

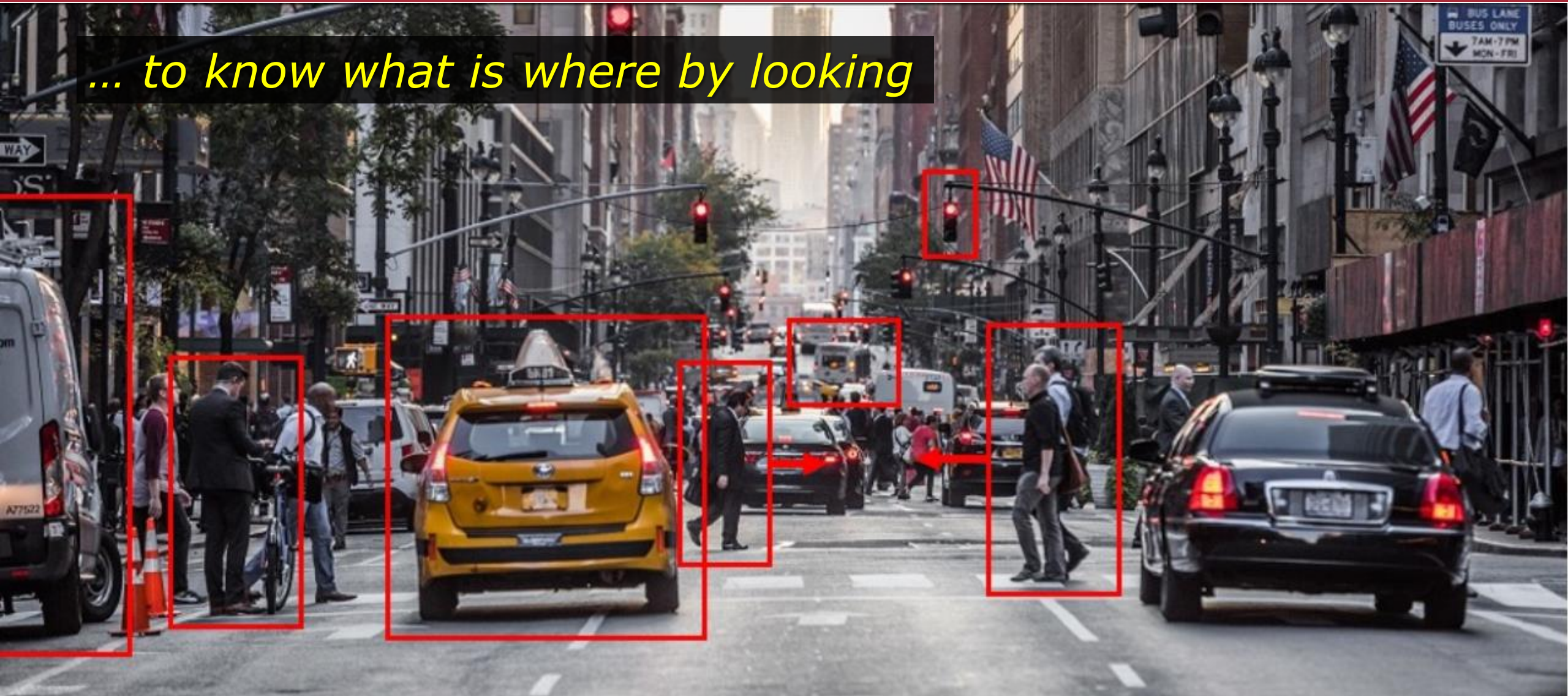


Computer Vision

CV Intro & CNNs

Computer Vision

... to know what is where by looking



The rise and impact of computer vision

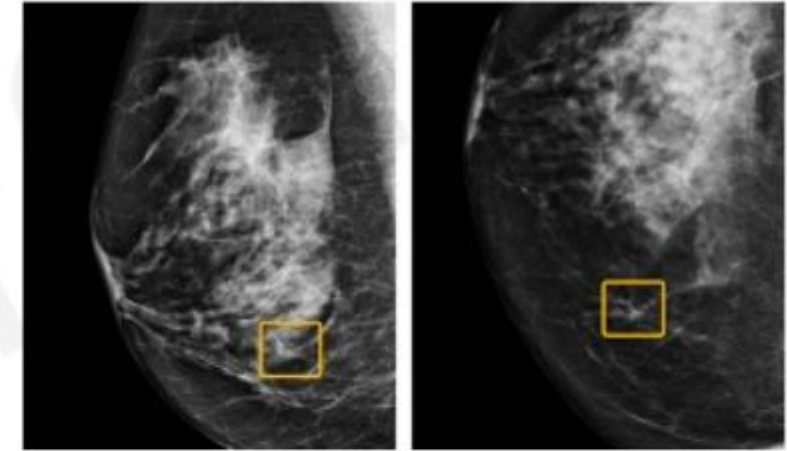
Robotics



Accessibility



Biology & Medicine



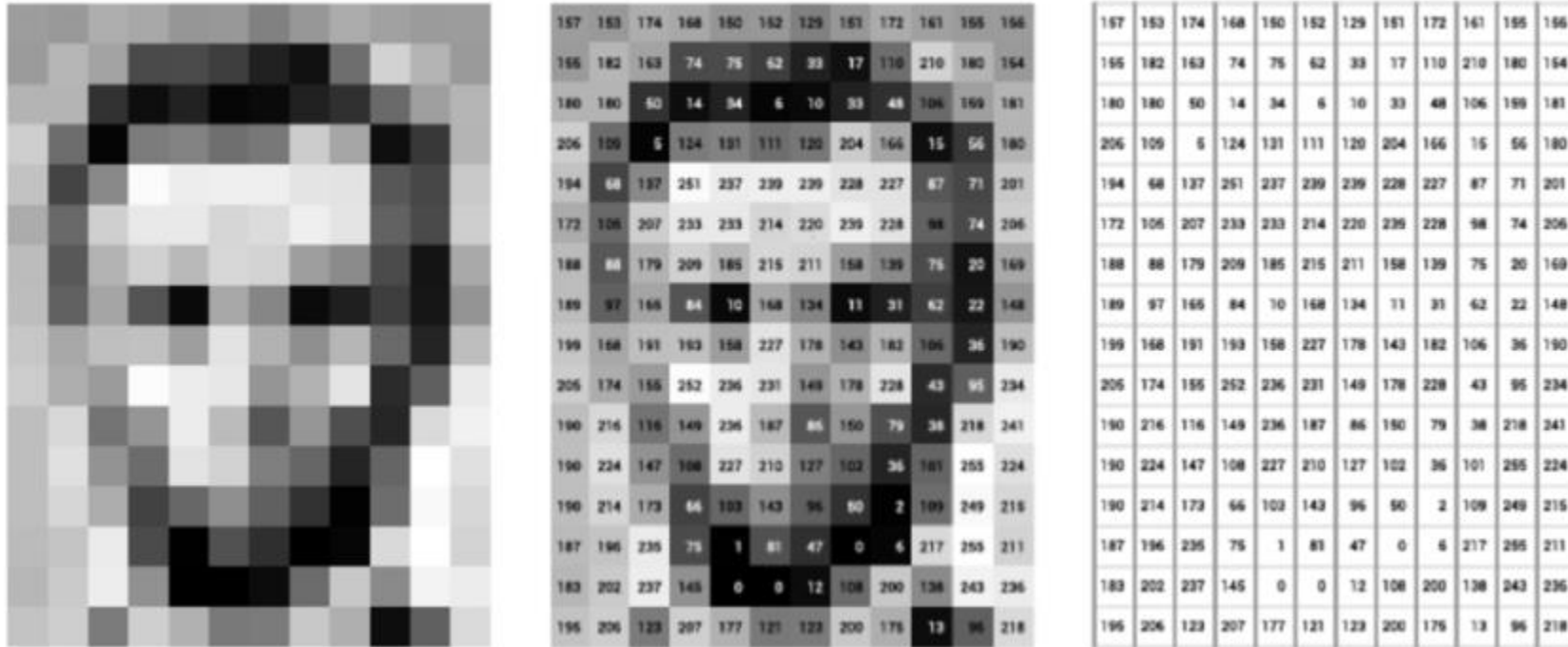
Autonomous driving



Mobile computing

What do computers see?

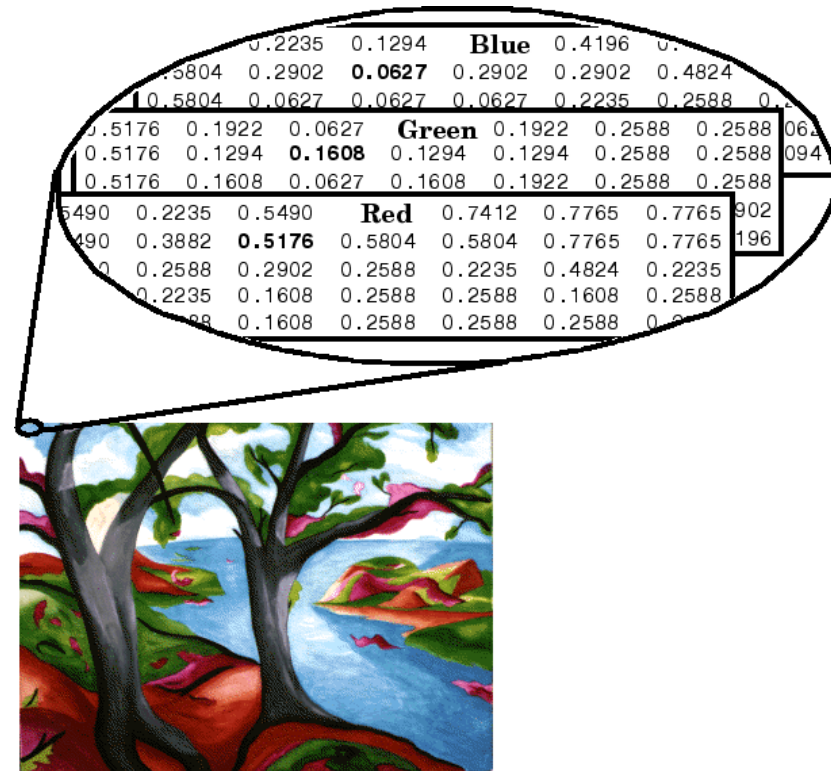
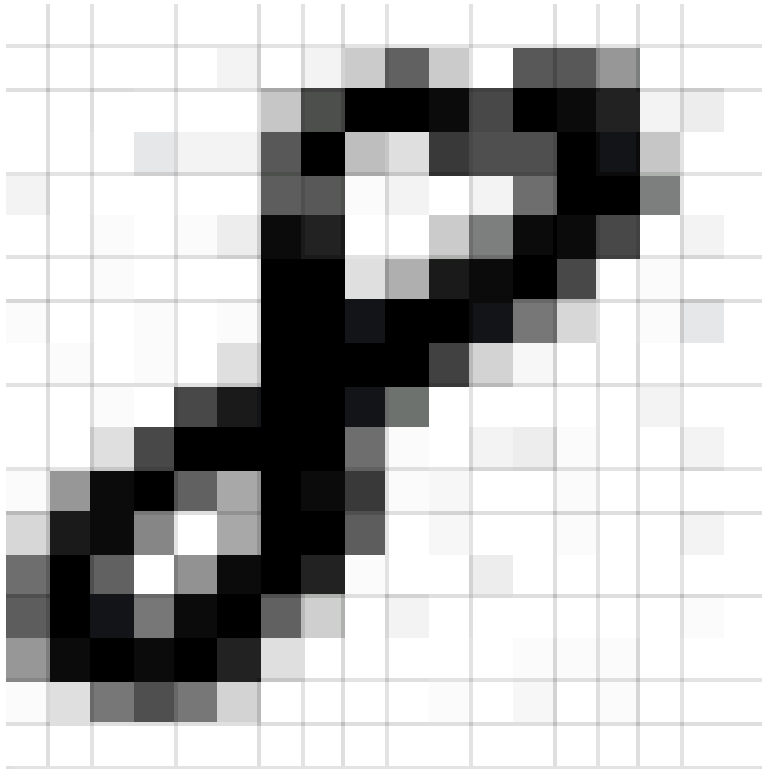
- Images are numbers



