



CS 4375-03

## Lecture 3-C: Convolutional Neural Networks

Feng Chen

University of Texas at Dallas

(Slides adapted from Andrew NG)

# Computer Vision Problems

## Image Classification



64x64

→ Cat? (0/1)

Nowadays, many real-world applications, such as virtual reality glamping and self-driving car, are driving by deep learning. For computer vision, convolutional neural networks (or CNN) is the most successful DL network for these applications.

One example is to build an image classification model to recognize cat in an image.

# Computer Vision Problems

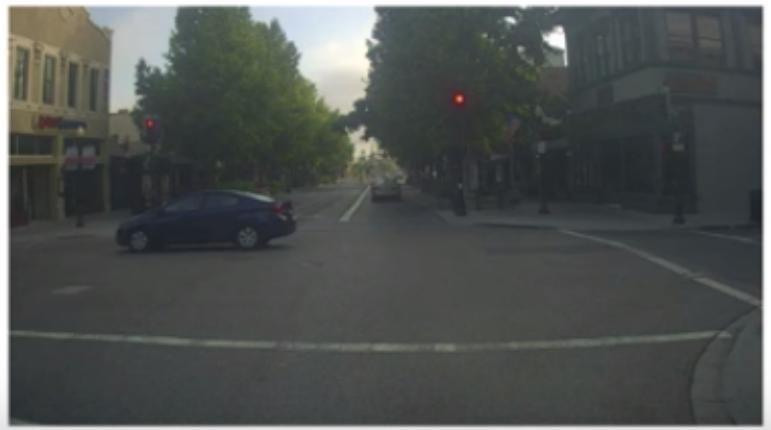
## Image Classification



→ Cat? (0/1)

64x64

## Object detection



Another example is object detection that is critical to build a self-driving car. In this image, you need to detect the location of your car and the locations of other cars in order to avoid them.

# Computer Vision Problems

Image Classification



→ Cat? (0/1)

64x64

Object detection



# Computer Vision Problems

Image Classification



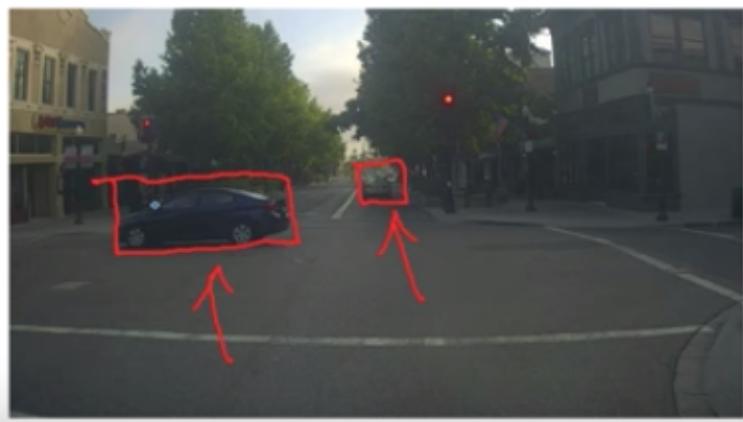
→ Cat? (0/1)

64x64

Neural Style Transfer



Object detection



The third example is to repaint an image in different style. For example, the left figure is a human face and the right figure is a paint by Picasso. We want to repaint the left figure using the style of the right figure.

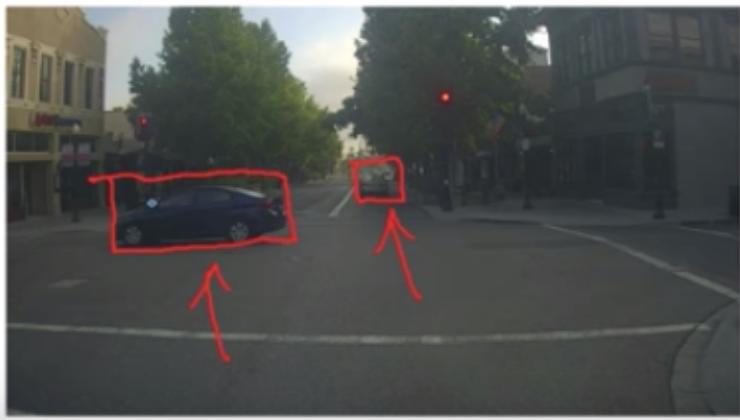
# Computer Vision Problems

Image Classification



→ Cat? (0/1)

Object detection



Neural Style Transfer



One of the challenges in computer vision is that the input can get really big. For example, the top cat image has the dimension 64 by 64 by 3 because for each pixel we have R, G, B, three colors. If we multiply them out, the input feature has the dimension 12288. That is not too bad for this image. However, if you work on a large image, such as the cat image shown on the bottom, it has 1000 pixel by 1000 pixel and just has 1 MB file size. The dimension of this image will be 1000 by 1000 by 3 because we have R, G, B, three channels, that is equal to 3 million pixels in total.

# Deep learning on large images

## Image Classification



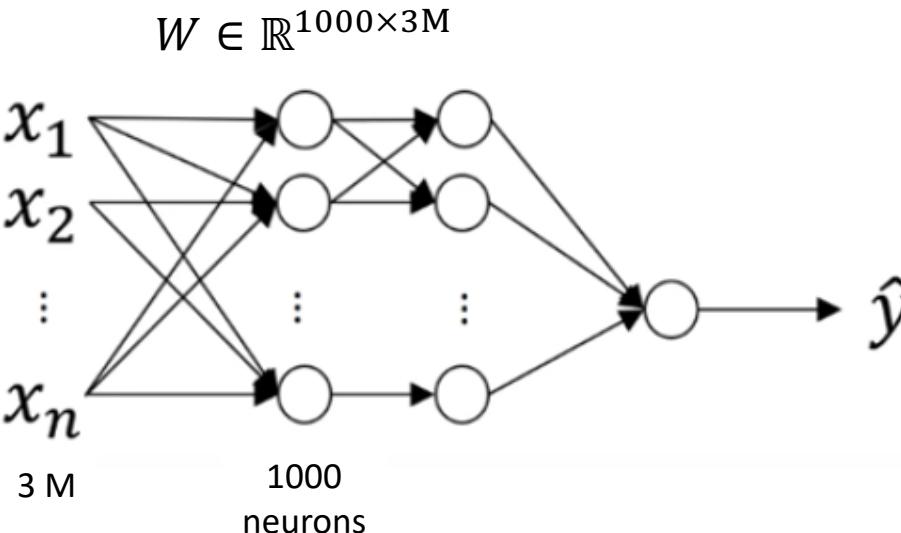
→ Cat? (0/1)

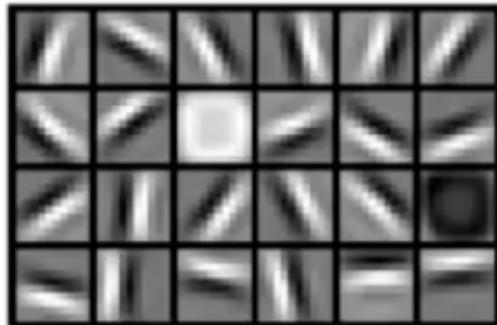
64x64x3 → 12288 variables



1000x1000x3 = 3 M

To consider 3 million features as the input to a fully connected neural network, then the weight matrix for the first latent layer will have the dimension 1000 by 3 million that is equal to 3 billion parameters in total. This is very large. With so many parameters, it will be very difficult to find enough training data to prevent the neural network from overfitting. In addition, both the computational and memory requirement are also infeasible in regular computers.





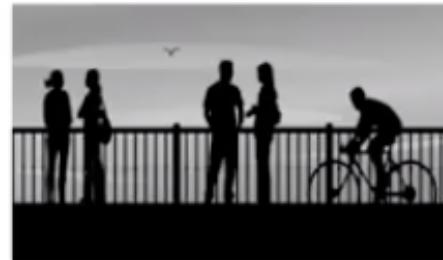
The first convolutional layer  
detects edges



The second convolutional layer detects  
partial objects



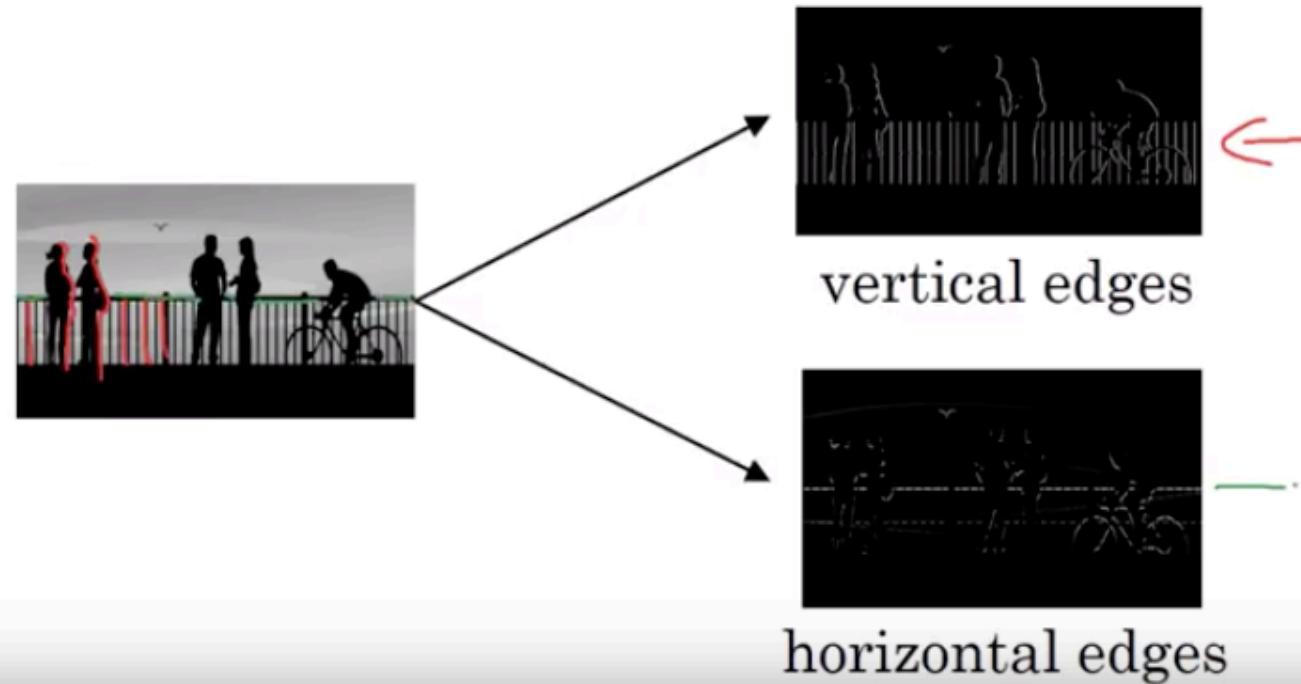
The third convolutional layer may detect  
complete objects like human faces in this case.



vertical edges

For computer vision applications, we often need to use high-resolution images. In order to do that, we need to implement a convolutional operator to extract meaningful features from the image, which is one of the fundamental building blocks of convolutional neural networks. In this page, the example on the top row illustrate convolutions using the example of edge detection. Here, I show three layers of convolutional neural networks for face recognition. The first convolutional layer detects edges, the second convolutional layer detects partial objects. The third convolutional layer may detect ...

For the second example, let us see how we can detect edges in the bottom image. We want computer to figure out what are objects in this image. The first thing you might do is to detect vertical edges shown on the bottom right for example. It has all those vertical lines around the railings and pedestrians.



We might also want to detect horizontal edges. In this image there is a very strong horizontal line around the railing and the bike. This example indicates the usefulness of edge detection. Now the question is: how can we detect vertical and horizontal edges from the input image like this one?

# Vertical edge detection

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	2	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

\*

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

=


Let us look at an example here. We have a six pixel by six pixel gray scale image. We can represent this image as six by six matrix, rather than six by six by three because there are no R, G, and B channels in this image.

Let us first consider the detection of vertical edges in this image. What you can do is to construct a three by three matrix. In terminology of convolutional neural networks, we also call this matrix a filter. I am going to construct the 3 by 3 filter that looks like this. Some papers also call this matrix a kernel instead of a filter. In the slides, I will use the filter terminology.

# Vertical edge detection

In order to calculate the convoluted matrix based on this filter. We take the 3 by 3 filter and place it on top of the three by three region of the original image.

We take the element wise product:  $3 \times 1 + 1 \times 1 + 2 \times 1 + 0 \times 0 + 5 \times 0 + 7 \times 0 + 1 \times -1 + 8 \times -1 + 2 \times -1 = -5$ .

We can add these numbers in any order.

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	2	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

\*

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

=

-5			

# Vertical edge detection

We now shift the block one step to the right and place the filter on top of the region.  
We do the same element-wise product and then do addition. We have...

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	2	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

\*

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

=

-5	<del>-4</del>		

# Vertical edge detection

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	2	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

\*

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

=

-5	-4	0	

# Vertical edge detection

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	2	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

\*

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

=

-5	-4	0	8

# Vertical edge detection

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	2	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

\*

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

=

-5	-4	0	8
<i>-10</i>			

# Vertical edge detection

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	2	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

\*

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

=

-5	-4	0	8
-10	<b>-2</b>		

# Vertical edge detection

After we calculate all the convoluted values, what we obtain is an edge detector. You can easily see why it is called edge detector on the next page.

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	2	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

\*

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

=

-5	-4	0	8
-10	-2	2	3
0	-2	-4	-7
-3	-2	-3	<b>-16</b>

# Vertical 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

=


Suppose the input image is shown on the left.

# Vertical 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



We can visualize both the input image and the filter on the bottom. -1 refers to black, 0 refers to grey, a positive value refers to white. We can clearly see there is an edge in the middle line of the input image.

$$\begin{matrix} & * & \begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix} & = & \begin{matrix} & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \end{matrix} \end{matrix}$$



# Vertical 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

This page shows the output matrix after we conduction the convolution based on this filter.

\*

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



# Vertical edge detection

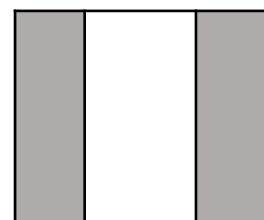
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0

\*

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



The bottom shows the visualization of the output matrix. The white region in the middle has clearly highlighted the middle edge region in the input image.

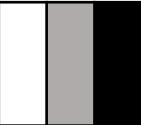
# Vertical edge detection

In this next page, I will use a different input image that also has a edge in the middle.

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

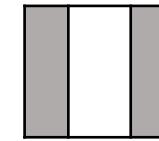


\*

$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array}$$


=

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



# Vertical edge detection

We want to check if the filter can still detect the edge region in the middle.

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

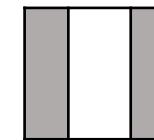


\*

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



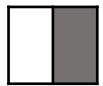
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10

\*

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

# Vertical 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



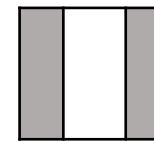
\*

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

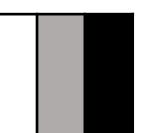


0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10



\*

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



# Vertical edge detection

This page shows the output matrix of of this new input image.

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



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	10	10	10



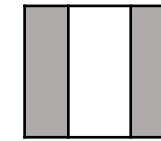
\*

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



\*

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

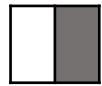


=

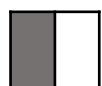
0	-30	-30	0
0	-30	-30	0
0	-30	-30	0
0	-30	-30	0

# Vertical 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



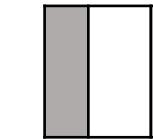
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10



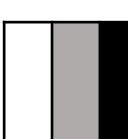
\*

$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array}$$
  

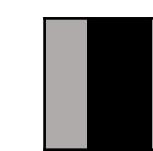

=

$$\begin{array}{|c|c|c|c|} \hline 0 & 30 & 30 & 0 \\ \hline 0 & 30 & 30 & 0 \\ \hline 0 & 30 & 30 & 0 \\ \hline 0 & 30 & 30 & 0 \\ \hline \end{array}$$
  


\*

$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array}$$
  


=

$$\begin{array}{|c|c|c|c|} \hline 0 & -30 & -30 & 0 \\ \hline 0 & -30 & -30 & 0 \\ \hline 0 & -30 & -30 & 0 \\ \hline 0 & -30 & -30 & 0 \\ \hline \end{array}$$
  


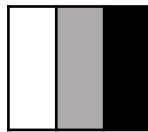
The visualization of the output image clearly indicates that it can detect the edge region in the middle.

To summarize, a filter with edge detection is able to highlight the edge region with a different color.

# Vertical and horizontal edge detection

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

Vertical



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

Horizontal



We have demonstrated that the left filter is able to detect vertical edges. In the following slides, I will show that the right filter is able to detect horizontal edges.

# horizontal edge detection

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



\*

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

Horizontal



\*

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

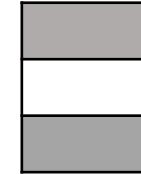
Horizontal



In this page, we have two different input images that both have horizontal edge in the middle. We apply the same filter for horizontal edge detection. The output matrices demonstrate that this filter is able to detection the horizontal edge regions in both input images as well.

=

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



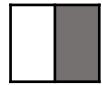
=

0	0	0	0
-30	-30	-30	-30
-30	-30	-30	-30
0	0	0	0

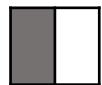


# horizontal edge detection

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



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	10	10	10



\*

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

Horizontal



=

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



\*

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

Horizontal



=

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



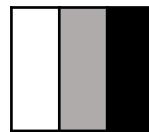
If we apply the filter to input images that do not have horizontal edges, but only vertical edges, the output matrix will be a zero matrix. That indicates that there are no edge regions in the input images.

# Vertical and horizontal edge detection

In this page, the input image has both horizontal and vertical edge regions. If we apply a filter for horizontal edge detection, we can then detect only the horizontal edge region in the input page as shown in the next page.

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

Vertical



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	0	10	10

\*

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



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

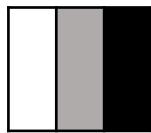
Horizontal



# Vertical and horizontal 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	0	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

# Learning to detect edges

*put a little bit more weight on the central pixel*

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

1	0	-1
2	0	-2
1	0	-1


*Sobel filter*

In the previous slides, I have demonstrated the filter shown on the left for vertical edge detection. We can obtain a different filter by putting a little more weight on the center pixels. This filter is called sobel filter that can better detect edges that are not directly vertical as shown on the example on the bottom.

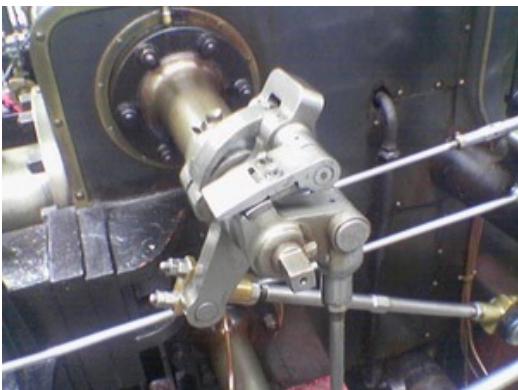
# Learning to detect edges

put a little bit more weight on the central pixel

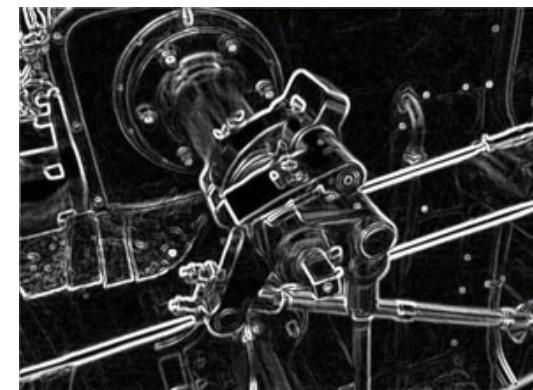
1	0	-1
1	0	-1
1	0	-1

1	0	-1
2	0	-2
1	0	-1


Sobel filter



\* Sobel filter =



# Learning to detect edges

put a little bit more weight on the central pixel

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

1	0	-1
2	0	-2
1	0	-1

Sobel filter

3	0	-3
10	0	-10
3	0	-3

Scharr filter

instead of only 1, 0, -1  
we can use other excessive  
numbers as well.

Instead of using only 1, 0 and -1, we can use other excessive numbers as well, such as 3, 10, 3, and then -3, -10, and -3. We obtain a well-known filter called Scharr filter that has swiftly different property.

# Learning to detect edges

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

1	0	-1
2	0	-2
1	0	-1

3	0	-3
10	0	-10
3	0	-3

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	2	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

( 6 × 6 )

With the rise of deep learning, one of things we learned is that, when you really want to detect edges in some complicated images maybe you can just learn the filter and treat the 9 numbers of this filter matrix as parameters, so that we can take the six by six image and convolute it with your 3 by 3 filter that gives you a good edge detector. We can then learn one of the filters that we have introduced, or learn even better ones directly from the data.

