

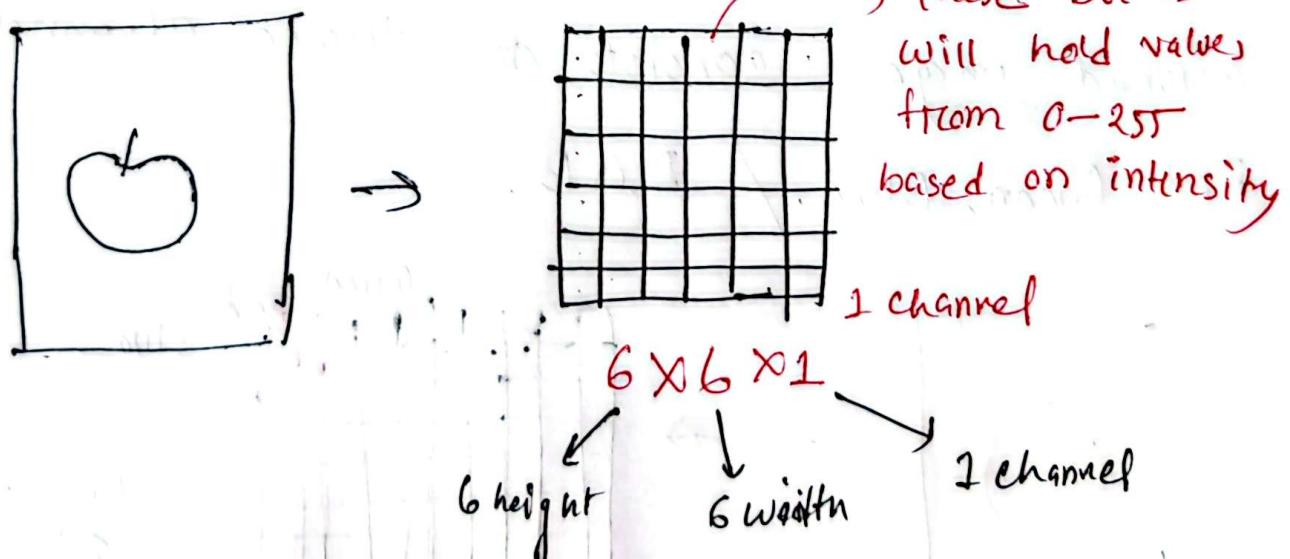
CNN

Convolutional Neural Network handles Image type datasets. Usual Artificial Neural Network may get complex due to the large number of features an image can have. So, before giving an image input to an ANN, we apply convolution layer on the input images to extract important patterns from the pictures.

Image:

Grey-scale Image →

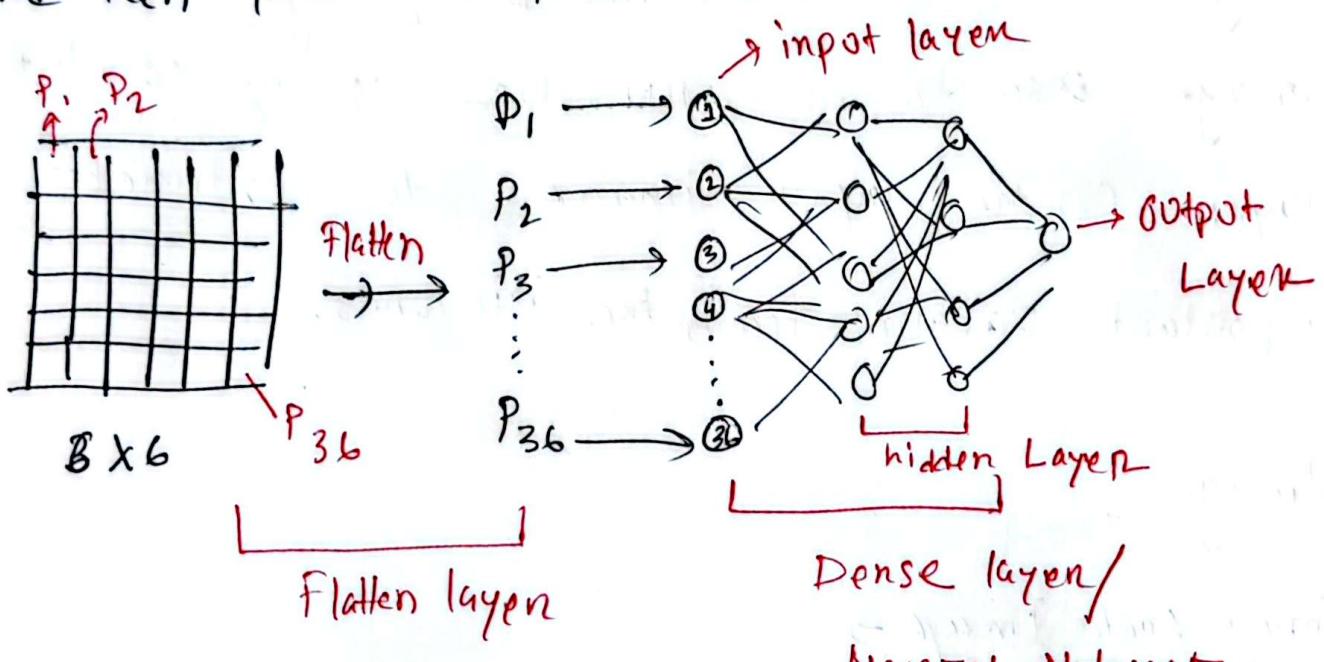
Black and white image has only one channel →



This image has 6×6 pixels = 36 pixels.

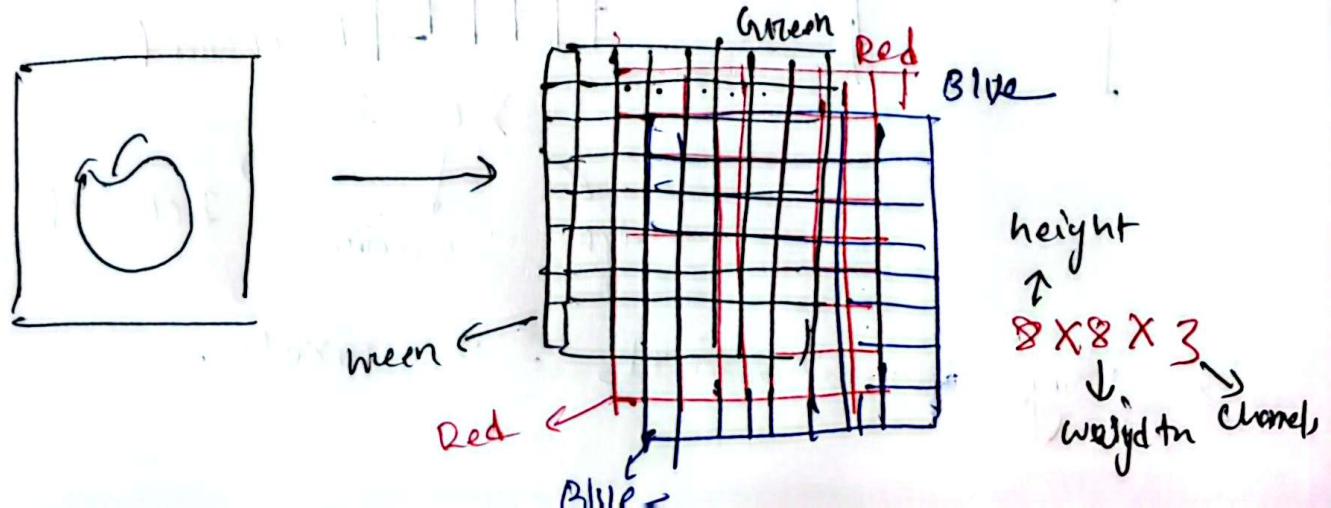
If we want to input this image inside an ANN architecture, the Neural Network will have 36 neurons in the input layer.

so, we will Flatten the image from 2D to 1D and then pass the pixel values inside NN.



