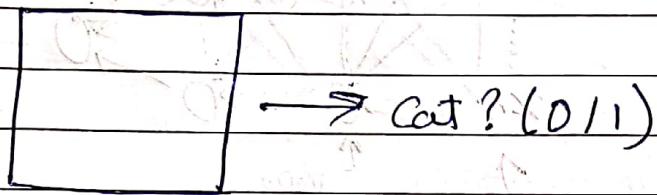


# Convolutional Neural Networks.

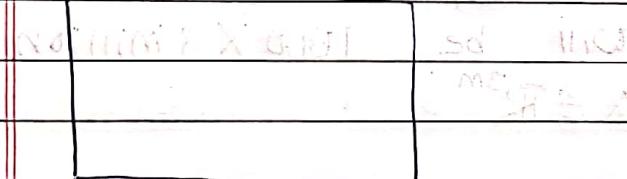
## \* Computer Vision.

Some Example of Computer Vision

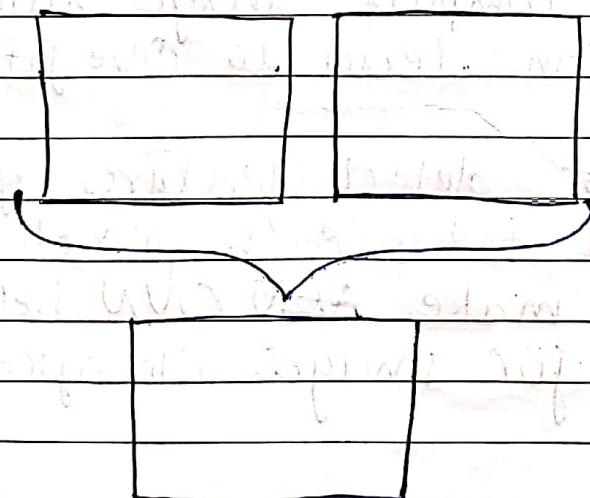
⇒ image classification.



⇒ object detection.

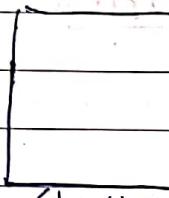


⇒ Neural Style Transfer



## \* Deep learning on large images.

Small image

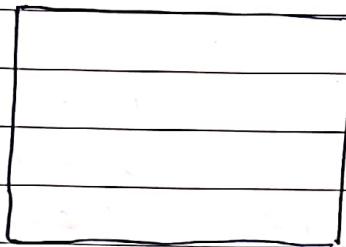


$\rightarrow$  cat? (0/1)

$64 \times 64 \times 3$

$$= \underline{12288}$$

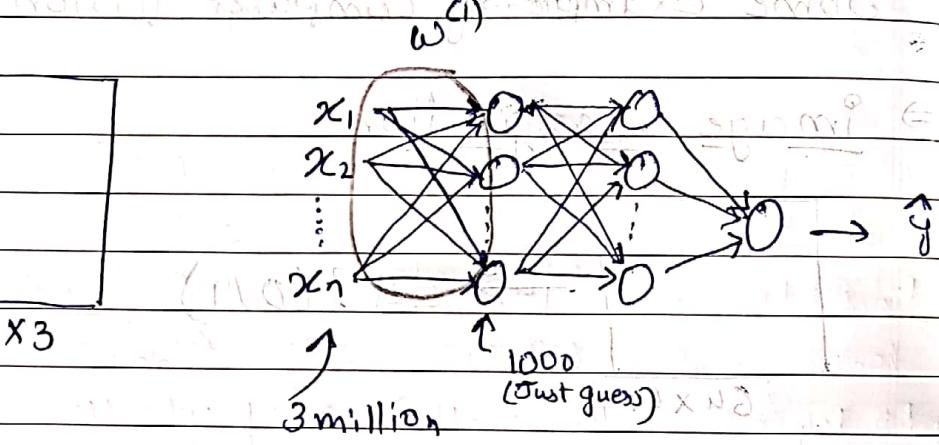
Large image



$1000 \times 1000 \times 3$

$$= 3 \text{ million}$$

$x \in \mathbb{R}^{3m}$



Then Total Number of weight  
on matrix  $w^{(1)}$

will be  $1000 \times 3 \text{ million}$

$$x \in \mathbb{R}^{3m}$$

$\therefore$  if we use Artificial Neural Network  
for solving image classification problem,  
model has to compute large number  
of parameter will can lead to Overfitting

$\rightarrow$  Also, ANN can not detect features of  
images, instead it takes each pixel  
as features, which make ~~ANN~~ CNN better  
than ANN for images Classification  
problems.

## \* Edge detection Examples.

→ Vertical edge detection (~~Bottom edge~~ ~~Top edge~~)

"Convolution"

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

 $\downarrow$ 

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

 $=$ 

-5	-4	0	8
-10	-2	2	3
0	-2	-4	-7
-3	-2	-3	-16

$3 \times 3$  filter       $4 \times 4$

6x6 (grey image)      Kernel)

4x4 matrix Solution:

$$\rightarrow (3 \times 1) + (1 \times 1) + (2 \times 1) + (0 \times 0) + (5 \times 0) + (7 \times 0) + (8 \times 1) + (2 \times -1) = -5$$

$$\rightarrow (0 \times 1) + (5 \times 1) + (7 \times 1) + (1 \times 0) + (8 \times 0) + (2 \times 0) + (2 \times -1) + (9 \times 1) + (5 \times -1) = -4$$

$$\vdots$$

$$\vdots$$

$$\rightarrow (1 \times 1) + (6 \times 1) + (2 \times 1) + (7 \times 0) + (2 \times 0) + (3 \times 0) + (7 \times -1) + (8 \times -1) + (9 \times -1) = -16$$

→ Horizontal Edge Detection (~~Bottom edge~~ ~~Top edge~~)

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
0	0	0	10	10	10
0	0	0	10	10	10

edge detection

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

filter
 $=$ 

0	0	0	0
30	10	-10	-30
30	10	-10	30
0	0	0	0

Input image      filter      result

PPB

# \* Learning to detect edges: convolution step - 7 \*

$\begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix}$	$\begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix}$
--	--	--	---

Prewitt

Top edge filter

Prewitt

Left edge filter

Prewitt

Bottom edge filter

Prewitt

Right edge filter

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

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

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

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

$$Z = (-\text{Sobel } H \times 8) + (\text{Sobel } V \times 8) + (\text{Sobel } D \times 8) + (\text{Sobel } A \times 8)$$

$$H = (1 \cdot \text{Top edge filter}) + (4 \cdot \text{Left edge filter}) + (6 \cdot \text{Bottom edge filter}) + (4 \cdot \text{Right edge filter})$$

$\begin{bmatrix} -3 & -10 & -3 \\ 0 & 0 & 0 \\ 3 & 10 & 3 \end{bmatrix}$	$\begin{bmatrix} -3 & 0 & 3 \\ -10 & 0 & 10 \\ -3 & 0 & 3 \end{bmatrix}$	$\begin{bmatrix} 3 & 10 & 3 \\ 0 & 0 & 0 \\ -3 & -10 & -3 \end{bmatrix}$	$\begin{bmatrix} 3 & 0 & -3 \\ 0 & 0 & -10 \\ 3 & 0 & -3 \end{bmatrix}$
--	--	--	---

Scharr

Top edge filter

Scharr

Left edge filter

Scharr

Bottom edge filter

Scharr

Right edge filter

\*

$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$
---	---

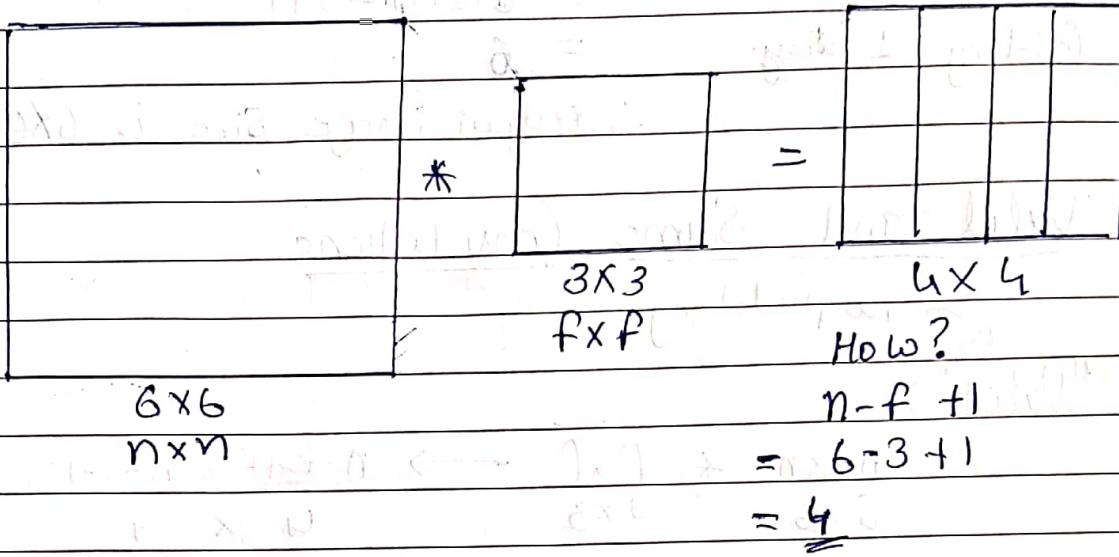
$w_1$	$w_2$	$w_3$
$w_4$	$w_5$	$w_6$
$w_7$	$w_8$	$w_9$

4x4

6x6

\* we treat these  
as parameters (it is sometimes  
Very Powerful idea)

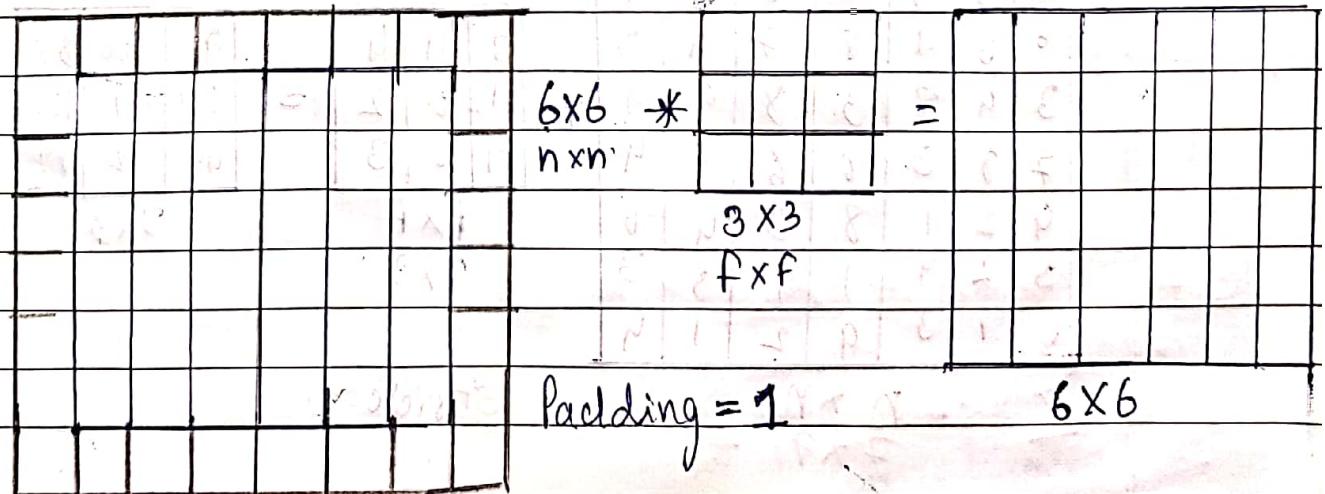
## \* Padding



→ two down sides of this

- 1) everytime we apply convolution operator, image shrink, ~~we~~ so we've gone from  $6 \times 6$  to  $4 \times 4$  then  $2 \times 2$ , which ~~is~~ we do not want every time, especially when image size is small.
- 2) Pixels on the corners around the edges are used much less in output, throwing away lots of information near the edge of the image.

