# CNNs
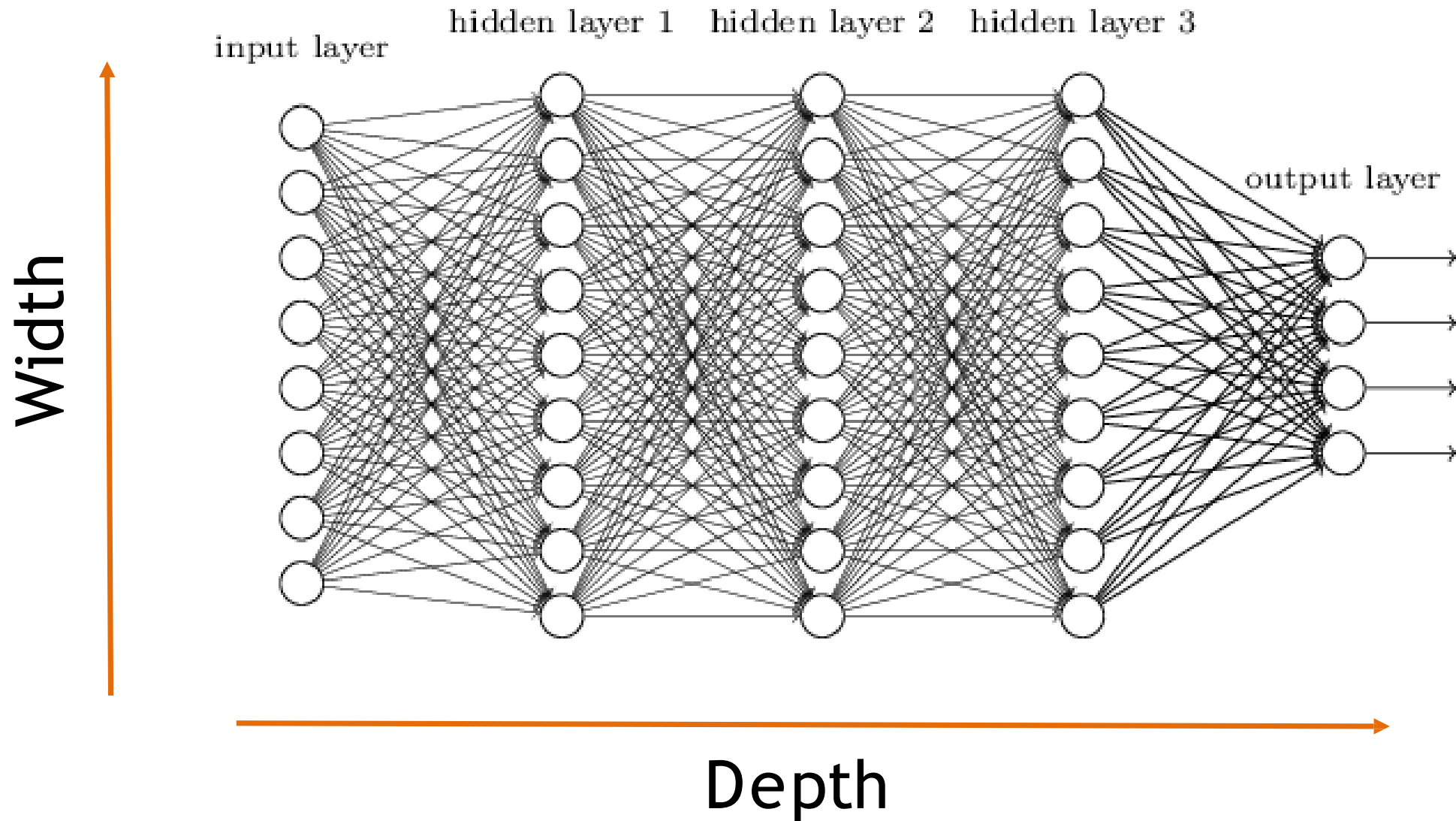
**Marcus Rüb**

**Hahn-Schickard Villingen-Schwenningen**
**Marcus.rueb@hahn-schickard.de**

# Fully Connected Neural Network
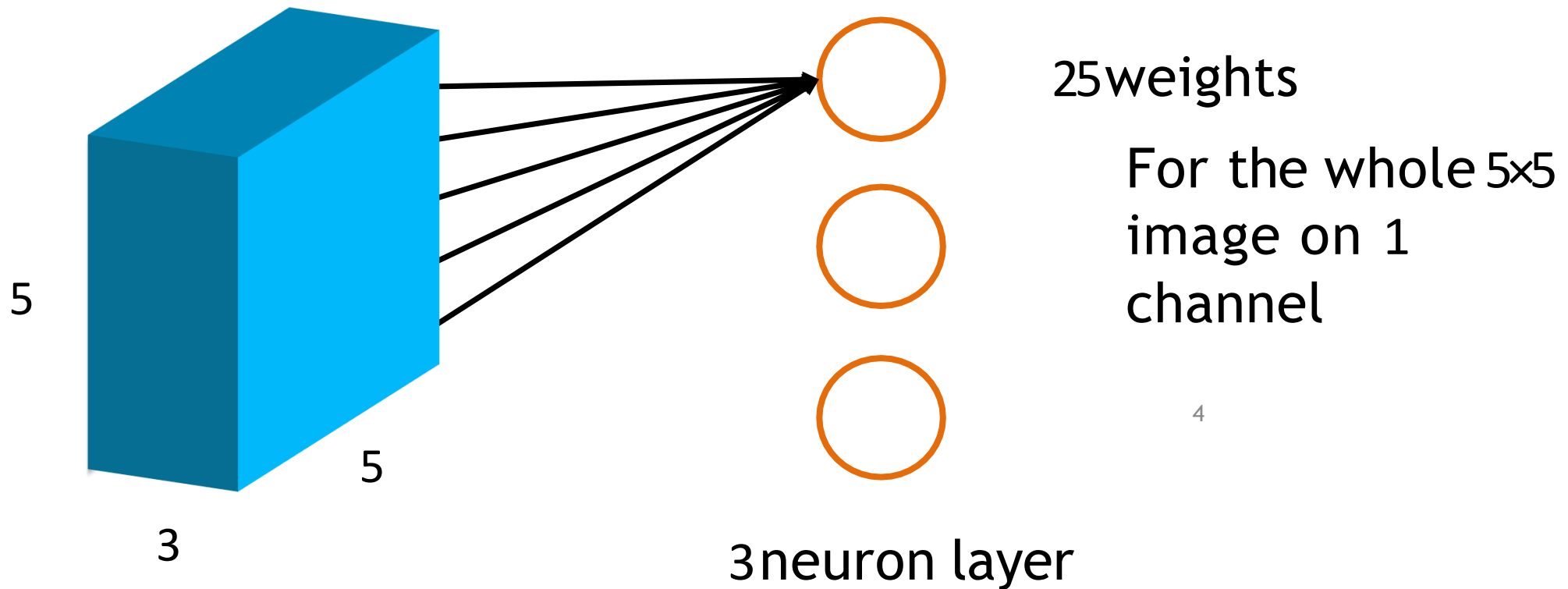
input layer

hidden layer 1    hidden layer 2    hidden layer 3

output layer

Width

Depth

# **Problems using FC Layers on Images**

- How to process a tiny image with FC layers



5

5

3

5weights

3

3neuron layer

# Problems using FC Layers on Images

- How to process a tiny image with FC layers



25 weights

For the whole 5×5 image on 1 channel

3 neuron layer

# Problems using FC Layers on Images

- How to process a tiny image with FC layers



75 weights

For the whole 5×5 image on the 3 channel

3 neuron layer

5

5

5

3

# Problems using FC Layers on Images

- How to process a tiny image with FC layers



75 weights

75 weights

75 weights

3 neuron layer

For the whole 5×5 image on the three channels per neuron

5

5

3

# Problems using FC Layers on Images

- How to process a normal image with FC layers



1000

1000

3

3neuron layer

# Problems using FC Layers on Images

- How to process a normal image with FC layers



1000

1000

3

$3\,billion$ weights

IMPRACTICAL

8

1000 neuron layer

# Why not simply more FC Layers?

We cannot make networks arbitrarily complex

