

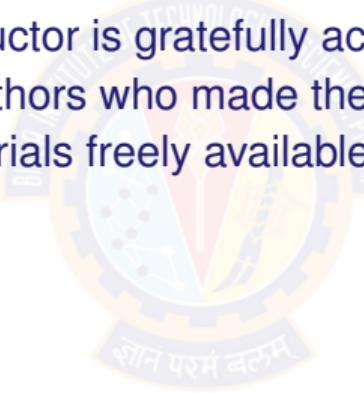
**BITS Pilani**  
Pilani | Dubai | Goa | Hyderabad

## DEEP LEARNING MODULE # 5 : CONVOLUTIONAL NEURAL NETWORK [CNN]

---

Seetha Parameswaran  
Asst Prof, BITS Pilani

The instructor is gratefully acknowledging  
the authors who made their course  
materials freely available online.

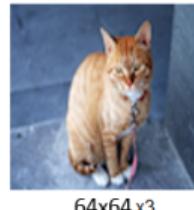


# IN THIS SEGMENT

- 1 COMPUTER VISION TASKS
- 2 IMAGE BASICS
- 3 IMAGE CLASSIFICATION CHALLENGES
- 4 CONVOLUTIONAL NEURAL NETWORK
- 5 CONVOLUTION
- 6 PUTTING IT TOGETHER



# IMAGE CLASSIFICATION



→ Cat? (0/1)



→ Cat? (0/1)



Classify into 3 classes

# FACE DETECTION



# OBJECT RECOGNITION

"SIFT" & Object Recognition, David Lowe, 1999

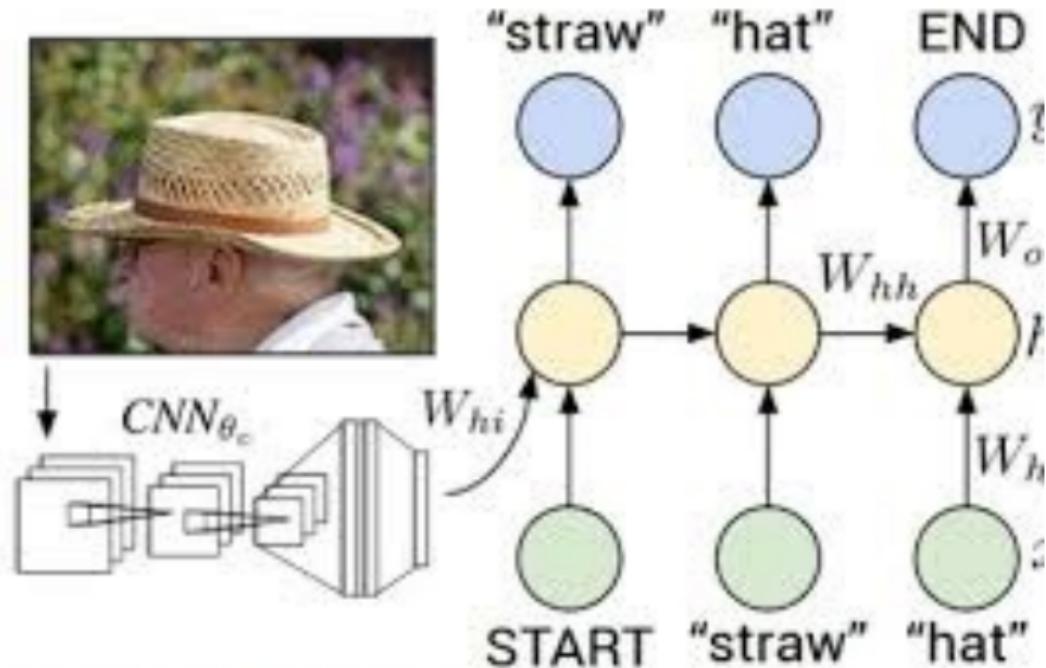


[Image](#) is public domain



[Image](#) is CC BY-SA 2.0

# IMAGE CAPTIONING



Andrej Karpathy, Fei-Fei Li

# IN THIS SEGMENT

- 1 COMPUTER VISION TASKS
- 2 IMAGE BASICS
- 3 IMAGE CLASSIFICATION CHALLENGES
- 4 CONVOLUTIONAL NEURAL NETWORK
- 5 CONVOLUTION
- 6 PUTTING IT TOGETHER



## DIGITAL GRayscale IMAGE

Pixel intensities =  $I(x, y)$



0	2	15	0	0	11	10	0	0	0	0	9	9	0	0
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0
0	10	16	115	238	255	244	245	243	250	249	255	222	103	10
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124
2	90	255	228	255	251	254	211	141	116	122	215	251	238	255
13	217	243	255	155	33	226	52	2	0	10	13	232	255	255
16	229	252	254	49	12	0	0	7	7	0	70	237	252	235
6	141	245	255	212	25	11	9	3	0	115	236	243	255	137
0	81	252	250	248	215	60	0	1	21	252	255	248	144	6
0	13	111	255	255	245	255	182	181	248	252	242	208	36	0
1	0	5	117	251	255	241	255	247	255	241	162	17	0	7
0	0	0	4	58	251	255	246	254	253	255	120	11	0	1
0	0	4	91	255	255	255	248	252	255	244	255	182	10	0
0	22	206	252	246	251	241	100	24	113	255	245	255	194	9
0	111	255	242	255	158	24	0	0	6	39	255	232	230	56
0	218	251	250	137	7	11	0	0	0	2	62	255	250	125
0	173	255	255	103	9	20	0	13	3	13	182	251	245	61
0	107	251	241	255	230	98	55	19	118	217	248	253	255	52
0	18	146	250	255	247	255	255	255	249	255	240	255	129	0
0	0	23	111	215	255	250	248	255	255	248	248	118	14	12
0	0	6	1	0	52	153	233	255	252	147	37	0	0	4
0	0	5	5	0	0	0	0	0	0	14	1	0	6	6

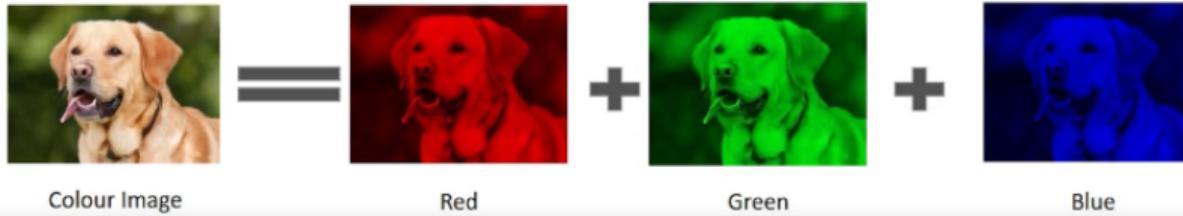
0	2	15	0	0	11	10	0	0	0	0	9	9	0	0	0
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0	29
0	10	16	119	238	255	244	245	243	250	249	255	222	103	10	0
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124	1
2	98	255	228	255	251	254	211	141	116	122	215	251	238	255	49
13	217	243	255	155	33	226	52	2	0	10	13	232	255	255	36
16	229	252	254	49	12	0	0	7	7	0	70	237	252	235	62
6	141	245	255	212	25	11	9	3	0	115	236	243	255	137	0
0	87	252	250	248	215	60	0	1	121	252	255	248	144	6	0
0	13	113	255	255	245	255	182	181	248	252	242	208	36	0	19
1	0	5	117	251	255	241	255	247	255	241	162	17	0	7	0
0	0	0	4	58	251	255	246	254	253	255	120	11	0	1	0
0	0	4	97	255	255	255	248	252	255	244	255	182	10	0	4
0	22	206	252	246	251	241	100	24	113	255	245	255	194	9	0
0	0	111	255	242	255	158	24	0	0	6	39	255	232	230	56
0	0	218	251	250	137	7	11	0	0	0	2	62	255	250	125
0	0	173	255	255	101	9	20	0	13	3	13	182	251	245	61
0	0	107	251	241	255	230	98	55	19	118	217	248	253	255	52
0	0	18	146	250	255	247	245	255	255	249	255	240	255	129	0
0	0	0	0	23	113	215	255	250	248	255	255	248	248	118	14
0	0	0	6	1	0	52	153	233	255	252	147	37	0	0	4
0	0	0	5	5	0	0	0	0	0	14	1	0	6	6	0

# DIGITAL COLOR IMAGE

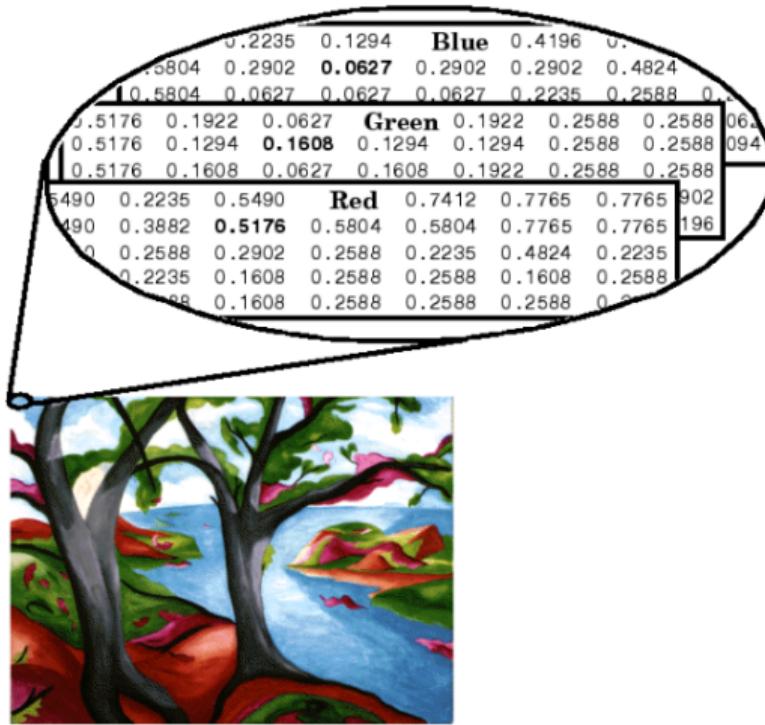


Colour Image

This image is composed of many colors and almost all colors can be generated from the three primary colors- **Red, Green, and Blue**. We can say that each colored image is composed of these three colors or 3 channels- Red, Green, and Blue-



# DIGITAL COLOR IMAGE



# IN THIS SEGMENT

- 1 COMPUTER VISION TASKS
- 2 IMAGE BASICS
- 3 IMAGE CLASSIFICATION CHALLENGES
- 4 CONVOLUTIONAL NEURAL NETWORK
- 5 CONVOLUTION
- 6 PUTTING IT TOGETHER



# IMAGE CLASSIFICATION - CHALLENGES

**Image Classification:** A core task in Computer Vision



This image by Nikita is  
licensed under CC-BY 2.0

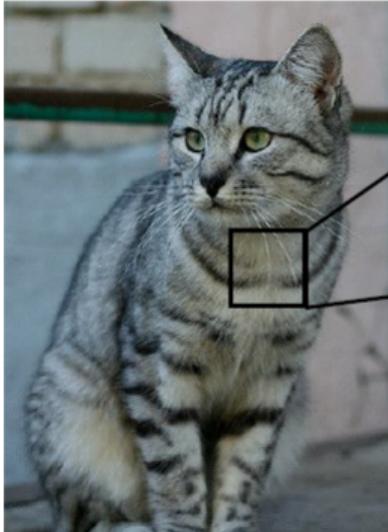
(assume given set of discrete labels)  
{dog, cat, truck, plane, ...}



cat

# IMAGE CLASSIFICATION - CHALLENGES

## The Problem: Semantic Gap



This image by [Nikita](#) is  
licensed under [CC-BY 2.0](#)

