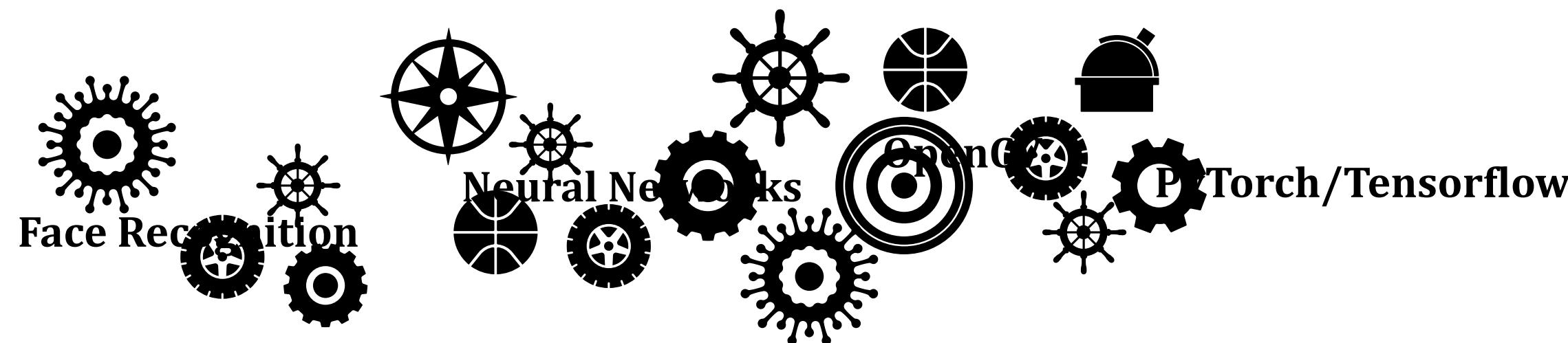
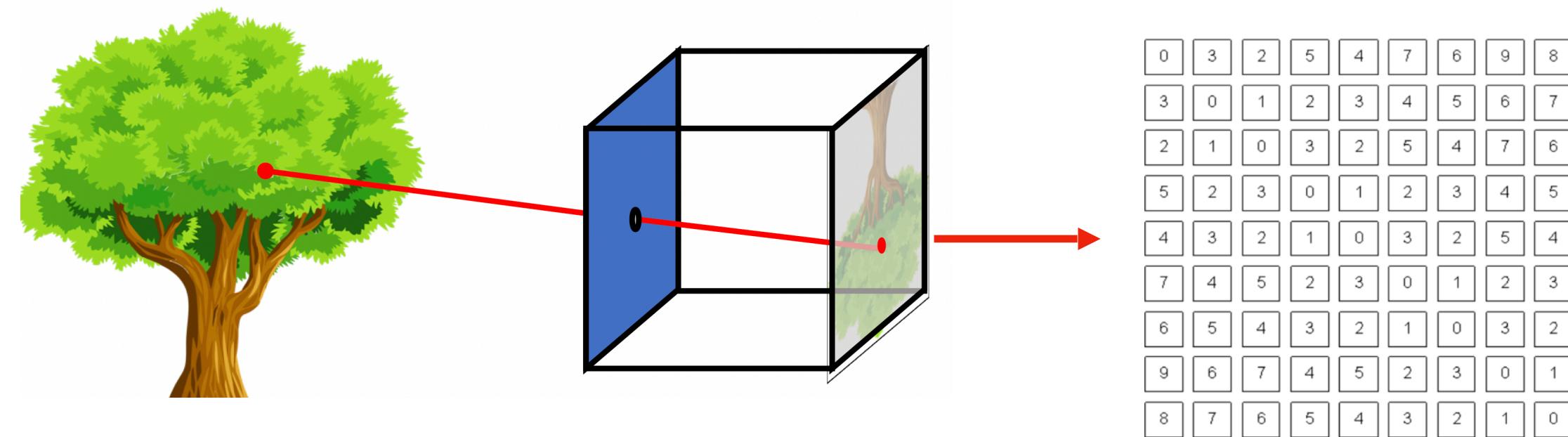


Computer Vision

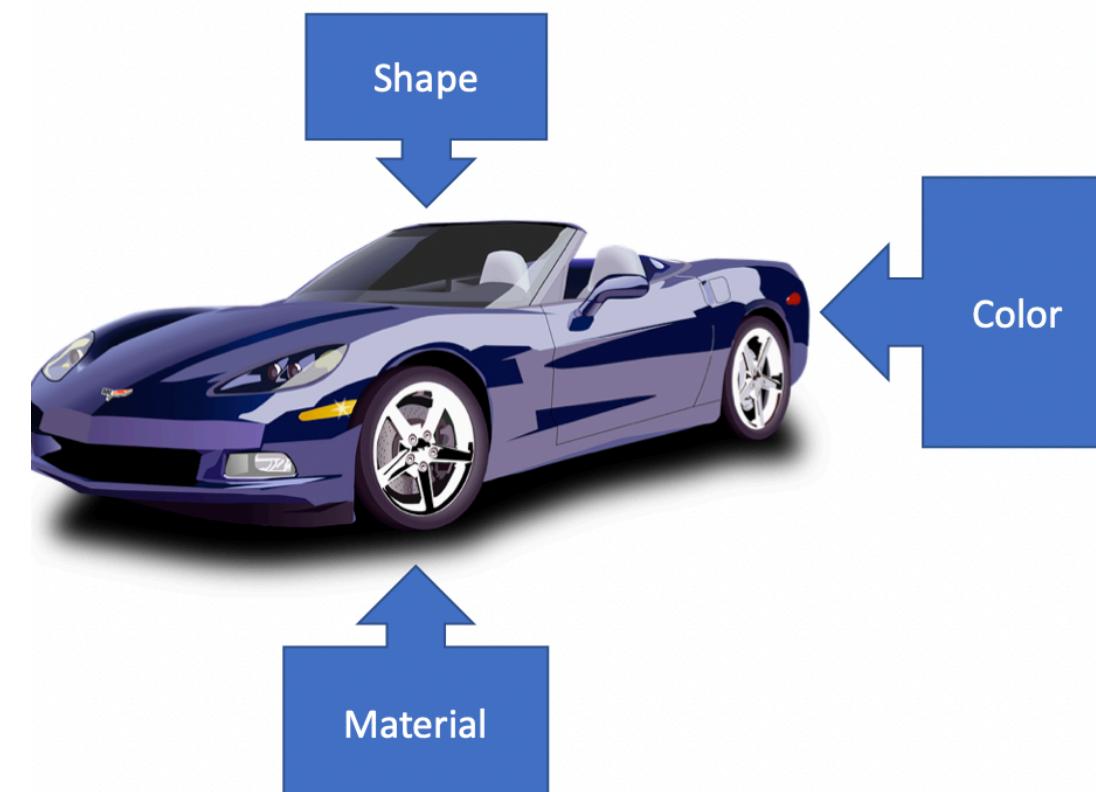


Face Recognition Neural Networks OpenCV PyTorch/Tensorflow
Object Detection Deep Learning Segmentation
 Image Classification

Computer Vision



- Computer vision process 2D matrix of picture element (pixel)
- In contrast, real world objects does not deal with pixels but with physical attributes such as shape, colour and material.

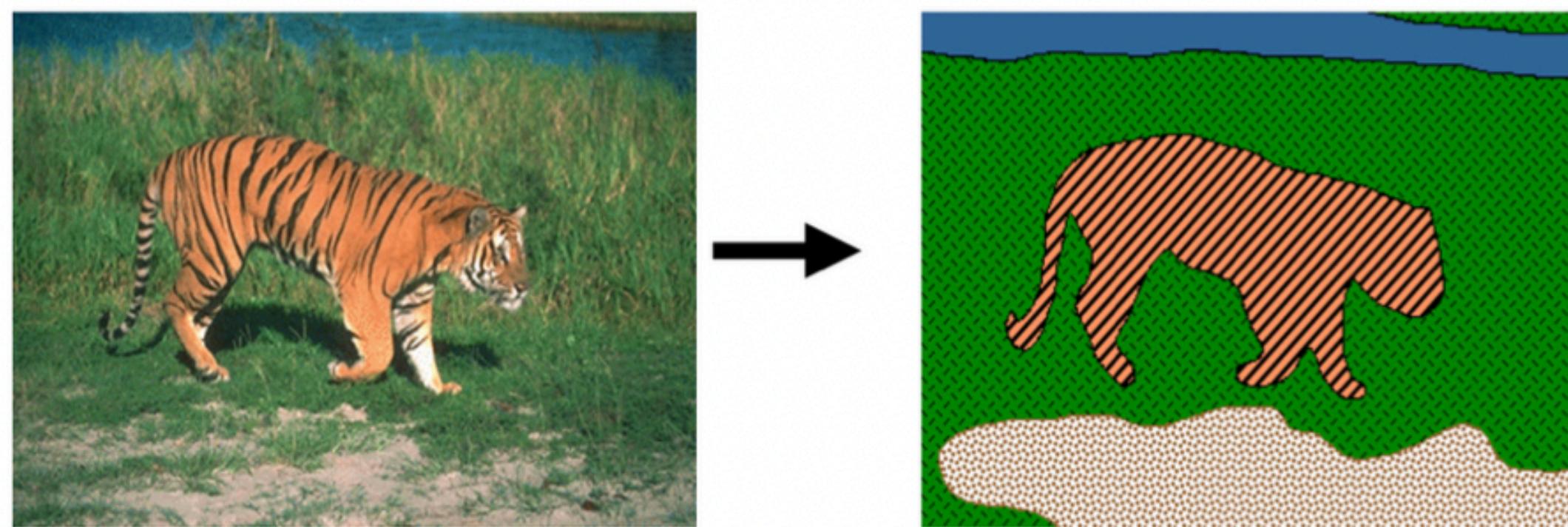


Not all 2D arrays are images



Computer Vision

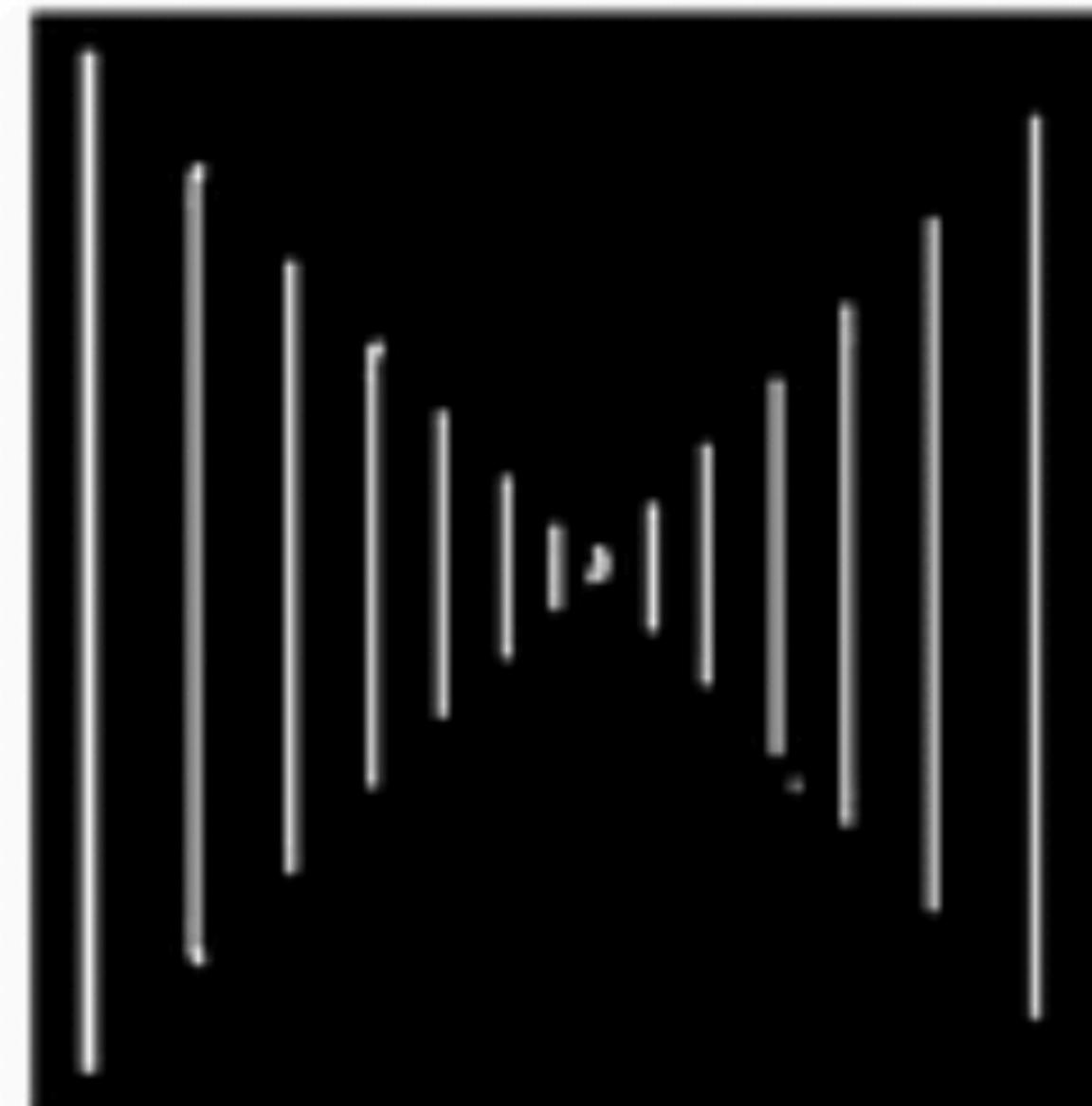
- Natural images are *not* arbitrary 2D arrays
- They have properties resulting from physics / math of image formation
- Convert from “pixels” to “objects”: which groups of pixels correspond to objects?



