

CSE 152: Computer Vision

Hao Su

Lecture 7: Neural Networks



Review of Filters: **From Linear to Non-linear**

Image filtering (Linear case)

$$f[\cdot, \cdot] \frac{1}{9}$$

1	1	1
1	1	1
1	1	1

$I[\cdot, \cdot]$

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

$h[\cdot, \cdot]$

0										

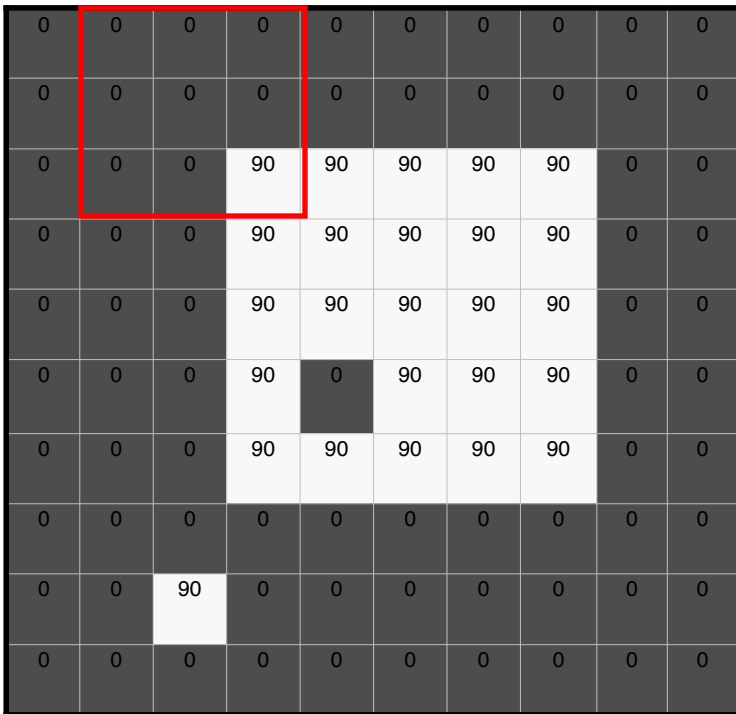
$$h[m, n] = \sum_{k, l} f[k, l] I[m + k, n + l]$$

Image filtering (Linear case)

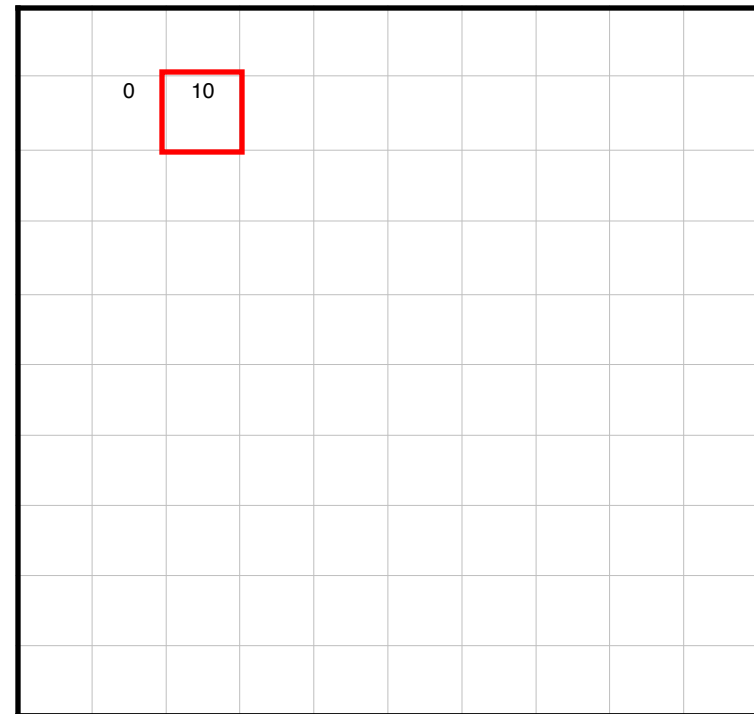
$$f[\cdot, \cdot] \frac{1}{9}$$

1	1	1
1	1	1
1	1	1

$I[\cdot, \cdot]$



$h[\cdot, \cdot]$



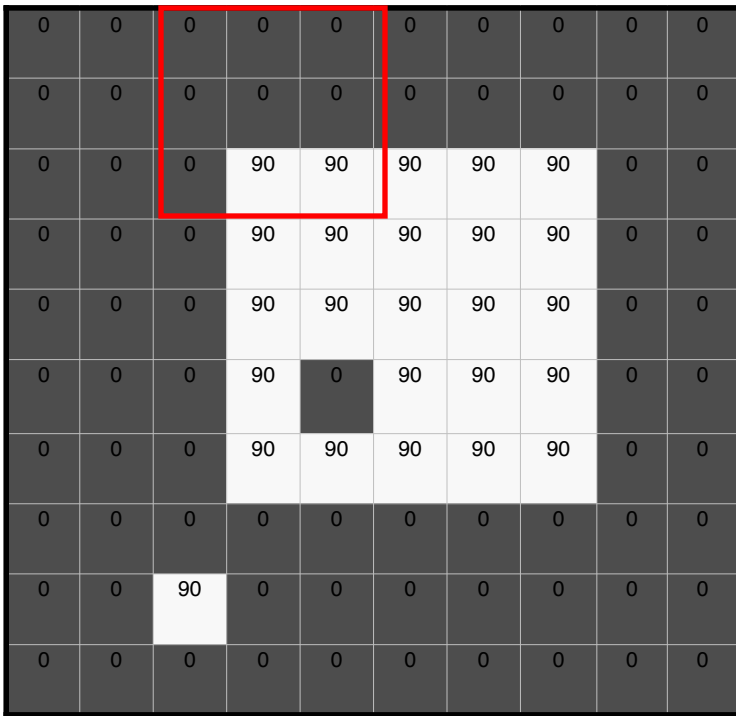
$$h[m, n] = \sum_{k, l} f[k, l] I[m + k, n + l]$$

Image filtering (Linear case)

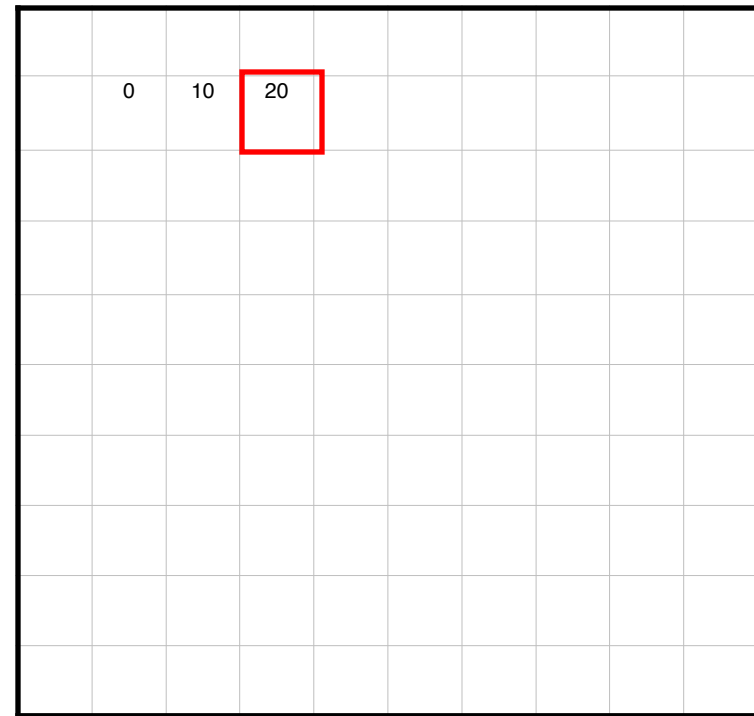
$$f[\cdot, \cdot] \frac{1}{9}$$

1	1	1
1	1	1
1	1	1

$I[\cdot, \cdot]$



$h[\cdot, \cdot]$



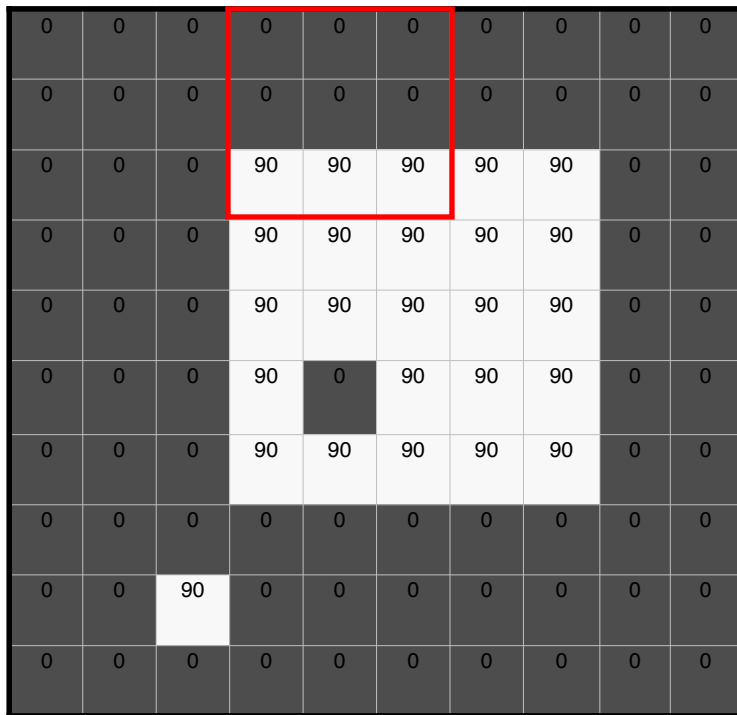
$$h[m, n] = \sum_{k, l} f[k, l] I[m + k, n + l]$$

Image filtering (Linear case)

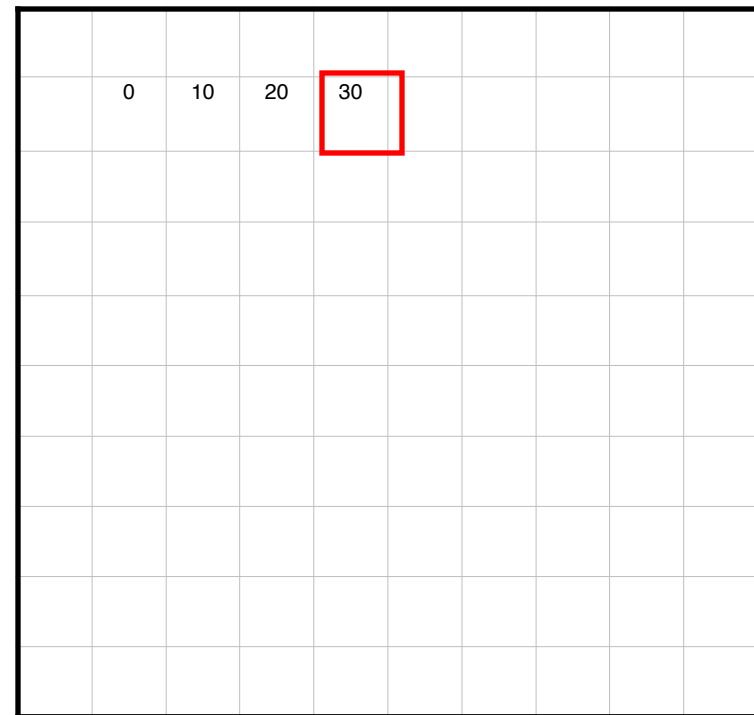
$$f[\cdot, \cdot] \frac{1}{9}$$

1	1	1
1	1	1
1	1	1

$I[\cdot, \cdot]$



$h[\cdot, \cdot]$



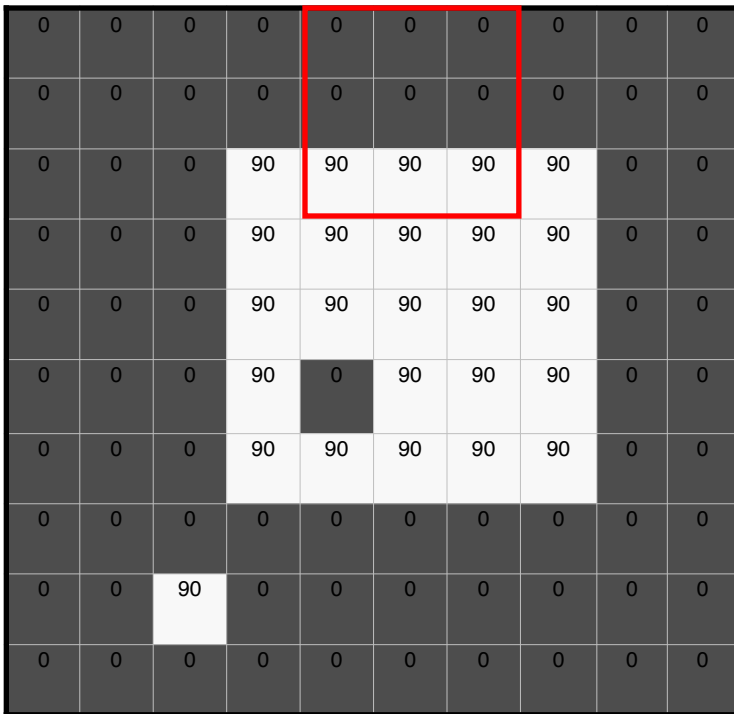
$$h[m, n] = \sum_{k, l} f[k, l] I[m + k, n + l]$$

Image filtering (Linear case)

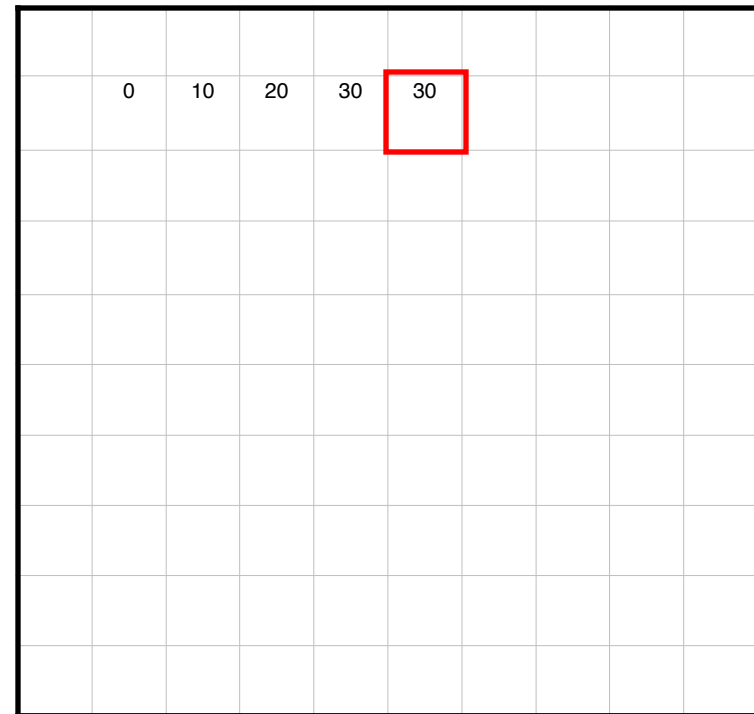
$$f[\cdot, \cdot] \frac{1}{9}$$

1	1	1
1	1	1
1	1	1

$I[\cdot, \cdot]$



$h[\cdot, \cdot]$



$$h[m, n] = \sum_{k, l} f[k, l] I[m + k, n + l]$$

Image filtering (Linear case) $f[\cdot, \cdot]$

1	1	1
1	1	1
1	1	1

$I[\cdot, \cdot]$

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

$h[\cdot, \cdot]$

	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	

$$h[m, n] = \sum_{k, l} f[k, l] I[m + k, n + l]$$

Reducing salt-and-pepper noise

3x3



5x5



7x7

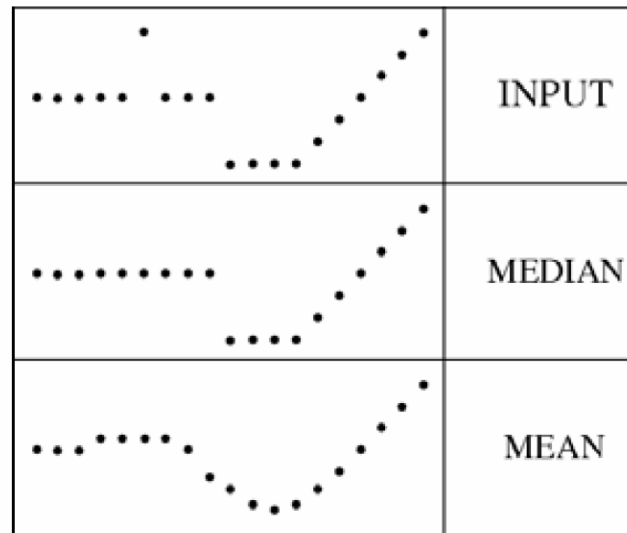


- What's wrong with the results?

Median filter (Non-linear)

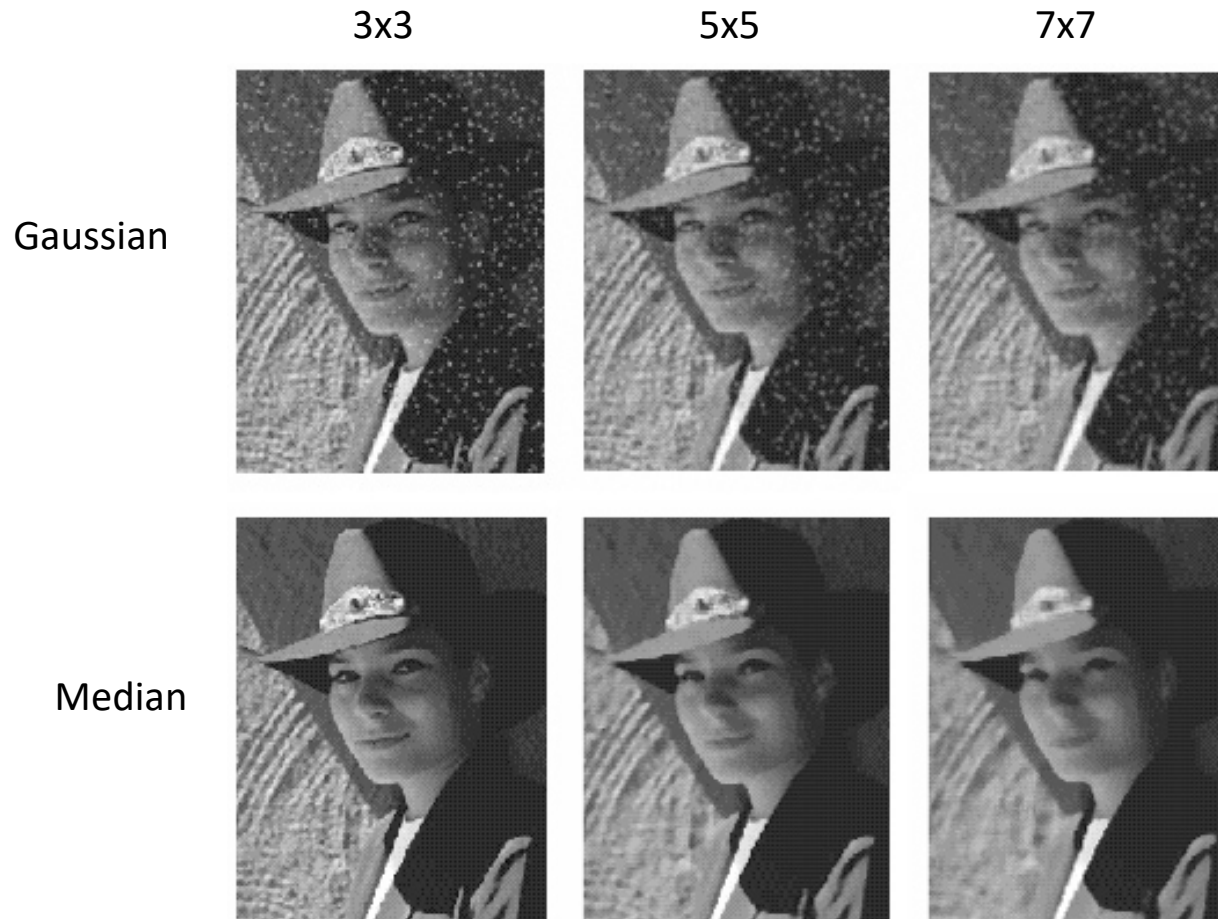
- What advantage does median filtering have over box filtering?
 - Robustness to outliers

filters have width 5 :



Source: K. Grauman

Median filter (Non-linear)



Gaussian vs. median filtering

3x3

5x5

7x7

Gaussian



Median



Neural Networks

A General Framework
from Linear to Non-linear

Image Classification: A core task in Computer Vision

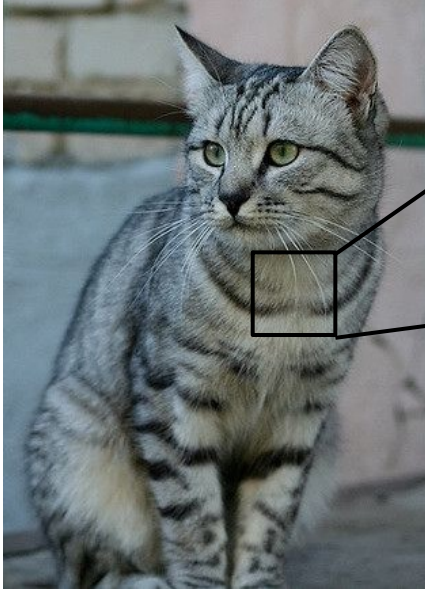


This image by Nikita is
licensed under [CC-BY 2.0](https://creativecommons.org/licenses/by/2.0/)

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

—————→ cat

The Problem: Semantic Gap



This image by Nikita is licensed under [CC-BY 2.0](https://creativecommons.org/licenses/by/2.0/)

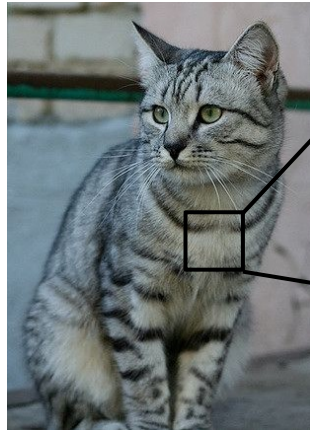
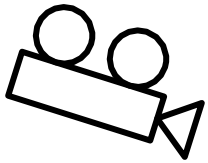
```
[[105 112 108 111 104 99 106 99 96 103 112 119 104 97 93 87]  
[ 91 98 102 106 104 79 98 103 99 105 123 136 110 105 94 85]  
[ 76 85 90 105 128 105 87 96 95 99 115 112 106 103 99 85]  
[ 99 81 81 93 120 131 127 100 95 98 102 99 96 93 101 94]  
[106 91 61 64 69 91 88 85 101 107 109 98 75 84 96 95]  
[114 108 85 55 55 69 64 54 64 87 112 129 98 74 84 91]  
[133 137 147 103 65 81 80 65 52 54 74 84 102 93 85 82]  
[128 137 144 140 109 95 86 70 62 65 63 63 60 73 86 101]  
[125 133 148 137 119 121 117 94 65 79 80 65 54 64 72 98]  
[127 125 131 147 133 127 126 131 111 96 89 75 61 64 72 84]  
[115 114 109 123 150 148 131 118 113 109 100 92 74 65 72 78]  
[ 89 93 90 97 106 147 131 118 113 114 113 109 106 95 77 60]  
[ 63 77 86 81 77 79 102 123 117 115 117 125 125 130 115 87]  
[ 62 65 82 89 78 71 80 101 124 126 119 101 107 114 131 119]  
[ 63 65 75 88 89 71 62 81 120 138 135 105 81 98 110 118]  
[ 87 65 71 87 106 95 69 45 76 130 126 107 92 94 105 112]  
[118 97 82 86 117 123 116 66 41 51 95 93 89 95 102 107]  
[164 146 112 80 82 120 124 104 76 48 45 66 88 101 102 109]  
[157 170 157 120 93 86 114 132 112 97 69 55 70 82 99 94]  
[130 128 134 161 139 100 109 118 121 134 114 87 65 53 69 86]  
[120 112 96 117 150 144 120 115 104 107 102 93 87 81 72 79]  
[123 107 96 86 83 112 153 149 122 109 104 75 80 107 112 99]  
[122 121 102 80 82 86 94 117 145 148 153 102 58 78 92 107]  
[122 164 148 103 71 56 78 83 93 103 119 139 102 61 69 84]]
```

What the computer sees

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

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

Challenges: Viewpoint variation



[105	112	108	111	104	99	106	99	96	103	112	119	104	97	93	87]
[91	98	102	106	104	79	98	103	99	105	123	136	110	105	94	85]
[76	85	90	105	128	105	87	96	95	99	115	112	106	103	99	85]
[99	81	81	93	120	131	127	100	95	98	102	99	96	93	101	94]
[106	91	61	64	69	91	88	85	101	107	109	98	75	84	96	95]
[114	100	05	55	55	69	64	54	64	87	112	125	98	74	84	91]
[133	137	147	103	65	81	80	65	52	54	74	84	102	93	85	82]
[128	137	144	140	109	95	86	70	62	65	63	63	60	73	86	101]
[125	133	148	137	119	121	117	94	65	79	80	65	54	64	72	90]
[127	125	131	147	133	127	126	131	111	96	89	75	61	64	72	84]
[115	114	109	123	150	148	131	118	113	109	100	92	74	65	72	78]
[89	93	98	97	100	147	131	118	113	114	113	100	106	95	77	80]
[63	77	86	81	77	79	102	123	117	115	117	125	125	130	115	87]
[62	65	82	89	78	71	80	101	124	126	119	101	107	114	131	119]
[63	65	75	88	69	71	62	81	120	130	135	105	81	90	110	118]
[87	65	71	97	106	95	69	65	76	130	126	107	92	94	105	112]
[118	97	82	86	117	123	116	66	41	51	95	93	89	95	102	107]
[164	146	112	80	82	120	124	104	76	48	45	66	88	101	102	109]
[157	170	157	120	63	86	114	132	112	97	69	55	70	82	90	94]
[130	128	134	161	139	100	109	118	121	134	114	87	65	53	69	86]
[120	112	96	117	150	144	120	115	104	107	102	93	87	81	72	79]
[123	107	96	86	83	112	153	149	122	109	104	75	80	107	112	99]
[122	121	102	80	82	86	94	117	145	148	153	102	58	78	92	107]
[122	164	148	103	71	56	78	83	93	103	119	139	102	61	69	84]

All pixels change when the camera moves!

