# CS540 Introduction to Artificial Intelligence
## Convolutional Neural Networks (II)

University of Wisconsin-Madison

**Spring 2025**
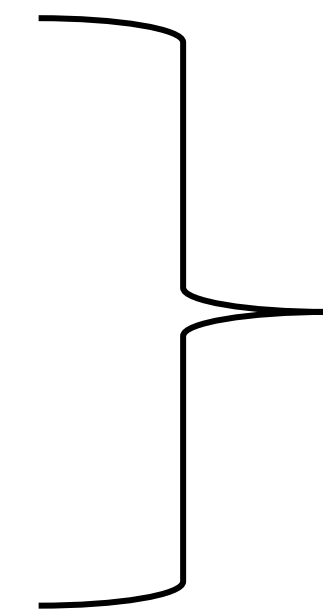
# Announcements

- **Homeworks**:
  - HW7 online, deadline on Monday **April 7ᵗʰ at 11:59 PM**

- Class roadmap and schedule:

| | |
|---|---|
| Machine Learning: Deep Learning II | ⎫ |
| Machine Learning: Deep Learning III | ⎬ Deep Learning |
| *Spring Recess March 22-30* | |

# Today's goals

- Review (some of) convolutional computations.

  - 2D convolutions, multiple input channels, pooling.

- Understand how convolutions are used as layers in a (deep) neural network.

- Build intuition for output of convolutional layers.

- Overview the evolution of deeper convolutional networks
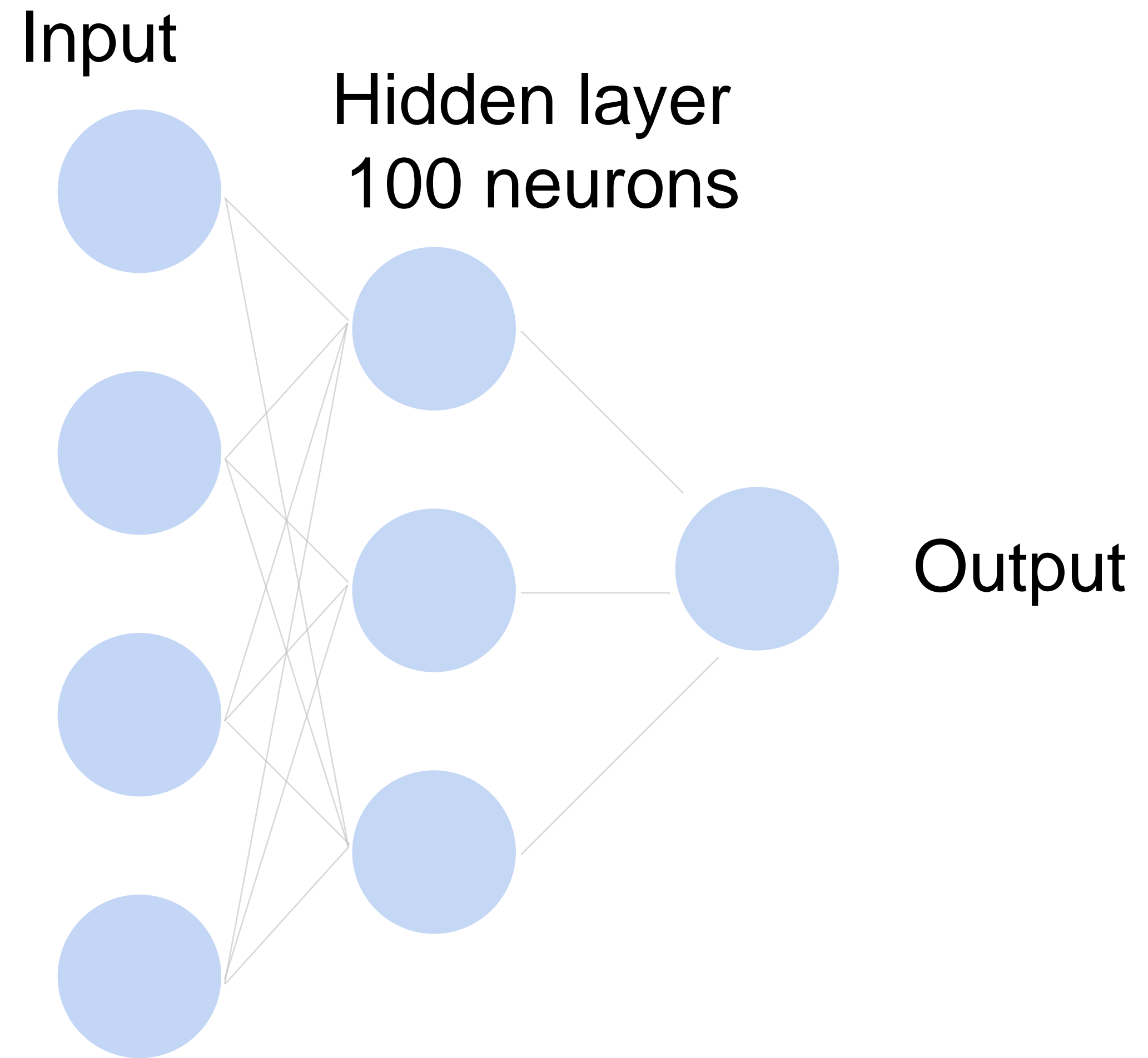
# How to classify

**Cats vs. dogs?**

Dual

# 12MP

wide-angle and
telephoto cameras

**36M** floats in a RGB image!

# Fully Connected Networks

Input

Hidden layer
100 neurons

**Cats vs. dogs?**

Output

36M elements x 100 = **3.6B** parameters!

# Review: 2-D Convolution

Input    Kernel    Output



$$0 \times 0 + 1 \times 1 + 3 \times 2 + 4 \times 3 = 19,$$
$$1 \times 0 + 2 \times 1 + 4 \times 2 + 5 \times 3 = 25,$$
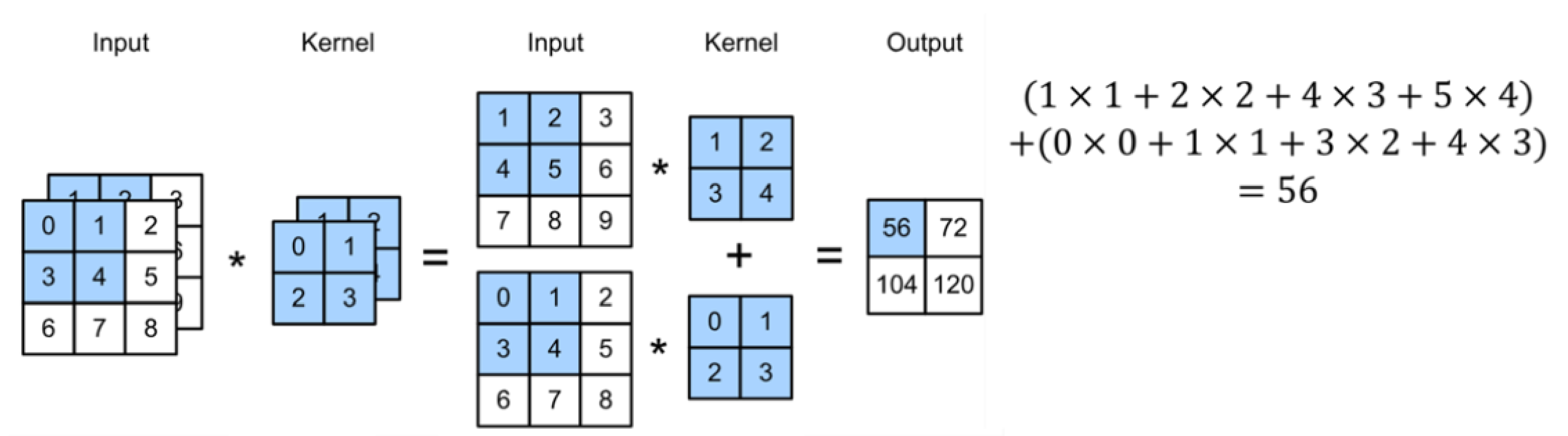$$3 \times 0 + 4 \times 1 + 6 \times 2 + 7 \times 3 = 37,$$
$$4 \times 0 + 5 \times 1 + 7 \times 2 + 8 \times 3 = 43.$$

(vdumoulin@ Github)

# Review: Multiple Input Channels

- Input and kernel can be 3D, e.g., an RGB image have 3 channels
- Have a kernel for each channel, and then sum results over channels