Colored Image →

Colored images consists of three channels → Red, Green, Blue. / RGB.



In case of three channel images, the number of pixel is usually same for each channel.

Each channel has $= 8 \times 8 = 64$ pixels.

There are total 3 channels,

so, total pixel $64 \times 3 = 192$ pixels.

For simplicity we have taken example of small dimension image. In real the dimension can be more and which will increase the number of pixels/features for the input layer. To handle this problem we add convolution to our image before training with Neural Network/dense layer.

Filters

We use filters to perform convolution on image. Filters are used to extract pattern from images.

For example →

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

This filter is used to extract horizontal edges.

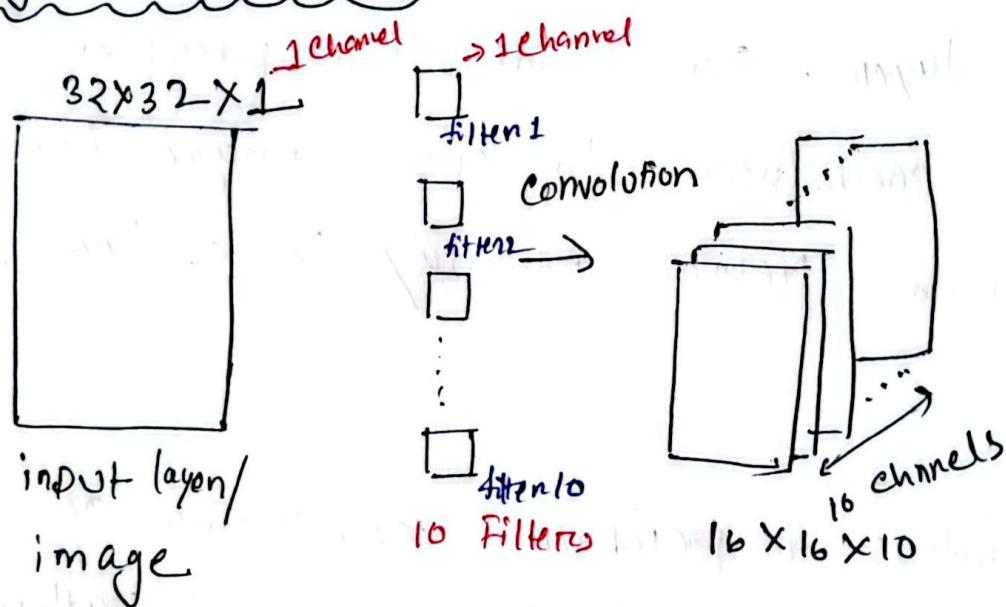
1	6	-1
1	0	-1
1	0	-1

This filter is used to extract vertical edges from an image.

Check the "Tea cup" image from the slide to understand it better.

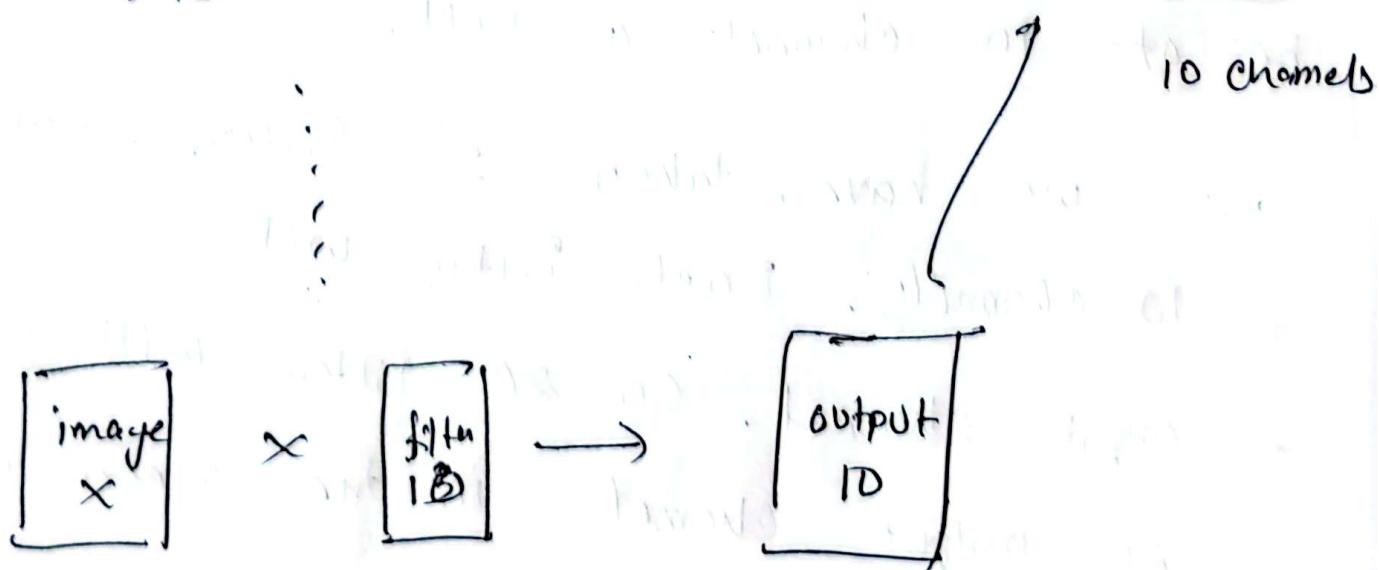
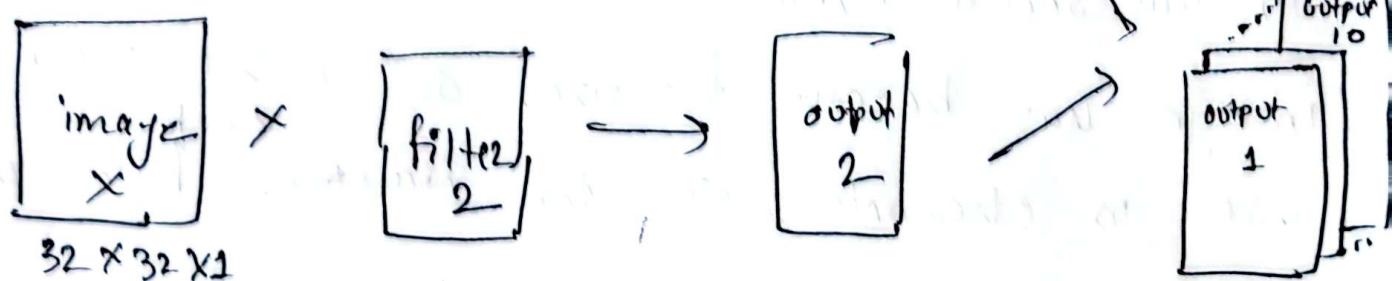
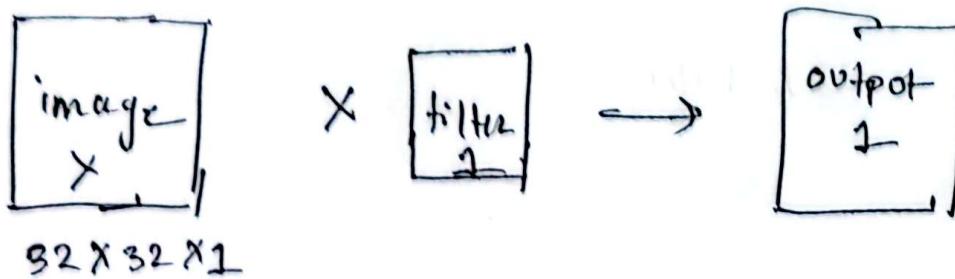
Convolution layers:

Single channel image \rightarrow

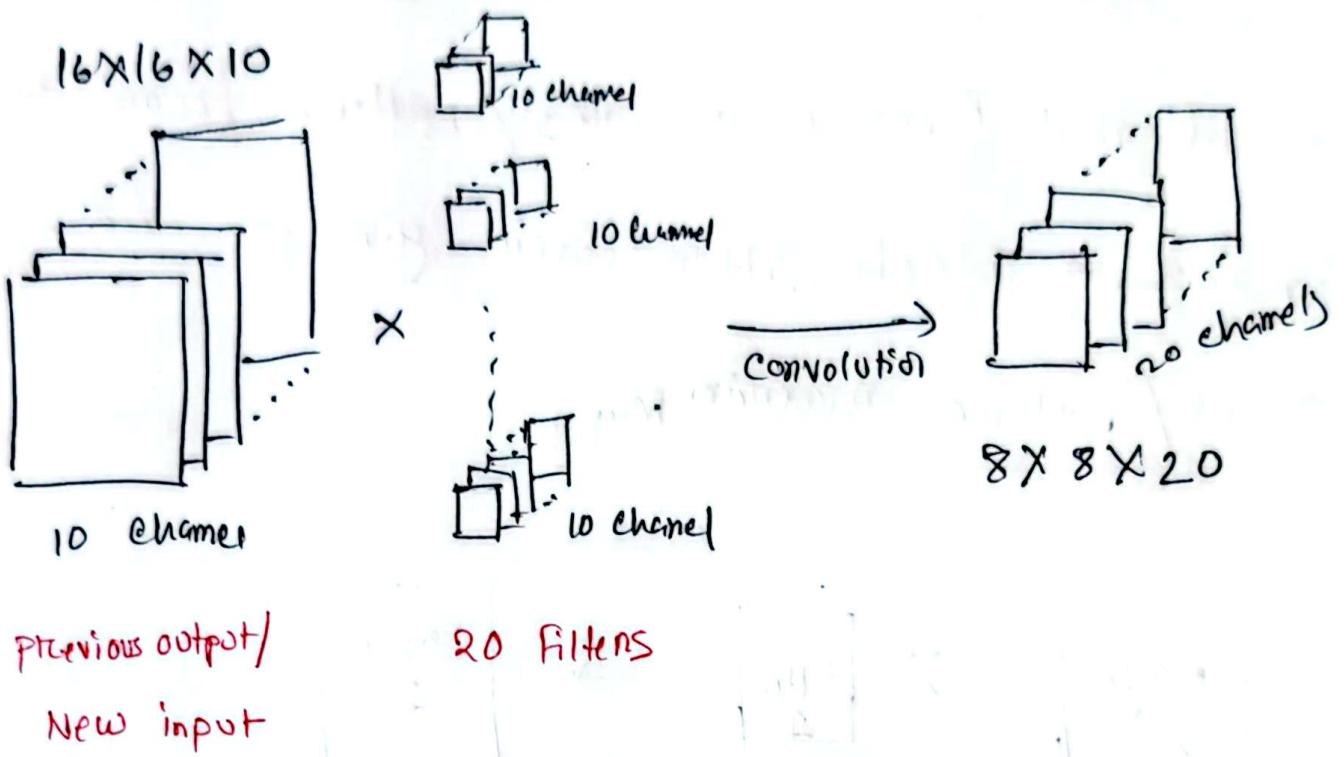


hence the image consists of 1 channel. So the filter will be of 1 channel as well. But we can use multiple feature to extract multiple different types of pattern (horizontal edge, vertical

edge etc). Hence 10 filter will extract 10 different types of features / pattern from the image. \Rightarrow Each filter will give us one output / feature matrix map



Now on this 10 channel output, we can perform another convolution operation.



In this second layer of convolution, the input image has become the size of $16 \times 16 \times 10$, with 10 channels. So, ~~the~~ ^{each} filter will be of 10 channels as well.

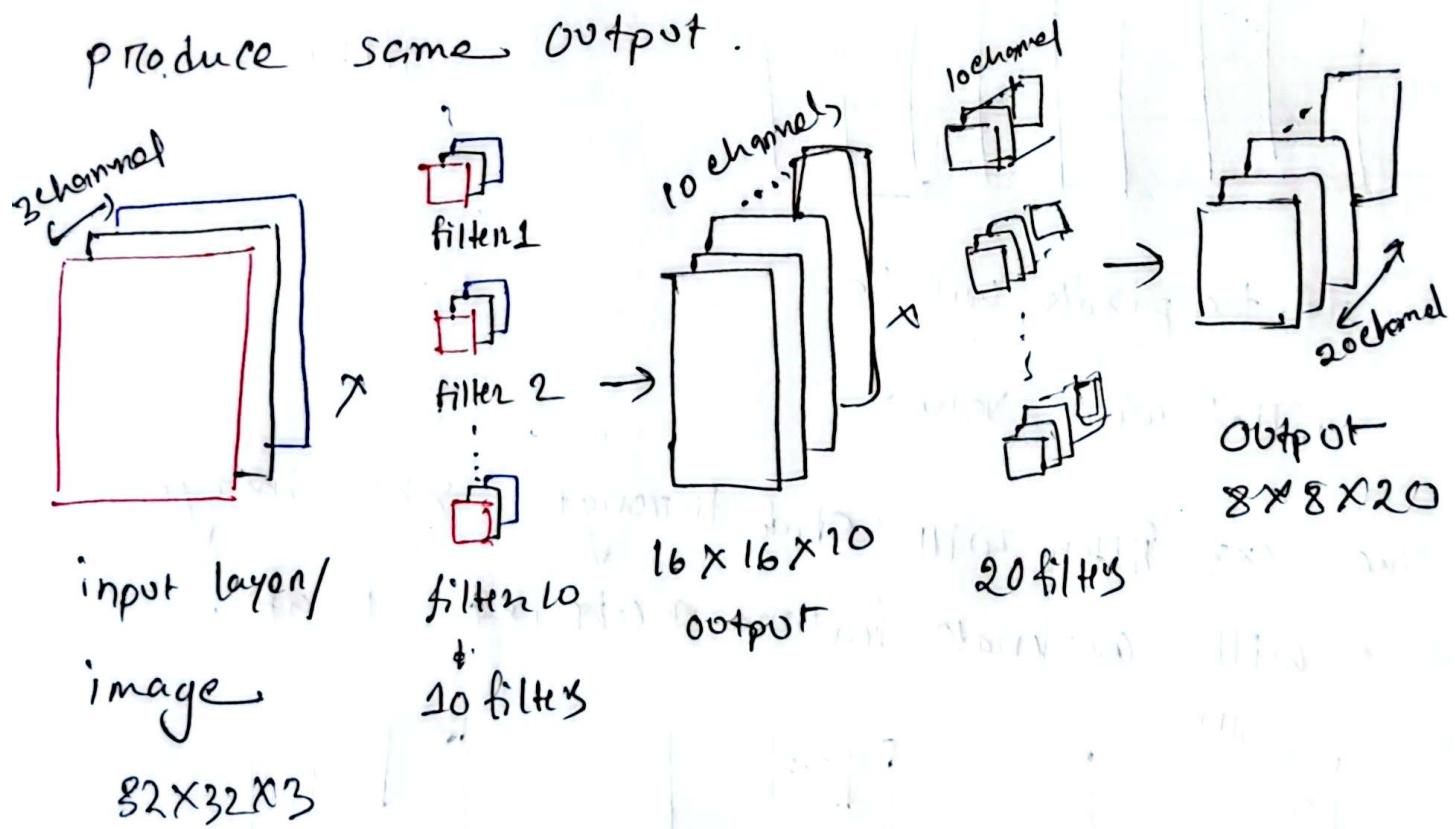
Hence we have taken 20 filters, each of 10 channels. Each filter will give us 1 output channel. So 20 filter will give us 20 output channel in the dimension of $8 \times 8 \times 20$.