<https://www.cs.cornell.edu/courses/cs4670/2018sp/lec02-imageformn.pdf>

the content are taken from here and there

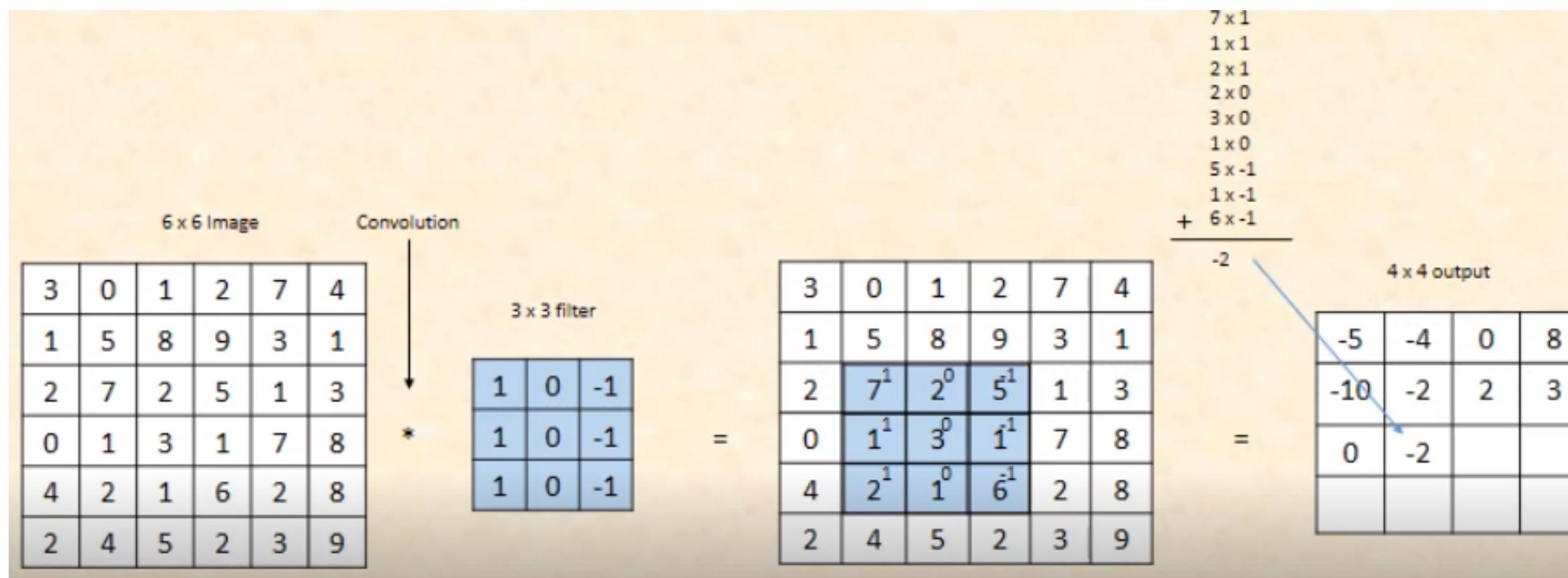
Computer Vision

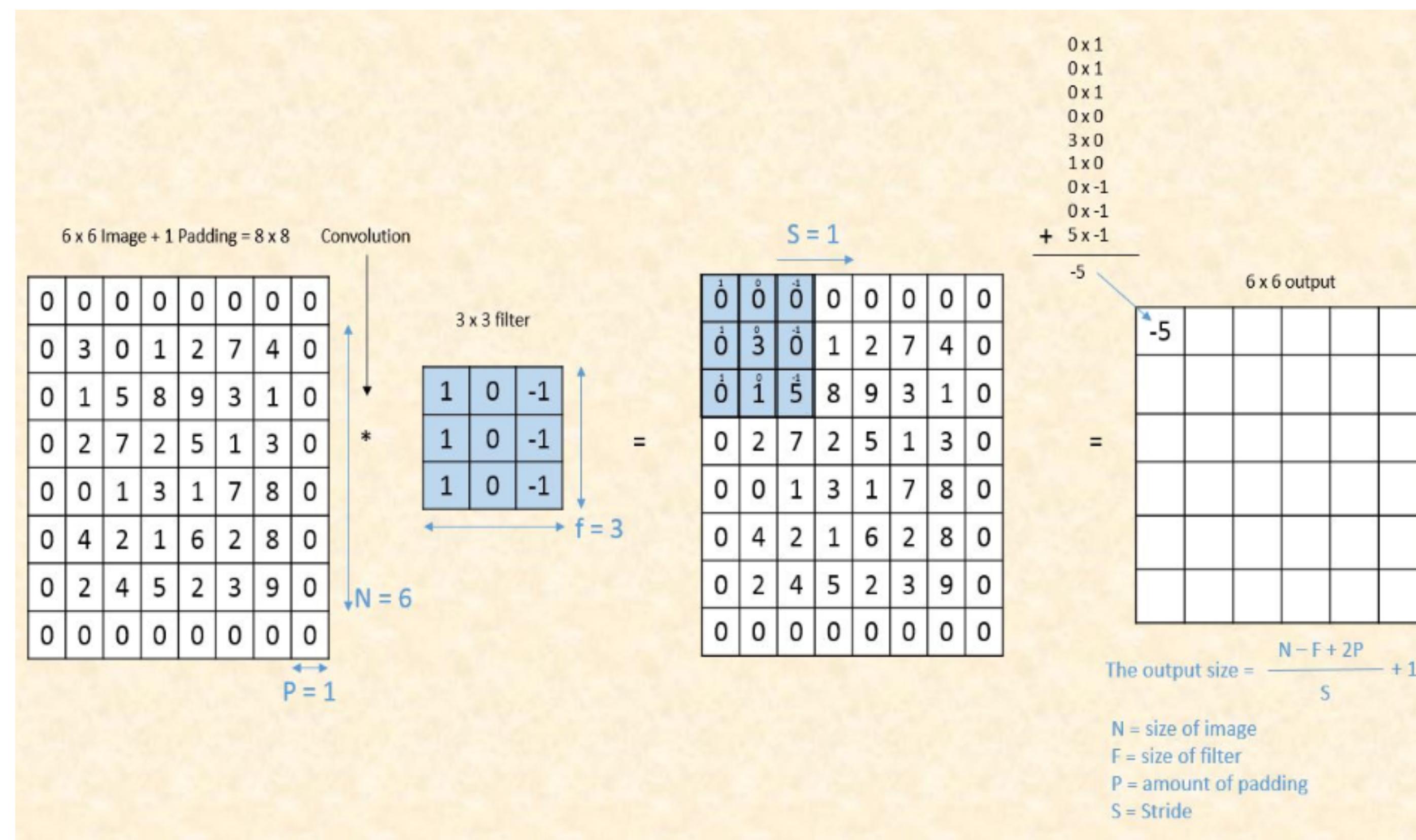


Mathematics

Convolution operation -

- convolving a 6×6 grayscale image with a 3×3 matrix called filter or kernel to produce a 4×4 matrix.
- dot product between the filter and the first 9 elements of the image matrix and fill the output matrix.





the content are taken from here and there

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

: :

Each filter detects a small pattern (3 x 3).

Convolution

stride=1

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

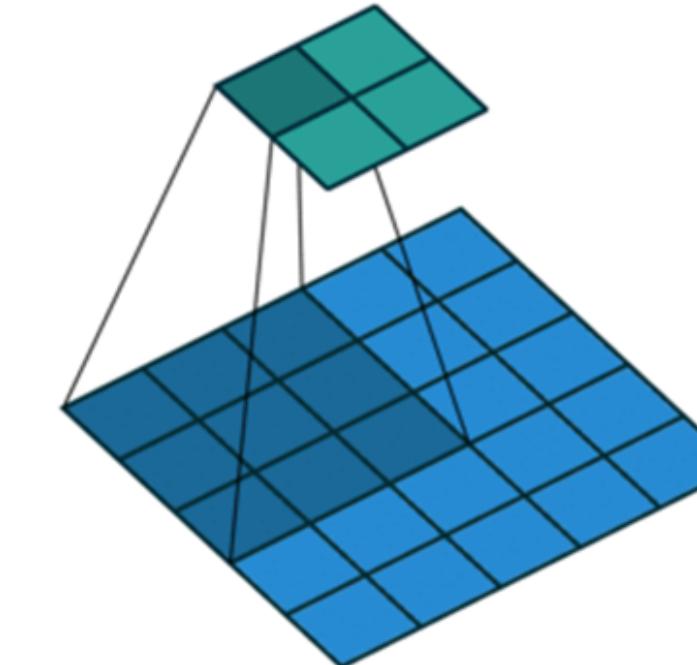
1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

Dot
product



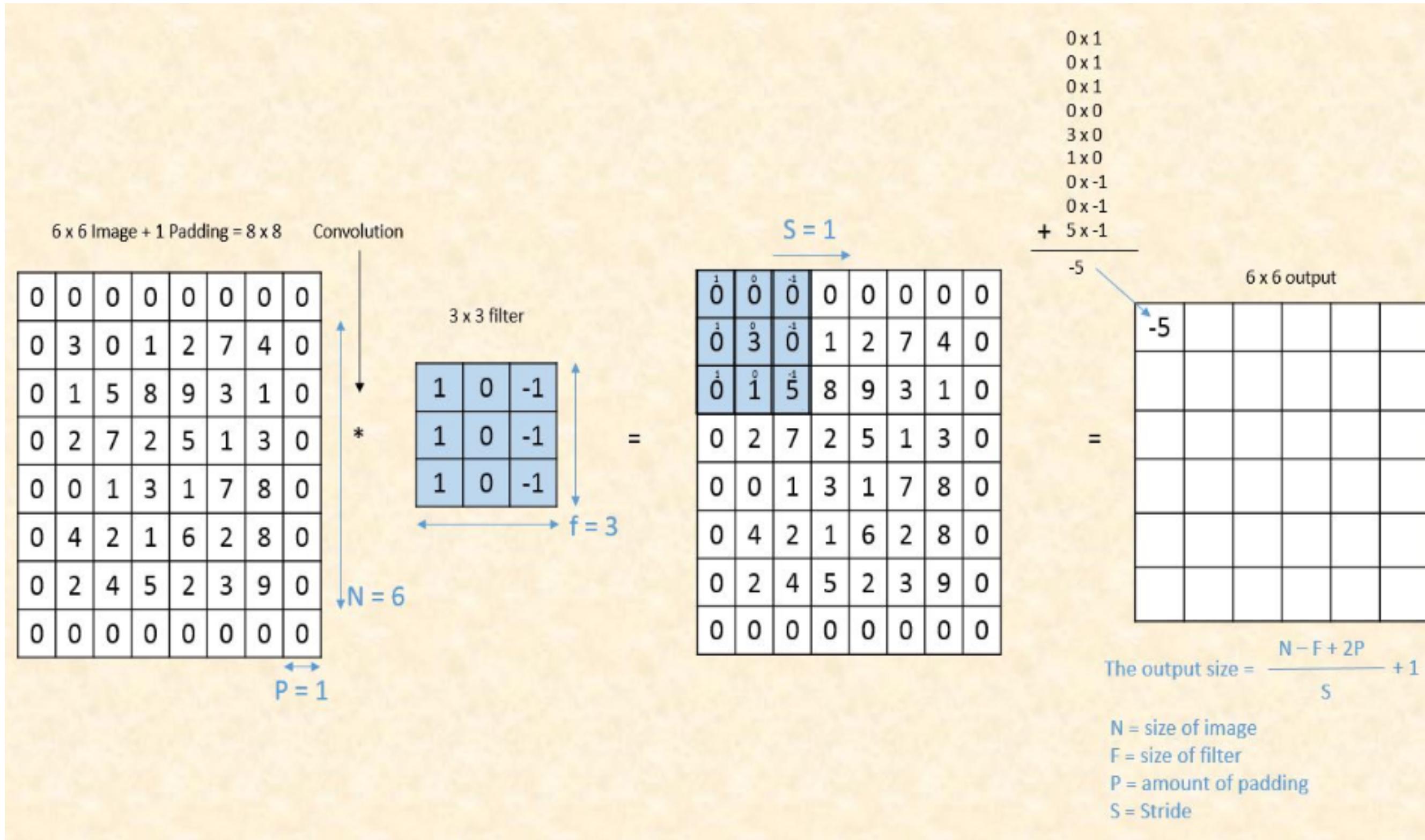
$$\begin{aligned} & 1*1 + 0*(-1) + 0*(-1) + \\ & 0*(-1) + 1*1 + 0*(-1) + \\ & 0*(-1) + 0*(-1) + 1*1 \end{aligned}$$



Challenges with convolution operation:

- Output may shrink.
- Data loss from image corners.

Solution - Padding (size of the output image is equal to input image size.)

