

Introduction to Computer Vision

CS 4243 Computer Vision & Pattern Recognition

Angela Yao

Today's Lecture

1. What is Computer Vision?
2. Why Is Computer Vision Difficult?
3. Digital Image Capture & Representations

What is Computer Vision?

Every Image Tells A Story



Goal of computer vision: perceive the “story” behind the picture.

Compute properties of the world

- 3D shape
- Names of people or objects
- What happened?
- When, how, why?

Can Computers Match Human Perception?



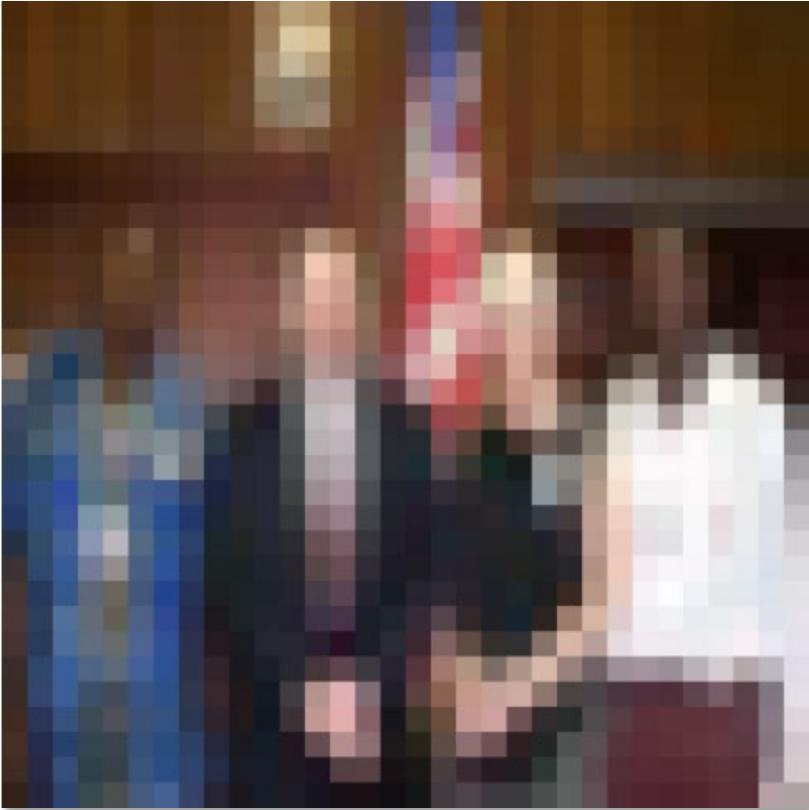
Yes and no (mainly no)

- computers can be better at “easy” things
- humans are much better at “hard” things

But huge progress has been made

- Accelerating in past 5 years due to deep learning
- What is considered “hard” keeps changing

Humans can tell a lot from very little ...



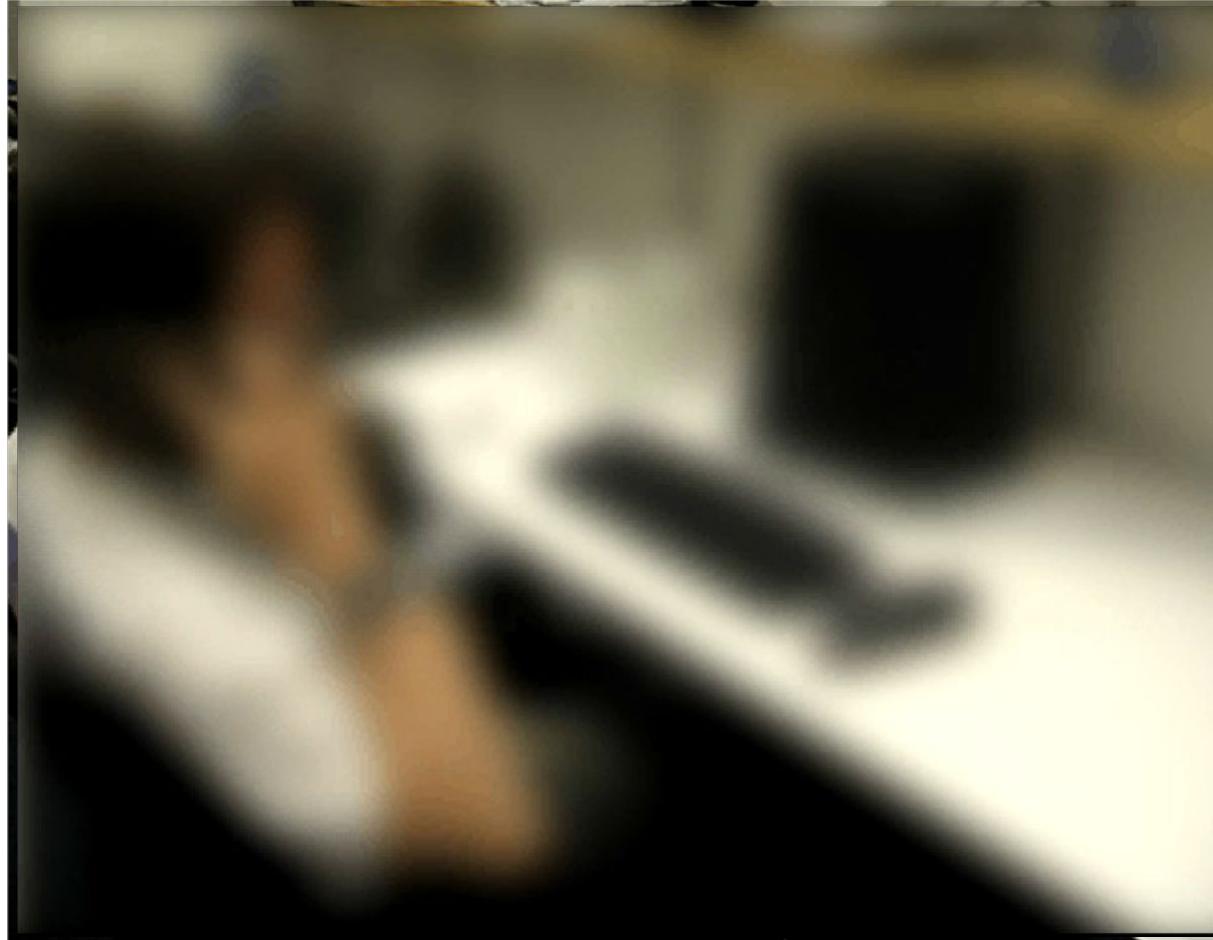
Source: "[80 million tiny images](#)" by Torralba, et al.

What do you see in this image?

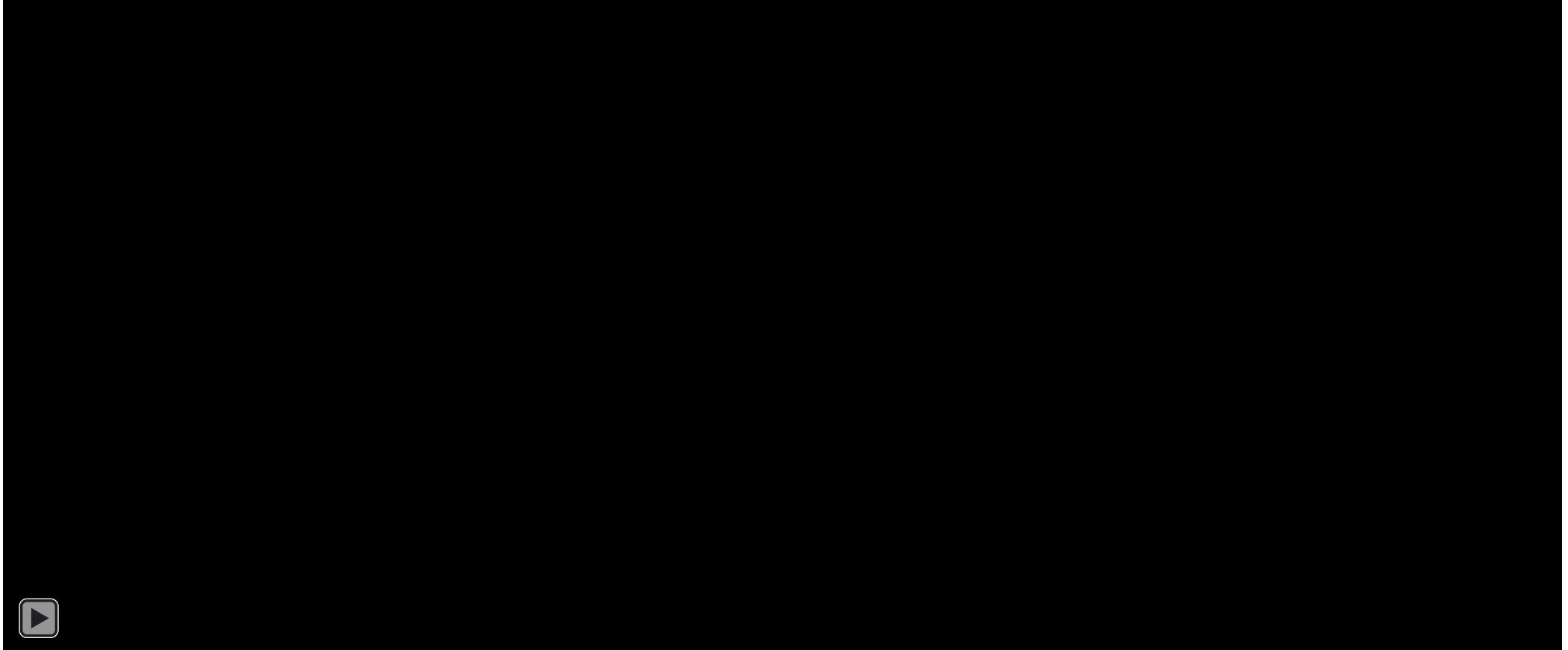
~900 pixels (32x32) is all we need to see 4 people standing in front of some cabinets and a flag(?).

Two are white and two are black.
Two are men and two are women?

But we rely heavily on prior knowledge ...

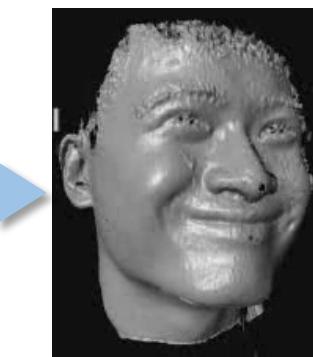
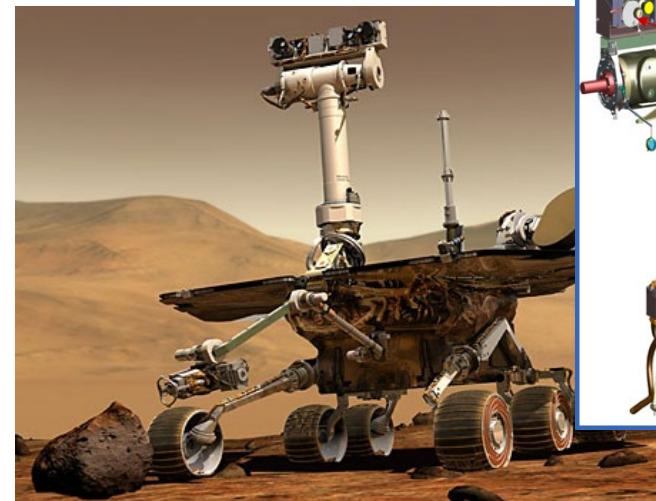


What is the goal of computer vision?



