

# Computer Vision II

Fundamentals of Artificial Intelligence

Instructor: Chenhui Chu

Email: [chu@i.kyoto-u.ac.jp](mailto:chu@i.kyoto-u.ac.jp)

Teaching Assistant: Youyuan Lin

E-mail: [youyuan@nlp.ist.i.kyoto-u.ac.jp](mailto:youyuan@nlp.ist.i.kyoto-u.ac.jp)

# Schedule

- 1. Overview of AI and this Course (4/14)
- 2. Introduction to Python (4/21)
- 3, 4. Mathematics Concepts I, II (4/28, 5/12)
- 5, 6. Regression I, II (5/19, 5/26)
- 7. Classification (6/2)
- 8. Introduction to Neural Networks (6/9)
- 9. Neural Networks Architecture and Backpropagation (6/16)
- 10. Fully Connected Layers (6/23)
- 11, 12, 13. **Computer Vision** I, II, III (6/30, 7/7, 7/14)
- 14. Natural Language Processing (7/17)

# Overview of This Course

11, 12, 13. Computer vision I,  
II, III

14. Natural language  
processing

## Deep Learning Applications



8. Neural network  
introduction

9. Architecture and  
backpropagation

10. Feedforward  
neural networks

## Deep Learning



5. Regression I

6. Regression II

7. Classification

## Basic Supervised Machine Learning



2. Python

3, 4. Mathematics concepts I, II

## Fundamental of Machine Learning

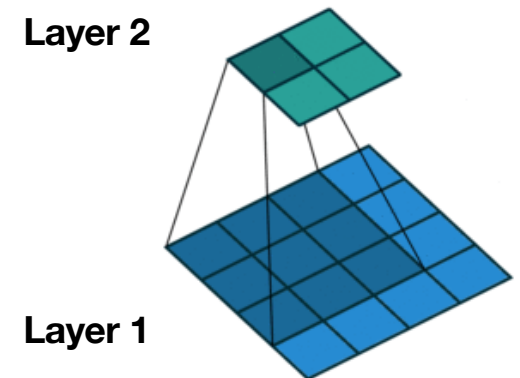
# Previously

- Last week, we looked at **convolutional layers**
- And we saw that they can be mathematically represented by an operation called the “**convolution operation**”
  - Just like Fully-Connected layers can be represented by Matrix-Multiplication
- We also saw that the “convolution operation” can be applied to images with certain kernels to produce “edge detection”

# Convolutional Layers

- Neurons are organized in 2-dimensional layers
- Neurons in 2 layers are only connected if they roughly belong to the same area of their respective layer
- Eg. The neuron in the top-left corner of layer 2 is only connected to the 9 neurons in the top-left corner of layer 1

This gives spatial information to the network

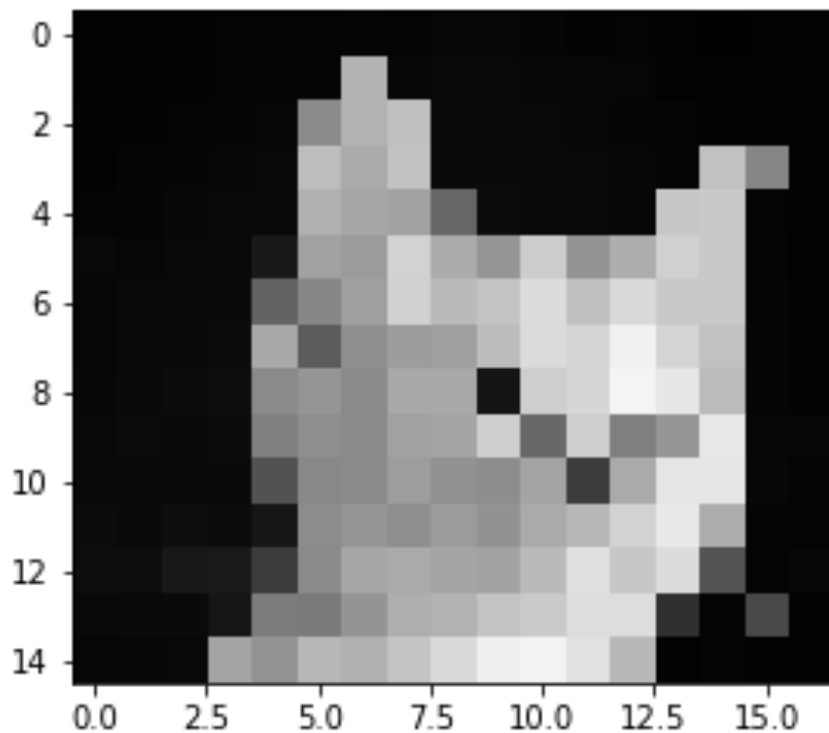


**Because all of the inputs of one neuron correspond to Neighboring pixels**

# An image as a 2D Array

- Greyscale image: each pixel has a grey value between 0 (black) and 1 (white)

Image with 18x20 pixels (greyscale)



Array with 18x20 numbers

0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0	0.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.5	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.8	0.5	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.6	0.4	0.0	0.0	0.0	0.0	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.1	0.6	0.6	0.8	0.7	0.6	0.8	0.6	0.7	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.4	0.5	0.6	0.8	0.7	0.8	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.7	0.4	0.6	0.6	0.6	0.7	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.6	0.5	0.7	0.7	0.1	0.8	0.8	1.0	0.9	0.7	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.5	0.6	0.5	0.6	0.6	0.8	0.4	0.8	0.5	0.6	0.9	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.3	0.5	0.5	0.6	0.6	0.6	0.6	0.2	0.7	0.9	0.9	0.0	0.0	0.0
0.0	0.0	0.1	0.0	0.1	0.5	0.6	0.6	0.6	0.6	0.7	0.7	0.8	0.9	0.7	0.0	0.0	0.0
0.1	0.1	0.1	0.1	0.2	0.5	0.7	0.7	0.6	0.6	0.7	0.9	0.8	0.9	0.3	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.5	0.6	0.7	0.7	0.8	0.8	0.9	0.9	0.2	0.0	0.3	0.0	0.0
0.0	0.0	0.0	0.6	0.6	0.7	0.7	0.8	0.9	0.9	0.9	0.9	0.7	0.0	0.0	0.0	0.0	0.0

# Convolution (1/4)

Input Array

0	0	0	0	0	0
0	0	1	1	0	0
0	1	1	1	0	0
0	1	1	1	1	0
0	1	1	1	1	0
0	1	1	1	1	1

\*

Kernel Array

2	0	2
0	1	0
-1	1	0

=

Output Array

1	1	1	-1
4	3	3	2
4	5	3	3
4	5	5	3

# Convolution (2/4)

- How we compute:

**Input Array**

0	0	0	0	0	0
0	0	1	1	0	0
0	1	1	1	0	0
0	1	1	1	1	0
0	1	1	1	1	0
0	1	1	1	1	1

**Kernel Array**

2	0	2
0	1	0
-1	1	0

=

**Output Array**

1	1	1	-1
4	3	3	2
4	5	3	3
4	5	5	3

$0 \times 2 + 0 \times 0 + 0 \times 2 + 0 \times 0 + 0 \times 1 + 1 \times 0 + 0 \times -1 + 1 \times 1 + 1 \times 0 = 1$



# Convolution (3/4)

- How we compute:

**Input Array**

0	0	0	0	0	0
0	0	1	1	0	0
0	1	1	1	0	0
0	1	1	1	1	0
0	1	1	1	1	0
0	1	1	1	1	1

**Kernel Array**

2	0	2
0	1	0
-1	1	0

**Output Array**

1	1	1	-1
4	3	3	2
4	5	3	3
4	5	5	3

$0 \times 2 + 0 \times 0 + 1 \times 2 + 0 \times 0 + 1 \times 1 + 1 \times 0 + 0 \times -1 + 1 \times 1 + 1 \times 0 = 4$

# Convolution (4/4)

- How we compute:

**Input Array**

0	0	0	0	0	0
0	0	1	1	0	0
0	1	1	1	0	0
0	1	1	1	1	0
0	1	1	1	1	0
0	1	1	1	1	1

**Kernel Array**

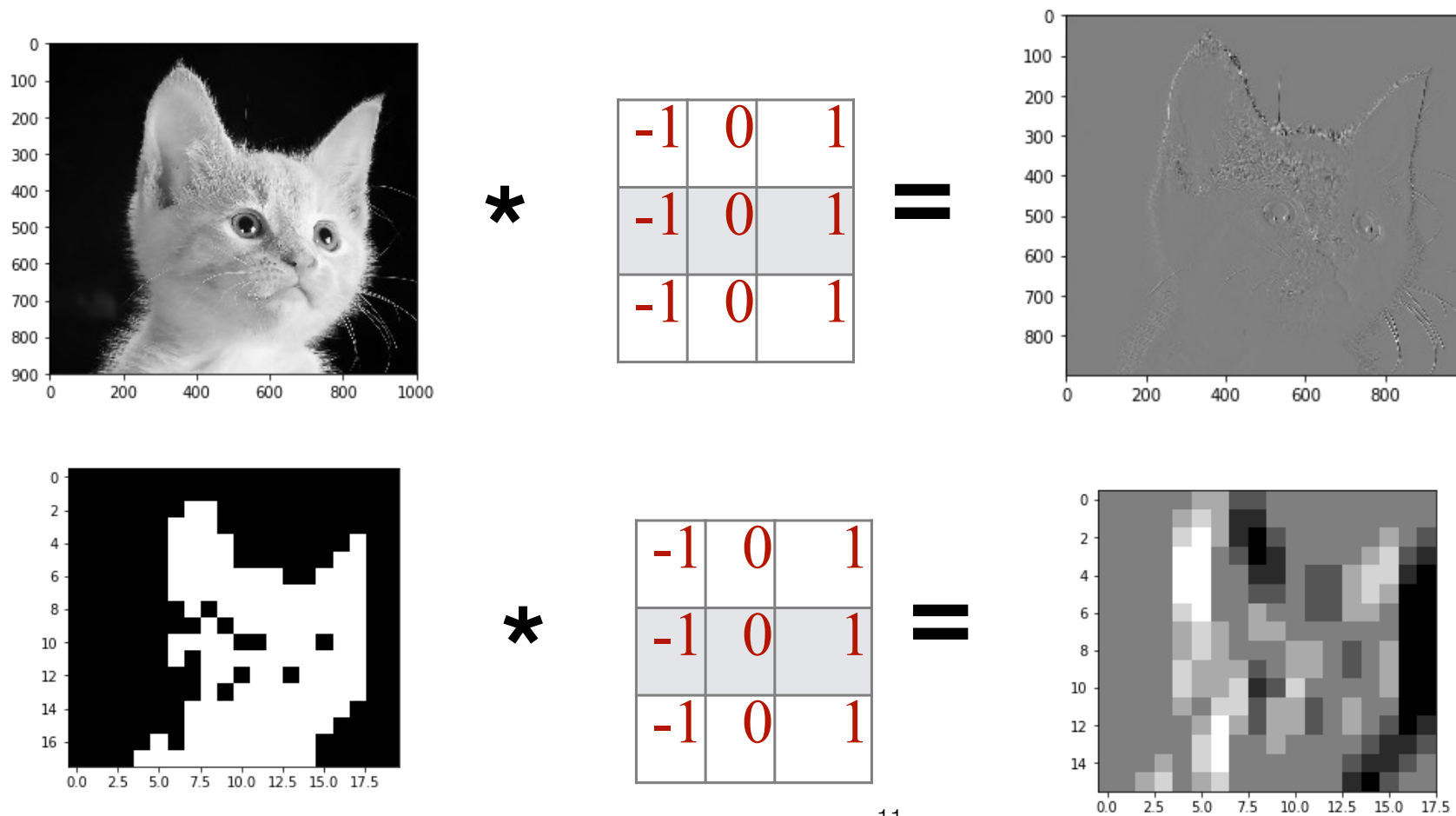
2	0	2
0	1	0
-1	1	0

**Output Array**

1	1	1	-1
4	3	3	2
4	5	3	3
4	5	5	3

$$0 \times 2 + 1 \times 0 + 1 \times 2 + 1 \times 0 + 1 \times 1 + 1 \times 0 + 1 \times -1 + 1 \times 1 + 1 \times 0 = 3$$

# Edge Detectors



# Today

- We:
  - Consider “***volume***” ***convolutions*** instead of the “flat” convolutions we just described
  - See one last type of layers: “***Max-Pooling***” ***layers***
  - Combine everything to create an ***Image Classifier***

