

CS 446/ECE 449: Machine Learning

Shenlong Wang

University of Illinois at Urbana-Champaign, 2024

Convolutional Neural Networks (CNNs)

Goals of this lecture

- Convolutional Neural Networks

Goals of this lecture

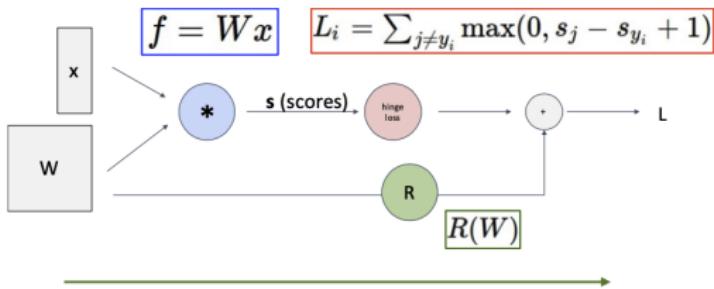
- Convolutional Neural Networks

Reading material

- I. Goodfellow et al.; Deep Learning; Chapters 6-9

Recall: Backpropagation

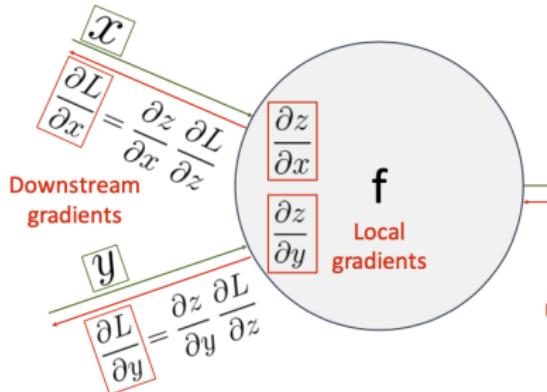
Represent complex expressions
as computational graphs



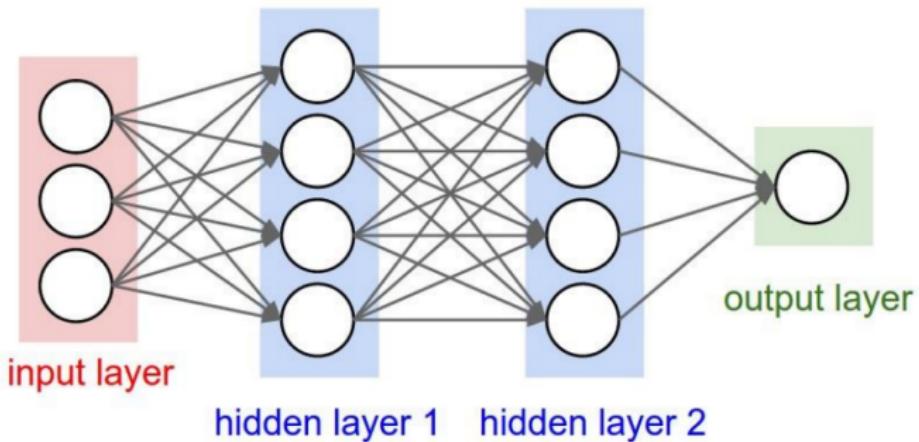
Forward pass computes outputs

Backward pass computes gradients

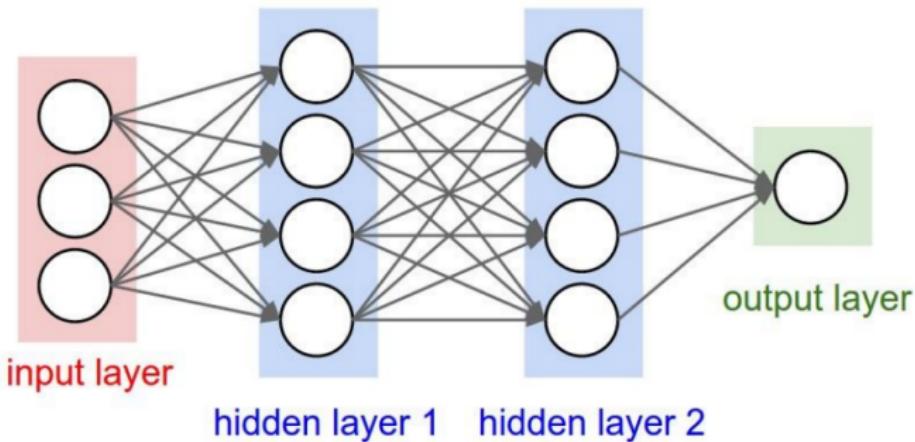
During the backward pass, each node in the graph receives **upstream gradients** and multiplies them by **local gradients** to compute **downstream gradients**.



Recall: Multi-layer Perceptron



Recall: Multi-layer Perceptron

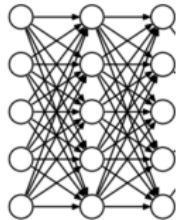


Problems:

- If x is a 500×500 image, z_1 is 1024-dimensional, what is the size of the W_1 ?
- If x is a 500×500 image, x is shifted by one pixel, do you expect the z_1 remain similar?