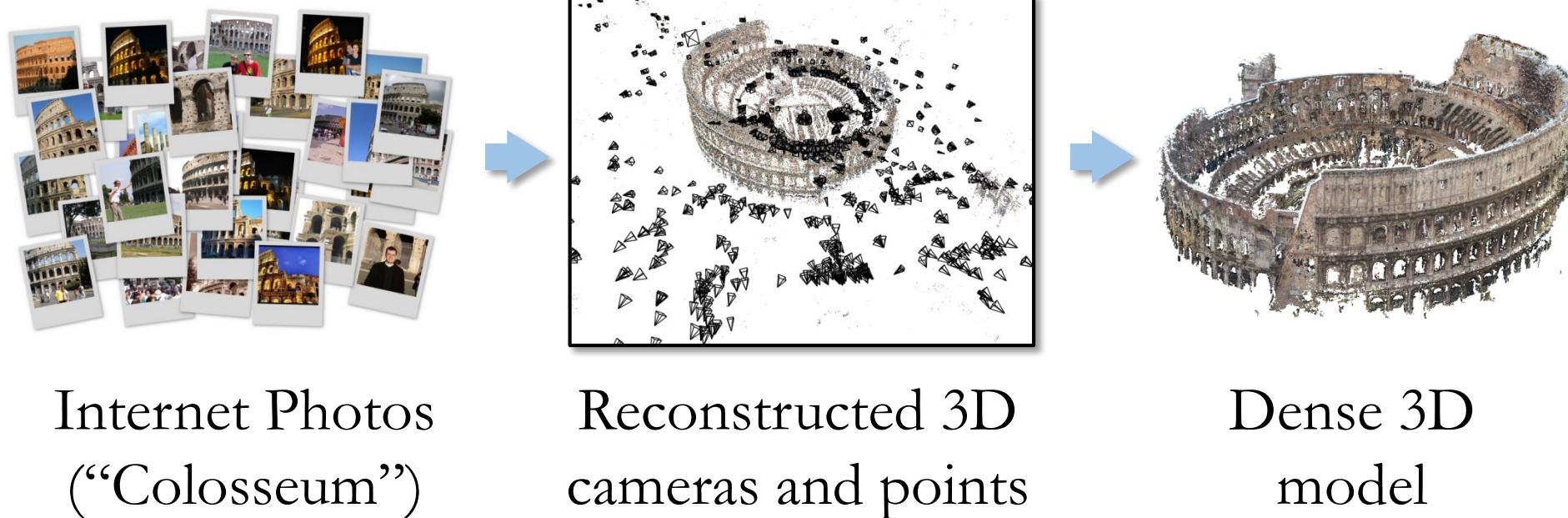
<https://www.hudsandguis.com/home/2011/01/02/terminator-the-evolution-of-machine-vision>

Compute the 3D shape of the world



01. Introduction to Computer Vision

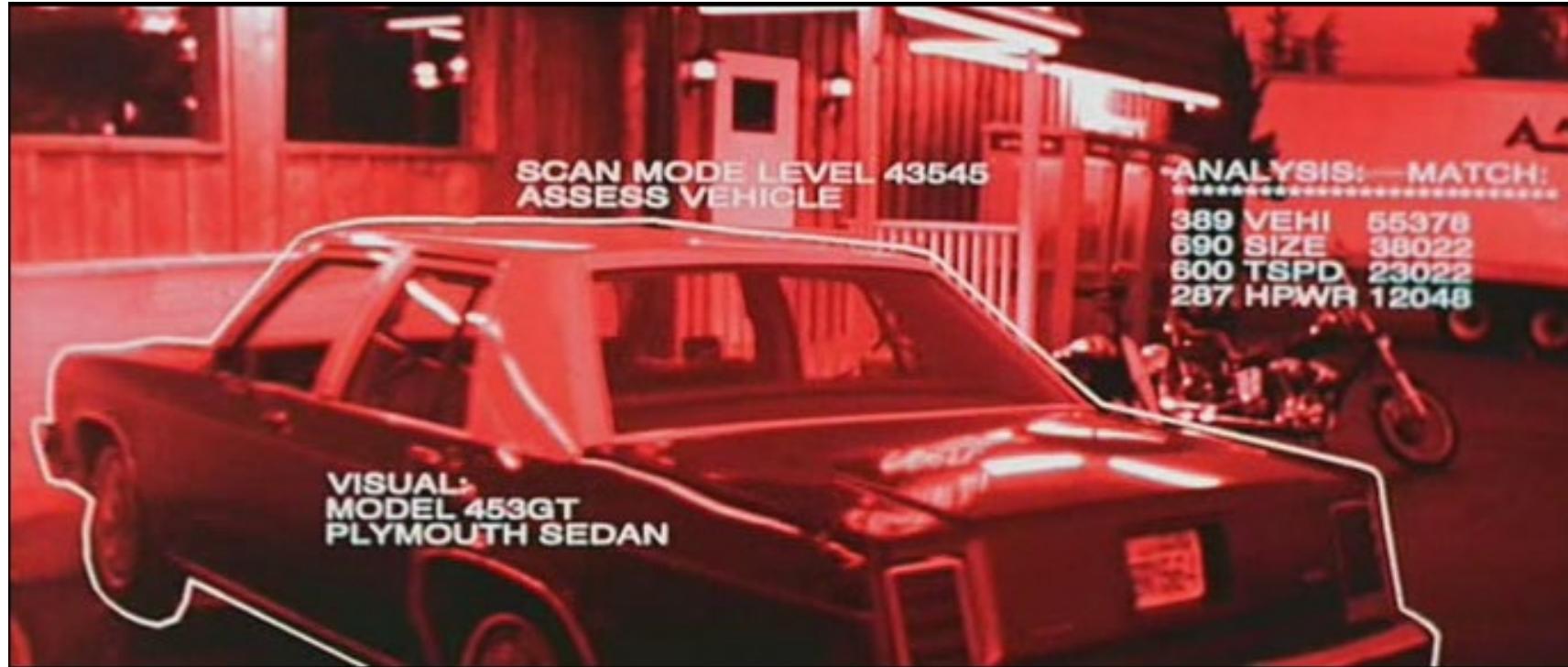
Compute the 3D shape of the world





01. Introduction to Computer Vision

Recognize Objects & People



Terminator 2 (1991)

Objects
Activities
Scenes
Locations
Text / writing
Faces
Gestures
Motions
Emotions...

Recognize Objects & People

sky

Cedar Point

ride

The Wicked Twister

Lake Erie

Ferris wheel

amusement park

tree

tree

water

ride

12 E
-12 E-

deck

tree

bench

tree

carousel

umbrellas

people sitting on ride

maxair

pedestrians

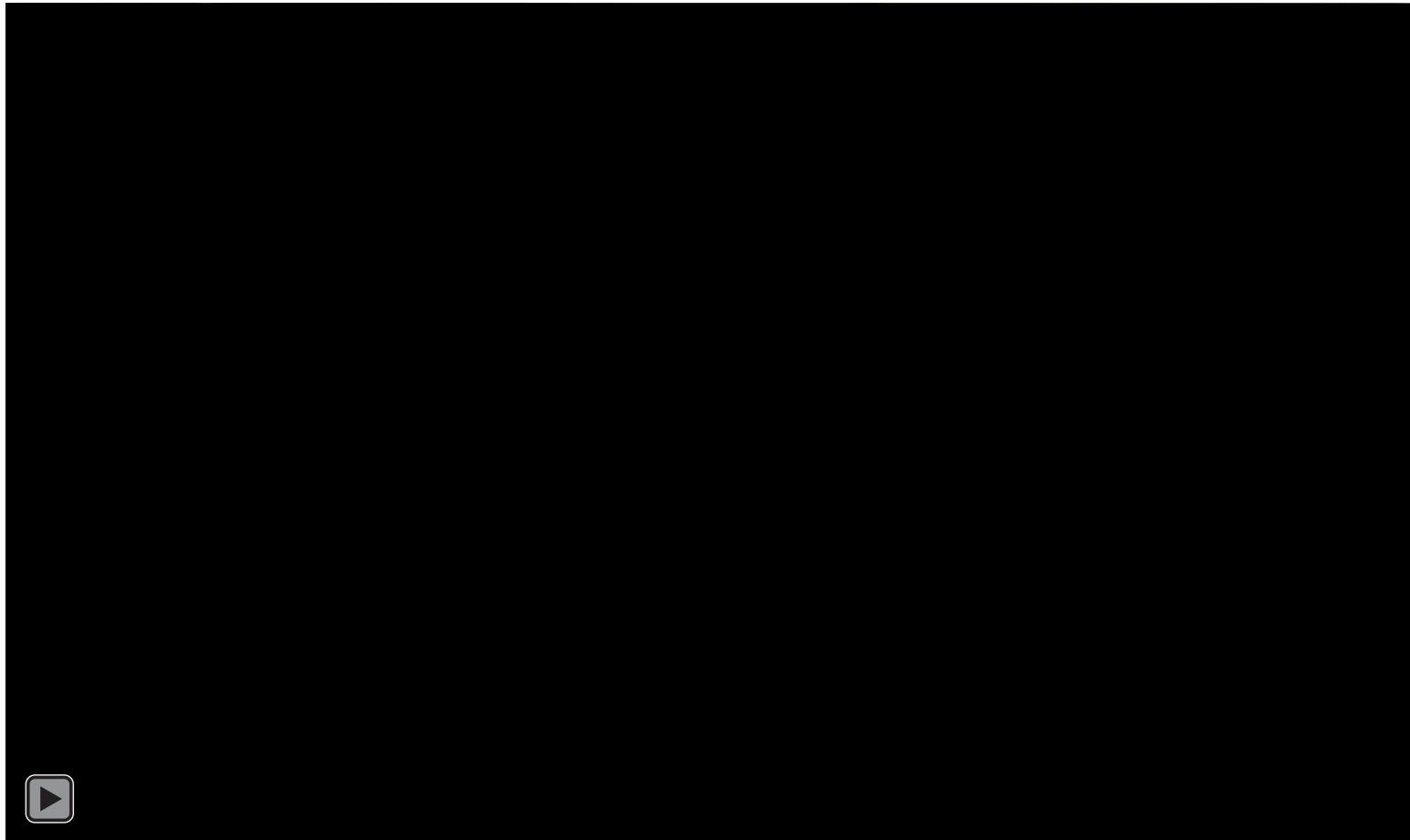


Check out the online web demo!

<https://cs.stanford.edu/people/karpathy/deepimagesent/generationdemo/>

“Enhance” Images

I cannot get the video clip to export correctly for the lecture video. Also the quality is awful... Just watch it on YouTube [[here](#)].



“Enhance” Images



Researchers warn peace sign photos could expose fingerprints

But the likelihood of anyone actually using images to recreate prints is pretty slim.



Jamie Rigg, @jmerigg
01.13.17 in Security

Comments

1721
Shares



01. Introduction to Computer Vision

Improve photos (“Computational Photography”)



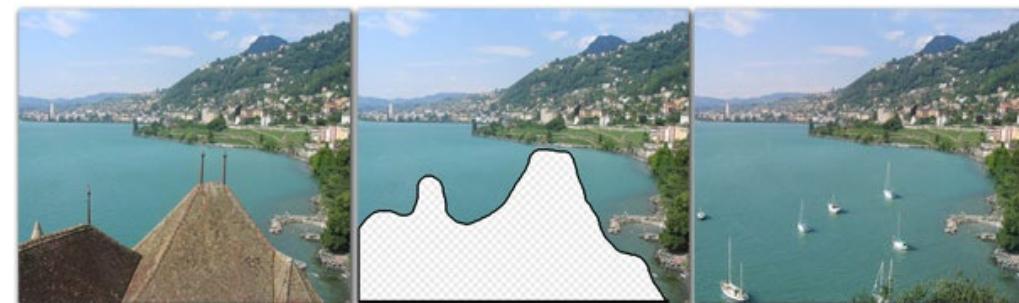
Super-resolution (source: 2d3)



Depth of field on cell phone camera
(source: [Google Research Blog](#))



Low-light photography
(credit: [Hasinoff et al., SIGGRAPH ASIA 2016](#))



Inpainting / image completion
(image credit: Hays and Efros)

01. Introduction to Computer Vision



Source images straight out of the camera (Camera: Nikon D7000; Lens: Nikon 10-24mm; f/11; 1/400s - 1/15s)



Tone-mapped HDR

<http://farbspiel-photo.com/wp-content/gallery/bna-images/hdr-before-and-after-izmir-harbor.jpg>

