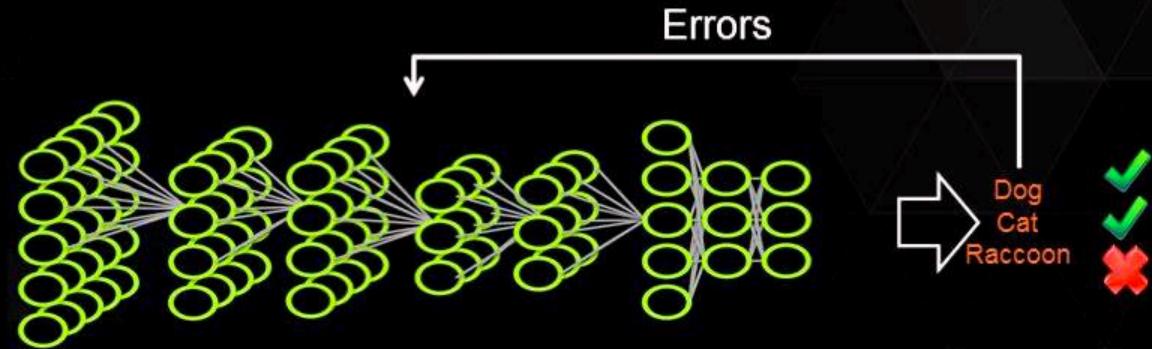
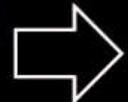


Deep Learning for Self Driving Cars

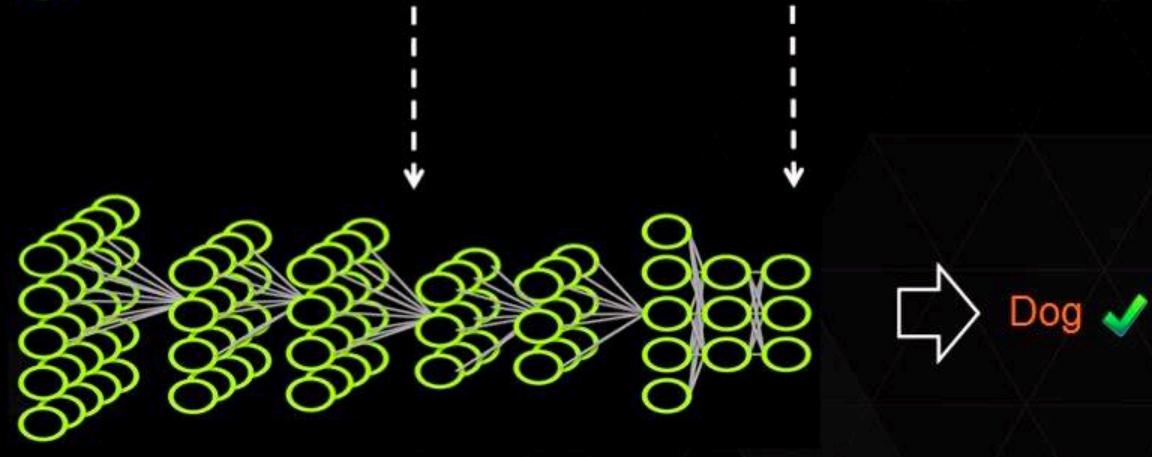
Neural Networks for Vision

DEEP LEARNING APPROACH

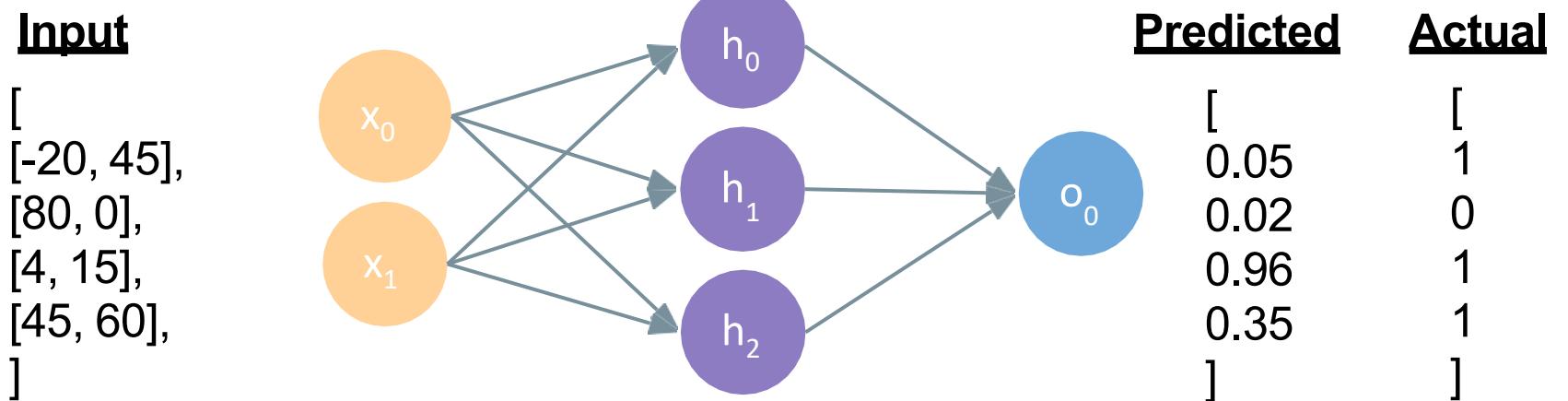
Train:



Deploy:



Total Loss



$$\text{total loss} := J(\theta) = \frac{1}{N} \sum_i^N loss(f(x^{(i)}; \theta), y^{(i)})$$

————— Predicted ————— Actual

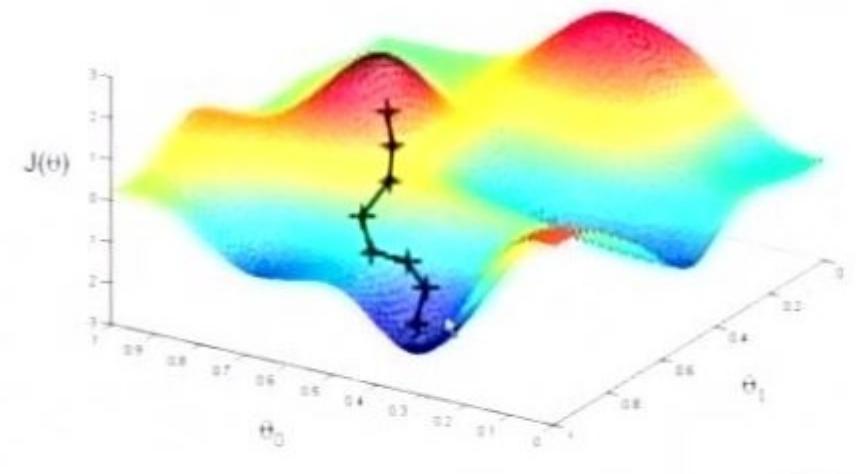
Why is it Stochastic Gradient Descent?

- Initialize θ randomly
- For N Epochs
 - For each training example (x, y) :

- Compute Loss Gradient:
- Update θ with update rule:

$$\frac{\partial J(\theta)}{\partial \theta}$$

Only an estimate of
true gradient!



$$\theta := \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$$

Minibatches Reduce Gradient Variance

- Initialize θ randomly
- For N Epochs
 - For each training **batch** $\{(x_0, y_0), \dots, (x_B, y_B)\}$:

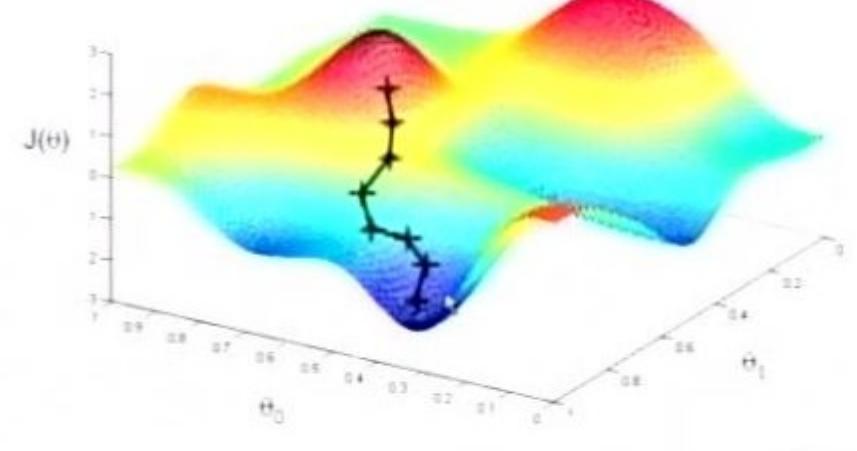
- Compute Loss Gradient:

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{1}{B} \sum_i^B \frac{\partial J_i(\theta)}{\partial \theta}$$

- Update θ with update rule:

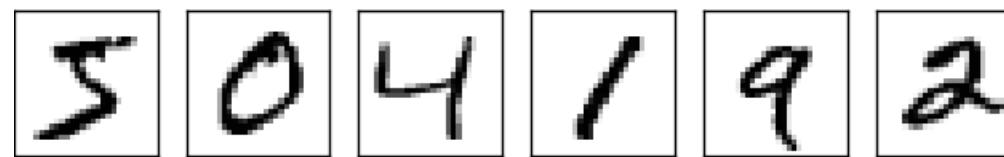
$$\theta := \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$$

More accurate estimate!

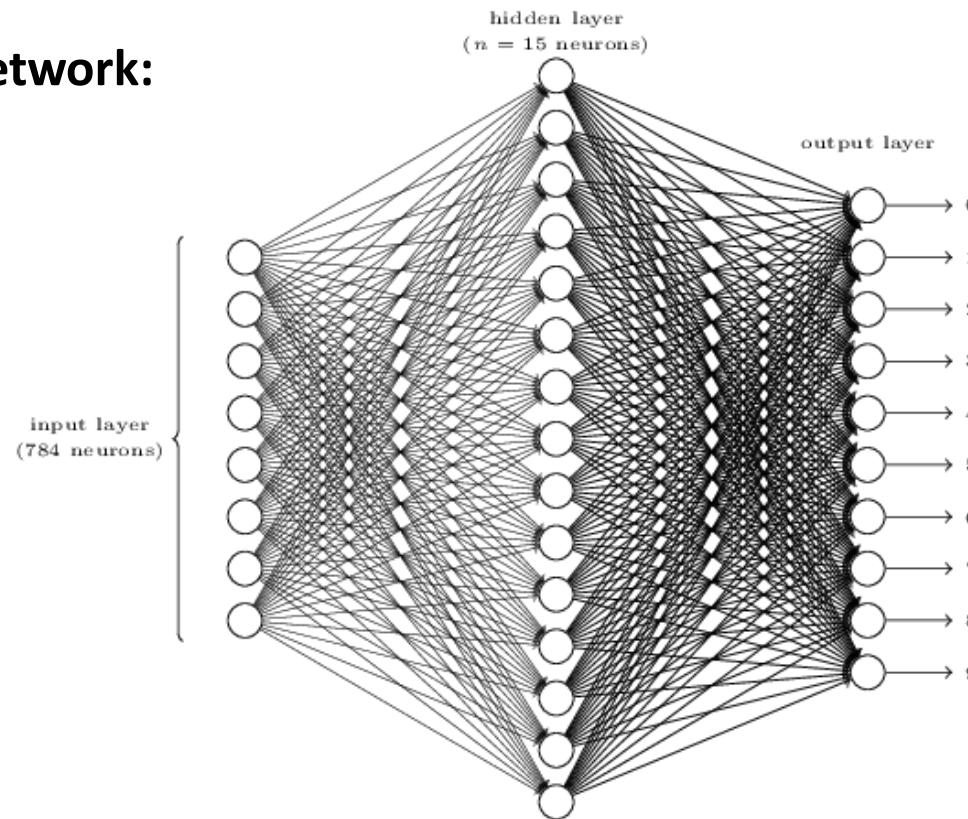


Classify and Image of a Number

Input:
(28x28)

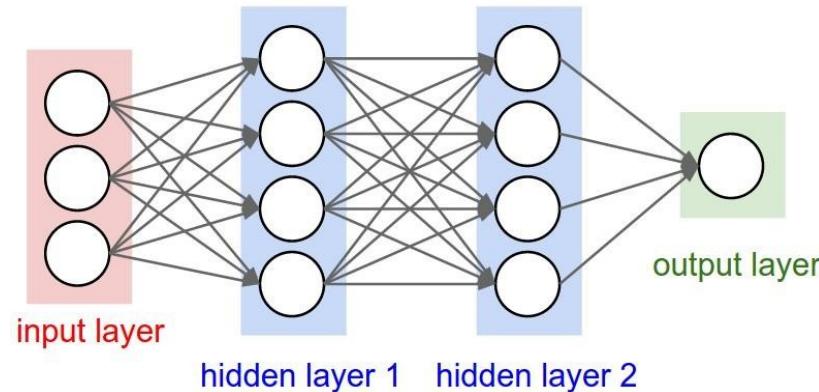


Network:

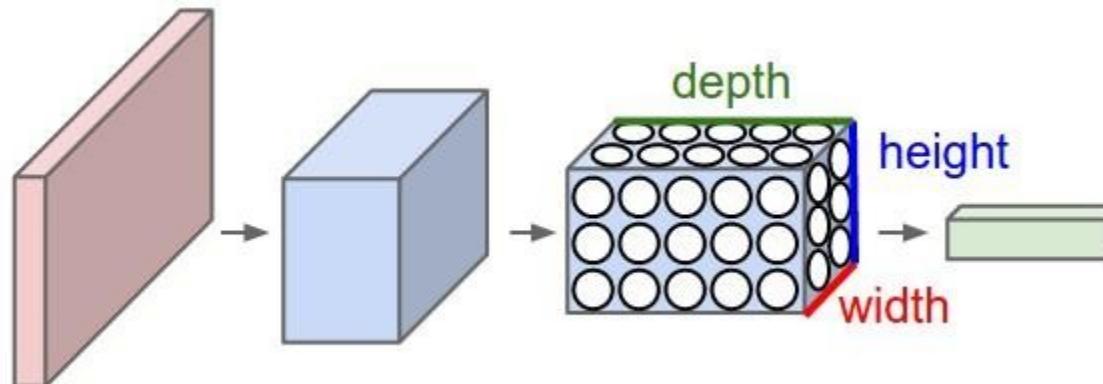


Convolutional Neural Networks

Regular neural network (fully connected):



Convolutional neural network:

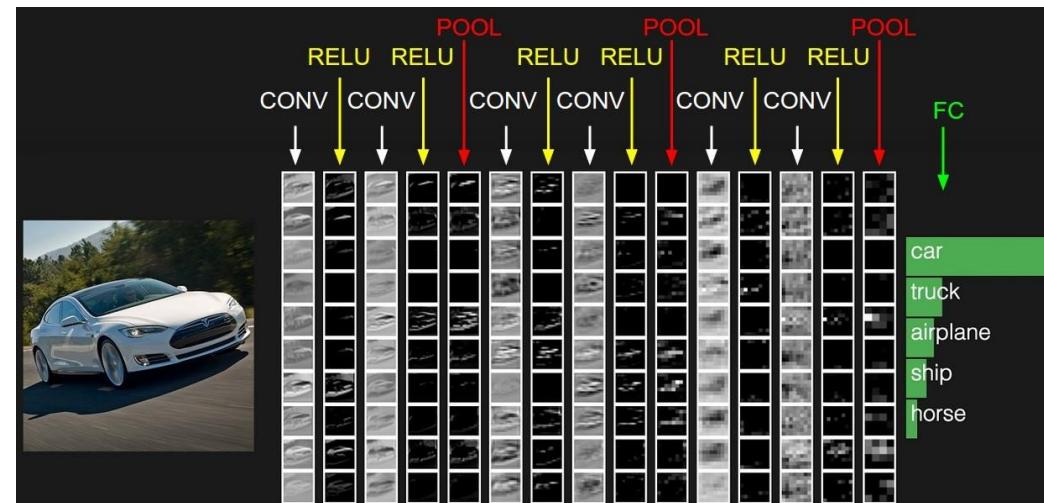


Each layer takes a 3d volume, produces 3d volume with some smooth function that may or may not have parameters.

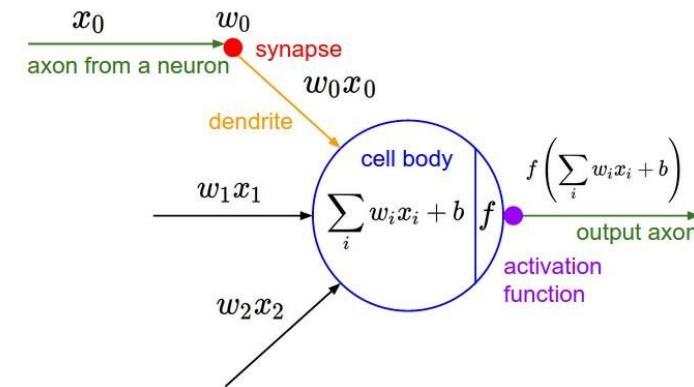
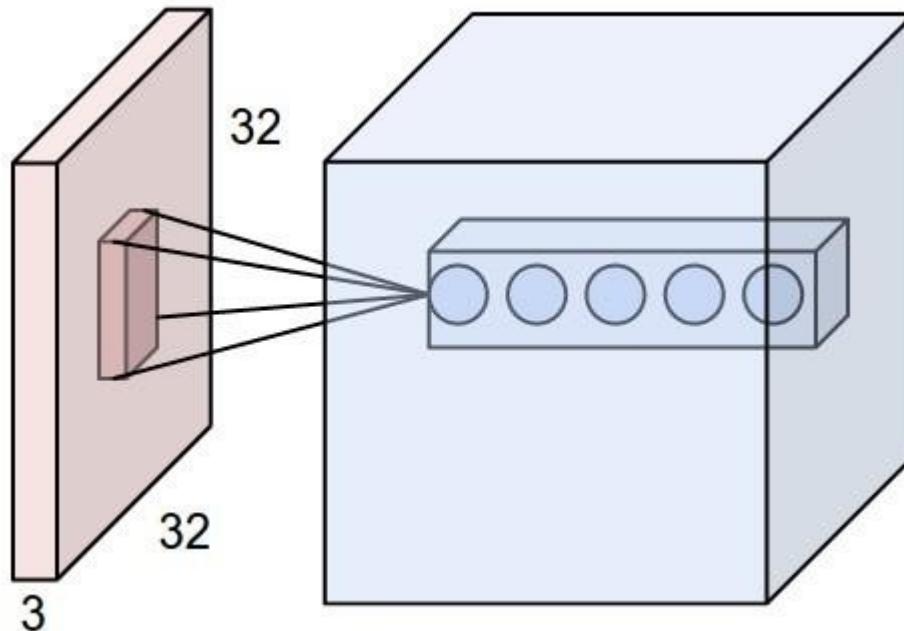
Convolutional Neural Networks: Layers

- **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.



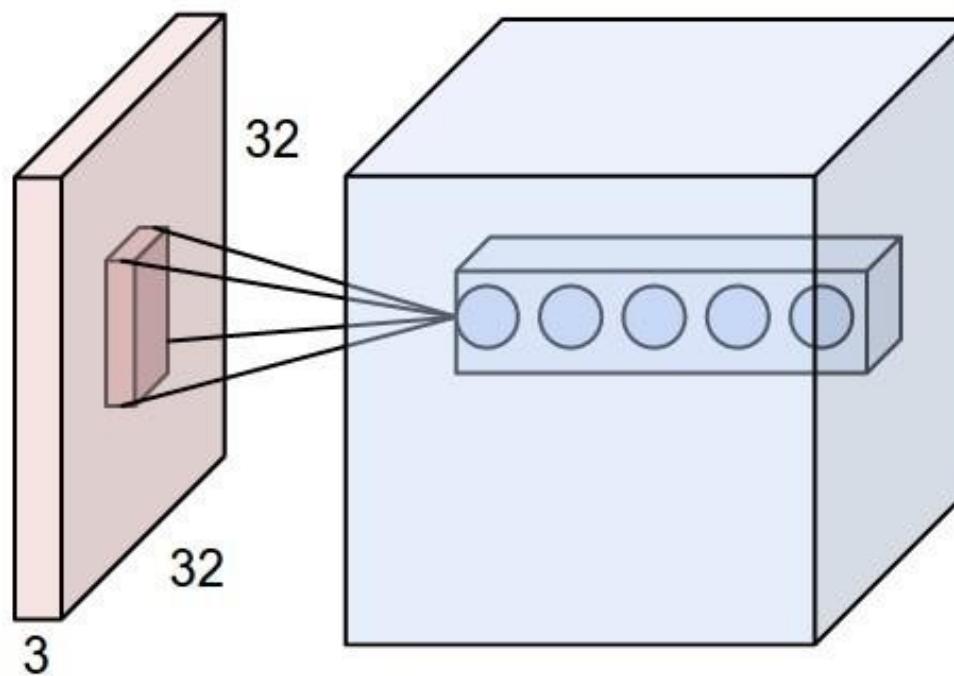
Dealing with Images: Local Connectivity



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

- Let f be the image and g be the kernel. The output of convolving f with g is denoted by $f * g$.

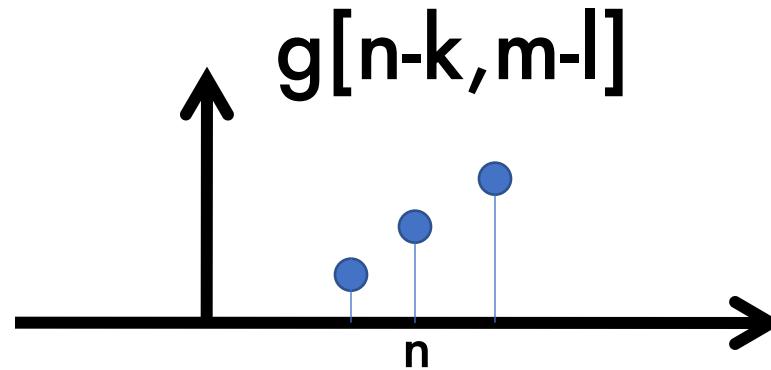
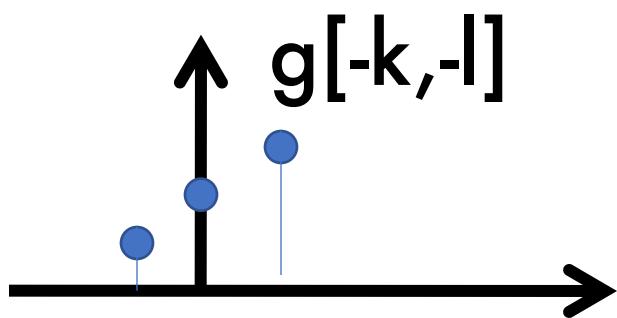
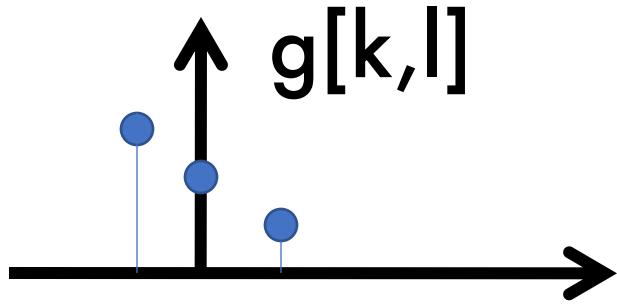
$$(f * g)[n, m] = \sum_{k,l} f[k, l] g[n - k, m - l]$$

- Weighted product of $f(k, l)$ by $g(-(k, l))$ computed at different locations n, m
- MATLAB: `conv2` vs. `filter2` (also `imfilter`)

Discrete convolution

$$(f * g)[n, m] = \sum_{k,l} f[k, l] g[n - k, m - l]$$

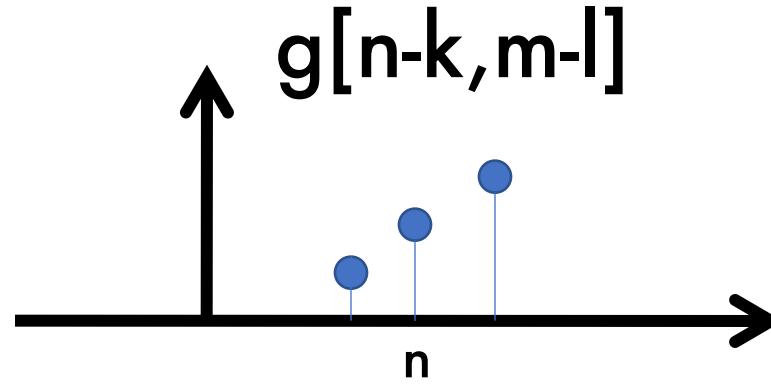
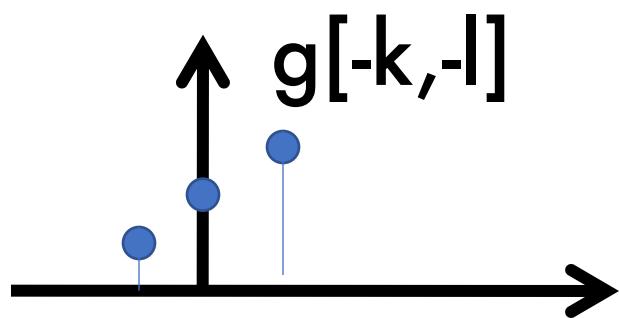
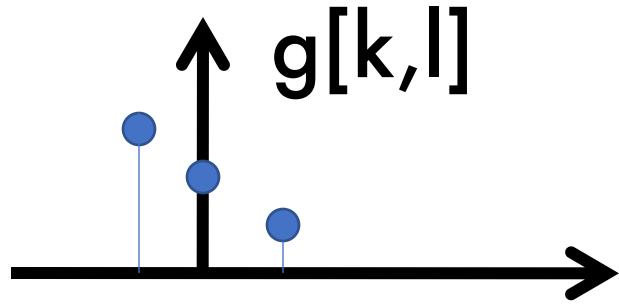
- Fold $g[k, l]$ about origin to form $g[-k, -l]$
- Shift the folded results by n, m to form $g[n - k, m - l]$
- Multiply $g[n - k, m - l]$ by $f[k, l]$
- Sum over all k, l
- Repeat for every n, m



