



Deep Learning

CS60010

Abir Das

Assistant Professor

Computer Science and Engineering Department
Indian Institute of Technology Kharagpur

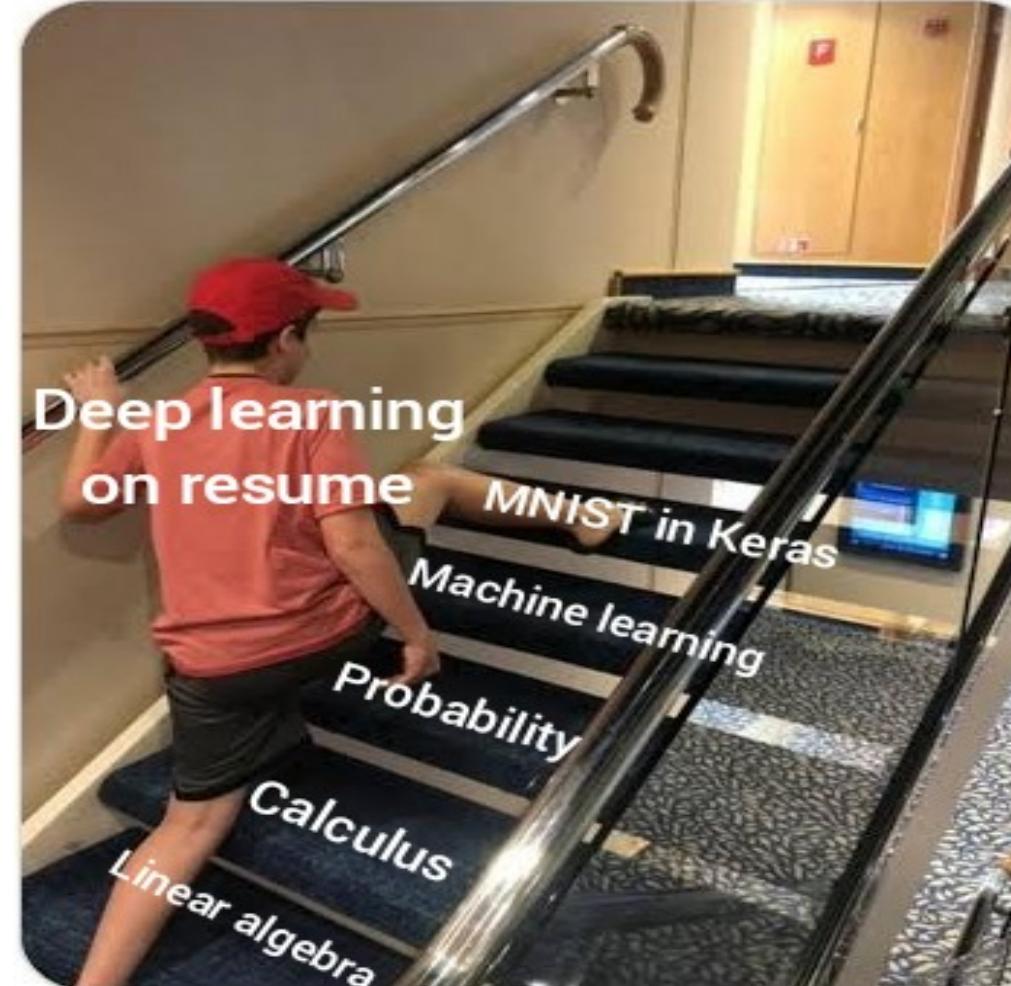
<http://cse.iitkgp.ac.in/~adas/>



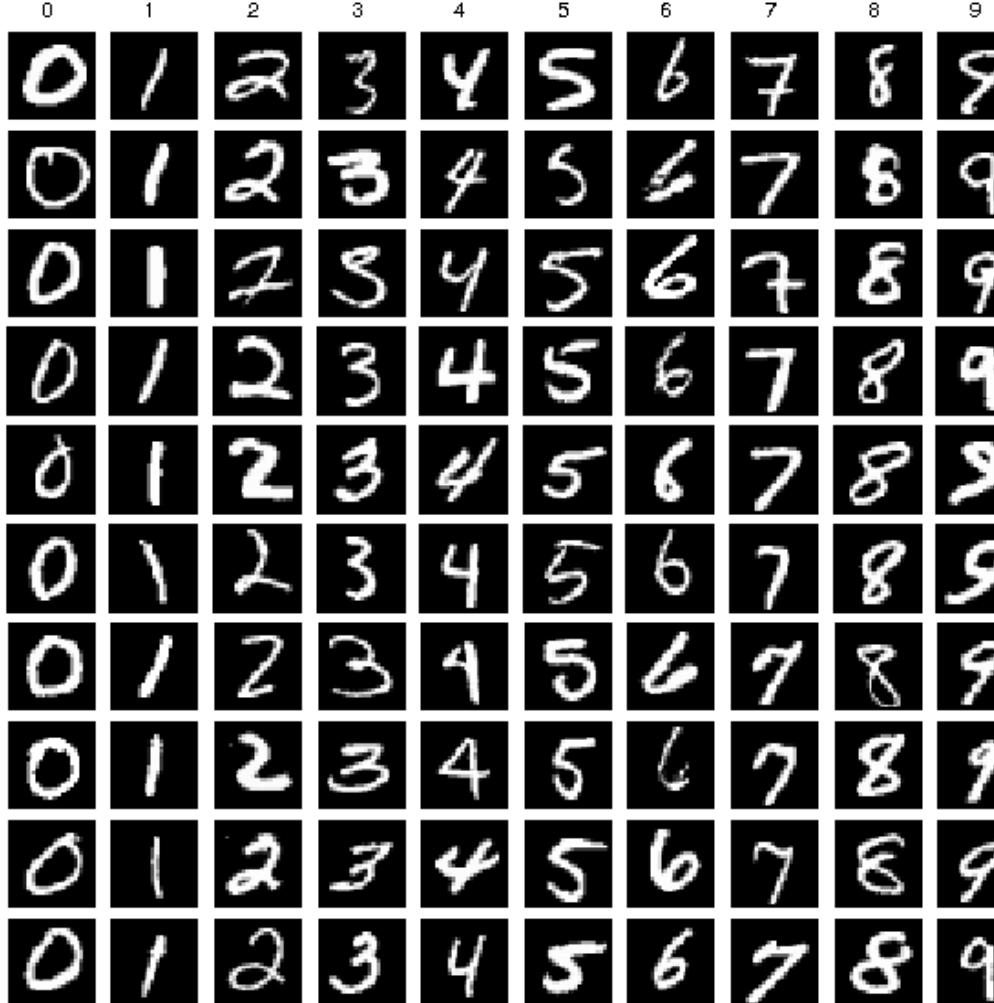
Agenda

The Building Blocks of Convolutional Neural Networks/CNNs/ConvNets

Importance of MNIST



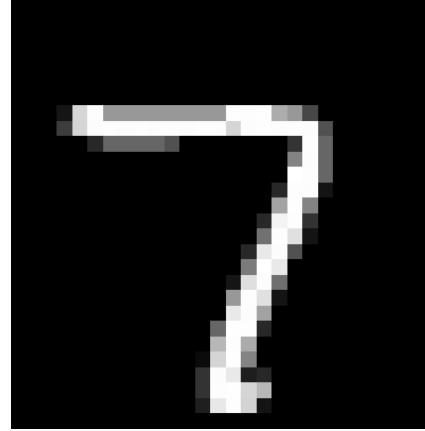
MNIST



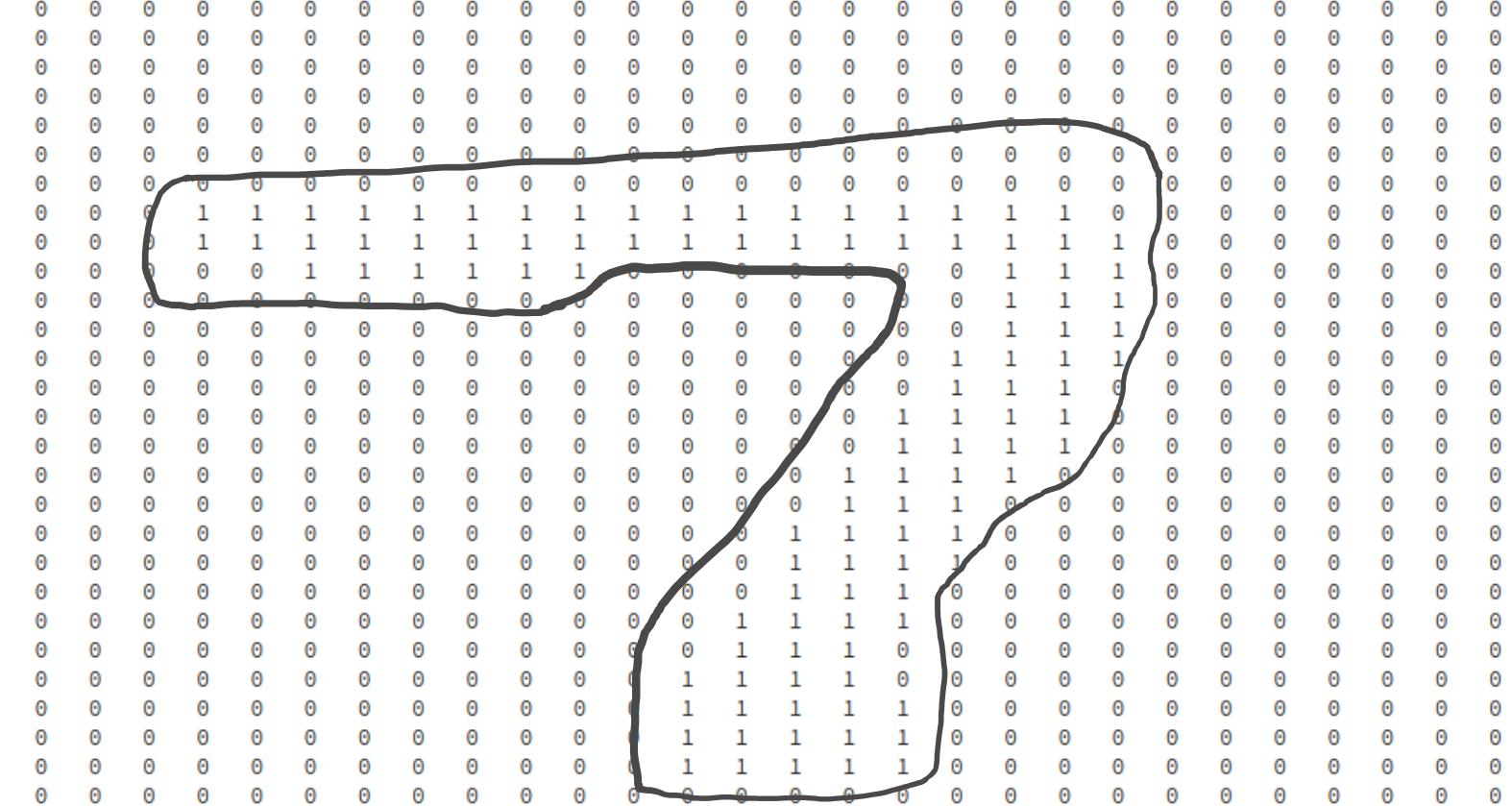
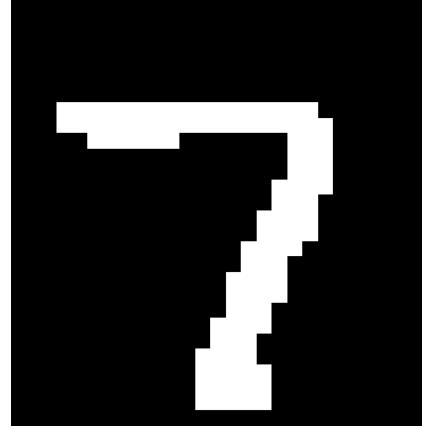
- database of handwritten digits
- 10 classes
 - Training set 60,000 images
