

COM SCI 188

# Intro to Robotics

## Lecture 7

Yuchen Cui  
Winter 2026

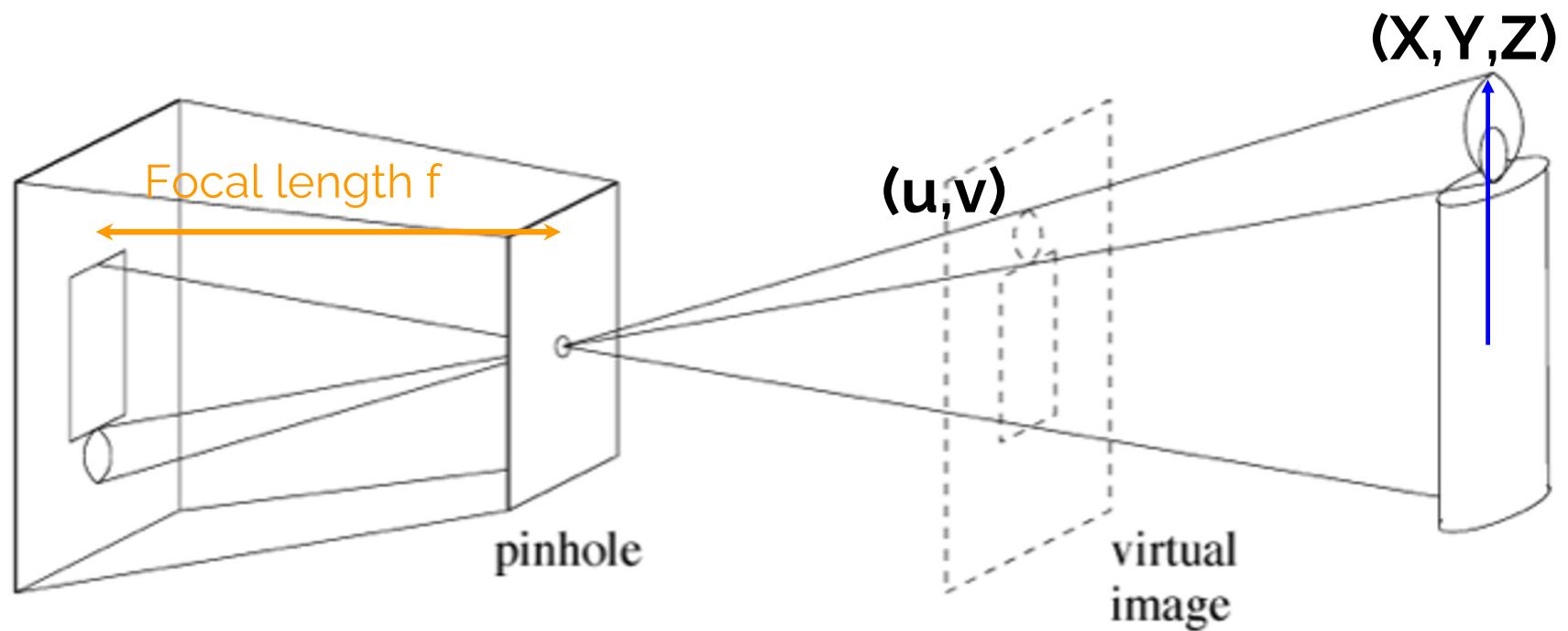
# Agenda

- Announcements
- Camera Model & Camera calibration
- Computer Vision Introduction
- Convolutional Neural Networks

# Announcements

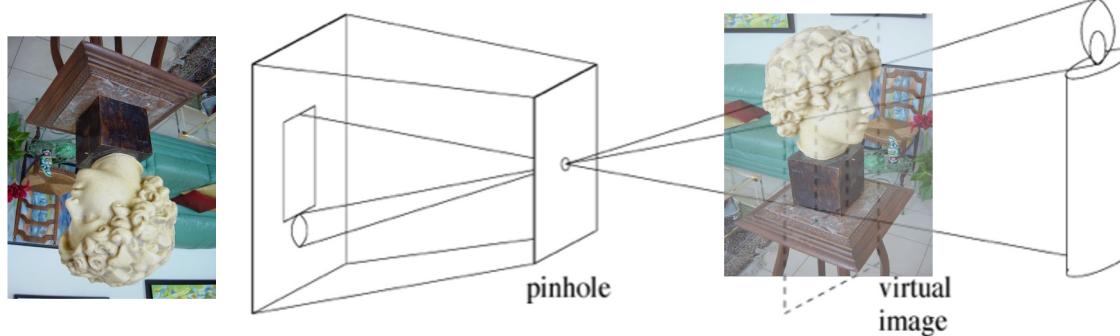
- Problem Set 2 due Friday
- Coding Assignment 2 out, due next Friday
- Midterm Exam (Week 7): Feb 19<sup>th</sup>

# Recap: pinhole camera model



# Recap: Camera Projection

**Image Unit:** Pixels



**World Unit:**  
Meters



If we let  $f$  to take care transform from word unit to image unit.

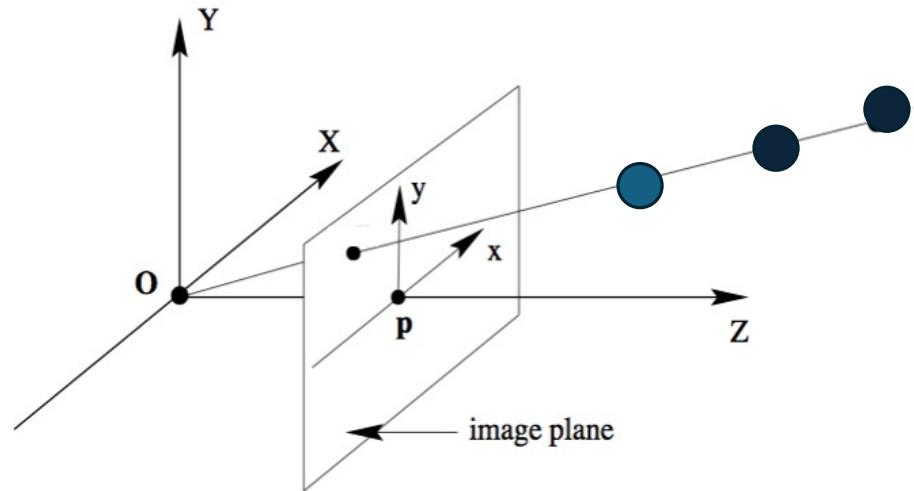
The unit of  $f$  is pixel

$$u = \frac{fX}{Z} \quad v = \frac{fY}{Z}$$

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \end{bmatrix}$$

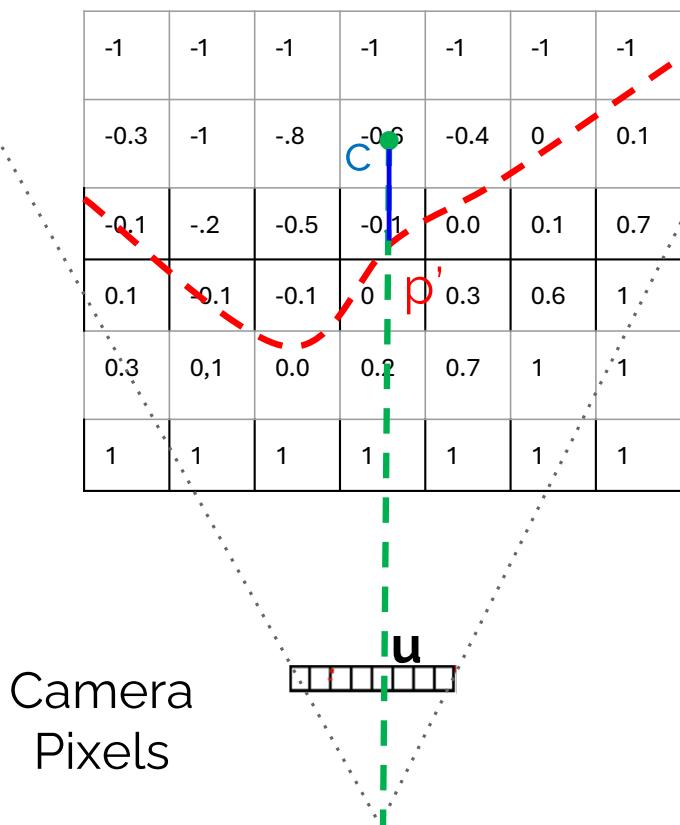
## Recap:

$$\lambda \begin{bmatrix} u \\ v \\ f \end{bmatrix} = \begin{bmatrix} X \\ Y \\ Z \end{bmatrix}$$



All the points on the same ray, projected to the same pixels [u,v]  
Only the point with smallest distance is observed in the image.

# Projective Sign Distance Function



Sign indicated occluded/free space with respect to camera

In this example:

Red: surface observed in camera

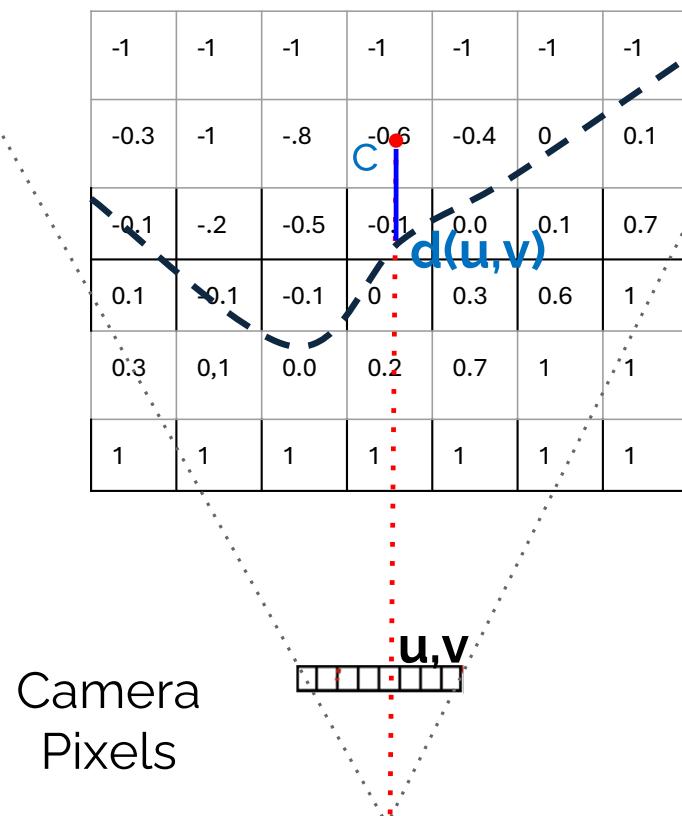
Green: camera projection ray

We take  $p'$  as the nearest point.

Sign:  $c$  is behind  $p'$   $\rightarrow$  occluded  $\rightarrow$  negative

Projective distance: Blue line

# Computing projective distance

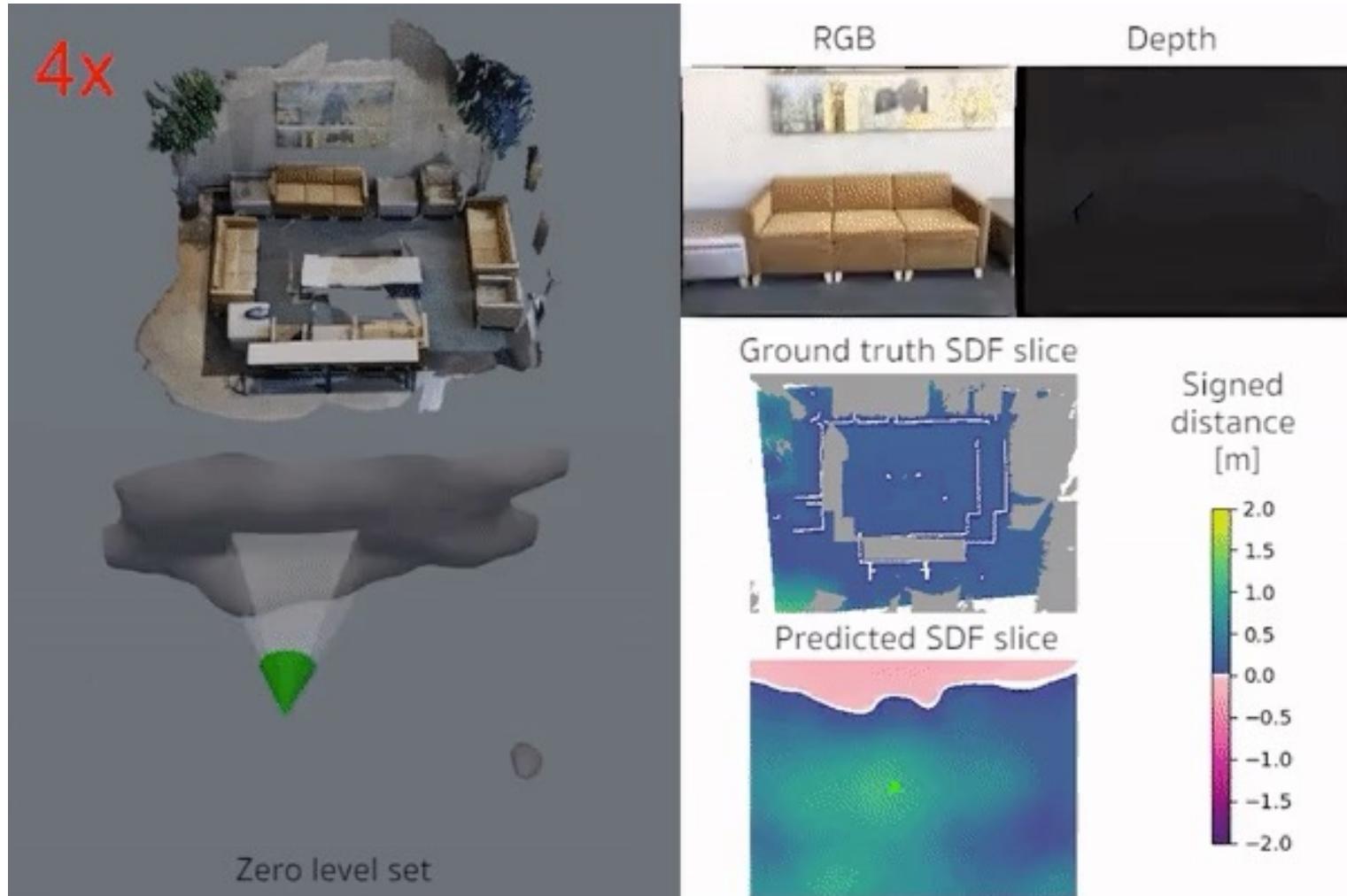


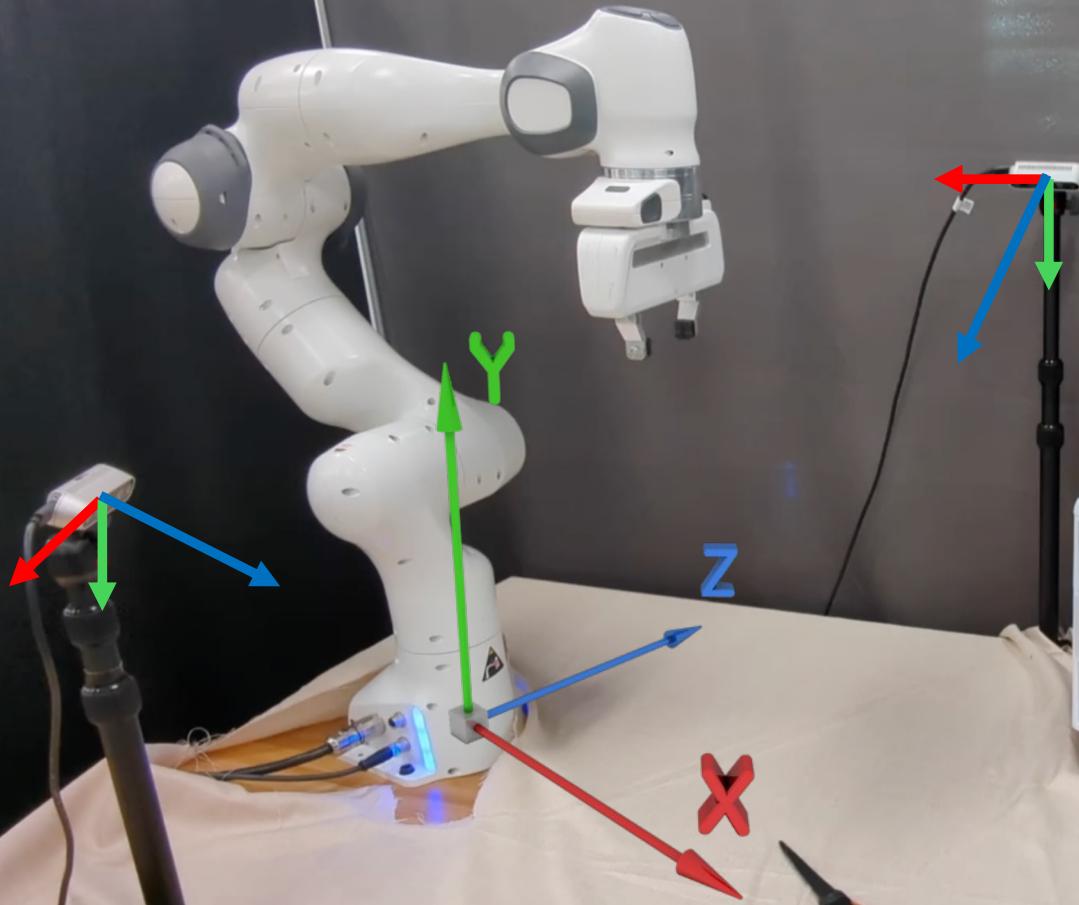
- 3D point location  $c(x, y, z)$
- Project  $c(x, y, z)$  to the 2D image, gives us 2D coordinate  $u, v$
- Depth reading on pixel is  $d(u, v)$

