

# Convolutional Neural Networks

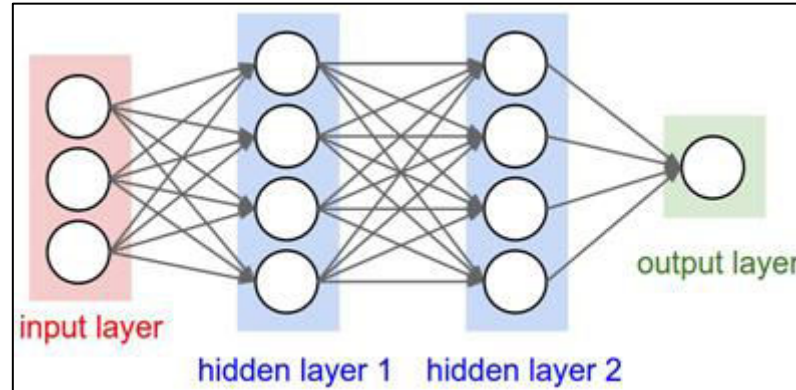
Radu Ionescu, Prof. PhD.  
raducu.ionescu@gmail.com

Faculty of Mathematics and Computer Science  
University of Bucharest

# Mini-batch SGD

Loop:

1. **Sample** a batch of data
2. **Forward** prop it through the graph, get loss
3. **Backprop** to calculate the gradients
4. **Update** the parameters using the gradient



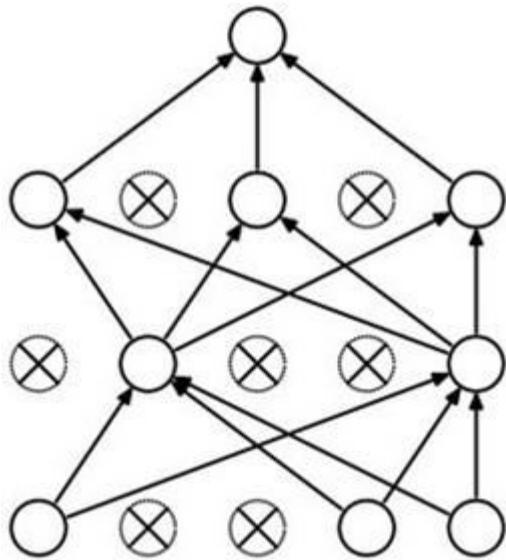
# Neural networks are universal function approximators

- **Universal Approximation Theorem:**

A feed-forward neural network with a hidden layer composed of a finite number of neurons can approximate any continuous function defined on a compact subset of  $\mathbb{R}^n$ .

- Although 2-layer neural networks (1-hidden layer) are universal function approximators, the hidden layer's width (number of neurons) can be exponentially large.
- In practice, we prefer deeper (with more layers) and thinner architectures

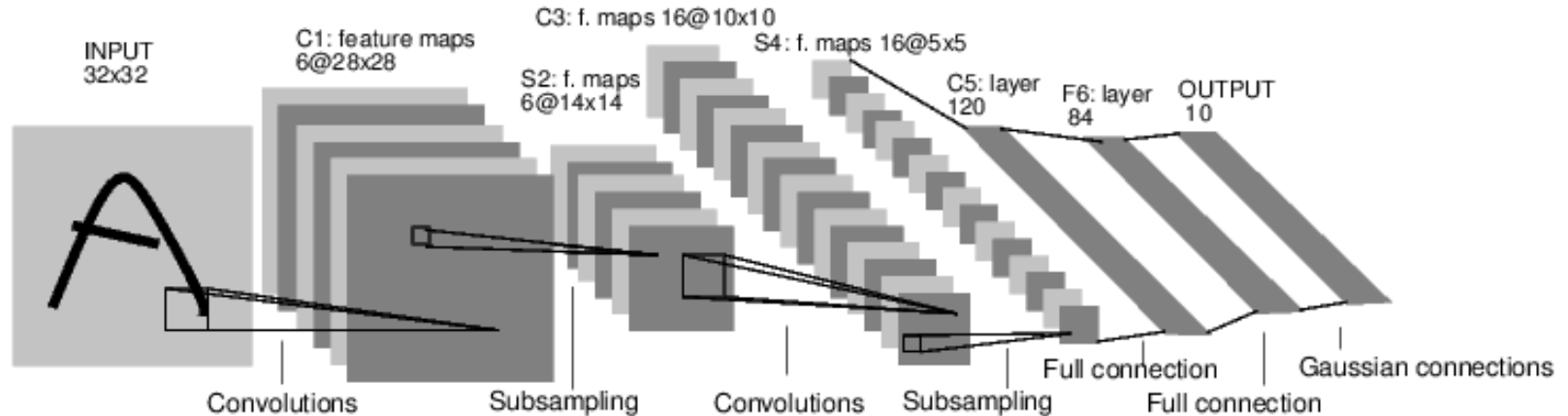
# Dropout



Forces the network to have a redundant representation.

