

# Lecture 10:

# Training Neural Networks

# Overview

## 1. One time setup

*activation functions, preprocessing, weight initialization, regularization, gradient checking*

## 2. Training dynamics

*babysitting the learning process, parameter updates, hyperparameter optimization*

## 3. Evaluation

*model ensembles*

# Evaluation:

## Model Ensembles

1. Train multiple independent models
2. At test time average their results

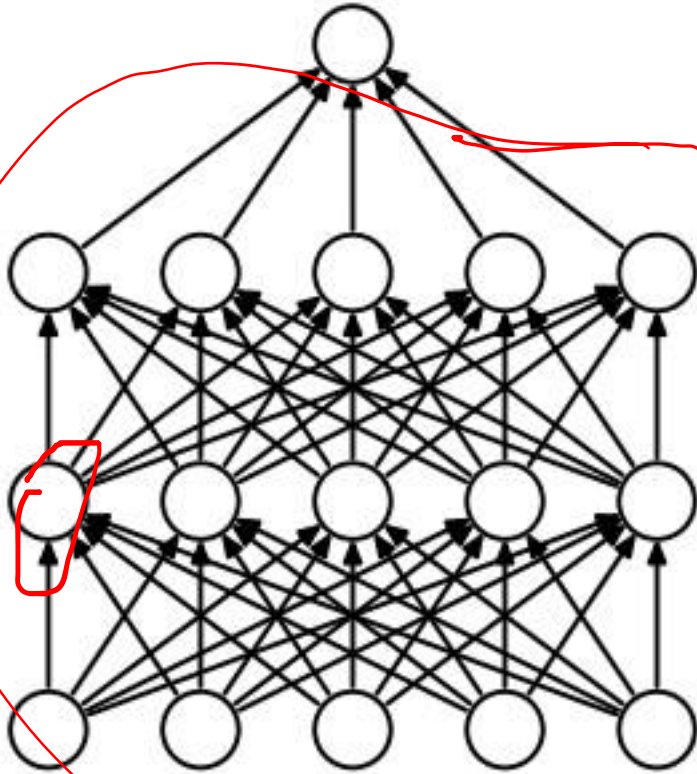
Enjoy 2% extra performance (?!!!)

## Fun Tips/Tricks:

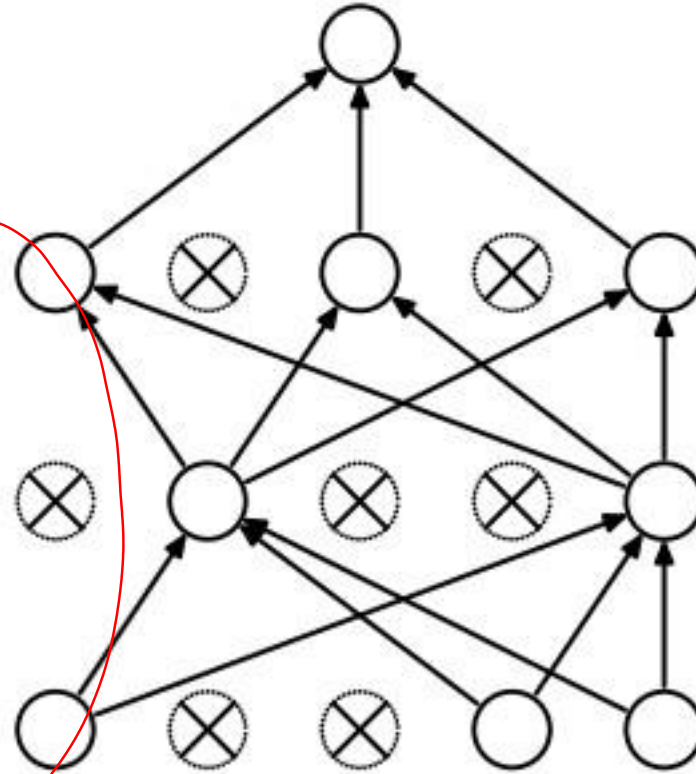
- can also get a small boost from averaging multiple model checkpoints of a single model.

# Regularization: Dropout

“randomly set some neurons to zero in the forward pass”



(a) Standard Neural Net



(b) After applying dropout.

*[Srivastava et al., 2014]*

```
p = 0.5 # probability of keeping a unit active. higher = less dropout
```

```
def train_step(X):
```

```
    """ X contains the data """
```

```
    # forward pass for example 3-layer neural network
```

```
    H1 = np.maximum(0, np.dot(W1, X) + b1)
```

```
    U1 = np.random.rand(*H1.shape) < p # first dropout mask
```

```
    H1 *= U1 # drop!
```

```
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
```

```
    U2 = np.random.rand(*H2.shape) < p # second dropout mask
```

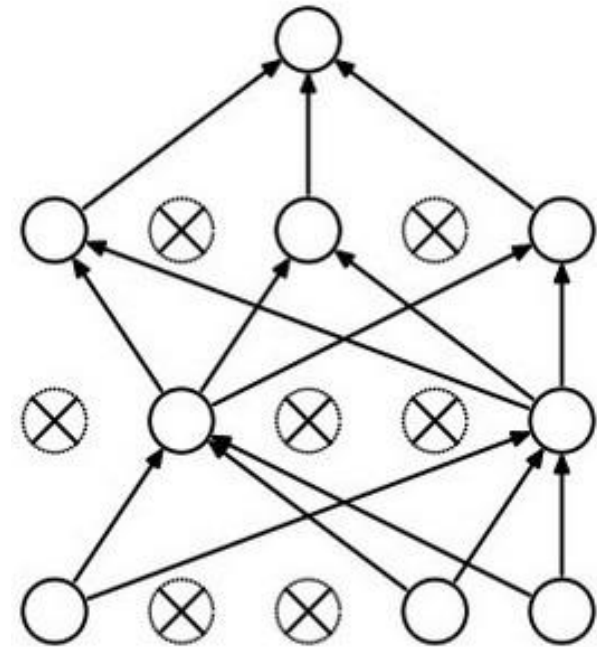
```
    H2 *= U2 # drop!
```

```
    out = np.dot(W3, H2) + b3
```

```
    # backward pass: compute gradients... (not shown)
```

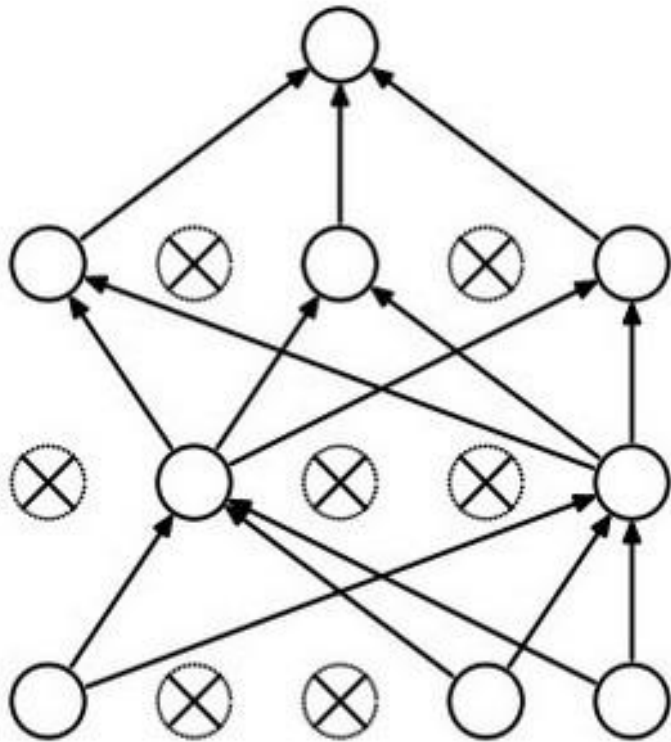
```
    # perform parameter update... (not shown)
```

Example forward pass with a 3-layer network using dropout



Waaaait a second...

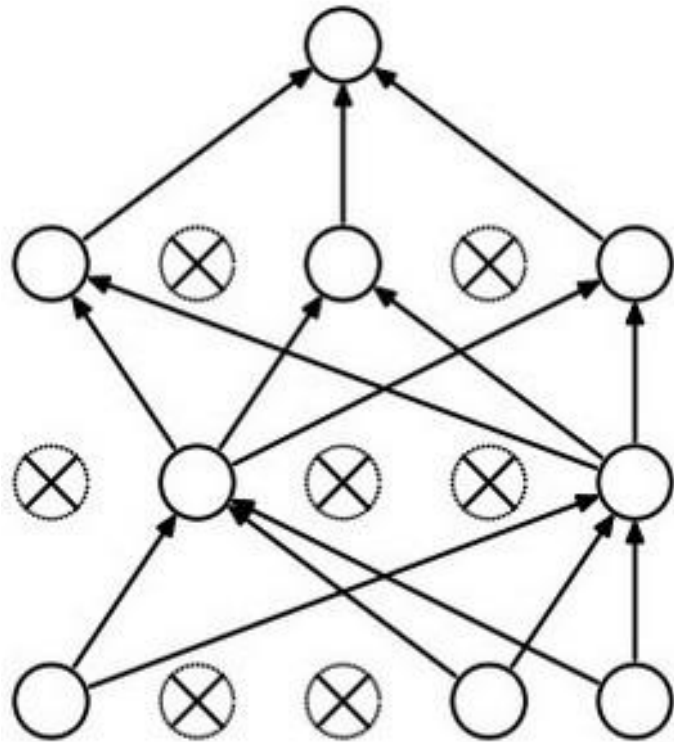
How could this possibly be a good idea?





# Waaaait a second...

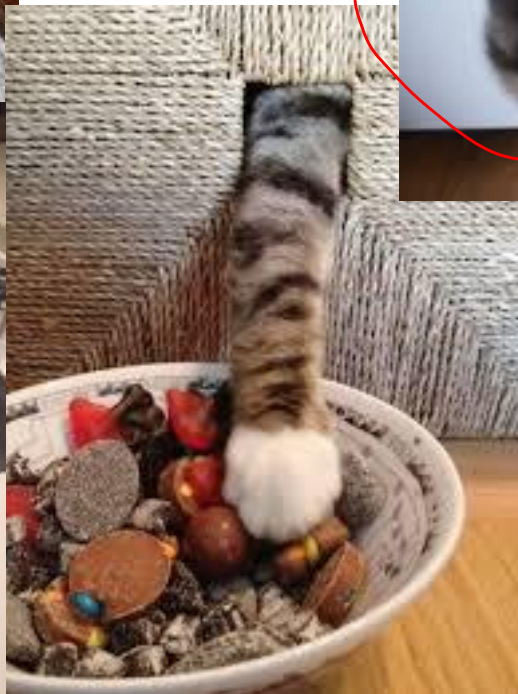
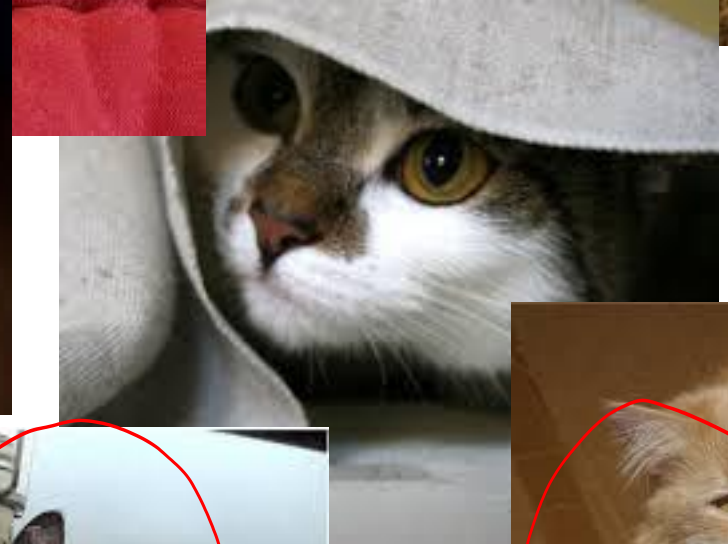
# How could this possibly be a good idea?



Forces the network to have a redundant representation.

