

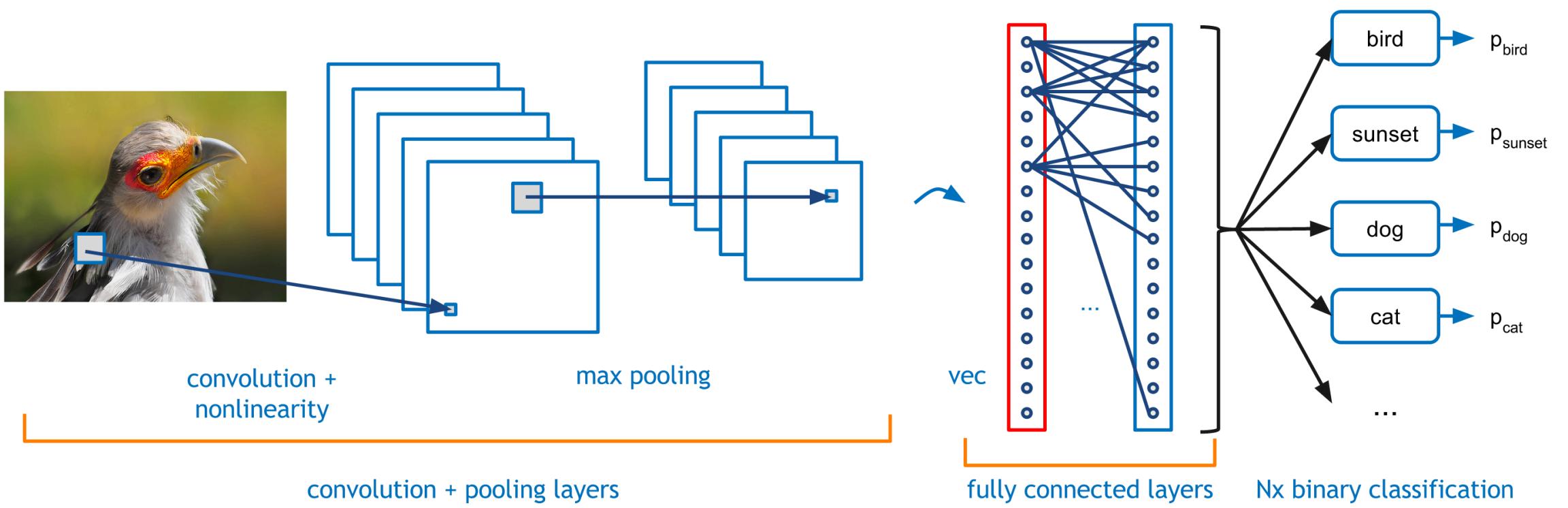
# Neural Networks MLPs and CNNs

Ardavan Pedram

It can be hard to explain  
the difference between  
the easy and the virtually  
impossible.



# High Level view

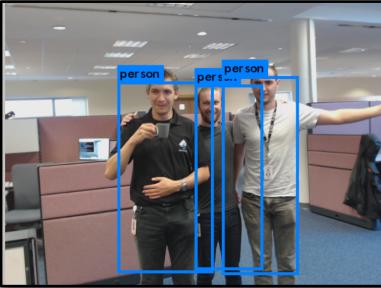


# Why Neural Networks

- **Conventional Computers**
  - Deductive Reasoning: Analyse the problem in terms of known rules
  - Code the analysis
  - Good for well defined problems
- **Neural Networks Learn From Examples**
  - Inductive Reasoning: Construct the rules
  - **Generalize to handle incomplete data**
  - No requirement for an explicit description of the problem
  - **Noise Resilient**



Image Classification



Object Detection



Semantic Segmentation

Computer Vision  
CNNs



Speaker  
Diarization



Speech  
Recognition

Speech Recognition  
RNNs, LSTMs



Translation



Sentiment Analysis

Natural Language Processing  
Sequence to sequence



Recommender

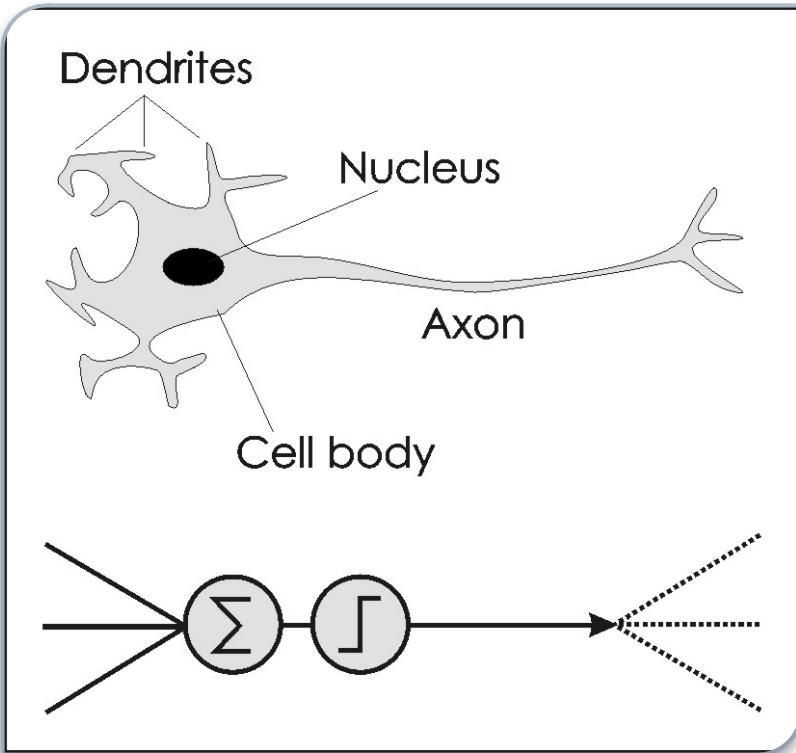


GamePlay

Many more emerging...

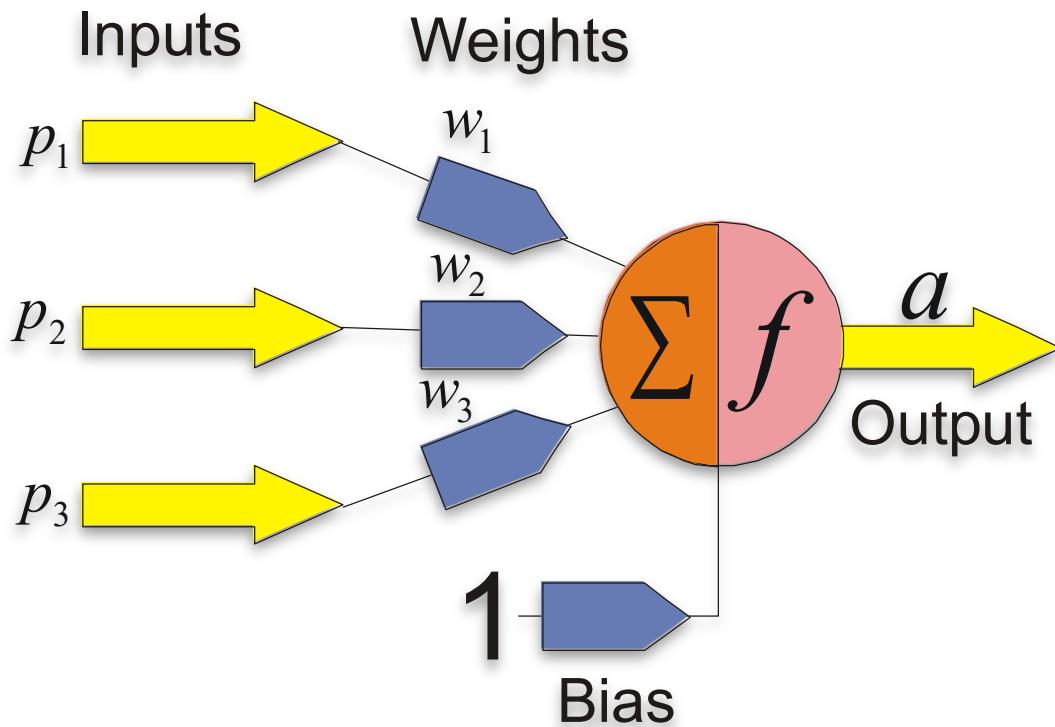
Others

# The Neuron



- A neuron's dendritic tree is connected to a thousand neighbouring neurons. When one of those neurons fire, a positive or negative charge is received by one of the dendrites. The strengths of all the received charges are added together through the processes of spatial and temporal summation

# Anatomy of Neuron

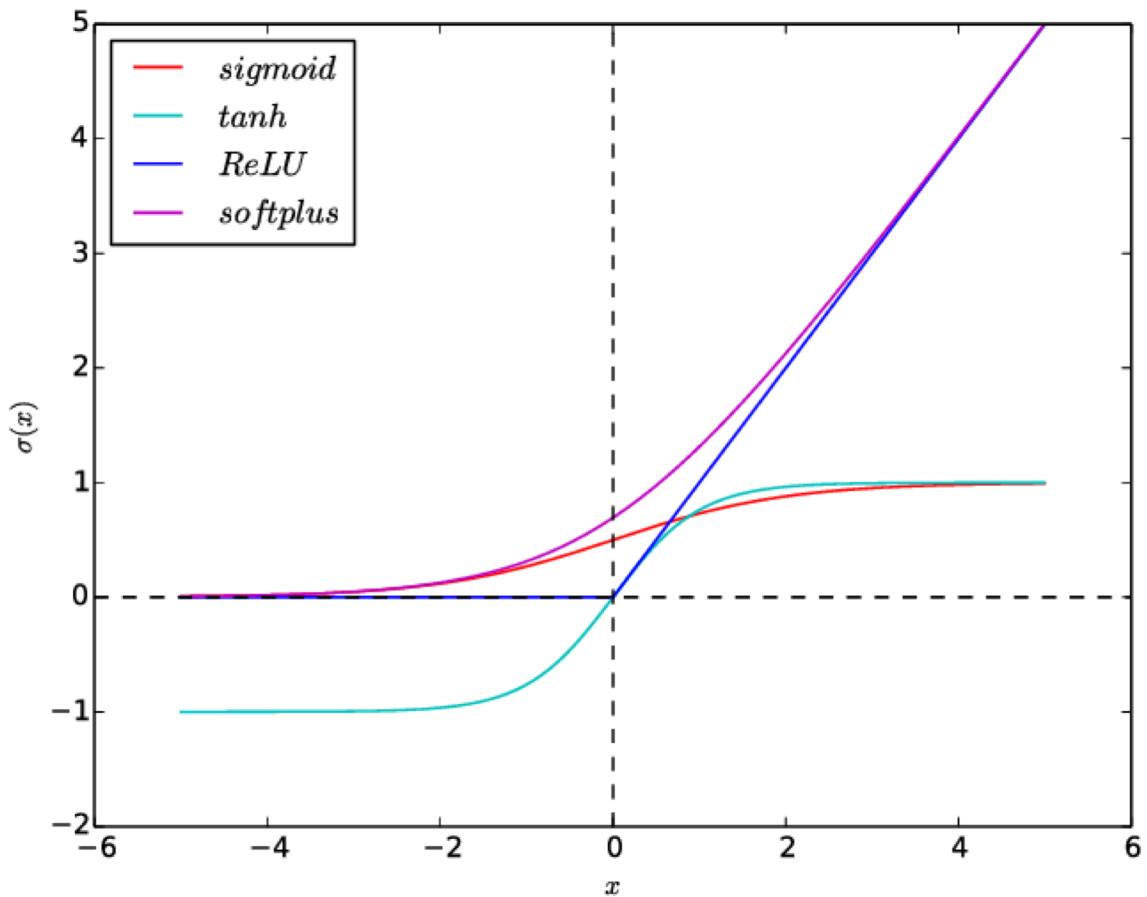


$$a = f(p_1 w_1 + p_2 w_2 + p_3 w_3 + b) = f\left(\sum p_i w_i + b\right)$$

- Basic Computation is Multiply Add
  - Sum of Inputs x weights
- Activation function can vary
  - Differentiable?

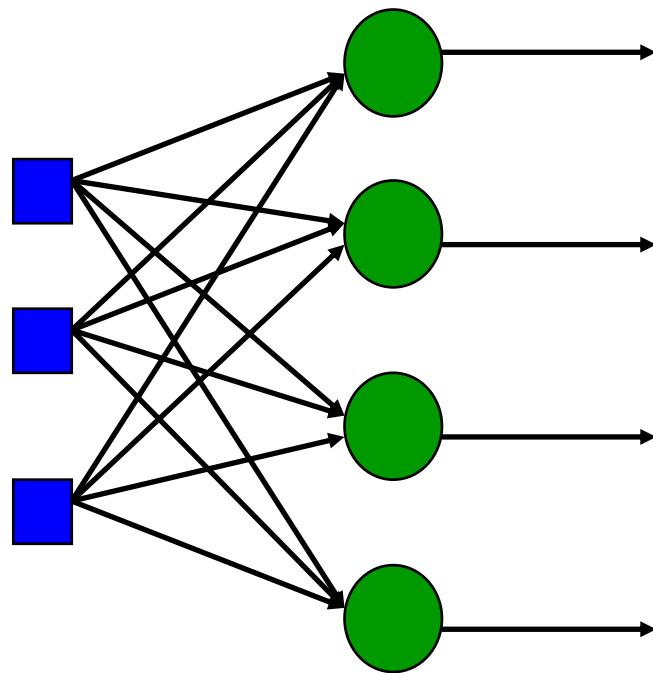
# What is the Activation Function $F$ ?

- Linear functions
  - limited because the output is simply proportional to the input.
- Non-linear functions
  - Sigmoid
  - Step



# Single Layer

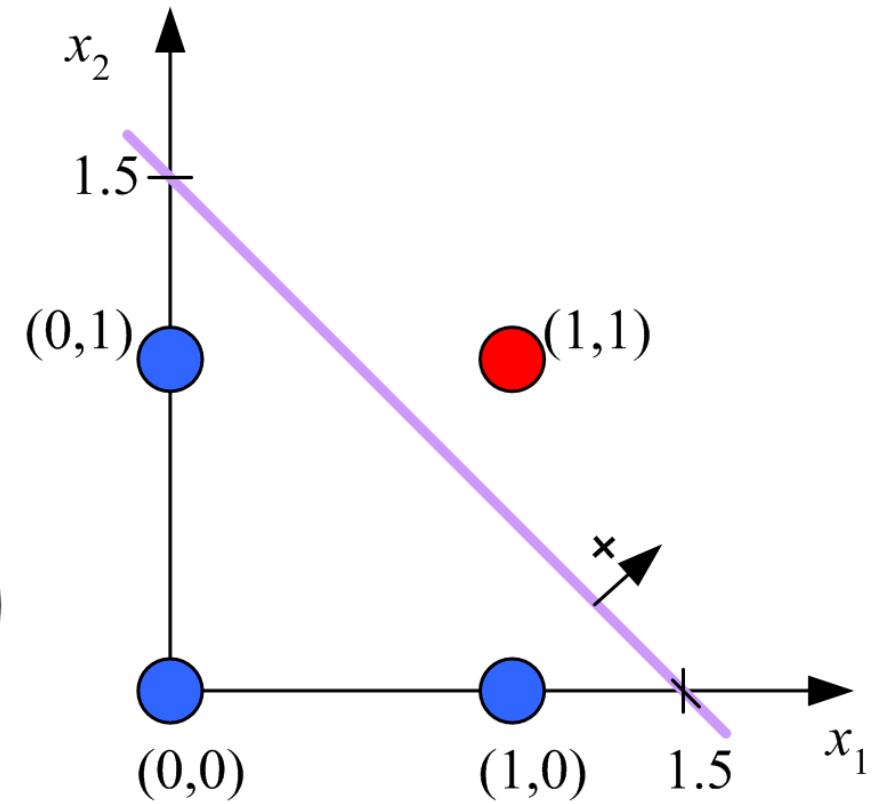
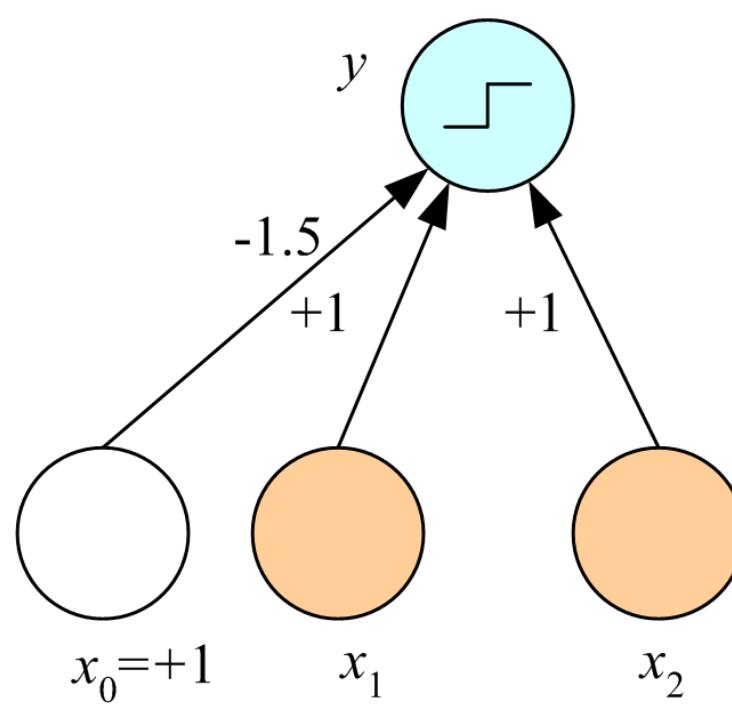
*Input layer  
of  
source nodes*



*Output layer  
of  
neurons*

# Example: Learning Boolean AND

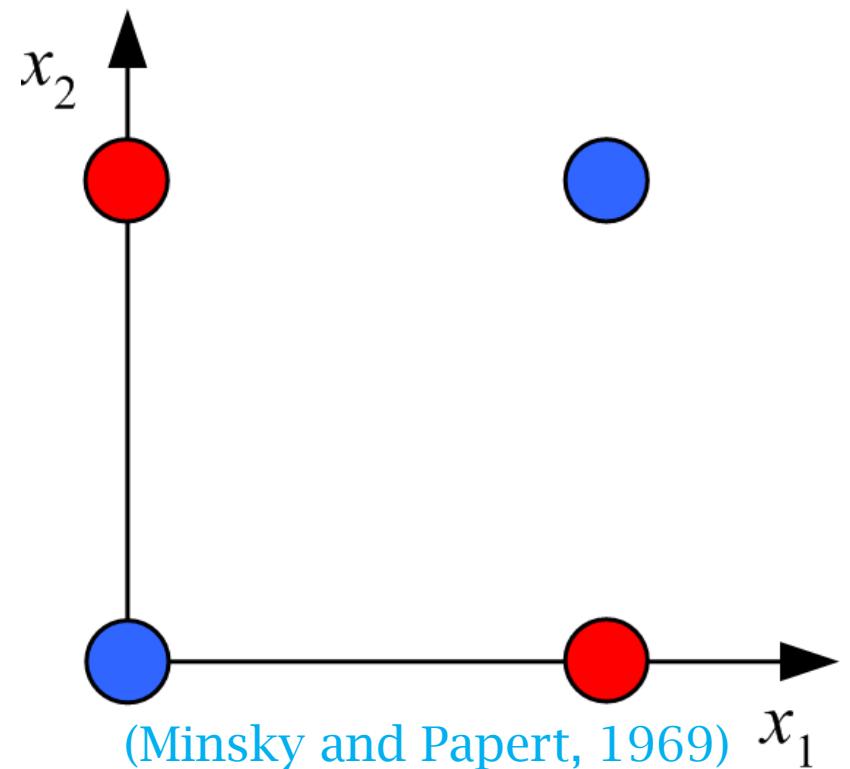
$x_1$	$x_2$	$r$
0	0	0
0	1	0
1	0	0
1	1	1



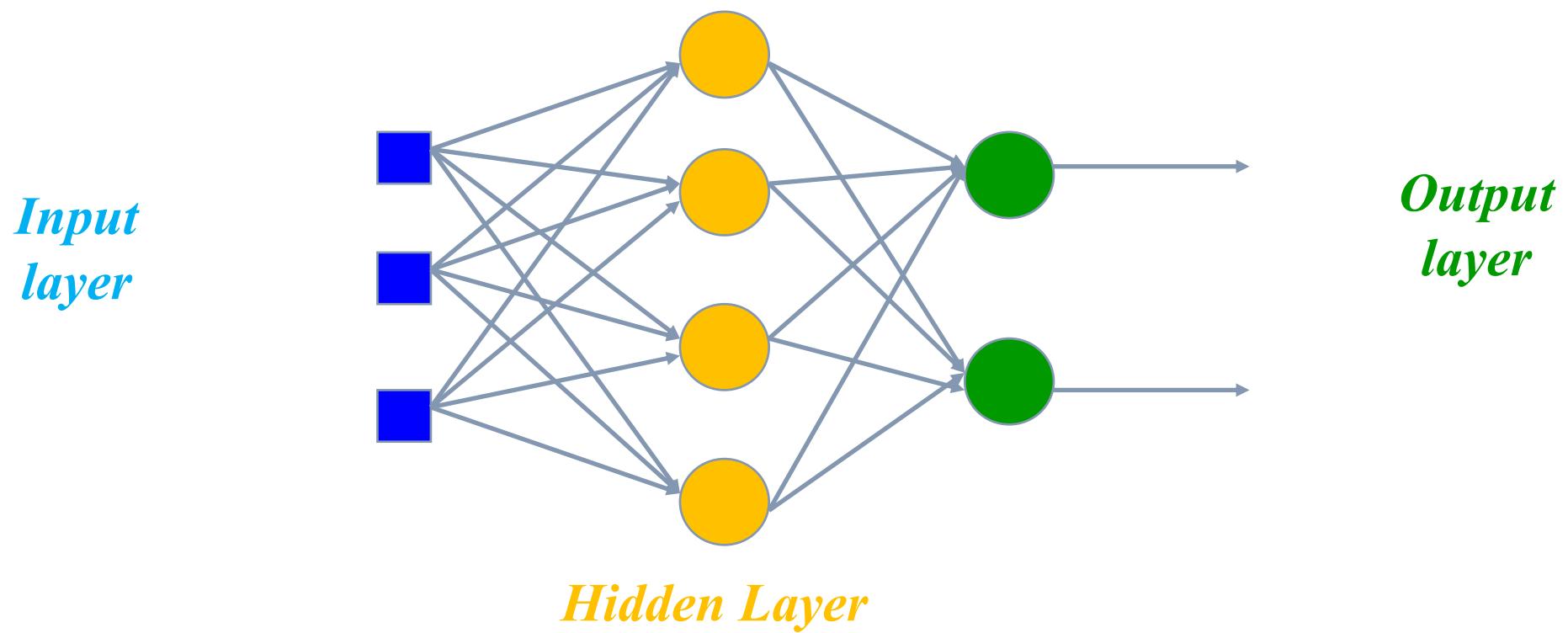
# What ABOUT XOR Function?

- Linearly non separable

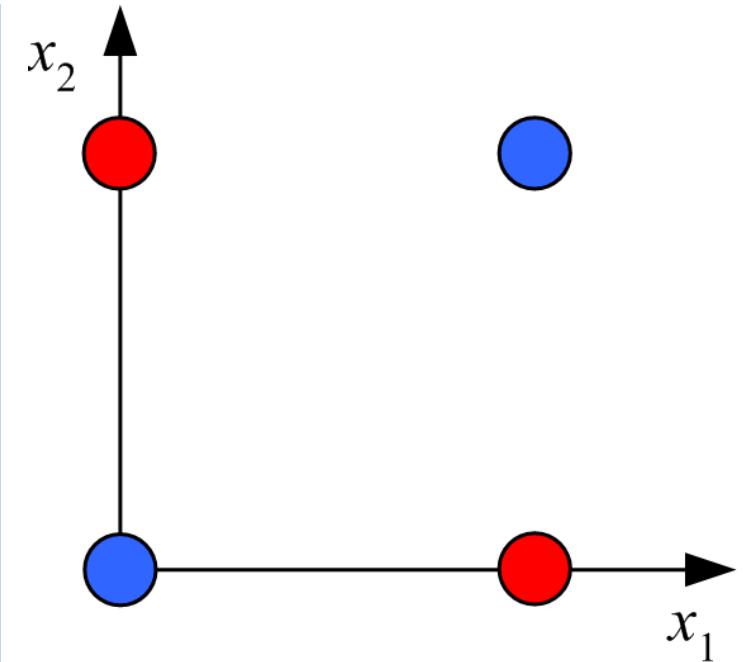
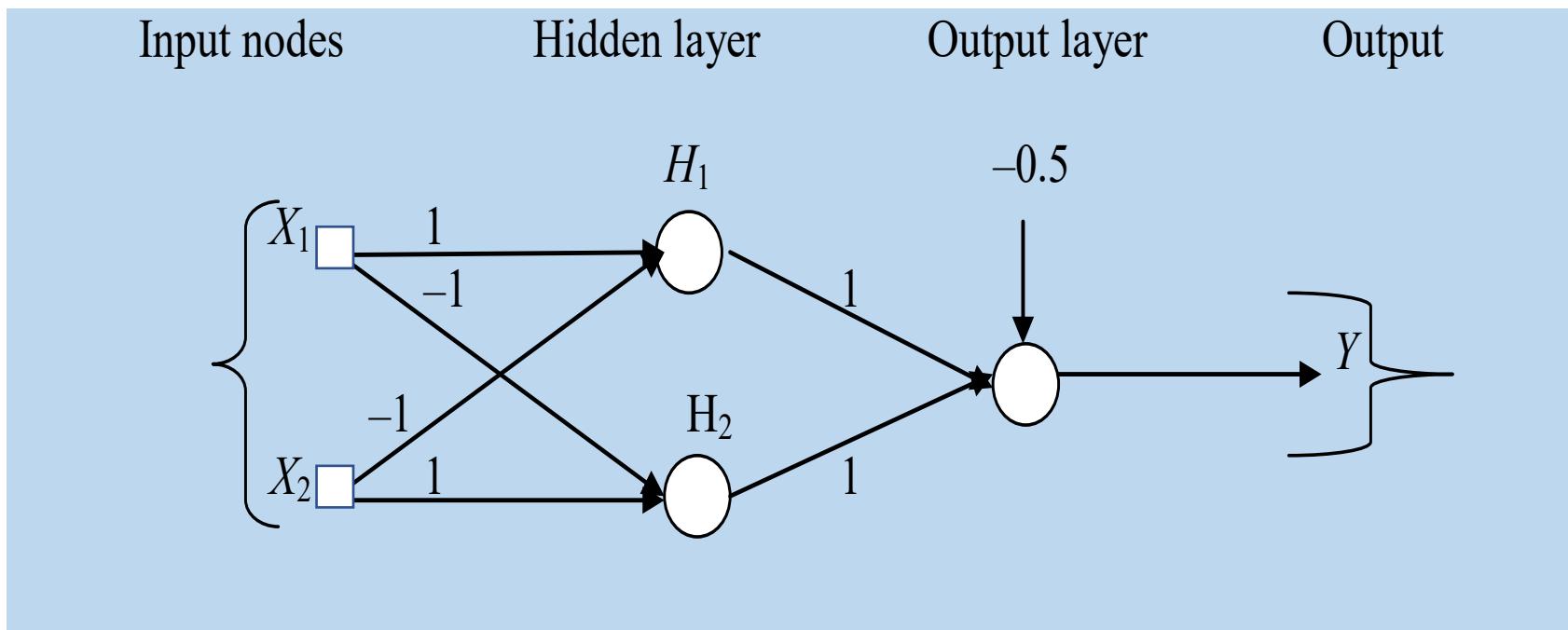
$x_1$	$x_2$	$r$
0	0	0
0	1	1
1	0	1
1	1	0