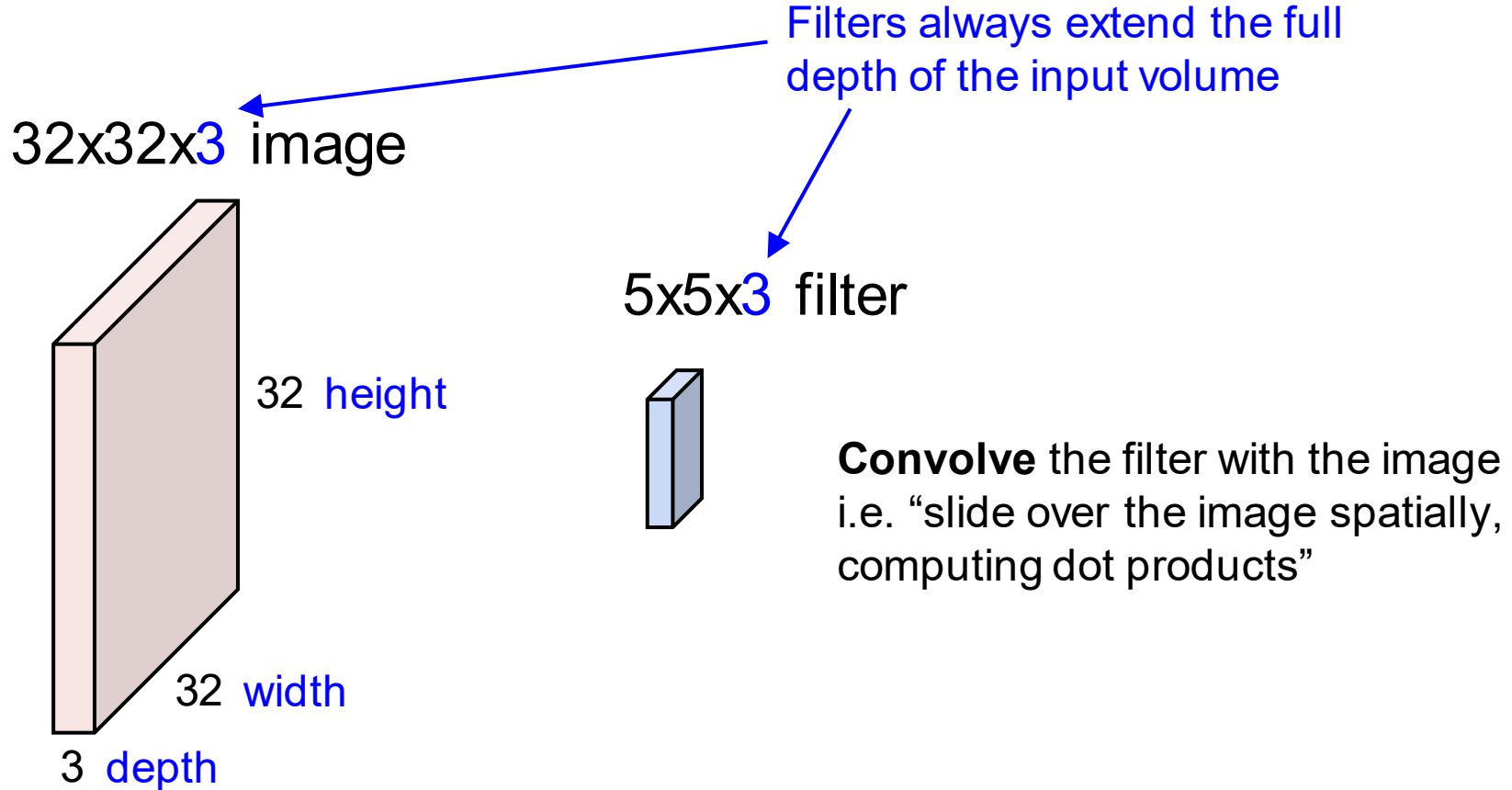


# Convolutional Neural Networks

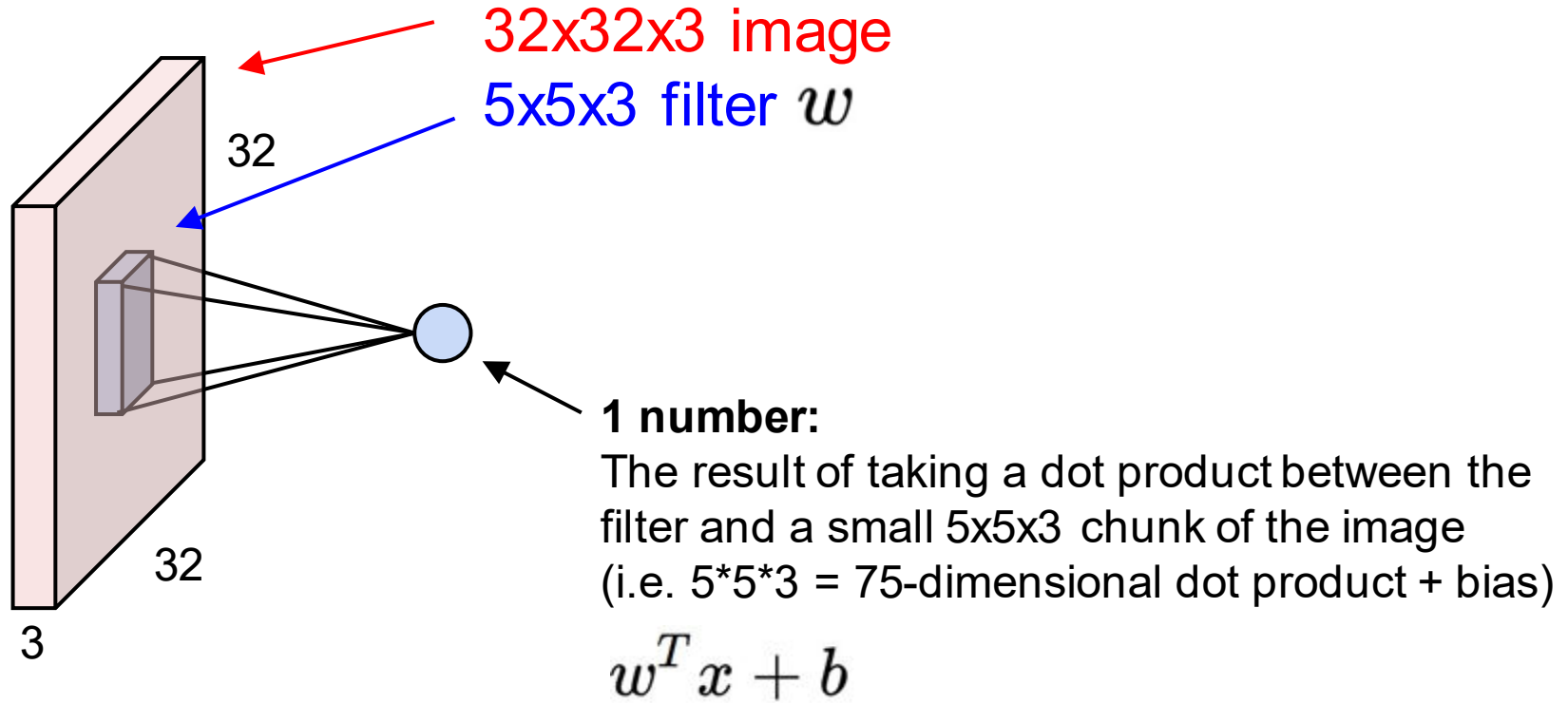


*[LeNet-5, LeCun 1998]*

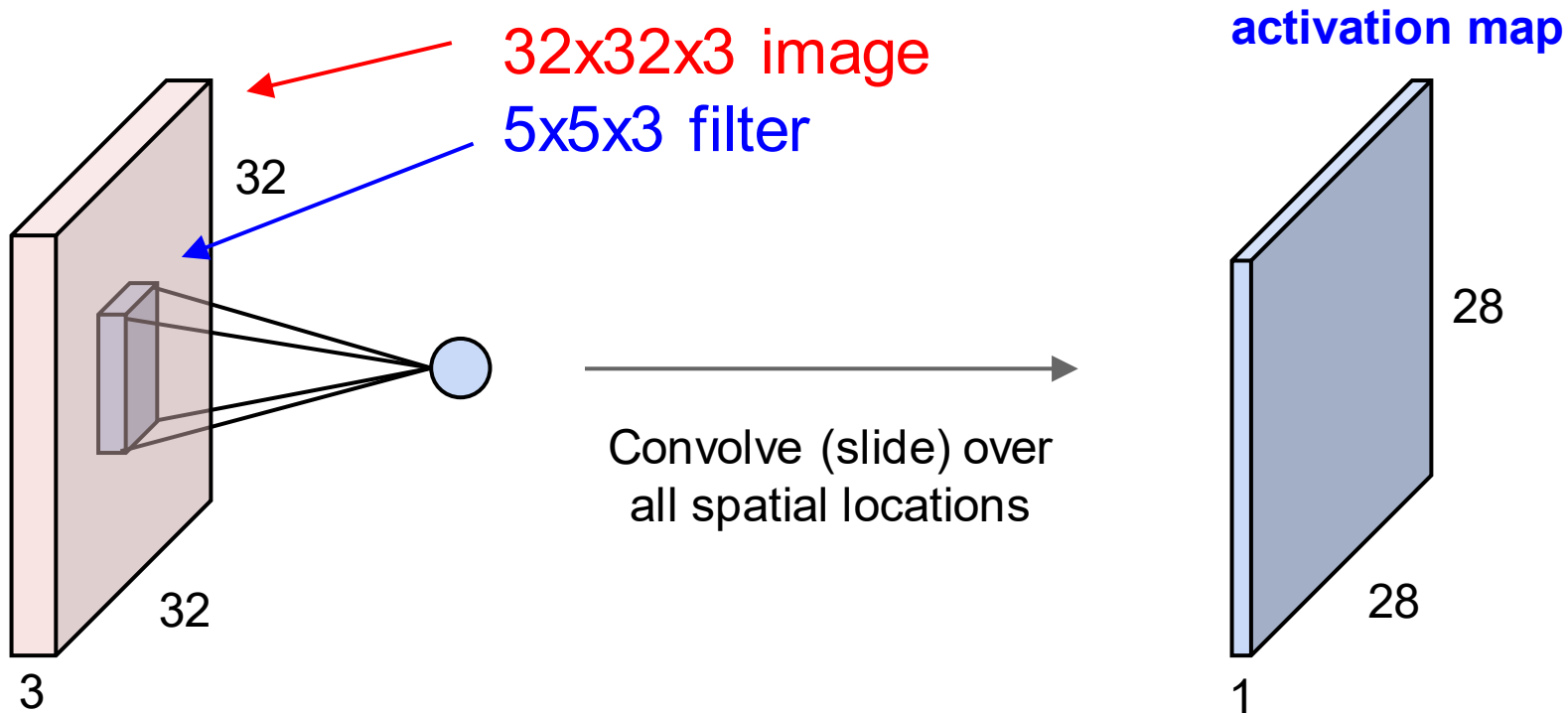
# Convolution Layer



# Convolution Layer

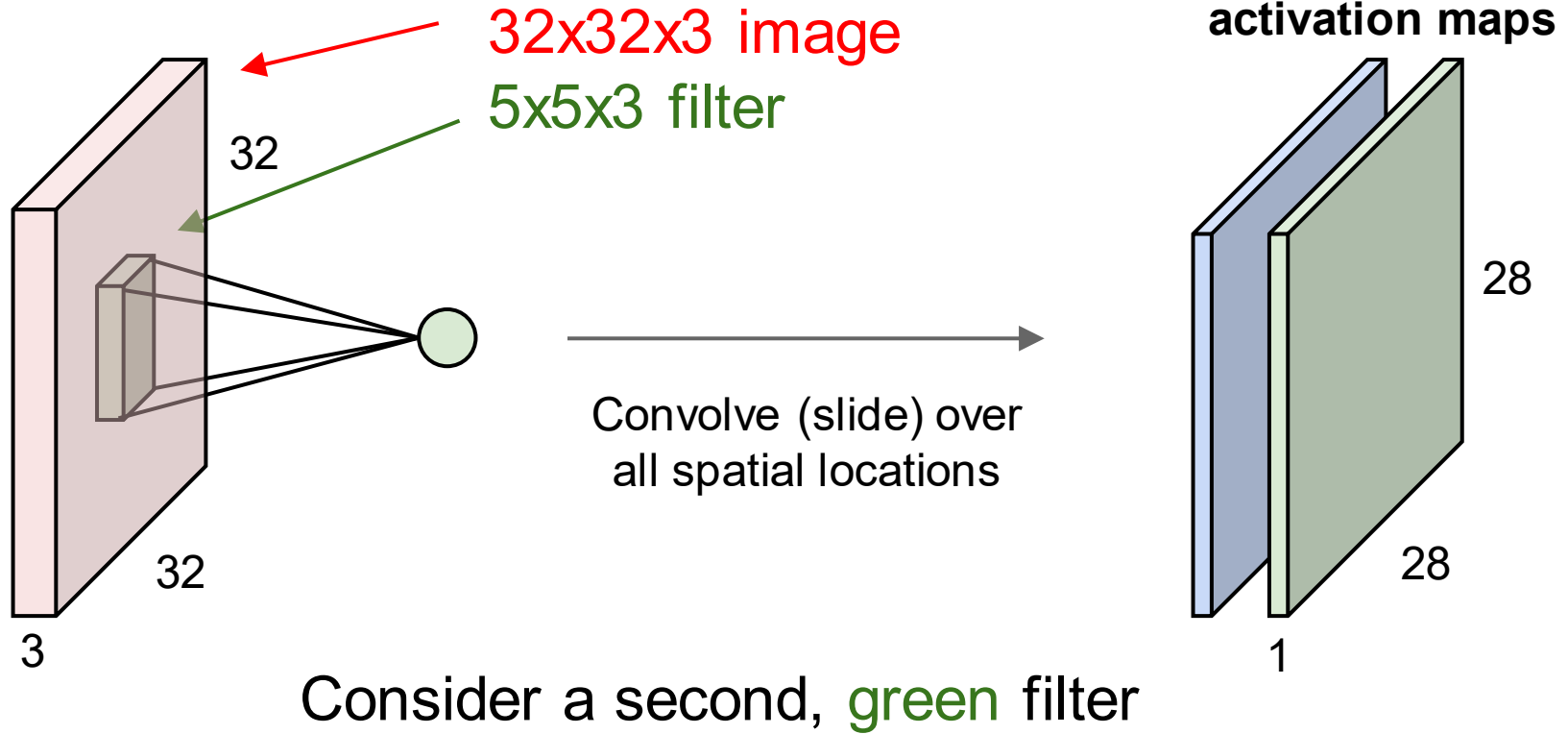


# Convolution Layer

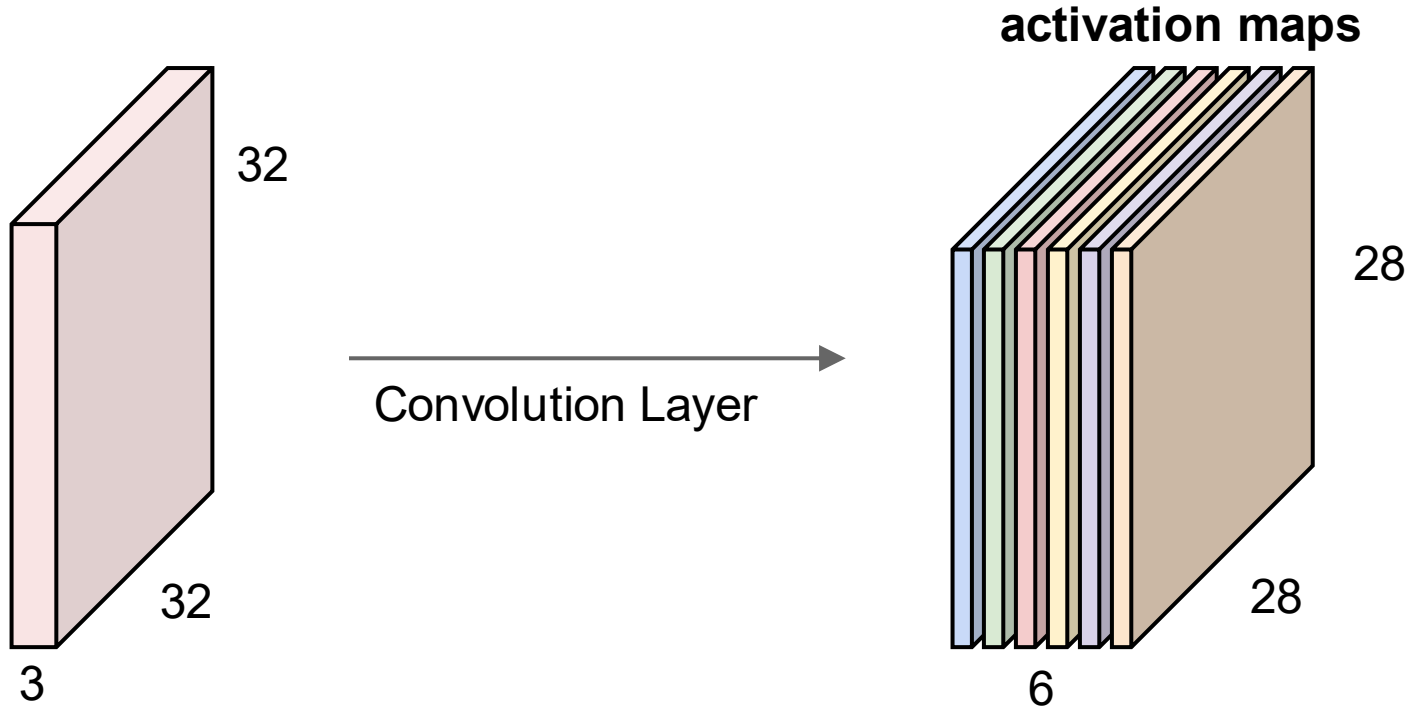




# Convolution Layer

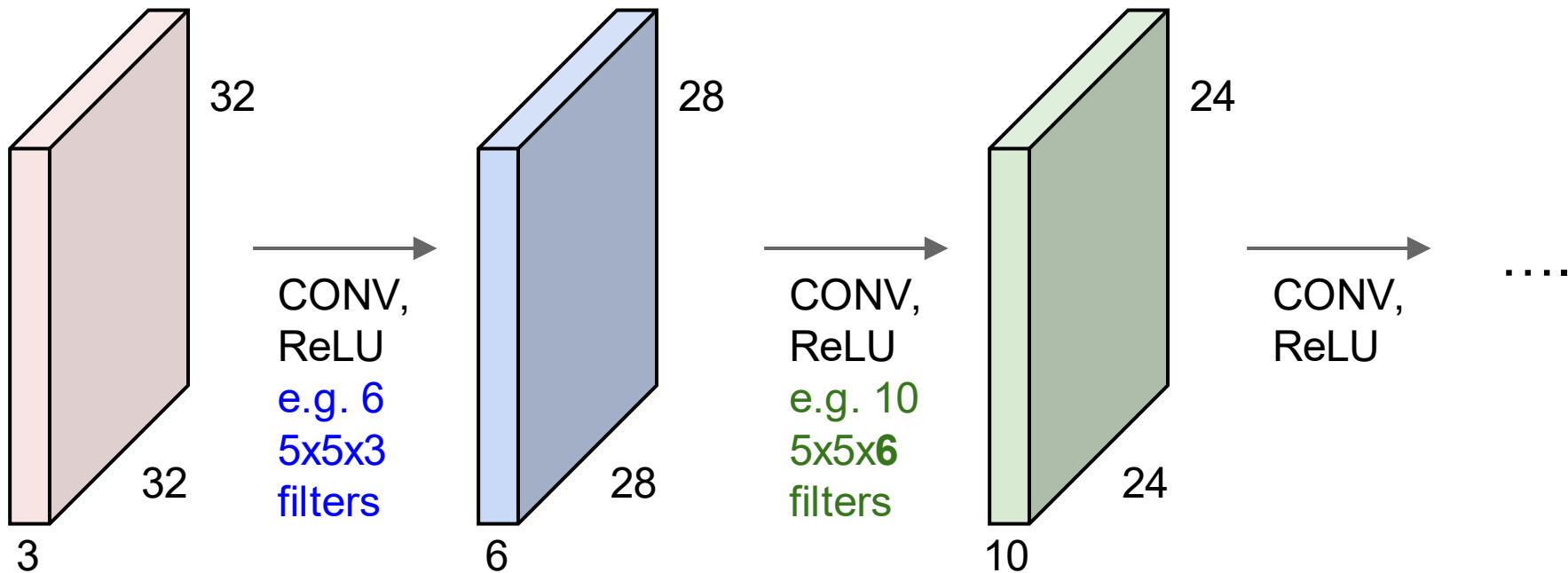


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

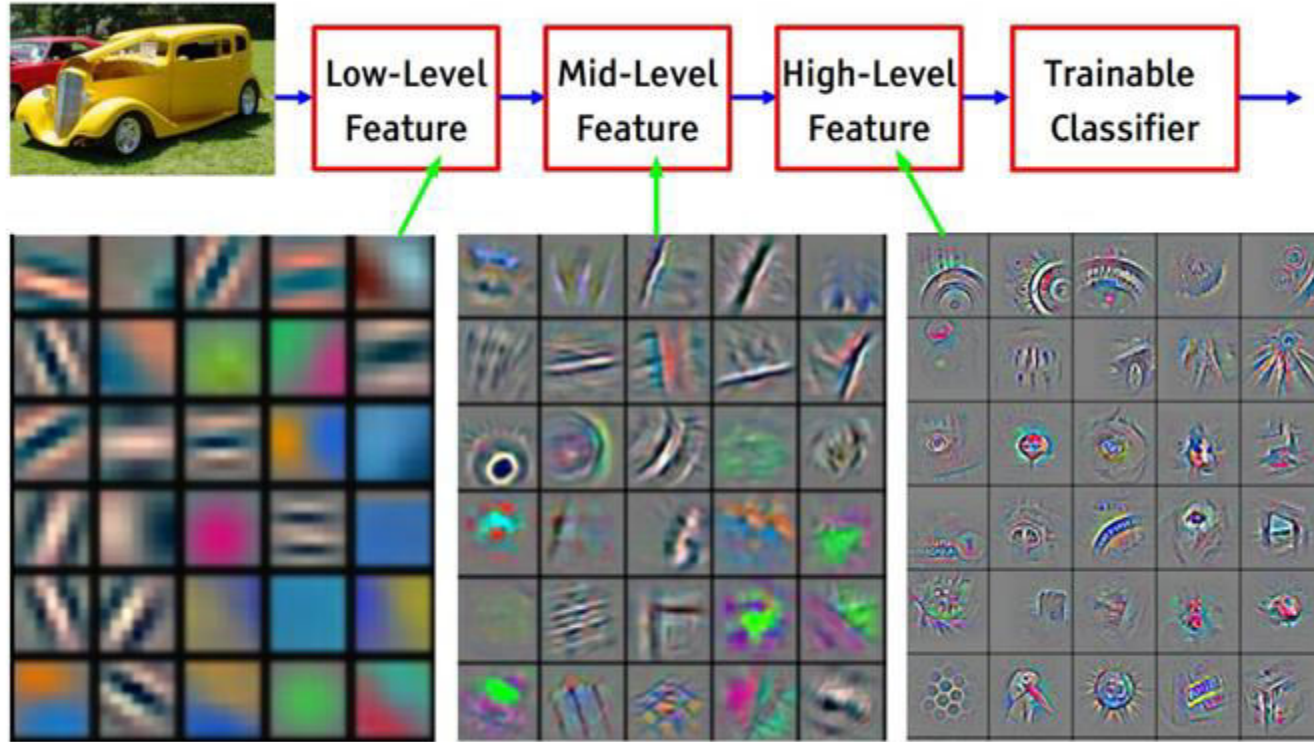


We stack these up to get a “new image” of size 28x28x6!

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



# Filters correspond to features / parts of the objects



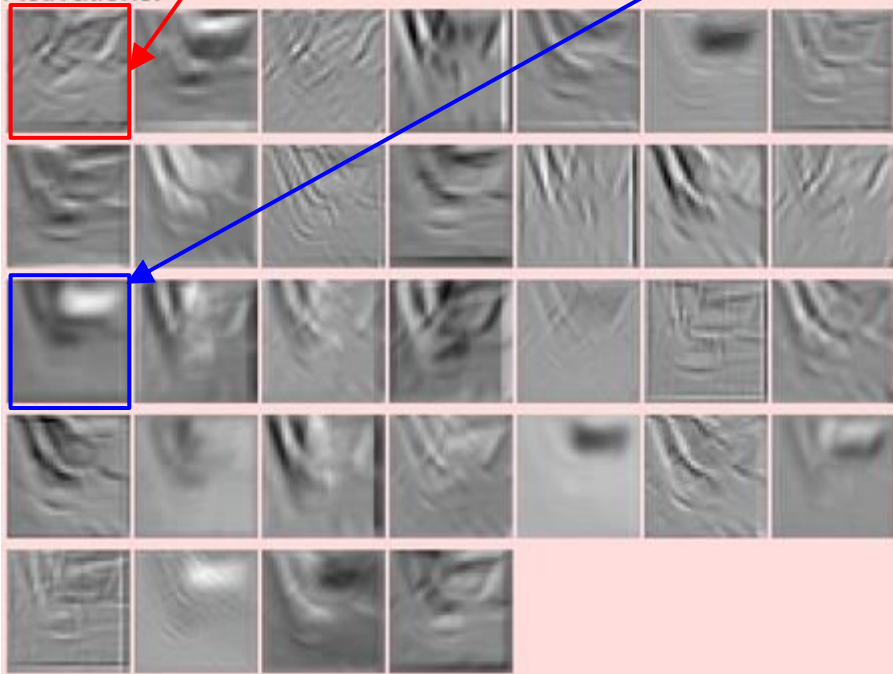
Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]



one filter =>  
one activation map

example 5x5 filters  
(32 total)

Activations:



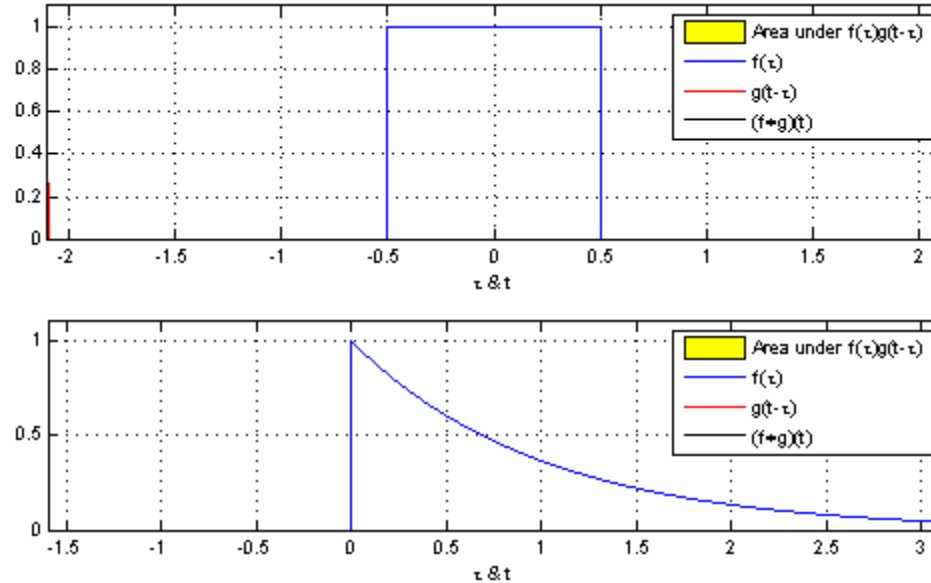
We call the layer convolutional  
because it is related to convolution  
of two signals:

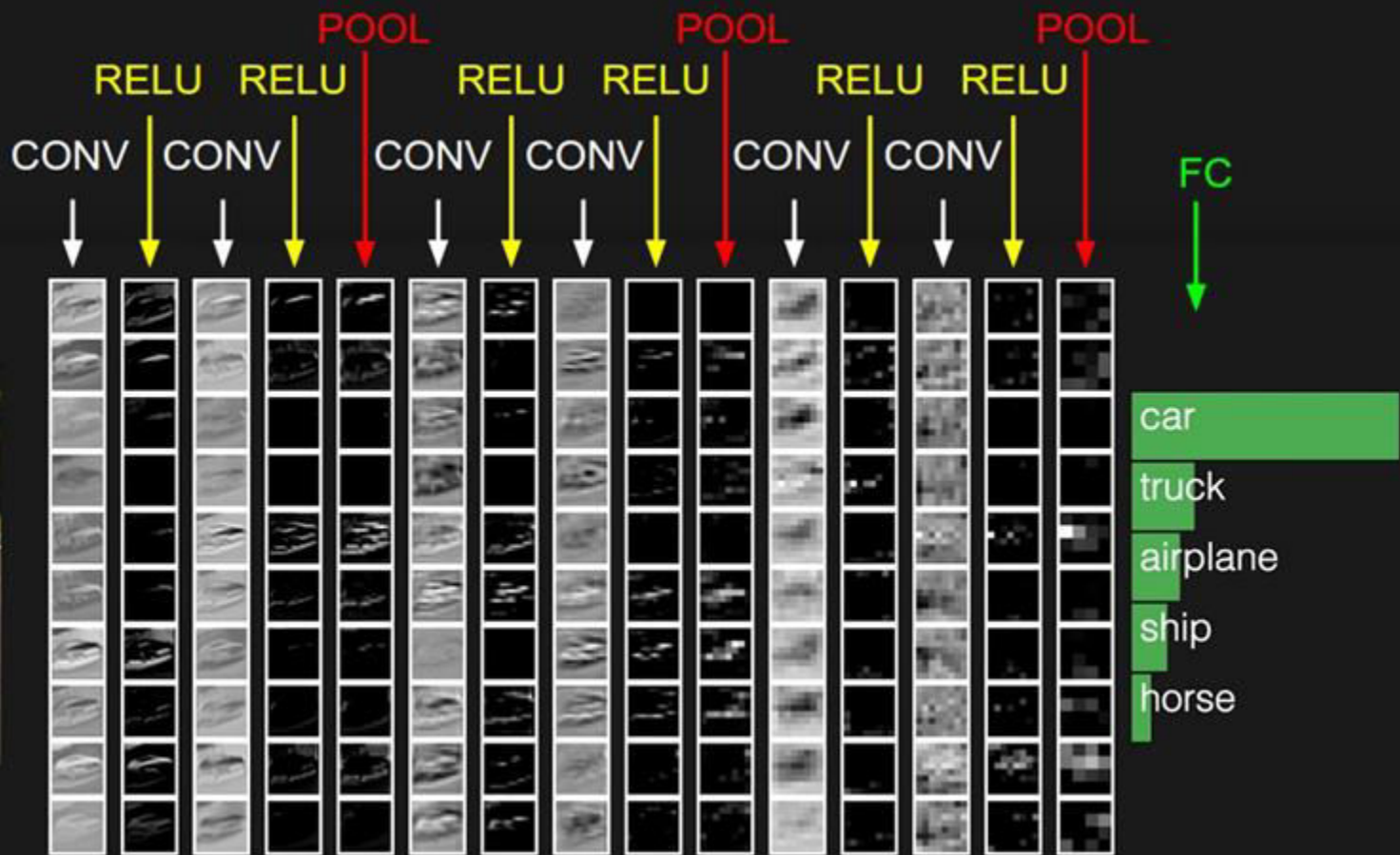
$$f[x,y] * g[x,y] = \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1,n_2] \cdot g[x-n_1,y-n_2]$$



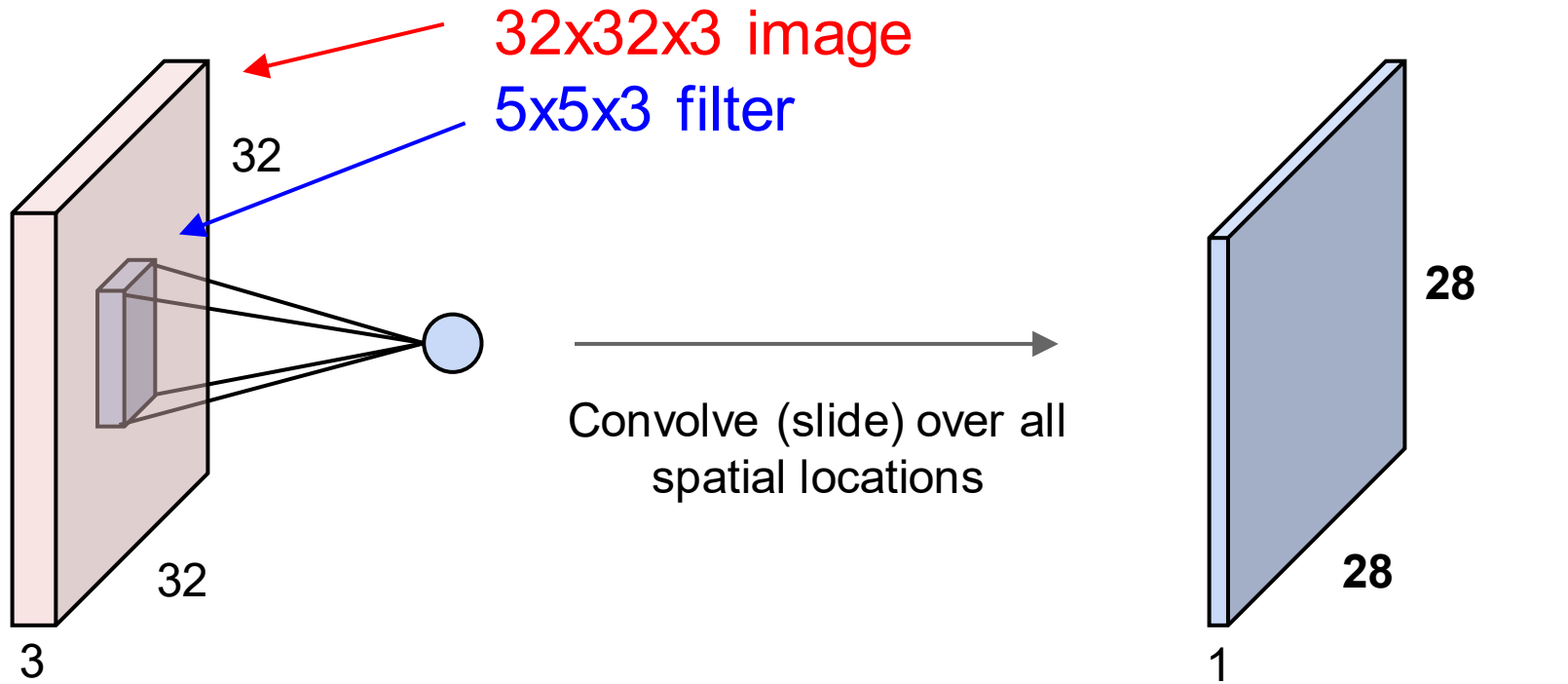
elementwise multiplication and sum of  
a filter and the signal (image)

# Maximum activation = highest response





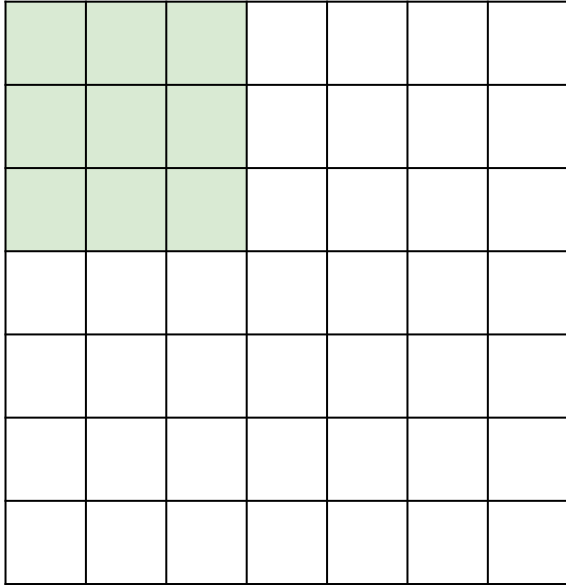
# Size of an activation map





# Size of an activation map

7

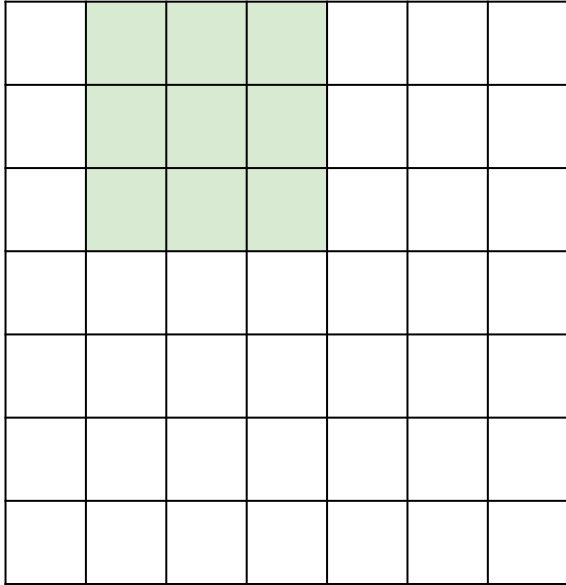


7

7x7 input (spatially, without depth)  
assume 3x3 filter

# Size of an activation map

7

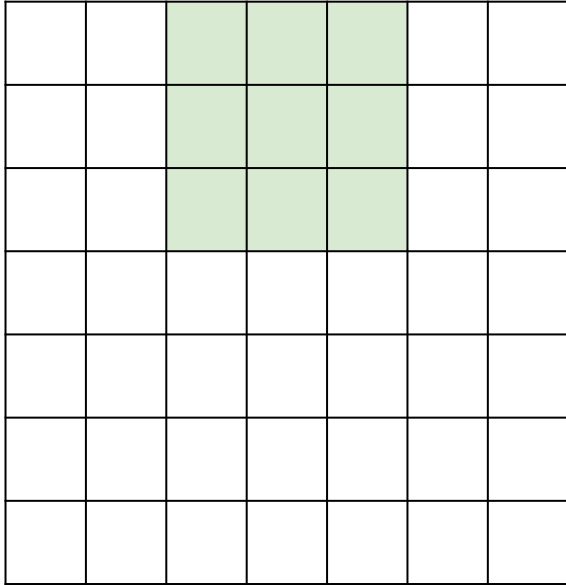


7

7x7 input (spatially, without depth)  
assume 3x3 filter

# Size of an activation map

7

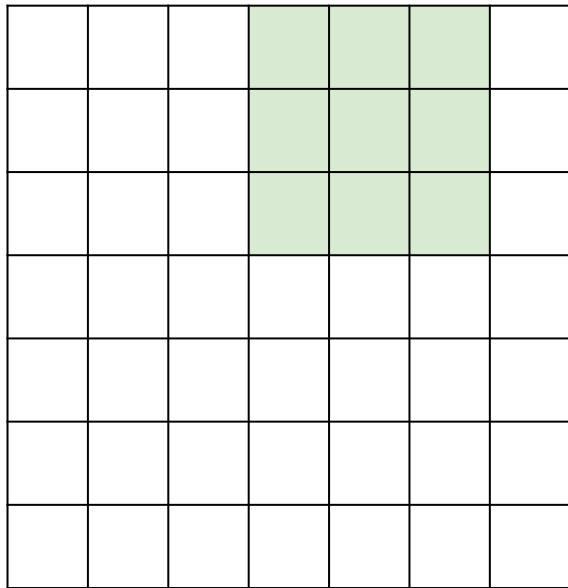


7

7x7 input (spatially, without depth)  
assume 3x3 filter

# Size of an activation map

7



7

7x7 input (spatially, without depth)  
assume 3x3 filter

# Size of an activation map

7

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

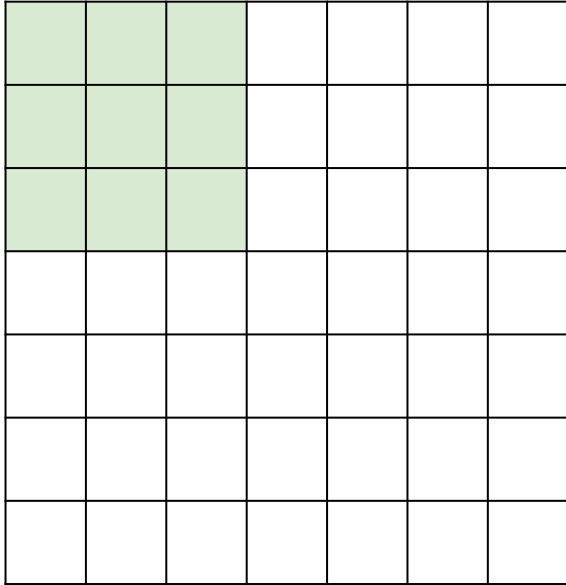
7

7x7 input (spatially, without depth)  
assume 3x3 filter

**=> 5x5 output**

# Size of an activation map

7

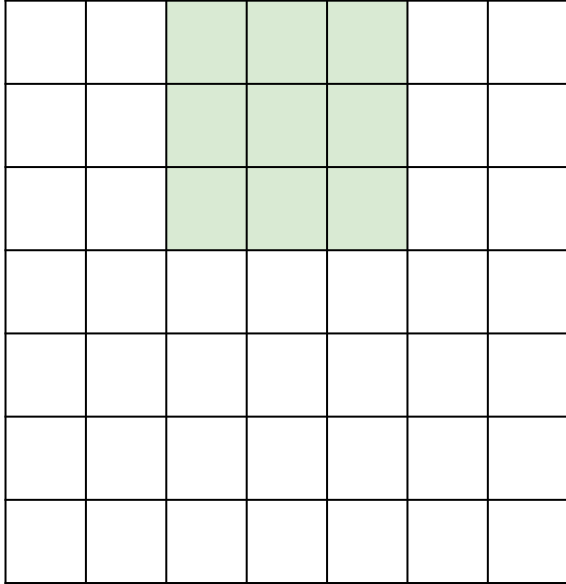


7

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

# Size of an activation map

7



7

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

# Size of an activation map

7

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

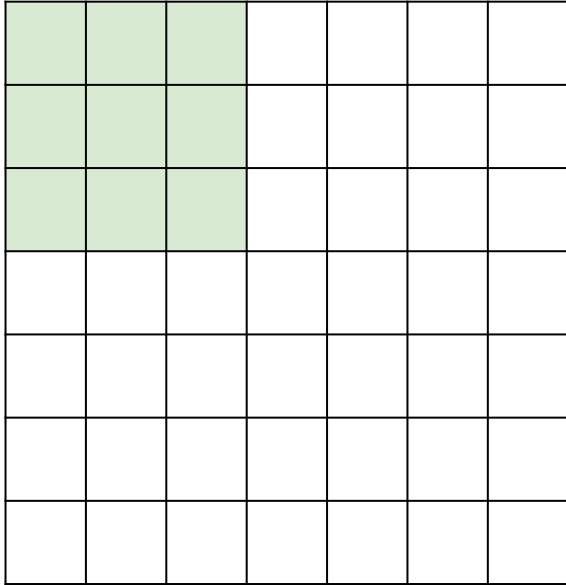
7

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



# Size of an activation map

7

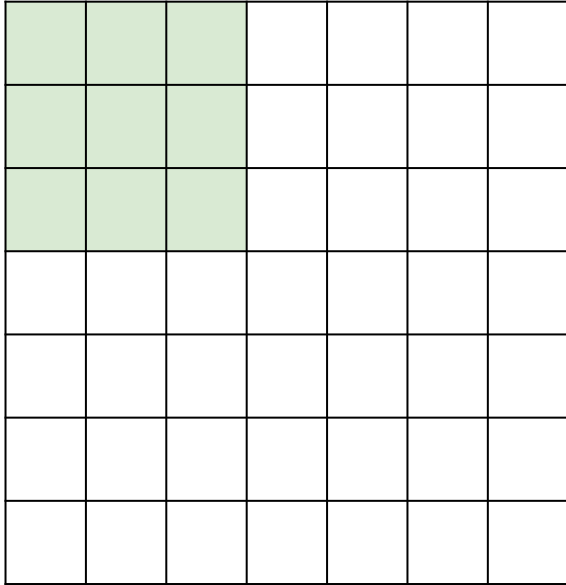


7

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

# Size of an activation map

7



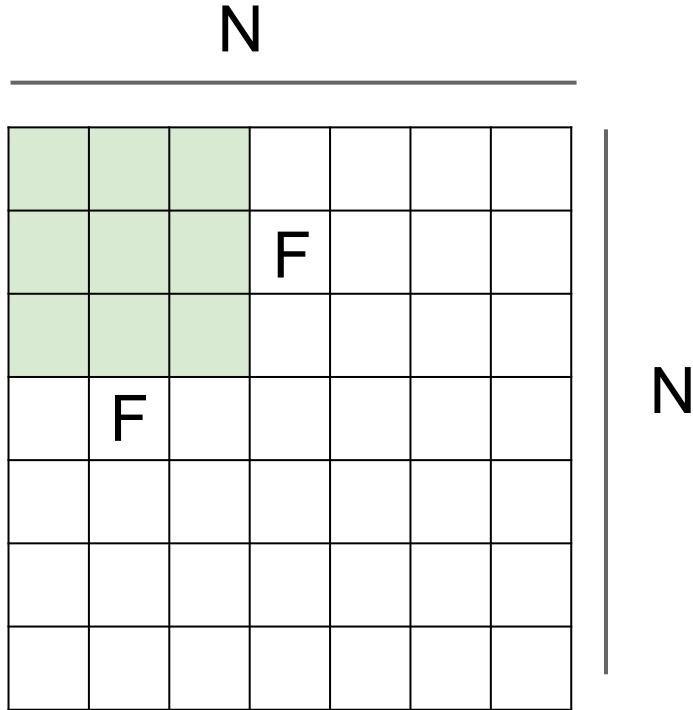
7

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

**It does not fit!**

We cannot apply a 3x3 filter  
on an image of 7x7 pixels  
using a stride of 3

# Size of an activation map



Output size:

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

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

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

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

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

# 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!**

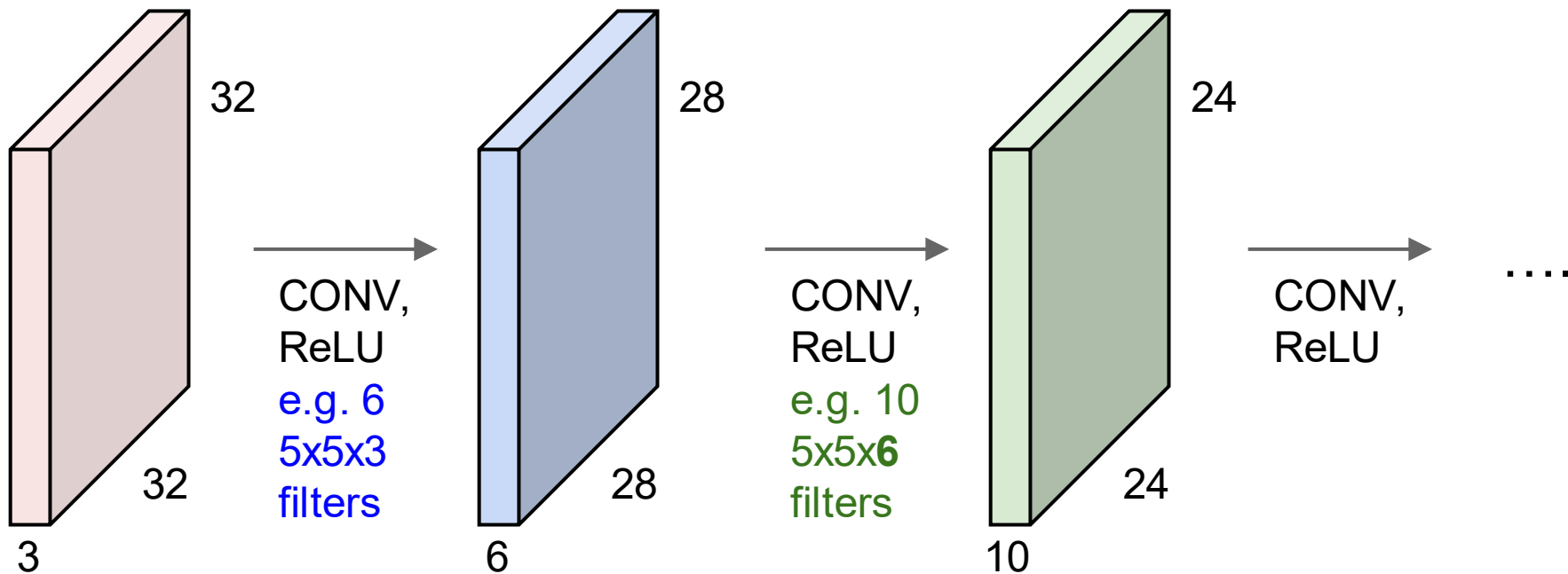
In general, common to see CONV layers with stride 1, filters of size  $F \times F$ , and zero-padding with  $(F-1)/2$ . (preserves the size of the input image / activation map)

e.g.  $F = 3 \Rightarrow$  zero pad with 1 (value 0)

$F = 5 \Rightarrow$  zero pad with 2 (value 0)

$F = 7 \Rightarrow$  zero pad with 3 (value 0)

- 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially!  
(32 => 28 => 24 ...)
- Shrinking too fast is not good, doesn't work well.

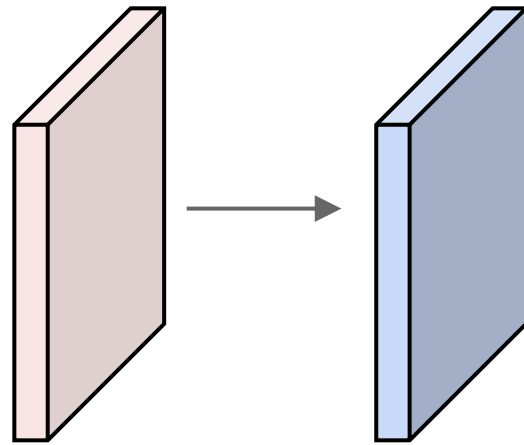


Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

Output volume size: ?





Examples time:

Input volume: **32x32x3**

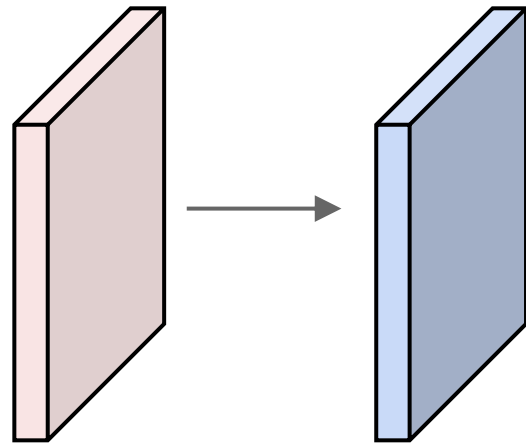
**10** **5x5** filters with stride **1**, pad **2**

Output volume size:

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

so the volume is

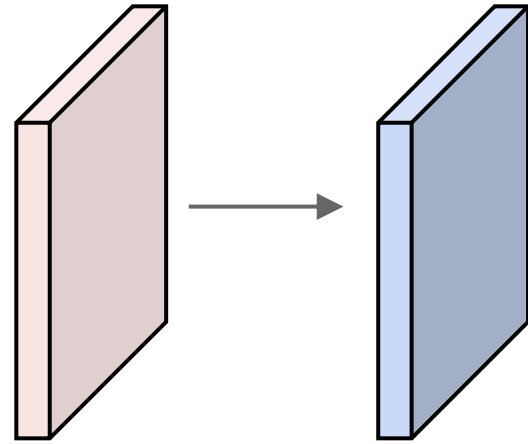
**32x32x10**



Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

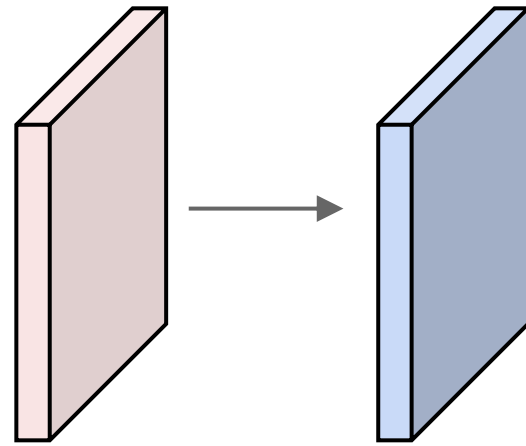


Number of parameters in this layer?

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)

=>  $76*10 = 760$

## Common settings:

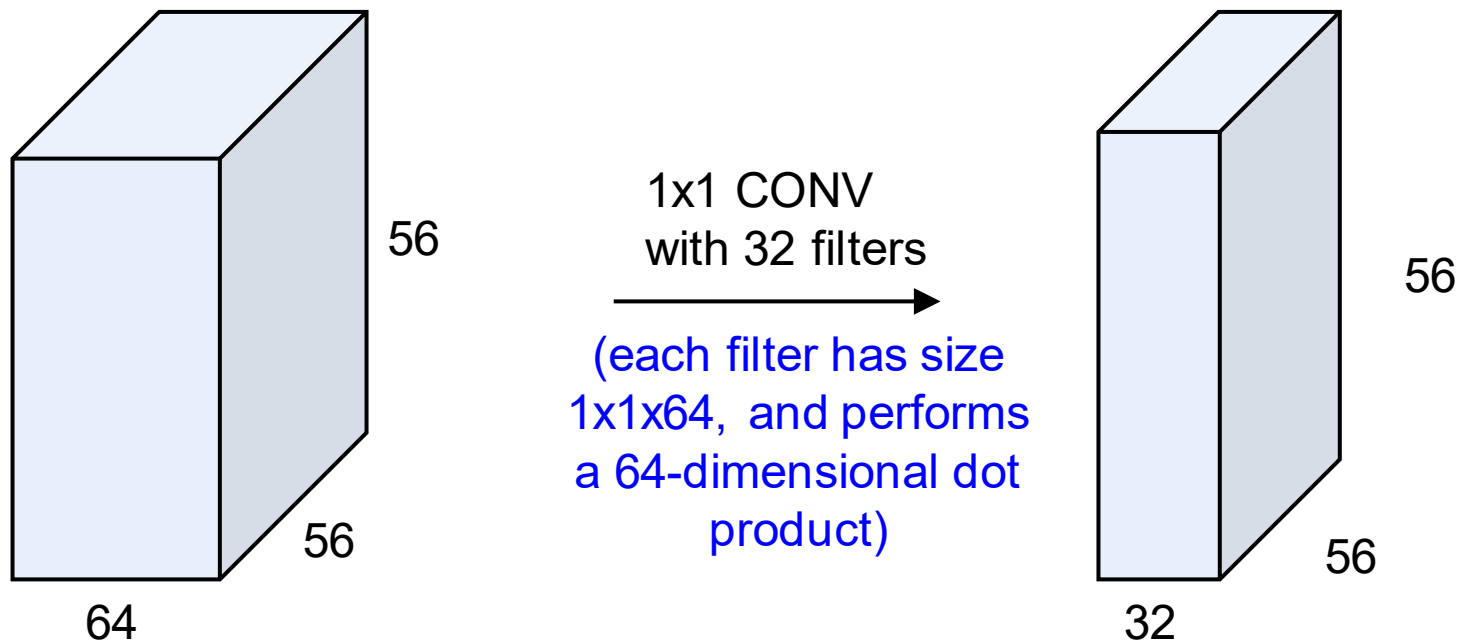
$K$  = (powers of 2, e.g. 32, 64, 128, 512)

- $F = 3, S = 1, P = 1$
- $F = 5, S = 1, P = 2$
- $F = 5, S = 2, P = ?$  (whatever fits)
- $F = 1, S = 1, P = 0$

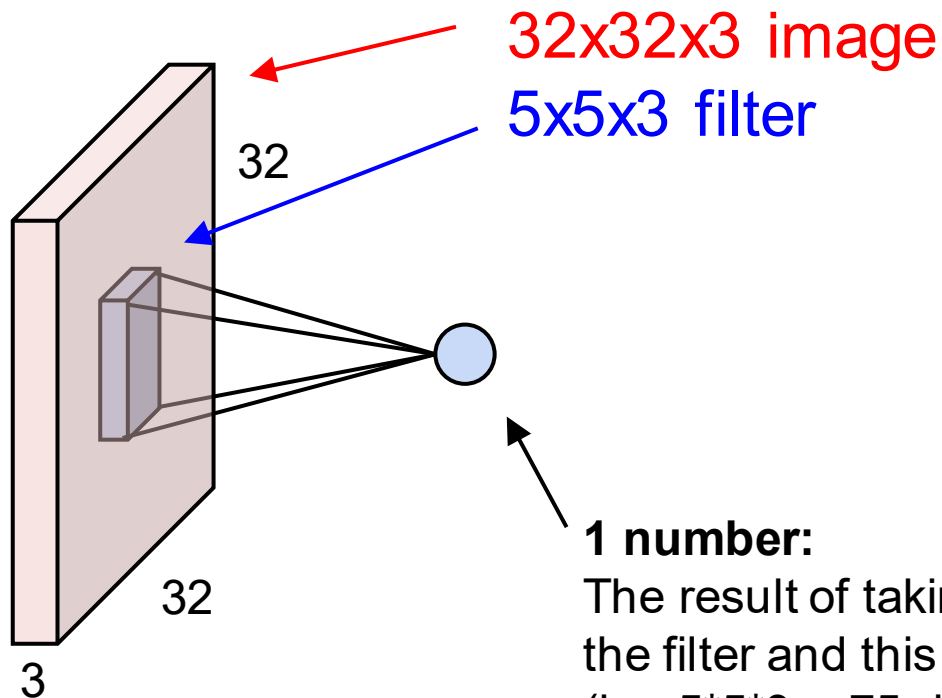
**Summary.** To summarize, the Conv Layer:

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
  - Number of filters  $K$ ,
  - their spatial extent  $F$ ,
  - the stride  $S$ ,
  - the amount of zero padding  $P$ .
- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F + 2P)/S + 1$
  - $H_2 = (H_1 - F + 2P)/S + 1$  (i.e. width and height are computed equally by symmetry)
  - $D_2 = K$
- With parameter sharing, it introduces  $F \cdot F \cdot D_1$  weights per filter, for a total of  $(F \cdot F \cdot D_1) \cdot K$  weights and  $K$  biases.
- In the output volume, the  $d$ -th depth slice (of size  $W_2 \times H_2$ ) is the result of performing a valid convolution of the  $d$ -th filter over the input volume with a stride of  $S$ , and then offset by  $d$ -th bias.