Challenges: Illumination



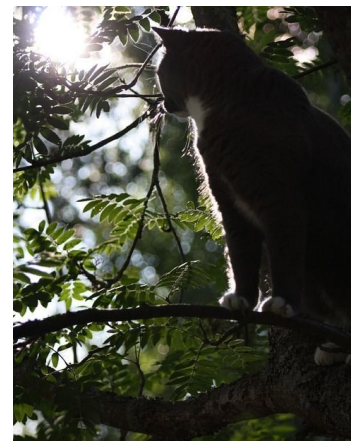
[This image](#) is [CC0 1.0](#) public domain



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

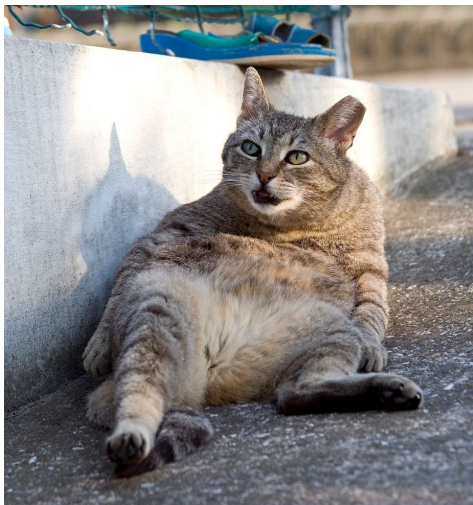


[This image](#) is [CC0 1.0](#) public domain



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

Challenges: Deformation



[This image](#) by [Umberto Salvagnin](#) is licensed under [CC-BY 2.0](#)



[This image](#) by [Umberto Salvagnin](#) is licensed under [CC-BY 2.0](#)



[This image](#) by [sare bear](#) is licensed under [CC-BY 2.0](#)



[This image](#) by [Tom Thai](#) is licensed under [CC-BY 2.0](#)

Challenges: Occlusion



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



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



This image by [jonsson](#) is licensed under [CC-BY 2.0](#)

Challenges: Background Clutter



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



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

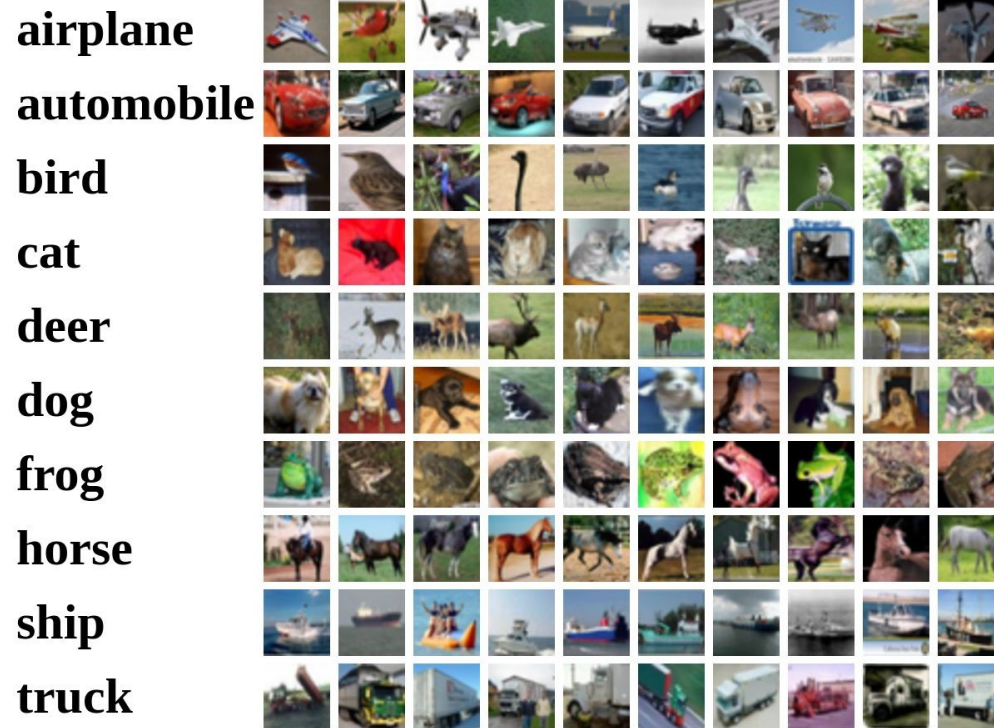
Challenges: Intraclass variation



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

Linear Classification

Recall CIFAR10



50,000 training images
each image is **32x32x3**

10,000 test images.

Parametric Approach

Image



Array of **32x32x3** numbers
(3072 numbers total)



10 numbers giving
class scores

↑
W

parameters
or weights

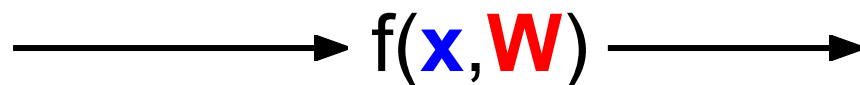
Parametric Approach: Linear Classifier

Image



Array of **32x32x3** numbers
(3072 numbers total)

$$f(x, W) = Wx$$



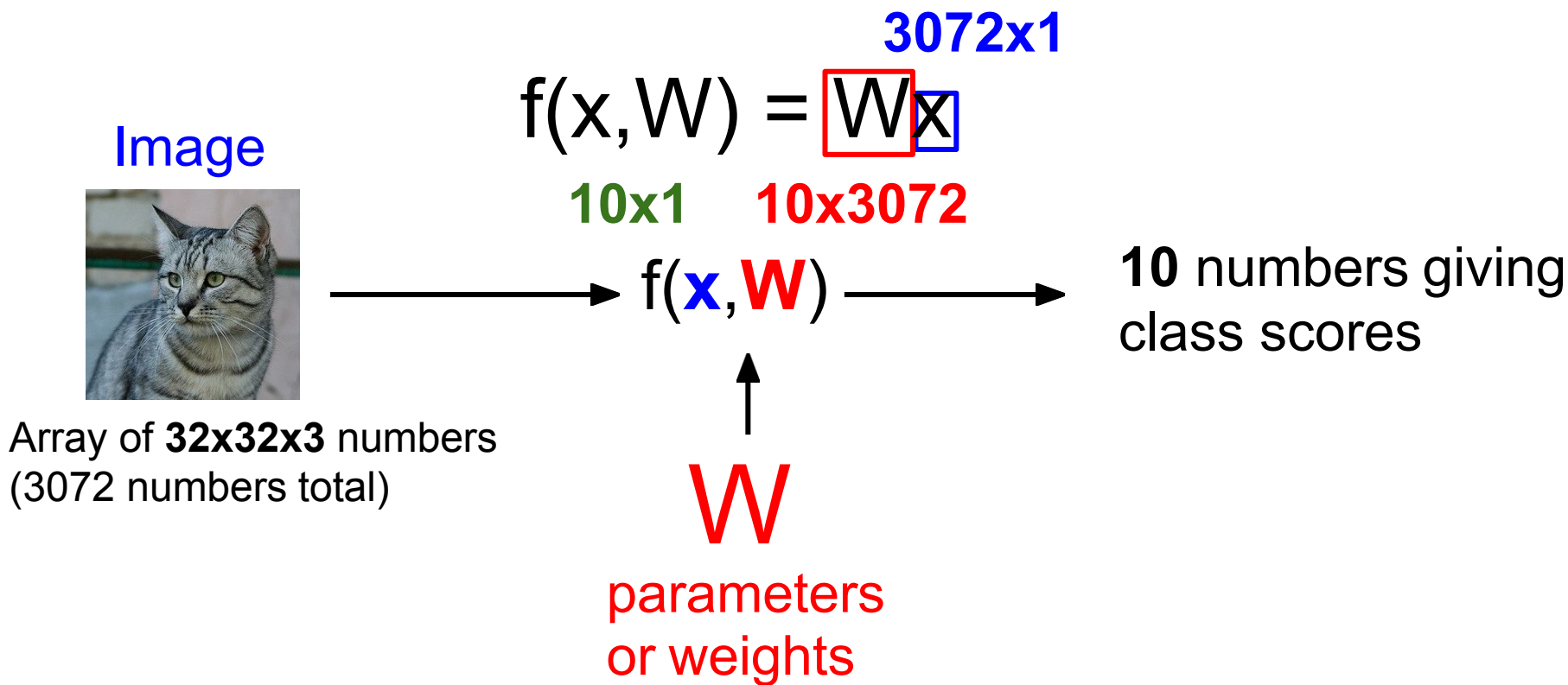
10 numbers giving
class scores



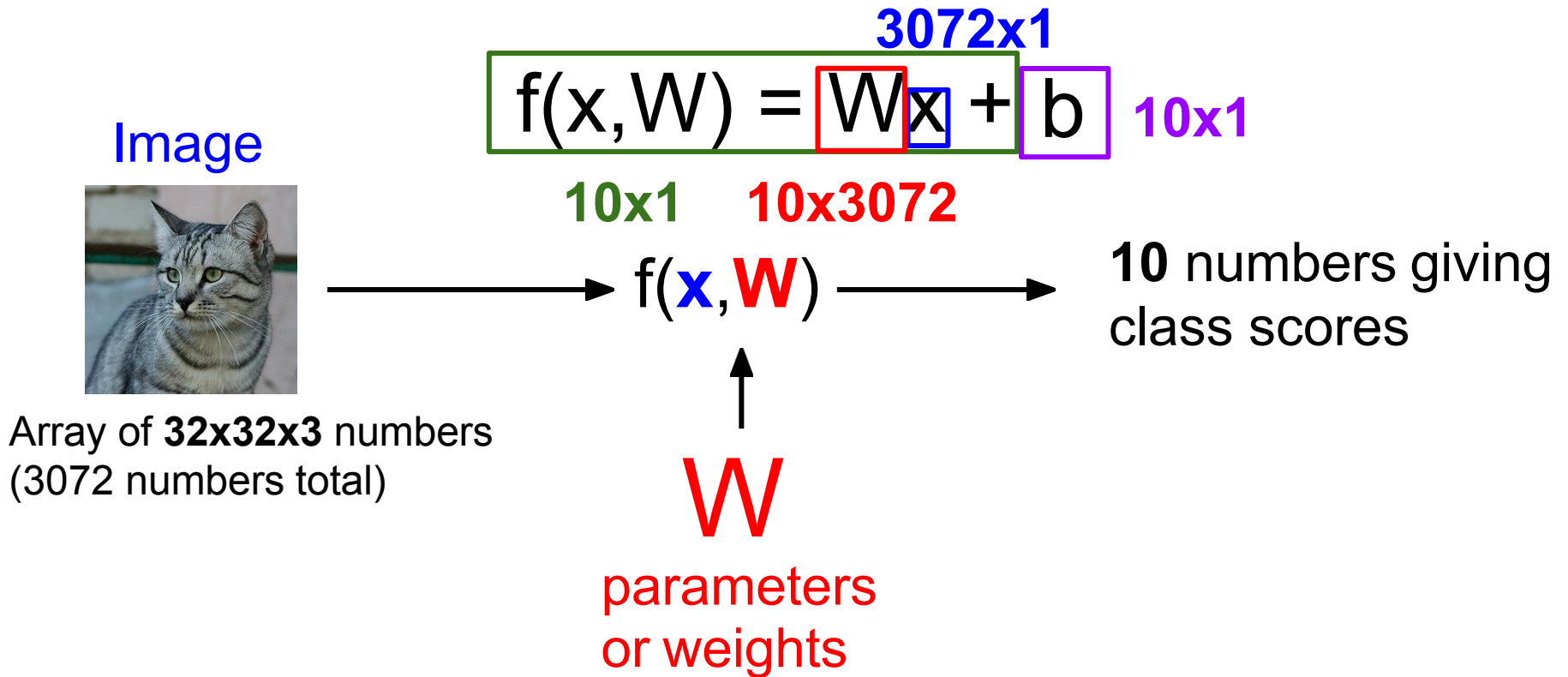
W

parameters
or weights

Parametric Approach: Linear Classifier

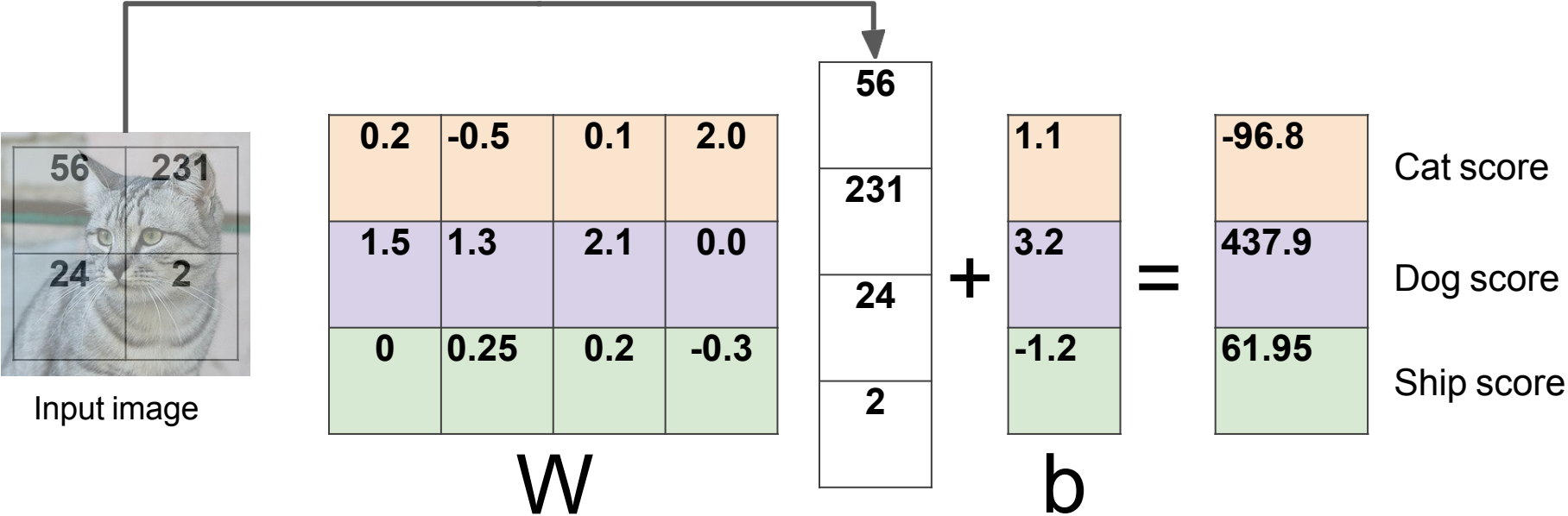


Parametric Approach: Linear Classifier



Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

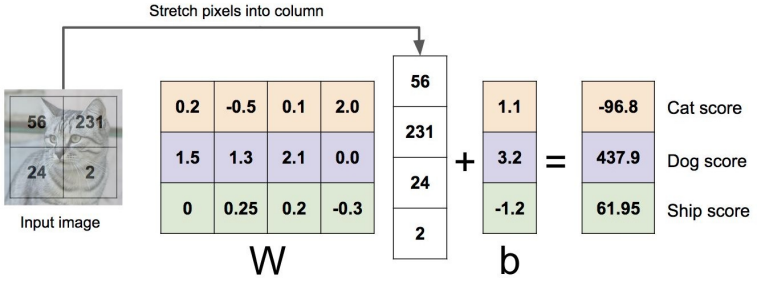
Stretch pixels into column



Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

Algebraic Viewpoint

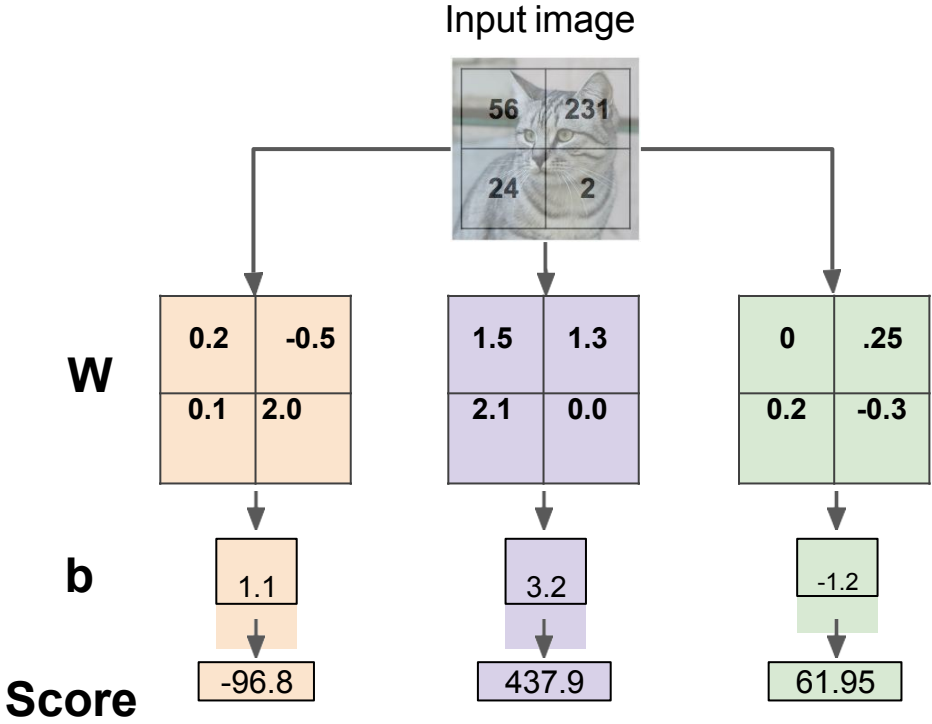
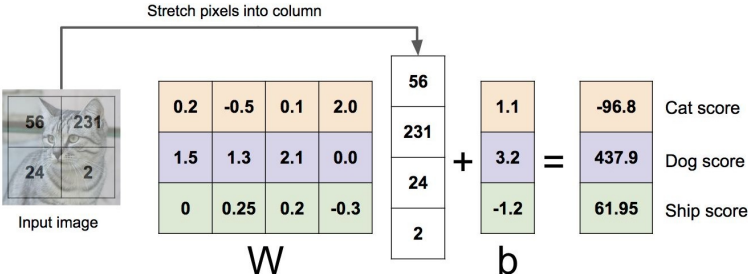
$$f(x,W) = Wx$$



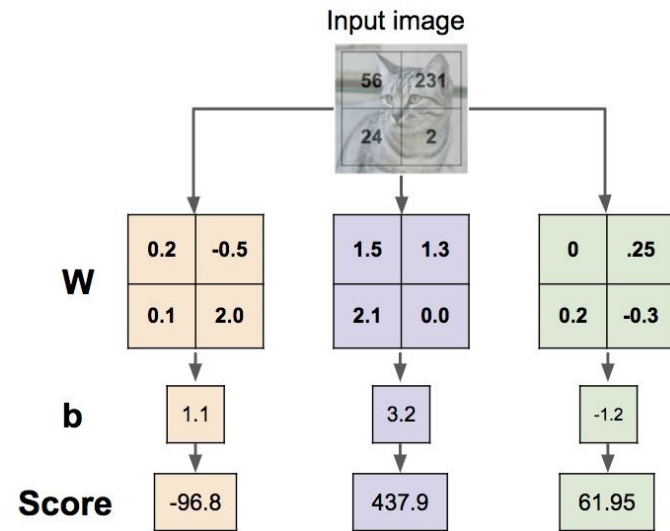
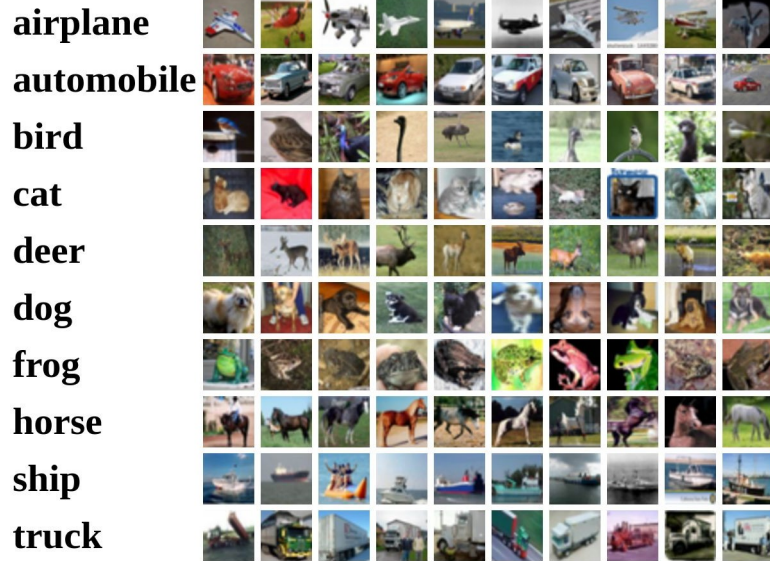
Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