[1] 185 112 106 111 104 99 106 99 96 103 112 119 104 97 93 87]
[1] 91 98 192 106 104 79 98 183 99 105 123 136 118 105 94 85]
[1] 76 85 98 105 128 105 87 96 95 99 115 112 106 103 99 85]
[1] 99 81 81 93 128 131 127 189 95 98 102 99 96 92 181 94]
[1] 106 91 61 64 69 91 88 85 181 187 109 98 75 84 96 95]
[1] 114 108 85 55 55 69 64 54 64 87 112 129 98 74 84 91]
[1] 133 137 147 103 65 81 88 65 52 54 74 84 102 93 85 82]
[1] 128 137 144 148 109 95 86 79 62 65 63 63 68 73 88 181]
[1] 125 133 148 137 119 121 117 94 65 79 86 65 54 64 72 98]
[1] 127 125 131 147 133 127 126 131 111 96 89 75 61 64 72 84]
[1] 115 114 109 123 158 148 131 118 113 109 106 92 74 65 72 78]
[1] 89 93 98 97 108 147 131 118 113 114 113 109 106 95 77 88]
[1] 63 77 86 81 77 79 102 123 117 115 117 125 125 138 115 87]
[1] 62 65 82 89 78 71 88 101 124 126 119 101 107 114 131 119]
[1] 63 65 75 88 89 71 62 81 120 138 135 105 81 90 118 118]
[1] 87 85 71 87 106 95 69 45 76 130 126 107 92 94 185 112]
[1] 118 97 82 86 117 122 116 66 41 51 95 93 89 95 182 187]
[1] 164 146 112 88 82 128 124 184 76 48 45 66 80 101 182 189]
[1] 157 179 157 128 93 86 114 132 112 97 69 55 78 82 99 94]
[1] 139 128 134 161 139 108 105 118 121 134 114 87 65 53 69 86]
[1] 128 112 96 117 158 144 128 115 184 197 102 93 87 81 72 79]
[1] 123 187 96 86 83 112 153 149 122 189 194 75 89 107 112 99]
[1] 122 121 192 98 82 85 94 117 145 148 153 102 58 78 92 147]
[1] 122 164 148 143 71 56 78 83 93 103 118 139 102 61 69 84]

What the computer sees

An image is just a big grid of numbers between [0, 255]:

e.g. 800 x 600 x 3  
(3 channels RGB)

# IMAGE CLASSIFICATION - CHALLENGES

## Challenges: Viewpoint variation



All pixels change when  
the camera moves!

This image by Nikita is  
licensed under [CC-BY 2.0](#)

# IMAGE CLASSIFICATION - CHALLENGES

## Challenges: Illumination



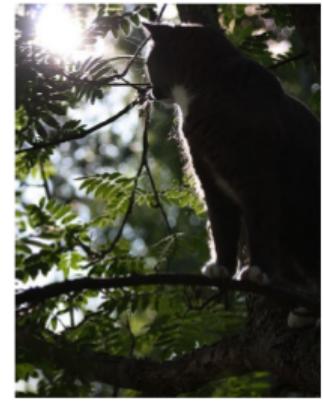
[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)



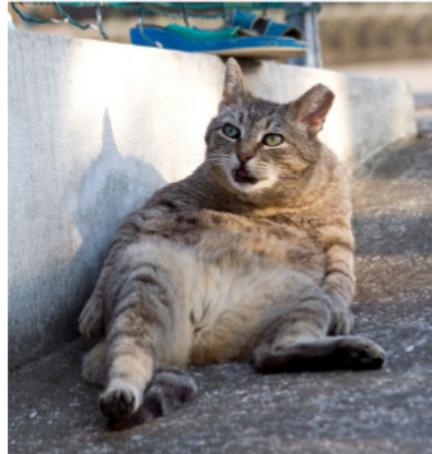
[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)

# IMAGE CLASSIFICATION - CHALLENGES

## Challenges: Deformation



This image by Umberto Salvagnin  
is licensed under CC-BY 2.0



This image by Umberto Salvagnin  
is licensed under CC-BY 2.0



This image by sare bear is  
licensed under CC-BY 2.0



This image by Tom Thai is  
licensed under CC-BY 2.0

# IMAGE CLASSIFICATION - CHALLENGES

## Challenges: Occlusion



[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)



[This image by jonsson is licensed under CC-BY 2.0](#)

# IMAGE CLASSIFICATION - CHALLENGES

## Challenges: Background Clutter



[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)

# IMAGE CLASSIFICATION - CHALLENGES

**Challenges:** Intraclass variation



This image is CC0 1.0 public domain

# IMAGE CLASSIFICATION - CHALLENGES

Pixelated Images



Larger Images



- A 50x50 Image requires 2500 input dimensions.
- A HD image 1920x1080 is about 2 million inputs. (2 Mega Pixels)

# IMAGE INVARIANCE - CHALLENGES

Object can appear anywhere in the image.

Image A

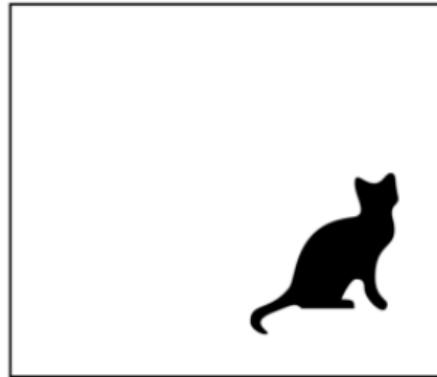
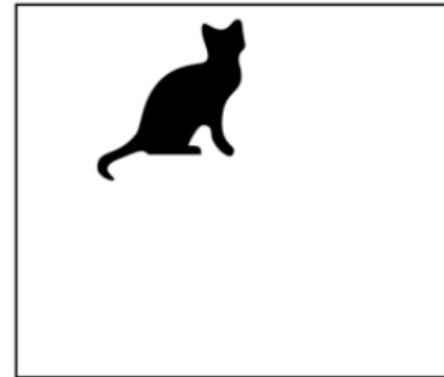


Image B

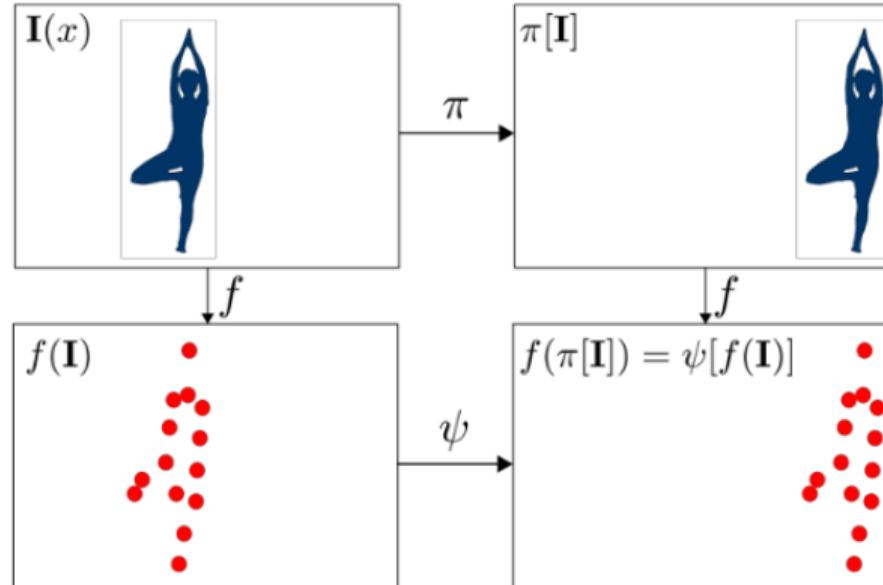


Model predicts “cat” for both images!

Source: <https://penkovsky.com/neural-networks/day5/>

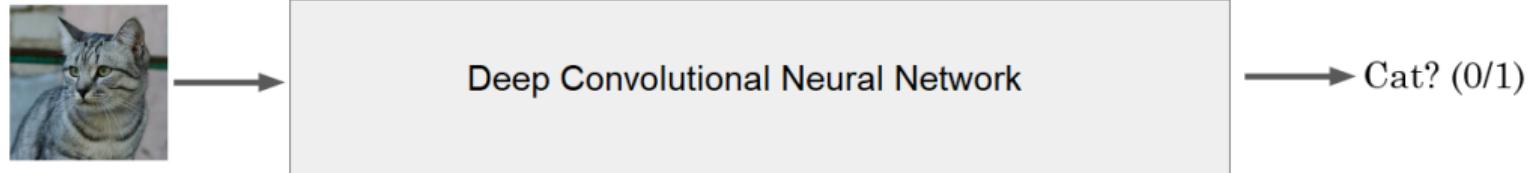
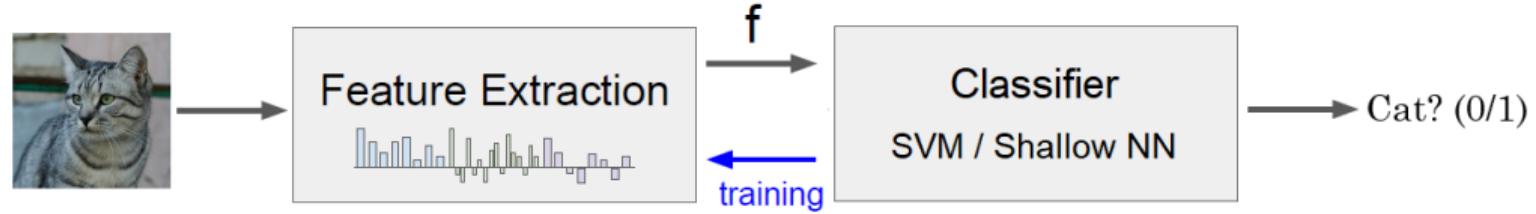
# IMAGE EQUIVARIANCE - CHALLENGES

Changes to the input predictably (equivalently) should not affect the output.



Source: <https://towardsdatascience.com/translational-invariance-vs-translational-equivariance-f9fbc8fcfa63a>

# IMAGE CLASSIFICATION - SOLUTION



# IN THIS SEGMENT

- 1 COMPUTER VISION TASKS
- 2 IMAGE BASICS
- 3 IMAGE CLASSIFICATION CHALLENGES
- 4 CONVOLUTIONAL NEURAL NETWORK
- 5 CONVOLUTION
- 6 PUTTING IT TOGETHER



# CONVOLUTIONAL NEURAL NETWORKS (CNN)

- Convolutional networks are simply neural networks that use convolution in place of general matrix multiplication in at least one of their layers.
- Special kind of neural network for processing data that has a known, grid-like topology.
- Eg: Time-series data – 1D grid taking samples at regular time intervals.
- Eg: Image data – 2D grid of pixels.
- Network employs a mathematical operation called **convolution**.

# EXAMPLE OF CNN – LENET-5

- Handwritten Character recognition
- 10-categories classification problem

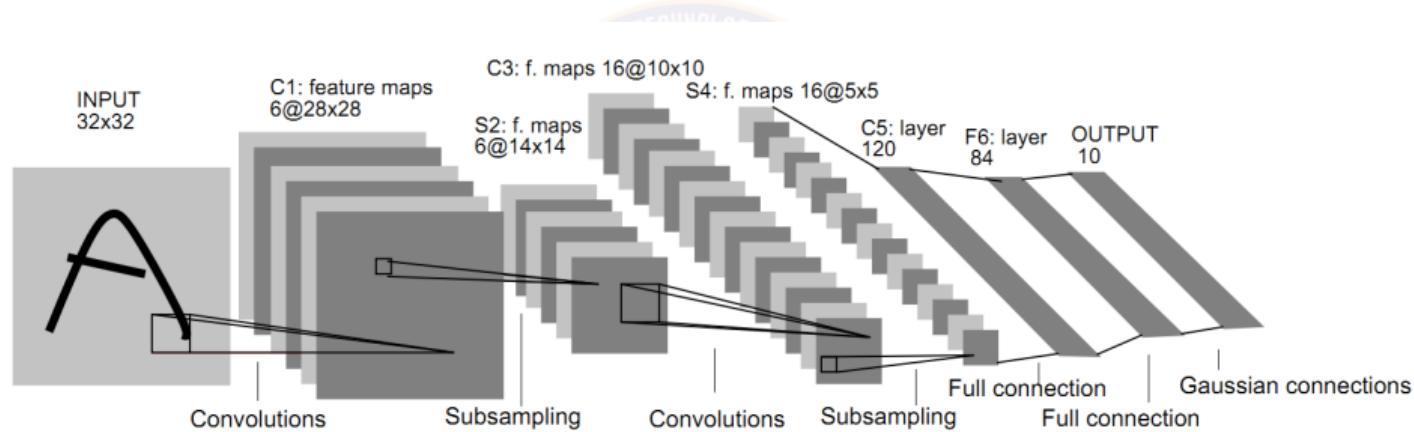


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

# WHY CNN FOR IMAGES

- Some patterns are much smaller than the whole image.



