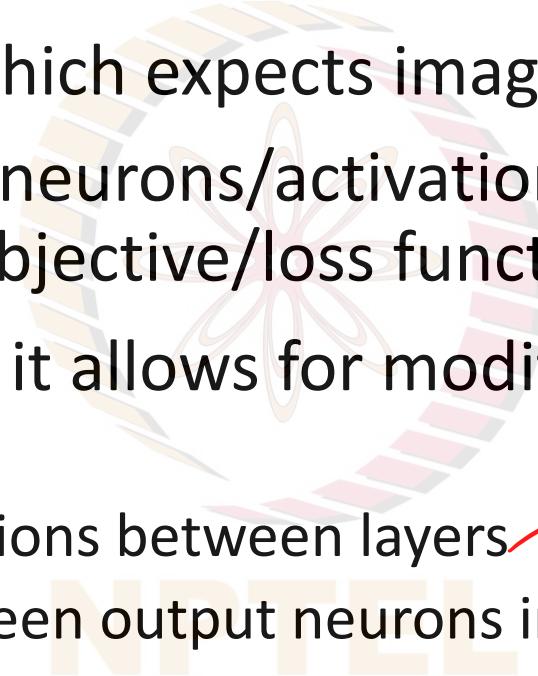

Convolutional Neural Networks



What are CNNs?

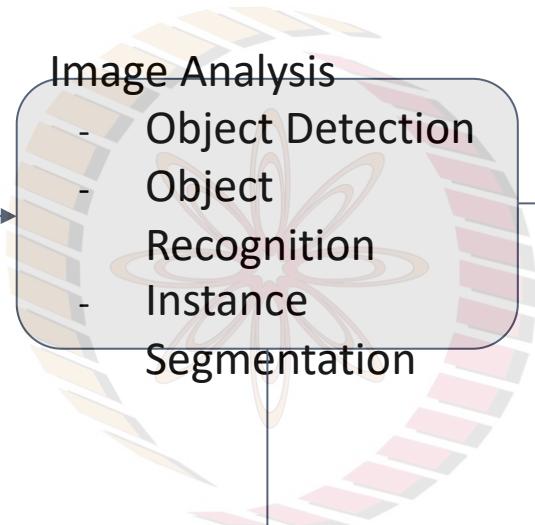
- A special type of ANNs which expects images as inputs
- Like ANNs they too have neurons/activations and weights (estimated during training) and an objective/loss function
- Since the input is images it allows for modifications to the architecture so that
 - There are sparse connections between layers 
 - Weights are shared between output neurons in the hidden layer 

Applications in Computer Vision & Image Analysis

- Image Recognition
- Object Detection and Localization
- Semantic segmentation
- Medical Image Analysis



Applications in Image Analysis

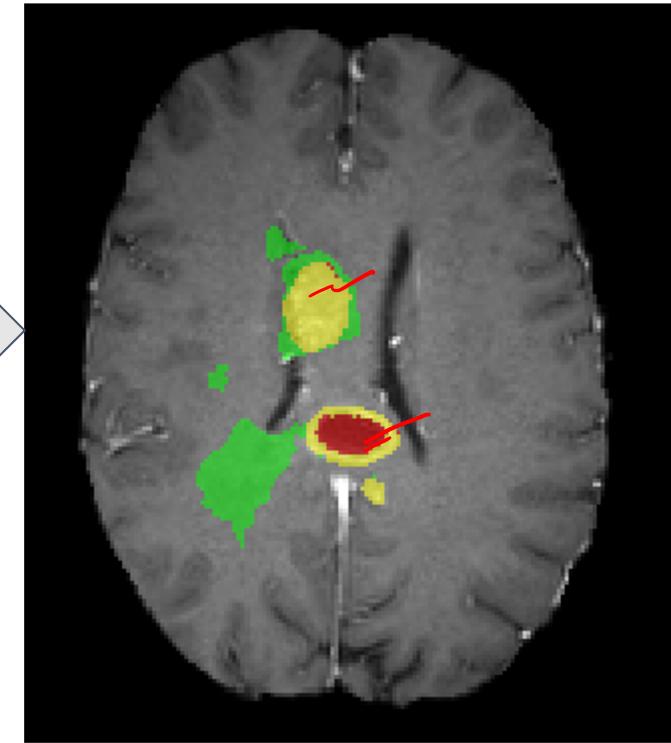
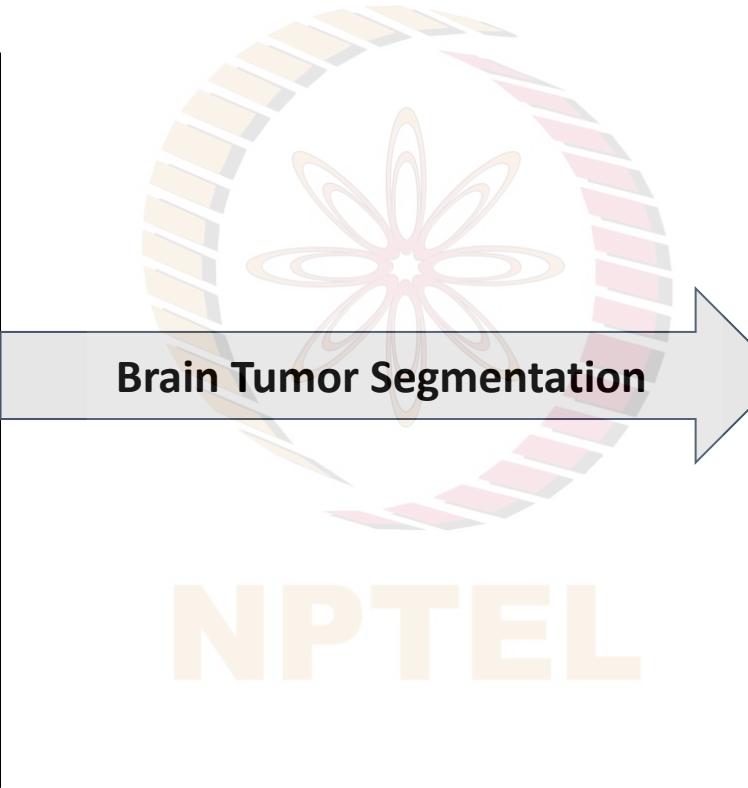
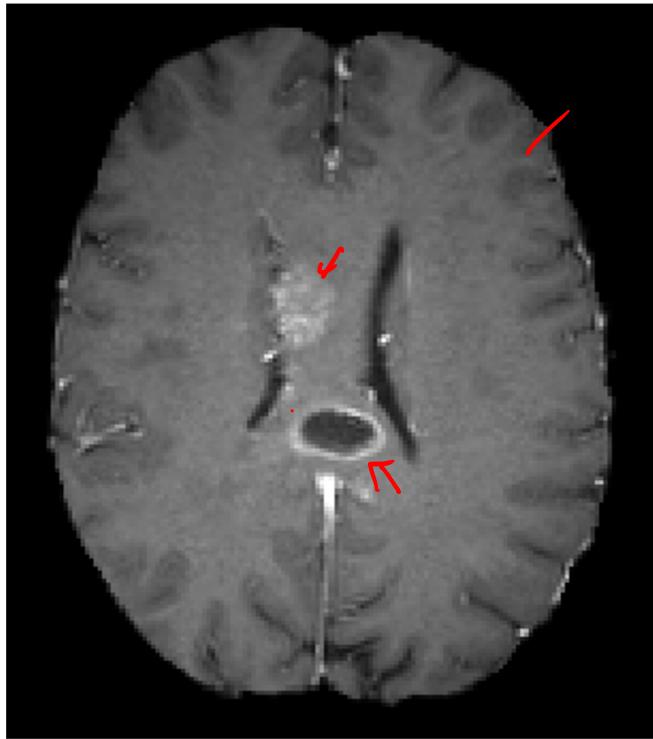


Herd	Goats	Herder
91%	73%	62%
Sheep	Livestoc	Grass
82%	k	59%
Pastur	71%	Farm
e	Herding	50%
80%	69%	1



- 1) <https://cloud.google.com/vision/>
- 2) <http://silverpond.com.au/object-detector.html>

Applications in Image Analysis

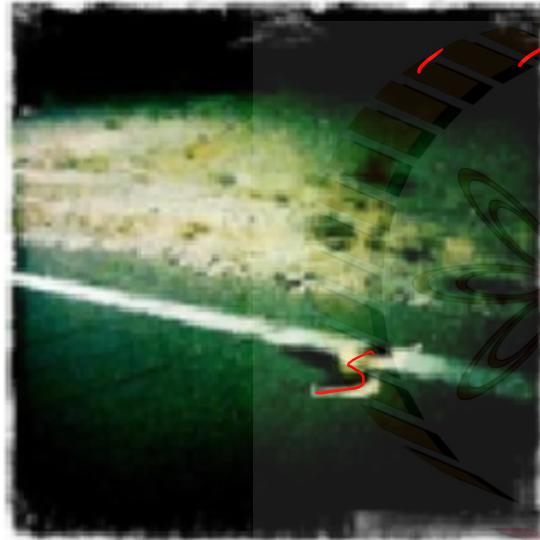


Preprocessed MR Image of
Brain

Tumor Segmentation

ImageNet Challenge

1000 classes, a million training images and report top-5 prediction accuracy. The top-5 error is now better than humans- Note humans are not trained! Human error: ~5.1%



rule, ruler ✓	sidewinder	hatchet✓	schipperke ✓
pencil box, pencil case ✓	maze, labyrinth	vase	schipperke
rubber eraser, rubber ✓	gar, garish	pitcher, ewer	groenendael
ballpoint, ballpoint pen ✓	valley, vale	coffeepot	doormat, welcome mat
pencil sharpener ✓	hammerhead	mask	teddy, teddy bear
carpenter's kit, tool kit	sea snake	cup	jigsaw puzzle

Progress in ImageNet challenge

Model	Top-1 Accuracy	Top-5 Accuracy	Parameters	Depth
Xception	0.79	0.945	22910480	126
VGG16	0.715	0.901	138357544	23
VGG19	0.727	0.91	143667240	26
ResNet50	0.759	0.929	25636712	168
InceptionV3	0.788	0.944	23851784	159
InceptionResNetV2	0.804	0.953	55873736	572
MobileNet	0.665	0.871	4253864	88
DenseNet121	0.745	0.918	8062504	121
DenseNet169	0.759	0.928	14307880	169
DenseNet201	0.77	0.933	20242984	201

The top-1 and top-5 accuracy refers to the model's performance on the ImageNet validation dataset (Source: <https://keras.io/applications/>)

Image Parametrization

- Gray scale images

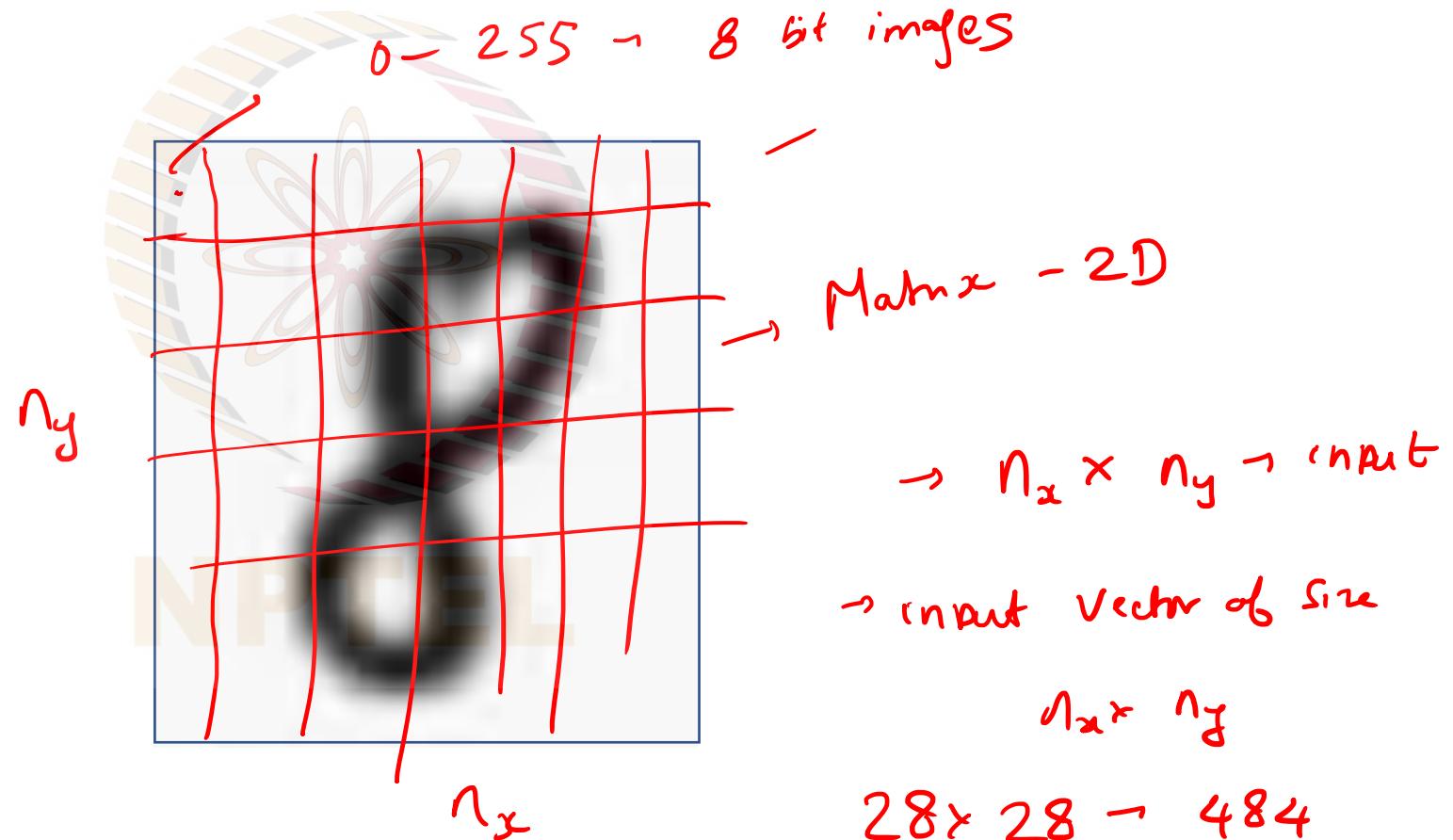
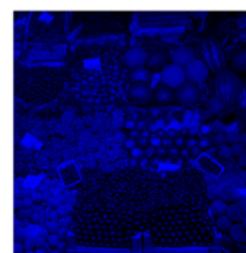
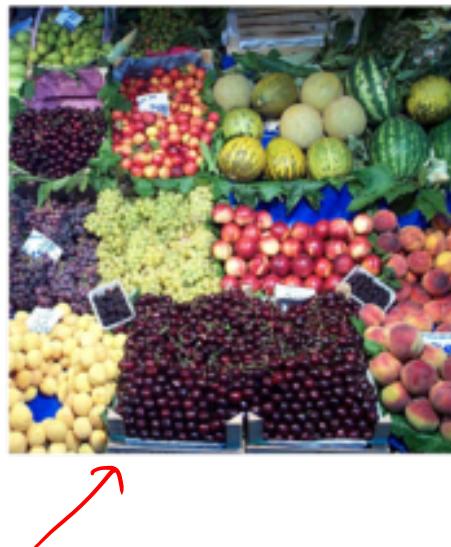


Image Parametrization

- RGB images



Why not regular Neural Networks?

