# Convolutional Neural Nets

* Convolution Neural Nets will be useful in Visual tasks (object Recognition)

  Ex- given a Image and asking is there any Car?., tiger? etc.,

* Some neurons in the Visual cortex That fire when presented lines at Specific angle.

  $V_1$: primary Visual Cortex → edges (which detects)

lly we have $V_2$, $V_3$.

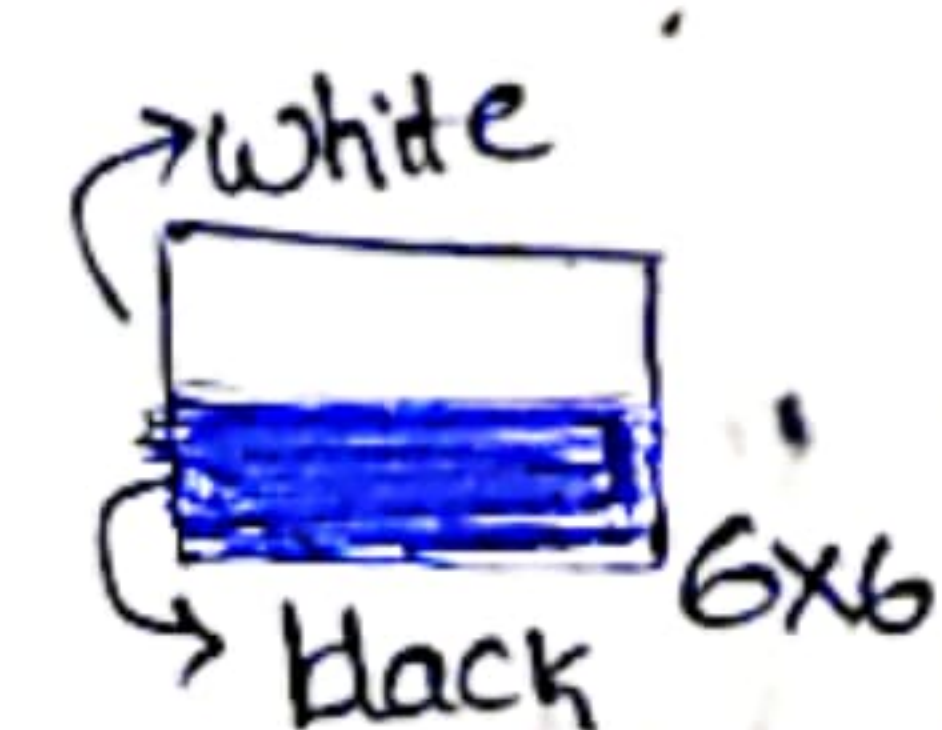## 60.2 Convolution : Edge Detection on Images:-

* Edge detection is $V_1$ (what primary Visual Cortex) will do

① lets assume we have an Image (Consider a grayscale Image (0 to 255) or (0 to 1)

So we do as follows:

① we do Convolution between input Image and Sobel edge detector / Kernel / filter / mask / Operator.

→ white

→ black 6×6

2D

| 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |
| 255 | 255 | 255 | 255 | 255 | 255 |
| 255 | 255 | 255 | 255 | 255 | 255 |
| 255 | 255 | 255 | 255 | 255 | 255 |

6×6

*

2D

| +1 | +2 | +1 |
|----|----|----|
| 0 | 0 | 0 |
| -1 | -2 | -1 |

3×3

So what Convolution will do is

 → Horz edge

| 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |
| 255 | 255 | 255 | 255 | 255 | 255 |
| 255 | 255 | 255 | 255 | 255 | 255 |
| 255 | 255 | 255 | 255 | 255 | 255 |

6X6

nxn

*

| +1 | +2 | +1 |
|----|----|----|
| 0  | 0  | 0  |
| -1 | -2 | -1 |

3X3

KxK

=

| 0 | 0 | 0 | 0 |
|---|---|---|---|
| -1020 | -1020 | -1020 | -1020 |
| +1020 | +1020 | +1020 | +1020 |
| 0 | 0 | 0 | 0 |

(n-k+1 x n-k+1)

| 0(+1) | 0(+2) | 0(+1) |
|-------|-------|-------|
| 0(0)  | 0(0)  | 0(0)  |
| 0(-1) | 0(-2) | 0(-1) |

0+0+0+ 0+0+0+
0+0+0 = 0

| 0(+1) | 0(+2) | 0(+1) |
|-------|-------|-------|
| 0(0)  | 0(0)  | 0(0)  |
| 255(-1) | 255(-2) | 255(-1) |

0+0+0+ 0+0+0
+(-255)+ (-510)
+ (-255)
= -1020

| 255(+1) | 255(+2) | 255(+1) |
|---------|---------|---------|
| 0(255)  | 0(255)  | 0(255)  |
| -1(255) | -2(255) | -1(255) |

255 + 510 + 255 + 0 + 0 + 0 +
-255 + (-510) + (-255)

= 0

$$\Box_{n\times n} * \Box_{K\times K} = \Box_{(n-k+1, n-k+1)}$$

so we got

 } Dark Shaded

Light Region {

by if we do for Vertical edge detection

Light

we get



 → Vertical edge

**60.3 Convolution : padding and Strides :-**

Convolution $6\times6 \rightarrow$ ...

Ex:- $6\times6 \rightarrow 8\times8$

**Padding**

Ex:-

Horizontal edge 1 →

Horz. edge 2 : →

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 255 | 255 | 255 | 255 | 255 | 255 | 0 |
| 0 | 255 | 255 | 255 | 255 | 255 | 255 | 0 |
| 0 | 255 | 255 | 255 | 255 | 255 | 255 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

$*$

| +1 | +2 | +1 |
|---|---|---|
| 0 | 0 | 0 |
| -1 | -2 | -1 |

$k \times k$

$\downarrow$

$P = 1$

$6\times6 \longrightarrow 8\times8$

and we done here

Zero padding.

$$\Rightarrow \quad n\times n \xrightarrow{\quad P=1 \quad} (n+2)\times(n+2) \Longrightarrow (n+2\times p) \times (n+(2\times p))$$

$$n\times n \xrightarrow[\text{padd } P]{\quad k\times k \quad} (n-k+2p+1) \times (n-k+2p+1)$$

By padding we observe some interesting effects on the

Operations.