$$(1 \times 1 + 2 \times 2 + 4 \times 3 + 5 \times 4)$$
$$+(0 \times 0 + 1 \times 1 + 3 \times 2 + 4 \times 3)$$
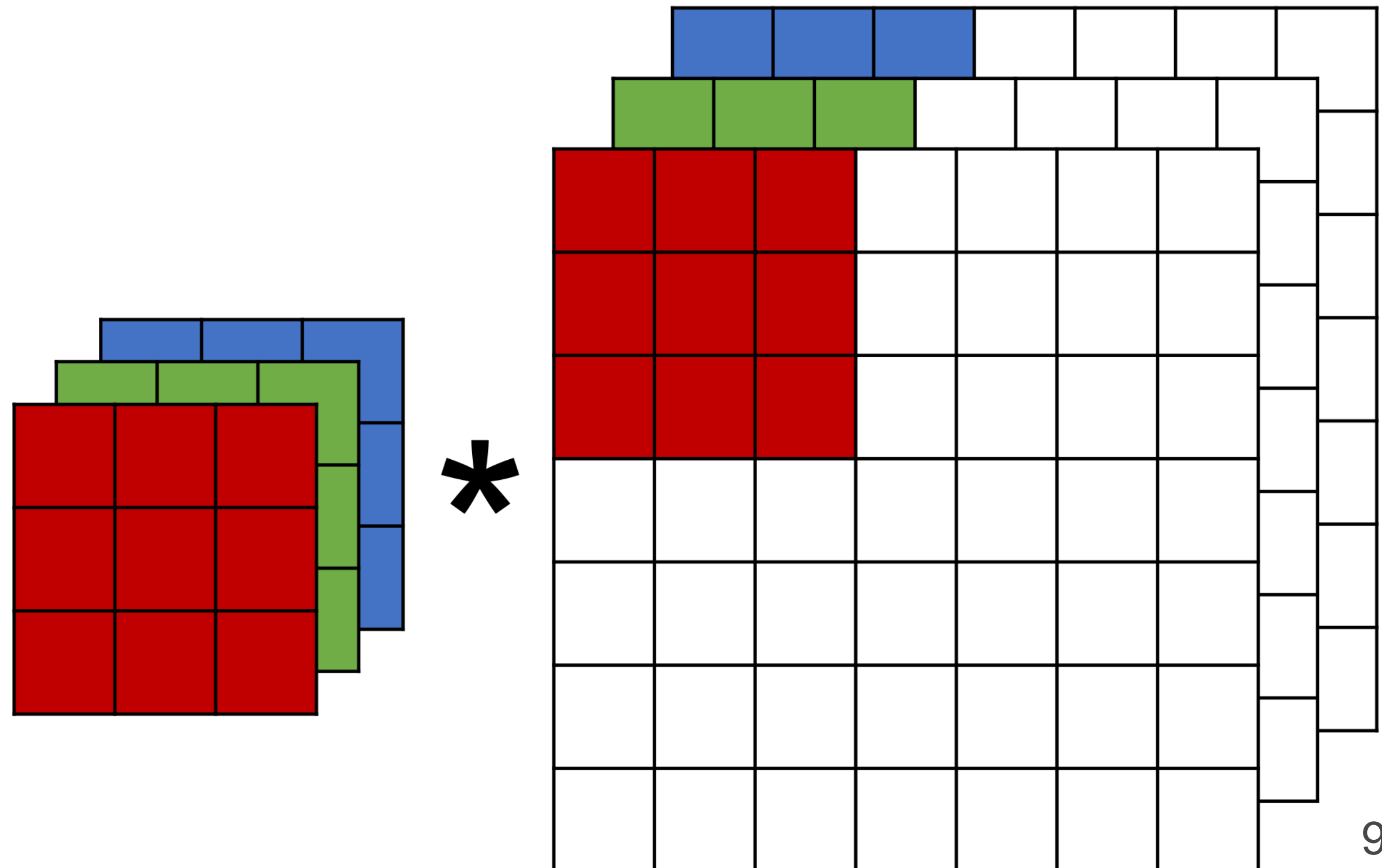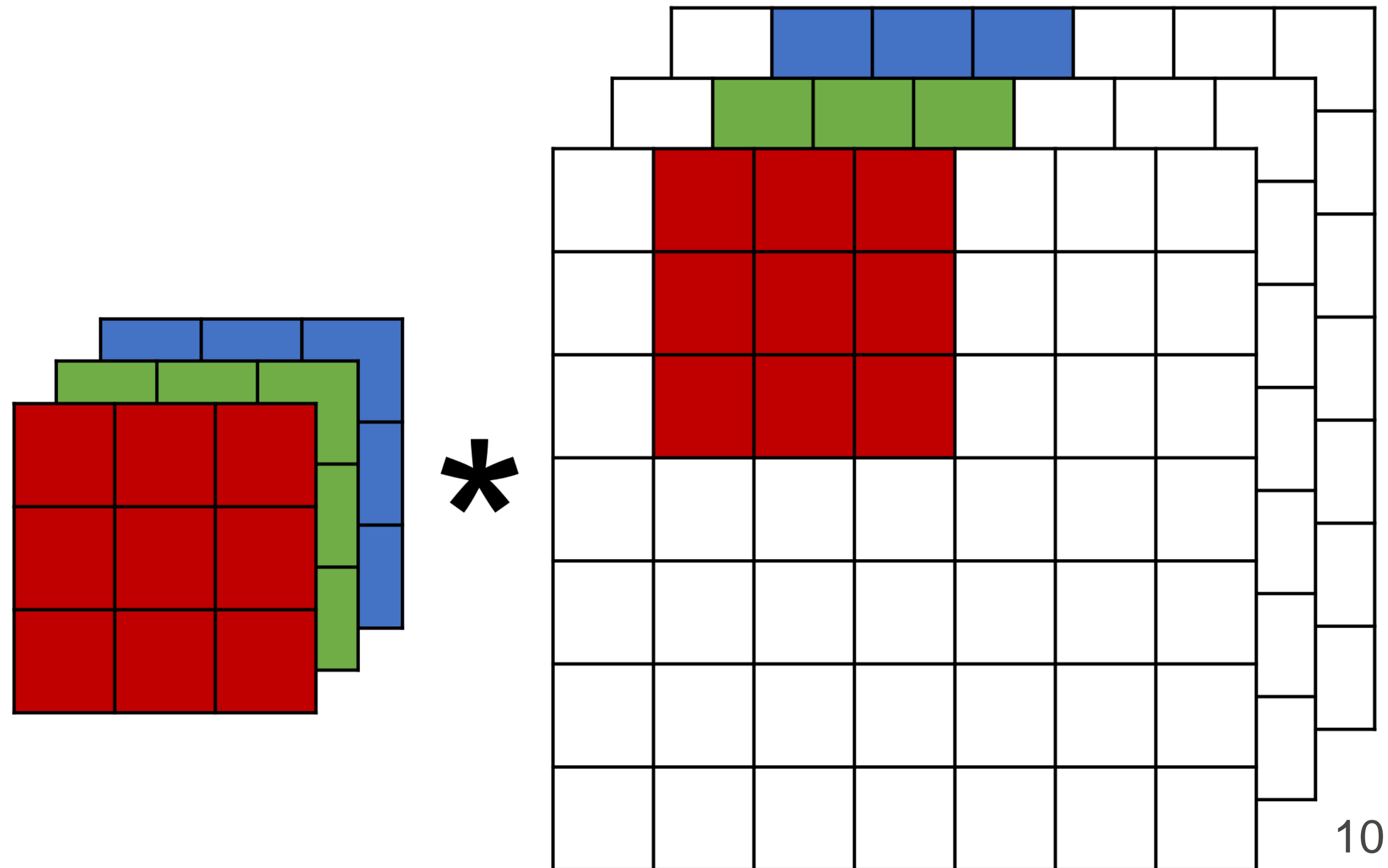$$= 56$$

# Review: Multiple Input Channels

- Input and kernel can be 3D, e.g., an RGB image have 3 channels

- Have a kernel for each channel, and then sum results over channels
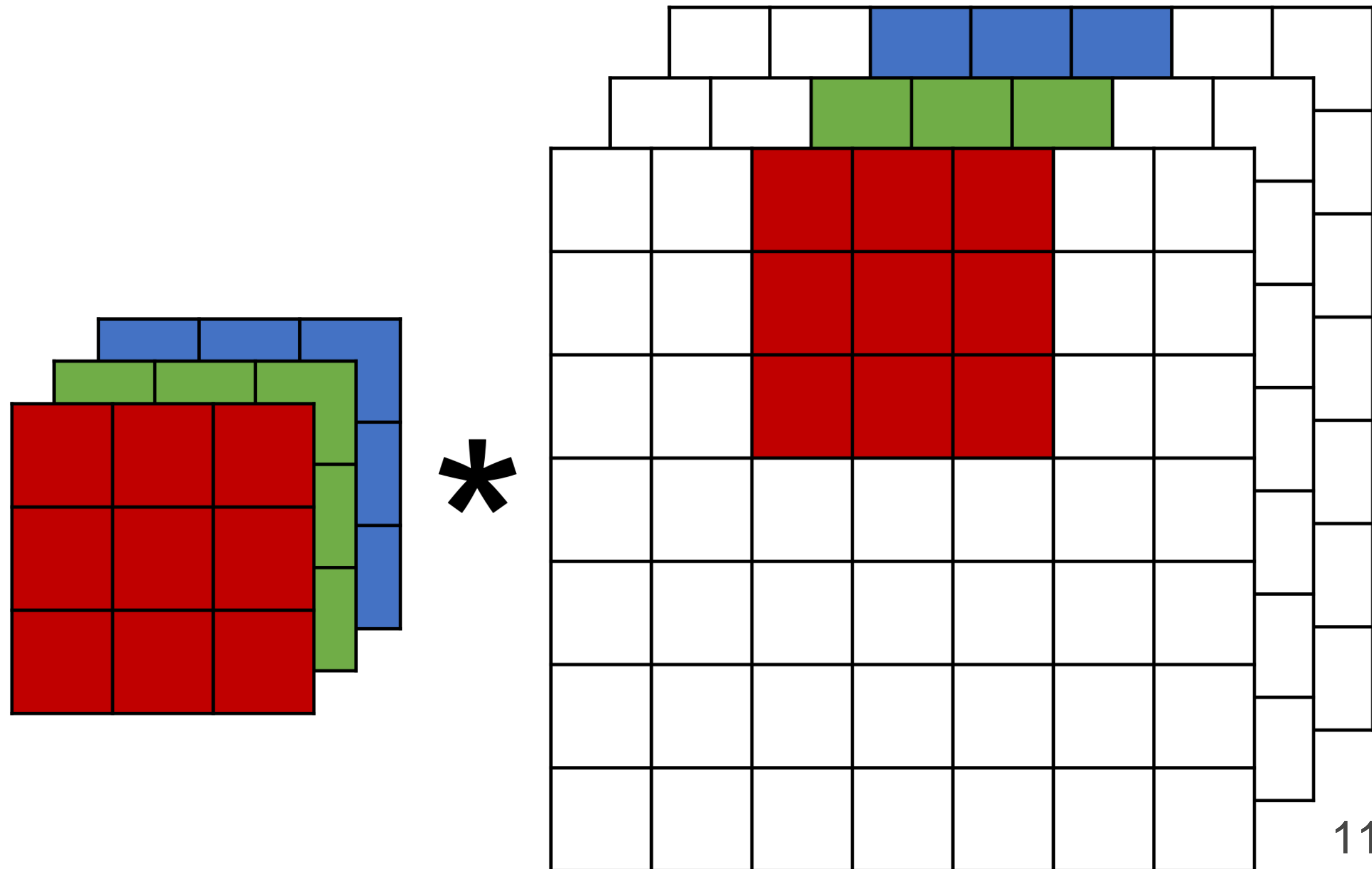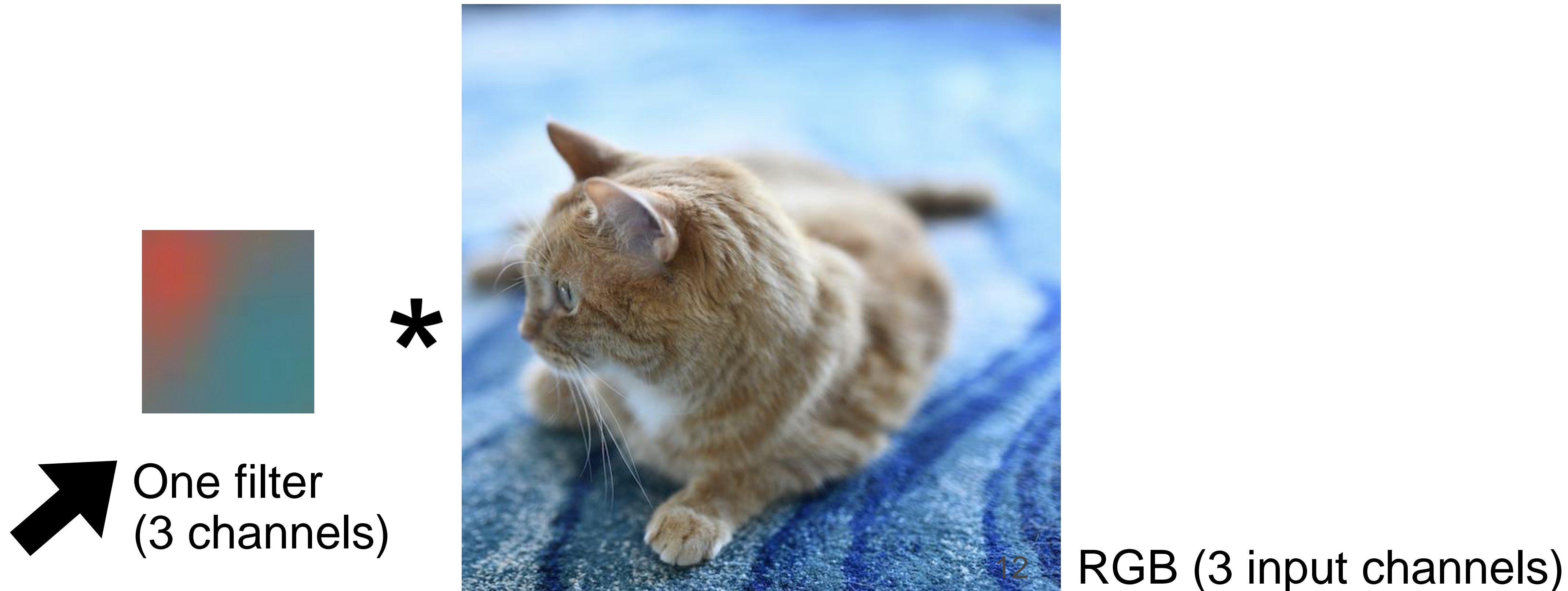
# Review: Multiple Input Channels