We can see the feature is reducing as we can pass this $8 \times 8 \times 20$ size output

into a flatten layer and then pass the values inside ANN/Dense layers for prediction.

Colored image:

The operation is same as grey scaled image, but the only difference is the input images consists of three images channels. So, we need 3 channels per filter for the first convolution layer. The other layer will produce same output.



Here ~~one~~ 1 filter produces 1 channel of feature map ~~matrix~~ for the output as well.

Convolution Operation:

Let's see,

1 channel:

Let's consider an image of $8 \times 8 \times 1$ and filter of 3×3 size.

Image

0.6	0.2	0.6	-0.6	0.02	-0.6
0.1	-0.2	-0.3	-0.1	0.2	0.3
-0.5	-0.1	-0.3	0.5	0.1	0.3

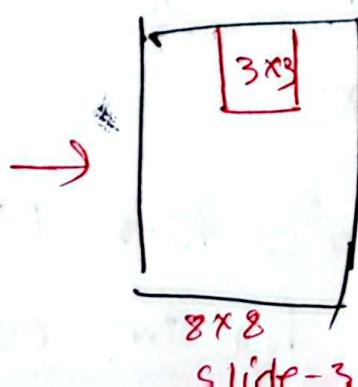
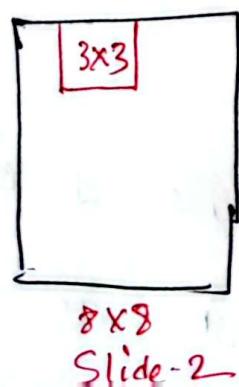
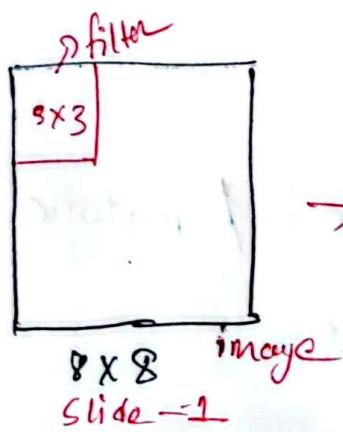
Filter

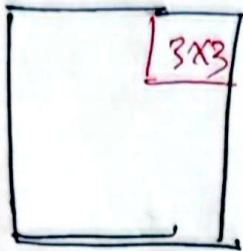
1	1	1
0	0	0
-1	-1	-1

$3 \times 3 \times 1$
filter

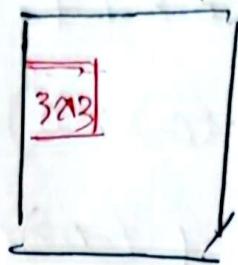
All the pixels will be
filled with values

The 3×3 filter will slide through 8×8 image
and will generate feature map.

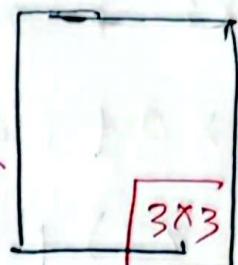




slide-4



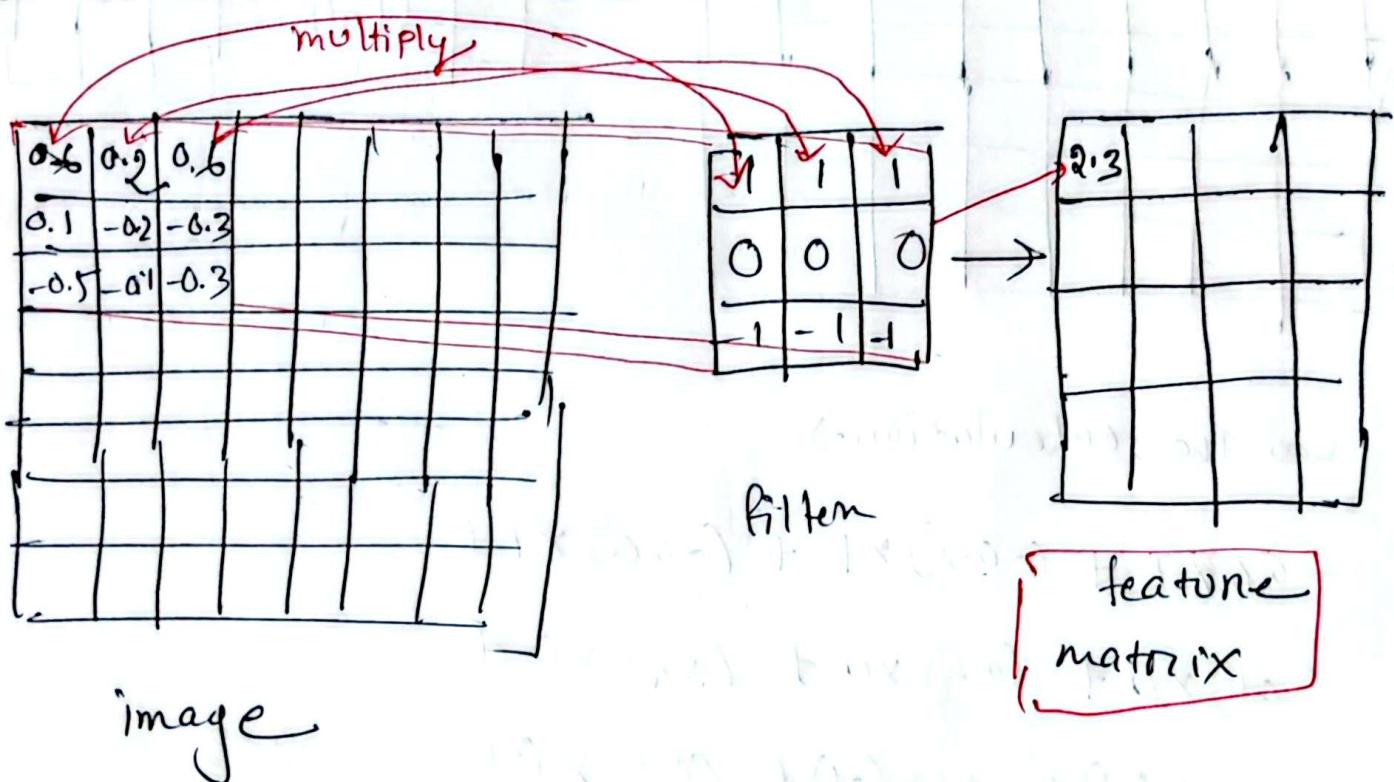
Slide - 5



last slide

we can fix the value of Stride (How many steps the filter ~~passes~~^{skips} to slide through the image) and the feature map will look different based on the value of Stride.

