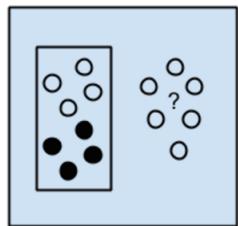


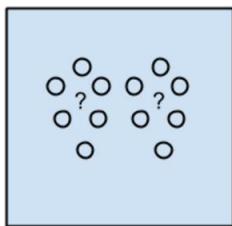


# Lecture 4: Computer Vision

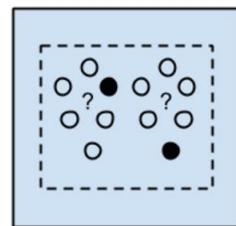
# Computer Vision is Deep Learning



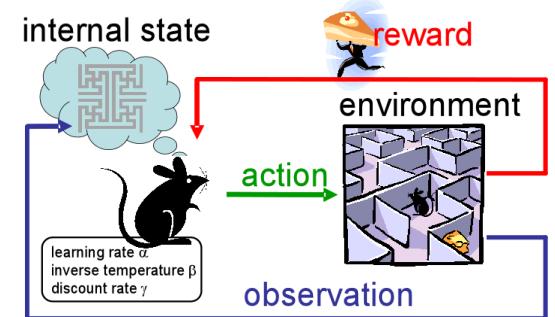
Supervised  
Learning



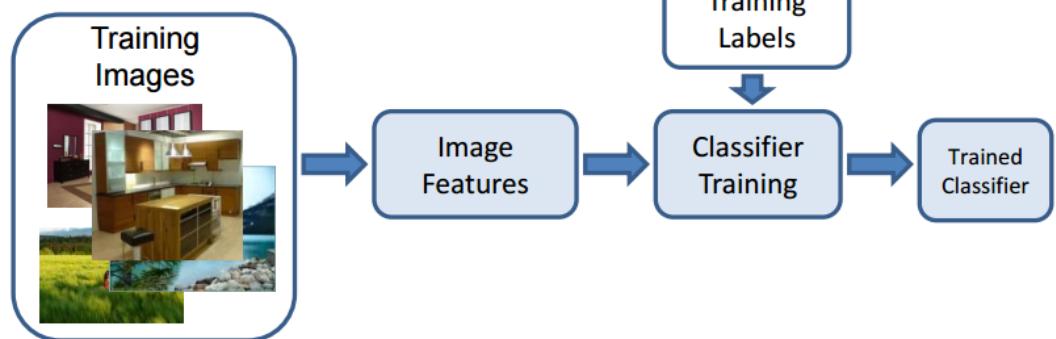
Unsupervised  
Learning



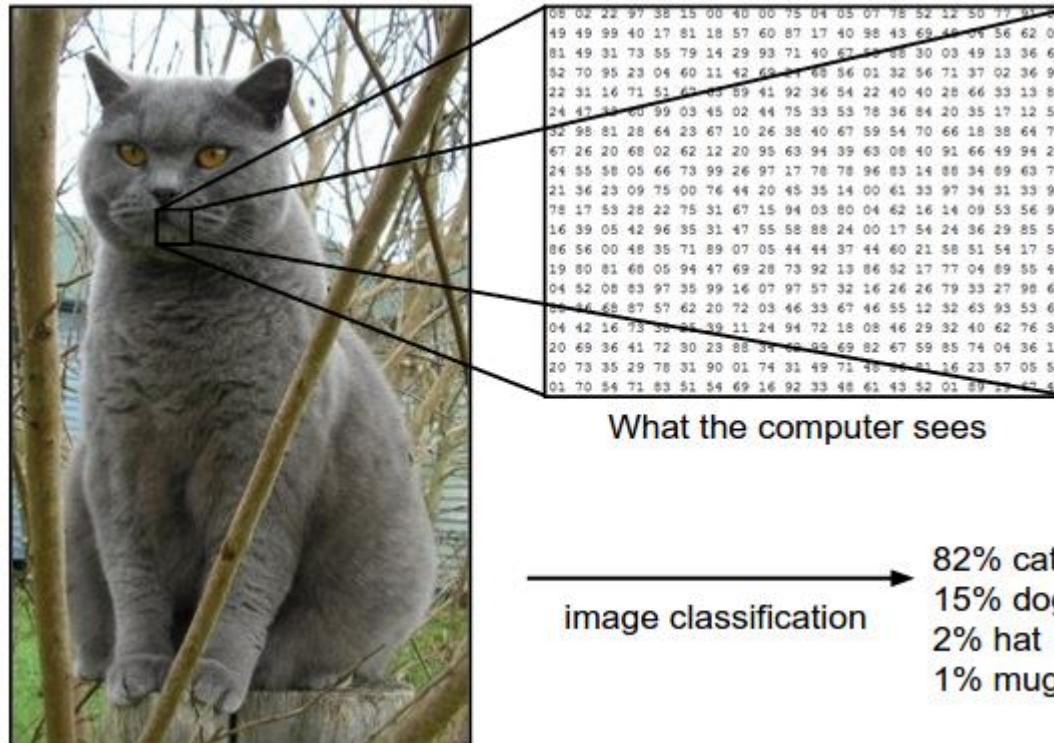
Semi-Supervised  
Learning



Reinforcement  
Learning



# Images are Numbers



- **Regression:** The output variable takes continuous values
- **Classification:** The output variable takes class labels
  - Underneath it may still produce continuous values such as probability of belonging to a particular class.

# Computer Vision with Deep Learning:

Our intuition about what's "hard" is flawed (in complicated ways)

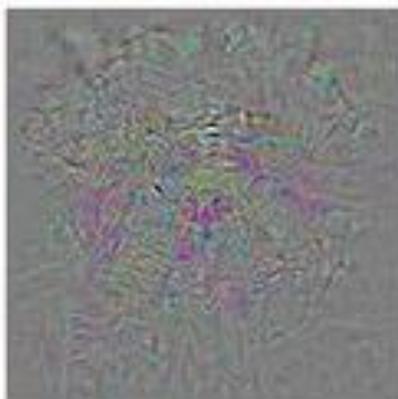
**Visual perception:** 540,000,000 years of data

**Bipedal movement:** 230,000,000 years of data

**Abstract thought:** 100,000 years of data



Prediction: Dog



+ Distortion

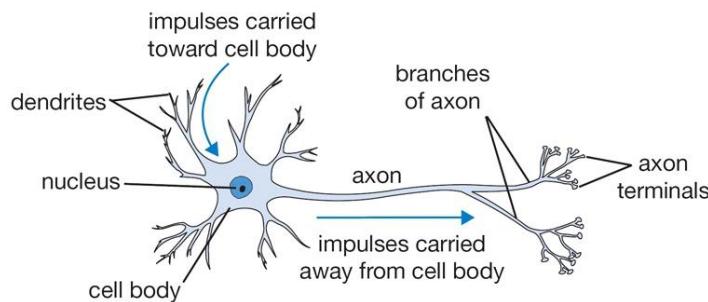


Prediction: Ostrich

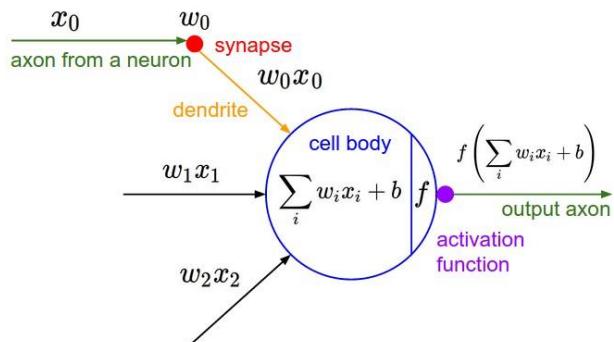
"Encoded in the large, highly evolved sensory and motor portions of the human brain is a **billion years of experience** about the nature of the world and how to survive in it.... Abstract thought, though, is a new trick, perhaps less than **100 thousand years old**. We have not yet mastered it. It is not all that intrinsically difficult; it just seems so when we do it."

- Hans Moravec, *Mind Children* (1988)

# Neuron: Biological Inspiration for Computation



- **Neuron:** computational building block for the brain
- **(Artificial) Neuron:** computational building block for the “neural network”



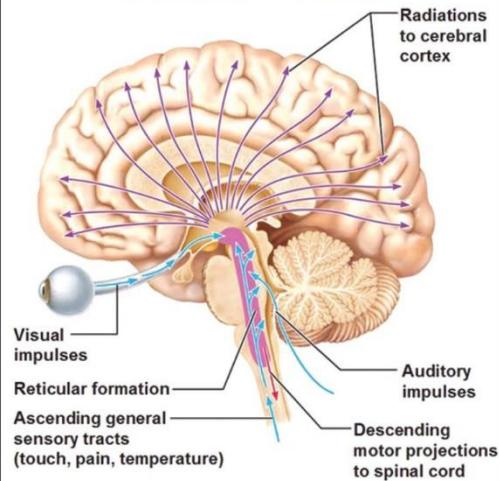
## Differences (among others):

- **Parameters:** Human brains have  $\sim 10,000,000$  times synapses than artificial neural networks.
- **Topology:** Human brains have no “layers”. Topology is complicated.
- **Async:** The human brain works asynchronously, ANNs work synchronously.
- **Learning algorithm:** ANNs use gradient descent for learning. Human brains use ... (we don't know)
- **Processing speed:** Single biological neurons are slow, while standard neurons in ANNs are fast.
- **Power consumption:** Biological neural networks use very little power compared to artificial networks
- **Stages:** Biological networks usually don't stop / start learning. ANNs have different fitting (train) and prediction (evaluate) phases.

## Similarity (among others):

- Distributed computation on a large scale.

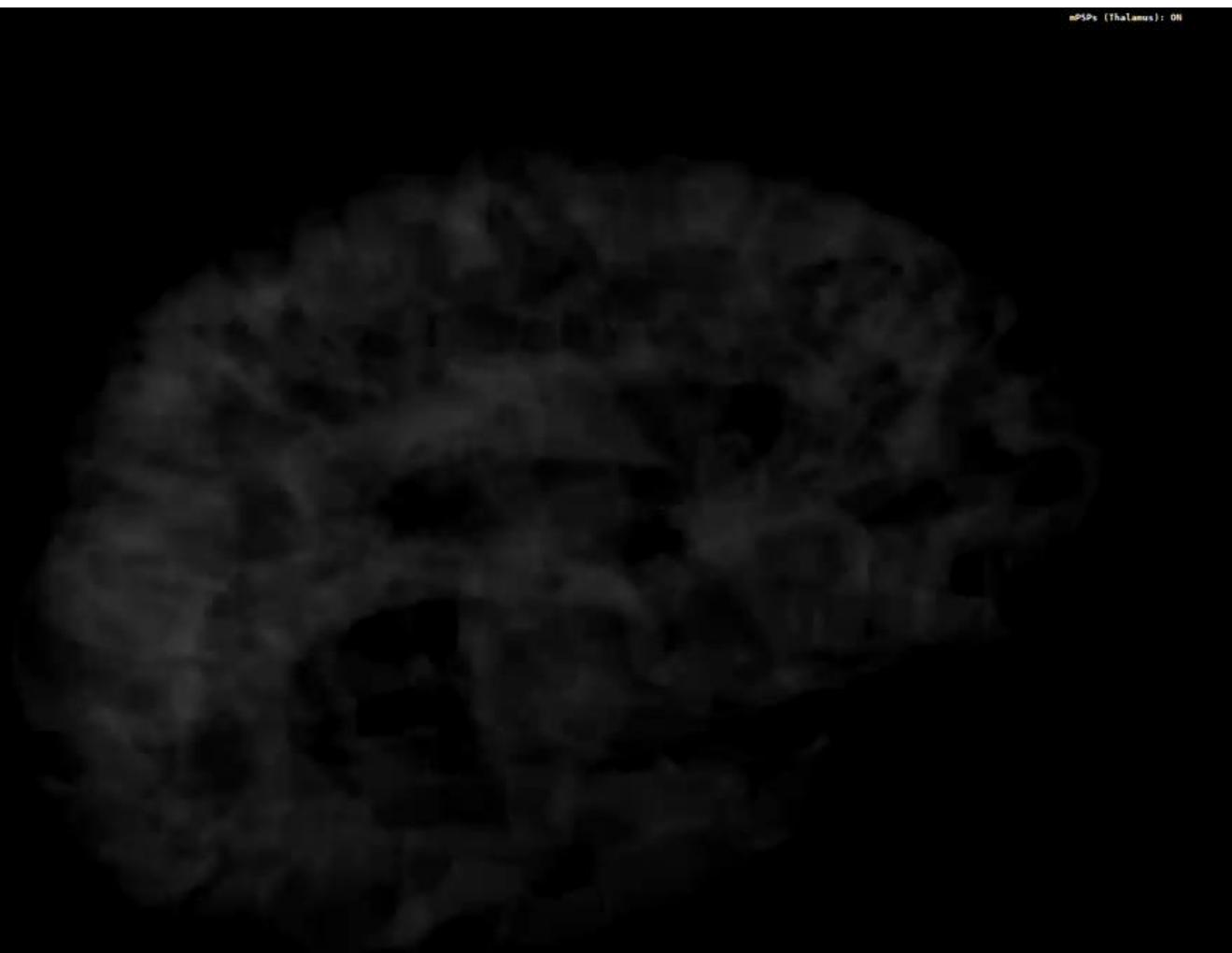
## The Reticular Formation



# Human Vision

Its structure is instructive and inspiring!

*Thalamocortical System Simulation: 8 million cortical neurons + 2 billion synapses:*

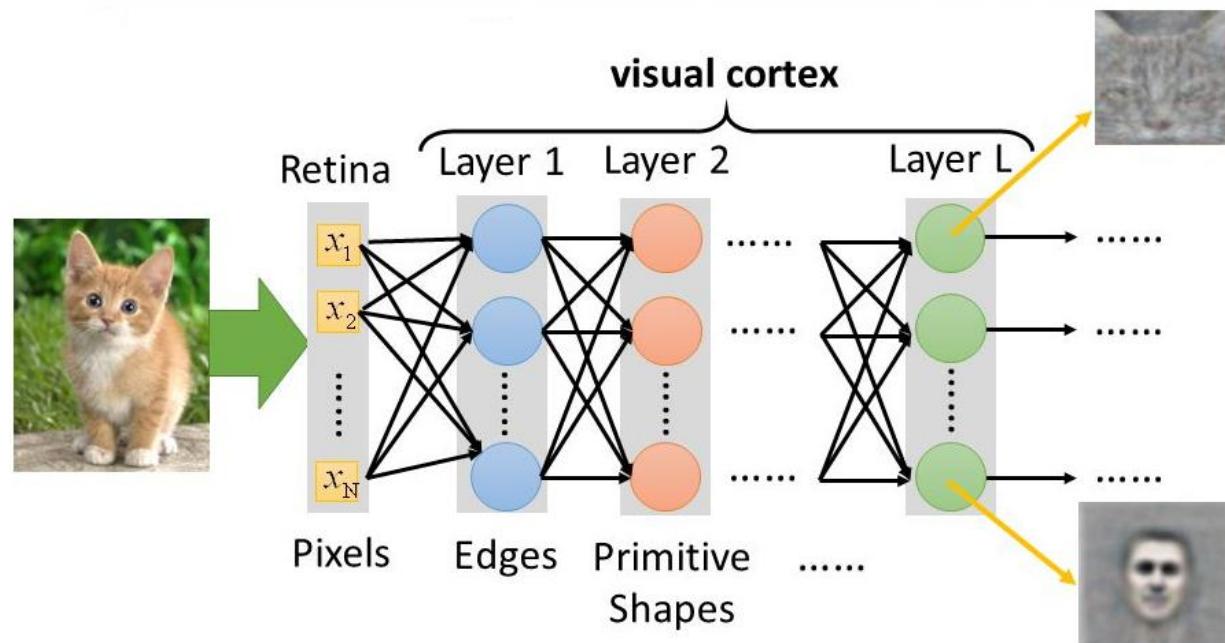
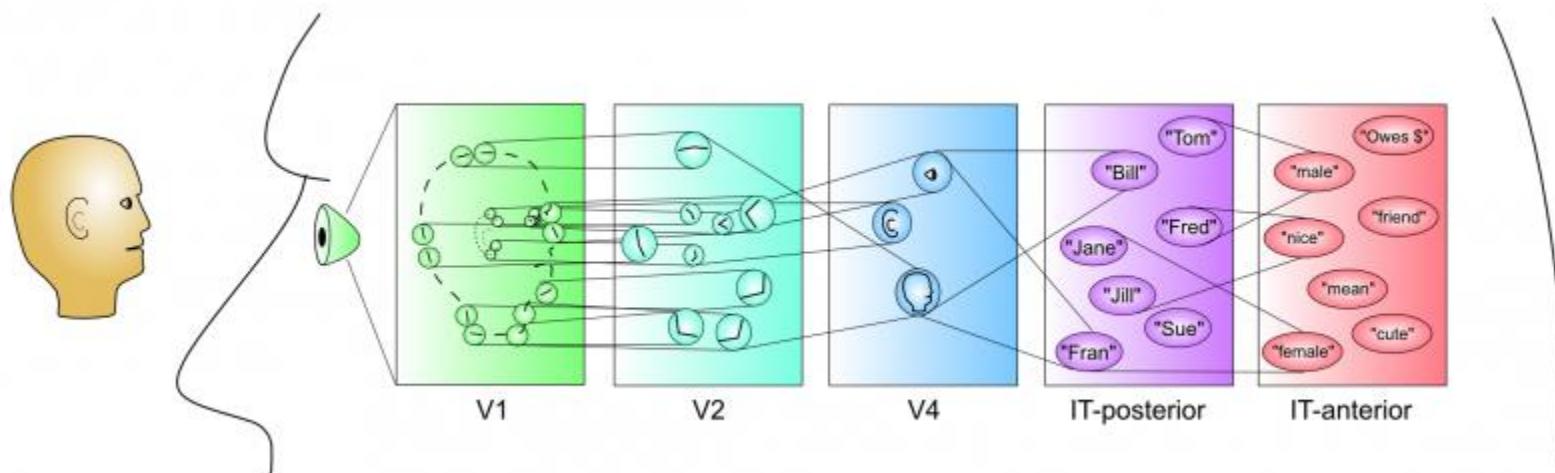


Retinal Ganglion Cell Activity



# Visual Cortex

(Its Structure is Instructive and Inspiring)



# Deep Learning is Hard: Illumination Variability



# Deep Learning is Hard: Pose Variability

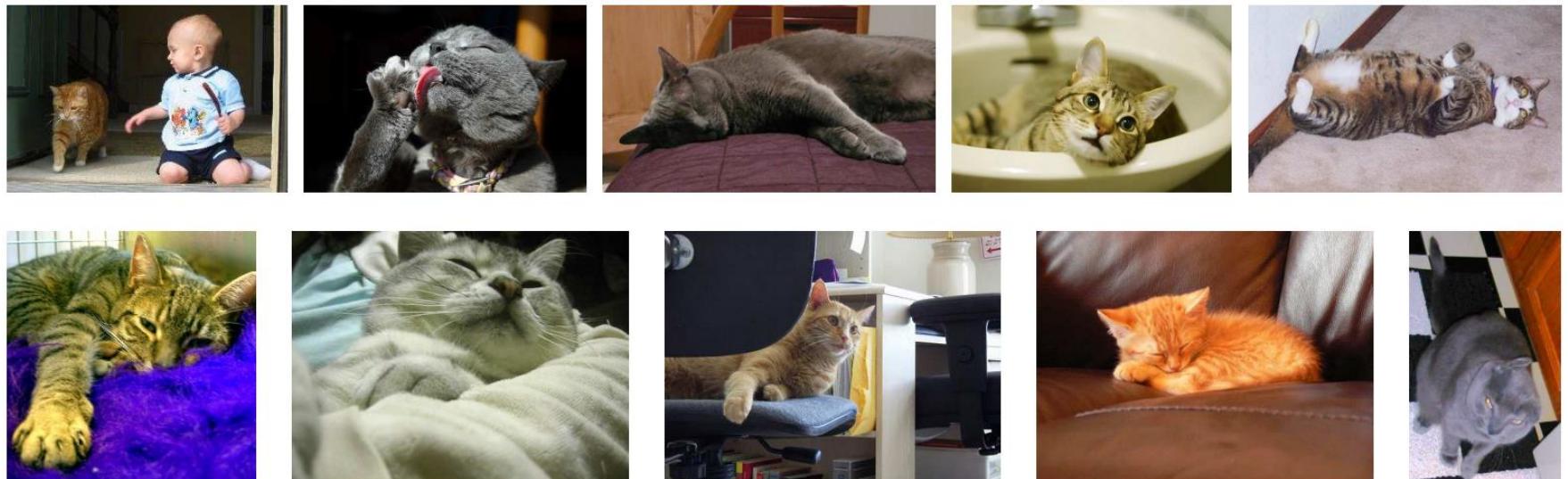


Figure 1. **The deformable and truncated cat.** Cats exhibit (al-

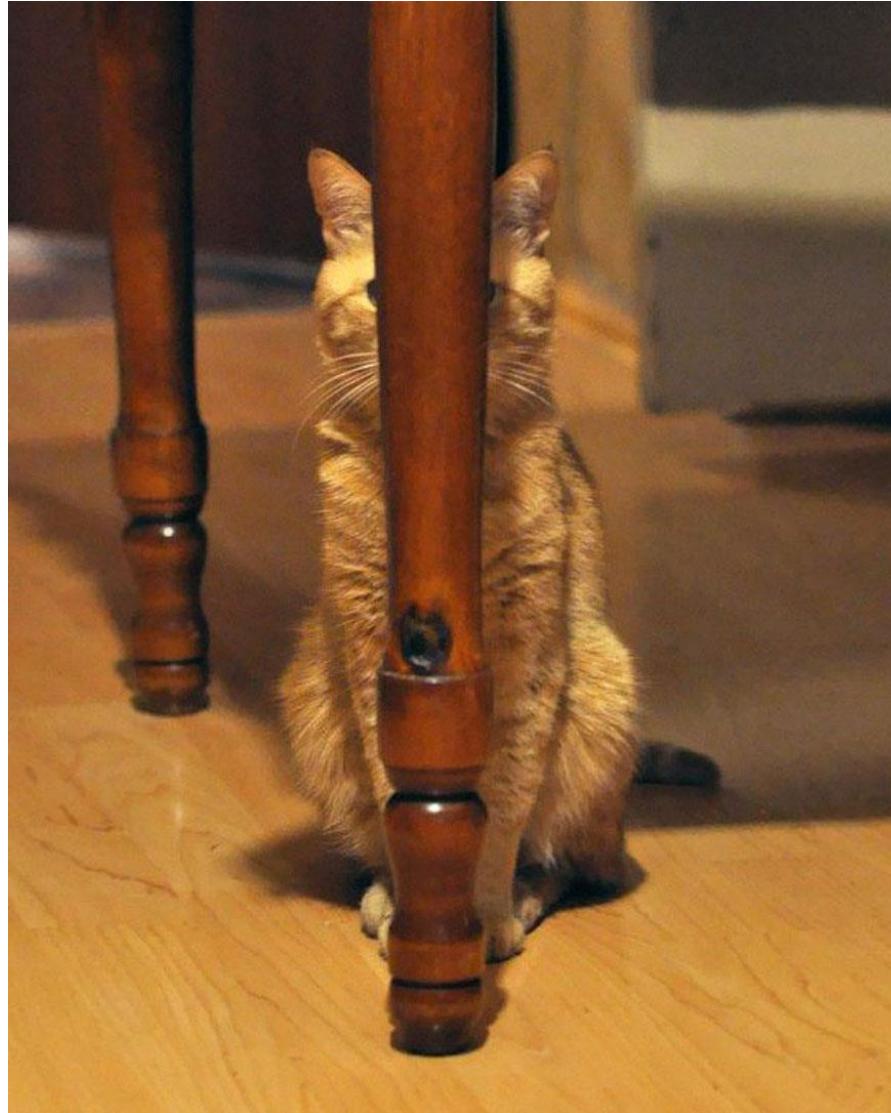
Parkhi et al. "The truth about cats and dogs." 2011.