- Why not just go deeper and get better?
    - No structure!!
    - It is just brute force!
    - Optimization becomes hard
    - Performance plateaus / drops!

- ## We want to restrict the degrees of freedom
  - We want a layer with structure
  - Weight sharing → using the same weights for different parts of the image
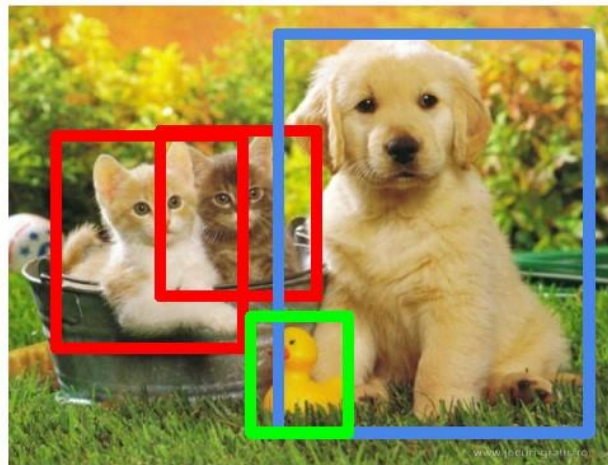
# Using CNNs in Computer Vision



| Classification | Classification + Localization | Object Detection | Instance Segmentation |
| --- | --- | --- | --- |
| CAT | CAT | CAT, DOG, DUCK | CAT, DOG, DUCK |

Single object

Multiple objects

[Li et al., CS231n Course Slides] Lecture 12: Detection and Segmentation

# Convolutions

# What are Convolutions?

$$f * g = \int^{\infty} f(\tau)g(t - \tau)d\tau$$

$f$ = red

$g$ = blue

$f * g$ = green

Convolution of two box functions

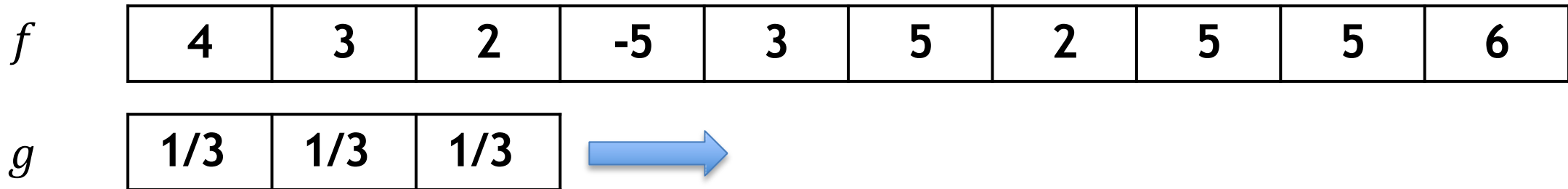Convolution of two Gaussians

Application of a filter to a function

— The 'smaller' one is typically called the filter kernel

# Discrete case: box filter

| $f$ | 4 | 3 | 2 | -5 | 3 | 5 | 2 | 5 | 5 | 6 |
|---|---|---|---|---|---|---|---|---|---|---|

| $g$ | 1/3 | 1/3 | 1/3 |
|---|---|---|---|

‘Slide’ filter kernel from left to right; at each position, compute a single value in the output data

