

Intro

- In this video you will learn about one more useful layer of neurons;
- We will build our first fully working neural network for images!

A color image input

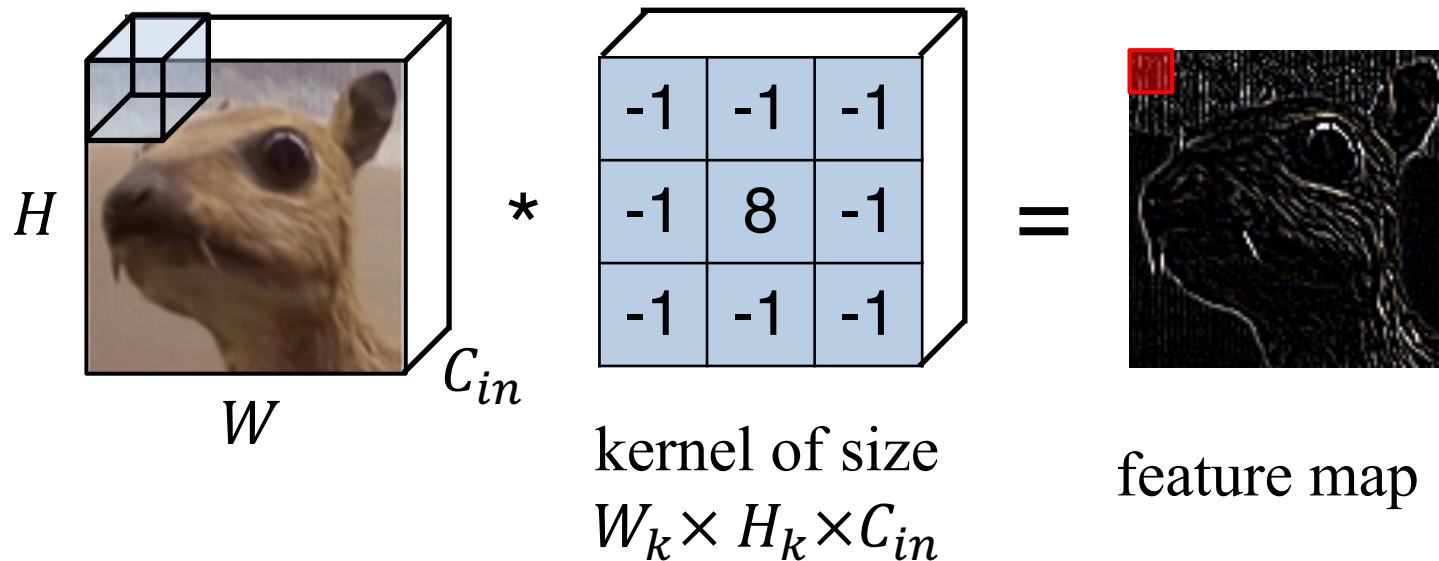
Let's say we have a color image as an input, which is $W \times H \times C_{in}$ **tensor** (multidimensional array), where

- W – is an image width,
- H – is an image height,
- C_{in} – is a number of input channels (e.g. 3 **R****G****B** channels).

A color image input

Let's say we have a color image as an input, which is $W \times H \times C_{in}$ **tensor** (multidimensional array), where

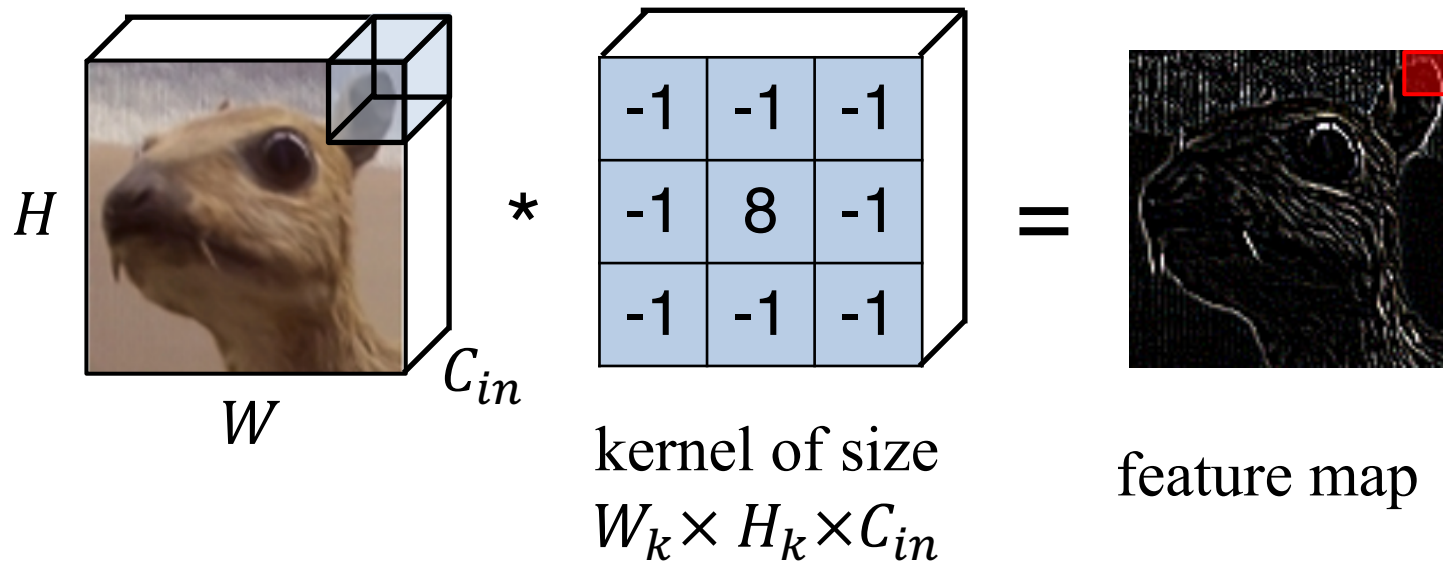
- W – is an image width,
- H – is an image height,
- C_{in} – is a number of input channels (e.g. 3 **R****G****B** channels).