\*

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

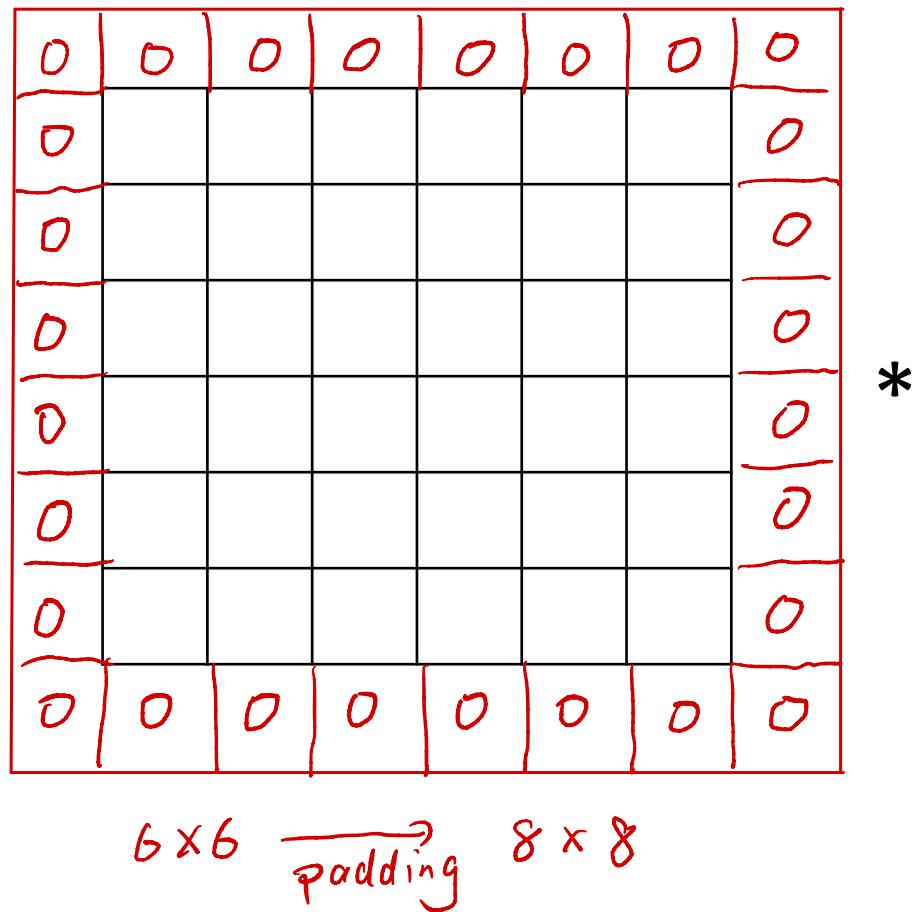
$$(3 \times 3) \\ w \in \mathbb{R}^{3 \times 3}$$

=


$$(4 \times 4)$$

Not matter what is the filter matrix, the output size will shrink the size by 2: from 6 to 4. We need an operation called padding that can keep the size the same.

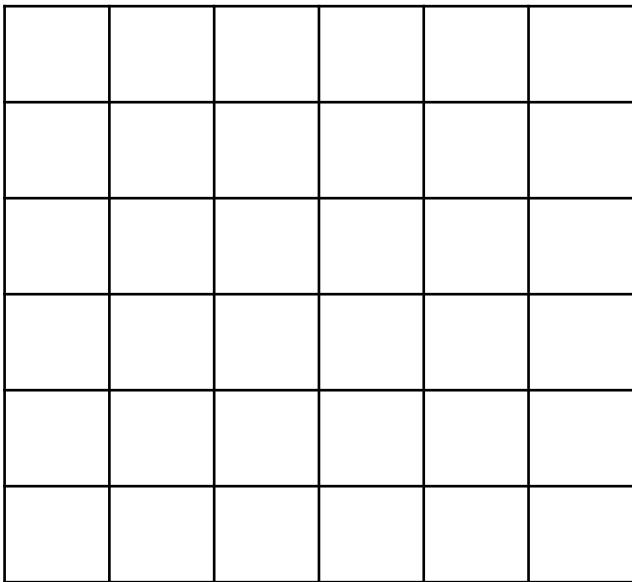
# Padding



You may have noticed that the input image will shrink after one convolution. For example, given an 6 by 6 image as input, the output image after one convolution will be 4 by 4. The image will get really small after several steps of convolutions. We want to avoid that. That means we do not want image to shrink every time we do convolution. To address this we can use padding. The idea is to pad the image with an additional border of zero pixels around the edges.

In this example, after padding the zero pixels, the size of the output image is now 6 by 6, the same size of the input image.

# Padding

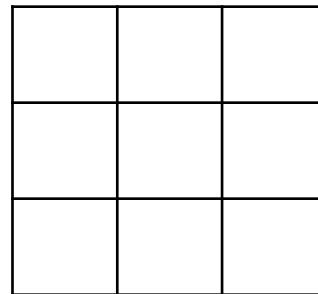


$6 \times 6$

$(n \times n)$

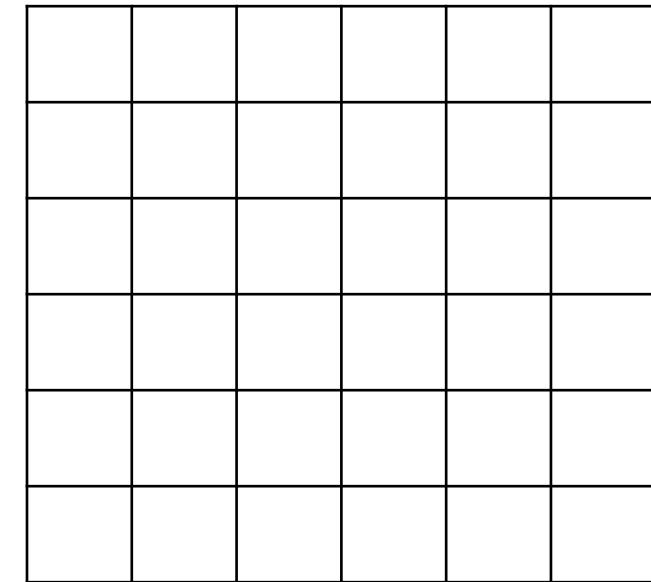
padding ( $p$ ) = 1

\*



$3 \times 3$   
 $(f \times f)$

=



$6 \times 6$

$\underbrace{(n+2p-f+1)}_{11} \times \underbrace{(n+2p-f+1)}_{11}$

$\underbrace{6+2-1-3+1}_{6}$

Suppose the input image is  $n$  by  $n$  and the filter is  $f$  by  $f$ , there is a simple formula to calculate the output image size:  $n + 2p - f + 1$ . Here  $p$  refers to the number of rows of zero padding. In this example,  $p$  is equal to 1. In this example, as  $n = 6$  and  $f = 3$ , using the formula we can calculate that the output image size  $6 + 2 - 3 + 1 = 6$ .

# Valid and Same convolutions

“convolution without padding” :

“convolution with padding ( $p$ )” :

“Same” : Pad so that output size is the same  
as the input size

In this previous page, given an input image of size 6 by 6 and a filter of size , by padding one row of zero pixels, the output image has the same input size. Now, given input image of any size, we need to a general rule to pad the input image such that the output size is the same as the input size.

# Valid and Same convolutions

Input Image:  $n \times n$ , Filter:  $f \times f$

“convolution without padding”:  $(n - f + 1) \times (n - f + 1)$

“convolution with padding ( $p$ )”:  $(n + 2 \cdot p - f + 1) \times (n + 2 \cdot p - f + 1)$

“Same”:

The question is to decide the value of  $f$ , such that  $n + 2p - f + 1$  is equal to  $n$ .

By solving this simple equation, we can find that the pad size  $p$  is only decided by the filter size  $f$ . If  $f$  is an odd number, than the pad size  $p$  should be equal to  $\frac{f-1}{2}$ . If  $f$  is an even number, than it is not possible to ensure the same input and output size. In this case, we take the floor of  $\frac{f-1}{2}$ , ensuring that the output size is no greater than the input size.

Pad so that output size is the same  
as the input size

$$\left\lfloor \frac{f-1}{2} \right\rfloor$$

decide the value of  $f$  such that

$$n + 2p - f + 1 = n$$

$$\Rightarrow 2p - f + 1 = 0$$

$$\Rightarrow p = \left\lfloor \frac{f-1}{2} \right\rfloor$$

$f$  is often an odd number. 37

# Strided convolution

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 2 & 3 & 7 & 4 & 6 & 2 & 9 \\ \hline 6 & 6 & 9 & 8 & 7 & 4 & 3 \\ \hline 3 & 4 & 8 & 3 & 8 & 9 & 7 \\ \hline 7 & 8 & 3 & 6 & 6 & 3 & 4 \\ \hline 4 & 2 & 1 & 8 & 3 & 4 & 6 \\ \hline 3 & 2 & 4 & 1 & 9 & 8 & 3 \\ \hline 0 & 1 & 3 & 9 & 2 & 1 & 4 \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline 3 & 4 & 4 \\ \hline 1 & 0 & 2 \\ \hline -1 & 0 & 3 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline \quad & \quad & \quad \\ \hline \quad & \quad & \quad \\ \hline \quad & \quad & \quad \\ \hline \end{array}$$

Strided convolution is another piece of the basic building block of convolutions as used in convolutional neural networks. Let me show you an example. Let's say you want to convolute this seven by seven image with a 3 by 3 filter shown in the middle. Instead of doing the usual way, we are going to do it with a stride of two. That means you take the element wise production as usual in the upper left. Then instead of stepping the shaded box over by one step, we are going to jump it over by two steps.

# Strided convolution

2	3	7	4	6	2	9
6	6	9	8	7	4	3
3	4	8	3	8	9	7
7	8	3	6	6	3	4
4	2	1	8	3	4	6
3	2	4	1	9	8	3
0	1	3	9	2	1	4

\*

3	4	4
1	0	2
-1	0	3

=


# Strided convolution

2	3	7	4	6	2	9
6	6	9	8	7	4	3
3	4	8	3	8	9	7
7	8	3	6	6	3	4
4	2	1	8	3	4	6
3	2	4	1	9	8	3
0	1	3	9	2	1	4

\*

3	4	4
1	0	2
-1	0	3

=

91		

# Strided convolution

2	3	7	4	6	2	9
6	6	9	8	7	4	3
3	4	8	3	8	9	7
7	8	3	6	6	3	4
4	2	1	8	3	4	6
3	2	4	1	9	8	3
0	1	3	9	2	1	4

\*

3	4	4
1	0	2
-1	0	3

=

91	100	

# Strided convolution

2	3	7	4	6	2	9
6	6	9	8	7	4	3
3	4	8	3	8	9	7
7	8	3	6	6	3	4
4	2	1	8	3	4	6
3	2	4	1	9	8	3
0	1	3	9	2	1	4

\*

3	4	4
1	0	2
-1	0	3

=

91	100	83

# Strided convolution

2	3	7	4	6	2	9
6	6	9	8	7	4	3
3	4	8	3	8	9	7
7	8	3	6	6	3	4
4	2	1	8	3	4	6
3	2	4	1	9	8	3
0	1	3	9	2	1	4

\*

3	4	4
1	0	2
-1	0	3

=

91	100	83
69		

# Strided convolution

2	3	7	4	6	2	9
6	6	9	8	7	4	3
3	4	8	3	8	9	7
7	8	3	6	6	3	4
4	2	1	8	3	4	6
3	2	4	1	9	8	3
0	1	3	9	2	1	4

\*

3	4	4
1	0	2
-1	0	3

=

91	100	83
69	91	

# Strided convolution

2	3	7	4	6	2	9
6	6	9	8	7	4	3
3	4	8	3	8	9	7
7	8	3	6	6	3	4
4	2	1	8	3	4	6
3	2	4	1	9	8	3
0	1	3	9	2	1	4

\*

3	4	4
1	0	2
-1	0	3

=

91	100	83
69	91	127

# Strided convolution

2	3	7	4	6	2	9
6	6	9	8	7	4	3
3	4	8	3	8	9	7
7	8	3	6	6	3	4
4	2	1	8	3	4	6
3	2	4	1	9	8	3
0	1	3	9	2	1	4

\*

3	4	4
1	0	2
-1	0	3

=

91	100	83
69	91	127
44		

# Strided convolution

2	3	7	4	6	2	9
6	6	9	8	7	4	3
3	4	8	3	8	9	7
7	8	3	6	6	3	4
4	2	1	8	3	4	6
3	2	4	1	9	8	3
0	1	3	9	2	1	4

\*

3	4	4
1	0	2
-1	0	3

=

91	100	83
69	91	127
44	72	

# Strided convolution ( $s = 2$ )

2	3	7	4	6	2	9
6	6	9	8	7	4	3
3	4	8	3	8	9	7
7	8	3	6	6	3	4
4	2	1	8	3	4	6
3	2	4	1	9	8	3
0	1	3	9	2	1	4

$7 \times 7$

$n \times n$

$7 \times 7$

\*

3	4	4
1	0	2
-1	0	3

$3 \times 3$

=

91	100	83
69	91	127
44	72	74

$3 \times 3$

\*

$f \times f$

=

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

\*

$3 \times 3$

=

$$\underbrace{\left\lfloor \frac{7+0-3}{2} + 1 \right\rfloor}_{\frac{11}{3}} \times \underbrace{\left\lfloor \frac{7+0-3}{2} + 1 \right\rfloor}_{\frac{11}{3}}$$

$\times$

$_{48}$

# Summary of notation: padding only

$n \times n$  image

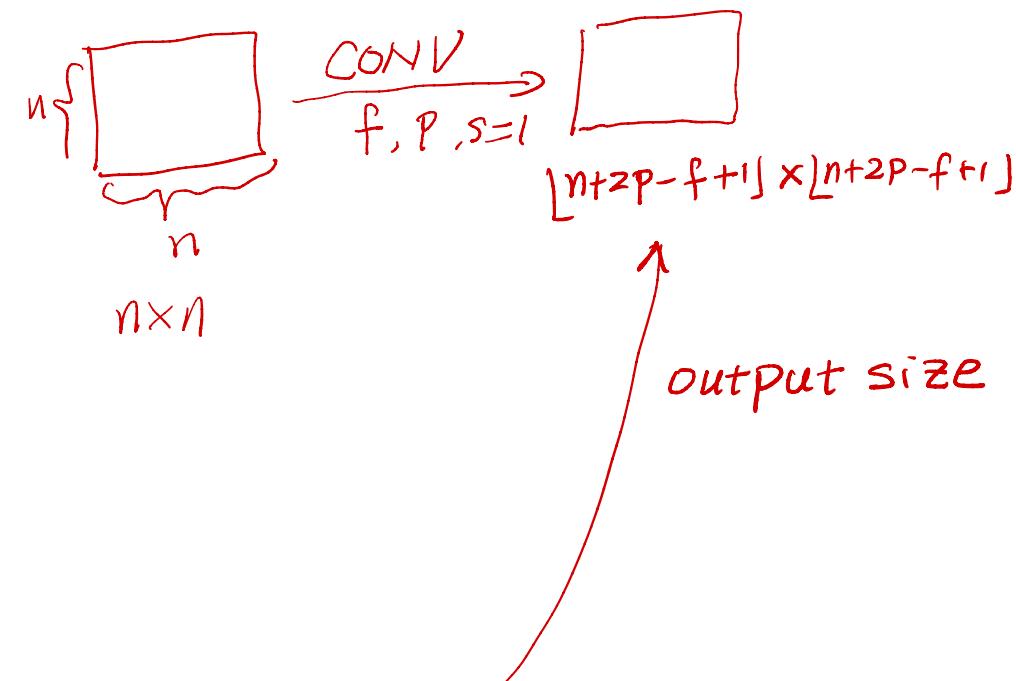
padding  $p$

output size

$f \times f$  filter

stride 1 (default)

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



Now I am ready to provide a summary of notations for padding only. Given a  $n$  by  $n$  input image and a  $f$  by  $f$  filter, suppose the padding size is  $p$  and the stride size is 1 (the default), then the output size is shown on the bottom.

# Summary of notation: stride only

$n \times n$  image

$f \times f$  filter

padding 0

stride  $s$

output size

$$\left\lfloor \frac{n + 2 \cdot 0 - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n + 2 \cdot 0 - f}{s} + 1 \right\rfloor$$

If padding  $p = 0$  and the stride is  $s$ , then the output size is shown on the new formula on the bottom.

# Summary of notation: padding + stride

$n \times n$  image

$f \times f$  filter

padding  $p$

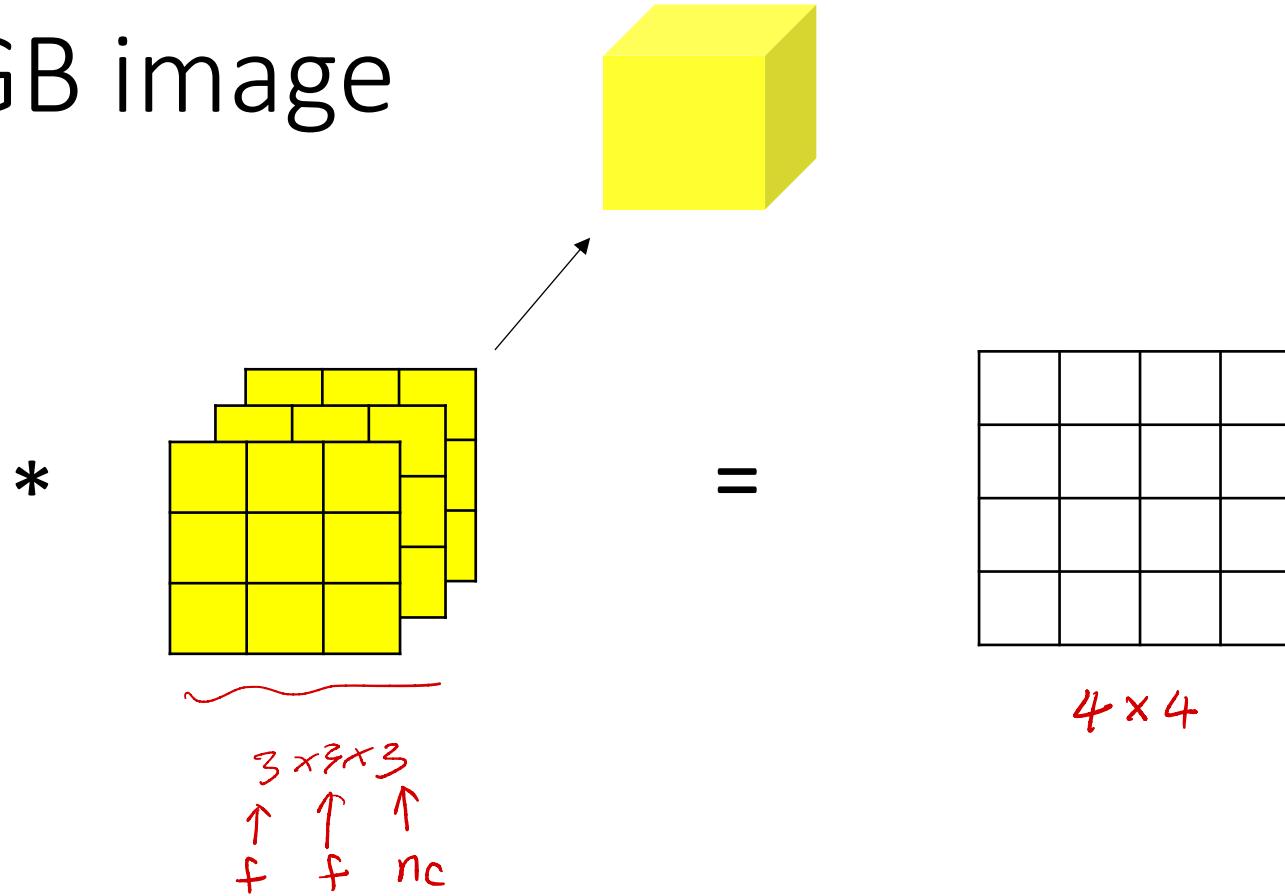
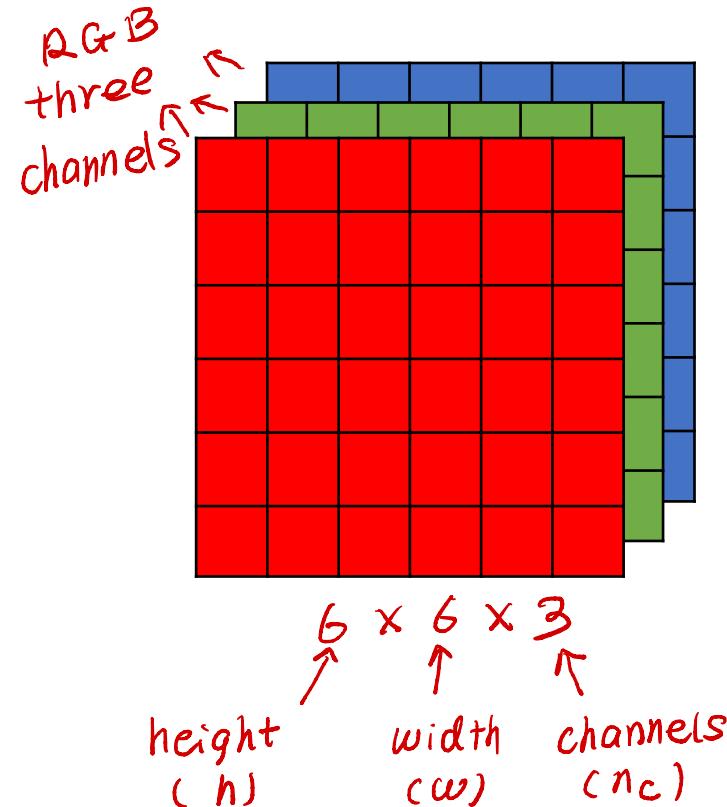
stride  $s$

output size

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

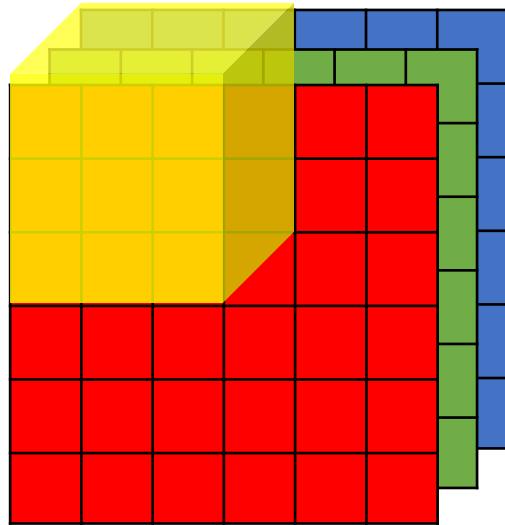
Considering both padding  $p$  and stride  $s$ , then we have the formula to calculate the outsize as:  
the food of  $(n+2p-f)/s + 1$ .