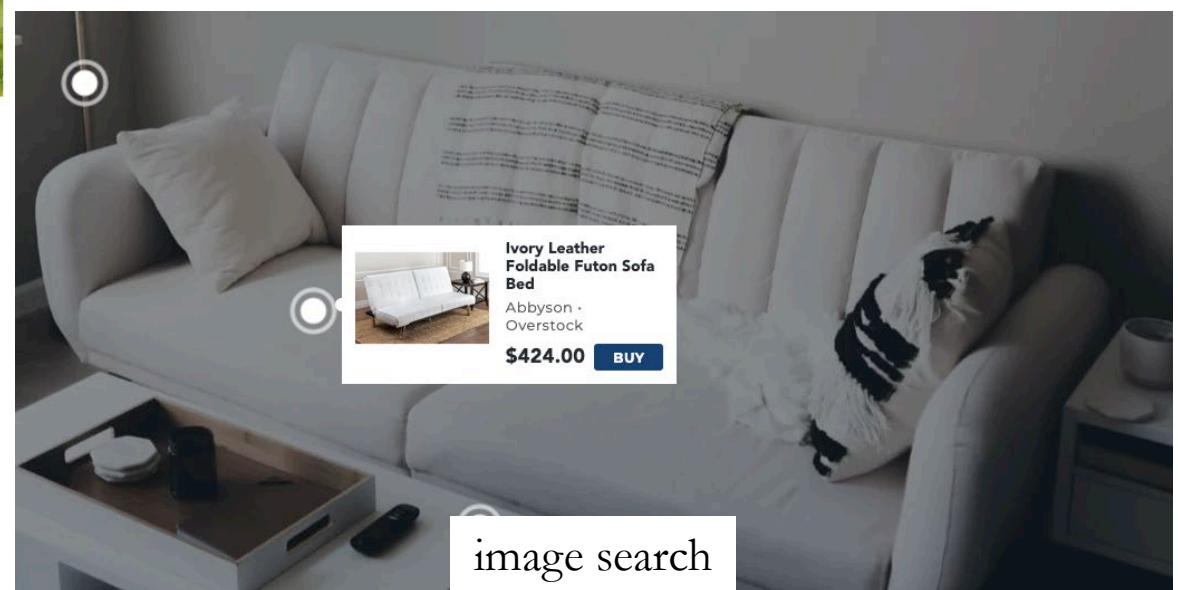
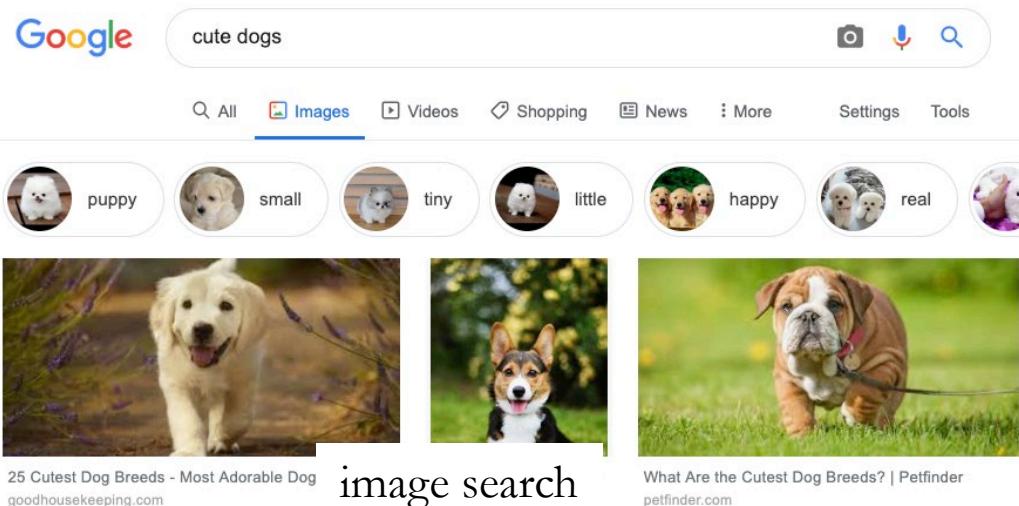


01. Introduction to Computer Vision

Final image after post-processing



Organization & Search



Goals of Computer Vision

1. **Measure**: properties of the 3D world from visual data
2. **Recognize**: algorithms and representations to allow a machine to recognize objects, people, scenes, activities.
3. **Generate**: enhance, complete and manipulate images and videos
4. **Organize**: algorithms to mine, search, and interact with visual data

Deep Learning

CS5242, CS 5260

CS4243
Recognition

Machine
Learning

Robotics

Human
Computer
Interaction

Virtual
Reality

Measure

Medical
Imaging

Optics

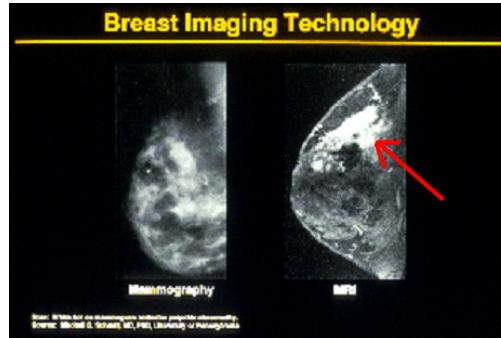
Computational
Photography

Graphics

Why Does Computer Vision Matter?



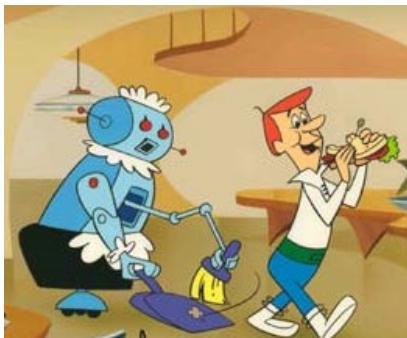
Safety



Health



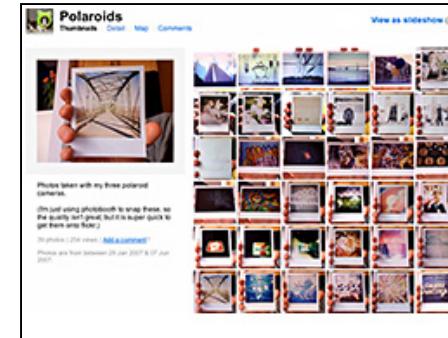
Security



Comfort



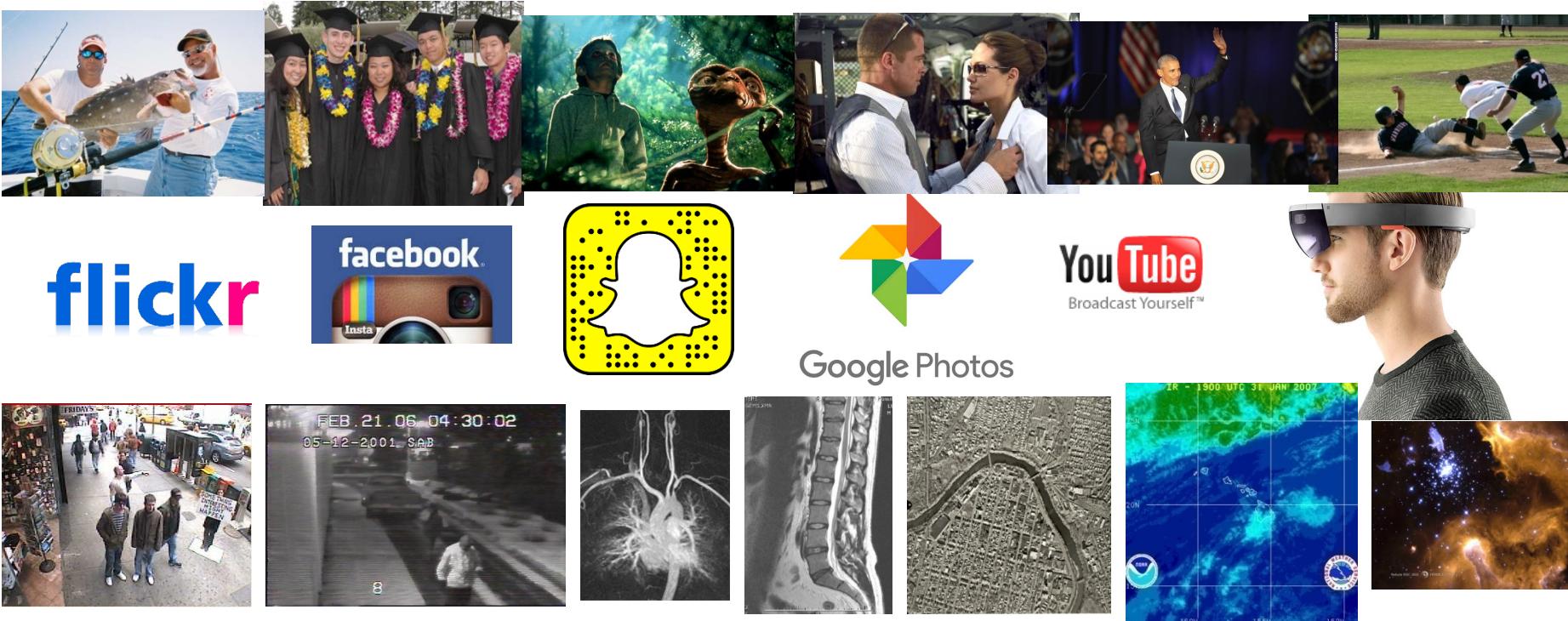
Entertainment



Access

Why Study Computer Vision?

Billions of images/videos captured per day



Huge number of useful applications...

Why Is Computer Vision Hard?

(After all, it comes so naturally to humans.)

Computer Vision is Easy?



Park or Bird? A National Park and Bird Identifying App Inspired by an xkcd Comic

by [E.D.W. Lynch](#) at 5:21 PM on October 21, 2014

[Facebook](#) [Twitter](#) [Pinterest](#) [Tumblr](#) [Flipboard](#) [More](#)

The screenshot shows a user interface for identifying photos. In the center, there's a dashed box containing a photo of a bird standing in front of a fire. Below this box is the text "EXAMPLE PHOTOS" followed by five small thumbnail images. To the right of the dashed box, there are two sections: "PARK or BIRD" and "PARK? ???". The "PARK or BIRD" section contains text about identifying national parks and birds using GPS and computer vision. The "PARK? ???" section contains text about identifying birds without GPS information. To the far right, there's a "BIRD? YES" section with the text "Dude, that is such a bird." At the bottom left, there's a "Photo credits" link.

screenshot via [Flickr](#), photo via [penguinthemagpie](#)

[Park or Bird?](#) is a web app by the [Flickr](#) Vision and Search team that can determine whether a photo depicts a bird or a National Park. The app stems from an [xkcd comic](#) about the often confusing difference between a simple computer task (determining whether a photo was shot in a National Park) and an extremely challenging task (determining whether a photo contains a bird).

Why is computer vision difficult?



Viewpoint variation



Illumination



Scale

Why is computer vision difficult?



Intra-class variation



Background clutter



Motion (Source: S. Lazebnik)



Occlusion

But there are lots of cues we can exploit...



NATIONAL GEOGRAPHIC.COM

© 2003 National Geographic Society. All rights reserved.

Bottom line

- Perception is an inherently ambiguous problem
 - Many different 3D scenes could have given rise to a particular 2D picture



- We often need to use prior knowledge about the structure of the world

Some historical Notes

1966: Papert assigns computer vision as an undergrad summer project

1960's: interpretation of synthetic worlds

1970's: some progress on interpreting selected images

1980's: ANNs come and go; shift toward geometry and increased mathematical rigor

1990's: face recognition; statistical analysis in vogue

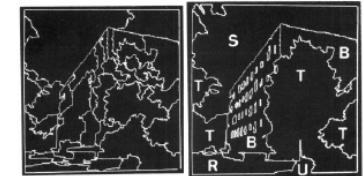
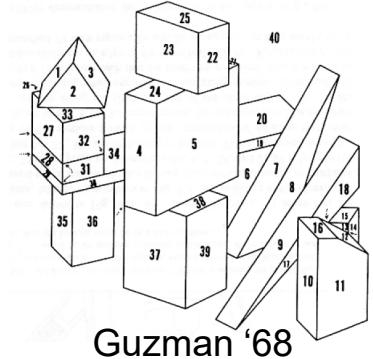
2000's: broader recognition; large annotated datasets available; video processing starts

2010's: Deep learning with CNNs

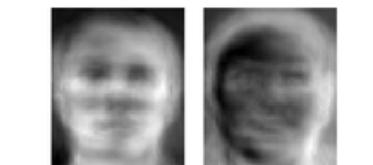
2020's: Widespread autonomous vehicles?

2030's: robot uprising?

Are we here now?



Ohta Kanade '78



Turk and Pentland '91

The state of Computer Vision and AI: we are really, really far.

Oct 22, 2012



Andrej Karpathy blog



01. Introduction

The picture above is funny.