Discrete convolution

$$(f * g)[n, m] = \sum_{k,l} f[k, l] g[n - k, m - l]$$

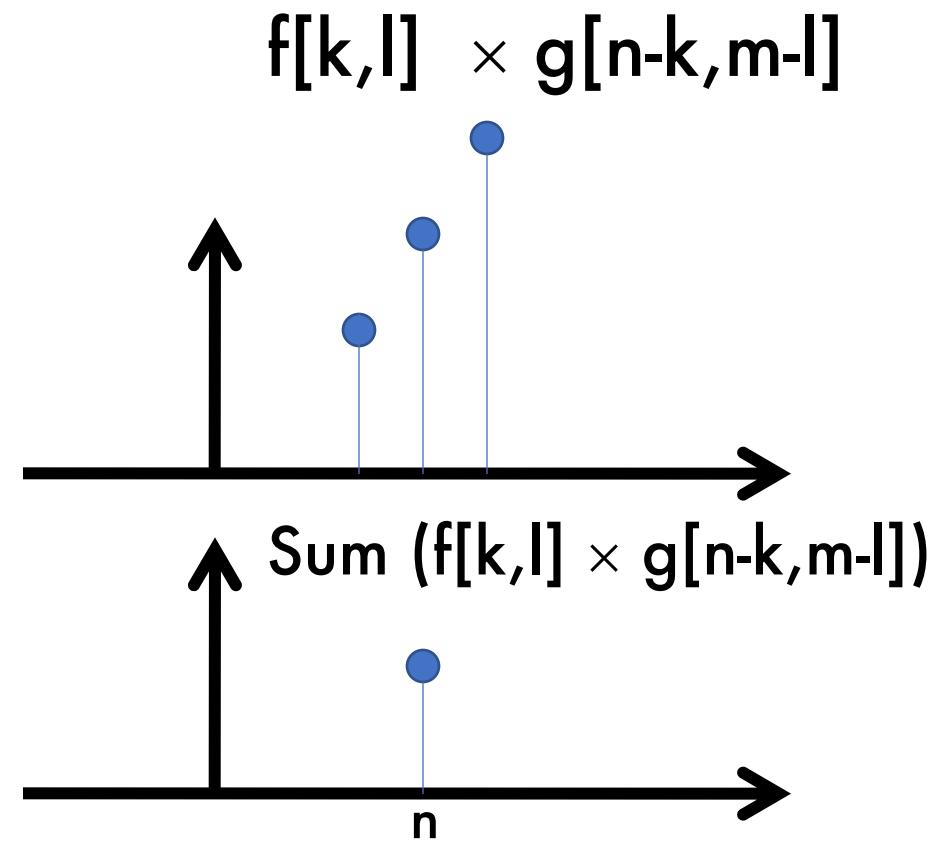
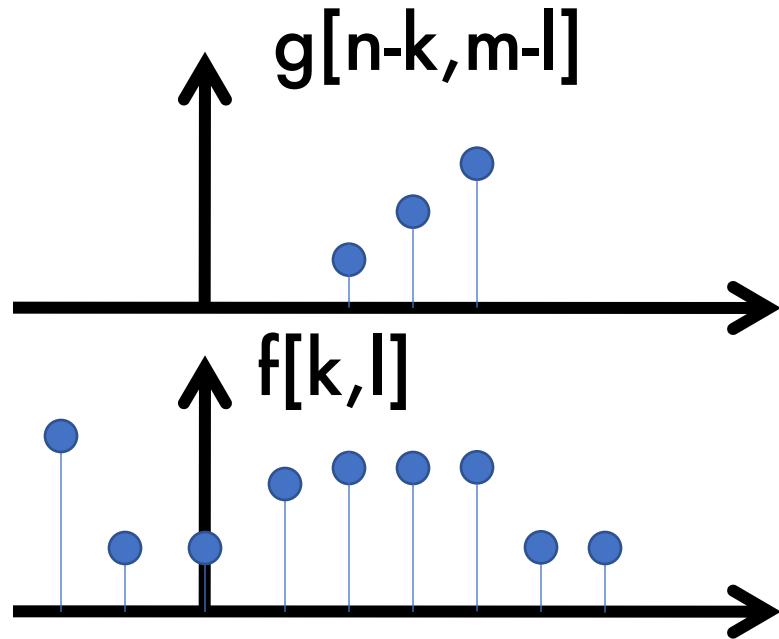
- Fold $g[k, l]$ about origin to form $g[-k, -l]$
- Shift the folded results by n, m to form $g[n - k, m - l]$
- Multiply $g[n - k, m - l]$ by $f[k, l]$
- Sum over all k, l
- Repeat for every n, m



Discrete convolution

$$(f * g)[n, m] = \sum_{k,l} f[k, l] g[n - k, m - l]$$

- Fold $g[k, l]$ about origin to form $g[-k, -l]$
- Shift the folded results by n, m to form $g[n - k, m - l]$
- Multiply $g[n - k, m - l]$ by $f[k, l]$
- Sum over all k, l
- Repeat for every n, m

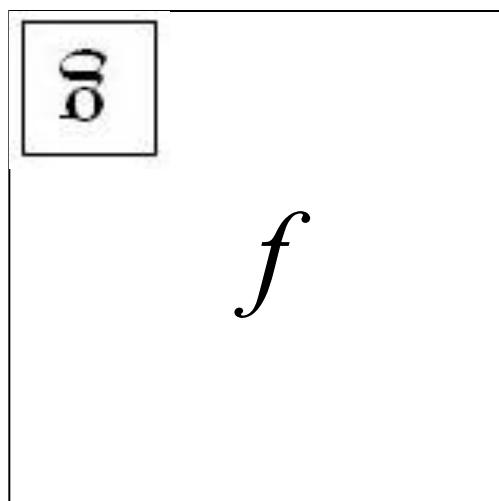


Convolution

- Let f be the image and g be the kernel. The output of convolving f with g is denoted by $f * g$.

$$(f * g)[n, m] = \sum_{k, l} f[k, l] g[n - k, m - l]$$

Weighted product of $f(k, l)$ by $g(-(k, l))$) computed at different locations n, m



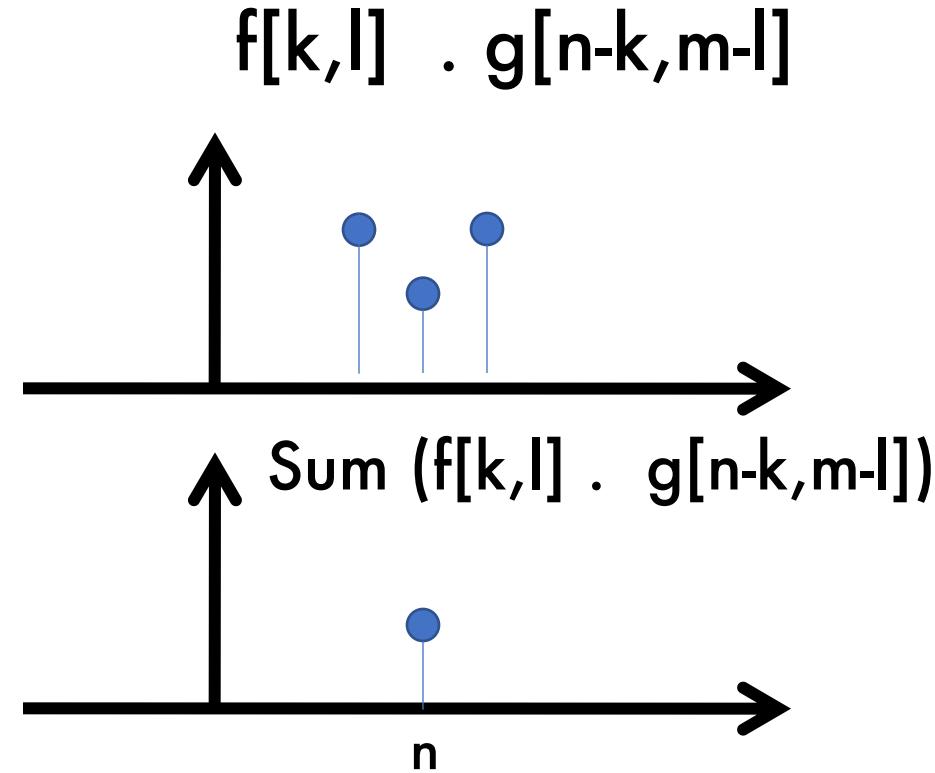
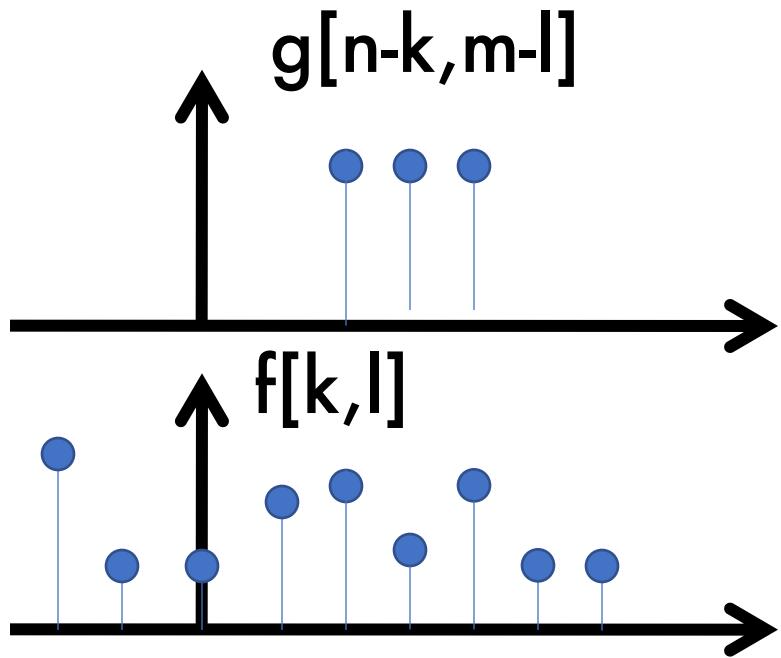
Box filter

- Kernel k with positive entries, that sum to 1.
- Notice: all weights are equal

$$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

Discrete convolution

- Fold $g[n,m]$ about origin to form $g[-k,-l]$
- Shift the folded results by n,m to form $g[n - k,m - l]$
- Multiply $g[n - k,m - l]$ by $f[k,l]$
- Sum over all k,l
- Repeat for every n,m

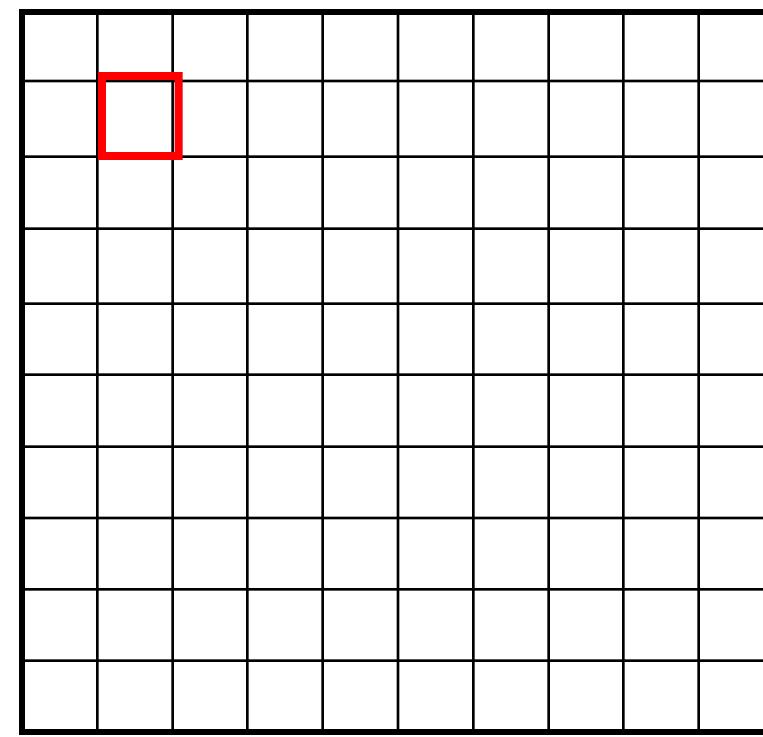


Box filter

$F[x, y]$

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	90	90	90	90	90	0	0	0	0
0	0	0	90	90	90	90	90	0	0	0	0
0	0	0	90	90	90	90	90	0	0	0	0
0	0	0	90	0	90	90	90	0	0	0	0
0	0	0	90	90	90	90	90	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0

$G[x, y]$



$$(f * g)[m, n] = \sum_{k,l} f[k, l] g[m - k, n - l]$$

Source: S. Seitz

Box filter

$F[x, y]$

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	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	90	0	90	90	90	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0

$G[x, y]$

0	10									

$$(f * g)[m, n] = \sum_{k,l} f[k, l] g[m - k, n - l]$$

Source: S. Seitz

Box filter

$F[x, y]$

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	90	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	90	0	0
0	0	0	90	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0

$G[x, y]$

			0	10	20					

$$(f * g)[m, n] = \sum_{k,l} f[k, l] g[m - k, n - l]$$

Source: S. Seitz

Box filter

$F[x, y]$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$G[x, y]$

$$(f * g)[m, n] = \sum_{k,l} f[k, l] g[m - k, n - l]$$

Source: S. Seitz

Box filter

$$F[x, y]$$

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	90	90	90	90	90	0	0	0	0
0	0	0	90	90	90	90	90	0	0	0	0
0	0	0	90	90	90	90	90	0	0	0	0
0	0	0	90	0	90	90	90	0	0	0	0
0	0	0	90	90	90	90	90	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0

$$G[x, y]$$

			0	10	20	30	30				

$$(f * g)[m, n] = \sum_{k,l} f[k, l] g[m - k, n - l]$$

Source: S. Seitz

Box filter

$$F[x, y]$$

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	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	90	0	90	90	90	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0

$$G[x, y]$$

	0	10	20	30	30	30	20	10	
	0	20	40	60	60	60	40	20	
	0	30	60	90	90	90	60	30	
	0	30	50	80	80	90	60	30	
	0	30	50	80	80	90	60	30	
	0	20	30	50	50	60	40	20	
	10	20	30	30	30	30	20	10	
	10	10	10	0	0	0	0	0	

$$(f * g)[m, n] = \sum_{k,l} f[k, l] g[m - k, n - l]$$

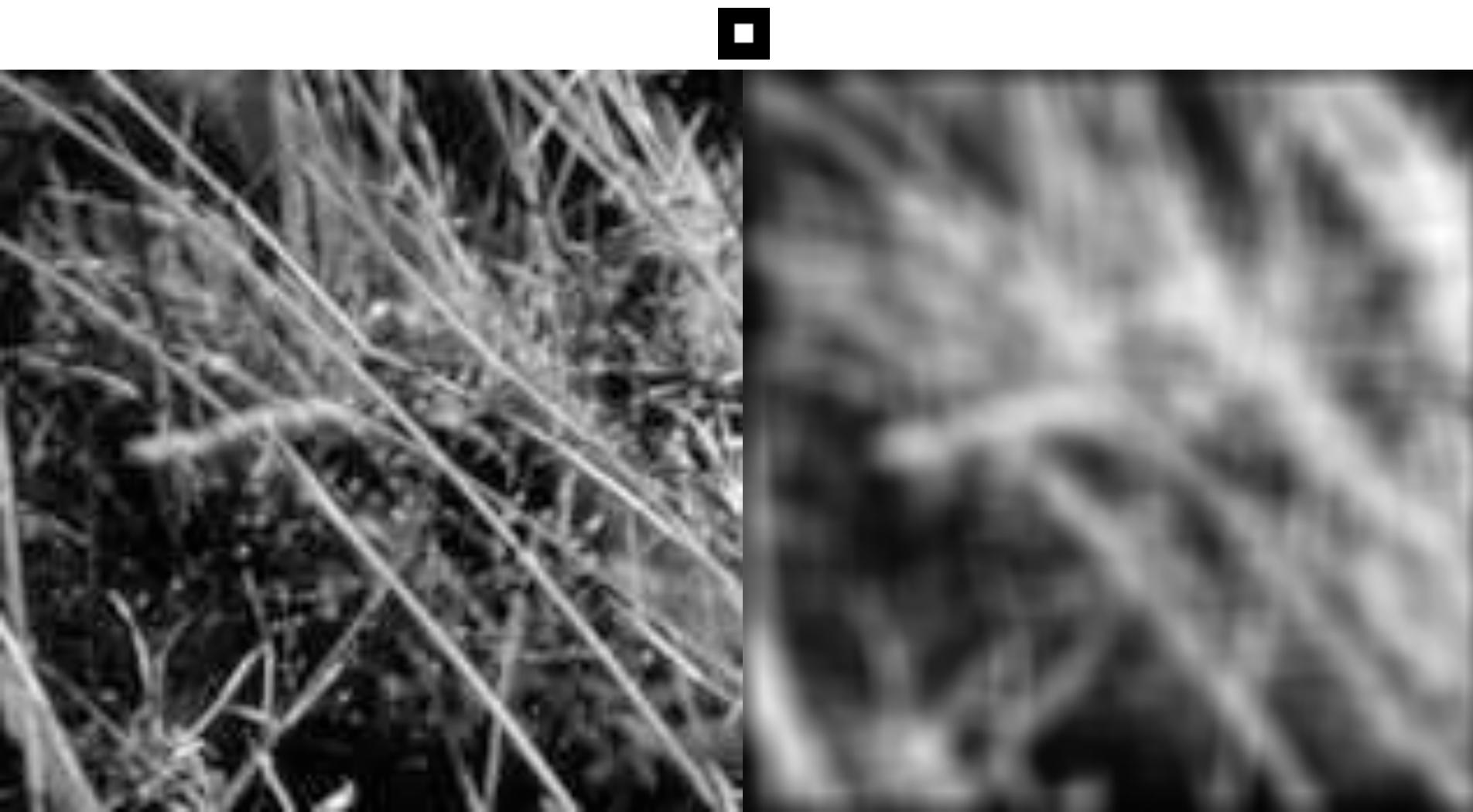
Source: S. Seitz

Box filter

- Replaces each pixel with an average of its neighborhood.
- Achieve smoothing effect (remove sharp features)

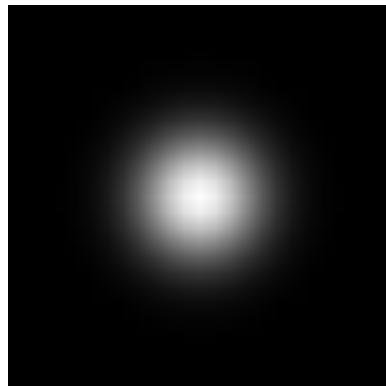
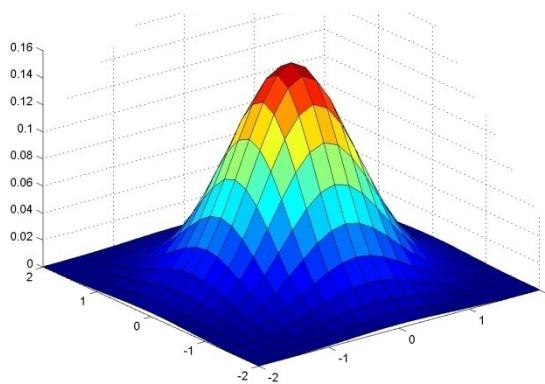
$$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

Example: Smoothing with a box filter



Smoothing with a Gaussian

- Weight contributions of neighboring pixels by nearness



0.003	0.013	0.022	0.013	0.003
0.013	0.059	0.097	0.059	0.013
0.022	0.097	0.159	0.097	0.022
0.013	0.059	0.097	0.059	0.013
0.003	0.013	0.022	0.013	0.003

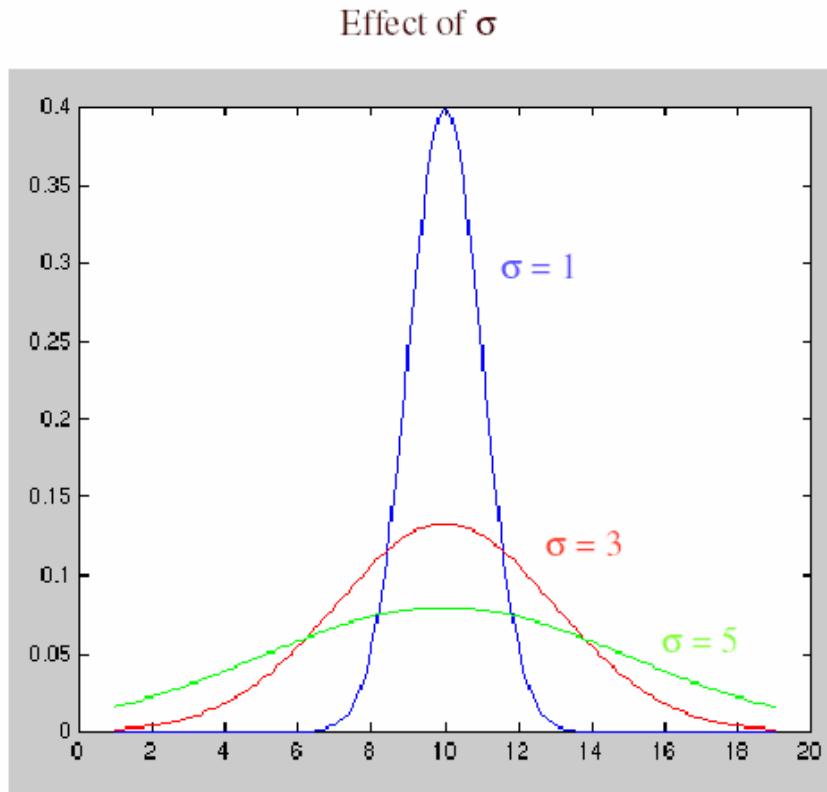
$5 \times 5, \sigma = 1$

$$G_\sigma = \frac{1}{2\pi\sigma^2} e^{-\frac{(x^2+y^2)}{2\sigma^2}}$$

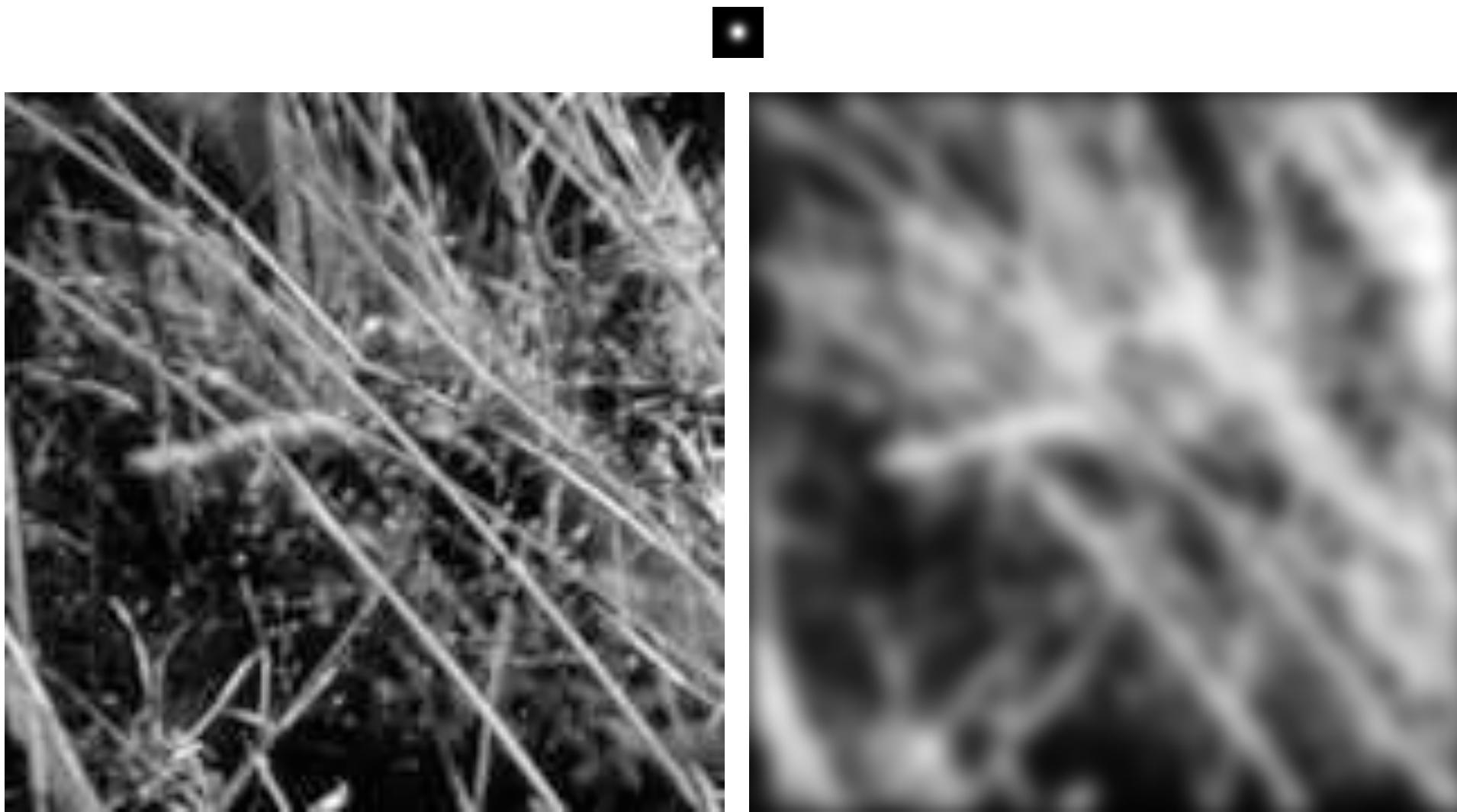
- Constant factor at front makes volume sum to 1 (can be ignored, as we should normalize weights to sum to 1 in any case).