$*$ By padding we can get another edge

○ Another type of padding is Same value padding.

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 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 | 250 | 250 | 250 | 250 | 250 | 250 | 250 |

But when it comes to here it get confuse

## Strides

| 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |
| 255 | 255 | 255 | 255 | 255 | 255 |
| 255 | 255 | 255 | 255 | 255 | 255 |
| 255 | 255 | 255 | 255 | 255 | 255 |

shifted by 1 = Stride of 1

Shifted by 2 = Stride of 2

$$n \times n \xrightarrow[K=K]{S=S} \left(\frac{n-k}{s}+1\right) \times \left(\frac{n-k}{s}+1\right)$$

$$\text{Ex: } 6 \times 6 \xrightarrow[K=3]{S=2} \left(\left\lfloor\frac{3}{2}\right\rfloor+1\right) \times \left(\left\lfloor\frac{3}{2}\right\rfloor+1\right)$$

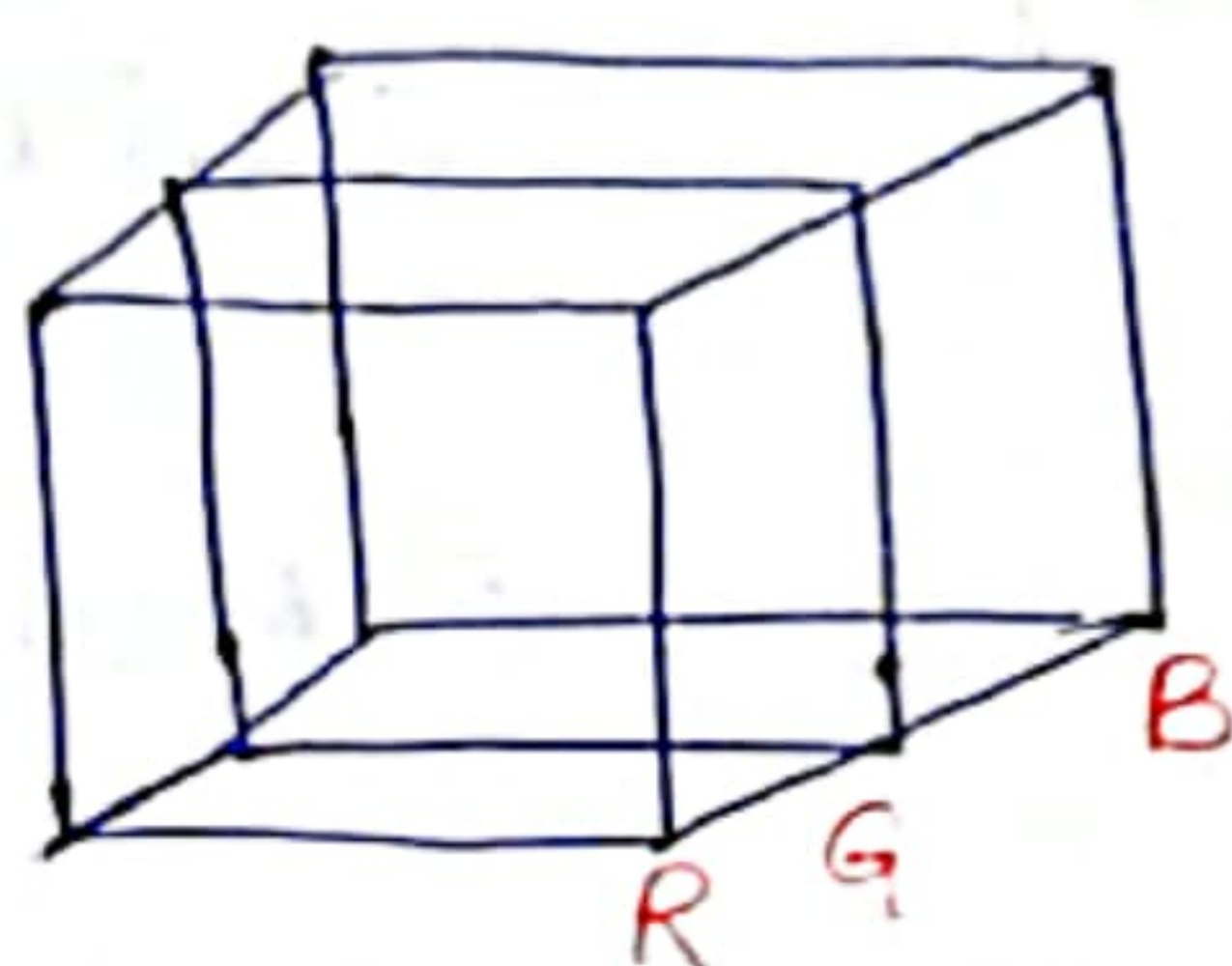$$6 \times 6 \xrightarrow[K=3]{S=2} \quad 2 \times 2 \quad (O/P)$$

$$\ast \quad n \times n \xrightarrow[\text{Padding } P]{K \times K} (n-k+2P+1) \times (n-k+2p+1)$$

$$n \times n \xrightarrow[P, S]{K \times K} \left(\left\lfloor\frac{n-k+2P}{S}\right\rfloor+1\right) \times \left(\left\lfloor\frac{n-k+2P}{S}\right\rfloor+1\right)$$

## 604 Convolution Over RGB Images :-

As we know each Image/pixel would be Represented using Three primary colours Red, Green, Blue.