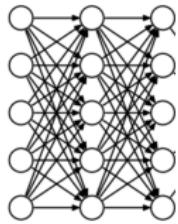
Components of a Fully-Connected Network

Fully-connected layers

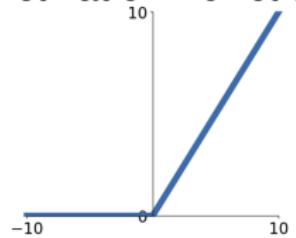


Components of a Fully-Connected Network

Fully-connected layers

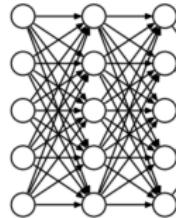


Activation Function

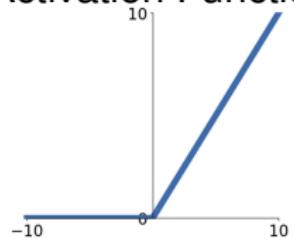


Components of a Convolutional Network

Fully-connected layers

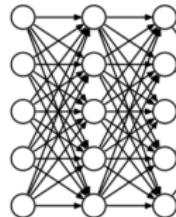


Activation Function

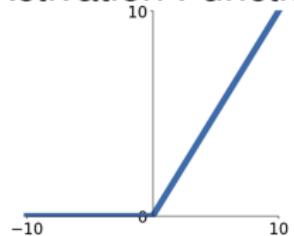


Components of a Convolutional Network

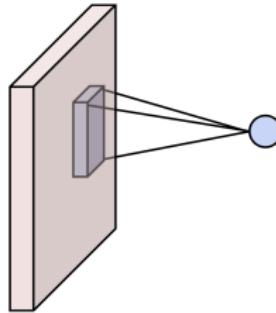
Fully-connected layers



Activation Function

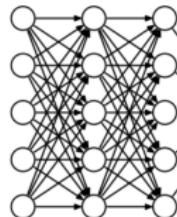


Convolution layers

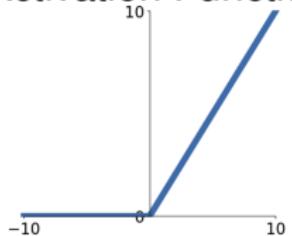


Components of a Convolutional Network

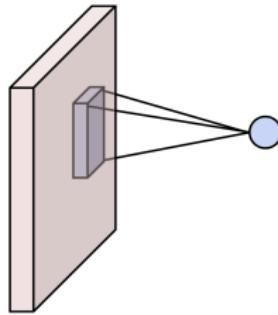
Fully-connected layers



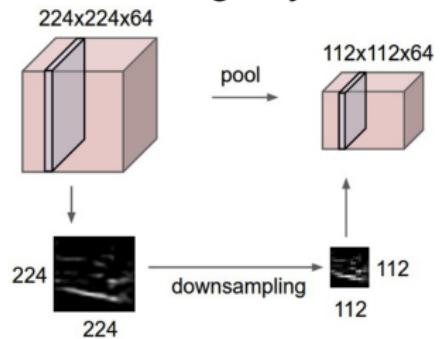
Activation Function



Convolution layers

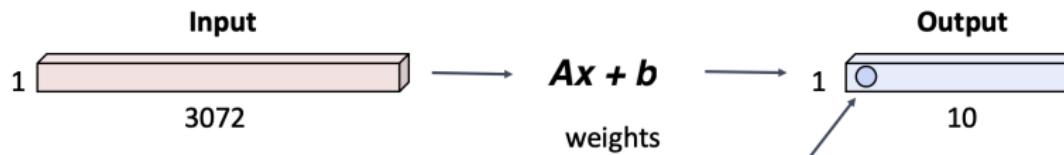


Pooling Layers



Fully-Connected Layer

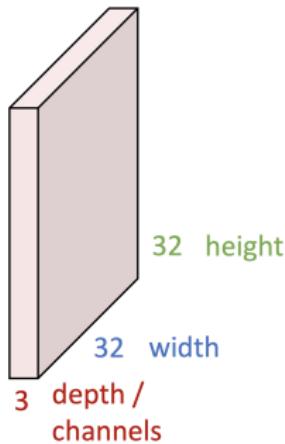
32x32x3 image -> stretch to 3072 x 1



1 number:
the result of taking a dot product between a row of W and the input (a 3072-dimensional dot product)