But for me it is also one of those examples that make me sad about the outlook for AI and for Computer Vision. What would it take for a computer to understand this image as you or I do? I challenge you to think explicitly of all the pieces of knowledge that have to fall in place for it to make sense. Here is my short attempt:

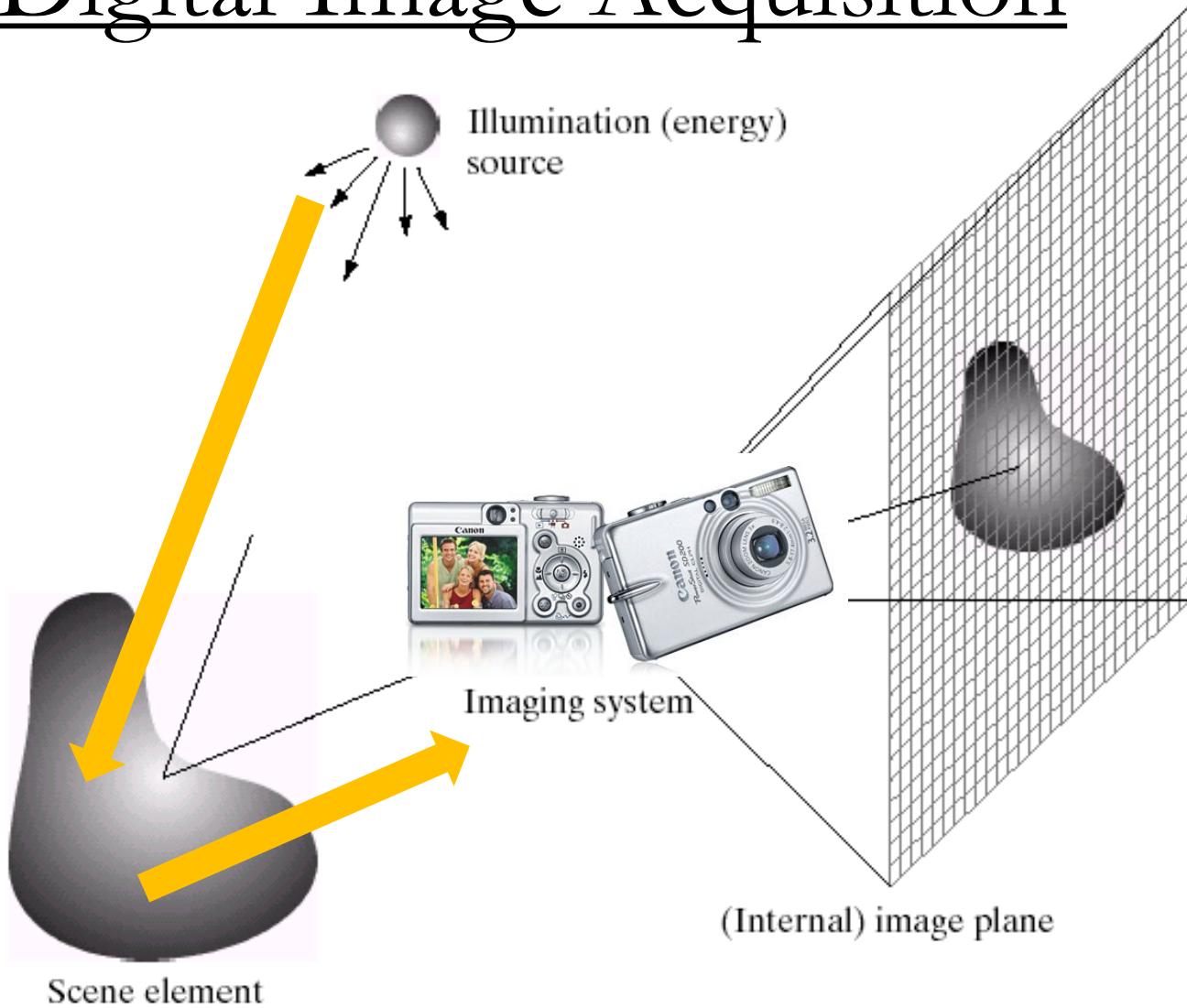
- You recognize it is an image of a bunch of people and you understand they are in a hallway
- You recognize that there are 3 mirrors in the scene so some of those people are "fake" replicas from different viewpoints.
- You recognize Obama from the few pixels that make up his face. It helps that he is in his suit and that he is surrounded by other people with suits.
- You recognize that there's a person standing on a scale, even though the scale occupies only very few white pixels that blend with the background. But, you've used the person's pose and knowledge of how people interact with objects to figure it out.
- You recognize that Obama has his foot positioned just slightly on top of the scale. Notice the language I'm using: It is in terms of the 3D structure of the scene, not the position of the leg in the 2D coordinate system of the image.
- You know how physics works: Obama is leaning in on the scale, which applies a force on it. Scale measures force that is applied on it, that's how it works => it will over-estimate the weight of the person standing on it.
- The person measuring his weight is not aware of Obama doing this. You derive this because you know his pose, you understand that the field of view of a person is finite, and you understand that he is not very likely to sense the slight push of Obama's foot.
- You understand that people are self-conscious about their weight. You also understand that he is reading off the scale measurement, and that shortly the over-estimated weight will confuse him because it will probably be much higher than what he expects. In other words, you reason about implications of the events that are about to unfold seconds after this photo was taken, and especially about the thoughts and how they will develop inside people's heads. You also reason about what pieces of information are available to people.
- There are people in the back who find the person's imminent confusion funny. In other words you are reasoning about state of mind of people, and their view of the state of mind of another person. That's getting frighteningly meta.
- Finally, the fact that the perpetrator here is the president makes it maybe even a little more funnier. You understand what actions are more or less likely to be undertaken by different people based on their status and identity.

<http://karpathy.github.io/2012/10/22/state-of-computer-vision/>

Digital Images

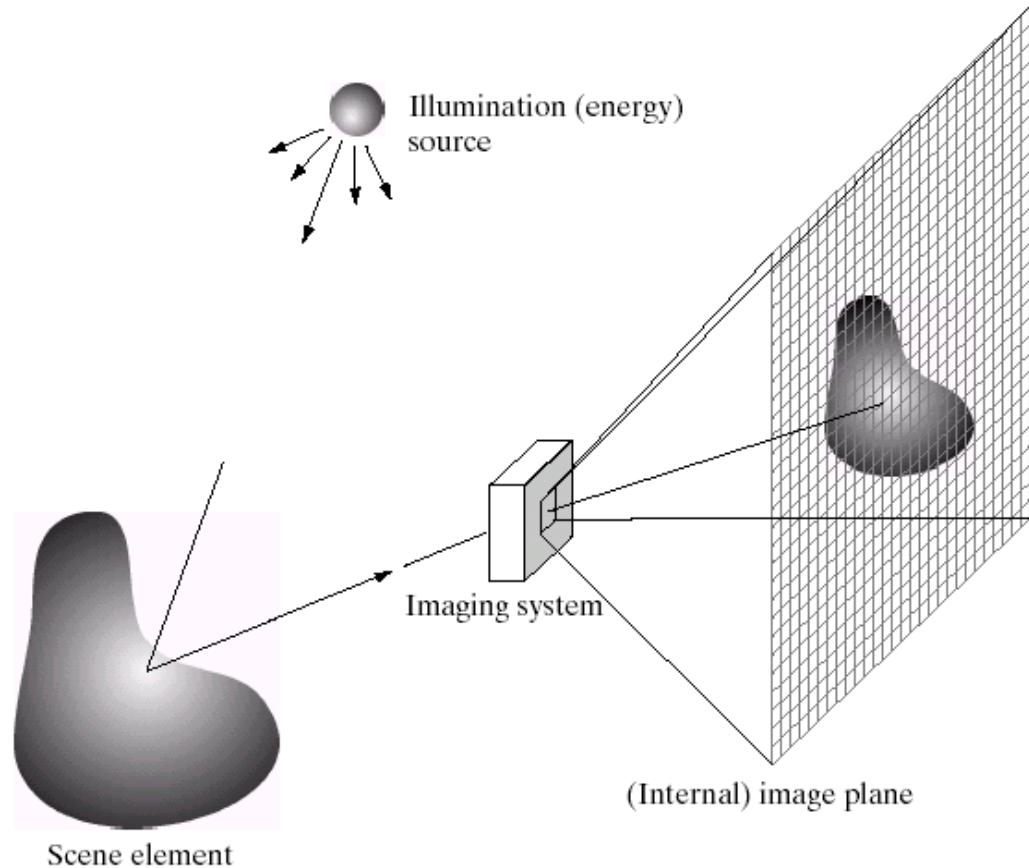
Imaging sensors, exposure time, spatial and intensity resolution

Digital Image Acquisition



1. Energy transfer from source (sun, light bulb) to scene to imaging system
2. Optics in imaging system focuses energy onto sensor (e.g. CCD).
3. Digital sensor measures amount of energy

Image Formation Model

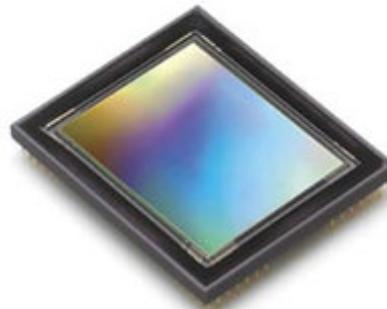


- $f(x, y)$ - 2D image function, where x, y are coordinates on the internal image plane
- Illumination component $i(x, y)$: amount of source illumination incident on the scene
- Reflectance component $r(x, y)$: amount of illumination reflected by objects in the scene

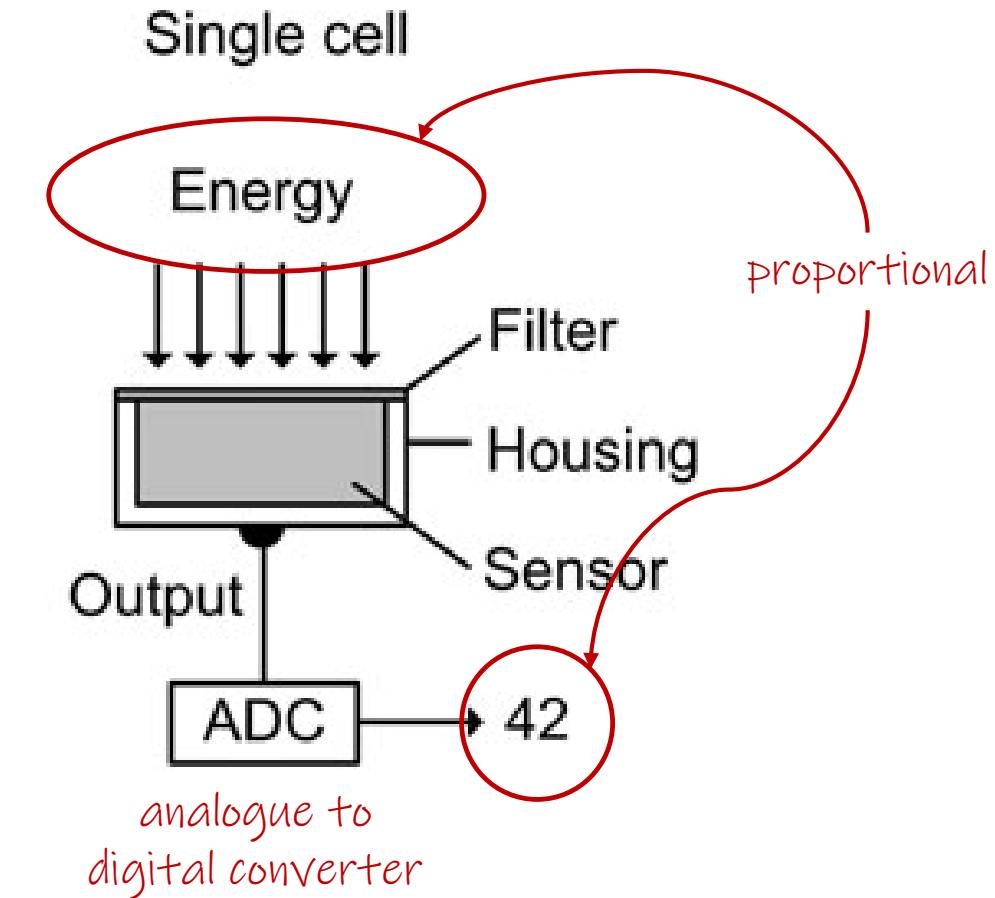
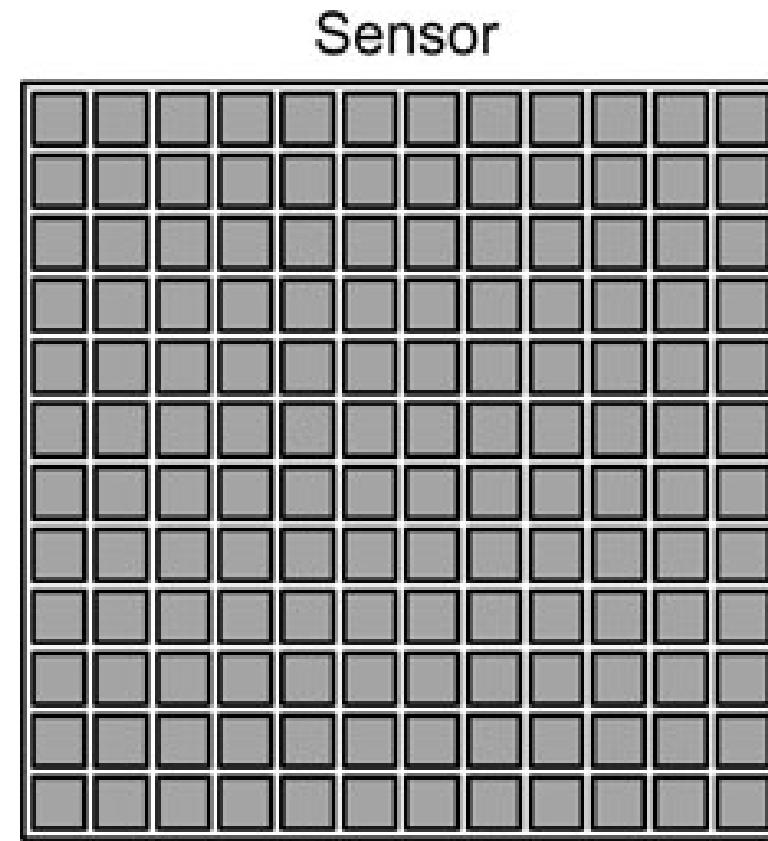
$$f(x, y) = i(x, y)r(x, y)$$

