

# Foundations of Convolutional Neural Networks

Pio

# Computer Vision Problems

## Image Classification



64x64

→ Cat? (0/1)

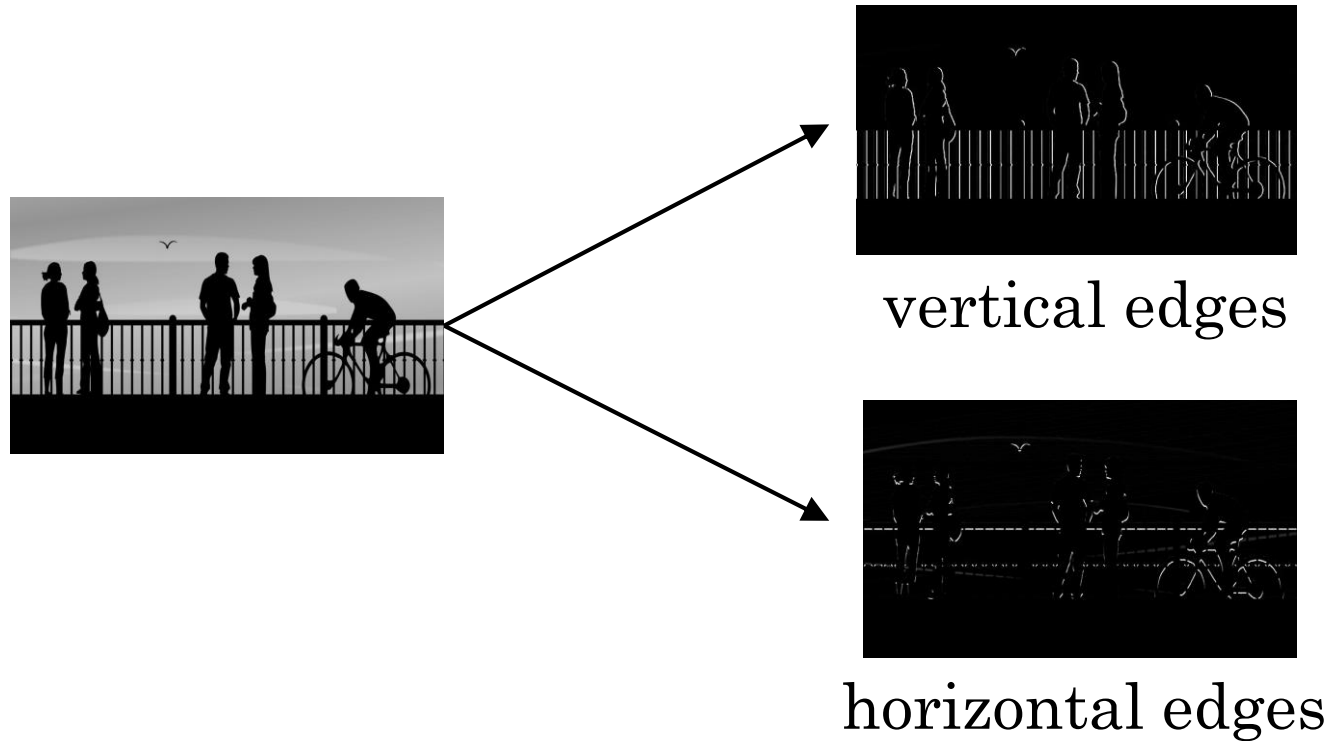
## Neural Style Transfer



## Object detection



# Computer Vision Problem



# Vertical edge detection

3 <sup>1</sup>	0 <sup>0</sup>	1 <sup>-1</sup>	2 <sup>-1</sup>	7 <sup>-0</sup>	4 <sup>-1</sup>
1 <sup>1</sup>	5 <sup>0</sup>	8 <sup>-1</sup>	9 <sup>-1</sup>	3 <sup>-0</sup>	1 <sup>-1</sup>
2 <sup>1</sup>	7 <sup>0</sup>	2 <sup>-1</sup>	5 <sup>-1</sup>	1 <sup>-0</sup>	3 <sup>-1</sup>
0 <sup>1</sup>	1 <sup>0</sup>	3 <sup>-1</sup>	1 <sup>-1</sup>	7 <sup>-0</sup>	8 <sup>-1</sup>
4	2	1	6	2	8
2	4	5	2	3	9