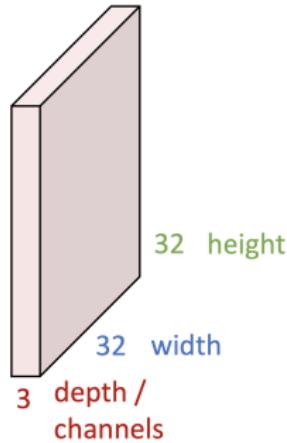
Convolution Layer

3x32x32 image: preserve spatial structure



Convolution Layer

$3 \times 32 \times 32$ image



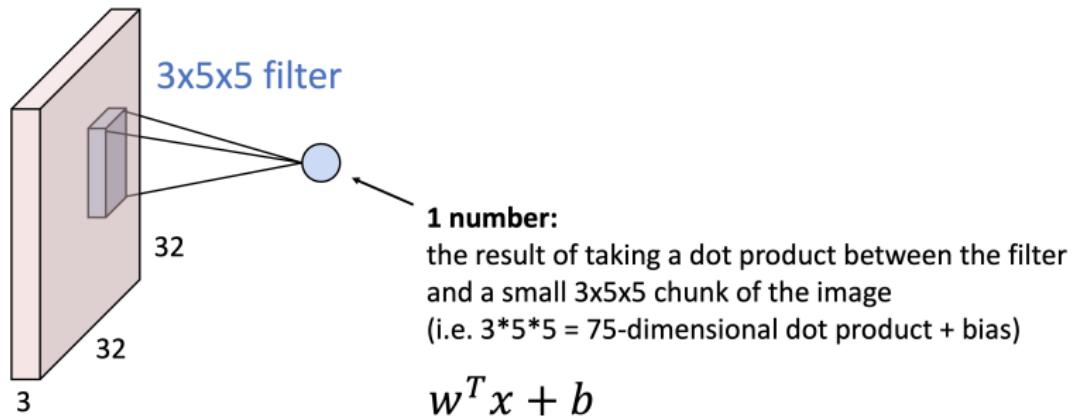
$3 \times 5 \times 5$ filter



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

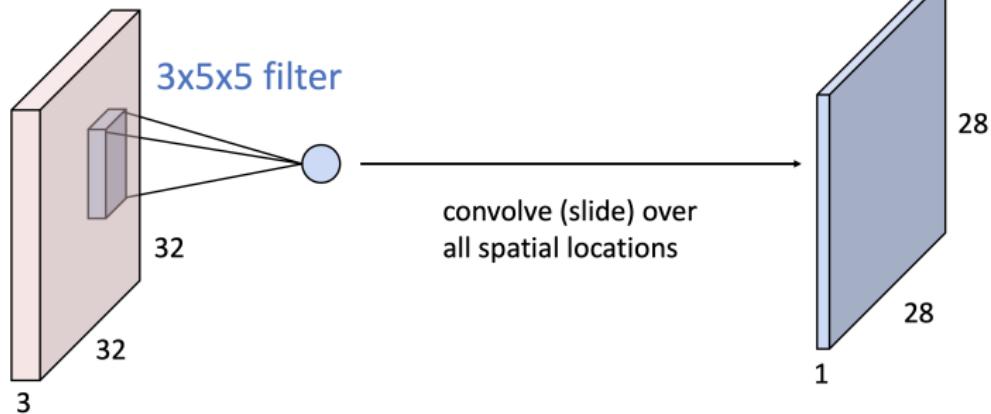
Convolution Layer

3x32x32 image



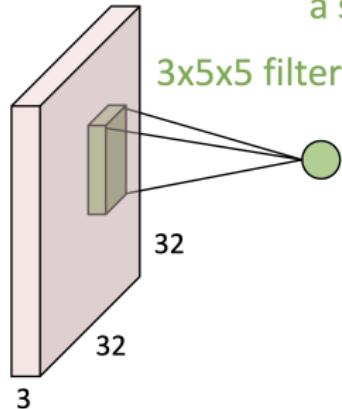
Convolution Layer

3x32x32 image



Convolution Layer

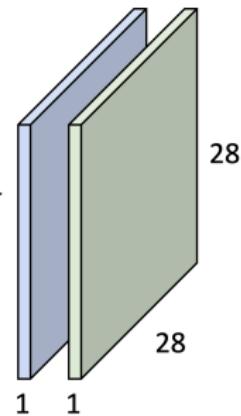
3x32x32 image



Consider repeating with
a second (green) filter:

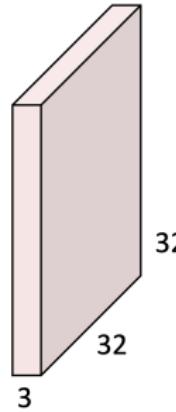
convolve (slide) over
all spatial locations

two 1x28x28
activation map



Convolution Layer

3x32x32 image



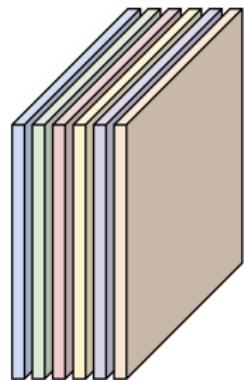
Consider 6 filters,
each 3x5x5



6x3x5x5
filters



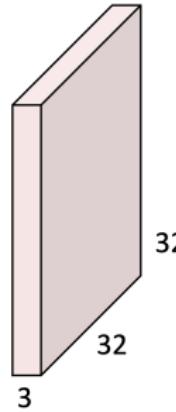
6 activation maps,
each 1x28x28



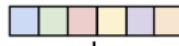
Stack activations to get a
6x28x28 output image!

Convolution Layer

3x32x32 image



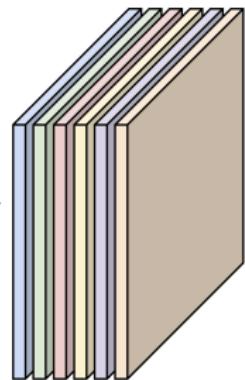
Also 6-dim bias vector:



6x3x5x5
filters

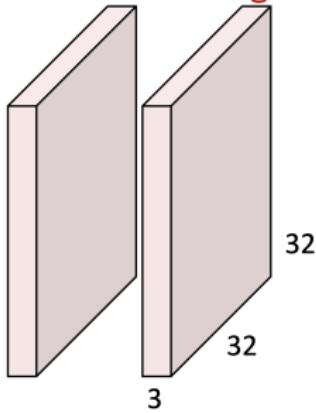


6 activation maps,
each 1x28x28

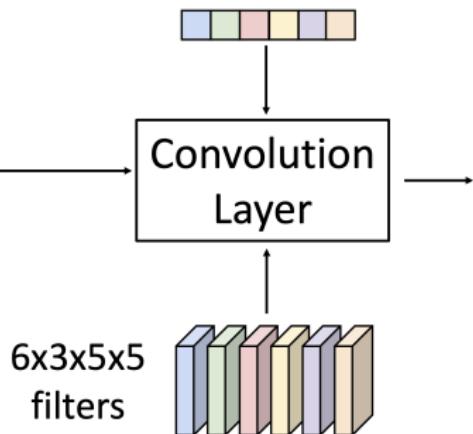


Stack activations to get a
6x28x28 output image!

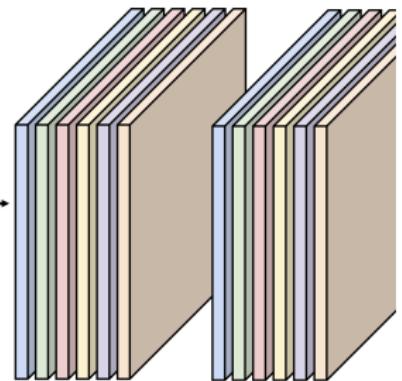
Convolution Layer
 $2 \times 3 \times 32 \times 32$
Batch of images