# Convolution on RGB image

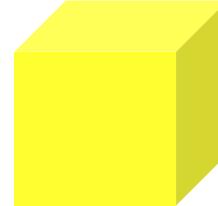
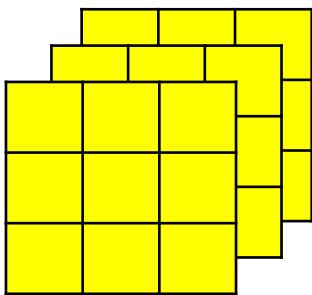


You have seen how we do convolution in 2D images. Here, I will show how to do convolution on 3D volumes or RGB image.

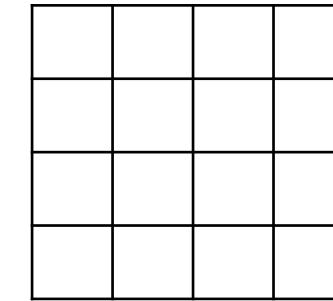
# Convolution on RGB image



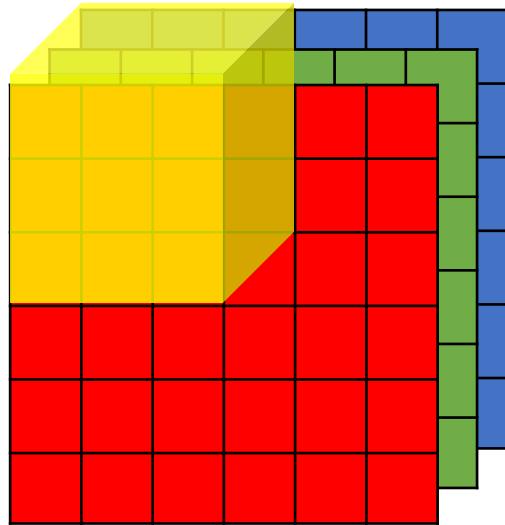
\*



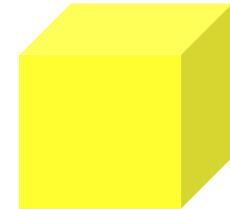
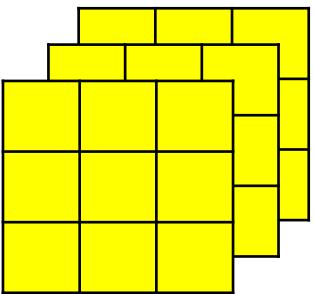
=



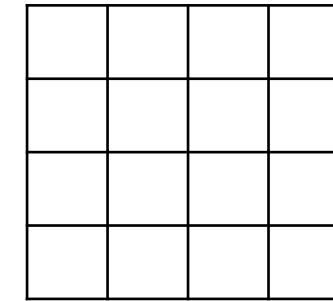
# Convolution on RGB image



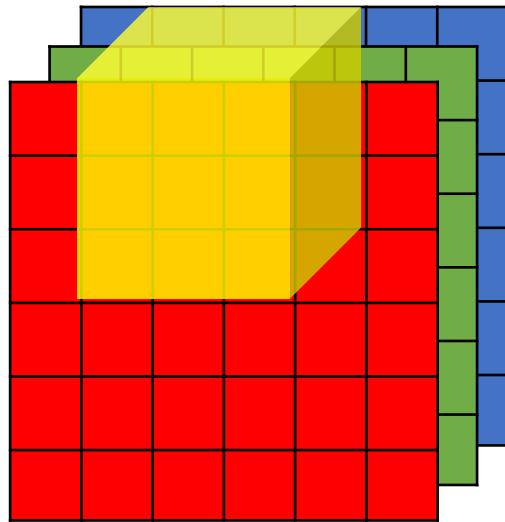
\*



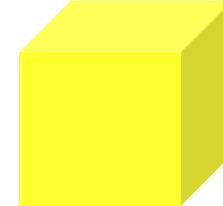
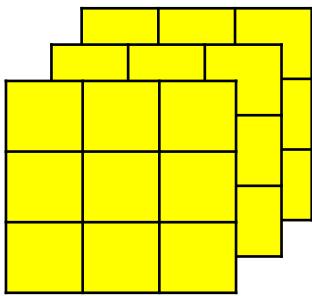
=



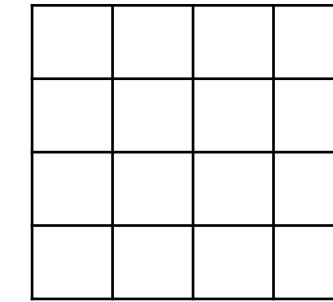
# Convolution on RGB image



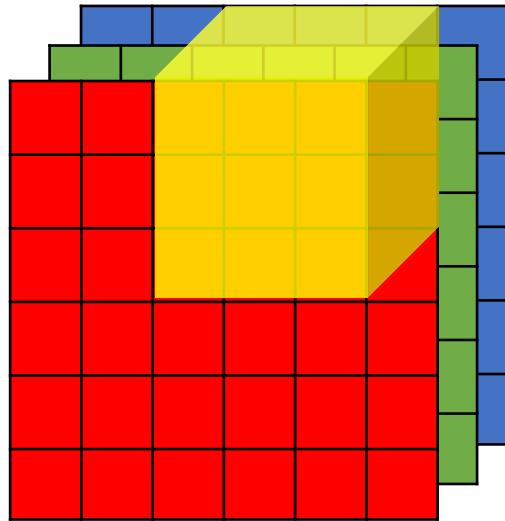
\*



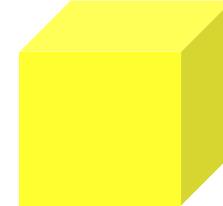
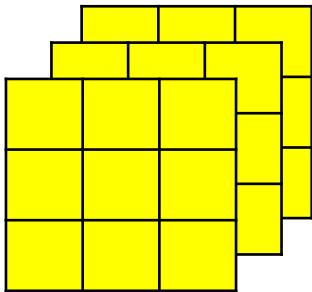
=



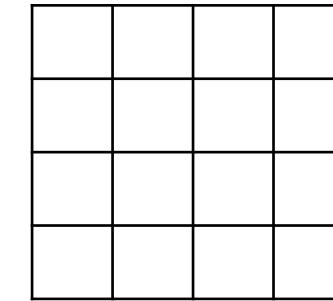
# Convolution on RGB image



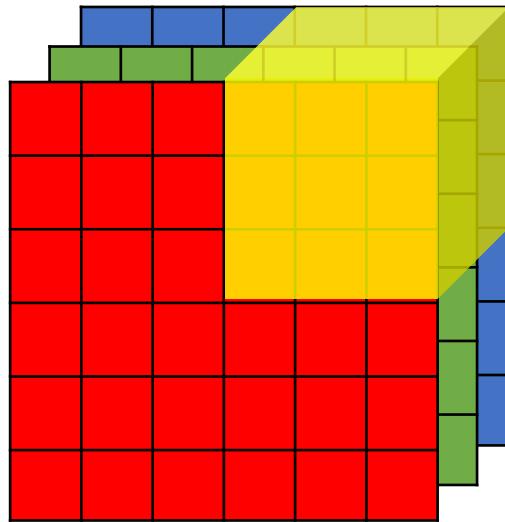
\*



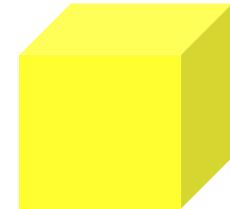
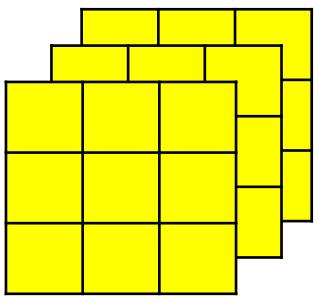
=



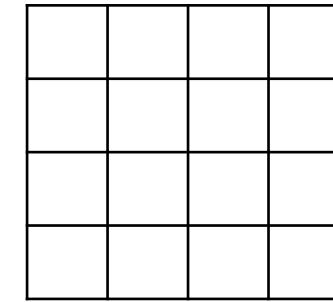
# Convolution on RGB image



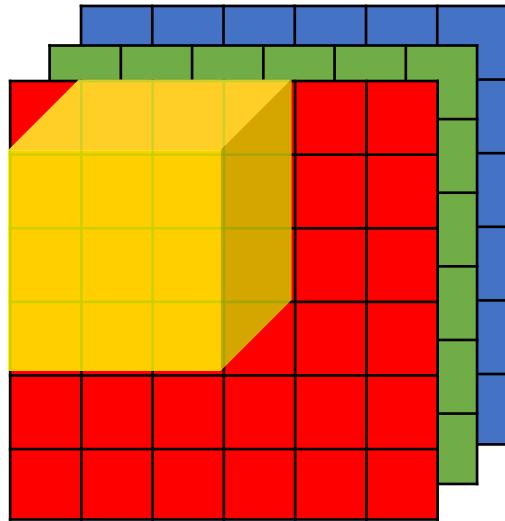
\*



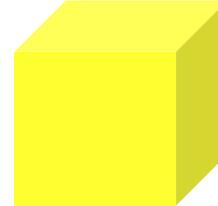
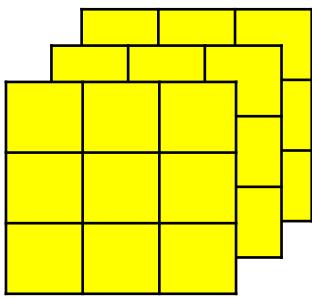
=



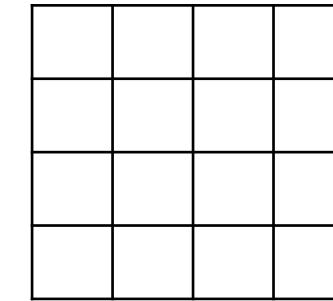
# Convolution on RGB image



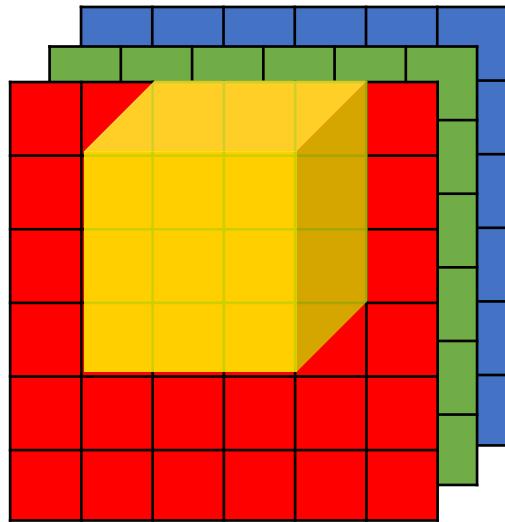
\*



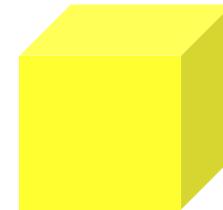
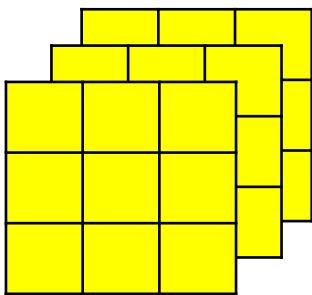
=



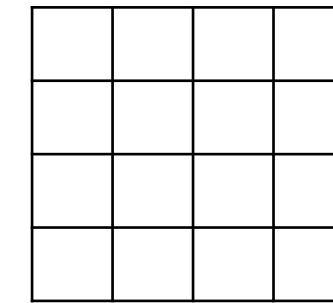
# Convolution on RGB image



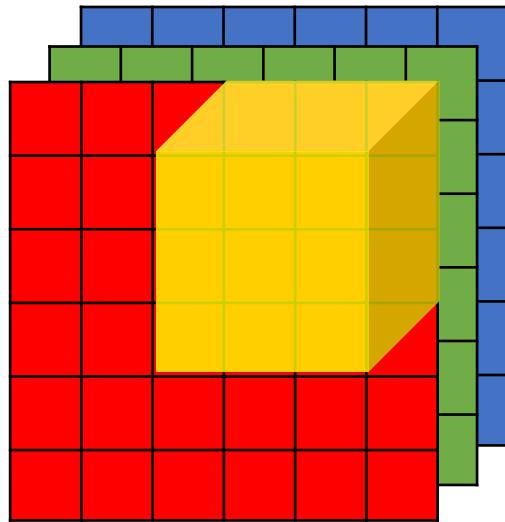
\*



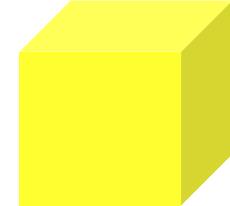
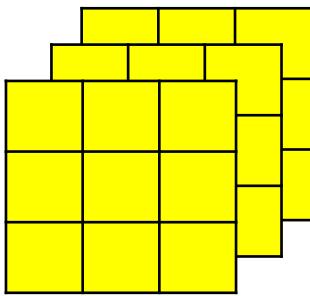
=



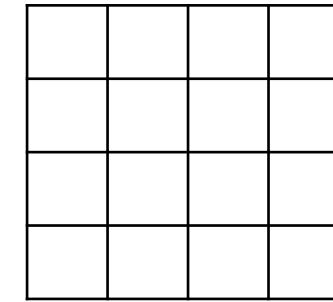
# Convolution on RGB image



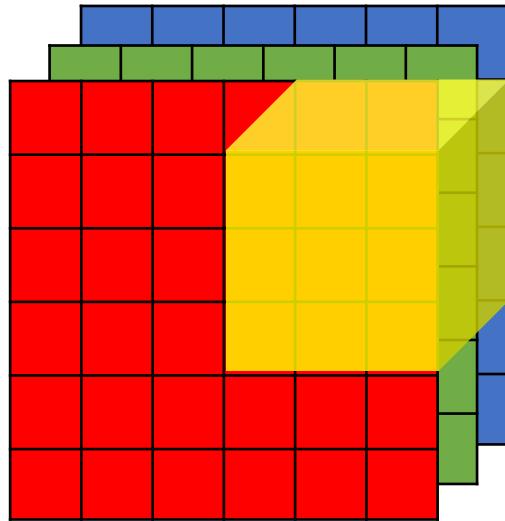
\*



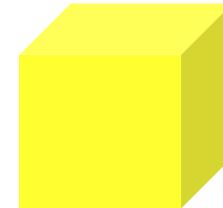
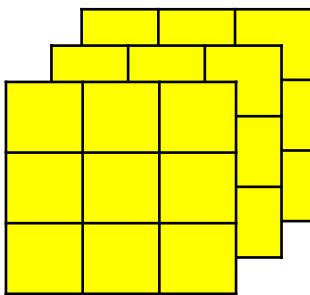
=



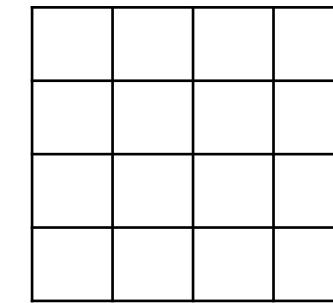
# Convolution on RGB image



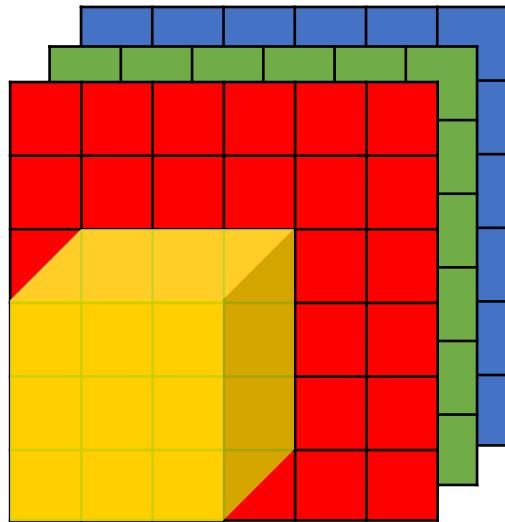
\*



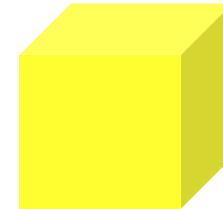
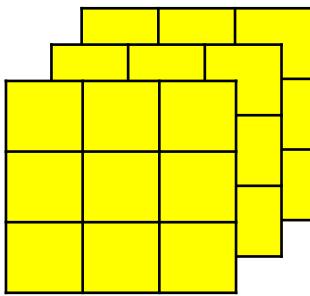
=



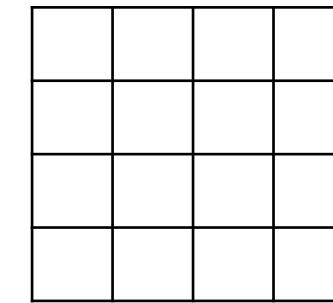
# Convolution on RGB image



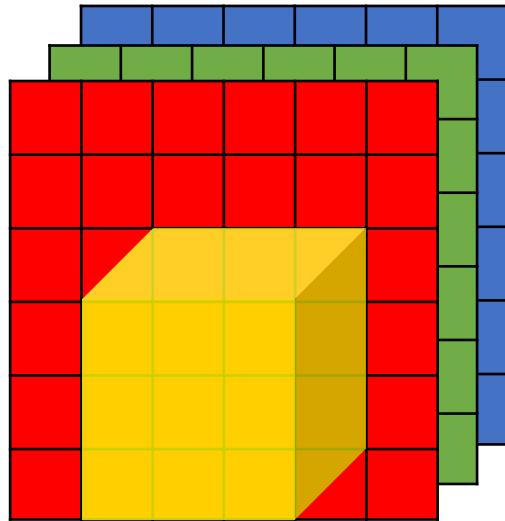
\*



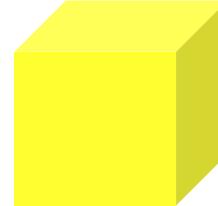
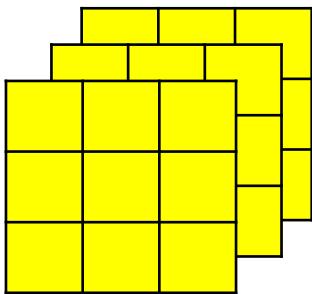
=



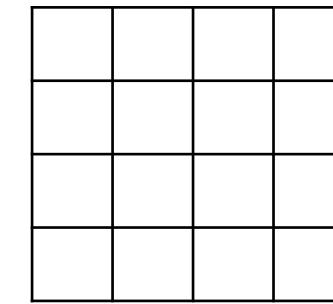
# Convolution on RGB image



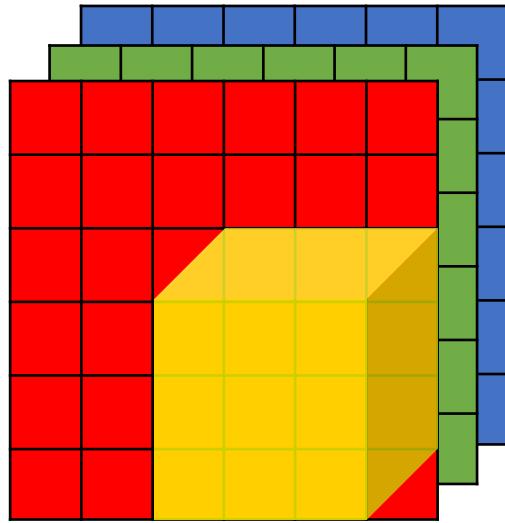
\*



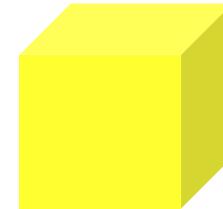
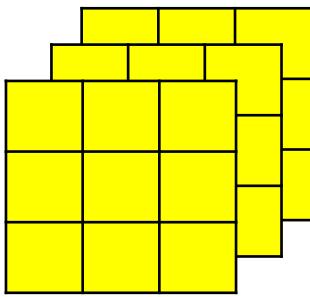
=