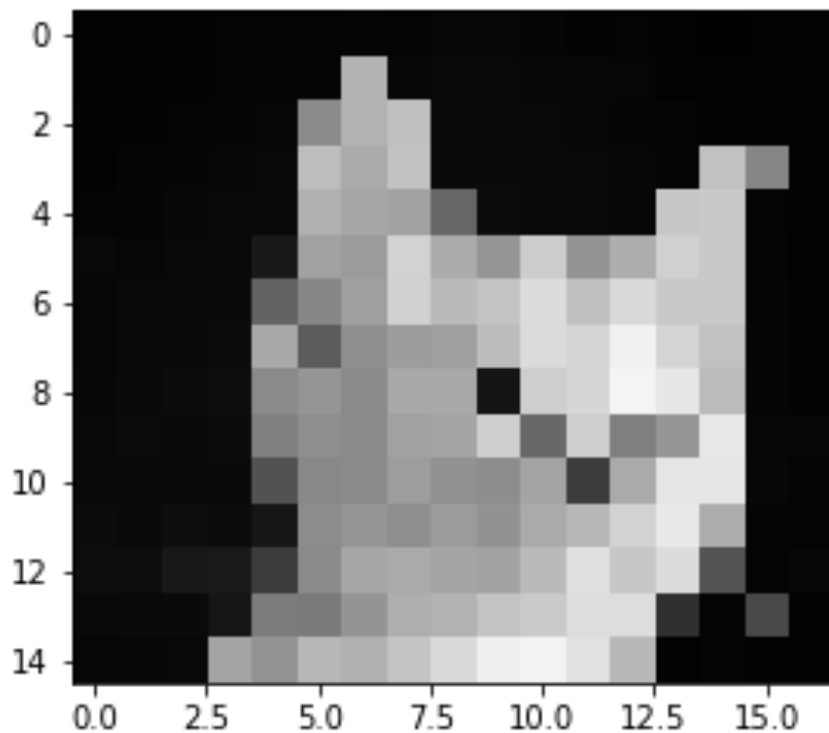
# Volume Convolution

- When we discussed convolution, we considered the input was a 2D array of numbers
  - This 2D array corresponds for example, to a Black & White image

# An Image as a 2D Array

- Greyscale image: each pixel has a grey value between 0 (black) and 1 (white)

Image with 18x20 pixels (greyscale)



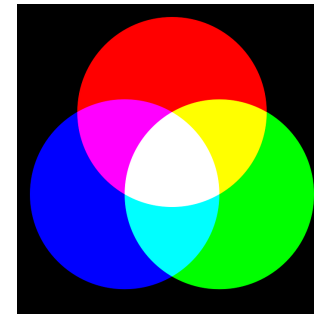
Array with 18x20 numbers

0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0	0.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.5	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.8	0.5	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.6	0.4	0.0	0.0	0.0	0.0	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.1	0.6	0.6	0.8	0.7	0.6	0.8	0.6	0.7	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.4	0.5	0.6	0.8	0.7	0.8	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.7	0.4	0.6	0.6	0.6	0.7	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.6	0.5	0.7	0.7	0.1	0.8	0.8	1.0	0.9	0.7	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.5	0.6	0.5	0.6	0.6	0.8	0.4	0.8	0.5	0.6	0.9	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.3	0.5	0.5	0.6	0.6	0.6	0.6	0.2	0.7	0.9	0.9	0.0	0.0	0.0
0.0	0.0	0.1	0.0	0.1	0.5	0.6	0.6	0.6	0.6	0.7	0.7	0.8	0.9	0.7	0.0	0.0	0.0
0.1	0.1	0.1	0.1	0.2	0.5	0.7	0.7	0.6	0.6	0.7	0.9	0.8	0.9	0.3	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.5	0.6	0.7	0.7	0.8	0.8	0.9	0.9	0.2	0.0	0.3	0.0	0.0
0.0	0.0	0.0	0.6	0.6	0.7	0.7	0.8	0.9	0.9	0.9	0.9	0.7	0.0	0.0	0.0	0.0	0.0

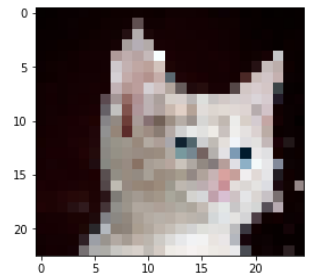
# Volume Convolution

- When we discussed convolution, we considered the input was a 2D array of numbers
  - This 2D array corresponds for example, to a Black & White image
- What about color images?
  - Color images can be represented by 3D arrays

# Color Images (1/5)

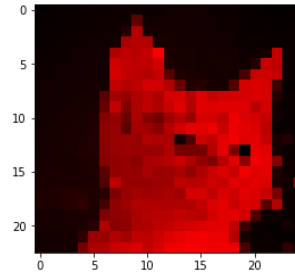


Superpose all channels



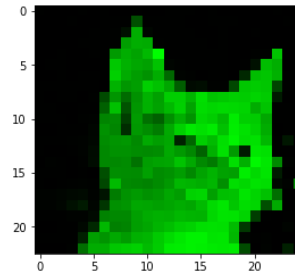
**R channel (red): 18x20 array**

0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0	0.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.5	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.8	0.5	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.6	0.4	0.0	0.0	0.0	0.0	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.1	0.6	0.6	0.8	0.7	0.6	0.8	0.6	0.7	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.4	0.5	0.6	0.8	0.7	0.8	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.7	0.4	0.6	0.6	0.6	0.7	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.6	0.5	0.7	0.7	0.1	0.8	0.8	1.0	0.9	0.7	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.5	0.6	0.5	0.6	0.6	0.8	0.4	0.8	0.5	0.6	0.9	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.3	0.5	0.5	0.6	0.6	0.6	0.6	0.2	0.7	0.9	0.9	0.0	0.0	0.0
0.0	0.0	0.1	0.0	0.1	0.5	0.6	0.6	0.6	0.7	0.7	0.8	0.9	0.7	0.0	0.0	0.0	0.0
0.1	0.1	0.1	0.1	0.2	0.5	0.7	0.7	0.6	0.6	0.7	0.9	0.8	0.9	0.3	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.5	0.6	0.7	0.7	0.8	0.8	0.9	0.9	0.2	0.0	0.3	0.0	0.0
0.0	0.0	0.0	0.6	0.6	0.7	0.7	0.8	0.9	0.9	0.9	0.9	0.7	0.0	0.0	0.0	0.0	0.0



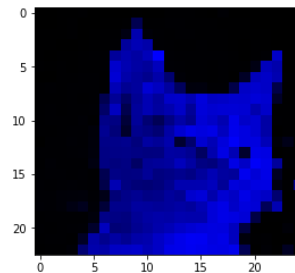
**G channel (green): 18x20 array**

0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0	0.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.5	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.8	0.5	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.6	0.4	0.0	0.0	0.0	0.0	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.1	0.6	0.6	0.8	0.7	0.6	0.8	0.6	0.7	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.4	0.5	0.6	0.8	0.7	0.8	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.7	0.4	0.6	0.6	0.6	0.7	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.6	0.5	0.7	0.7	0.1	0.8	0.8	1.0	0.9	0.7	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.5	0.6	0.5	0.6	0.6	0.8	0.4	0.8	0.5	0.6	0.9	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.3	0.5	0.5	0.6	0.6	0.6	0.6	0.2	0.7	0.9	0.9	0.0	0.0	0.0
0.0	0.0	0.1	0.0	0.1	0.5	0.6	0.6	0.6	0.6	0.7	0.7	0.8	0.9	0.7	0.0	0.0	0.0
0.1	0.1	0.1	0.1	0.2	0.5	0.7	0.7	0.6	0.6	0.7	0.9	0.8	0.9	0.3	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.5	0.6	0.7	0.7	0.8	0.8	0.9	0.9	0.2	0.0	0.3	0.0	0.0
0.0	0.0	0.0	0.6	0.6	0.7	0.7	0.8	0.9	0.9	0.9	0.9	0.7	0.0	0.0	0.0	0.0	0.0



**B channel (blue): 18x20 array**

0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0	0.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.5	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.8	0.5	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.6	0.4	0.0	0.0	0.0	0.0	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.1	0.6	0.6	0.8	0.7	0.6	0.8	0.6	0.7	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.4	0.5	0.6	0.8	0.7	0.8	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.7	0.4	0.6	0.6	0.6	0.7	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.6	0.5	0.7	0.7	0.1	0.8	0.8	1.0	0.9	0.7	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.5	0.6	0.5	0.6	0.6	0.8	0.4	0.8	0.5	0.6	0.9	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.3	0.5	0.5	0.6	0.6	0.6	0.6	0.2	0.7	0.9	0.9	0.0	0.0	0.0
0.0	0.0	0.1	0.0	0.1	0.5	0.6	0.6	0.6	0.6	0.7	0.7	0.8	0.9	0.7	0.0	0.0	0.0
0.1	0.1	0.1	0.1	0.2	0.5	0.7	0.7	0.6	0.6	0.7	0.9	0.8	0.9	0.3	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.5	0.6	0.7	0.7	0.8	0.8	0.9	0.9	0.2	0.0	0.3	0.0	0.0
0.0	0.0	0.0	0.6	0.6	0.7	0.7	0.8	0.9	0.9	0.9	0.9	0.7	0.0	0.0	0.0	0.0	0.0



In a computer, a color image is usually represented in the RGB format

We separate the color of the images in Red, Green and Blue Components

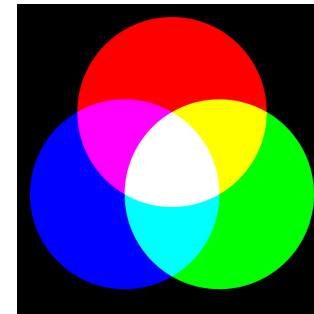
We have then 3 images that can each be represented by a 2D array



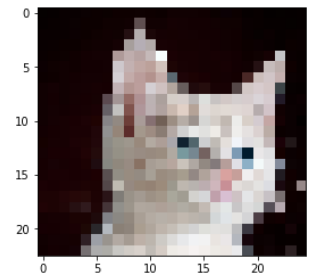
# Color Images (2/5)

- In a computer, a color image is usually represented in the **RGB format**
- We separate the color of the images in Red, Green and Blue Components
- We have then 3 images that can each be represented by a 2D array
  - One color image = three 2D arrays

# Color Images (3/5)

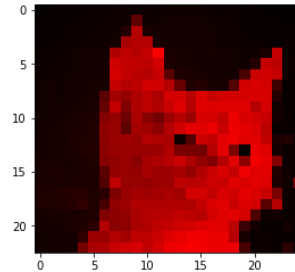


Superpose all channels



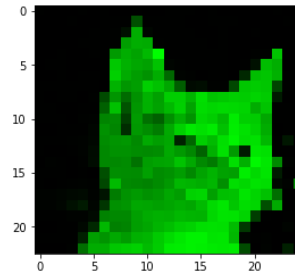
**R channel (red): 18x20 array**

0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0	0.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.5	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.8	0.5	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.6	0.4	0.0	0.0	0.0	0.0	0.0	0.8	0.8	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.1	0.6	0.6	0.8	0.7	0.6	0.8	0.6	0.7	0.8	0.8	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.4	0.5	0.6	0.8	0.7	0.8	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.7	0.4	0.6	0.6	0.6	0.7	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.6	0.5	0.7	0.7	0.1	0.8	0.8	1.0	0.9	0.7	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.5	0.6	0.5	0.6	0.6	0.8	0.4	0.8	0.5	0.6	0.9	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.3	0.5	0.5	0.6	0.6	0.6	0.6	0.2	0.7	0.9	0.9	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.1	0.0	0.1	0.5	0.6	0.6	0.6	0.7	0.7	0.8	0.9	0.7	0.0	0.0	0.0	0.0	0.0	0.0
0.1	0.1	0.1	0.1	0.2	0.5	0.7	0.7	0.6	0.6	0.7	0.9	0.8	0.9	0.3	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.5	0.6	0.7	0.7	0.8	0.8	0.9	0.9	0.2	0.0	0.3	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.6	0.6	0.7	0.7	0.8	0.9	0.9	0.9	0.9	0.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0



