

**DSE 3121 DEEP LEARNING**

## **Convolutional Neural Networks**

**Dr. Rohini Rao & Dr. Abhilash K Pai**

Dept. of Data Science and Computer Applications

MIT Manipal

# The Convolution Operation - 1D

- Convolution is a linear operation on two functions of a real-valued argument, where one function is applied over the other to yield element-wise dot products.
- Example: Consider a discrete signal ' $x_t$ ' which represents the position of a spaceship at time ' $t$ ' recorded by a laser sensor.
- Now, suppose that this sensor is noisy.
- To obtain a less noisy measurement we would like to average several measurements.
- Considering that, the most recent measurements are more important, we would like to take a weighted average over ' $x_t$ '. The new estimate at time ' $t$ ' is computed as follows:

$$s_t = \sum_{a=0}^{\infty} x_{t-a} w_{-a} = (x * w)_t$$

convolution  
↑  
input      Filter/Mask/Kernel



$x_0$



$x_1$



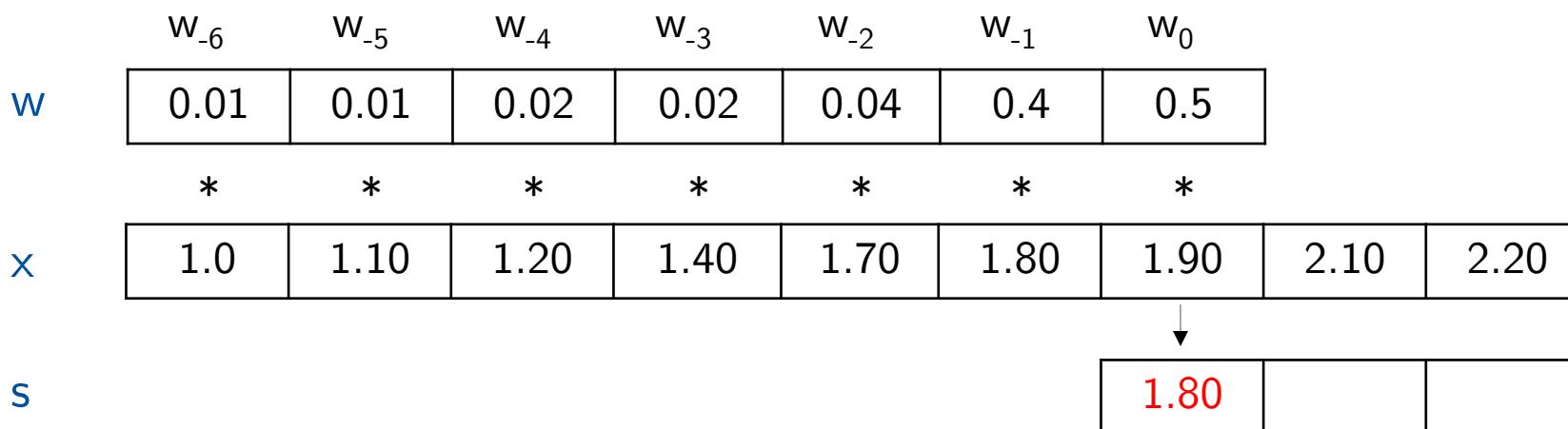
$x_2$

# The Convolution Operation - 1D

- In practice, we would sum only over a small window.

For example:  $s_t = \sum_{a=0}^6 x_{t-a} w_{-a}$

- We just slide the filter over the input and compute the value of  $s_t$  based on a window around  $x_t$



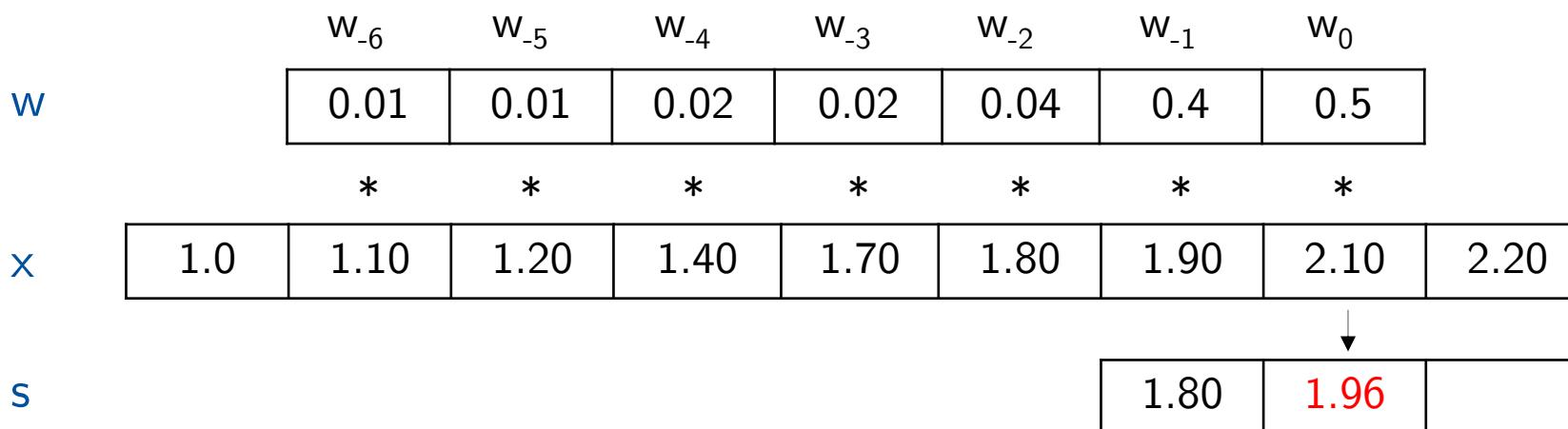
Content adapted from : CS7015 Deep Learning, Dept. of CSE, IIT Madras

# The Convolution Operation - 1D

- In practice, we would sum only over a small window.

For example:  $s_t = \sum_{a=0}^6 x_{t-a} w_{-a}$

- We just slide the filter over the input and compute the value of  $s_t$  based on a window around  $x_t$



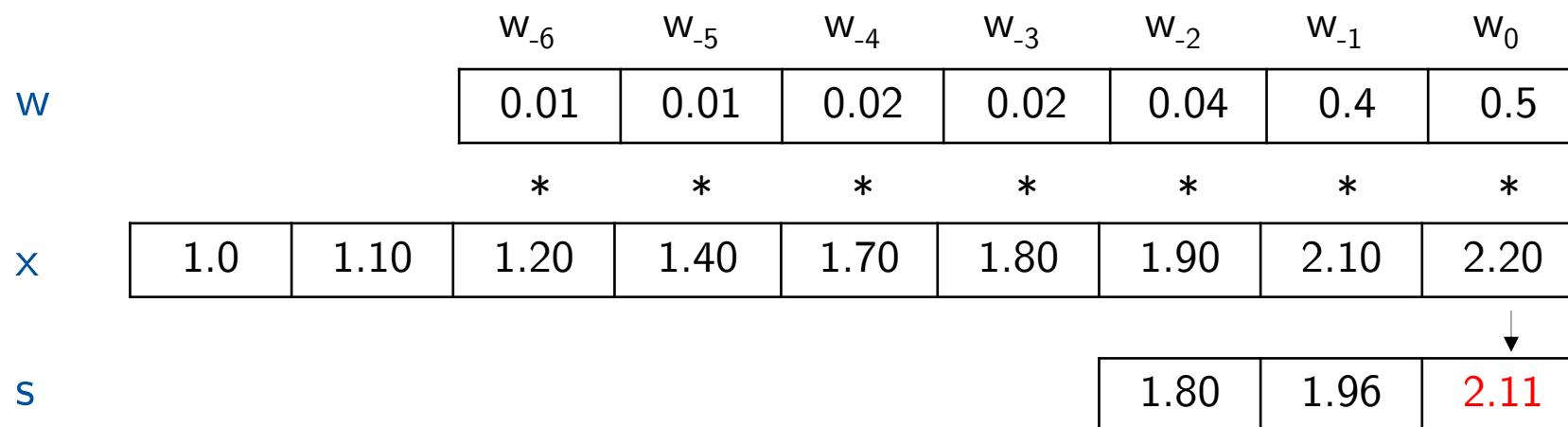
Content adapted from : CS7015 Deep Learning, Dept. of CSE, IIT Madras

# The Convolution Operation - 1D

- In practice, we would sum only over a small window.

$$\text{For example: } s_t = \sum_{a=0}^6 x_{t-a} w_{-a}$$

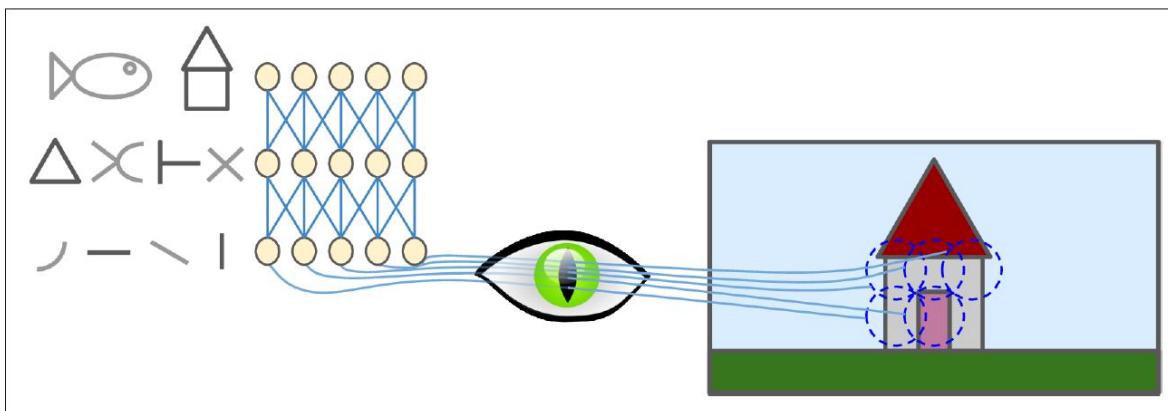
- We just slide the filter over the input and compute the value of  $s_t$  based on a window around  $x_t$



- Use cases of 1-D convolution : Audio signal processing, stock market analysis, time series analysis etc.

Content adapted from : CS7015 Deep Learning, Dept. of CSE, IIT Madras

# Local Receptive Field



Source: Aurelien Geron, "Hands-On Machine Learning with Scikit-Learn , Keras & Tensorflow, O'Reilly Publications

Figure 14-1. Local receptive fields in the visual cortex

- Convolutional neural networks (CNNs) emerged from the study of the brain's visual cortex, and used in image recognition since the 1980s
- **David H. Hubel and Torsten Wiesel** performed experiments on cats giving crucial insights on the structure of the visual cortex
- Important Findings
  - many neurons in the visual cortex have a small local receptive field
  - The receptive fields of different neurons may overlap, and together they tile the whole visual field.
  - Some neurons react only to images of horizontal lines, while others react only to lines with different orientations
  - Some neurons have larger receptive fields, and they react to more complex patterns that are combinations of the lower-level patterns.
  - These observations led to the idea that the higher-level neurons are based on the outputs of neighboring lower-level neurons

# Convolutional Neural Network

- Studies of the visual cortex inspired the neocognitron evolved into what we now call convolutional neural networks. An
- important milestone was a 1998 paper by Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner, which introduced the LeNet-5 architecture, widely used to recognize handwritten check numbers
- Building blocks: convolutional layers and pooling layers.

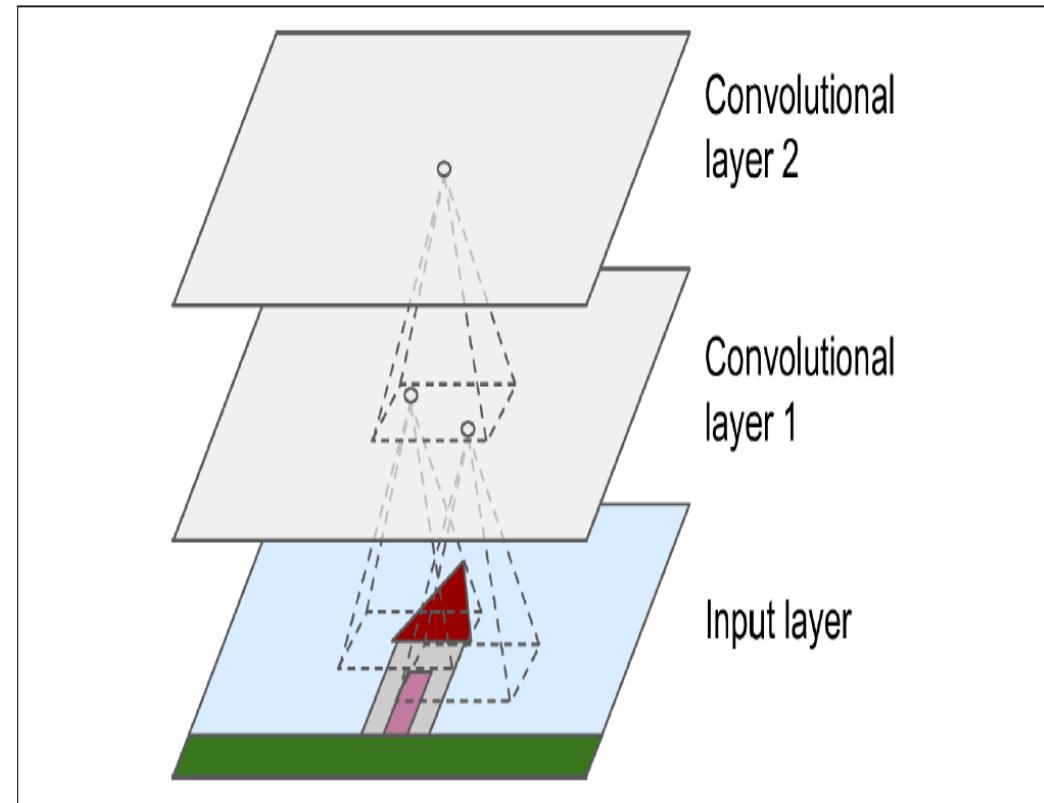
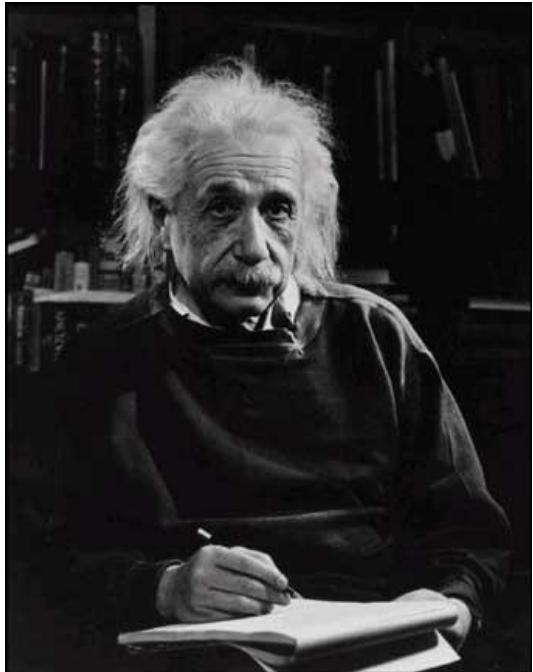


Figure 14-2. CNN layers with rectangular local receptive fields

Source: Aurelien Geron, "Hands-On Machine Learning with Scikit-Learn , Keras & Tensorflow, O'Reilly Publications

# Convolution in 2-D using Images : What is an Image?



What we see

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

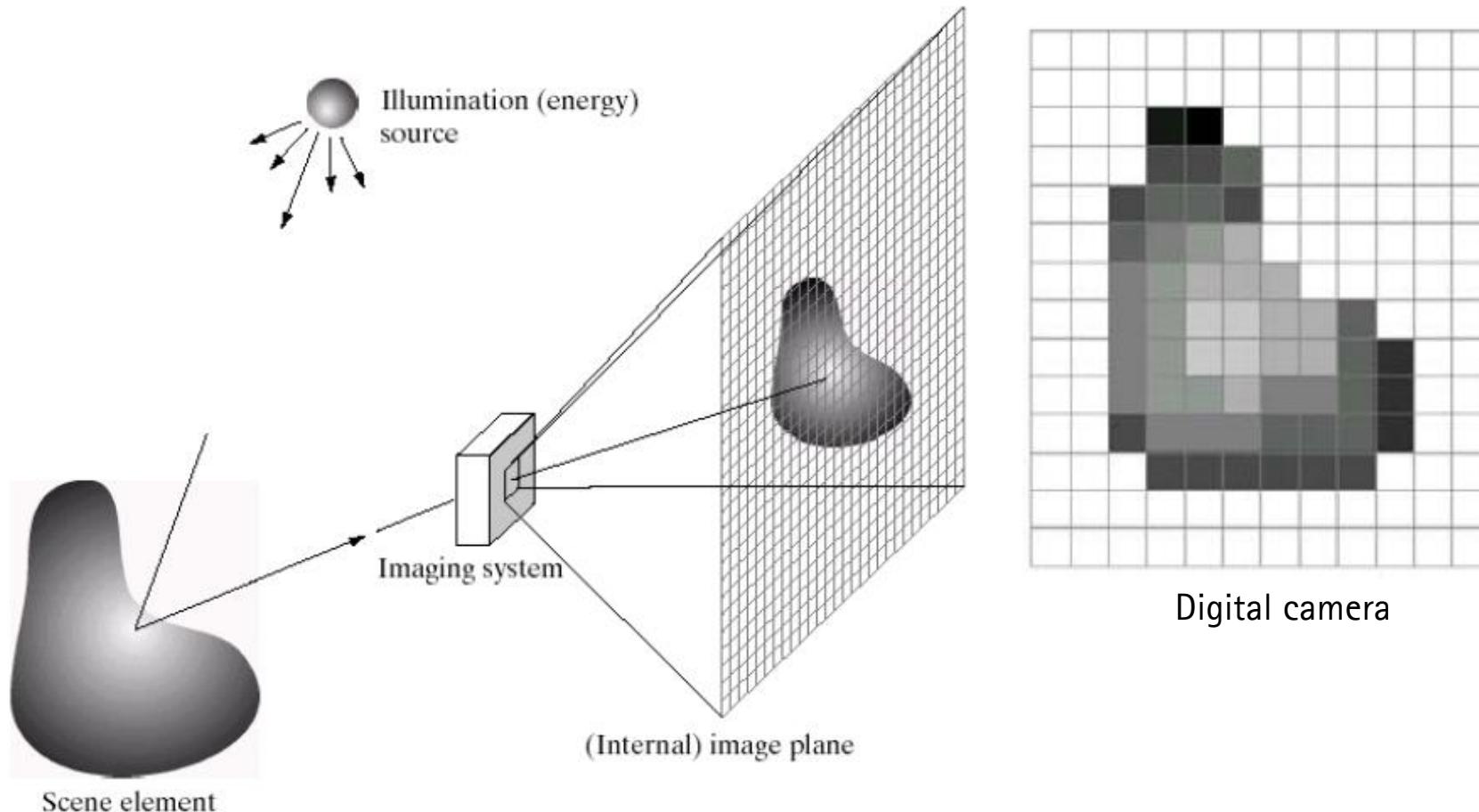
What a computer sees

# Convolution in 2-D using Images : What is an Image?

- An image can be represented mathematically as a function  $f(x,y)$  which gives the intensity value at position  $(x,y)$ , where,  $f(x,y) \in \{0,1,\dots,I_{\max-1}\}$  and  $x,y \in \{0,1,\dots,N-1\}$ .
- Larger the value of  $N$ , more is the clarity of the picture (larger resolution), but more data to be analyzed in the image.
- If the image is a **Gray-scale** (8-bit per pixel) image, then it requires  $N^2$  Bytes for storage.
- If the image is color - **RGB**, each pixel requires 3 Bytes of storage space.

$N$  is the resolution of the image and  $I_{\max}$  is the level of discretized brightness value.

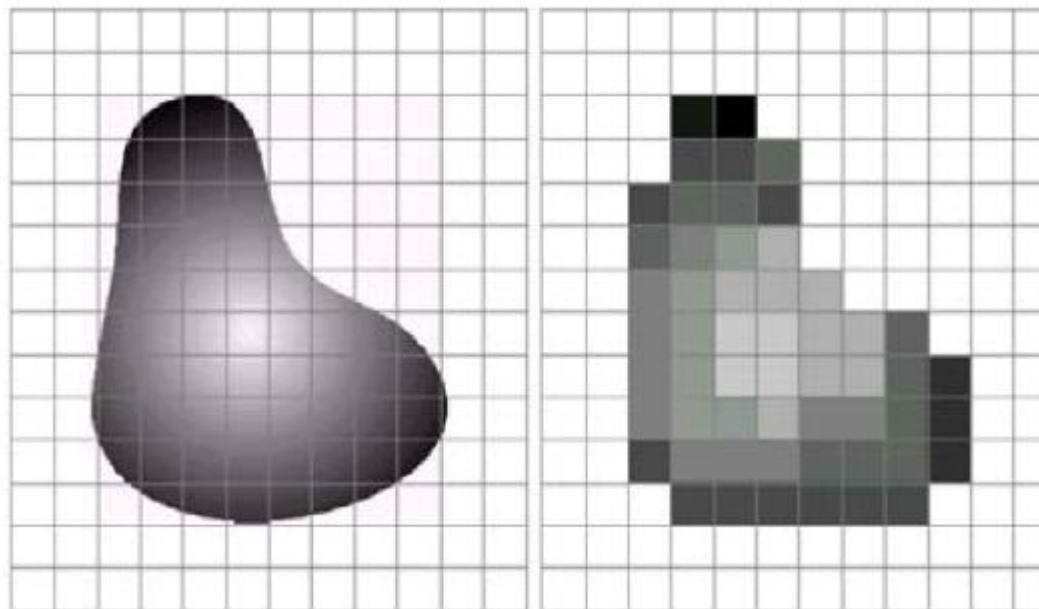
# Convolution in 2-D using Images : What is an Image?



[Source: D. Hoiem]

# Convolution in 2-D using Images : What is an Image?

- **Sample** the 2-D space on a regular grid.
- **Quantize** each sample, i.e., the photons arriving at each active cell are integrated and then digitized.



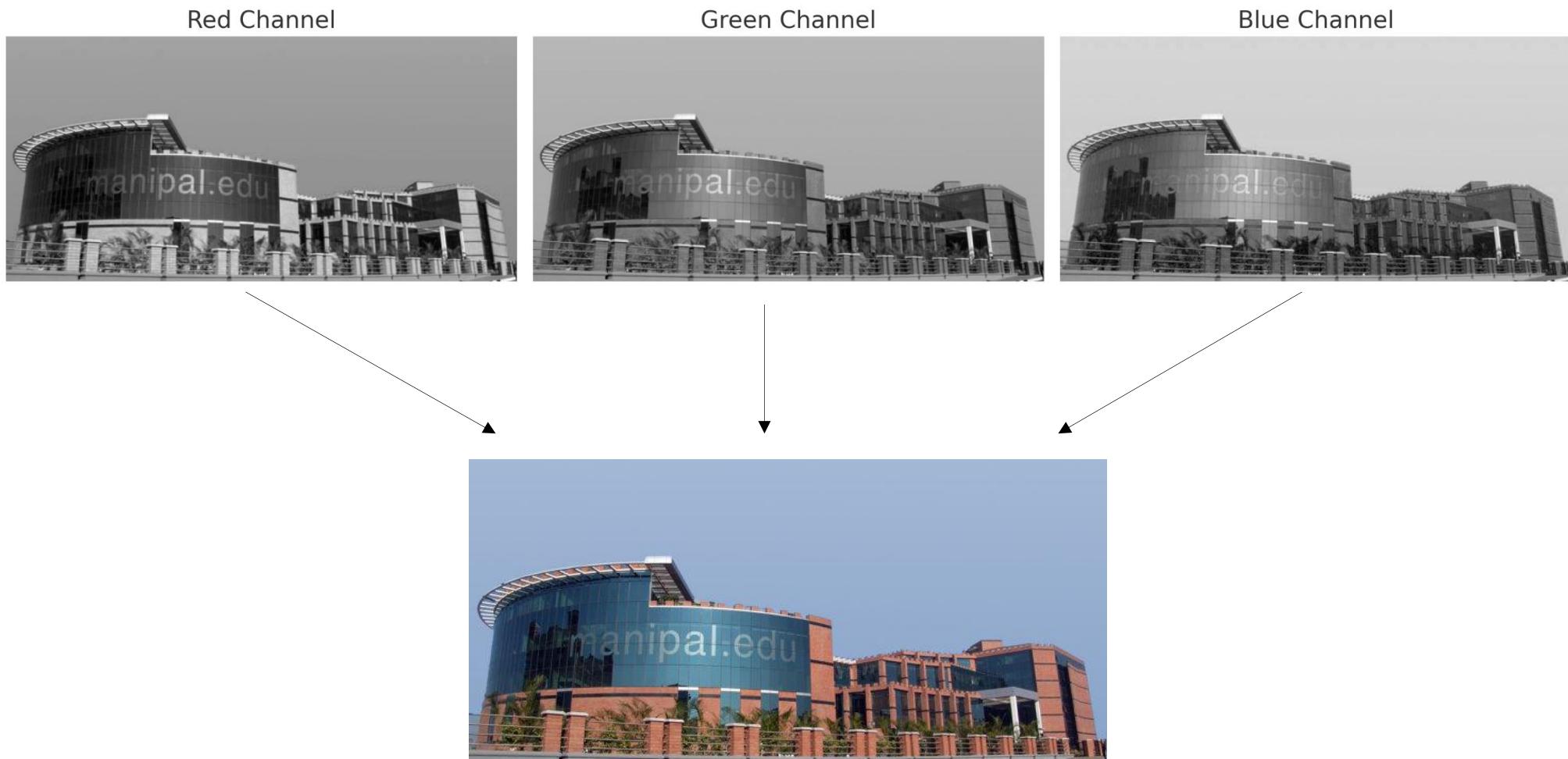
[Source: D. Hoiem]

# Convolution in 2-D using Images : What is an Image?

- A grid (matrix) of intensity values.