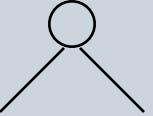
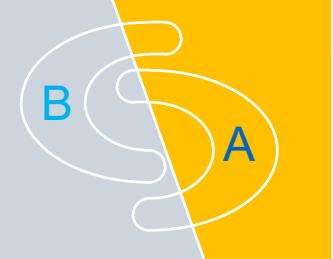
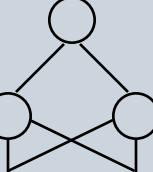
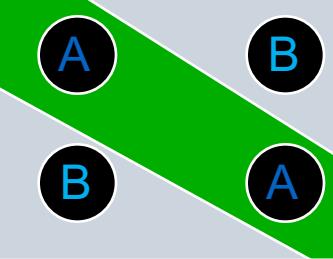
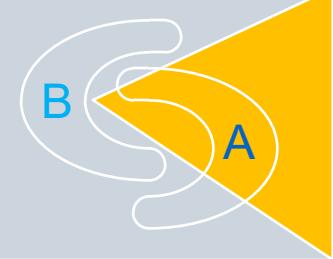
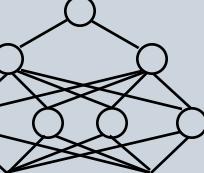
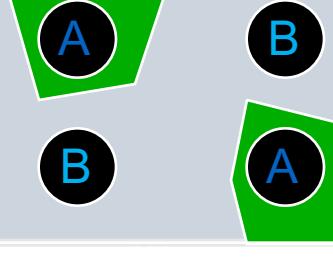
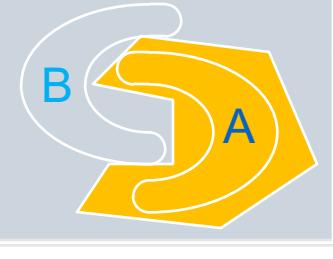
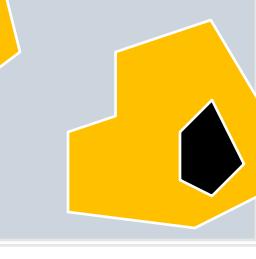
# Hidden Layer(s)



# XOR Function

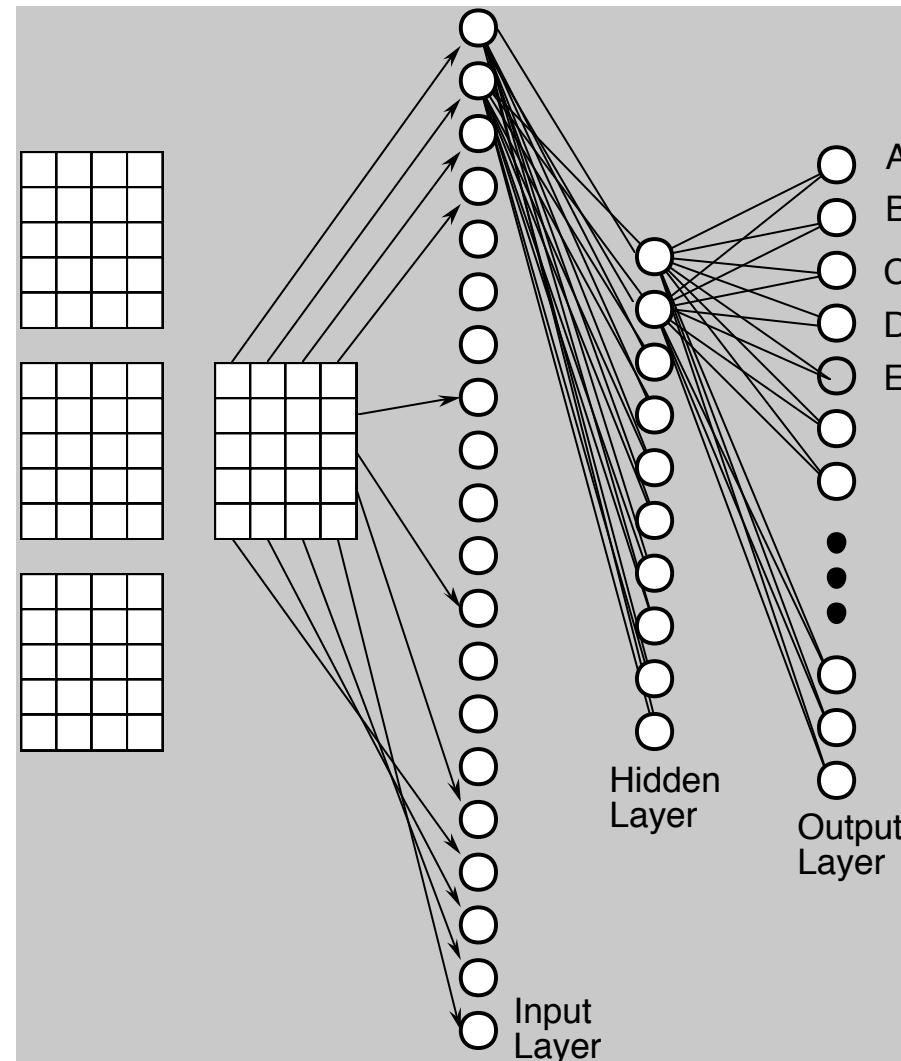


# Non-Linearly Separable Problems

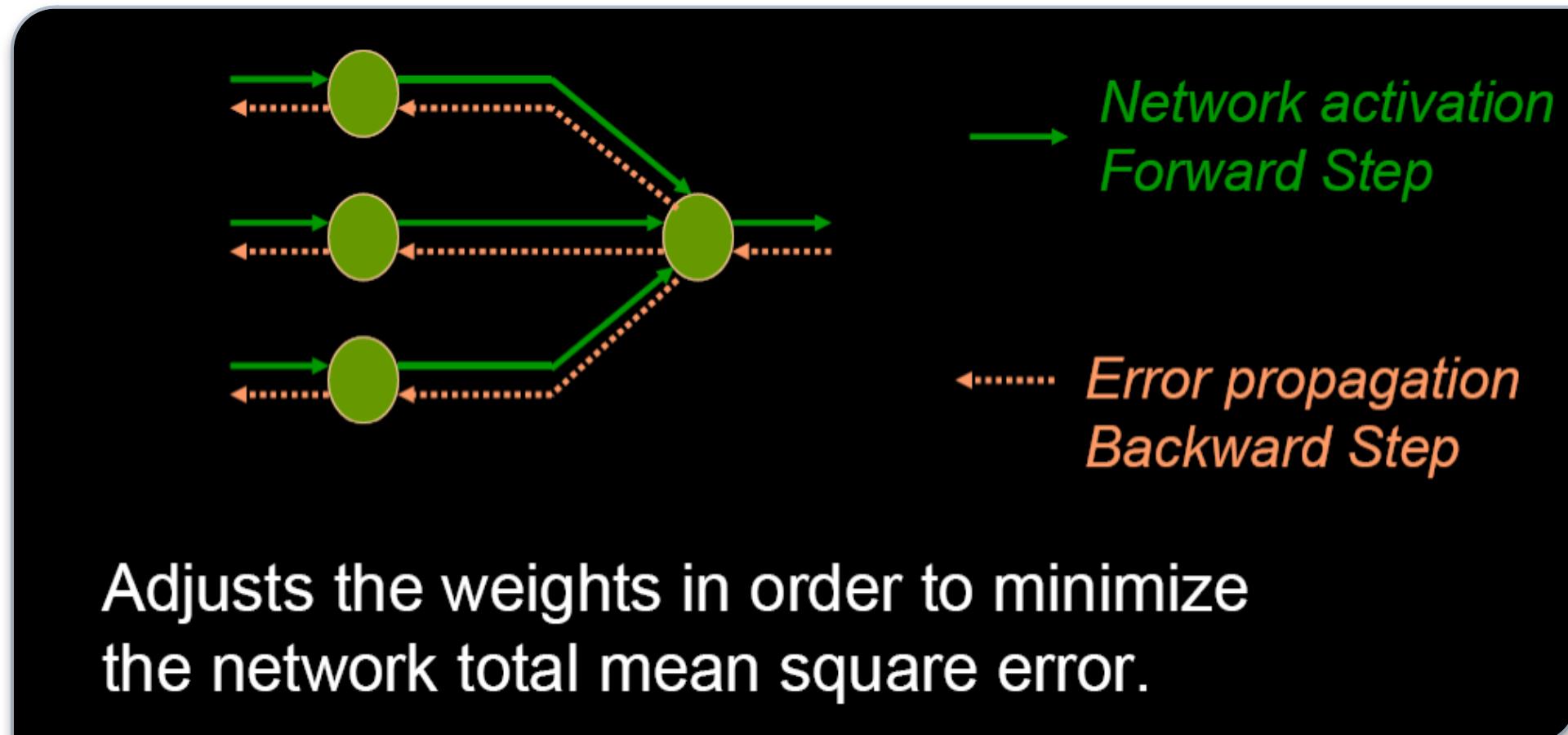
Structure	<i>Types of Decision Regions</i>	<i>Exclusive-OR Problem</i>	<i>Classes with Meshed regions</i>	<i>Most General Region Shapes</i>
<i>Single-Layer</i> 	<i>Half Plane Bounded By Hyperplane</i>			
<i>Two-Layer</i> 	<i>Convex Open Or Closed Regions</i>			
<i>Three-Layer</i> 	<i>Arbitrary (Complexity Limited by No. of Nodes)</i>			

# Neural network for OCR

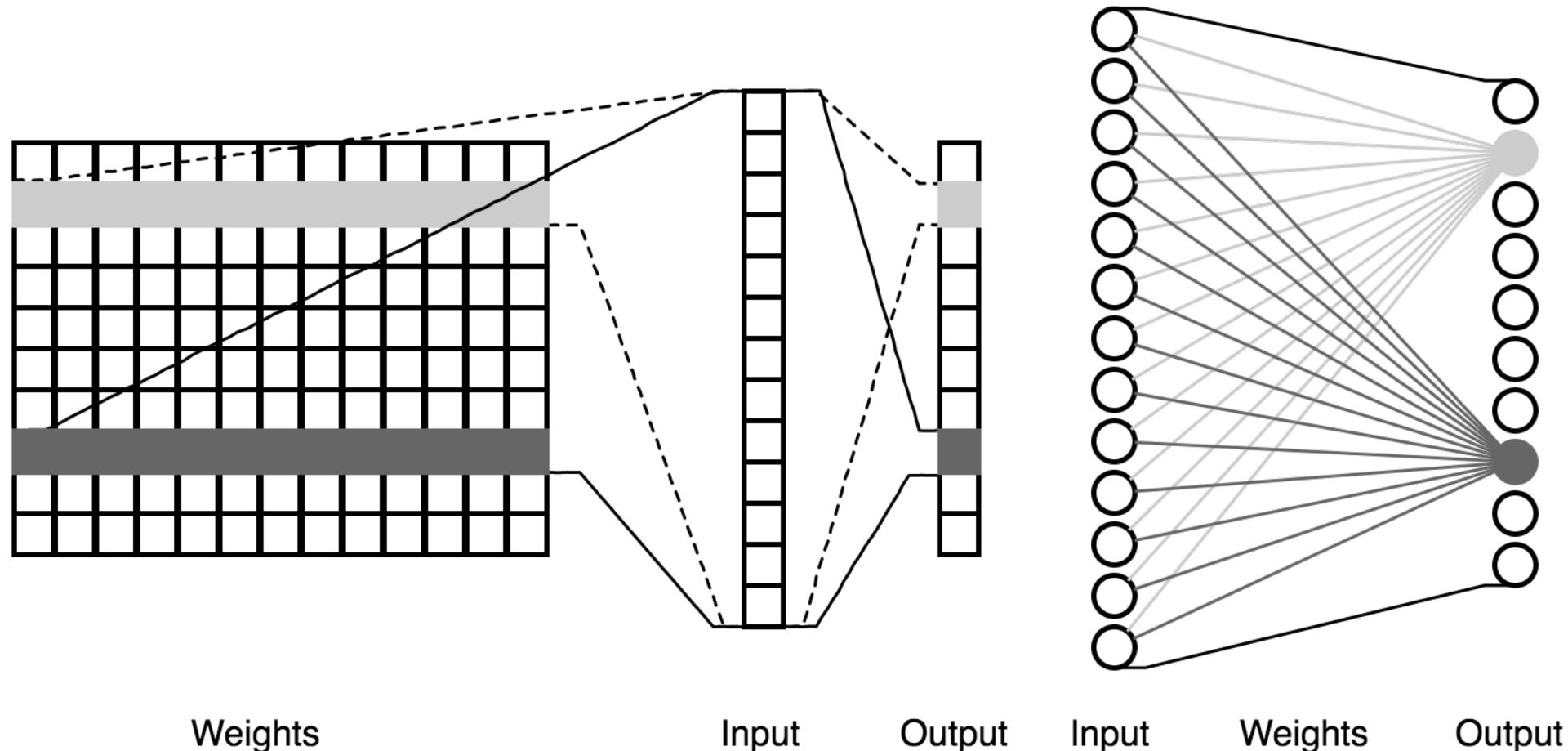
- Feedforward network
- Trained using Back-propagation