Choosing kernel width

- Rule of thumb: set filter half-width to about 3σ



Smoothing with a Gaussian



Gaussian noise

- Mathematical model: sum of many independent factors
- Assumption: independent, zero-mean noise

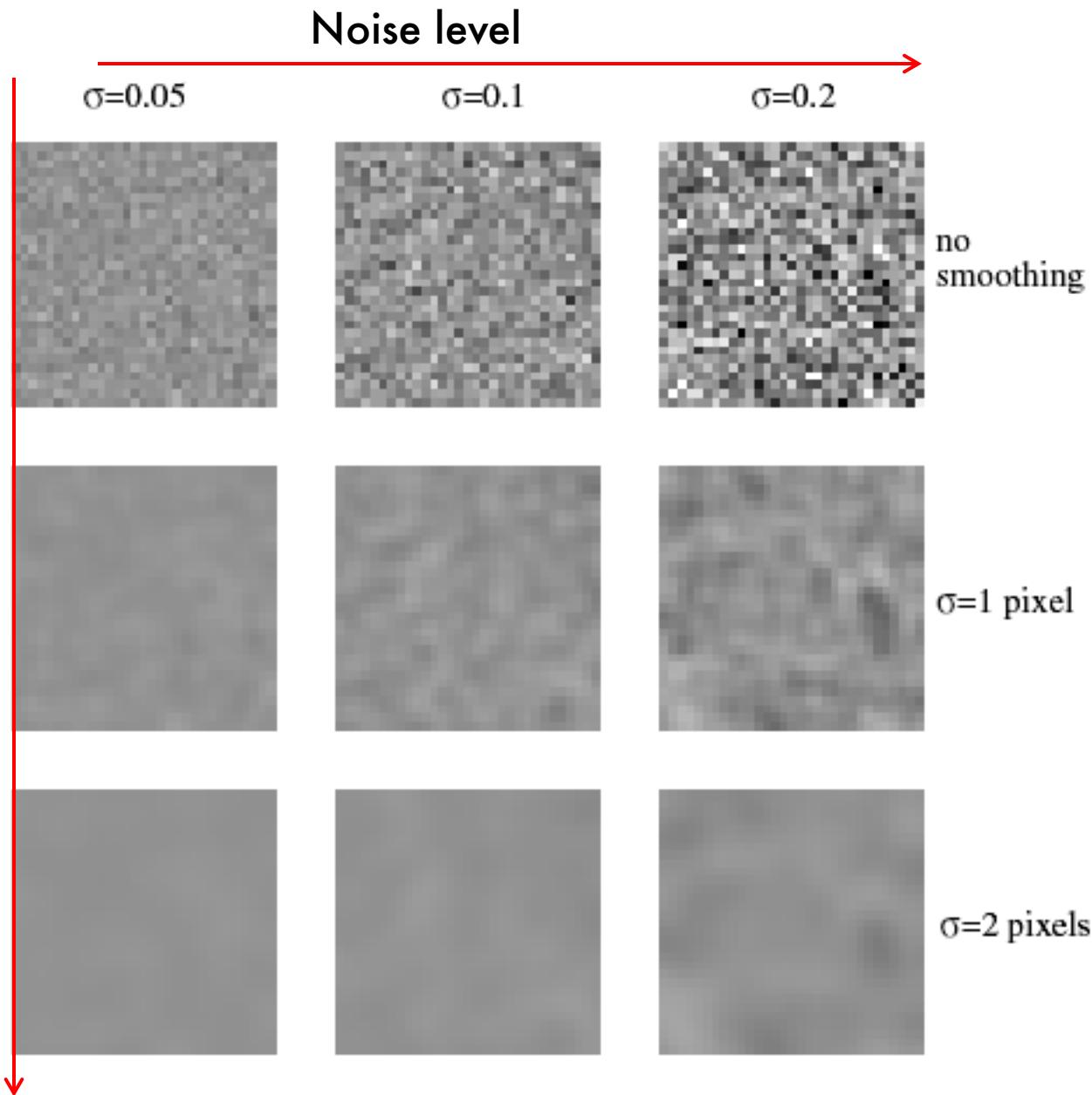
Image
Noise



$$f(x, y) = \overbrace{\hat{f}(x, y)}^{\text{Ideal Image}} + \overbrace{\eta(x, y)}^{\text{Noise process}}$$

Gaussian i.i.d. ("white") noise:
 $\eta(x, y) \sim \mathcal{N}(\mu, \sigma)$

Smoothing with a Gaussian



**Smoothing
reduces pixel
noise:**

Each row shows smoothing with Gaussians of different width; each column shows different amounts of Gaussian noise.

Gaussian filters

- Remove “high-frequency” components from the image (low-pass filter)
- Convolution with self is another Gaussian
 - $F * G(\sigma) * G(\sigma) = F * G(\sqrt{\sigma})$
- *Separable* kernel
 - Factors into product of two 1D Gaussians $G(x,y) = G(x) G(y)$
 - Why this is useful?
can convolve all rows, then all columns

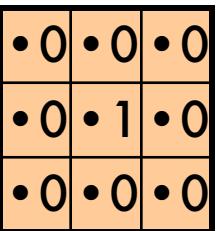
Separability of the Gaussian filter

$$\begin{aligned} G_\sigma(x, y) &= \frac{1}{2\pi\sigma^2} \exp^{-\frac{x^2 + y^2}{2\sigma^2}} \\ &= \left(\frac{1}{\sqrt{2\pi}\sigma} \exp^{-\frac{x^2}{2\sigma^2}} \right) \left(\frac{1}{\sqrt{2\pi}\sigma} \exp^{-\frac{y^2}{2\sigma^2}} \right) \end{aligned}$$

The 2D Gaussian can be expressed as the product of two functions, one a function of x and the other a function of y

In this case, the two functions are the (identical) 1D Gaussian

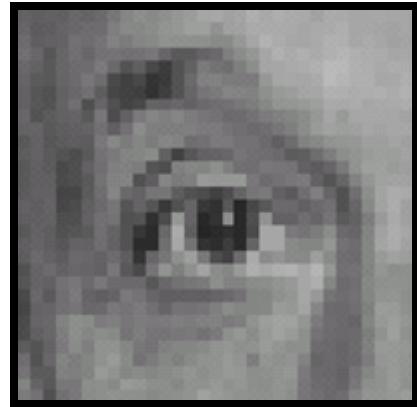
Convolution: Properties

- Commutative: $a * b = b * a$
 - Conceptually no difference between filter and signal
- Associative: $a * (b * c) = (a * b) * c$
 - Often apply several filters one after another: $((a * b_1) * b_2) * b_3$
 - This is equivalent to applying one filter: $a * (b_1 * b_2 * b_3)$
- Distributes over addition: $a * (b + c) = (a * b) + (a * c)$
- Scalars factor out: $ka * b = a * kb = k(a * b)$
- Identity: unit impulse $e =$
 $a * e = a$ 

Convolution: Properties

- **Linearity:** $\text{filter}(f_1 + f_2) = \text{filter}(f_1) + \text{filter}(f_2)$
- **Shift invariance:** $\text{filter}(\text{shift}(f)) = \text{shift}(\text{filter}(f))$
(same behavior regardless of pixel location)
- Theoretical result: any linear shift-invariant operator can be represented as a convolution

Other examples of convolution



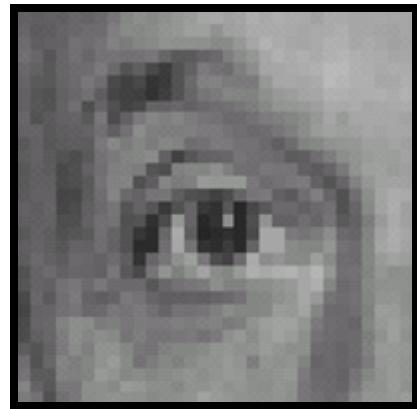
Original

$$\begin{array}{|c|c|c|} \hline \bullet & 0 & 0 \\ \hline 0 & \bullet & 1 \\ \hline 0 & 0 & 0 \\ \hline \end{array}$$

=

?

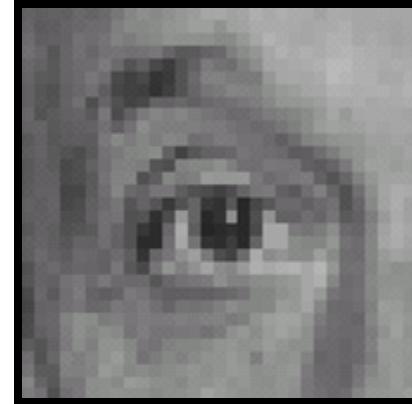
Other examples of convolution



Original

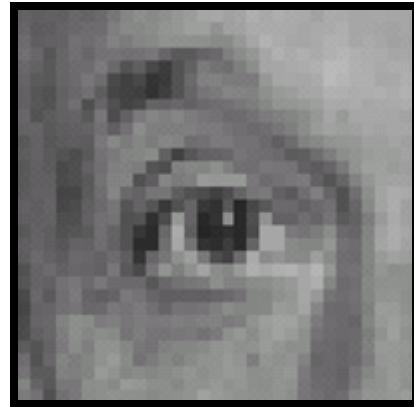
$$\begin{array}{|c|c|c|} \hline \bullet & 0 & \bullet & 0 \\ \hline 0 & \bullet & 1 & 0 \\ \hline \bullet & 0 & \bullet & 0 \\ \hline \end{array}$$

=



Filtered
(no change)

Other examples of convolution



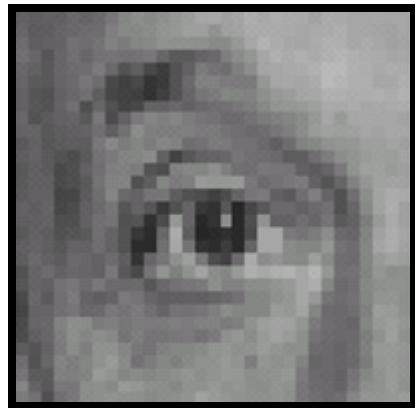
Original

$$\begin{array}{|c|c|c|} \hline \bullet & 0 & 0 \\ \hline 0 & \bullet & 0 \\ \hline 0 & 0 & \bullet \\ \hline \end{array}$$

=

?

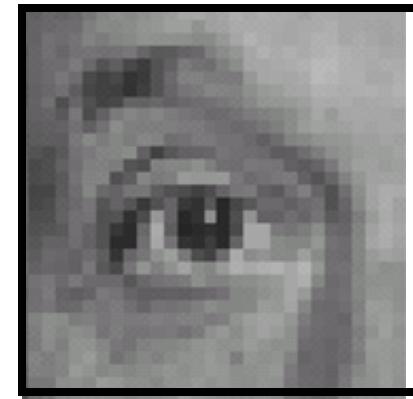
Practice with linear filters



Original

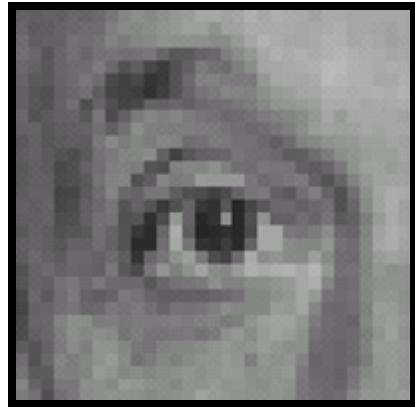
$$\begin{array}{|c|c|c|} \hline \bullet & 0 & 0 \\ \hline 0 & \bullet & 0 \\ \hline 0 & 0 & \bullet \\ \hline \end{array}$$

=



Shifted left
By 1 pixel

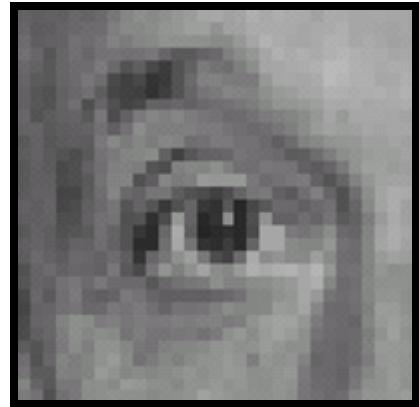
Other examples of convolution



Original

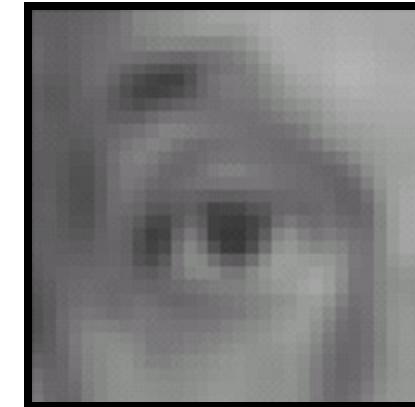
$$\frac{1}{9} \begin{array}{|c|c|c|} \hline \bullet & 1 & \bullet & 1 & \bullet & 1 \\ \hline \bullet & 1 & \bullet & 1 & \bullet & 1 \\ \hline \bullet & 1 & \bullet & 1 & \bullet & 1 \\ \hline \end{array} = ?$$

Other examples of convolution



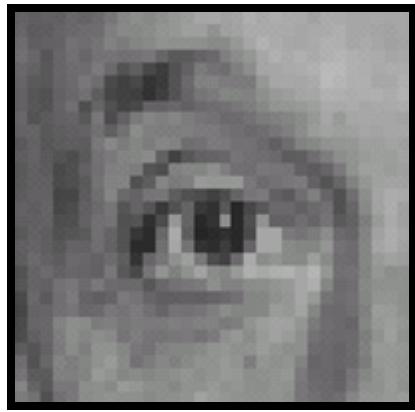
Original

$$\frac{1}{9} \begin{array}{|c|c|c|} \hline \bullet & 1 & \bullet & 1 & \bullet & 1 \\ \hline \bullet & 1 & \bullet & 1 & \bullet & 1 \\ \hline \bullet & 1 & \bullet & 1 & \bullet & 1 \\ \hline \end{array} =$$



Blur (with a
box filter)

Other examples of convolution



$$\begin{array}{|c|c|c|} \hline \bullet & 0 & 0 \\ \hline 0 & \bullet & 2 \\ \hline 0 & 0 & \bullet \\ \hline \end{array}$$

$$- \frac{1}{9} \begin{array}{|c|c|c|} \hline \bullet & 1 & 1 \\ \hline 1 & \bullet & 1 \\ \hline 1 & 1 & \bullet \\ \hline \end{array}$$

= ?

(Note that filter sums to 1)

Original

$$\begin{array}{|c|c|c|} \hline \bullet & 0 & 0 \\ \hline 0 & \bullet & 1 \\ \hline 0 & 0 & \bullet \\ \hline \end{array}$$

+

$$\begin{array}{|c|c|c|} \hline \bullet & 0 & 0 \\ \hline 0 & \bullet & 1 \\ \hline 0 & 0 & \bullet \\ \hline \end{array}$$

-

$$\frac{1}{9} \begin{array}{|c|c|c|} \hline \bullet & 1 & 1 \\ \hline 1 & \bullet & 1 \\ \hline 1 & 1 & \bullet \\ \hline \end{array}$$

details

Source: D. Lowe

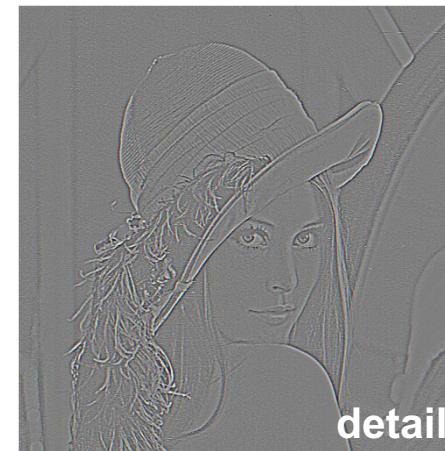
- What does blurring take away?



-



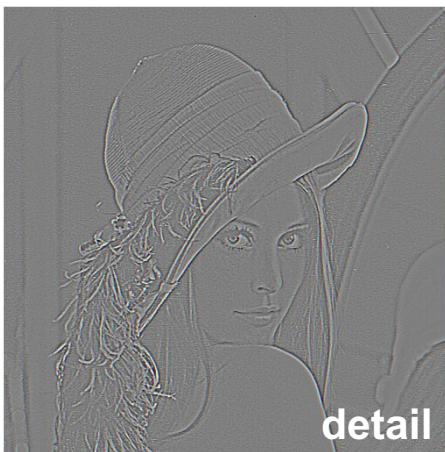
=



- Let's add it back:



+ a



=



Other examples of convolution



$$\begin{array}{|c|c|c|} \hline \bullet & 0 & 0 \\ \hline 0 & \bullet & 2 \\ \hline 0 & 0 & \bullet \\ \hline \end{array}$$

$$- \frac{1}{9} \begin{array}{|c|c|c|} \hline \bullet & 1 & 1 \\ \hline 1 & \bullet & 1 \\ \hline 1 & 1 & \bullet \\ \hline \end{array} =$$

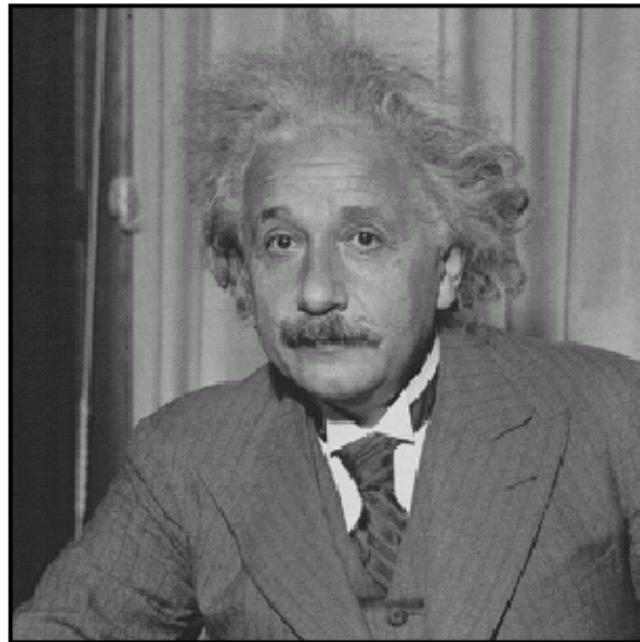


Original

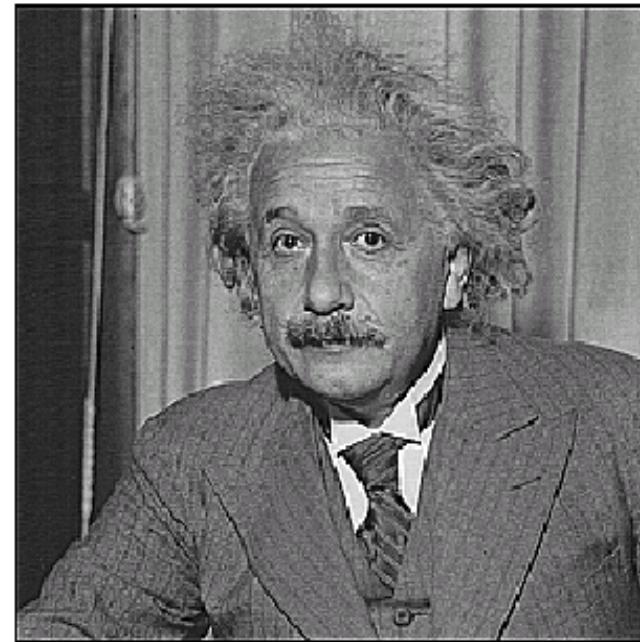
Sharpening filter

- Accentuates differences with local average

Sharpening

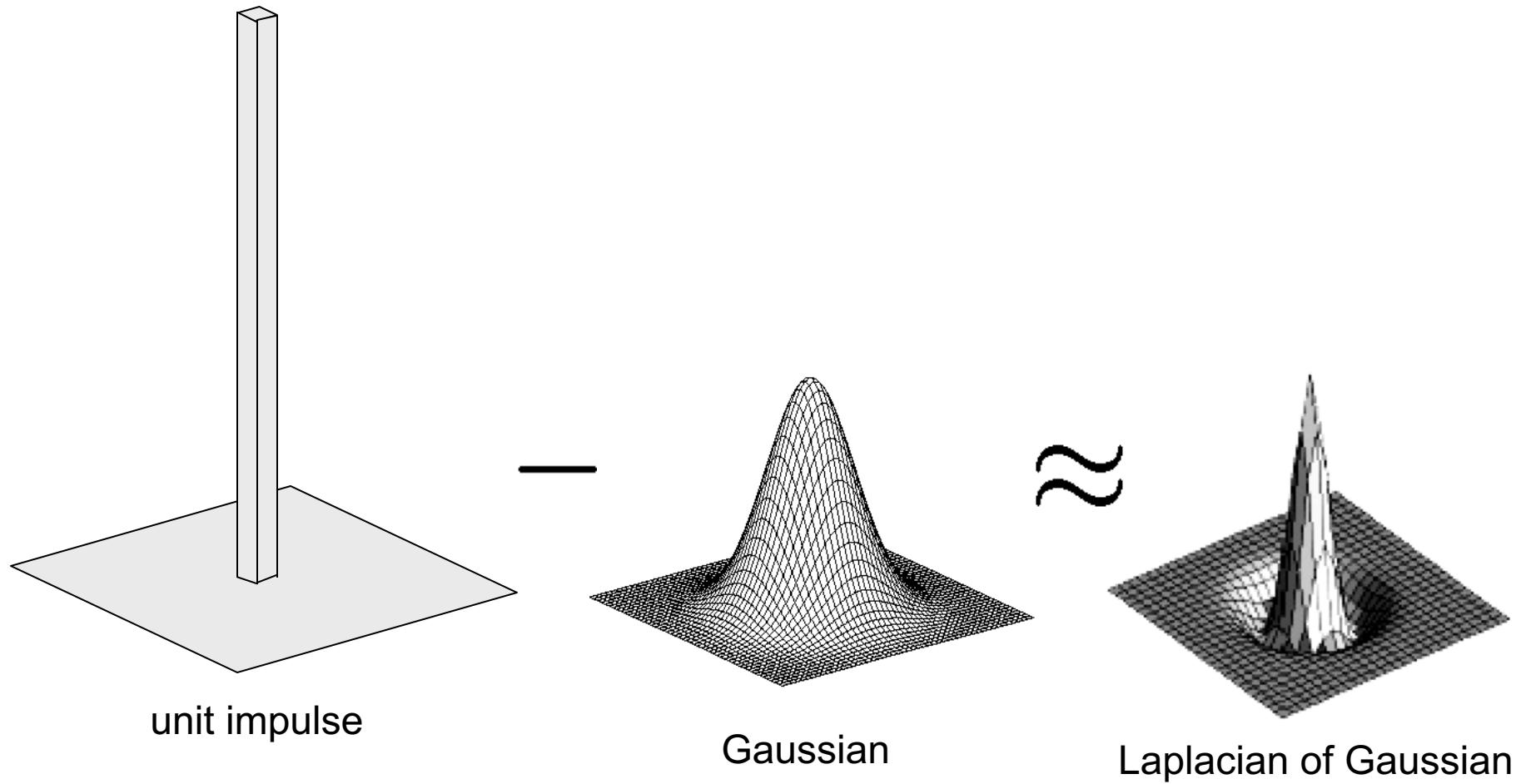


before



after

Sharpening



$$f * ((1 + \alpha) e^{-\frac{x^2}{2}} - g)$$

Impulse function

Differentiation by convolution

- Recall, for 2D function, $f(x,y)$:
- We could approximate this as

$$\frac{\partial f}{\partial x} = \lim_{\varepsilon \rightarrow 0} \left(\frac{f(x + \varepsilon, y) - f(x, y)}{\varepsilon} \right)$$

$$\frac{\partial f}{\partial x} \approx \frac{f(x_{n+1}, y) - f(x_n, y)}{\Delta x}$$

- This is linear and shift invariant \rightarrow convolution

2D Kernel

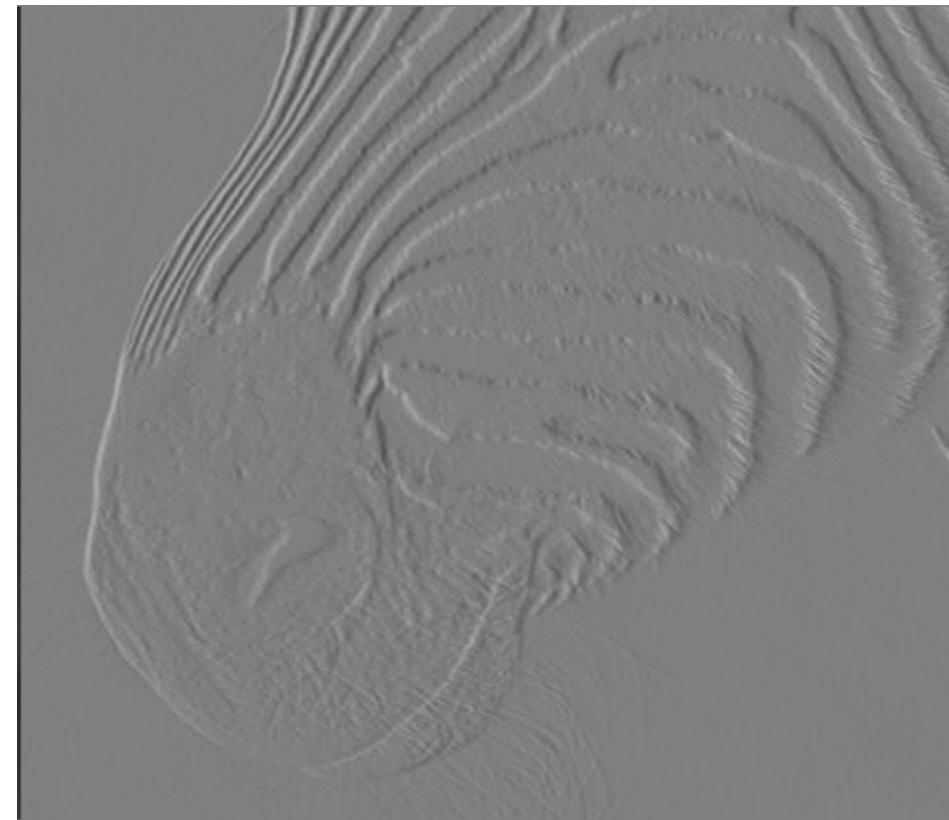
$$\begin{matrix} -1 & 1 \\ -1 & 1 \end{matrix}$$

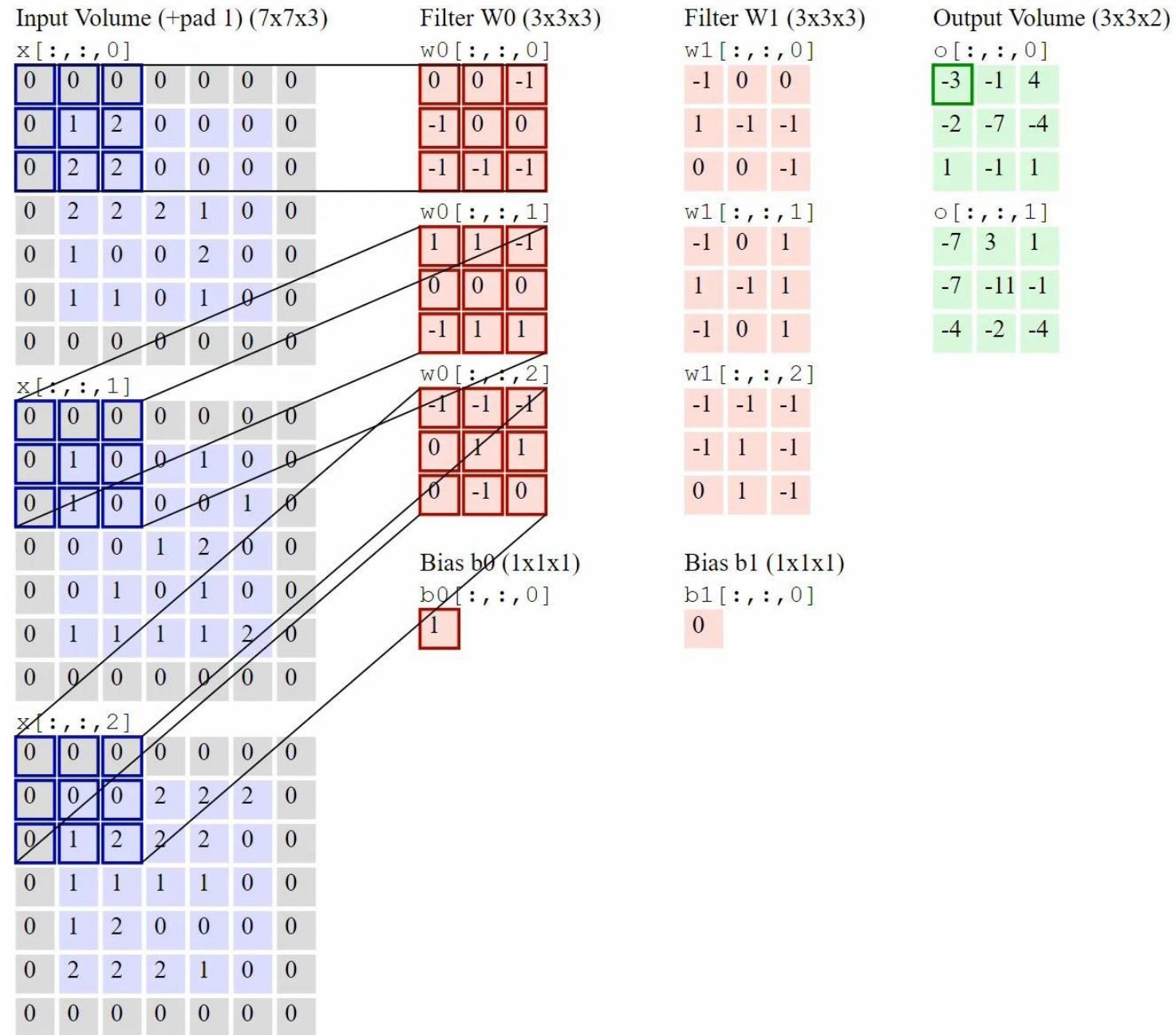
Rudimentary edge detector!

