

Convolutional Neural Networks

Eren Erdal Aksoy



HALMSTAD
UNIVERSITY

Introduction: What is an Image?

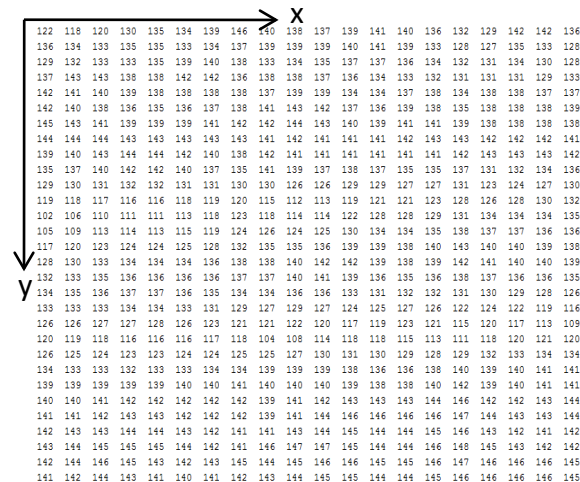
- An image is a 2D function: $f(x, y)$, where x and y are spatial coordinates, and the **amplitude** of the function is called **intensity** or **grey level**.
- An images is just a matrix of numbers, e.g. [0, 255].
- Images are composed of picture elements: **pixels**

Tasks in Computer Vision

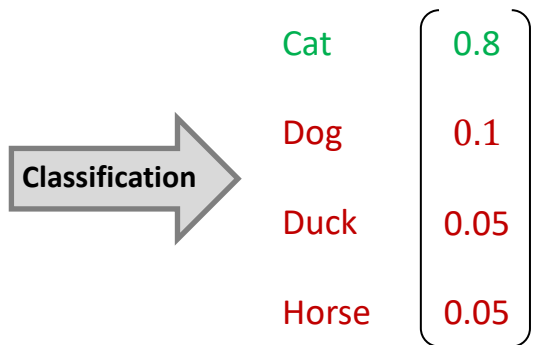
- **Classification:** output variable takes class label. Can produce probability of belonging to a particular class.
- **Regression:** output variable takes continuous value



What we see!

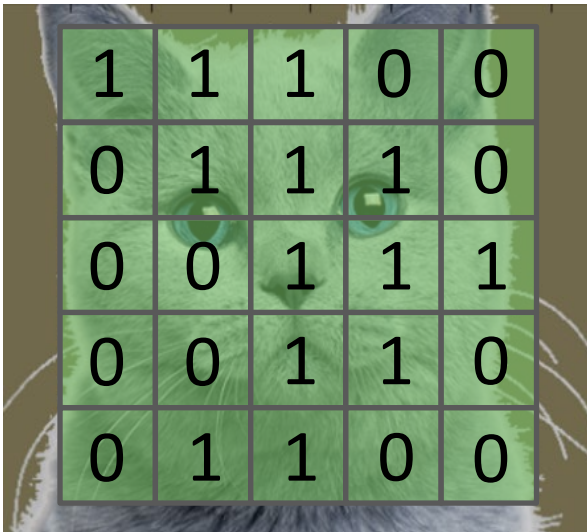


What computers see!



Introduction: What is an Image Processing?

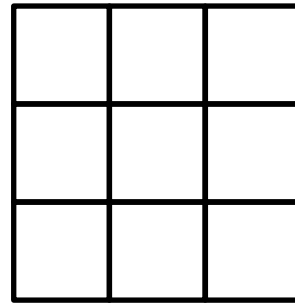
- **Image processing** is a method to perform some operations on an image, in order to get an enhanced image or to extract some useful information from it.
- For instance, we slide the **3x3 filter** over the image, element-wise multiply, and add the outputs



Image

1	0	1
0	1	0
1	0	1

Filter (kernel)



Result

$$I(x, y) * h = \sum_{i=-a}^a \sum_{j=-b}^b I(x - i, y - j) \cdot h(i, j)$$



Inner Product

Image Convolution

a process of combining image **pixels** with a certain matrix weight to identify specific features of the image, such as edge detection, sharpening, blurring,

Introduction: Image Convolution

- **Image processing** is a method to perform some operations on an image, in order to get an enhanced image or to extract some useful information from it.
- For instance, we slide the **3x3 filter** over the image, element-wise multiply, and add the outputs

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

Image

1	0	1
0	1	0
1	0	1

Filter (kernel)

$$I(x, y) * h = \sum_{i=-a}^a \sum_{j=-b}^b I(x - i, y - j) \cdot h(i, j)$$

Result

Introduction: Image Convolution

- **Image processing** is a method to perform some operations on an image, in order to get an enhanced image or to extract some useful information from it.
- For instance, we slide the **3x3 filter** over the image, element-wise multiply, and add the outputs

1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

Image

1	0	1
0	1	0
1	0	1

Filter (kernel)

$$I(x, y) * h = \sum_{i=-a}^a \sum_{j=-b}^b I(x - i, y - j) \cdot h(i, j)$$

4		

Result

$$\begin{aligned} & (1 \times 1) + (1 \times 0) + (1 \times 1) + \\ & (0 \times 0) + (1 \times 1) + (1 \times 0) + \\ & (0 \times 1) + (0 \times 0) + (1 \times 1) = 4 \end{aligned}$$

Introduction: Image Convolution

- **Image processing** is a method to perform some operations on an image, in order to get an enhanced image or to extract some useful information from it.
- For instance, we slide the **3x3 filter** over the image, element-wise multiply, and add the outputs

1	1 _{x1}	1 _{x0}	0 _{x1}	0
0	1 _{x0}	1 _{x1}	1 _{x0}	0
0	0 _{x1}	1 _{x0}	1 _{x1}	1
0	0	1	1	0
0	1	1	0	0

Image

1	0	1
0	1	0
1	0	1

Filter (kernel)

$$I(x, y) * h = \sum_{i=-a}^a \sum_{j=-b}^b I(x - i, y - j) \cdot h(i, j)$$

4	3	

Result

$$\begin{aligned}
 & (1 \times 1) + (1 \times 0) + (0 \times 1) + \\
 & (1 \times 0) + (1 \times 1) + (1 \times 0) + \\
 & (0 \times 1) + (1 \times 0) + (1 \times 1) = 3
 \end{aligned}$$

Introduction: Image Convolution

- **Image processing** is a method to perform some operations on an image, in order to get an enhanced image or to extract some useful information from it.
- For instance, we slide the **3x3 filter** over the image, element-wise multiply, and add the outputs

1	1	1 _{x1}	0 _{x0}	0 _{x1}
0	1	1 _{x0}	1 _{x1}	0 _{x0}
0	0	1 _{x1}	1 _{x0}	1 _{x1}
0	0	1	1	0
0	1	1	0	0

Image

1	0	1
0	1	0
1	0	1

Filter (kernel)

$$I(x, y) * h = \sum_{i=-a}^a \sum_{j=-b}^b I(x - i, y - j) \cdot h(i, j)$$

4	3	4

Result

$$\begin{aligned}
 & (1 \times 1) + (0 \times 0) + (0 \times 1) + \\
 & (1 \times 0) + (1 \times 1) + (0 \times 0) + \\
 & (1 \times 1) + (1 \times 0) + (1 \times 1) = 4
 \end{aligned}$$

Introduction: Image Convolution

- **Image processing** is a method to perform some operations on an image, in order to get an enhanced image or to extract some useful information from it.
- For instance, we slide the **3x3 filter** over the image, element-wise multiply, and add the outputs

1	1	1	0	0
0 _{x1}	1 _{x0}	1 _{x1}	1	0
0 _{x0}	0 _{x1}	1 _{x0}	1	1
0 _{x1}	0 _{x0}	1 _{x1}	1	0
0	1	1	0	0

Image

1	0	1
0	1	0
1	0	1

Filter (kernel)

$$I(x, y) * h = \sum_{i=-a}^a \sum_{j=-b}^b I(x - i, y - j) \cdot h(i, j)$$

4	3	4
2		

Result

$$\begin{aligned} & (0 \times 1) + (1 \times 0) + (1 \times 1) + \\ & (0 \times 0) + (0 \times 1) + (1 \times 0) + \\ & (0 \times 1) + (0 \times 0) + (1 \times 1) = 2 \end{aligned}$$

Introduction: Image Convolution

- **Image processing** is a method to perform some operations on an image, in order to get an enhanced image or to extract some useful information from it.
- For instance, we slide the **3x3 filter** over the image, element-wise multiply, and add the outputs

1	1	1	0	0
0	1 _{x1}	1 _{x0}	1 _{x1}	0
0	0 _{x0}	1 _{x1}	1 _{x0}	1
0	0 _{x1}	1 _{x0}	1 _{x1}	0
0	1	1	0	0

Image

1	0	1
0	1	0
1	0	1

Filter (kernel)

$$I(x, y) * h = \sum_{i=-a}^a \sum_{j=-b}^b I(x - i, y - j) \cdot h(i, j)$$

4	3	4
2	4	

Result

$$\begin{aligned} & (1 \times 1) + (1 \times 0) + (1 \times 1) + \\ & (0 \times 0) + (1 \times 1) + (1 \times 0) + \\ & (0 \times 1) + (1 \times 0) + (1 \times 1) = 4 \end{aligned}$$

Introduction: Image Convolution

- **Image processing** is a method to perform some operations on an image, in order to get an enhanced image or to extract some useful information from it.
- For instance, we slide the **3x3 filter** over the image, element-wise multiply, and add the outputs

1	1	1	0	0
0	1	1 _{x1}	1 _{x0}	0 _{x1}
0	0	1 _{x0}	1 _{x1}	1 _{x0}
0	0	1 _{x1}	1 _{x0}	0 _{x1}
0	1	1	0	0

Image

1	0	1
0	1	0
1	0	1

Filter (kernel)

$$I(x, y) * h = \sum_{i=-a}^a \sum_{j=-b}^b I(x - i, y - j) \cdot h(i, j)$$

4	3	4
2	4	3

Result

$$\begin{aligned} & (1 \times 1) + (1 \times 0) + (0 \times 1) + \\ & (1 \times 0) + (1 \times 1) + (1 \times 0) + \\ & (1 \times 1) + (1 \times 0) + (0 \times 1) = 3 \end{aligned}$$

Introduction: Image Convolution

- **Image processing** is a method to perform some operations on an image, in order to get an enhanced image or to extract some useful information from it.
- For instance, we slide the **3x3 filter** over the image, element-wise multiply, and add the outputs

1	1	1	0	0
0	1	1	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0 _{x0}	0 _{x1}	1 _{x0}	1	0
0 _{x1}	1 _{x0}	1 _{x1}	0	0

Image

1	0	1
0	1	0
1	0	1

Filter (kernel)

$$I(x, y) * h = \sum_{i=-a}^a \sum_{j=-b}^b I(x - i, y - j) \cdot h(i, j)$$

4	3	4
2	4	3
2		

Result

$$\begin{aligned} & (0 \times 1) + (0 \times 0) + (1 \times 1) + \\ & (0 \times 0) + (0 \times 1) + (1 \times 0) + \\ & (0 \times 1) + (1 \times 0) + (1 \times 1) = 2 \end{aligned}$$

Introduction: Image Convolution

- **Image processing** is a method to perform some operations on an image, in order to get an enhanced image or to extract some useful information from it.
- For instance, we slide the **3x3 filter** over the image, element-wise multiply, and add the outputs

1	1	1	0	0
0	1	1	1	0
0	0 _{x1}	1 _{x0}	1 _{x1}	1
0	0 _{x0}	1 _{x1}	1 _{x0}	0
0	1 _{x1}	1 _{x0}	0 _{x1}	0

Image

1	0	1
0	1	0
1	0	1

Filter (kernel)

$$I(x, y) * h = \sum_{i=-a}^a \sum_{j=-b}^b I(x - i, y - j) \cdot h(i, j)$$

4	3	4
2	4	3
2	3	

Result

$$\begin{aligned} & (0 \times 1) + (1 \times 0) + (1 \times 1) + \\ & (0 \times 0) + (1 \times 1) + (1 \times 0) + \\ & (1 \times 1) + (1 \times 0) + (0 \times 1) = 3 \end{aligned}$$

Introduction: Image Convolution

- **Image processing** is a method to perform some operations on an image, in order to get an enhanced image or to extract some useful information from it.
- For instance, we slide the **3x3 filter** over the image, element-wise multiply, and add the outputs

1	1	1	0	0
0	1	1	1	0
0	0	1 _{x1}	1 _{x0}	1 _{x1}
0	0	1 _{x0}	1 _{x1}	0 _{x0}
0	1	1 _{x1}	0 _{x0}	0 _{x1}

Image

1	0	1
0	1	0
1	0	1

Filter (kernel)

$$I(x, y) * h = \sum_{i=-a}^a \sum_{j=-b}^b I(x - i, y - j) \cdot h(i, j)$$

4	3	4
2	4	3
2	3	4

Result

$$\begin{aligned} & (1 \times 1) + (1 \times 0) + (1 \times 1) + \\ & (1 \times 0) + (1 \times 1) + (0 \times 0) + \\ & (1 \times 1) + (0 \times 0) + (0 \times 1) = 4 \end{aligned}$$

Introduction: Feature Extraction

- Each **filter** identifies a **specific set of features** of the image



Original



Sharpen



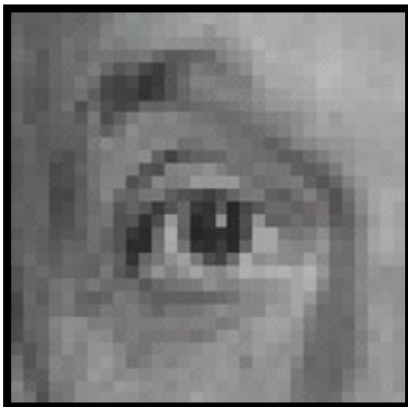
Edge Detect



Strong Edge Detect

How to define mask coefficients?

Depends on what the filter is supposed to do!