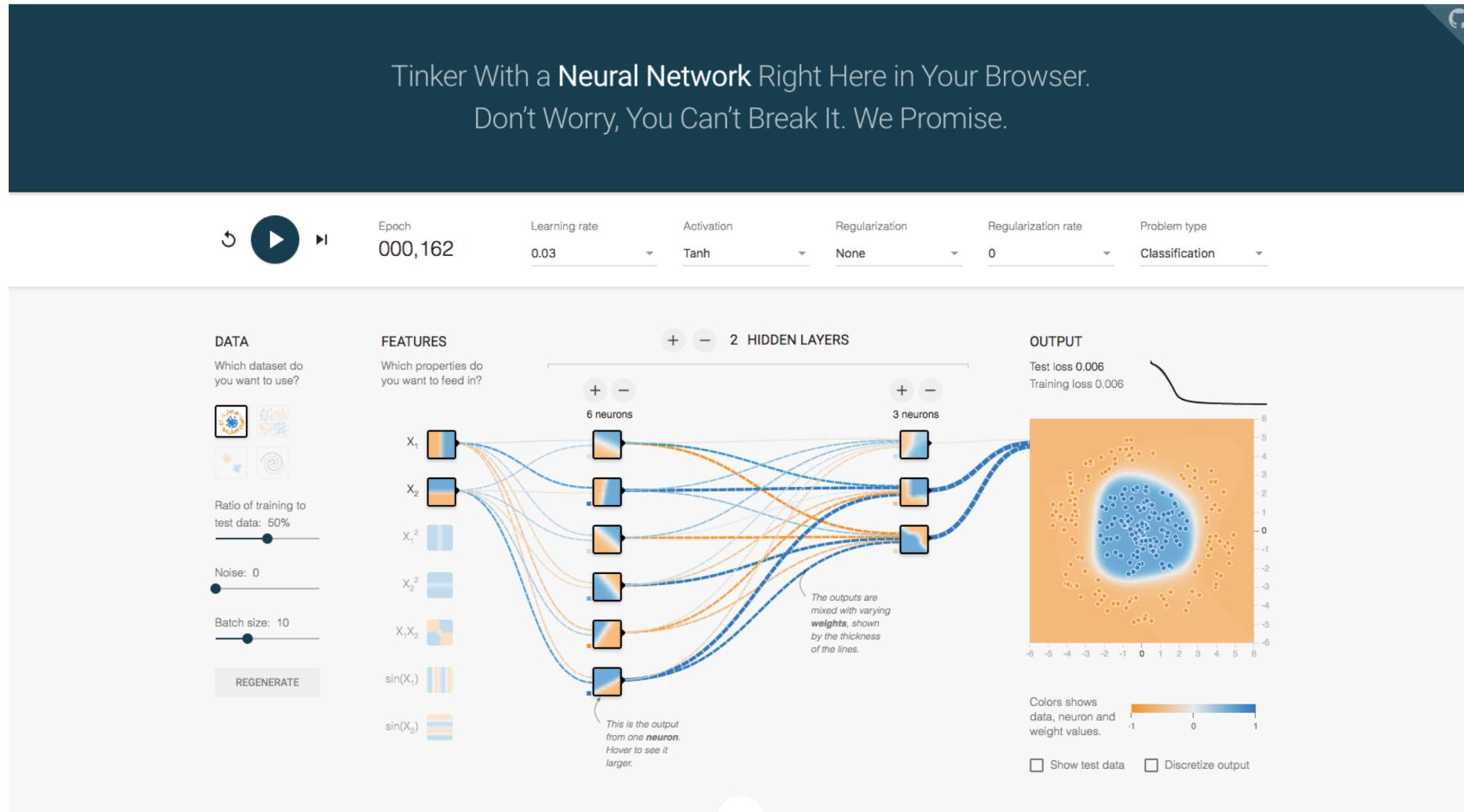
# Learning: Back Propagation

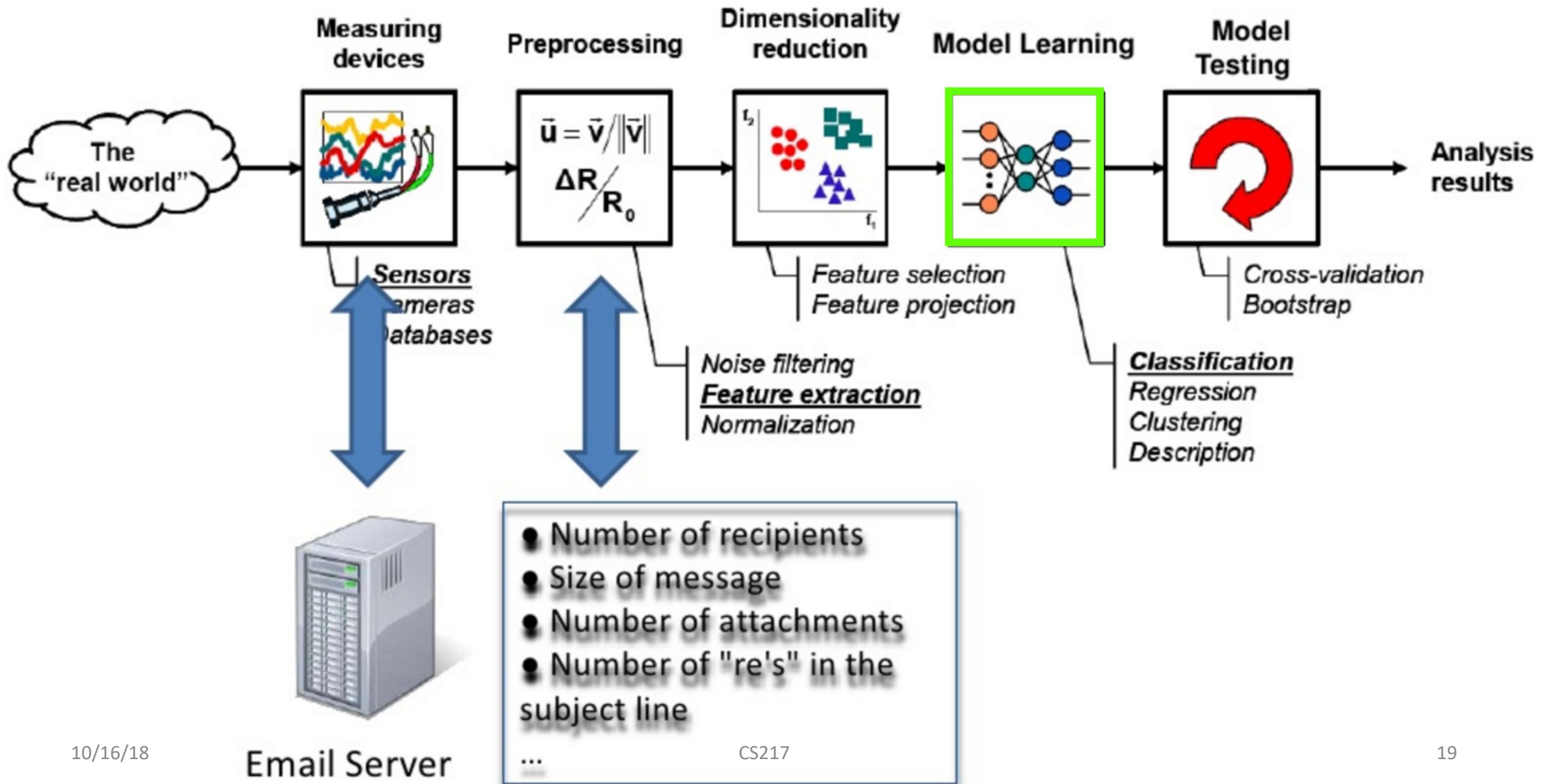


# Forward Propagation is GEMV

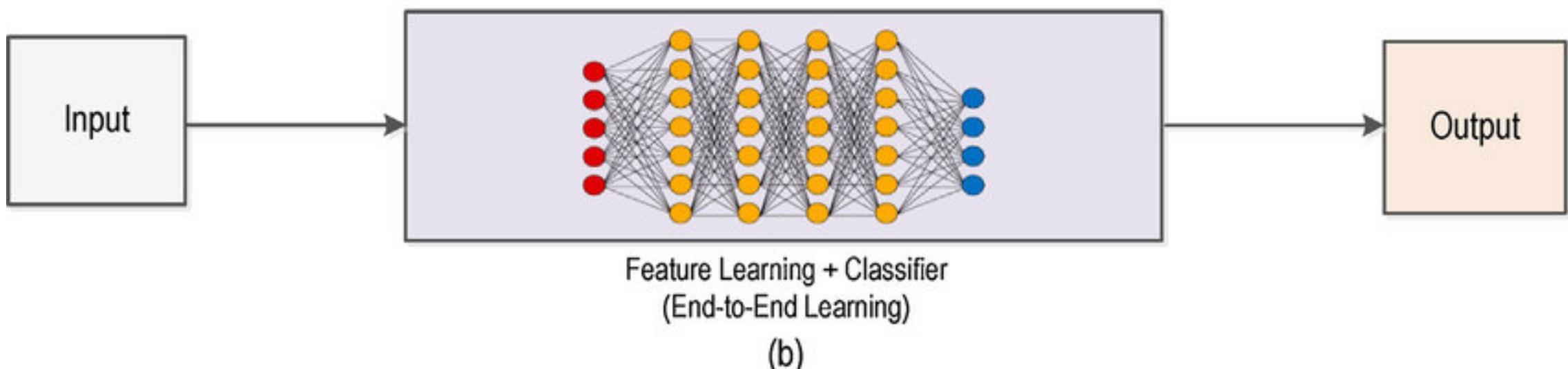


# Playground.tensorflow.com



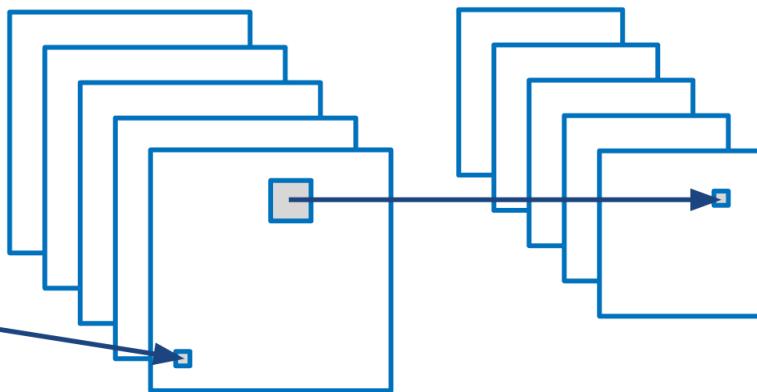


# Shallow & Deep Learning

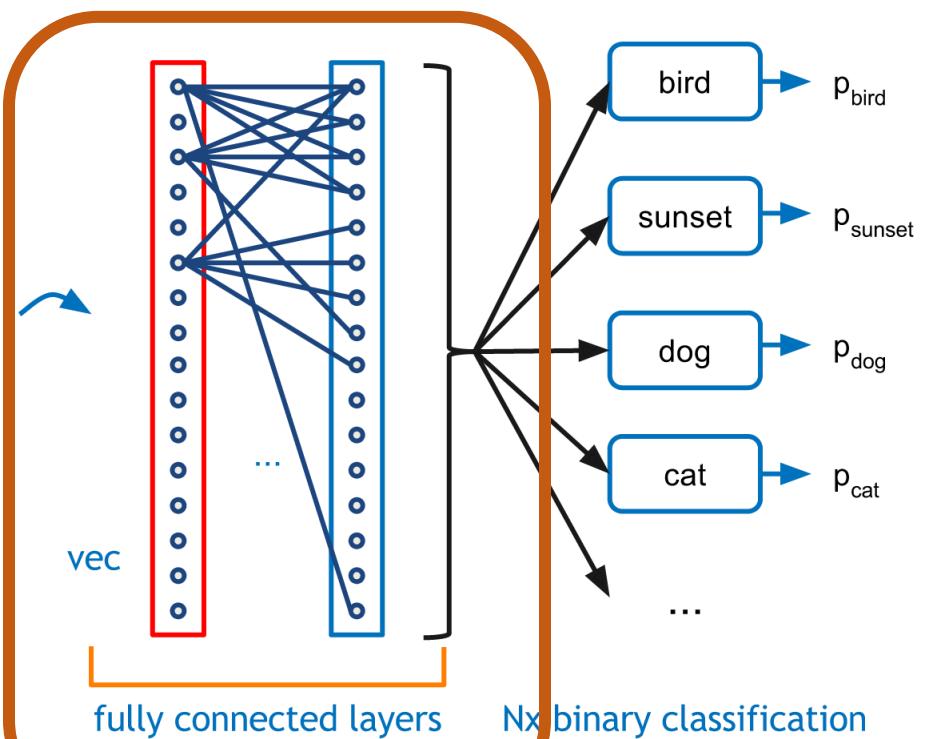




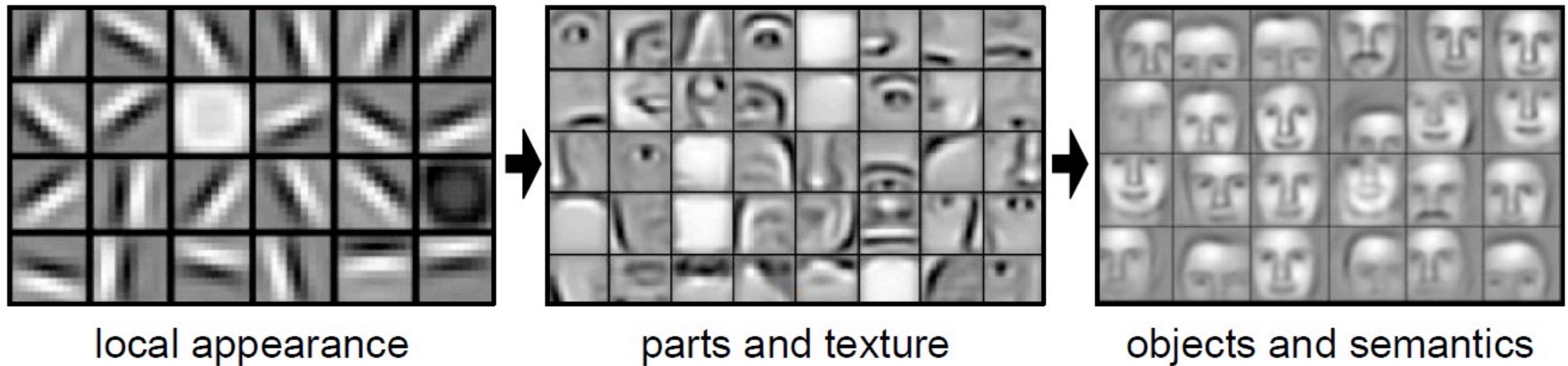
convolution +  
nonlinearity



convolution + pooling layers



# Features in End-to-end learning



**Learn the representation from data**

*[figure credit H. Lee]*

# Computer Vision Tasks

**Classification**



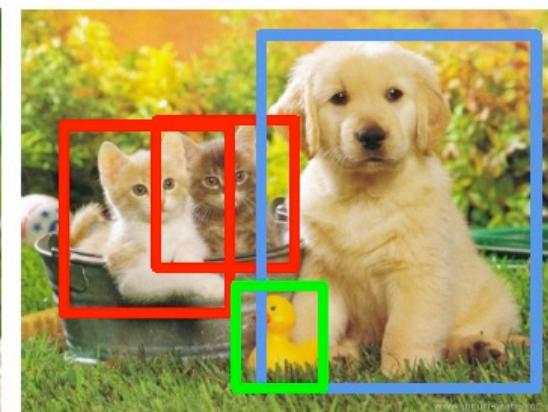
CAT

**Classification + Localization**



CAT

**Object Detection**



CAT, DOG, DUCK

**Instance Segmentation**



CAT, DOG, DUCK

Single object

Multiple objects

# Challenges

Viewpoint variation



Scale variation



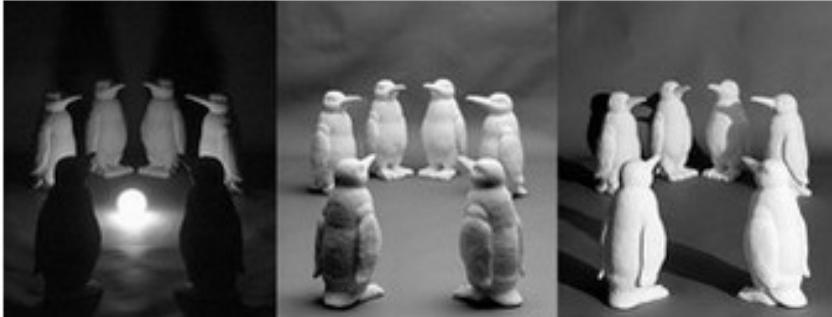
Deformation



Occlusion



Illumination conditions



Background clutter

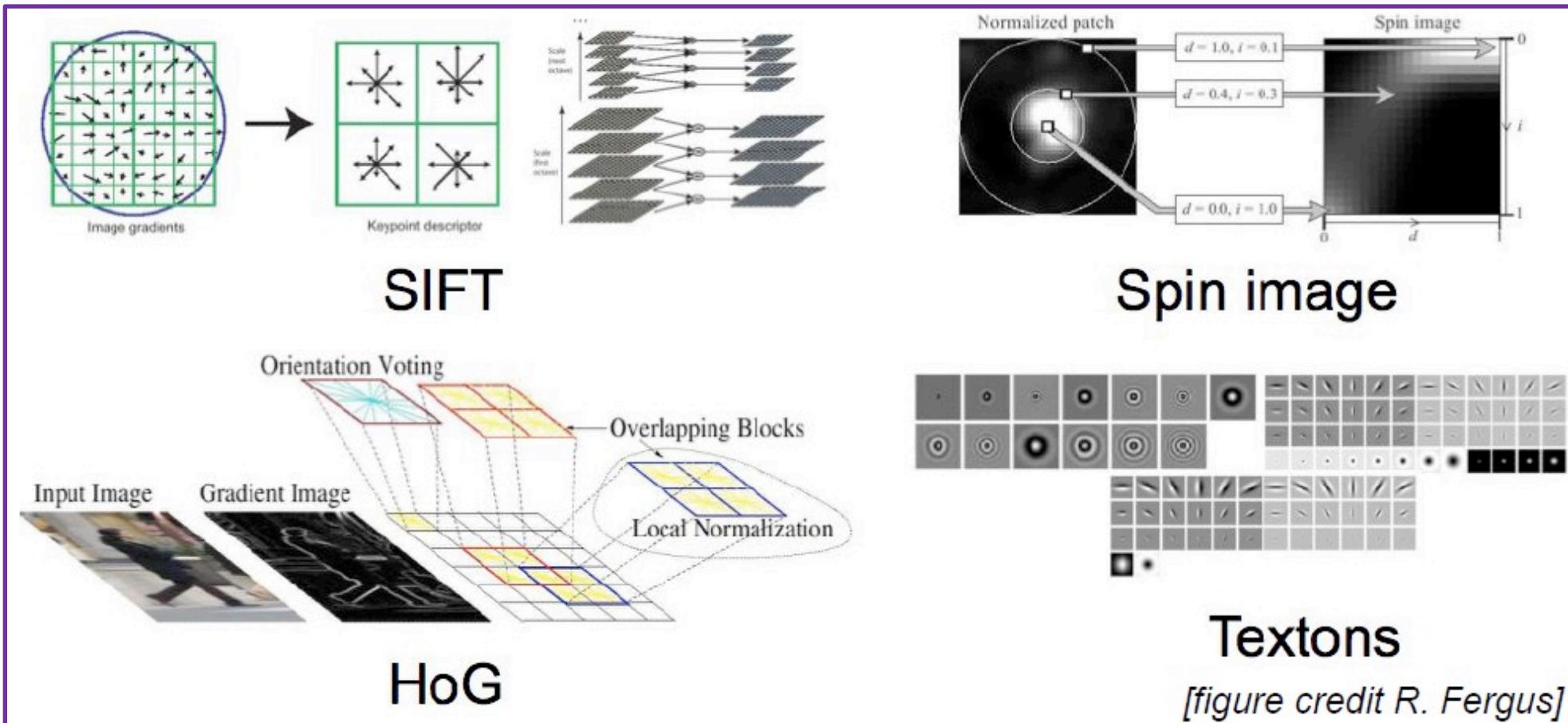


Intra-class variation



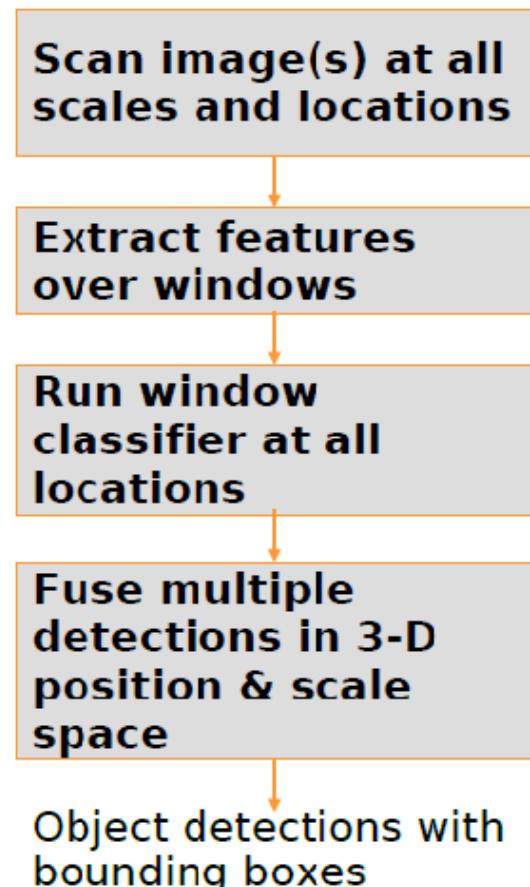
# Hand Engineered Features

- Years of Vision algorithms

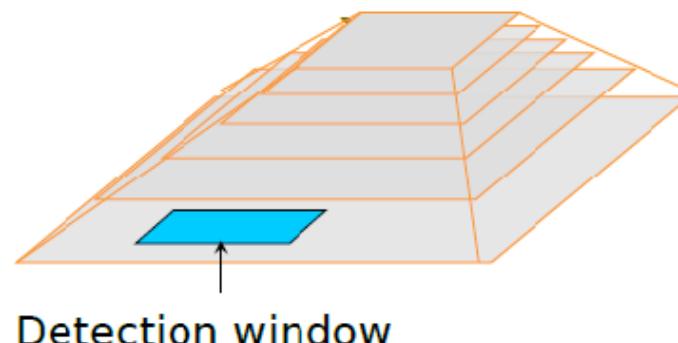


# Image Scanning Detectors

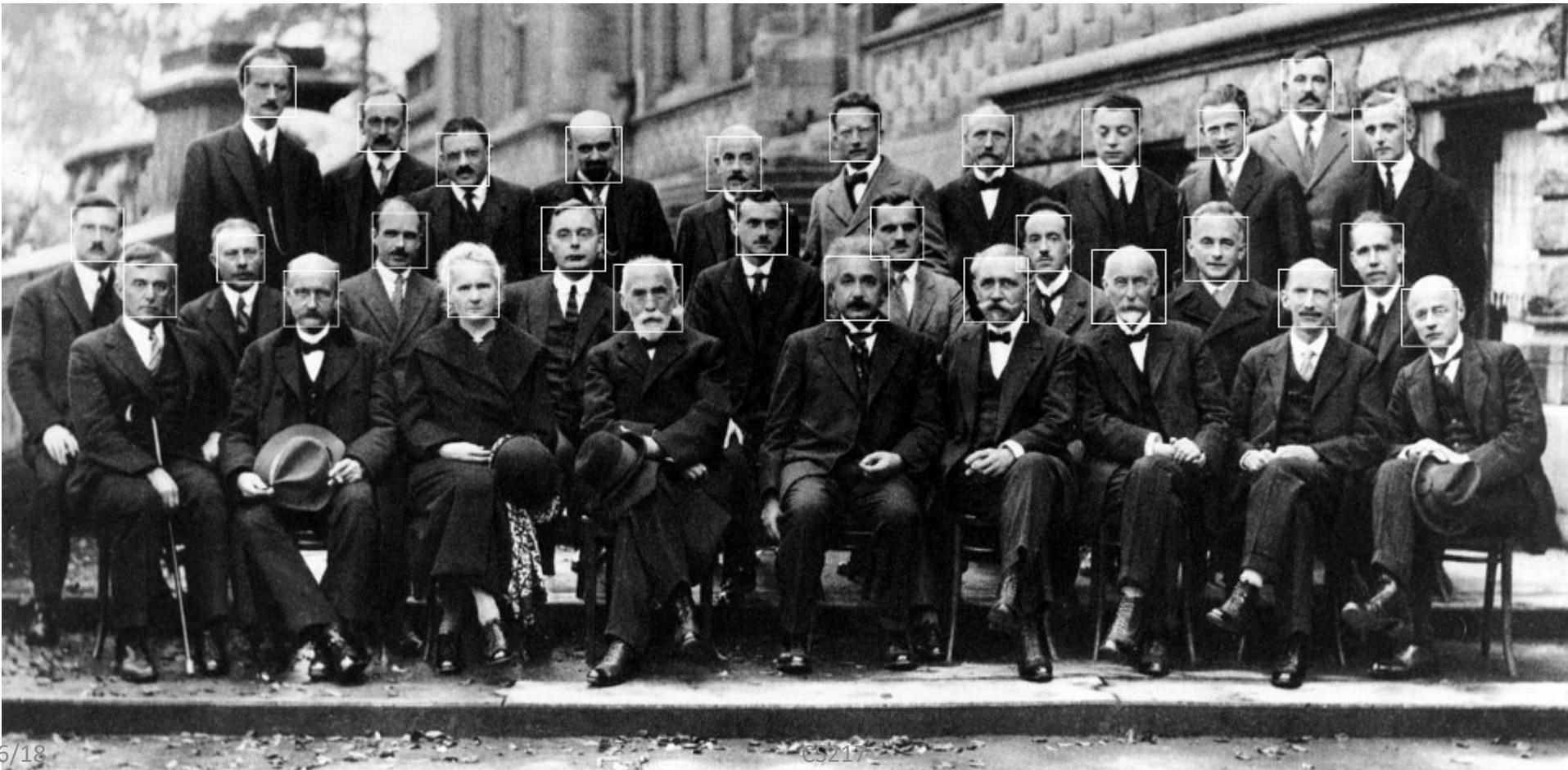
## Detection Phase



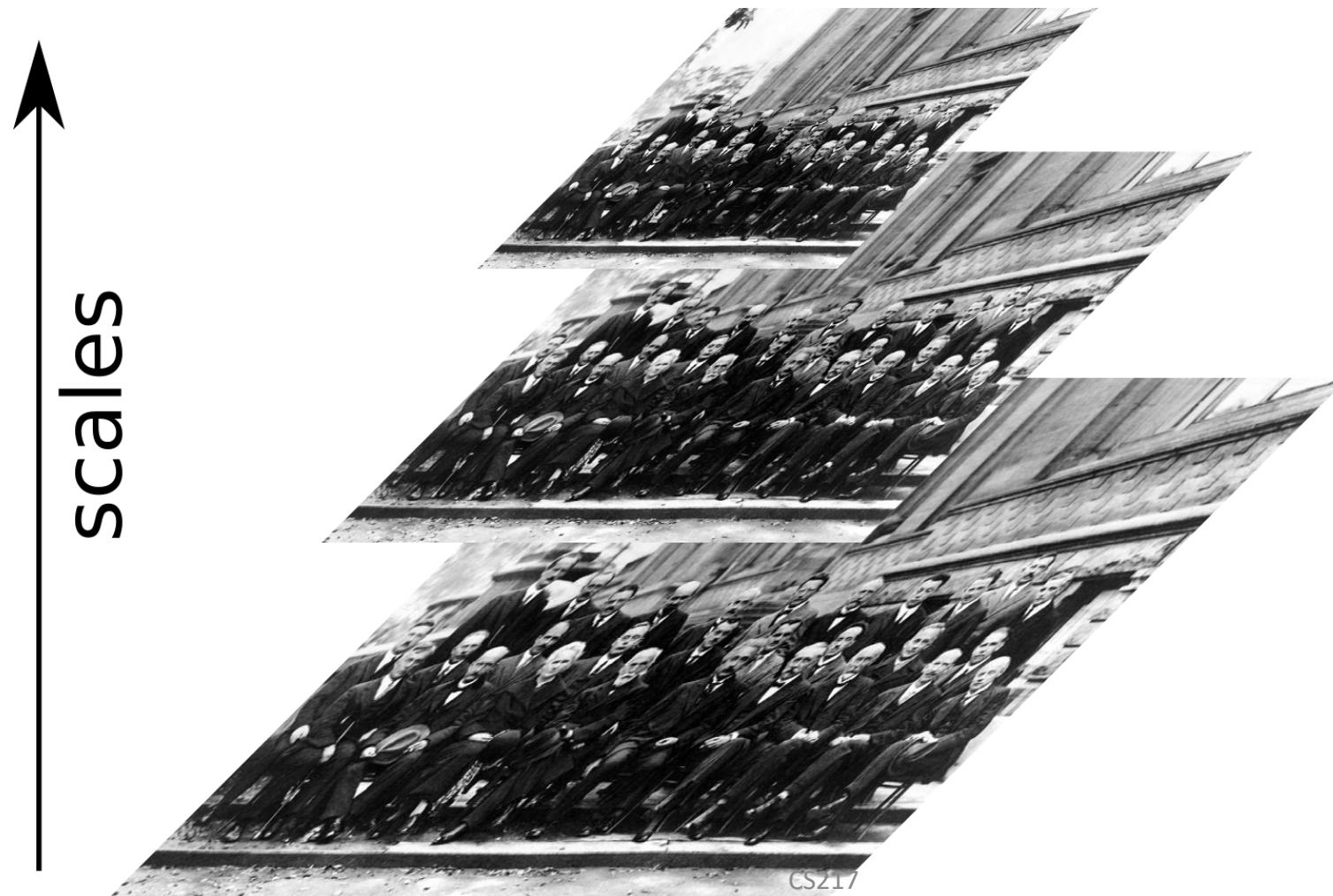
Scale-space pyramid



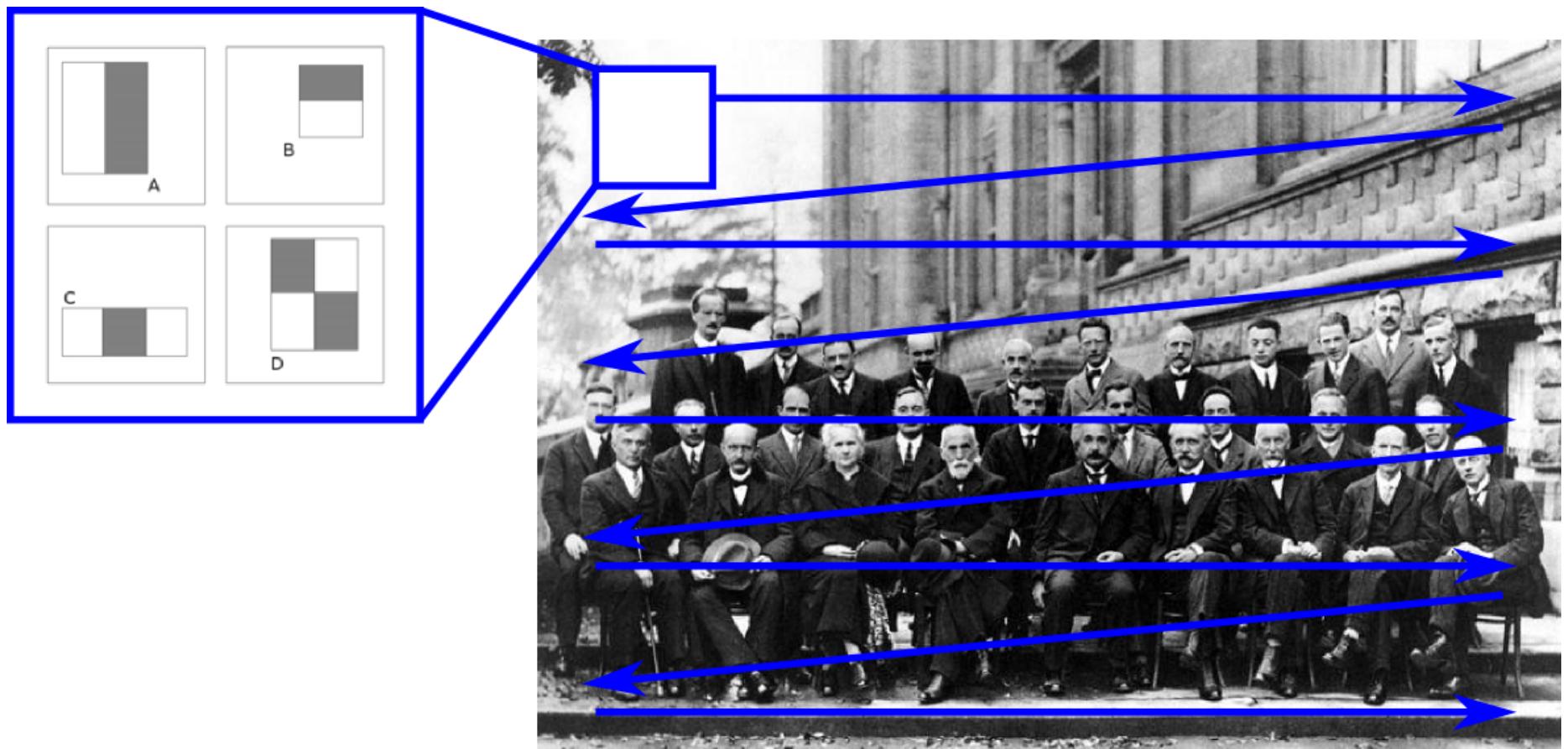
# Viola Jones Face Detection

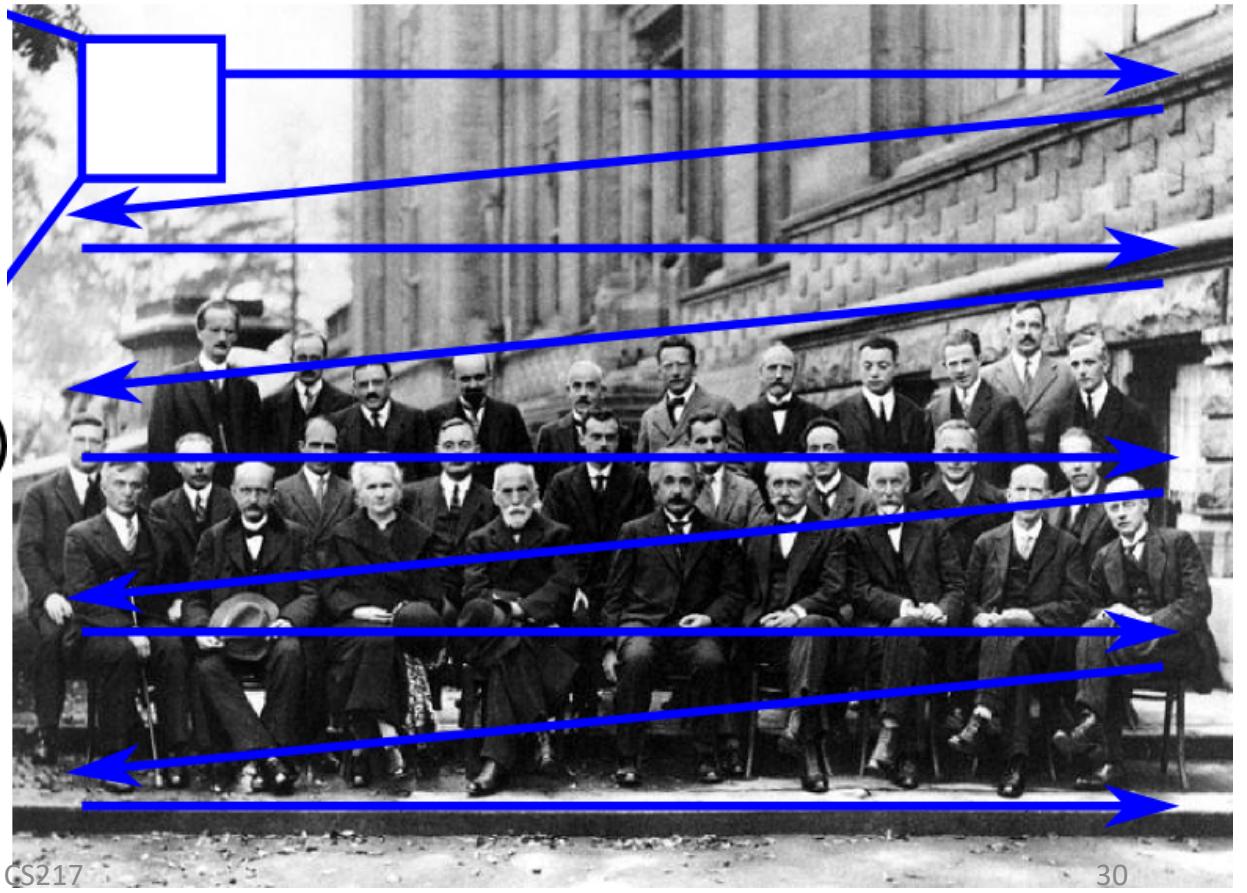
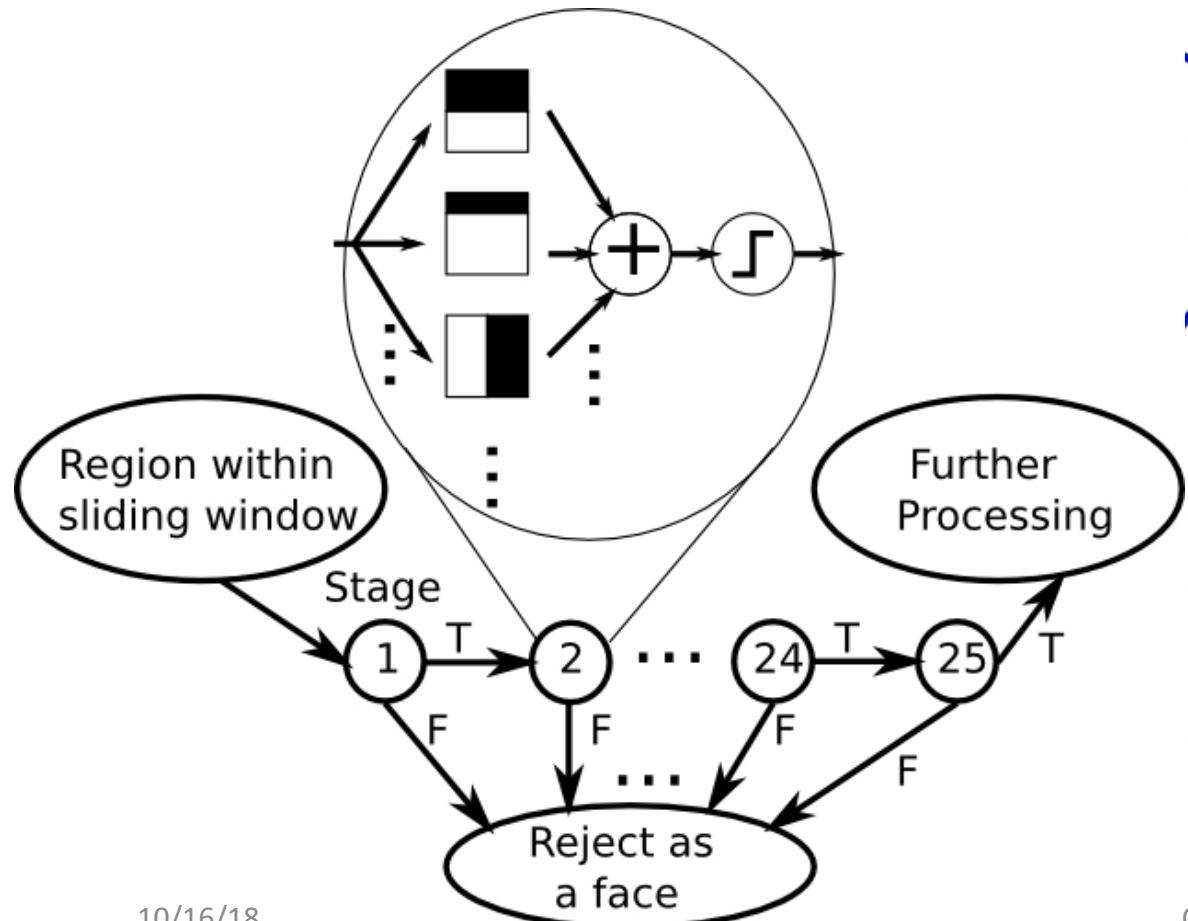