Algebraic Viewpoint

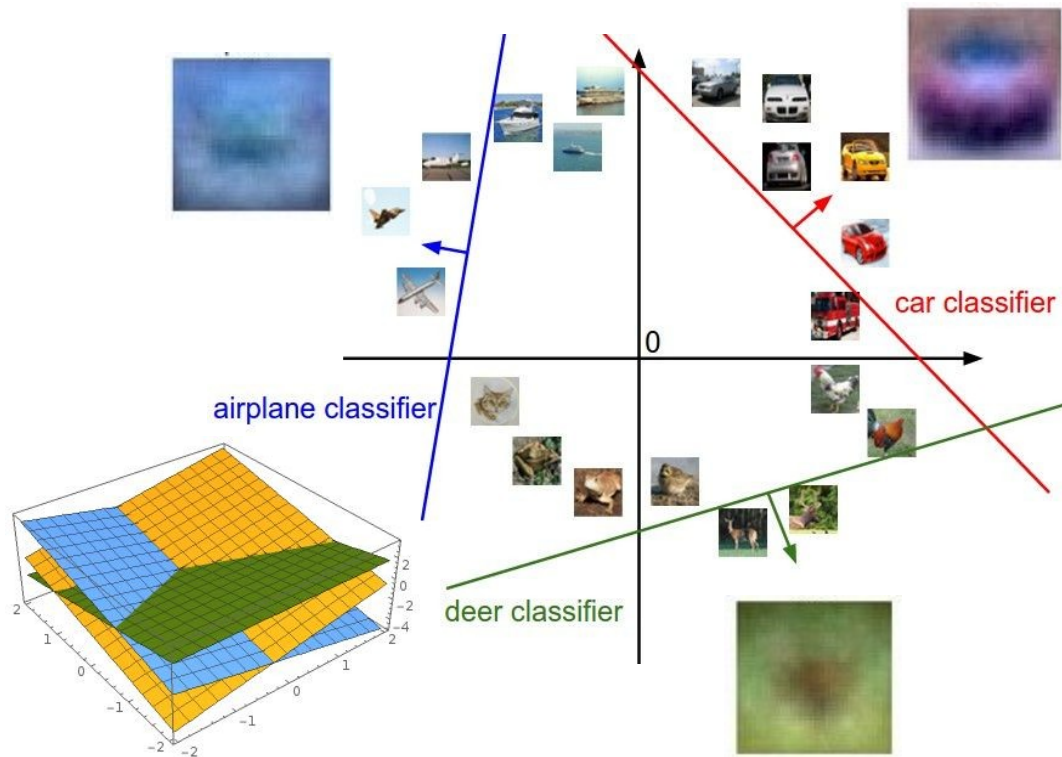
$$f(x,W) = Wx$$



Interpreting a Linear Classifier



Interpreting a Linear Classifier: Geometric Viewpoint



$$f(x, W) = Wx + b$$

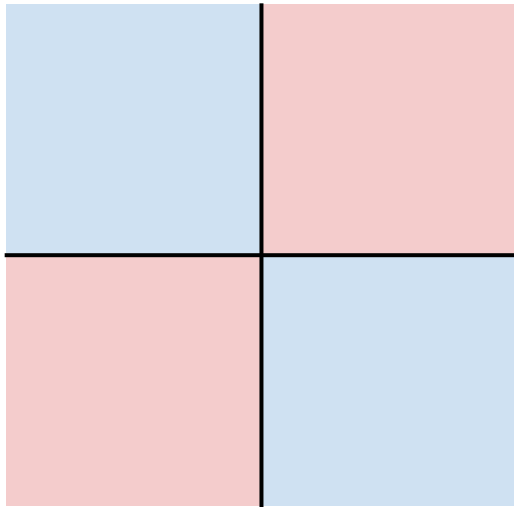


Array of **32x32x3** numbers
(3072 numbers total)

Hard cases for a linear classifier

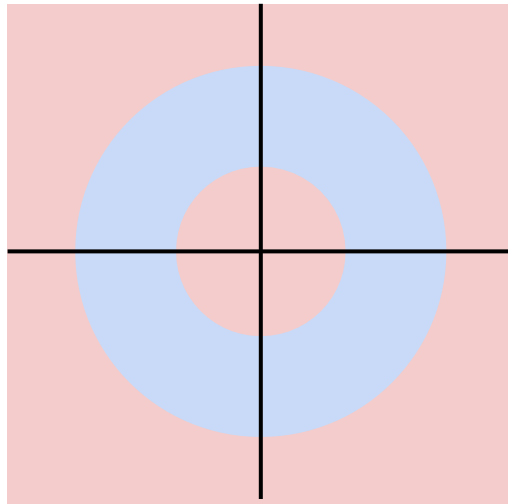
Class 1:
First and third quadrants

Class 2:
Second and fourth quadrants



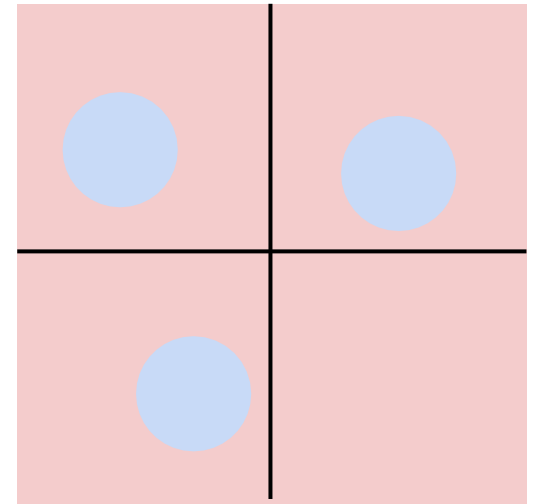
Class 1:
 $1 \leq \text{L2 norm} \leq 2$

Class 2:
Everything else



Class 1:
Three modes

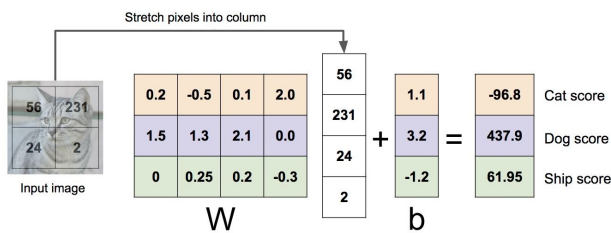
Class 2:
Everything else



Linear Classifier: Three Viewpoints

Algebraic Viewpoint

$$f(x, W) = Wx$$



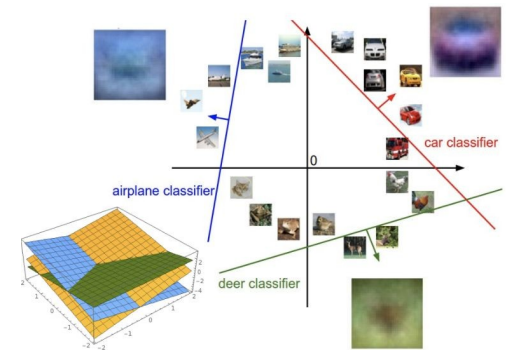
Visual Viewpoint

One template per class

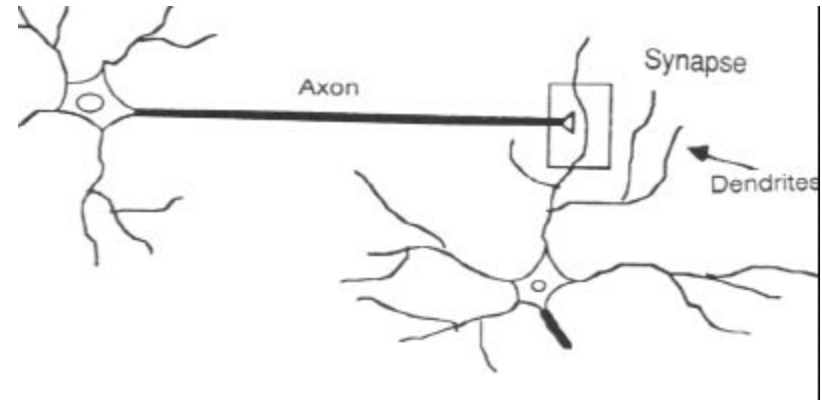
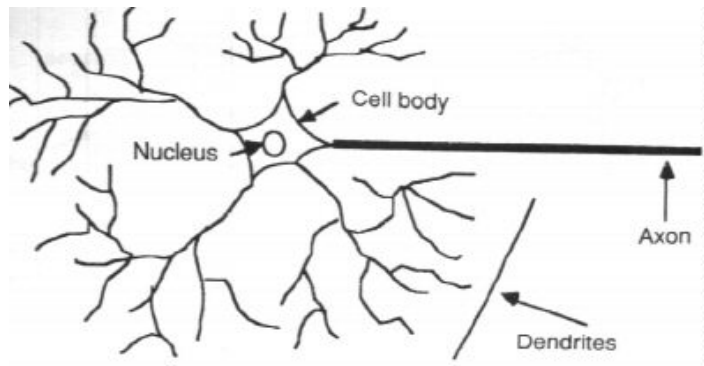


Geometric Viewpoint

Hyperplanes cutting up space

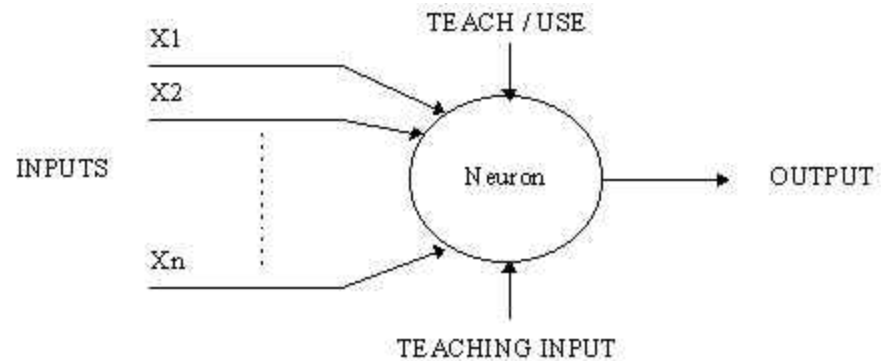


How the Human Brain learns



- In the human brain, a typical neuron collects signals from others through a host of fine structures called *dendrites*.
- The neuron sends out spikes of electrical activity through a long, thin strand known as an *axon*, which splits into thousands of branches.
- At the end of each branch, a structure called a *synapse* converts the activity from the axon into electrical effects that inhibit or excite activity in the connected neurons.

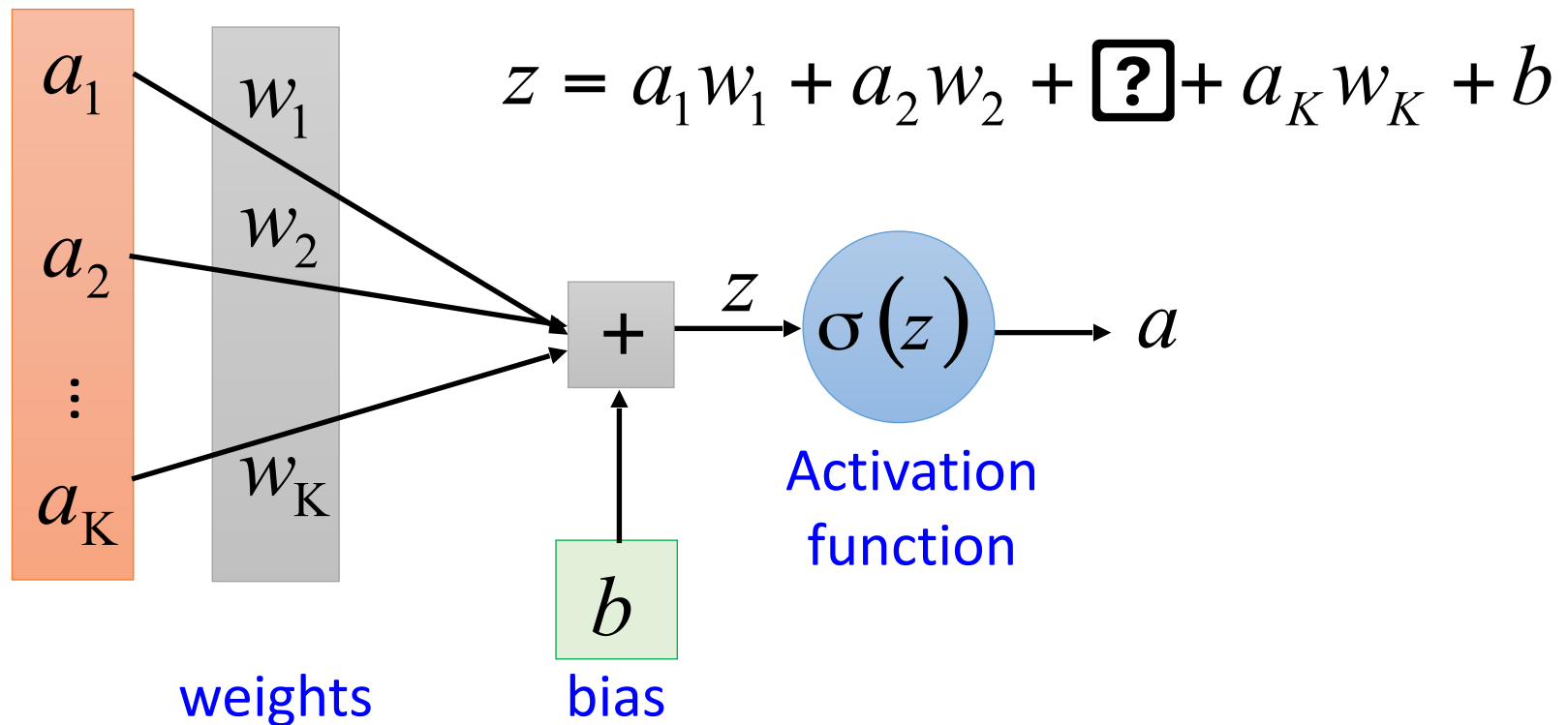
A Simple Neuron



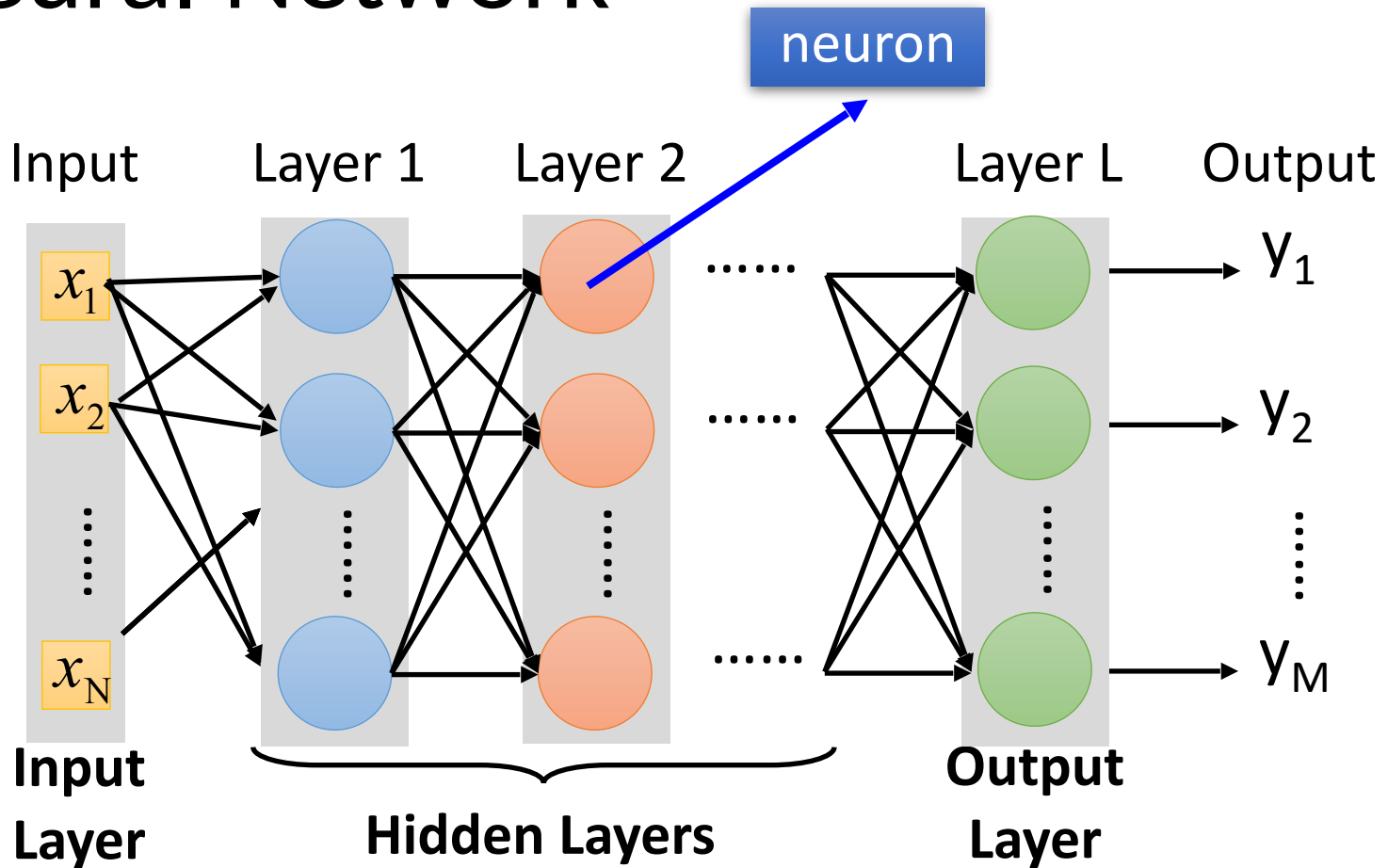
- An artificial neuron is a device with many inputs and one output.

Element of Neural Network

Neuron $f: R^K \rightarrow R$

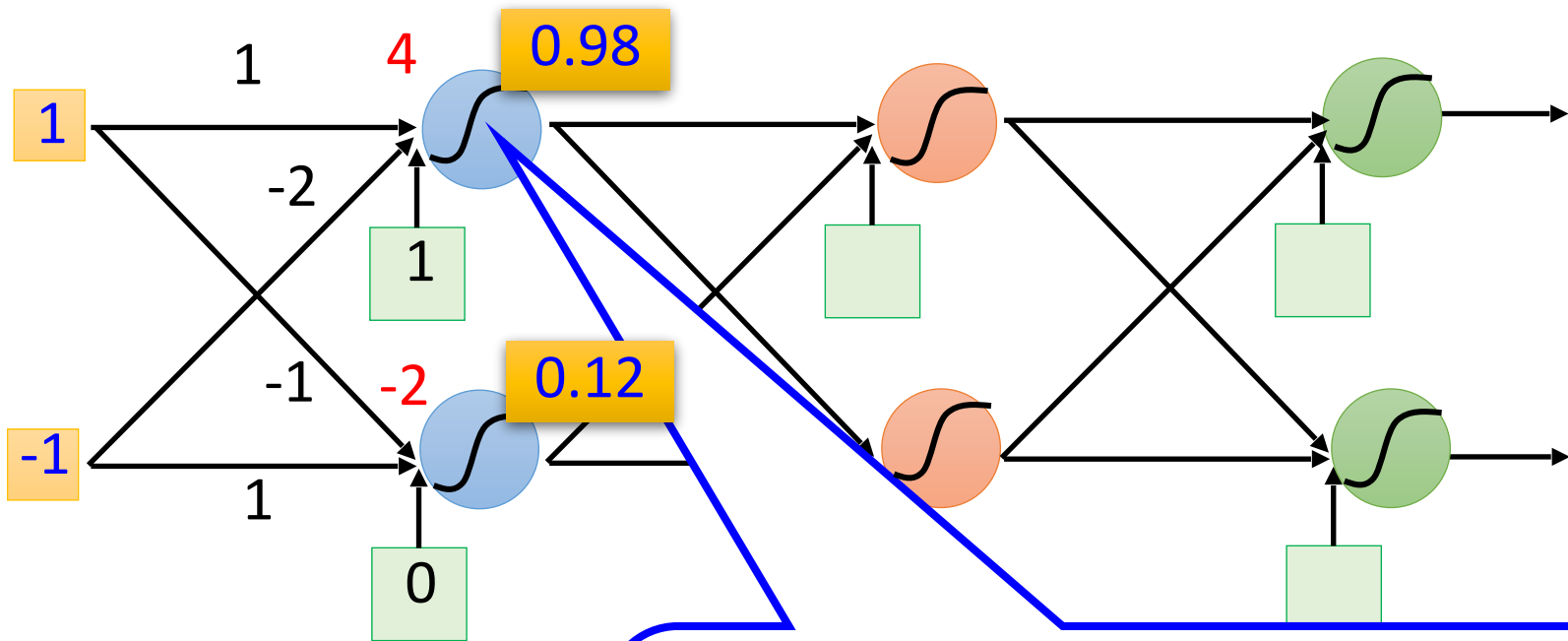


Neural Network



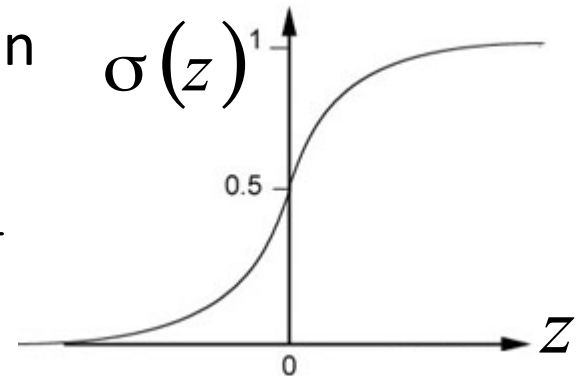
Deep means many hidden layers

Example of Neural Network



Sigmoid Function

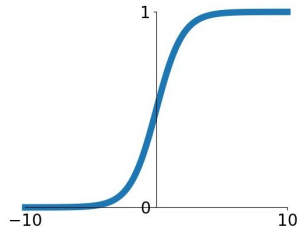
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



Activation functions

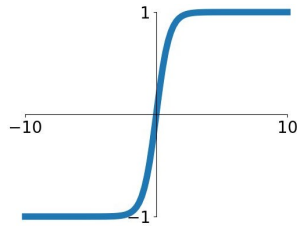
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



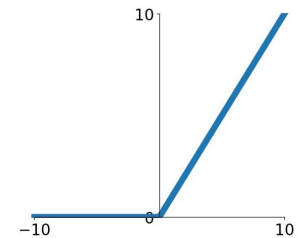
tanh

$$\tanh(x)$$



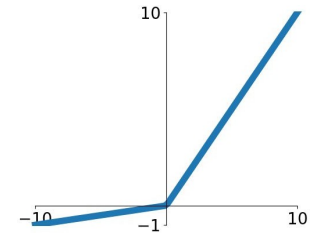
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

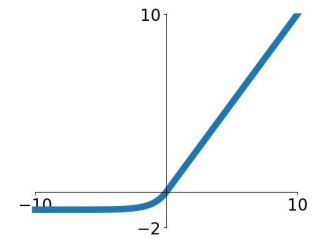


Maxout

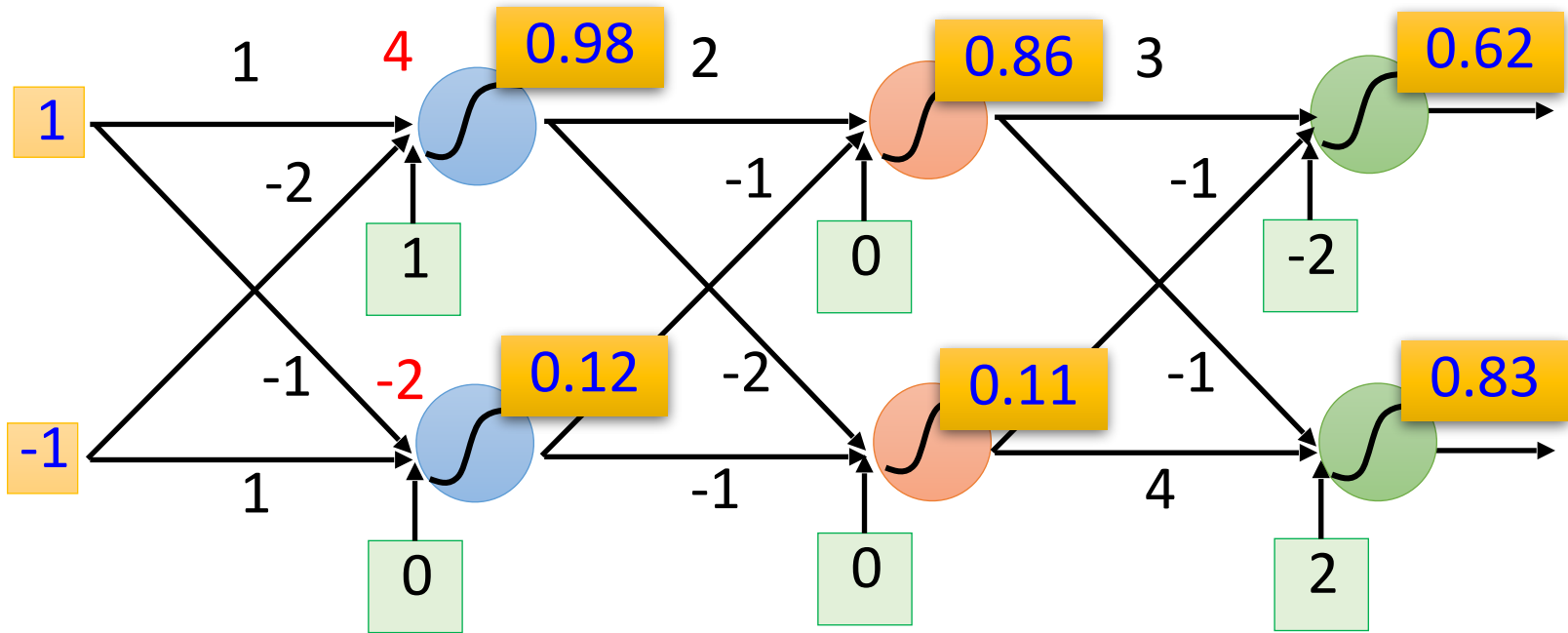
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

