

7th March, 2023

1/9

Introduction to Computer Vision: CNNs

DSML : Computer Vision

Class starts
@ 9:05 pm.



**What normal people see
when they walk on street**



**What Computer Vision
folks see**



WHO WOULD WIN?



**STATE OF THE ART
NEURAL NETWORK**



ONE NOISY BOI

About me :

Name : Dhruv. Jawali.

Qualifications : * B.Tech (CSA) from NIT Goa.
* PhD, Indian Institute of Science .
(learning Filters, Filterbanks, Wavelets
and Multiscale Representations.)

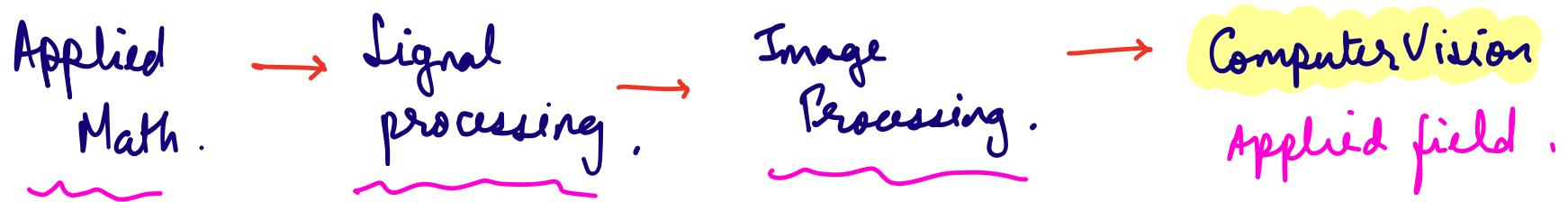
Work Experience : * Software engineer, Samsung Research.
* Deep Learning Researcher, Neurospin.ai.
* ML Engineer and Instructor @ Scales.

Interests : ❤
→ Coding , C.V. Research (3D)
→ Problem Solving
→ Teaching.

Agenda for today:

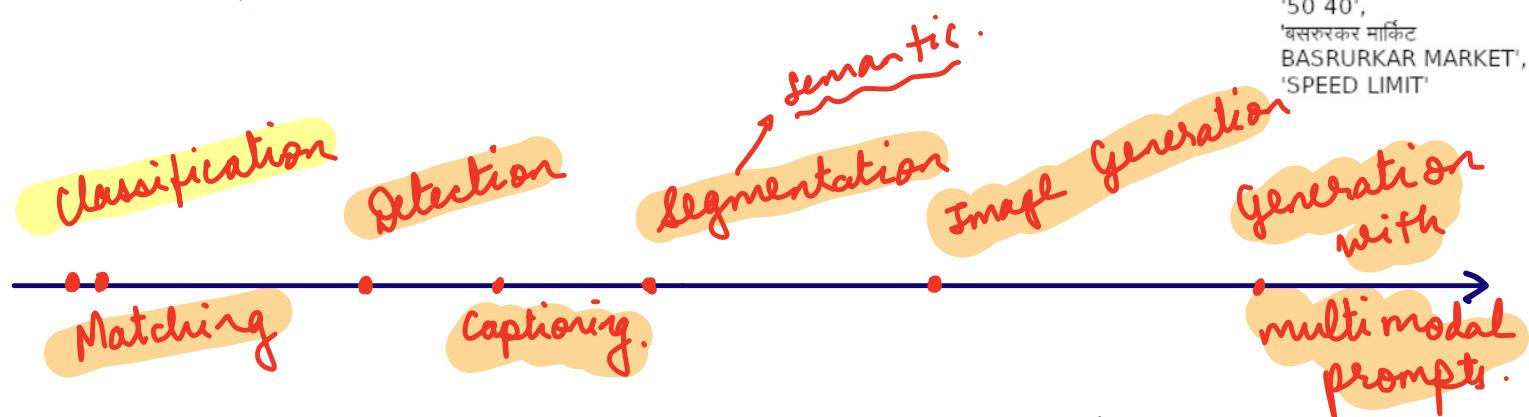
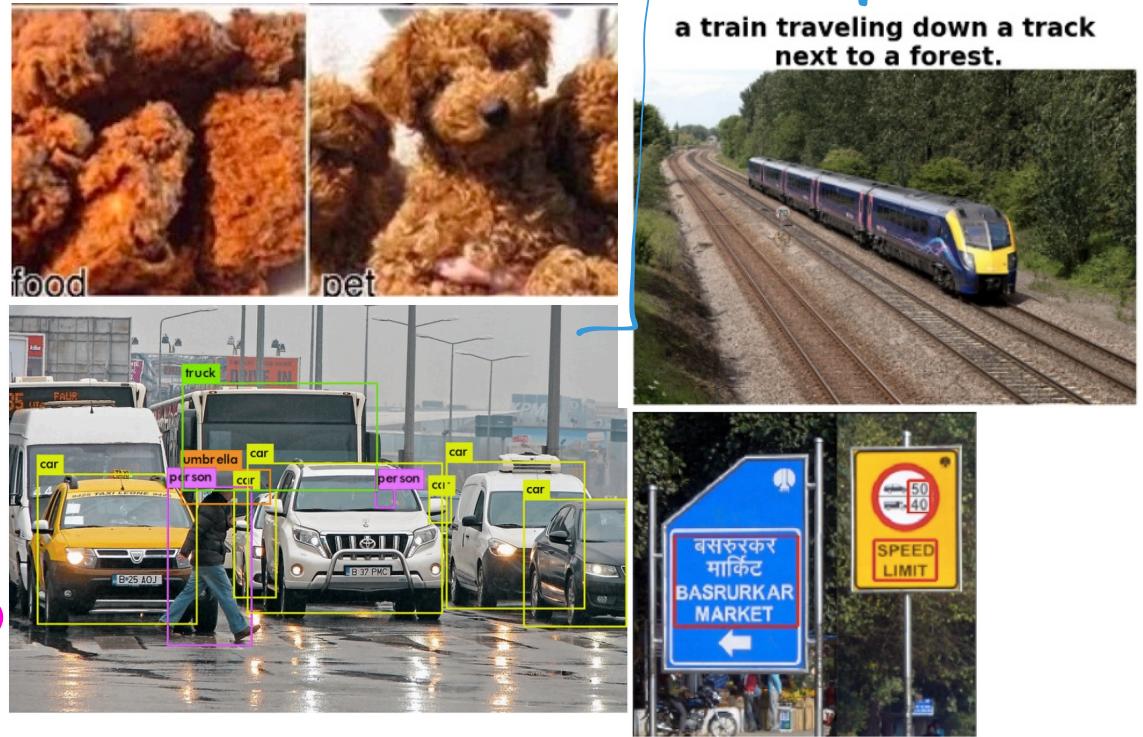
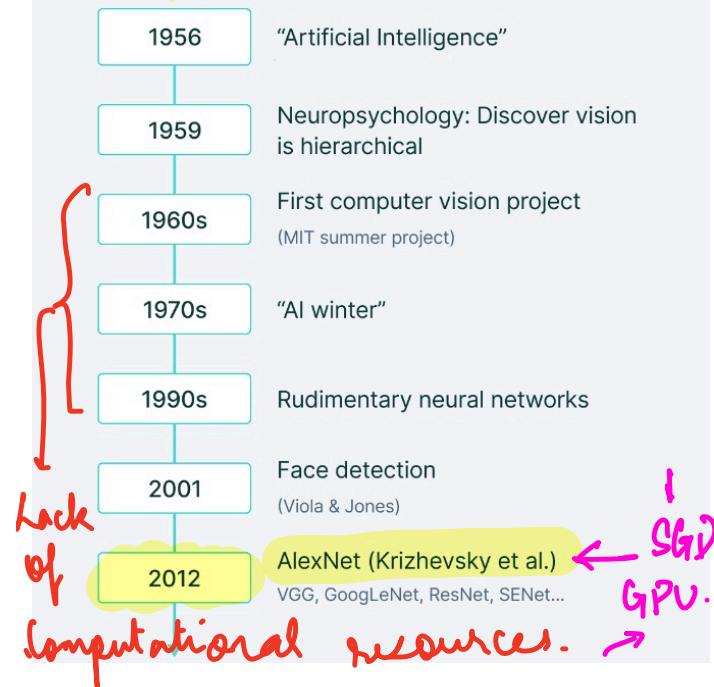
- * Introduction to the Field.
- * Inspiration for CNNs : Neuroscience!
- * Traditional Neural Networks vs. CNNs.
- * Deep Dive:
 - Convolution.
 - Pooling
 - CNN architectures.