- $0 \leq i(x, y) < \infty$ - no upper bound
- $0 \leq r(x, y) \leq 1$
0 for total absorption, 1 for total reflectance

Imaging System Sensors



CCD sensor



Exposure Time

Exposure is the amount of time that incident light can reach the imaging sensor.

HWQ: Why does motion blur occur and what adjustments should we make to compensate?



Correctly exposed



Over exposed

exposure time too long
too much light reaches sensor



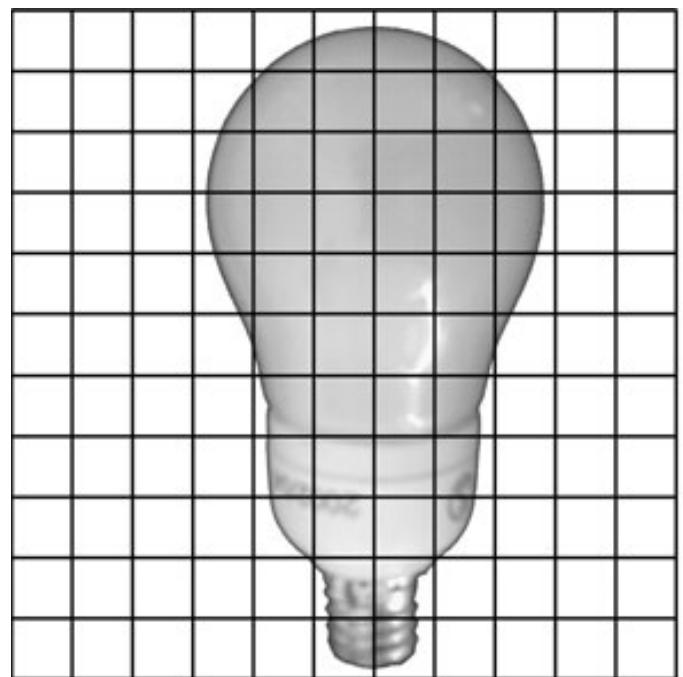
Under exposed

exposure time too short
not enough light reaches sensor

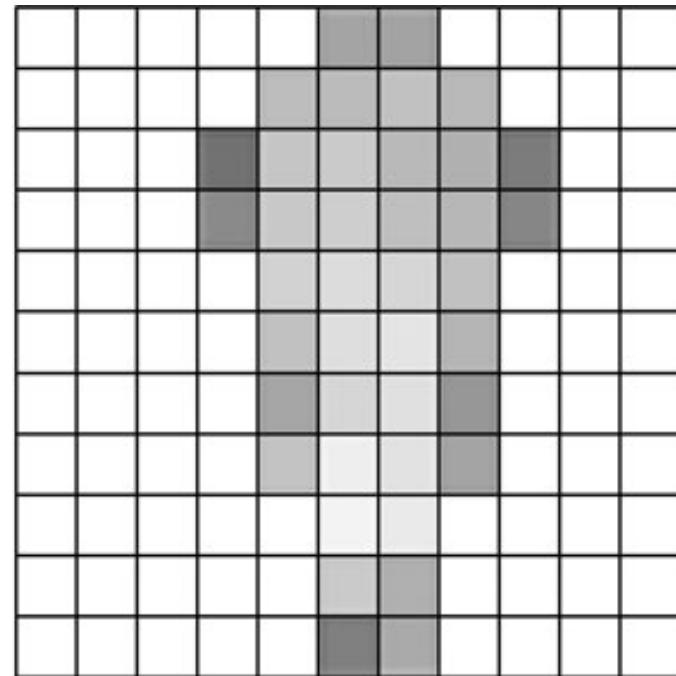


Motion blur

Image Formation



continuous projection
onto imaging sensor

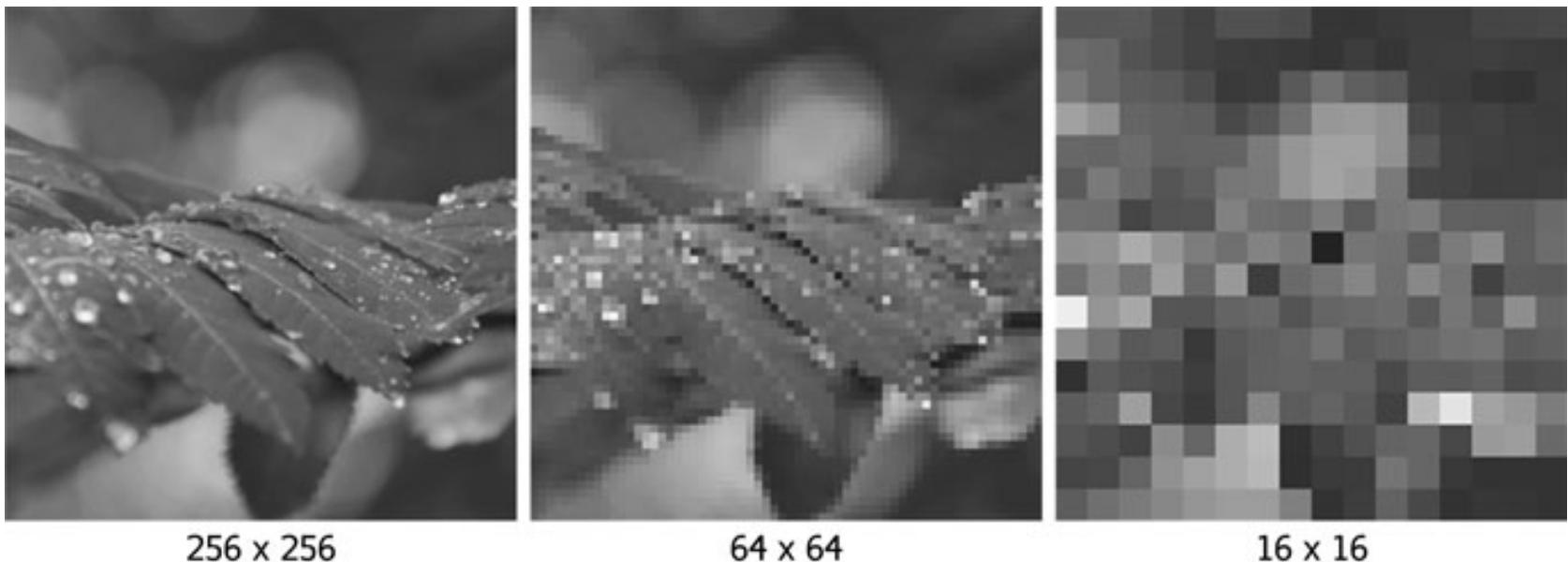


digital sampling

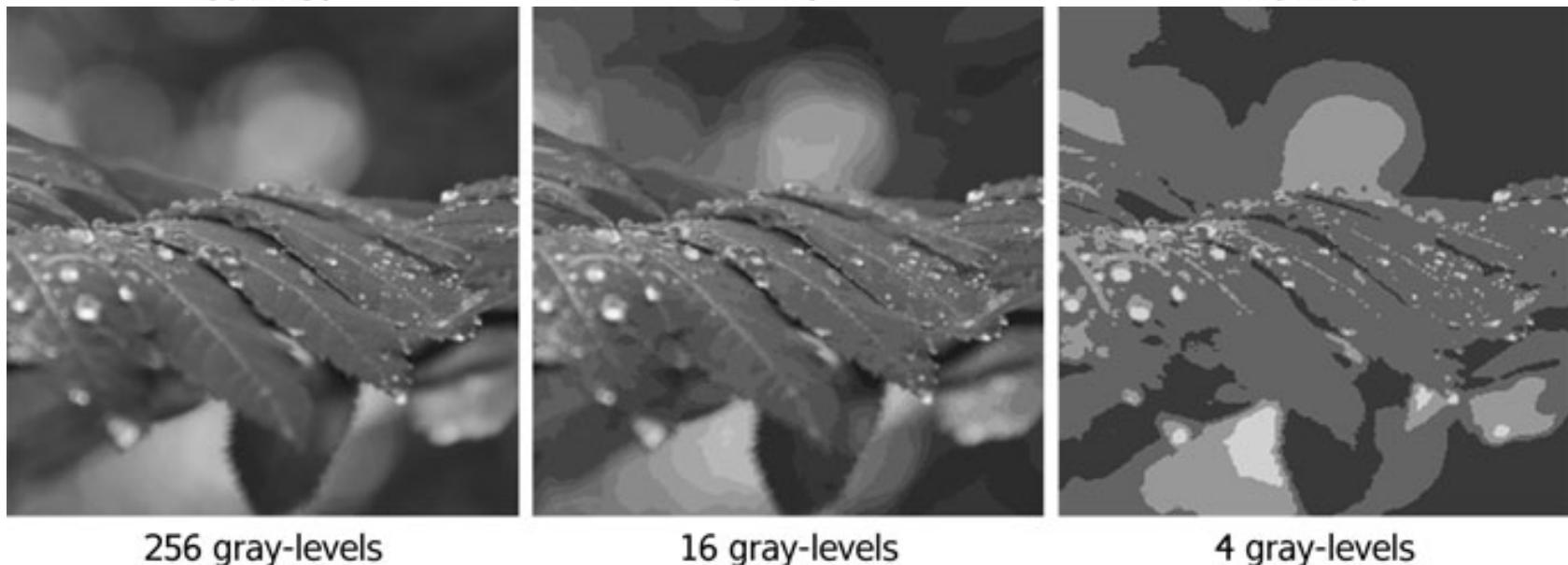
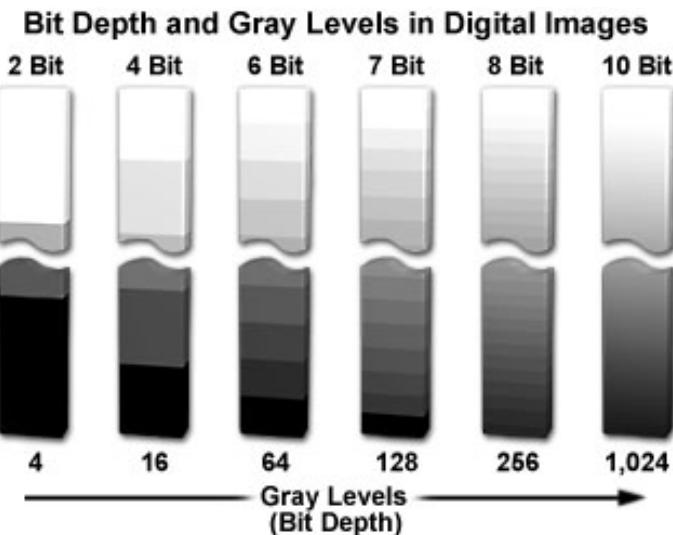
spatial resolution
intensity resolution