6 x 6

convolution  
\*

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

3 x 3

filter / kernel

=

0	-2	-4	-7
-3	-2	-3	-16

4 x 4

# Vertical edge detection

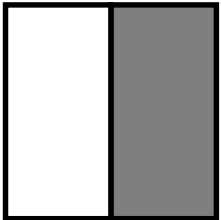
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0

\*

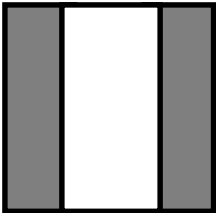
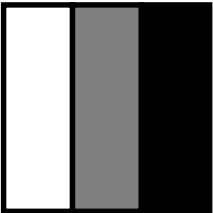
1	0	-1
1	0	-1
1	0	-1

=

0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0




\*




# Vertical edge detection examples

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




\*

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




=

0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	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	10	10	10




\*

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



=

0	-30	-30	0
0	-30	-30	0
0	-30	-30	0
0	-30	-30	0



# Vertical and Horizontal Edge Detection

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

Vertical

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

Horizontal

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

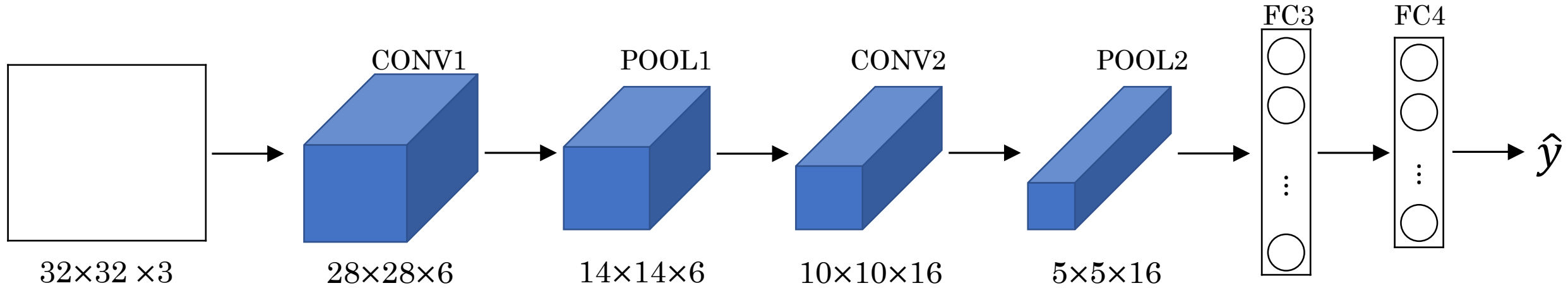
\*

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

=

0	0	0	0
30	10	-10	-30
30	10	-10	-30
0	0	0	0

# Types of layer in a convolutional network:



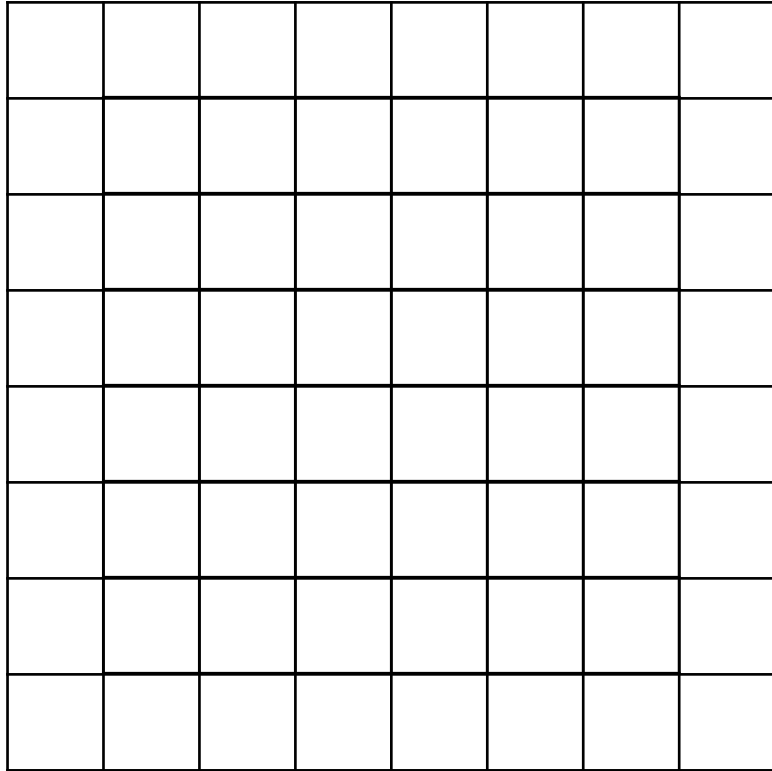
- Convolution (CONV)
- Pooling (POOL)
- Fully connected (FC)