→ To Represent an image R, G, B's will be layered as a matrix to form as a Image
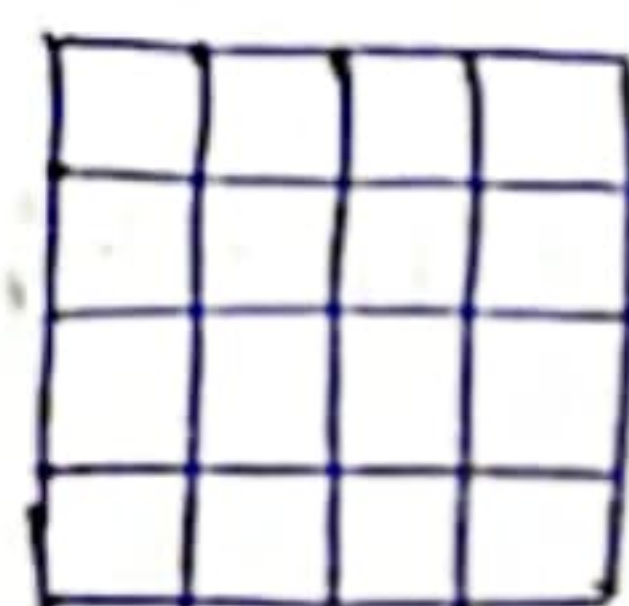


[R, G, B]

each R, G, B will be a Channel.
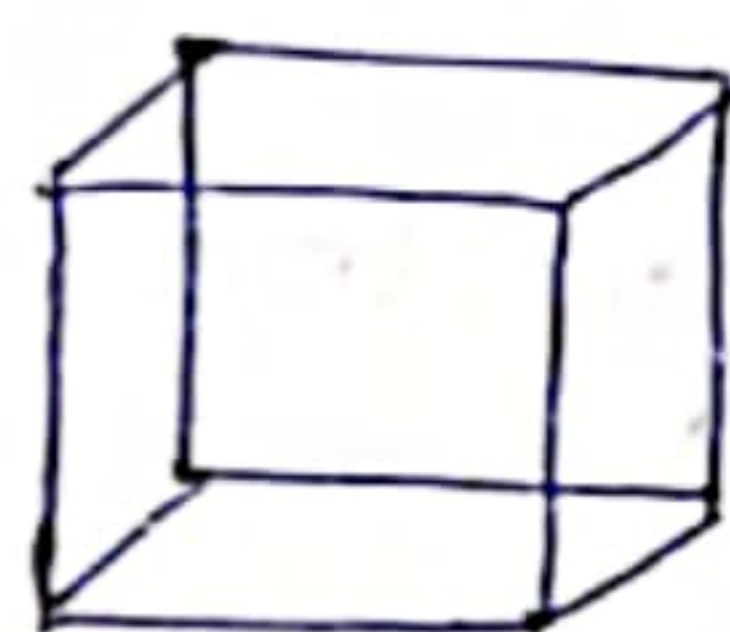
→ each dimension would be Represented as as a Tensor
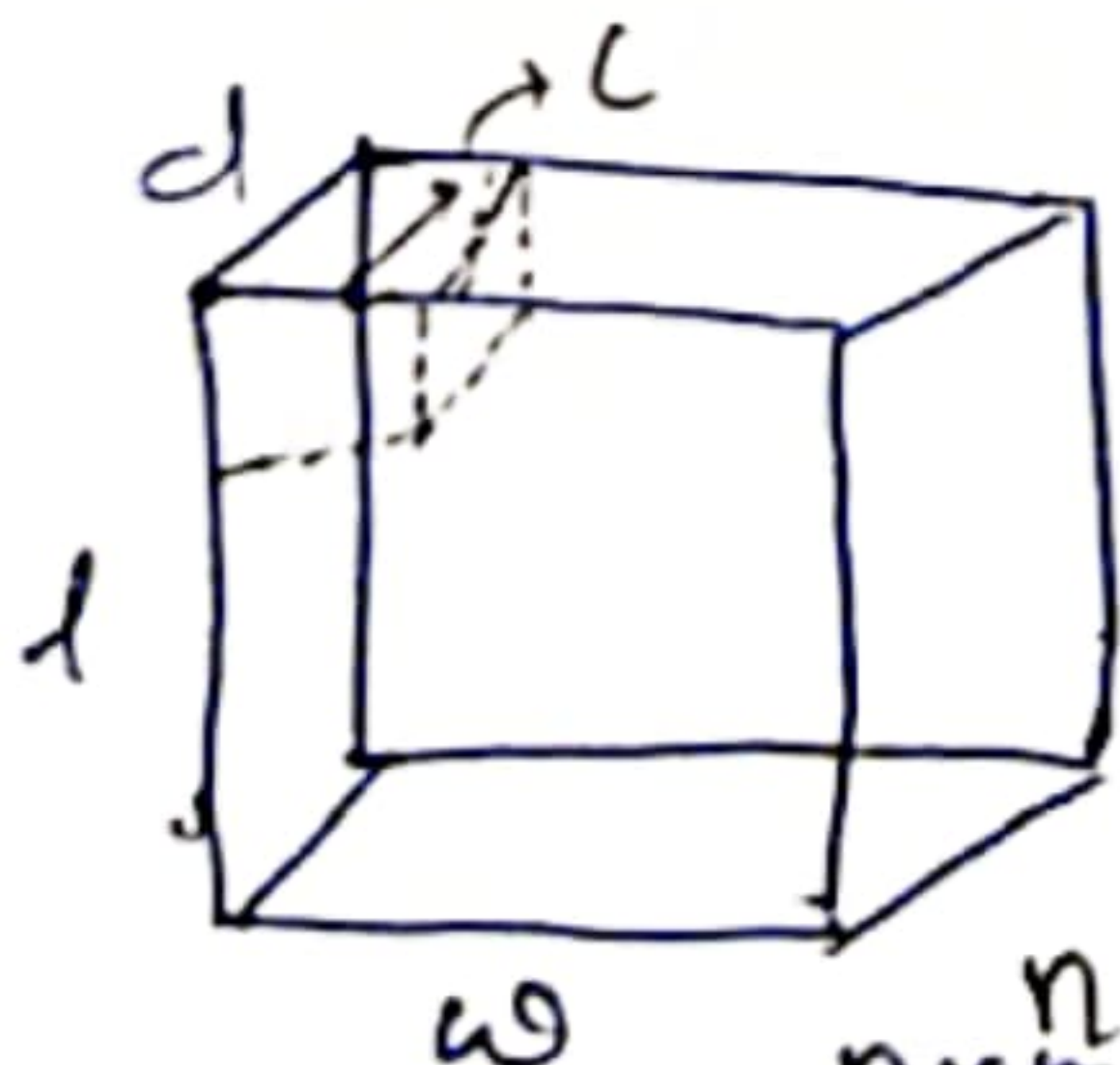