 - Test set of 10,000 images
- Greyscale images of size 28x28
- Often treated as ‘Hello World’ for any ML/DL practitioner

Image taken from Researchgate

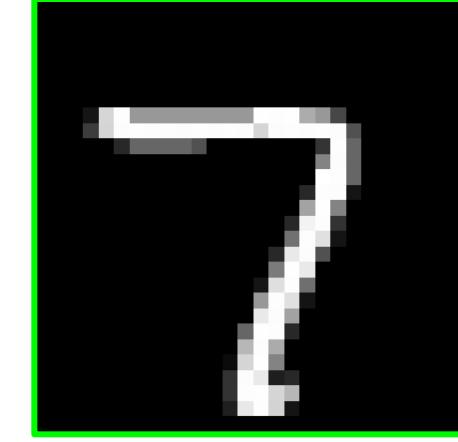
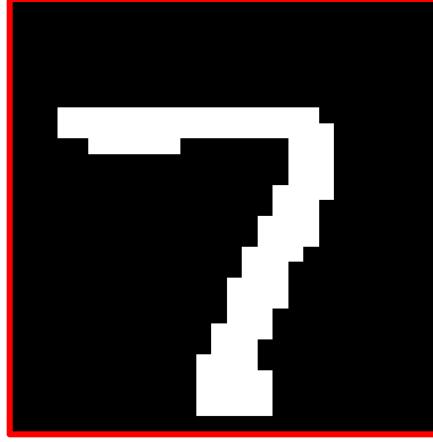
Example from MNIST



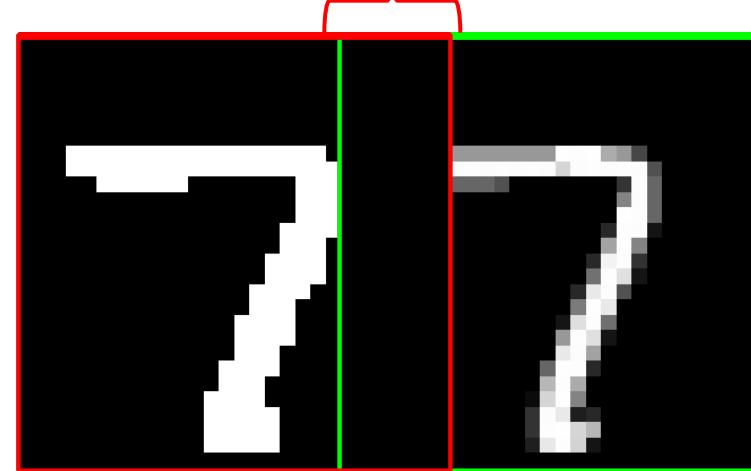
7 Like Window?



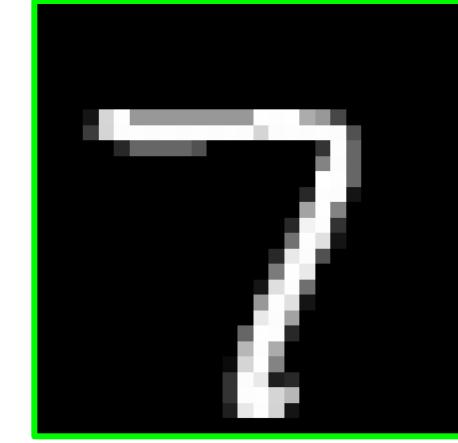
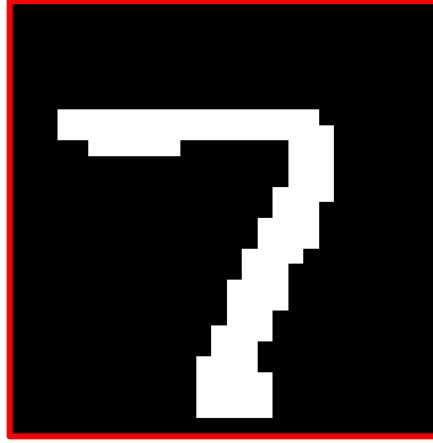
Filtering



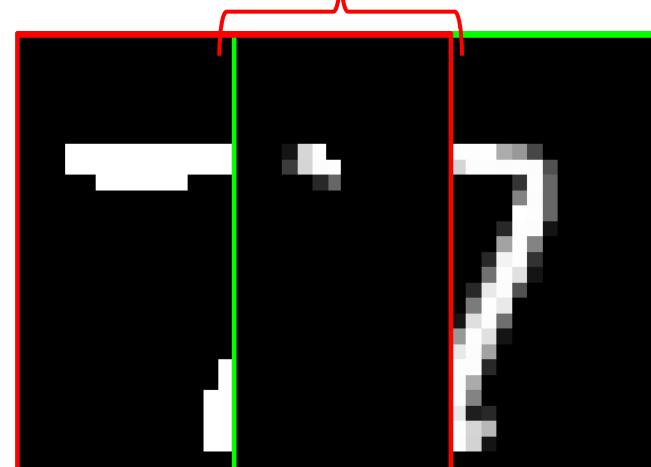
0



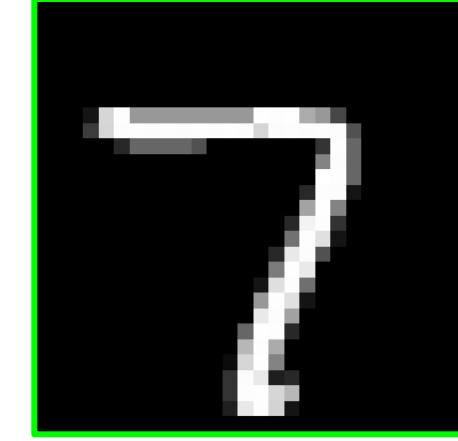
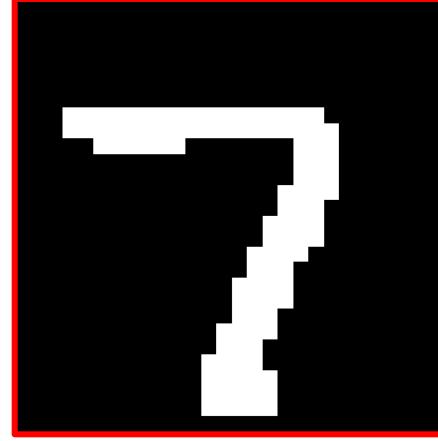
Filtering