# IMAGE PYRAMID

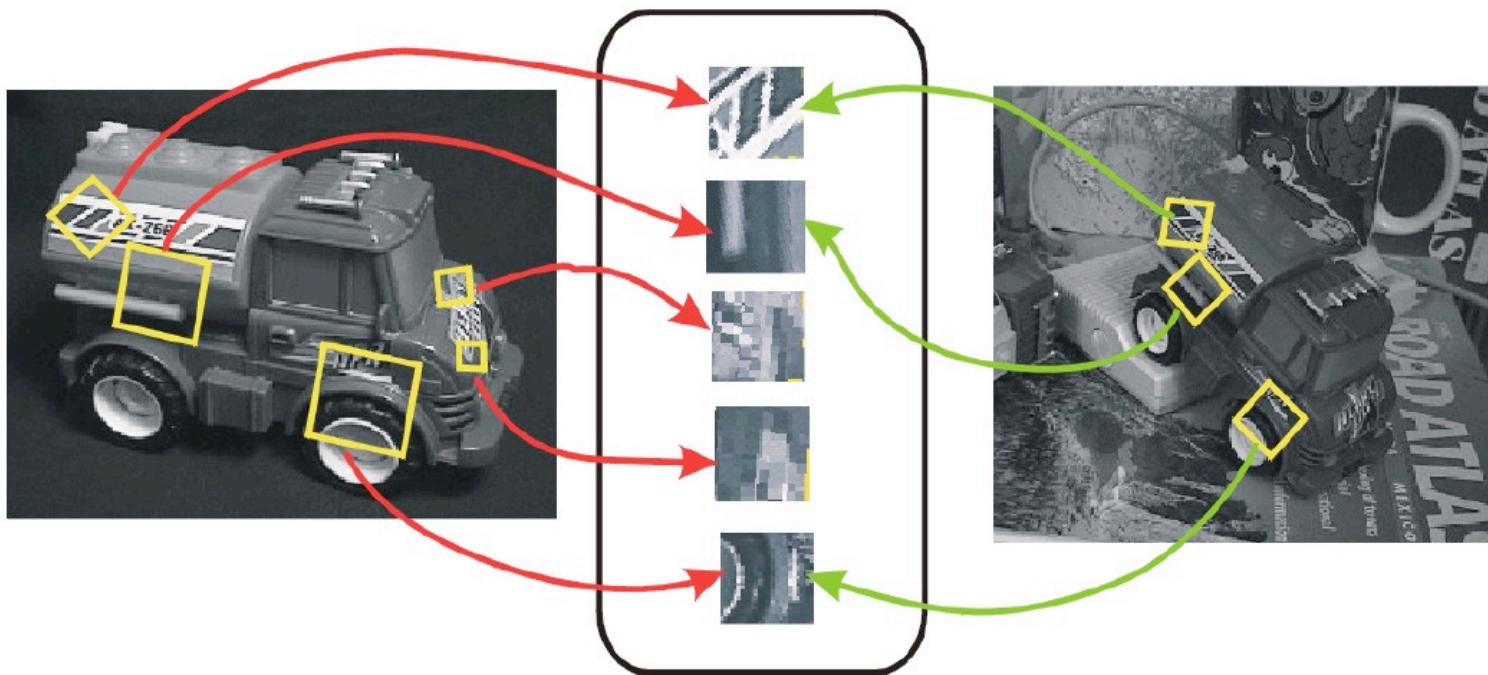


# SLIDING WINDOW and COMPARE





# Feature Descriptors, Detection and Matching



- Encode distinctive local structure at a collection of image points for matching between images
- Modest changes in viewing conditions
  - scale
  - orientation
  - Contrast

# Descriptors

- A good Tutorial:

<https://gilscvblog.com/2013/08/18/a-short-introduction-to-descriptors/>

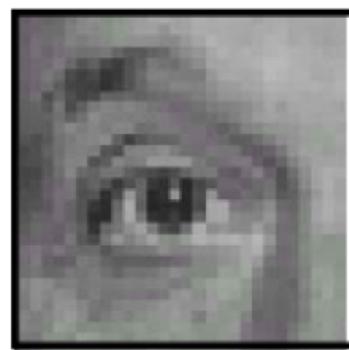
- Scale Invariant Feature Transformation (SIFT)
- Histogram of Oriented Gradients (HOG)

- Feature point detection
- Feature descriptors
- Potential applications –
  - Image panoramas
  - image matching and registration
  - Long range motion (tracking by detection)
  - Stereoscopic vision / 3D reconstruction
  - Object recognition

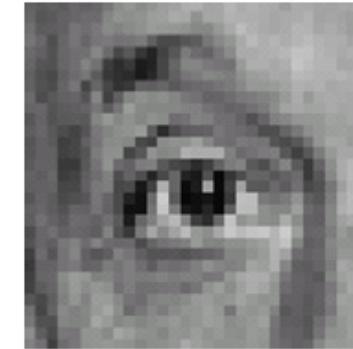
# Convolution and Filtering



$$\begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{matrix}$$



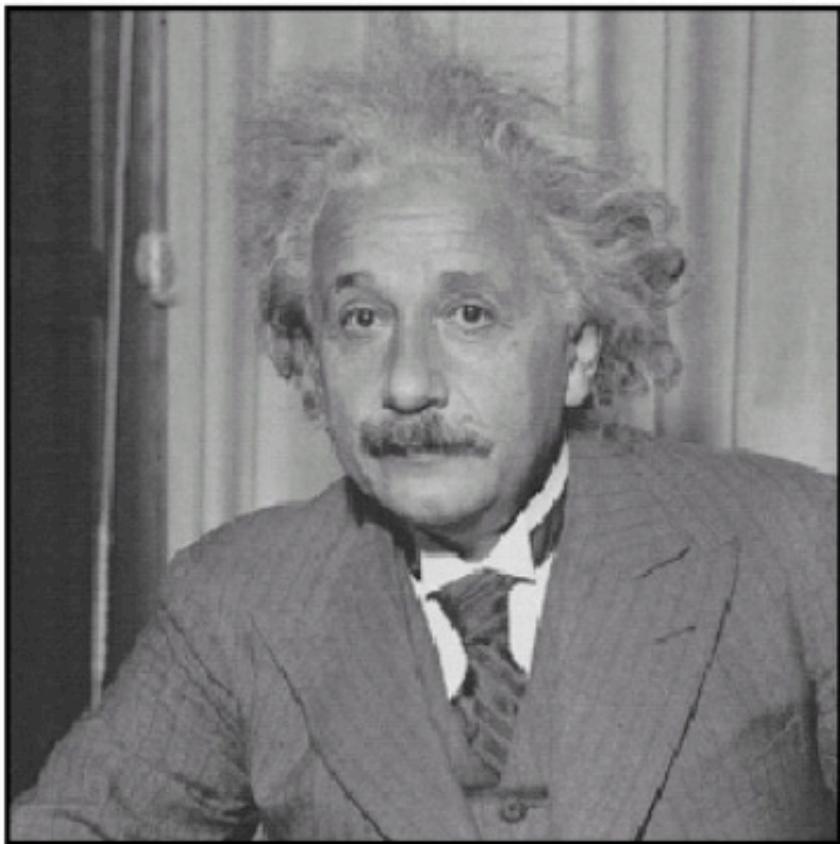
$$* \left( \begin{matrix} 0 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{matrix} - \frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix} \right) =$$



Sharpening filter  
(accentuates edges)

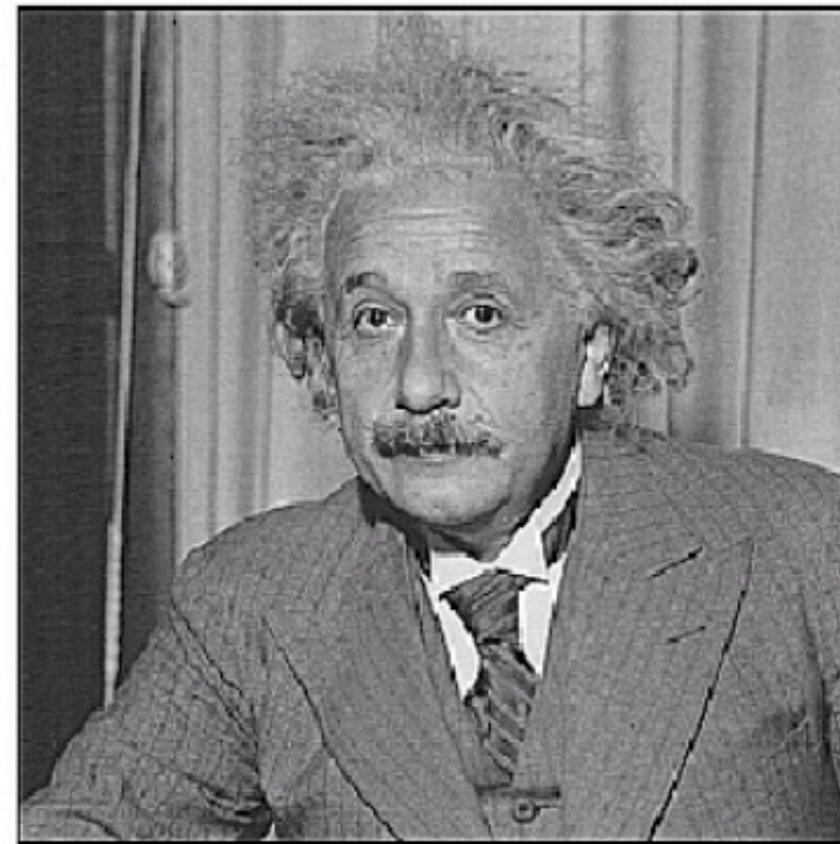
- <http://www.cs.toronto.edu/~urtasun/courses/CV/lecture02.pdf>

# CONVOLUTION AND FILTERING



**before**

10/16/18

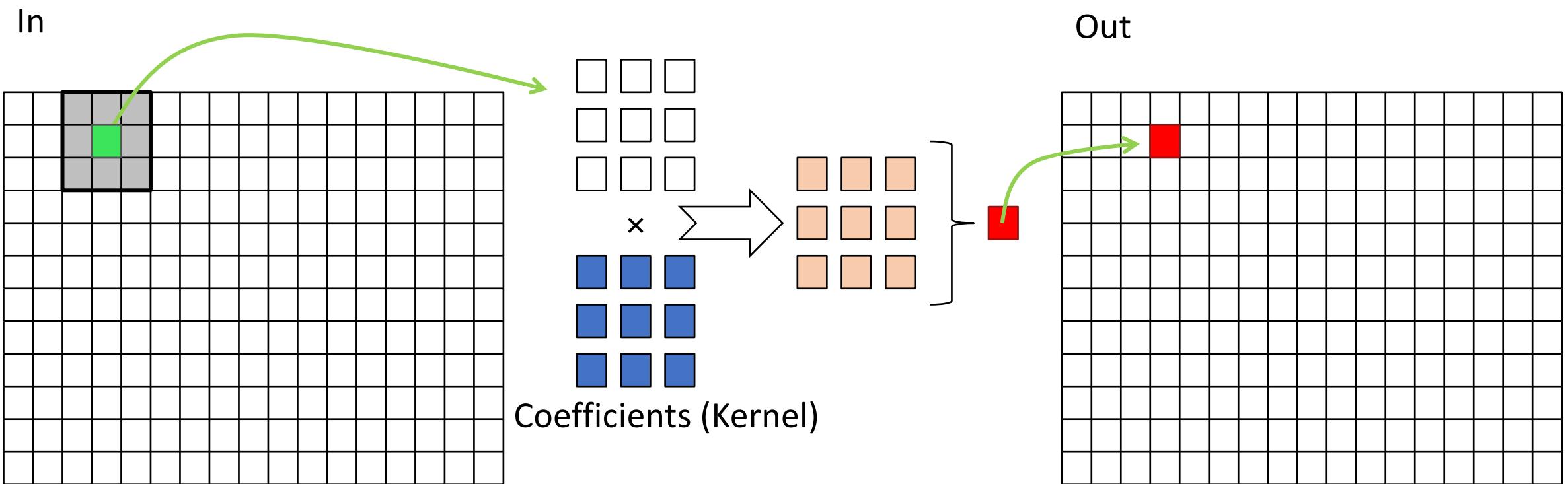


**after**

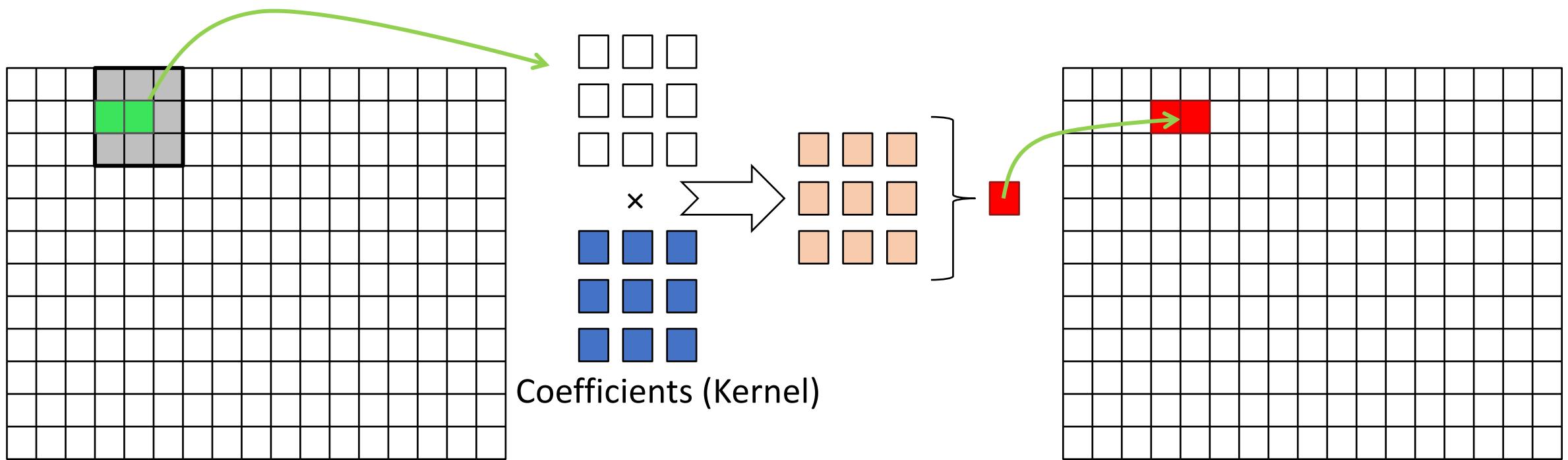
CS217

34

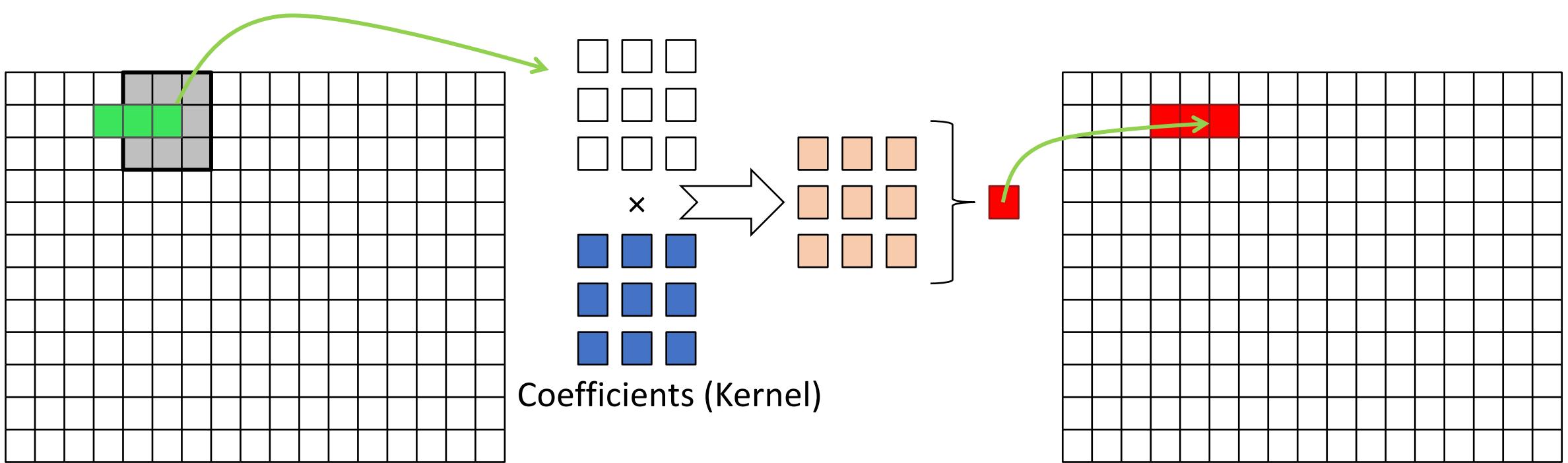
# Convolution Operation



# Convolution Operation



# Convolution Operation

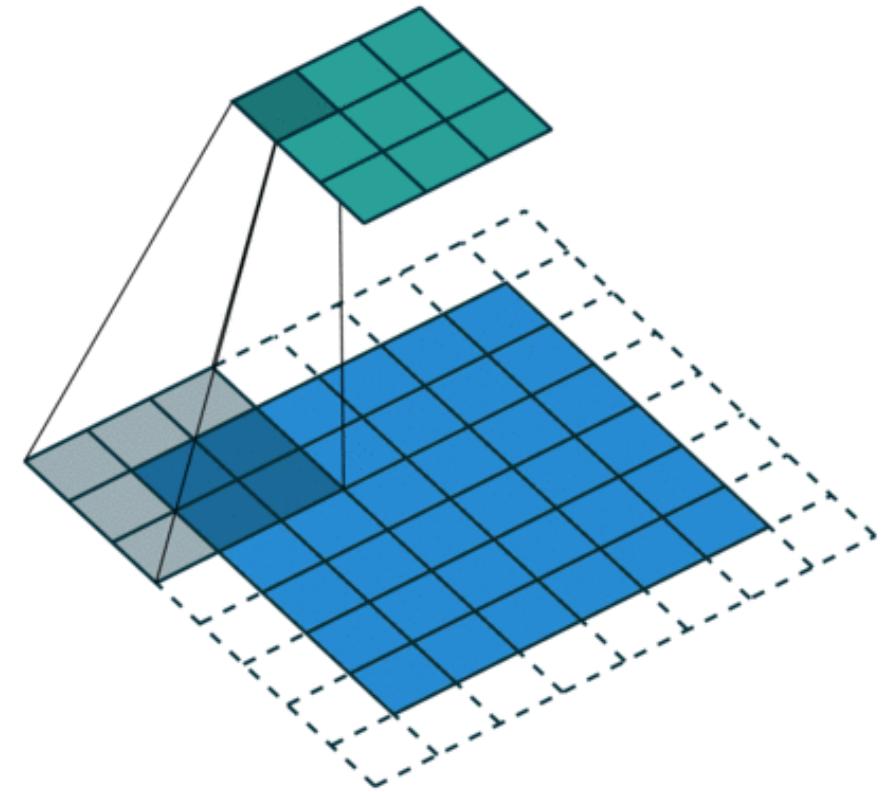


# Convolution Operation

3 <sub>0</sub>	3 <sub>1</sub>	2 <sub>2</sub>	1	0
0 <sub>2</sub>	0 <sub>2</sub>	1 <sub>0</sub>	3	1
3 <sub>0</sub>	1 <sub>1</sub>	2 <sub>2</sub>	2	3
2	0	0	2	2
2	0	0	0	1

12	12	17
10	17	19
9	6	14

No zero padding, unit strides



Zero padding, non-unit strides

Chihuahua or Muffin?



10/16/18

That Is The  
Question



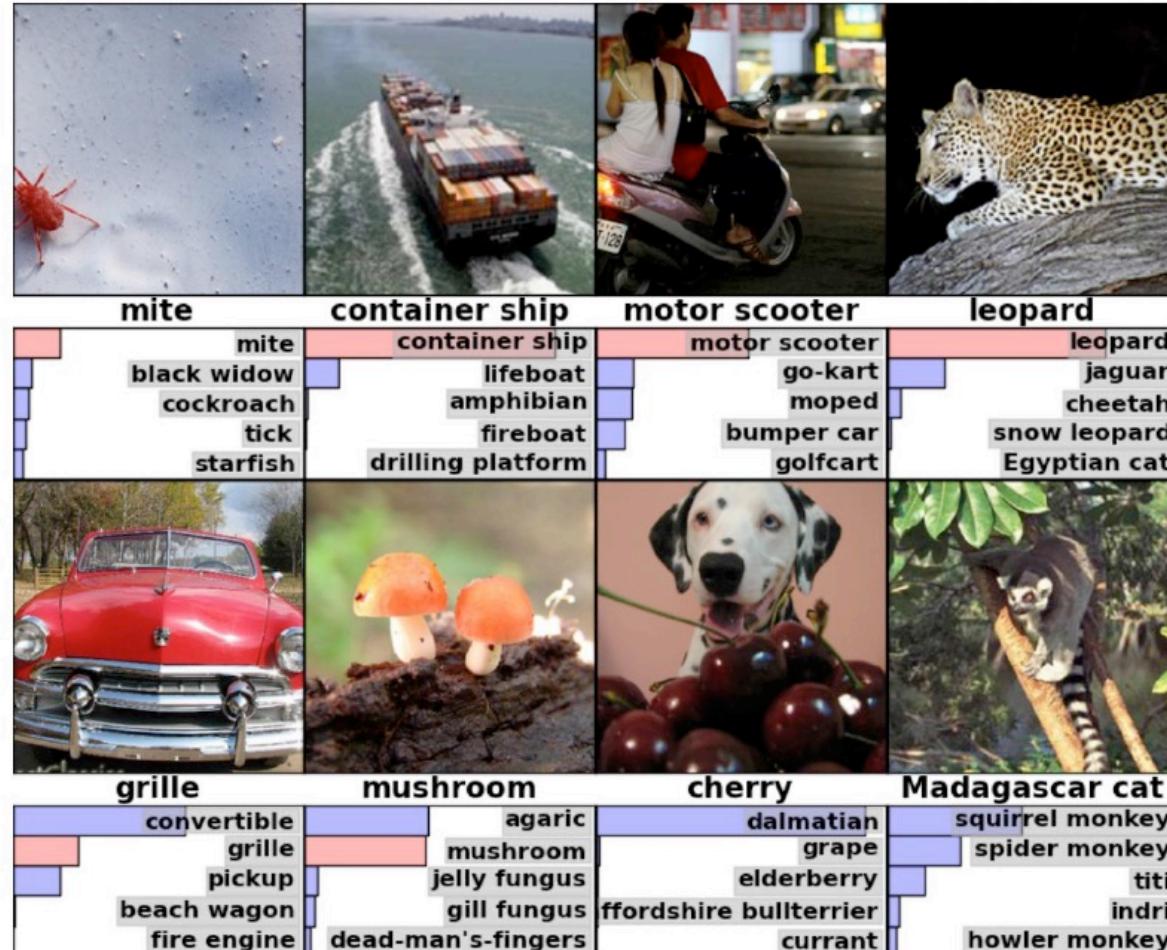
CS217

39

# ImageNET challenge

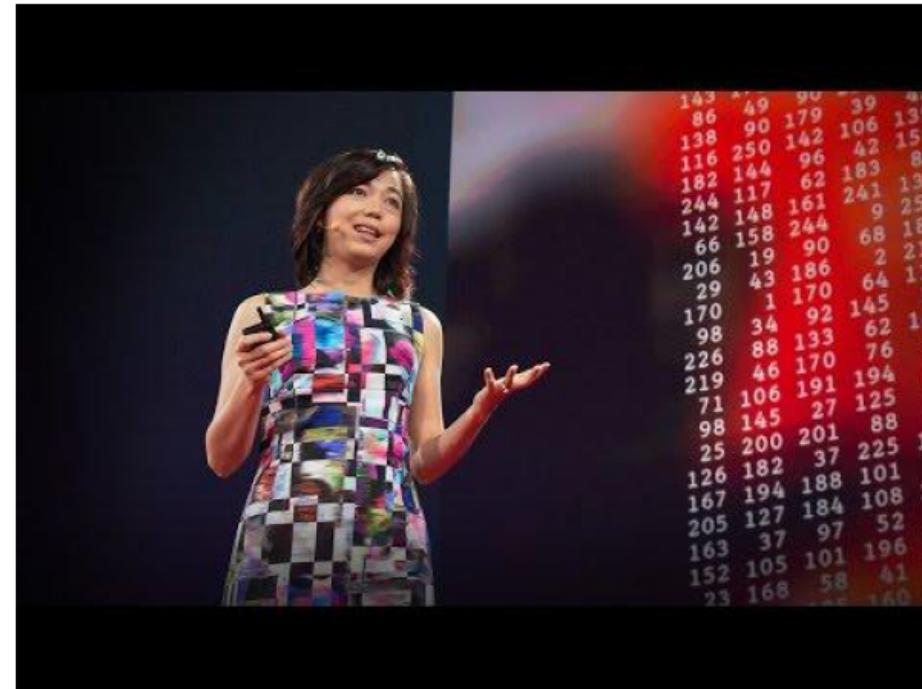
IMAGENET

- 1,000 object classes (categories).
- Images:
  - 1.2 M train
  - 100k test.