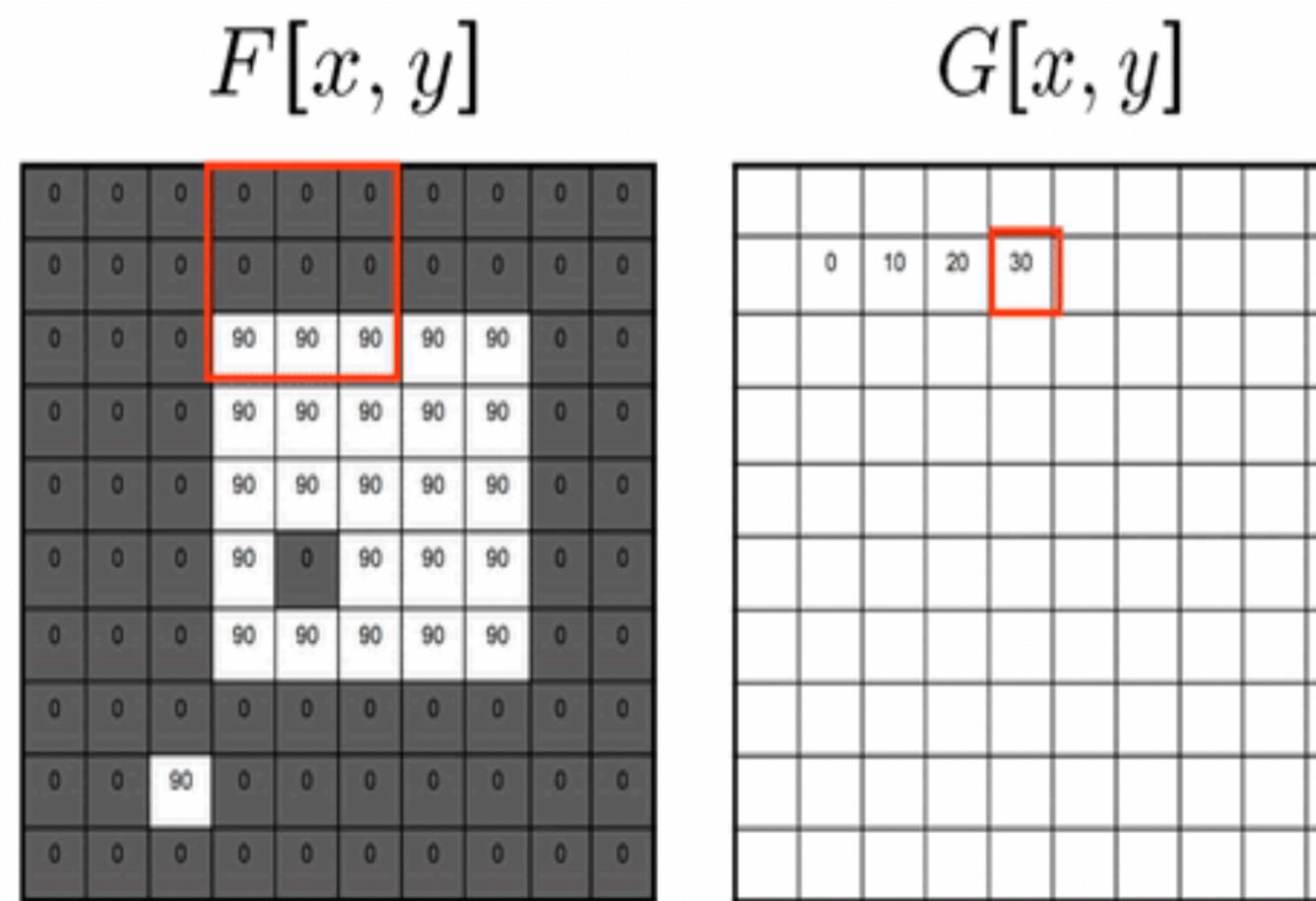


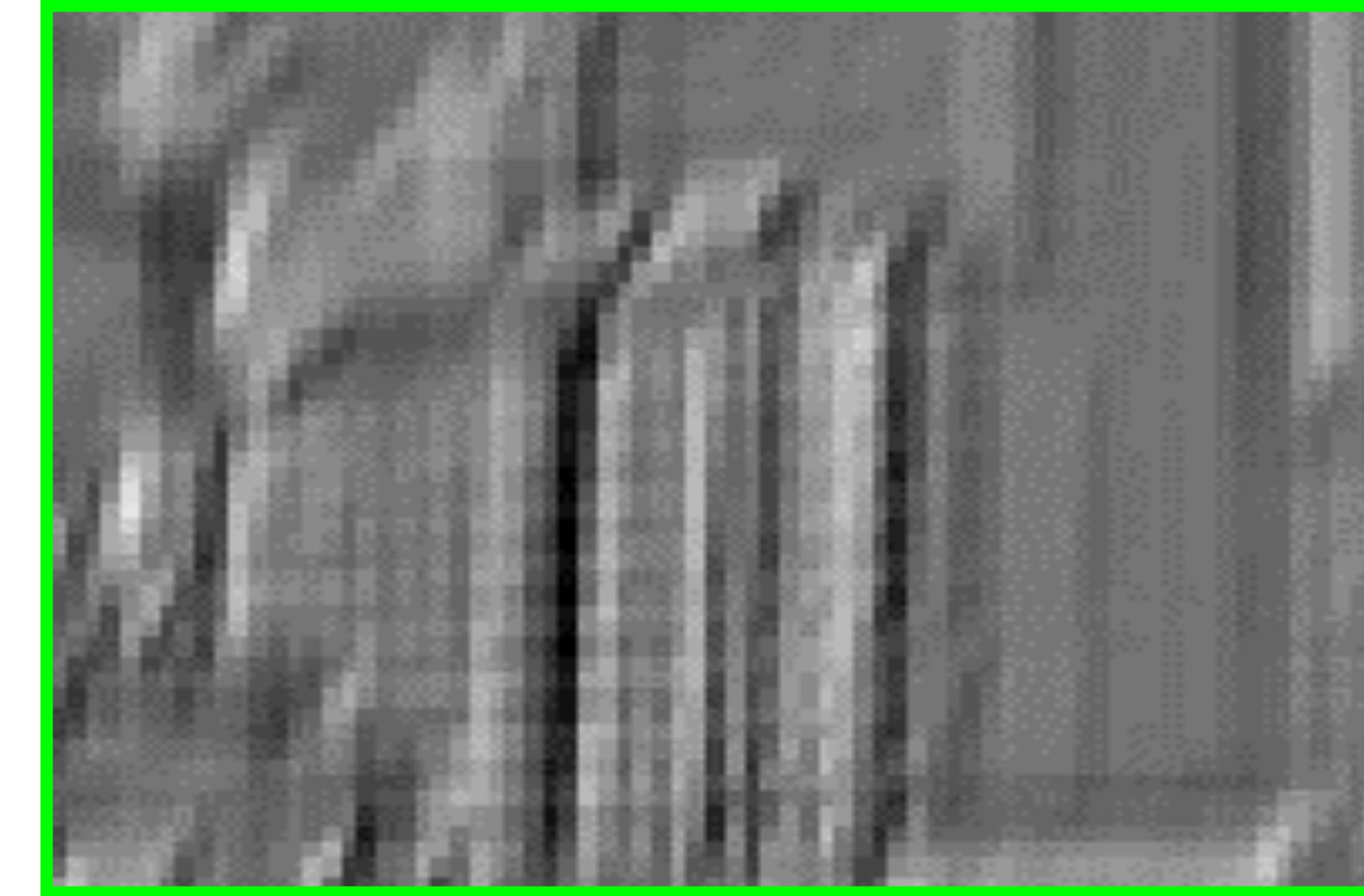
the content are taken from here and there

Computer Vision

1. Filtering and Convolution

Filtering operations such as blurring, sharpening and noise reduction are applied to images using convolution. Convolution involves sliding a filter or kernel over an image and performing mathematical operations on each pixel. This process enables various improvements such as anti-aliasing, edge detection and texture removal.





- Weighted moving sum generating features

Biological Vision

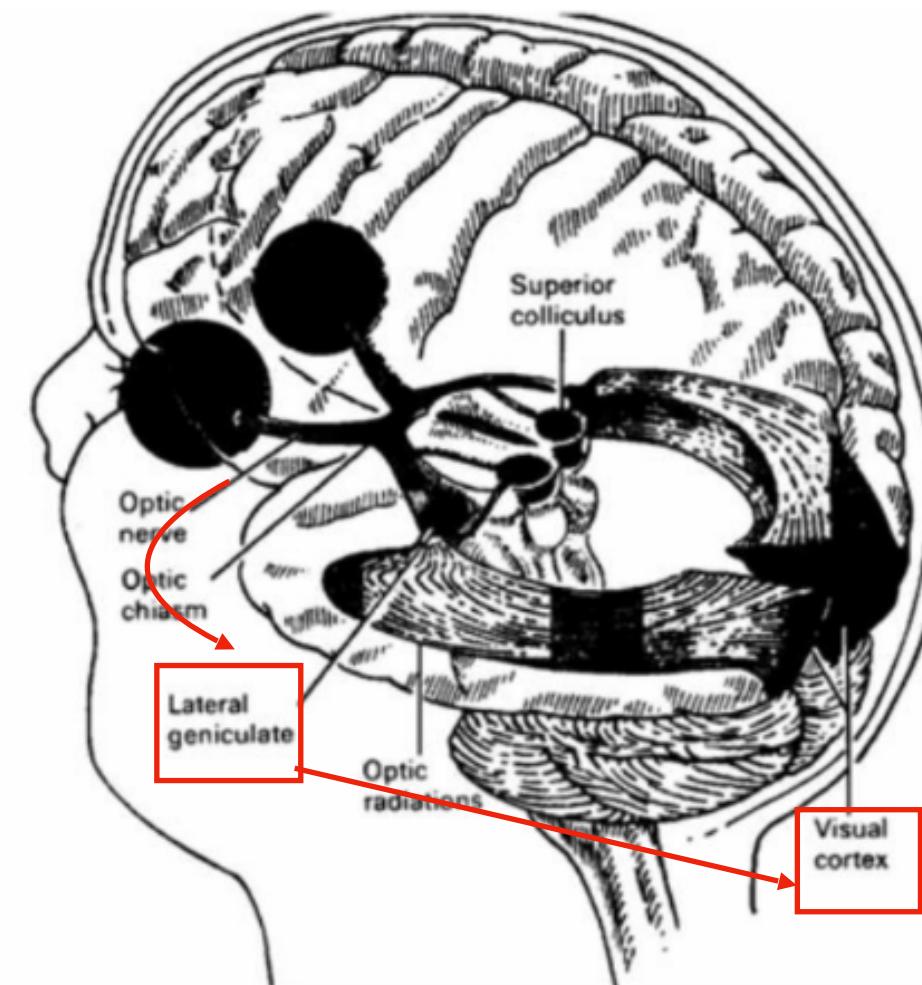
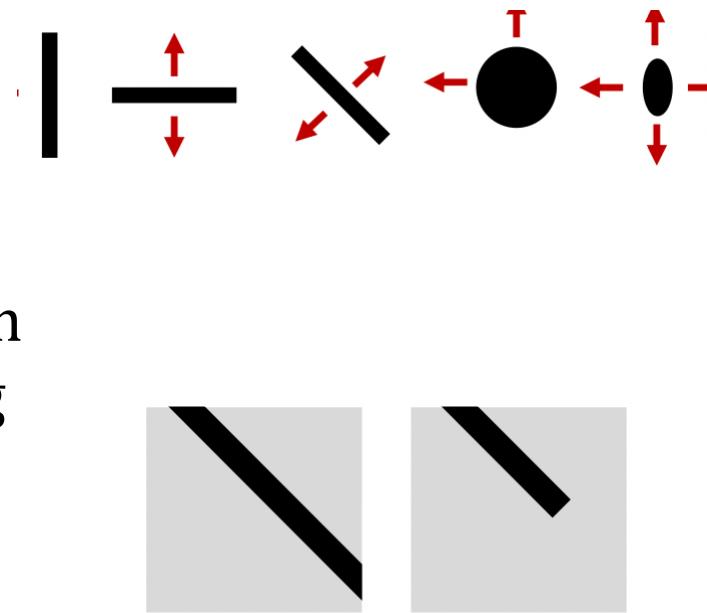
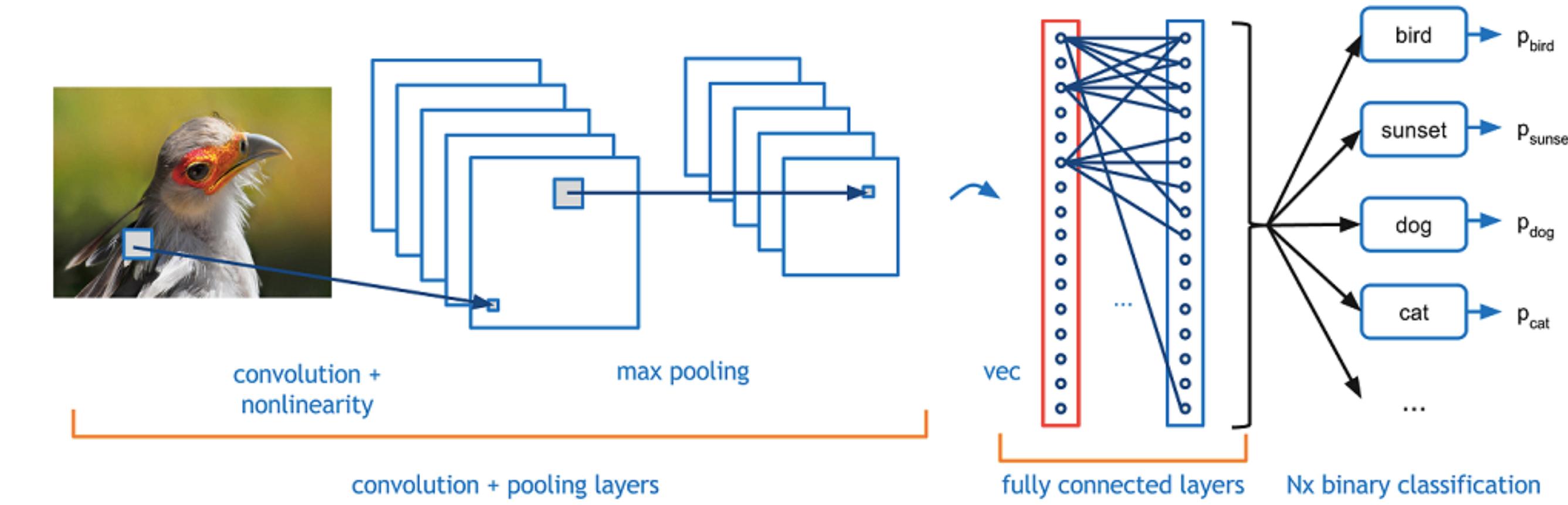
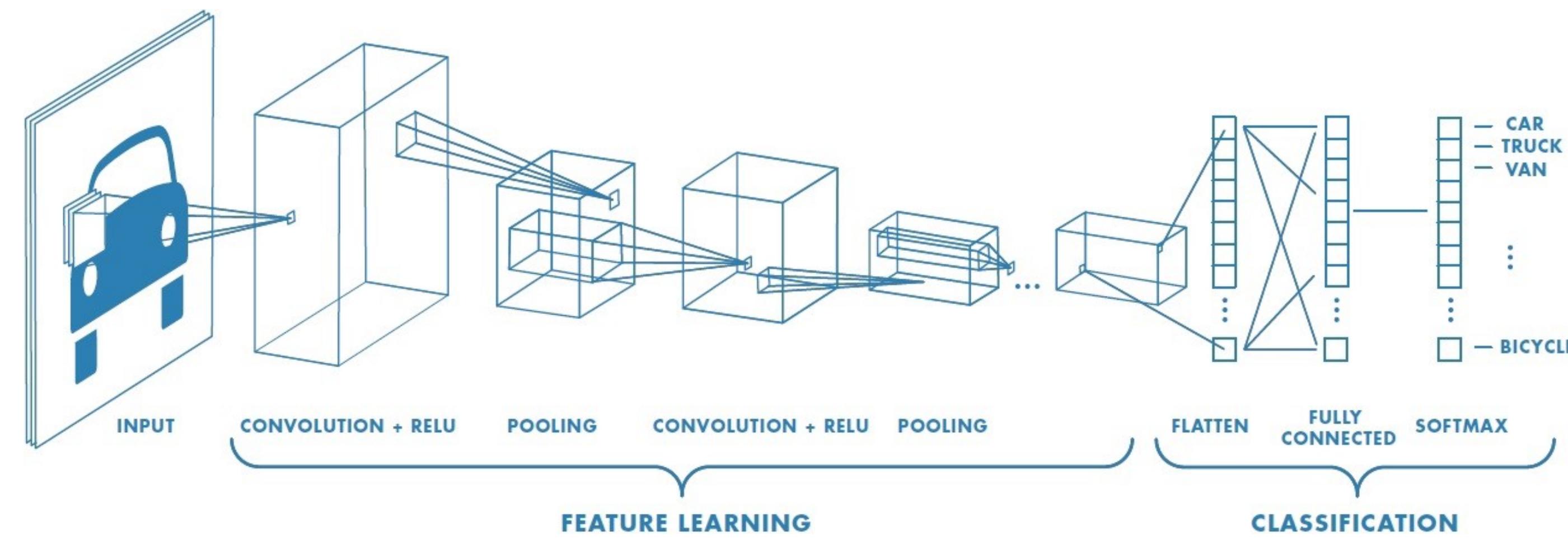


