

**EE 417 Introduction to Computer Vision**  
**/ EE 569 3D Vision**

Instructor: Associate Prof. Mehmet Keskinöz

Communication Theory & Technologies (CTT) Group,  
Electronics Engineering Program,  
Computer Science and Engineering,  
Cyber-Security Program  
Faculty of Engineering and Natural Sciences

Email: keskinoz@sabanciuniv.edu

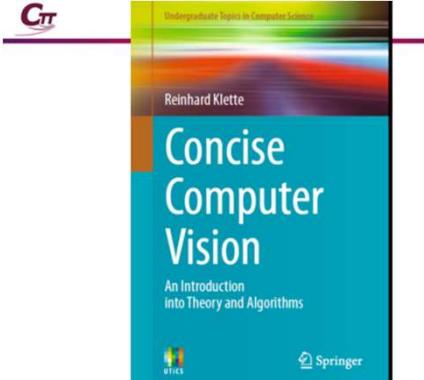
. Sabancı .  
Universitesi

**Assoc. Professor Mehmet Keskinöz**

Email: [keskinoz@sabanciuniv.edu](mailto:keskinoz@sabanciuniv.edu)  
URL: <http://people.sabanciuniv.edu/keskinoz/>

- Education:
  - BS from EE Department of Boğaziçi University, 1994
  - MS (in 1997) and PhD (in 2001) from Electrical and Computer Engineering Carnegie Mellon University (CMU).
- Research Interests
  - High Rate and Robust Communication Receiver Design for Wired and Wireless Technologies such as ADSL, 3G , 4G, 802.11, Wireless Sensor Networks, MIMO etc.
  - Network Coding: Optimal Power Allocation
  - Big Data Analytics& Storage& Communications
  - Information and Decision Fusion; Biometric Fusion
  - Distributed Estimation and Detection
  - Digital Multimedia and Biometric Security; Biometric Hashing and Digital Watermarking Technologies; 2-D Bar Code Design and Its Applications.
  - Control Algorithm Design for Target Detection, Localization and Navigation
  - Correlation Filter Pattern Recognition Theory and Applications ;
  - Error Control Codes (Turbo , LDPC etc) and Their Applications

**Text:** Concise Computer Vision: An Introduction into Theory and Algorithms, Springer, Series:Undergraduate Topics in Computer Science, by Reinhard Klette, 2014. ISBN: 978-1-4471-6319-0



**Reference Books:**  
Computer Vision: Algorithms and Applications, R. Szeliski, Springer, 2010.  
“Fundamentals of Digital Image Processing,” by Anil K. Jain, Englewood Cliffs, NJ : Prentice Hall, 1989. ISBN: 0133361659

**What is Vision?**

“Vision is the act of knowing what is where by looking.” --Aristotle

Special emphasis: relationship between 3D world and a 2D image. Location and identity of objects.

**What is Computer Vision?**

**It is related, but not equivalent to:**

- Photogrammetry
- Image Processing
- Artificial Intelligence

## Computer Vision

- The ability of computers to see
  - Image Understanding
  - Machine Vision
  - Robot Vision
  - Image Analysis
  - Video Understanding

## Why study Computer Vision?

- Images and movies are everywhere
- Fast-growing collection of useful applications
  - building representations of the 3D world from pictures
  - automated surveillance (who's doing what)
  - movie post-processing
  - face finding
- Various deep and attractive scientific mysteries
  - how does object recognition work?
- Greater understanding of human vision

## Goals and Objectives

- Introduce the fundamental problems of computer vision.
- Introduce the main concepts and techniques used to solve those problems.
- Enable one to implement vision algorithms
- Enable one to make sense of the vision literature

## A picture is worth a thousand words.



**A picture is worth 1000 words.**

**A video is worth 1000 sentences?**



<http://marsrovers.jpl.nasa.gov/gallery/press/opportunity/20040125a.html>  
JPL Mars' Panorama captured by the Opportunity

- Rich info. from visual data
- Examples of images around us
  - natural photographic images; artistic and engineering drawings
  - scientific images (satellite, medical, etc.)
- “Motion pictures” => video
  - movie, TV program; family video; surveillance and highway/ferry camera

## Why Do We Process Images?

- Enhancement and restoration
  - Remove artifacts and scratches from an old photo/movie
  - Improve contrast and correct blurred images
- Composition (for magazines and movies), **Display, Printing ...**
- Transmission and storage
  - images from oversea via Internet, or from a remote planet
- Information analysis and automated recognition
  - Providing “human vision” to machines
- Medical imaging for diagnosis and exploration
- Security, forensics and rights protection
  - Encryption, hashing, digital watermarking, digital fingerprinting ...



# Digital Data or Media

Image      Speech      Text      Video



Ls. (.A.dministra  
tors-.O.n.Line)  
People who mon  
itor and oversee  
online chat roo  
ms and online ga  
mes. In other w  
ords, the "onlin  
e police.".....A



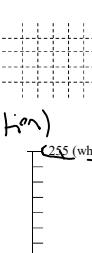
Digitized Image

Analog Image

They are all a sequence (or array) of numbers with finite precision!

