# Deep Learning
# CS60010

**Abir Das**

Assistant Professor
Computer Science and Engineering Department
Indian Institute of Technology Kharagpur

http://cse.iitkgp.ac.in/~adas/

# Announcements

- Class Test 1 is graded.
- Solutions are also uploaded.
- Login to the moodle server where you gave the exam. In gradebooks you will get to see the marks.
- Ignore the percentages. Just the marks obtained will be taken to compute your grades at the end. The percentages shown here are meaningless.
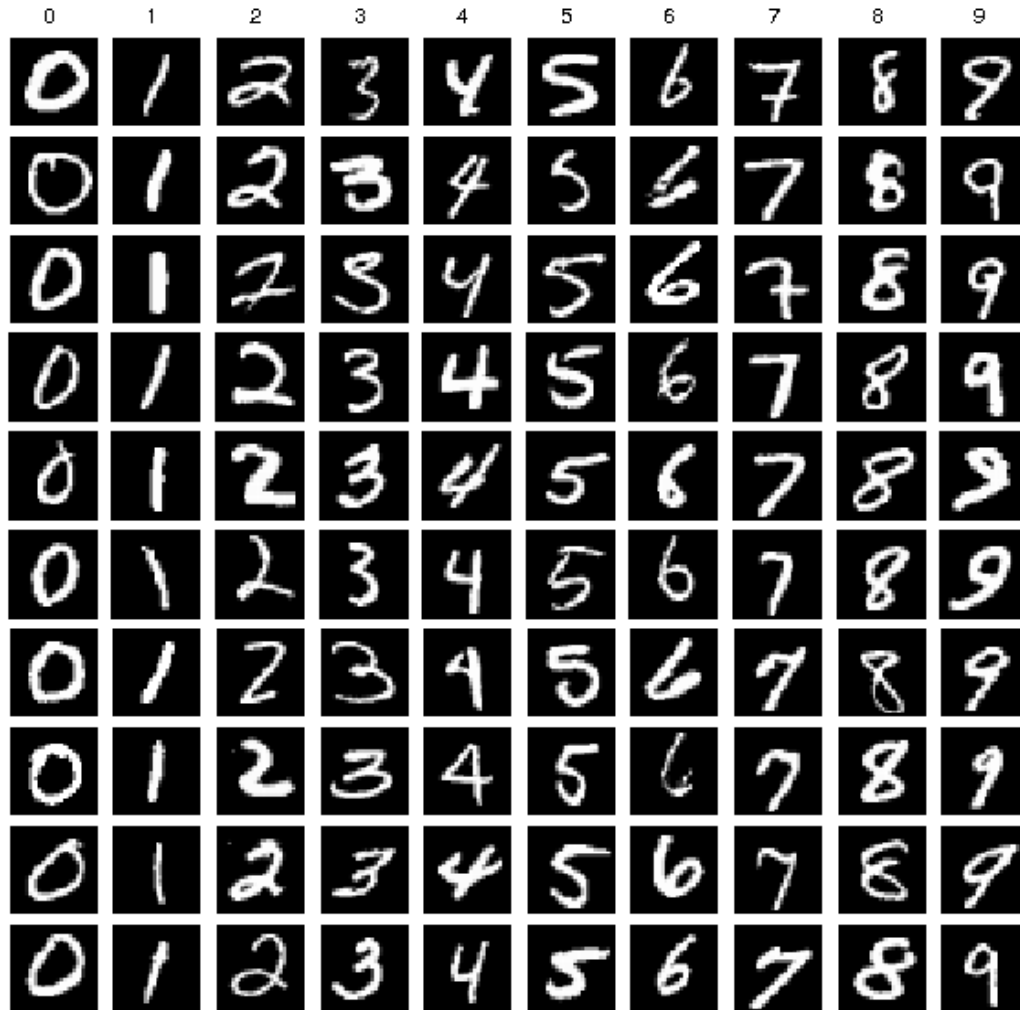
# Agenda