Original

0	0	0
0	0	0
0	0	0

Filter



Filtered
(Deleting)

Introduction: Feature Extraction

- Each **filter** identifies a **specific set of features** of the image



Original



Sharpen



Edge Detect



Strong Edge Detect

How to define mask coefficients?

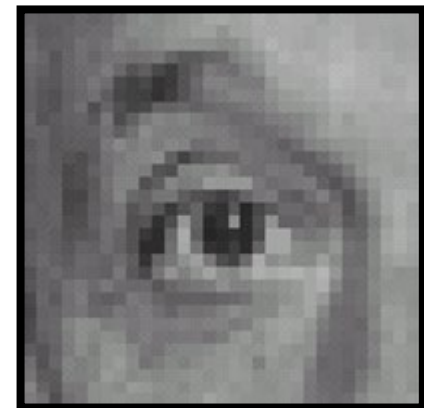
Depends on what the filter is supposed to do!



Original

0	0	0
0	1	0
0	0	0

Filter



Filtered
(No change)

Introduction: Feature Extraction

- Each **filter** identifies a **specific set of features** of the image



Original



Sharpen



Edge Detect



Strong Edge Detect

How to define mask coefficients?

Depends on what the filter is supposed to do!



Original

0	0	0
0	0	1
0	0	0

Filter



Shifted *left*
by 1 pixel

Manual Feature Extraction

- Features should be **robust to image variations**

Domain Knowledge

Define Features

Detect Features to Classify

Viewpoint variation



Scale variation



Deformation



Occlusion



Illumination conditions



Background clutter



Intra-class variation



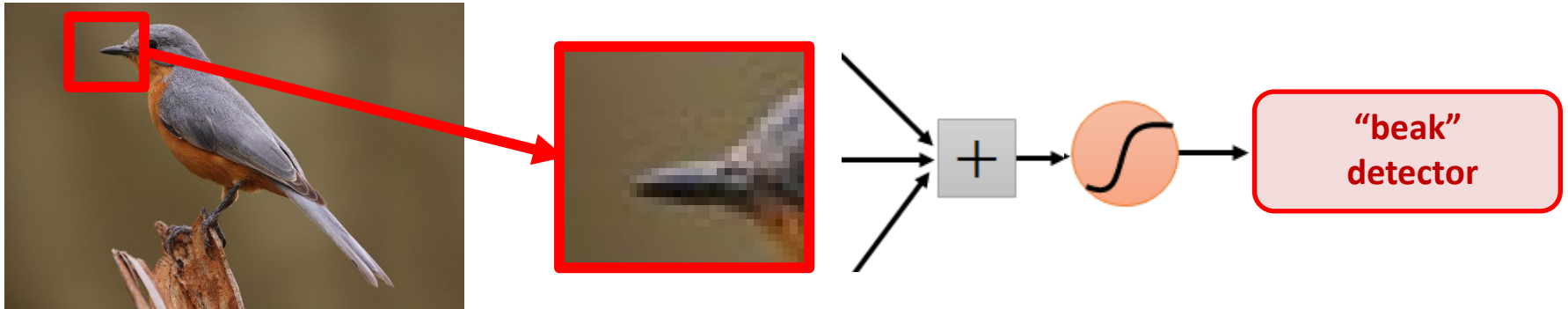
Can we learn a *hierarchy of features* **robust to image variations** directly from the data instead of hand engineering?

Solution: Convolutional Neural Networks

Why CNN for Image?

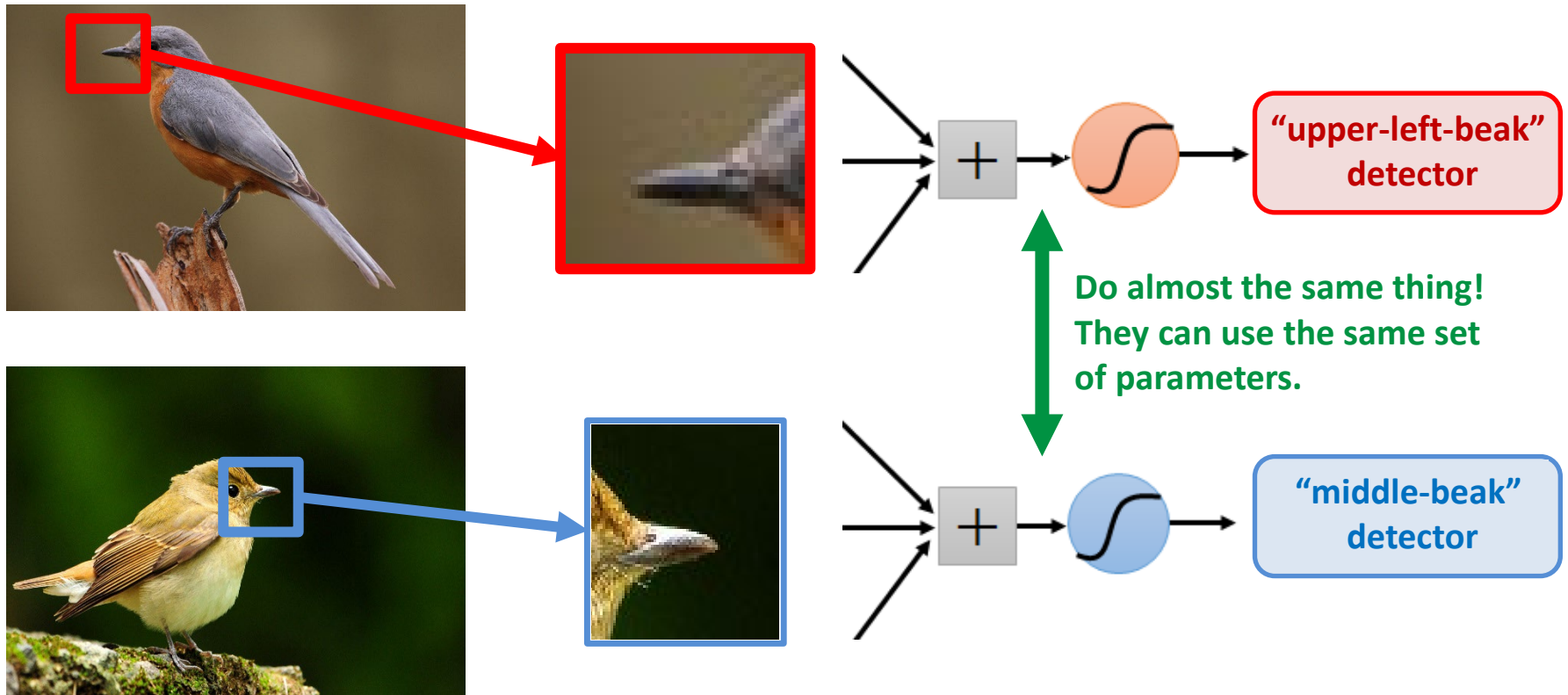
1. Some patterns are much smaller than the whole image

- A neuron **does not have to see the whole image** to discover the pattern
- Connecting to a small region with **fewer parameters**



Why CNN for Image?

1. Some patterns are much smaller than the whole image
 - A neuron **does not have to see the whole image** to discover the pattern
 - Connecting to a small region with **fewer parameters**
2. The very same patterns appear in different regions



Why CNN for Image?

1. **Some patterns are much smaller than the whole image**
 - A neuron **does not have to see the whole image** to discover the pattern
 - Connecting to a small region with **fewer parameters**
2. **The very same patterns appear in different regions**
3. **Subsampling the pixels will not change the object**
 - We can **subsample** the pixels to make the image **smaller**
 - **Fewer parameters** for the network to process the image

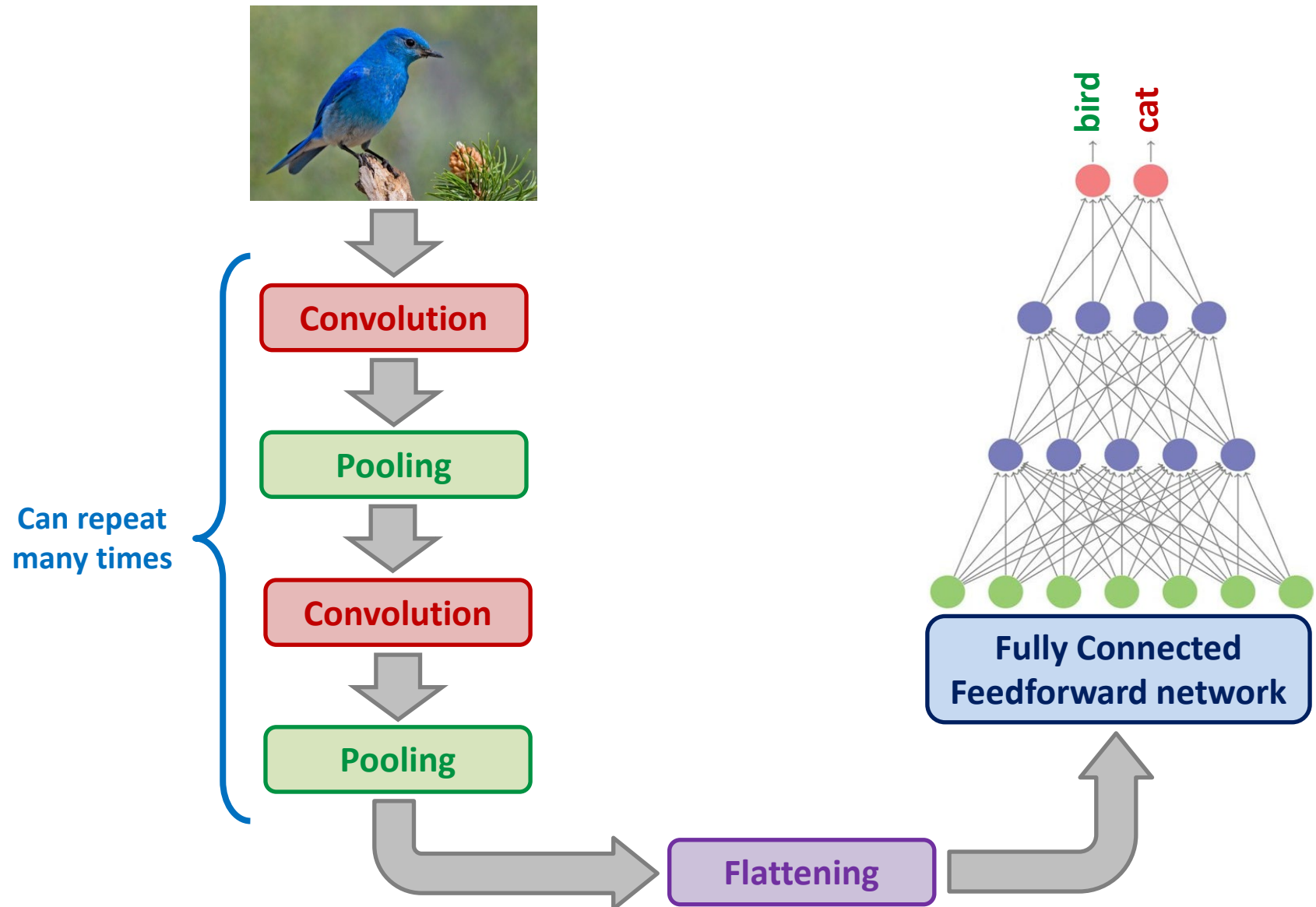


bird

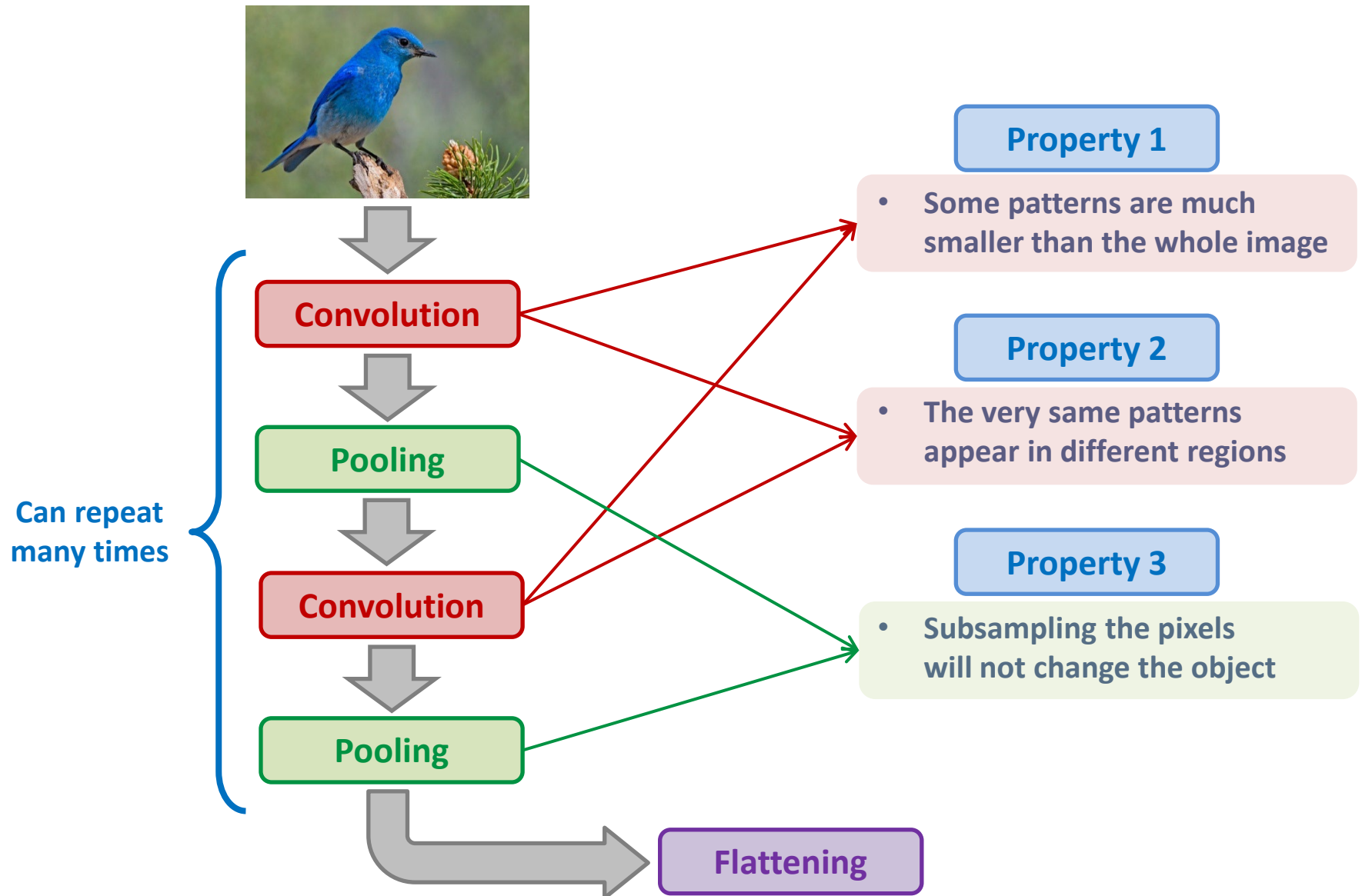


bird

The Whole CNN



The Whole CNN



CNN: Convolution

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

6 x 6 Image

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

⋮

Each filter detects a small pattern (3 x 3).