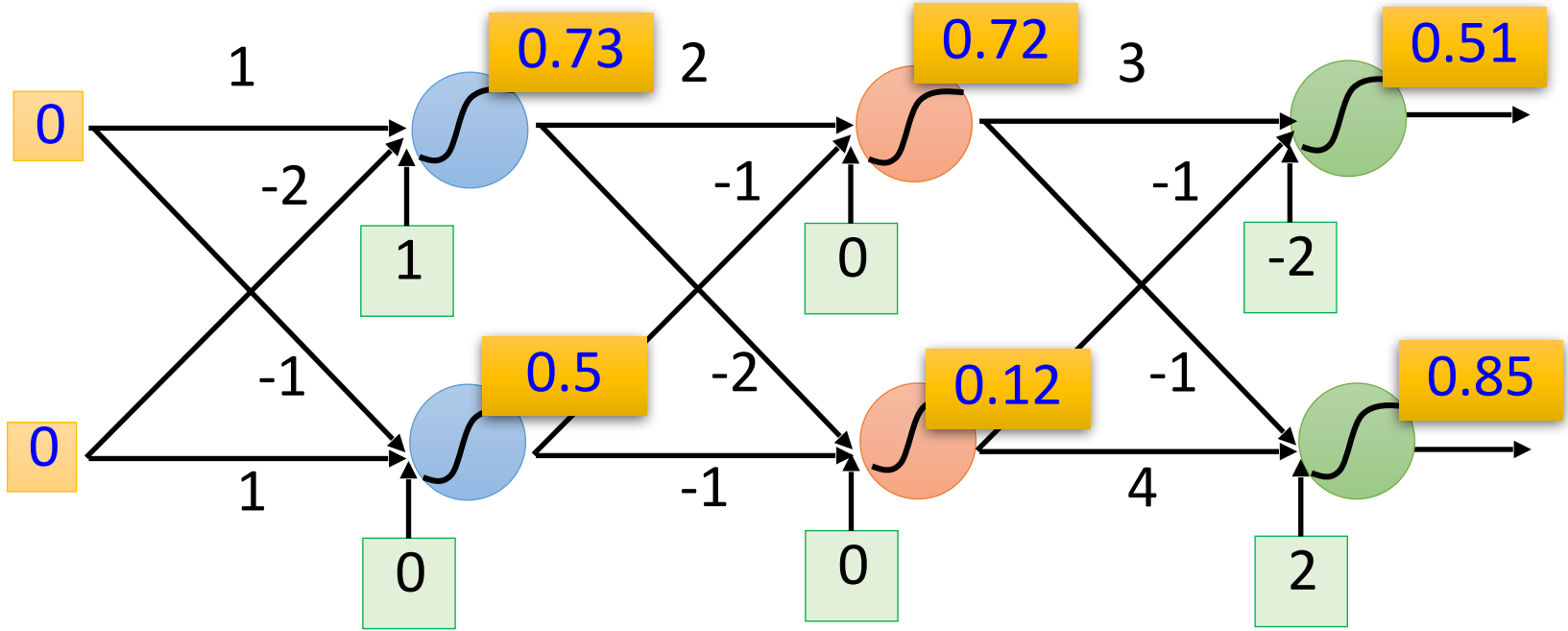
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Example of Neural Network



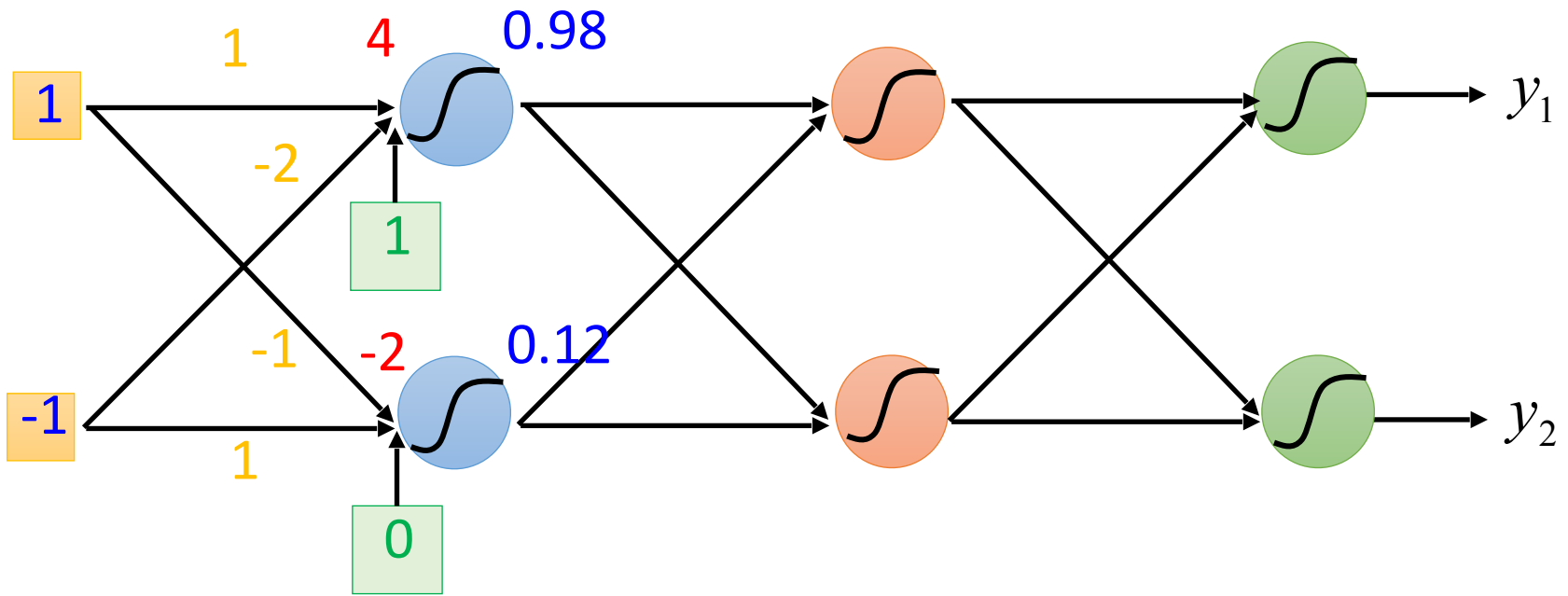
Example of Neural Network



$$f: \mathbb{R}^2 \rightarrow \mathbb{R}^2 \quad f\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix}\right) = \begin{bmatrix} 0.62 \\ 0.83 \end{bmatrix} \quad f\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}\right) = \begin{bmatrix} 0.51 \\ 0.85 \end{bmatrix}$$

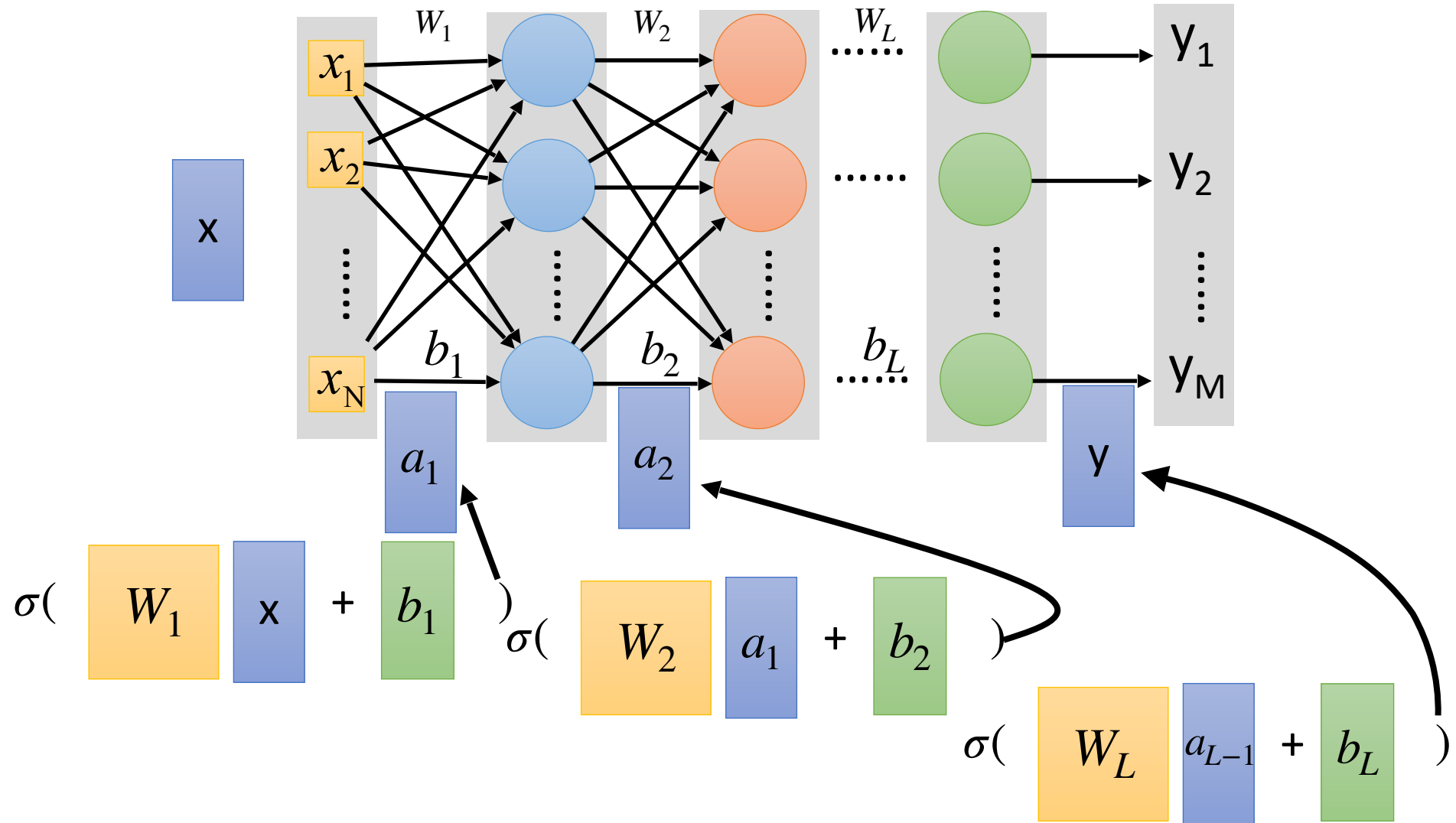
Different parameters define different function

Matrix Operation

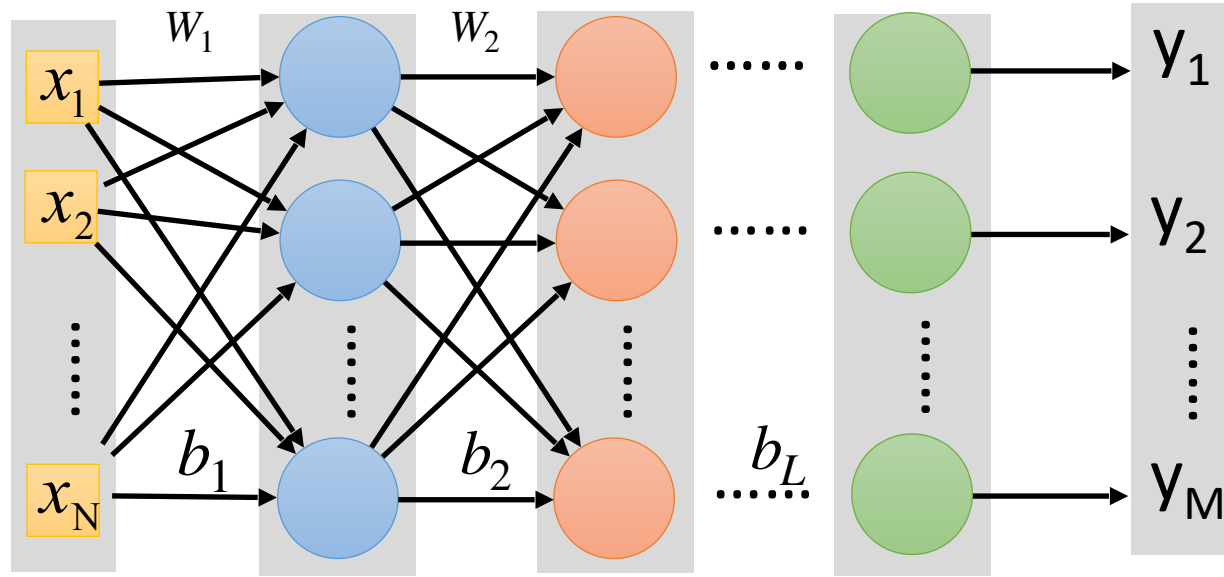


$$\sigma\left(\underbrace{\begin{bmatrix} 1 & -2 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix}}_{\begin{bmatrix} 4 \\ -2 \end{bmatrix}} \right) = \begin{bmatrix} 0.98 \\ 0.12 \end{bmatrix}$$

Neural Network



Neural Network



$$\mathbf{y} = f(\mathbf{x})$$

Using parallel computing techniques to speed up matrix operation

$$\sigma(W_L \cdots \sigma(W_2 \sigma(W_1 x + b_1) + b_2) + \cdots b_L)$$

Softmax

- Softmax layer as the output layer

Ordinary Layer

$$z_1 \longrightarrow \sigma \longrightarrow y_1 = \sigma(z_1)$$

$$z_2 \longrightarrow \sigma \longrightarrow y_2 = \sigma(z_2)$$

$$z_3 \longrightarrow \sigma \longrightarrow y_3 = \sigma(z_3)$$

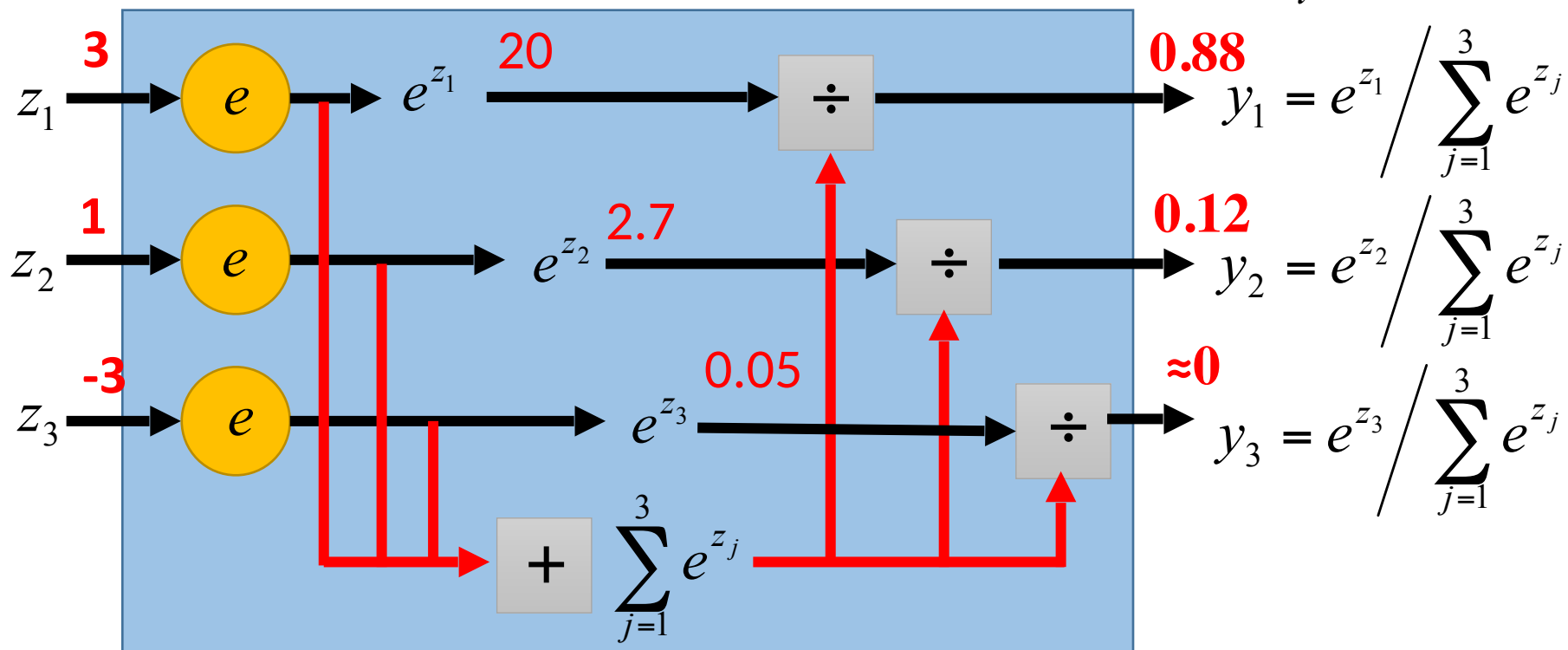
In general, the output of network can be any value.

May not be easy to interpret

Softmax

- Softmax layer as the output layer

Softmax Layer

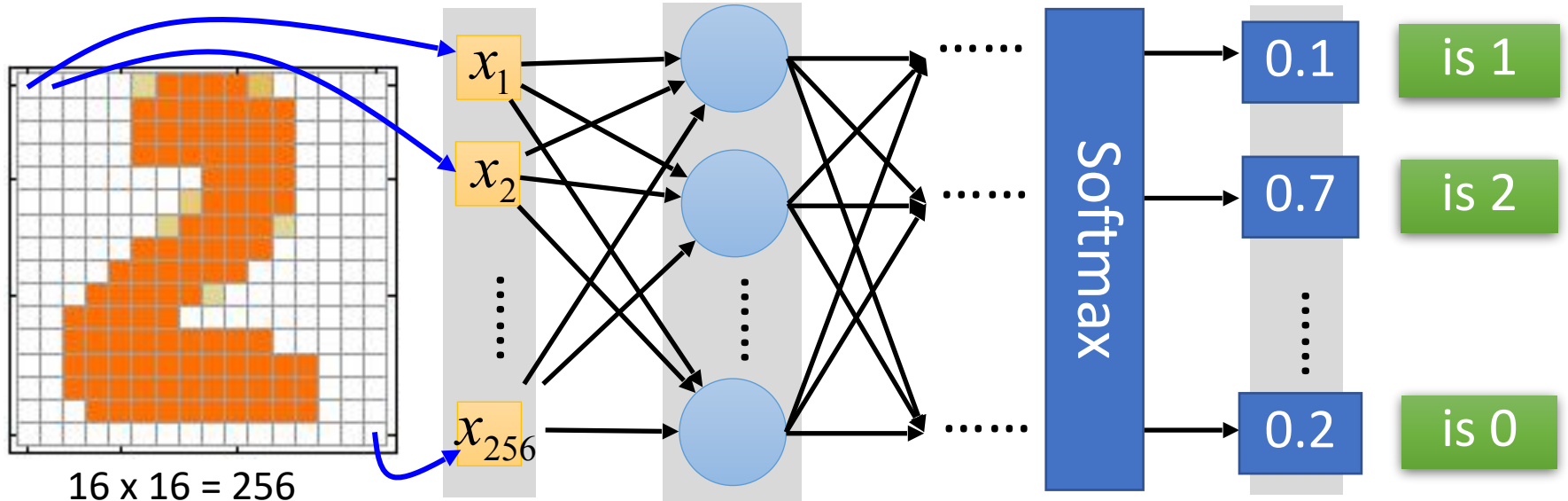


Probability:

- $1 > y_i > 0$
- $\sum_i y_i = 1$

How to set network parameters

$$\theta = \{W_1, b_1, \dots, W_n, b_n\}$$



16 x 16 = 256

Ink \rightarrow 1

No ink \rightarrow 0

Set the network parameters θ such that

Input x_1 has the maximum value

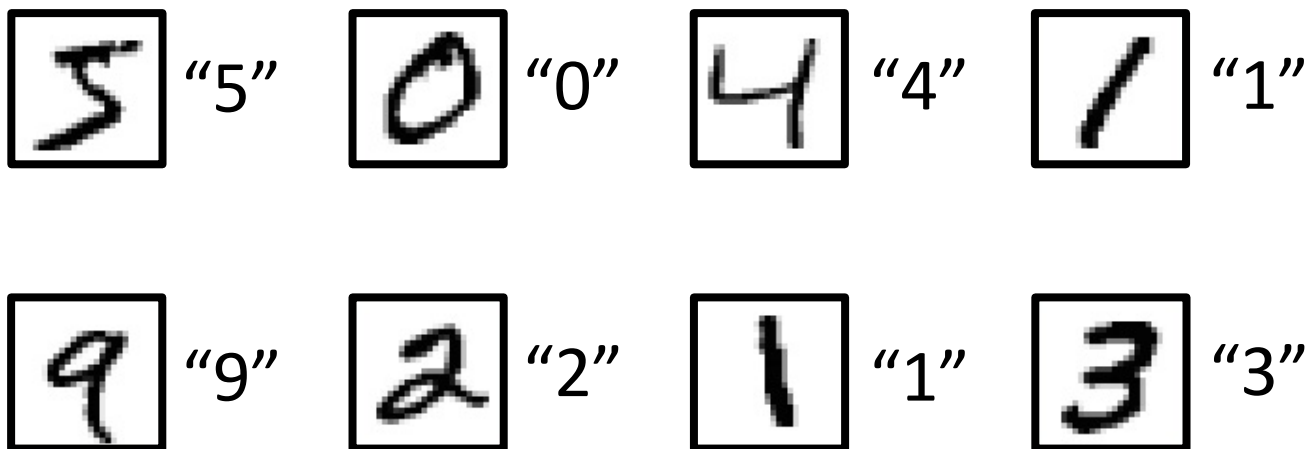
Input x_2 has the maximum value

How to let the neural network achieve this



Training Data

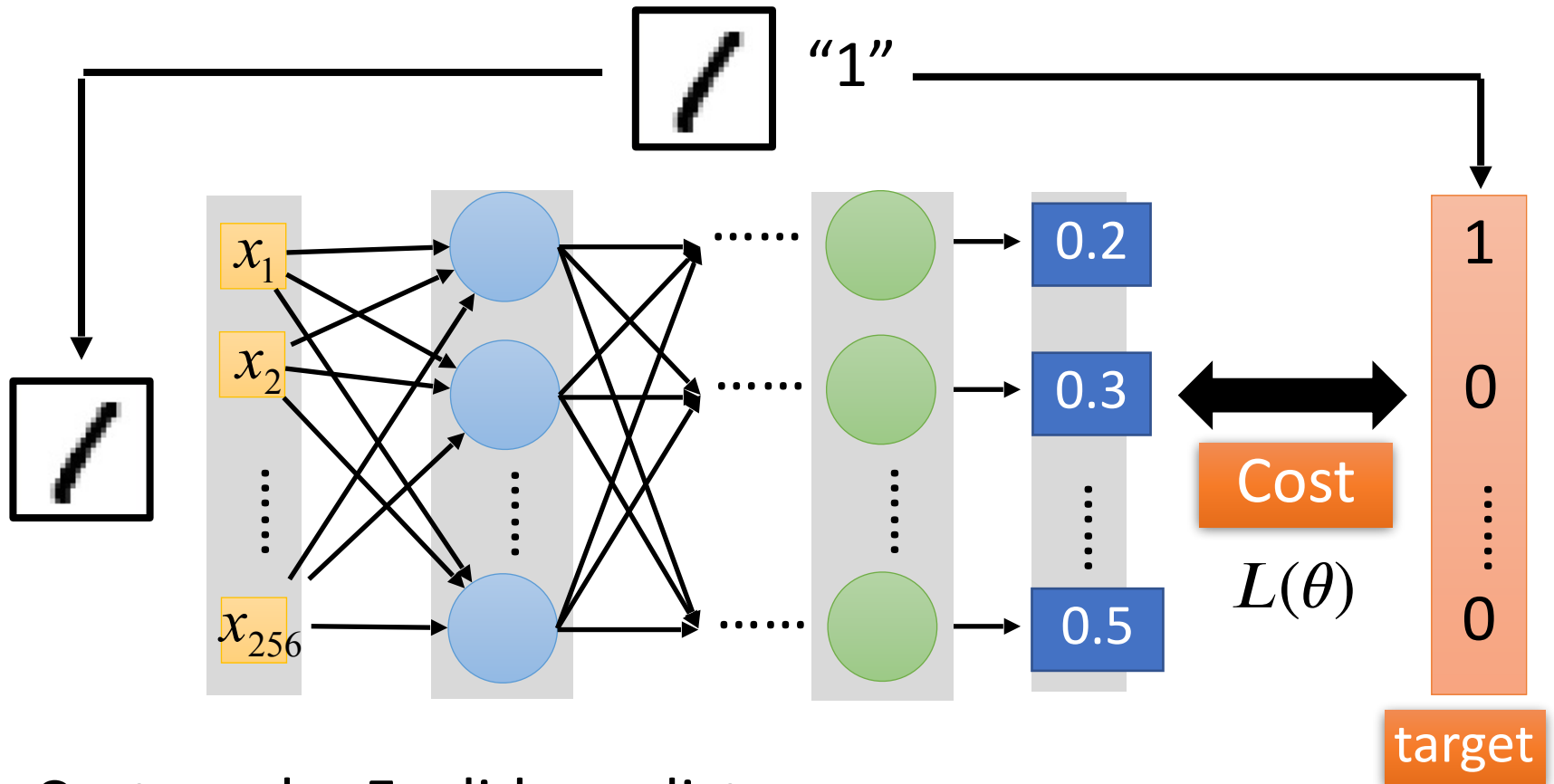
- Preparing training data: images and their labels



Using the training data to find
the network parameters.