- ANNs take a vector of inputs and product as output another hidden layer vector fully connected to the input-
 - For small image sizes the number of weights/parameters to be estimated are not large – but consider a 224x224x3 image - RGB images.
 - A single neuron in the output layer will have 224x224x3 weights coming into it.
 - A ‘Volume’ image input like RGB images will lead to an explosion in the number of weights – Requires more memory, computations and data
- CNNs exploit the structure of images
 - leading to sparse connections between input and output neurons
 - parameter sharing between output neurons

$$\begin{aligned} 30 \times 30 \times 3 &\rightarrow 2700 \text{ input neurons} \\ [256 \times 256 \times 3] &\rightarrow 10^5 \text{ neurons} \\ &\rightarrow 10^8 \rightarrow \text{Weights} \end{aligned}$$

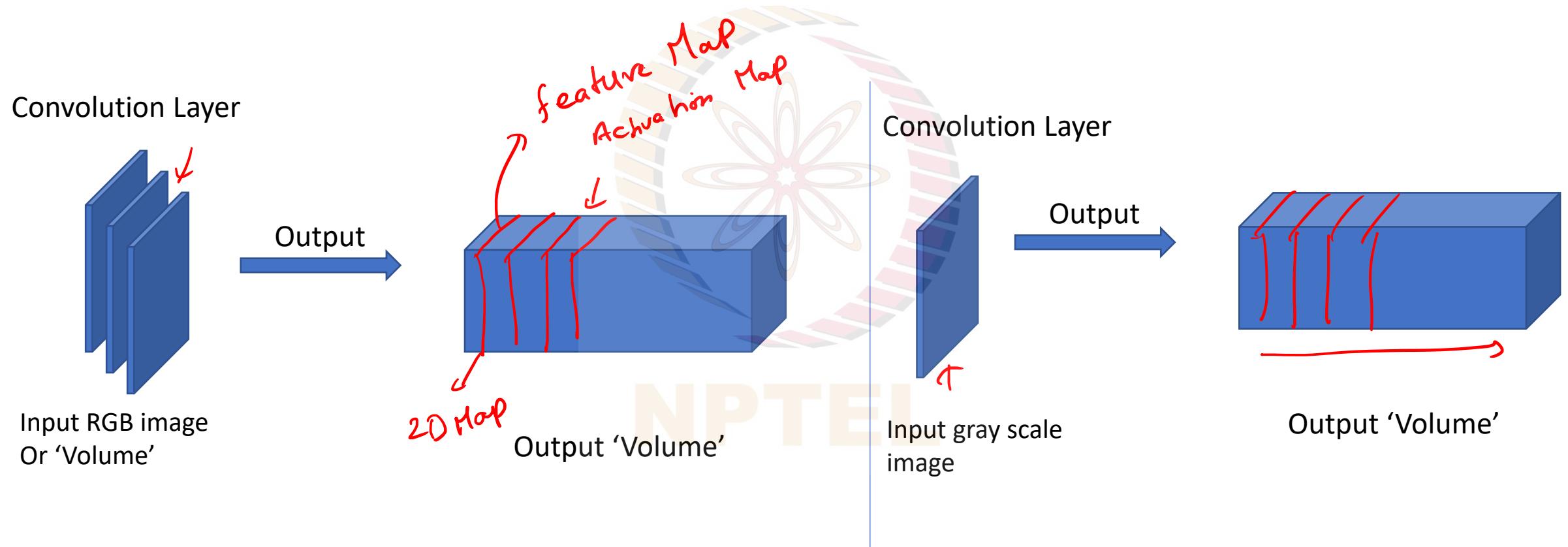
CNN layers and operations

- Like regular ANNs, CNNs stack a sequence of layers followed by an output layer – Classification layer
- Each layer performs a 'Convolution' operation (Conv Layer) or a 'Pooling' operation (Pooling Layer)
- The Conv Layer and Pooling Layers are alternatively stacked – leading to a series of fully connected layers followed by an output layer

CNN Layers & Operations

- CNNs take as input an image ‘Volume’ – For e.g. RGB volume or sub-volume of a 3D medical image- and each layer outputs another 3D volume following a convolution or pooling
- Contrast to ANN where the output from each layer is another vector of Neurons

CNN layers

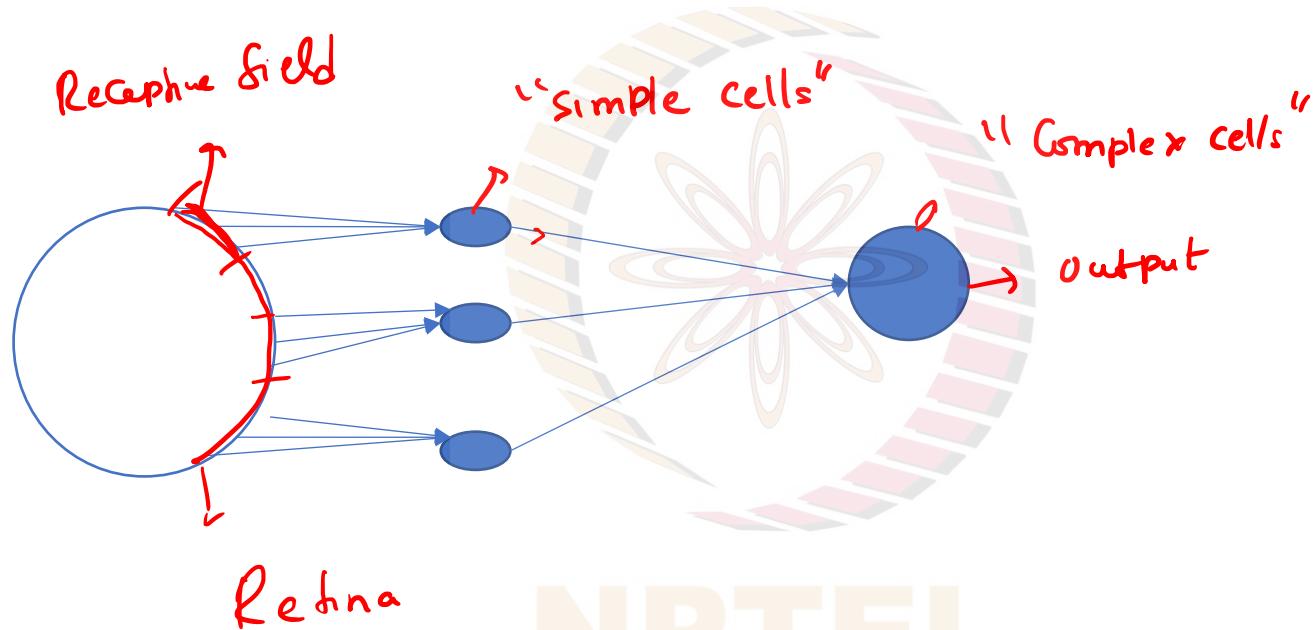


Why Convolutions? Biological Inspiration

- In 1962 Hubel and Wiesel did a series of studies measuring activations in the neuronal cells of cat visual cortex.
- Simple Cells- Respond to edges, lines etc in the receptive field of the cells- linear response
- Complex Cells- Take as input, activations from the simple cells and produces a rectified output- invariant to translation

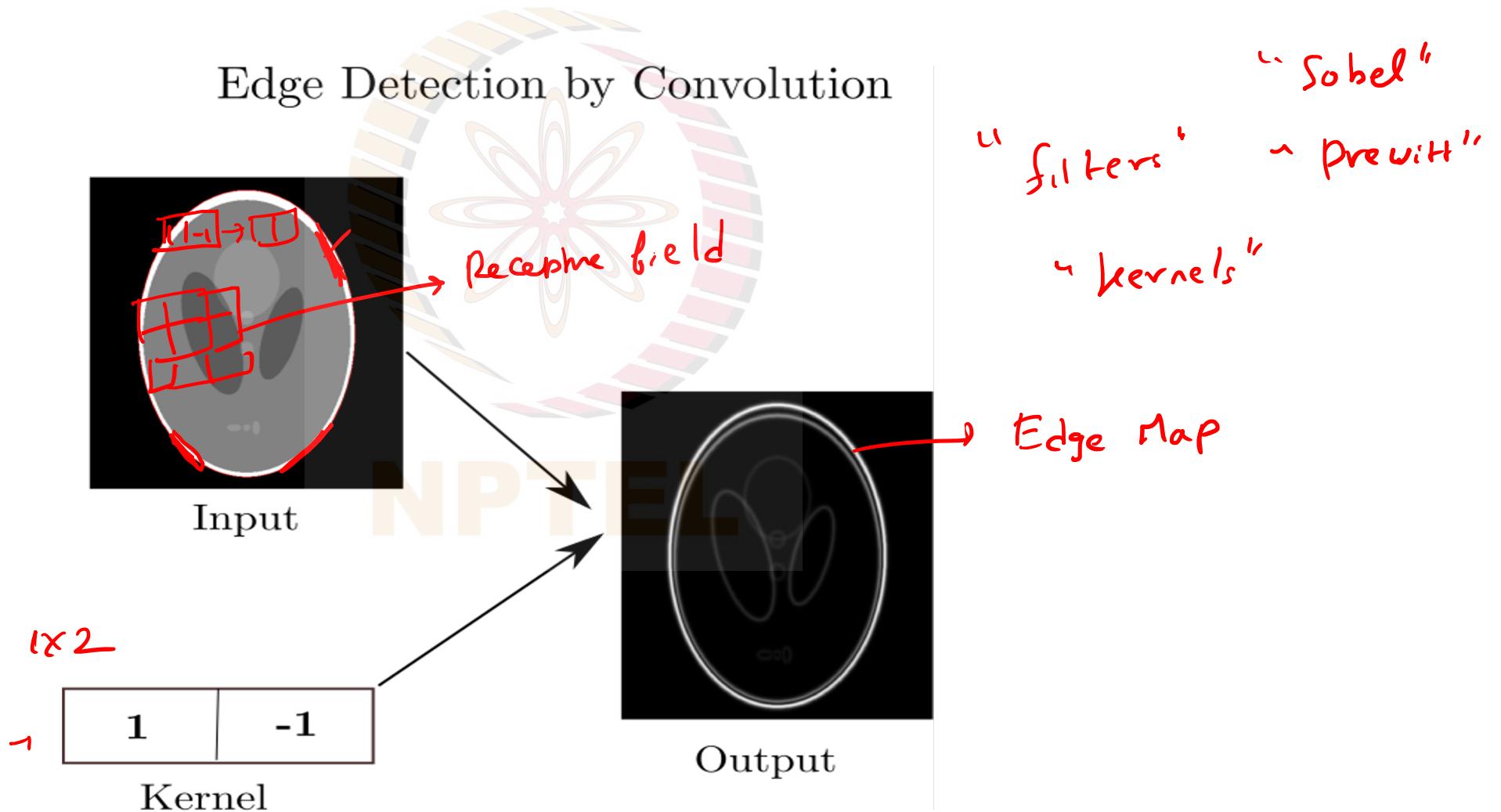
The logo for NPTEL (National Programme on Technology Enhanced Learning) is displayed as a watermark. It consists of the letters "NPTEL" in a bold, sans-serif font, with each letter in a different color: N is orange, P is yellow, T is light blue, E is red, and L is green. The letters are slightly overlapping and have a three-dimensional effect.

Biological Inspiration



NPTEL

Why Convolutions?



Convolutions

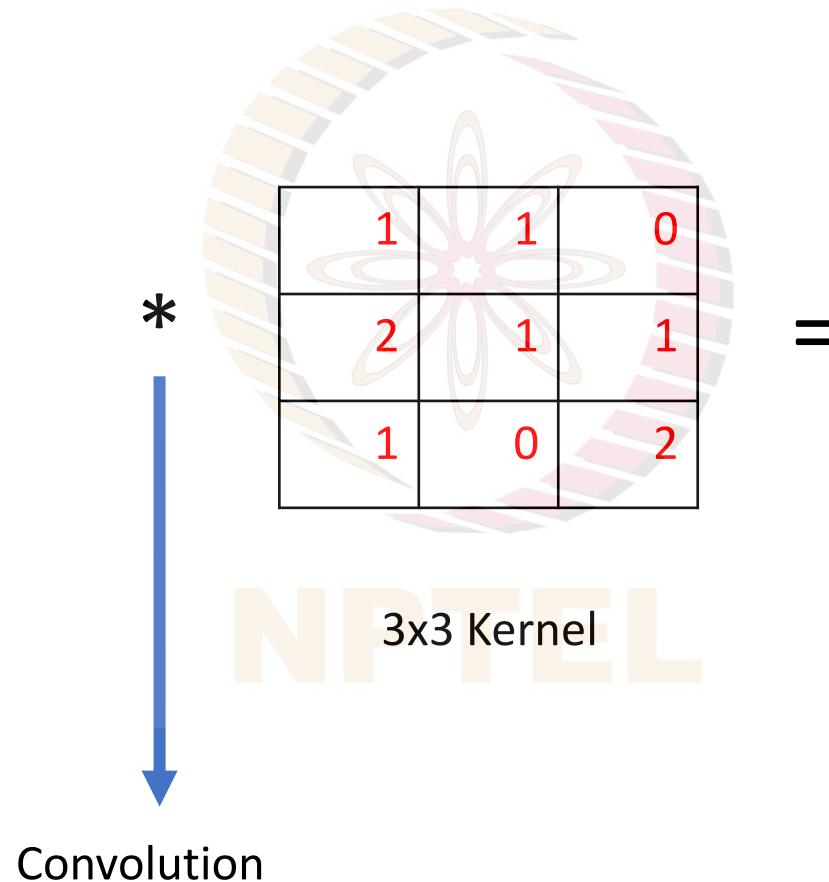
- Every output neuron is connected a small neighbourhood in the input through a weight matrix – Filter or kernel.
- We can define multiple kernels for every conv layer each giving rise to an output.
- Each filter is moved around the input giving rise to one 2D output. The outputs corresponding to each filter are stacked giving rise to an output volume.