Computer Vision as a field of study



Typical problems in Computer Vision

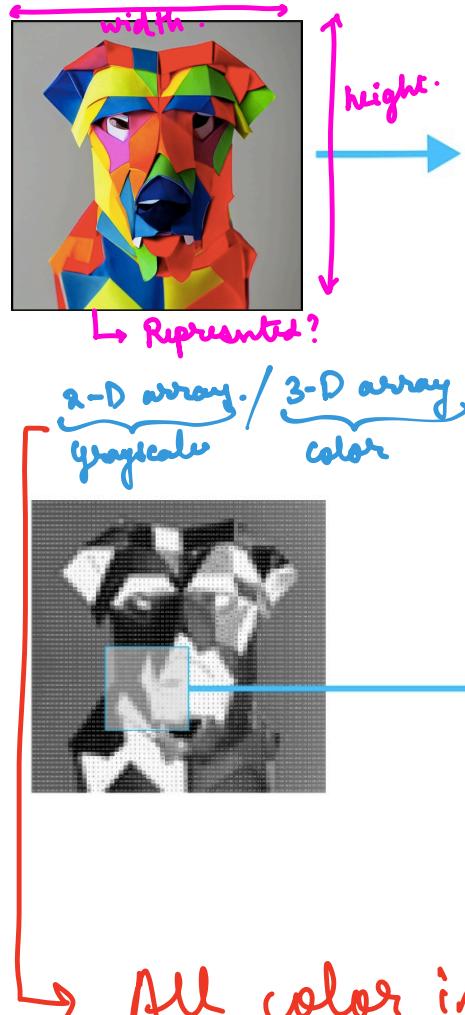
A brief history of Computer Vision



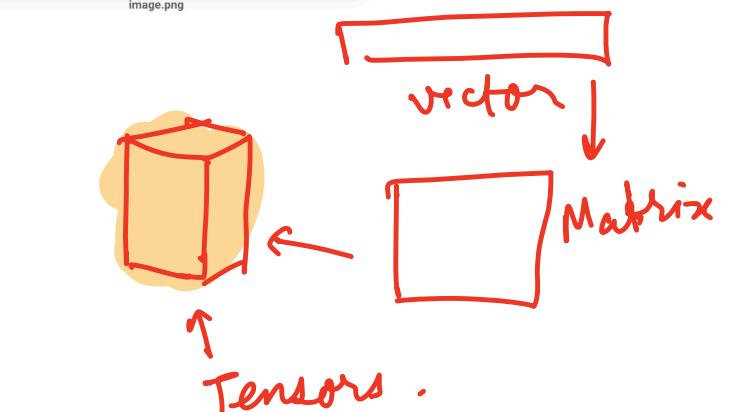
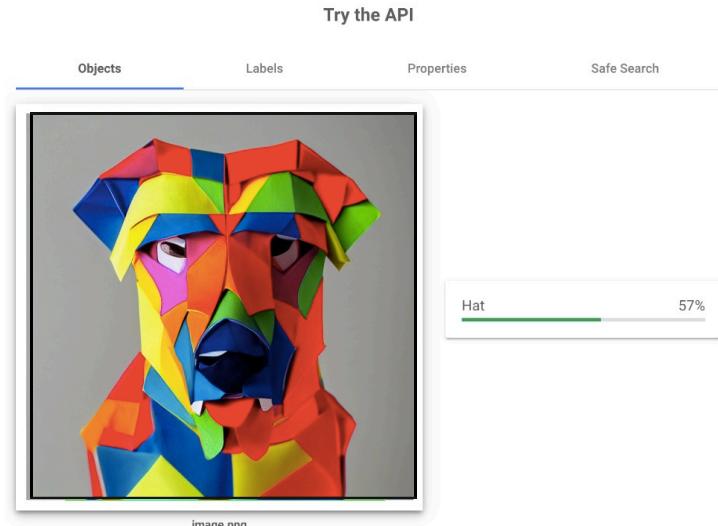
Increasing order of Difficulty

How does a computer "see"?

0-255 → Values in pixels.



158	176	246	246	251	241	235	242	254	249	244	253	248	255	127	0
159	172	243	247	249	239	240	251	255	185	220	255	249	244	27	4
160	168	239	248	247	250	253	252	246	109	247	250	255	160	4	28
161	164	237	248	248	249	249	255	199	15	234	255	254	97	27	3
162	163	235	250	248	249	246	255	122	0	188	255	195	24	0	4
162	162	233	252	249	249	251	250	44	0	139	255	62	0	8	6
163	158	228	254	249	246	255	188	0	0	93	185	0	0	0	0
161	165	236	252	249	246	255	190	0	0	38	68	13	50	78	87
160	224	253	247	249	248	249	251	58	0	12	25	55	86	100	67
207	255	251	249	255	247	247	255	189	0	8	32	0	0	0	0
255	251	255	145	144	255	244	248	253	58	0	7	12	12	9	5
255	248	251	46	0	192	255	241	255	112	0	3	1	3	3	3
248	255	205	3	0	22	229	250	255	167	0	8	1	4	3	2
243	255	154	0	12	0	66	251	253	209	5	12	10	5	5	3
245	255	182	16	0	7	0	116	255	232	30	0	3	5	1	5
250	252	227	155	25	7	2	0	169	255	57	8	34	4	1	4