# Deep Learning is Hard: Intra-Class Variability



Parkhi et al. "Cats and dogs." 2012.

# Occlusion



# Occlusion



# Occlusion

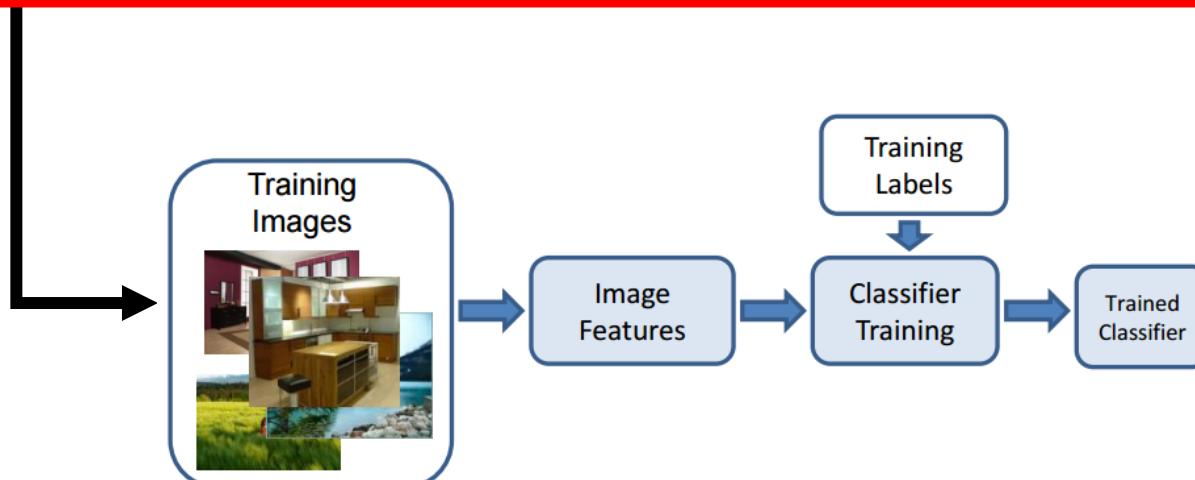
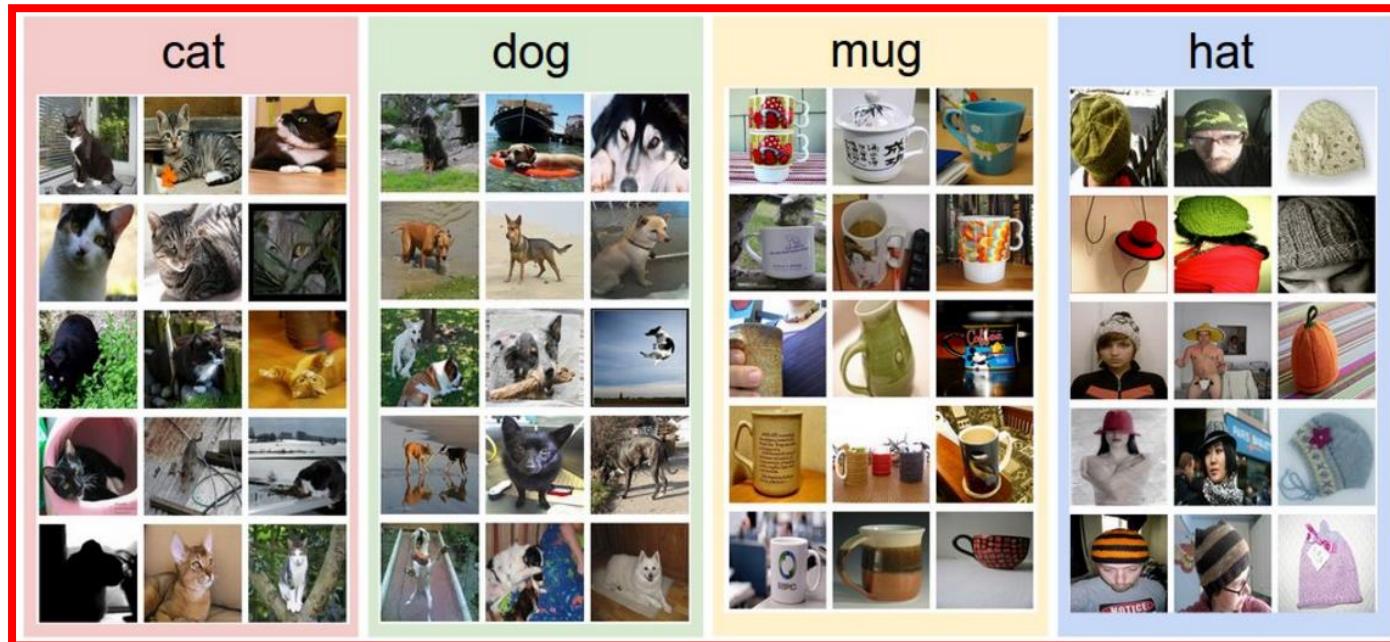


# Philosophical Ambiguity: “Image Classification” is not (yet) “Understanding”



10  
Cats

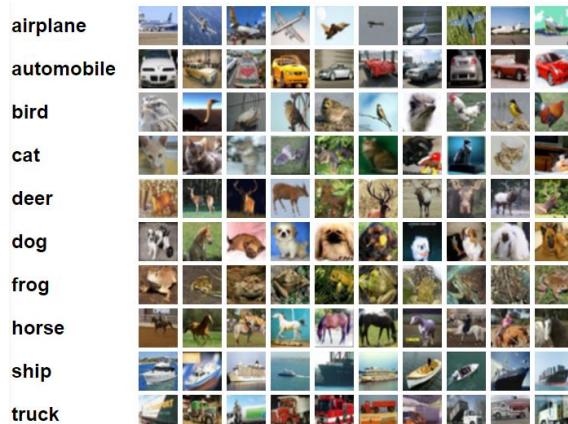
# Image Classification Pipeline



# Famous Computer Vision Datasets

9 3 1 9 7 2 4 5 1 0 3 2 4 3 7 5 9 0 3 4 9  
3 0 2 4 2 9 4 8 3 2 0 1 3 5 3 5 7 4 6 8 5  
2 8 2 3 2 3 8 2 4 9 8 2 9 1 3 9 1 1 1 9 9  
6 3 6 9 0 3 6 0 3 0 1 1 3 9 3 1 5 0 4 9 6  
3 3 8 0 7 0 5 6 9 8 8 4 1 4 4 4 6 9 5 3 3  
1 1 9 5 8 0 4 3 7 7 5 0 5 4 2 0 9 8 1 2 4  
9 5 0 0 5 1 1 1 7 4 7 7 2 6 5 1 8 2 4 1 1  
0 2 1 6 1 7 0 9 5 6 3 2 6 6 7 1 5 2 3 2  
9 4 3 2 1 0 0 2 0 8 7 4 0 9 7 9 3 6 9 3 4  
5 5 1 6 6 2 7 6 7 5 6 6 5 8 1 6 8 7 1 0 5  
7 1 7 5 9 2 3 9 4 3 0 4 5 8 0 0 4 0 4 6 6  
1 6 7 9 6 4 1 1 4 1 3 1 2 3 4 8 1 5 5 0 7  
0 1 6 1 6 7 5 5 5 6 6 8 8 1 7 2 8 3 7 6 5  
6 4 6 8 7 7 1 3 0 7 3 8 6 9 1 6 7 3 6 4 8  
7 9 7 3 1 3 9 7 9 3 6 2 4 9 2 1 4 5 0 3 8

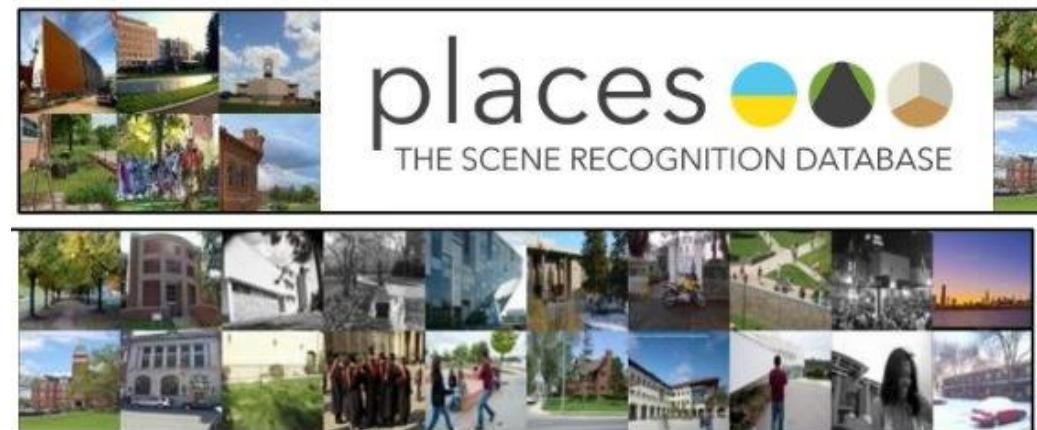
**MNIST:** handwritten digits



**CIFAR-10(0):** tiny images

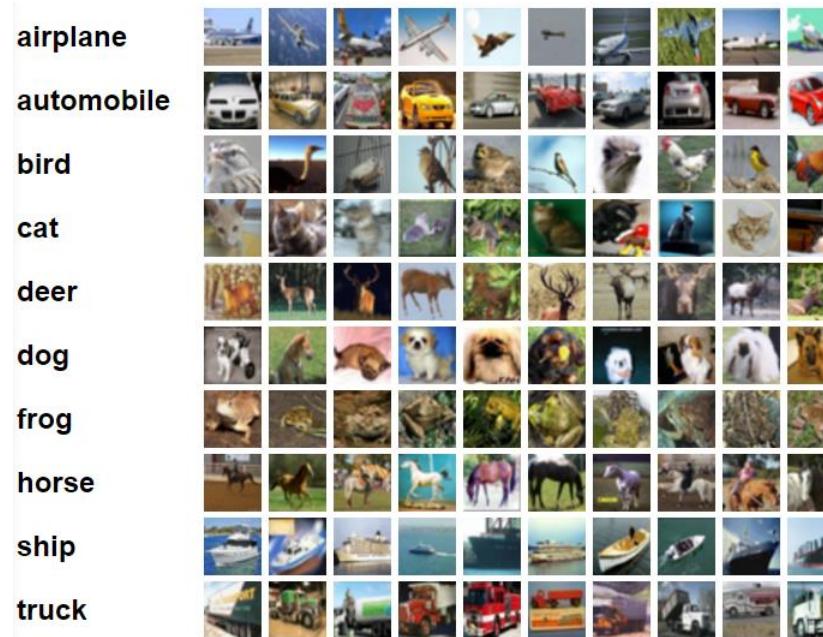
offspring printer housing animal weight drop egg white  
teacher housing computer album garage dorm flower television  
register gallery court key structure light date spread  
horse king fireplace church press market lighter  
restaurant counter cup pack  
hotel road paper side site door coffee  
sport screen coast tree file tower hill can camp fish bathroom  
sky plant wine fox house school stock film  
bread weapon table top man car gun  
cloud cover cover top man car gun  
spring range leash van suite mirror seat fly  
descent fruit dog bed shop people study  
kitchen train kit roll goal  
engine camera memoriesieve cell kid center  
chain stone boat tea overall sleeve step  
dinner home room office bar watch  
apple girl flat rule hall  
flag bank cross chair mine castle ocean  
radio valley t-shirt club  
beach support level line street golf  
base library stage video food building  
tool material player leg shirt desk vehicle  
football hospital match equipment cell phone mountain  
short circuit bridge scale gas pedal microphone recording crowd telephone

**ImageNet:** WordNet hierarchy



**Places:** natural scenes

# Let's Build an Image Classifier for CIFAR-10



test image				training image				pixel-wise absolute value differences			
56	32	10	18	10	20	24	17	46	12	14	1
90	23	128	133	8	10	89	100	82	13	39	33
24	26	178	200	12	16	178	170	12	10	0	30
2	0	255	220	4	32	233	112	2	32	22	108

→ 456

# Let's Build an Image Classifier for CIFAR-10

$$\begin{array}{c} \text{test image} \\ \left| \begin{array}{cccc} 56 & 32 & 10 & 18 \\ 90 & 23 & 128 & 133 \\ 24 & 26 & 178 & 200 \\ 2 & 0 & 255 & 220 \end{array} \right| - \begin{array}{cccc} 10 & 20 & 24 & 17 \\ 8 & 10 & 89 & 100 \\ 12 & 16 & 178 & 170 \\ 4 & 32 & 233 & 112 \end{array} = \begin{array}{cccc} 46 & 12 & 14 & 1 \\ 82 & 13 & 39 & 33 \\ 12 & 10 & 0 & 30 \\ 2 & 32 & 22 & 108 \end{array} \rightarrow 456 \end{array}$$



## Accuracy

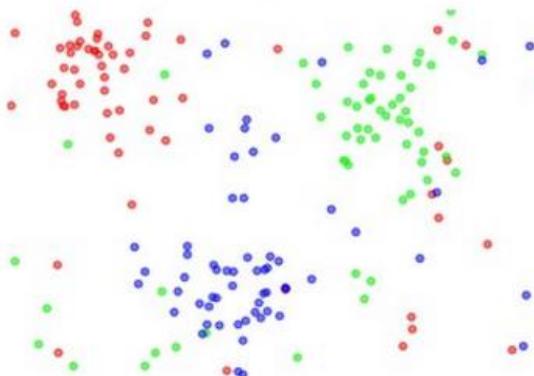
Random: **10%**

Our image-diff (with L1): **38.6%**

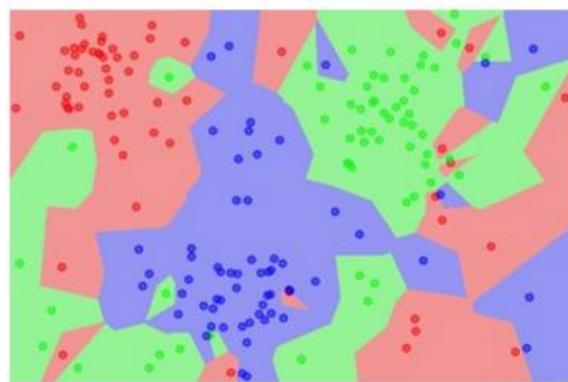
Our image-diff (with L2): **35.4%**

# K-Nearest Neighbors: Generalizing the Image-Diff Classifier

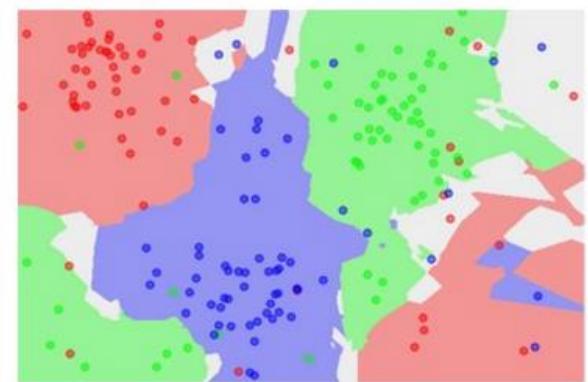
the data



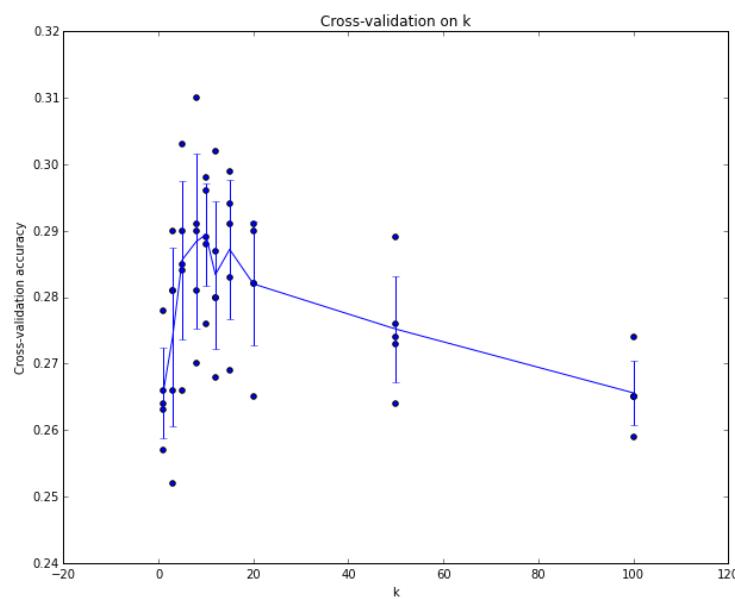
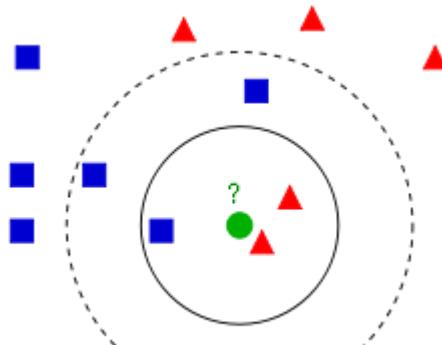
NN classifier



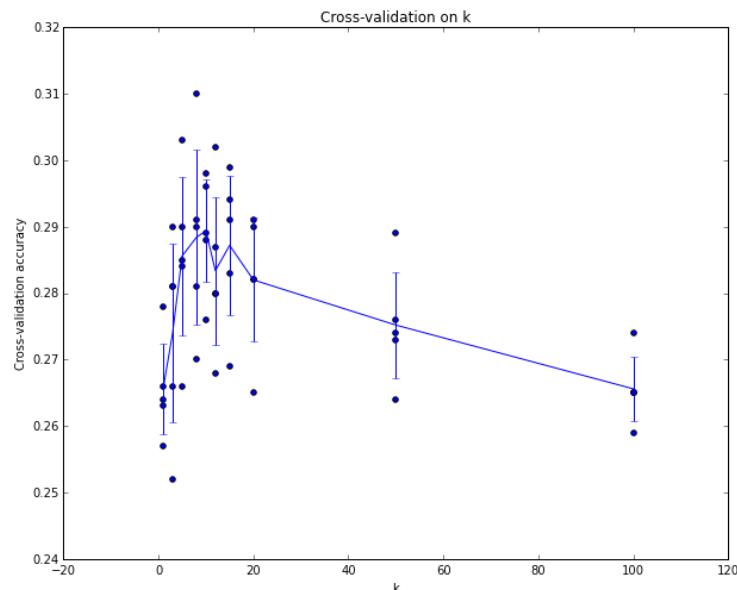
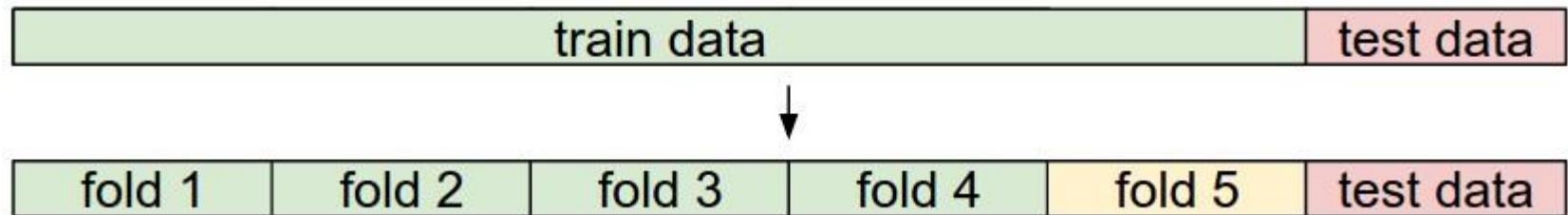
5-NN classifier



Tuning (hyper)parameters:



# K-Nearest Neighbors: Generalizing the Image-Diff Classifier



## Accuracy

Random: **10%**

Training and testing on the same data: **35.4%**

7-Nearest Neighbors: **~30%**

Human: **~95%**

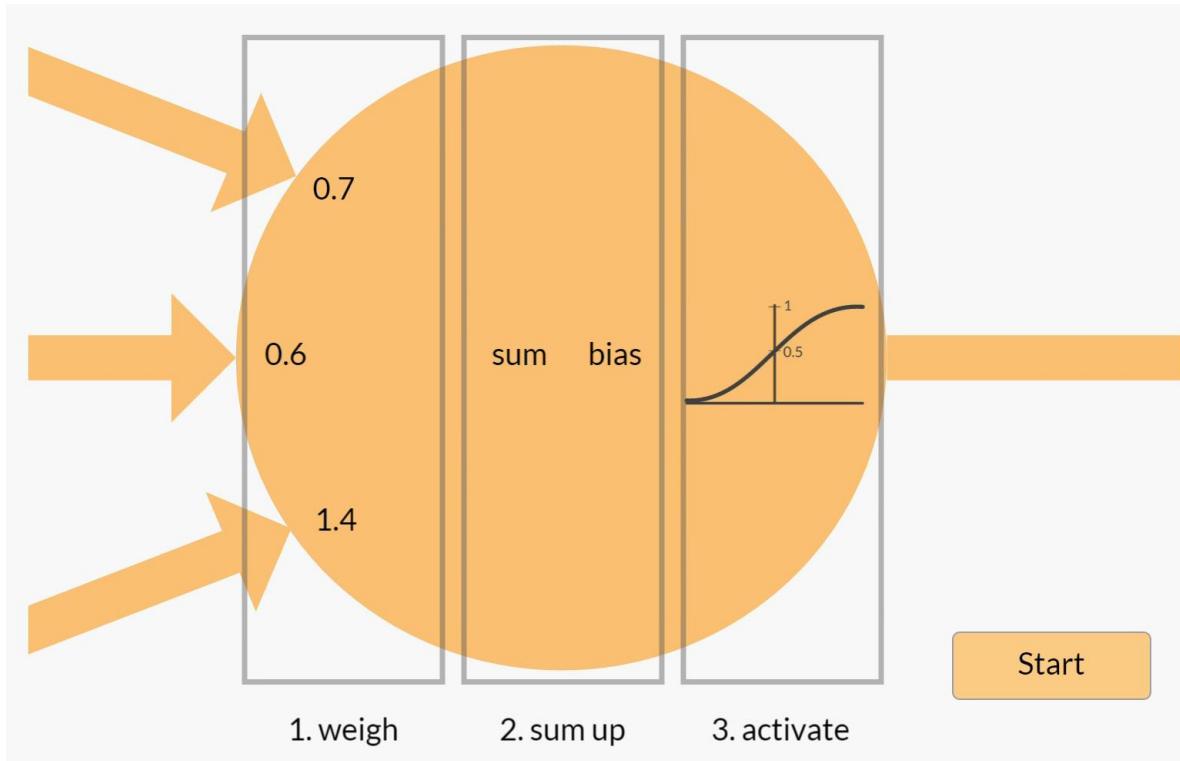
...