**CIT** Sampling and Quantization

- Computer handles “discrete” data.
- Sampling
  - Sample the value of the image at the nodes of a regular grid on the image plane.
  - A pixel (picture element) at  $(i, j)$  is the image intensity value at grid point indexed by the integer coordinate  $(i, j)$ .
- Quantization (rounding off operation)
  - Is a process of transforming a real valued sampled image to one taking only a finite number of distinct values.
  - Each sampled value in a 256-level grayscale image is represented by 8 bits.



77 = 3.14159

C<sub>π</sub>

Why Digital?

- “Exactness”
  - Perfect reproduction without degradation
  - Perfect duplication of processing result
- Convenient & powerful computer-aided processing
  - Can perform sophisticated processing through computer hardware or software
  - Even kindergartners can do some!
- Easy storage and transmission
  - 1 CD can store hundreds of family photos!
  - Paperless transmission of high quality photos through network within seconds

CIT

Image Sampling

- The sampling theorem applies to 2D signal (images) too.

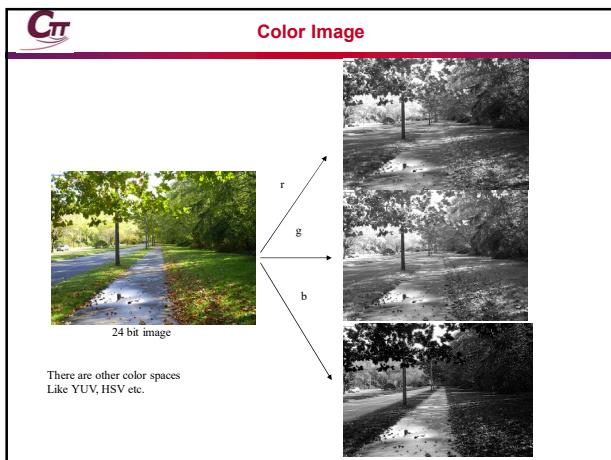
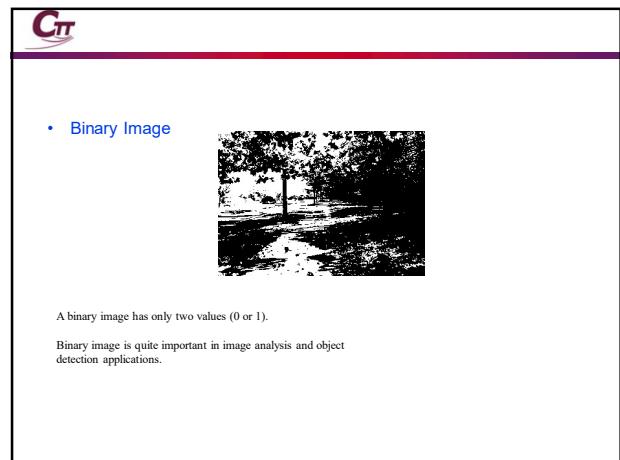
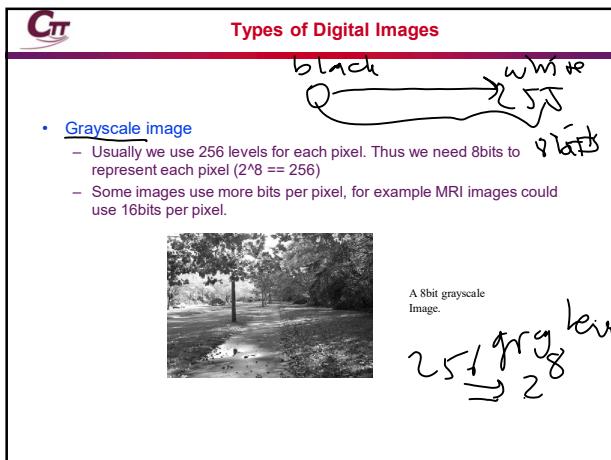
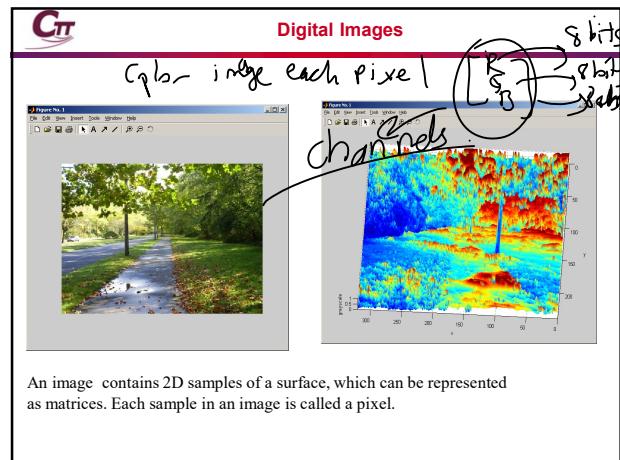
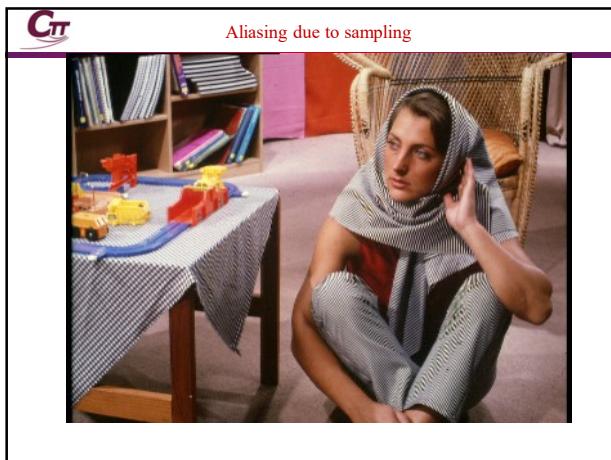
A portrait of a woman wearing a hat, sampled onto a 5x5 grid of blue dots. The grid is centered on her face, with dots at each corner and in the middle of each side.

