

CAP 5415: Computer Vision



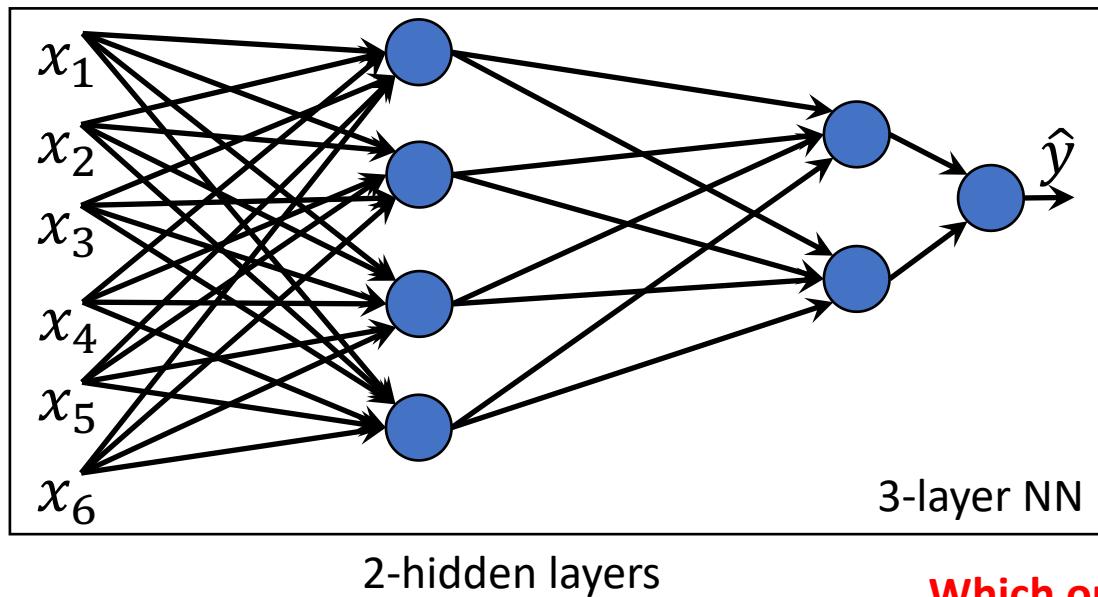
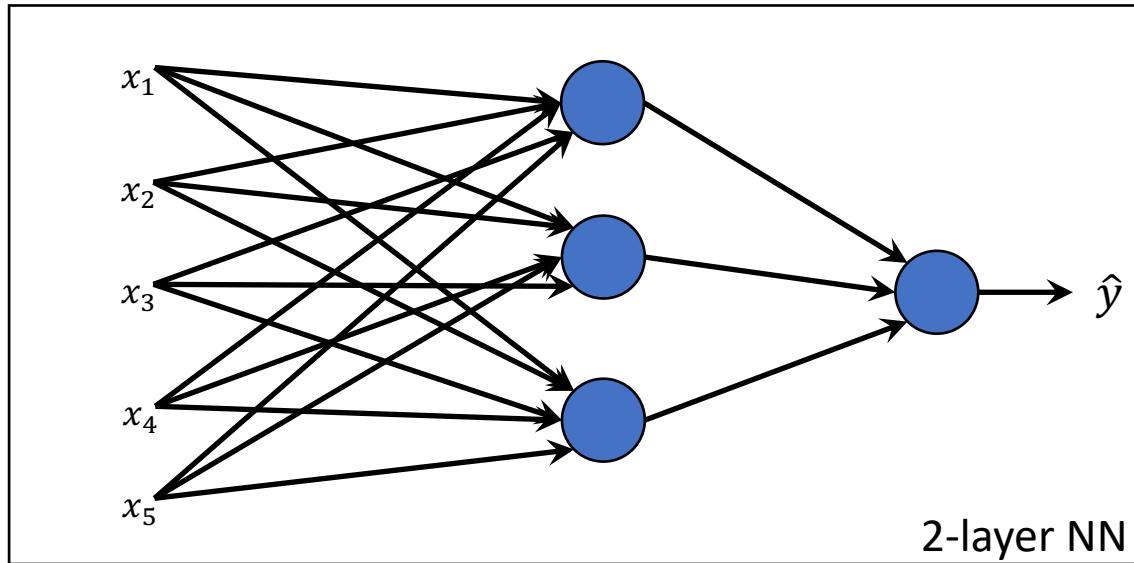
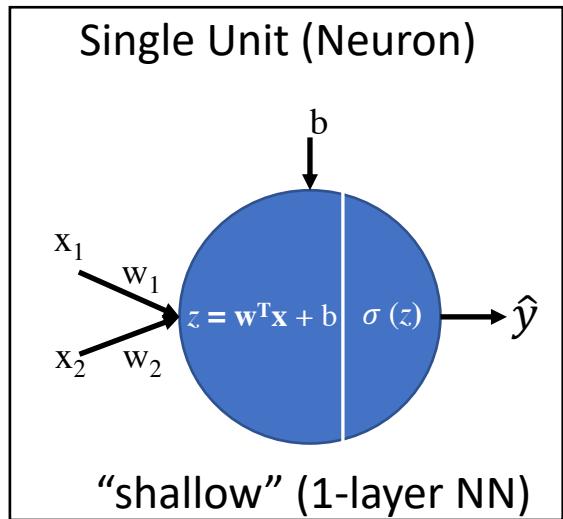
Intro to Convolutional Neural Networks



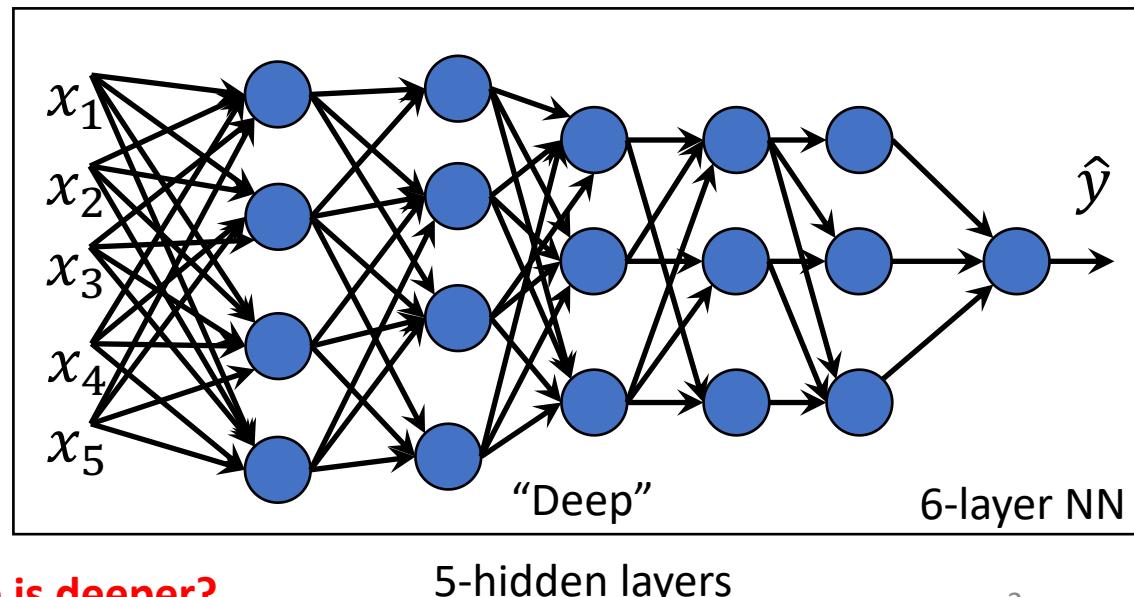
Dr. Sedat Ozer



Deep vs. Shallow



Which one is deeper?

