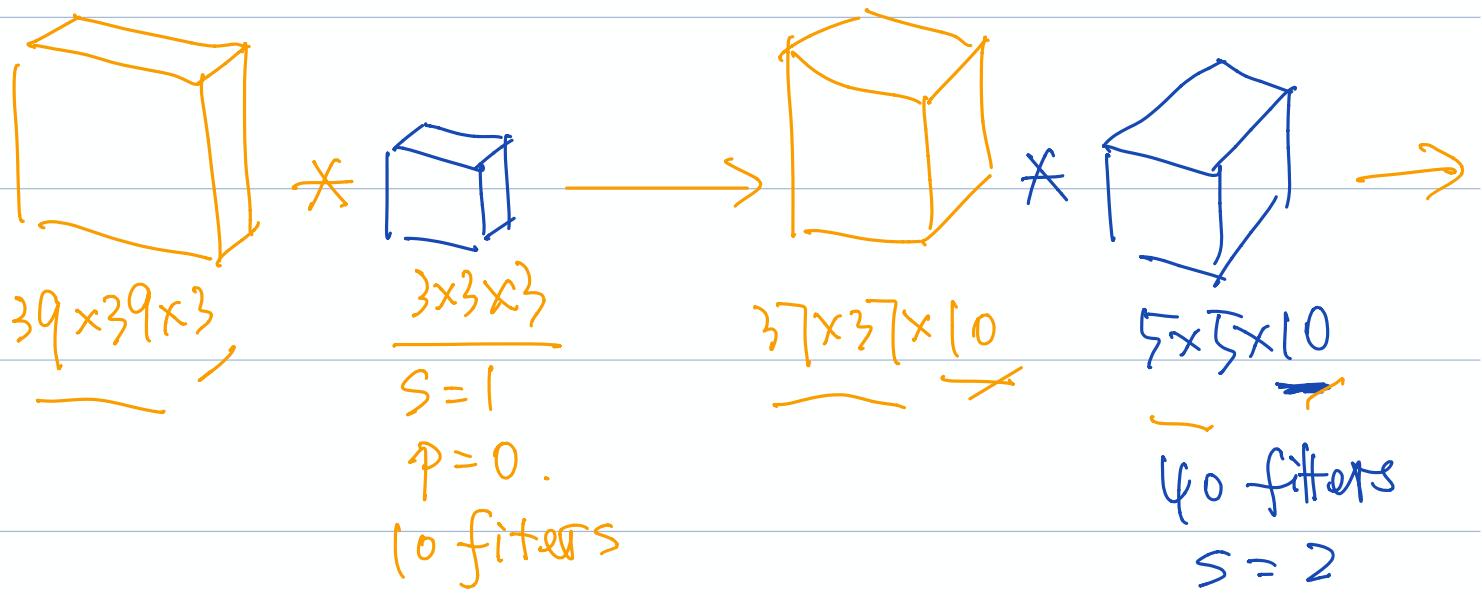
ReLU function definition:

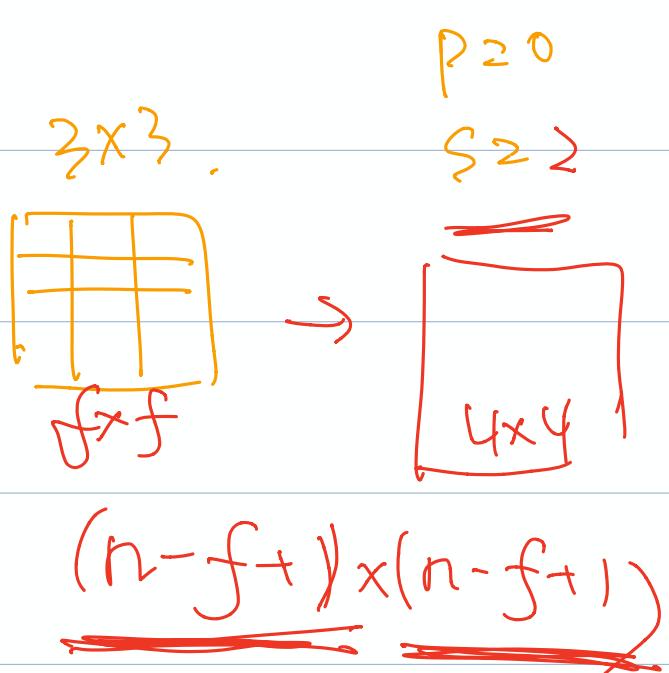
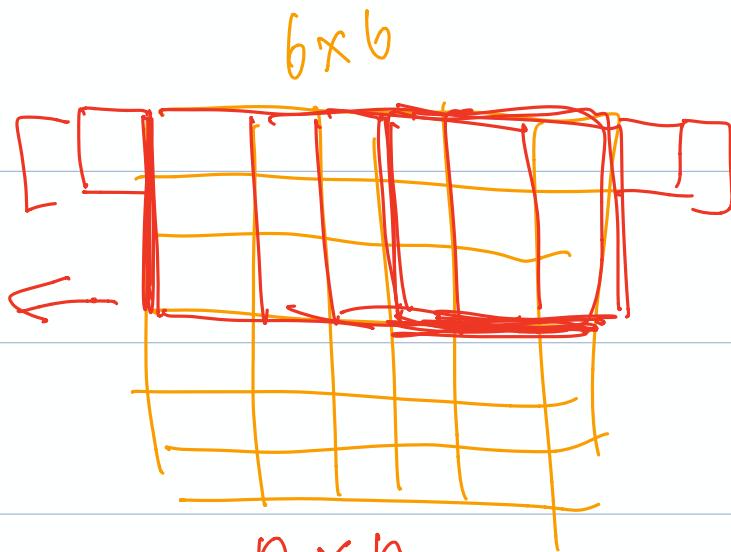
$$\begin{cases} 0 & x \leq 0 \\ x & x > 0 \end{cases}$$

$\frac{27}{11}$   
 10 filters  $3 \times 3 \times 3$

$$27 \times 10 + 10 = 280$$







$$(n+2p)-f+1$$

$$6-3+1=4$$

$$p=1$$

$$8 \times 8$$

$$3 \times 3$$

$$8-3+1=6$$

$$\left\lfloor \frac{8-3}{2}+1 \right\rfloor = 2,5+1$$

$$\left\lfloor \frac{n+2p-f}{s}+1 \right\rfloor$$