The Building Blocks of Convolutional Neural Networks/CNNs/ConvNets
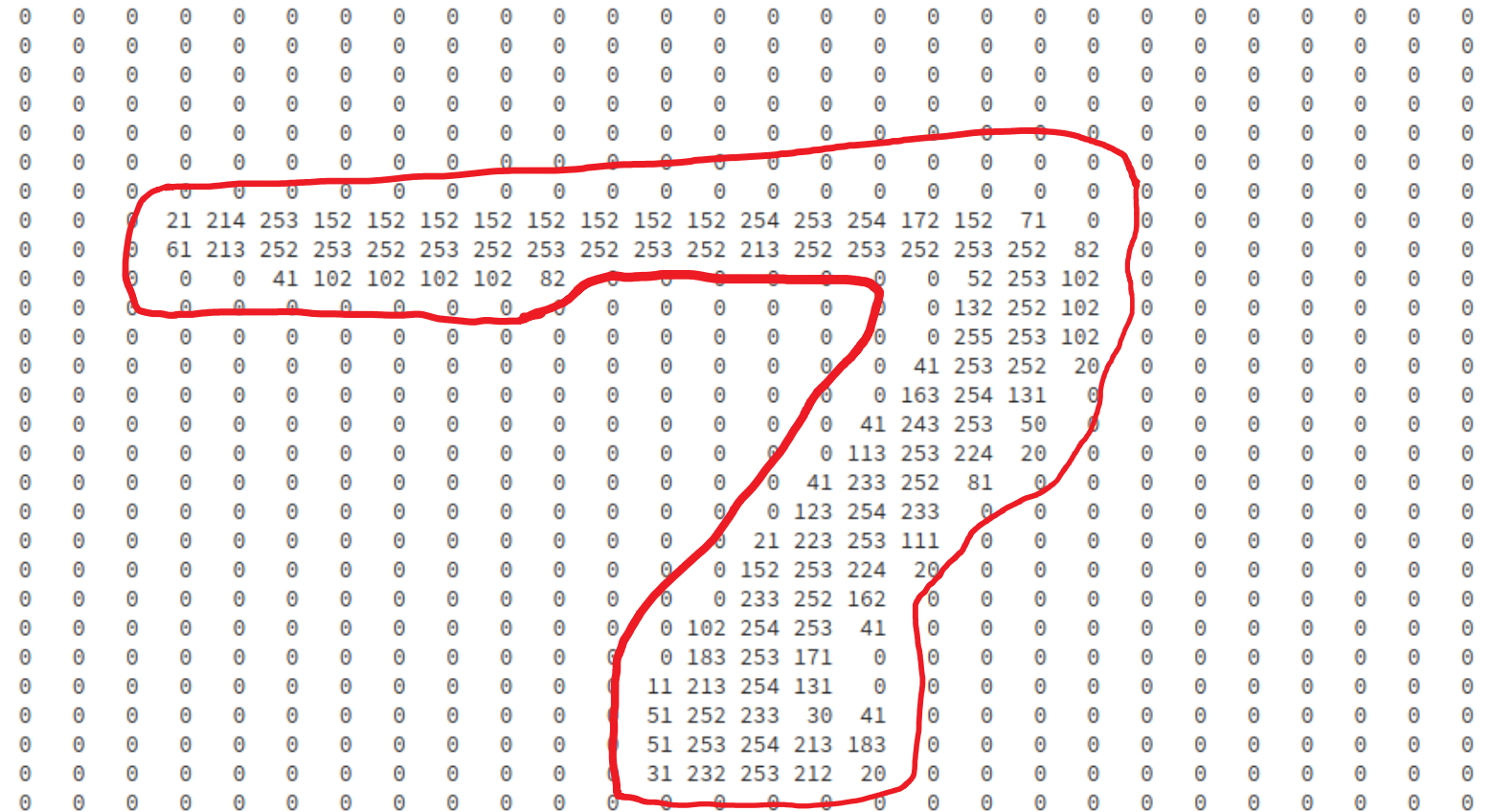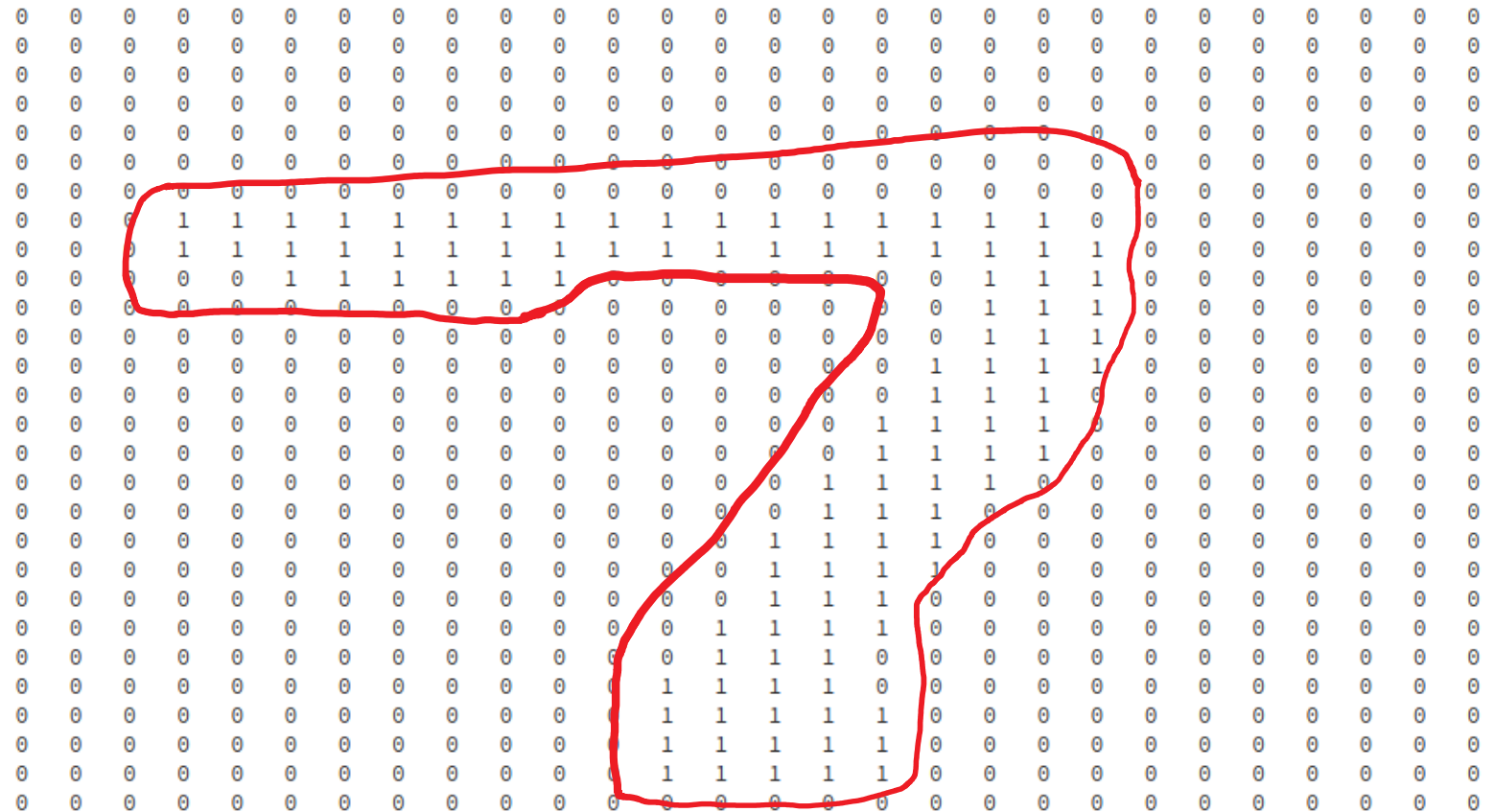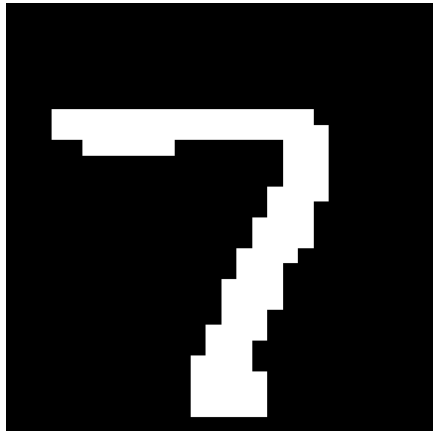
# Importance of MNIST