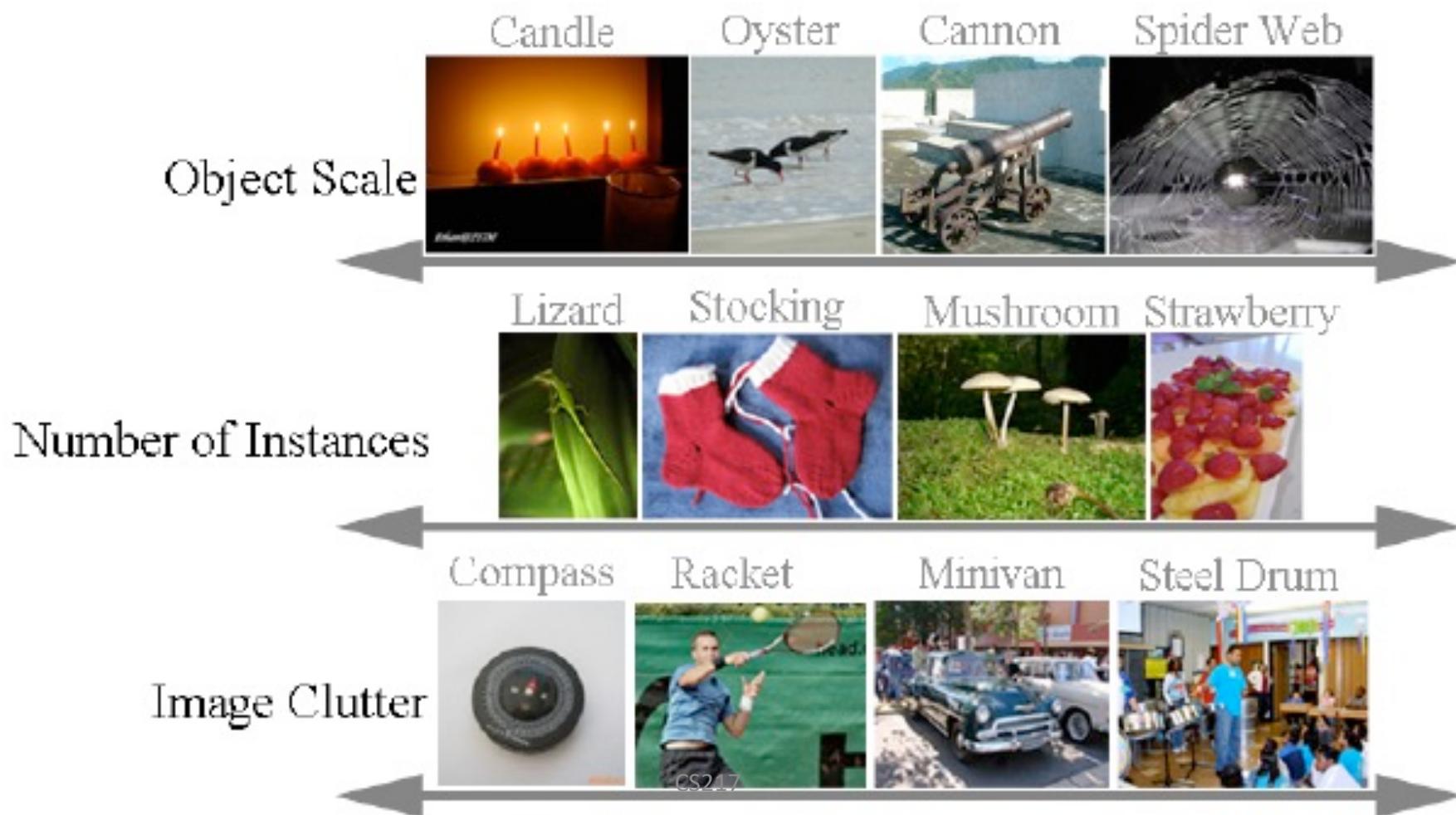


# Imagenet dataset

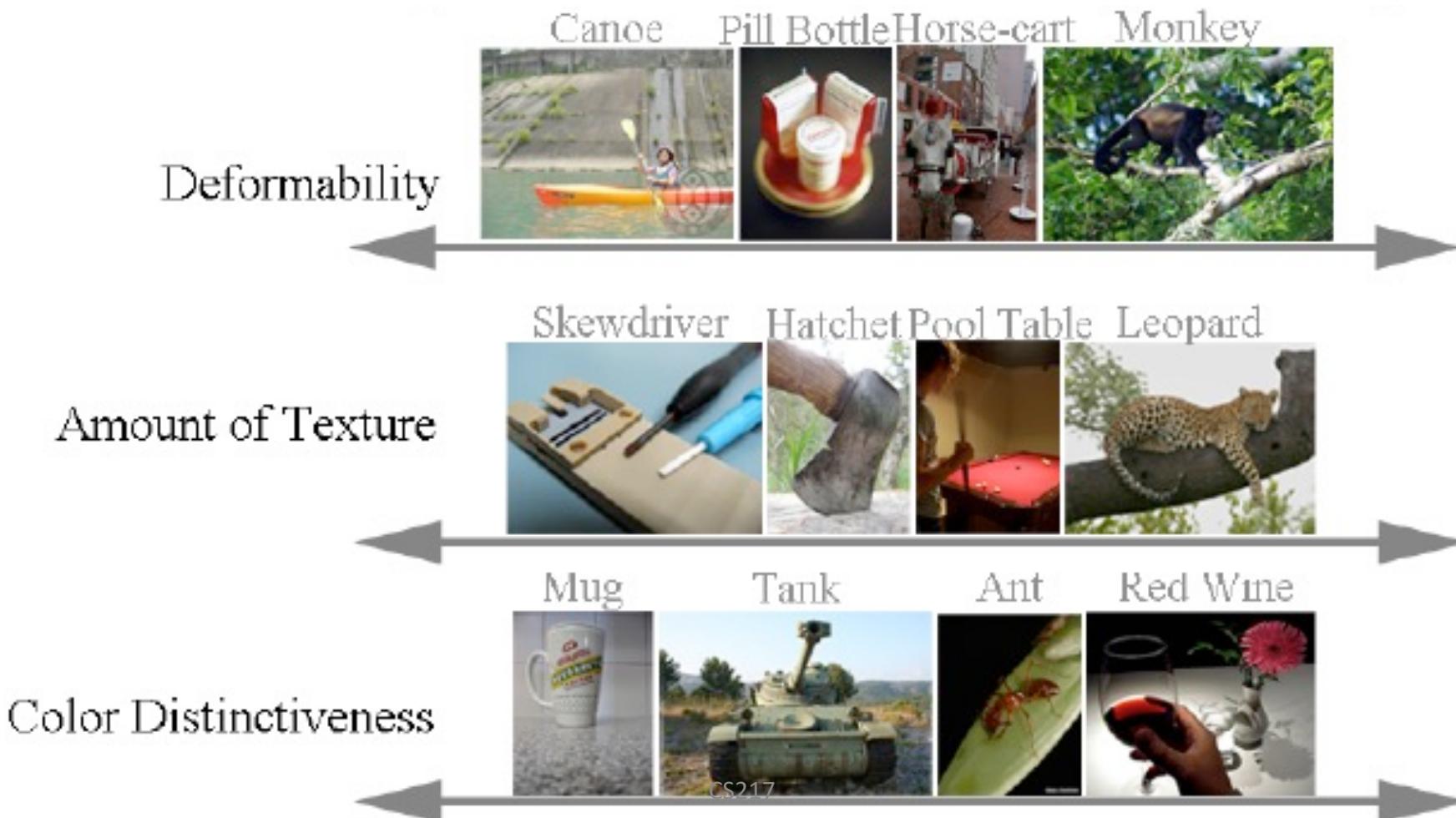
Li Fei-Fei, ["How we're teaching computers to understand pictures"](#) TEDTalks 2014.



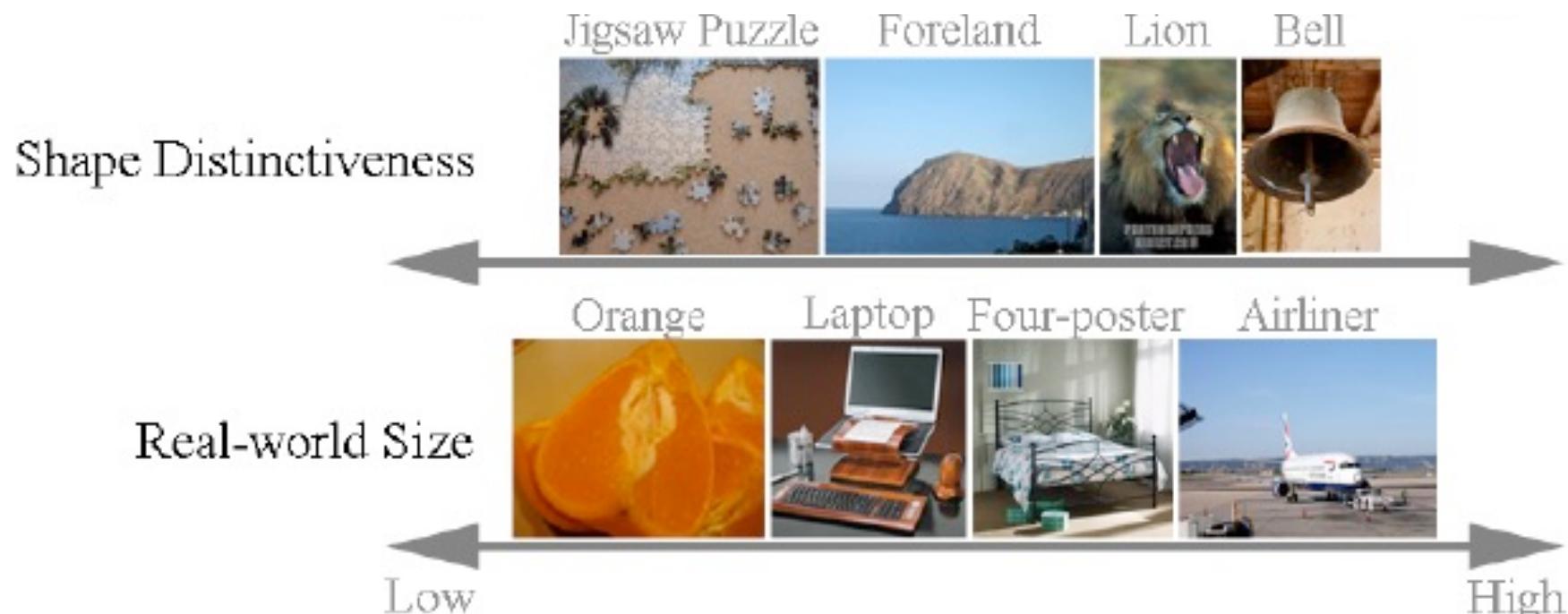
# IMAGENET



# IMAGENET

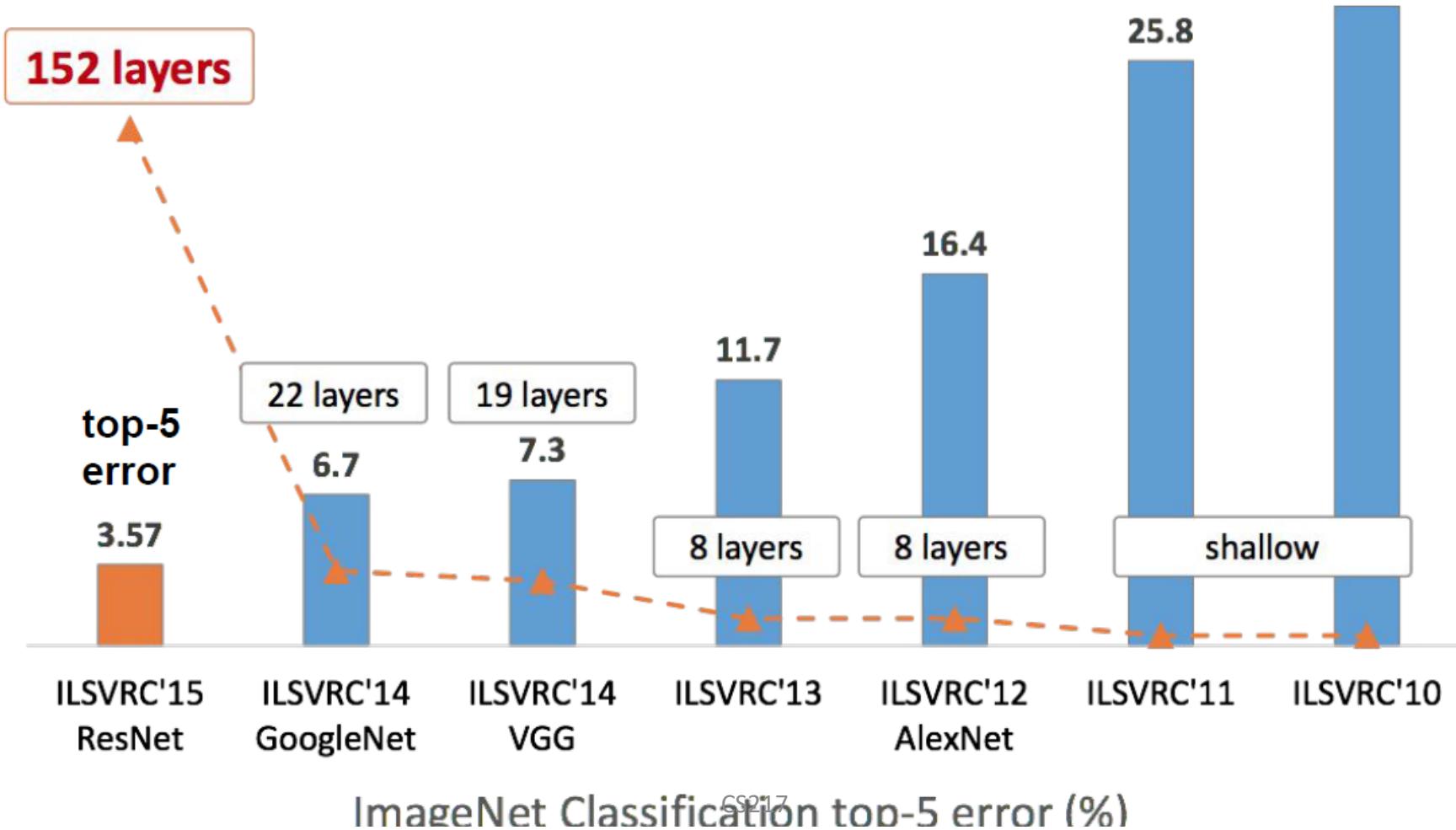


# IMAGENET



# IMAGENET CHALLENGE WINNERS

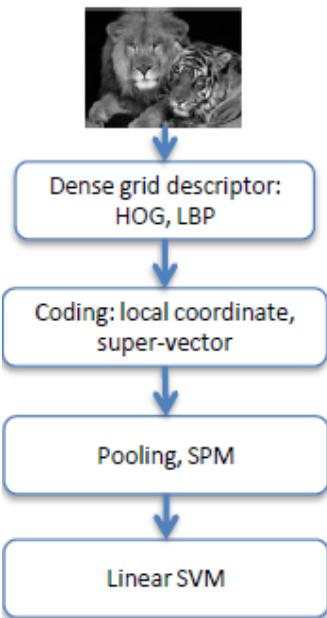
## SHALLOW vs. DEEP



# DEEP Models Beat Shallow Models

## Year 2010

NEC-UIUC

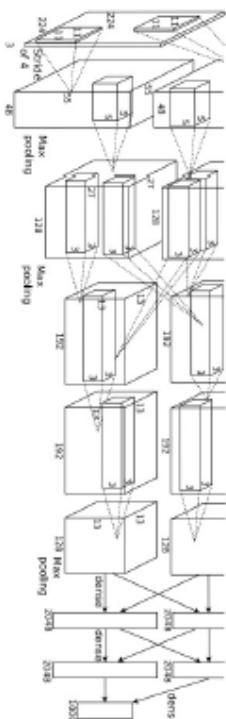


[In CVPR 2011]

10/16/18

Year 2012

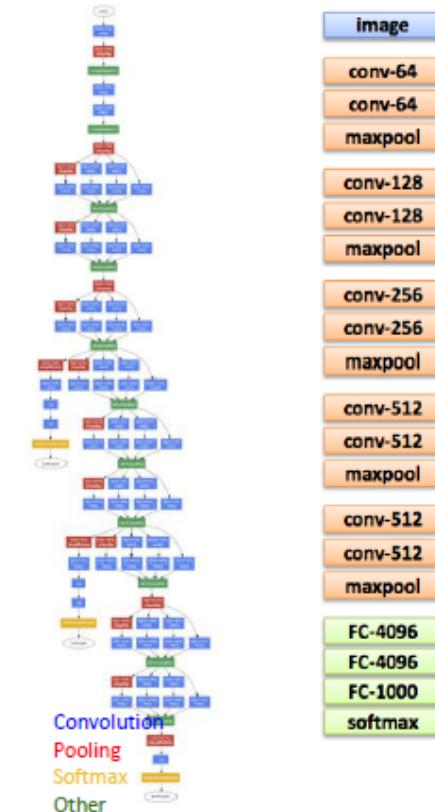
## SuperVision



[Krizhevsky NIPS 2012]

Year 2014

## GoogLeNet      VGG

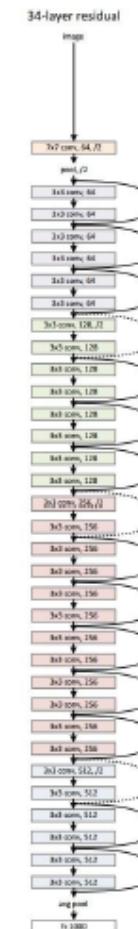


[Szegedy arxiv 2014]

CS217

**Year 2015**

MSRA



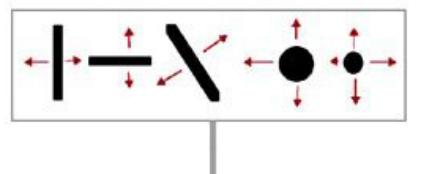
46



**WE NEED TO GO**

**DEEPER**

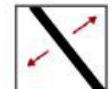
memegenerator.net



**Simple Cells: Response to light orientation**

**Complex Cells: Response to light orientation & movement**

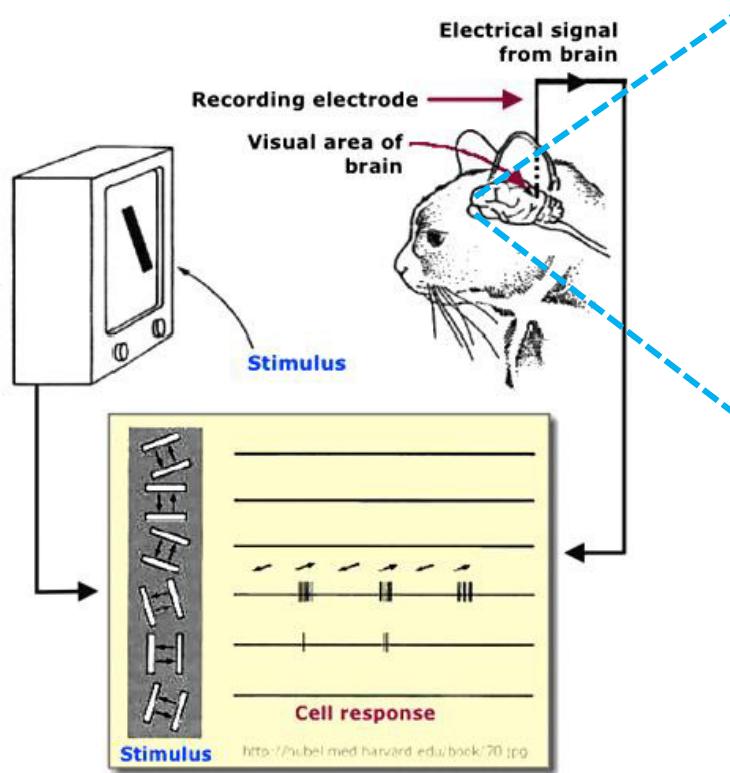
**Hypercomplex Cells: Response to movement with an end point**



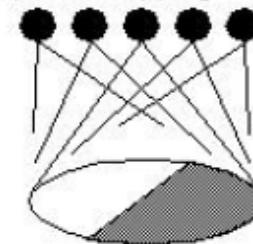
No response



Response (end point)



Hubel & Weisel  
topographical mapping



### featural hierarchy



hyper-complex cells



complex cells



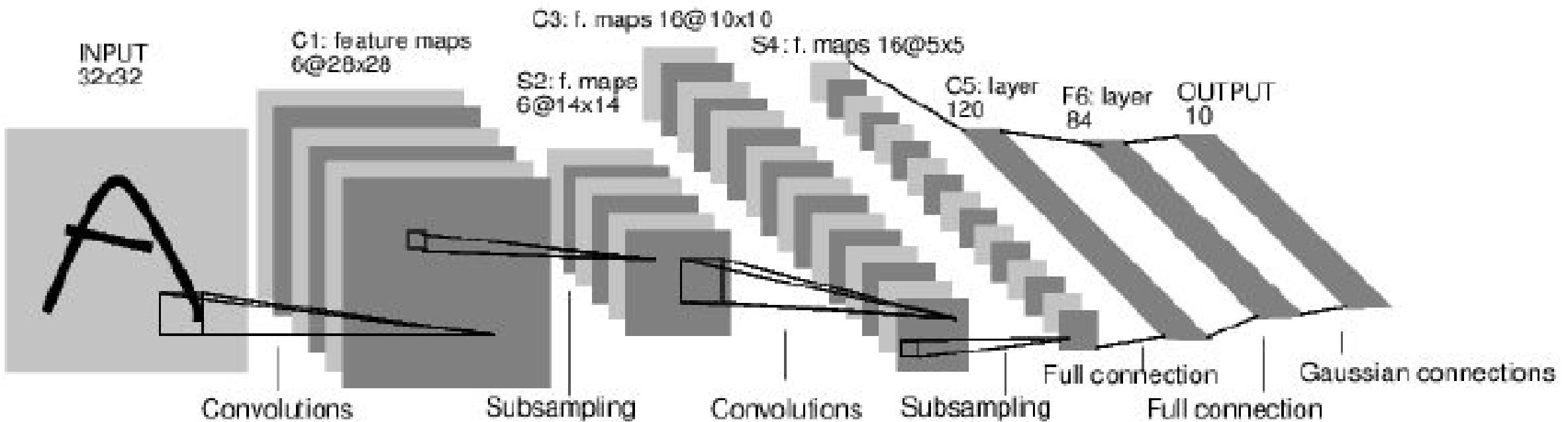
simple cells

- (□) high level
- (○) mid level
- (●) low level

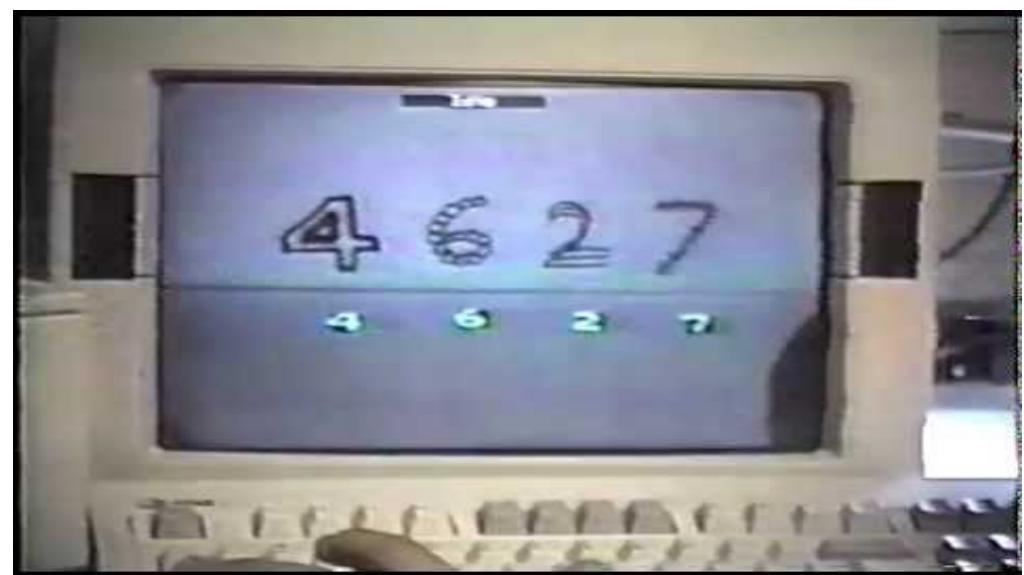
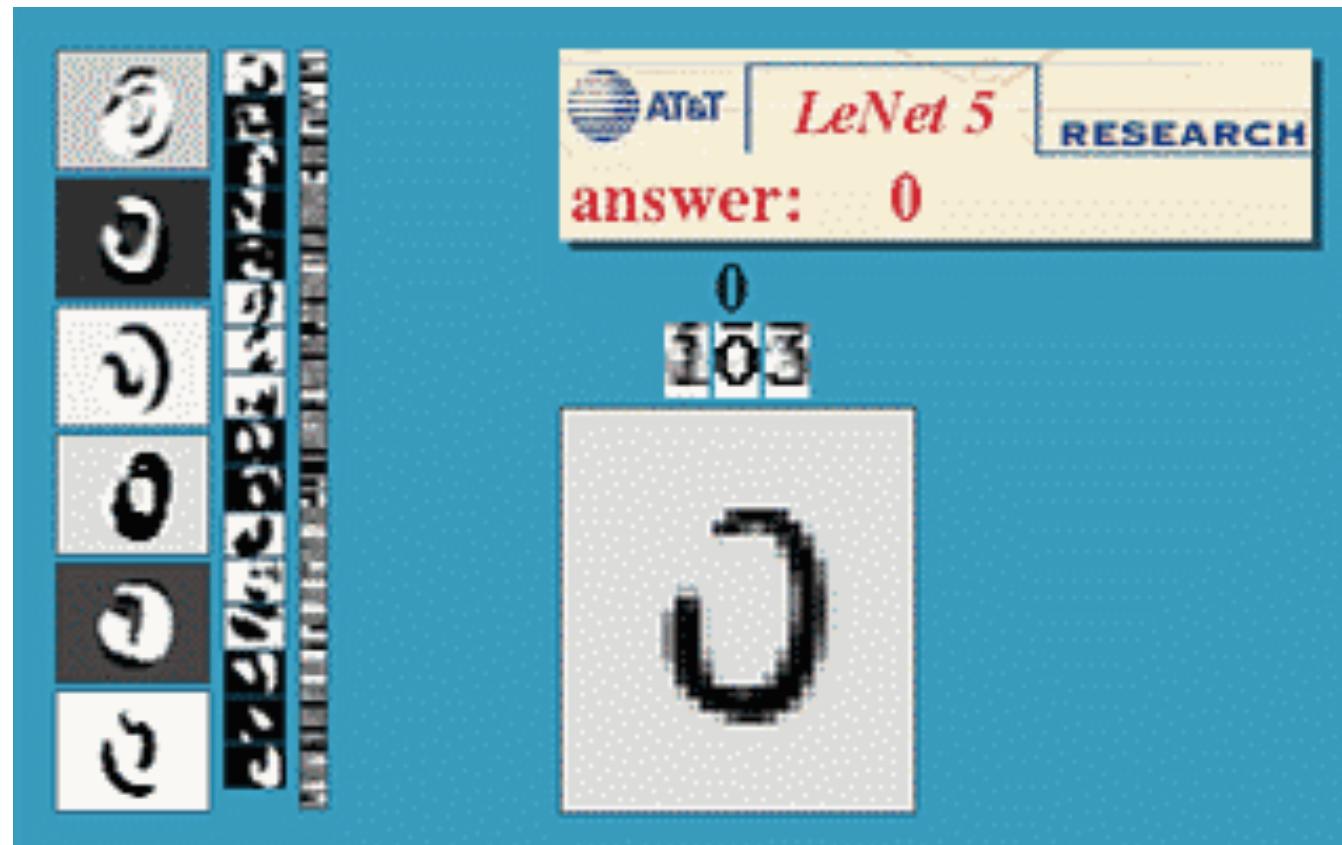
Hubel & Wiesel, 1959