# Convolution Parameters

- Padding
- Stride
- Multiple Filters

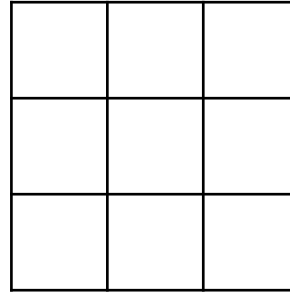
# Padding



$6 \times 6 \rightarrow 8 \times 8$

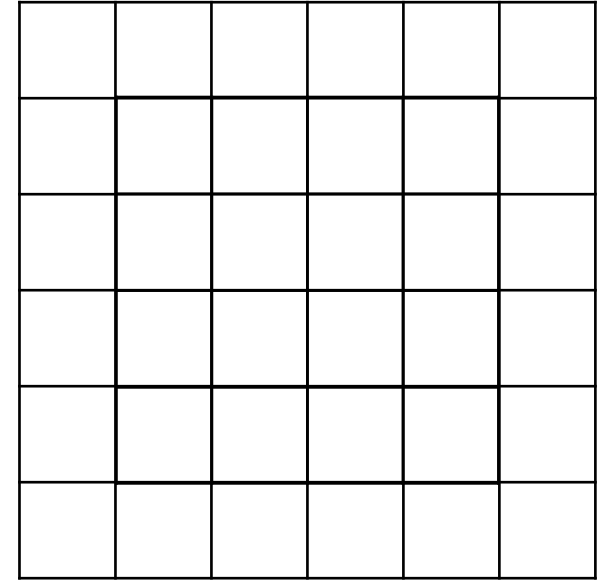
Padding  $p = 1$

\*



$3 \times 3$

=



$4 \times 4 \rightarrow 6 \times 6$

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

Valid

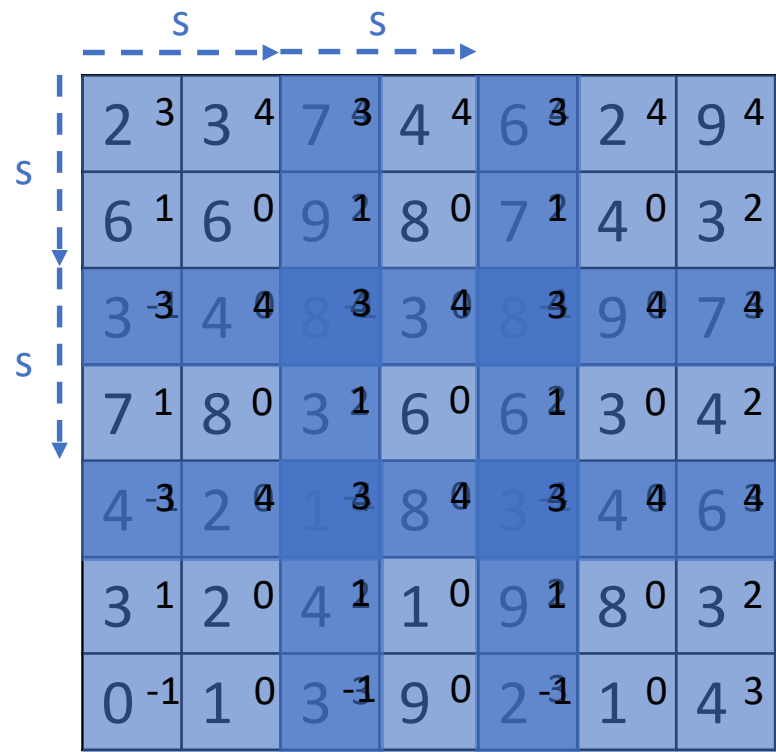
Convolution

$(n + 2p - f + 1)$   
 $\times (n + 2p - f + 1)$

Same

Convolution

# Strided convolution



$n \times n$  image  
padding  $p$   
stride  $s = 2$

\*

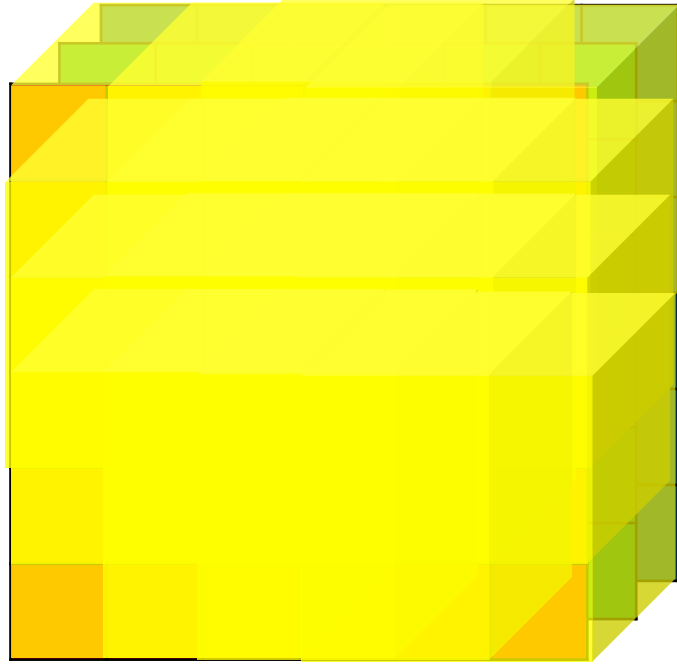
3	4	4
1	0	2
-1	0	3

$f \times f$  filter

=