# Discrete case: box filter

| f | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 4 | 3 | 2 | -5 | 3 | 5 | 2 | 5 | 5 | 6 |

| g | | | |
|---|---|---|---|
| | 1/3 | 1/3 | 1/3 |

| f *g | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 3 | | | | | | | | |

$$4 \cdot \frac{1}{3} + 3 \cdot \frac{1}{3} + 2 \cdot \frac{1}{3} = 3$$

# Discrete case: box filter

| $f$ | 4 | 3 | 2 | -5 | 3 | 5 | 2 | 5 | 5 | 6 |
|---|---|---|---|---|---|---|---|---|---|---|

| $g$ | 1/3 | 1/3 | 1/3 | |
|---|---|---|---|---|

| $f * g$ | | 3 | 0 | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|

$$3 \cdot \frac{1}{3} + 2 \cdot \frac{1}{3} + (-5) \cdot \frac{1}{3} = 0$$

16

# What areConvolutions?

## Discrete case: box filter

$f$

| 4 | 3 | 2 | -5 | 3 | 5 | 2 | 5 | 5 | 6 |
|---|---|---|----|---|---|---|---|---|---|

$g$

| 1/3 | 1/3 | 1/3 |
|-----|-----|-----|

$f * g$

|   | 3 | 0 | 0 |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|

$$2 \cdot \frac{1}{3} + (-5) \cdot \frac{1}{3} + 3 \cdot \frac{1}{3} = 0$$

# What areConvolutions?

## Discrete case: box filter

$f$

| 4 | 3 | 2 | -5 | 3 | 5 | 2 | 5 | 5 | 6 |
|---|---|---|---|---|---|---|---|---|---|

$g$

| 1/3 | 1/3 | 1/3 |
|-----|-----|-----|

$f*g$

| | 3 | 0 | 0 | 1 | | | | | |
|---|---|---|---|---|---|---|---|---|---|

$$(-5) \cdot \frac{1}{3} + 3 \cdot \frac{1}{3} + 5 \cdot \frac{1}{3} = 1$$

18

# Discrete case: box filter

$f$

| 4 | 3 | 2 | -5 | 3 | 5 | 2 | 5 | 5 | 6 |
|---|---|---|----|---|---|---|---|---|---|

$g$

| 1/3 | 1/3 | 1/3 |
|-----|-----|-----|

$f*g$

| | 3 | 0 | 0 | 1 | 10/3 | | | | |
|---|---|---|---|---|------|---|---|---|---|

$$3 \cdot \frac{1}{3} + 5 \cdot \frac{1}{3} + 2 \cdot \frac{1}{3} = \frac{10}{3}$$

19

# What areConvolutions?

## Discrete case: box filter

| $f$ | 4 | 3 | 2 | -5 | 3 | 5 | 2 | 5 | 5 | 6 |
|-----|---|---|---|----|---|---|---|---|---|---|

| $g$ | | | | | | 1/3 | 1/3 | 1/3 | | |
|-----|---|---|---|---|---|-----|-----|-----|---|---|

| $f*g$ | | 3 | 0 | 0 | 1 | 10/3 | 4 | | | |
|-------|---|---|---|---|---|------|---|---|---|---|

$$5 \cdot \frac{1}{3} + 2 \cdot \frac{1}{3} + 5 \cdot \frac{1}{3} = 4$$

# Discrete case: box filter

$f$

| 4 | 3 | 2 | -5 | 3 | 5 | 2 | 5 | 5 | 6 |
|---|---|---|----|---|---|---|---|---|---|

$g$

| 1/3 | 1/3 | 1/3 |
|-----|-----|-----|

$f * g$

| | 3 | 0 | 0 | 1 | 10/3 | 4 | 4 | | |
|---|---|---|---|---|------|---|---|---|---|

$$2 \cdot \frac{1}{3} + 5 \cdot \frac{1}{3} + 5 \cdot \frac{1}{3} = 4$$

21

**What areConvolutions?**

# Discrete case: box filter

| $f$ | | 4 | 3 | 2 | -5 | 3 | 5 | 2 | 5 | 5 | 6 |

| $g$ | | | | | | | | | 1/3 | 1/3 | 1/3 |

| $f * g$ | | | 3 | 0 | 0 | 1 | 10/3 | 4 | 4 | 16/3 | |

$$5 \cdot \frac{1}{3} + 5 \cdot \frac{1}{3} + 6 \cdot \frac{1}{3} = \frac{16}{3}$$

# Discrete case: box filter

| 4 | 3 | 2 | -5 | 3 | 5 | 2 | 5 | 5 | 6 |
|---|---|---|----|---|---|---|---|---|---|

| 1/3 | 1/3 | 1/3 |
|-----|-----|-----|

| ?? | 3 | 0 | 0 | 1 | 10/3 | 4 | 4 | 16/3 | ?? |
|----|---|---|---|---|------|---|---|------|----|

## What to do at boundaries?

# Discrete case: box filter

| 4 | 3 | 2 | -5 | 3 | 5 | 2 | 5 | 5 | 6 |
|---|---|---|----|---|---|---|---|---|---|

| 1/3 | 1/3 | 1/3 |
|-----|-----|-----|

| ?? | 3 | 0 | 0 | 1 | 10/3 | 4 | 4 | 16/3 | ?? |
|----|---|---|---|---|------|---|---|------|----|