A color image input

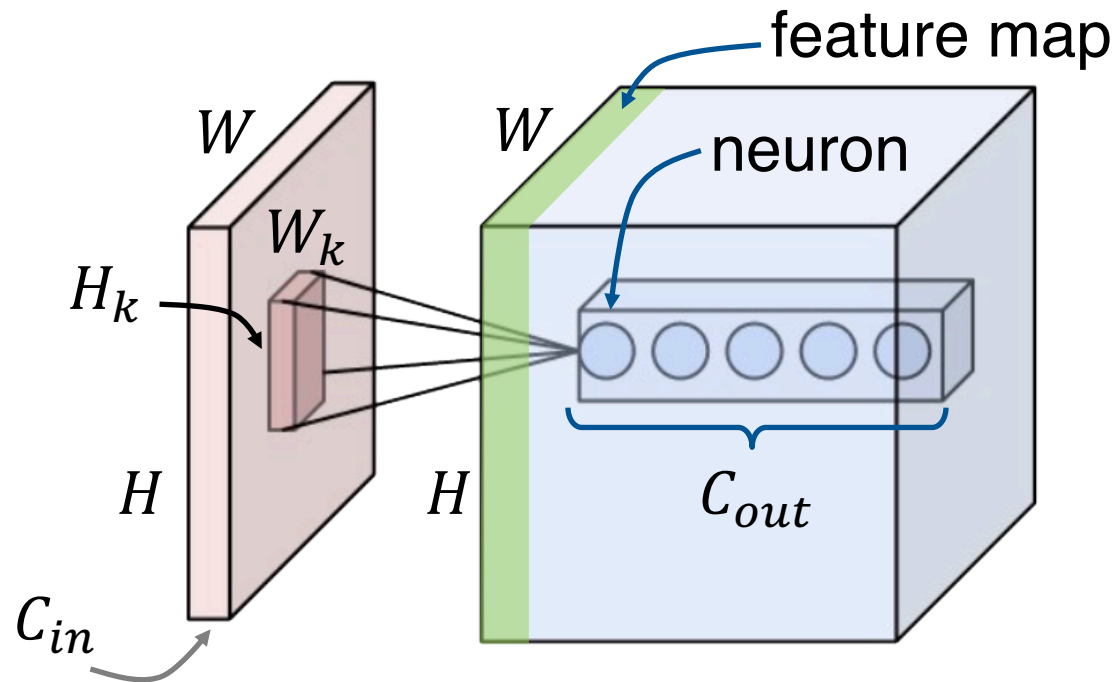
Let's say we have a color image as an input, which is $W \times H \times C_{in}$ **tensor** (multidimensional array), where

- W – is an image width,
- H – is an image height,
- C_{in} – is a number of input channels (e.g. 3 **R****G****B** channels).