Image Credit: The Three R's of Computer Vision,
Jitendra Malik
UC Berkeley

simple visual forms,
edges, corners
Pathway
of
Information
Processing
high-level object
descriptions, faces, objects

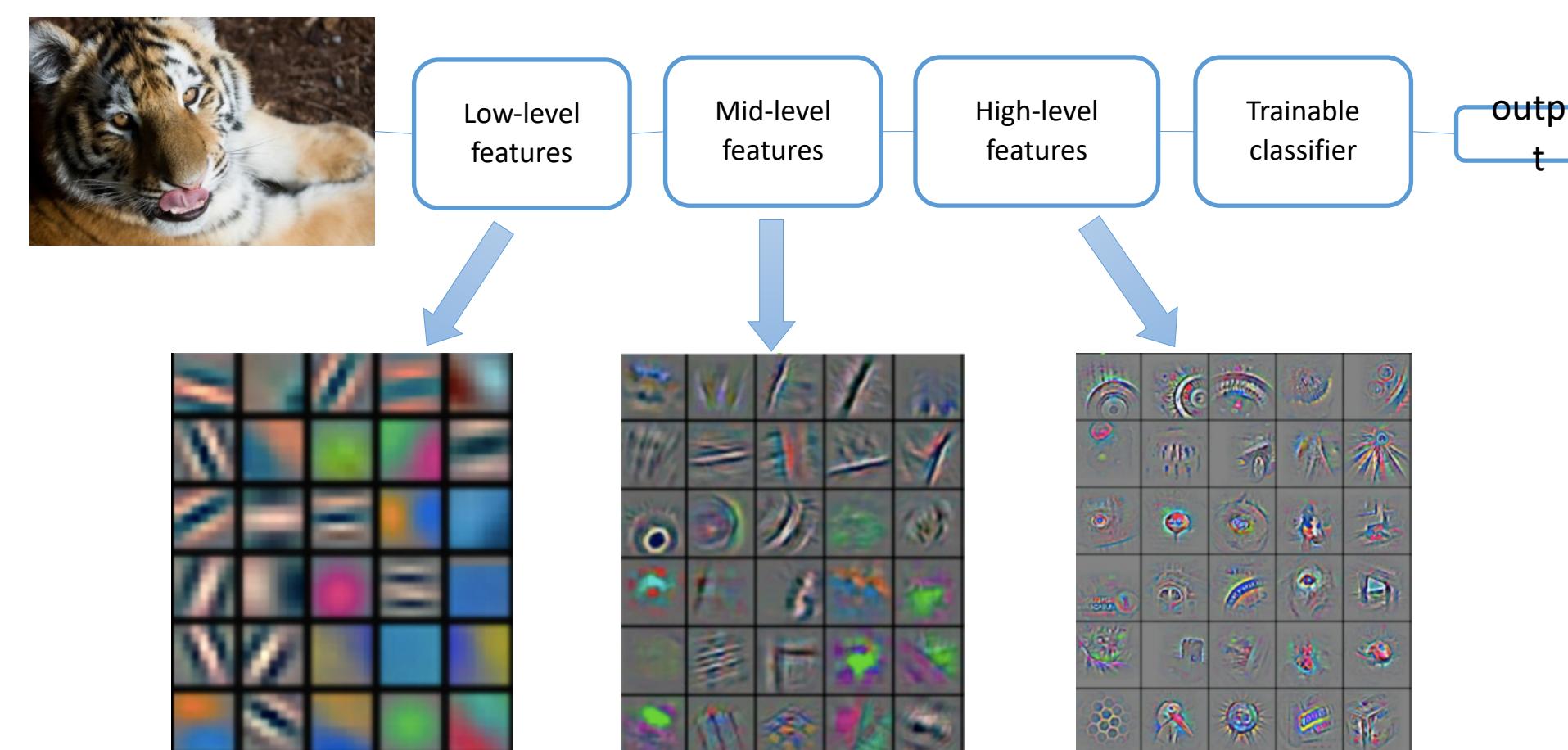




the content are taken from here and there

DEEP LEARNING

- Deep learning (a.k.a. representation learning) seeks to learn rich hierarchical representations (i.e. features) automatically through multiple stage of feature learning process.



Convolutional Neural Networks

- Each layer in a CNN applies a different set of filters, typically hundreds or thousands of them, and combines the results, feeding the output into the next layer in the network.
- During training, a CNN automatically learns the values for these filters.
- In the context of image classification, our CNN may learn to:
 - ✓ Detect edges from raw pixel data in the first layer.
 - ✓ Use these edges to detect shapes (i.e., “blobs”) in the second layer.
 - ✓ Use these shapes to detect higher-level features such as facial structures, parts of a car, etc. in the highest layers of the network.
- The last layer in a CNN uses these higher-level features to make predictions regarding the contents of the image.

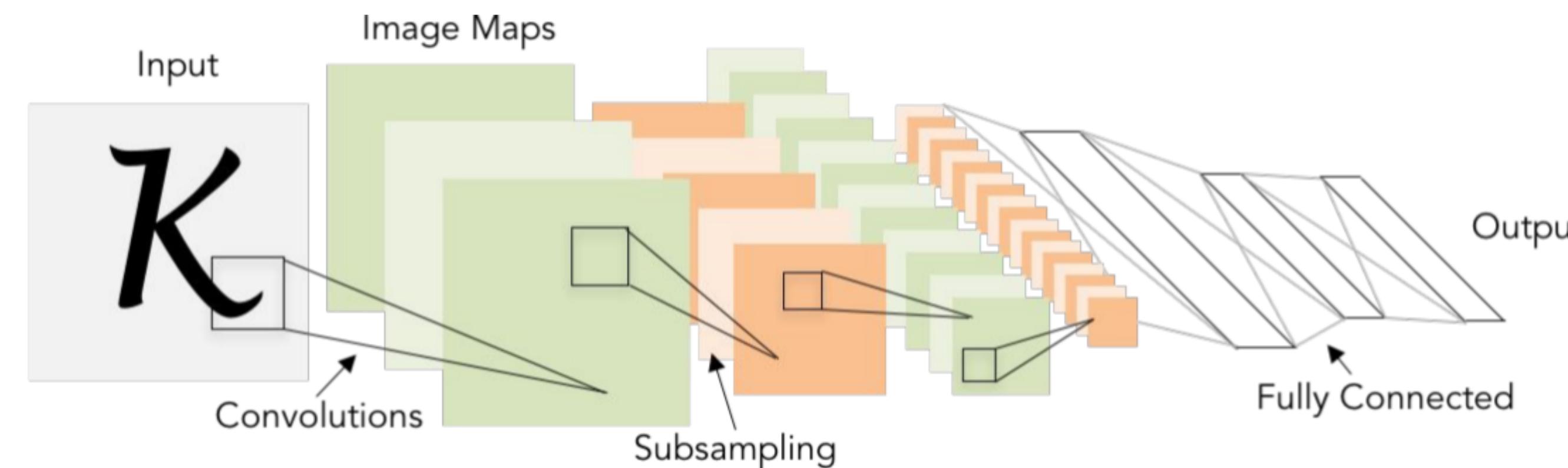
Convolutional Neural Networks

CNN consist of three layers -

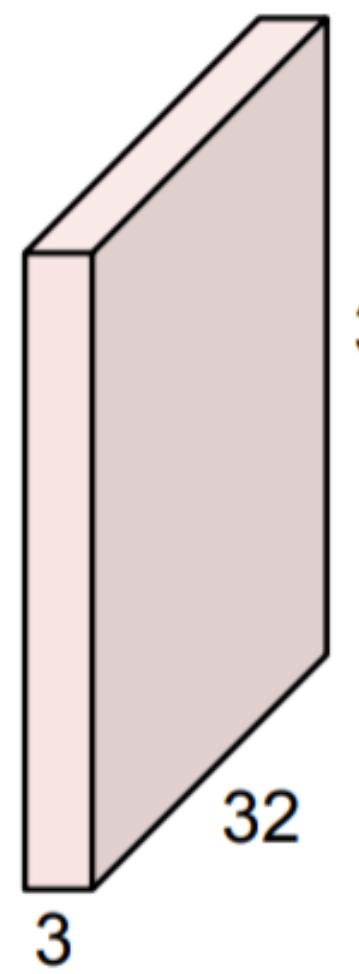
1. Convolution Layer,
2. Pooling Layer, and
3. Fully Connected Layer.

ConvNet architecture is formed stacking these layers.

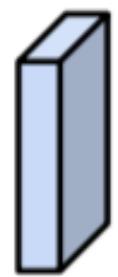
- The output can be a softmax layer indicating whether there is a cat or something else. You can also have a sigmoid layer to give you a probability of the image being a cat.