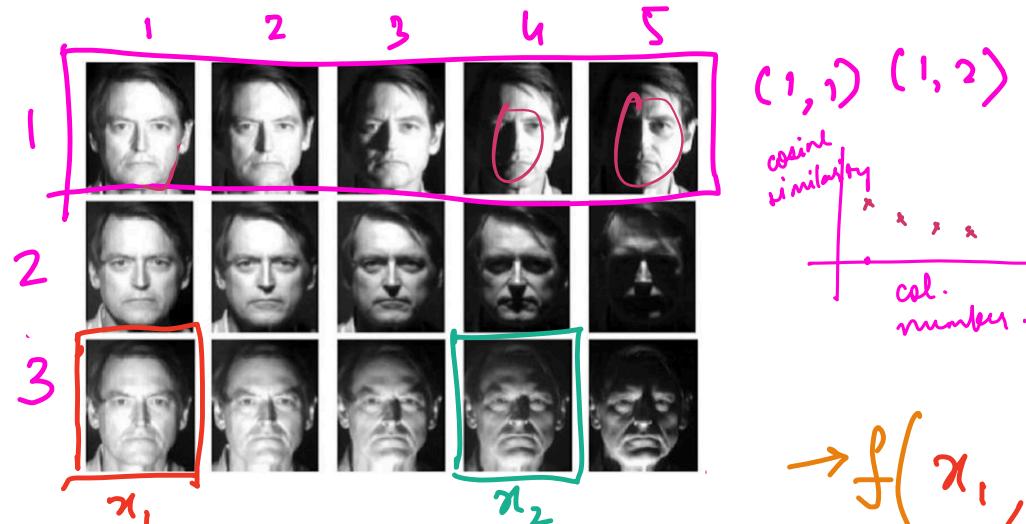


→ All color images are 3-D arrays.

(height , width , channels).

$$\begin{array}{ccc} - & - & - \end{array} \rightarrow (256)^3$$

What makes a Computer vision task difficult?



$$\rightarrow f(x_1, x_2)$$

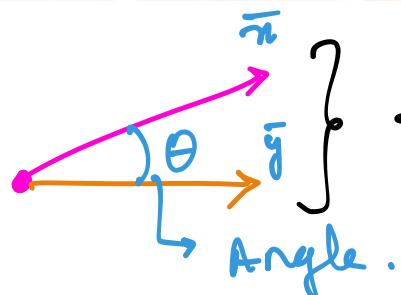
↳ occlusions (loss from 3D \rightarrow 2D).



"How to get a similarity measure so that all

→ Different "deformation" images are pointing to the same obj?"

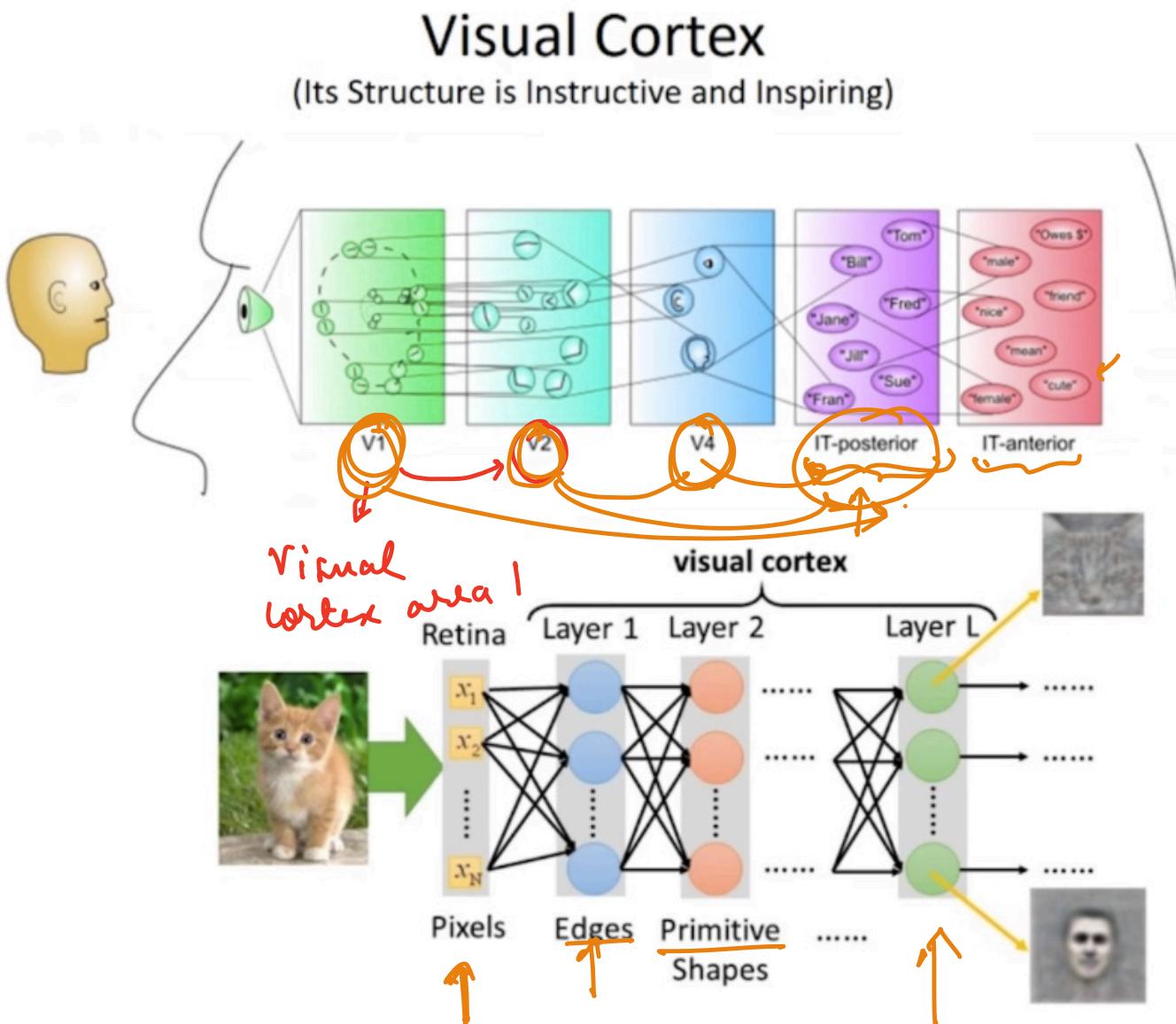
Math:



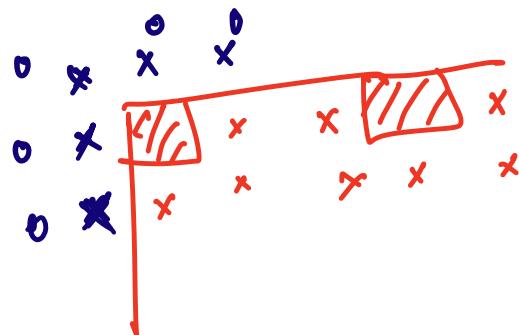
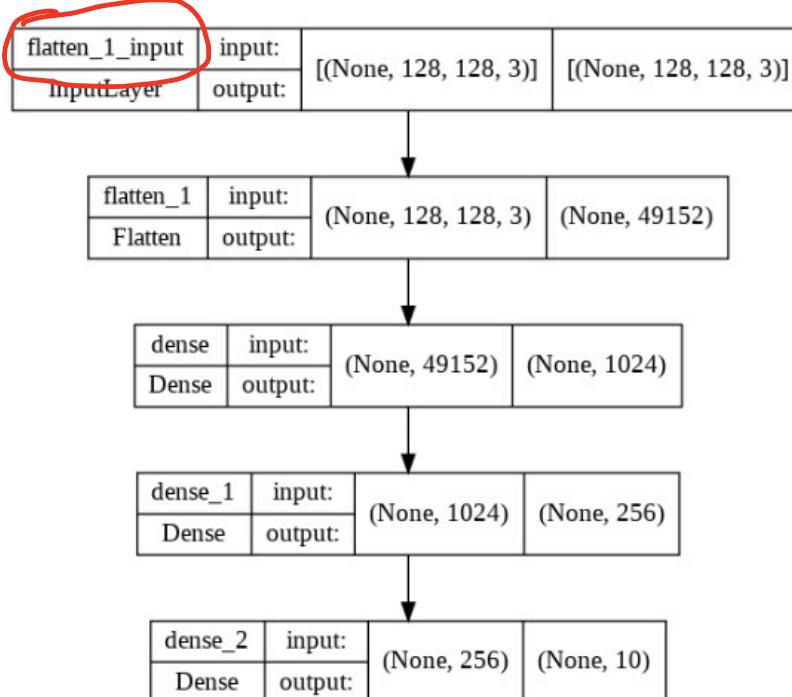
similarity measures:

$$\cos(\theta) = \frac{\bar{x}^T \bar{y}}{\|\bar{x}\| \cdot \|\bar{y}\|}$$

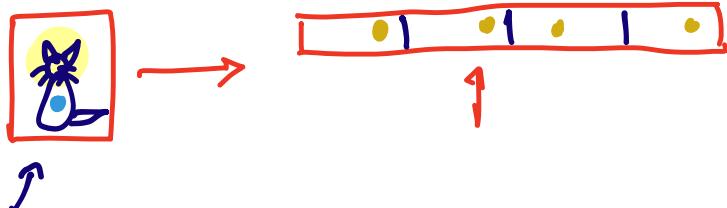
Inspiration : The visual cortex :



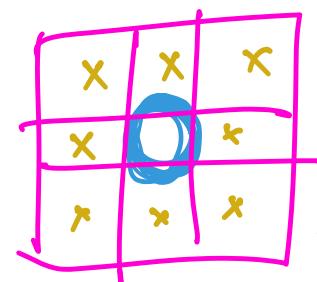
Deep dive: Using Neural Networks for Classification.



