

# CSCI567 Machine Learning (Fall 2017)

Prof. Fei Sha

U of Southern California

Lecture on Sept. 19, 2017

# Outline

- 1 Administration
- 2 Review of last lecture
- 3 Convolution neural networks

# Outline

- 1 Administration
- 2 Review of last lecture
- 3 Convolution neural networks

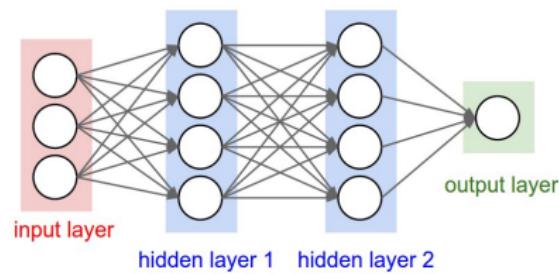
# Administrative stuff

- Thursday lecture: to be taught by Zhiyun Lu due to my traveling
-

# Outline

- 1 Administration
- 2 Review of last lecture
- 3 Convolution neural networks

# Neural nets



# Key steps (essentially, chain rule in calculus)

To compute

$$\frac{\partial \ell}{\partial w_{ji}} = \frac{\partial \ell}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}} = z_i \frac{\partial \ell}{\partial a_j}$$

as  $w_{ji}$  affects only  $a_j$

Passthru the nonlinear transfer function

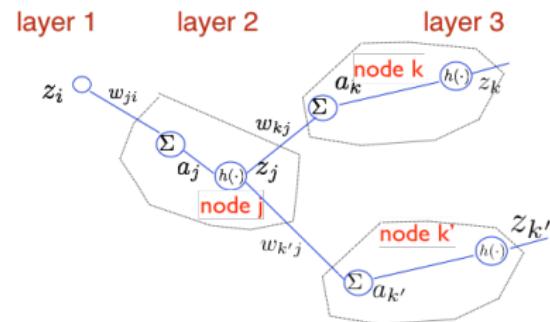
$$\frac{\partial \ell}{\partial a_j} = \frac{\partial \ell}{\partial z_j} \frac{\partial z_j}{\partial a_j} = h'(a_j) \frac{\partial \ell}{\partial z_j}$$

Collect errors from next layer

$$\frac{\partial \ell}{\partial z_j} = \sum_k \frac{\partial \ell}{\partial a_k} \frac{\partial a_k}{\partial z_j} = \sum_k \frac{\partial \ell}{\partial a_k} w_{kj}$$

Recursion

$$\frac{\partial \ell}{\partial a_j} = h'(a_j) \sum_k \frac{\partial \ell}{\partial a_k} w_{kj}$$



# Optimization tricks

- Stochastic gradient descent with minimatch, momentum
- Early stopping, data augmentation, injecting noise for overcoming overfitting

# Acknowledgements

The lecture slides borrow *heavily* from the following sources:

- Stanford Course Cs231n: Convolutional Neural Networks for Visual Recognition Spring 2017 (taught by Prof. Fei-Fei Li and her students).  
The website <http://cs231n.stanford.edu/> provides also very valuable notes, demo codes, and pointers to useful information
- Dr. Ian Goodfellow's lectures on neural networks and deep learning.  
The website <http://deeplearningbook.org> provides very valuable lecture notes, pointers to video-tapped lectures and other useful information

# Outline

- 1 Administration
- 2 Review of last lecture
- 3 Convolution neural networks
  - Motivation
  - Architecture
  - Algorithm

# Image Classification: A core task in Computer Vision



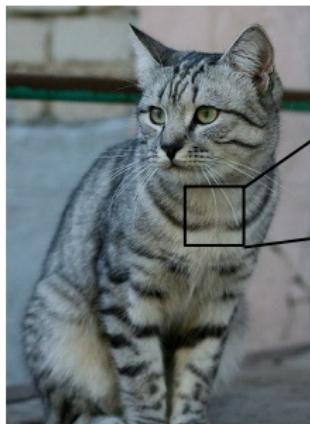
This image by nikita is  
licensed under CC-BY 2.0

(assume given set of discrete labels)  
{dog, cat, truck, plane, ...}



cat

# The Problem: Semantic Gap



This image by nikita is  
licensed under CC-BY 2.0

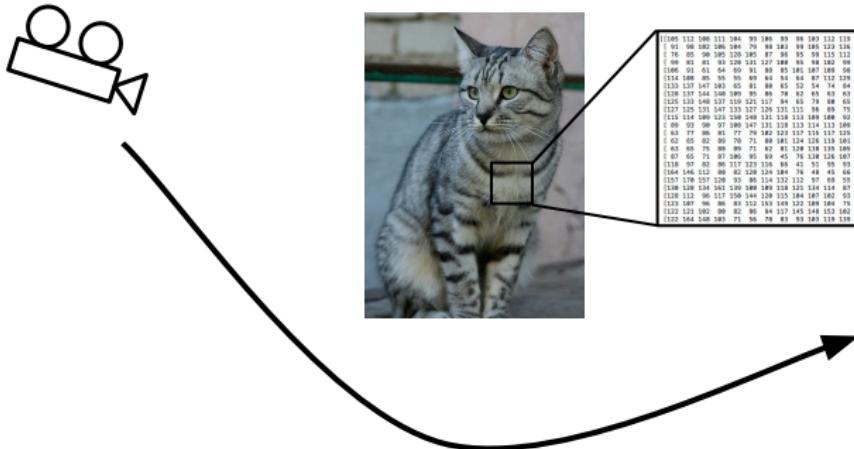
[1] [85 112 180 111 184 99 186 99 86 185 112 119 184 97 93 87]
[1] [76 85 98 105 120 185 87 96 95 99 115 112 106 183 99 85]
[1] [99 81 81 93 128 131 127 180 95 98 182 99 96 93 101 94]
[1] [86 91 61 64 69 91 88 85 181 187 109 98 75 84 96 95]
[1] [114 180 85 95 59 62 64 54 87 120 128 98 74 84 96 95]
[1] [133 117 120 183 65 81 80 52 74 64 80 82 75 85 82]
[1] [128 137 144 148 180 95 86 78 62 65 63 63 68 73 86 101]
[1] [125 133 148 137 119 121 117 94 65 79 88 65 54 64 72 98]
[1] [127 125 131 147 133 127 126 131 111 96 89 75 61 64 72 84]
[1] [115 114 189 123 158 140 131 118 113 189 180 92 74 65 77 78]
[1] [89 111 112 180 111 184 99 186 99 86 185 112 119 184 97 93 87]
[1] [63 77 86 81 77 79 182 123 117 115 117 125 125 138 115 87]
[1] [62 65 82 82 78 71 81 181 124 126 119 181 187 114 131 119]
[1] [63 65 75 88 89 71 62 81 128 138 135 185 81 98 118 118]
[1] [87 65 71 87 104 95 69 45 76 139 126 187 92 94 105 112]
[1] [118 117 120 183 65 81 80 52 74 64 80 82 75 85 82]
[1] [164 146 112 88 82 128 124 184 76 48 45 66 88 183 182 189]
[1] [157 178 157 128 93 86 114 132 112 97 69 55 78 82 99 94]
[1] [138 128 134 161 139 188 109 118 121 134 114 87 65 53 69 86]
[1] [128 120 96 88 83 131 152 124 119 121 134 114 87 65 53 69 86]
[1] [123 121 112 180 111 184 99 186 99 86 117 125 148 153 181 79 92 187]
[1] [122 164 148 183 71 56 78 83 93 183 119 139 182 61 69 84]

What the computer sees

An image is just a big grid of numbers between [0, 255]:

e.g. 800 x 600 x 3  
(3 channels RGB)

## Challenges: Viewpoint variation



All pixels change when  
the camera moves!

This image by Nikita is  
licensed under CC-BY 2.0

# Challenges: Illumination



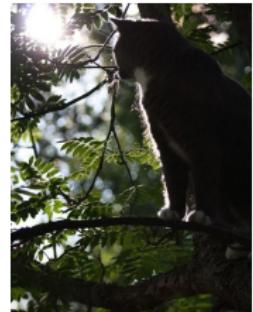
[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)

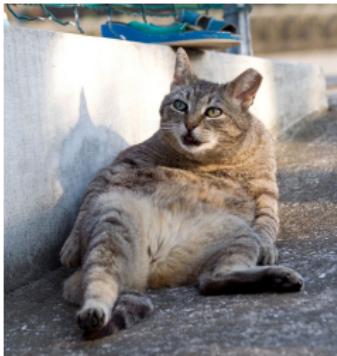


[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)

# Challenges: Deformation



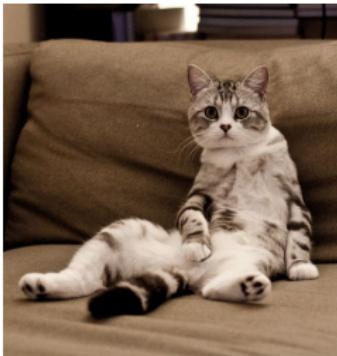
This image by Umberto Salvagnin  
is licensed under CC-BY 2.0



This image by Umberto Salvagnin  
is licensed under CC-BY 2.0



This image by save bear is  
licensed under CC-BY 2.0



This image by Tom Thai is  
licensed under CC-BY 2.0

# Challenges: Occlusion



[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)



This image by [jonsson](#) is licensed under [CC-BY 2.0](#)

## Challenges: Background Clutter



[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)

## Challenges: Intraclass variation



[This image](#) is CC0 1.0 public domain

# Fundamental problems in vision and learning

## The key challenge

How to train a visual recognition model to be resilient to all those variations?

