



# Convolutional Neural Networks

Convolutions, pooling and CNNs. Neural architectures  
for computer vision.

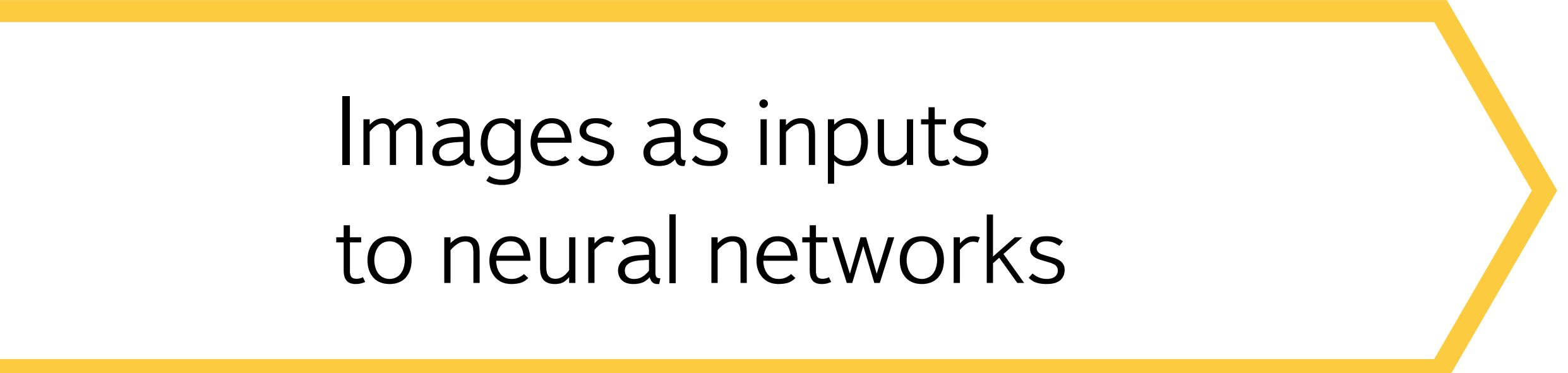
Fourth Machine Learning in High Energy Physics Summer School,  
MLHEP 2018, August 6--12

Alexey Artemov<sup>1,2</sup>

<sup>1</sup>Skoltech <sup>2</sup>National Research University Higher School of Economics

# Lecture overview

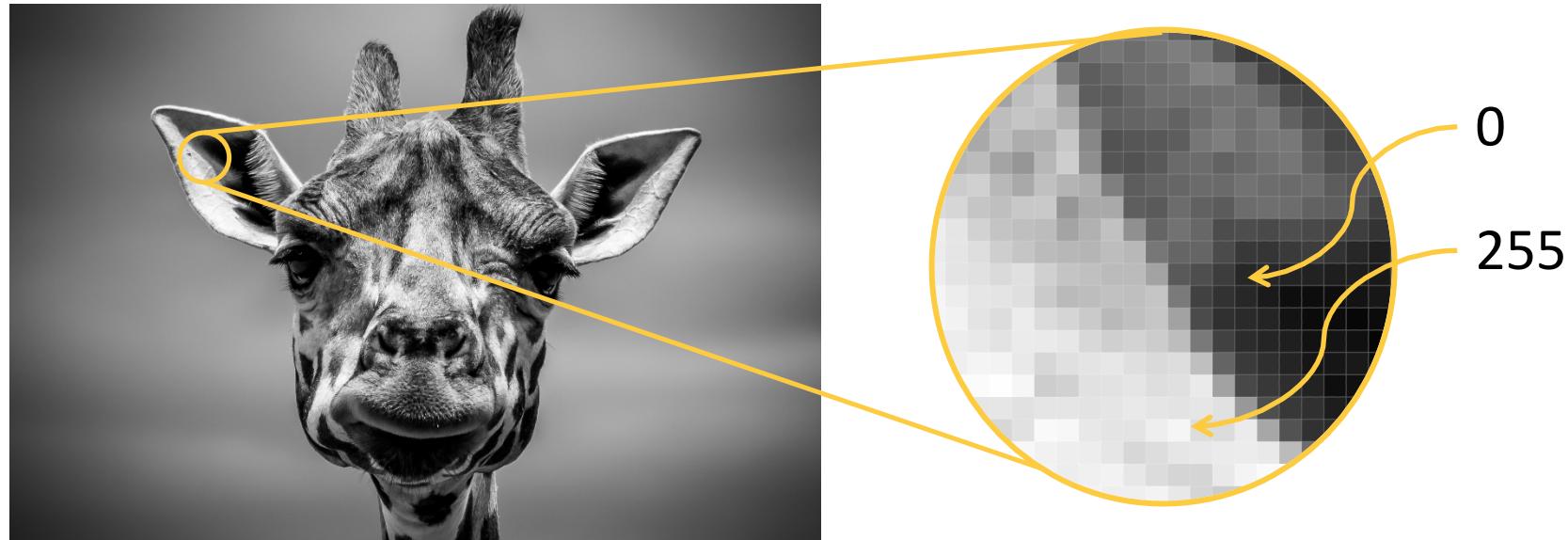
- Digital images and processing by neural networks
- Image processing operations: convolutions and pooling
- Convolutional neural networks from scratch
- Modern computer vision architectures: AlexNet, VGG, Inception and ResNets



Images as inputs  
to neural networks

# Digital representation of an image

- Grayscale image is a matrix of pixels (picture elements)
- Dimensions of this matrix are called image resolution (e.g.  $300 \times 300$ )
- Each pixel stores its brightness (or intensity) ranging from 0 to 255, 0 intensity corresponds to black color:



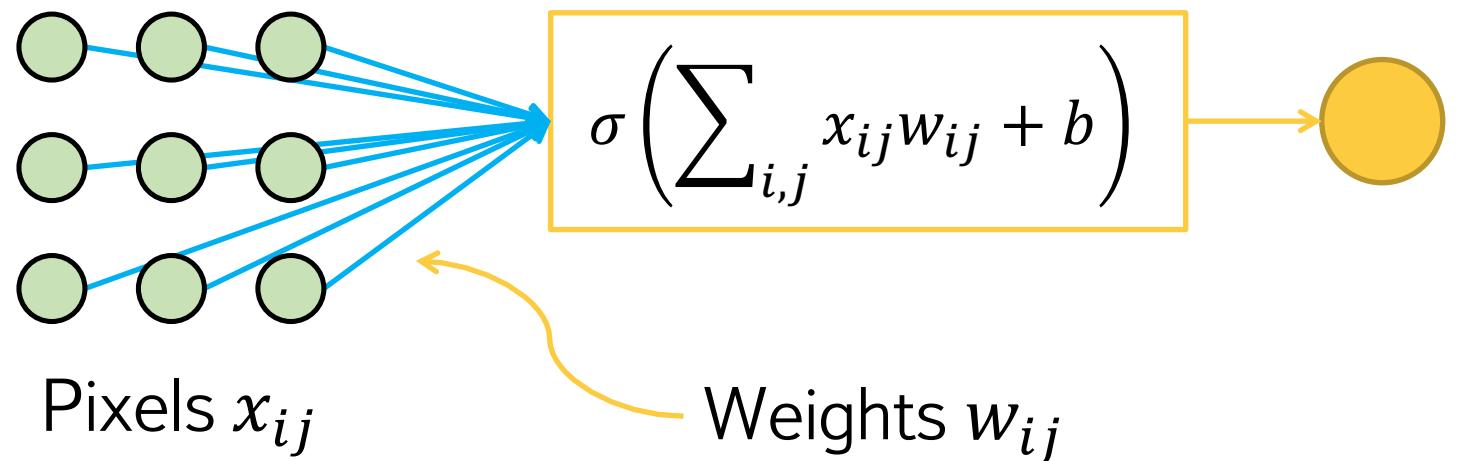
# Image as a neural network input

- Normalize input pixels:  $x_{norm} = \frac{x}{255} - 0.5$

# Image as a neural network input

- Normalize input pixels:  $x_{norm} = \frac{x}{255} - 0.5$

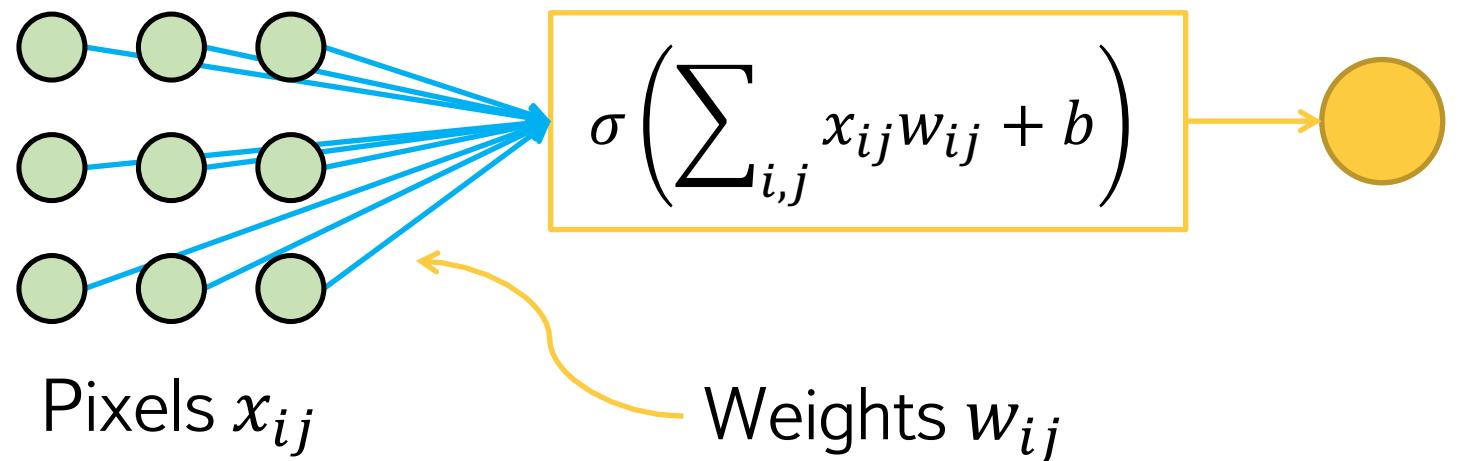
- Maybe MLP will work?



# Image as a neural network input

- Normalize input pixels:  $x_{norm} = \frac{x}{255} - 0.5$

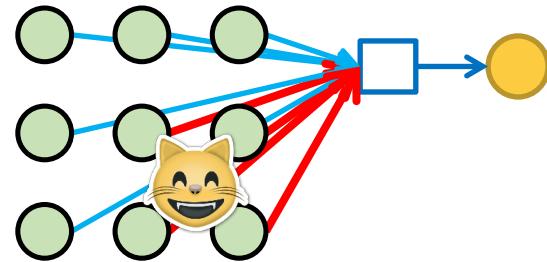
- Maybe MLP will work?



- Actually, no!

# Why not MLP?

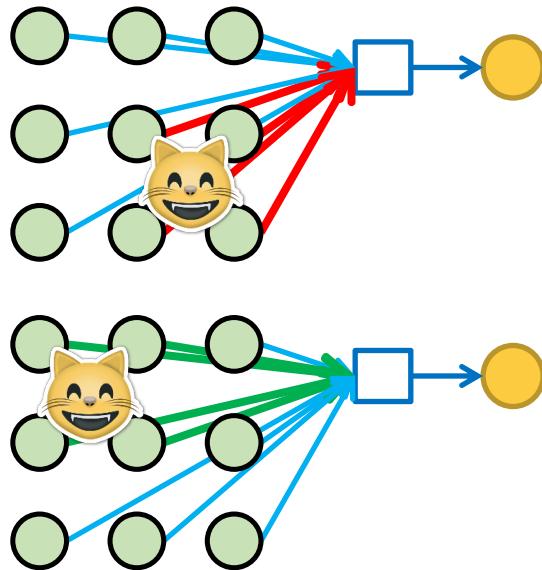
- Let's say we want to train a "cat detector"



On this training image red weights  $w_{ij}$  will change a little bit to better detect a cat

# Why not MLP?

- Let's say we want to train a "cat detector"

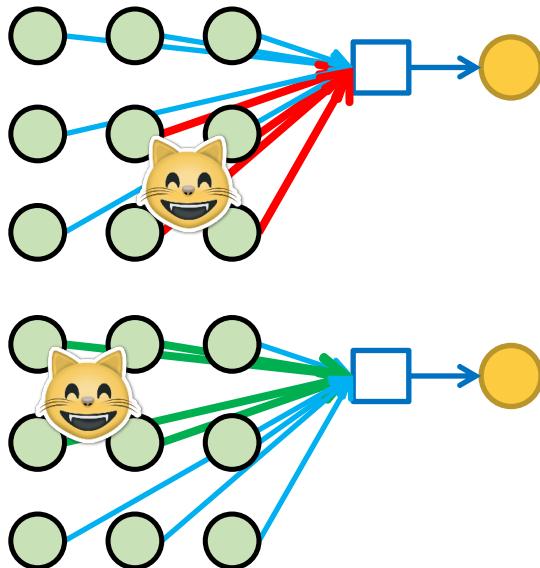


On this training image red weights  $w_{ij}$  will change a little bit to better detect a cat

On this training image green weights  $w_{ij}$  will change...