Cost

Given a set of network parameters θ , each example has a cost value.



Cost can be Euclidean distance or cross entropy of the network output and target

Soft-entropy Loss

The score of label category is larger than other categories:

$$score_{\text{label}} > score_j \text{ for any } j \neq \text{label}$$

How to set up a loss for this goal?

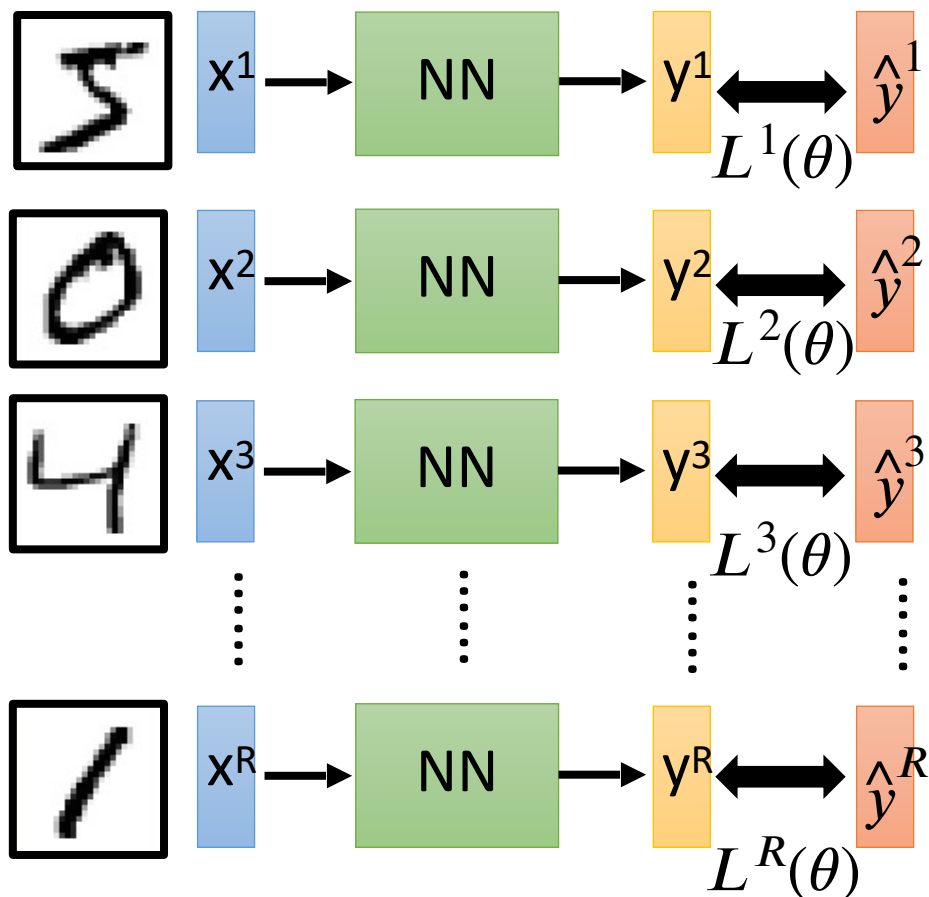
Soft-entropy Loss

$$\text{Let } score_{label} = \frac{e^{f(x, W)_{label}}}{\sum_j e^{f(x, W)_j}}$$

If we set $\ell(f(x; W, b), label) = -\log score_{label}$

Total Cost

For all training data ...



Total Cost:

$$C(\theta) = \sum_{r=1}^R L^r(\theta)$$

How bad the network parameters θ is on this task

Find the network parameters θ^* that minimize this value

Gradient Descent


Error Surface

Assume there are only two parameters w_1 and w_2 in a network.


$$\theta = \{w_1, w_2\}$$

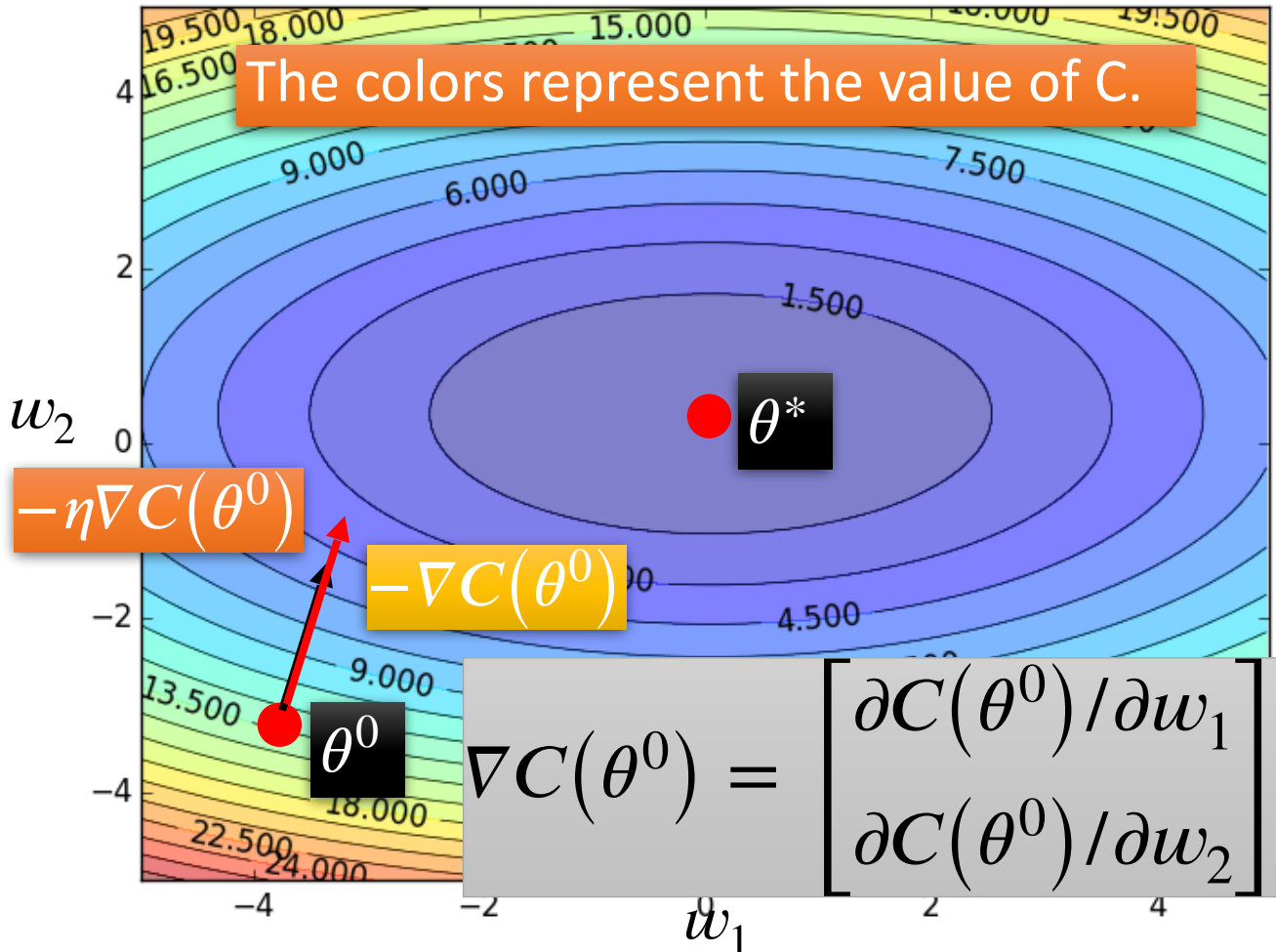
Randomly pick a starting point θ^0

Compute the negative gradient at θ^0

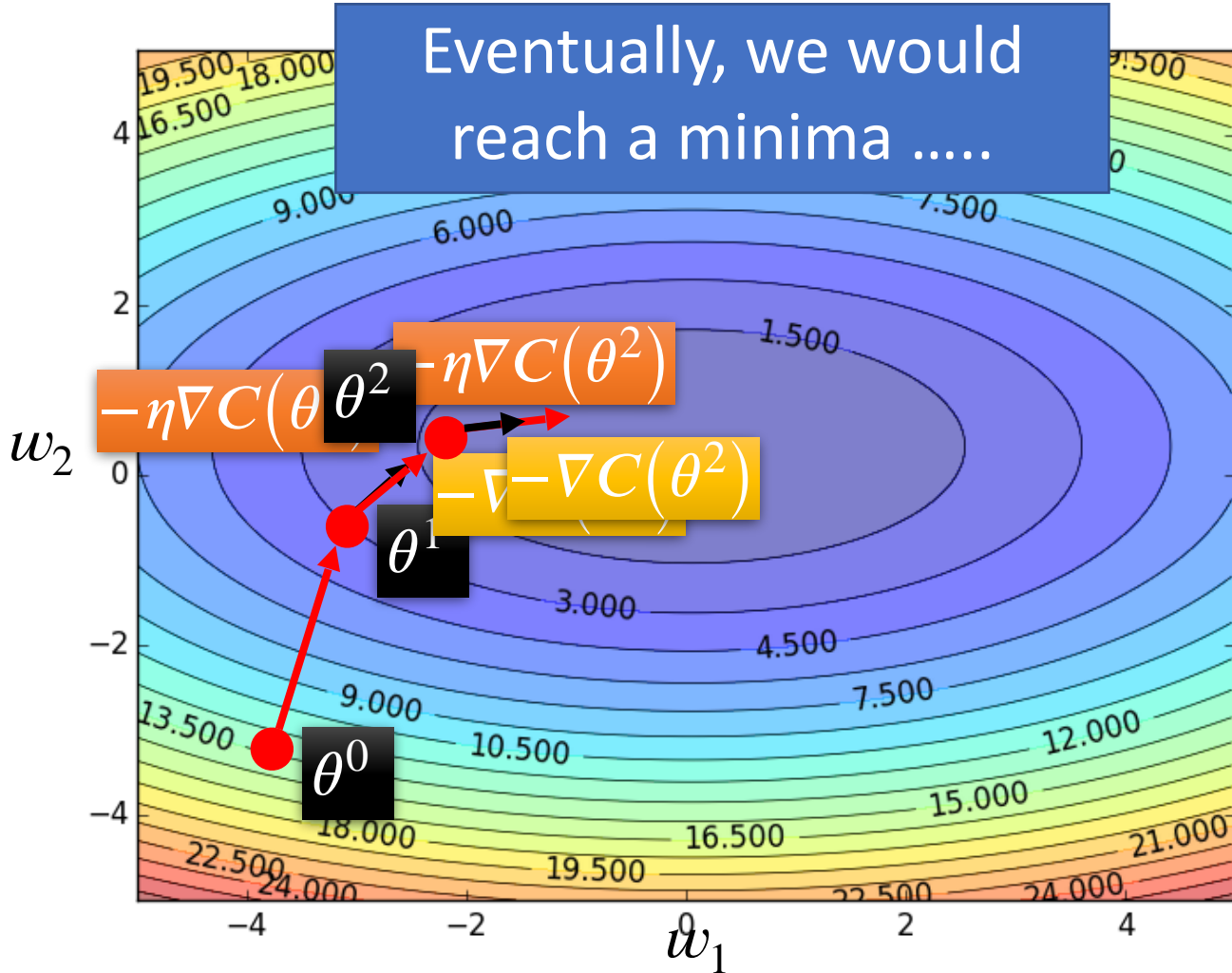
 $-\nabla C(\theta^0)$

Times the learning rate η

 $-\eta \nabla C(\theta^0)$




Gradient Descent




Randomly pick a starting point θ^0

Compute the negative gradient at θ^0

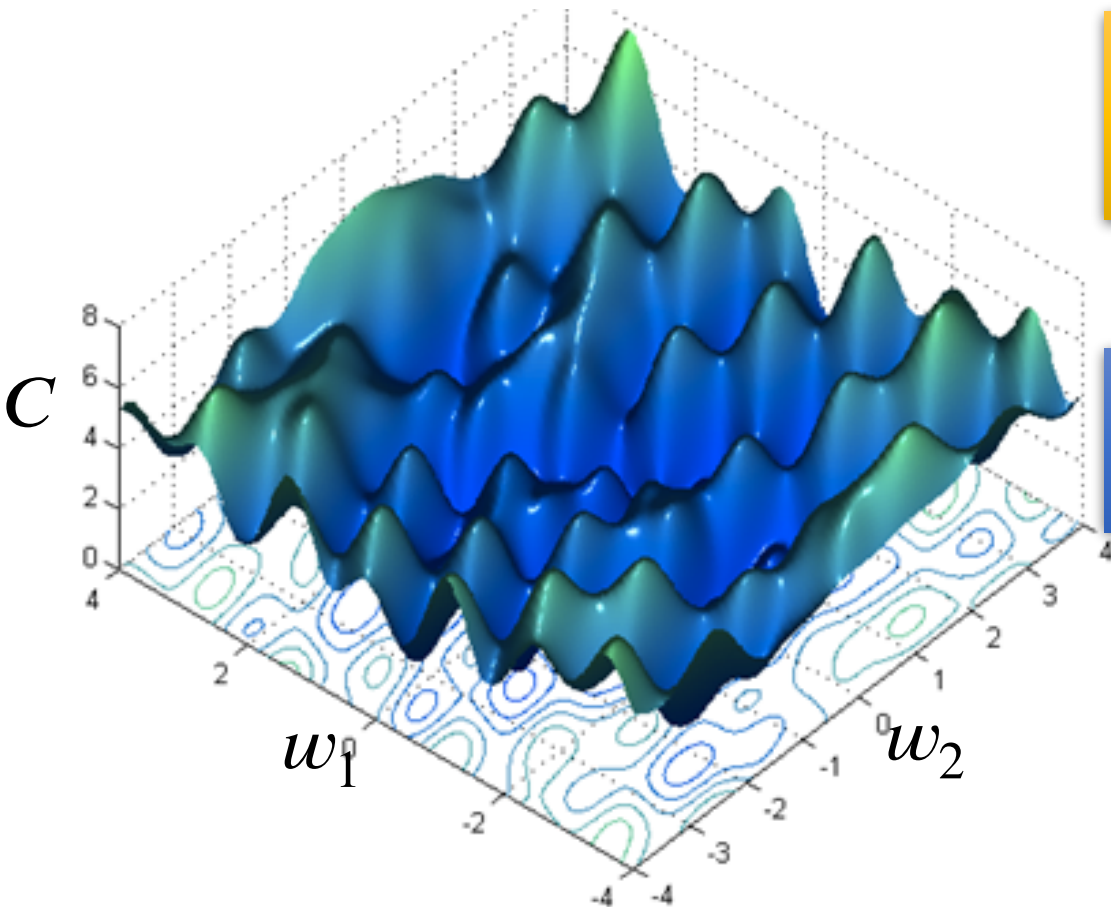
 $-\nabla C(\theta^0)$

Times the learning rate η

 $-\eta \nabla C(\theta^0)$

Local Minima

- Gradient descent never guarantee global minima



Different initial
point θ^0

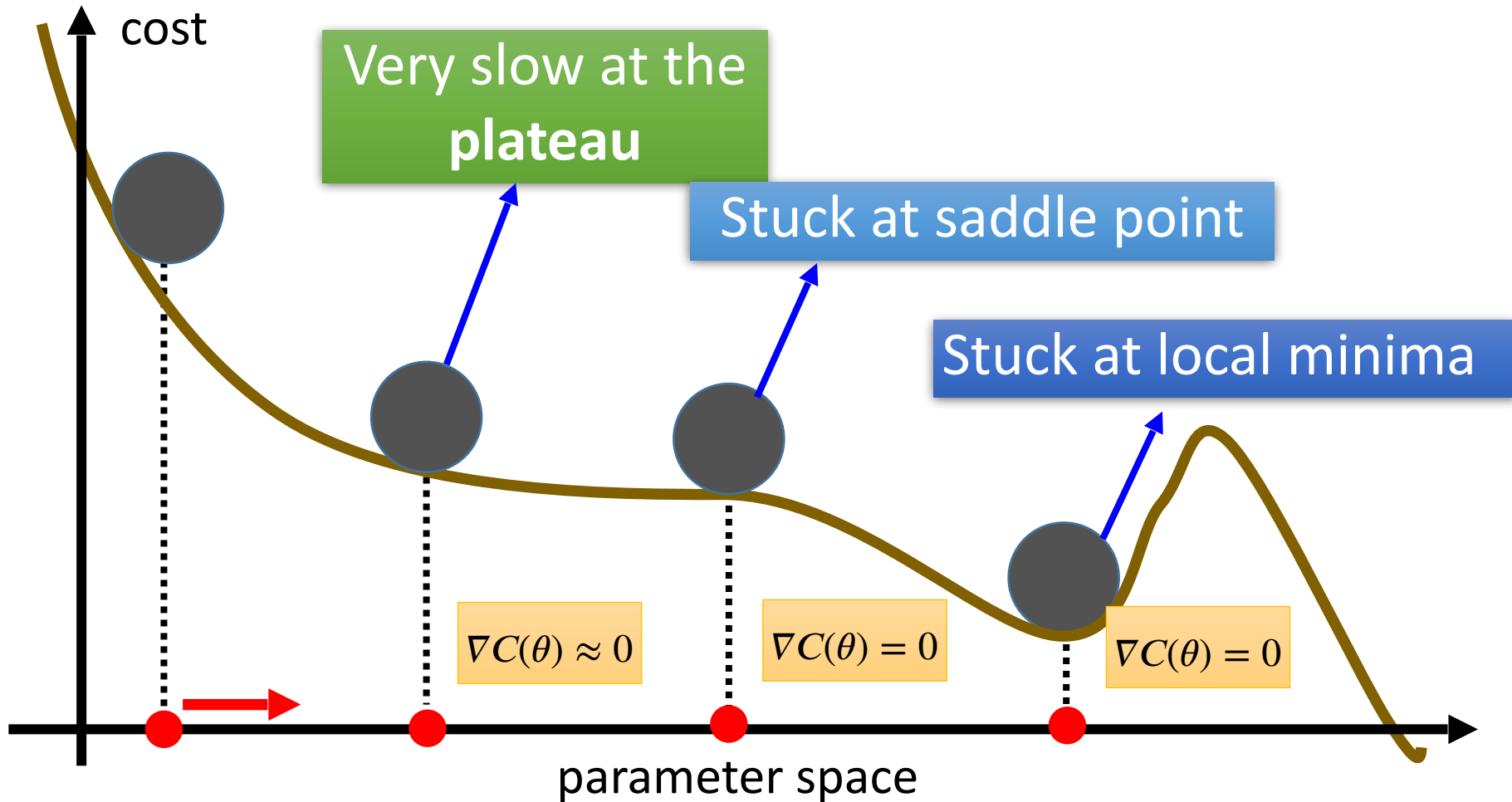


Reach different minima,
so different results

Who is Afraid of Non-Convex
Loss Functions?

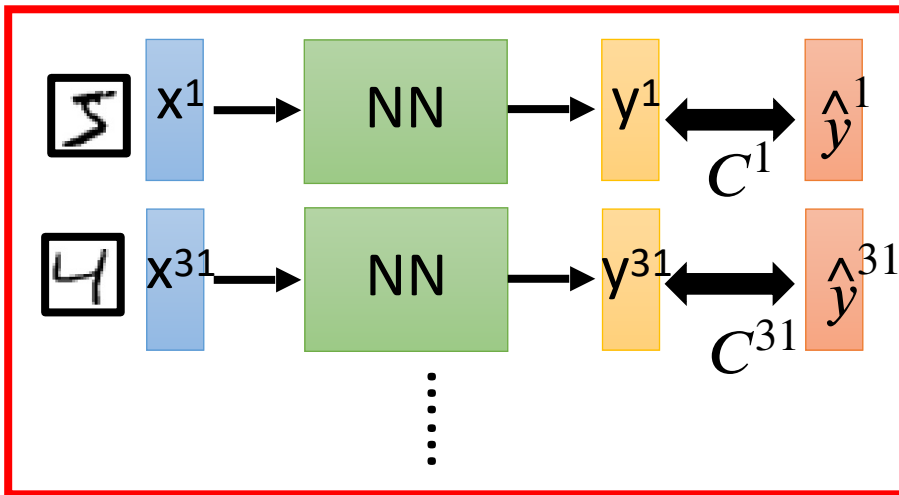
[http://videolectures.net/
eml07_lecun_wia/](http://videolectures.net/eml07_lecun_wia/)

Besides local minima

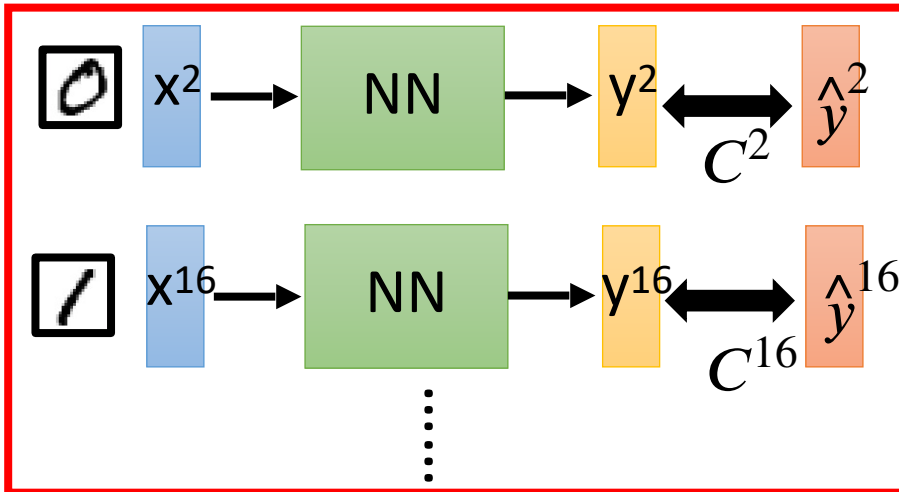


Mini-batch

Mini-batch



Mini-batch



➤ Randomly initialize θ^0

➤ Pick the 1st batch

$$C = C^1 + C^{31} + \dots$$

$$\theta^1 \leftarrow \theta^0 - \eta \nabla C(\theta^0)$$

➤ Pick the 2nd batch

$$C = C^2 + C^{16} + \dots$$

$$\theta^2 \leftarrow \theta^1 - \eta \nabla C(\theta^1)$$

⋮

➤ Until all mini-batches have been picked

one epoch

Repeat the above process

Backpropagation

- A network can have millions of parameters.
 - Backpropagation is the way to compute the gradients efficiently (not today)
 - Ref: http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS_2015_2/Lecture/DNN%20backprop.ecm.mp4/index.html
- Many toolkits can compute the gradients automatically

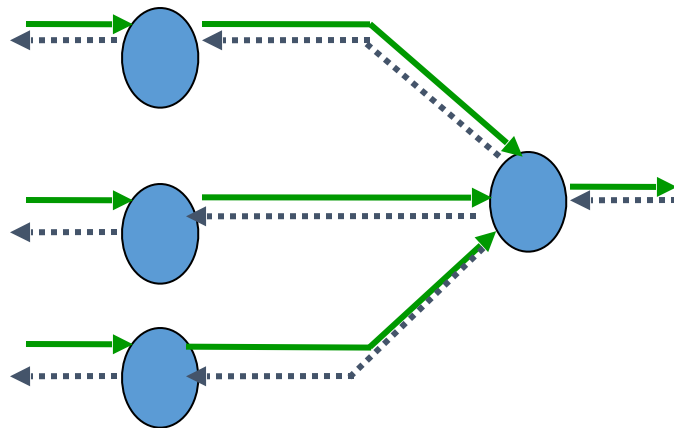
theano



Ref: http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS_2015_2/Lecture/Theano%20DNN.ecm.mp4/index.html

Back Propagation

- Back-propagation training algorithm



*Network activation
Forward Step*

*Error propagation
Backward Step*

- Backprop adjusts the weights of the NN in order to minimize the network total error.