Convolutional Neural Networks: **~97.75%**

# Reminder: Weighing the Evidence

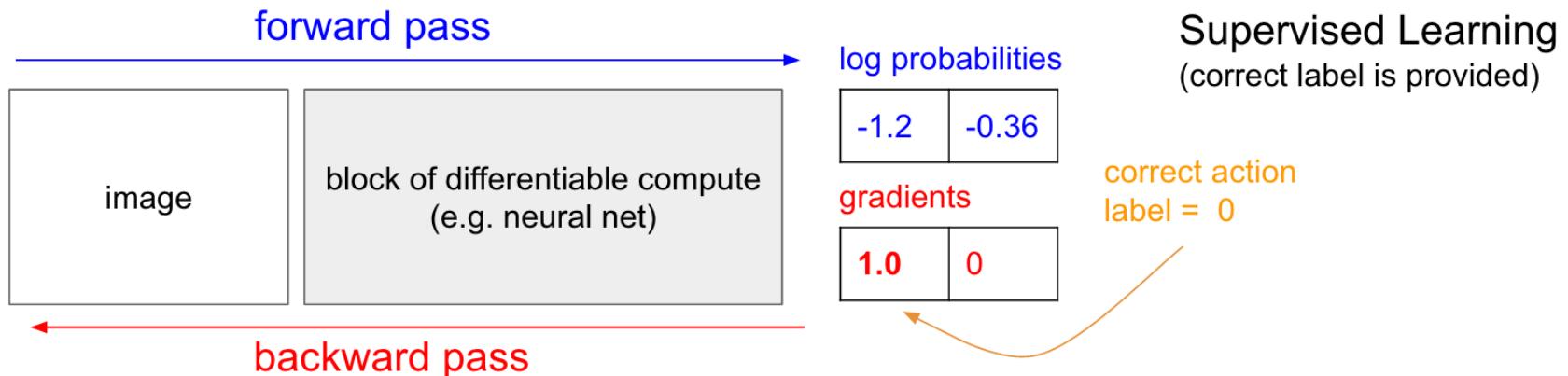
Evidence

Decisions



$$\text{output} = \begin{cases} 0 & \text{if } \sum_j w_j x_j \leq \text{threshold} \\ 1 & \text{if } \sum_j w_j x_j > \text{threshold} \end{cases}$$

# Reminder: “Learning” is Optimization of a Function

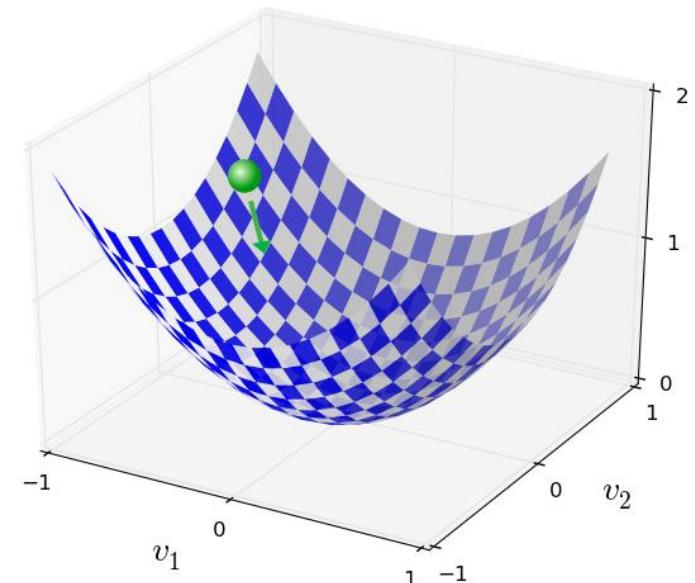


Ground truth for “6”:

$$y(x) = (0, 0, 0, 0, 0, 0, 1, 0, 0, 0)^T$$

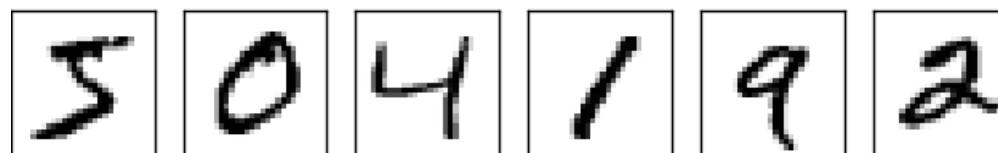
“Loss” function:

$$C(w, b) \equiv \frac{1}{2n} \sum_x \|y(x) - a\|^2$$

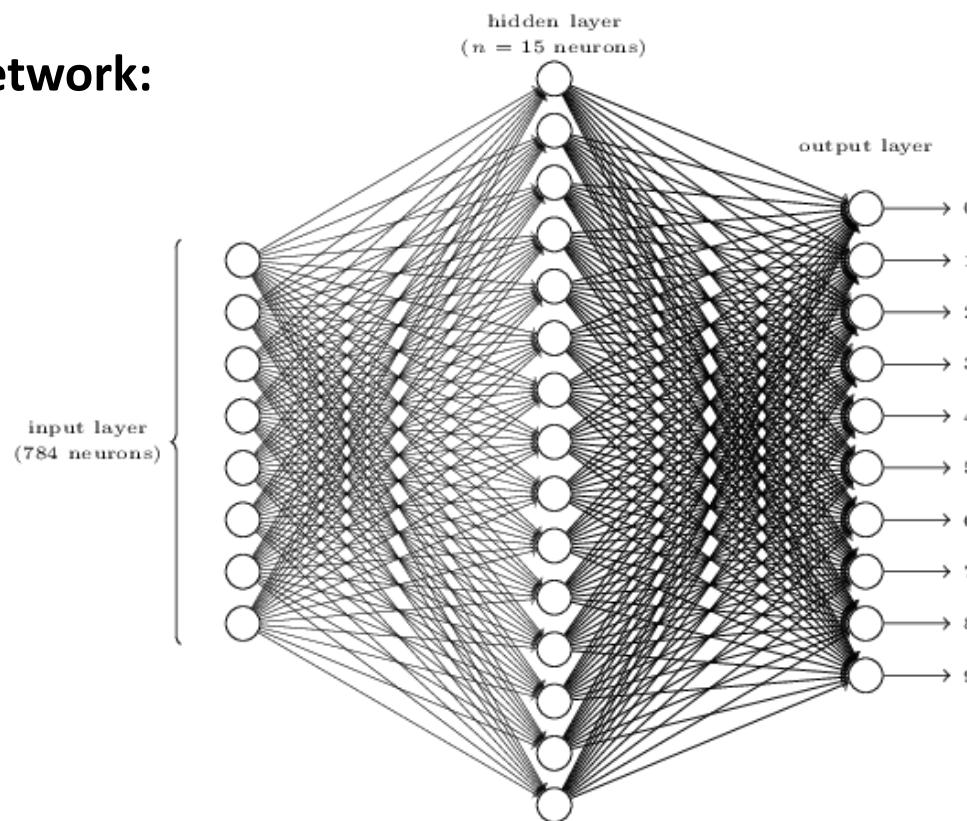


# Classify and Image of a Number

**Input:**  
(28x28)

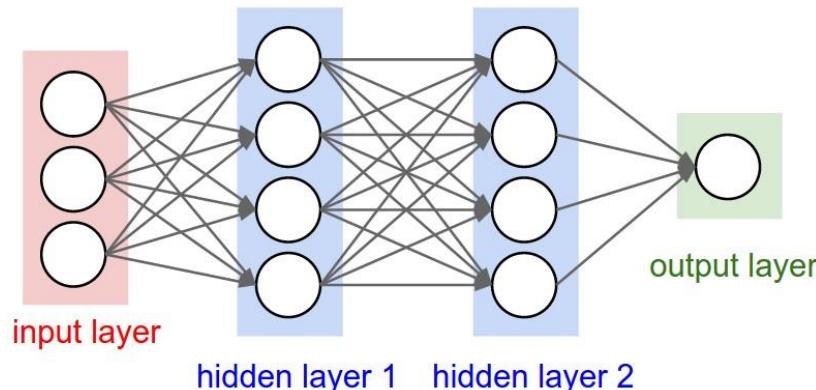


**Network:**

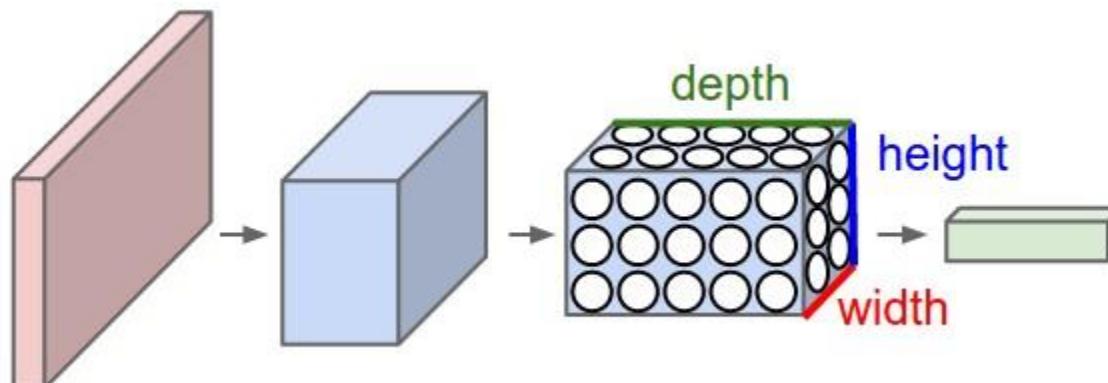


# Convolutional Neural Networks

Regular neural network (fully connected):



Convolutional neural network:

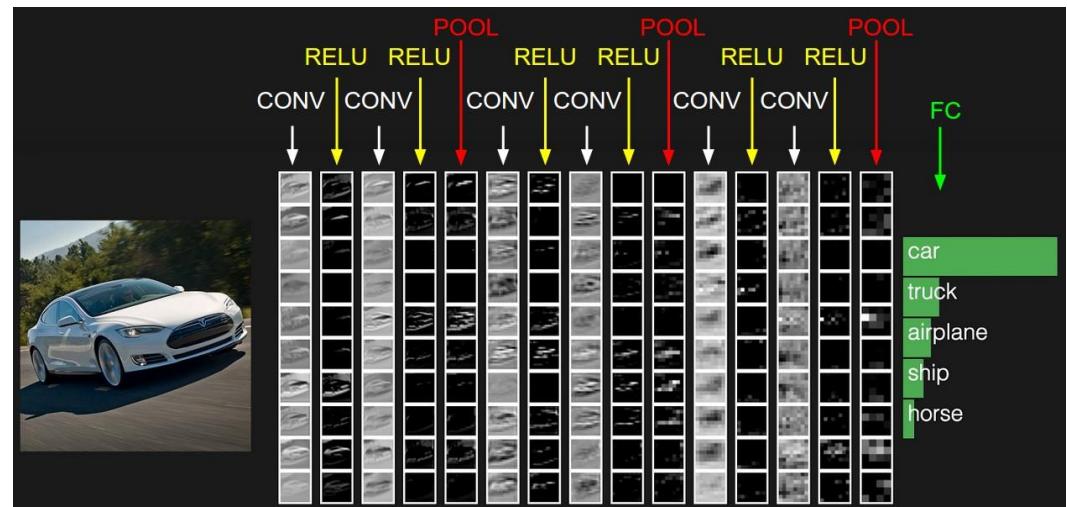


Each layer takes a 3d volume, produces 3d volume with some smooth function that may or may not have parameters.

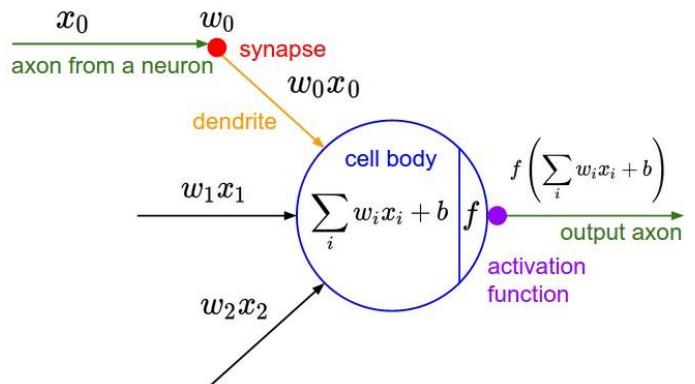
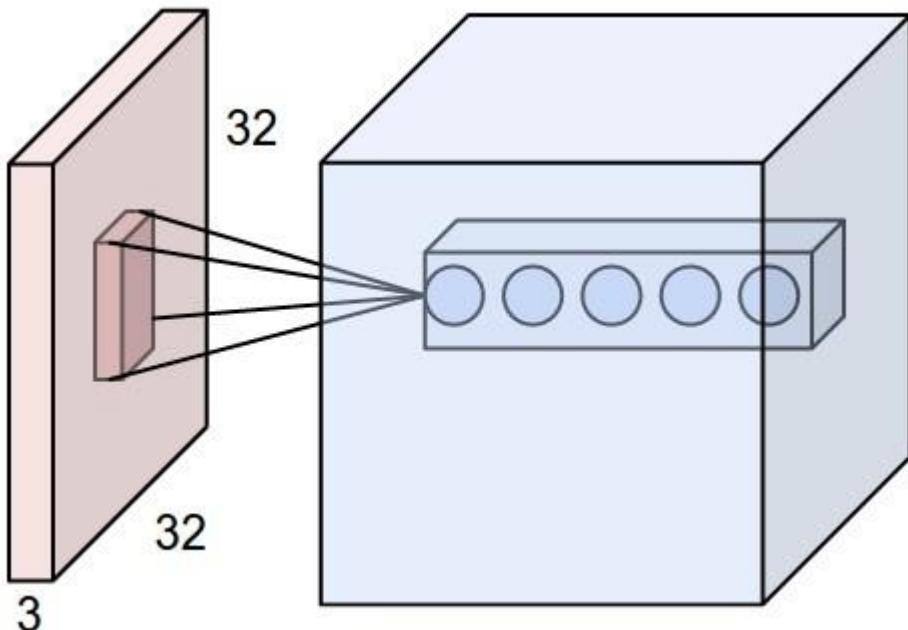
# Convolutional Neural Networks: Layers

- **INPUT** [32x32x3] will hold the raw pixel values of the image, in this case an image of width 32, height 32, and with three color channels R,G,B.
- **CONV** layer will compute the output of neurons that are connected to local regions in the input, each computing a dot product between their weights and a small region they are connected to in the input volume. This may result in volume such as [32x32x12] if we decided to use 12 filters.
- **RELU** layer will apply an elementwise activation function, such as the  $\max(0,x)$  thresholding at zero. This leaves the size of the volume unchanged ([32x32x12]).
- **POOL** layer will perform a downsampling operation along the spatial dimensions (width, height), resulting in volume such as [16x16x12].
- **FC** (i.e. fully-connected) layer will compute the class scores, resulting in volume of size [1x1x10], where each of the 10 numbers correspond to a class score, such as among the 10 categories of CIFAR-10. As with ordinary Neural Networks and as the name implies, each neuron in this layer will be connected to all the numbers in the previous volume.

Layers highlighted in blue have learnable parameters.



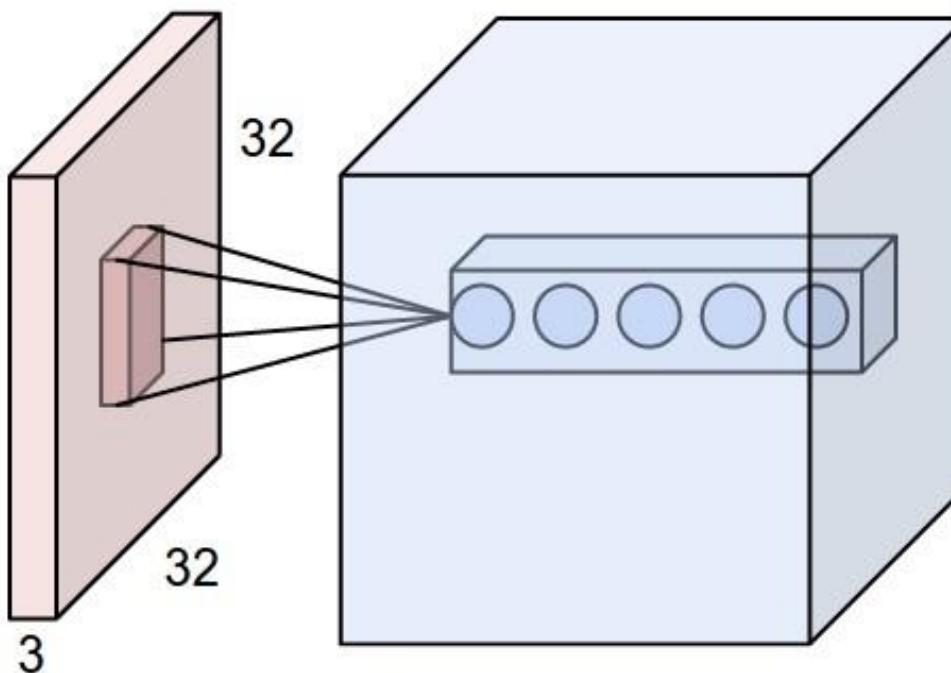
# Dealing with Images: Local Connectivity



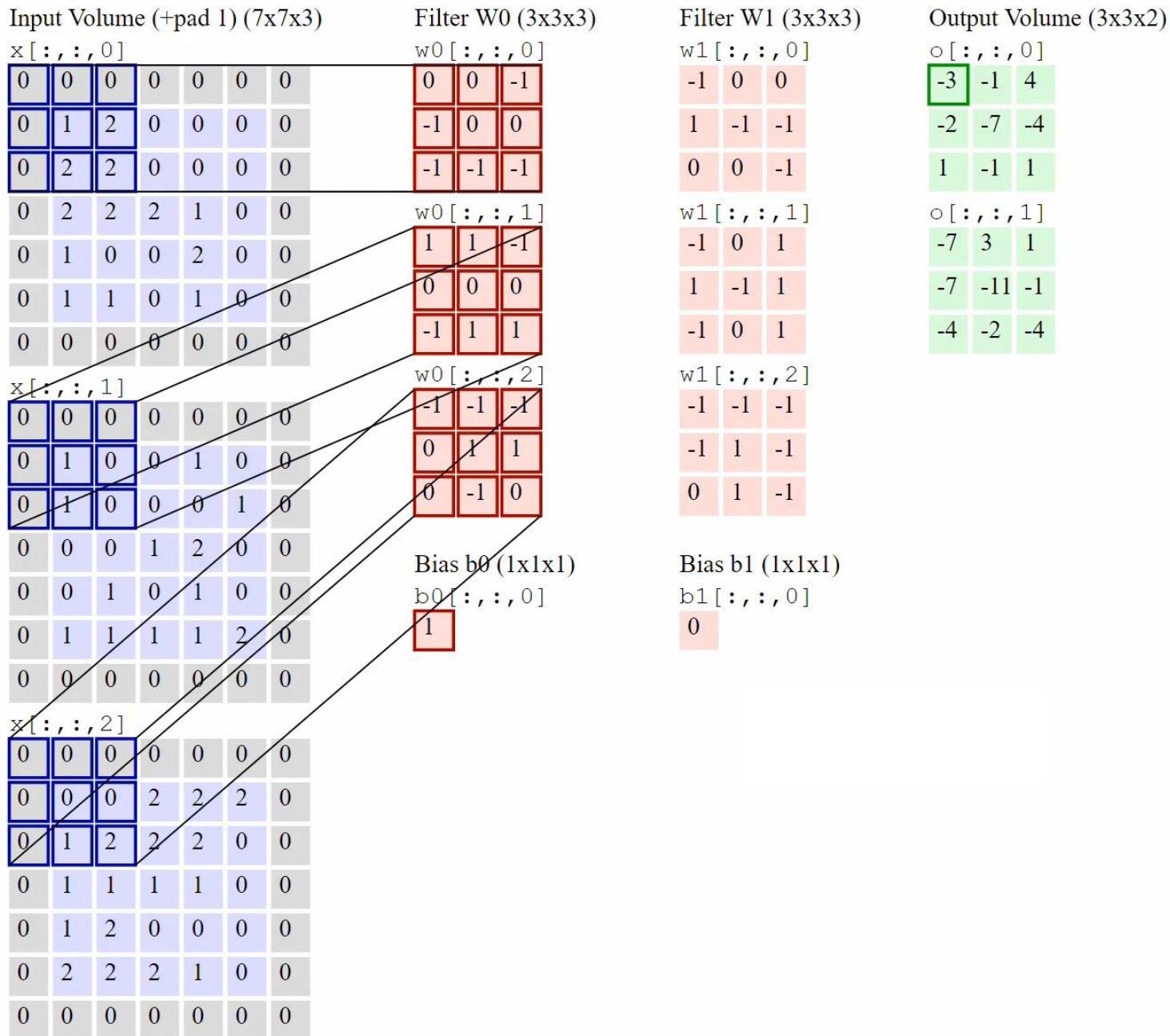
Same neuron. Just more focused (narrow “receptive field”).

