# What is "mode"?

↓ i/p
N = 5
K = 3
N-k+1
= 5-3+R
= 3
o/p
shrunk in size

```
convolve2d(A, np.fliplr(np.flipud(w)), mode='valid')
```

- The movement of the filter is bounded by the edges of the image. The output is therefore always smaller than the input. **Input Image**

Filter?

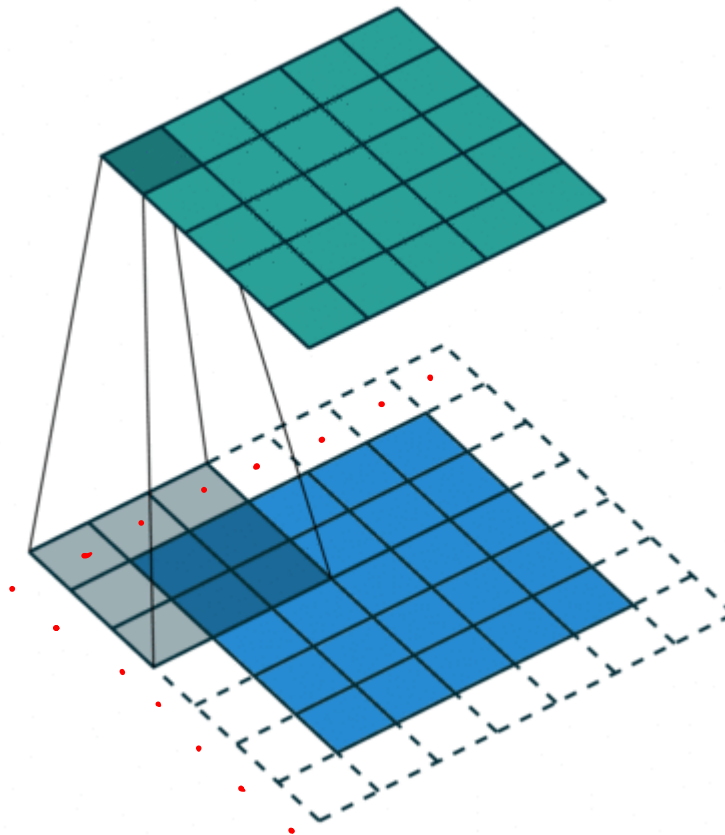| 0 | 1 | 2 |
|---|---|---|
| 0 | 1 | 2 |
| 0 | 1 | 2 |

Input Image:

| 3₀ | 3₁ | 2₂ | 1 | 0 |
|----|----|----|---|---|
| 0₂ | 0₂ | 1₀ | 3 | 1 |
| 3₀ | 1₁ | 2₂ | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

**Output Image**

| 12.0 | 12.0 | 17.0 |
|------|------|------|
| 10.0 | 17.0 | 19.0 |
| 9.0 | 6.0 | 14.0 |

# Padding ( "same" mode)

- What if we want the output to be of the **same size** as the input image?
- Then we can use padding i.e. add imaginary zeros around the input.



$$\text{Input} = N = 5 \, (\cdots \to 7)$$

$$\text{Filter} = K = 3$$

$$o/p = 5$$

$$N - K + 1$$
$$= 7 - 3 + 1$$
$$= 4 + 1$$
$$\to \quad = 5$$

Trainer: Dr. Darshan Ingle.

# Even more padding! ( "full " mode)

- We could extend the filter further and still get non-zero outputs.

- This is ~~not very common these days.~~

- "full" padding"
  - i.  Input Length = N
  - ii.  Kernel length = K
  - iii.  Output length = N+K-1

5

3

S+2+1

S+3-1

= 8-1

= 7

# Summary of Modes



It's here...

It's finally here!

- Input length = N
- Kernel length = K

| Mode | Output Size | Usage |
| --- | --- | --- |
| Valid | N - K + 1 | Typical |
| Same | N | Typical |
| Full | N + K - 1 | Atypical |

# A new perspective on Convolution

SO NOW WHAT?

Sliding Pattern Finder that passes through the entire image.

Input

\* Filter

=

The pattern is found here!

Trainer: Dr. Darshan Ingle.

# How to view Convolution as Matrix Multiplication?

- Lets see 1D Convolution first because its easy to understand and 2D Convolution is just going to be a generalization of this.
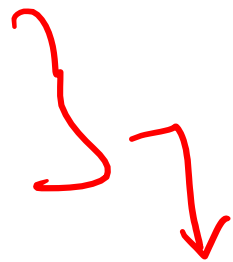
Input image: a = [$a_1$, $a_2$, $a_3$, $a_4$]

Filter: w = [$w_1$, $w_2$]