Property 1

Note that these **filter indices** are the network parameters to be **learned**.
We do not predefine them like in conventional image processing methods.

CNN: Convolution

Stride is how far the filter moves in every step along one direction.

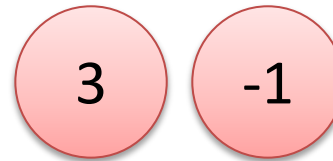
Stride=1

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

6 x 6 Image

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1



CNN: Convolution

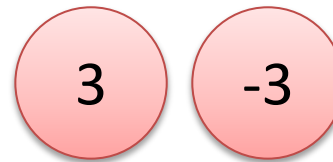
What if Stride is 2?

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

6 x 6 Image

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1



Unless otherwise stated,
stride is by default 1

CNN: Convolution

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

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

6 x 6 Image

A feature map, a.k.a. Activation Map, is the result of applying a convolution across an image

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

Feature Map

CNN: Convolution

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

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

6 x 6 Image

Why do we have two highest scores there?

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

Feature Map

CNN: Convolution

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

This filter detects
left diagonal edges

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

6 x 6 Image

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

Property 2

The very same patterns
appear in different regions

CNN: Convolution

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

6 x 6 Image

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

Follow the same process for **every** filter

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

Feature Maps

CNN: Convolution

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

This filter detects
vertical edges

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

6 x 6 Image

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

For example,
if we had **6 3x3 filters**, we'll get **6 separate feature maps**.
We stack these up to get a “**new image**” of size **4x4x6**

Feature Maps are 4 x 4 images

CNN: Spatial Dimensions

- A closer look at spatial dimensions

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

Input size: **6x6** (spatially)

Filter size: **3x3**

Stride: **1**

Output image size: **4x4**

- **Note that** applying convolution repeatedly will **shrink feature maps spatially!**
- **Shrinking** too **fast** is **not good**, doesn't work well!
- To **preserve size spatially**, it is common to **zero-pad** the border

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

Padding: **1**

Output Image Size

$$Width_{out} = \frac{Width_{in} - Width_{filter} + 2 * PaddingSize}{Stride} + 1$$

$$Height_{out} = \frac{Height_{in} - Height_{filter} + 2 * PaddingSize}{Stride} + 1$$

$$Dimension_{out} = Number_{filter}$$

CNN: Spatial Dimensions

- For example;

Input Image: 32x32x3

Filter size: 5x5

Filter number: 10

Stride: 1

Padding: 2

Output Image Size

$$Width_{out} = \frac{Width_{in} - Width_{filter} + 2 * PaddingSize}{Stride} + 1 = \frac{32 - 5 + 2 * 2}{1} + 1 = 32$$

$$Height_{out} = \frac{Height_{in} - Height_{filter} + 2 * PaddingSize}{Stride} + 1 = \frac{32 - 5 + 2 * 2}{1} + 1 = 32$$

$$Dimension_{out} = Number_{filter} = 10$$

=> Output Image Size: 32x32x10

Number of parameters in this layer?

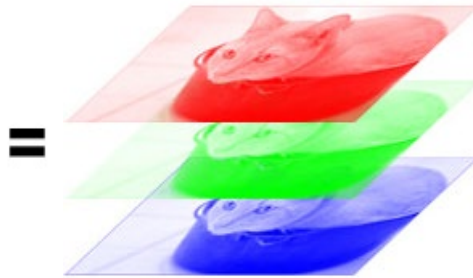
Each filter has 5x5x3 +1 = 76 parameters (+1 for bias)

Since there are 10 filters, in total 76*10 = 760 parameters

CNN: Convolution for Multi-channel Images



Colorful
Image



Red, Green, Blue
Image Channels



Input

4	9	2	5	8	3
5	6	2	4	0	3
2	4	5	4	5	2
5	6	5	4	7	8
5	7	7	9	2	1
5	8	5	3	8	4

$6 \times 6 \times 3$

$*$

Filter 1

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

$3 \times 3 \times 3$

Filter 2

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

$3 \times 3 \times 3$

$=$

Feature Map 1

4×4

$=$

Feature Map 2

4×4

Output

$4 \times 4 \times 2$

CNN: Convolution for Multi-channel Images

0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...

Input Channel #1 (Red)

0	0	0	0	0	0	...
0	167	166	167	169	169	...
0	164	165	168	170	170	...
0	160	162	166	169	170	...
0	156	156	159	163	168	...
0	155	153	153	158	168	...
...

Input Channel #2 (Green)

0	0	0	0	0	0	...
0	163	162	163	165	165	...
0	160	161	164	166	166	...
0	156	158	162	165	166	...
0	155	155	158	162	167	...
0	154	152	152	157	167	...
...

Input Channel #3 (Blue)

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

Kernel Channel #1



308

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

Kernel Channel #2



-498

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



164

+

+

+ 1 = -25



Bias = 1

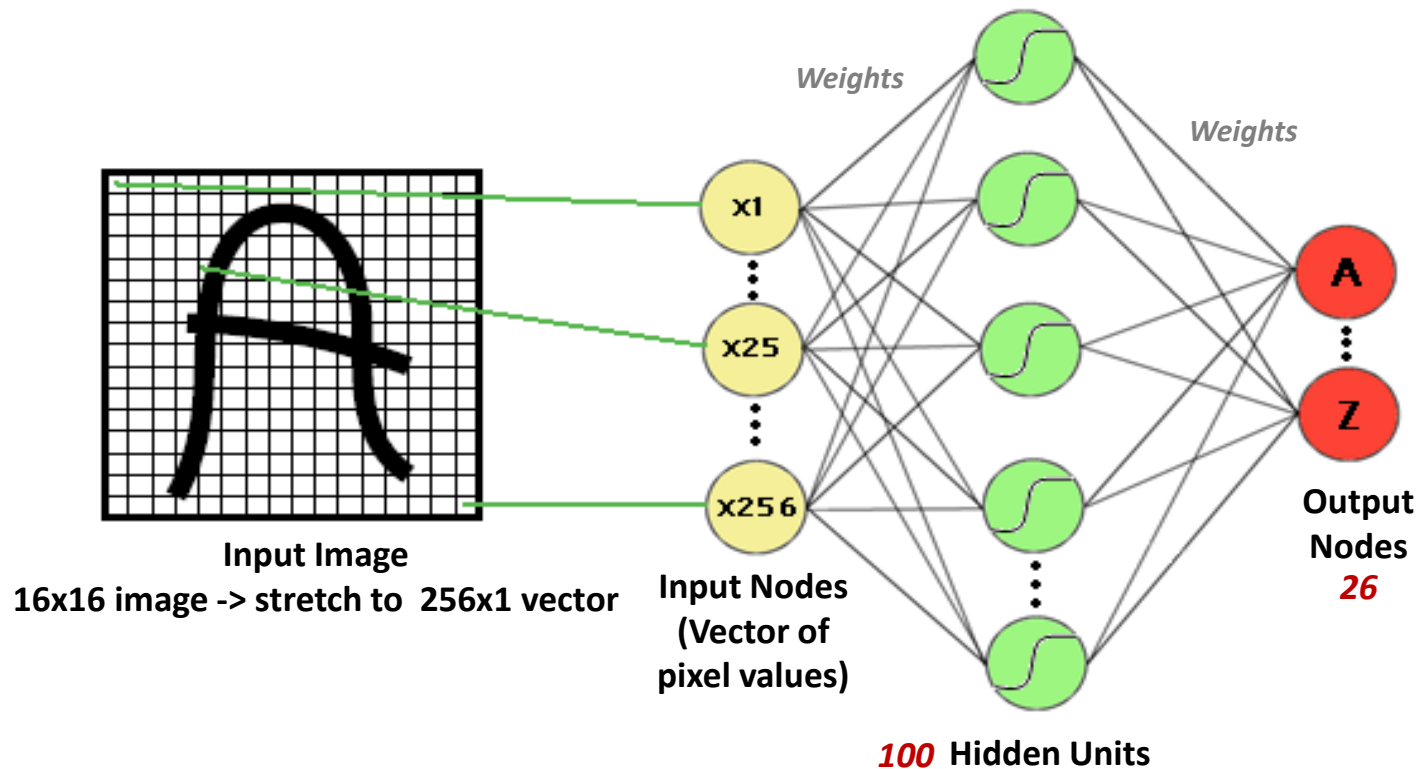
Click [here](#) to see the gif animation!

Output

-25				...
				...
				...
				...
...

Convolution versus Fully Connected

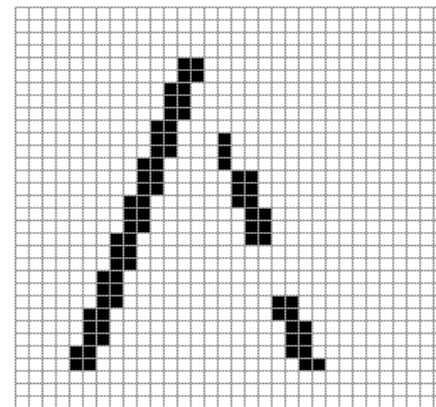
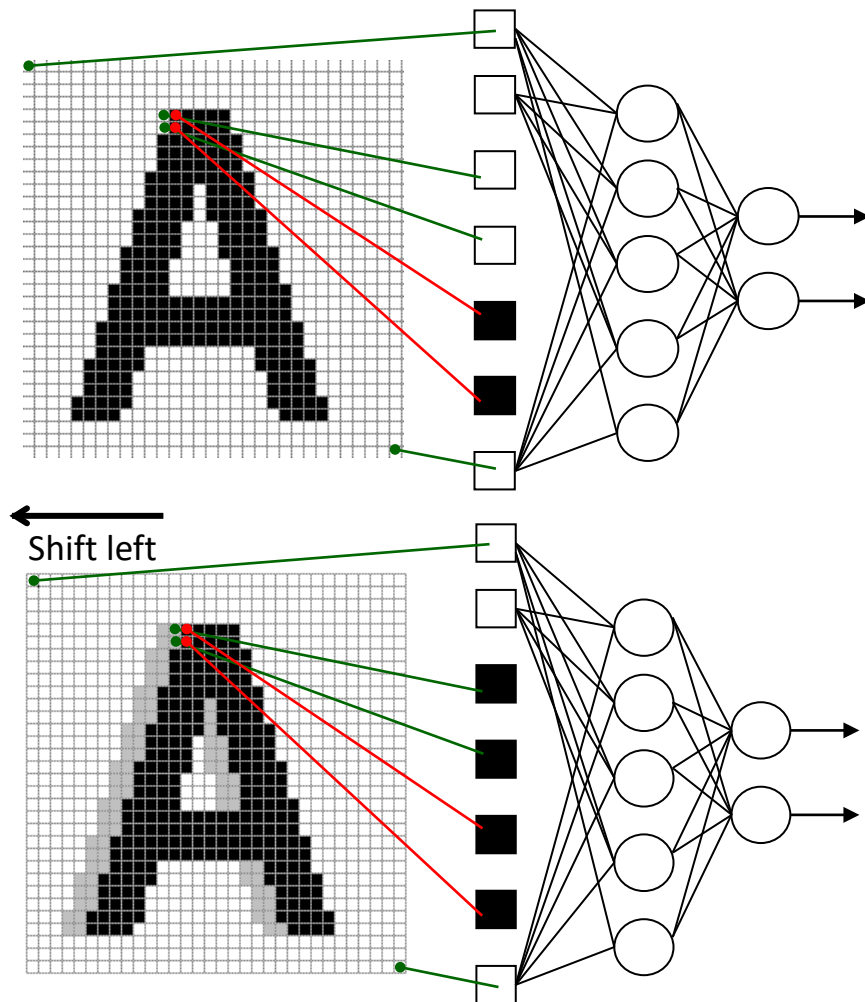
- In a **fully-connected** mode,
 - a neuron in the hidden layer is connected to all neurons in the input layer
 - **no spatial information** is preserved (i.e. if we shuffle the image pixels, there is no difference!)
 - the number of **trainable parameters** becomes extremely **large**



Total trainable parameter number:
 $(256 \times 100) + 100 + (26 \times 100) + 26 = 28326$

Convolution versus Fully Connected

- In a **fully-connected** mode,
 - the number of **trainable parameters** becomes extremely **large**
 - there is **little or no invariance** to shifting, scaling, and other forms of distortion