# Why not MLP?

- Let's say we want to train a "cat detector"



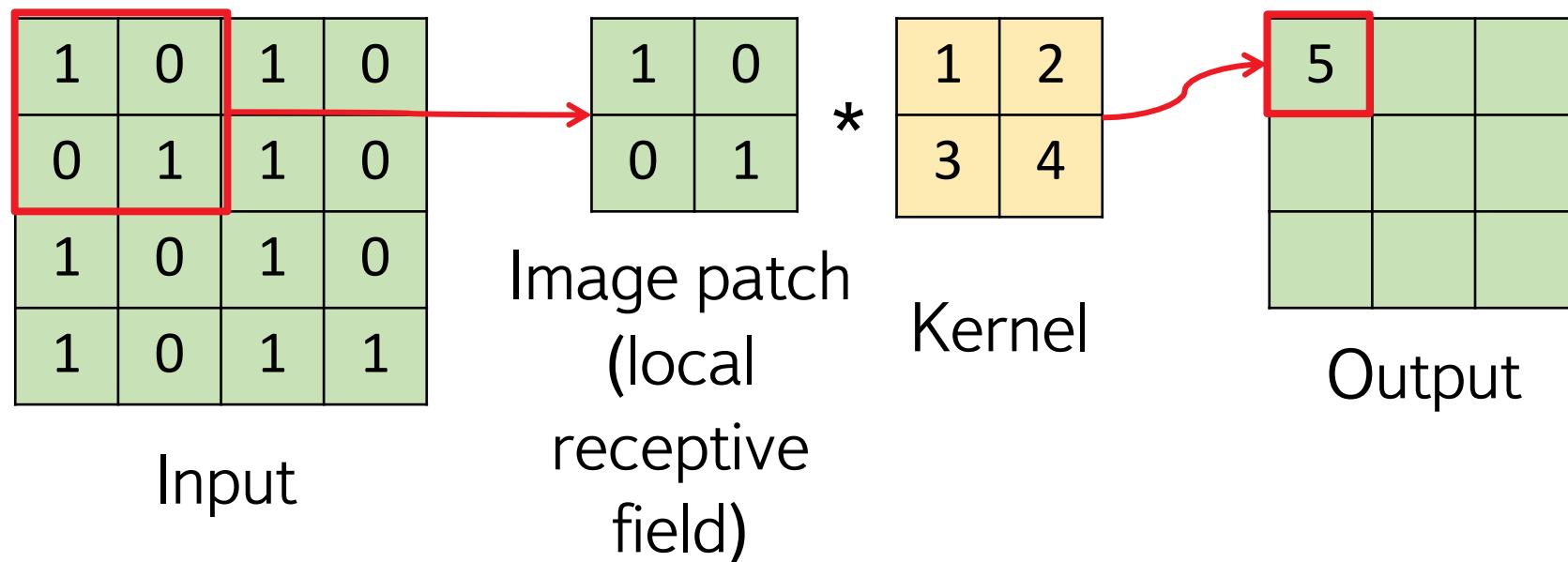
- We learn the same "cat features" in different areas and don't fully utilize the training set!
- What if cats in the test set appear in different places?

On this training image **red** weights  $w_{ij}$  will change a little bit to better detect a cat

On this training image **green** weights  $w_{ij}$  will change...

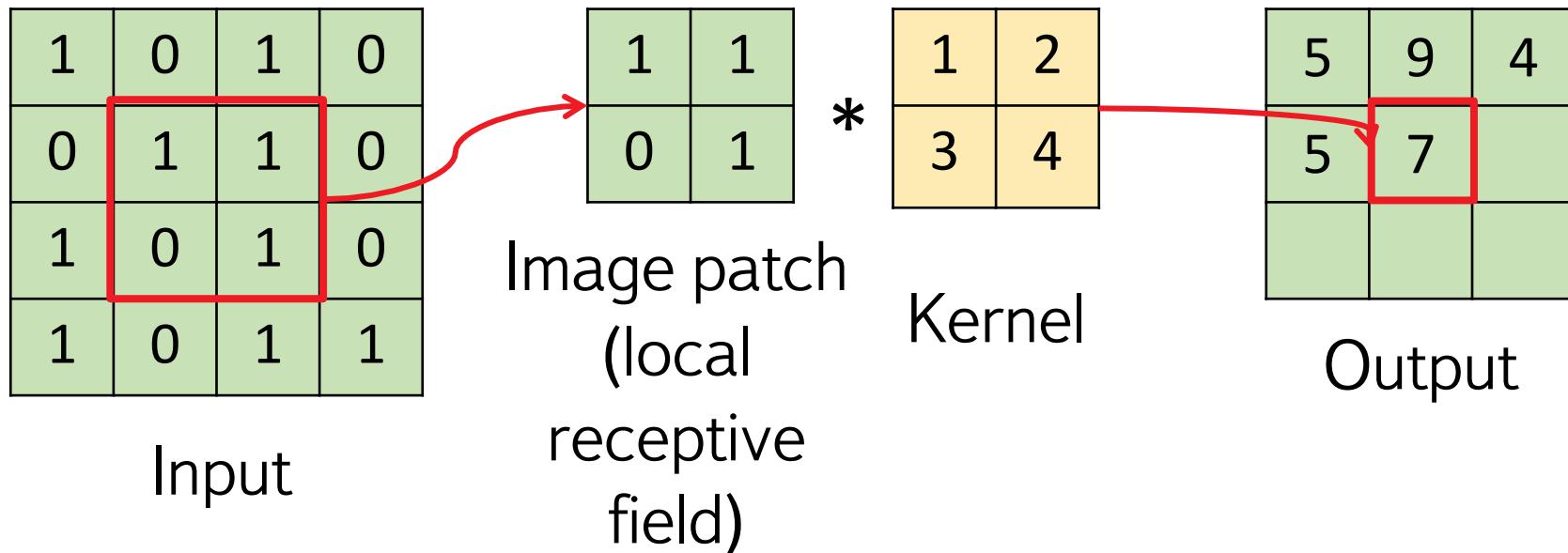
# Convolutions will help!

- Convolution is a dot product of a kernel (or filter) and a patch of an image (local receptive field) of the same size



# Convolutions will help!

- Convolution is a dot product of a kernel (or filter) and a patch of an image (local receptive field) of the same size



# Convolutions have been used for a while



Original  
image

$$\text{Original image} * \begin{matrix} \text{Kernel} \\ \begin{array}{|c|c|c|} \hline -1 & -1 & -1 \\ \hline -1 & 8 & -1 \\ \hline -1 & -1 & -1 \\ \hline \end{array} \end{matrix} = \begin{matrix} \text{Edge detection} \\ \text{Resulting image showing edges} \end{matrix}$$

The diagram illustrates the convolution process for edge detection. On the left is the original image of a dog's head. In the center is a 3x3 kernel matrix with values: -1, -1, -1 in the top row; -1, 8, -1 in the middle row; and -1, -1, -1 in the bottom row. To the right of the kernel is an equals sign. To the right of the equals sign is the resulting image, which shows the edges of the dog's face highlighted in white against a dark background, demonstrating edge detection.

Sums up to 0 (black color)  
when the patch is a solid fill

# Convolutions have been used for a while



Original  
image

$$\text{Original image} * \begin{matrix} \text{Kernel} \\ \begin{array}{ccc} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{array} \end{matrix} = \text{Edge detection}$$
A processed image showing the edges of the squirrel's head highlighted in white against a dark background.

$$\text{Original image} * \begin{matrix} \text{Kernel} \\ \begin{array}{ccc} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{array} \end{matrix} = \text{Sharpening}$$
A processed image where the edges of the squirrel's head appear more prominent and the overall image has a slightly increased contrast.

Doesn't change an image for solid fills  
Adds a little intensity on the edges

# Convolutions have been used for a while



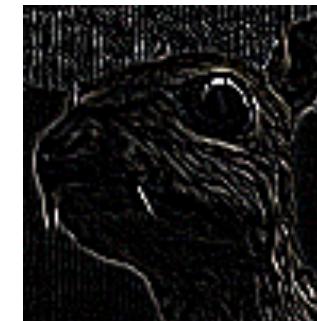
Original  
image

Kernel

-1	-1	-1
-1	8	-1
-1	-1	-1

\*

=



Edge  
detection

0	-1	0
-1	5	-1
0	-1	0

\*

=

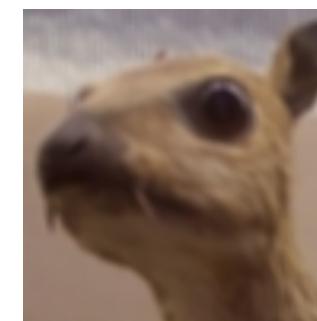


Sharpening

1	1	1
1	1	1
1	1	1

$\ast \frac{1}{9}$

=



Blurring

# Convolution is similar to correlation

$$\begin{array}{|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 1 & 0 \\ \hline 0 & 0 & 0 & 1 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 0 & 1 & 0 \\ \hline 0 & 0 & 2 \\ \hline \end{array}$$

Input                      Kernel                      Output

# Convolution is similar to correlation

$$\begin{array}{|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 1 & 0 \\ \hline 0 & 0 & 0 & 1 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 0 & 1 & 0 \\ \hline 0 & 0 & 2 \\ \hline \end{array}$$

Input

Kernel

Output

$$\begin{array}{|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 1 \\ \hline 0 & 0 & 1 & 0 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 0 & 0 & 1 \\ \hline 0 & 1 & 0 \\ \hline \end{array}$$

Input

Kernel

Output

# Convolution is similar to correlation

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

Input

1	0
0	1

\*

=

0	0	0
0	1	0
0	0	2

Max = 2

Simple  
classifier

Max = 1

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

Input

1	0
0	1

\*

=

0	0	0
0	0	1
0	1	0

Output

# Convolution is translation equivariant

$$\begin{array}{|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 1 & 0 \\ \hline 0 & 0 & 0 & 1 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 0 & 1 & 0 \\ \hline 0 & 0 & 2 \\ \hline \end{array}$$

Input                    Kernel                    Output

# Convolution is translation equivariant

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

Input

1	0
0	1

\*

=

0	0	0
0	1	0
0	0	2

at each position an

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

Input

1	0
0	1

\*

=

2	0	0
0	1	0
0	0	0

Output

# Convolution is translation equivariant

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

Input

1	0
0	1

\*

=

0	0	0
0	1	0
0	0	2

Max = 2

Didn't  
change

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

Input

1	0
0	1

\*

=

2	0	0
0	1	0
0	0	0

Max = 2

Output

# Convolutional layer in neural network

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

Input 3x3  
image with  
zero **padding**  
(grey area)

multiplication summed ("do")

$\sigma(w_6 + w_8 + w_9 + b)$	...	...
...	...	...
...	...	...