Output image: b = a * w = [$a_1w_1 + a_2w_2$, $a_2w_1 + a_3w_2$, $a_3w_1 + a_4w_2$]

$$a = (a_1, a_2, a_3, a_4)$$

$$w = (w_1, w_2)$$

CONVOLUTION

$$b = a * w = (a_1w_1 + a_2w_2, a_2w_1 + a_3w_2, a_3w_1 + a_4w_2)$$

# 1D Convolution in general

$$b_i = \sum_{i'=1}^{K} a_{i+i'} w_{i'}$$

- **Note:** It is very same as 2D Convolution's equation, just without 2nd index
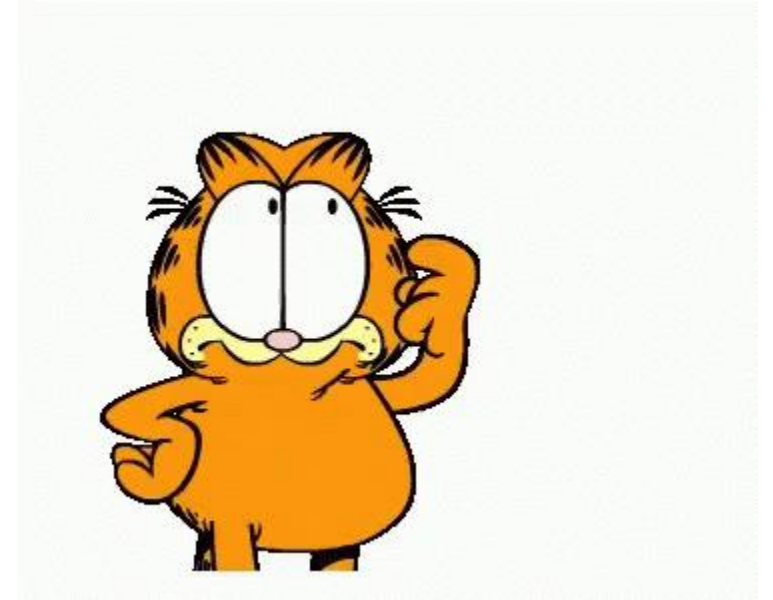
# Matrix Multiplication

*ID Array*

- Can we do the same operation using Matrix Multiplication.

- YES!!!

- HOW????????

*Matrix Multiplication*

*Kernel & I/p*

$$\begin{pmatrix} a_1 w_1 + a_2 w_2 \\ a_2 w_1 + a_3 w_2 \\ a_3 w_1 + a_4 w_2 \end{pmatrix} = \begin{pmatrix} w_1 & w_2 & 0 & 0 \\ 0 & w_1 & w_2 & 0 \\ 0 & 0 & w_1 & w_2 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{pmatrix}$$

b                              w                              a

# So what do we understand?

- By repeating the same filter again and again, we can do convolution without actually doing convolution.

$$\begin{pmatrix} a_1w_1 + a_2w_2 \\ a_2w_1 + a_3w_2 \\ a_3w_1 + a_4w_2 \end{pmatrix} = \begin{pmatrix} w_1 & w_2 & 0 & 0 \\ 0 & w_1 & w_2 & 0 \\ 0 & 0 & w_1 & w_2 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{pmatrix}$$

Trainer: Dr. Darshan Ingle.

# Problem with Matrix Multiplication?

# Problem with Matrix Multiplication?

① filter gets repeated multiple times

② It takes up a lot more space

③ Filter $w$ was vector of size 2.

In Matrix mult, $w$ becomes a 2D Matrix of size __3×4.__

$$\begin{pmatrix} a_1 w_1 + a_2 w_2 \\ a_2 w_1 + a_3 w_2 \\ a_3 w_1 + a_4 w_2 \end{pmatrix} = \begin{pmatrix} w_1 & w_2 & 0 & 0 \\ 0 & w_1 & w_2 & 0 \\ 0 & 0 & w_1 & w_2 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{pmatrix}$$

# Lets replace Matrix Multiplication with Convolution

MATRIX Mult. is BIG & SLOW.

∴ We will go for Convolution which is SMALL & FAST.

# Parameter Sharing / Weight Sharing

$$a = w^T \cdot x$$

whr  $a = $ o/p Activation

$X = $ i/p feature vector

$w = $ weight matrix.

Why What if I use the same weights $(w_1, w_2)$ again & again?

↳ ADV:

① less parameters

② use lesser RAM

③ Computation bcomes far more efficient.

∴ Convolution ~~takes~~ saves both SPACE & TIME by using less weights.

# Why do this?

Consider a fully connected ANN

Ex: MNIST dataset

28×28 = 784-sized input vector

What if the image size is slightly larger & it is a color image?

Ex: CIFAR-10 : 32×32×3 = 3072 -sized input vector.

Modern CNN such VGG : 224×224
  If we are using full weight matrix, the #i/p features = 1,50,528 features

Modern HD image : 1280×720 = 2.8 million features.

Too Large for a NN.

# Why do this?

Convolution $\equiv$ Corelation.

FILTER $\equiv$ Pattern Finder.

We want the same filter to look at all locations on the image.
↳ This is called as Translational Invariance.

# Translational Invariance

- Suppose we are building a Dog vs cat classifier



Cat          Cat

If we use a fully connected NN, then the NN has to learn the weights for each of these positions separately.

& yet after that, NN wont generalize that well bcz if we come across the same cat in a different position, NN shall fail to recognize it.

# Translational Invariance



Cat



Cat

∴ Having a Shared pattern finder makes more sense.