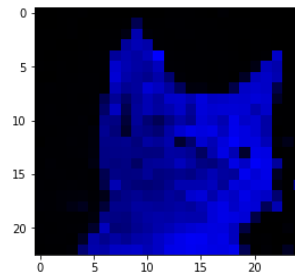
**G channel (green): 18x20 array**

0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0	0.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.5	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.8	0.5	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.6	0.4	0.0	0.0	0.0	0.0	0.0	0.8	0.8	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.1	0.6	0.6	0.8	0.7	0.6	0.8	0.6	0.7	0.8	0.8	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.4	0.5	0.6	0.8	0.7	0.8	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.7	0.4	0.6	0.6	0.6	0.7	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.6	0.5	0.7	0.7	0.1	0.8	0.8	1.0	0.9	0.7	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.5	0.6	0.5	0.6	0.6	0.8	0.4	0.8	0.5	0.6	0.9	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.3	0.5	0.5	0.6	0.6	0.6	0.6	0.2	0.7	0.9	0.9	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.1	0.0	0.1	0.5	0.6	0.6	0.6	0.6	0.7	0.7	0.8	0.9	0.7	0.0	0.0	0.0	0.0	0.0
0.1	0.1	0.1	0.1	0.2	0.5	0.7	0.7	0.6	0.6	0.7	0.9	0.8	0.9	0.3	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.5	0.6	0.7	0.7	0.8	0.8	0.9	0.9	0.2	0.0	0.3	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.6	0.6	0.7	0.7	0.8	0.9	0.9	0.9	0.9	0.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0



**B channel (blue): 18x20 array**

0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0	0.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.5	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.8	0.5	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.7	0.7	0.6	0.4	0.0	0.0	0.0	0.0	0.0	0.8	0.8	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.1	0.6	0.6	0.8	0.7	0.6	0.8	0.6	0.7	0.8	0.8	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.4	0.5	0.6	0.8	0.7	0.8	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.7	0.4	0.6	0.6	0.6	0.7	0.9	0.8	0.9	0.8	0.8	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.6	0.5	0.7	0.7	0.1	0.8	0.8	1.0	0.9	0.7	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.5	0.6	0.5	0.6	0.6	0.8	0.4	0.8	0.5	0.6	0.9	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.3	0.5	0.5	0.6	0.6	0.6	0.6	0.2	0.7	0.9	0.9	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.1	0.0	0.1	0.5	0.6	0.6	0.6	0.6	0.7	0.7	0.8	0.9	0.7	0.0	0.0	0.0	0.0	0.0
0.1	0.1	0.1	0.1	0.2	0.5	0.7	0.7	0.6	0.6	0.7	0.9	0.8	0.9	0.3	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.1	0.5	0.5	0.6	0.7	0.7	0.8	0.8	0.9	0.9	0.2	0.0	0.3	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.6	0.6	0.7	0.7	0.8	0.9	0.9	0.9	0.9	0.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0



In a computer, a color image is usually represented in the RGB format

We separate the color of the images in Red, Green and Blue Components

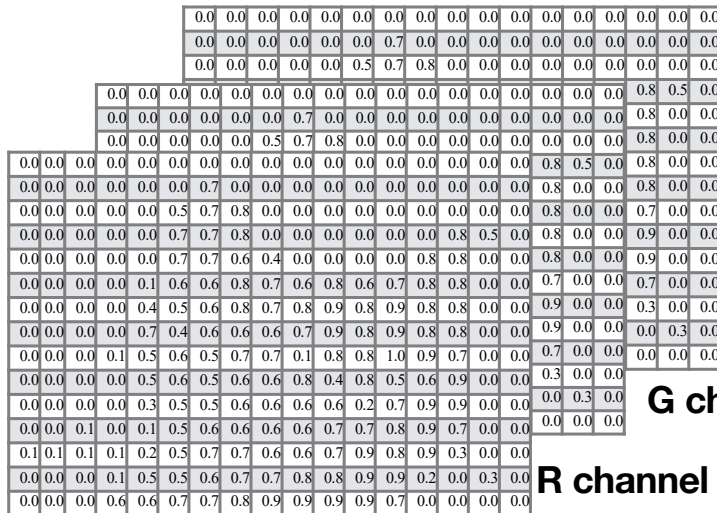
We have then 3 images that can each be represented by a 2D array

# Color Images (4/5)

- In a computer, a color image is usually represented in the ***RGB format***
- We separate the color of the images in Red, Green and Blue Components
- We have then 3 images that can each be represented by a 2D array
  - One color image = three 2D arrays
  - One color image = one 3D array

# Color Images (5/5)

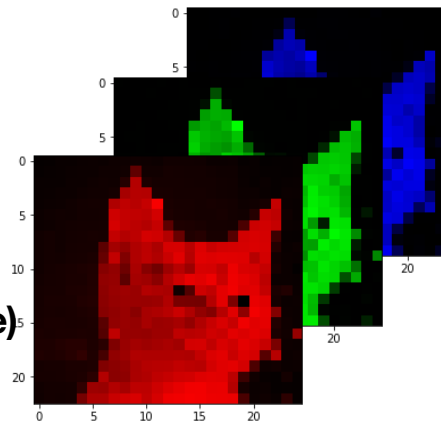
## An 18x20x3 array



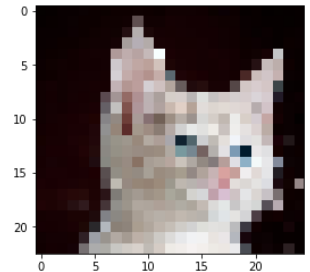
**B channel (blue)**

### G channel (green)

### R channel (red)



Superpose all channels



**In a computer, a color image is usually represented in the RGB format**

## We separate the color of the images in Red, Green and Blue Components

**We have then 3 images that can each be represented by a 2D array**

# Volume Convolution

- When we discussed convolution, we considered the input was a 2D array of numbers
  - A 2D array corresponds for example, to a Black&White image
- What about color images?
  - Color images can be represented by 3D arrays
  - In Image processing, it is a convention to call the 3rd dimension the ***“channel”*** dimension
- How do we apply convolutions to 3D arrays?

# “Flat” Convolution

- How we compute:

**Input Array**

0	0	0	0	0	0
0	0	1	1	0	0
0	1	1	1	0	0
0	1	1	1	1	0
0	1	1	1	1	0
0	1	1	1	1	1

**Kernel Array**

2	0	2
0	1	0
-1	1	0

**Output Array**

1	1	1	-1
4	3	3	2
4	5	3	3
4	5	5	3

$$0 \times 2 + 1 \times 0 + 1 \times 2 + 1 \times 0 + 1 \times 1 + 1 \times 0 + 1 \times -1 + 1 \times 1 + 1 \times 0 = 3$$

# Volume Convolution (1/9)

- We now consider we have a Kernel Array whose 3rd dimension is the same as the input
- Computation is the same as for 2D input, but we sum across channels
- Result is a 2D array

**3D Input Array**

0	0	0	0	1	0	0	0
0	0	2	0	2	0	0	0
0	0	2	3	0	0	0	0
0	0	1	1	0	0	1	3
0	1	1	1	0	0	0	0
0	1	1	4	1	0	2	1
0	1	1	1	1	0	1	Channel 3
0	1	1	1	1	1	1	Channel 2
0	1	1	1	1	1	1	Channel 1

**3D Kernel Array**

2	0	2
6	1	4
1	2	-1
0	1	2
-1	1	3

**\***

**=**

**2D Output Array**

1	1	1	-1
4	3	3	2
4	5	3	3
4	5	5	3

$$0 \times 1 + 2 \times 2 + 3 \times -1 + 0 \times 0 + 1 \times 1 + 1 \times 2 + 1 \times -1 + 1 \times 1 + 1 \times 3 + 2 \times 6 + 0 \times 1 + 2 \times 4 + \dots = 1$$

# Volume Convolution (2/9)

- If we apply a 3D Kernel to a 3D input, we get a 2D array
- Can we get a 3D output array?
  - Yes, by using more than one kernel



# Volume Convolution (3/9)

- Also, we can convolve the input with more than one kernel at a time to produce a 3D output with more than one channel

3D Kernel Array X2

**3D Input Array**

		0	0	0	1	0	0	
	0	2	0	2	0	0	0	
0	0	2	3	0	0	0	0	
0	0	1	1	0	0	0	0	
0	1	1	1	0	0	0	3	
0	1	1	4	1	0	0	0	
0	1	1	1	1	0	0	1	Channel 3
0	1	1	1	1	0	0	1	Channel 2
0	1	1	1	1	1	0	1	Channel 1

**\***

			-1	1			1	
		2	3			1	2	
	-1	1			2	0	1	
	0	3			1	1		
	2	0			-1			
		2	0			2		
		6	1			4	0	
	1	2			-1	0	0	
	0	1			2	3		
	-1	1			3			

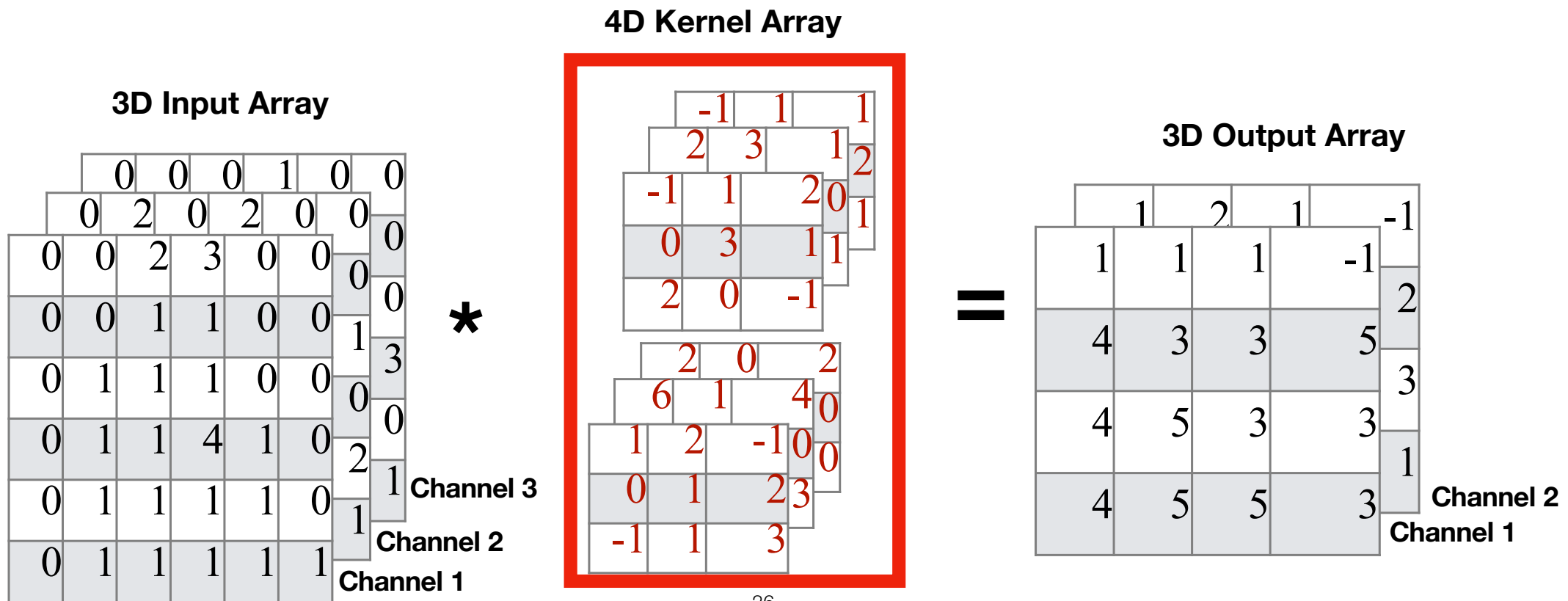
=

**3D Output Array**

		1	2	1			-1	
	1	1		1			-1	
	4	3		3			5	
	4	5		3			3	
	4	5		5			3	
								Channel 2
								Channel 1

# Volume Convolution (4/9)

- Also, we can convolve the input with more than one kernel at a time to produce a 3D output with more than one channel



# Volume Convolution (5/9)

- The values of the 4D kernel is what we are going to learn when we train our Convolutional Network

Learnable 4D Kernel

**3D Input Array**

		0	0	0	1	0	0	
	0	2	0	2	0	0	0	
0	0	2	3	0	0	0	0	
0	0	1	1	0	0	0	0	
0	1	1	1	0	0	0	0	
0	1	1	4	1	0	0	0	
0	1	1	1	1	0	0	0	
0	1	1	1	1	1	0	0	