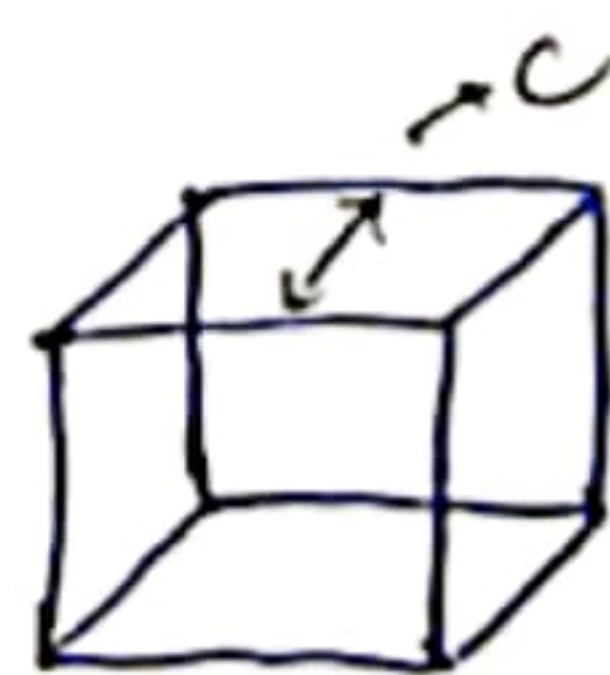1D-Tensor    2D-Tensor    3D-Tensor.

each 3D tensor ⟶ $n \times m \times c$ Channels

# as we seen Convolution in 2D matrix → Component wise mul & add.

to when we do convolution over 3D tensors. (Colour Images)



$d$   $c$

conv *

$c$

$K \times K \times c$

Thex $c's$ should be same.

⟶ $(n-k+1) \times$
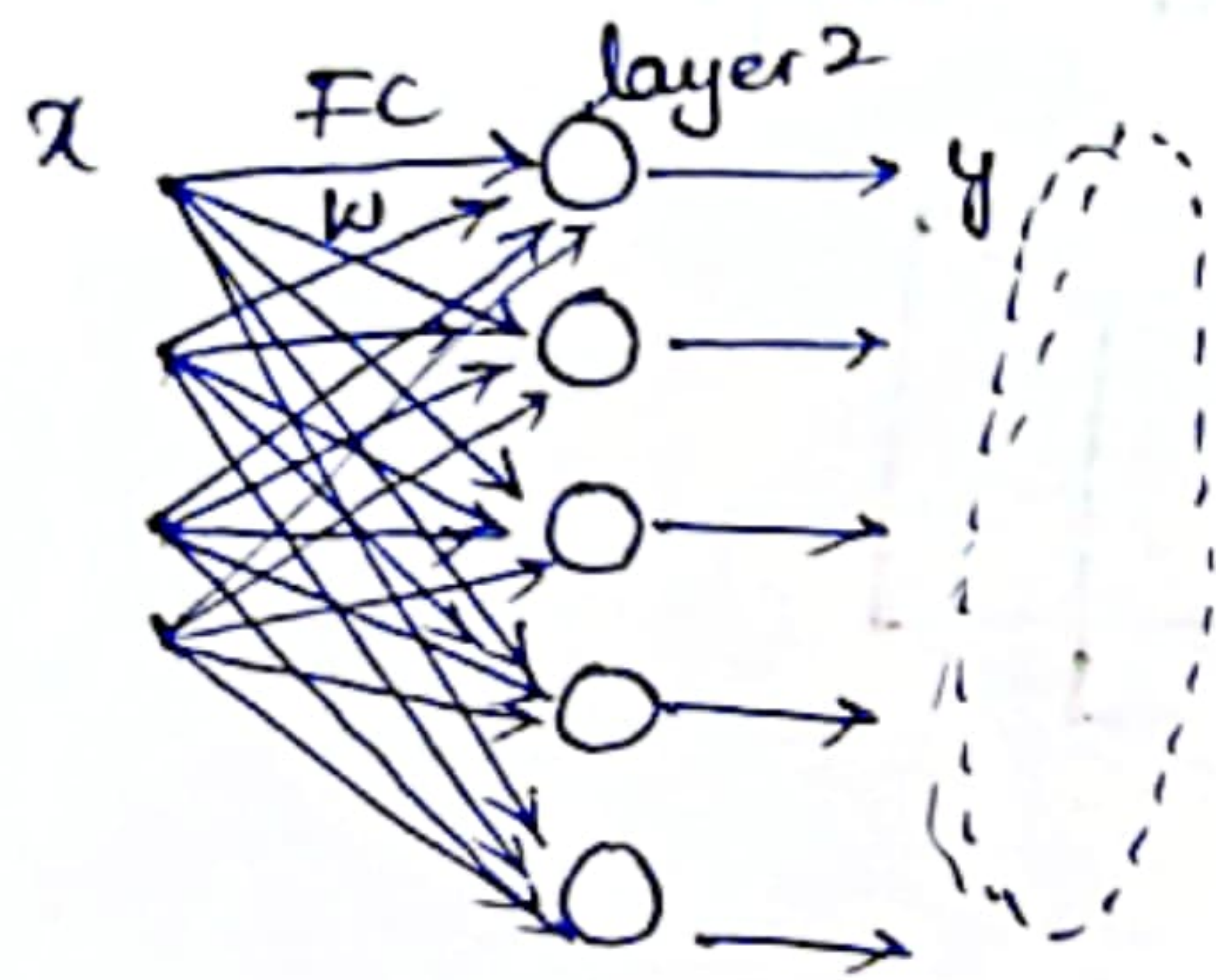$(n-k+1) \times 1 (1 chan$

$c →$ no of chanels
⟶ R, G, B (depth)

$n \times m \times c$
$l \times w \times d$

O/p we get a two dimensional Image

## 60.5 Convolution layer

* As we seen MLP layer. where we have inputs/ previous layer O/p's.



So as we are doing

1. $W^T x = z$ (we are finding $z$)

2. $ReLU(z) = y$ (applying ReLU on top of it).