## What to do at boundaries?

24

### Option 1: Shrink

| 3 | 0 | 0 | 1 | 10/3 | 4 | 4 | 16/3 |
|---|---|---|---|------|---|---|------|

# Discrete case: box filter

| 0 | 4 | 3 | 2 | -5 | 3 | 5 | 2 | 5 | 5 | 6 | 0 |
|---|---|---|---|----|---|---|---|---|---|---|---|

| 1/3 | 1/3 | 1/3 |
|-----|-----|-----|

| ?? | 3 | 0 | 0 | 1 | 10/3 | 4 | 4 | 16/3 | ?? |
|----|---|---|---|---|------|---|---|------|----|

## What to do at boundaries?

$$0 \cdot \frac{1}{3} + 4 \cdot \frac{1}{3} + 3 \cdot \frac{1}{3} = \frac{7}{3}$$

25

Option 2: Pad (often 0's)

| 7/3 | 3 | 0 | 0 | 1 | 10/3 | 4 | 4 | 16/3 | 11/3 |
|-----|---|---|---|---|------|---|---|------|------|

# Convolutions on Images

Image 5×5

| -5 | 3 | 2 | -5 | 3 |
|----|---|---|----|---|
| 4 | 3 | 2 | 1 | -3 |
| 1 | 0 | 3 | 3 | 5 |
| -2 | 0 | 1 | 4 | 4 |
| 5 | 6 | 7 | 9 | -1 |

Kernel 3×3

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

Output 3×3

| 6 | | |
|---|---|---|
| | | |
| | | |

$5·3+ (−1)·3+ (−1)·2+ (−1)·0+ (−1)·4$
$=15−9=6$

# Convolutions on Images

Image 5×5

| -5 | 3 | 2 | -5 | 3 |
|----|----|----|----|----|
| 4 | 3 | 2 | 1 | -3 |
| 1 | 0 | 3 | 3 | 5 |
| -2 | 0 | 1 | 4 | 4 |
| 5 | 6 | 7 | 9 | -1 |

Kernel 3×3

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

Output 3×3

| 6 | 1 | |
|----|----|----|
| | | |
| | | |

$5 \cdot 2 + (-1) \cdot 2 + (-1) \cdot 1 + (-1) \cdot 3 + (-1) \cdot 3 = 10 - 9 = 1$

# Convolutions on Images

Image 5×5

| -5 | 3 | 2 | -5 | 3 |
|----|---|---|----|---|
| 4 | 3 | 2 | 1 | -3 |
| 1 | 0 | 3 | 3 | 5 |
| -2 | 0 | 1 | 4 | 4 |
| 5 | 6 | 7 | 9 | -1 |

Kernel 3×3

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

Output 3×3

| 6 | 1 | 8 |
|---|---|---|
| | | |
| | | |

$5 \cdot 1 + (-1) \cdot (-5) + (-1) \cdot (-3) + (-1) \cdot 3$
$+ (-1) \cdot 2$
$= 5 + 3 = 8$

# Convolutions on Images

Image 5×5

| -5 | 3 | 2 | -5 | 3 |
|----|---|---|----|---|
| 4  | 3 | 2 | 1  | -3 |
| 1  | 0 | 3 | 3  | 5 |
| -2 | 0 | 1 | 4  | 4 |
| 5  | 6 | 7 | 9  | -1 |

Kernel 3×3

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

Output 3×3

| 6  | 1 | 8 |
|----|---|---|
| -7 |   |   |
|    |   |   |

$5 \cdot 0 + (-1) \cdot 3 + (-1) \cdot 0 + (-1) \cdot 1 + (-1) \cdot 3$
$= 0 - 7 = -7$

# Convolutions on Images

**Image 5×5**

| -5 | 3 | 2 | -5 | 3 |
|----|---|---|----|---|
| 4 | 3 | 2 | 1 | -3 |
| 1 | 0 | 3 | 3 | 5 |
| -2 | 0 | 1 | 4 | 4 |
| 5 | 6 | 7 | 9 | -1 |

**Kernel 3×3**

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

**Output 3×3**

| 6 | 1 | 8 |
|---|---|---|
| -7 | 9 | |
| | | |

$5 \cdot 3 + (-1) \cdot 2 + (-1) \cdot 3 + (-1) \cdot 1 + (-1) \cdot 0$
$= 15 - 6 = 9$

# Convolutions on Images

**Image 5×5**

| -5 | 3 | 2 | -5 | 3 |
|----|---|---|----|---|
| 4  | 3 | 2 | 1  | -3 |
| 1  | 0 | 3 | 3  | 5 |
| -2 | 0 | 1 | 4  | 4 |
| 5  | 6 | 7 | 9  | -1 |

**Kernel 3×3**

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

**Output 3×3**

| 6  | 1 | 8 |
|----|---|---|
| -7 | 9 | 2 |
|    |   |   |

$5 \cdot 3 + (-1) \cdot 1 + (-1) \cdot 5 + (-1) \cdot 4 + (-1) \cdot 3$
$= 15 - 13 = 2$

# Convolutions on Images

**Image 5×5**

| -5 | 3 | 2 | -5 | 3 |
|----|---|---|----|---|
| 4 | 3 | 2 | 1 | -3 |
| 1 | 0 | 3 | 3 | 5 |
| -2 | 0 | 1 | 4 | 4 |
| 5 | 6 | 7 | 9 | -1 |

**Kernel 3×3**

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

**Output 3×3**

| 6 | 1 | 8 |
|---|---|---|
| -7 | 9 | 2 |
| -5 | | |

$5 \cdot 0 + (-1) \cdot 0 + (-1) \cdot 1 + (-1) \cdot 6$
$+ (-1) \cdot (-2)$
$= -5$