Sampling on a grid

A 5x5 grid of blue dots. Two blue triangles, representing sinc functions, are drawn over the dots at the bottom center and the right center. These triangles overlap several dots, illustrating how multiple samples can be taken from a single sinc function, leading to aliasing.

Sampling problem

A photograph of a woman sitting cross-legged on a carpeted floor. She is wearing a patterned headscarf and matching pants. Her hands are resting on her knees. In front of her is a low wooden table covered with a blue and white checkered cloth. On the table are several colorful plastic toy cars and trucks. In the background, there is a bookshelf filled with books and a small bed with a pink blanket. The overall atmosphere is domestic and suggests a moment of quiet in a family home.

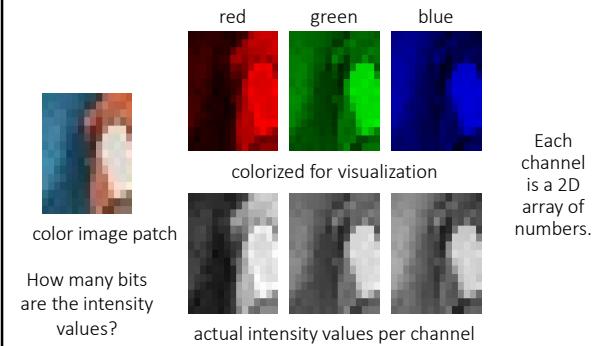


## What is an image?



A (color) image is a 3D tensor of numbers.

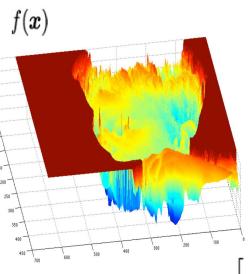
## What is an image?



## What is an image?



grayscale image  
What is the range of the image function  $f$ ?



$$\text{domain } \mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$$

A (grayscale) image is a 2D function.



## RGB Primaries and Color Representation

- Use red, green, blue light to represent a large number of visible colors
- The contribution from each primary is normalized to [0, 1]

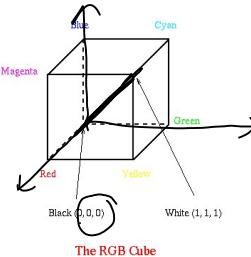
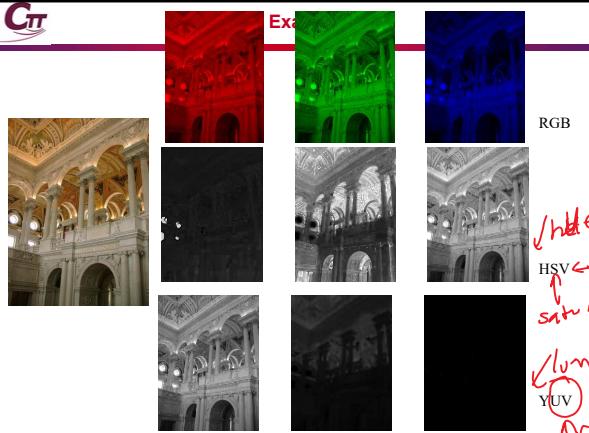


FIGURE 6.8 RGB 24-bit color cube.

Color-cube figures: left figure is from B.Liu ELE330 S'01 lecture notes @ Princeton, right figure is from slides at Gonzalez/ Woods DIP book website



## Example



## Color Coordinates Used in TV Transmission

- Facilitate sending color video via 6MHz mono TV channel
- YIQ for NTSC (National Television Systems Committee) transmission system

- Use receiver primary system ( $R_N$ ,  $G_N$ ,  $B_N$ ) as TV receivers standard
- Transmission system use (Y, I, Q) color coordinate
  - Y ~ luminance, I & Q ~ chrominance

$$\begin{bmatrix} Y \\ I \\ Q \end{bmatrix} = \begin{bmatrix} 0.299 & 0.587 & 0.114 \\ 0.596 & -0.275 & -0.321 \\ 0.212 & -0.523 & 0.311 \end{bmatrix} \begin{bmatrix} R_N \\ G_N \\ B_N \end{bmatrix}, \quad \begin{bmatrix} Y \\ U \\ V \end{bmatrix} = \begin{bmatrix} 0.299 & 0.587 & 0.114 \\ -0.147 & -0.289 & 0.436 \\ 0.615 & -0.515 & -0.100 \end{bmatrix} \begin{bmatrix} R_P \\ G_P \\ B_P \end{bmatrix}.$$