9 output neurons (**feature map**)  
with only 10 parameters

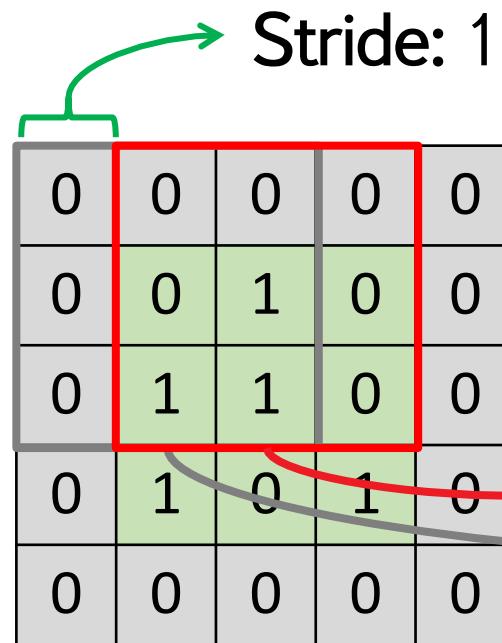
Shared bias:

$b$

Shared kernel:

$w_1$	$w_2$	$w_3$
$w_4$	$w_5$	$w_6$
$w_7$	$w_8$	$w_9$

# Convolutional layer in neural network



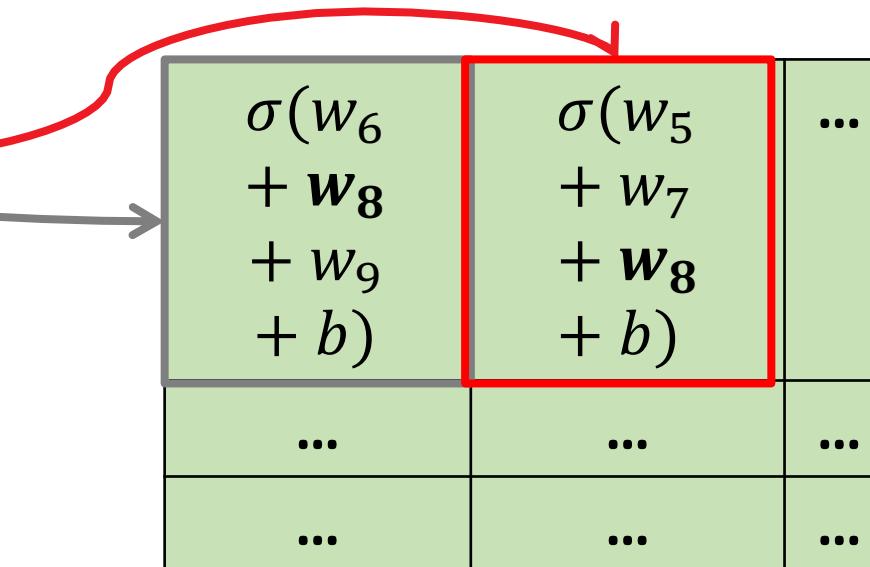
Input 3x3  
image with  
zero **padding**  
(grey area)

Shared bias:

$b$

Shared kernel:

$w_1$	$w_2$	$w_3$
$w_4$	$w_5$	$w_6$
$w_7$	$w_8$	$w_9$

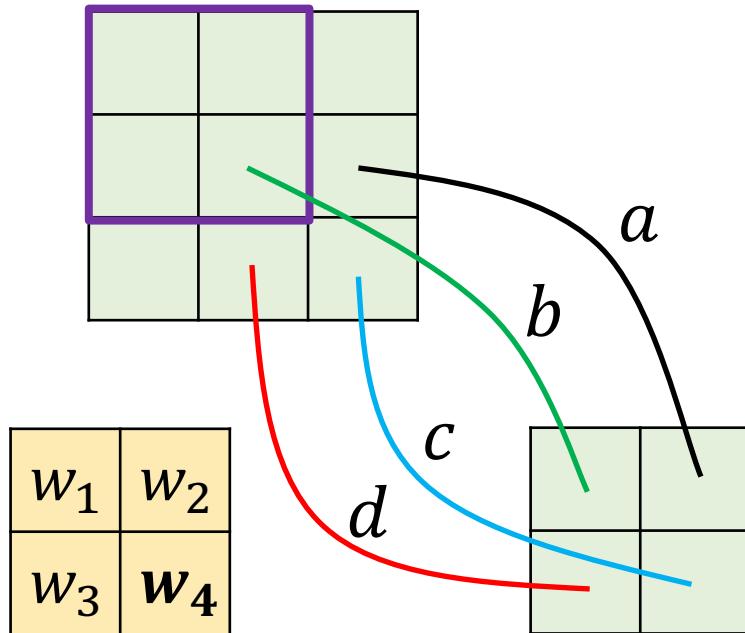


9 output neurons (**feature map**)  
with only 10 parameters

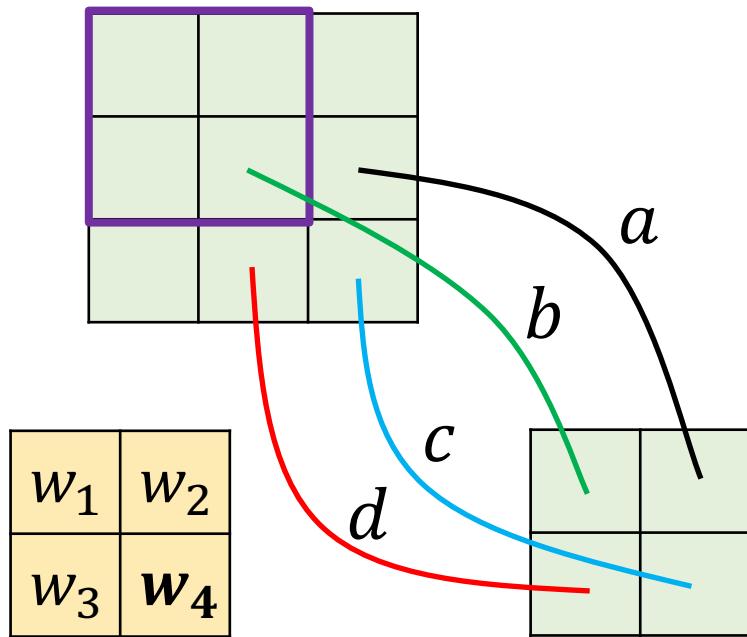
weights are \*Sha

# Backpropagation for CNN

Gradients are first calculated as if the kernel weights were not shared:



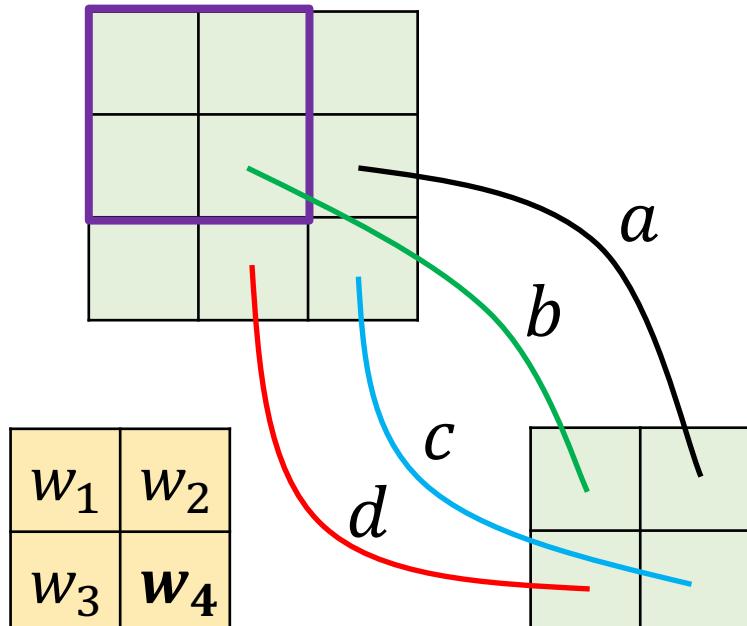
# Backpropagation for CNN



Gradients are first calculated as if the kernel weights were not shared:

$$a = a - \gamma \frac{\partial L}{\partial a} \quad b = b - \gamma \frac{\partial L}{\partial b}$$
$$c = c - \gamma \frac{\partial L}{\partial c} \quad d = d - \gamma \frac{\partial L}{\partial d}$$

# Backpropagation for CNN

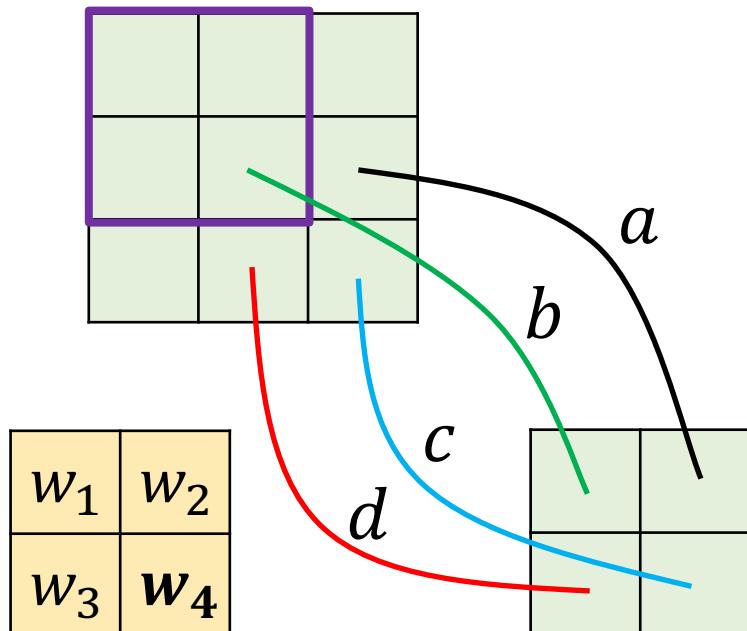


Gradients are first calculated as if the kernel weights were not shared:

$$a = a - \gamma \frac{\partial L}{\partial a} \quad b = b - \gamma \frac{\partial L}{\partial b}$$
$$c = c - \gamma \frac{\partial L}{\partial c} \quad d = d - \gamma \frac{\partial L}{\partial d}$$

$$w_4 = w_4 - \gamma \left( \frac{\partial L}{\partial a} + \frac{\partial L}{\partial b} + \frac{\partial L}{\partial c} + \frac{\partial L}{\partial d} \right)$$

# Backpropagation for CNN



Gradients are first calculated as if the kernel weights were not shared:

$$a = a - \gamma \frac{\partial L}{\partial a} \quad b = b - \gamma \frac{\partial L}{\partial b}$$
$$c = c - \gamma \frac{\partial L}{\partial c} \quad d = d - \gamma \frac{\partial L}{\partial d}$$

$$w_4 = w_4 - \gamma \left( \frac{\partial L}{\partial a} + \frac{\partial L}{\partial b} + \frac{\partial L}{\partial c} + \frac{\partial L}{\partial d} \right)$$

Gradients of the same shared weight are summed up!

# Convolutional vs fully connected layer

- In convolutional layer, the same kernel is used for every output neuron, this way we share parameters of the network and train a better model;

# Convolutional vs fully connected layer

- In convolutional layer, the same kernel is used for every output neuron, this way we share parameters of the network and train a better model;
- 300x300 input, 300x300 output, 5x5 kernel – 26 parameters in convolutional layer and  $8.1 \times 10^9$  parameters in fully connected layer (each output is a perceptron);

# Convolutional vs fully connected layer

- In convolutional layer, the same kernel is used for every output neuron, this way we share parameters of the network and train a better model;
- 300x300 input, 300x300 output, 5x5 kernel – 26 parameters in convolutional layer and  $8.1 \times 10^9$  parameters in fully connected layer (each output is a perceptron);
- Convolutional layer can be viewed as a special case of a fully connected layer when all the weights outside the local receptive field of each neuron equal 0 and kernel parameters are shared between neurons.

# Intermediate summary

- We've introduced a convolutional layer which works better than fully connected layer for images: it has fewer parameters and acts the same for every patch of input.
- This layer will be used as a building block for larger neural networks!

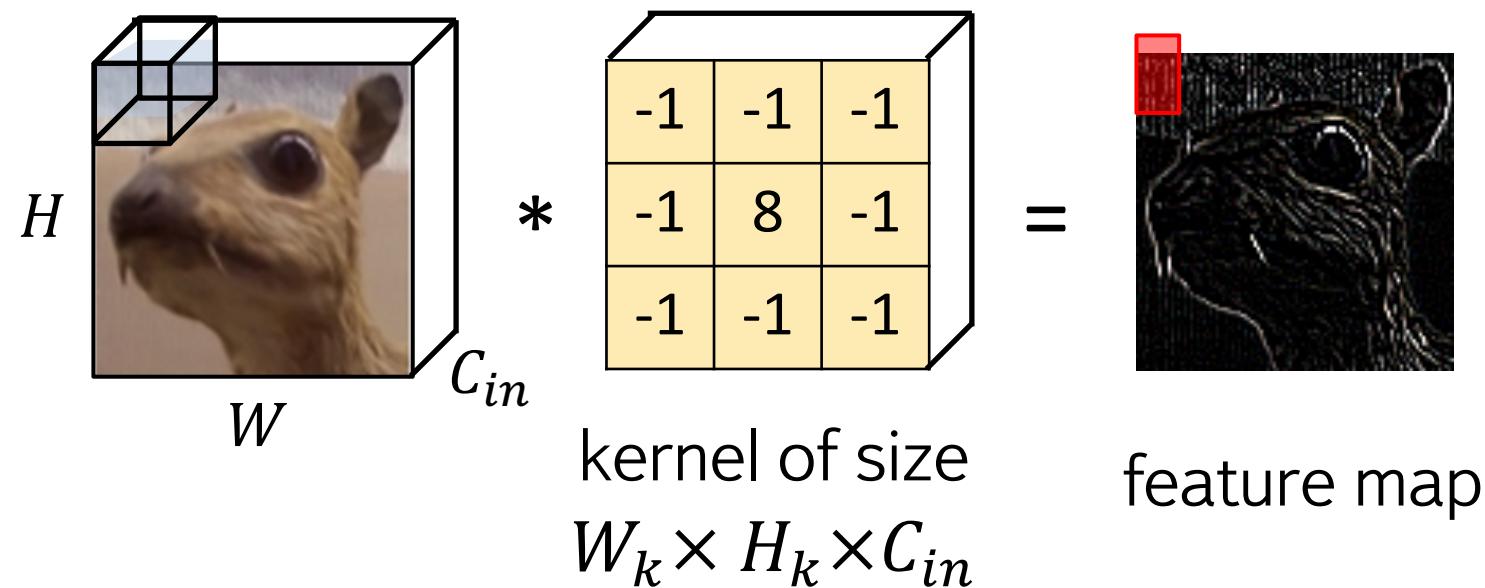
# Building convolutional neural networks for vision