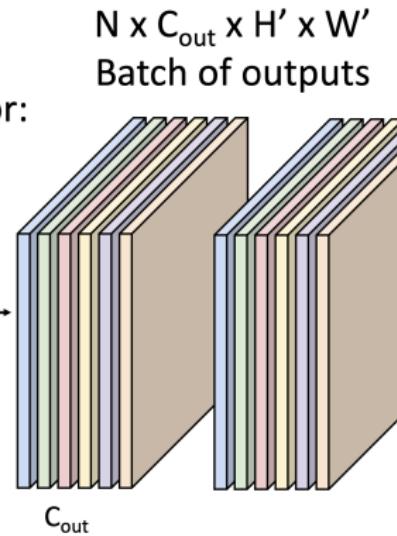
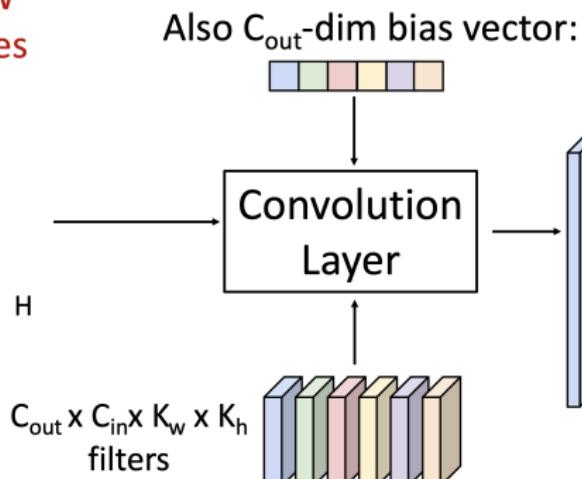
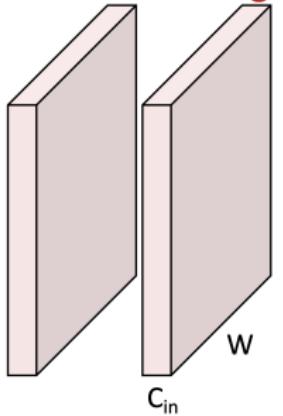
Also 6-dim bias vector:



$2 \times 6 \times 28 \times 28$
Batch of outputs

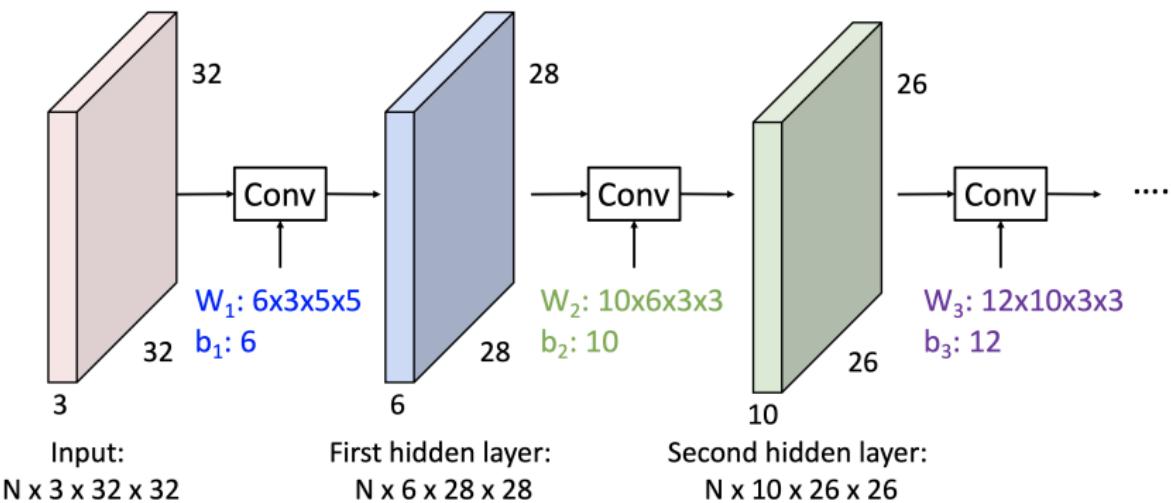


Convolution Layer
 $N \times C_{in} \times H \times W$
Batch of images



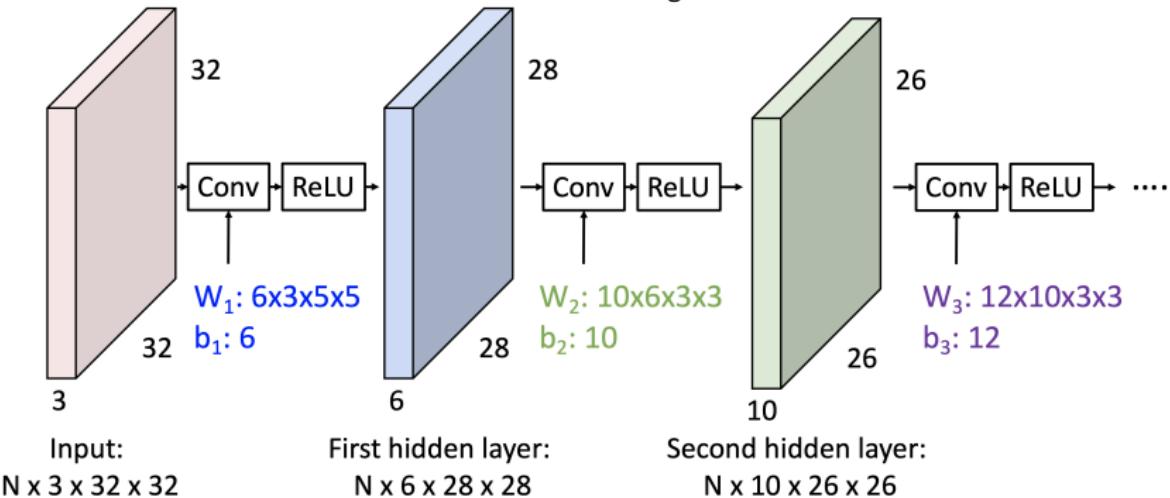
Stacking Convolutions

Q: What happens if we stack two convolution layers?