# MNIST



- database of handwritten digits
- 10 classes
  - Training set 60,000 images
  - Test set of 10,000 images
- Greyscale images of size 28x28
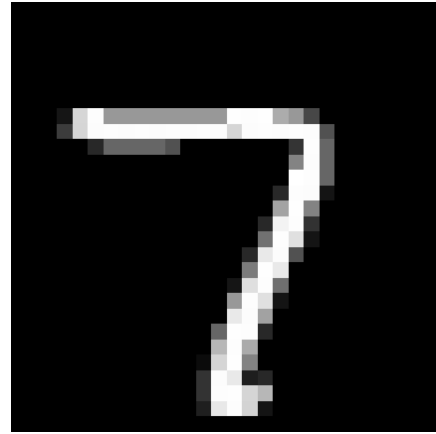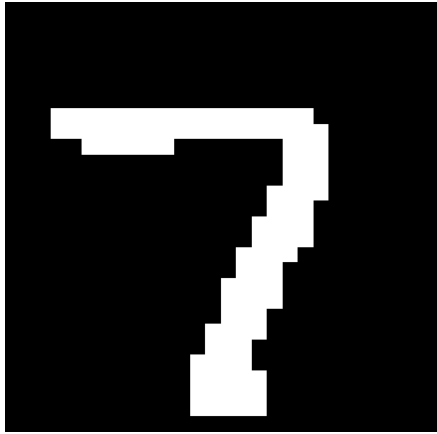- Often treated as `Hello World' for any ML/DL practioner
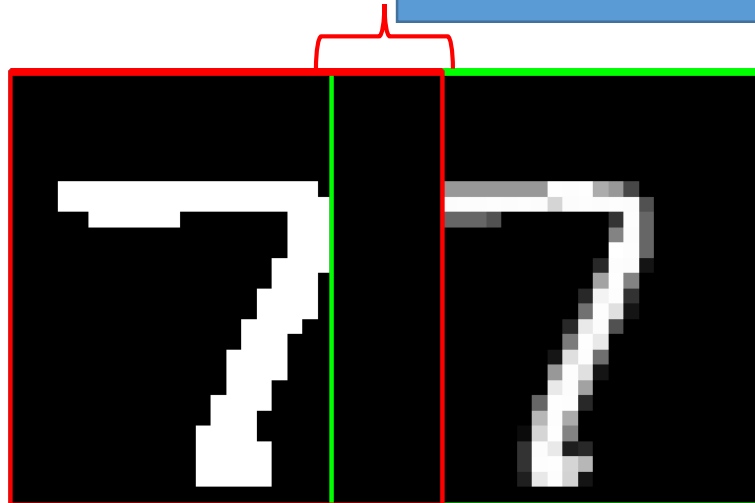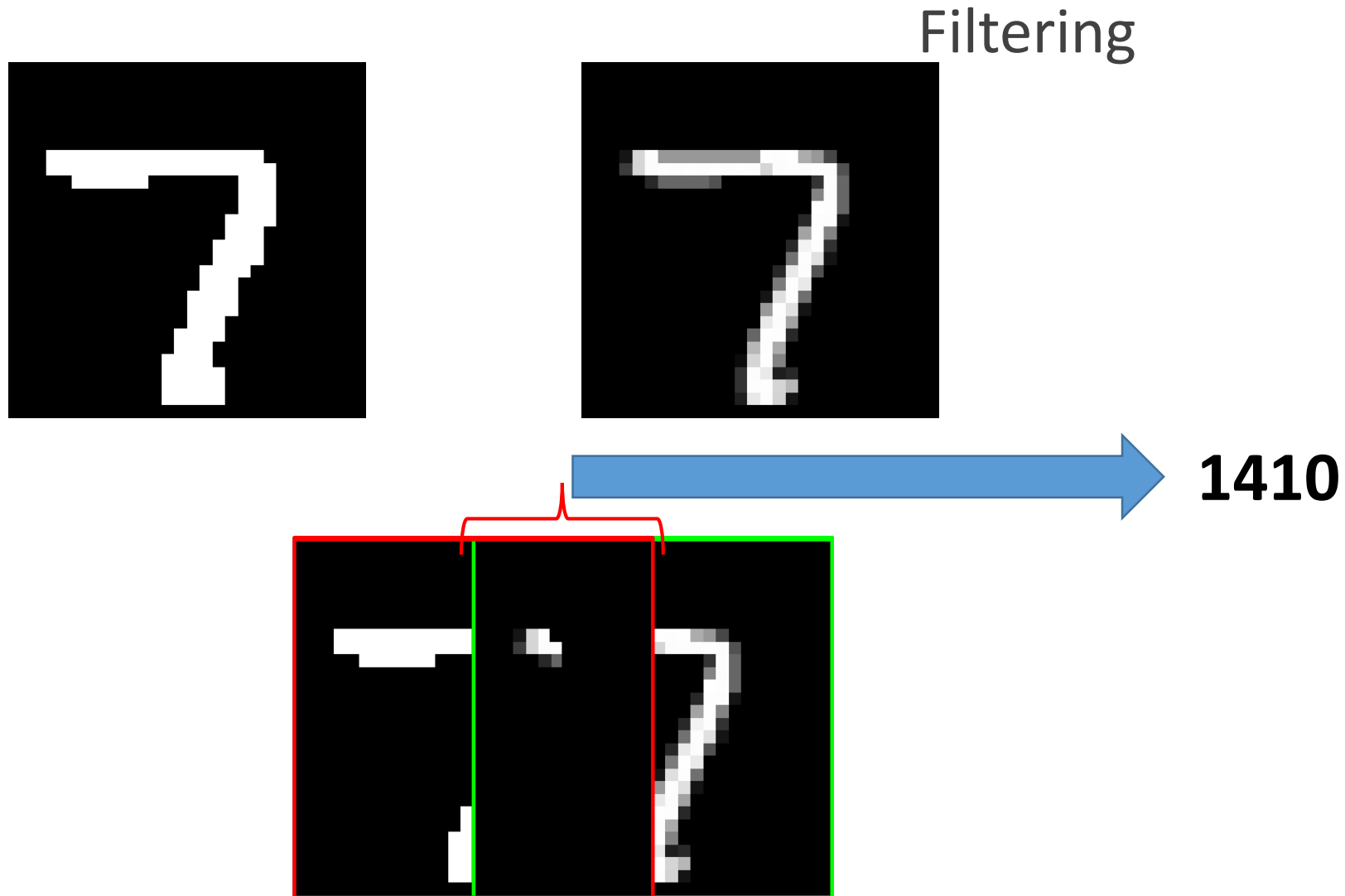
Image taken from Researchgate

