

# CSCI567 Machine Learning (Spring 2021)

Sirisha Rambhatla

University of Southern California

Feb 17, 2021

# Outline

- 1 Logistics
- 2 Review of last lecture
- 3 Convolutional neural networks (ConvNets/CNNs)

# Outline

- 1 Logistics
- 2 Review of last lecture
- 3 Convolutional neural networks (ConvNets/CNNs)

# Logistics

- Sign-up with your group members for the project!

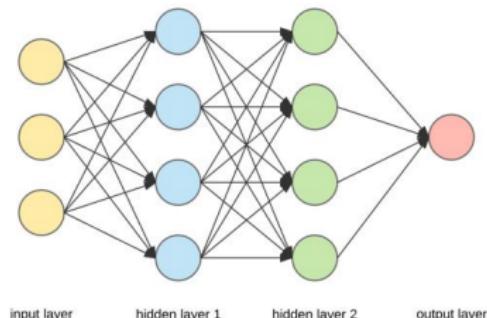
# Outline

- 1 Logistics
- 2 Review of last lecture
- 3 Convolutional neural networks (ConvNets/CNNs)

# Math formulation of neural nets

An L-layer neural net can be written as

$$\mathbf{f}(\mathbf{x}) = \mathbf{h}_L (\mathbf{W}_L \mathbf{h}_{L-1} (\mathbf{W}_{L-1} \cdots \mathbf{h}_1 (\mathbf{W}_1 \mathbf{x})))$$



To ease notation, for a given input  $\mathbf{x}$ , define recursively

$$\mathbf{o}_0 = \mathbf{x}, \quad \mathbf{a}_\ell = \mathbf{W}_\ell \mathbf{o}_{\ell-1}, \quad \mathbf{o}_\ell = \mathbf{h}_\ell(\mathbf{a}_\ell) \quad (\ell = 1, \dots, L)$$

where

- $\mathbf{W}_\ell \in \mathbb{R}^{D_\ell \times D_{\ell-1}}$  is the weights for layer  $\ell$
- $D_0 = D, D_1, \dots, D_L$  are numbers of neurons at each layer
- $\mathbf{a}_\ell \in \mathbb{R}^{D_\ell}$  is input to layer  $\ell$
- $\mathbf{o}_\ell \in \mathbb{R}^{D_\ell}$  is output to layer  $\ell$
- $\mathbf{h} : \mathbb{R}^{D_\ell} \rightarrow \mathbb{R}^{D_\ell}$  is activation functions at layer  $\ell$

# Backprop = SGD for neural nets

The **backpropagation** algorithm (**Backprop**)

Initialize  $\mathbf{W}_1, \dots, \mathbf{W}_L$  (all  $\mathbf{0}$  or randomly). Repeat:

- ① randomly pick one data point  $n \in [N]$
- ② **forward propagation**: for each layer  $\ell = 1, \dots, L$ 
  - compute  $\mathbf{a}_\ell = \mathbf{W}_\ell \mathbf{o}_{\ell-1}$  and  $\mathbf{o}_\ell = \mathbf{h}_\ell(\mathbf{a}_\ell)$   $(\mathbf{o}_0 = \mathbf{x}_n)$
- ③ **backward propagation**: for each  $\ell = L, \dots, 1$ 
  - compute

$$\frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_\ell} = \begin{cases} \left( \mathbf{W}_{\ell+1}^T \frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_{\ell+1}} \right) \circ \mathbf{h}'_\ell(\mathbf{a}_\ell) & \text{if } \ell < L \\ 2(\mathbf{h}_L(\mathbf{a}_L) - \mathbf{y}_n) \circ \mathbf{h}'_L(\mathbf{a}_L) & \text{else} \end{cases}$$

- update weights

$$\mathbf{W}_\ell \leftarrow \mathbf{W}_\ell - \eta \frac{\partial \mathcal{E}_n}{\partial \mathbf{W}_\ell} = \mathbf{W}_\ell - \eta \frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_\ell} \mathbf{o}_{\ell-1}^T$$

*Think about how to do the last two steps properly!*

# Outline

- 1 Logistics
- 2 Review of last lecture
- 3 Convolutional neural networks (ConvNets/CNNs)
  - Motivation
  - Architecture

# Acknowledgements

Not much math, a lot of empirical intuitions

## Acknowledgements

Not much math, a lot of empirical intuitions

The materials borrow heavily from the following sources:

- Stanford Course Cs231n: <http://cs231n.stanford.edu/>
- Dr. Ian Goodfellow's lectures on deep learning:  
<http://deeplearningbook.org>

Both website provides tons of useful resources: notes, demos, videos, etc.

# Image Classification: A core task in Computer Vision



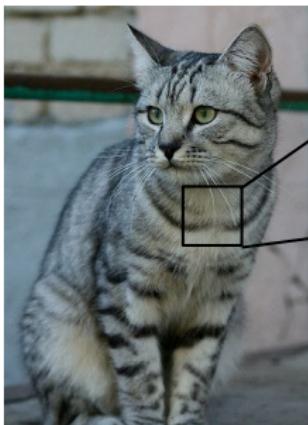
This image by nikita is  
licensed under CC-BY 2.0

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



cat

# The Problem: Semantic Gap



|  |
|--|
| [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] [106 91 61 64 69 91 88 85 181 187 109 98 75 84 96 95]            |
| [1] [114 180 85 95 59 61 68 54 87 120 126 98 74 84 96 95]            |
| [1] [133 137 144 183 86 81 88 52 74 64 80 95 82 83 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 82 74 65 77 78]  |
| [1] [89 101 111 121 147 133 127 126 131 111 96 89 75 61 64 72 84]    |
| [1] [63 77 86 81 77 79 182 123 117 115 117 125 139 115 87]           |
| [1] [62 65 82 82 78 71 81 181 124 126 119 101 187 114 131 119]       |
| [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 97 107 111 121 126 124 124 124 124 124 124 124 124 124 124] |
| [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] [128 122 96 117 158 144 120 115 184 187 182 93 87 81 73 79]      |
| [1] [123 128 111 101 88 83 111 152 144 144 144 144 144 144 144 144]  |
| [1] [122 164 148 183 71 56 78 83 93 103 119 139 182 61 69 84]        |

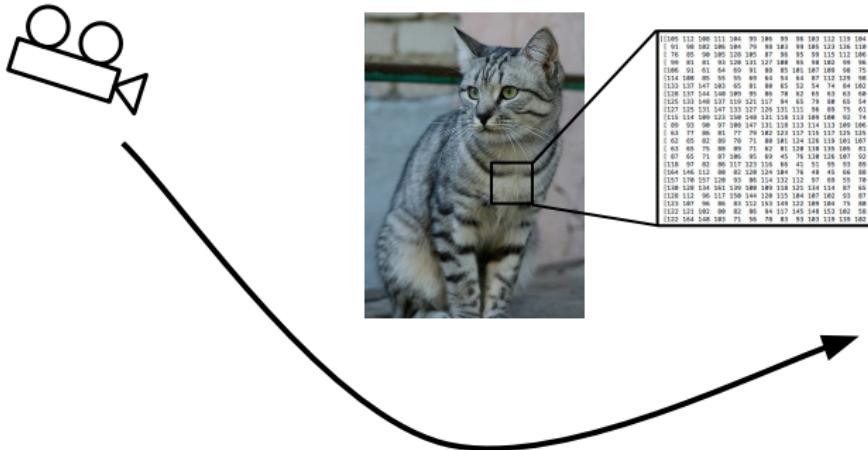
What the computer sees

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

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

This image by nikita is  
licensed under CC-BY 2.0

## Challenges: Viewpoint variation



All pixels change when  
the camera moves!