↖



I am going to lose
positional information
if I flatten the
image.



← neighbourhood
of a pixel .

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

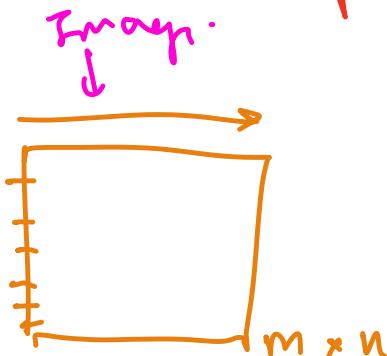
kernel.

$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$
$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$
$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$

$\underbrace{\qquad\qquad\qquad}_{3 \times 3}$

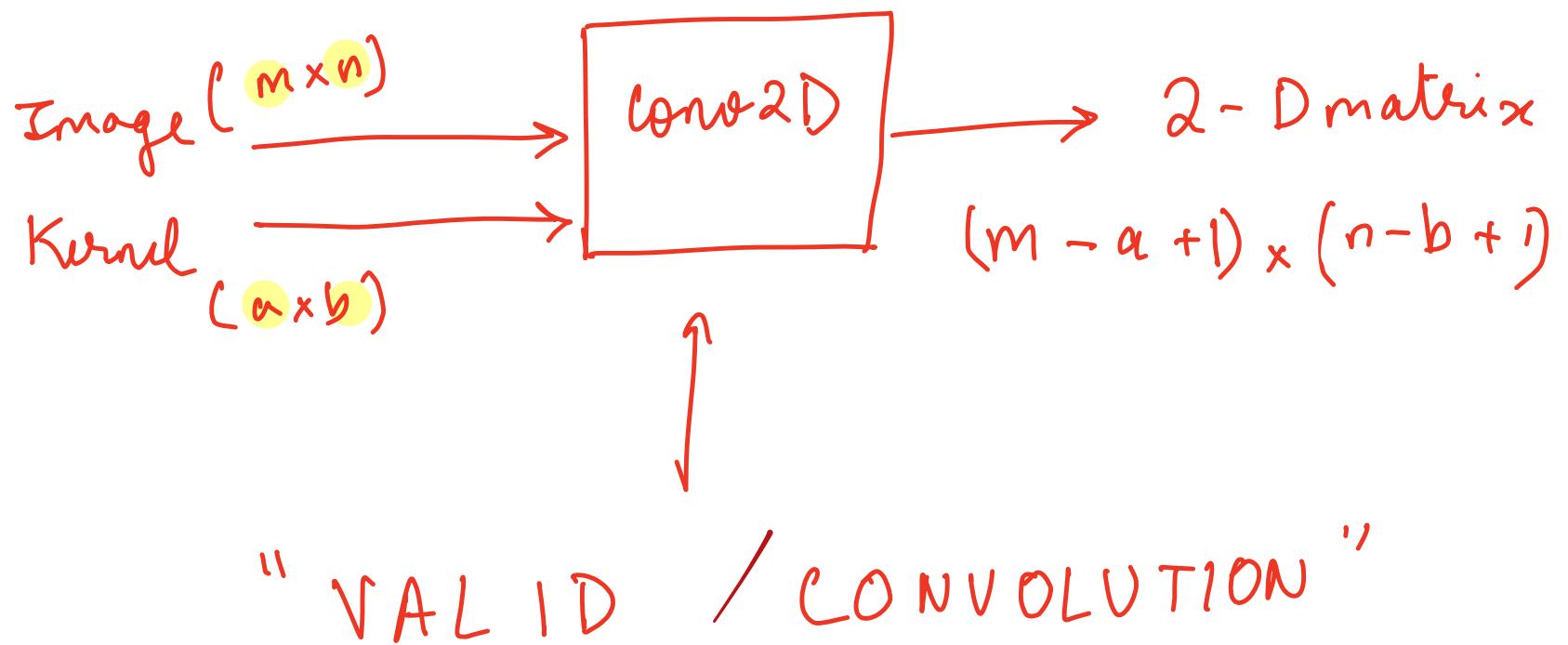
$\leftarrow \underline{6 \times 6}$.

No. of shifts : 4 horizontal,
4 vertical.



$m - a + 1$ (vertical).
 $n - b + 1$ (horizontal).