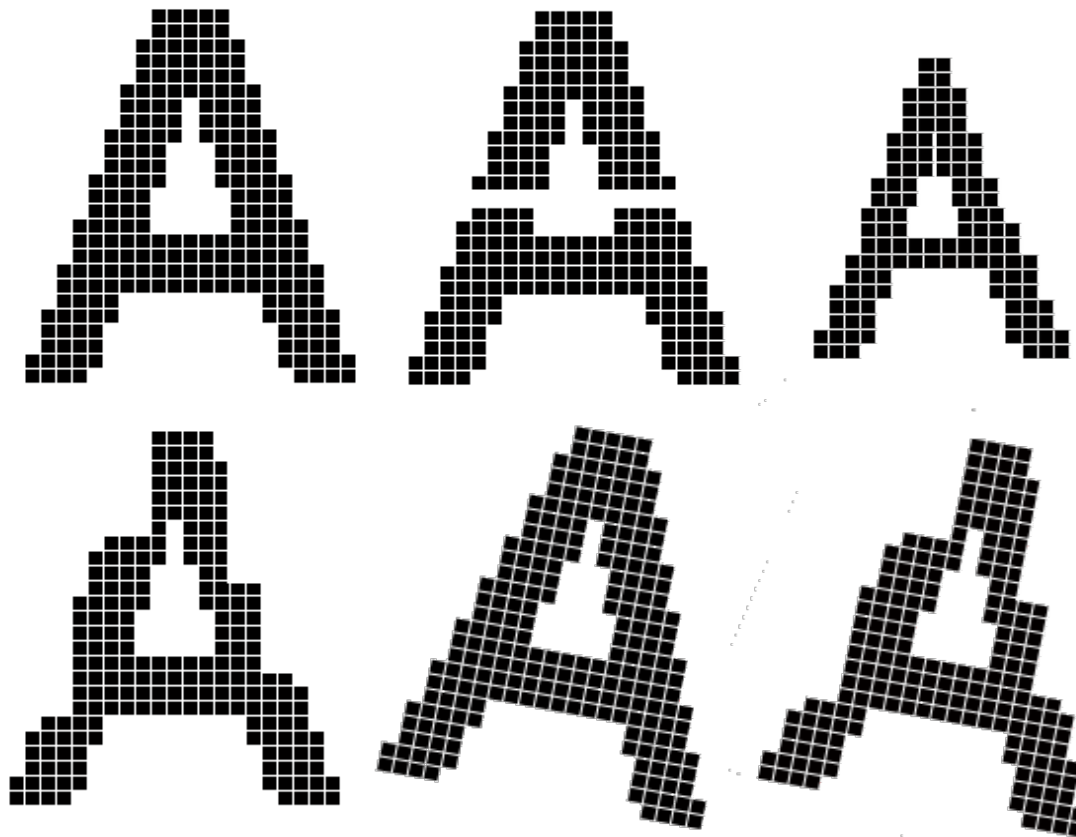


154 input change
from 2 shift left

77 : black to white
77 : white to black

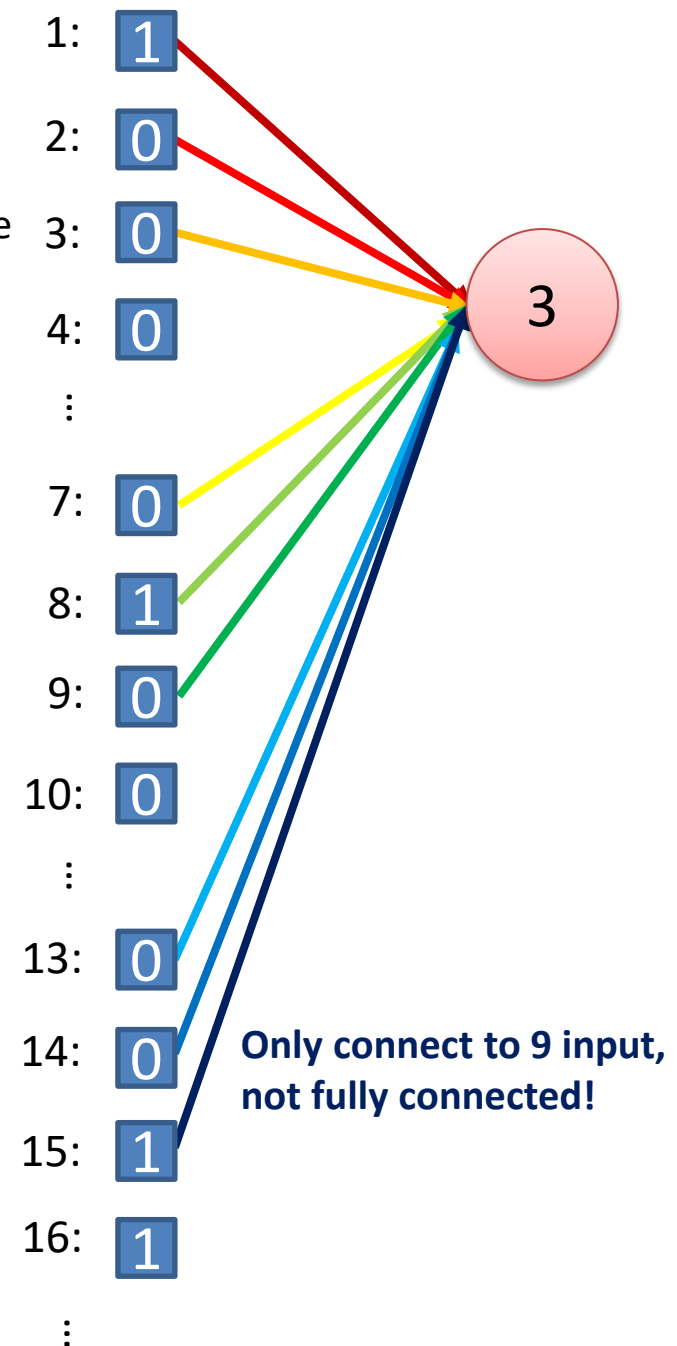
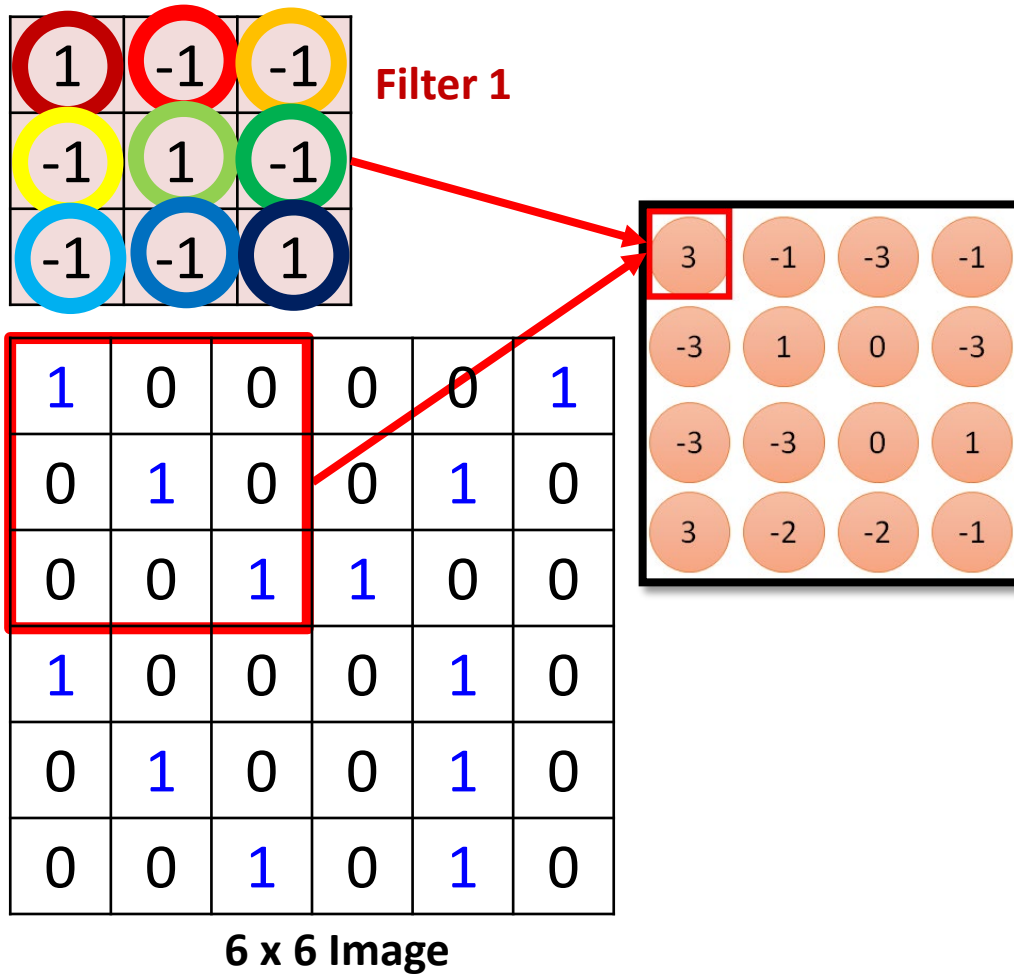
Convolution versus Fully Connected

- In a **fully-connected** mode,
 - the number of **trainable parameters** becomes extremely **large**
 - there is ***little or no invariance*** to shifting, scaling, and other forms of distortion
 - the **topology** of the input data is completely **ignored**



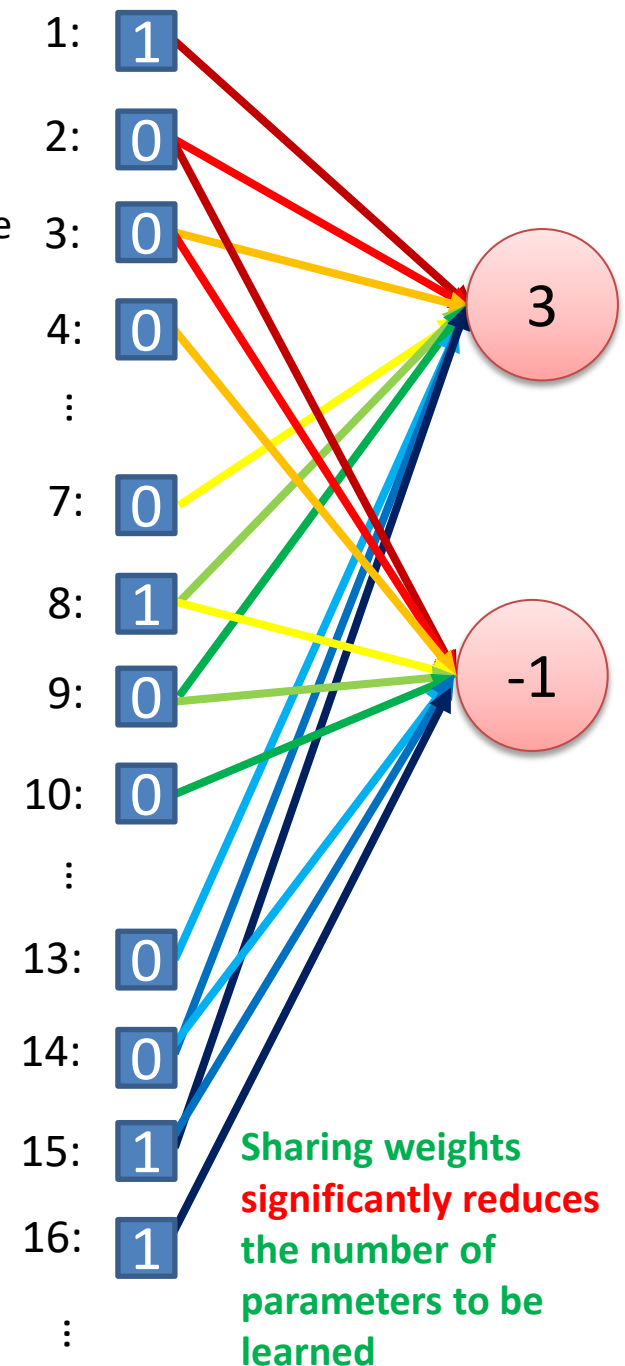
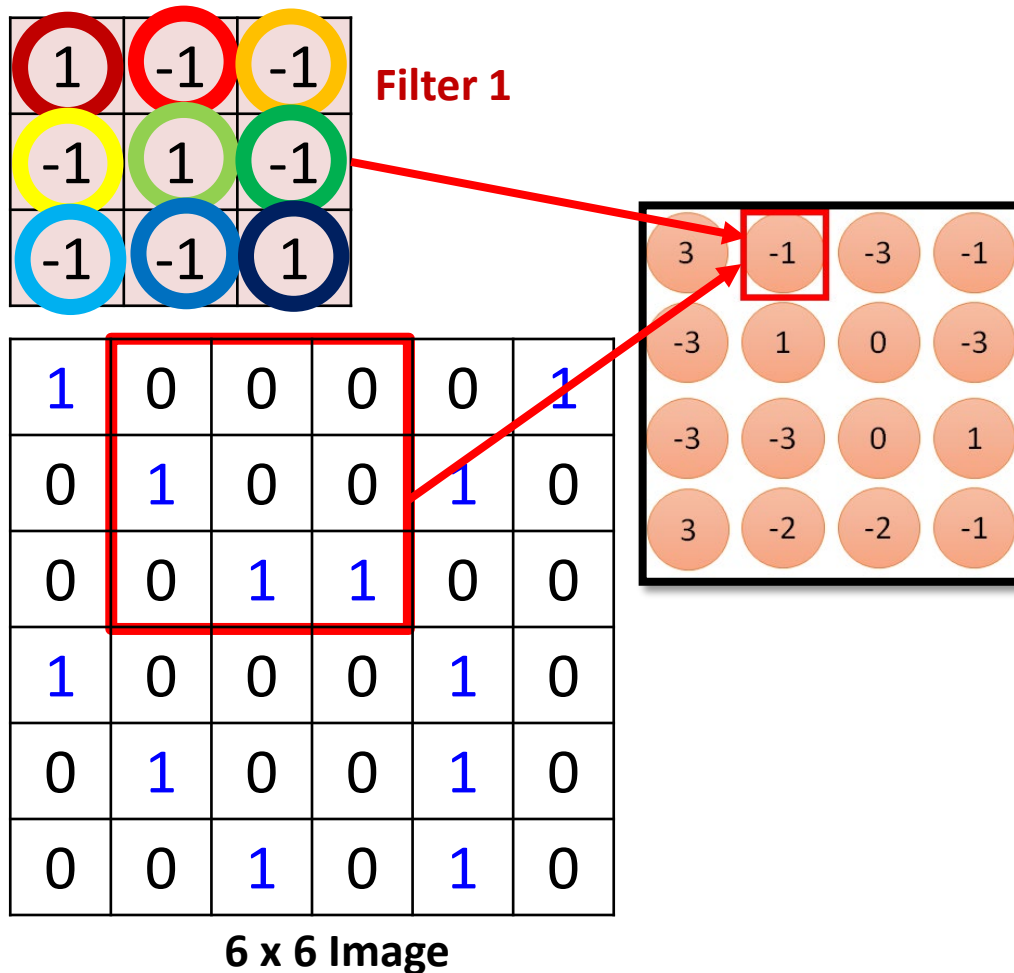
Convolution versus Fully Connected