# 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 **RGB** channels).

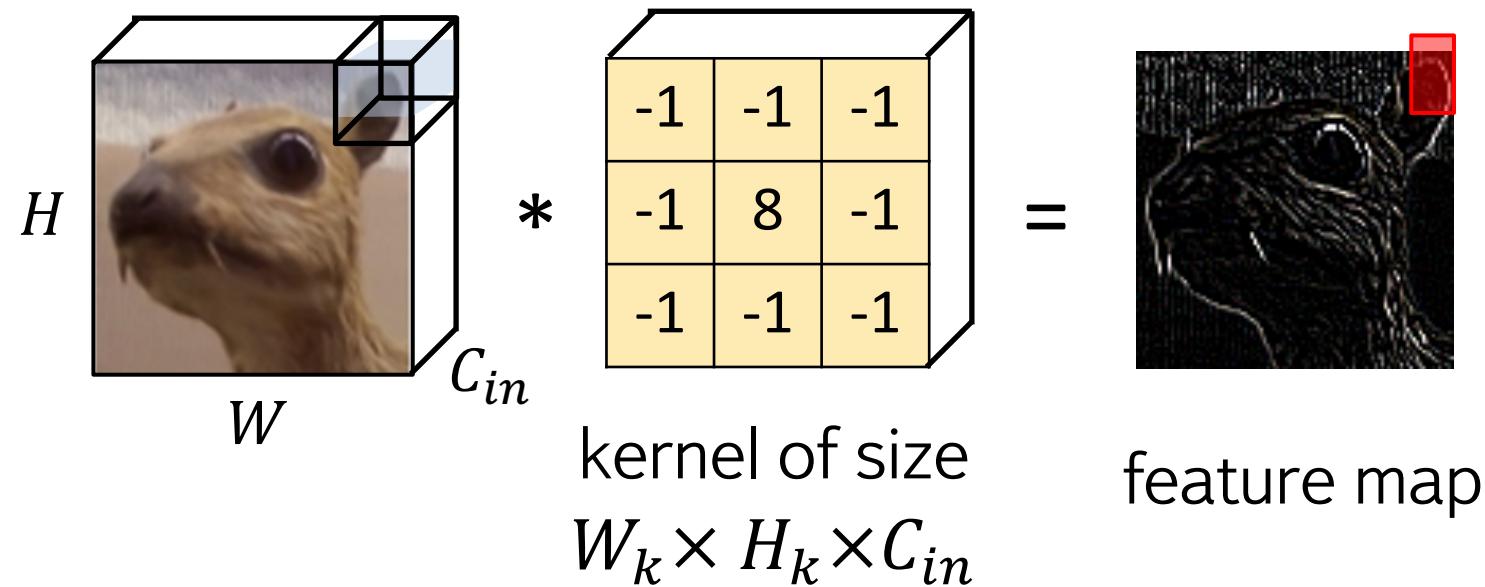
# A color image input

- 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 **RGB** channels).



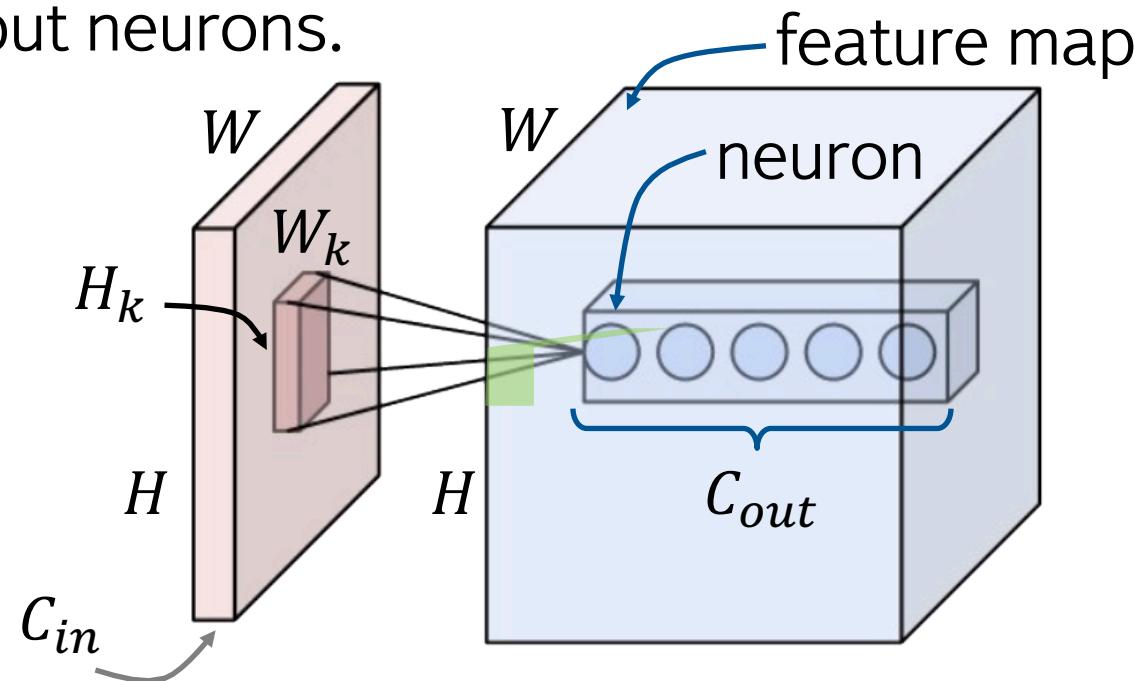
# A color image input

- 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 **RGB** channels).



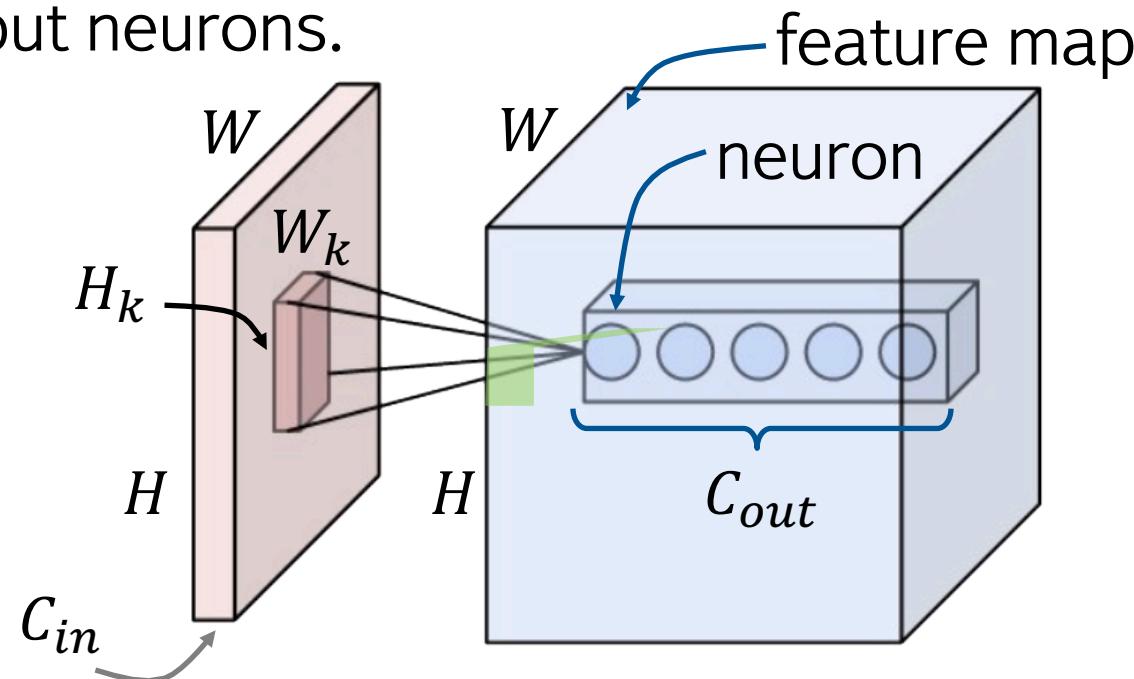
# One kernel is not enough!

- 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.



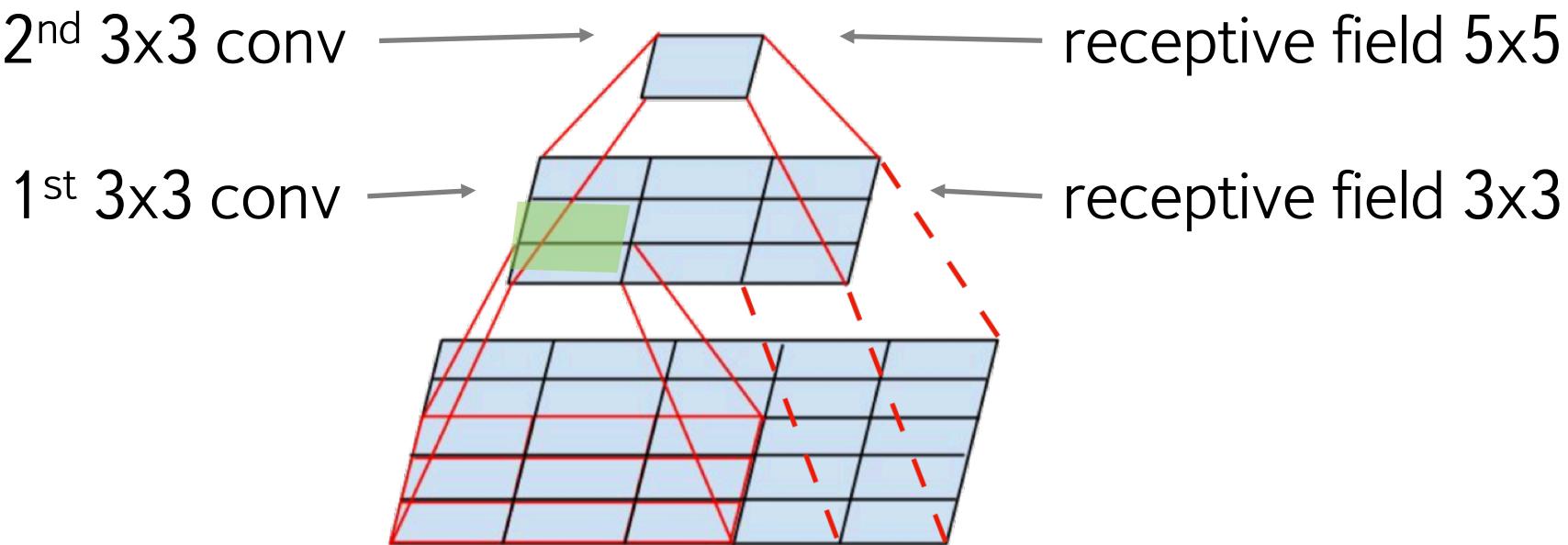
- Using  $(W_k * H_k * C_{in} + 1) * C_{out}$  parameters.

# One convolutional layer is not enough!

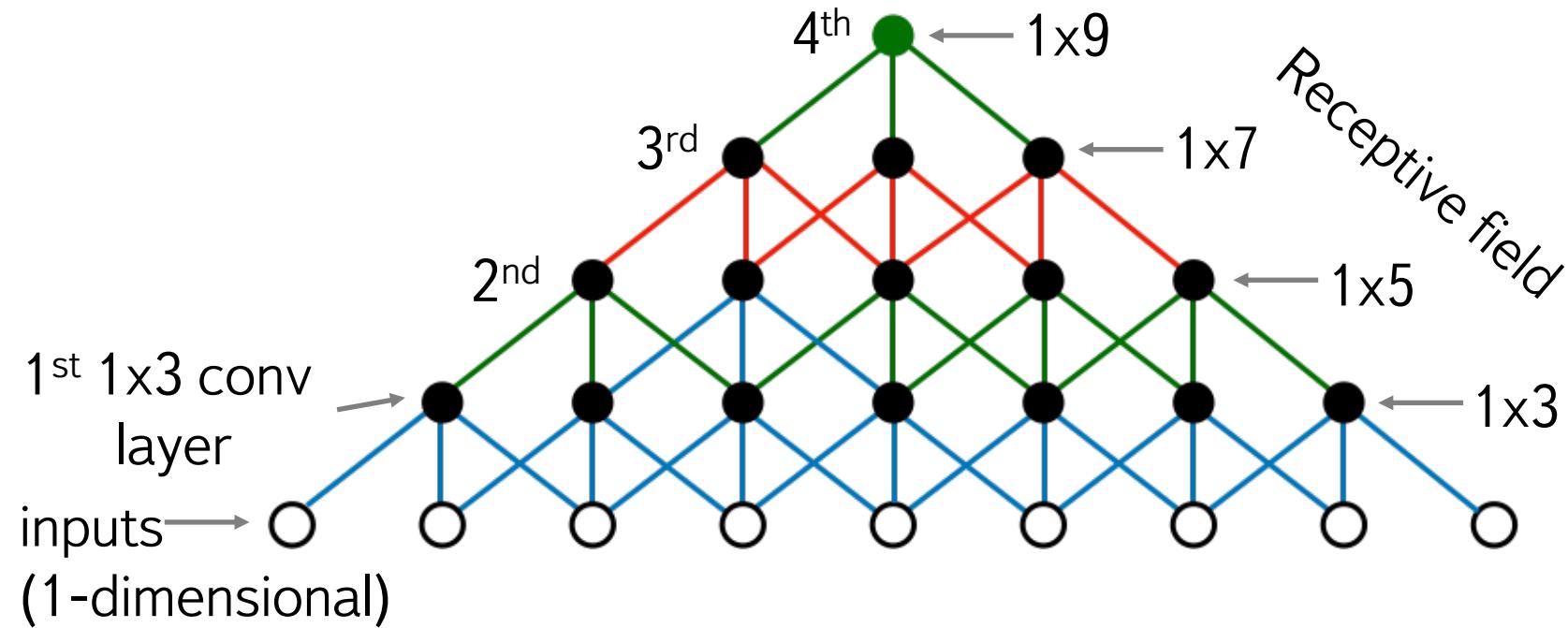
- Let's say neurons of the 1st convolutional layer look at the patches of the image of size 3x3.
- What if an object of interest is bigger than that?
- We need a 2<sup>nd</sup> convolutional layer on top of the 1<sup>st</sup>!