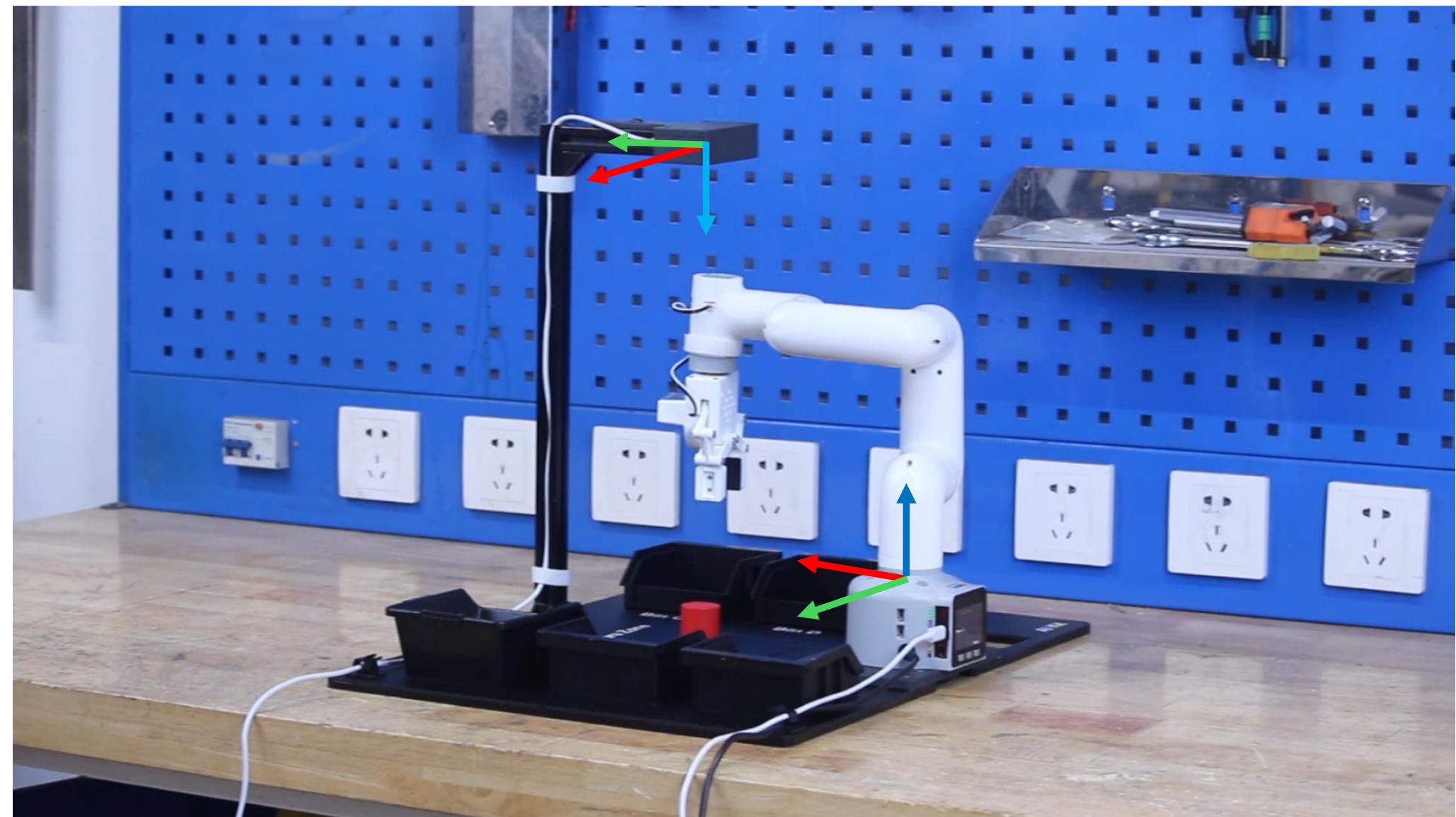
Signed projective distance

$$d_{\text{proj}} = d(u, v) - z$$

# iSDF: Real-Time Neural Signed Distance Fields for Robot Perception







# Recap: Camera Projection

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \mathbf{K} \text{ Intrinsic Matrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{array}{c|c} I & t \\ \hline \mathbf{0} & 1 \end{array} \times \begin{array}{c|c} R & \mathbf{0} \\ \hline \mathbf{0} & 1 \end{array} \text{ rotation} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

translation

Camera Extrinsic Matrix [R|t]

Camera Parameter  
Camera Projection Matrix

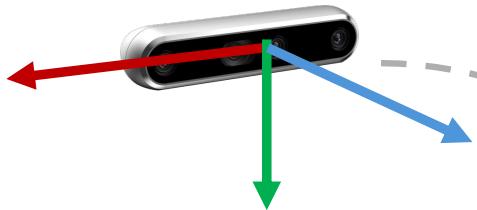
$$\mathbf{x}^P = \mathbf{P} \mathbf{X}$$

# Recap: Camera Projection

$$\mathbf{x} = \mathbf{K} \begin{bmatrix} \mathbf{R} | \mathbf{t} \end{bmatrix} \mathbf{X}$$

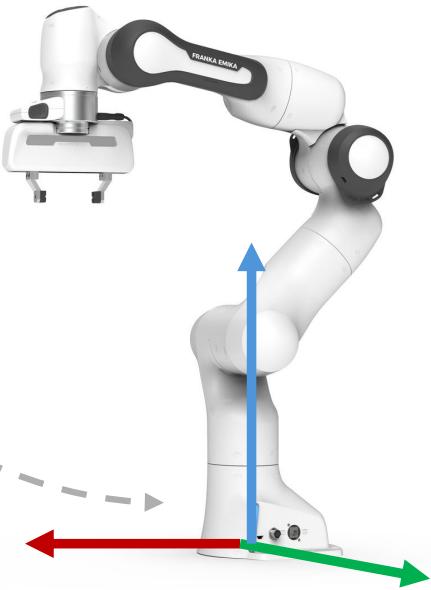
- Map a 3D point  $\mathbf{X}$  into a 2D coordinate in image  $\mathbf{x}$
- How to describe its *pose* in the world? (extrinsic matrix)
- How to describe its internal parameters? (intrinsic matrix)

# Camera Calibration



**Camera/Image frame:**  
reference frame for *visual observations*

**Robot base frame:**  
reference frame for *actions*



**Hand-Eye Calibration**

# Camera Calibration

**Goal:** express pixels in robot frame!

**Option 1:** solve for camera parameters

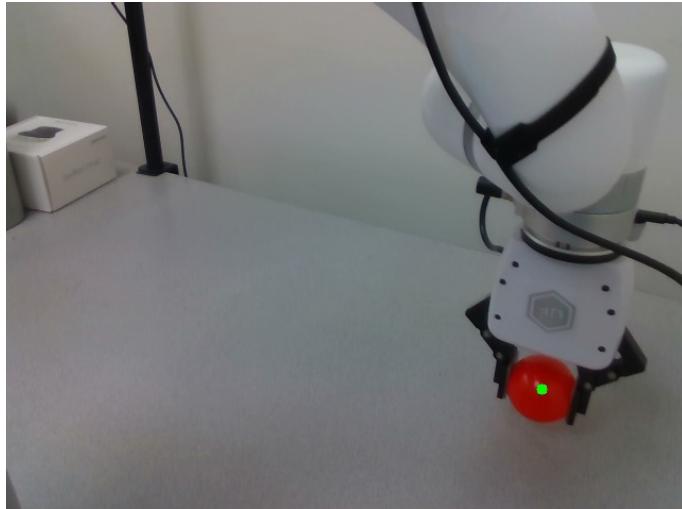
- intrinsic (focal length, principle point, pixel size)
- extrinsic (rotation angles, translation)

**Option 2:** solve for projection matrix

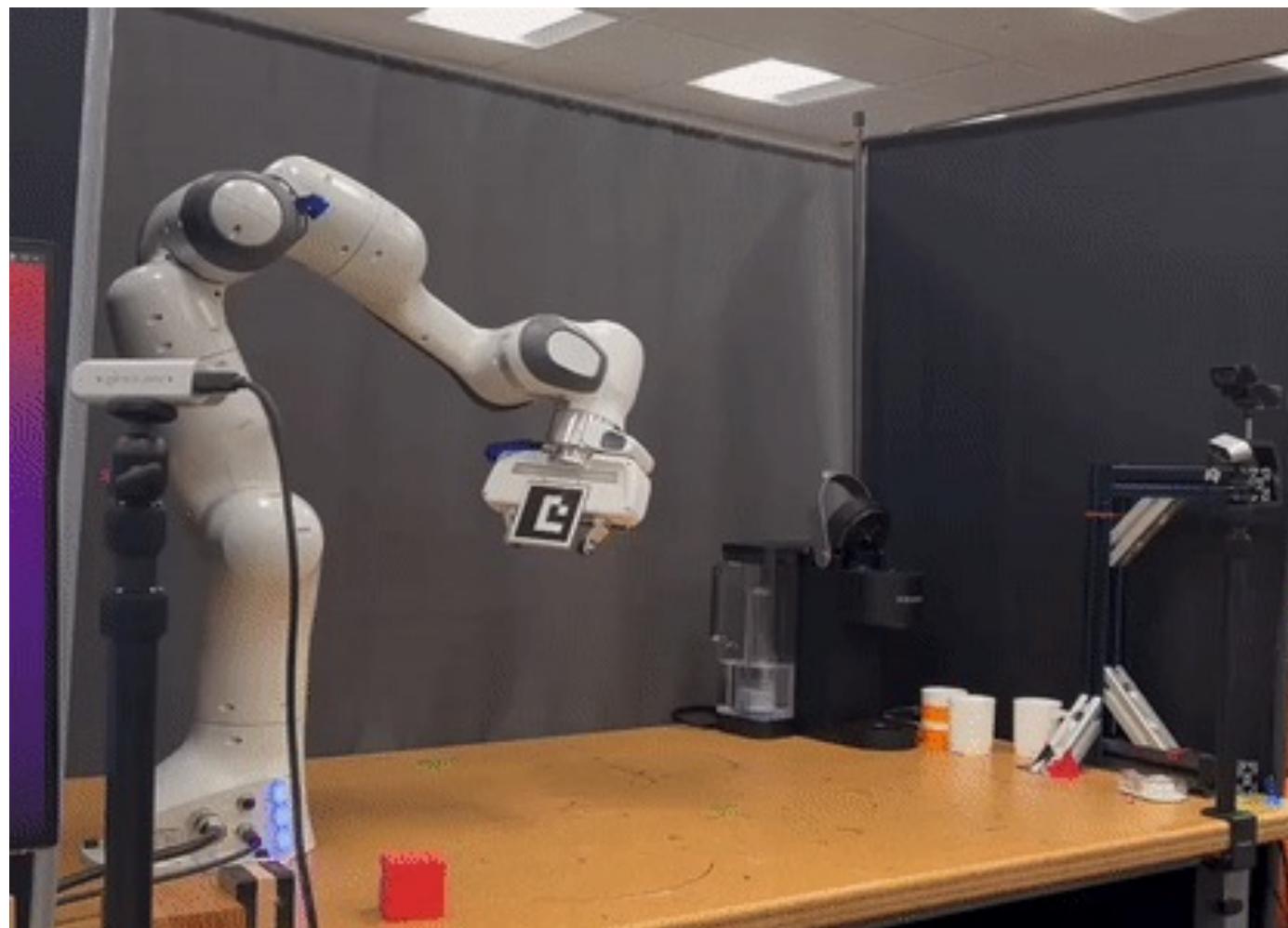
$$\mathbf{X} = \begin{bmatrix} wx \\ wy \\ w \end{bmatrix} = \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} = \mathbf{P}\mathbf{X}$$

**Robot:** move\_to(x=0.3, y=-0.2, z = 0.1)

**Camera:**



During the hand-eye calibration process, we move robot to a known position in its base frame, and captures the image with camera and detects the robot's hand's position (approximated with the red ball)



Video credit: Priya Sundaresan

# Problem Setup: solve for transformation

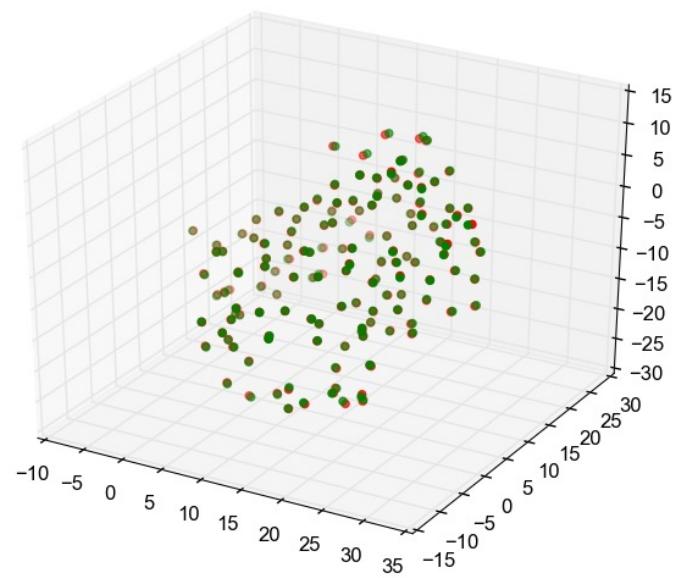
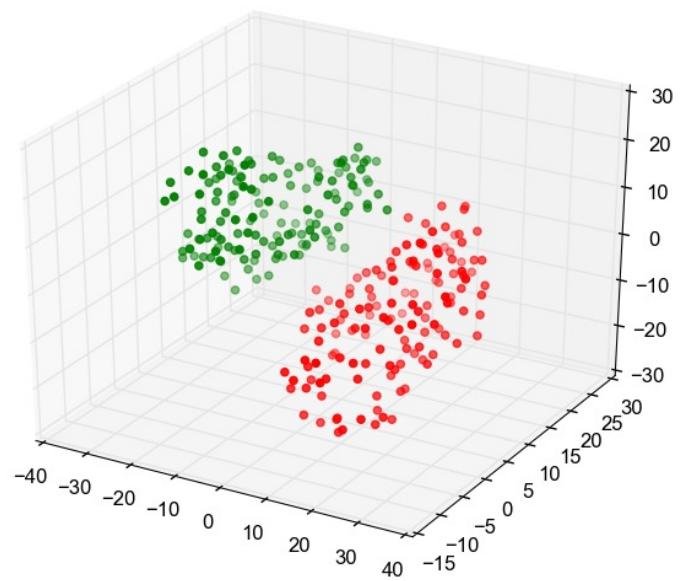
You're given:

- A set of points in frame A:  $\{\mathbf{p}_i^A\}$
- Corresponding points in frame B:  $\{\mathbf{p}_i^B\}$

You want to find:

- A **rotation matrix**  $R \in SO(3)$
- A **translation vector**  $t \in \mathbb{R}^3$   
such that:

$$\mathbf{p}_i^B \approx R\mathbf{p}_i^A + t$$



# The Kabsch-Umeyama Algorithm

- 1. Compute Centroids**
- 2. Center the Points**
- 3. Compute Cross-Covariance Matrix**
- 4. Apply SVD (Singular Value Decomposition)**
- 5. Construct the Rotation Matrix**
- 6. Compute the Translation Vector**

## **Output the Transformation**

Combine rotation and translation into a rigid body transformation

# The Kabsch-Umeyama Algorithm

**1. Compute centroids:**

$$\bar{\mathbf{p}}^A = \frac{1}{N} \sum_{i=1}^N \mathbf{p}_i^A, \quad \bar{\mathbf{p}}^B = \frac{1}{N} \sum_{i=1}^N \mathbf{p}_i^B$$

---

**2. Center the points:**

$$\mathbf{q}_i^A = \mathbf{p}_i^A - \bar{\mathbf{p}}^A, \quad \mathbf{q}_i^B = \mathbf{p}_i^B - \bar{\mathbf{p}}^B$$

---

# The Kabsch-Umeyama Algorithm

**3. Compute covariance matrix:**

$$H = \sum_{i=1}^N \mathbf{q}_i^A (\mathbf{q}_i^B)^\top$$

---

**4. Singular Value Decomposition (SVD):**

$$H = U \Sigma V^\top$$

# The Kabsch-Umeyama Algorithm

**5. Compute rotation matrix  $R$ :**

$$R = VU^\top$$

If  $\det(R) < 0$ , correct for reflection:

$$\begin{aligned} V[:, 3] &\leftarrow -V[:, 3] \quad (\text{flip the third column of } V) \\ R &= VU^\top \end{aligned}$$

---

**6. Compute translation vector  $t$ :**

$$\mathbf{t} = \bar{\mathbf{p}}^B - R\bar{\mathbf{p}}^A$$

---

# Computer Vision Basics

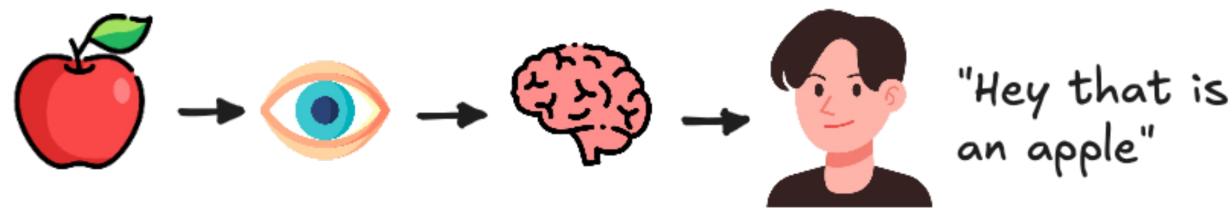
## Image Processing + CNN Survival Kit

Slides Adapted from Stanford EE227/CS227A: Robot Perception

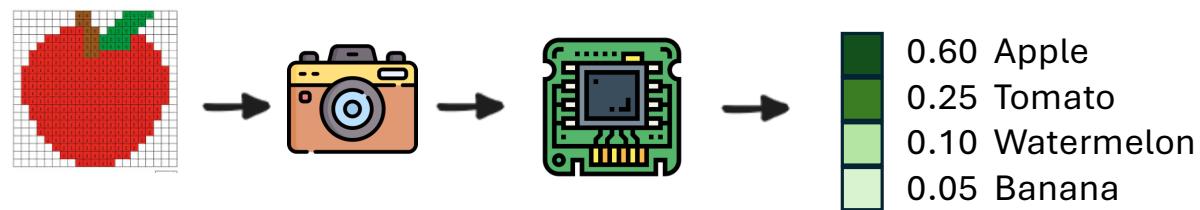
credit: Shuran Song

What is computer vision?