This image by Nikita is  
licensed under CC-BY 2.0

## Challenges: Illumination



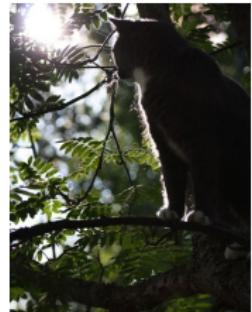
[This image is CC0 1.0 public domain](#)



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

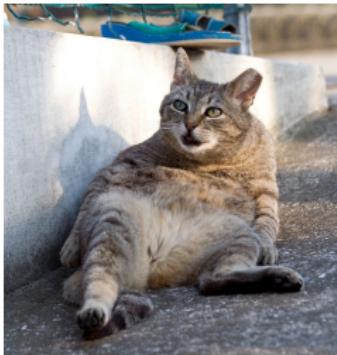


[This image is CC0 1.0 public domain](#)



[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

## The key challenge

How to train a model that can tolerate all those variations?

# Fundamental problems in vision

## The key challenge

How to train a model that can tolerate all those variations?

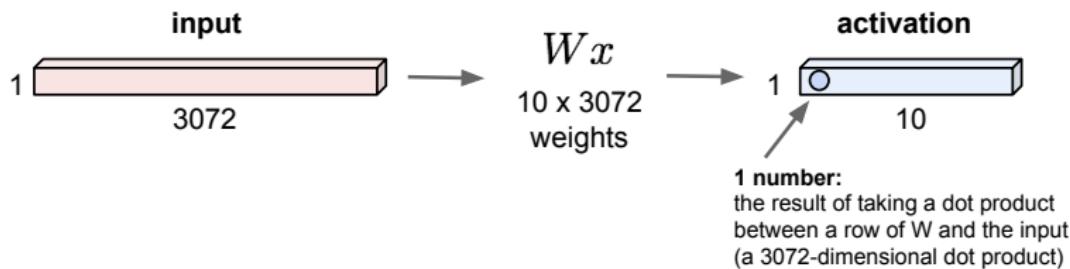
## Main ideas

- need a lot of data that exhibits those variations
- need more specialized models to capture the invariance

# Issues of standard NN for image inputs

## Fully Connected Layer

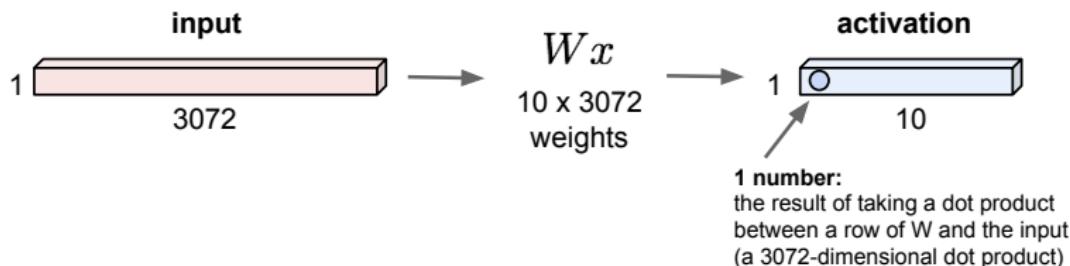
32x32x3 image  $\rightarrow$  stretch to 3072 x 1



# Issues of standard NN for image inputs

## Fully Connected Layer

32x32x3 image  $\rightarrow$  stretch to 3072 x 1



# Solution: Convolutional Neural Net (ConvNet/CNN)

A special case of fully connected neural nets

# Solution: Convolutional Neural Net (ConvNet/CNN)

A special case of fully connected neural nets

- usually consist of **convolution layers**, ReLU layers, **pooling layers**, and regular fully connected layers

# Solution: Convolutional Neural Net (ConvNet/CNN)

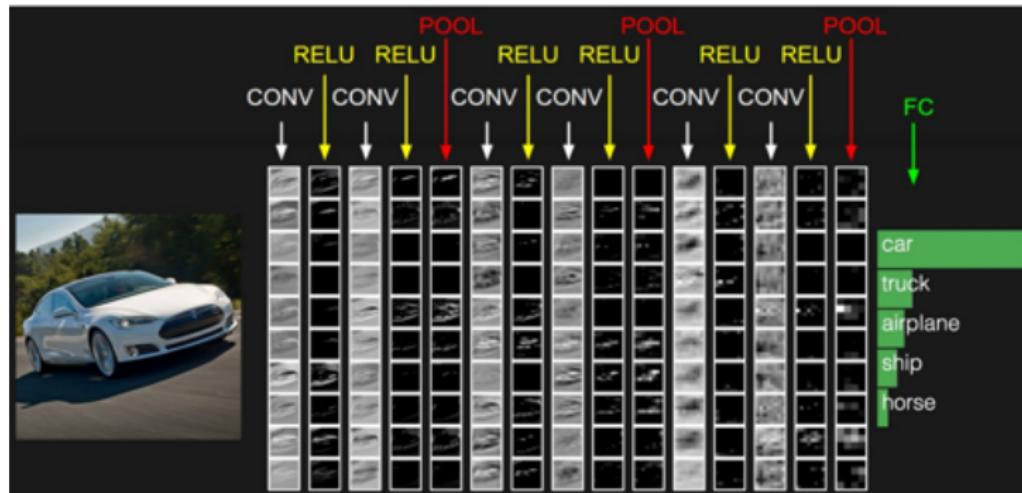
A special case of fully connected neural nets

- usually consist of **convolution layers**, ReLU layers, **pooling layers**, and regular fully connected layers
- key idea: *learning from low-level to high-level features*

# Solution: Convolutional Neural Net (ConvNet/CNN)

A special case of fully connected neural nets

- usually consist of **convolution layers**, ReLU layers, **pooling layers**, and regular fully connected layers
- key idea: *learning from low-level to high-level features*