- In a **convolution** mode, spatial structure is preserved. We have
 - **fewer parameters!**
 - **shared weights!**



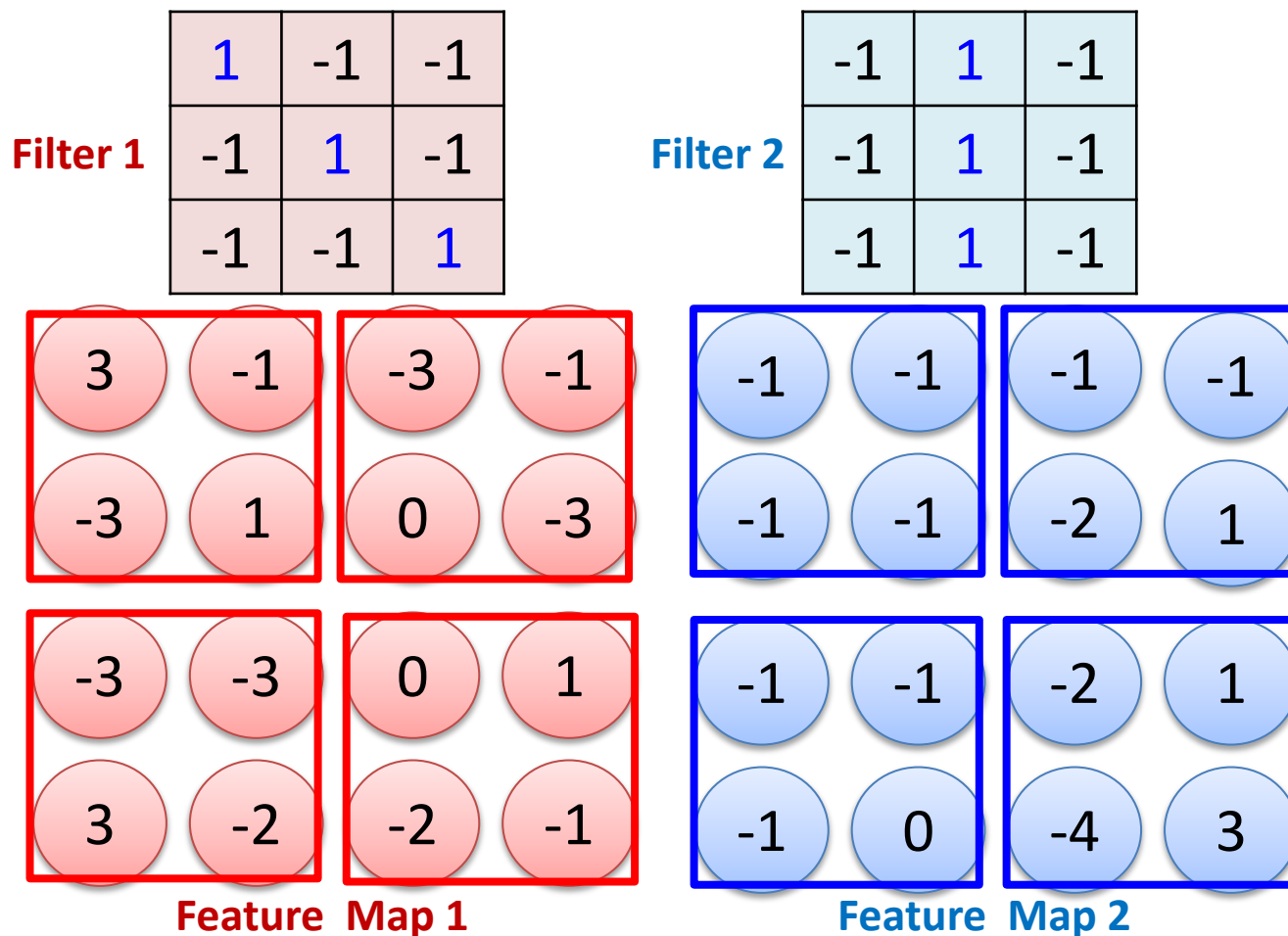
Convolution versus Fully Connected

- In a **convolution** mode, spatial structure is preserved. We have
 - **fewer parameters!**
 - **shared weights!**



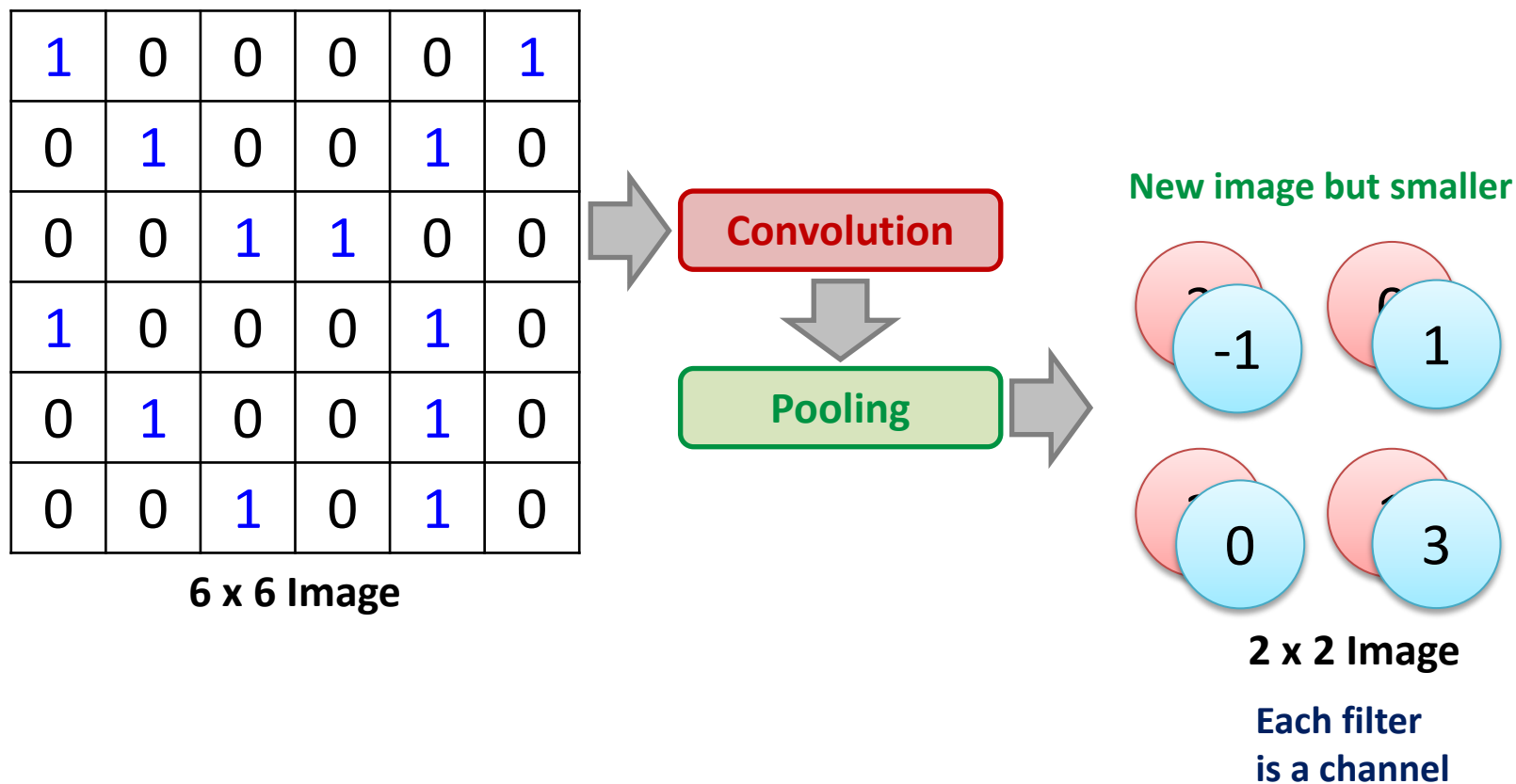
CNN: Pooling

- Pooling layer makes the representation **smaller**, more **manageable**, and **robust to position variations**
- Pooling operates over each feature map (a.k.a. activation map) **independently**
 - **Max-pooling**: takes a max response over spatial locations
 - **Average-pooling**: takes an average response over spatial locations