# WHY CNN FOR IMAGES

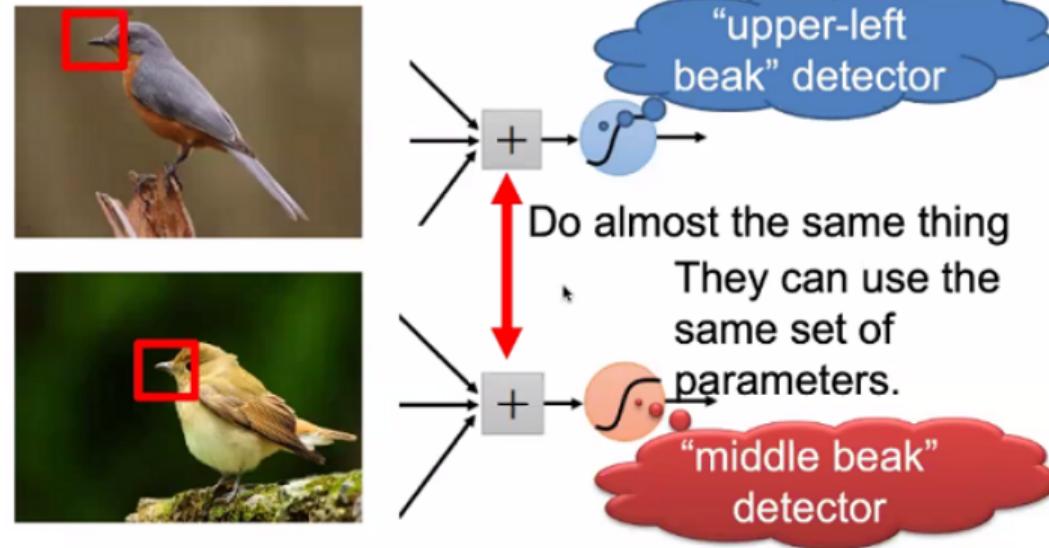
- Some patterns are much smaller than the whole image.

A neuron does not have to see the whole image to discover the pattern.



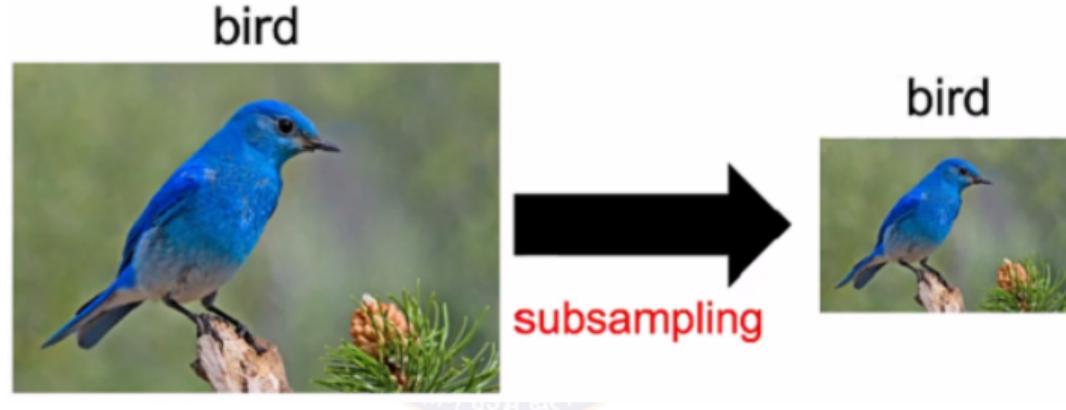
# WHY CNN FOR IMAGES

- Same pattern may appear in different regions.



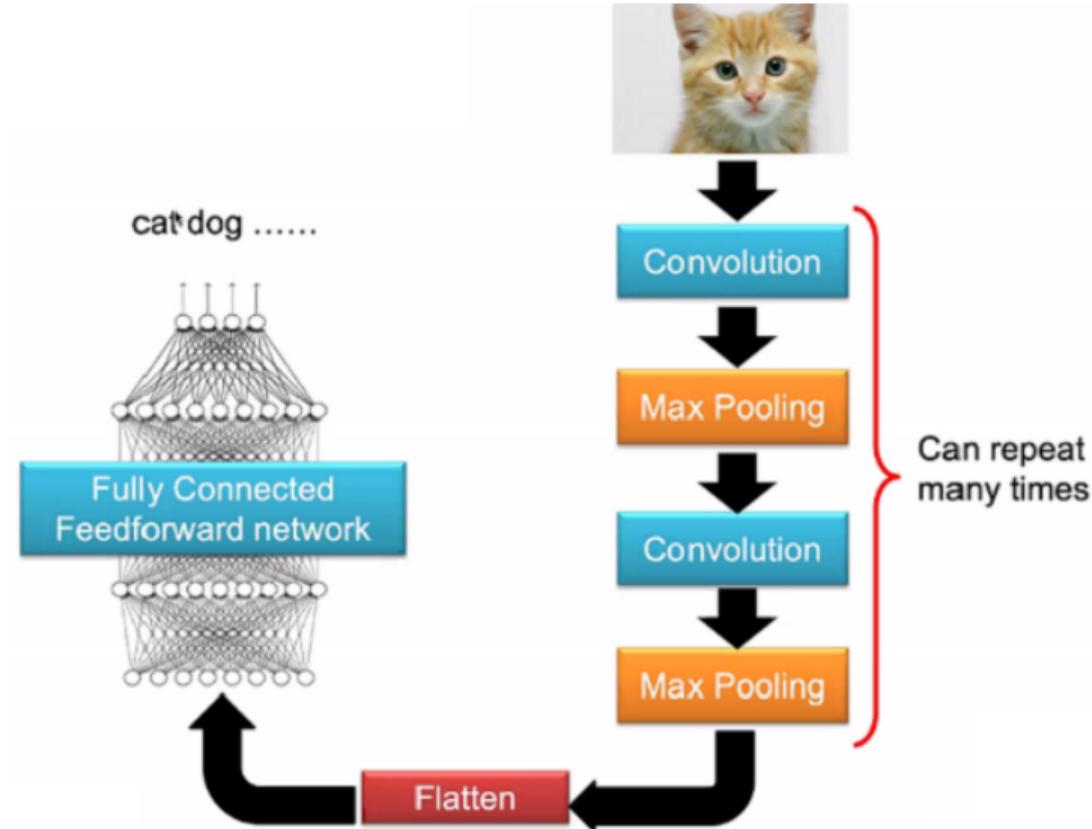
# WHY CNN FOR IMAGES

- Sub-sampling the pixels will not change the object.



- We can sub-sample the image to make it smaller.
- Less parameters for the network to process the image.

# THE WHOLE CNN



# THE WHOLE CNN

## Property 1

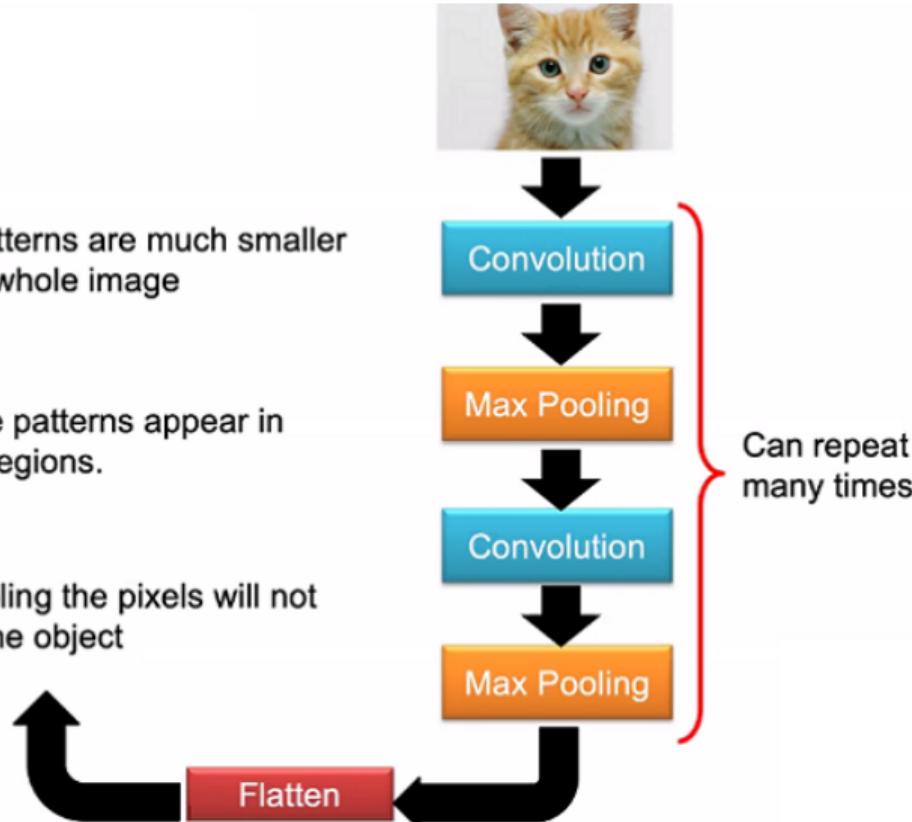
- Some patterns are much smaller than the whole image

## Property 2

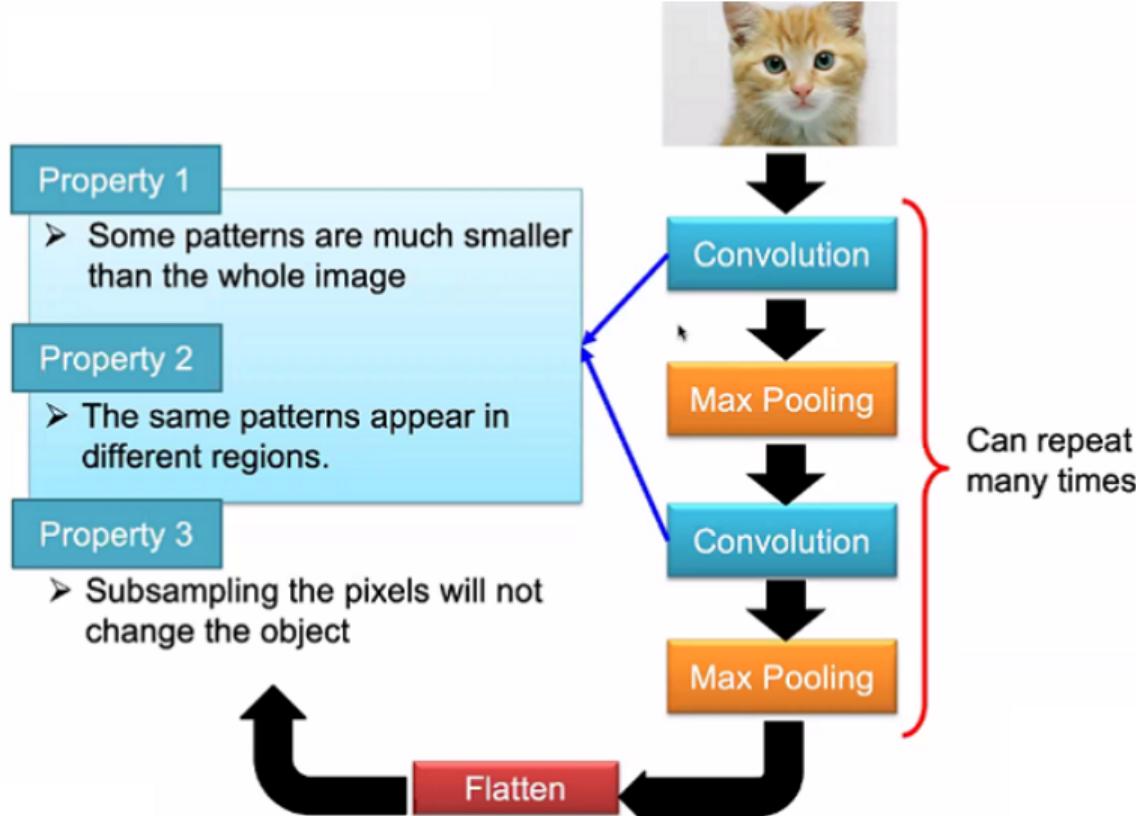
- The same patterns appear in different regions.