An image is just a matrix of numbers $[0, 255]$
i.e. $1080 \times 1080 \times 1$ for a grayscale image

What do computers see?

- Images as matrices



An image is just a matrix of numbers [0, 255]
i.e. 1080 x 1080 x 3 for an RGB image

High level feature detection

- Lets identify key features in each image category



Nose,
Eyes,
Mouth



Wheels,
License Plate,
Headlights



Door,
Windows,
Steps

Manual feature extraction

Domain Knowledge

Define Features

Detect features to
classify

Viewpoint variation



Scale variation



Deformation



Occlusion



Illumination conditions



Background clutter



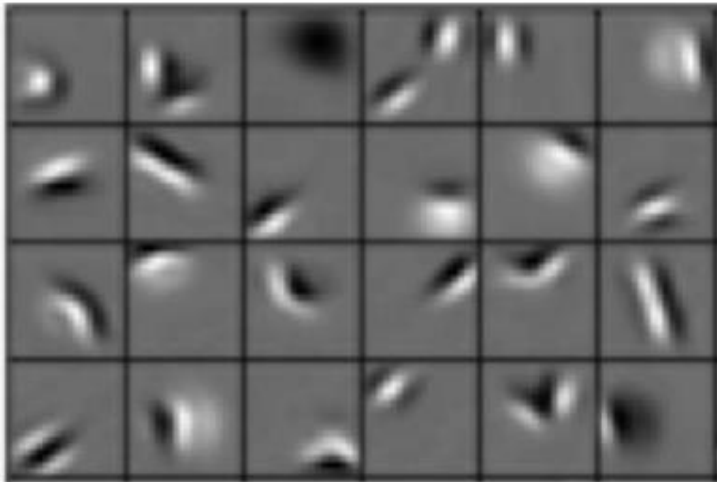
Intra-class variation



Learning feature representations

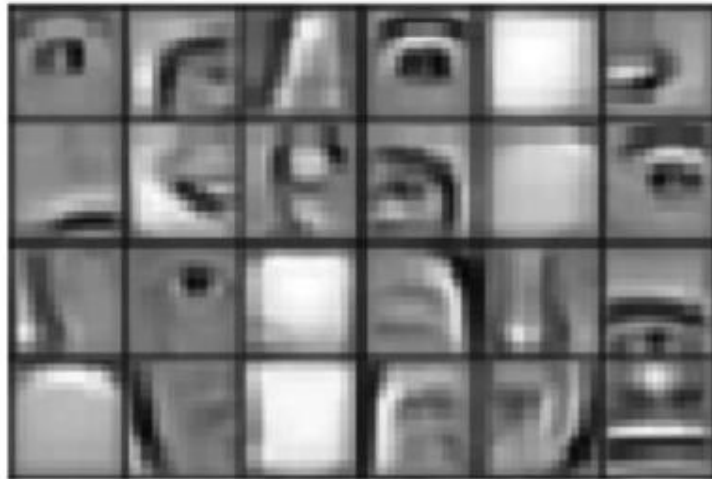
- Can we learn a **hierarchy of features** directly from the data instead of hand engineering?

Low level features



Edges, dark spots

Mid level features



Eyes, ears, nose

High level features



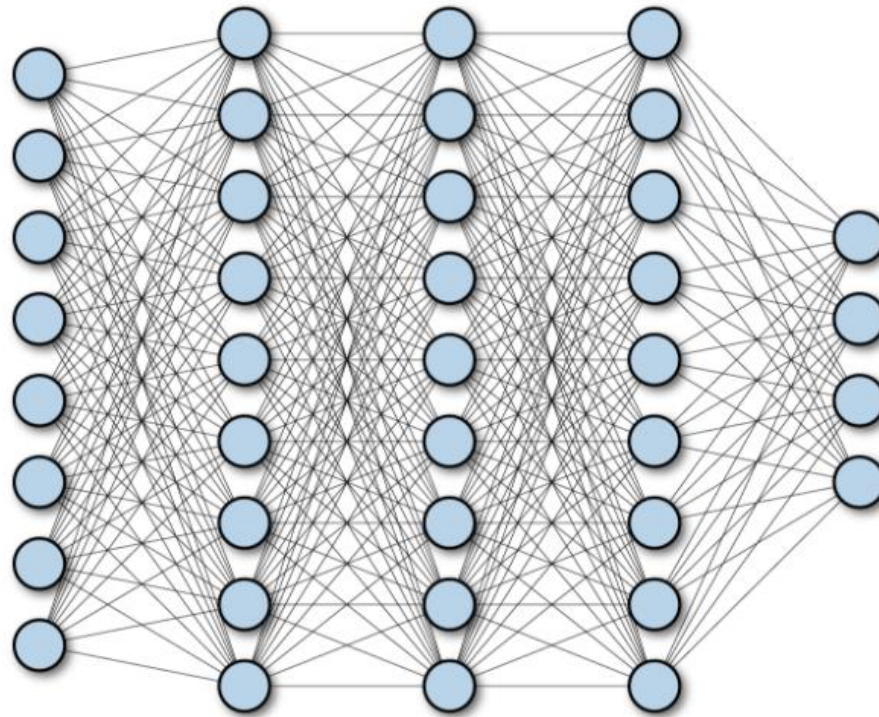
Facial structure

Learning visual features

- Fully connected neural network

Input:

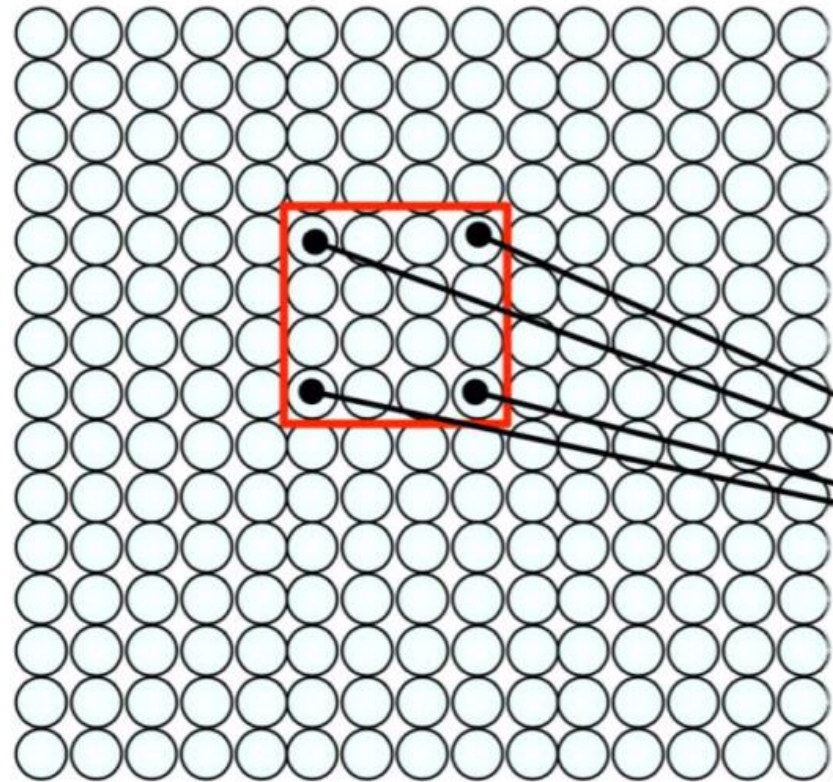
- 2D image
- Vector of pixel values



Learning visual features

- Using spatial structure

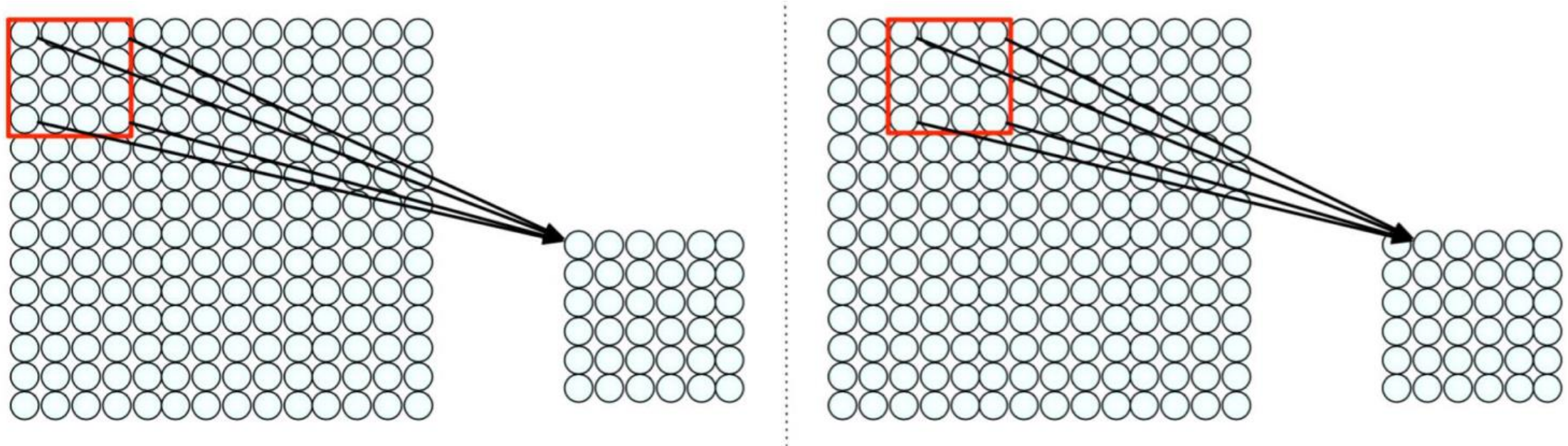
Input: 2D image.
Array of pixel values



Idea: connect patches of input
to neurons in hidden layer.
Neuron connected to region of
input. Only “sees” these values.