# Convolutions on Images

**Image 5×5**

| -5 | 3 | 2 | -5 | 3 |
|----|---|---|----|---|
| 4 | 3 | 2 | 1 | -3 |
| 1 | 0 | 3 | 3 | 5 |
| -2 | 0 | 1 | 4 | 4 |
| 5 | 6 | 7 | 9 | -1 |

**Kernel 3×3**

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

**Output 3×3**

| 6 | 1 | 8 |
|---|---|---|
| -7 | 9 | 2 |
| -5 | -9 | |

$5 \cdot 1 + (-1) \cdot 3 + (-1) \cdot 4 + (-1) \cdot 7 + (-1) \cdot 0$
$= 5 - 14 = -9$

33

# Convolutions on Images

Image 5×5

| -5 | 3 | 2 | -5 | 3 |
|----|---|---|----|---|
| 4 | 3 | 2 | 1 | -3 |
| 1 | 0 | 3 | 3 | 5 |
| -2 | 0 | 1 | 4 | 4 |
| 5 | 6 | 7 | 9 | -1 |

Kernel 3×3

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

Output 3×3

| 6 | 1 | 8 |
|---|---|---|
| -7 | 9 | 2 |
| -5 | -9 | 3 |

$5 \cdot 4 + (-1) \cdot 3 + (-1) \cdot 4 + (-1) \cdot 9 + (-1) \cdot 1$
$= 20 - 17 = 3$

- Each kernel gives us a different image filter

Input

Edge detection

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

Box mean

$$\frac{1}{9}\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

Sharpen

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

Gaussian blur

$$\frac{1}{16}\begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$

LET'S LEARN THESE FILTERS!

# Convolutions on RGB Images

width  height  depth

image 32×32×3

Depth dimension *must* match; i.e., filter extends the full depth of the input

filter 5×5×3

32

5

5

3

Convolve filter with image
i.e., 'slide' over it and:
- apply filter at each location
- dot products

32

3

Images have depth: e.g. RGB -> 3 channels

# Convolutions on RGB Images

32×32×3 image (pixels $X$)

5×5×3 filter (weights vector $w$)



1 number at a time:
equal to dot product between filter weights $w$ and $x_i$−th chunk of the image. Here: 5·5·3=75-dim dot product +bias

$$z_4 = w^t x_i + b$$

(5×5×3 ×1)        (5×5×3)×1        1

# Convolutions on R G B Images

Activation map
(also feature map)

32×32×3 image

5×5×3 filter

32

5

5

3

Convolve

Slide over all spatial locations $x_i$
and compute all output $z_i$;
w/o padding, there are
28×28 locations

28

28

1

32

3

# Convolution Layer

# Convolution Layer

32×32×3 image

Activation maps

5×5×3 filter

32

5

32

5

3

Convolve

Let's apply a different filter with different weights!

28

28

1 1

# Convolution "Layer"

32×32×3 image

Activation maps

32

32

3

Convolve

28

28

5

Let's apply **five** filters, each with different weights!

# Convolution Layer

- A basic layer is defined by
  - Filter width and height (depth is implicitly given)
  - Number of different filter banks (#weight sets)


- Each filter captures a different image characteristic

# Different Filters



Low-Level Feature

- Each filter captures different image characteristics:
  - Horizontal edges
  - Vertical edges
  - Circles
  - Squares
  - ...

[Zeiler & Fergus, ECCV'14] Visualizing and Understanding Convolutional Networks

# Dimensions of a Convolution Layer

# Convolution Layers: Dimensions

1

Image 7×7

Input:       7×7
Filter:      3×3
Output:      5×5

# Convolution Layers: Dimensions

2

Image 7×7

Input:          7×7
Filter:         3×3
Output:         5×5

# Convolution Layers: Dimensions

| Input: | 7×7 |
|---|---|
| Filter: | 3×3 |
| Output: | 5×5 |

# Convolution Layers: Dimensions

4

Image 7×7

Input:          7×7
Filter:         3×3
Output:        5×5

# Convolution Layers: Dimensions

5

Image 7×7

Input:          7×7
Filter:         3×3
Output:         5×5

# Convolution Layers: Stride

With a stride of 1

Image 7×7

| | | | | | | |
|---|---|---|---|---|---|---|
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |

Input:        7×7
Filter:       3×3
Stride:        1
Output:       5×5

Stride of $S$: apply filter every $S$-th spatial location; i.e. subsample the image

# Convolution Layers: Stride

With a stride of 2

Image 7×7

Input:          7×7
Filter:         3×3
Stride:          2
Output:         3×3

# Convolution Layers: Stride

With a stride of 2

Image 7×7

| | | | | | | |
|---|---|---|---|---|---|---|
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |

Input:      7×7
Filter:     3×3
Stride:       2
Output:     3×3

# Convolution Layers: Stride

With a stride of 2



Image 7×7

Input:          7×7
Filter:         3×3
Stride:          2
Output:         3×3

# Convolution Layers: Stride

With a stride of 3

| Input: | 7×7 |
|---|---|
| Filter: | 3×3 |
| Stride: | 3 |
| Output: | ?×? |

Image 7×7

# Convolution Layers: Stride

With a stride of 3

Image 7×7



Input:          7×7
Filter:         3×3
Stride:          3
Output:         ?×?

# Convolution Layers: Stride

With a stride of 3

Image 7×7

Input:          7×7
Filter:         3×3
Stride:          3
Output:         ?×?

Does not really fit (remainder left)
→ Illegal stride for input & filter size!