## Human Vision



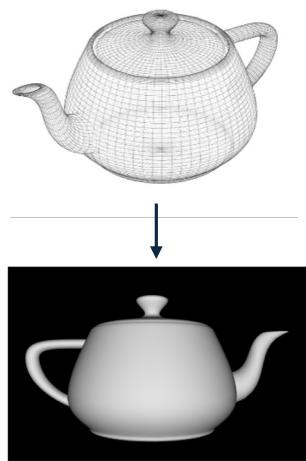
## Computer Vision



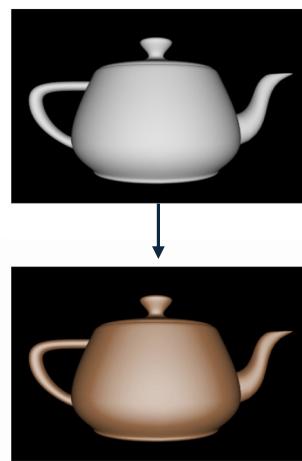
# Computer Vision: Quick Overview

What is computer vision?

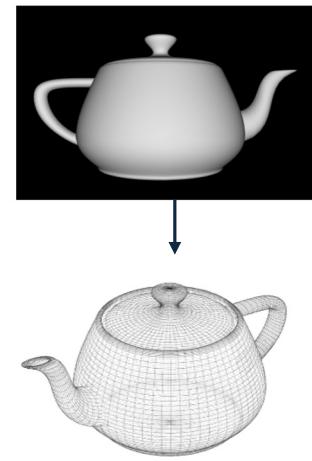
---



**Computer Graphics**  
Models to Images



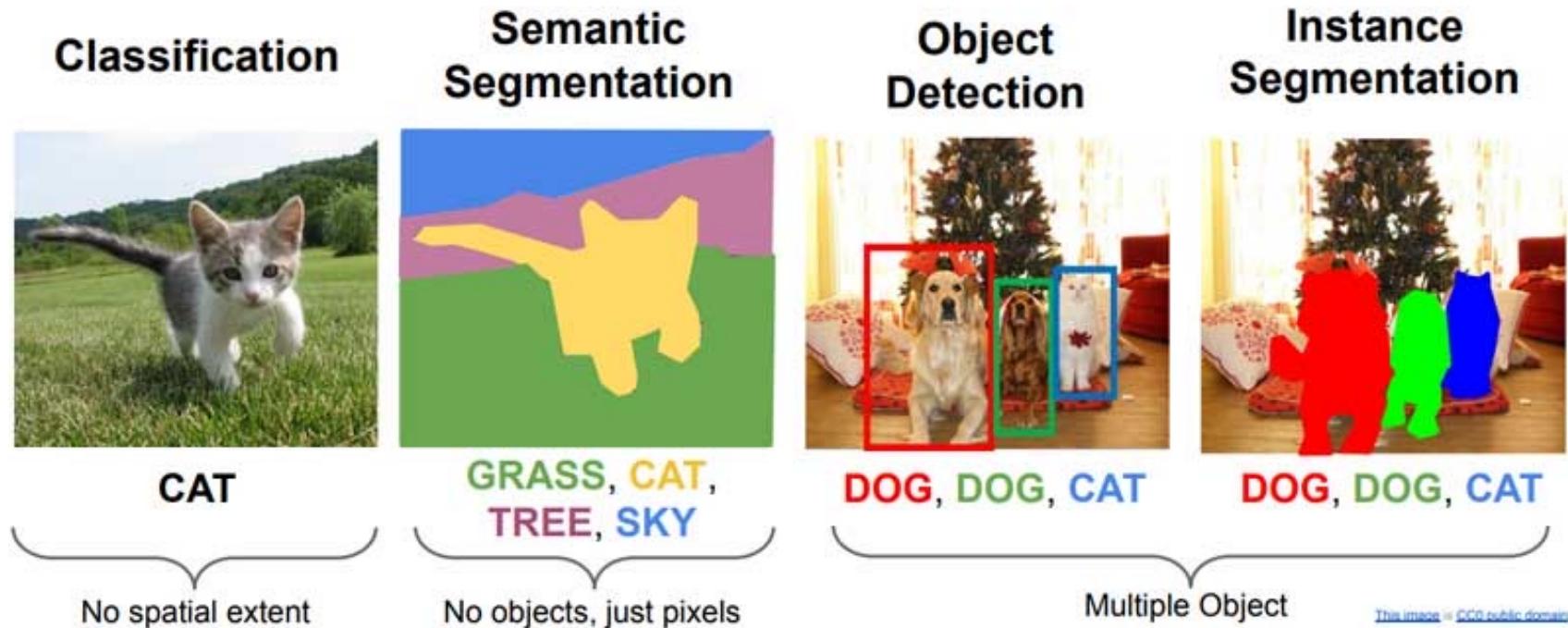
**Comp. Photography**  
Images to Images



**Computer Vision**  
Images to Models

# Computer Vision: Quick Overview

Common 2D computer vision tasks



# What Makes 2D Computer Vision Hard?



**Variation:** same cat different poses, view points ...

# What Makes 2D Computer Vision Hard?



**Variation:** different cat with different shape, color, texture...

# What Makes 2D Computer Vision Hard?



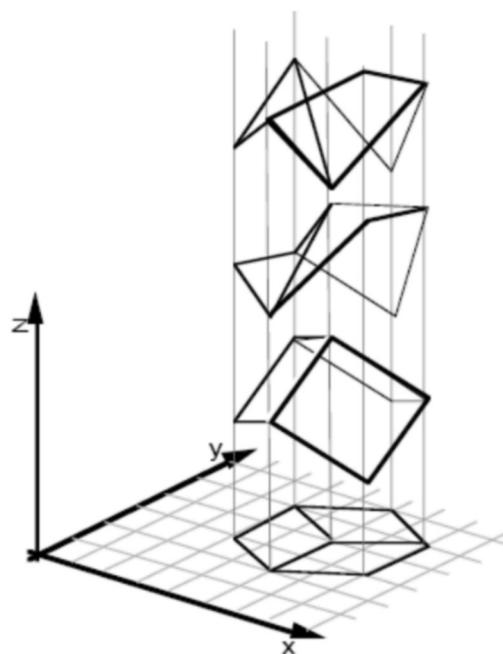
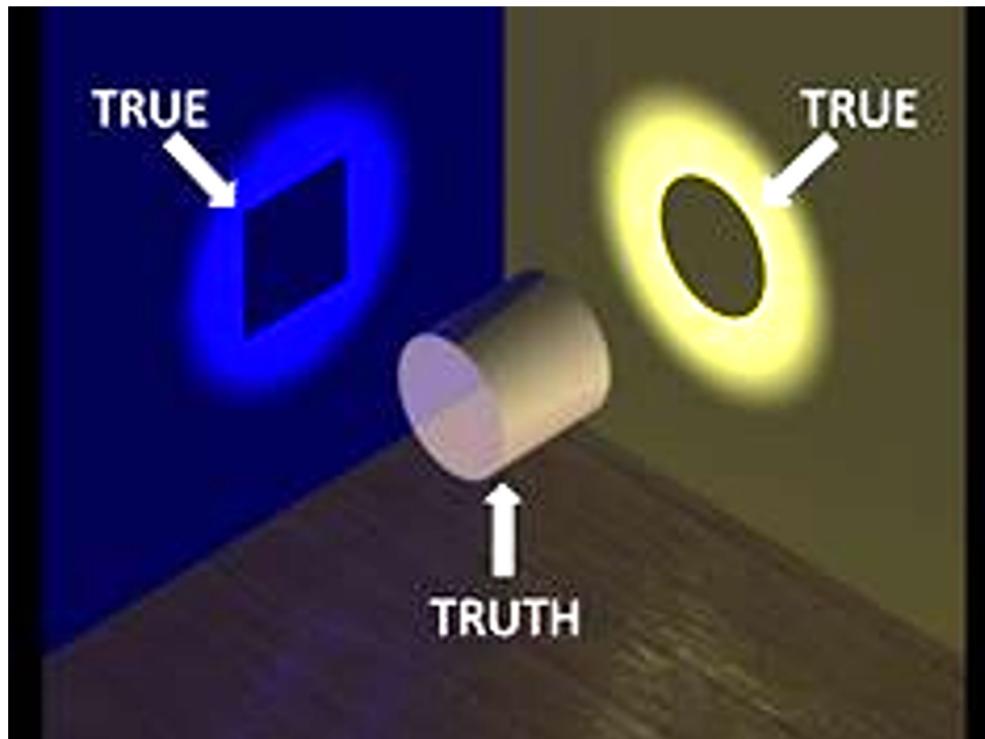
**Illumination**

# What Makes 2D Computer Vision Hard?



**Occlusion:** partial observation

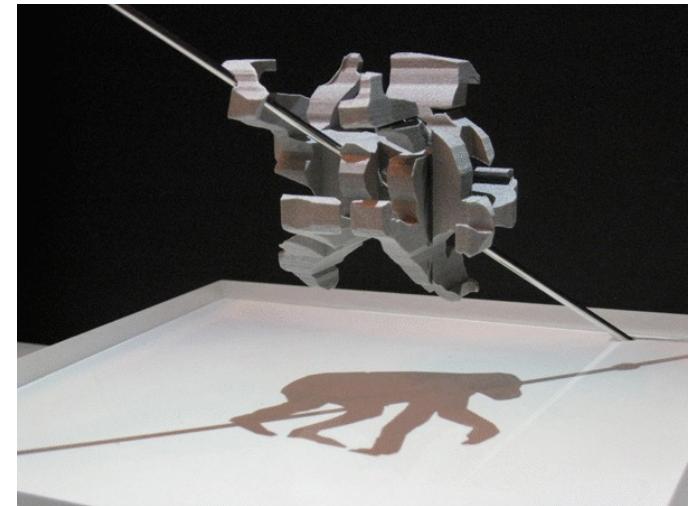
# What Makes 2D Computer Vision Hard?



[Sinha and Adelson 1993]

**Ambiguity:** Ill-posed Problem

# What Makes 2D Computer Vision Hard?



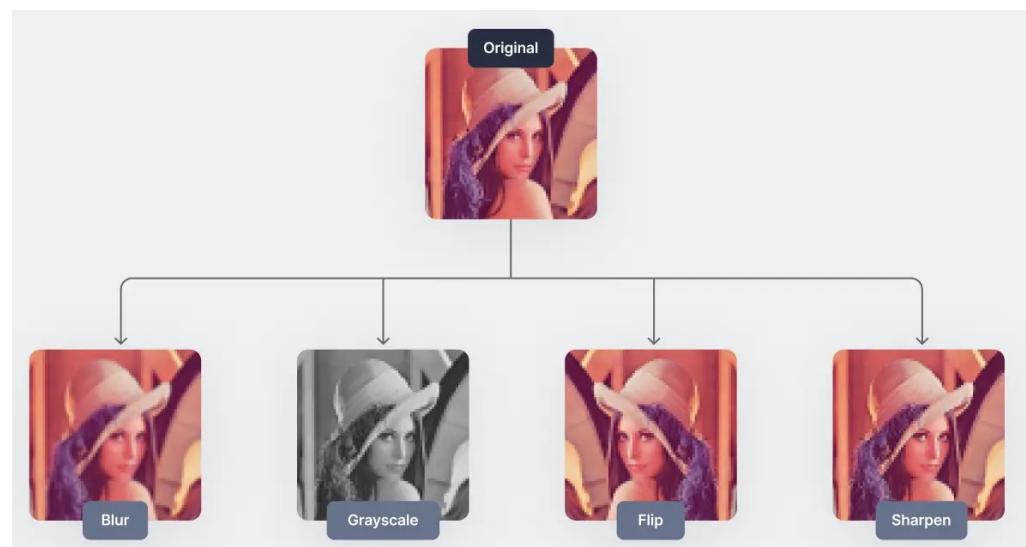
More examples:  
<https://www.jvmuntean.com/gallery-1#photos>)

# Basic Image Processing

There are many operations we can apply on a image.

Some examples:

- Resize
- Color manipulation
- Contrast
- Brightness
- Gamma
- Filtering
- Denosing
- ...

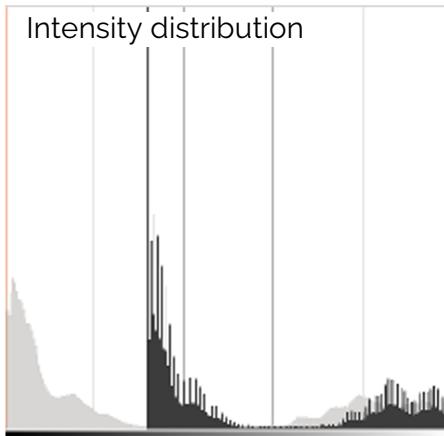


# Brightness

(overall intensity of an image)

Adjust brightness:  $g(x) = f(x) + \beta$

Increasing (/ decreasing) the  $\beta$  value will add (/ subtract) a constant value to every pixel. Pixel values outside of the  $[0 ; 1]$  range will be saturated (i.e. a pixel value higher (/ lesser) than 1 (/ 0) will be clamped to 1 (/ 0)).



In light gray, histogram of the original image,  
in dark gray when brightness = 0.15



$f(x)$ : input image pixel value

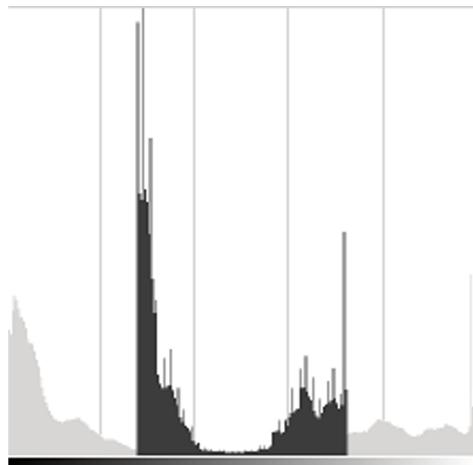


$g(x)$ : output image pixel value

# Contrast

Adjust contrast:  $g(x) = \alpha f(x)$

The  $\alpha$  parameter will modify how the intensity distribution spread. If  $\alpha < 1$ , the color levels will be compressed and the result will be an image with less contrast.