Channel 3  
Channel 2  
Channel 1

\*

			$\Theta_{19}$	$\Theta_{20}$	$\Theta_{21}$			
		$\Theta_{10}$	$\Theta$	$\Theta_{12}$				
	$\Theta_1$	$\Theta_2$	$\Theta_3$	15	7			
	$\Theta_4$	$\Theta_5$	$\Theta_6$	18				
	$\Theta_7$	$\Theta_8$	$\Theta_9$					
			$\Theta_{44}$	$\Theta_{45}$	$\Theta_{46}$			
		$\Theta_{37}$	$\Theta$	$\Theta_{39}$	9			
	$\Theta_{28}$	$\Theta_{29}$	$\Theta_{30}$	2	1			
	$\Theta_{31}$	$\Theta_{32}$	$\Theta_{33}$	13				
	$\Theta_{34}$	$\Theta_{35}$	$\Theta_{36}$					

=

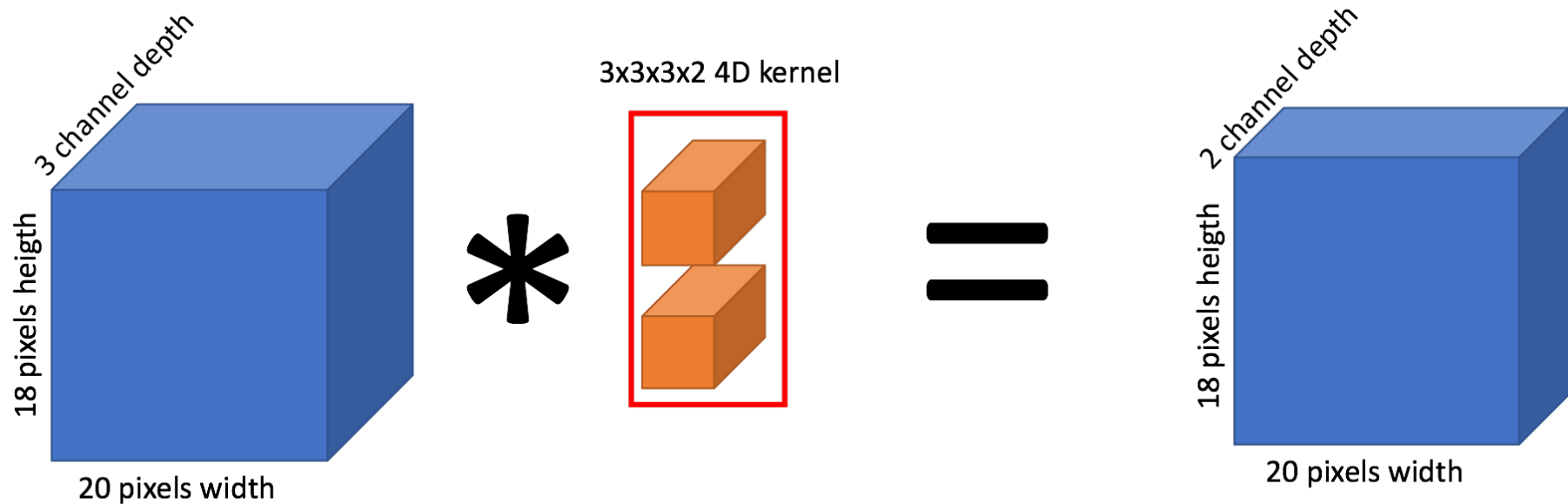
**3D Output Array**

		1	2	1				
	1	1	1		-1			
4	3	3		5				
4	5	3		3				
4	5	5		3				

Channel 2  
Channel 1

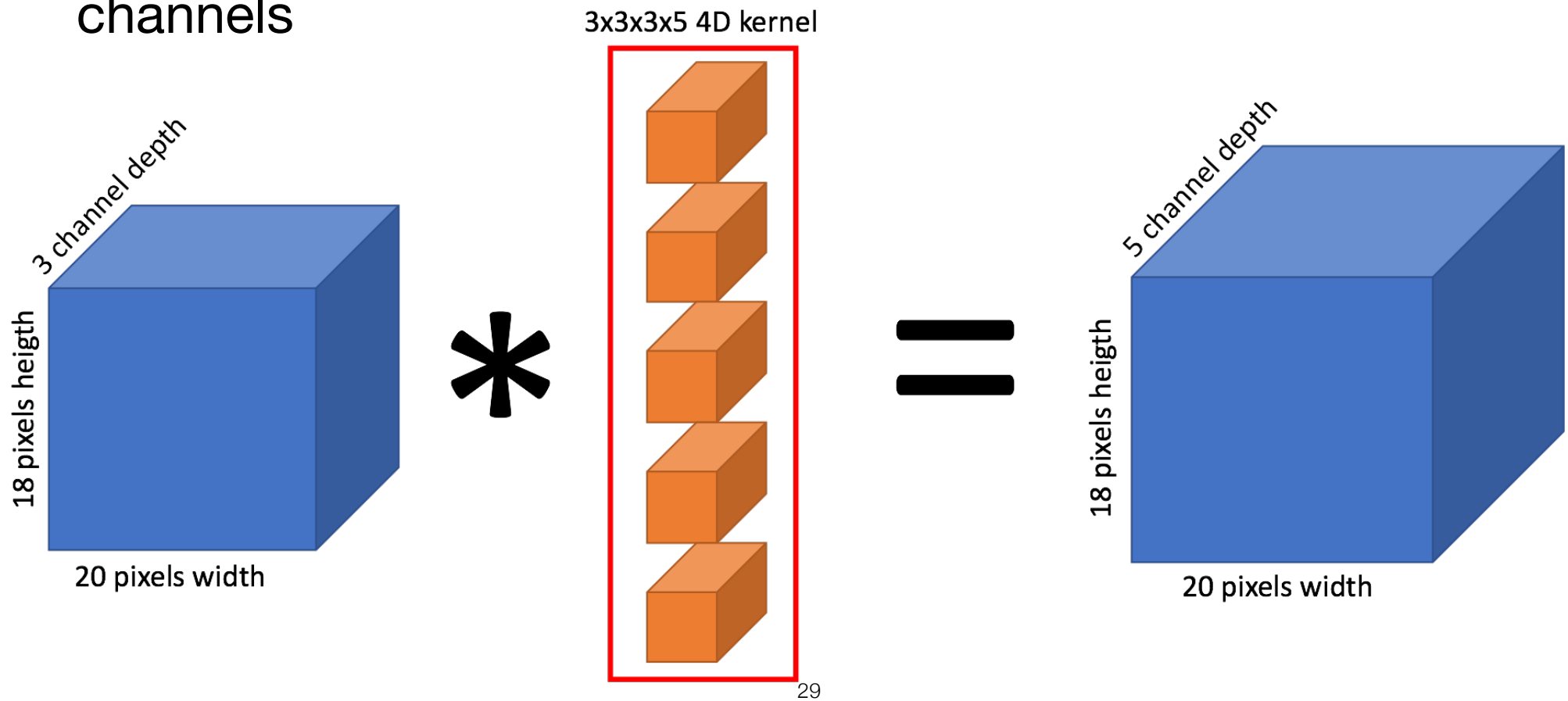
# Volume Convolution (6/9)

- For visualization, it can be interesting to forget the numbers, and just look at 3D arrays as if they were 3D shapes



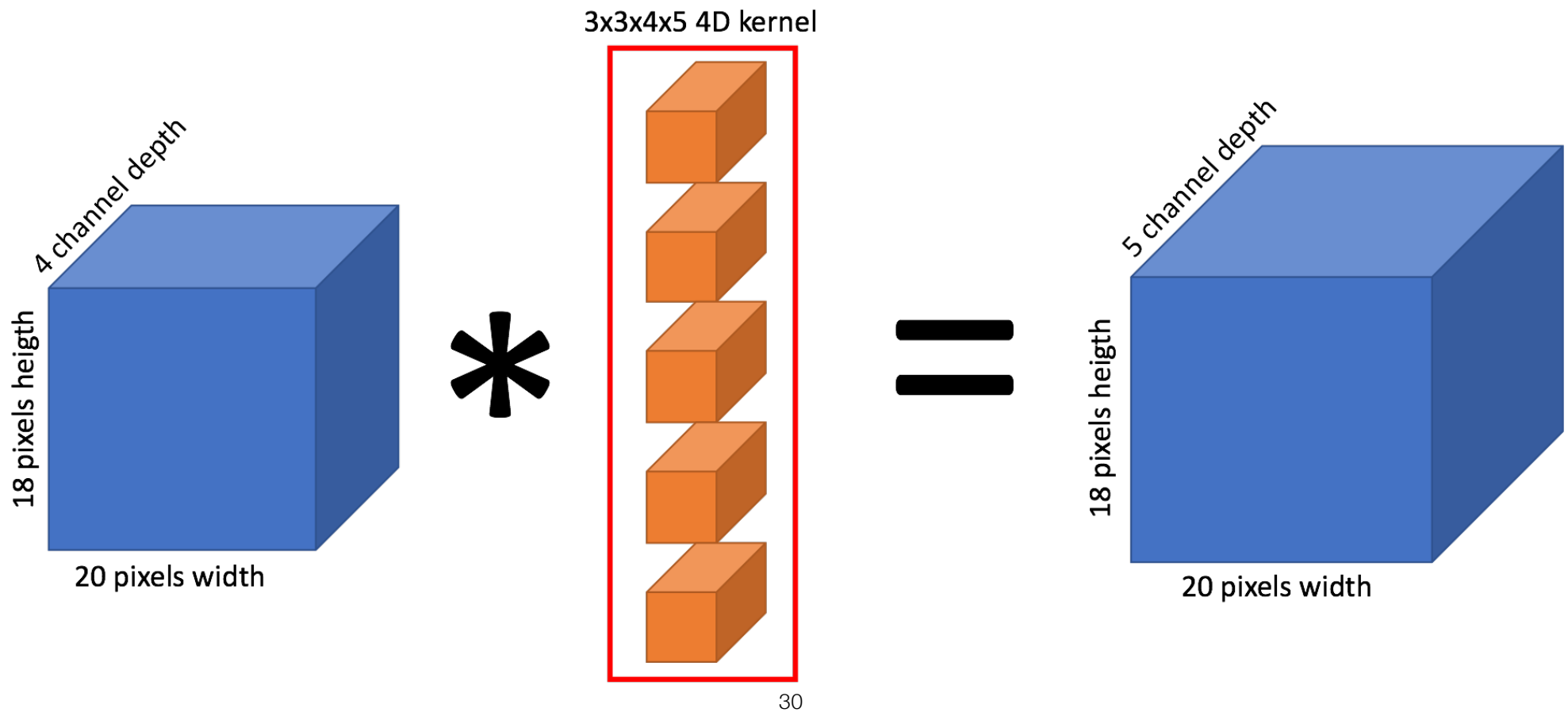
# Volume Convolution (7/9)

- We can generalize to any number of input and output channels



# Volume Convolution (8/9)

- We can generalize to any number of input and output channels



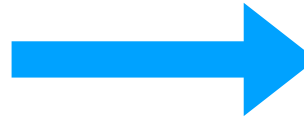
# Volume Convolution (9/9)

- In short the most important thing to remember:
  - Convolutions can take a 3D array as input
  - Can produce a 3D array as output
  - The number of *channels* (ie. third dimension) of input and output array can be different
- Now, we are going to see one last operation: ***“Max Pooling”***

# Max Pooling (1/4)

- Another operation commonly used in CNN: Max Pooling
- Simply takes the max of an area of the input
- Used to reduce the size of the input

1	1	1	-1
4	3	3	2
4	5	3	3
4	5	5	3



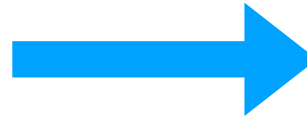
4	3
5	5



# Max Pooling (2/4)

- Another operation commonly used in CNN: Max Pooling
- Simply takes the max of an area of the input
- Used to reduce the size of the input