Next: Convolutional Neural Networks

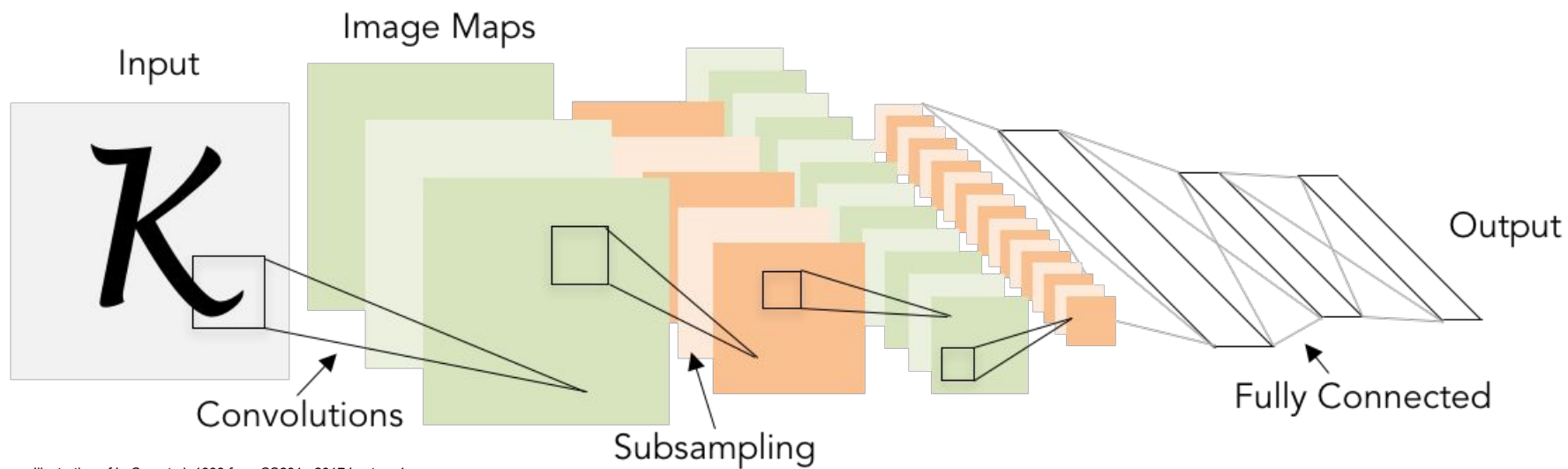
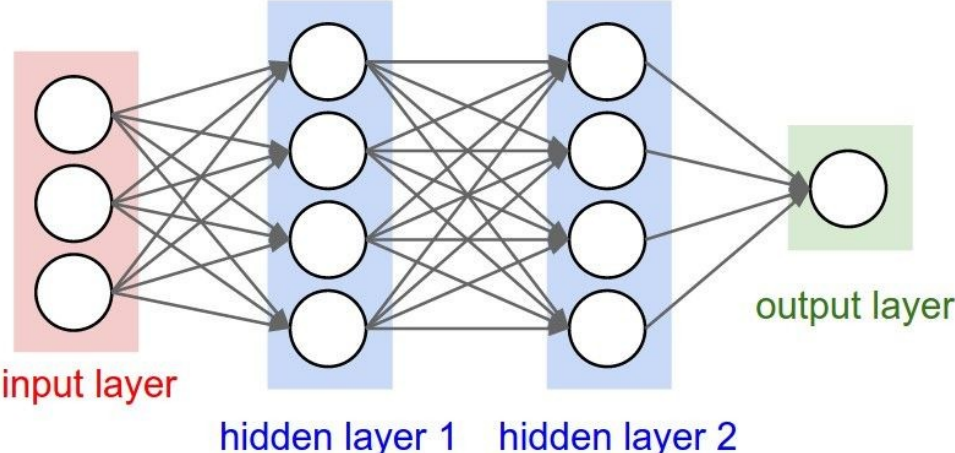
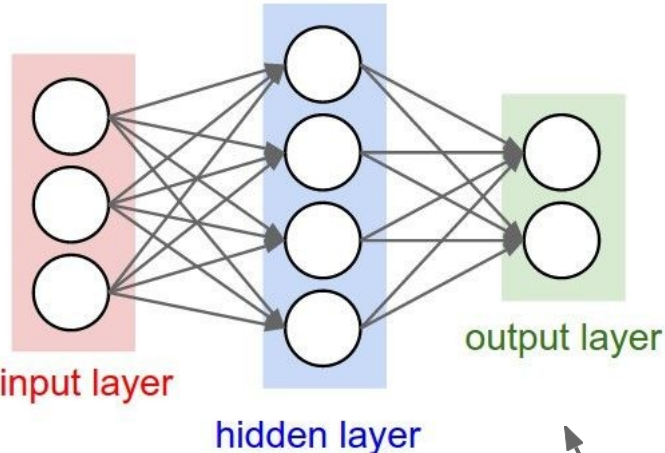


Illustration of LeCun et al. 1998 from CS231n 2017 Lecture 1

Neural networks: Architectures



"2-layer Neural Net", or
"1-hidden-layer Neural Net"

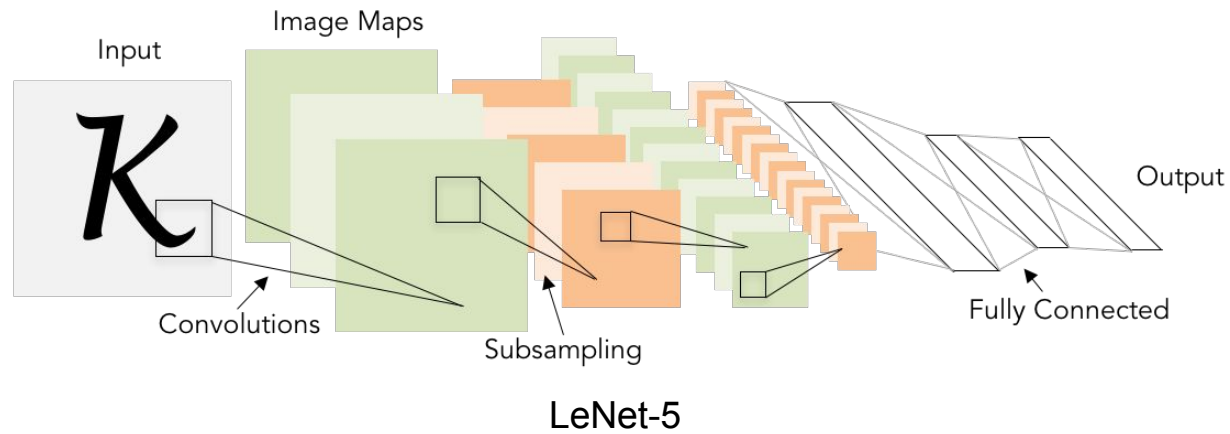
"3-layer Neural Net", or
"2-hidden-layer Neural Net"

"Fully-connected" layers

Gradient-based learning applied to document recognition

[LeCun, Bottou, Bengio, Haffner 1998]

A bit of history:



A bit of history:

ImageNet Classification with Deep Convolutional Neural Networks

[Krizhevsky, Sutskever, Hinton, 2012]

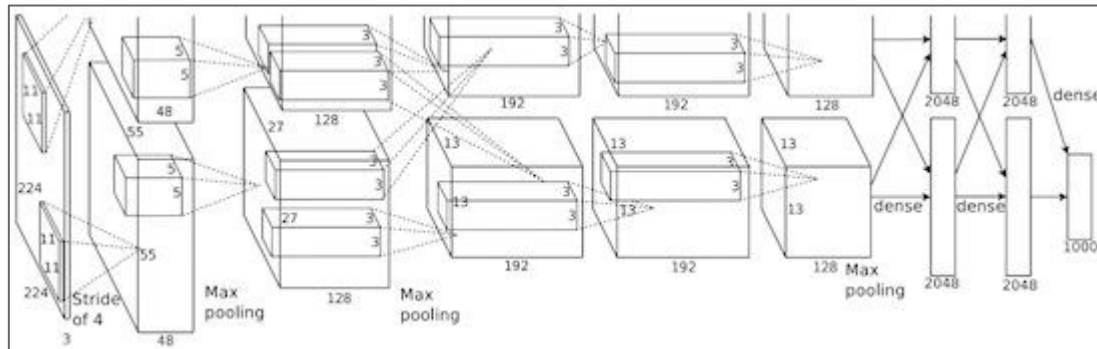


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

“AlexNet”

Fast-forward to today: ConvNets are everywhere



self-driving cars

Photo by Lane McIntosh. Copyright CS231n 2017.



[This image](#) by GBPublic_PR is licensed under [CC-BY 2.0](#)

NVIDIA Tesla line

(these are the GPUs on rye01.stanford.edu)

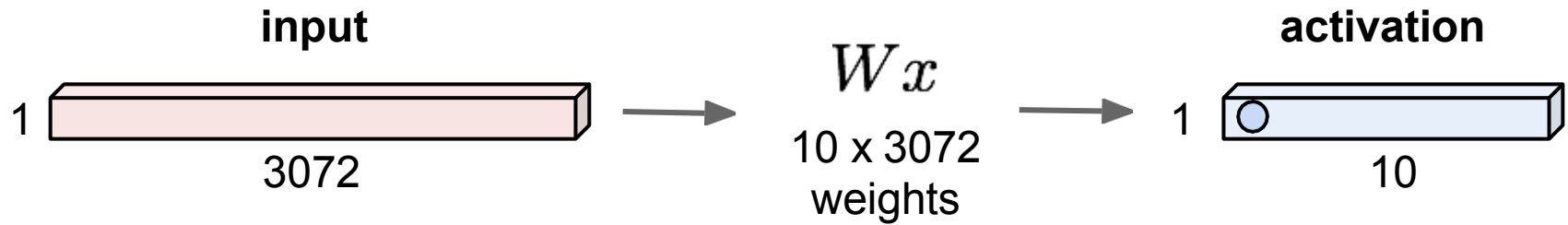
Note that for embedded systems a typical setup would involve NVIDIA Tegras, with integrated GPU and ARM-based CPU cores.

Convolutional Neural Networks

(First without the brain stuff)

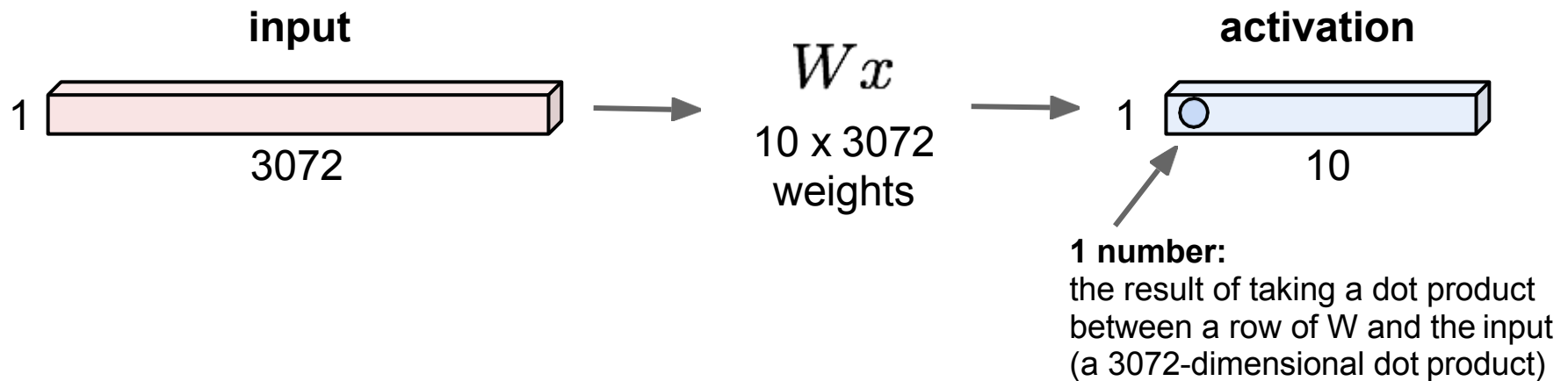
Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



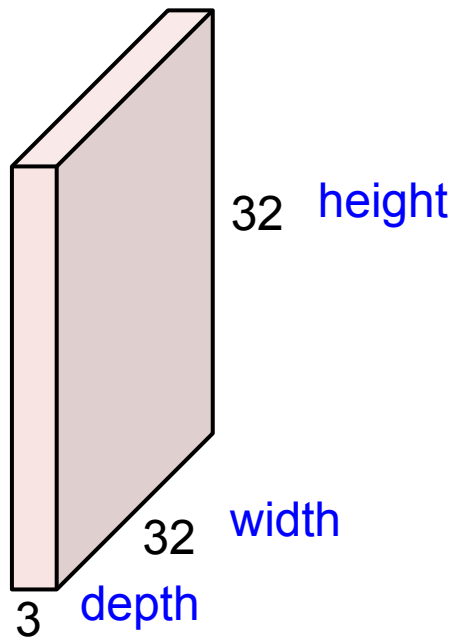
Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



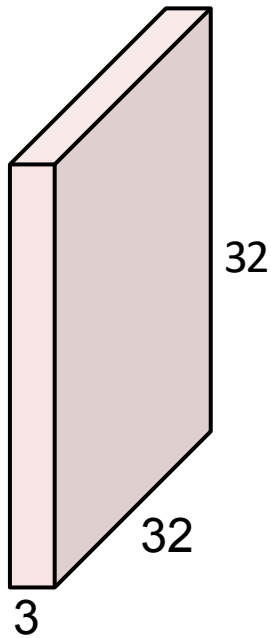
Convolution Layer

32x32x3 image -> preserve spatial structure



Convolution Layer

32x32x3 image

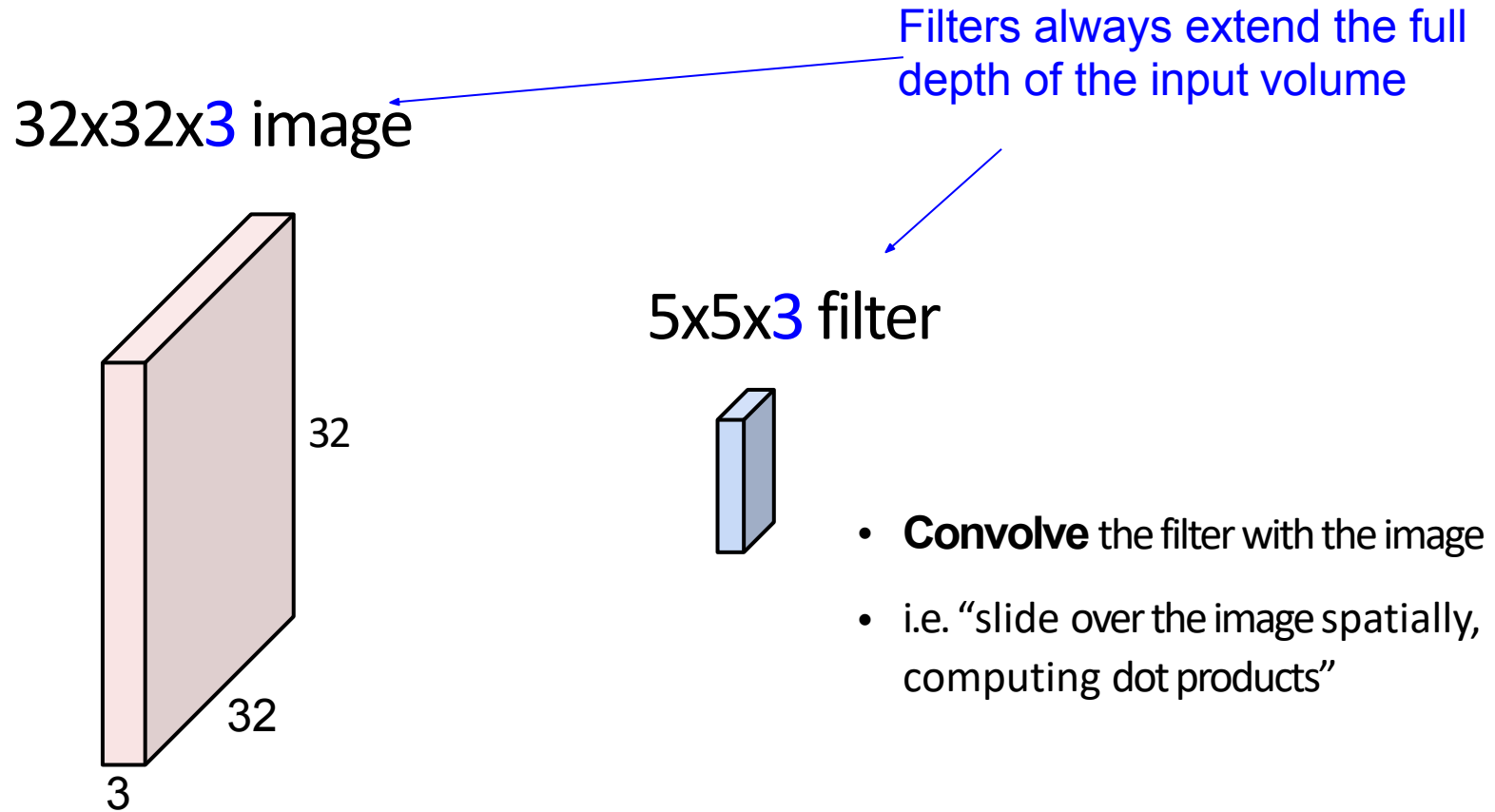


5x5x3 filter

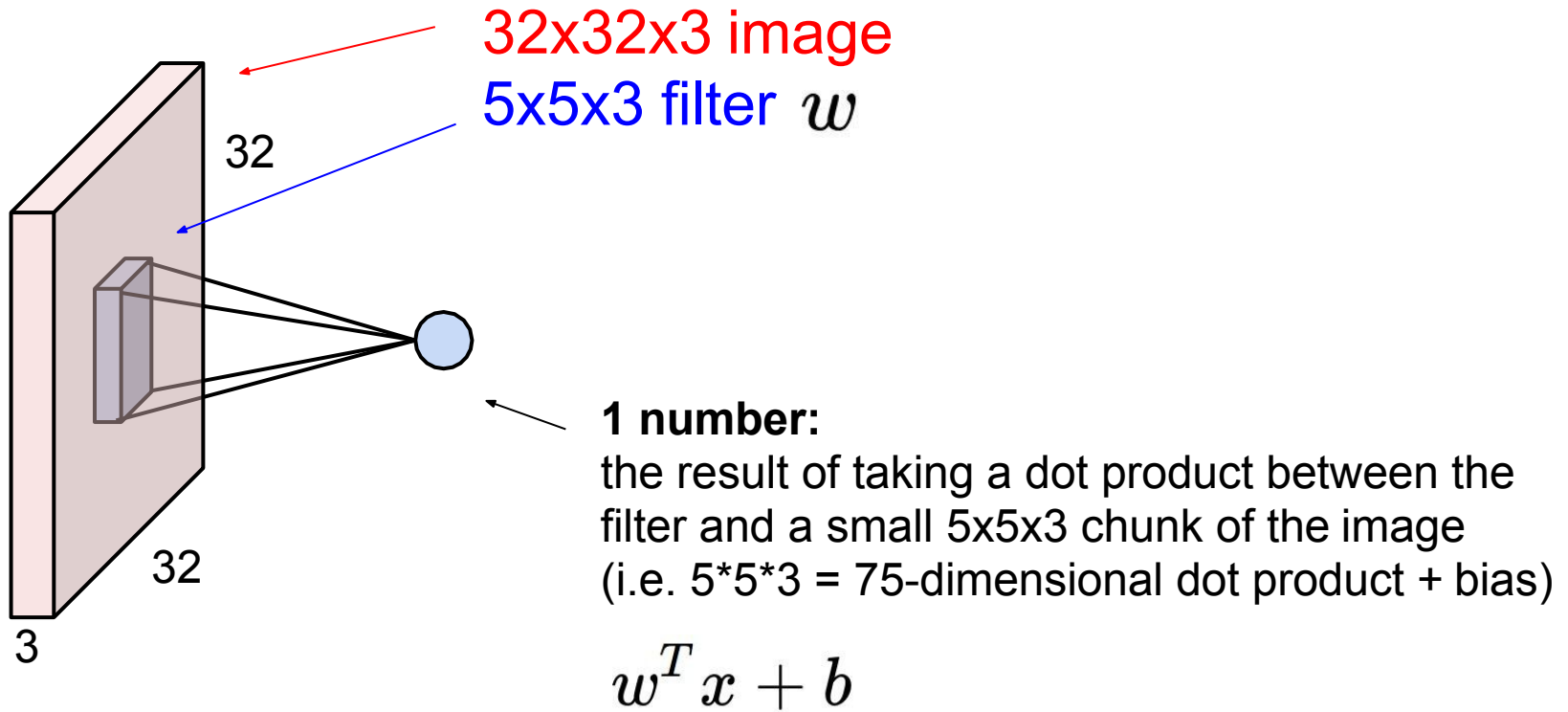


- **Convolve** the filter with the image
- i.e. “slide over the image spatially, computing dot products”