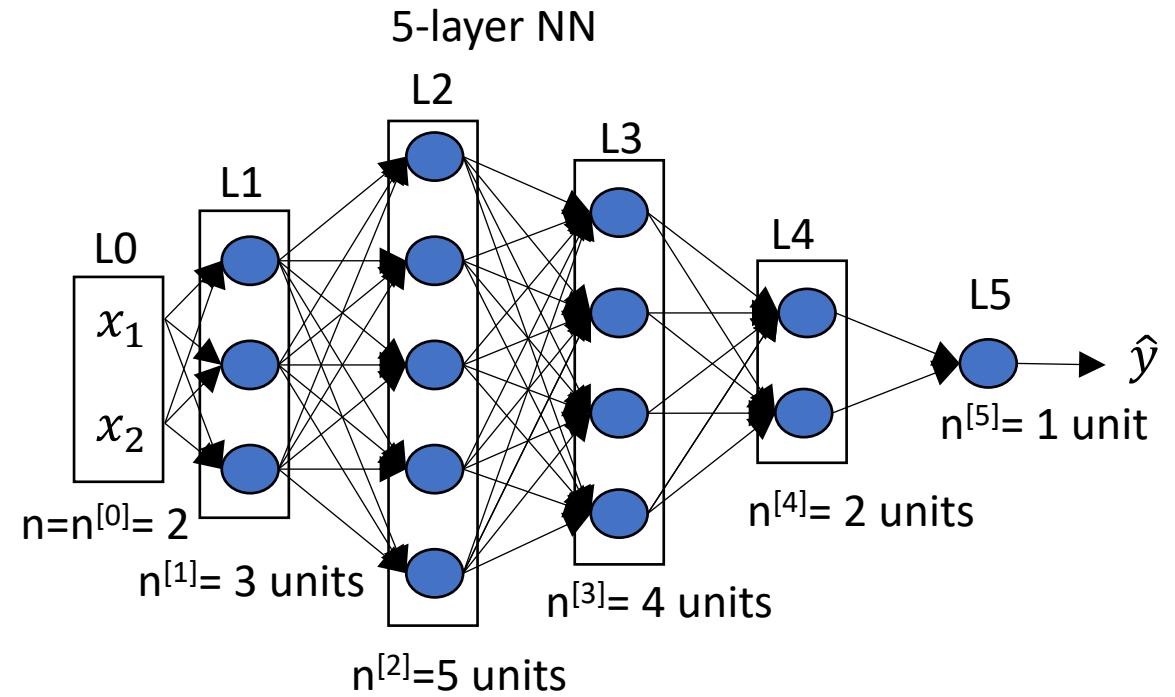


A deep NN: Notation

$$A^{[0]} = X = \begin{bmatrix} & | & | & | \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} \dots & \mathbf{x}^{(m)} \\ & | & | & | \end{bmatrix} \quad (n^{[0]} \times m)$$

$$A^{[l]} = \begin{bmatrix} & | & | & | \\ \mathbf{a}^{[l](1)} & \mathbf{a}^{[l](2)} \dots & \mathbf{a}^{[l](m)} \\ & | & | & | \end{bmatrix} \quad (n^{[l]} \times m)$$

$(Z^{[l]}$ has the same shape as $A^{[l]}$)



$$Z^{[l]} = W^{[l]} A^{[l-1]} + \mathbf{b}^{[l]}$$

$$A^{[l]} = g(Z^{[l]})$$

Dims

	Dims
$\mathbf{z}^{[l]}$:	$(n^{[l]} \times 1)$
$\mathbf{a}^{[l]}$:	$(n^{[l]} \times 1)$
$\mathbf{W}^{[l]}$:	$(n^{[l]} \times n^{[l-1]})$

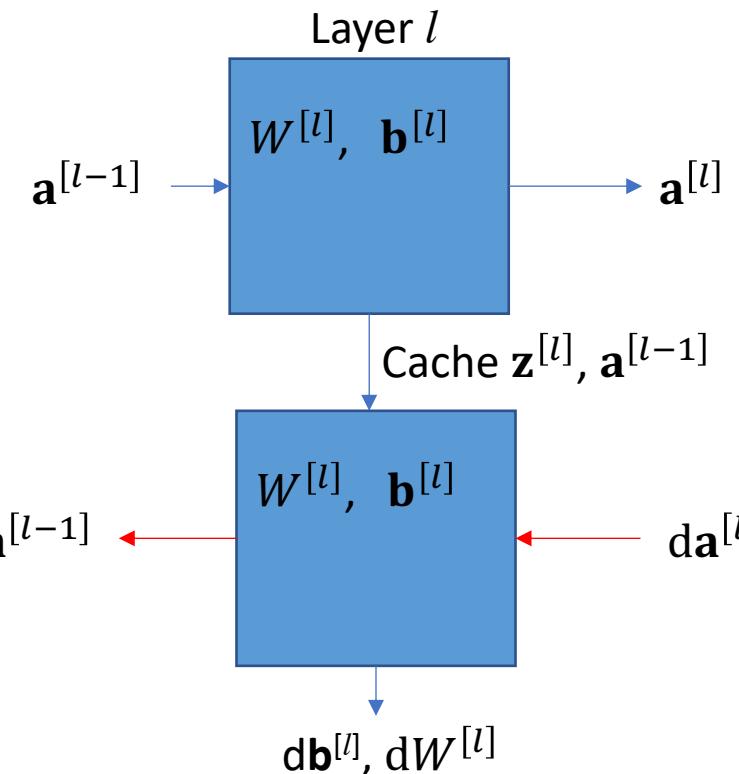
Hyperparameters

- We have parameters: $W^{[l]}, \mathbf{b}^{[l]}$ for each layer (actual parameters)
- We also have additional parameters affecting the performance of our algorithm: Hyperparameters. I.e., the hyperparameters are the parameters that affect our actual parameters. Examples:
 - Learning rate,
 - number of iterations,
 - total number of hidden layers,
 - hidden unit numbers in each layer,
 - Choice of activation function,
 - Momentum,
 - mini batch size,
 - regularization parameters,
 - choice of cost function, choice of the optimization technique, ...

Forward vs. Backward computation

- Forward direction computes the output and the loss.
- Backward direction computes the derivatives to find the update directions for the parameters: weights & biases.

Forward direction:



Backward direction:

Element-wise product

$$dz^{[l]} = da^{[l]} * g^{[l]'}(z^{[l]})$$

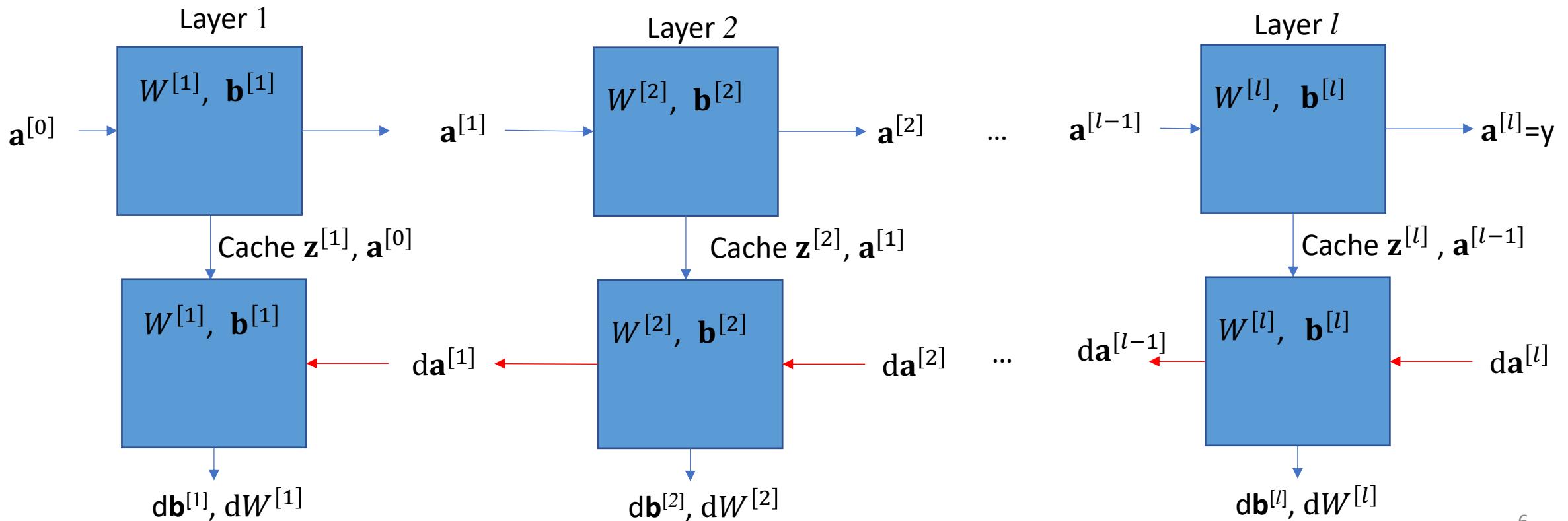
$$dW^{[l]} = dz^{[l]} a^{[l-1]}$$

$$db^{[l]} = dz^{[l]}$$

$$da^{[l-1]} = W^{[l]T} dz^{[l]}$$

Forward vs. Backward computation

$$\begin{aligned} \mathbf{d}\mathbf{z}^{[l]} &= \mathbf{d}\mathbf{a}^{[l]} * g'^{[l]}(\mathbf{z}^{[l]}) \\ \mathbf{d}W^{[l]} &= \mathbf{d}\mathbf{z}^{[l]} \mathbf{a}^{[l-1]} \\ \mathbf{d}\mathbf{b}^{[l]} &= \mathbf{d}\mathbf{z}^{[l]} \\ \mathbf{d}\mathbf{a}^{[l-1]} &= W^{[l]T} \mathbf{d}\mathbf{z}^{[l]} \end{aligned}$$



Backward Propagation

Element-wise multiplication

$$d\mathbf{z}^{[l]} = d\mathbf{a}^{[l]} * g^{[l]'}(\mathbf{z}^{[l]})$$

$$dW^{[l]} = d\mathbf{z}^{[l]} \mathbf{a}^{[l-1]}$$

$$d\mathbf{b}^{[l]} = d\mathbf{z}^{[l]}$$

$$d\mathbf{a}^{[l-1]} = W^{[l]T} d\mathbf{z}^{[l]}$$

For a single sample

Element-wise multiplication

$$d\mathbf{Z}^{[l]} = d\mathbf{A}^{[l]} * g^{[l]'}(\mathbf{z}^{[l]})$$

$$dW^{[l]} = \frac{1}{m} d\mathbf{Z}^{[l]} \mathbf{A}^{[l-1]T}$$

$$d\mathbf{b}^{[l]} = \frac{1}{m} \text{np.sum}(d\mathbf{Z}^{[l]}, axis=1, keepdims=True)$$

$$d\mathbf{A}^{[l-1]} = W^{[l]T} d\mathbf{Z}^{[l]}$$

For m data samples (“Vectorized” versions)

Convolutional Neural Networks

Story so far...