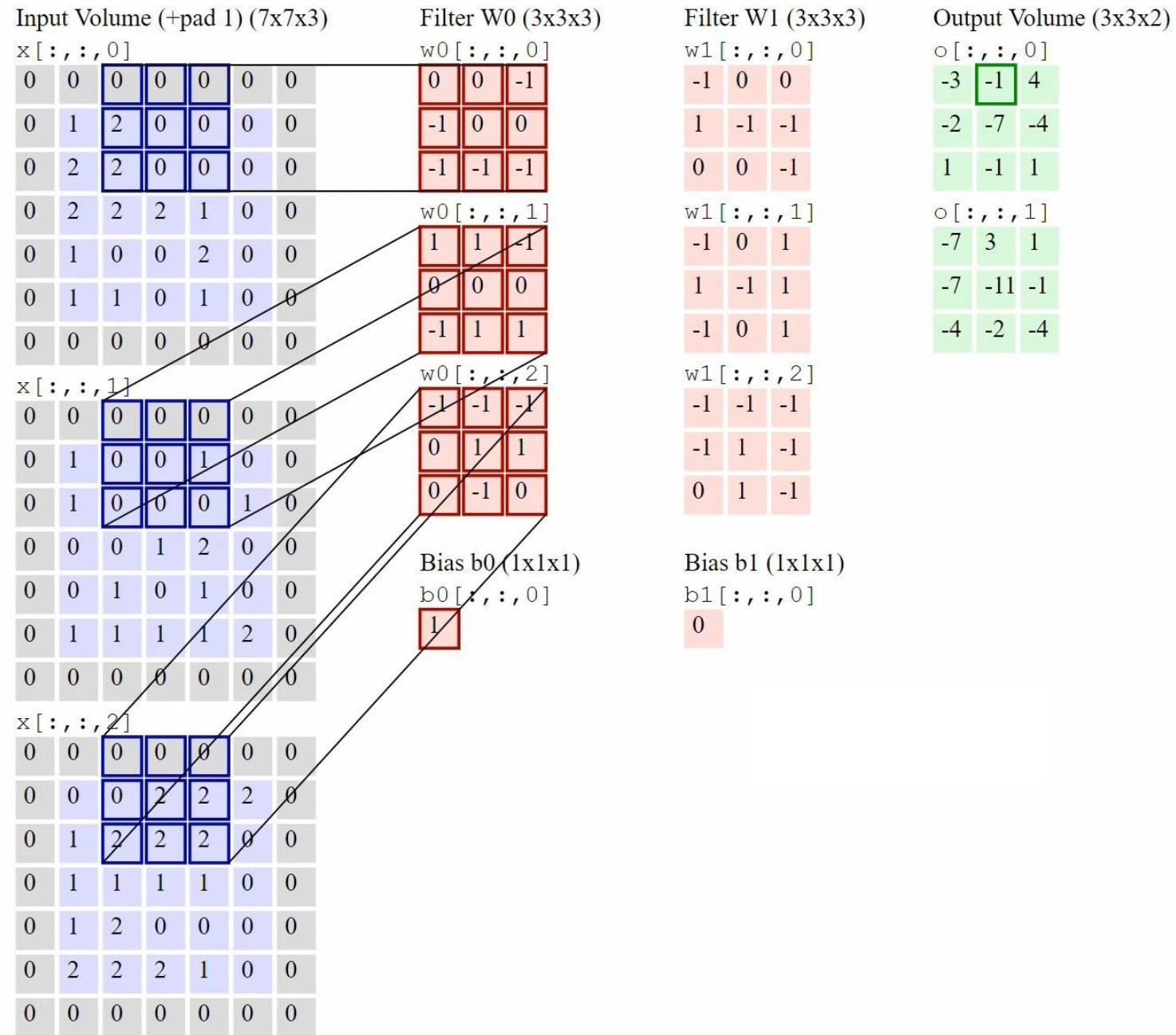
Differentiation by convolution

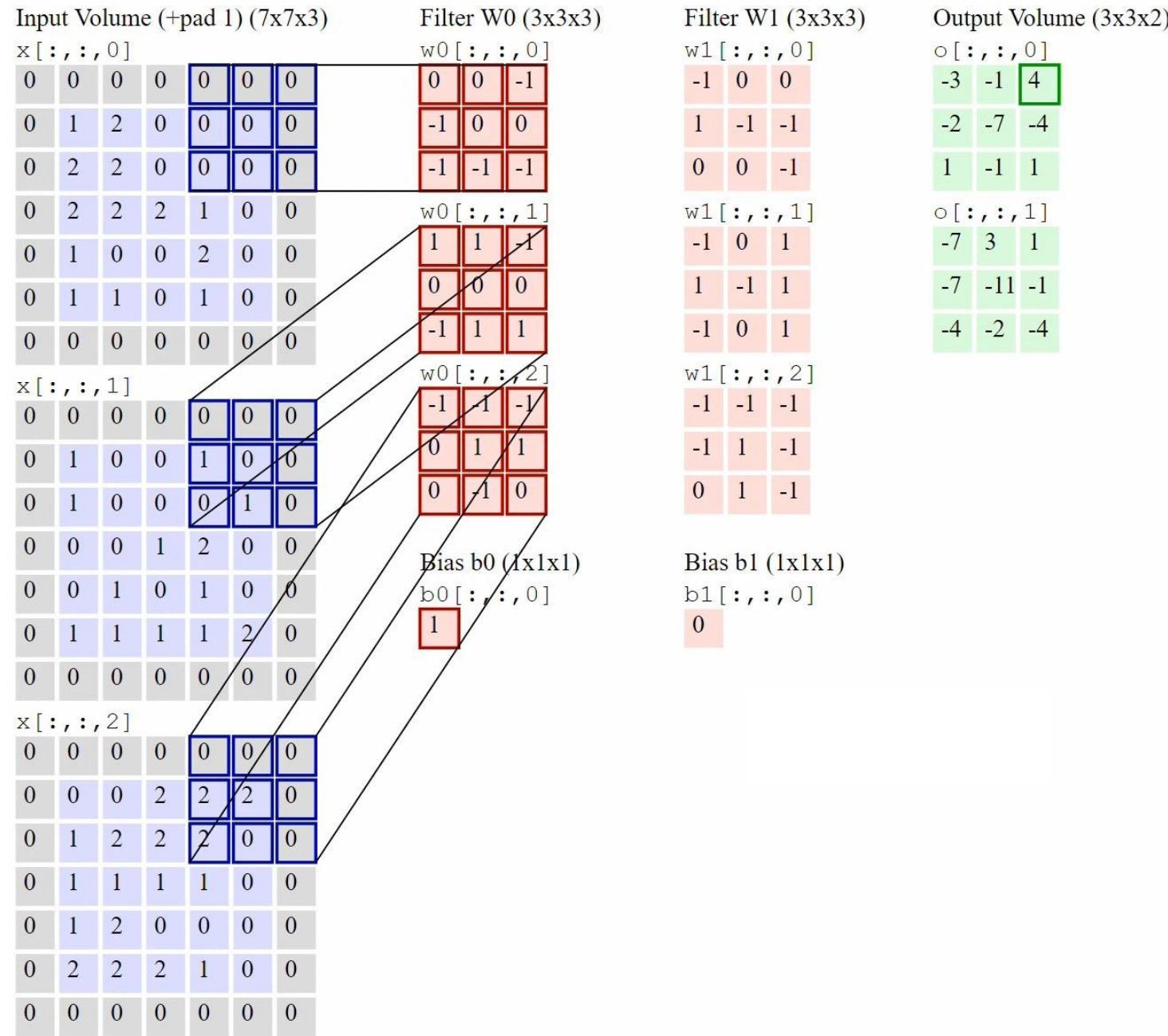
Directional Derivatives Of A Binary Image			
Original Image			
2D Kernel	2D Kernel	2D Kernel	2D Kernel
0 0 0	-1 0 1	-2 -1 0 1 2	-3 -2 -1 0 1 2 3
0 1 0	-1 0 1	-2 -1 0 1 2	-3 -2 -1 0 1 2 3
0 0 0	-1 0 1	-2 -1 0 1 2	-3 -2 -1 0 1 2 3
		-2 -1 0 1 2	-3 -2 -1 0 1 2 3
		-2 -1 0 1 2	-3 -2 -1 0 1 2 3
		-2 -1 0 1 2	-3 -2 -1 0 1 2 3