In terms of our nonlinear basis functions, those functions must be able to map the images to an invariant representation (with respect to pose, illumination, etc). But *how to get such representations?*

## Main ideas

- We will need a lot of data that exhibits those variations – otherwise the models cannot automatically conjure up invariances.
- We will need powerful models to capture the invariances.

# How biology does it?

A bit of history:

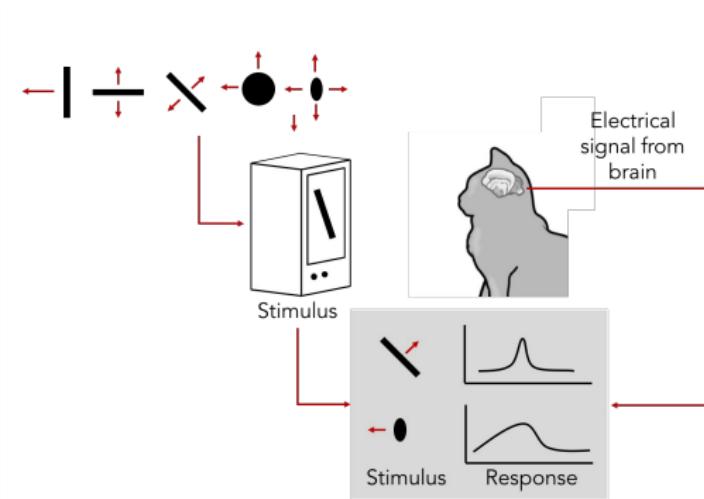
**Hubel & Wiesel,**  
**1959**

RECEPTIVE FIELDS OF SINGLE  
NEURONES IN  
THE CAT'S STRIATE CORTEX

**1962**

RECEPTIVE FIELDS, BINOCULAR  
INTERACTION  
AND FUNCTIONAL ARCHITECTURE IN  
THE CAT'S VISUAL CORTEX

**1968...**



Cat image by CNX OpenStax is licensed under CC BY 4.0; changes made

# Hierarchical organization

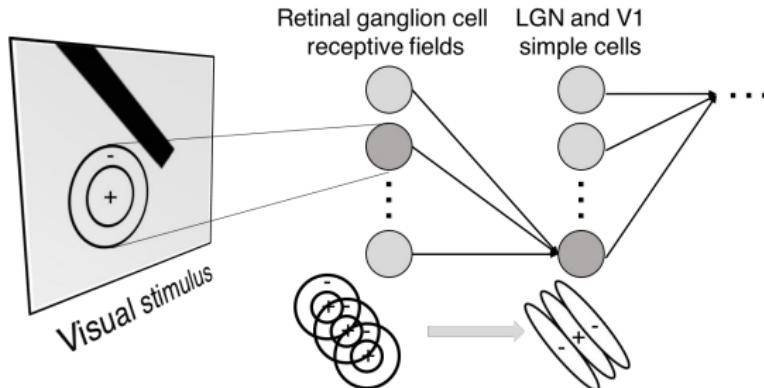
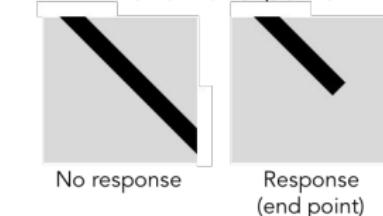


Illustration of hierarchical organization in early visual pathways by Lane McIntosh, copyright CS231n 2017

**Simple cells:**  
Response to light orientation

**Complex cells:**  
Response to light orientation and movement

**Hypercomplex cells:**  
response to movement with an end point

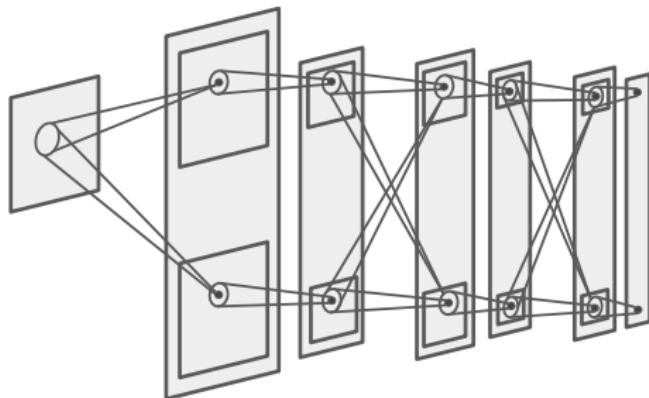


# An attempt to mimic the biological pipeline

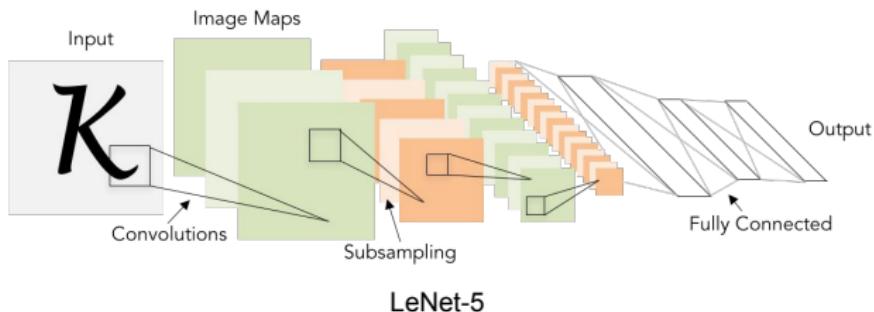
A bit of history:

## Neocognitron [Fukushima 1980]

"sandwich" architecture (SCSCSC...)  
simple cells: modifiable parameters  
complex cells: perform pooling



# A bit of history: Gradient-based learning applied to document recognition [LeCun, Bottou, Bengio, Haffner 1998]



# A bit of history: ImageNet Classification with Deep Convolutional Neural Networks [Krizhevsky, Sutskever, Hinton, 2012]

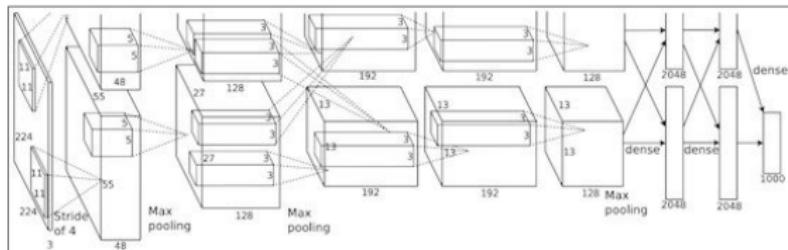


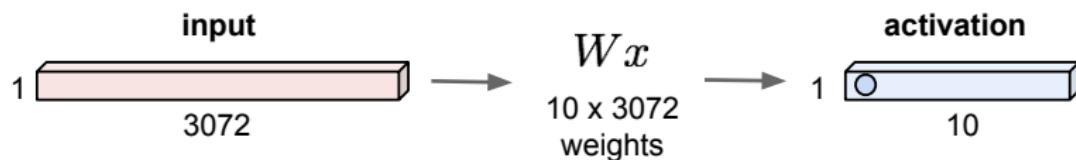
Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

“AlexNet”

# What is wrong with standard NN for image as input?

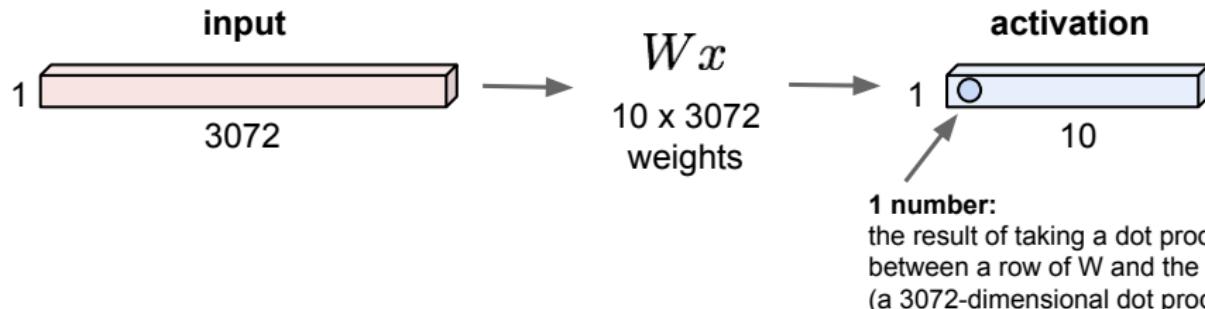
## Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