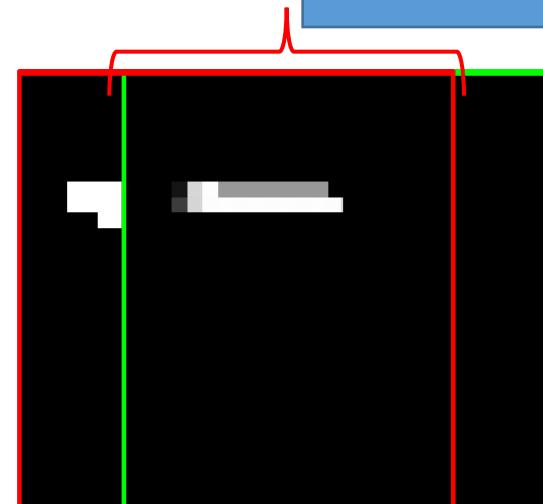
1410



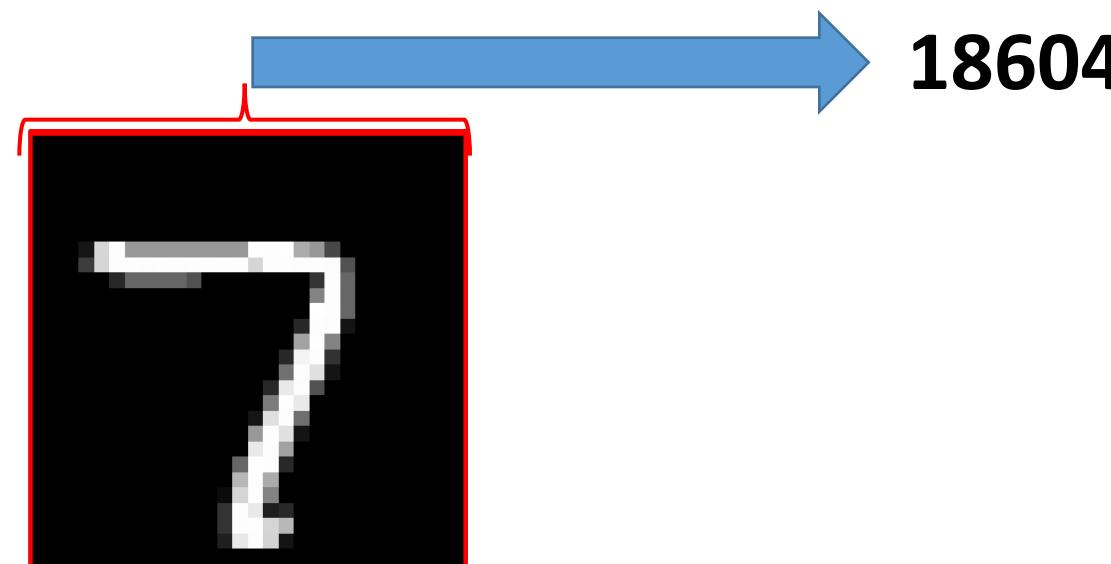
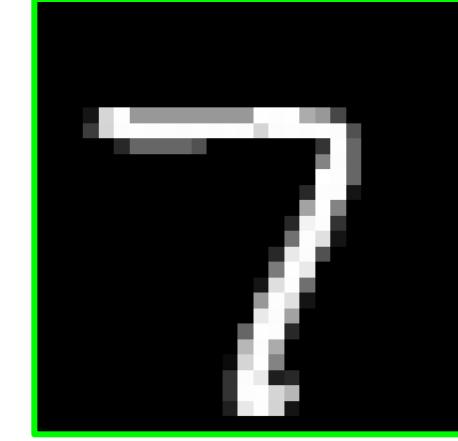
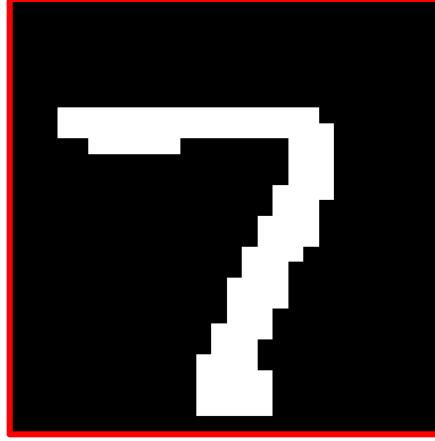
Filtering



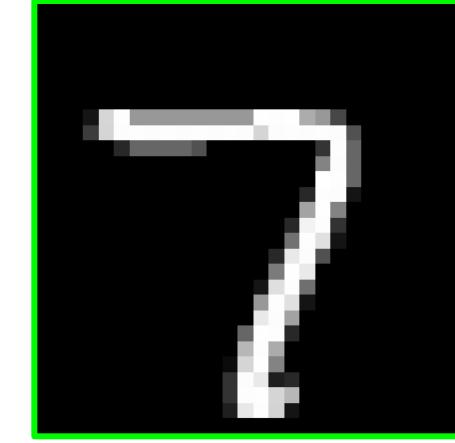
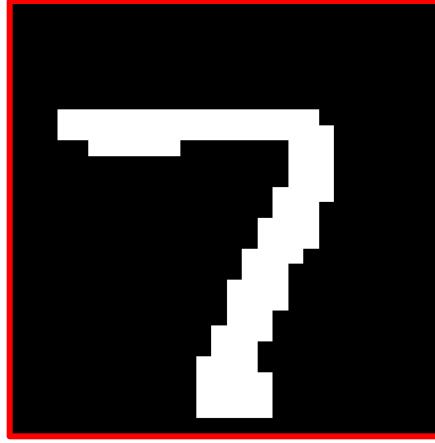
4098



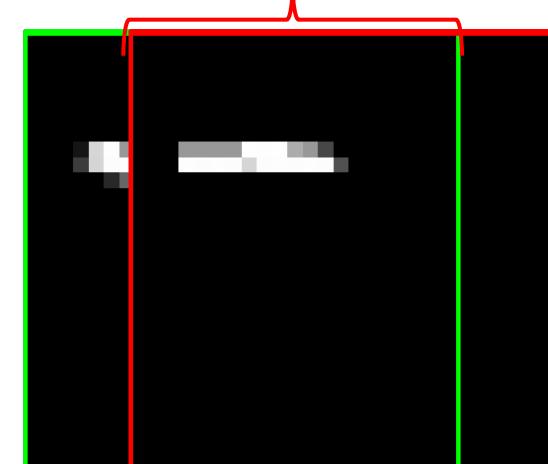
Filtering

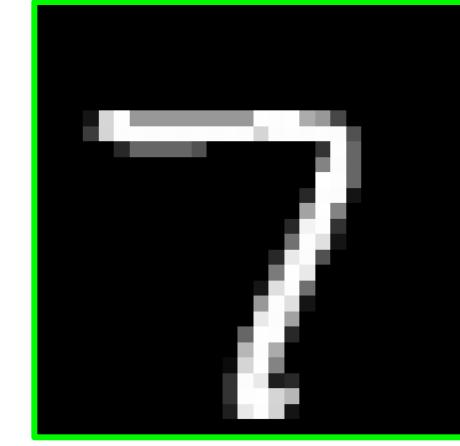
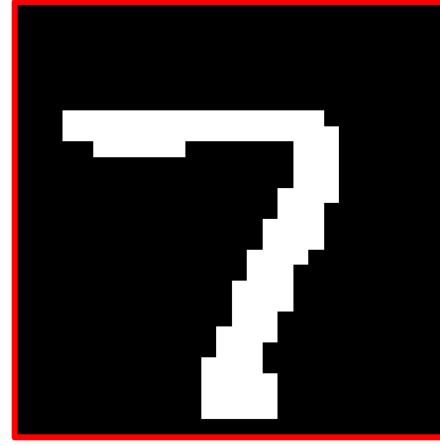


Filtering



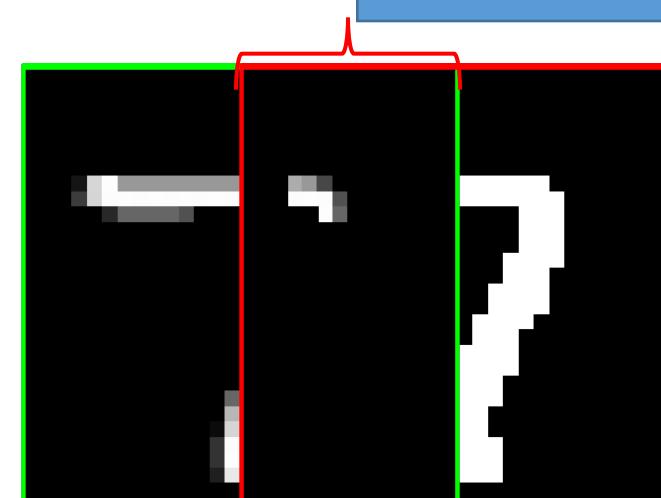
4331



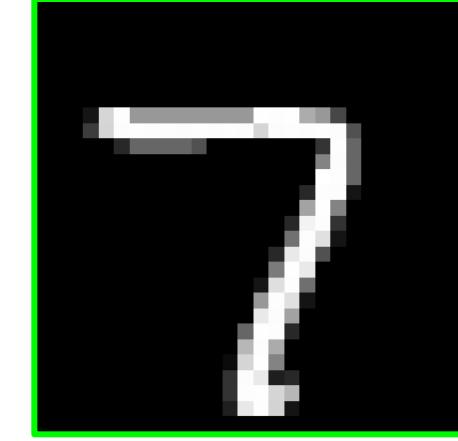
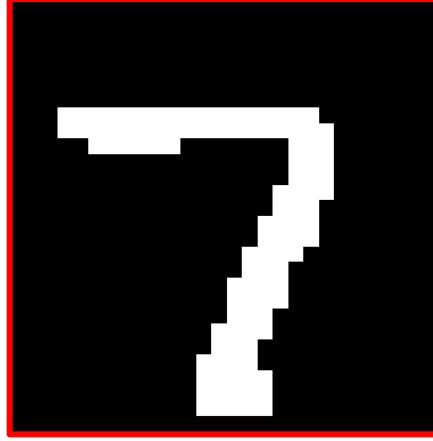