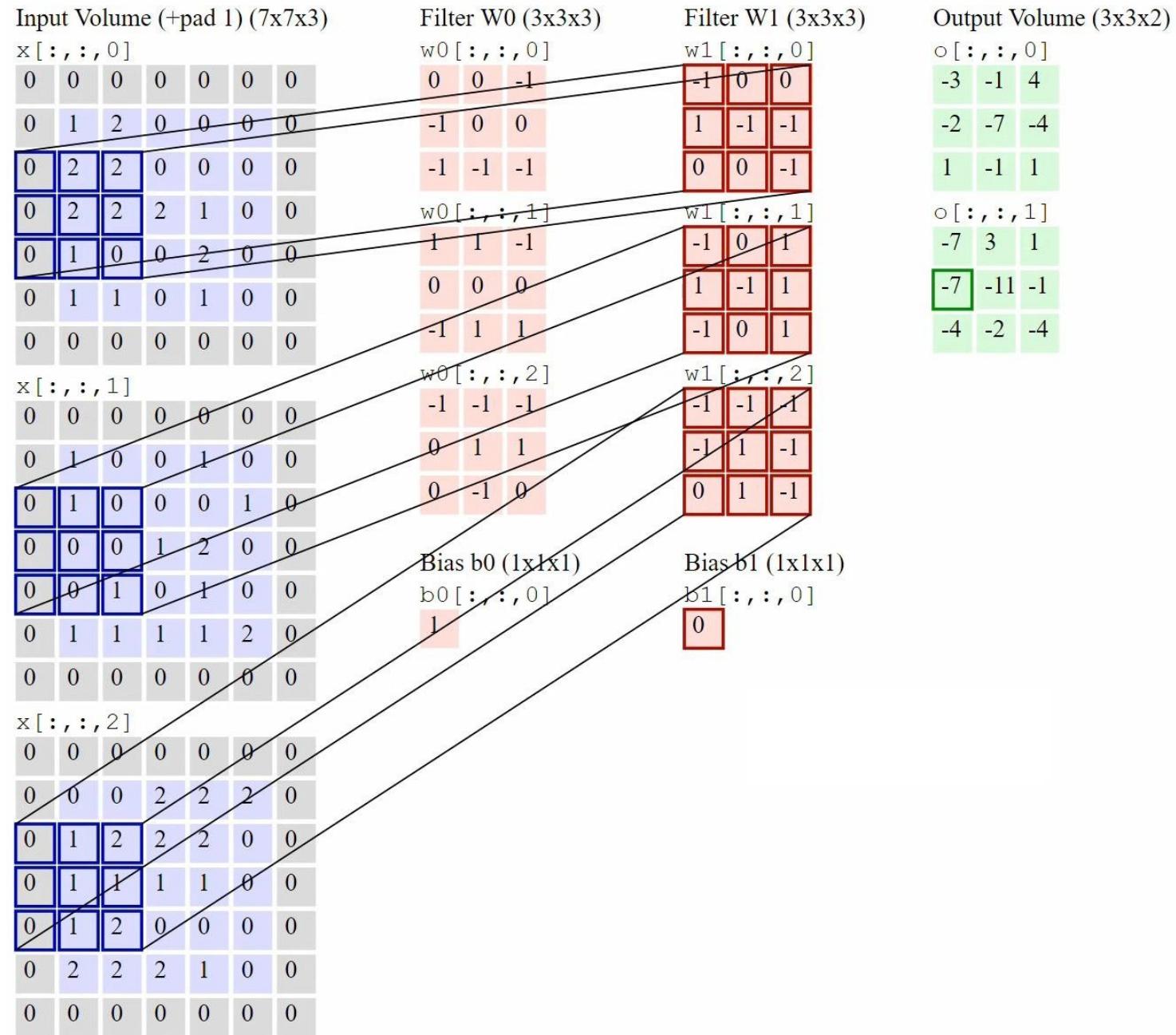
Differentiation by convolution

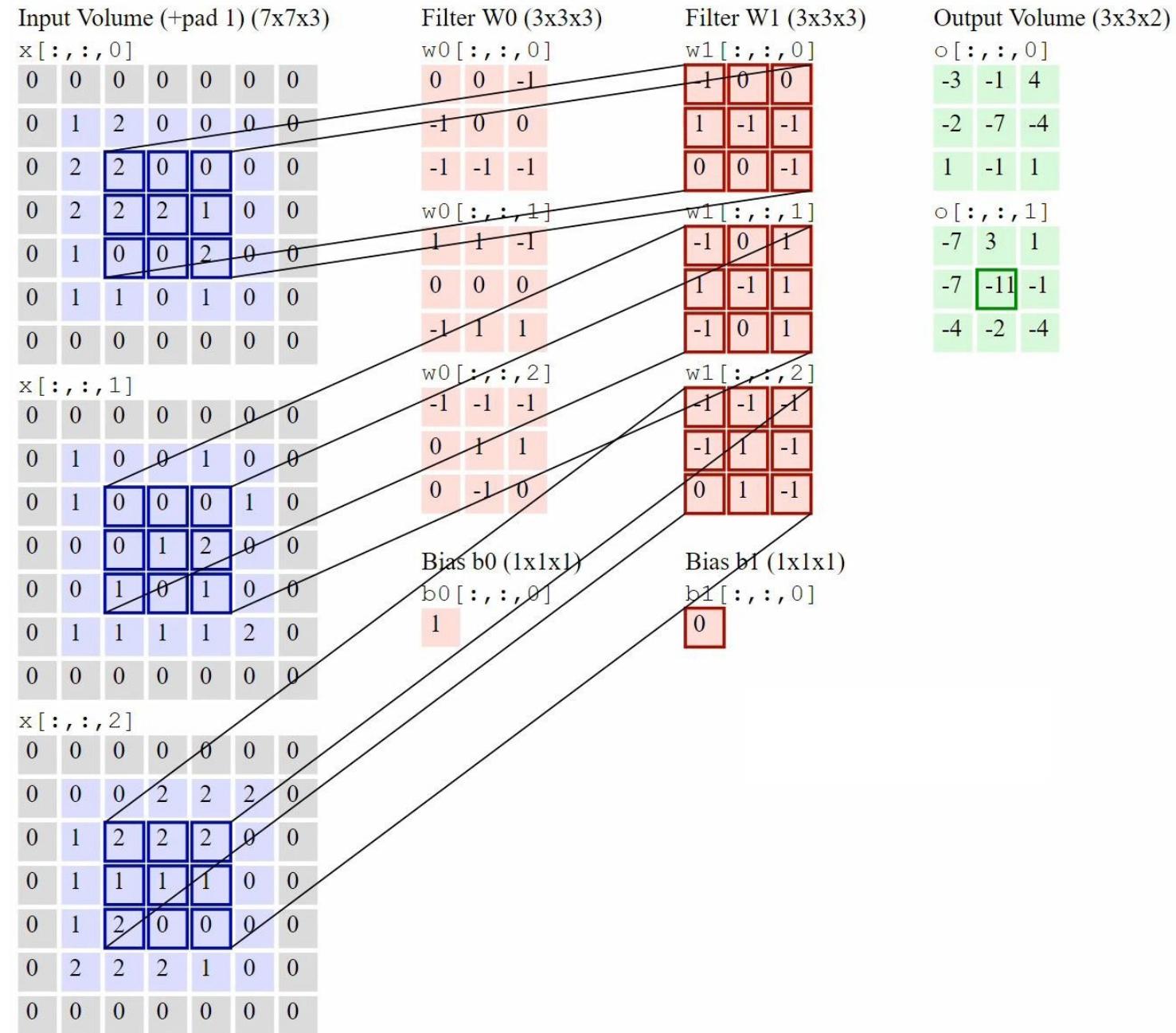


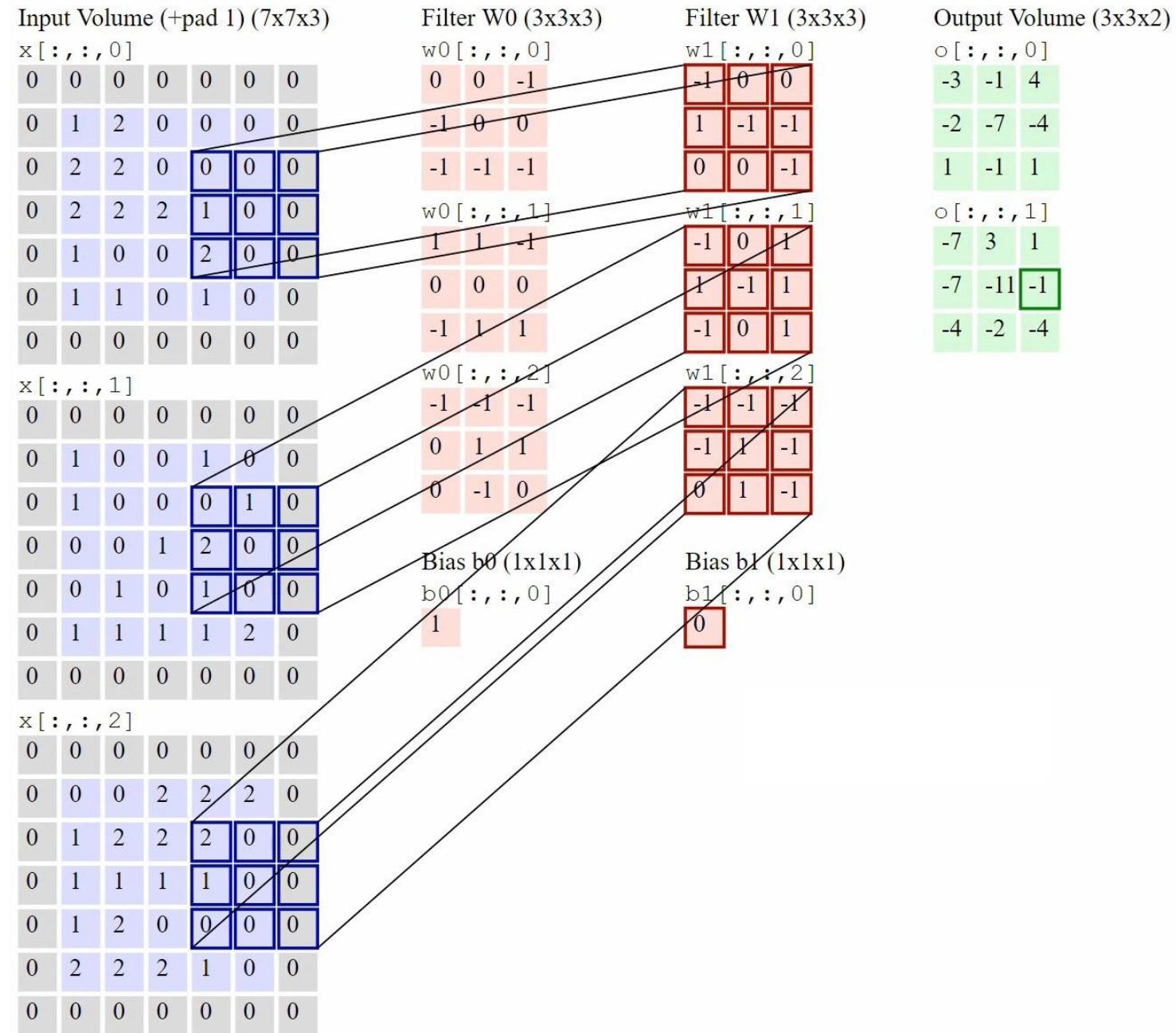


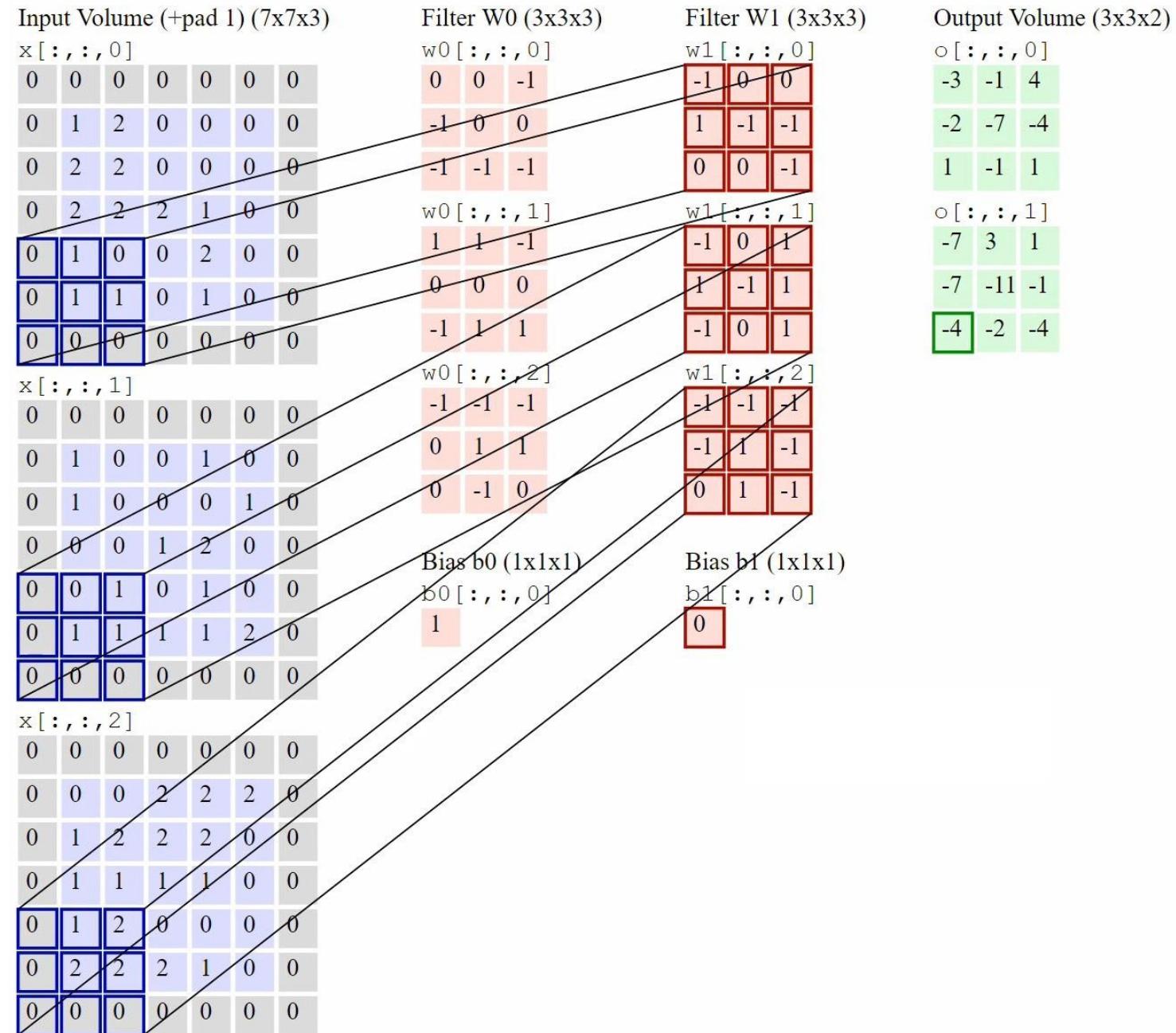


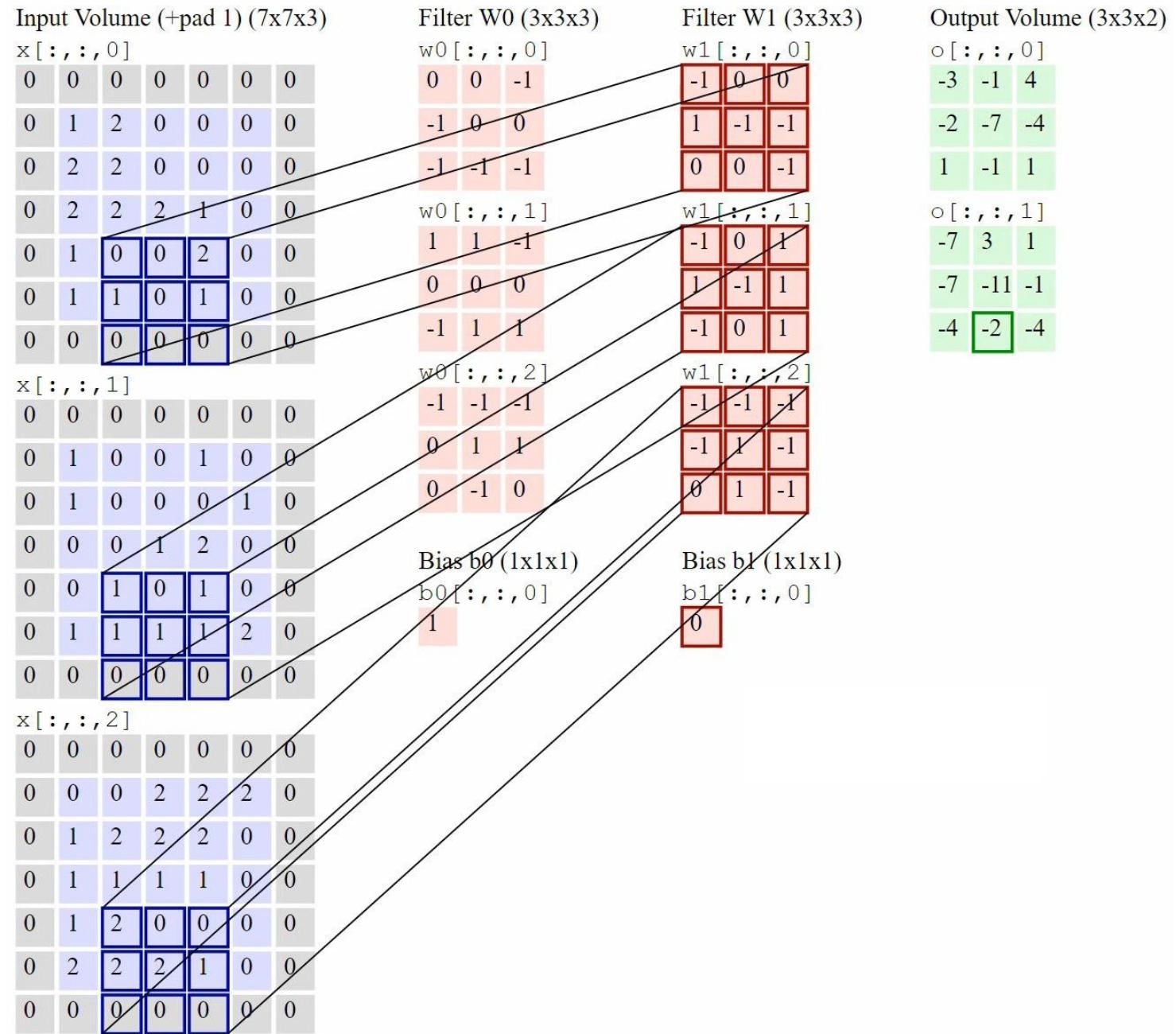


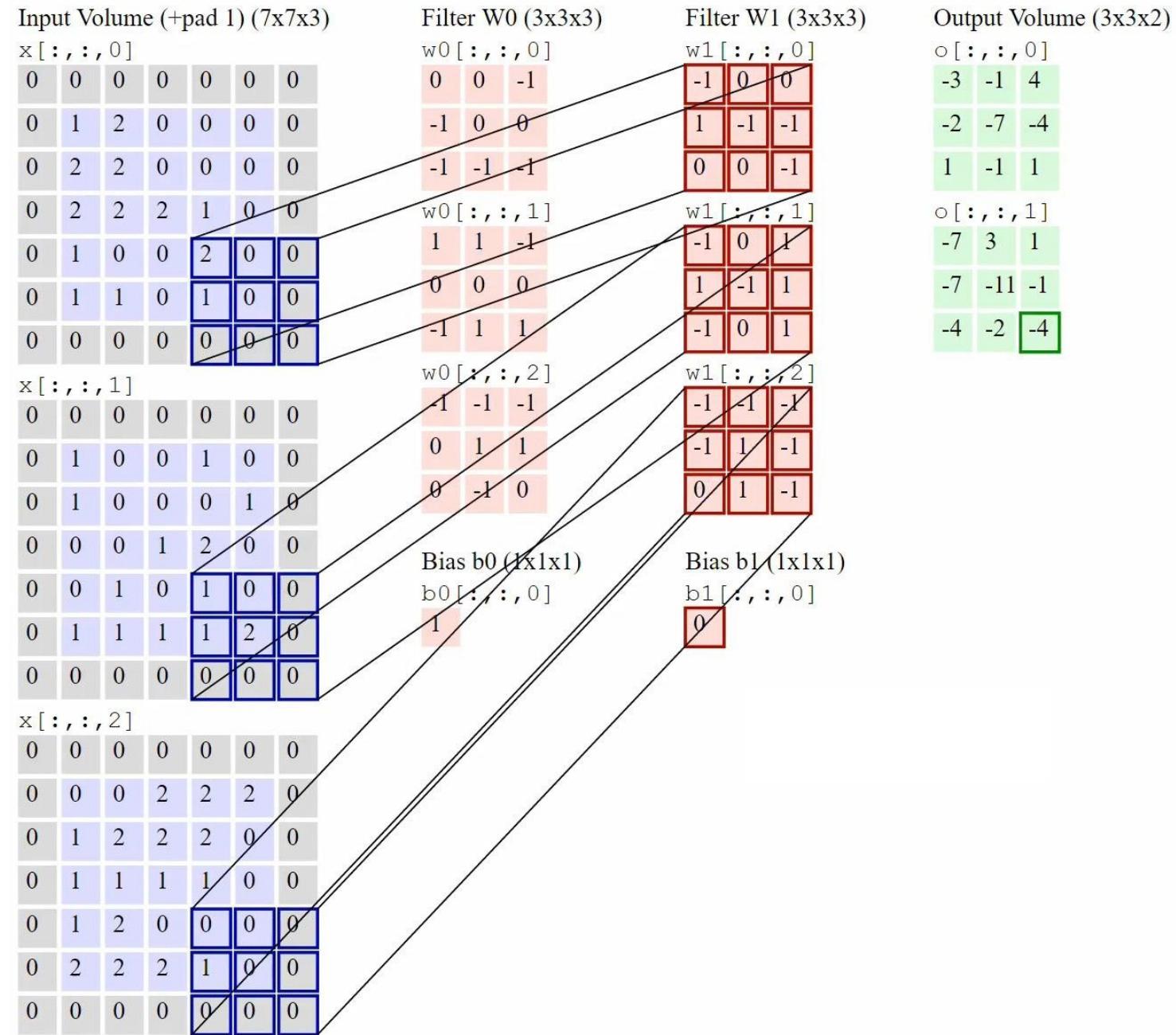




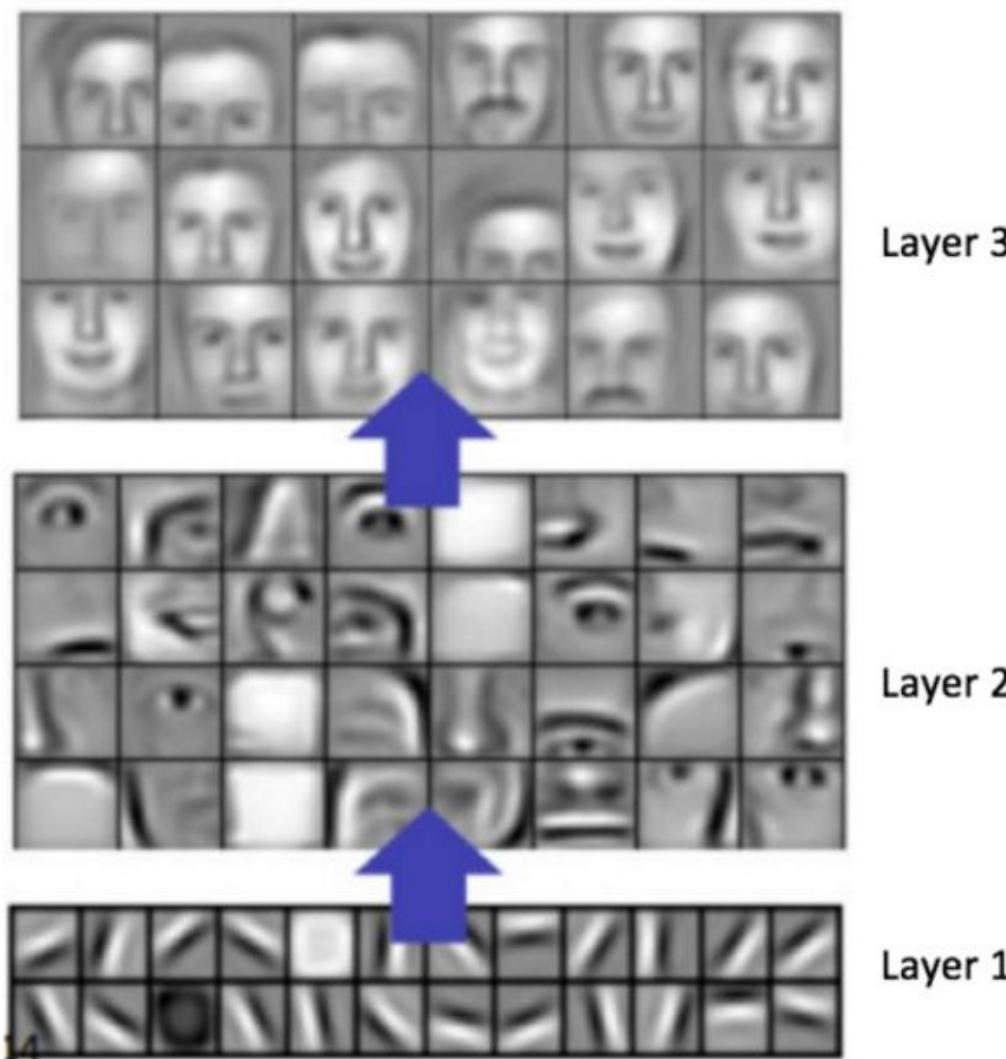




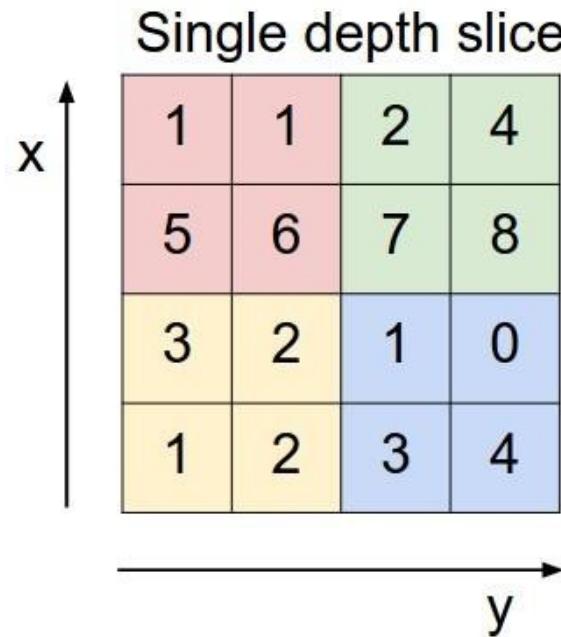




Convolution: Representation Learning



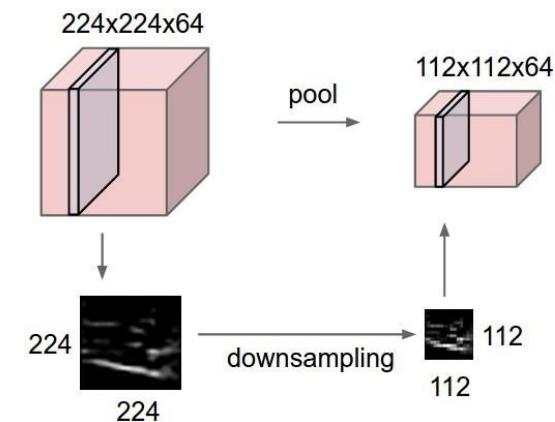
ConvNets: Pooling



max pool with 2x2 filters
and stride 2

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

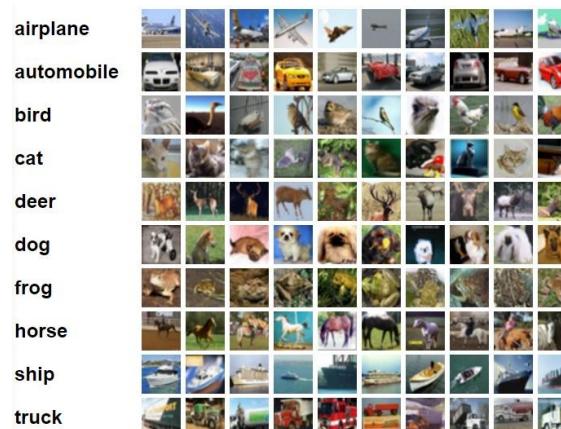
6	8
3	4



Famous Computer Vision Datasets



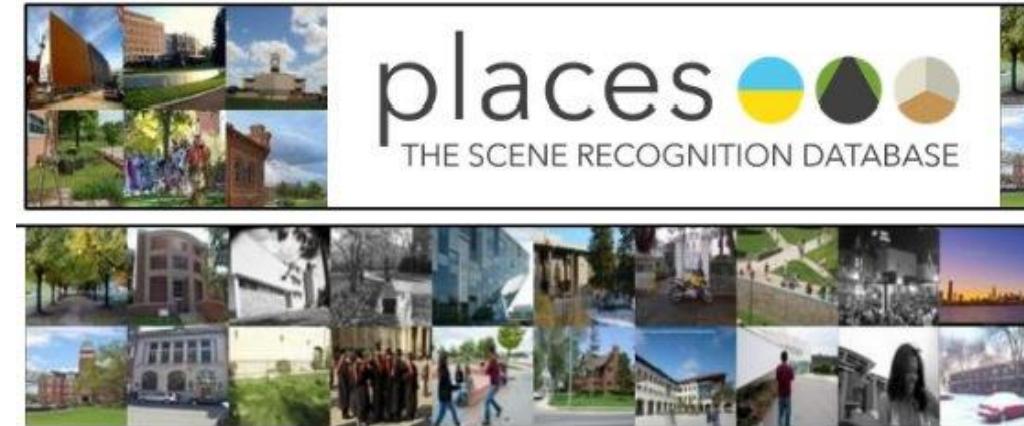
MNIST: handwritten digits



CIFAR-10(0): tiny images

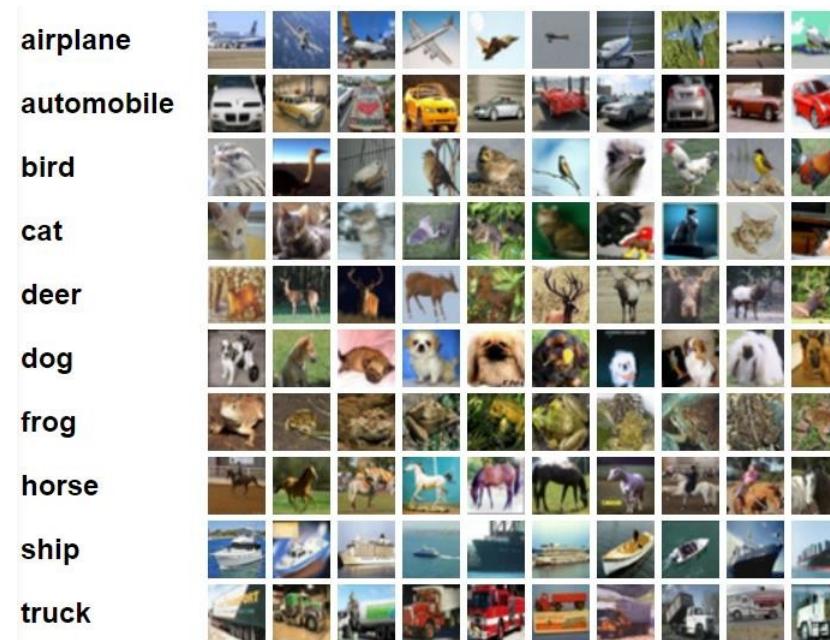
printer housing animal weight drop headquarters egg white
offspring teacher computer album garage daffodil flower television
register measure gallery court key structure light date spread breakfast
king fireplace horse church press market lighter
restaurant counter hotel road Paper cup side site door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

ImageNet: WordNet hierarchy



Places: natural scenes

Let's Build an Image Classifier for CIFAR-10



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

Let's Build an Image Classifier for CIFAR-10

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



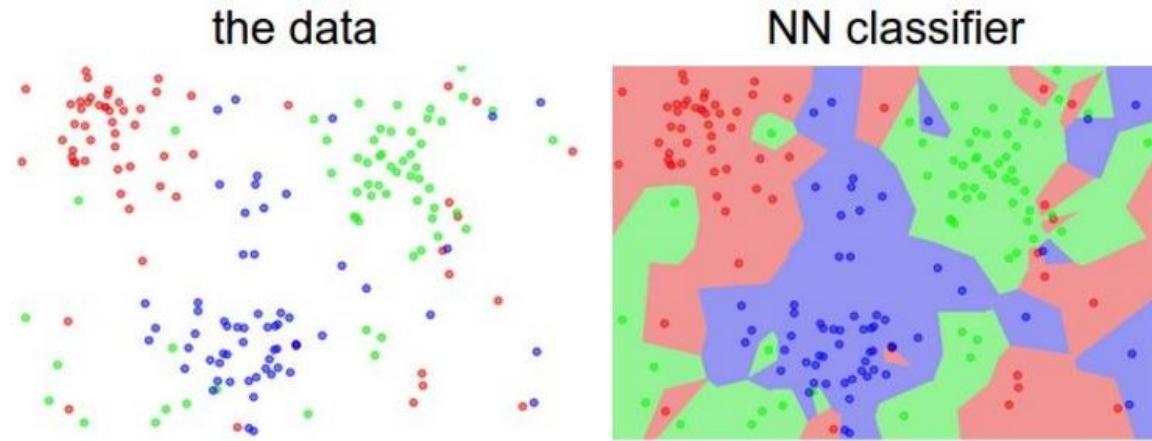
Accuracy

Random: **10%**

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

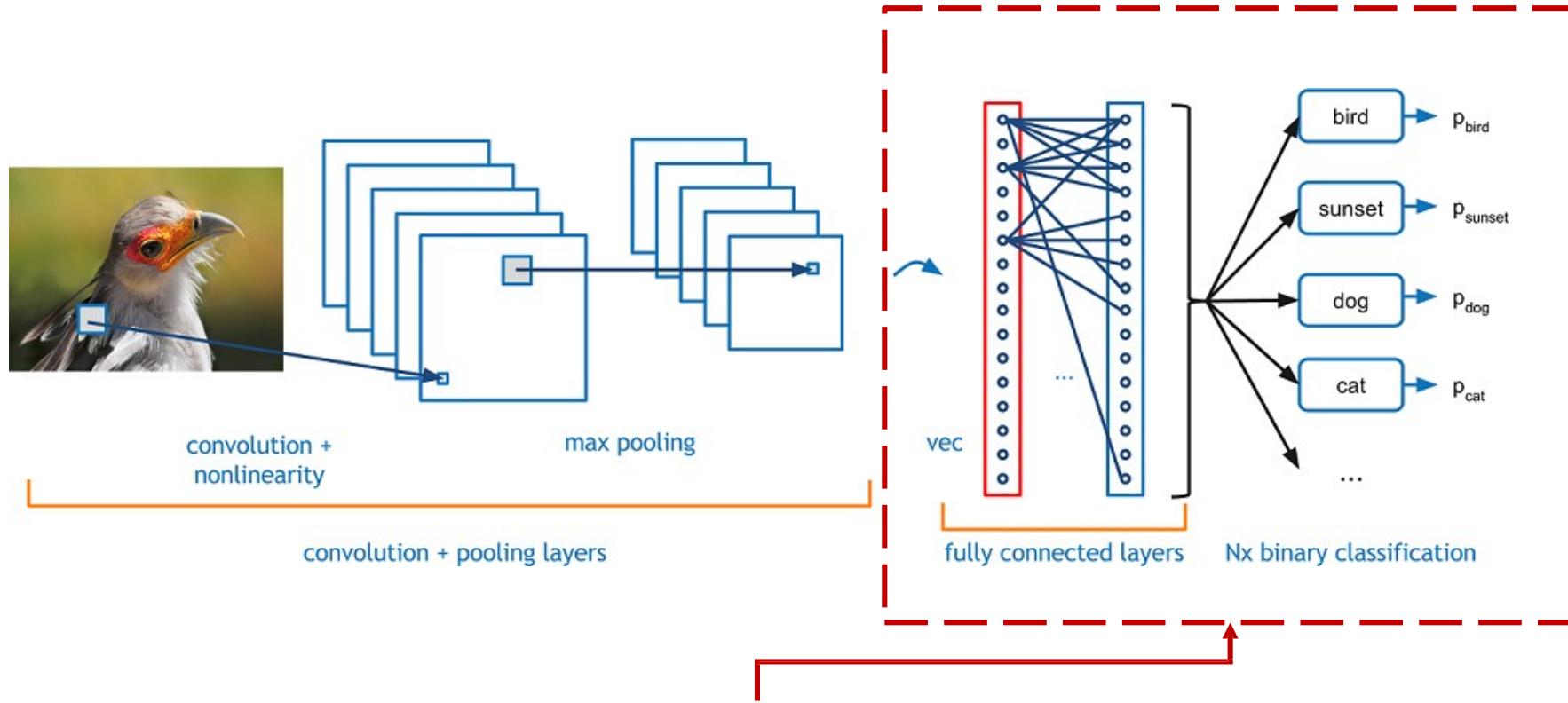
Our image-diff (with L2): **35.4%**

K-Nearest Neighbors: Generalizing the Image-Diff Classifier



- Accuracy
- Random: 10%
- Training and testing on the same data: 35.4%
- 7-Nearest Neighbors: ~30%
- Human: ~94%
- ...
- Convolutional Neural Networks: ~95%

Same Architecture, Many Applications

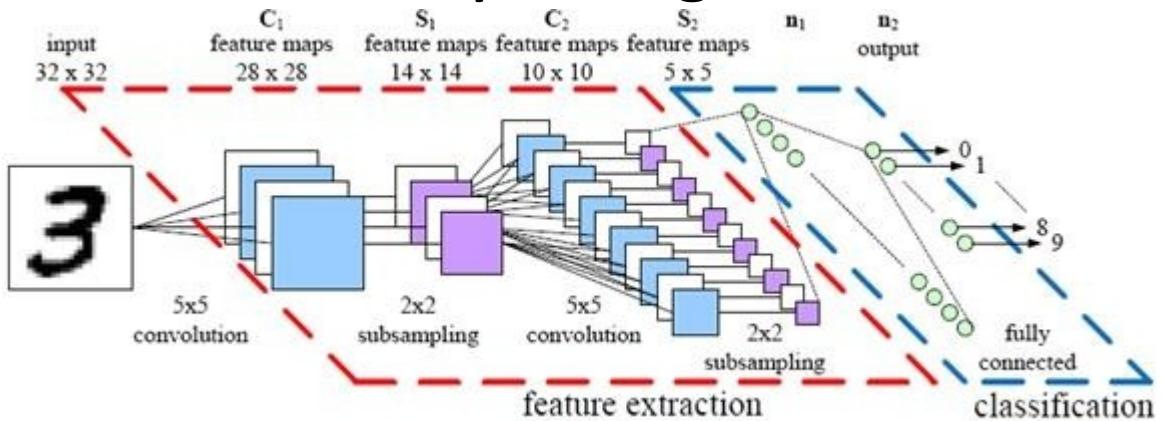


This part might look different for:

- 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

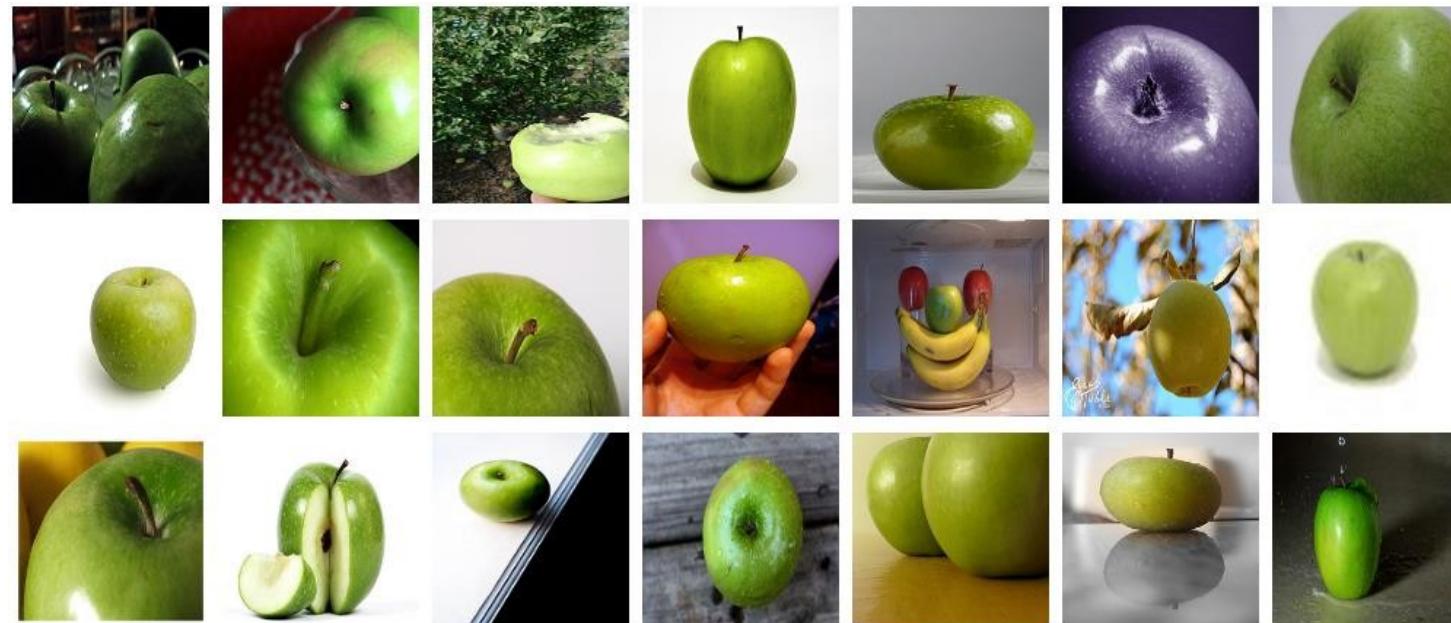
Case Study: ImageNet



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

What is ImageNet?