## Property 3

- Subsampling the pixels will not change the object



# THE WHOLE CNN



# THE WHOLE CNN

## The whole CNN

### Property 1

- Some patterns are much smaller than the whole image

### Property 2

- The same patterns appear in different regions.

### Property 3

- Subsampling the pixels will not change the object



Convolution



Max Pooling



Convolution



Max Pooling



Can repeat  
many times

# IN THIS SEGMENT

- 1 COMPUTER VISION TASKS
- 2 IMAGE BASICS
- 3 IMAGE CLASSIFICATION CHALLENGES
- 4 CONVOLUTIONAL NEURAL NETWORK
- 5 CONVOLUTION
- 6 PUTTING IT TOGETHER



# CONVOLUTION LAYER

- Perform Convolution operation.
- In mathematics, convolution is a mathematical operation on two functions that produces a third function that expresses how the shape of one is modified by the other. [Wikipedia]
- The term convolution refers to both the result function and to the process of computing it. [Wikipedia]
- We will assume that the kernel is square matrix.
- We will assume that the kernel is centered on the pixel of interest. (One reason why kernel is odd matrix.)
- Convolution operation

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

# 1-D CONVOLUTION

- Suppose we track position of an object using a sensor or GPS device at discrete time intervals.
- Assume that the sensor is noisy.
- To obtain a less noisy estimate, we would like to take a weighted average.
- We would only sum over a small window or a weight array (filter) ( $w$ ).

$$s_t = \sum_{a=0}^k x_{t-a} w_a$$

- Slide the filter over the input and convolve.
- The input and the kernel are one dimensional.

# 1-D CONVOLUTION

$$S = x_5w_0 + x_4w_1 + x_3w_2 + x_2w_3 + x_1w_4 + x_0w_5$$

X	1.00	1.10	1.20	1.40	1.70	1.80	1.90	2.10	2.20	2.40	2.50
---	------	------	------	------	------	------	------	------	------	------	------

W	0.01	0.02	0.03	0.04	0.4	0.5					
---	------	------	------	------	-----	-----	--	--	--	--	--

S			1.80								
---	--	--	------	--	--	--	--	--	--	--	--

W	0.01	0.02	0.03	0.04	0.4	0.5					
---	------	------	------	------	-----	-----	--	--	--	--	--

S		1.80	1.96								
---	--	------	------	--	--	--	--	--	--	--	--

W	0.01	0.02	0.03	0.04	0.4	0.5					
---	------	------	------	------	-----	-----	--	--	--	--	--

S		1.80	1.96	2.11							
---	--	------	------	------	--	--	--	--	--	--	--

# 2D CONVOLUTION OPERATION

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  $6 \times 6$

These values inside the matrix are to be learned by the network. We call these as parameters or weights.

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

Kernel / Filter 1

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

Kernel / Filter 2

Property 1: Each filter detects a small pattern ( $3 \times 3$ )

# CONVOLUTION OPERATION

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

Start at top left corner.

Perform element wise product and add all the values.

Output

3

# CONVOLUTION OPERATION

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

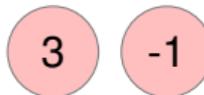
Input Image  $6 \times 6$

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

Filter 1

Slide by one pixel towards right.  
Perform element wise product and  
add all the values.

Output



# CONVOLUTION OPERATION

stride  $s = 1$

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

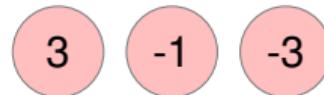
Input Image  $6 \times 6$

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

Filter 1

Slide by one pixel towards right.  
Perform element wise product and  
add all the values.

Output



# CONVOLUTION OPERATION

stride  $s = 1$

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

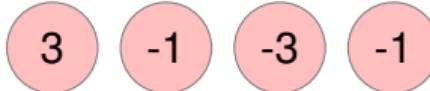
Input Image  $6 \times 6$

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

Filter 1

Slide by one pixel towards right.  
Perform element wise product and  
add all the values.

Output



# CONVOLUTION OPERATION

stride  $s = 1$

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

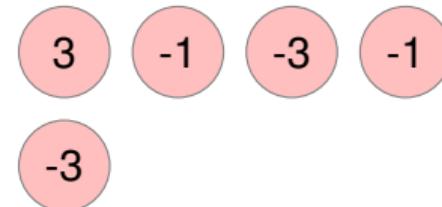
Input Image  $6 \times 6$

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

Filter 1

Jump to next row, left most pixel.  
The filter has to completely  
superimpose the pixels of input.  
NOT PARTIALLY.

Output



# CONVOLUTION OPERATION

stride  $s = 1$

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

After sliding through the entire image.

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

# CONVOLUTION OPERATION

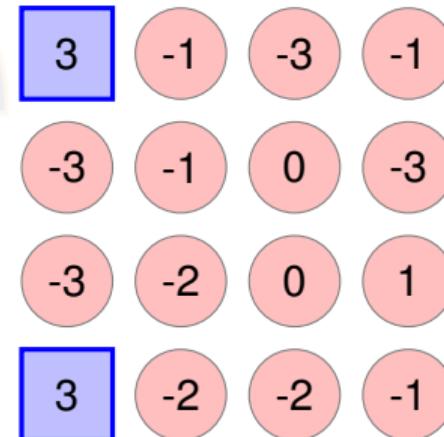
stride  $s = 1$

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



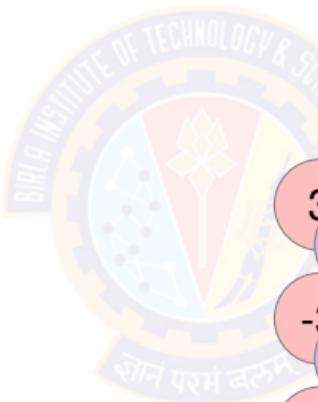
Property 2 is satisfied.

# CONVOLUTION OPERATION

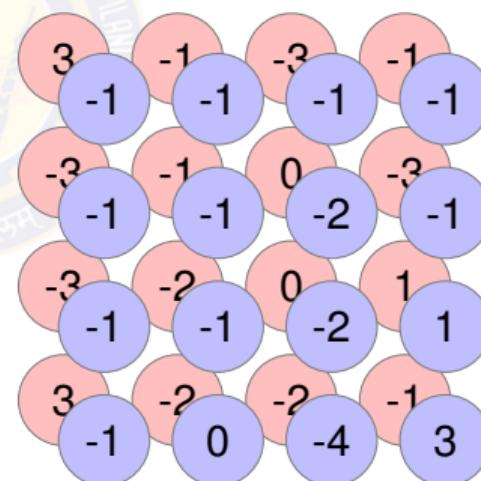
stride  $s = 1$

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  $6 \times 6$



Filter 1			Filter 2		
1	-1	-1	-1	1	-1
-1	1	-1	-1	1	-1
-1	-1	1	-1	1	-1



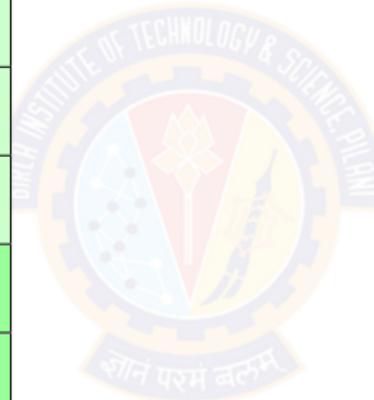
Feature  
Maps

# CONVOLUTION OPERATION

stride  $s = 3$

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  $6 \times 6$



Filter1

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

Output

3	-1
3	-1

# CONVOLUTION WITH STRIDING

- Stride means skip by how many pixels.

Input Size :  $n \times n$

Stride :  $s$

Kernel Size :  $f \times f$

Output Size :  $\left\lfloor \frac{n-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n-f}{s} + 1 \right\rfloor$

- The filter must lie completely inside the input.

# PADDING

- When filter is convolved with the input image
  - ▶ Image shrinks –  $6 \times 6$  image became a  $4 \times 4$  image.
  - ▶ Pixels in the corners are used only once, when compared the pixels in the middle. The data in the corners are thrown away.
- To resolve, apply padding to the input image before convolution.
- Padding means append zeros around the boundary of the input.

# CONVOLUTION WITH PADDING

Input Image with padding  $8 \times 8$

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

Filter1

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

# PADDING

- Two kinds of padding

- Valid Padding

- Padding is not applied.

$$p = 0$$

- Input size  $\neq$  output size.
    - Output shrinks when compared to the input.

- Same Padding

- Padding is applied.

$$p = \frac{f - 1}{2}$$

- Input size = output size.
    - We maintain the size of input and output as the same.

# PADDING

- How much to pad?

$$\text{Padding : } p = \frac{f - 1}{2}$$

Kernel  $3 \times 3 : p = 1$

Kernel  $5 \times 5 : p = 2$

Kernel  $7 \times 7 : p = 3$

- Padding preserves the input size.