# Convolution Layers: Dimensions

Input width of

$N$

Input height of

$N$

Filter width of $F$

Filter height of $F$

Input:    $N \times N$
Filter:    $F \times F$
Stride:    $S$

Output:    $\left(\frac{N-F}{S}+1\right) \times \left(\frac{N-F}{S}+1\right)$

$N = 7, F = 3, S = 1: \dfrac{7-3}{1}+1=5$

$N = 7, F = 3, S = 2: \dfrac{7-3}{2}+1=3$

$N = 7, F = 3, S = 3: \dfrac{7-3}{3}+1=2.\bar{3}$

Fractions are illegal

# Convolution Layers: Dimensions

Input Image



Conv + ReLU

32

5 filters

5×5×3

32

3

28

Conv + ReLU

8 filters

5×5×5

28

5

24

Conv + ReLU

12 filters

5×5×8

24

8

20

12

Shrinking down so quickly (32→28→24→20) is typically not a good idea...

## Why padding?

- Sizes get small too quickly
- Corner pixel is only used once

Image 7x7+zero padding

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|
| 0 | | | | | | | | 0 |
| 0 | | | | | | | | 0 |
| 0 | | | | | | | | 0 |
| 0 | | | | | | | | 0 |
| 0 | | | | | | | | 0 |
| 0 | | | | | | | | 0 |
| 0 | | | | | | | | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

# Why padding?

- Sizes get small too quickly
- Corner pixel is only used once

# Convolution Layers: Padding

Image 7×7 + zero padding

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Input ($N \times N$):     7×7
Filter ($F \times F$):     3×3
Padding ($P$):     1
Stride ($S$):     1
Output     7×7

Most common is 'zero' padding

Output Size:

$$\left( \left\lfloor \frac{N + 2 \cdot P - F}{S} \right\rfloor + 1 \right) \times \left( \left\lfloor \frac{N + 2 \cdot P - F}{S} \right\rfloor + 1 \right)$$

$\lfloor \ \rfloor$ denotes the floor operator (as in practice an integer division is performed)

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Image 7×7+zero padding

DL: Prof. Niessner, Prof. Leal-Taixé

## Types of convolutions:

- Valid convolution: using no padding

- Same convolution: output=input size

Set padding to $P = \frac{F-1}{2}$

# Convolution Layers: Dimensions

## Example

Input image: 32×32×3
10 filters 5×5
Stride 1
Pad 2          Depth of 3 is implicitly given

Output size is:

$$\frac{32+2\cdot2-5}{1}+1=32$$

i.e. 32×32×10

32

10 filters
5×5×3

32

3

Remember

Output: $\left(\left\lfloor\frac{N+2\cdot P-F}{S}\right\rfloor+1\right)\times\left(\left\lfloor\frac{N+2\cdot P-F}{S}\right\rfloor+1\right)$

# Convolution Layers: Dimensions

## Example

Input image: 32×32×3

10 filters 5×5

Stride 1

Pad 2

Output size is:

$$\frac{32 + 2 \cdot 2 - 5}{1} + 1 = 32$$

i.e. 32×32×10

32

10 filters
5×5×3

32

3

**Remember**

Output: $\left( \left\lfloor \frac{N + 2 \cdot P - F}{S} \right\rfloor + 1 \right) \times \left( \left\lfloor \frac{N + 2 \cdot P - F}{S} \right\rfloor + 1 \right)$

# Example

Input image: 32×32×3
10 filters 5 × 5
Stride 1
Pad 2



32
10 filters
5×5×3
32
3

Number of parameters (weights):
Each filter has 5 × 5 × 3 + 1 = 76 params
-> 76 · 10 = 760 parameters in layer

(+1 for bias)

# Convolutional Neural Network (CNN)

# ConvNet is concatenation of Conv Layers and activations

Input Image

Conv + ReLU

Conv + ReLU

Conv + ReLU

32
28
24
20

5 filters
8 filters
12 filters

5×5×3
5×5×5
5×5×8

32
28
24
12

3
5
8

# CNN Learned Filters



Low-Level Feature → Mid-Level Feature → High-Level Feature → Trainable Classifier

[Zeiler & Fergus, ECCV'14] Visualizing and Understanding Convolutional Networks

# Pooling

# Pooling Layer



224x224x64

pool →

112x112x64

224

224

downsampling

112

112

Lecture 5: Convolutional Neural Networks

70

# Pooling Layer: MaxPooling

Single depth slice of input

| | | | |
|---|---|---|---|
| 3 | 1 | 3 | 5 |
| 6 | 0 | 7 | 9 |
| 3 | 2 | 1 | 4 |
| 0 | 2 | 4 | 3 |

Max pool with
2×2 filters and stride 2

➡

'Pooled' output

| | |
|---|---|
| 6 | 9 |
| 3 | 4 |

- Conv Layer = 'Feature Extraction'
  - Computes a feature in a given region

- Pooling Layer = 'Feature Selection'
  - Picks the strongest activation in a region

- Input is a volume of size $W_{in} \times H_{in} \times D_{in}$
- Two hyperparameters
  - Spatial filter extent $F$  ⎱ Filter count $K$ and padding $P$
  - Stride $S$                  ⎰ make no sense here
- Output volume is of size $W_{out} \times H_{out} \times D_{out}$

  - $W_{out} = \dfrac{W_{in} - F}{S} + 1$

  - $H_{out} = \dfrac{H_{in} - F}{S} + 1$

  - $D_{out} = D_{in}$

- Does not contain parameters; e.g. it's fixed function

# Pooling Layer

- Input is a volume of size $W_{in} \times H_{in} \times D_{in}$
- Two hyperparameters
    - Spatial filter extent $F$
    - Stride $S$