\* Two solve both of this problems, ~~we~~ use padding  
→ add addition border to the image.



$$\begin{aligned}
 n \times n &= 6 \times 6 & \therefore n+2p-f+1 \\
 f \times f &= 3 \times 3 & = 6+2(1)-3+1 \\
 \text{Padding} &= 1 \quad \text{filter} & = 6
 \end{aligned}$$

$\therefore$  output image size is  $6 \times 6$

### \* Valid and Same Convolutions.

$\leftarrow$  no padding

"Valid":

$$\begin{matrix}
 n \times n & * & f \times f & \rightarrow & n-f+1 \times n-f+1 \\
 6 \times 6 & & 3 \times 3 & & 4 \times 4
 \end{matrix}$$

"Same": Pad so that output size is the same as the input size.

$$n+2p-f+1 \times n+2p-f+1$$

$$n+2p-f+1 = n \Rightarrow p = \frac{f-1}{2}$$

$\Rightarrow$  filter is usually odd matrix number  
matrix ;  $3 \times 3, 5 \times 5, 7 \times 7$ .

### \* Strided Convolutions

2	3	7	4	6	2	9	3	9	4	91	100	83	
6	6	9	8	7	4	3	*	1	0	2	69	91	127
3	4	8	3	8	9	7	-1	0	3	44	72	74	
7	8	3	6	6	3	9							
4	2	1	8	3	4	6							
3	2	4	1	9	8	3							
0	1	3	9	2	1	4							

$$\begin{matrix}
 n \times n & = & 7 \times 7 & \text{Stride} = 2 \\
 7 \times 7 & & &
 \end{matrix}$$

Padding  $\rightarrow P$   
Stride  $\rightarrow S$

$$= \left[ \frac{n+2P-f+1}{S} \right] \times \left[ \frac{n+2P-f+1}{S} \right]$$

$$= \left[ \frac{7+2(0)-3+1}{2} \right] \times \left[ \frac{7+2(0)-3+1}{2} \right] = 3 \times 3$$

$\rightarrow$  Stride is the number of pixels shifts over the input matrix when the stride is 1 then we move the filter to 1 pixel at a time. when the stride is 2 then we move the filter to 2 pixel at a time and so on.

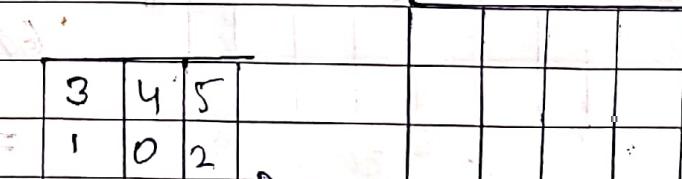
$\rightarrow$  if  $\left[ \frac{n+2P-f+1}{S} \right]$  part is not integer but fraction,

then, ~~the~~ round-off value will be counted

$$\text{i.e } \lfloor z \rfloor = \text{floor}(z); \text{ here } z = \frac{n+2P-f+1}{S}$$

★ Technical note on Cross-Correlation Vs. convolution.

3	4	5
1	0	2
-1	0	7

$\star$   $=$  

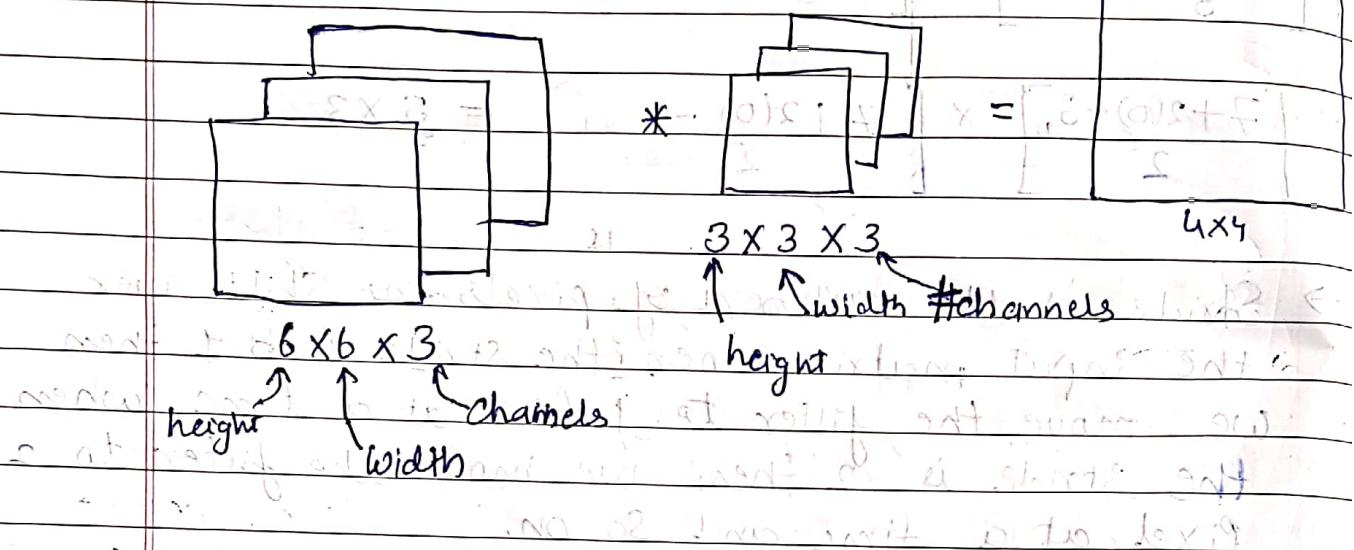
7	2	5
9	0	4
1	1	3

This operation called cross-correlation or convolution operator.

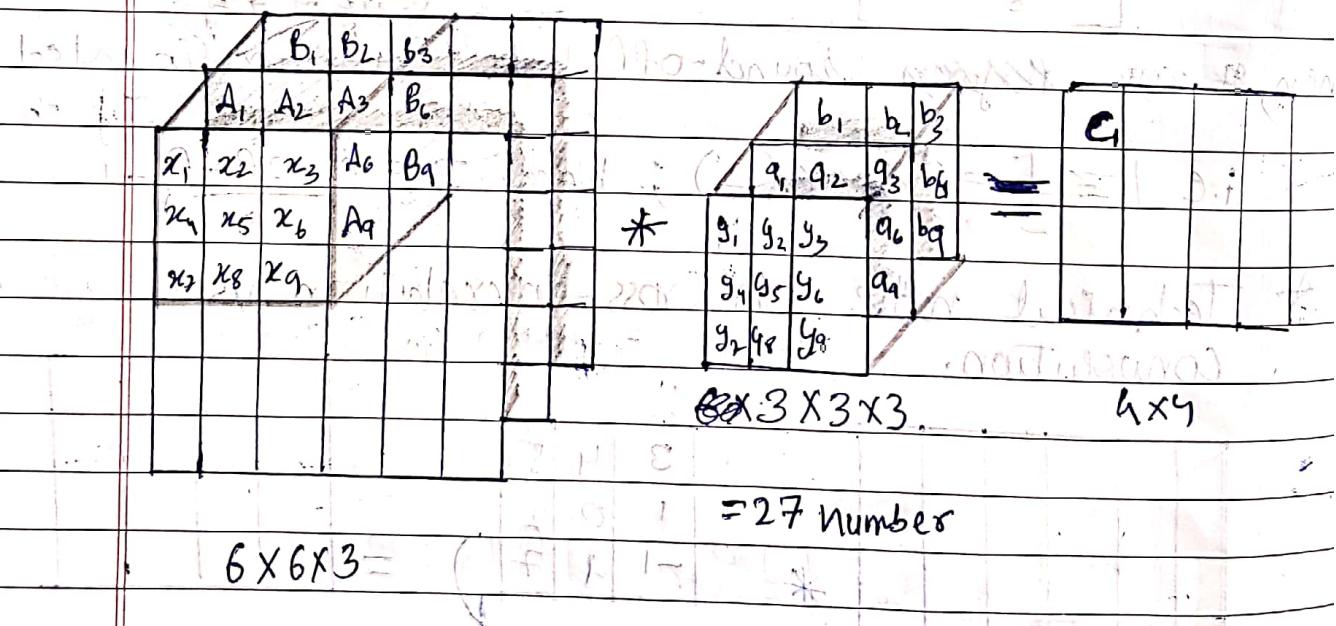
$$\text{As: } (A * B) * C = A * (B * C)$$

## \* Convolutions over Volumes.

\* Convolutions on RGB images



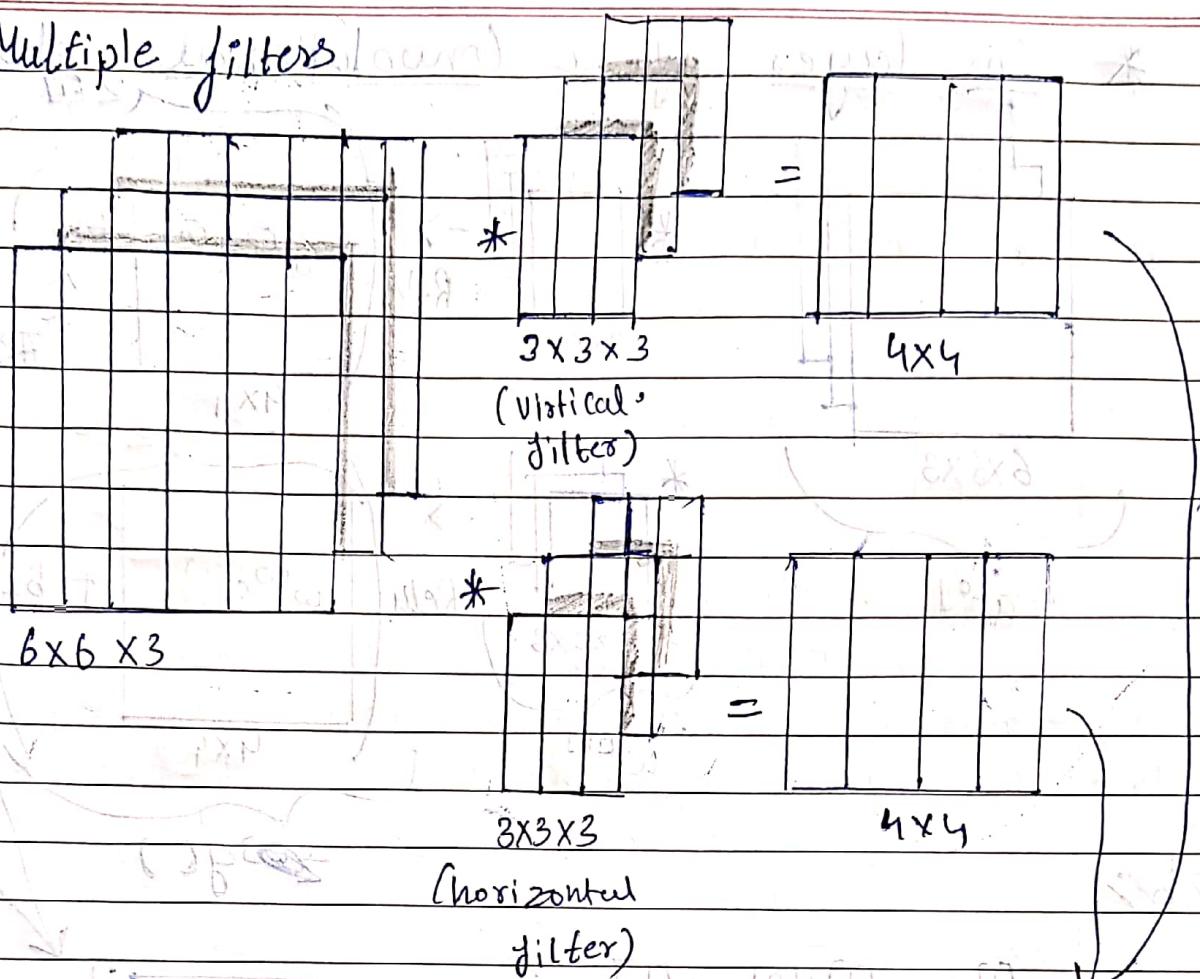
=> How it works?