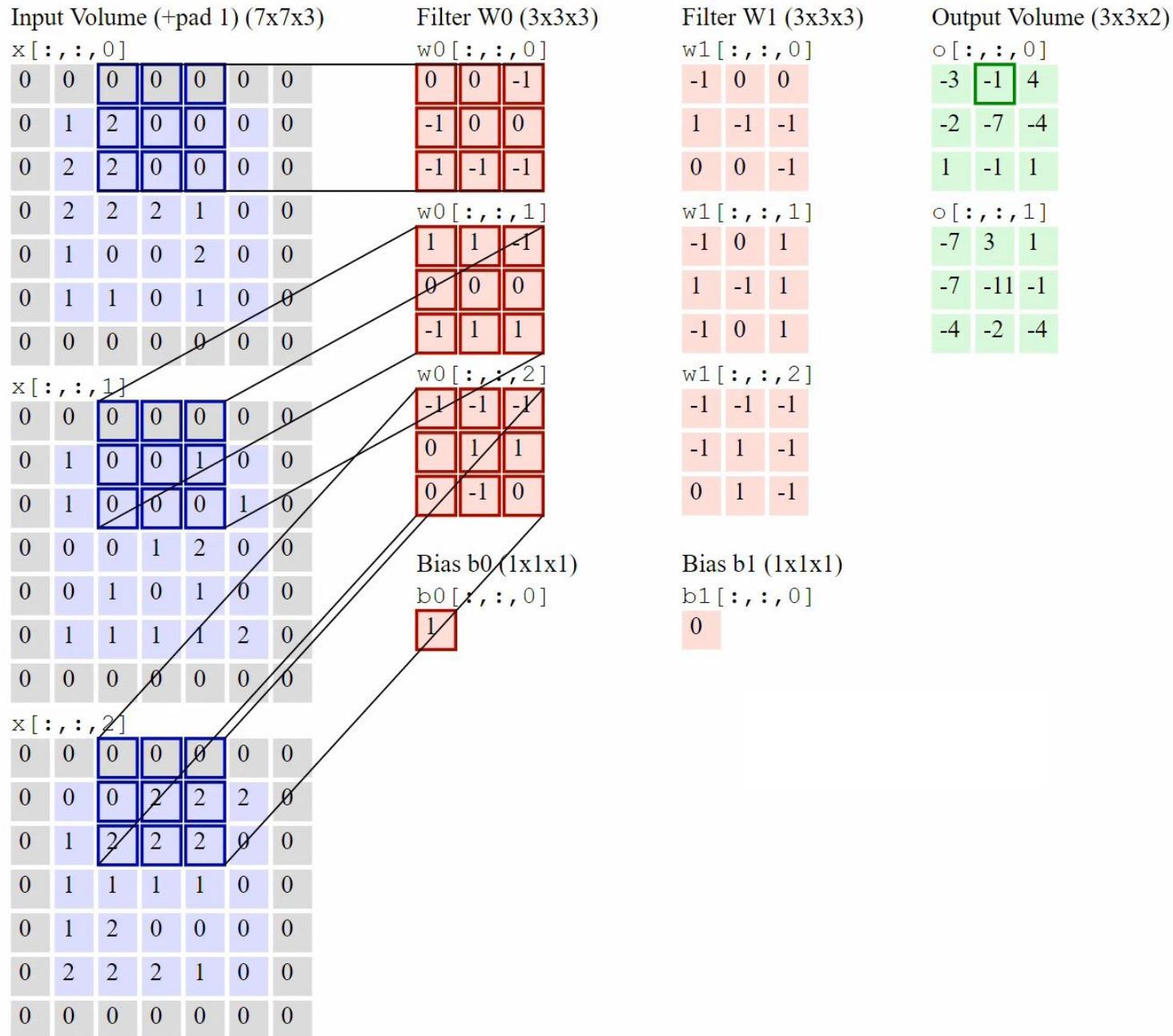
The parameters on each filter are spatially “shared”  
(if a feature is useful in one place, it’s useful elsewhere)

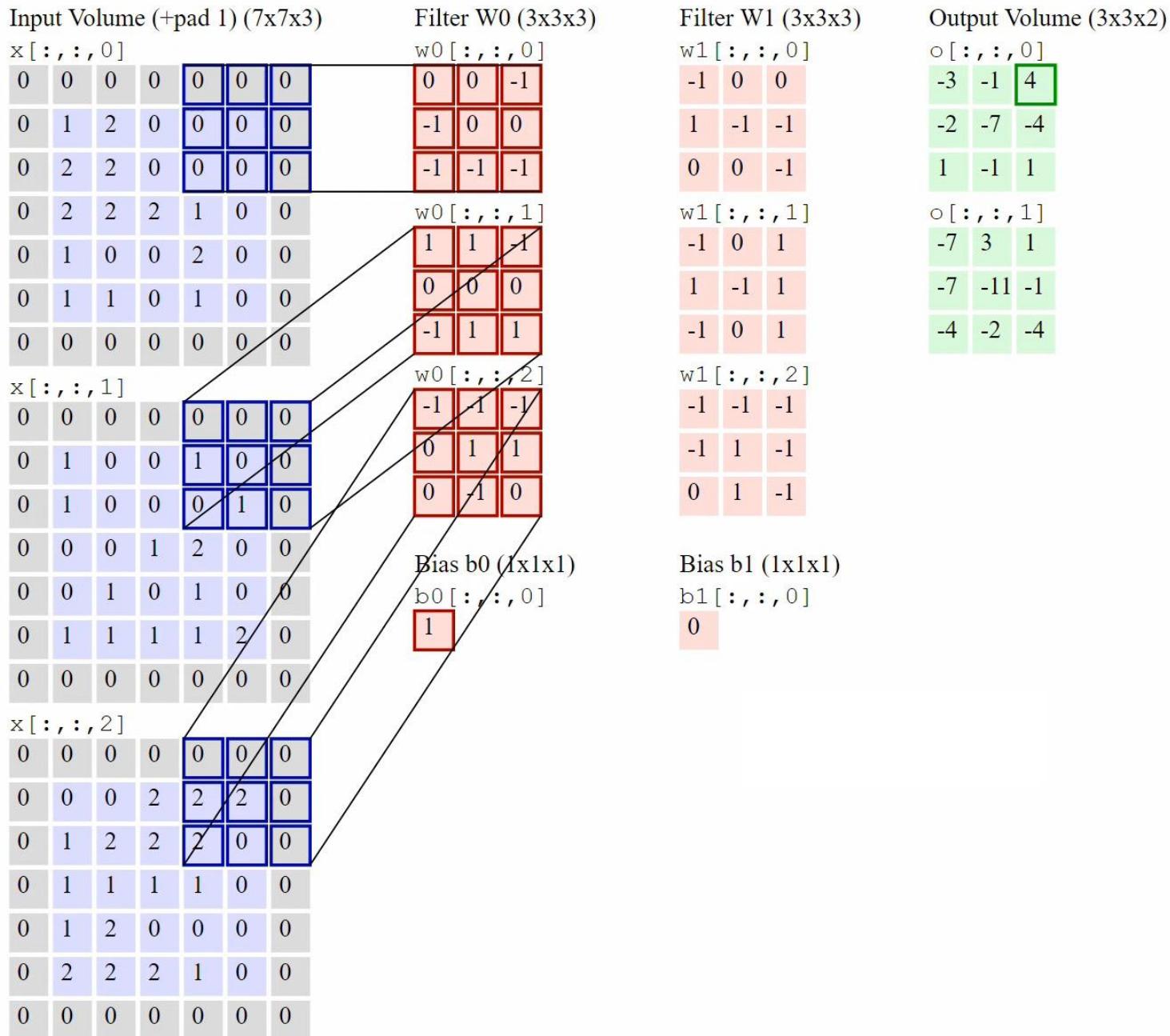
# ConvNets: Spatial Arrangement of Output Volume

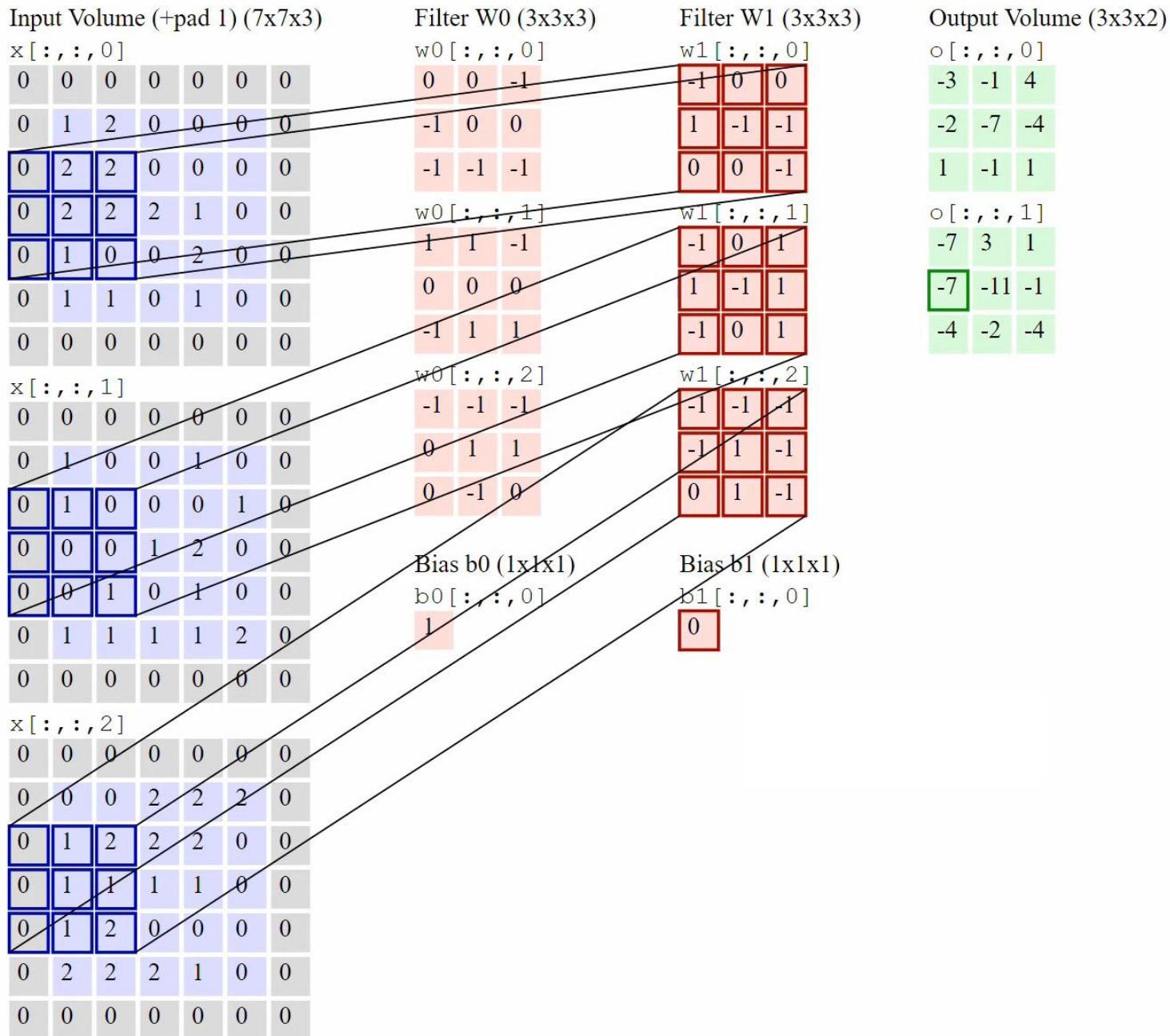


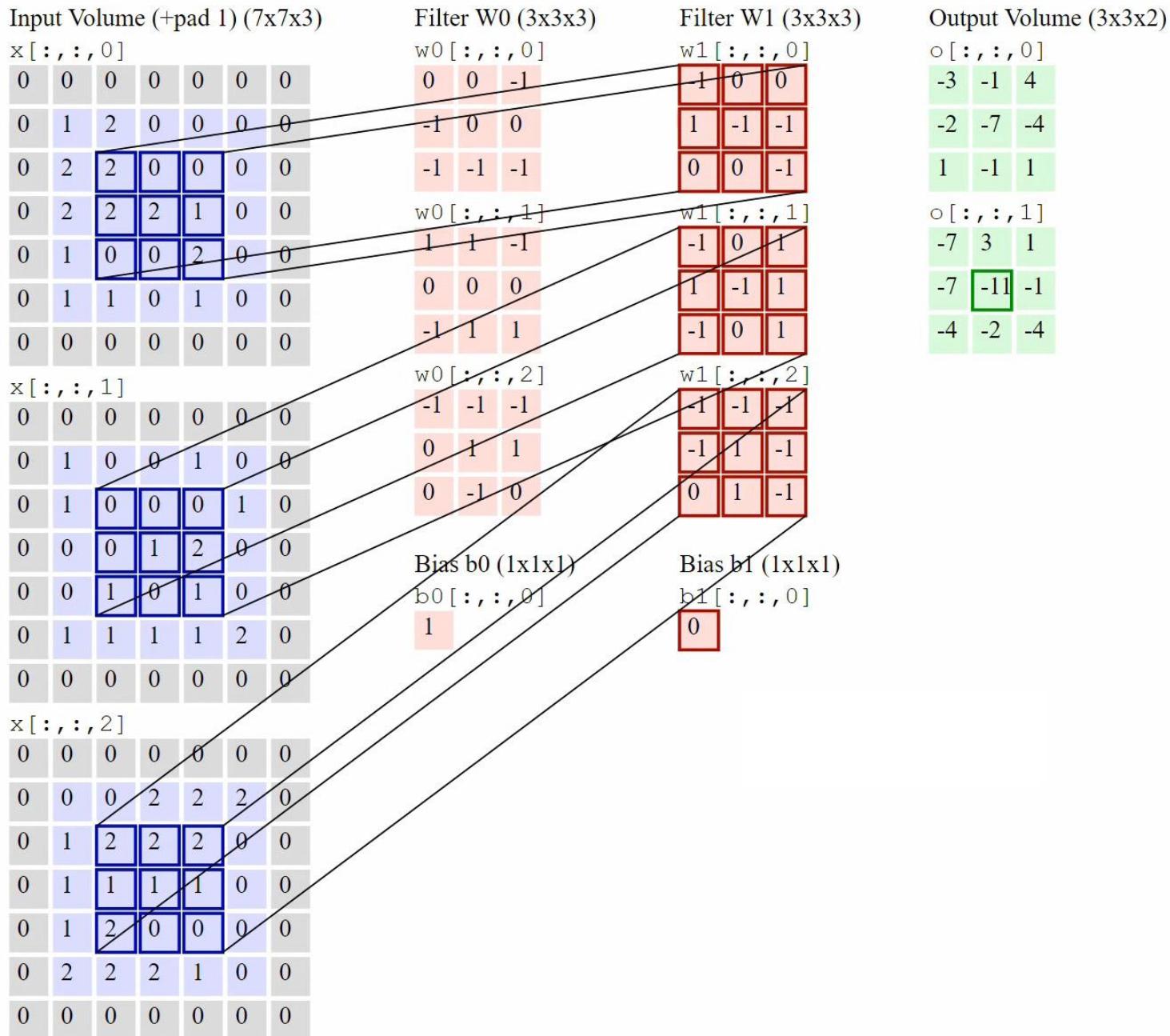
- **Depth:** number of filters
- **Stride:** filter step size (when we “slide” it)
- **Padding:** zero-pad the input