* Convolution layers are biologically inspired.

as we seen in edge detectors

$$1 \text{ edge detector} \longrightarrow 1 \text{ Kernel}$$

lly for multiple edge detectors $\longrightarrow$ multiple Kernels.

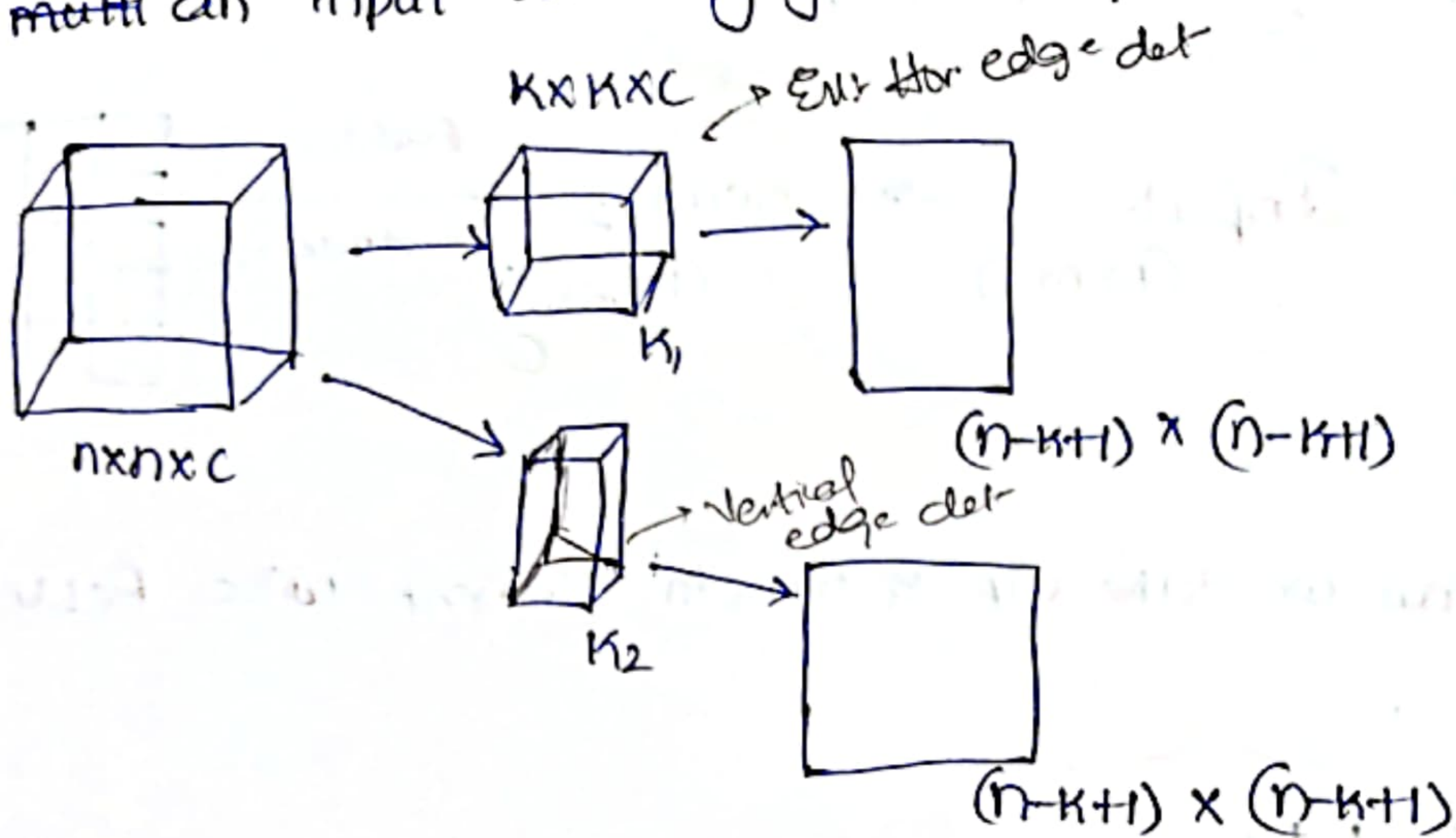* As we are using Sobel in Classical Image processing while in CNN we learn Kernel matrix as of back propagation.

As in MLP ! we learn weights w's $\longrightarrow$ $w \cdot x$.

where in CNN! we learn Kernel matrices $\longrightarrow$ $I_p * W_{K \times K}$    ← Kernel

Convolution $\longrightarrow$ 2D matrix $\longrightarrow$ Component wise mul (add)

Conv $\longrightarrow$ 3D tensor $\longrightarrow$ " " "

* for multi an input we may get multiple kernels

$K \times K \times C \longrightarrow$ Ext Hor edge det

$n \times n \times c$

$K_1$

$(n-k+1) \times (n-k+1)$

$\longrightarrow$ Vertical edge det

$K_2$

$(n-k+1) \times (n-k+1)$

∴ for each convolution layer.