Btw, 1x1 convolutional layers make perfect sense

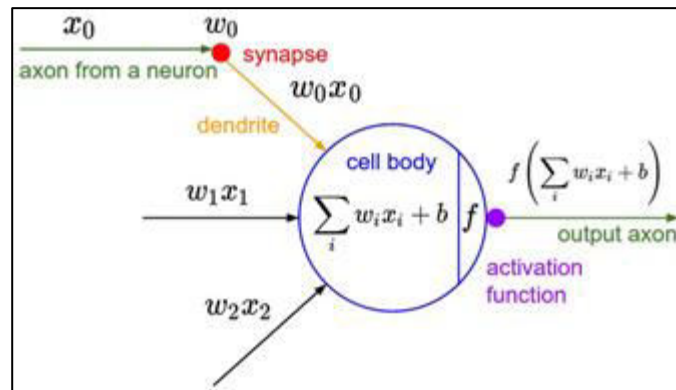


# The brain/neuron view of CONV Layer



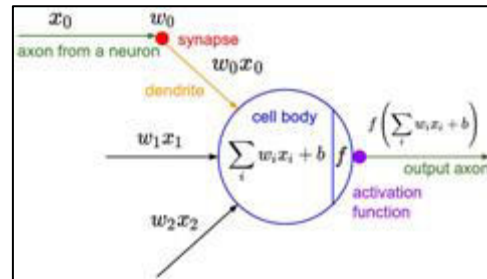
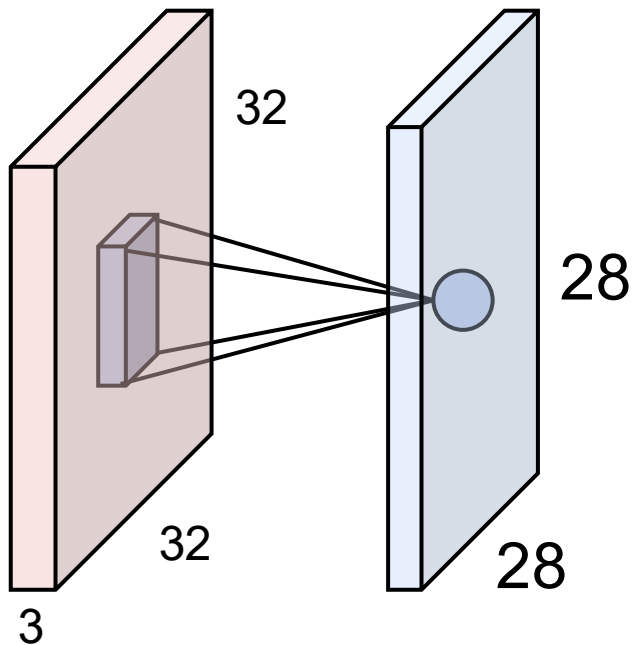
The result of taking a dot product between the filter and this part of the image (i.e.  $5 \times 5 \times 3 = 75$ -dimensional dot product + bias)

$$w^T x + b$$



It's just a neuron with local connectivity...

# The brain/neuron view of CONV Layer

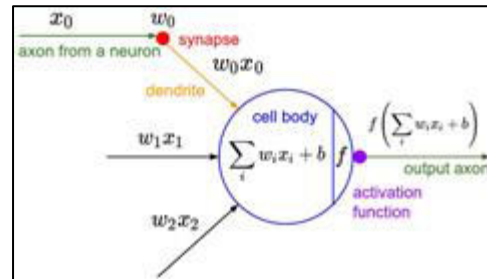
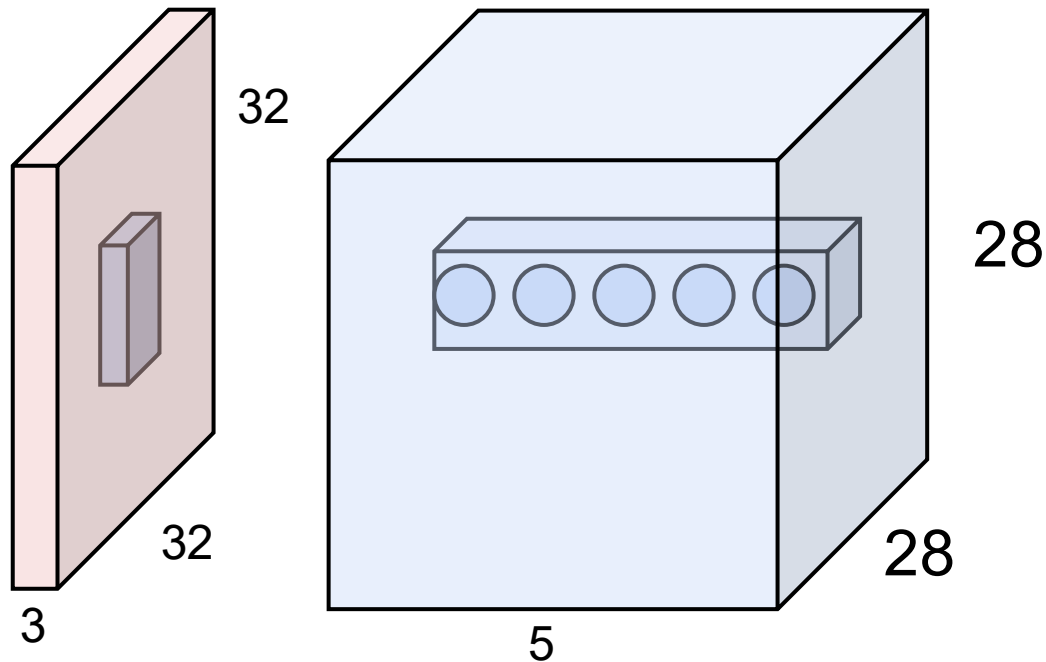


An activation map is a 28x28 sheet of neuron outputs:

1. Each is connected to a small region in the input
2. All of them share parameters

5x5 filter => 5x5 receptive field for each neuron

## The brain/neuron view of CONV Layer

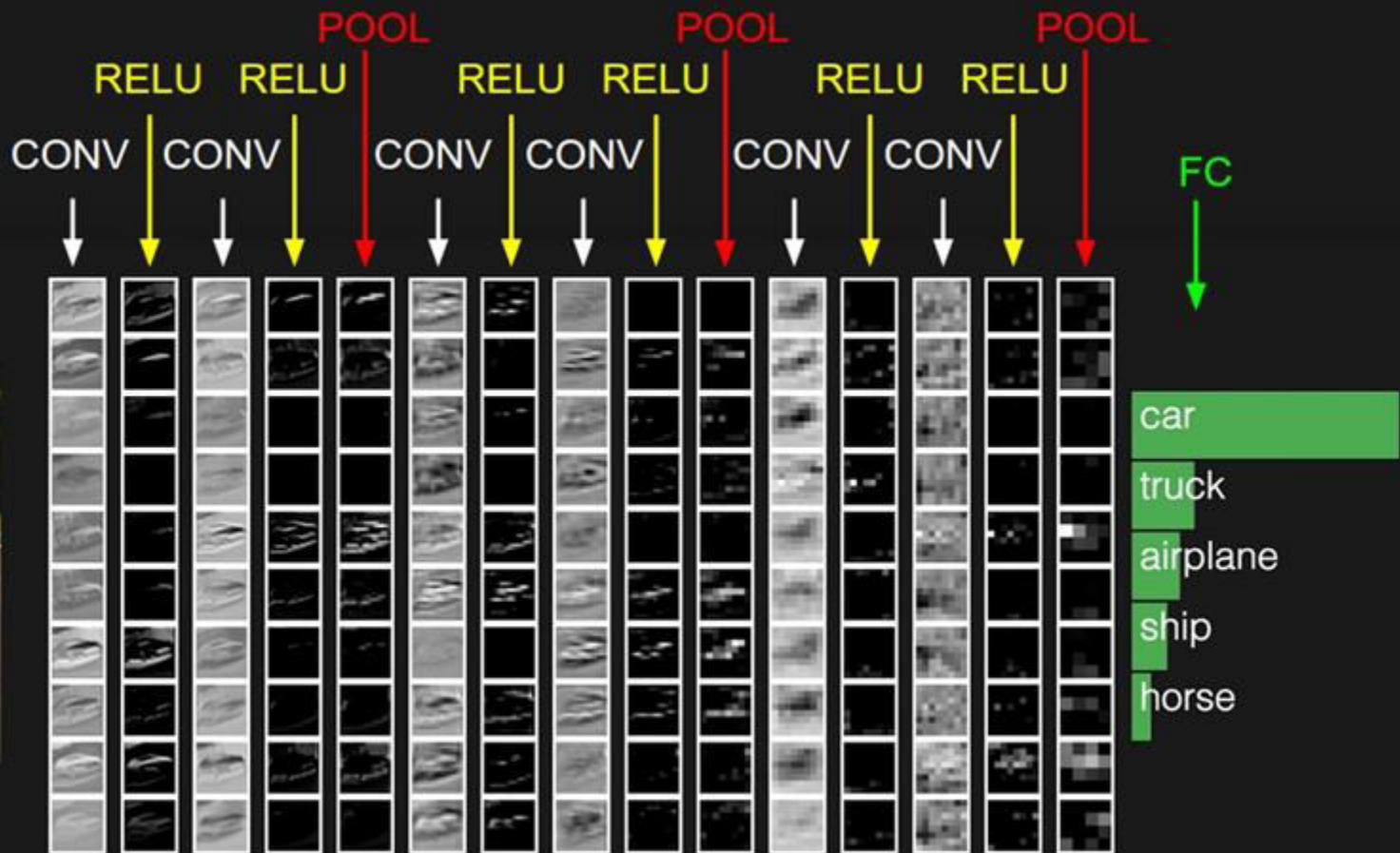


E.g. with 5 filters,  
CONV layer consists of  
neurons arranged in a 3D grid  
(28x28x5)

There will be 5 different neurons all looking at the same region in the input volume

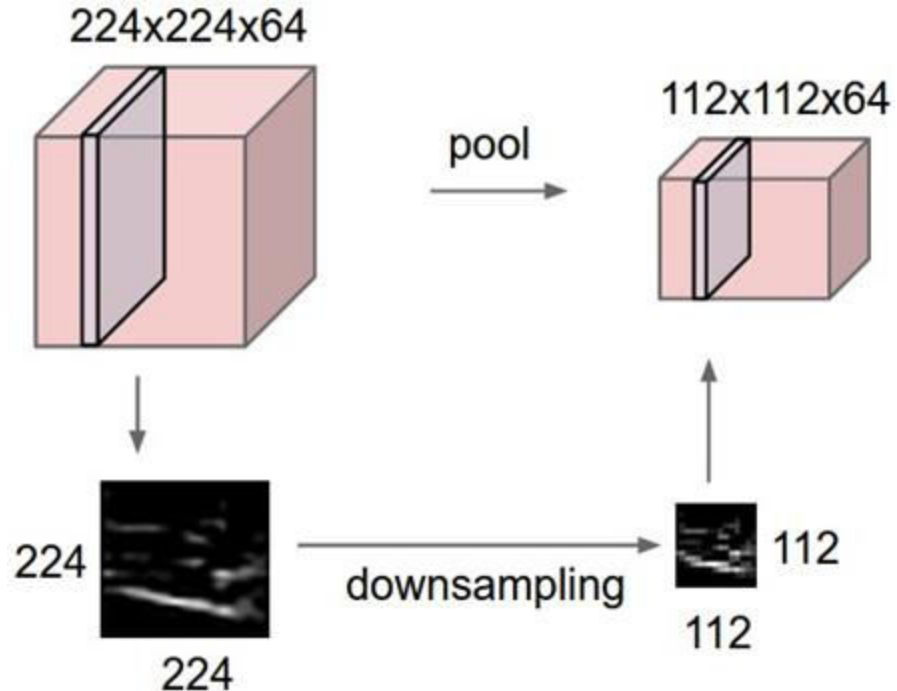


## Two more layers to go: POOL/FC

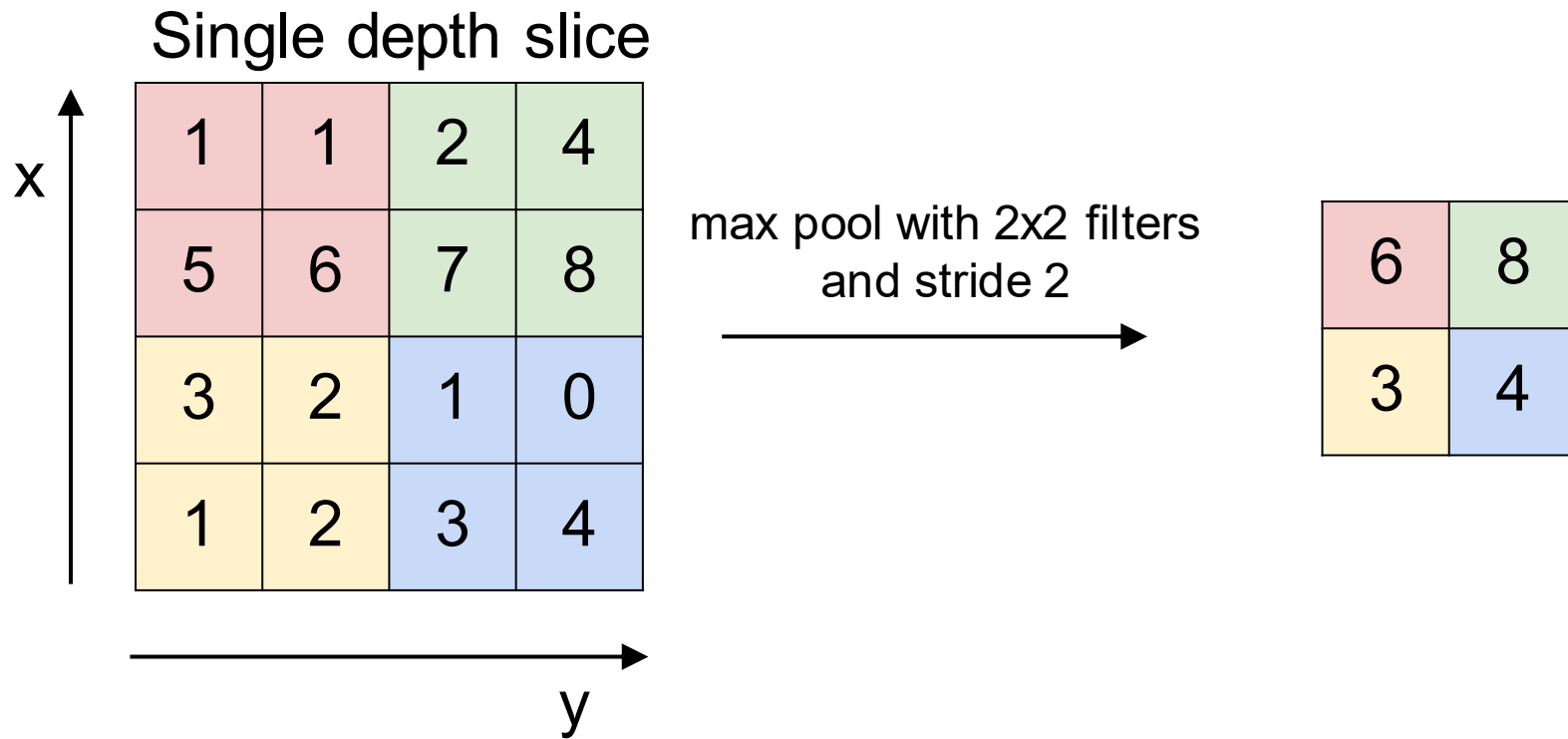


# Pooling layer

- Makes the representations smaller and more manageable
- Operates over each activation map independently:



# MAX Pooling



## Common settings:

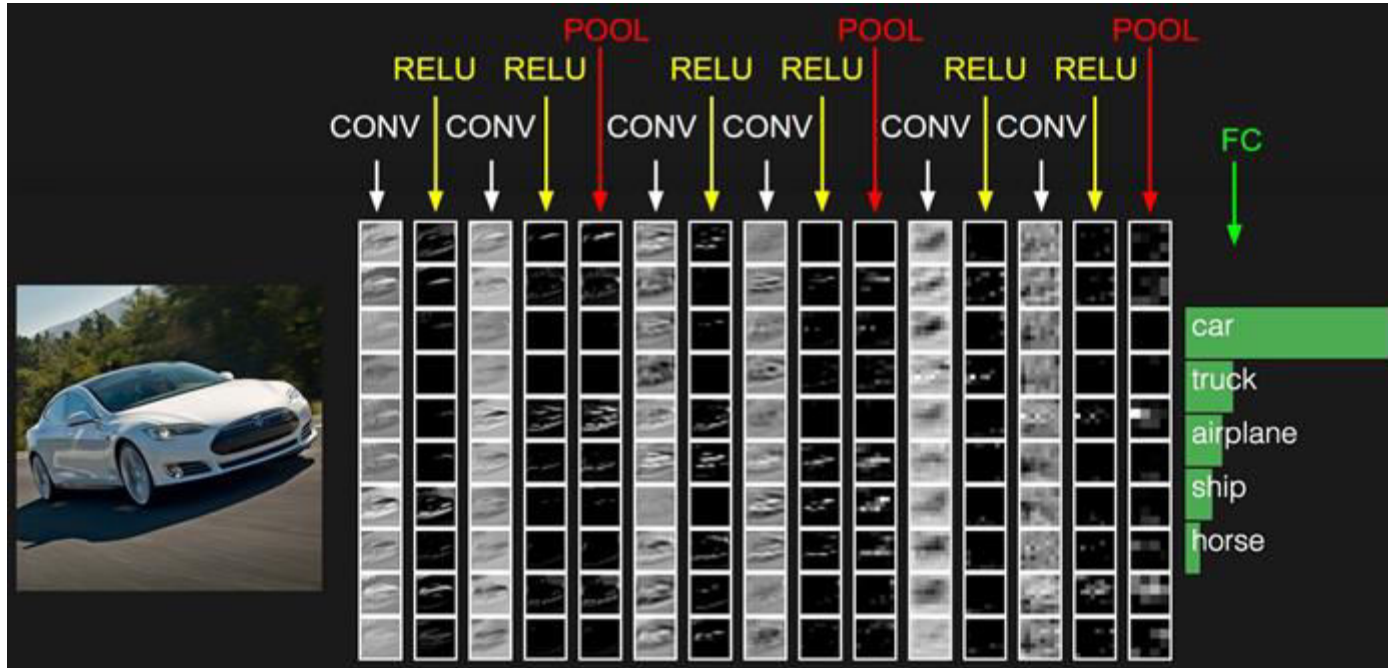
$$F = 2, S = 2$$

$$F = 3, S = 2$$

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
  - their spatial extent  $F$ ,
  - the stride  $S$ ,
- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F)/S + 1$
  - $H_2 = (H_1 - F)/S + 1$
  - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

