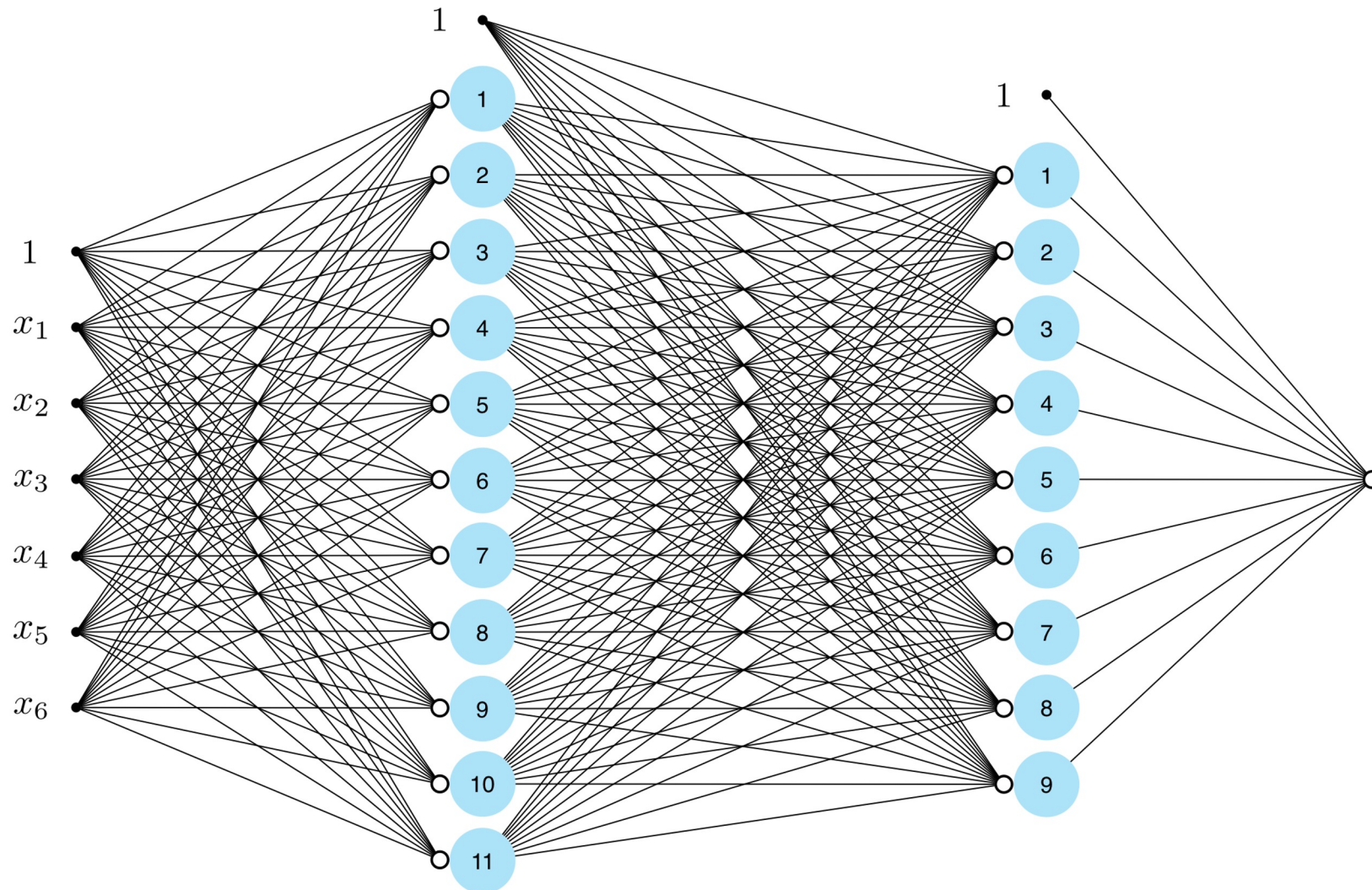


Convolutional Neural Networks

Tony T. Luo
July 2023

Multi-Layer Perceptron: fully connected neural network

A two-layer MLP network



Python implementation (using Numpy)

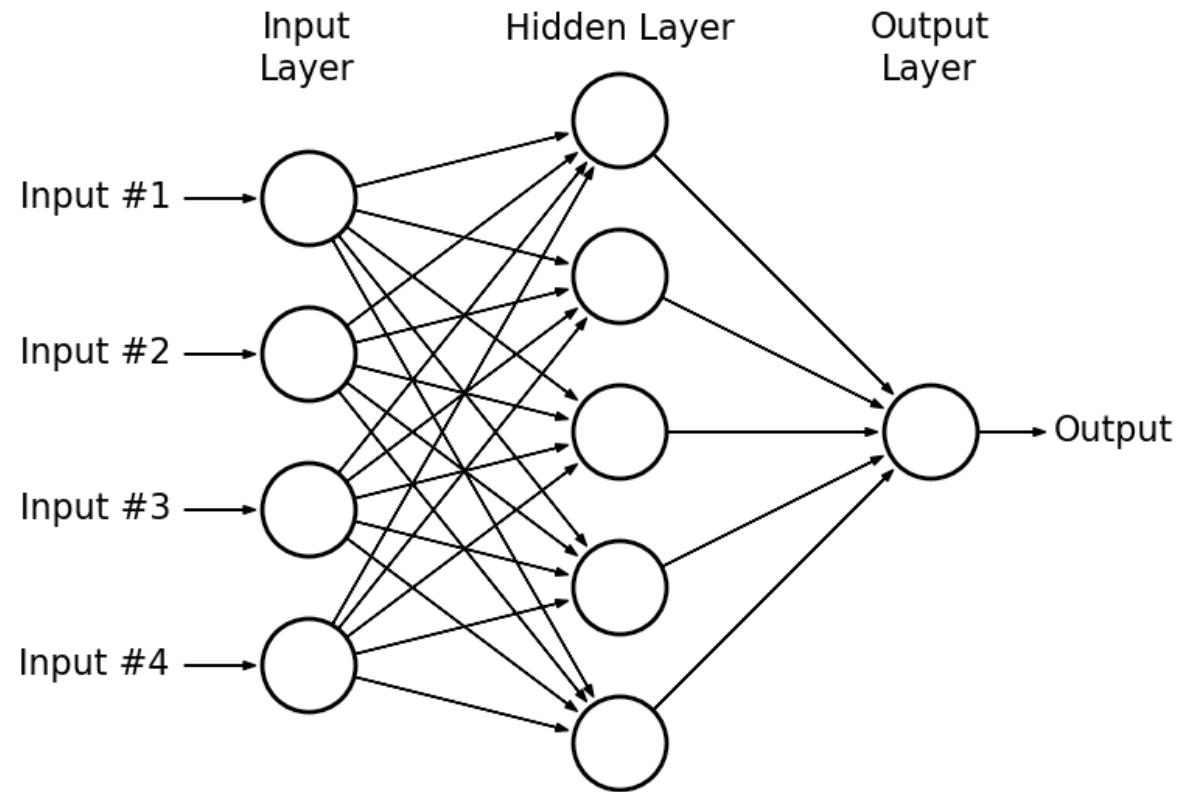
```
1  # neural network feature transformation
2  def feature_transforms(a, w):
3
4      # loop through each layer
5      for W in w:
6
7          # compute inner-product with current layer weights
8          a = W[0] + np.dot(a.T, W[1:])
9
10         # pass through activation
11         a = activation(a).T
12
13     return a
```

- Notice bias and activation function

Convolutional Neural Networks

Why CNN (why not MLP)?