$$C_1 = (x_1 * y_1) + (A_1 * a_1) + (B_1 * b_1)$$

$$C_1 = \sum_{i=1}^9 (x_i * y_i) + \sum_{i=1}^9 (A_i * a_i) + \sum_{i=1}^9 (B_i * b_i)$$

→ Multiple filters (local)



Summary:

$$n \times n \times n_c * f \times f \times f_c$$

$$6 \times 6 \times 3 \quad 3 \times 3 \times 3$$

have to be same



4x4x2

$$\rightarrow n - f + 1 \times n - f + 1 \times n_f$$

$$4 \times 4 \times 2$$

number of filters

# filters

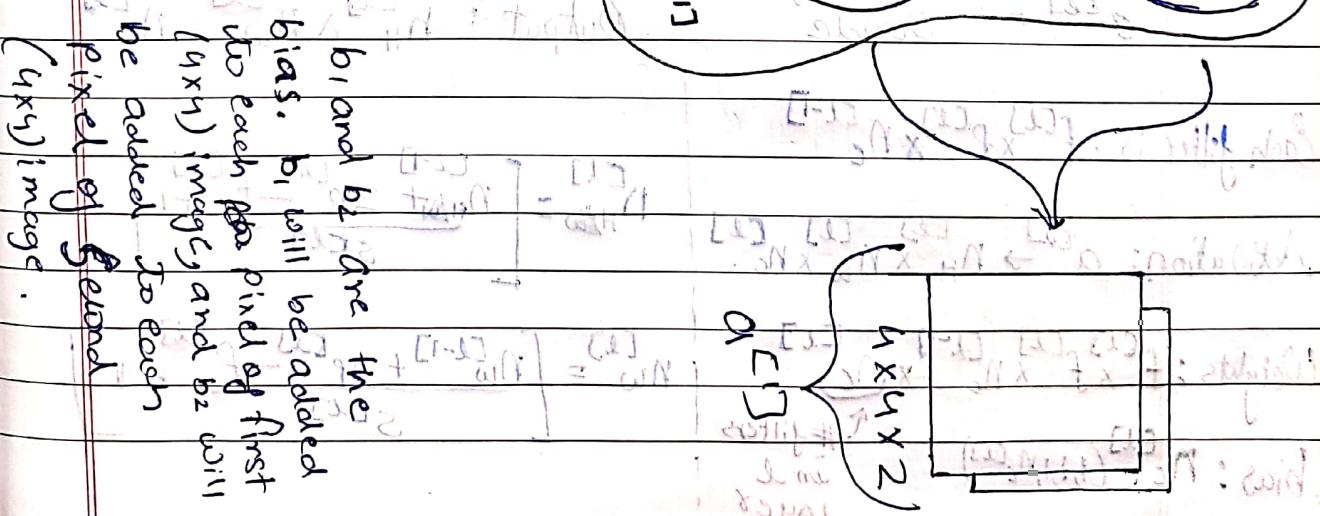
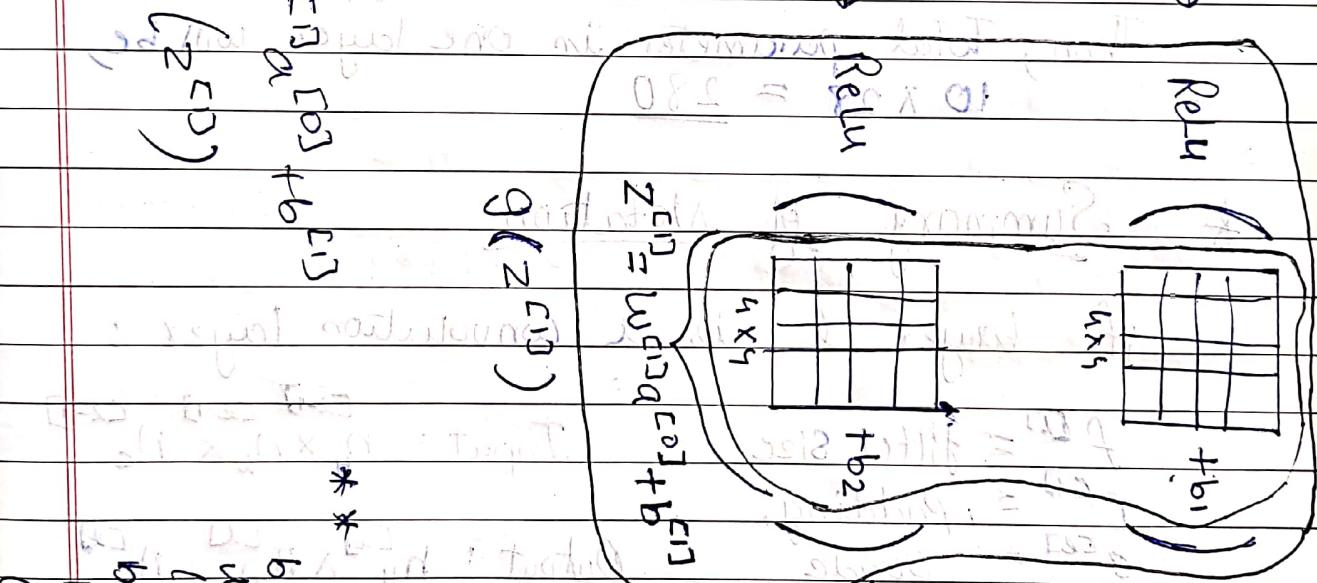
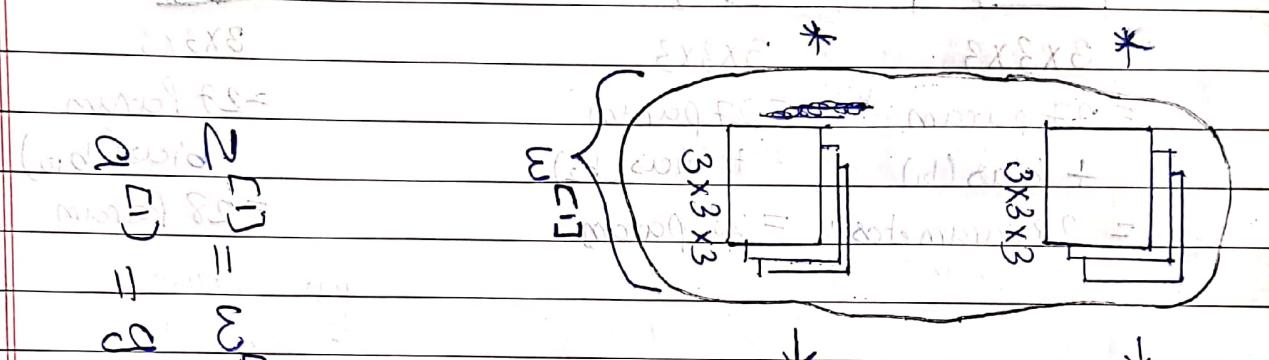
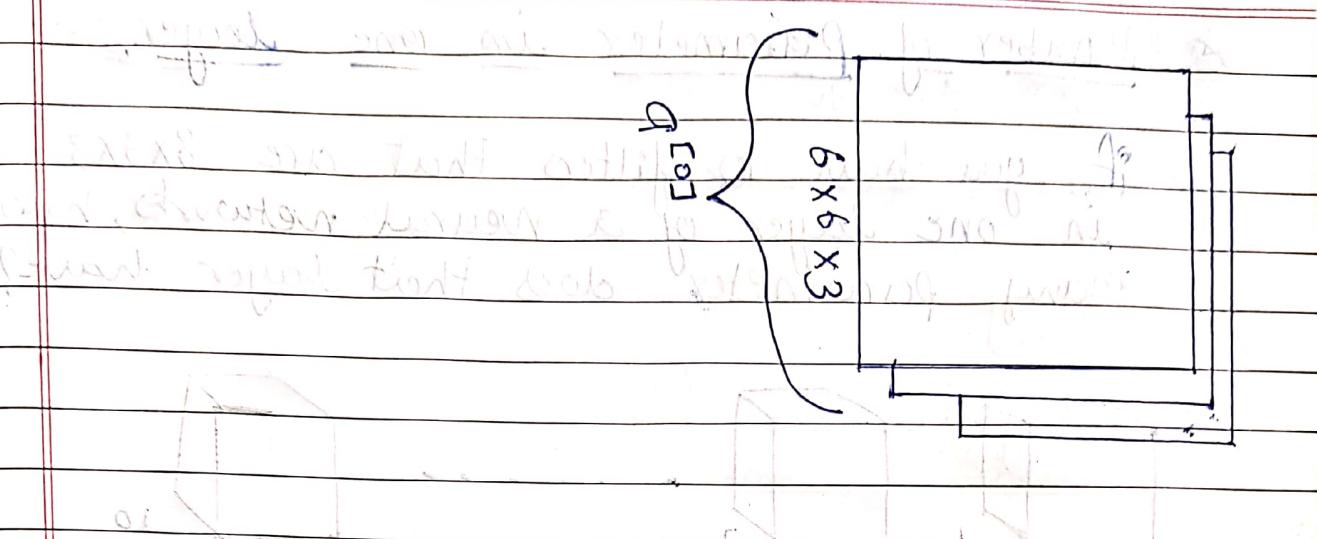
SPD GDI DINA

## Example of a layer

\* One layer of a Convolutional Network.

classmate

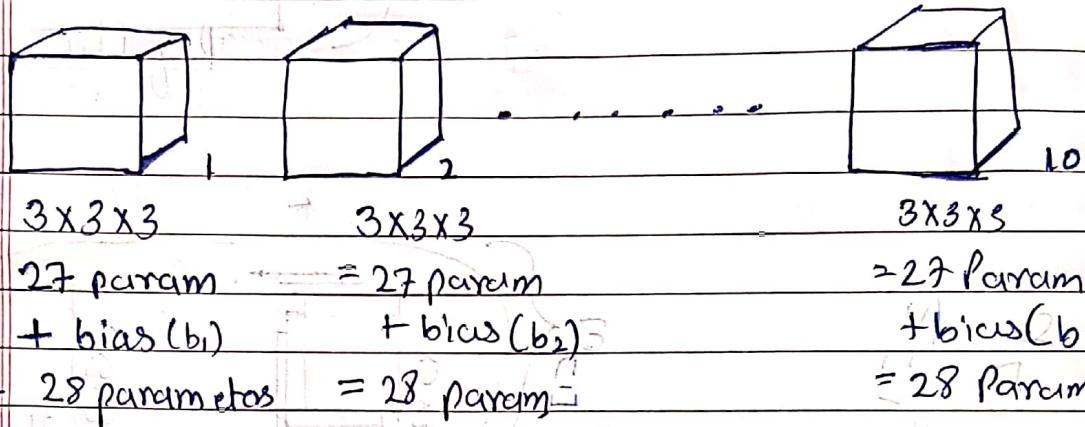
Date \_\_\_\_\_



$b_1$  and  $b_2$  are the bias.  $b_1$  will be added to each pixel of first  $(4 \times 4)$  image, and  $b_2$  will be added to each pixel of second  $(4 \times 4)$  image.