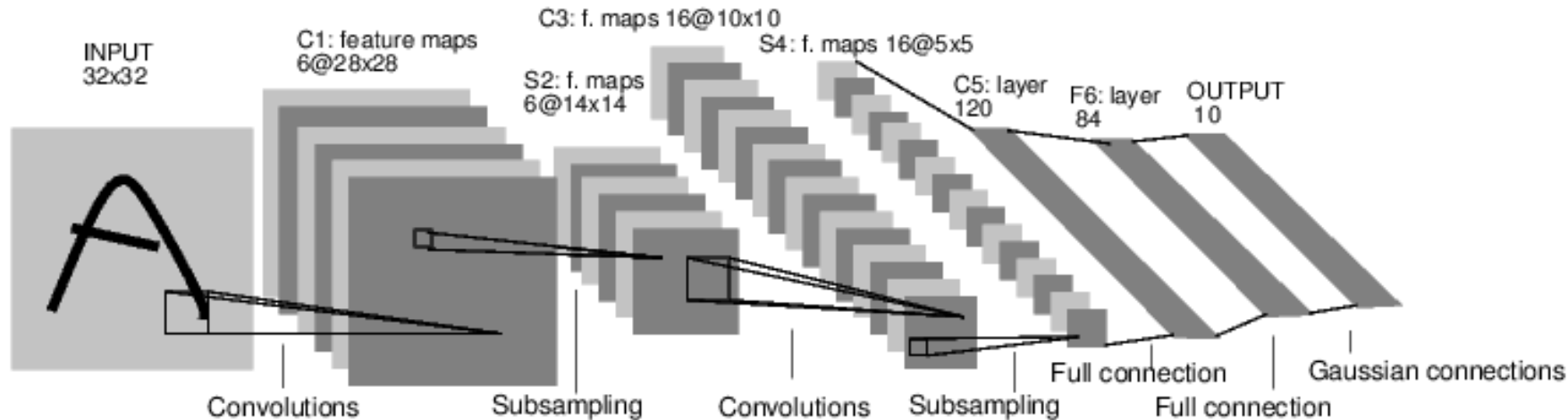
# Fully-Connected Layer (FC layer)

- Contains neurons that connect to the entire input volume, just as in ordinary Neural Networks



# Case Study: LeNet-5

[LeCun et al., 1998]



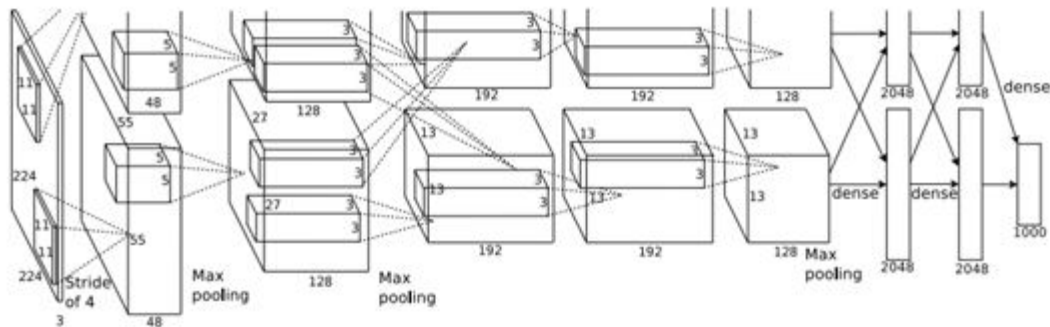
Conv filters were 5x5, applied at stride 1

Subsampling (Pooling) layers were 2x2, applied at stride 2

i.e. architecture is [CONV-POOL-CONV-POOL-CONV-FC]

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

**First layer (CONV1):** 96 11x11 filters applied at stride 4

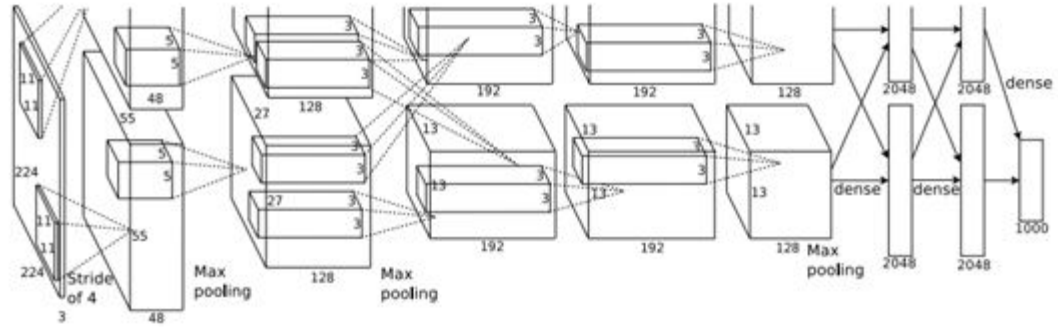
=>

Q: What is the output volume size?

Hint:  $(227-11)/4+1 = 55$

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

**First layer (CONV1):** 96 11x11 filters applied at stride 4

=>

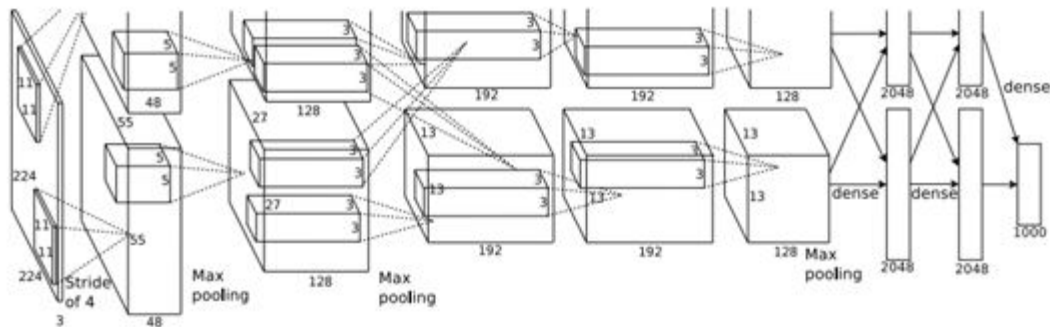
Output volume **[55x55x96]**

Q: What is the total number of parameters in this layer?



# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

**First layer (CONV1):** 96 11x11 filters applied at stride 4

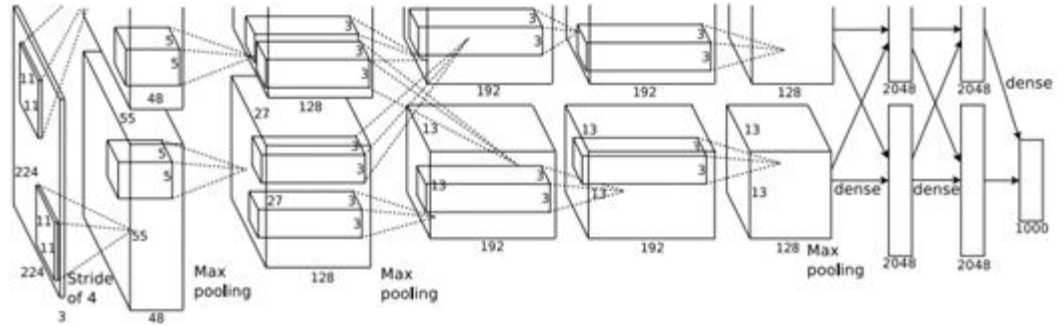
=>

Output volume **[55x55x96]**

Parameters:  $(11*11*3)*96 = \mathbf{35K}$

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

After CONV1: 55x55x96

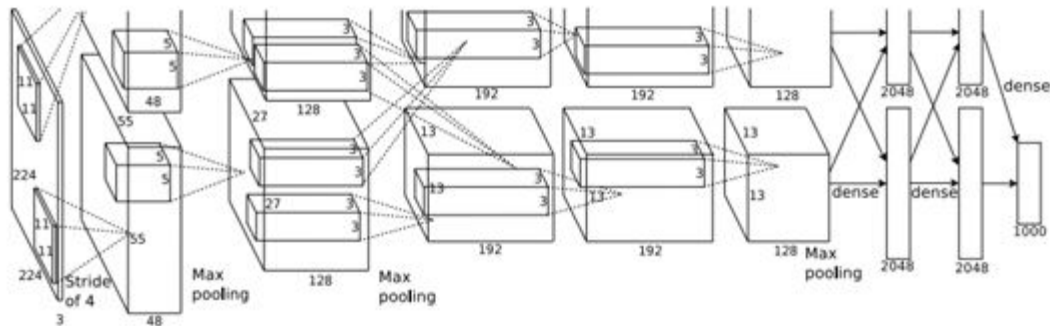
**Second layer (POOL1):** 3x3 filters applied at stride 2

Q: What is the output volume size?

Hint:  $(55-3)/2+1 = 27$

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

After CONV1: 55x55x96

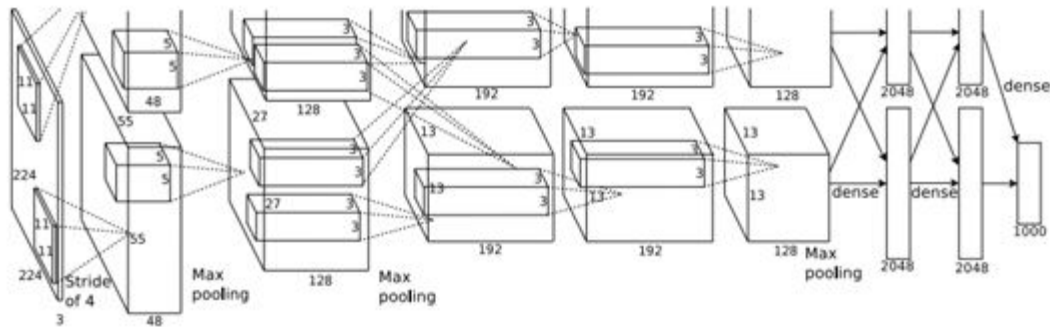
**Second layer (POOL1):** 3x3 filters applied at stride 2

Output volume: 27x27x96

Q: What is the number of parameters in this layer?

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

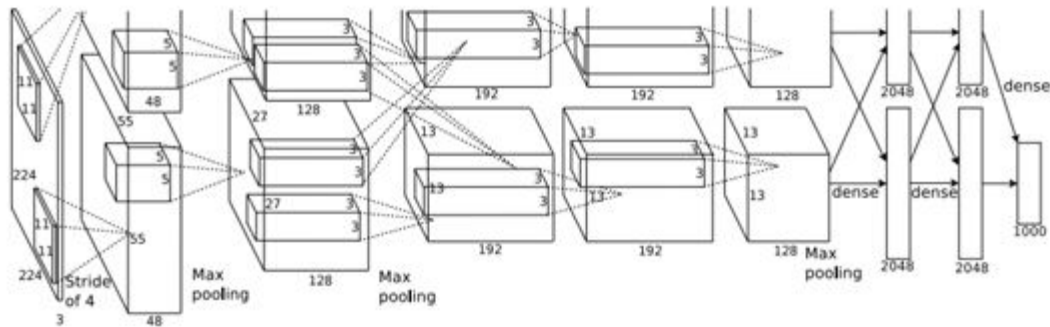
After CONV1: 55x55x96

**Second layer (POOL1):** 3x3 filters applied at stride 2

Output volume: 27x27x96

Parameters: 0

[Krizhevsky et al. 2012]



Input: 227x227x3 images

After CONV1: 55x55x96

After POOL1: 27x27x96

■ ■ ■

# Case Study: AlexNet

[Krizhevsky et al. 2012]

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1**: 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1**: 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2**: 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3**: 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4**: 384 3x3 filters at stride 1, pad 1

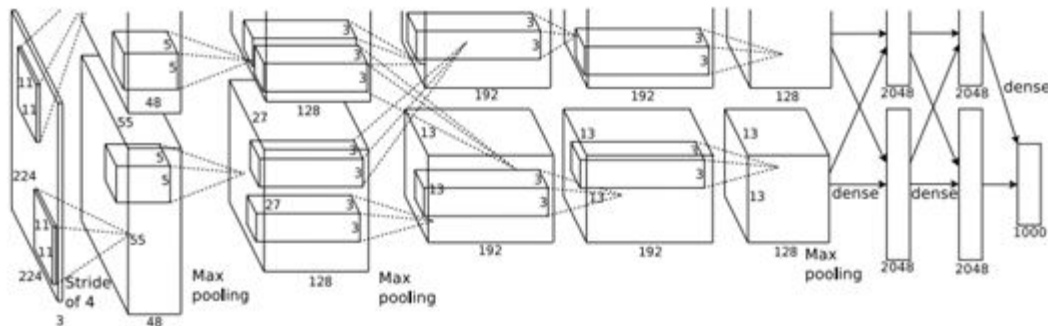
[13x13x256] **CONV5**: 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3**: 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



## Details/Retrospectives:

- First use of ReLU
- Used Norm layers (not common anymore)
- Heavy data augmentation
- Dropout 0.5
- Batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when validation accuracy plateaus
- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% -> 16.4%

# Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Only 3x3 CONV stride 1, pad 1  
and 2x2 MAX POOL stride 2

Best model

11.2% top 5 error in ILSVRC 2013

=>

7.3% top 5 error

| ConvNet Configuration       |                        |                               |  |  |   |
|-----------------------------|------------------------|-------------------------------|--|--|---|
| A                           | A-LRN                  | B                             | C  | D  | E   |
| 11 weight layers            | 11 weight layers       | 13 weight layers              | 16 weight layers                           | 16 weight layers                           | 19 weight layers  |
| input (224 × 224 RGB image) |                        |                               |  |  |   |
| conv3-64                    | conv3-64<br>LRN        | conv3-64<br><b>conv3-64</b>   | conv3-64<br>conv3-64                       | conv3-64<br>conv3-64                       | conv3-64<br>conv3-64                                    |
| maxpool                     |                        |                               |  |  |   |
| conv3-128                   | conv3-128              | conv3-128<br><b>conv3-128</b> | conv3-128<br>conv3-128                     | conv3-128<br>conv3-128                     | conv3-128<br>conv3-128                                  |
| maxpool                     |                        |                               |  |  |   |
| conv3-256<br>conv3-256      | conv3-256<br>conv3-256 | conv3-256<br>conv3-256        | conv3-256<br>conv3-256<br><b>conv1-256</b> | conv3-256<br>conv3-256<br><b>conv3-256</b> | conv3-256<br>conv3-256<br>conv3-256<br><b>conv3-256</b> |
| maxpool                     |                        |                               |  |  |   |
| conv3-512<br>conv3-512      | conv3-512<br>conv3-512 | conv3-512<br>conv3-512        | conv3-512<br>conv3-512<br><b>conv1-512</b> | conv3-512<br>conv3-512<br><b>conv3-512</b> | conv3-512<br>conv3-512<br>conv3-512<br><b>conv3-512</b> |
| maxpool                     |                        |                               |  |  |   |
| conv3-512<br>conv3-512      | conv3-512<br>conv3-512 | conv3-512<br>conv3-512        | conv3-512<br>conv3-512<br><b>conv1-512</b> | conv3-512<br>conv3-512<br><b>conv3-512</b> | conv3-512<br>conv3-512<br>conv3-512<br><b>conv3-512</b> |
| maxpool                     |                        |                               |  |  |   |
| FC-4096                     |                        |                               |  |  |   |
| FC-4096                     |                        |                               |  |  |   |
| FC-1000                     |                        |                               |  |  |   |
| soft-max                    |                        |                               |  |  |   |