Problems with MLP (all FC layers)

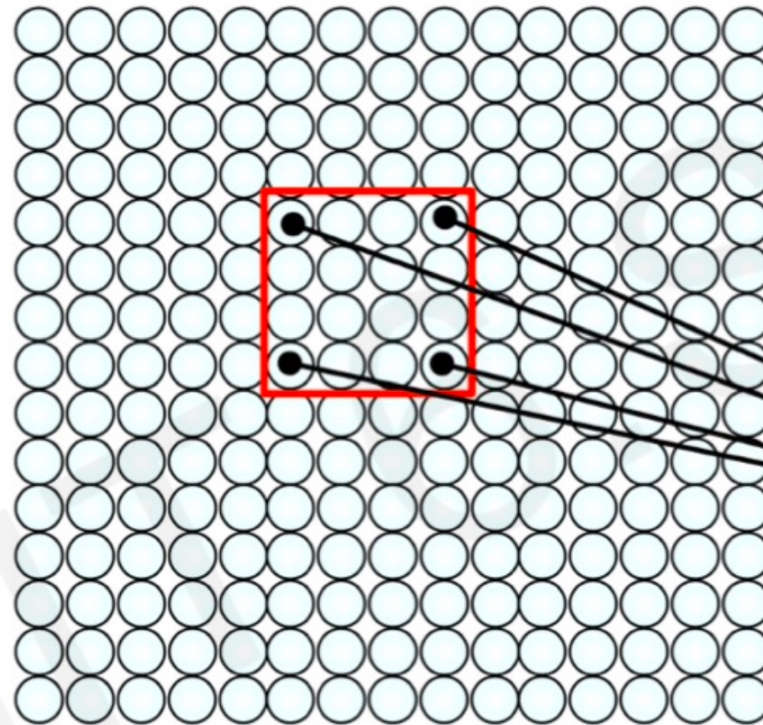


- Too many parameters
 - $U_1 * U_2 * U_3 \dots$
- Does not preserve spatial information
 - Every neuron sees the same inputs!

Preserving spatial structure

Each neuron only sees a “patch”

Input: 2D image.
Array of pixel values

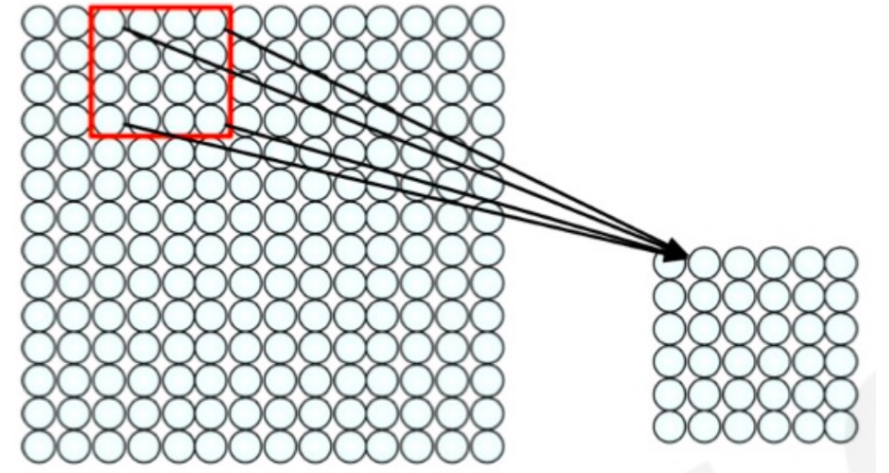


Idea: connect patches of input
to neurons in hidden layer.

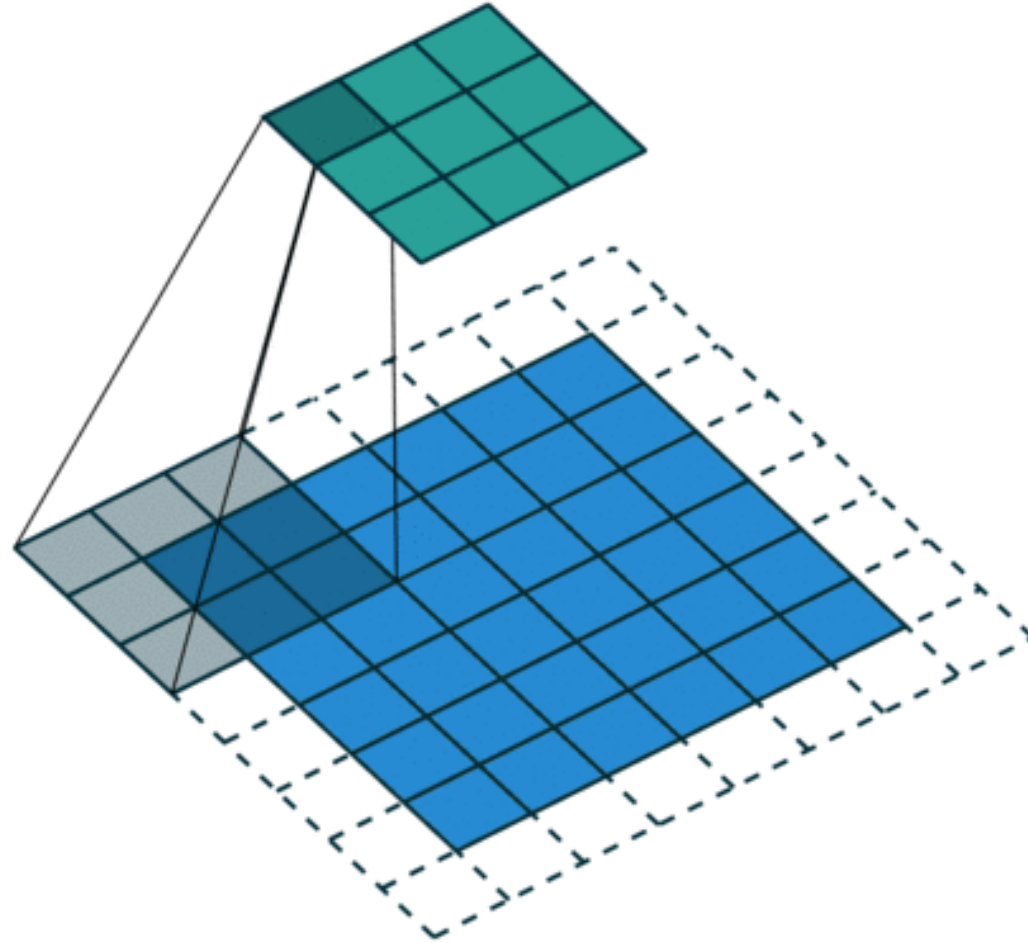
Neuron connected to region of
input. Only “sees” these values.

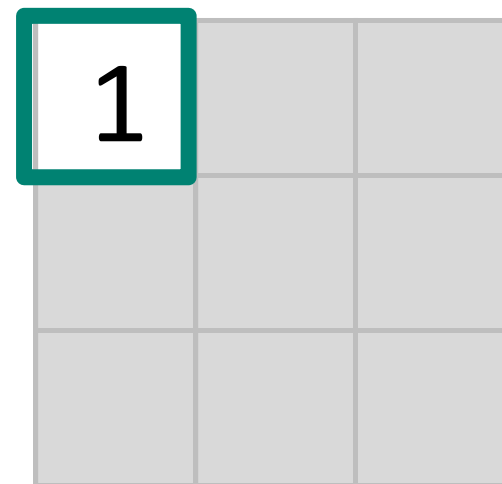
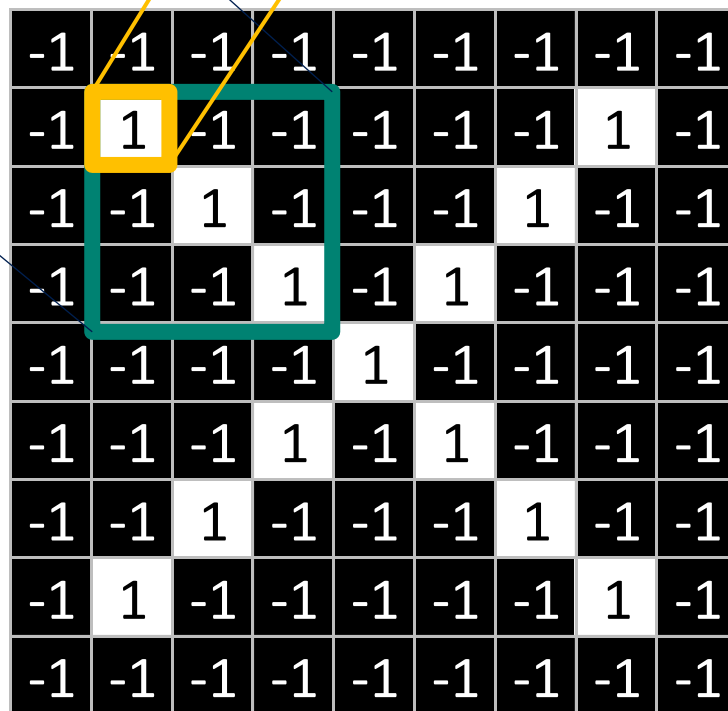
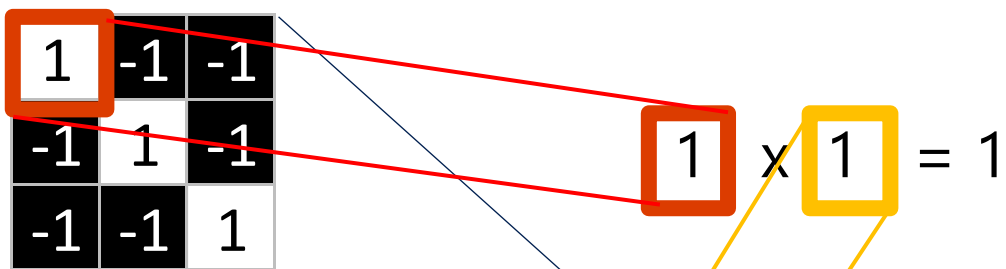
How to map a patch to a single neuron?