255	255	255	255	255	255	255	255	255	255	255	255	255	255	255
255	255	255	255	255	255	255	255	255	255	255	255	255	255	255
255	255	255	0	0	255	255	255	255	255	255	255	255	255	255
255	255	255	75	75	75	255	255	255	255	255	255	255	255	255
255	255	75	95	95	75	255	255	255	255	255	255	255	255	255
255	255	96	127	145	175	255	255	255	255	255	255	255	255	255
255	255	127	145	175	175	175	255	95	255	255	255	255	255	255
255	255	127	145	200	200	175	175	95	255	255	255	255	255	255
255	255	127	145	145	175	127	127	95	47	255	255	255	255	255
255	255	127	145	145	175	127	127	95	47	255	255	255	255	255
255	255	74	127	127	127	95	95	95	47	255	255	255	255	255
255	255	255	74	74	74	74	74	74	255	255	255	255	255	255
255	255	255	255	255	255	255	255	255	255	255	255	255	255	255
255	255	255	255	255	255	255	255	255	255	255	255	255	255	255

# Convolution in 2-D using Images : What is an Image?



# The Convolution Operation - 2D

- Images are good examples of 2-D inputs.
- A 2-D convolution of an **Image 'I'** using a **filter 'K'** of size '**m x n**' is now defined as (looking at previous pixels):

$$S_{ij} = (I * K)_{ij} = \sum_{a=0}^{m-1} \sum_{b=0}^{n-1} I_{i-a,j-b} K_{a,b}$$

- In practice, one of the way is to look at the succeeding pixels:

$$S_{ij} = (I * K)_{ij} = \sum_{a=0}^{m-1} \sum_{b=0}^{n-1} I_{i+a,j+b} K_{a,b}$$

# The Convolution Operation - 2D

- Another way is to consider center pixel as reference pixel, and then look at its surrounding pixels:

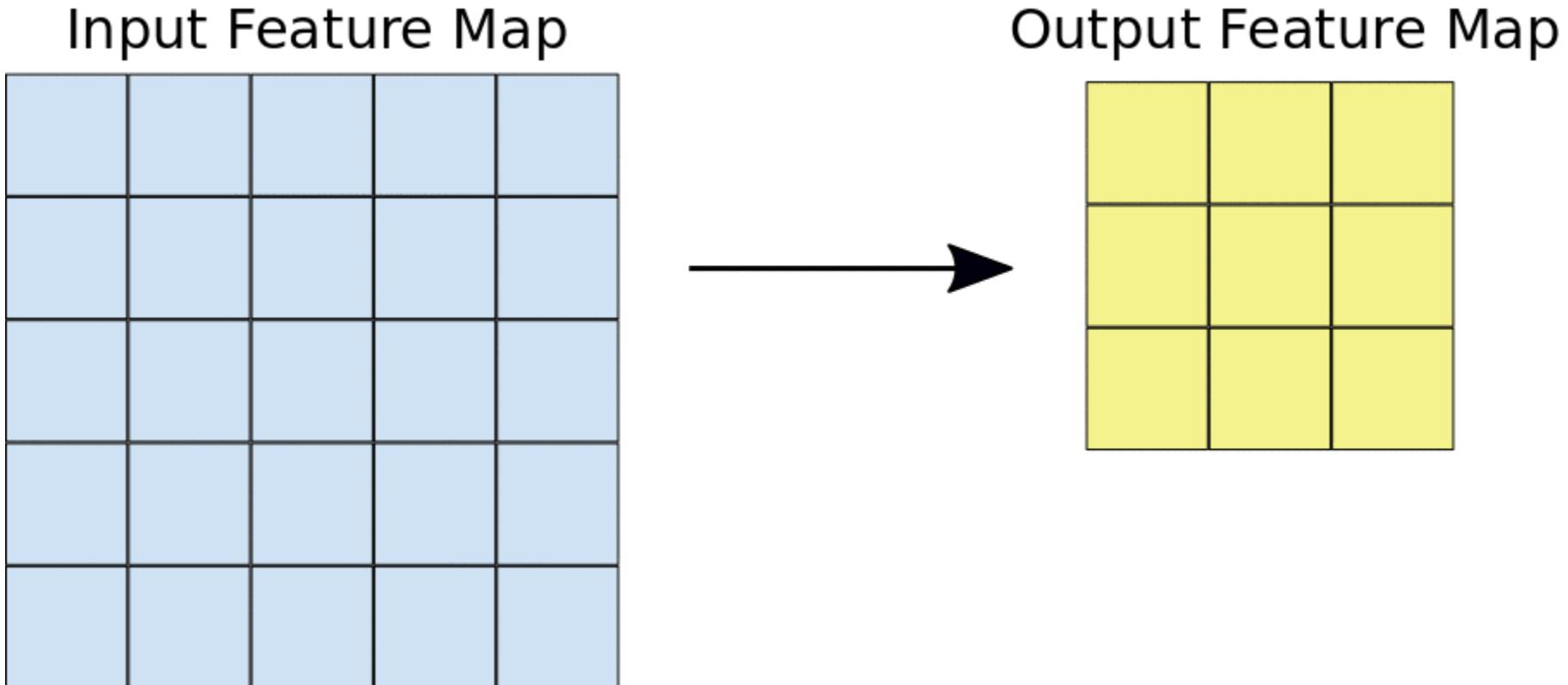
$$S_{ij} = (I * K)_{ij} = \sum_{a=[-m/2]}^{[m/2]} \sum_{b=[-n/2]}^{[n/2]} I_{i-a,j-b} * K_{(m/2)+a, (n/2)+b}$$

Pixel of interest

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

Content adapted from : CS7015 Deep Learning, Dept. of CSE, IIT Madras

# The Convolution Operation - 2D



Source: <https://developers.google.com/>

# The Convolution Operation - 2D

Input Image

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

Convolutional Filter

1	0	0
1	1	0
0	0	1

Source: <https://developers.google.com/>

# The Convolution Operation - 2D

3x1	5x0	2x0	8	1
9x1	7x1	5x0	4	3
2x0	0x0	6x1	1	6
6	3	7	9	2
1	4	9	5	1

3+0+0+9+7+0+0+0+6



Output Feature Map

25	18	17
18	22	14
20	15	23

Source: <https://developers.google.com/>

# The Convolution Operation - 2D



$$* \frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} =$$



Smoothening Filter

# The Convolution Operation - 2D



$$* \begin{matrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{matrix} =$$



Sharpening Filter

# The Convolution Operation - 2D



$$\begin{matrix} * & \begin{matrix} 1 & 1 & 1 \\ 1 & -8 & 1 \\ 1 & 1 & 1 \end{matrix} & = \end{matrix}$$

Filter for edge  
detection



# The Convolution Operation - 2D

[But what is a convolution? \(youtube.com\)](https://www.youtube.com)

# The Convolution Operation – 2D : Various filters (edge detection)

## Prewitt

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

$S_x$

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

$S_y$



Input image



After applying  
Horizontal edge  
detection filter

## Sobel

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

$S_x$

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

$S_y$



After applying  
Vertical edge  
detection filter

## Laplacian

0	1	0
1	-4	1
0	1	0

## Roberts

0	1
-1	0

$S_x$

1	0
0	-1

$S_y$

# The Convolution Operation - 2D

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

Filter 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

Input image:  $6 \times 6$

# The Convolution Operation - 2D

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

Filter 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

Input image:  $6 \times 6$

# The Convolution Operation - 2D

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

Filter 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

Dot  
product  


Input image:  $6 \times 6$

# The Convolution Operation - 2D

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

Filter 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

Dot  
product



3

Input image:  $6 \times 6$

# The Convolution Operation - 2D

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

Filter 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

Dot  
product

3

Input image:  $6 \times 6$

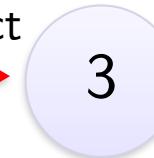
Note: Stride is the number of “unit” the kernel is shifted per slide over rows/ columns

# The Convolution Operation - 2D

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

Dot  
product



3

Input image:  $6 \times 6$

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

Filter 1

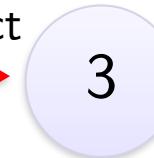
Note: Stride is the number of “unit” the kernel is shifted per slide over rows/ columns

# The Convolution Operation - 2D

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

Dot  
product



3

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

Filter 1

Input image:  $6 \times 6$

Note: Stride is the number of “unit” the kernel is shifted per slide over rows/ columns

# The Convolution Operation - 2D

stride=1

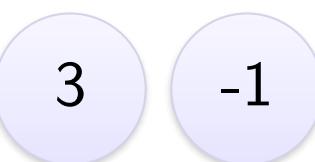
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

Input image:  $6 \times 6$

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

Filter 1

Dot  
product



Note: Stride is the number of “unit” the kernel is shifted per slide over rows/ columns

# The Convolution Operation - 2D

If stride=2

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

Filter 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

Input image:  $6 \times 6$

Note: Stride is the number of “unit” the kernel is shifted per slide over rows/ columns

# The Convolution Operation - 2D

If stride=2

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

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

Filter 1

3

Input image:  $6 \times 6$

Note: Stride is the number of “unit” the kernel is shifted per slide over rows/ columns

# The Convolution Operation - 2D

If stride=2

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

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

Filter 1

3

Input image:  $6 \times 6$

Note: Stride is the number of “unit” the kernel is shifted per slide over rows/ columns

# The Convolution Operation - 2D

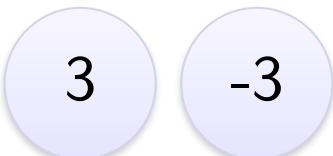
If stride=2

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

Input image:  $6 \times 6$

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

Filter 1



Note: Stride is the number of “unit” the kernel is shifted per slide over rows/ columns

# The Convolution Operation - 2D

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

Input image:  $6 \times 6$

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

Filter 1

# The Convolution Operation - 2D

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

Input image:  $6 \times 6$

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

Filter 1

# The Convolution Operation - 2D

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

3

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

Filter 1

Input image:  $6 \times 6$

# The Convolution Operation - 2D

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

Input image:  $6 \times 6$

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

Filter 1

3

# The Convolution Operation - 2D

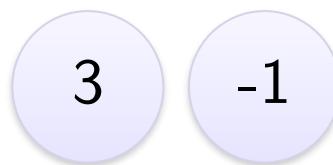
stride=1

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

Input image:  $6 \times 6$

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

Filter 1



# The Convolution Operation - 2D

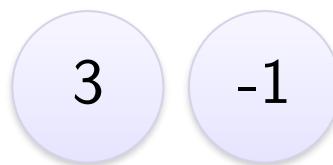
stride=1

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

Input image:  $6 \times 6$

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

Filter 1



# The Convolution Operation - 2D

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

Input image:  $6 \times 6$

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

Filter 1



# The Convolution Operation - 2D

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

Input image:  $6 \times 6$

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

Filter 1



# The Convolution Operation - 2D

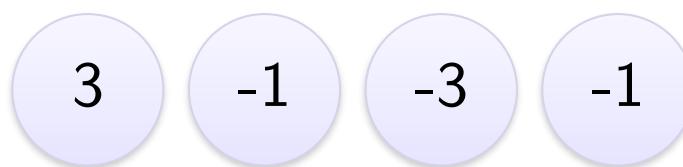
stride=1

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

Input image:  $6 \times 6$

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

Filter 1



# The Convolution Operation - 2D

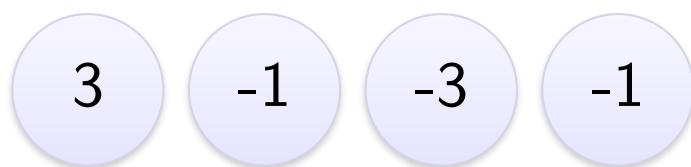
stride=1

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

Input image:  $6 \times 6$

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

Filter 1



# The Convolution Operation - 2D

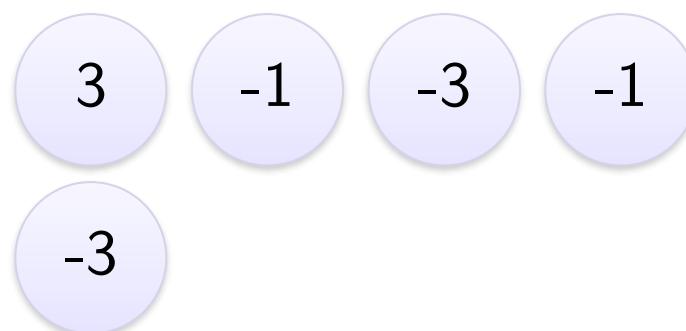
stride=1

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

Input image:  $6 \times 6$

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

Filter 1



# The Convolution Operation - 2D

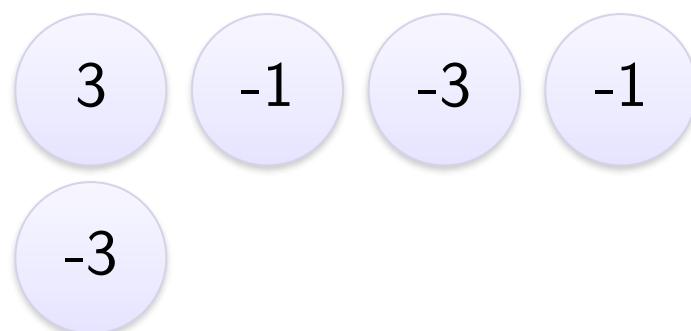
stride=1

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

Input image:  $6 \times 6$

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

Filter 1



# The Convolution Operation - 2D

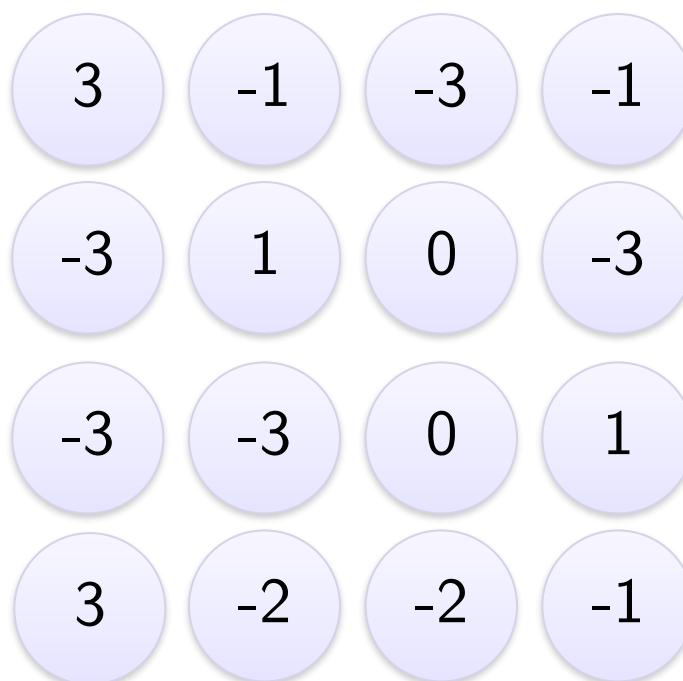
stride=1

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

Input image:  $6 \times 6$

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

Filter 1



# The Convolution Operation - 2D

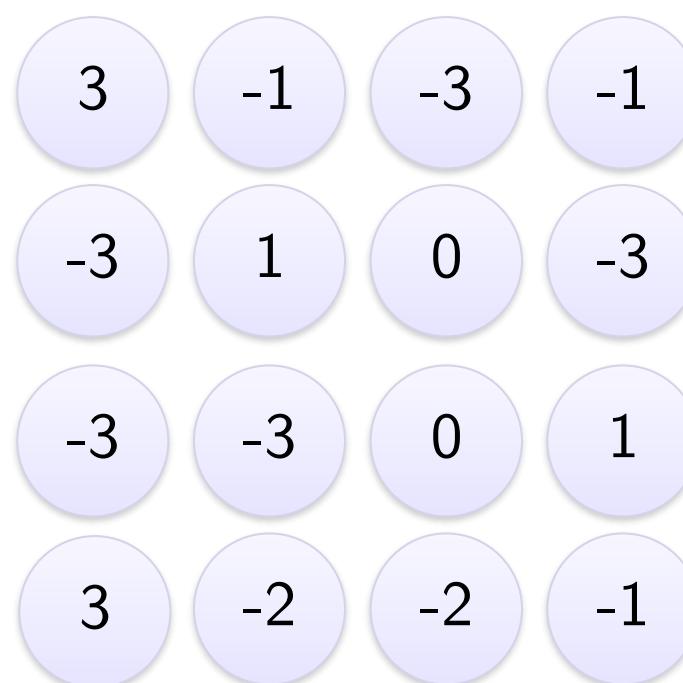
stride=1

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

Input image:  $6 \times 6$

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

Filter 1



$4 \times 4$  Feature Map

# The Convolution Operation - 2D

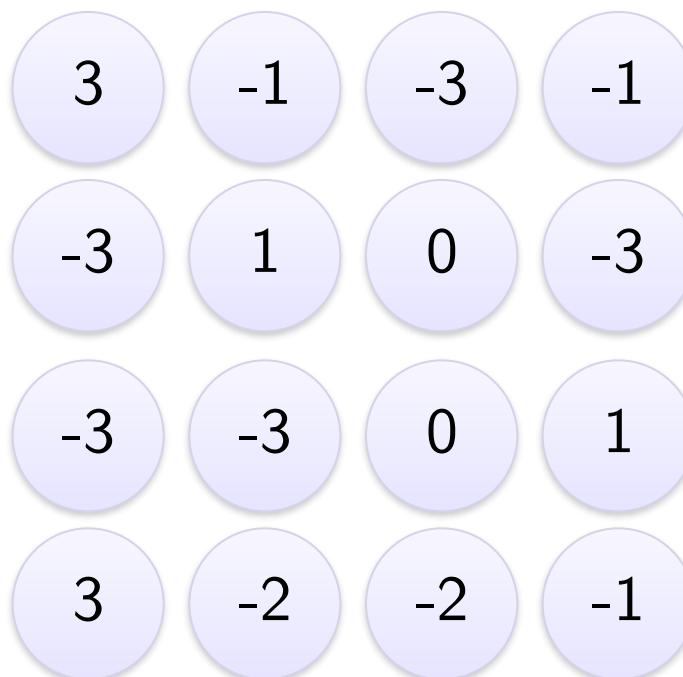
stride=1

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

Input image:  $6 \times 6$

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

Filter 2



# The Convolution Operation - 2D

stride=1

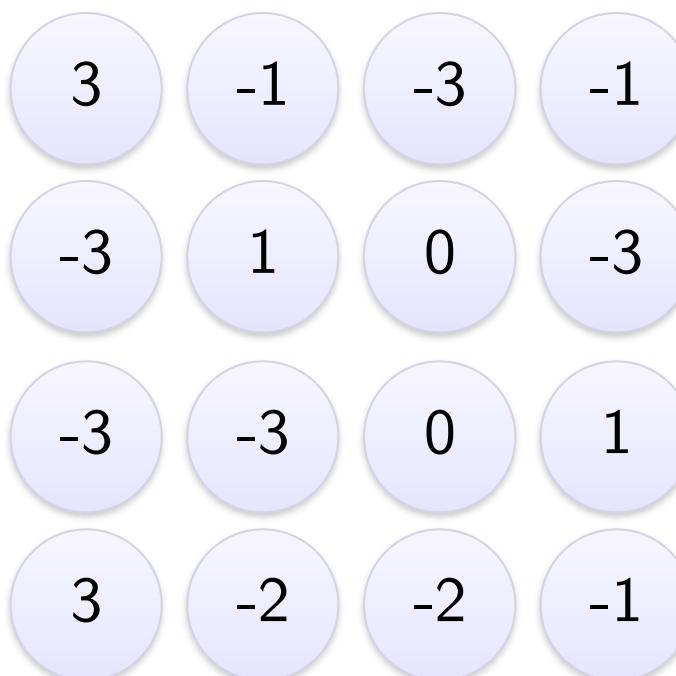
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

Input image:  $6 \times 6$

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

Filter 2

Repeat for each filter!



# The Convolution Operation - 2D

stride=1

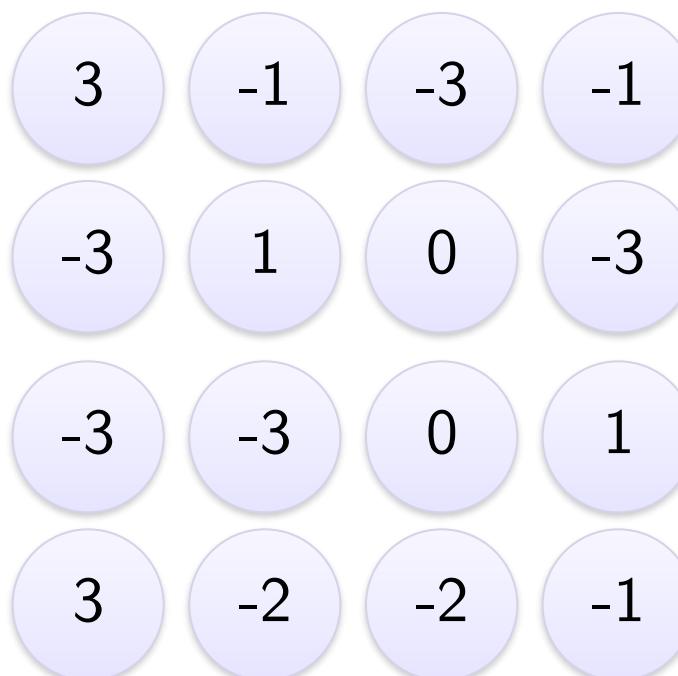
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

Input image:  $6 \times 6$

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

Filter 2

Repeat for each filter!



# The Convolution Operation - 2D

stride=1

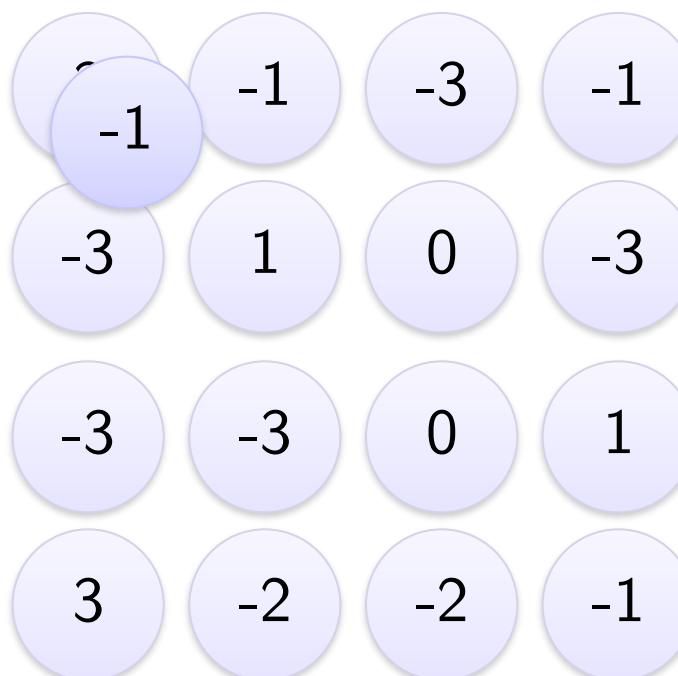
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

Input image:  $6 \times 6$

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

Filter 2

Repeat for each filter!



# The Convolution Operation - 2D

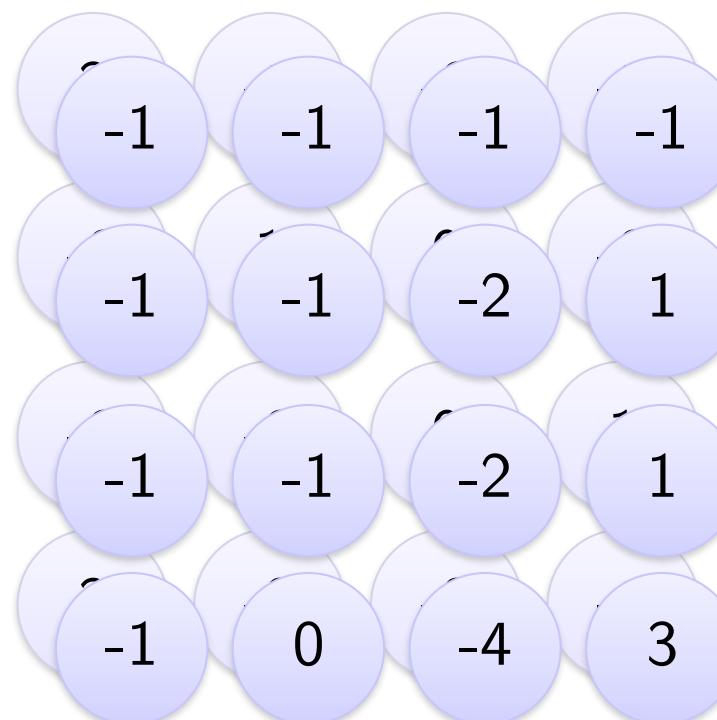
stride=1

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

Input image:  $6 \times 6$

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

Filter 2



Repeat for each filter!

# The Convolution Operation - 2D

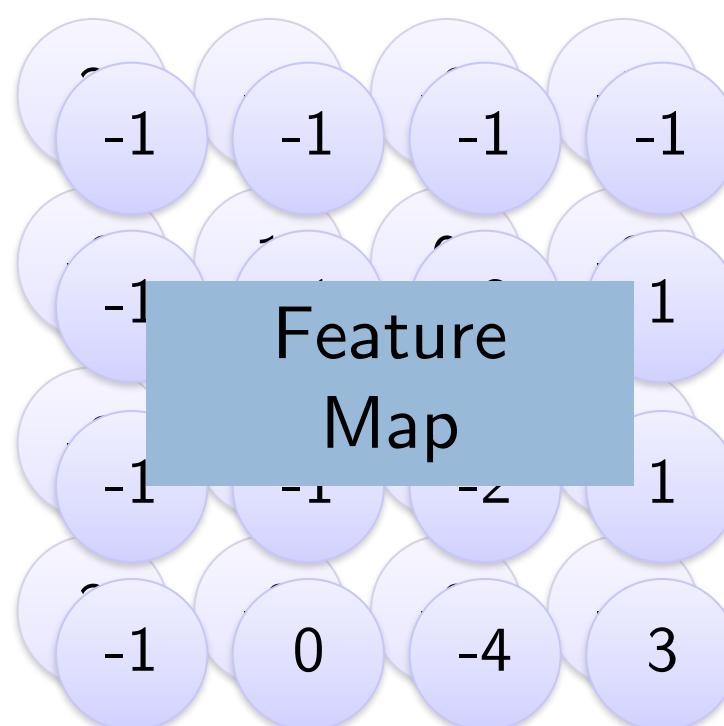
stride=1

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

Input image:  $6 \times 6$

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

Filter 2



Repeat for each filter!