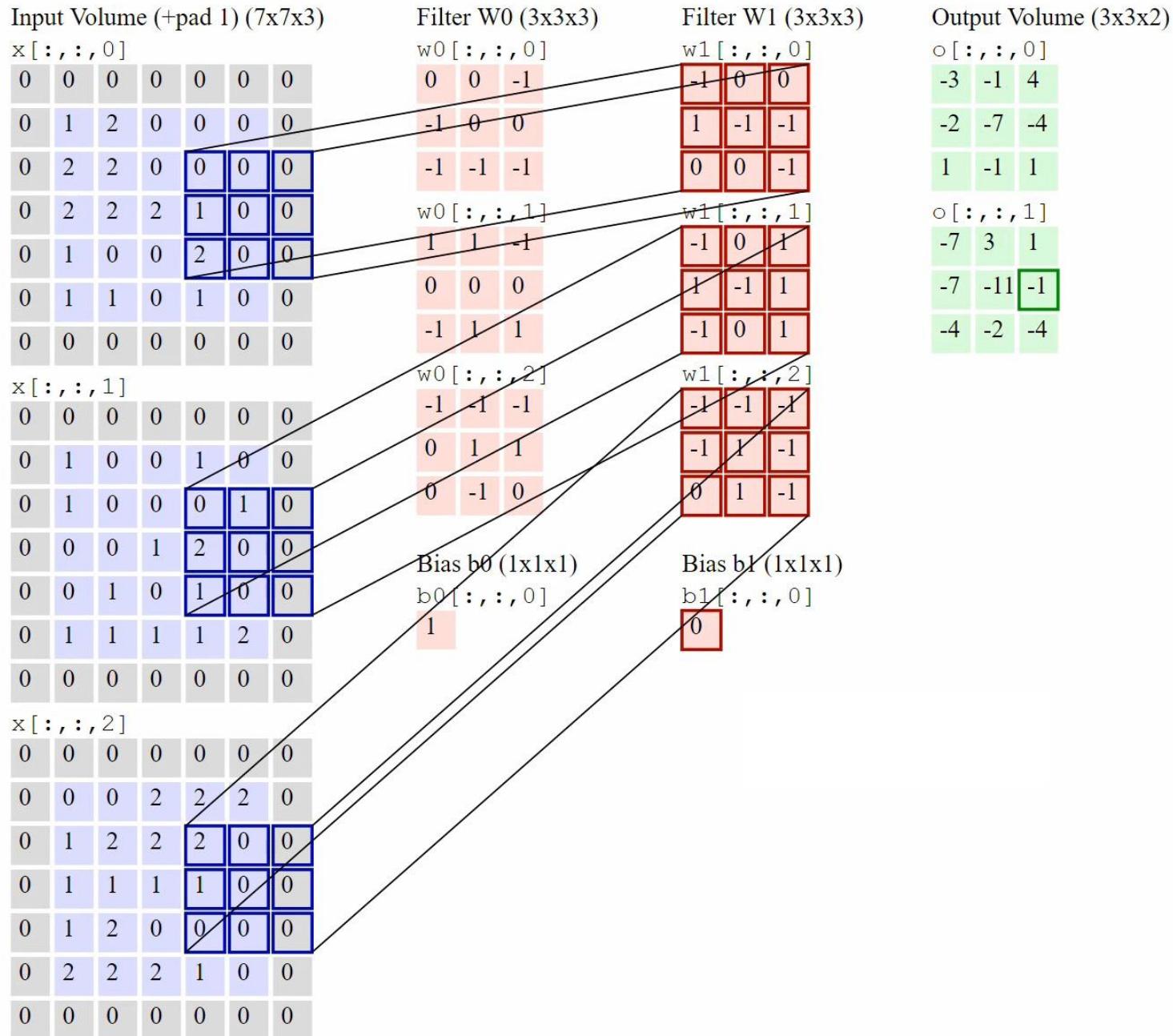


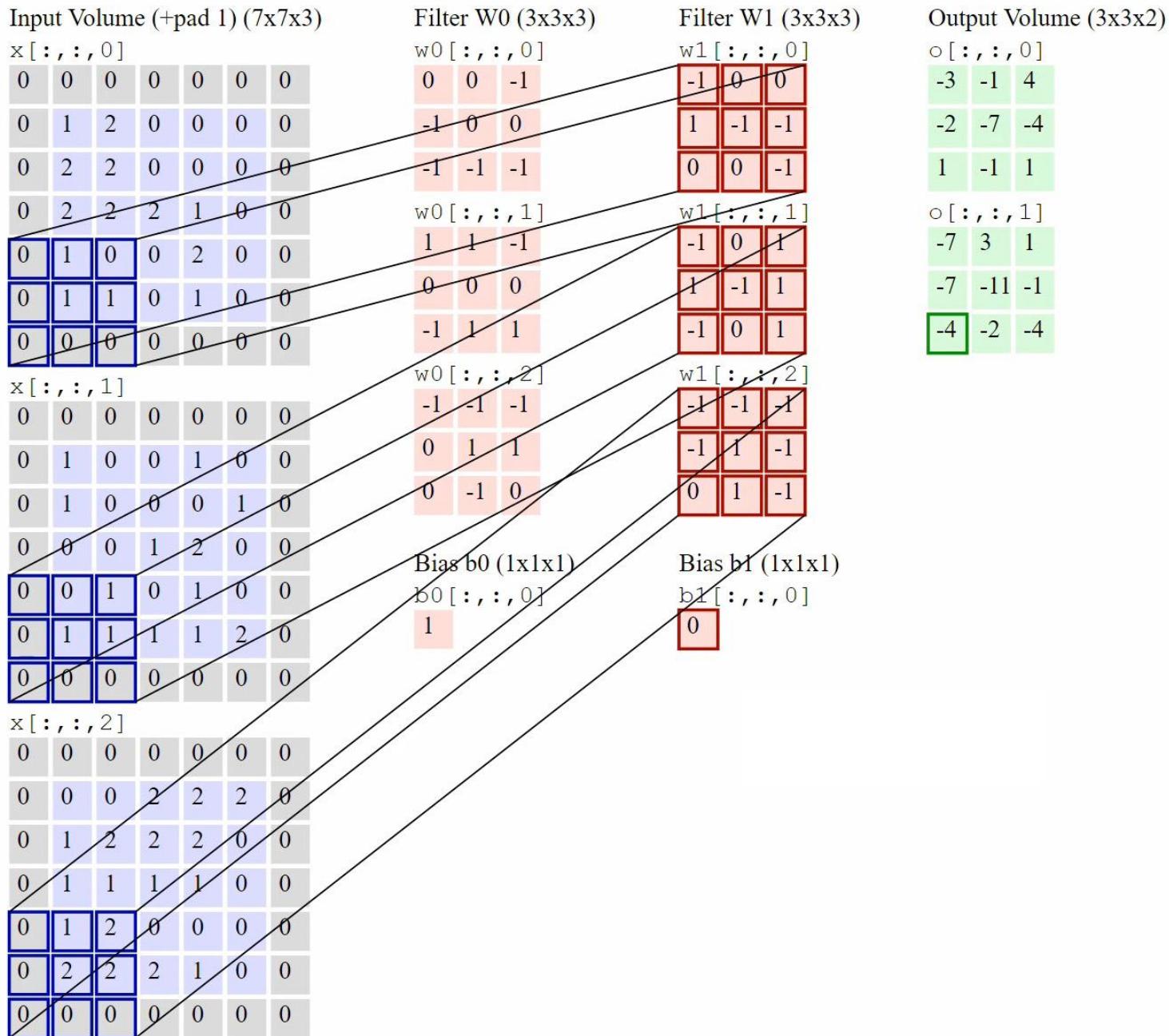


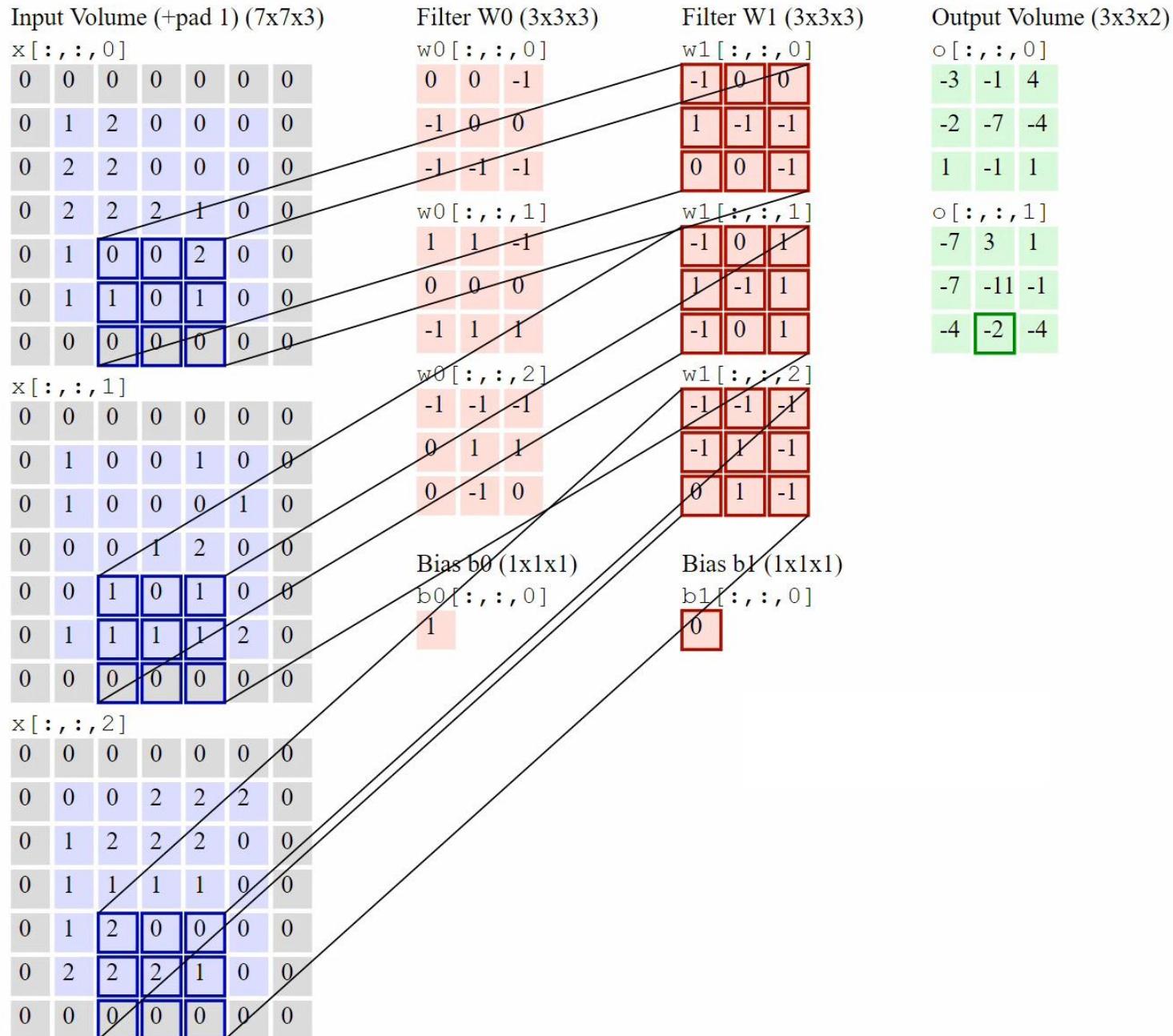


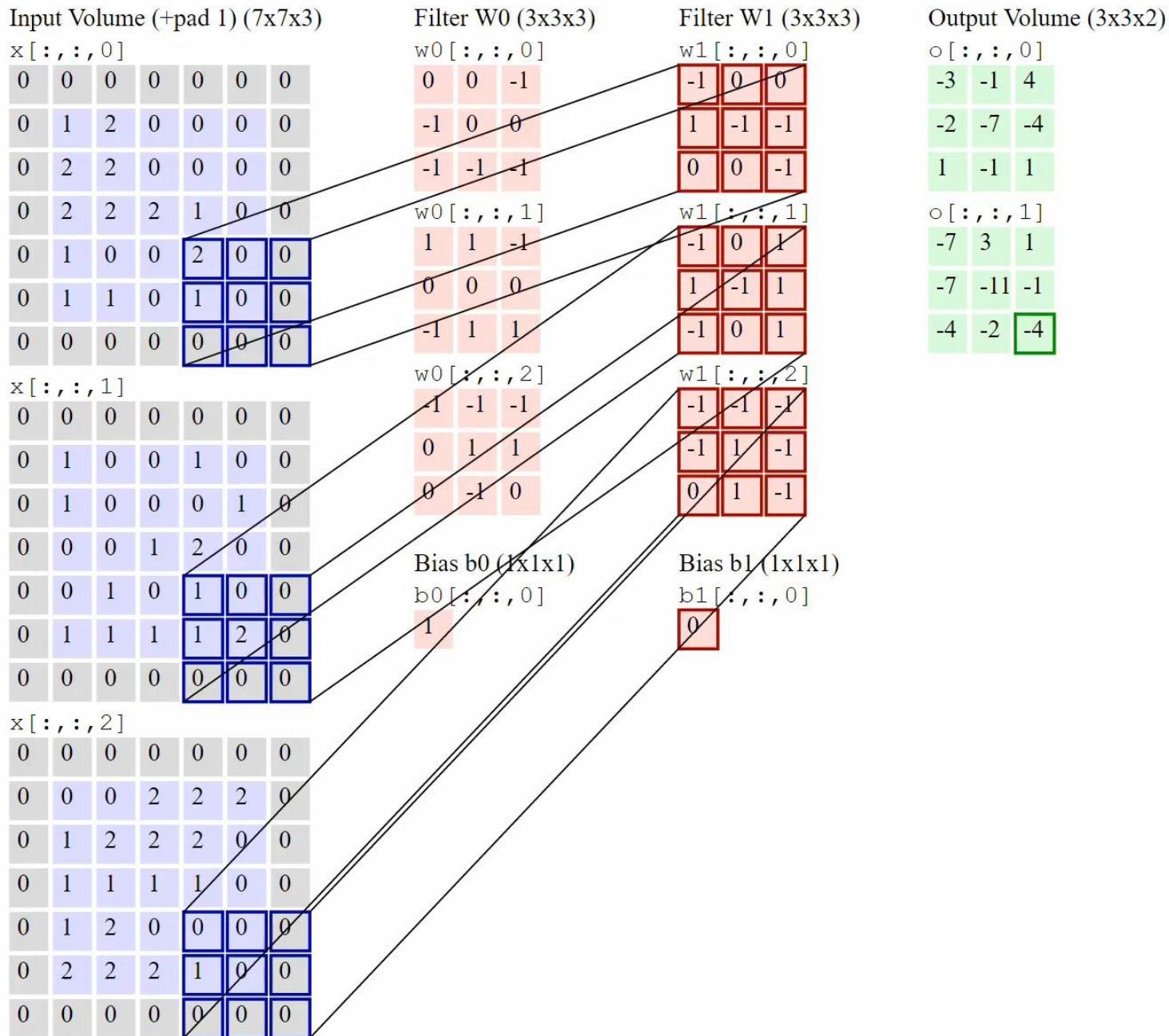












# Convolution

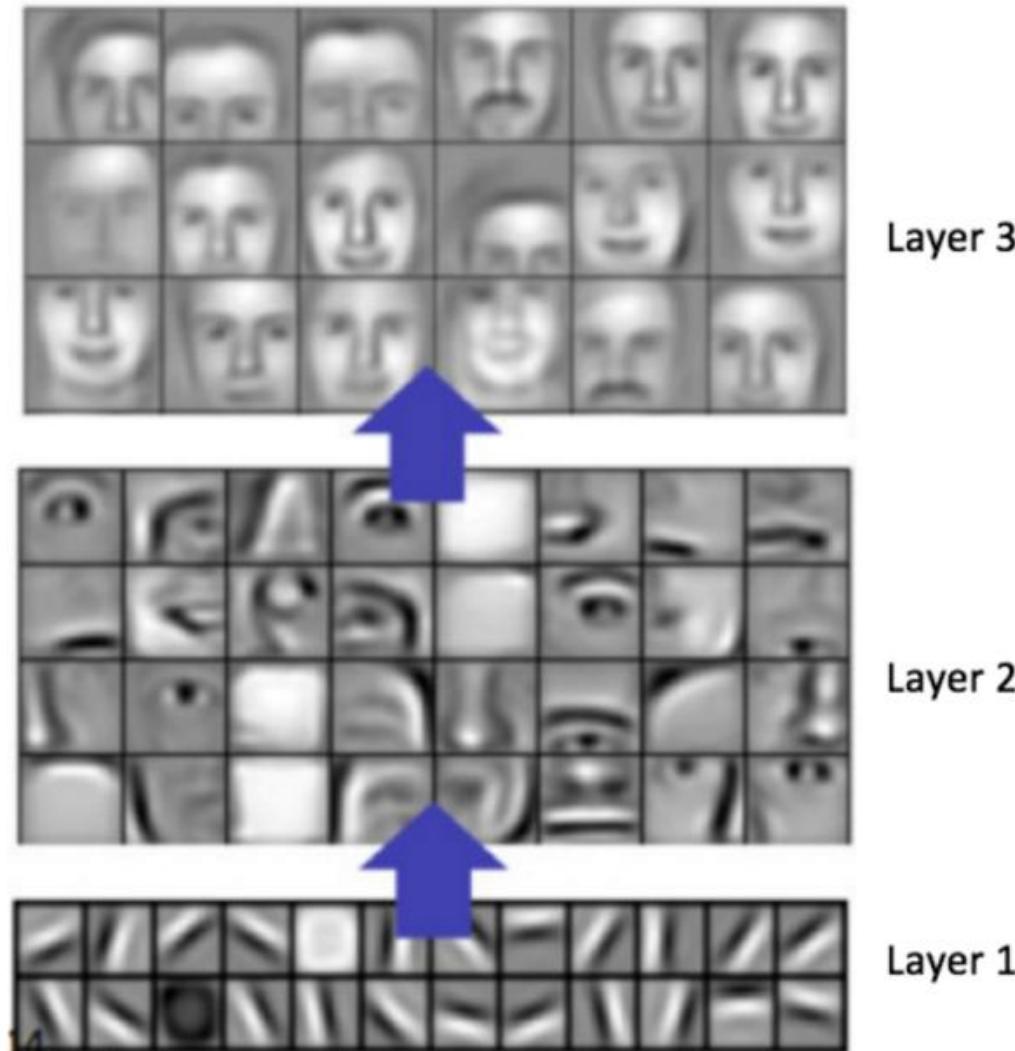
Operation	Filter	Convolved Image
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	

# Convolution

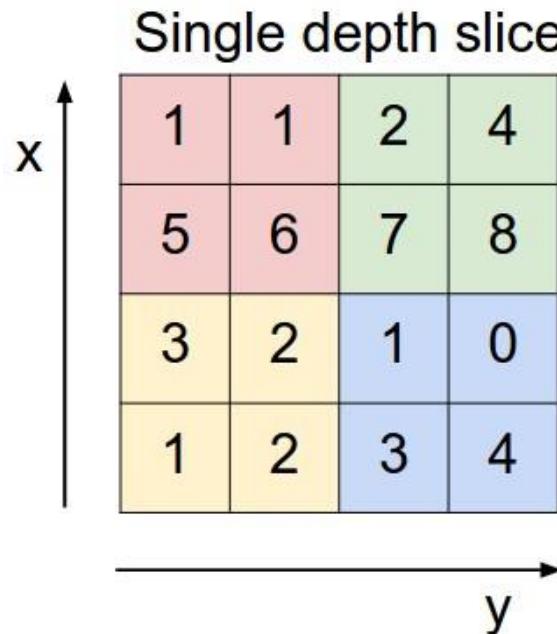


Input

# Convolution: Representation Learning



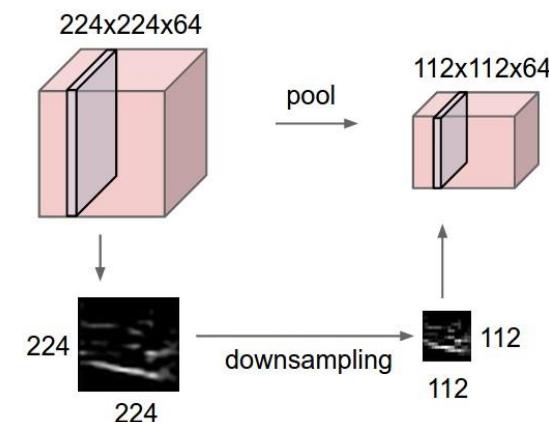
# ConvNets: Pooling



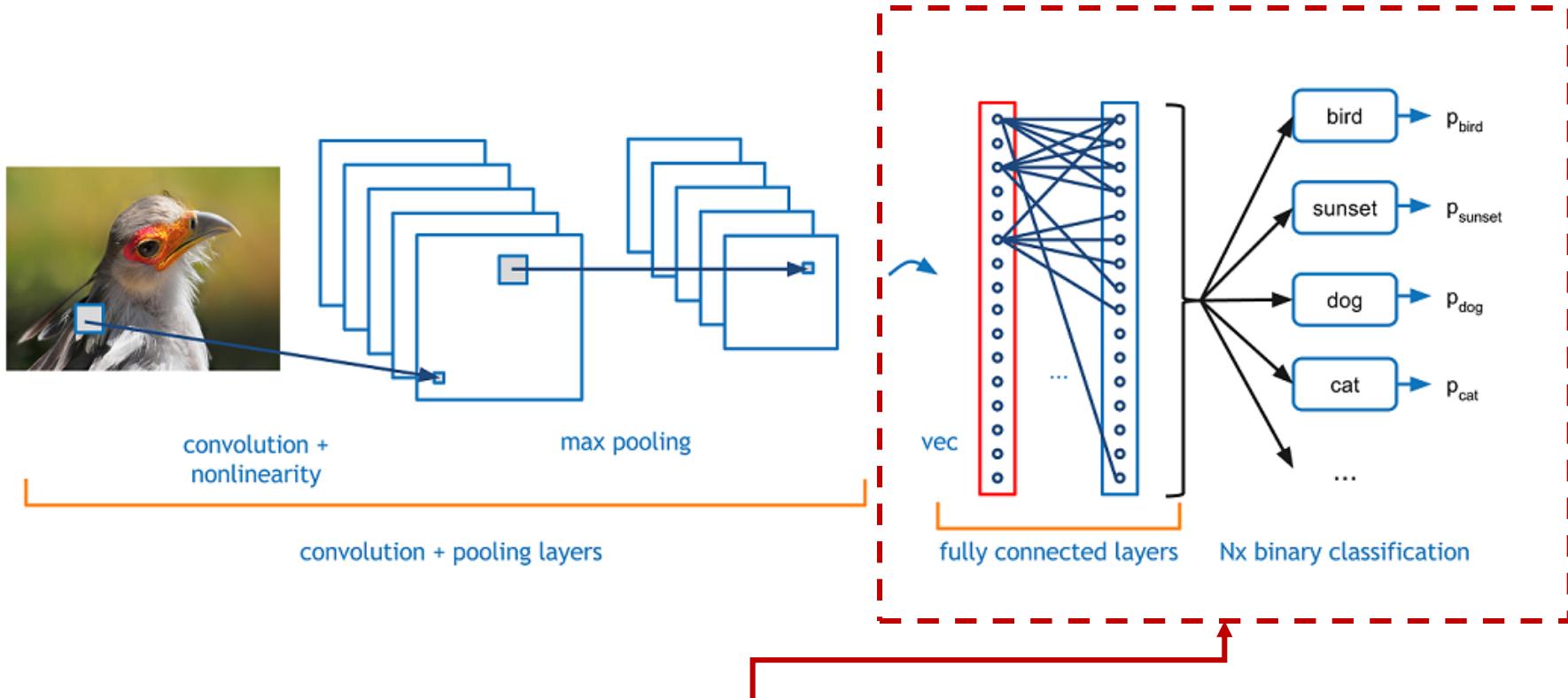
max pool with 2x2 filters  
and stride 2

The result of applying a 2x2 max pooling filter with stride 2 to the input grid. The resulting 2x2 grid contains the maximum values from each 2x2 receptive field: the top-left cell contains 6 (from the pink cells), the top-right cell contains 8 (from the light green cells), the bottom-left cell contains 3 (from the yellow cells), and the bottom-right cell contains 4 (from the light blue cells).