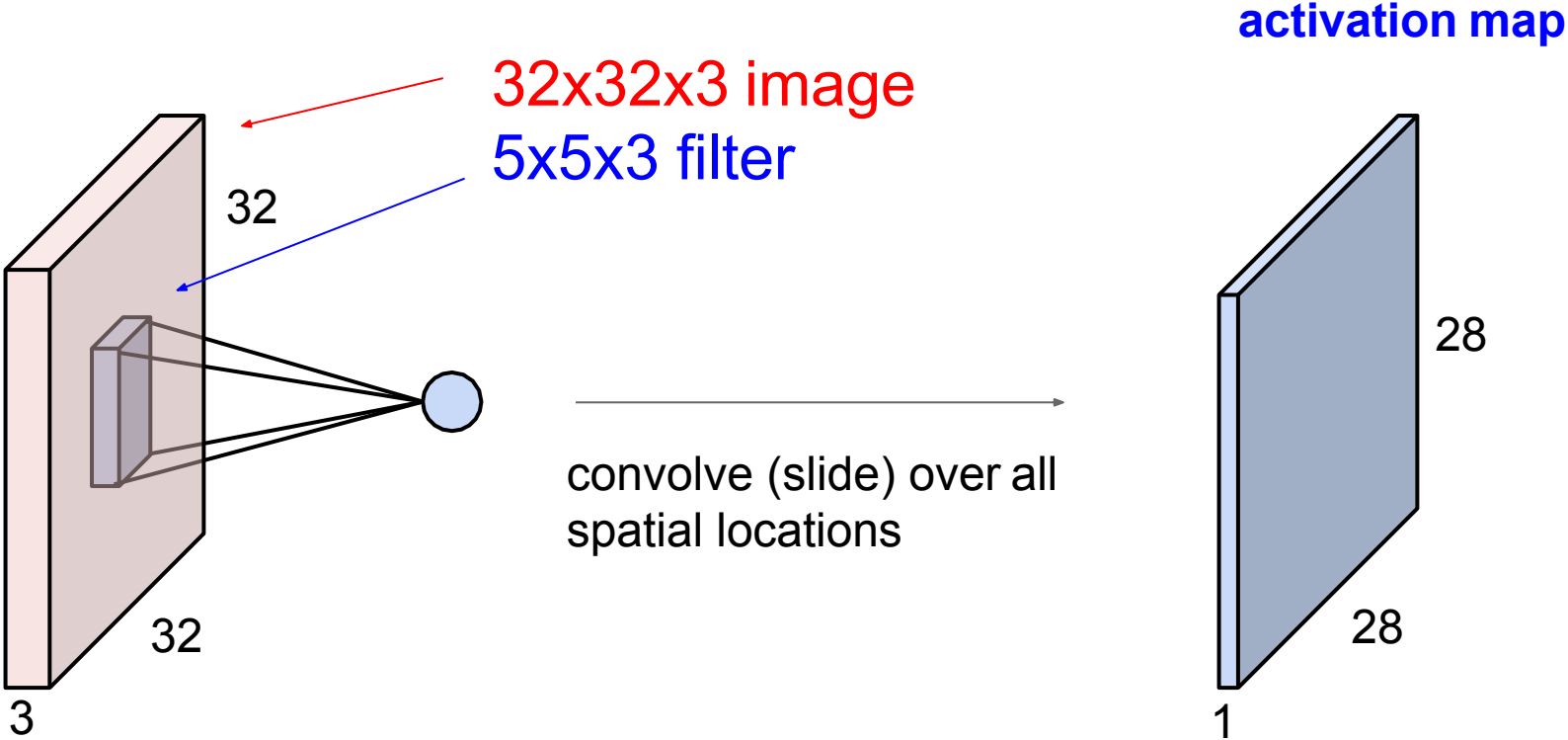
Convolution Layer



Convolution Layer

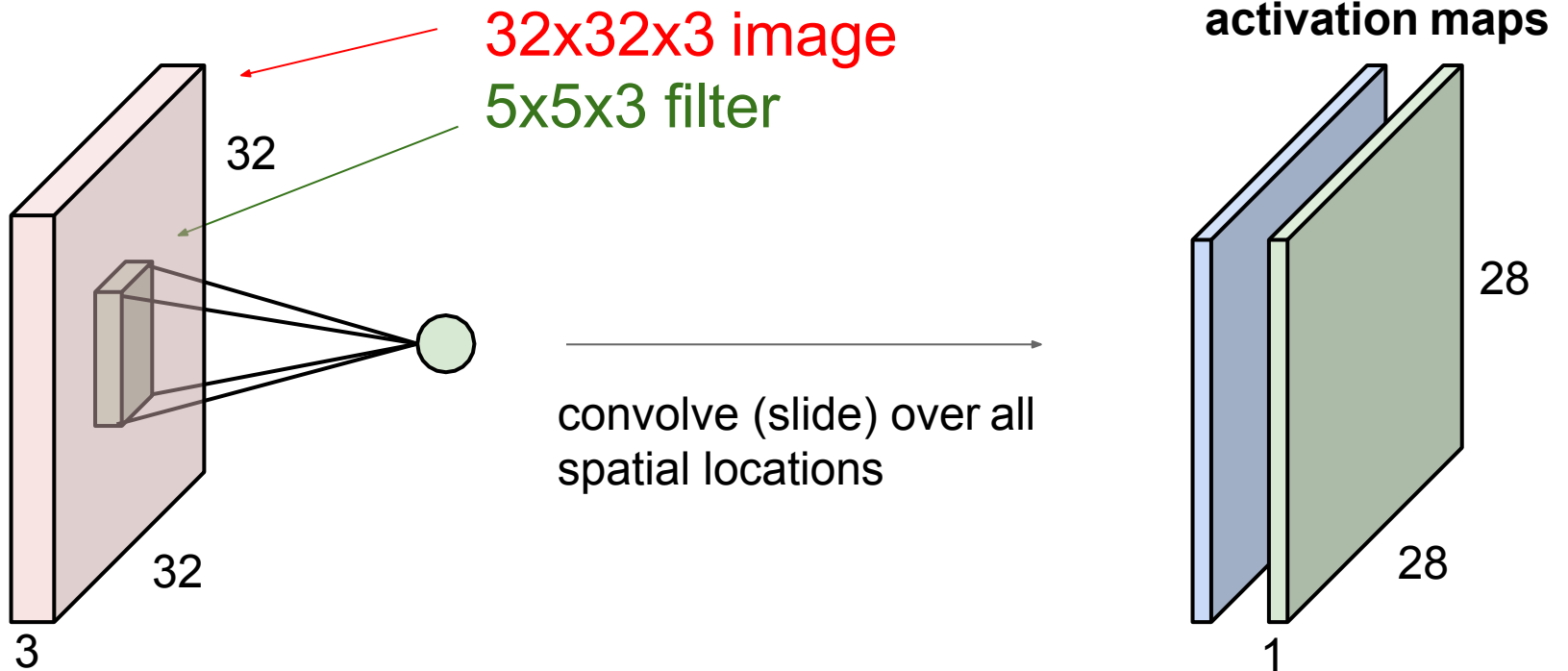


Convolution Layer

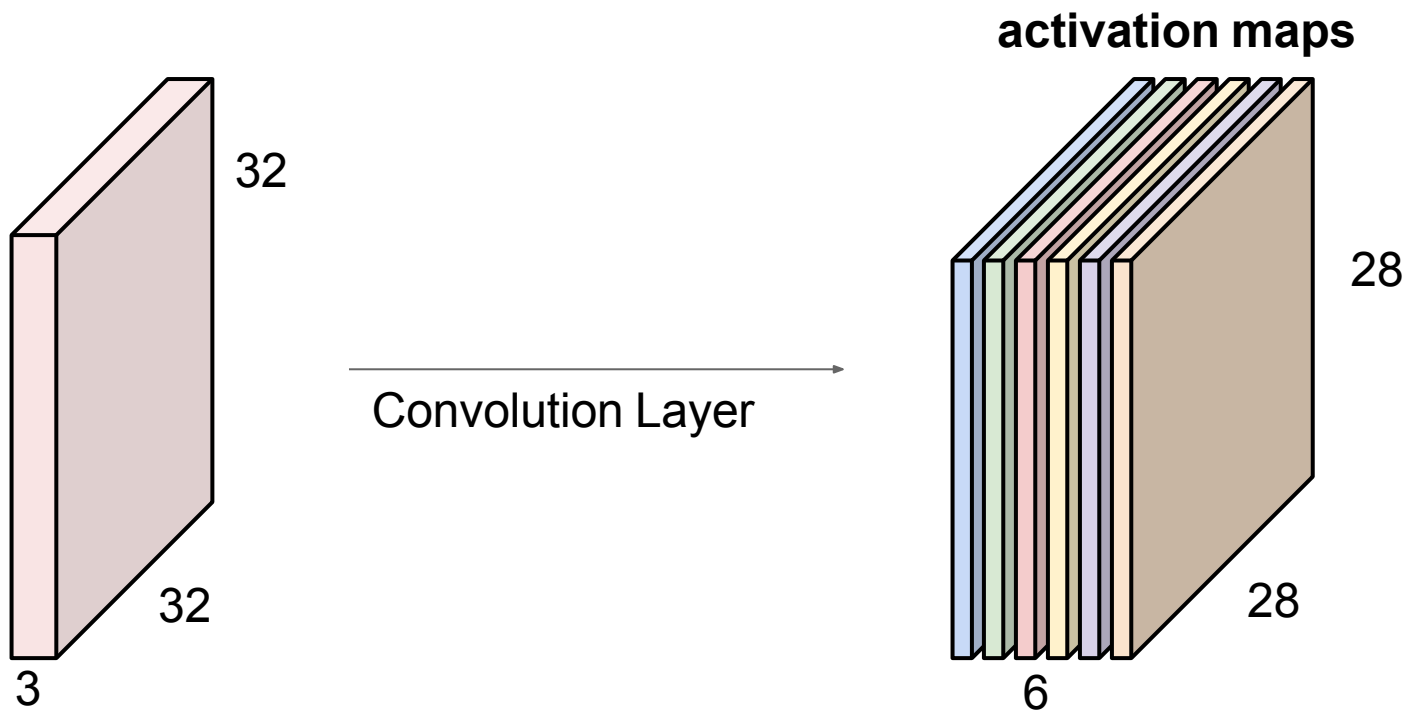


Convolution Layer

consider a second, **green** filter

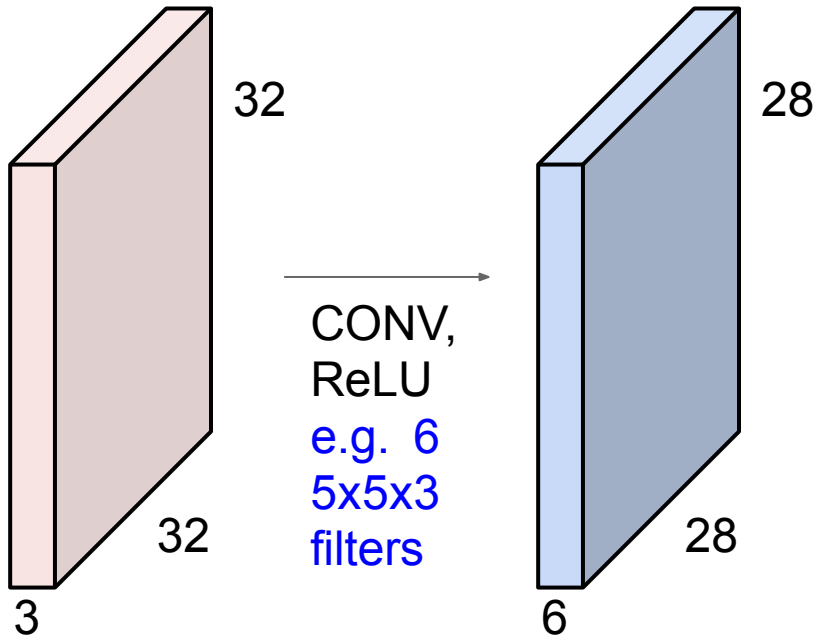


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

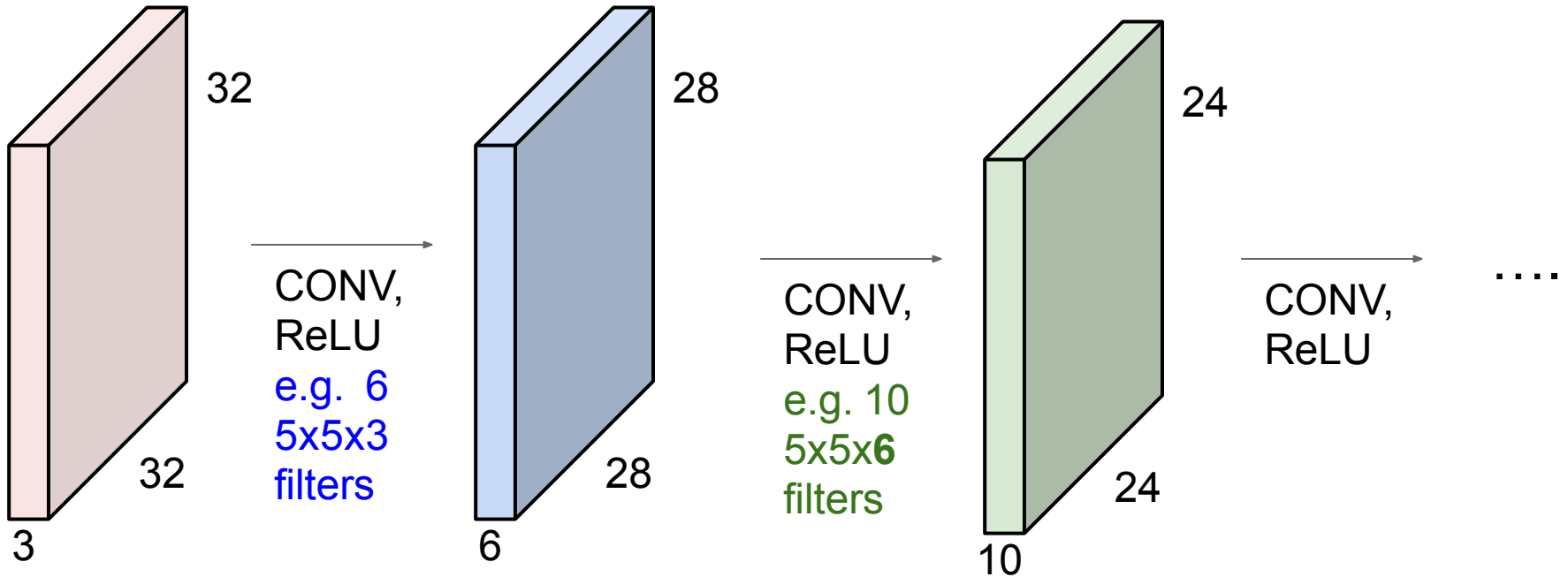


We stack these up to get a “new image” of size 28x28x6!

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



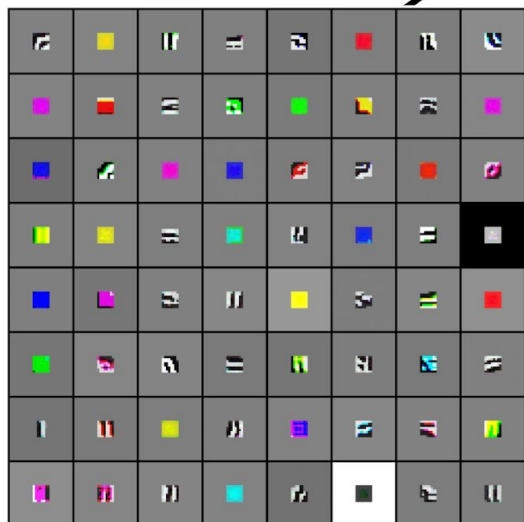
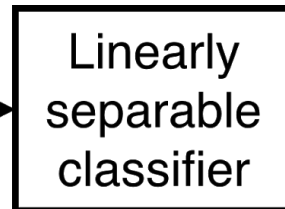
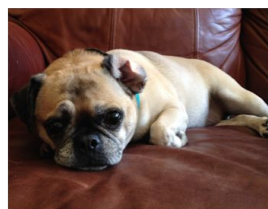
Preview: ConvNet is a sequence of Convolution Layers, interspersed with activation functions



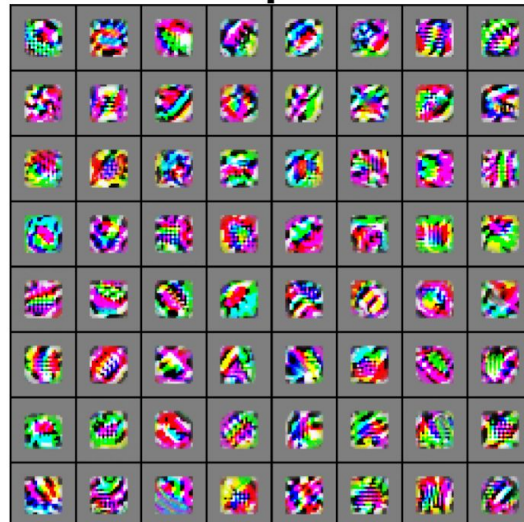
Preview

[Zeiler and Fergus 2013]

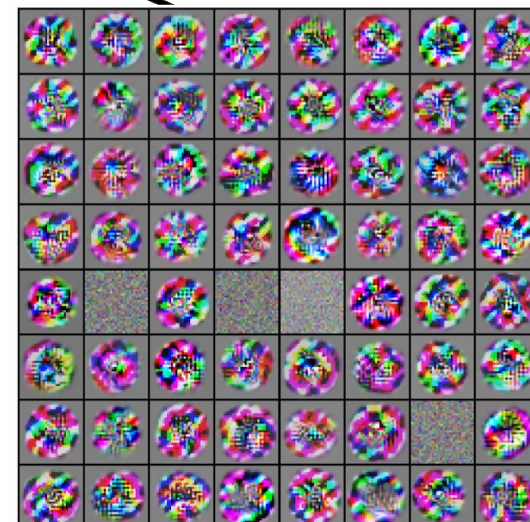
Visualization of VGG-16 by Lane McIntosh. VGG-16 architecture from [Simonyan and Zisserman 2014].



VGG-16 Conv1_1

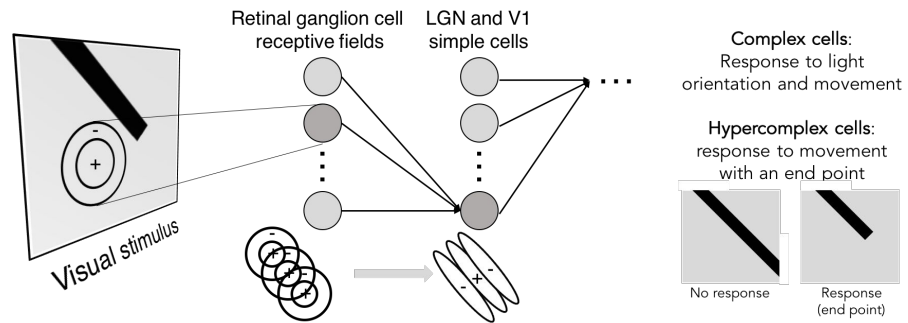
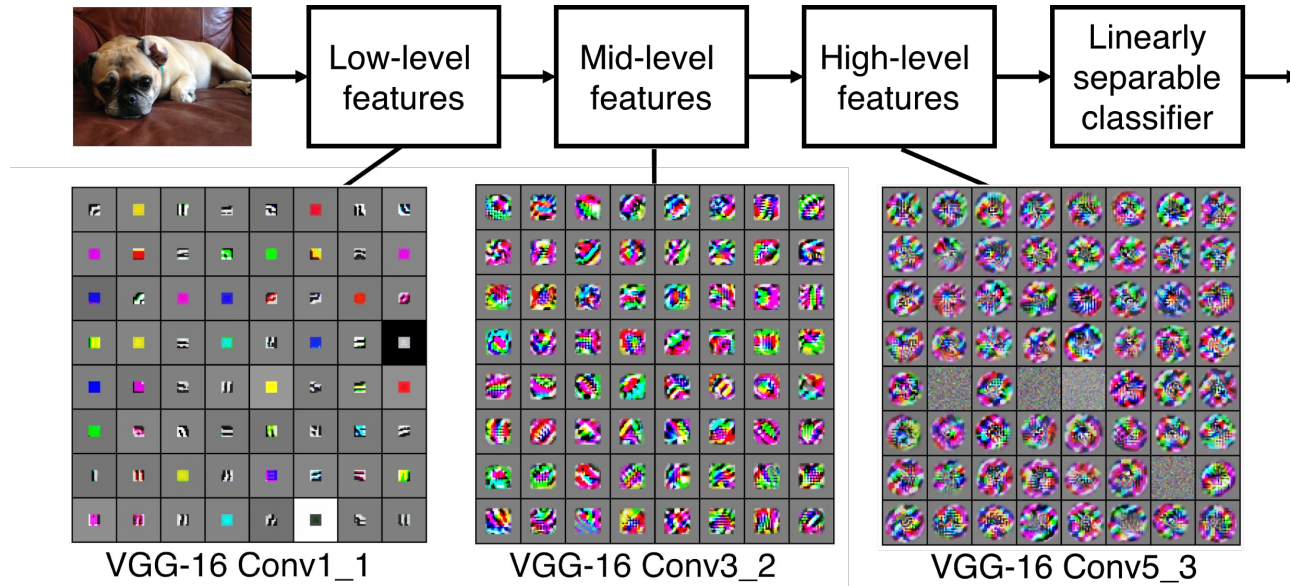