# Example from MNIST

# 7 Like Window?

# Filtering



**0**

# Filtering



**1410**

# Filtering



**4098**

# Filtering



**18604**

# Filtering

**4331**

# Filtering



**1589**

# Filtering



**0**

# Classification by Matching Filters

Filters

Test Images



But what if the test image is a little
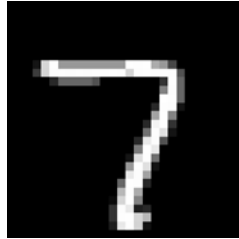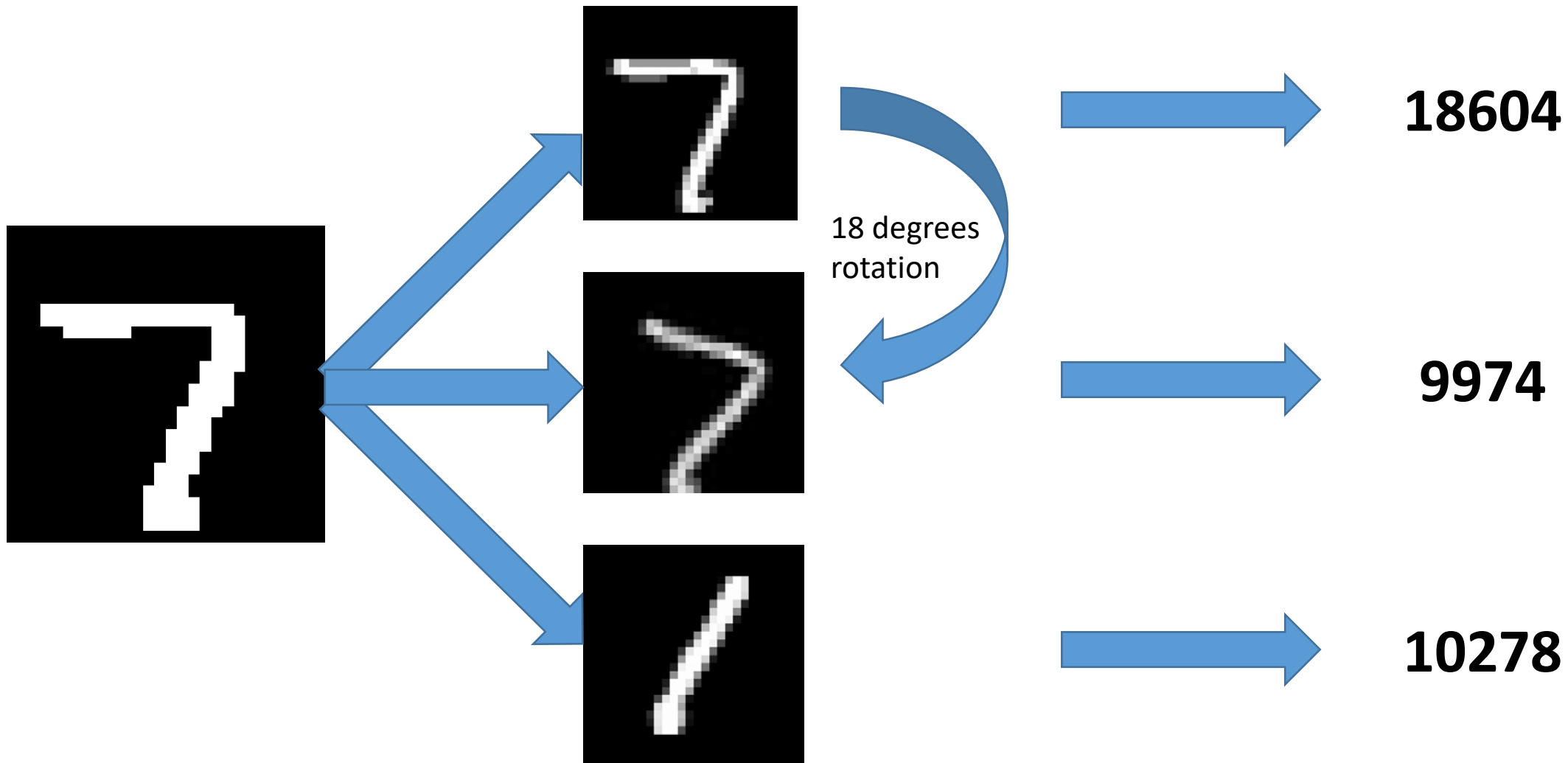Rotated          (or)
Skewed           (or)
Zoomed out   and so on