CS217

# Le-NET 5: 1989

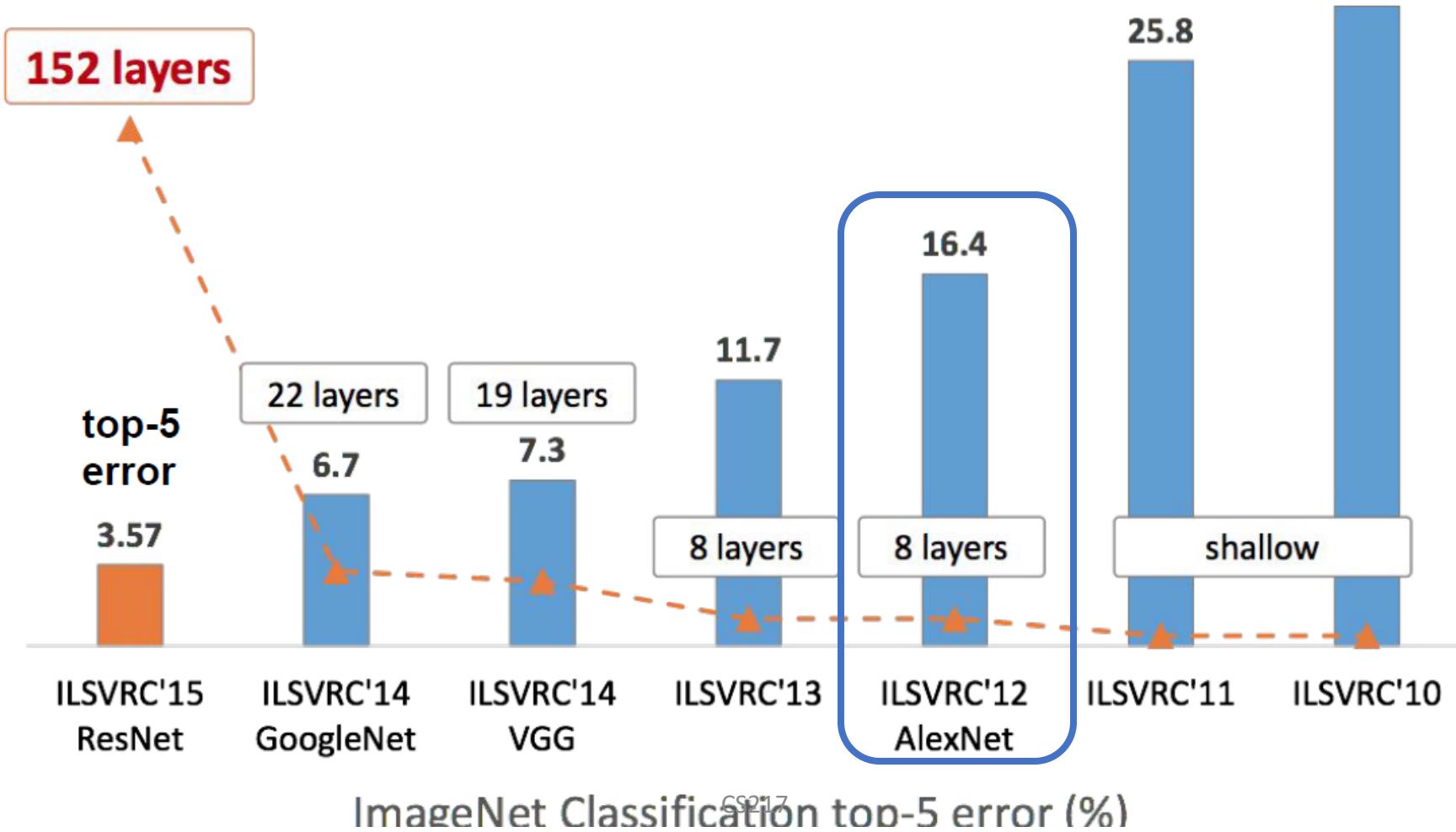


# LE-NET 5

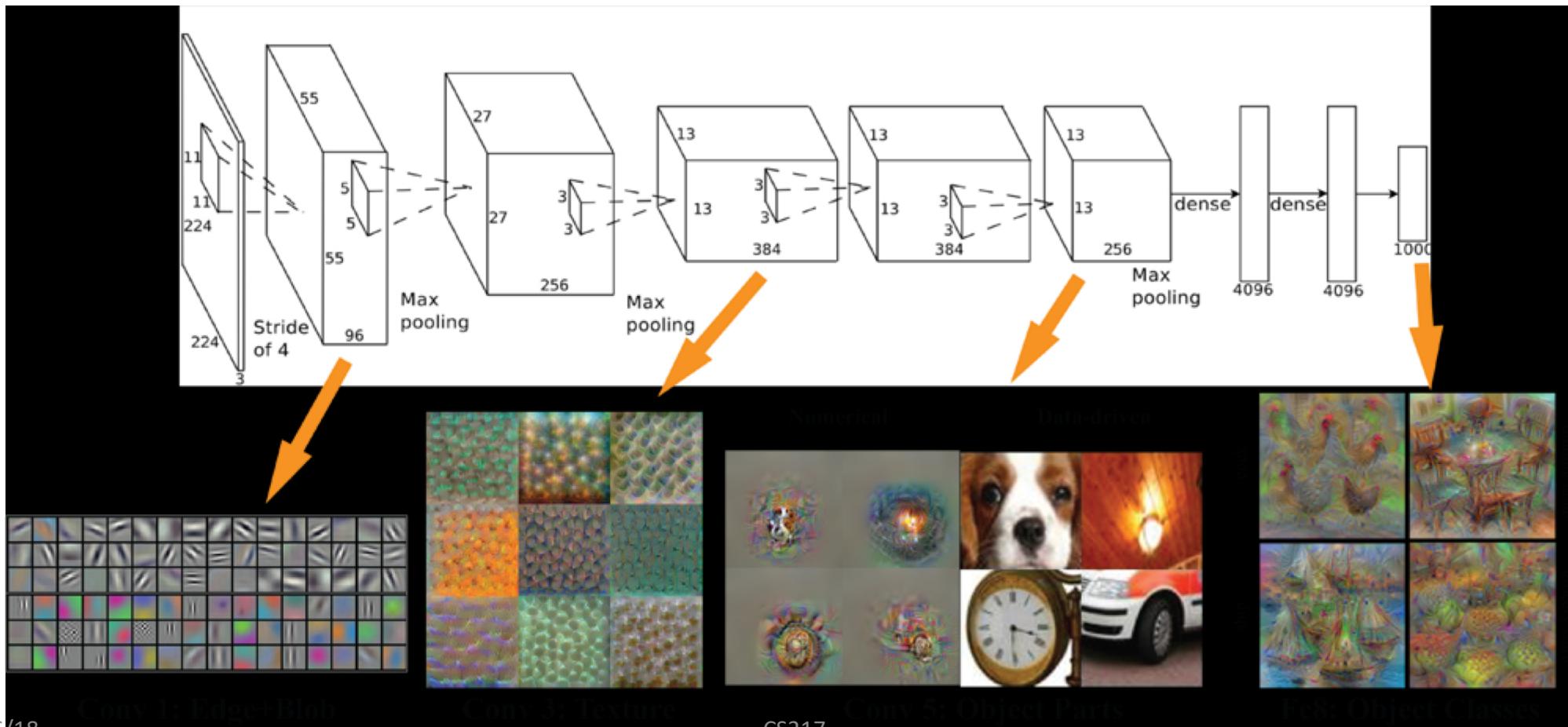


# IMAGENET CHALLENGE WINNERS

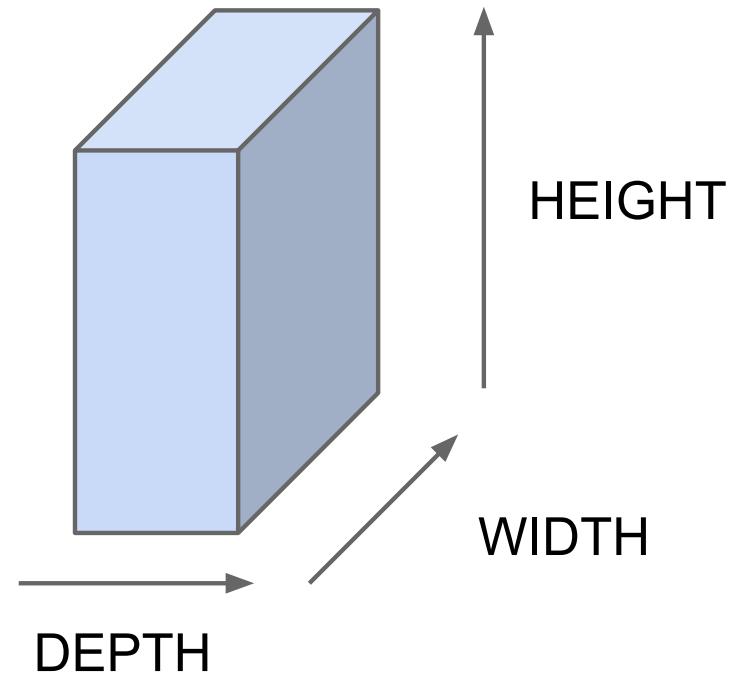
## SHALLOW vs. DEEP



# ALEX NET ARCHITECTURE



All Neural Net  
activations  
arranged in 3  
**dimensions:**

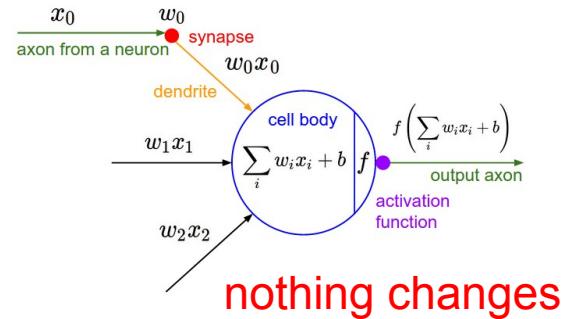
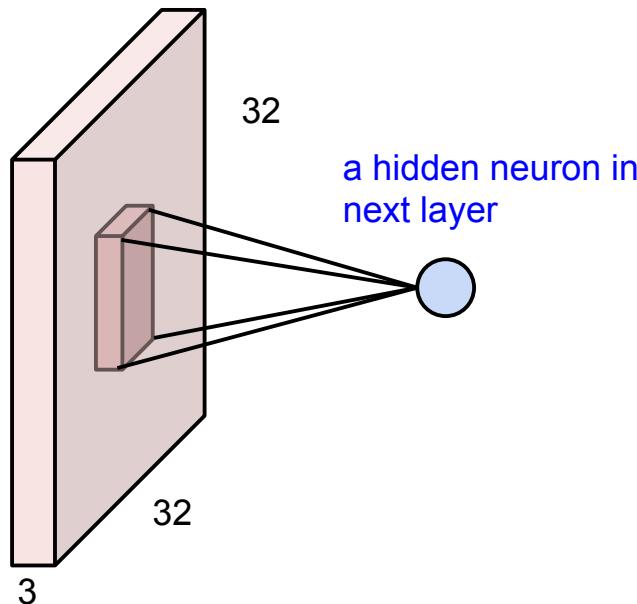


For example, a CIFAR-10 image is a  $32 \times 32 \times 3$  volume  
32 width, 32 height, 3 depth (RGB channels)

Source: [Fei-Fei Li & Andrej Karpathy]

# Local Neighborhoods

Convolutional Neural Networks  
are just Neural Networks BUT:  
**1. Local connectivity**



Source: [Fei-Fei Li & Andrej Karpathy]

# Local Neighborhoods

Convolutional Neural Networks  
are just Neural Networks BUT:  
**1. Local connectivity**

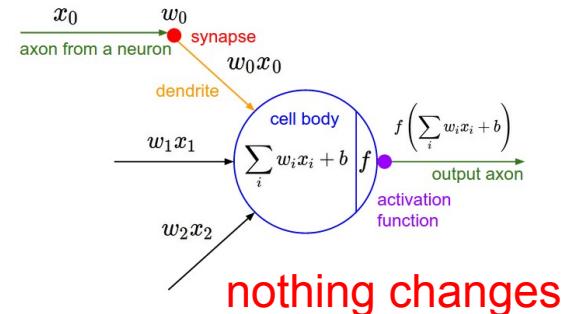
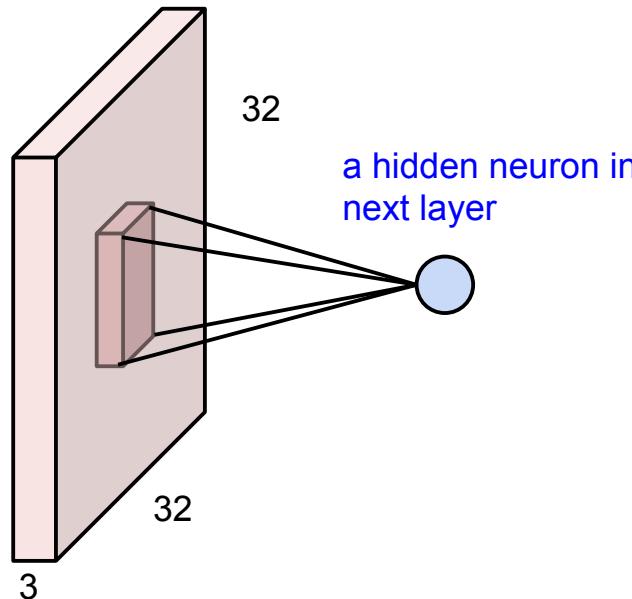


image: 32x32x3 volume

**before:** full connectivity: 32x32x3 weights

**now:** one neuron will connect to, e.g. 5x5x3 chunk and only have 5x5x3 weights.

note that connectivity is:

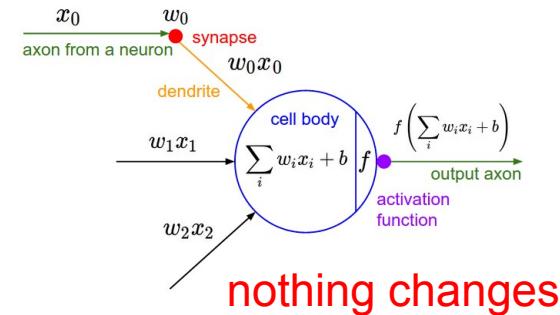
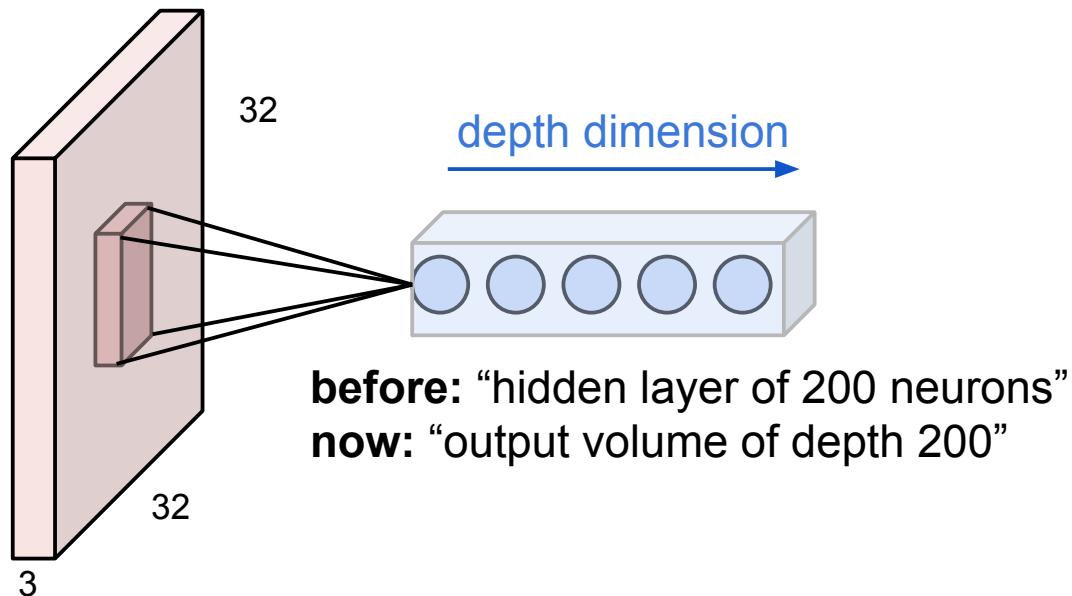
- local in space (5x5 inside 32x32)
- but full in depth (all 3 depth channels)

Source: [Fei-Fei Li & Andrej Karpathy]

# Kernels (Filters) Create Depth

Convolutional Neural Networks  
are just Neural Networks BUT:

## 1. Local connectivity

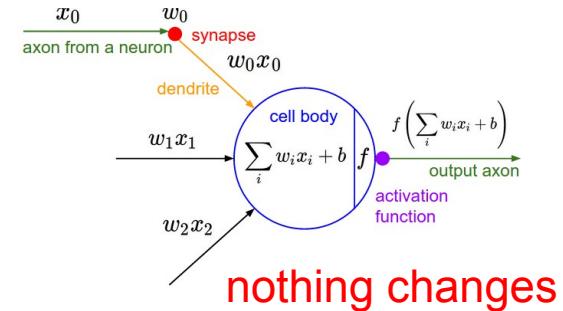
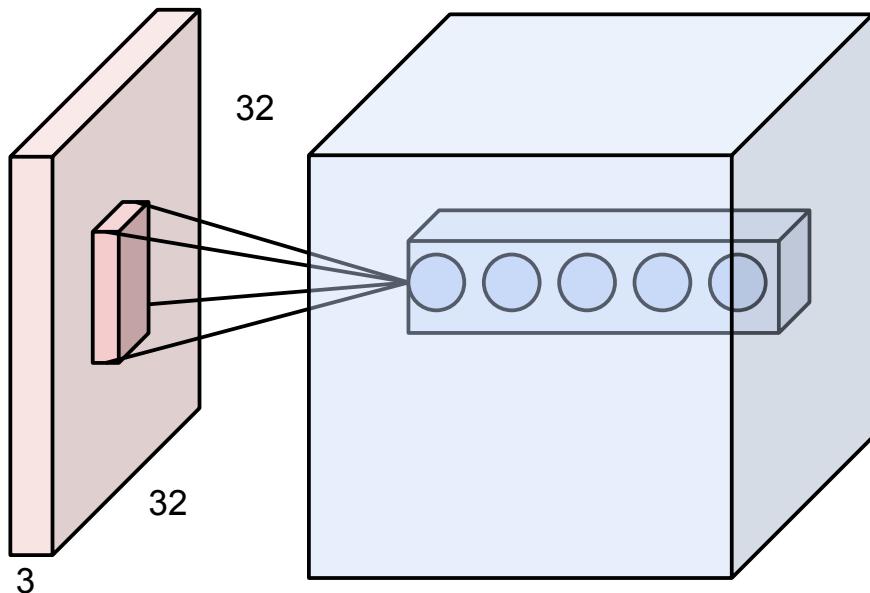


Multiple neurons all  
looking at the same  
region of the input  
volume, stacked  
along depth.

Source: [Fei-Fei Li & Andrej Karpathy]

# Kernels (Filters) Create Depth

Convolutional Neural Networks  
are just Neural Networks BUT:  
**1. Local connectivity**



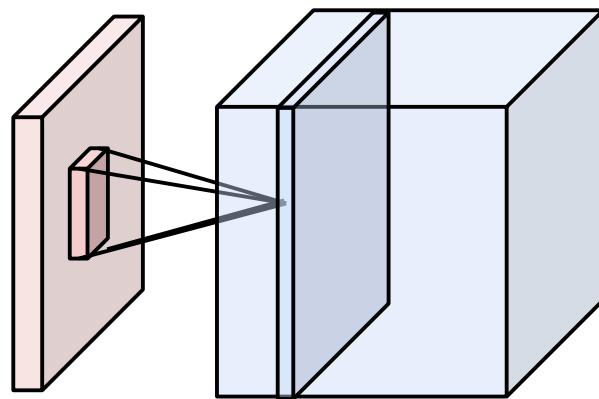
These form a single  
[ $1 \times 1 \times \text{depth}$ ]  
“depth column” in the  
output volume

Source: [Fei-Fei Li & Andrej Karpathy]

# Convolutional Layers

These layers are called **Convolutional Layers**

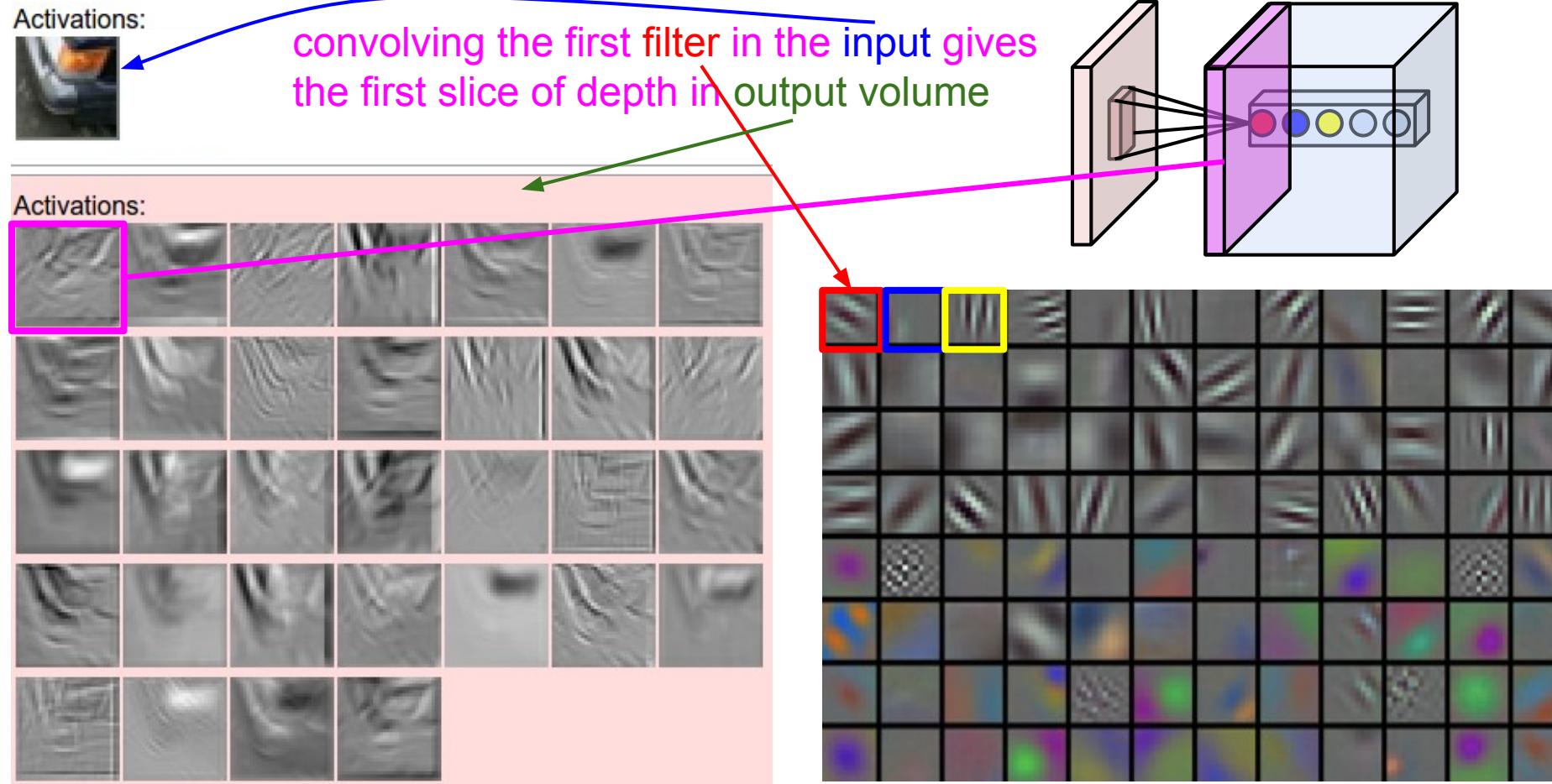
1. Connect neurons only to local receptive fields
2. Use the same neuron weight parameters for neurons in each “depth slice” (i.e. across spatial positions)



one activation map (a depth slice),  
computed with one set of weights

Source: [Fei-Fei Li & Andrej Karpathy]

# Filters, Convolutions, Depth



Input Volume (+pad 1) (7x7x3)							Filter W0 (3x3x3)		
							w0[:, :, 0]		
0	0	0	0	0	0	0	-1	0	1
0	0	1	2	0	0	0	-1	1	-1
0	2	1	2	1	1	0	-1	1	0
0	2	0	2	0	1	0	w0[:, :, 1]		
0	0	2	1	2	0	0	-1	1	-1
0	1	1	0	2	2	0	0	-1	1
0	0	0	0	0	0	0	1	0	0
x[:, :, 1]							w0[:, :, 2]		
0	0	0	0	0	0	0	1	1	-1
0	1	2	2	2	1	0	0	1	1
0	0	1	2	2	0	0	1	-1	-1
0	2	1	2	1	1	0	Bias b0 (1x1x1)		
0	0	0	0	1	1	0	b0[:, :, 0]	1	
0	2	1	2	0	0	0	Bias b1 (1x1x1)		
0	0	0	0	0	0	0	b1[:, :, 0]	0	
x[:, :, 2]							toggle movement		
0	0	0	0	0	0	0			
0	0	1	1	2	0	0			
0	2	0	2	1	0	0			
0	2	2	2	2	1	0			
0	0	1	2	0	2	0			
0	0	0	10	16	18	0			

Filter W1 (3x3x3)			Output Volume (3x3x2)		
w1[:, :, 0]			o[:, :, 0]		
-1	0	1	2	4	3
1	0	-1	4	0	-1
-1	0	-1	3	-1	4
w1[:, :, 1]			o[:, :, 1]		
1	-1	0	1	-6	-7
-1	-1	0	-5	-11	-5
-1	1	0	2	-2	1
w1[:, :, 2]					
1	1	0			
0	-1	1			
-1	-1	-1			

- Input 7x7 with 3 Channels
- Kernels 3x3 with 3 Channels
- 2 Kernels (Filters) W0 and W1
- Output has 2 Channels (# of Kernels)

$$(1 \times -1 + 2 \times 1) +$$

$$(1 \times -1 + 2 \times 1) +$$

$$(1 \times 1 + 2 \times -1) +$$

$$1$$

$$= 2$$

Input Volume (+pad 1) (7x7x3)							Filter W0 (3x3x3)		
							w0[:, :, 0]		
0	0	0	0	0	0	0	-1	0	1
0	0	1	2	0	0	0	-1	1	-1
0	2	1	2	1	1	0	-1	1	0
0	2	0	2	0	1	0	w0[:, :, 1]	-1	1
0	0	2	1	2	0	0	-1	1	1
0	1	1	0	2	2	0	0	-1	1
0	0	0	0	0	0	0	w0[:, :, 2]	1	0
x[:, :, 1]	0	0	0	0	0	0	1	1	-1
0	1	2	2	2	1	0	0	1	1
0	0	1	2	2	0	0	1	-1	-1
0	2	1	2	1	1	0	1	1	-1
0	0	0	0	1	1	0	0	-1	1
x[:, :, 2]	0	0	0	0	0	0	0	-1	-1
0	0	0	0	0	0	0	Bias b0 (1x1x1)	b0[:, :, 0]	1
0	0	0	0	0	0	0			
10/16/18	0	0	0	0	0	0			

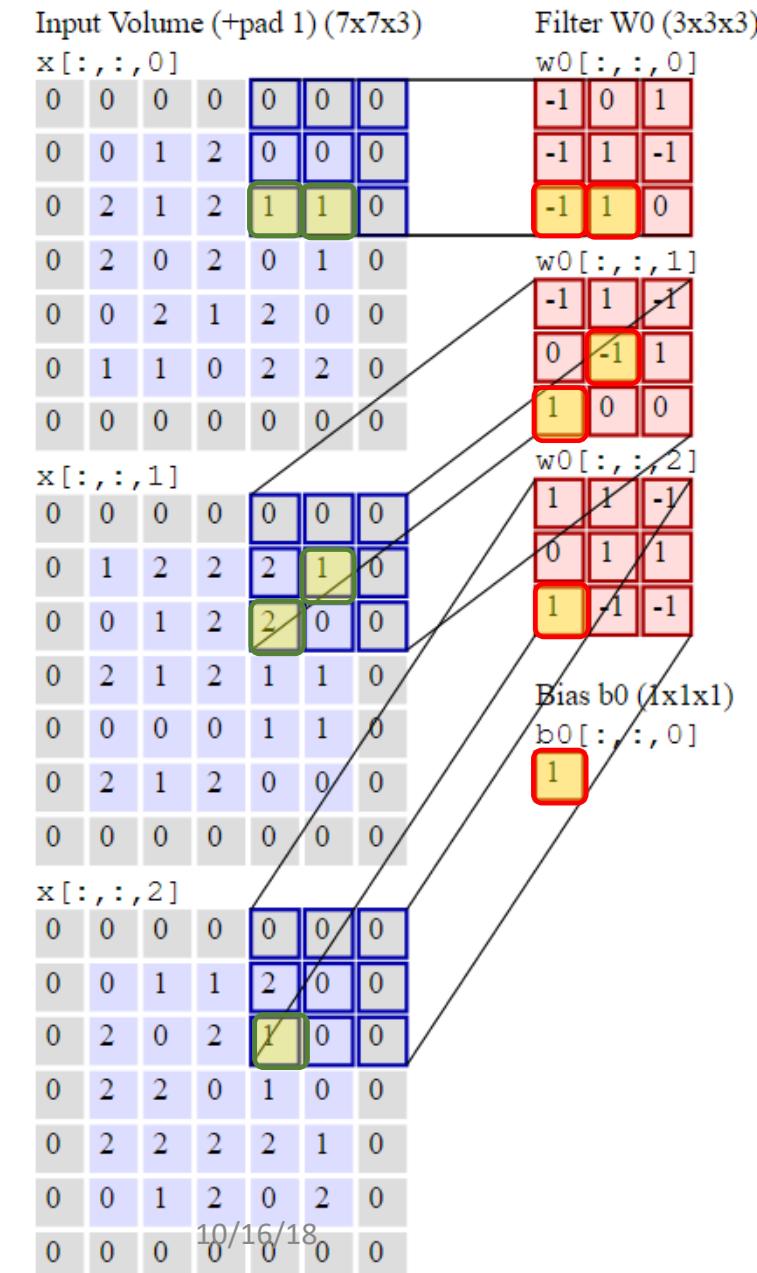
Filter W1 (3x3x3)			Output Volume (3x3x2)		
w1[:, :, 0]			o[:, :, 0]		
-1	0	1	2	4	3
1	0	-1	4	0	-1
-1	0	-1	3	-1	4
w1[:, :, 1]			o[:, :, 1]		
1	-1	0	-4	-6	-7
-1	-1	0	-5	-11	-5
-1	1	0	2	-2	1
w1[:, :, 2]					
1	1	0			
0	-1	1			
-1	-1	-1			
Bias b1 (1x1x1)					
b1[:, :, 0]					
0					

toggle movement

# Example

- Input 7x7 with 3 Channels
- Kernels 3x3 with 3 Channels
- 2 Kernels (Filters) W0 and W1
- Output has 2 Channels (# of Kernels)

$$\begin{aligned}
 & (1 \times -1 + 2 \times 1 + 1 \times -1 + 2 \times 1) + \\
 & (2 \times 1 + 2 \times 1 + 1 \times 1 + 2 \times -1 + 2 \times -1) + \\
 & (1 \times 1 + 2 \times 1 + 2 \times -1 + 1 \times -1) + \\
 & 1 \\
 & = 4
 \end{aligned}$$