The logo for NPTEL (National Programme on Technology Enhanced Learning) is displayed as a watermark. It consists of the word "NPTEL" in a bold, sans-serif font, with each letter in a different color: N is orange, P is red, T is blue, E is green, and L is yellow. Behind the letters are several thin, curved lines of the same colors, creating a stylized, radiating effect.

Convolution

1	0	2	2	1
0	2	1	1	3
7	0	1	2	1
5	1	3	2	2
2	3	6	1	5

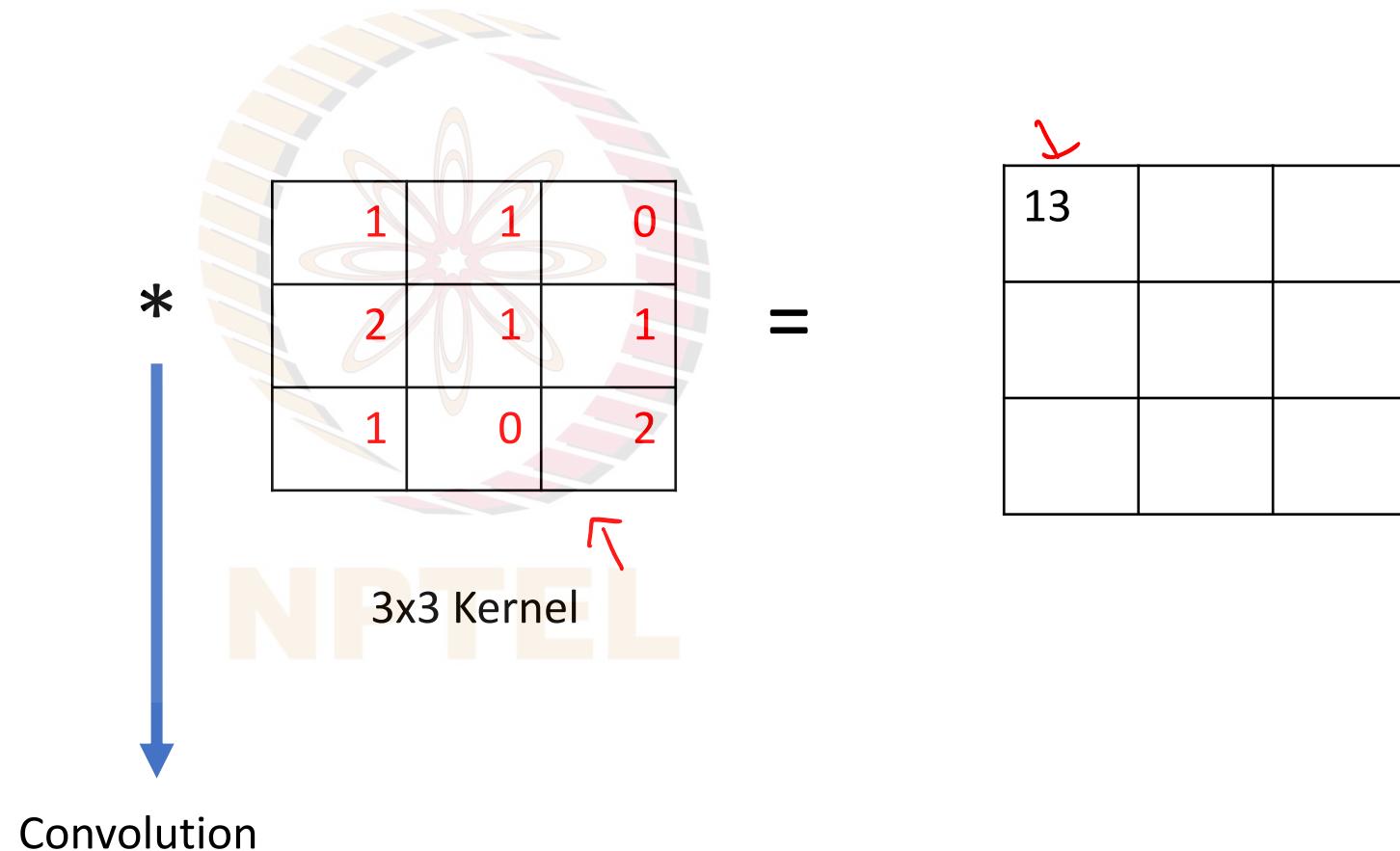
Image [5X5]



Convolution

1 1	0 1	2 0	2	1
0 2	2 1	1 1	1	3
7 1	0 0	1 2	2	1
5	1	3	2	2
2	3	6	1	5

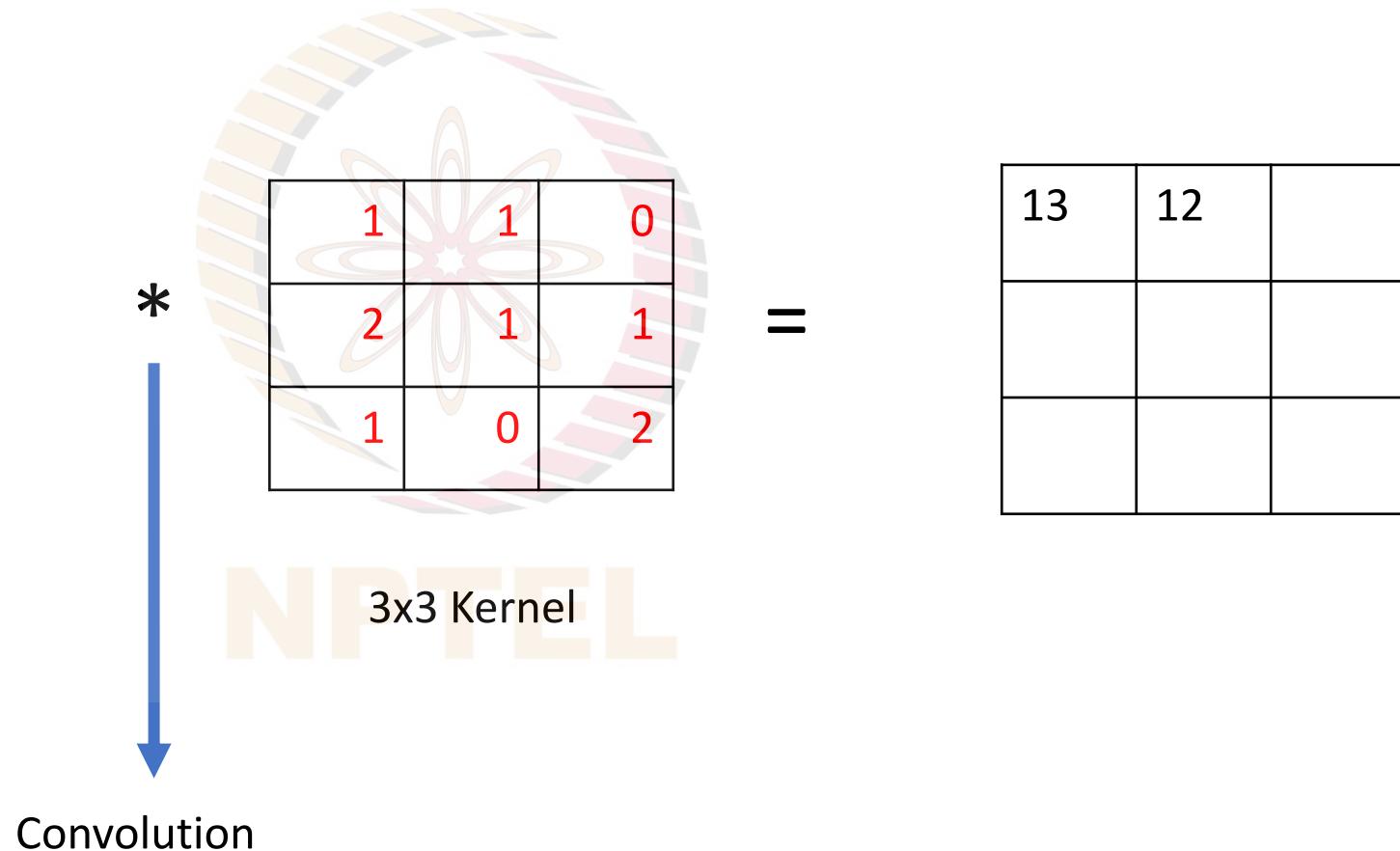
Image [5X5]



Convolution

1	0	1	2	1	2	0	1
0	2	2	1	1	1	1	3
7	0	1	1	0	2	2	1
5	1	3	2	2	2		
2	3	6	1	5			

Image [5X5]



Convolution

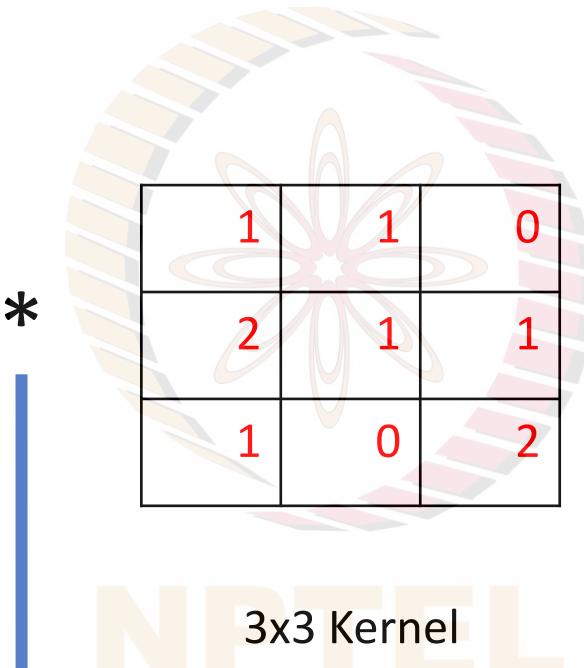
1	0	2	2	1	0
0	2	1	2	1	1
7	0	1	1	2	0
5	1	3	2	2	
2	3	6	1	5	

Image [5X5]