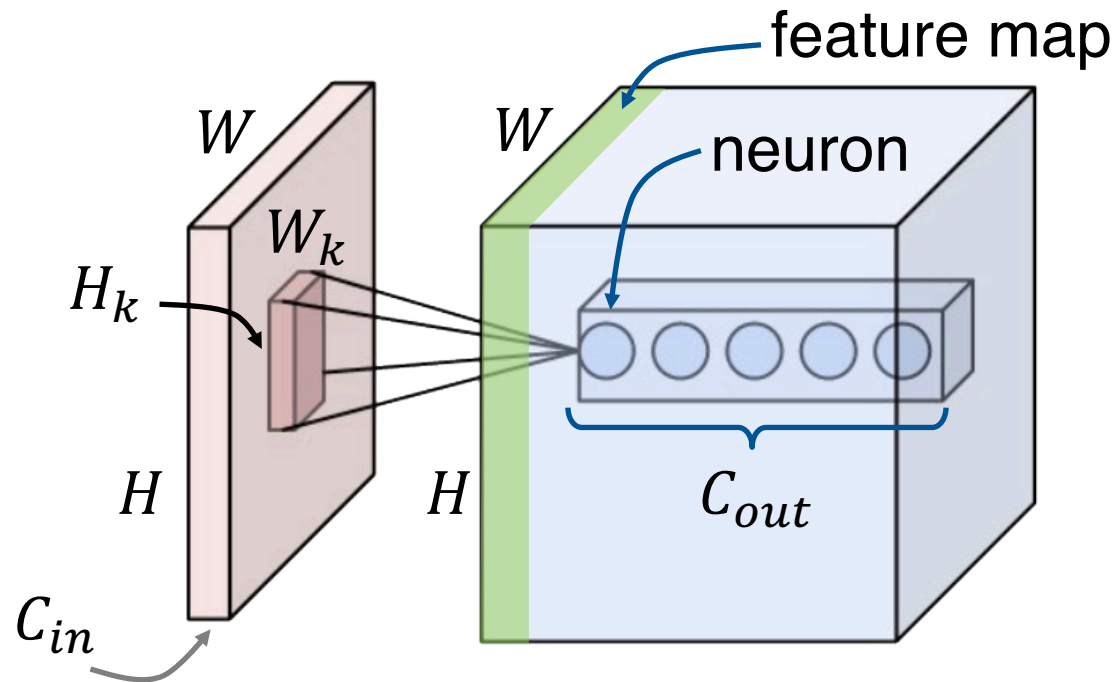
One kernel is not enough!

- We want to train C_{out} kernels of size $W_k \times H_k \times C_{in}$.
- Having a stride of 1 and enough zero padding we can have $W \times H \times C_{out}$ output neurons.



One kernel is not enough!

- We want to train C_{out} kernels of size $W_k \times H_k \times C_{in}$.
- Having a stride of 1 and enough zero padding we can have $W \times H \times C_{out}$ output neurons.



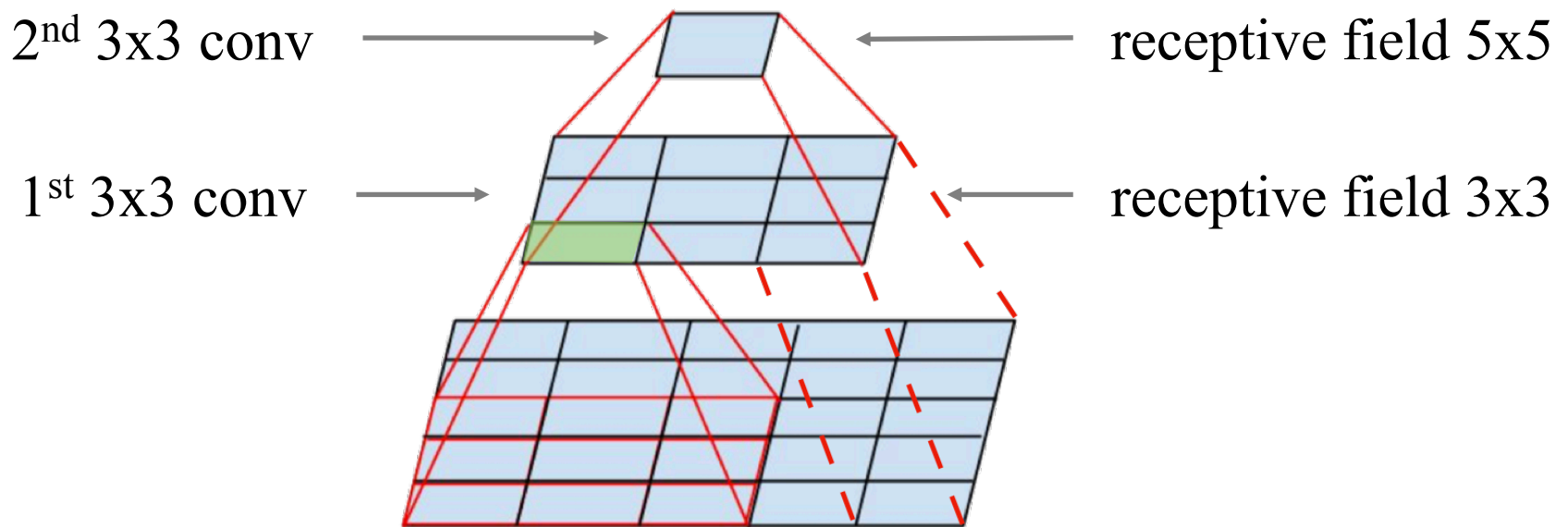
- Using $(W_k * H_k * C_{in} + 1) * C_{out}$ parameters.

One convolutional layer is not enough!