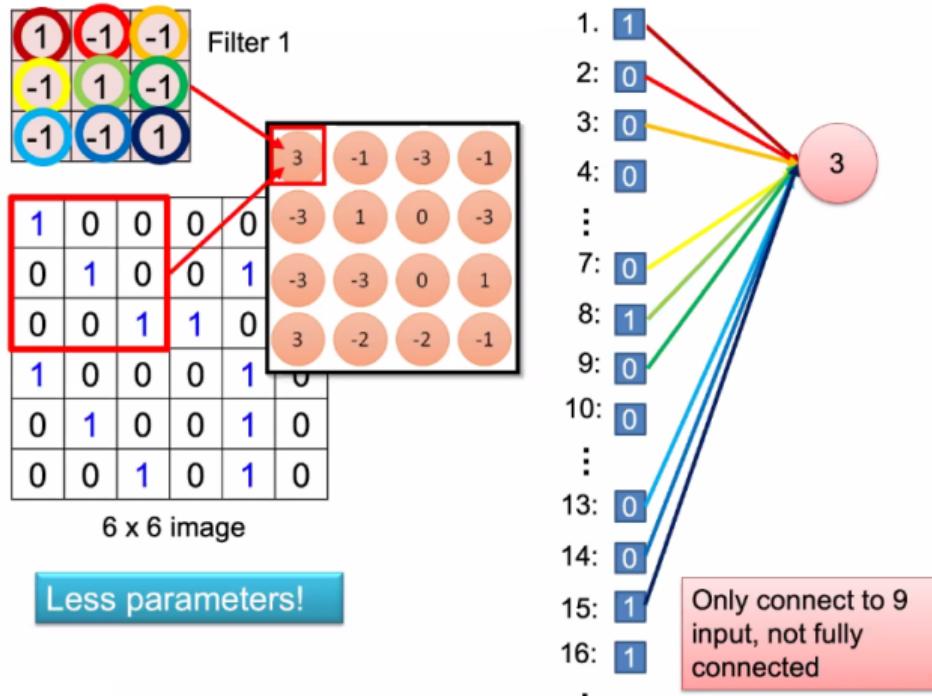
Input Size :  $n \times n$

Padding :  $p$

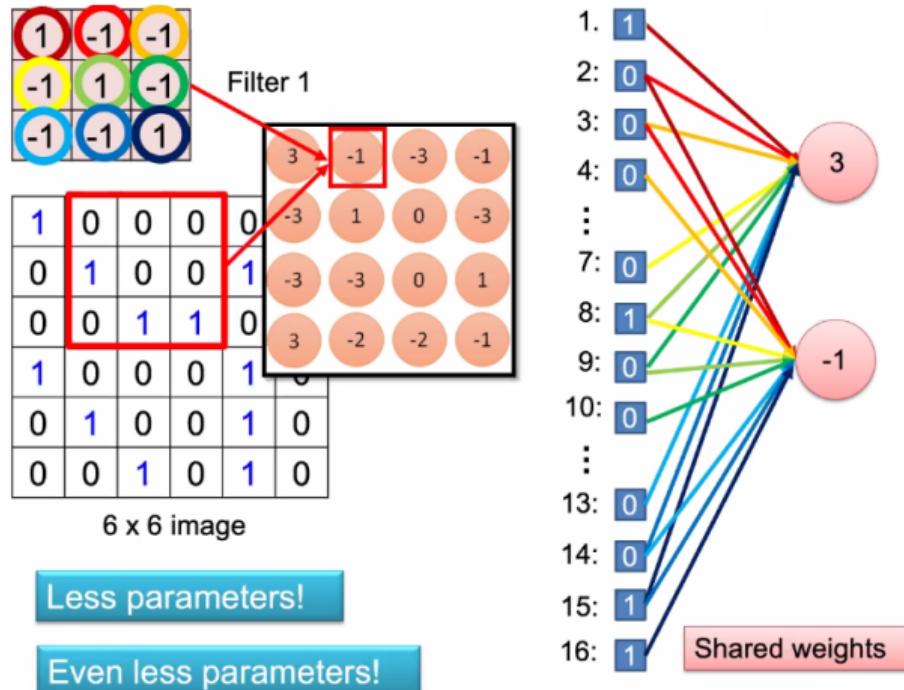
Kernel Size :  $f \times f$

Output Size :  $(n + 2p - f + 1) \times (n + 2p - f + 1)$

# CONVOLUTIONAL LAYER

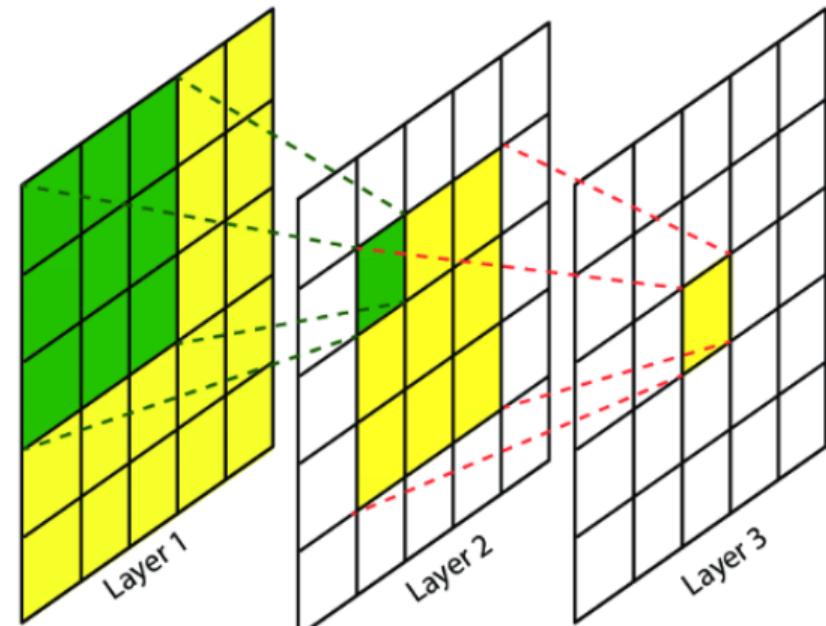


# CONVOLUTIONAL LAYER

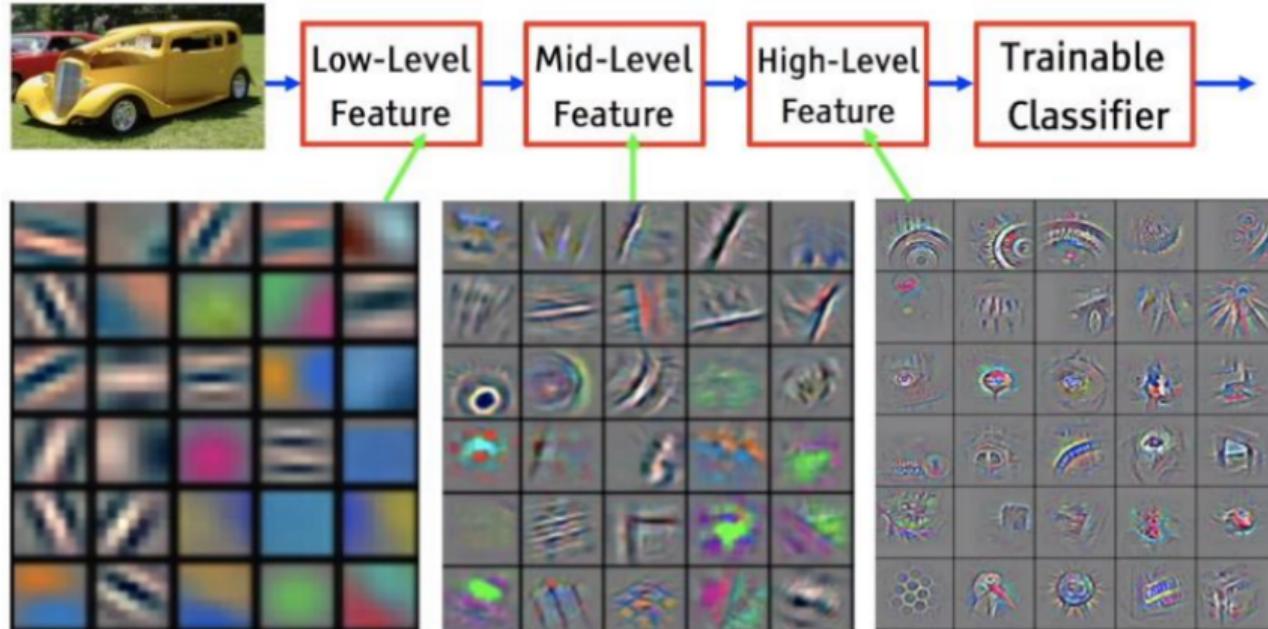


# ACTIVATIONS OR RECEPTIVE FIELDS

- Effective area “seen” by a neuron
- What is the effective number of pixels seen at each layer?
  - Input image  $x$
  - $h_1 = \text{Conv2d}(x, k = 3)$
  - $h_2 = \text{Conv2d}(h_1, k = 3)$
  - $h_3 = \text{Conv2d}(h_2, k = 3)$
  - ...
  - $h_n = \text{Conv2d}(h_{n-1}, k = 3)$
- Homework: What happens if the kernel size was 5?



# LOCAL TO GLOBAL



# PARAMETER SIZE IN CONVOLUTION LAYER

- Input image:  $B \times W_{in} \times H_{in} \times C_{in}$ 
  - ▶ Input data has  $B$  samples
  - ▶ Input or Image size is  $W_{in} \times H_{in} \times C_{in}$
- Convolution filter
  - ▶ Weight parameters:  $C_{in} \times C_{out} \times K_W \times K_H$
  - ▶ Bias parameters:  $C_{out}$
- Output activation:  $B \times W_{out} \times H_{out} \times C_{out}$ 
  - ▶  $C_{out}$  filters of shape  $C_{in} \times K_W \times K_H$  operate on each sample

# PARAMETER SIZE IN CONVOLUTION LAYER

- Each filter gives one 2D output.
- K filter will give K such 2D outputs.
- The number of parameters depend on the number of kernels and the kernel size.

Number of kernels in layer ( $I - 1$ ) :  $K_{I-1}$

Kernel size :  $f \times f$

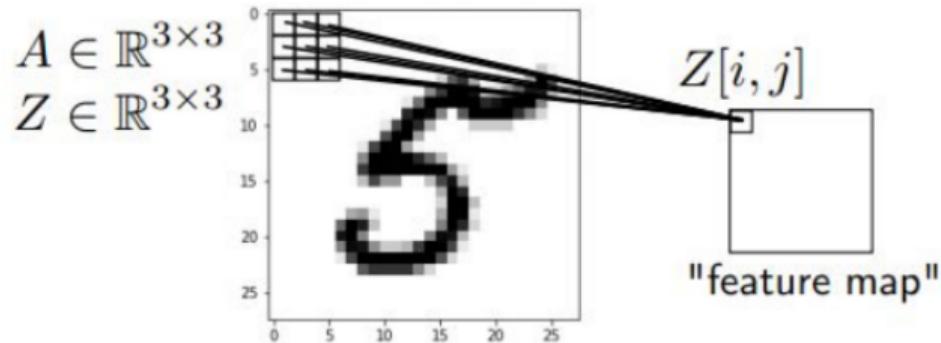
Number of kernels in layer ( $I$ ) :  $K_I$

Number of Parameters in layer ( $I$ ) :  $K_{I-1} * f * f * K_I$

# CROSS-CORRELATION VS CONVOLUTION

- Convolution in Deep Learning is actually cross-correlation.
- Cross-correlation is sliding dot product over the images.
- Cross-correlation

$$Z[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k K[u, v] A[i + u, j + v]$$



# CROSS-CORRELATION VS CONVOLUTION

## Cross-Correlation

$$Z[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k K[u, v] A[i+u, j+v]$$

(Note the looping direction given in red.)

1) -1,-1	2) -1,0	3) -1,1
4) 0,-1	5) 0,0	6) 0,1
7) 1,-1	8) 1,0	9) 1,1

## Convolution

$$Z[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k K[u, v] A[i-u, j-v]$$

(Flipping the kernel.)

9) -1,-1	8) -1,0	7) -1,1
6) 0,-1	5) 0,0	4) 0,1
3) 1,-1	2) 1,0	1) 1,1

# CROSS-CORRELATION VS CONVOLUTION

- Convolution in Deep Learning is actually cross-correlation.
- Convolution (as in Signal Processing) has nice association property.

$$(A * B) * C = A * (B * C)$$

- In DL, we usually don't care about that ( as opposed to traditional computer vision and signal processing applications.)
- Cross-correlation is easier to implement.
- Maybe the term "convolution" for cross-correlation because popular, because "Cross-correlation Neural Network" sounds weird.

# TRANSPOSED CONVOLUTION

## Regular Convolution:

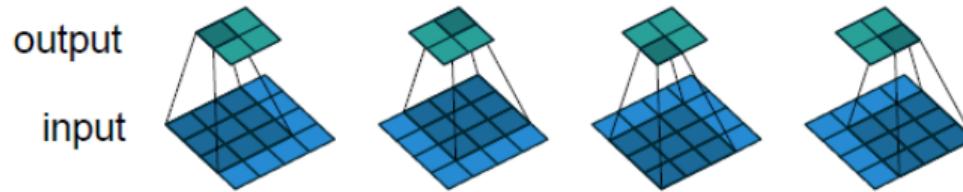
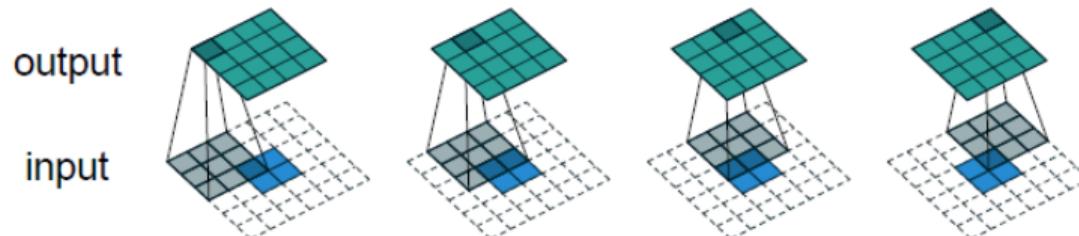


Figure 2.1: (No padding, unit strides) Convolving a  $3 \times 3$  kernel over a  $4 \times 4$  input using unit strides (i.e.,  $i = 4$ ,  $k = 3$ ,  $s = 1$  and  $p = 0$ ).

## Transposed Convolution (emulated with direct convolution):



Dumoulin, Vincent, and Francesco Visin. "[A guide to convolution arithmetic for deep learning](#)." *arXiv preprint arXiv:1603.07285* (2016).

# DILATED CONVOLUTION

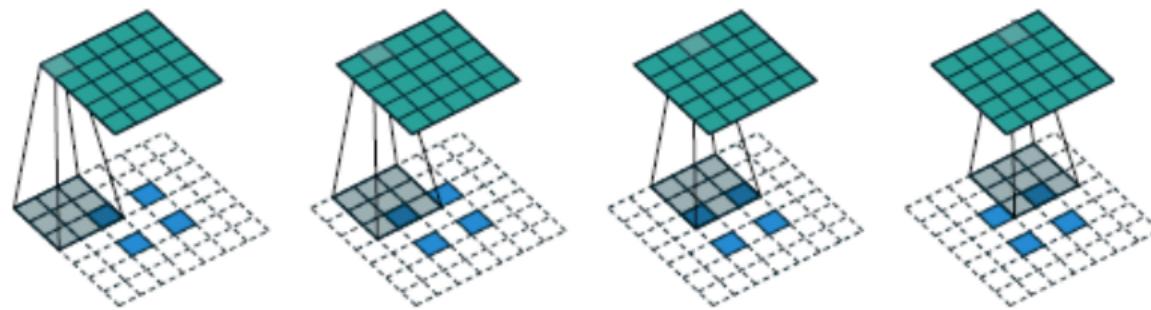


Figure 4.5: The transpose of convolving a  $3 \times 3$  kernel over a  $5 \times 5$  input using  $2 \times 2$  strides (i.e.,  $i = 5, k = 3, s = 2$  and  $p = 0$ ). It is equivalent to convolving a  $3 \times 3$  kernel over a  $2 \times 2$  input (with 1 zero inserted between inputs) padded with a  $2 \times 2$  border of zeros using unit strides

Dumoulin, Vincent, and Francesco Visin. "[A guide to convolution arithmetic for deep learning](#)." arXiv preprint arXiv:1603.07285 (2016).