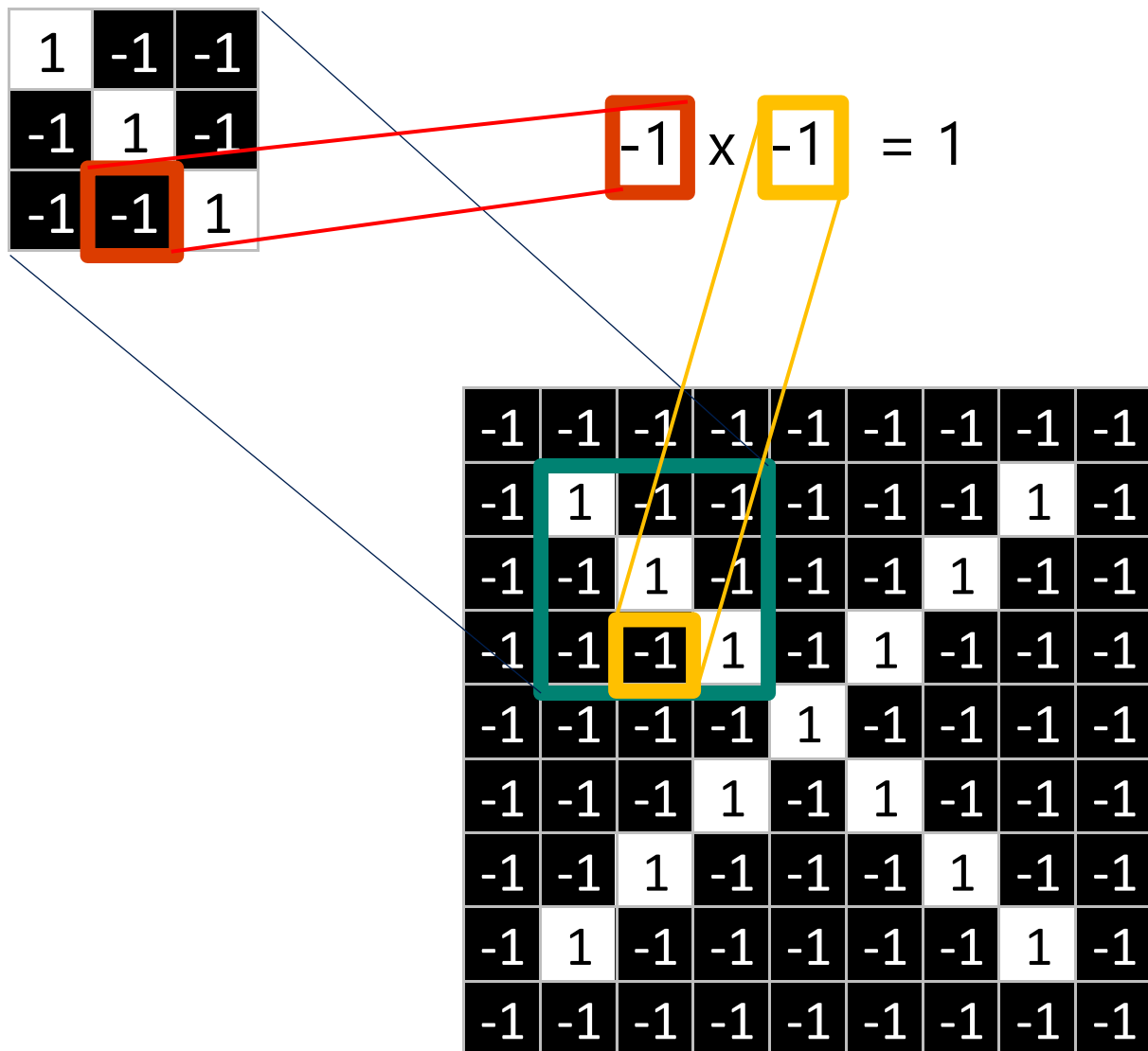
- Basic idea: weighted sum
- Weights are called a “**filter**” or “**kernel**”
- $\text{patch} \cdot \text{filter} = \text{output hidden neuron}$
 - Element-wise multiply
- This operation is called “**convolution**”
- Each filter corresponds to a certain “**feature**” of input
- Use multiple filters to extract multiple features



Convolution Operation



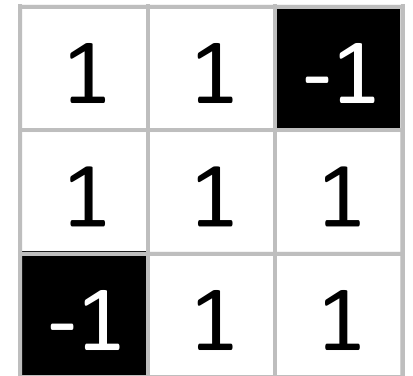
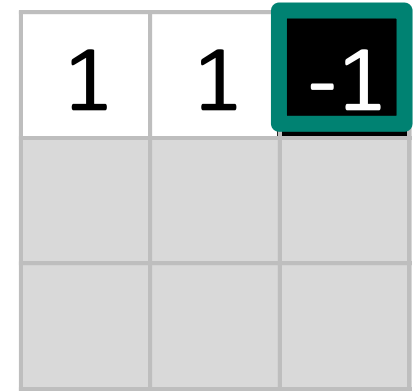
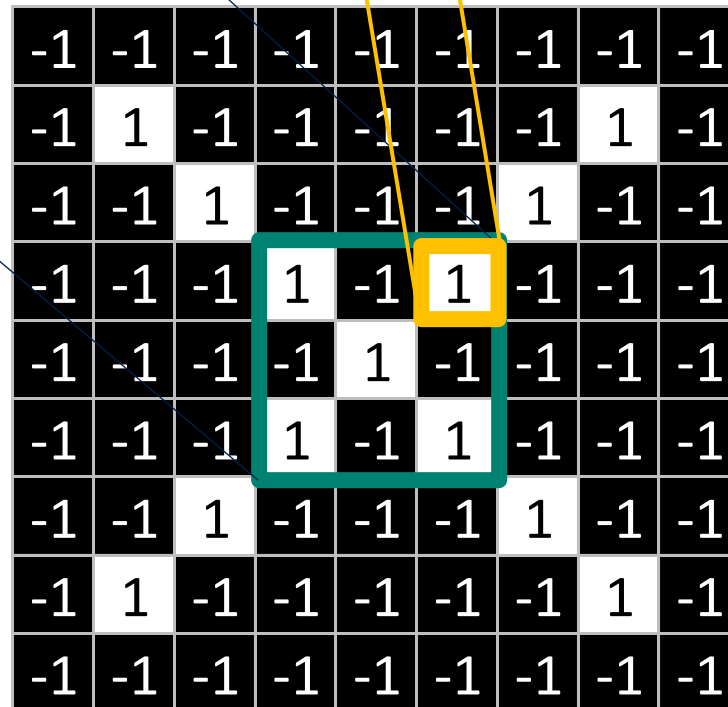
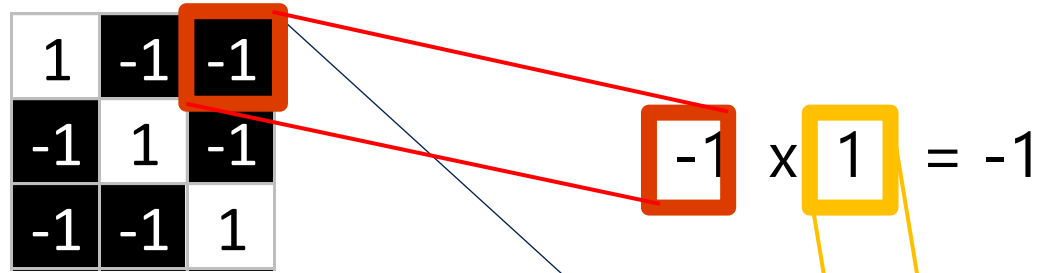