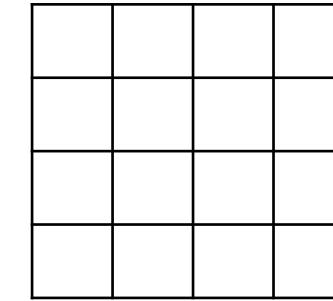
# Convolution on RGB image



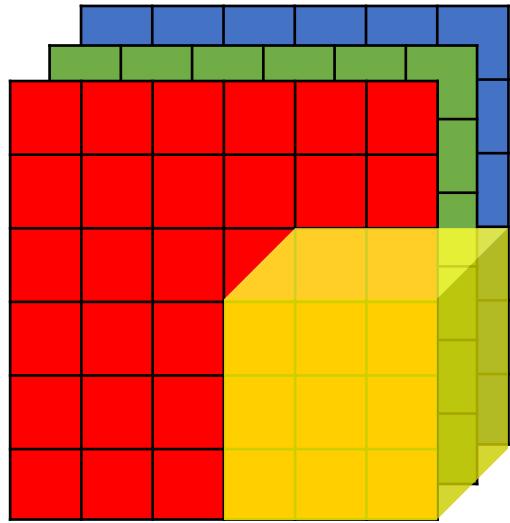
\*



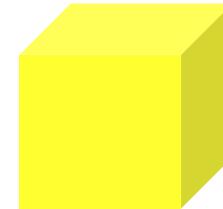
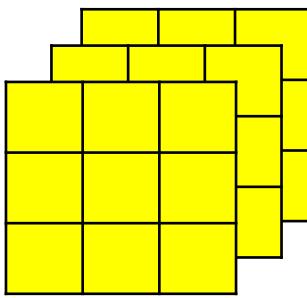
=



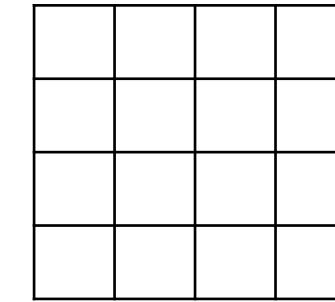
# Convolution on RGB image



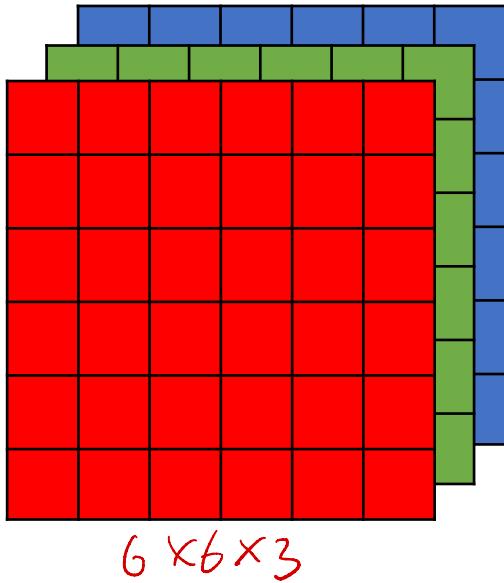
\*



=

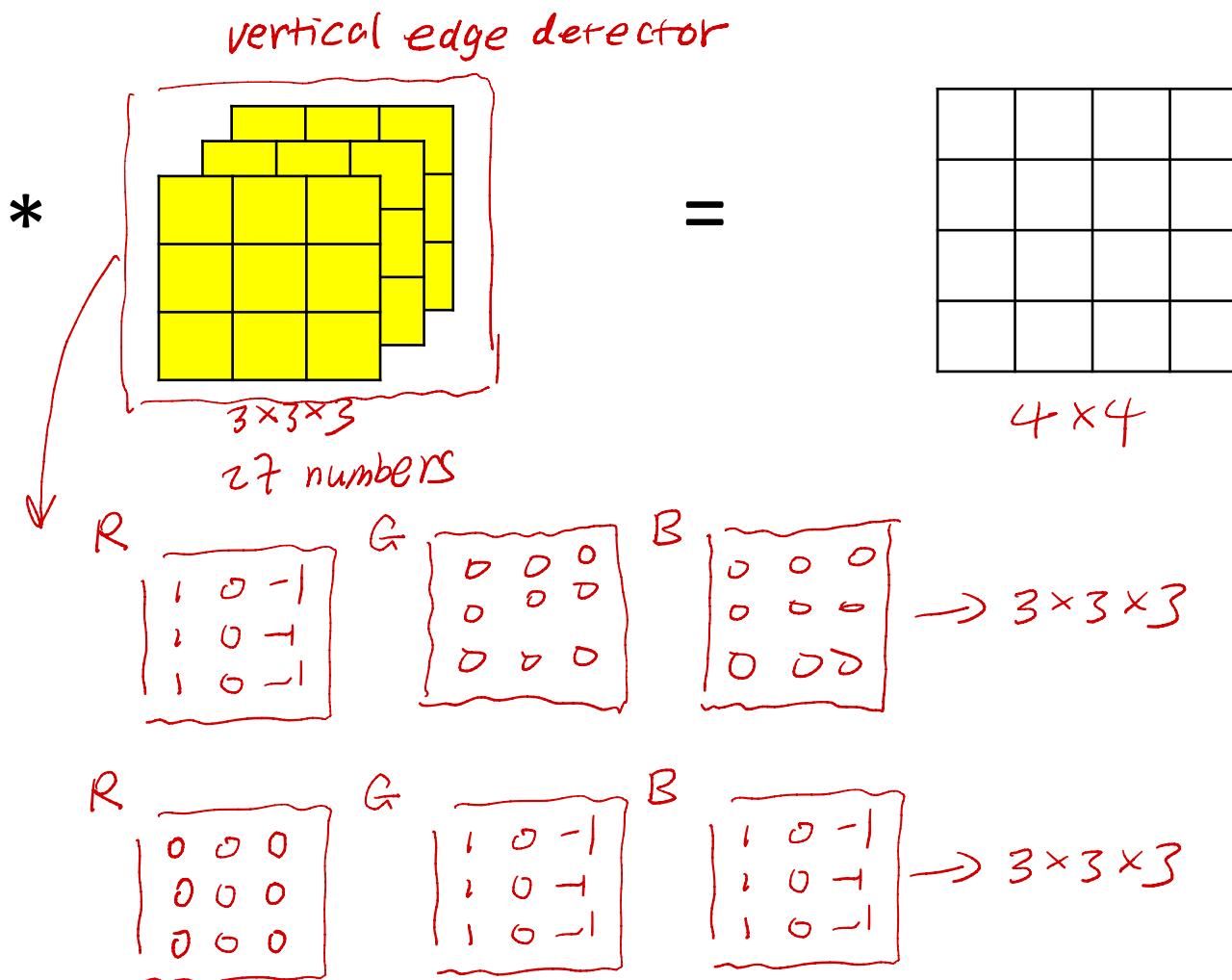


# Multiple filters

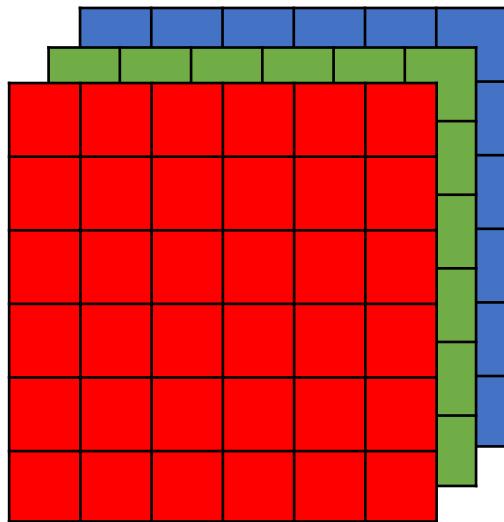


$6 \times 6 \times 3$

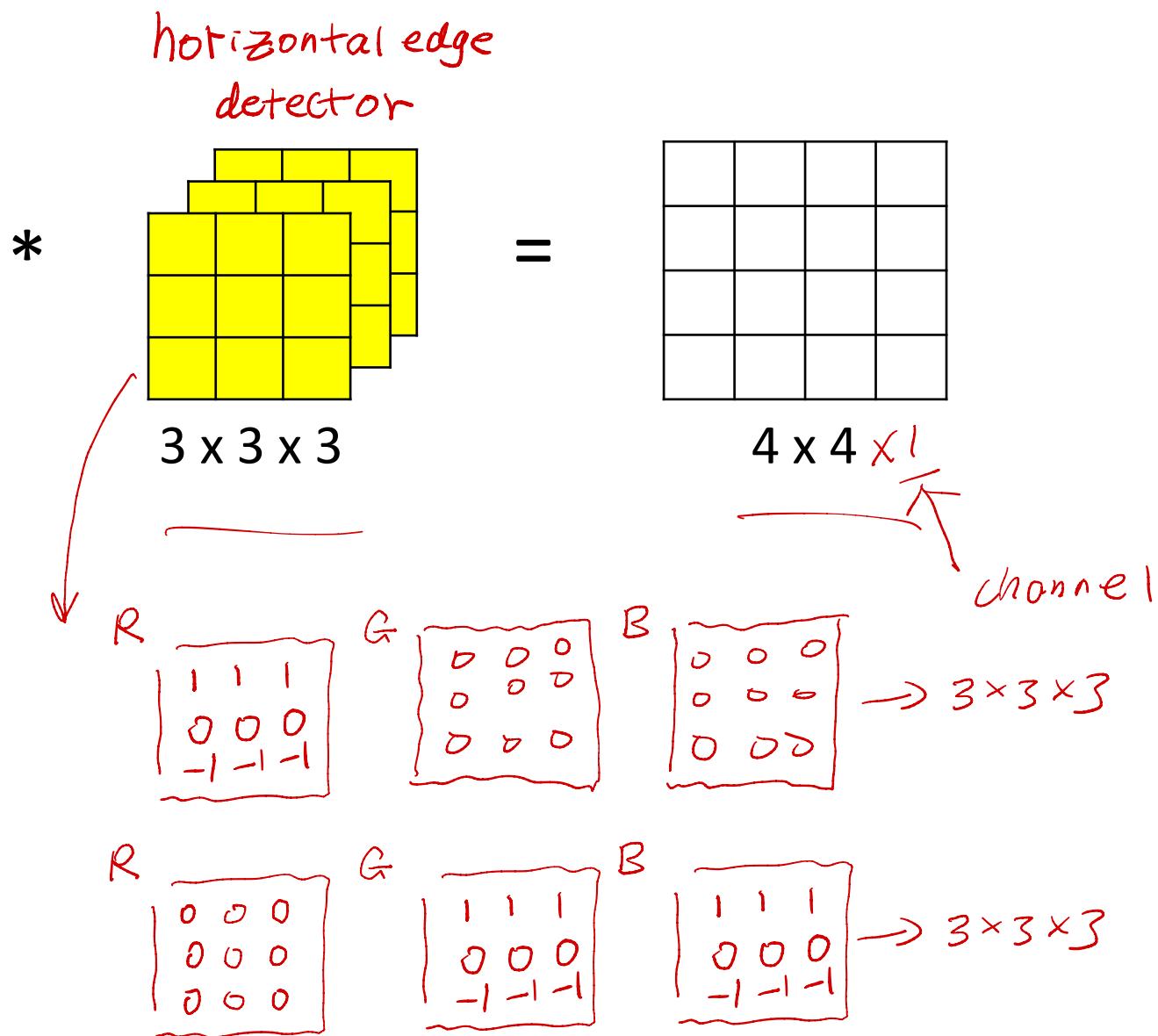
$$\begin{array}{c} \text{input} \\ \text{6x6x3} \end{array} * \begin{array}{c} \text{filter} \\ \text{3x3x3} \end{array} = \begin{array}{c} \text{output} \\ \text{4x4} \end{array}$$



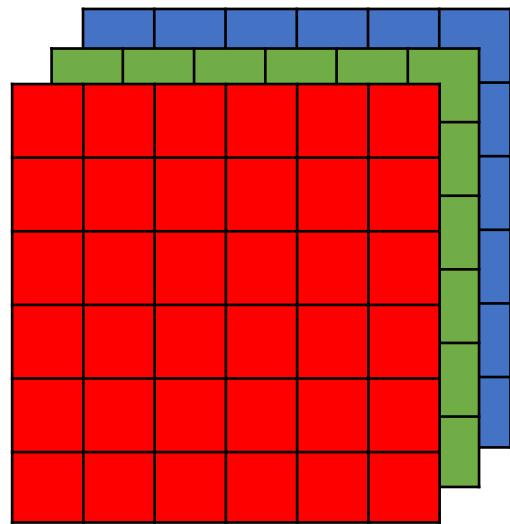
# Multiple filters



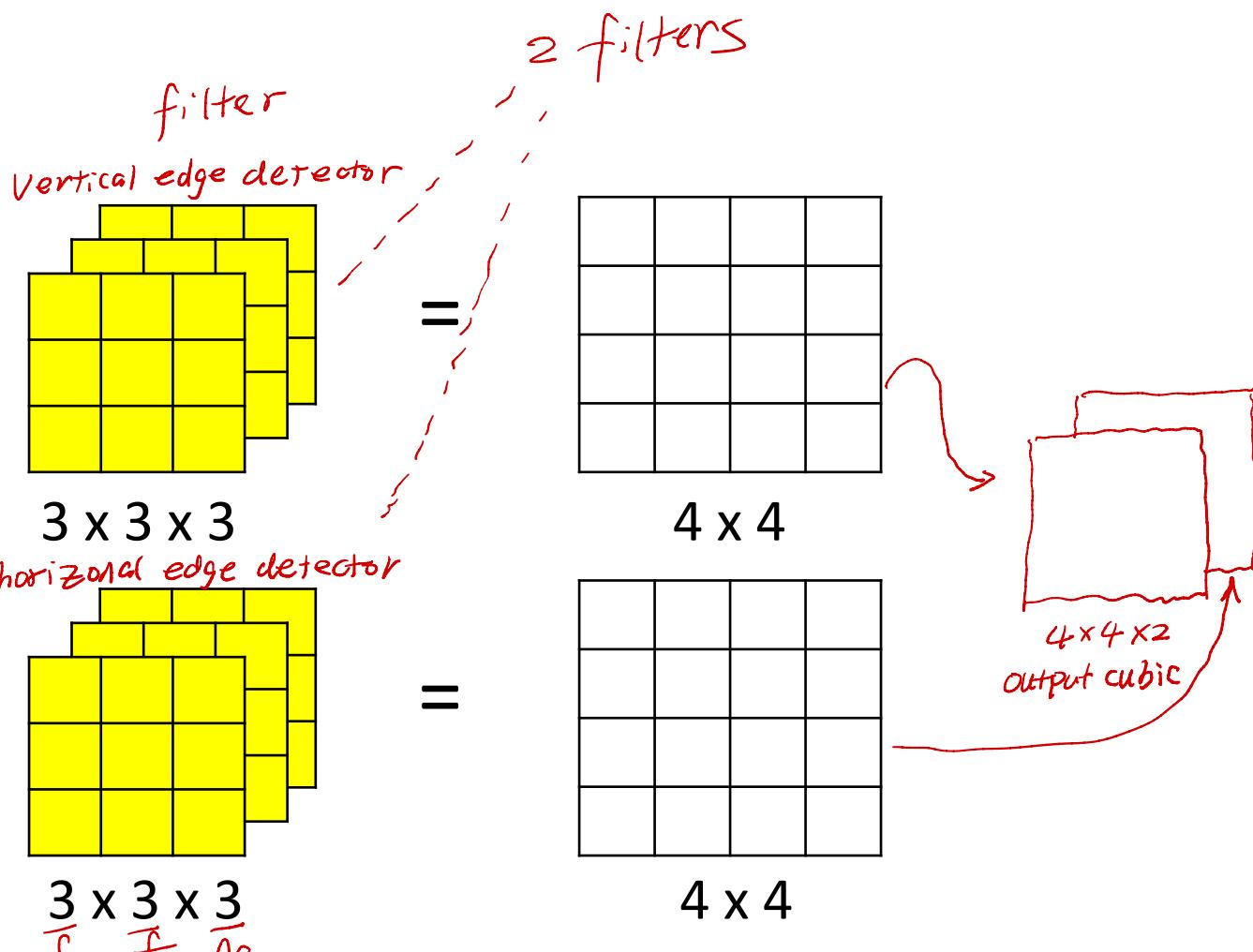
$6 \times 6 \times 3$



# Multiple filters



$6 \times 6 \times 3$   
 $\bar{f}$   
 $n$   
 $\uparrow$   
 $\# \text{ channels (depths)}$   
 $n_c$

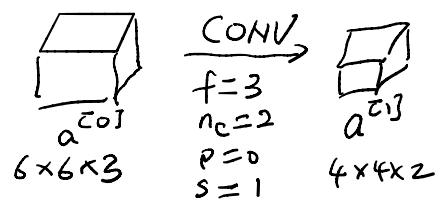
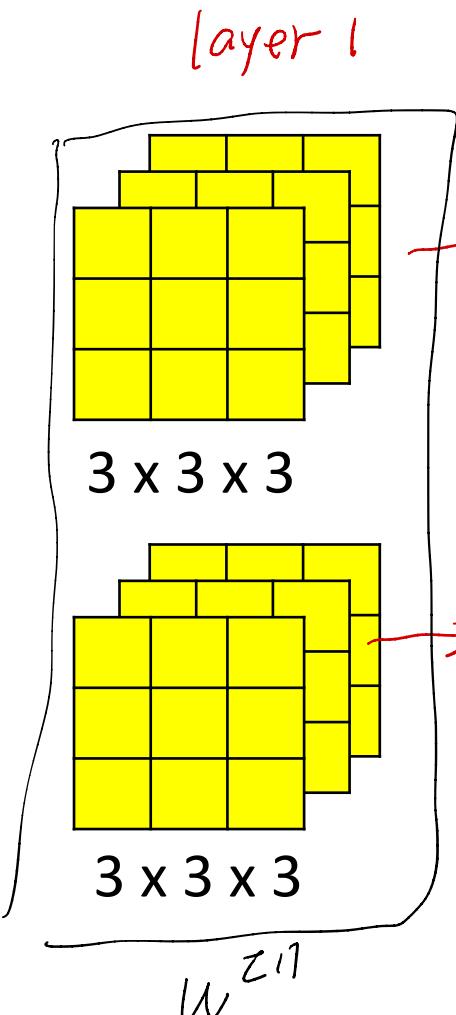
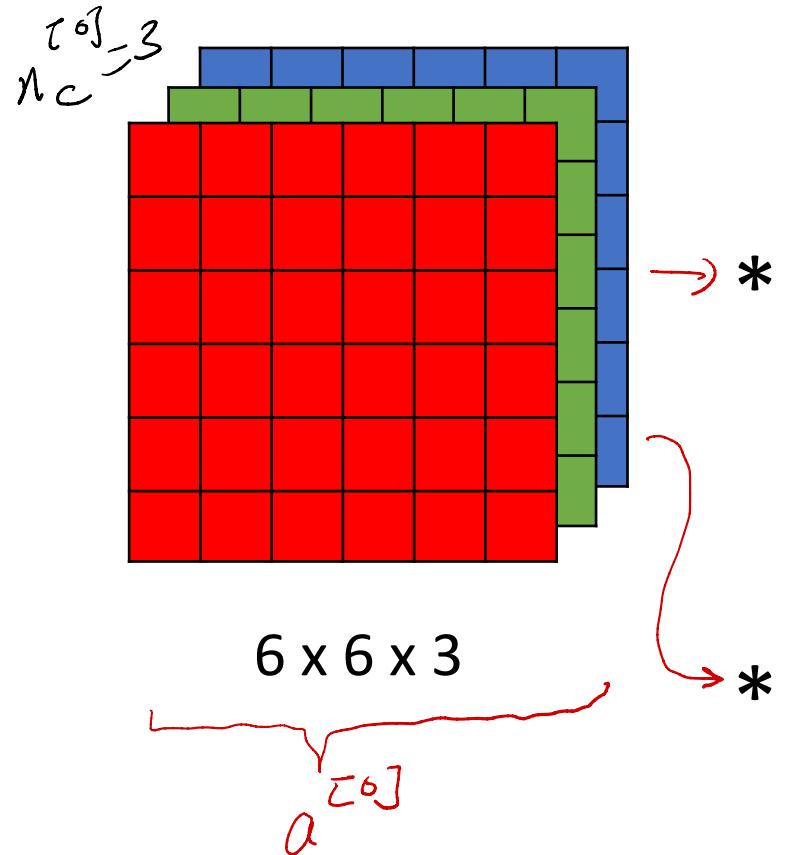


summary of dimensions

$$n \times n \times n_c * f \times f \times n_c \rightarrow (n-f+1) \times (n-f+1) \times n_c' \\ 6 \times 6 \times 3 \quad 3 \times 3 \times 3 \quad 4 \times 4 \times 2 \quad \# \text{ filters}$$

$$3 \times 3 \times 3 \times 2 + 2 = 9 \cdot 6 + 2 = 54 + 2 = 56 \text{ # parameters}$$