# The Convolution Operation - 2D

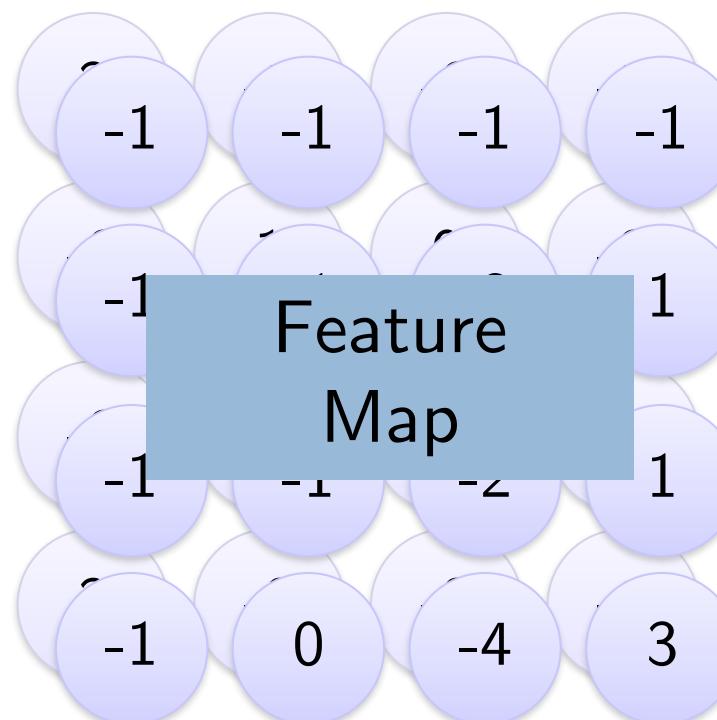
stride=1

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

Input image:  $6 \times 6$

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

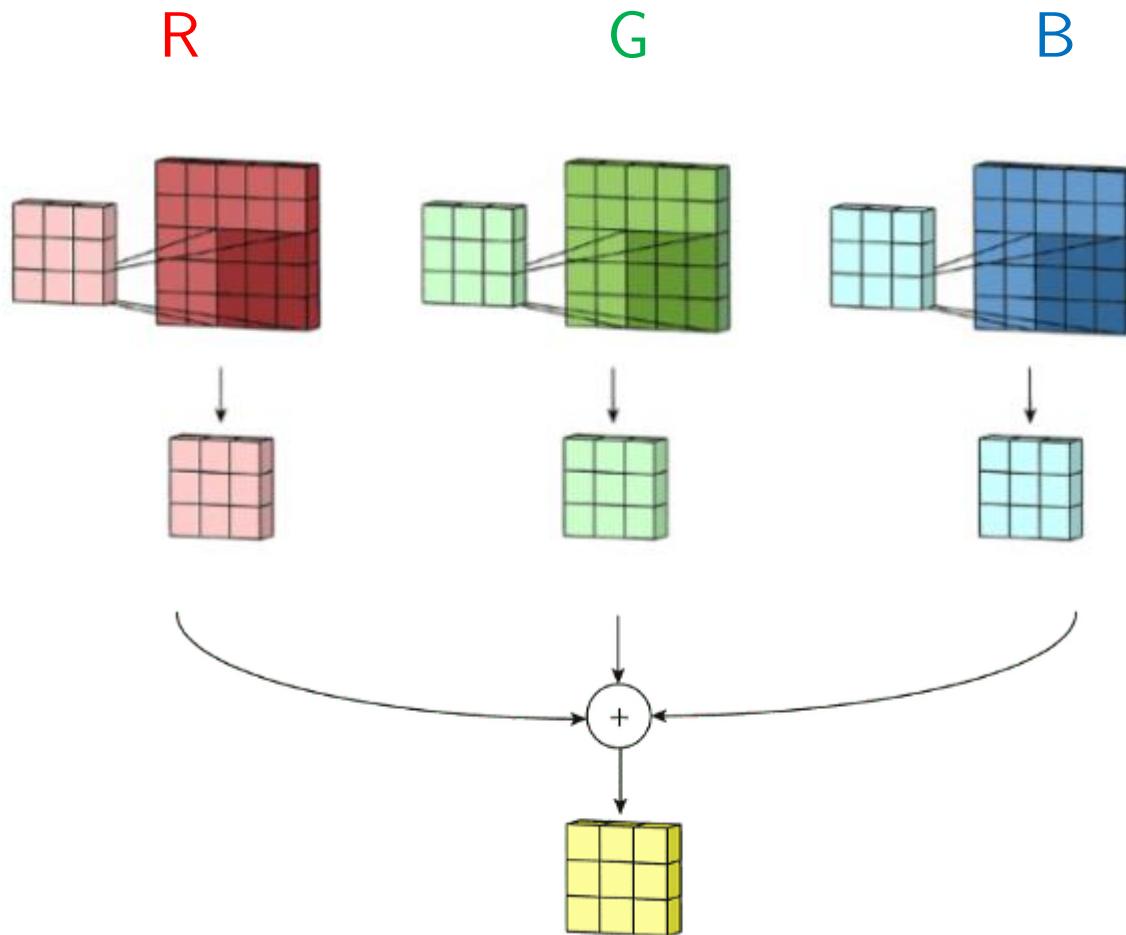
Filter 2



Repeat for each filter!

Two  $4 \times 4$  images  
Forming  $4 \times 4 \times 2$  matrix

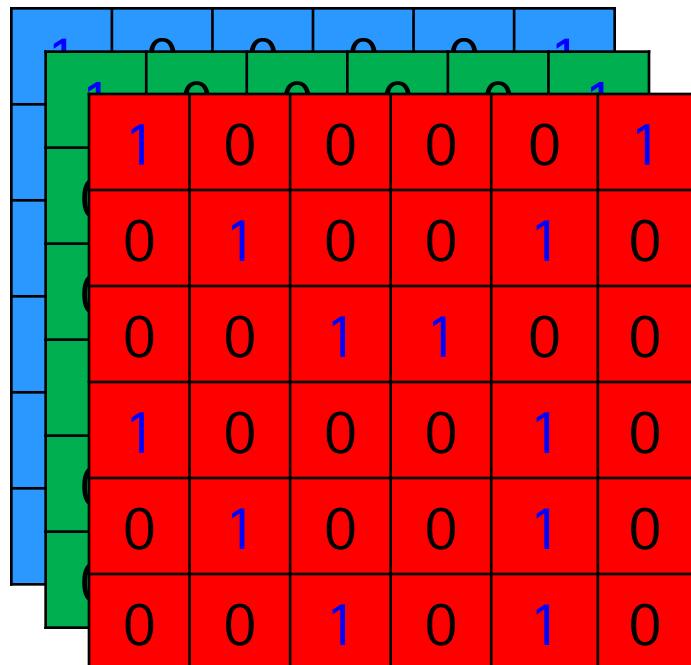
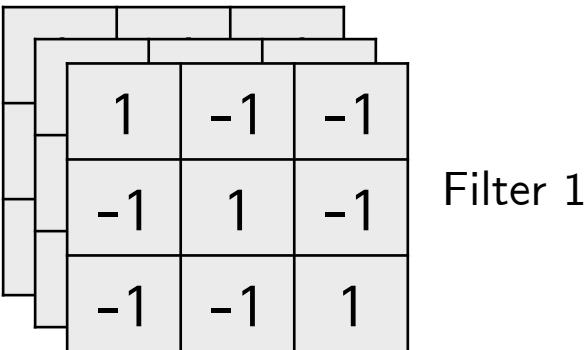
# The Convolution Operation –RGB Images



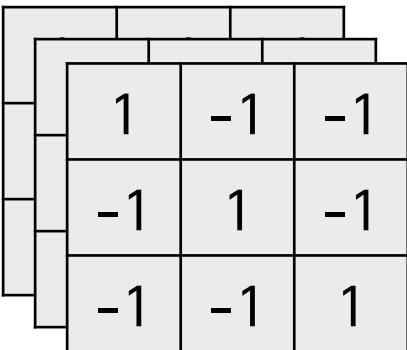
Apply the filter to R, G, and B channels of the image and combine the resultant feature maps to obtain a 2-D feature map.

Source: Intuitively Understanding Convolutions for Deep Learning | by Irhum Shafkat | Towards Data Science

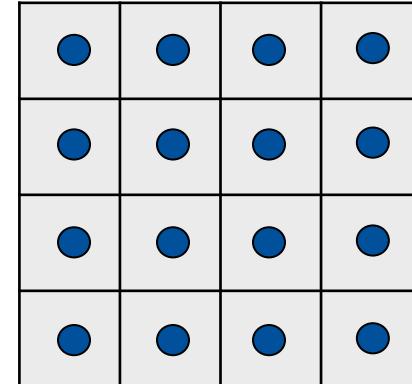
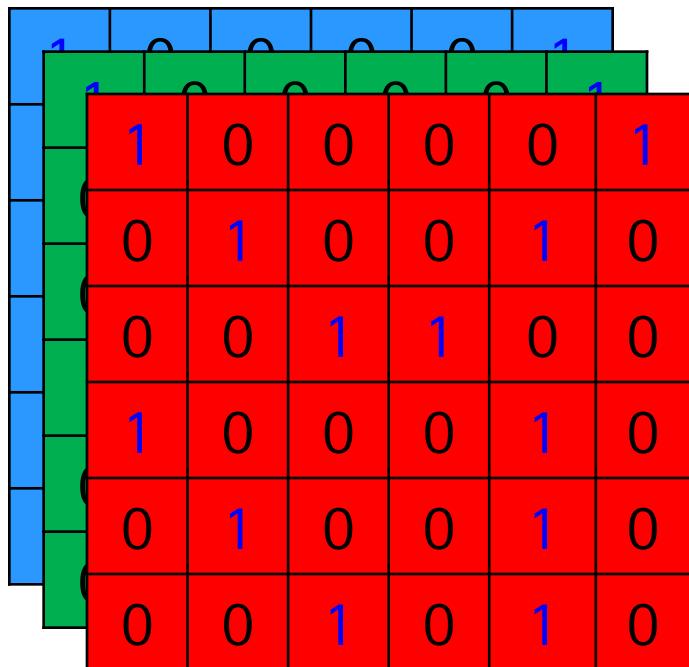
# The Convolution Operation –RGB Images multiple filters



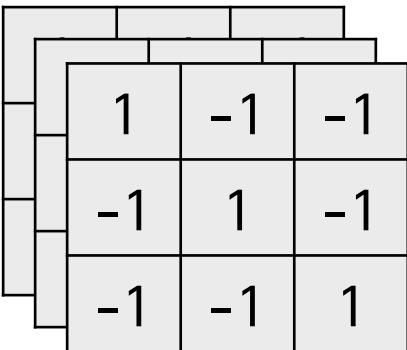
# The Convolution Operation –RGB Images multiple filters



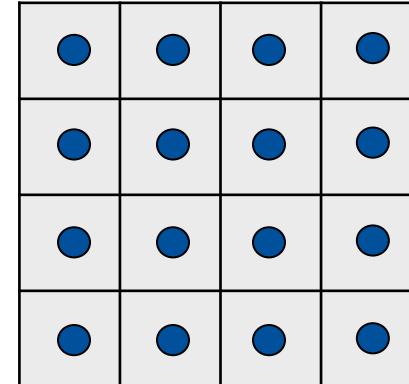
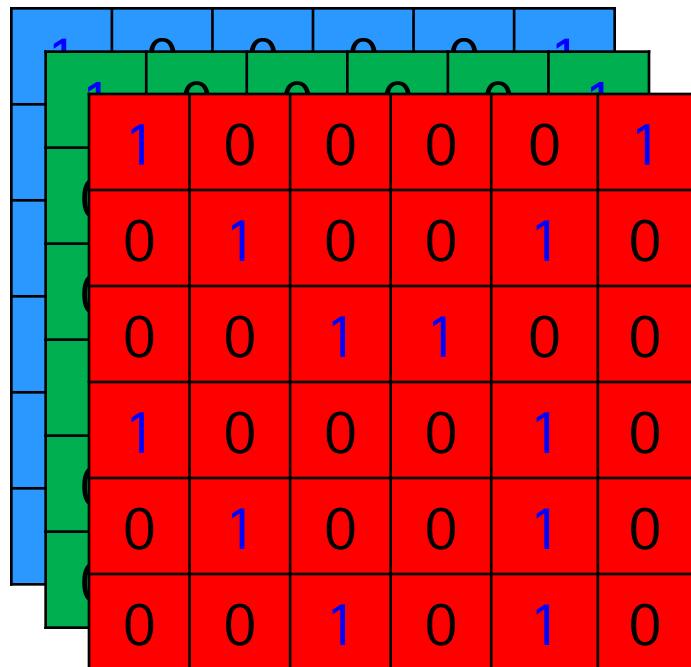
Filter 1



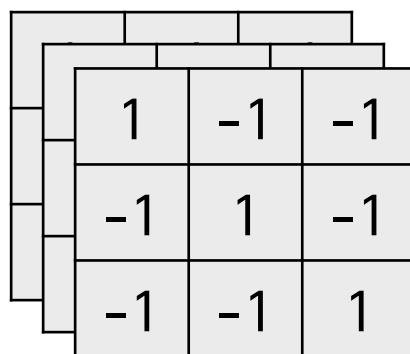
# The Convolution Operation –RGB Images multiple filters



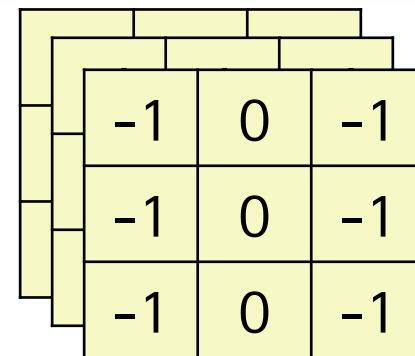
Filter 1



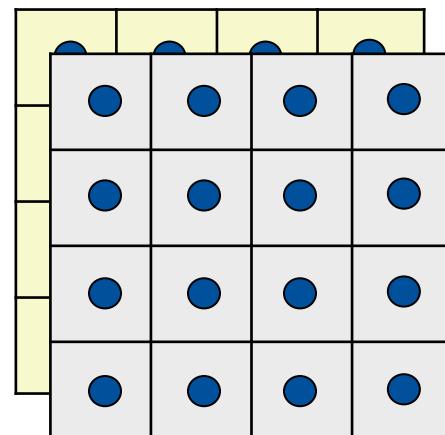
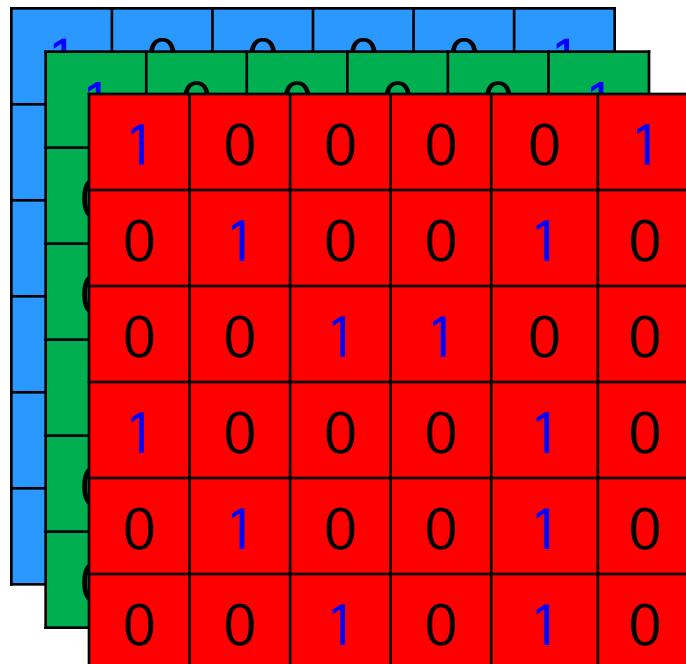
# The Convolution Operation –RGB Images multiple filters



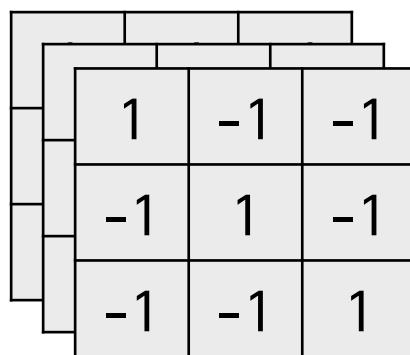
Filter 1



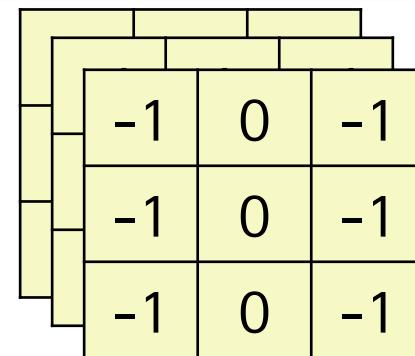
Filter 2



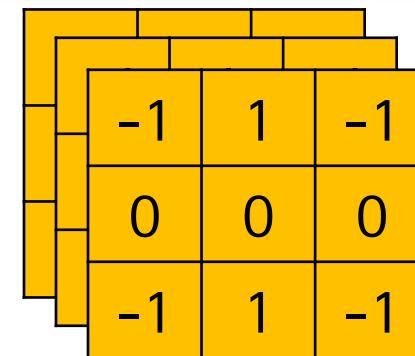
# The Convolution Operation –RGB Images multiple filters



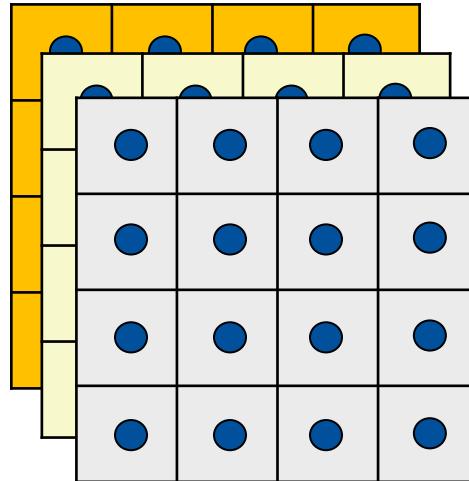
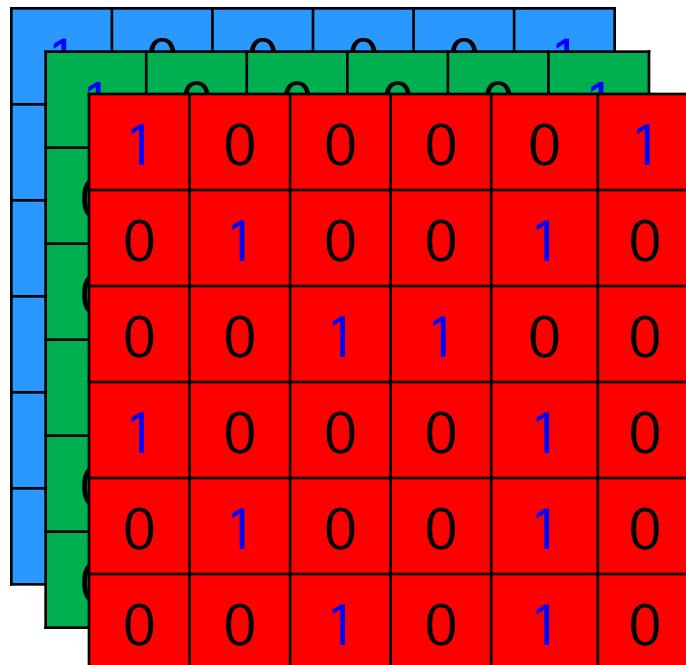
Filter 1



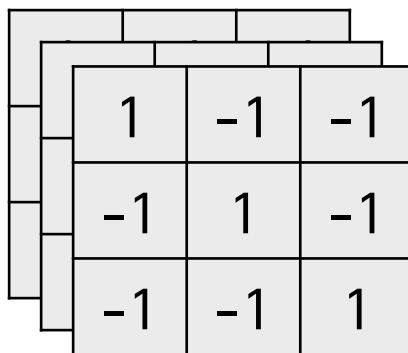
Filter 2



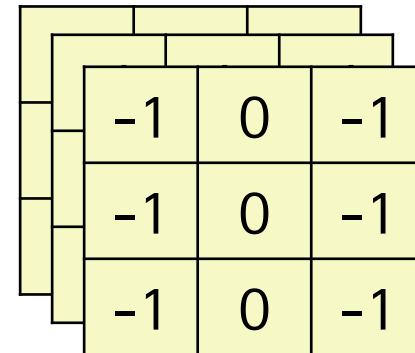
Filter K



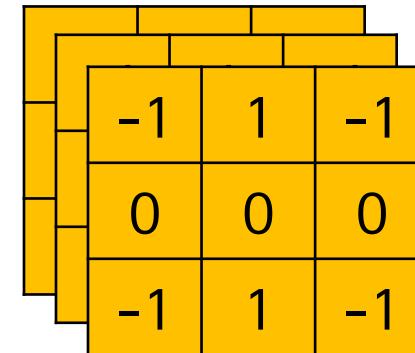
# The Convolution Operation –RGB Images multiple filters



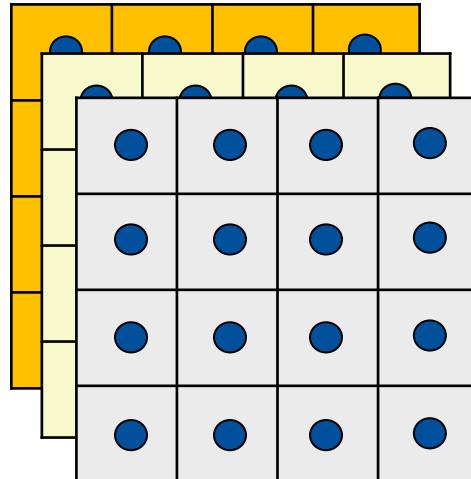
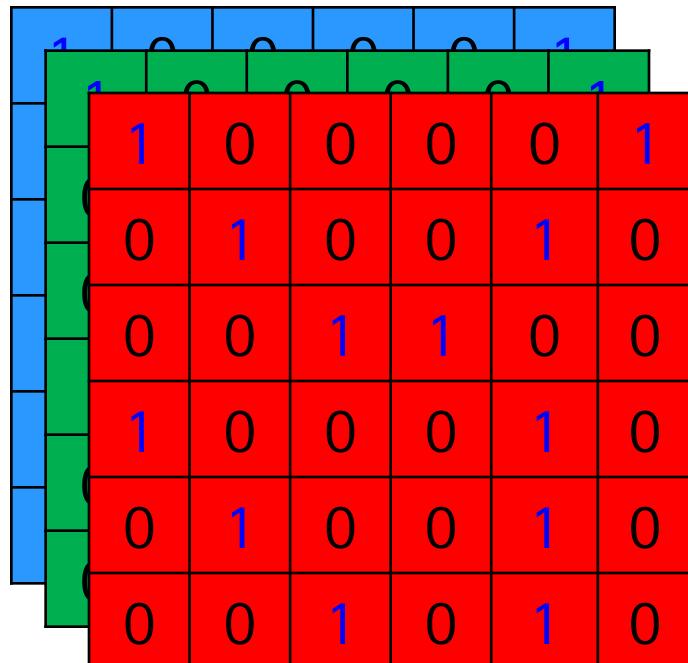
Filter 1



Filter 2

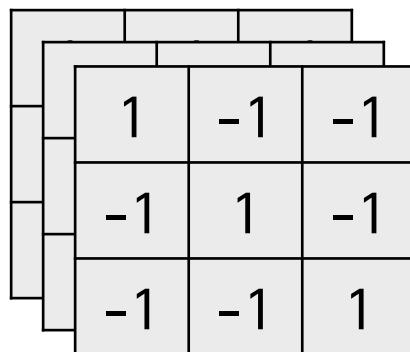


Filter K

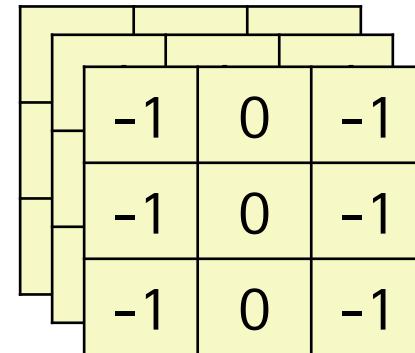


K-filters = K-Feature Maps

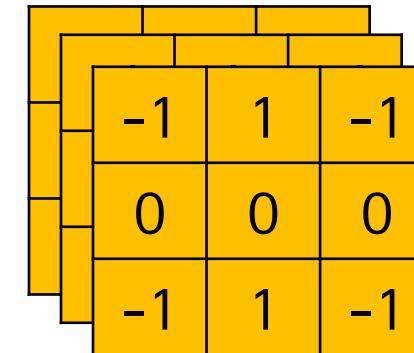
# The Convolution Operation –RGB Images multiple filters