Table 2: Number of parameters (in millions).

| Network              | A,A-LRN | B   | C   | D   | E   |
|----------------------|---------|-----|-----|-----|-----|
| Number of parameters | 133     | 133 | 134 | 138 | 144 |

(not counting biases)

INPUT: [224x224x3] memory:  $224*224*3=150K$  params: 0  
CONV3-64: [224x224x64] memory:  $224*224*64=3.2M$  params:  $(3*3*3)*64 = 1,728$   
CONV3-64: [224x224x64] memory:  $224*224*64=3.2M$  params:  $(3*3*64)*64 = 36,864$   
POOL2: [112x112x64] memory:  $112*112*64=800K$  params: 0  
CONV3-128: [112x112x128] memory:  $112*112*128=1.6M$  params:  $(3*3*64)*128 = 73,728$   
CONV3-128: [112x112x128] memory:  $112*112*128=1.6M$  params:  $(3*3*128)*128 = 147,456$   
POOL2: [56x56x128] memory:  $56*56*128=400K$  params: 0  
CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*128)*256 = 294,912$   
CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*256)*256 = 589,824$   
CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*256)*256 = 589,824$   
POOL2: [28x28x256] memory:  $28*28*256=200K$  params: 0  
CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*256)*512 = 1,179,648$   
CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*512)*512 = 2,359,296$   
CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*512)*512 = 2,359,296$   
POOL2: [14x14x512] memory:  $14*14*512=100K$  params: 0  
CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$   
CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$   
CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$   
POOL2: [7x7x512] memory:  $7*7*512=25K$  params: 0  
FC: [1x1x4096] memory: 4096 params:  $7*7*512*4096 = 102,760,448$   
FC: [1x1x4096] memory: 4096 params:  $4096*4096 = 16,777,216$   
FC: [1x1x1000] memory: 1000 params:  $4096*1000 = 4,096,000$

**TOTAL memory:**  $24M * 4 \text{ bytes} \sim 93MB / \text{image}$  (only forward!  $\sim x2$  for backward)

**TOTAL params:** 138M parameters

| ConvNet Configuration     |                  |                  |    |
|---------------------------|------------------|------------------|----|
| B                         | C                | D                |    |
| 13 weight layers          | 16 weight layers | 16 weight layers | 19 |
| put (224 × 224 RGB image) |                  |                  |    |
| conv3-64                  | conv3-64         | conv3-64         | cc |
| <b>conv3-64</b>           | conv3-64         | conv3-64         | cc |
| maxpool                   |                  |                  |    |
| conv3-128                 | conv3-128        | conv3-128        | co |
| <b>conv3-128</b>          | conv3-128        | conv3-128        | co |
| maxpool                   |                  |                  |    |
| conv3-256                 | conv3-256        | conv3-256        | co |
| conv3-256                 | conv3-256        | conv3-256        | co |
|                           | <b>conv1-256</b> | <b>conv3-256</b> | co |
|                           |                  | <b>conv3-256</b> | co |
| maxpool                   |                  |                  |    |
| conv3-512                 | conv3-512        | conv3-512        | co |
| conv3-512                 | conv3-512        | conv3-512        | co |
|                           | <b>conv1-512</b> | <b>conv3-512</b> | co |
|                           |                  | <b>conv3-512</b> | co |
| maxpool                   |                  |                  |    |
| conv3-512                 | conv3-512        | conv3-512        | co |
| conv3-512                 | conv3-512        | conv3-512        | co |
|                           | <b>conv1-512</b> | <b>conv3-512</b> | co |
|                           |                  | <b>conv3-512</b> | co |
| maxpool                   |                  |                  |    |
| FC-4096                   |                  |                  |    |
| FC-4096                   |                  |                  |    |
| FC-1000                   |                  |                  |    |
| soft-max                  |                  |                  |    |



(not counting biases)

INPUT: [224x224x3] memory:  $224*224*3=150K$  params: 0  
CONV3-64: [224x224x64] memory:  $224*224*64=3.2M$  params:  $(3*3*3)*64 = 1,728$   
CONV3-64: [224x224x64] memory:  $224*224*64=3.2M$  params:  $(3*3*64)*64 = 36,864$   
POOL2: [112x112x64] memory:  $112*112*64=800K$  params: 0  
CONV3-128: [112x112x128] memory:  $112*112*128=1.6M$  params:  $(3*3*64)*128 = 73,728$   
CONV3-128: [112x112x128] memory:  $112*112*128=1.6M$  params:  $(3*3*128)*128 = 147,456$   
POOL2: [56x56x128] memory:  $56*56*128=400K$  params: 0  
CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*128)*256 = 294,912$   
CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*256)*256 = 589,824$   
CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*256)*256 = 589,824$   
POOL2: [28x28x256] memory:  $28*28*256=200K$  params: 0  
CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*256)*512 = 1,179,648$   
CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*512)*512 = 2,359,296$   
CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*512)*512 = 2,359,296$   
POOL2: [14x14x512] memory:  $14*14*512=100K$  params: 0  
CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$   
CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$   
CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$   
POOL2: [7x7x512] memory:  $7*7*512=25K$  params: 0  
FC: [1x1x4096] memory: 4096 params:  $7*7*512*4096 = 102,760,448$   
FC: [1x1x4096] memory: 4096 params:  $4096*4096 = 16,777,216$   
FC: [1x1x1000] memory: 1000 params:  $4096*1000 = 4,096,000$

Note:

Most memory is in  
early CONV

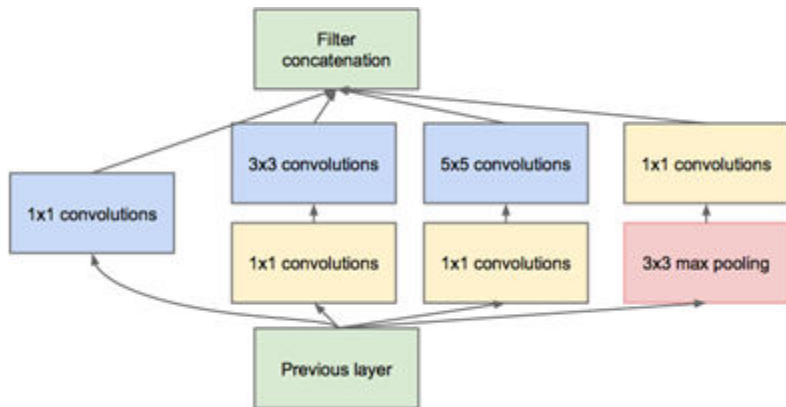
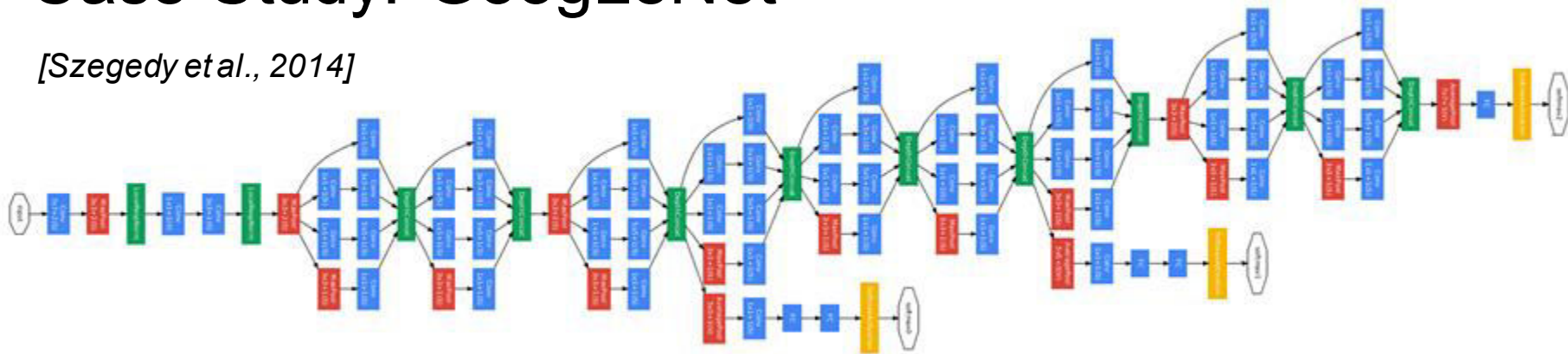
Most params are  
in late FC

TOTAL memory:  $24M * 4 \text{ bytes} \sim 93MB / \text{image}$  (only forward!  $\sim x2$  for backward)

TOTAL params: 138M parameters

# Case Study: GoogLeNet

[Szegedy et al., 2014]



Inception module

ILSVRC 2014 winner (6.7% top 5 error)

# Case Study: GoogLeNet

| type           | patch size/<br>stride | output<br>size | depth | #1×1 | #3×3<br>reduce | #3×3 | #5×5<br>reduce | #5×5 | pool<br>proj | params | ops  |
|----------------|-----------------------|----------------|-------|------|----------------|------|----------------|------|--------------|--------|------|
| convolution    | 7×7/2                 | 112×112×64     | 1     |      |                |      |                |      |              | 2.7K   | 34M  |
| max pool       | 3×3/2                 | 56×56×64       | 0     |      |                |      |                |      |              |        |      |
| convolution    | 3×3/1                 | 56×56×192      | 2     |      | 64             | 192  |                |      |              | 112K   | 360M |
| max pool       | 3×3/2                 | 28×28×192      | 0     |      |                |      |                |      |              |        |      |
| inception (3a) |                       | 28×28×256      | 2     | 64   | 96             | 128  | 16             | 32   | 32           | 159K   | 128M |
| inception (3b) |                       | 28×28×480      | 2     | 128  | 128            | 192  | 32             | 96   | 64           | 380K   | 304M |
| max pool       | 3×3/2                 | 14×14×480      | 0     |      |                |      |                |      |              |        |      |
| inception (4a) |                       | 14×14×512      | 2     | 192  | 96             | 208  | 16             | 48   | 64           | 364K   | 73M  |
| inception (4b) |                       | 14×14×512      | 2     | 160  | 112            | 224  | 24             | 64   | 64           | 437K   | 88M  |
| inception (4c) |                       | 14×14×512      | 2     | 128  | 128            | 256  | 24             | 64   | 64           | 463K   | 100M |
| inception (4d) |                       | 14×14×528      | 2     | 112  | 144            | 288  | 32             | 64   | 64           | 580K   | 119M |
| inception (4e) |                       | 14×14×832      | 2     | 256  | 160            | 320  | 32             | 128  | 128          | 840K   | 170M |
| max pool       | 3×3/2                 | 7×7×832        | 0     |      |                |      |                |      |              |        |      |
| inception (5a) |                       | 7×7×832        | 2     | 256  | 160            | 320  | 32             | 128  | 128          | 1072K  | 54M  |
| inception (5b) |                       | 7×7×1024       | 2     | 384  | 192            | 384  | 48             | 128  | 128          | 1388K  | 71M  |
| avg pool       | 7×7/1                 | 1×1×1024       | 0     |      |                |      |                |      |              |        |      |
| dropout (40%)  |                       | 1×1×1024       | 0     |      |                |      |                |      |              |        |      |
| linear         |                       | 1×1×1000       | 1     |      |                |      |                |      |              | 1000K  | 1M   |
| softmax        |                       | 1×1×1000       | 0     |      |                |      |                |      |              |        |      |

Fun features:

- Only 5 million params!  
(Removes FC layers completely)


**Compared to AlexNet:**

- 12x less params
- 2x more computations
- 6.67% (vs. 16.4%)

# Case Study: ResNet

[He et al., 2015]


ILSVRC 2015 winner (3.6% top 5 error)



## MSRA @ ILSVRC & COCO 2015 Competitions

- **1st places** in all five main tracks
  - ImageNet Classification: “Ultra-deep” (quote Yann) **152-layer** nets
  - ImageNet Detection: **16%** better than 2nd
  - ImageNet Localization: **27%** better than 2nd
  - COCO Detection: **11%** better than 2nd
  - COCO Segmentation: **12%** better than 2nd

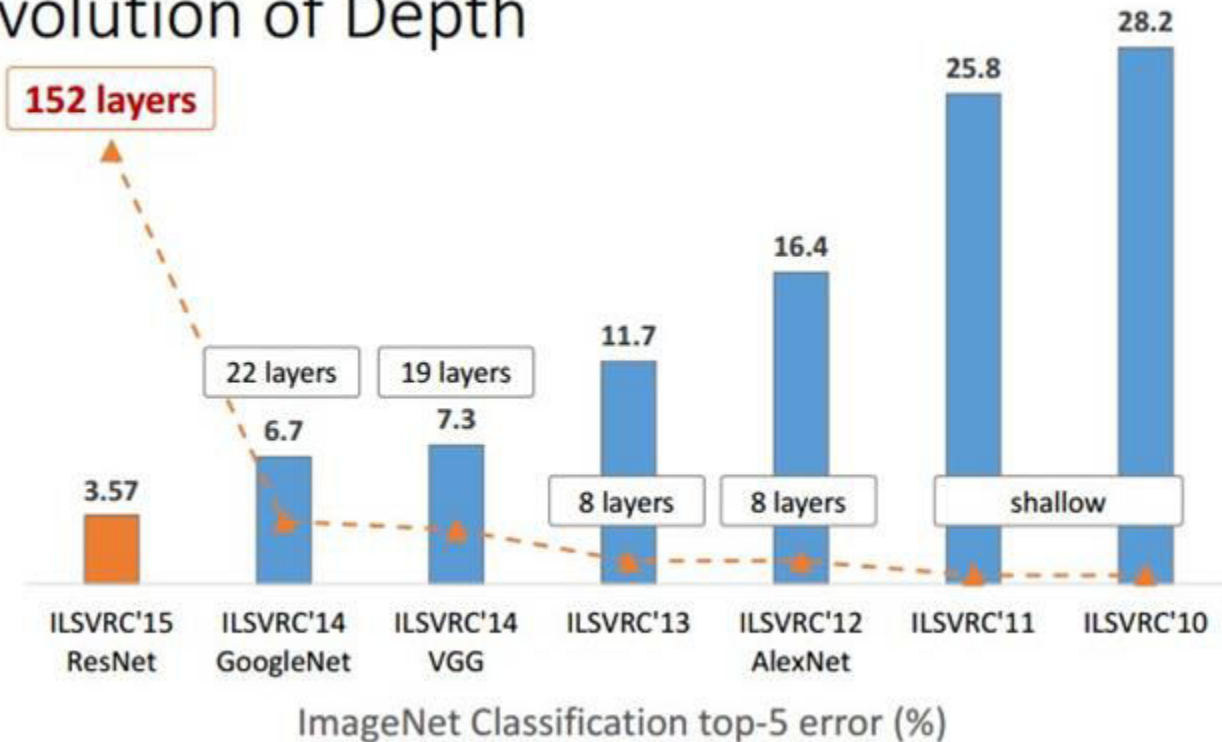
\*improvements are relative numbers



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. “Deep Residual Learning for Image Recognition”. arXiv 2015.