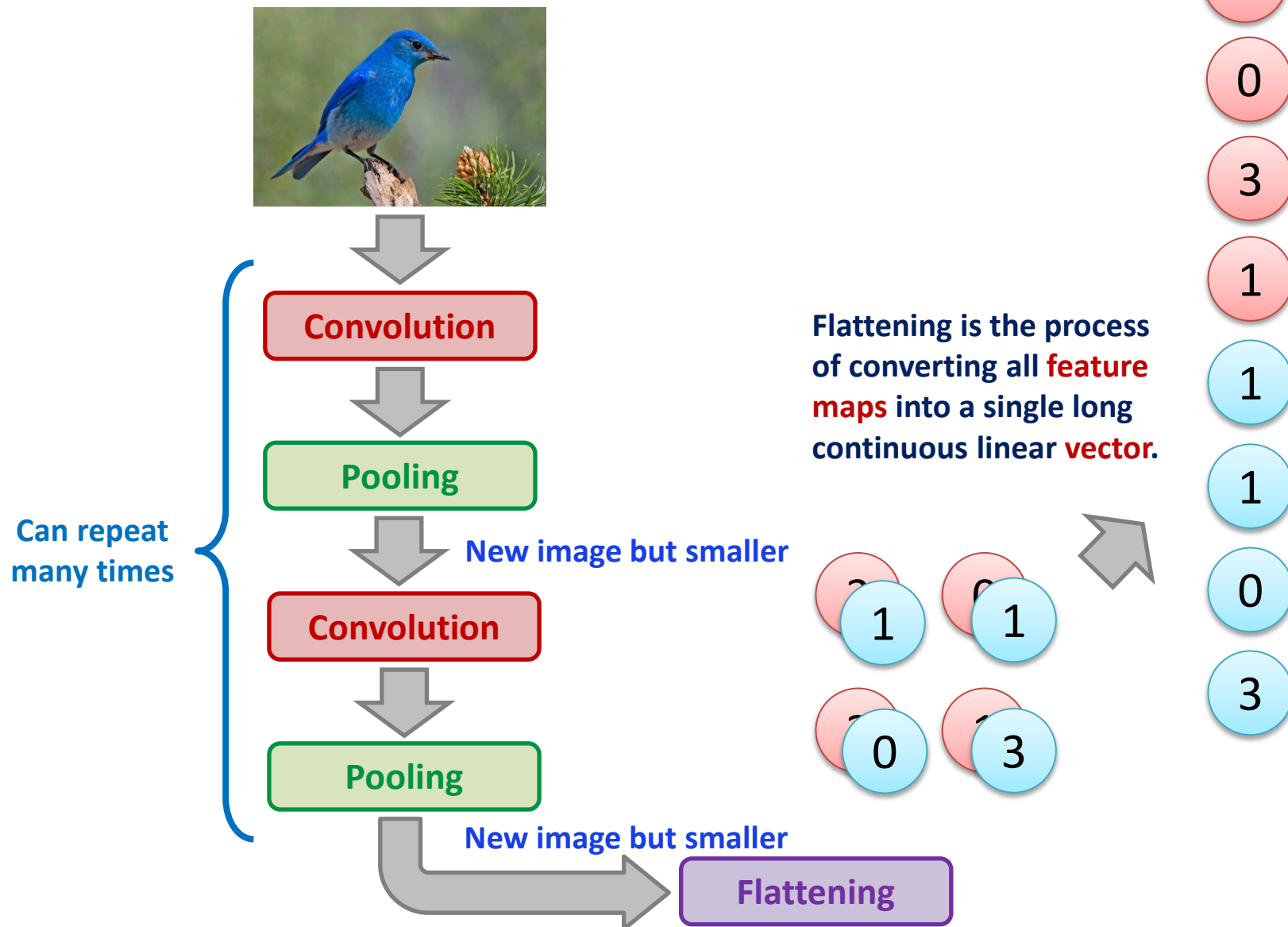


CNN: Pooling

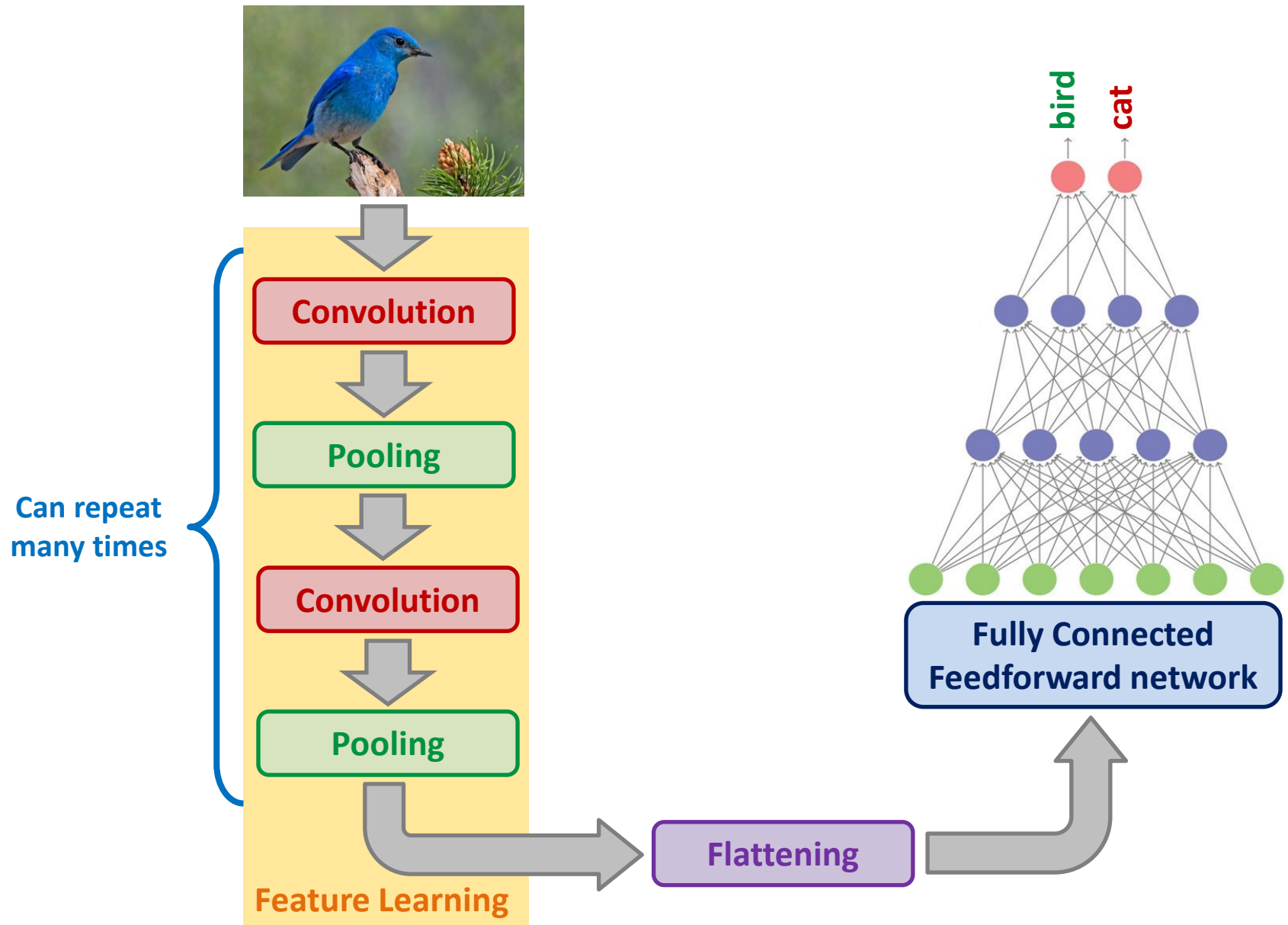
- Pooling layer makes the representation **smaller**, more **manageable**, and **robust to position variations**
- Pooling operates over each feature map (a.k.a. activation map) **independently**
 - **Max-pooling**: takes a max response over spatial locations
 - **Average-pooling**: takes an average response over spatial locations



CNN: Flattening



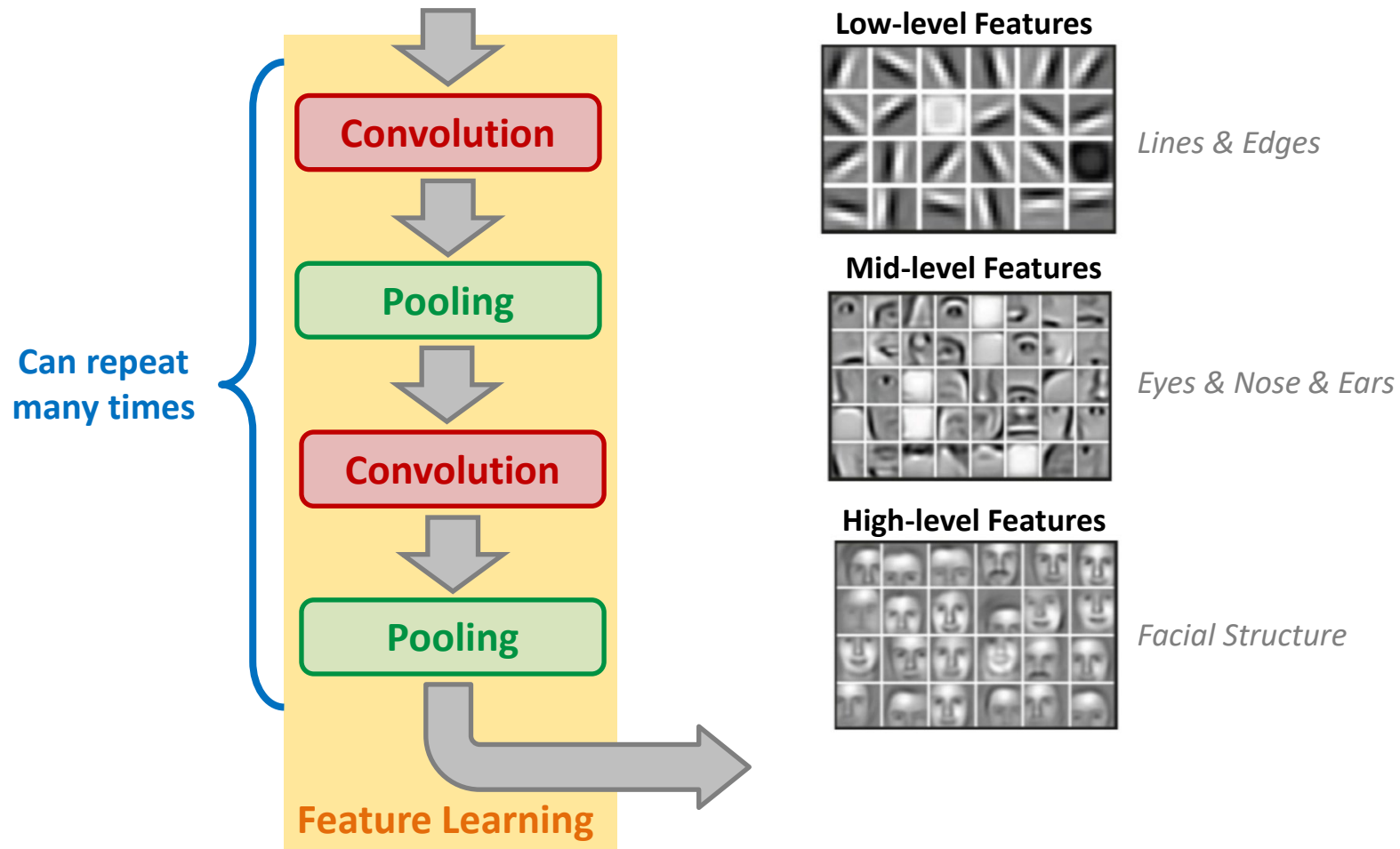
The Whole CNN



The Whole CNN

■ Feature Learning

- Learn a hierarchy of **features** in input image through **convolution**
- Introduce **non-linearity** through activation function (**Conv + ReLU**)
- Reduce dimensionality and preserve spatial invariance with **pooling**

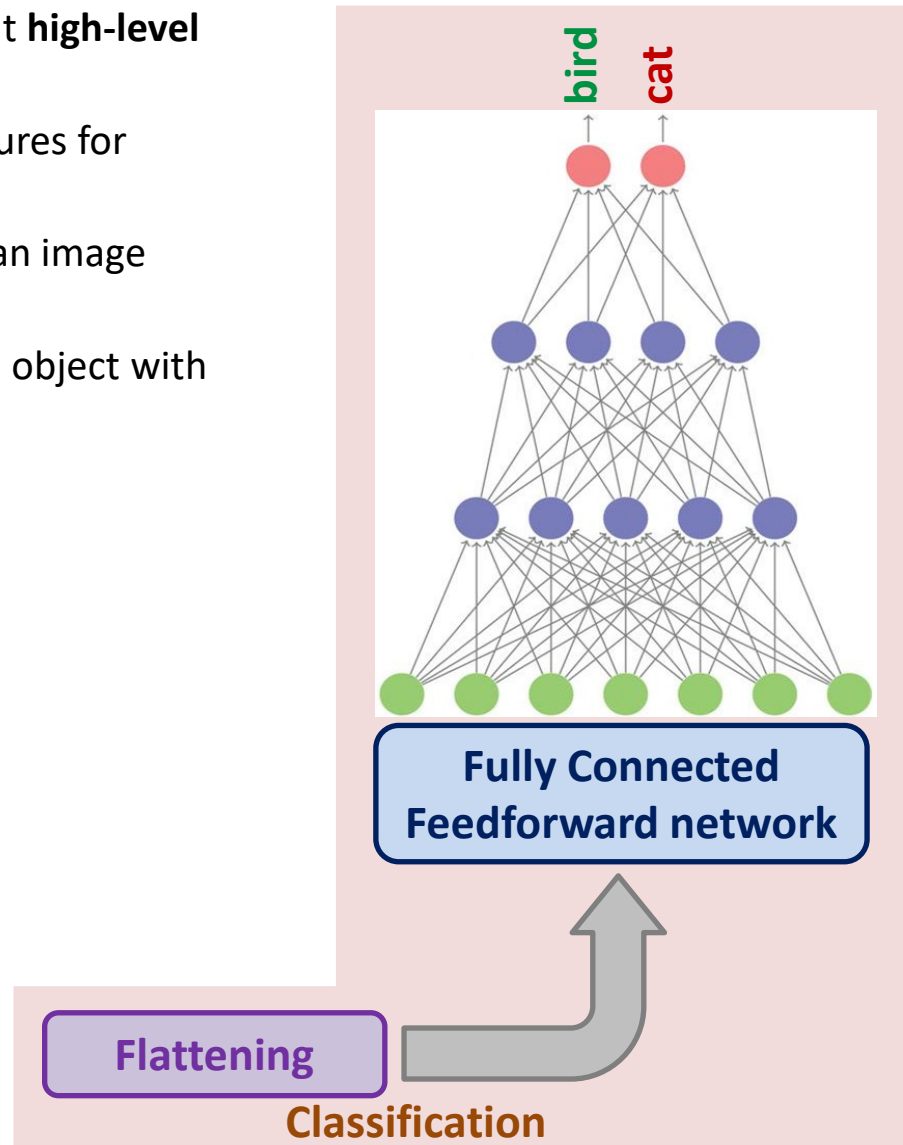


The Whole CNN

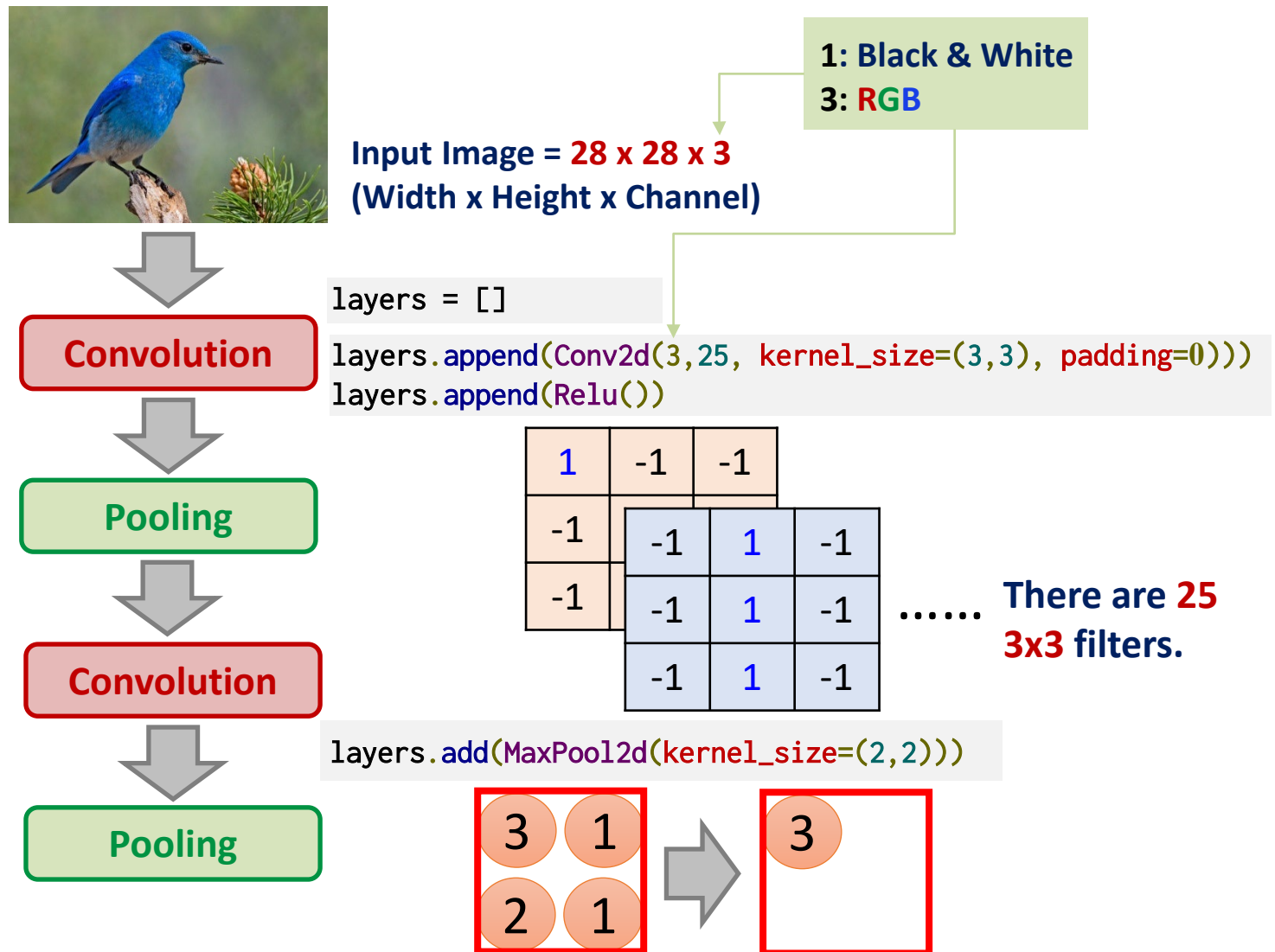
■ Classification

- Convolution and Pooling layers output **high-level features** of the input
- Fully connected layer uses these features for **classifying** input image
- Express output as the **probability** of an image belonging to a particular class
- Apply **Softmax function** to classify an object with probabilistic values between **0 and 1**