# Effect of Slight Rotation
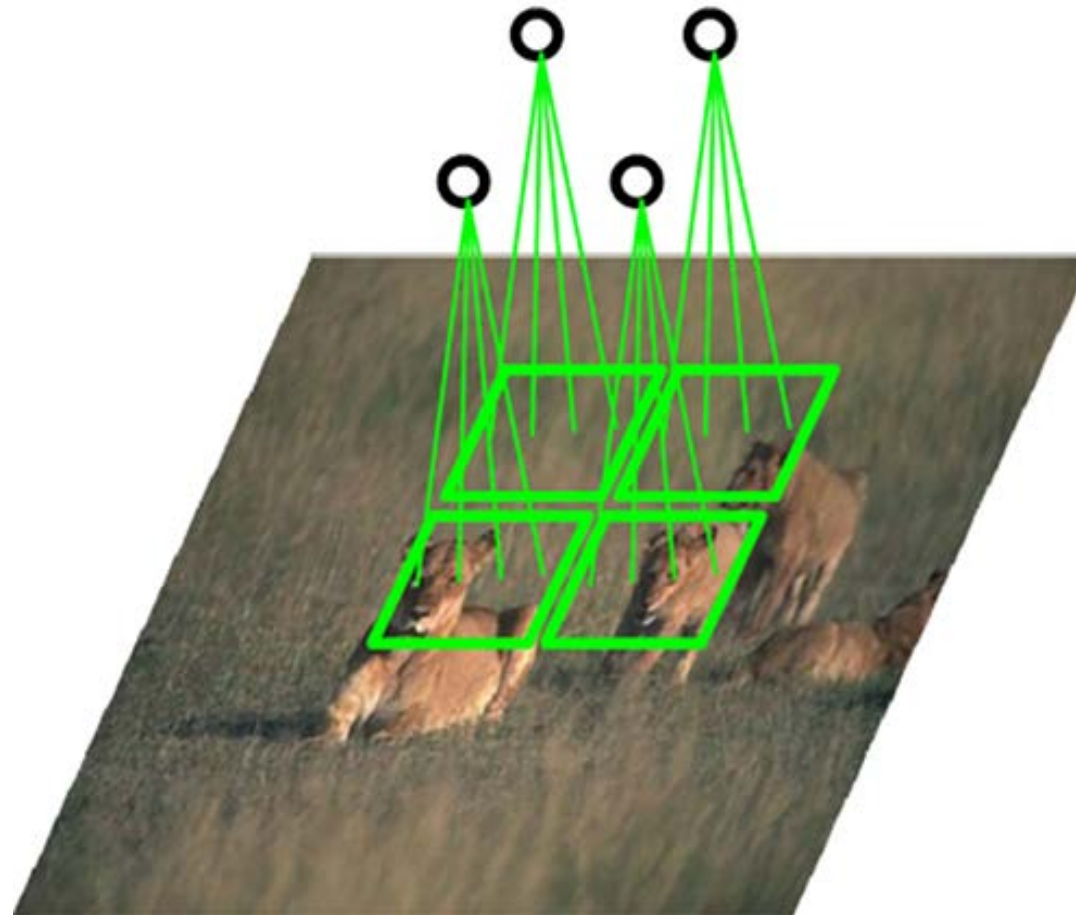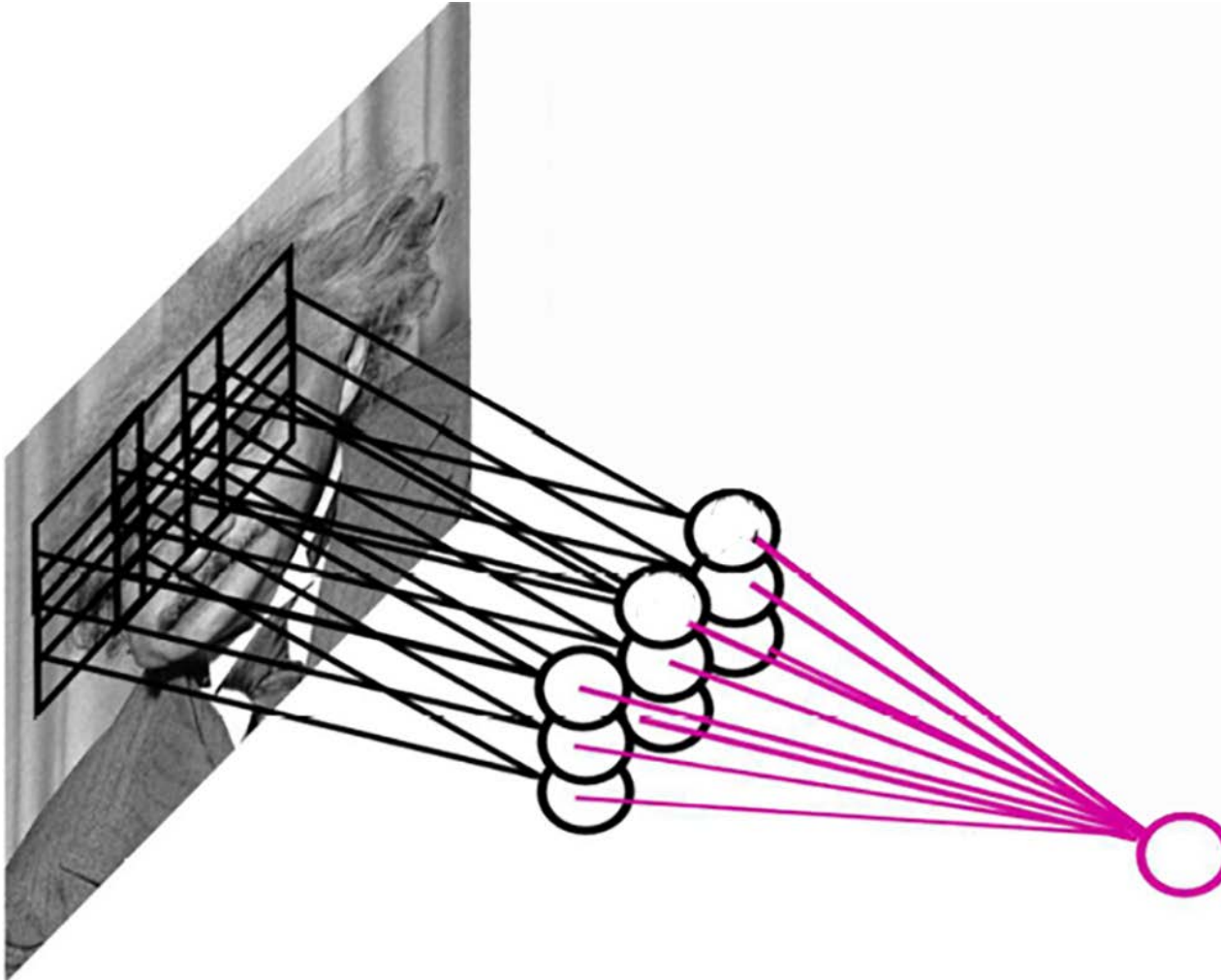


18 degrees rotation

18604

9974

10278

# Convnets (Fukushima, LeCun, Hinton)

# Convnets (Fukushima, LeCun, Hinton)

# Pooling



By "pooling" (e.g., max or average) filter responses, we gain robustness to the exact spatial location of features.

Image courtesy: Nando de Freitas, University of Oxford

# Convolution



Image courtesy:Vincent Dumoulin

