# Convolutional Neural Networks
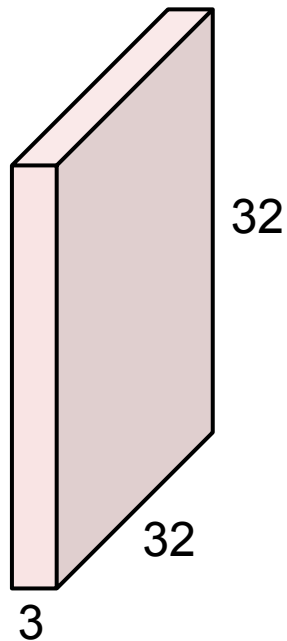
# Convolution Layer

32x32x3 image -> preserve spatial structure

32 height

32 width

3 depth

# Convolution Layer

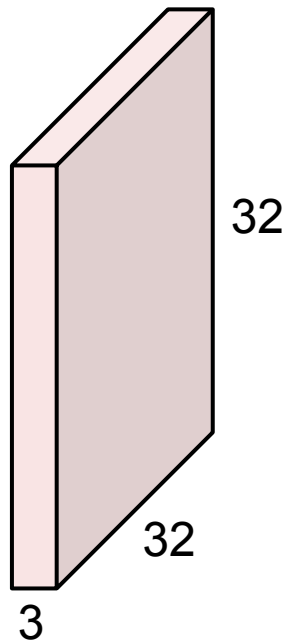32x32x3 image



32

32

3

5x5x3 filter



**Convolve** the filter with the image
i.e. "slide over the image spatially,
computing dot products"

# Convolution Layer

32x32x**3** image

**Filters always extend the full depth of the input volume**

32

32

**3**

5x5x**3** filter

**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"
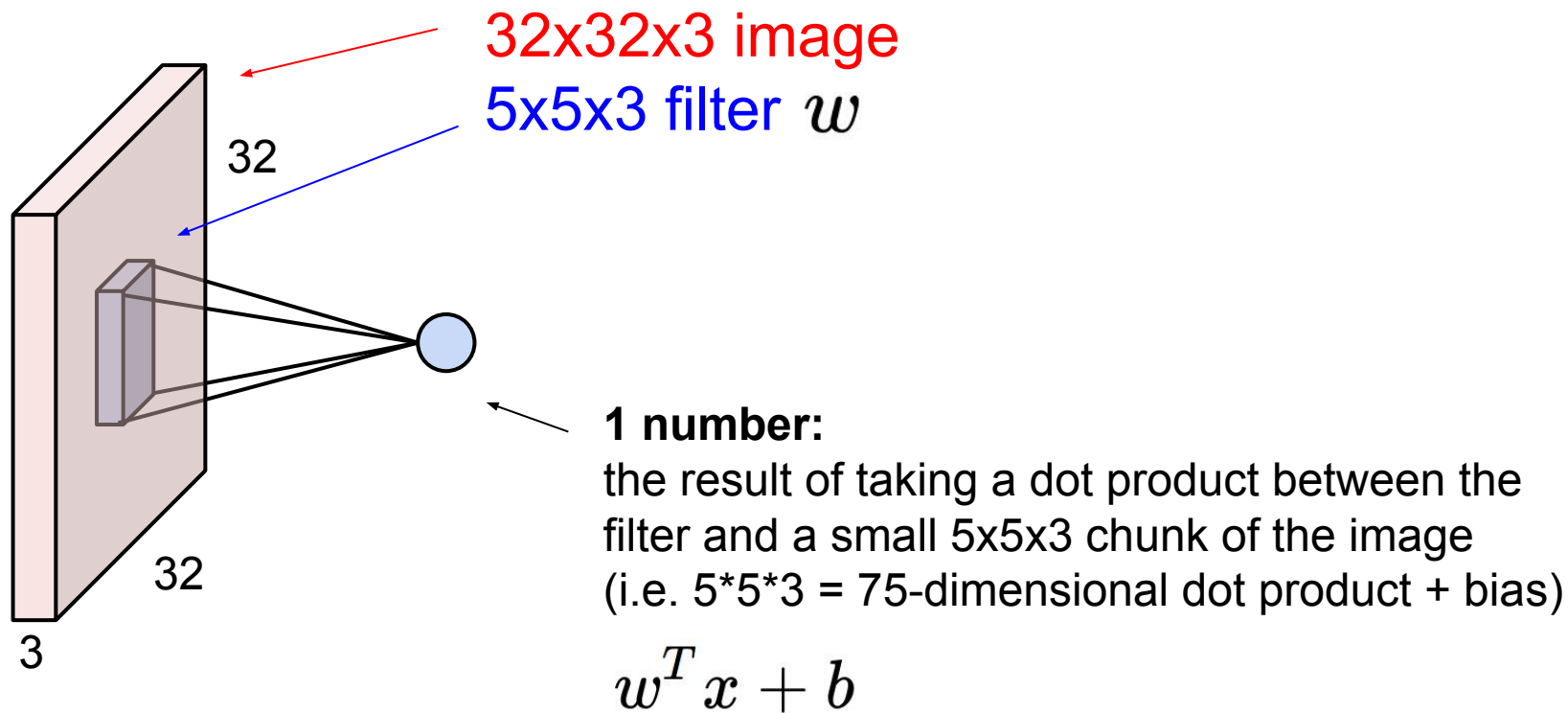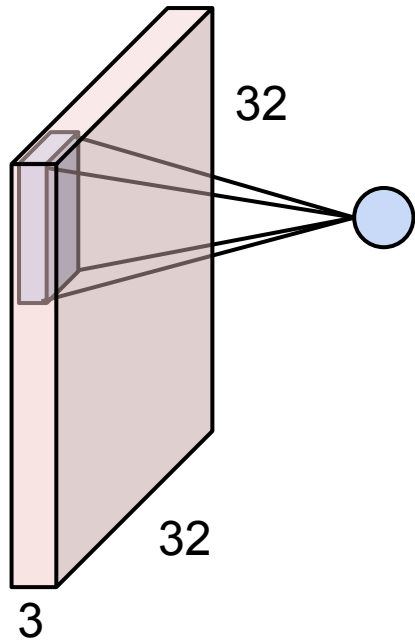
# Convolution Layer