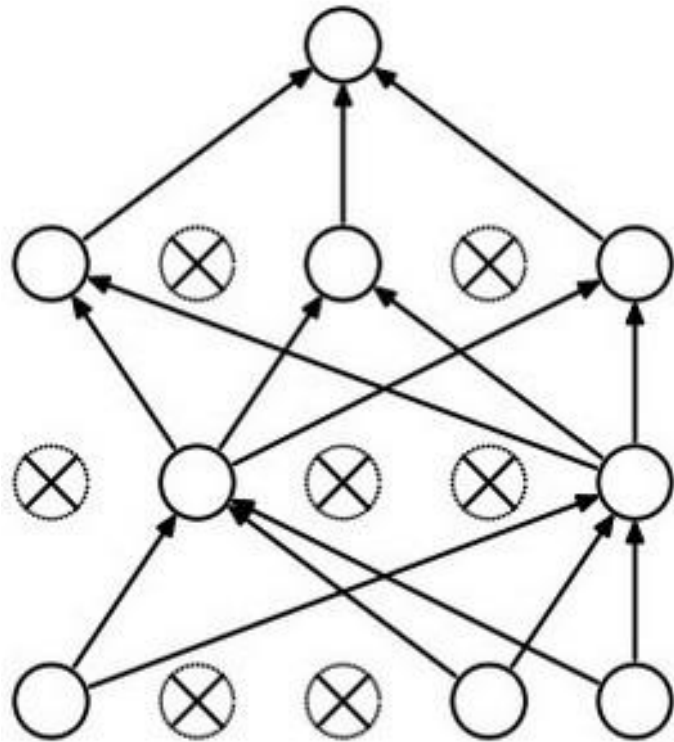


# Training with occlusions?



Waaaait a second...

How could this possibly be a good idea?



Another interpretation:

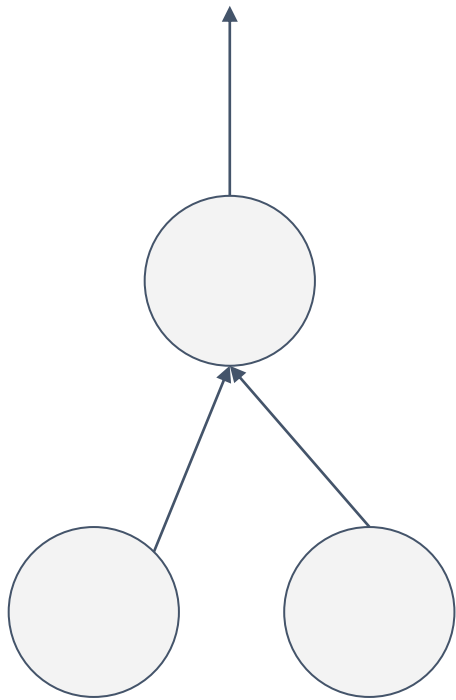
Dropout is training a large ensemble of models (that share parameters).

Each binary mask is one model, gets trained on only ~one datapoint.

# At test time....

Can in fact do this with a single forward pass! (approximately)

Leave all input neurons turned on (no dropout).

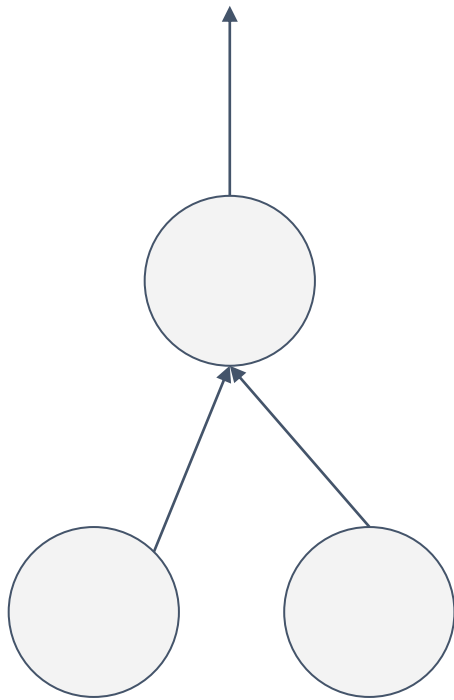


(this can be shown to be an approximation to evaluating the whole ensemble)

# At test time....

Can in fact do this with a single forward pass! (approximately)

Leave all input neurons turned on (no dropout).



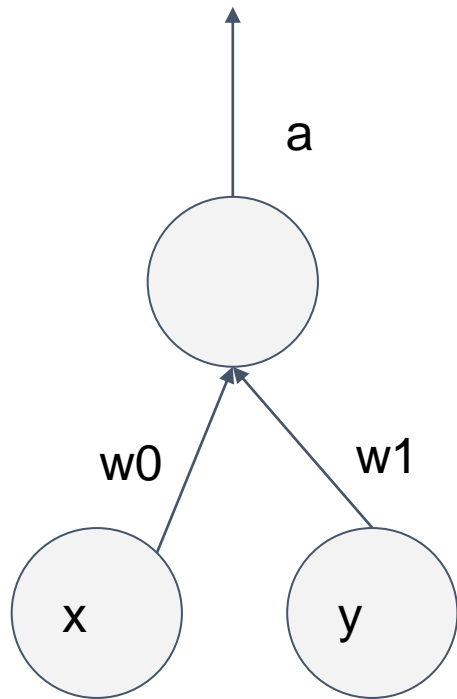
Q: Suppose that with all inputs present at test time the output of this neuron is  $x$ .

What would its output be during training time, in expectation? (e.g. if  $p = 0.5$ )

# At test time....

Can in fact do this with a single forward pass! (approximately)

Leave all input neurons turned on (no dropout).



during test:  $a = w_0 * x + w_1 * y$

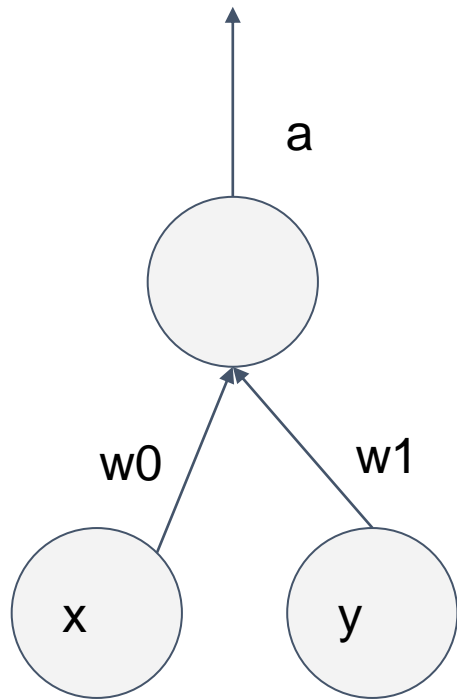
during train:

$$\begin{aligned} E[a] &= \frac{1}{4} * (\underbrace{w_0 * 0 + w_1 * 0}_{w_0 * 0 + w_1 * y} \\ &\quad \underbrace{w_0 * 0 + w_1 * y}_{w_0 * x + w_1 * 0} \\ &\quad \underbrace{w_0 * x + w_1 * y}) \\ &= \frac{1}{4} * (2 w_0 * x + 2 w_1 * y) \\ &= \frac{1}{2} * (w_0 * x + w_1 * y) \end{aligned}$$

# At test time....

Can in fact do this with a single forward pass! (approximately)

Leave all input neurons turned on (no dropout).



during test:  $a = w_0 * x + w_1 * y$

during train:

$$\begin{aligned} E[a] &= \frac{1}{4} * (w_0 * 0 + w_1 * 0 \\ &\quad w_0 * 0 + w_1 * y \\ &\quad w_0 * x + w_1 * 0 \\ &\quad w_0 * x + w_1 * y) \\ &= \frac{1}{4} * (2 w_0 * x + 2 w_1 * y) \\ &= \frac{1}{2} * (w_0 * x + w_1 * y) \end{aligned}$$

With  $p=0.5$ , using all inputs in the forward pass would inflate the activations by 2x from what the network was “used to” during training!  
=> Have to compensate by scaling the activations back down by  $\frac{1}{2}$

# We can do something approximate analytically

```
def predict(X):  
    # ensembled forward pass  
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations  
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations  
    out = np.dot(W3, H2) + b3
```

At test time all neurons are active always

=> We must scale the activations so that for each neuron:

output at test time = expected output at training time



# Dropout Summary

```
""" Vanilla Dropout: Not recommended implementation (see notes below) """
```

```
p = 0.5 # probability of keeping a unit active. higher = less dropout
```

```
def train_step(X):
```

```
    """ X contains the data """
```

```
    # forward pass for example 3-layer neural network
```

```
    H1 = np.maximum(0, np.dot(W1, X) + b1)
```

```
    U1 = np.random.rand(*H1.shape) < p # first dropout mask
```

```
    H1 *= U1 # drop!
```

```
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
```

```
    U2 = np.random.rand(*H2.shape) < p # second dropout mask
```

```
    H2 *= U2 # drop!
```

```
    out = np.dot(W3, H2) + b3
```

```
    # backward pass: compute gradients... (not shown)
```

```
    # perform parameter update... (not shown)
```

```
def predict(X):
```

```
    # ensembled forward pass
```

```
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations
```

```
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations
```

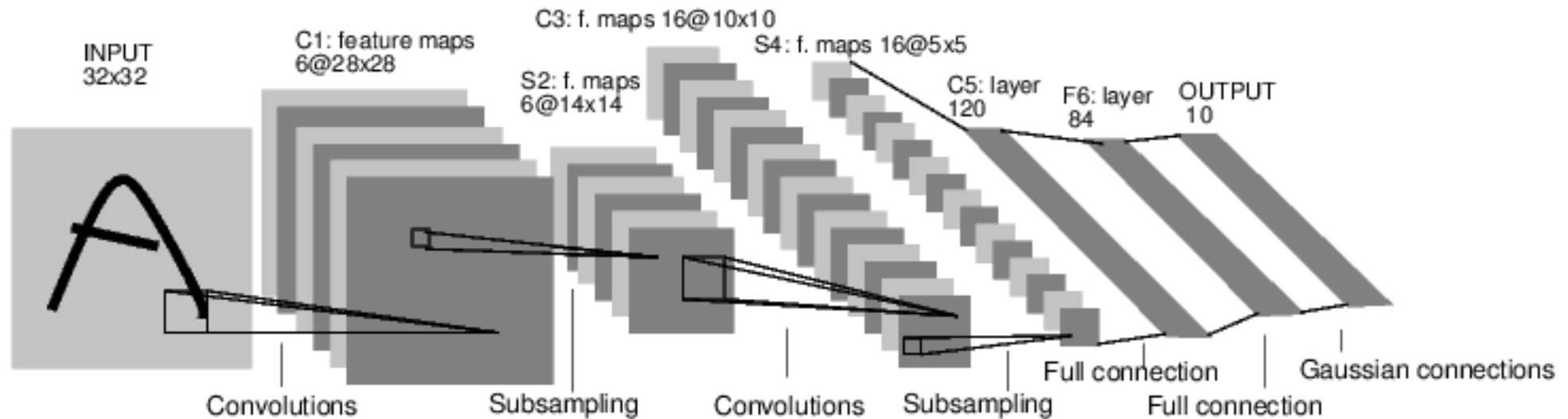
```
    out = np.dot(W3, H2) + b3
```

drop in forward pass

scale at test time

# Lecture 11:

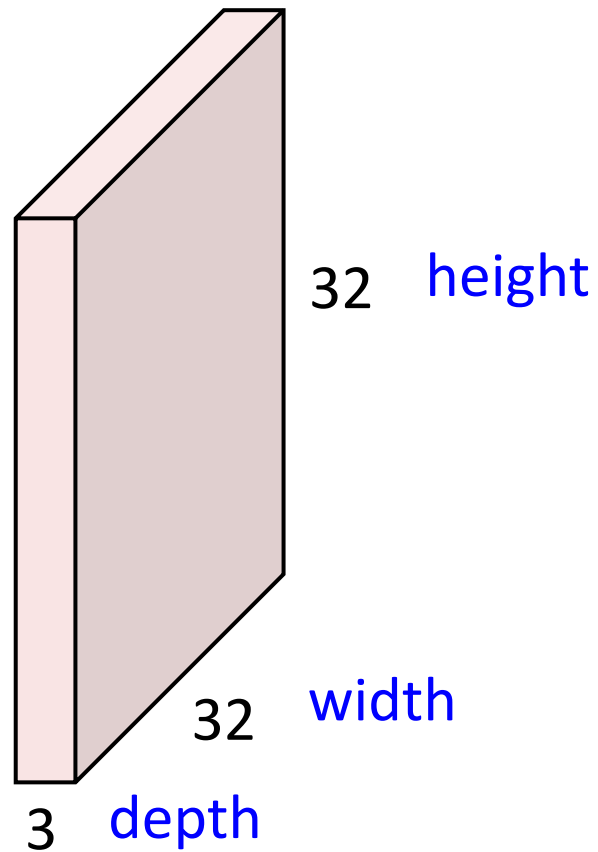
# Convolutional Neural Networks



[LeNet-5, LeCun 1980]