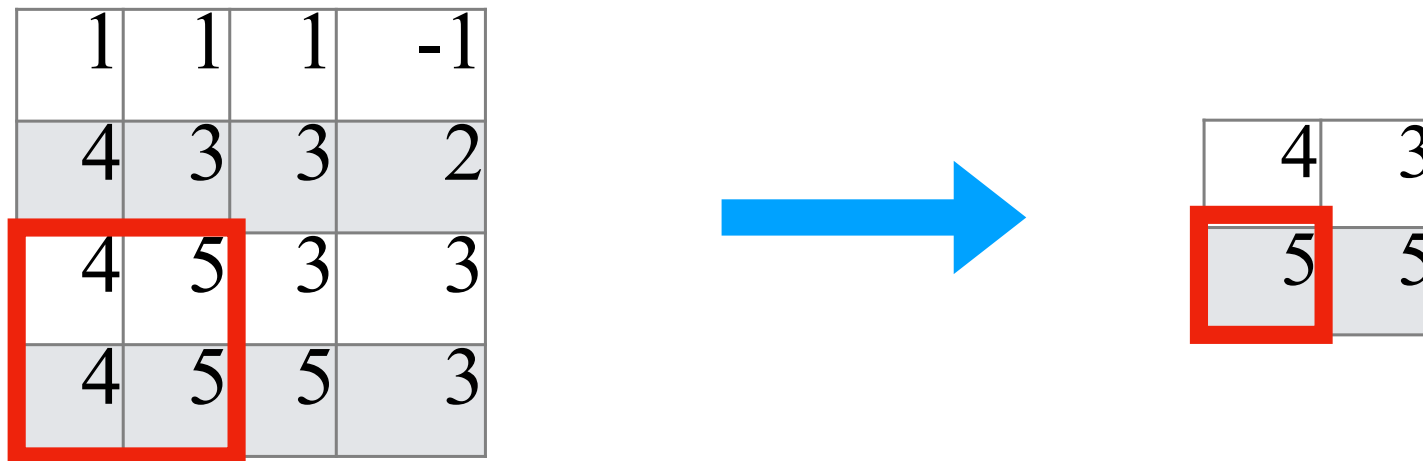
1	1	1	-1
4	3	3	2
4	5	3	3
4	5	5	3



4	3
5	5

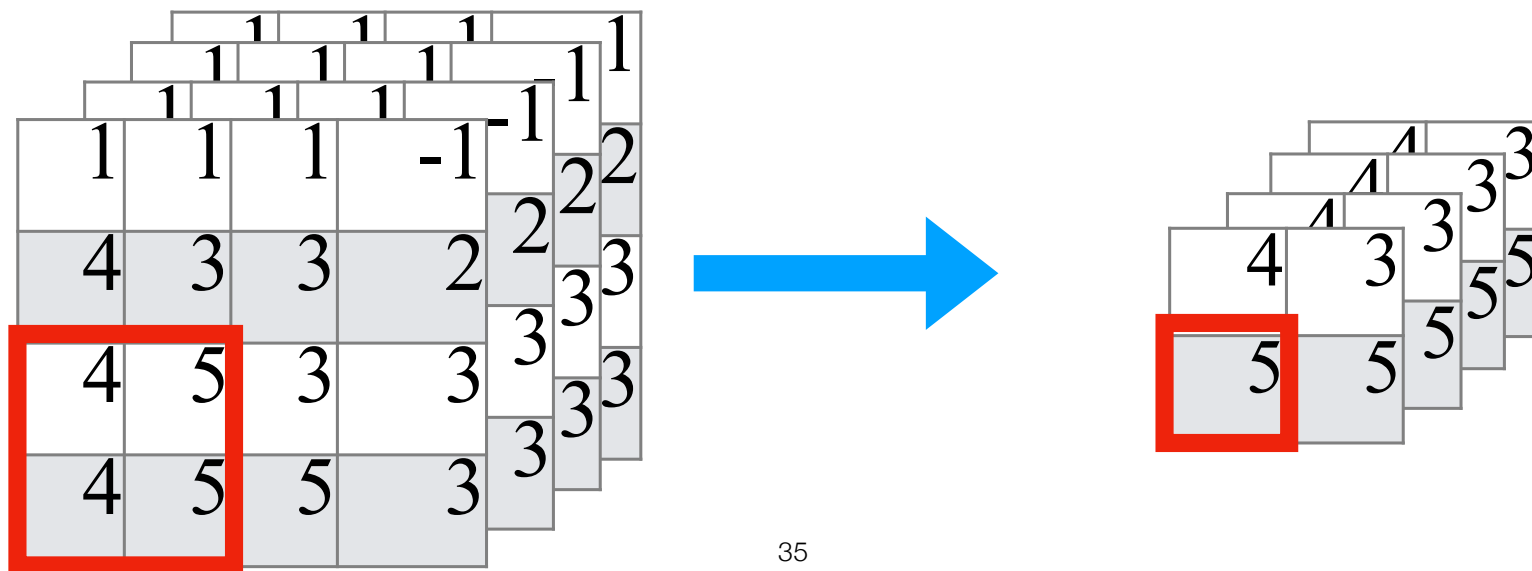
# Max Pooling (3/4)

- Another operation commonly used in CNN: Max Pooling
- Simply takes the max of an area of the input
- Used to reduce the size of the input



# Max Pooling (4/4)

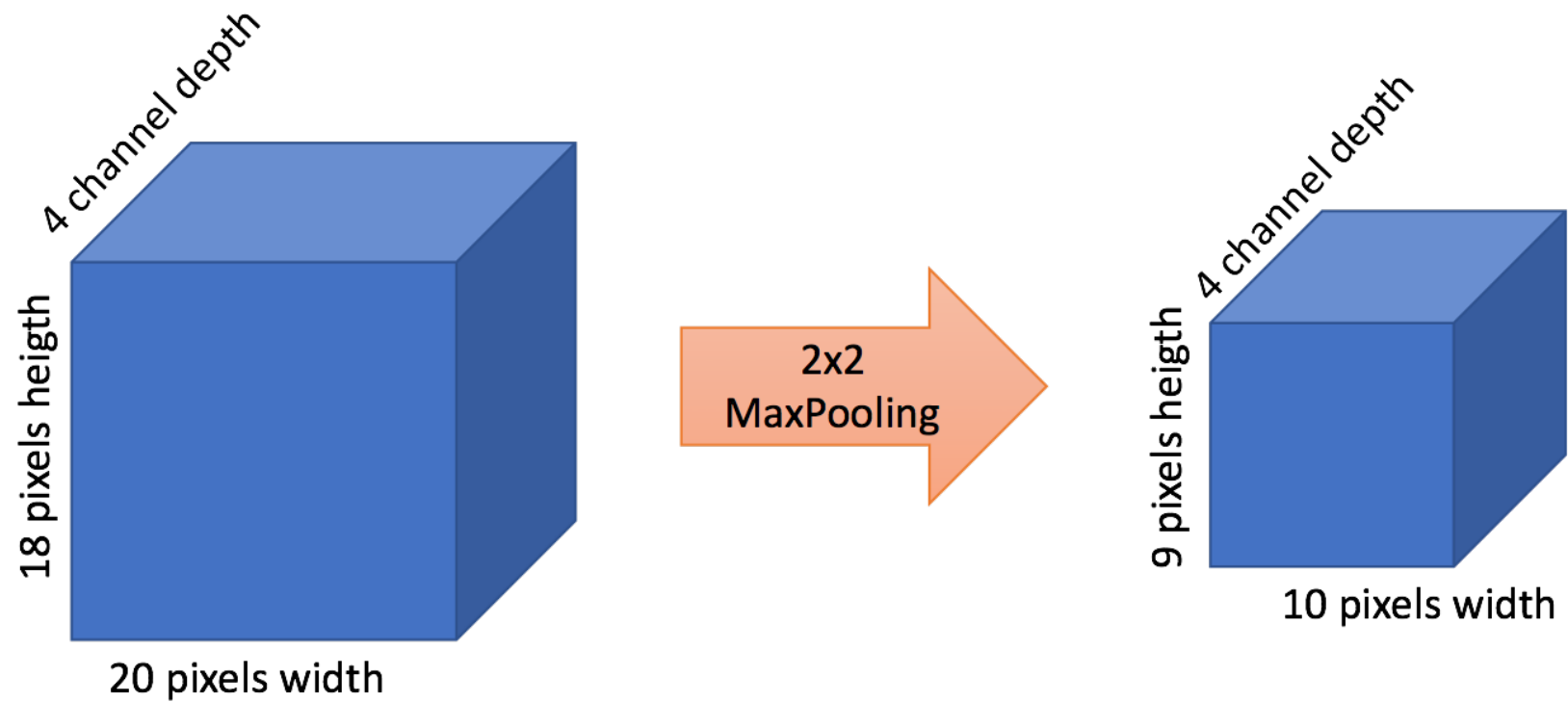
- Another operation commonly used in CNN: Max Pooling
- Simply takes the max of an area of the input
- Used to reduce the size of the input



# “Volume” Max Pooling

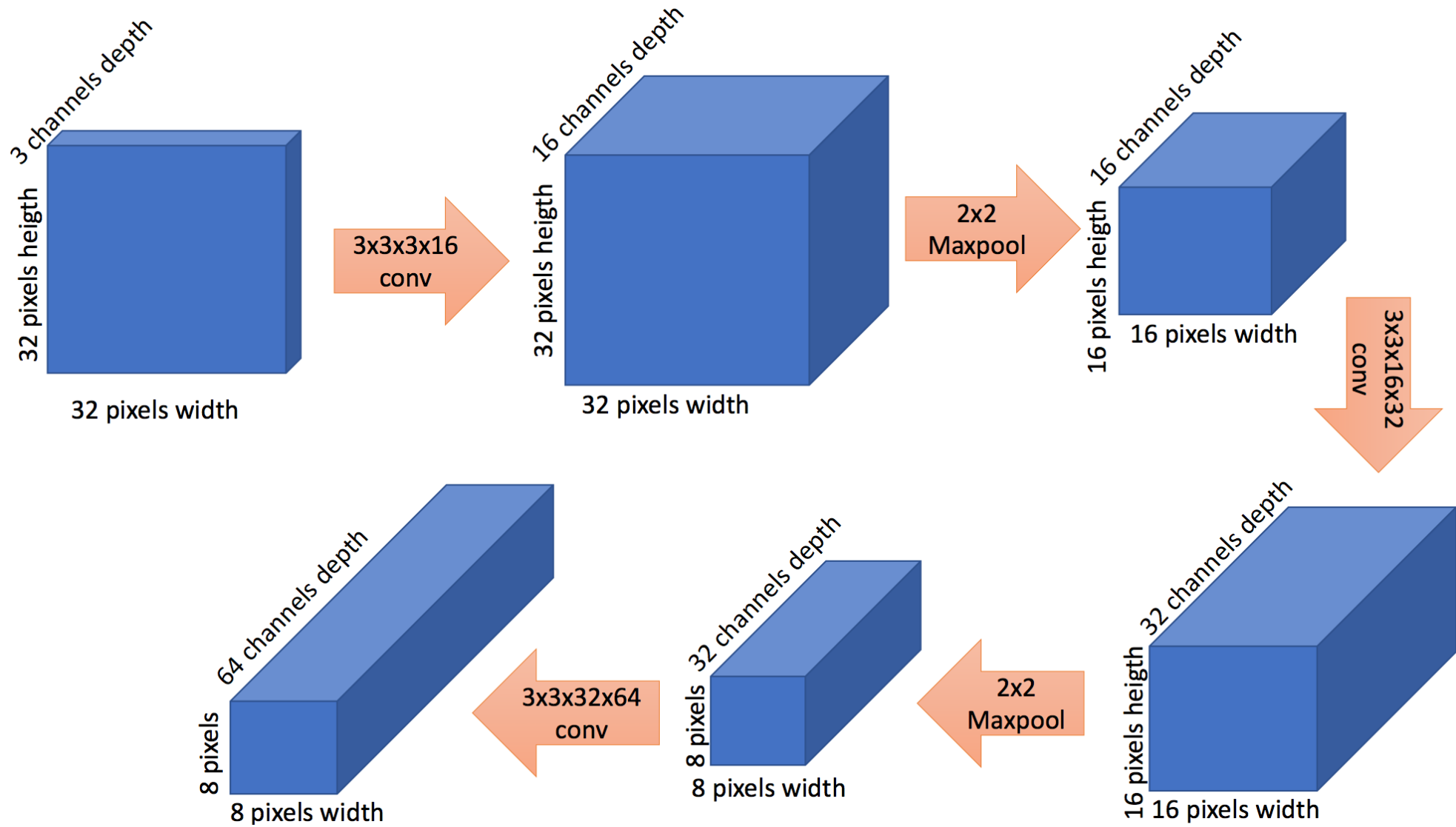
- Max-Pooling can be applied to a 3D array as well
- It is applied to each channel separately
  - The input is a 3D array
  - The output is a 3D array with the same number of channels as the input
  - But with other dimensions divided by 2