32x32x3 image

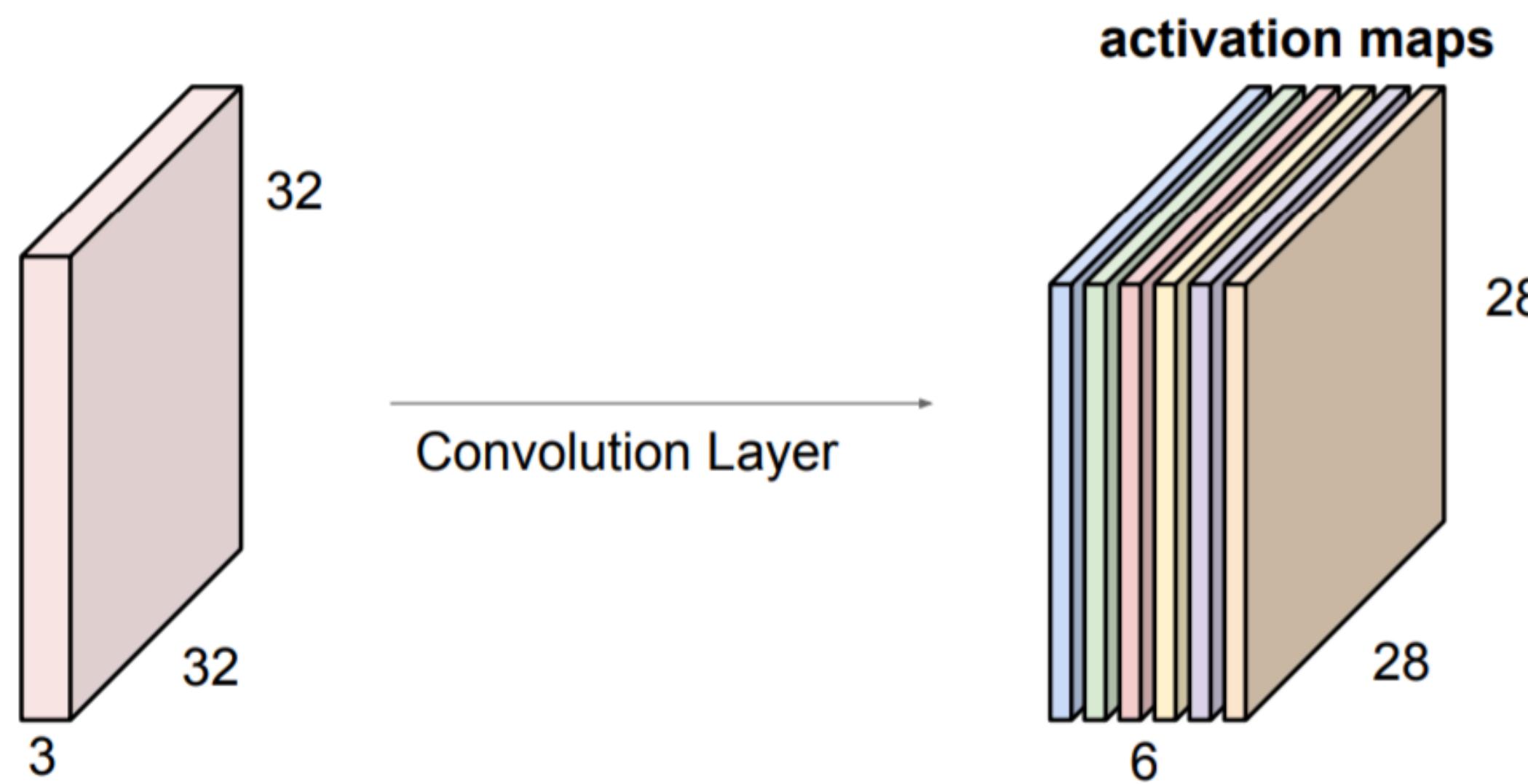


5x5x3 filter

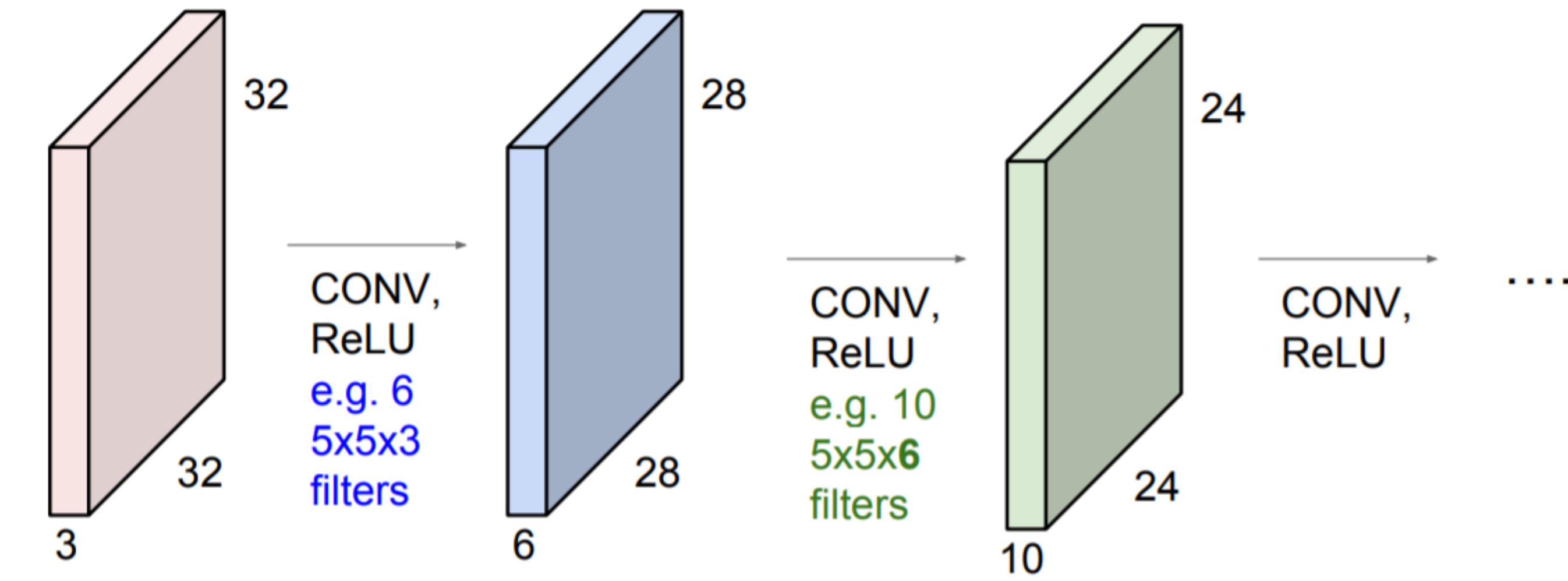


Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

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



Ex: we have $6 -> 5 \times 5$ filters, we'll get 6 separate activation maps. “new image” of size $28 \times 28 \times 6$!



the content are taken from here and there

Convolution over volume

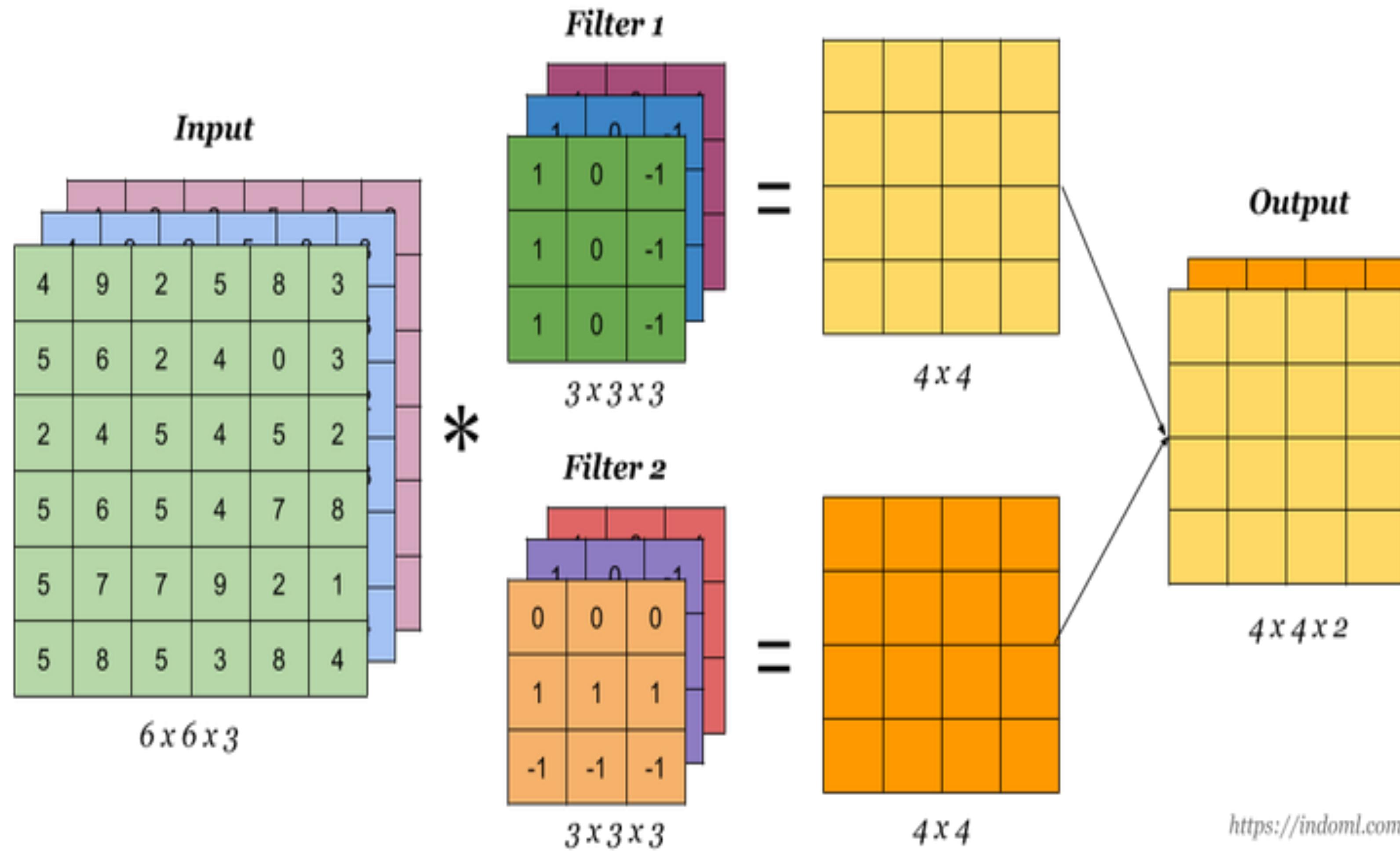


depth of the filter = depth of image
(necessary condition)

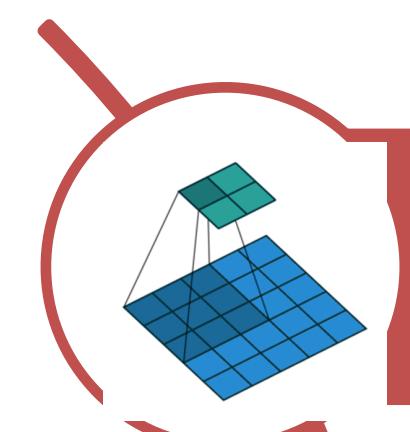
- For every filter one feature map is created
- Depth of the output image represents the number of filters

With every increasing layer size of image decreases and number of filter increases

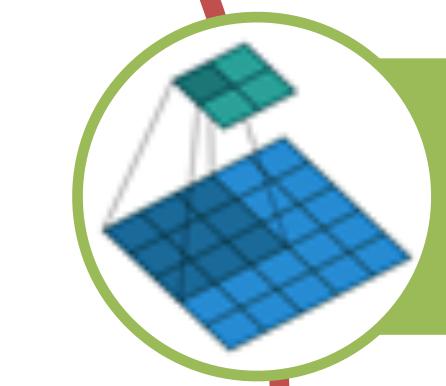
Convolution over volume



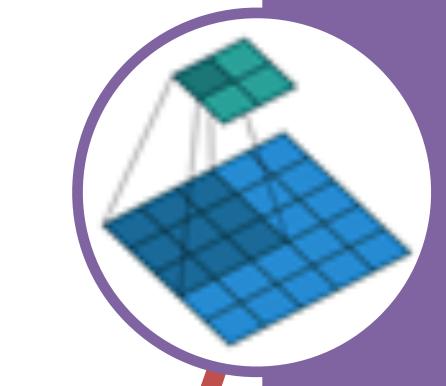
Strides



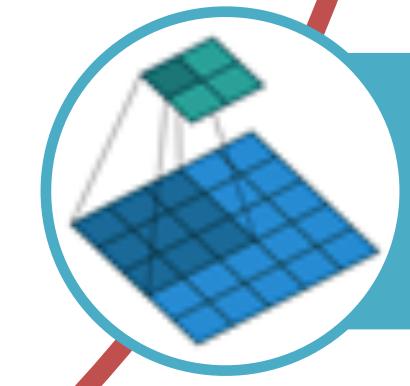
Used to reduce Dimensionality



Decides the shift along axis

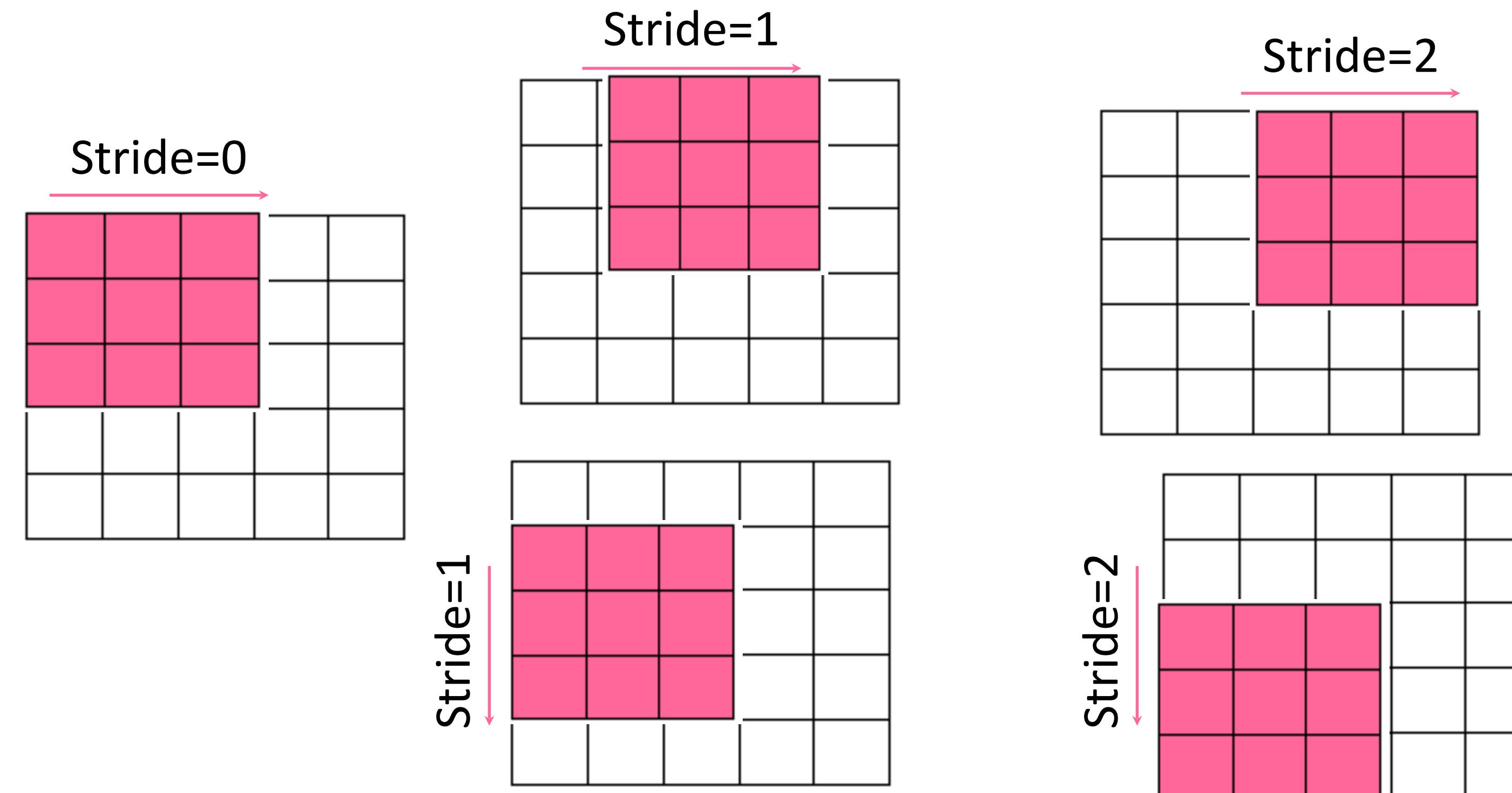


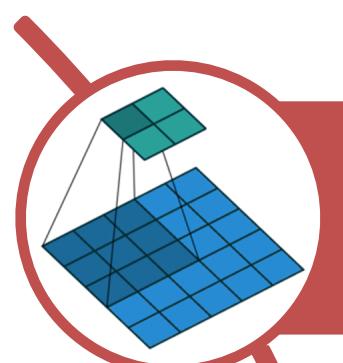
Smaller strides identifies the pattern better



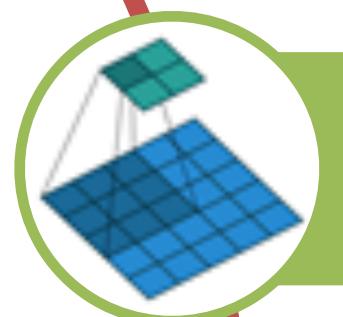
Used for pattern recognition

Stride- Visualisation

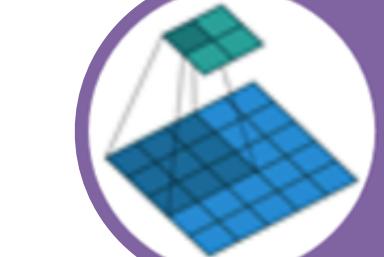




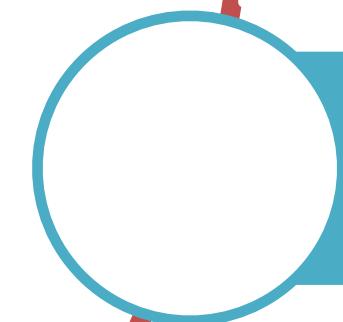
Used to avoid reduction in size



Retaining the original size



Protect the information along
the border

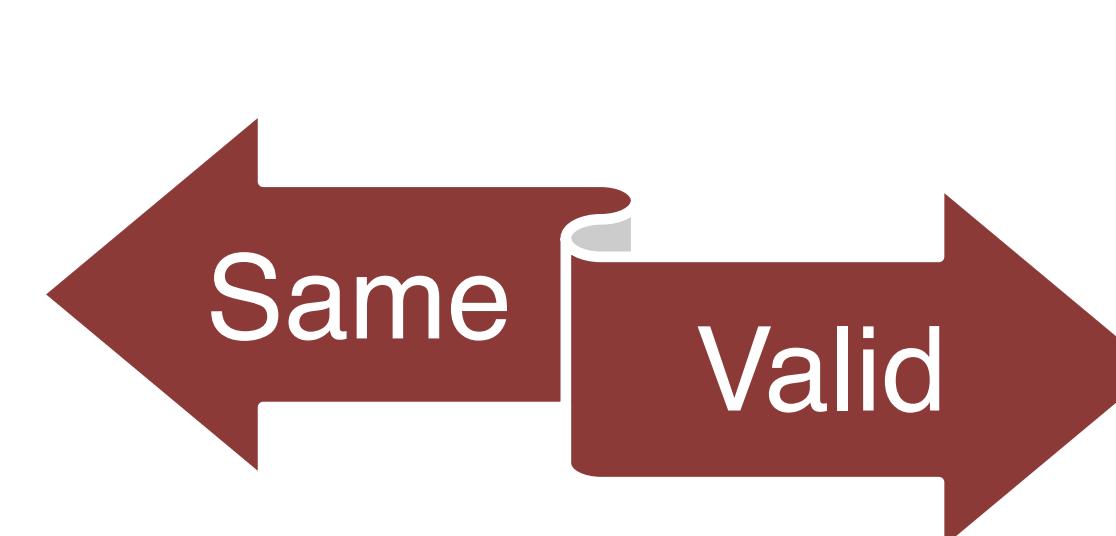


Input : image ($h * w * d$), filter($f_h * f_w * d$)

Types of padding

Same

after padding the size of image is unaffected.