# Convolution on Color Images

Color: 3D     : H × W × C     (Ht × Width × Color)



Original Color Image

Matlab RGB Matrix

# Convolution on Color Images



2-D

3-D

Note: same size in depth dimension

Input (5×5)

filter 2×2

4×4

o/p Image

N−k+1

5−2+1

= 4

# Convolution on Color Images

2-D (2-D "dot product") - a grayscale pattern-finder

$$(A * w)_{ij} = \sum_{i'=1}^{K} \sum_{j'=1}^{K} A(i + i', j + j') w(i', j')$$

E.g. if the filter is looking for red circles it won't match a green circle

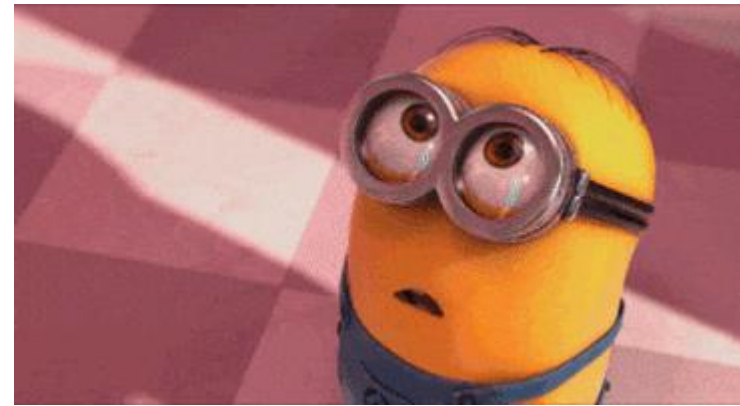3-D (3-D "dot product") - a color pattern-finder

$$(A * w)_{ij} = \sum_{c=1}^{3} \sum_{i'=1}^{K} \sum_{j'=1}^{K} A(i + i', j + j', c) w(i', j', c)$$

Trainer: Dr. Darshan Ingle.

# What more? 3D

- Input Image: H x W x 3, which is a 3D vector

- Kernel: K x K x 3, which is a 3D vector

I/P = 3D
O/P = 2D

- Output Image: (H-K+1) x (W-K+1), which is a 2D vector

- We know that Neural Networks have a repeating structures (i.e. they have "uniformity")

- Ex: For a Dense layer, if the input is 1D, then the output is also 1D and hence can be fed from one layer to another.

- But, in our case, we see that the output is 2D and input is 3D.

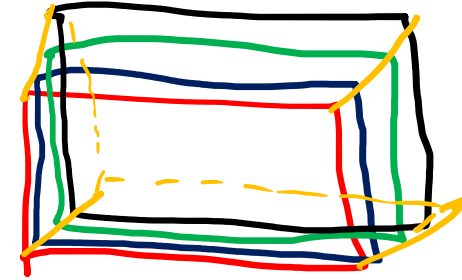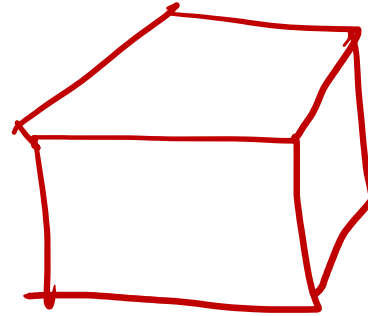- Question: So how do we solve this?
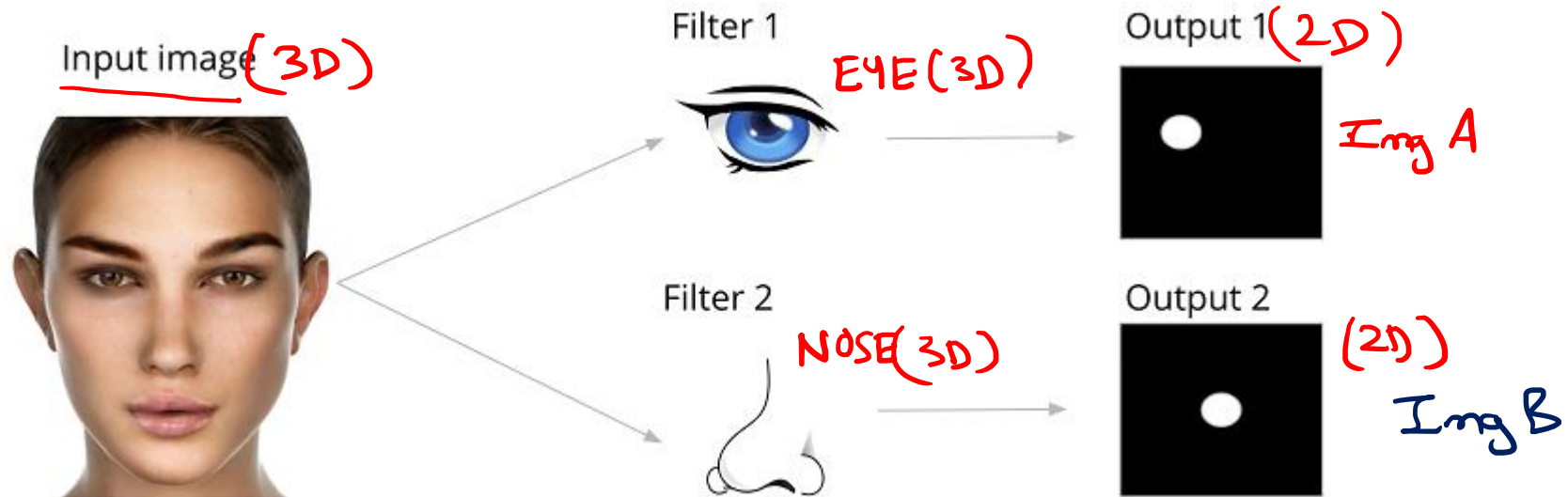
- Answer: Lets see this now.