- Input and kernel can be 3D, e.g., an RGB image have 3 channels

- Have a kernel for each channel, and then sum results over channels

# Review: Multiple Input Channels

- Input and kernel can be 3D, e.g., an RGB image have 3 channels

- Have a kernel for each channel, and then sum results over channels



11

# Multiple Input Channels

- Input and kernel can be 3D, e.g. RGB image has 3 channels
- Also call each 3D kernel a "**filter**", which produces only **one** output channel (due to summation over channels)



One filter
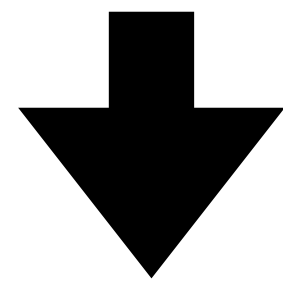(3 channels)

RGB (3 input channels)

# Multiple filters (in one layer)

- Apply multiple filters on the input
- Each filter may learn different features about the input
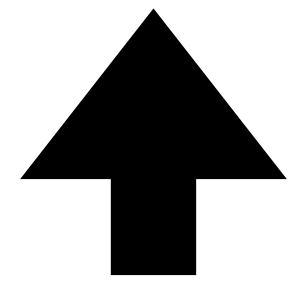- Each filter (3D kernel) produces one output channel



**A different filter**
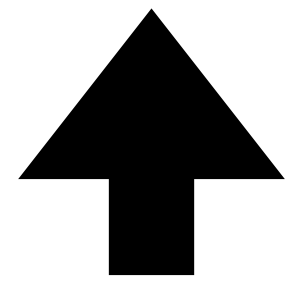
RGB (3 input channels)
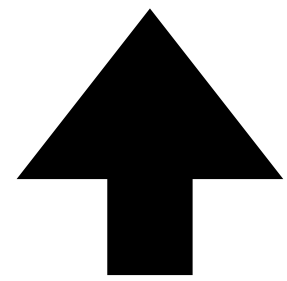
# Output shape

**Kernel/filter size**

$$\lfloor (n_h - k_h + p_h + s_h)/s_h \rfloor \times \lfloor (n_w - k_w + p_w + s_w)/s_w \rfloor$$

**Input size**   **Pad**   **Stride**

14
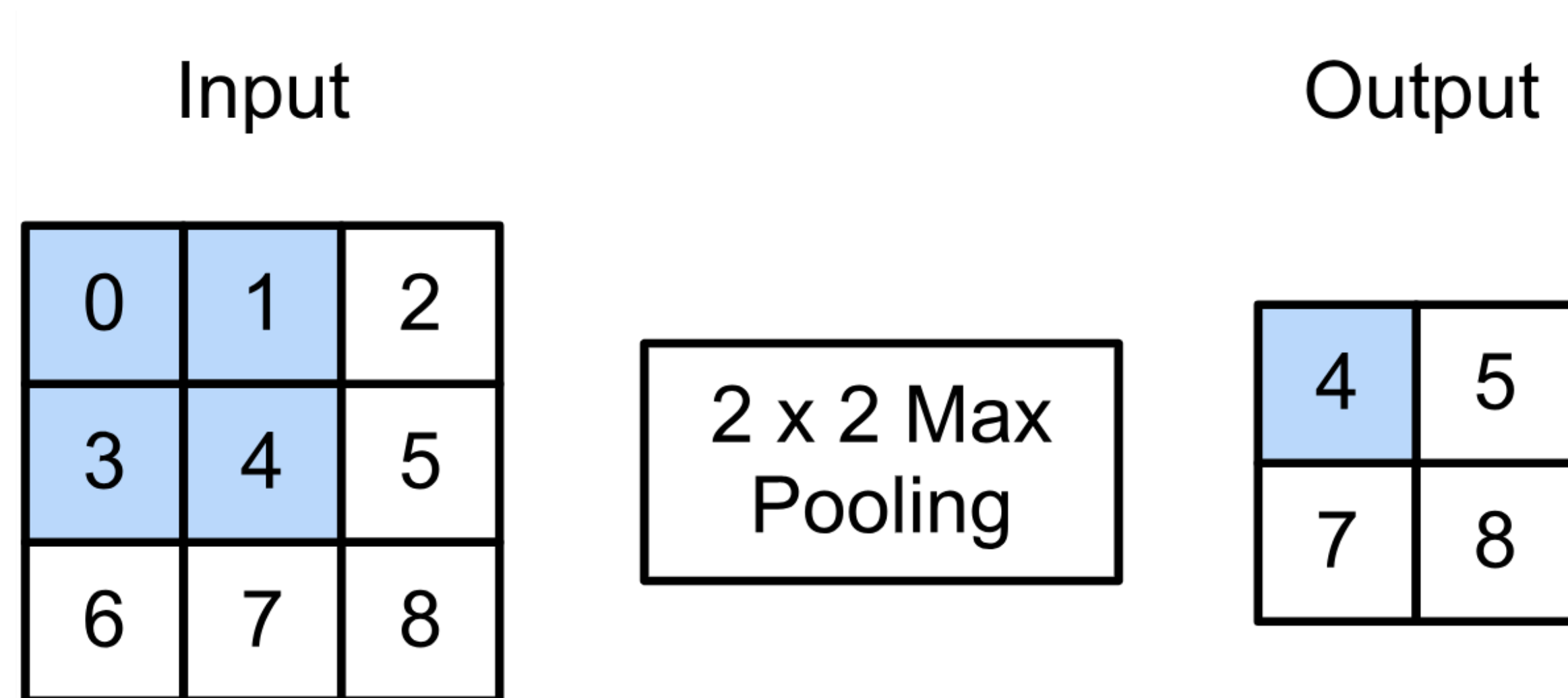
Consider a convolution layer with 16 filters. Each filter has a size of 11x11x3, a stride of 2x2. Given an input image of size 22x22x3, if we don't allow a filter to fall outside of the input, what is the output size?