1	1	1
1	1	1
1	1	

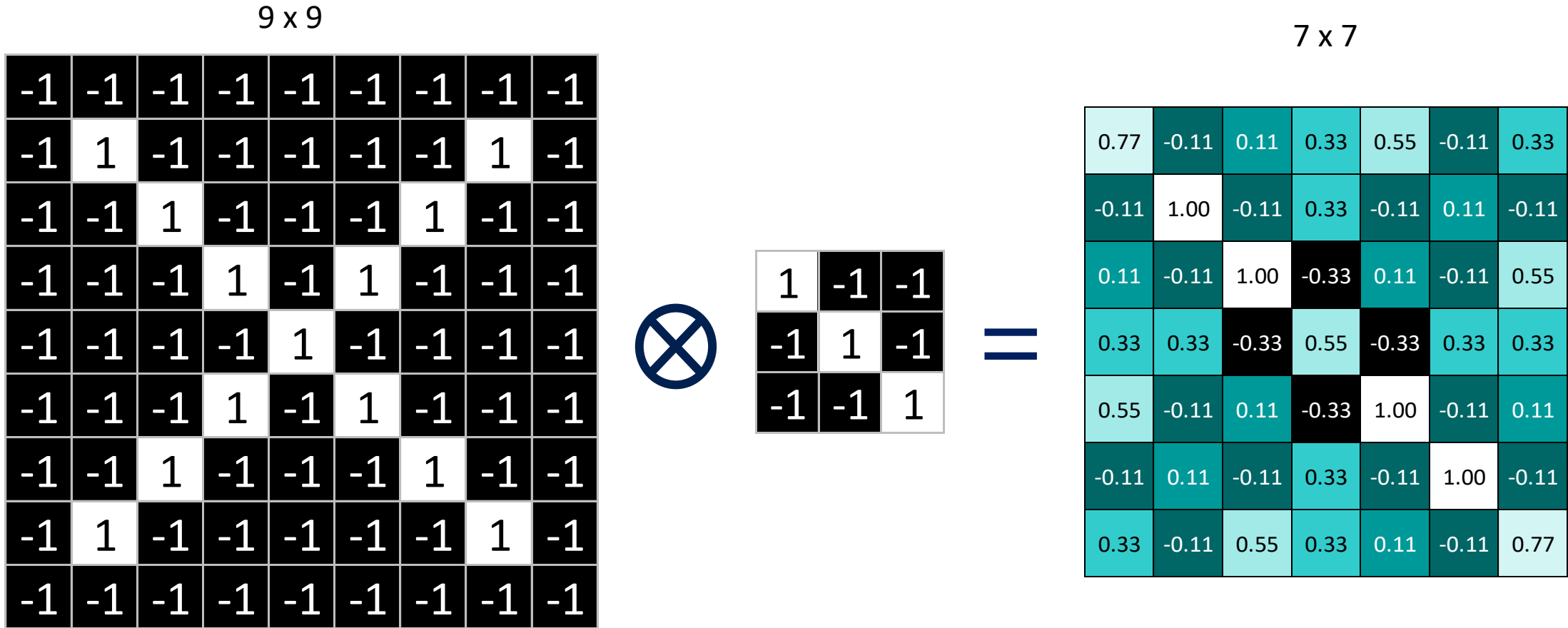
Sum: 9 \rightarrow 9/9=1
(single neuron)

Move on to another patch



Sum: 5 \rightarrow $5/9 = 0.55$

After convolution over every patch...



Feature map

Assuming stride (step size) = 1

Each channel has a filter

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

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

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

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



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

=

0.77	-0.11	0.11	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77

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



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

=

0.33	-0.55	0.11	-0.11	0.11	-0.55	0.33
-0.55	0.55	-0.55	0.33	-0.55	0.55	-0.55
0.11	-0.55	0.55	-0.77	0.55	-0.55	0.11
-0.11	0.33	-0.77	1.00	-0.77	0.33	-0.11
0.11	-0.55	0.55	-0.77	0.55	-0.55	0.11
-0.55	0.55	-0.55	0.33	-0.55	0.55	-0.55
0.33	-0.55	0.11	-0.11	0.11	-0.55	0.33