6	8
3	4



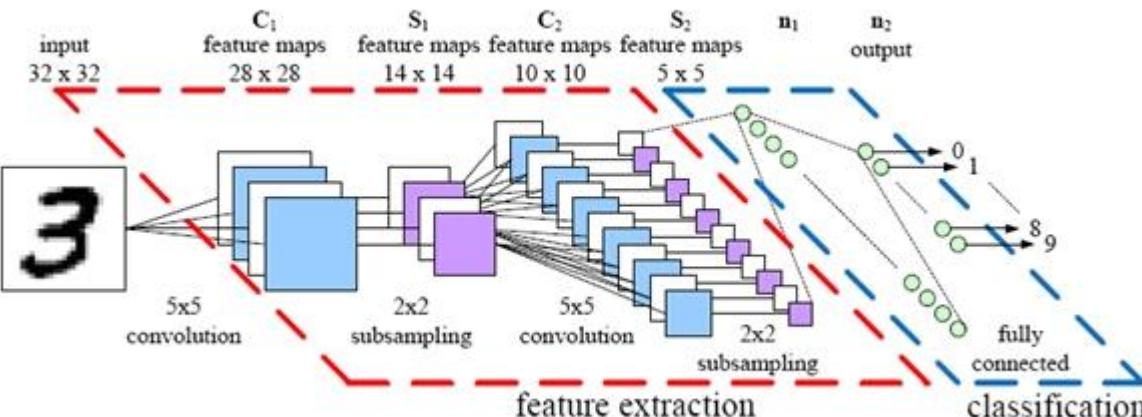
# Same Architecture, Many Applications



This part might look different for:

- Different image classification **domains**
- Image captioning with **recurrent neural networks**
- Image object localization with **bounding box**
- Image segmentation with **fully convolutional networks**
- Image segmentation with **deconvolution layers**

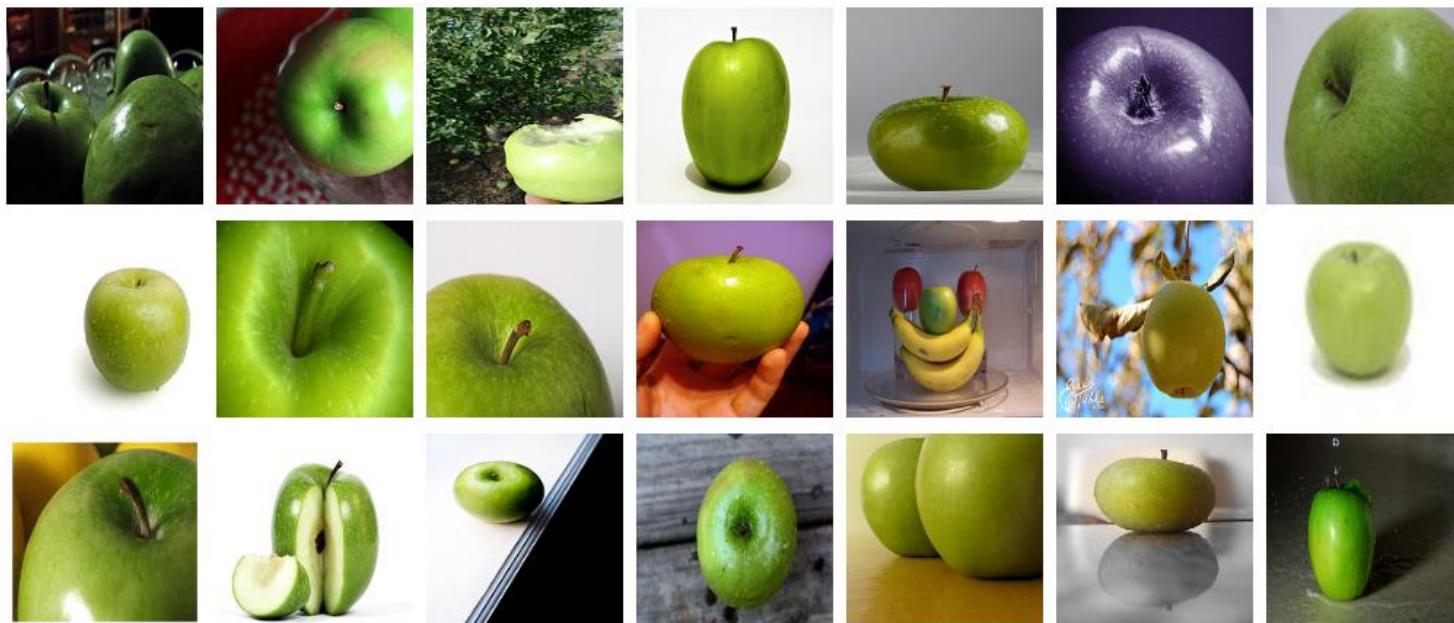
# Object Recognition Case Study: ImageNet



mite	container ship	motor scooter	leopard
mite	container ship	motor scooter	leopard
black widow	lifeboat	go-kart	jaguar
cockroach	amphibian	moped	cheetah
tick	fireboat	bumper car	snow leopard
starfish	drilling platform	golfcart	Egyptian cat

# What is ImageNet?

- **ImageNet:** dataset of 14+ million images (21,841 categories)
- Let's take the high level category of **fruit** as an example:
  - Total 188,000 images of fruit
  - There are 1206 Granny Smith apples:

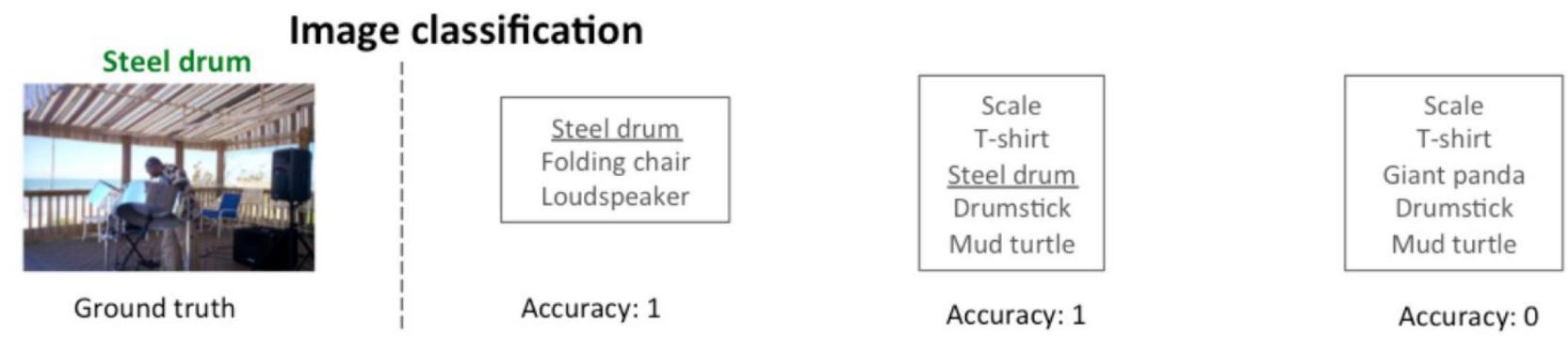


# What is ImageNet?

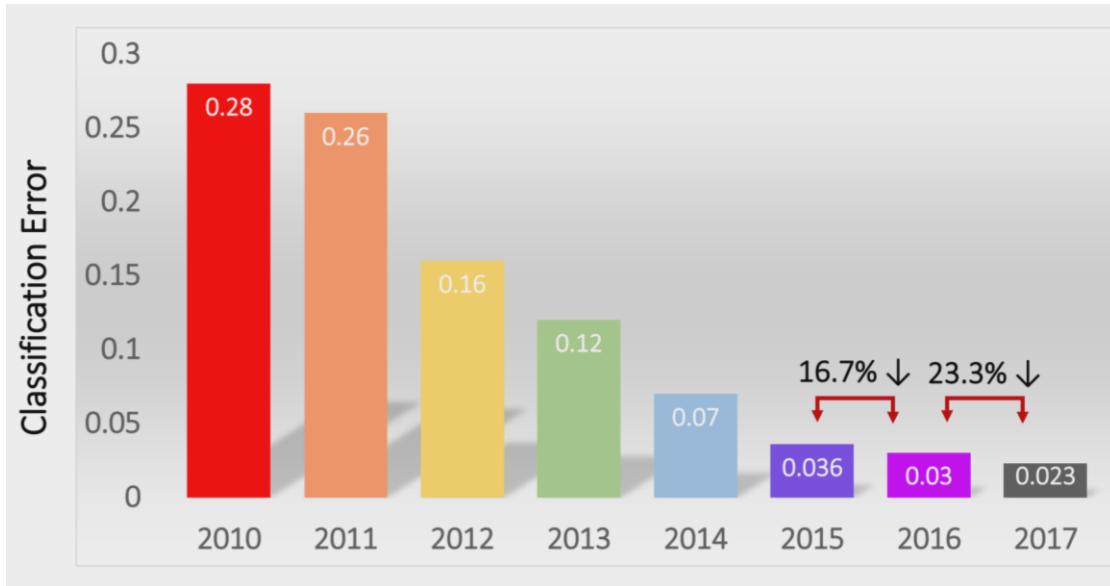
- Dataset → • **ImageNet**: dataset of 14+ million images
- Competition → • **ILSVRC**: ImageNet Large Scale Visual Recognition Challenge
- Networks → • AlexNet (2012)  
• ZFNet (2013)  
• VGGNet (2014)  
• GoogLeNet (2014)  
• ResNet (2015)  
• CUIImage (2016)  
• SENet (2017)

# ILSVRC Challenge Evaluation for Classification

- Top 5 error rate:
  - You get 5 guesses to get the correct label



- ~20% reduction in accuracy for Top 1 vs Top 5
- Human annotation is a binary task: “apple” or “not apple”

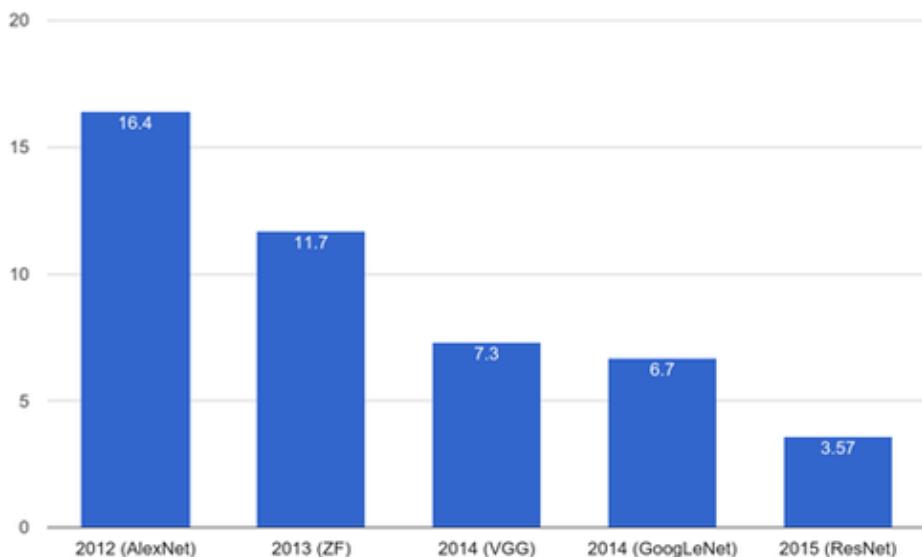


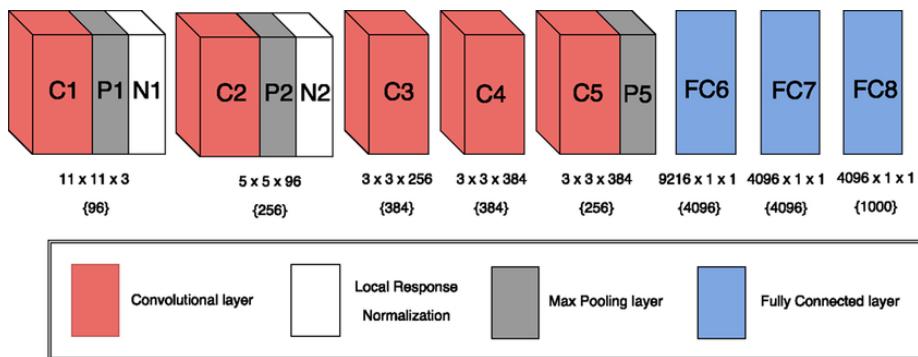
- Human error: 5.1%
  - Surpassed in 2015

- **AlexNet (2012): First CNN (15.4%)**
  - 8 layers
  - 61 million parameters
- **ZFNet (2013): 15.4% to 11.2%**
  - 8 layers
  - More filters. Denser stride.
- **VGGNet (2014): 11.2% to 7.3%**
  - Beautifully uniform:  
3x3 conv, stride 1, pad 1, 2x2 max pool
  - 16 layers
  - 138 million parameters
- **GoogLeNet (2014): 11.2% to 6.7%**
  - Inception modules
  - 22 layers
  - 5 million parameters  
(throw away fully connected layers)
- **ResNet (2015): 6.7% to 3.57%**
  - More layers = better performance
  - 152 layers
- **CUIImage (2016): 3.57% to 2.99%**
  - Ensemble of 6 models
- **SENet (2017): 2.99% to 2.251%**
  - Squeeze and excitation block: network is allowed to adaptively adjust the weighting of each feature map in the convolutional block.

- **AlexNet (2012): First CNN (15.4%)**
  - 8 layers
  - 61 million parameters
- **ZFNet (2013): 15.4% to 11.2%**
  - 8 layers
  - More filters. Denser stride.
- **VGGNet (2014): 11.2% to 7.3%**
  - Beautifully uniform:  
3x3 conv, stride 1, pad 1, 2x2 max pool
  - 16 layers
  - 138 million parameters
- **GoogLeNet (2014): 11.2% to 6.7%**
  - Inception modules
  - 22 layers
  - 5 million parameters  
(throw away fully connected layers)
- **ResNet (2015): 6.7% to 3.57%**
  - More layers = better performance
  - 152 layers
- **CUIImage (2016): 3.57% to 2.99%**
  - Ensemble of 6 models

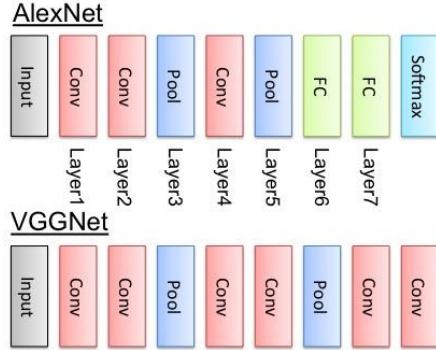
ImageNet Classification Error (Top 5)





- **AlexNet (2012): First CNN (15.4%)**
  - 8 layers
  - 61 million parameters
- **ZFNet (2013): 15.4% to 11.2%**
  - 8 layers
  - More filters. Denser stride.
- **VGGNet (2014): 11.2% to 7.3%**
  - Beautifully uniform:  
3x3 conv, stride 1, pad 1, 2x2 max pool
  - 16 layers
  - 138 million parameters
- **GoogLeNet (2014): 11.2% to 6.7%**
  - Inception modules
  - 22 layers
  - 5 million parameters  
(throw away fully connected layers)
- **ResNet (2015): 6.7% to 3.57%**
  - More layers = better performance
  - 152 layers
- **CUIImage (2016): 3.57% to 2.99%**
  - Ensemble of 6 models

Krizhevsky et al. "Imagenet classification with deep convolutional neural networks." Advances in neural information processing systems. 2012.

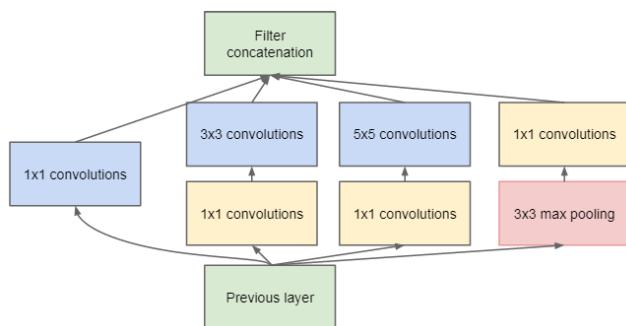
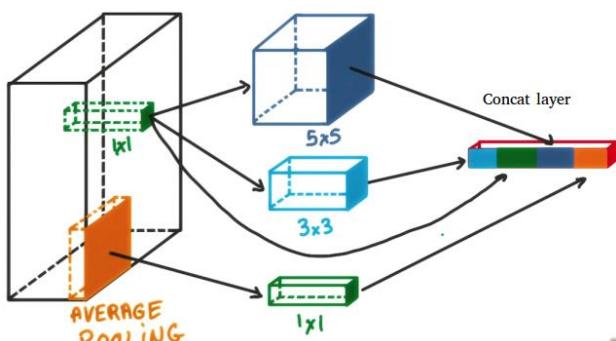
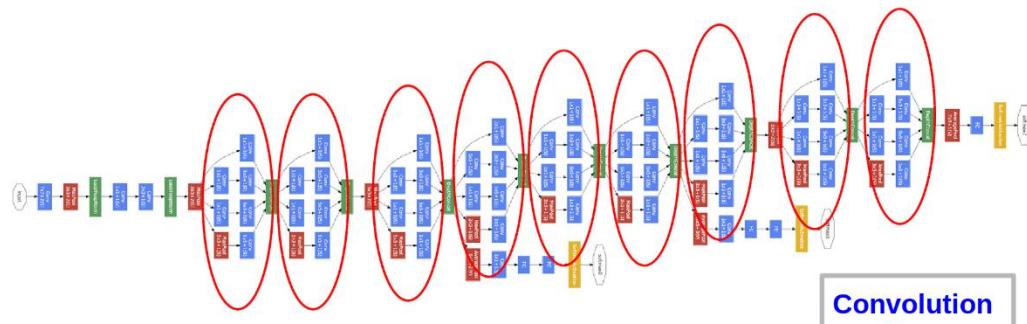


Legend:

- Input : Image input
- Conv : Convolutional layer
- Pool : Max-pooling layer
- FC : Fully-connected layer
- Softmax : Softmax layer

- **AlexNet (2012): First CNN (15.4%)**
  - 8 layers
  - 61 million parameters
- **ZFNet (2013): 15.4% to 11.2%**
  - 8 layers
  - More filters. Denser stride.
- **VGGNet (2014): 11.2% to 7.3%**
  - Beautifully uniform:  
3x3 conv, stride 1, pad 1, 2x2 max pool
  - 16 layers
  - 138 million parameters
- **GoogLeNet (2014): 11.2% to 6.7%**
  - Inception modules
  - 22 layers
  - 5 million parameters  
(throw away fully connected layers)
- **ResNet (2015): 6.7% to 3.57%**
  - More layers = better performance
  - 152 layers
- **CUIImage (2016): 3.57% to 2.99%**
  - Ensemble of 6 models

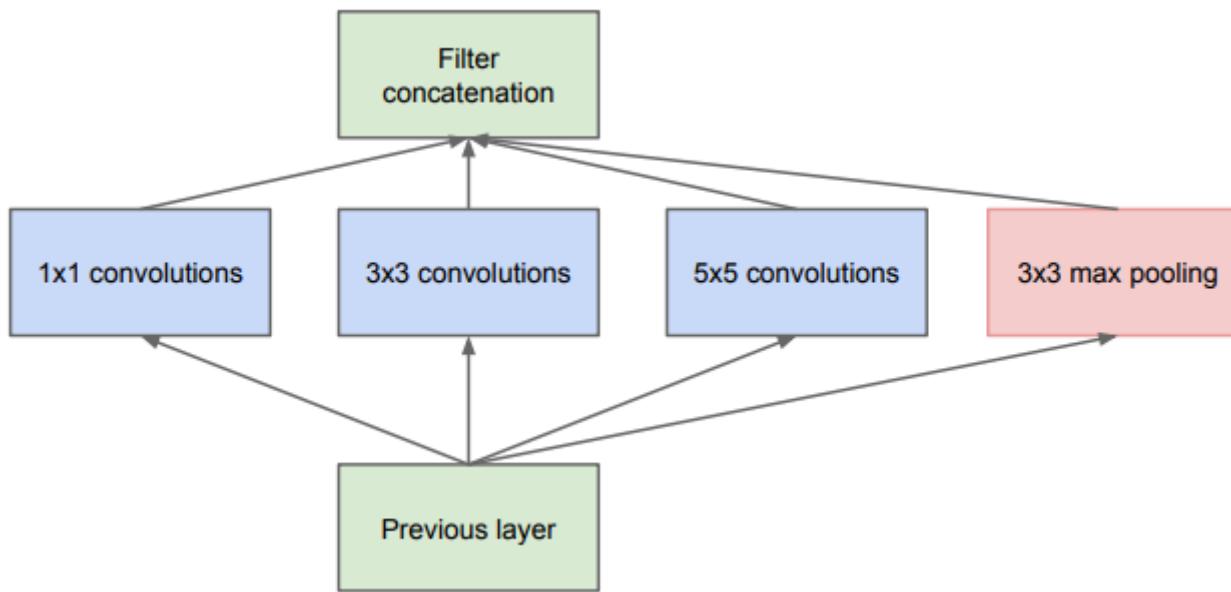
Simonyan et al. "Very deep convolutional networks for large-scale image recognition." 2014.



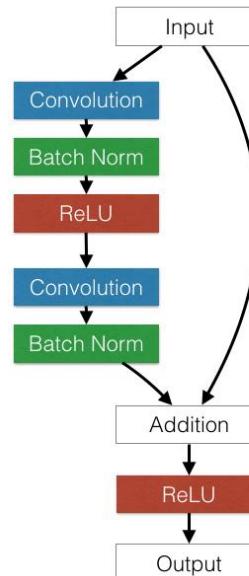
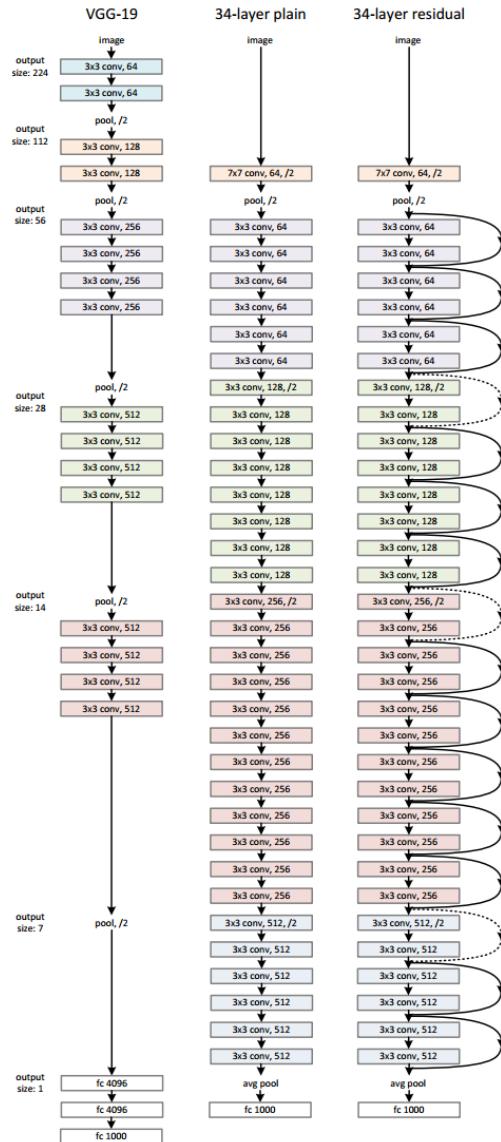
- **AlexNet (2012): First CNN (15.4%)**
  - 8 layers
  - 61 million parameters
- **ZFNet (2013): 15.4% to 11.2%**
  - 8 layers
  - More filters. Denser stride.
- **VGGNet (2014): 11.2% to 7.3%**
  - Beautifully uniform:  
3x3 conv, stride 1, pad 1, 2x2 max pool
  - 16 layers
  - 138 million parameters
- **GoogLeNet (2014): 11.2% to 6.7%**
  - Inception modules
  - 22 layers
  - 5 million parameters  
(throw away fully connected layers)
- **ResNet (2015): 6.7% to 3.57%**
  - More layers = better performance
  - 152 layers
- **CUIImage (2016): 3.57% to 2.99%**
  - Ensemble of 6 models

Szegedy et al. "Going deeper with convolutions." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2015.