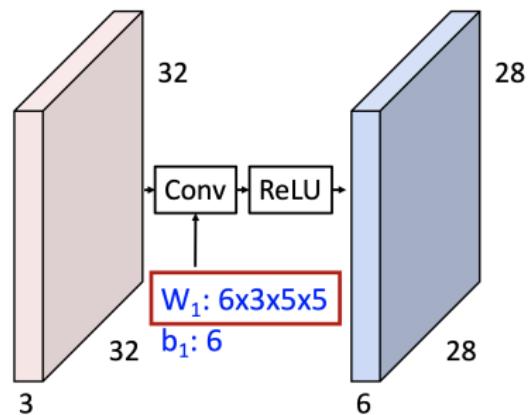


Stacking Convolutions

Q: What happens if we stack (Recall $y=W_2W_1x$ is two convolution layers? a linear classifier)
A: We get another convolution!



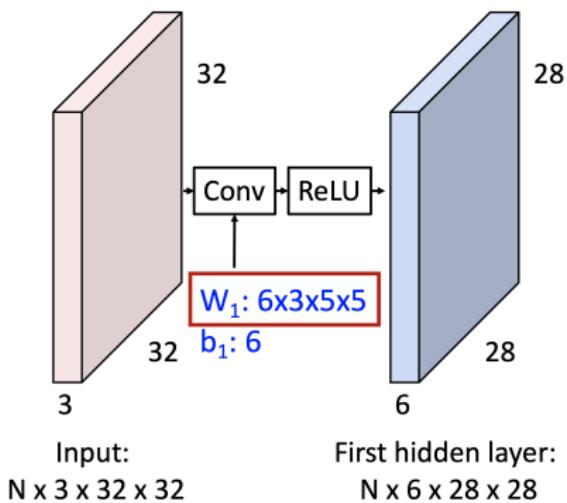
What do convolutional filters learn?



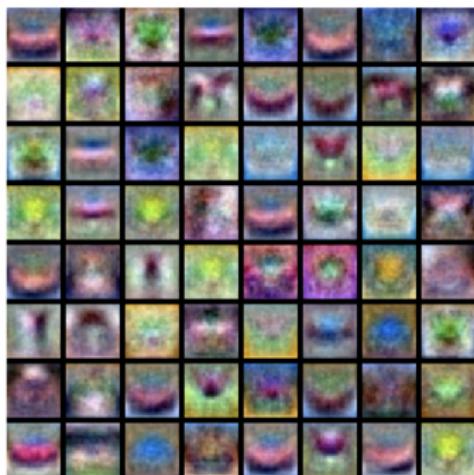
Linear classifier: One template per class



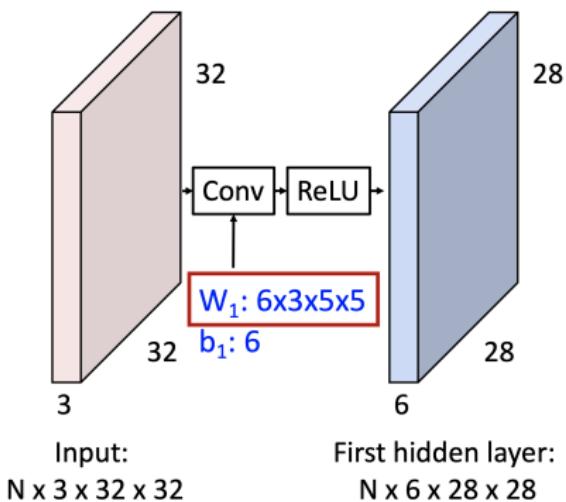
What do convolutional filters learn?



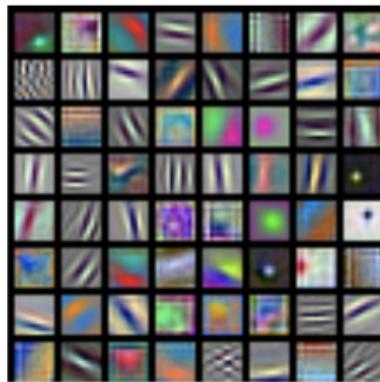
MLP: Bank of whole-image templates



What do convolutional filters learn?

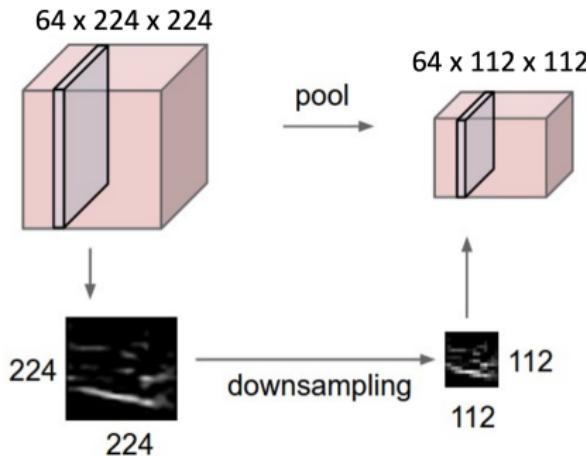


First-layer conv filters: local image templates
(Often learns oriented edges, opposing colors)



AlexNet: 64 filters, each $3 \times 11 \times 11$

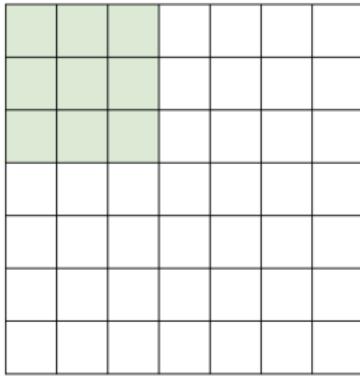
Pooling Layers: Another way to downsample