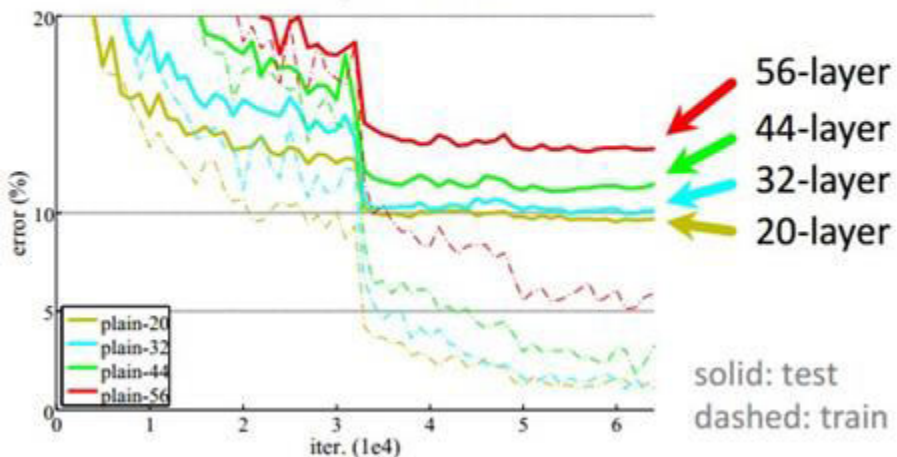
Slide after Kaiming He

# Revolution of Depth

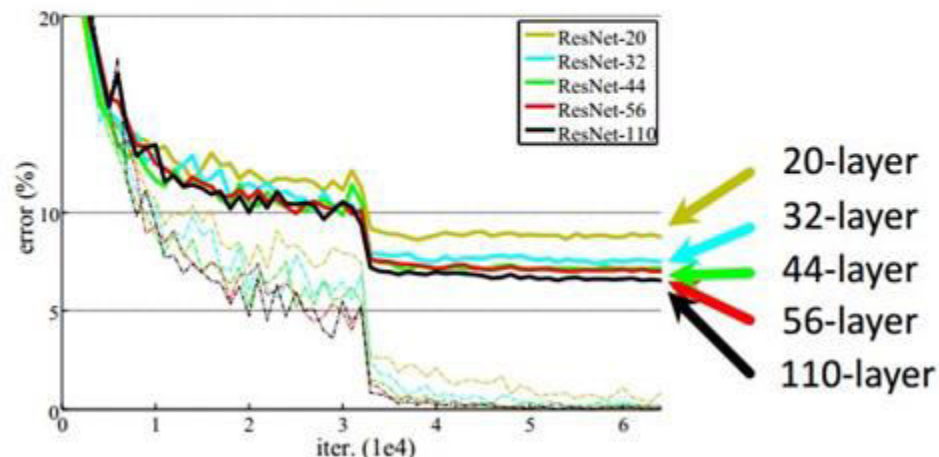


# CIFAR-10 experiments

CIFAR-10 plain nets



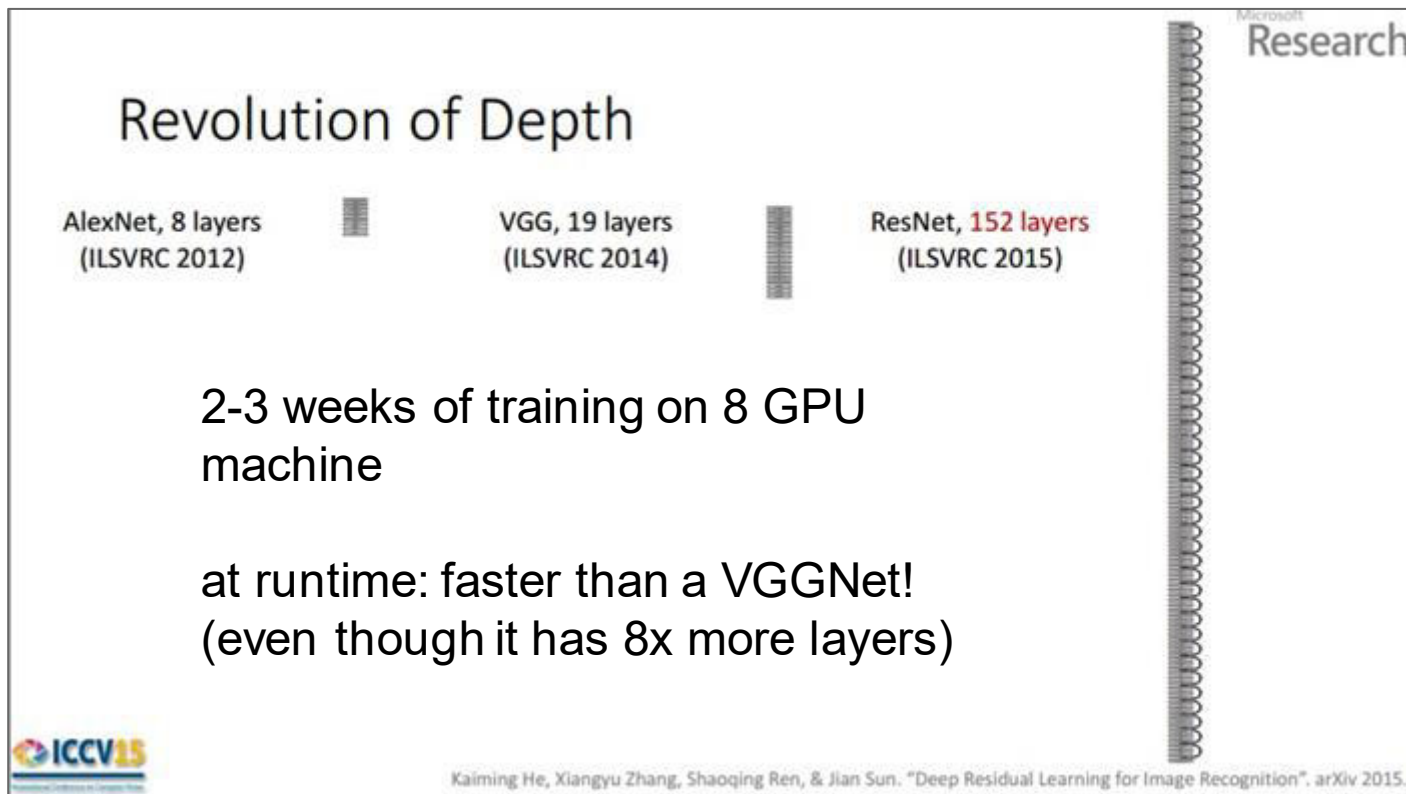
CIFAR-10 ResNets



# Case Study: ResNet

[He et al., 2015]

ILSVRC 2015 winner (3.6% top 5 error)



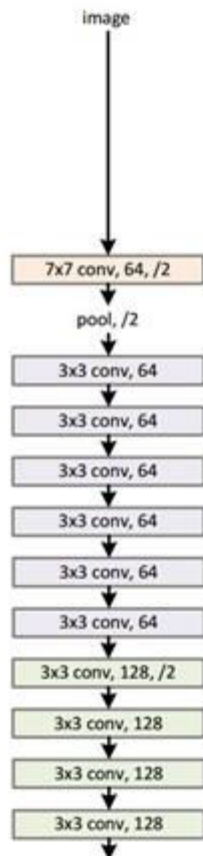
Slide after Kaiming He



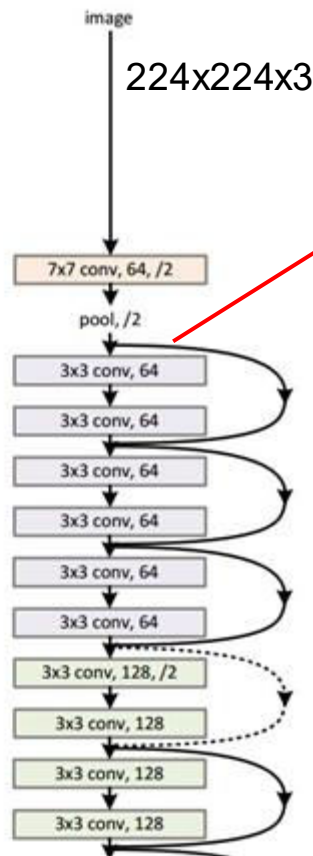
# Case Study: ResNet

[He et al., 2015]

34-layer plain



34-layer residual

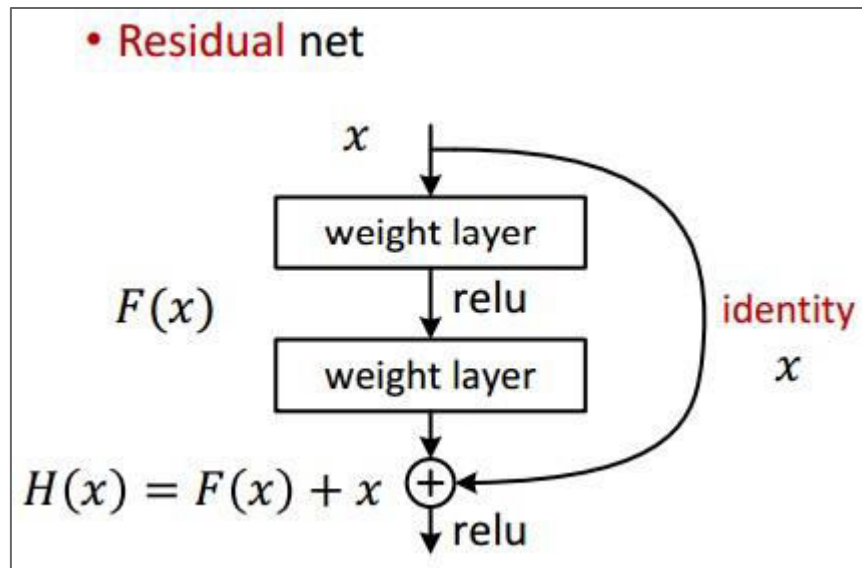
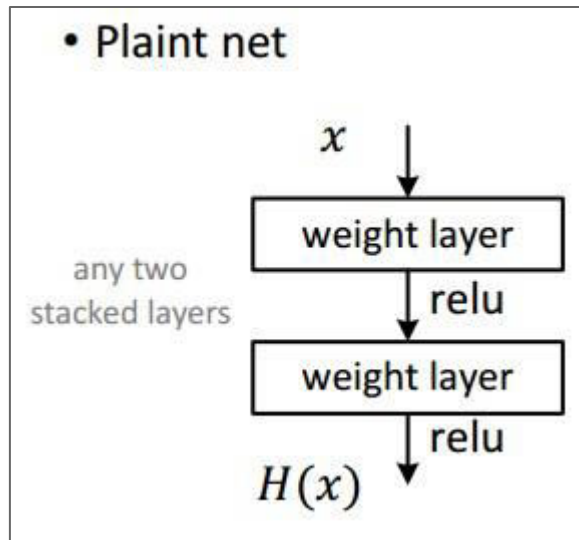


Spatial dimension  
only  $56 \times 56$ !



# Case Study: ResNet

[He et al., 2015]



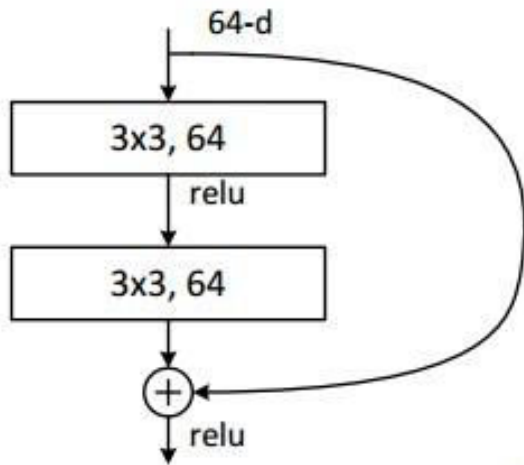
# Case Study: ResNet

*[He et al., 2015]*

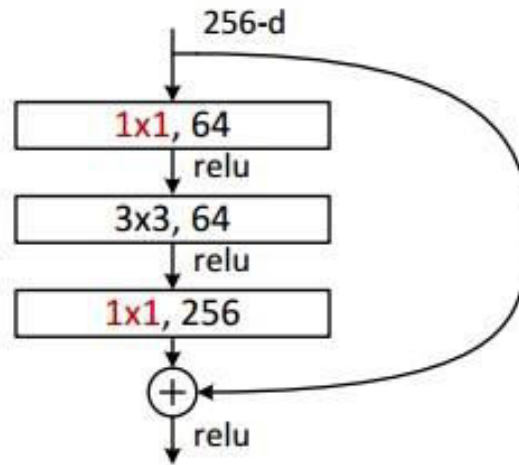
- Batch Normalization after every CONV layer
- Xavier/2 initialization
- SGD + momentum (0.9)
- Learning rate: 0.1, divided by 10 when validation error plateaus
- Mini-batch size 256
- Weight decay of  $1e-5$
- No dropout used!

# Case Study: ResNet

[He et al., 2015]



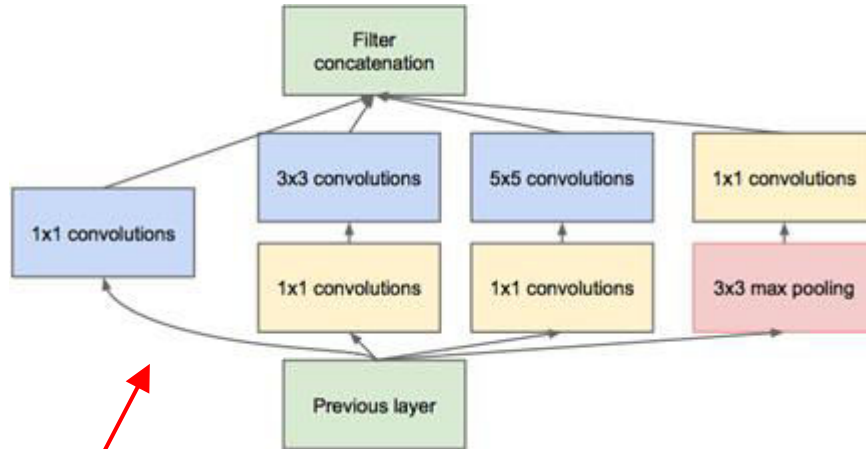
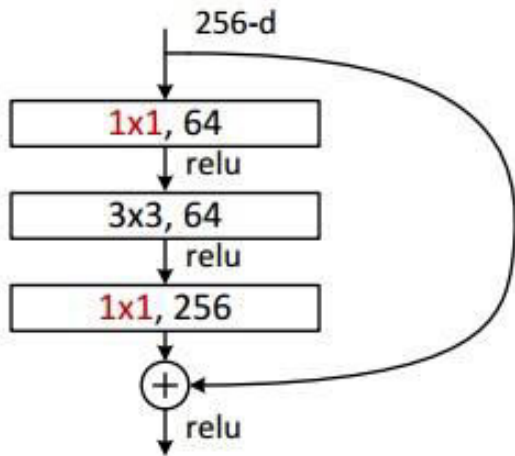
all-3x3



**bottleneck**  
(for ResNet-50/101/152)

# Case Study: ResNet

[He et al., 2015]

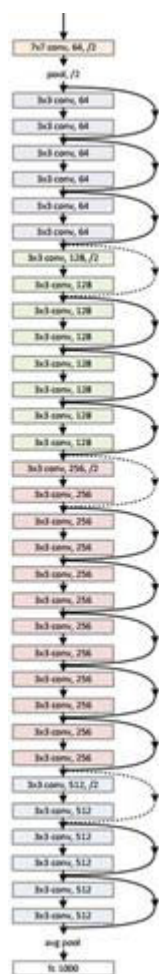


This trick is also used in GoogLeNet

# Case Study: ResNet

[He et al., 2015]

| layer name | output size | 18-layer  | 34-layer  | 50-layer  | 101-layer  | 152-layer  |
|------------|-------------|---|---|---|--|--|
| conv1      | 112×112     | 7×7, 64, stride 2   |   |   |  |  |
| conv2_x    | 56×56       | 3×3 max pool, stride 2  |   |   |  |  |
|            |             | $\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$   | $\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$   | $\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$    | $\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$     | $\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$     |
| conv3_x    | 28×28       | $\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$ | $\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$  | $\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$   | $\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$   |
| conv4_x    | 14×14       | $\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$ | $\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$ | $\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$ | $\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$ |
| conv5_x    | 7×7         | $\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$  | $\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$  |
|            | 1×1         | average pool, 1000-d fc, softmax  |   |   |  |  |
| FLOPs      |             | $1.8 \times 10^9$   | $3.6 \times 10^9$   | $3.8 \times 10^9$   | $7.6 \times 10^9$  | $11.3 \times 10^9$   |



# Summary

- ConvNets stack CONV, POOL, FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Typical architectures look like

**$[(\text{CONV-RELU})^N\text{-POOL?}]^M\text{-(FC-RELU)}^K, \text{SOFTMAX}$**

where N is usually up to  $\sim 5$ , M is large,  $0 \leq K \leq 2$ .

- ... but recent advances such as ResNet/GoogLeNet challenge this paradigm