# Max Pooling



# Image Classifier

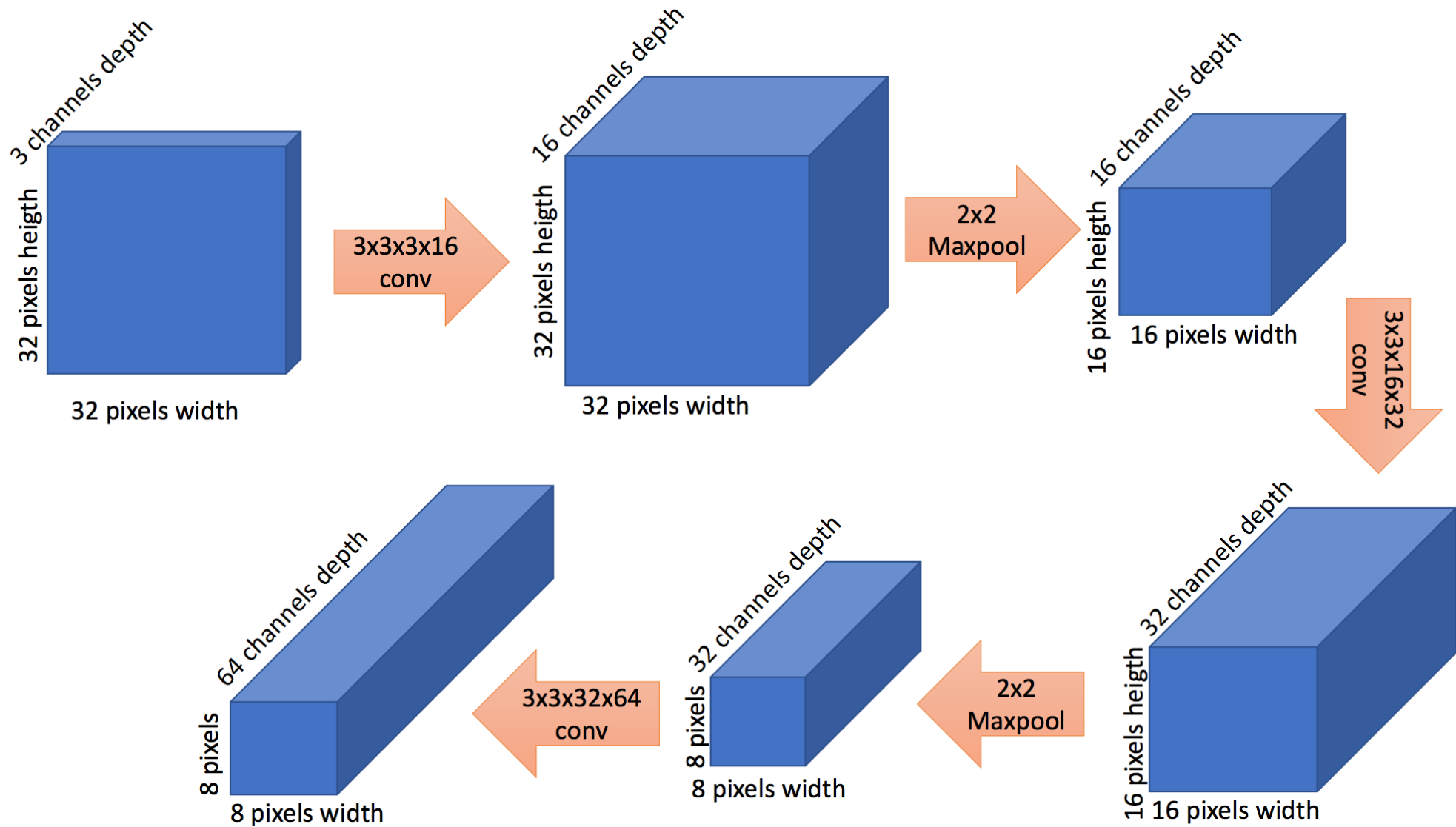
- Finally, we can look at how we build a full image classifier
- A typical modern image classifier is a Multi-Layered Neural Network
  - Input image is sent to a convolutional Layer
  - The result is sent to a Max-Pooling layer
  - The result is sent to a Convolutional Layer
  - And so on....



# Image Classifier

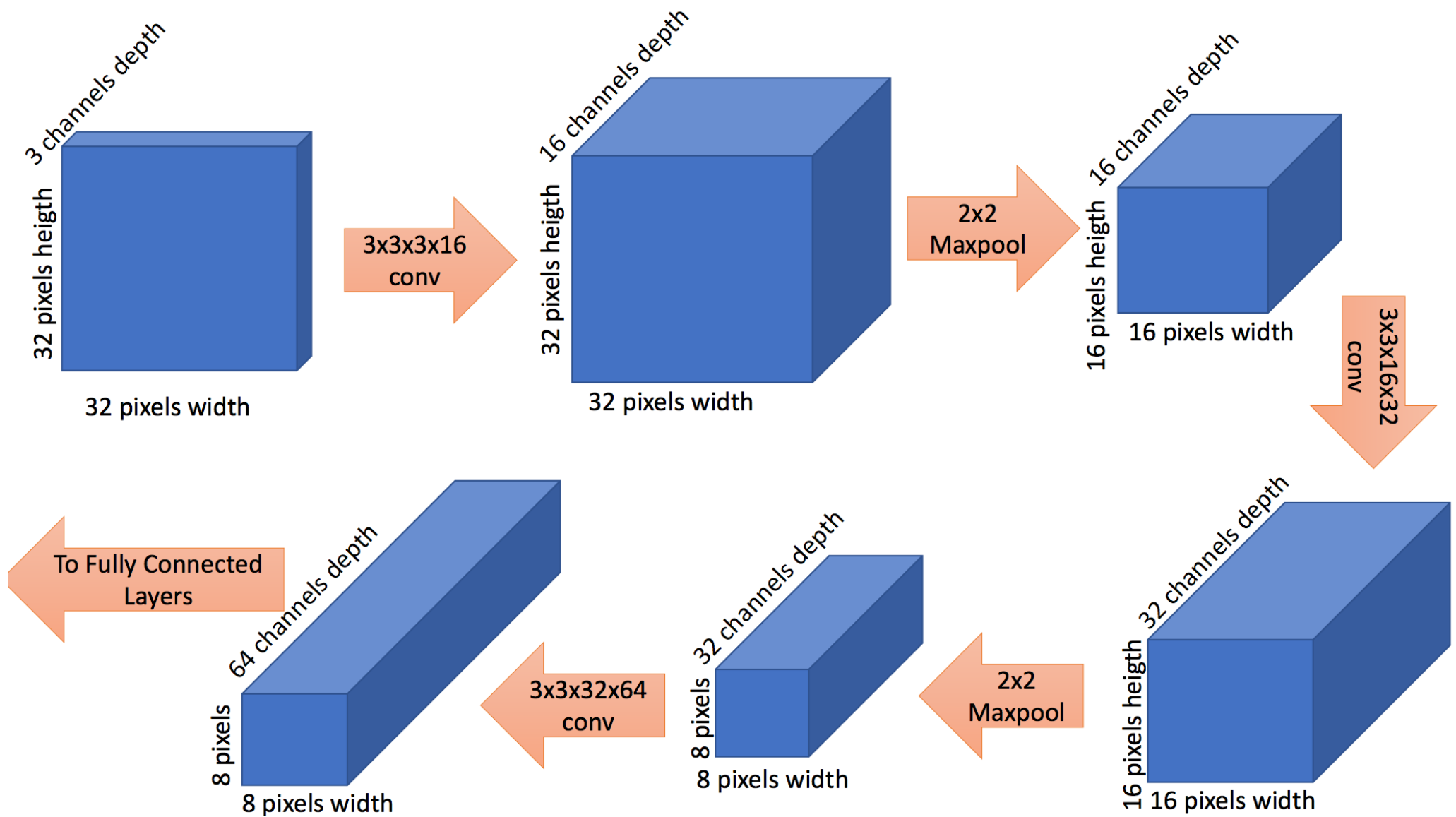
- Finally, we can look at how we build a full image classifier
- A typical modern image classifier is a Multi-Layered Neural Network
- Each Convolutional Layer + Max-Pooling Layer produce a 3D array that is “narrower” and “deeper” than the input



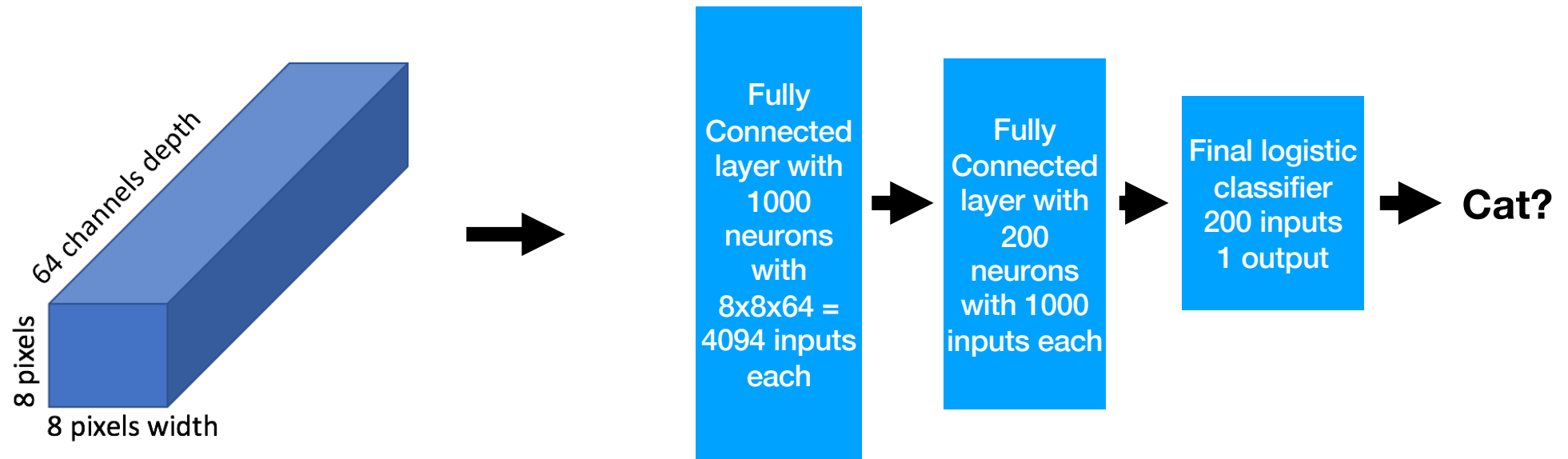


# Image Classifier

- Finally, we can look at how we build a full image classifier
- A typical modern image classifier is a Multi-Layered Neural Network
- Each Convolutional Layer + Max-Pooling Layer produce a 3D array that is “narrower” and “deeper” than the input
- At the end, we send the “deep and narrow” 3D array to a Fully-Connected Classifier