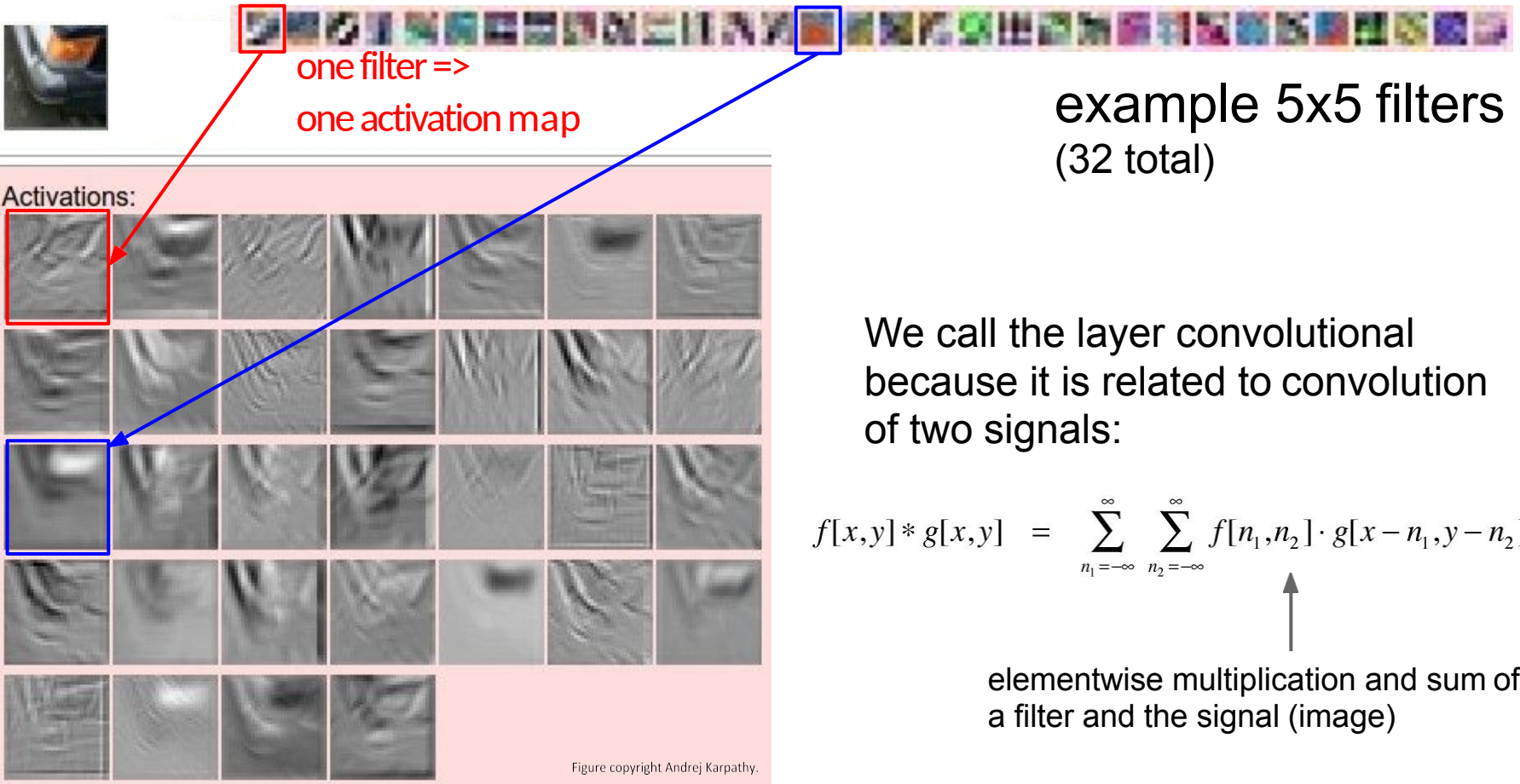
VGG-16 Conv3_2



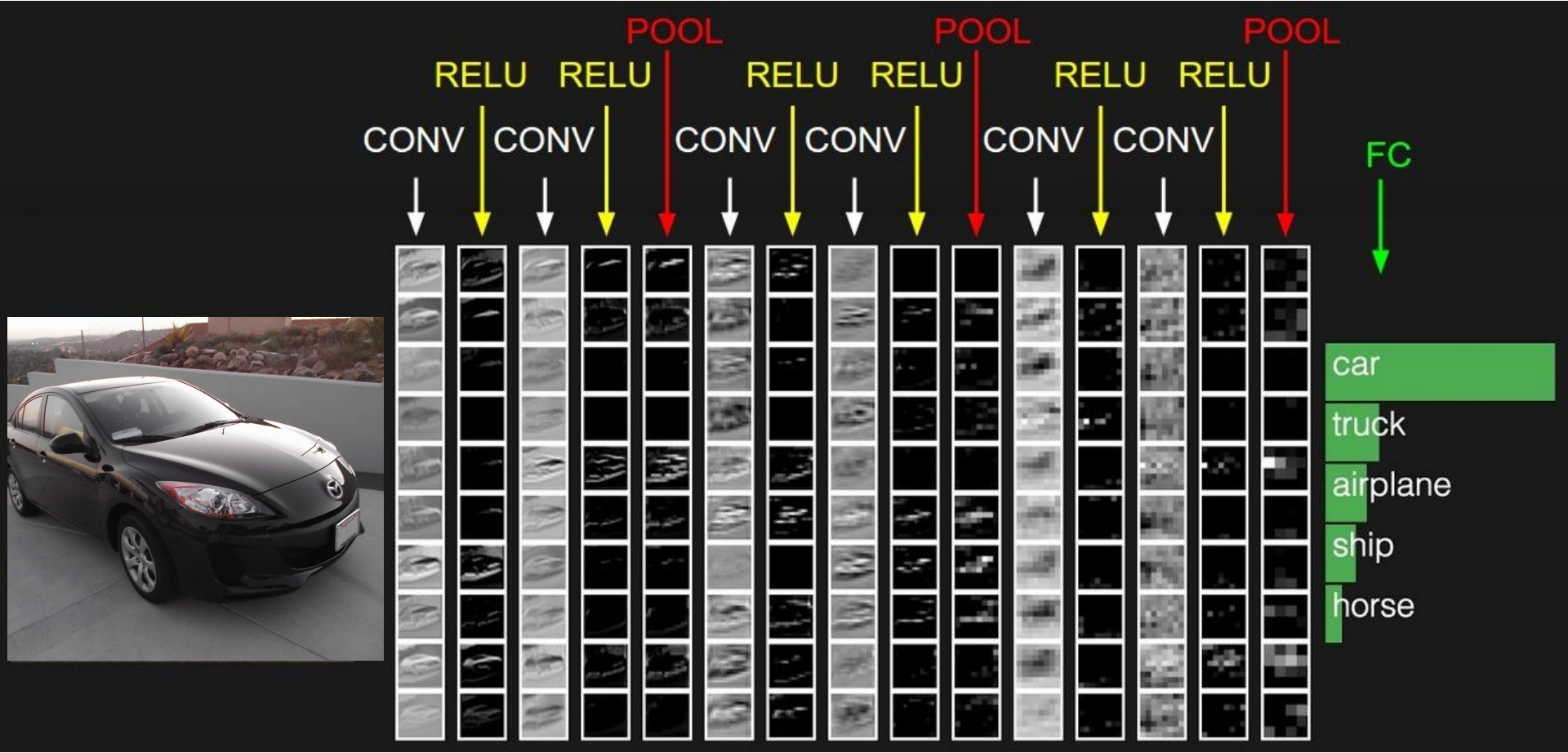
VGG-16 Conv5_3

Preview

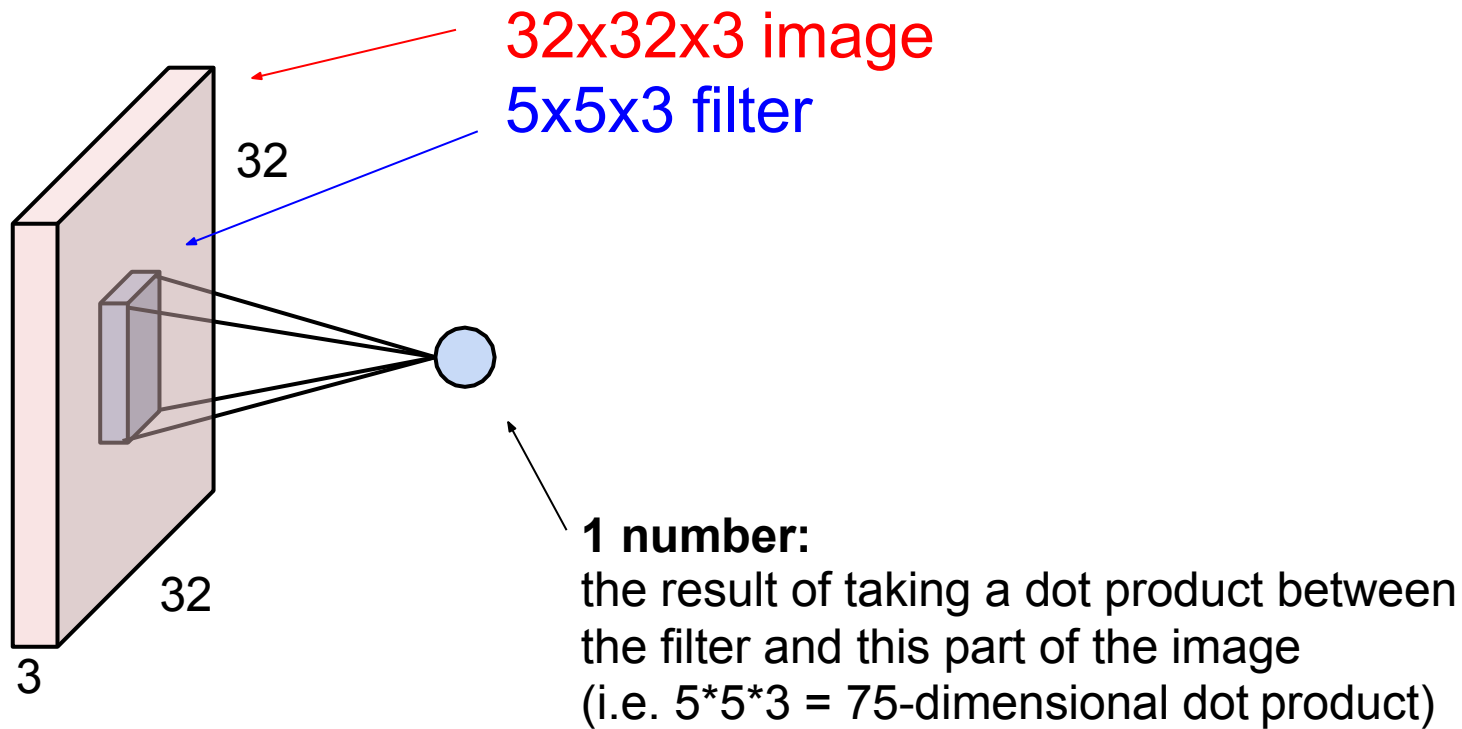




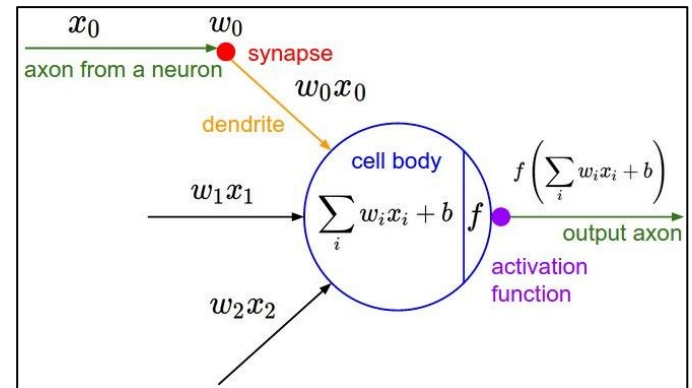
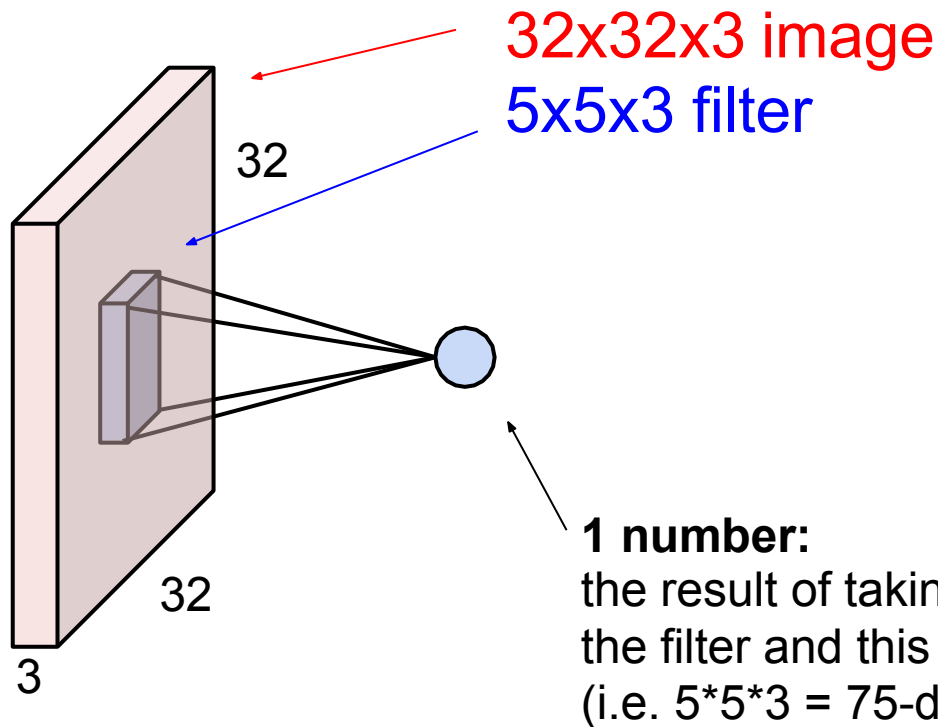
Preview:



The brain/neuron view of CONV Layer

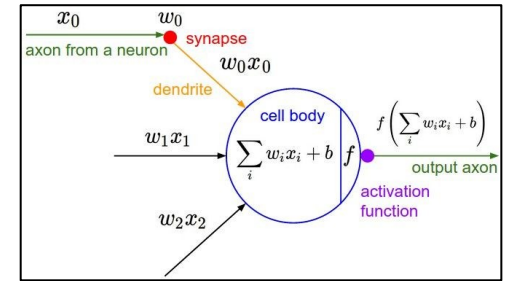
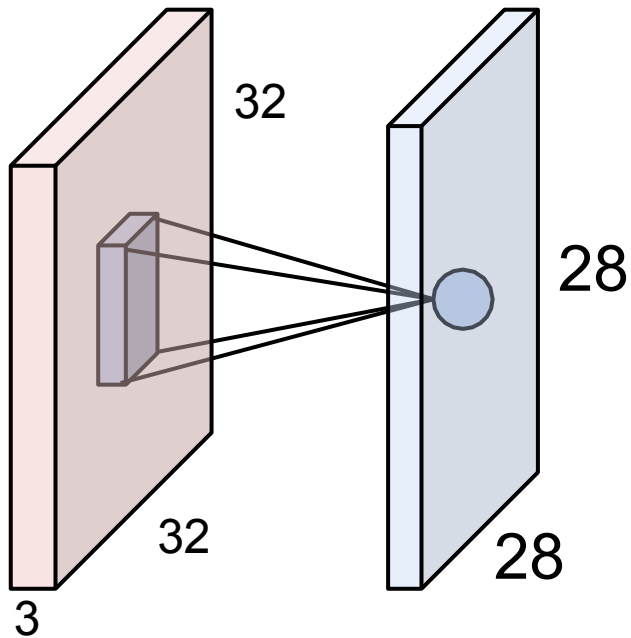


The brain/neuron view of CONV Layer



It's just a neuron with local connectivity...

The brain/neuron view of CONV Layer

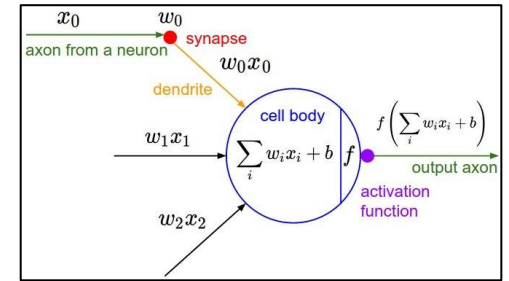
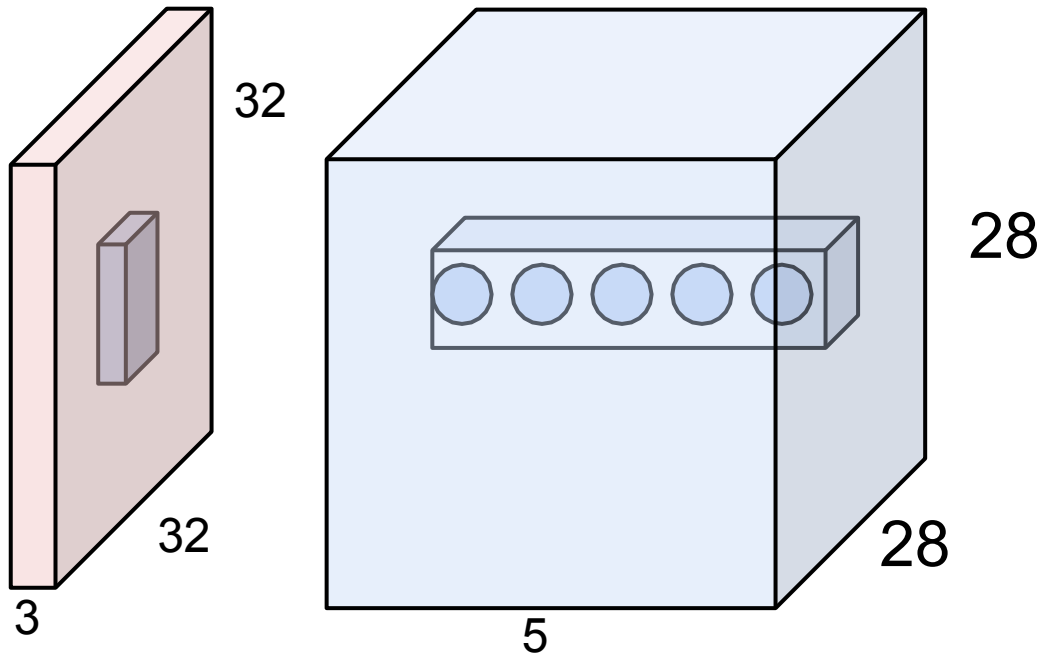


An activation map is a 28x28 sheet of neuron outputs:

1. Each is connected to a small region in the input
2. All of them share parameters

“5x5 filter” -> “5x5 receptive field for each neuron”

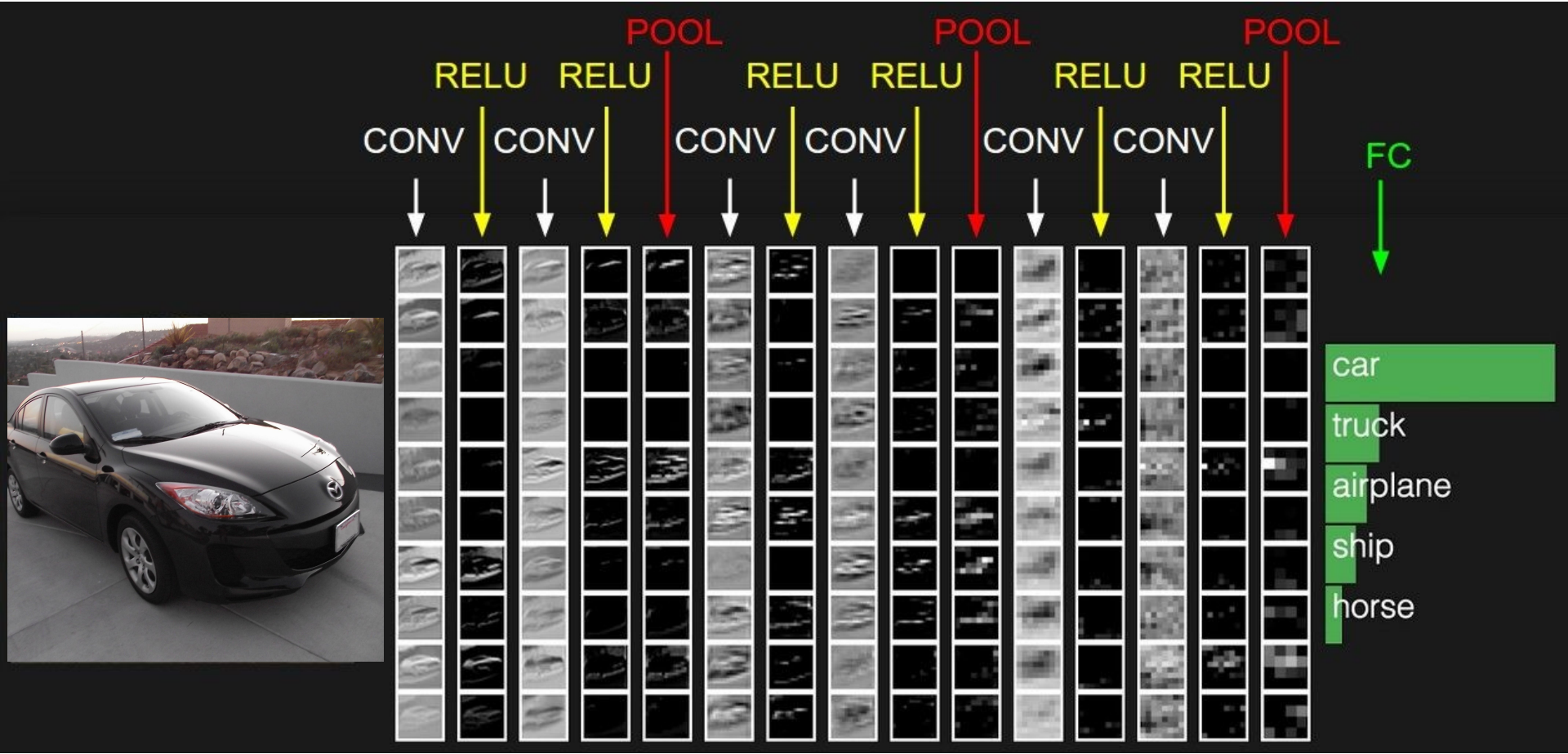
The brain/neuron view of CONV Layer



E.g. with 5 filters, CONV layer consists of neurons arranged in a 3D grid (28x28x5)

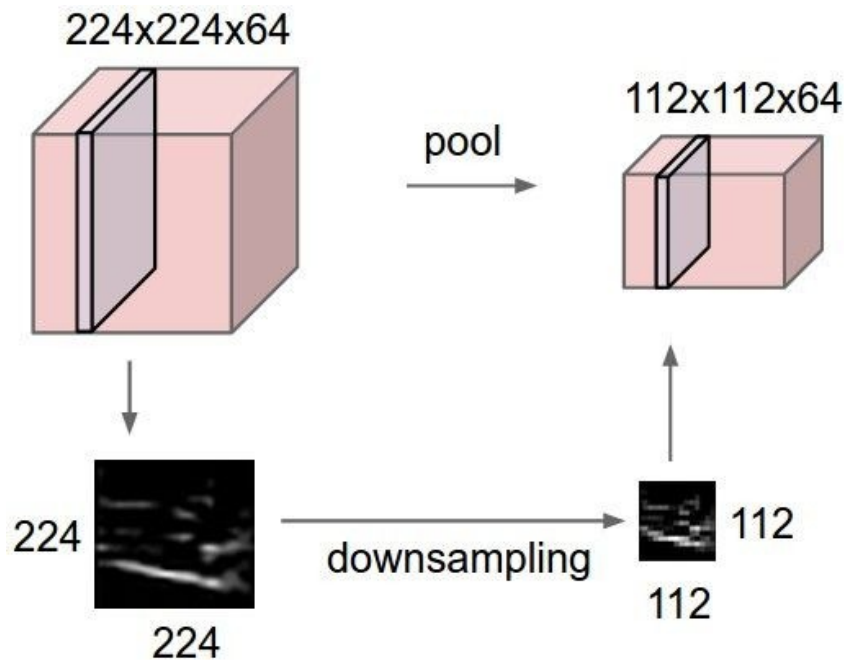
There will be 5 different neurons all looking at the same region in the input volume

two more layers to go: POOL/FC



Pooling layer

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



Max Pooling

Single depth slice

x ↑

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

→ y

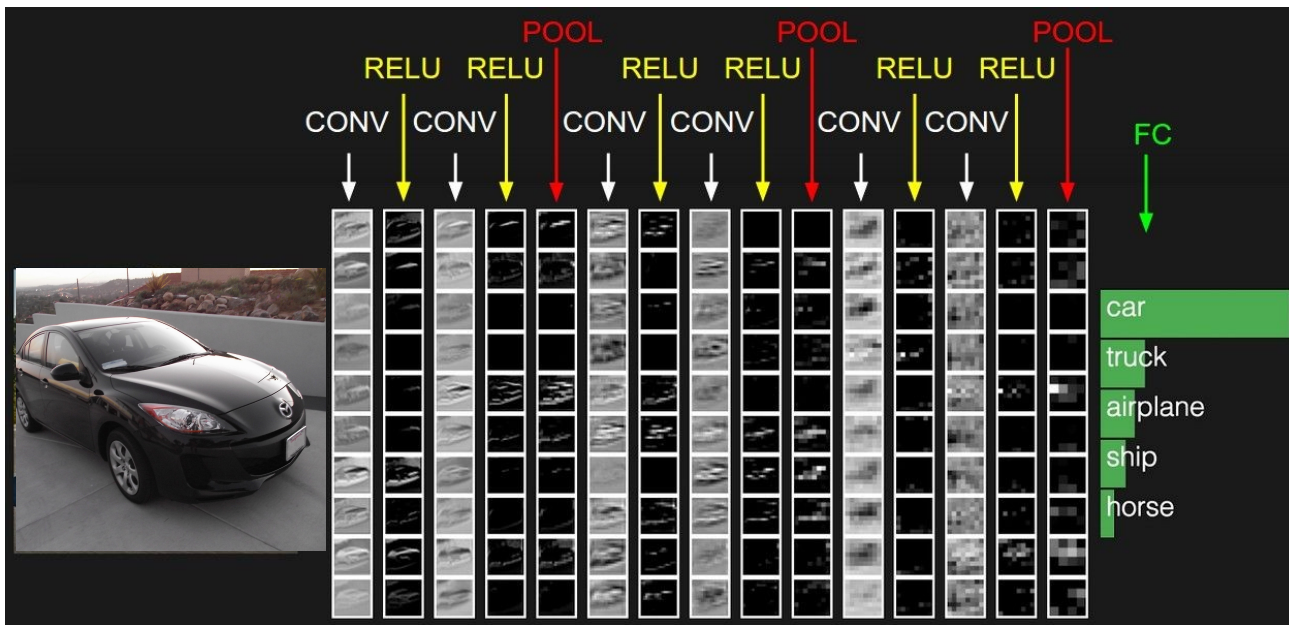
max pool with 2x2 filters
and stride 2



6	8
3	4

Fully Connected Layer (FC layer)

- Contains neurons that connect to the entire input volume, as in ordinary Neural Networks



Summary

- ConvNets stack CONV, POOL, FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Typical architectures look like
[(CONV-RELU)*N-POOL?]*M-(FC-RELU)*K, SOFTMAX
where N is usually up to ~5, M is large, $0 \leq K \leq 2$.
 - but recent advances such as ResNet/GoogLeNet challenge this paradigm