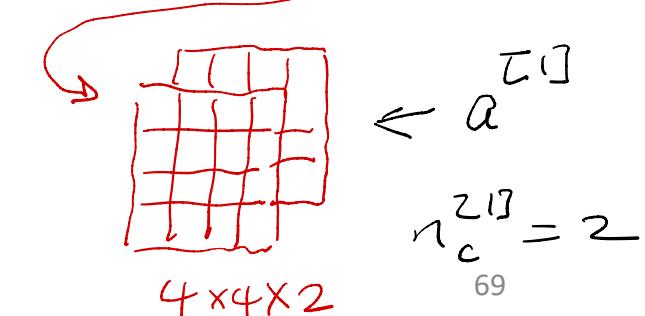
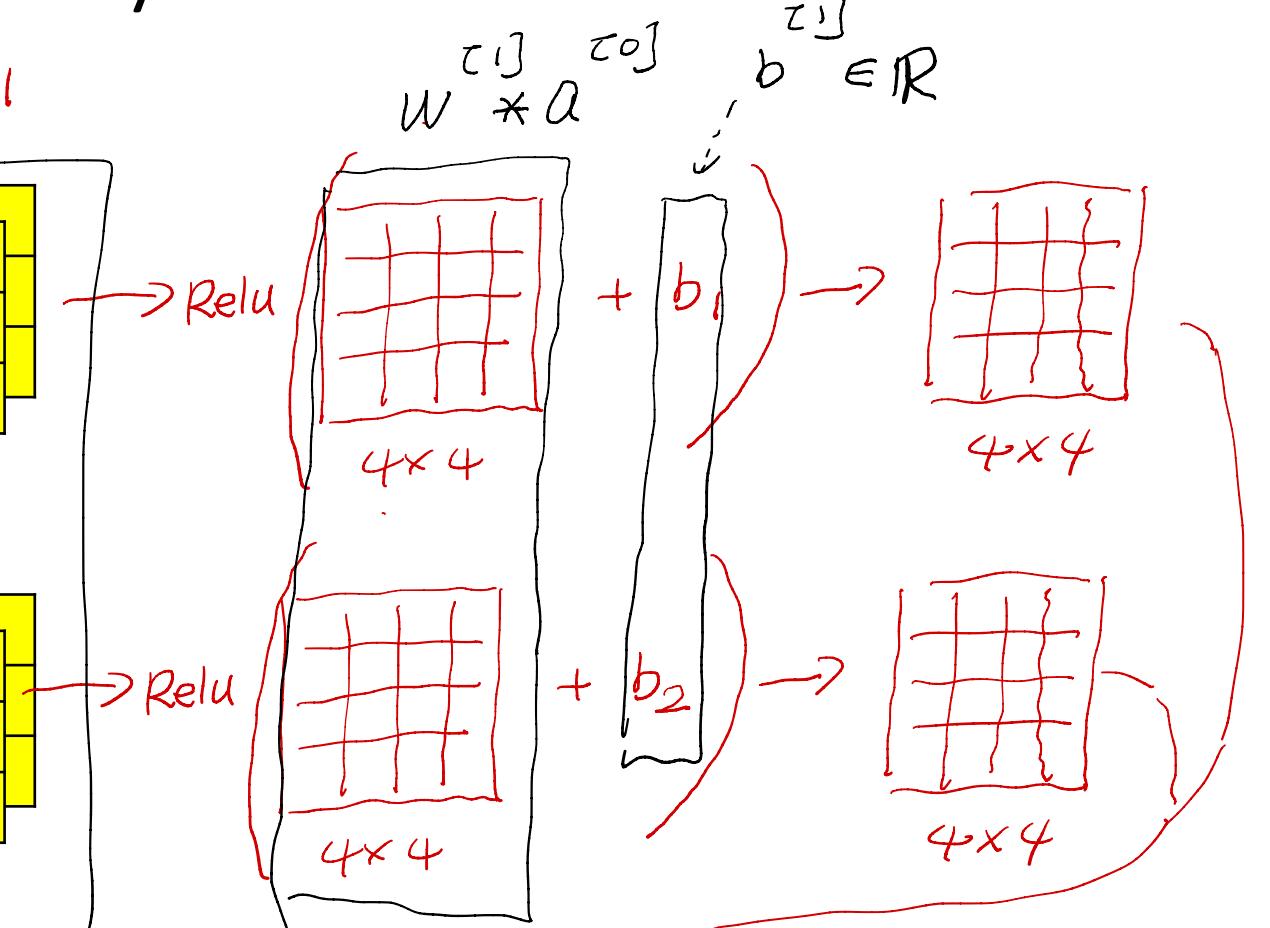
# Example of a convolution layer



$$z^{(1)} = W^{(1)} * a^{(0)} + b^{(1)}$$

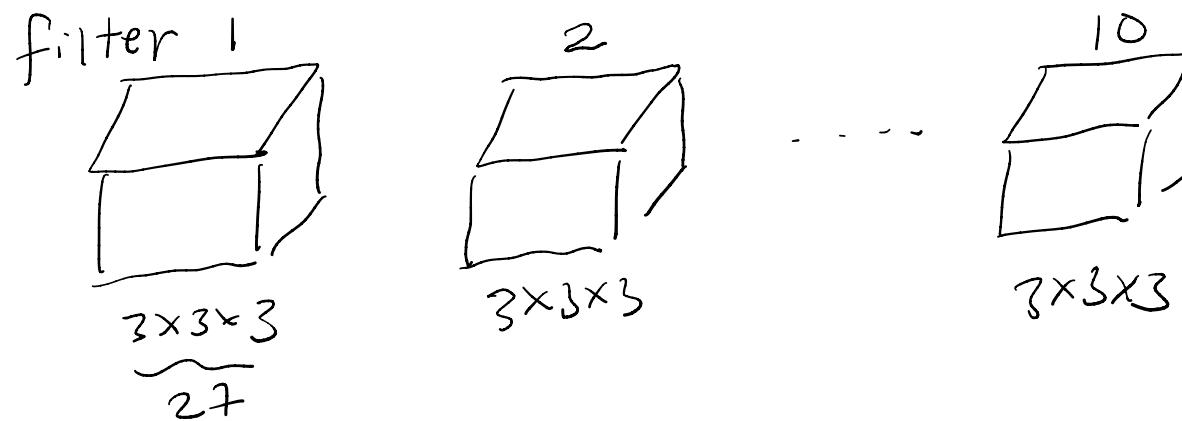
$$a^{(1)} = g(z^{(1)})$$

Now I am ready to introduce a full picture of a convolution layer.



# Number of parameters in one layer

If you have 10 filters that have  $3 \times 3 \times 3$  in one layer of a neural network, how many parameters does that layer have?



$$\underbrace{27 \times 10}_{\text{bias}} + \underbrace{10}_{\text{bias}} = 270 + 10 = 280$$

# Summary of notation

If layer  $\ell$  is a convolution layer:

**INPUT:**  $n_H^{[\ell-1]} \times n_W^{[\ell-1]} \times n_C^{[\ell-1]}$

Hyper Parameters:

$f^{[\ell]}$  = filter size

$p^{[\ell]}$  = padding

$s^{[\ell]}$  = stride

$n_C^{[\ell]}$  = number of filters

Each filter is  $= f^{[\ell]} \times f^{[\ell]} \times n_C^{[\ell-1]}$

Activations (Output)  $= n_H^{[\ell]} \times n_W^{[\ell]} \times n_C^{[\ell]}$

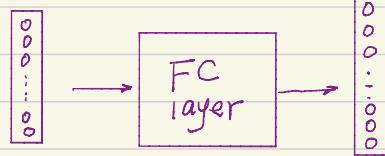
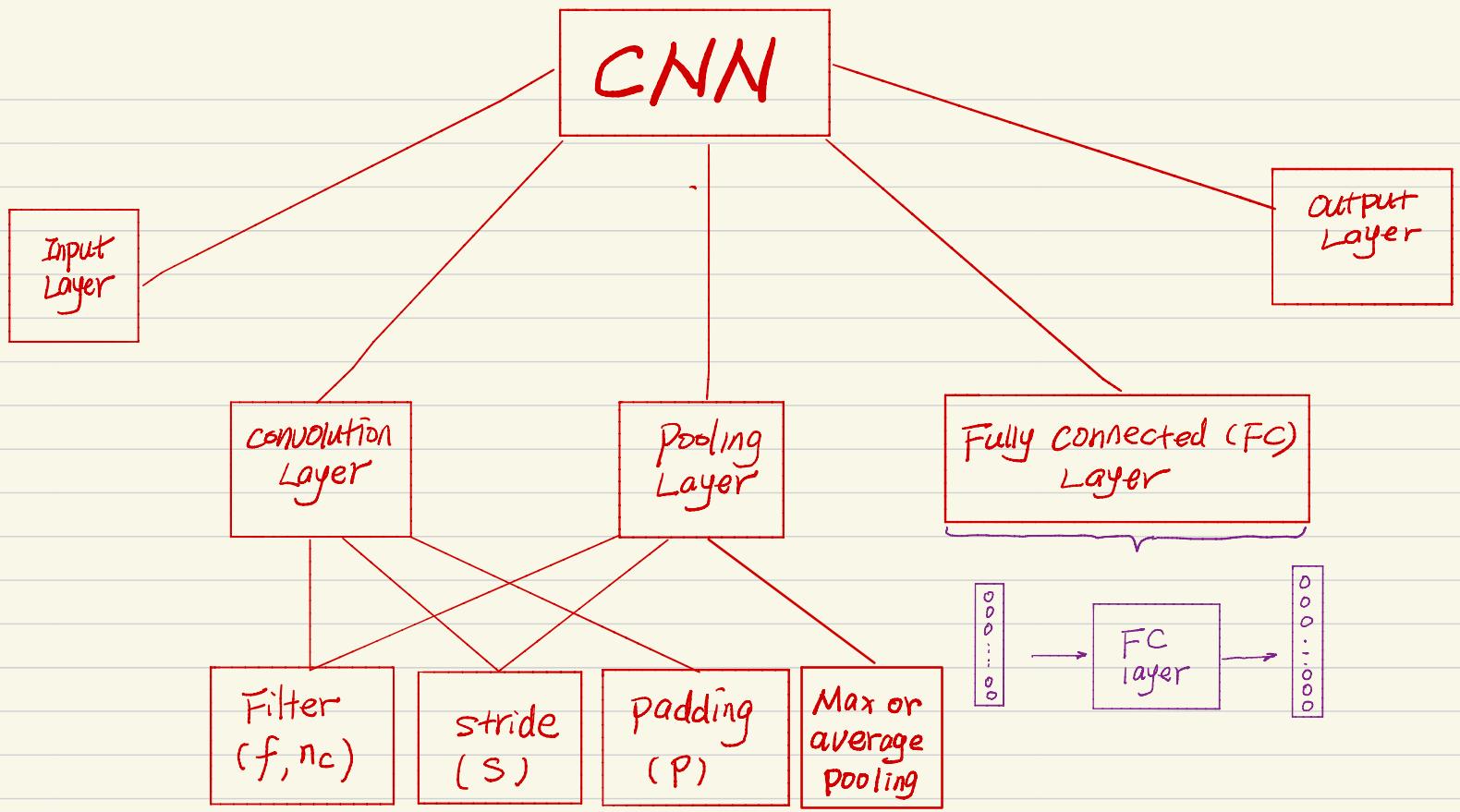
# weights  $= f^{[\ell]} \times f^{[\ell]} \times n_C^{[\ell-1]} \times \underbrace{n_C^{[\ell]}}_{\# \text{ filters in layer } \ell}$

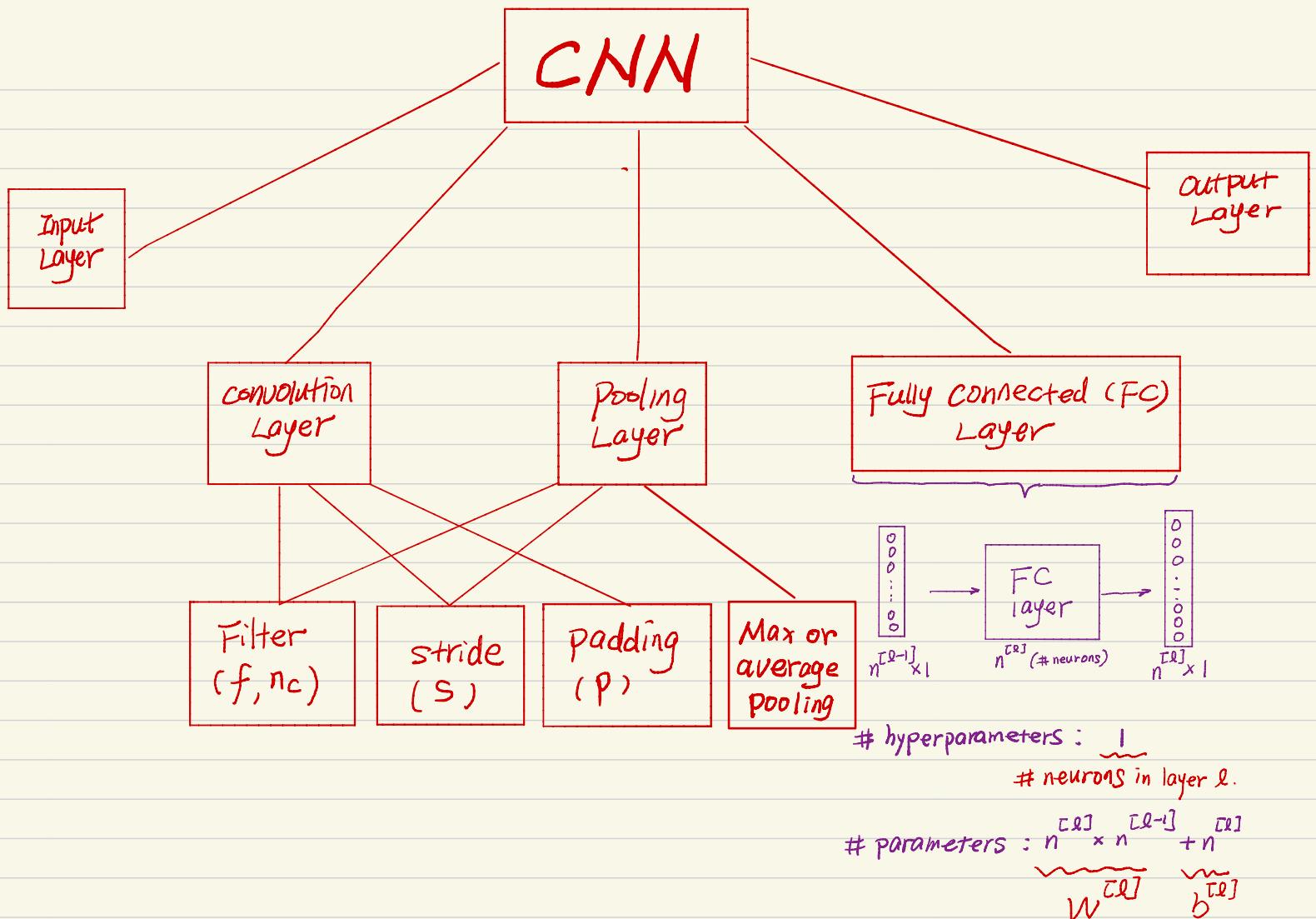
# bias  $= n_C^{[\ell]}$

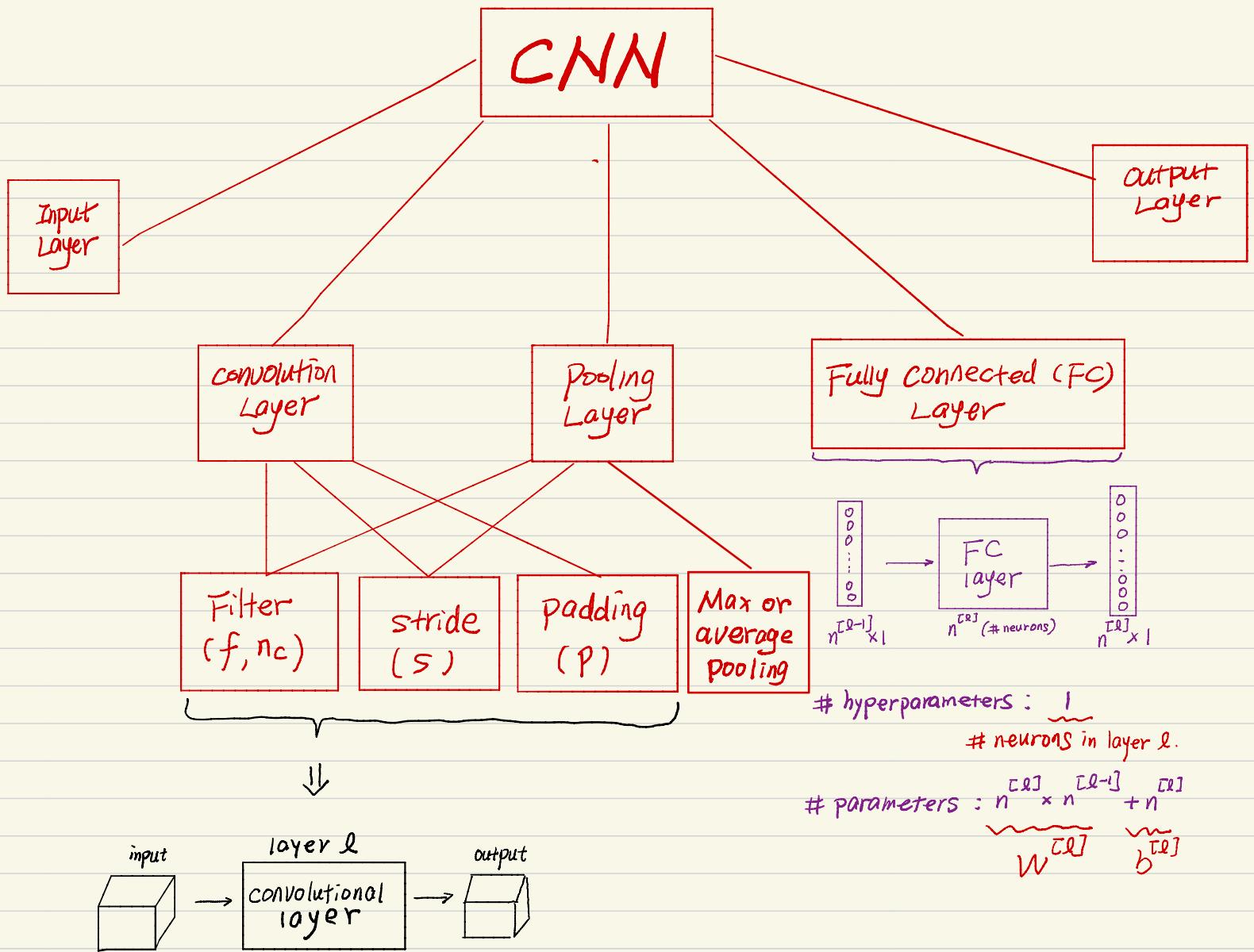
$$\text{OUTPUT} = n_H^{[\ell]} = \left\lfloor \frac{n_H^{[\ell-1]} + 2p^{[\ell]} - f^{[\ell]}}{s^{[\ell]}} + 1 \right\rfloor$$

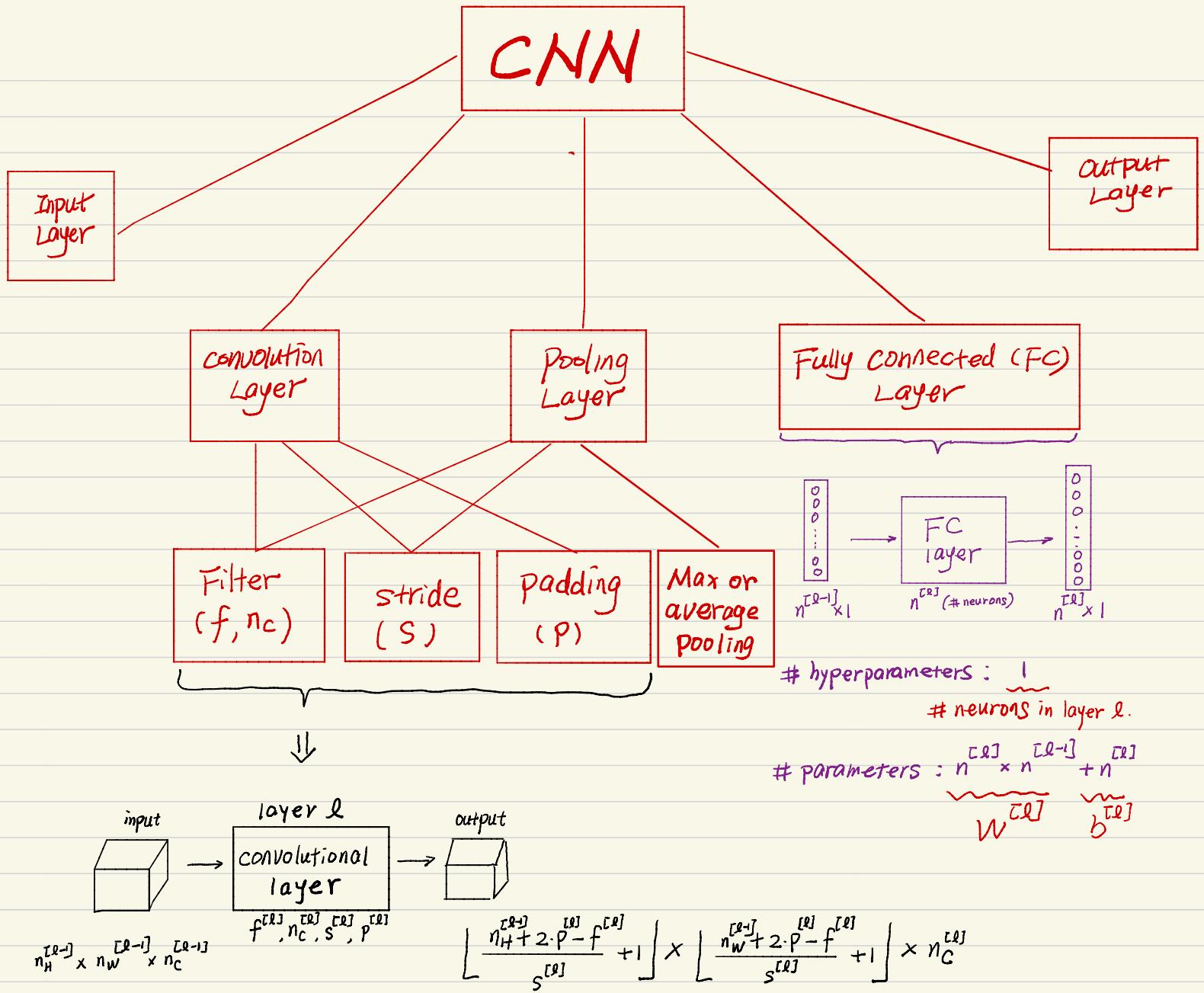
$$n_W^{[\ell]} = \left\lfloor \frac{n_W^{[\ell-1]} + 2p^{[\ell]} - f^{[\ell]}}{s^{[\ell]}} + 1 \right\rfloor$$