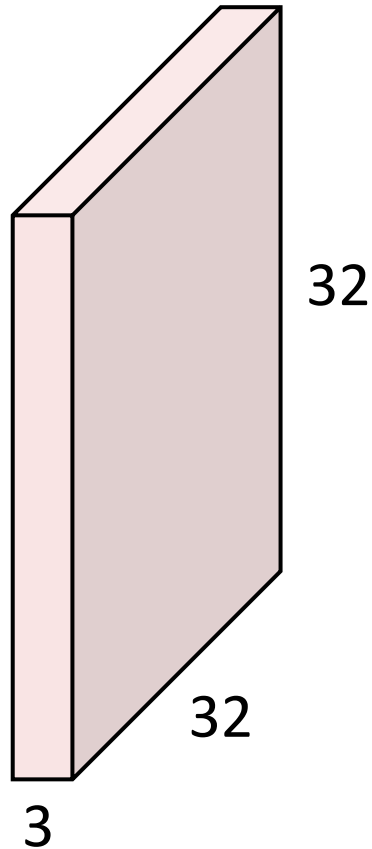
# Convolution Layer

32x32x3 image



# Convolution Layer

32x32x3 image



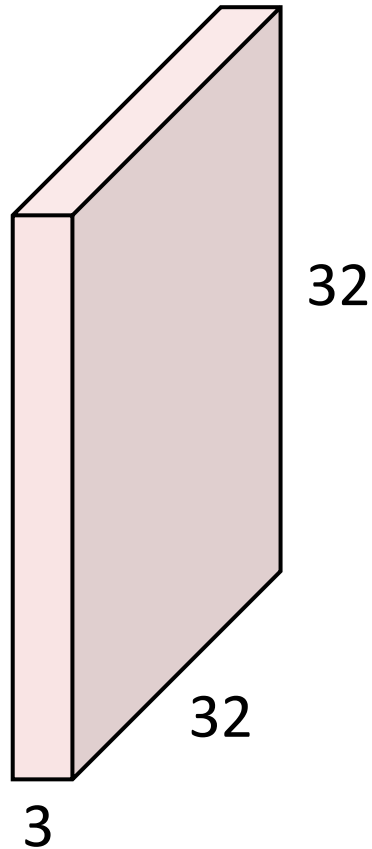
5x5x3 filter



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

# Convolution Layer

32x32x3 image



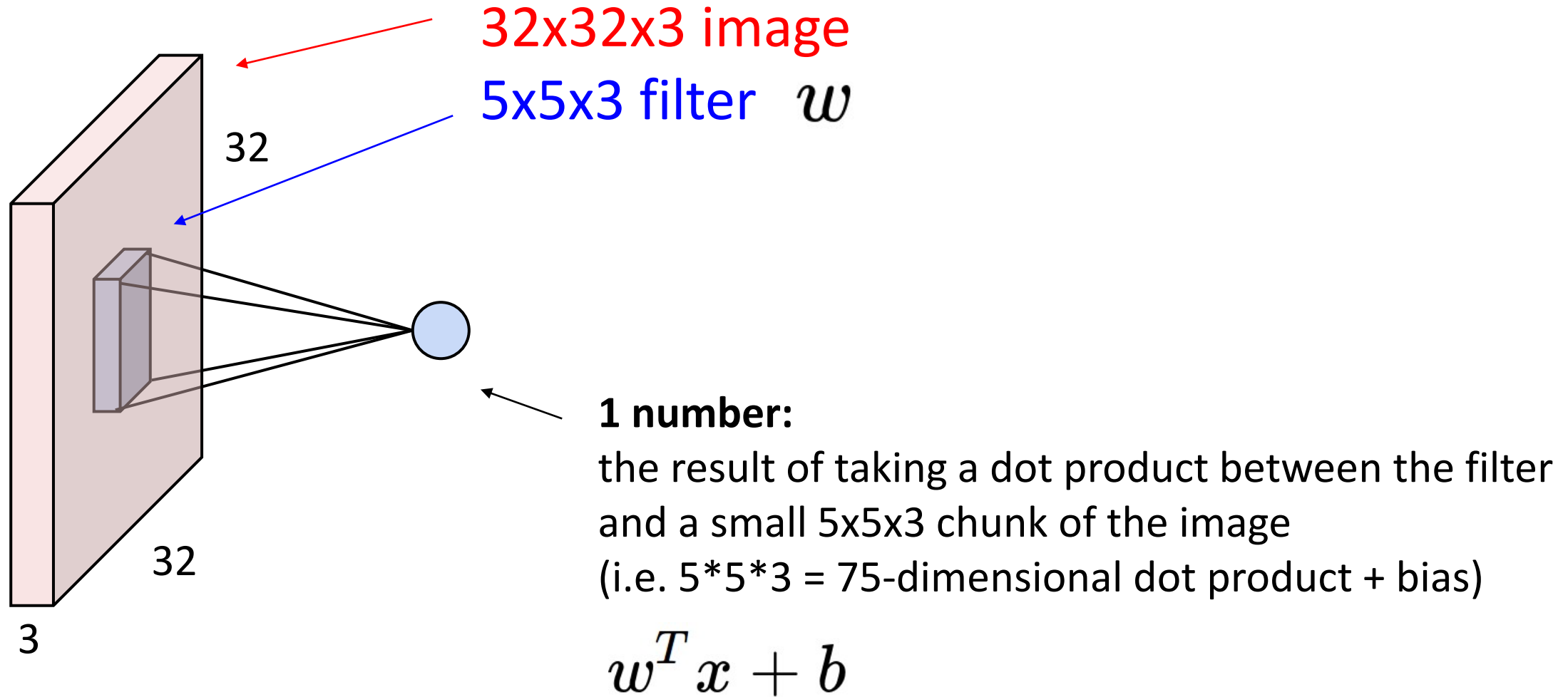
Filters always extend the full depth of the input volume