Common settings:
$$F = 2, S = 2$$
$$F = 3, S = 2$$

- Output volume is of size $W_{out} \times H_{out} \times D_{out}$

    - $W_{out} = \dfrac{W_{in} - F}{S} + 1$

    - $H_{out} = \dfrac{H_{in} - F}{S} + 1$

    - $D_{out} = D_{in}$

- Does not contain parameters; e.g. it's fixed function

# Pooling Layer: Average Pooling

Single depth slice of input

| 3 | 1 | 3 | 5 |
|---|---|---|---|
| 6 | 0 | 7 | 9 |
| 3 | 2 | 1 | 4 |
| 0 | 2 | 4 | 3 |

Average pool with
2×2 filters and stride 2

'Pooled' output

| 2.5 | 6 |
|---|---|
| 1.75 | 3 |

- Typically used deeper in the network

# C N N Prototype



[Li et al., CS231n Course Slides]
Lecture 5: Convolutional Neural Networks

# Final Fully-Connected Layer

- Same as what we had in 'ordinary' neural networks
  – Make the final decision with the extracted features from the convolutions
  – One or two FC layers typically

# **Convolutions vs Fully-Connected**

- In contrast to fully-connected layers, we want to restrict the degrees of freedom
  - FC is somewhat brute force
  - Convolutions are structured

- Sliding window to with the same filter parameters to extract image features
  - Concept of weight sharing
  - Extract same features independent of location

# Receptive field

- Spatial extent of the connectivity of a convolutional filter

5x5 input

3x3 filter

=

3x3 output

- Spatial extent of the connectivity of a convolutional filter



5x5 input

3x3 filter

=

3x3 output

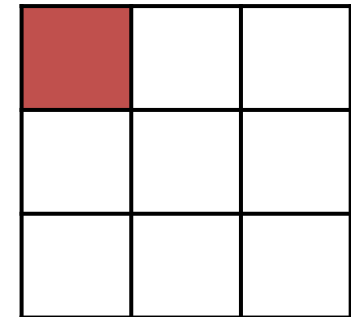3x3 receptive field = 1 output pixel is connected to 9 input pixels

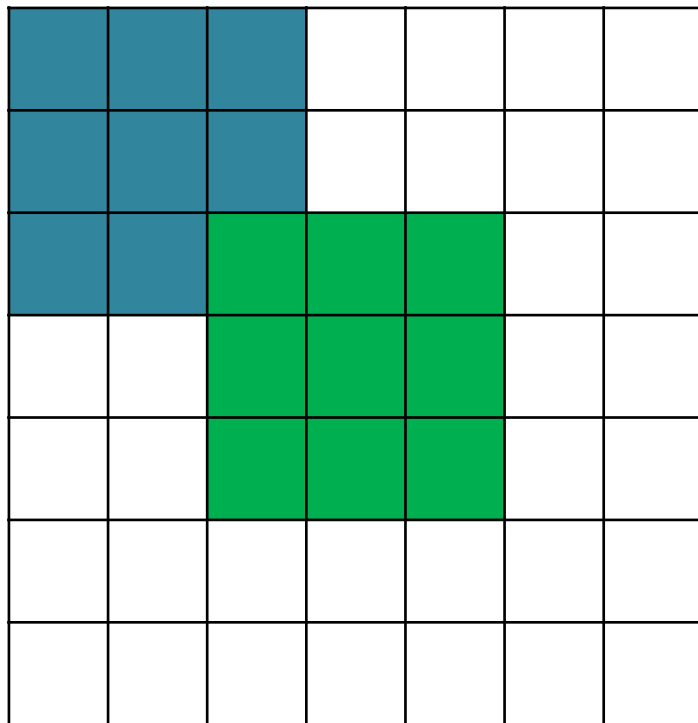- **Spatial extent of the connectivity of aconvolutional filter**



7x7 input

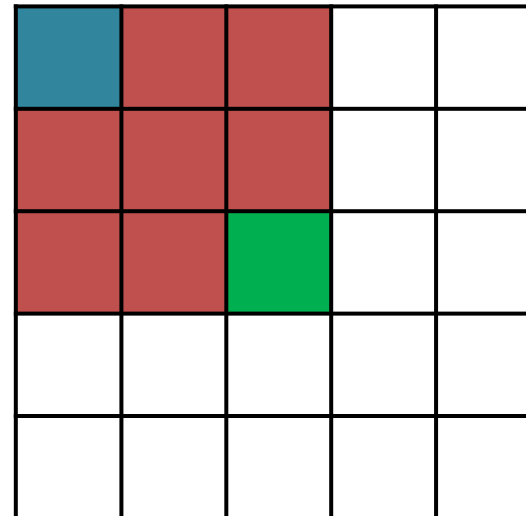3x3 receptive field = 1output pixel is connected to 9 input pixels
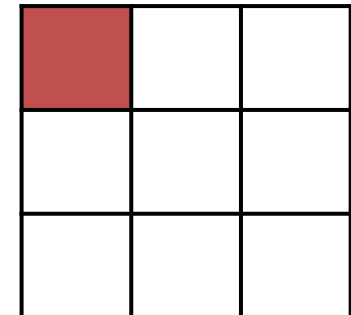
# Receptive Field