2D
I/P = N (2D)
Filter = k
O/P : N-k+1
(2D)

$$(A * w)_{ij} = \sum_{c=1}^{3} \sum_{i'=1}^{K} \sum_{j'=1}^{K} A(i + i', j + j', c) w(i', j', c)$$
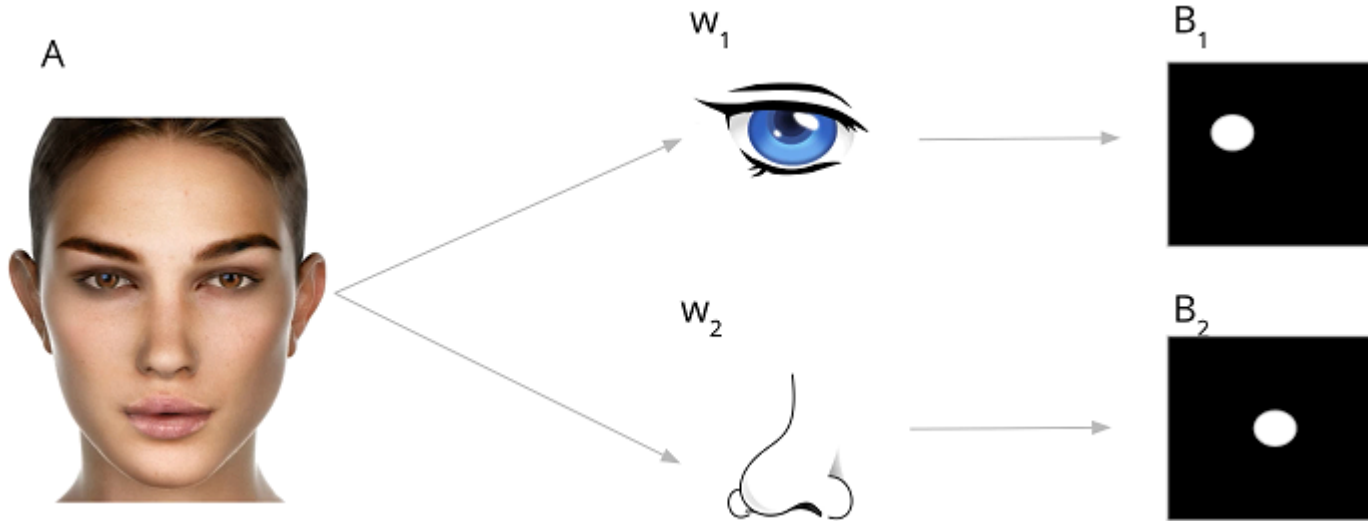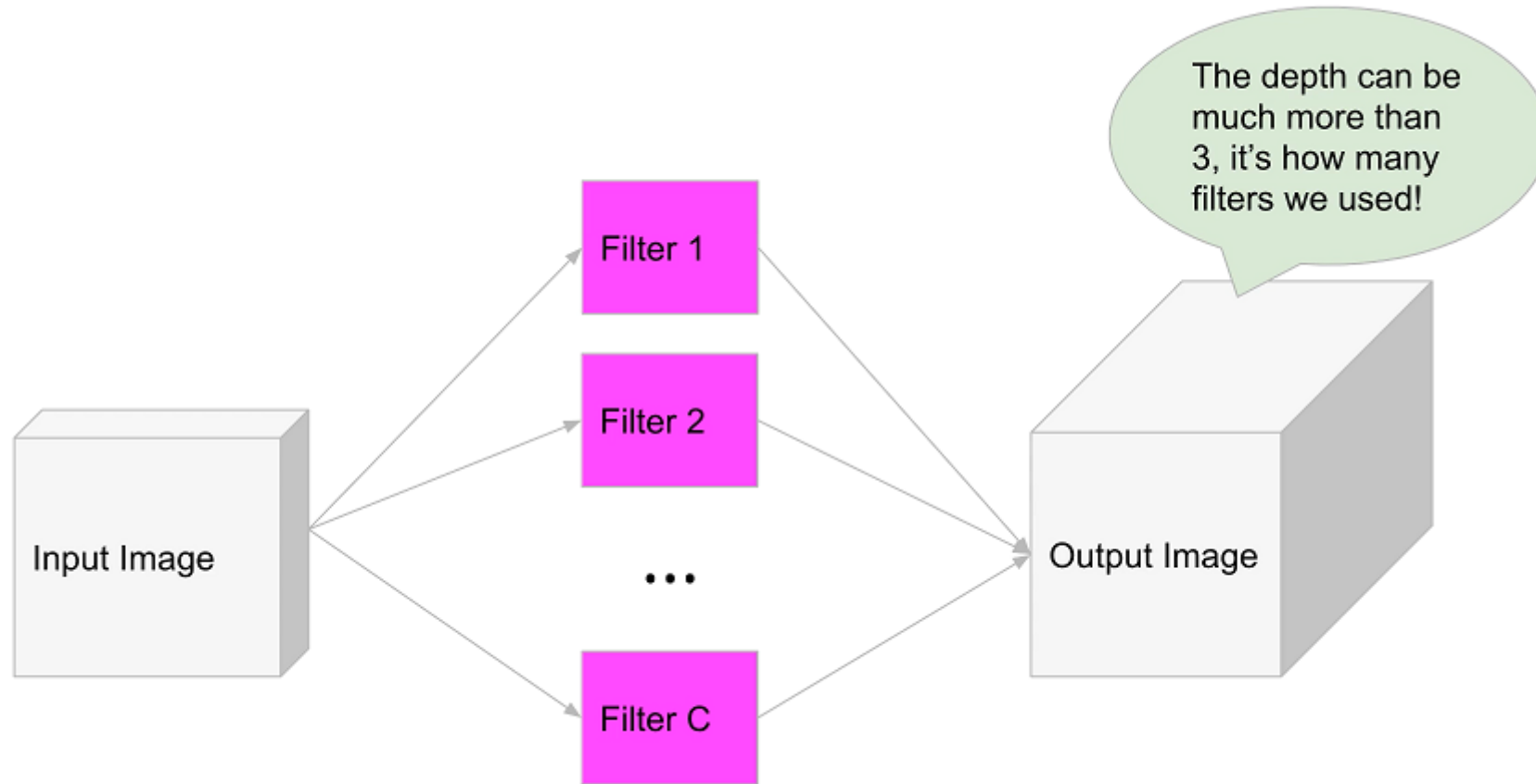
# Multiple Features