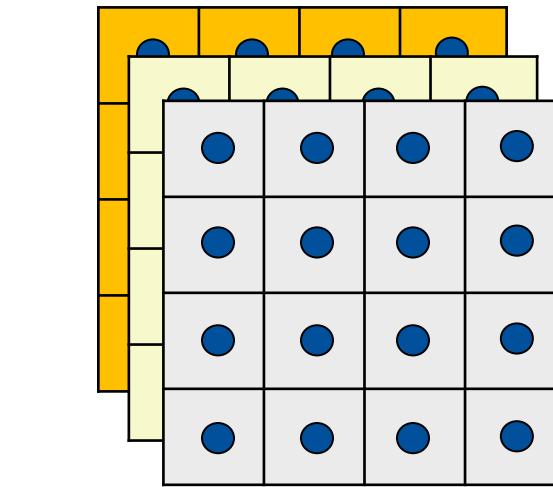
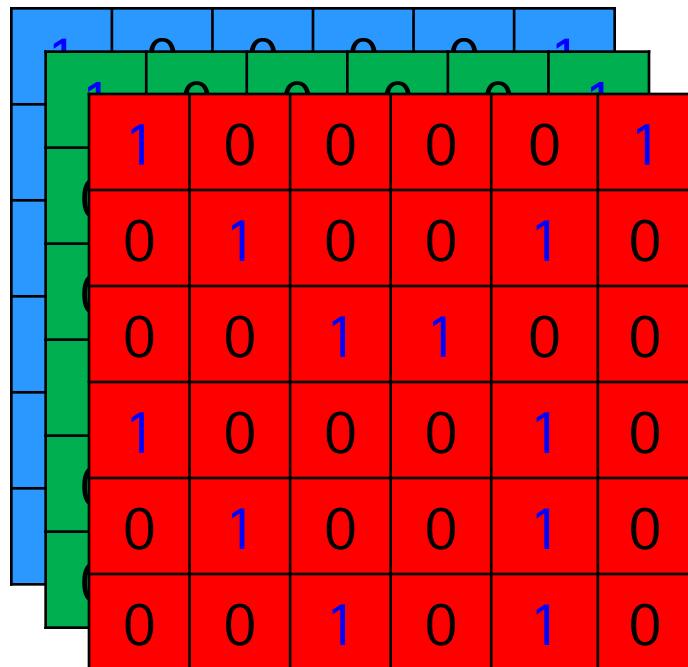
Filter 1



Filter 2



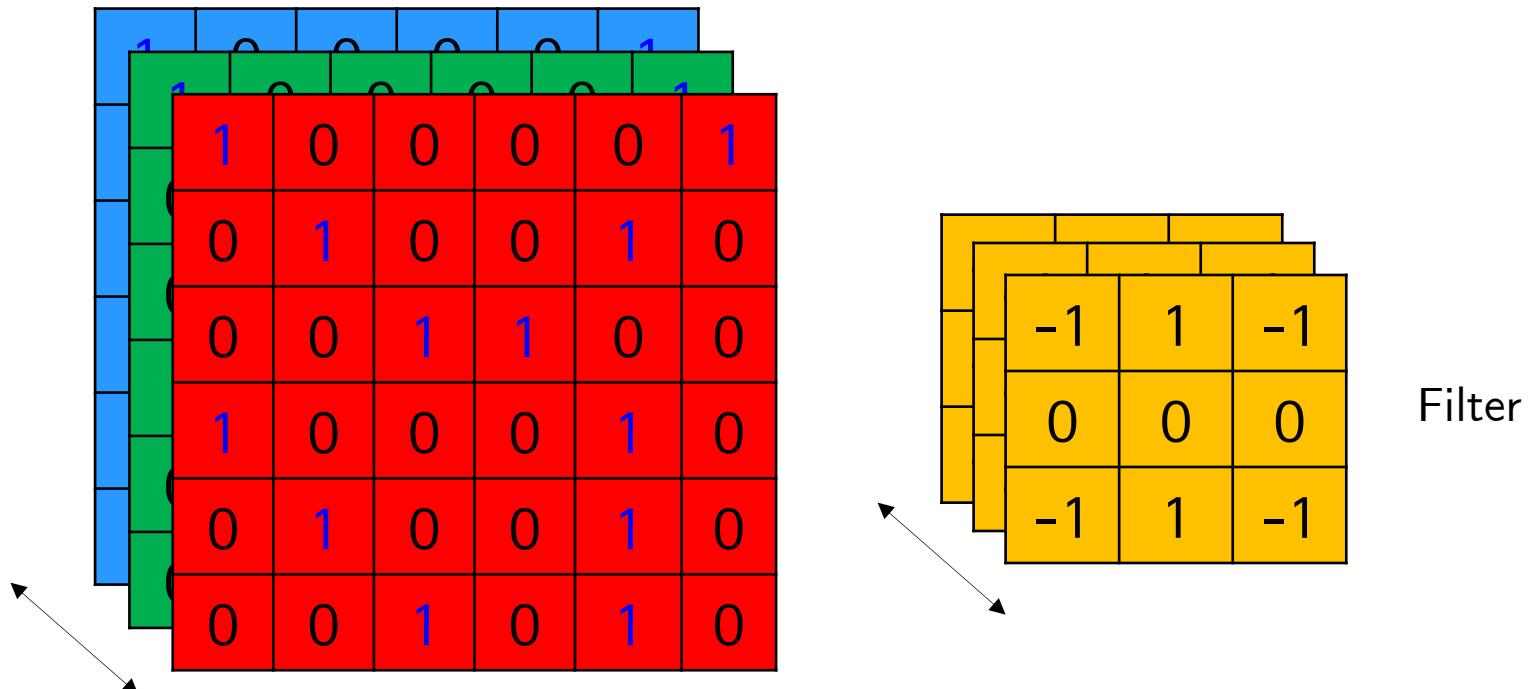
Filter K



K-filters = K-Feature Maps

Depth of feature map = No. of feature maps = No. of filters

# The Convolution Operation : Terminologies



1. Depth of an Input Image = No. of channels in the Input Image = Depth of a filter
2. Assuming square filters, Spatial Extent ( $F$ ) of a filter is the size of the filter

# The Convolution Operation : Zero Padding




conv<sub>3x3</sub>

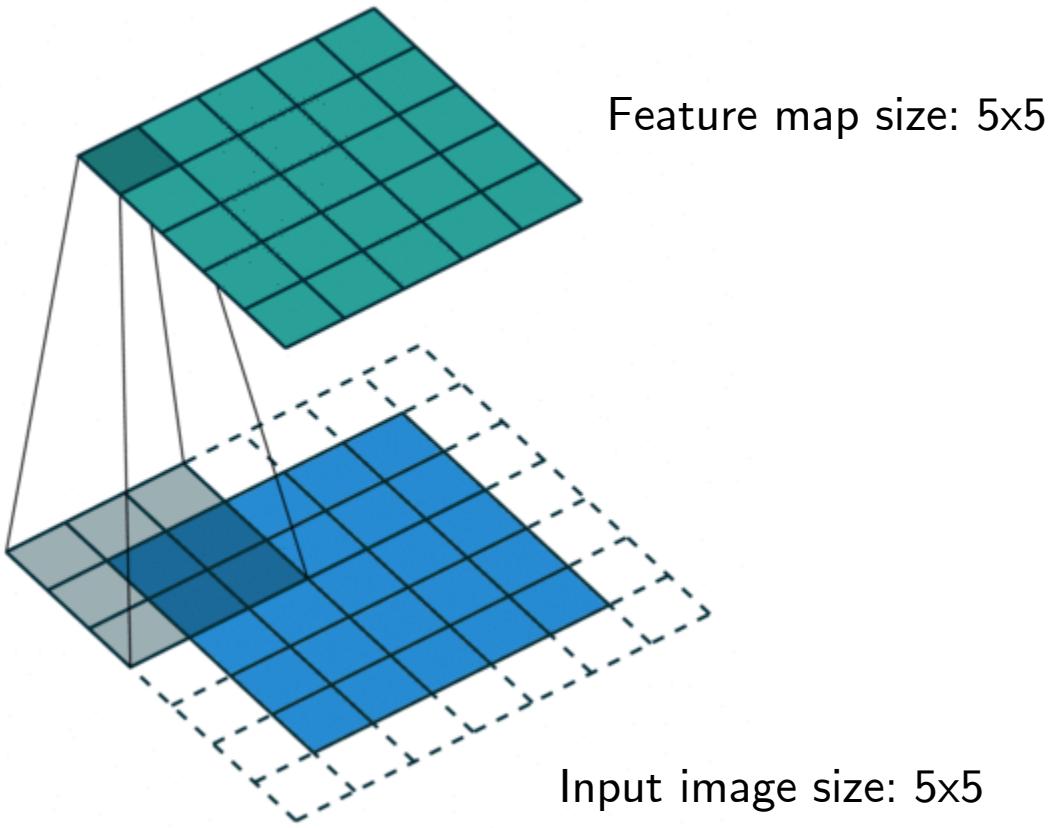
4x4


2x2

0	0	0	0	0	0
0	35	19	25	6	0
0	13	22	16	53	0
0	4	3	7	10	0
0	9	8	1	3	0
0	0	0	0	0	0

Pad Zeros and then convolve to obtain a feature map with dimension = input image dimension

# The Convolution Operation : Zero Padding



Source: [Intuitively Understanding Convolutions for Deep Learning | by Irhum Shafkat | Towards Data Science](#)

# The Convolution Operation: Summary

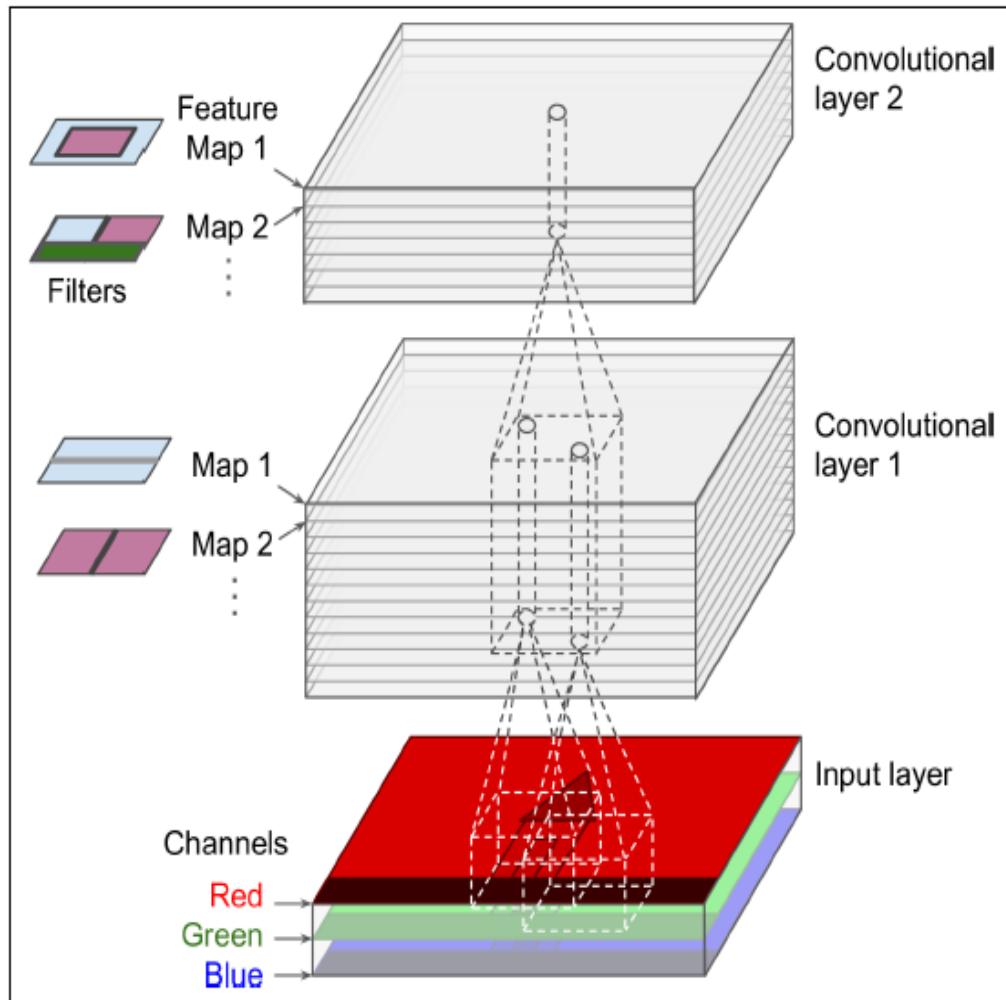
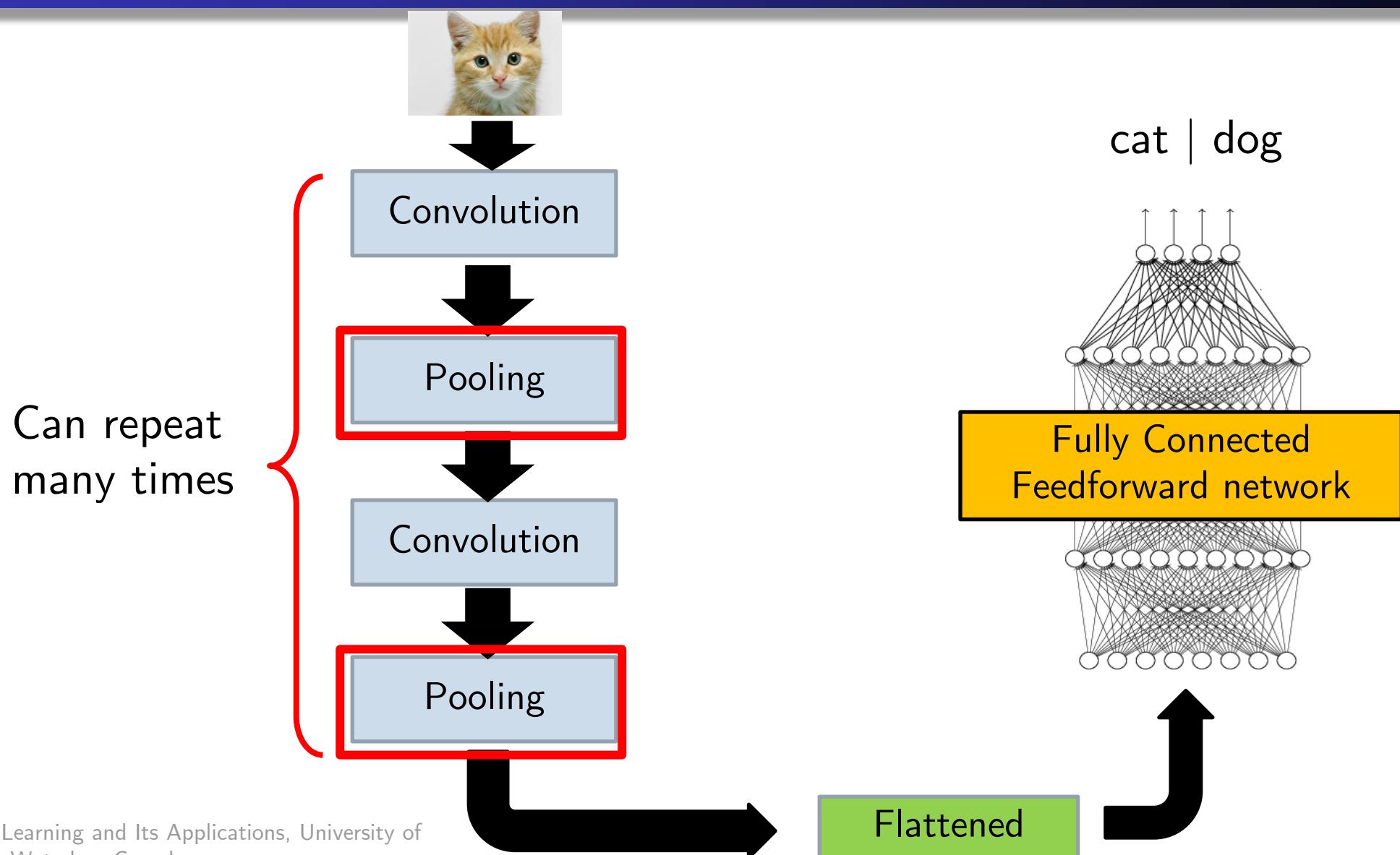


Figure 14-6. Convolution layers with multiple feature maps, and images with three color channels

- Convolutional Layer simultaneously applies multiple trainable filters to its inputs, thus detecting multiple features
  - Each convolutional layer has multiple filters & outputs 1 feature map per filter
  - 1 neuron per pixel in each feature map
  - all neurons share same parameters (weights and biases)
- Padding
  - Same- uses 0 padding if necessary
  - Valid – no 0 padding , could ignore some rows /cols at the bottom right

```
conv = keras.layers.Conv2D(filters=32, kernel_size=3, strides=1,  
                           padding="SAME", activation="relu")
```

# Convolutional Neural Network (CNN) : At a glance



Source: CS 898: Deep Learning and Its Applications, University of Waterloo, Canada.

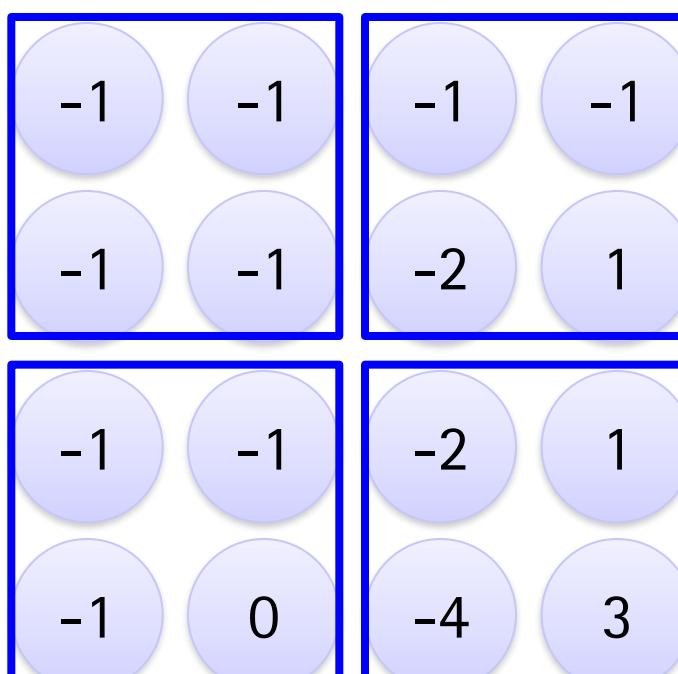
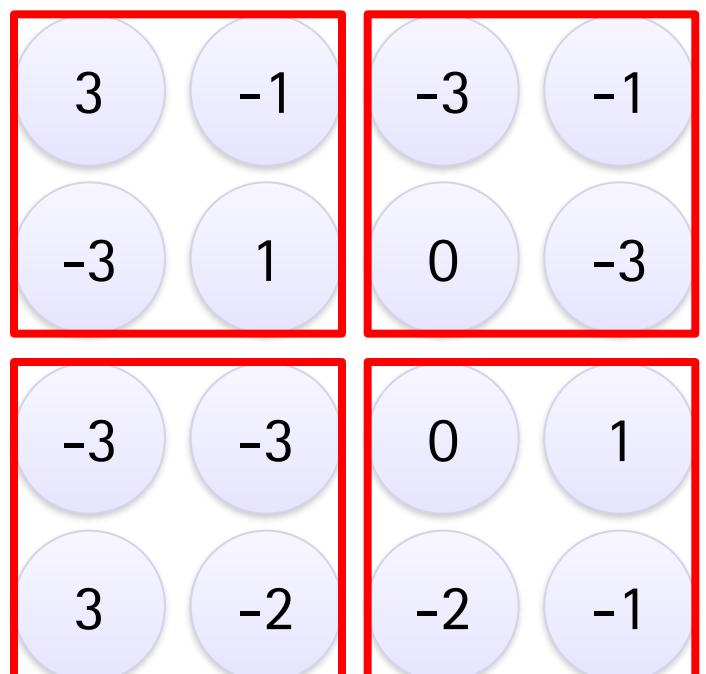
# Pooling

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

Filter 1

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

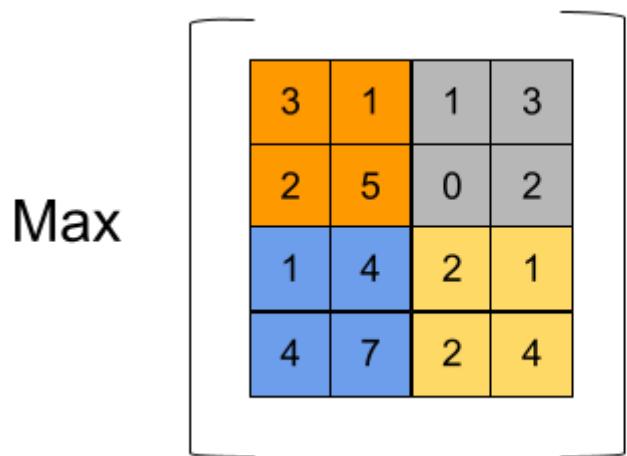
Filter 2



- Max Pooling
- Average Pooling

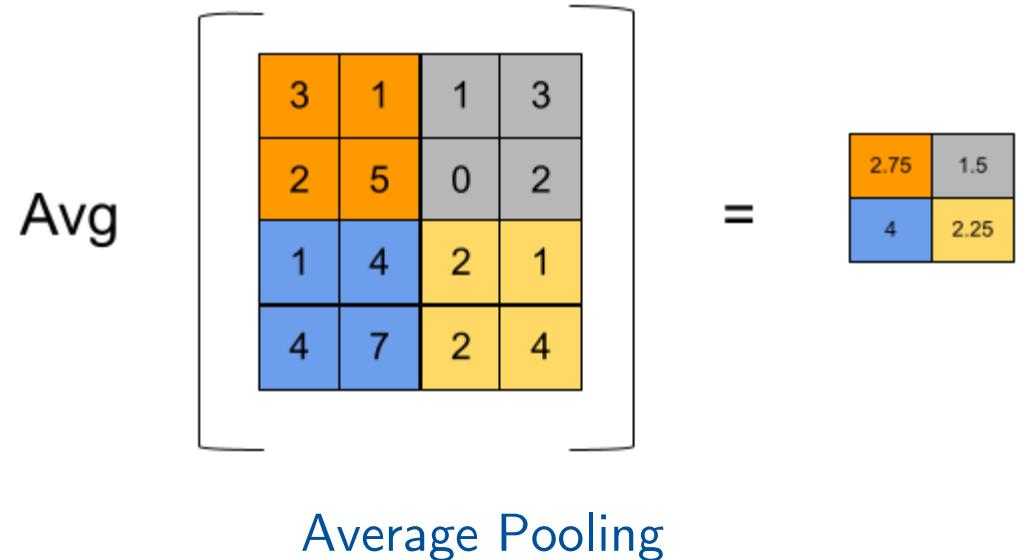
Source: CS 898: Deep Learning and Its Applications, University of Waterloo, Canada.

# Pooling



$$= \begin{matrix} 5 & 3 \\ 7 & 4 \end{matrix}$$

Max. Pooling

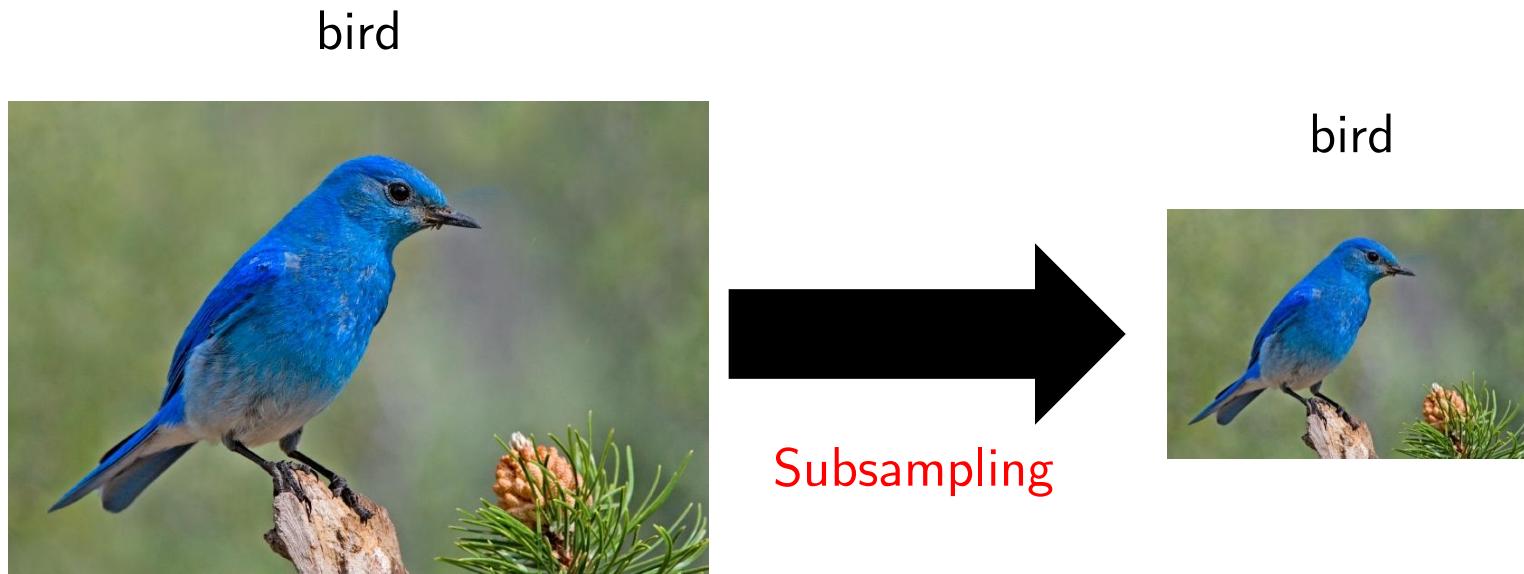


Average Pooling

Stride ?