Resolution

Spatial Resolution



Intensity Resolution

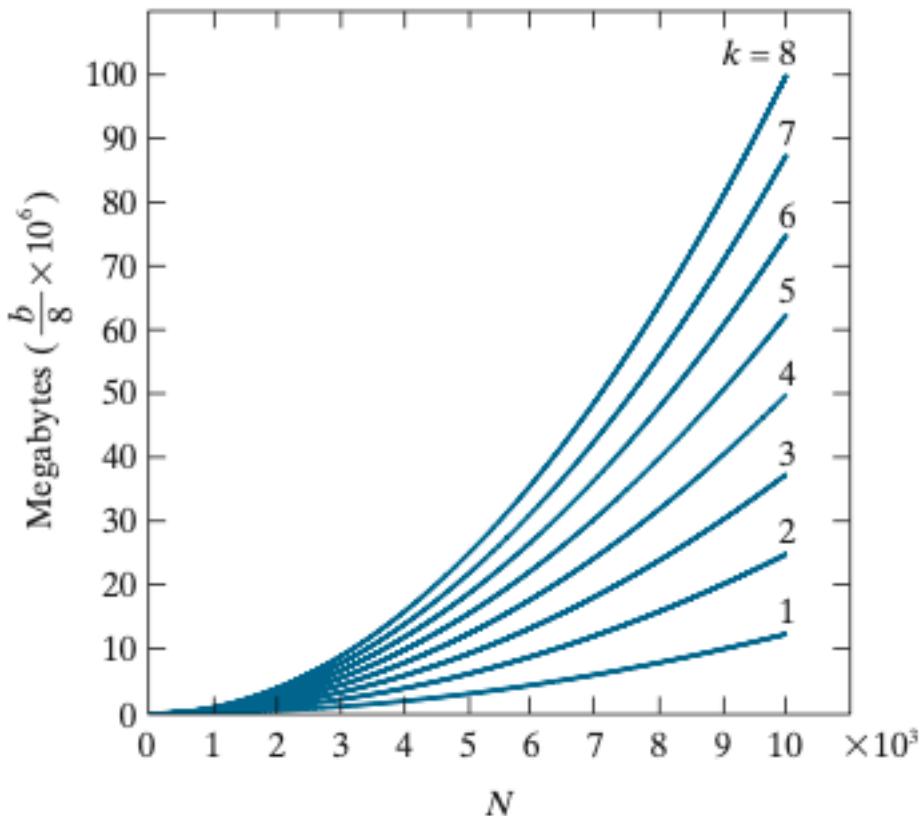


we almost always work with
8 bits, i.e. from 0-255

some visible false contouring

Resolution Determines File Size

Number of megabytes required to store images for various values of N and k .



- Spatial and intensity resolution determines the file size and also transmission time if we were to send images over the internet
- An image of size $M \cdot N$ with a bit-depth of k requires b bits to store, where
$$b = M \cdot N \cdot k$$
- b is the storage requirement for a raw bitmap (.bmp). Compressed formats (.png, .jpg) can reduce file sizes

Colour Representations

RGB images, Bayer patterns, colour to greyscale,
HSV colour space, colour-based applications

Colour Images



Compare

- Sky vs. building in the center
- Top of the fountain vs. clouds
- [0,0,0] is black, [255,255,255] is white

[R, G, B] :[120, 30, 24]



red channel



green channel



blue channel

Images in Python

row ↓ column →

0.92	0.93	0.94	0.97	0.62	0.37	0.85	0.97	0.93	0.92	0.99
0.95	0.89	0.82	0.89	0.56	0.31	0.75	0.92	0.81	0.95	0.91
0.89	0.72	0.51	0.55	0.51	0.42	0.57	0.41	0.49	0.91	0.92
0.96	0.95	0.88	0.94	0.56	0.46	0.91	0.87	0.90	0.97	0.95
0.71	0.81	0.81	0.87	0.57	0.37	0.80	0.88	0.89	0.79	0.85
0.49	0.62	0.60	0.58	0.50	0.60	0.58	0.50	0.61	0.45	0.33
0.86	0.84	0.74	0.58	0.51	0.39	0.73	0.92	0.91	0.49	0.74
0.96	0.67	0.54	0.85	0.48	0.37	0.88	0.90	0.94	0.82	0.93
0.69	0.49	0.56	0.66	0.43	0.42	0.77	0.73	0.71	0.90	0.99
0.79	0.73	0.90	0.67	0.33	0.61	0.69	0.79	0.73	0.93	0.97
0.91	0.94	0.89	0.49	0.41	0.78	0.78	0.77	0.89	0.99	0.93
0.89	0.73	0.58	0.58	0.73	0.72	0.77	0.73	0.71	0.91	0.91
0.79	0.73	0.90	0.67	0.33	0.61	0.69	0.79	0.73	0.93	0.97
0.91	0.94	0.89	0.49	0.41	0.78	0.78	0.77	0.89	0.99	0.93
0.89	0.73	0.58	0.58	0.73	0.72	0.77	0.73	0.71	0.91	0.91
0.79	0.73	0.90	0.67	0.33	0.61	0.69	0.79	0.73	0.93	0.97
0.91	0.94	0.89	0.49	0.41	0.78	0.78	0.77	0.89	0.99	0.93

[0,255] → [0,1] for arithmetic convenience

represented as a matrix

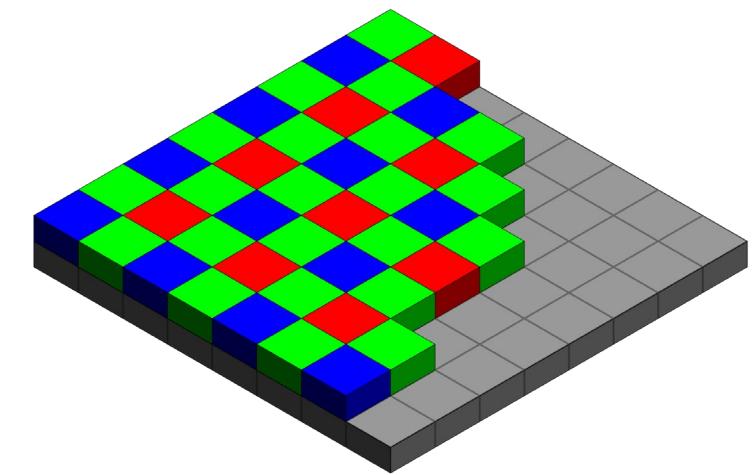
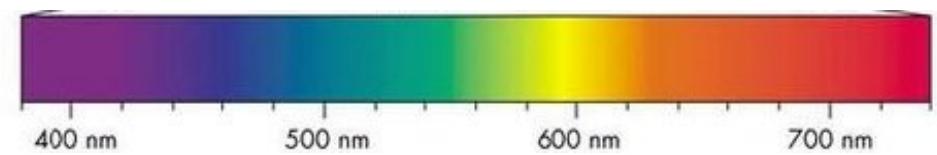
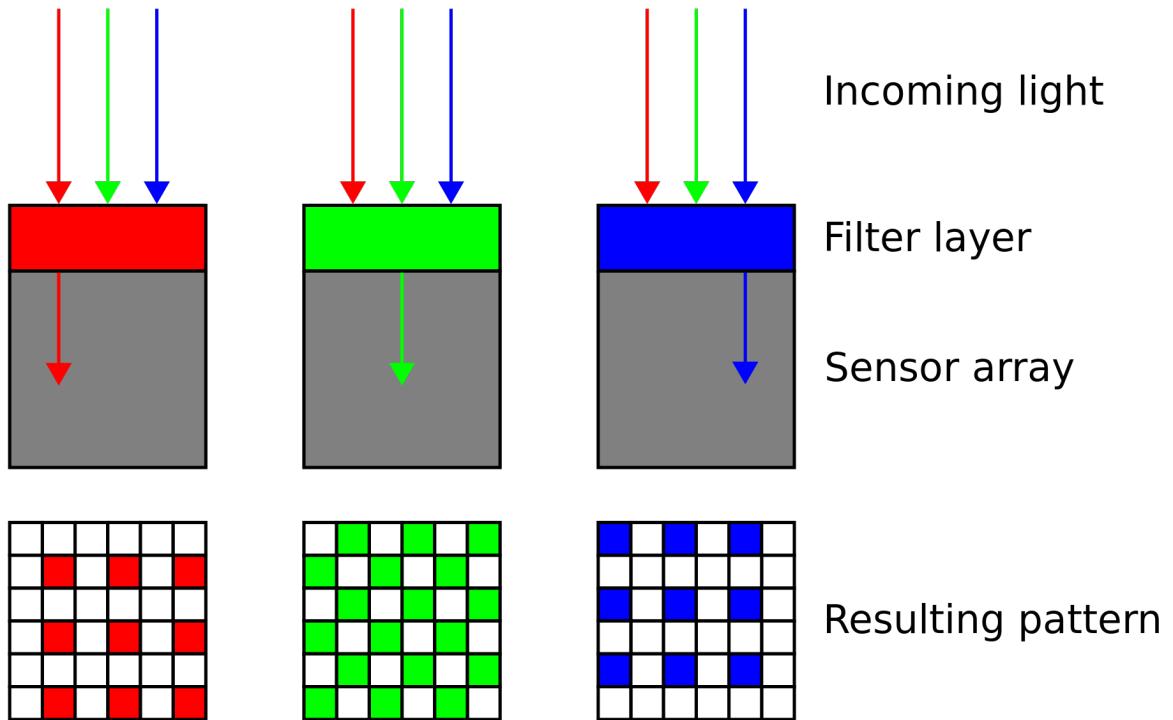
N x M RGB image “im”

$im(0,0,0)$ = top-left pixel value in R-channel

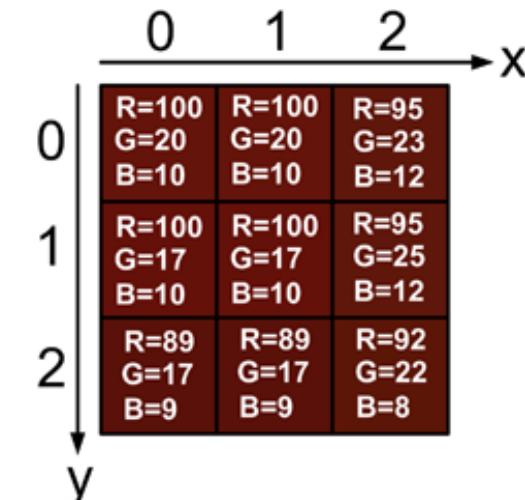
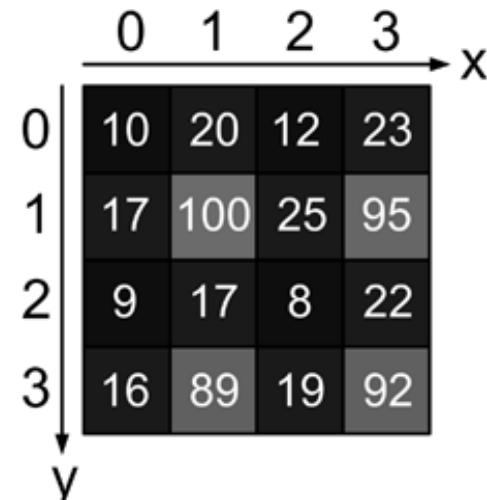
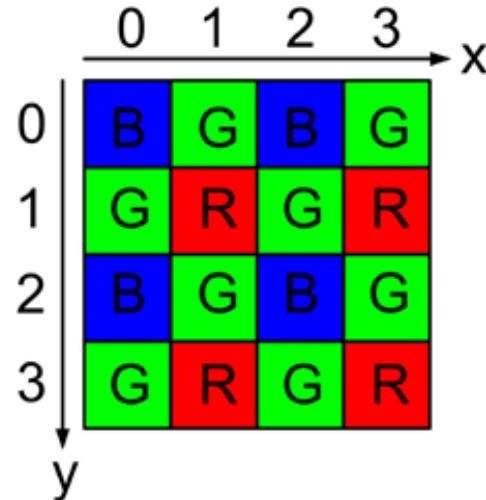
$im(y, x, b)$ = y pixels down, x pixels to right in the b^{th} channel