- YUV (YCbCr) for PAL and digital video
  - Y ~ luminance, Cb and Cr ~ chrominance

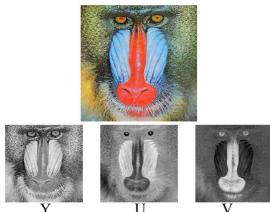
## Color System in Video

- YUV was used in PAL (an analog video standard) and also used for digital video.
- Y is the luminance component (brightness)  

$$Y = 0.299 R + 0.587 G + 0.144 B$$
- U and V are color components  

$$U = B - Y$$
  

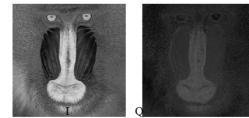
$$V = R - Y$$



## Color Image

- YIQ is the color standard in NTSC.

$$\begin{bmatrix} Y' \\ I \\ Q \end{bmatrix} = \begin{bmatrix} 0.299 & 0.587 & 0.144 \\ 0.595879 & -0.274133 & -0.321746 \\ 0.211205 & -0.523083 & 0.311878 \end{bmatrix} \begin{bmatrix} R' \\ G' \\ B' \end{bmatrix}$$

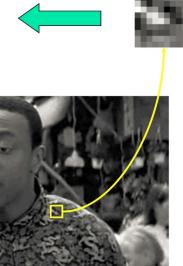


- YCbCr: A color system used in JPEG.

$$\begin{bmatrix} Y' \\ C_b \\ C_r \end{bmatrix} = \begin{bmatrix} 65.481 & 128.553 & 24.966 \\ -37.797 & -74.203 & 112 \\ 112 & -93.786 & -18.214 \end{bmatrix} \begin{bmatrix} R' \\ G' \\ B' \end{bmatrix} + \begin{bmatrix} 16 \\ 128 \\ 128 \end{bmatrix}$$

- 2-D array of numbers (intensity values, gray levels)
- Gray levels 0 (black) to 255 (white)
- Color image is 3 2-D arrays of numbers
  - Red
  - Green
  - Blue
- Resolution (number of rows and columns)
  - 128X128
  - 256X256
  - 512X512
  - 640X480

58	89	116	97	88	83	83	81
23	57	75	89	88	89	89	81
26	57	62	58	75	81	81	75
21	98	65	1	22	80	81	84
147	65	114	73	48	98	73	81
115	98	155	114	48	28	48	73
74	81	73	94	73	81	89	89
45	82	92	58	73	81	82	82
81	103	106	97	89	82	82	82
89	86	89	97	81	78	82	97



## Video

- Sequence of frames
- 30 frames per second

## Video Clip



**C<sub>π</sub>**

## Sequence of Images

**C<sub>π</sub>**

## Why is Computer Vision Hard?

### Visualizing Images

Recall two ways of visualizing an image

Intensity pattern	2d array of numbers

We “see it” at this level      Computer works at this level

**C<sub>π</sub>**

## Why is Computer Vision Hard?

We are trying to infer things about the world from an array of numbers

Shoulder of a cow...

problems: too local; lack of context.  
But wait, it's even worse than that...

**C<sub>π</sub>**

## Why is Computer Vision Hard?

If we already know the geometry, surface material and lighting conditions, it is well-understood how to generate the value at each pixel. [this is Computer Graphics]

But this confluence of factors contributing to each pixel can not be easily decomposed. The process can not be inverted.

**C<sub>π</sub>**

## More Difficulties

Object appearance changes with respect to viewpoint

**C<sub>π</sub>**

## Effects of Lighting

Object appearance also varies with respect to lighting magnitude and direction

**C<sub>TR</sub>**

## Why is Computer Vision Hard?

We are trying to infer things about the world from an array of numbers

Putdata: /home/camps/covgray.jpg															
File	146	161	185	159	165	172	166	142	143	141	149	154	152	149	158
	148	154	152	149	158	173	164	147	144	141	149	157	167	161	151
	147	146	145	148	157	160	151	139	140	138	149	157	167	161	151
	147	149	157	167	167	159	139	129	130	132	149	157	167	161	151
	148	155	167	176	163	150	135	131	131	131	155	167	176	163	150
	138	146	152	155	148	139	133	133	133	134	146	152	155	148	139
	131	132	132	131	132	133	131	127	130	132	132	132	132	131	132
	133	132	129	127	134	141	134	122	125	127	127	127	127	127	127
	129	127	128	128	131	132	130	127	128	127	127	127	127	127	127
	129	127	128	128	131	132	130	128	130	129	129	129	129	129	129

Shoulder  
of a cow...

problems: too local; lack of context.  
But wait, it's even worse than that...

**C<sub>TR</sub>**

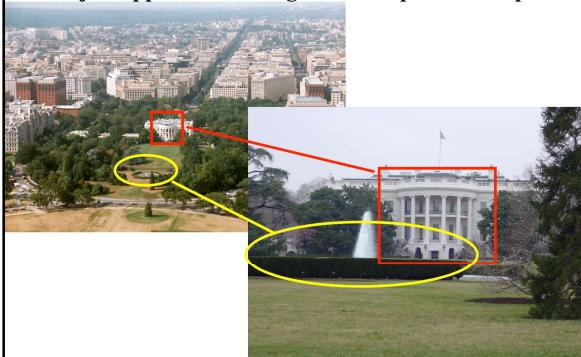
## Why is Computer Vision Hard?