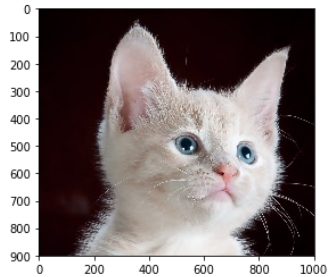


# Final Fully Connected Classifier

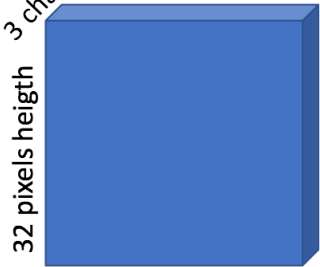


# Image Classifier

Input Image as a 3D array



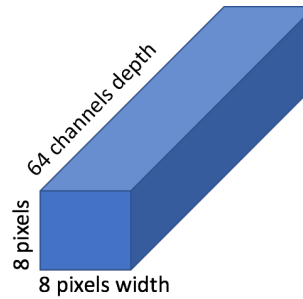
3 channels depth



32 pixels width



Alternate Convolutional Layers  
and Max Pooling Layers



8 pixels

8 pixels width

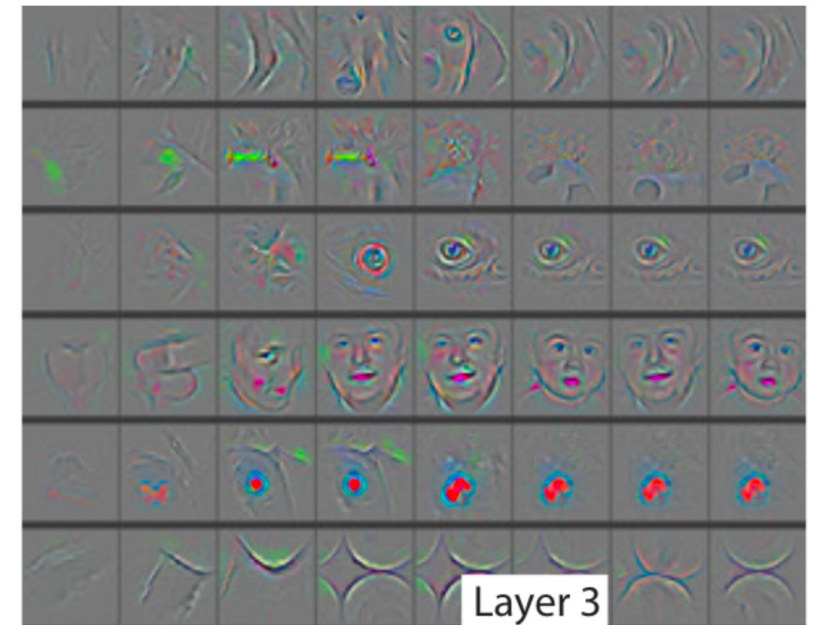
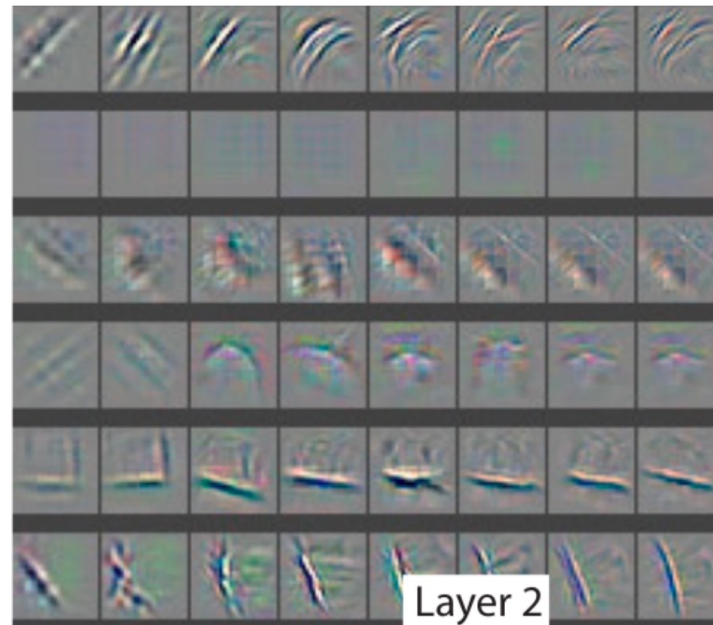
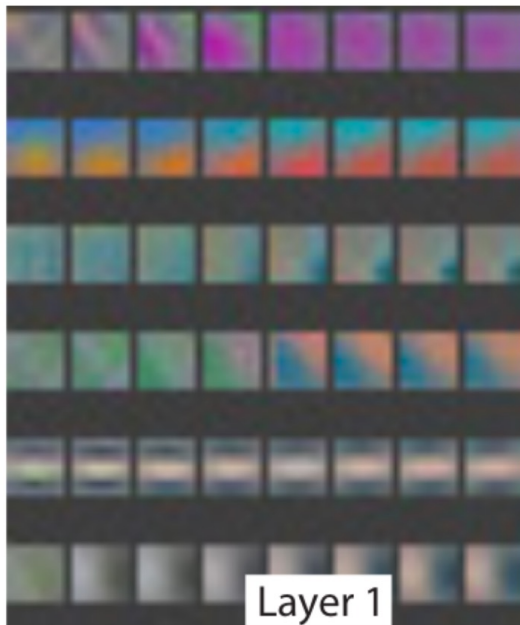


Classifier with Fully Connected  
Layers



**Cat?**

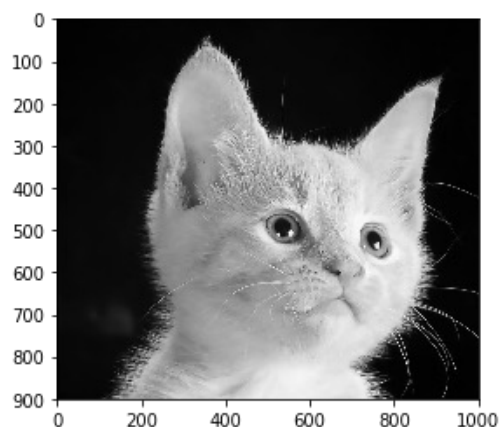
# What Do the Convolution Kernels Learn to Recognize?



From Zeiler&Fergus “Visualizing and Understanding Convolutional Networks”

# Learnable Kernels and Edge Detectors

In practice, we will be  
learning these parameters  
from examples:



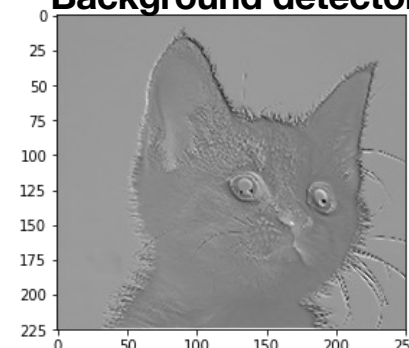
\*

$\Theta_1$	$\Theta_2$	$\Theta_3$
$\Theta_4$	$\Theta_5$	$\Theta_6$
$\Theta_7$	$\Theta_8$	$\Theta_9$

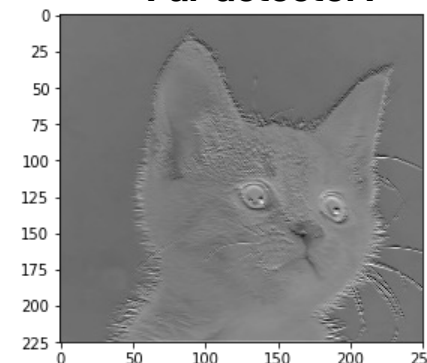
=

- We know some kernels can emphasize some edges in an image (edge detectors)
- Maybe by training the parameters, we discover kernels that can emphasize interesting aspects of the image?

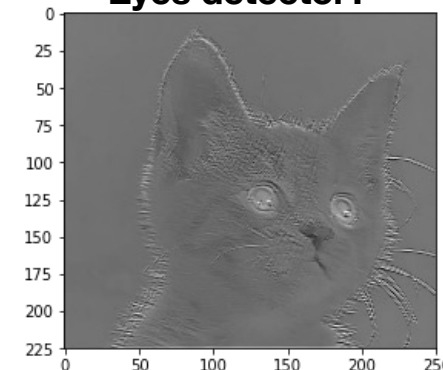
Background detector?



Fur detector?



Eyes detector?



# About Parameters Size (1/3)

- We said that if we have a Fully-Connected Layer with  $n$  neurons and  $m$  input, it contains  $n \times (m+1)$  parameters
- How many parameters in a Max-Pooling Layer ?
- How many parameters in a Convolutional Layer?



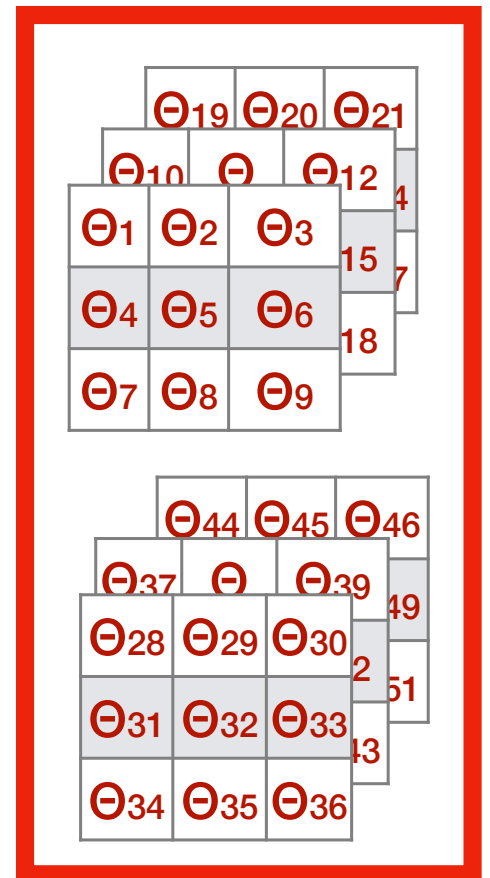
# About Parameters Size (2/3)

- We said that if we have a Fully-Connected Layer with  $n$  neurons and  $m$  input, it contains  $n \times (m+1)$  parameters
- How many parameters in a Max-Pooling Layer ?
  - 0 parameters
- How many parameters in a Convolutional Layer?
  - It depends on the 4D kernel size

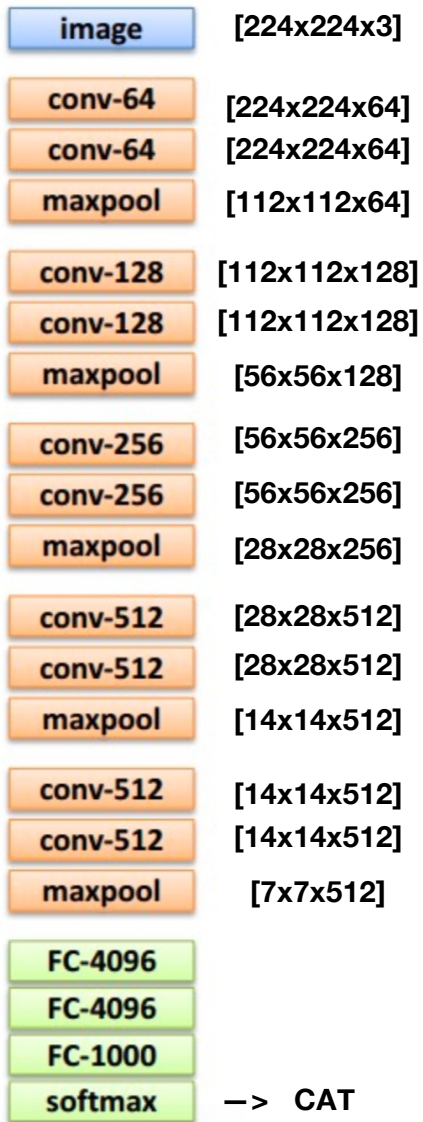
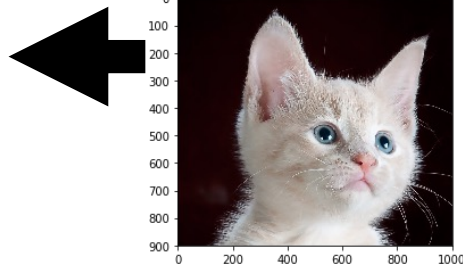
# About Parameters Size (3/3)

- Kernel of size  $k \times k$  with  $C_i$  input channels and  $C_o$  output channels
- It means we have  $C_o$  neurons connected to  $k \times k \times C_i$  inputs for every location in the image
  - $(k \times k \times C_i + 1) \times C_o$  parameters
- For the kernel on the right:  $k=3$   $C_i=3$   $C_o=2$ 
  - $(3 \times 3 \times 3 + 1) \times 2 = 56$  parameters

Learnable 4D Kernel

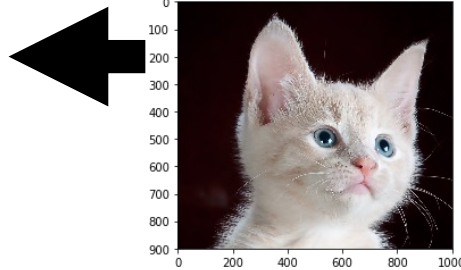
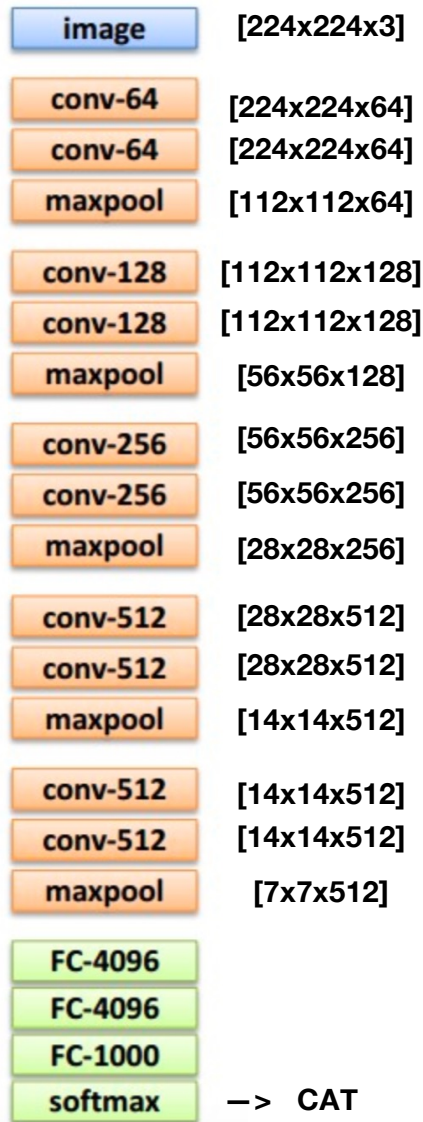