for an input $n \times n \times c$

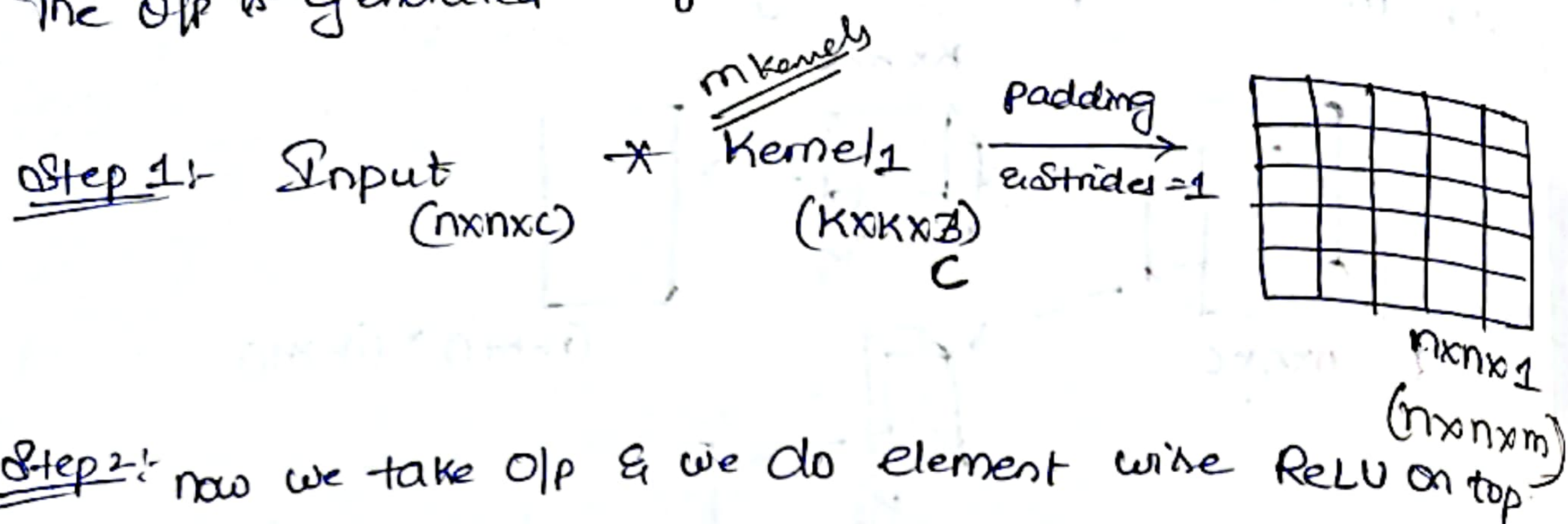we get kernel Ma O/p of $n \times n \times m$ (where m is no of kernels)

$n \times n \times c$
$\xrightarrow[\begin{array}{c} Padding \\ S=1 \end{array}]{\begin{array}{c} (K \times K \times c) \\ 10 \text{ multiple kernels} \end{array}}$

$\omega$  $I_1$  $I_2$  $I_3$

$n \times n \times m$
$\hookrightarrow$ no of kernels

$\downarrow$
no of
tensor

(which is hyperparameter)

for an 3D, RGB Image

what padding, S

$C=3$

are also hyperparameter

to $l, \omega$ (length of O/p, & width of O/p) $\longrightarrow$ Padding, Strides & kernel.
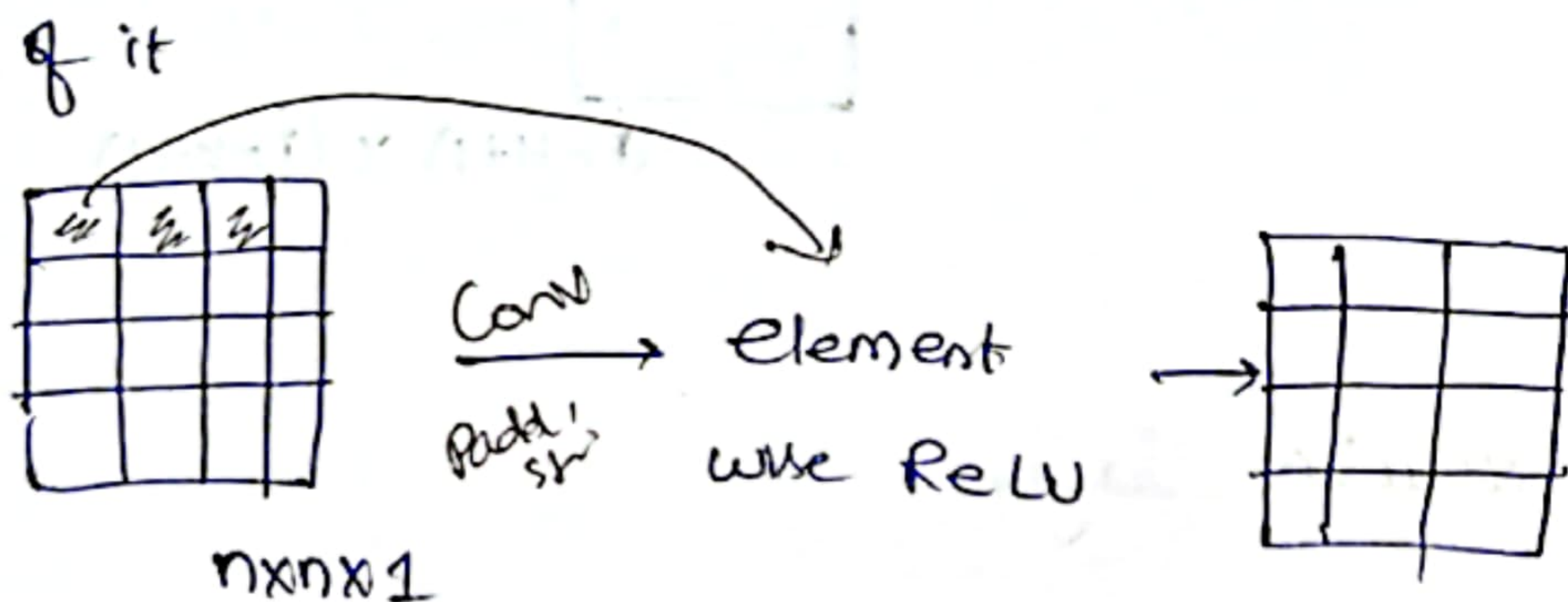
$d \longrightarrow (m) \# of$ Kernels.

$\to$ AO for one layer of Convolution layer we have 3 hyper-parameters $k$, # kernels, padding, & strides.
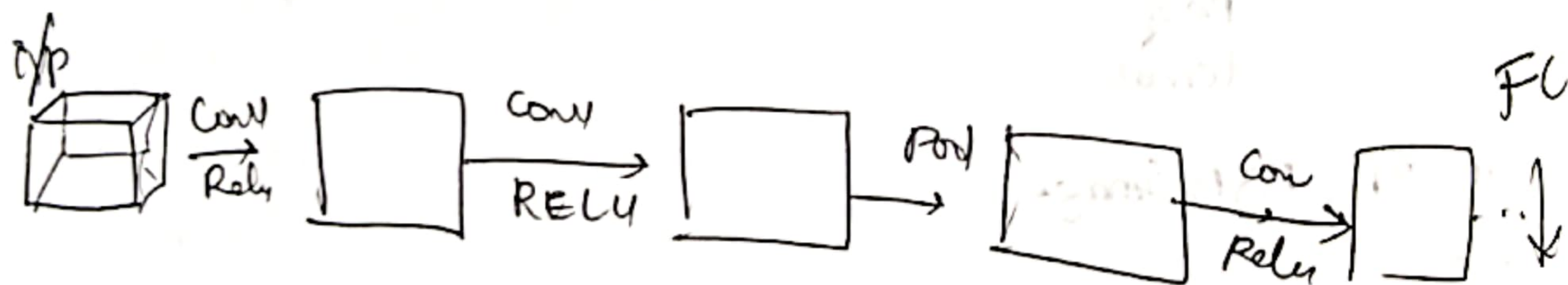
The O/P is generated as follows:

**Step 1:-** Input $\overset{m\ kernels}{*}$ Kernel$_1$ $\xrightarrow[\text{& stride} = 1]{\text{padding}}$
(n×n×c)      (K×K×3)
                   C

n×n×1
(n×n×m)

**Step 2:-** now we take O/P & we do element wise ReLU on top
of it

n×n×1 $\xrightarrow[\text{Padd!}\ \text{str}]{\text{Conv}}$ element wise ReLU $\longrightarrow$

n×n×c $\longleftrightarrow$ (n×n×1)

There size maynot be same it vary depends un padding & striding.

o/p $\xrightarrow[\text{Relu}]{\text{Conv}}$ $\xrightarrow[\text{RELU}]{\text{Conv}}$ $\xrightarrow{\text{Pool}}$ $\xrightarrow[\text{Relu}]{\text{Conv}}$ FC

we do apply multiple layers of conv layers

## 6.6 max pooling :-

* Max pooling is one of the layer which we add in our model.

* Max pooling makes the model location-Invariant, Scale-Invariant, Rotational-Invariant

To understand max pooling

Imagine we have 4x4 Image



$$\xrightarrow[\quad S=2 \quad]{\text{max pooling} \atop K=2}$$

| 6 | 8 |
|---|---|
| 3 | 4 |

Now if we do max pooling with Kernel/filter = 2 & strides = 2

by if we do with $\xrightarrow[\text{max pooling}]{K=2 \: \& \: S=1}$

| 6 | 7 | 8 |
|---|---|---|
| 6 | 6 | 8 |
| 3 | 3 | 4 |

*if we apply another level of max pooling

| 6 | 8 |
|---|---|
| 3 | 4 |

$\xrightarrow[2\times2]{\text{max pooling}}$

| 8 |
|---|

1x1

and this max value (8) will be the max value among all the values of convolution matrix input.

This is why we can say that it is location Invariance.

so we can easily detect image if we even irrespective of its location as the image seeks max activation value.

max-pooling is extremely used in convolers.

→ we Can also do Avg/mean pooling where Rather than max value it Considers Avg/mean value.

## 60.7 CNN Training : Optimization

→ As we seen in MLP :

when eve we are applying / using backpropagation which is the key idea in mlp for all algorithms like SGD, Adagrad, Adam. etc all are differentiable.

that means

$$x \longrightarrow \hat{y} \rightarrow L(y.\hat{y}) \rightarrow \frac{\partial L}{\partial w} \rightarrow \nabla_w L.$$

lly In Case of

Conv layer ⟶ we do Convolution ; ReLU

⌐ (on matrices)

⟶ which do

┌ - - - - - - - - - - - - - - ┐
Elementwise + addition → diff
Mul to
└ - - - - - - - - - - - - - - ┘
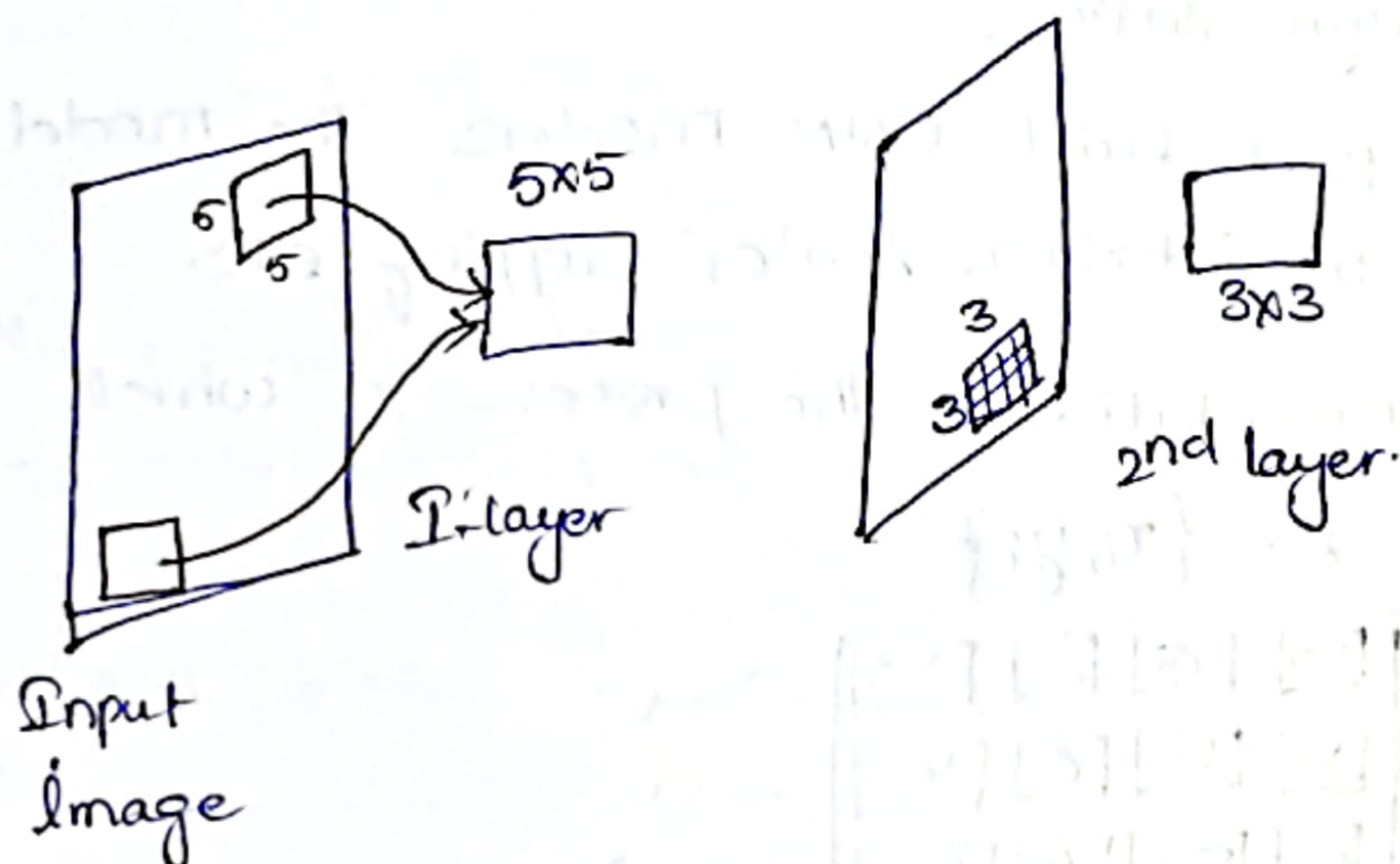
lly for max pooling