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



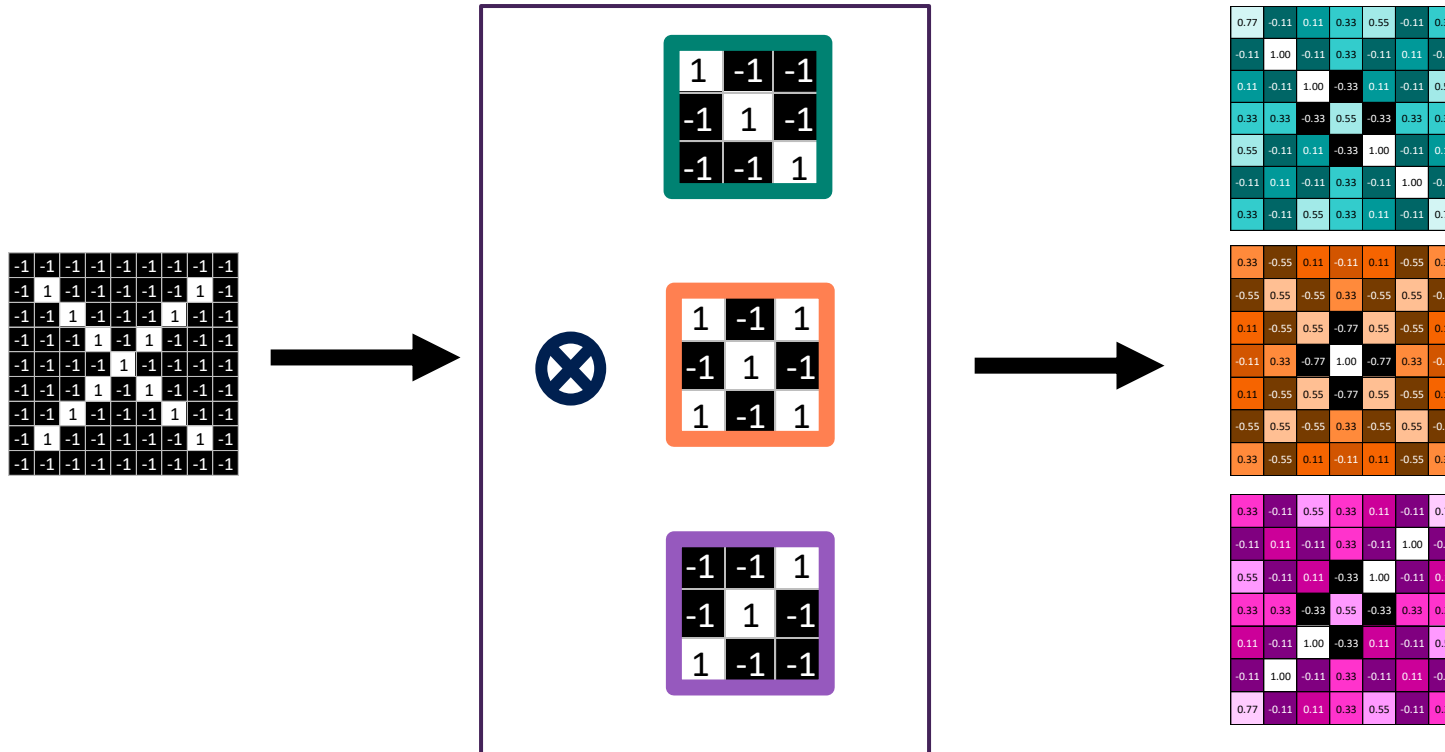
-1	-1	1
-1	1	-1
1	-1	-1

=

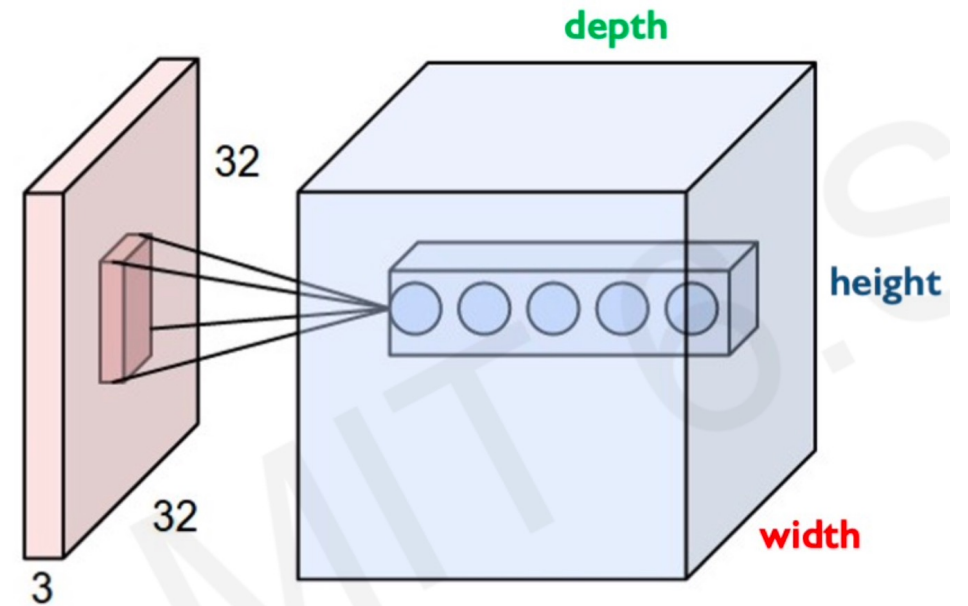
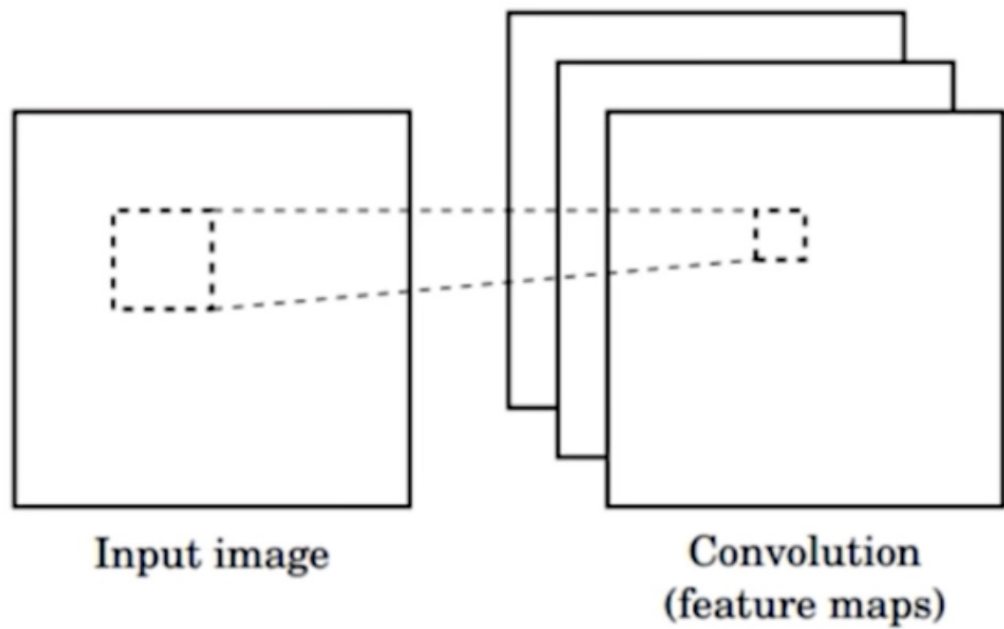
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.77	-0.11	0.11	0.33	0.55	-0.11	0.33

Output of a convolution layer

- A **stack** of feature maps

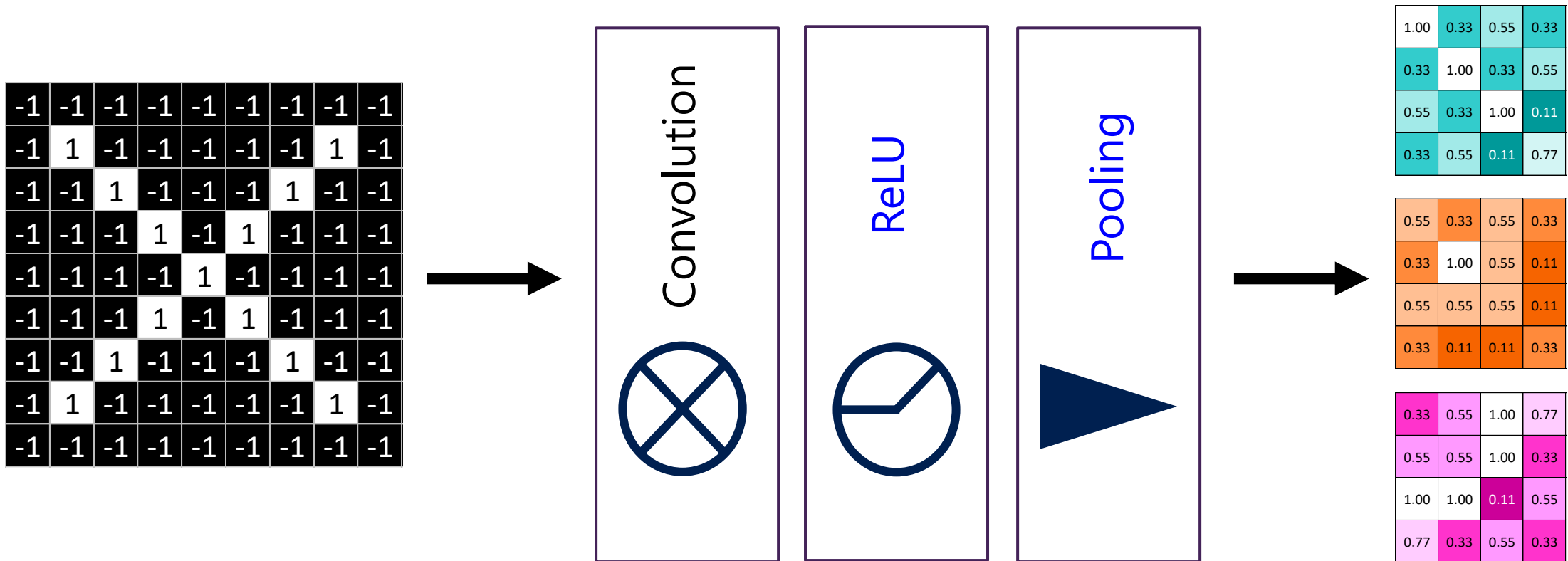


- Stacked together (depth increases)