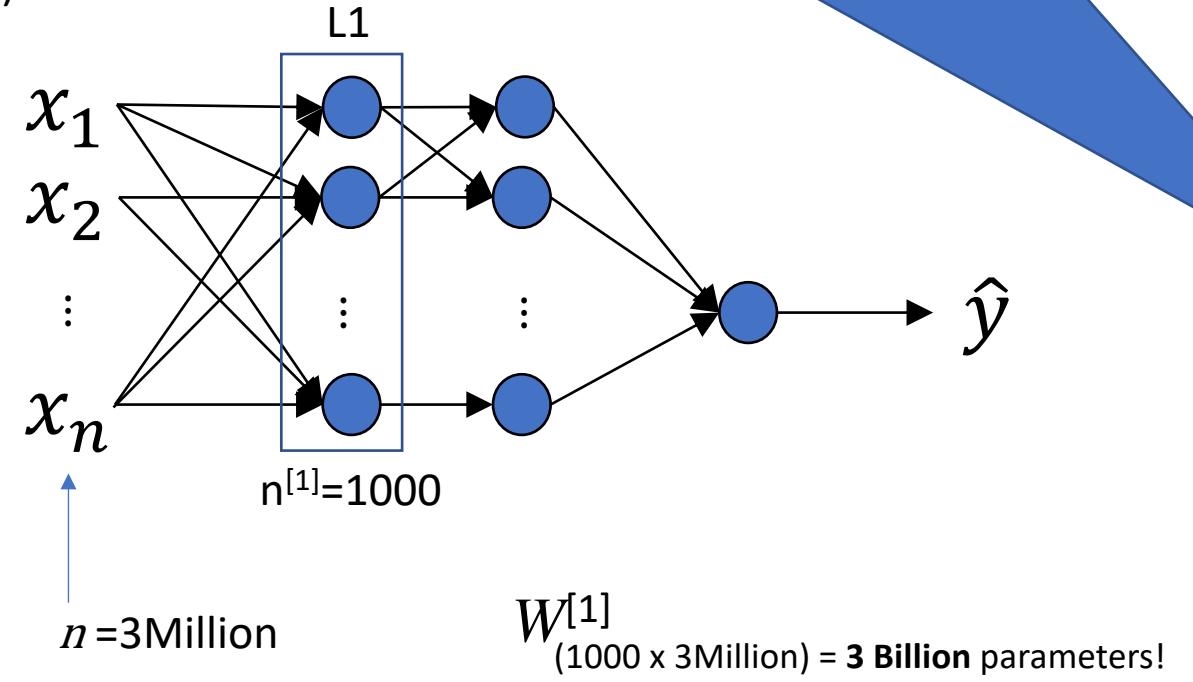
- Up until now, we learned how to
- Lets consider the case where the



A high-resolution image: $1000 \times 1000 \times 3$
 $= (3\text{Million} \times 1)$

As an alternative approach:
Use a convolutional NN !!

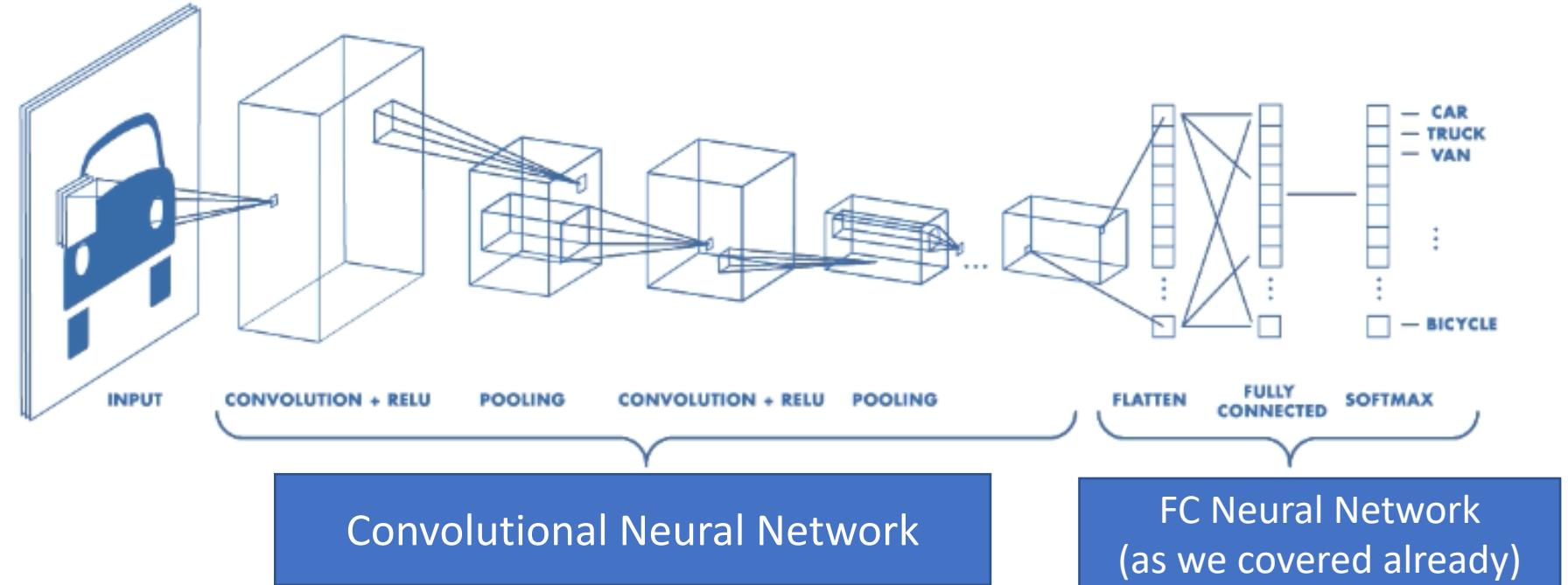
Problem: Is there a person in this picture?
(answer: yes/no)



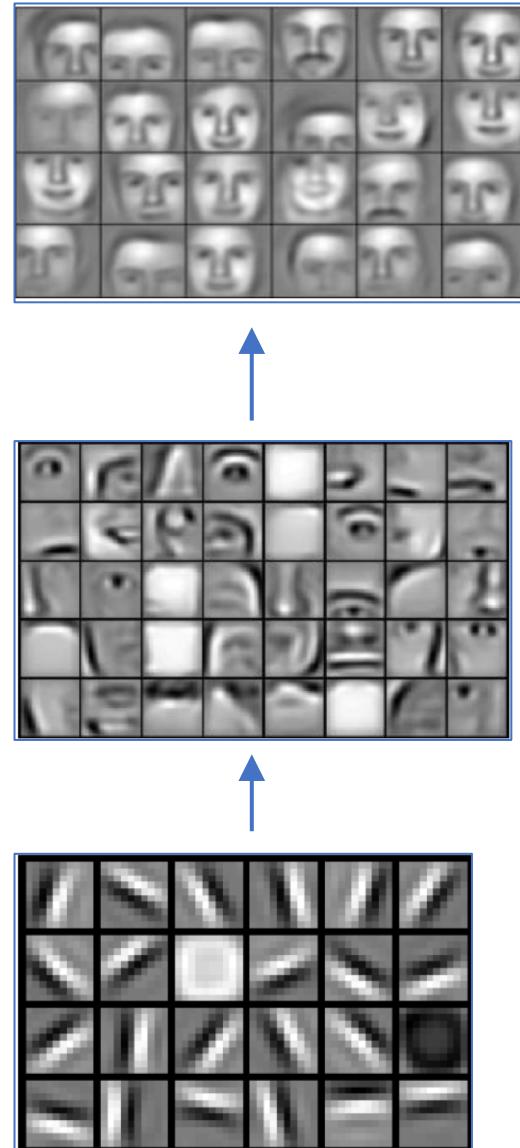
Convolutional Neural Networks (CNN)

Includes many additional components - operations. Most common ones are:

- Filtering,
- Padding,
- Stride,
- Pooling.