# One convolutional layer is not enough!

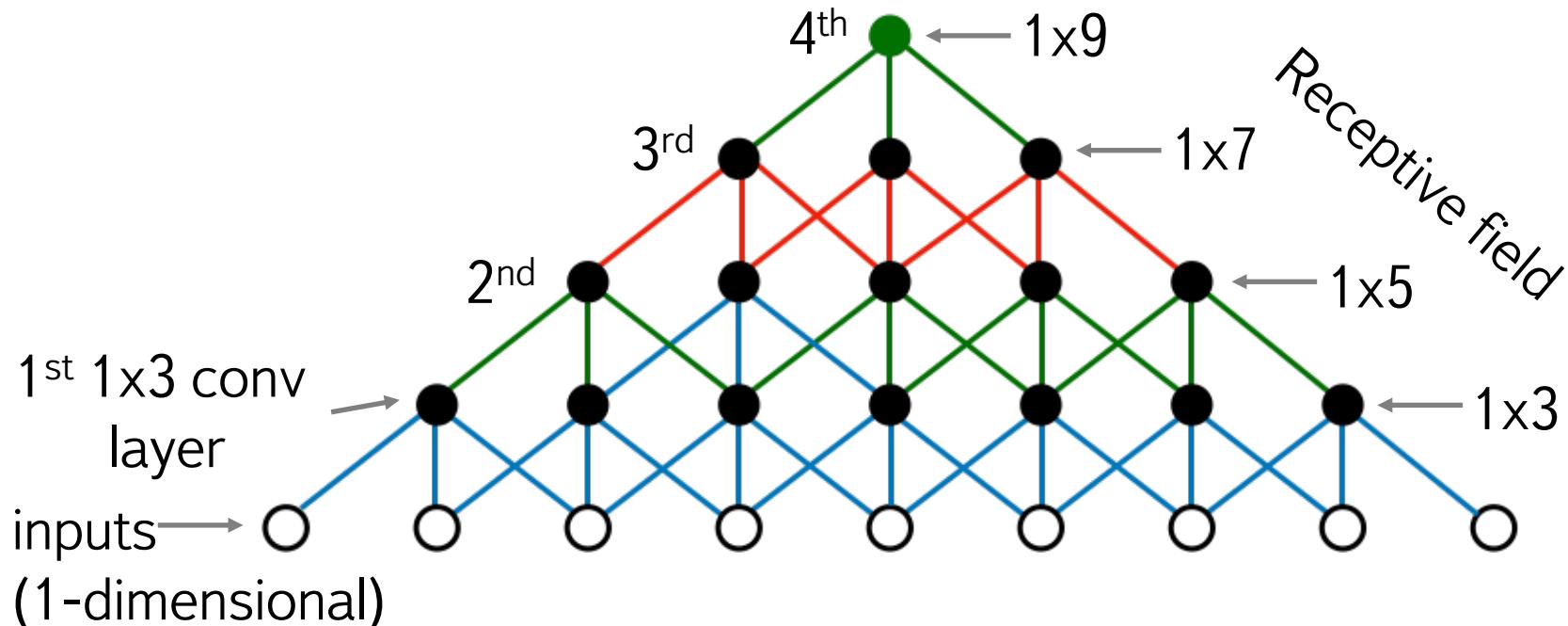
- Let's say neurons of the 1st convolutional layer look at the patches of the image of size 3x3.
- What if an object of interest is bigger than that?
- We need a 2<sup>nd</sup> convolutional layer on top of the 1<sup>st</sup>!



# Receptive field after N convolutional layers



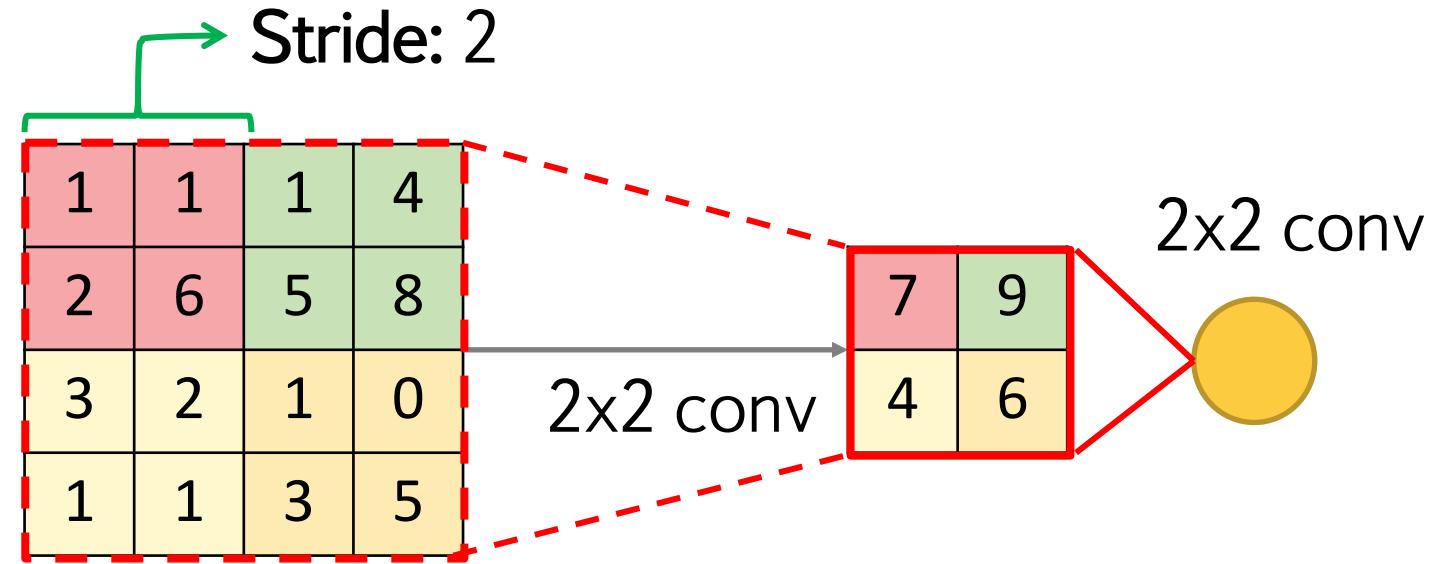
# Receptive field after $N$ convolutional layers



- If we stack  $N$  convolutional layers with the same kernel size 3x3 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!

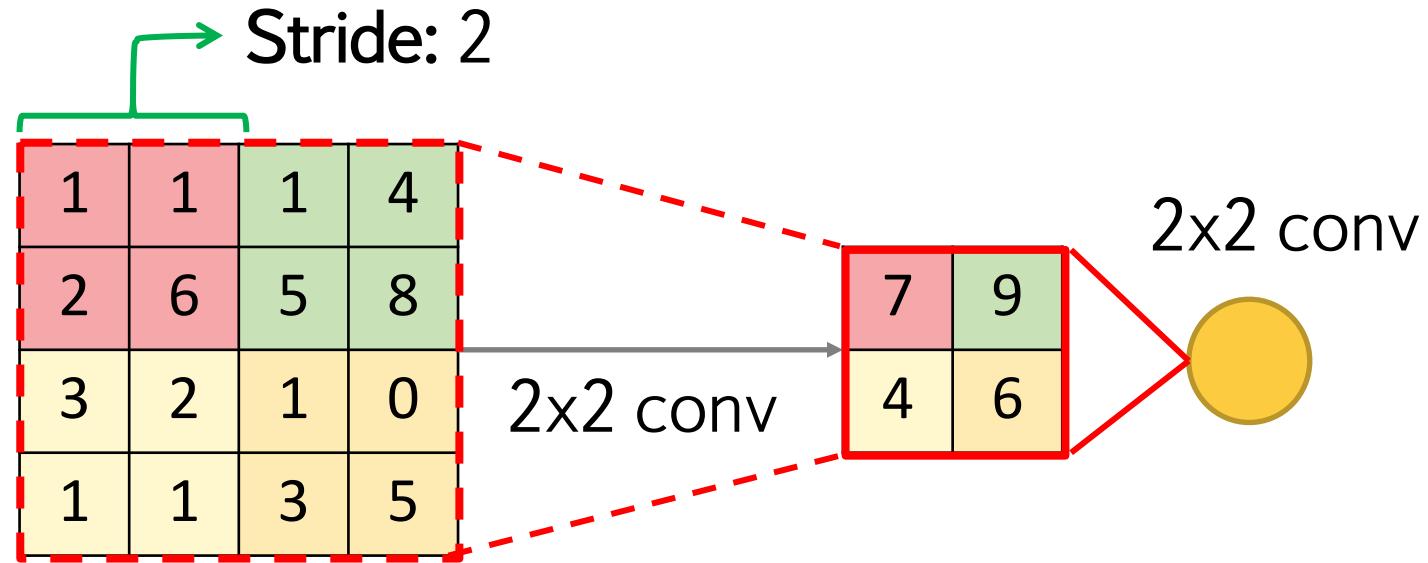
# We need to grow receptive field faster!

- 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

\*

=

0	0	0
0	1	0
0	0	2

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

\*

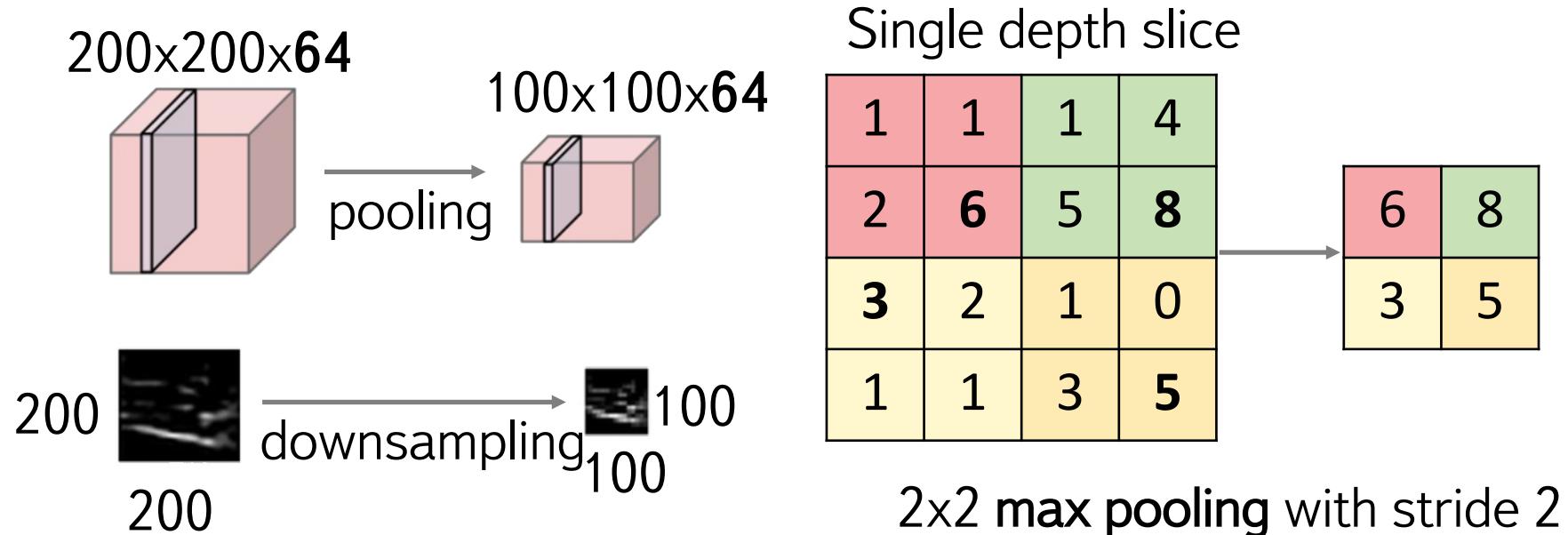
=

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.