$im(N-1, M-1, 2)$ = bottom-right pixel in B-channel

Colour Capture – Bayer Filter



Demosaicing Bayer Patterns



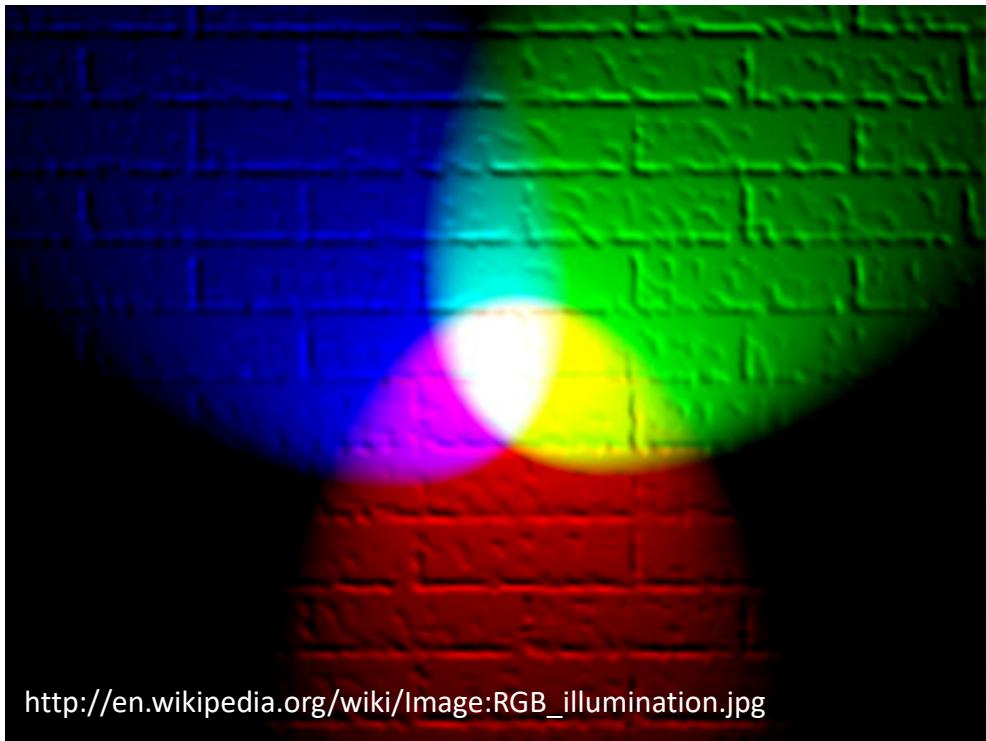
$$g(x, y) = \begin{cases} [R, G, B]_B = [f(x + 1, y + 1), f(x + 1, y), f(x, y)] \\ [R, G, B]_{GB} = [f(x, y + 1), f(x, y), f(x - 1, y)] \\ [R, G, B]_{GR} = [f(x + 1, y), f(x, y), f(x, y - 1)] \\ [R, G, B]_R = [f(x, y), f(x - 1, y), f(x - 1, y - 1)] \end{cases}$$

sensitive to blue

sensitive to green followed by
blue / red in row

sensitive to red

RGB Colour Model



http://en.wikipedia.org/wiki/Image:RGB_illumination.jpg

- “default” colour space
- **additive** model used in electronic displays & photography
- based on **trichromatic** human perception of colours

Colour to Greyscale

greyscale
intensity

Weight factor

$$I = W_R \cdot R + W_G \cdot G + W_B \cdot B$$

$$W_R + W_G + W_B = 1$$

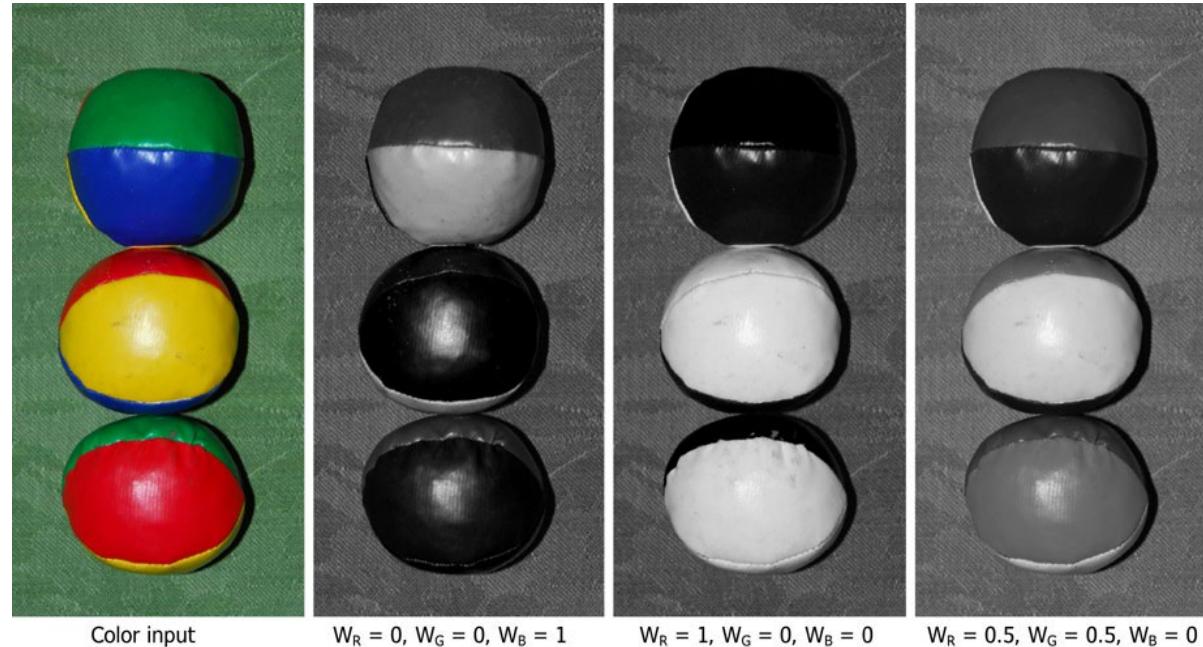
keeps values
within [0, 255]

$$W_R = W_G = W_B = \frac{1}{3}$$

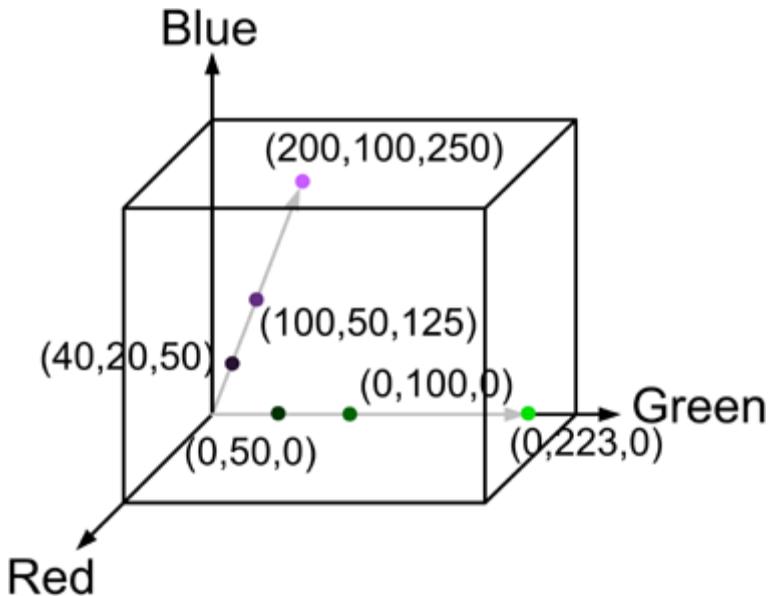
default, but can be adjusted
depending on application

$$W_R = 0.299, \quad W_G = 0.587, \quad W_B = 0.114$$

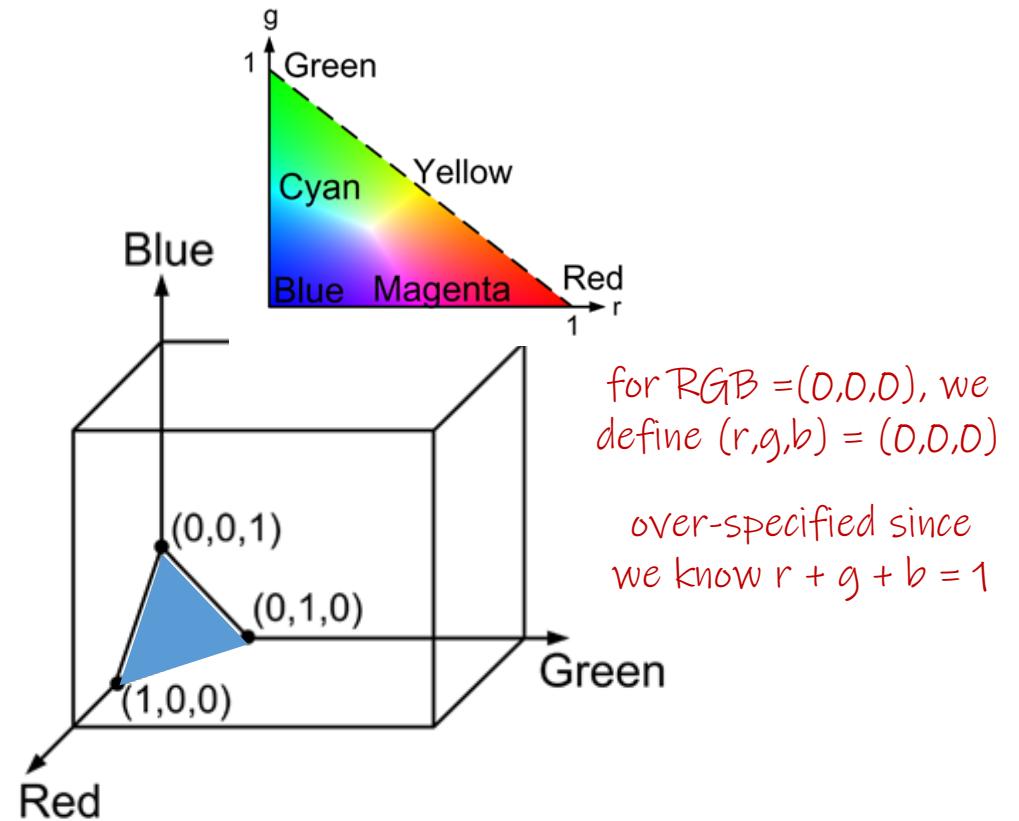
good compromise for
visualization purposes



Normalized RGB



$(0, 50, 0)$, $(0, 100, 0)$, $(0, 223, 0)$ all lie on the same vector spanned by $(0, 255, 0)$ → 3 shades of **green**
same colour, different illumination



$$(r, g, b) = \left(\frac{R}{R+G+B}, \frac{G}{R+G+B}, \frac{B}{R+G+B} \right)$$

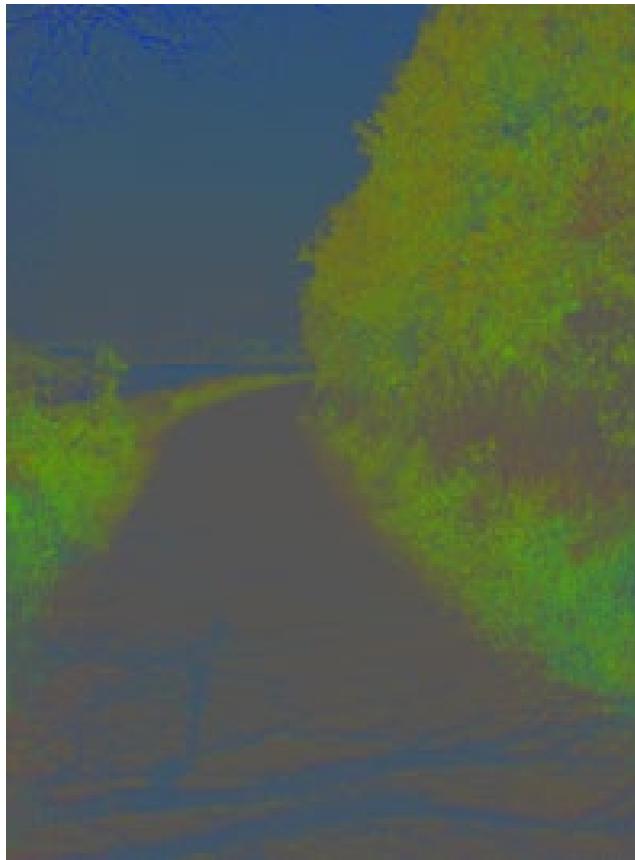
$$(R, G, B) \Leftrightarrow \underbrace{(r, g)}_{\text{true "colour"}}, I \quad I = \frac{R+G+B}{3}$$

Why separate colour and intensity?