3	5	9	1	10
13	2	4	6	11
16	24	9	13	1
7	1	6	8	3
8	4	9	1	9

<https://images.app.goo.gl/WtzAyvDBrRbzSZ3P7>

Valid

no padding done.

0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	3	5	9	1	10	0	0
0	0	13	2	4	6	11	0	0
0	0	16	24	9	13	1	0	0
0	0	7	1	6	8	3	0	0
0	0	8	4	9	1	9	0	0
0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0

$$\begin{matrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 3 & 3 & 4 & 4 & 7 & 0 & 0 & 0 \\ 0 & 9 & 7 & 6 & 5 & 8 & 2 & 0 & 0 \\ 0 & 6 & 5 & 5 & 6 & 9 & 2 & 0 & 0 \\ 0 & 7 & 1 & 3 & 2 & 7 & 8 & 0 & 0 \\ 0 & 0 & 3 & 7 & 1 & 8 & 3 & 0 & 0 \\ 0 & 4 & 0 & 4 & 3 & 2 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{matrix} * \begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix} = \begin{matrix} -10 & -13 & 1 & & & \\ -9 & 3 & 0 & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{matrix}$$

$6 \times 6 \rightarrow 3 \times 3$

<https://images.app.goo.gl/UPri8nVvzJcrUxD9>

Padding Numerical

- If the image matrix is 6×6
- If the size of the filter is $3 \times 3 \times 2$,
- If stride is 1 and padding is 1
- The output is calculated as

$$\left(\frac{h + (2 * p) - fh}{s} + 1 \right) * \left(\frac{w + (2 * p) - fw}{s} + 1 \right) * d$$

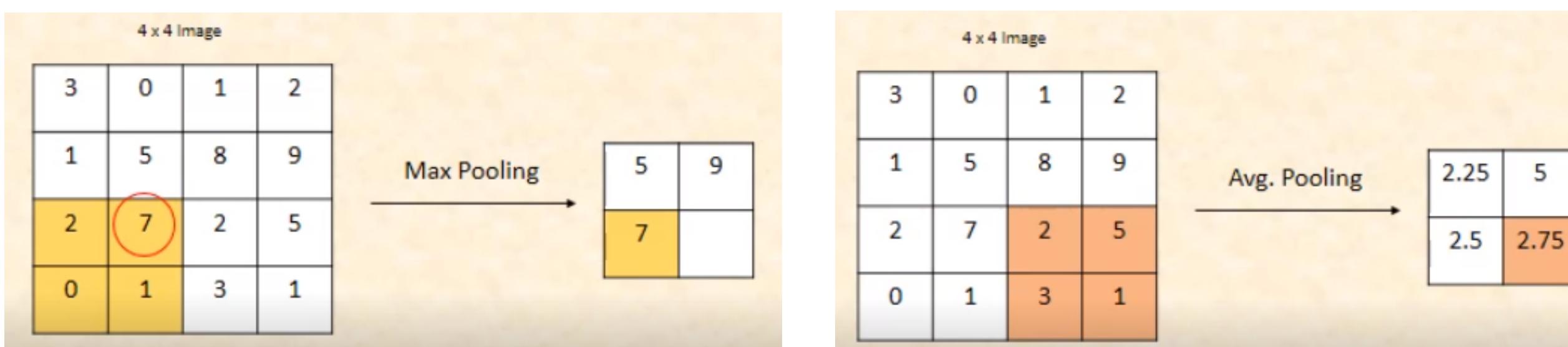
- So, output image is

$$\left(\frac{6 + (2 * 1) - 3}{1} + 1 \right) * \left(\frac{6 + (2 * 1) - 3}{1} + 1 \right) * 2$$

$$= 6 * 6$$

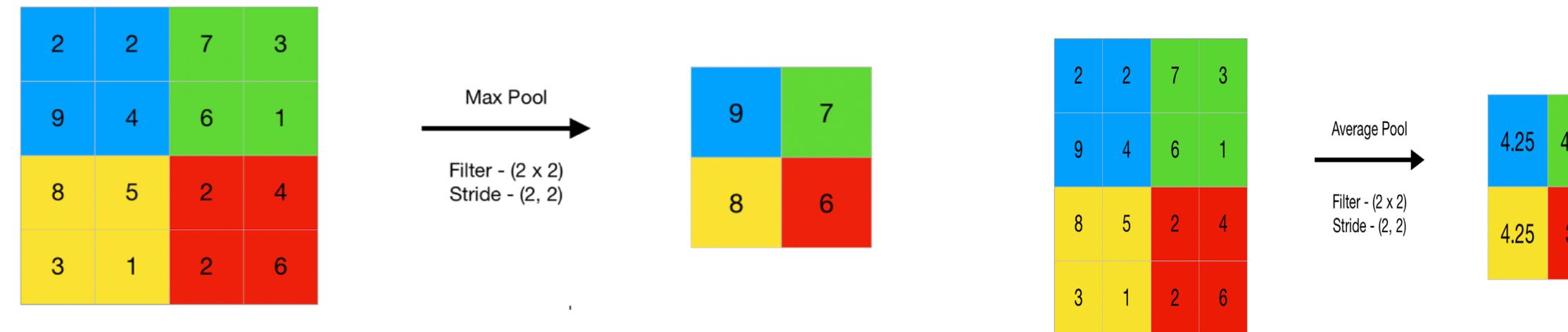
Pooling Layer:

- Reduce number of parameters.
- In pooling layer, we have two hyperparameters filter size and stride which are fixed only once.



the content are taken from here and there

- Pooling layers are used to reduce the dimensions of the feature maps. Reduces the number of parameters to learn and the amount of computation performed in the network.
- The pooling layer summarises the features present in a region of the feature map generated by a convolution layer. So, further operations are performed on summarised features instead of precisely positioned features generated by the convolution layer. This makes the model more robust to variations in the position of the features in the input image.
- Average pooling computes the average of the elements present in the region of feature map covered by the filter. Thus, while max pooling gives the most prominent feature in a particular patch of the feature map, average pooling gives the average of features present in a patch.