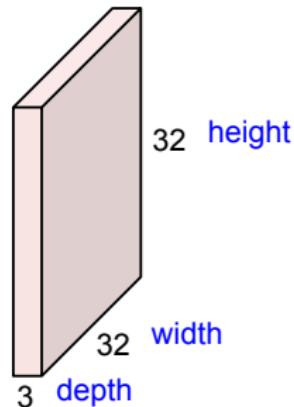


# Convolution layer

Arrange neurons as a **3D volume** naturally

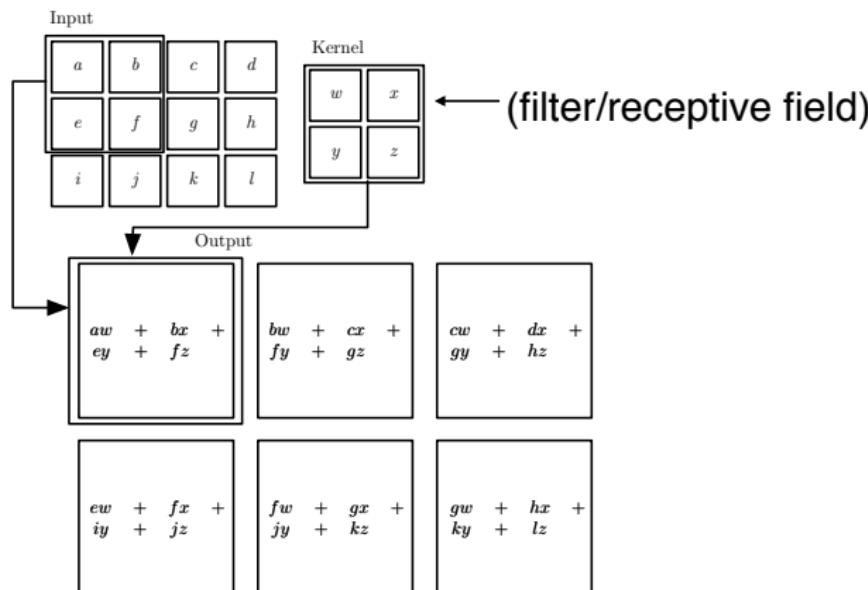
## Convolution Layer

32x32x3 image -> preserve spatial structure



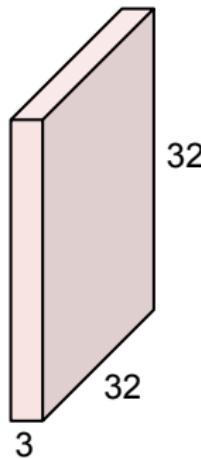
# Convolution

## 2D Convolution



# Convolution Layer

32x32x3 image



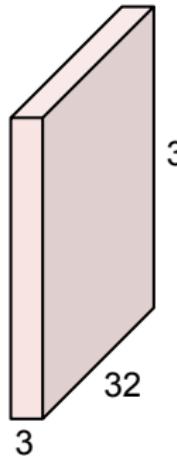
5x5x3 filter



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

# Convolution Layer

32x32x3 image



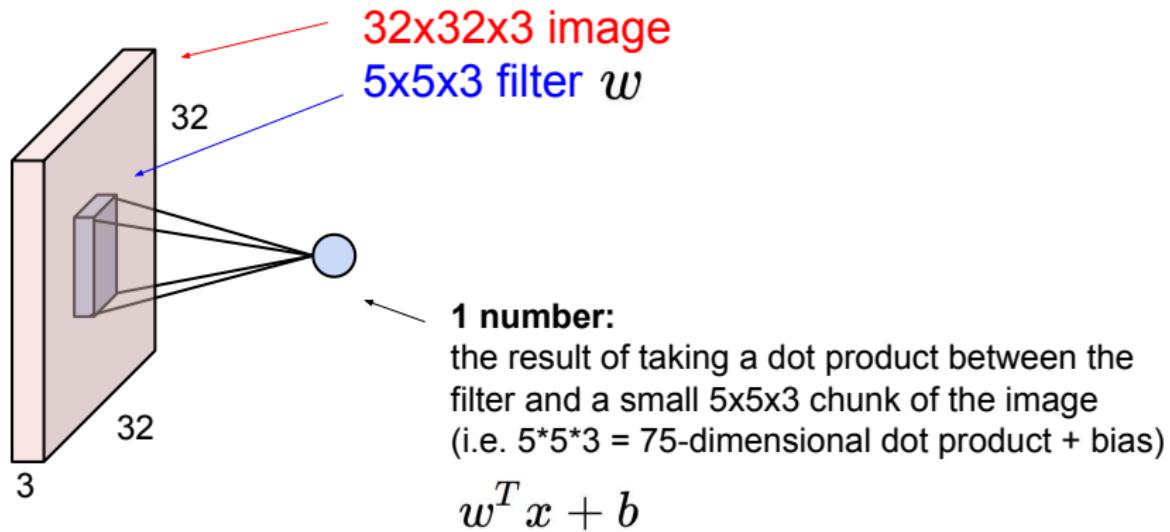
5x5x3 filter



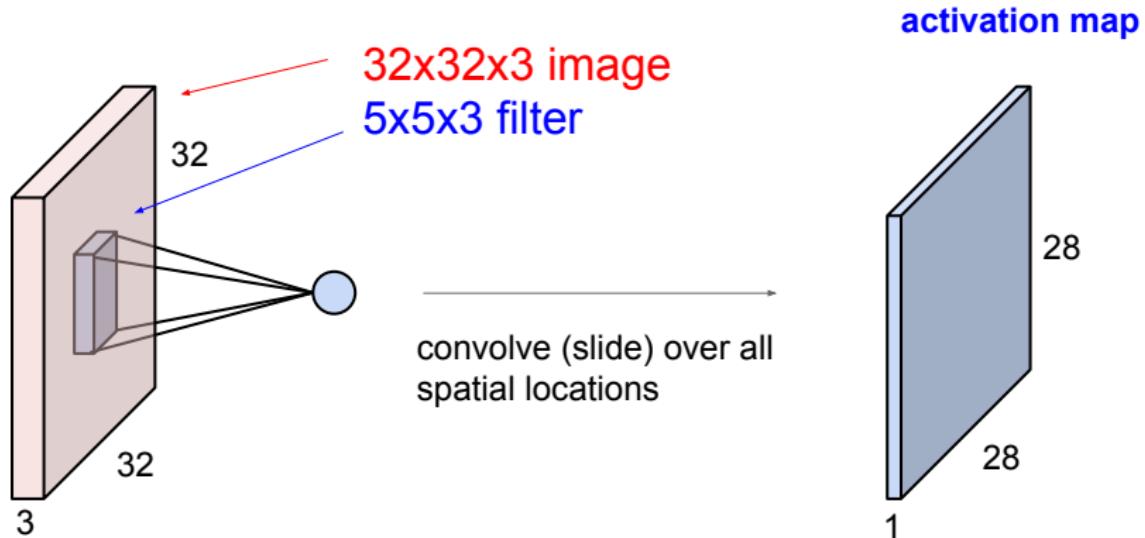
Filters always extend the full depth of the input volume