If we already know the geometry, surface material and lighting conditions, it is well-understood how to generate the value at each pixel. [this is Computer Graphics]

But this confluence of factors contributing to each pixel can not be easily decomposed. The process can not be inverted.

## More Difficulties

Object appearance changes with respect to viewpoint



## Effects of Lighting



Object appearance also varies with respect to lighting magnitude and direction

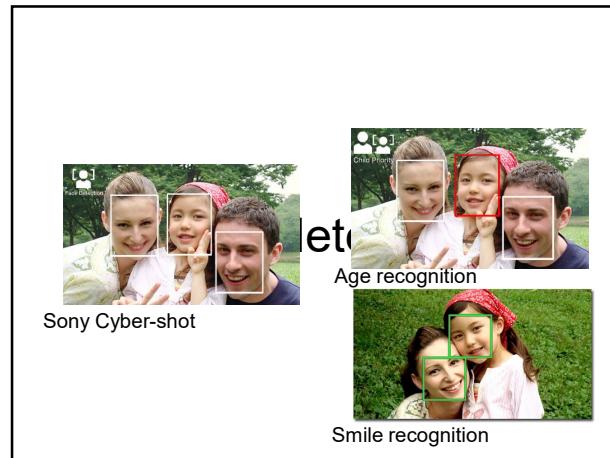
*plus we might have occlusion*



## Applications of computer vision

### Automated visual inspection





## Face makeovers

Creating your own new look is easy

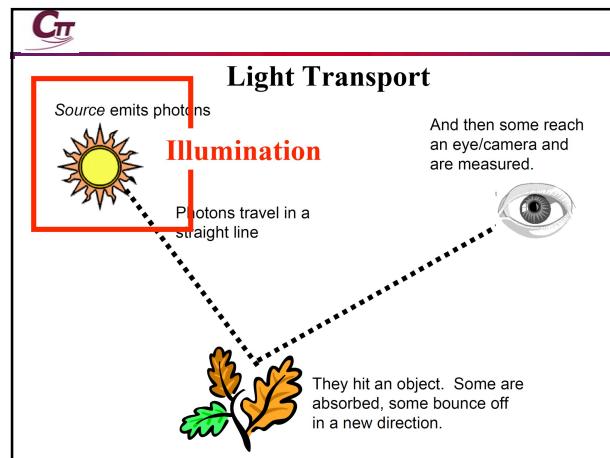
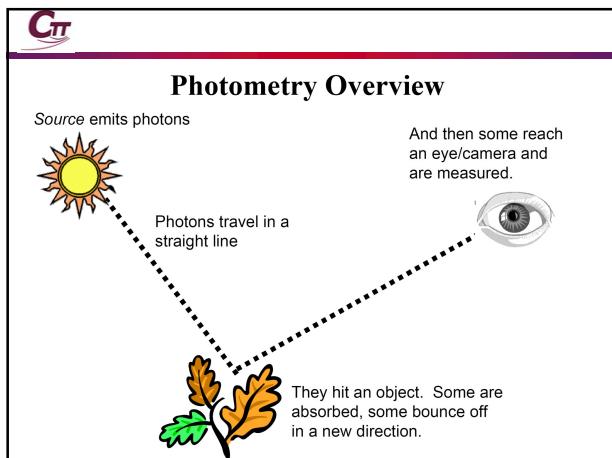
1. Upload your photo
2. Apply some makeup
3. Choose a hairstyle