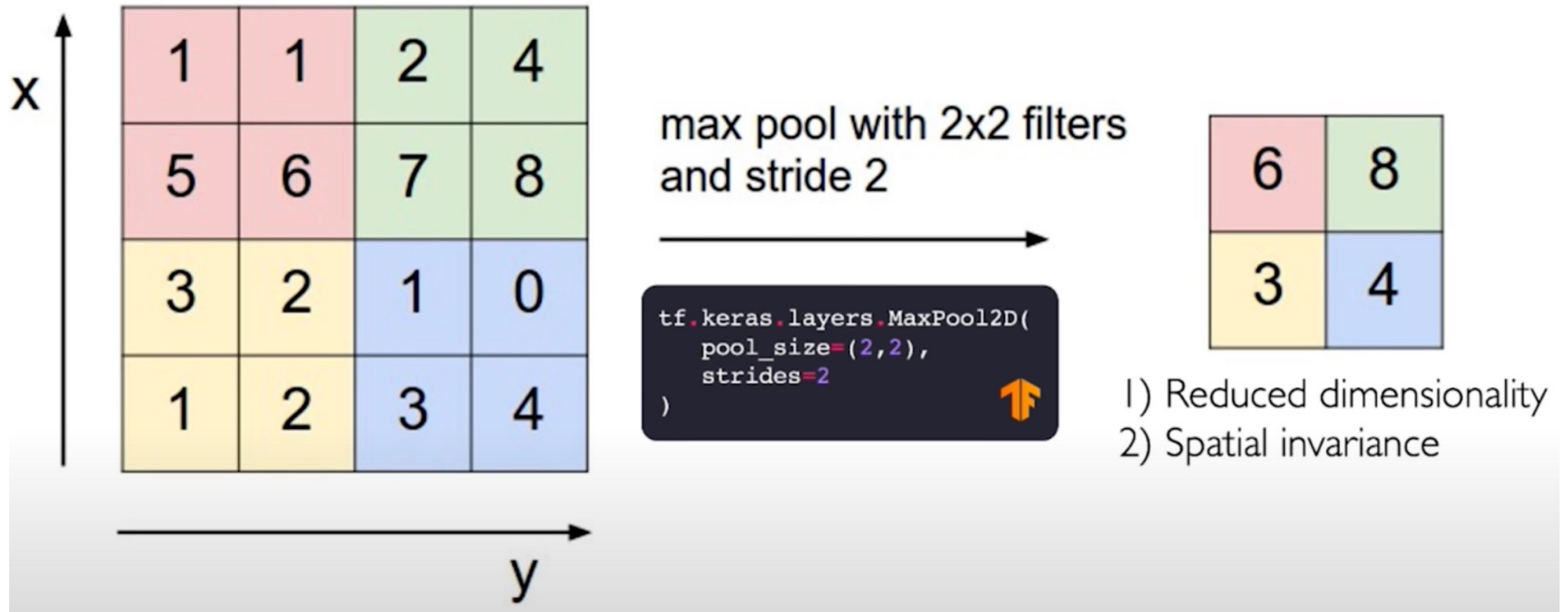


Pipeline

- Add *activation function* to each convolution
 - *Sigmoid, Relu, Tanh, etc.*
- Add *pooling* to downsample feature map

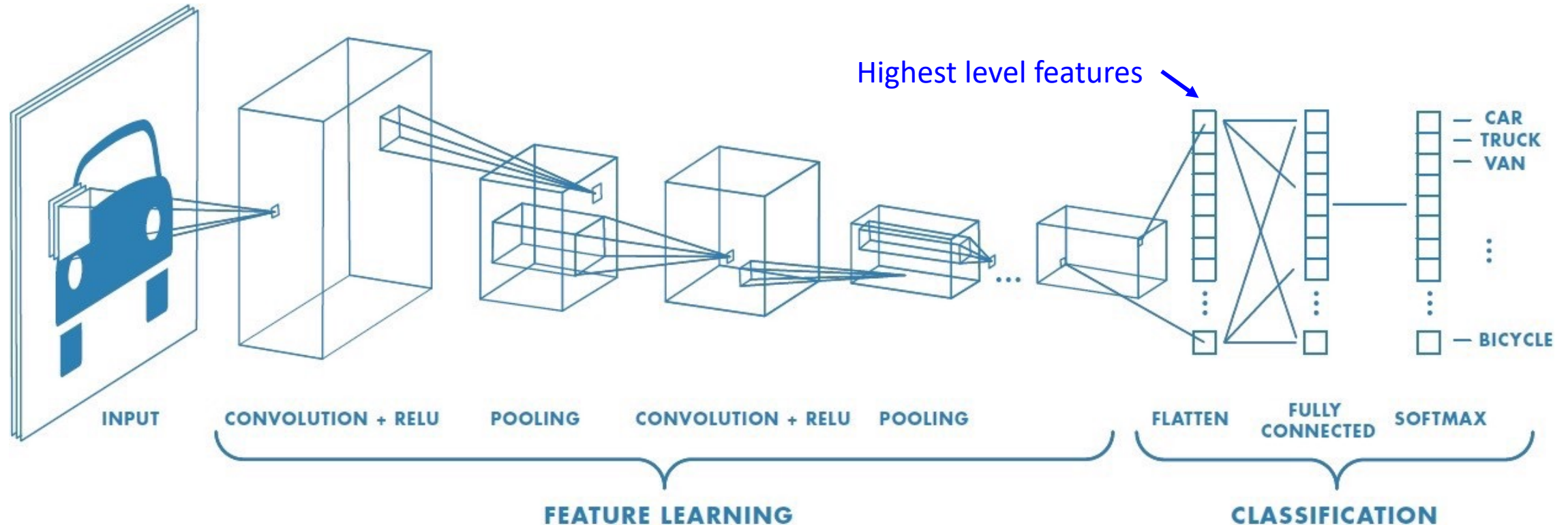


Pooling (max pooling, avg pooling)



Feature map

CNN: concatenates many layers together (feature extraction) + FC (classification/regression)



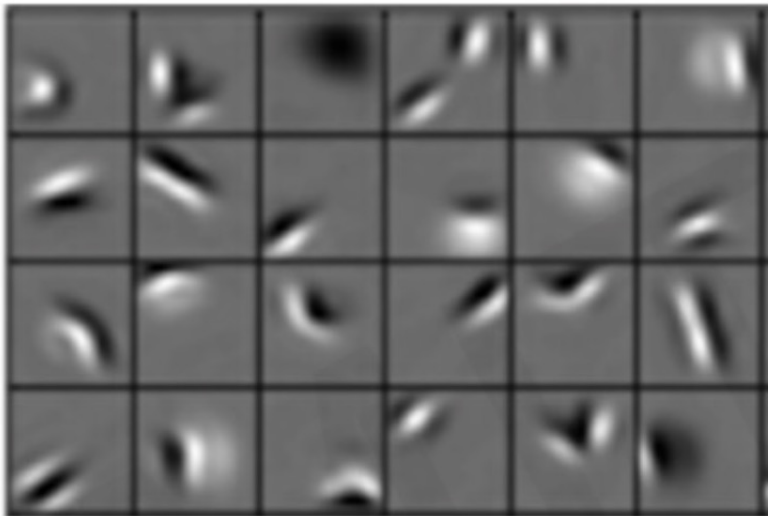
- Convolution + activation + pooling
- Convolution increases depth by <# of filters> folds
- Output of FC: scores of the input belonging to each class

SOFTMAX FUNCTION

$$s(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

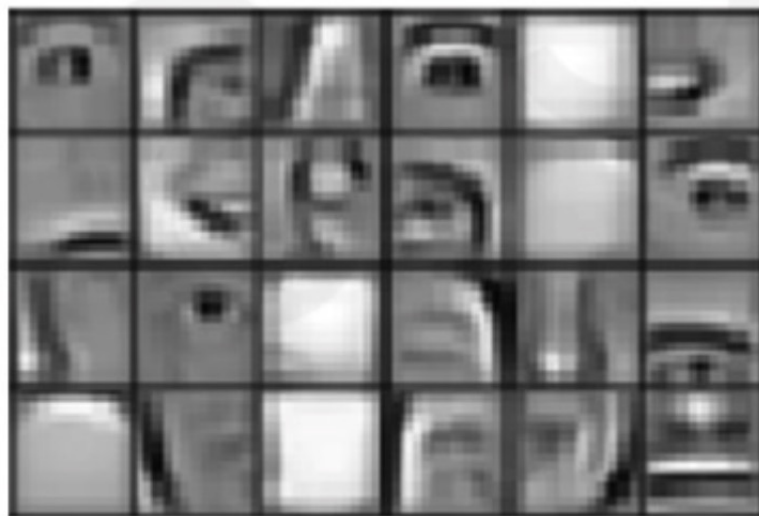
Intuition of “features”