Filter W0 (3x3x3)

w0[:, :, 1]	-1 1 1	0 -1 1	1 0 0
-1 1 1	0 0 0	0 0 0	0 0 0
0 -1 1	0 0 0	0 0 0	0 0 0

Filter W1 (3x3x3)

w1[:, :, 1]	1 -1 0	-4 -6 -7	-5 -11 -5
1 -1 0	0 0 0	0 0 0	0 0 0
-4 -6 -7	0 0 0	0 0 0	0 0 0

Output Volume (3x3x2)

o[:, :, 1]	1 -1 0	-4 -6 -7	-5 -11 -5
1 -1 0	0 0 0	0 0 0	0 0 0
-4 -6 -7	0 0 0	0 0 0	0 0 0

Bias b0 (1x1x1)

b0[:, :, 0]	0
-------------	---

Bias b1 (1x1x1)

b1[:, :, 0]	0
-------------	---

toggle movement

# Example

- Input 7x7 with 3 Channels
- Kernels 3x3 with 3 Channels
- 2 Kernels (Filters) W0 and W1
- Output has 2 Channels (# of Kernels)

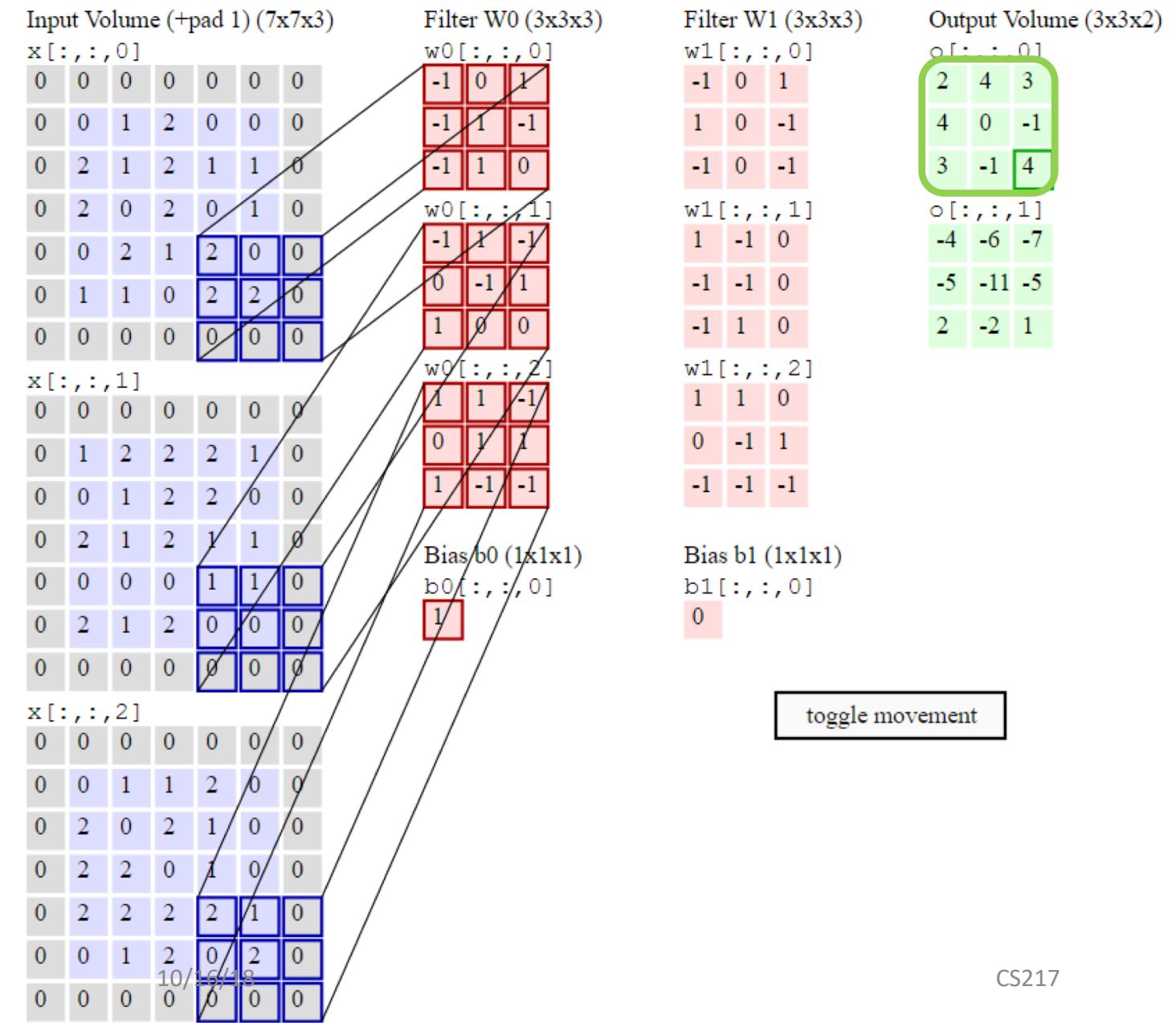
$$(1 \times 1 + 1 \times 1) +$$

$$(1 \times 1 + 2 \times 1) +$$

$$(1 \times 1) +$$

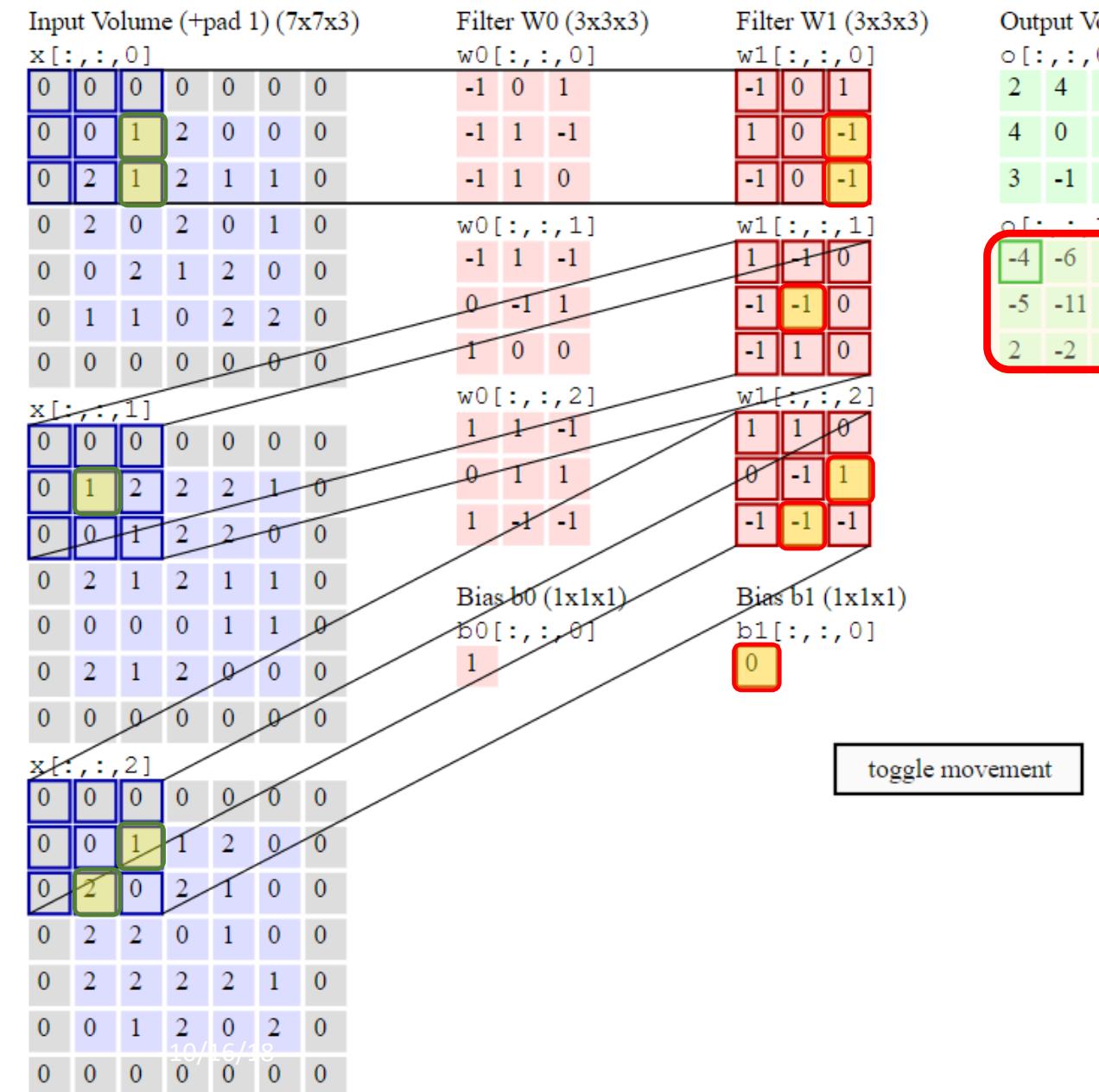
$$1$$

$$= 3$$



- Computed 1st Channel of Output
  - Use Filter W1 with the same input volume

CS217



- Input 7x7 with 3 Channels
- Kernels 3x3 with 3 Channels
- Now with W1
- Output's 2nd Channel

$$(1 \times -1 + 1 \times -1) +$$

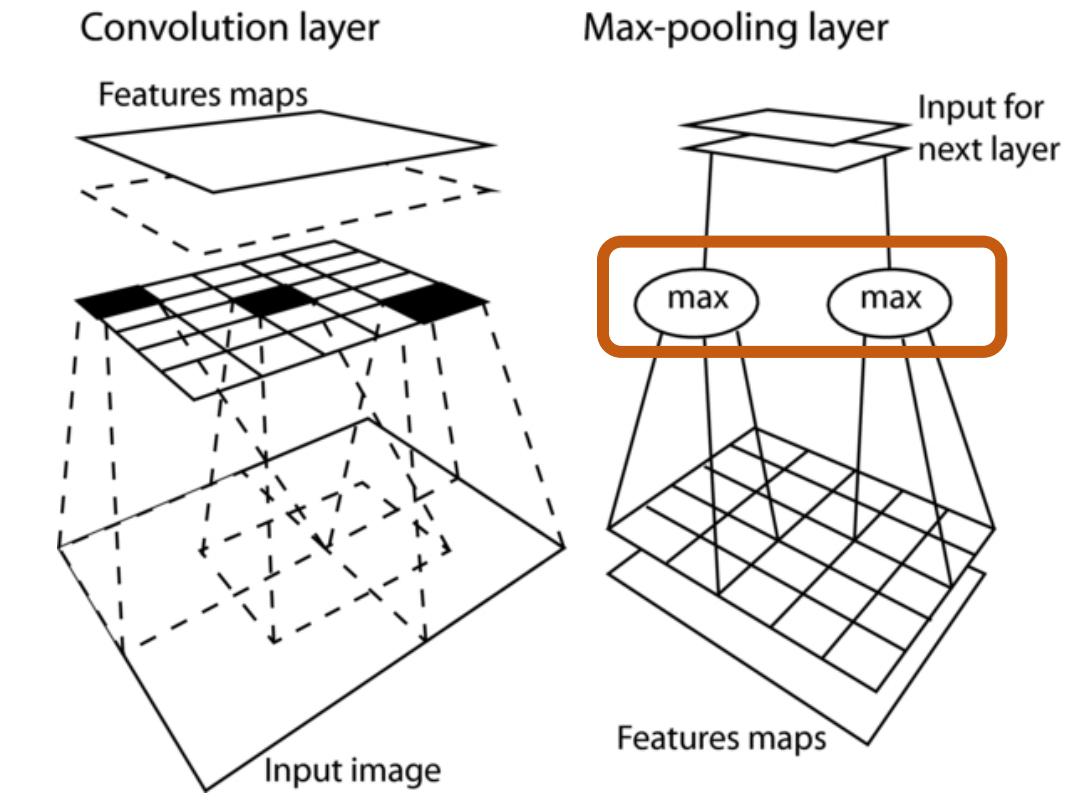
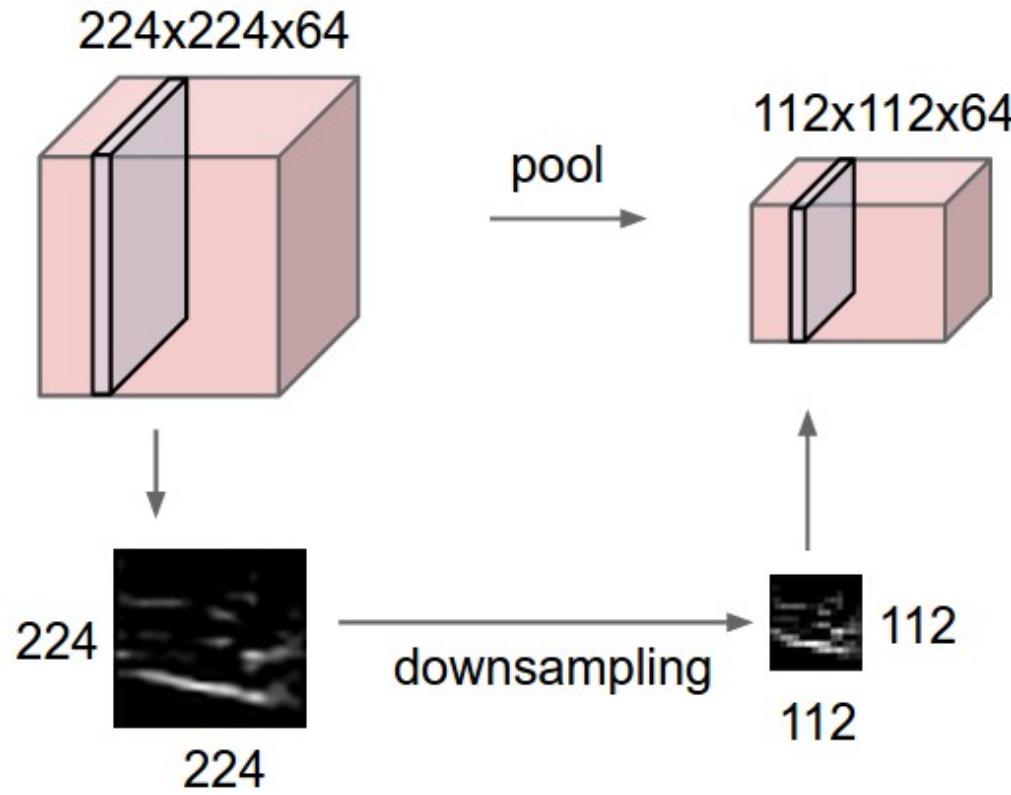
$$(1 \times -1) +$$

$$(1 \times 1 + 2 \times -1) +$$

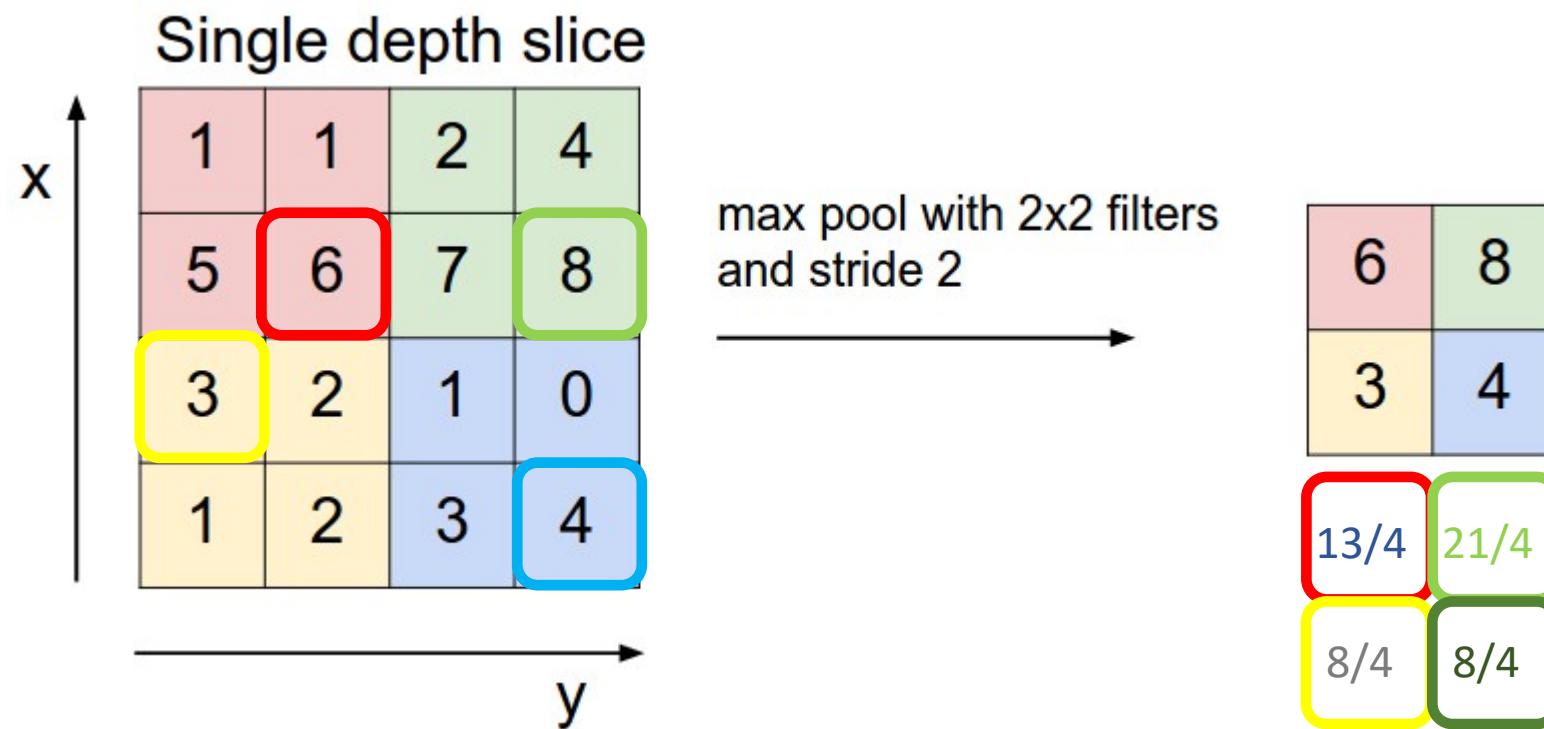
$$0$$

$$= -4$$

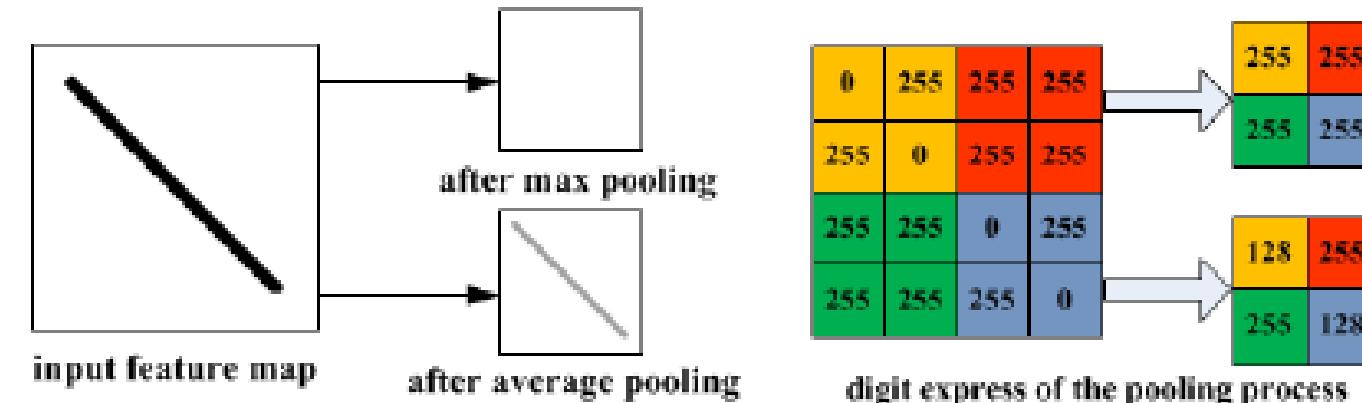
# Pooling Layer



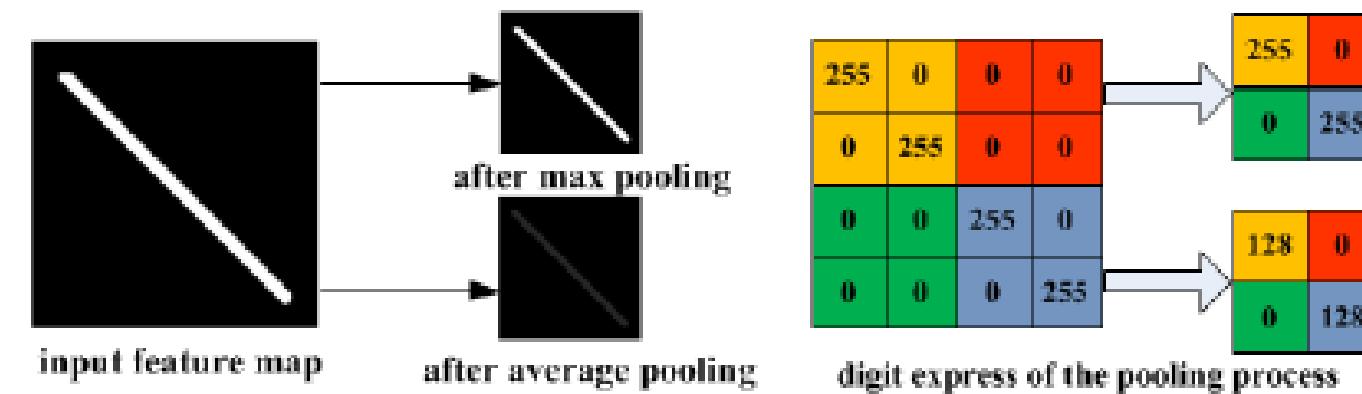
# Max Pooling



# Pooling



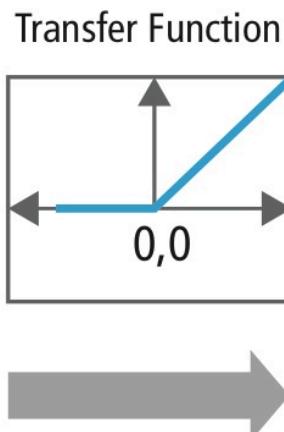
(a) Illustration of max pooling drawback



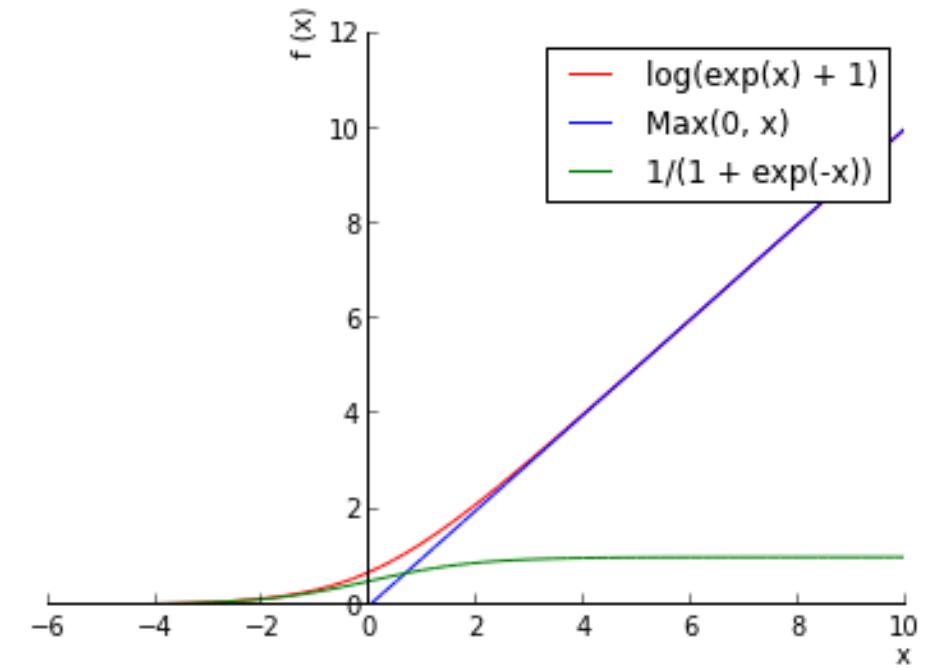
(b) Illustration of average pooling drawback

# Rectified Linear Units

15	20	-10	35
18	-110	25	100
20	-15	25	-10
101	75	18	23



15	20	0	35
18	0	25	100
20	0	25	0
101	75	18	23



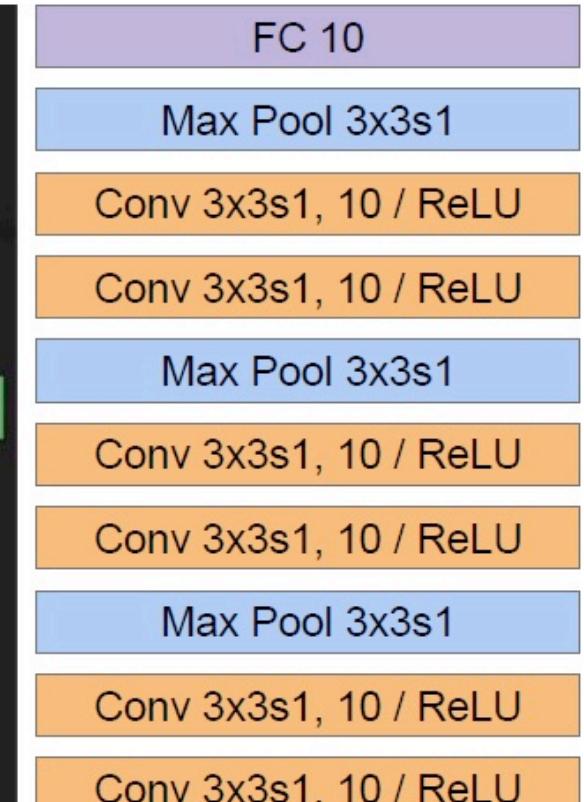
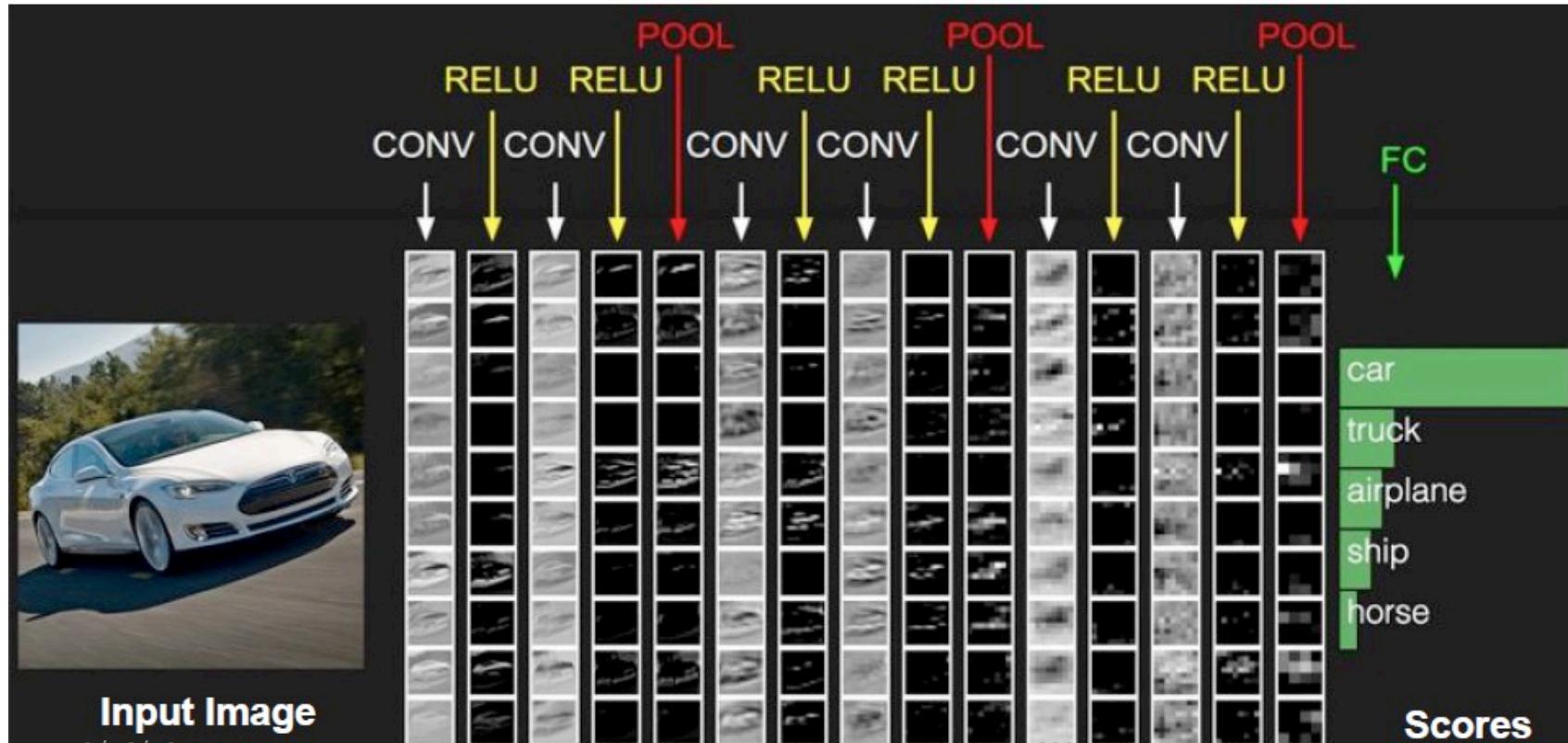
EASY.....