Let us calculate the first slide \rightarrow

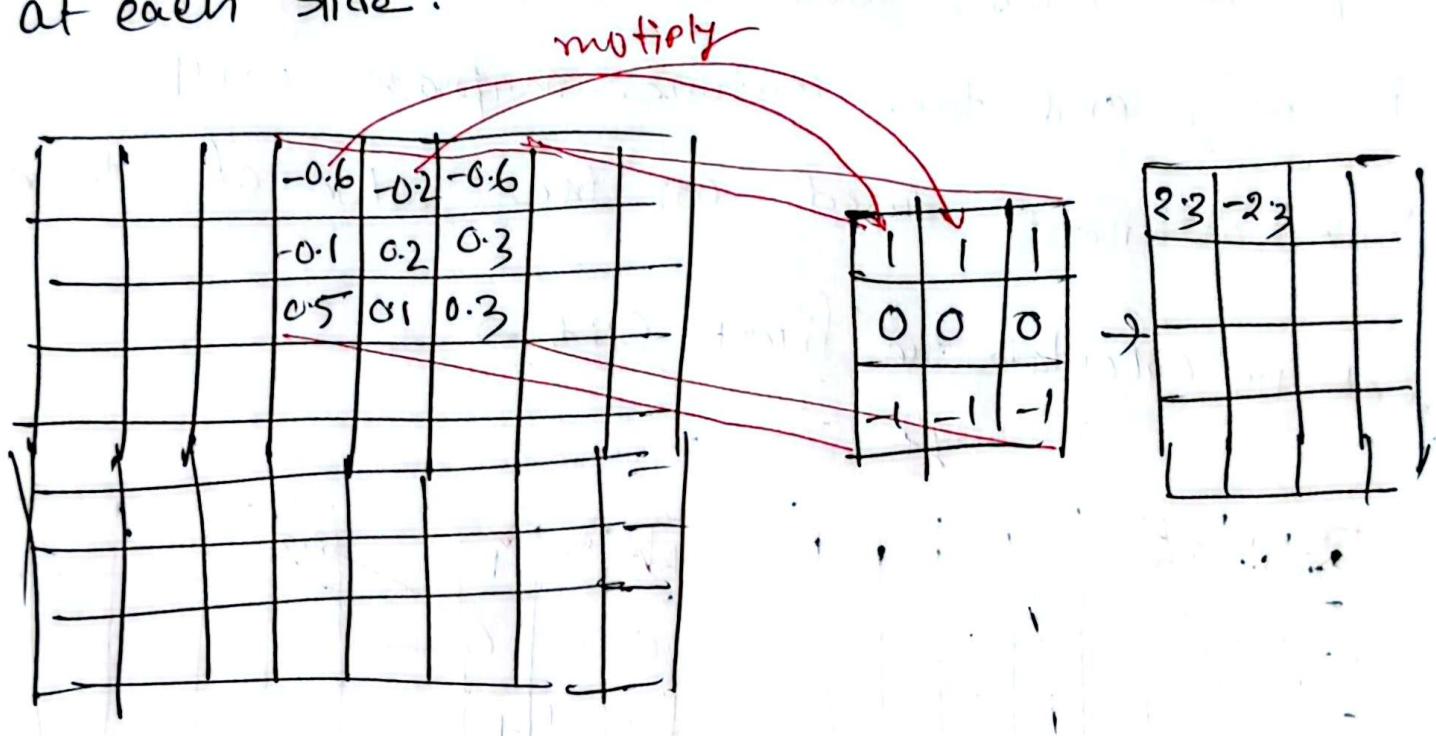


here we multiply each pixel with its corresponding position's filter's value and add all the multiplications

$$\begin{aligned}
 & 0.6 \times 1 + 0.2 \times 1 + 0.6 \times 1 + \\
 & 0.1 \times 0 + (-0.2 \times 0) + (-0.3 \times 0) + \\
 & -0.5 \times (-1) + (-0.1) \times (-1) + (-0.3) \times (-1) = 2.3
 \end{aligned}$$

second slide!

Let us consider we jump/skip three pixels at each slide.



see the calculation \rightarrow

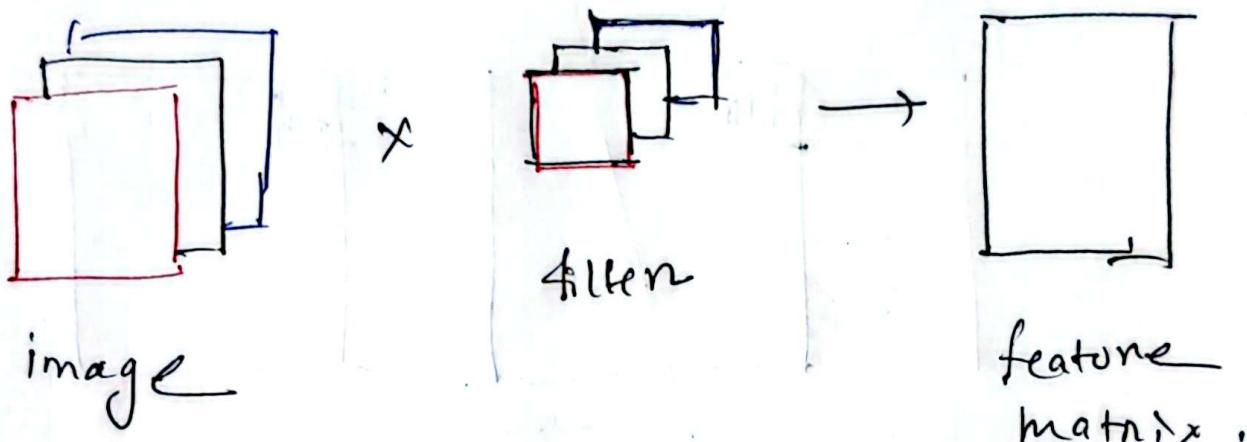
$$\begin{aligned}
 & -0.6 \times 1 + (-0.2) \times 1 + (-0.6) \times 1 + \\
 & (-0.1) \times 0 + (0.2) \times 0 + (0.3) \times 0 + \\
 & 0.5 \times (-1) + 0.1 \times (-1) + 0.3 \times (-1) = -2.3
 \end{aligned}$$

this how we slide through the entire image and generate the feature maps for one filter.

Three channel!

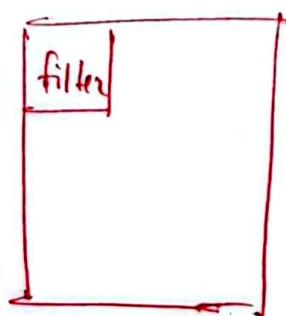
We know for colored image (3 channel)

We use three channel filter.