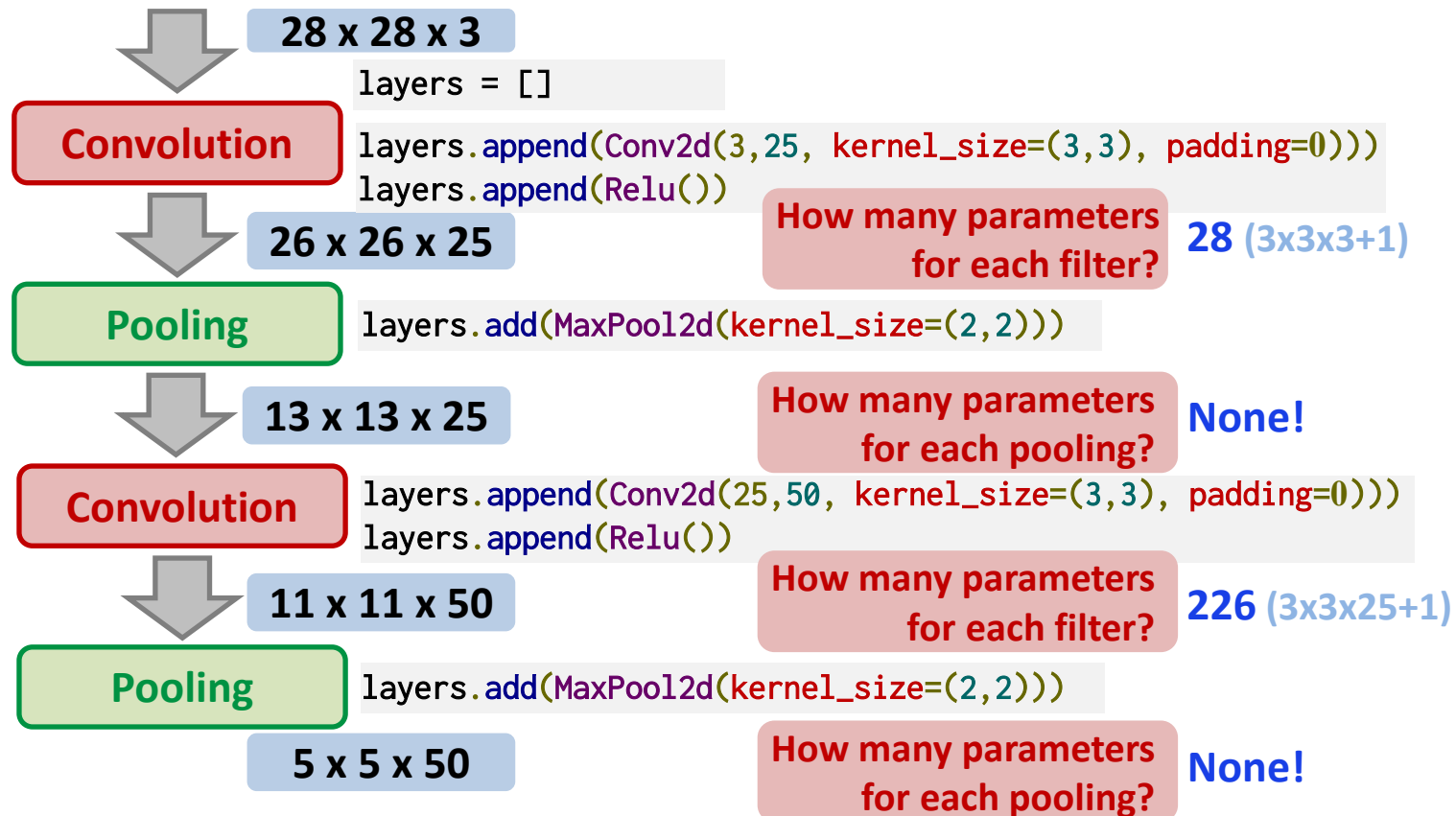
$$\text{softmax}(z_i) = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$



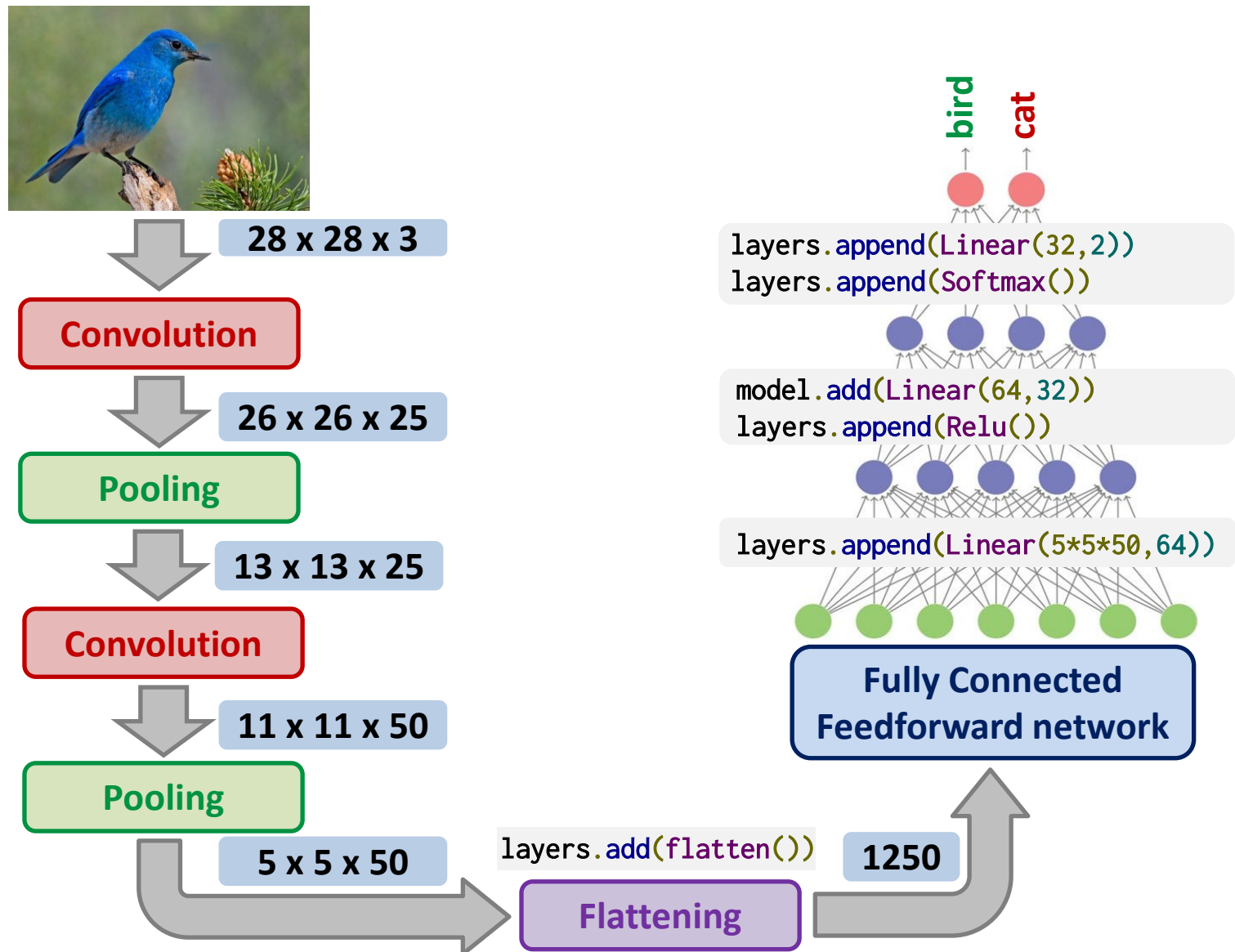
The Whole CNN in PyTorch



The Whole CNN in PyTorch

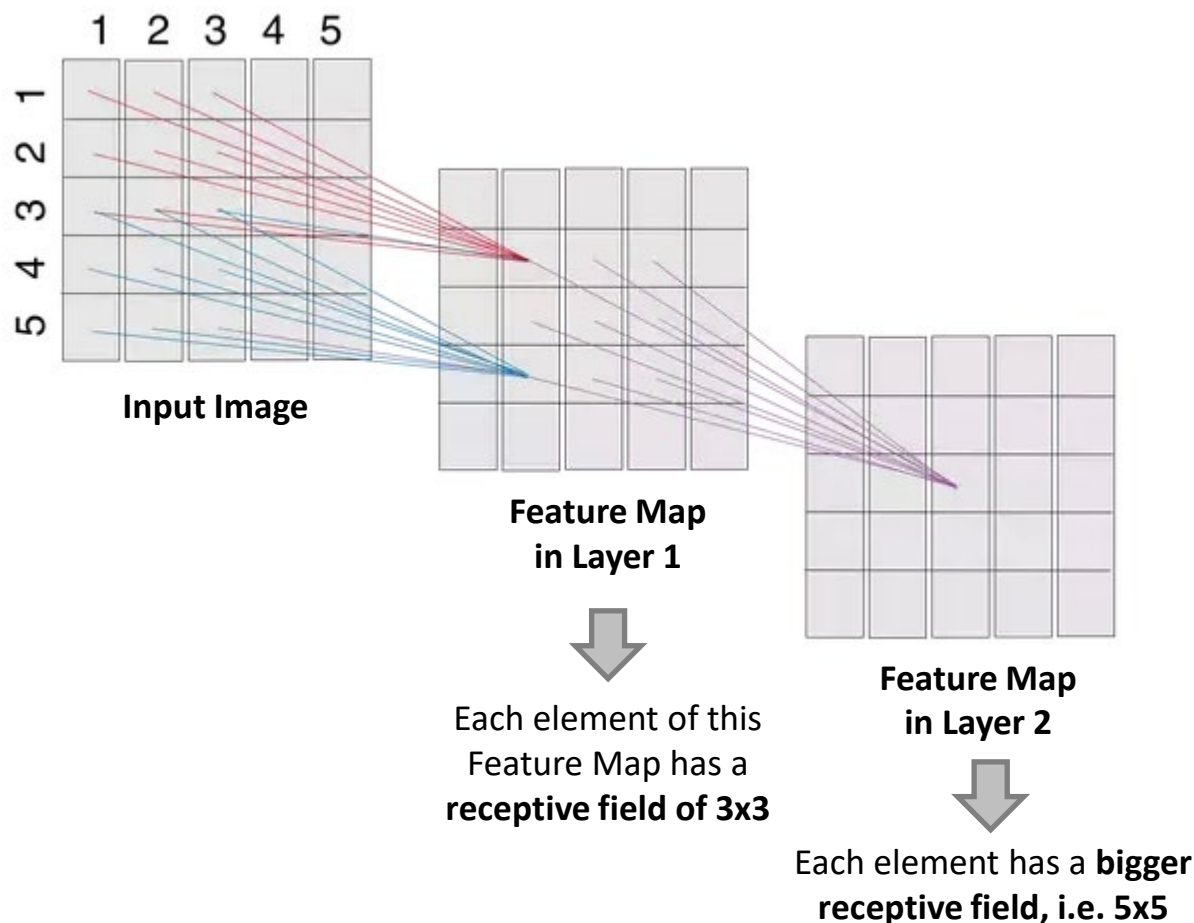


The Whole CNN in PyTorch



CNN: Receptive Fields

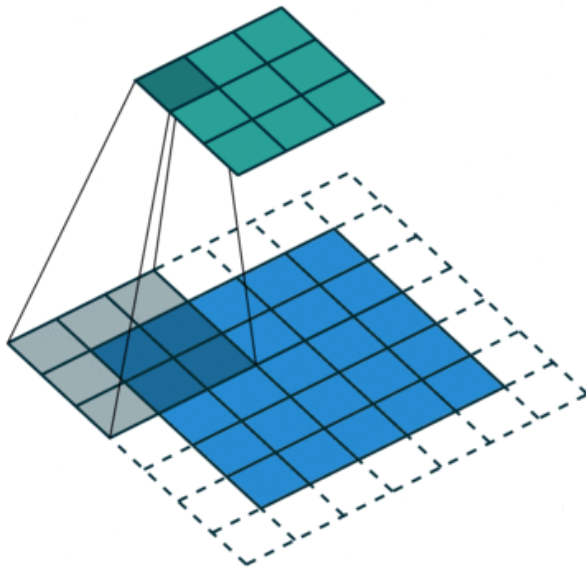
- The **receptive field** is defined as the region in the input space that a particular CNN's feature is looking at (i.e. be affected by)
- With the network going **deeper**, the size of the **receptive field** is becoming **larger**.



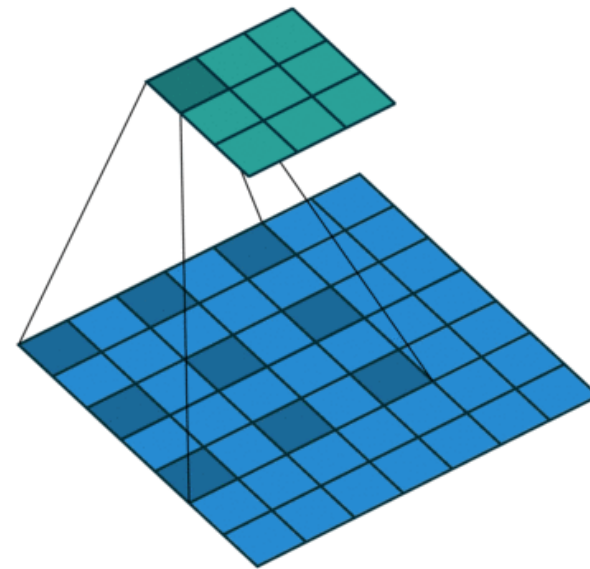
The **receptive field** of vanilla convolutions **grows linearly** with the number of layers (**3x3, 5x5, 7x7**, etc.).