- **ImageNet**: dataset of 14+ million images (21,841 categories)
 - Links to images not images
- 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?

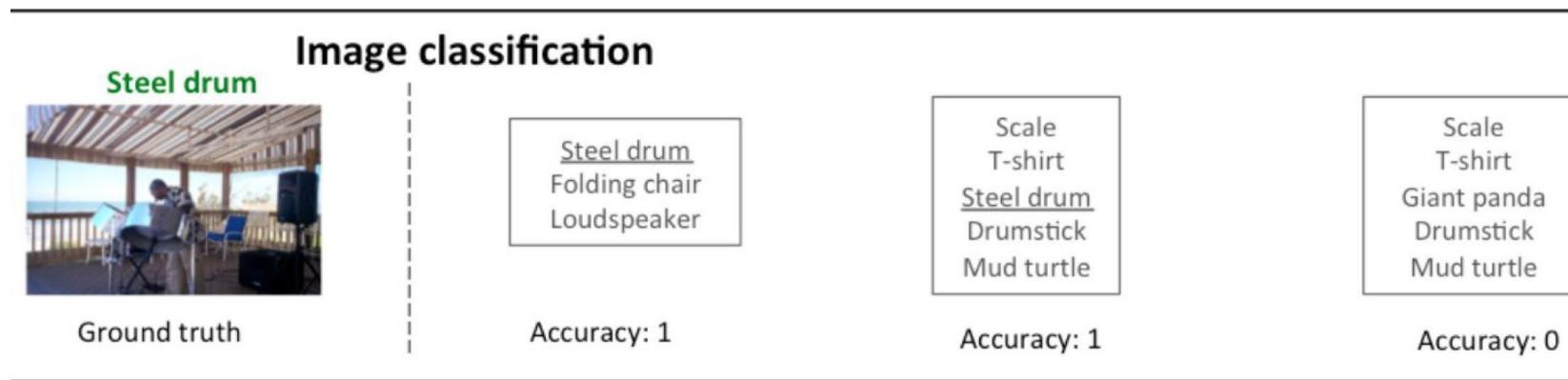
Dataset → • **ImageNet**: dataset of 14+ million images

Competition → • **ILSVRC**: ImageNet Large Scale Visual Recognition Challenge

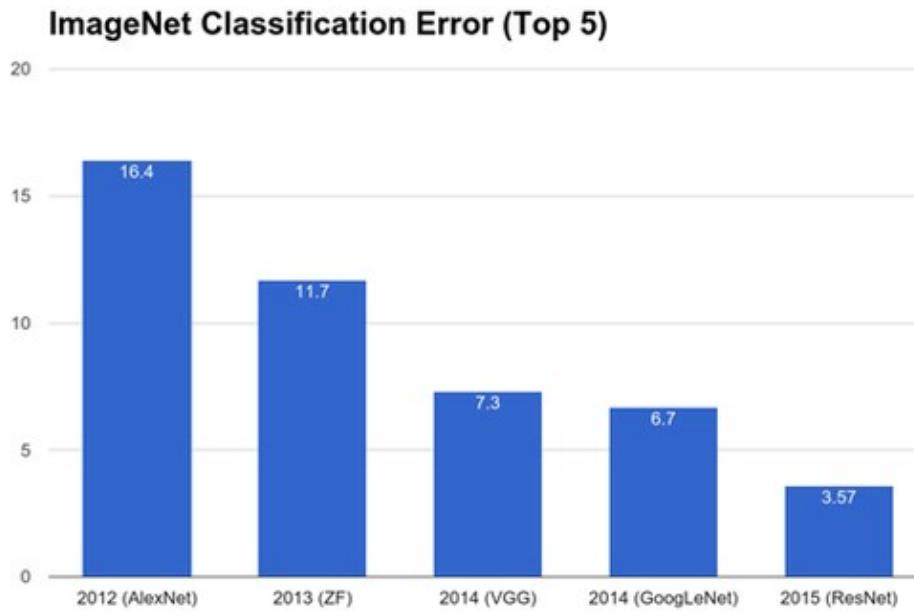
Networks → • AlexNet (2012)
• ZFNet (2013)
• VGGNet (2014)
• GoogLeNet (2014)
• ResNet (2015)
• CUIImage (2016)

ILSVRC Challenge Evaluation for Classification

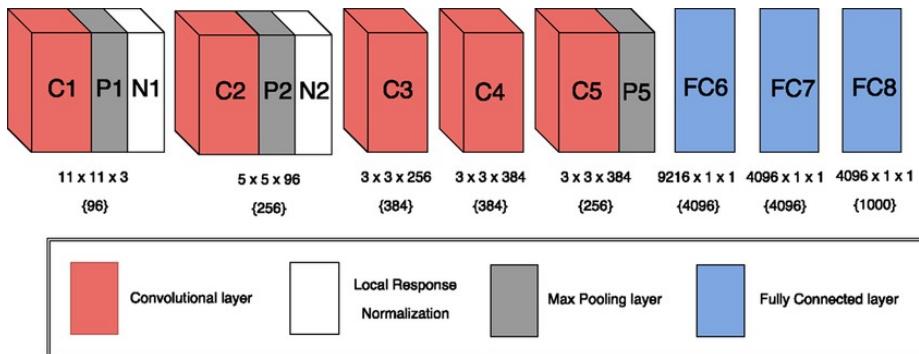
- Top 5 error rate:
 - You get 5 guesses to get the correct label



- ~20% reduction in accuracy for Top 1 vs Top 5
 - Example: In 2012 AlexNet achieved
- Human annotation is a binary task: “apple” or “not apple”

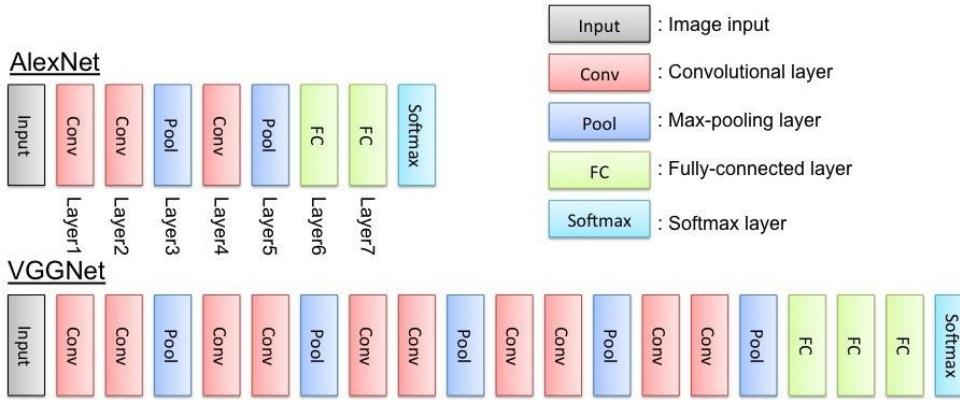


- **AlexNet (2012): First CNN (15.4%)**
 - 8 layers
 - 61 million parameters
- **ZFNet (2013): 15.4% to 11.2%**
 - 8 layers
 - More filters. Denser stride.
- **VGGNet (2014): 11.2% to 7.3%**
 - Beautifully uniform:
3x3 conv, stride 1, pad 1, 2x2 max pool
 - 16 layers
 - 138 million parameters
- **GoogLeNet (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



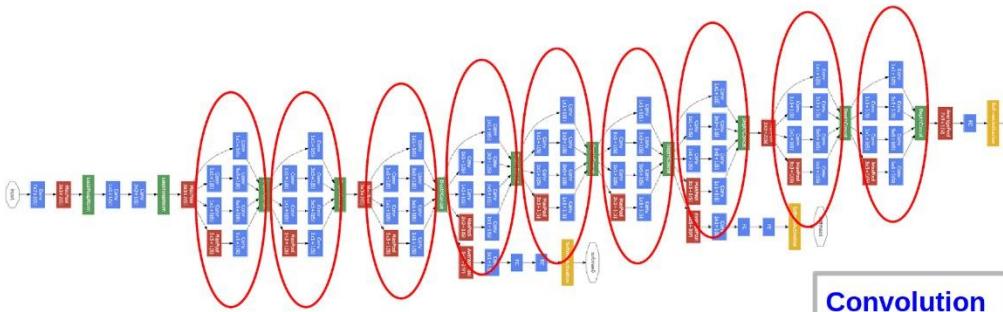
- **AlexNet (2012): First CNN (15.4%)**
 - 8 layers
 - 61 million parameters
- **ZFNet (2013): 15.4% to 11.2%**
 - 8 layers
 - More filters. Denser stride.
- **VGGNet (2014): 11.2% to 7.3%**
 - Beautifully uniform:
3x3 conv, stride 1, pad 1, 2x2 max pool
 - 16 layers
 - 138 million parameters
- **GoogLeNet (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
- **CUImage (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.

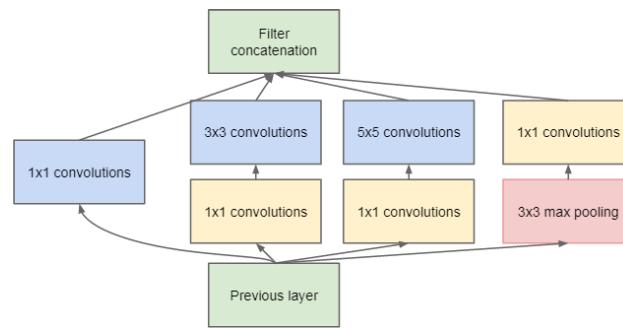
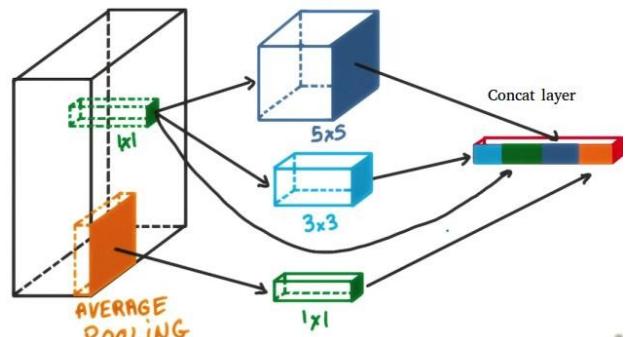


- **AlexNet (2012): First CNN (15.4%)**
 - 8 layers
 - 61 million parameters
- **ZFNet (2013): 15.4% to 11.2%**
 - 8 layers
 - More filters. Denser stride.
- **VGGNet (2014): 11.2% to 7.3%**
 - Beautifully uniform:
3x3 conv, stride 1, pad 1, 2x2 max pool
 - 16 layers
 - 138 million parameters
- **GoogLeNet (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.

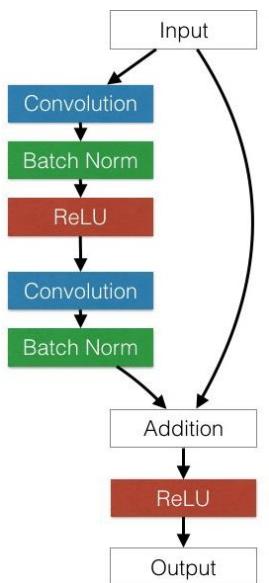
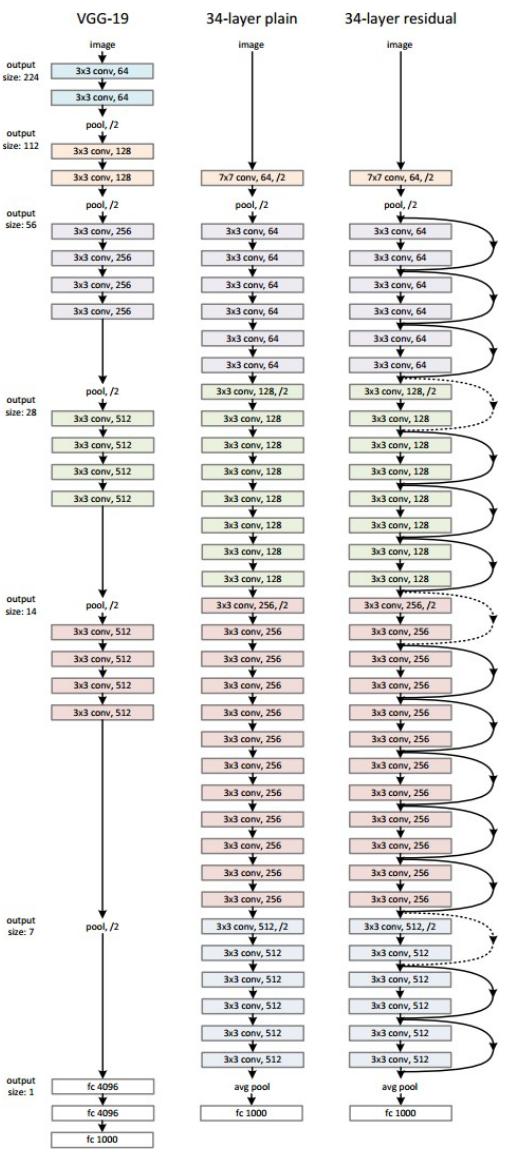


Convolution
Pooling
Softmax
Concat/Normalize



- **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

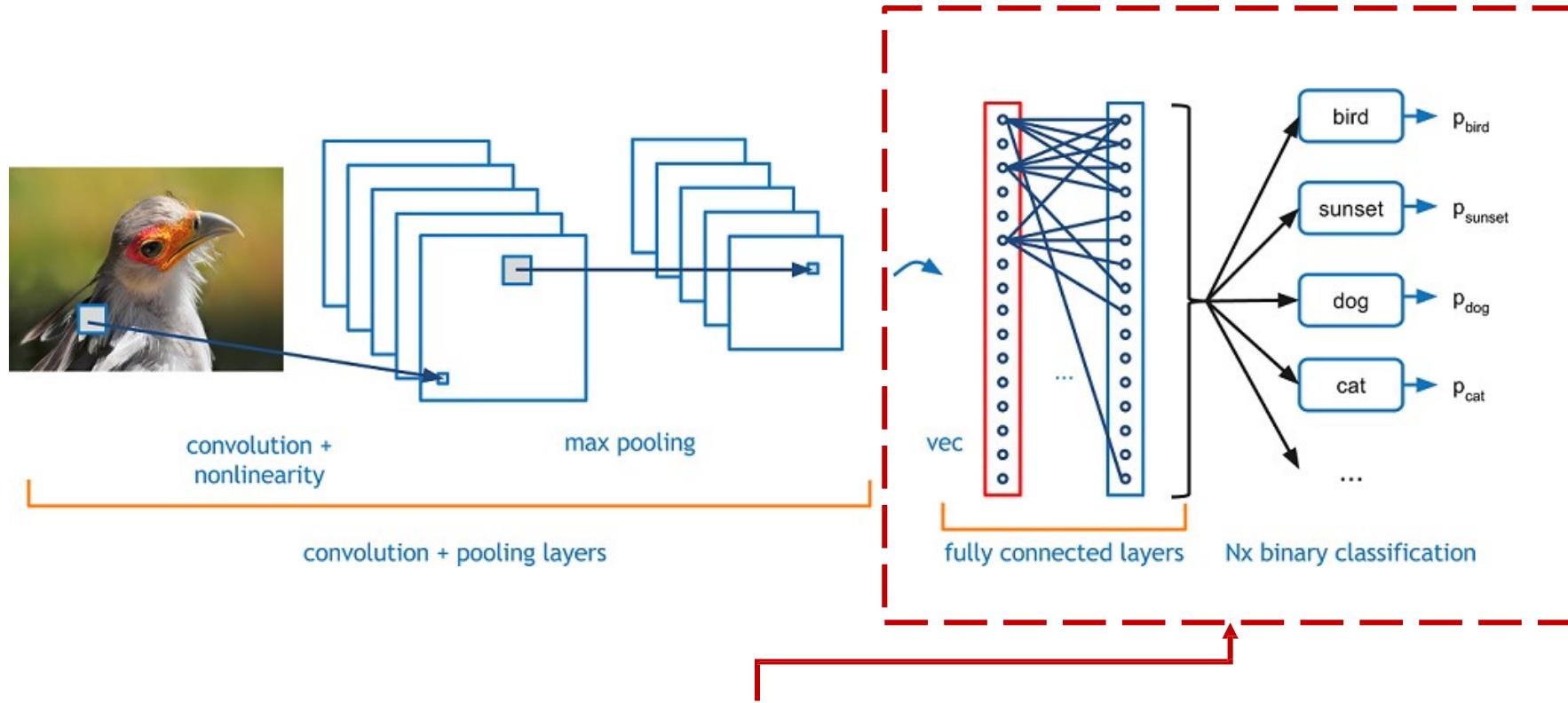
Szegedy et al. "Going deeper with convolutions." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2015.



- **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:
 3×3 conv, stride 1, pad 1, 2×2 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
- **CUImage (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.

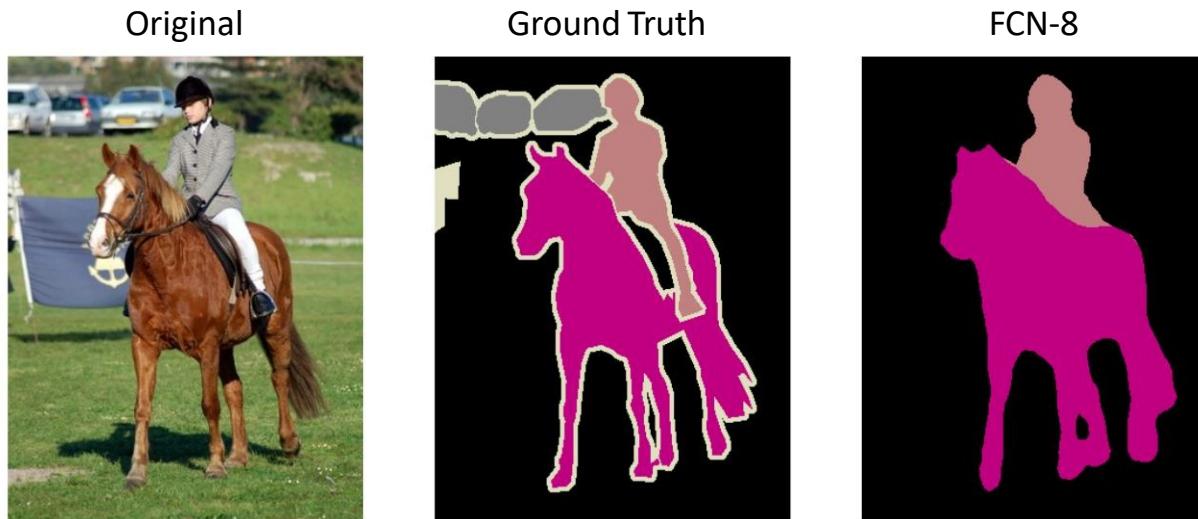
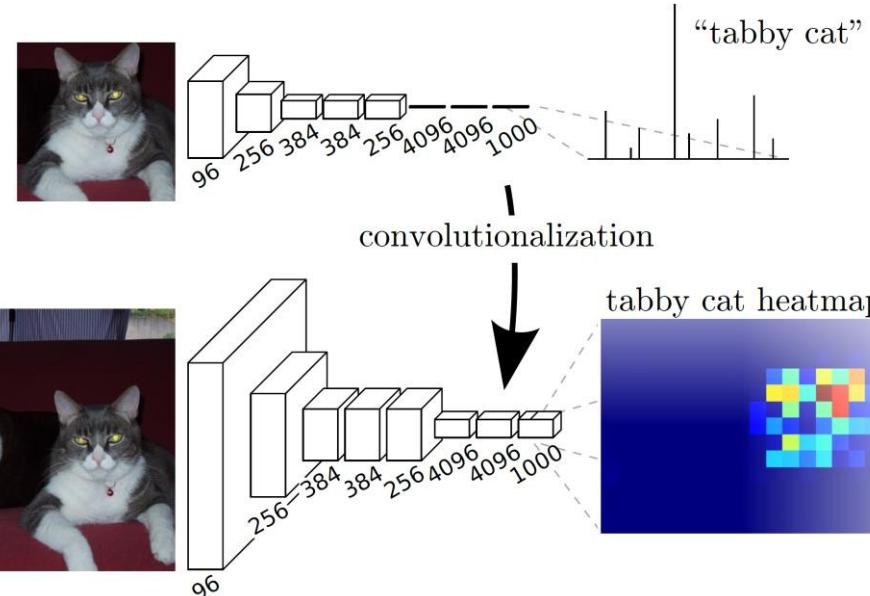
Same Architecture, Many Applications



This part might look different for:

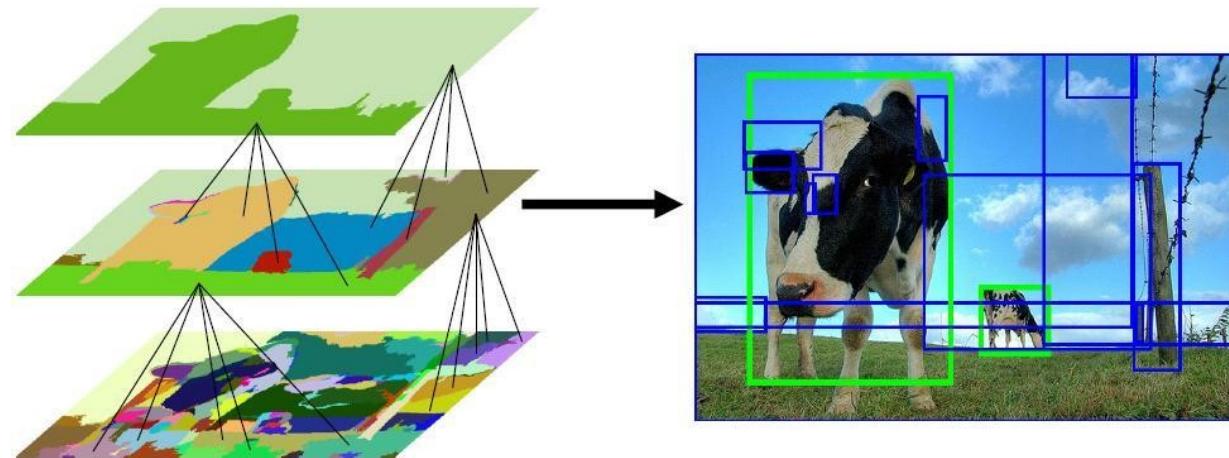
- Different image classification **domains**
- Image captioning with **recurrent neural networks**
- Image object localization with **bounding box**
- Image segmentation with **fully convolutional networks**
- Image segmentation with **deconvolution layers**

Segmentation

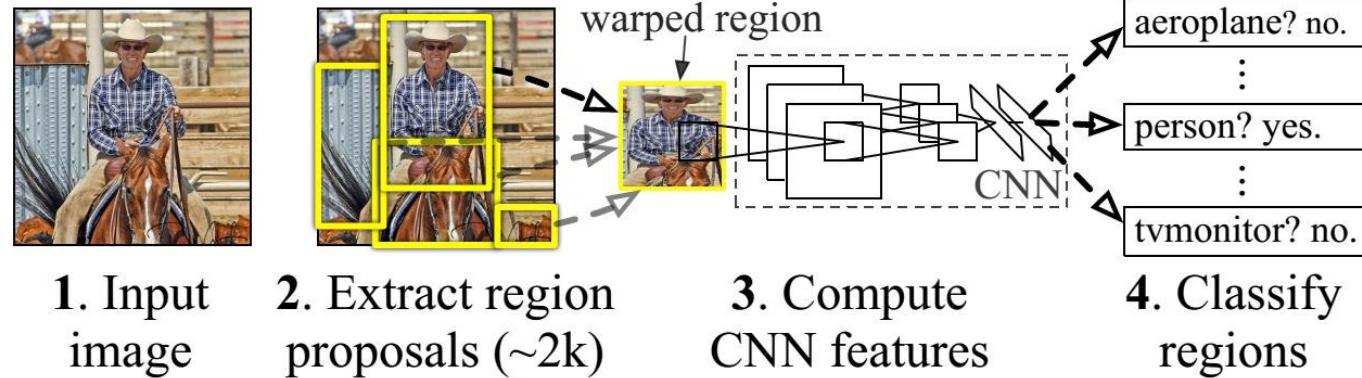




Object Detection



R-CNN: *Regions with CNN features*





Training Data



Testing Data



Minibatches Reduce Gradient Variance

- Initialize θ randomly
- For N Epochs
 - For each training **batch** $\{(x_0, y_0), \dots, (x_B, y_B)\}$:

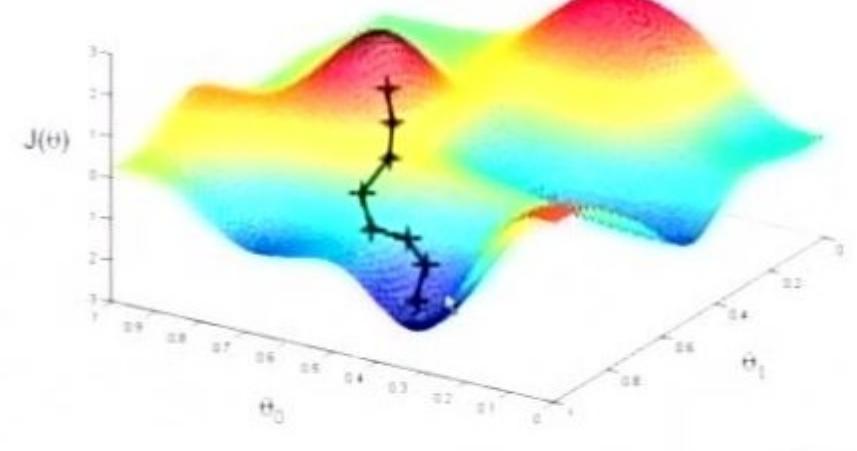
- Compute Loss Gradient:

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{1}{B} \sum_i^B \frac{\partial J_i(\theta)}{\partial \theta}$$

- Update θ with update rule:

$$\theta := \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$$

More accurate estimate!