RIGHT?



Conv 3x3s1, 10 / ReLU

Type: Conv Kernel Size: 3x3 Stride: 1 Channels: 10 Activation: ReLU



# DEMO

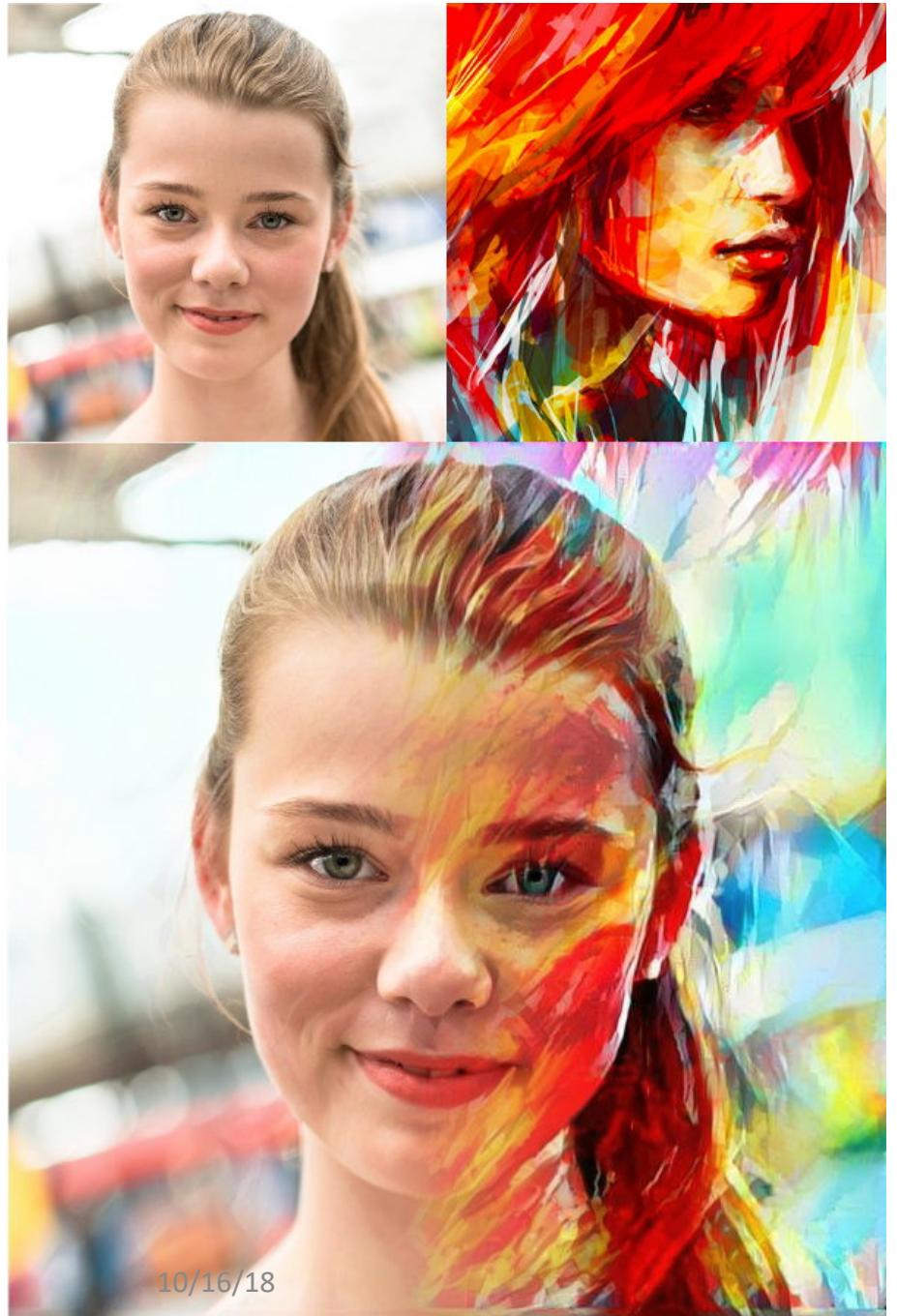
Classify MNIST digits with a  
Convolutional Neural Network



Classify CIFAR-10 with  
Convolutional Neural Network

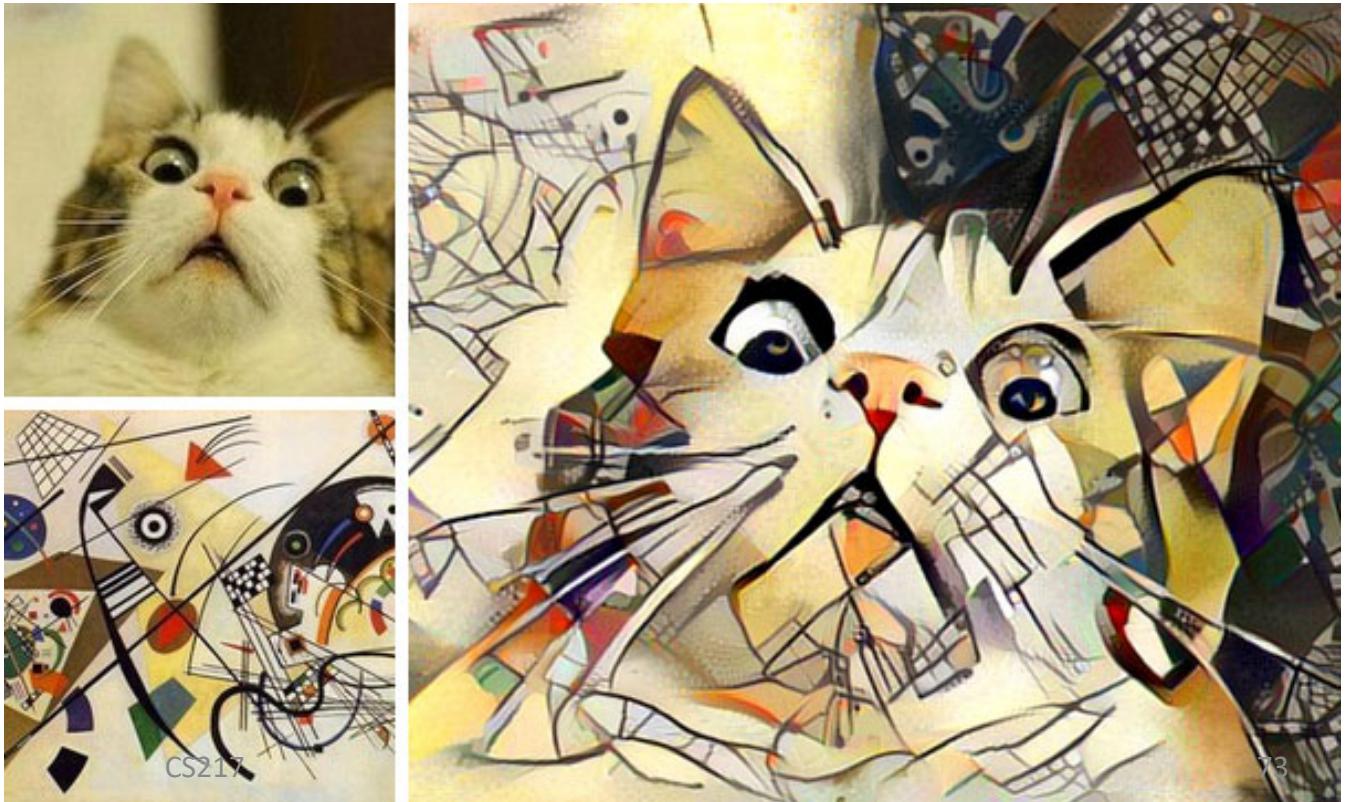


- <http://cs.stanford.edu/people/karpathy/convnetjs/>



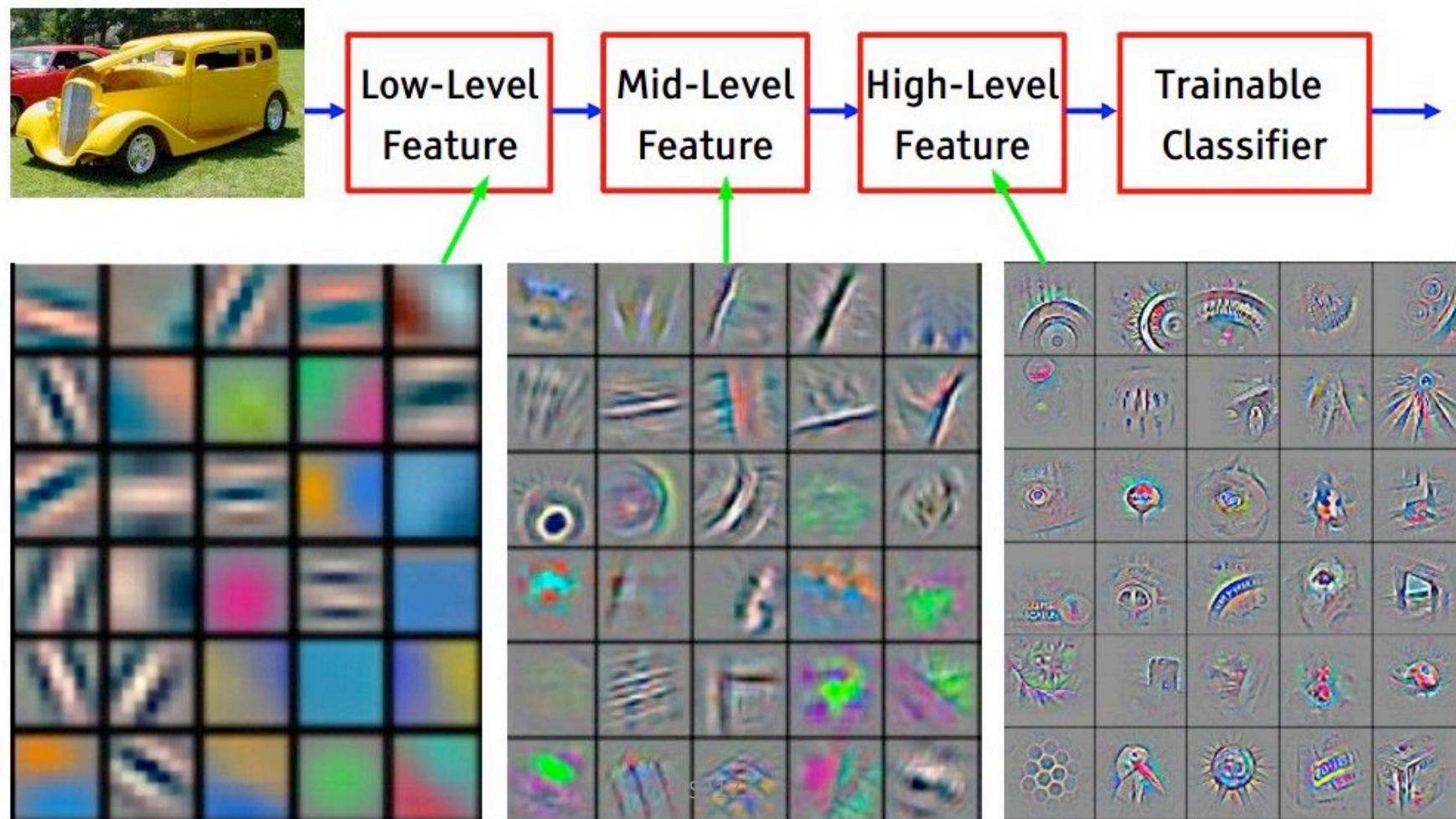
10/16/18

# Next: Visualizing CNNs and more



73

# Features and Layers



# Deep Visualization

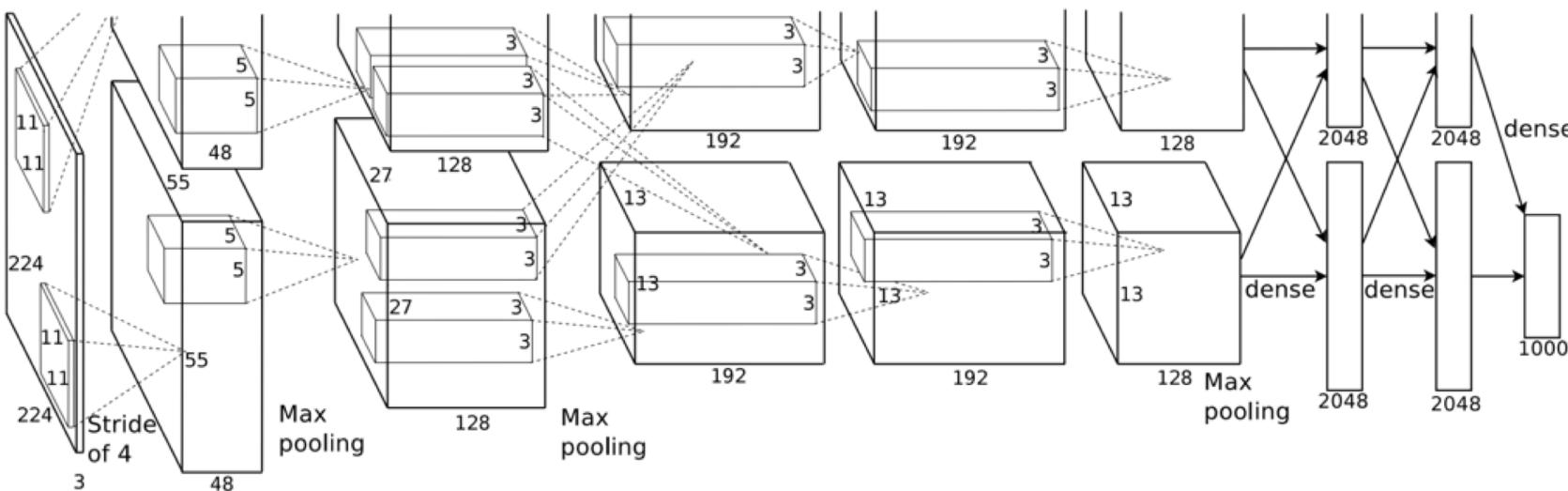




- Watch the video for a nice recap of CNNs and what we learnt so far

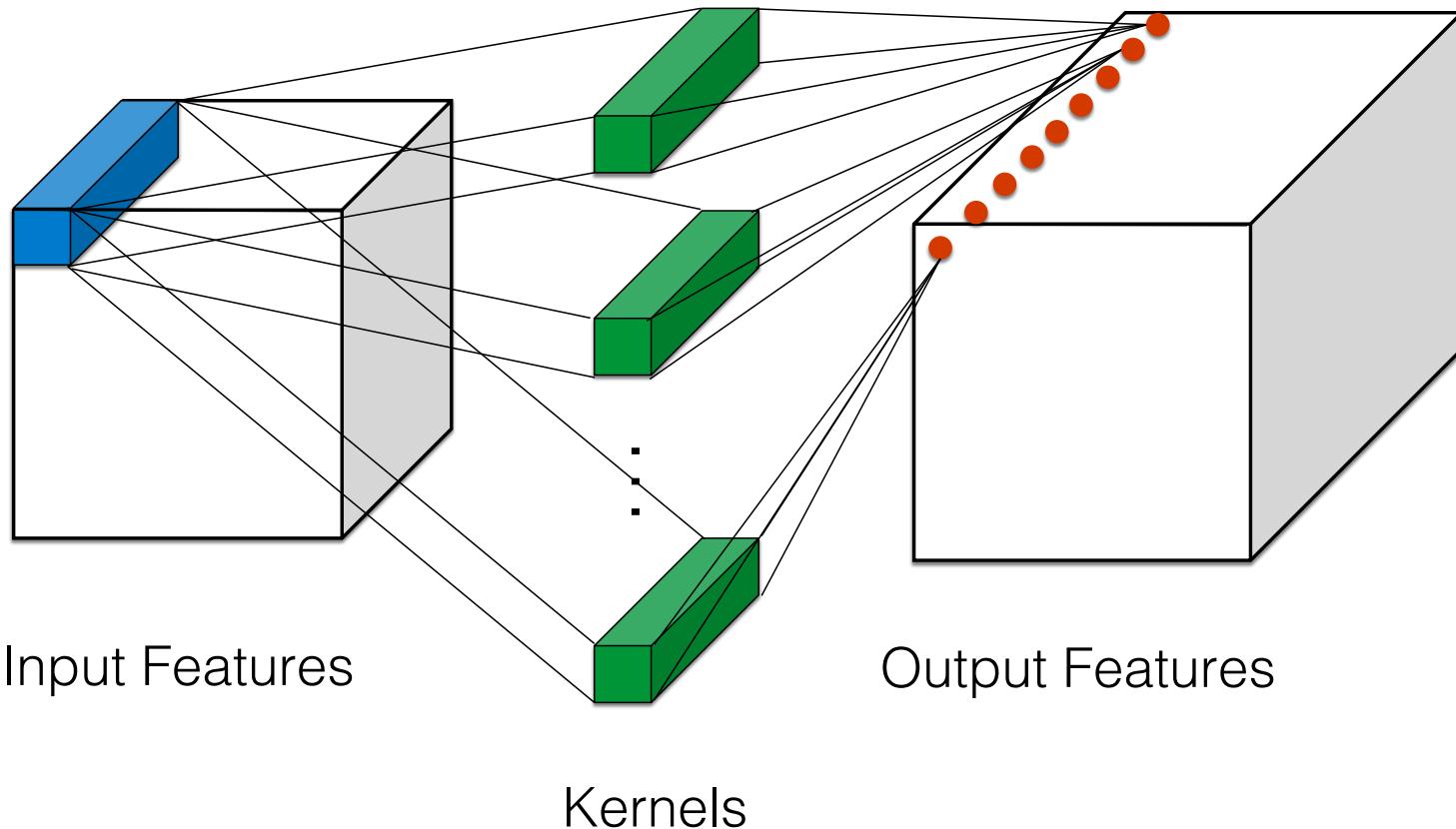
- [https://archive.org/details/Redwood\\_Center\\_2016\\_03\\_01\\_Leon\\_Gatys](https://archive.org/details/Redwood_Center_2016_03_01_Leon_Gatys)

- Four dimensions in each layer
  - Image Dimensions (X,Y)
  - # Channels (C)
  - # Kernels (K)
  - Filter dimensions ( $F_h$ ,  $F_w$ )



Alex Krizhevsky, Ilya Sutskever, Geoffrey E. Hinton, ImageNet Classification with Deep Convolutional Neural Networks

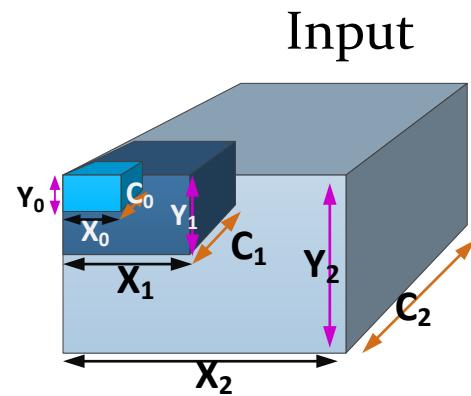
# Convolution Layer



# Optimize CNNs for Hardware

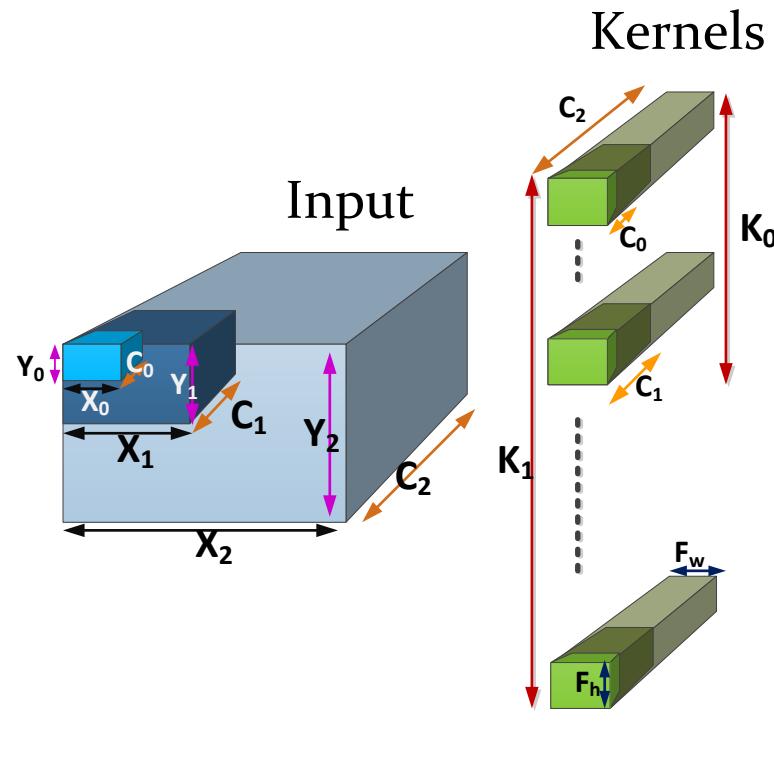
- Turn it to GEMM/GEMV?
- Design Accelerators
- Take advantage of Locality
- A desired Datapath
  - Data Movement
  - PEs
  - Systolic/Non Systolic
- Take advantage of sparsity
- Change Nature of Computation (Winograd)

# Multilevel Blocking of Convolutional Layer



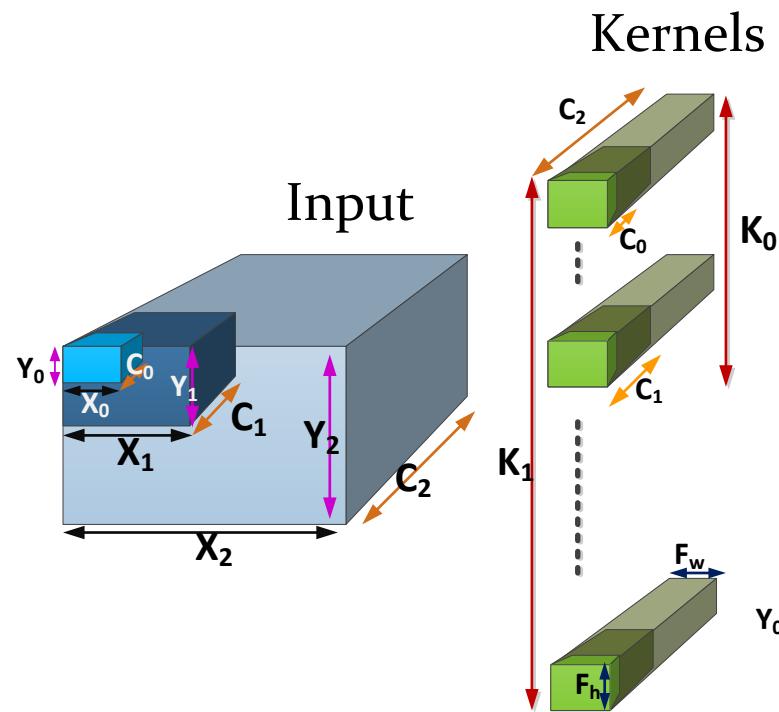
- Input Image Dimensions:
  - Image Dimensions ( $X, Y$ )
  - # Channels ( $C$ )
  - # Kernels ( $K$ )
  - Filter dimensions ( $F_h, F_w$ )
- Nested blocking

# Multilevel Blocking of Convolutional Layer



- Kernel Dimensions:
  - Image Dimensions ( $X, Y$ )
  - # Channels ( $C$ )
  - # Kernels ( $K$ )
  - Filter dimensions ( $F_h, F_w$ )
- Nested blocking

# Multilevel Blocking of Convolutional Layer



- Output Dimensions:
  - Image Dimensions ( $X, Y$ )
  - # Channels (C)
  - **# Kernels (K)**

- Nested blocking

# Multilevel Blocking of Convolutional Layer