# Fully Connected Layer

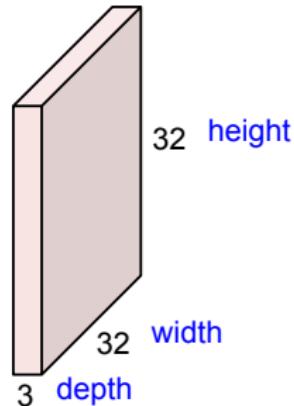
32x32x3 image -> stretch to 3072 x 1



# Convolution with filters to preserve spatial structure

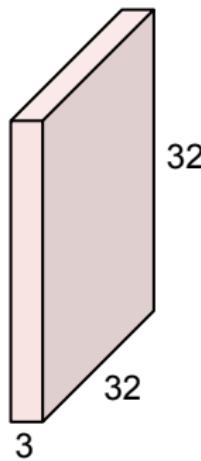
## Convolution Layer

32x32x3 image -> preserve spatial structure



# Convolution Layer

32x32x3 image



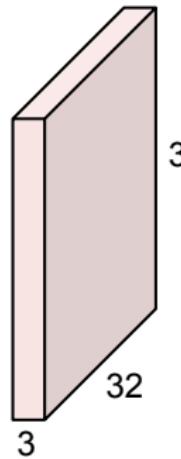
5x5x3 filter



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

# Convolution Layer

32x32x3 image



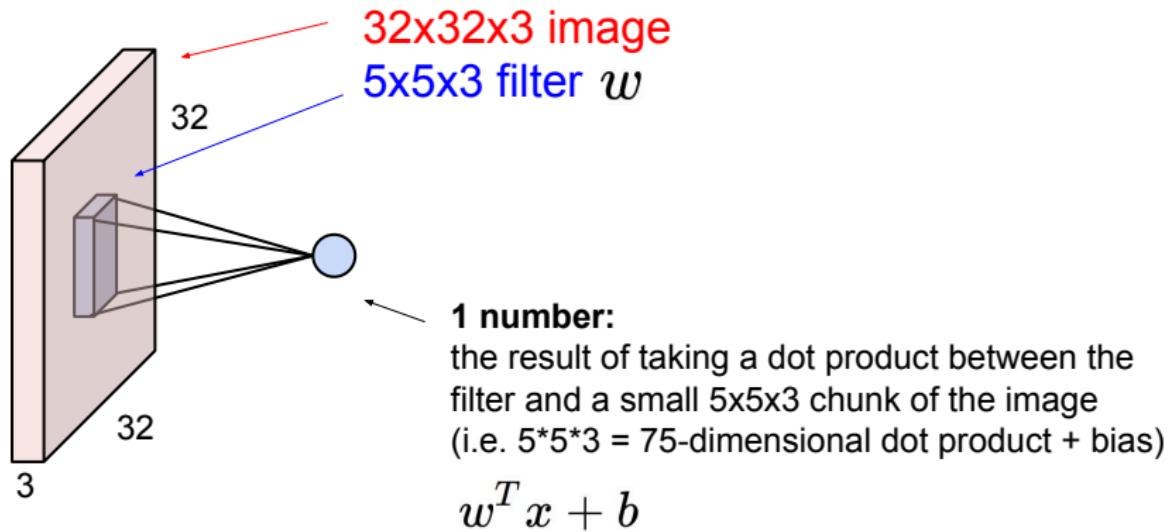
5x5x3 filter



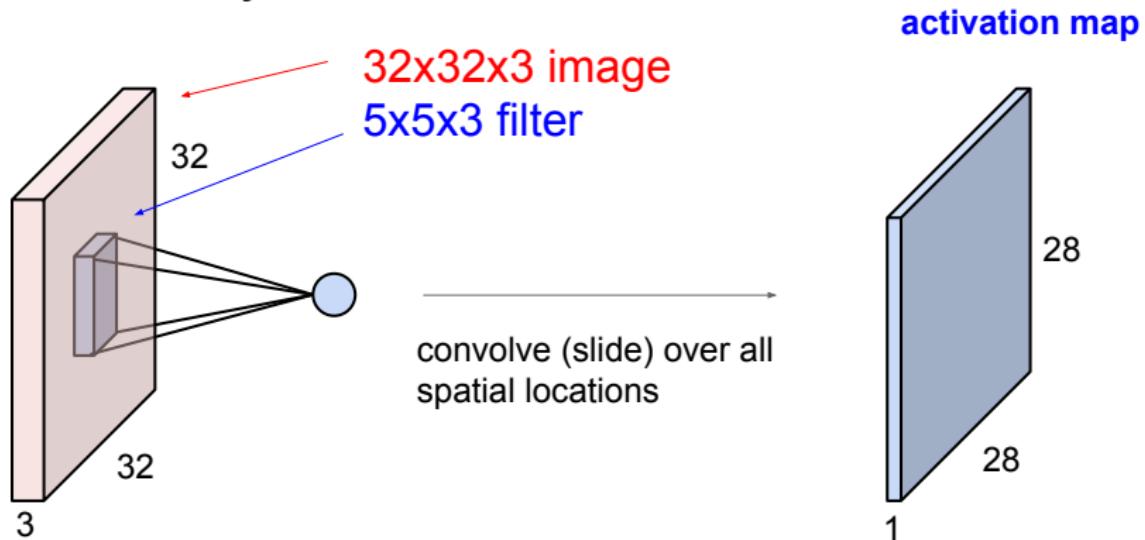
Filters always extend the full depth of the input volume