In light gray, histogram of the original image, in darker gray when contrast  $\alpha < 1$

Note:  $f(x)$  is normalized to have zero mean,  $g(x)$  will be converted back to  $(0,1)$  before visualization



$f(x)$ : input image pixel value



$g(x)$ : output image pixel value

# Gamma Correction

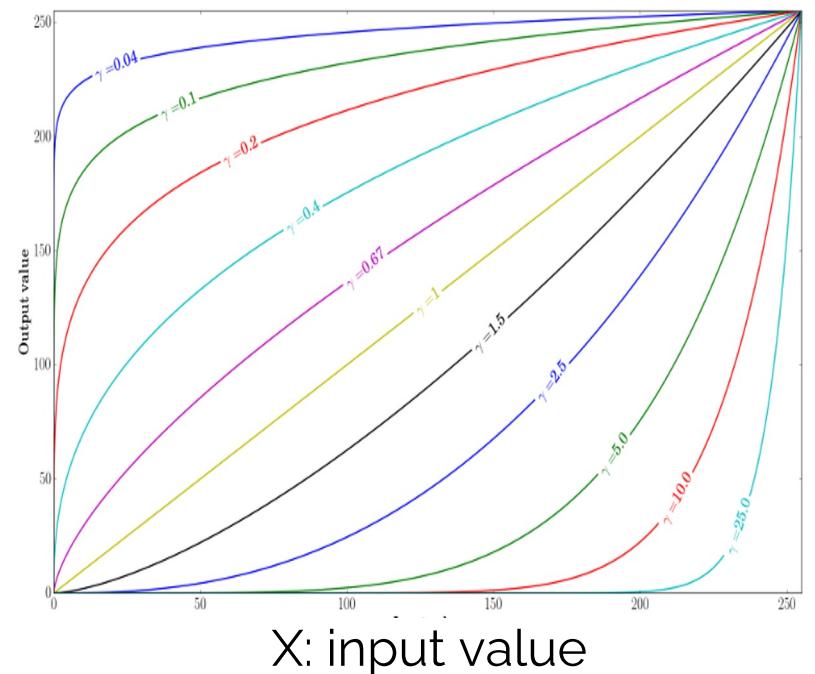
Gamma correction can be used to correct the brightness of an image by using a non linear transformation between the input values and the mapped output values:

$$g(\mathbf{x}) = [f(\mathbf{x})]^{1/\gamma}$$

As this relation is nonlinear, the effect will not be the same for all the pixels and will depend to their original value.

When  $\gamma < 1$ , the original dark regions will be brighter and the histogram will be shifted to the right.  
(reducing contrast)

Y: output value



# Gamma Correction



The following image has been corrected with:  $\alpha=1.3$  and  $\beta=40$ . The overall brightness has been improved but the clouds are now greatly saturated due to clipping.



The following image has been corrected with: gamma correction  $\gamma=0.4$ .

# Image Filtering

Image filtering: compute function of local neighborhood at each position

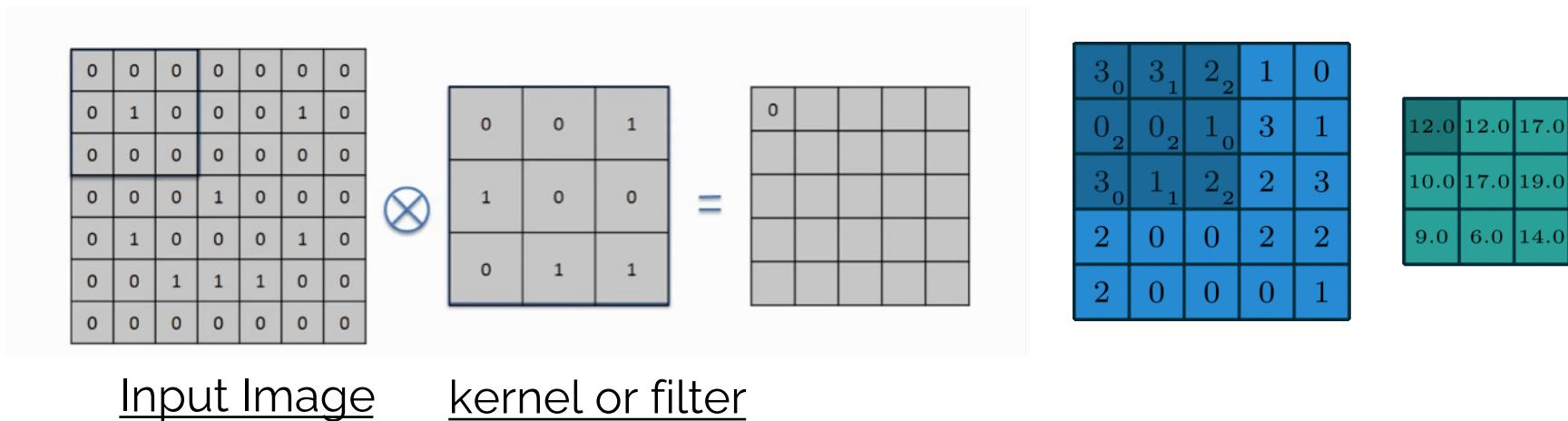
- Filtering can be used for:
  - Enhance images
    - Denoise, increase contrast, etc.
  - Extract information from images
    - Texture, edges, distinctive points, etc.
  - Detect patterns
    - Template matching
  - Deep Convolutional Networks



As seen on Instagram...

# Convolution operation

- Sliding the kernel (small 2D matrix) over the 2D image (Big 2D matrix)
- Performing elementwise multiplication with the part of the input it is currently on
- Sum up the result into a single output pixel
- Slide to the next location



# Padding

What about near the edge?

- the filter window falls off the edge of the image

Boundary padding

- Zero (smaller value around boundary)
- Circular/warp (best for panorama)
- Replicate
- Symmetric



Zero Pad



Circular



Replicate



Symmetric

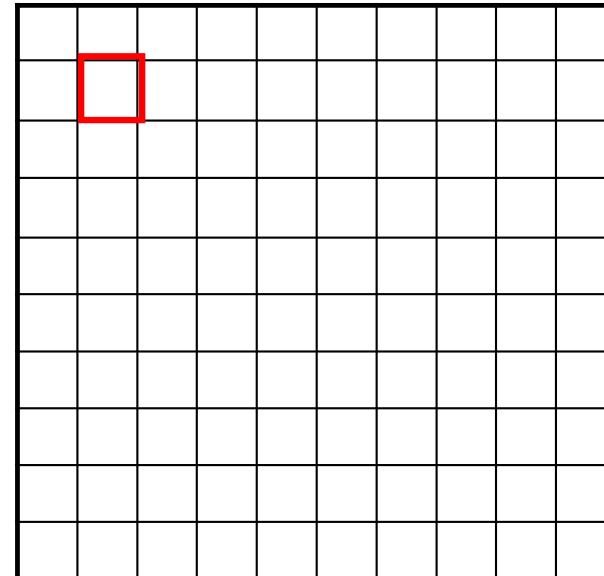
# Image Filtering (step by step)

$$g[\cdot, \cdot] \frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

$$f[.,.]$$

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

$$h[.,.]$$



Dot product at each position

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

Credit: S. Seitz

# Image Filtering

$f[.,.]$

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

Dot product at each position

$$g[\cdot, \cdot] \frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

$h[.,.]$

0	10									

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

Credit: S. Seitz

# Image Filtering

$$f[.,.]$$

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

Dot product at each position

$$g[\cdot, \cdot] \frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

$$h[.,.]$$

0	10	20							

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

Credit: S. Seitz

# Image Filtering

$f[.,.]$

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

Dot product at each position

$$g[\cdot, \cdot] \frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

$h[.,.]$

0	10	20	30								

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

Credit: S. Seitz

# Image Filtering

 $f[.,.]$ 

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

Dot product at each position

$$g[\cdot, \cdot] \frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

 $h[.,.]$ 

			0	10	20	30	30		

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

Credit: S. Seitz

# Image Filtering

$$f[.,.]$$

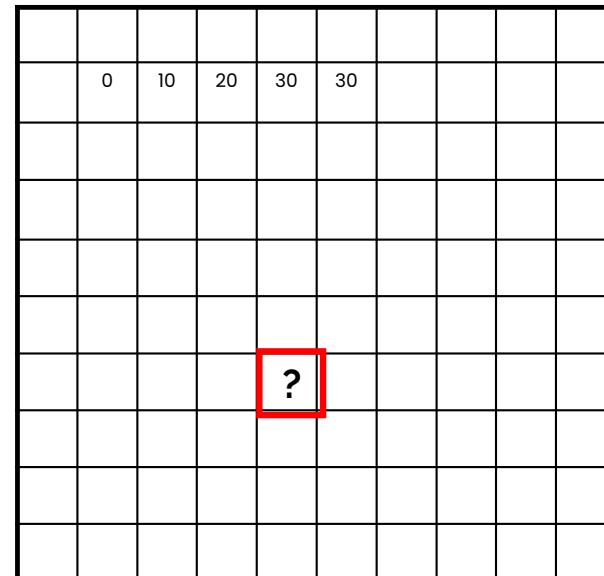
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0

Dot product at each position

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

$$g[\cdot, \cdot] \frac{1}{9} \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

$$h[.,.]$$



Credit: S. Seitz

# Image Filtering

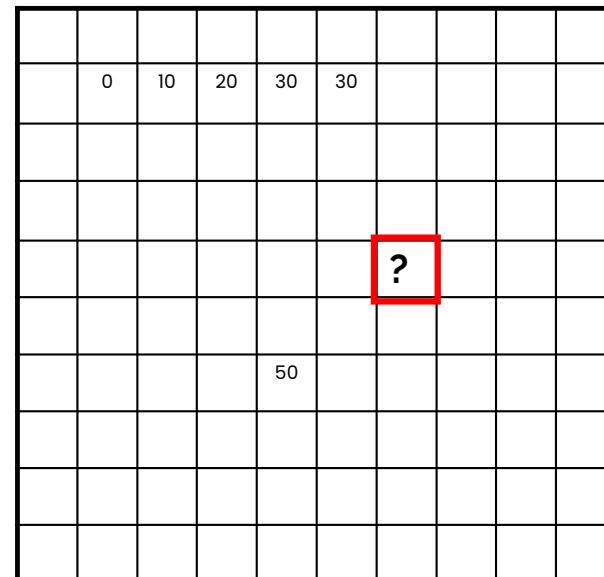
$$f[.,.]$$

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	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	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

Dot product at each position

$$g[\cdot, \cdot] \frac{1}{9} \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

$$h[.,.]$$



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

Credit: S. Seitz

# Image Filtering

$$f[.,.]$$

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

$$g[\cdot, \cdot]$$

$$\begin{array}{|ccc|}\hline & 1 & 1 & 1 \\ \hline 1 & & & \\ \hline 9 & 1 & 1 & 1 \\ \hline & 1 & 1 & 1 \\ \hline\end{array}$$

$$h[.,.]$$

	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	

Dot product at each position

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

Credit: S. Seitz

# Image Filtering: Box filter

What does it do?

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

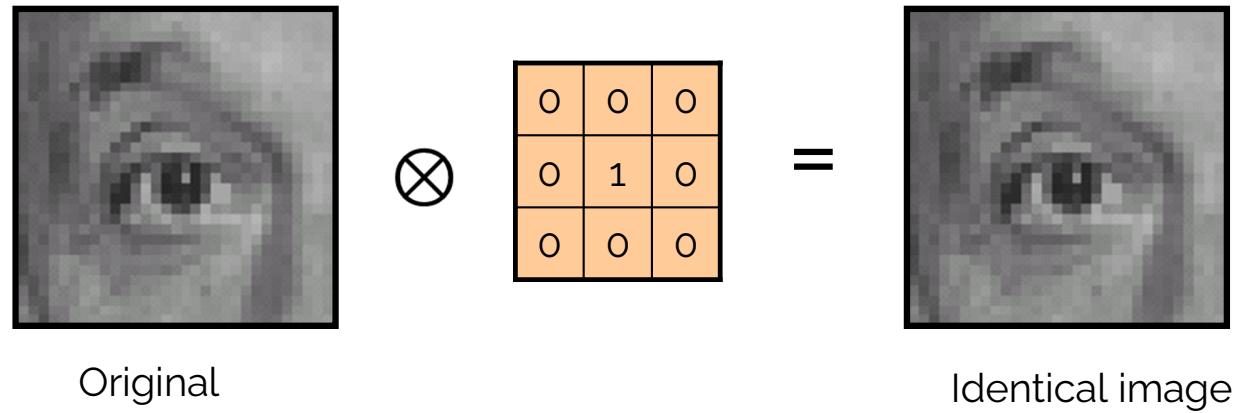
$$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

Slide credit: David Lowe  
(UBC)

# Smoothing with Box filter

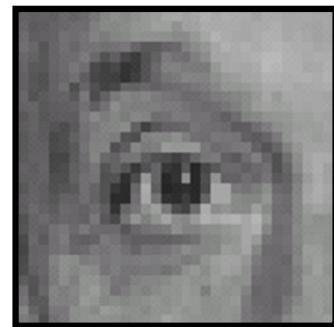


# Practice with Linear filter

$$\text{Original} \otimes \begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 0 & 1 & 0 \\ \hline 0 & 0 & 0 \\ \hline \end{array} = \text{Identical image}$$


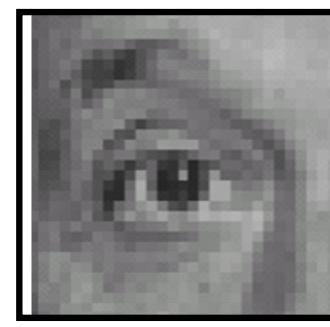
Source: D. Lowe

# Practice with Linear filter



0	0	0
1	0	0
0	0	0

=

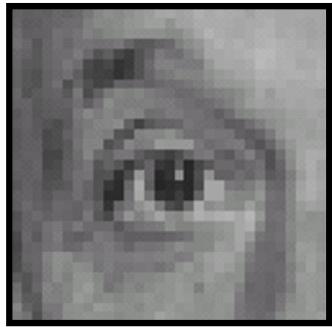


Original

Shifted right  
By 1 pixel

Source: D. Lowe

# Practice with Linear filter

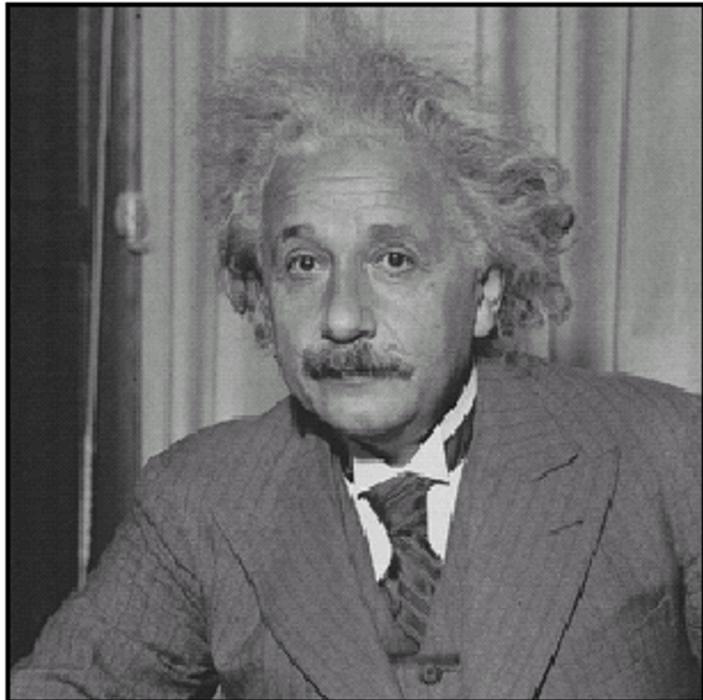

$$\otimes \left( \begin{array}{|ccc|} \hline 0 & 0 & 0 \\ \hline 0 & 2 & 0 \\ \hline 0 & 0 & 0 \\ \hline \end{array} - \frac{1}{9} \begin{array}{|ccc|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array} \right) = \text{Sharpened Image}$$


Original

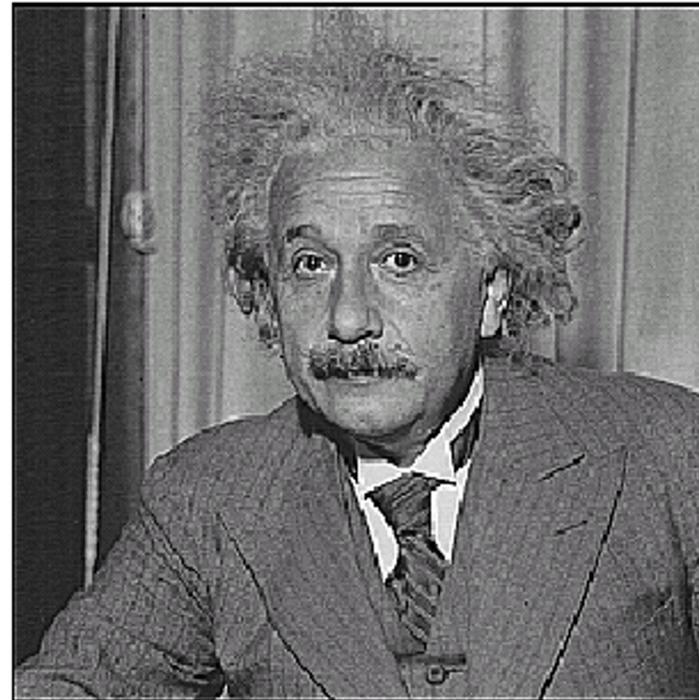
**Sharpening filter**  
(Amplifies the difference)

Source: D. Lowe

# Sharpening



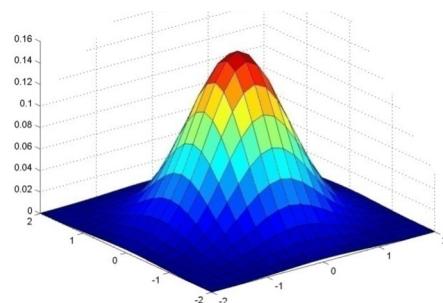
**before**



**after**

Source: D. Lowe

# Other Filters: Gaussian Kernel



Approximated by:

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

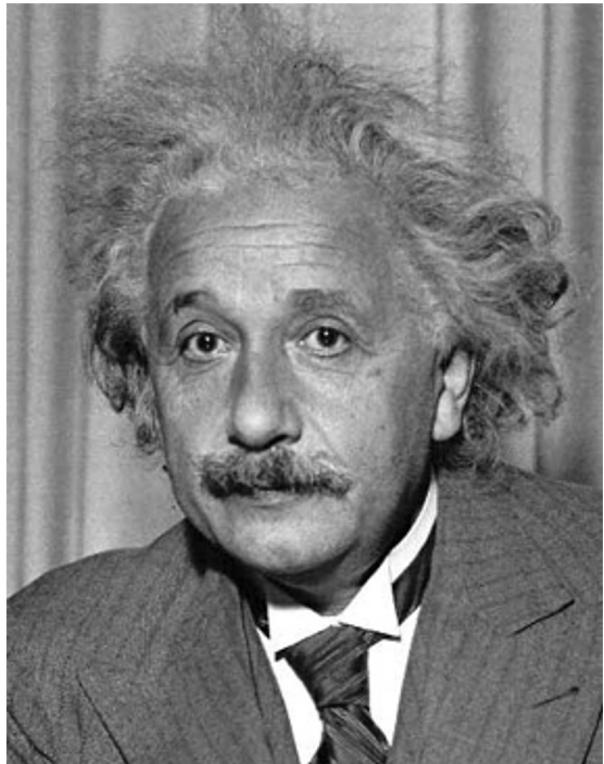
$5 \times 5, \sigma = 1$

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

Constant factor at front  
makes matrix sum to 1

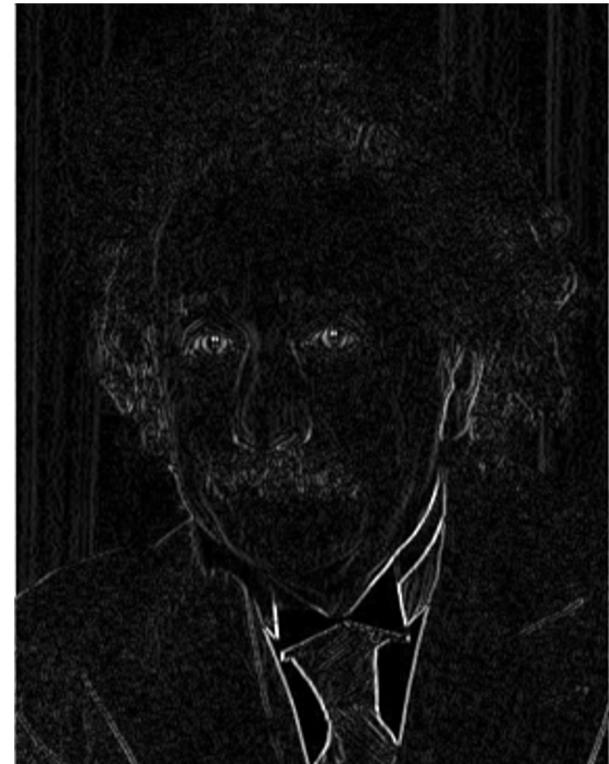
Source: C. Rasmussen

## Other Filters: Edge



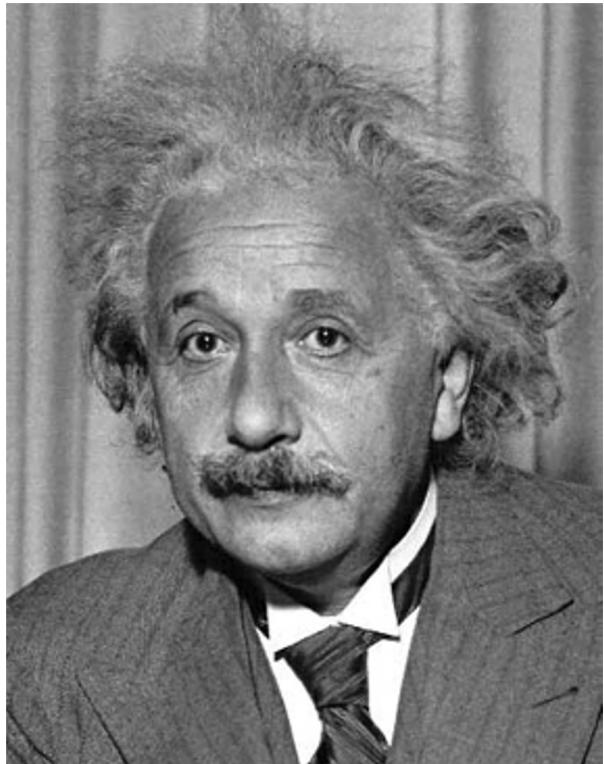
1	0	-1
2	0	-2
1	0	-1

Vertical  
Sobel



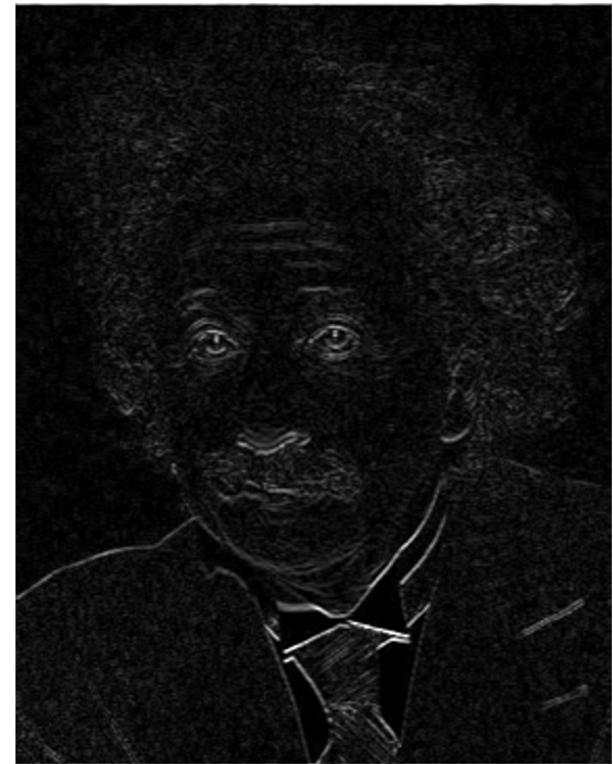
Vertical Edge  
(absolute value)

## Other Filters: Edge



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

Horizontal  
Sobel

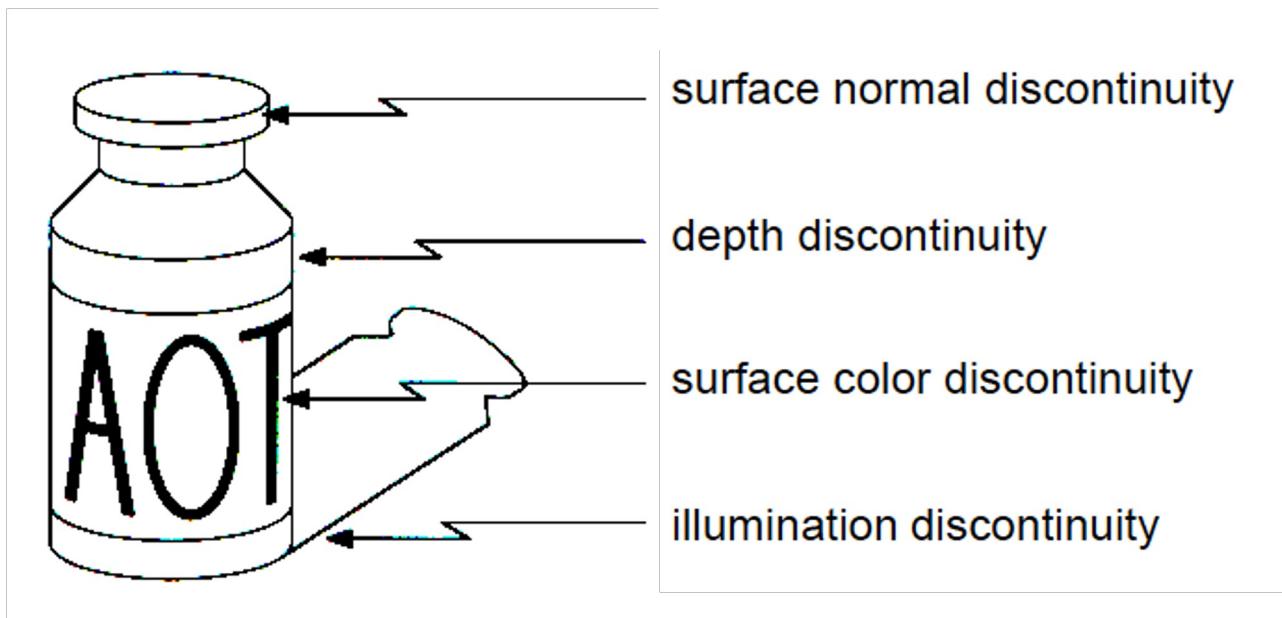


Horizontal Edge  
(absolute value)

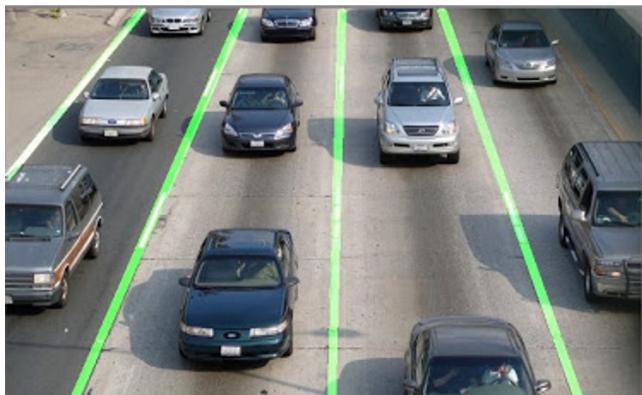
# Edge Detection

**Goal:** Identify sudden changes (discontinuities) in an image

Why? Intuitively, edges carry a lot of the semantic and shape information from the image



# Edge Detection in Robotics



Extract information Lane detection

geometry and camera viewpoint

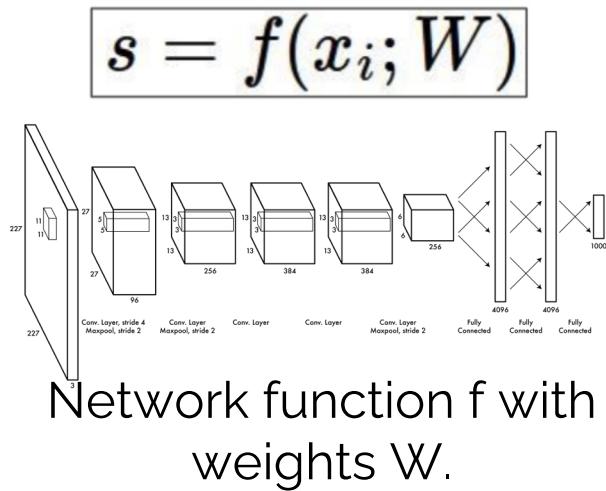
## Convolutional Neural Network (CNN)

Can we automatically *learn* which filter to apply  
to get most useful information?

# Convnet Building Blocks



Input Image:  
 $x_i$



cat  
car  
frog

**3.2**  
**5.1**  
**-1.7**

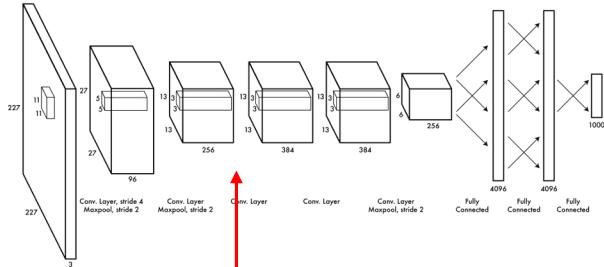
Output:  $s$  is a vector of classification scores

# Convnet Building Blocks



Input Image:  
 $x_i$

$$s = f(x_i; W)$$



Network function  $f$  with  
weights  $W$ .

Network Layers (function approximator):  
Convolutional layers  
Pooling layers  
Fully connected layers

cat  
car  
frog

**3.2**  
**5.1**  
**-1.7**

**1**  
0  
0

Groundtruth  
Label

Losses (objectives):  
Regression  
Classification

# Convolutional Layer

**Input:** an image

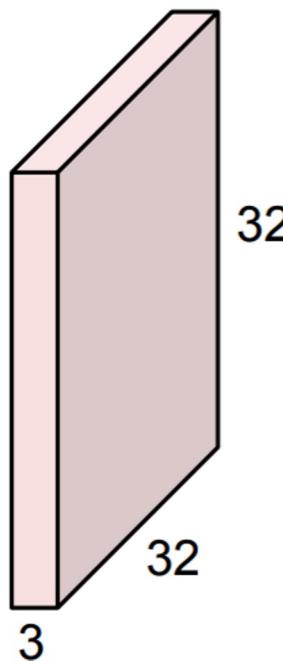
**Processing:** convolution with multiple filters

**Output:** an image, # channels = # filters, weighted sum of input

```
torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, padding_mode='zeros')
```