# REMEMBER IMAGE INVARIANCE CHALLENGE

Object can appear anywhere in the image.

Image A

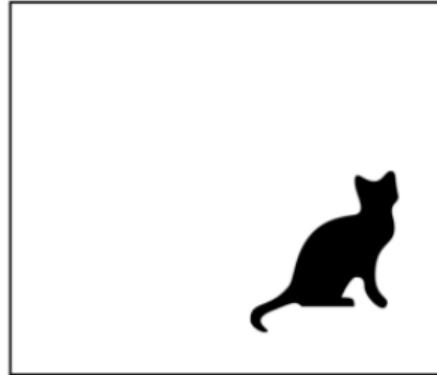
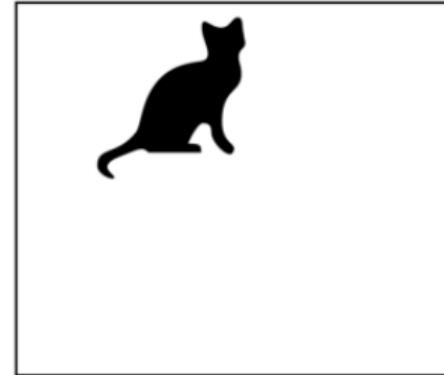


Image B

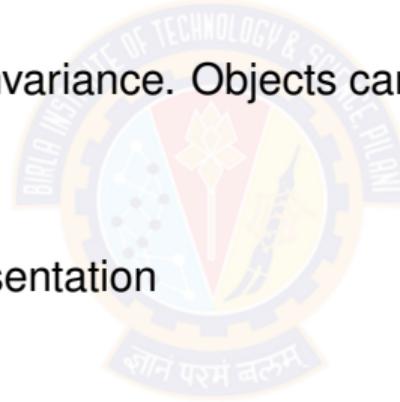


Model predicts “cat” for both images!

Source: <https://penkovsky.com/neural-networks/day5/>

# POOLING (POOL) LAYER

- Helps to achieve image invariance. Objects can appear anywhere in the Image and still get detected.
- Used for sub-sampling
- Reduce the size of representation
- Speed up computation



# POOLING

- Apply a filter of size  $f$  and stride  $s$ . Padding is not applied.
- The output dimensions will be

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

- Fixed computation
- Gradient descent is not applied as there are not parameters to be learned.
- Pooling is applied on each of the channels.

# POOLING LAYER

- Two types
  - ① Max pooling
  - ② Average pooling

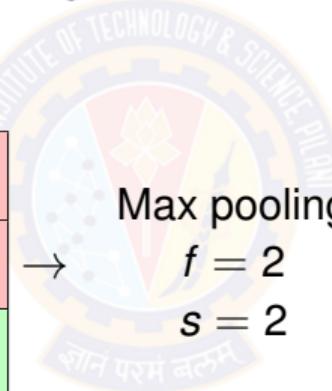


# MAX POOLING

- Take the maximum of the sub region that is considered.

Input  $4 \times 4$

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



Output

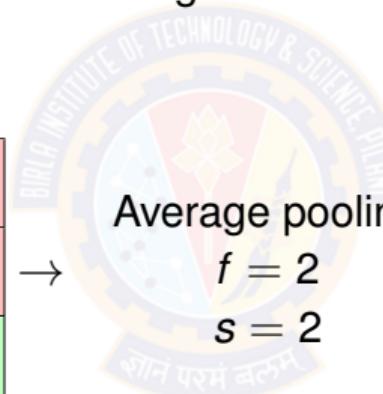
$2 \times 2$

6	9
7	8

# AVERAGE POOLING

- Compute the average of the sub region that is considered.

Input $4 \times 4$			
2	4	7	4
6	6	9	8
3	4	8	3
7	5	3	6



Average pooling

$$f = 2$$

$$s = 2$$

Output

$$2 \times 2$$

4.5	7
4.75	5

# CONVOLUTION + MAX POOLING

Input Image  $6 \times 6$

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

Filter 1

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

Filter 2

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

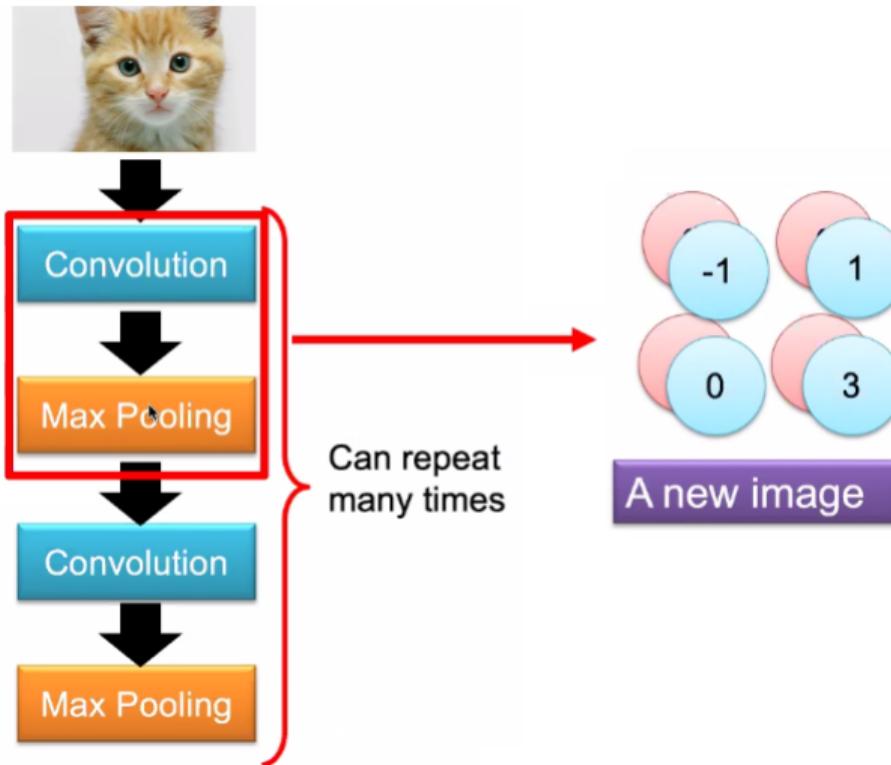
Convolve s=1

3	-1	-3	-1
-3	-1	0	-3
-3	-2	0	1
3	-2	-2	-1
-1	-1	-1	-1
-1	-1	-2	-1
-1	-1	-2	1
-1	0	-4	3

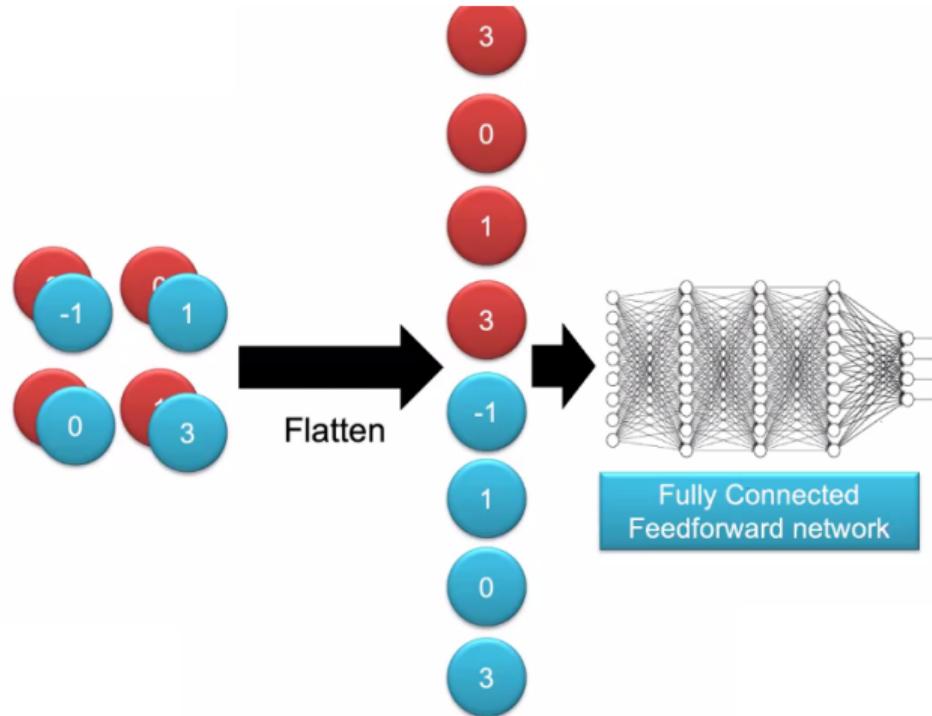
Maxpool f=2, s=2

3	0
3	1
-1	-1
0	3

# CONVOLUTION + MAX POOLING



# FLATTEN



# IMAGE IS HIGHLY NON-LINEAR

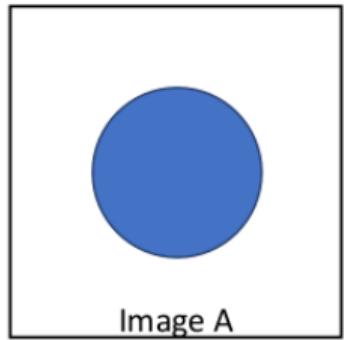


Image A

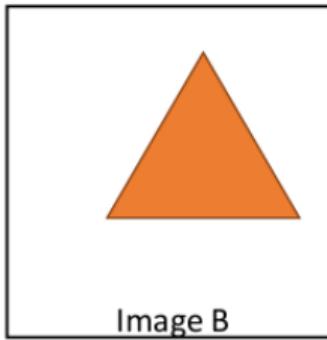
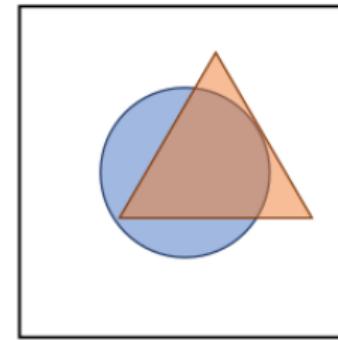
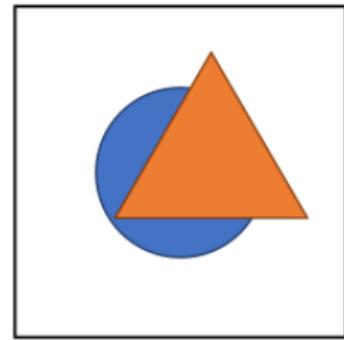