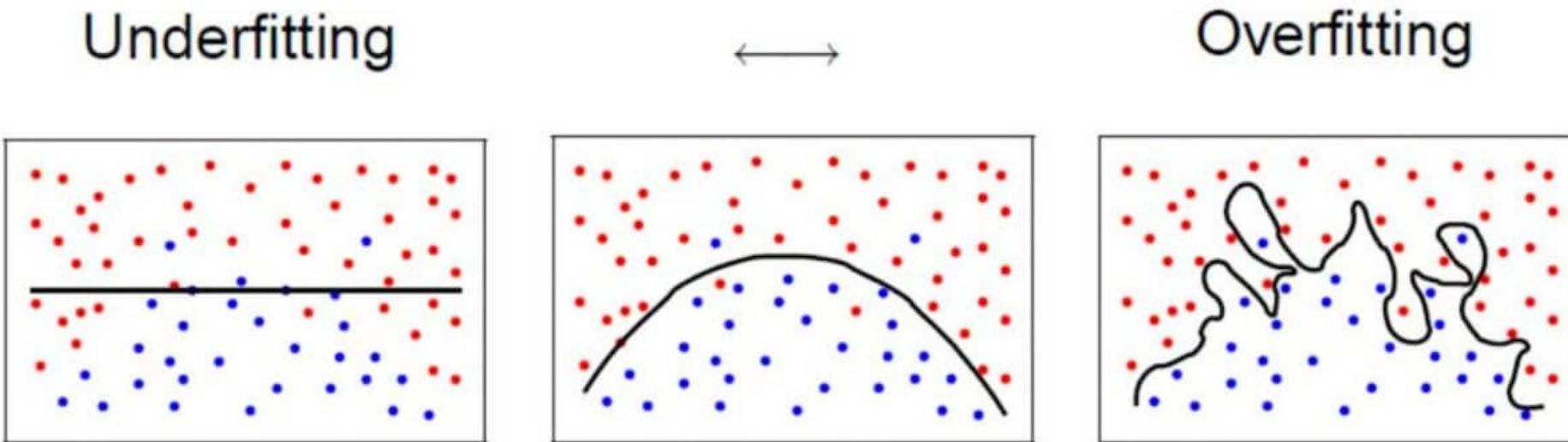
Advantages of Minibatches

- More accurate estimation of gradient
 - Smoother convergence
 - Allows for larger learning rates
- Minibatches lead to fast training!
 - Can parallelize computation + achieve significant speed increases on GPU's

Training Neural Networks In Practice: Fighting Overfitting

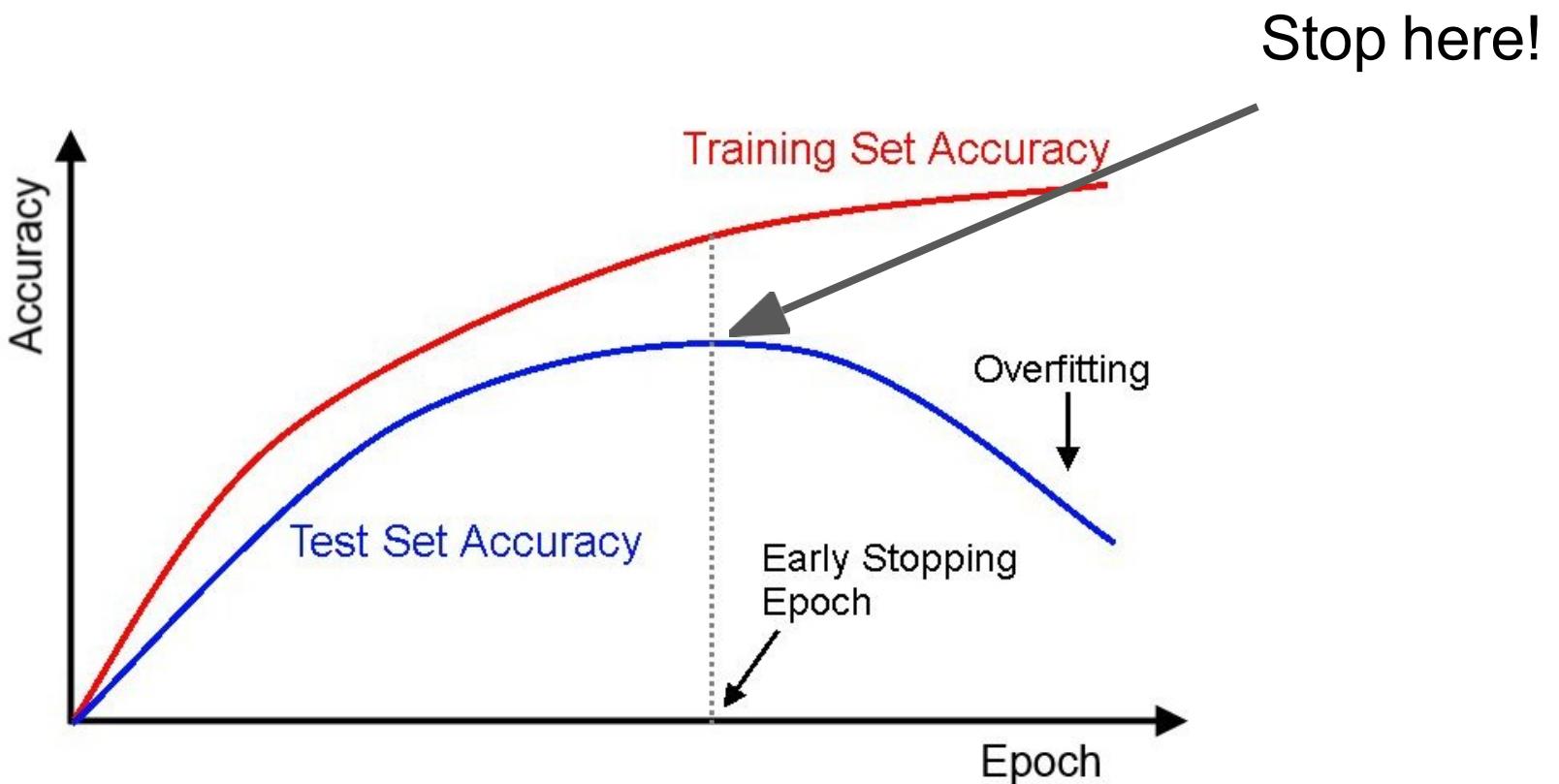
(Stolen from N. Locascio)

The Problem of Overfitting



(Stolen from N. Locascio)

Regularization II: Early Stopping



(Stolen from N. Locascio)

End-To-End

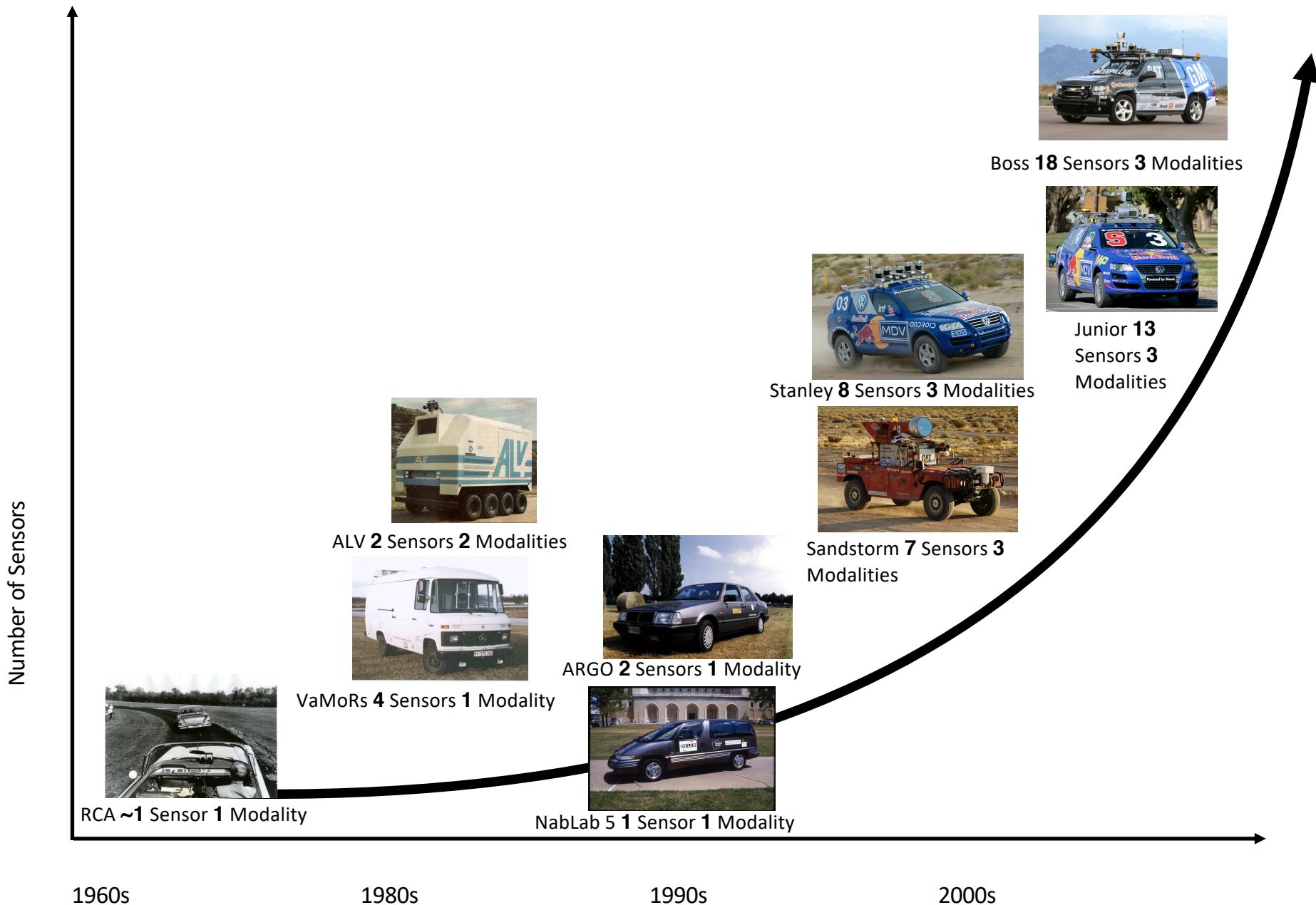


DAVE 2 Driving a Lincoln

- A convolutional neural network
- Trained by human drivers
- Learns perception, path planning, and control
"pixel in, action out"
- Front-facing camera is the only sensor

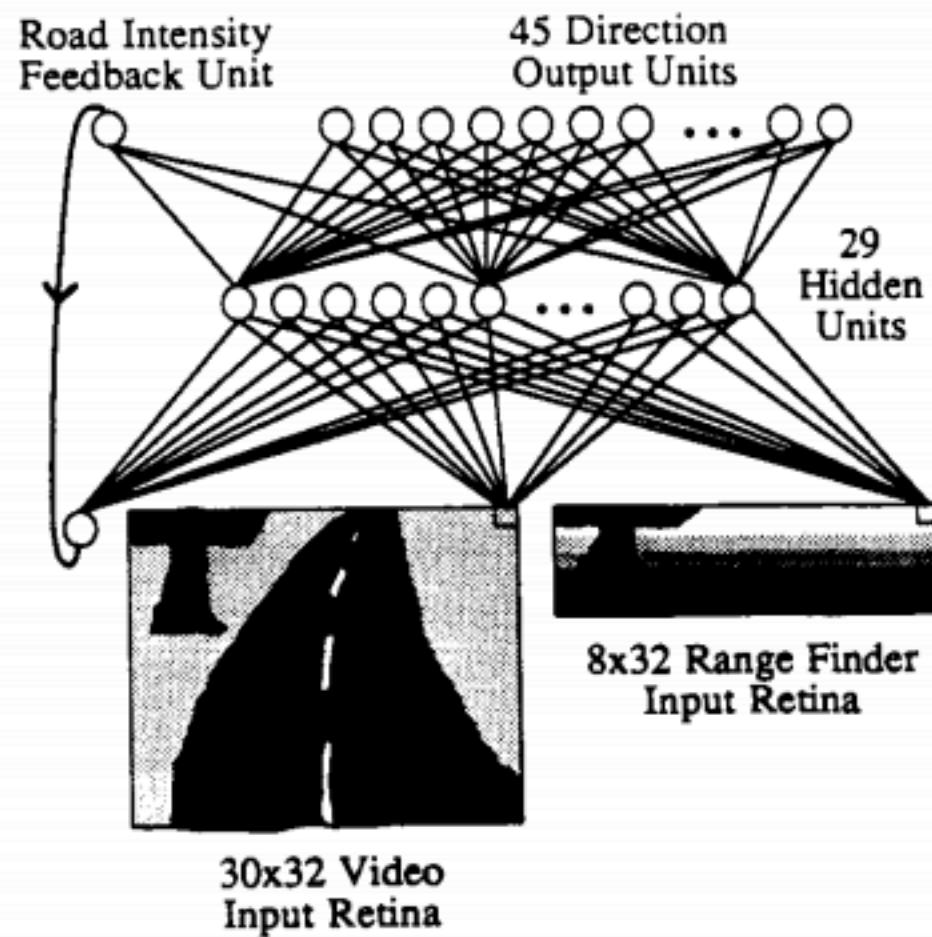
DAVE-1





1989

ALVINN Architecture



Neural Network-Based Autonomous Driving

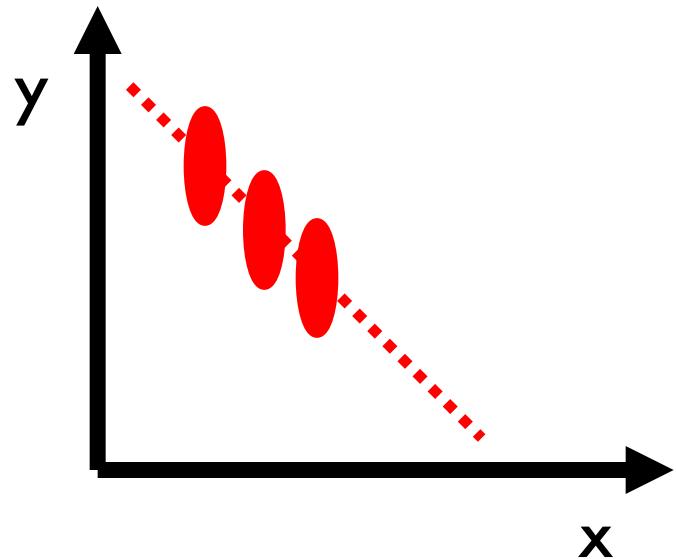
23 November 1992

[Courtesy of Dean Pomerleau]

Hough transform

P.V.C. Hough, *Machine Analysis of Bubble Chamber Pictures*, Proc. Int. Conf. High Energy Accelerators and Instrumentation, 1959

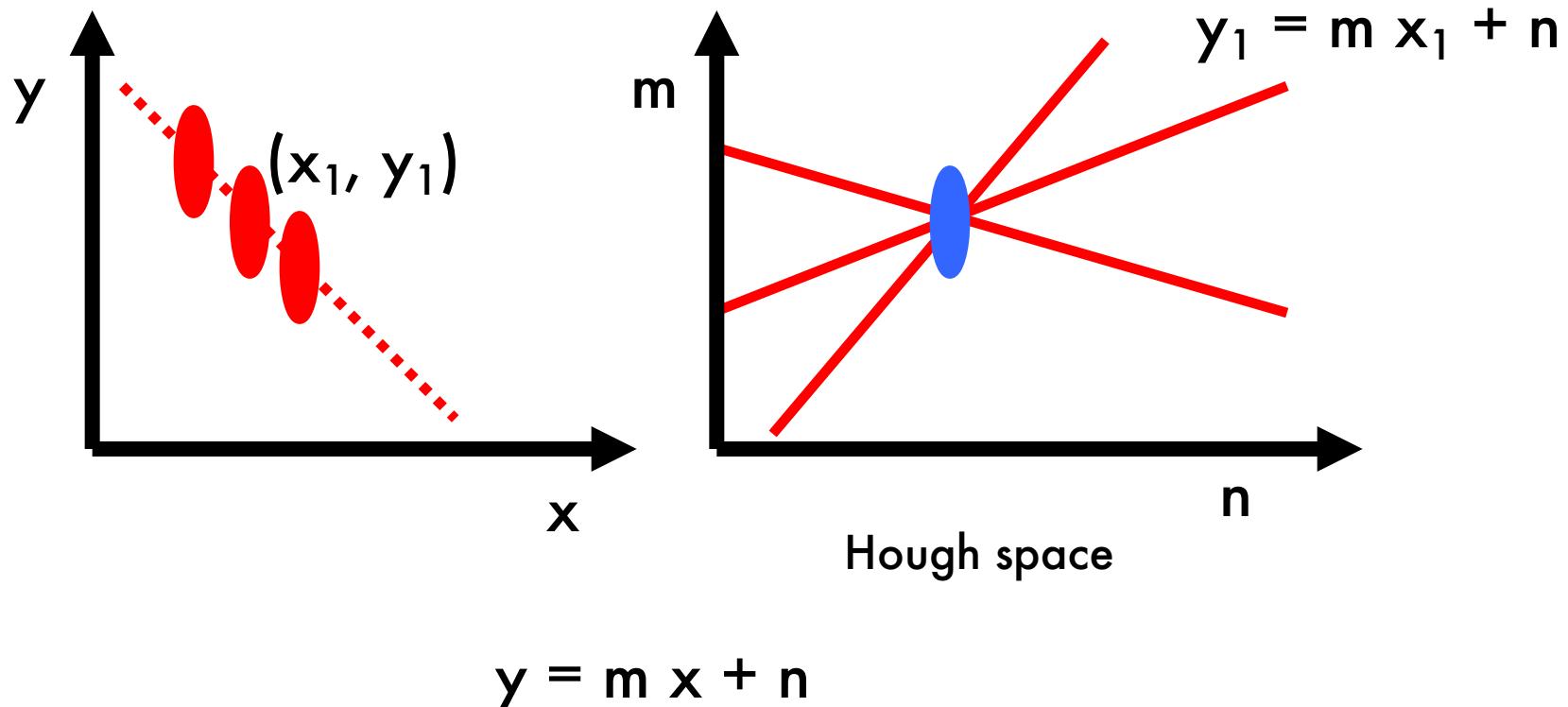
Given a set of points, find the curve or line that explains the data points best



Hough transform

P.V.C. Hough, *Machine Analysis of Bubble Chamber Pictures*, Proc. Int. Conf. High Energy Accelerators and Instrumentation, 1959

Given a set of points, find the curve or line that explains the data points best

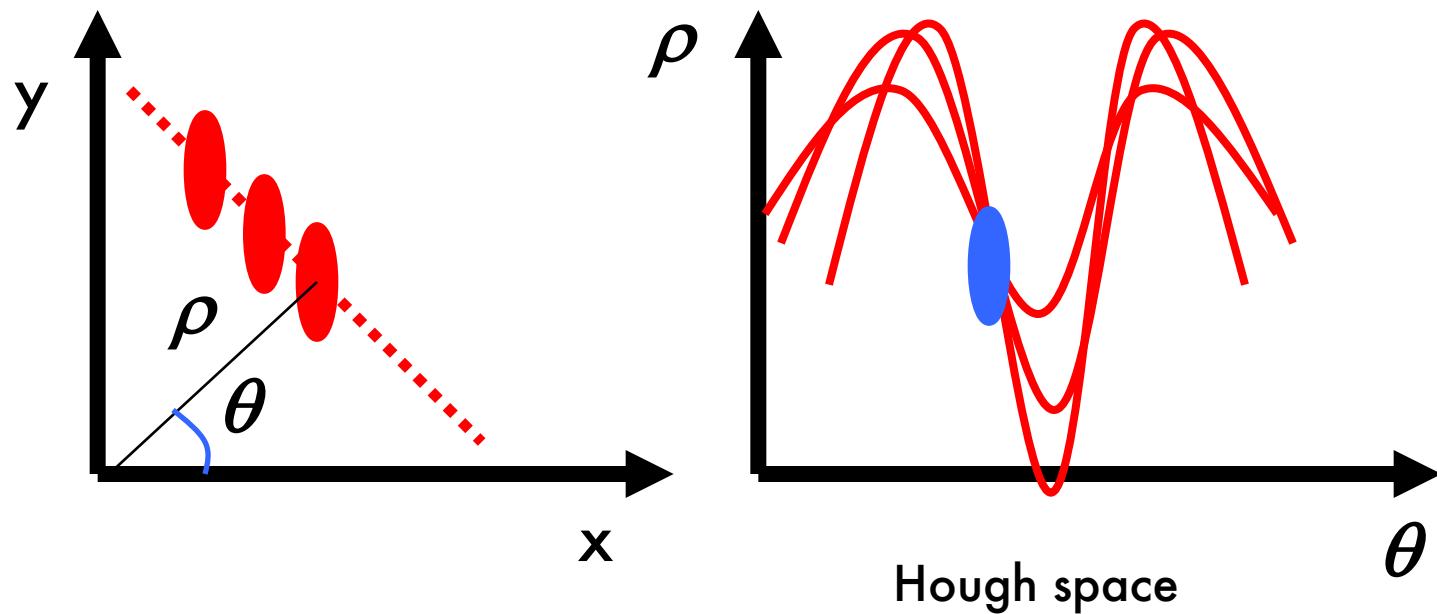


Hough transform

P.V.C. Hough, *Machine Analysis of Bubble Chamber Pictures*, Proc. Int. Conf. High Energy Accelerators and Instrumentation, 1959

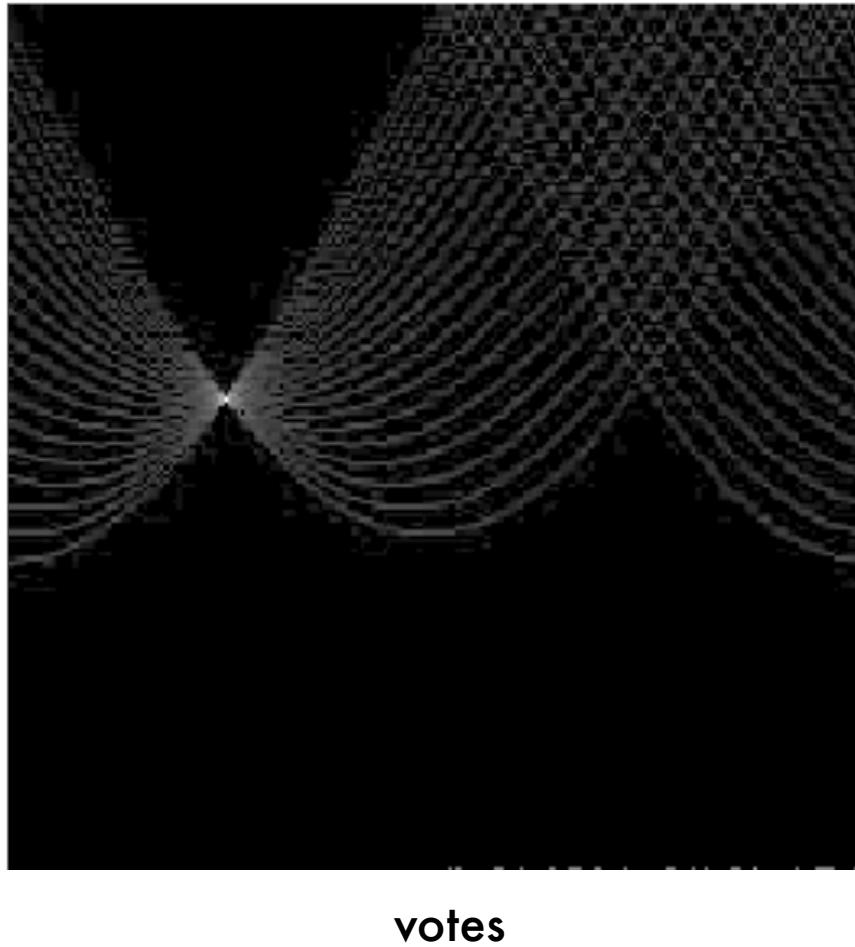
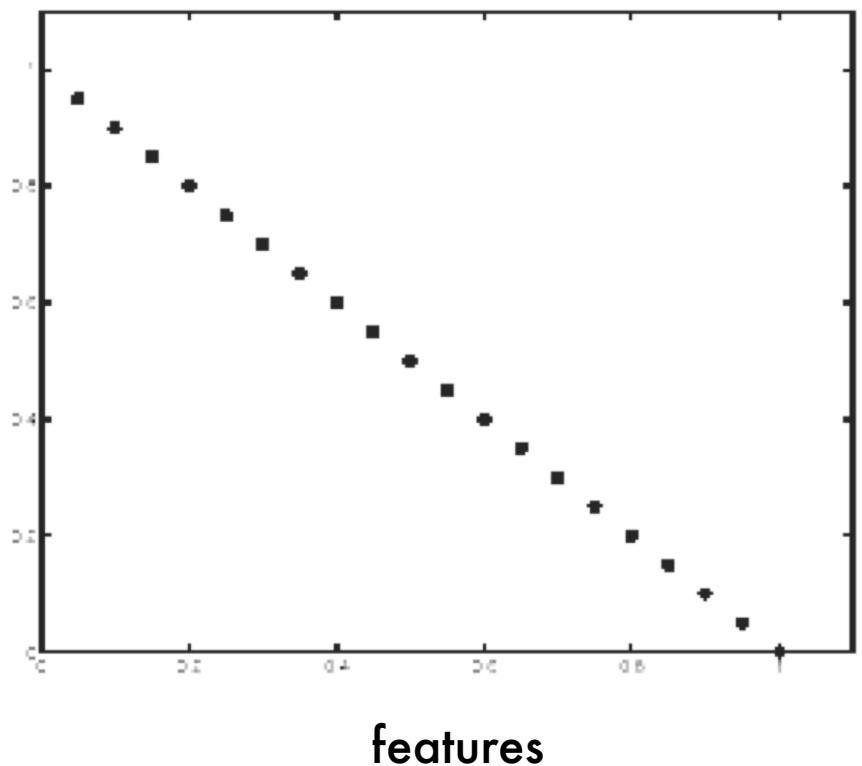
Issue : parameter space $[m,n]$ is unbounded...