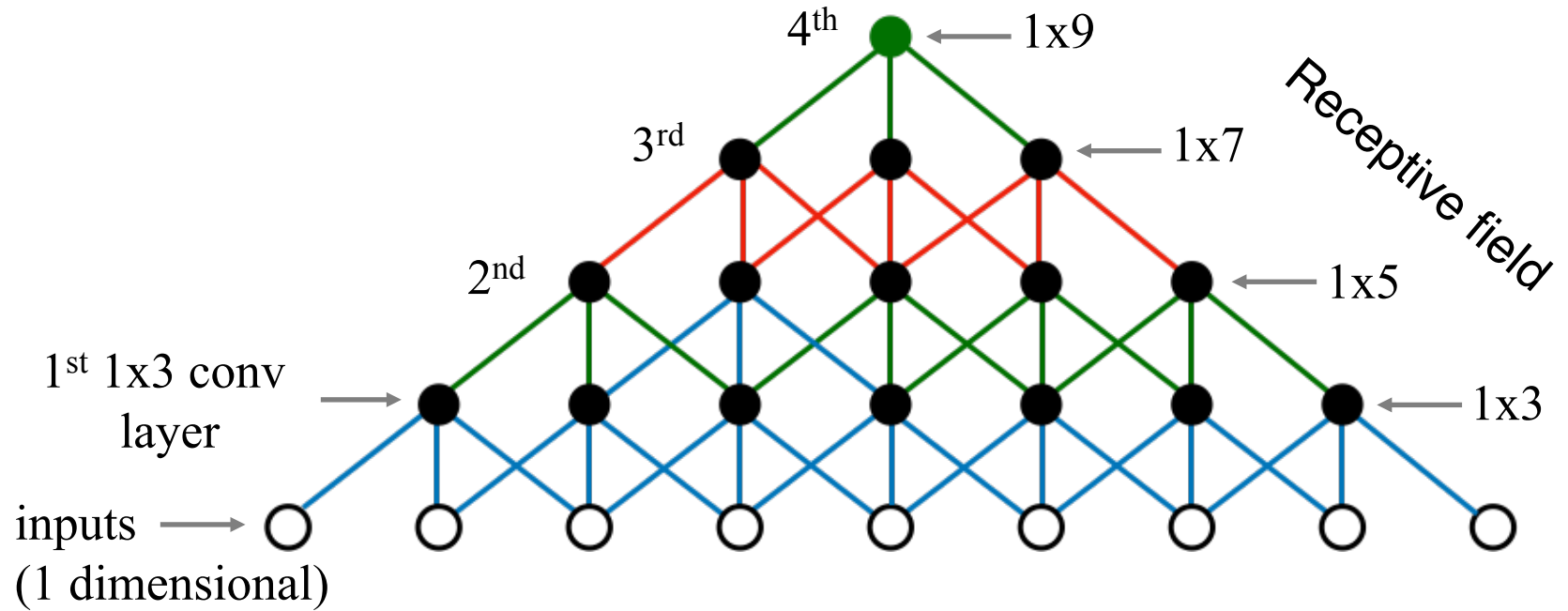
- Let's say neurons of the 1st convolutional layer look at the patches of the image of size 3x3.
- What if an object of interest is bigger than that?
- We need a 2nd convolutional layer on top of the 1st!

One convolutional layer is not enough!

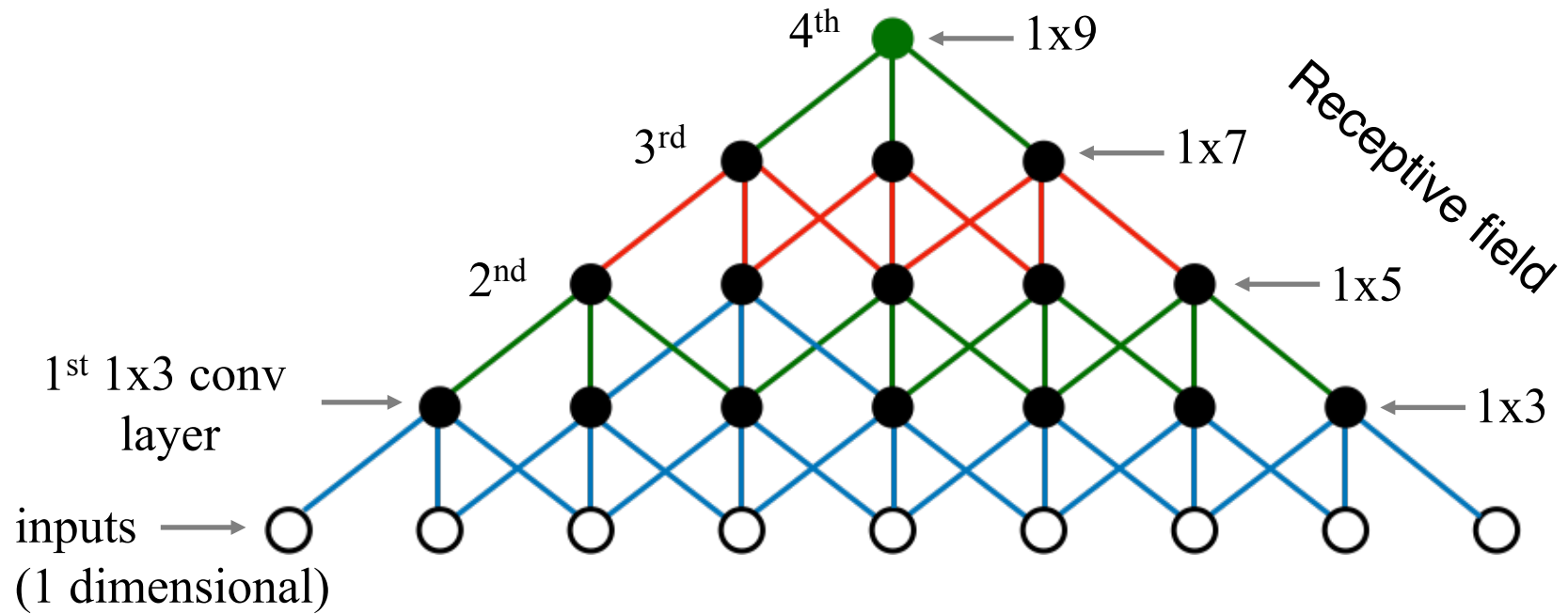
- Let's say neurons of the 1st convolutional layer look at the patches of the image of size 3x3.
- What if an object of interest is bigger than that?
- We need a 2nd convolutional layer on top of the 1st!