Low level features



Edges, dark spots

Mid level features



Eyes, ears, nose

High level features

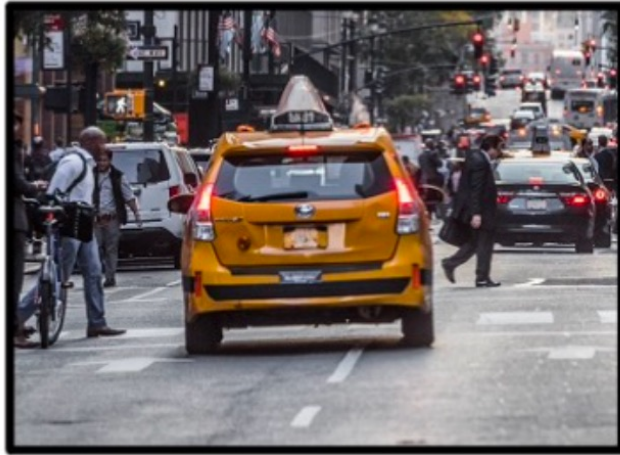


Facial structure

CNN is Versatile

- Image classification
- Object detection
- Segmentation
- ... (probabilistic control / robotics)

Object detection



Output:
taxi: (x_1, y_1, w_1, h_1)

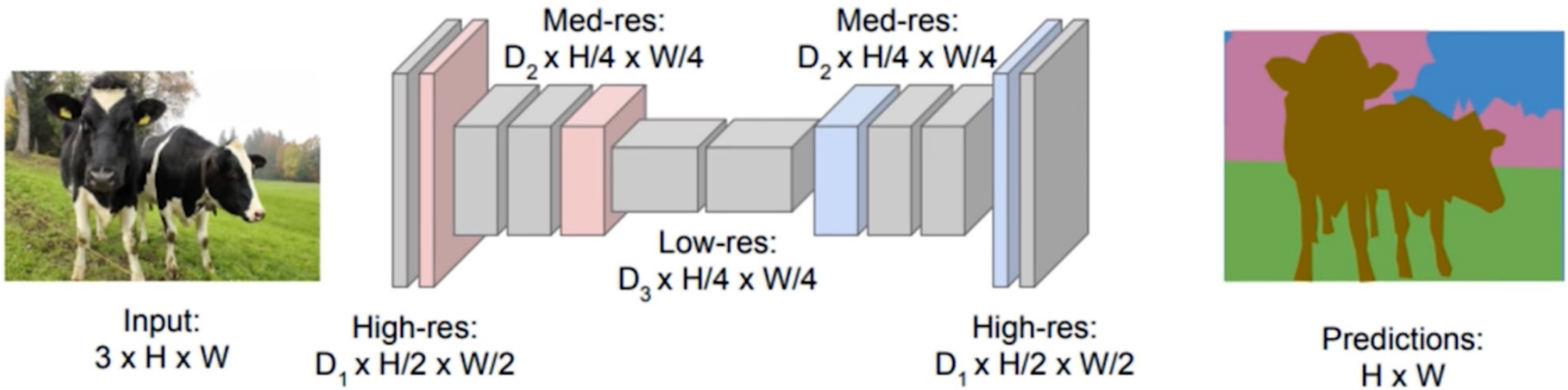
Multi-OD



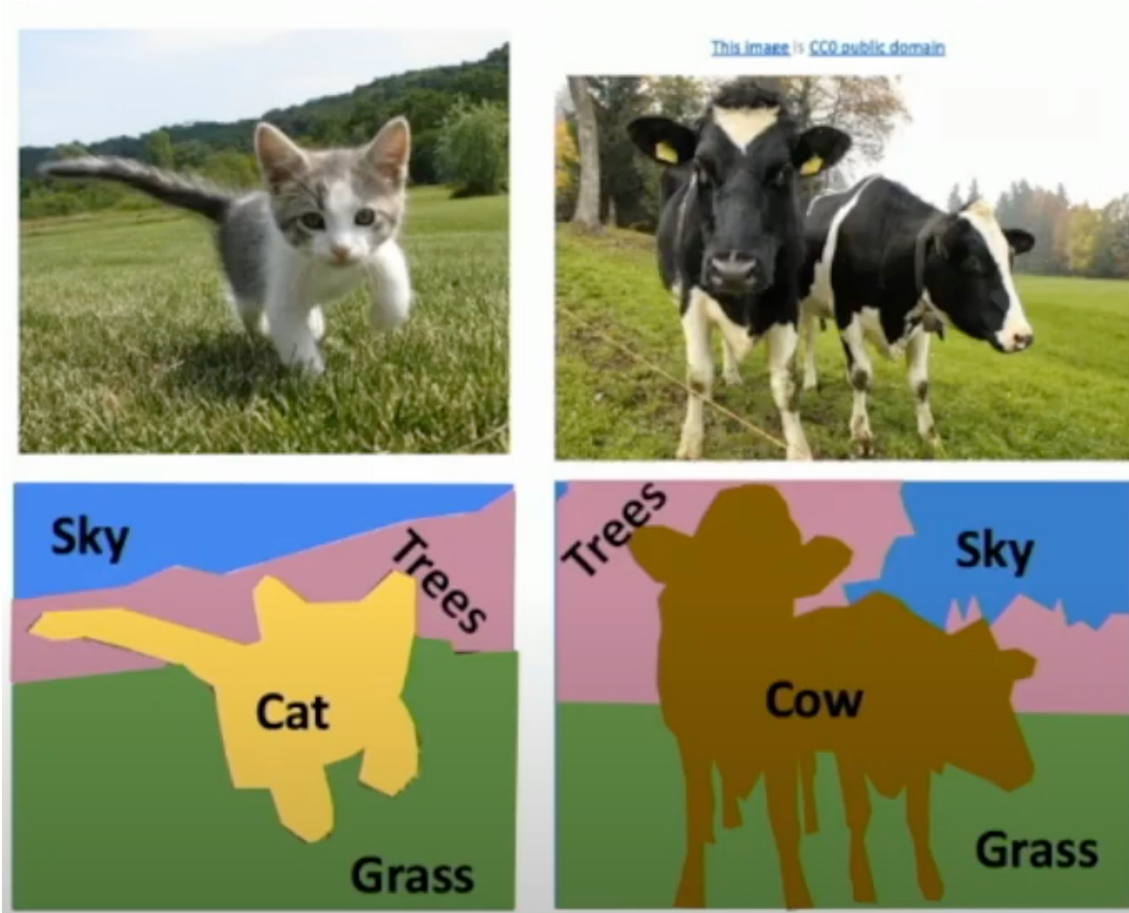
Output:
taxi: (x_1, y_1, w_1, h_1)
person: (x_2, y_2, w_2, h_2)
person: (x_3, y_3, w_3, h_3)
....