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

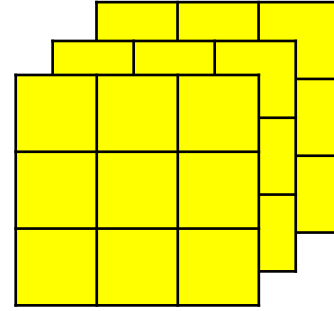
# Convolutions on RGB image



6 x 6 x 3

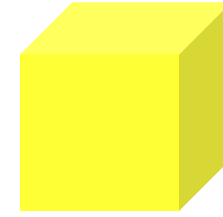
height x width x channels

\*

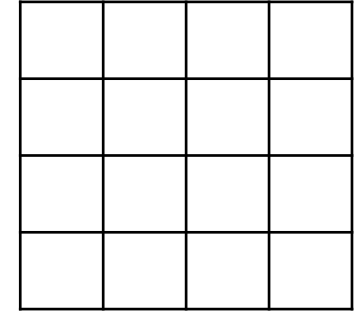


3 x 3 x 3

f x f x channels

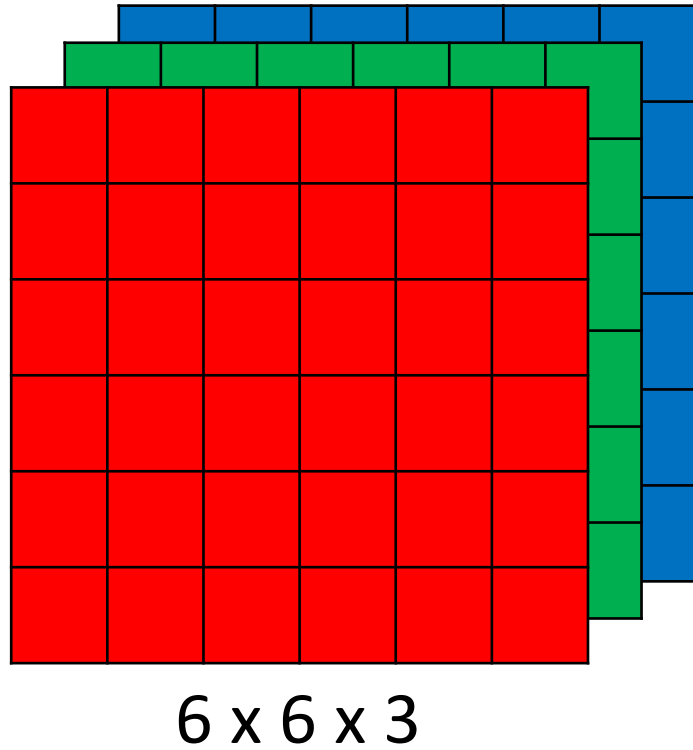


=

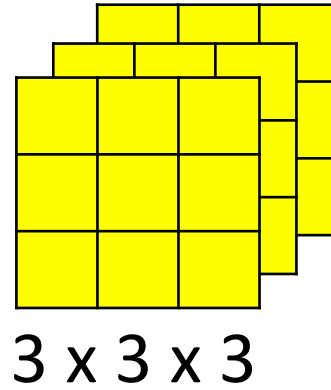


4 x 4

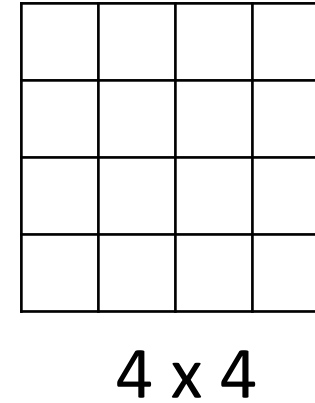
# Multiple filters



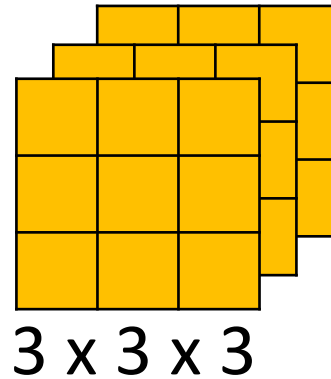
\*



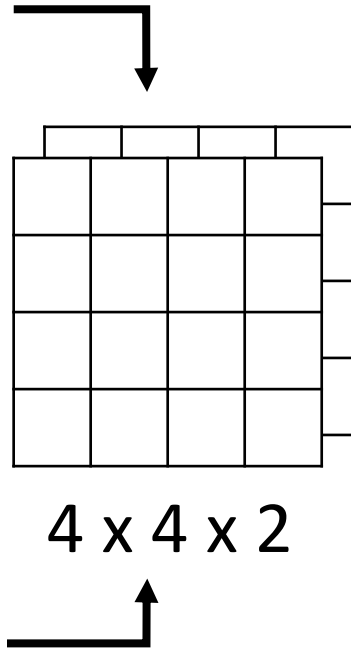
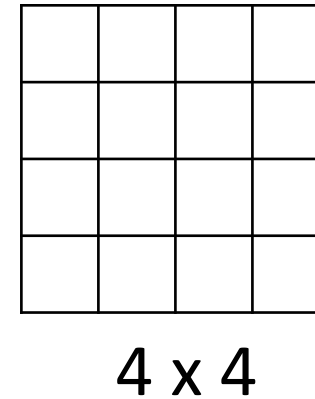
=



\*



=



# Summary of Convolutions

Hyperparameters:

$n_H$  : image height

$n_W$  : image width

$n_C$  : # of image channels (# of filters)

$f$  : filter size

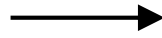
$s$  : stride

$p$  : padding

$$\left\lfloor \frac{n_H + 2p - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_W + 2p - f}{s} + 1 \right\rfloor \times n_C$$

# Pooling layer: Max pooling

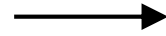
1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2



9	2
6	3

# Pooling layer: Max pooling

1	3	2	1	3
2	9	1	1	5
1	8	2	1	2
8	3	8	1	0
5	6	1	2	9






# Pooling layer: Average pooling

1	3	2	1
2	9	1	1
1	4	2	3
5	6	1	2



3.75	1.25
4	2

# Summary of pooling

Hyperparameters:

$f$  : filter size

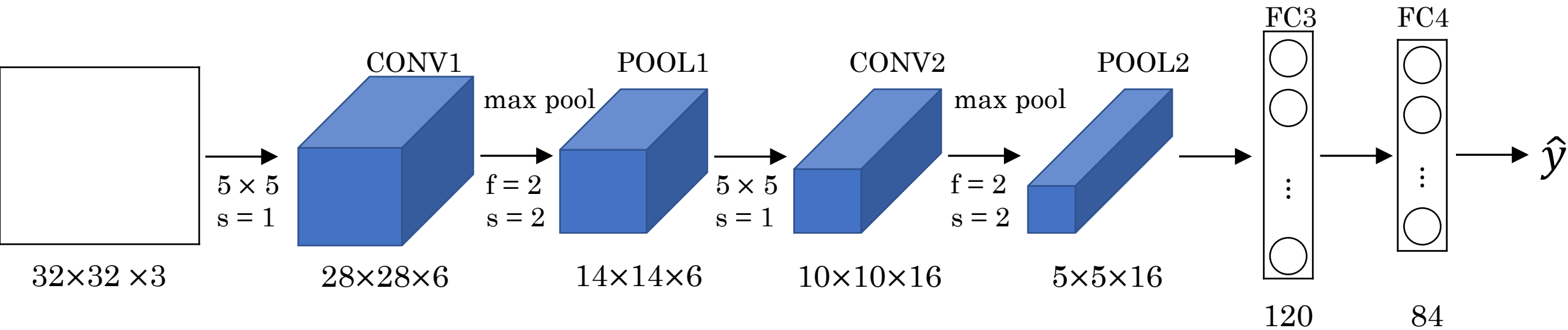
$s$  : stride

Max or average pooling

~~$p$  : padding~~

$$\left\lfloor \frac{n_H - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_W - f}{s} + 1 \right\rfloor \times n_C$$

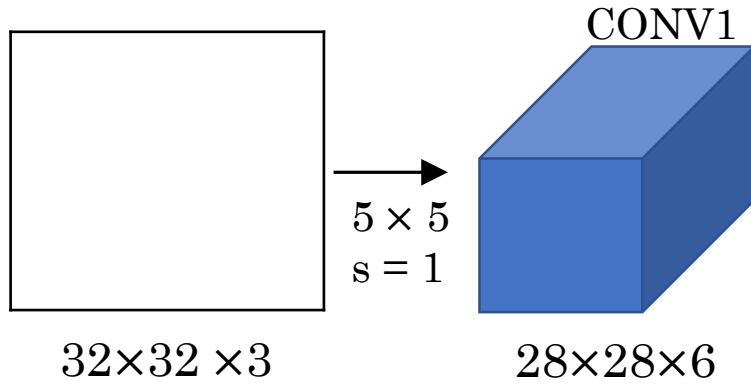
# Neural network example (LeNet – 5)



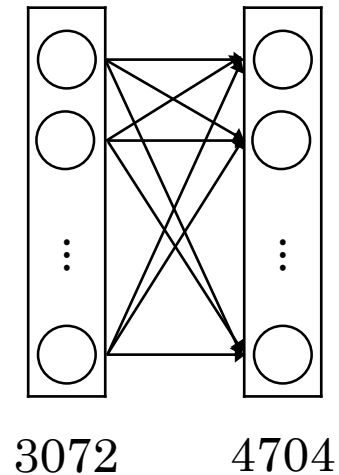
	Activation shape	Activation Size	# parameters
Input:	(32,32,3)	3,072	0
CONV1 (f=5, s=1)	(28,28,6)	4704	156
POOL1	(14,14,6)	1176	0
CONV2 (f=5, s=1)	(10,10,16)	1,600	416
POOL2	(5,5,16)	400	0
FC3	(120,1)	120	48,001
FC4	(84,1)	84	10,081
Softmax	(10,1)	10	841

$$\begin{matrix} n_H, n_W & \downarrow \\ n_C & \uparrow \end{matrix}$$

# Why convolutions



	Activation shape	Activation Size	# parameters
Input:	(32,32,3)	3,072	0
CONV1 (f=5, s=1)	(28,28,6)	4,704	156



$$3072 \times 4704 \approx 14\text{M}$$

# Why convolutions

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

 $*$ 

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

 $=$ 

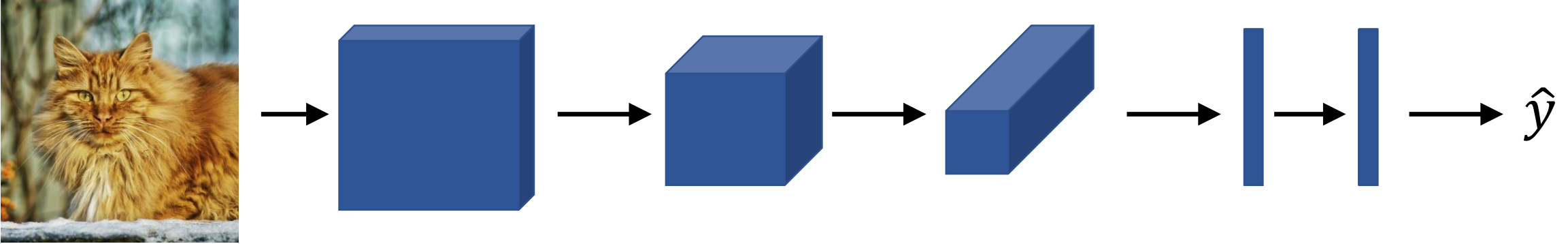
0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0