Hyperparameters:
Kernel Size
Stride
Pooling function

Stride

7

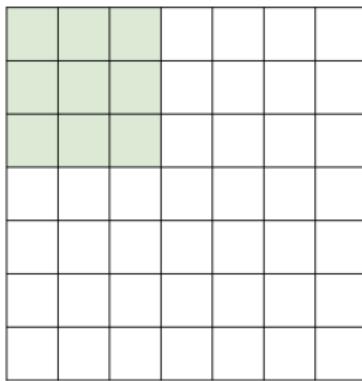


7x7 input (spatially)
assume 3x3 filter

7

Stride

7

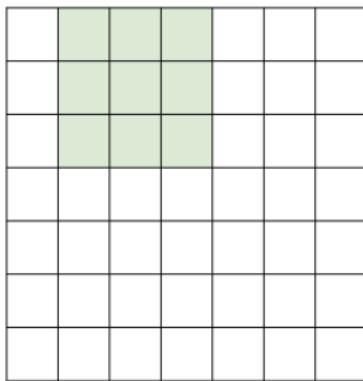


7x7 input (spatially)
assume 3x3 filter

7

Stride

7

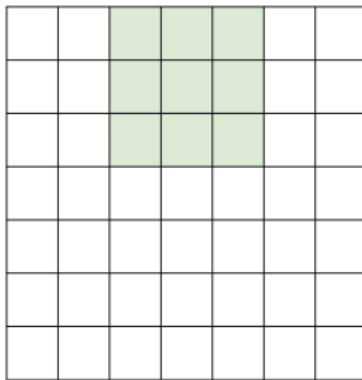


7x7 input (spatially)
assume 3x3 filter

7

Stride

7

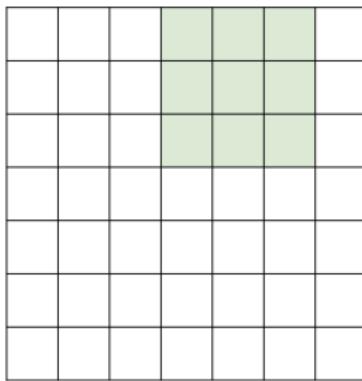


7x7 input (spatially)
assume 3x3 filter

7

Stride

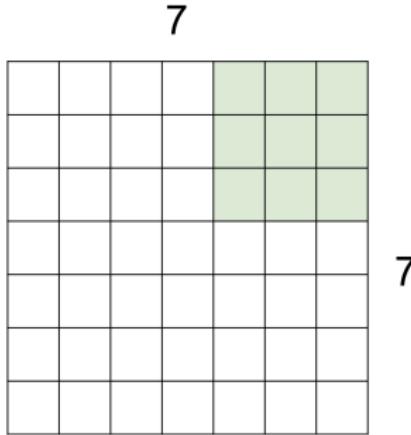
7



7x7 input (spatially)
assume 3x3 filter

7

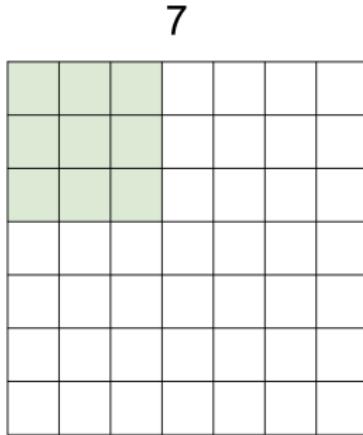
Stride



7x7 input (spatially)
assume 3x3 filter

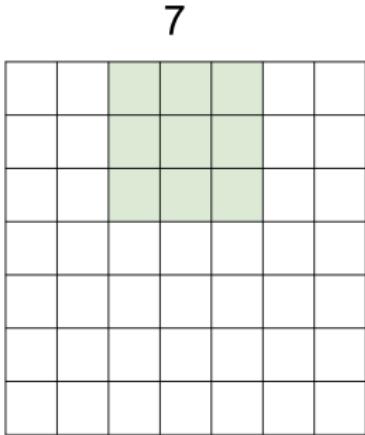
=> 5x5 output

Stride



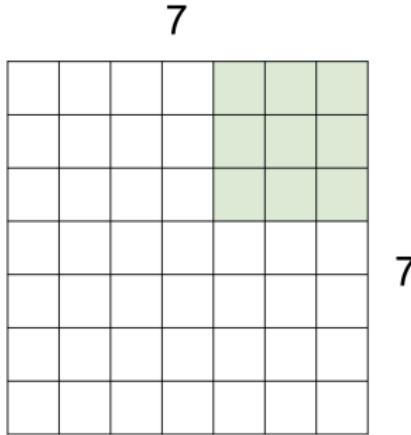
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

Stride



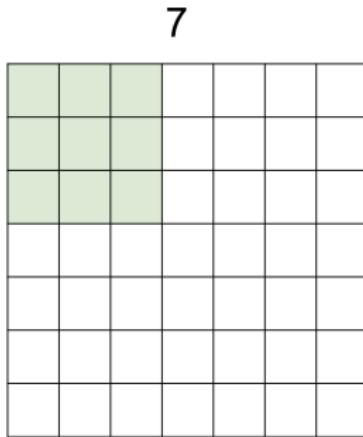
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

Stride



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

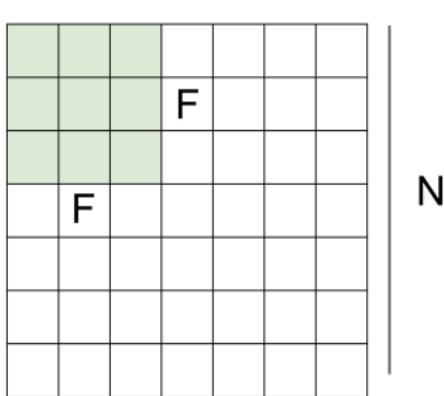
Stride



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

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

Stride



Output size:

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

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

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

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

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

Padding

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:)

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

Padding

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

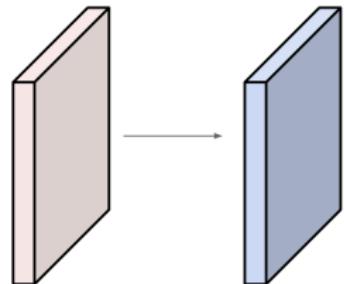
$F = 5 \Rightarrow$ zero pad with 2

$F = 7 \Rightarrow$ zero pad with 3

Padding

Examples time:

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



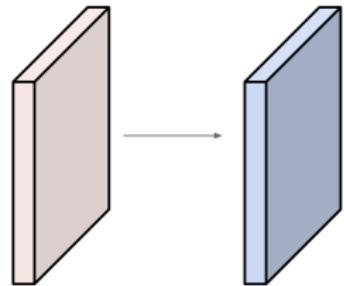
Output volume size:
 $(32+2*2-5)/1+1 = 32$ spatially, so
32x32x10

Padding

Examples time:

Input volume: **32x32x3**

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



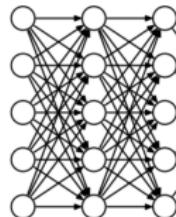
Number of parameters in this layer?