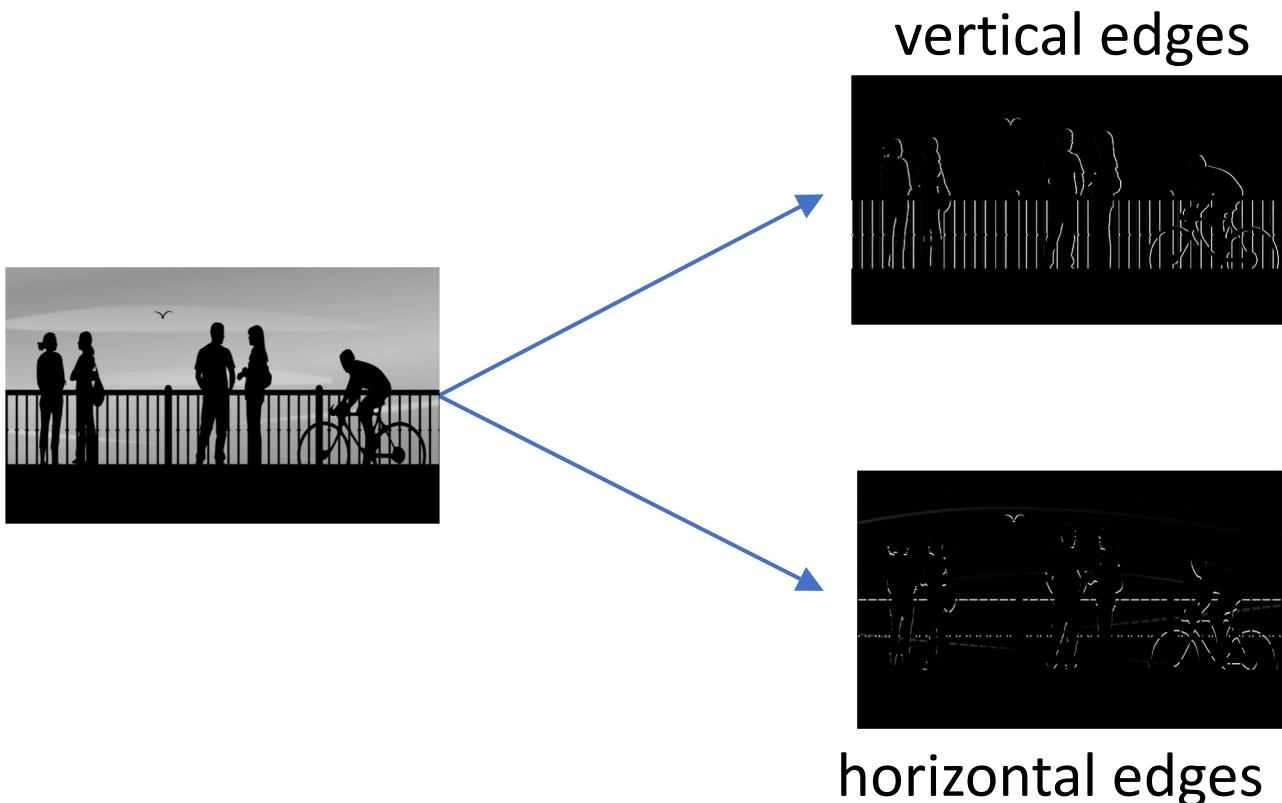
Features ?



Edge Detection

- To explain the Convolutional NNs, we will look at the edge detection example first.