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

# Convolution Layer

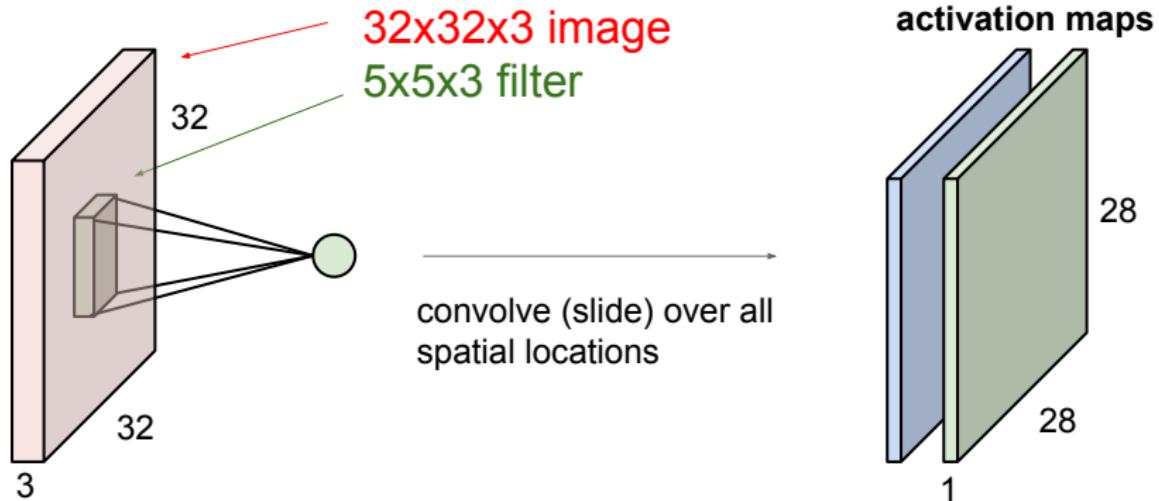


# Convolution Layer

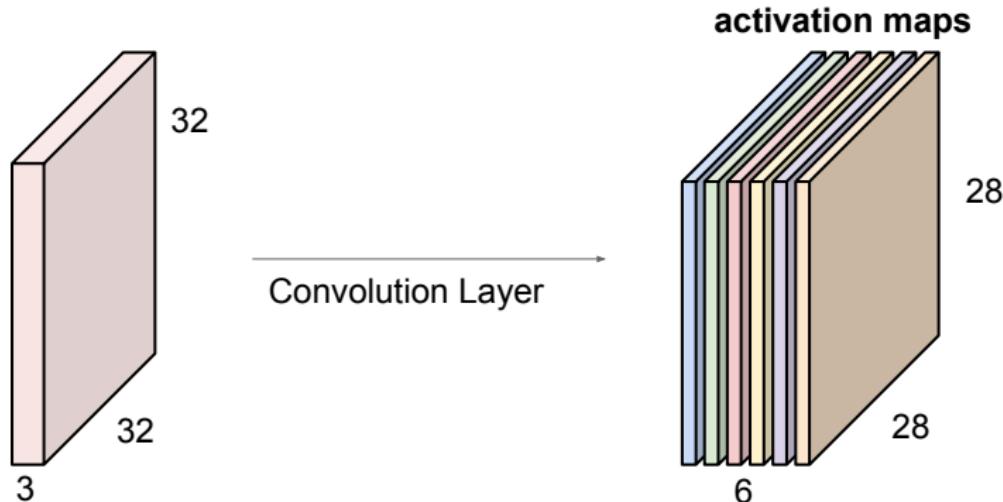


# Convolution Layer

consider a second, green filter

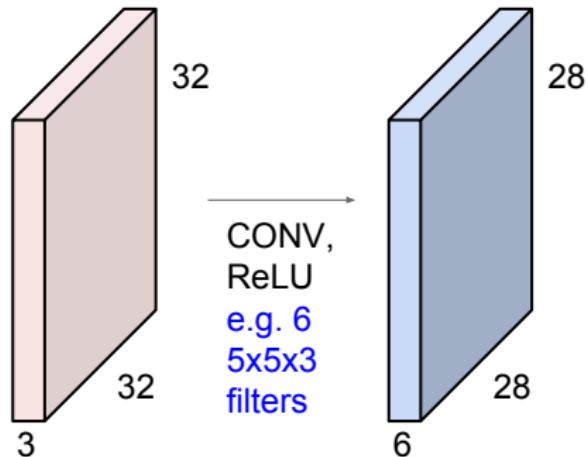


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

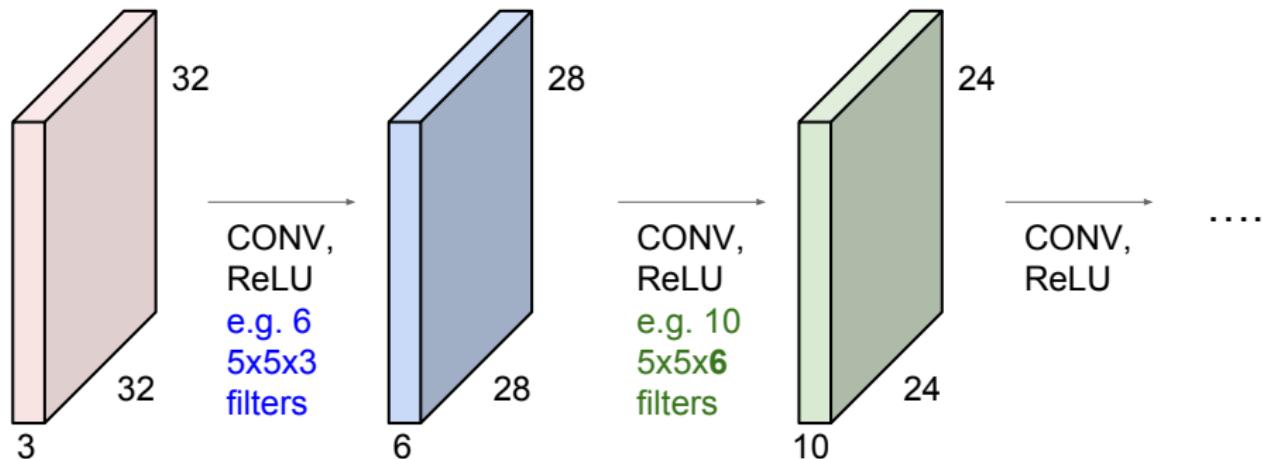


We stack these up to get a “new image” of size  $28 \times 28 \times 6$ !