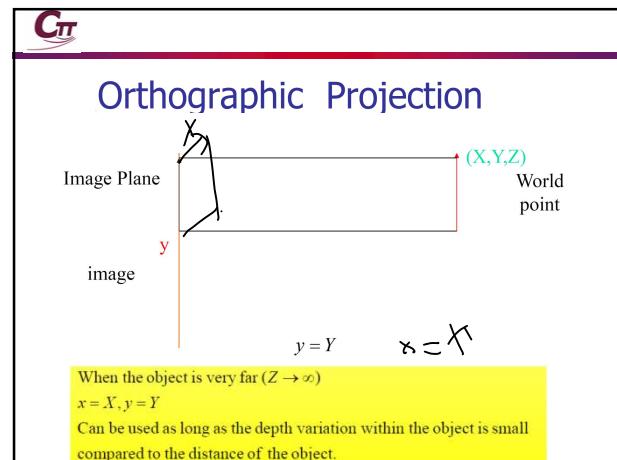
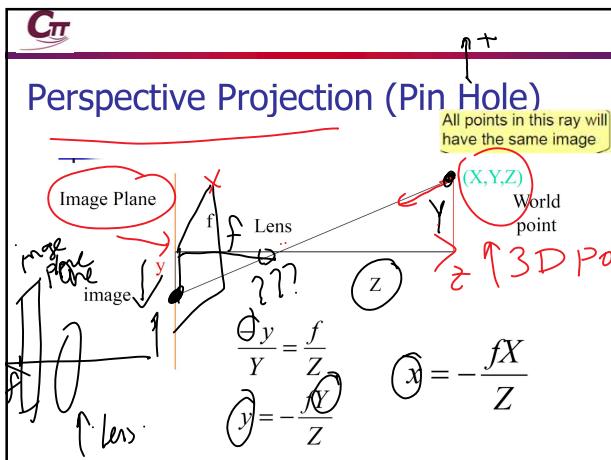
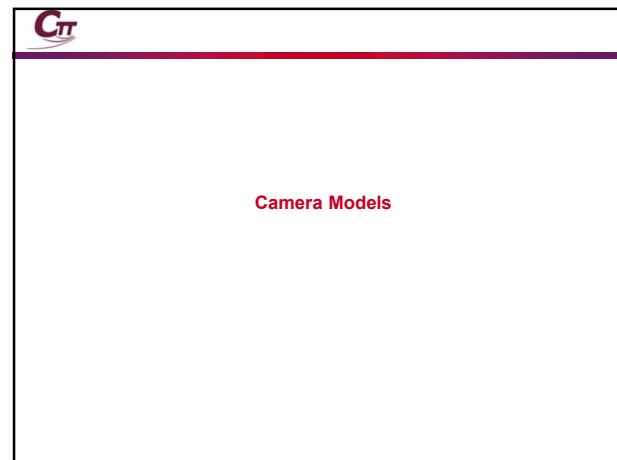
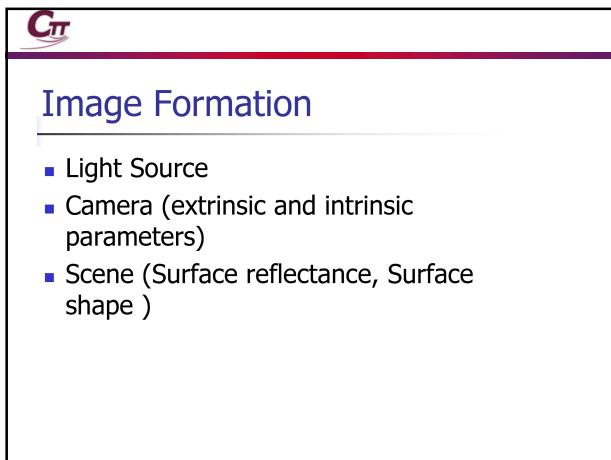
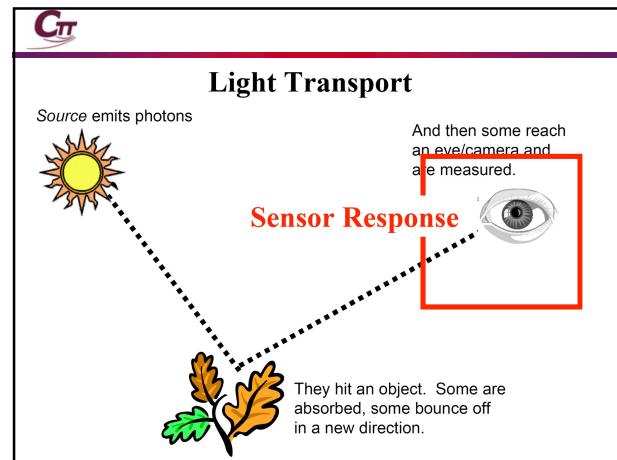
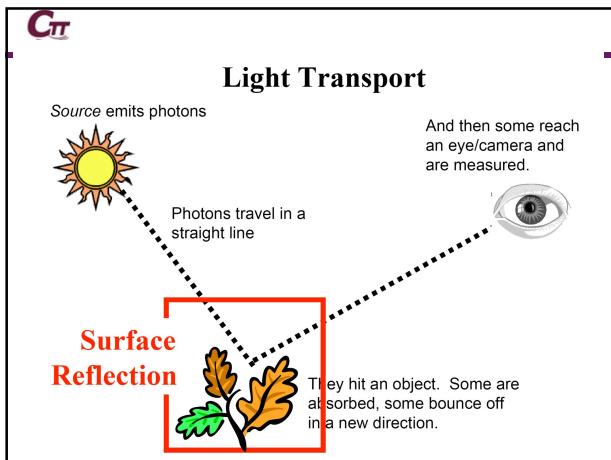
try it now!

TODAY'S FEATURED MAKEOVER: rtyukljkj, Lkmyjny by tazza

TODAY'S FEATURED ADVICE QUESTION: which look is better?

Ask your beauty beauty question. Our community will answer it for you!





**Plane to Plane Mappings**

Rigid, Similarity, Affine, & Projective Mappings  
Homography Estimation  
Image Warping

**Image Warping**

Objective: Change appearance of image by performing geometric transformation, i.e., change the position of a point in the image to a new position.

Example:

**Camera Projection Models**

Projection Models  
Intrinsic (lens) Parameters  
Extrinsic (pose) Parameters  
Camera Calibration

focal length,  
pixelsizes etc.  
rotation matrix  
translation

**Application: Eyevision System  
CAMERA CALIBRATION!**

**Eyevision : SuperBowl XXXV**

**Human Visual Perception**

## Information Processing by Human Observer

The diagram illustrates the visual process: an image is input into an eye, which then produces a perceived image. A dashed arrow labeled "understanding of content" points from the perceived image to the text below.