→

Convolution

*



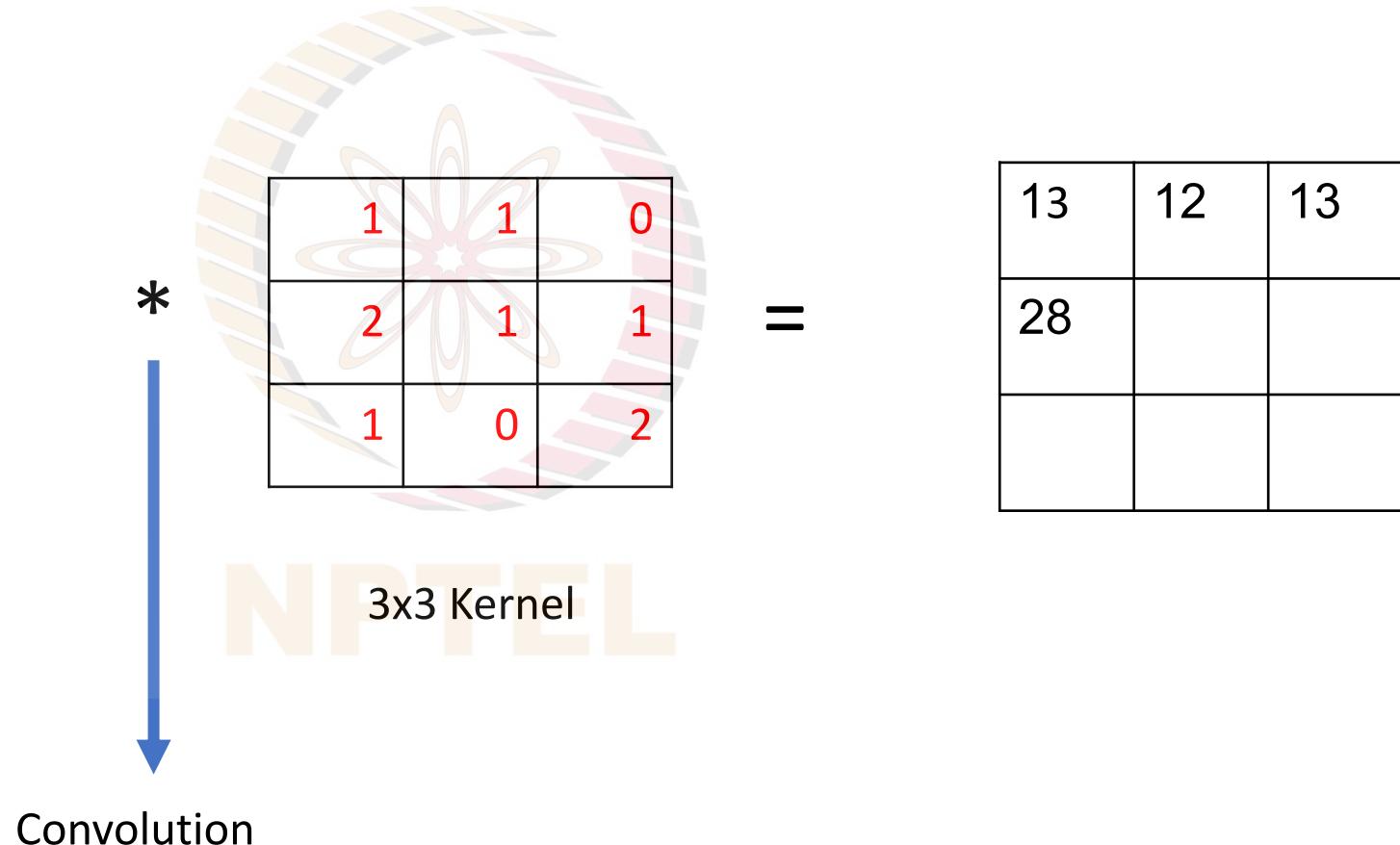
=

13	12	13

Convolution

1	0	2	2	1
0 1	2 1	1 0	1	3
7 2	0 1	1 1	2	1
5 1	1 0	3 2	2	2
2	3	6	1	5

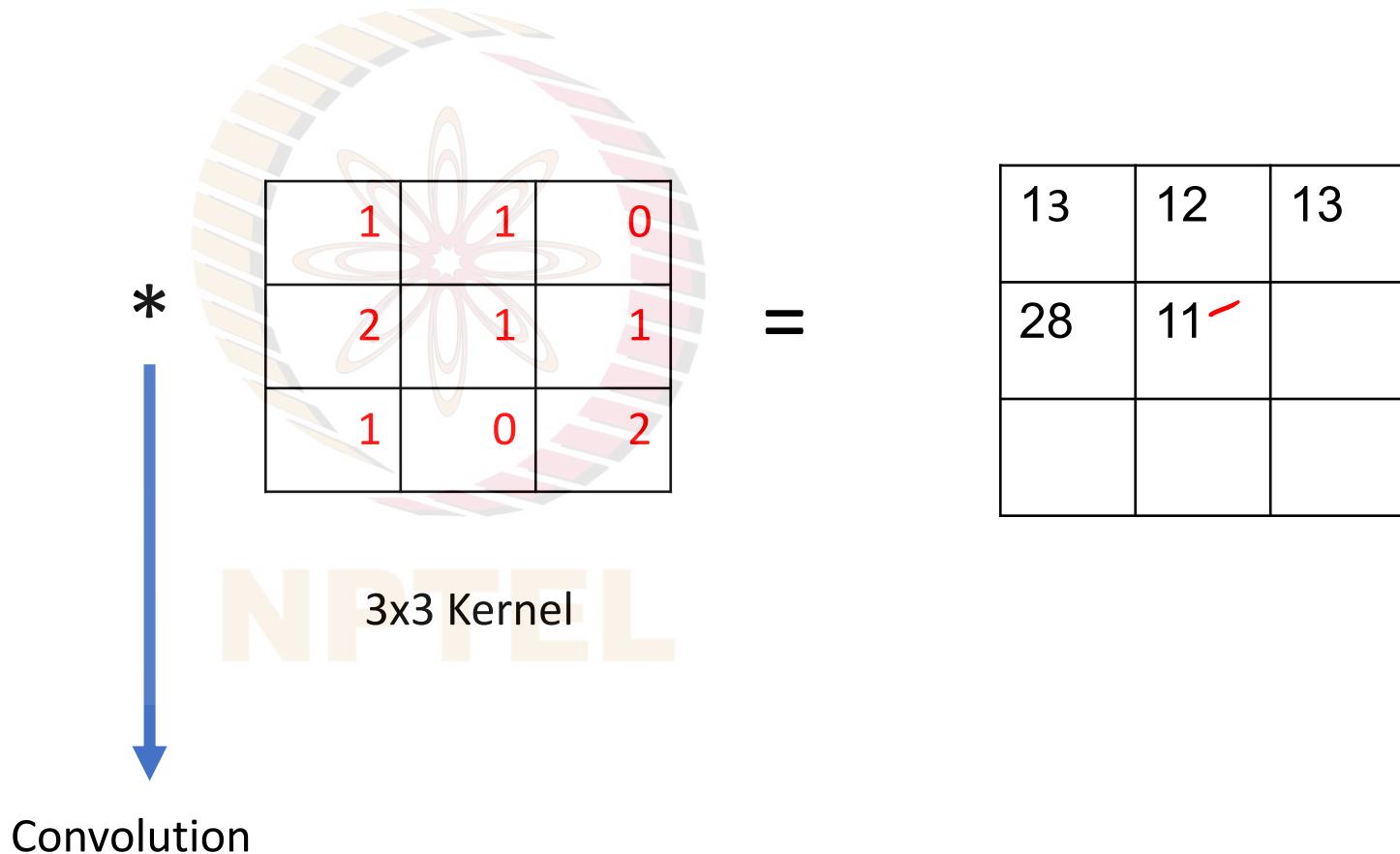
Image [5X5]



Convolution

1	0	2	2	1
0	2	1	1	1
7	0	2	1	1
5	1	1	3	0
2	3	6	1	5

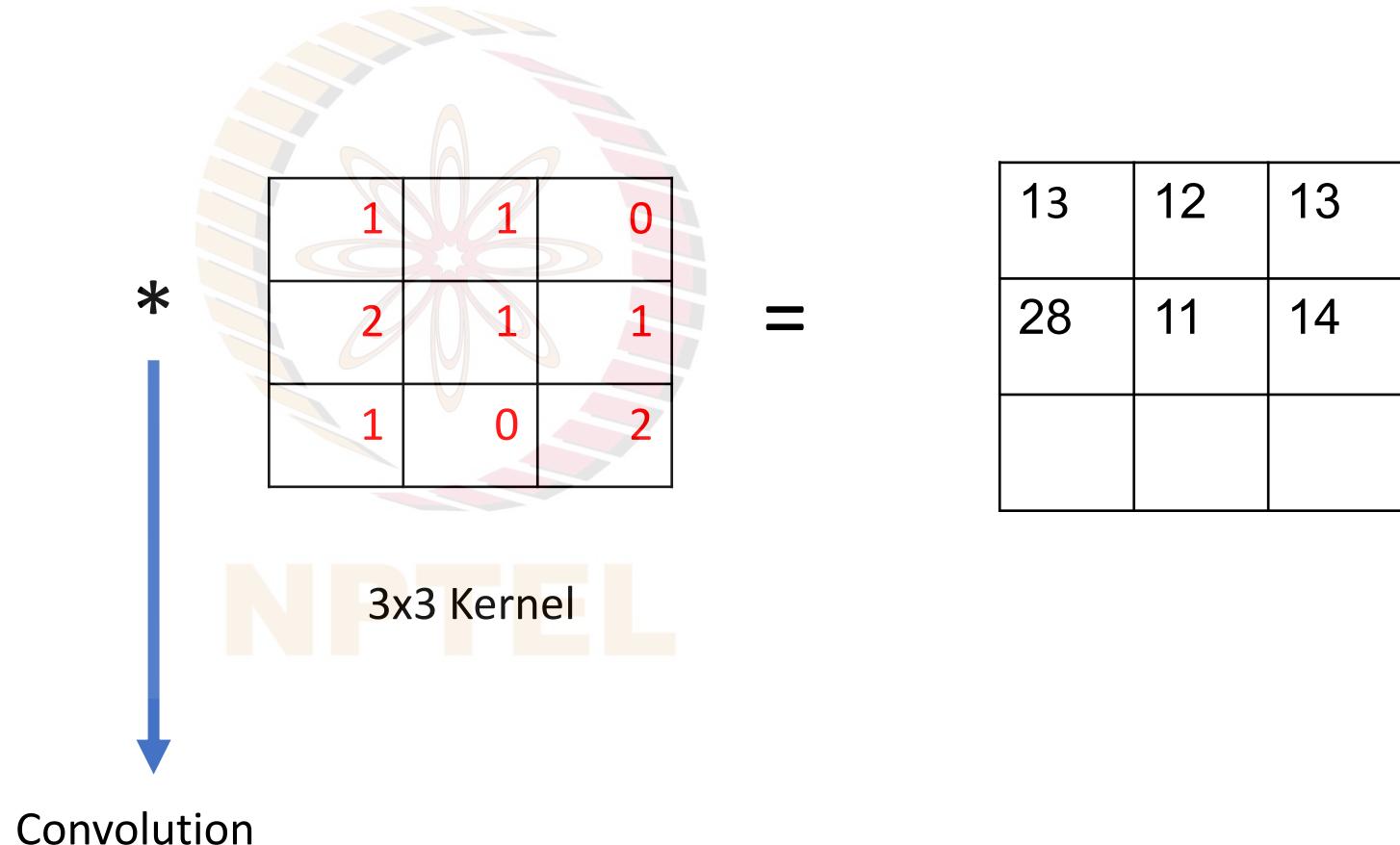
Image [5X5]



Convolution

1	0	2	2	1
0	2	1 1	1 1	3 0
7	0	1 2	2 1	1 1
5	1	3 1	2 0	2 2
2	3	6	1	5

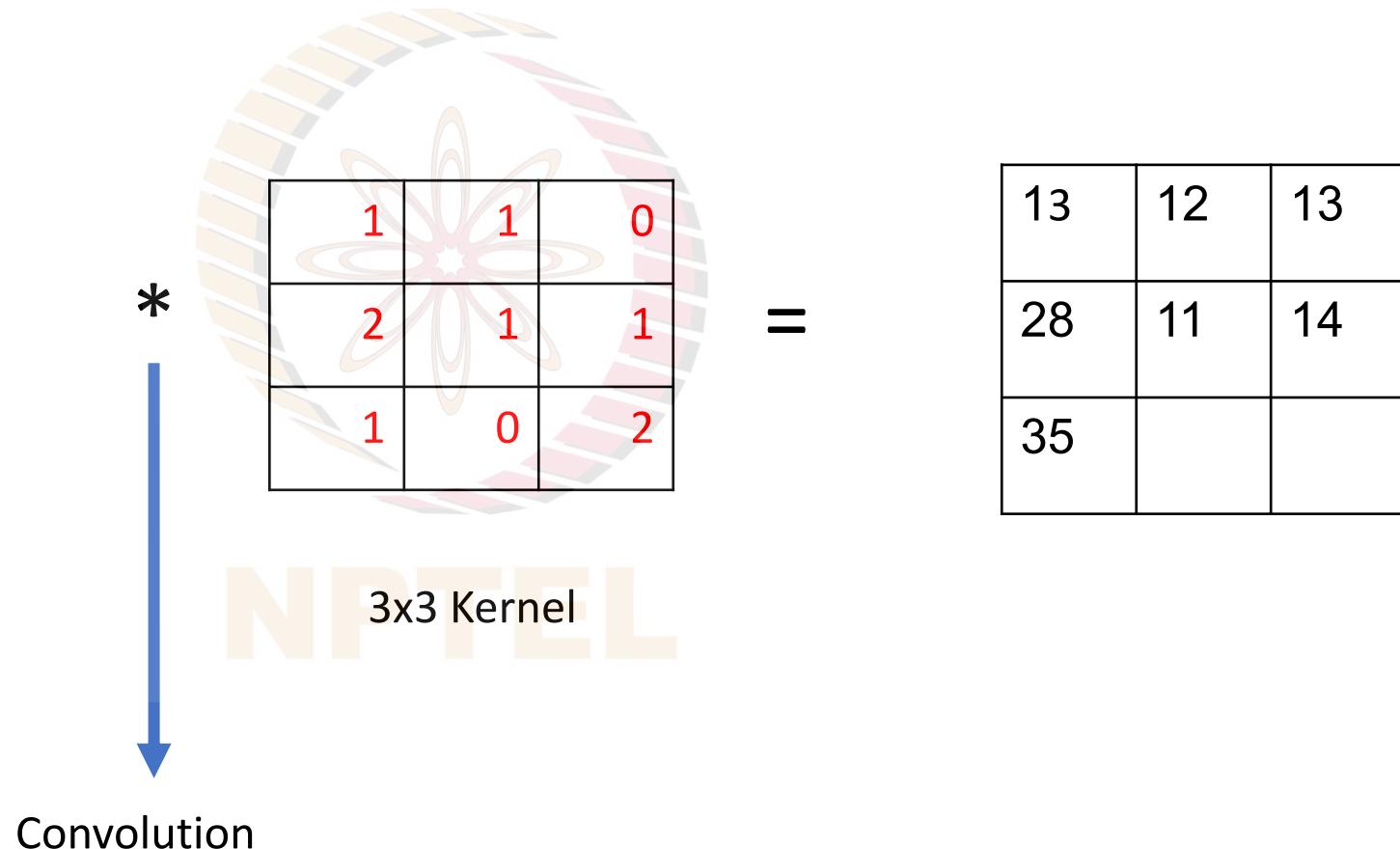
Image [5X5]