Filtering

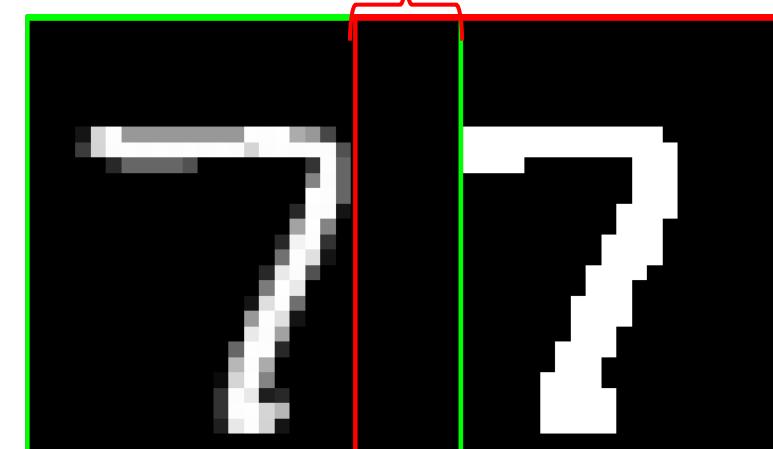
1589



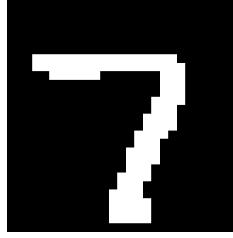
Filtering



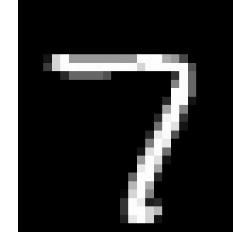
0



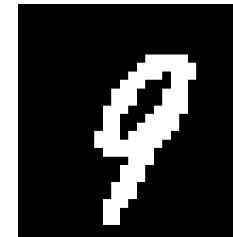
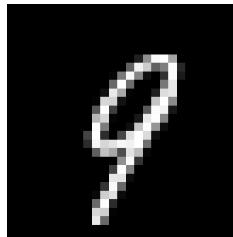
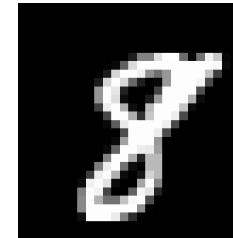
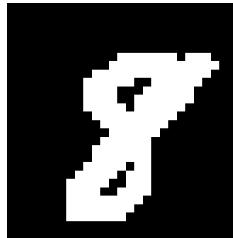
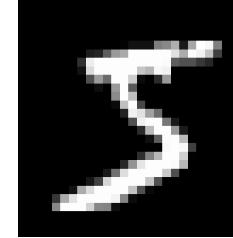
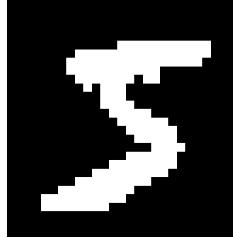
Filters



Test Images

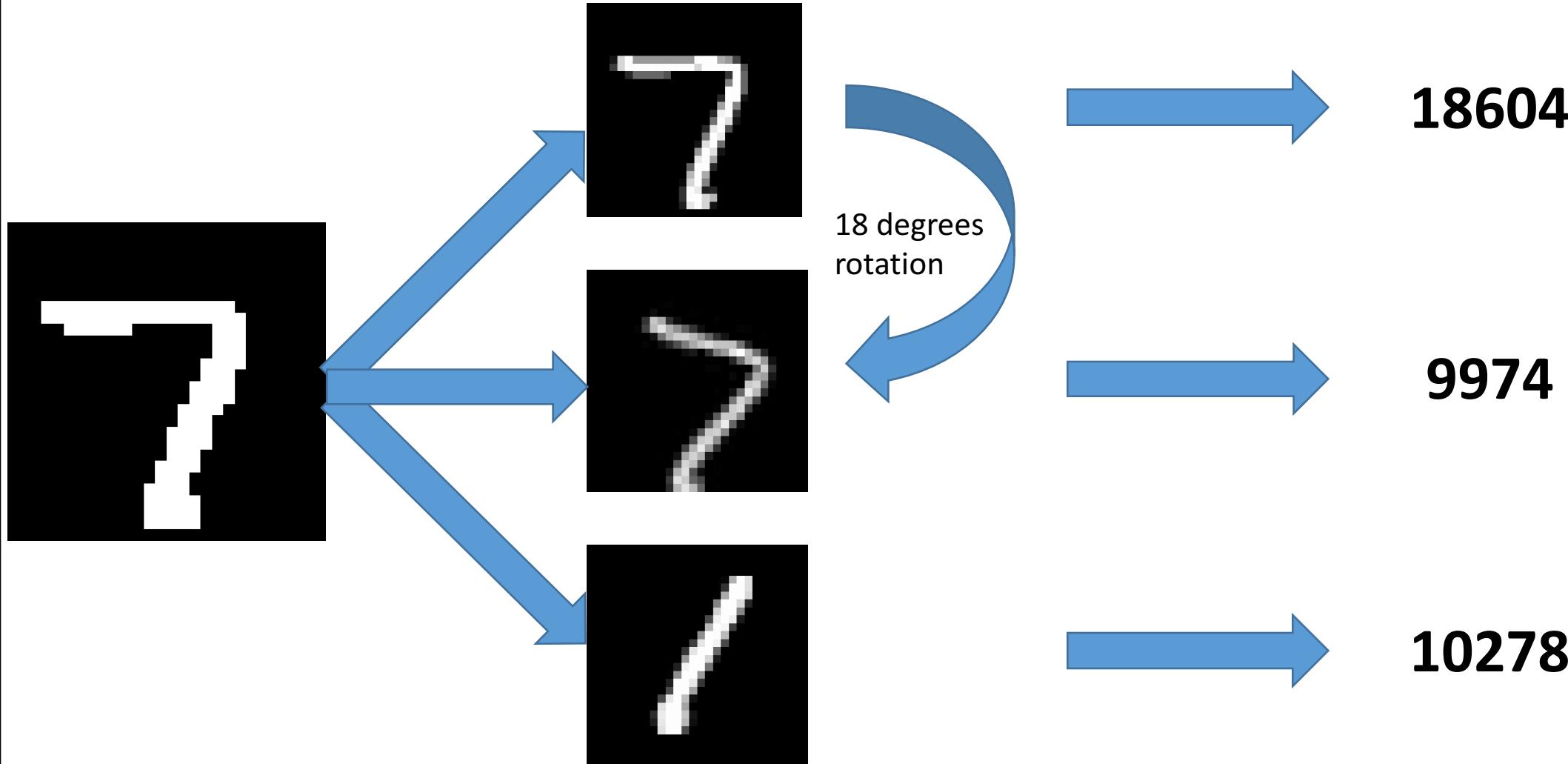


Classification by Matching Filters

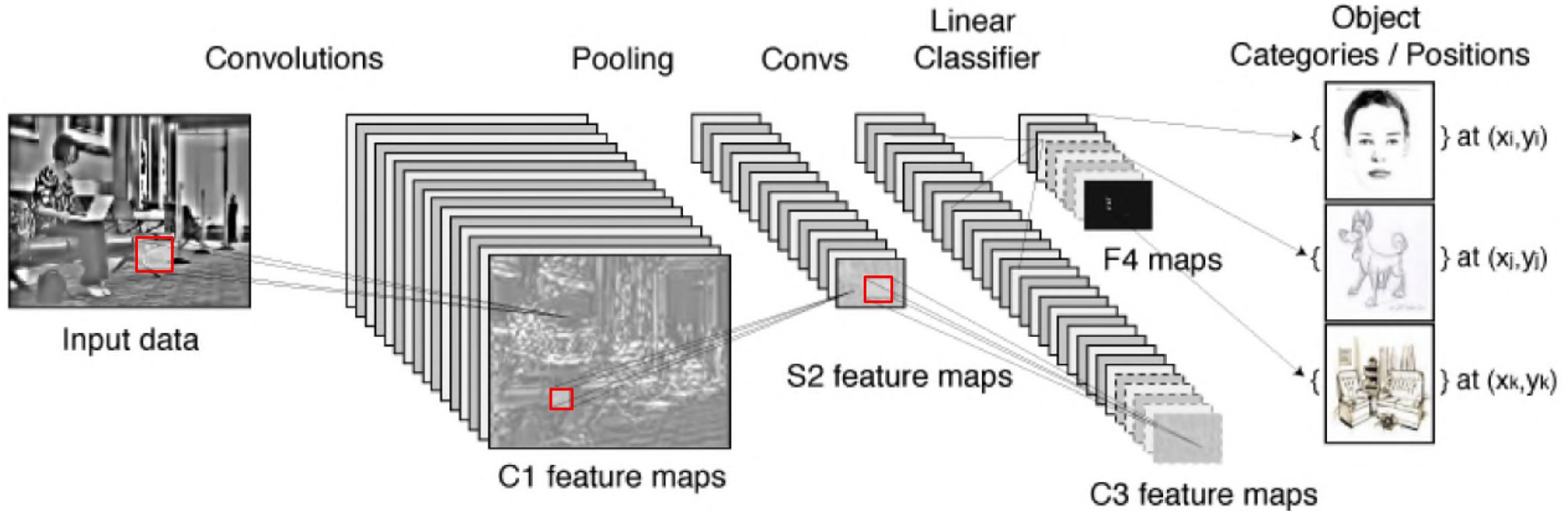


But what if the test image is a little
Rotated (or)
Skewed (or)
Zoomed out and so on