# Convolution with Zero Padding

# Convolution with Strides

Convolving a 3x3 kernel over a 5x5 input using 2x2 strides

Image courtesy:Vincent Dumoulin

# Convolution with Strides and Zero Padding

Convolving a 3x3 kernel over a 5 x5 input using 1 x1 strides and half padding



Image courtesy:Vincent Dumoulin

# Inputs Generally have Multiple Channels



Image courtesy: Vincent Dumoulin

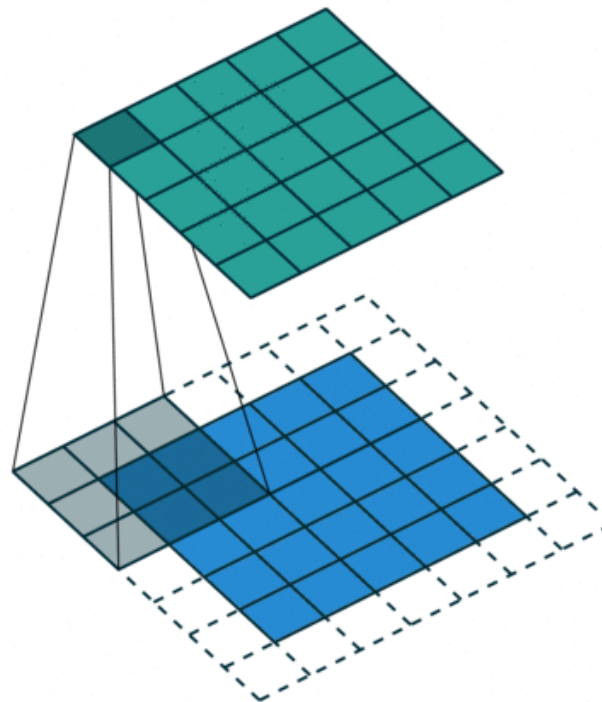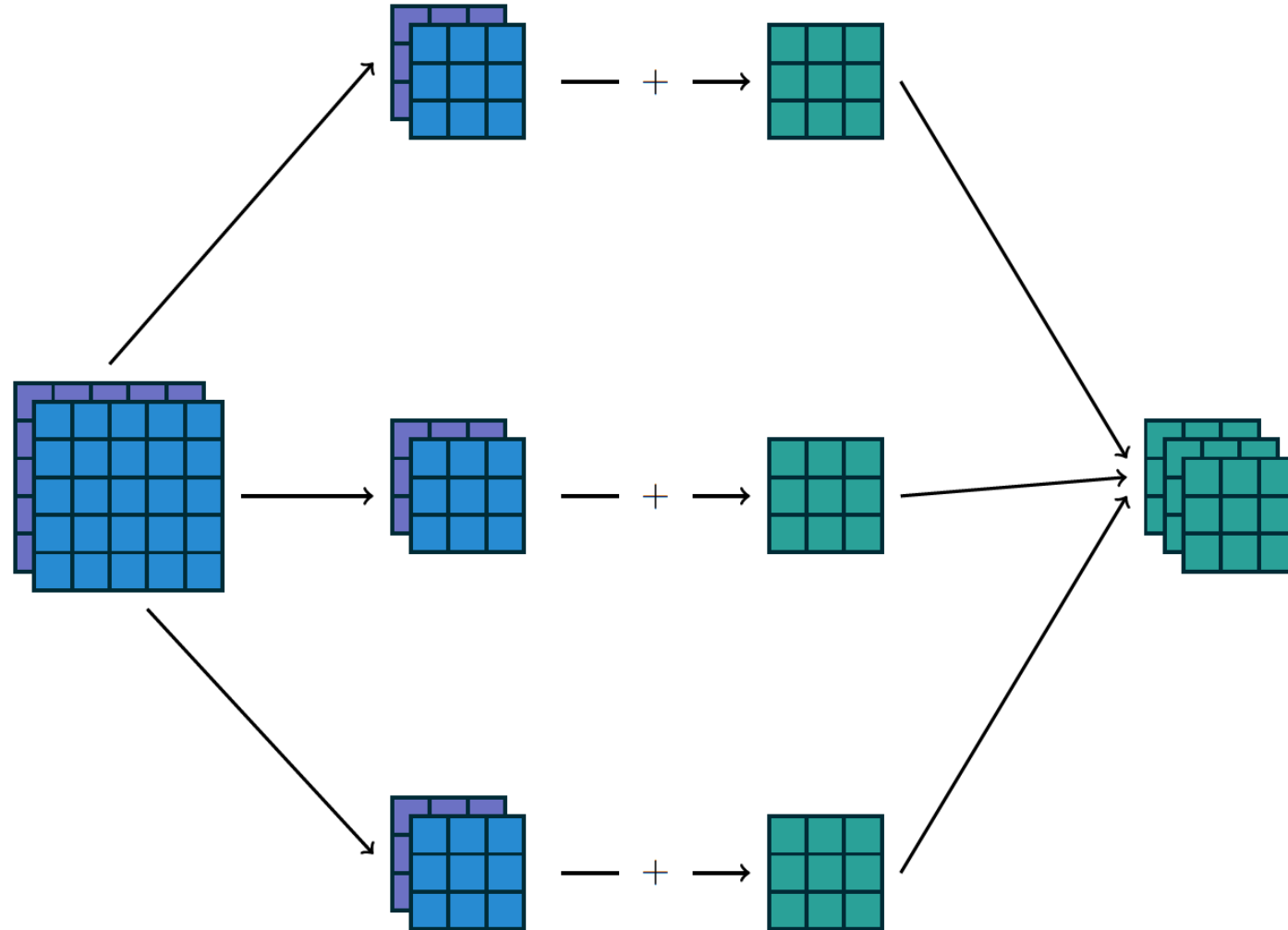# Convolution Arithmetic

