# Introduction to Vision Models
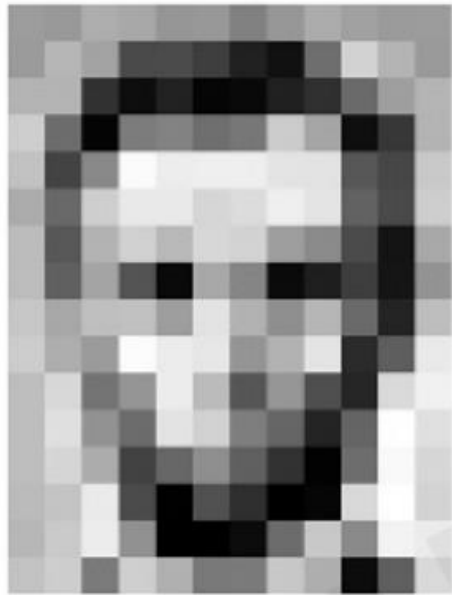
CNN, VIT AND BEYOND…!

# What is an image?

**How computers sees an image?**



Input Image

Pixel Representation

classification

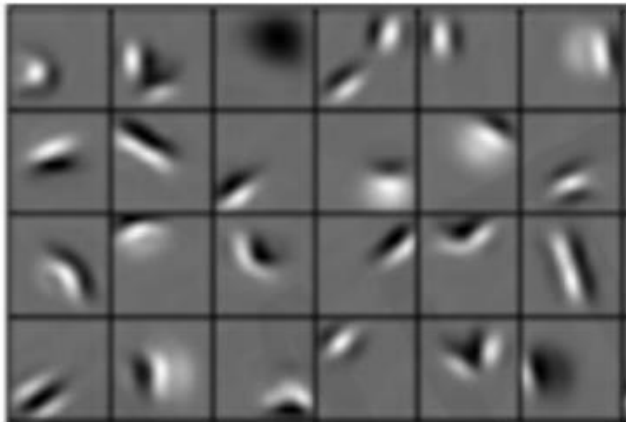| | |
|---|---|
| Lincoln | 0.8 |
| Washington | 0.1 |
| Jefferson | 0.05 |
| Obama | 0.05 |

**Local features within images?**



Low level features — Edges, dark spots

Mid level features — Eyes, ears, nose

High level features — Facial structure

How to process such data with NN?



**Input:**
- 2D image
- Vector of pixel values

But this is not how our brain Sees!

# How Brain actually sees things?

# Getting spatial views with filters!



Connect patch in input layer to a single neuron in subsequent layer.
Use a sliding window to define connections.