# 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!

6	8
3	5

Maximum = 8

7	8
3	5

Maximum = 8

- 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!

6	8
3	5

Maximum = 8

7	8
3	5

Maximum = 8

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

6	8
3	5

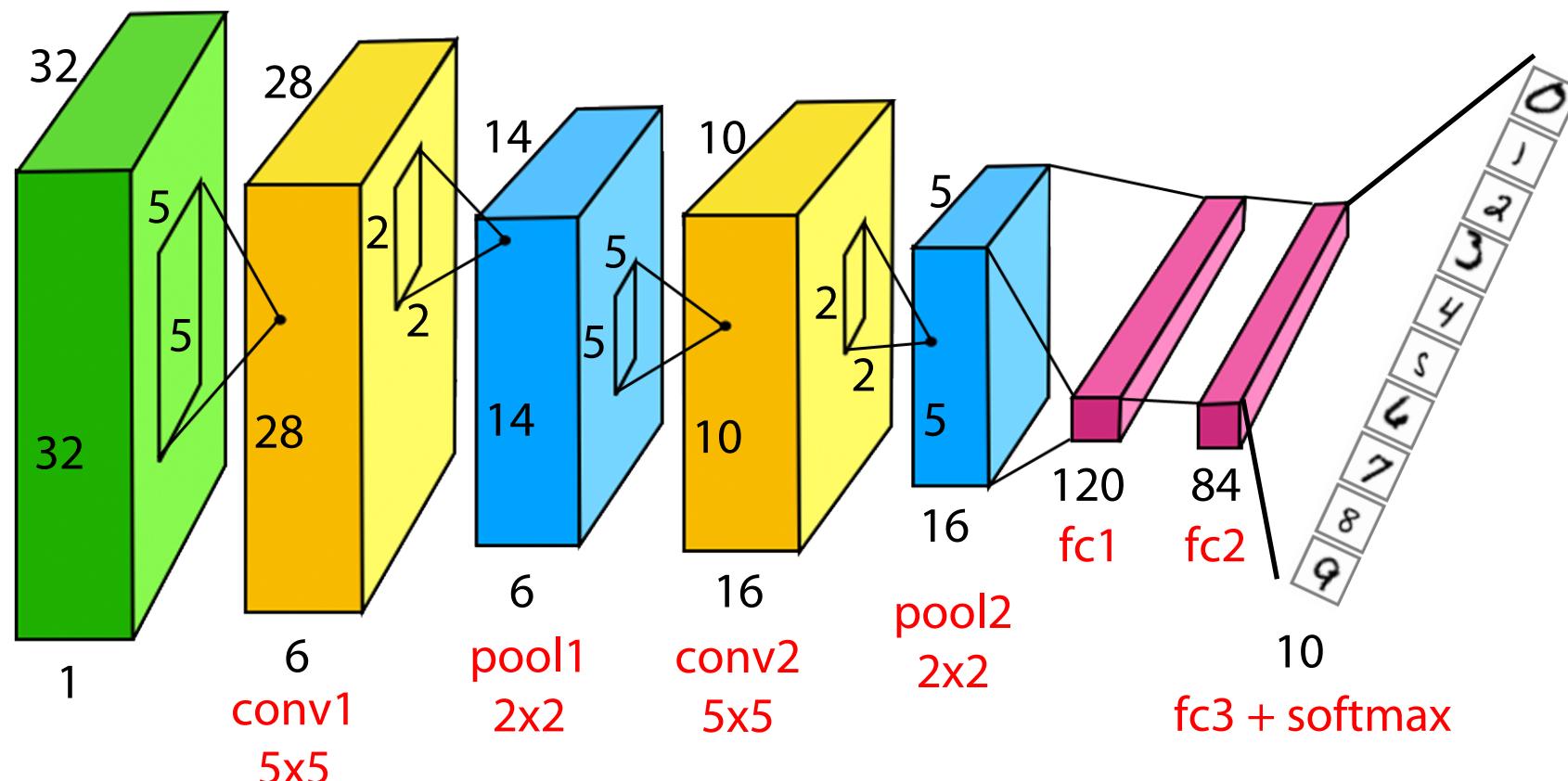
Maximum = 8

7	9
3	5

Maximum = 9

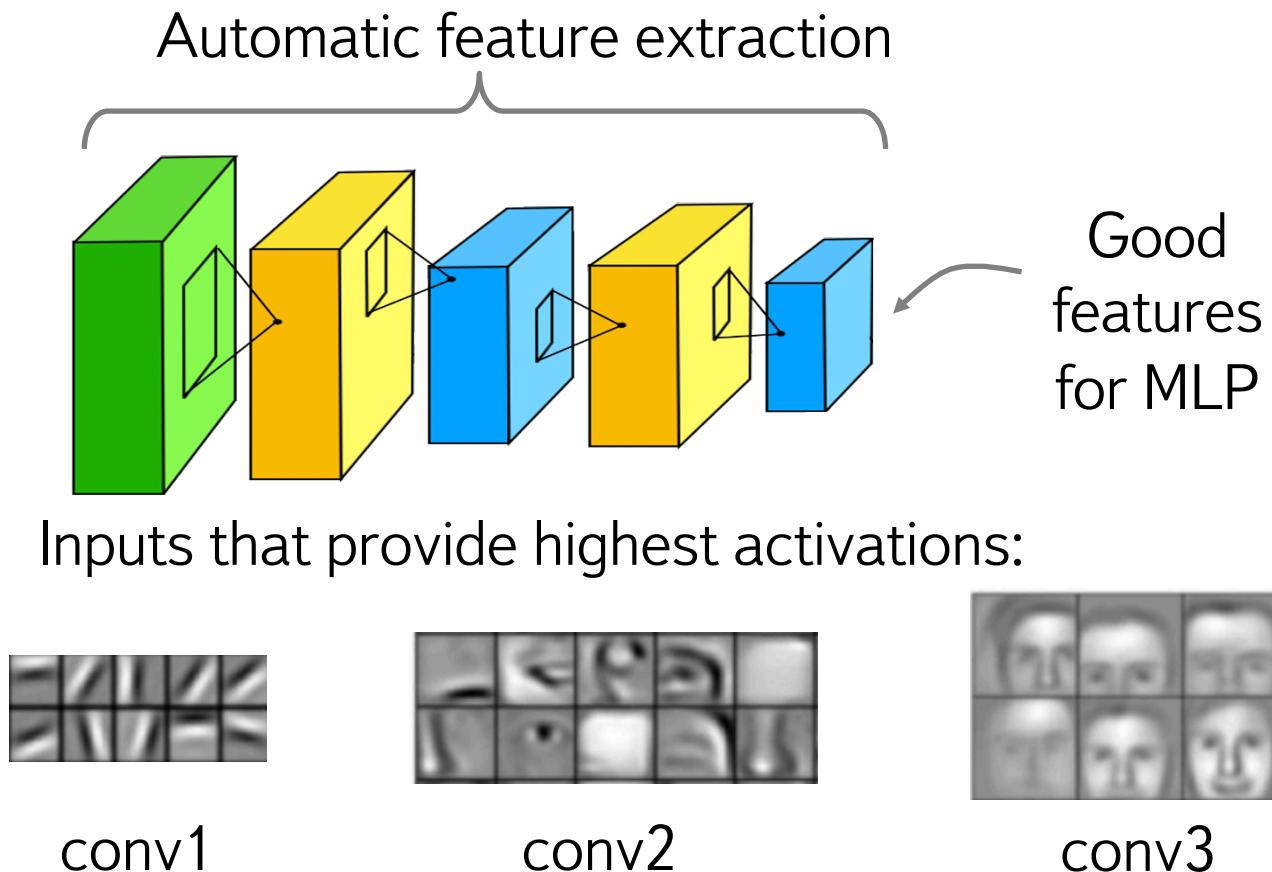
# Putting it all together into a simple CNN

- 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.



# Summary

- Using convolutional, pooling and fully connected layers we've built our first network for handwritten digits recognition!

# Neural architectures for computer vision

# ImageNet classification dataset



flamingo

cock

ruffed grouse

quail

partridge



Egyptian cat

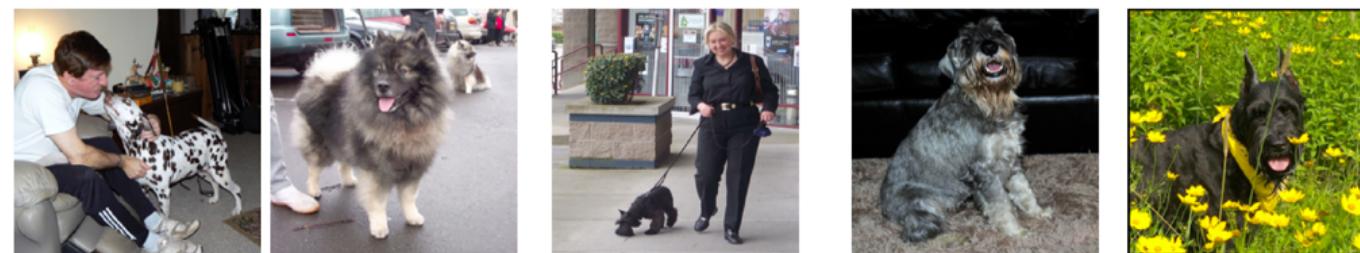
Persian cat

Siamese cat

tabby

lynx

...



dalmatian

keeshond

miniature schnauzer

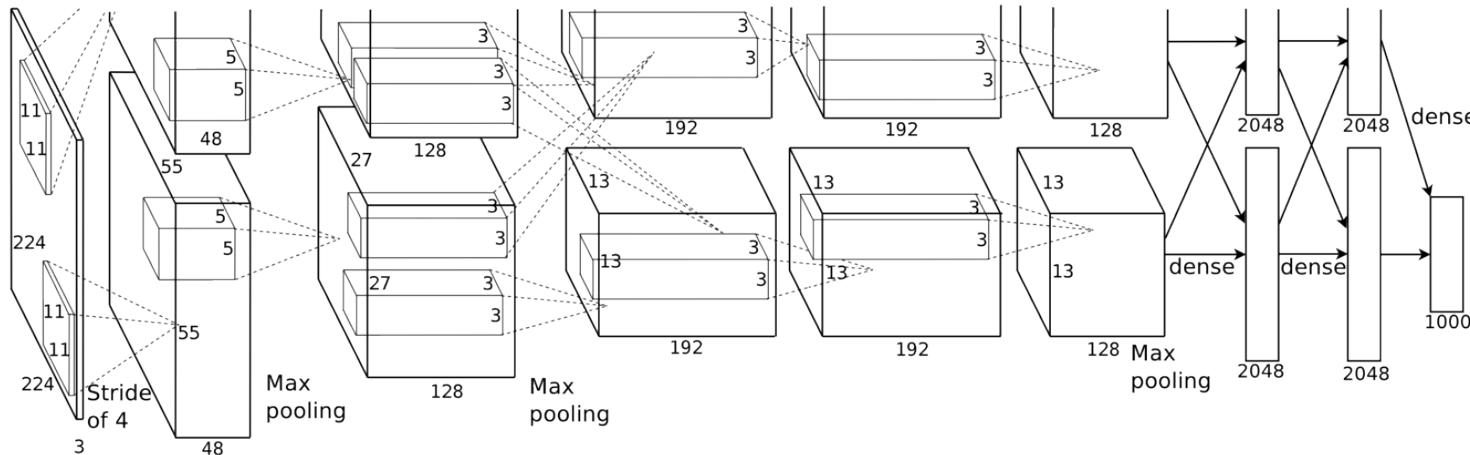
standard schnauzer

giant schnauzer

...