# Example: VGG Network (1/2)



- “Very Deep Convolutional Networks for Large-Scale Image Recognition” Simonyan and Zisserman, 2015
- Best Image Classifier in 2015
- Today the best models have much more layers (> 100)
- Could you compute the number of parameters in each layer?
- All kernels are of size 3 (k=3)
- Conv-64 mean convolutional layer with 64 output channels
- FC-4096 means Fully Connected Layer with 4096 neurons
- Size of input can be guessed from the size of output of the previous layer

# Example: VGG Network (2/2)



- Number of parameters for each layer:
- Conv-64 (1st):  $(3 \times 3 \times 3 + 1) \times 64 = 1792$
- Conv-64 (2nd):  $(3 \times 3 \times 64 + 1) \times 64 = 36\,928$
- Conv-128 (1st):  $(3 \times 3 \times 64 + 1) \times 128 = 73\,856$
- Conv-128 (2nd):  $(3 \times 3 \times 128 + 1) \times 128 = 147\,584$
- Conv-256 (1st):  $(3 \times 3 \times 128 + 1) \times 256 = 295\,168$
- Conv-256 (2nd):  $(3 \times 3 \times 256 + 1) \times 256 = 590\,080$
- Conv-512:  $(3 \times 3 \times 256 + 1) \times 512 = 1\,180\,160$
- 3x Conv-512 (2nd step) :  $(3 \times 3 \times 512 + 1) \times 512 = 2\,359\,808 \times 3$
- FC-4096 (1st):  $(7 \times 7 \times 512 + 1) \times 4096 = 102\,764\,544$
- FC-4096 (2nd) :  $(4096 + 1) \times 4096 = 16\,781\,312$
- FC-1000 :  $(4096 + 1) \times 1000 = 4\,097\,000$
- Total:  $1792 + 36\,928 + 73\,856 + 147\,584 + 295\,168 + 590\,080 + 1\,180\,160 + 2\,359\,808 \times 3 + 102\,764\,544 + 16\,781\,312 + 4\,097\,000 = 133\,047\,848$  parameters

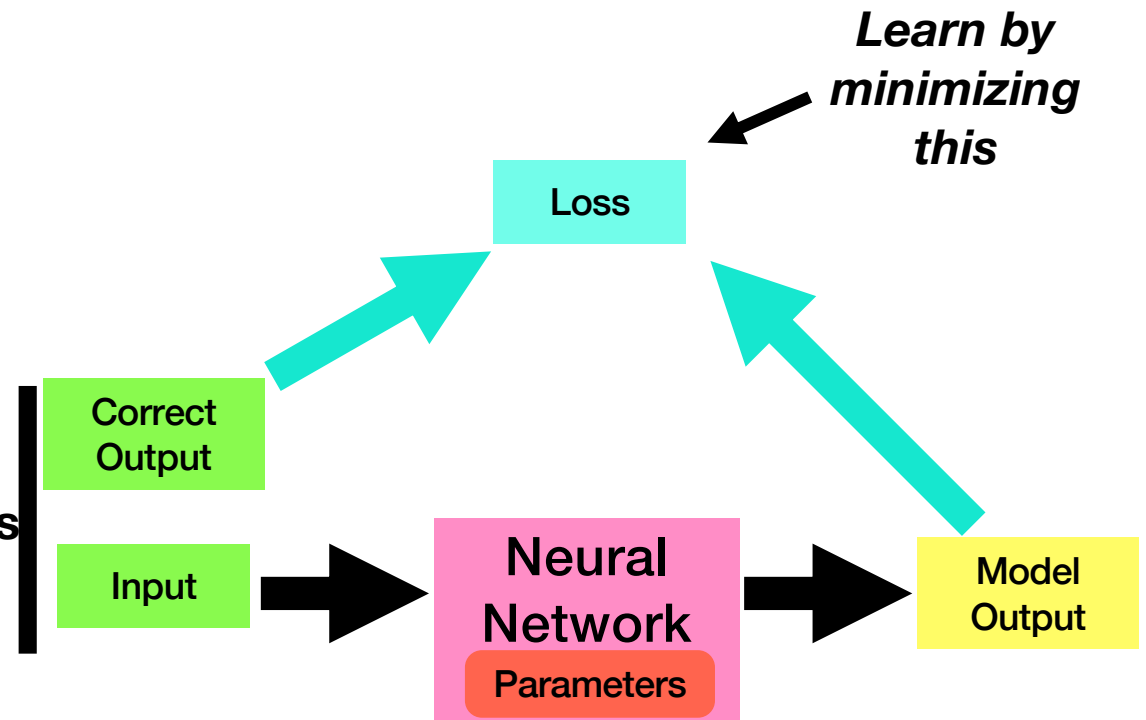
# Supervised Learning (1/2)

- In supervised learning, we usually have:

- A **MODEL**: a “parameterized” function that takes input and produce output

- A **Loss**: A function that compute how different the **Examples** model output is from the correct output

- **Examples** of input and correct output



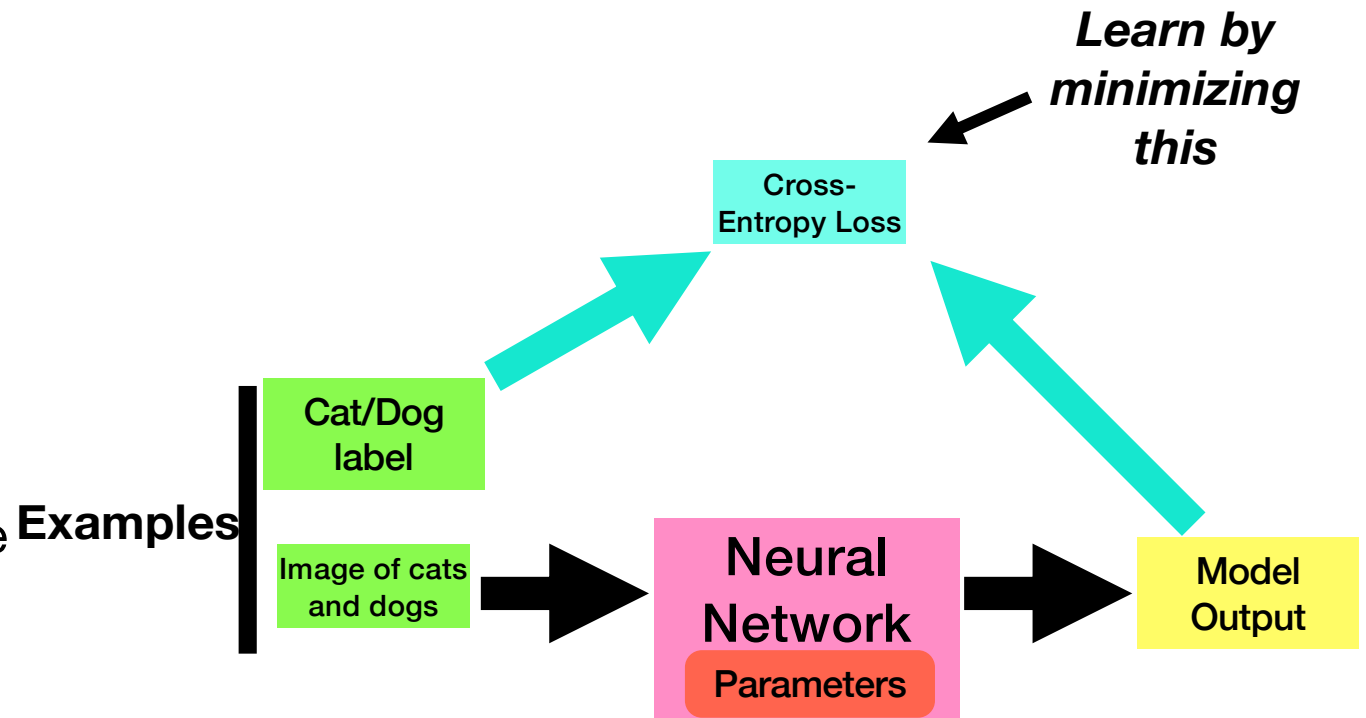
# Supervised Learning (2/2)

- In supervised learning, we usually have:

- A **MODEL**: a “parameterized” function that takes input and produce output

- A **Loss**: A function that compute how different the model output is from the correct output

- **Examples** of input and correct output



# Actual Training

- Train your network with labeled images and you are good:



**Not Cat**



**Cat**



**Cat**



**Not Cat**

**How to get 5  
images for the price  
of one?**

# Data Augmentation (1/4)

**How to get 5  
images for  
the price of  
one?**



**A cat**



# Data Augmentation (2/4)

How to get 5  
images for  
the price of  
one?



**A cat**

**Mirroring**



**Still a cat**

# Data Augmentation (3/4)

How to get 5  
images for  
the price of  
one?



**A cat**



**Mirroring**

**Still a cat**



**Color change**

**Still a cat**

# Data Augmentation (4/4)

How to get 5  
images for  
the price of  
one?



A cat

Cropping



Still a cat

Mirroring



Still a cat

Color change



Still a cat

Rotation



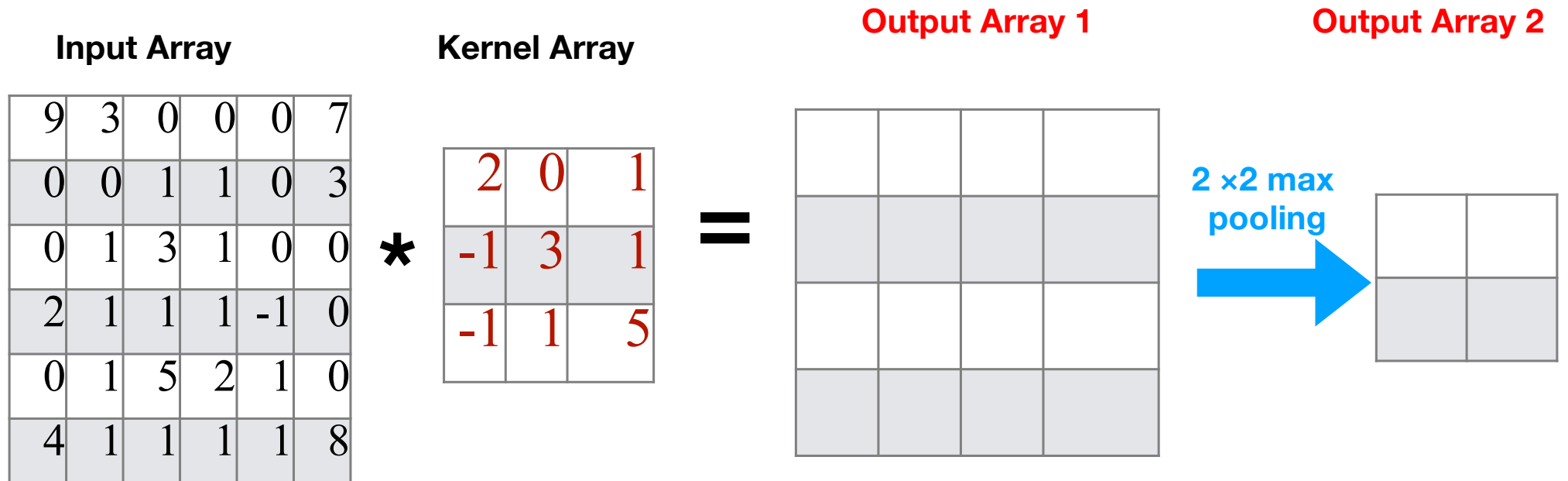
Still a cat

And any combination  
of these

# Next Time

- We will learn how to do image recognition with Chainer

# Exercise



# Report

- Submit the report of **exercise in pdf** via Panda
- Submission due: **next lecture**
- Name the pdf file as **student id\_name**

# Final Report

- Content: Describe a real problem in your study that can be solved by the models (either basic supervised machine learning or deep learning) introduced in this course together with the model structure and the reason for using the model. It would be better if you implemented the model and conducted experiments. The report should have no less than 2,000 words in English. Figures & diagrams are allowed to assist your illustration.
- Submit your final report in pdf named as [student id\_name] via Panda by August 1st, 2025.