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

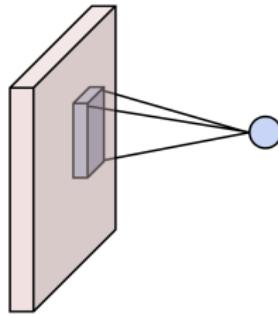
$$\Rightarrow 76 * 10 = 760$$

Components of a Convolutional Network

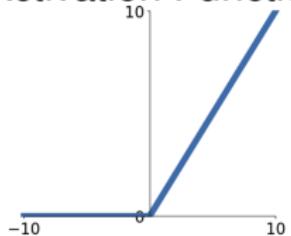
Fully-connected layers



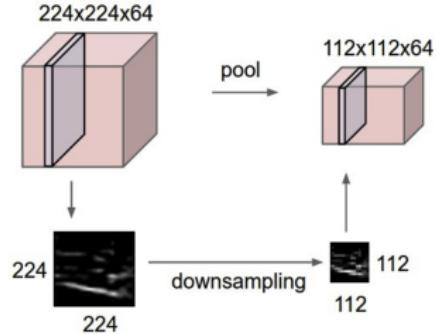
Convolution layers



Activation Function



Pooling Layers

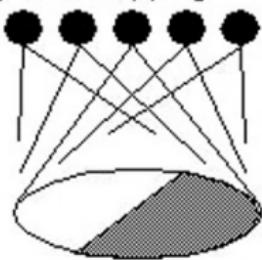


Backstory

- Biological inspiration: D. Hubel and T. Wiesel (1959, 1962, Nobel Prize 1981)
- Visual cortex consists of a hierarchy of simple, complex, and hyper-complex cells

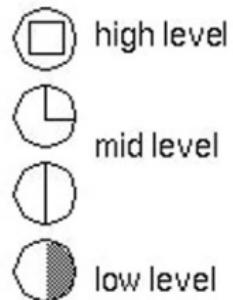
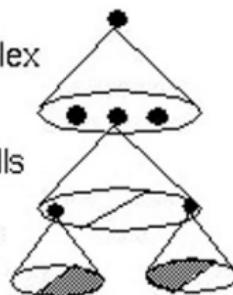
Hubel & Weisel

topographical mapping



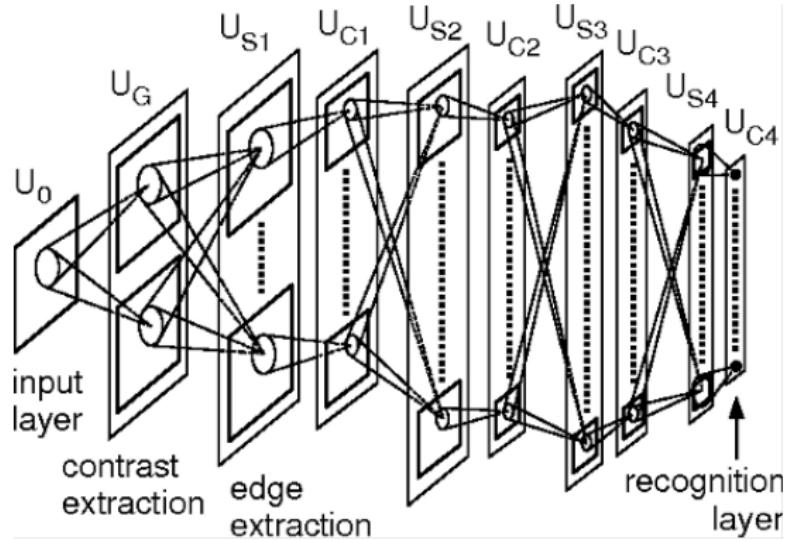
featural hierarchy

hyper-complex
cells
complex cells
simple cells



Backstory

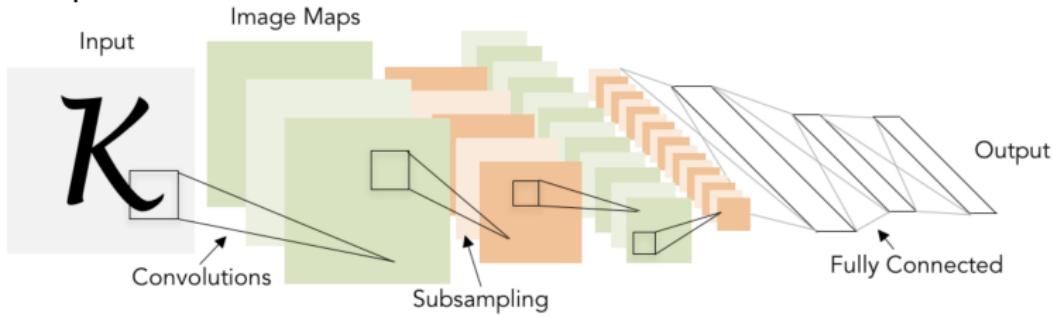
- Neocognitron (Kunihiko Fukushima in 1979)
- Inspired by Hubel and Wiesel, and inspired later CNNs