# AlexNet (2012)

- First deep convolutional neural net for ImageNet
- Significantly reduced top 5 error from 26% to 15%

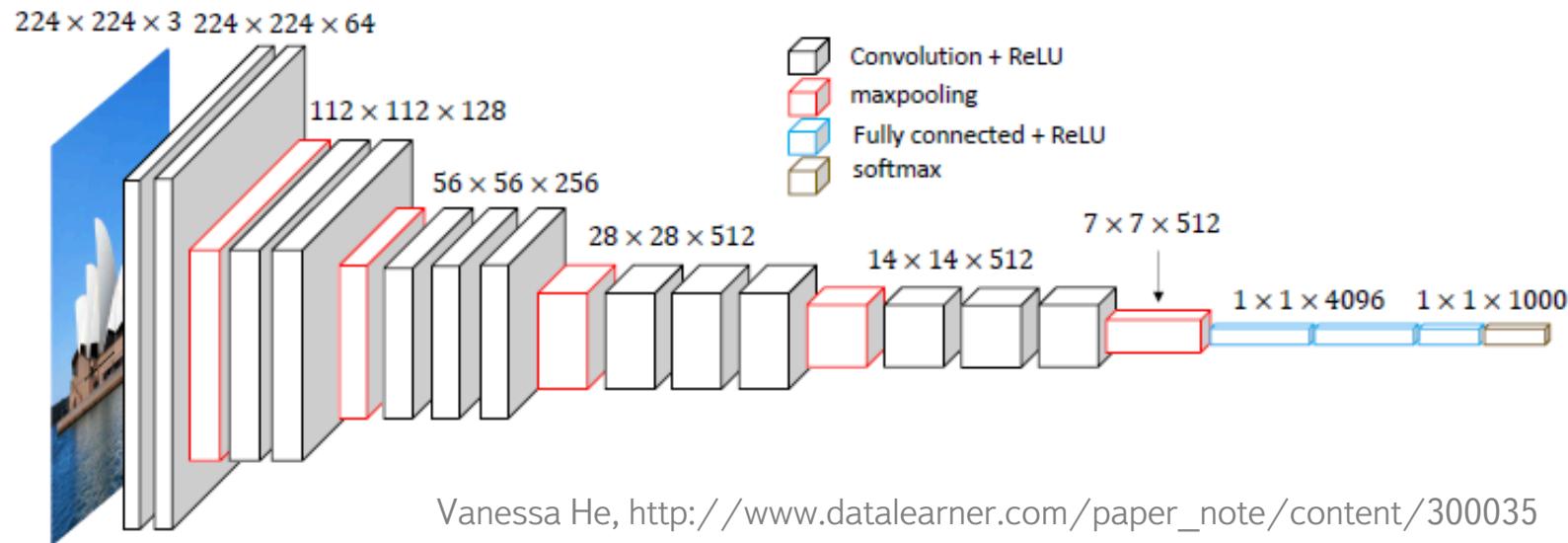


Alex Krizhevsky, <https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>

- 11x11, 5x5, 3x3 convolutions, max pooling, dropout, data augmentation, ReLU activations, SGD with momentum
- 60 million parameters
- Trains on 2 GPUs for 6 days

# VGG (2015)

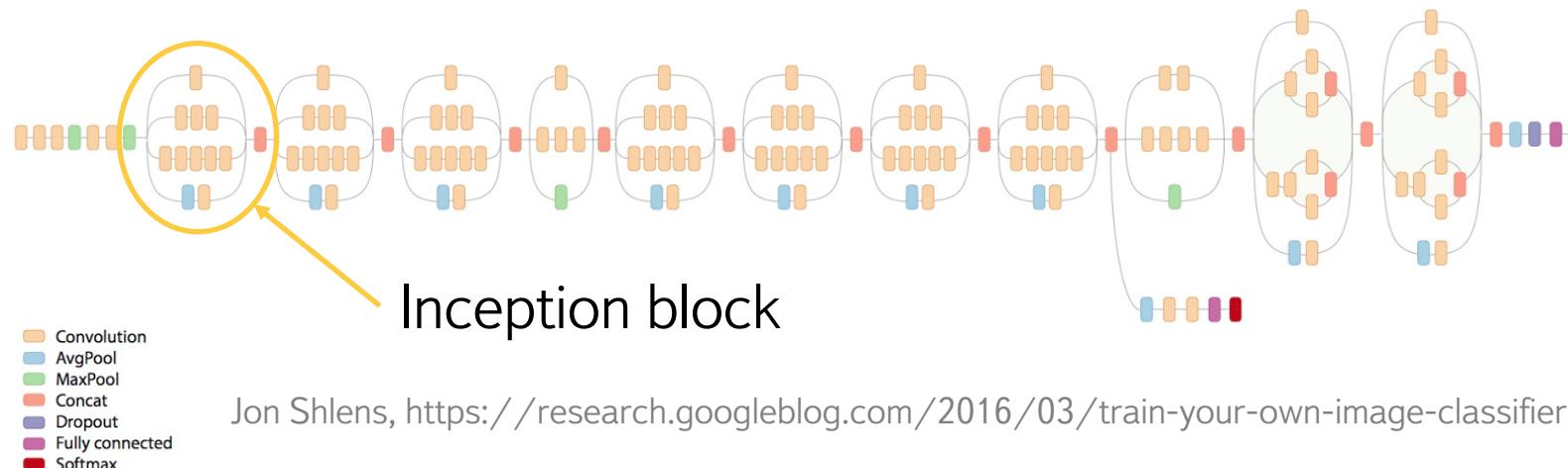
- Similar to AlexNet, only 3x3 convolutions, but lots of filters!
- ImageNet top 5 error: 8.0% (single model)



- Training similar to AlexNet with additional multi-scale cropping.
- 138 million parameters
- Trains on 4 GPUs for 2-3 weeks

# Inception V3 (2015)

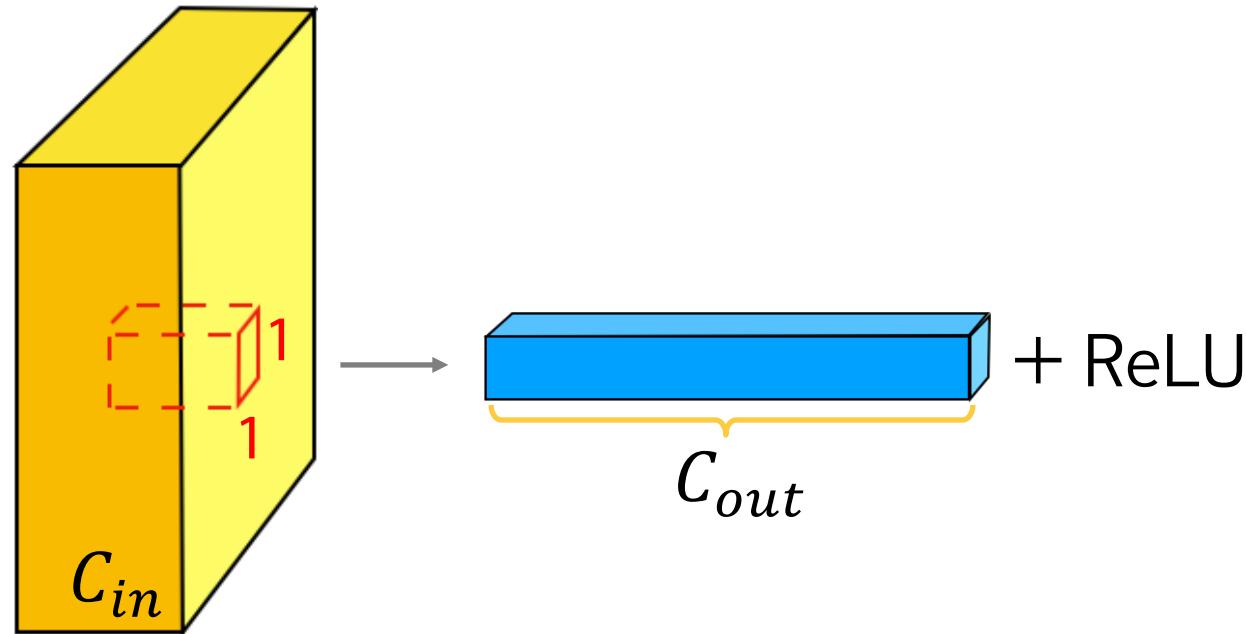
- Similar to AlexNet? Not quite, uses Inception block introduced in GoogLeNet (a.k.a. Inception V1)
- ImageNet top 5 error: 5.6% (single model), 3.6% (ensemble)



- Batch normalization, image distortions, RMSProp
- 25 million parameters!
- Trains on 8 GPUs for 2 weeks

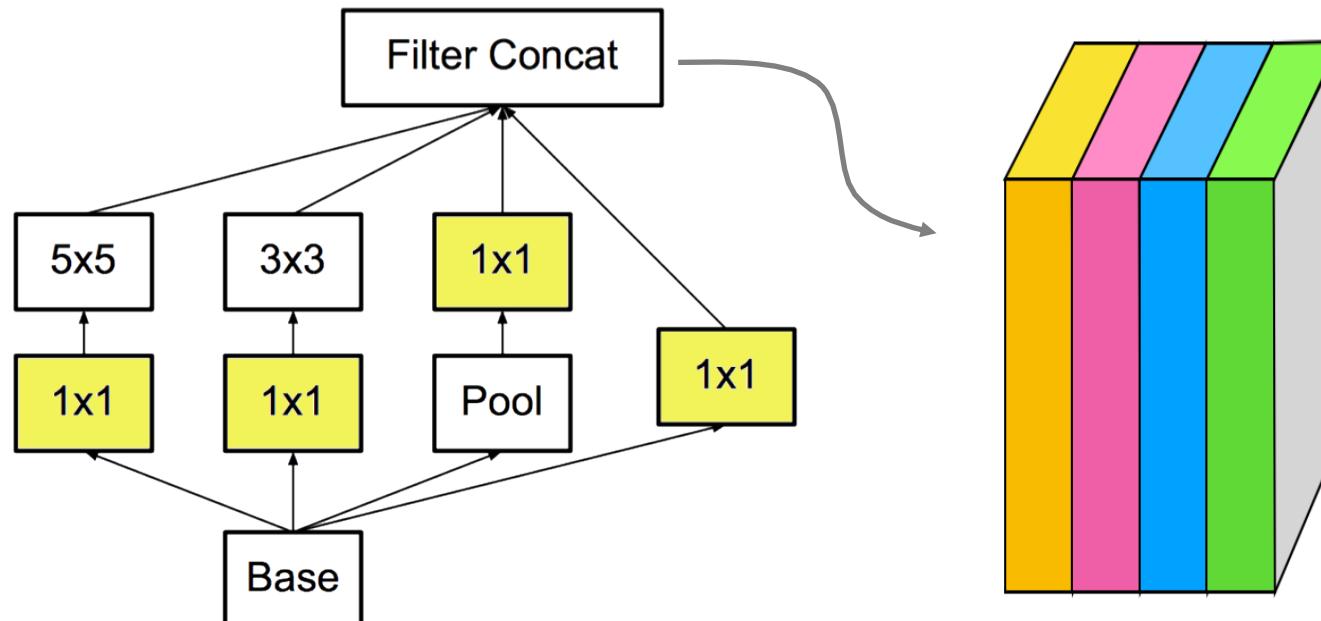
# 1x1 convolutions

- Such convolutions capture interactions of input channels in one “pixel” of feature map
- They can reduce the number of channels not hurting the quality of the model, because different channels can correlate
- Dimensionality reduction with added ReLU activation



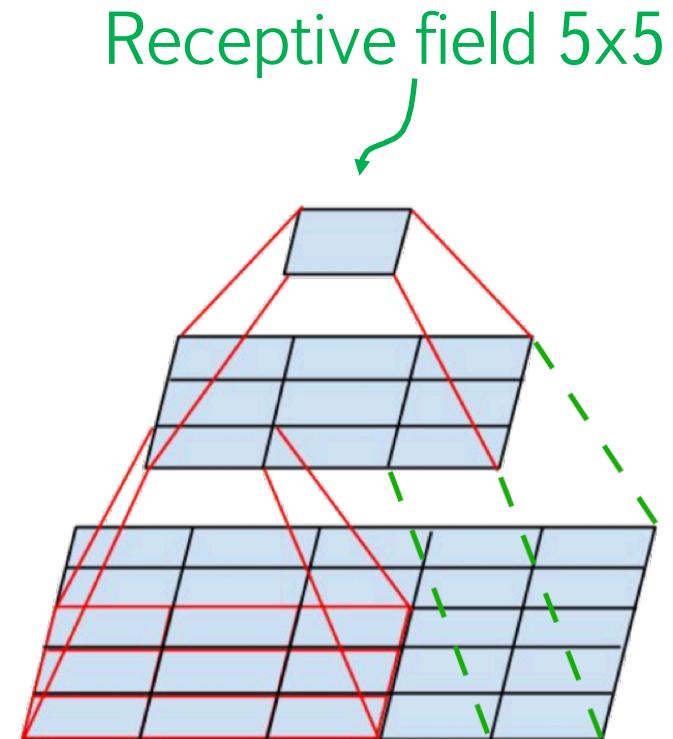
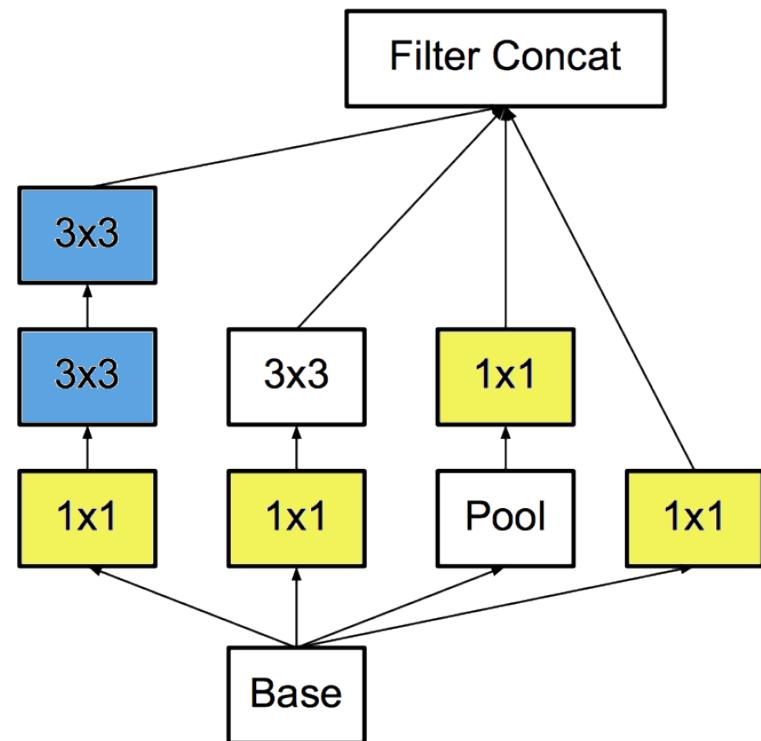
# Basic Inception block