Receptive field after N convolutional layers



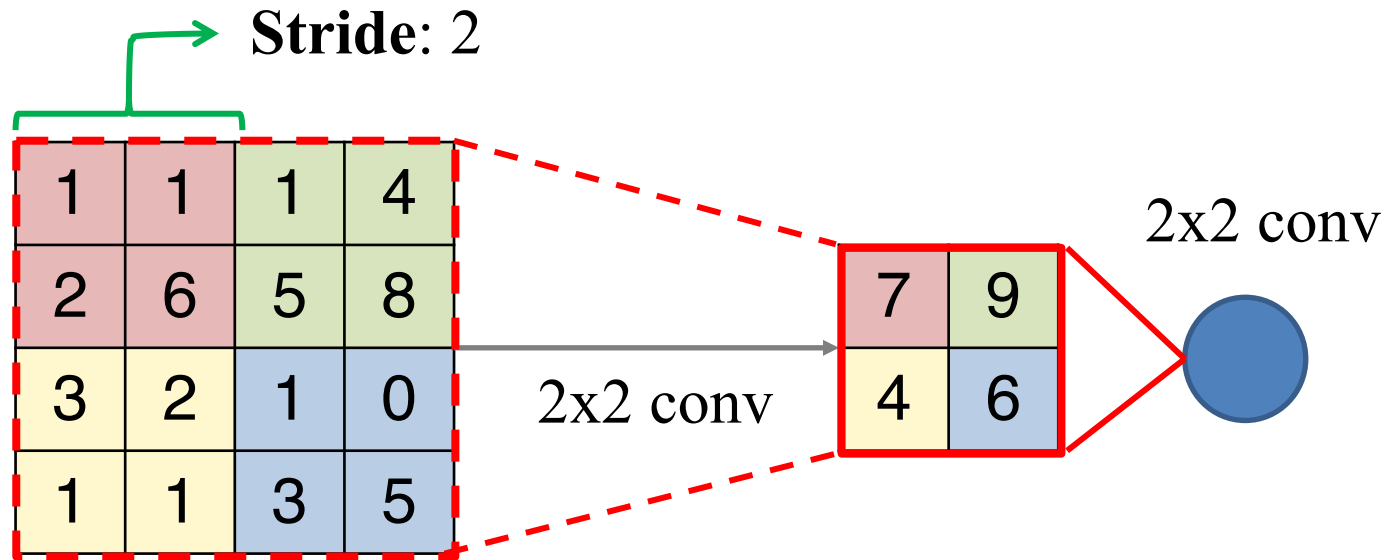
Receptive field after N convolutional layers



- If we stack N convolutional layers with the same kernel size 3×3 the receptive field on N -th layer will be $2N + 1 \times 2N + 1$.
- It looks like we need to stack a lot of convolutional layers!
To be able to identify objects as big as the input image **300x300** we will need **150** convolutional layers!