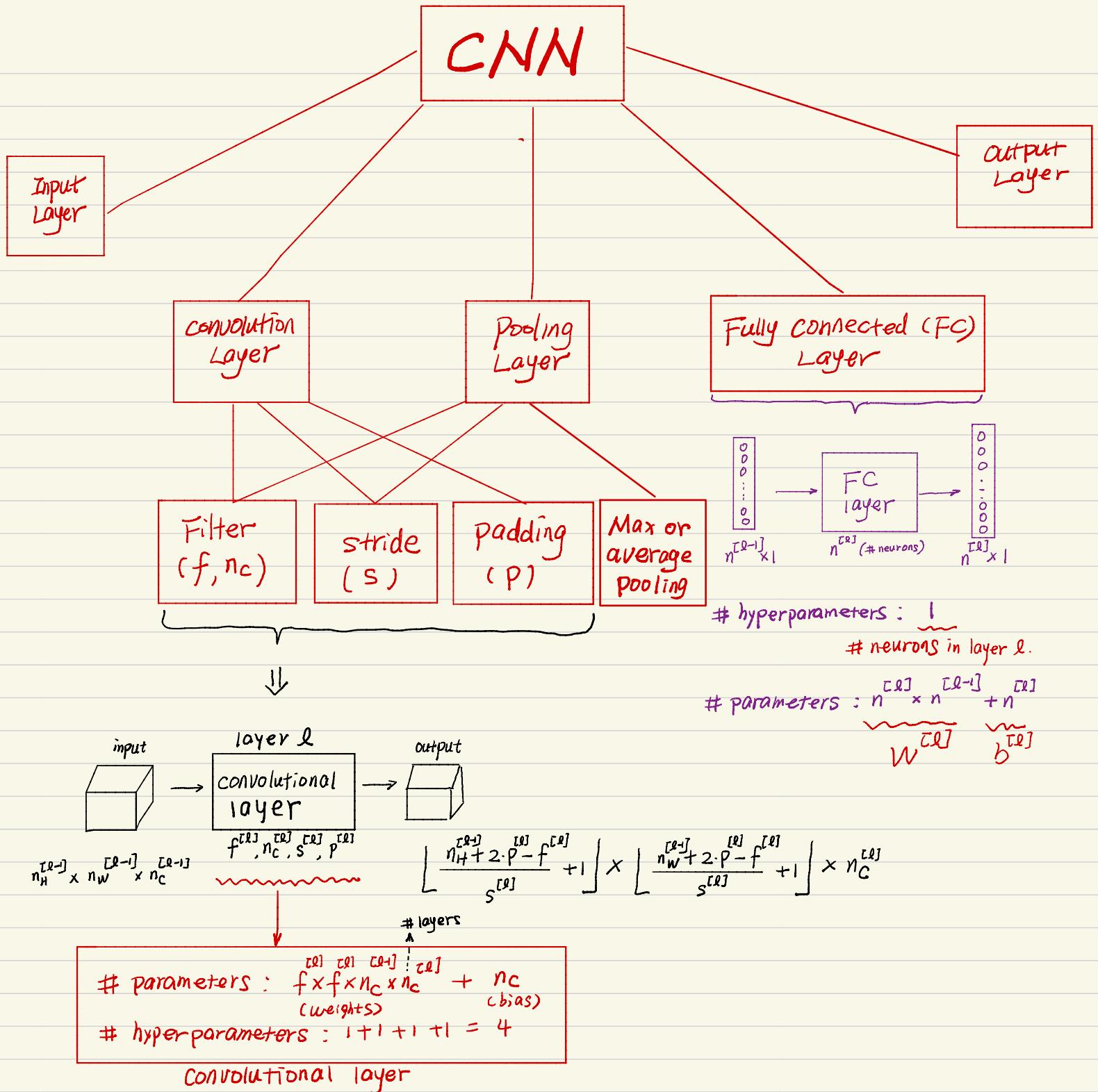
$\overbrace{n_C^{[\ell]}}^{\# \text{ channels of the output volume}}$  = Number of filters











# Pooling layer: Max pooling

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

4x4

# Pooling layer: Max pooling

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

$4 \times 4$


output of this specific pooling is  
 $2 \times 2$  matrix

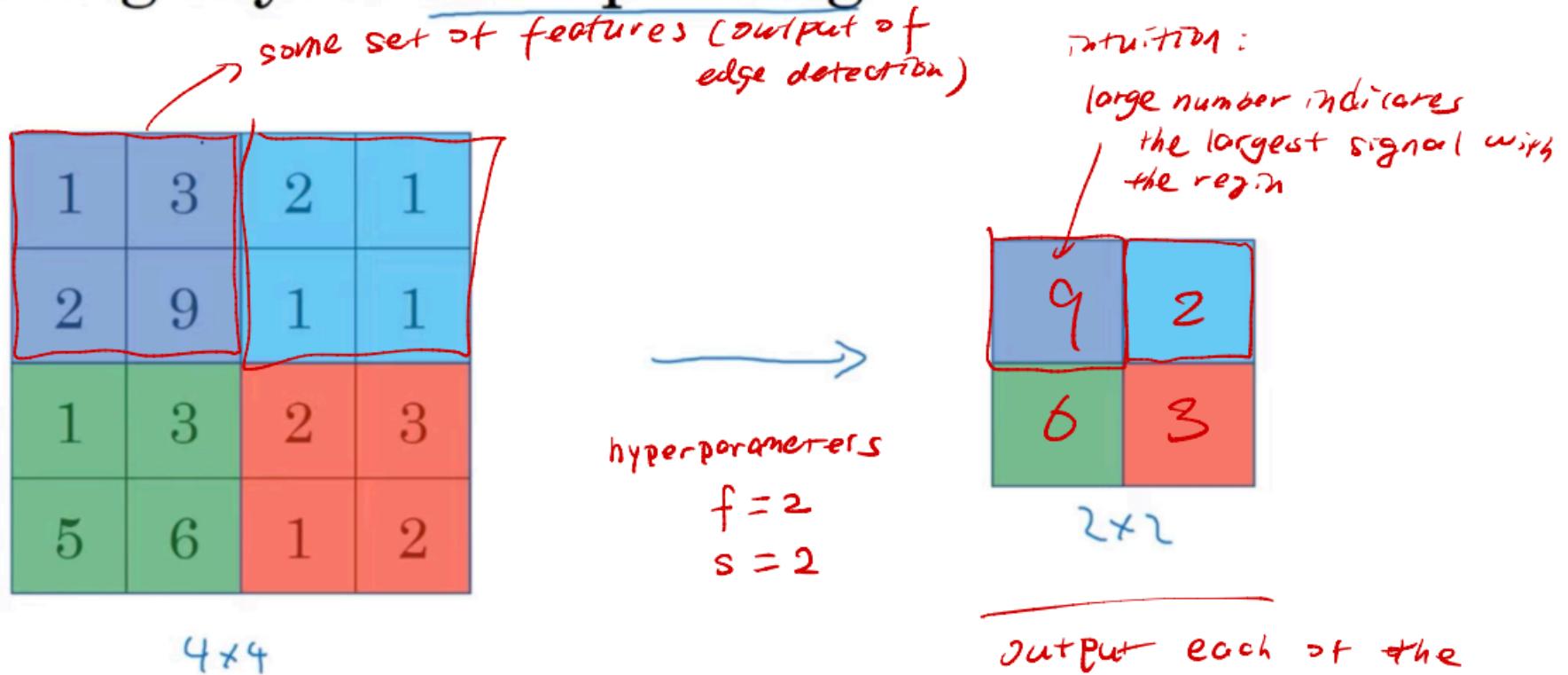
# Pooling layer: Max pooling

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




$4 \times 4$   
break into different regions  
with different colors

# Pooling layer: Max pooling



Output each of the output will be max of the corresponding shaded region

# Pooling layer: Max pooling

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0
5	6	1	2	9

5x5




$$\begin{aligned} & \boxed{f=3} \quad 3 \times 3 \\ & s=1 \quad \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \\ & = \left\lfloor \frac{5+0-3}{1} + 1 \right\rfloor = 3 \end{aligned}$$

an example with different hyperparameters

# Pooling layer: Max pooling

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0
5	6	1	2	9

$5 \times 5$



9		

$$\boxed{f=3 \\ s=1}$$

$$3 \times 3 \\ \lfloor \frac{n+2p-f}{s} + 1 \rfloor \\ = \lfloor \frac{5+0-3}{1} + 1 \rfloor = 3$$

an example with different hyperparameters

# Pooling layer: Max pooling

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0
5	6	1	2	9

5x5



9	9	

$$\begin{cases} f = 3 \\ s = 1 \end{cases}$$

$$\begin{aligned} & 3 \times 3 \\ & \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \\ & = \left\lfloor \frac{5+0-3}{1} + 1 \right\rfloor = 3 \end{aligned}$$

# Pooling layer: Max pooling

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0
5	6	1	2	9

5 × 5



9	9	5

$$\begin{aligned} & \boxed{\begin{array}{l} f=3 \\ s=1 \end{array}} \quad 3 \times 3 \\ & \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \\ & = \left\lfloor \frac{5+0-3}{1} + 1 \right\rfloor = 3 \end{aligned}$$

# Pooling layer: Max pooling

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0
5	6	1	2	9

5 × 5



9	9	5
9		

$$\begin{aligned} & \boxed{\begin{array}{l} f=3 \\ s=1 \end{array}} \quad 3 \times 3 \\ & \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \\ & = \left\lfloor \frac{5+0-3}{1} + 1 \right\rfloor = 3 \end{aligned}$$

# Pooling layer: Max pooling

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0
5	6	1	2	9

5 × 5



9	9	5
9	9	

$$\begin{aligned} & \boxed{\begin{array}{l} f=3 \\ s=1 \end{array}} \quad 3 \times 3 \\ & \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \\ & = \left\lfloor \frac{5+0-3}{1} + 1 \right\rfloor = 3 \end{aligned}$$

# Pooling layer: Max pooling

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0
5	6	1	2	9

5 × 5



9	9	5
9	9	5

$$\begin{array}{|l|} \hline f=3 \\ s=1 \\ \hline \end{array} \quad \begin{array}{l} 3 \times 3 \\ \lfloor \frac{n+2p-f}{s} + 1 \rfloor \\ = \lfloor \frac{5+0-3}{1} + 1 \rfloor = 3 \end{array}$$

# Pooling layer: Max pooling

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0
5	6	1	2	9

5 × 5



9	9	5
9	9	5
8		

$$\begin{array}{|l|} \hline f=3 \\ s=1 \\ \hline \end{array} \quad \begin{array}{l} 3 \times 3 \\ \lfloor \frac{n+2p-f}{s} + 1 \rfloor \\ = \lfloor \frac{5+0-3}{1} + 1 \rfloor = 3 \end{array}$$

# Pooling layer: Max pooling

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0
5	6	1	2	9

5 × 5



9	9	5
9	9	5
8	6	

$$\begin{array}{|l|l|} \hline f=3 & 3 \times 3 \\ s=1 & \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \\ \hline \end{array}$$
$$= \left\lfloor \frac{5+0-3}{1} + 1 \right\rfloor = 3$$

# Pooling layer: Max pooling

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0
5	6	1	2	9

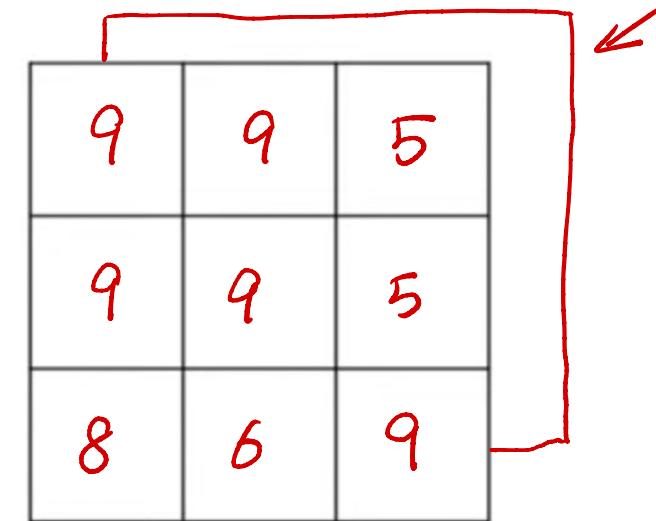
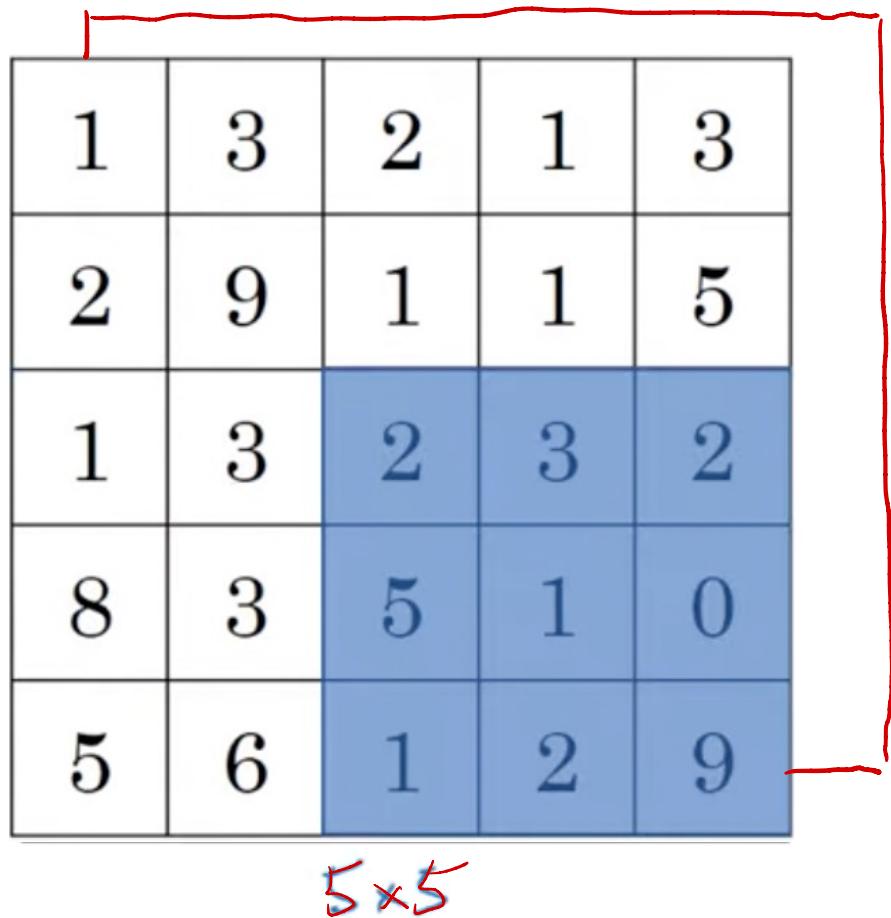
5 × 5



9	9	5
9	9	5
8	6	9

$$\begin{array}{|l|} \hline f=3 \\ s=1 \\ \hline \end{array} \quad \begin{array}{l} 3 \times 3 \\ \lfloor \frac{n+2p-f}{s} + 1 \rfloor \\ = \lfloor \frac{5+0-3}{1} + 1 \rfloor = 3 \end{array}$$

# Pooling layer: Max pooling



$$\boxed{\begin{array}{l} f=3 \\ s=1 \end{array}} \quad \begin{array}{l} 3 \times 3 \\ \lfloor \frac{n+2p-f}{s} + 1 \rfloor \\ = \lfloor \frac{5+0-3}{1} + 1 \rfloor = 3 \end{array}$$

# Pooling layer: Average pooling

1	3	2	1
2	9	1	1
1	4	2	3
5	6	1	2



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


# Pooling layer: Average pooling

1	3	2	1
2	9	1	1
1	4	2	3
5	6	1	2



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

3.75	1.25
4	2

# Summary of pooling

Hyperparameters:

f : filter size

Popular settings

$f=2, s=2$

s : stride

$f=3, s=2$

Max or average pooling

(typically, no padding)

NO parameters to learn.

INPUT:  $n_H \times n_W \times n_C$



=

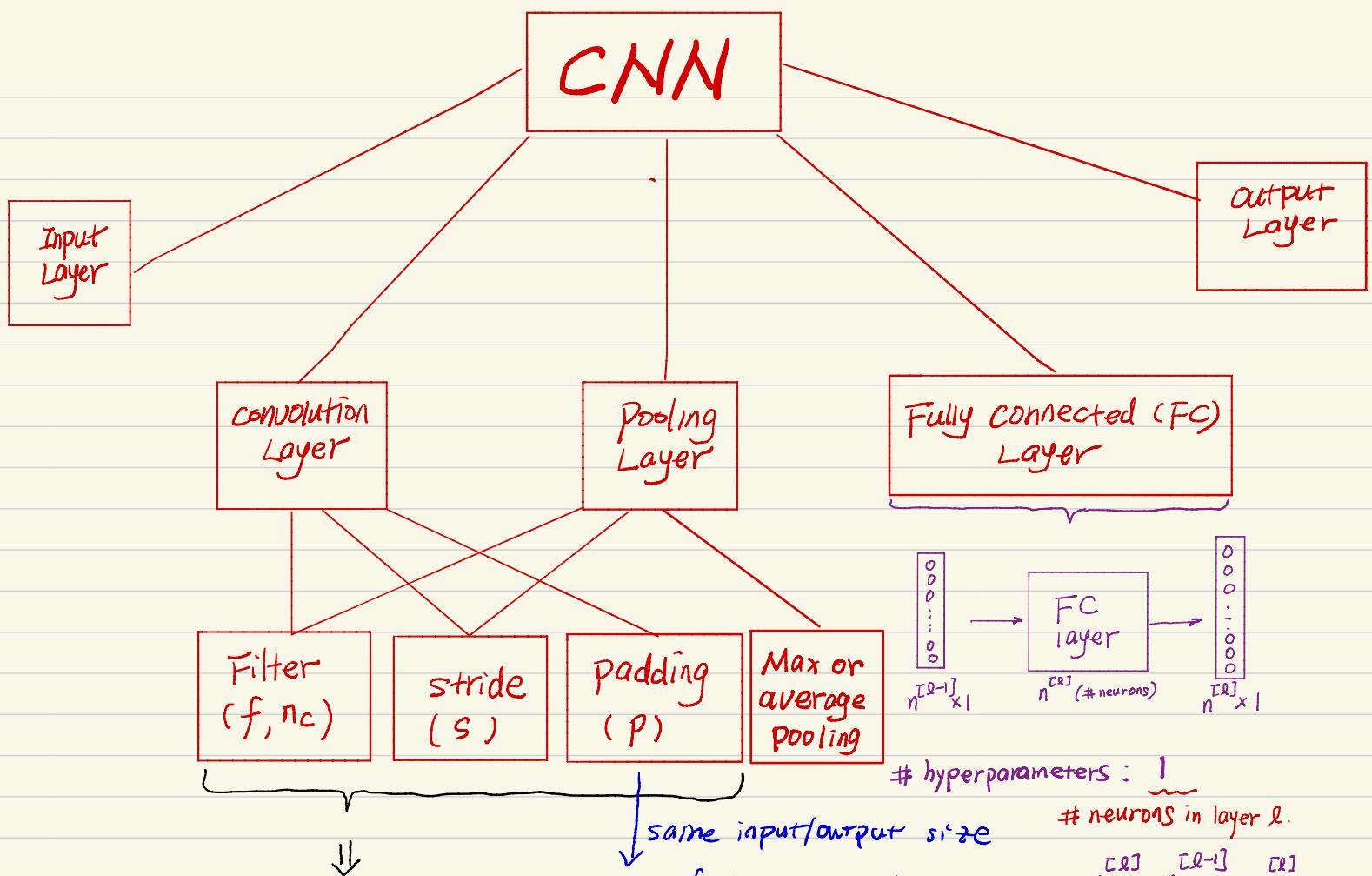
OUTPUT:

$$\left\lfloor \frac{n_H-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_W-f}{s} + 1 \right\rfloor \times n_C$$

equal

Pooling is applied  
to each channel  
separately.