we do differentiation as given by

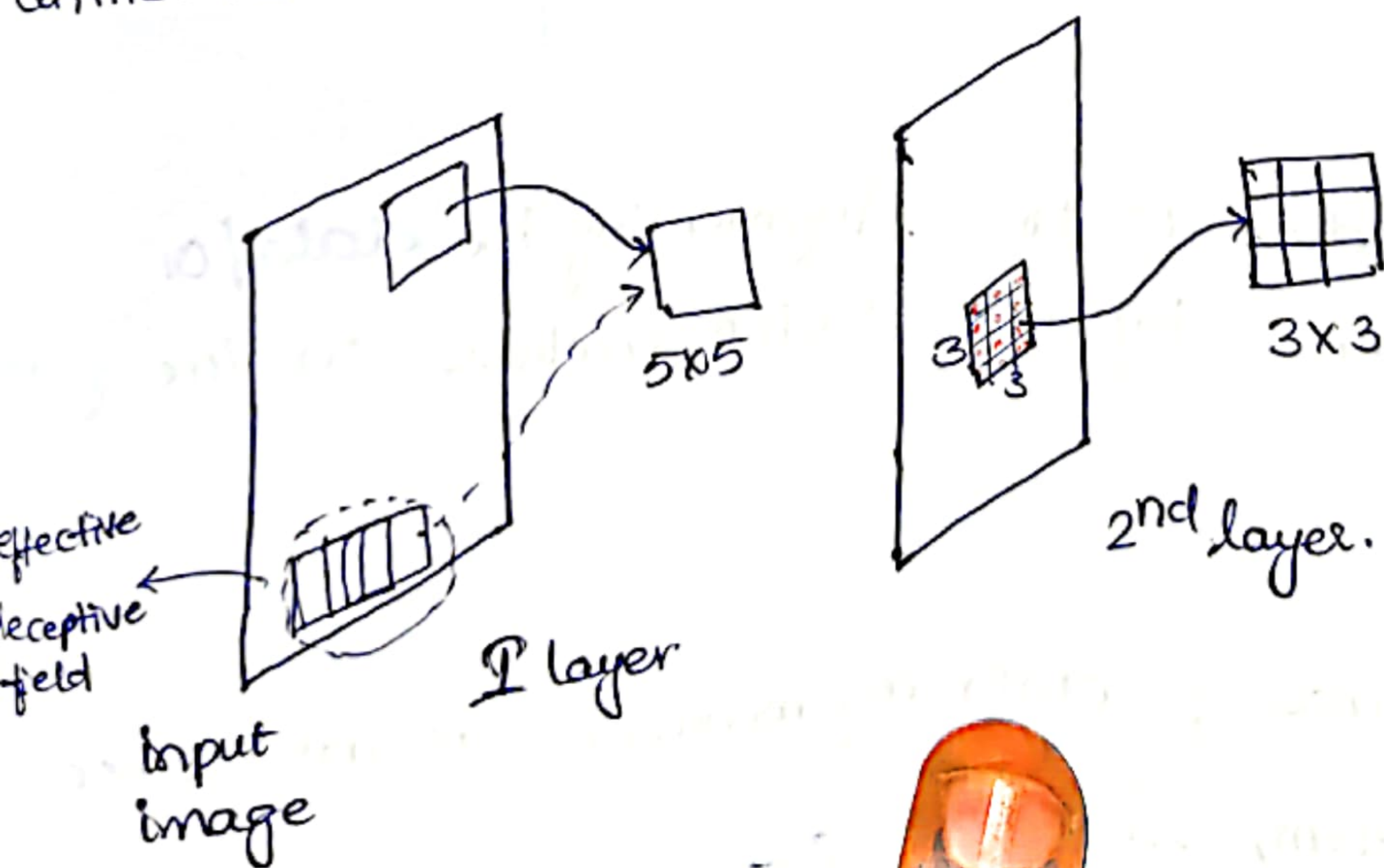max value $\xrightarrow{\text{diff}}$ 1

non-max value ⟶ 0

60    __Receptive fields and Effective Receptive fields :-__

Receptive field is the Region of Image of which the Convolution of Kernel at a given time t. is been convoluted with



Input Image

where as effective Receptive field would be when we have 2, 3 --- layers (multi layer CNN) the effective Region of the Original Image, whose pixels directly or indirectly Contributed to get O/p kernel

## 60.9 Example CNN: LeNet [1998]:-
← LeNet printout →

## 60.10 ImageNet Dataset
← Wikipedia →

## 60.11 Data Augmentation:-

while we apply a Build CNN models The model needs to be Robust to Rotation, Scale, Cropping etc.

so Data Augmentation is The process in which for a given Data $D = \{x_i, y_i\}$



So tries to do all possible Operations like flip, hor-shift, Vertical shift, Rotation, Zoom, Shear etc., Such that the model become Robust and easy to detect the image.

→ So the Core idea is to Augmenting The data/or adding the data by doing Such Operations on the given image.

→ so the main use of data augmentation is invariance to rotation, Zoom, Shift etc.,

and we Can also transform a small datasets → larger datasets.