Reminders / Announcements

- Project Proposal has been graded
 - o Goal: to give you detailed feedback. Grading is "light" (Min: 90%, Max: 99%)
 - <u>Compute:</u> Google Colab Pro is free for 1 year: https://blog.google/outreach-initiatives/education/colab-higher-education/
 - Oct 22: PyTorch tutorial in class (taught by Yu Liu)
- Midterm Exam is on 10/20
 - o In class; closed-book; 1 hour; everything up to and including 10/15 lecture
 - More details in the next class.
- Homework 1 is being graded (ETA ~1 week)
- Homework 2 will be released this weekend (and due ~ Nov 3)

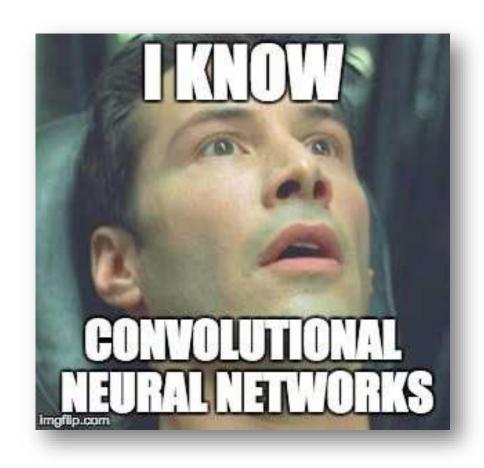
Your Project Topics are Cool!

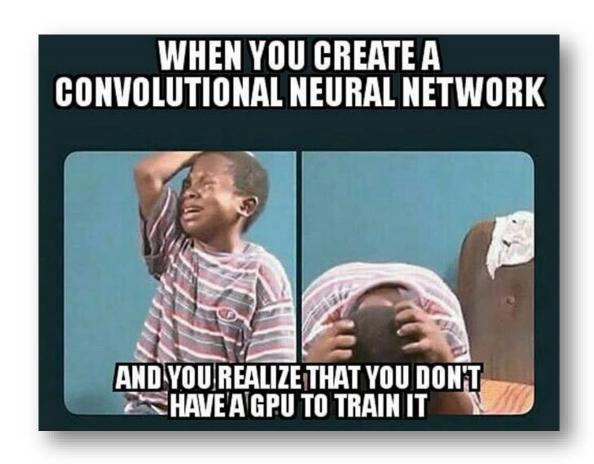


- Self-supervised learning for medical imaging
- Combined stereo-mono depth estimation
- Interactive Video Generation
- Interactive Image Segmentation
- Anomaly Detection in Astrophysics
- Enhancing Robustness of Al Watermarking
- Fish Growth Monitoring via Vision

- Calorie/Nutrition Estimation from Images
- Vision-based Parking Monitoring
- Photogrammetry
- Detection of "Al-Generated" Images
- Embodied Visual Category Discovery
- Adversarial Defense in Event-Based Vision
- Paraboloid CNNs

Lecture 10 Convolutional Neural Networks







| 10 BREAKTHROUGH | TECHNOLOGIES 2013

CONNECT

Introduction The 10 Technologies Past Years

Deep Learning

With massive amounts of computational power, machines can now recognize objects and translate speech in real time. Artificial intelligence is finally getting smart.

Temporary Social Media

Messages that quickly self-destruct could enhance the privacy of online communications and make people freer to be spontaneous.

Prenatal DNA Sequencing

Reading the DNA of fetuses will be the next frontier of the genomic revolution. But do you really want to know about the genetic problems or musical aptitude of your unborn child?

Additive Manufacturing

Skeptical about 3-D printing? GE, the world's largest manufacturer, is on the verge of using the technology to make jet parts.

Baxter: The Blue-Collar Robot

Rodney Brooks's newest creation is easy to interact with, but the complex innovations behind the robot show just how hard it is to get along with people.

Memory Implants

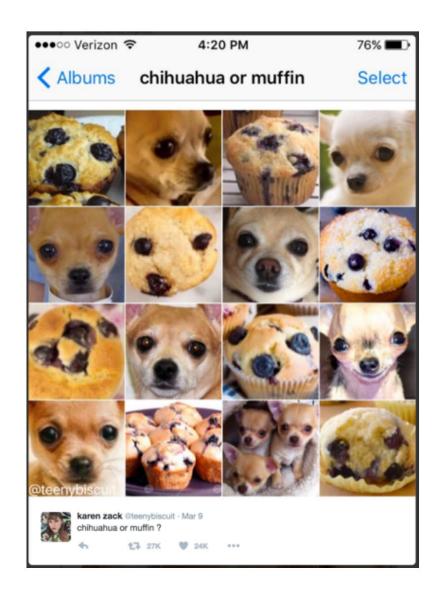
Smart Watches

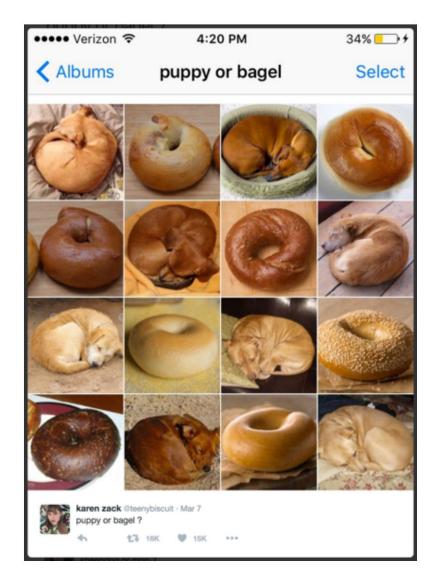
Ultra-Efficient Solar

Big Data from

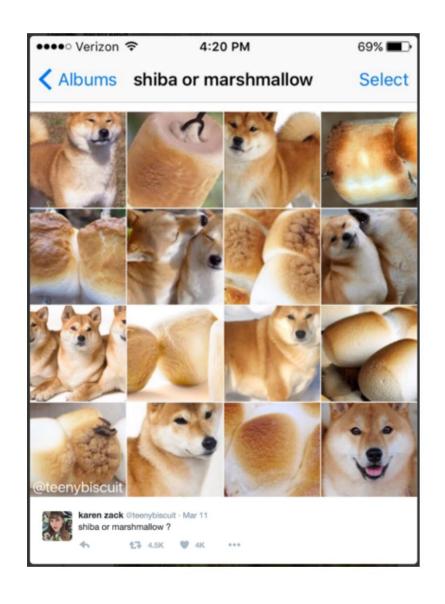
Supergrids

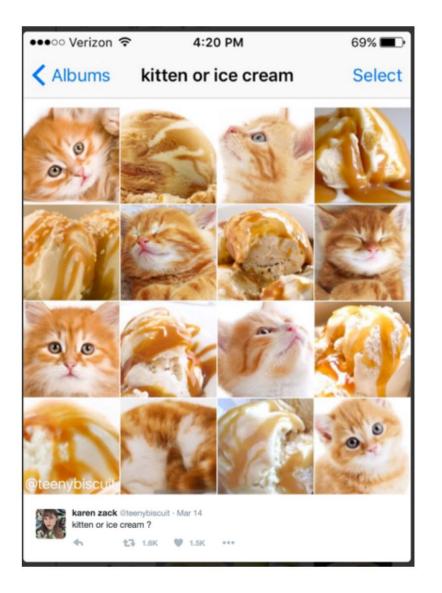
(Unrelated) Dog vs Food





(Unrelated) Dog vs Food





CNNs in 2012: "SuperVision" (aka "AlexNet")

"AlexNet" — Won the ILSVRC2012 Challenge

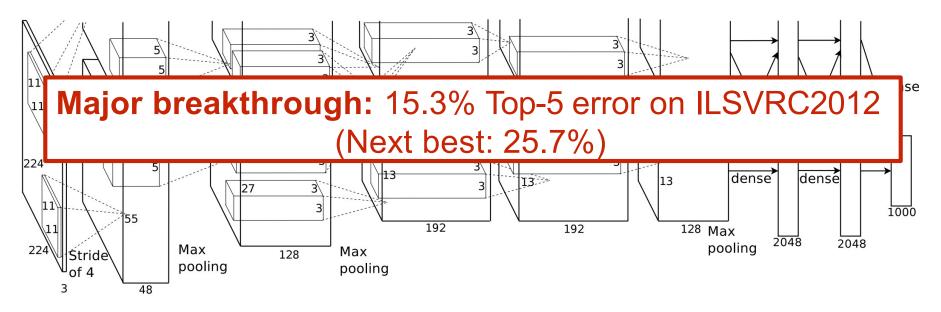
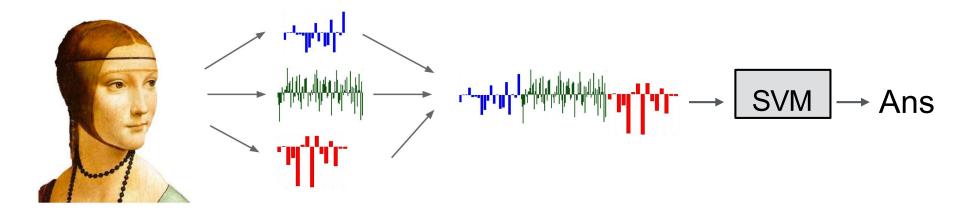


Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–1000.

[Krizhevsky, Sutskever, Hinton. NIPS 2012]

Recap: Before Deep Learning



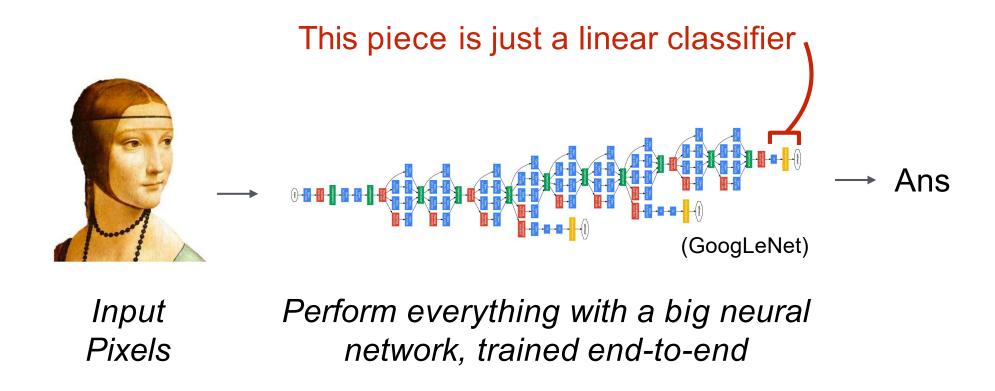
Input Pixels Extract Features

Concatenate into a vector **x**

Linear Classifier

Figure: Karpathy 2016

The last layer of (most) CNNs are linear classifiers

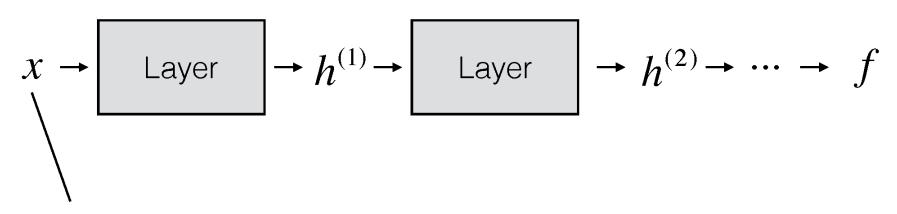


Key: perform enough processing so that by the time you get to the end of the network, the classes are linearly separable

ConvNets

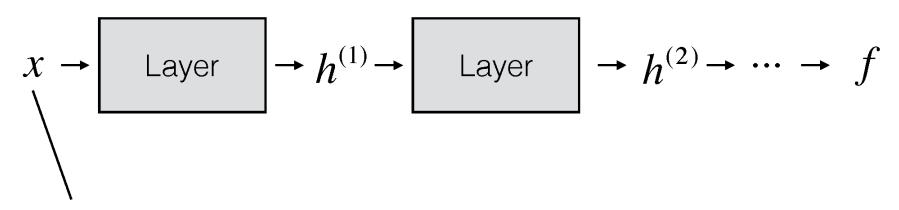
They're just neural networks with 3D activations and weight sharing

What shape should the activations have?



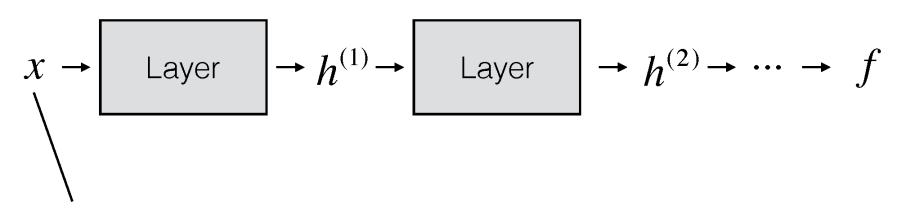
- The input is an image, which is 3D (RGB channel, height, width)

What shape should the activations have?



- The input is an image, which is 3D (RGB channel, height, width)
- We could flatten it to a 1D vector, but then we lose structure

What shape should the activations have?

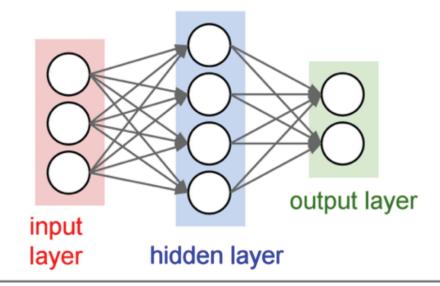


- The input is an image, which is 3D (RGB channel, height, width)
- We could flatten it to a 1D vector, but then we lose structure
- What about keeping everything in 3D?