eye, nose, lips, eyebrows

Input image (3D)

Filter 1

EYE (3D)

Output 1 (2D)

Img A

Filter 2

NOSE (3D)

Output 2 (2D)

Img B

# Multiple Features

- Consider we have an image A with dimensions: H x W x 3 **3D**

- If we use the "same" mode, and apply one filter, then we get output, lets say B1 with dimensions: H x W **2D** *eye*

- If we use the "same" mode, and apply one more filter, then we get output, lets say B2 with dimensions: H x W **2D** *nose*

- If we stack up B1 and B2 together, we get a new output image, lets say B with dimensions: H x W x 2, i.e. a 3D image as our original image



A    W₁    B₁

W₂    B₂

# Multiple Features



Trainer: Dr. Darshan Ingle.

# Convolution in Deep Neural Network

- Lets vectorize this operation, we don't need to do each color convolution separately

$$B = A \overset{convolve}{*} w$$

A = Input Img, w = filter

B = Output Img

$$shape(A) = H \times W \times C_1$$

$$shape(w) = C_1 \times K \times K \times C_2$$

$$shape(B) = H \times W \times C_2$$

$$B(i,j,c) = \sum_{i'=1}^{K} \sum_{j'=1}^{K} \sum_{c'=1}^{C_1} A(i+i',j+j',c')w(c',i',j',c)$$



SPLOID GIF

# So What is this 3$^{rd}$ dimension in the Output image?

3rd is certainly not the color.

I/p to NN: H × W × 3

After subsequent convolutions, we get H × W × (arbitrary #)

3rd dimension = feature maps.

# Convolution Layer

$$\text{Conv layer}: \ \underline{\sigma}(W \overset{\swarrow \text{convolution}}{*} x + \underline{b})$$

$$\text{Dense layer}: \ \underline{\sigma}(W^{T} \cdot x + \underline{b})$$

# Shape of the bias

- In a Dense layer, if $W^T x$ is a vector of size M, $b$ is also a vector of size M

- In a Conv layer, $b$ **does not** have the same shape as $W * x$ (a 3-D image)

- Technically, this is not allowed by the rules of matrix arithmetic

- But the rules of broadcasting (in Numpy code) allow it

- If $W * x$ has the shape H x W x $C_2$, then $b$ is a vector of size $C_2$
  - One scalar per feature map

$$Conv\ layer : \sigma(W * x + b)$$
$$Dense\ layer : \sigma(W^T x + b)$$

# How much do we save? <u>CONVOLUTION</u>



Has convolution saved my time & space?

Input Img: $32 \times 32 \times \underline{3}$

Filter : $3 \times 5 \times 5 \times 64$ (i.e. 64 feature maps)

<u>Output</u> Img: $28 \times 28 \times 64$ (assuming mode = 'valid')

why 28? $N - k + 1$

$$M - k + 1 = 32 - 5 + 1 = 28$$

$\therefore$ # of parameters (ignore the bias) $= 3 \times 5 \times 5 \times 64 = 4800$

# How much do we save?  MATRIX MULTIPLICATION

Flattened Input Img = $32 \times 32 \times 3 = 3072$

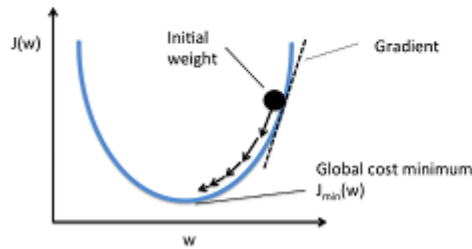—"— O/p Img = $28 \times 28 \times 64 = 50176$