the content are taken from here and there

Why Pooling

- Subsampling pixels will not change the object

bird



Subsampling

bird

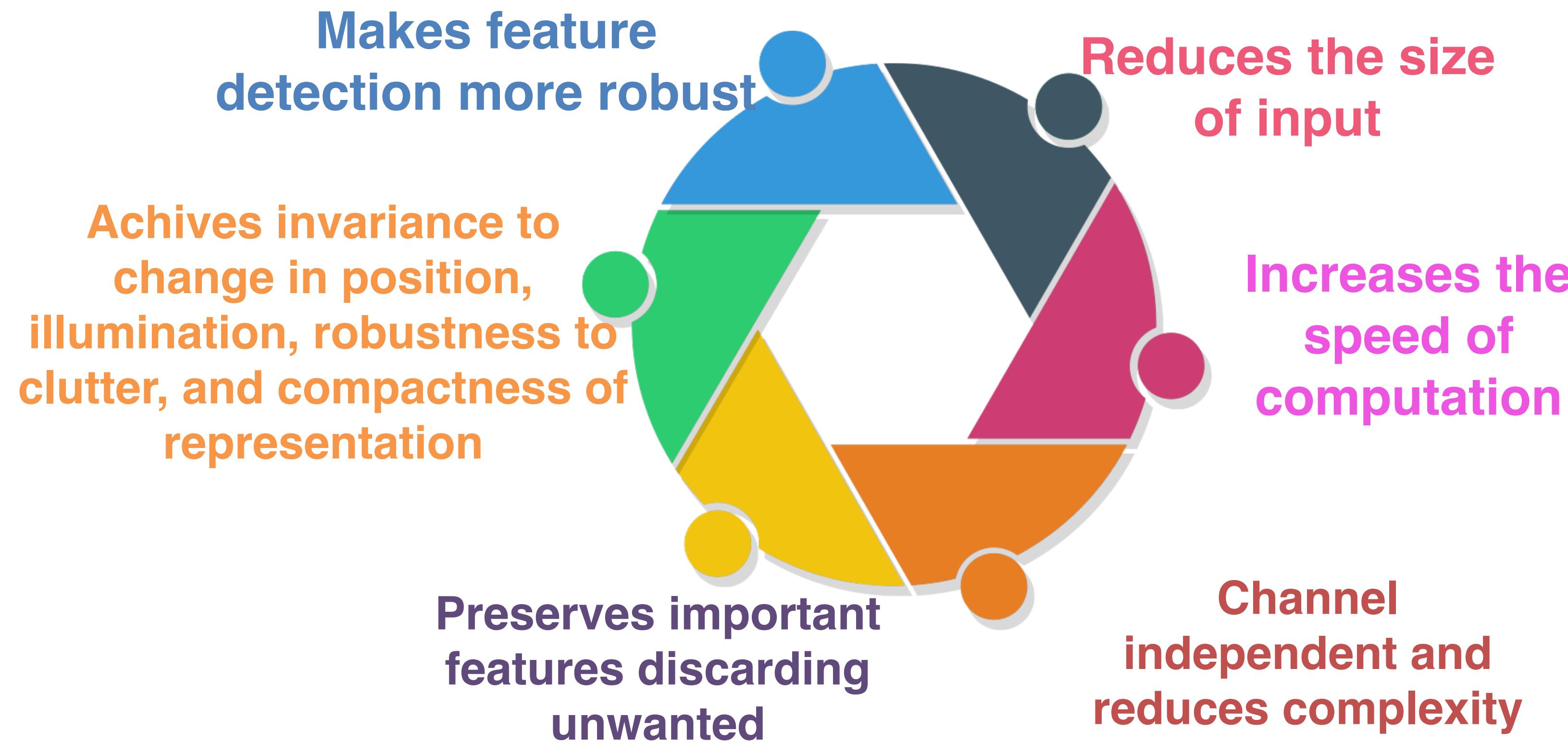


We can subsample the pixels to make image smaller



fewer parameters to characterize the image

Importance of pooling



Types of Pooling

Max

Take the max value
in each block

Max

Average

Average

average all values in
each block

Max Pooling

29	15	28	184
0	100	70	38
12	12	7	2
12	12	45	6

2 x 2
pool size

100	184
12	45

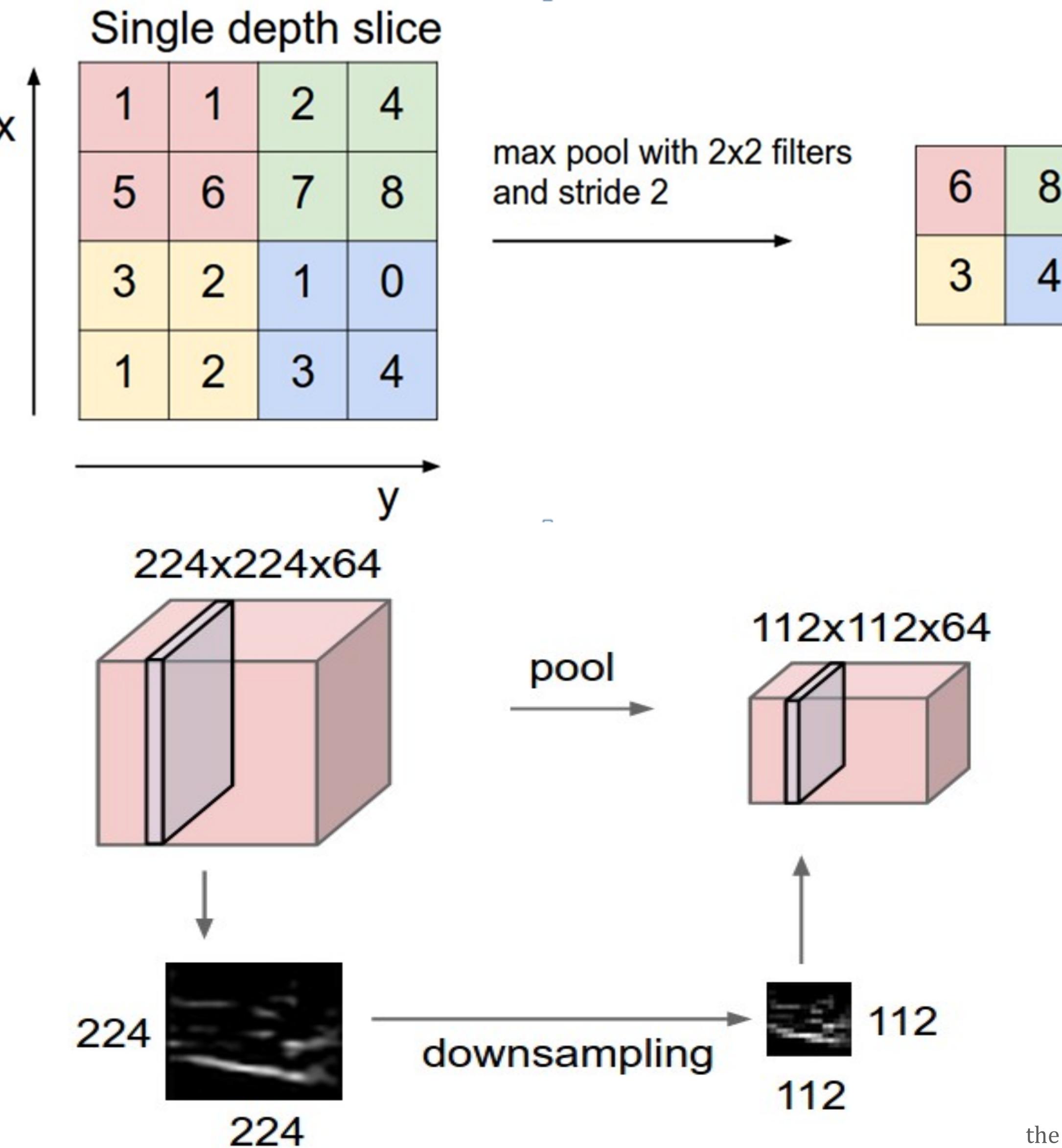
Average Pooling

31	15	28	184
0	100	70	38
12	12	7	2
12	12	45	6

2 x 2
pool size

36	80
12	15

Visualisation of Pooling

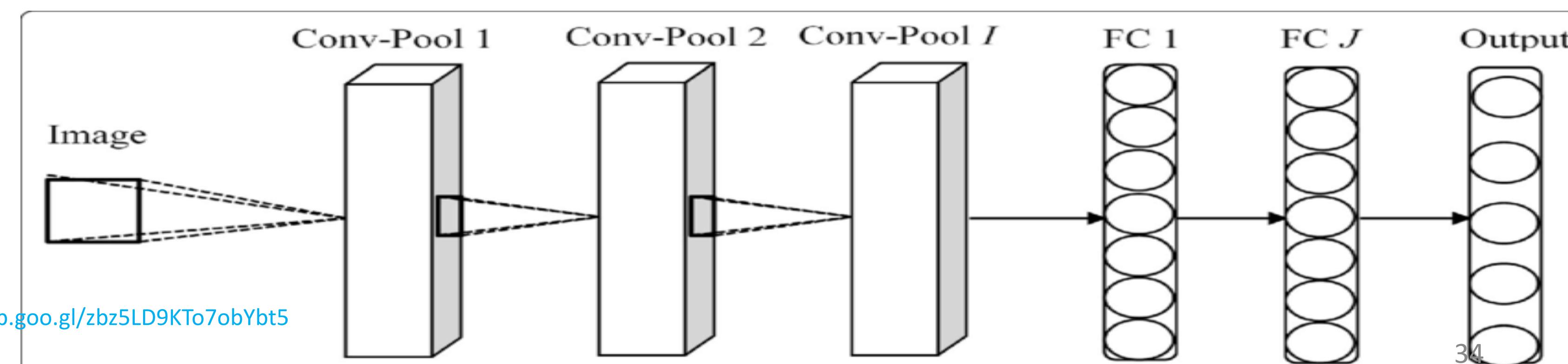


Flattening or FC layer

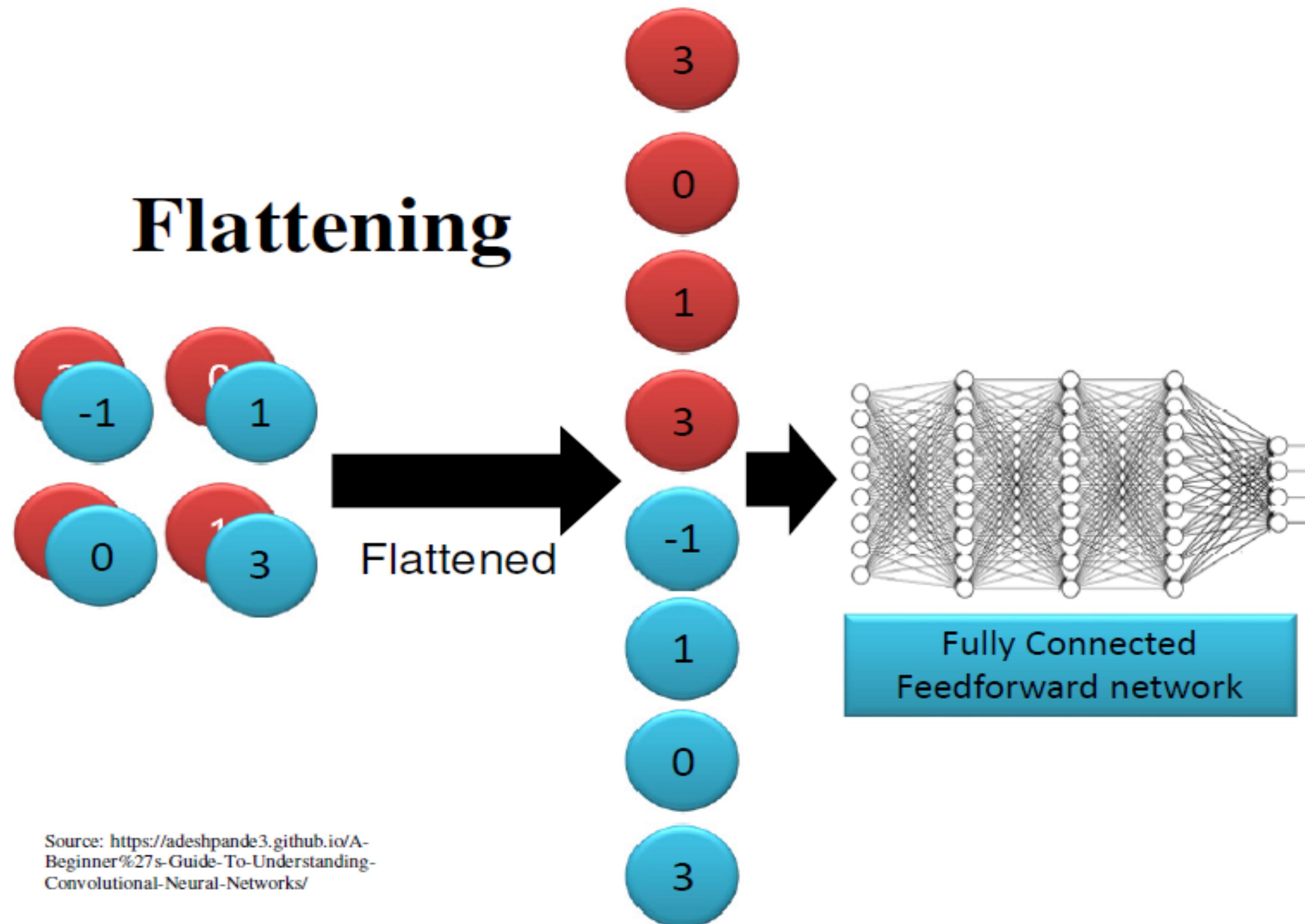


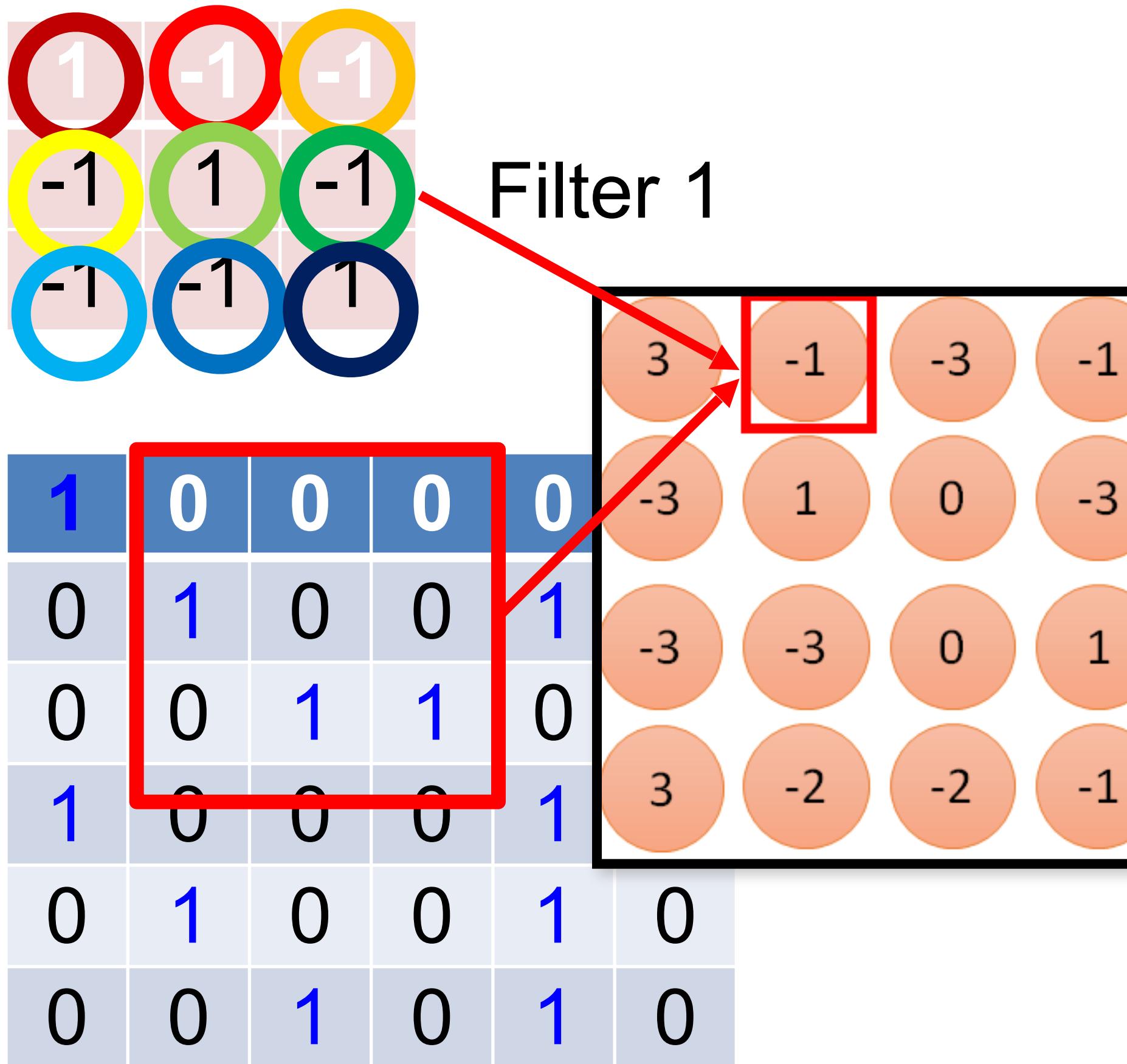
Hidden layers are called classification layer

Number of neurons in the output layer is equal to the number of classes.