Image B



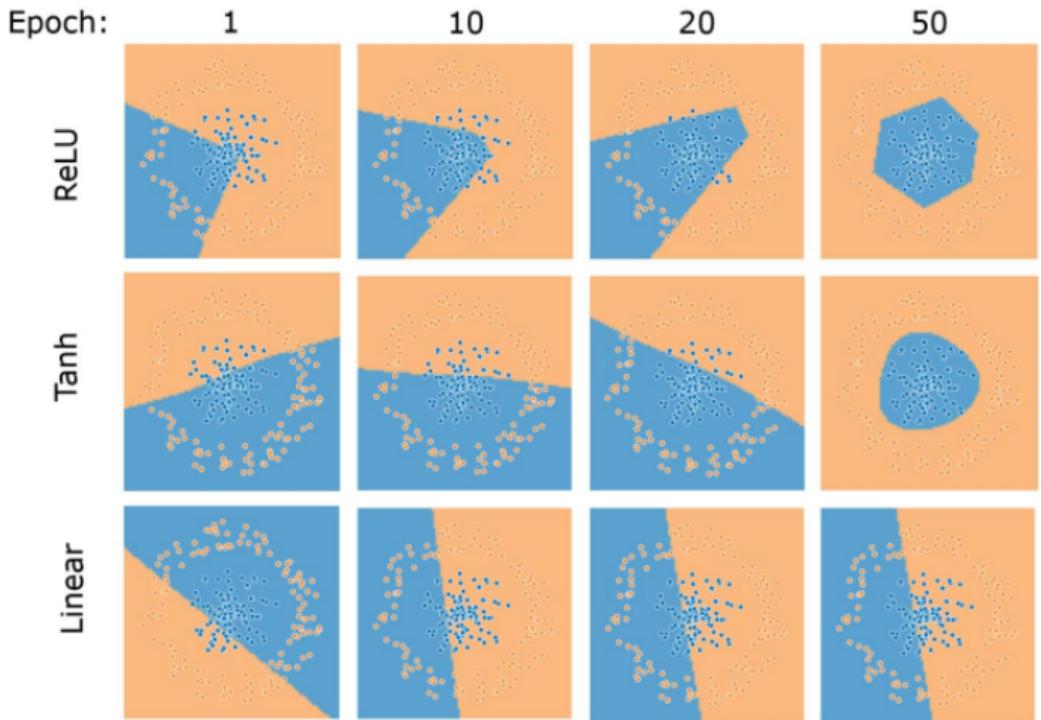
Linear combination  
 $A+B$



Actual combination  
 $A+B$

# RECTIFIED LINEAR UNIT (ReLU)

- Very fast
- Gradients
- Piece-wise linear approximation



# RELU IN IMAGES

Ground truth image



Standard fully-connected net

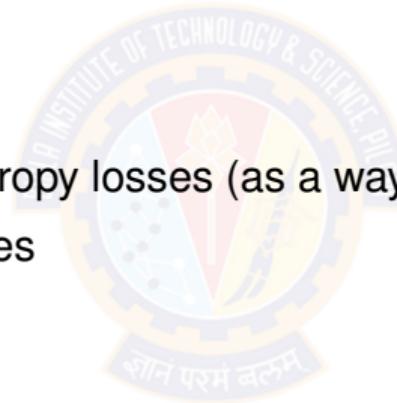


With Positional Encoding



# LOSS FOR TRAINING

- Mean Squared Error
- Binary Cross Entropy
- Multi class Cross Entropy
- Class-weighted cross-entropy losses (as a way to counter imbalance)
- L1 (MAE), L2 (MSE) losses
- Focal, Dice, Huber loss
- Margin ranking losses
- Generalized Intersection-over-Union (GloU) loss
- Connectionist Temporal Classification (CTC) loss



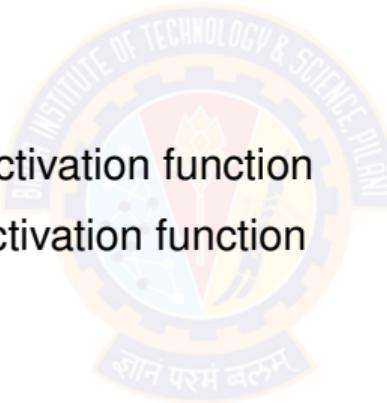
# REGULARISATION AND OPTIMIZATION

- Drop out
- Batch Normalisation
- RMSProp, Adam or any other optimizer can be used.



# CLASSIFICATION

- $p \in [0, 1]$  - use Sigmoid activation function
- k classes - use softmax activation function



# IN THIS SEGMENT

- 1 COMPUTER VISION TASKS
- 2 IMAGE BASICS
- 3 IMAGE CLASSIFICATION CHALLENGES
- 4 CONVOLUTIONAL NEURAL NETWORK
- 5 CONVOLUTION
- 6 PUTTING IT TOGETHER



# PUTTING EVERYTHING TOGETHER

- ① Identify the task as classification or regression. This helps in identify the output size.
- ② Collect the  $N$  images. This helps in identifying input size as  $W \times H \times C$ .
- ③ Feature extraction
  - ▶ Decide on the number of Convolution layers. This can be derived by looking at how large or small pattern we are interested in.
  - ▶ Decide on the kernel size and number of kernels in each layer. Smaller patterns in images require smaller kernels. Number of kernels can be derived from the non-linearity.
  - ▶ Use ReLU activation function for better approximation.
- ④ Classification
  - ▶ Use Fully connected layers for classification.
  - ▶ The number of (dense) layers depends on the non-linearity.
  - ▶ The last (output) layer activation is decided based on Step 1.
- ⑤ Train and Test on the available dataset.

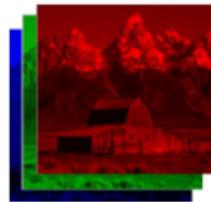
# CONVOLUTIONAL NEURAL NETWORK



- Input: 1 or 3 images
  - Grey scale or color
  - Will assume color to be generic

# CONVOLUTIONAL NEURAL NETWORK

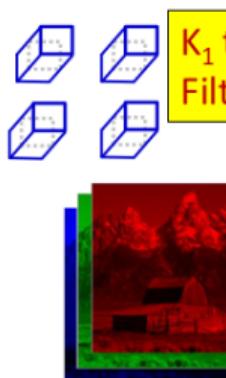
- Input is convolved with a set of  $K_1$  filters or kernels.
  - ▶ Typically,  $K_1$  is a power of 2 e.g 2,4,8,16,32,...
  - ▶ Filters are typically  $5 \times 5 \times 3$ ,  $3 \times 3 \times 3$  or  $1 \times 1 \times 3$  (better notation).
  - ▶ Typical stride is 1 or 2.



$I \times I$  image

# CONVOLUTIONAL NEURAL NETWORK

- Input is convolved with a set of  $K_1$  filters or kernels.
  - ▶ Typically,  $K_1$  is a power of 2 e.g 2,4,8,16,32,...
  - ▶ Filters are typically 5x5, 3x3 or 1x1.
  - ▶ Filters should be small enough to capture fine features, particularly important for scaled down images.



Parameters to choose:  $K_1$ ,  $L$  and  $S$

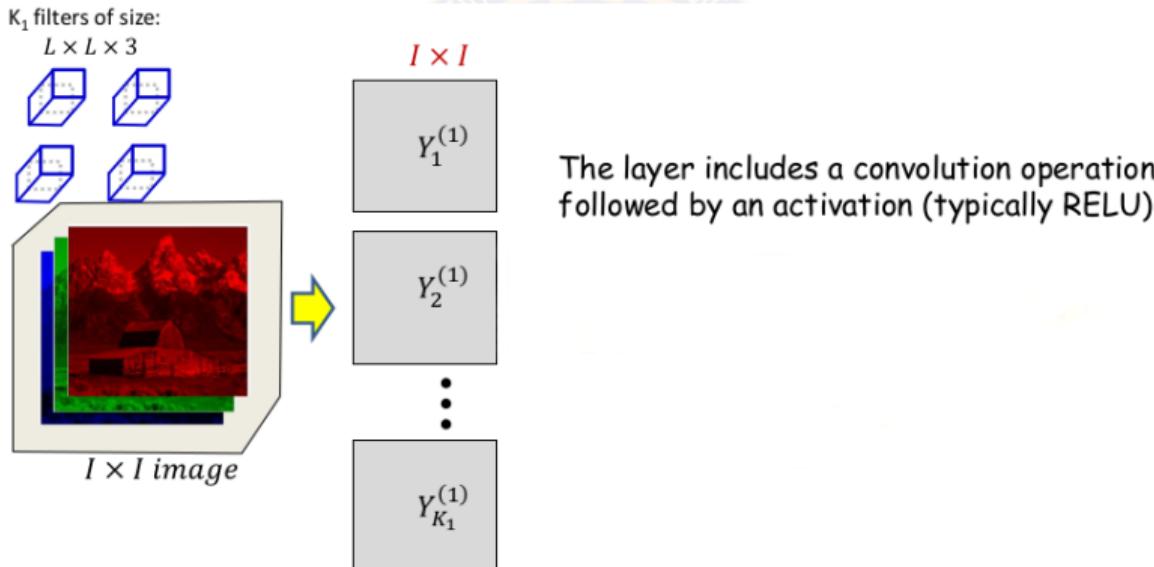
1. Number of filters  $K_1$
2. Size of filters  $L \times L \times 3 + \text{bias}$
3. Stride of convolution  $S$

Total number of parameters:  $K_1(3L^2 + 1)$

# CONVOLUTIONAL NEURAL NETWORK

- First Convolutional layer:

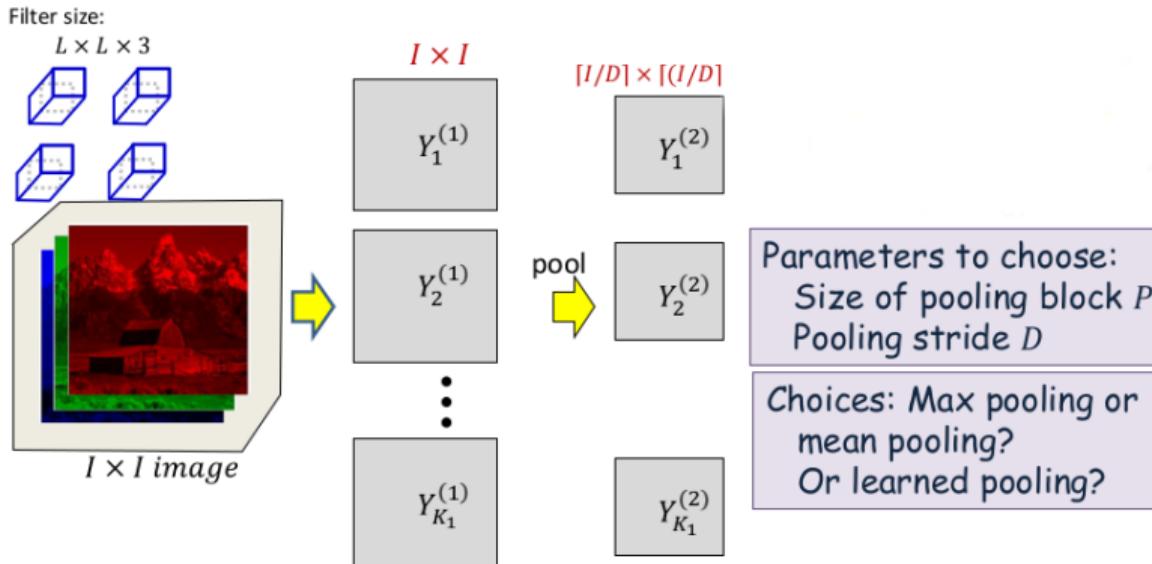
- ▶ Several convolutional filters, which are 3D.
- ▶ Convolution is followed typically by a ReLU activation.
- ▶ Each filter creates a singly 2D feature map.



# CONVOLUTIONAL NEURAL NETWORK

- First Pooling layer:

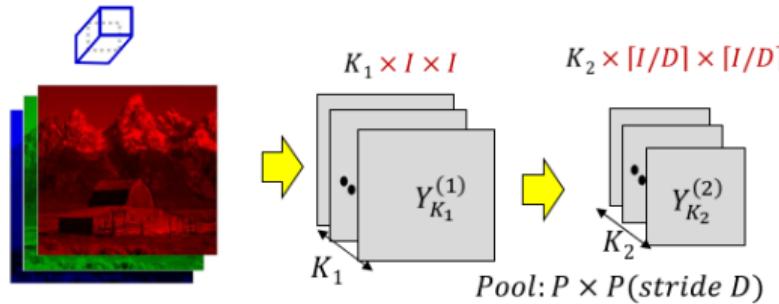
- ▶ For each  $P \times P$  block of each map, pool down to a single value.