Weights Matrix = $3072 \times 50176 = 1,54,140,672 \approx 154$ million

Compared to Convolution, we see that it takes $\sim 32000$ times more parameters.
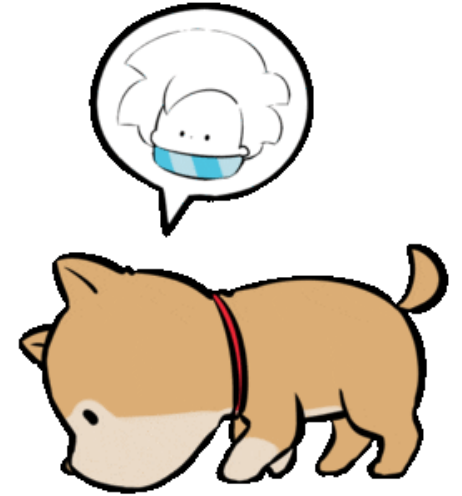
$\therefore$ Huge RAM. $\Rightarrow$ Suboptimal

# How are Convolution filters found?

- Since convolution is just a part of some neural network layer, it's easy to conceive of how the filters will be found

- Initially, we looked at convolution as an image *modifier* (blur, edge)

- Now, we see it as a pattern finder / shared-parameter matrix multiplication / feature transformer

- In other words, W will be found the same as before, automatically!