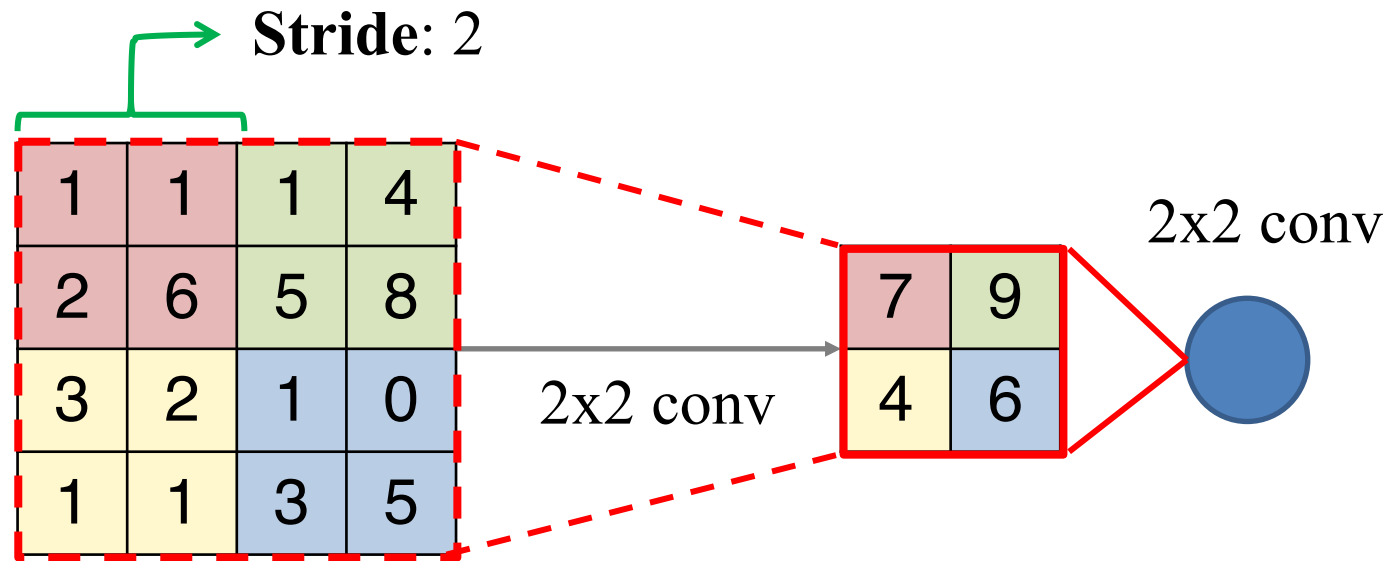
We need to grow receptive field faster!

We can increase a **stride** in our convolutional layer to reduce the output dimensions!



We need to grow receptive field faster!

We can increase a **stride** in our convolutional layer to reduce the output dimensions!



Further convolutions will effectively **double** their receptive field!

How do we maintain translation invariance?

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

Input

*

1	0
0	1

Kernel

=

0	0	0
0	1	0
0	0	2

Output

Max = 2



Didn't
change



Max = 2

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

Input

*

1	0
0	1

Kernel

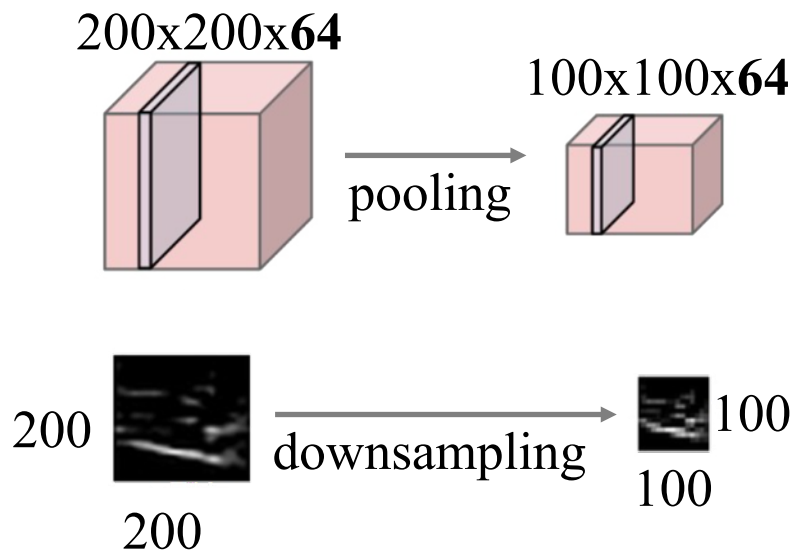
=

2	0	0
0	1	0
0	0	0

Output

Pooling layer will help!

This layer works like a convolutional layer but doesn't have kernel, instead it calculates **maximum** or **average** of input patch values.



Single depth slice

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

The diagram shows the result of a 2x2 max pooling operation with a stride of 2. The output is a 2x2 grid of values: 6, 8, 3, and 5. The values 6 and 8 are in red boxes, 3 is in a yellow box, and 5 is in a blue box.

6	8
3	5

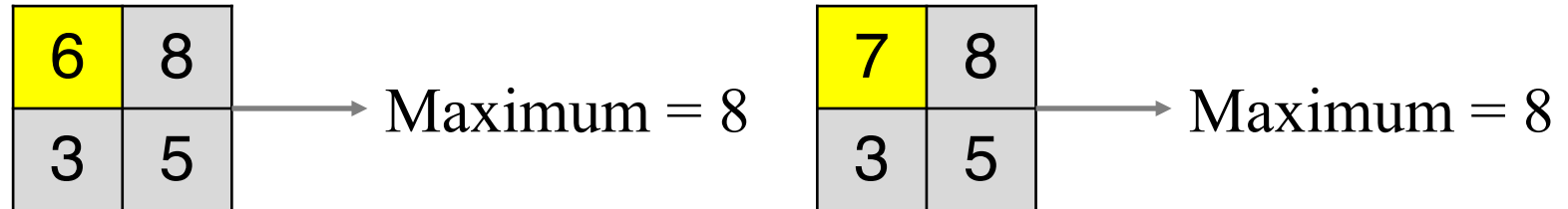
2x2 **max pooling** with stride 2

Backpropagation for max pooling layer

Strictly speaking: maximum is not a differentiable function!