- 11x11x16
- 6x6x16
- 7x7x16
- 5x5x16

# Pooling Layer

# 2-D Max Pooling

- Returns the maximal value in the sliding window



Input

| 0 | 1 | 2 |
|---|---|---|
| 3 | 4 | 5 |
| 6 | 7 | 8 |

2 x 2 Max Pooling

Output

| 4 | 5 |
|---|---|
| 7 | 8 |

$$max(0,1,3,4) = 4$$

# Average Pooling

- Max pooling: the strongest pattern signal in a window

- Average pooling: replace max with mean in max pooling

  - The average signal strength in a window

Max pooling

Average pooling

# Convolutional Neural Network Architecture
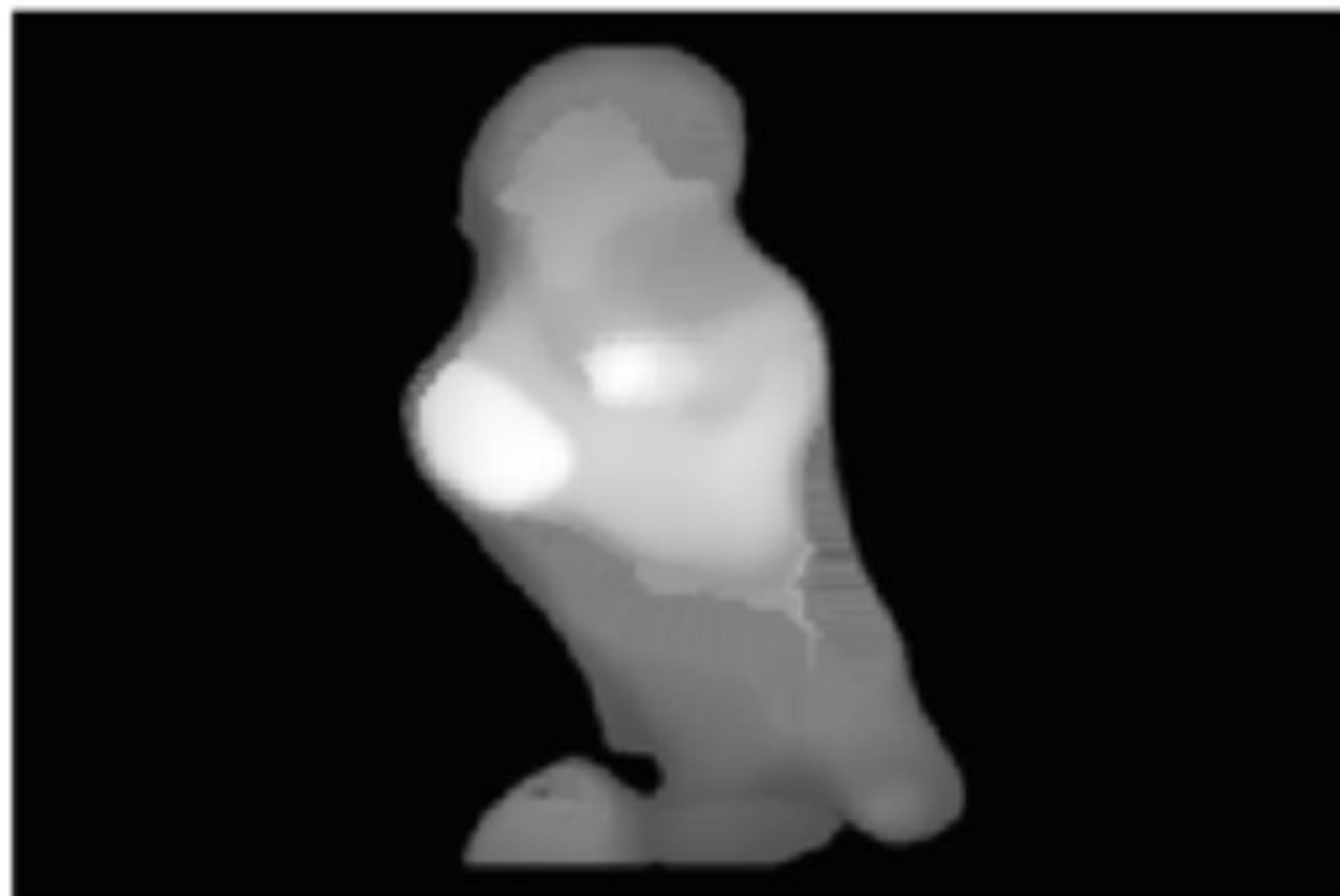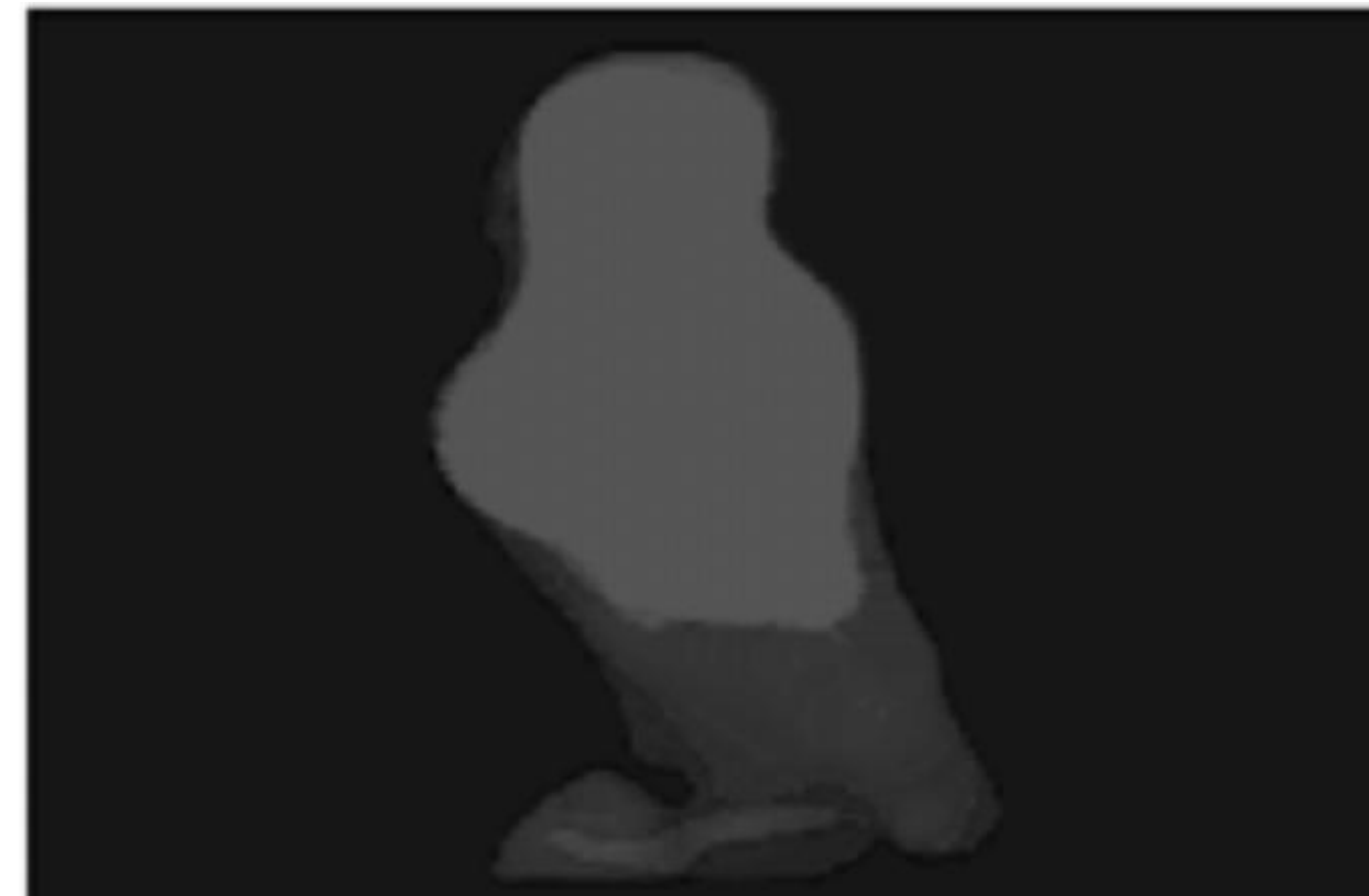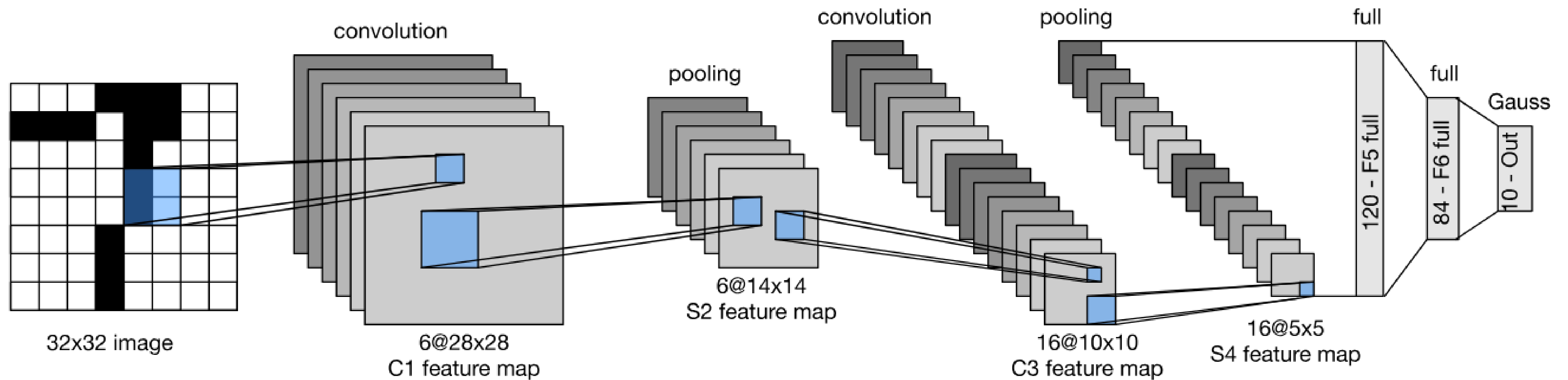


Gradient-based learning applied to document recognition, by Y. LeCun, L. Bottou, Y. Bengio and P. Haffner

# Convolutional Neural Network Intuition

Early layers recognize simple visual features, later layers recognize more complex visual features.

Suppose we want to classify pictures of cats or dogs. How would you do this?

Look for features of cats or dogs in the image and use for decision.

- Example: cats have cat-like faces, dogs have dog-like faces.

- How do you determine what is a "cat-like" face vs a "dog-like" face?