ConvNets

They're just neural networks with 3D activations and weight sharing

before:



(1D vectors)

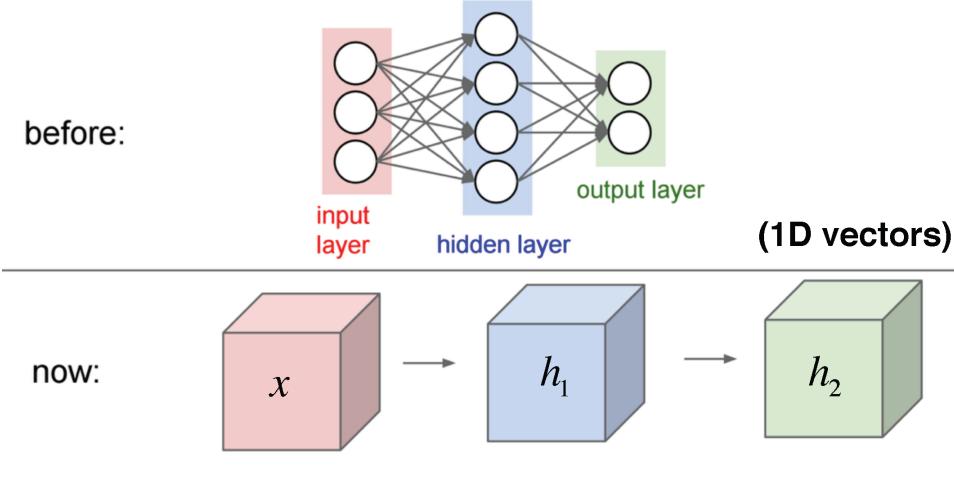
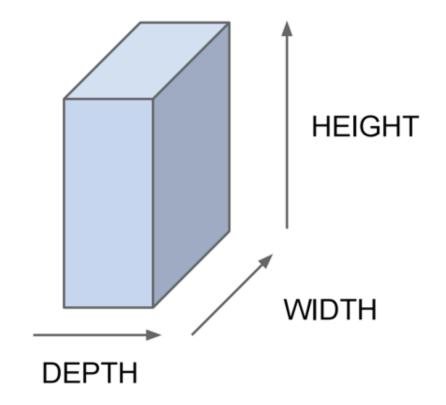


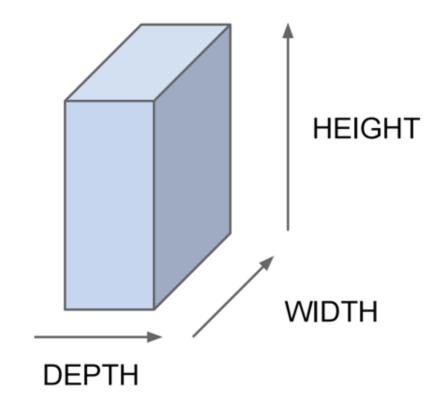
Figure: Andrej Karpathy

(3D arrays)

All Neural Net activations arranged in 3 dimensions:

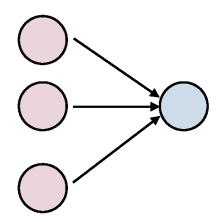


All Neural Net activations arranged in 3 dimensions:

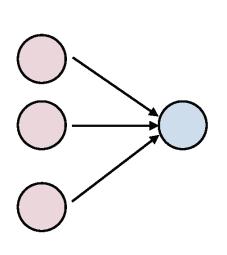


For example, a CIFAR-10 image is a 3x32x32 volume (3 depth — RGB channels, 32 height, 32 width)

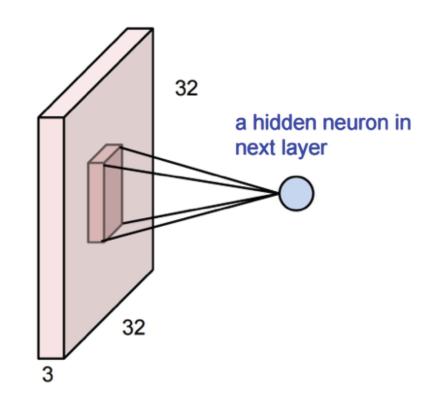
1D Activations:

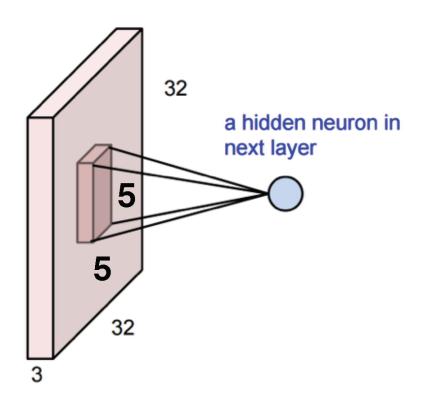


1D Activations:

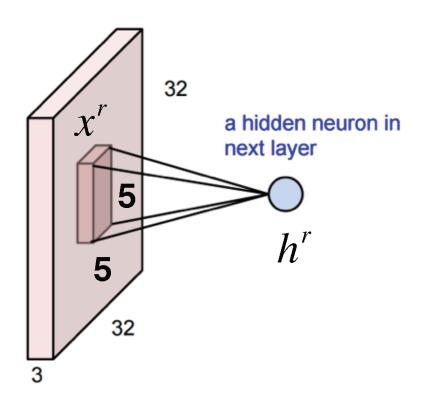


3D Activations:



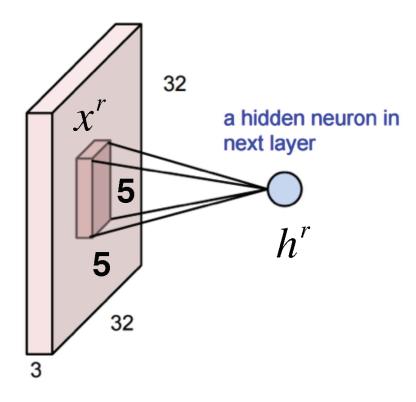


- The input is 3x32x32
- This neuron depends on a 3x5x5 chunk of the input
- The neuron also has a 3x5x5 set of weights and a bias (scalar)



Example: consider the region of the input " x^r "

With output neuron h^r

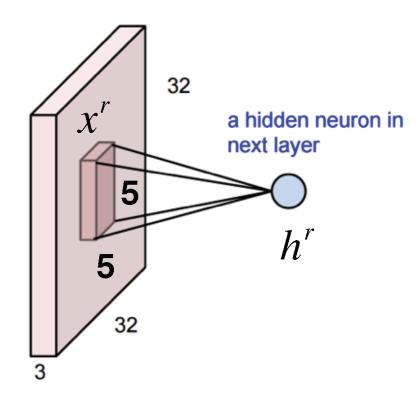


Example: consider the region of the input " x^{r} "

With output neuron h^r

Then the output is:

$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$



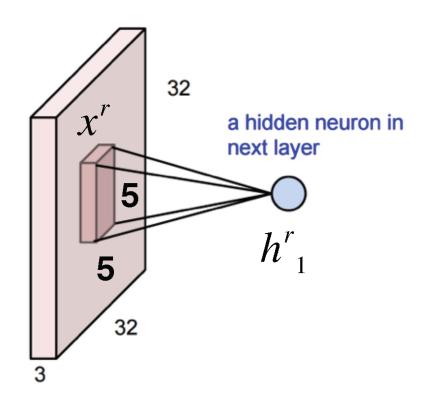
Example: consider the region of the input " x^{r} "

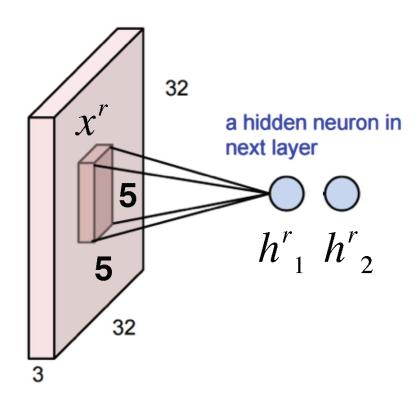
With output neuron h^r

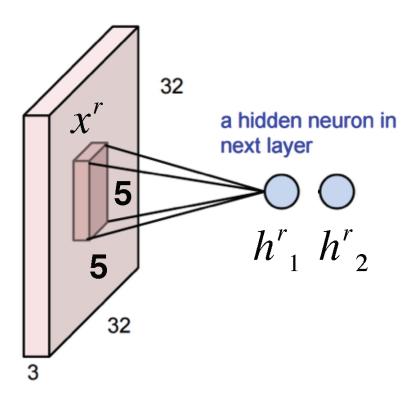
Then the output is:

$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$

Sum over 3 axes



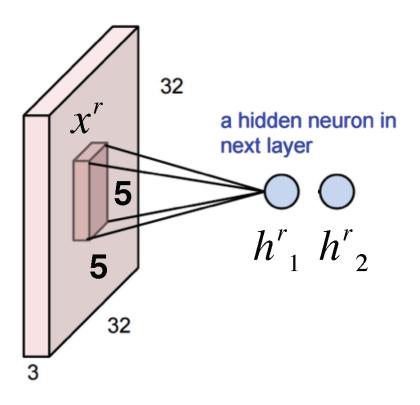




With 2 output neurons

$$h^{r}_{1} = \sum_{ijk} x^{r}_{ijk} W_{1ijk} + b_{1}$$

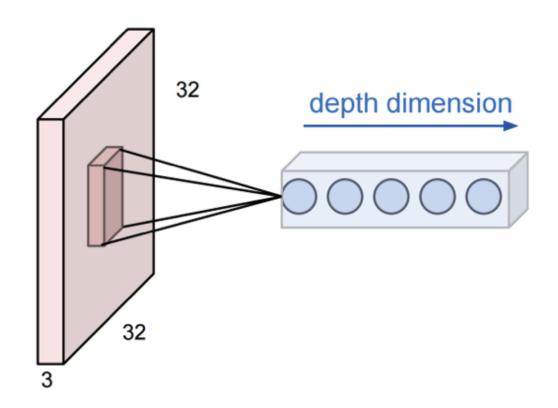
$$h^{r}_{2} = \sum_{ijk} x^{r}_{ijk} W_{2ijk} + b_{2}$$

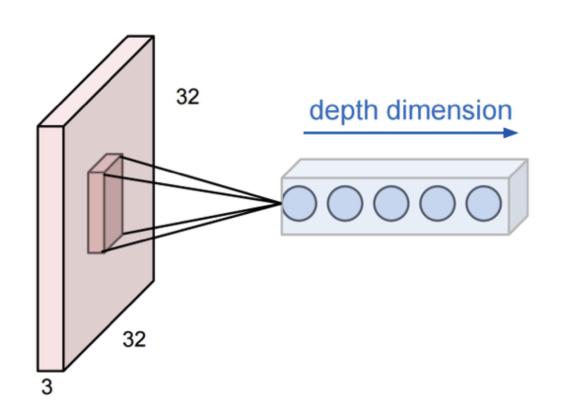


With 2 output neurons

$$h^r_{1} = \sum_{ijk} x^r_{ijk} W_{1ijk} + b_{1}$$

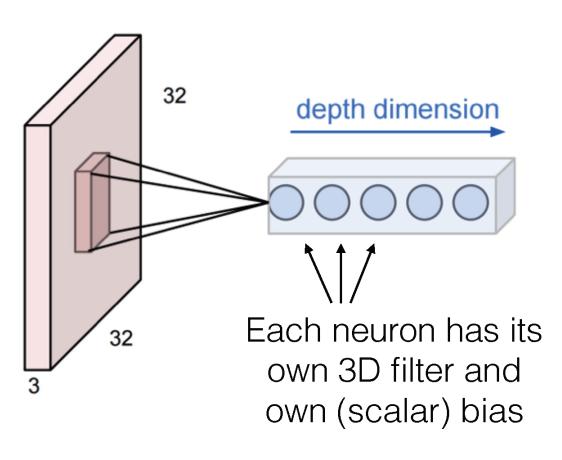
$$h^r_2 = \sum_{ijk} x^r_{ijk} W_{2ijk} + b_2$$





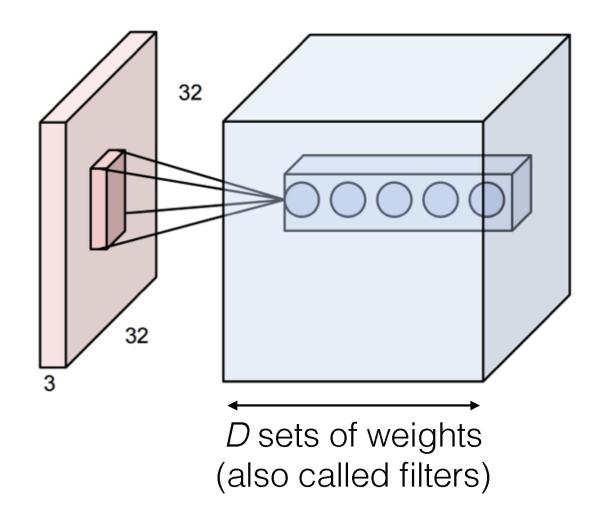
We can keep adding more outputs

These form a column in the output volume: [depth x 1 x 1]

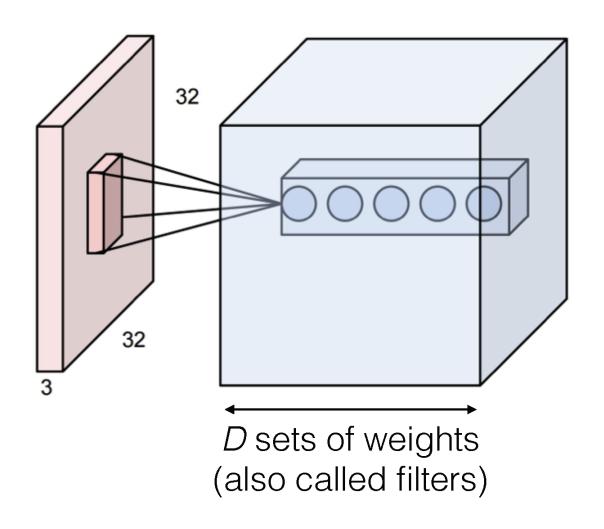


We can keep adding more outputs

These form a column in the output volume: [depth x 1 x 1]