Flattening or FC layer

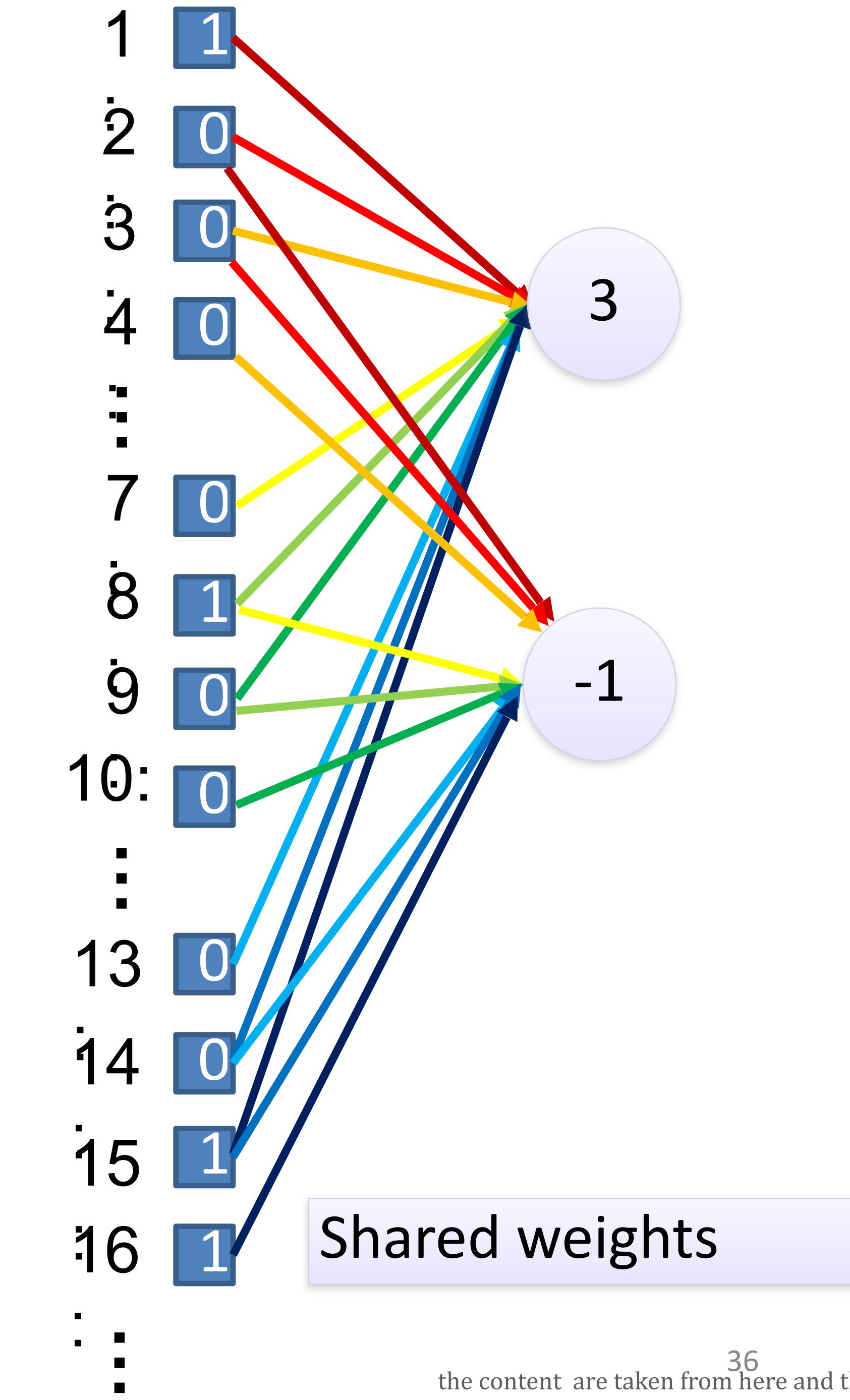




6 x 6 image

Fewer parameters

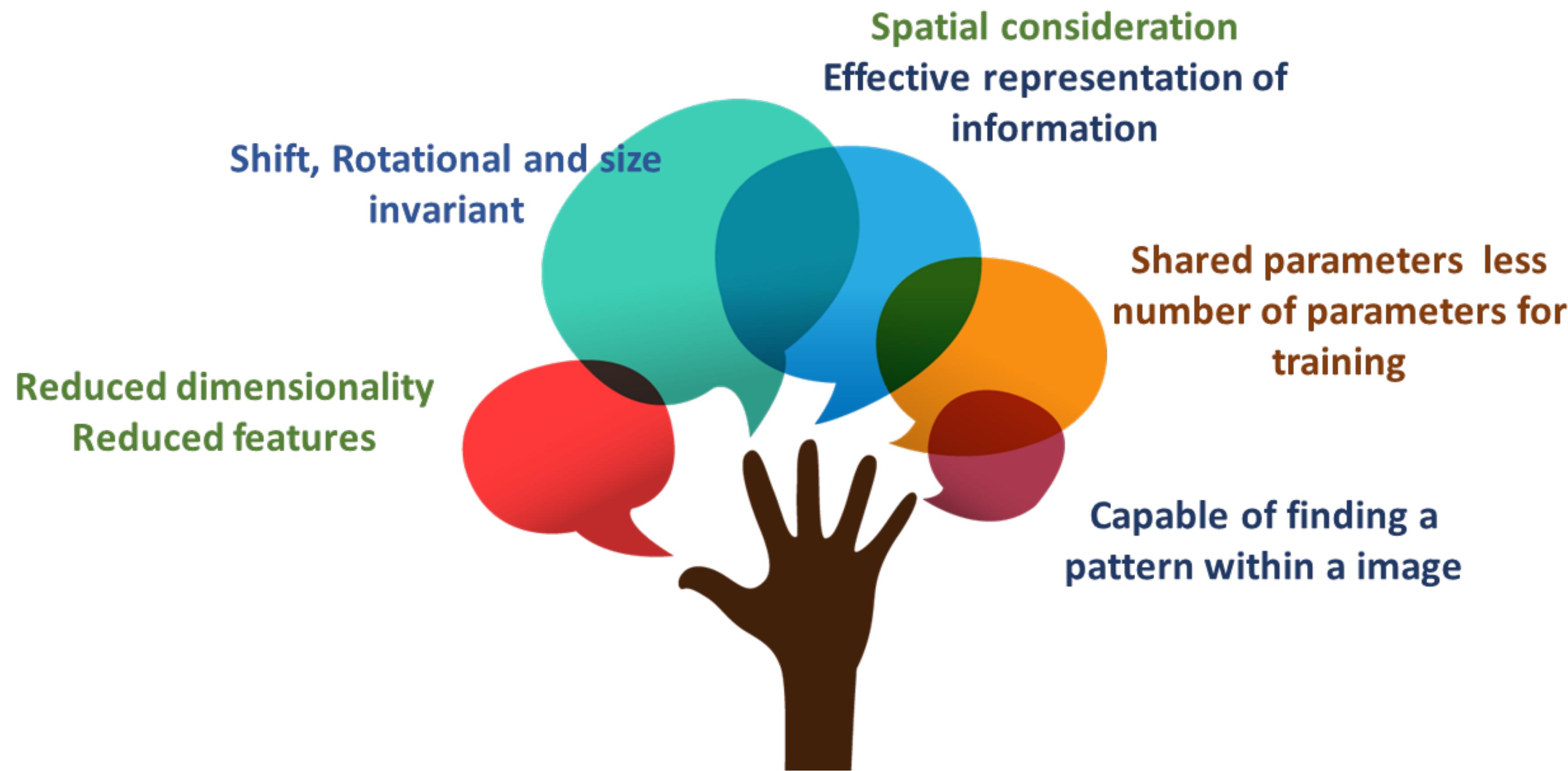
Even fewer parameters



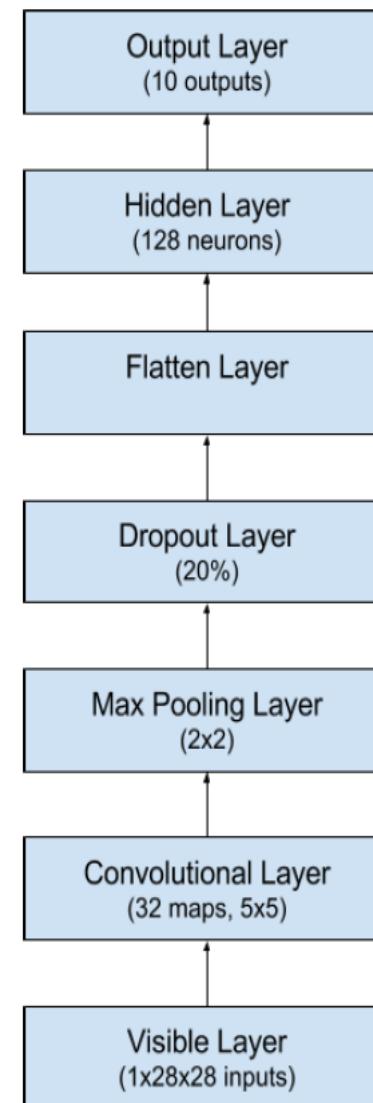
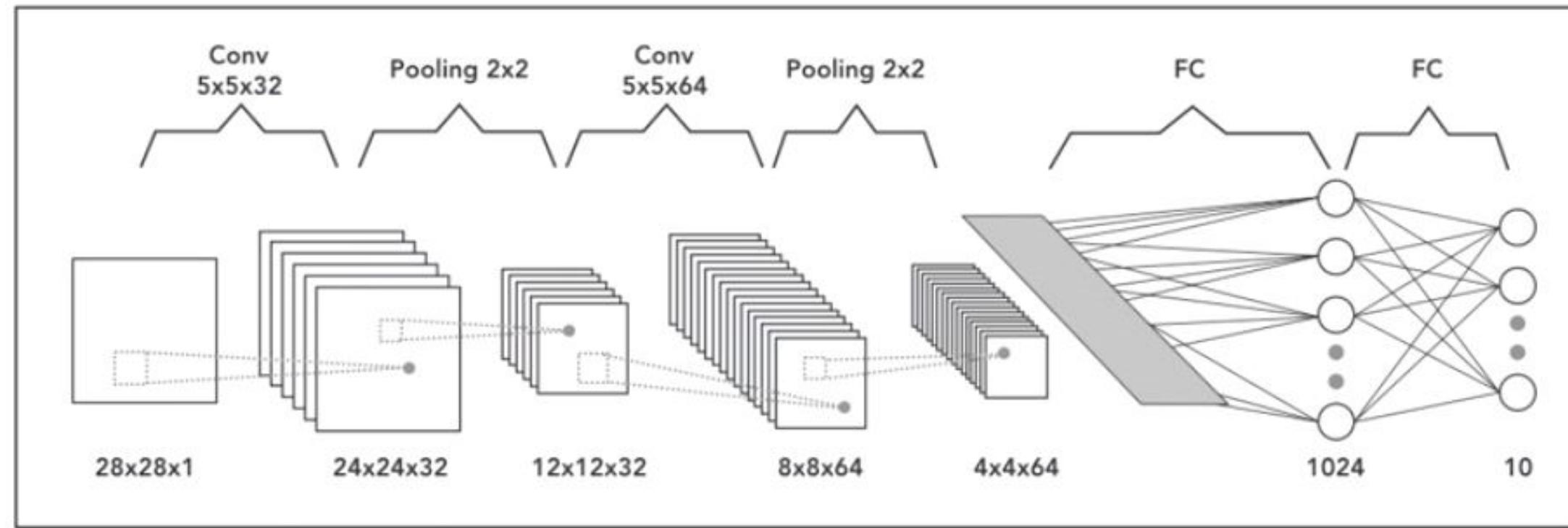
CNN Overview

- Better than fully connected network as it considers the spatial aspect .
- Use shared weights to reduce the number of parameters and hence number of trainable parameters are reasonable
- Represent a small region with fewer parameters hence reducing dimensionality
- The kernels (shared weights) are employed for learning the patterns, giving importance to only those features that are important
- The patterns in the input data will be further refined
- Little or no invariance to shifting, scaling, and other forms of distortion
- Network is capable of finding a object/ feature in any part of the image.

Convolution Layer



Architecture:



CNN Architecture for MNIST

the content are taken from here and there

Summary of CNN

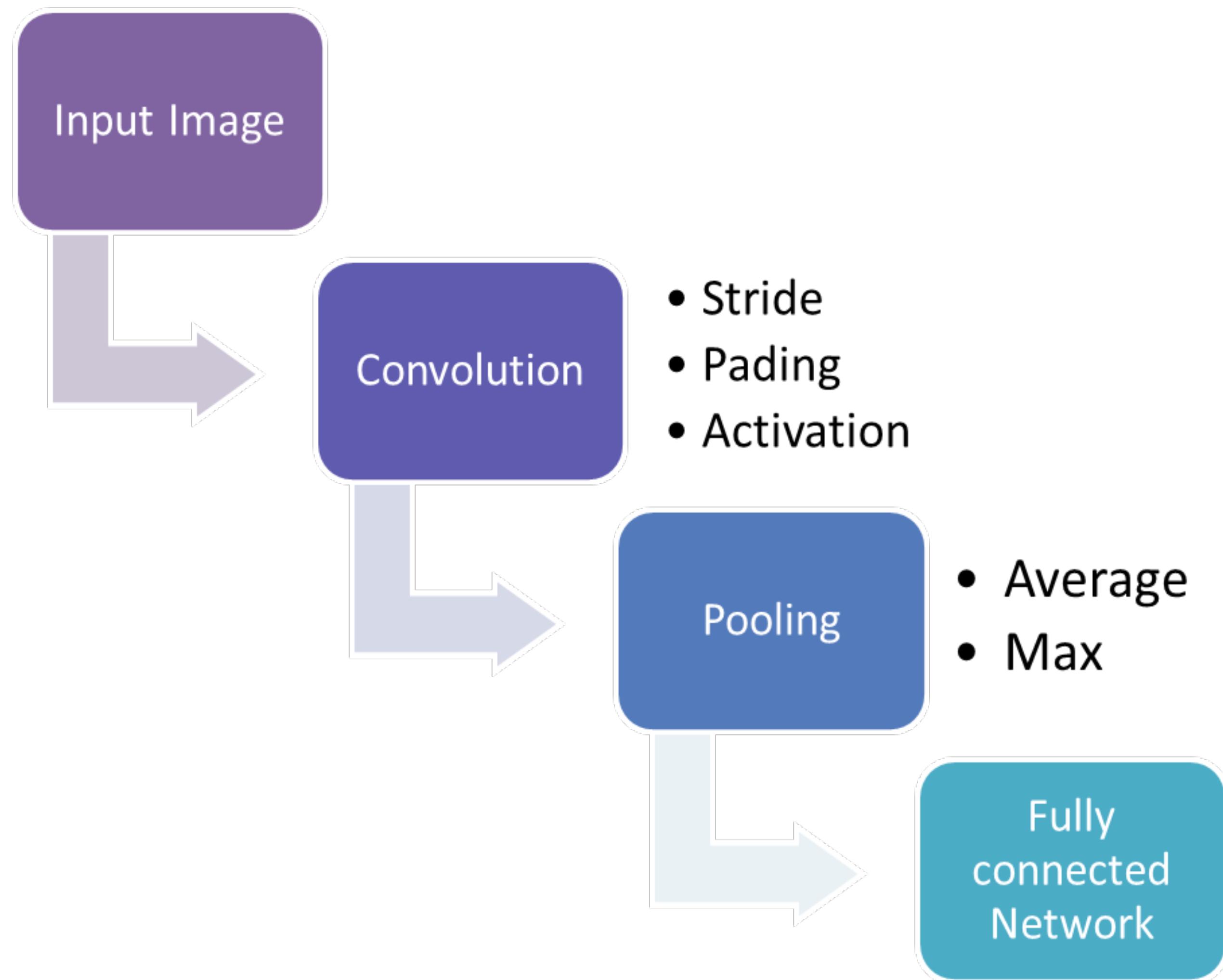
- CNN helps in analyzing the imaginary.
- Q -> What human see vs what system see?
- Size of image reducing with increase in stride value.
- Padding the input image with zero - retains depth.
- Kernel/filter that extracts useful info (edges, color, ...)
- Pooling helps in reducing the spatial size of the image.
- Output volume is controlled by 3 parameters (no. of filters, stride, zero padding) - $([W-F+2P]/S) + 1$.
- Output layer in CNN is a FC layer.

CNN: INPUT => CONV => RELU => FC => SOFTMAX

CNN that accepts an input, applies a convolution layer, then an activation layer, then a fully-connected layer, and, finally, a softmax classifier to obtain the output classification probabilities.

Architecture of CNN

- A typical CNN has 4 layers



Quiz

Ex:

Q1. Consider input image of $32 \times 32 \times 3$ and we have $10 \rightarrow 5 \times 5$ filters with stride 1, pad 2.

What is output volume?

Q2 . Number of parameters in this layer?

Key

- $(32+2*2-5)/1+1 = 32$ spatially, so $32 \times 32 \times 10$.
- each filter has $5 \times 5 \times 3 + 1 = 76$ params (+1 for bias) $\Rightarrow 76 \times 10 = 760$