Look for features of "cat-like" faces and "dog-like" faces.
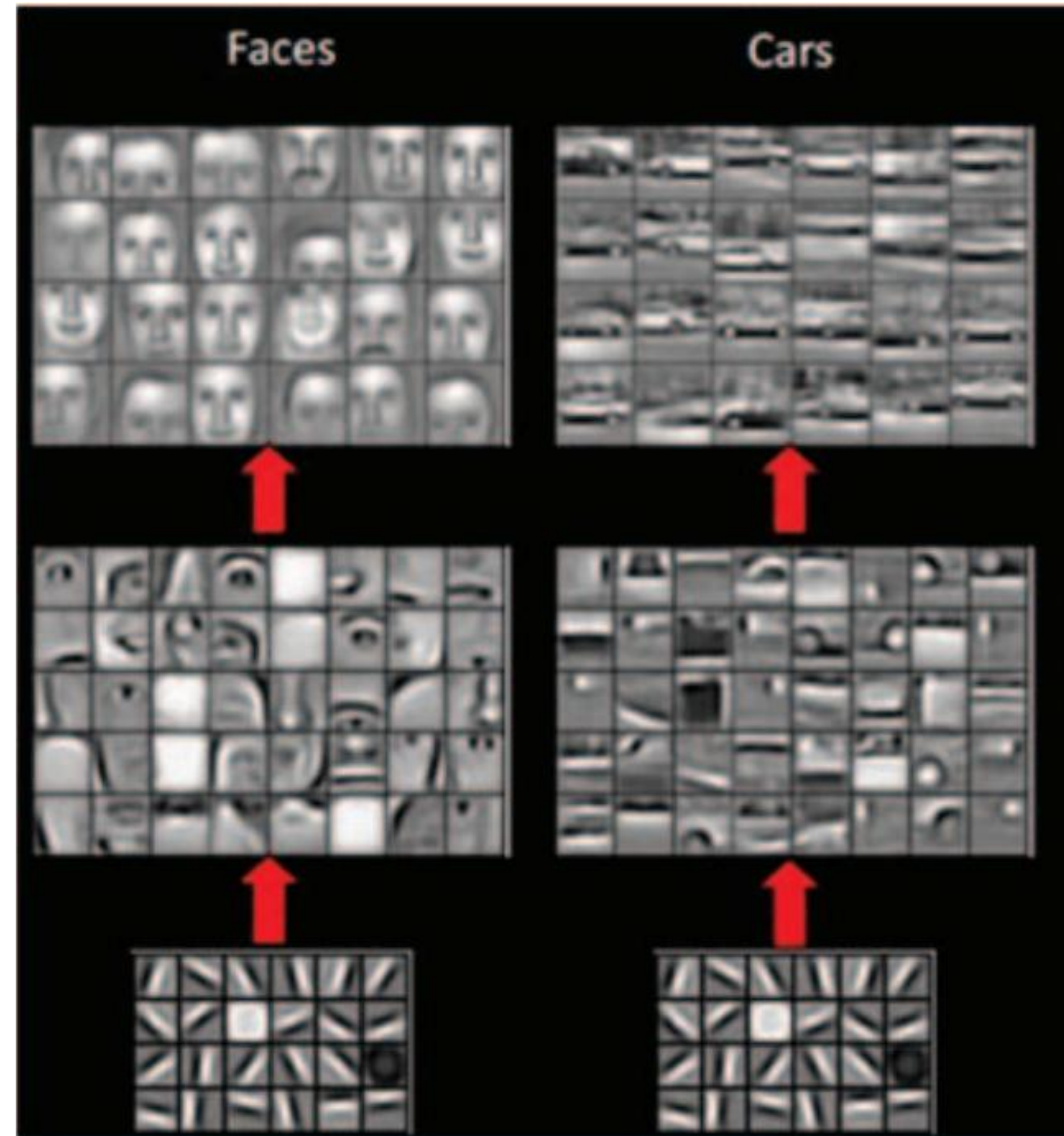
- Example: Dogs have longer snouts.

- How do you determine what is a long snout?

# Feature Learning

Later layers recognize complete objects

Middle layers recognize parts of objects

Early layers recognize simple patterns



Adaptive Neuron Apoptosis for Accelerating Deep Learning on Large Scale Systems. Seigel et al. 2016.

# Convolutional Neural Networks Examples

# Evolution of neural net architectures



LeNet
(1989)

AlexNet
(2012)

VGG
(2014)

Inception Net
(2014)

ResNet
(2015)

DenseNet
(2016)

c : Channel-wise concatenation

# LeNet Architecture
## (first convolutional neural net; 1989)



convolution

pooling

convolution

pooling

full

full

Gauss

32x32 image

6@28x28
C1 feature map

6@14x14
S2 feature map

16@10x10
C3 feature map

16@5x5
S4 feature map

120 - F5 full

84 - F6 full

10 - Out

34

# Handwritten Digit Recognition

# MNIST

- Centered and scaled
- 50,000 training data
- 10,000 test data
- 28 x 28 images
- 10 classes

Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, 1998 Gradient-based learning applied to document recognition

# LeNet Architecture

Gradient-based learning applied to document recognition,
by Y. LeCun, L. Bottou, Y. Bengio and P. Haffner