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

# Convolution Layer

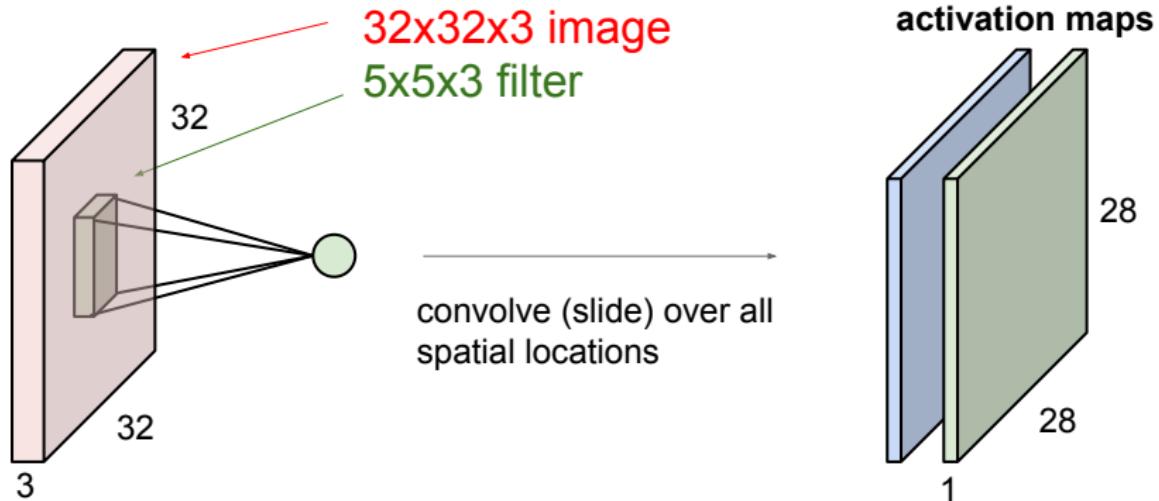


# Convolution Layer

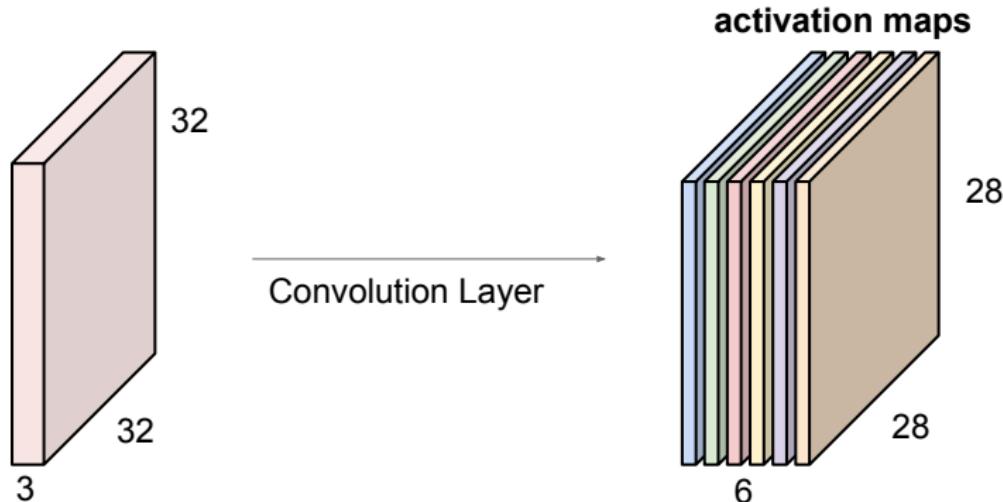


# Convolution Layer

consider a second, green filter

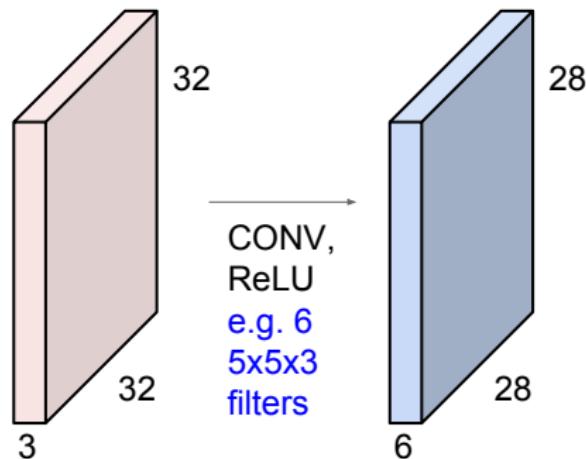


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

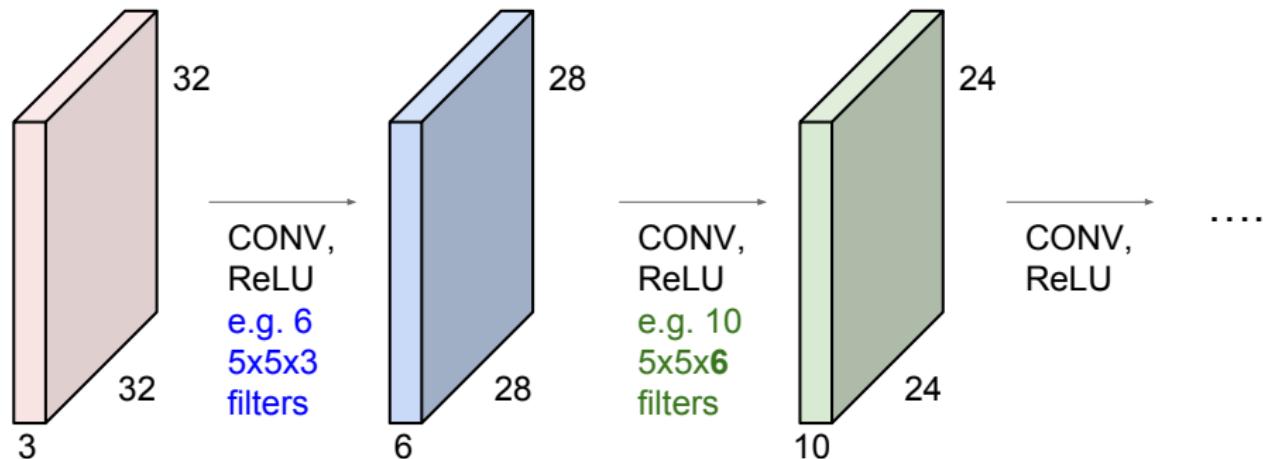


We stack these up to get a “new image” of size  $28 \times 28 \times 6$ !

**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

[Zeiler and Fergus 2013]

Visualization of VGG-16 by Lane McIntosh. VGG-16 architecture from [Simonyan and Zisserman 2014].