# Inception Module



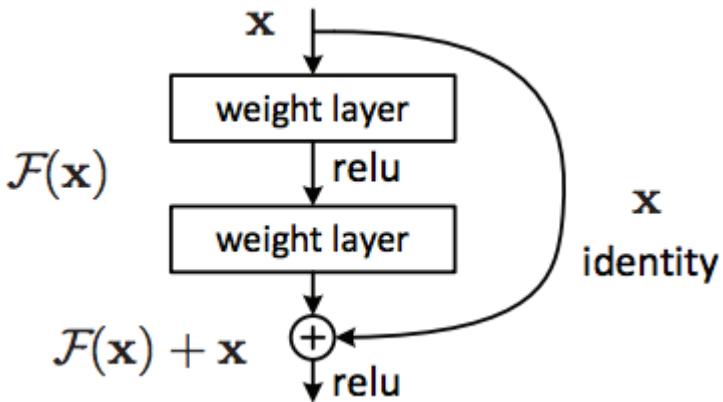
- **Process:** do different size convolutions, and concatenate
- Convolution sizes:
  - Smaller convolutions: local features
  - Larger convolutions: high-abstraction features
- **Result:** Fewer parameters and better performance



- **AlexNet (2012): First CNN (15.4%)**
  - 8 layers
  - 61 million parameters
- **ZFNet (2013): 15.4% to 11.2%**
  - 8 layers
  - More filters. Denser stride.
- **VGGNet (2014): 11.2% to 7.3%**
  - Beautifully uniform:  
3x3 conv, stride 1, pad 1, 2x2 max pool
  - 16 layers
  - 138 million parameters
- **GoogLeNet (2014): 11.2% to 6.7%**
  - Inception modules
  - 22 layers
  - 5 million parameters  
(throw away fully connected layers)
- **ResNet (2015): 6.7% to 3.57%**
  - More layers = better performance
  - 152 layers
- **CUIImage (2016): 3.57% to 2.99%**
  - Ensemble of 6 models

He, Kaiming, et al. "Deep residual learning for image recognition." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2016.

# Residual Block

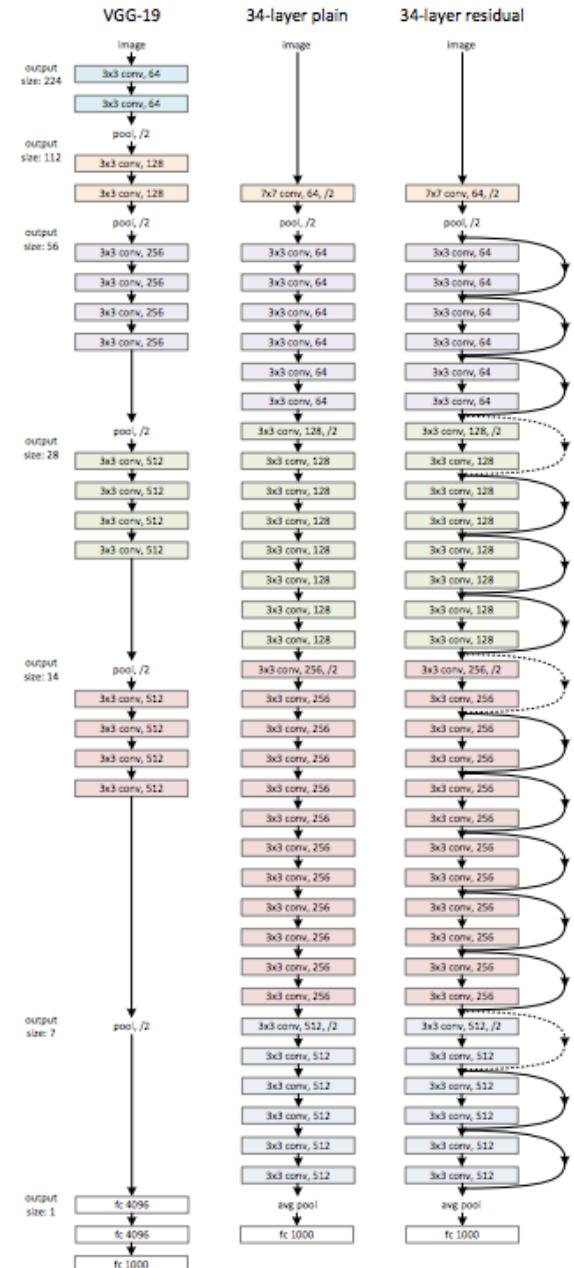


- **Initial Observation:**

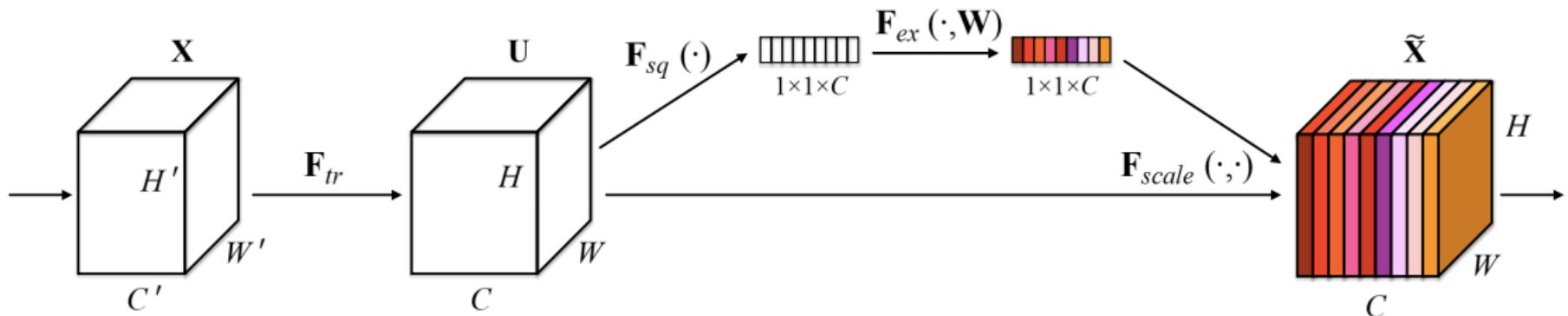
- Network depth often increases representation power, but is harder to train.

- **Residual Block:**

- Repeat a simple network block (think: RNN)
- Pass input along without transformation: help ensure that each layer learns something new

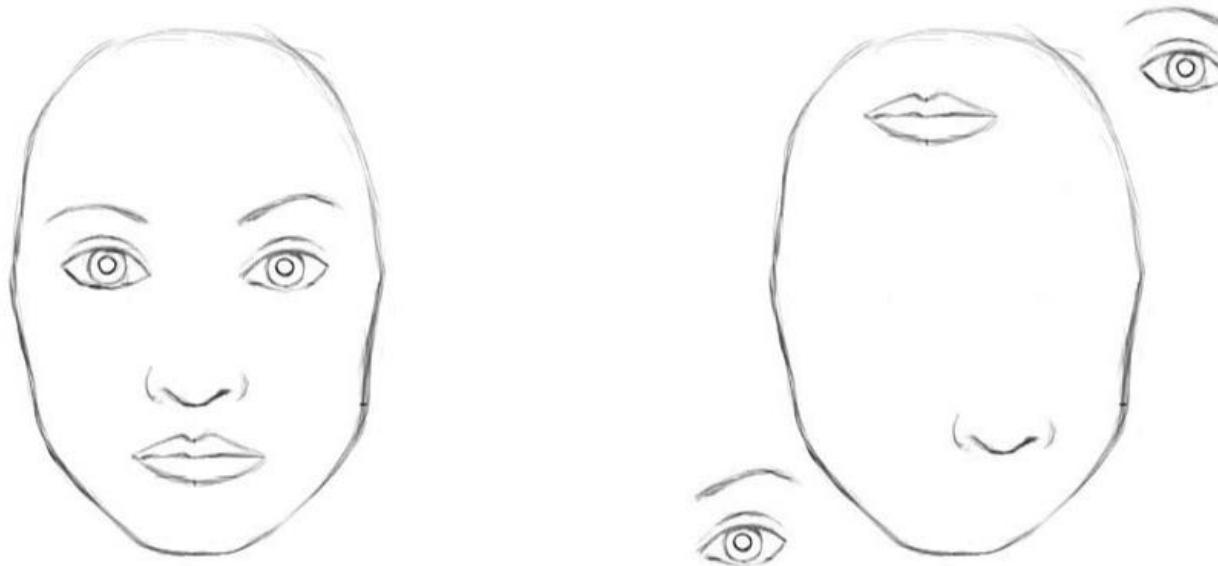


# SENet: Squeeze-and-Excitation Networks



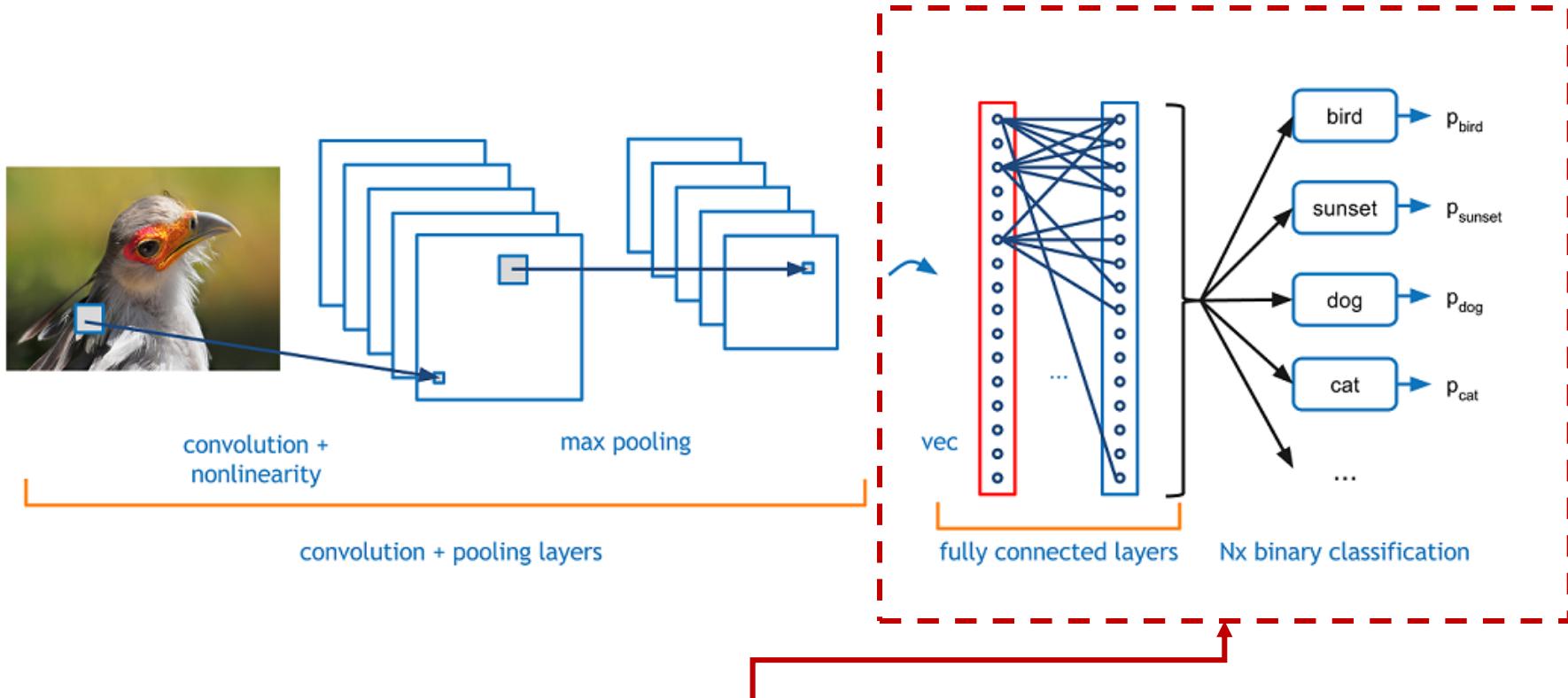
- **Content-aware channel weighting:** Add parameters to each channel of a convolutional block so that the network can adaptively adjust the weighting of each feature map
- This approach is simple and can be added to any model
  - **Takeaway for thought:** Parameterize everything (that's cost-effective) including higher-order hyper-parameters.

# Capsule Networks (Hinton)



- A CNN see both images as the same. The problem:
  - *Internal data representation of a convolutional neural network does not take into account important spatial hierarchies between simple and complex objects.*
- See upcoming online-only lecture on capsule networks.

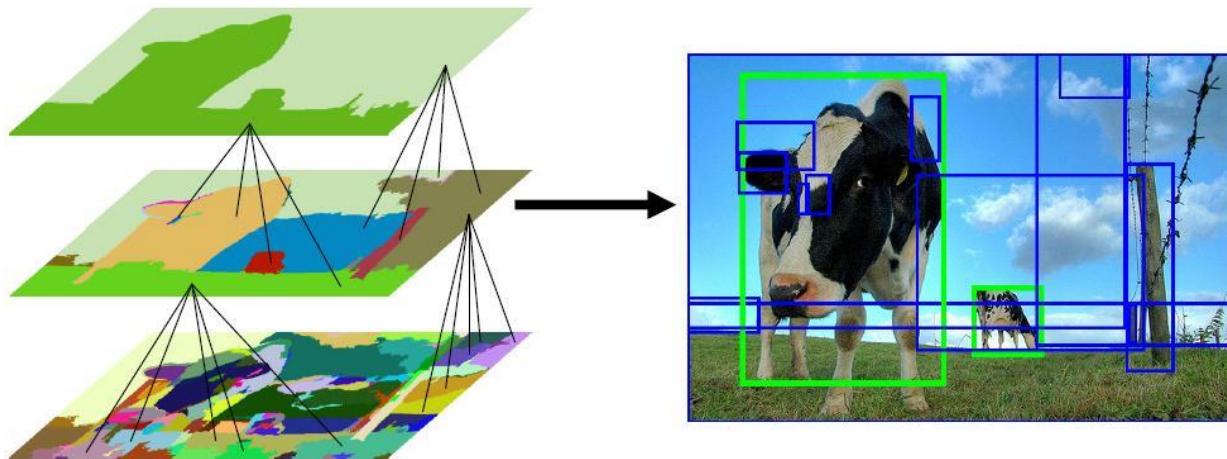
# Same Architecture, Many Applications



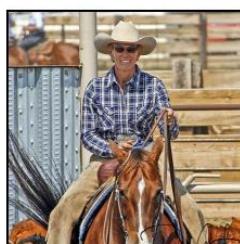
This part might look different for:

- Different image classification **domains**
- Image captioning with **recurrent neural networks**
- Image object localization with **bounding box**
- Image segmentation with **fully convolutional networks**
- Image segmentation with **deconvolution layers**

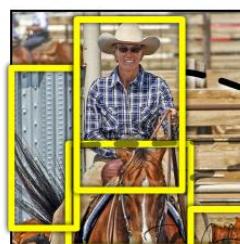
# Object Detection



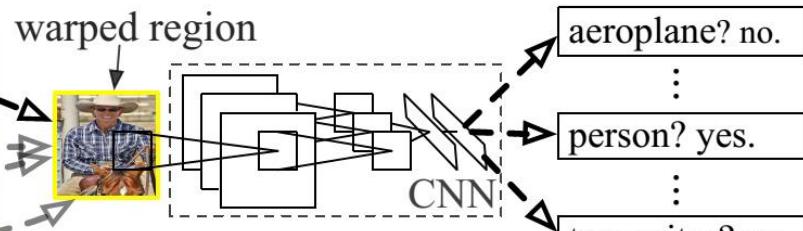
## R-CNN: *Regions with CNN features*



1. Input image



2. Extract region proposals (~2k)



3. Compute CNN features

4. Classify regions

# Fully Convolutional Networks

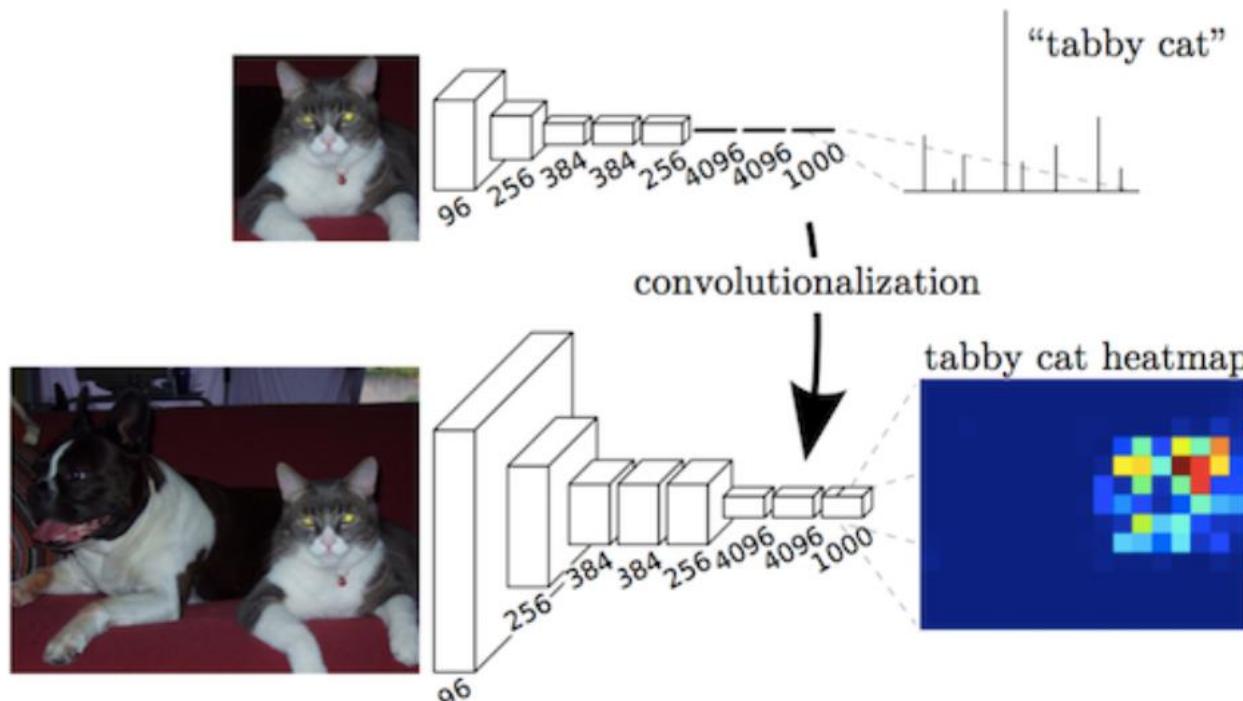
- **Goal:** Classify every pixel in an image.
- **Difficulty:** Hard
- **Why?**
  - When precise boundaries of objects matter (medical, driving)
  - Useful for fusing with other sensors (LIDAR)



# FCN (Nov 2014)

Paper: “Fully Convolutional Networks for Semantic Segmentation”

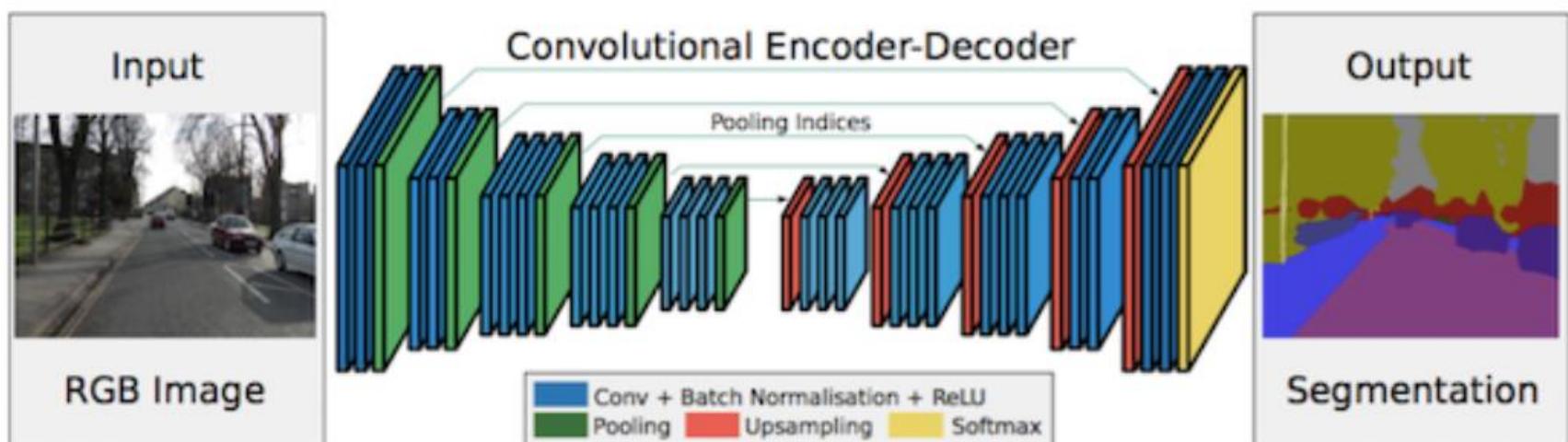
- Repurpose Imagenet pretrained nets
- Upsample using deconvolution
- Skip connections to improve coarseness of upsampling



# SegNet (Nov 2015)

Paper: "SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation"

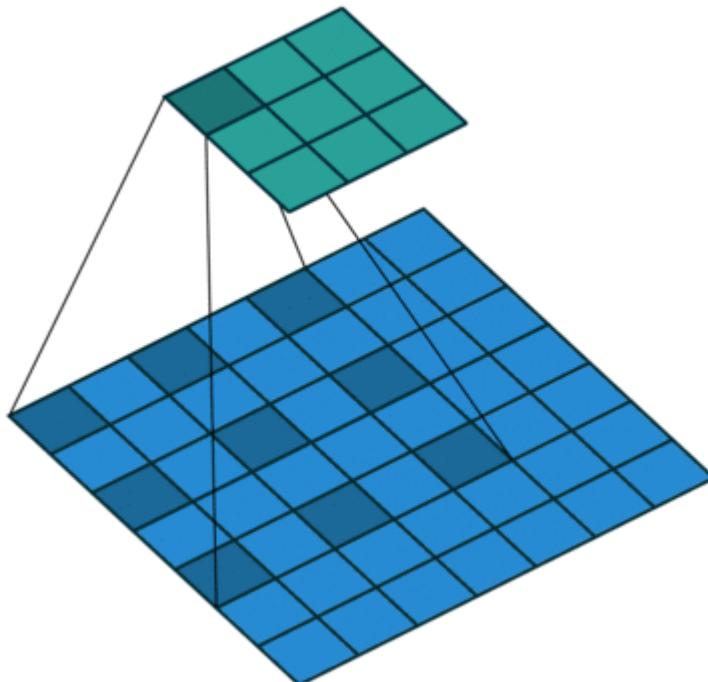
- Maxpooling indices transferred to decoder to improve the segmentation resolution.