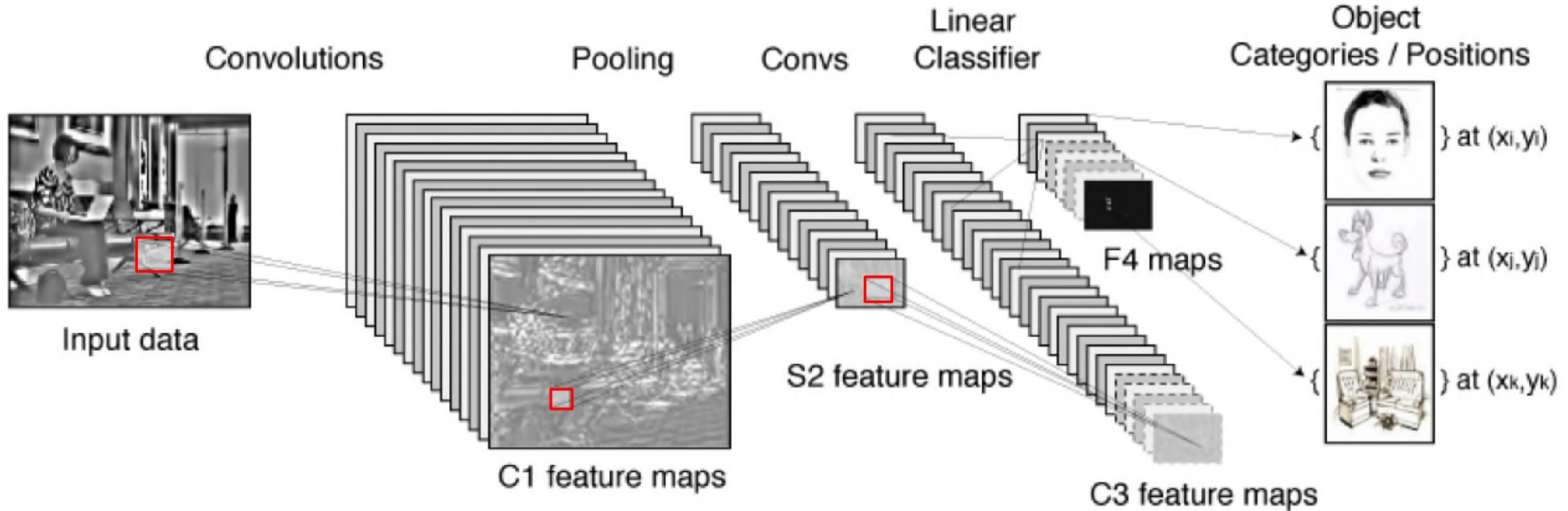
Effect of Slight Rotation



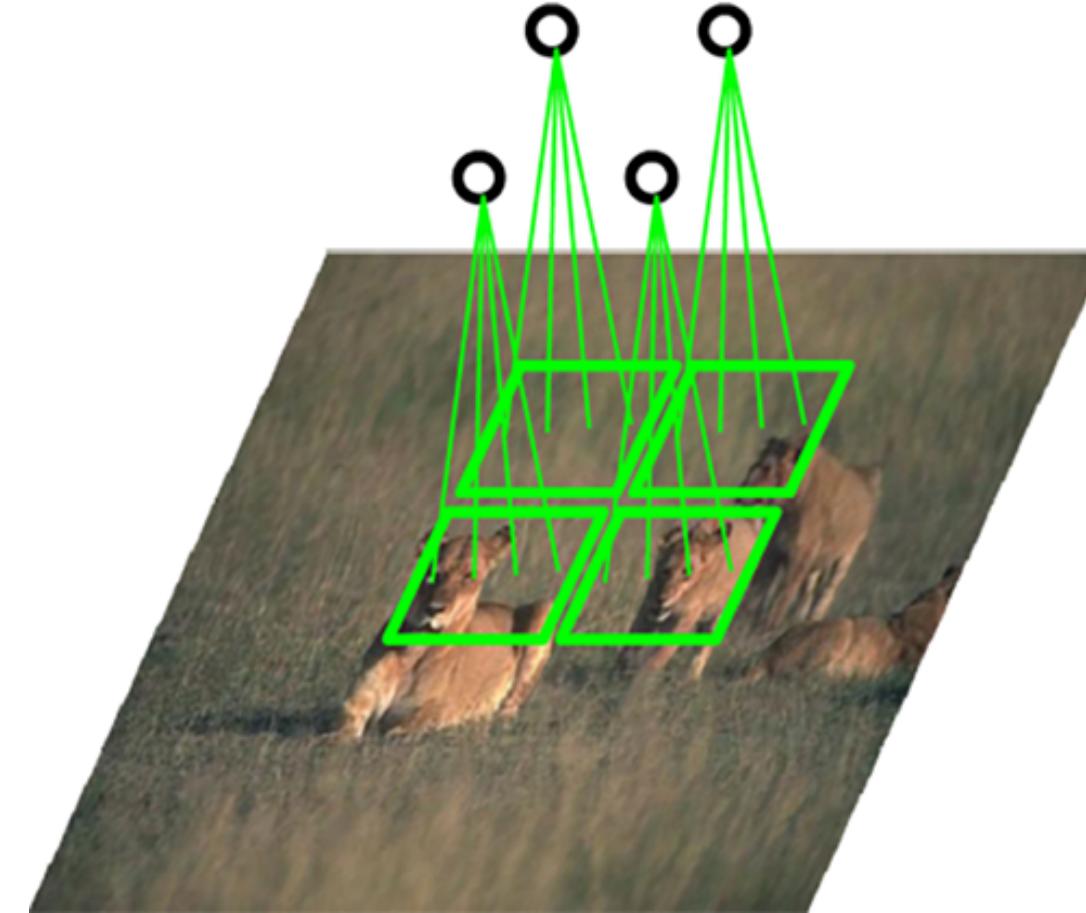
Convnets (Fukushima, LeCun, Hinton)



Convnets (Fukushima, LeCun, Hinton)

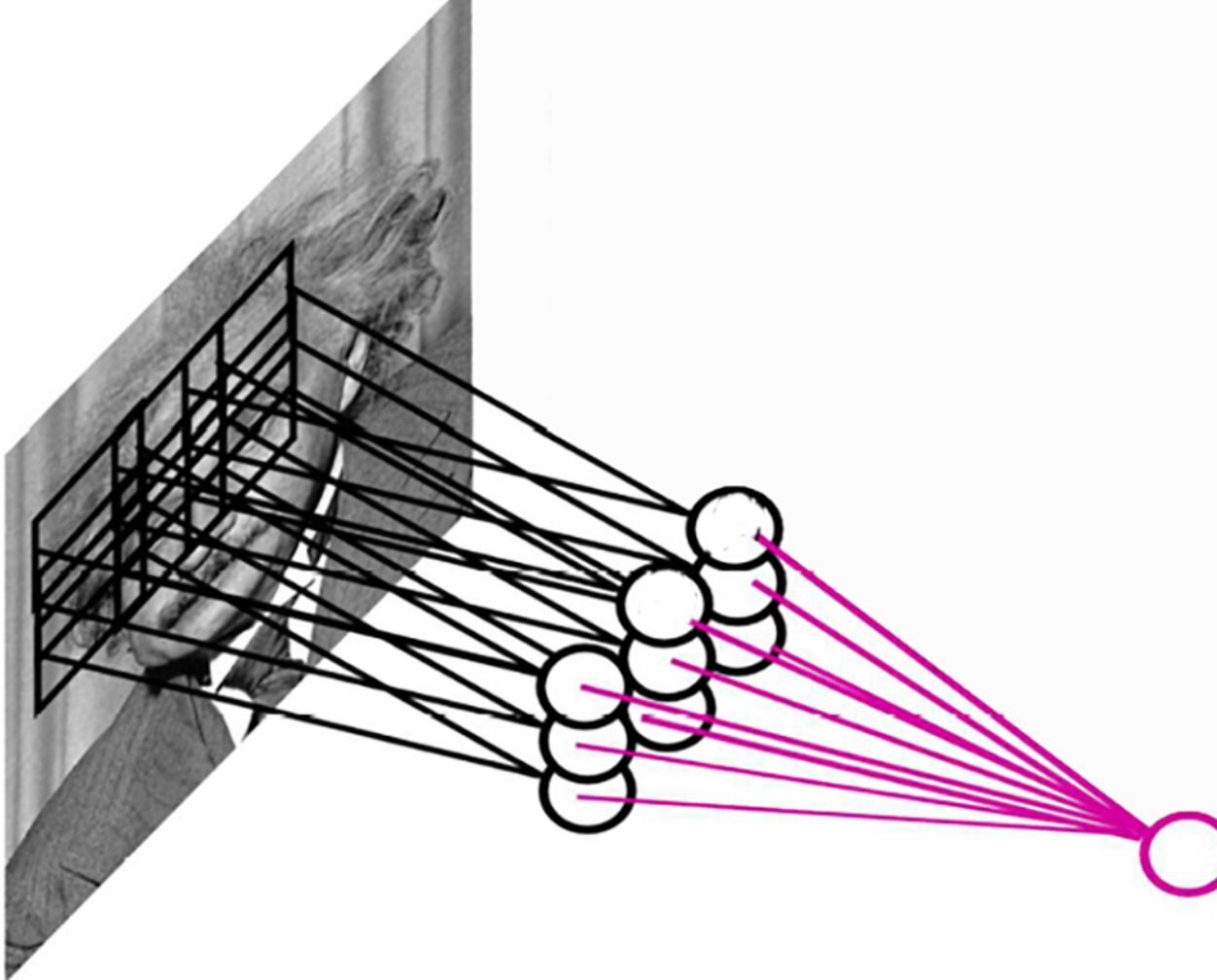


Convnets (Fukushima, LeCun, Hinton)



Slide courtesy: Nando de Freitas, University of Oxford

Pooling



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

Convolution

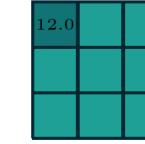
Filter

0	1	2
2	2	0
0	1	2

Image

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

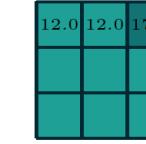
3 ₀	3 ₁	2 ₂	1	0
0 ₂	0 ₂	1 ₀	3	1
3 ₀	1 ₁	2 ₂	2	3
2	0	0	2	2
2	0	0	0	1