## \* Number of Parameters in one layer.

If you have 10 filters that are  $3 \times 3 \times 3$  in one layer of a neural network, how many parameters does that layer have?



Then, Total parameter in one layer will be,  
 $10 \times 28 = 280$

## \* Summary of Notation.

If layer  $l$  is a convolution layer :

$f^{[l]}$  = filter size

$p^{[l]}$  = padding

$s^{[l]}$  = stride

Input :  $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$

Output :  $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

Each filter is :  $f^{[l]} \times p^{[l]} \times n_C^{[l-1]}$

Activation:  $a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

Weights :  $f^{[l]} \times f^{[l]} \times n_C^{[l-1]} \times n_C^{[l]}$   
 $\uparrow \# \text{filters}$

bias :  $n_C^{[l]} - (1, 1, n_C^{[l]})$   
 $\uparrow \text{in } l \text{ layer}$

$$n_H^{[l]} = \left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]} + 1}{s^{[l]}} \right\rfloor$$

$$n_W^{[l]} = \left\lfloor \frac{n_W^{[l-1]} + 2p^{[l]} - f^{[l]} + 1}{s^{[l]}} \right\rfloor$$

# A Simple Convolution networks Example

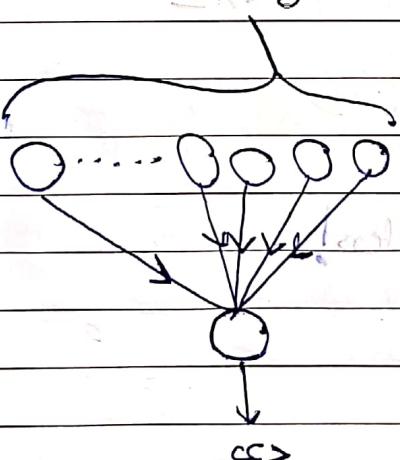
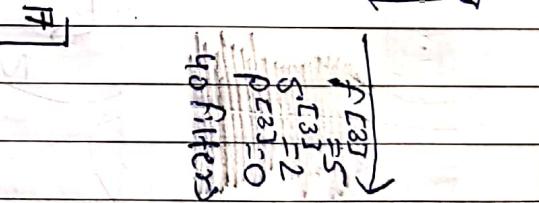
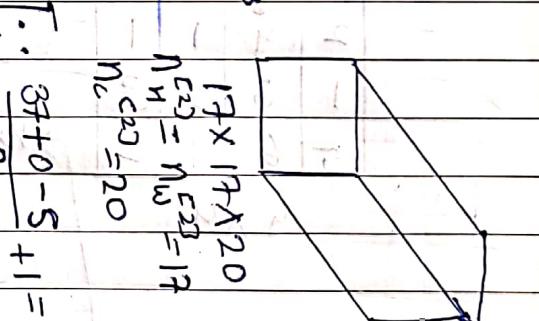
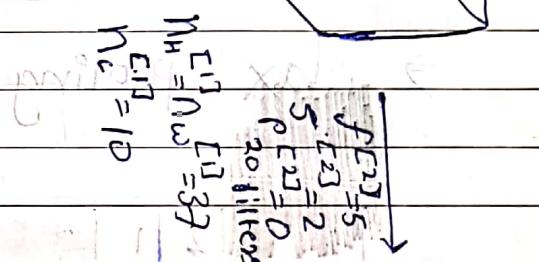
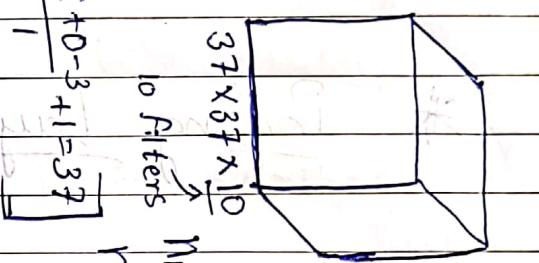
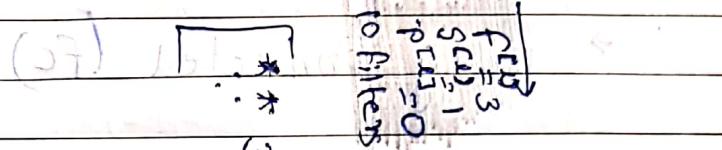
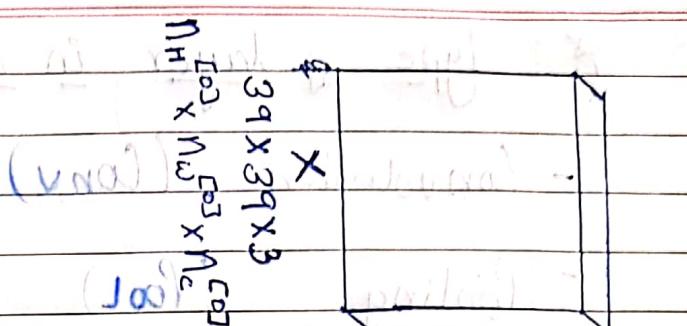
classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

## \* Example ConvNet

Input (Input image)  $\rightarrow$   $(39 \times 39 \times 3)$



\* Type of layer in a Convolutional network:

- Convolution (Conv)
- Pooling (Pool)
- Fully Connected (Fc)

### \* Pooling Layer

→ Max pooling: Calculate largest value in each patch of feature map. filter size  $(2 \times 2)$

1	3	2	1	
2	9	1	1	
1	3	2	3	
5	6	1	2	

→ hyperparameters:  $2 \times 2$

$$4 \times 4 \quad f=2$$

$$s=2$$

No parameters!

filter:  $3 \times 3$ 

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

$n=5$   
 $f=3$   
 $s=1$

$3 \times 3 \times n_c$

$$5 \times 5 \times n_c$$

$$\therefore n + 2p - f + 1$$

$$= \underline{\underline{3}}$$

Input

Output

(1, 3, 1) (1, 3, 1)

⇒ Average Pooling: Calculate Average Value in each patch of feature map.

1	3	2	1	avg	→	3.75   1.25	
2	9	1	1			4   2	
1	3	2	3	avg	$f=2$	(1, 0.5)	8-0-4
5	6	1	2		$s=2$		

## Summary of pooling

Hyperparameters:

$$n_H \times n_W \times n_c$$

f: filter size

$$\frac{n_H - f + 1}{s} \times \frac{n_W - f + 1}{s} \times n_c$$

s: stride

Max or average pooling.

p: padding (very very rare)

(most common method is zero padding)

No parameters to learn.

# Convolutional neural network Example

classmate

Data  
Page

	Activation Shape	Activation Size	# Parameters
Input :	(32, 32, 3)	3, 072	0
Conv1 ( $f=5, s=1$ )	(28, 28, 6)	4704	456
Pool1 ( $f=3, s=2$ )	(14, 14, 6)	1176	0
Conv2 ( $f=5, s=1$ )	(10, 10, 16)	1600	2416
Pool2 ( $f=3, s=2$ )	(5, 5, 16)	400	0
FC-3	(120, 1)	120	48120
FC-4	(84, 1)	84	10164
Softmax	(10, 1)	10	924

Rule :  $(f * f * \# \text{Previous channels} + \text{bias}(1)) * \# \text{Filters}$

⇒ Pooling layers : do not have parameters

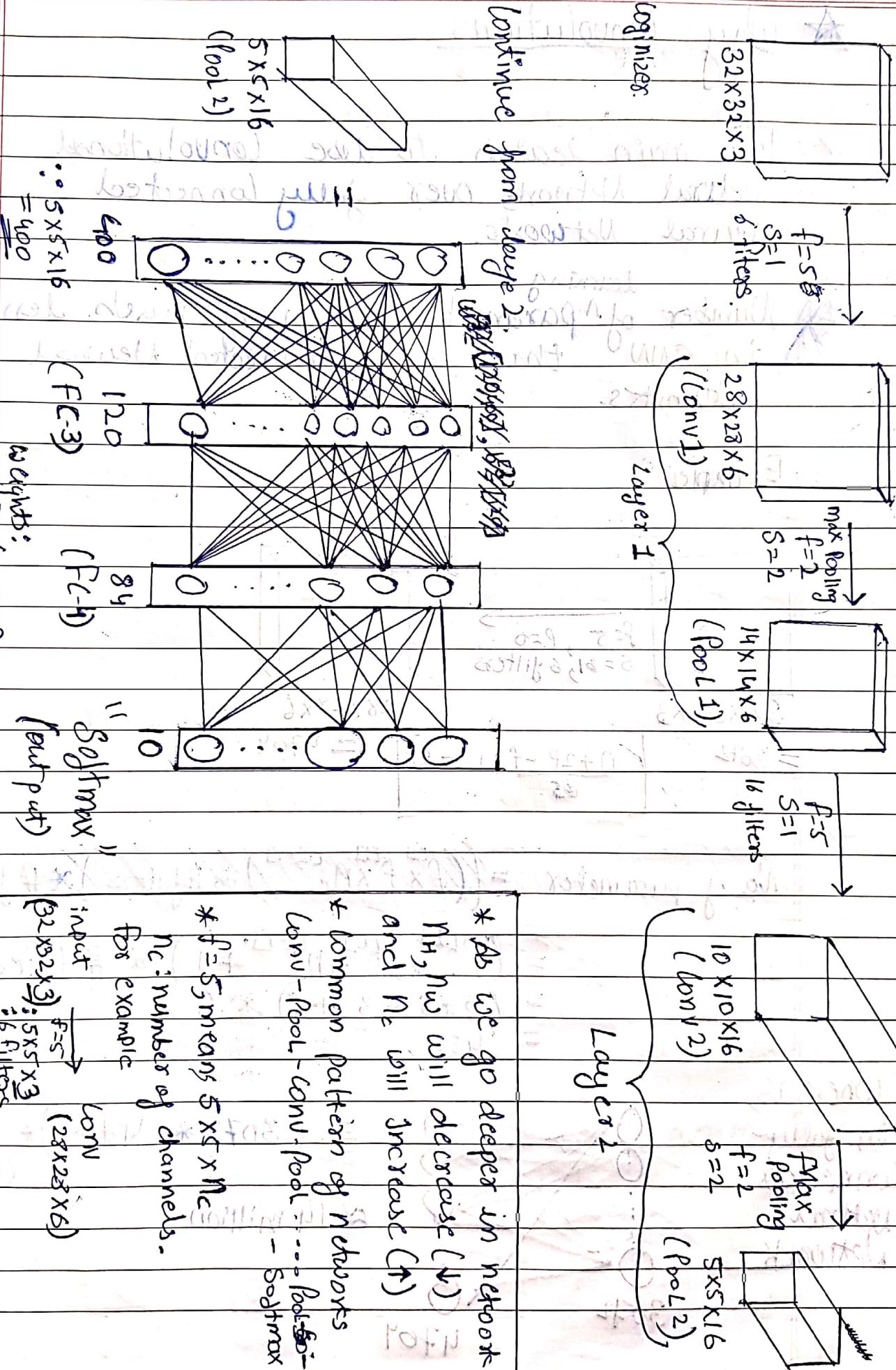
for fully connected layer

Rule:

- $(\# \text{Nodes in layer}) * (\# \text{Nodes in previous layer}) + \# \text{Nodes in layer} (\text{bias})$

# (LeNet-5)

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_  
consider  
digit Recognizer  
problem

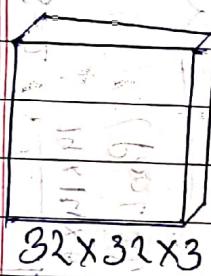


## ★ Why Convolutions

→ Two main reason to use Convolutional Neural Networks over fully connected neural networks.

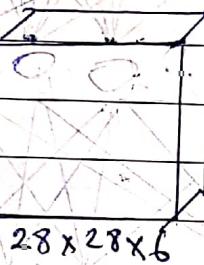
~~learning~~ Number of parameters are very much less in CNN than fully connected Neural Networks.

### Example



$$f=5, p=0 \\ s=1, 6 \text{ filters}$$

$32 \times 32 \times 3$



$28 \times 28 \times 6$

$$= 3072$$

$$n + 2p - f + 1 = 28$$

$\Theta 5$

$$= 4704$$

No. of parameter

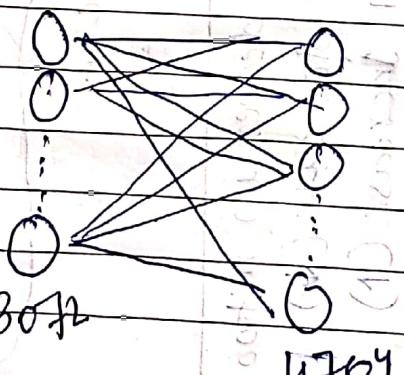
$$= (f^{[L]} * f^{[L]} * n_c^{[L-1]} + 1) * \# \text{filters} * \# \text{filters}$$

$$= (5 * 5 * 3 + 1) * 6$$

$$= 456$$

Whereas,

In fully-connected neural network



$$= (3072 * 4704) + 4704$$

(bias)

$\approx 14 \text{ million}$

1) parameter Sharing : A feature detector (Such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.

2) Sparsity of Connections : In each layer, each output value depends only on a small number of input.

\* Translation Invariance :

Using CNN methodology means that if image of cat is shifted by a few pixels, it'll still be identified as a cat.

\* Putting it together.

Training set  $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$

$x^{(i)}$  is image,  $y^{(i)}$  is label (output)

$$\text{Cost} \rightarrow J = \frac{1}{m} \sum_{i=1}^m L(y^{(i)}, \hat{y}^{(i)})$$

$\hat{y}$  is predicted label (output)

use gradient descent to optimize parameters to reduce  $J$

"x1" 100	0	0
0	0	0

$$\Delta J = \frac{\partial J}{\partial \theta} = \frac{\partial J}{\partial \theta_1} = \frac{\partial J}{\partial \theta_2} = \dots$$

## \* Case Studies

### \* why look at Case Studies? (historical)

↳ historical when technology is new  
↳ how it's developed

### \* Outline

↳ neural nets & architectural optimizations (2)

#### Classic networks:

- LeNet-5

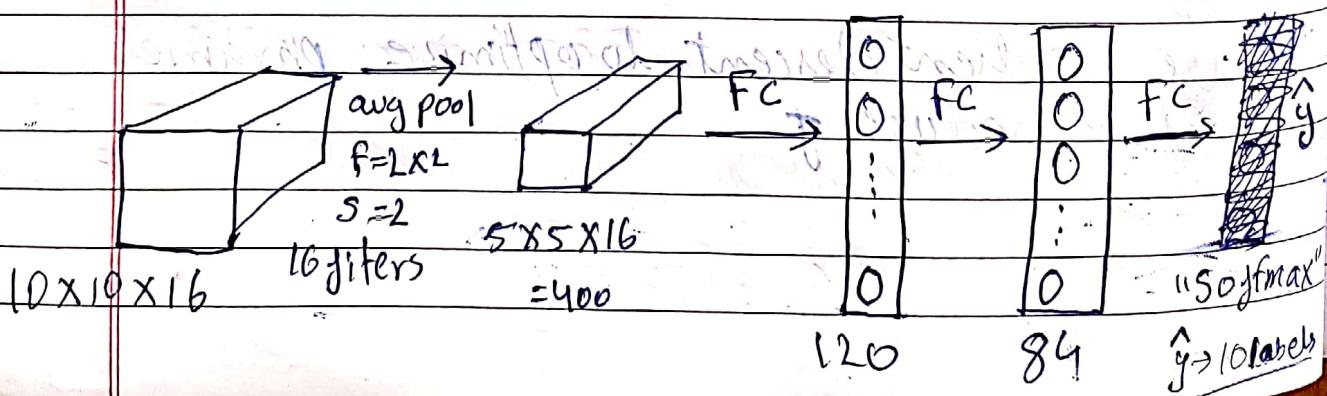
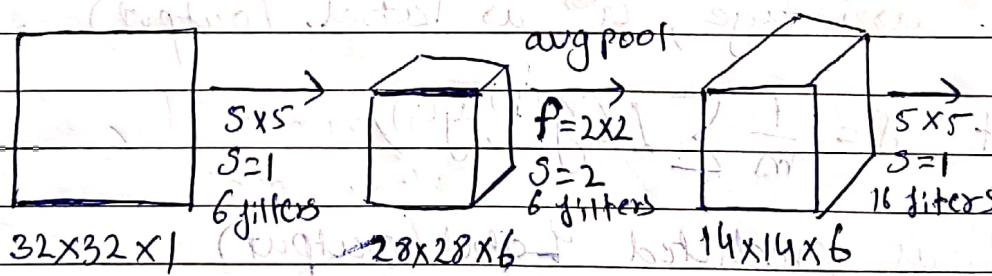
- AlexNet

- VGG

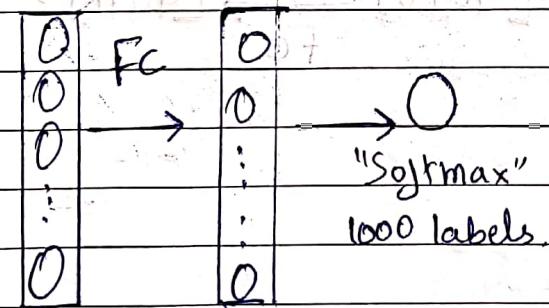
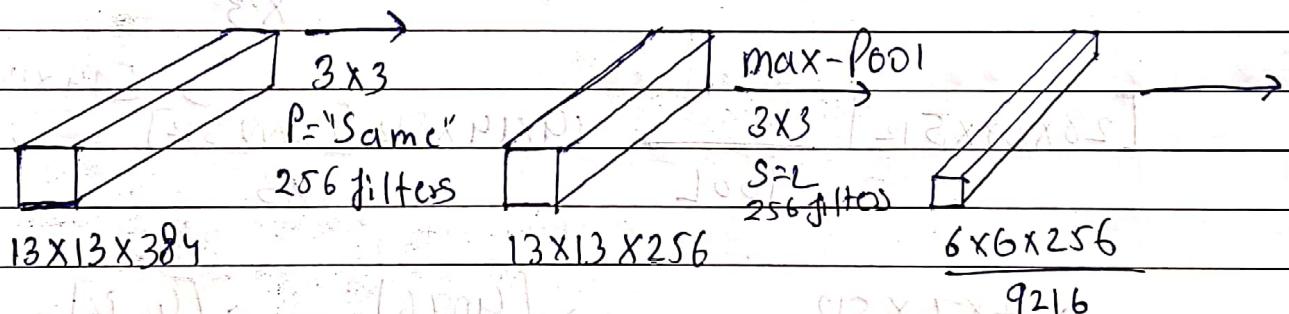
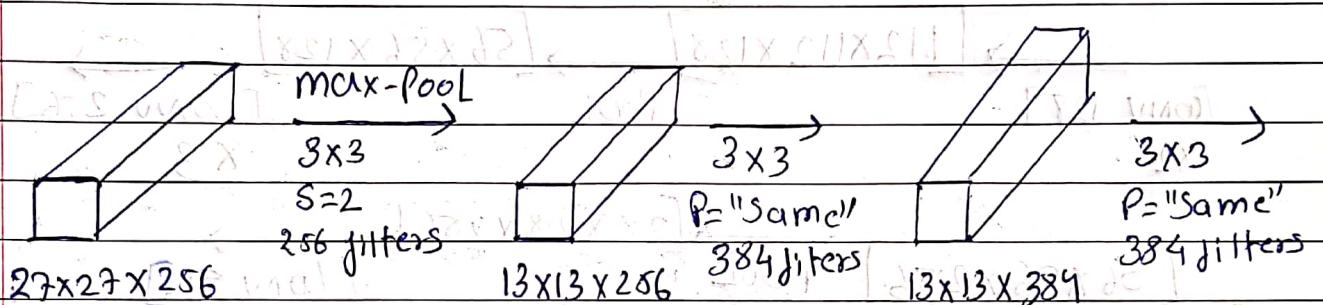
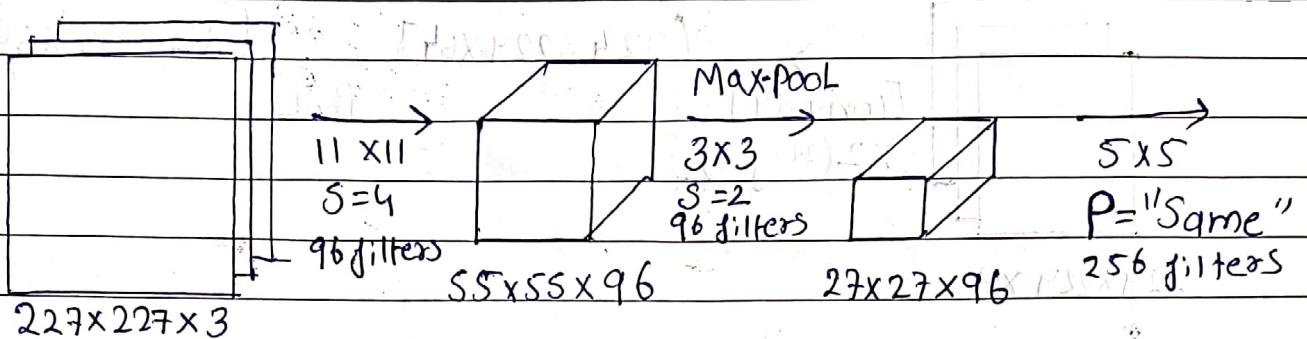
- ResNet (Residual Network, 152 layers)