# Dilated Convolutions (Nov 2015)

Paper: "Multi-Scale Context Aggregation by Dilated Convolutions"

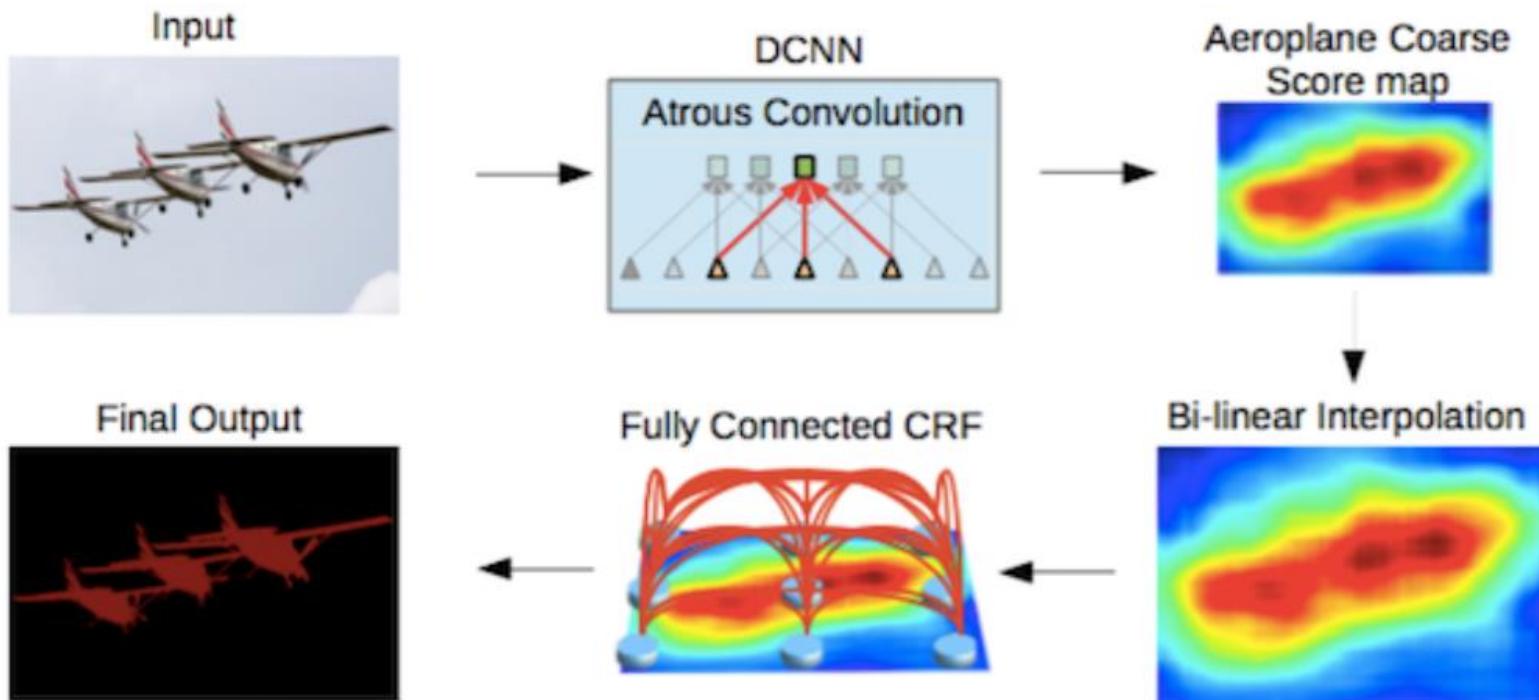
- Since pooling decreases resolution:
  - Added “dilated convolution layer”
- Still interpolate up from 1/8 of original image size



# DeepLap v1, v2 (Jun 2016)

Paper: "DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs"

- Added fully-connected Conditional Random Fields (CRFs) – as a post-processing step
  - Smooth segmentation based on the underlying image intensities



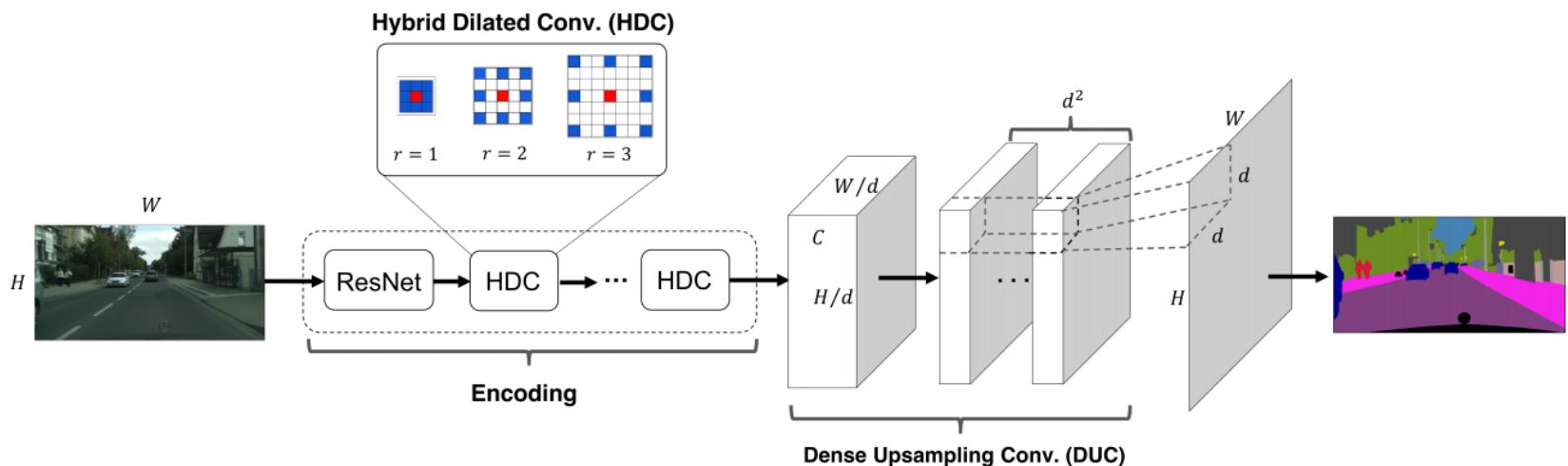
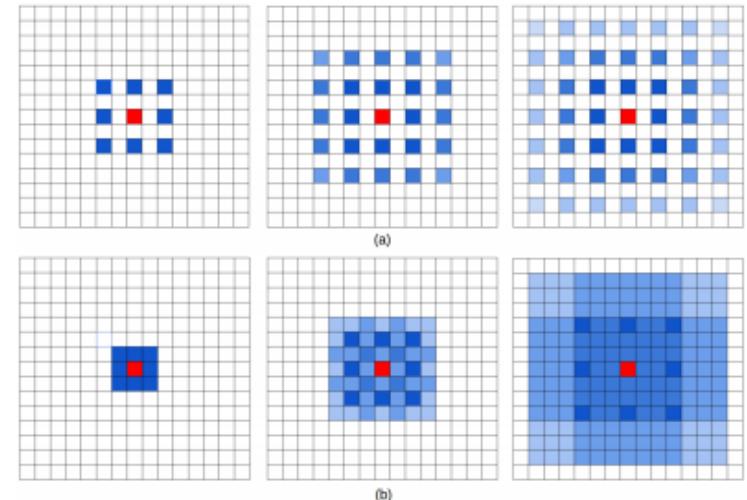
# Key Aspects of Segmentation

- **Fully convolutional networks (FCNs)** - replace fully-connected layers with convolutional layers
  - Deeper, updated models (now ResNet) consistent with ImageNet Challenge object classification tasks.
- **Conditional Random Fields (CRFs)** to capture both local and long-range dependencies within an image to refine the prediction map.
- **Dilated convolution** (aka Atrous convolution) – maintain computational cost, increase resolution of intermediate feature maps

# ResNet-DUC (Nov 2017)

Paper: "Understanding Convolution for Semantic Segmentation"

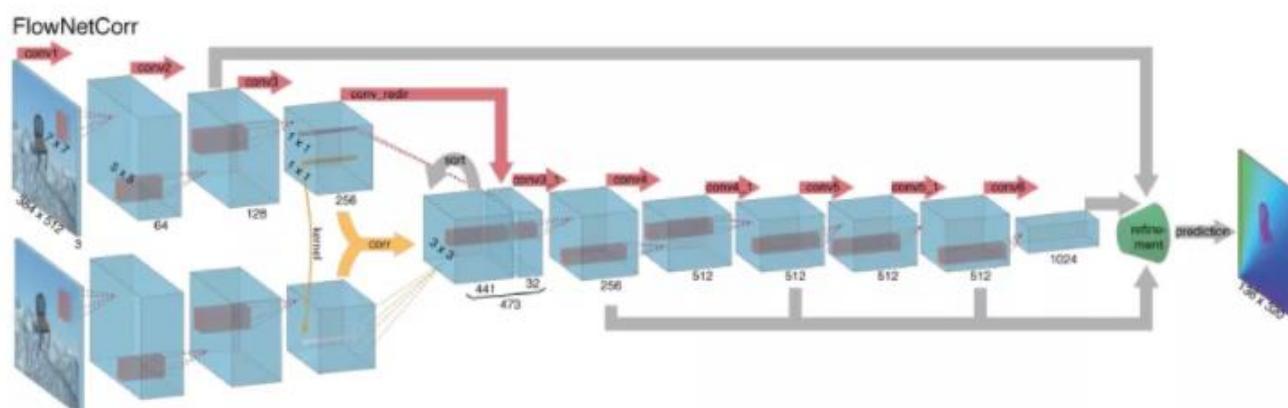
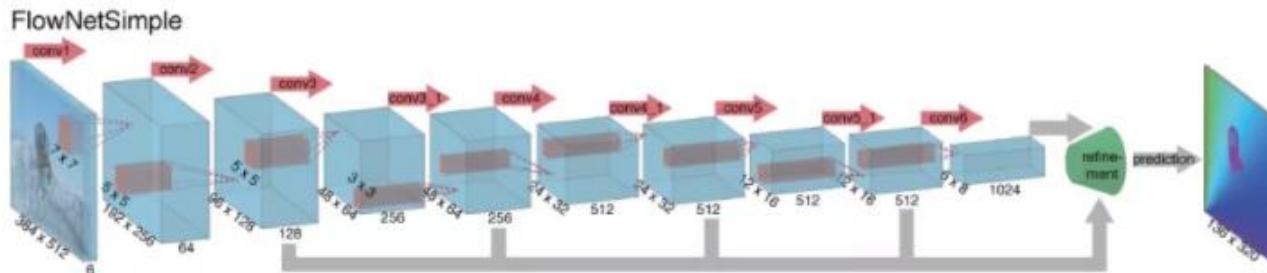
- Dense upsampling convolution (DUC) instead of bilinear upsampling
  - **Learnable:** Learn the upscaling filters
- Hybrid dilated convolution (HDC)
  - Use a different dilation rate



# FlowNet (May 2015)

Paper: "FlowNet: Learning Optical Flow with Convolutional Networks"

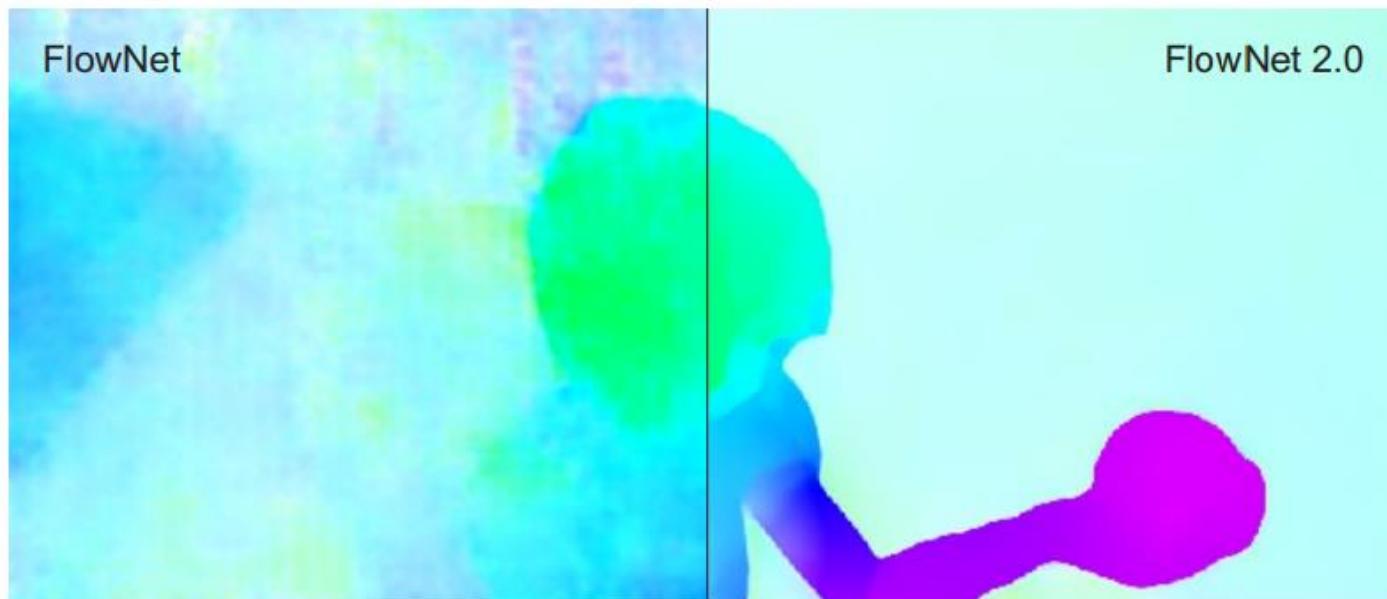
- Learn flow from image-pair, end to end.
  - FlowNetS – stacks two images as input
  - FlowNetC – convolute separately, combine with correlation layer



# FlowNet 2.0 (Dec 2016)

Paper: "FlowNet 2.0: Evolution of Optical Flow Estimation with Deep Networks"

- Stack FlowNetS and FlowNetC
- Improvement over FlowNet
  - Smooth flow fields
  - Preserves fine-motion detail
  - Runs at 8-140fps
- Observations:
  - Stacking networks as an approach
  - Order of training dataset matters



# SegFuse: Dynamic Driving Scene Segmentation



[cars.mit.edu/segfuse](https://selfdrivingcars.mit.edu/segfuse)

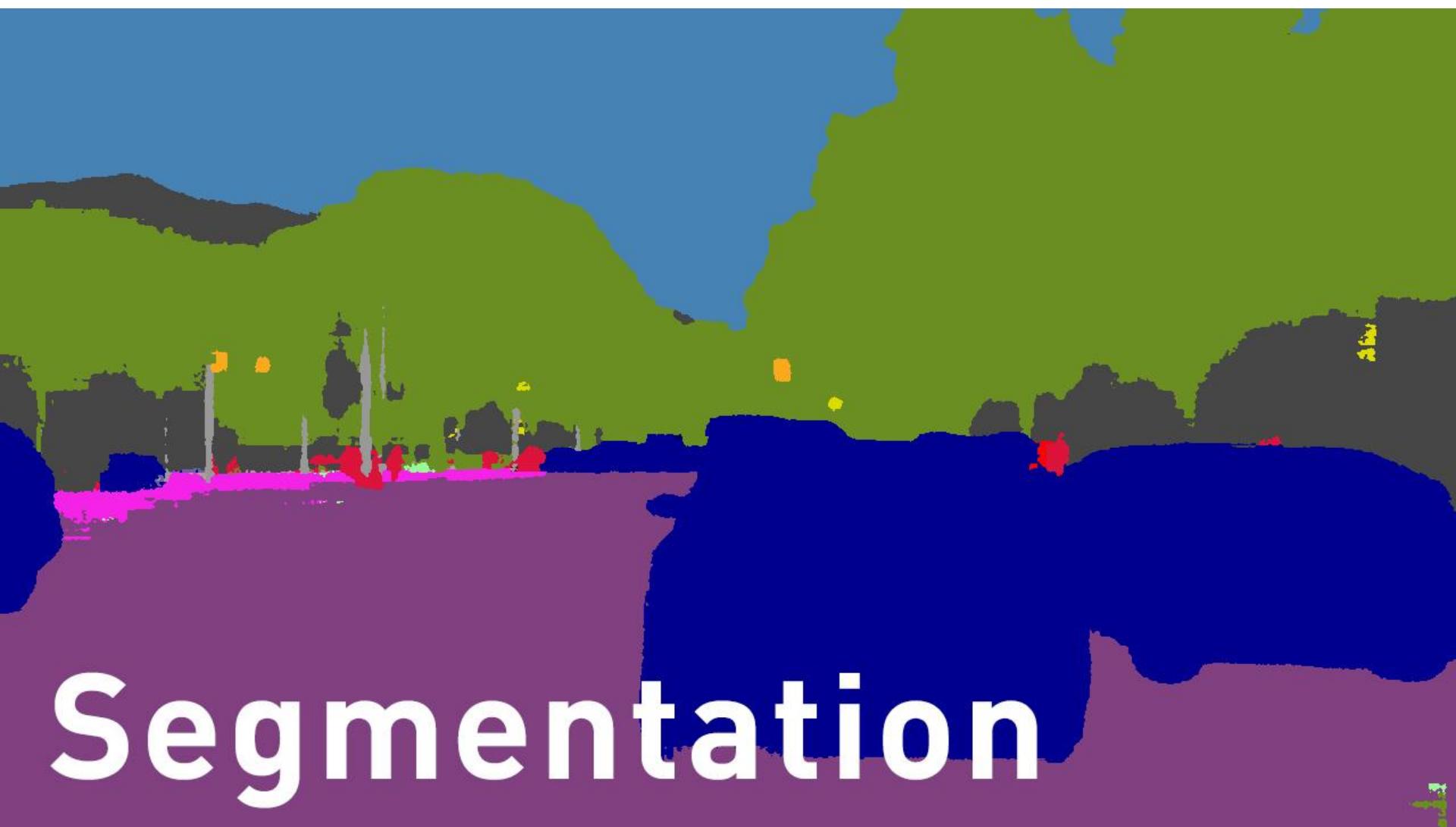
# SegFuse: Dynamic Driving Scene Segmentation



# Ground Truth

[cars.mit.edu/segfuse](http://cars.mit.edu/segfuse)

# SegFuse: Dynamic Driving Scene Segmentation



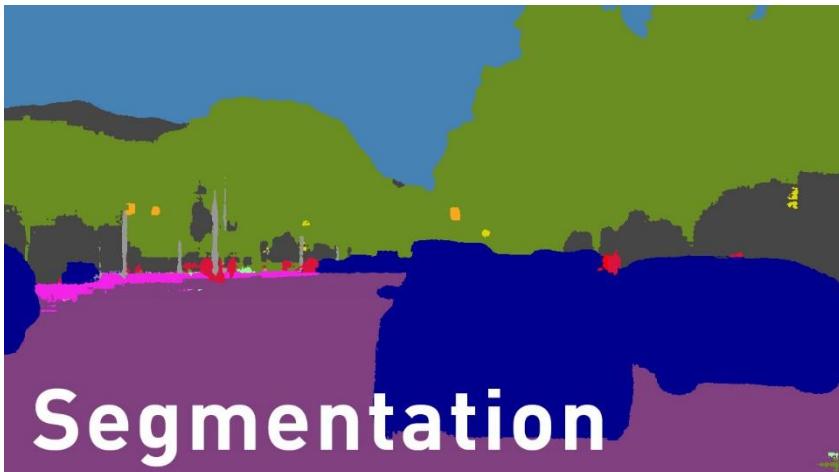
[cars.mit.edu/segfuse](http://cars.mit.edu/segfuse)

# SegFuse: Dynamic Driving Scene Segmentation

# Optical Flow

[cars.mit.edu/segfuse](http://cars.mit.edu/segfuse)

# SegFuse: Dynamic Driving Scene Segmentation



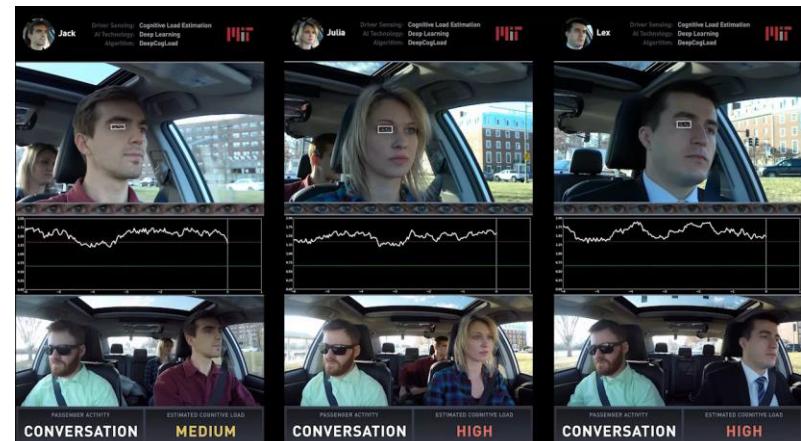
[cars.mit.edu/segfuse](http://cars.mit.edu/segfuse)

# Thank You

*Tomorrow: Waymo*



*Next lecture: Deep Learning for Human Sensing*



## Upcoming online-only lectures:

- Capsule networks
- Generative adversarial networks