- Visual perception**
  - Concerns how an image is perceived by a human observer
    - preliminary processing by eye
    - further processing by brains
  - Important to develop image fidelity measures
    - How to evaluate DIP/DVP algorithms & systems

## Eye Anatomy

The top part shows the external features of the eye: Iris, Pupil, Sclera, and Conjunctiva. The bottom part shows a cross-section of the eye with labels for: Sclera, Iris, Cornea, Pupil, Lens, Choroid, Optic nerve, Vitreous, Macula, and Retina. A small diagram at the bottom left shows light rays passing through the eye to form an image on the retina. The source of the diagrams is cited as <http://www.stukeseye.com/Anatomy.asp>.

## Eye Versus Camera

The diagram compares the components of a camera with those of the human eye:

Camera components	Eye components
Lens	Lens, cornea
Shutter	Iris, pupil
Film	Retina
Cable to transfer images	Optic nerve send the info to the brain

**FIGURE 2.1** Simplified diagram of a cross section of the human eye.

This detailed cross-section diagram of the human eye includes the following labels: Cornea, Iris, Anterior chamber, Lens, Ciliary muscle, Ciliary fibers, Visual axis, Vitreous humor, Fovea, Blind spot, Retina, Sclera, Choroid, Nerve & sheath, and Ciliary body. A dashed line indicates the visual axis. The figure is noted to be from slides at Gonzalez/Woods DIP book website (Chapter 2).

- Cross section illustration
- Retina ~ the "film" in eyes to hold our visual sensors

## Two Types of Photoreceptors at Retina

- Rods**
  - Long and thin
  - Large quantity (~ 100 million)
  - Provide scotopic vision (i.e., dim light vision or at low illumination); night vision
  - Only extract luminance information and provide a general overall picture
- Cones**
  - Short and thick, densely packed in fovea (center of retina)
  - Much fewer (~ 6.5 million) and less sensitive to light than rods
  - Provide photopic vision (i.e., bright light vision or at high illumination); day vision
  - Help resolve fine details as each cone is connected to its own nerve end
  - Responsible for color vision
  - There are three types of cones : Red, Green and blue cones (each type has different frequency response)
- Mesopic vision**
  - provided at intermediate illumination by both rod and cones

## Human Visual Perception (Monochrome Vision)

**Figure 2.6**

Figure is from slides at Gonzalez/Woods DIP book website (Chapter 2)

- Light is an electromagnetic wave
  - with wavelength of 350nm to 780nm stimulating human visual response
- Expressed as spectral energy distribution  $I(\lambda)$ 
  - The range of light intensity levels that human visual system can adapt is huge: ~ on 10 orders of magnitude ( $10^{10}$ ) but not simultaneously
  - Brightness adaptation: small intensity range to discriminate simultaneously

**Luminance vs. Brightness**

- Luminance (or intensity)**
  - Independent of the luminance of surroundings
$$L(x, y) = \int_0^{\infty} I(x, y, \lambda) V(\lambda) d\lambda$$

$I(x, y, \lambda)$  -- spatial light distribution  
 $V(\lambda)$  -- relative luminous efficiency function of visual system  
 (bell shape; different for scotopic vs. photopic vision; highest for green wavelength, second for red, and least for blue )
- Brightness**
  - Perceived luminance  $\rightarrow$  Depends on surrounding luminance

Same lum.  
Different brightness

Different lum.  
Similar brightness

**Contrast and Weber's Law**

- From psychovisual research
  - HVS more sensitive to luminance contrast than absolute luminance
  - Eye-brain response to the % changes in intensity is approx. constant
- Weber's Law:  $|L_s - L_o| / L_s \approx \text{const (0.02)}$ 
  - Luminance of an object ( $L_o$ ) is set to be just noticeable from surround luminance ( $L_s$ )
  - For just-noticeable luminance difference (jnd)  $\Delta L$ :  
 Define  $\Delta C \equiv \Delta L / L \approx d(\log L) \approx 0.02$   
 equal increments in log luminance are perceived as equally different
- Empirical Luminance-to-Contrast models
  $C = 50 \log_{10} L$  (logarithmic law, widely used)  
 $L \in [1, 100]$  and  $C \in [0, 100]$

**Contrast and Weber's Ratio**

- The ability of the eye to discriminate between changes in light intensity at any specific adaptation level is different
- Weber ratio :  $\frac{\Delta I}{I}$ ,  $\Delta I$  is the increment of illumination discriminable 50% of the time with background illumination  $I$ 
  - Small  $\frac{\Delta I}{I}$  : good brightness discrimination
  - Large  $\frac{\Delta I}{I}$  : poor brightness discrimination

**Contrast and Weber's Ratio**

- Brightness discrimination is poor at low levels of illumination
- At low levels of illumination, vision is carried out by activity of rods, at high levels, it is by cones

**FIGURE 2.7**

(a) An example showing that perceived brightness is not a simple function of intensity. The relative vertical positions between the two profiles in (b) have no spatial significance; they were chosen for clarity.