32x32x3 image

5x5x3 filter $w$

32

32

3

**1 number:**
the result of taking a dot product between the filter and a small 5x5x3 chunk of the image (i.e. 5*5*3 = 75-dimensional dot product + bias)
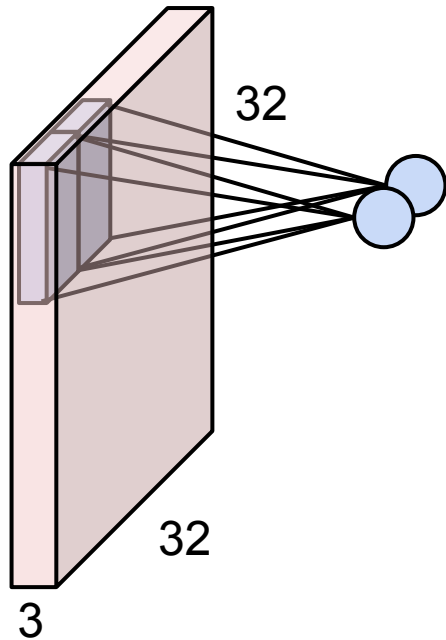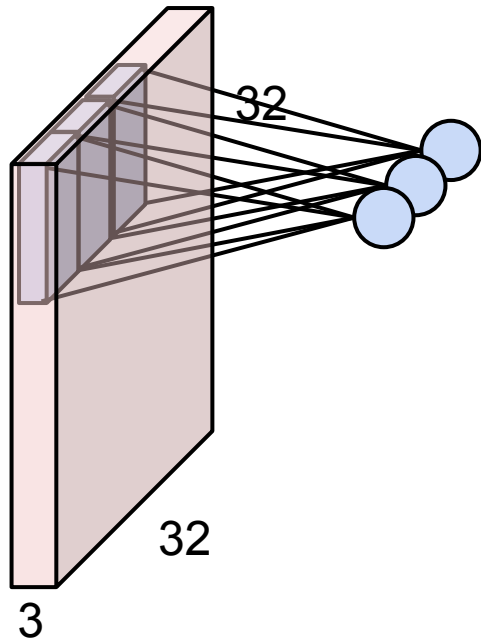
$$w^T x + b$$
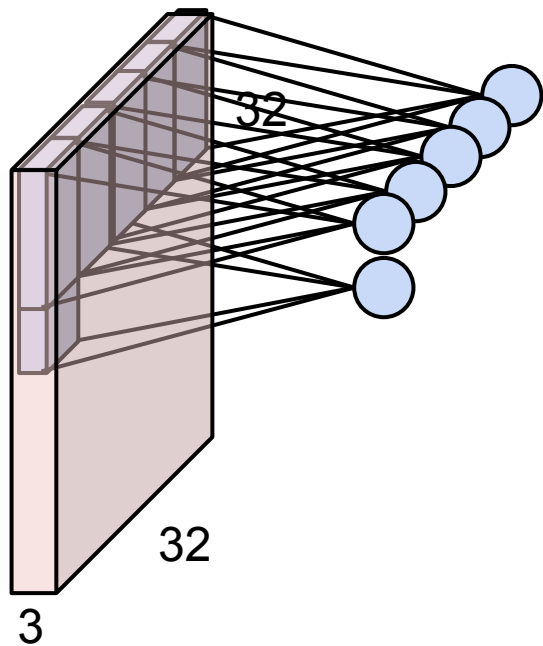
# Convolution Layer