# CNN

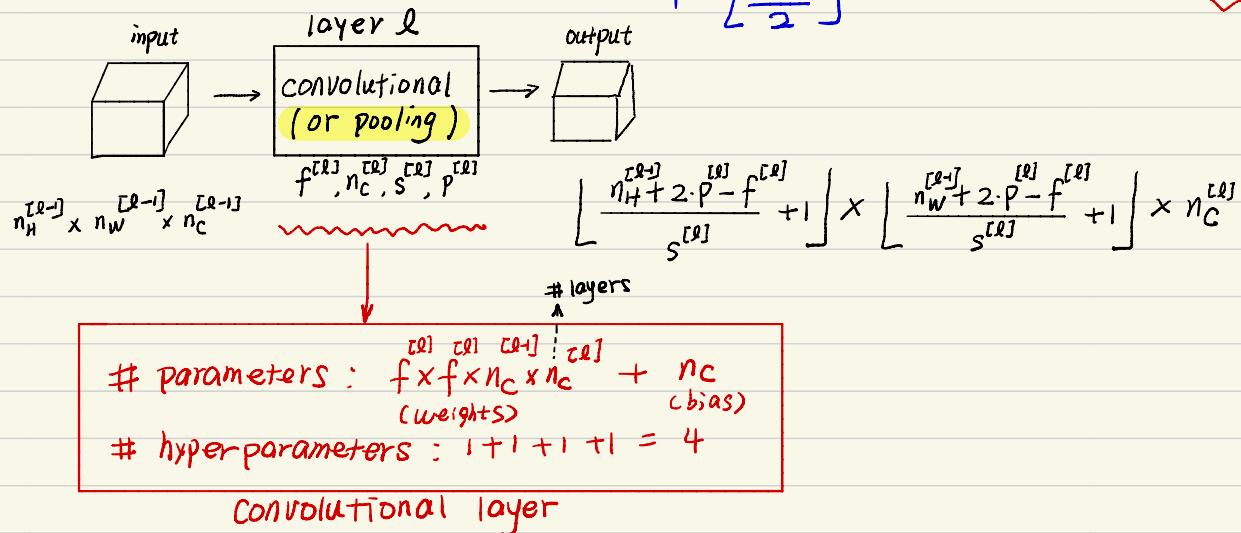


# hyperparameters : 1

# neurons in layer  $l$ .

$$\# \text{parameters} : n^{[l-1]} \times n^{[l-1]} + n^{[l]}$$

$$W^{[l]} \quad b^{[l]}$$



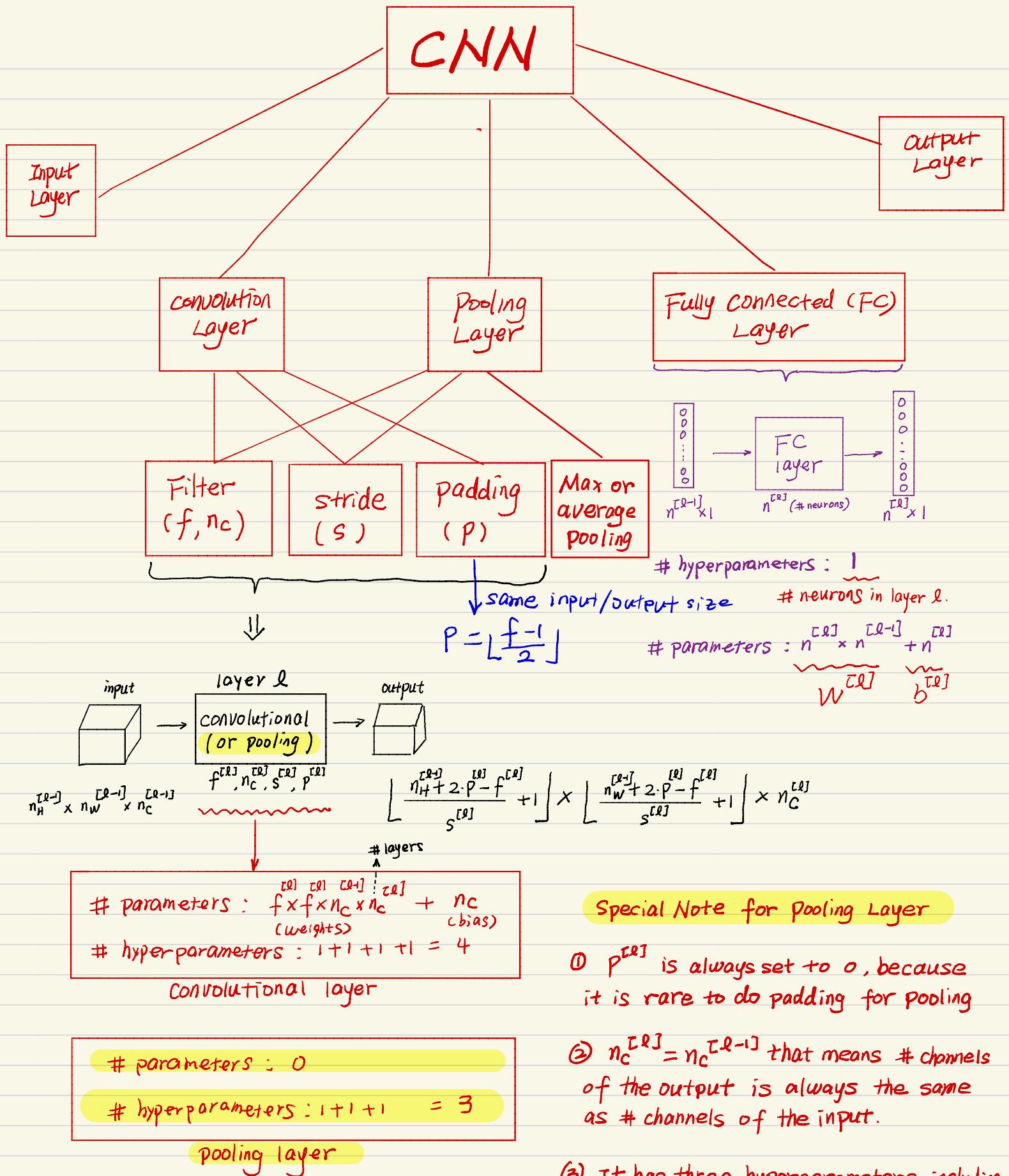
$$\# \text{parameters} : f \times f \times n_C^{[l-1]} \times n_C^{[l]} + n_C^{[l]}$$

$$\# \text{hyperparameters} : 1 + 1 + 1 + 1 = 4$$

# parameters : 0

# hyperparameters : 1 + 1 + 1 = 3

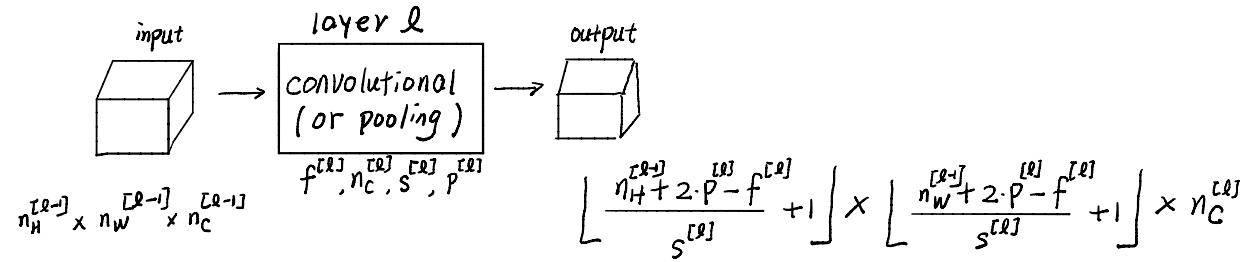
pooling layer



# Example CNN network - 1

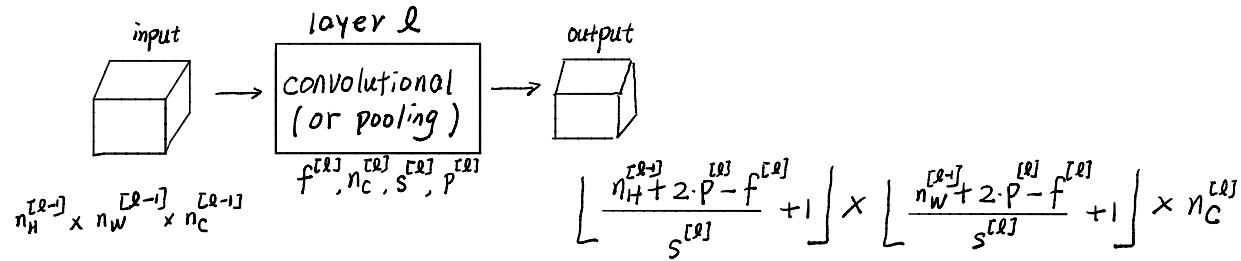
INPUT (39x39x3)
CONV1 ( $f=3, n_C = 10, s=1, p=0$ )
CONV2 ( $f=5, n_C = 20, s=2, p=0$ )
CONV3 ( $f=5, n_C = 40, s=2, p=0$ )
softmax (10,1)

# Example CNN network - 1

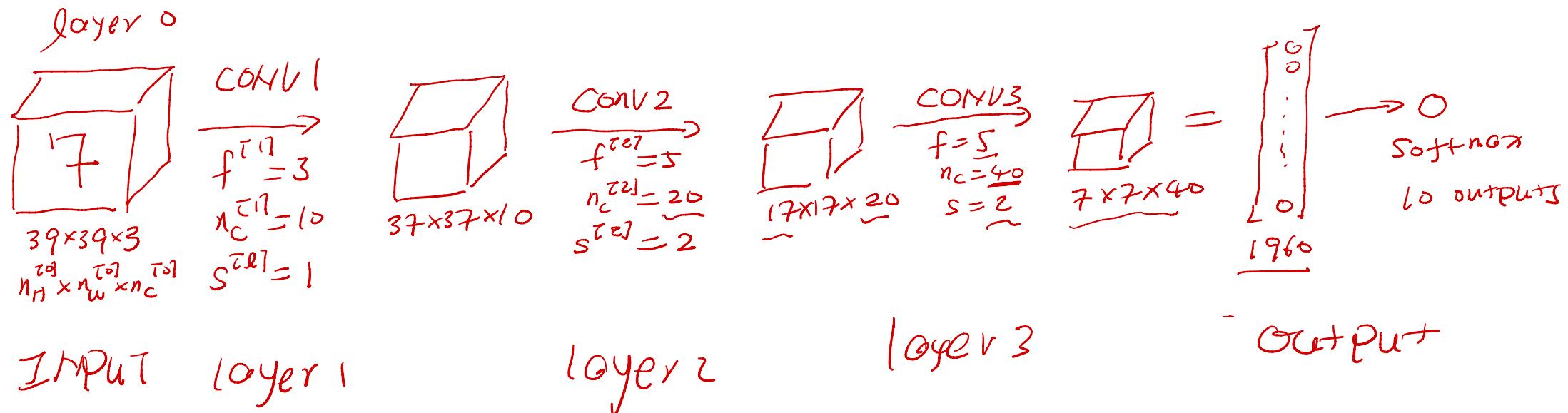


INPUT (39x39x3)
CONV1 ( $f=3, n_C = 10, s=1, p=0$ )
CONV2 ( $f=5, n_C = 20, s=2, p=0$ )
CONV3 ( $f=5, n_C = 40, s=2, p=0$ )
softmax (10,1)

# Example CNN network - 1

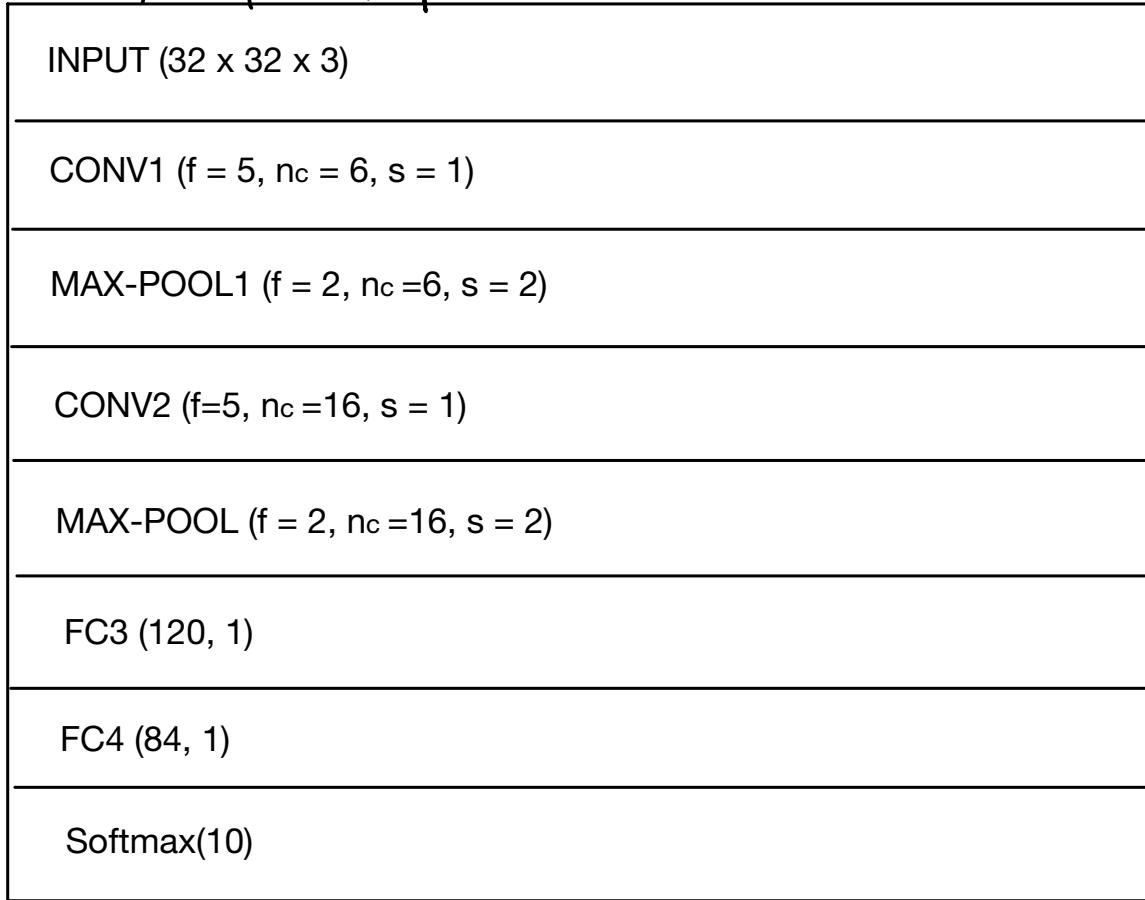


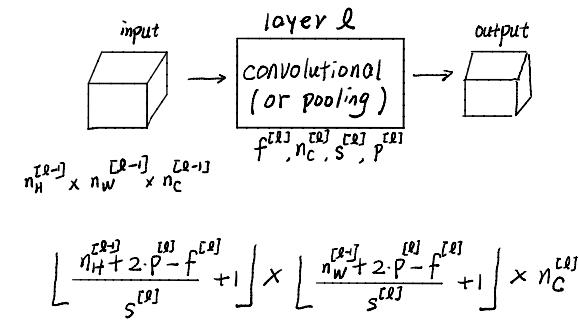
INPUT (39x39x3)
CONV1 ( $f=3, n_c = 10, s=1, p=0$ )
CONV2 ( $f=5, n_c = 20, s=2, p=0$ )
CONV3 ( $f=5, n_c = 40, s=2, p=0$ )
softmax (10,1)



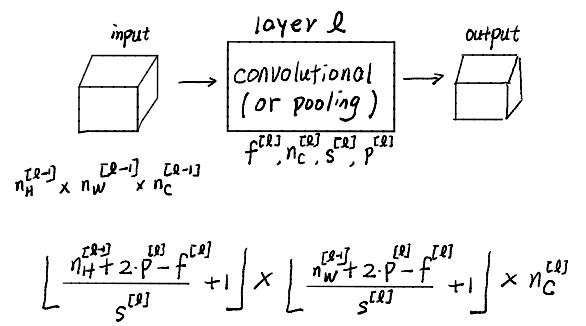
# Example CNN network - 2 (LeNet)

By default,  $P = 0$ .

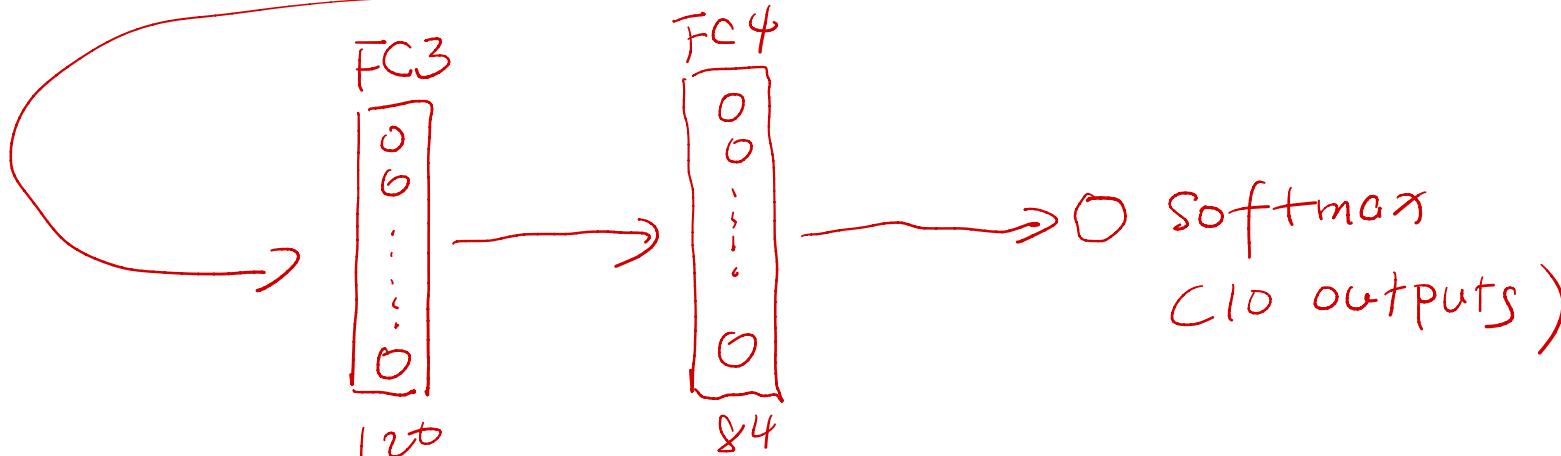
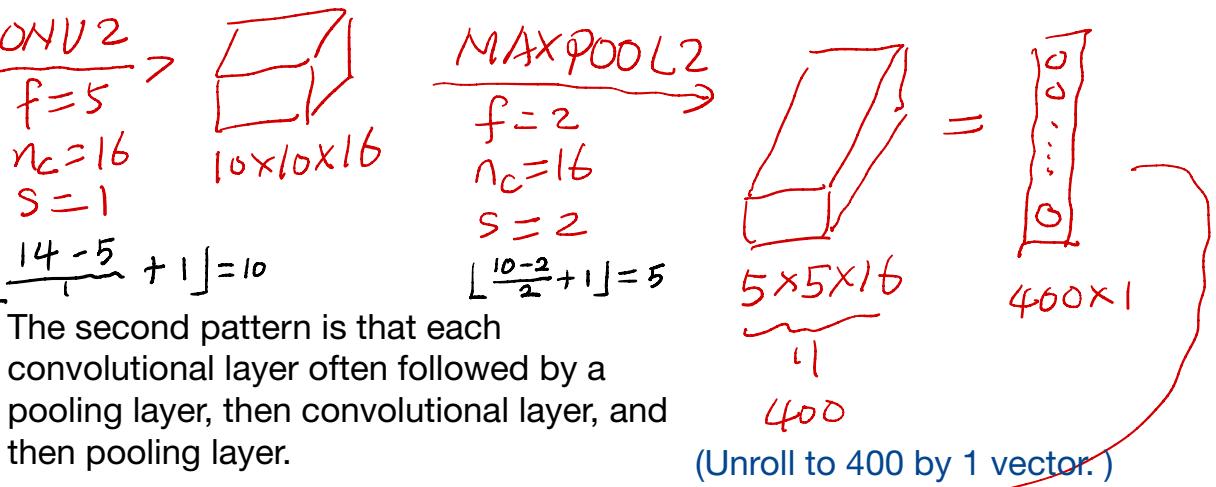
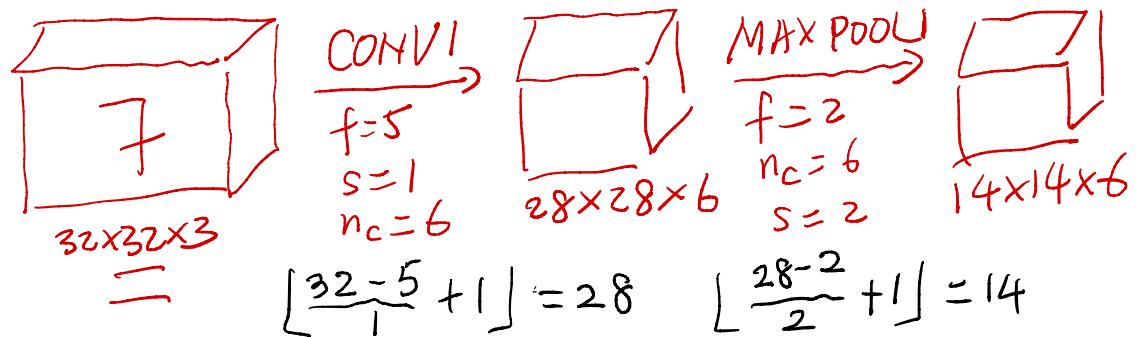




# Example CNN network - 2 (LeNet)



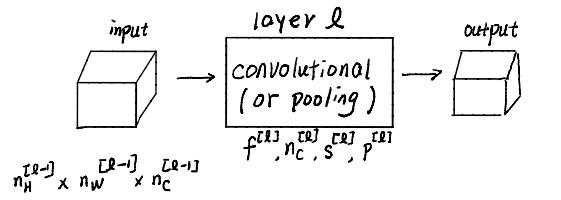
# Example CNN network - 2 (LeNet)



The last pattern is that the last several layers are always fully connected layers with number of neurons decreasing from 120, to 84, and finally 10 outputs for ten class classification in this example network.