# Convolutional Layer

3x 32 x32 image

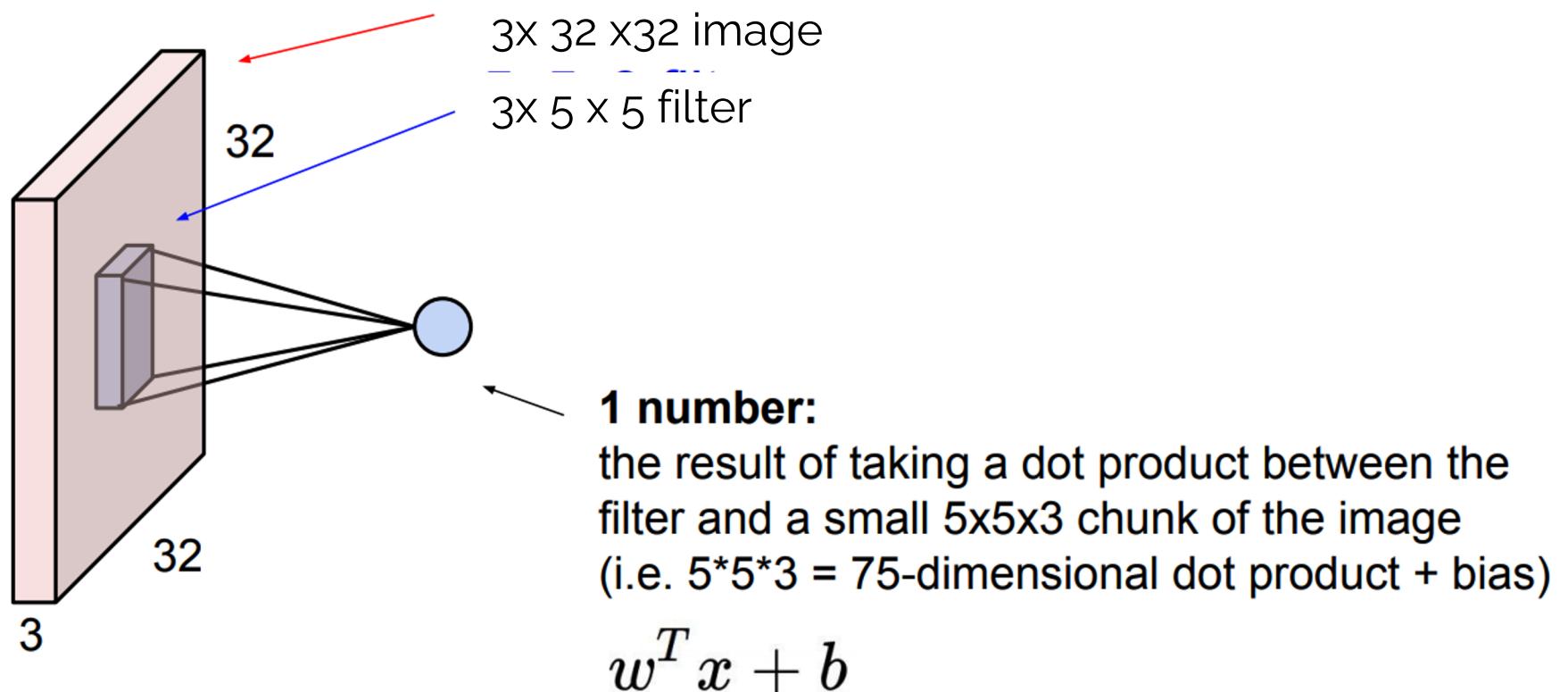


3x 5 x 5 filter

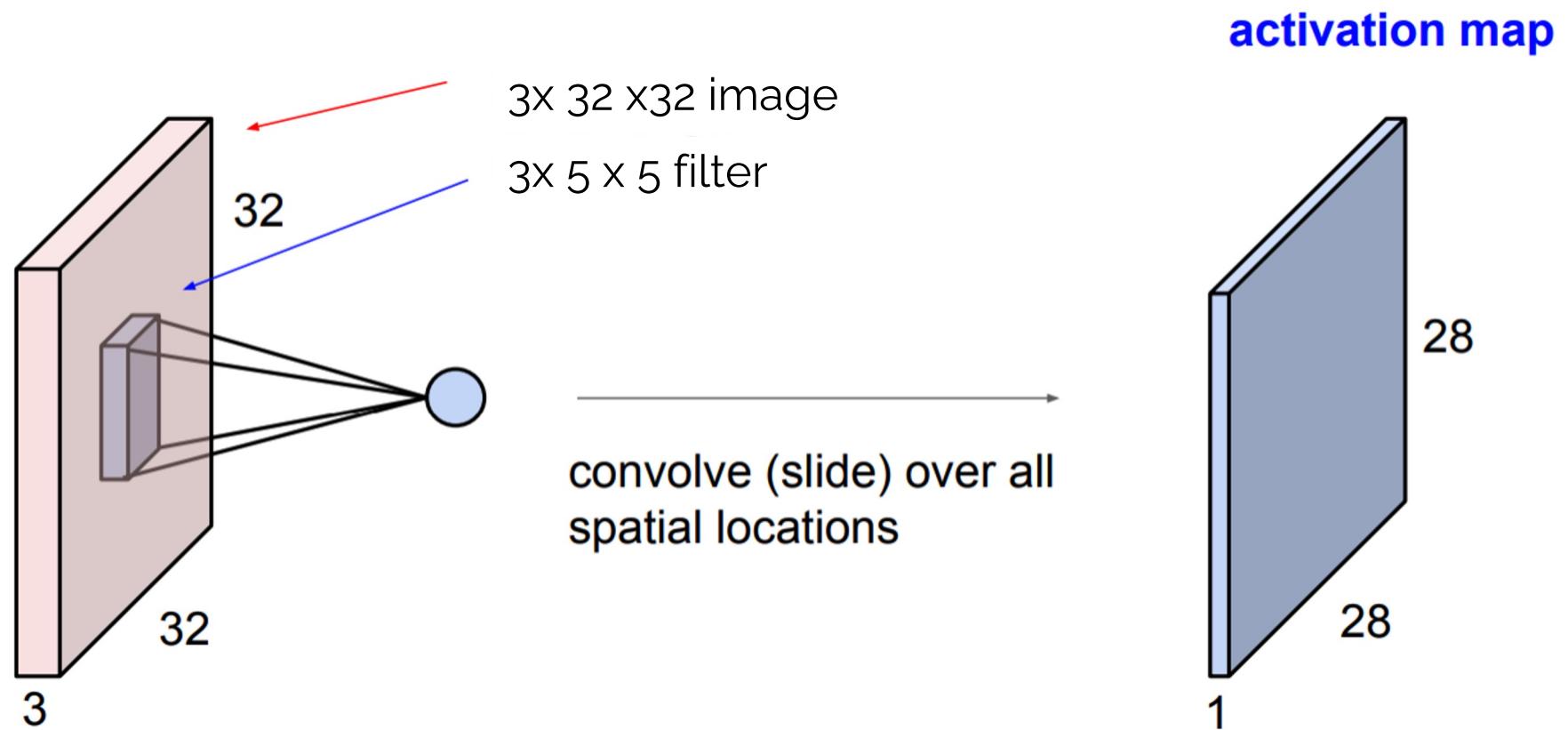


**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

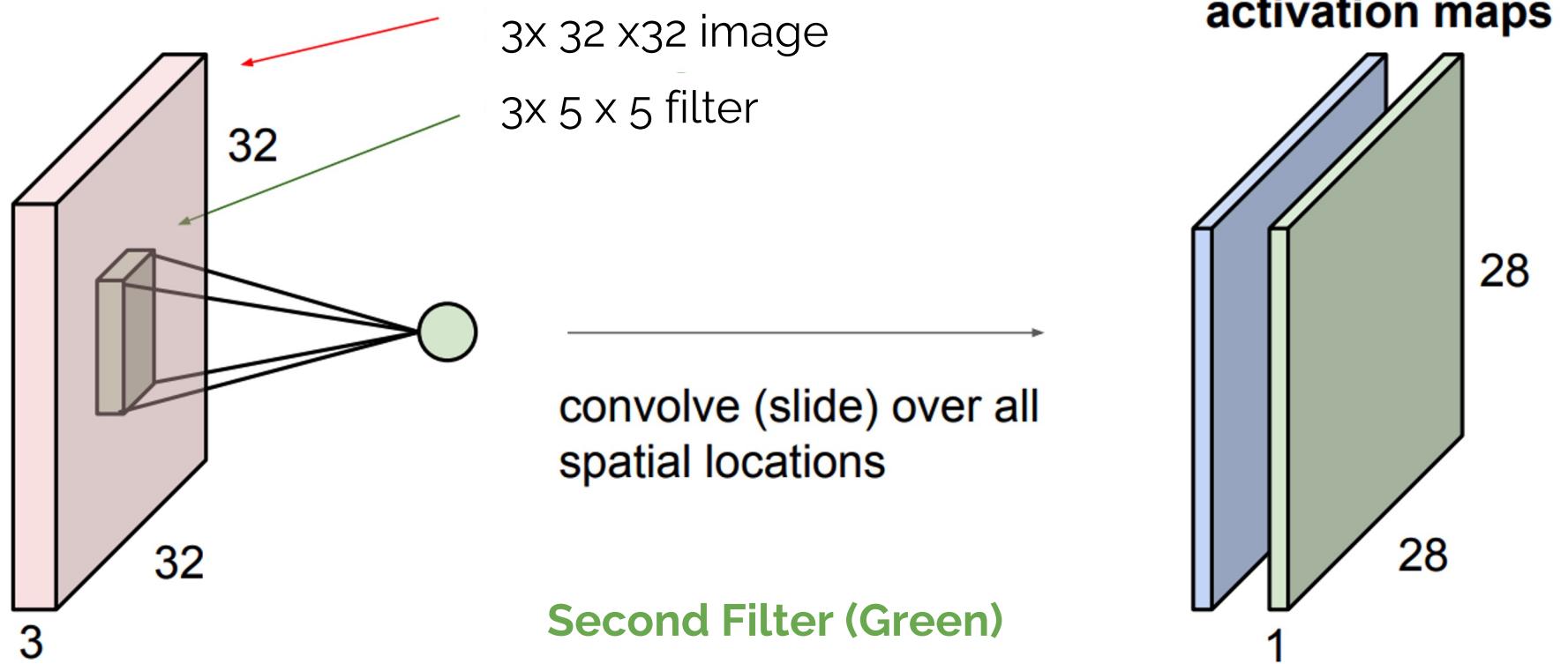
# Convolutional Layer



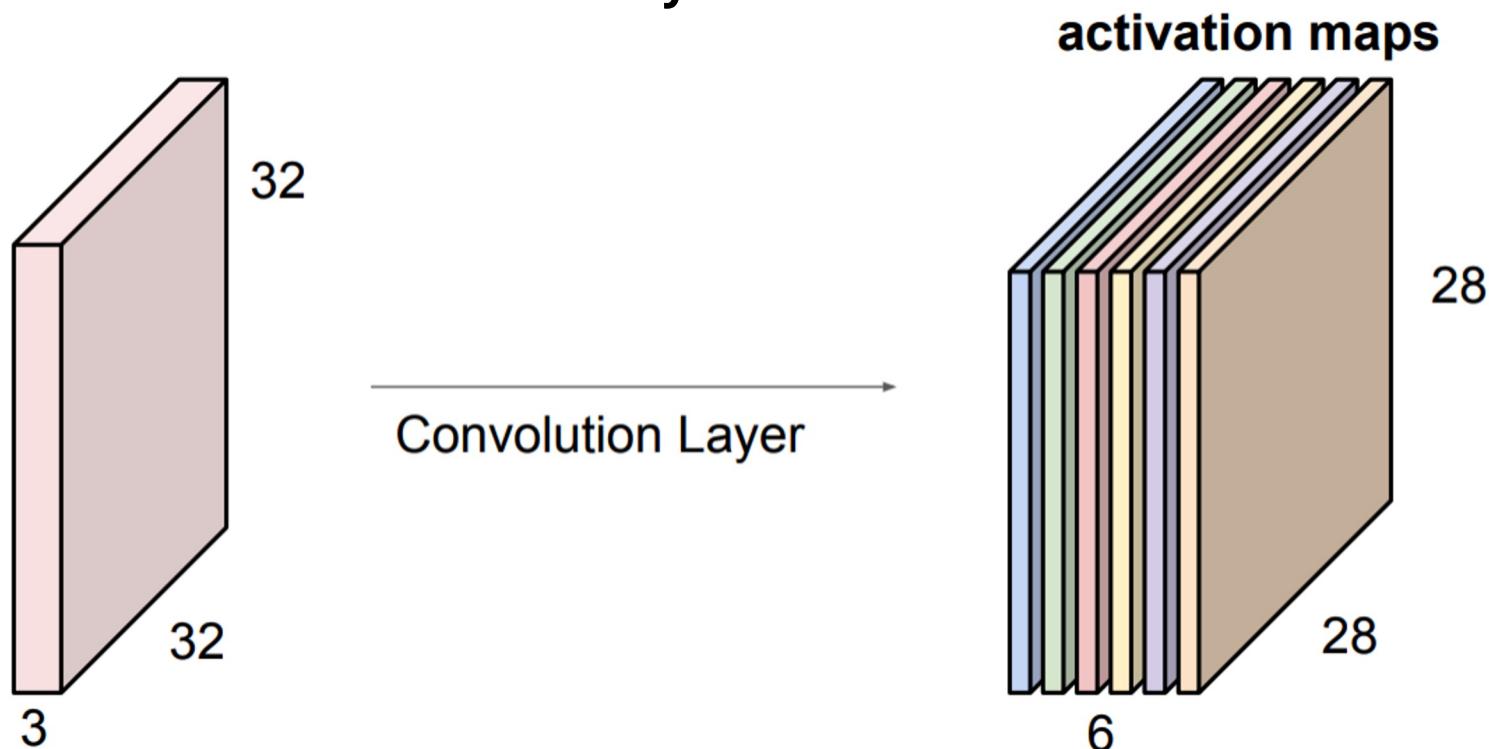
# Convolutional Layer



# Convolutional Layer



# Convolutional Layer



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!

# Kernel size

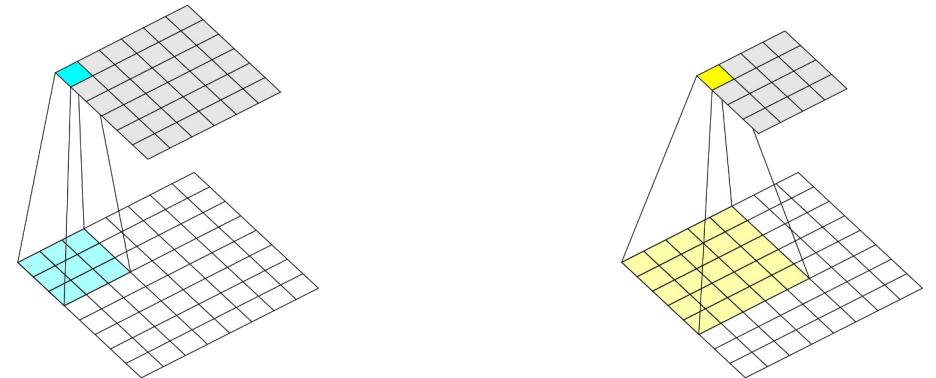
How big the filter for a layer is

Typically  $1 \times 1 \leftrightarrow 7 \times 7$

**Larger** kernel aggregate information from **larger** local region (**bigger receptive field**)

$1 \times 1$  is just linear combination of channels in previous image (no spatial processing)

Filters usually have same number of channels as input image.

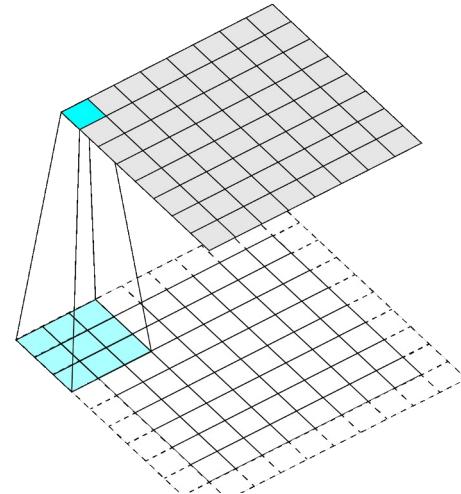
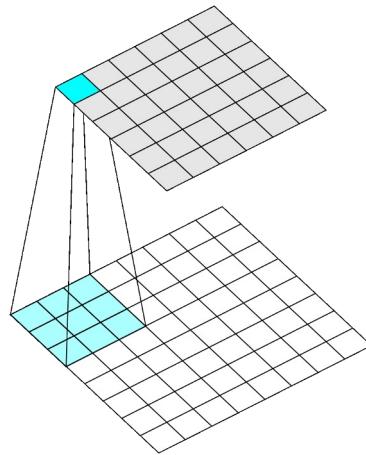


# Padding

Convolutions have problems on edges

Do nothing: output a little smaller than input

Pad: add extra pixels on edge

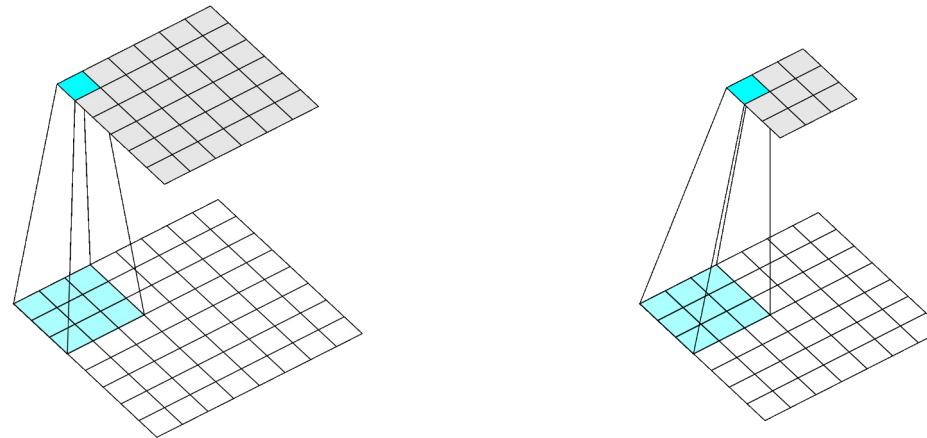


# Stride

How far to move filter between applications

We've done stride 1 convolutions up until now, approximately preserves image size

Could move filter further, downsample image



# Input - output size

- Input  $C_{in} \times H_{in} \times W_{in}$  image
- Conv2d(K,N,S,P): Number of kernel K, kernel size [NxN], stride: s, padding: p
- What's the output feature size ( $C_{out}, H_{out}, W_{out}$ ) ?

$$H_{out} = \left\lfloor \frac{H_{in} + 2 * padding - (N-1) - 1}{stride} + 1 \right\rfloor$$

$$W_{out} = \left\lfloor \frac{W_{in} + 2 * padding - (N-1) - 1}{stride} + 1 \right\rfloor$$

$$C_{out} = K$$

- Ex: Input: 3x200x300 Conv2d (18,3,1,1) Output size?

# Input - output size

- Input  $C_{in} \times H_{in} \times W_{in}$  image
- Conv2d(K,N,S,P): Number of kernel K, kernel size [NxN], stride: s, padding: p
- What's the output feature size ( $C_{out}, H_{out}, W_{out}$ ) ?

$$H_{out} = \left\lfloor \frac{H_{in} + 2 * padding - (N-1) - 1}{stride} + 1 \right\rfloor$$

$$W_{out} = \left\lfloor \frac{W_{in} + 2 * padding - (N-1) - 1}{stride} + 1 \right\rfloor$$

$$C_{out} = K$$

- Ex: Input: 3x200x300 Conv2d (18,3,1,1) Output size?
- How about Conv2d (18,5,1,2)?