Convolution

1	0	2	2	1
0	2	1	1	3
7 1	0 1	1 0	2	1
5 2	1 1	3 1	2	2
2 1	3 0	6 2	1	5

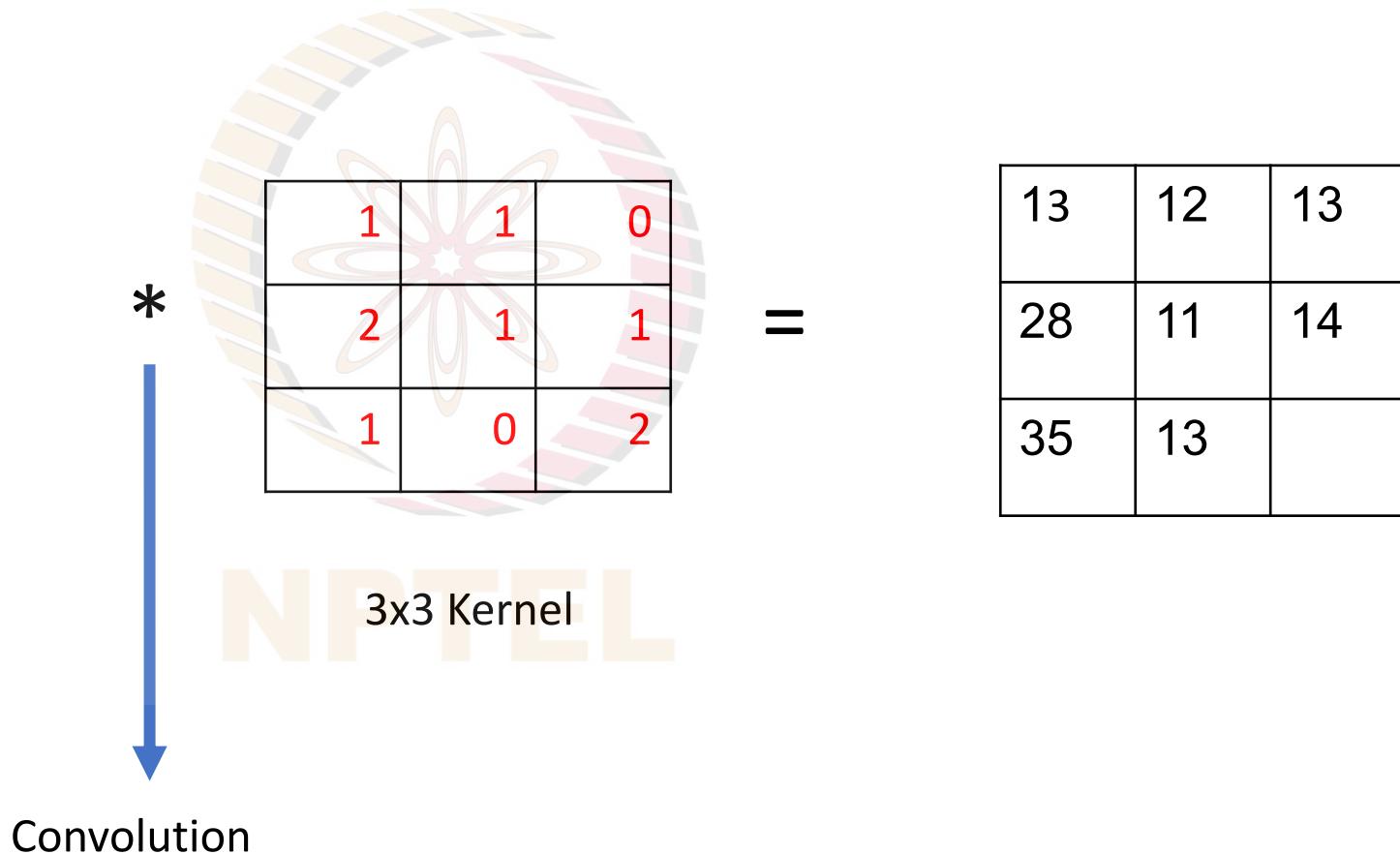
Image [5X5]



Convolution

1	0	2	2	1
0	2	1	1	3
7	0 1	1 1	2 0	1
5	1 2	3 1	2 1	2
2	3 1	6 0	1 2	5

Image [5X5]



Convolution

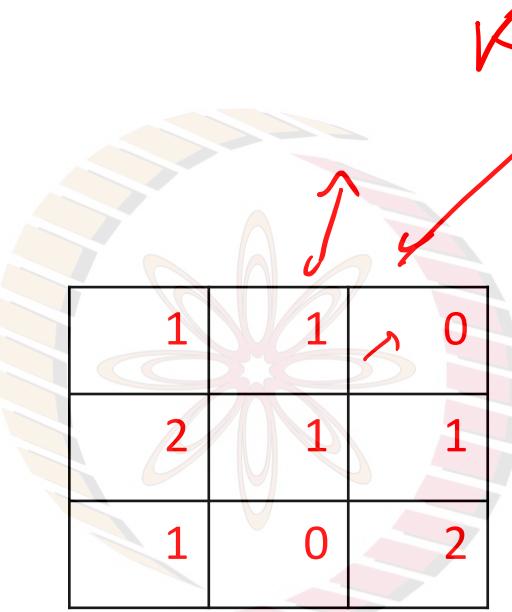
1	0	2	2	1
0	2	1	1	3
7	0	1	1	2
5	1	3	2	2
2	3	6	1	1

Image [5X5]

↓

Convolution

*



3x3 Kernel

$$n_x \times n_y , f_x, f_y$$

=

13	12	13
28	11	14
35	13	29

↓

$$\left. \begin{array}{l} n_x - f_x + 1 \\ n_y - f_y + 1 \end{array} \right\}$$

K' Feature Maps

Pooling

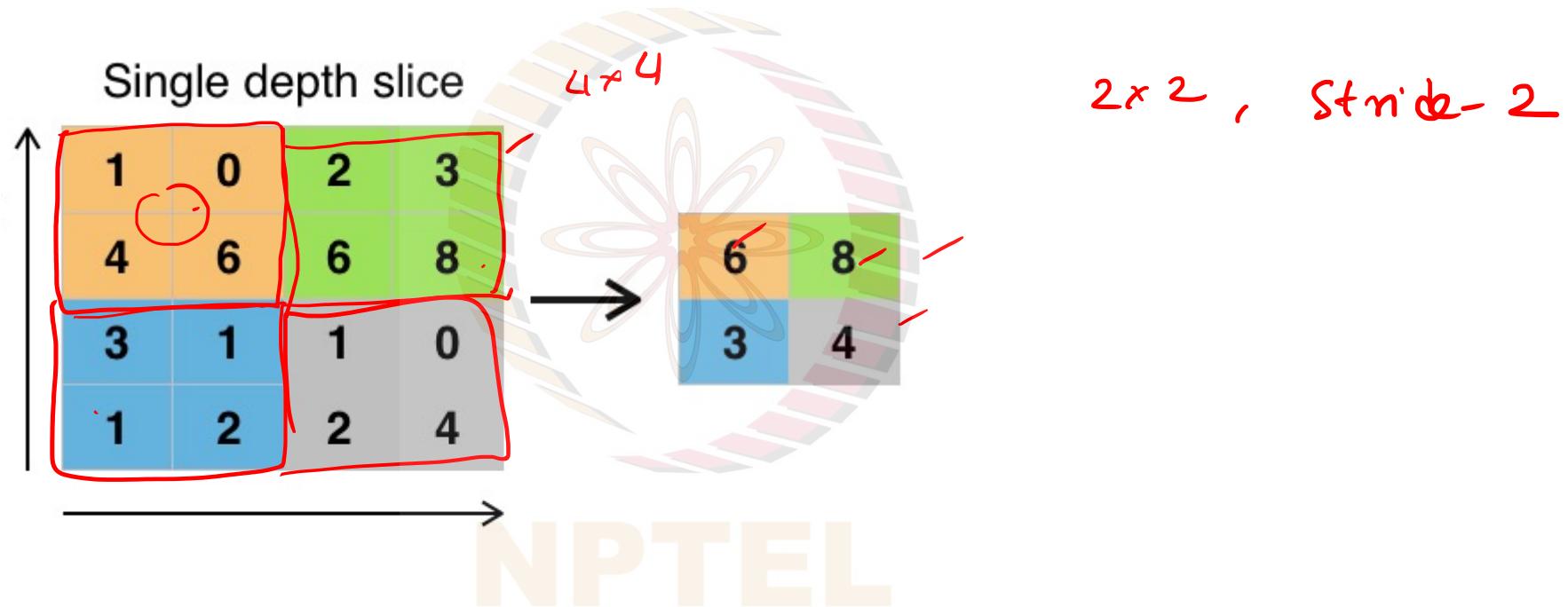
- Provides translational invariance by subsampling
- Reduces size of the feature maps
- Average Pooling and Max Pooling are commonly used

$$2 \text{ } 78 \times 256 \rightarrow \frac{32 \times 32}{= 11}$$

↑
factor 8 reduction

NPTEL

Max Pooling



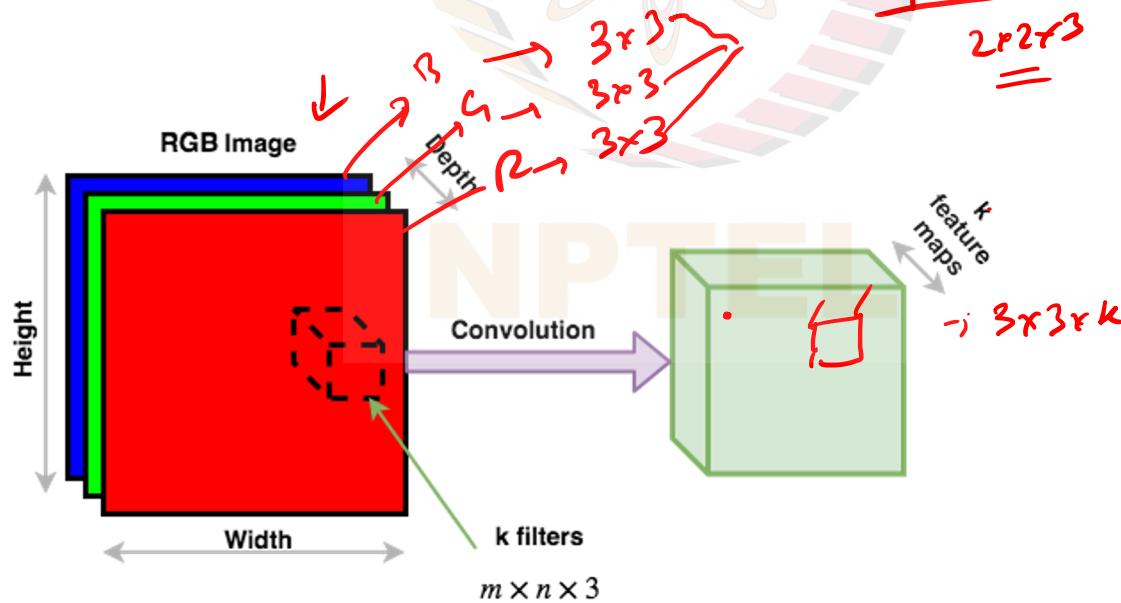
https://upload.wikimedia.org/wikipedia/commons/e/e9/Max_pooling.png

By Aphex34 [CC BY-SA 4.0 (<https://creativecommons.org/licenses/by-sa/4.0>)], from Wikimedia Commons

Volume Convolutions

Every layer of a ConvNet has the same API:

- Takes a 3D volume of numbers
- Outputs a 3D volume of number



$$\begin{array}{|c|c|c|} \hline & 1 & 1 \\ \hline 1 & & \\ \hline & -1 & 0 \\ \hline \end{array} \quad 2 \times 2 \times 3 = 27 + 1$$

$(3 \times 3 \times 3)$
filter \rightarrow
 $27 + 1$

RGB

3 channels

$n_x \times n_y \times 3$

$3 \times 3 \times 5 \rightarrow 4^5 + 1$

Size of Output Volume

- Size of the Output Volume or Feature map depend on
 - Size of Input Feature Map
 - Kernel size
 - Zero padding
 - **Stride**

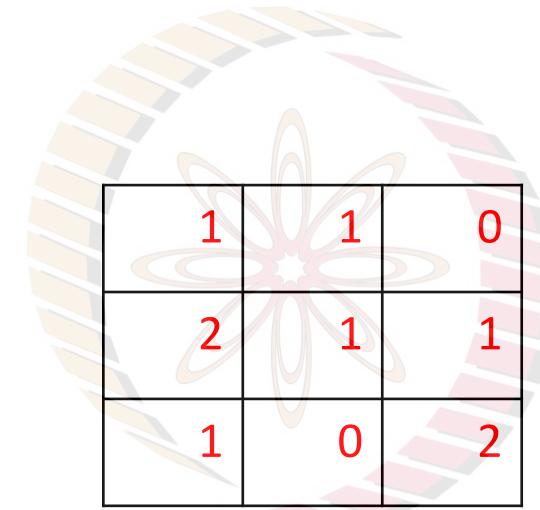


Padded Convolution

2	1	1
0	1	2
1	3	2



Image [3X3]



Convolution

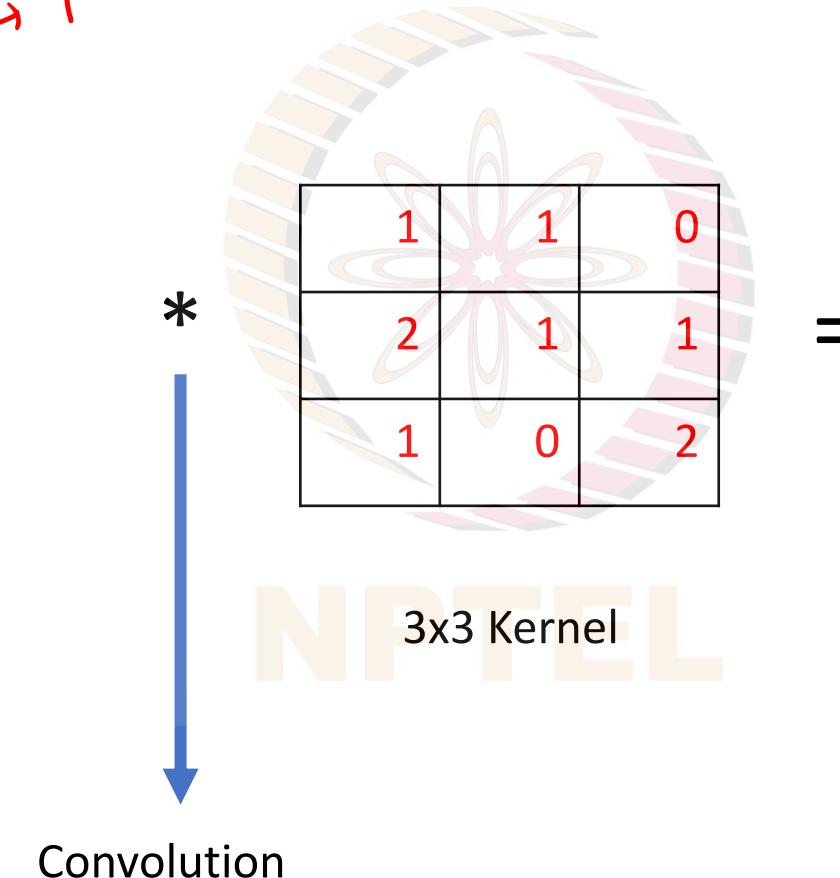
3x3 Kernel

(x1)