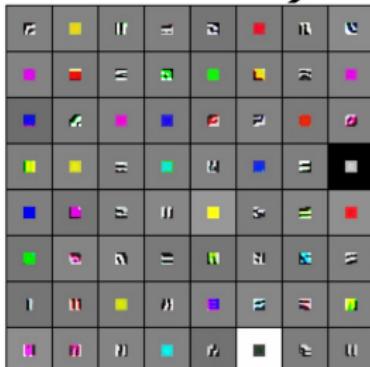


Low-level features

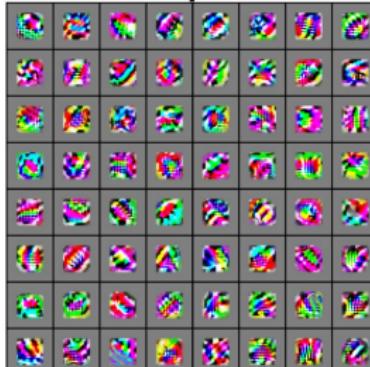
Mid-level features

High-level features

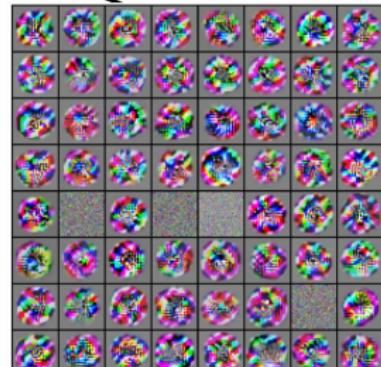
Linearly  
separable  
classifier



VGG-16 Conv1\_1

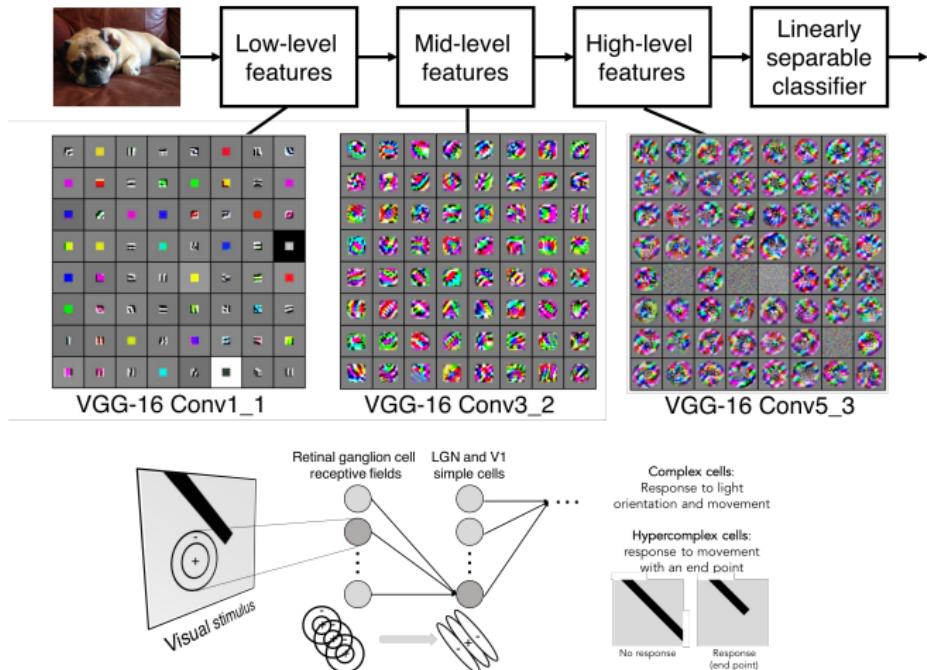


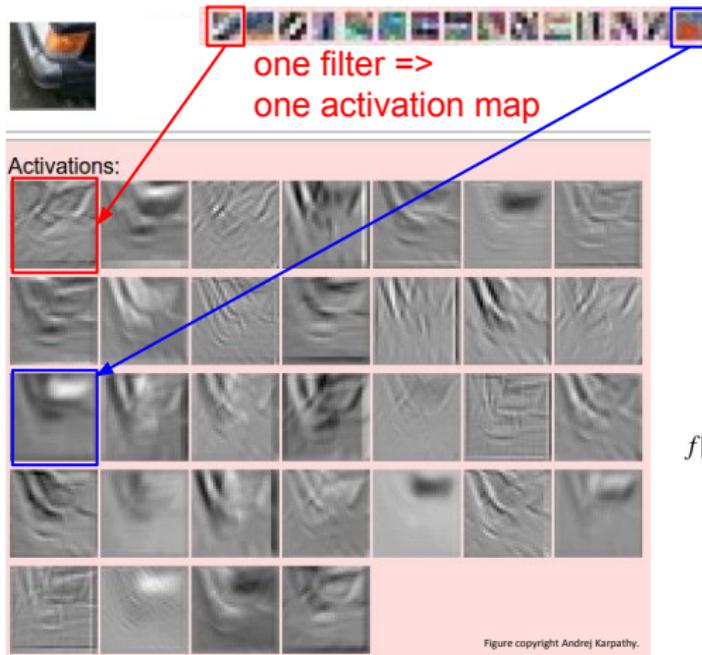
VGG-16 Conv3\_2



VGG-16 Conv5\_3

## Preview





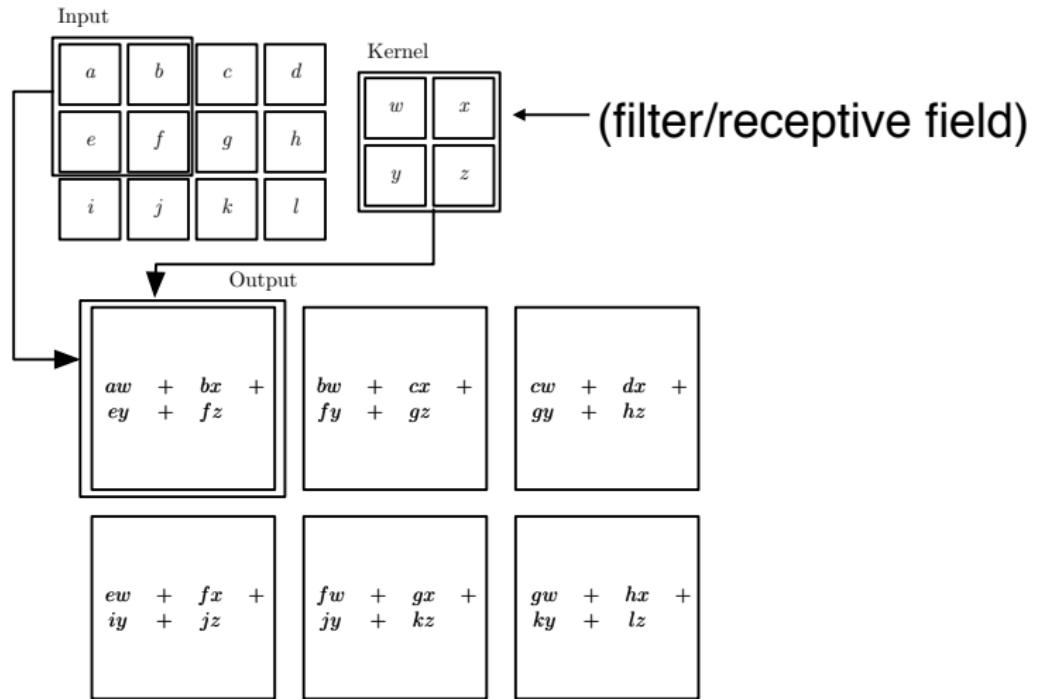
example 5x5 filters  
(32 total)

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

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

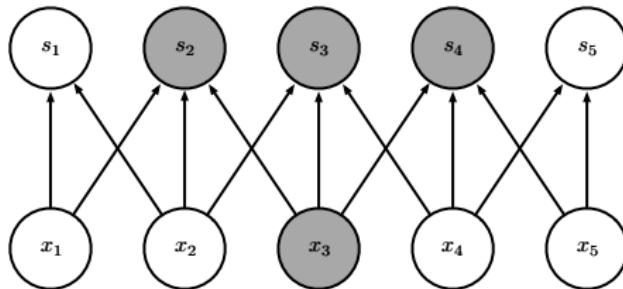
↑  
elementwise multiplication and sum of a filter and the signal (image)

# 2D Convolution

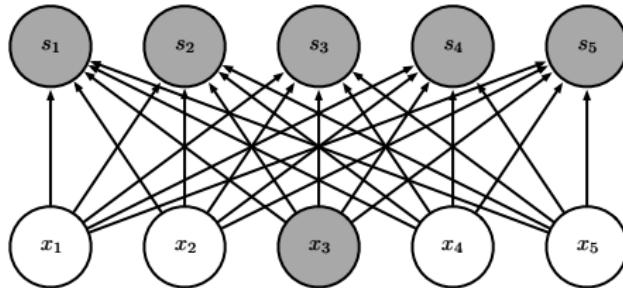


# Local Receptive Field Leads to Sparse Connectivity (affects less)

Sparse  
connections  
due to small  
convolution  
kernel

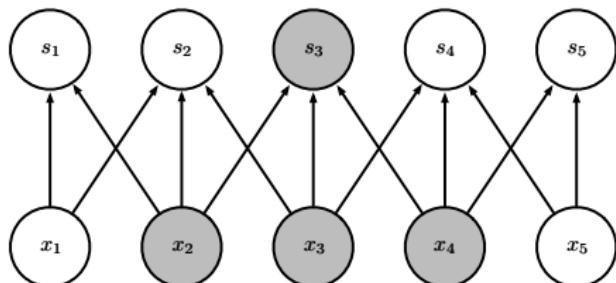


Dense  
connections



# Sparse connectivity: being affected by less

Sparse connections due to small convolution kernel



Dense connections

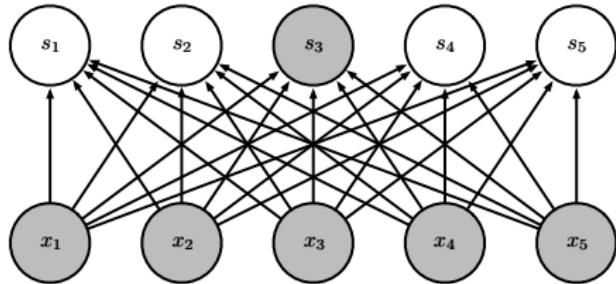


Figure 9.3

However, “depth” by stacking layers creates larger receptive fields

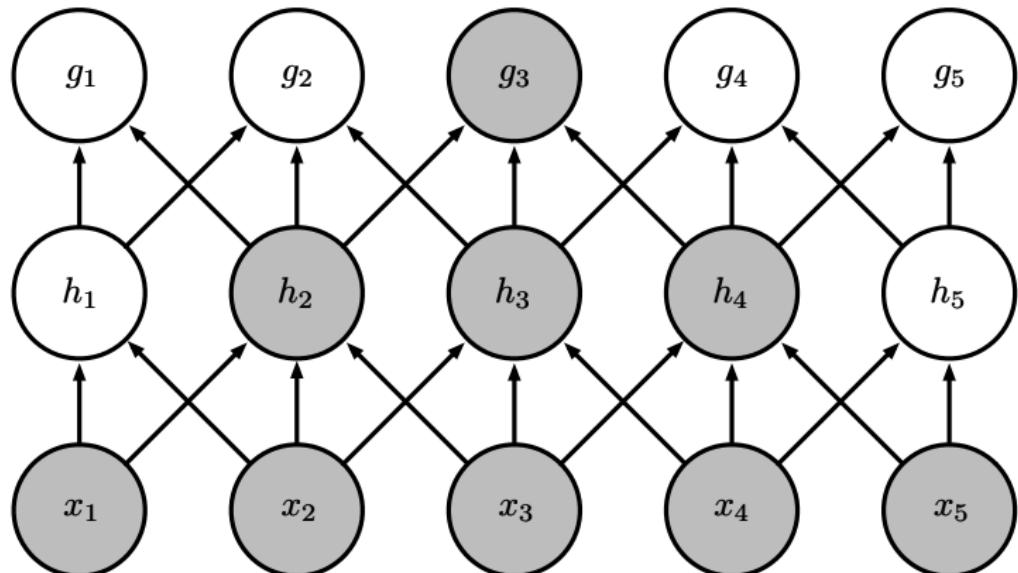
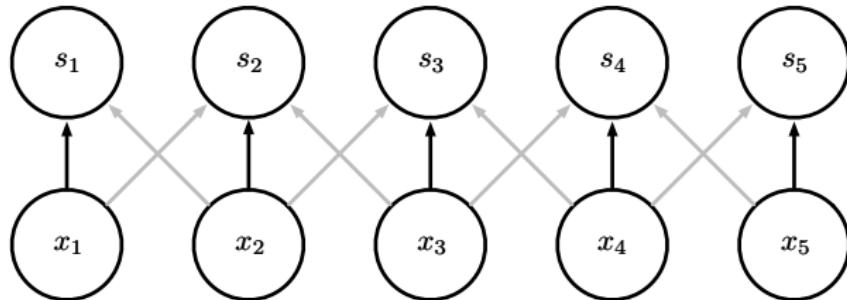


Figure 9.4

# Parameter Sharing

Convolution  
shares the same  
parameters  
across all spatial  
locations



Traditional  
matrix  
multiplication  
does not share  
any parameters

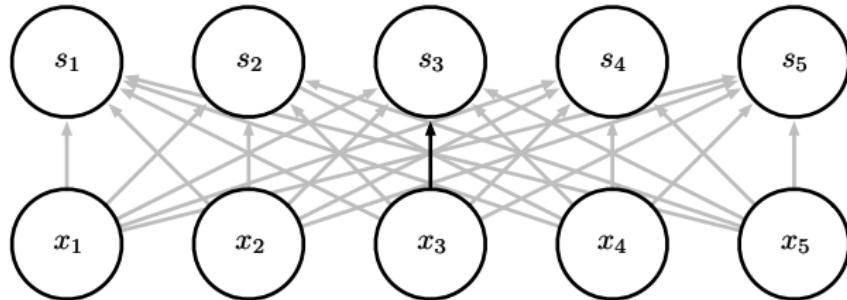
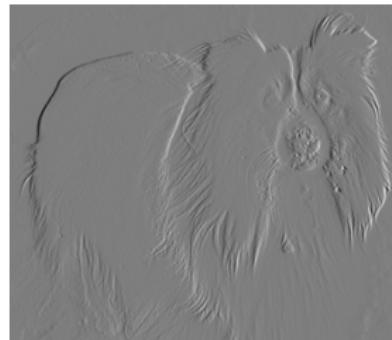