# Pooling Layer

**Input:** an image

**Processing:** pool pixel values over region

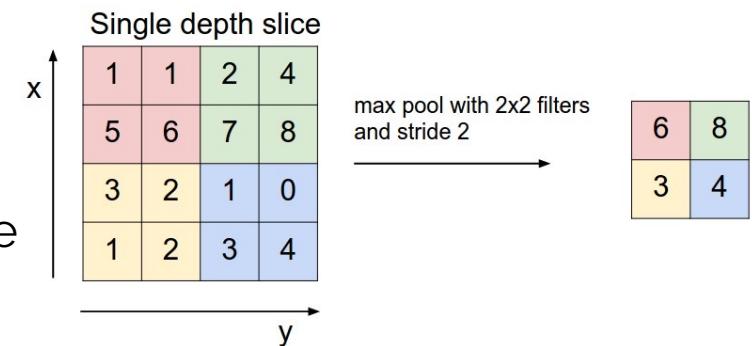
**Output:** an image, shrunk by a factor of the stride

Hyperparameters:

What kind of pooling? Average, mean, max, min

How big of stride? Controls downsampling

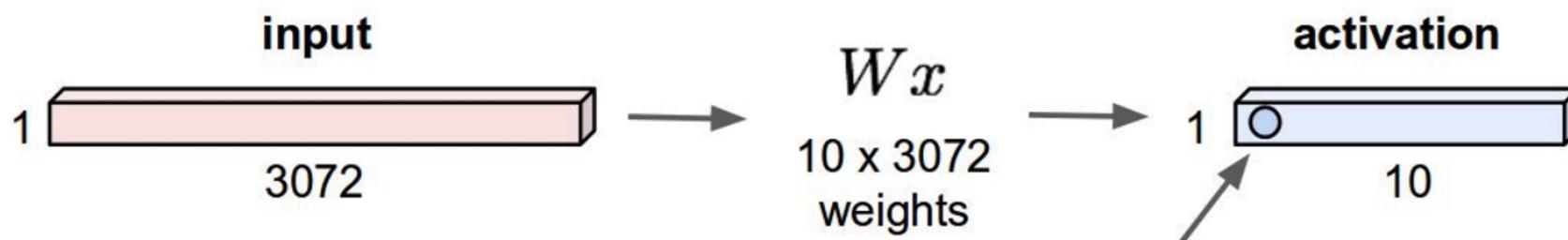
How big of region? Usually not bigger than stride



Most common: maxpooling 2x2 stride of 2 or 3x3 stride of 3

# Fully Connected Layer (InnerProduct)

Every input neuron connects to every output neuron



Often used to go from image feature map -> final output or map image features to a single vector

Eliminates spatial information

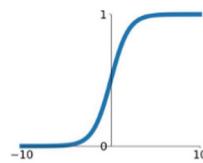
# Activation Layer

Used to increase non-linearity of the network without affecting receptive fields of the conv layers.

Most commonly used activation layer: ReLU

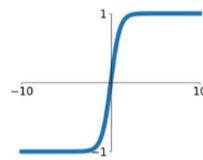
**Sigmoid**

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



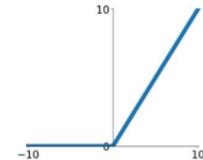
**tanh**

$$\tanh(x)$$



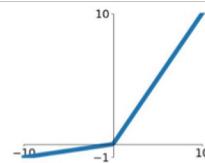
**ReLU**

$$\max(0, x)$$



**Leaky ReLU**

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

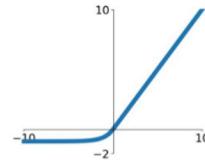


**Maxout**

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

**ELU**

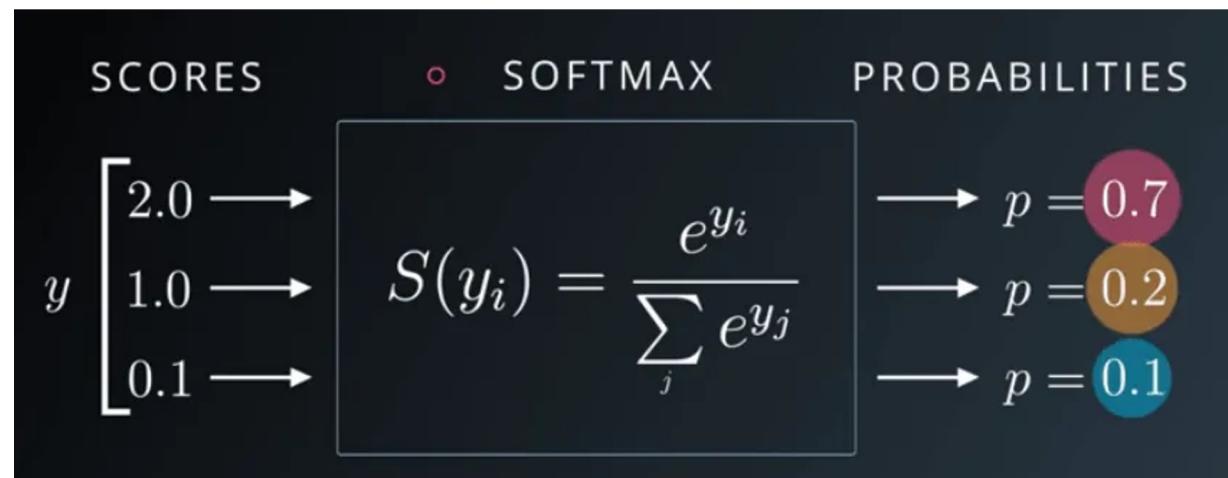
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



# Activation Layer: Softmax

A special kind of activation layer, usually at the end of fully connected layer and before cross-entropy loss.

Can be viewed as a fancy normalization function that produce a probability distribution vector.



# Loss Function

---

A loss function scores how far off a prediction is from the desired “target” output.

- `loss(prediction,target)` returns a number called “the loss”
- Big Loss = Bad Error
- Small Loss = Minor Error
- Zero Loss = No Error

# Loss Function

---

## Regression

### **Predict real-valued output**

What temperature will it be tomorrow?

Loss function: squared error, L1 error ...

## Classification

### **Predict category output**

Will it be sunny tomorrow?

Often predicts continuous values which are probabilities of different labels, use the one with max probabilities as predicted class.

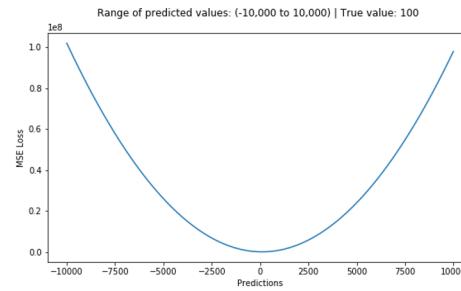
Loss function: cross-entropy

# Regression

---

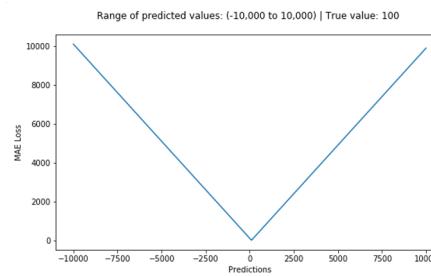
**L2 Loss**

$$MSE = \frac{\sum_{i=1}^n (y_i - y_i^p)^2}{n}$$



**L1 Loss**

$$MAE = \frac{\sum_{i=1}^n |y_i - y_i^p|}{n}$$



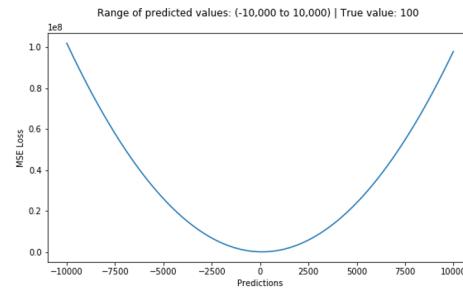
L1 loss is more robust to outliers, but its derivatives around zero are not continuous, making it inefficient to find the solution. L2 loss is sensitive to outliers, but gives a more stable and closed form solution.

# Regression

---

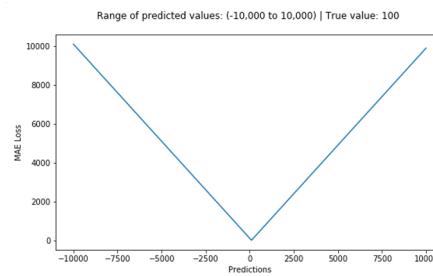
**L2 Loss**

$$MSE = \frac{\sum_{i=1}^n (y_i - y_i^p)^2}{n}$$



**L1 Loss**

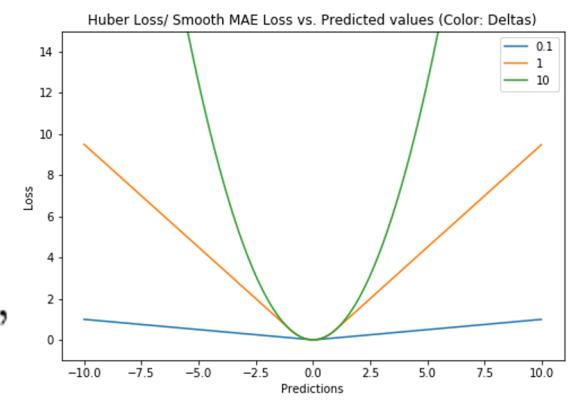
$$MAE = \frac{\sum_{i=1}^n |y_i - y_i^p|}{n}$$



**Smooth L1**

$$L_\delta(y, f(x)) = \begin{cases} \frac{1}{2}(y - f(x))^2 & \text{for } |y - f(x)| \leq \delta, \\ \delta |y - f(x)| - \frac{1}{2}\delta^2 & \text{otherwise.} \end{cases}$$

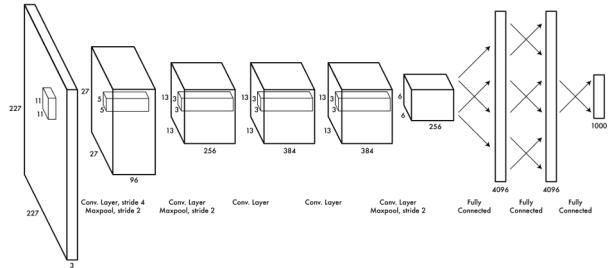
**L2 when  $\delta \sim 0$  and L1 when  $\delta \sim \infty$  (large numbers.)**



# Classification: CrossEntropy



Input Image:  
 $x_i$



Network function  $f$  with  
weights  $W$ .

$$s = f(x_i; W)$$

cat  
car  
frog

**3.2**  
5.1  
-1.7

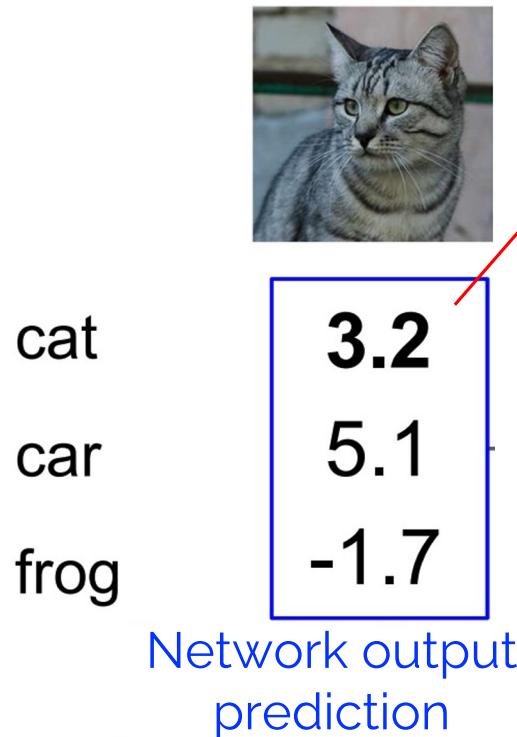
**1**  
0  
0

$s$  is a vector of  
classification  
scores

Groundtruth  
Label

# Classification: CrossEntropy

---

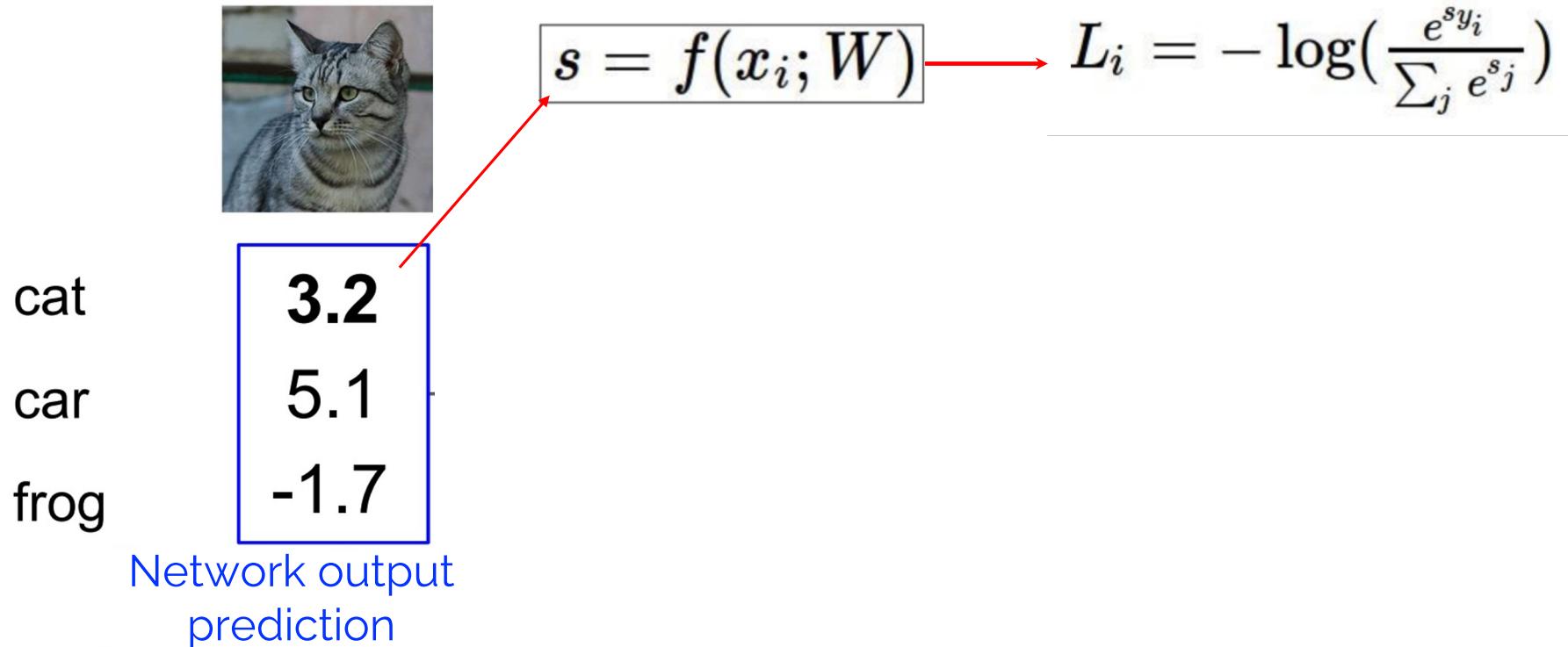


$$s = f(x_i; W)$$

scores = unnormalized log probabilities of the classes.