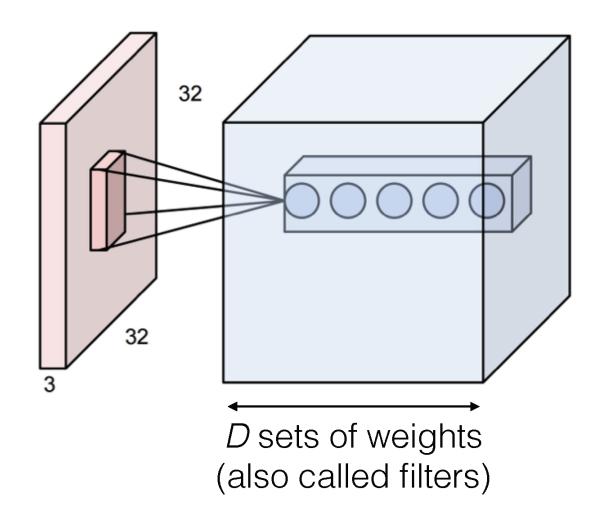
Now repeat this across the input

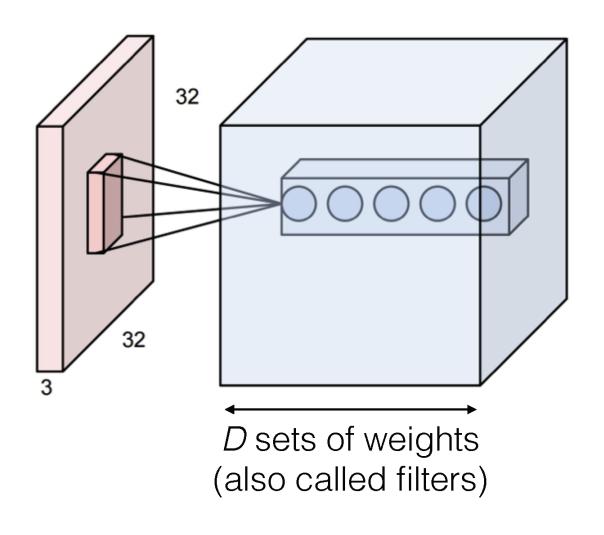


Now repeat this across the input

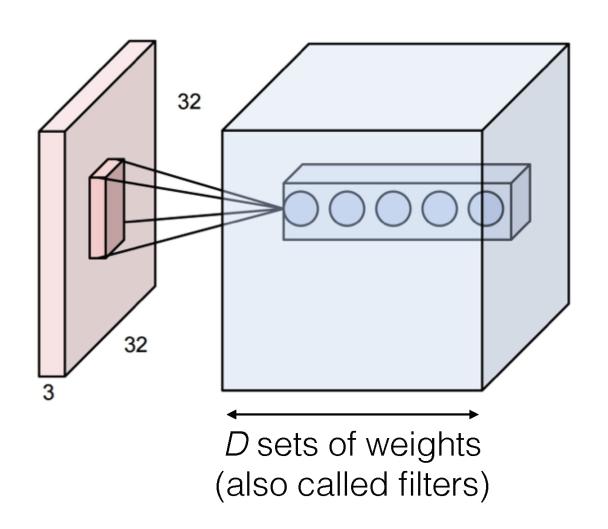
Weight sharing:

Each filter shares the same weights (but each depth index has its own set of weights)



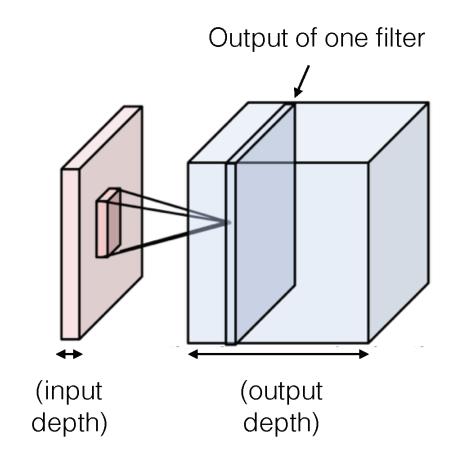


With weight sharing, this is called convolution



With weight sharing, this is called convolution

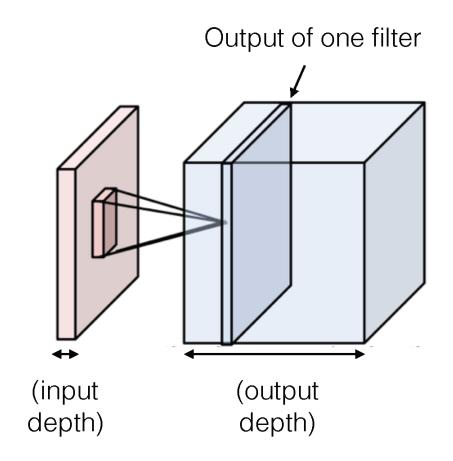
Without weight sharing, this is called a locally connected layer



One set of weights gives one slice in the output

To get a 3D output of depth *D*, use *D* different filters

In practice, ConvNets use many filters (~64 to 1024)



One set of weights gives one slice in the output

To get a 3D output of depth *D*, use *D* different filters

In practice, ConvNets use many filters (~64 to 1024)

All together, the weights are **4** dimensional: (output depth, input depth, kernel height, kernel width)

Fixed Filters



original

Vertical Sobel filter:

1	2	1		
0	0	0		
-1	-2	-1		

CONVOLUTION



vertical Sobel filter

Learned **Filters**



original

 W_{23}

CONVOLUTION

Learned **Features**

We can unravel the 3D cube and show each layer separately:



Figure: Andrej Karpathy

We can unravel the 3D cube and show each layer separately:

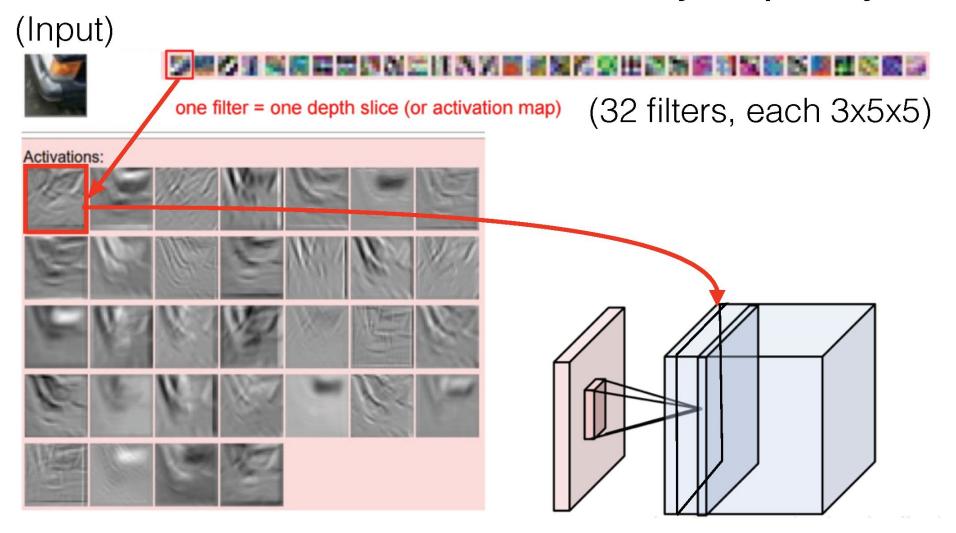


Figure: Andrej Karpathy

We can unravel the 3D cube and show each layer separately:

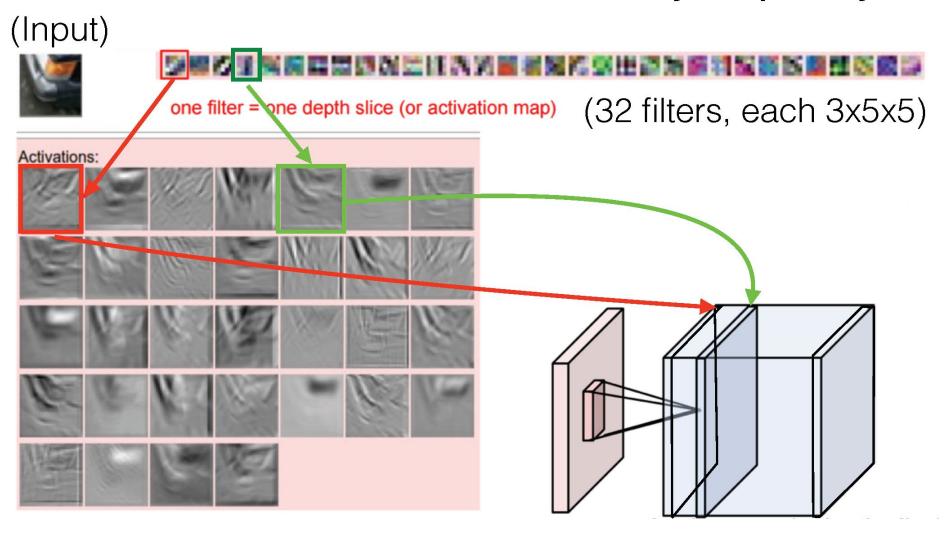
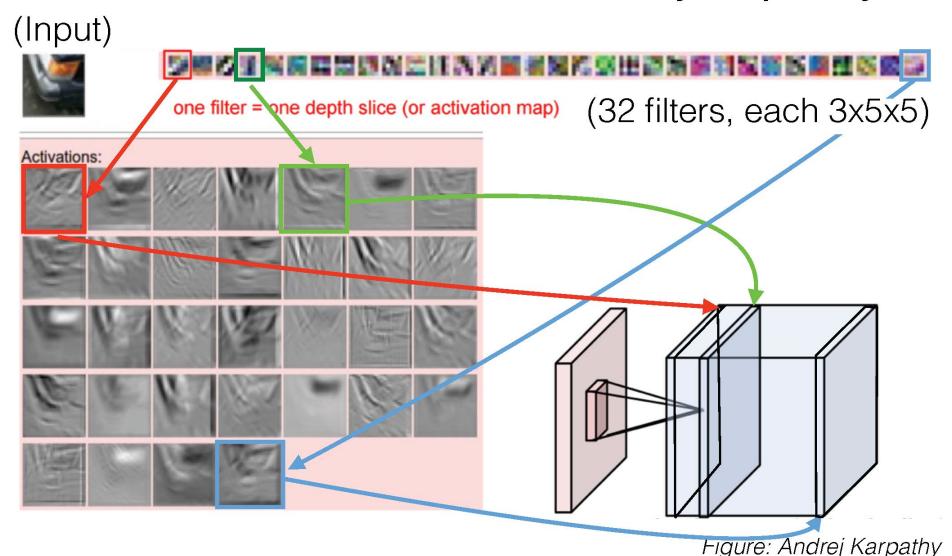
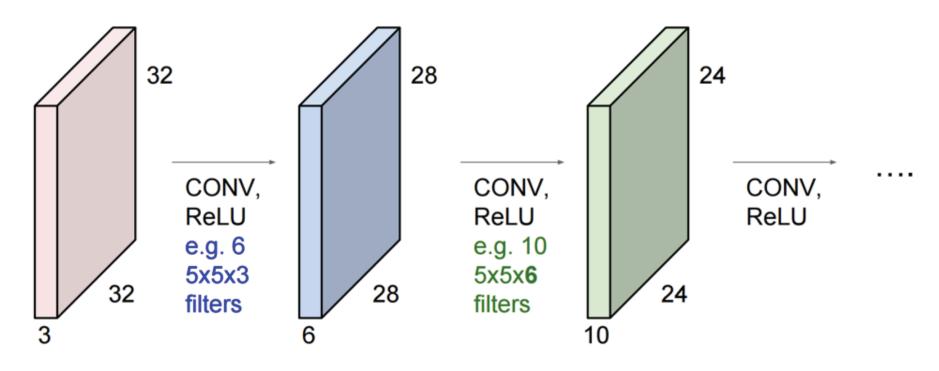


Figure: Andrej Karpathy

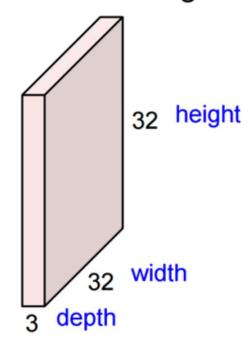
We can unravel the 3D cube and show each layer separately:



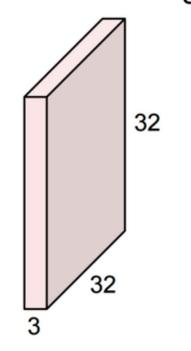
A **ConvNet** is a sequence of convolutional layers, interspersed with activation functions (and possibly other layer types)