32

32

3

# Convolution Layer



32

32

3

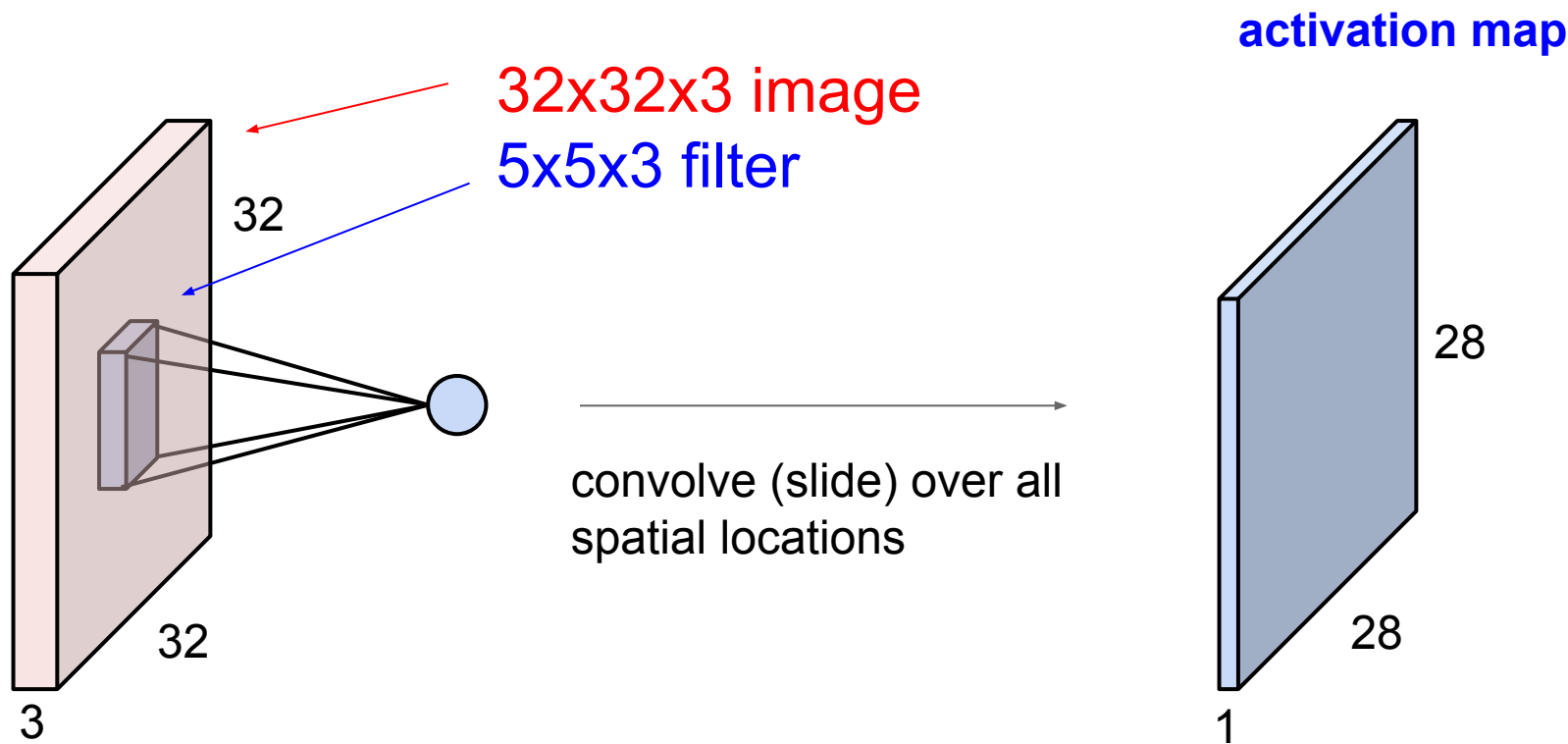# Convolution Layer

# Convolution Layer
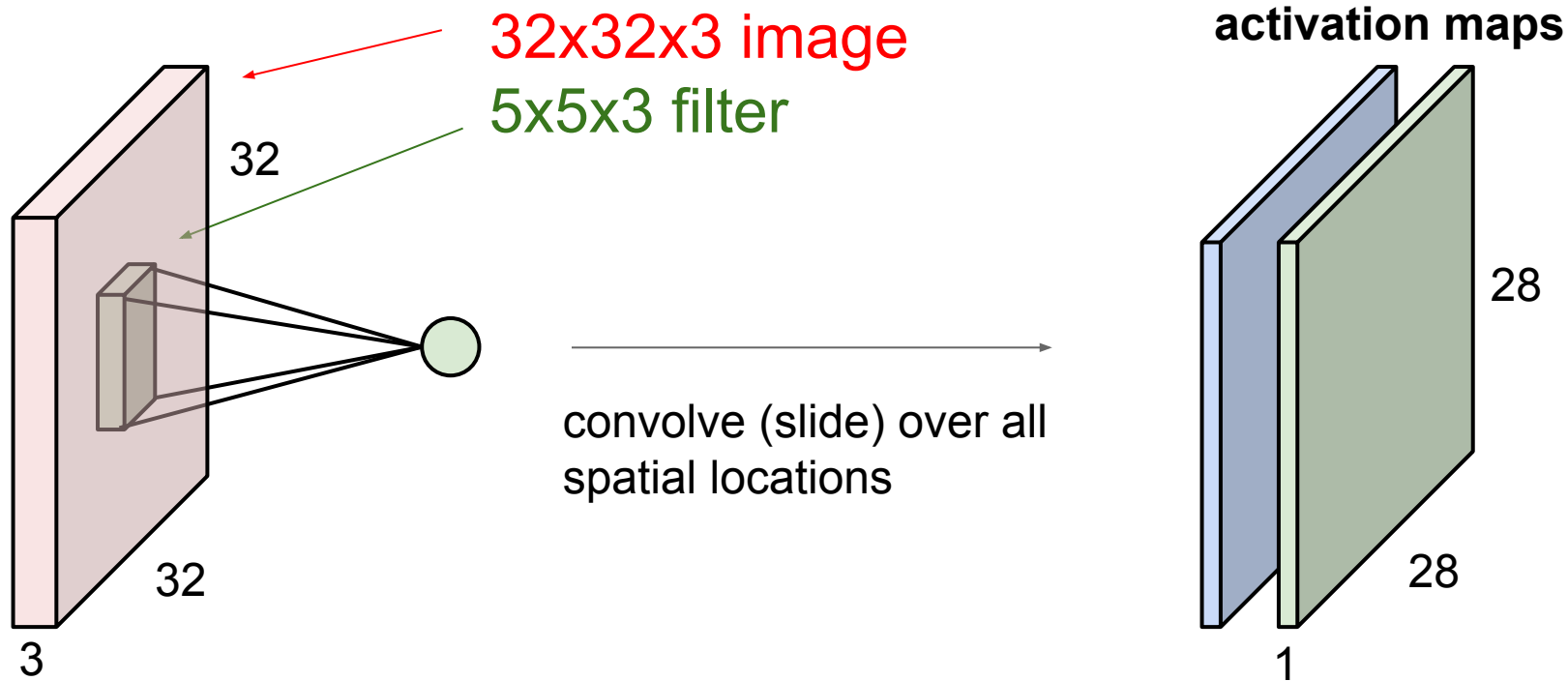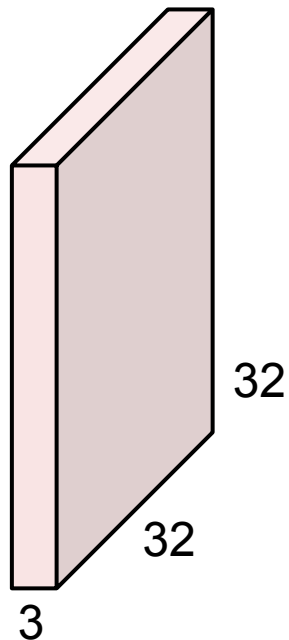
# Convolution Layer



32x32x3 image

5x5x3 filter

activation map

32

32

3

convolve (slide) over all
spatial locations

28

28

1

# Convolution Layer

consider a second, green filter



32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all
spatial locations

**activation maps**

28

28

1
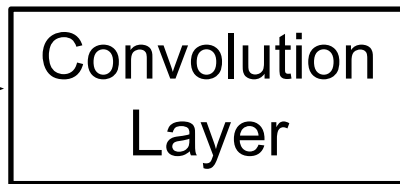
# Convolution Layer

3x32x32 image

Consider 6 filters,
each 3x5x5

6 activation maps,
each 1x28x28

32

32

3

Convolution
Layer

6x3x5x5
filters

Stack activations to get a
6x28x28 output image!