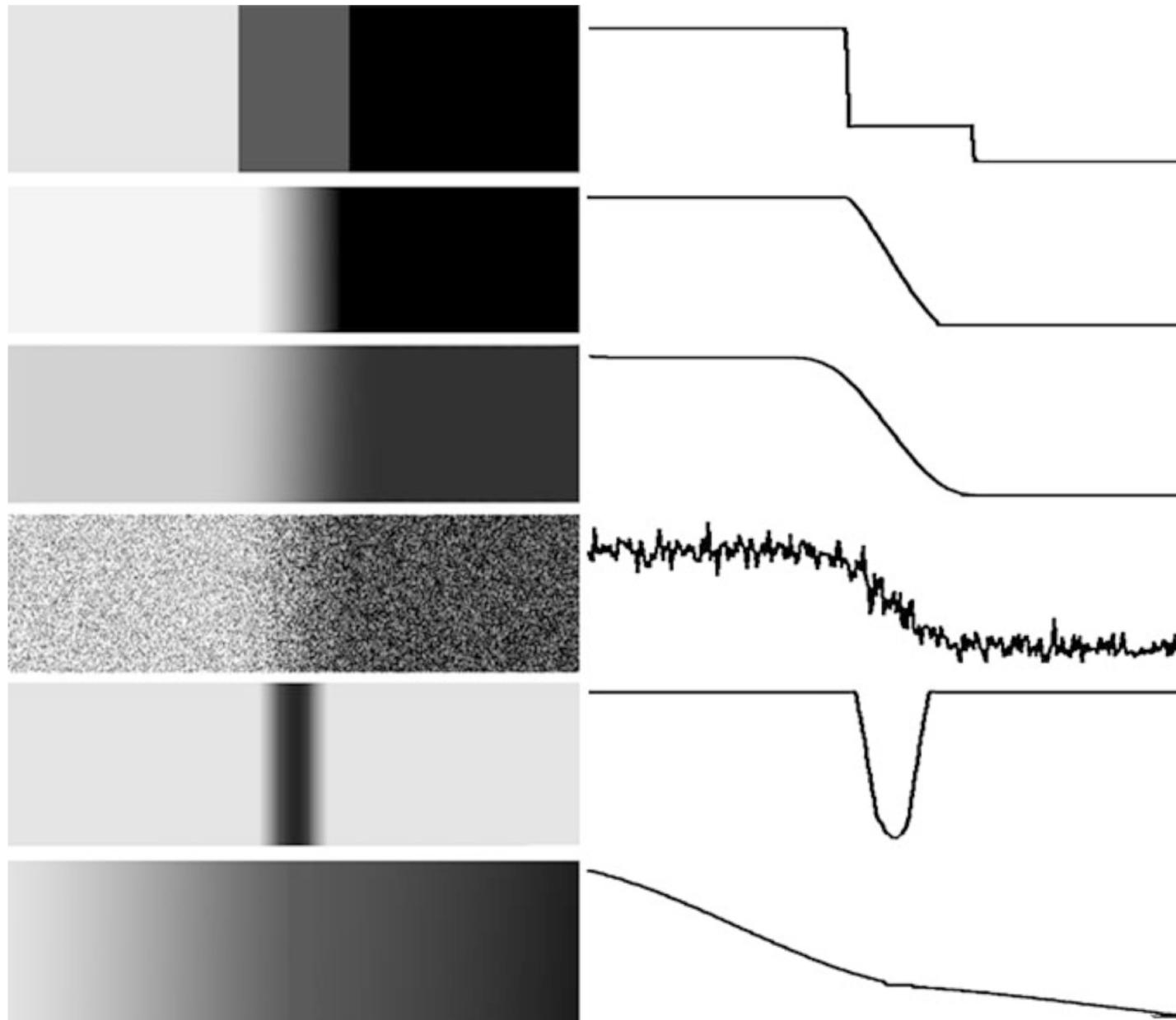
Edges



EDGES

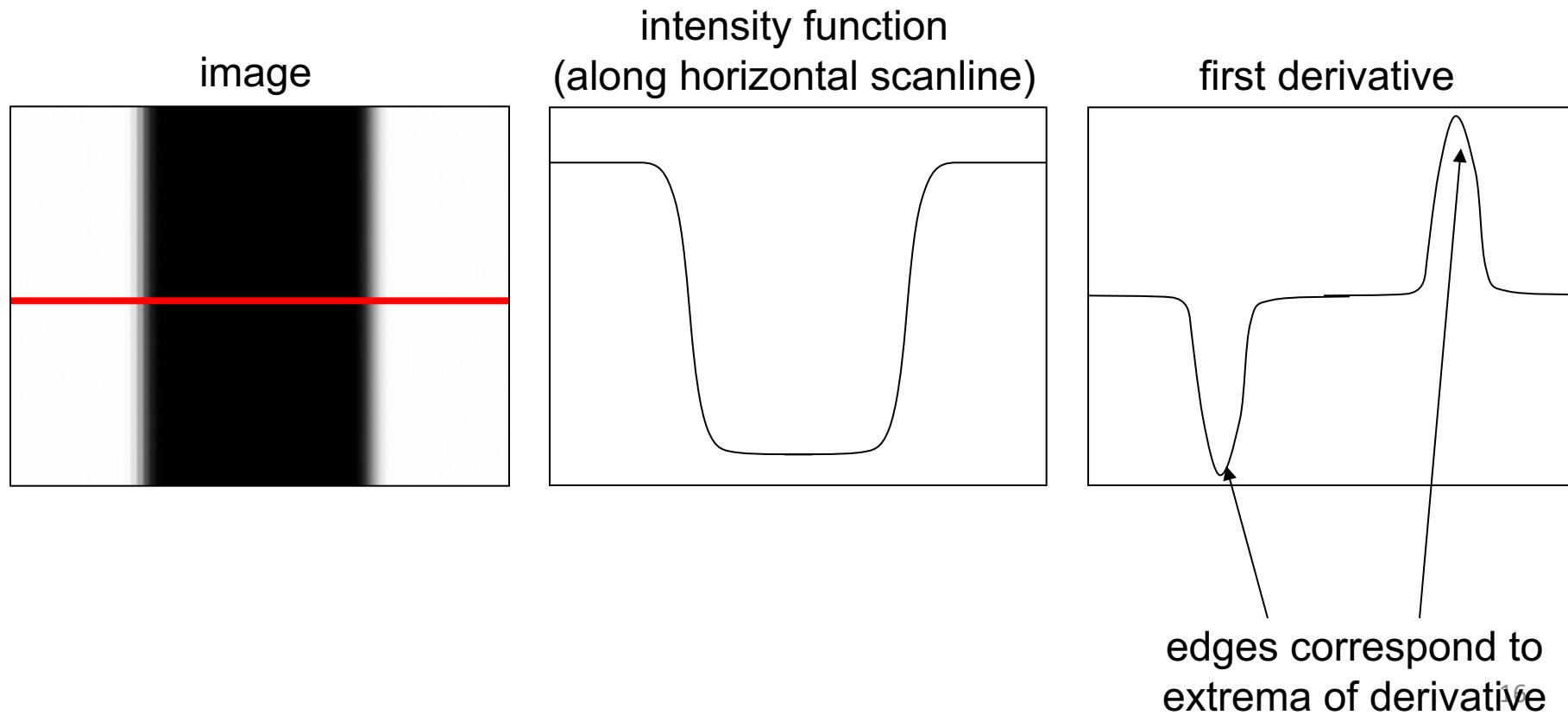
- **Discontinuities** in images are features that are often useful for initializing an image analysis procedure.
- Edges are important information for understanding an image; by moving “non-edge” data we also **simplify** the data.

Edge Models



Characterizing Edges

- An edge is a place of rapid change in the image intensity function



Derivatives and Average

- **Derivative:** rate of change
- Examples:
 - Speed is a rate of change of a distance, $X=V.t$
 - Acceleration is a rate of change of speed, $V=a.t$

Derivative

$$\frac{df}{dx} = \lim_{\Delta x \rightarrow 0} \frac{f(x) - f(x - \Delta x)}{\Delta x} = f'(x) = f_x$$

Example: $y = x^2 + x^4$

$$\frac{dy}{dx} = 2x + 4x^3$$

Discrete Derivative

$$\frac{df}{dx} = \lim_{\Delta x \rightarrow 0} \frac{f(x) - f(x - \Delta x)}{\Delta x} = f'(x)$$

$$\frac{df}{dx} = \frac{f(x) - f(x - 1)}{1} = f'(x)$$

$$\frac{df}{dx} = f(x) - f(x - 1) = f'(x)$$

Discrete Derivative / Finite Difference

$$\frac{df}{dx} = f(x) - f(x-1) = f'(x)$$

Backward difference

$$\frac{df}{dx} = f(x) - f(x+1) = f'(x)$$

Forward difference

$$\frac{df}{dx} = f(x+1) - f(x-1) = f'(x)$$

Central difference

Example: Finite Difference

$$f(x) = 10 \quad 15 \quad 10 \quad 10 \quad 25 \quad 20 \quad 20 \quad 20$$
$$f'(x) = 0 \quad \begin{matrix} [-1 & 1] \\ 5 \end{matrix} \quad \begin{matrix} [-1 & 1] \\ -5 \end{matrix} \quad \begin{matrix} [-1 & 1] \\ 0 \end{matrix} \quad \begin{matrix} [-1 & 1] \\ 15 \end{matrix} \quad \begin{matrix} [-1 & 1] \\ -5 \end{matrix} \quad \begin{matrix} [-1 & 1] \\ 0 \end{matrix} \quad \begin{matrix} [-1 & 1] \\ 0 \end{matrix}$$

Derivative Masks

Backward difference $[-1 \quad 1]$

Forward difference $[1 \quad -1]$

Central difference $[-1 \quad 0 \quad 1]$

Derivative in 2-D

Given function

$$f(x, y)$$

The diagram consists of two rectangular boxes. The left box is labeled "Image Axis 1" and the right box is labeled "Image Axis 2". Two blue arrows point from these boxes towards the variables x and y respectively in the mathematical expression $f(x, y)$.

Gradient vector

$$\nabla f(x, y) = \begin{bmatrix} \frac{\partial f(x, y)}{\partial x} \\ \frac{\partial f(x, y)}{\partial y} \end{bmatrix} = \begin{bmatrix} f_x \\ f_y \end{bmatrix}$$

Gradient magnitude

$$|\nabla f(x, y)| = \sqrt{f_x^2 + f_y^2}$$

Gradient direction

$$\theta = \tan^{-1} \frac{f_x}{f_y}$$

Derivative of Images

Derivative masks (filters): $f_x \Rightarrow \frac{1}{3} \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}$ $f_y \Rightarrow \frac{1}{3} \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$

$$I = \begin{bmatrix} 10 & 10 & 20 & 20 & 20 \\ 10 & 10 & 20 & 20 & 20 \\ 10 & 10 & 20 & 20 & 20 \\ 10 & 10 & 20 & 20 & 20 \\ 10 & 10 & 20 & 20 & 20 \end{bmatrix}$$
$$I_x = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & \boxed{10} & \boxed{10} & 0 & 0 \\ 0 & 10 & 10 & 0 & 0 \\ 0 & 10 & 10 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

(Vertical) Edge Detection

4 ¹	1 ⁰	4 ⁻¹	0 ⁻¹	1 ⁻⁰	2 ⁻¹
2 ¹	7 ⁰	4 ⁻¹	3 ⁻¹	0 ⁻⁰	1 ⁻¹
0 ¹	6 ⁰	5 ⁻¹	4 ⁻¹	2 ⁻⁰	2 ⁻¹
5 ¹	2 ⁰	1 ⁻¹	0 ⁻¹	4 ⁻⁰	2 ⁻¹
1	5	3	5	6	6
3	2	2	1	7	5

Convolution

*

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

=

-7	7	10	2
-3	8	4	2
-3	-4	-3	-1
3	3	-11	-7

Edge Detection (1)

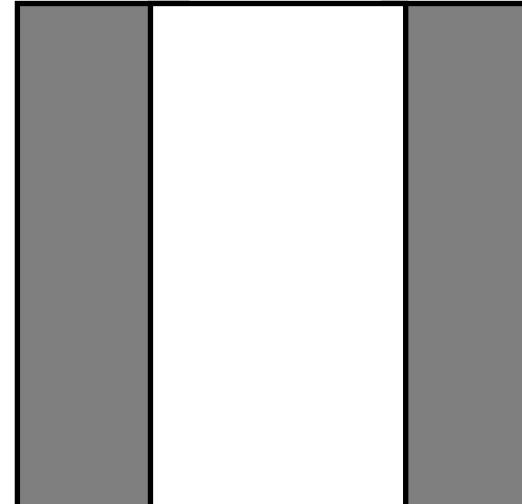
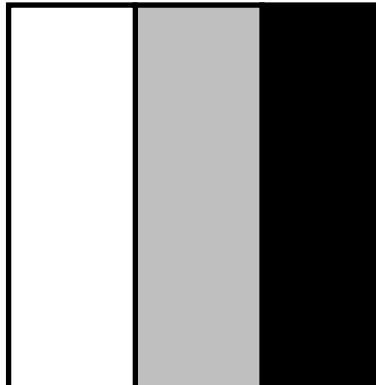
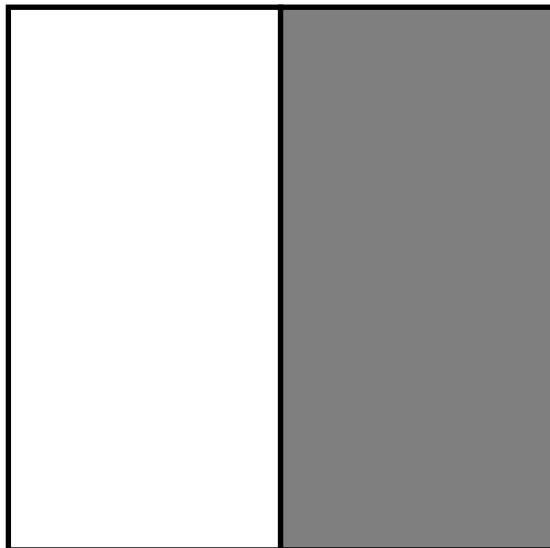
20	20	20	0	0	0
20	20	20	0	0	0
20	20	20	0	0	0
20	20	20	0	0	0
20	20	20	0	0	0
20	20	20	0	0	0

*

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

=

0	60	60	0
0	60	60	0
0	60	60	0
0	60	60	0



Edge Detection (2)