- Still gradient descent, `model.fit()`

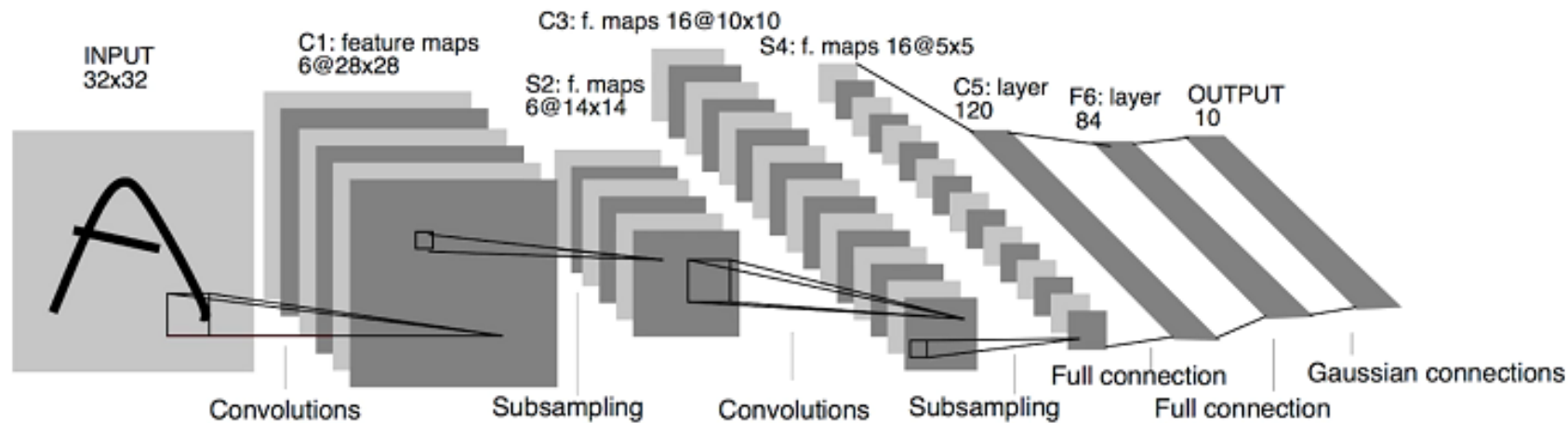$$Conv\ layer: \sigma(W * x + b)$$

$$Dense\ layer: \sigma(W^T x + b)$$

# CNN Architecture
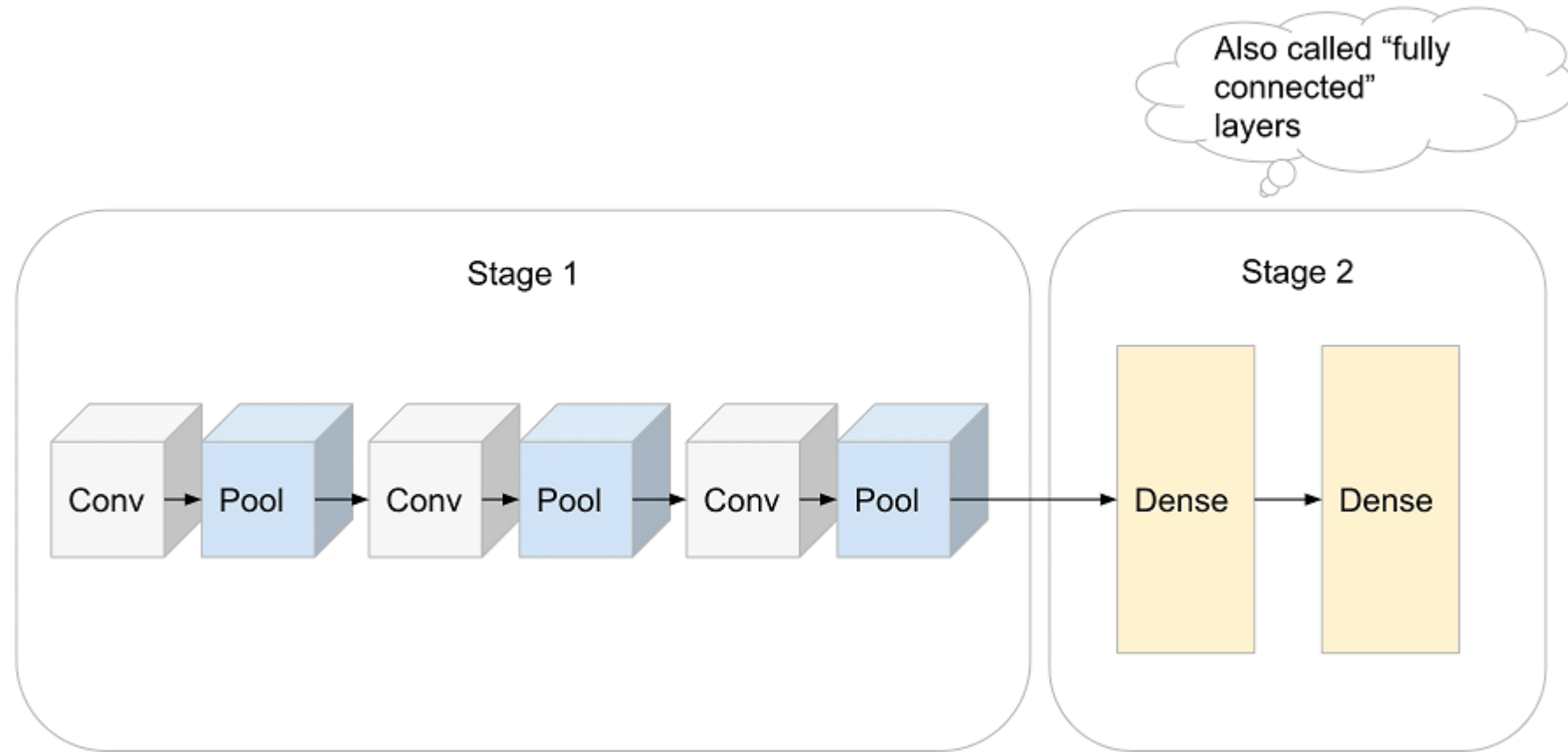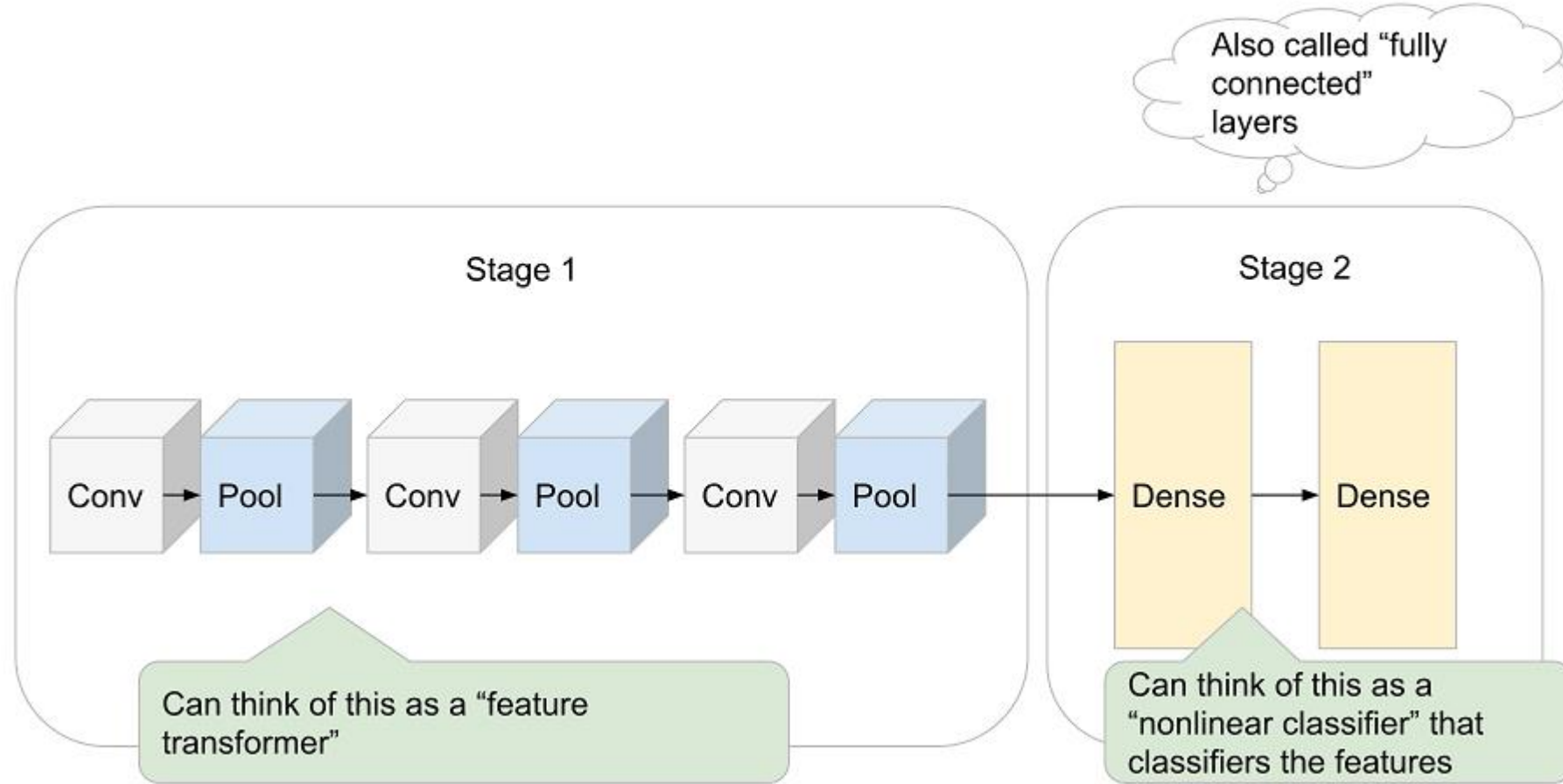
Trainer: Dr. Darshan Ingle.