# Why Pooling ?

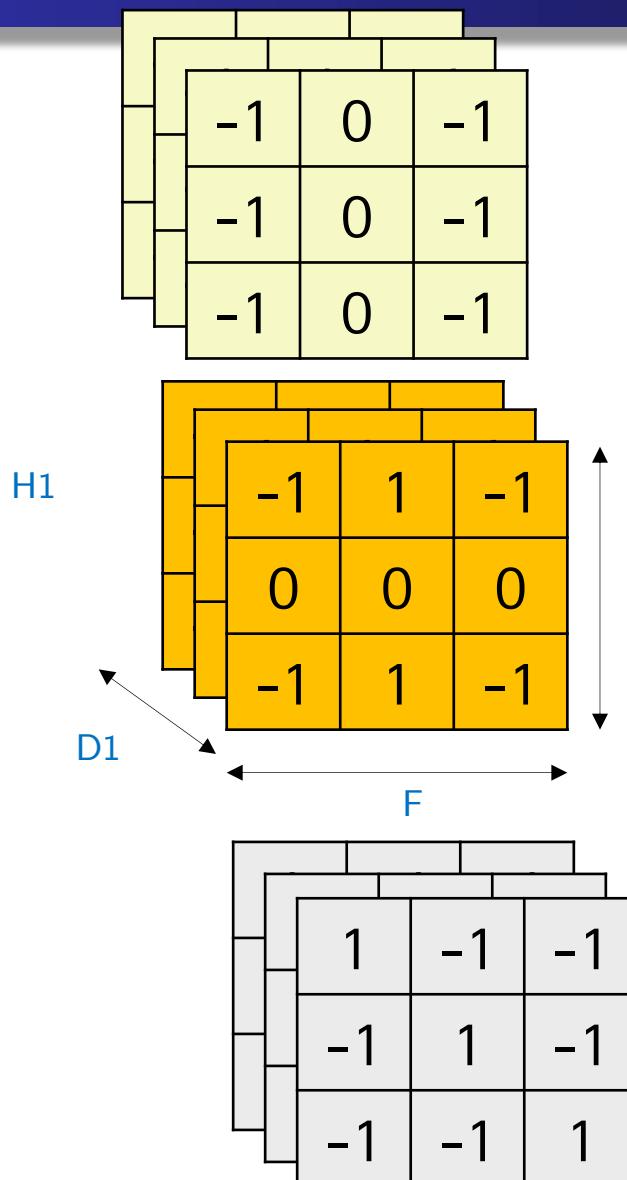
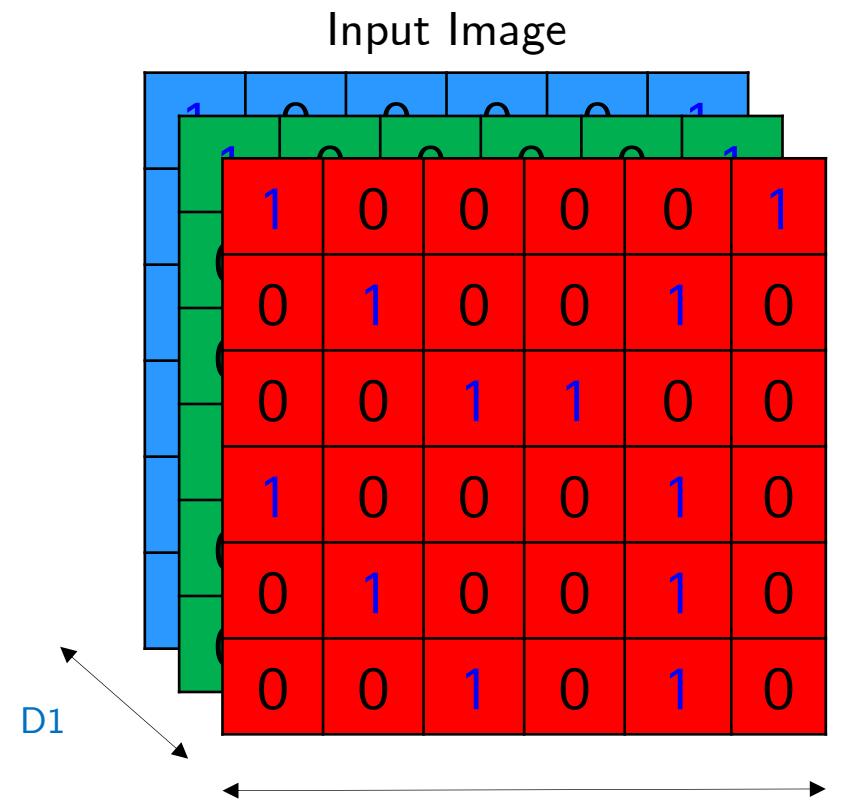
- Subsampling pixels will not change the object



- We can subsample the pixels to make image smaller
- Therefore, fewer parameters to characterize the image

Source: CS 898: Deep Learning and Its Applications, University of Waterloo, Canada.

# Relation between i/p size, feature map size, filter size



$$W_2 = \frac{W_1 - F + 2P}{S} + 1$$

$$H_2 = \frac{H_1 - F + 2P}{S} + 1$$

$$D_2 = K$$

# Connections between layers in a CNN

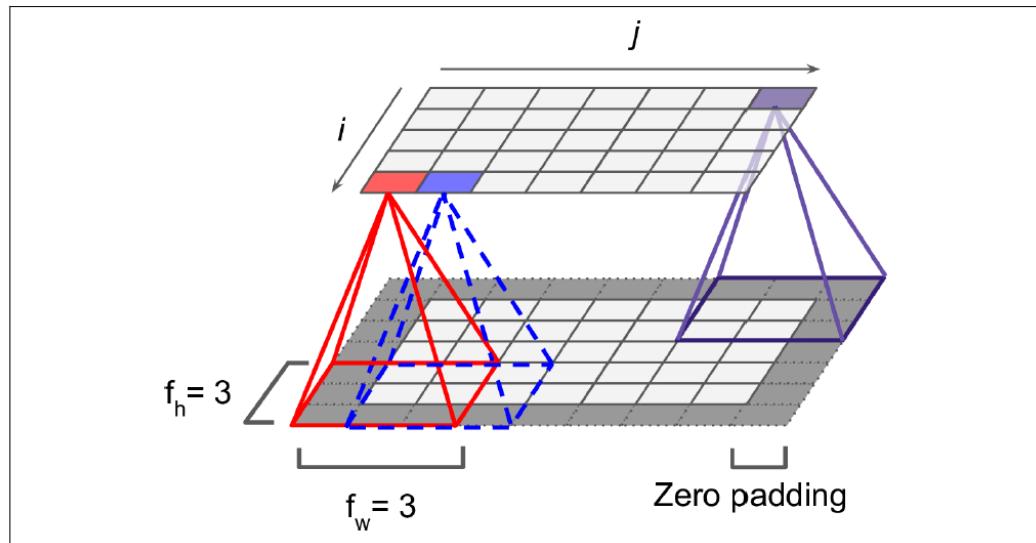


Figure 14-3. Connections between layers and zero padding

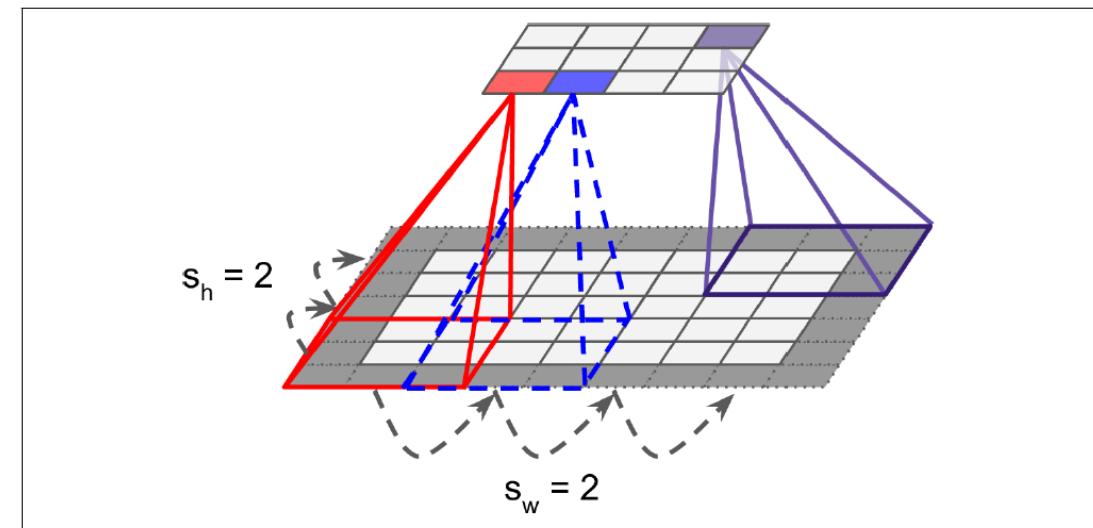


Figure 14-4. Reducing dimensionality using a stride of 2

Source: Aurelien Geron, "Hands-On Machine Learning with Scikit-Learn , Keras & Tensorflow, O'Reilly Publications

A neuron located in row  $i$ , column  $j$  of a given layer is connected to the outputs of the neurons in the previous layer located in

- rows  $i$  to  $i + f_h - 1$ , columns  $j$  to  $j + f_w - 1$ , where  $f_h$  and  $f_w$  are the height and width of the receptive field

The shift from one receptive field to the next is called the stride.

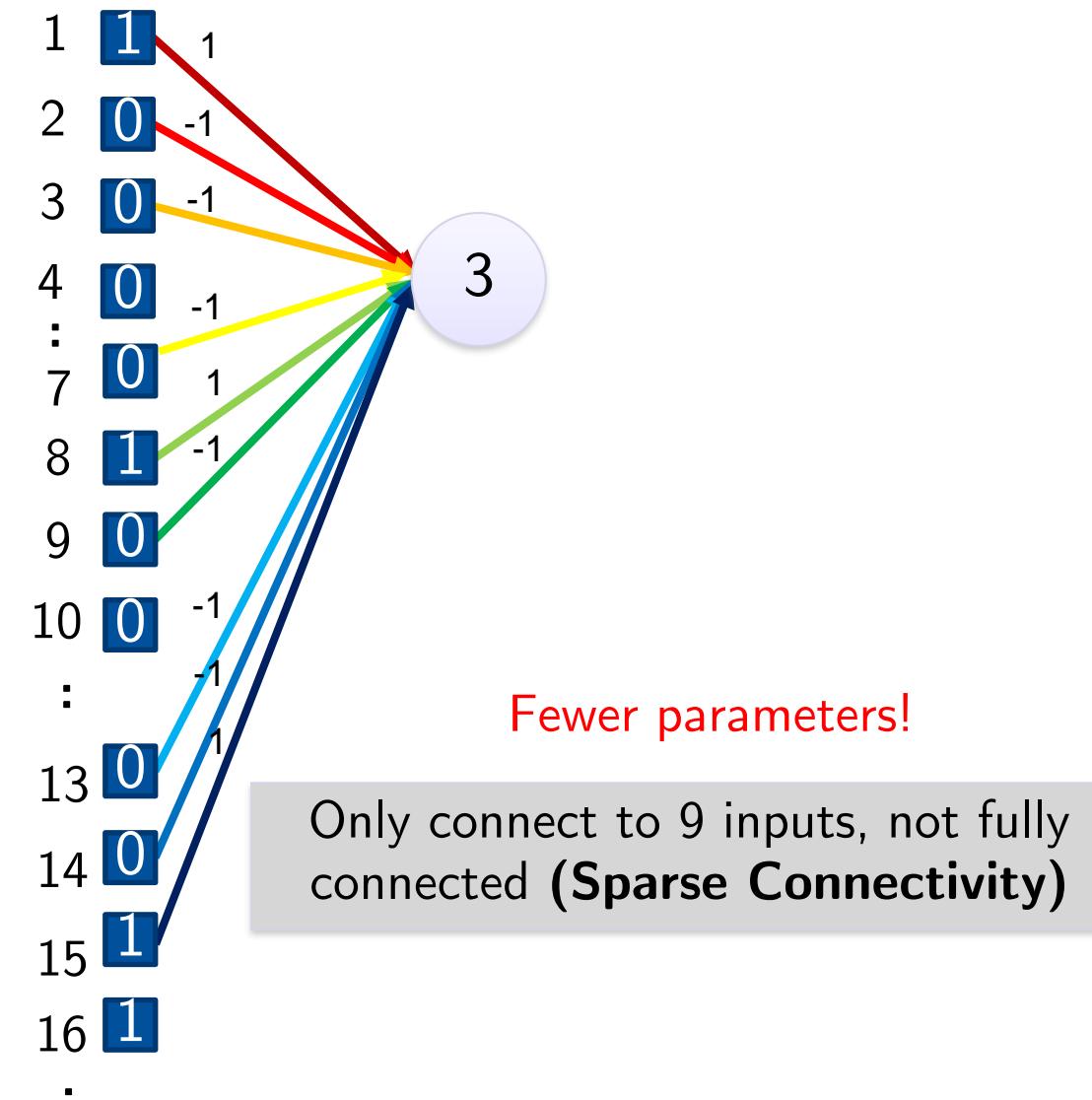
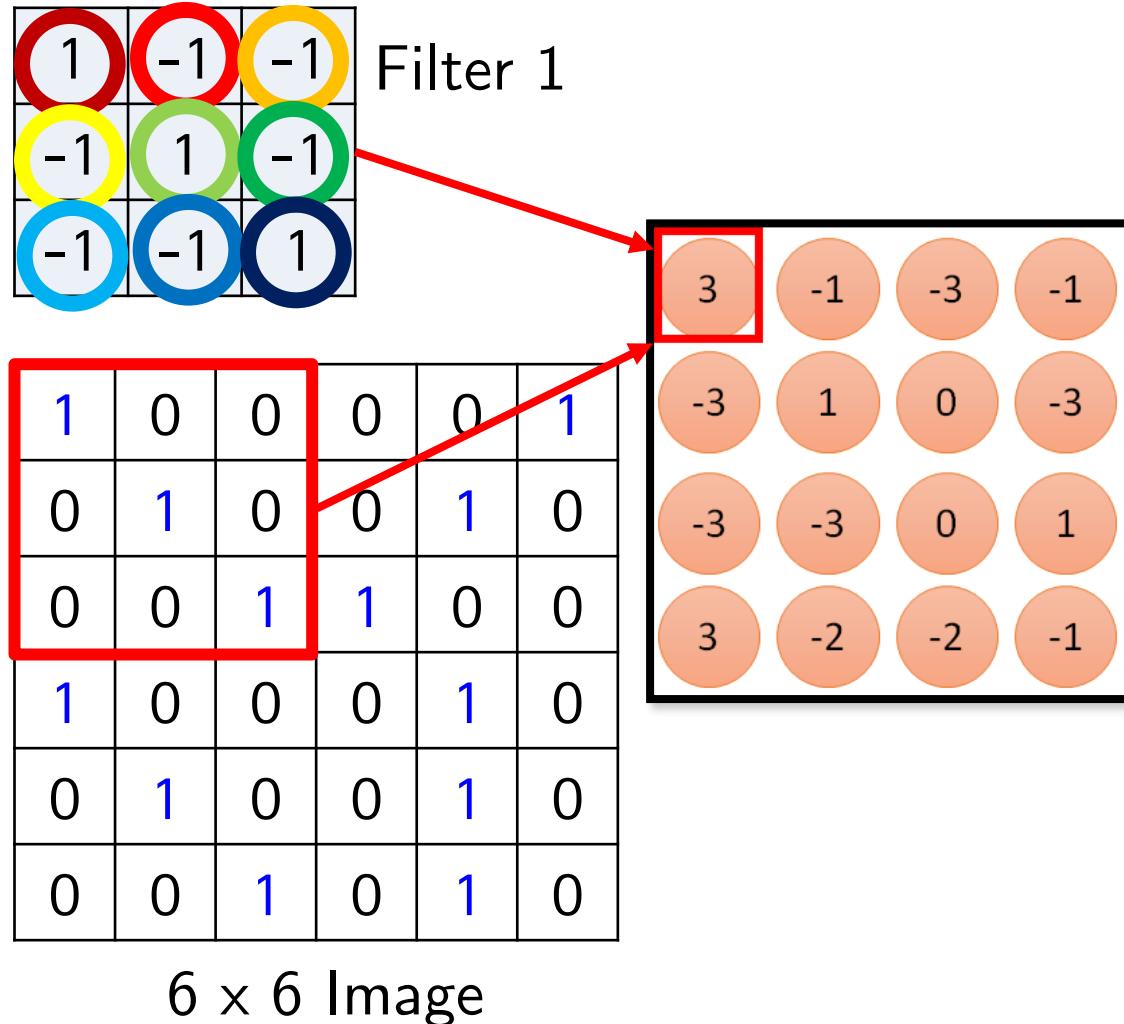
A neuron located in row  $i$ , column  $j$  in the upper layer is connected to the outputs of the neurons in the previous layer located in

- rows  $i \times s_h$  to  $i \times s_h + f_h - 1$ , columns  $j \times s_w$  to  $j \times s_w + f_w - 1$ ,
- where  $s_h$  and  $s_w$  are the vertical and horizontal strides.

# Important properties of CNN

- Sparse Connectivity
- Shared weights
- Equivariant representation

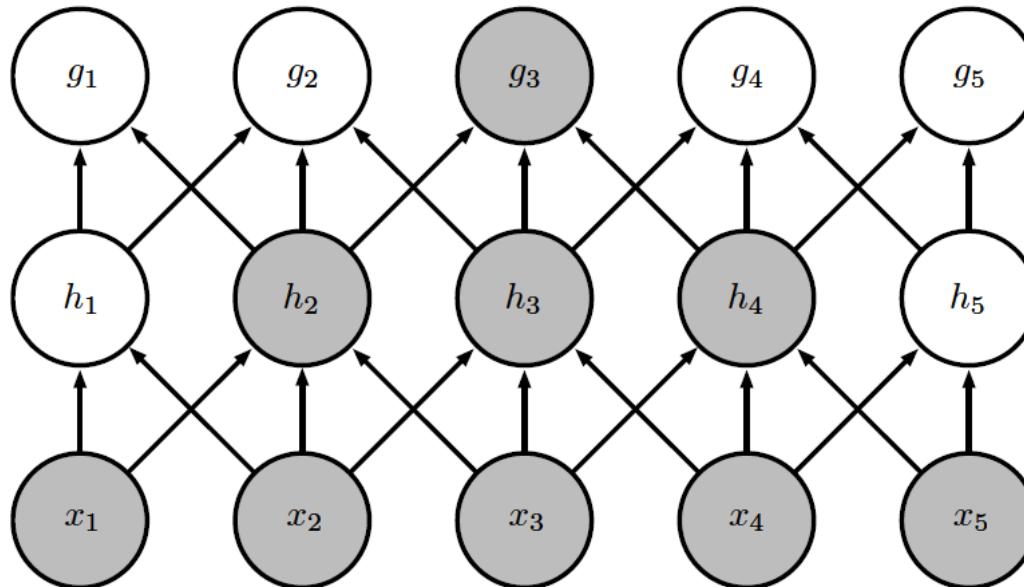
# Properties of CNN



Source: CS 898: Deep Learning and Its Applications, University of Waterloo, Canada.

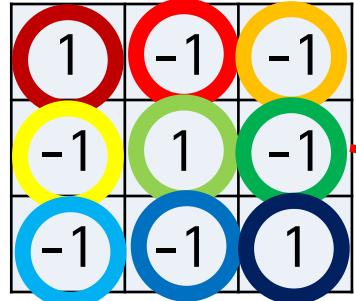
# Properties of CNN

Is sparse connectivity good?



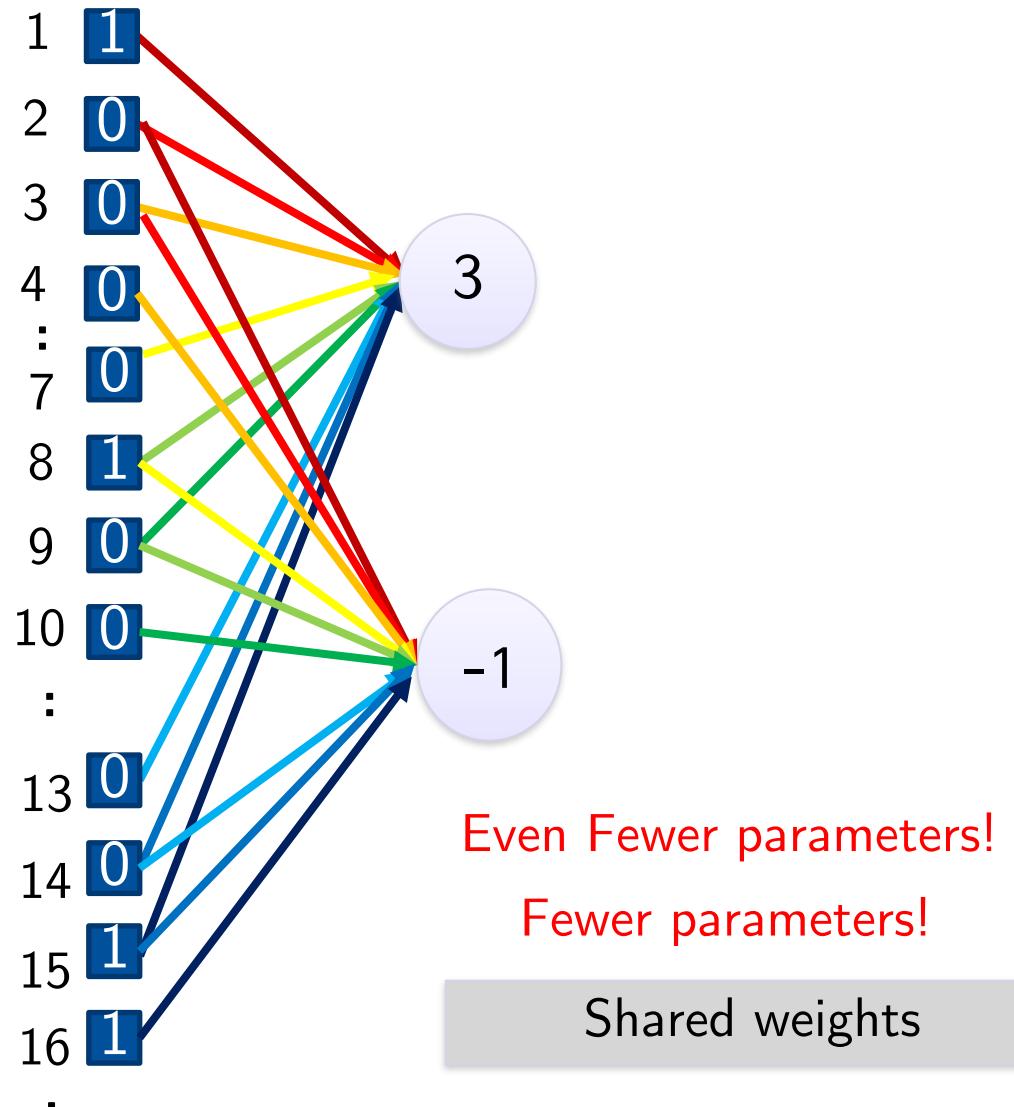
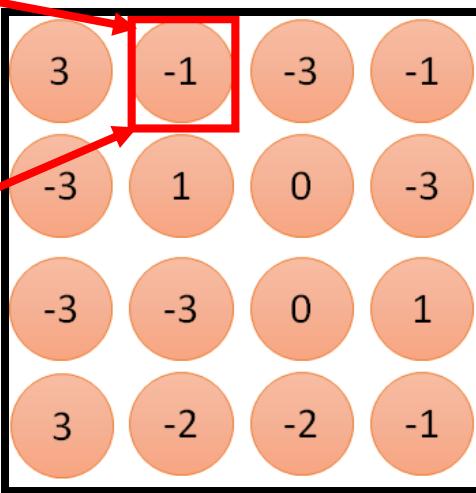
Ian Goodfellow et al. 2016

# Properties of CNN



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

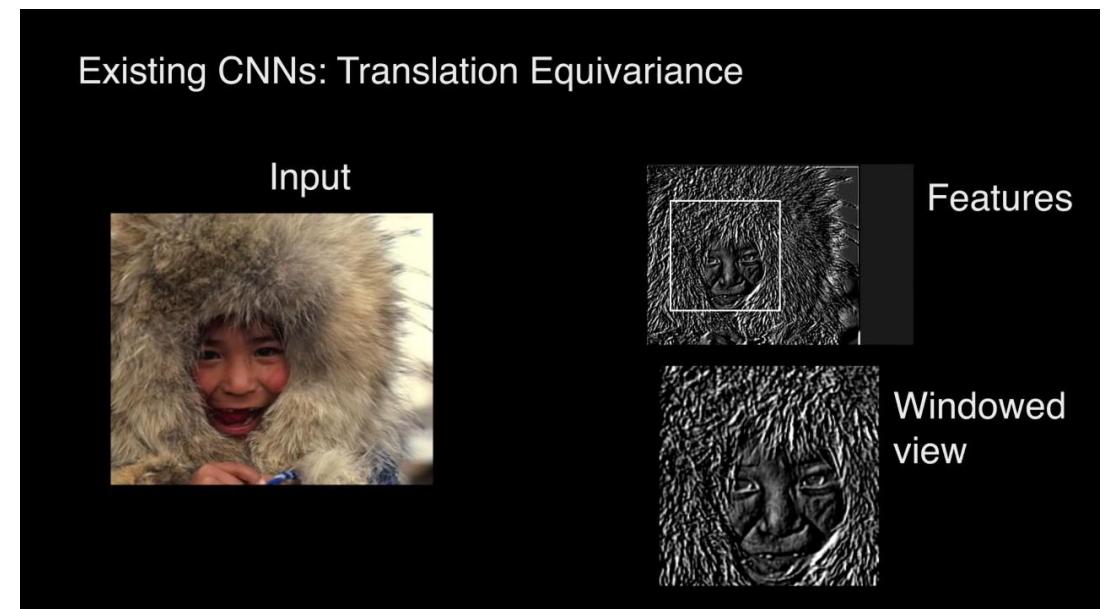
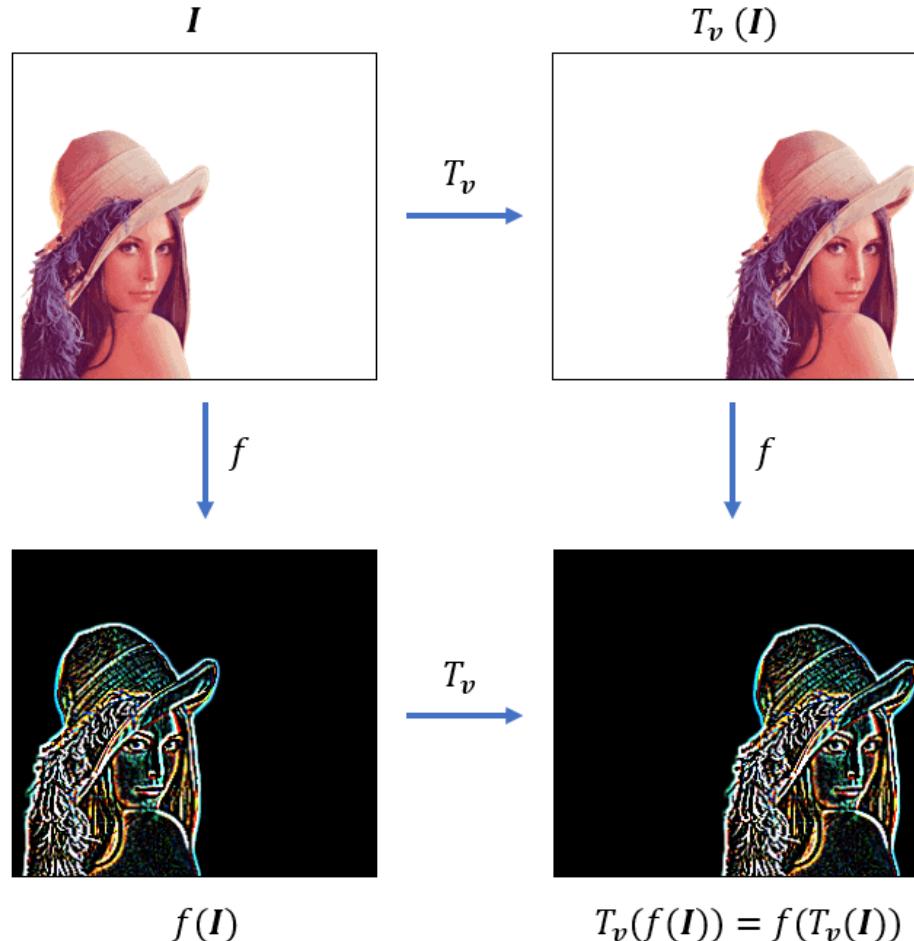
6 x 6 Image



Source: CS 898: Deep Learning and Its Applications, University of Waterloo, Canada.