Backpropagation for max pooling layer

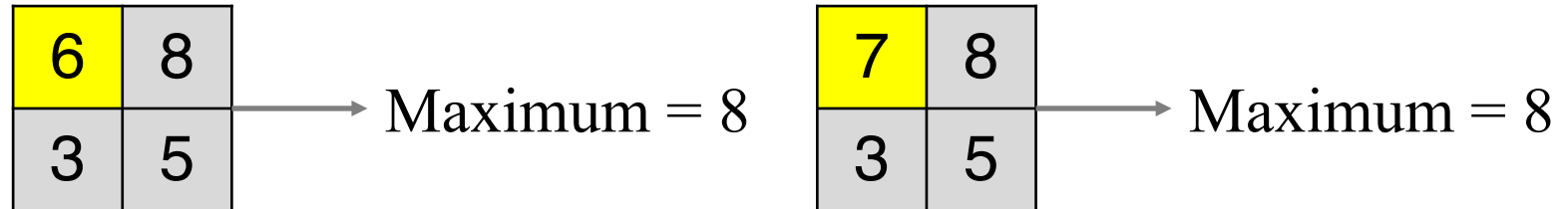
Strictly speaking: maximum is not a differentiable function!



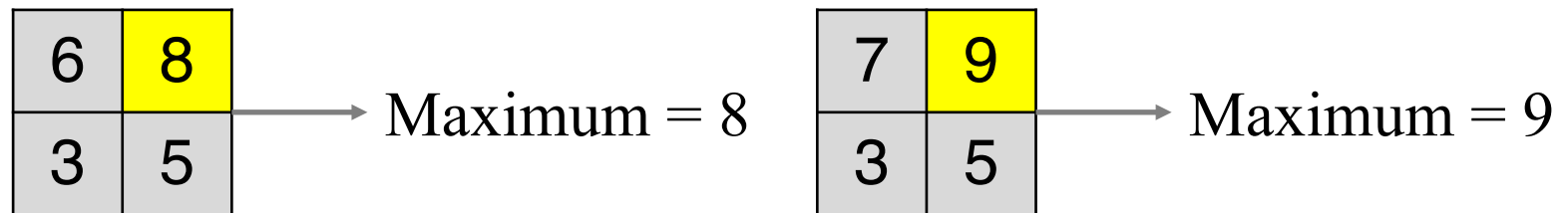
There is no gradient with respect to non maximum patch neurons,
since changing them slightly
does not affect the output.

Backpropagation for max pooling layer

Strictly speaking: maximum is not a differentiable function!