Slide inspiration: Justin Johnson

# Convolution Layer



3x32x32 image

Also 6-dim bias vector:

6 activation maps, each 1x28x28

32

32

3

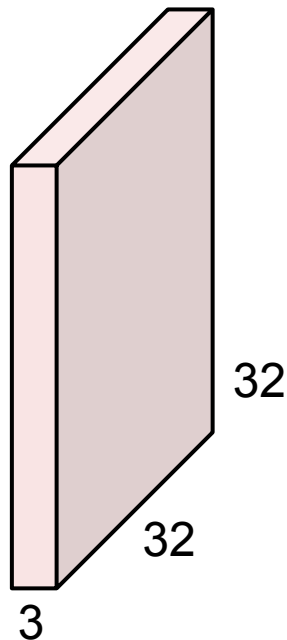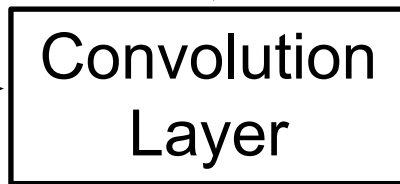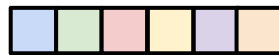6x3x5x5 filters

Convolution Layer

Stack activations to get a 6x28x28 output image!

Slide inspiration: Justin Johnson

# Convolution Layer



3x32x32 image

Also 6-dim bias vector:

28x28 grid, at each point a 6-dim vector

32
32
3

Convolution Layer

6x3x5x5 filters

Stack activations to get a 6x28x28 output image!

Slide inspiration: Justin Johnson

# Convolution Layer



2x3x32x32
Batch of images

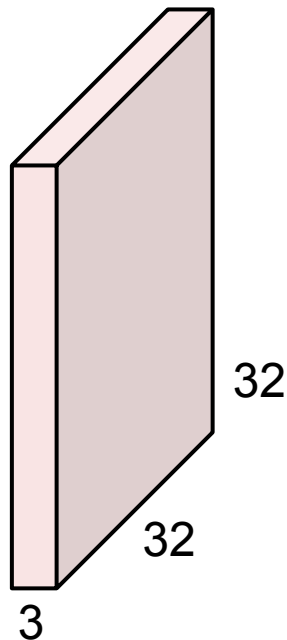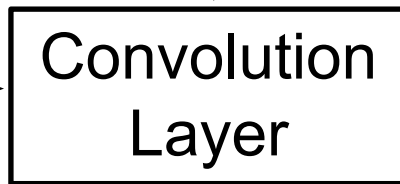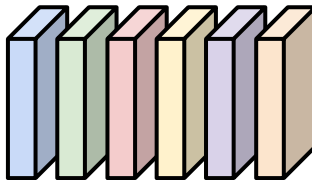Also 6-dim bias vector:

2x6x28x28
Batch of outputs

32

32

3

6x3x5x5
filters

Convolution Layer

Slide inspiration: Justin Johnson

# Convolution Layer

N x $C_{in}$ x H x W
Batch of images

N x $C_{out}$ x H' x W'
Batch of outputs

Also $C_{out}$-dim bias vector:



Convolution
Layer

H

W

$C_{in}$

$C_{out}$ x $C_{in}$ x $K_w$ x $K_h$
filters

$C_{out}$

Slide inspiration: Justin Johnson

**Preview:** ConvNet is a sequence of Convolution Layers



32

28

CONV
e.g. 6
5x5x3
filters

32

28

3

6

**Preview:** ConvNet is a sequence of Convolution Layers



32

32

3

CONV

e.g. 6
5x5x3
filters

28

28

6

CONV

e.g. 10
5x5x**6**
filters

24

24

10

CONV

....

**Preview:** ConvNet is a sequence of Convolution Layers, interspersed with activation functions



32
32
3

CONV
ReLU
e.g. 6
5x5x3
filters

28
28
6

CONV
ReLU
e.g. 10
5x5x**6**
filters

24
24
10

CONV
ReLU

....

A closer look at spatial dimensions:

7



7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:

7
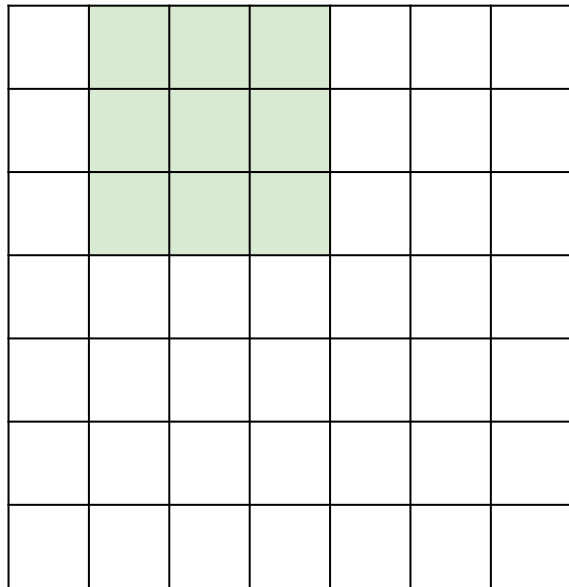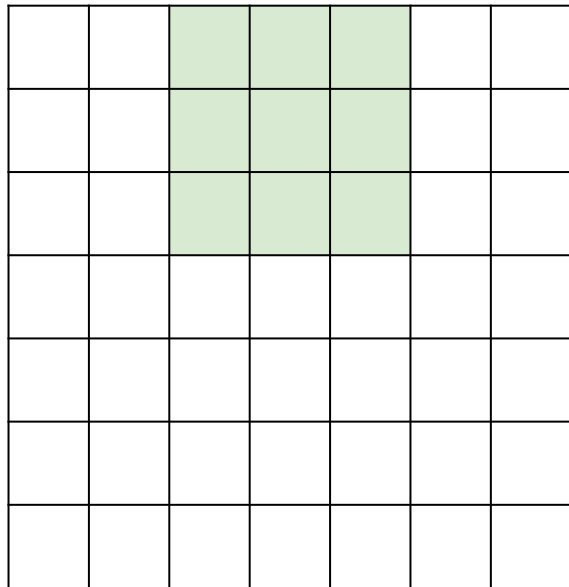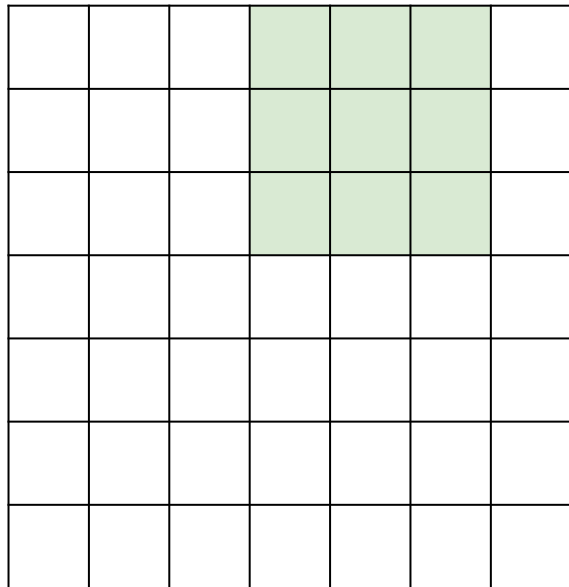


7

7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:

7



7

7x7 input (spatially)
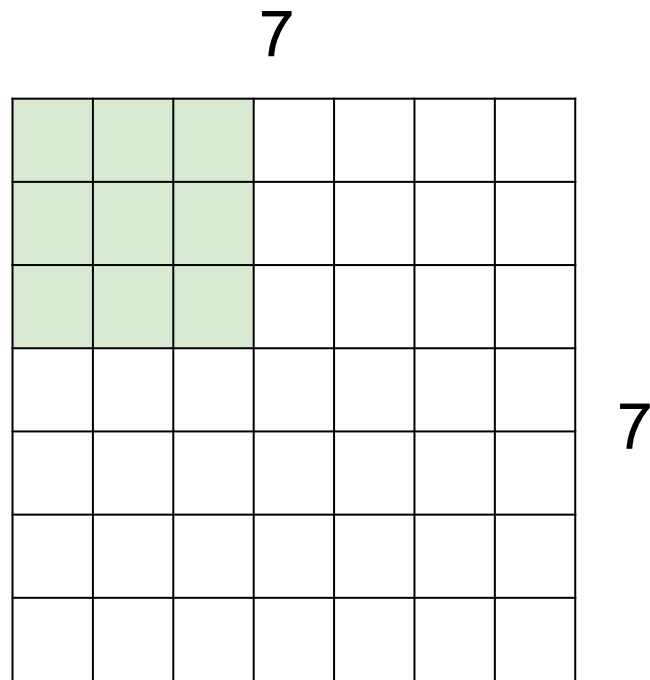assume 3x3 filter

A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter

**=> 5x5 output**
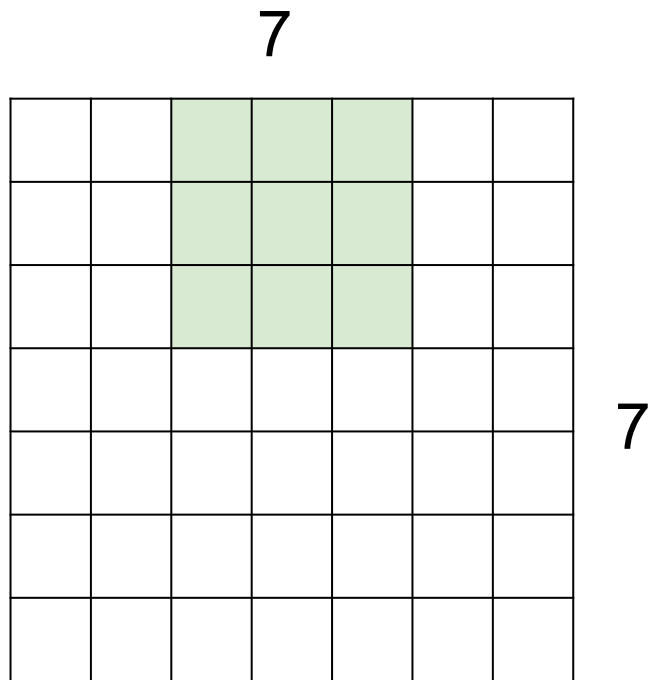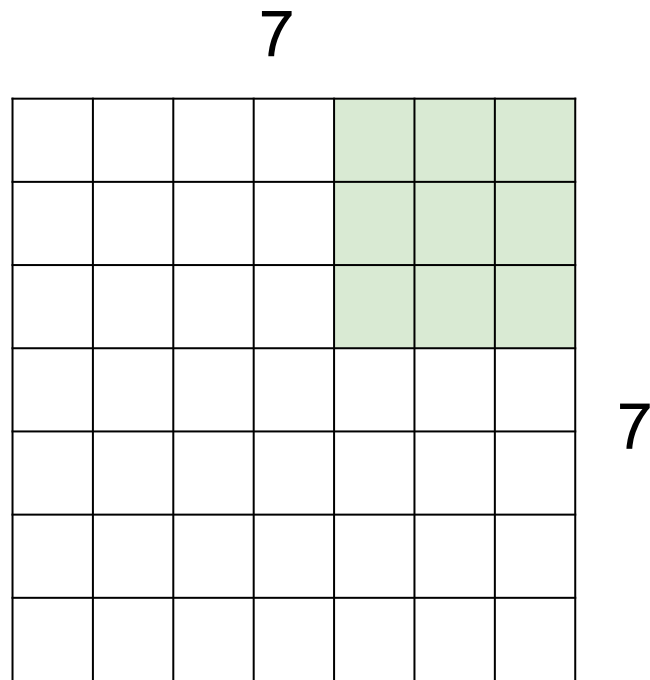
A closer look at spatial dimensions:
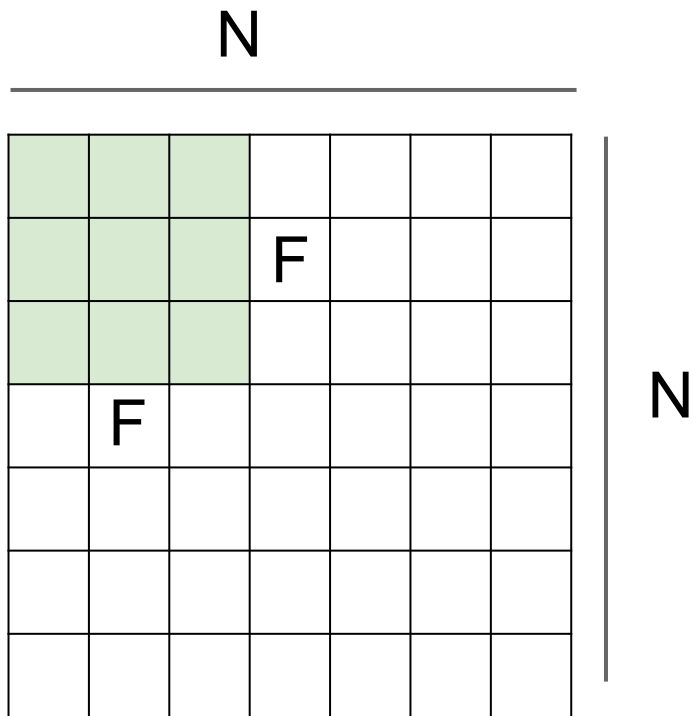
7



7

7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

A closer look at spatial dimensions:



7

7

7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**
**=> 3x3 output!**

Output size:
**(N - F) / stride + 1**

e.g. N = 7, F = 3:
stride 1 => (7 - 3)/1 + 1 = 5
stride 2 => (7 - 3)/2 + 1 = 3
stride 3 => (7 - 3)/3 + 1 = 2.33 :\

# In practice: Common to zero pad the border

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

e.g. input 7x7
**3x3** filter, applied with **stride 1**
**pad with 1 pixel** border => what is the output?

(recall:)
(N - F) / stride + 1

# In practice: Common to zero pad the border



e.g. input 7x7
**3x3** filter, applied with **stride 1**
**pad with 1 pixel** border => what is the output?

**7x7 output!**

(recall:)
(N + 2P - F) / stride + 1

# In practice: Common to zero pad the border



e.g. input 7x7
**3x3** filter, applied with **stride 1**
**pad with 1 pixel** border => what is the output?