**Parameter sharing:** A feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.

**Sparsity of connections:** In each layer, each output value depends only on a small number of inputs.

# Putting it together

Training set  $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$ .



$$\text{Cost } J = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Use gradient descent to optimize parameters to reduce  $J$

# Quiz

- What do you think applying this filter to a grayscale image will do?

$$\begin{bmatrix} 0 & 1 & -1 & 0 \\ 1 & 3 & -3 & -1 \\ 1 & 3 & -3 & -1 \\ 0 & 1 & -1 & 0 \end{bmatrix}$$

- a) Detect vertical edges
- b) Detect horizontal edges
- c) Detect image contrast
- d) Detect 45 degree edges



- Because pooling layers do not have parameters, they do not affect the backpropagation (derivatives) calculation.
  - a) True
  - ☒ b) False

- In lecture we talked about “parameter sharing” as a benefit of using convolutional networks. Which of the following statements about parameter sharing in ConvNets are true? (Check all that apply.)

- ☒ It allows gradient descent to set many of the parameters to zero, thus making the connections sparse.
- ☐ It allows parameters learned for one task to be shared even for a different task (transfer learning).
- ☐ It reduces the total number of parameters, thus reducing overfitting.
- ☒ It allows a feature detector to be used in multiple locations throughout the whole input image/input volume.

- In lecture we talked about “sparsity of connections” as a benefit of using convolutional layers. What does this mean?
  - a) Regularization causes gradient descent to set many of the parameters to zero.
  - b) Each activation in the next layer depends on only a small number of activations from the previous layer.**
  - c) Each layer in a convolutional network is connected only to two other layers
  - d) Each filter is connected to every channel in the previous layer.



# Valid and Same convolutions

→ no padding

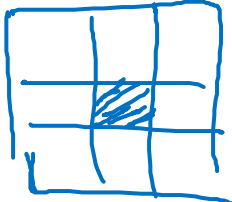
“Valid”:  $n \times n \quad * \quad f \times f \quad \rightarrow \quad \frac{n-f+1}{1} \times n-f+1$   
 $6 \times 6 \quad * \quad 3 \times 3 \quad \rightarrow \quad 4 \times 4$

“Same”: Pad so that output size is the same as the input size.

$$\begin{aligned} n+2p-f+1 \times n+2p-f+1 \\ \cancel{n+2p-f+1} = \cancel{n} \Rightarrow \boxed{p = \frac{f-1}{2}} \\ 3 \times 3 \quad p = \frac{3-1}{2} = 1 \quad \left| \begin{array}{l} 5 \times 5 \\ f=5 \end{array} \right. \end{aligned}$$

$f$  is usually odd

1x1  
3x3  
5x5  
7x7



$p=2$

# Technical note on cross-correlation vs. convolution

Convolution in math textbook:

2 <sup>7</sup>	3 <sup>2</sup>	7 <sup>5</sup>	4	6	2
6 <sup>9</sup>	6 <sup>0</sup>	9 <sup>4</sup>	8	7	4
3 <sup>-1</sup>	4 <sup>1</sup>	8 <sup>3</sup>	3	8	9
7	8	3	6	6	3
4	2	1	8	3	4
3	2	4	1	9	8

3	4	5
1	0	2
-1	9	7

7	2	5
9	0	4
-1	1	3


$$(A * B) * C = A * (B * C)$$

# Learning to detect edges

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	-3

Scharr Filter

- You have an input volume that is  $32 \times 32 \times 16$ , and apply max pooling with a stride of 2 and a filter size of 2. What is the output volume?
  - a)  $16 \times 16 \times 8$
  - b)  $15 \times 15 \times 16$
  - c)  $32 \times 32 \times 8$
  - d)  $16 \times 16 \times 16$