# Equivariance to translation

- A function  $f$  is equivariant to a function  $g$  if  $f(g(x)) = g(f(x))$  or if the output changes in the same way as the input.

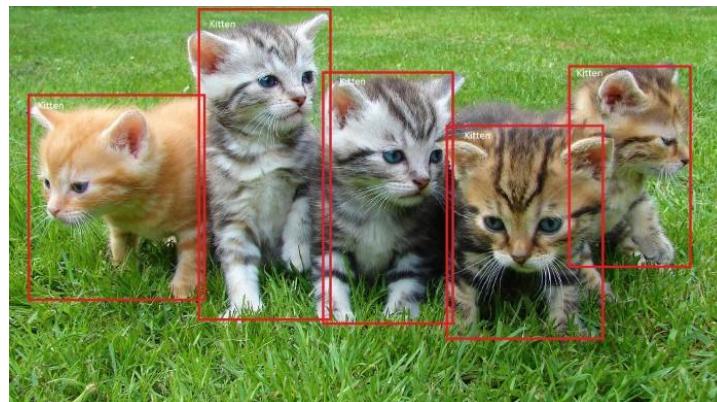


Source: <https://fabianfuchsml.github.io/equivariance1of2/>

Source: [Translation Invariance and Equivariance in Computer Vision | Baeldung on Computer Science](#)

# Equivariance to translation

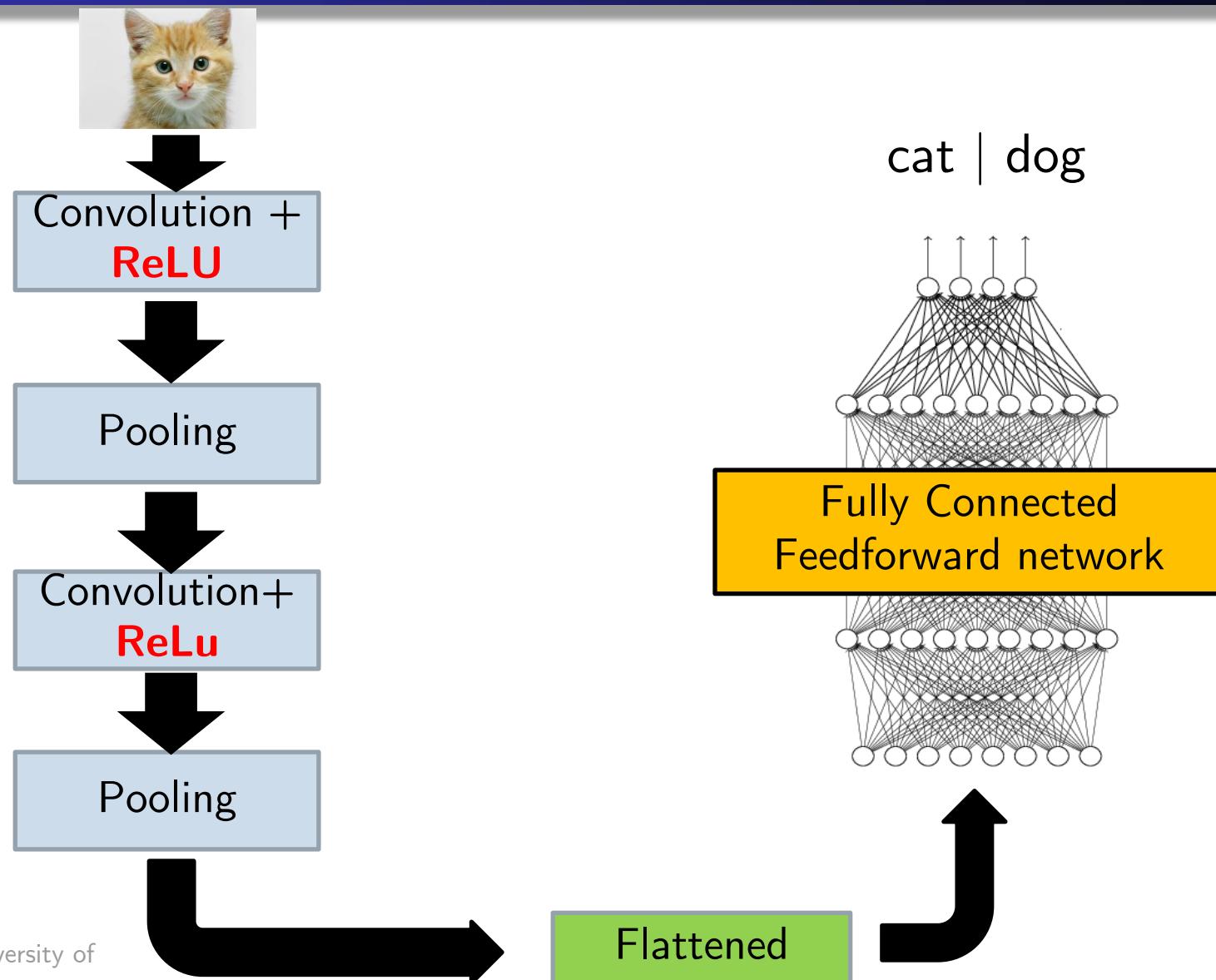
- A function **f** is equivariant to a function **g** if  $f(g(x)) = g(f(x))$  or if the output changes in the same way as the input.
- This is achieved by the concept of weight sharing.
- As the same weights are shared across the images, hence if an object occurs in any image, it will be detected irrespective of its position in the image.



# CNN vs Fully Connected NN

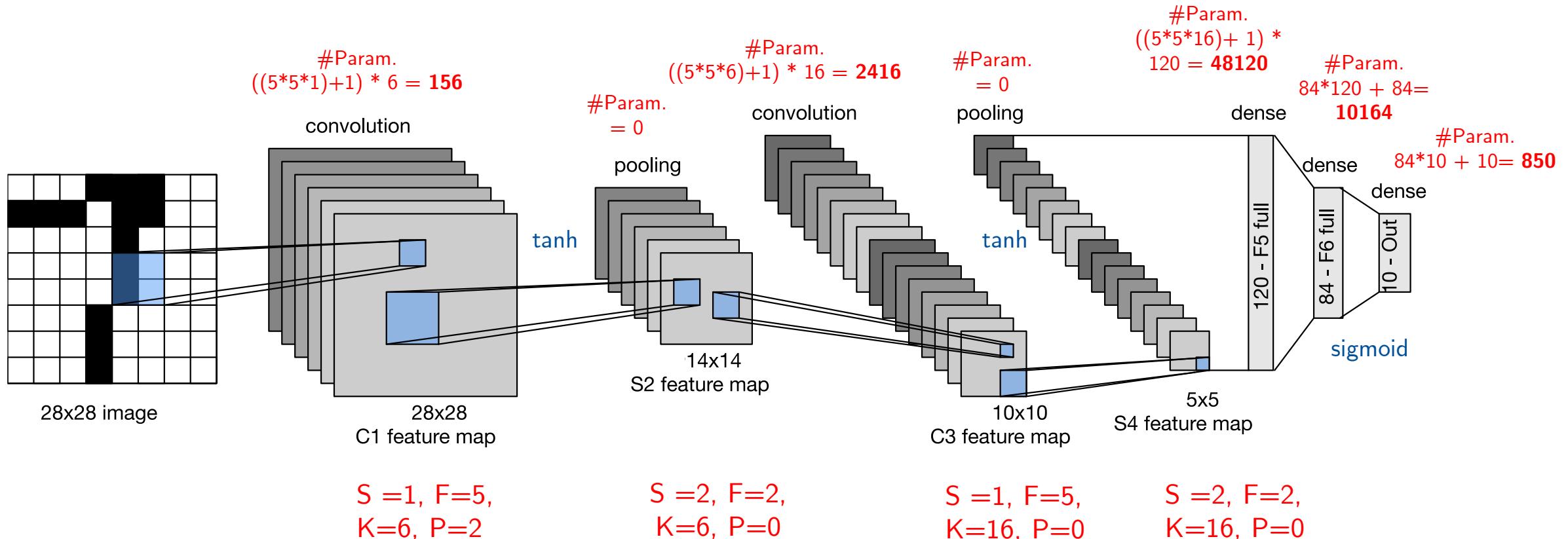
- A CNN compresses the fully connected NN in two ways:
  - Reducing the number of connections
  - Shared weights
- Max pooling further reduces the parameters to characterize an image.
- CNN captures the neighbourhood relationships between pixel.

# Convolutional Neural Network (CNN) : Non-linearity with activation



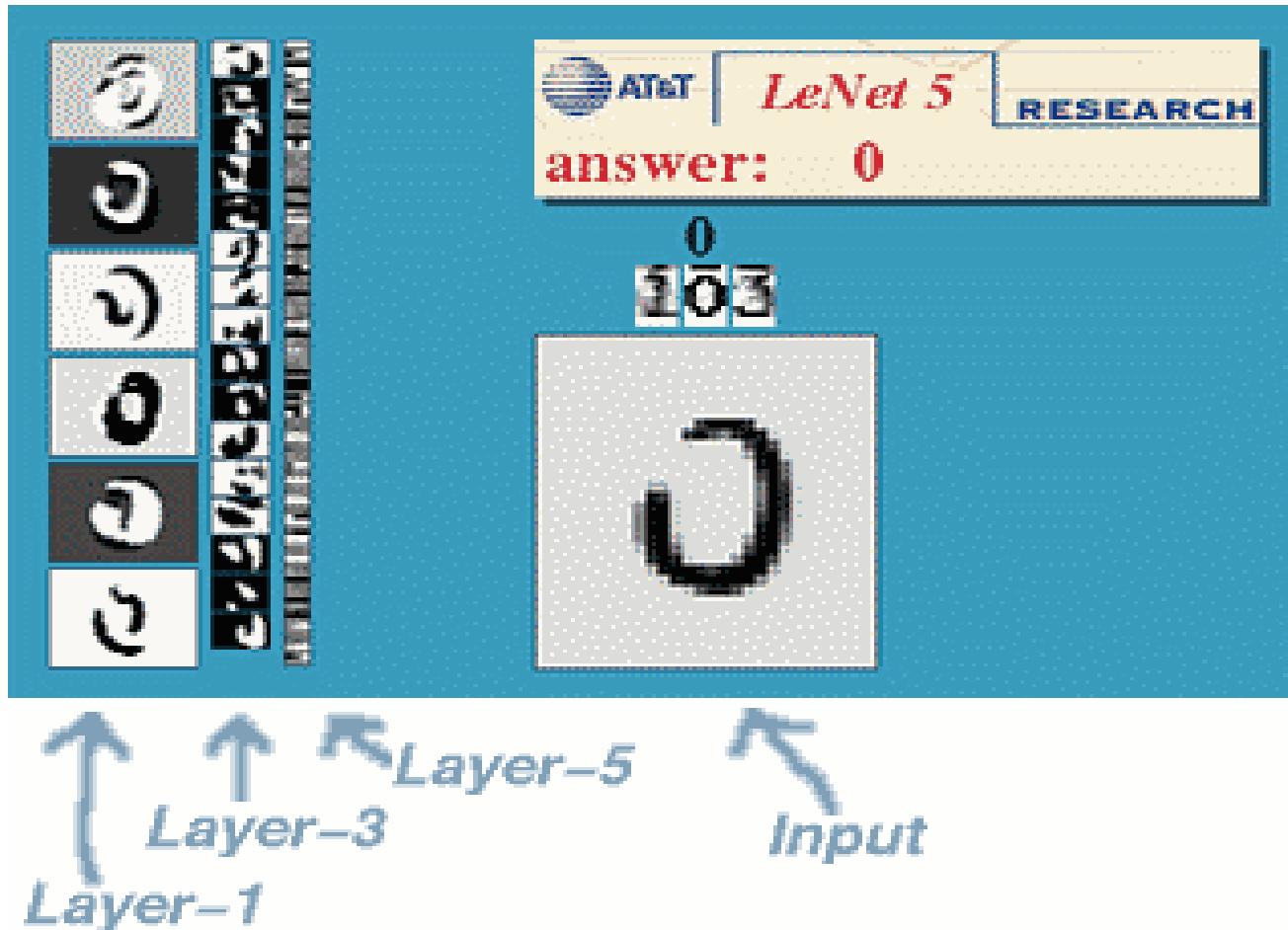
Source: CS 898: Deep Learning and Its Applications, University of Waterloo, Canada.

# LeNet-5 Architecture for handwritten text recognition



LeCun, Y., Bottou, L., Bengio, Y., Haffner, P., & others. (1998). Gradient-based learning applied to document recognition.  
*Proceedings of the IEEE*, 86(11), 2278–2324.

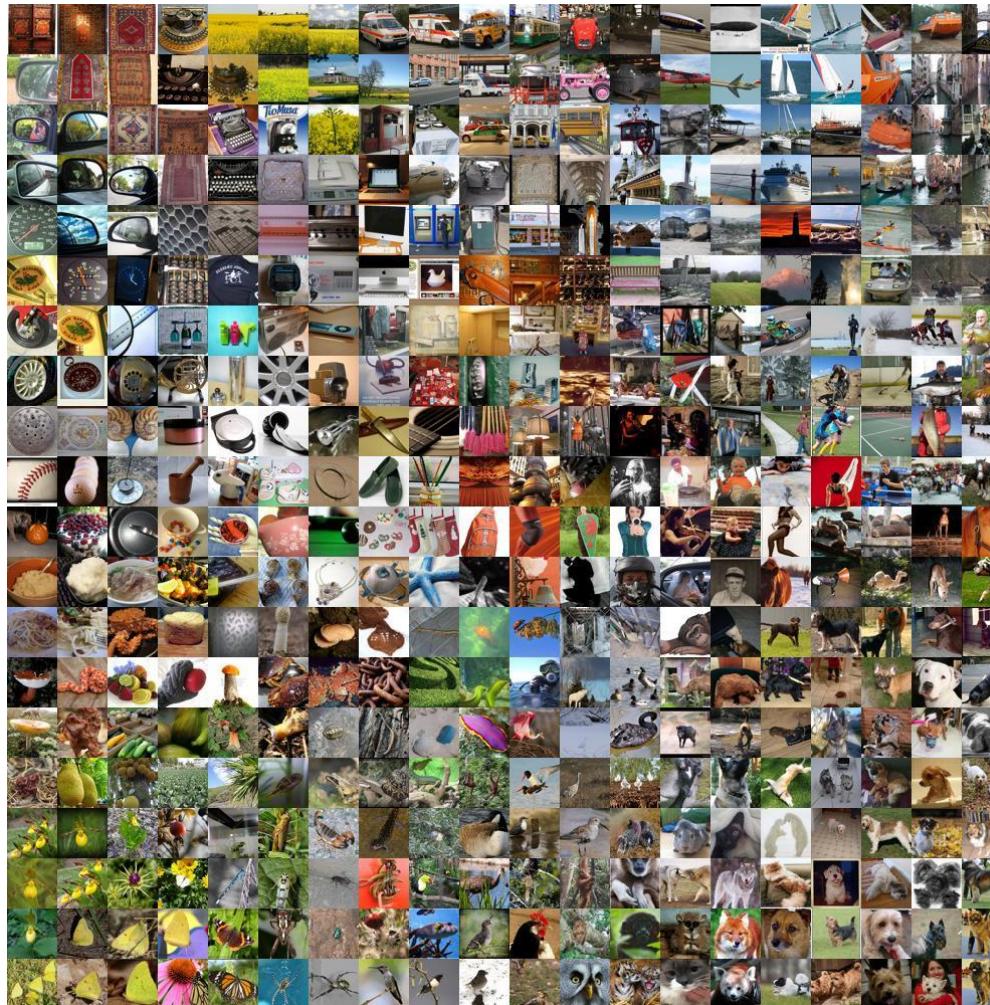
# LeNet-5 Architecture for handwritten number recognition



Source: <http://yann.lecun.com/>

# ImageNet Dataset

More than 14 million images.

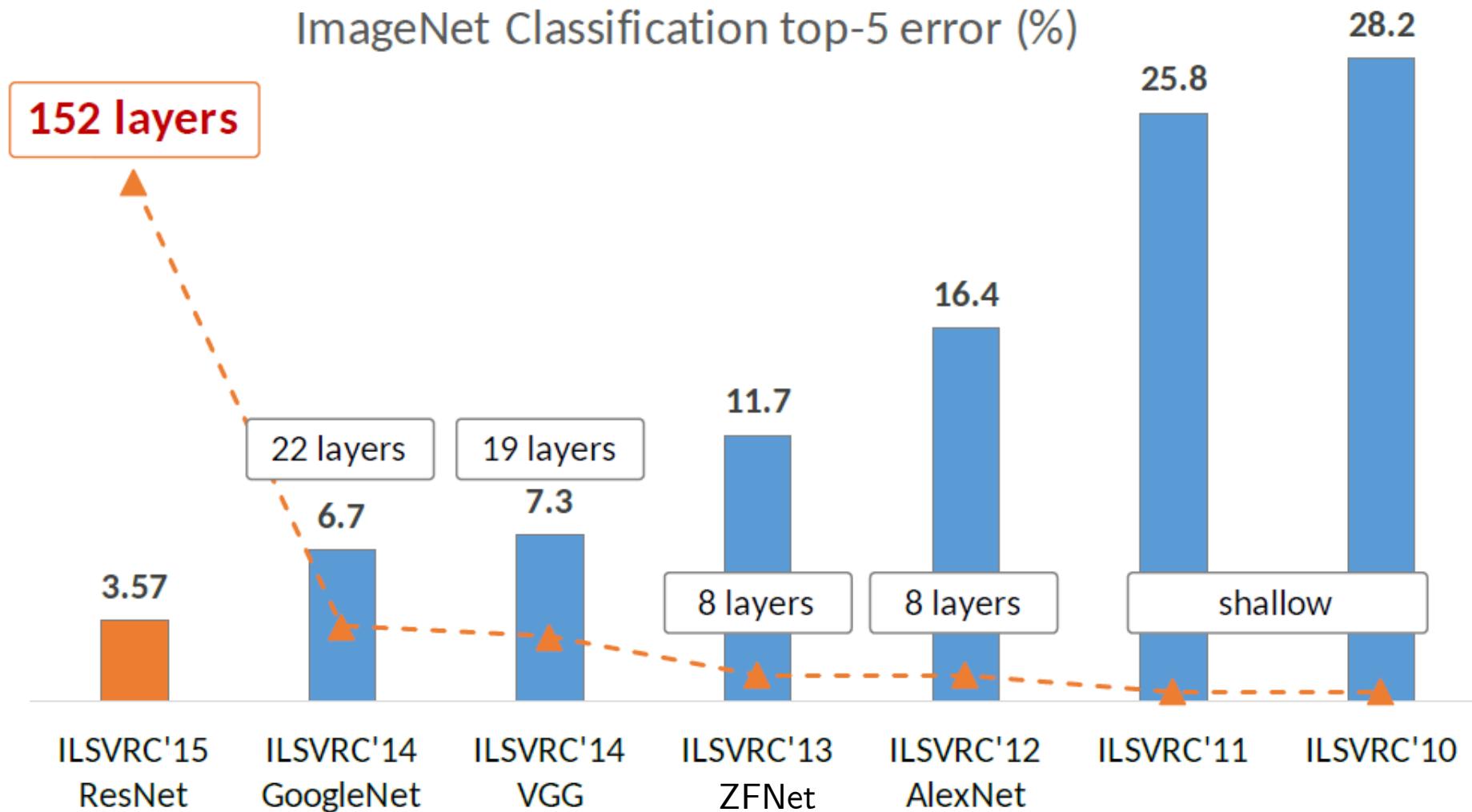


22,000 Image categories

Deng, Jia, et al. "Imagenet: A large-scale hierarchical image database." *IEEE conference on computer vision and pattern recognition*. IEEE, 2009.

# ImageNet Large Scale Visual Recognition Challenge

- 1000 ImageNet Categories



## ImageNet Classification with Deep Convolutional Neural Networks

**Alex Krizhevsky**  
University of Toronto  
[kriz@cs.utoronto.ca](mailto:kriz@cs.utoronto.ca)

**Ilya Sutskever**  
University of Toronto  
[ilya@cs.utoronto.ca](mailto:ilya@cs.utoronto.ca)

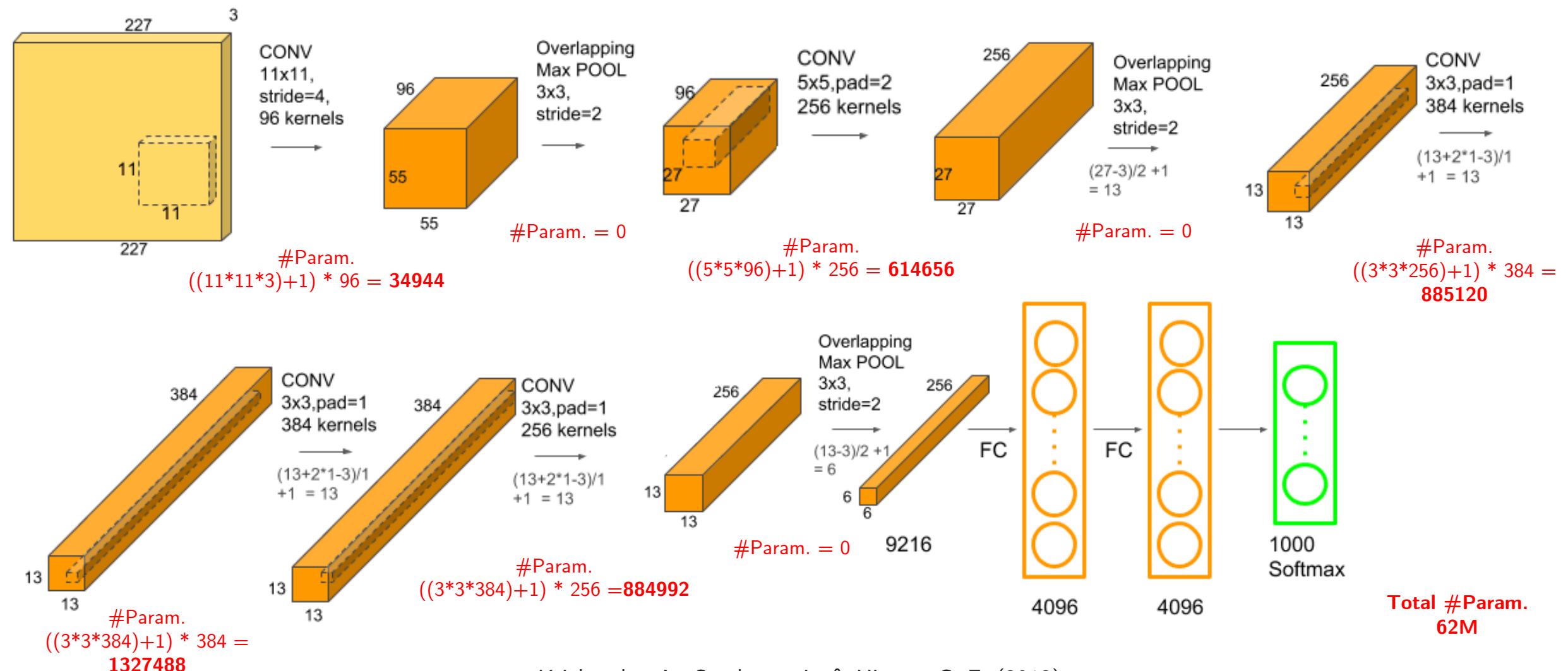
**Geoffrey E. Hinton**  
University of Toronto  
[hinton@cs.utoronto.ca](mailto:hinton@cs.utoronto.ca)

### Abstract

We trained a large, deep convolutional neural network to classify the 1.2 million high-resolution images in the ImageNet LSVRC-2010 contest into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 37.5% and 17.0% which is considerably better than the previous state-of-the-art. The neural network, which has 60 million parameters and 650,000 neurons, consists of five convolutional layers, some of which are followed by max-pooling layers, and three fully-connected layers with a final 1000-way softmax. To make training faster, we used non-saturating neurons and a very efficient GPU implementation of the convolution operation. To reduce overfitting in the fully-connected layers we employed a recently-developed regularization method called “dropout” that proved to be very effective. We also entered a variant of this model in the ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3%, compared to 26.2% achieved by the second-best entry.

- Used **ReLU** activation function instead of sigmoid and tanh.
- Used **data augmentation** techniques that consisted of image translations, horizontal reflections, and patch extractions.
- Implemented **dropout** layers.