**7x7 output!**
in general, common to see CONV layers with stride 1, filters of size FxF, and zero-padding with (F-1)/2. (will preserve size spatially)
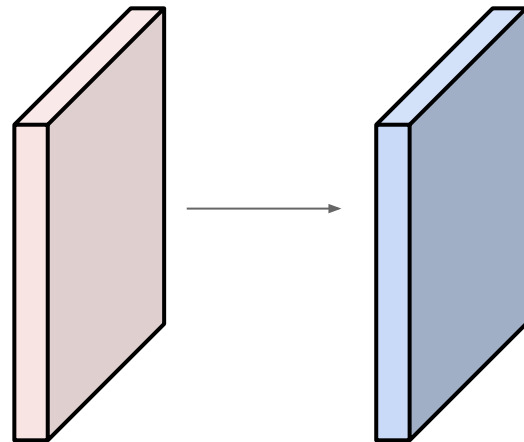e.g. F = 3 => zero pad with 1
       F = 5 => zero pad with 2
       F = 7 => zero pad with 3

Examples time:

Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2



Output volume size: ?

Examples time:
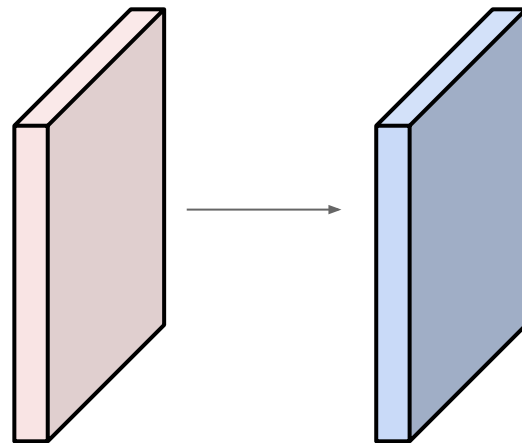
Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2

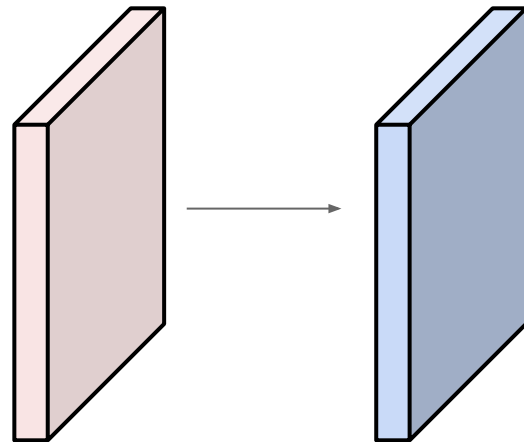Output volume size:
(32+2*2-5)/1+1 = 32 spatially, so
**32x32x10**

Examples time:

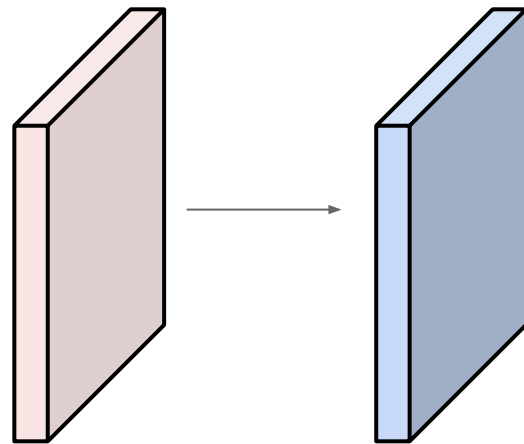Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2



Number of parameters in this layer?

Examples time:

Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2



Number of parameters in this layer?
each filter has 5*5*3 + 1 = 76 params        (+1 for bias)
=> 76*10 = **760**

# Convolution layer: summary

Let's assume input is $W_1$ x $H_1$ x C

Conv layer needs 4 hyperparameters:
- Number of filters **K**
- The filter size **F**
- The stride **S**
- The zero padding **P**

This will produce an output of $W_2$ x $H_2$ x K
where:
- $W_2 = (W_1 - F + 2P)/S + 1$
- $H_2 = (H_1 - F + 2P)/S + 1$

Number of parameters: $F^2CK$ and K biases

# Convolution layer: summary

Let's assume input is $W_1$ x $H_1$ x C
Conv layer needs 4 hyperparameters:
- Number of filters **K**
- The filter size **F**
- The stride **S**
- The zero padding **P**

This will produce an output of $W_2$ x $H_2$ x K
where:
- $W_2 = (W_1 - F + 2P)/S + 1$
- $H_2 = (H_1 - F + 2P)/S + 1$

Number of parameters: $F^2CK$ and K biases

K = (powers of 2, e.g. 32, 64, 128, 512)
- F = 3, S = 1, P = 1
- F = 5, S = 1, P = 2
- F = 5, S = 2, P = ? (whatever fits)
- F = 1, S = 1, P = 0

# Example: CONV layer in PyTorch

Conv2d

Conv layer needs 4 hyperparameters:
- Number of filters **K**
- The filter size **F**
- The stride **S**
- The zero padding **P**