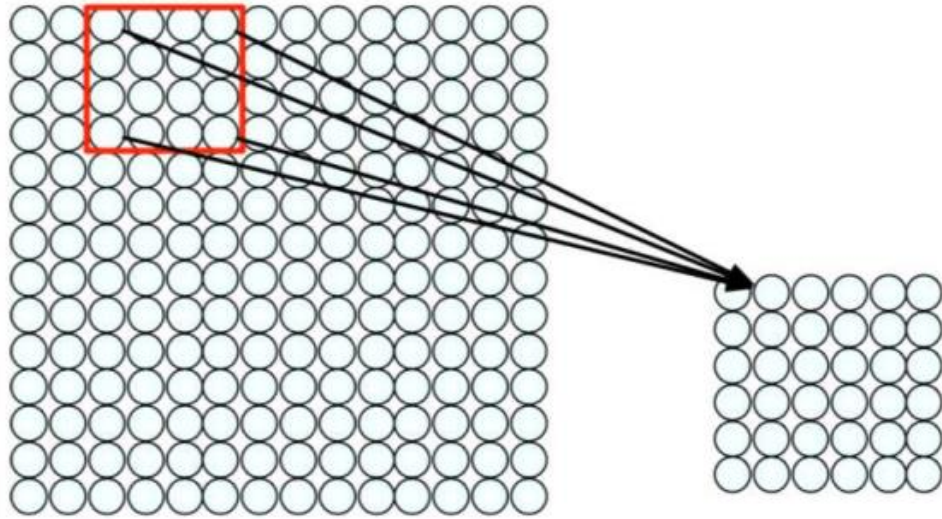
Learning visual features

- Using spatial structure



Connect patch in input layer to a single neuron in subsequent layer.
Use a sliding window to define connections.

Feature extraction with convolution



- Filter of size 4x4 : 16 different weights
- Apply this same filter to 4x4 patches in input
- Shift by 2 pixels for next patch

This “patchy” operation is **convolution**

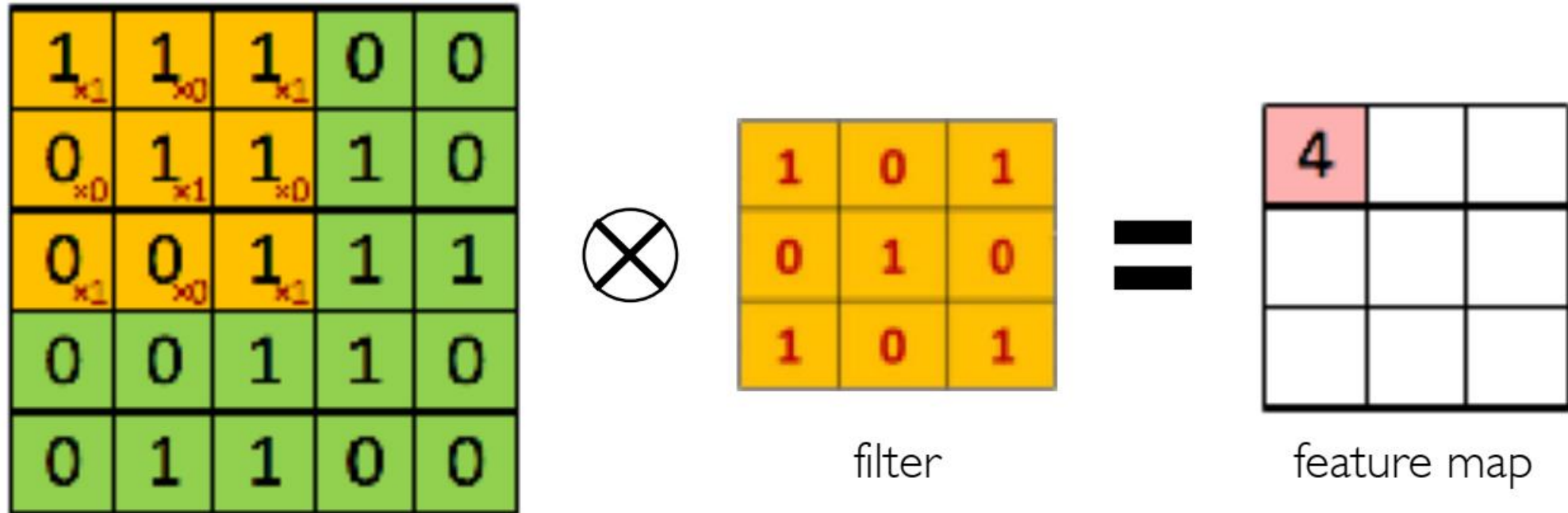
1) Apply a set of weights – a filter – to extract **local features**

2) Use **multiple filters** to extract different features

3) **Spatially share** parameters of each filter

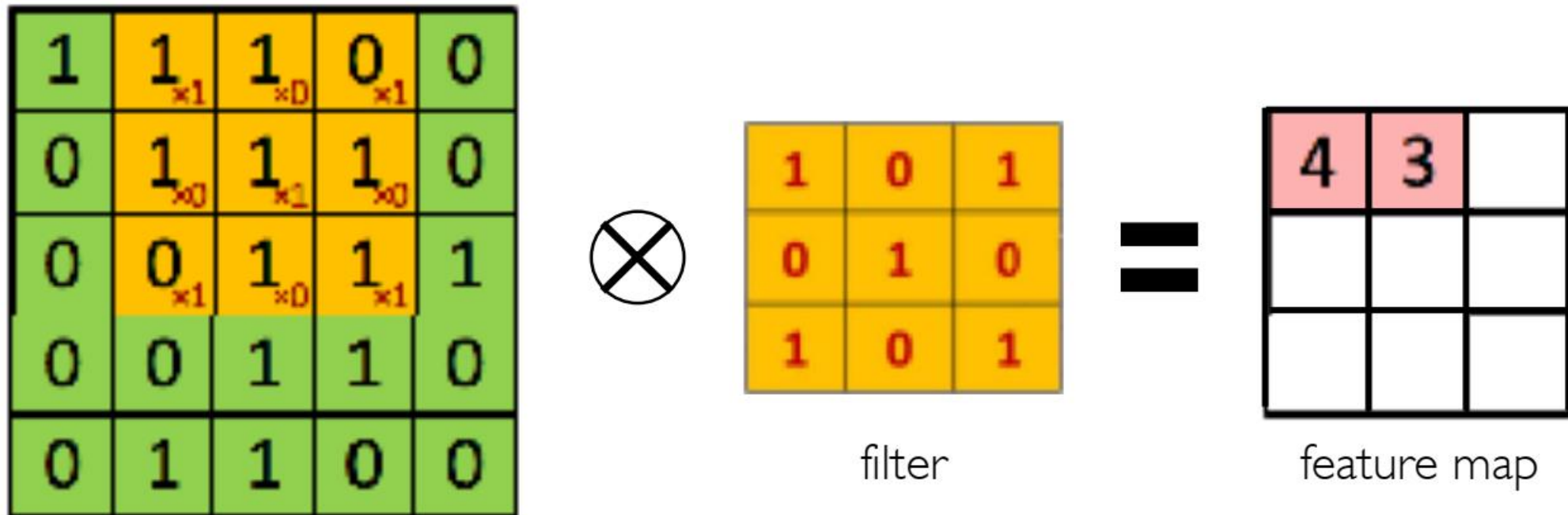
The convolution operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



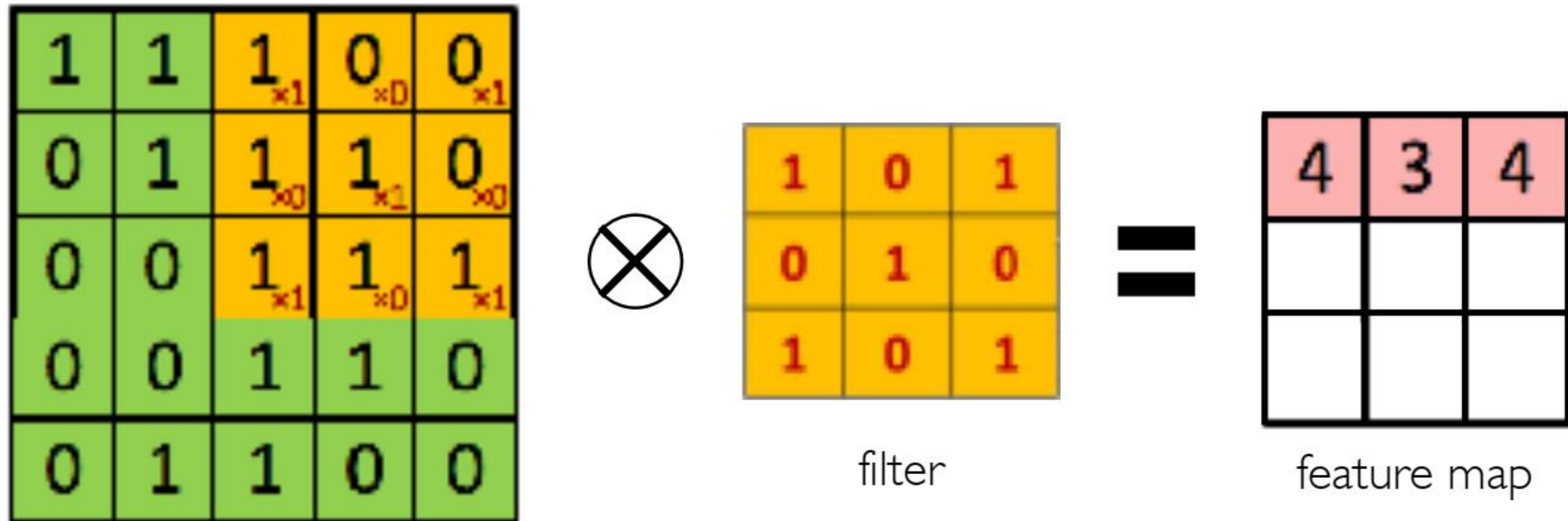
The convolution operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