(For simplicity we are assuming square image and filter/kernel)

Image width = image height = $w$

Filter width = Filter height = $k$

Stride = $s$

Output size = $\left\lfloor \dfrac{w-k}{s} \right\rfloor + 1$

Padding = $p$ → This means image dimension becomes $w + 2p$

So, output size = $\left\lfloor \dfrac{w+2p-k}{s} \right\rfloor + 1$
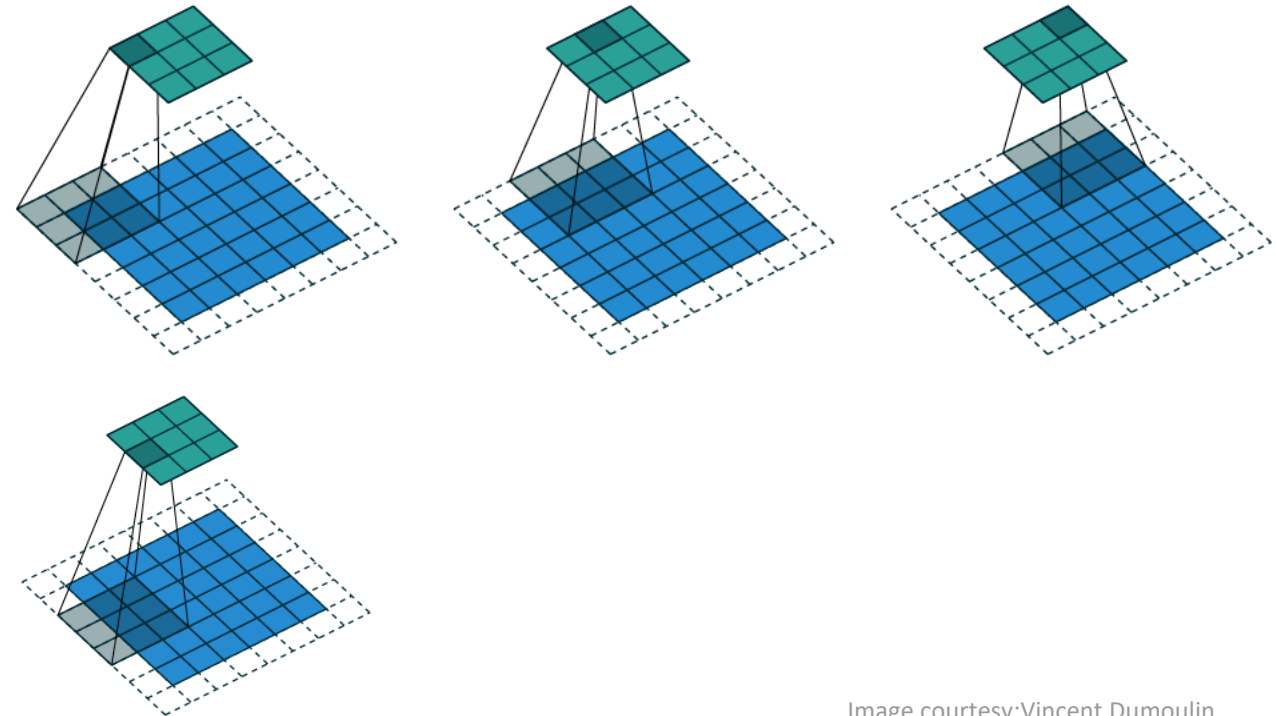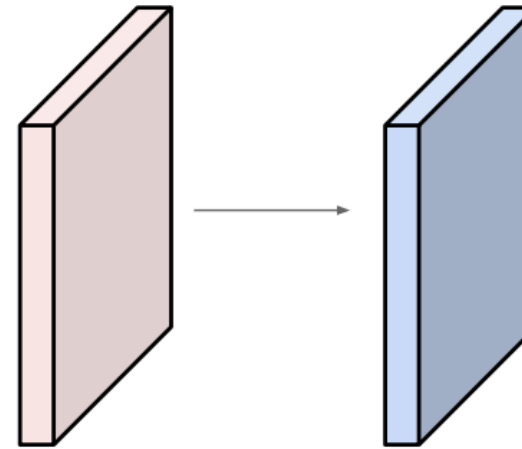
Note the box function

$w = 6, p = 1, k = 3, s = 2$



Image courtesy:Vincent Dumoulin

# Convolution Arithmetic

Input volume: 32x32x3 [w, h, c].
64 filters of size 3x3 [k, k] with
stride 2 [s], pad 1 [p]



What is the output feature map size?