Intensity carries more “information”



Original RGB
Image



Only color shown -
Constant intensity



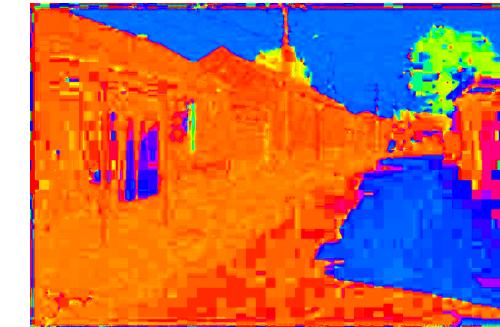
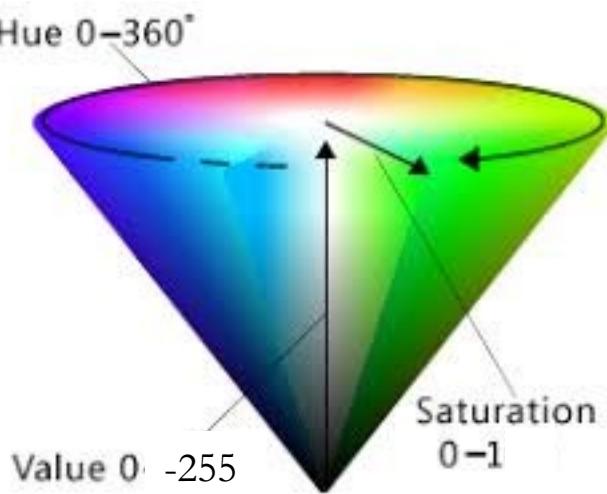
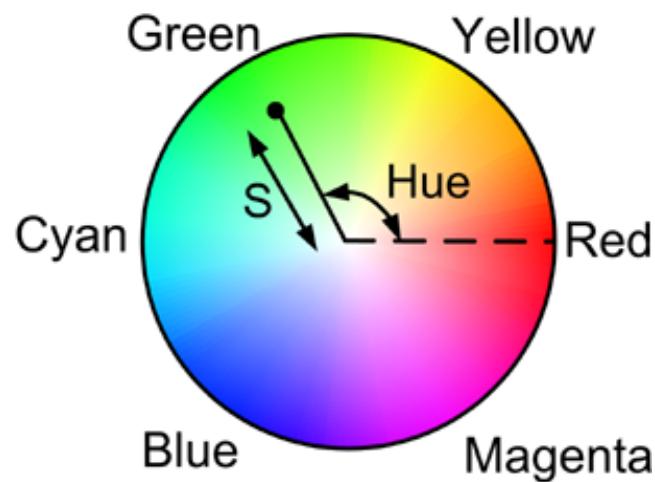
Only intensity shown
– Constant color

HSV Colour Space

Hue: “pure” colour (0 to 360)

Saturation: “purity”, mixing pure colour (1) with white light (0)

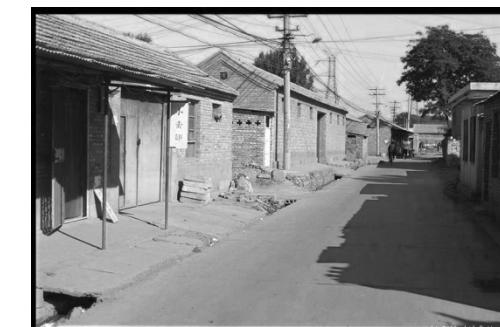
Value: achromatic mixing from black (0) to white (255)



H
(S=1, V=1)



S
(H=1, V=1)



V
(H=1, S=0)

HSV Colour Space

$$H = \begin{cases} \frac{G-B}{V-\min\{R,G,B\}} \cdot 60^\circ, & \text{if } V = R \text{ and } G \geq B; \\ \left(\frac{B-R}{V-\min\{R,G,B\}} + 2\right) \cdot 60^\circ, & \text{if } G = V; \\ \left(\frac{R-G}{V-\min\{R,G,B\}} + 4\right) \cdot 60^\circ, & \text{if } B = V; \\ \left(\frac{R-B}{V-\min\{R,G,B\}} + 5\right) \cdot 60^\circ, & \text{if } V = R \text{ and } G < B \end{cases} \quad H \in [0^\circ, 360^\circ[$$

$$S = \frac{V - \min\{R, G, B\}}{V} \quad S \in [0, 1]$$

$$V = \max\{R, G, B\} \quad V \in [0, 255]$$

What is Colour Good for?

Indexing and Retrieval!

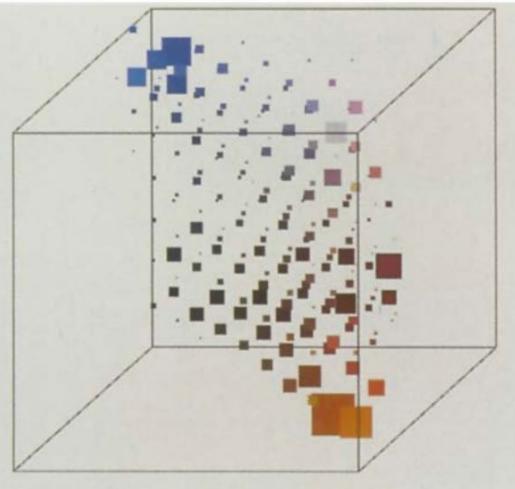
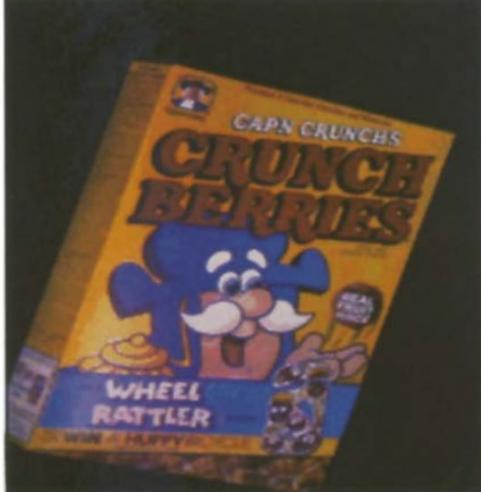
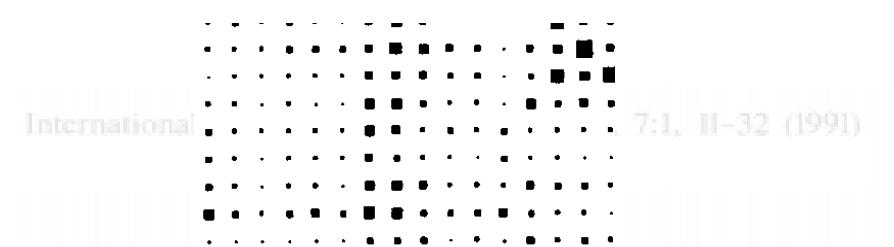


Fig. 1. Left: Image of a Crunchberries cereal box. Right: Three dimensional color histogram of the Crunchberries image with the black background substrated.

need to be developed which run in real time and subserve the task of finding the location of a known object. Color can be successfully used for both tasks.

This article demonstrates that color histograms of multicolored objects provide a robust, efficient cue for indexing into a large database of models. It shows that color histograms are stable object representations in the presence



International

7:1, II-32 (1991)

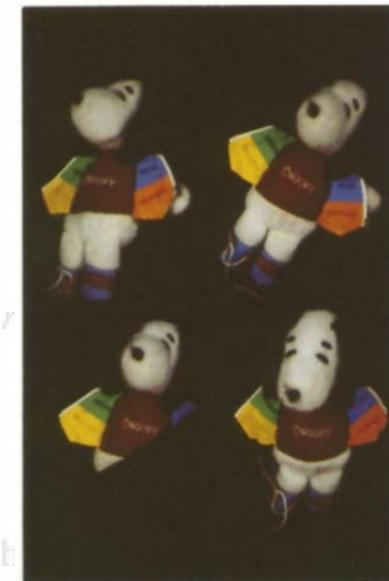


Fig. 2. Four views of Snoopy.

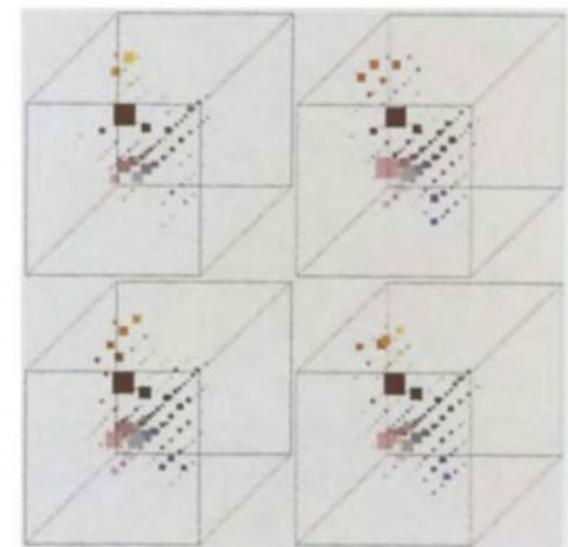


Fig. 3. Histograms of the four views of Snoopy.

Colour-Based Segmentation

HWQ: Why is green preferred for special effects?

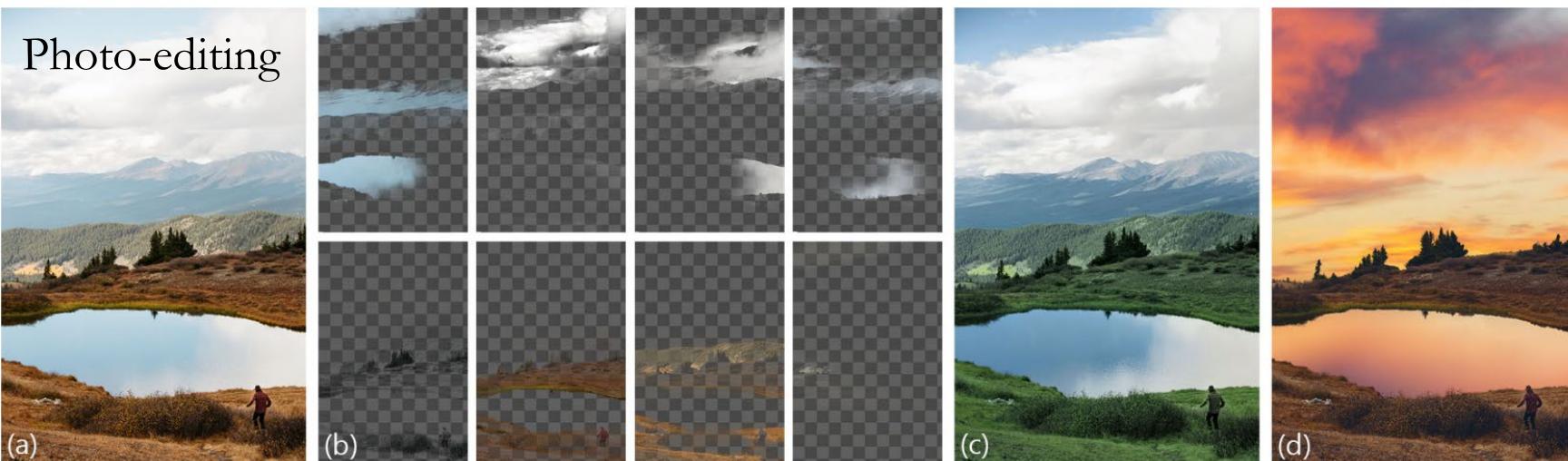
Green-screen for virtual backgrounds



Skin detection



Photo-editing



Summary

1. Computer vision aims to automatically interpret images (and video):
 - Real-world measurements, semantics & recognition
 - Manipulation & generation, organization & search
2. Visual perception is inherently ambiguous, algorithms have a hard time generalizing to all the challenges of the real visual world
3. Digital images
 - Exposure time & gain are two critical parameters of image capture
 - Quantization determines resolution in space and intensity
4. Colour Representations
 - Colour: Bayer mosaic, different colour spaces (RGB, HSV)
 - Intensity more informative than colour
 - Colour can be used for indexing & retrieval, segmentation (e.g. skin)