# LeNet in Pytorch

```python
def __init__(self):
    super(LeNet5, self).__init__()
    # Convolution (In LeNet-5, 32x32 images are given as input. Hence padding of 2 is done below)
    self.conv1 = torch.nn.Conv2d(in_channels=1, out_channels=6, kernel_size=5, stride=1, padding=2, bias=True)
    # Max-pooling
    self.max_pool_1 = torch.nn.MaxPool2d(kernel_size=2)
    # Convolution
    self.conv2 = torch.nn.Conv2d(in_channels=6, out_channels=16, kernel_size=5, stride=1, padding=0, bias=True)
    # Max-pooling
    self.max_pool_2 = torch.nn.MaxPool2d(kernel_size=2)
    # Fully connected layer
    self.fc1 = torch.nn.Linear(16*5*5, 120)    # convert matrix with 16*5*5 (= 400) features to a matrix of 120 features (columns)
    self.fc2 = torch.nn.Linear(120, 84)        # convert matrix with 120 features to a matrix of 84 features (columns)
    self.fc3 = torch.nn.Linear(84, 10)         # convert matrix with 84 features to a matrix of 10 features (columns)
```
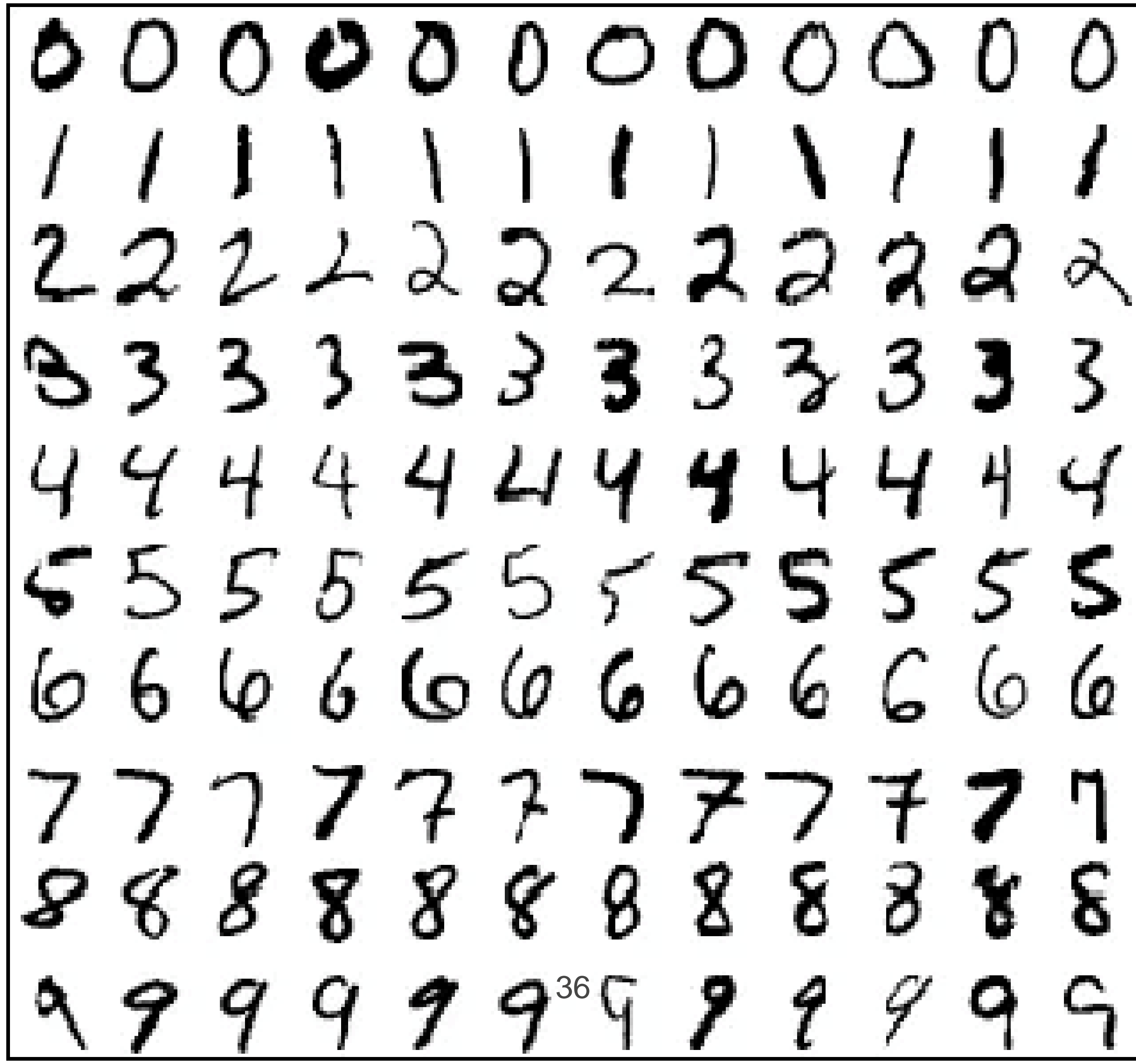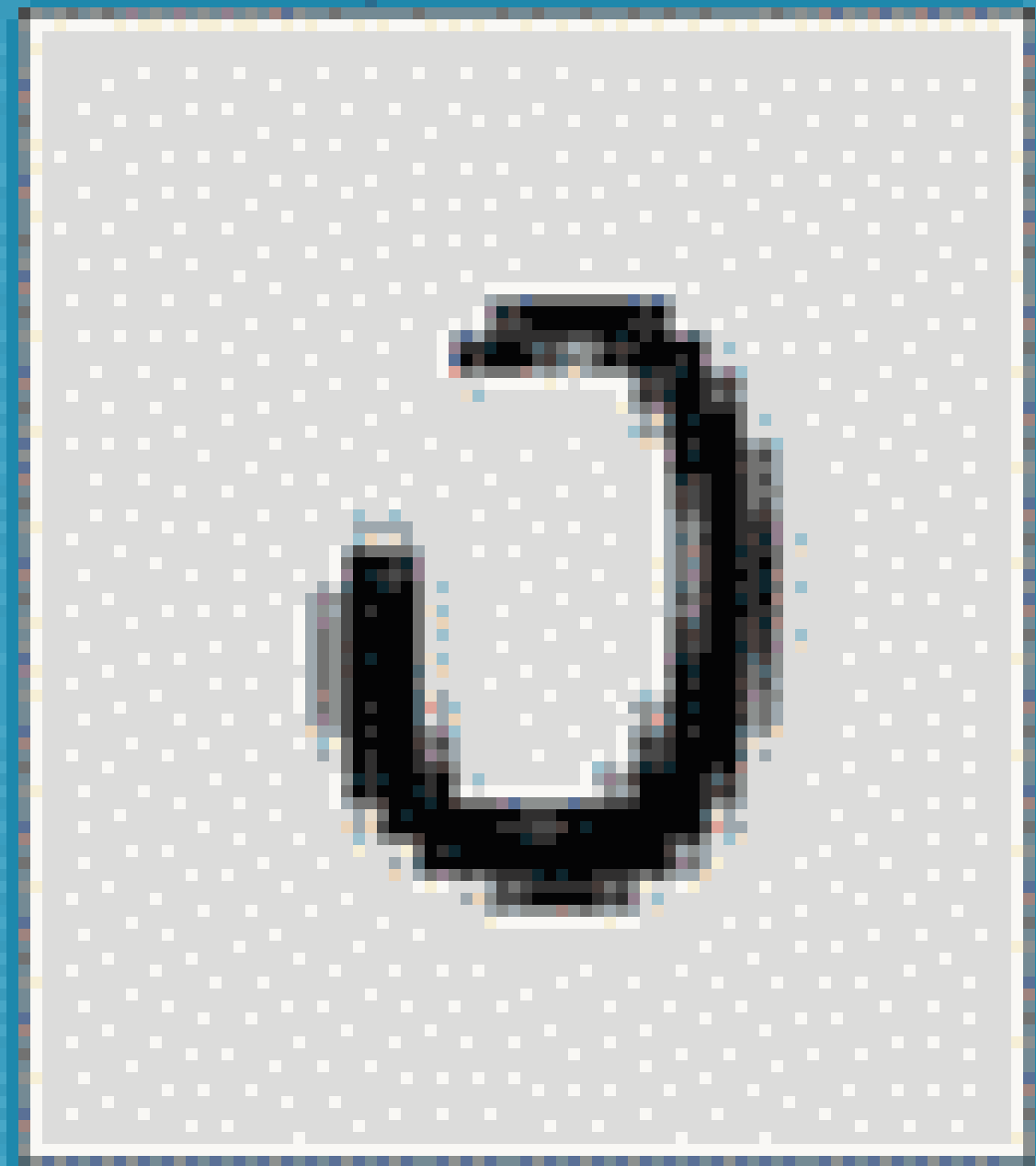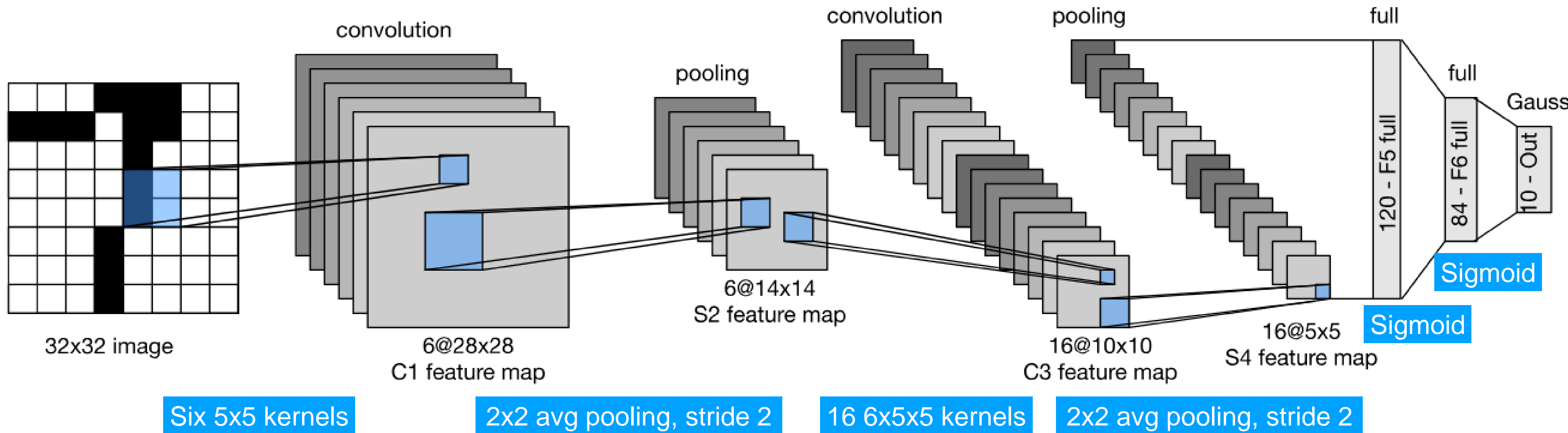
https://github.com/bollakarthikeya/LeNet-5-PyTorch/blob/master/lenet5_gpu.py

```python
def forward(self, x):
    # convolve, then perform ReLU non-linearity
    x = torch.nn.functional.relu(self.conv1(x))
    # max-pooling with 2x2 grid
    x = self.max_pool_1(x)
    # convolve, then perform ReLU non-linearity
    x = torch.nn.functional.relu(self.conv2(x))
    # max-pooling with 2x2 grid
    x = self.max_pool_2(x)
    # first flatten 'max_pool_2_out' to contain 16*5*5 columns
    # read through https://stackoverflow.com/a/42482819/7551231
    x = x.view(-1, 16*5*5)
    # FC-1, then perform ReLU non-linearity
    x = torch.nn.functional.relu(self.fc1(x))
    # FC-2, then perform ReLU non-linearity
    x = torch.nn.functional.relu(self.fc2(x))
    # FC-3
    x = self.fc3(x)

    return x
```

# LeNet in Pytorch

# AlexNet



Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*.
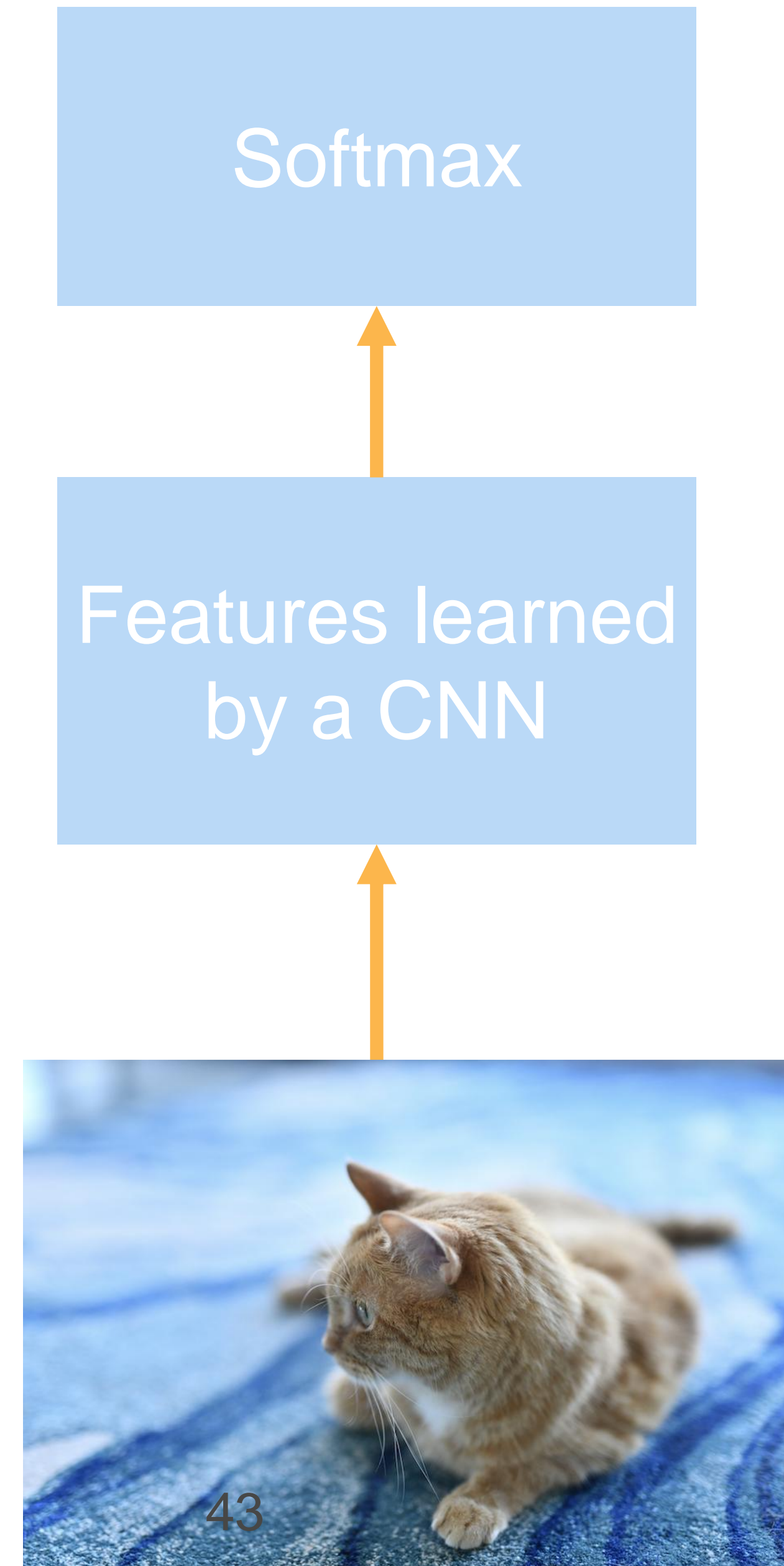
41