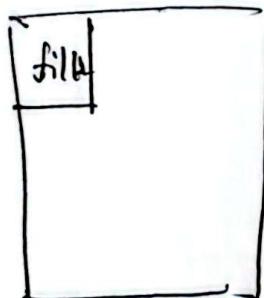


Hence for 1 filter we will get 1 feature matrix as well.

Let us see the first slide for colored image →



Red
channel



Green
channel



Blue
channel

Each channel ^{in image} has its own channel in the filter. We will multiply the value like below

and sum all the multiplication up to generate ONE number for the first slide.

2x2

R1	R2	R3	
R4	R5	R6	
R7	R8	R9	
.	.	.	
			11111111

Red channel
X image

G1	G2	G3	
G4	G5	G6	
G7	G8	G9	
.	.	.	
			11111111

Green channel
Image

B1	B2	B3	
B4	B5	B6	
B7	B8	B9	
.	.	.	
			11111111

Blue channel
Image

F1	F2	F3
F4	F5	F6
F7	F8	F9

Red channel
filter

Convolution ↓

$$\begin{aligned} & R_1 F_1 + R_2 F_2 + R_3 F_3 \\ & + R_4 F_4 + R_5 F_5 + R_6 F_6 \\ & + R_7 F_7 + R_8 F_8 + R_9 F_9 \end{aligned}$$

F1	F2	F3
F4	F5	F6
F7	F8	F9

green
channel
filter

$$G_1 F_1 + G_2 F_2 +$$

$$+ \dots + G_9 F_9$$

F1	F2	F3
F4	F5	F6
F7	F8	F9

Blue channel
filter



$$B_1 F_1 + B_2 F_2 +$$

$$+ \dots + B_9 F_9$$

Let us assume the summation = 120.

so the feature map matrix will look like →

120		

Second Slide →

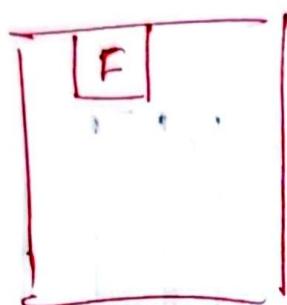


image
red channel

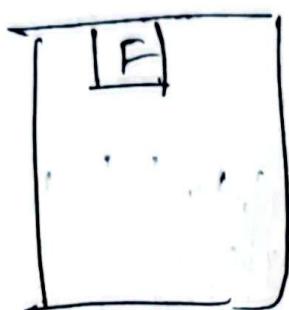


image
green
channel

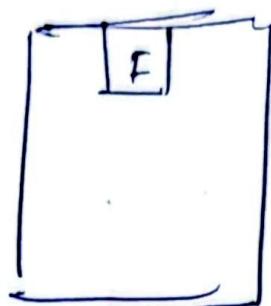
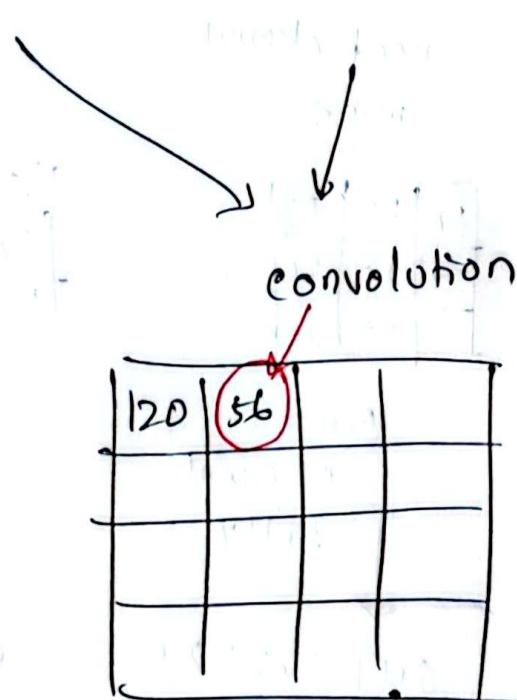


image Blue
channel



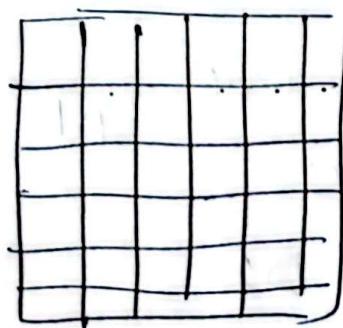
Now we need to know how to calculate the dimension of this feature map. Before that we will learn about padding layer.

convolution operation

Padding Layer:

In general the feature map generated from the feature convolution operation has size less than the original input. In order to keep the size equal, the input image/feature map is padded with zeros around the input/feature map. It is called zero padding.

Padding.



Actual Image/
Feature map

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

1 zero padded
layer.

We can add multiple layer of padding. Sometimes the outer most row/column of the Actual image / feature map gets copied into the padding for more

accuracy. The padding layer can be added before convolution, but it is not mandatory.

Calculating the dimension of feature map:

The generated feature map's dimension can vary based on the input image size, filter size, stride and padding.

$$\text{Output dimension (width)} = \boxed{\frac{W - F + 2P}{S} + 1}$$

W = width of input image

F = width/^{height} of filter

P = Number of zero padding

S = stride

H = Height of input image

$$\text{Output dimension (Height)} = \frac{H - F + 2P}{S} + 1$$

If input image is of $\rightarrow W \times H \times D_1$

Filter $\rightarrow F \times F \times D_2$

Number of filters $\Rightarrow K$

$$\boxed{\text{The output/feature map dimension} = W_2 \times H_2 \times K}$$

$$\text{where } W_2 = \frac{W_1 - F \times 2P}{S} + 1 \text{ and } H_2 = \frac{H_1 - F \times 2P}{S} + 1$$

Convolution operation as Matrix Multiplication

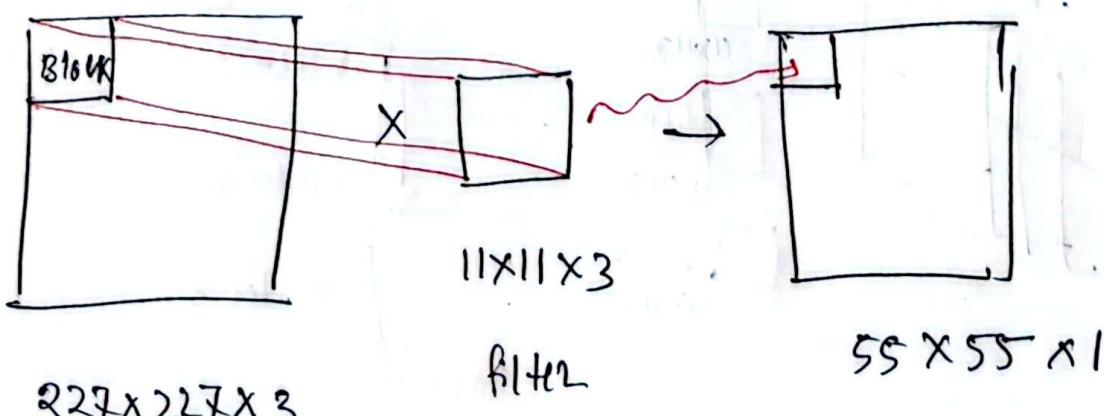
Let us consider an ~~image~~^{filter} of $11 \times 11 \times 3$ and ~~filter of~~ image of $227 \times 227 \times 3$. If

$$P = 0 \text{ and } S = 4,$$

$$\text{the dimension of feature map} = \frac{w-f+2p}{s} + 1 \\ = \frac{227 - 11 + 2 \cdot 0}{4} + 1 \\ = 55$$