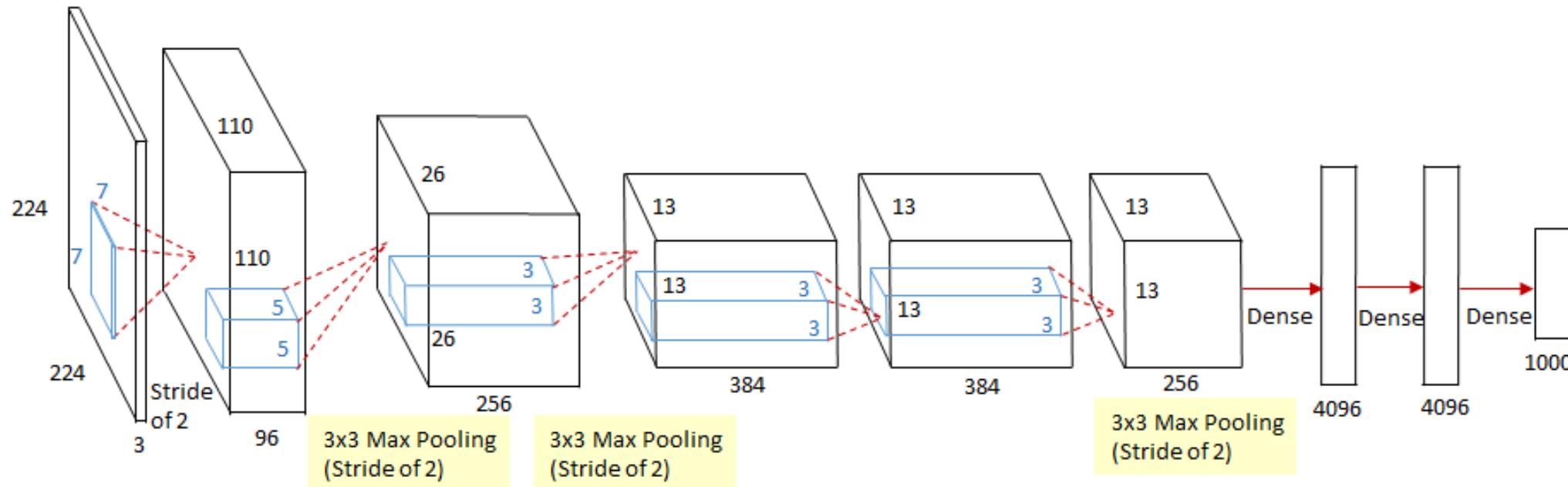
# AlexNet Architecture



Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012).

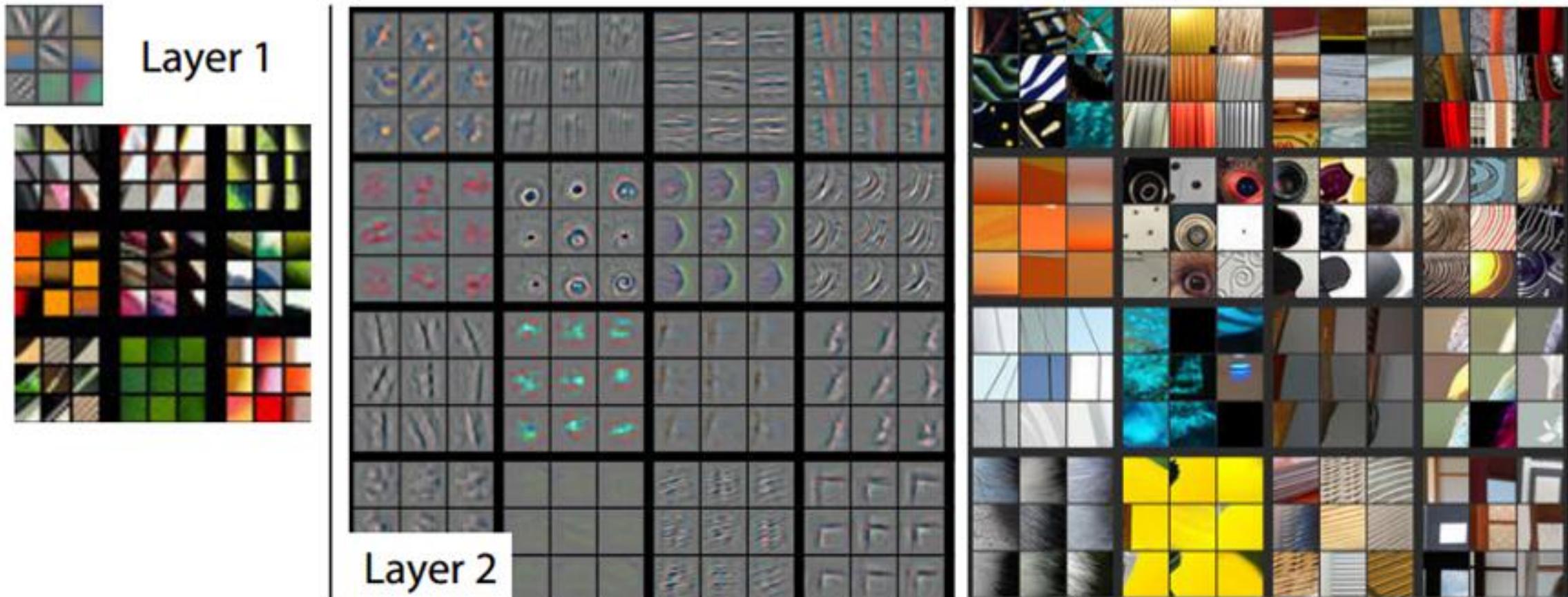
Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25.

# ZFNet Architecture (2013)



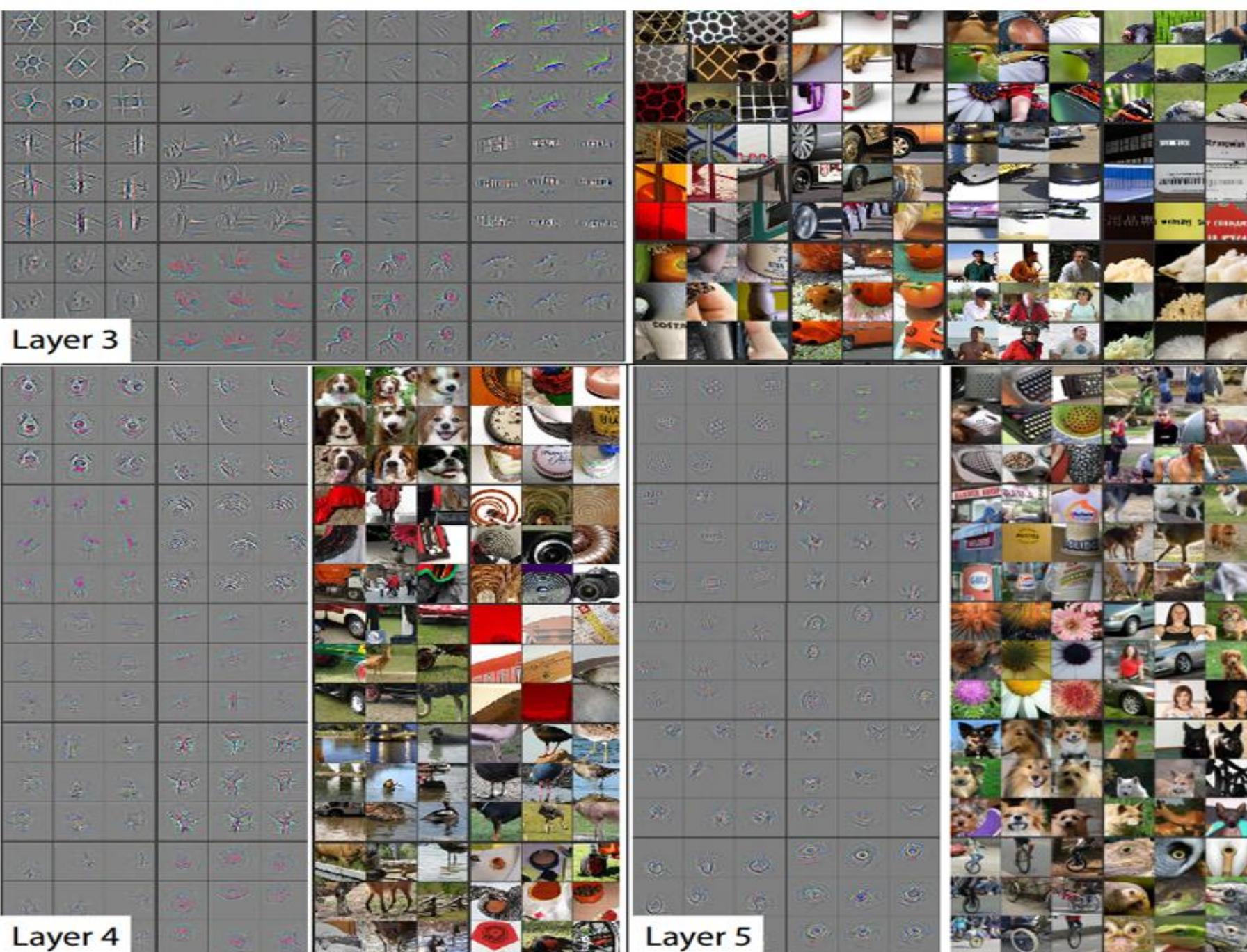
- Used filters of size 7x7 instead of 11x11 in AlexNet
- Used Deconvnet to visualize the intermediate results.

Zeiler, M. D., & Fergus, R. (2013). Visualizing and understanding convolutional networks. In *European conference on computer vision* (pp. 818-833). Springer, Cham.



Visualizations of Layer 1 and 2. Each layer illustrates 2 pictures, one which shows the filters themselves and one that shows what part of the image are most strongly activated by the given filter. For example, in the space labeled Layer 2, we have representations of the 16 different filters (on the left)

[Visualizing and Understanding Deep Neural Networks by Matt Zeiler - YouTube](#)



Visualizations of Layers 3, 4, and 5

# VGGNet Architecture (2014)

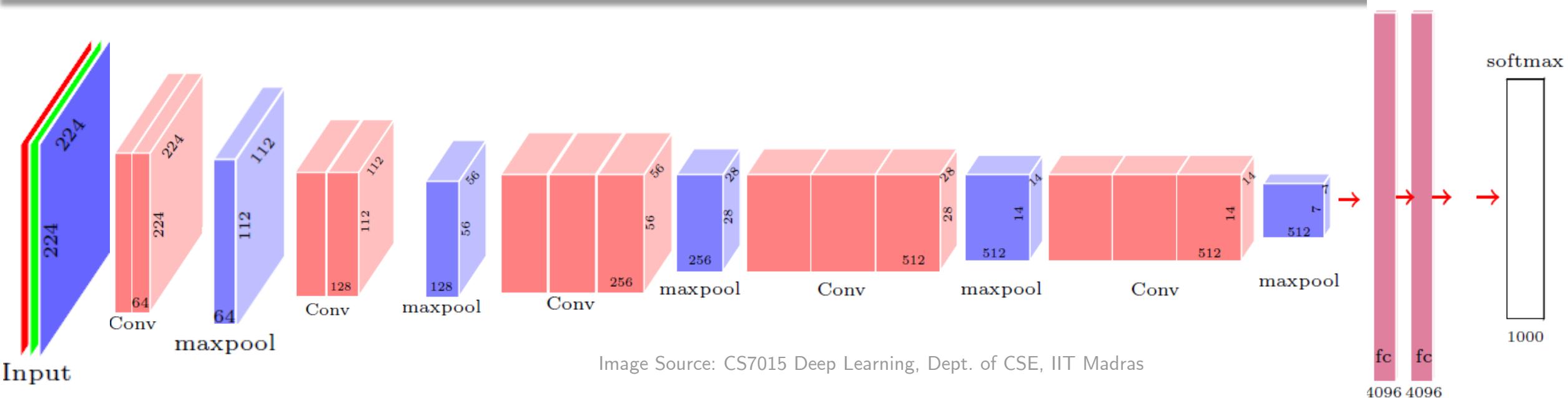
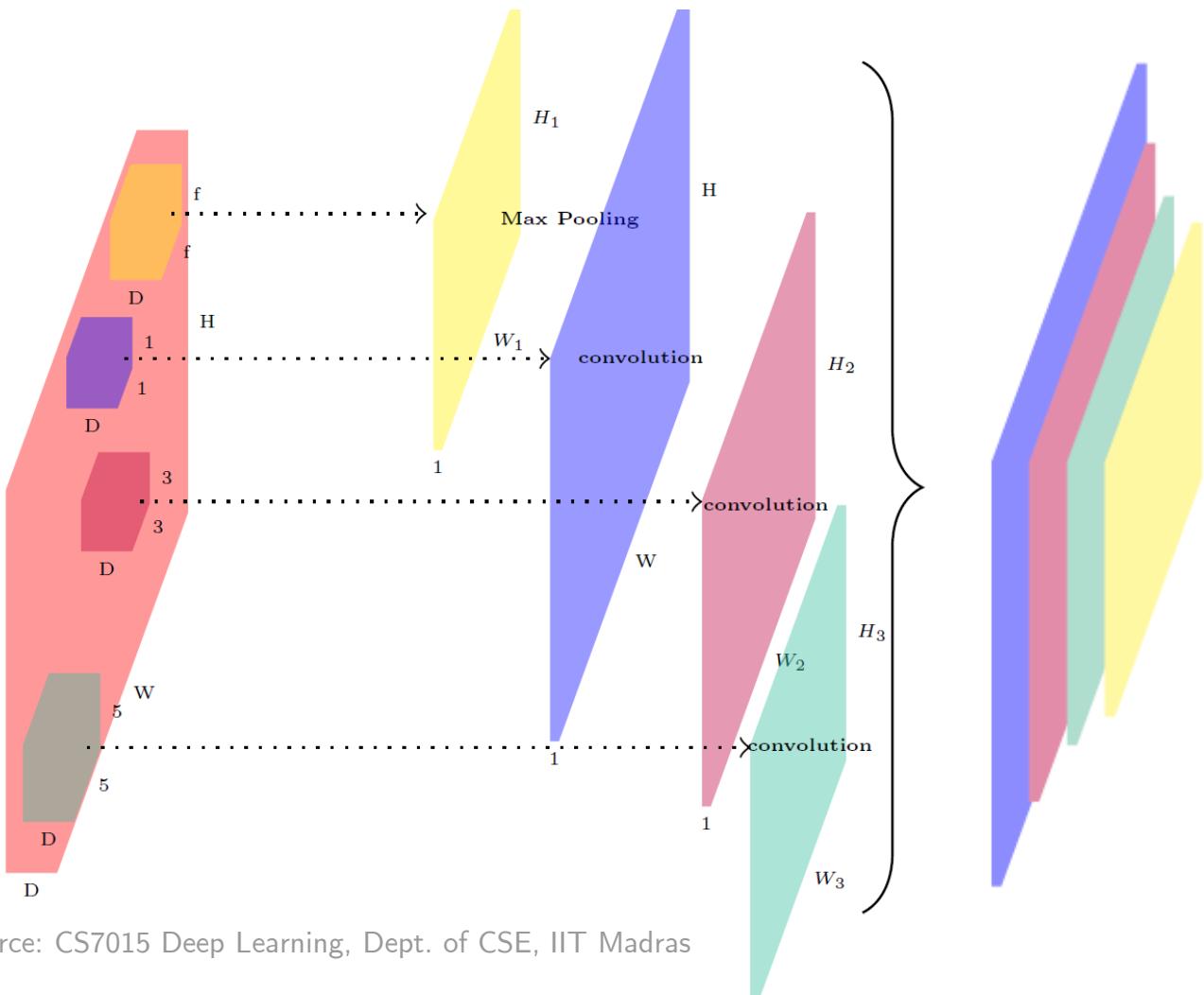


Image Source: CS7015 Deep Learning, Dept. of CSE, IIT Madras

- Used filters of size 3x3 in all the convolution layers.
- 3 conv layers back-to-back have an effective receptive field of 7x7.
- Also called VGG-16 as it has 16 layers.
- This work reinforced the notion that convolutional neural networks have to have a deep network of layers in order for this hierarchical representation of visual data to work

Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition , International Conference on Learning Representations (ICLR14)

# GoogleNet Architecture (2014)

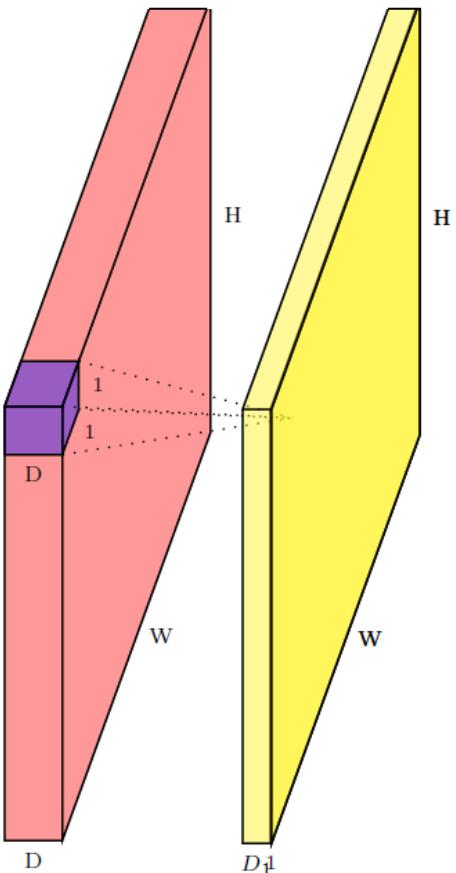


Source: CS7015 Deep Learning, Dept. of CSE, IIT Madras

- Most of the architectures discussed till now apply either of the following after each convolution operation:
  - Max Pooling
  - 3x3 convolution
  - 5x5 convolution
- Idea: Why cant we apply them all together at the same time and concatenate the feature maps.
- Problem: This will result in large number of computations.
- Specifically, each element of the output required  **$O(F \times F \times D)$**  computations

Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., ... & Rabinovich, A. (2015). Going deeper with convolutions.  
In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR'15)*

# GoogleNet Architecture (2014)

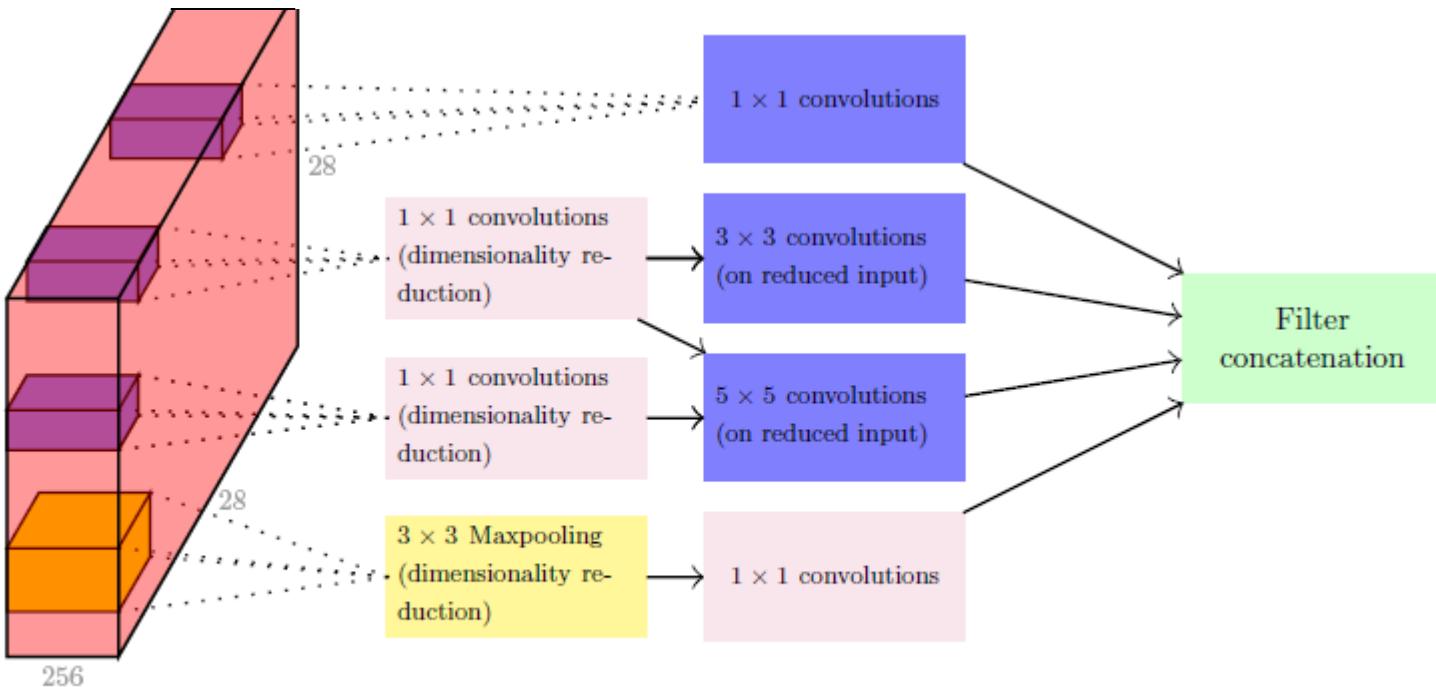


- Solution: Apply 1x1 convolutions
- 1x1 convolution aggregates along the depth allowing dimensionality reduction across depth.
- So, if we apply  $D_1$  1x1 convolutions ( $D_1 < D$ ), we will get an output of size  $W \times H \times D_1$
- So, the total number of computations will reduce to  **$O(F \times F \times D_1)$**
- We could then apply subsequent 3x3, 5x5 filters on this reduced output

Source: CS7015 Deep Learning, Dept. of CSE, IIT Madras

Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., ... & Rabinovich, A. (2015). Going deeper with convolutions.  
In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR'15)*

# GoogleNet Architecture (2014)



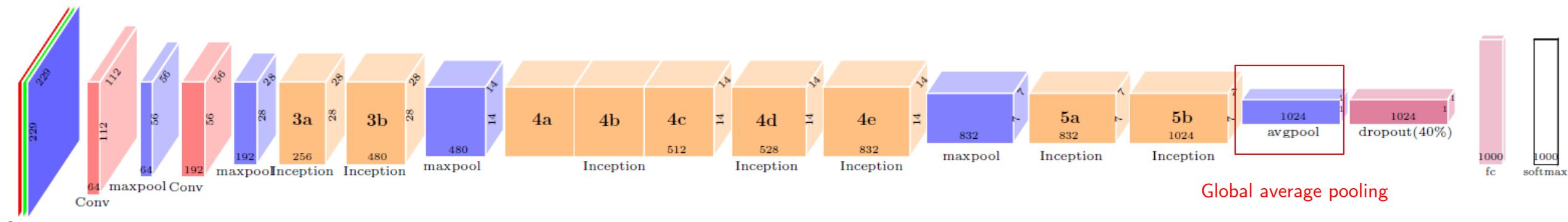
The Inception module

- Also, we might want to use different dimensionality reductions (applying 1x1 convolutions of different sizes) before the 3x3 and 5x5 filters.
- We can also add the maxpooling layer followed by 1x1 convolution.
- After this, we concatenate all these layers.
- This is called the **Inception module**.
- **GoogleNet** contains many such inception modules.

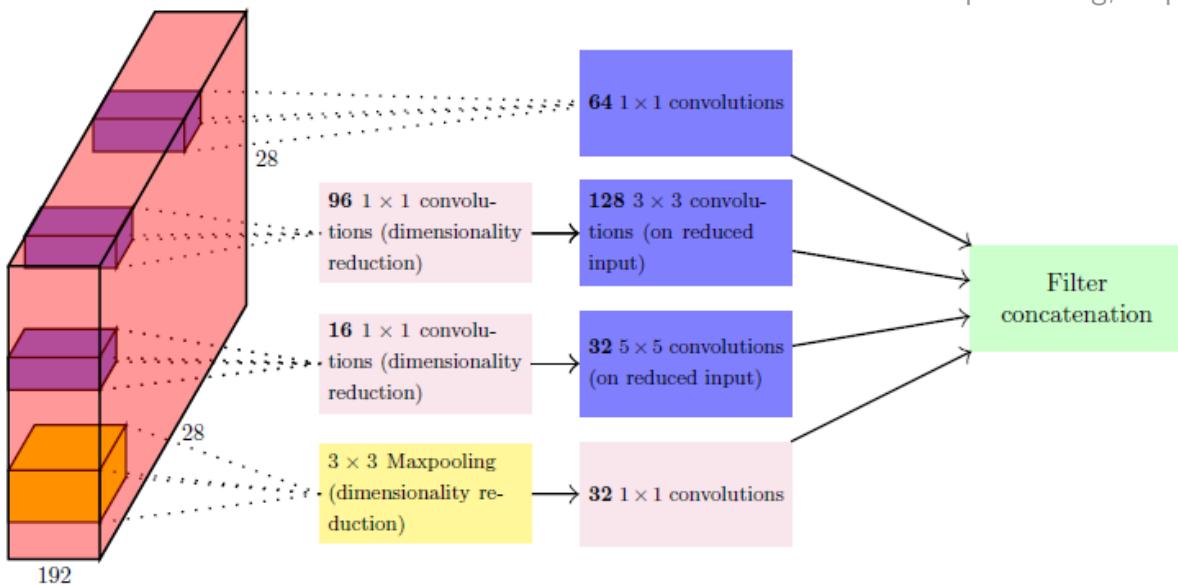
Source: CS7015 Deep Learning, Dept. of CSE, IIT Madras

Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., ... & Rabinovich, A. (2015). Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR'15)*

# GoogleNet Architecture (2014)



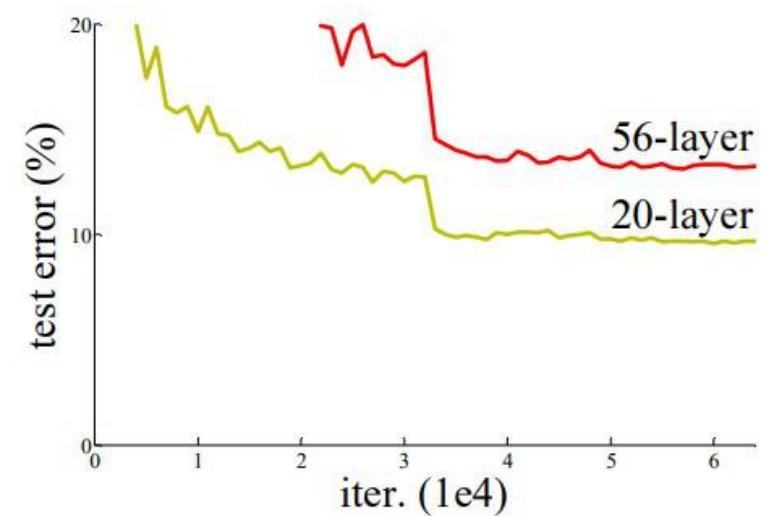
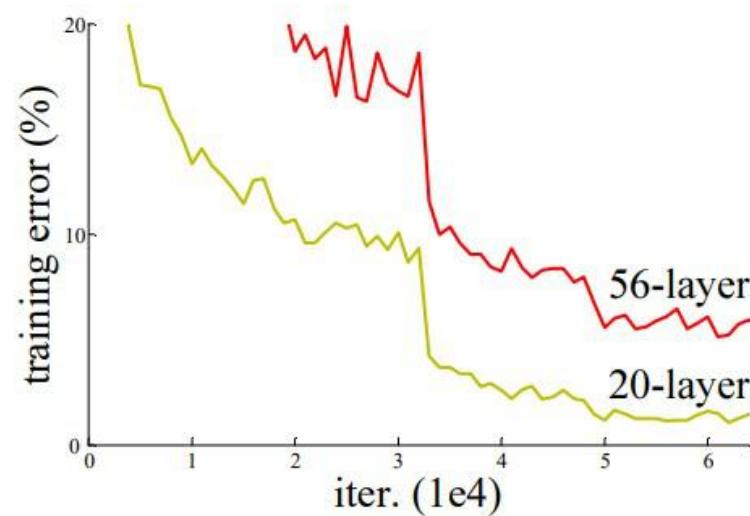
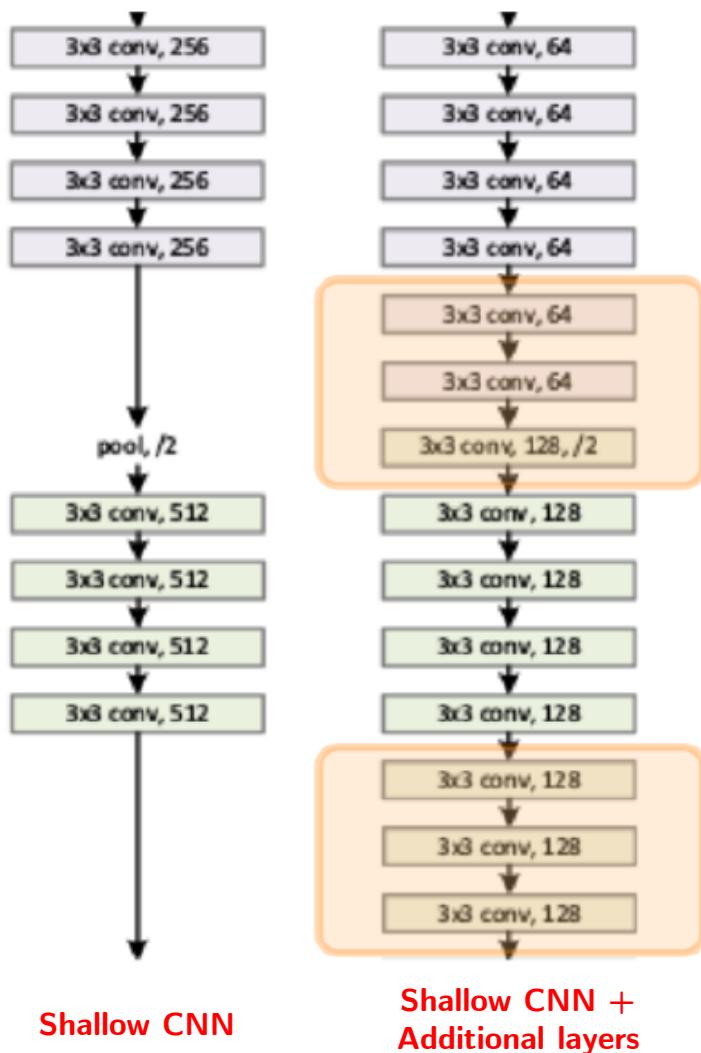
Source: CS7015 Deep Learning, Dept. of CSE, IIT Madras



- 12 times less parameters and 2 times more computations than AlexNet (which is acceptable, given the performance)
- Used **Global Average Pooling** instead of Flattening.
  - If we flatten the final max pool output, the fully connected layer will have  $[1 \times 50176]$  as the input.
  - With global average pooling, the fully connected layer will now have  $[1 \times 1024]$  as the input, which will result in a significant reduction in the no. of parameters (allowing to add more Inception blocks)

Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., & Rabinovich, A. (2015). Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR'15)*

# ResNet Architecture (2015)

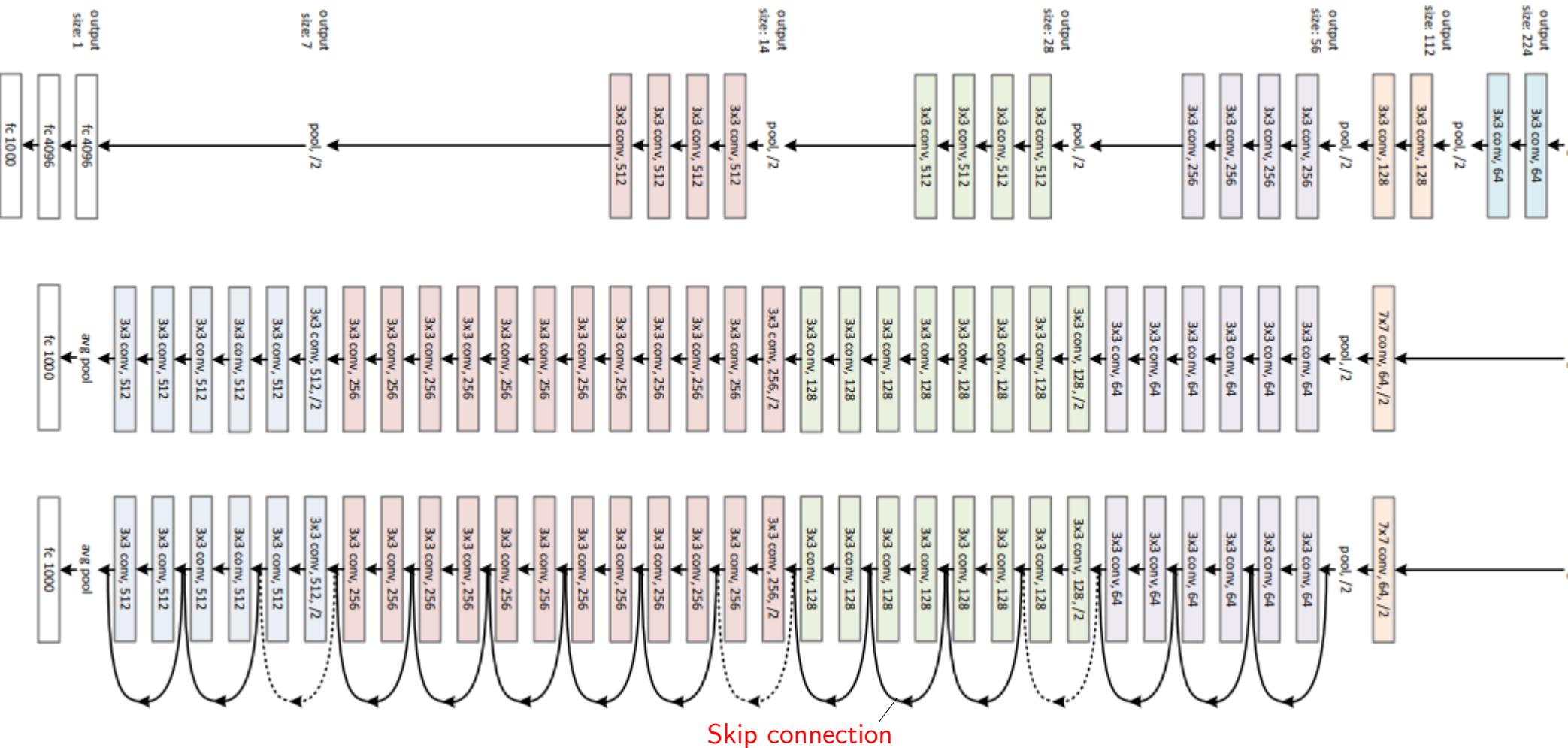


**Effect of increasing layers of shallow CNN when experimented over the CIFAR dataset**

Source: [Residual Networks \(ResNet\) - Deep Learning - GeeksforGeeks](#)

He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770-778).

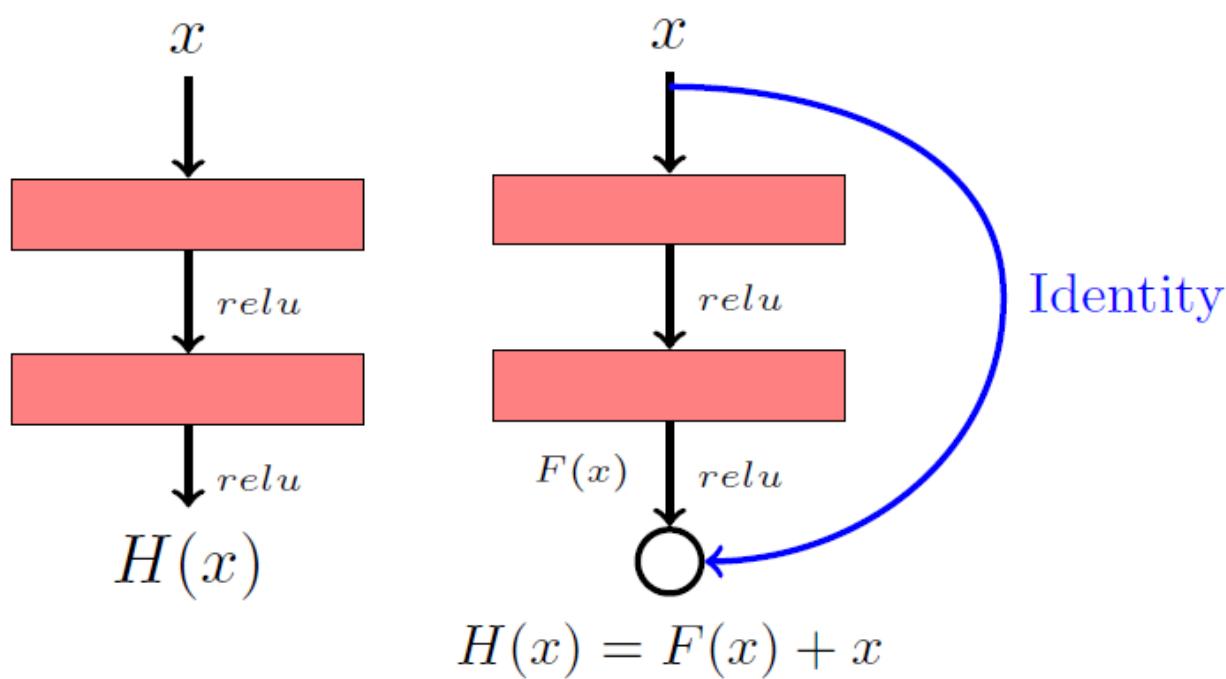
# ResNet Architecture (2015)



Source: [Residual Networks \(ResNet\) - Deep Learning - GeeksforGeeks](#)

He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770-778).

# ResNet Architecture (2015): Skip Connection – The Residual Block



Source: CS7015 Deep Learning, Dept. of CSE, IIT Madras

- It **mitigates the Vanishing/Exploding Gradient Problem** (usually occurs in Deep Networks) by allowing gradient to flow through multiple paths.
- The identity shortcuts act as a form of **regularization**. They can prevent overfitting by reducing the complexity of the mapping the network has to learn, leading to better performance on unseen data.
- Because the network learns the **residual** (the difference between input and output) rather than the full transformation, it **simplifies the learning process** and helps the network focus on refining the output instead of relearning the input information.

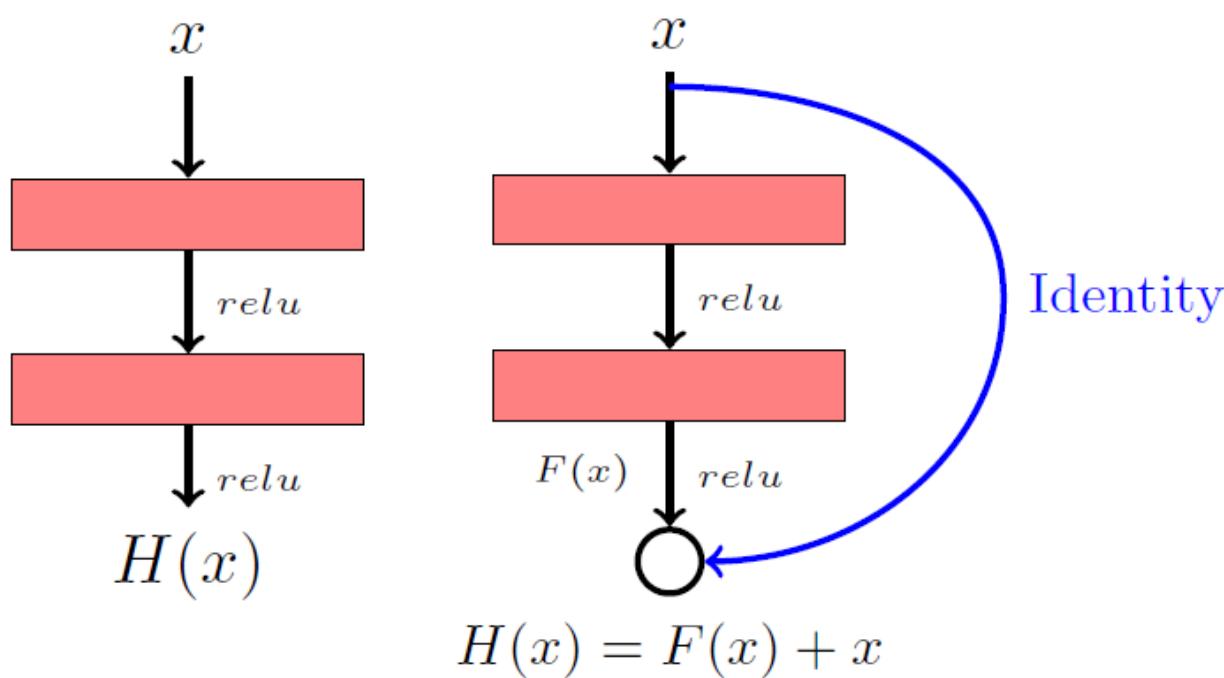
He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770-778).

## Bag of tricks used in ResNet:

- Batch Normalization after every CONV layer
- Xavier/2 initialization
- SGD + Momentum(0.9)
- Learning rate: 0.1, divided by 10 when validation error plateaus
- Mini-batch size 256
- No dropout used

He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition.  
In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770-778).

# ResNet Architecture (2015): Skip Connection – The Residual Block



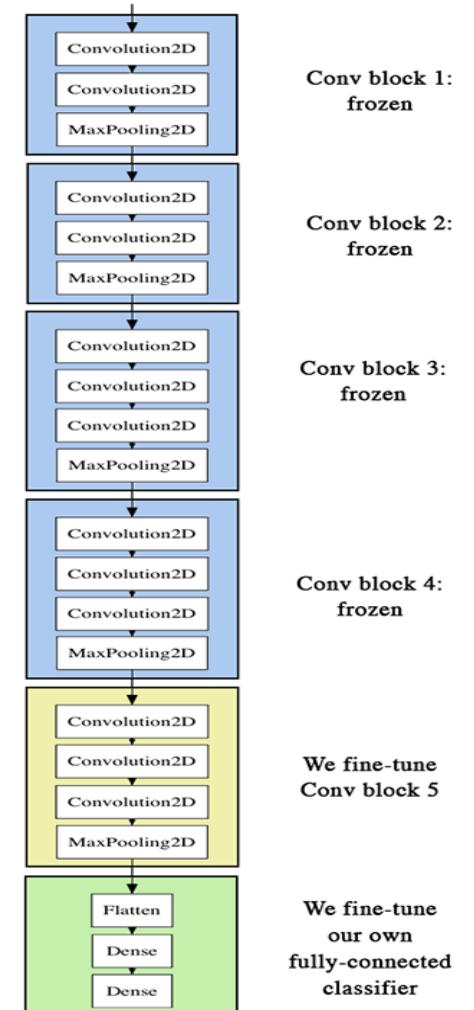
Source: CS7015 Deep Learning, Dept. of CSE, IIT Madras

- It **mitigates the Vanishing/Exploding Gradient Problem** (usually occurs in Deep Networks) by allowing gradient to flow through multiple paths.
- The identity shortcuts act as a form of **regularization**. They can prevent overfitting by reducing the complexity of the mapping the network has to learn, leading to better performance on unseen data.
- Because the network learns the **residual** (the difference between input and output) rather than the full transformation, it **simplifies the learning process** and helps the network focus on refining the output instead of relearning the input information.

He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770-778).

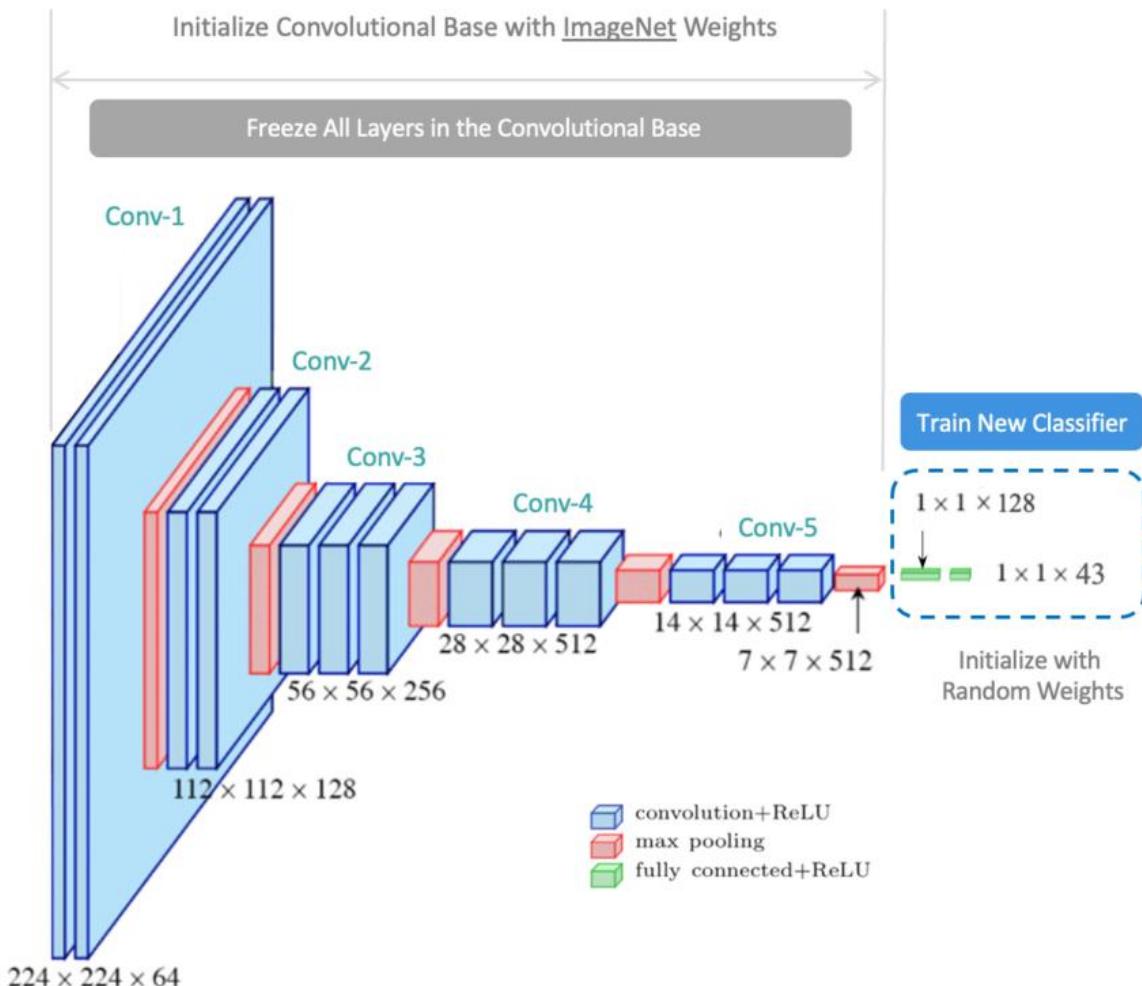
# Leveraging a Pre-trained model

- A pretrained model has been previously trained on a large dataset for a large scale image classification task
- The spatial hierarchy of features learnt by the pretrained model acts as a generic model of the visual world
- Two ways to use a pretrained model:
  - **(1) Feature extraction**
    - Use the convolution base of pre trained model
      - First few layers pick, local , generic features like edges, colors and textures
      - Layers higher up pick abstract concepts like – ear,eye
  - **(2) Transfer Learning/Fine –tuning a pretrained model**
    - Add custom network on top of already trained base network
    - Freeze the base network
    - Train the custom network
    - Unfreeze some layers in base network
    - Jointly train both unfrozen layers and the custom network



Source: Aurelien Geron, "Hands-On Machine Learning with Scikit-Learn , Keras & Tensorflow, O'Reilly Publications

# Transfer Learning vs Fine-Tuning

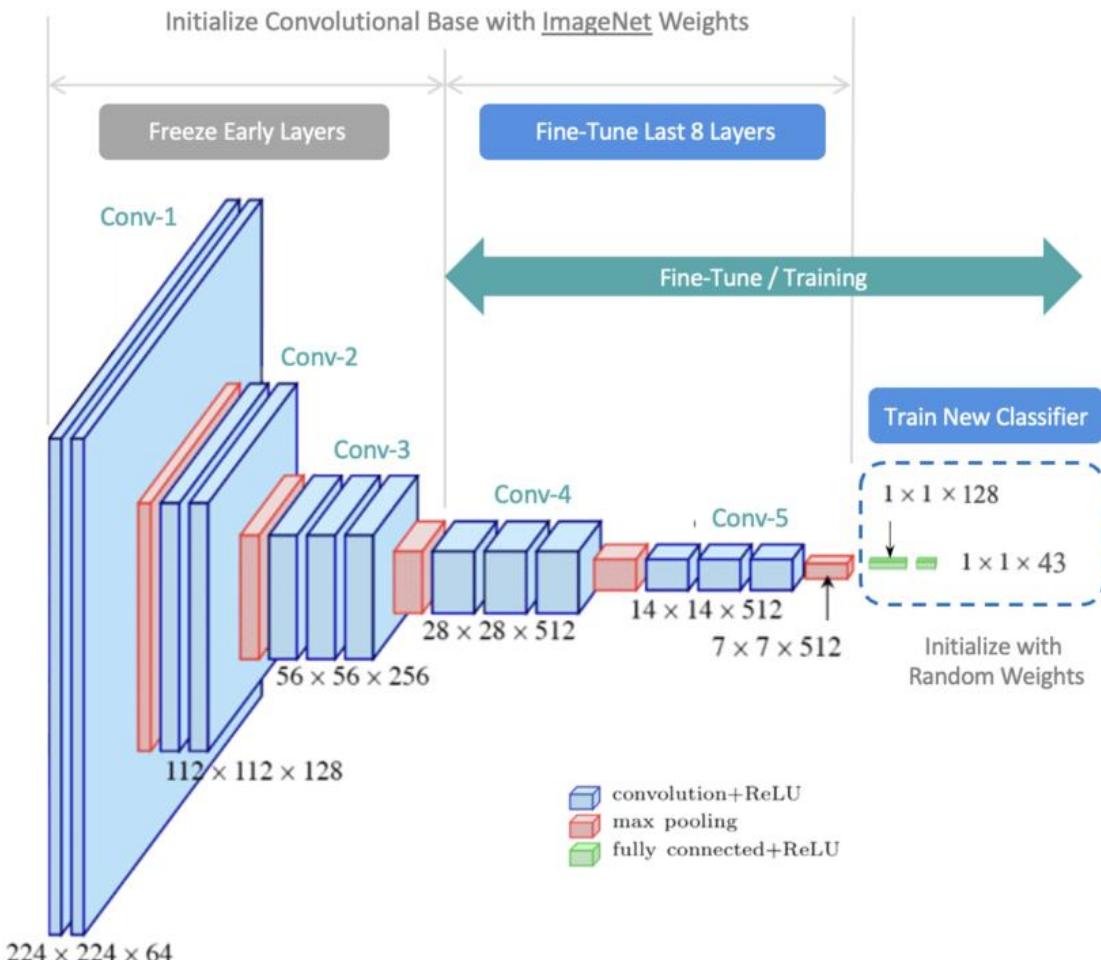


## Transfer Learning:

- Freeze the feature extraction layers of the pre-trained model and train only the final classifier layer.
- Allows you to quickly study how a pre-trained model can be customized for a new dataset.
- However, sometimes retraining the classifier isn't enough. This is where Fine-Tuning can be very beneficial.

Source: [Unlock the Power of Fine-Tuning Pre-Trained Models in TensorFlow & Keras \(learnopencv.com\)](#)

# Transfer Learning vs Fine-Tuning



## Fine-tuning:

- Similar to transfer learning, but with the additional step of selectively training some of the pre-trained model's layers.
- Freeze the first several layers of the feature extractor but allow the last few layers to be trained further.
- The idea is that the first several layers in the feature extractor represent generic, low-level features (e.g., edges, corners, and arcs) that are fundamental building blocks necessary to support many classification tasks.

Source: [Unlock the Power of Fine-Tuning Pre-Trained Models in TensorFlow & Keras \(learnopencv.com\)](http://learnopencv.com)