32x32x3 image



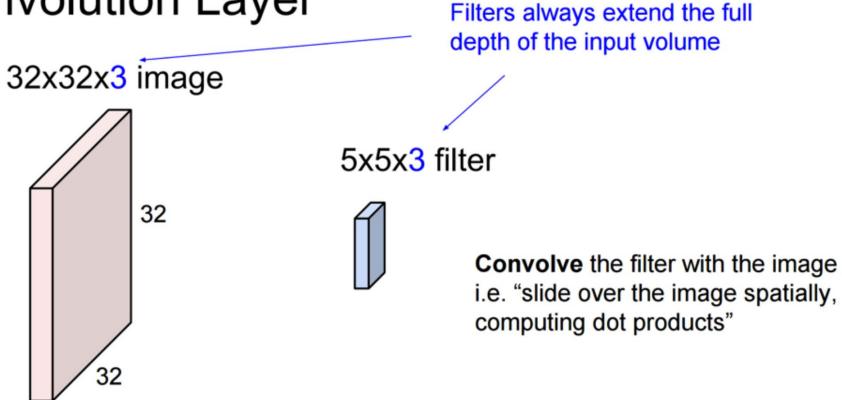
32x32x3 image

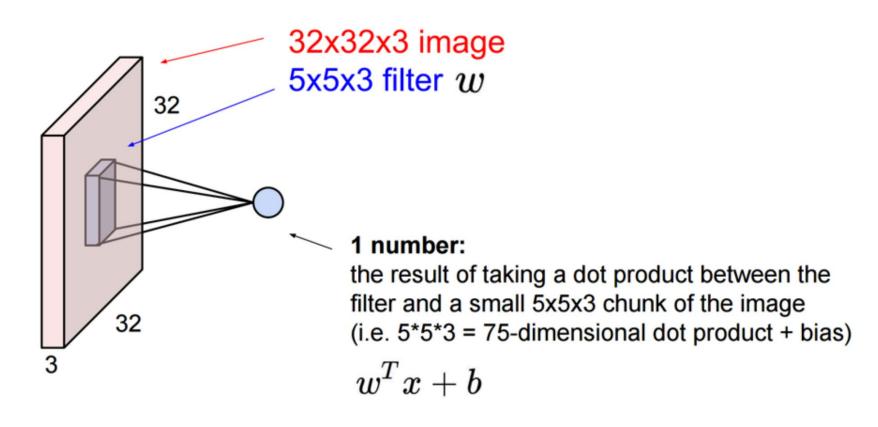


5x5x3 filter



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

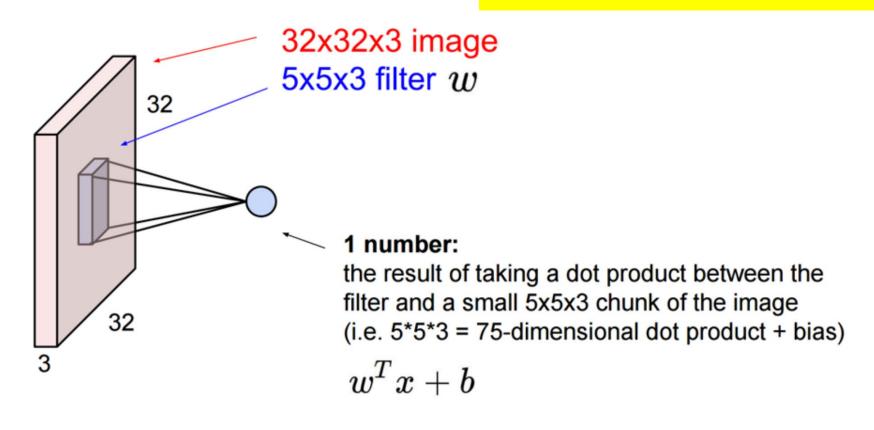


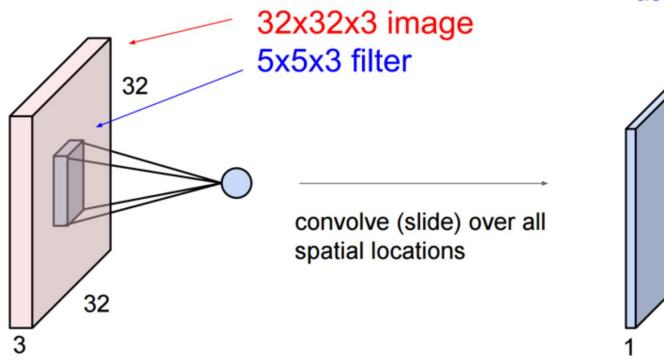


What will the output size be?

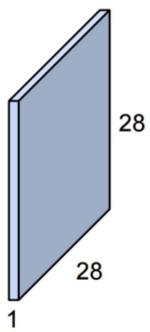
You will need to make some assumptions ...

Convolution Layer



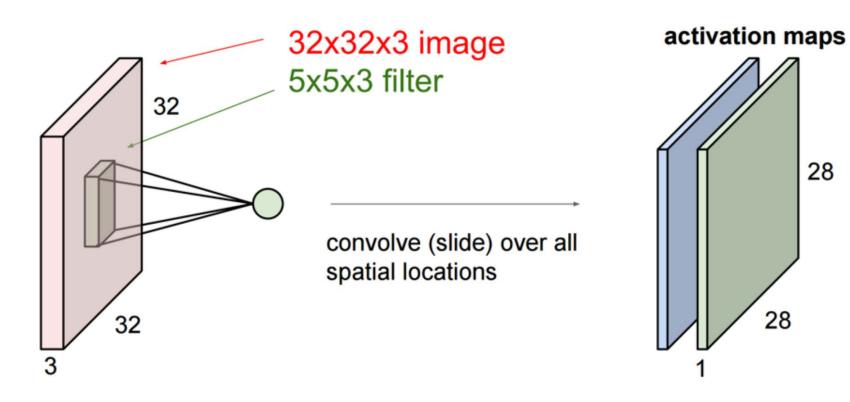


activation map

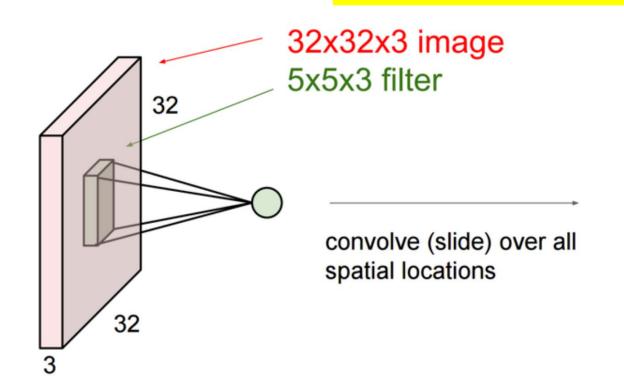


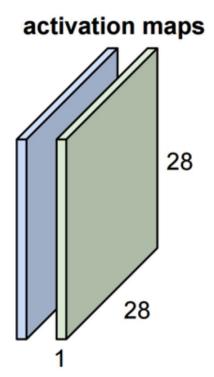
Consider a second filter ...

28

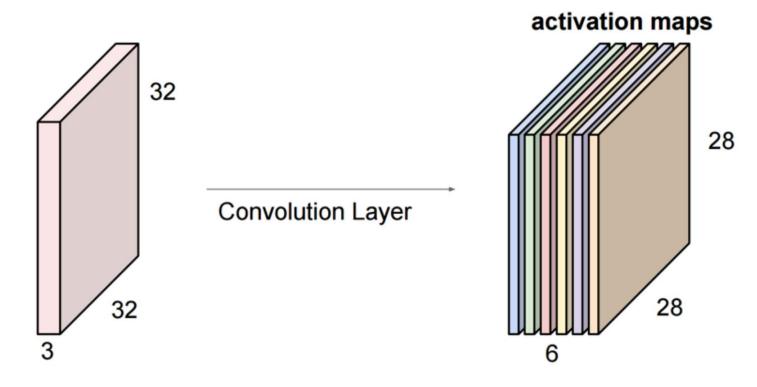


What will the output size be if we have 6 filters?

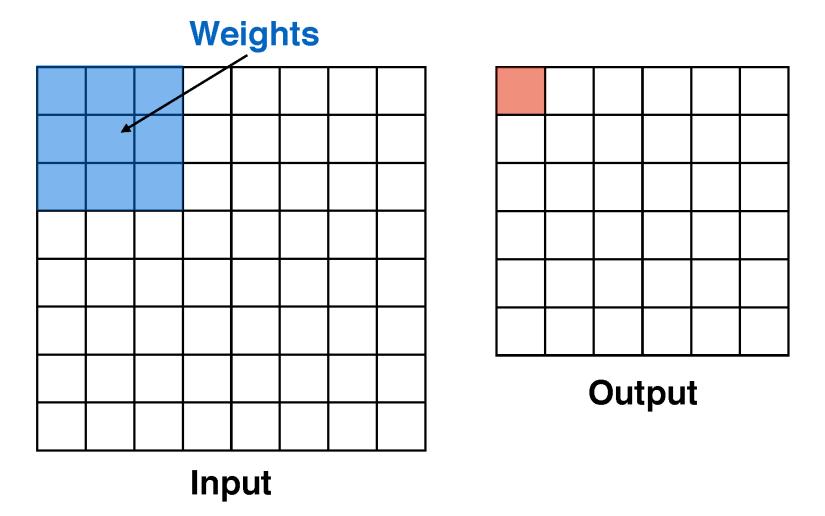


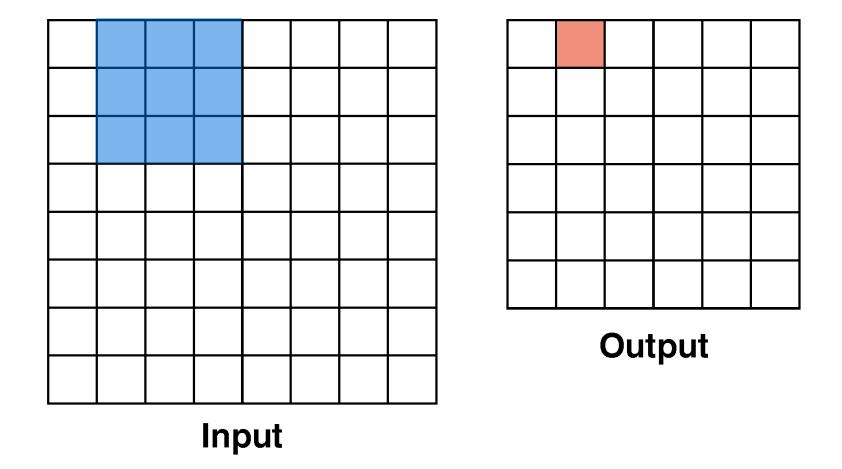


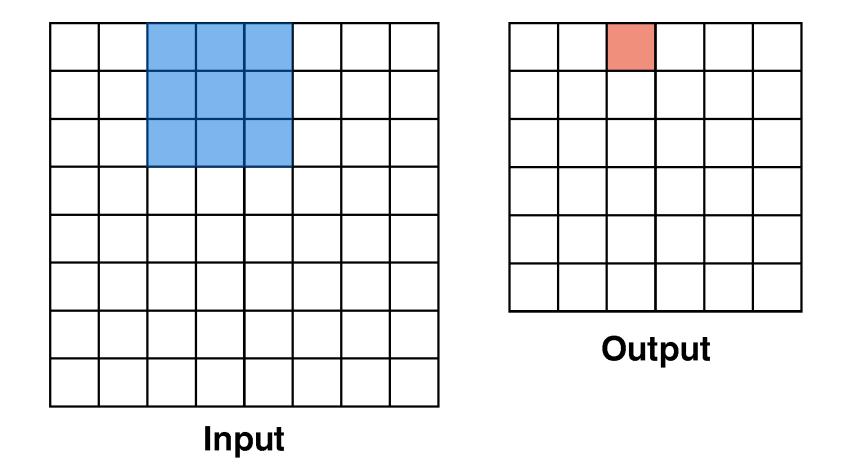
For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

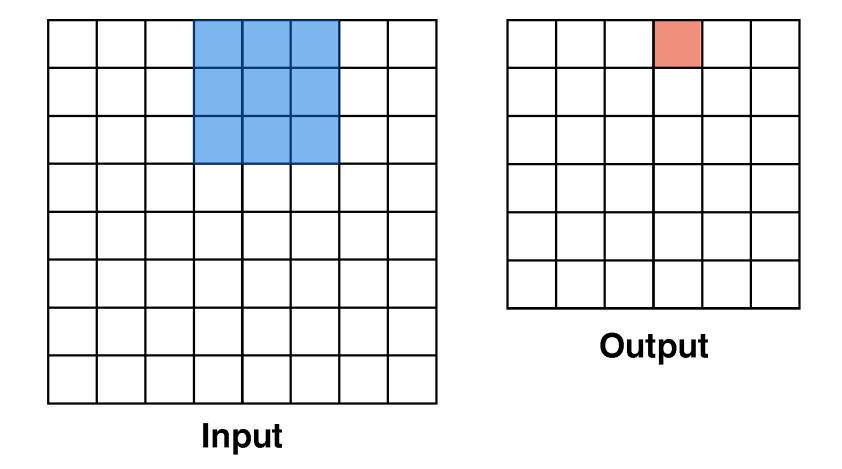


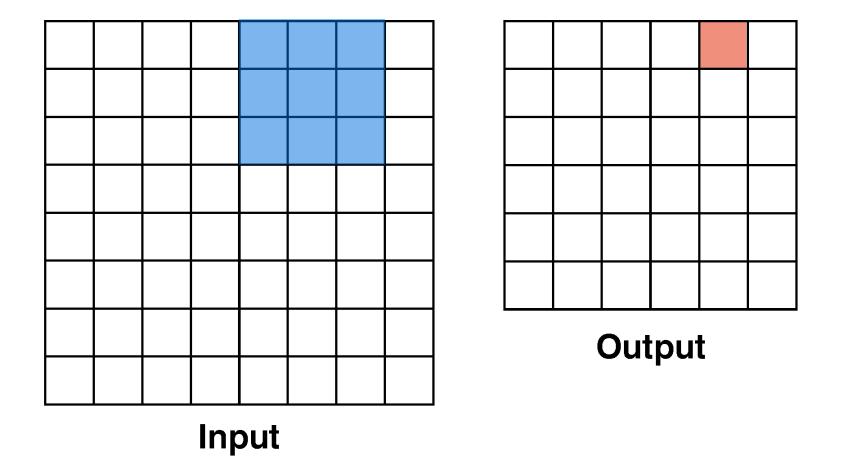
We stack these up to get a "new image" of size 28x28x6!