- All operations inside a block use stride 1 and enough padding to output the same spatial dimensions ( $W \times H$ ) of feature map.
- 4 different feature maps are concatenated on depth at the end



# Replace 5x5 convolutions

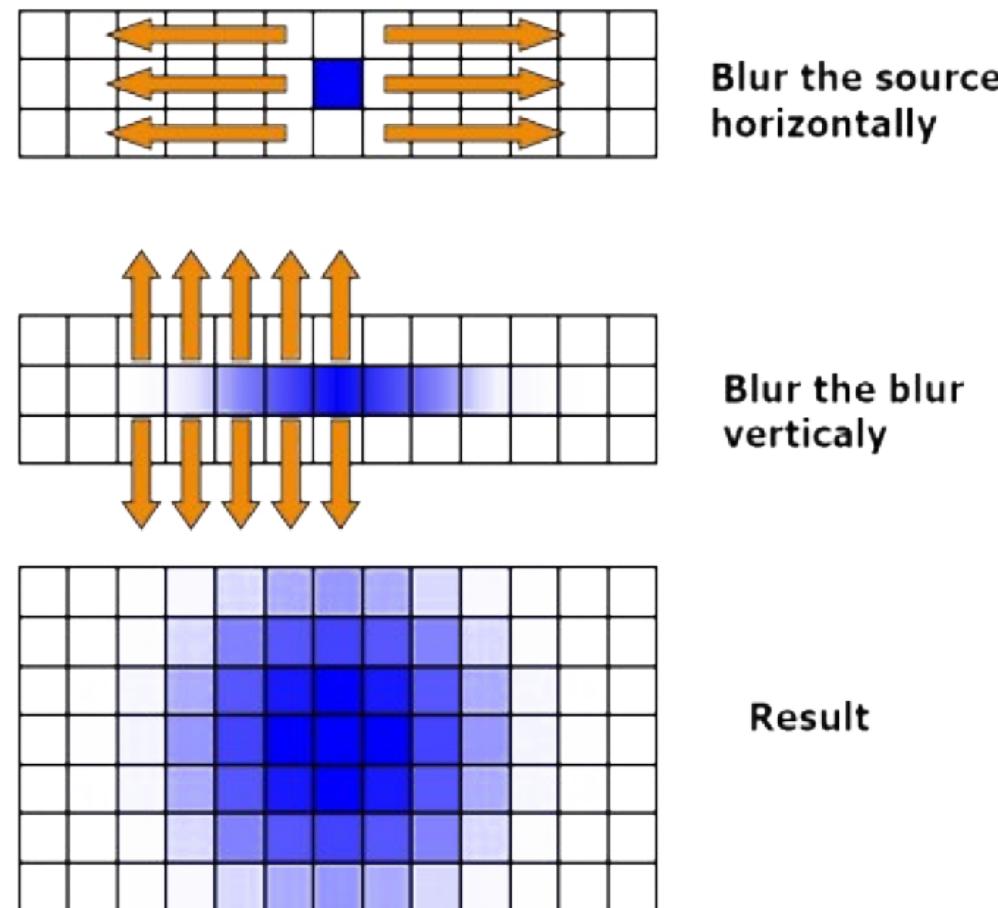
- 5x5 convolutions are expensive! Let's replace them with two layers of 3x3 convolutions which have an effective receptive field of 5x5.



Christian Szegedy, <https://arxiv.org/pdf/1512.00567.pdf>

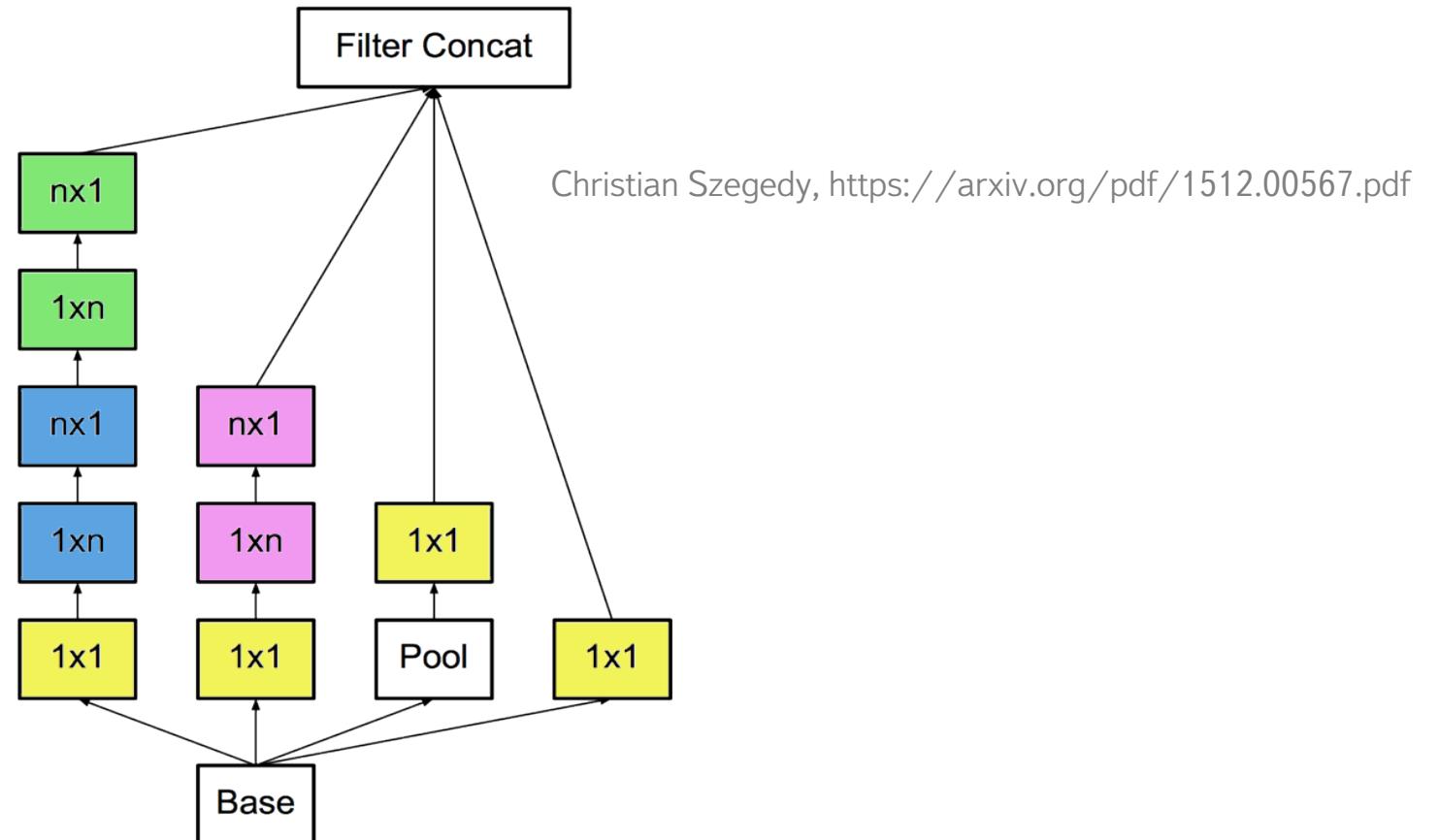
# Filter decomposition

- It's known that a Gaussian blur filter can be decomposed in two 1 dimensional filters:



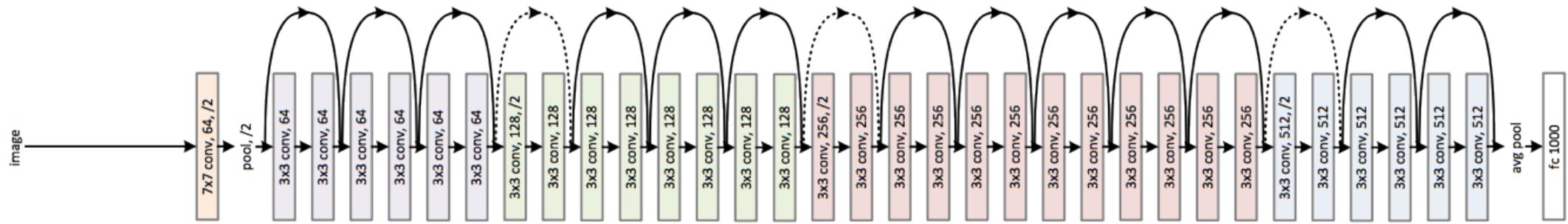
# Filter decomposition in Inception block

- 3x3 convolutions are currently the most expensive parts!
- Let's replace each 3x3 layer with 1x3 layer followed by 3x1 layer.



# ResNet (2015)

- Introduces residual connections
- ImageNet top 5 error: 4.5% (single model), 3.5% (ensemble)

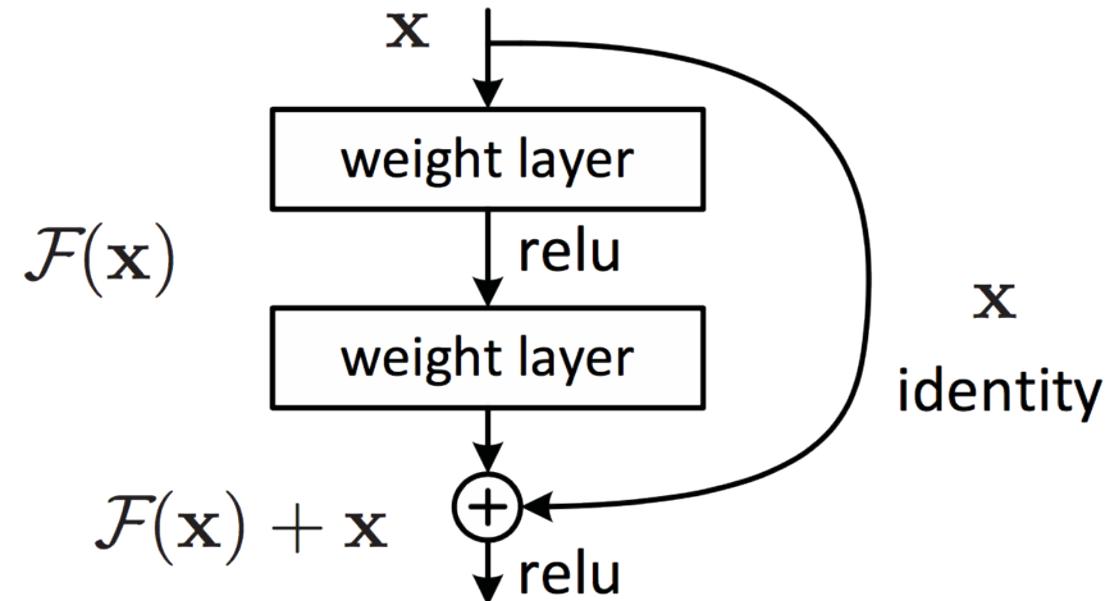


Kaiming He, <https://arxiv.org/pdf/1512.03385.pdf>

- 152 layers, few 7x7 convolutional layers, the rest are 3x3, batch normalization, max and average pooling.
- 60 million parameters
- Trains on 8 GPUs for 2-3 weeks.

# Residual connections

- We create output channels adding a small delta  $F(x)$  to original input channels  $x$ :



Kaiming He, <https://arxiv.org/pdf/1512.03385.pdf>

- This way we can stack thousands of layers and gradients do not vanish thanks to residual connections

# Summary

- By stacking more convolution and pooling layers you can reduce the error! Like in AlexNet or VGG.
- But you cannot do that forever, you need to utilize new kind of layers like Inception block or residual connections.
- You've probably noticed that one needs a lot of time to train her neural network!