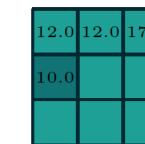
3	3 ₀	2 ₁	1 ₂	0
0	0 ₂	1 ₂	3 ₀	1
3	1 ₀	2 ₁	2 ₂	3
2	0	0	2	2
2	0	0	0	1



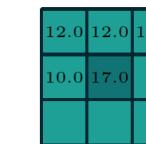
3	3	2 ₀	1 ₁	0 ₂
0	0	1 ₂	3 ₂	1 ₀
3	1	2 ₀	2 ₁	3 ₂
2	0	0	2	2
2	0	0	0	1



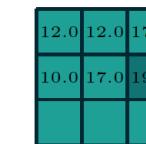
3	3	2	1	0
0 ₀	0 ₁	1 ₂	3	1
3 ₂	1 ₂	2 ₀	2	3
2 ₀	0 ₁	0 ₂	2	2
2	0	0	0	1



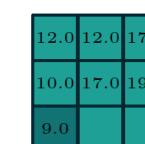
3	3	2	1	0
0	0 ₀	1 ₁	3 ₂	1
3	1 ₂	2 ₂	2 ₀	3
2	0 ₀	0 ₁	2 ₂	2
2	0	0	0	1



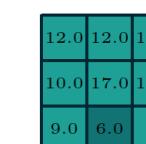
3	3	2	1	0
0	0	1 ₀	3 ₁	1 ₂
3	1	2 ₂	2 ₂	3 ₀
2	0	0 ₀	2 ₁	2 ₂
2	0	0	0	1



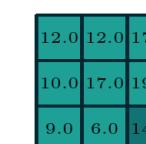
3	3	2	1	0
0	0	1 ₃	3	1
3 ₀	1 ₁	2 ₂	2	3
2 ₂	0 ₂	0 ₀	2	2
2 ₀	0 ₁	0 ₂	0	1



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



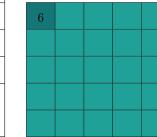
3	3	2	1	0
0	0	1	3	1
3	1 ₂	2 ₂	2 ₀	3 ₀
2	0 ₂	0 ₂	2 ₀	2
2	0 ₀	0 ₁	0 ₂	1



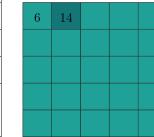
Convolution with Zero Padding

Filter
0 1 2
2 2 0
0 1 2

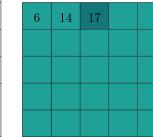
0 ₀	0 ₁	0 ₂	0	0	0	0
0 ₂	3 ₂	3 ₀	2	1	0	0
0 ₀	0 ₁	0 ₂	1	3	1	0
0	3	1	2	2	3	0
0	2	0	0	2	2	0
0	2	0	0	0	1	0
0	0	0	0	0	0	0



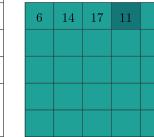
0	0 ₀	0 ₁	0 ₂	0	0	0
0	3 ₂	3 ₂	2 ₀	1	0	0
0	0 ₀	0 ₁	1 ₂	3	1	0
0	3	1	2	2	3	0
0	2	0	0	2	2	0
0	2	0	0	0	1	0
0	0	0	0	0	0	0



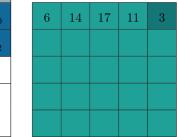
0	0	0 ₀	0 ₁	0 ₂	0	0
0	3	3 ₂	2 ₂	1 ₂	0 ₀	0
0	0	0 ₀	1 ₁	3 ₂	1	0
0	3	1	2	2	3	0
0	2	0	0	2	2	0
0	2	0	0	0	1	0
0	0	0	0	0	0	0



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



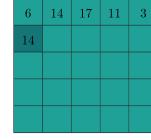
0	0	0	0 ₀	0 ₁	0 ₂	0
0	3	3	2 ₁	1 ₂	0 ₀	0
0	0	0	1 ₀	3 ₀	1 ₁	0 ₂
0	3	1	2	2	3	0
0	2	0	0	2	2	0
0	2	0	0	0	1	0
0	0	0	0	0	0	0



Image

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

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



...

...

...

...