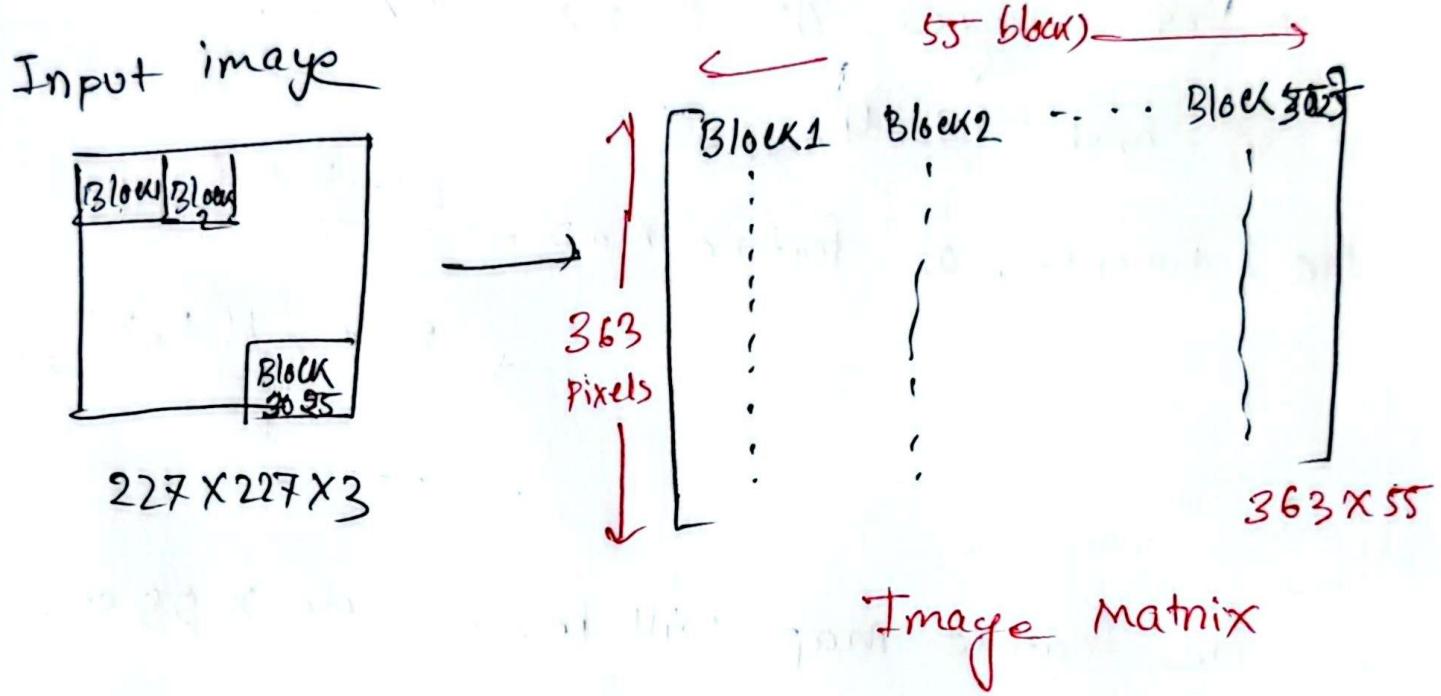
so, the feature map will be of $55 \times 55 \times 1$.

so, we are basically dividing the image into $55 \times 55 = 3025$ blocks using the $11 \times 11 \times 3$ filter.