# CONVOLUTIONAL NEURAL NETWORK

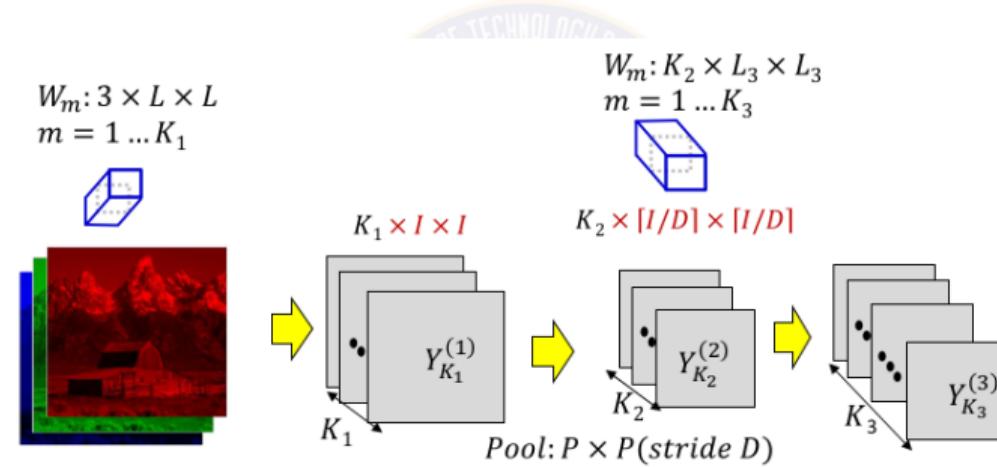
- Better drawing

$$W_m: 3 \times L \times L$$
$$m = 1 \dots K_1$$



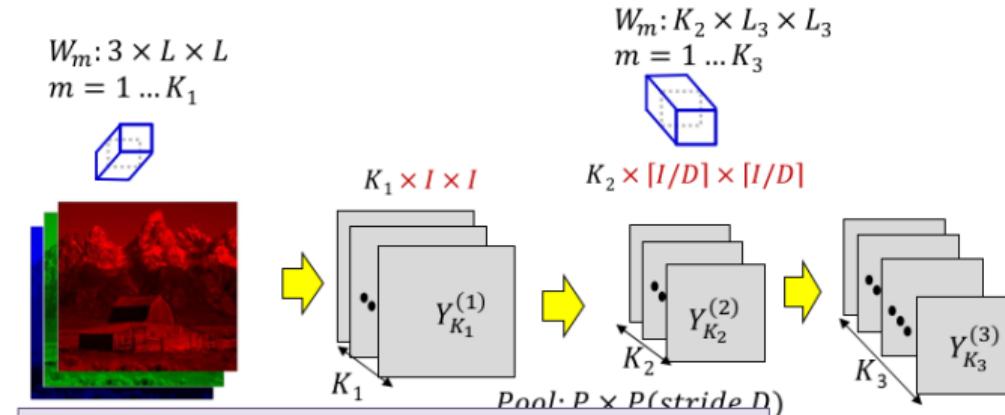
# CONVOLUTIONAL NEURAL NETWORK

- Second Convolutional layer"  $K_2$  3D filters resulting in  $K_2$  feature maps.



# CONVOLUTIONAL NEURAL NETWORK

- Second Convolutional layer



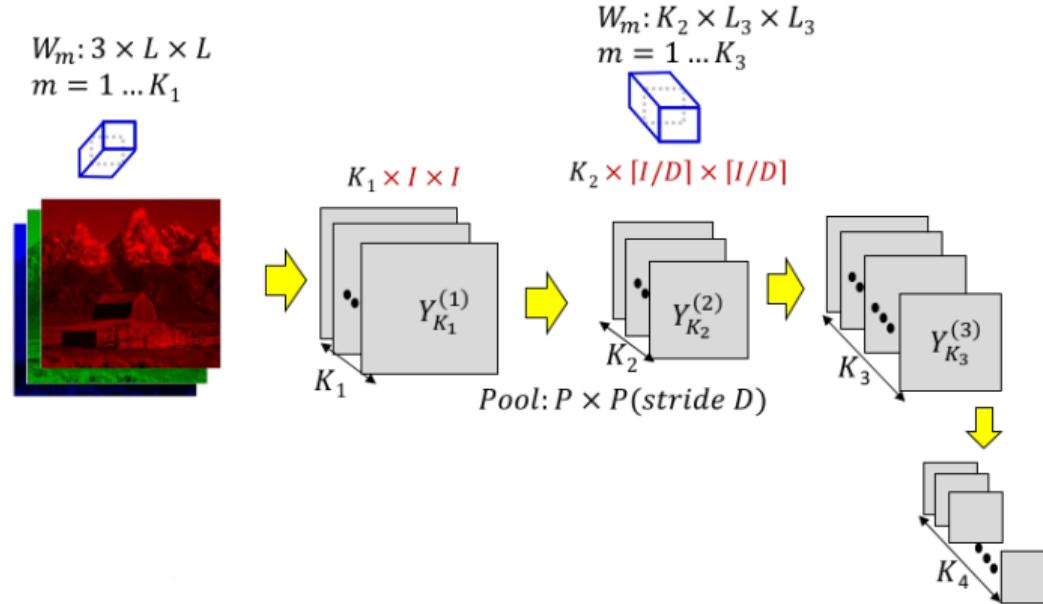
Parameters to choose:  $K_3, L_3$  and  $S_3$

1. Number of filters  $K_3$
2. Size of filters  $L_3 \times L_3 \times K_2 + bias$
3. Stride of convolution  $S_3$

Total number of parameters:  $K_3(K_2L_3^2 + 1)$   
All these parameters must be learned

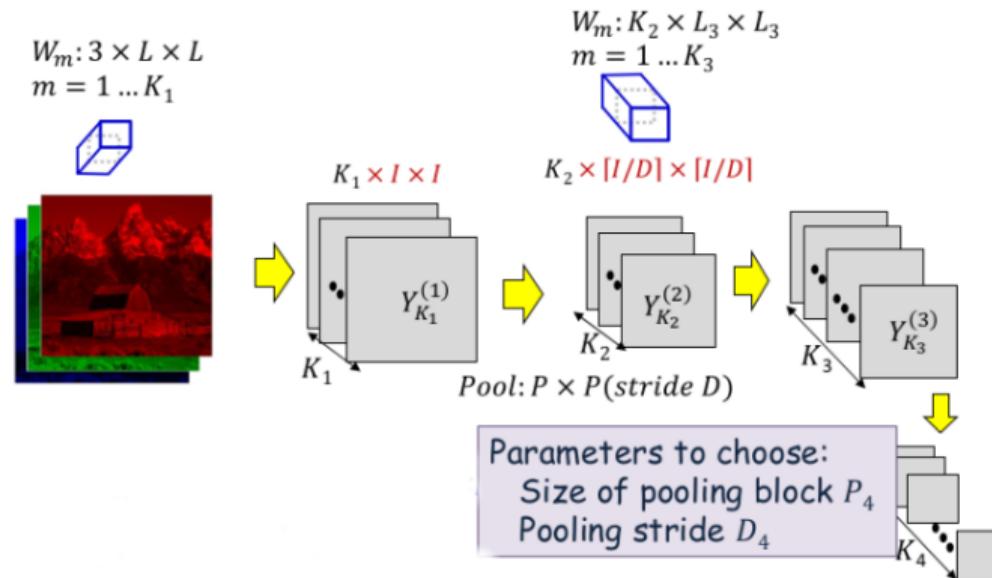
# CONVOLUTIONAL NEURAL NETWORK

- **Second Pooling layer:**  $K_2$  pooling operations results in  $K_2$  reduced 2D maps.



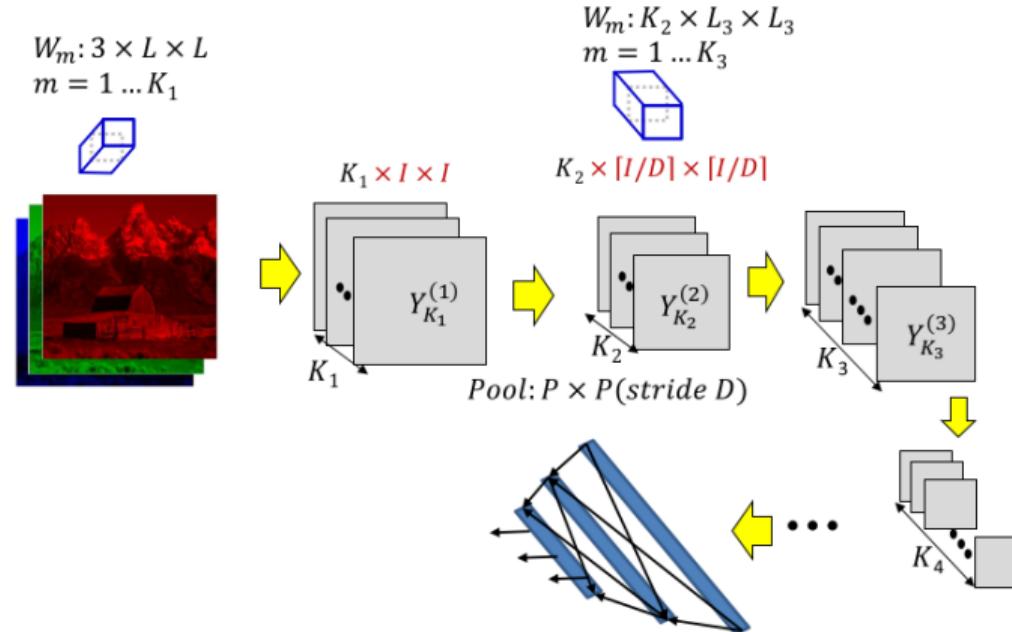
# CONVOLUTIONAL NEURAL NETWORK

- Second Pooling layer



# CONVOLUTIONAL NEURAL NETWORK

- **Fully connected layer:** This continues for several layers until the final convolved feature maps are fed to MLP.



# CONVOLUTIONAL NEURAL NETWORK

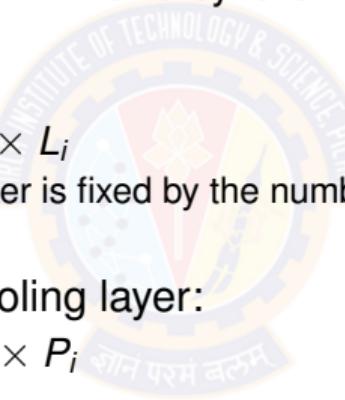
## Size of Layers

- Each convolution layer with stride 1 typically maintains the size of the image.
  - ▶ With appropriate zero padding.
  - ▶ If performed without zero padding it will decrease the size of the input.
- Each convolution layer will generally increase the number of maps from the previous layer.
  - ▶ Increasing layers reduces the amount of information lost by subsequent downsampling.
- Each pooling layer with stride D decreases the size of the maps by a factor of D.
- Filters within a layer must all be the same size, but sizes may vary with layer. Similarly for pooling D may vary with layer.
- In general the number of convolutional filters increases with layers.

# CONVOLUTIONAL NEURAL NETWORK

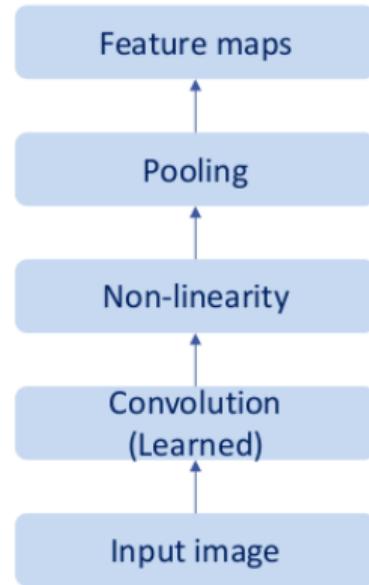
## Parameters to choose (design choices)

- Number of convolutional and downsampling layers
  - ▶ And arrangement (order in which they follow one another)
- For each convolution layer:
  - ▶ Number of filters  $K_i$
  - ▶ Spatial extent of filter  $L_i \times L_i$ 
    - ★ The “depth” of the filter is fixed by the number of filters in the previous layer  $K_{i-1}$ .
  - ▶ The stride  $S_i$
- For each downsampling/pooling layer:
  - ▶ Spatial extent of filter  $P_i \times P_i$
  - ▶ The stride  $D_i$
- For the final MLP:
  - ▶ Number of layers
  - ▶ Number of neurons in each layer



# CNN-Recap

- Neural network with specialized connectivity structure
- Feed-forward:
  - Convolve input
  - Non-linearity (rectified linear)
  - Pooling (local max)
- Supervised training
- Train convolutional filters by back-propagating error
- Convolution over time



## References

- ① Deep Learning with Python by Francois Chollet.  
<https://livebook.manning.com/book/deep-learning-with-python/>
- ② Deep Learning by Ian Goodfellow, Yoshua Bengio, Aaron Courville  
<https://www.deeplearningbook.org/>

**Thank You!**