Convolution (2D)



2 more types : (a) SAME TYPE CONV.
(b) FULL TYPE CONV.

kernel

0 6 0
0 0 0
0 0

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

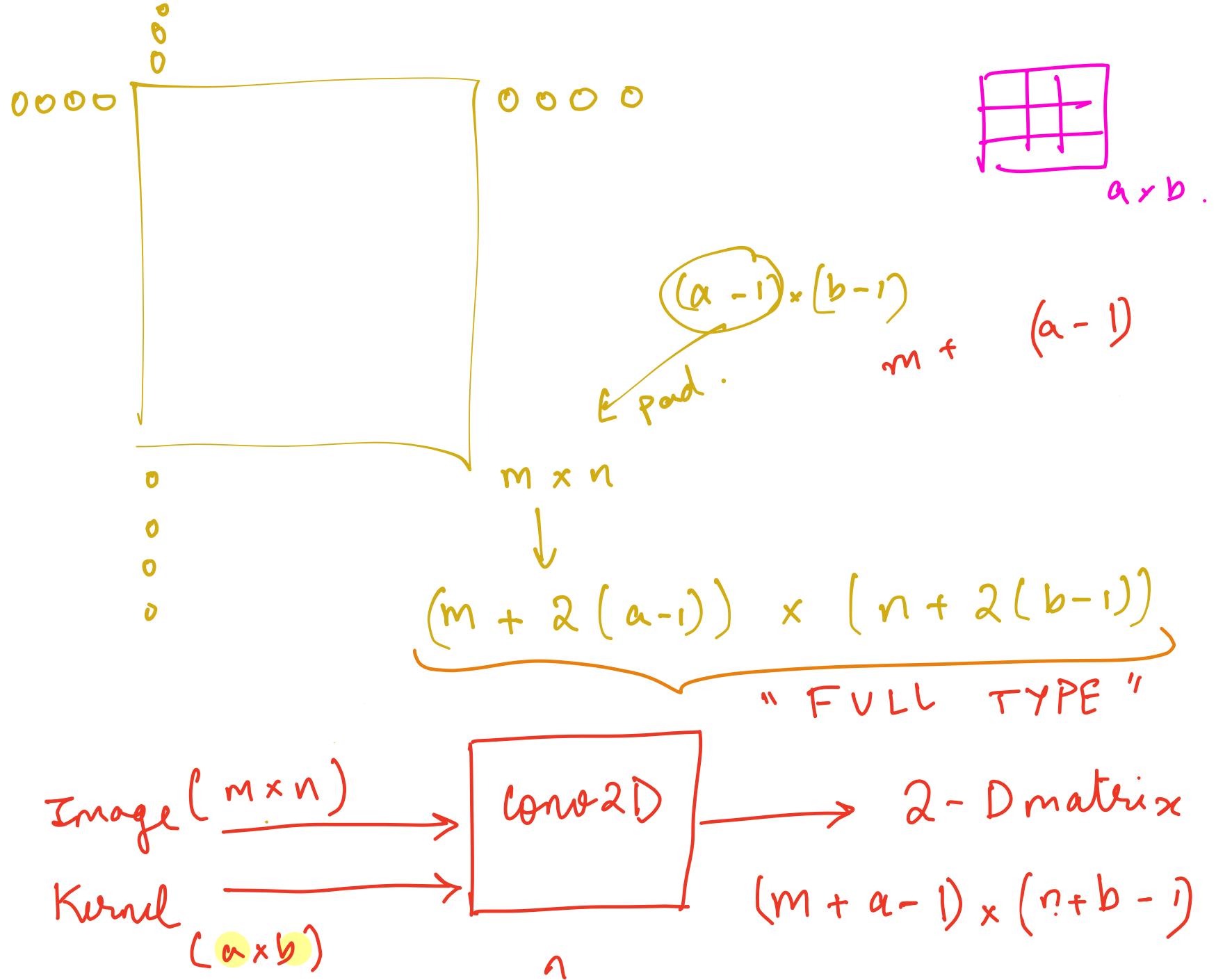
$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$
$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$
$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$

$$3 \times 3 - 2 \times 2$$
$$5 \times 5 - 4 \times 4$$
$$a \times b - (a-1)(b-1)$$

2 rows,
2 cols.

0 0 0 0 . - - .
0 6 0 0 - - - .