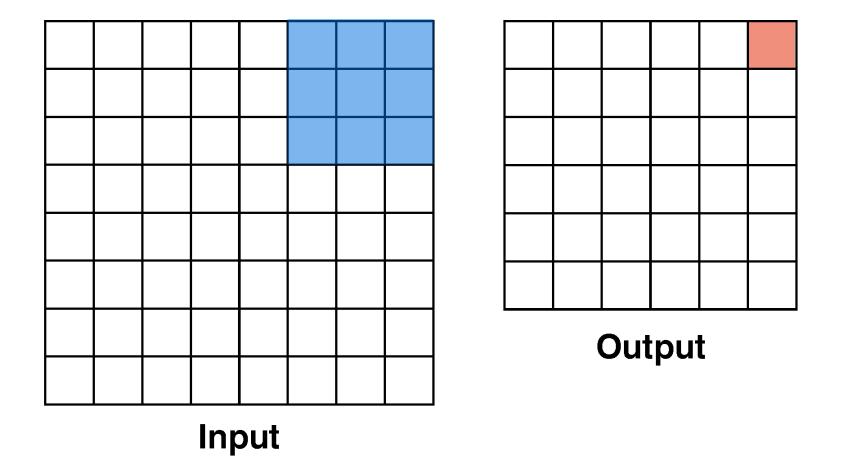




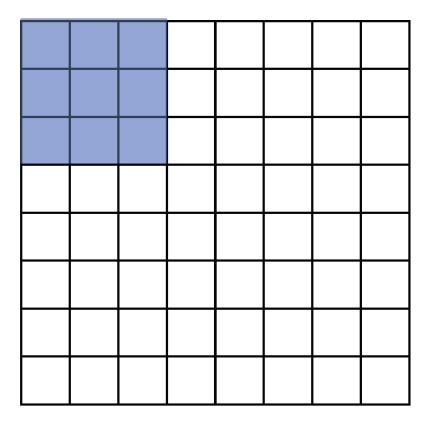








During convolution, the weights "slide" along the input to generate each output

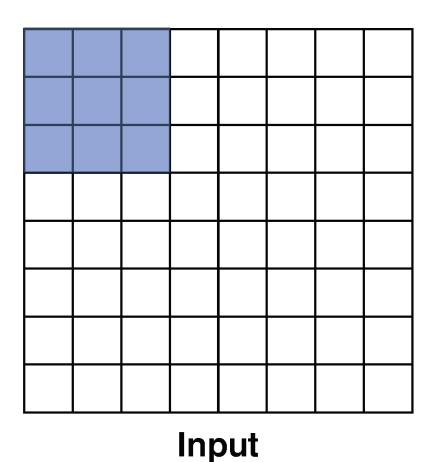


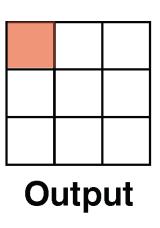
Recall that at each position, we are doing a **3D** sum:

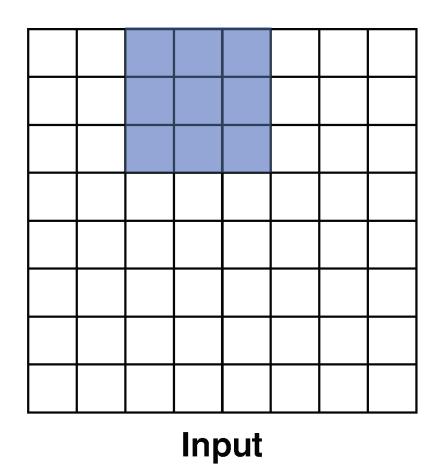
$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$

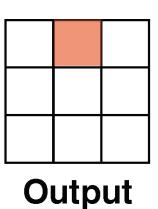
(channel, row, column)

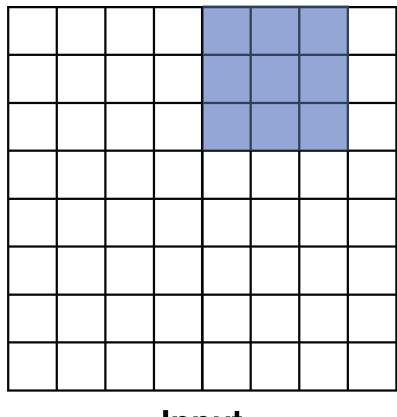
Input

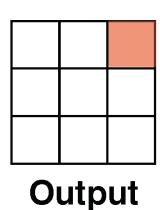




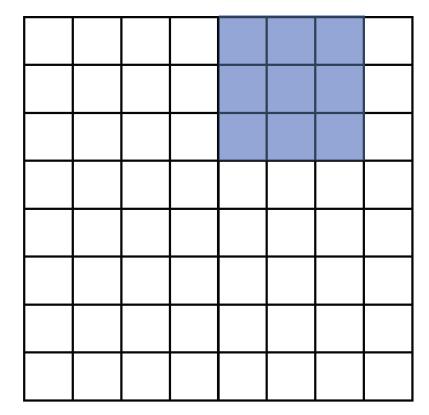




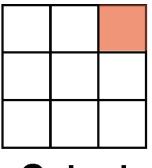




Input



Input

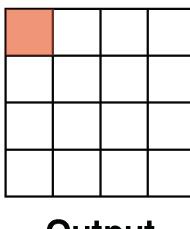


Output

- Notice that with certain strides, we may not be able to cover all of the input
- The output is also half the size of the input

We can also pad the input with zeros.

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

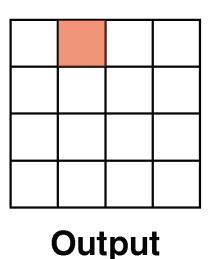


Output

Input

We can also pad the input with zeros.

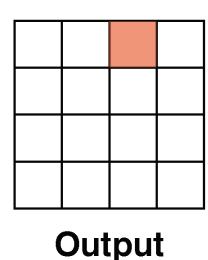
0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0



Input

We can also pad the input with zeros.

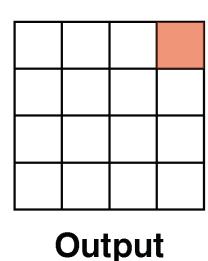
0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0



Input

We can also pad the input with zeros.

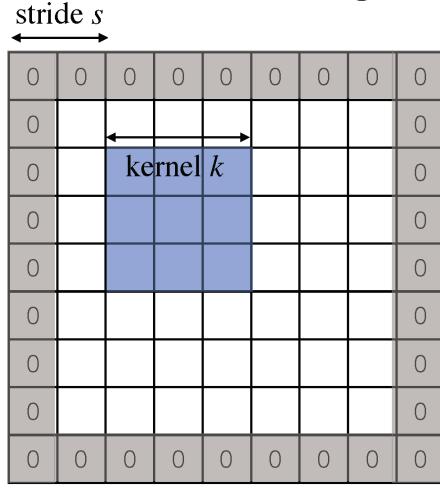
0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0



Input

Convolution:

How big is the output?



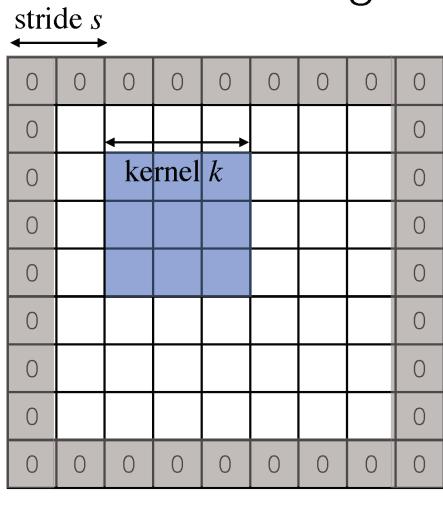
In general, the output has size:

$$w_{\text{out}} = \left[\frac{w_{\text{in}} + 2p - k}{s} \right] + 1$$

$$p \leftarrow \text{width } w_{\text{in}} \qquad p \rightarrow$$

Convolution:

How big is the output?



width $w_{\rm in}$

Example: k=3, s=1, p=1

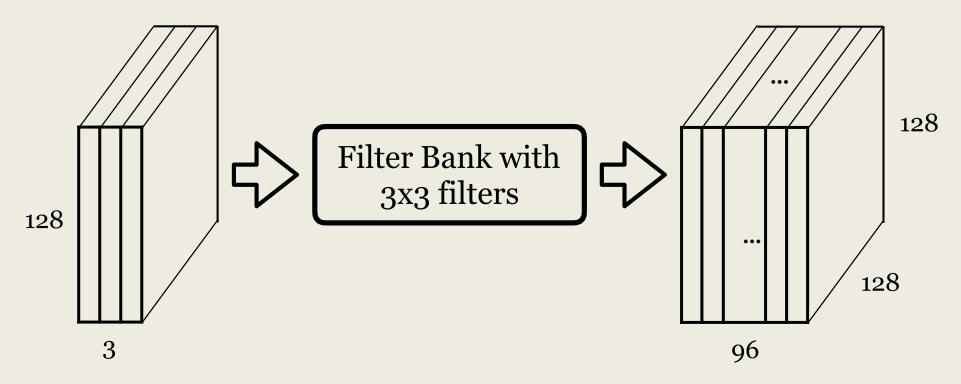
$$w_{\text{out}} = \left[\frac{w_{\text{in}} + 2p - k}{s} \right] + 1$$

$$= \left[\frac{w_{\text{in}} + 2 - 3}{1} \right] + 1$$

$$= w_{\text{in}}$$

VGGNet [Simonyan 2014] uses filters of this shape

Knowledge Check ...



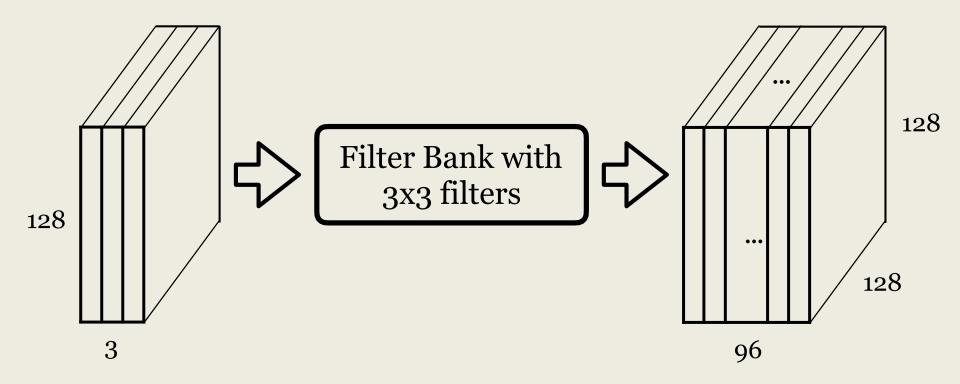
How many parameters does each filter have?

(a) 9 (b) 27

(b) 27 (c) 96

(d) 864

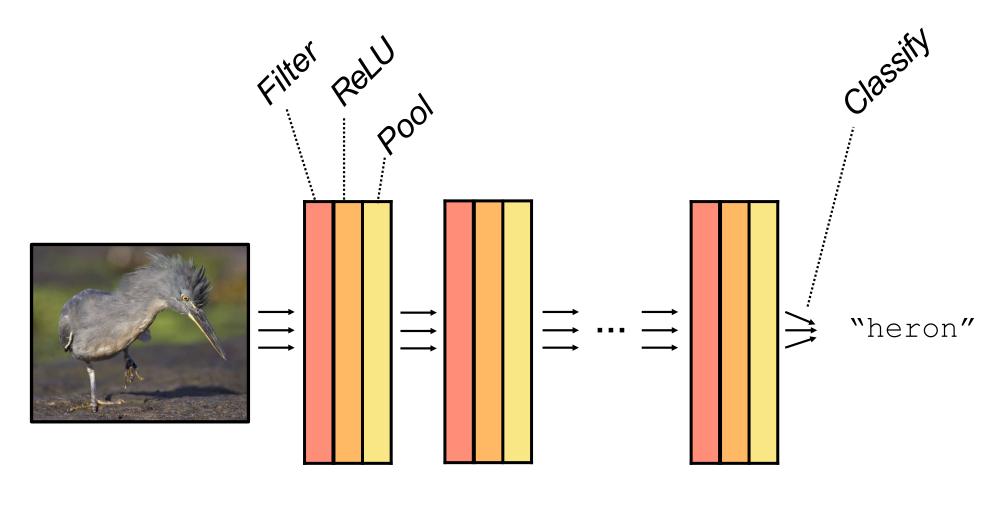
Knowledge Check ...



How many filters are in the bank?

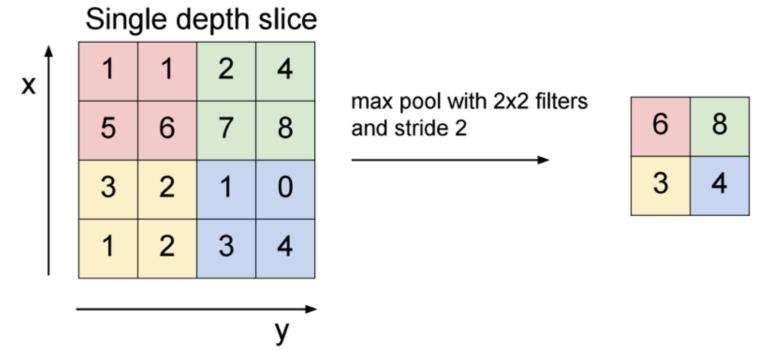
(a) 3 (b) 27 (c) 96 (d) can't say

Computation in a neural net



$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

Max Pooling

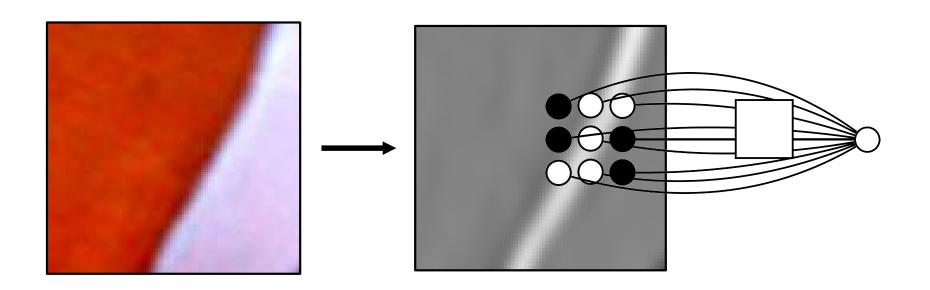


What's the backprop rule for max pooling?

- In the forward pass, store the index that took the max
- The backprop gradient is the input gradient at that index

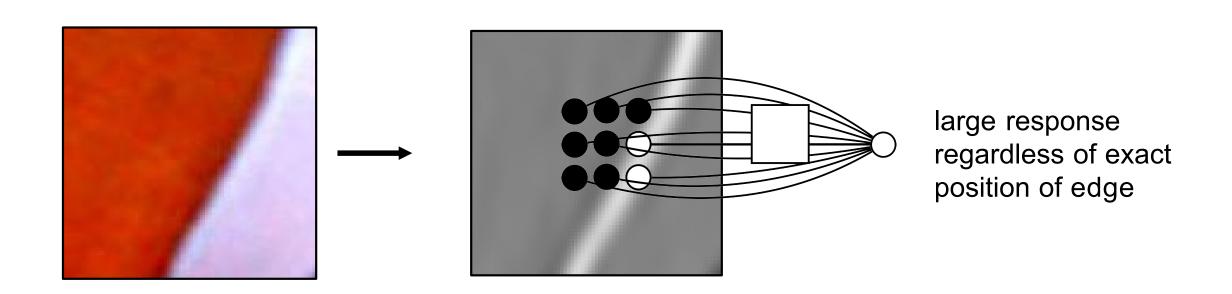
Pooling — Why?