Figure 9.5

# A simple example why filtering is useful



Input

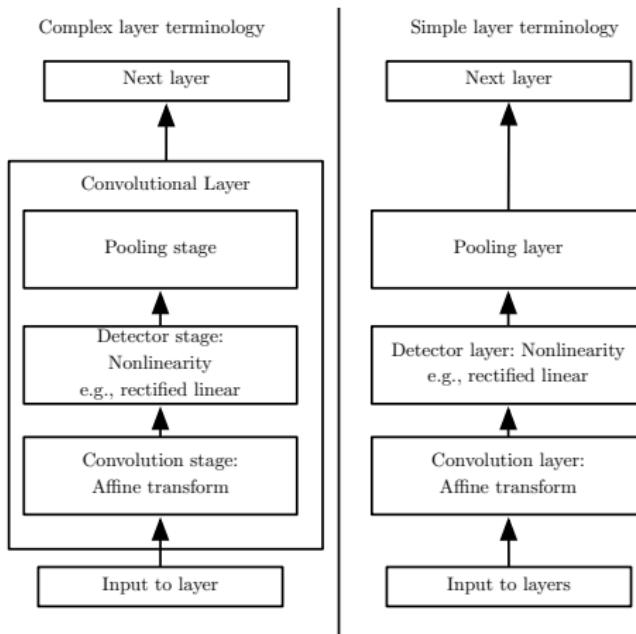


Output

1	-1
---	----

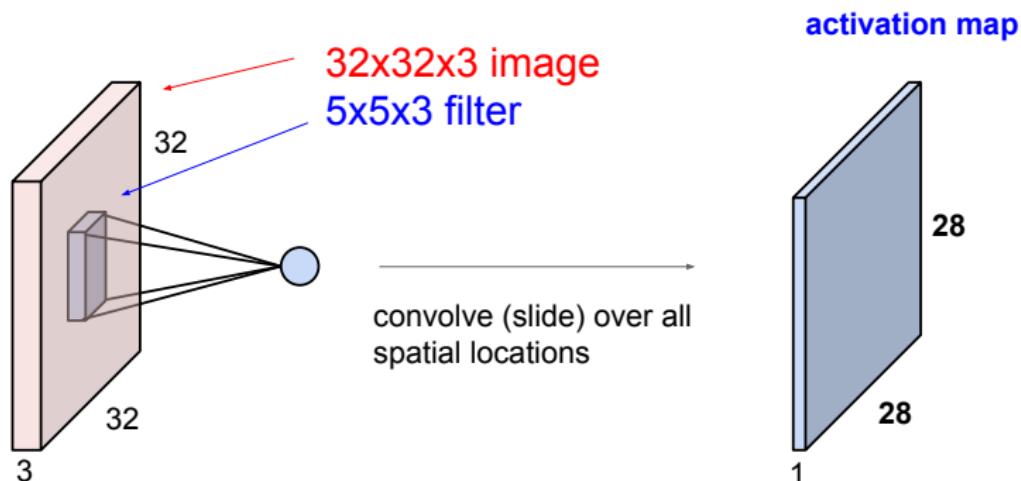
Kernel

# Compounding layers as convolution neural nets (CNNs)



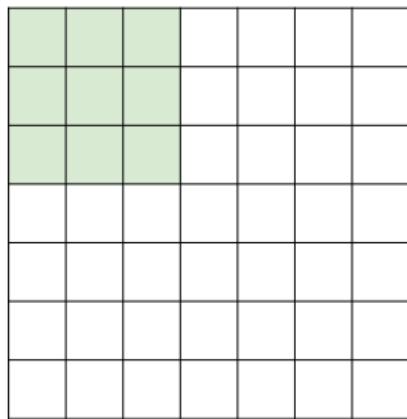
# Stride and padding for convolution

A closer look at spatial dimensions:



A closer look at spatial dimensions:

7

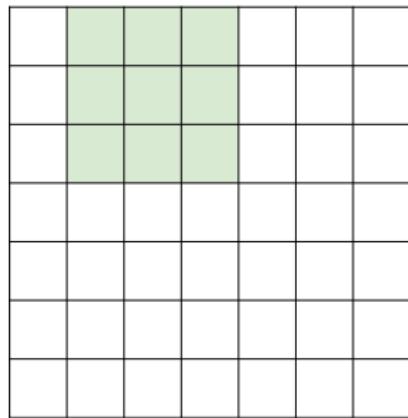


7x7 input (spatially)  
assume 3x3 filter

7

A closer look at spatial dimensions:

7

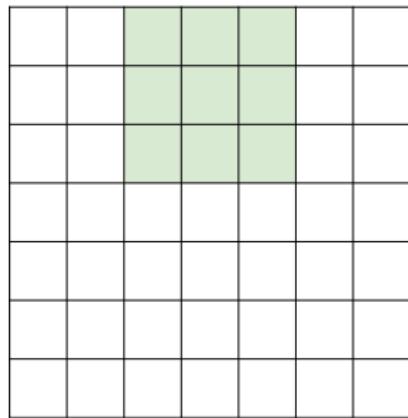


7x7 input (spatially)  
assume 3x3 filter

7

A closer look at spatial dimensions:

7

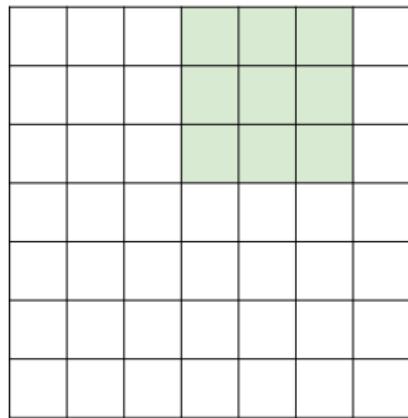


7

7x7 input (spatially)  
assume 3x3 filter

A closer look at spatial dimensions:

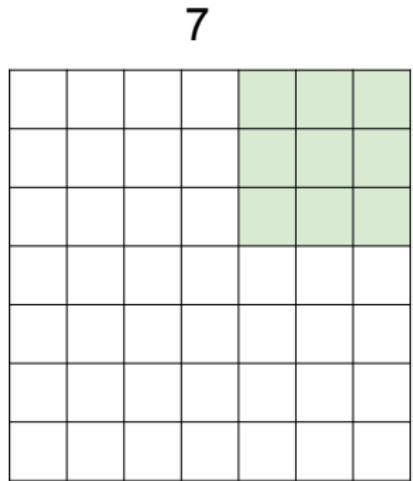
7



7x7 input (spatially)  
assume 3x3 filter

7

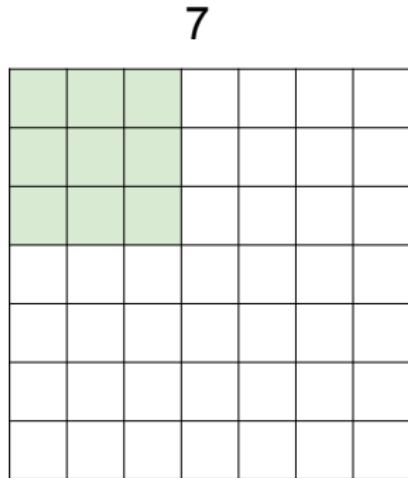
A closer look at spatial dimensions:



7x7 input (spatially)  
assume 3x3 filter

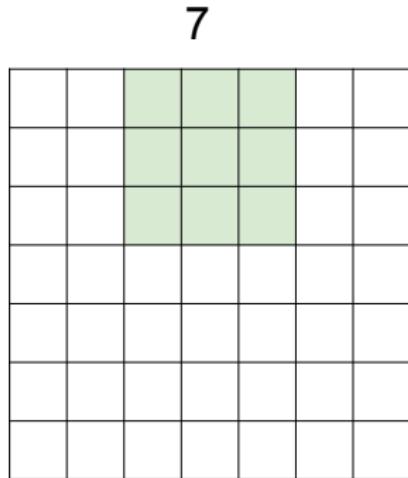
=> **5x5 output**

A closer look at spatial dimensions:



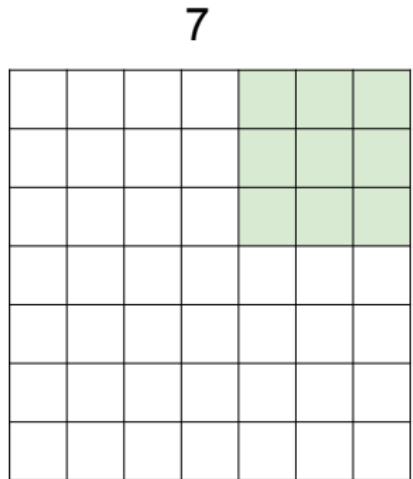
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

A closer look at spatial dimensions:



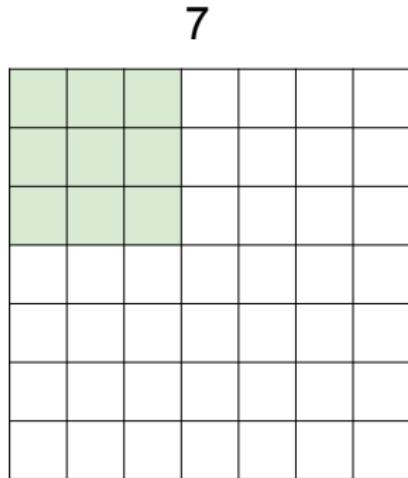
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

A closer look at spatial dimensions:



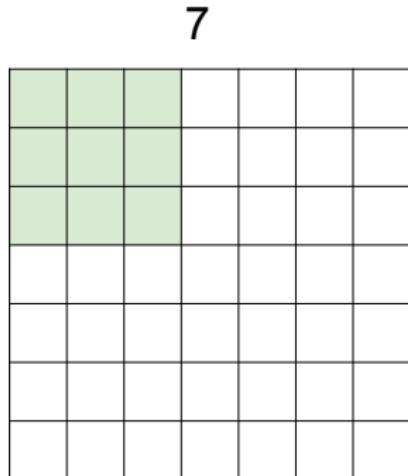
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**  
**=> 3x3 output!**

A closer look at spatial dimensions:



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

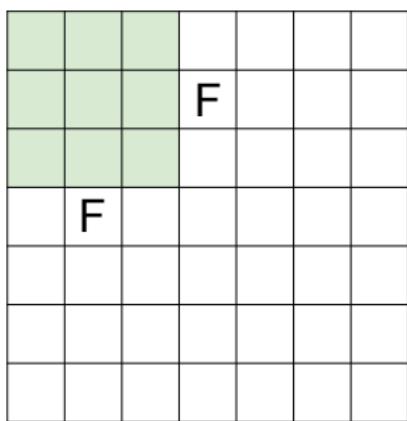
A closer look at spatial dimensions:



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

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

N



N

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

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

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

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

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

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 FxF, and zero-padding with  $(F-1)/2$ . (will preserve size spatially)

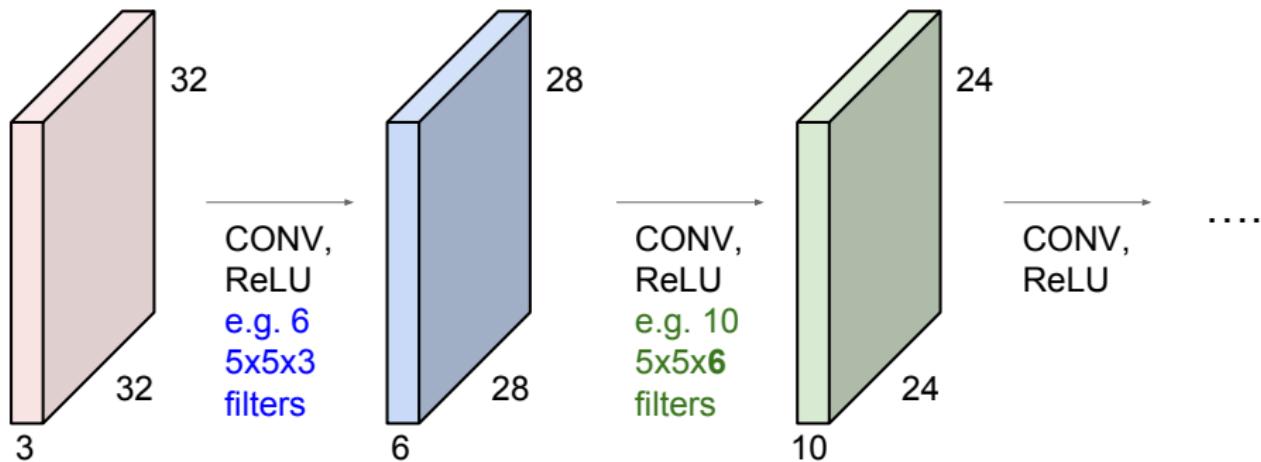
e.g.  $F = 3 \Rightarrow$  zero pad with 1

$F = 5 \Rightarrow$  zero pad with 2

$F = 7 \Rightarrow$  zero pad with 3

## Remember back to...

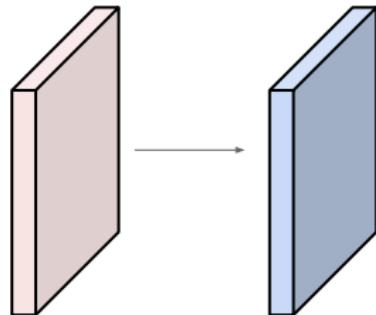
E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially!  
(32  $\rightarrow$  28  $\rightarrow$  24 ...). Shrinking too fast is not good, doesn't work well.



Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

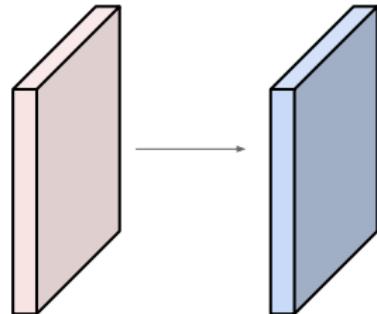


Output volume size: ?

Examples time:

Input volume: **32x32x3**

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



Output volume size:

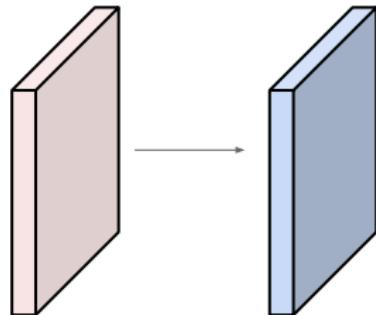
$(32+2*2-5)/1+1 = 32$  spatially, so

**32x32x10**

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

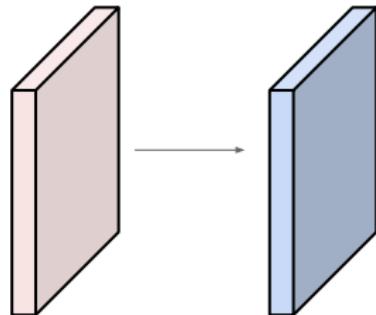


Number of parameters in this layer?

Examples time:

Input volume: **32x32x3**

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



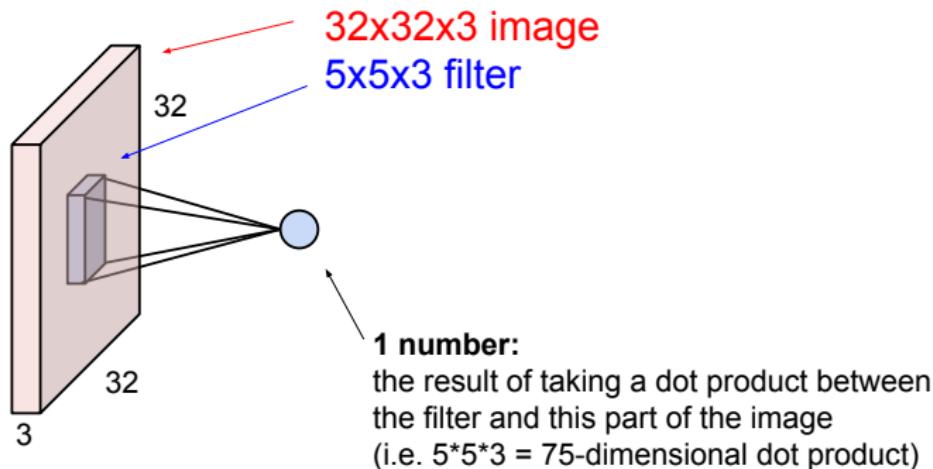
Number of parameters in this layer?

each filter has  $5 \times 5 \times 3 + 1 = 76$  params      (+1 for bias)

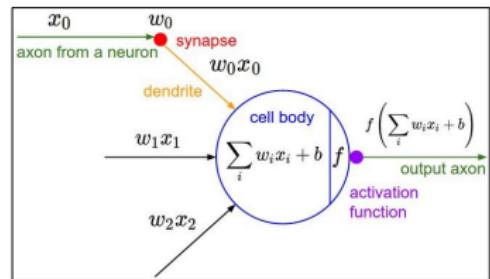
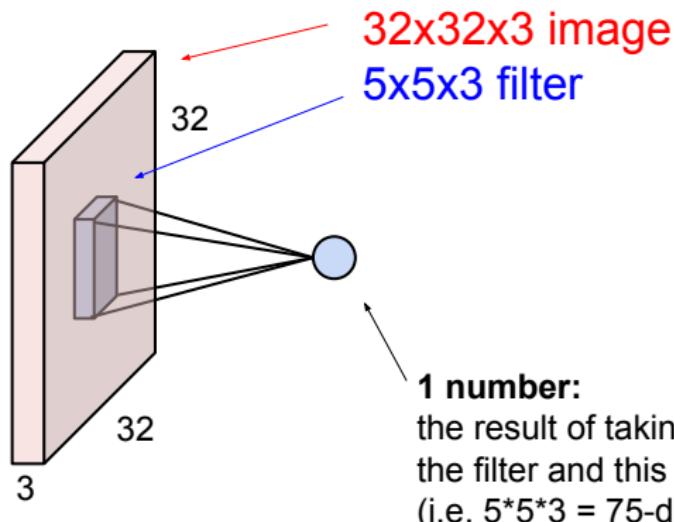
$$\Rightarrow 76 \times 10 = 760$$