- Inception

$\Rightarrow$  LeNet-5 (LeCun et al., 1998. Research Paper)



⇒ AlexNet [Krizhevsky et al., 2012. ImageNet classification with deep CNN]



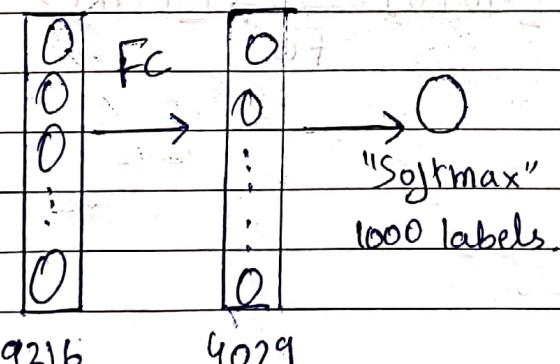
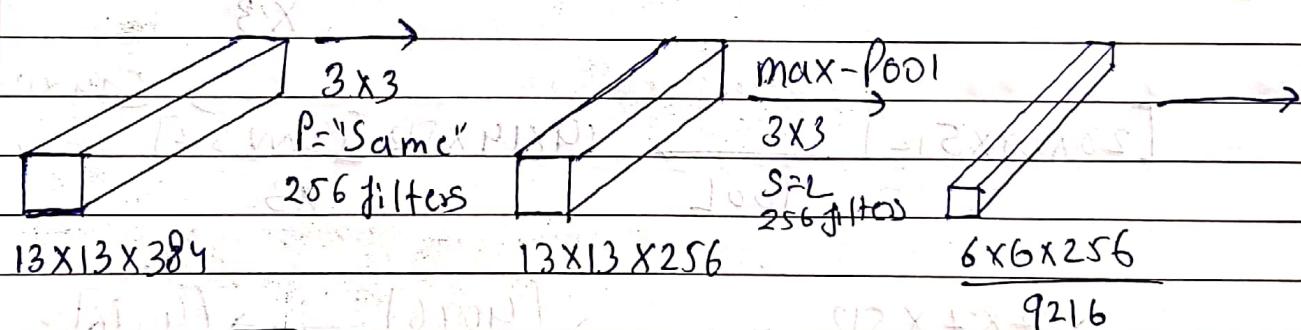
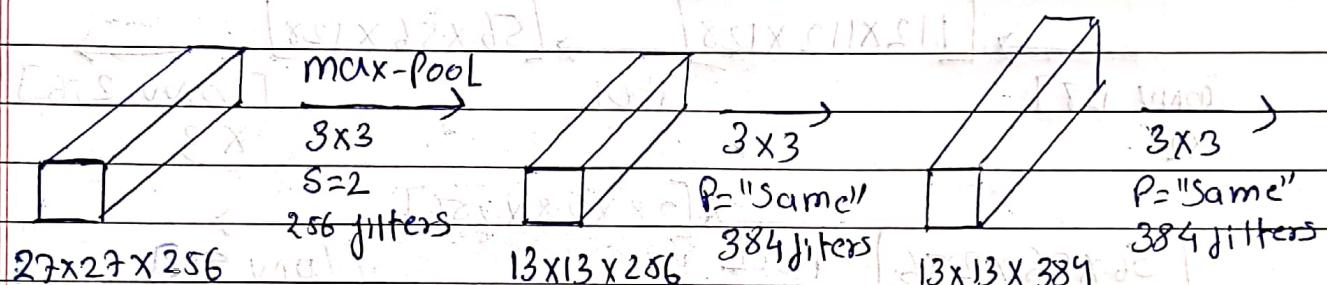
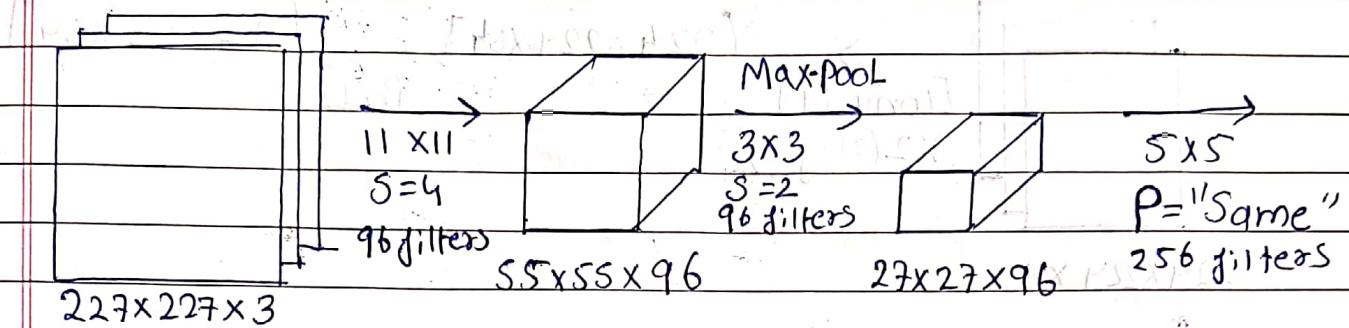
→ Similar to Lenet but much bigger ( $\approx 60^M$  parameters)  
→ ReLU as activation

289 training steps

(↓) well with rough up

(↑), N

⇒ AlexNet [Krizhevsky et al., 2012. ImageNet classification with deep CNN]

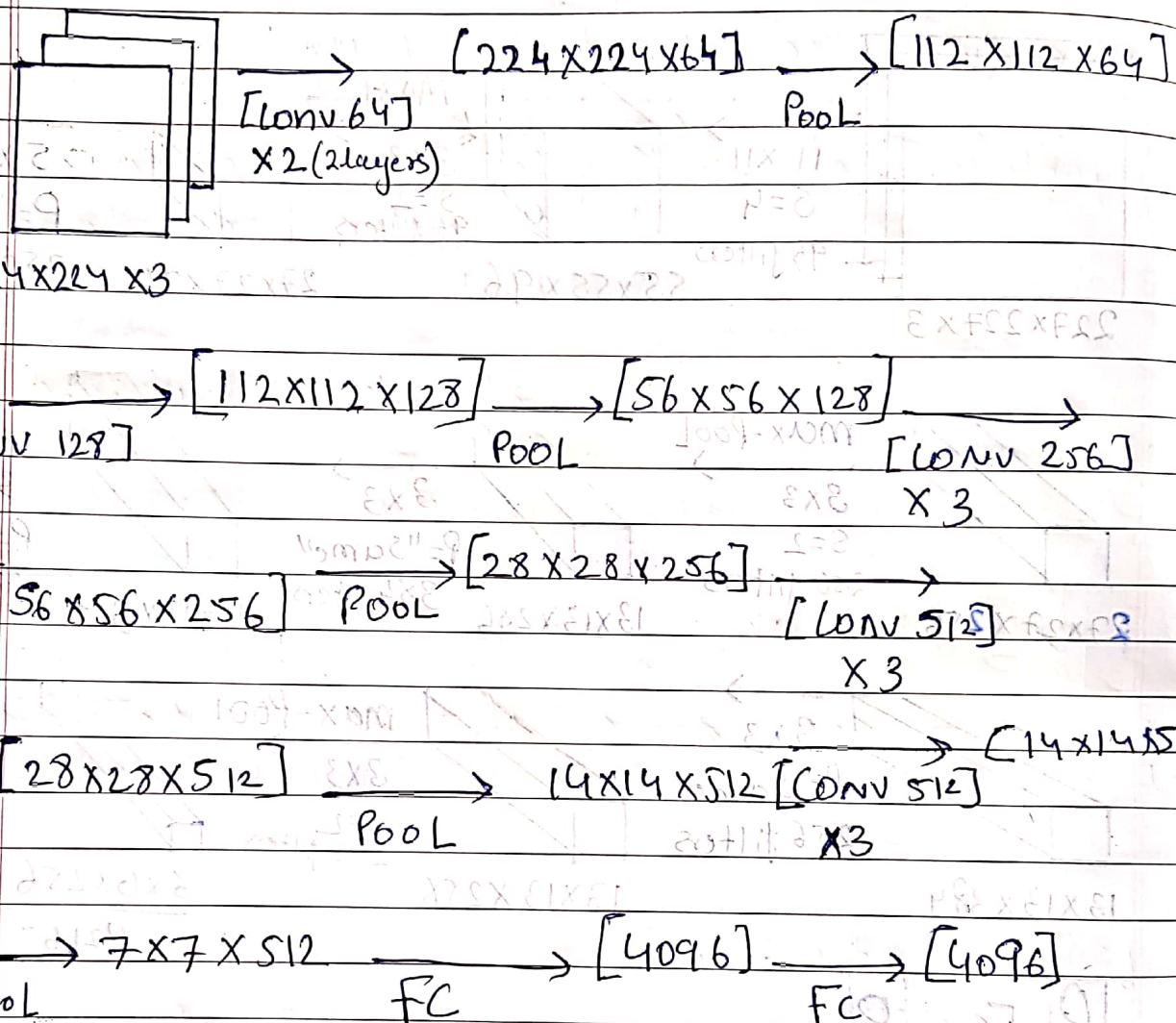


→ Similar to LeNet but  
much bigger ( $\approx 60^M$  parameters),  
→ ReLu as activation

(↓) Full with depth 08

(↑) A

→ VGG-16 [Simonyan & Zisserman 2015. Very deep CNN for Large Scale Image Recognition  
(Conv = 3x3 filters, S=1, P="Same") (Max-pool = 2x2, S=2)



→ Total  $\approx 138M$  parameters.

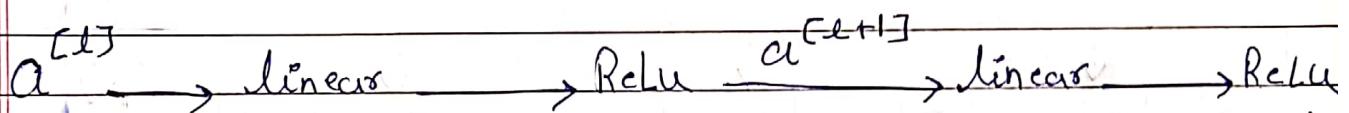
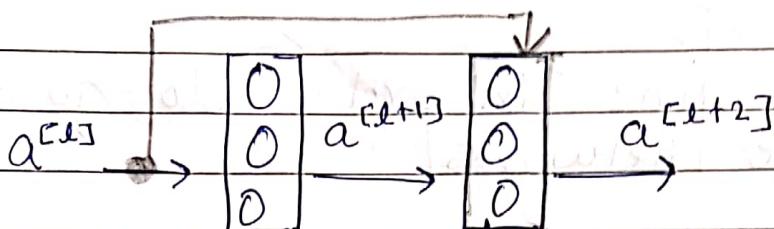
→ go deeper ~~↑~~  $n_h, n_w (\downarrow)$

$n_c (\uparrow)$

Caser studies:

## Residual Networks (ResNet)

Residual block.