Padded Convolution

0	0	0	0	0
0	2	1	1	0
0	0	1	2	0
0	1	3	2	0
0	0	0	0	0

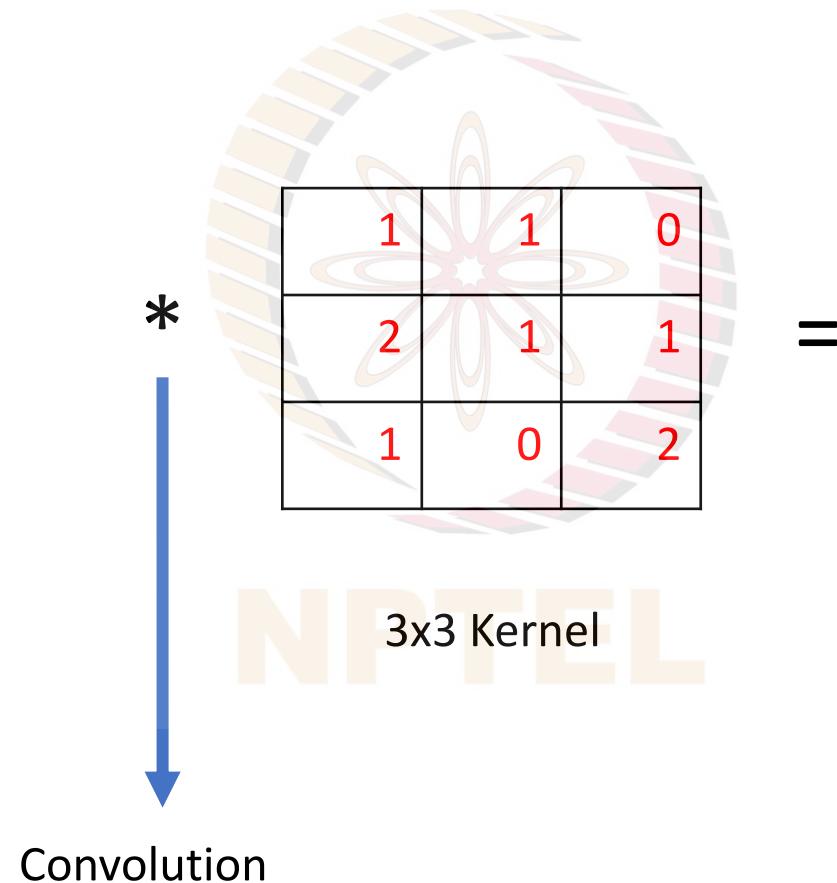
$P=1$



Padded Convolution

0	0	0	0	0
1	1	0	0	0
0	2	2	1	1
1	1	1	1	0
0	1	0	0	1
1	2	0	0	0
0	1	3	2	0
0	0	0	0	0

Image [3X3]

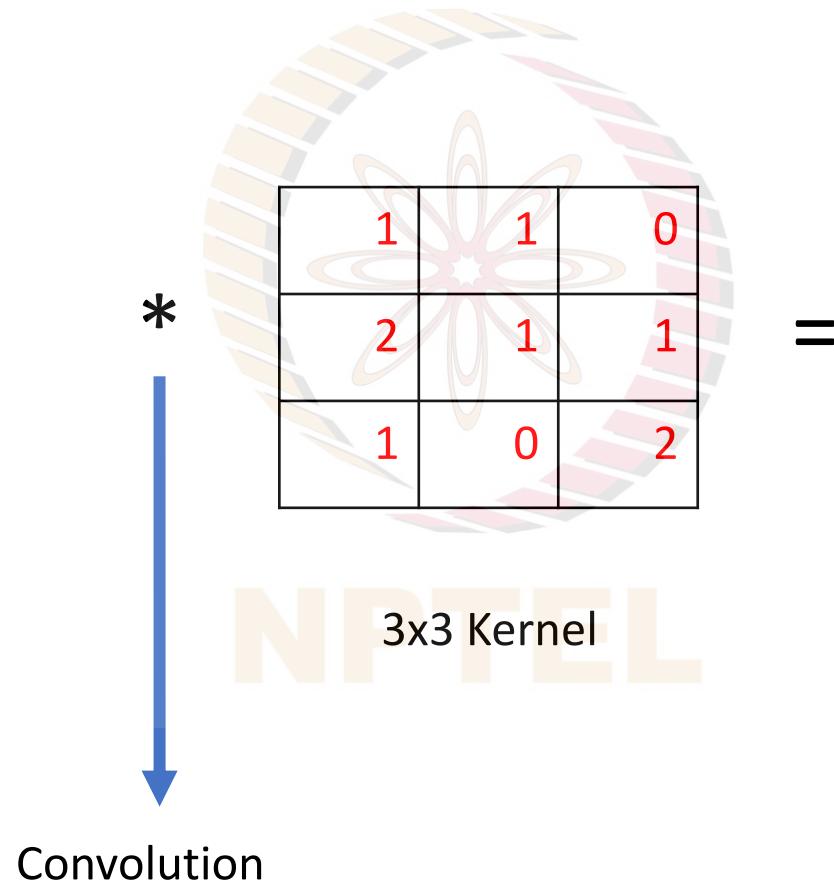


5		

Padded Convolution

0	0	1	0	1	0	0
0	2	2	1	1	1	1
0	0	1	1	0	2	2
0	1	3	2	2	0	0
0	0	0	0	0	0	0

Image [3X3]

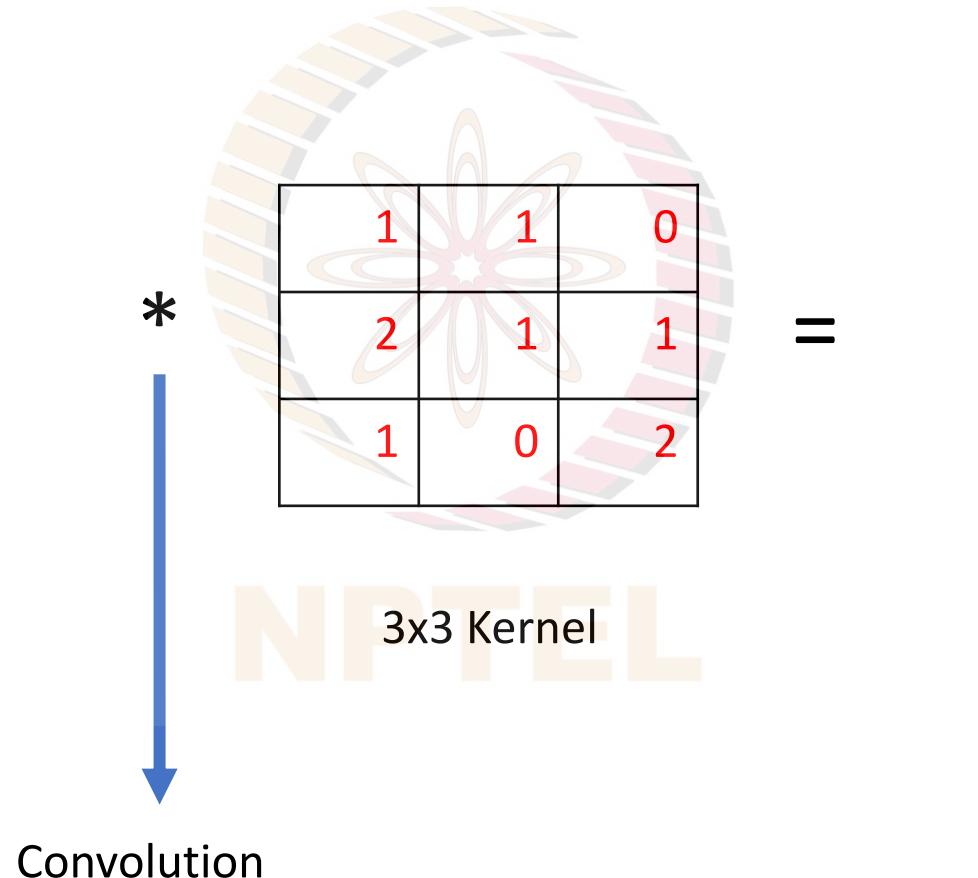


5	10	

Padded Convolution

0	0	0	0	0	0		
0	2	1	2	1	1	0	1
0	0	1	1	2	0	0	2
0	1	3	2	2	0		
0	0	0	0	0	0		

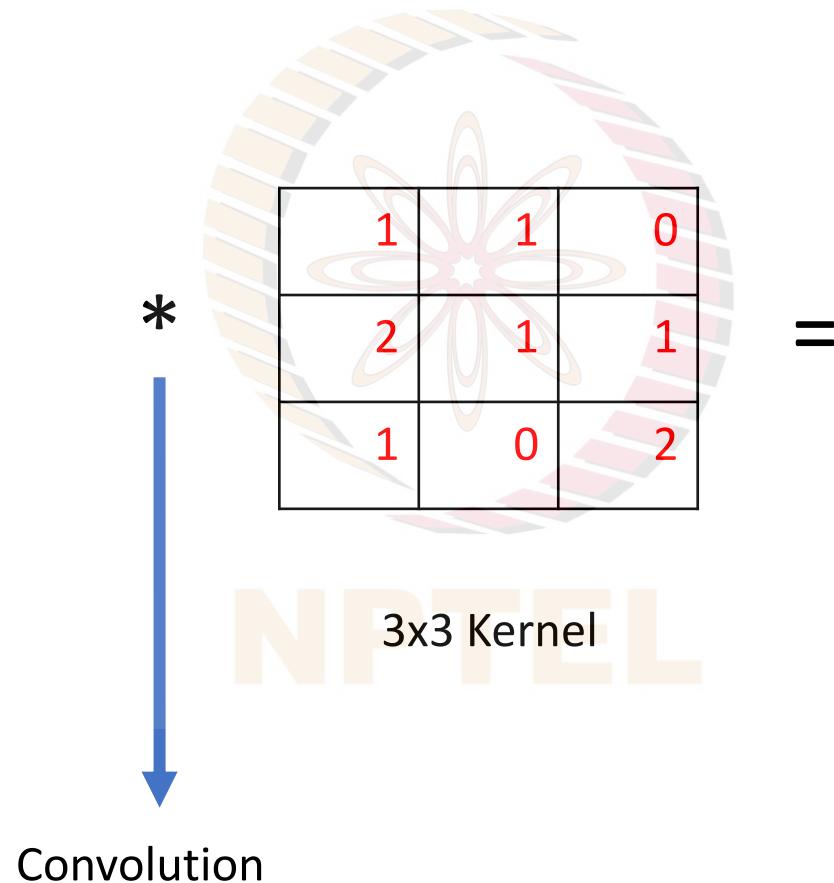
Image [3X3]



Padded Convolution

0	0	0	0	0
0	2	1	1	0
1	1	0	1	0
0	2	0	1	1
2	1	1	1	1
0	1	1	0	3
1	0	3	2	2
0	0	0	0	0

Image [3X3]

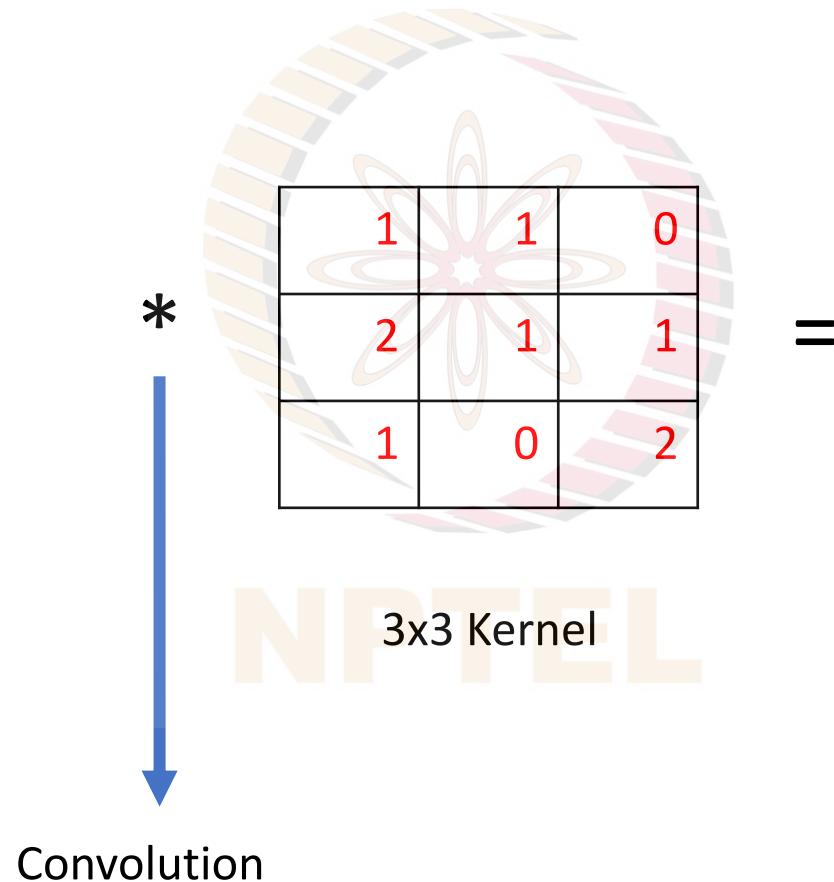


5	10	4
9		

Padded Convolution

0	0	0	0	0
0	2 1	1 1	1 0	0
0	0 2	1 1	2 1	0
0	1 1	3 0	2 2	0
0	0	0	0	0

Image [3X3]