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



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



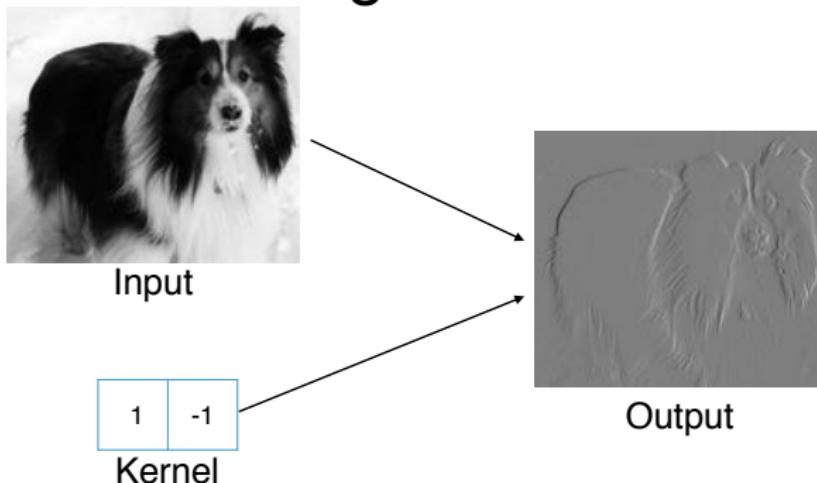
## Why convolution makes sense?

Main idea: **if a filter is useful at one location, it should be useful at other locations.**

## Why convolution makes sense?

Main idea: **if a filter is useful at one location, it should be useful at other locations.**

### A simple example why filtering is useful



# Connection to fully connected NNs

A convolution layer is a special case of a fully connected layer:

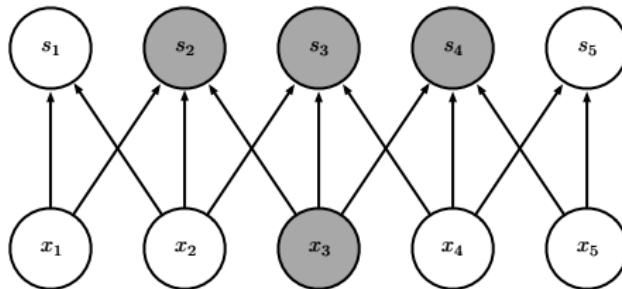
# Connection to fully connected NNs

A convolution layer is a special case of a fully connected layer:

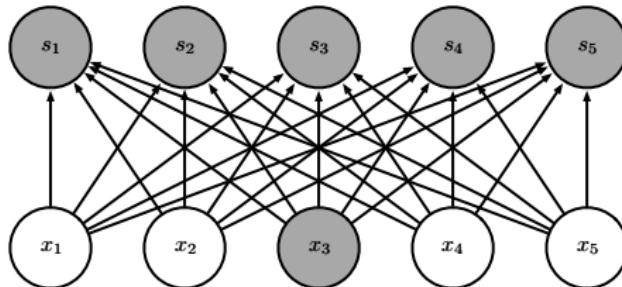
- filter = weights with **sparse connection**

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

Sparse  
connections  
due to small  
convolution  
kernel

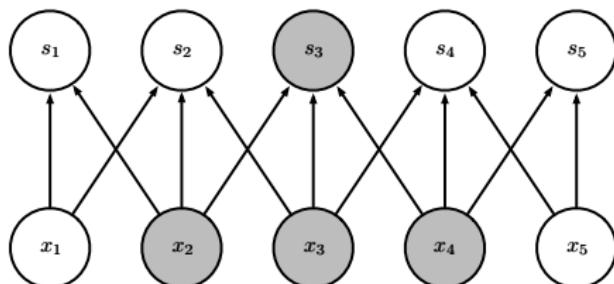


Dense  
connections



# Sparse connectivity: being affected by less

Sparse connections due to small convolution kernel



Dense connections

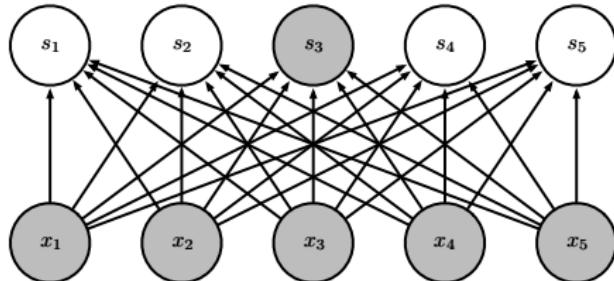


Figure 9.3

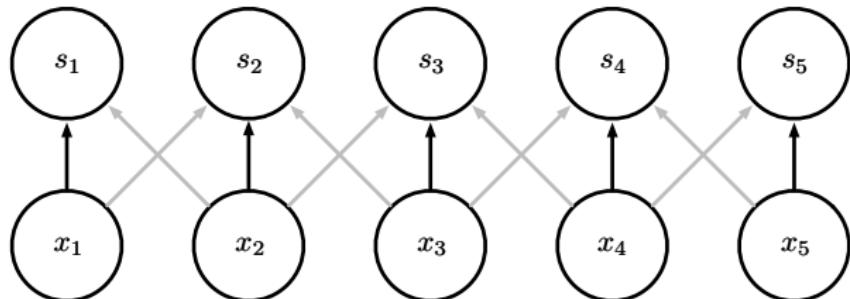
# Connection to fully connected NNs

A convolution layer is a special case of a fully connected layer:

- filter = weights with **sparse connection**
- **parameters sharing**

# Parameter Sharing

Convolution  
shares the same  
parameters  
across all spatial  
locations



Traditional  
matrix  
multiplication  
does not share  
any parameters

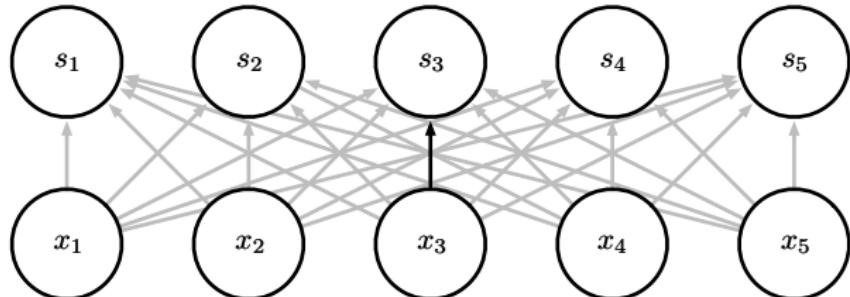


Figure 9.5

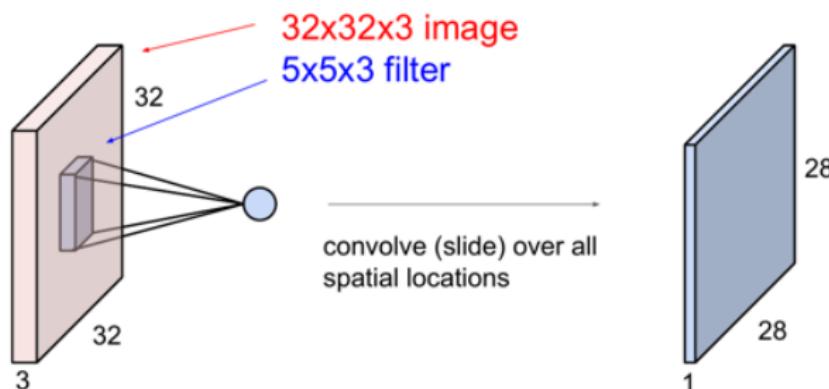
## Connection to fully connected NNs

A convolution layer is a special case of a fully connected layer:

- filter = weights with **sparse connection**
- **parameters sharing**

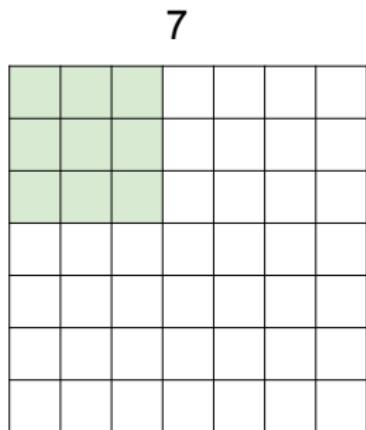
*Much less parameters!* Example (ignore bias terms):

- FC:  $(32 \times 32 \times 3) \times (28 \times 28) \approx 2.4M$
- CNN:  $5 \times 5 \times 3 = 75$



# Spatial arrangement: stride and padding

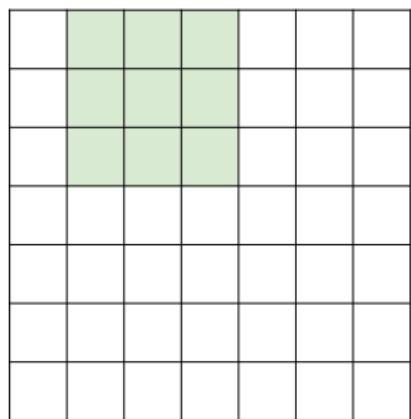
A closer look at spatial dimensions:



7x7 input (spatially)  
assume 3x3 filter

A closer look at spatial dimensions:

7

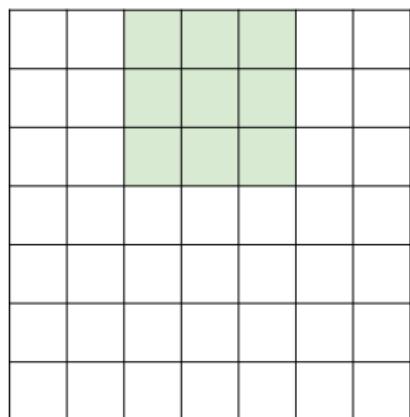


7x7 input (spatially)  
assume 3x3 filter

7

A closer look at spatial dimensions:

7

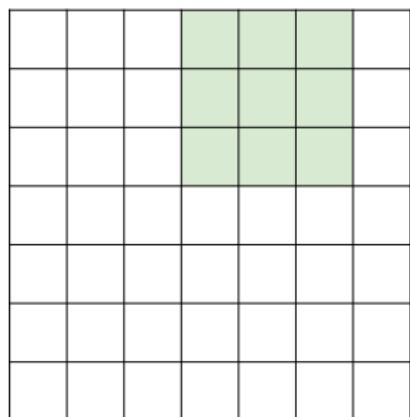


7x7 input (spatially)  
assume 3x3 filter

7

A closer look at spatial dimensions:

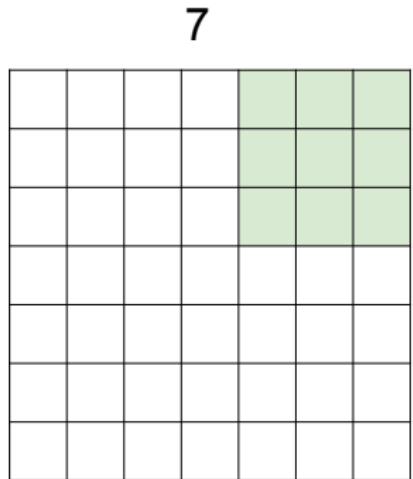
7



7x7 input (spatially)  
assume 3x3 filter

7

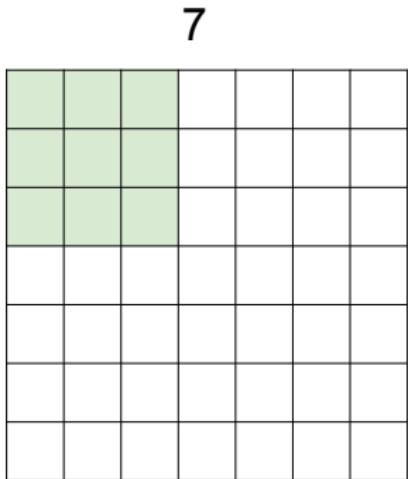
A closer look at spatial dimensions:



7x7 input (spatially)  
assume 3x3 filter

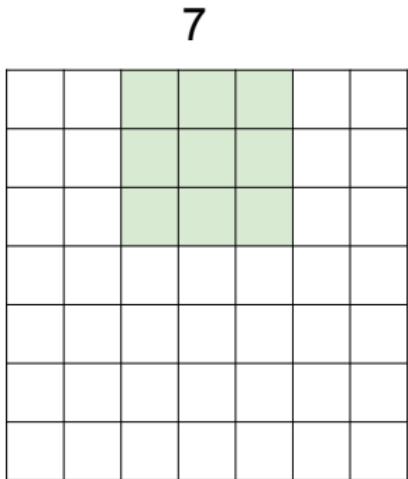
=> **5x5 output**

A closer look at spatial dimensions:



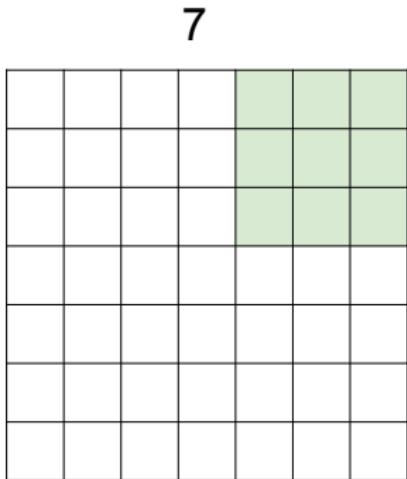
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

A closer look at spatial dimensions:



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

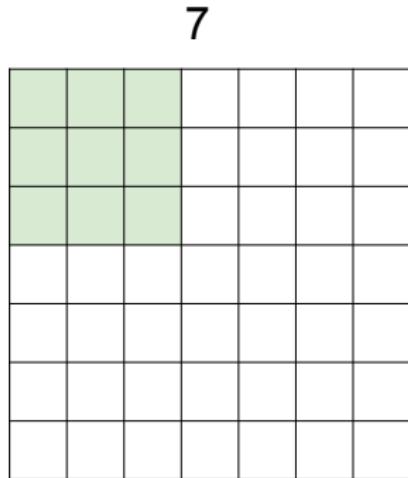
A closer look at spatial dimensions:



7

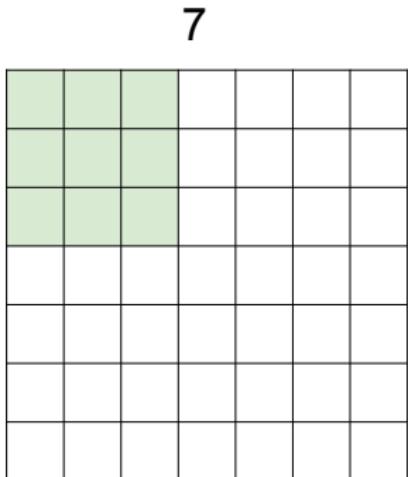
7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**  
**=> 3x3 output!**

A closer look at spatial dimensions:



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

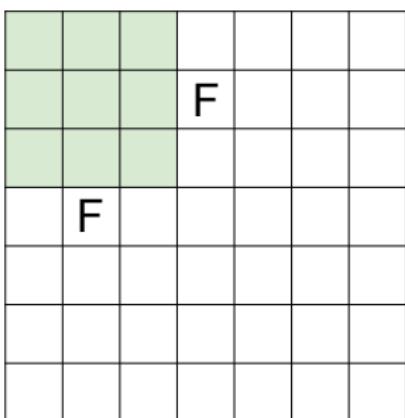
A closer look at spatial dimensions:



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)

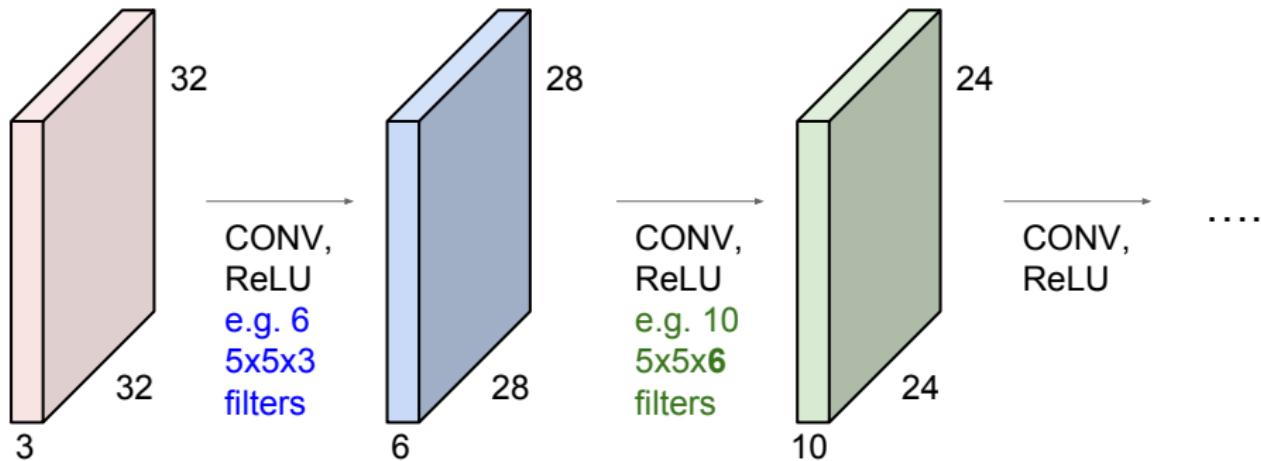
e.g.  $F = 3 \Rightarrow$  zero pad with 1

$F = 5 \Rightarrow$  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.



## Summary for convolution layer

**Input:** a volume of size  $W_1 \times H_1 \times D_1$

## Summary for convolution layer

**Input:** a volume of size  $W_1 \times H_1 \times D_1$

**Hyperparameters:**

- $K$  filters of size  $F \times F$
- stride  $S$
- amount of zero padding  $P$  (for one side)

## Summary for convolution layer

**Input:** a volume of size  $W_1 \times H_1 \times D_1$

**Hyperparameters:**

- $K$  filters of size  $F \times F$
- stride  $S$
- amount of zero padding  $P$  (for one side)

**Output:** a volume of size  $W_2 \times H_2 \times D_2$  where

- $W_2 =$
- $H_2 =$
- $D_2 =$

## Summary for convolution layer

**Input:** a volume of size  $W_1 \times H_1 \times D_1$

**Hyperparameters:**

- $K$  filters of size  $F \times F$
- stride  $S$
- amount of zero padding  $P$  (for one side)

**Output:** a volume of size  $W_2 \times H_2 \times D_2$  where

- $W_2 = (W_1 + 2P - F)/S + 1$
- $H_2 =$
- $D_2 =$

## Summary for convolution layer

**Input:** a volume of size  $W_1 \times H_1 \times D_1$

**Hyperparameters:**

- $K$  filters of size  $F \times F$
- stride  $S$
- amount of zero padding  $P$  (for one side)

**Output:** a volume of size  $W_2 \times H_2 \times D_2$  where

- $W_2 = (W_1 + 2P - F)/S + 1$
- $H_2 = (H_1 + 2P - F)/S + 1$
- $D_2 =$

## Summary for convolution layer

**Input:** a volume of size  $W_1 \times H_1 \times D_1$

**Hyperparameters:**

- $K$  filters of size  $F \times F$
- stride  $S$
- amount of zero padding  $P$  (for one side)

**Output:** a volume of size  $W_2 \times H_2 \times D_2$  where

- $W_2 = (W_1 + 2P - F)/S + 1$
- $H_2 = (H_1 + 2P - F)/S + 1$
- $D_2 = K$

## Summary for convolution layer

**Input:** a volume of size  $W_1 \times H_1 \times D_1$

**Hyperparameters:**

- $K$  filters of size  $F \times F$
- stride  $S$
- amount of zero padding  $P$  (for one side)