Pooling across spatial locations achieves stability w.r.t. small translations:



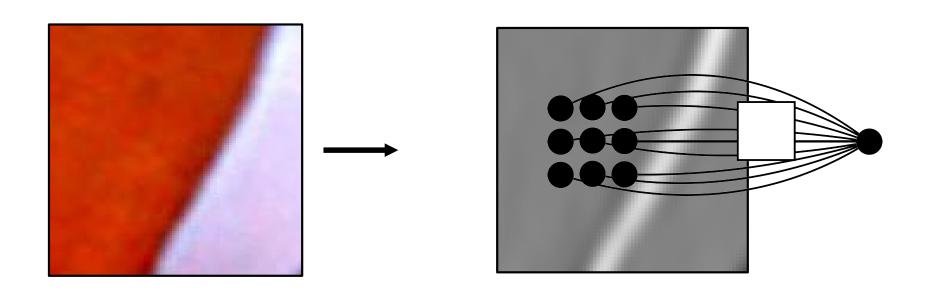
Pooling — Why?

Pooling across spatial locations achieves stability w.r.t. small translations:



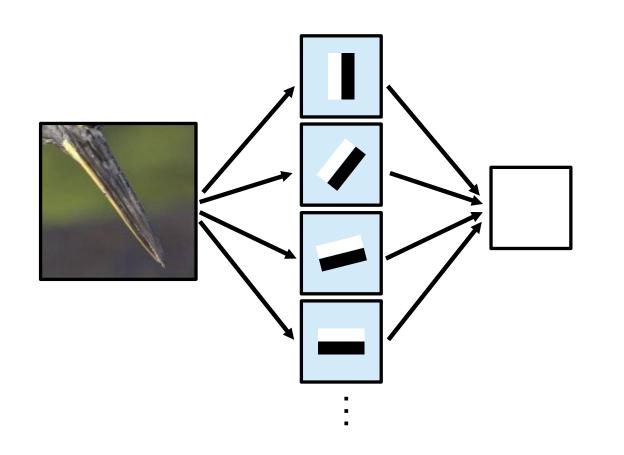
Pooling — Why?

Pooling across spatial locations achieves stability w.r.t. small translations:



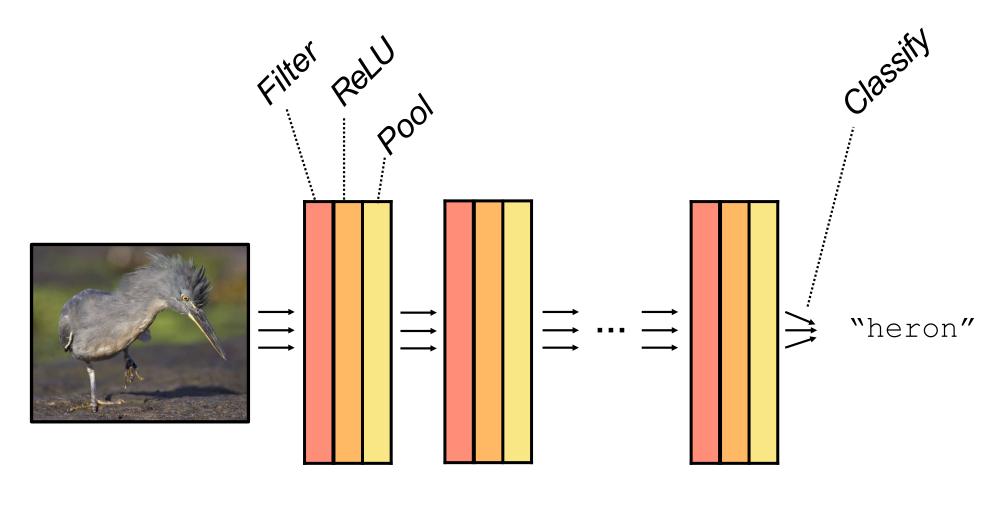
Pooling across channels — Why?

Pooling across feature channels (filter outputs) can achieve other kinds of invariances:



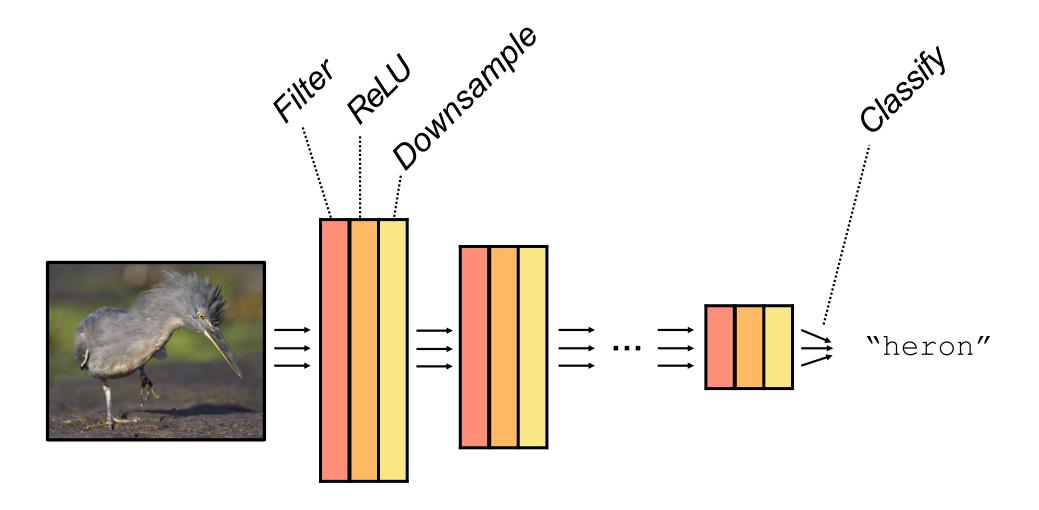
for any edge, regardless of its orientation

Computation in a neural net



$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

Computation in a neural net



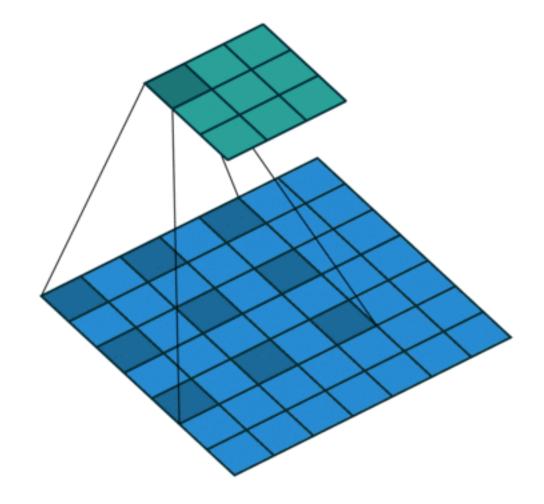
$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

Pooling vs Downsampling

Dilated Convolutions

Allows increasing the receptive field of the convolutional layer

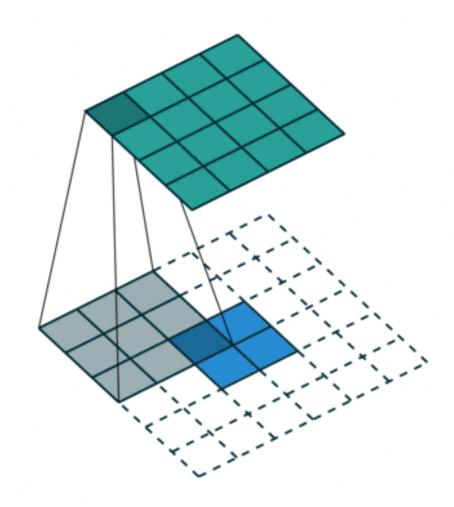
Useful for looking at larger spatial context without looking at every pixel



Transposed Convolution

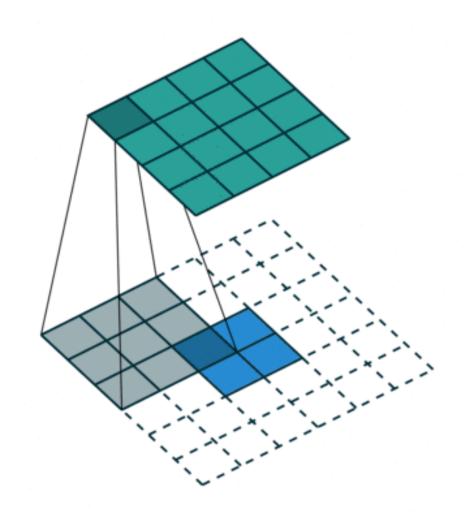
The transposed convolution a.k.a

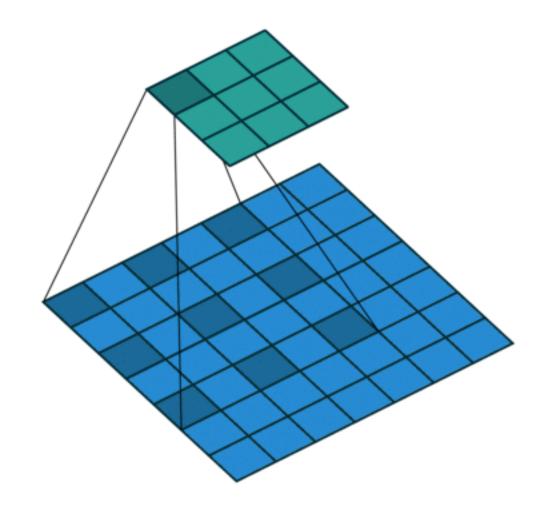
- deconvolution layer
- fractionally strided convolution



Transposed Conv Conv

vs. Dilated



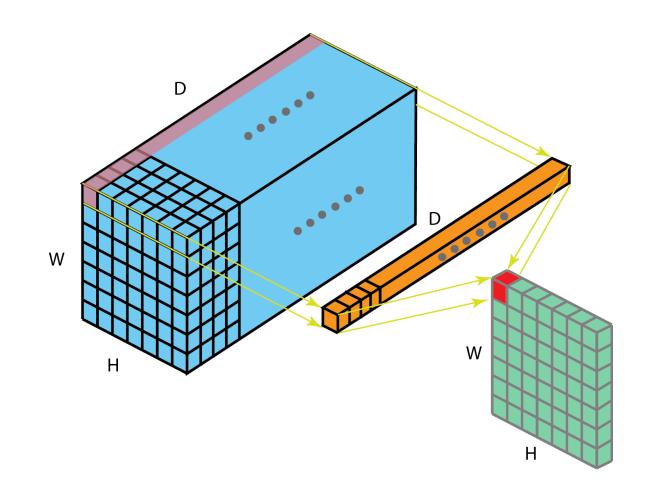


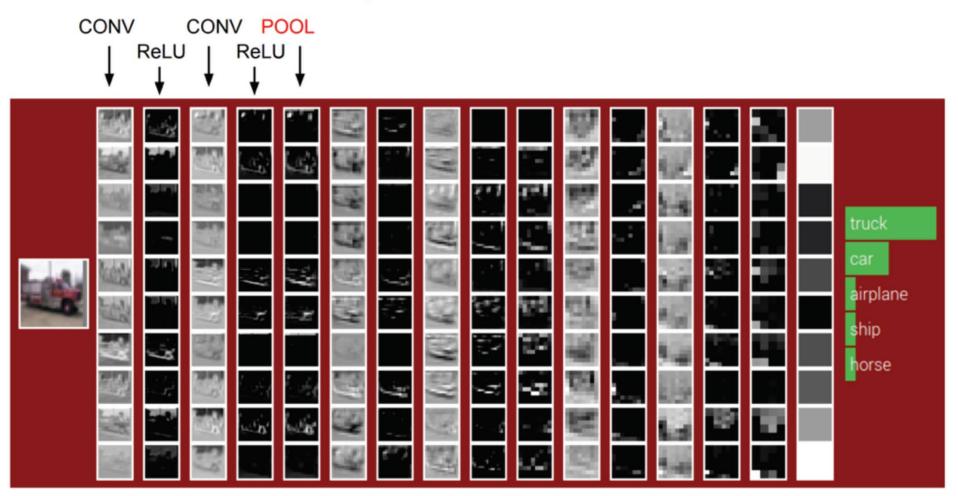
1x1 convolution

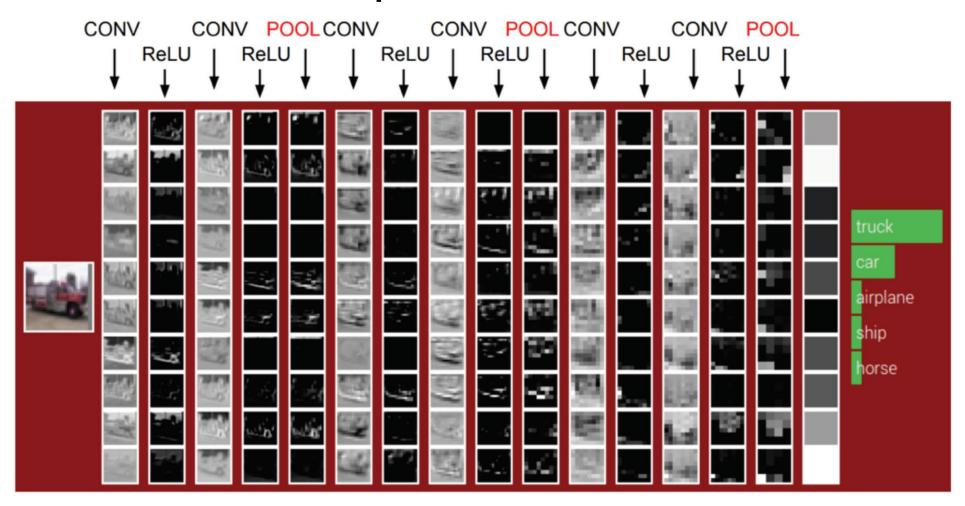
How is this not just multiplication?