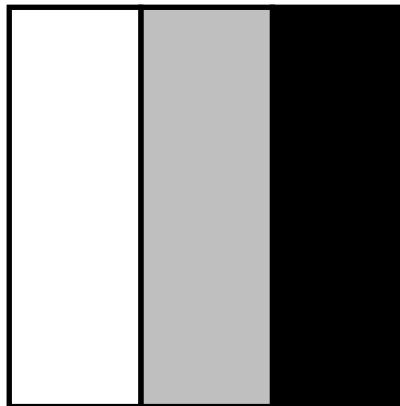
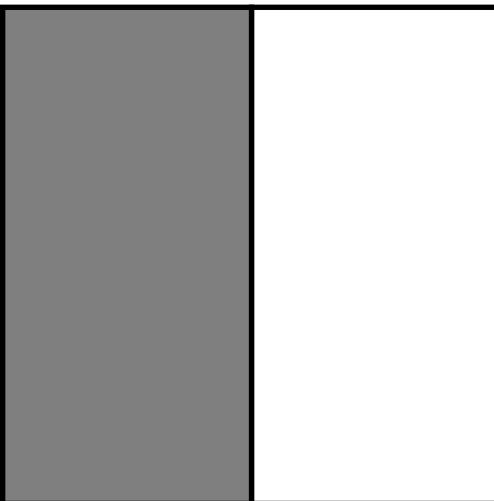
0	0	0	20	20	20
0	0	0	20	20	20
0	0	0	20	20	20
0	0	0	20	20	20
0	0	0	20	20	20
0	0	0	20	20	20

*

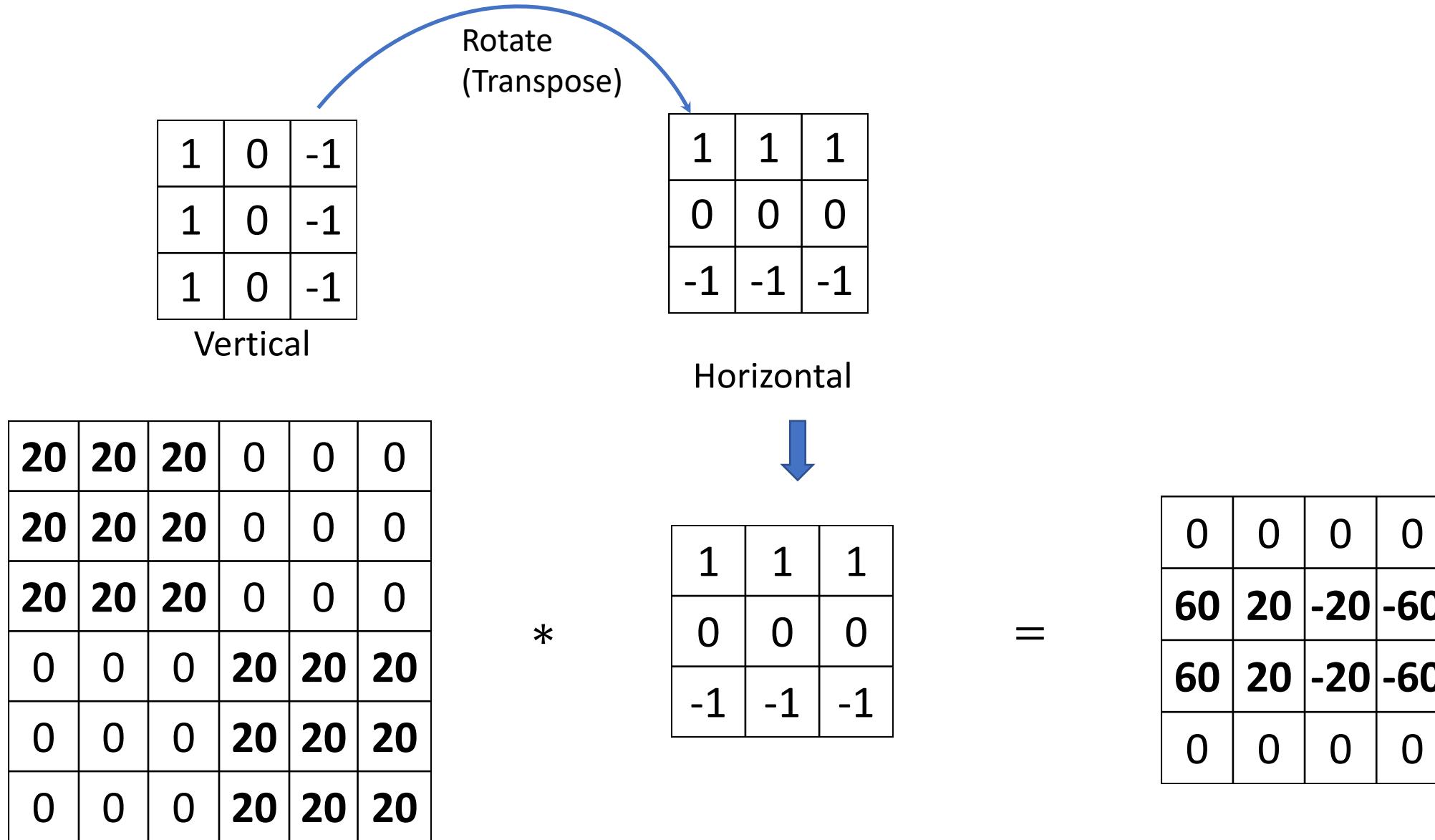
1	0	-1
1	0	-1
1	0	-1

=

0	-60	-60	0
0	-60	-60	0
0	-60	-60	0
0	-60	-60	0



Vertical vs. Horizontal Edge Detection



More Filters

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

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

Sobel filter

3	0	-3
10	0	-10
3	0	-2

Scharr filter

Which filter to pick? (Why to pick! Learn it)

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

6x6

*

Convolution

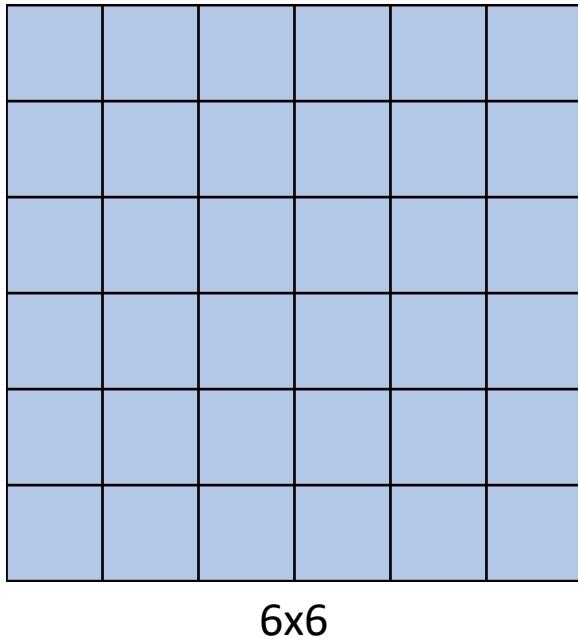
3×3

=

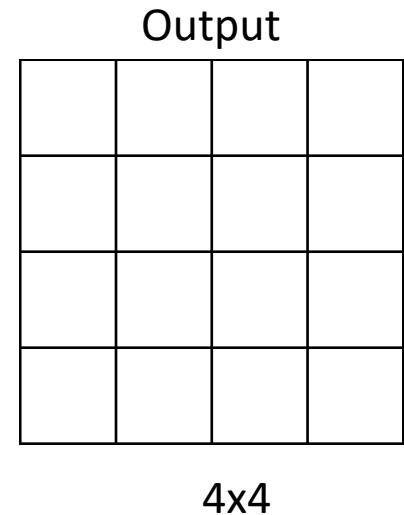
4x4

Padding

Remember the original example:

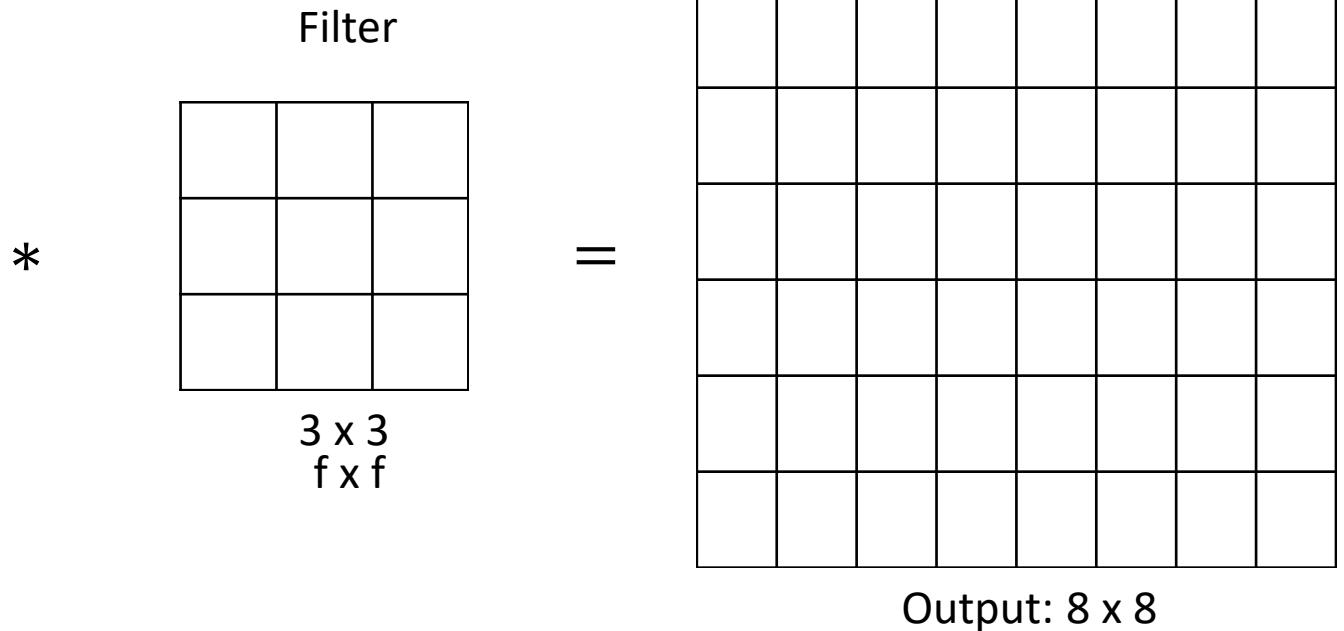
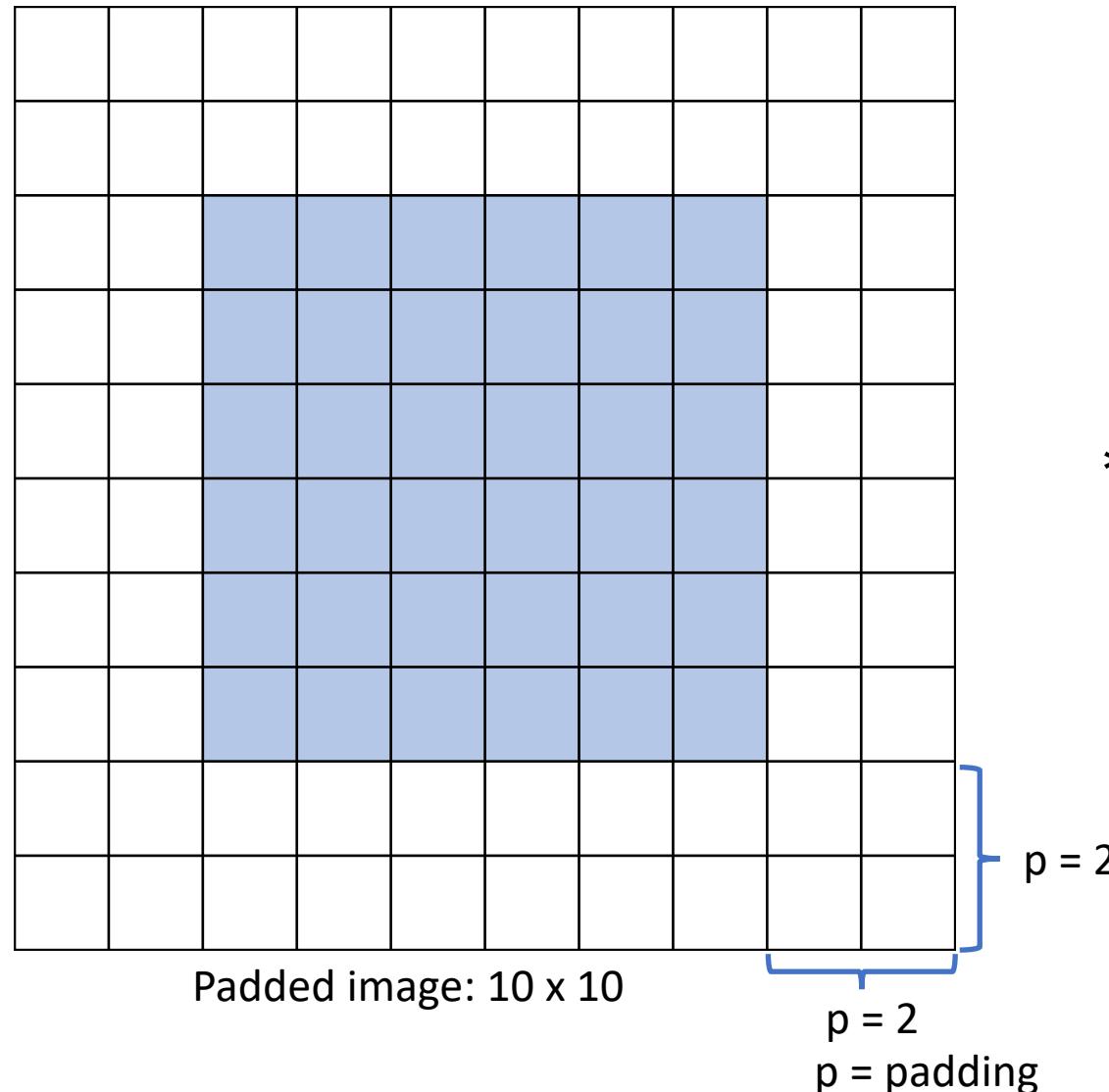


$$\begin{matrix} \text{*} & \text{Filter} \\ & \begin{matrix} \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \end{matrix} \\ & 3 \times 3 \end{matrix}$$



Padding

Original image: $6 \times 6 = n \times n$



Final output dims = $(n-f+2p+1) \times (n-f+2p+1)$

“Valid” & “Same” Convolutions

“Valid”: if no padding!

“Same”: Pad such that the output size matches the input size.

$$P = (f-1)/2$$

Strided convolution

Stride = 2

Stride = 2

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

*

3	4	4
1	0	2
-1	0	3

=

91	100	83
69	91	127
44	72	74

Output Size (Output “Shape”)

Image dims: $n \times n$

Filter dims: $f \times f$

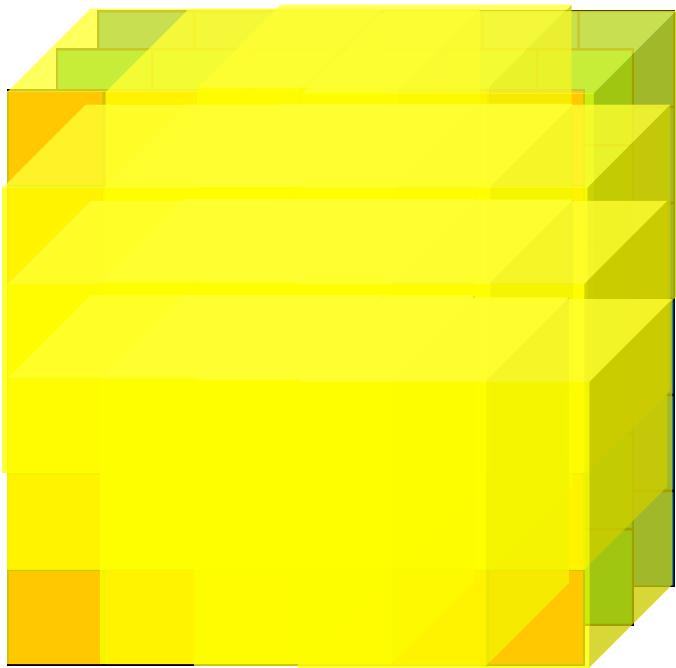
Padding: p

Stride: s

$$\text{Output dims: } \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor$$

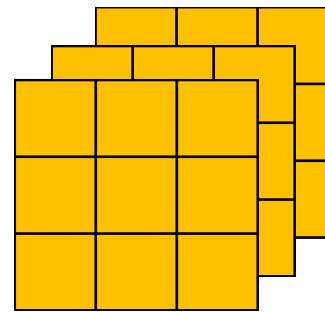
Floor (round to the lowest integer)