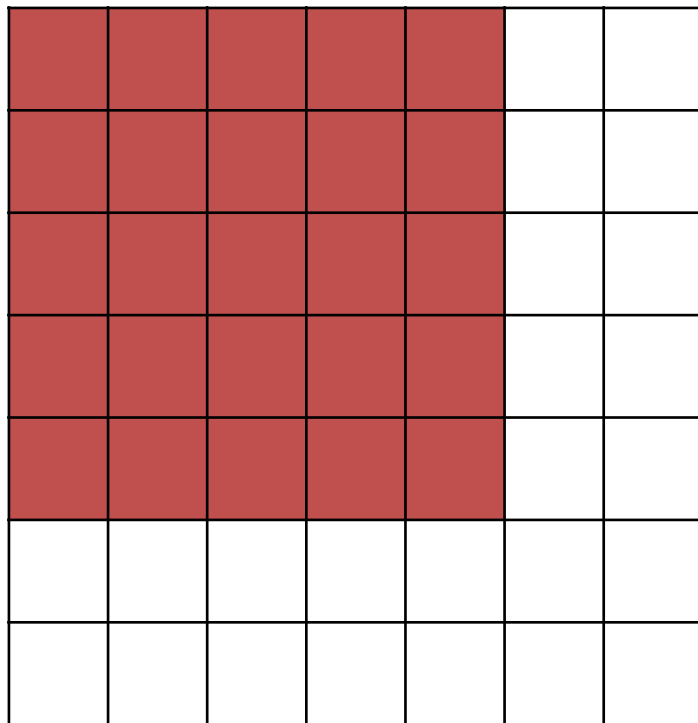- **Spatial extent of the connectivity of a convolutional filter**



7x7 input

3x3 receptive field = 1 output pixel is connected to 9 input pixels

# Receptive Field

- **Spatial extent of the connectivity of a convolutional filter**

7x7 input

3x3 receptive field = 1 output pixel is connected to 9 input pixels
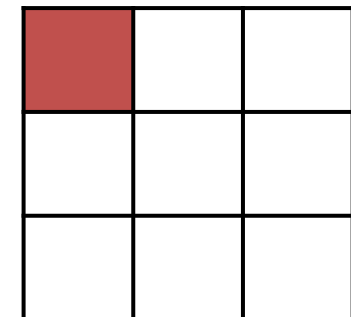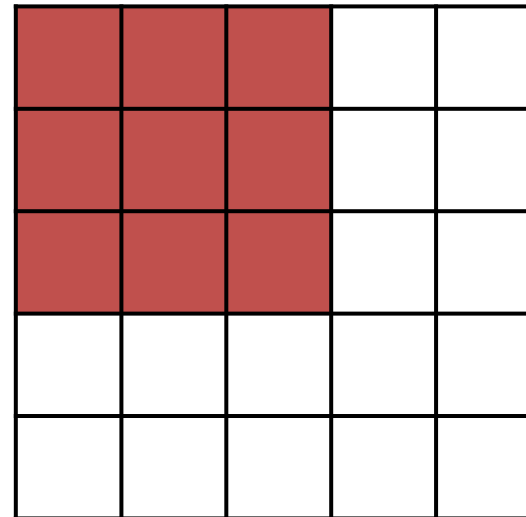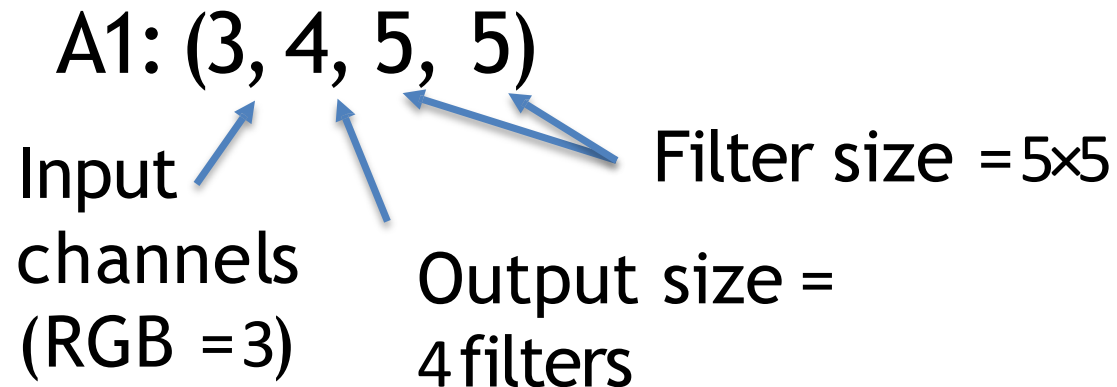
- **Spatial extent of the connectivity of aconvolutional filter**

7x7 input

5x5 receptive field on the original input:
one output value is connected to 25 input pixels

# Hands-on

**https://colab.research.google.com/github/tensorflow/docs/blob/master/site/en/tutorials/keras/classification.ipynb#scrollTo=2tRmdq_8CaXb**

- You are given a convolutional layer with 4 filters, kernel size 5, stride 1, and no padding that operates on an RGB image.

- Q1: What are the dimensions and the shape of its weight tensor?

  ☐ A1: (3, 4, 5,5)
  ☐ A2: (4, 5,5)
  ☐ A3: depends on the width and height of the image

- You are given a convolutional layer with 4 filters, kernel size 5, stride 1, and no padding that operates on an RGB image.

- Q1: What are the dimensions and the shape of its weight tensor?

A1: (3, 4, 5, 5)

Input channels (RGB = 3)

Output size = 4 filters

Filter size = 5×5

- You are given a convolutional layer with 4 filters, kernel size 5, stride 1, and no padding that operates on an RGB image.

- Q1: What are the dimensions and the shape of its weight tensor?

  ☐ A1: (3, 4, 5, 5)
  ☐ A2: (4, 5, 5)
  ☐ A3: depends on the width and height of the image