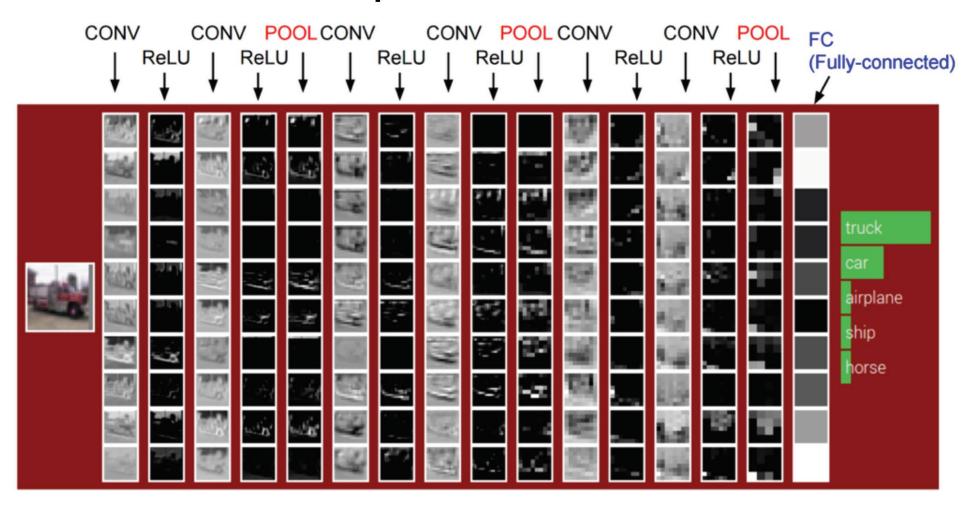
Multiplications followed by a RELU activation

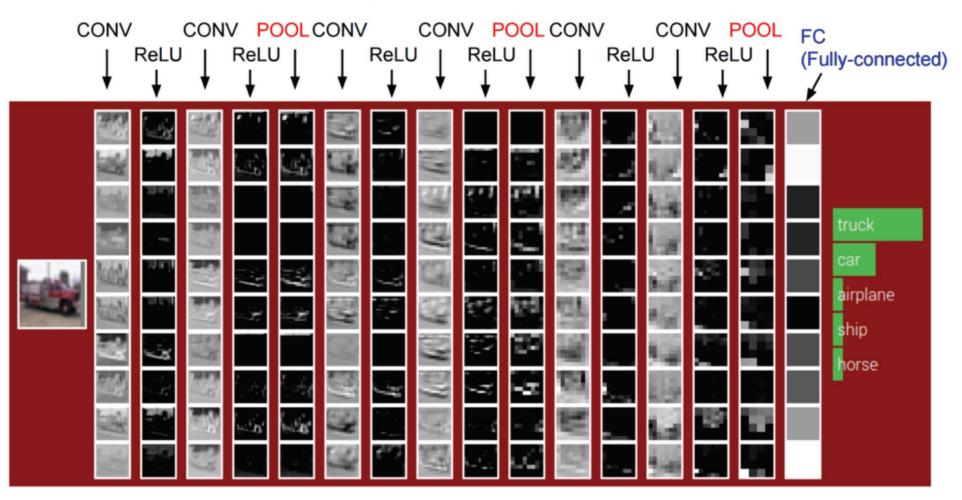
Good for dimensionality reduction, efficient storage





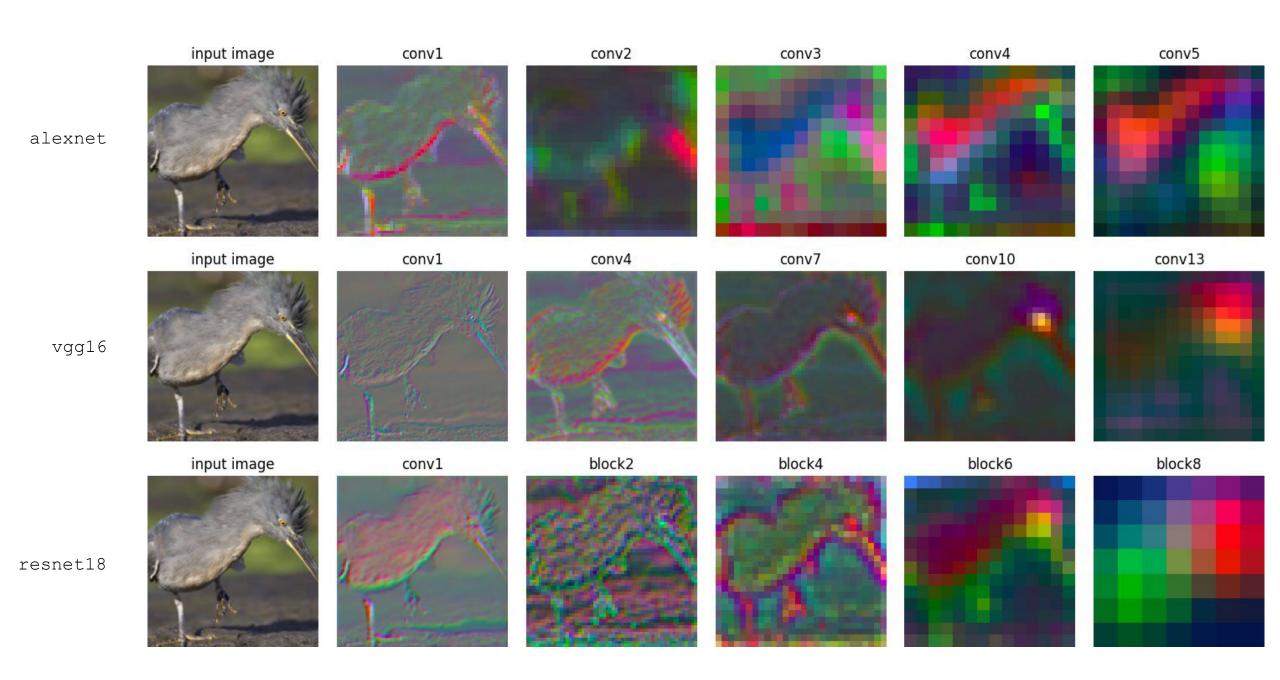






10x3x3 conv filters, stride 1, pad 1 2x2 pool filters, stride 2

Figure: Andrej Karpathy



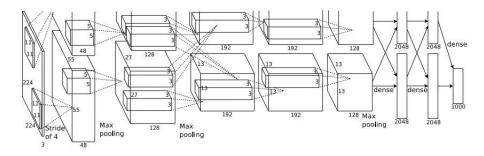
Layer Visualizations



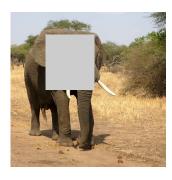
Which pixels matter: Saliency via Occlusion

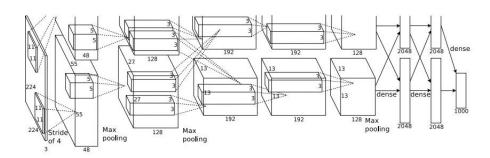
Mask part of the image before feeding to CNN, check how much predicted probabilities change





P(elephant) = 0.95



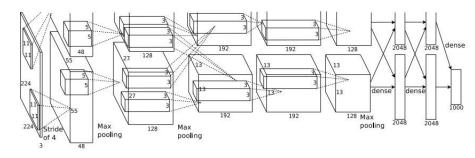


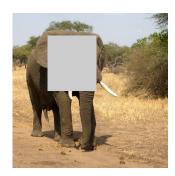
P(elephant) = 0.75

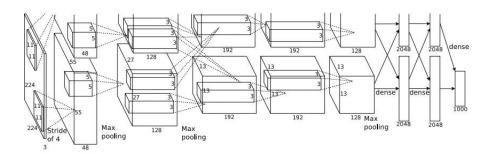
Which pixels matter: Saliency via Occlusion

Mask part of the image before feeding to CNN, check how much predicted probabilities change







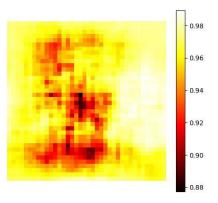


Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014

Boat image is CC0 public domain

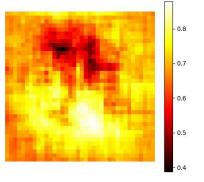
Elephant image is CC0 public domain
Go-Karts image is CC0 public domain





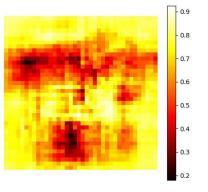
African elephant, Loxodonta africana





go-kart





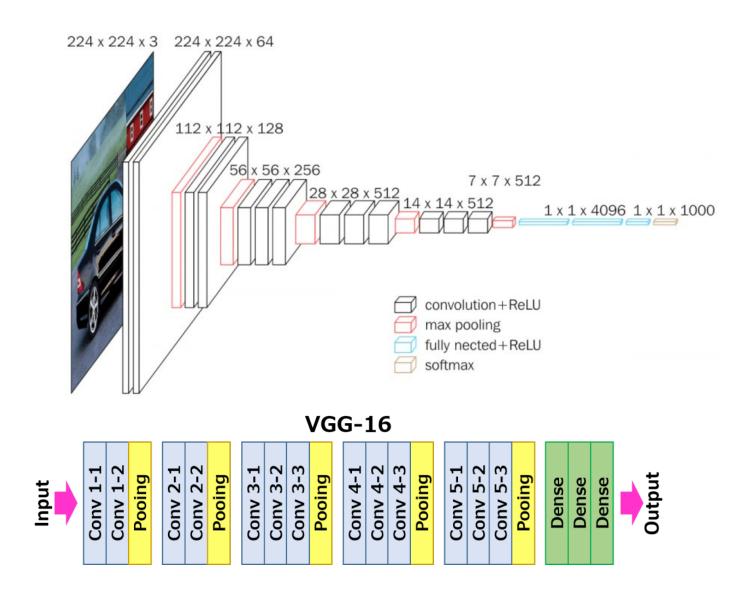
Common Architectures

Common Architectures

(We will revisit these in detail when we study specific tasks like object detection, image generation, etc.)

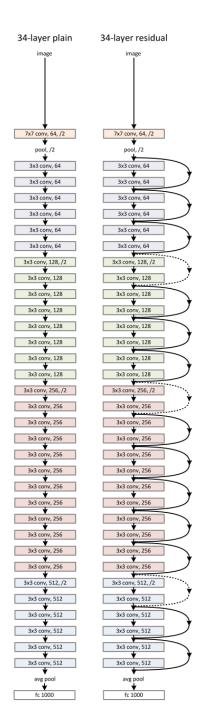
VGG

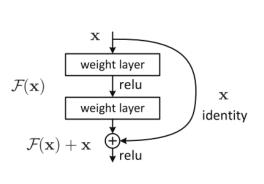
- Simonyan and Zisserman,
 "Very Deep Convolutional Networks for Large-Scale Image Recognition"
- Used to be very common (before ResNets)



ResNet

- He, Kaiming; Zhang, Xiangyu; Ren, Shaoqing; Sun, Jian (2016). "Deep Residual Learning for Image Recognition" (PDF). Proc. Computer Vision and Pattern Recognition (CVPR), IEEE.
- Deep networks with more layers does not always mean better performance (vanishing gradient problem)
- Residual blocks = has skip connections
- Skipped layers train faster at the beginning, then later are filled out



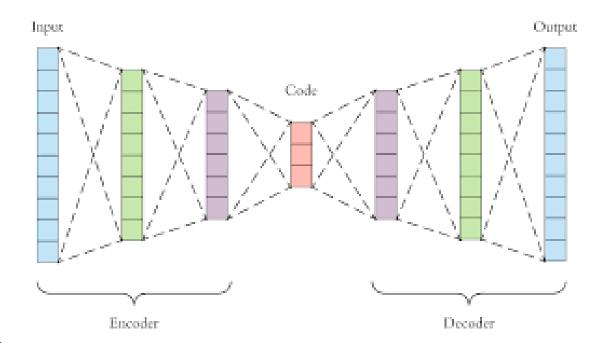


Autoencoder

 Can be done with either fully connected or convolutional layers

 Idea is to reduce the input to a bottleneck or latent code, then reconstruct it again

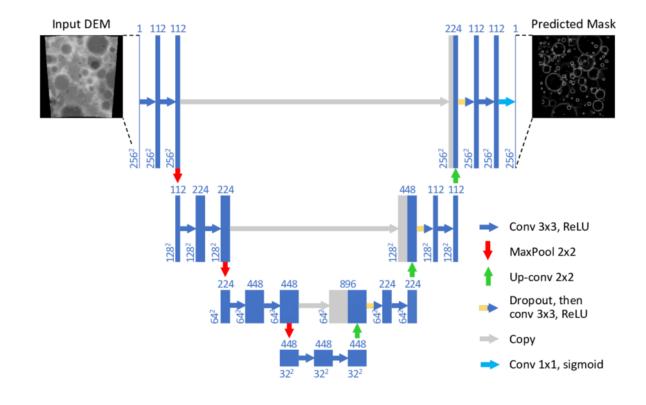
 Sometimes can be used to train a feature extractor by enforcing the output = input, and then use the first part of the network as a feature extractor

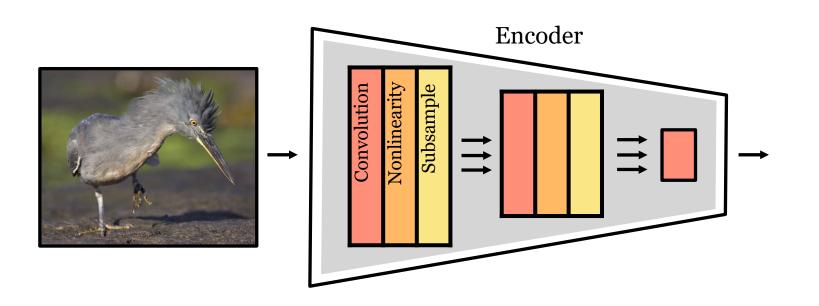


U-Net

Common architecture for image reconstruction tasks

 Features skip connections and transposed convolutions (up-conv)