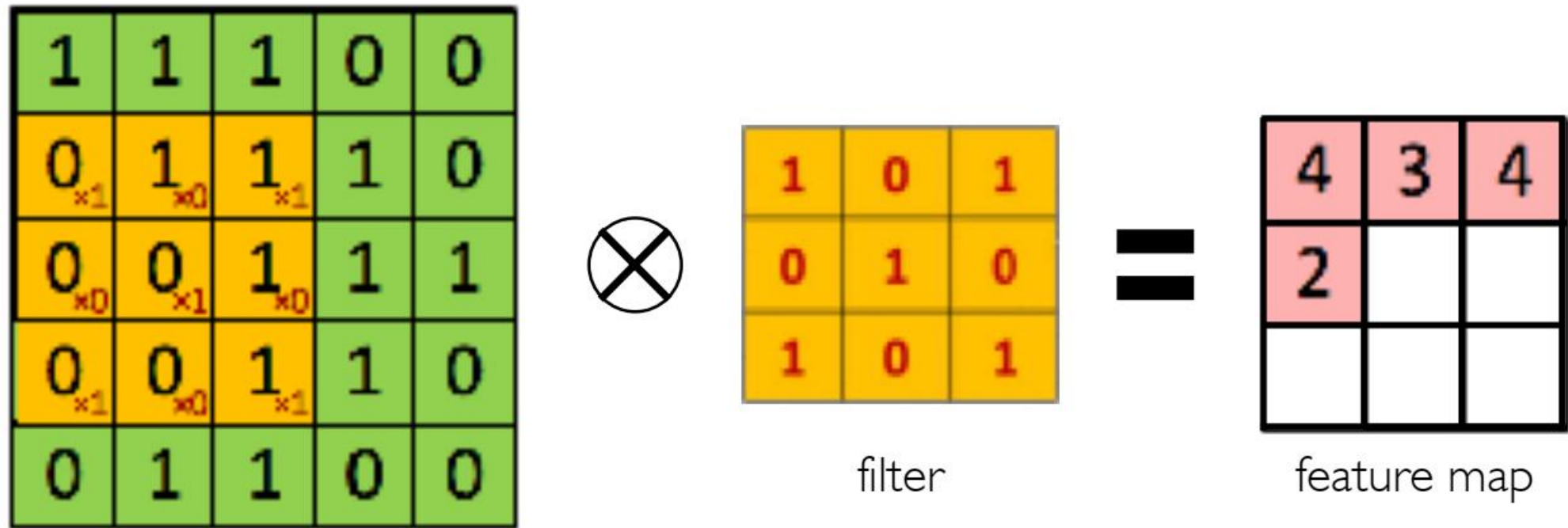
The convolution operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



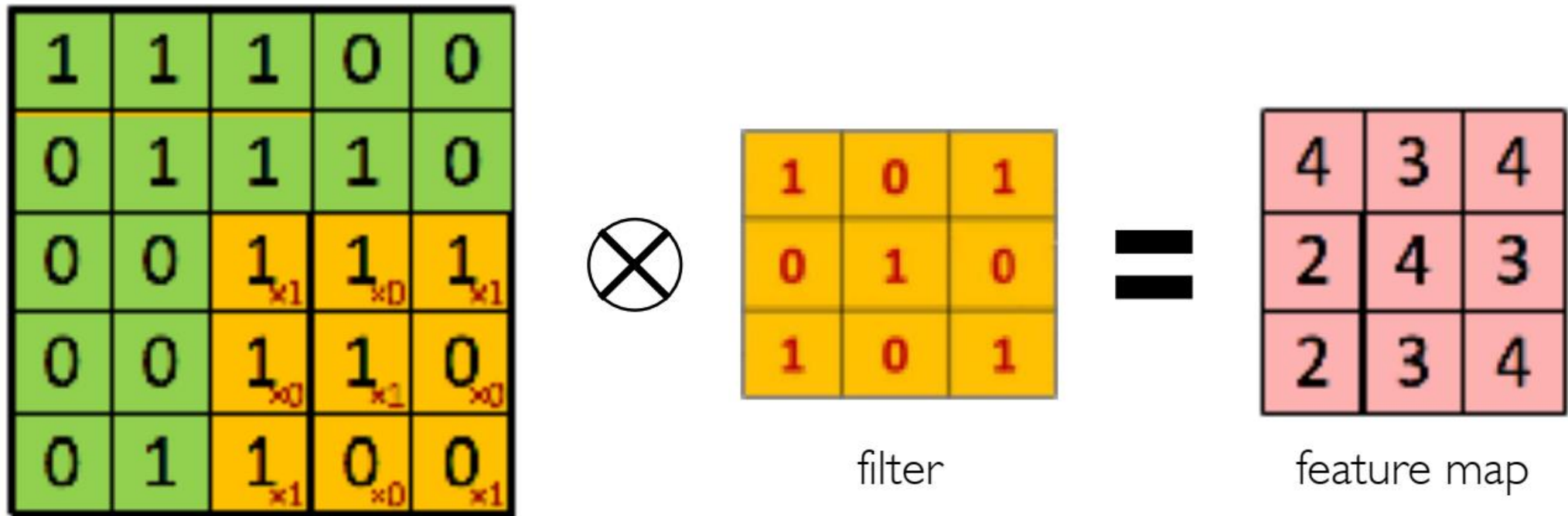
The convolution operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



The convolution operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



Purpose of Convolutions

- What does this filter do?

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0



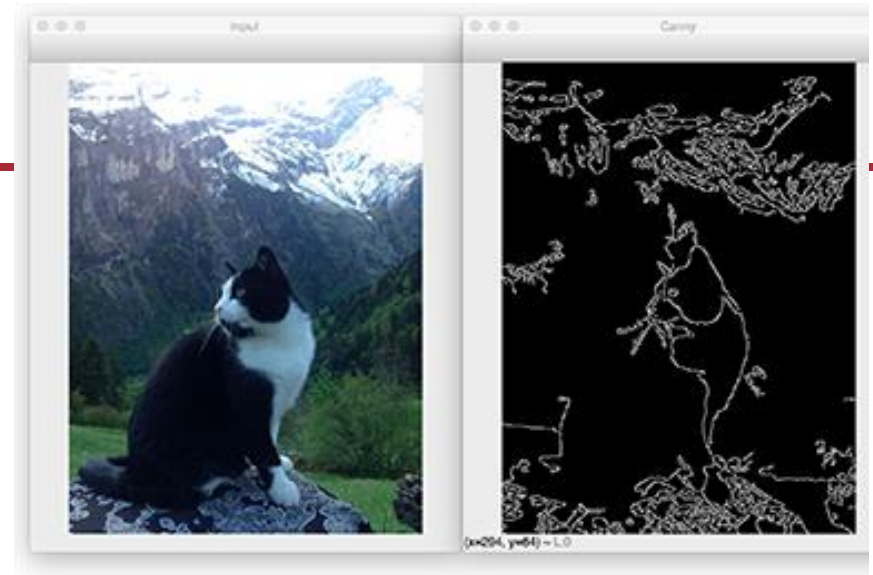
*

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



=



0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0



Purpose of Convolutions

- Different filters can be used to extract various features



Edge detection


$$* \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix} =$$


Kernel

The diagram illustrates the edge detection process. It shows an original image of a squirrel head being convolved with a 3x3 kernel. The kernel is a Laplacian filter, represented by the matrix: $\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$. The result is a grayscale image showing only the edges of the squirrel head. An arrow points from the word 'Kernel' to the matrix.

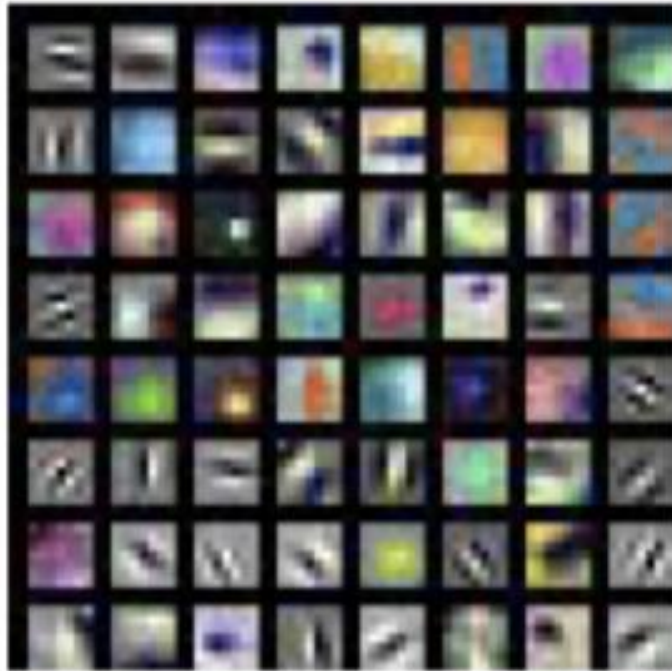
Sharpen


$$* \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix} =$$


The diagram illustrates the sharpening process. It shows the same original image of a squirrel head being convolved with a different 3x3 kernel. The kernel is a sharpening filter, represented by the matrix: $\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$. The result is a grayscale image where the edges are emphasized, making the squirrel head appear sharper.

Where does the “Deep Learning” part come in?

- Like dense fully-connected layers, we can just learn these filters!