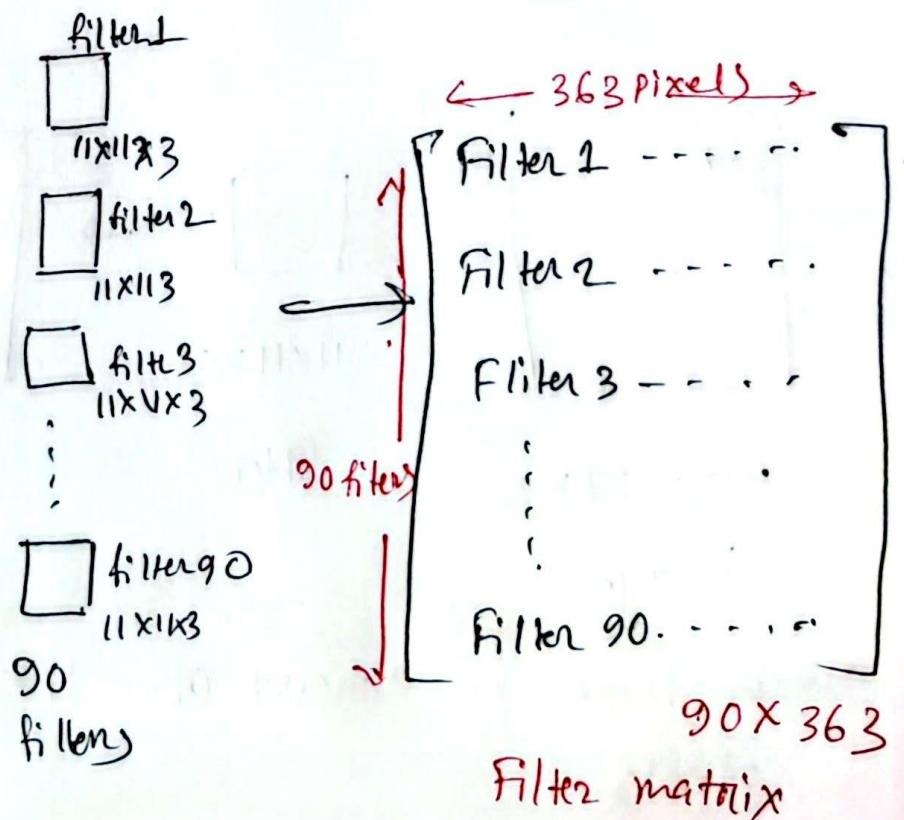


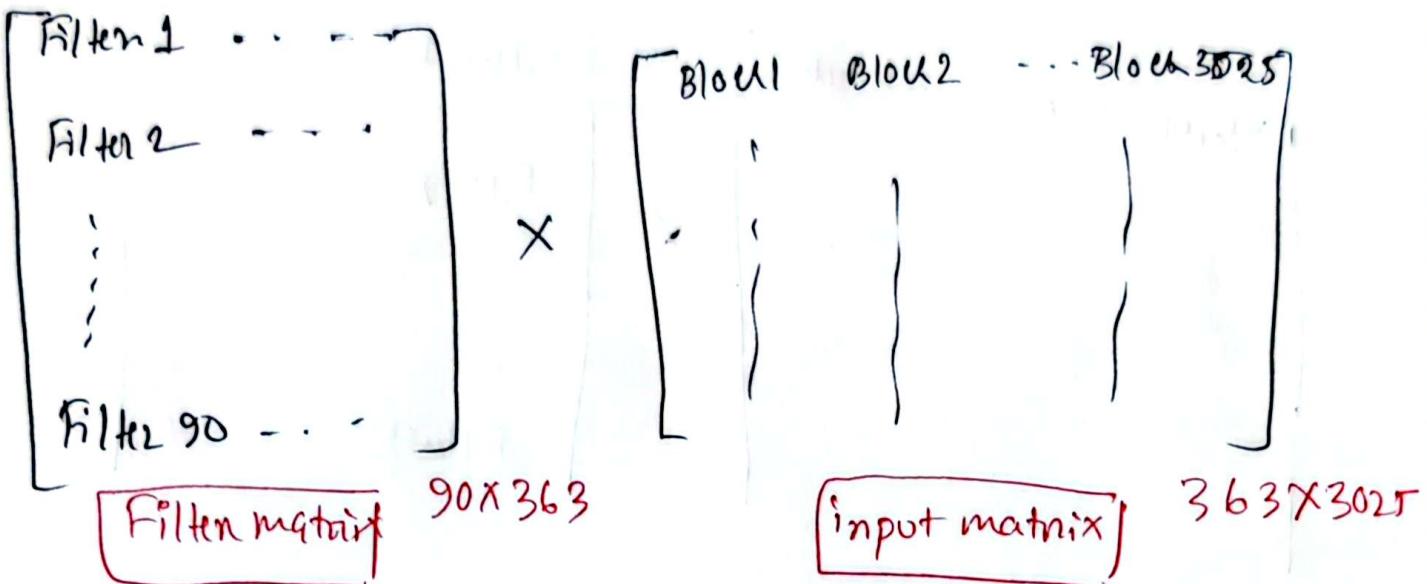
Each block is consist of $11 \times 11 \times 3$ pixels = 363 pixels.

We can represent each image as 3025 blocks where each block has 363 pixels.



If there are total 90 filters used in the convolution layer, the filter matrix will be off





$=$

The result of the multiplication is shown as a vertical column of boxes labeled "Feature map 1", "Feature map 2", and "Feature map 90". A red box highlights this column and its dimensions 90×3025 . To the right of the column, handwritten text states: "each feature map of 3025 pixels". Below the column, a red bracket groups the three feature maps with the label "90 filter generates 90 feature maps".

$$\therefore \text{output feature map matrix} = \text{filter matrix} \times \text{input matrix}$$

$$= (90 \times 363) \times (363 \times 3025)$$

$$= 90 \times 3025$$

$$= 90 \times \underbrace{55 \times 55}_{\text{feature map dimension}}$$

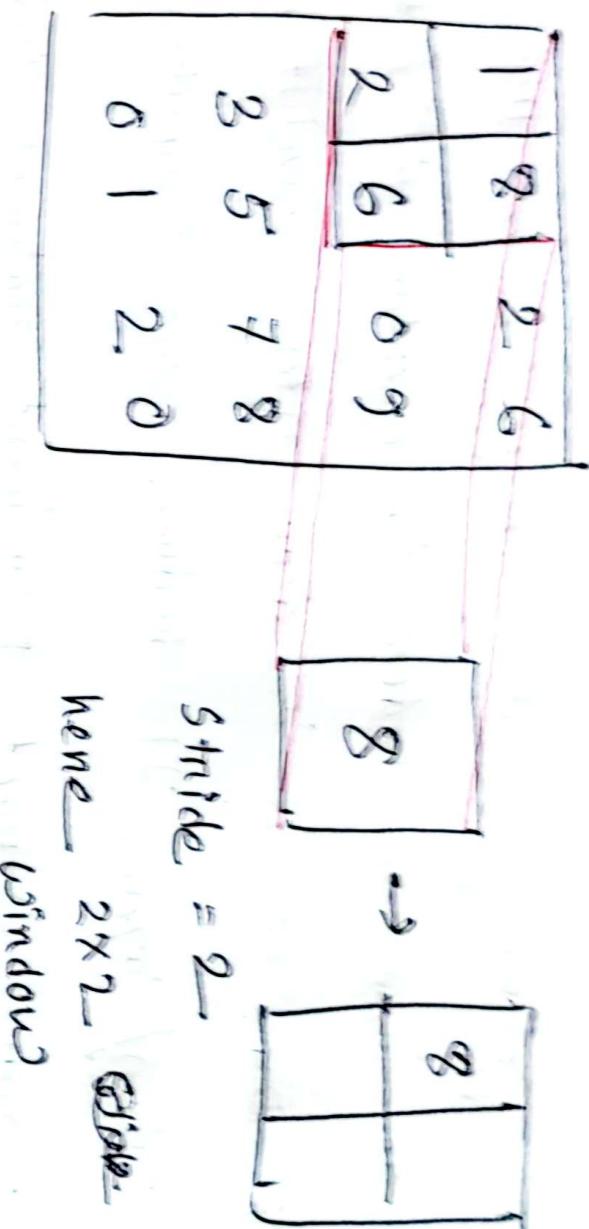
90 channels.

Pooling Layer:

After generating feature maps, we can use pooling layer to reduce the dimension of feature map ~~and~~. So it is done after convolution operation/ layer. It reduces the computational complexity, extract prominent features and reduces dimension.

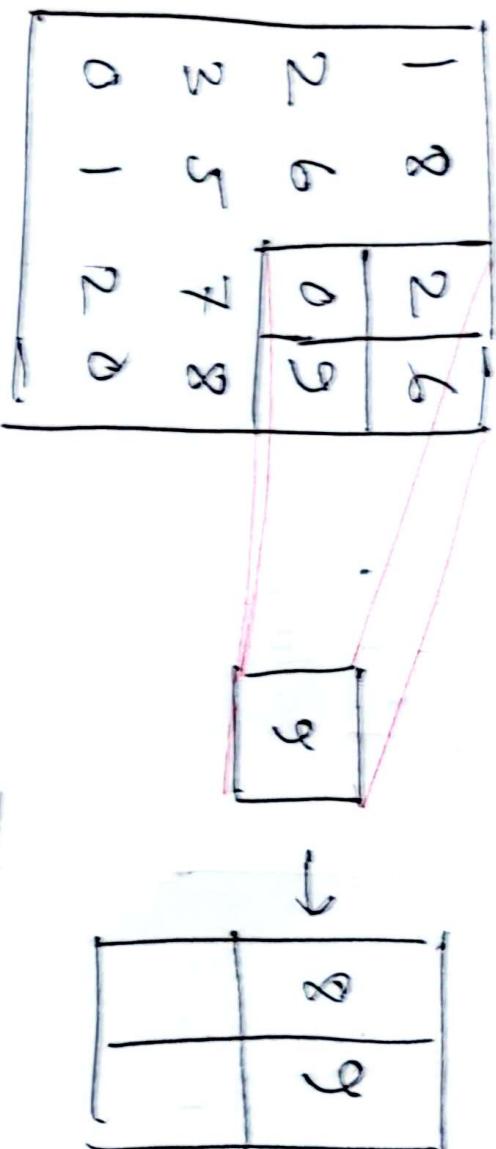
We can do two types of pooling!

Max Pooling!



The 2×2 window extracts the maximum value from the feature matrix. Here among 1, 8, 7, 6, $\rightarrow 8$ is extracted.

If stride = 2 \rightarrow second slide is \rightarrow



Similarly we will get \rightarrow

8	9
5	8

Average pooling:

Instead of pooling the max value, it take the average of all values inside the window.

and and

Normalization:

Normalization scales the input images¹ into a specific range. Please read the slide to understand the types of Normalization.