Convolution with Strides

Convolving a 3×3 kernel over a 5×5 input using 2×2 strides

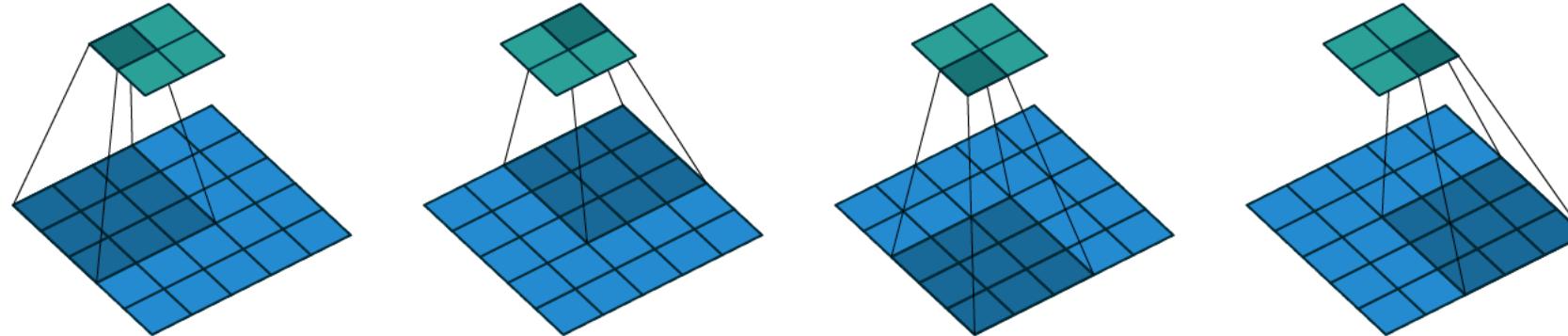


Image courtesy: Vincent Dumoulin

Convolution with Strides and Zero Padding

Convolving a 3×3 kernel over a 5×5 input using 1×1 strides and half padding

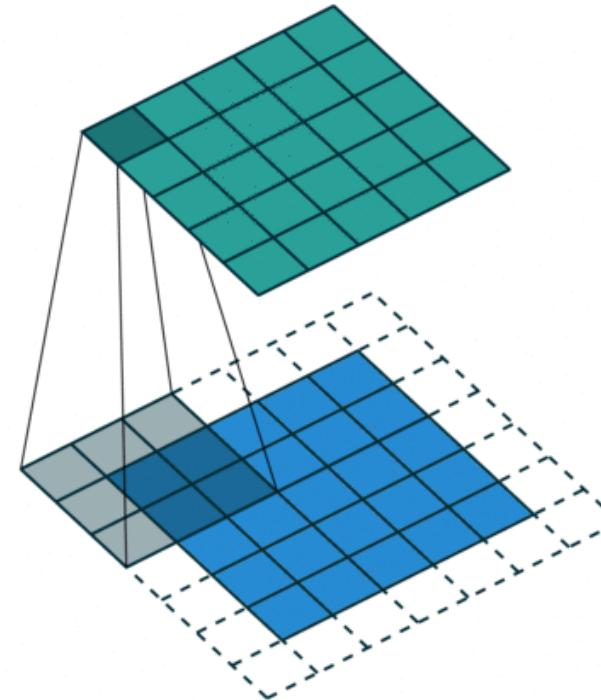


Image courtesy: Vincent Dumoulin

Inputs Generally have Multiple Channels

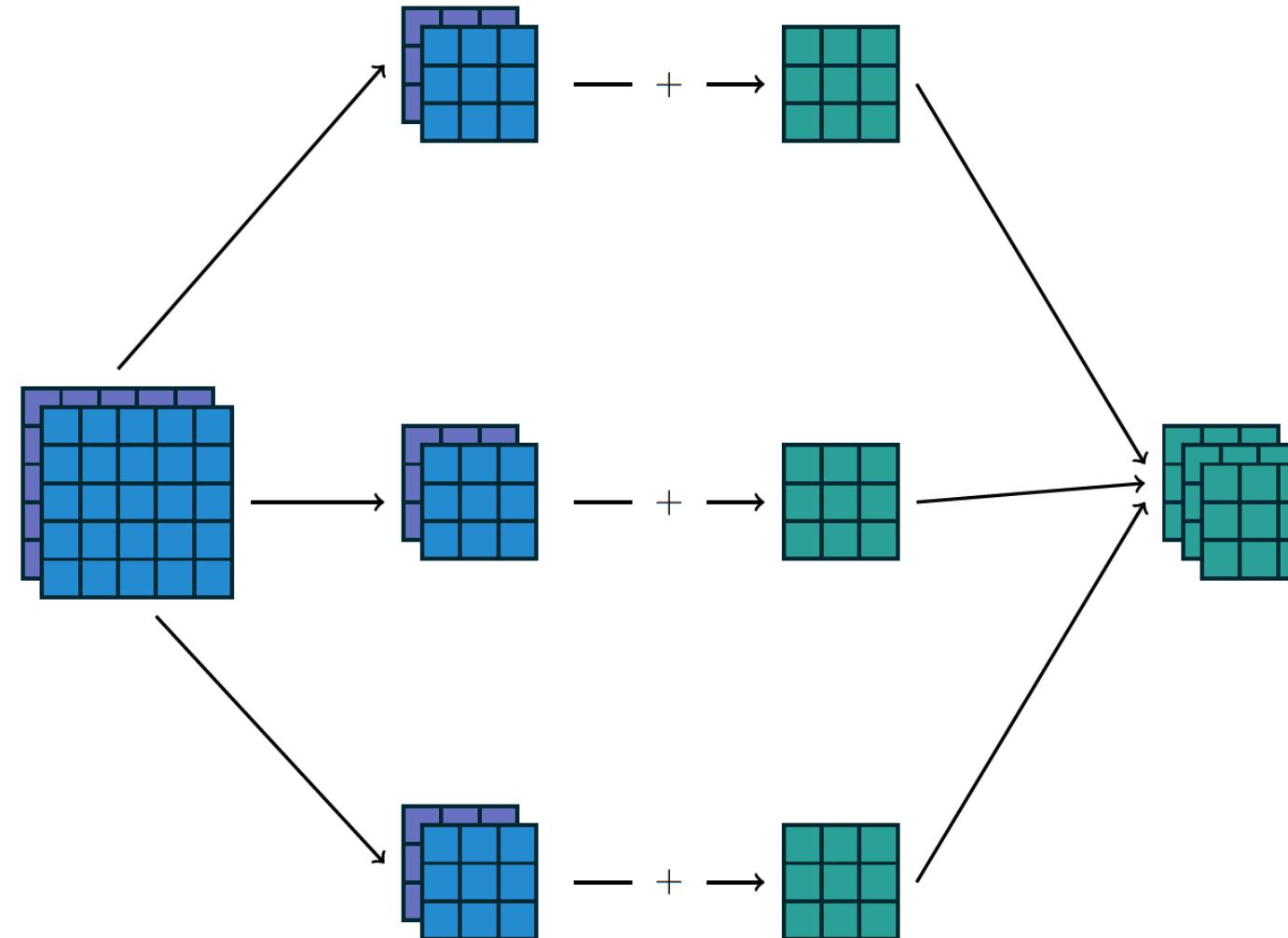


Image courtesy: Vincent Dumoulin

Convolution Arithmetic

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

Image width = image height = w

Filter width = Filter height = k

Stride = s

$$\text{Output size} = \left\lceil \frac{w-k}{s} \right\rceil + 1$$

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

$$\text{So, output size} = \left\lceil \frac{w+2p-k}{s} \right\rceil + 1$$

Note the box function

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

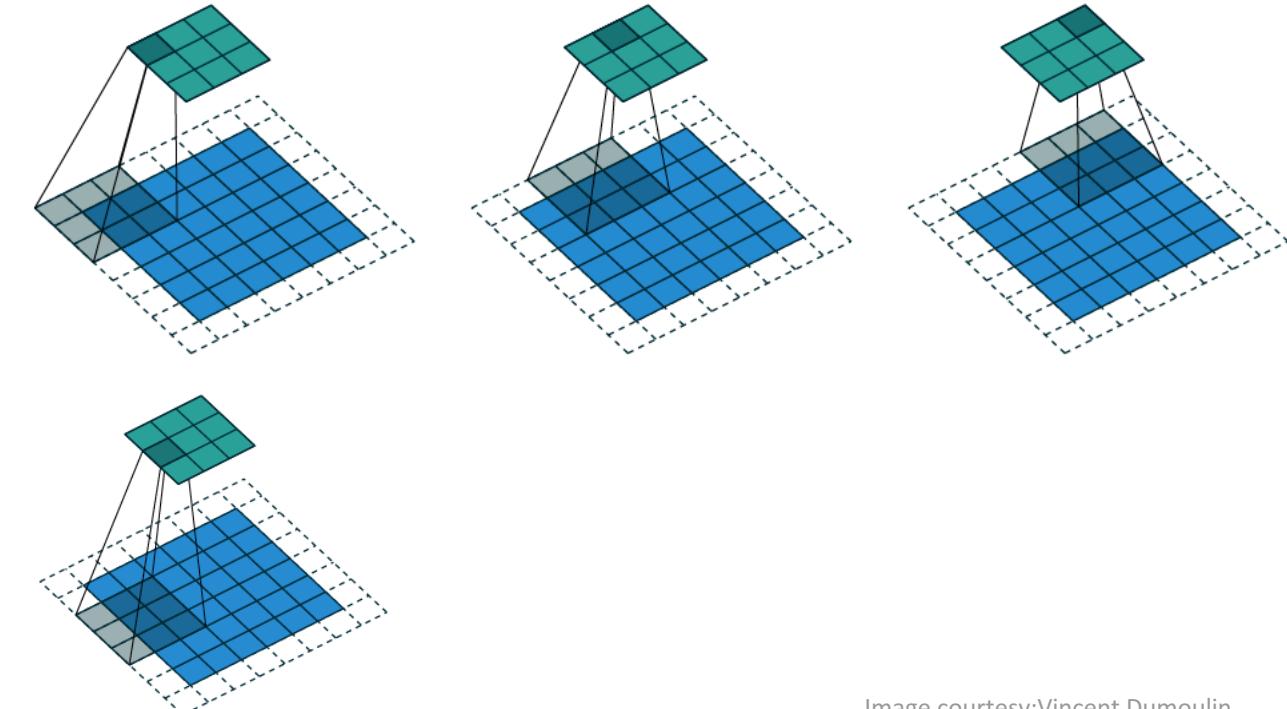
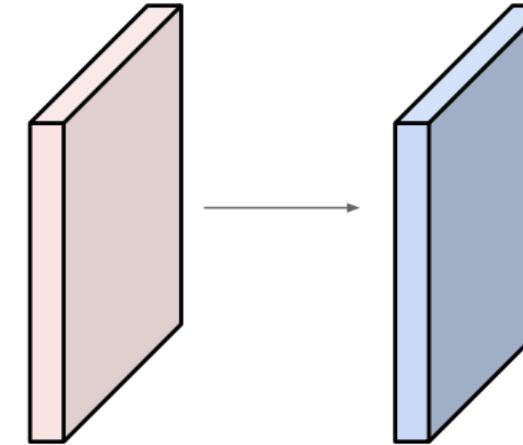


Image courtesy: Vincent Dumoulin

Convolution Arithmetic

Input volume: $32 \times 32 \times 3$ [w, h, c].

64 filters of size 3×3 [k, k] with
stride 2 [s], pad 1 [p]



What is the output feature map size?

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

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

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

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

Image courtesy: Andrej Karpathy

Pooling

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

3.0		

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

3.0	3.0	

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

3.0	3.0	3.0

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

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

3.0	3.0	3.0

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

3.0	3.0	3.0

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

3.0	3.0	3.0

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

3.0	3.0	3.0

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

3.0	3.0	3.0

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

3.0	3.0	3.0

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

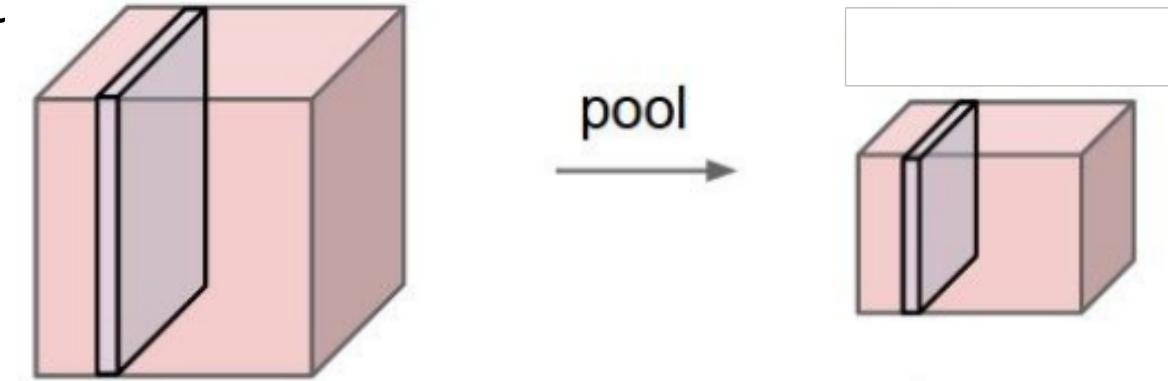
$$\text{Output size} = \left\lceil \frac{w-k}{s} \right\rceil + 1$$

Input volume: $32 \times 32 \times 3$ [w, h, c].

Max-pooling kernel of size 2×2 [k, k] with stride 2 [s]

What is the output feature map size?

$$\left\lceil \frac{32-2}{2} \right\rceil + 1 = 16 \quad \text{So, } 16 \times 16 \times 3 \text{ [w, h, c]}$$

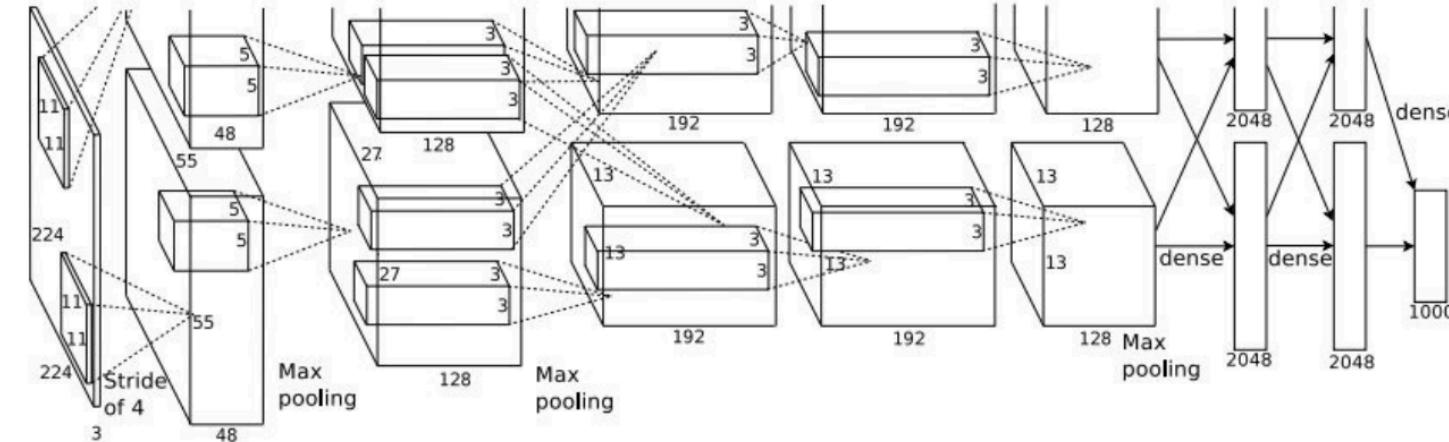


And What is the number of parameters in this pooling layer?