- Use a polar representation for the parameter space



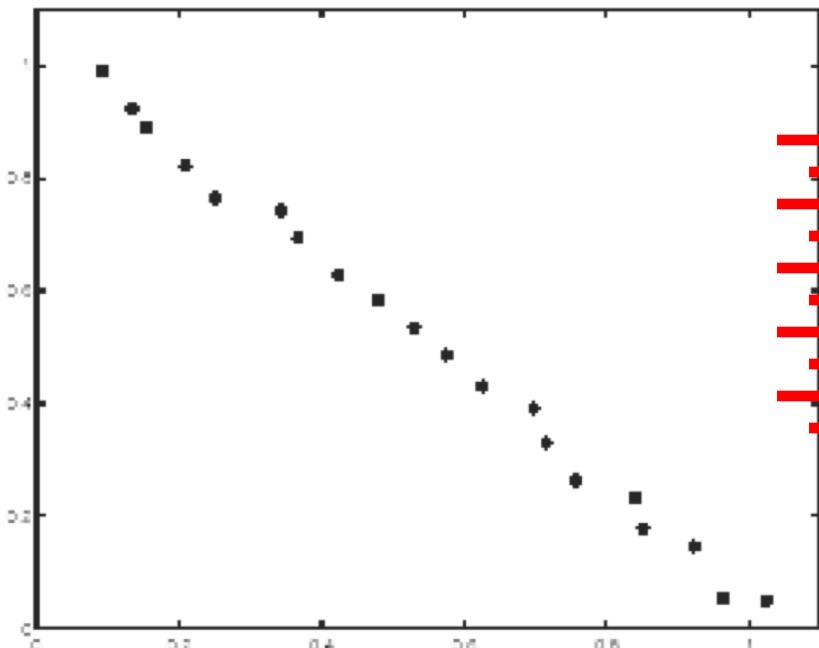
$$x \cos \theta + y \sin \theta = \rho$$

Hough transform - experiments

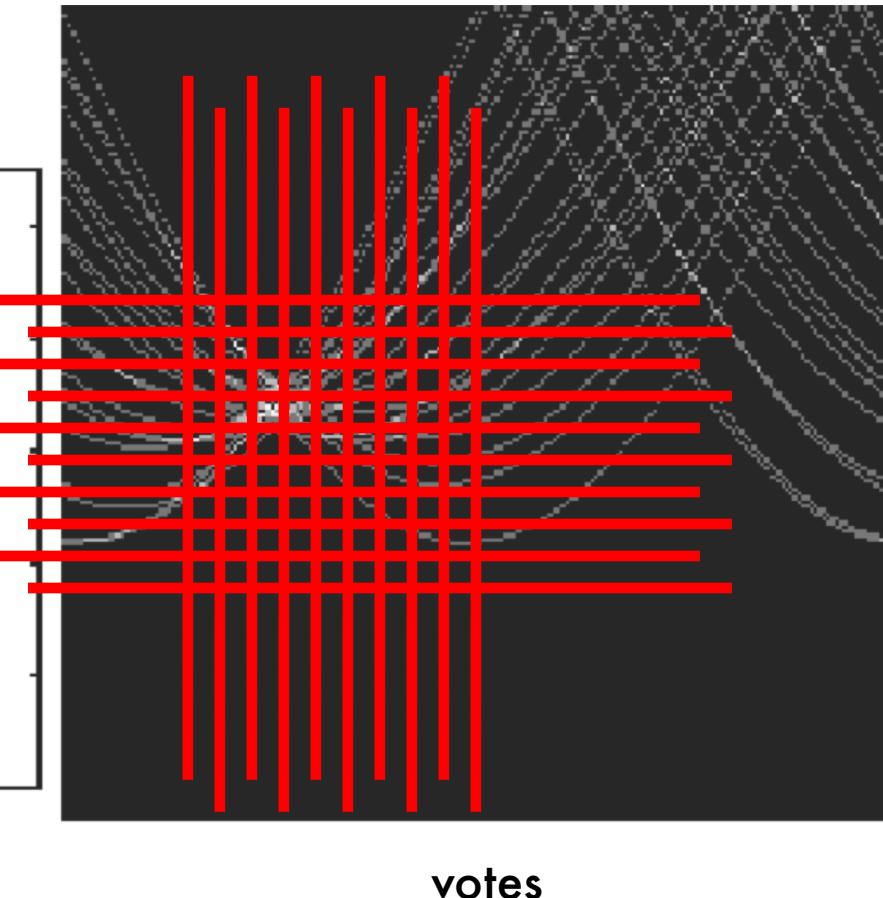


Hough transform - experiments

Noisy data



features



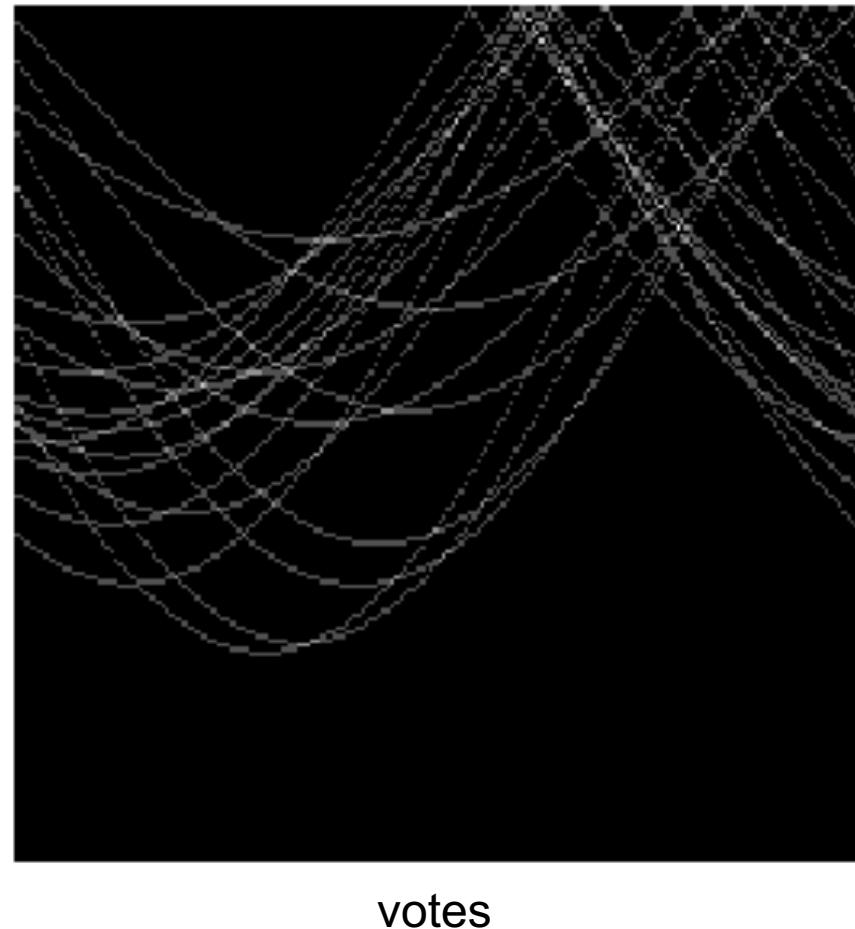
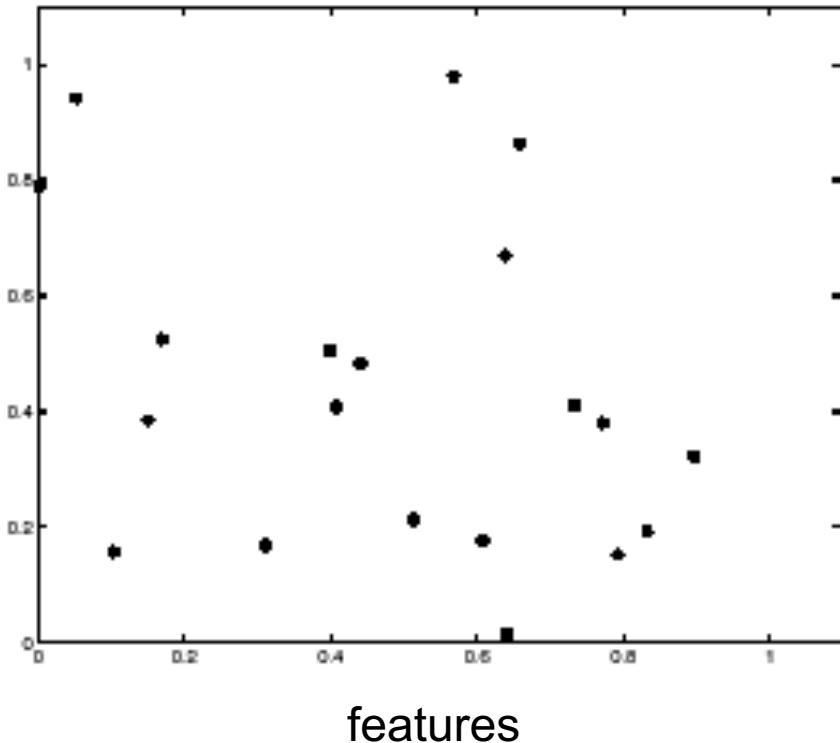
votes

How to compute the intersection point?

IDEA: introduce a grid a count intersection points in each cell

Issue: Grid size needs to be adjusted...

Hough transform - experiments



Issue: spurious peaks due to uniform noise

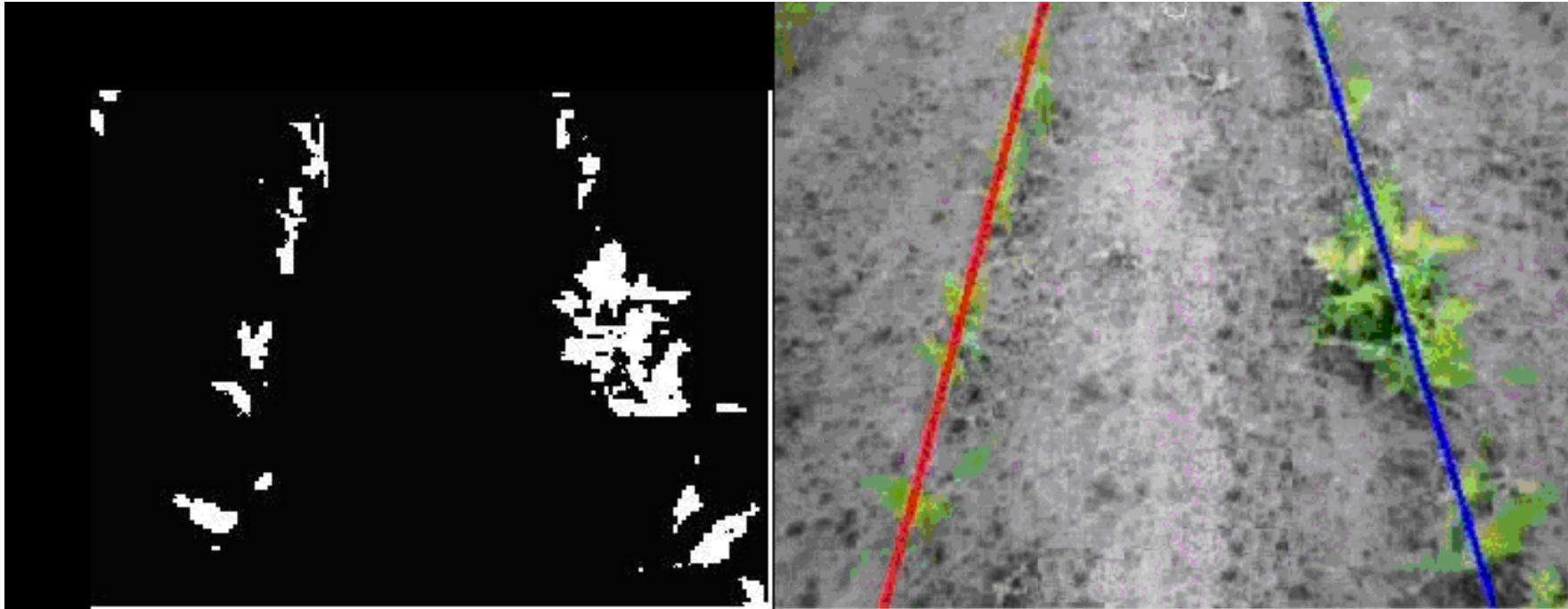
Hough transform - conclusions

- All points are processed independently, so can cope with occlusion/outliers
- Some robustness to noise: noise points unlikely to contribute consistently to any single bin

Bad:

- Spurious peaks due to uniform noise
- Trade-off noise-grid size (hard to find sweet point)

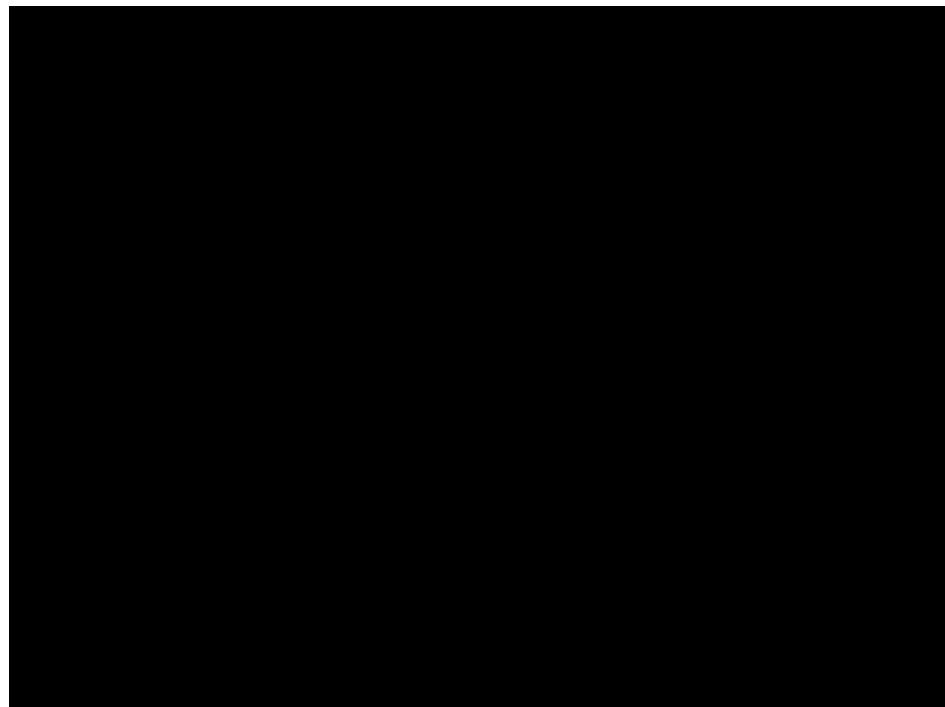
Hough transform - experiments



Courtesy of TKK Automation Technology Laboratory



NavLab 97



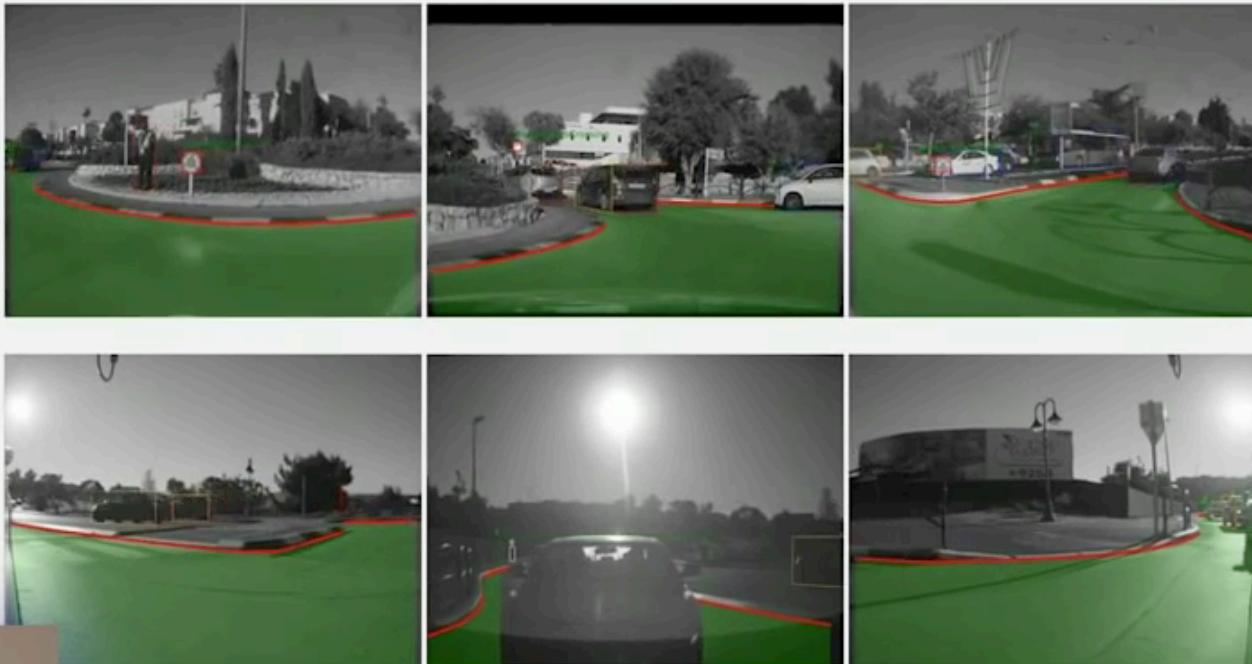


Componet





Free Space through Pixel Labeling



Our Vision. Your Safety.™



MOBILEYE®

Our Vision. Your Safety.™



Hybrid?



Learning Affordance for Direct Perception in Autonomous Driving

Chenyi Chen Ari Seff Alain Kornhauser Jianxiong Xiao

Princeton University



Controversy

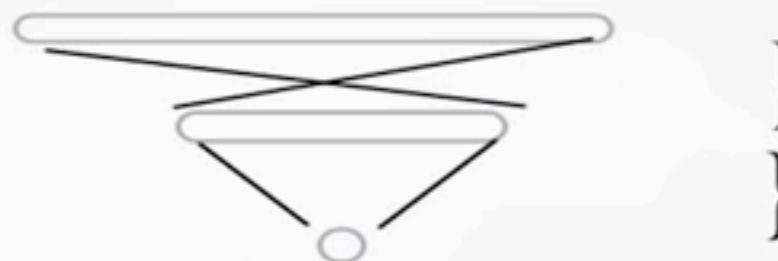
DNN Failure

Example: the n-bit parity problem

$x \in \{0, 1\}^n$ input nodes

$\frac{n}{2} + 1$ hidden nodes

$y \in \{0, 1\}$ output node



} weights are hard-wired

} weights found by system
of linear equations

Setiono, 1997

- Even though there exists weights that solve the n-bit parity problems, “learning” them using the available training techniques does not work for $n > 30$.
- This failure to train a DNN holds true also for over-subscribed architectures.

DNN Failure

Example: Learning Arithmetic Operations

Hoshen & Peleg, 2015

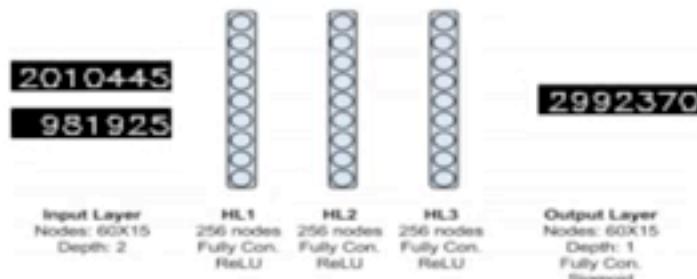


Figure 2. A diagram showing the construction of a neural network with 3 hidden layers able to perform addition using visual data. Two pictures are used as input and one picture as output. The network is fully connected and uses ReLU units in the hidden layers and sigmoid in the output layer. The hidden layers have 256 units each.

	Example A	Example B	Failure Example
Input Picture 1	981925	3570002	3750668
Input Picture 2	2010445	1216536	3643197
Network Output Picture	2992370	4786538	7393855
Ground Truth Picture	2992370	4786538	7393865

Figure 1. Input and output examples from our neural network trained for addition. The first two examples show a typical correct. The last example shows a rare failure case.

- DNN failed on the task of multiplication - whatever architecture they used they were unsuccessful in training the DNN.

DNN Failure

Example: Pentomino Dataset



Figure 1: Different classes of Pentomino shapes used in our dataset.



(a) sprites, not all same type



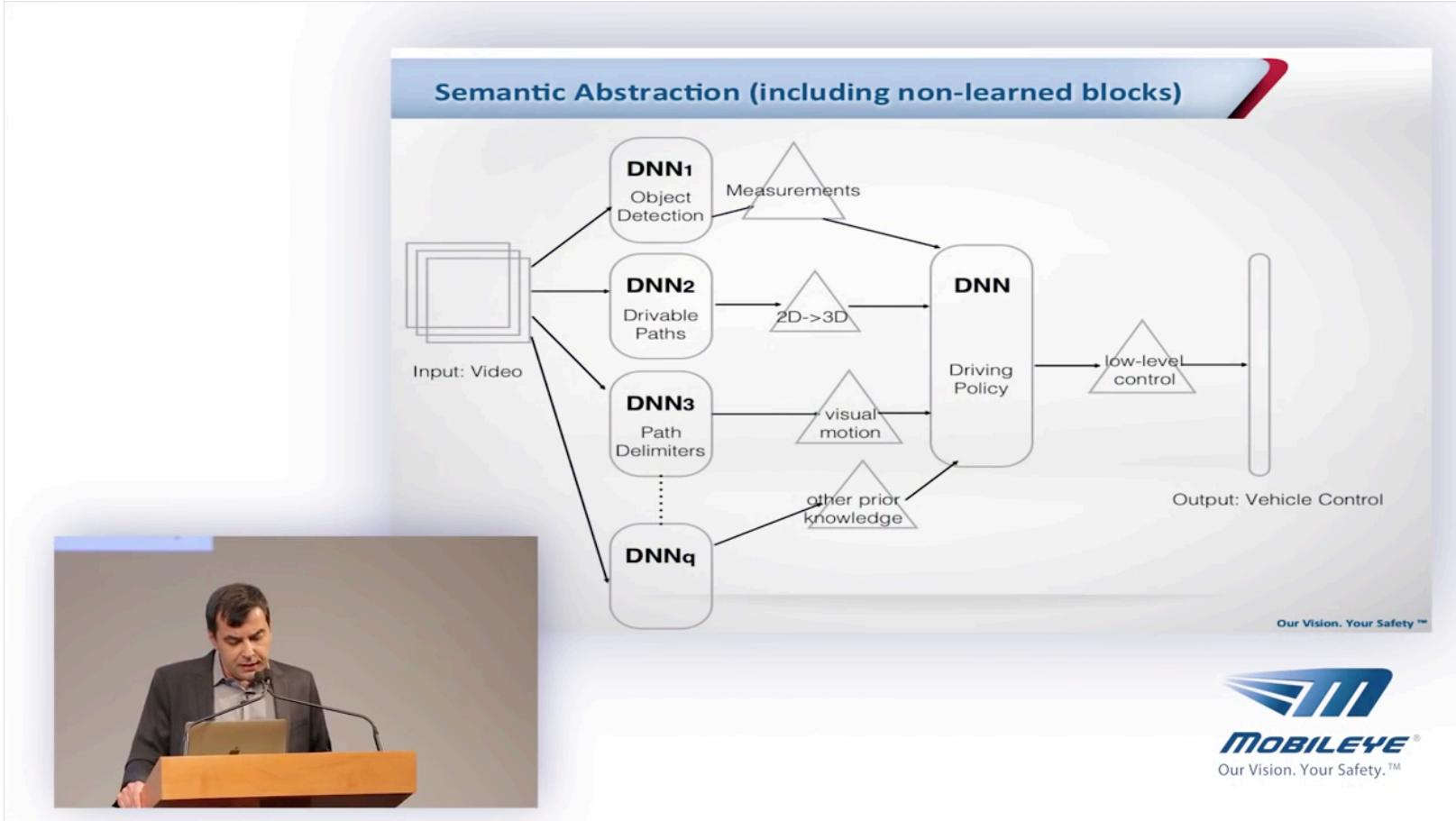
(b) sprites, all same type

Figure 2: Left: (a): An example image from the dataset which has *a different sprite type* in it. Right (b): An example image from the dataset that has only one type of Pentomino object in it, but with different orientations and scales.

Gulcehre & Bengio, 2016

- Different part types and can appear following some 2D geometric transformations
- Task: find out whether all parts are of the same class or not.

- DNN failed on the task when end-to-end was concerned.
- DNN succeeded when the task was broken down into first finding the category of each part and then making a decision whether all categories are the same.



What is Wrong with “End-to-End”?

- It is all about “corner cases” - rare events.
 - The probability of rare events decreases exponentially with end-to-end architecture.
- The importance of “Prior Knowledge”
 - DNNs could fail miserably on seemingly simple problems.



Our Vision. Your Safety.™



DNNs for Control?

SE3-Nets: Learning Rigid Body Motion using Deep Neural Networks

Arunkumar Byravan and Dieter Fox
University of Washington, Seattle

Supplementary video

This video has five sections: “Approach” & four sections of “Results”

Deep
Mind
Human
Models