# Classification: CrossEntropy

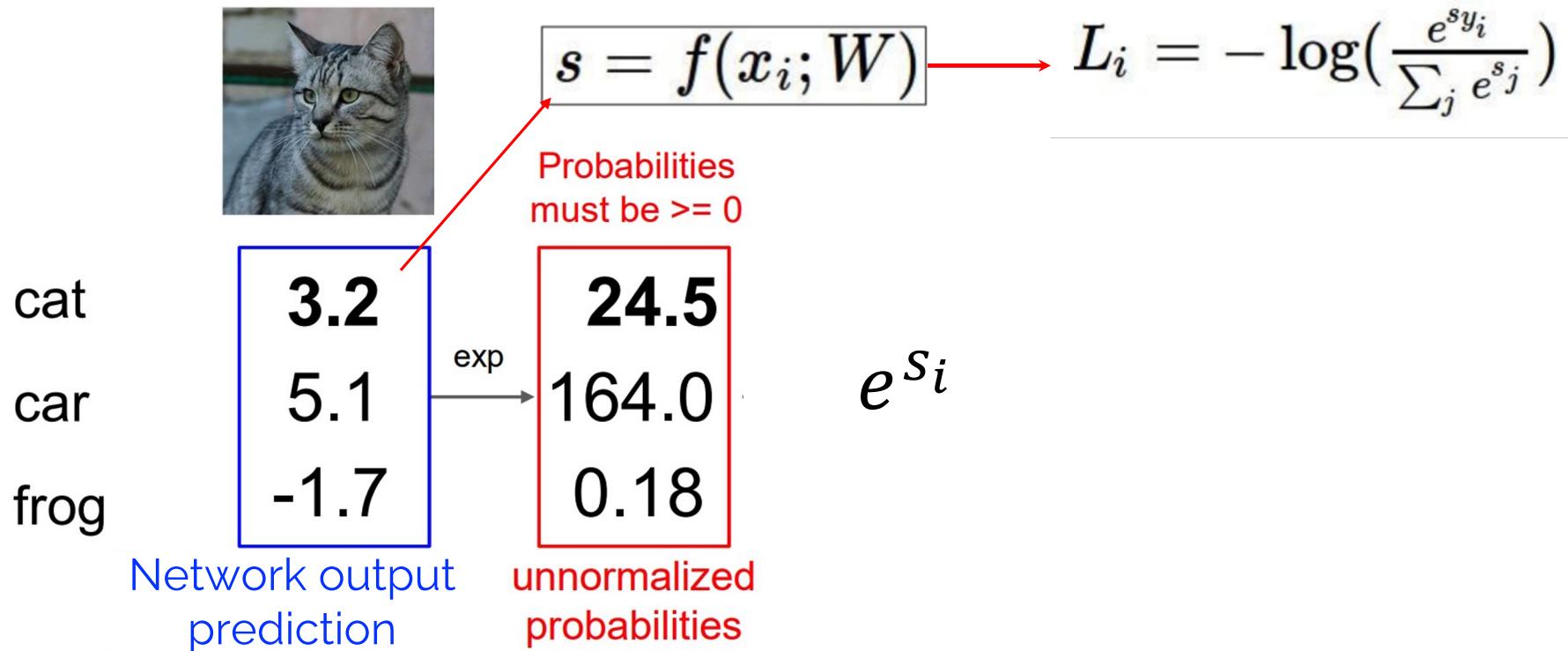
---



Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

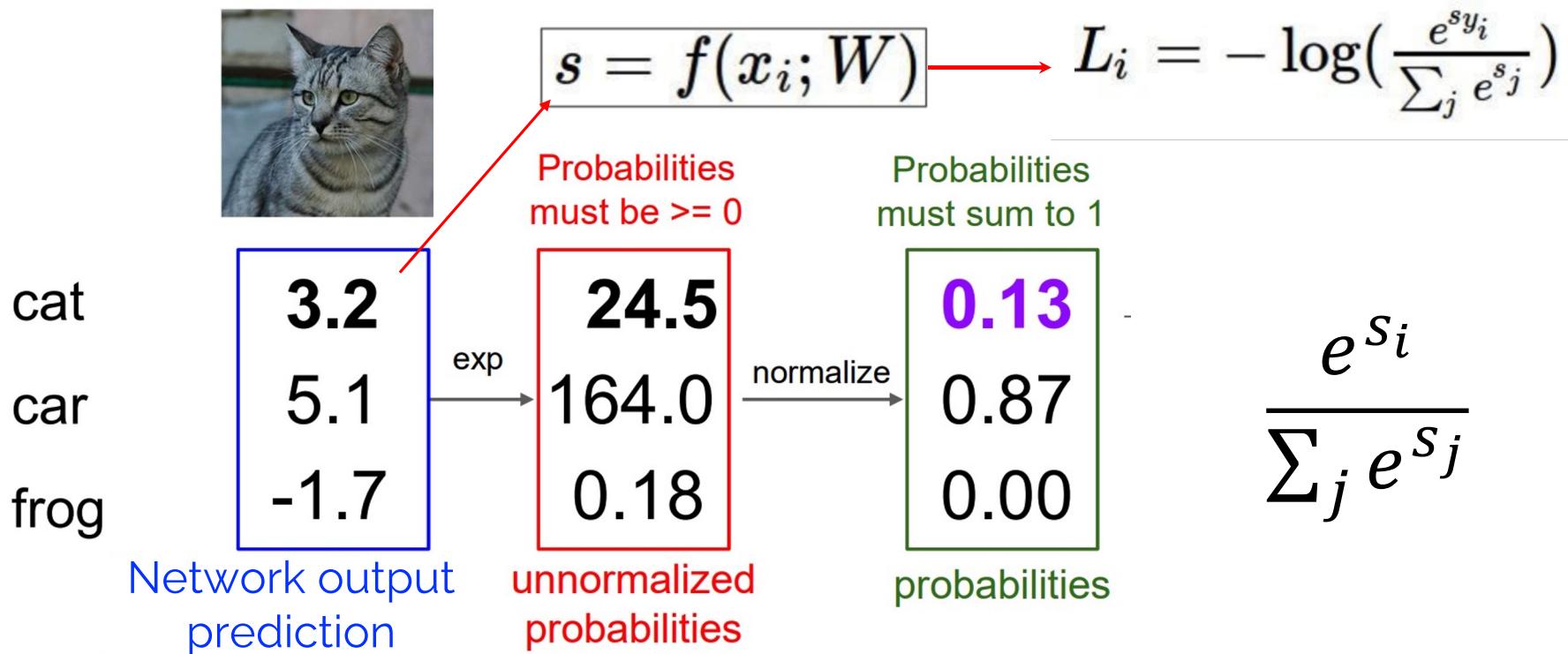
# Classification: CrossEntropy

---



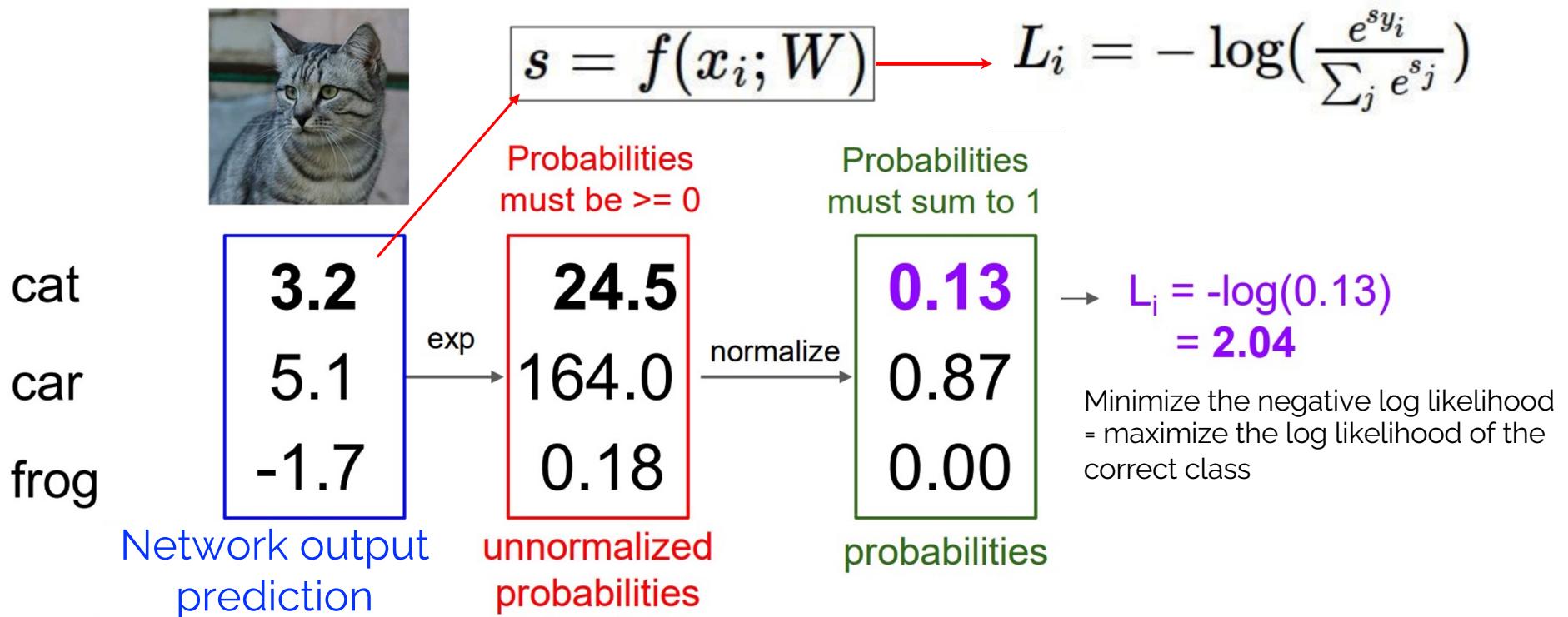
Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

# Classification: CrossEntropy



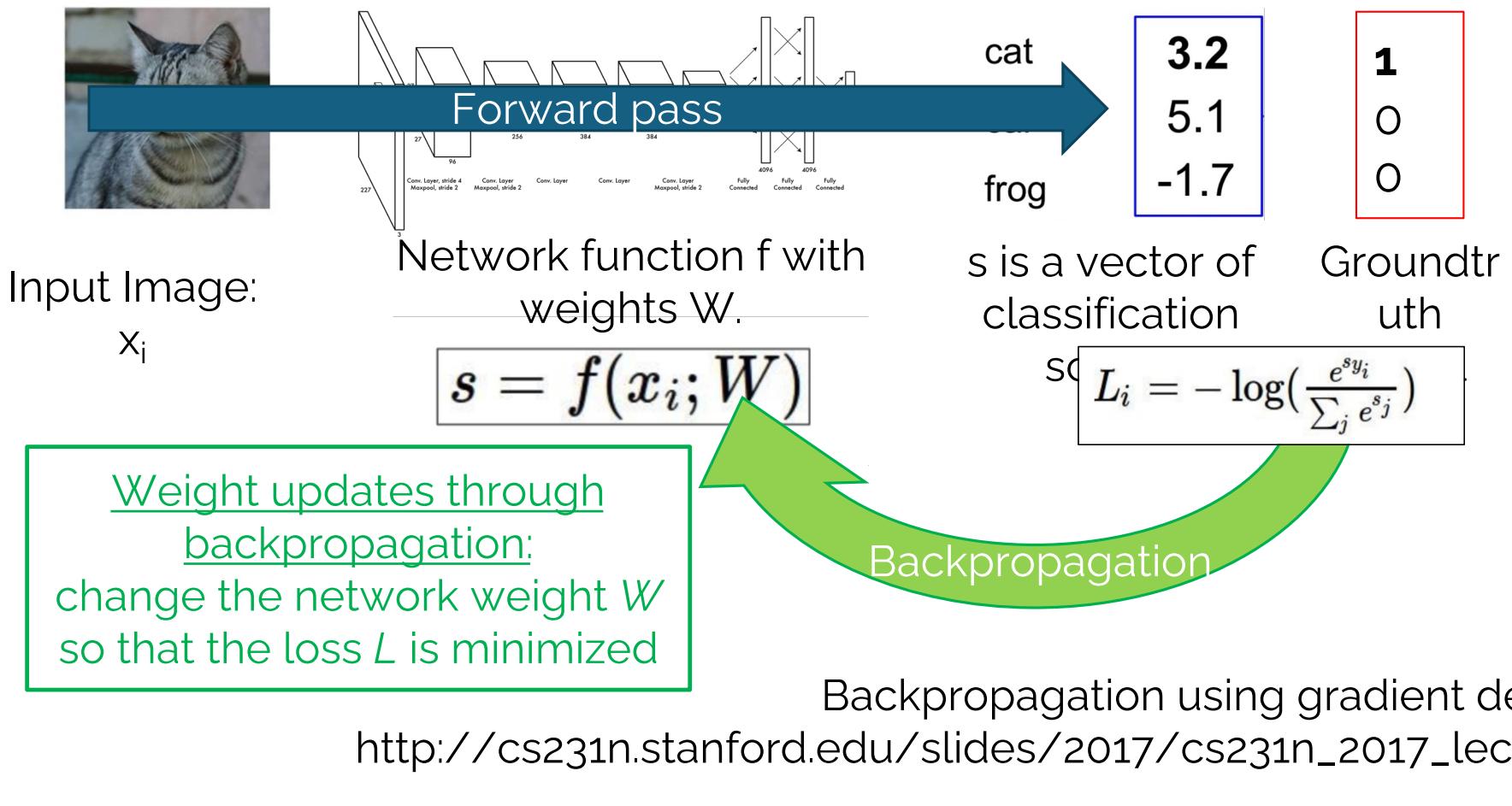
Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

# Classification: CrossEntropy



Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

# Training the network



# Summary of Convnet Building Blocks

## **Convolutional layers:**

Used convolution operation to extract features  
Convolution with stride > 1 will downsample image

## **Pooling layers:**

Used to downsample feature maps, make processing more efficient

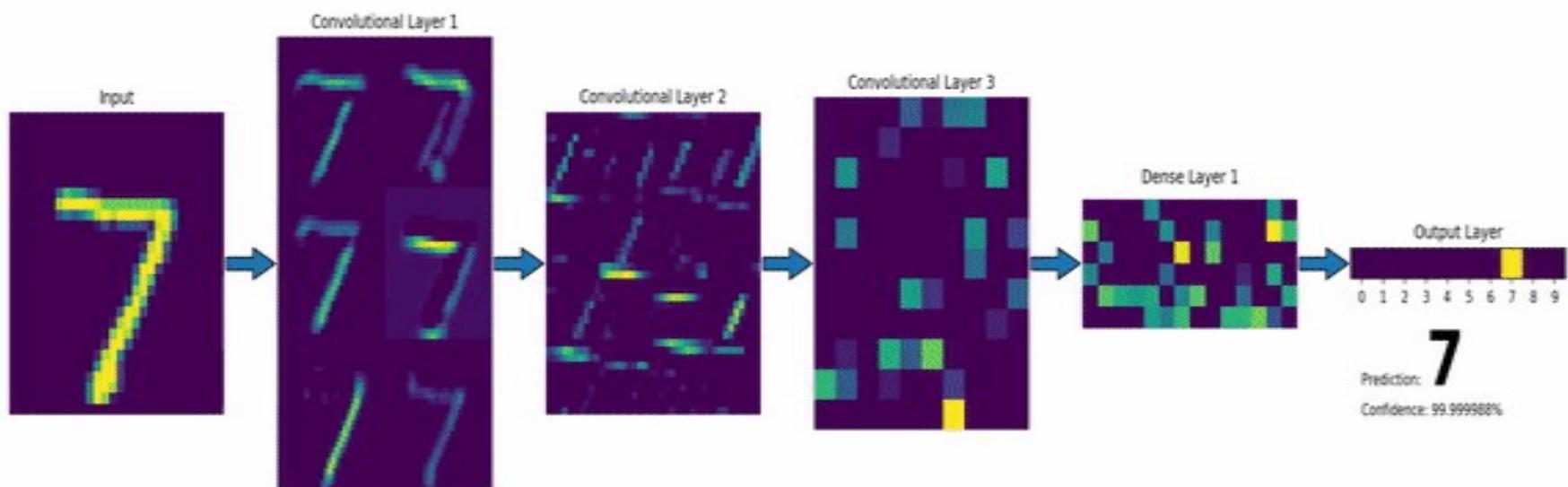
## **Fully connected layers:**

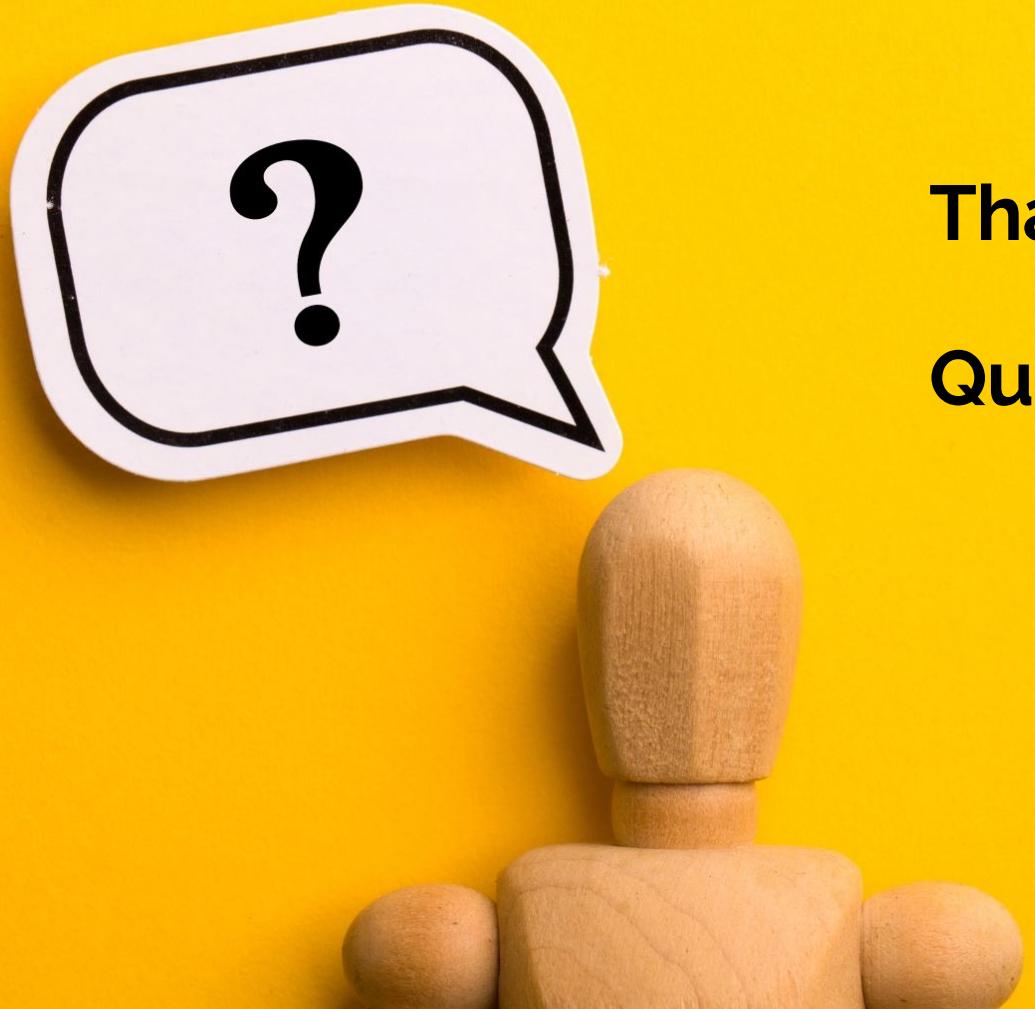
Often used as last layer, to map image features -> prediction  
No spatial information  
Inefficient: lots of weights, no weight sharing

**Activation layers:** increase non-linearity of the network

**Loss function:** classification and regression

# MNIST Example





**That's it for today!**

**Questions?**