Convolutions on RGB image



$6 \times 6 \times 3$

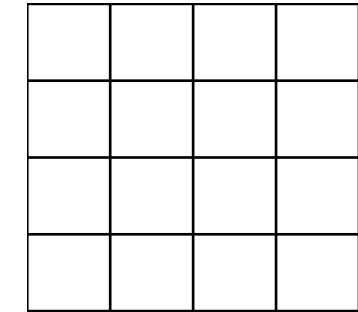
*



$3 \times 3 \times 3$

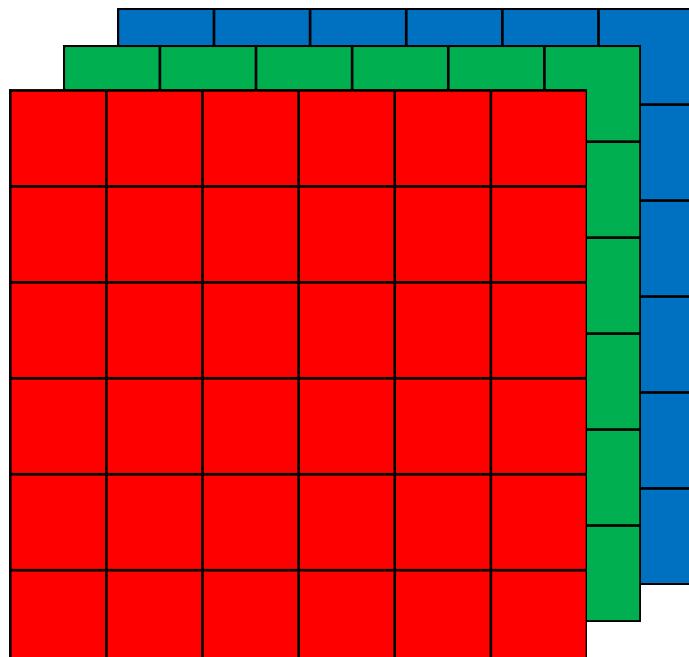


=



4×4

Multiple filters



6 x 6 x 3

Number of channels

$n \times n \times n_c$

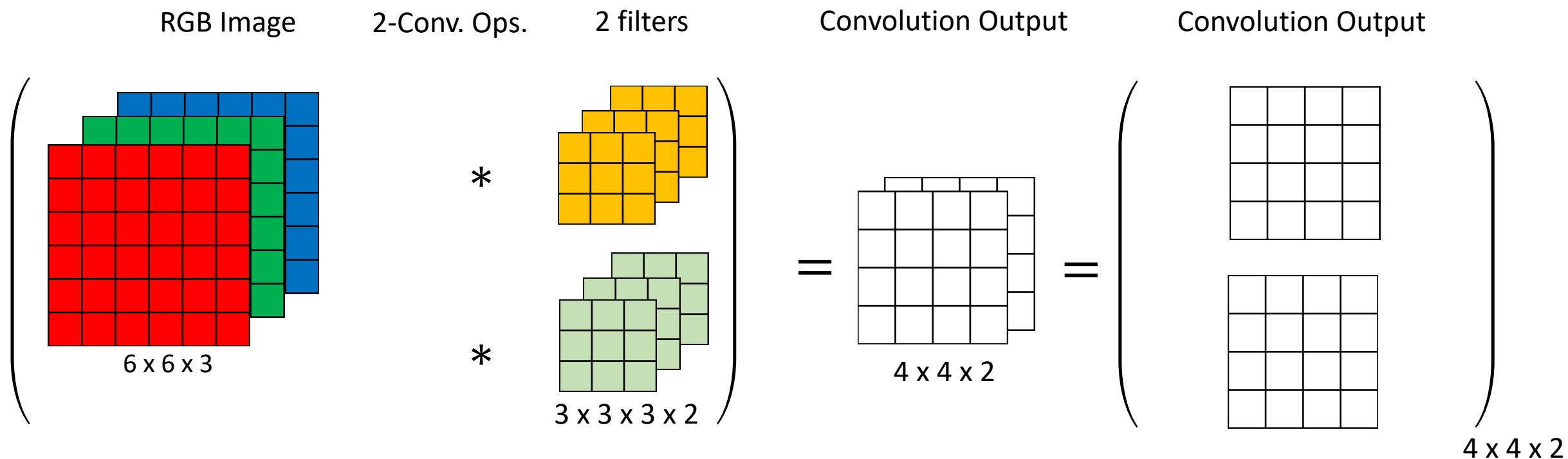
$$\begin{matrix} * & \begin{matrix} \text{3} \times 3 \times 3 \end{matrix} & = & \begin{matrix} \text{4} \times 4 \end{matrix} \\ \begin{matrix} \text{3} \times 3 \times 3 \end{matrix} & \begin{matrix} \text{4} \times 4 \end{matrix} & & \end{matrix}$$
A diagram illustrating a convolution operation. On the left, a 3x3x3 cube of yellow blocks is shown. An asterisk (*) is placed above it, followed by the text "3 x 3 x 3". To its right is an equals sign (=). To the right of the equals sign is a 4x4 grid of empty white boxes, with the text "4 x 4" below it. This represents a 3x3 input feature map being processed by a 4x4 kernel.

$$\begin{matrix} * & \begin{matrix} \text{3} \times 3 \times 3 \end{matrix} & = & \begin{matrix} 4 & 4 \end{matrix} \end{matrix}$$

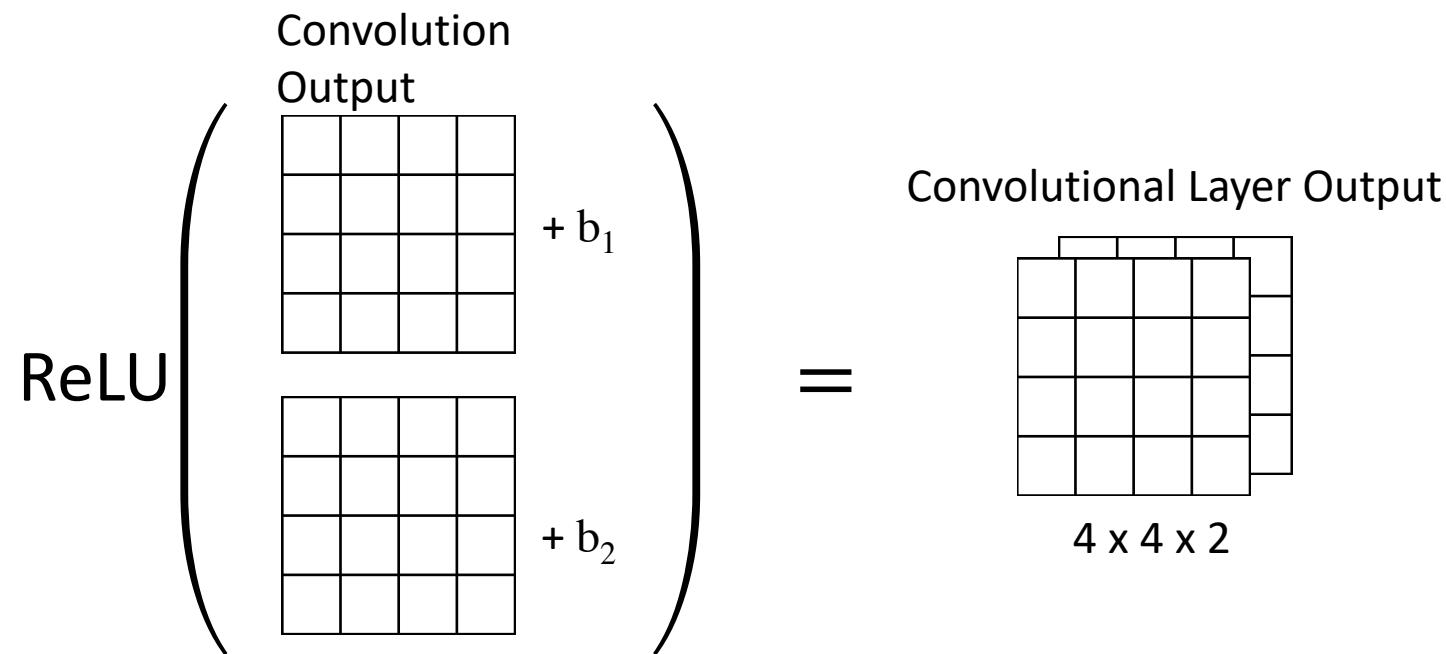
The diagram illustrates a mathematical operation involving tensors. On the left, a 3x3x3 tensor is shown as a stack of three 3x3 matrices. A black asterisk (*) is placed to its left, indicating scalar multiplication. To the right of the equals sign (=), a 4x4 matrix is shown as a stack of four 4x4 matrices. This visualizes how a 3D tensor can be converted into a 2D matrix through a scaling operation.

$$* \quad f \times f \times n_c \times 2 = (n - f + 1) \times (n - f + 1) \times 2$$

Multiple filters



A Convolutional Layer of CNN (ConvNet)



Convolutional Layer Output = $\text{ReLU} ((RGBImage * filter) + b)$

Output of a neuron in a FC NN: $a = \text{ReLU} (\mathbf{w}^T \mathbf{x} + b)$ in FC NN

Similar notation!
(but not equal)

Input – Output Dimensions for l^{th} Convolutional Layer

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



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

Padding size: $p^{[l]} \times p^{[l]}$

Stride size: $s^{[l]} \times s^{[l]}$

Filter size: $f^{[l]} \times f^{[l]} \times n_C^{[l-1]} \times n_C^{[l]} = \text{One_FilterDims} \times n_C^{[l]}$

(Number of Weights = Filter size), Bias size: $n_C^{[l]}$

Output dims (for one image) : $\left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor \times \left\lfloor \frac{n_W^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor \times n_C^{[l]}$

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

The output dimensions for m input samples $A^{[l]}$: $m \times n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

This lecture contains material from Andrew Ng and Ulas Bagci