And What is the number of parameters in this convolution layer?

$$\left\lfloor \frac{32 + 2 * 1 - 3}{2} \right\rfloor + 1 = 16$$

So, $16 \times 16 \times 64$ [w, h, c]

$64 \times 3 \times 3 \times 3$ [c_out, w, h, c_in] = 1728

# Pooling



3x3 max-pooling on
5x5 input with
1x1 stride

Image courtesy:Vincent Dumoulin

# Pooling Arithmetic

(For simplicity we are assuming square input and max pooling kernel)

Input width = Input height = $w$

Filter width = Filter height = $k$

Stride = $s$
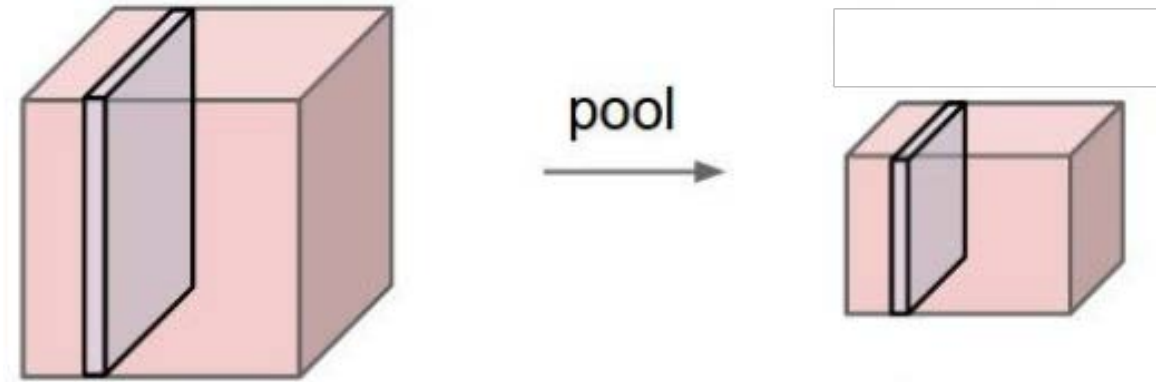
Output size = $\left\lfloor \dfrac{w-k}{s} \right\rfloor + 1$

Input volume: 32x32x3 [w, h, c].
Max-pooling kernel of size 2x2 [k, k] with stride 2 [s]
What is the output feature map size?



And What is the number of parameters in this pooling layer?

0

$\left\lfloor \dfrac{32-2}{2} \right\rfloor + 1 = 16$    So, $16 \times 16 \times 3$ [w, h, c]

Image courtesy: Andrej Karpathy

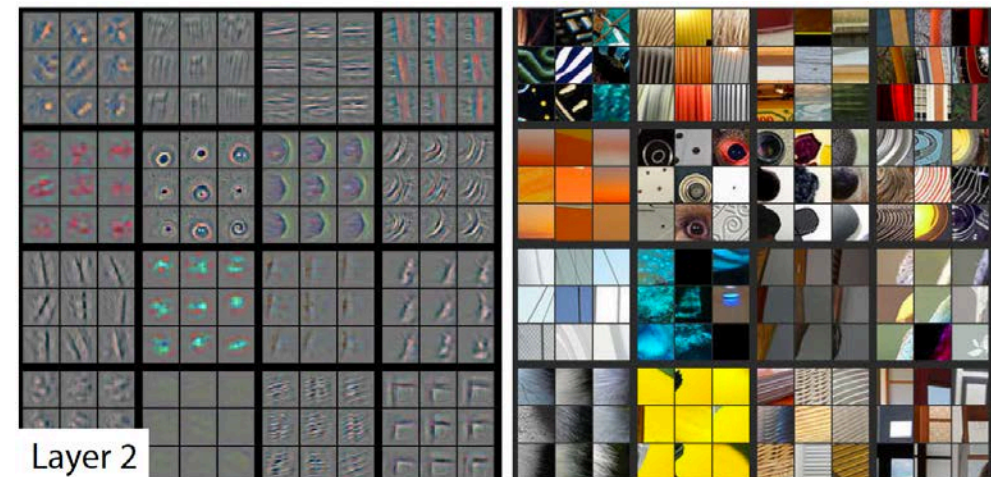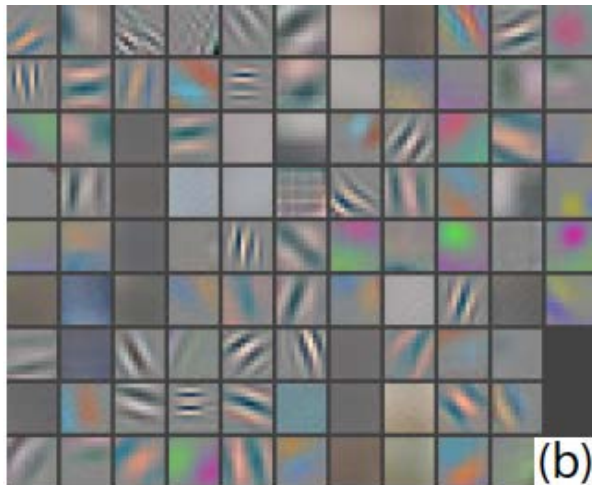# Visualizations

AlexNet (2012)



First layer (CONV1): 96  11x11 filters



(b)

Layer 2

Image courtesy: Zeiler, Fergus, 2013

# Visualizations



Layer 3

Layer 4

Image courtesy: Zeiler, Fergus, 2013