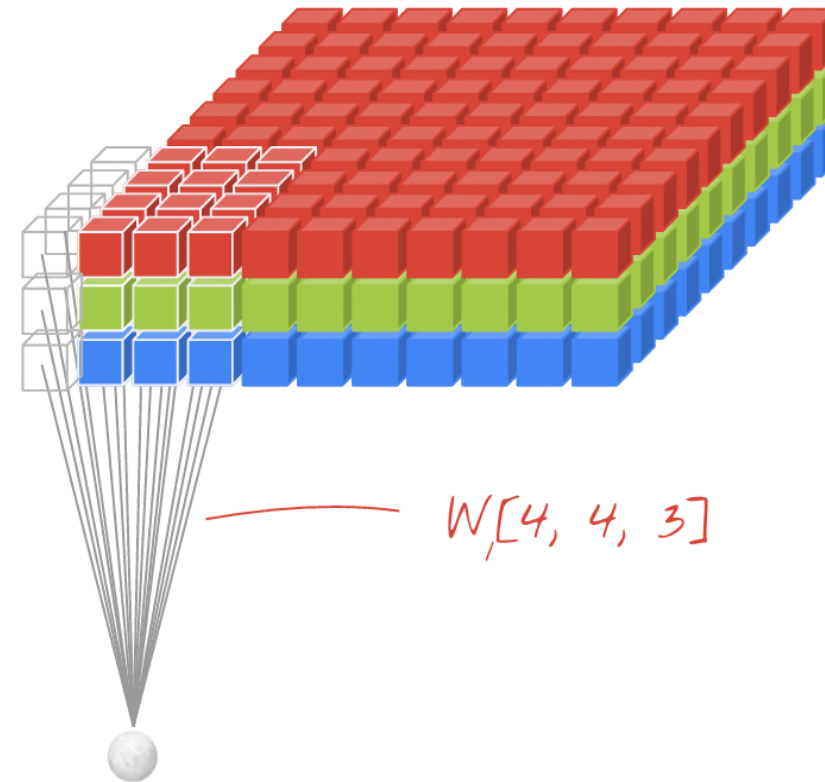
ResNet-18:
 $64 \times 3 \times 7 \times 7$



ResNet-101:
 $64 \times 3 \times 7 \times 7$

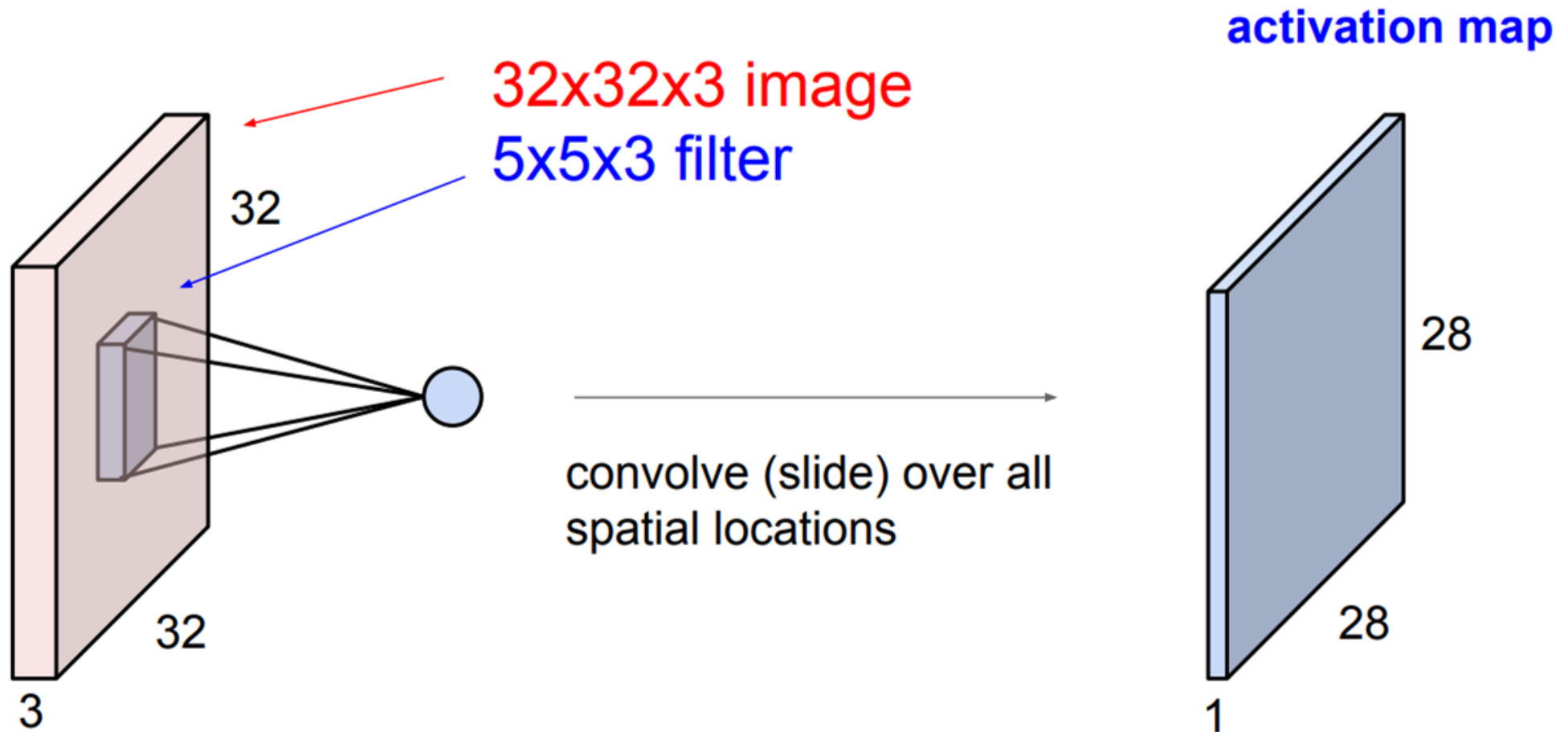
Filters (3D)

- Steps:
 - Compute the dot product for each channel (**same as 2D**)
 - Sum over each channel
- **Note:** The depth of the filter is always the same as the depth of the input image



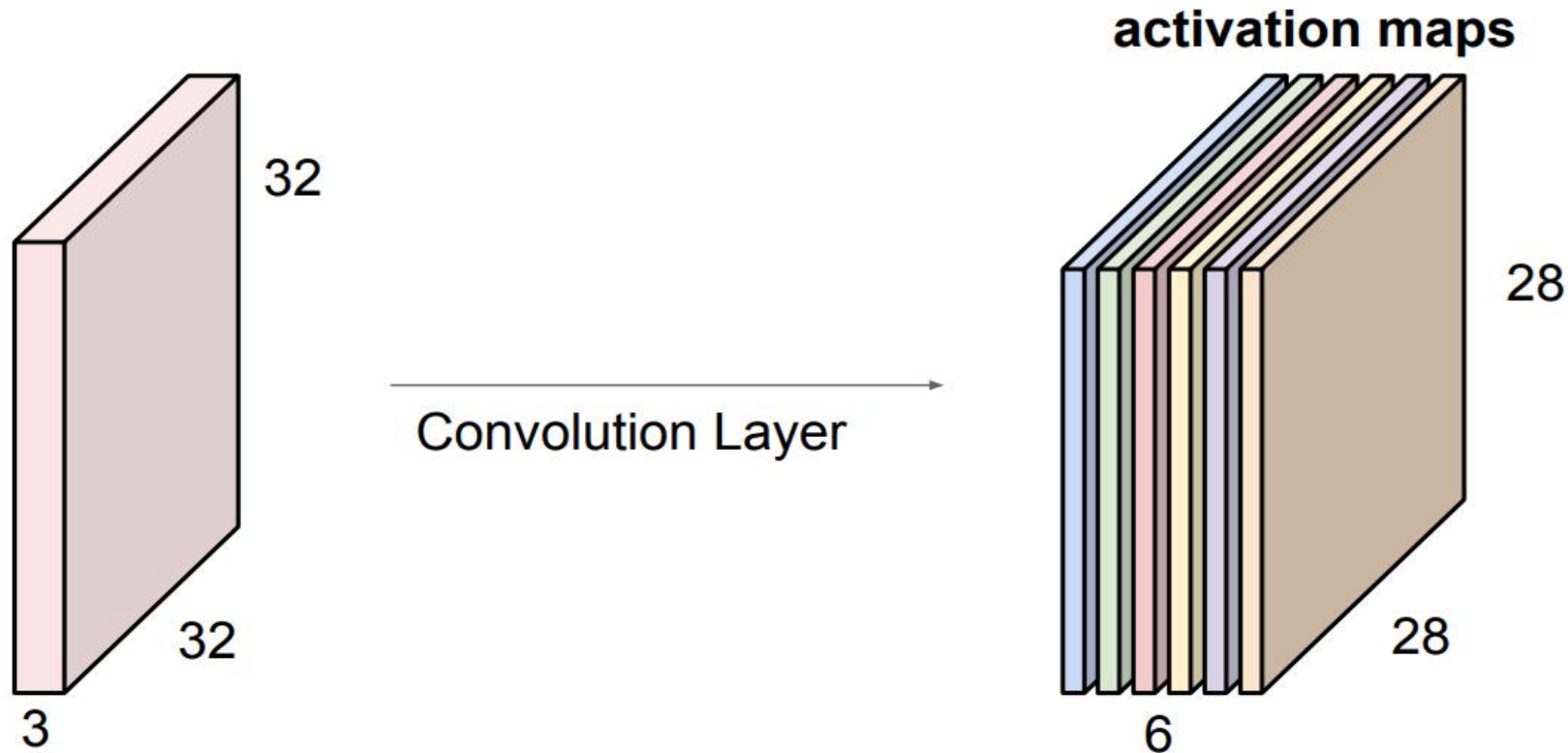
⚠ W_1 and W_2 are distinct $4 \times 4 \times 3$ filters

Convolutional Layers



Convolutional Layers

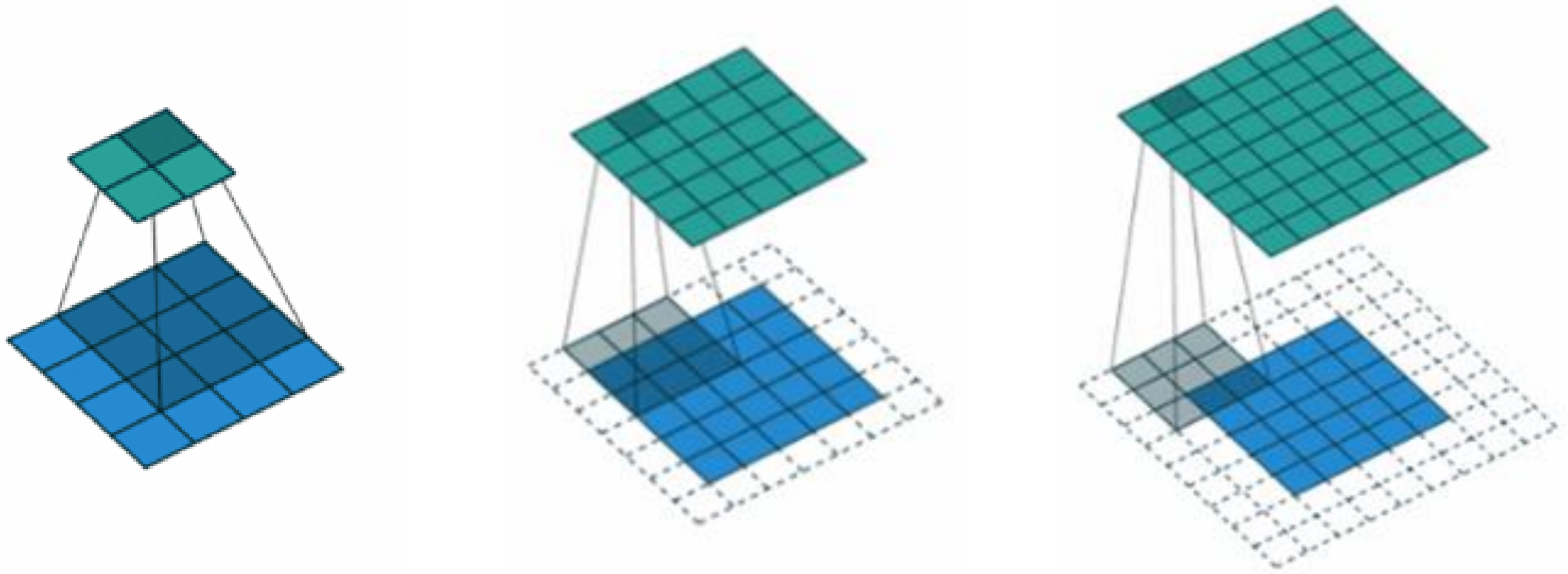
For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a “new image” of size 28x28x6!