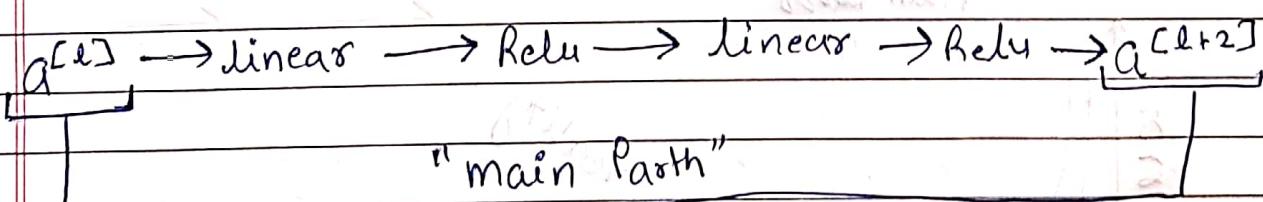


$$z^{[l+1]} = W^{[l+1]} a^{[l]} + b^{[l+1]}$$

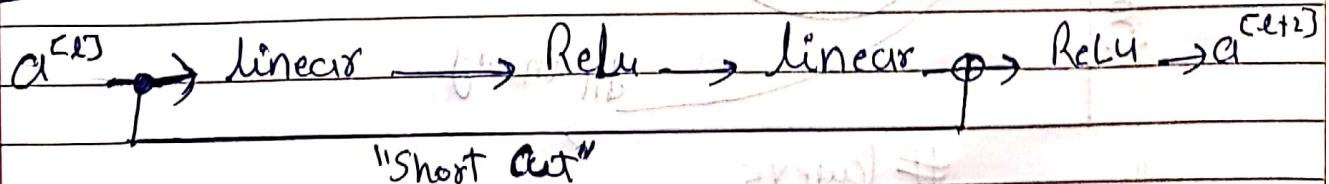
$$z^{[l+2]} = W^{[l+2]} a^{[l+1]} + b^{[l+2]}$$

$$a^{[l+1]} = g(z^{[l+1]})$$

$$a^{[l+2]} = g(z^{[l+2]})$$



In ResNet



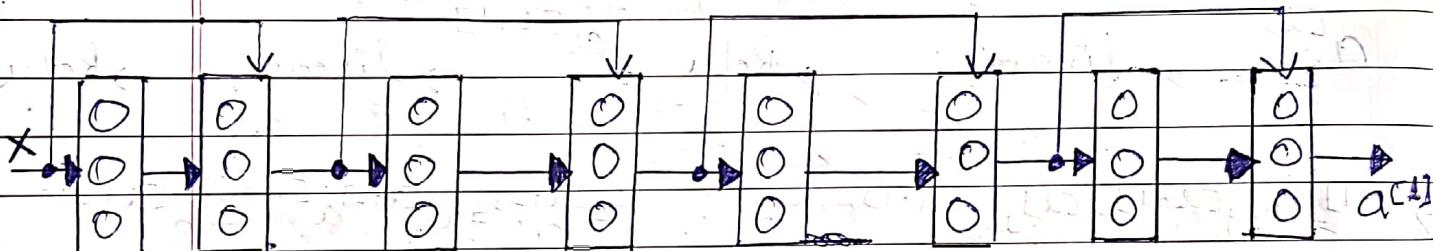
therefore Now equation

$$a^{[l+2]} = g(z^{[l+2]} + a^{[lT]})$$

This is called "Short cut" or "Skip Connection"

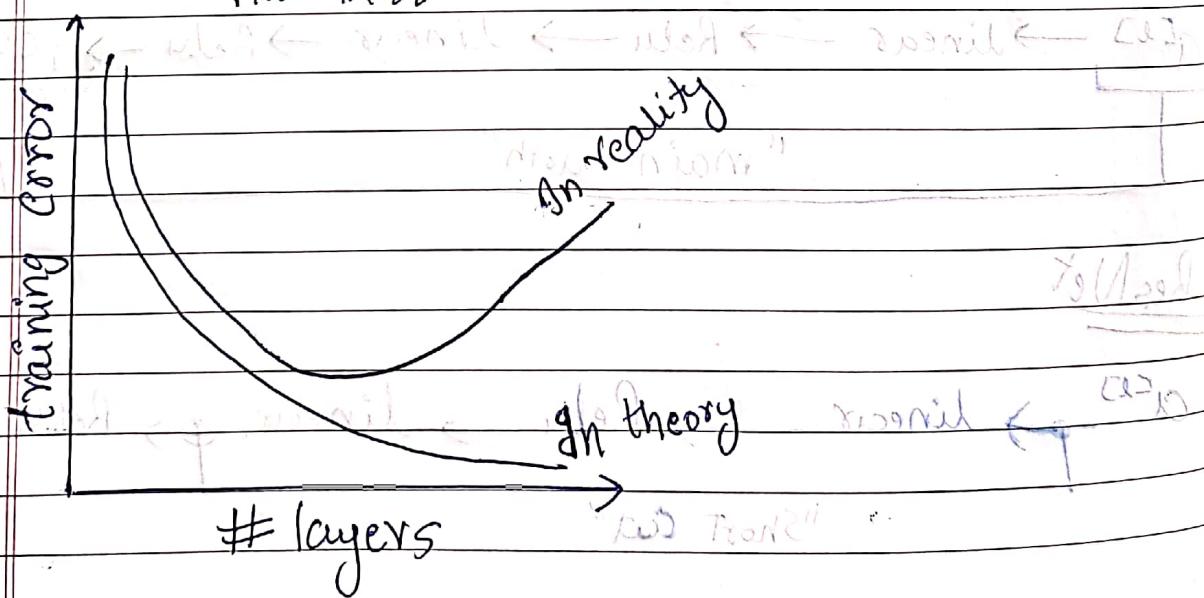
This is done in order to go deeper into the networks

## Residual Networks.



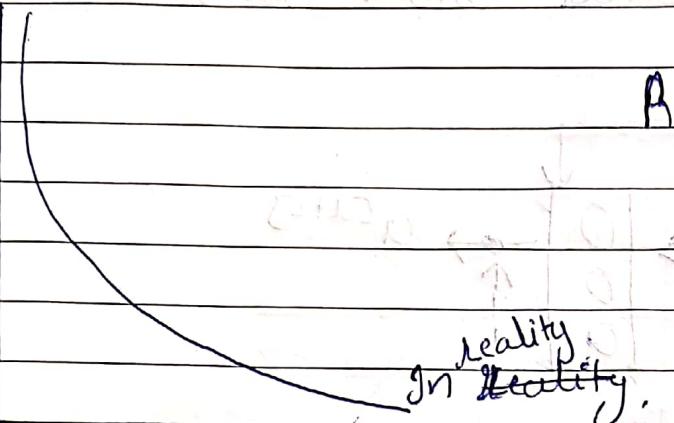
"Plain": without skipping layers.

"Plain"



# ResNet

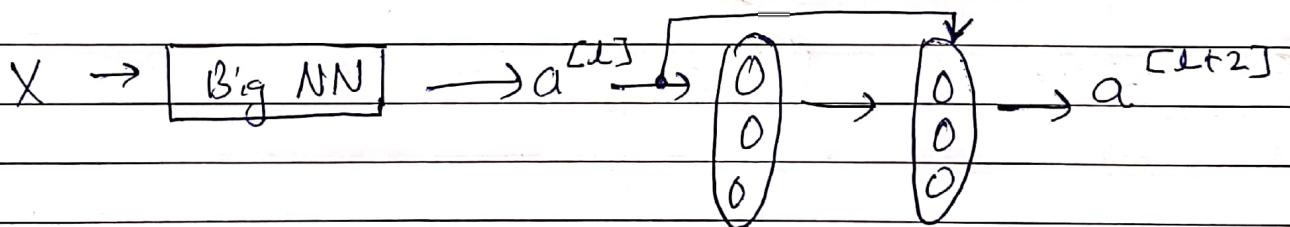
Training error



ResNet allows ~~it~~ to go deeper as it uses "short cut" or "skip connection" method without ↑ training error.

## Q Why ResNets Work So Well?

$$X \rightarrow \boxed{\text{Big NN}} \rightarrow a^{[l+2]}$$



ReLU  $a > 0$

$$a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

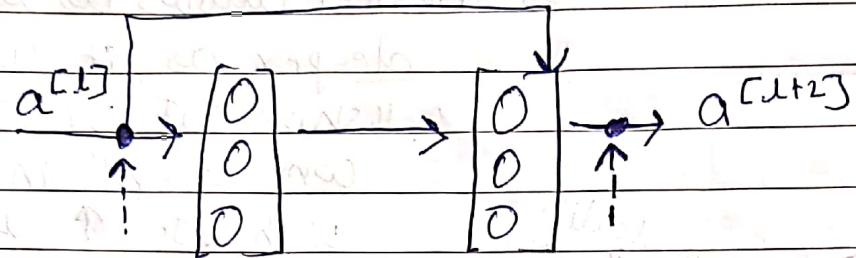
$$= g(w^{[l+2]} a^{[l+1]} + b^{[l+2]} + a^{[l]})$$

if  $w^{[l+2]} = 0, b^{[l+2]} = 0$

then,  $g = (a^{[l]})$

$$a^{[l+2]} = a^{[l]}$$

It shows that, identity function is easy for Residual block to learn.



"Same" convolution means dimension is same at input point and output point.