0 0
0 0



0 6
0 0

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

$$m + (a - i) + (a - 1) = m + 2(a - 1)$$

$a - 1$

0 0 0

0 0

0 0

0 0

0 0

0 0

0 0

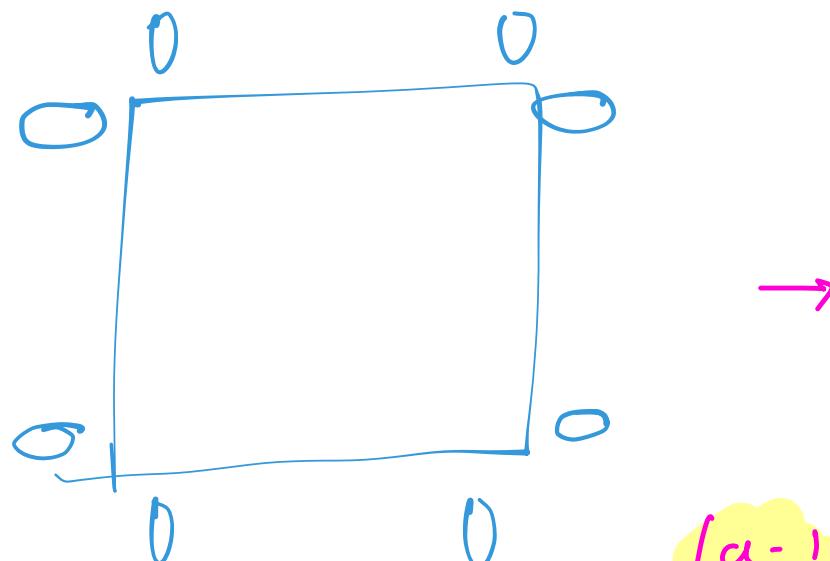
0 0

"Same" convolution \leftarrow To do same convolution,
we must have
odd length!

Image
 $m \times n$

Kernel
 $(a \times b)$

$(p, q) \leftarrow$ padding



$$m + 2 \cdot p - (a - 1) = m$$

$$n + 2 \cdot q - (b - 1) = n$$

$$p = \left\lceil \frac{a-1}{2} \right\rceil$$

$$q = \left\lceil \frac{b-1}{2} \right\rceil$$

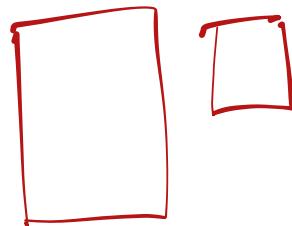
Summary:

- ① The CNN architecture is inspired from the Visual cortex.
- ② We use CNNs because flattening results in loss of **positional info**.
- ③ The solution is to use the "Convolution" operation.
- ④ The 3 main types are: $(m \times n \rightarrow \text{Image})$ $(a \times b \rightarrow \text{Kernel})$
 - if (i) Valid - $m - a + 1 \times n - b + 1$
 - if (ii) Same - $m \times n$
 - if (iii) Full - $m + a - 1 \times n + b - 1$

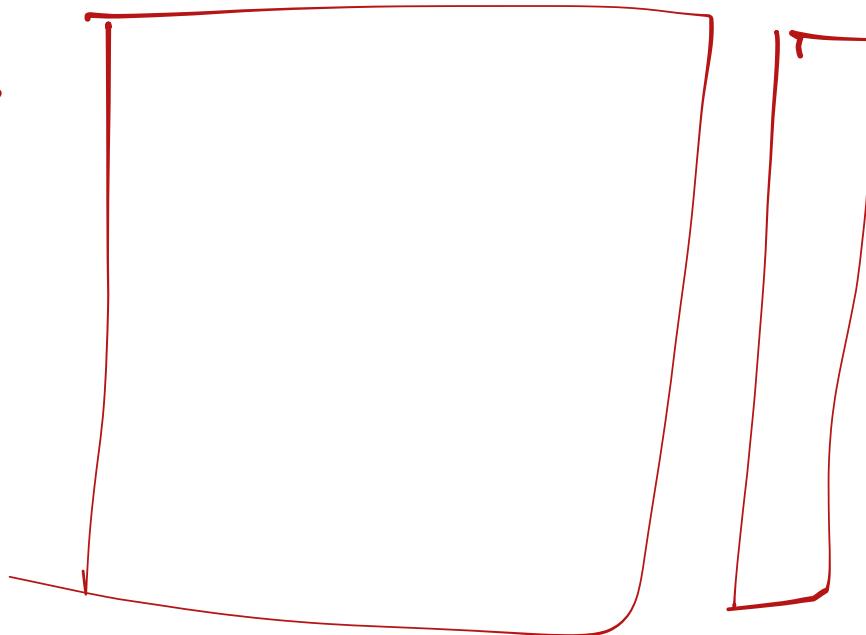


Topics :

- (a) Vision Transformer .
- (b) (Cross encoders) (NLP)
- (c) (Bi- Encoders)



$\xrightarrow{2D\text{-cons}}$



$\xrightarrow{\text{Toeplitz matrix.}}$

$$\begin{bmatrix} 1 & 2 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 2 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 2 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 2 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 2 & 1 \end{bmatrix}$$

$$\xrightarrow{W^T b \equiv x_1 * k_1}$$

```
import numpy as np
from scipy.signal import convolve2d

a = np.random.randn(640, 480)
b = np.random.randn(4, 4)
c = convolve2d(a, b)

print(c.shape)
```

(643, 483)