Convolution Layers – Padding

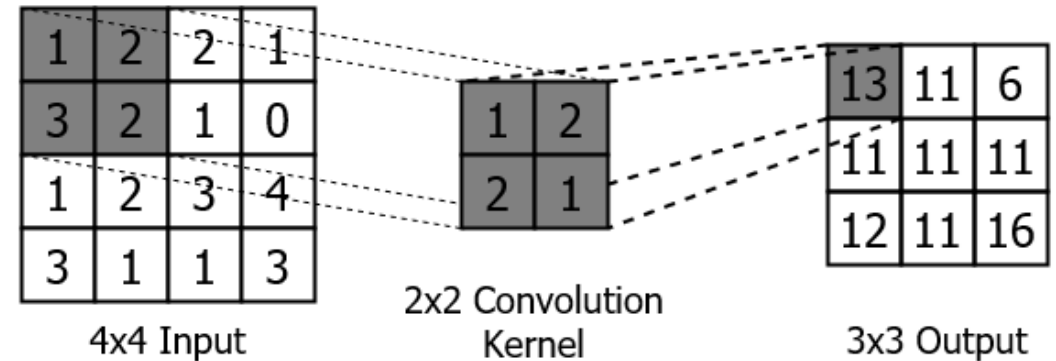
- Convolving an image with a filter results in a block with a smaller height and width – what if we want the height and width as before?



Convolution Layers – Same vs Valid Padding

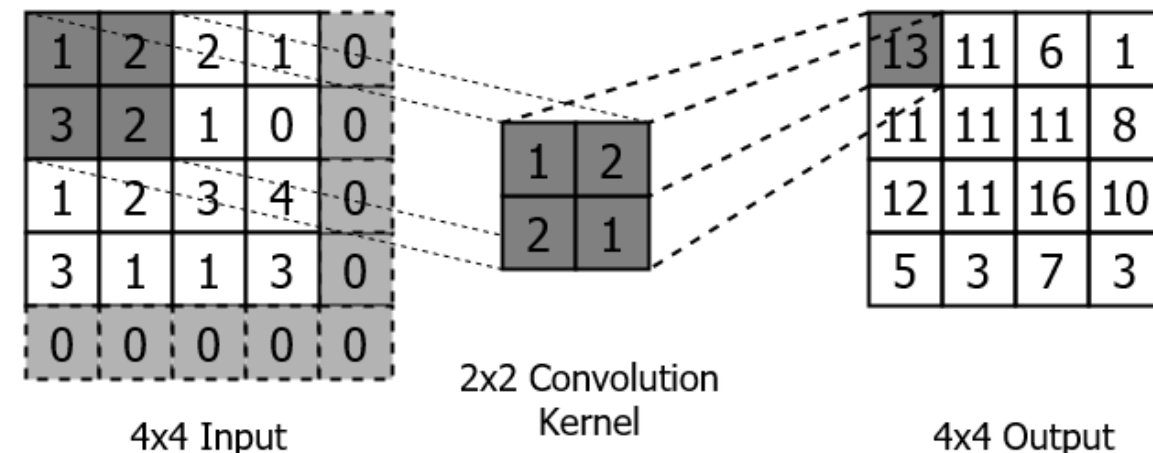
- Valid padding: no padding

Without Padding (Padding=Valid)



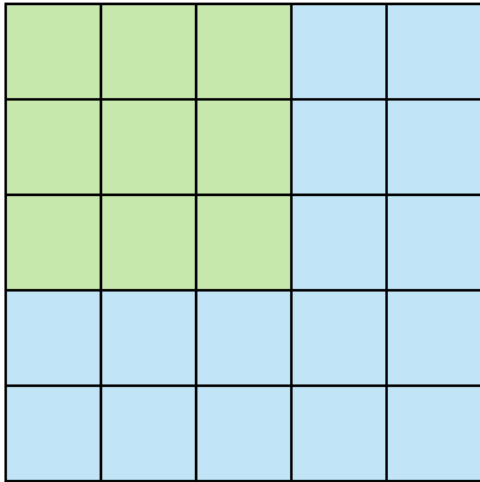
- Same padding: padding with 0s (or possibly some other constant value) to preserve the spatial dimensions of the output

With Padding (Padding=Same)

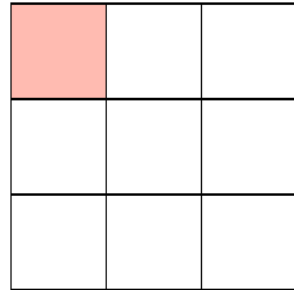


Convolution Layers – Stride

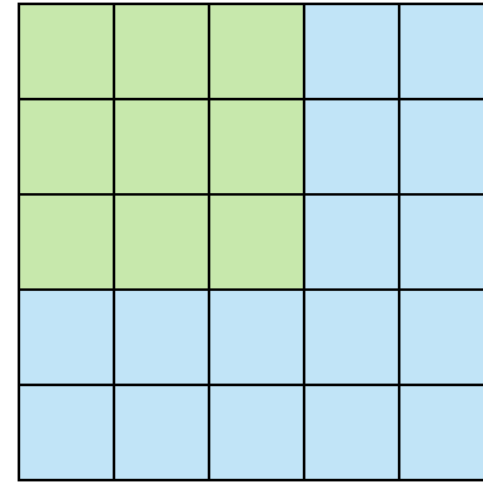
- The number of pixels to slide the filter by (both horizontally and vertically):
- A stride of 1 will shift the filter every pixel
- A stride of 2 will shift the filter every 2 pixels



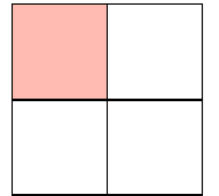
Stride 1



Feature Map



Stride 2



Feature Map

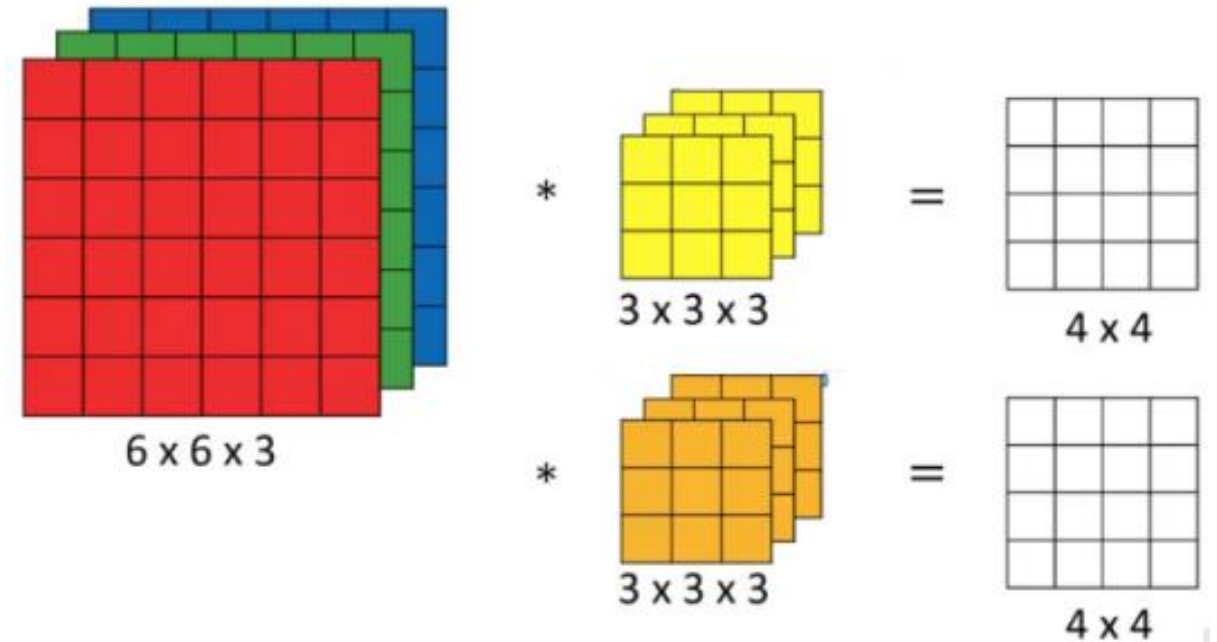
Output Dimensions

- $\text{ceil}[(W - F + 2P) / S] + 1$

- W: Input Dimension
- F: Kernel / Filter Size
- P: Padding size
- S: Strides

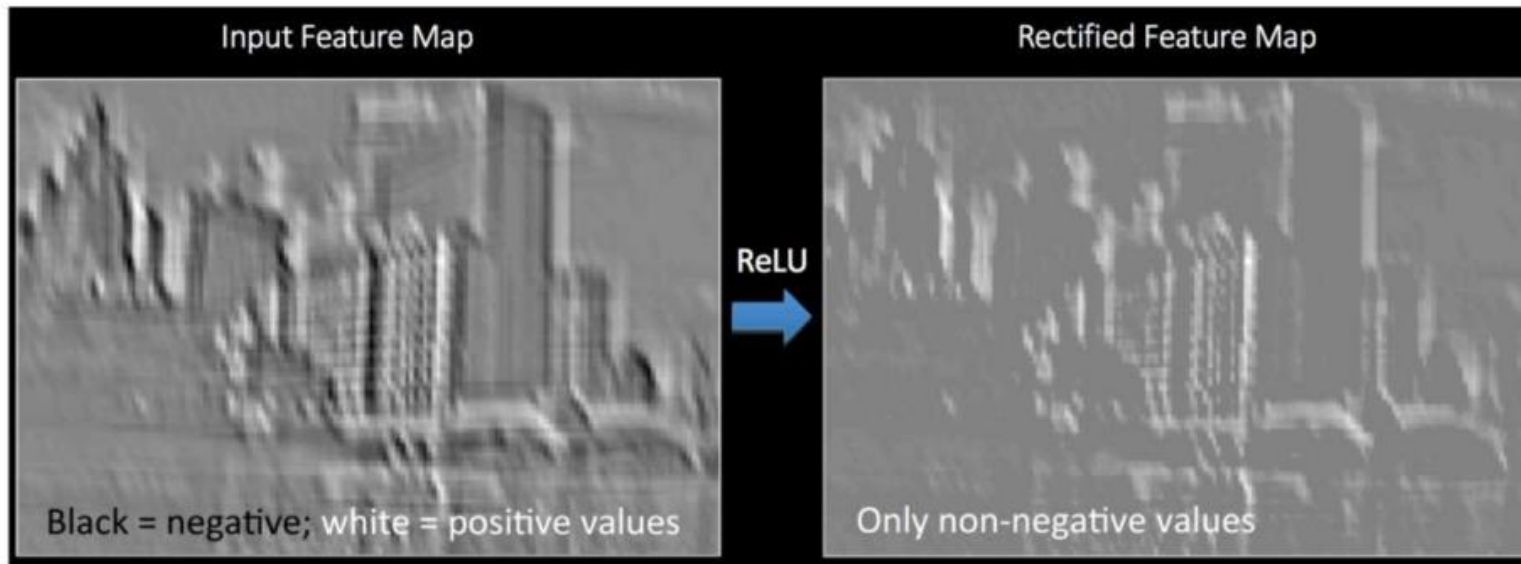
- For the figure on the right:

- Assume no padding and a stride of 1
- $W' = \text{ceil}[(6 - 3 + 2 * 0) / 1] + 1 = 4$
- $H' = \text{ceil}[(6 - 3 + 2 * 0) / 1] + 1 = 4$

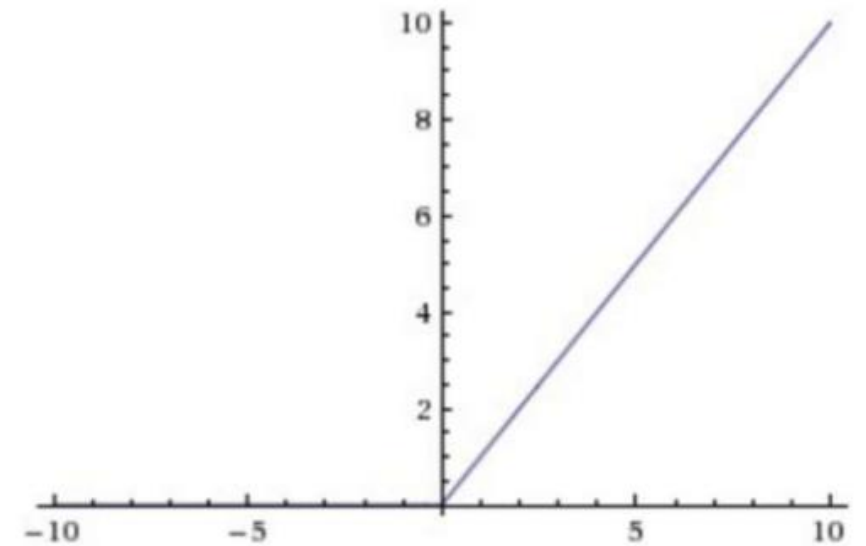


Introducing non-linearity

- Apply after every convolution operation (i.e., after convolutional layers)
- ReLU: pixel-by-pixel operation that replaces all negative values by zero. **Non-linear operation!**



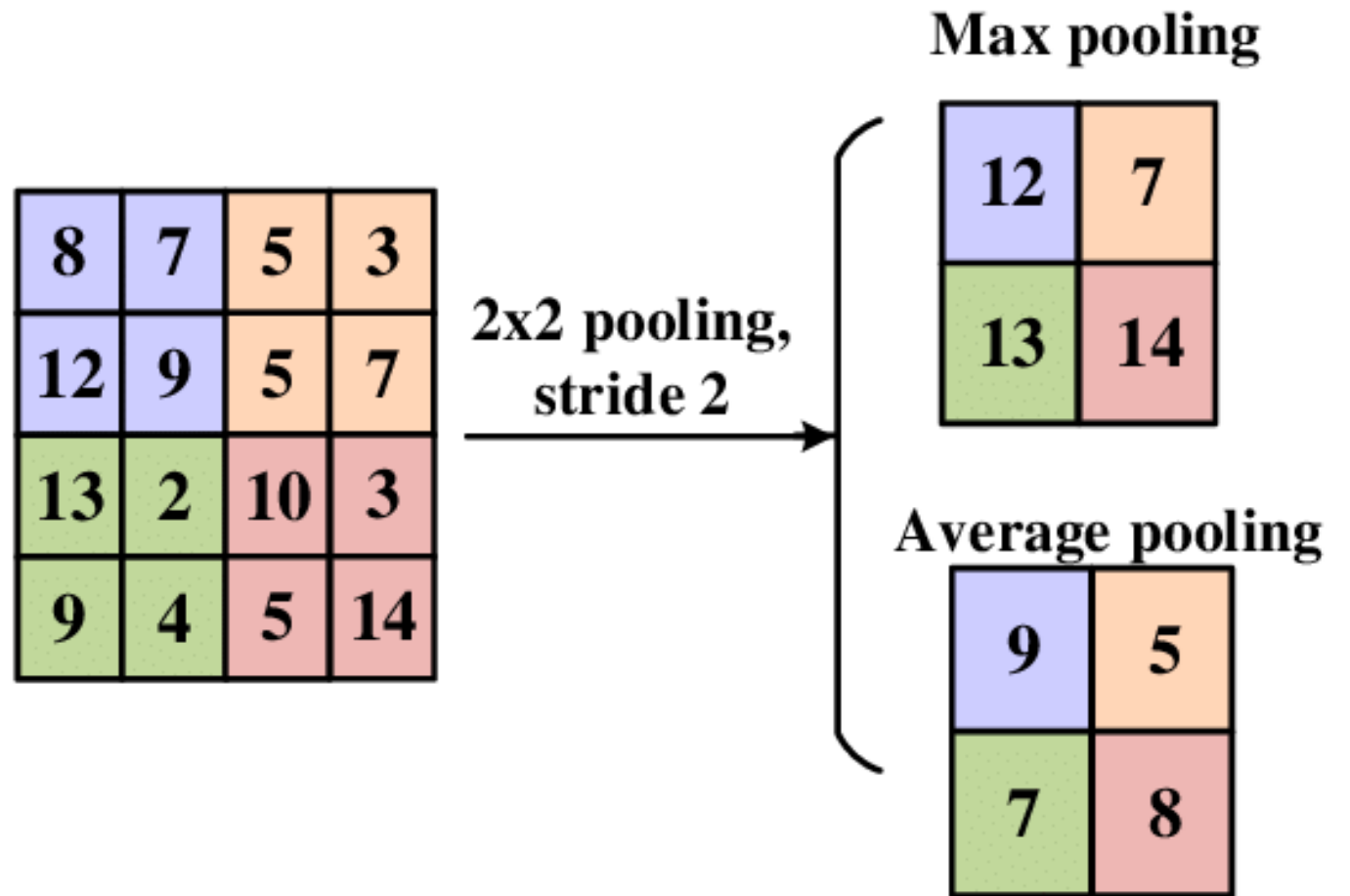
Rectified Linear Unit (ReLU)



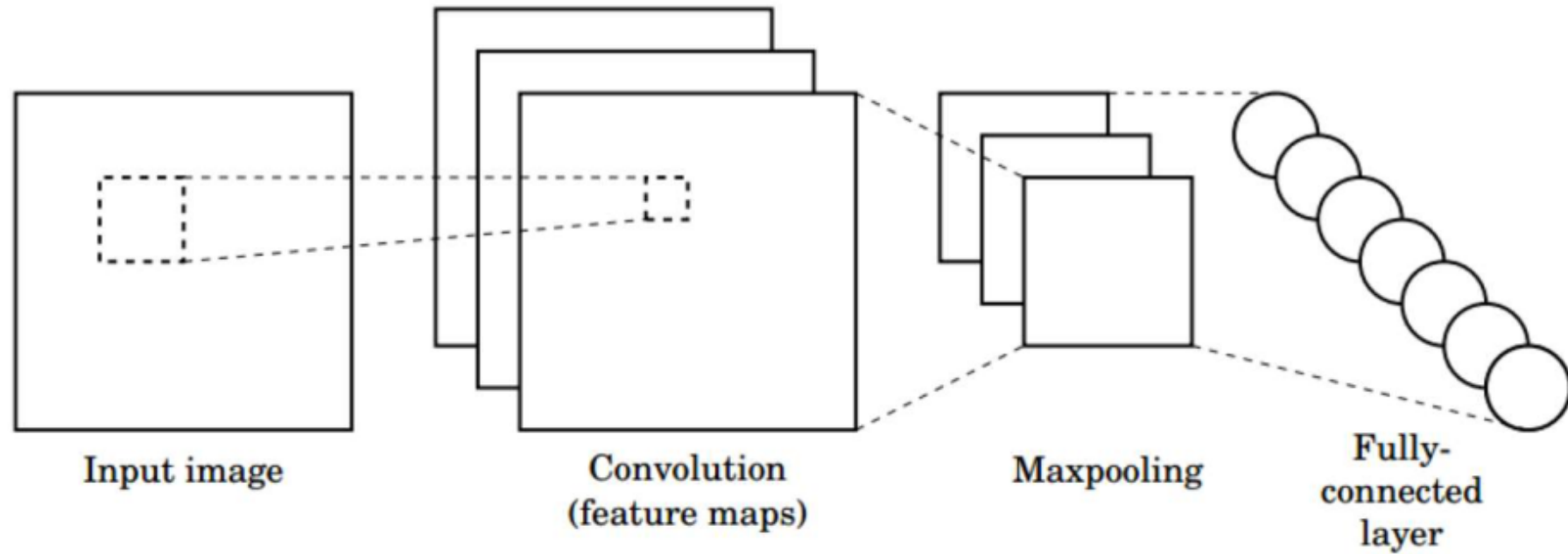
$$g(z) = \max(0, z)$$

Pooling Layers

- Reduces output size
- Applied to each channel independently
- Neighboring features may be similar
 - Doesn't remove too much information
- Max pooling takes the max
- Average pooling takes the average