Convolutional Networks

Classic architecture: [Conv, ReLU, Pool] x N, flatten, [FC, ReLU] x N, FC

Example: LeNet-5



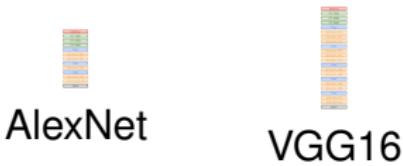
Lecun et al, "Gradient-based learning applied to document recognition", 1998

CNN Architectures



AlexNet

CNN Architectures



AlexNet

VGG16

CNN Architectures



AlexNet



VGG16



VGG19

CNN Architectures



AlexNet



VGG16



VGG19



GoogLeNet

CNN Architectures



AlexNet



VGG16



VGG19



GoogLeNet



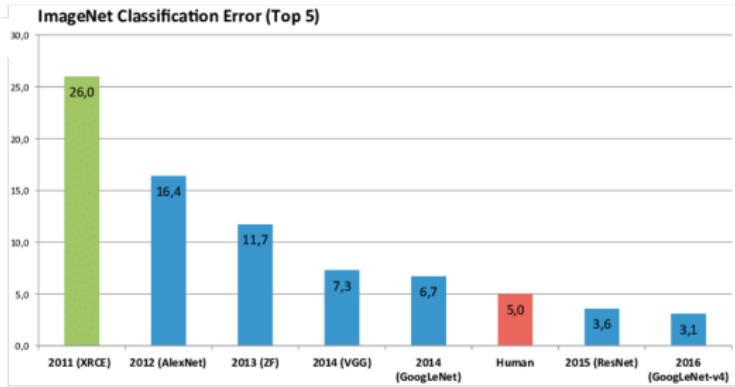
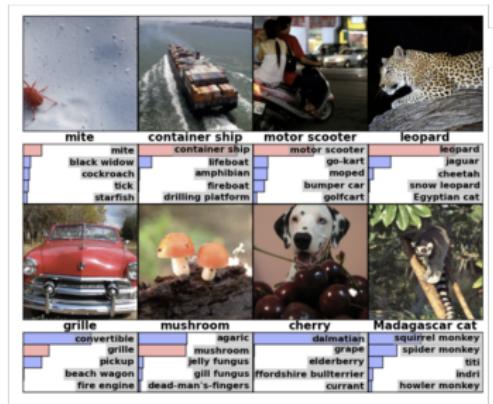
ResNet

CNN Architectures

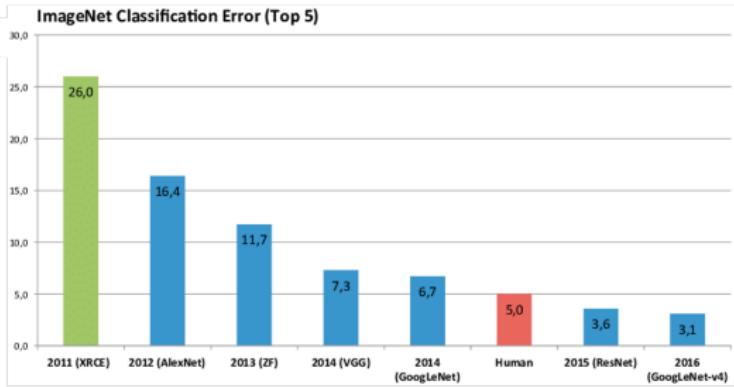
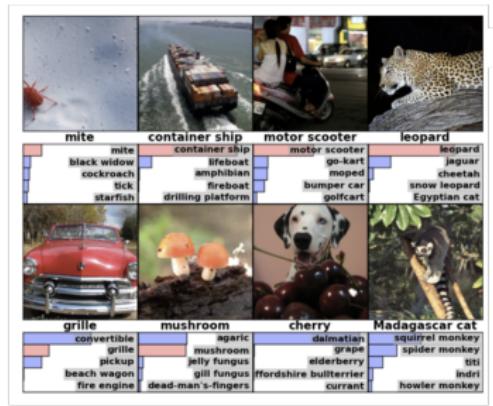


Mo

ImageNet



ImageNet





Which Architectures should I use?

- For most problems you should use an off-the-shelf architecture (with proven record after many “graduate student descent”); don’t try to design your own.

Which Architectures should I use?

- For most problems you should use an off-the-shelf architecture (with proven record after many “graduate student descent”); don’t try to design your own.
- If you just care about accuracy, **ResNet-50** or **ResNet-101** are great choices

Which Architectures should I use?

- For most problems you should use an off-the-shelf architecture (with proven record after many “graduate student descent”); don’t try to design your own.
- If you just care about accuracy, **ResNet-50** or **ResNet-101** are great choices
- If you want an efficient network (real-time, run on mobile, etc), try **MobileNets** and **ShuffleNets**

Which Architectures should I use?

- For most problems you should use an off-the-shelf architecture (with proven record after many “graduate student descent”); don’t try to design your own.
- If you just care about accuracy, **ResNet-50** or **ResNet-101** are great choices
- If you want an efficient network (real-time, run on mobile, etc), try **MobileNets** and **ShuffleNets**
- Start from pretrained models.
https://pytorch.org/serve/model_zoo.html

Demo time: [https://huggingface.co/spaces/Xenova/
object-detection-webcam_testing](https://huggingface.co/spaces/Xenova/object-detection-webcam_testing)

Important topics of this lecture

Important topics of this lecture

- Getting to know CNNs and its variants

Important topics of this lecture

- Getting to know CNNs and its variants
- Getting to know their use

Important topics of this lecture

- Getting to know CNNs and its variants
- Getting to know their use
- Contrasting CNNs to MLPs

Important topics of this lecture

- Getting to know CNNs and its variants
- Getting to know their use
- Contrasting CNNs to MLPs

Important topics of this lecture

- Getting to know CNNs and its variants
- Getting to know their use
- Contrasting CNNs to MLPs

Next up:

Recurrent neural networks