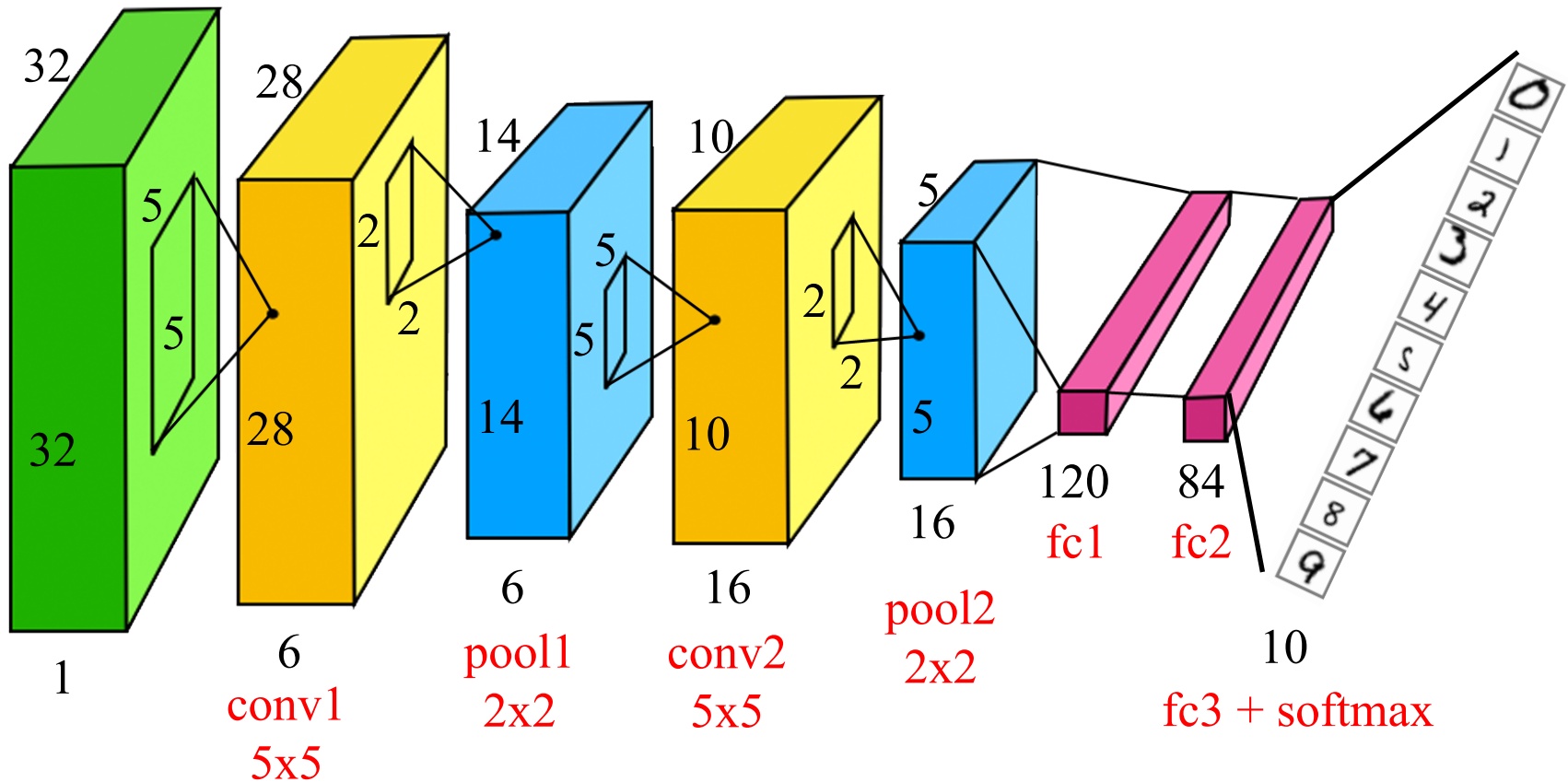
There is no gradient with respect to non maximum patch neurons, since changing them slightly does not affect the output.



For the maximum patch neuron we have a gradient of 1.

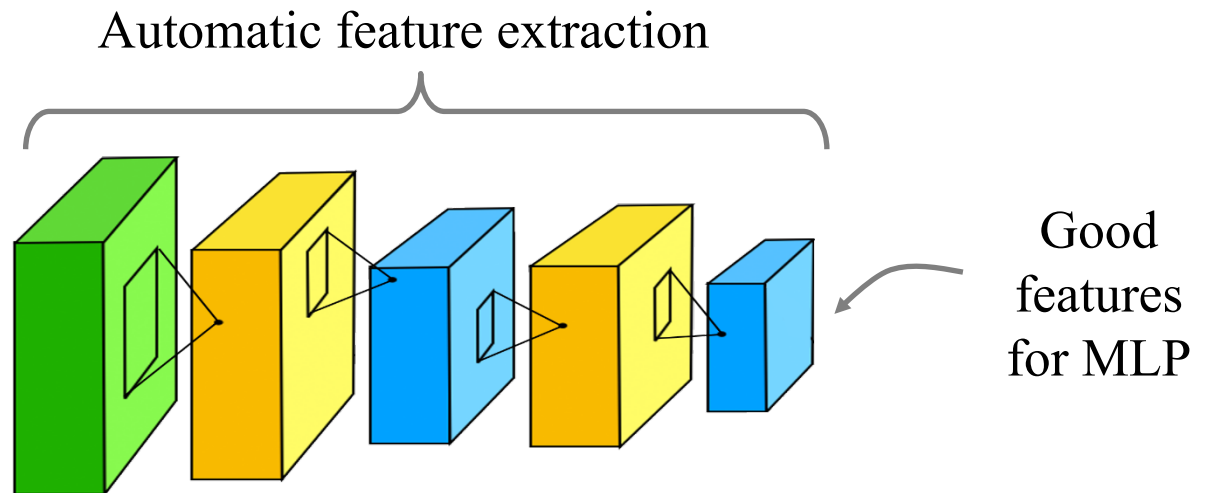
Putting it all together into a simple CNN

LeNet-5 architecture (1998) for handwritten digits recognition on MNIST dataset:



Learning deep representations

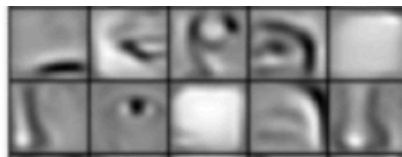
Neurons of deep convolutional layers learn complex representations that can be used as features for classification with MLP.



Inputs that provide highest activations:



conv1



conv2



conv3

Summary

- Using convolutional, pooling and fully connected layers we've built our first network for handwritten digits recognition!
- In the next video we'll overview tips and tricks that are utilized in modern neural network architectures.