# Example CNN network – 3 (AlexNet)

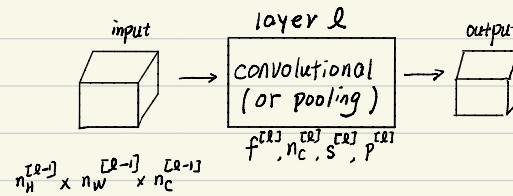
INPUT (227x227x3)
CONV1 ( $f=11, n_c = 96, s=4$ )
MAX-POOL1 ( $f=3, s=2$ )
CONV2 ( $f=5, n_c = 256, s=1$ ), same
MAX-POOL2 ( $f=3, s=2$ )
CONV3 ( $f=3, n_c = 384, s=1$ ), same
CONV4 ( $f=3, n_c = 384, s=1$ ), same
CONV5 ( $f=3, n_c = 256, s=1$ ), same
MAX-POOL3 ( $f=3, s=2$ )
FC1 (4096)
FC2 (4096)
Softmax (1000)



$$P = \left\lfloor \frac{f-1}{2} \right\rfloor$$

$$\left\lfloor \frac{n_h^{l-1} + 2 \cdot P - f^{l-1}}{s^{l-1}} + 1 \right\rfloor \times \left\lfloor \frac{n_w^{l-1} + 2 \cdot P - f^{l-1}}{s^{l-1}} + 1 \right\rfloor \times n_c^l$$

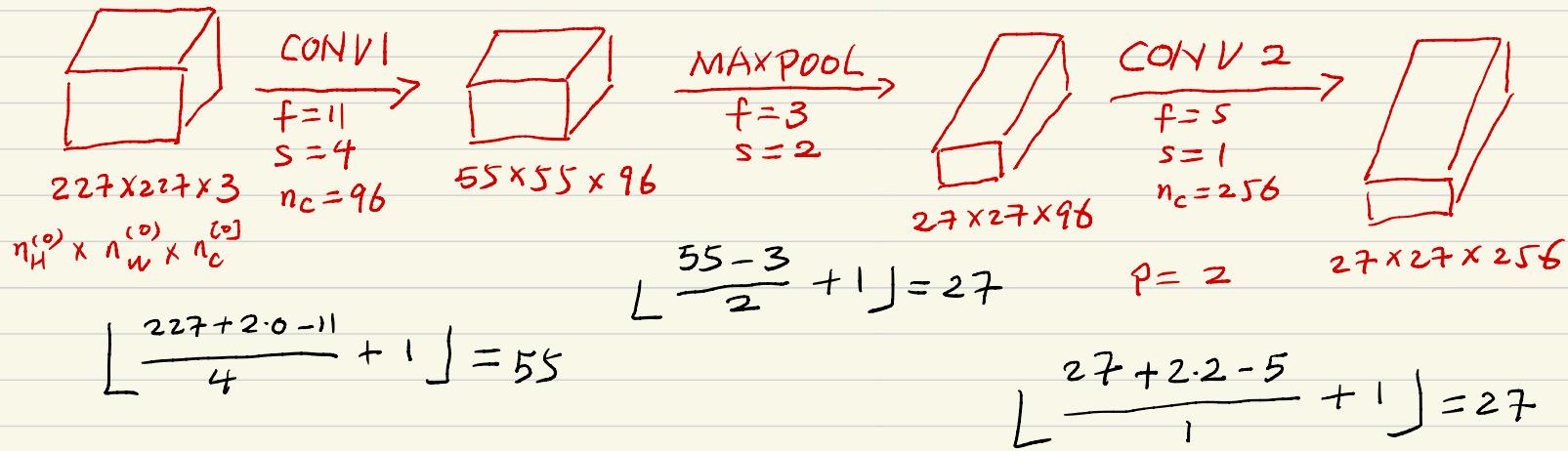
"same" means that the input and output have the same width and height.



$$P = \left\lfloor \frac{f-1}{2} \right\rfloor$$

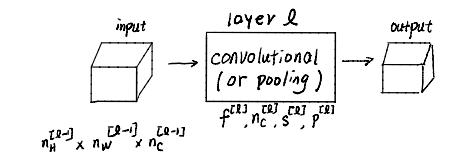
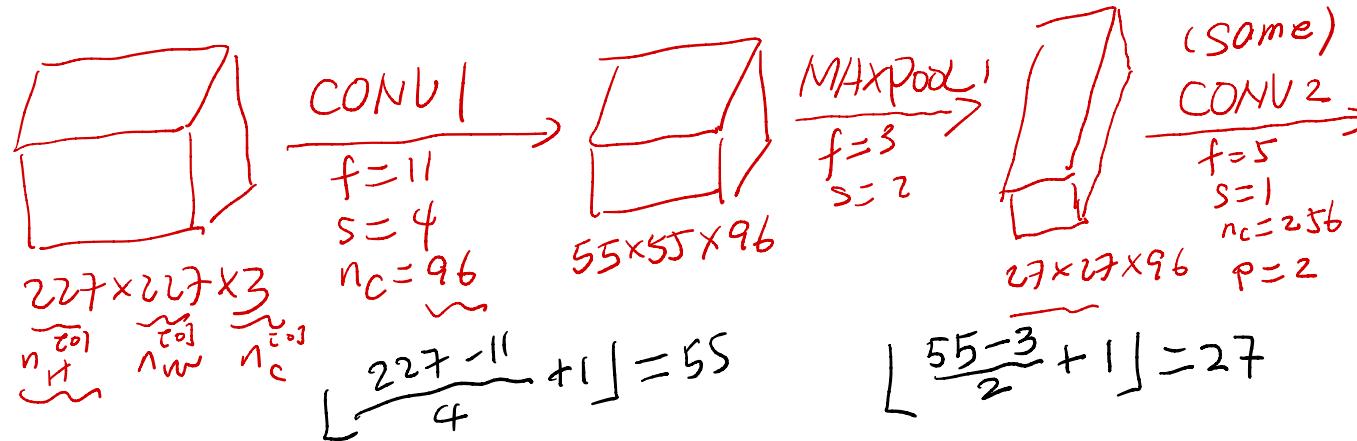
$$\left\lfloor \frac{n_H^{l-1} + 2 \cdot P - f^{l,ij}}{s^{l,ij}} + 1 \right\rfloor \times \left\lfloor \frac{n_W^{l-1} + 2 \cdot P - f^{l,ij}}{s^{l,ij}} + 1 \right\rfloor \times n_C^{l,ij}$$

INPUT (227x227x3)	
CONV1 ( $f=11, n_c = 96, s=4$ )	$P = \left\lfloor \frac{f-1}{2} \right\rfloor$
MAX-POOL1 ( $f=3, s=2$ )	
CONV2 ( $f=5, n_c = 256, s=1$ , same)	
MAX-POOL2 ( $f=3, s=2$ )	
CONV3 ( $f=3, n_c = 384, s=1$ , same)	
CONV4 ( $f=3, n_c = 384, s=1$ , same)	
CONV5 ( $f=3, n_c = 256, s=1$ , same)	
MAX-POOL3 ( $f=3, s=2$ )	
FC1 (4096)	
FC2 (4096)	
Softmax (1000)	

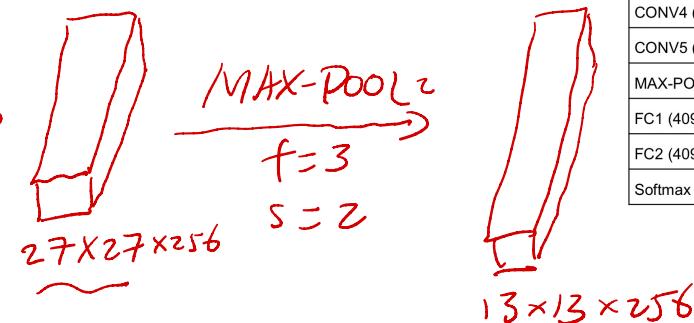


# Alexnet

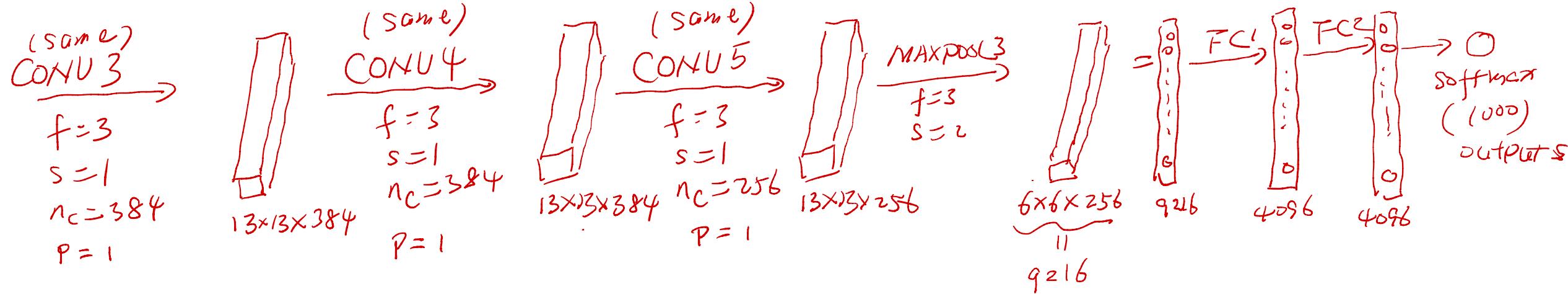
## Example CNN network - 3



$$\left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor \times \left\lfloor \frac{n_W^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor \times n_C^{[l]}$$



INPUT ( $227 \times 227 \times 3$ )
CONV1 ( $f=11, n_c = 96, s=4$ )
MAX-POOL1 ( $f=3, s=2$ )
CONV2 ( $f=5, n_c = 256, s=1$ , same)
MAX-POOL2 ( $f=3, s=2$ )
CONV3 ( $f=3, n_c = 384, s=1$ , same)
CONV4 ( $f=3, n_c = 384, s=1$ , same)
CONV5 ( $f=3, n_c = 256, s=1$ , same)
MAX-POOL3 ( $f=3, s=2$ )
FC1 (4096)
FC2 (4096)
Softmax (1000)



$$\left\lfloor \frac{f-1}{2} \right\rfloor = \left\lfloor \frac{3-1}{2} \right\rfloor = 1$$

# Example CNN network – 4 (VGG-16)

**CONV = 3x3 filter, s = 1, same**

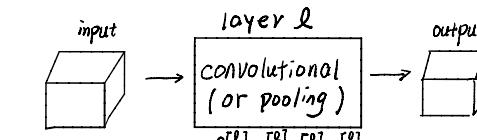
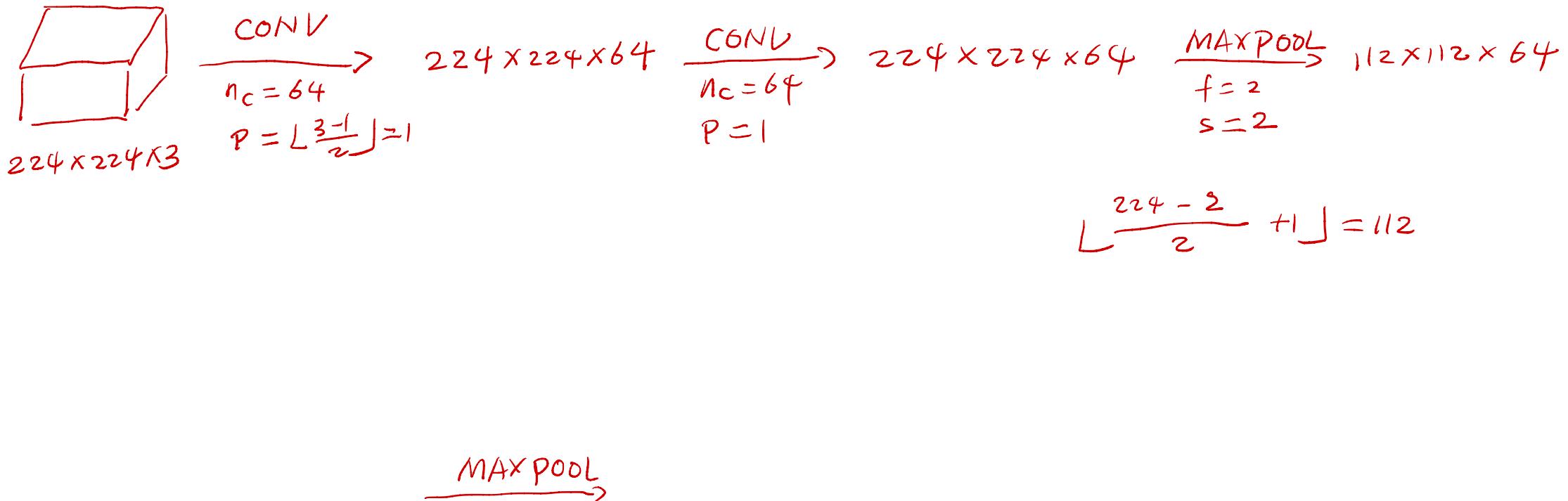
**MAX-POOL = 2x2 filter, s = 2**

Simonyan, Karen, and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition." *arXiv preprint arXiv:1409.1556* (2014). (**44,780** citations).

INPUT (224x224x3)
CONV ( $n_c = 64$ ) X 2
MAX-POOL
CONV ( $n_c = 128$ ) X 2
MAX-POOL
CONV ( $n_c = 256$ ) X 3
MAX-POOL
CONV ( $n_c = 256$ ) X 3
MAX-POOL
CONV ( $n_c = 256$ ) X 3
MAX-POOL
FC (4096)
FC (4096)
Softmax (1000,1)

$\text{CONV} = 3 \times 3$  filter,  $S=1$ ,  $\text{MAXPOOL} = 2 \times 2$ ,  $S=2$

## Example CNN network - 4

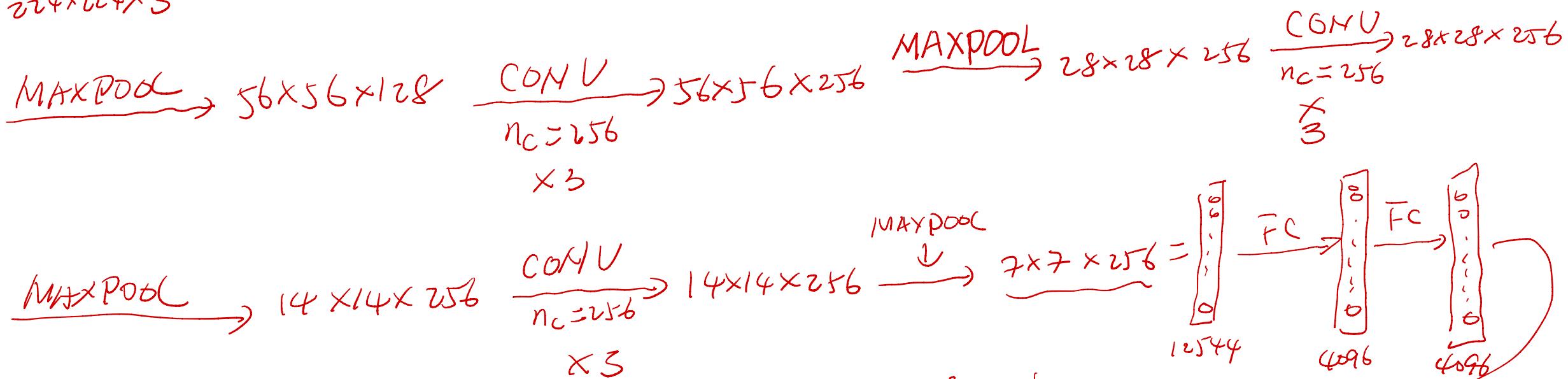
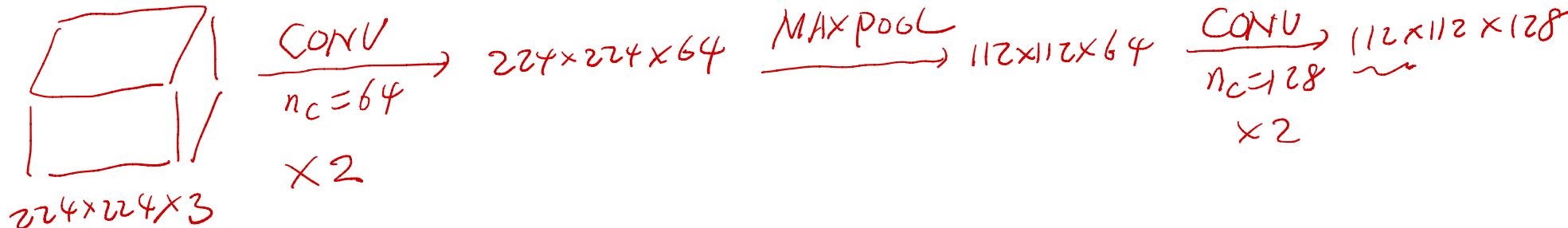


INPUT ( $224 \times 224 \times 3$ )
CONV ( $n_c = 64$ ) X 2
MAX-POOL
CONV ( $n_c = 128$ ) X 2
MAX-POOL
CONV ( $n_c = 256$ ) X 3
MAX-POOL
CONV ( $n_c = 256$ ) X 3
MAX-POOL
CONV ( $n_c = 256$ ) X 3
MAX-POOL
FC (4096)
FC (4096)
Softmax (1000, 1)

UGG-16

CONV = 3x3 filter, S=1, same MAXPOOL = 2x2, S=2

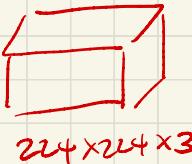
## Example CNN network - 4



INPUT (224x224x3)
CONV ( $n_c = 64$ ) X 2
MAX-POOL
CONV ( $n_c = 128$ ) X 2
MAX-POOL
CONV ( $n_c = 256$ ) X 3
MAX-POOL
CONV ( $n_c = 256$ ) X 3
MAX-POOL
CONV ( $n_c = 256$ ) X 3
MAX-POOL
FC (4096)
FC (4096)
Softmax (1000,1)

# Neural network example

	Activation shape	Activation Size	# parameters
Input:	(32,32,3)	3,072	0
CONV1 (f=5, s=1)	(28,28,8)	6,272	<del>208</del> 608
POOL1	(14,14,8)	1,568	0
CONV2 (f=5, s=1)	(10,10,16)	1,600	<del>416</del> 3216
POOL2	(5,5,16)	400	0
FC3	(120,1)	120	<del>48,001</del> 48120
FC4	(84,1)	84	<del>10,081</del> 10164
Softmax	(10,1)	10	<del>841</del> 850



CONV1

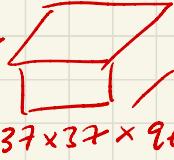
$$\begin{array}{l} f=4 \\ s=2 \\ n_c=96 \end{array}$$



$$110 \times 110 \times 96$$

MAX-POOL1

$$\begin{array}{l} f=3 \\ s=3 \end{array}$$



$$37 \times 37 \times 96$$

CONV2

$$\begin{array}{l} f=3 \\ s=2 \\ n_c=384 \end{array}$$



$$18 \times 18 \times 384$$

$$\left\lfloor \frac{224+0-4}{2} + 1 \right\rfloor = 111$$

$$\left\lfloor \frac{111+0-3}{3} + 1 \right\rfloor = \left\lfloor \frac{108}{3} + 1 \right\rfloor = 17 + 1$$

$$\left\lfloor \frac{37+0-3}{2} + 1 \right\rfloor = 18$$

$$= 36 + 1 = 37$$

MAX-POOL  
 $f=3, s=1$

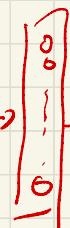


CONV3

$$\begin{array}{l} f=6 \\ s=2 \\ n_c=256 \end{array}$$



$$8 \times 8 \times 256$$

 $=$ 


$\text{softmax}$   
(class output)

$$\left\lfloor \frac{18+0-3}{1} + 1 \right\rfloor = 16 \quad p=2 \quad \left\lfloor \frac{16-6+2 \times 2}{2} + 1 \right\rfloor = 8$$