0

Visualizations

AlexNet (2012)



First layer (CONV1): 96 11x11 filters

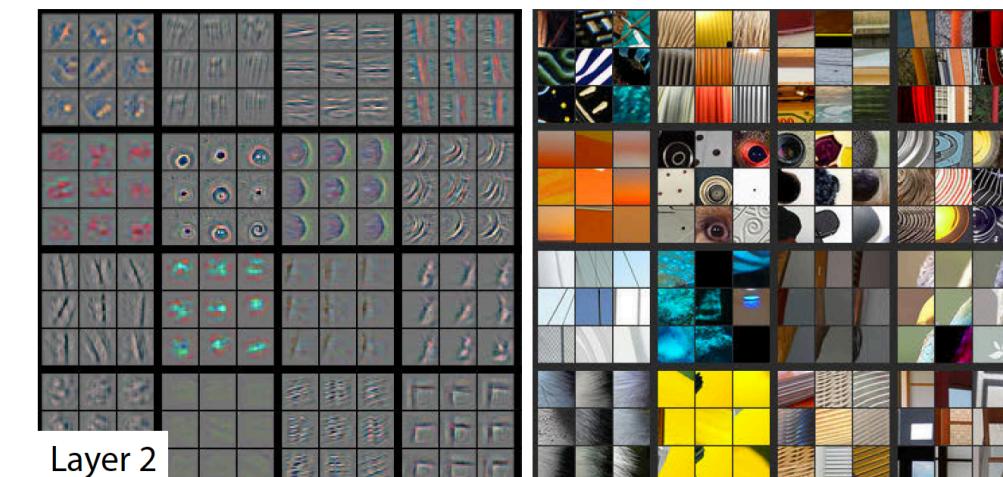
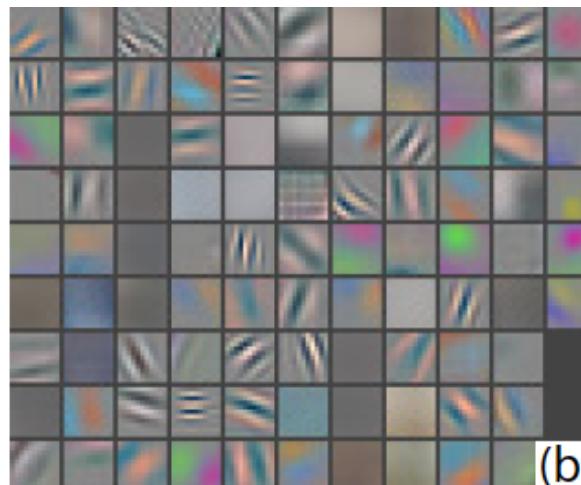


Image courtesy: Zeiler, Fergus, 2013

Visualizations

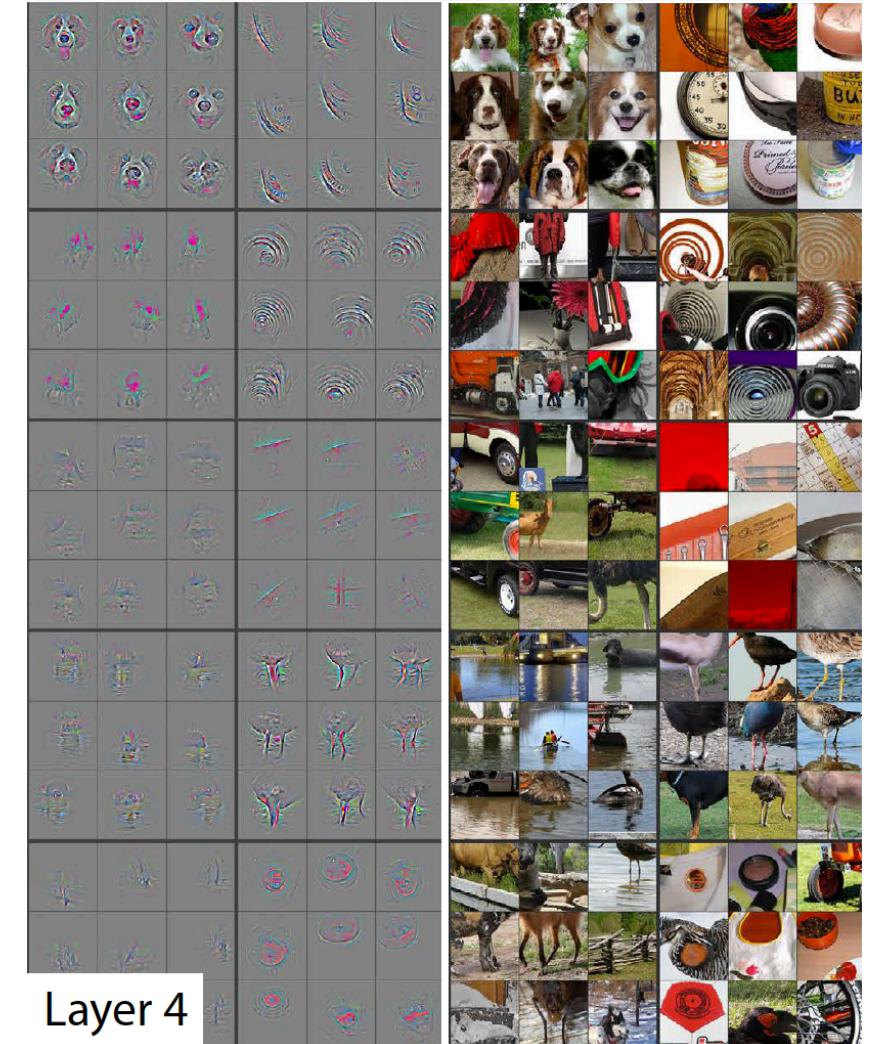
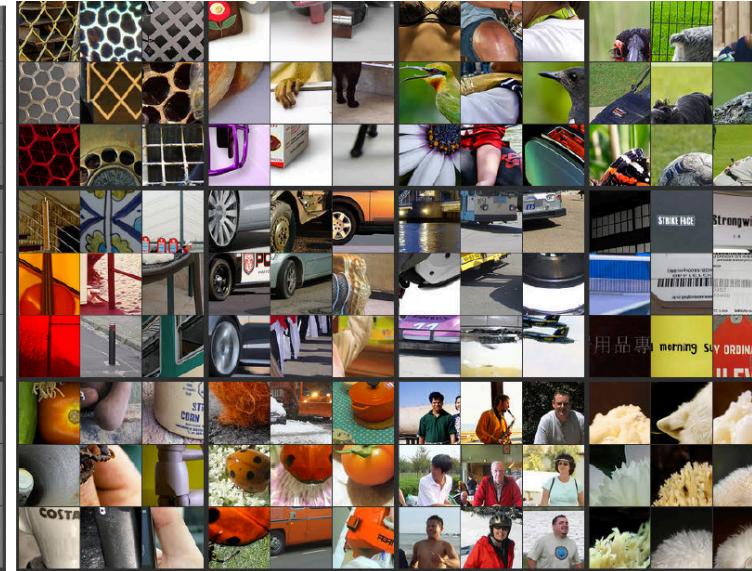
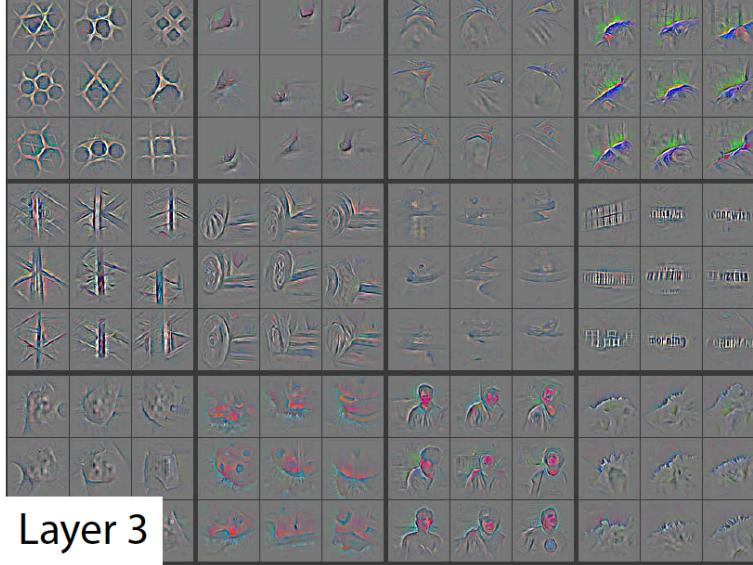


Image courtesy: Zeiler, Fergus, 2013
30

<https://distill.pub/2018/building-blocks/>

The Building Blocks of Interpretability

Interpretability techniques are normally studied in isolation.

We explore the powerful interfaces that arise when you combine them—and the rich structure of this combinatorial space.

CHOOSE AN INPUT IMAGE



For instance, by combining feature visualization (*what is a neuron looking for?*) with attribution (*how does it affect the output?*), we can explore how the network decides between labels like **Labrador retriever** and **tiger cat**.



Several floppy ear detectors seem to be important when distinguishing dogs, whereas pointy ears are used to classify "tiger cat".

CHANNELS THAT MOST SUPPORT ...

LABRADOR RETRIEVER ← → TIGER CAT