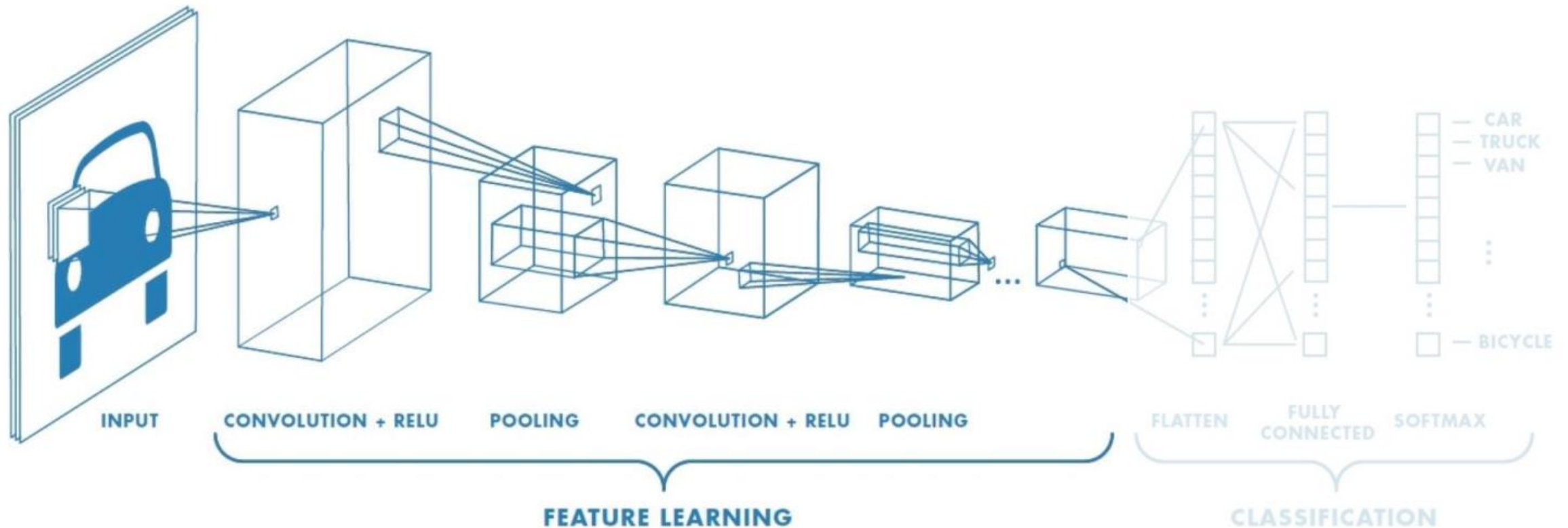
CNNs



1. **Convolution:** Apply filters with learned weights to generate feature maps.
2. **Non-linearity:** Often ReLU.
3. **Pooling:** Downsampling operation on each feature map.

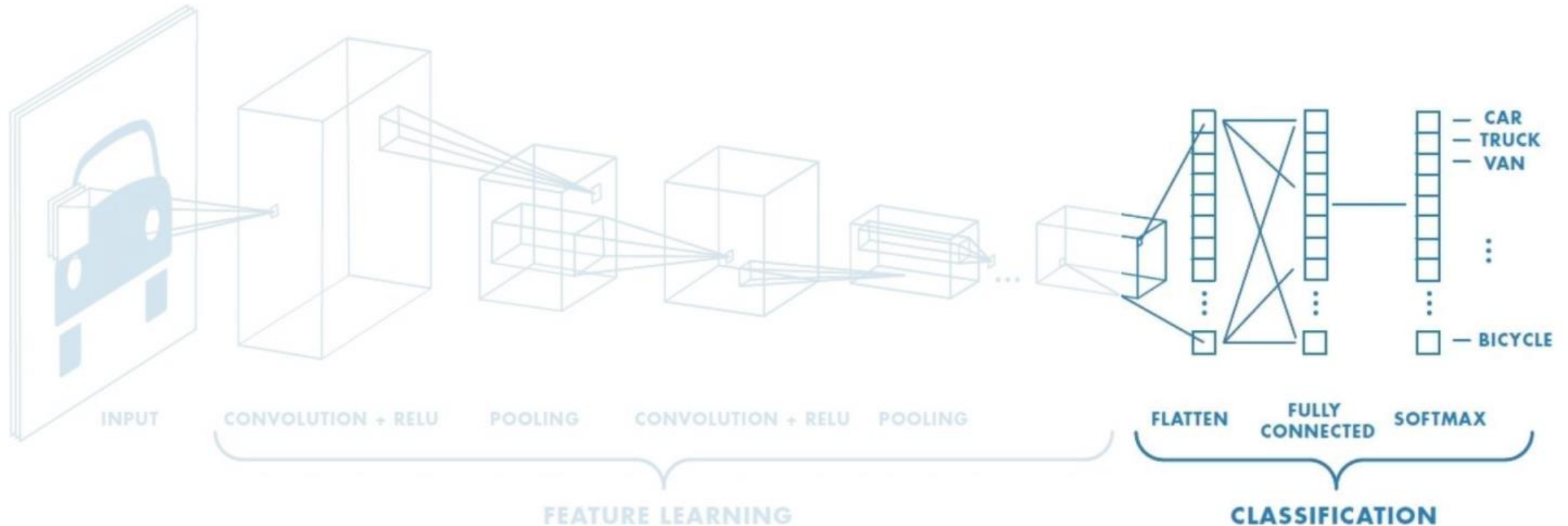
Train model with image data.
Learn weights of filters in convolutional layers.

CNNs for classification: *Feature learning*



1. Learn features in input image through **convolution**
2. Introduce **non-linearity** through activation function (real-world data is non-linear!)
3. Reduce dimensionality and preserve spatial invariance with **pooling**

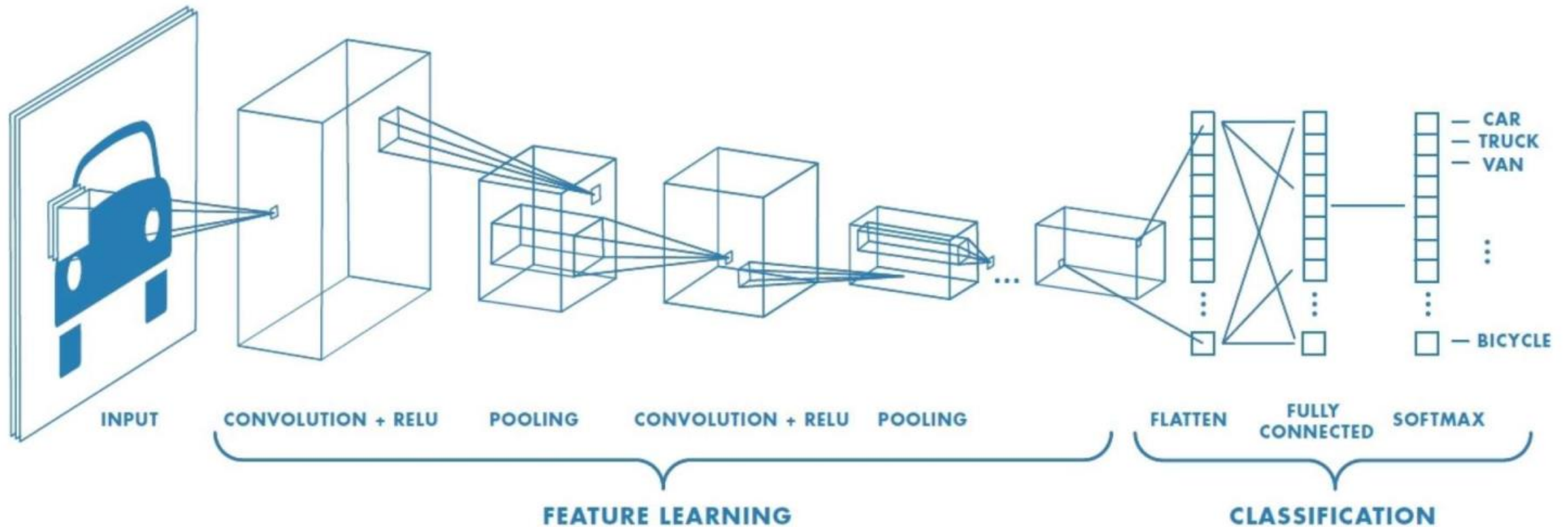
CNNs for classification: *Class probabilities*



- CONV and POOL layers output high-level features of input
- Fully connected layer uses these features for classifying input image
- Express output as **probability** of image belonging to a particular class

$$\text{softmax}(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

CNNs: Training with backpropagation

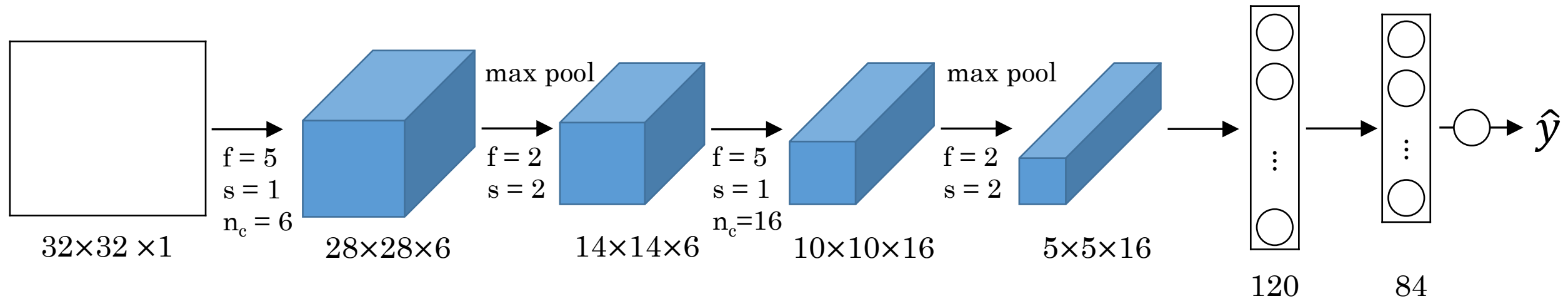


Learn weights for convolutional filters and fully connected layers

Backpropagation: cross-entropy loss

$$J(\theta) = \sum_i y^{(i)} \log(\hat{y}^{(i)})$$

CNN Example



CNN Example

	Activation Shape	Activation Size	# of parameters
Input	(32,32,3)	3072	
CONV 1 (f=5,s=1,n _c = 8)	(28,28,8)	6272	$(5*5*3 + 1) * 8 = 608$
POOL 1	(14,14,8)	1568	
CONV 2 (f=5,s=1, n _c = 16)	(10,10,16)	1600	$(5*5*8 + 1) * 16 = 3216$
POOL 2	(5,5,16)	400	
FC 3	(120,1)	120	$400 * 120 + 120 = 48120$
FC 4	(84,1)	84	$120 * 84 + 84 = 10164$
Softmax	(10,1)	10	$84*10 + 10 = 850$

Thanks