Deng et al. 2009

# **AlexNet**

- AlexNet won ImageNet competition in 2012

- Deeper and bigger LeNet

- Paradigm shift for computer vision

# AlexNet Architecture

AlexNet
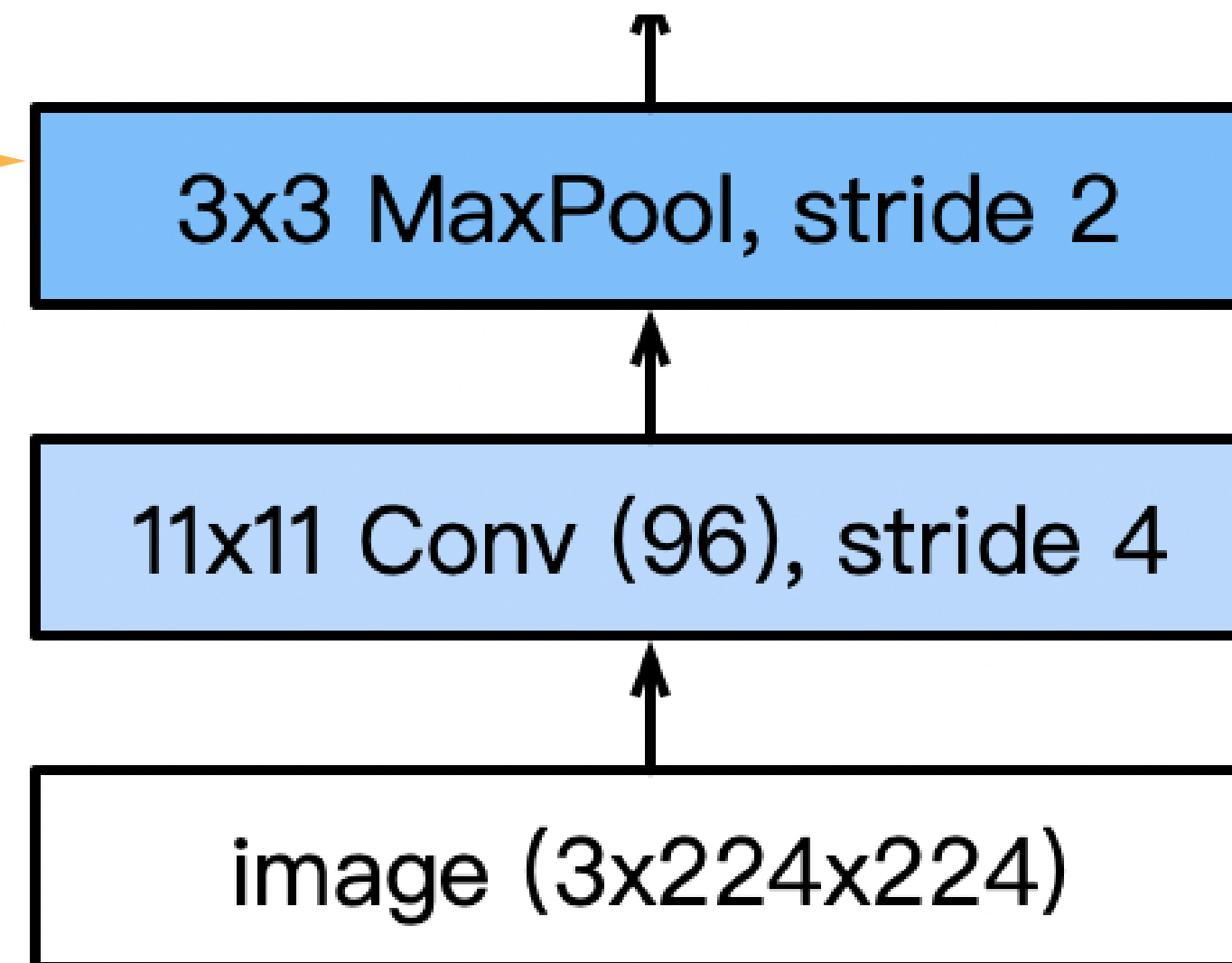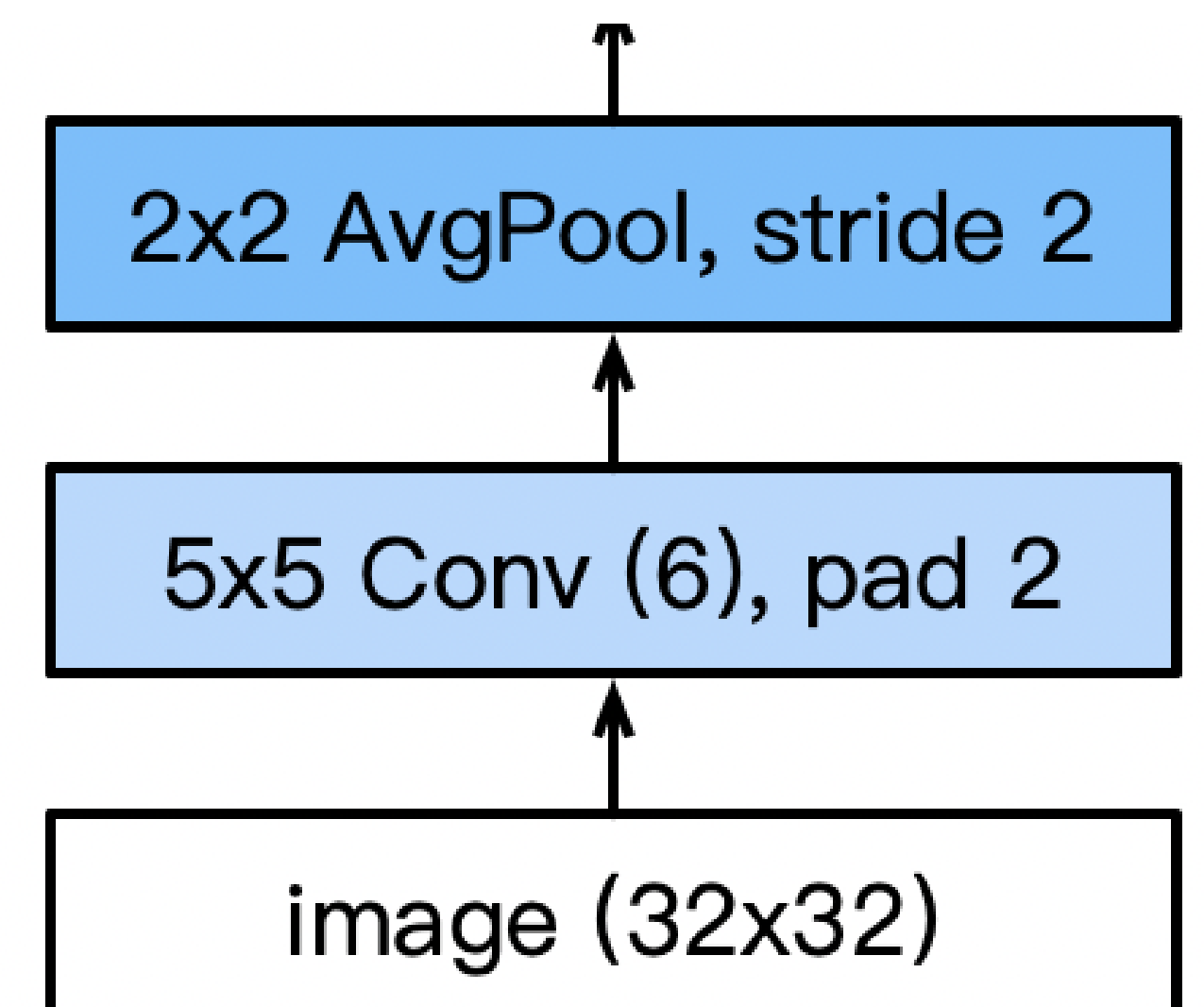
LeNet

Larger pool size

3x3 MaxPool, stride 2

2x2 AvgPool, stride 2

11x11 Conv (96), stride 4

5x5 Conv (6), pad 2

image (3x224x224)

image (32x32)

Larger kernel size, stride because of the increased image size, and more output channels.

# AlexNet Architecture



AlexNet

3x3 MaxPool, stride 2

3x3 Conv (384), pad 1

3x3 Conv (384), pad 1

3x3 Conv (384), pad 1

3x3 MaxPooling, stride 2

5x5 Conv (256), pad 2

3 additional convolutional layers

More output channels.