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

| 1 | -1 | 1 |
|---|----|----|
| -1 | 1 | -1 |
| 1 | -1 | 1 |

| -1 | -1 | 1 |
|----|----|----|
| -1 | 1 | -1 |
| 1 | -1 | -1 |

Input

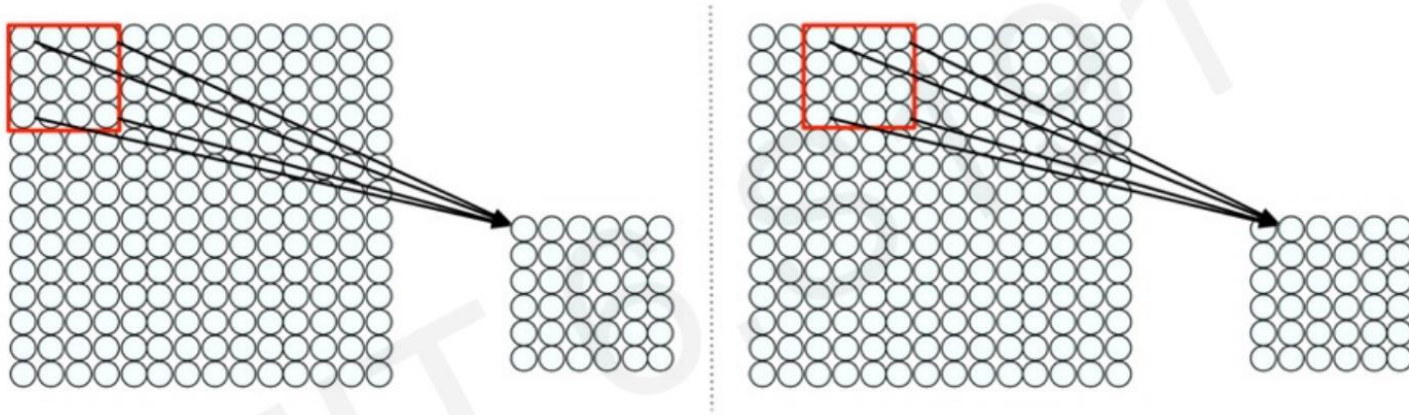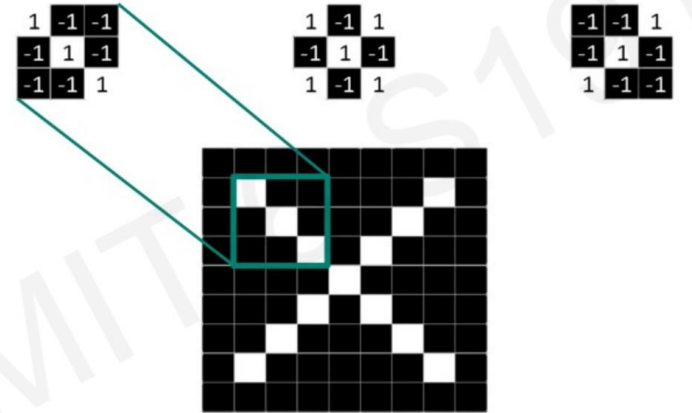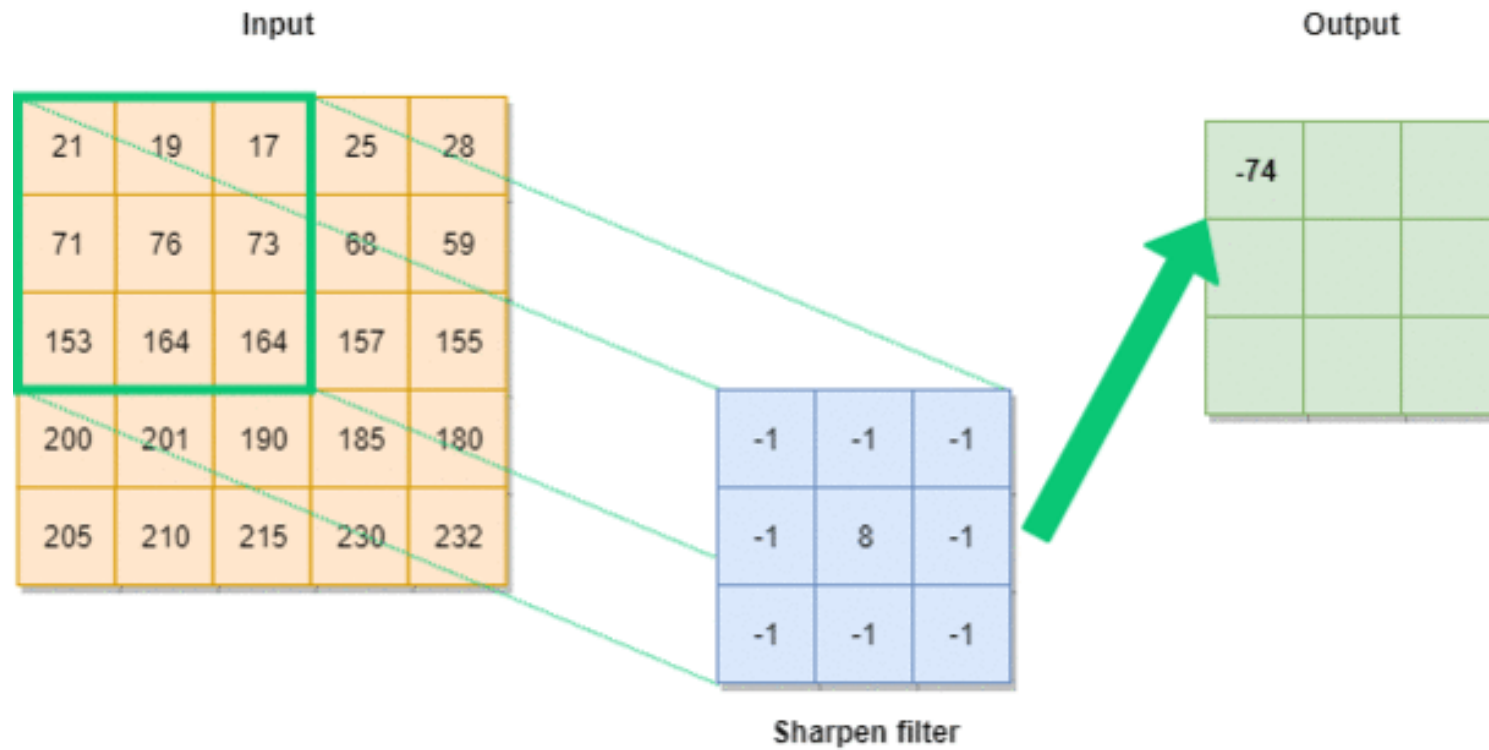| 21 | 19 | 17 | 25 | 28 |
|---|---|---|---|---|
| 71 | 76 | 73 | 68 | 59 |
| 153 | 164 | 164 | 157 | 155 |
| 200 | 201 | 190 | 185 | 180 |
| 205 | 210 | 215 | 230 | 232 |

Sharpen filter

| -1 | -1 | -1 |
|---|---|---|
| -1 | 8 | -1 |
| -1 | -1 | -1 |

Output

| -74 | | |
|---|---|---|
| | | |
| | | |

AIGeekProgrammer.com © 2019

filter 1
1x2x2

map 1
map 2
map 3
map 4
map 5
map 6
map 7
map 8

second CNN layer with 16 filters

filter a
8x2x2

map 1
map 2
map 3
map 4
map 5
map 6
map 7
map 8

map a
map b
map c
map d
map e
map f
map g
map h
map i
map j
map k
map l
map m
map n
map o
map p

Faces    Cars    Elephants    Chairs

Convolution Neural Network (CNN)