CNN: Dilated Convolutions

- Vanilla convolutions struggle to integrate global context since the effective receptive field of units can only grow linearly with layers.
- The dilation operation supports the exponential expansion of the receptive field without loss of resolution or coverage.
- The receptive field **grows exponentially** while the number of parameters **grows linearly**.
- Dilated convolutions have **gridding artifacts**
 - Neighboring output units are computed from completely separated sets of input units.



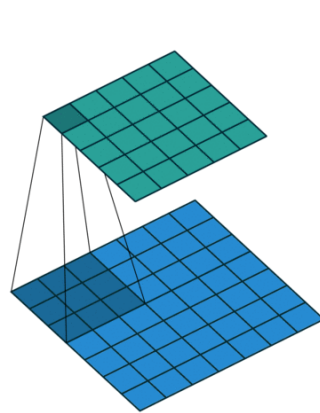
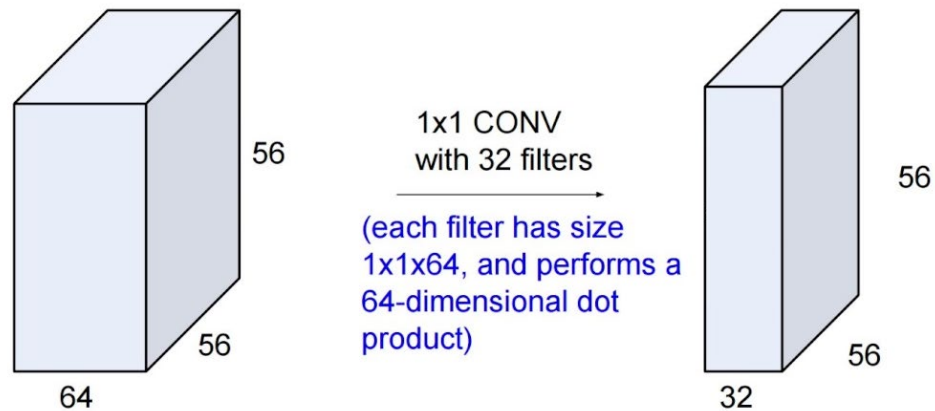
Standard convolution
(the receptive field is 5x5)



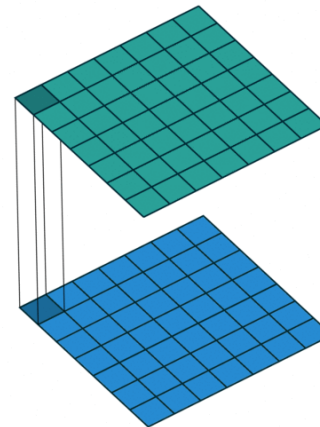
Dilated convolution
(the receptive field is larger, i.e. 7x7)

CNN: 1x1 Convolution

- 1x1 convolution leads to an **increase/decrease in the dimension** of the number of feature maps
- Useful in semantic segmentation to merge different layers

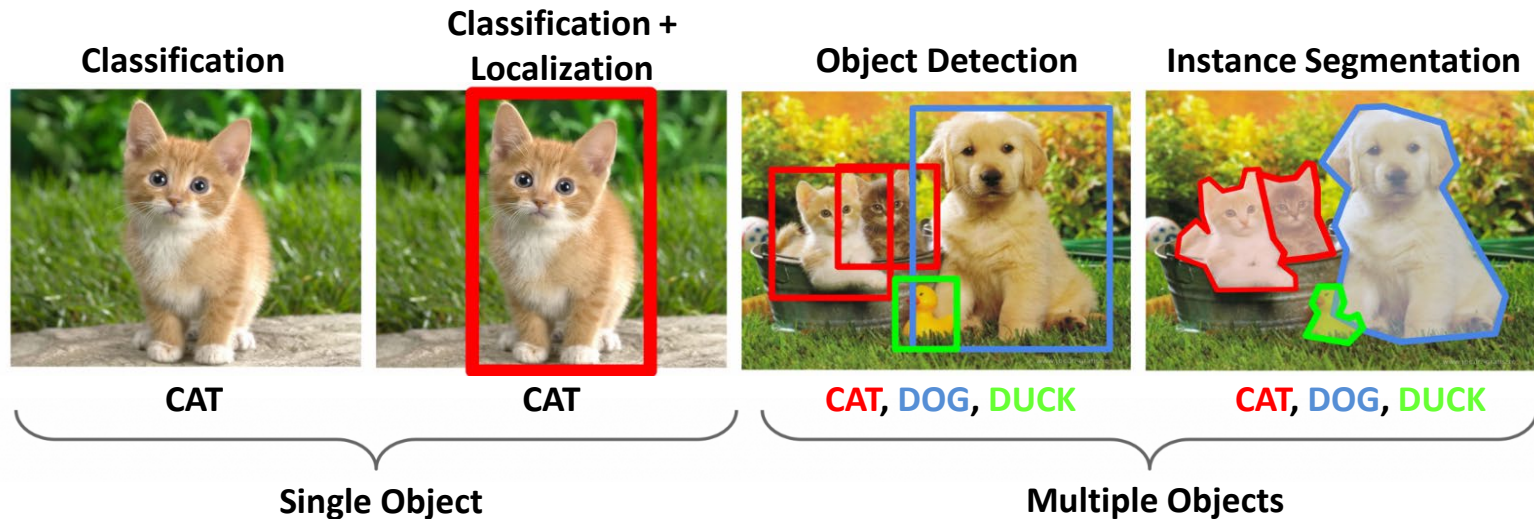


**Convolution with
kernel of size 3x3**



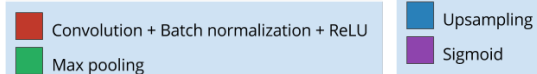
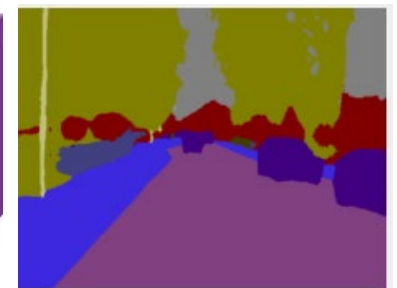
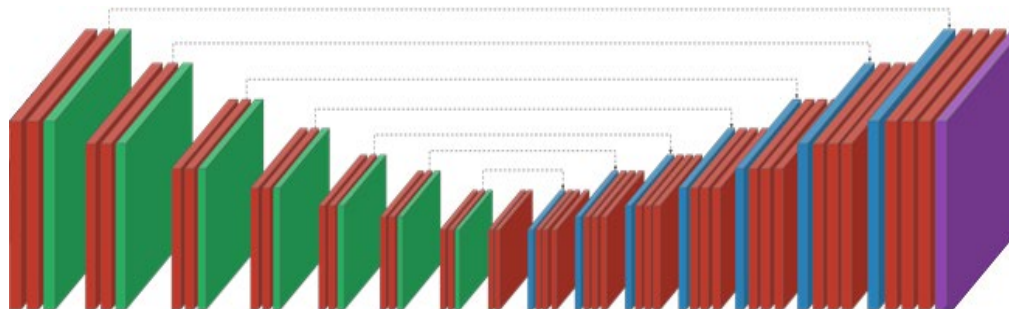
**Convolution with
kernel of size 1x1**

Computer Vision Tasks for CNNs



■ Semantic Image Segmentation with CNNs

- Convolutional Encoder-Decoder structure
- Not instance-aware
- Skip connections



Transfer Learning with CNNs

- **Common “Wisdom”:** You need a lot of data to train a CNN
- **Solution: Transfer learning**
 - taking a model trained on a similar task that has lots of data and adopting this model to the task that may not have enough training data
 - This strategy is Pervasive!

Train on ImageNet



Why on ImageNet?

- Convenience, lots of **data**
- We know how to **train these well**

Note that, for some other tasks we would need to start with something else (e.g., videos for optical flow)

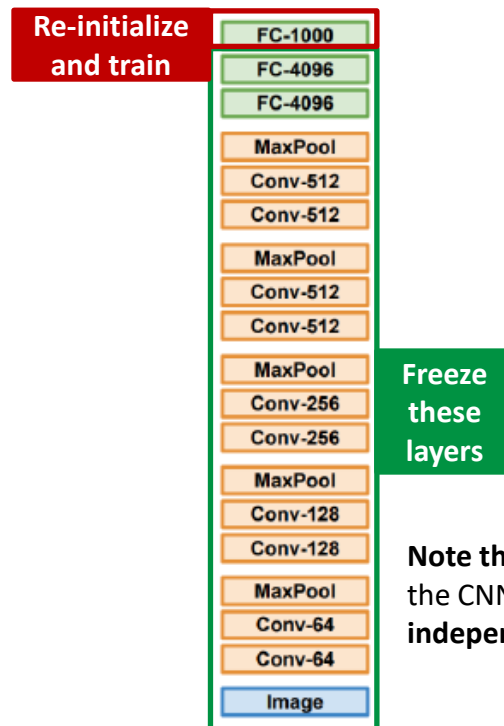
Transfer Learning with CNNs

- **Common “Wisdom”:** You need a lot of data to train a CNN
- **Solution: Transfer learning**
 - taking a model trained on a similar task that has lots of data and adopting this model to the task that may not have enough training data
 - This strategy is Pervasive!

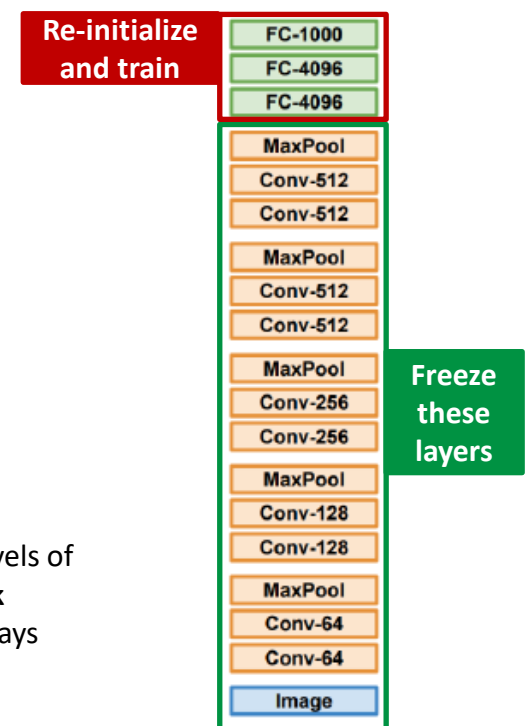
Train on **ImageNet**



Small dataset with C classes



Larger dataset



Note that lower levels of the CNN are at **task independent** anyways

Conclusion

- A Convolutional Neural Network (CNN) is a special kind of **multi-layer feed-forward neural network** that can extract topological properties from an image.
- CNNs are **neurobiologically** motivated by the findings of locally sensitive and orientation-selective nerve cells in the visual cortex.
- CNNs can **recognize visual patterns** with extreme variability.
- If a **neuron** in the feature map **fires**, this corresponds to a **match with the template**.
- **Filter indices** in CNNs **are learned** directly by the network.
- Pooling (**subsampling**) layers reduce the spatial resolution of each feature map
 - By reducing the **spatial resolution** of the feature map, a **certain degree** of **shift** and **distortion invariance** is achieved.
- Like almost every other neural network, CNNs **are trained** with a version of the **back-propagation algorithm**.

Discussion

- Which of the following is the correct order for the typical CNN operation?

- Convolution -> max pooling -> flattening -> full connection



- Max pooling -> convolution -> flattening -> full connection
 - Flattening -> max pooling -> convolution -> full connection
 - None

- What are the benefits of using convolutional layers instead of having fully connected ones for image processing tasks?

- Convolutional layers use
 - spatial context by only assigning weights to nearby pixels
 - translation invariance
 - lot fewer parameters due to weight sharing.

- What is the resulting image size if we apply a 3x3 convolution to an image with a size of 150x150? (no padding!)

- 148x148

Discussion

- What is the resulting image size if we apply a pooling operation (with the size of 2x2) to an image with a size of 150x150?

- 75x75

- Assume that we have an input tensor with the size of 128x128x16. How many parameters would a single 1x1 convolutional filter have, including the bias?

- 16 (depth channel only) + 1 (bias) = 17

- Assume that we have a binary classification task of categorizing images as *dog* or *not dog*. We implement a CNN with a single output neuron. Let's assume that the output of this neuron is z . The final output of our network, \hat{y} is then given by:

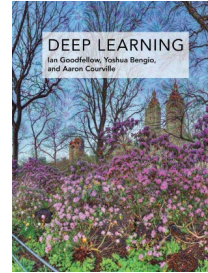
$$\hat{y} = \textit{sigmoid}(\textit{ReLU}(z))$$

We also have a threshold of 0.5 to classify all images as a dog, if their final value is $\hat{y} \geq 0.5$. What problem are we going to face in this case?

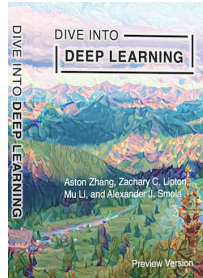
- Since we first use ReLU and then sigmoid all the predictions will be positive. Note that ReLU response will be ***always*** ≥ 0 . When you pass this to Sigmoid ($\mathcal{S}(x) = \frac{1}{1+e^{-x}}$), the output will ***always be*** ≥ 0.5 . Therefore, $\textit{sigmoid}(\textit{ReLU}(z)) \geq 0.5 \forall z$.

Reading Material

- Deep Learning Book
 - Chapter [9.1](#), [9.2](#), [9.3](#), [9.4](#)



- Dive into Deep Learning Book
 - Chapter [7](#)



- Video
 - Please check this video explaining how CNN features are learned in detail
<https://www.youtube.com/watch?v=N6wn8zMRIVE>
- Blogs
 - <https://medium.com/@subham.tiwari186/understanding-convolutional-neural-networks-cnns-72b62b17dd3a>

- CNN Explainer:
 - <https://poloclub.github.io/cnn-explainer/>