5x5x3 filter

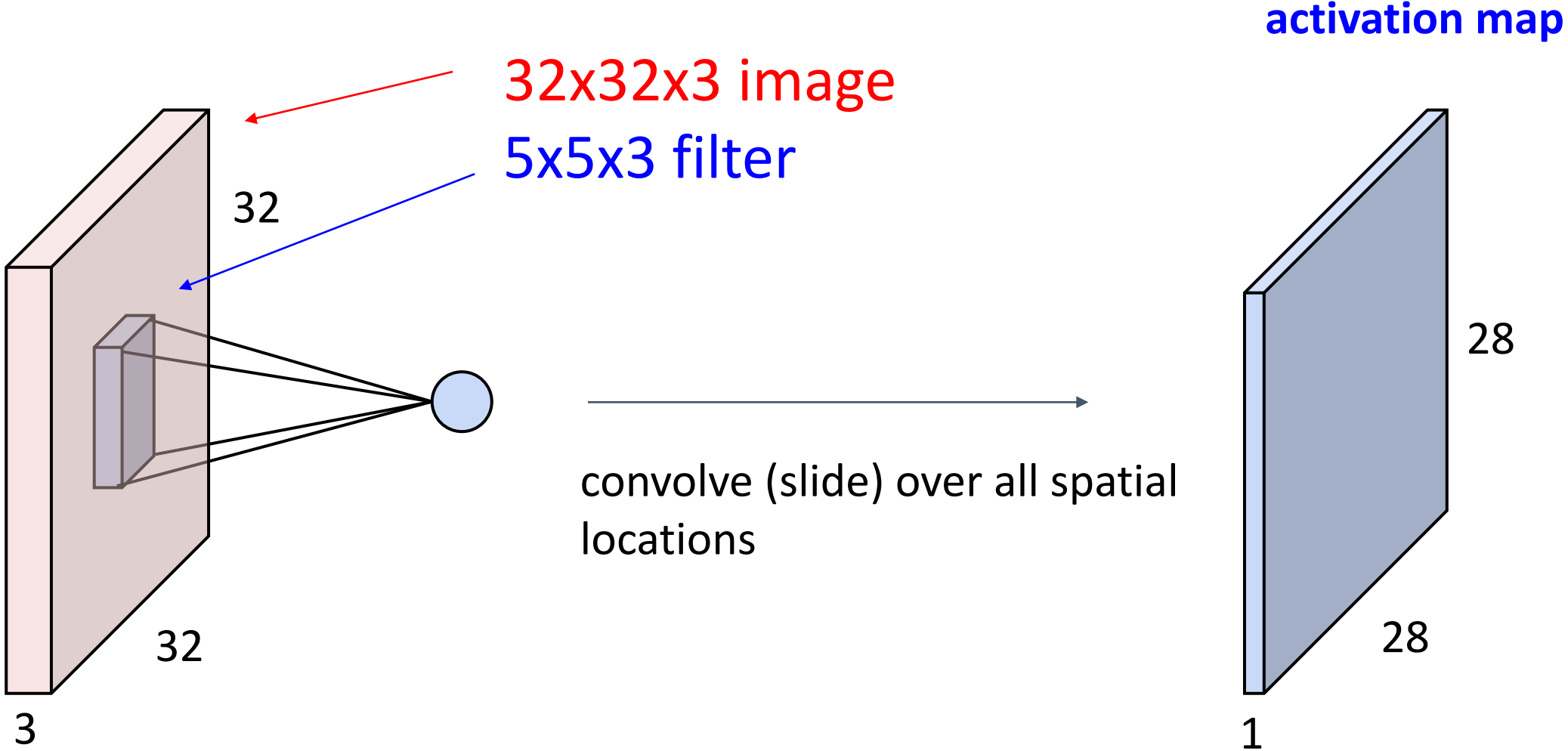


**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

# Convolution Layer

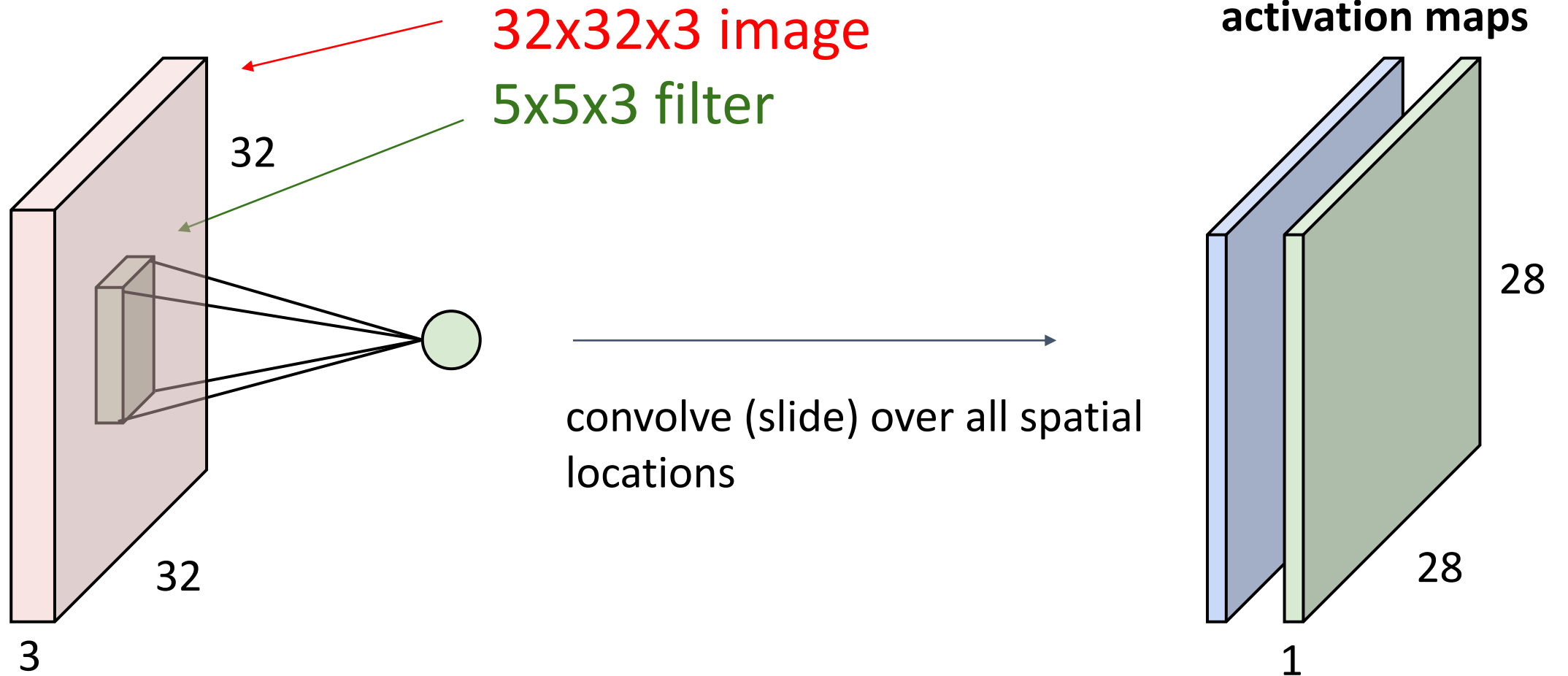


# Convolution Layer



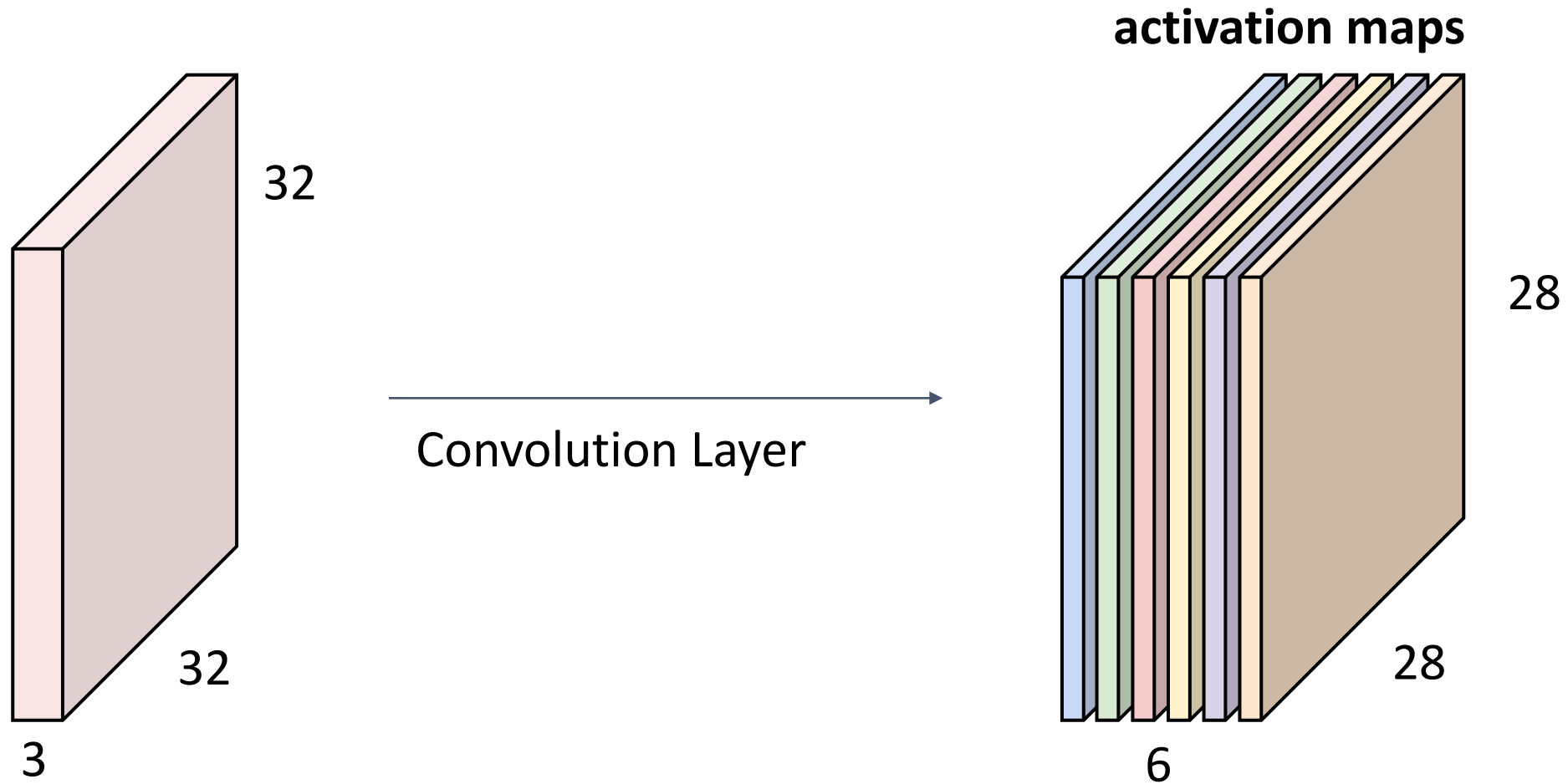
# Convolution Layer

consider a second, green filter



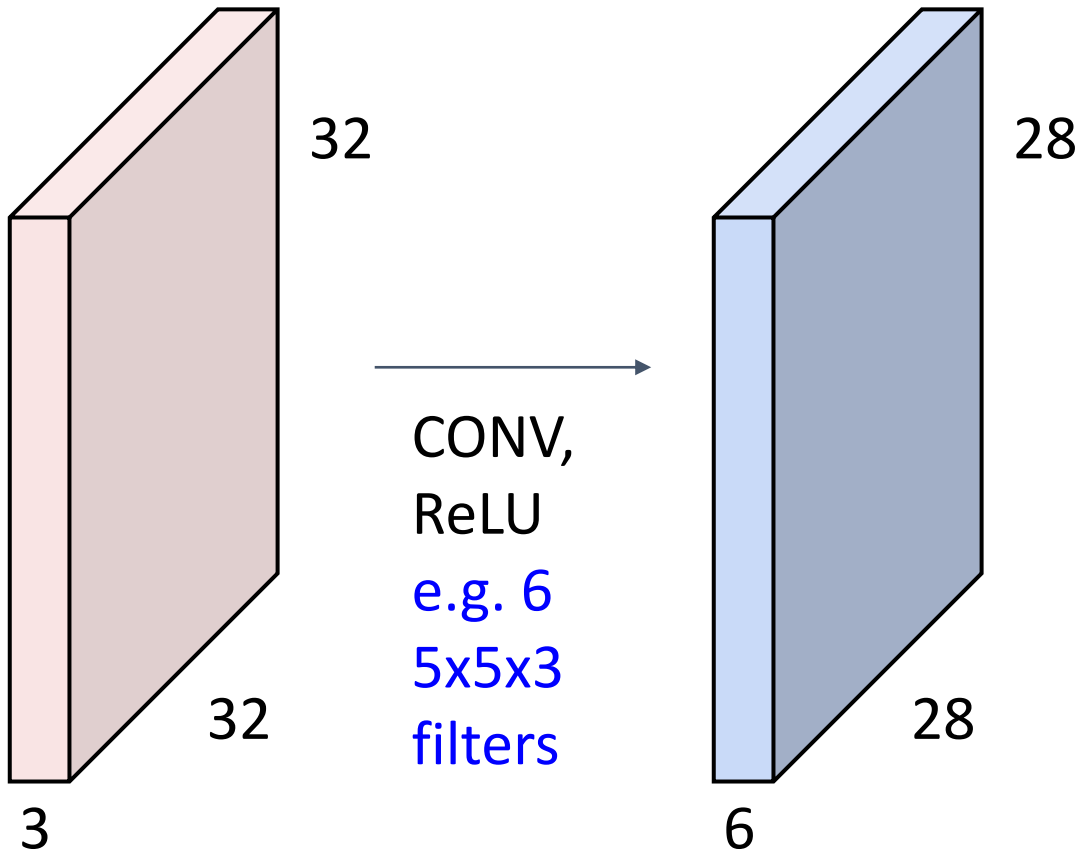


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

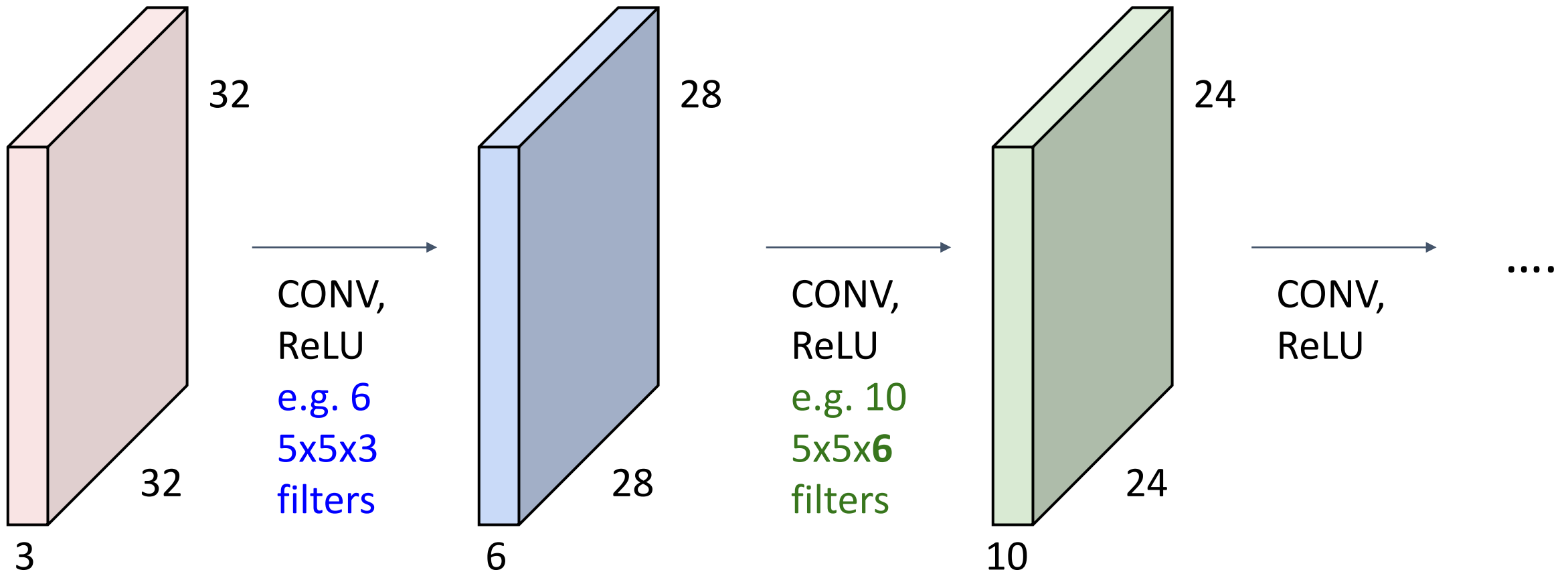


We stack these up to get a “new image” of size 28x28x6!