Figure is from slides at Gonzalez/Woods DIP book website (Chapter 2)

- Visual system tends to undershoot or overshoot around the boundary of regions of different intensities
- Demonstrates that the perceived brightness is not a simple function of light intensity

### C<sub>IT</sub> Spatial Frequency

- Spatial frequency measures how fast the image intensity changes in the image plane
- Spatial frequency can be completely characterized by the variation frequencies in two orthogonal directions (e.g horizontal and vertical)
  - $f_x$ : cycles/horizontal unit distance
  - $f_y$ : cycles/vertical unit distance
- It can also be specified by magnitude and angle of change

$$f_m = \sqrt{f_x^2 + f_y^2}, \theta = \arctan(f_y / f_x)$$

### C<sub>IT</sub> Sinusoids Temporal frequency

Diagram illustrating sinusoidal signals:

- A sinusoidal wave  $x(t) = A \sin(\omega t)$  with amplitude  $A$  and period  $T$ .
- The period  $T = \frac{1}{f}$  where  $f$  is the frequency.
- The frequency  $f = \frac{1}{T}$ .
- Handwritten note:  $\omega = 2\pi f$

### C<sub>IT</sub>

Figure 2.1 Two-dimensional sinusoidal signals: (a)  $(f_x, f_y) = (5, 0)$ ; (b)  $(f_x, f_y) = (5, 10)$ . The horizontal and vertical units are the width and height of the image, respectively. Therefore,  $f_x = 5$  means that there are five cycles along each row.

### C<sub>IT</sub> Visual Angle and Spatial Frequency

- Angular Spatial Frequency
  - Measures the extent of spatial transition in unit of "cycles per visual degree"
- Visibility thresholds
  - Eyes are most sensitive to medium spatial freq. and least sensitive to high frequencies
    - similar to a band-pass filter
  - More sensitive to horizontal and vertical changes than other orientations

Diagram illustrating angular spatial frequency as a band-pass filter.

### C<sub>IT</sub> Visibility Thresholds at Various Frequencies

- Measuring visibility threshold using sinusoidal grating of varying contrast and frequencies
- Human visual system's frequency response is similar to a band-pass filter

From Jain's Fig.3.7 (pp55) w.r.t. radial spatial freq.

### C<sub>IT</sub> Visibility Threshold at Various Spatial Frequency

- Visibility threshold at different spatial frequency
  - Eyes are most sensitive to mid frequencies, and least sensitive to high frequencies
  - Most sensitive to horizontal and vertical ones than other orientations

 Monochrome Vision Models	 Image Fidelity Criteria
<ul style="list-style-type: none"> <li>• Human Visual System (HVS) <ul style="list-style-type: none"> <li>– from experiment with sinusoidal grating of varying contrast</li> <li>– similar to a band-pass filter <ul style="list-style-type: none"> <li>• most sensitive to mid frequencies</li> <li>• least sensitive to high frequencies</li> </ul> </li> <li>– also depends on the orientation of grating <ul style="list-style-type: none"> <li>• most sensitive to horizontal and vertical ones</li> </ul> </li> </ul> </li> <li>• Overall monochrome vision models <ul style="list-style-type: none"> <li>– How light is transformed by eye to brightness information</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• Subjective measures <ul style="list-style-type: none"> <li>– Examination by human viewers</li> <li>– Goodness scale: excellent, good, fair, poor, unsatisfactory</li> <li>– Impairment scale: unnoticeable, just noticeable, ...</li> <li>– Comparative measures <ul style="list-style-type: none"> <li>• with another image or among a group of images</li> </ul> </li> </ul> </li> <li>• Objective (Quantitative) measures <ul style="list-style-type: none"> <li>– Mean square error and variations</li> <li>– Pro: <ul style="list-style-type: none"> <li>• Simple, less dependent on human subjects, &amp; easy to handle mathematically</li> </ul> </li> <li>– Con: <ul style="list-style-type: none"> <li>• Not always reflect human perception</li> </ul> </li> </ul> </li> </ul>

 Mean-square Criterion
<ul style="list-style-type: none"> <li>• Average (or sum) of squared difference of pixel luminance between two images <math display="block">\mathcal{E}_1 = E[ u - u' ^2] \quad (\text{mean square error})</math> <math display="block">\mathcal{E}_2 = \frac{1}{MN} \sum_{m=1}^M \sum_{n=1}^N  u(m, n) - u'(m, n) ^2 \quad (\text{average squared error})</math> <math display="block">\mathcal{E}_3 = \frac{1}{MN} \sum_{m=1}^M \sum_{n=1}^N E[ u(m, n) - u'(m, n) ^2] \quad (\text{average mean square error})</math> </li> <li>• Signal-to-noise ratio (SNR) <ul style="list-style-type: none"> <li>– <math>\text{SNR} = 10 \log_{10} (\sigma_s^2 / \sigma_e^2)</math> in unit of decibel (dB) <ul style="list-style-type: none"> <li>• <math>\sigma_s^2</math> image variance</li> <li>• <math>\sigma_e^2</math> variance of noise or error</li> </ul> </li> <li>– <math>\text{PSNR} = 10 \log_{10} (A^2 / \sigma_e^2)</math> <ul style="list-style-type: none"> <li>• A is peak-to-peak value</li> <li>• PSNR is about 12-15 dB higher than SNR</li> </ul> </li> </ul> </li> </ul>