# Adding more filters

The brain/neuron view of CONV Layer

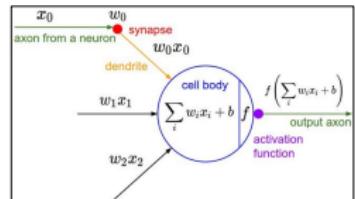
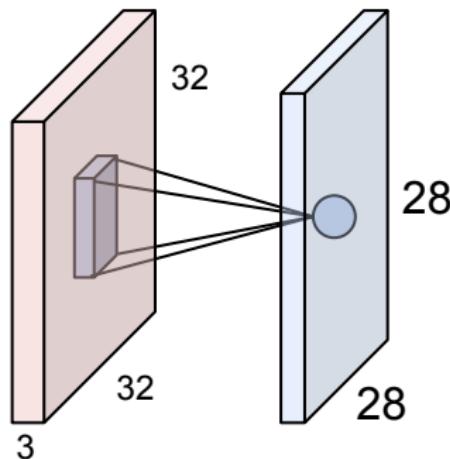


## The brain/neuron view of CONV Layer



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

## The brain/neuron view of CONV Layer

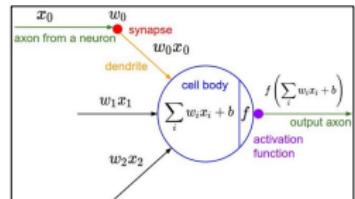
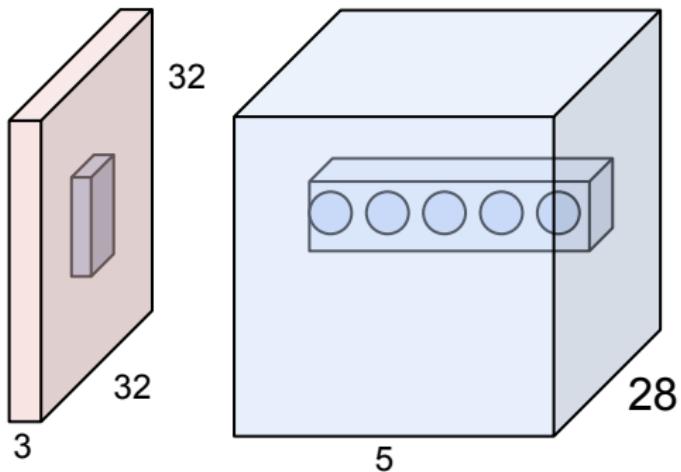


An activation map is a  $28 \times 28$  sheet of neuron outputs:

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

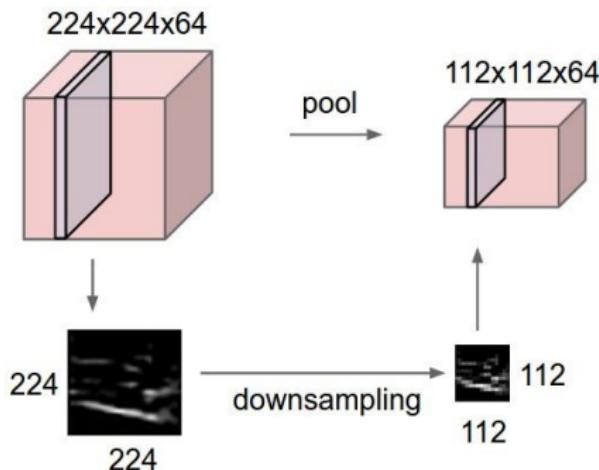


E.g. with 5 filters,  
CONV layer consists of  
neurons arranged in a 3D grid  
( $28 \times 28 \times 5$ )

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

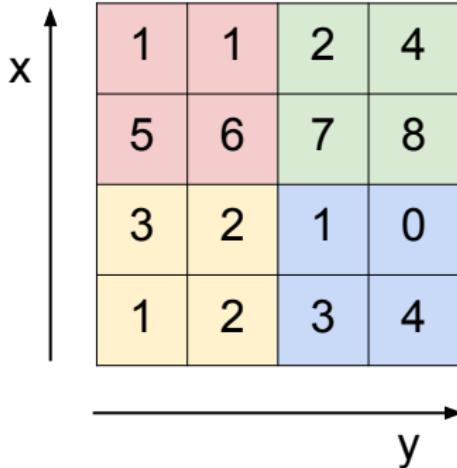
# Pooling layer

- makes the representations smaller and more manageable
- operates over each activation map independently:



# MAX POOLING

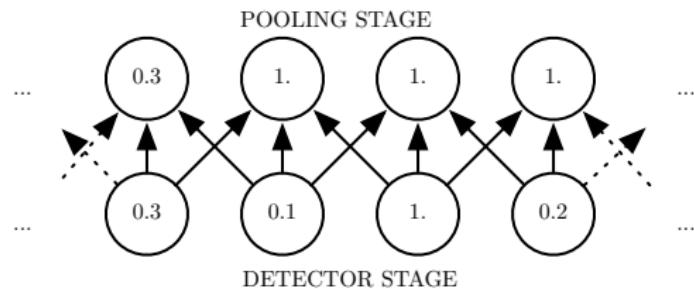
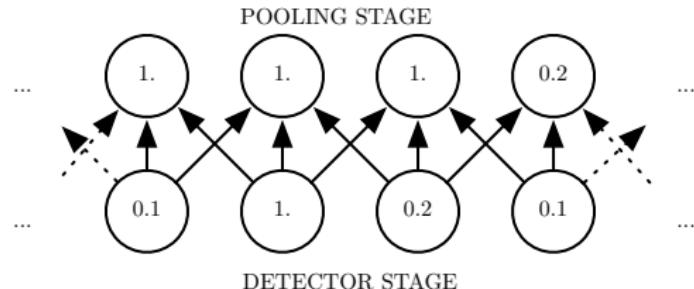
Single depth slice



max pool with 2x2 filters  
and stride 2

6	8
3	4

# Max Pooling and Invariance to Translation



# Cross-Channel Pooling and Invariance to Learned Transformations

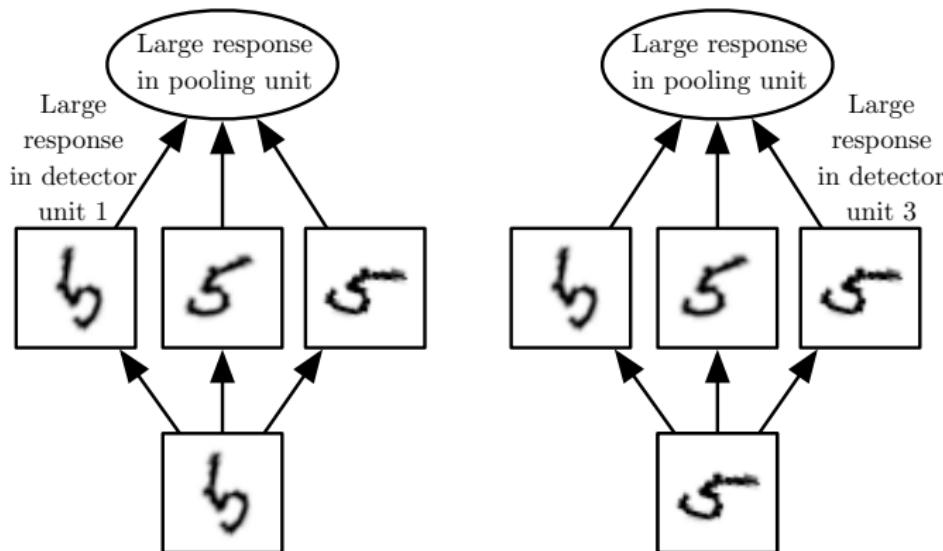


Figure 9.9

# Pooling with Downsampling

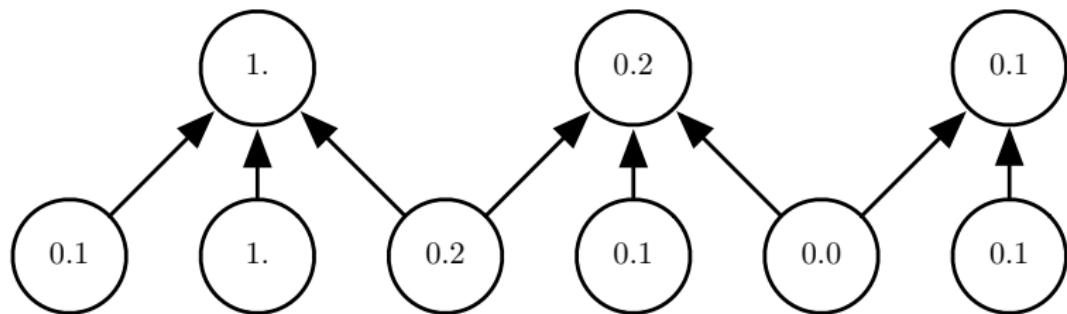


Figure 9.10

# Putting things together to classify

## 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  
 **$[(CONV-RELU)^*N-POOL?]^*M-(FC-RELU)^*K, SOFTMAX$**   
where N is usually up to ~5, M is large,  $0 \leq K \leq 2$ .
  - but recent advances such as ResNet/GoogLeNet challenge this paradigm

# Demo

<http://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

# How do we learn the parameters of the filters from data?

## Error backpropagation