In order to make "Same" Convolution, multiple \$w\_s\$ to \$a^{[l]}

Example if  $a^{[l+2]} \rightarrow \mathbb{R}^{256 \times 256}$

and  $a^{[l]} \rightarrow \mathbb{R}^{128 \times 128}$

then do,  $w_s a^{[l]}$  where  $w_s \rightarrow \mathbb{R}^{256 \times 128}$

$$(W_{1,1} + W_{1,2}) \cdot p = 1515$$

$$(1515 + 1515) \cdot 1 = 3030$$

$$3030 \approx 3030$$

$$(1515) = 0$$

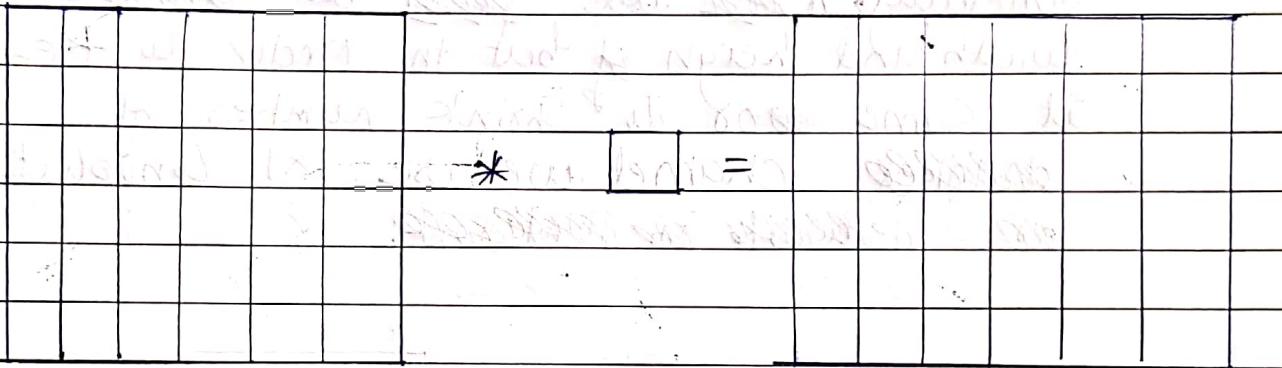
$$1515 \approx 1515$$

# Network in Network

## 1x1 Convolutions

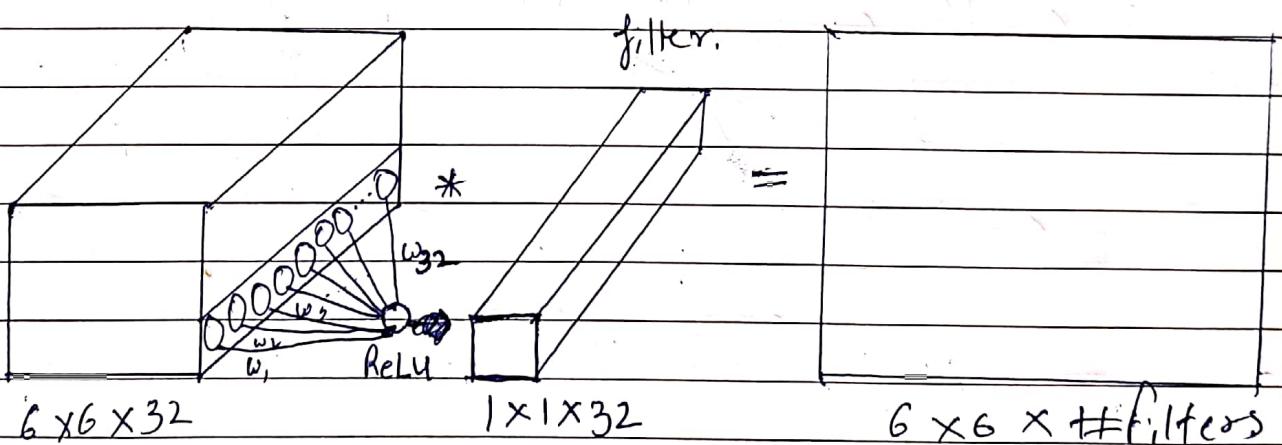
Why does a 1x1 convolution do?

It's like a fully connected layer with shared weights.



$6 \times 6 \times 1$

$6 \times 6 \times \# \text{Filters}$



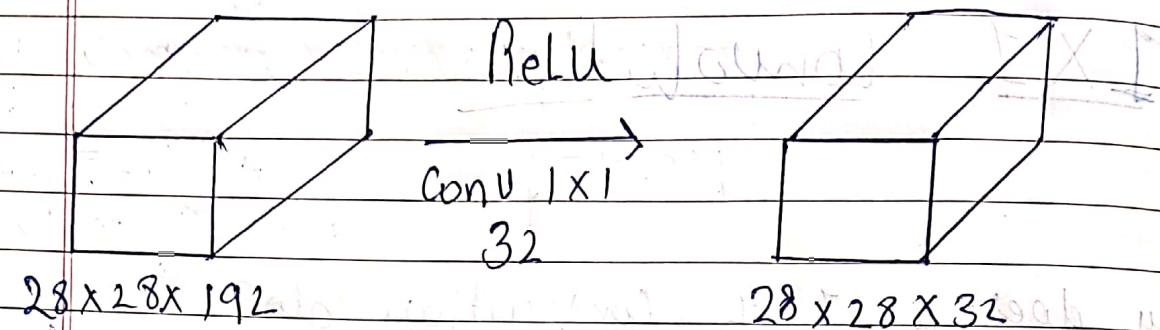
here,  $w_1$  = first channel of filter ( $1 \times 1 \times 1$ )

$w_2$  = Second channel of filter ( $1 \times 1 \times 1$ )

$w_3$  = Third channel of filter ( $1 \times 1 \times 1$ )

$w_{32}$  = 32<sup>th</sup> channel of filter.

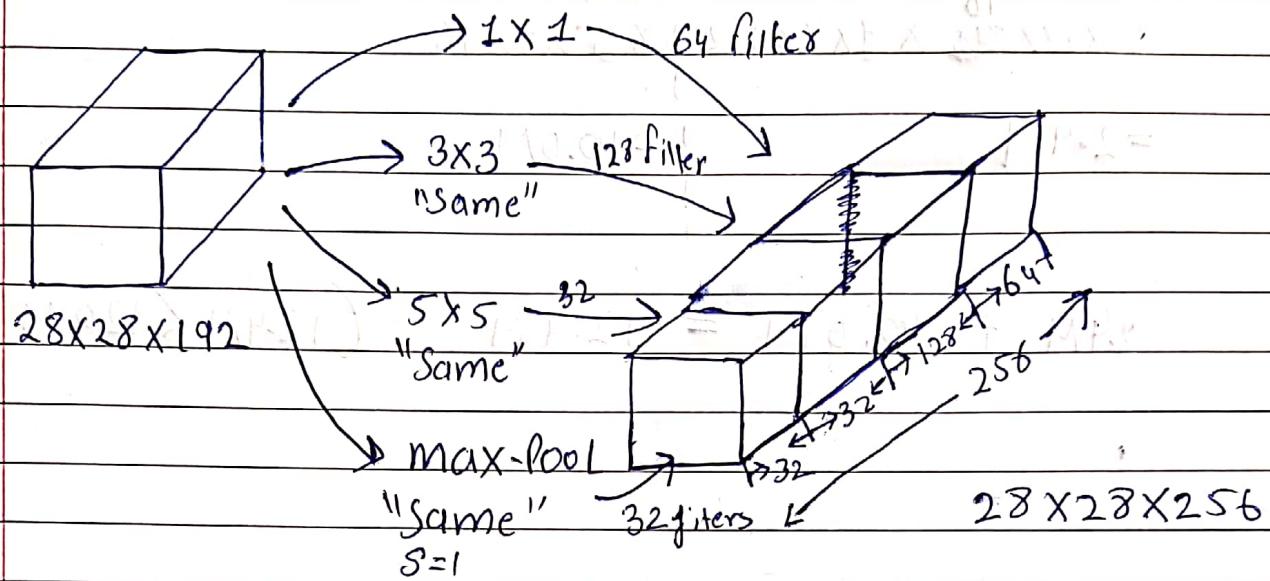
## Example Using $1 \times 1$ Convolutions



conv or Pool layer

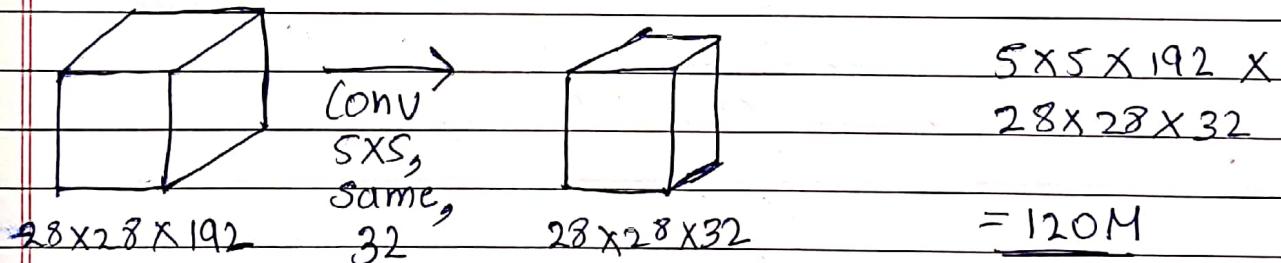
Generally, we ~~use~~ to shrink width and height of but in order to keep it same ~~as~~ its number of ~~channel~~ channel we use  $1 \times 1$  convolutions ~~to~~ ~~reduce~~ ~~in~~ ~~problem~~.

## Inception Network motivation.

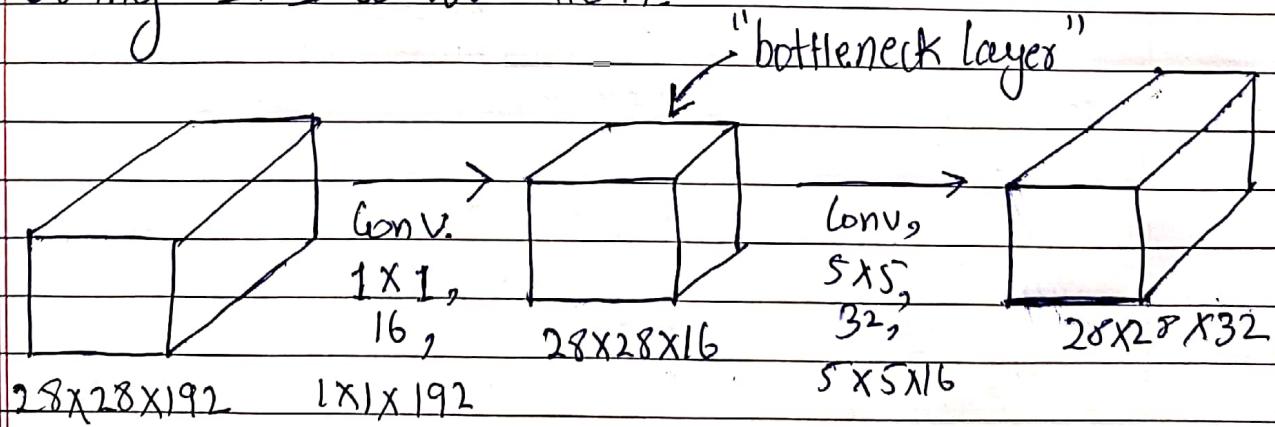


⇒ The problem of Computational Cost.

As it is  $32$  filter



Using  $1 \times 1$  convolution.



Let's look at computational cost here.

$$28 \times 28 \times \frac{16}{16} \times 1 \times 1 \times 192 \times 5 \times 5 \times 16 \\ = 2.4M$$

$$1 \times 28 \times 28 \times 32$$

$$1 \times 10.0M$$

$$2.4M + 10.0M = 12.4M \quad (12.4M < 120M)$$

~~28x28x8x8x2~~

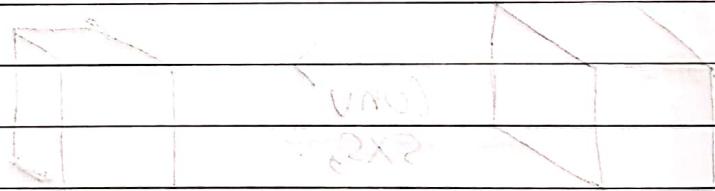
~~192x5x5x16~~

~~1=8~~

softmax with softmax

softmax with softmax

softmax with softmax



Host =

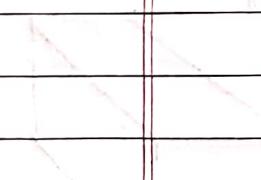
~~(28x8x8x2~~

~~softmax~~

~~softmax~~

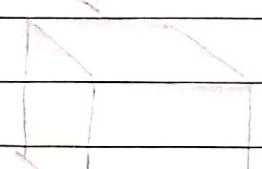
~~Computational cost~~

~~Computational cost~~



~~softmax~~

~~softmax~~



~~softmax~~

~~softmax~~

~~softmax~~

~~softmax~~

~~softmax~~

~~softmax~~

~~softmax~~

~~softmax~~