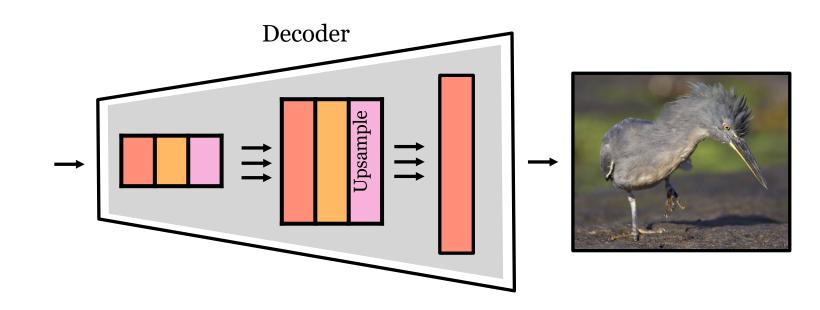


Image-to-image

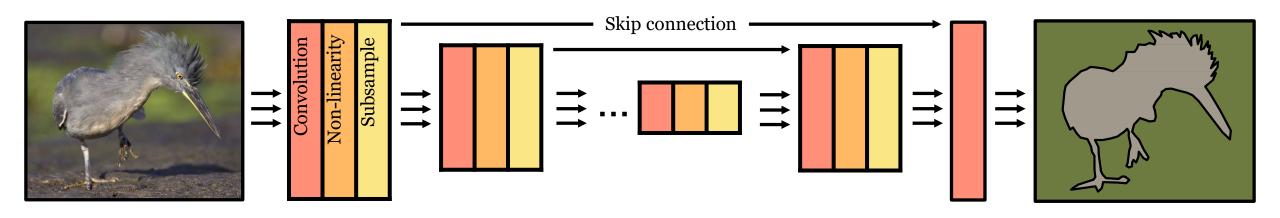
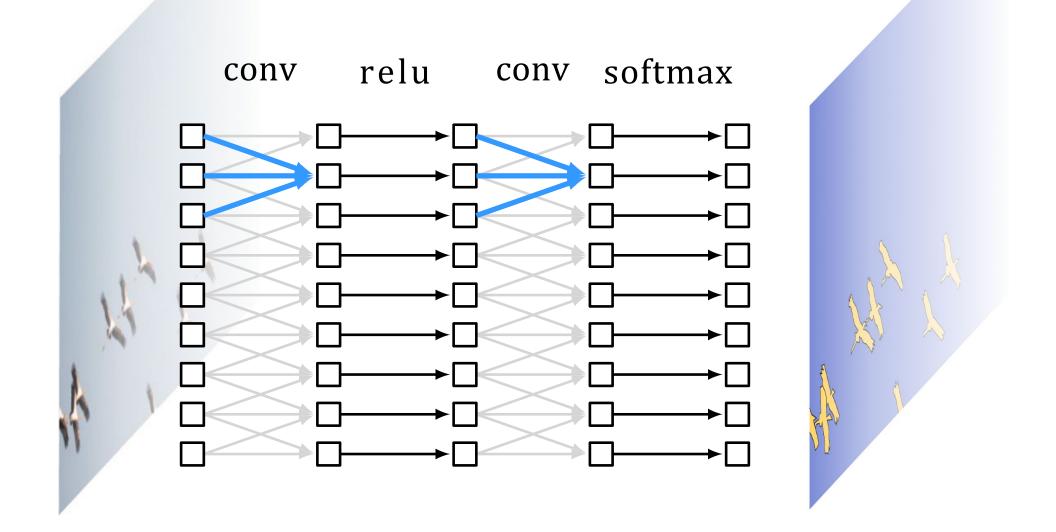
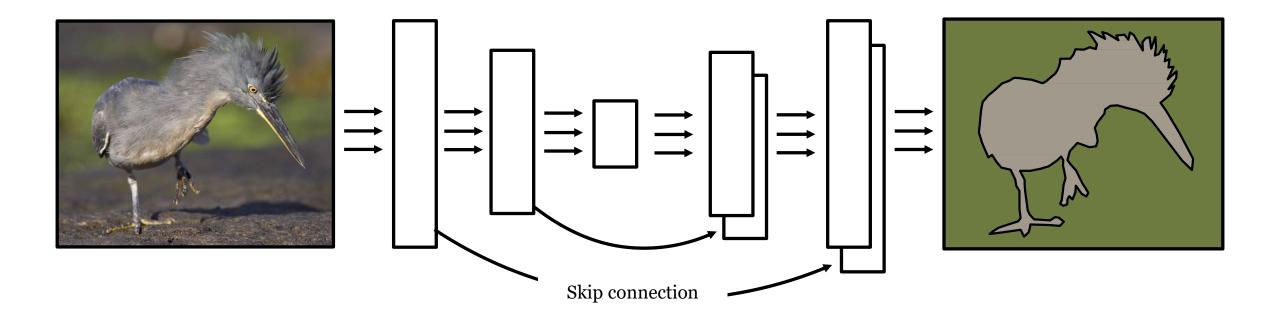


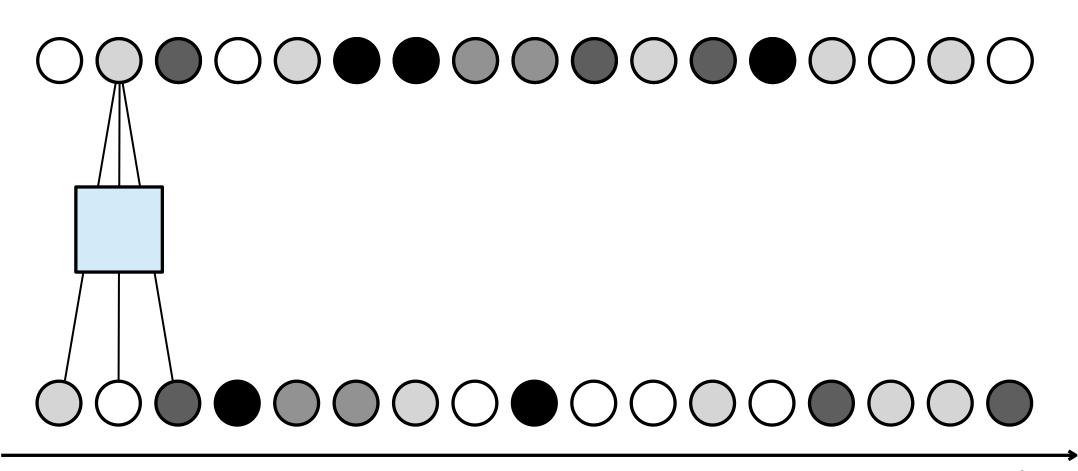
Image-to-image



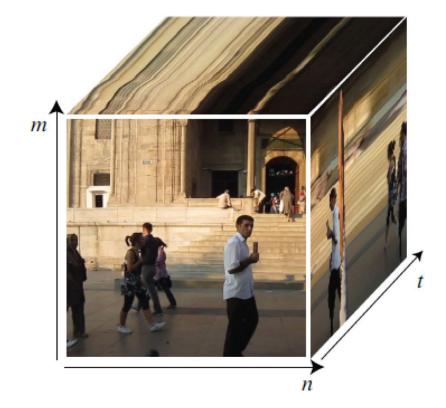
U-net



Convolutions in time







Training ConvNets

How do you actually train these things?

Roughly speaking:

Gather labeled data

Find a ConvNet architecture

Minimize the loss



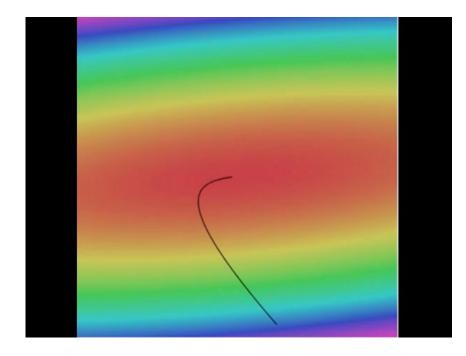




Recap

Learning network parameters through optimization





```
# Vanilla Gradient Descent
while True:
  weights_grad = evaluate_gradient(loss_fun, data, weights)
  weights += - step_size * weights_grad # perform parameter update
```

<u>Landscape image</u> is <u>CC0 1.0</u> public domain <u>Walking man image</u> is <u>CC0 1.0</u> public domain

Mini-batch Gradient Descent

Loop:

- 1. Sample a batch of training data (~100 images)
- 2. Forwards pass: compute loss (avg. over batch)
- 3. Backwards pass: compute gradient
- 4. Update all parameters

Note: usually called "stochastic gradient descent" even though SGD has a batch size of 1

Training a convolutional neural network

- Split and preprocess your data
- Choose your network architecture
- Initialize the weights
- Find a learning rate and regularization strength
- Minimize the loss and monitor progress
- Fiddle with knobs

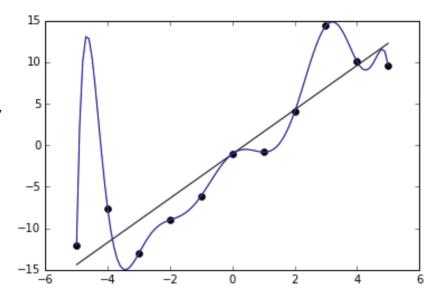
Recall: Overfitting

Overfitting

Overfitting: modeling noise in the training set instead of the "true" underlying relationship

Underfitting: insufficiently modeling the relationship in the training set

General rule: models that are "bigger" or have more capacity are more likely to overfit



[Image: https://en.wikipedia.org/wiki/File:Overfitted_Data.png]

Regularization

Regularization reduces overfitting:

$$L = L_{\text{data}} + L_{\text{reg}} \qquad \qquad L_{\text{reg}} = \lambda \frac{1}{2} ||W||_2^2$$

$$\lambda = 0.001 \qquad \qquad \lambda = 0.01$$

$$\lambda = 0.1$$

[Andrej Karpathy http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html]

Example Regularizers

L2 regularization

$$L_{\text{reg}} = \lambda \frac{1}{2} ||W||_2^2$$

(L2 regularization encourages small weights)

L1 regularization

$$L_{\text{reg}} = \lambda ||W||_1 = \lambda \sum_{ij} |W_{ij}|$$

(L1 regularization encourages sparse weights: weights are encouraged to reduce to exactly zero)

$$L_{\text{reg}} = \lambda_1 ||W||_1 + \lambda_2 ||W||_2^2$$

(combine L1 and L2 regularization)

Max norm

Clamp weights to some max norm

$$||W||_2^2 \le c$$

"Weight decay"

Regularization is also called "weight decay" because the weights "decay" each iteration:

$$L_{\text{reg}} = \lambda \frac{1}{2} ||W||_2^2 \longrightarrow \frac{\partial L}{\partial W} = \lambda W$$

Gradient descent step:

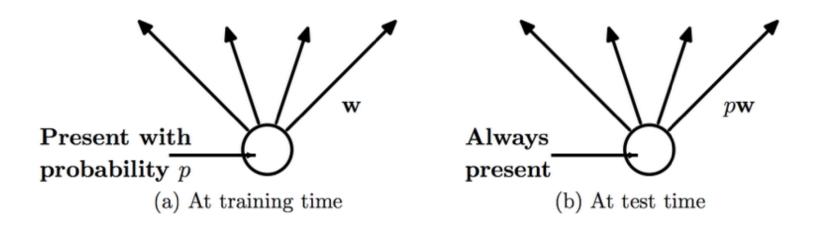
$$W \leftarrow W - \alpha \lambda W - \frac{\partial L_{\text{data}}}{\partial W}$$

Weight decay: $\alpha\lambda$ (weights always decay by this amount)

Note: biases are sometimes excluded from regularization

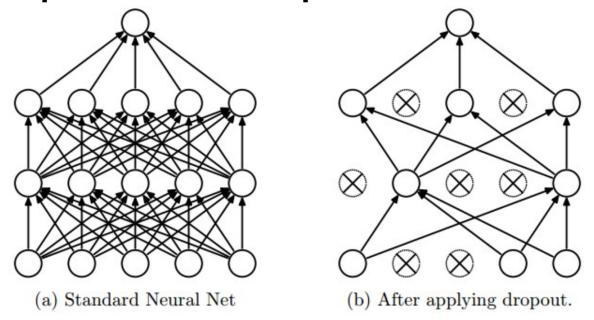
Dropout

Simple but powerful technique to reduce overfitting:



Dropout

Simple but powerful technique to reduce overfitting:



Note: Dropout can be interpreted as an approximation to taking the geometric mean of an ensemble of exponentially many models

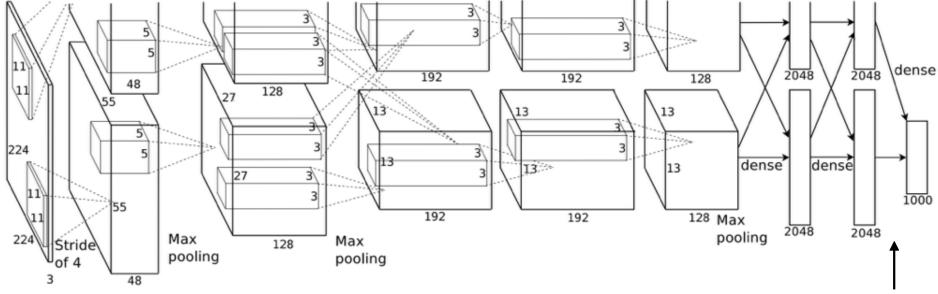
[Srivasta et al, "Dropout: A Simple Way to Prevent Neural Networks from Overfitting", JMLR 2014]

Dropout

Case study: [Krizhevsky 2012]

"Without dropout, our network exhibits substantial overfitting."

Dropout here 2048 2048 128



But not here — why?

[Krizhevsky et al, "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012]

Summary of things to fiddle

- Network architecture
- Learning rate, decay schedule, update type
- Regularization (L2, L1, maxnorm, dropout, ...)
- Loss function (softmax, SVM, ...)
- Weight initialization

Neural network parameters