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

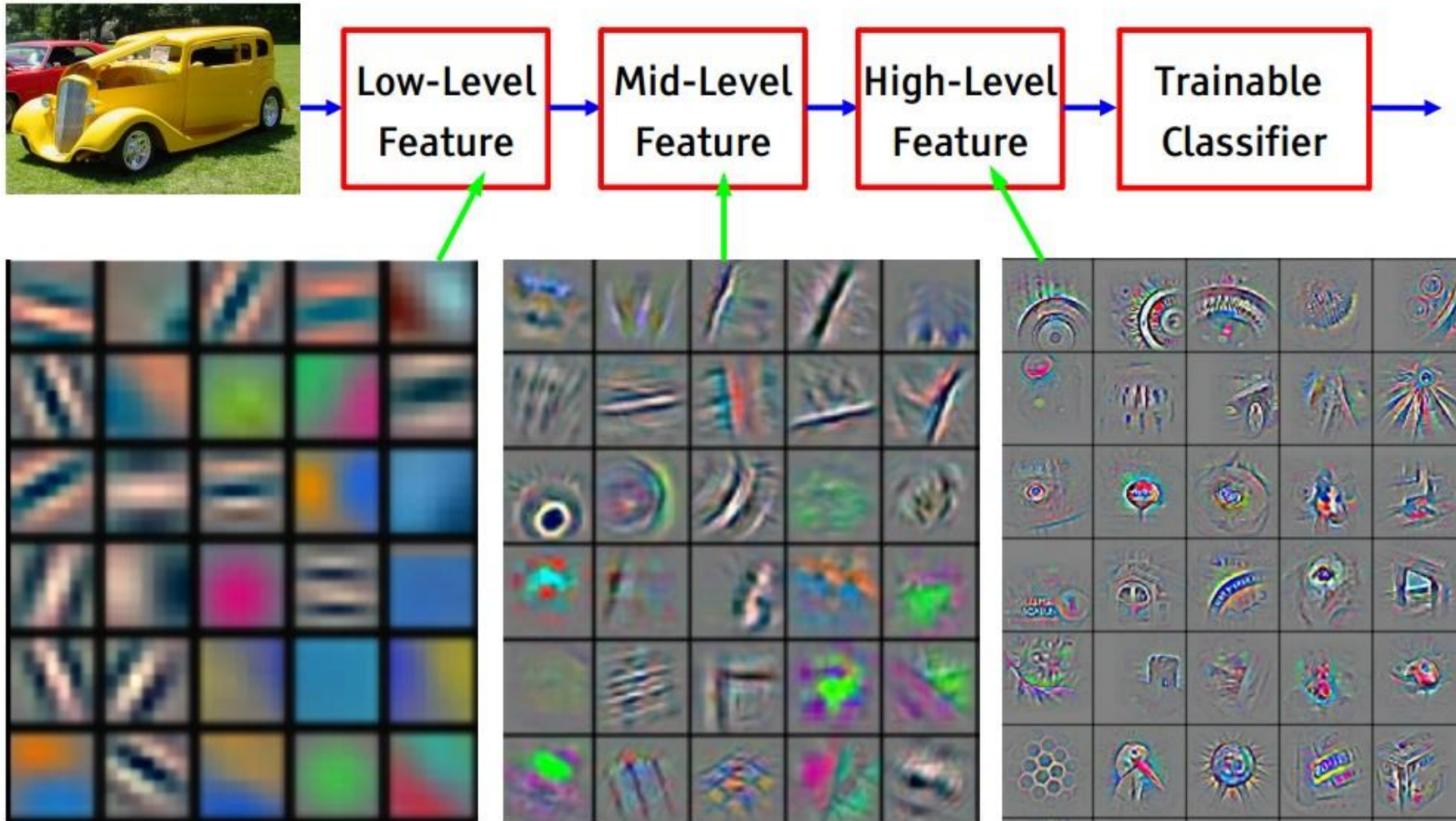


**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



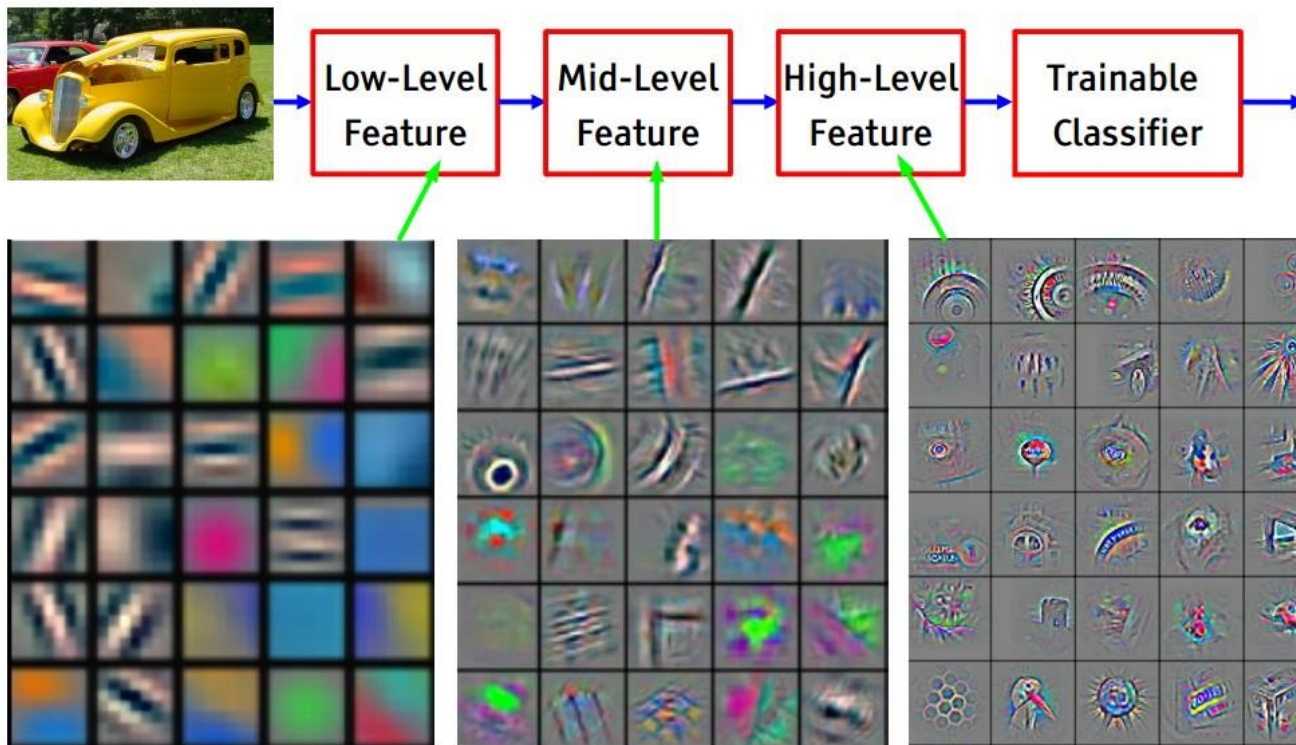
# Preview

[From recent Yann LeCun slides]



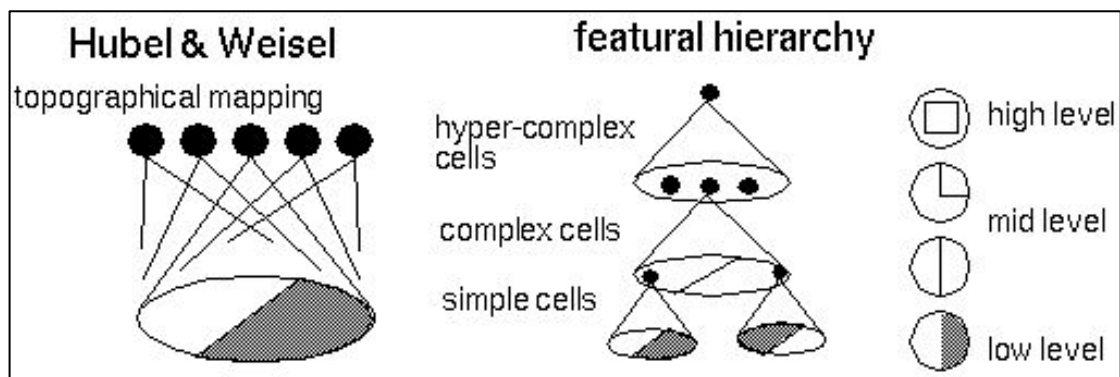
Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

# Preview



[From recent Yann LeCun slides]

Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

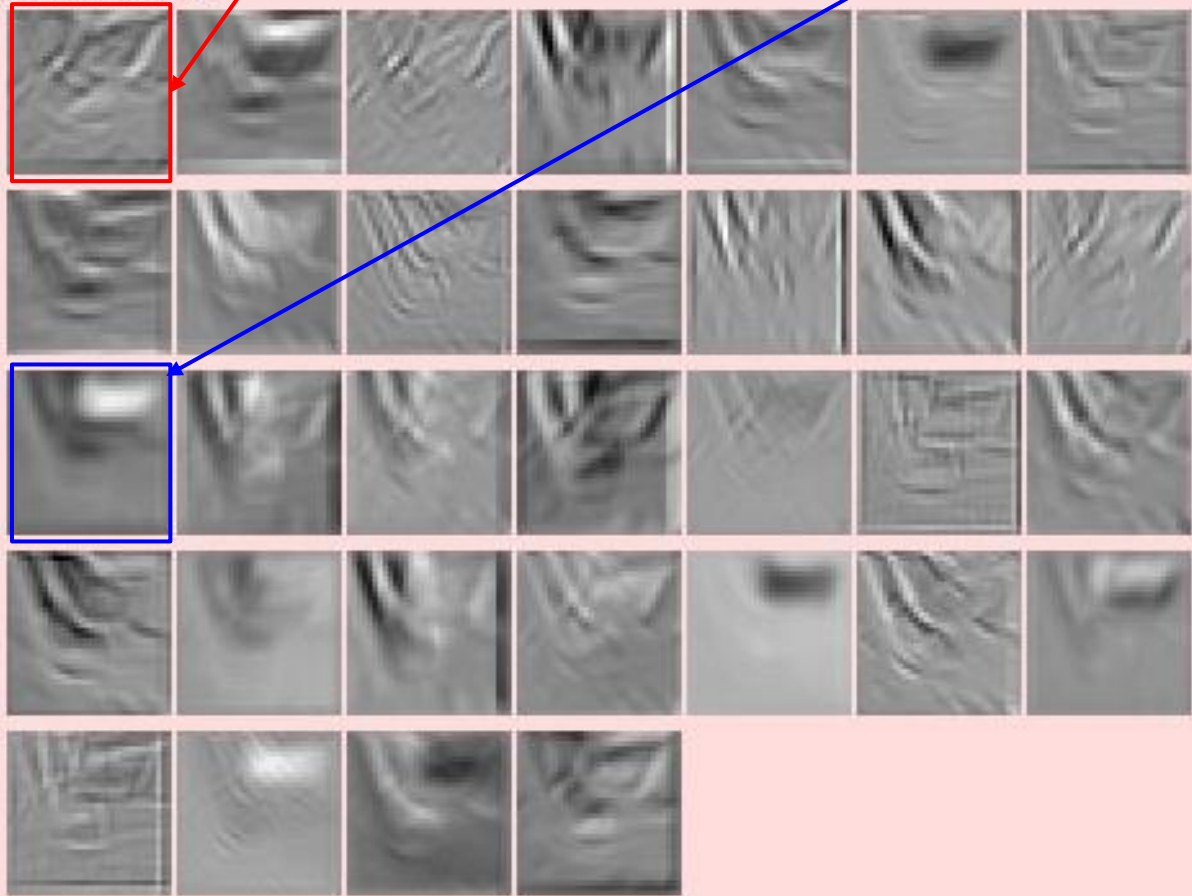




one filter =>  
one activation map

example 5x5 filters  
(32 total)

Activations:

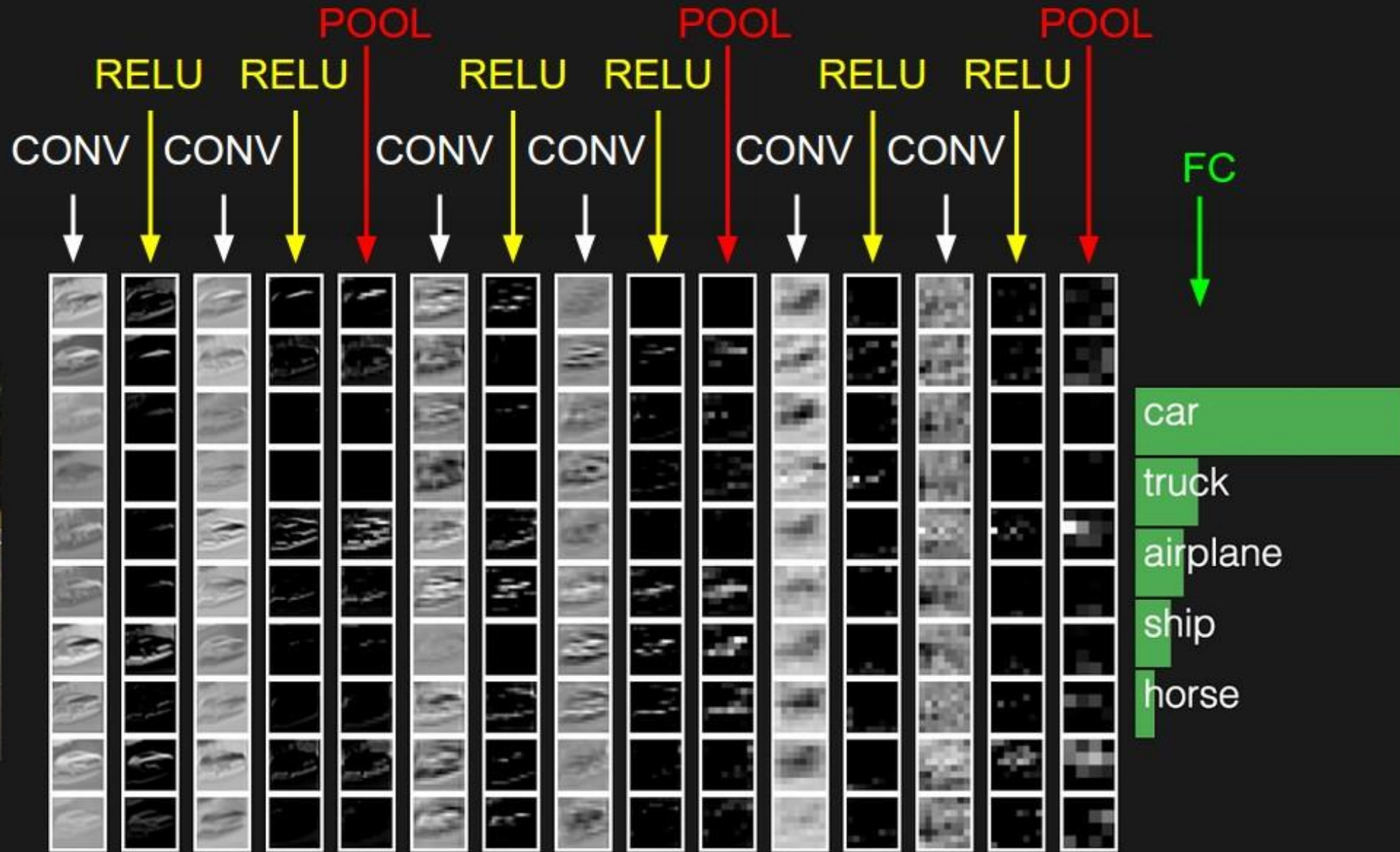


We call the layer convolutional because it is related to convolution of two signals:

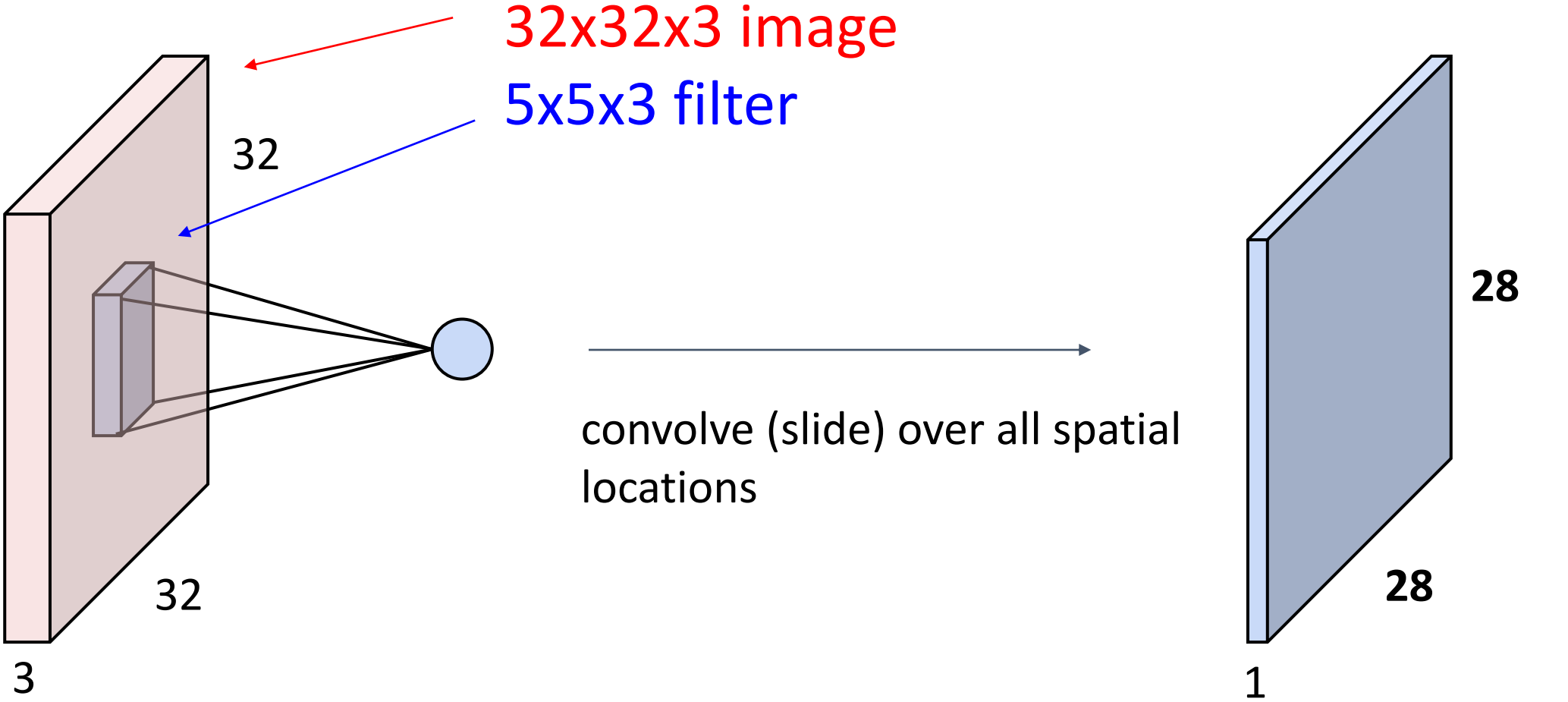
$$f[x,y] * g[x,y] = \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1,n_2] \cdot g[x-n_1,y-n_2]$$

↑  
elementwise multiplication and sum of a filter and the signal (image)

preview:



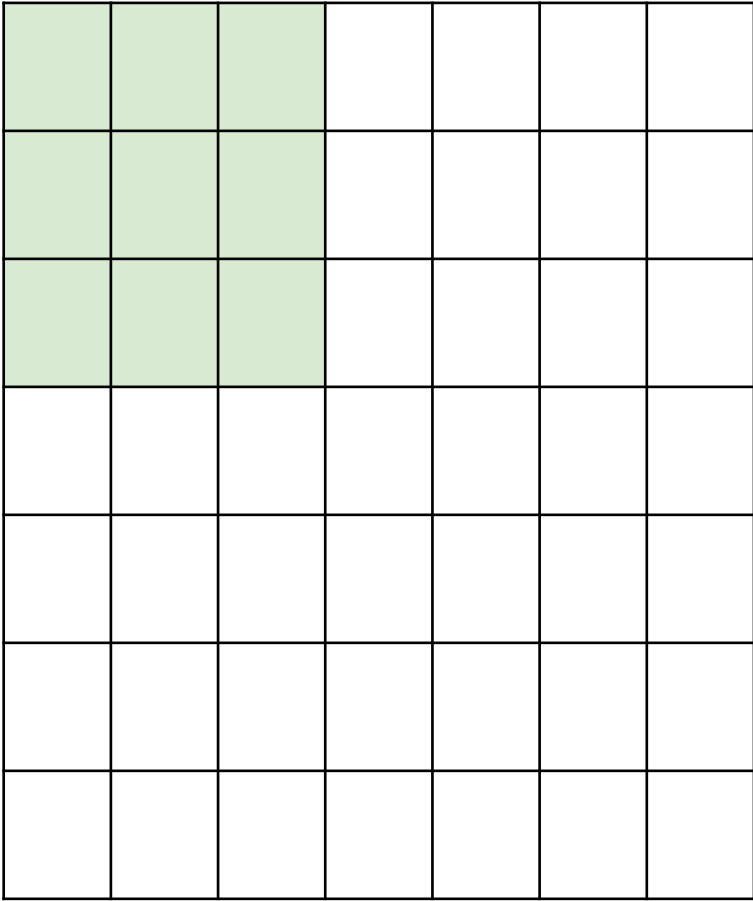
# A closer look at spatial dimensions:





# A closer look at spatial dimensions:

7

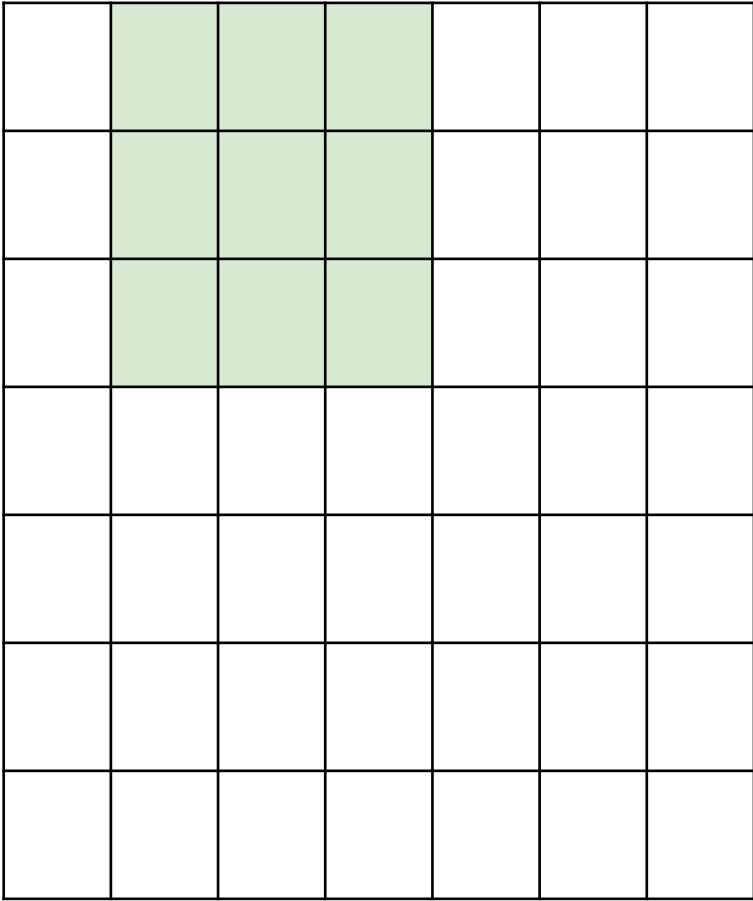


7x7 input (spatially)  
assume 3x3 filter

7

# A closer look at spatial dimensions:

7

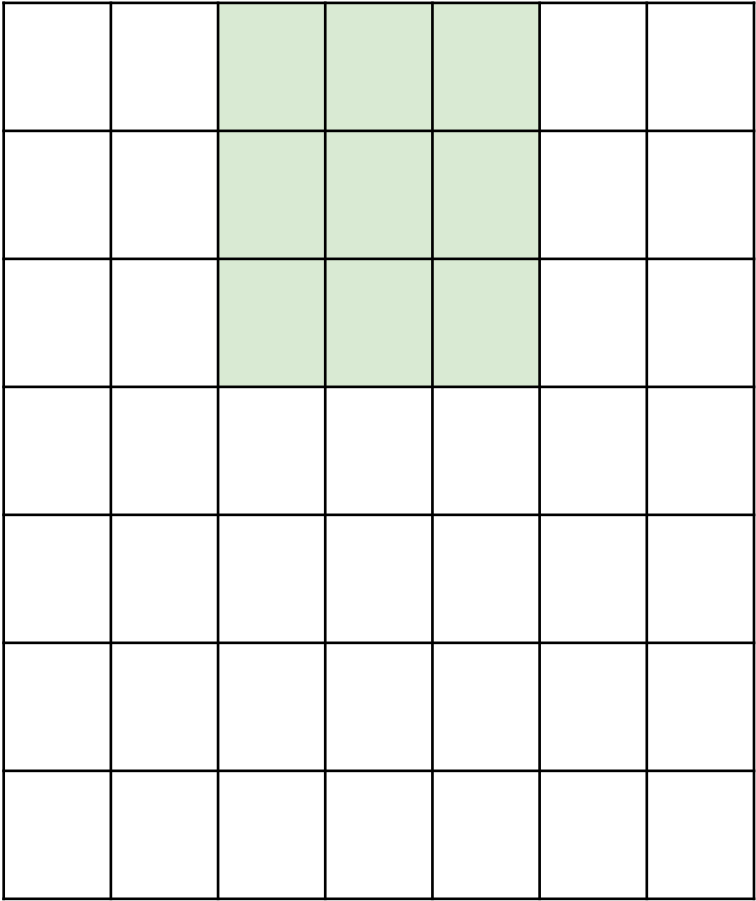


7x7 input (spatially)  
assume 3x3 filter

7

# A closer look at spatial dimensions:

7

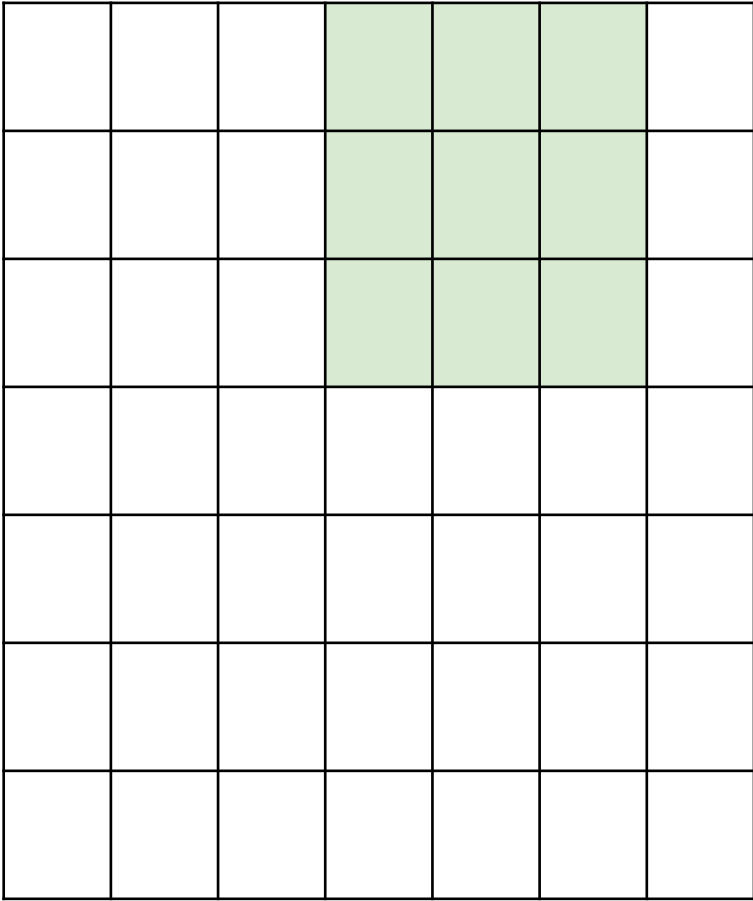


7x7 input (spatially)  
assume 3x3 filter

7

# A closer look at spatial dimensions:

7

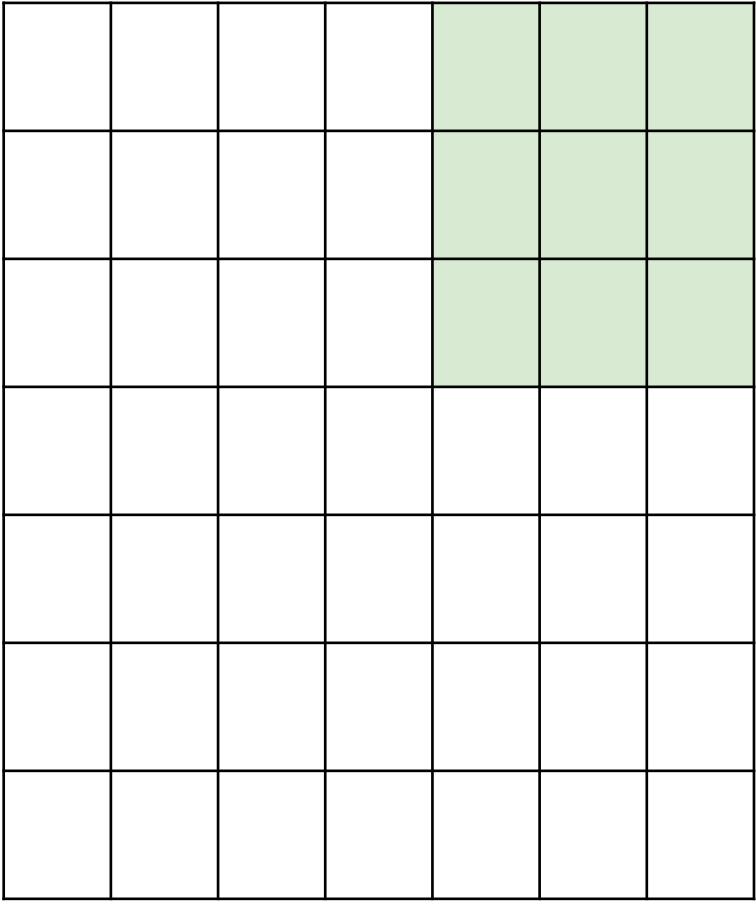


7x7 input (spatially)  
assume 3x3 filter

7

# A closer look at spatial dimensions:

7



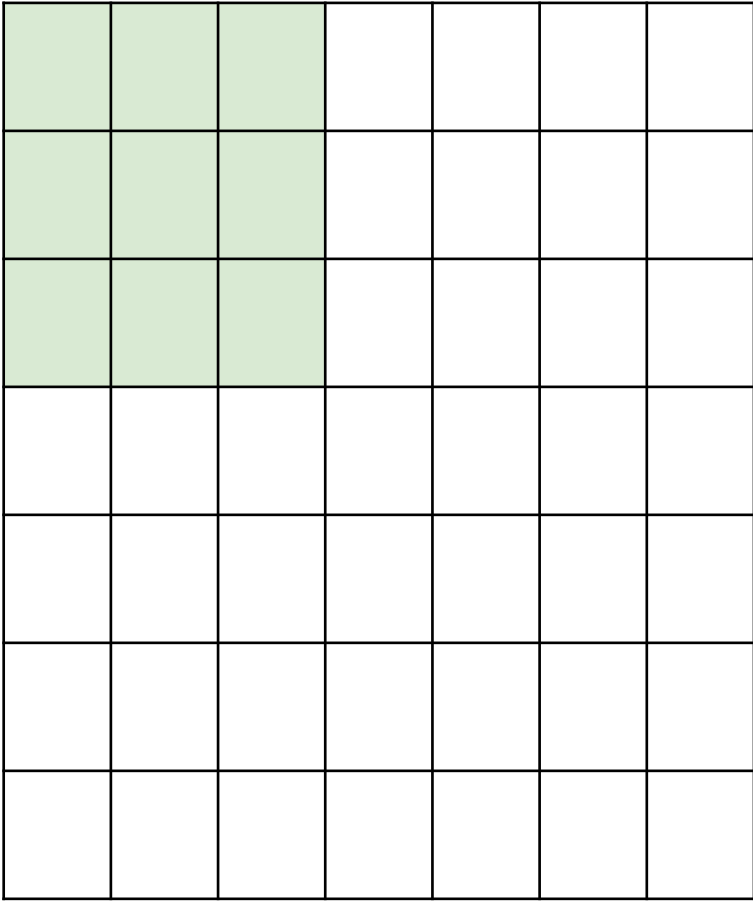
7x7 input (spatially)  
assume 3x3 filter

=> **5x5 output**

7

# A closer look at spatial dimensions:

7

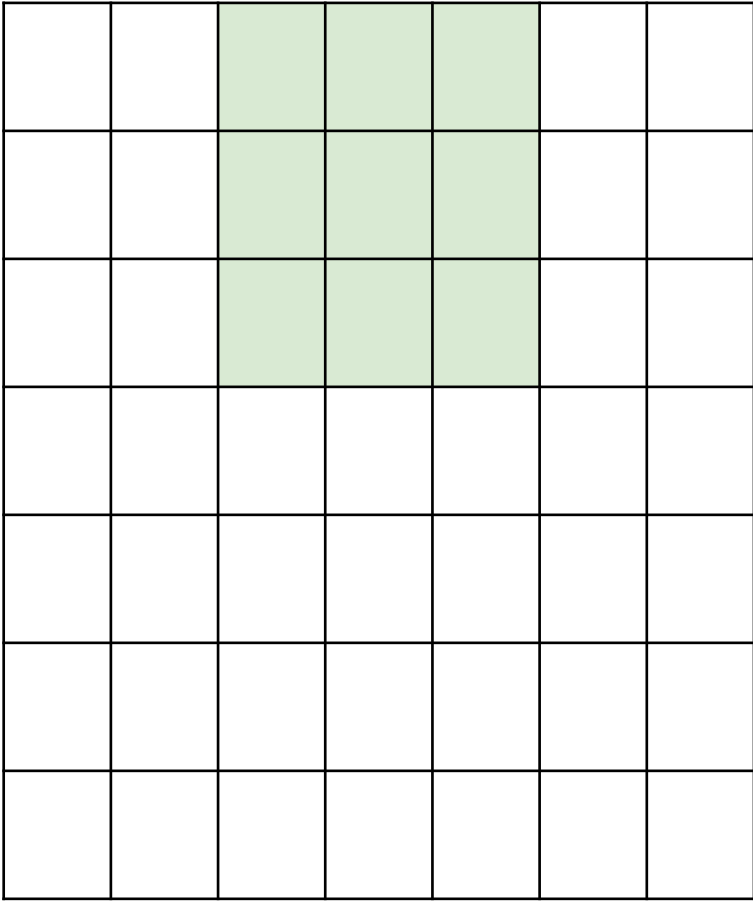


7

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

# A closer look at spatial dimensions:

7

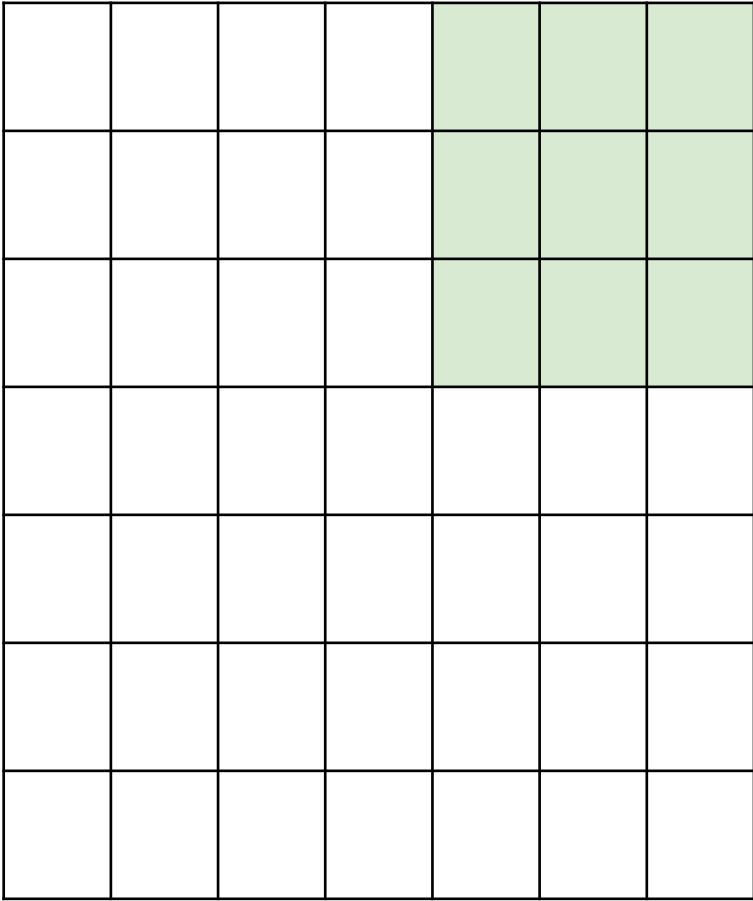


7

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

# A closer look at spatial dimensions:

7



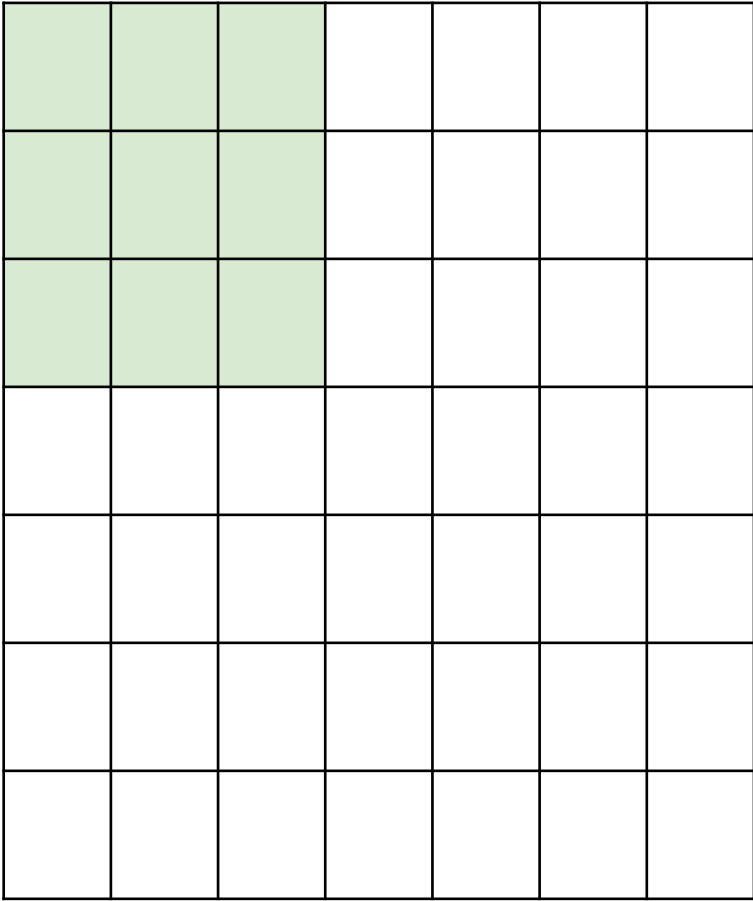
7

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



# A closer look at spatial dimensions:

7

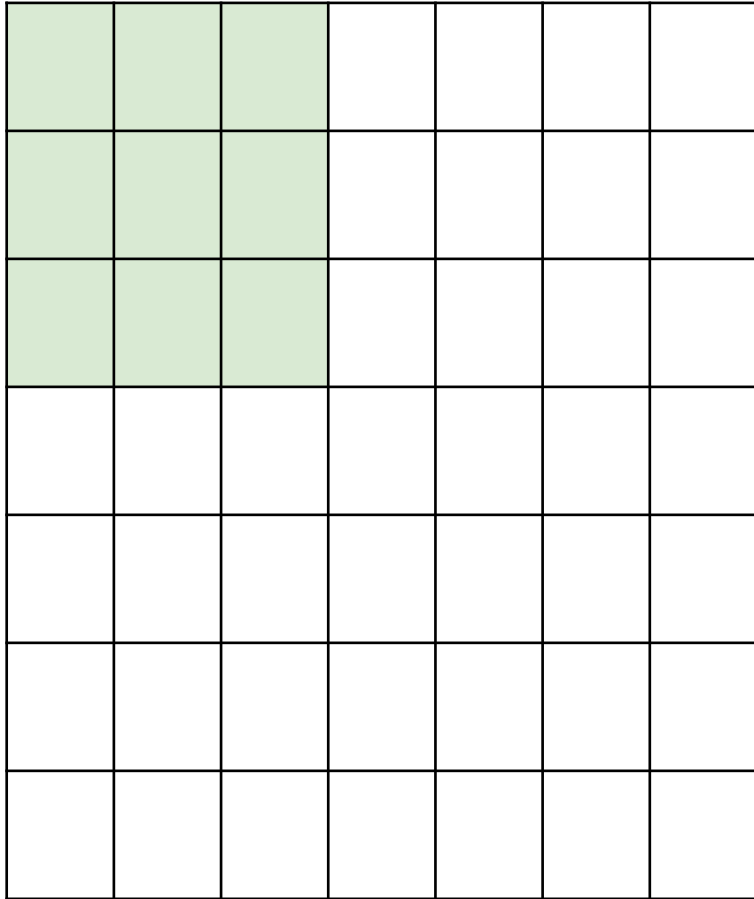


7

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

A closer look at spatial dimensions:

7

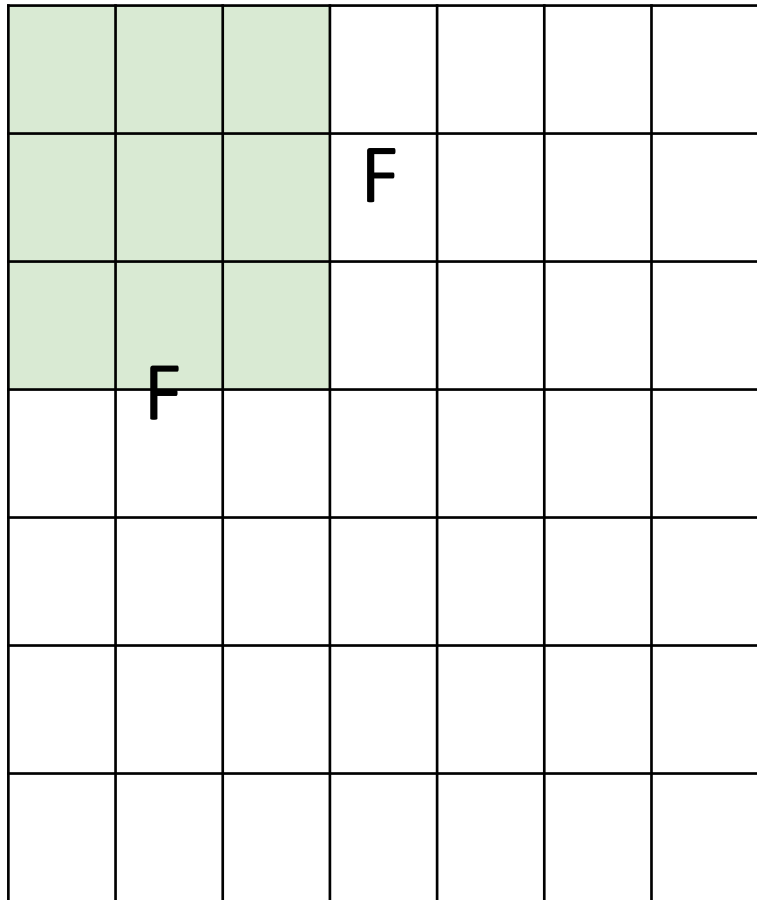


7

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

**doesn't fit!**  
cannot apply 3x3 filter on 7x7  
input with stride 3.

N



Output size:

$$(N - F) / \text{stride} + 1$$

e.g.  $N = 7, F = 3$ :

$$\text{stride } 1 \Rightarrow (7 - 3) / 1 + 1 = 5$$

$$\text{stride } 2 \Rightarrow (7 - 3) / 2 + 1 = 3$$

$$\text{stride } 3 \Rightarrow (7 - 3) / 3 + 1 = 2.33 \text{ :}\backslash$$

# 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:)

$$(N - F) / \text{stride} + 1$$

# 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?

**7x7** output!

# 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?

**7x7 output!**

in general, common to see CONV layers with stride 1, filters of size  $F \times F$ , and zero-padding with  $(F-1)/2$ . (will preserve size spatially)

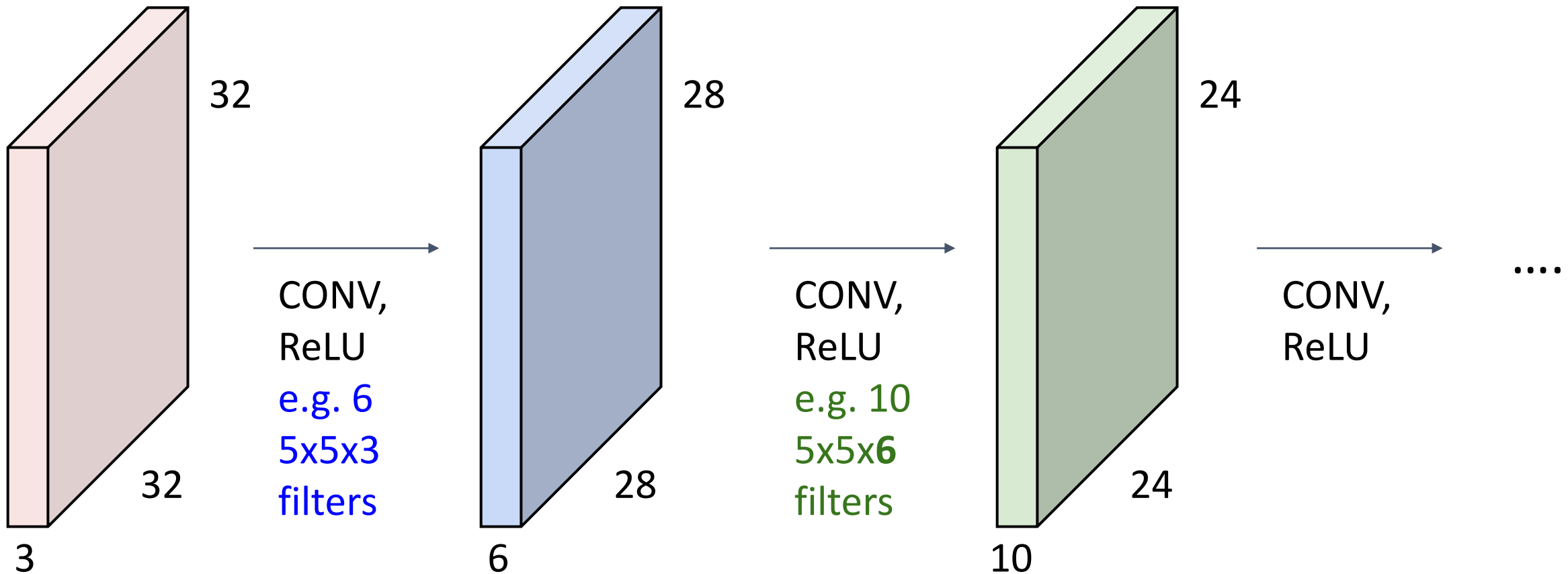
e.g.  $F = 3 \Rightarrow$  zero pad with 1

$F = 5 \Rightarrow$  zero pad with 2

$F = 7 \Rightarrow$  zero pad with 3

## Remember back to...

E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially! (32  $\rightarrow$  28  $\rightarrow$  24 ...). Shrinking too fast is not good, doesn't work well.

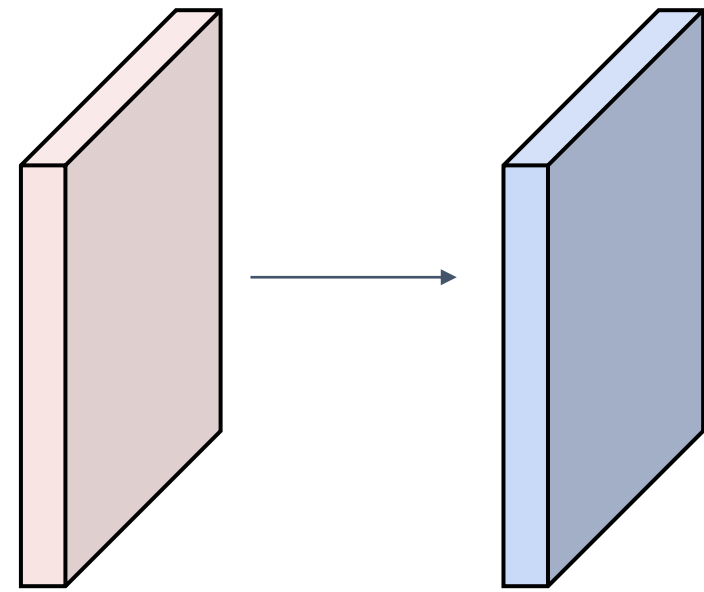


Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

Output volume size: ?





Examples time:

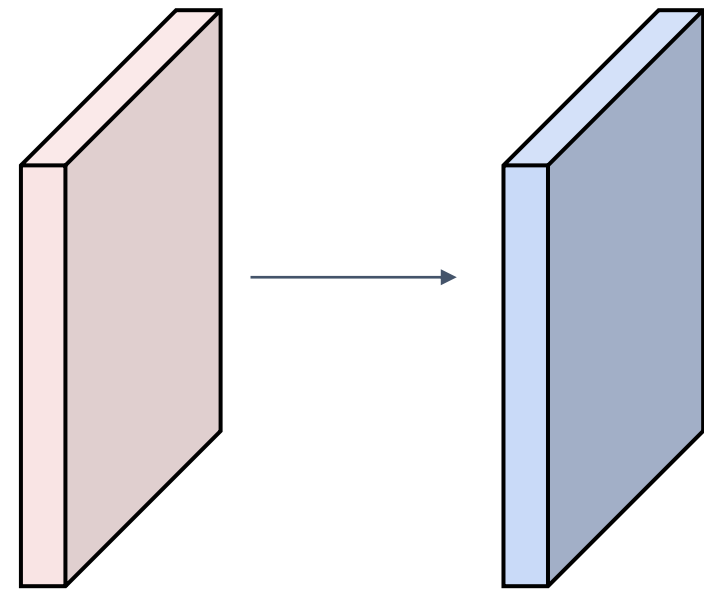
Input volume: **32x32x3**

**10** **5x5** filters with stride **1**, pad **2**

Output volume size:

$(32+2*2-5)/1+1 = 32$  spatially, so

**32x32x10**

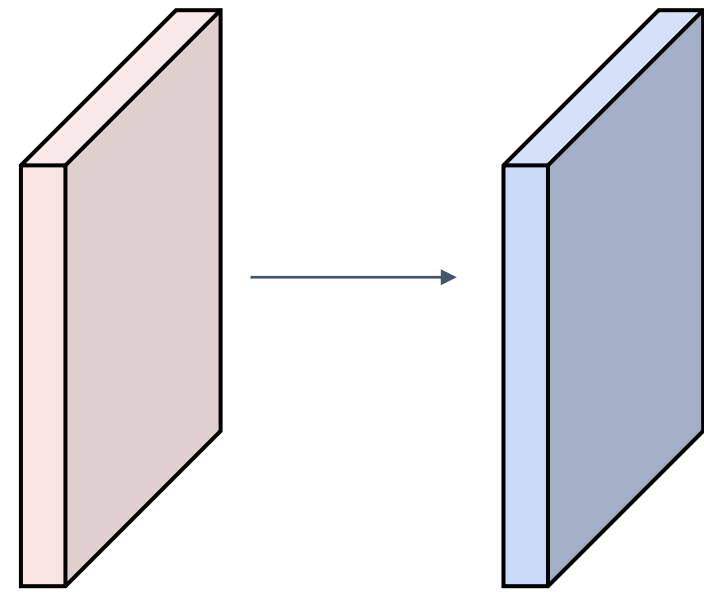


Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

Number of parameters in this layer?



Examples time:

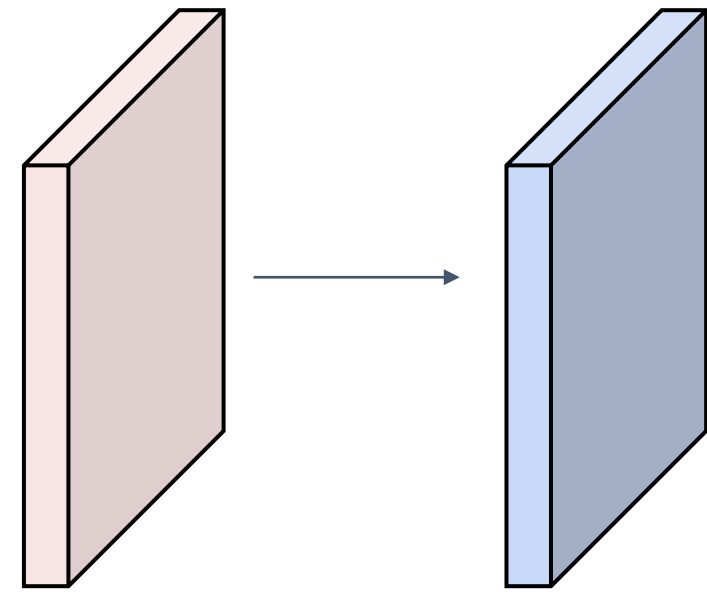
Input volume: **32x32x3**

**10** **5x5** filters with stride 1, pad 2

Number of parameters in this layer?

each filter has  $5*5*3 + 1 = 76$  params

=>  $76*10 = 760$



(+1 for bias)