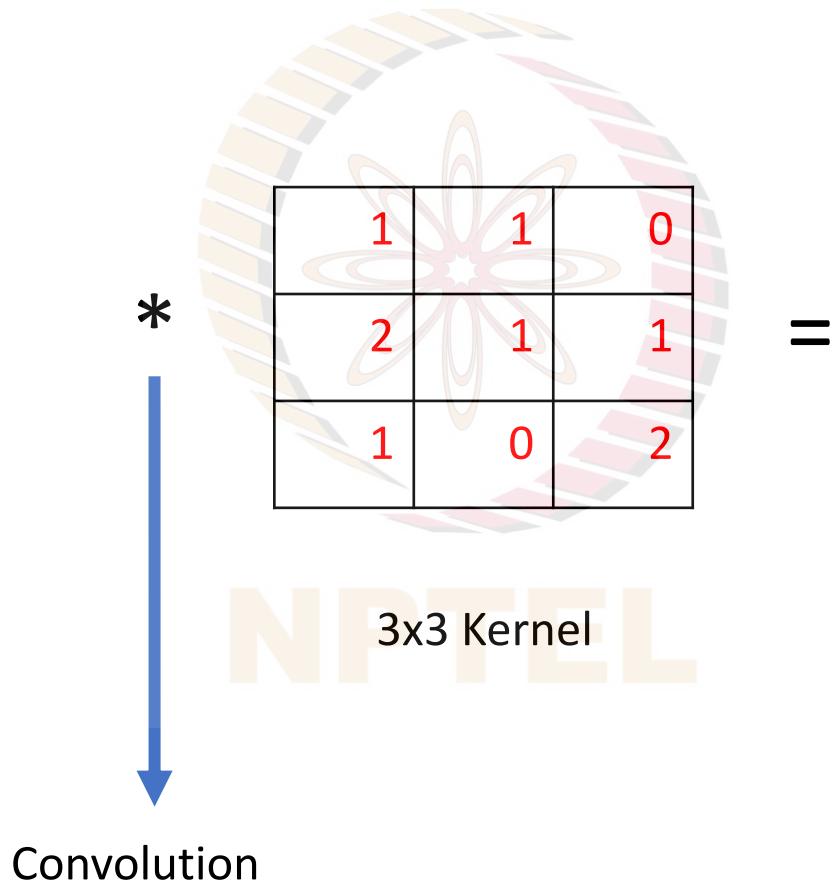


5	10	4
9	11	

Padded Convolution

0	0	0	0	0	0
0	2	1	1	1	0
0	0	1	2	2	1
0	1	3	1	2	0
0	0	0	0	0	0

Image [3X3]

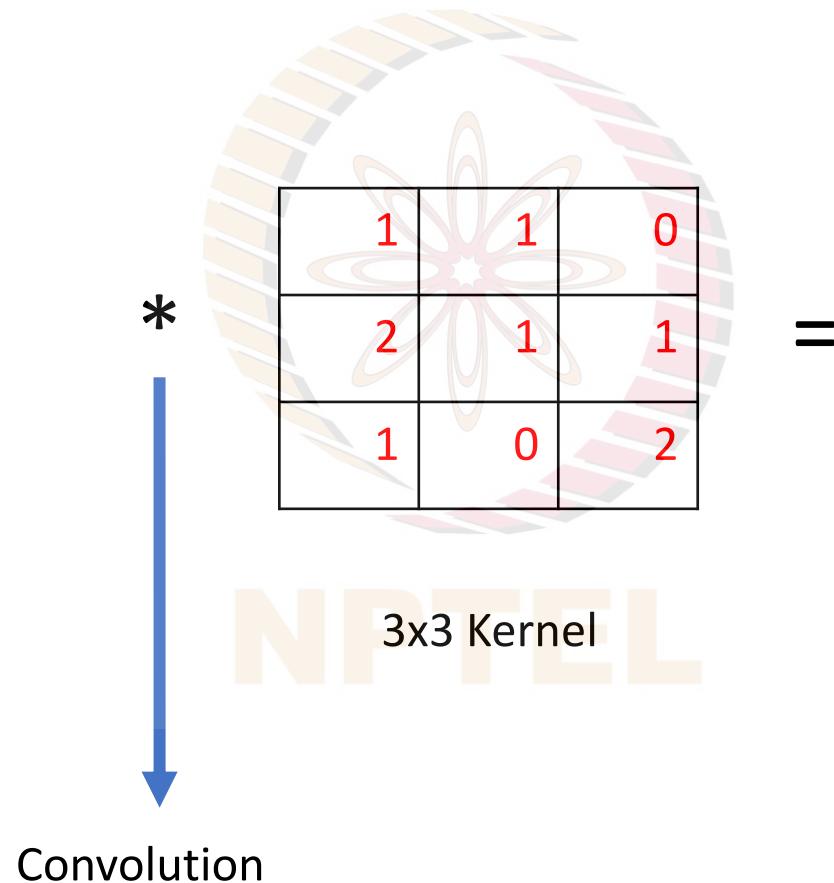


5	10	4
9	11	9

Padded Convolution

0	0	0	0	0
0	2	1	1	0
0	1	0	1	0
0	2	1	1	3
0	1	0	0	2

Image [3X3]

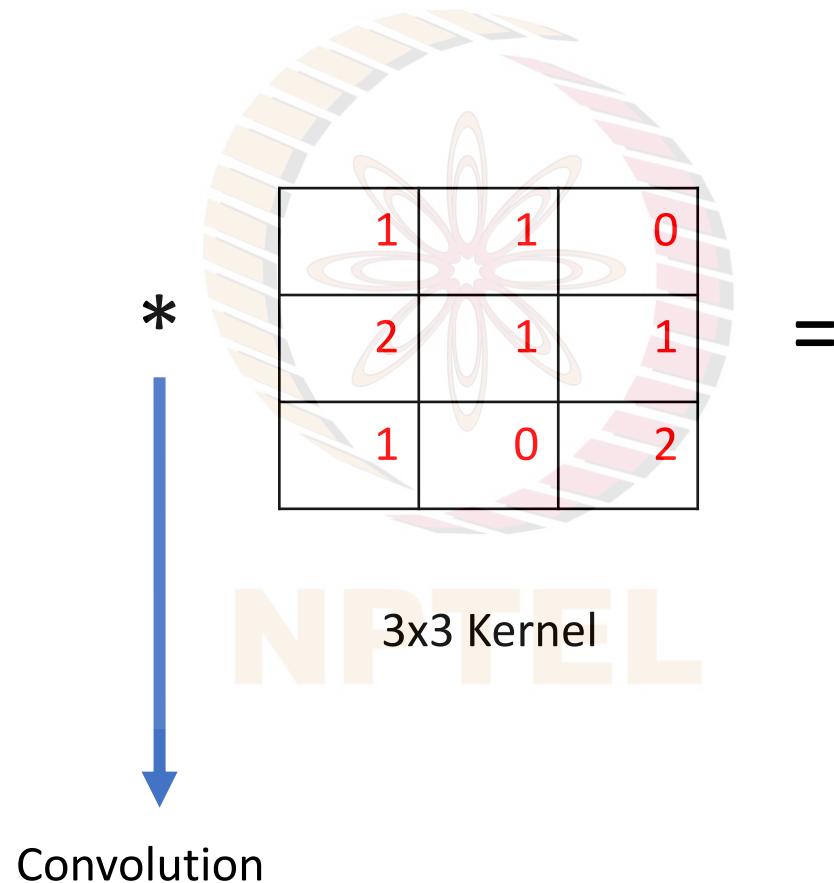


Convolution

Padded Convolution

0	0	0	0	0
0	2	1	1	0
0	0	1	1	2
0	1	2	3	1
0	0	1	0	2

Image [3X3]

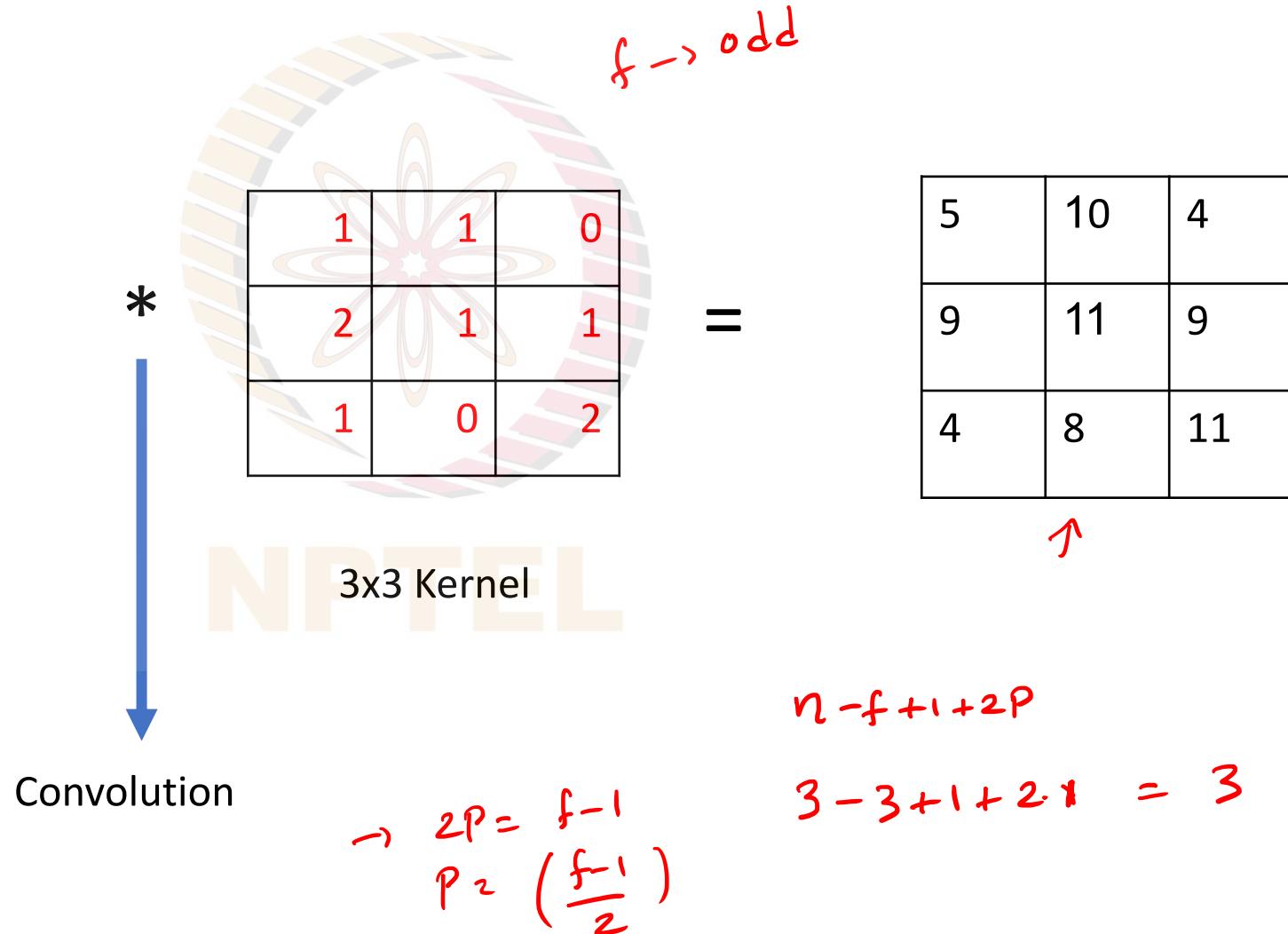


5	10	4
9	11	9
4	8	0

Padded Convolution

0	0	0	0	0
0	2	1	1	0
0	0	1 1	2 1	0 0
0	1	3 2	2 1	0 1
0	0	0 1	0 0	0 2

Image [3X3]