**Output:** a volume of size  $W_2 \times H_2 \times D_2$  where

- $W_2 = (W_1 + 2P - F)/S + 1$
- $H_2 = (H_1 + 2P - F)/S + 1$
- $D_2 = K$

**#parameters:**  $(F \times F \times D_1 + 1) \times K$  weights

## Summary for convolution layer

**Input:** a volume of size  $W_1 \times H_1 \times D_1$

**Hyperparameters:**

- $K$  filters of size  $F \times F$
- stride  $S$
- amount of zero padding  $P$  (for one side)

**Output:** a volume of size  $W_2 \times H_2 \times D_2$  where

- $W_2 = (W_1 + 2P - F)/S + 1$
- $H_2 = (H_1 + 2P - F)/S + 1$
- $D_2 = K$

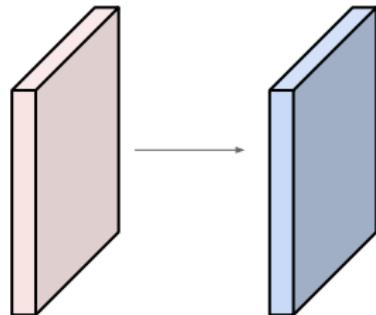
**#parameters:**  $(F \times F \times D_1 + 1) \times K$  weights

**Common setting:**  $F = 3, S = P = 1$

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

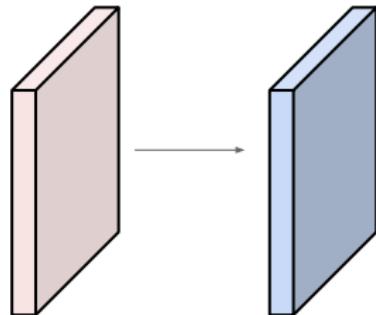


Output volume size: ?

Examples time:

Input volume: **32x32x3**

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



Output volume size:

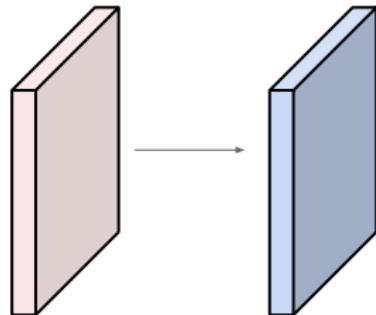
$(32+2*2-5)/1+1 = 32$  spatially, so

**32x32x10**

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

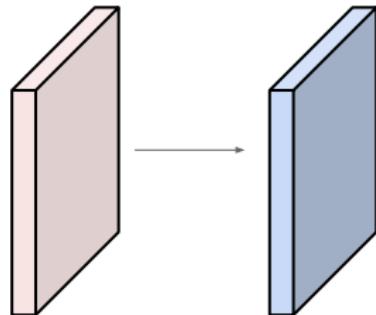


Number of parameters in this layer?

Examples time:

Input volume: **32x32x3**

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



Number of parameters in this layer?

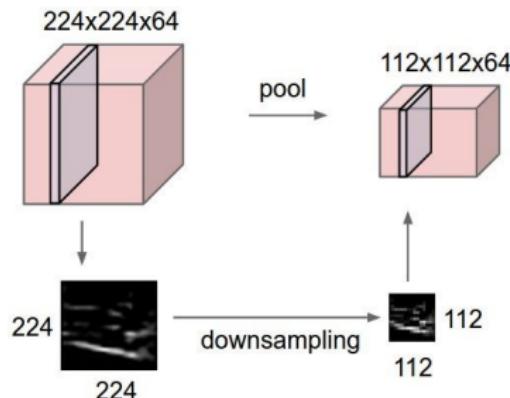
each filter has  $5*5*3 + 1 = 76$  params      (+1 for bias)

$$\Rightarrow 76 * 10 = 760$$

# Another element: pooling

## Pooling layer

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



# Pooling

Similar to a filter, except

- depth is always 1

# Pooling

Similar to a filter, except

- depth is always 1
- different operations: average, L2-norm, max

# Pooling

Similar to a filter, except

- depth is always 1
- different operations: average, L2-norm, max
- no parameters to be learned

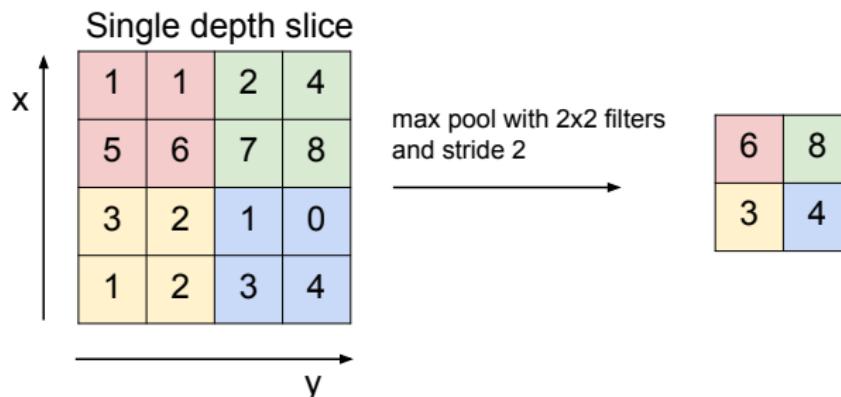
# Pooling

Similar to a filter, except

- depth is always 1
- different operations: average, L2-norm, max
- no parameters to be learned

**Max pooling** with  $2 \times 2$  filter and stride 2 is very common

## MAX POOLING



# Putting everything together

## Typical architecture for CNNs:

Input → [[Conv → ReLU]\*N → Pool?] \*M → [FC → ReLU]\*Q → FC

# Putting everything together

## Typical architecture for CNNs:

Input → [[Conv → ReLU]\*N → Pool?] \*M → [FC → ReLU]\*Q → FC

Common choices:  $N \leq 5$ ,  $Q \leq 2$ ,  $M$  is large

# Putting everything together

## Typical architecture for CNNs:

Input → [[Conv → ReLU]\*N → Pool?] \*M → [FC → ReLU]\*Q → FC

Common choices:  $N \leq 5$ ,  $Q \leq 2$ ,  $M$  is large

**Well-known CNNs:** LeNet, AlexNet, ZF Net, GoogLeNet, VGGNet, etc.

All achieve excellent performance on image classification tasks.

# How to train a CNN?

*How do we learn the filters/weights?*

# How to train a CNN?

*How do we learn the filters/weights?*

Essentially the same as FC NNs: apply **SGD/backpropagation**