# CNN Architecture

- Now that you understand how exactly a convolution layer works including the bias term and activation function we can now consider the architecture of a convolutional neural network and why it's that way.

- So as a little bit of a history lesson, modern CNN is essentially all originated from the same model, the LeNet.

- This is named after Yann LeCun, one of the original Deep Learning pioneers, along with Geoff Hinton and Yoshua Bengio.
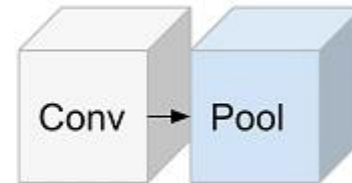
# Typical CNN

# Typical CNN

# Pooling means Downsampling

i.e. output a Smaller image from a bigger one.

eg: If i/p image = $100 \times 100$

then if pool_size = 2

then my o/p image = $50 \times 50$
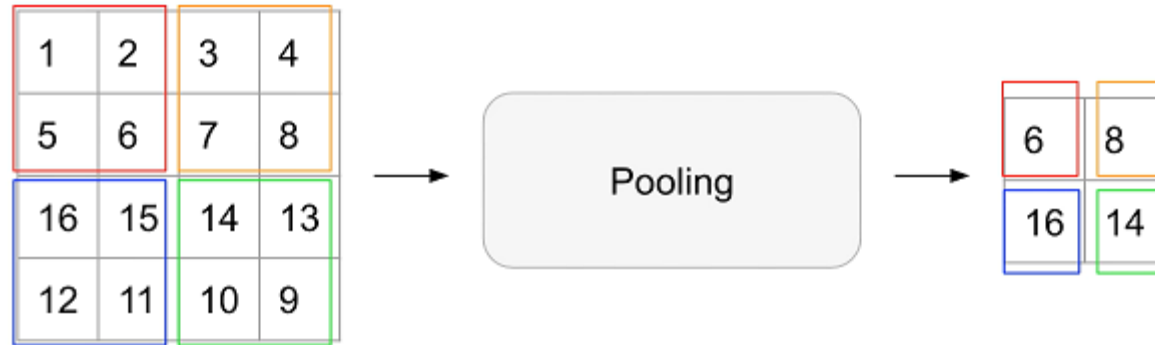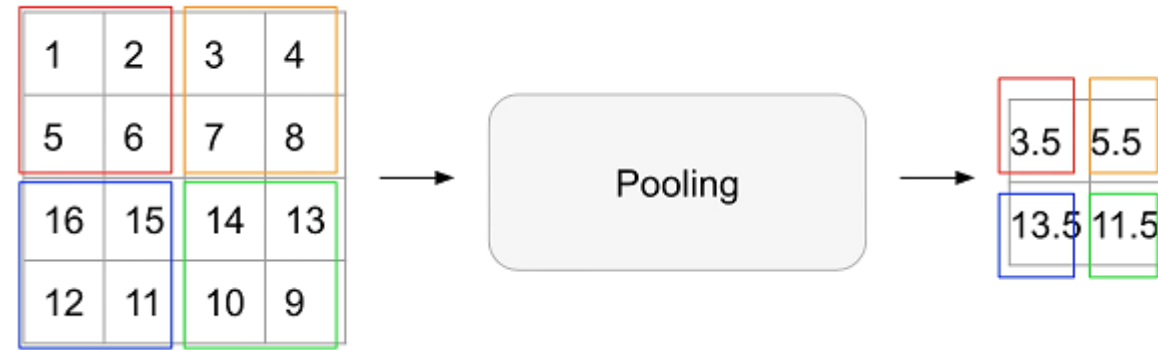
∴ [Downsample by 2.]

POOLING

# Types of Pooling

- There are two types of Pooling:

1. Max pooling

2. Average pooling

- Which one to use is a Hyperparameter choice.

# Max Pooling

# Average Pooling

# Why to use Pooling?



- Practical: If we shrink the image, we have less data to process.

- Translational Invariance: I don't care where in the image the feature occurred, I just care that it did.