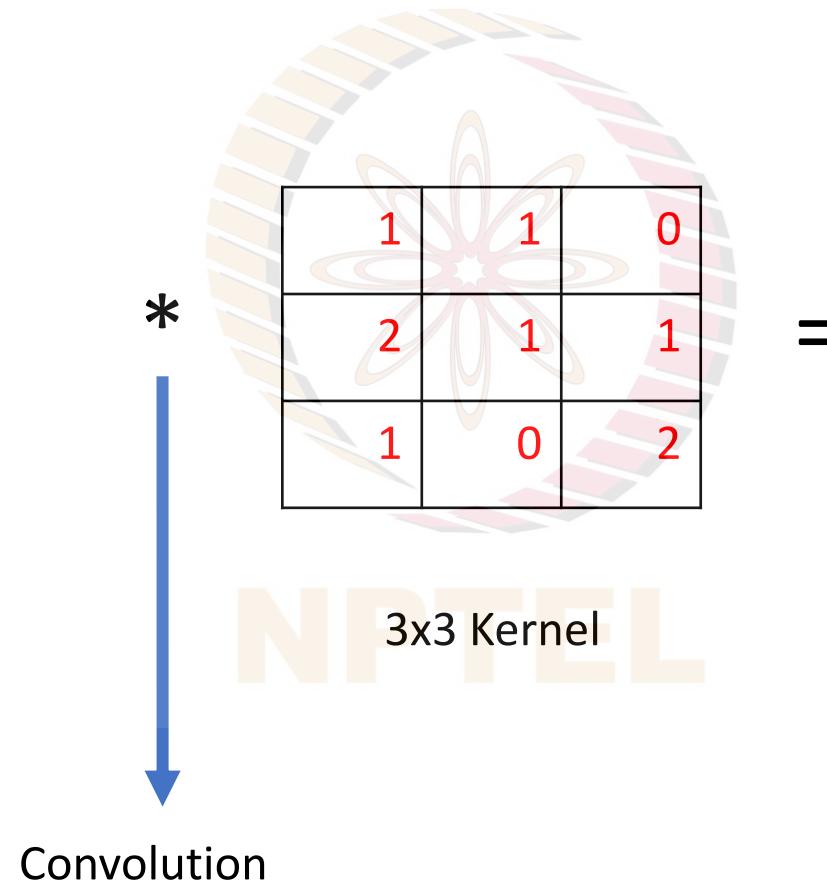
Strided Convolutions



Strided Convolution – Stride = 1

1	0	2	2	1
0	2	1	1	3
7	0	1	2	1
5	1	3	2	2
2	3	6	1	5

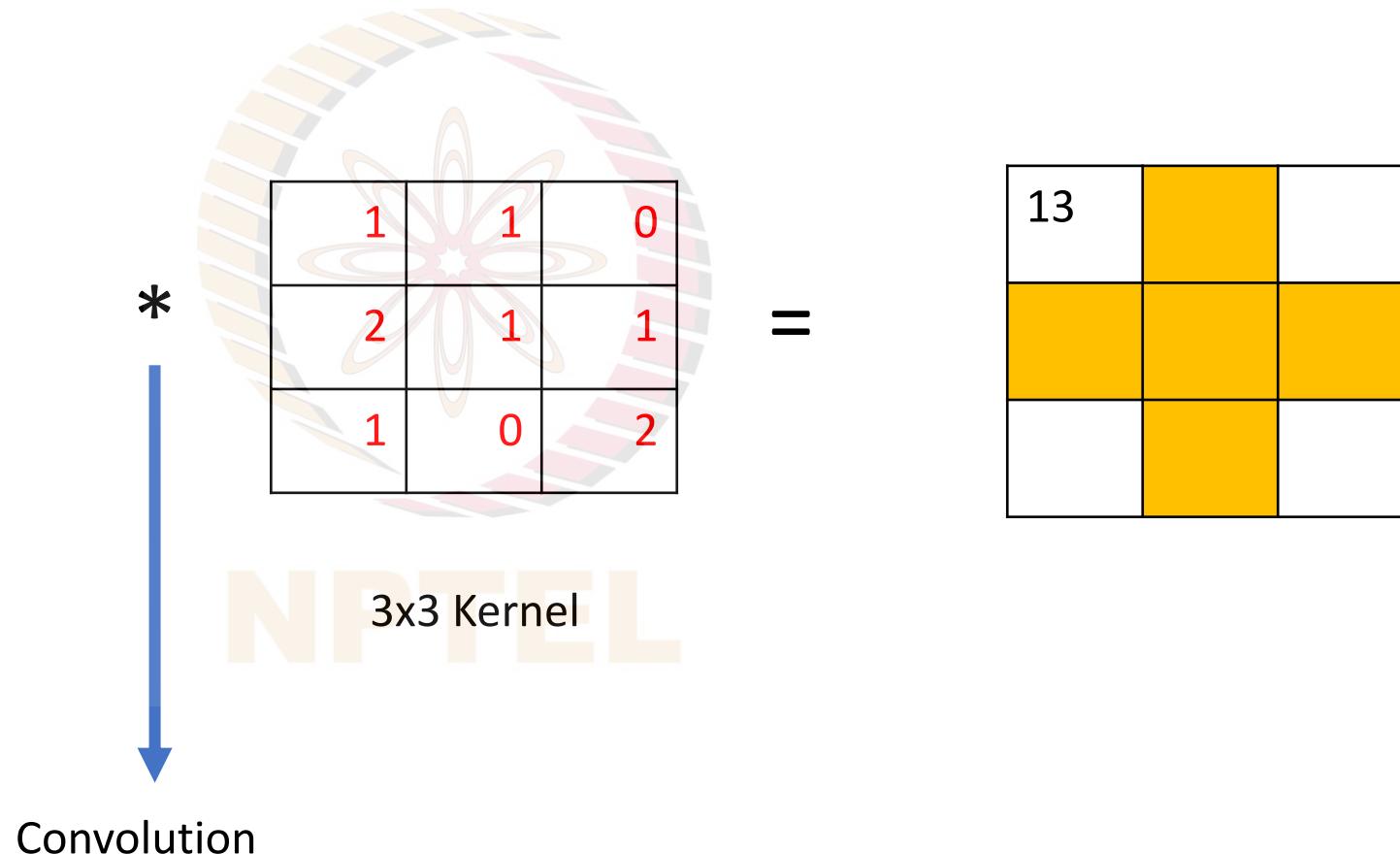
Image [5X5]



Strided Convolution

1 1	0 1	2 0	2	1
0 2	2 1	1 1	1	3
7 1	0 0	1 2	2	1
5	1	3	2	2
2	3	6	1	5

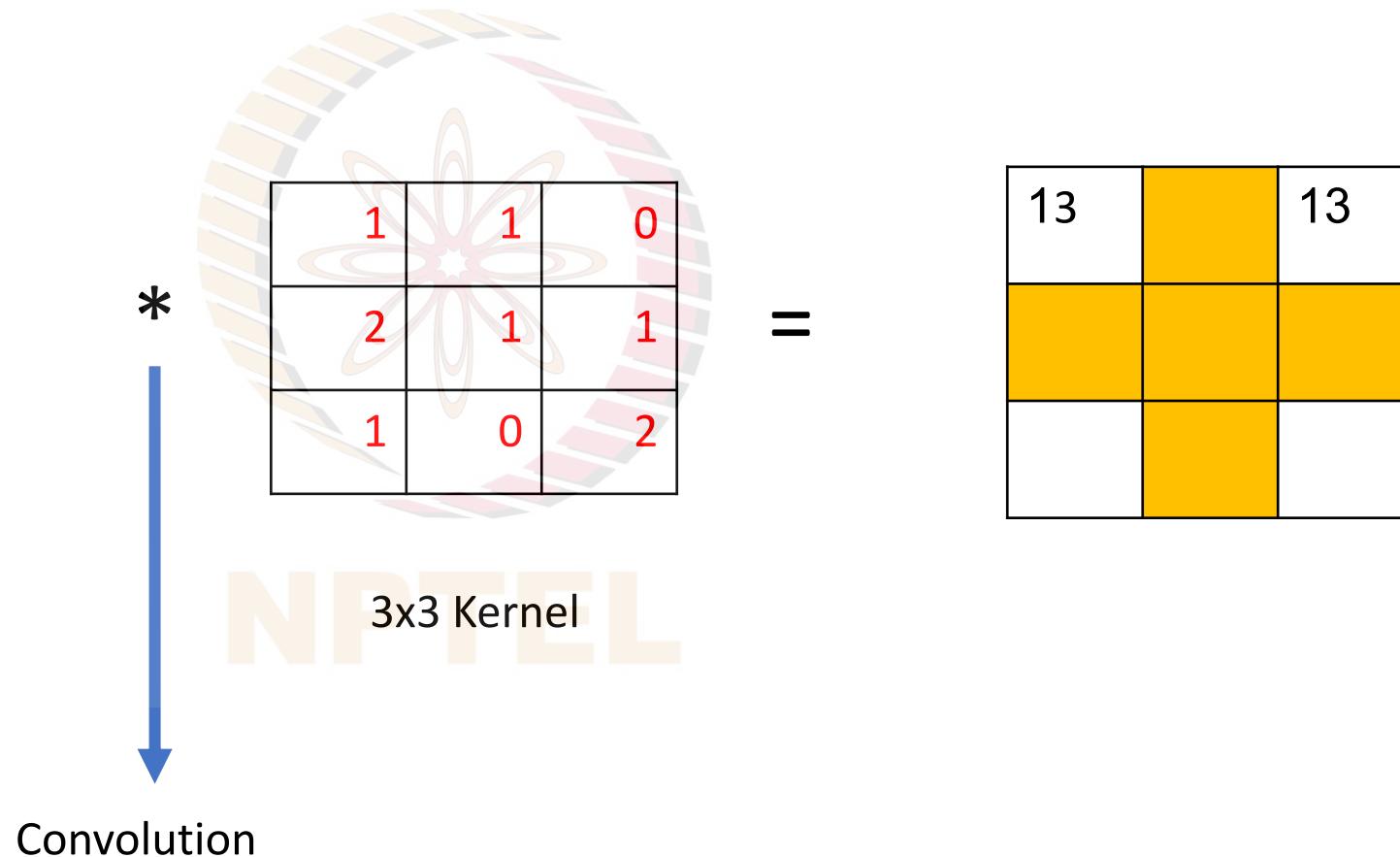
Image [5X5]



Strided Convolution

1	0	2	2	1
0	2	1	2	1
7	0	1	1	2
5	1	3	2	2
2	3	6	1	5

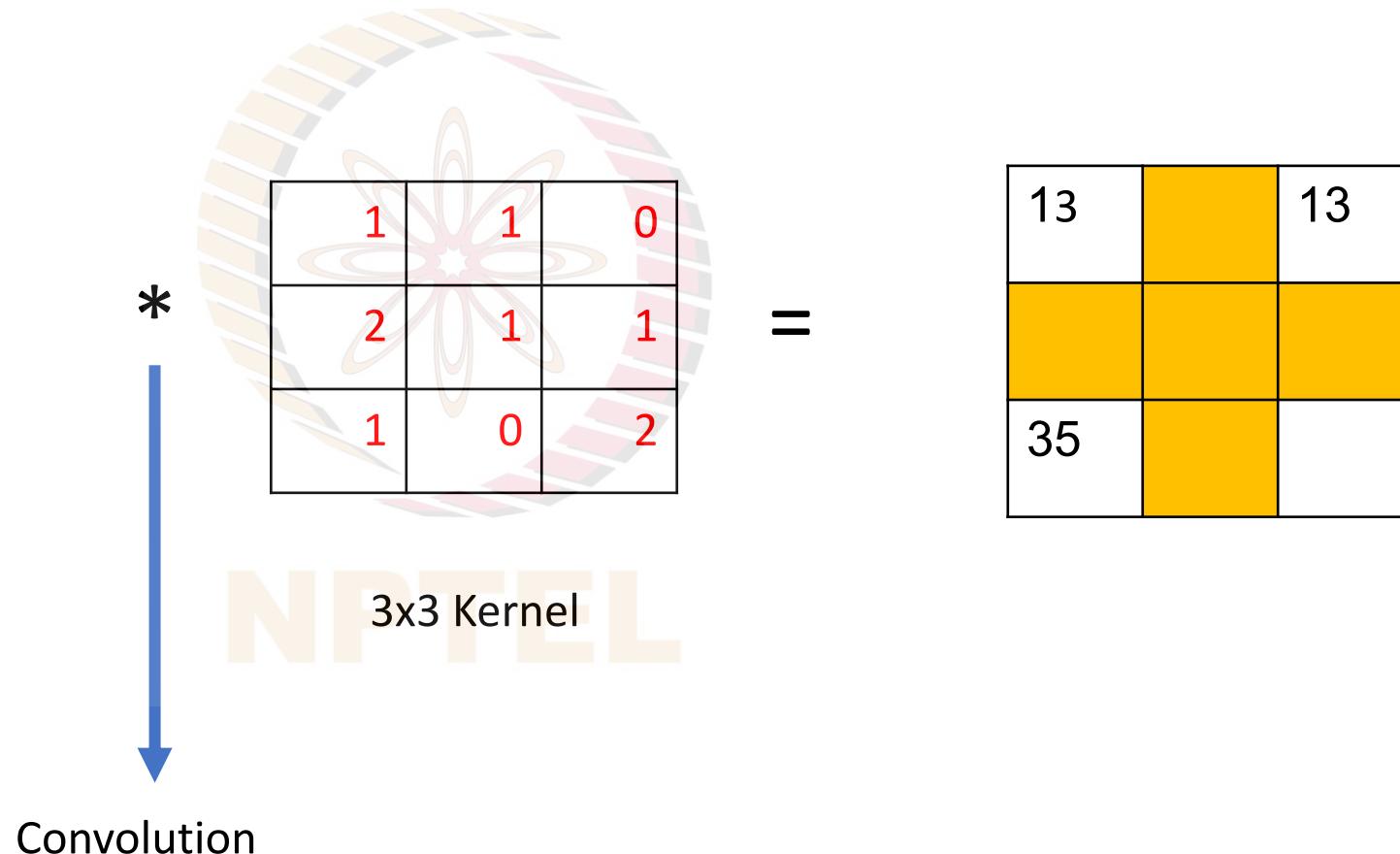
Image [5X5]



Strided Convolution

1	0	2	2	1
0	2	1	1	3
7 1	0 1	1 0	2	1
5 2	1 1	3 1	2	2
2 1	3 0	6 2	1	5

Image [5X5]

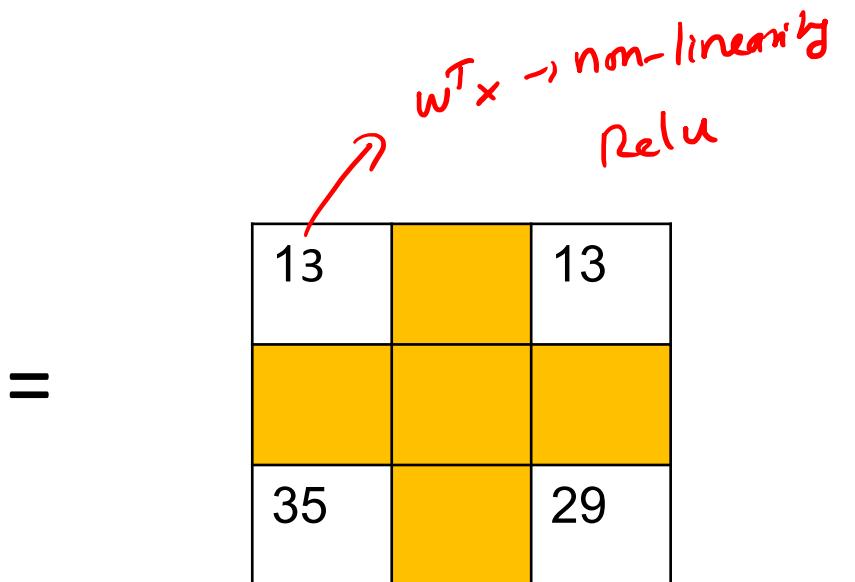
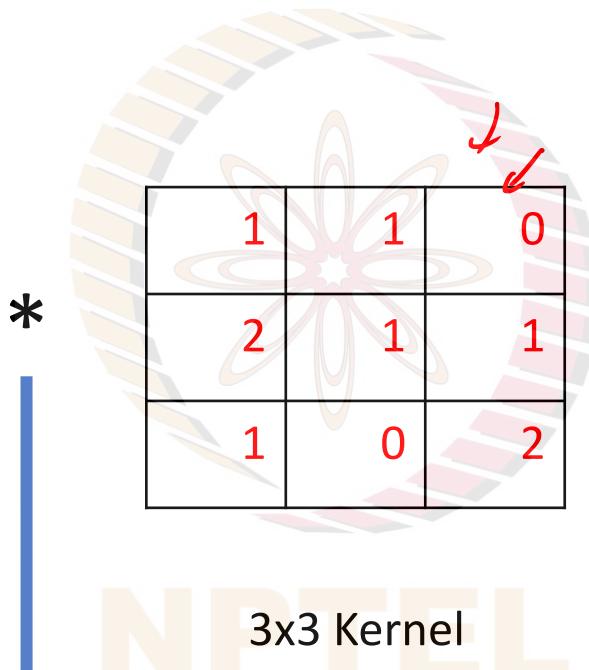


Strided Convolution

1	0	2	2	1
0	2	1	1	3
7	0	1	1	2 1
5	1	3 2	2 1	2 1
2	3	6 1	1 0	5 2

Image [5X5]

Convolution



$$\frac{n-f+2P}{s} + 1$$
$$\frac{5-3+0}{2} \rightarrow \frac{2}{2} = 1$$