Semantic segmentation:

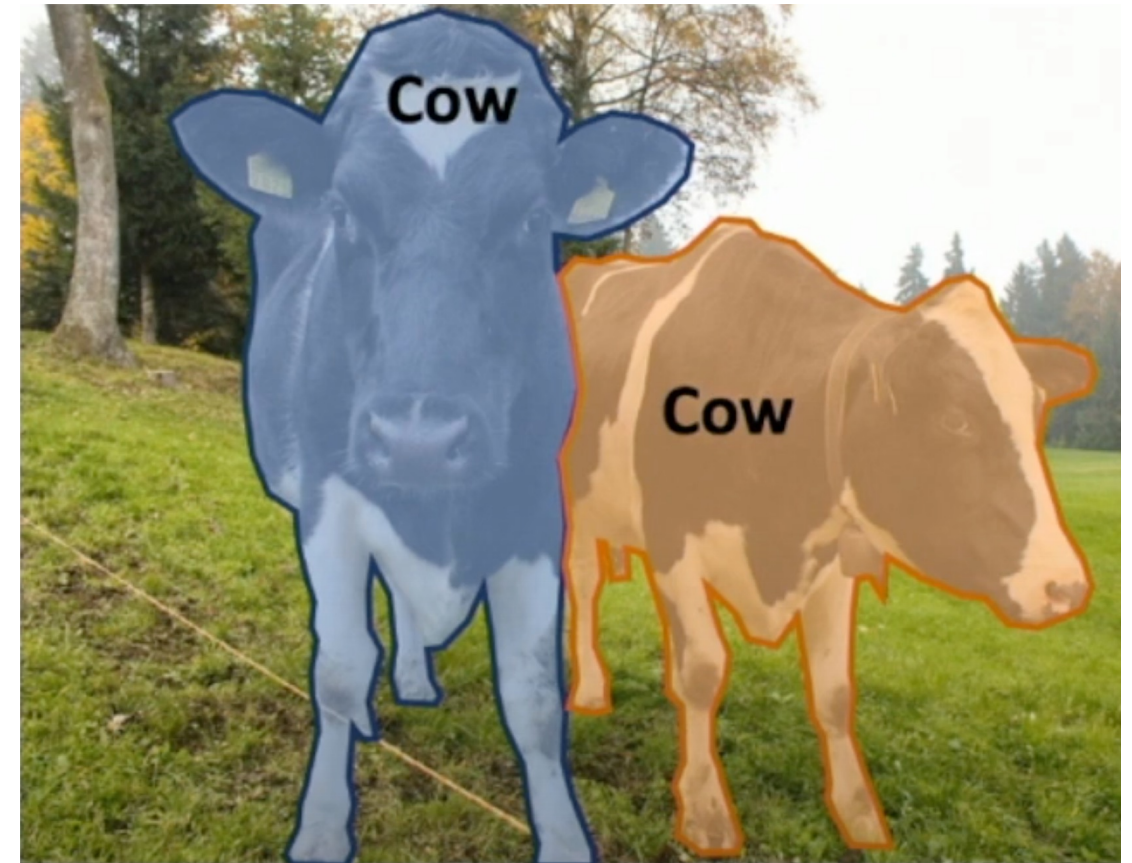
Assign a class label (membership) to every pixel



Semantic segmentation (pixels)

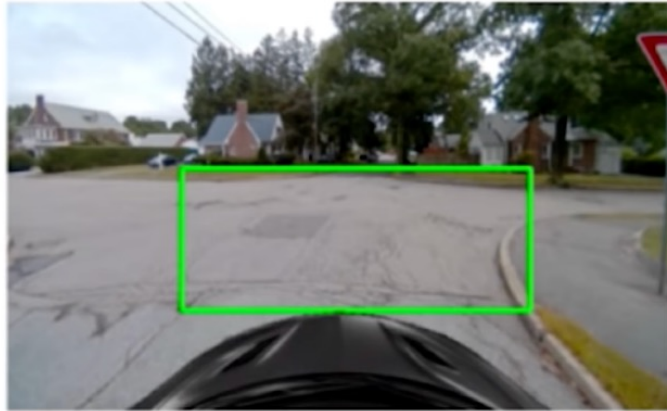


Instance segmentation (objects)

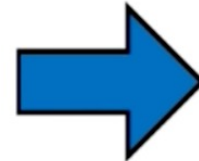
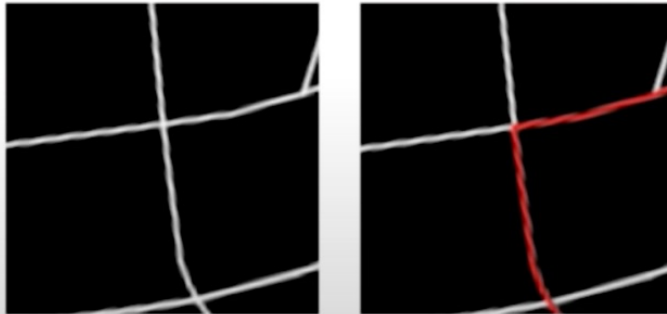


Robotic control

Raw Perception
 I
(ex. camera)



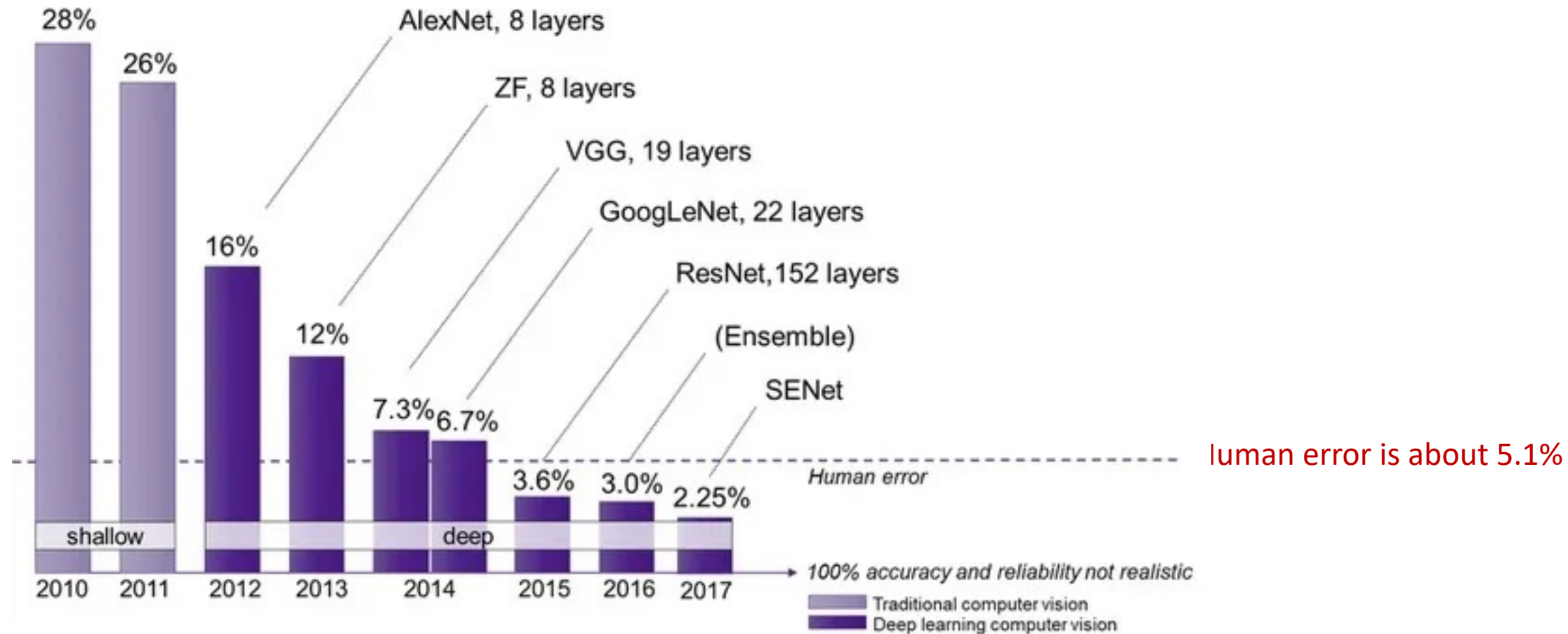
Coarse Maps
 M
(ex. GPS)



Possible Control Commands



Anecdote: ImageNet competition (ILSVRC)



Now let's build our own CNN in Pytorch

Drawer slide

- Output dimension = $(W - F + 2P) / S + 1$
- W: width & height
- F: filter size
- P: padding
- S: stride

Coding Demo:

1_mlp_Day1.ipynb

2_cnn_Day1.ipynb

3_transfer_Day1.ipynb

4_save_load_Day1.ipynb

5_PINNs_Day1.ipynb