LeNet

2x2 AvgPool, stride 2

5x5 Conv (16)

45

# AlexNet Architecture



AlexNet

LeNet

1000 classes output

Increase hidden size from 120 to 4096

| AlexNet |
| --- |
| Dense (1000) |
| Dense (4096) |
| Dense (4096) |

| LeNet |
| --- |
| Dense (10) |
| Dense (84) |
| Dense (120) |

# More Differences…

- Change activation function from sigmoid to ReLu (no more vanishing gradient)



sigmoid

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

Saturating gradients

# More Differences…

- Change activation function from sigmoid to ReLu (no more vanishing gradient)
- Data augmentation

# Complexity

| | #parameters | |
|---|---|---|
| | **AlexNet** | **LeNet** |
| **Conv1** | 35K | 150 |
| **Conv2** | 614K | 2.4K |
| **Conv3-5** | 3M | |
| **Dense1** | 26M | 0.048M |
| **Dense2** | 16M | 0.01M |
| **Total** | 46M | 0.06M |

Dense (1000)

↑

Dense (4096)

↑

Dense (4096)

↑

Max Pooling

↑

3x3 Conv (384)

↑

3x3 Conv (384)

↑

3x3 Conv (384)

↑

Max Pooling

↑

5x5 Conv (256)

↑

Max Pooling

↑

11x11 Conv (96), stride 4

↑

image (224x224)

# Complexity

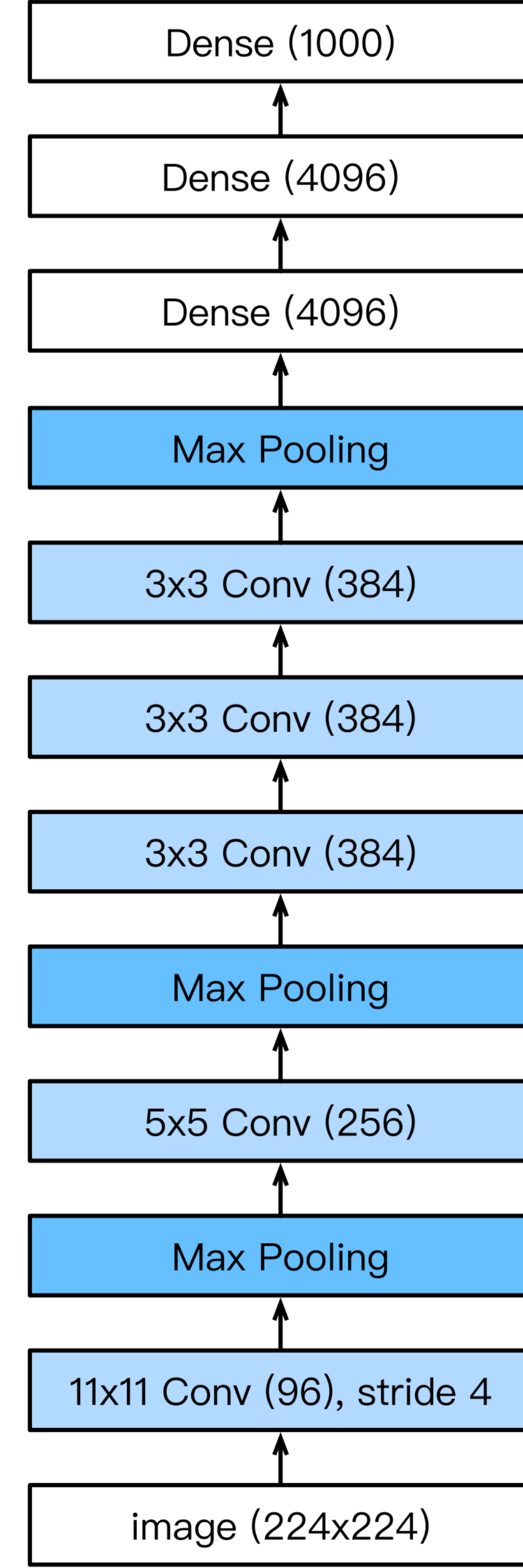| | #parameters | |
|---|---|---|
| | **AlexNet** | **LeNet** |
| **Conv1** | 35K | 150 |
| **Conv2** | 614K | 2.4K |
| **Conv3-5** | 3M | |
| **Dense1** | 26M | 0.048M |
| **Dense2** | 16M | 0.01M |
| **Total** | 46M | 0.06M |

11x11x3x96=35k

50

| Dense (1000) |
|---|
| Dense (4096) |
| Dense (4096) |
| Max Pooling |
| 3x3 Conv (384) |
| 3x3 Conv (384) |
| 3x3 Conv (384) |
| Max Pooling |
| 5x5 Conv (256) |
| Max Pooling |
| 11x11 Conv (96), stride 4 |
| image (224x224) |

ImageNet Top-5 Classification Error (%)

# AlexNet

3x3 MaxPool, stride 2

11x11 Conv (96), stride 4

image (3x224x224)

Each Conv1 kernel is 3x11x11, can be visualized as an RGB patch:



(c)

[Visualizing and Understanding Convolutional Networks. M Zeiler & R Fergus 2013]

Which of the following are true about AlexNet? Select all that apply.

A. AlexNet contains 8 conv/fc layers. The first five are convolutional layers.

B. The last three layers are fully connected layers.

C. some of the convolutional layers are followed by [max-pooling](layers) (layers).

D. AlexNet achieved excellent performance in the 2012 ImageNet challenge.

Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks.

*Advances in neural information processing systems* (pp. 1097–1105).

# Which of the following are true about AlexNet? Select all that apply.

A. AlexNet contains 8 conv/fc layers. The first five are convolutional layers.

B. The last three layers are fully connected layers.

C. some of the convolutional layers are followed by max-pooling (layers).

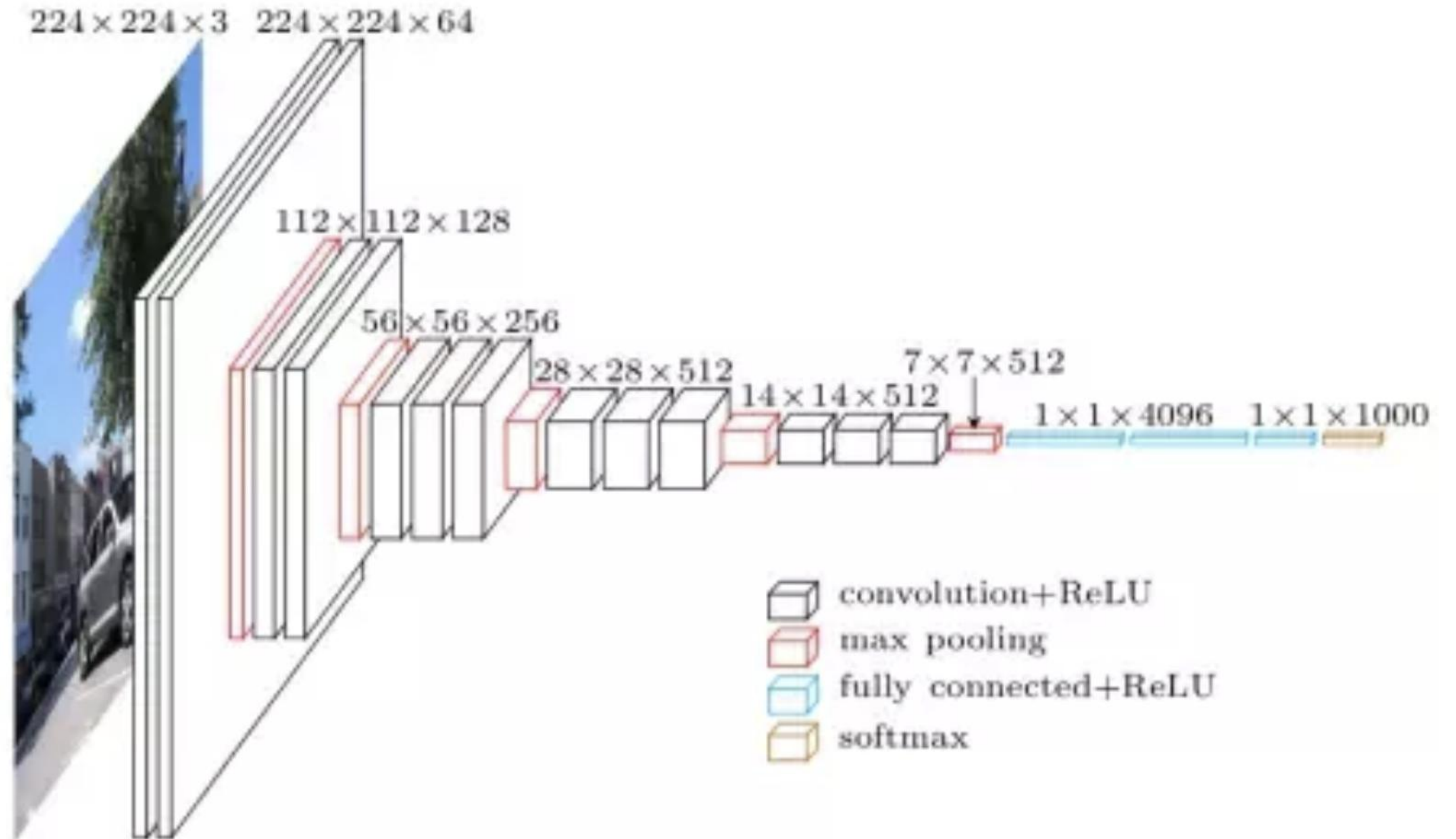D. AlexNet achieved excellent performance in the 2012 ImageNet challenge.

All options are true!

Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks.

*Advances in neural information processing systems* (pp. 1097–1105).

# VGG



VGG Block: Multiple convolution layers followed by pooling.

# Progress

- LeNet (1995)
  - 2 convolution + pooling layers
  - 2 hidden dense layers
- AlexNet
  - Bigger and deeper LeNet
  - ReLu, preprocessing
- VGG
  - Bigger and deeper AlexNet (repeated VGG blocks)

# Which of the following statement is True for the success of deep models?

- Better design of the neural networks
- Large scale training dataset
- Available computing power
- All of the above

# Which of the following statement is True for the success of deep models?

- Better design of the neural networks
- Large scale training dataset
- Available computing power
- All of the above

# Suggested Reading

Example using PyTorch:

https://pytorch.org/tutorials/beginner/blitz/cifar10_tutorial.html

# Summary of today

- Reviewed (some of) convolutional computations.

  - 2D convolutions, multiple input channels, pooling.

- Shown how convolutions are used as layers in a (deep) neural network.

- Built intuition for output of convolutional layers.

- Overviewed the evolution of deeper convolutional networks

# **Acknowledgement**: