# Computer Vision (CSE3010)

**Dr. Susant Kumar Panigrahi**
**Assistant Professor**
**School of Electrical & Electronics Engineering**

# Module-1 Syllabus

**Digital Image Formation And Low Level Processing**:

- Overview and State-of-the-art, Fundamentals of Image Formation, Transformation: Orthogonal, Euclidean, Affine, Projective, Fourier Transform,

- Convolution and Filtering, Image Enhancement, Restoration, Histogram Processing.

# Module-2 Syllabus

**Depth Estimation And Multi-Camera Views:**

Depth Estimation and Multi-Camera Views: Perspective, Binocular Stereopsis: Camera and Epipolar Geometry; Homography, Rectification, DLT, RANSAC, 3-D reconstruction framework; Auto-calibration. apparel.

# Module-3 Syllabus

**Feature Extraction And Image Segmentation:**

- **Feature Extraction**: Edges - Canny, LOG, DOG; Line detectors (Hough Transform), Corners - Harris and Hessian Affine, Orientation Histogram, SIFT, SURF, HOG, GLOH, Scale-Space Analysis- Image Pyramids and Gaussian derivative filters, Gabor Filters and DWT.

- **Image Segmentation:** Region Growing, Edge Based approaches to segmentation, Graph-Cut, Mean-Shift, MRFs, Texture Segmentation; Object detection.

# Module-4 Syllabus

**Pattern Analysis And Motion Analysis:**

• **Pattern Analysis**: Clustering: K-Means, K-Medoids, Mixture of Gaussians, Classification: Discriminant Function, Supervised, Un-supervised, Semi-supervised; Classifiers: Bayes, KNN, ANN models;

• **Dimensionality Reduction**: PCA, LDA, ICA; Non-parametric methods. Motion Analysis: Background Subtraction and Modelling, Optical Flow, KLT, Spatio-Temporal Analysis, Dynamic Stereo; Motion parameter estimation.

# Module-5 Syllabus

**Shape From X:**

Light at Surfaces; Phong Model; Reflectance Map;

Albedo estimation; Photometric Stereo; Use of Surface Smoothness

Constraint; Shape from Texture, color, motion and edges.

**Guest Lecture on Contemporary Topics**

## Text Books

1. Richard Szeliski, Computer Vision: Algorithms and Applications, Springer-Verlag London Limited 2011.
2. Computer Vision: A Modern Approach, D. A. Forsyth, J. Ponce, Pearson Education, 2003.

## Reference Book(s):

1. R.C. Gonzalez and R.E. Woods, Digital Image Processing, Addison- Wesley, 1992.
2. Richard Hartley and Andrew Zisserman, Multiple View Geometry in Computer Vision, Second Edition, Cambridge University Press, March 2004.
3. K. Fukunaga; Introduction to Statistical Pattern Recognition, Second Edition, Academic Press, Morgan Kaufmann, 1990.

## Required Tools/Software/IDLE:

1. Python/jupyter-notebook/google-colab
2. OpenCV
3. MATLAB

# Indicative List of Experiments:

1. Implement image preprocessing and Edge
2. Implement camera calibration methods
3. Implement Projection
4. Determine depth map from Stereo pair
5. Construct 3D model from Stereo pair
6. Implement Segmentation methods
7. Construct 3D model from defocus image
8. Construct 3D model from Images
9. Implement optical flow method
10. Implement object detection and tracking from video
11. Face detection and Recognition
12. Object detection from dynamic Background for Surveillance
13. Content based video retrieval
14. Construct 3D model from single image

# Computer Vision
# Unit – 01
# Low-Level Vision (Digital Image Formation And Low Level Processing)

# VISION: HUMANS ARE VISUAL CREATURES

- Half of the human brain is directly or indirectly devoted to processing visual information.
  - The eye's retina, which contains 150 million light-sensitive rod and cone cells

- The brain can identify images seen for as little as 13 milliseconds.
  - Helps the brain as it decides where to focus the eyes
  - Deciding where to move the eyes can take 100 to 140 ms, so very high-speed understanding must occur before that.

- At least 65 % of people are "visual learners"

- Humans have a remarkable ability to remember pictures.
  - More than 2000 pictures with at least 90 % accuracy.
- What our eyes see can influence what we hear, which is called the "McGurk Effect".

# Why we want Machines to Emulate Human Vision?

- ✓ To engage machines to perform our mundane works.

- ✓ Human visions are more qualitative rather quantitative.

- ✓ Human vision is limited.

**Human Perception VS Machine Vision**

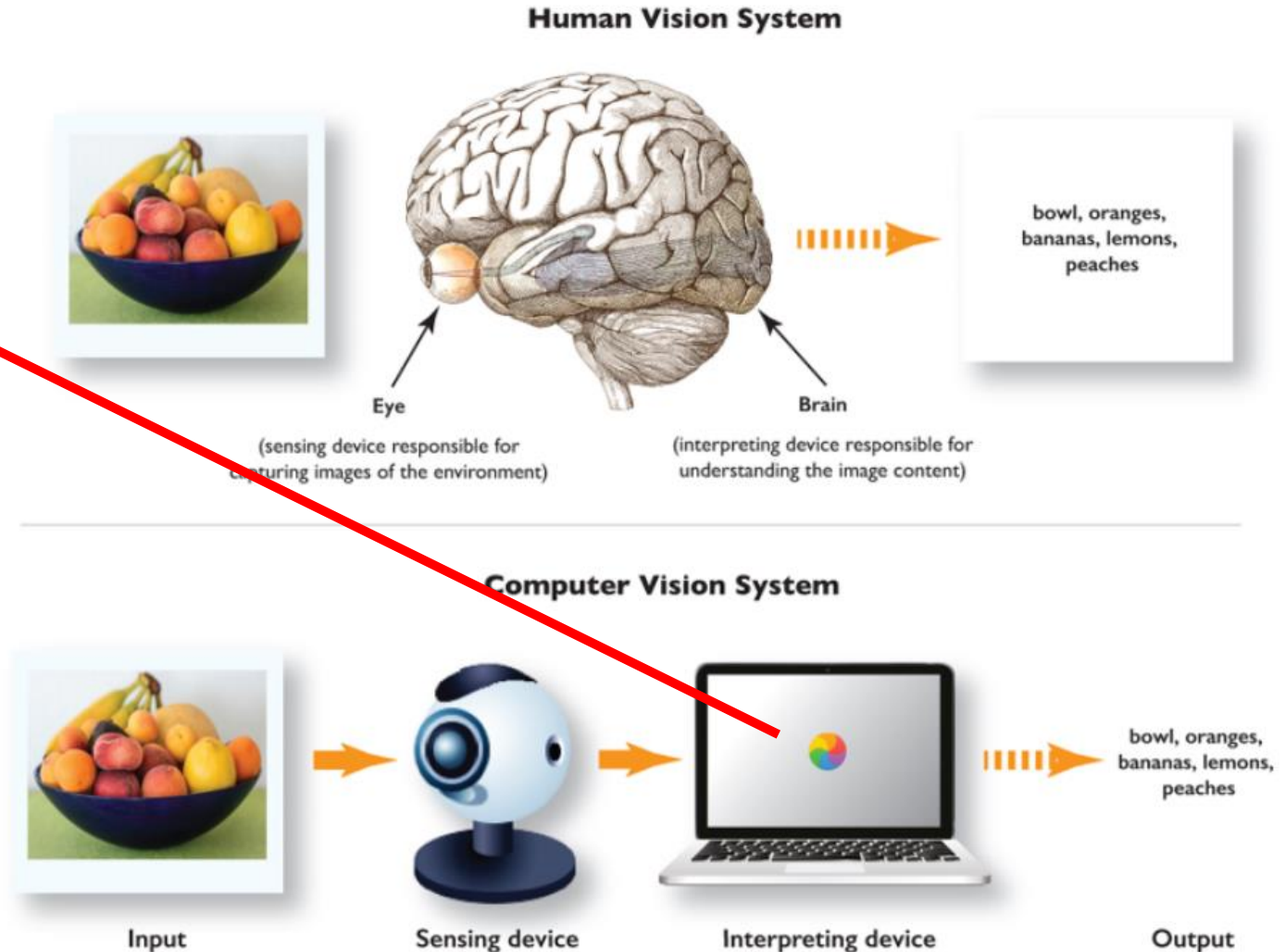**Limited vs Entire EM spectrum**

# HUMAN VS COMPUTER VISION

Interpreting Device or Vision Software: To come up with some symbolic description of the scene. [Bowl, Orange, Lemon, Banana and Peach].

More sophisticated algorithms may also illustrate "How fresh the fruits are what color/shape/size the bowl is"

Computer Vision*: where humans teach computers to see and interpret the world around them.

**Human Vision System**

bowl, oranges, bananas, lemons, peaches

Eye
(sensing device responsible for capturing images of the environment)

Brain
(interpreting device responsible for understanding the image content)

**Computer Vision System**

bowl, oranges, bananas, lemons, peaches

Input

Sensing device

Interpreting device

Output

*a subfield of Deep Learning and Artificial Intelligence

# But, what really is Computer Vision?

VISION is.....

... Automating human visual processes.
-- David Mahar

Computer vision that emulates human vision.

..... An information processing task.

..... Inverting image formation.
--- Berthold Horn

..... Inverse of computer graphics.

...... REALY USEFUL.

# *Why Computer Vision Difficult?*

- Inverse Problem

**Perceives** the "Story" behind the picture/video

Loss of dimensionality; while capturing image using sensing devices.

Seeks to recover some unknowns given insufficient information to fully specify the problem

## Challenges:

1. Viewpoint Variation: Input image may align in different directions that leads the computer vision system to predict inaccurate results.

2. Scale Variation: Images captured closer to camera looks bigger and vice versa. Variation in size or scale affects the decision taking capabilities of the system.

3. Deformation: System may learn from perfect image and depicts a particular perception about the shape, size and other features. In real-world shape may change that leads to inaccuracy when shape of the object is deformed.

4. Inter-class variation: Objects of same class may come in different shape, size, colour and texture; but the algorithm need to identify them as single class.

5. **Scale Variation:** Incomplete information due to occlusion results in inaccurate interpretation.



6. **Illumination Variation:** Same image captured under different illumination levels.

# Computer Vision System: How Does CV Works?

Computer Vision System:

1) Image Formation: Cameras to obtain visual data,

2) Low-level image processing:

3) Low level vision

4) Middle-level Vision

5) High-level vision: Image Understanding

6) Decision making

# Vision deals with Images…!

A picture worth more than thousand words.

# Scope and Challenges in CV

- ✓ Vision is a hard problem.

- ✓ Vision is multi-disciplinary.

- ✓ Considerable progress has been made.

- ✓ Many successful real-world applications.

- ✓ Computer vision as part of AI

# IMAGE ACQUISITION: IMAGING SYSTEM

**Camera + Scanner → Digital Camera: Get images into computer**

# IMAGE ACQUISITION: IMAGING SYSTEM



An example of digital image acquisition process. (a) Energy (illumination) Source. (b) An element of a scene. (c) Imaging System. (d) Projection of the scene on to the imaging plane. (e) Digitized Image

**Formulation**:

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

Where:
$i(x, y)$ = Illumination (Energy) source.

$$0 < i(x, y) < \infty$$

$r(x, y)$ = Reflectance component of the scene.
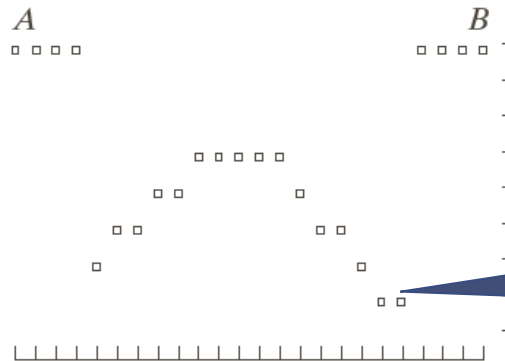
$$0 \leq r(x, y) \leq 1$$

# DIGITAL IMAGE

- Reflectance of some scenes

  - 0.01 for black velvet

  - 0.65 for stainless steel

  - 0.80 for flat-white wall paint

  - 0.90 for silver-plated metal

  - 0.93 for snow

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

Where:
$i(x, y)$ = Illumination (Energy) source.

$$0 < i(x, y) < \infty$$

$r(x, y)$ = Reflectance component of the scene.

$$0 \leq r(x, y) \leq 1$$

# DIGITAL IMAGE

- Image: Two dimensional function $f(x, y)$

  Where $x$ and $y$ spatial(plane) coordinates

- Intensity or gray level: The amplitude of $f$ at any pair of coordinates $(x, y)$
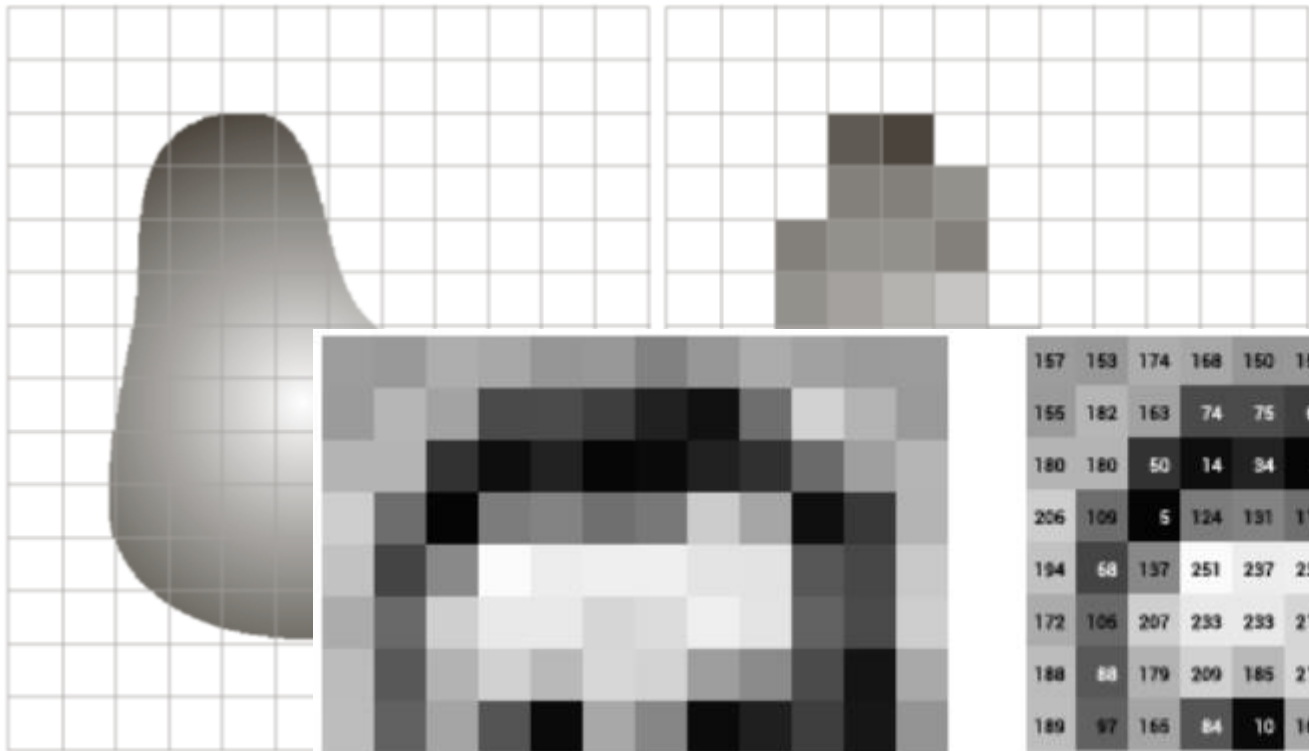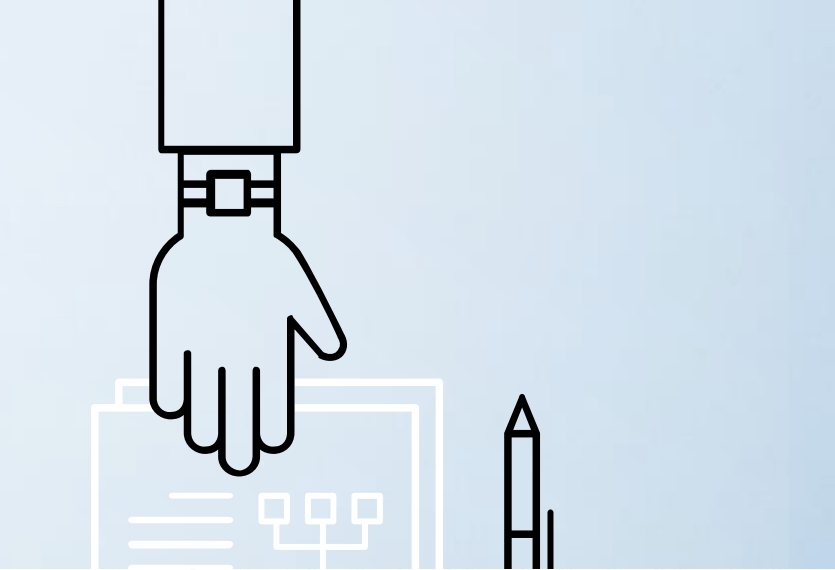
- Image as function

# IMAGE SAMPLING & QUANTIZATION



Digitizing the coordinate values

Digitizing the amplitude values

Generating a digital image. (a) Continuous image. (b) A scan line from A to B in continuous image, used to illustrate the concept of sampling and quantization. (c) Sampling and quantization. (d) Digital scan line.
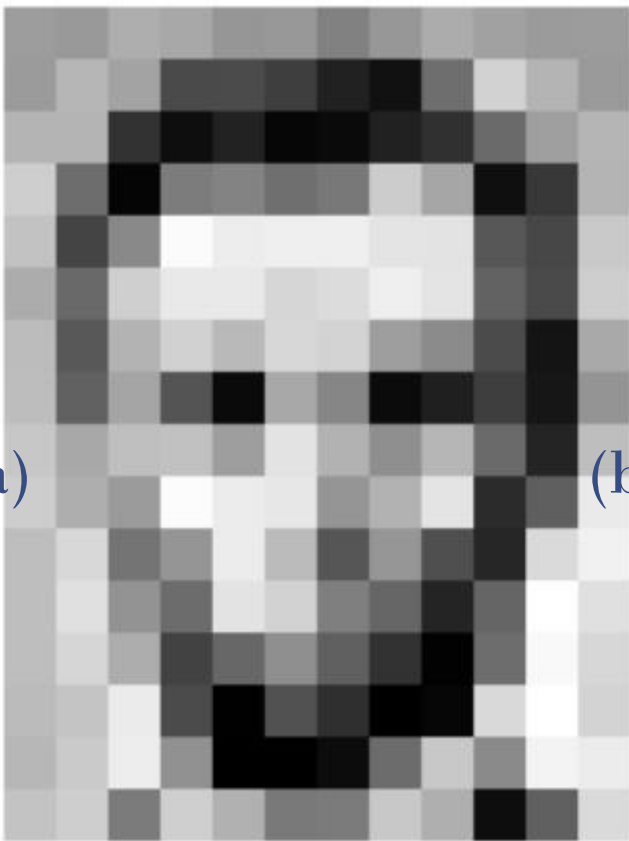
# IMAGE SAMPLING & QUANTIZATION



(a)

(b)

(a) Continuous image

# DIGITAL IMAGE REPRESENTATION

$$f(x, y) = \begin{bmatrix} f(0, 0) & f(0, 1) & \cdots & f(0, N-1) \\ f(1, 0) & f(1, 1) & \cdots & f(1, N-1) \\ \vdots & \vdots & & \vdots \\ f(M-1, 0) & f(M-1, 1) & \cdots & f(M-1, N-1) \end{bmatrix}$$



(0, 0)

*y*

*x*

A monochrome image and the convention used to represent rows

(*x*) and columns (*y*).

**2D Sampling**:
- Represented in a matrix with *M* rows and *N* columns.
- Equally spaced samples in two dimension.

  [ Digitizing the coordinate values are called **sampling**]

*Quantization:*
- Digitizing the amplitude values are called Quantization.
- The amplitude values depend on the "data type" 0r "class" by which it is represented.
- The following are few popular representations:
  - ✓ *uint8*
  - ✓ *uint16*
  - ✓ *unit32*
  - ✓ *double*
  - ✓ *single*

# DIGITAL IMAGE REPRESENTATION

- Discrete intensity interval $[0, L-1], L = 2^k$ ; L = Dynamic Range

- The number $b$ of bits required to store a $M \times N$ digitized image

$$b = M \times N \times k$$

Number of storage bits for various values of $N$ and $k$.

| $N/k$ | 1 ($L = 2$) | 2 ($L = 4$) | 3 ($L = 8$) | 4 ($L = 16$) | 5 ($L = 32$) | 6 ($L = 64$) | 7 ($L = 128$) | 8 ($L = 256$) |
|---|---|---|---|---|---|---|---|---|
| 32 | 1,024 | 2,048 | 3,072 | 4,096 | 5,120 | 6,144 | 7,168 | 8,192 |
| 64 | 4,096 | 8,192 | 12,288 | 16,384 | 20,480 | 24,576 | 28,672 | 32,768 |
| 128 | 16,384 | 32,768 | 49,152 | 65,536 | 81,920 | 98,304 | 114,688 | 131,072 |
| 256 | 65,536 | 131,072 | 196,608 | 262,144 | 327,680 | 393,216 | 458,752 | 524,288 |
| 512 | 262,144 | 524,288 | 786,432 | 1,048,576 | 1,310,720 | 1,572,864 | 1,835,008 | 2,097,152 |
| 1024 | 1,048,576 | 2,097,152 | 3,145,728 | 4,194,304 | 5,242,880 | 6,291,456 | 7,340,032 | 8,388,608 |
| 2048 | 4,194,304 | 8,388,608 | 12,582,912 | 16,777,216 | 20,971,520 | 25,165,824 | 29,369,128 | 33,554,432 |
| 4096 | 16,777,216 | 33,554,432 | 50,331,648 | 67,108,864 | 83,886,080 | 100,663,296 | 117,440,512 | 134,217,728 |
| 8192 | 67,108,864 | 134,217,728 | 201,326,592 | 268,435,456 | 335,544,320 | 402,653,184 | 469,762,048 | 536,870,912 |

# TYPES DIGITAL IMAGE

- Common image formats include:
    - 1 sample per point (B&W or Grayscale)
    - 3 samples per point (Red, Green, and Blue)
    - 4 samples per point (Red, Green, Blue, and "Alpha", a.k.a. Opacity)



- For most of this course we will focus on grey-scale and colour images

# TYPES OF DIGITAL IMAGES



A grayscale (monochrome) image and the pixel values in a 6 × 6 neighborhood.



A binary image and the pixel values in a 6 × 6 neighborhood.

**MAX-MIN Pixel Intensity**:

Gray-scale Image:

$$min = 0$$
$$max = 2^b - 1$$

b = Number of bits required to represent the pixel intensity values.

**MAX-MIN Pixel Intensity**:

Binary Image:

$$min = 0$$
$$max = 1$$

[Only two intensity values can be represented]

# COLOR IMAGE



(a) Color (RGB) image

(b) Red-Channel

(c) Green-Channel of the color image

(d) Blue-Channel

A color/RGB image and the pixel values in a 3 × 3 neighborhood.

Partial Pixel Region tool readout:

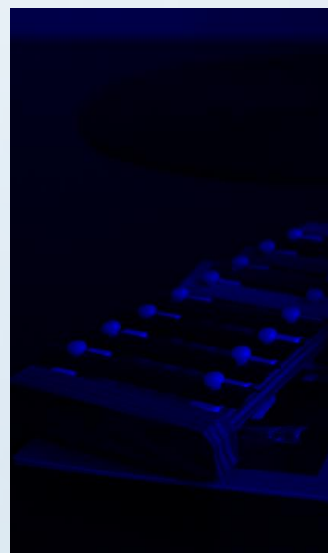| | G: 7 | G: 7 | G:10 | G:12 |
|---|---|---|---|---|
| 0 | B:13 | B:20 | B:16 | B:16 |
| :14 | | | | |
| :13 | R: 6 | R: 4 | R: 6 | R:16 |
| : 6 | G: 6 | G:10 | G: 6 | G: 8 |
| :16 | B:16 | B:17 | B:17 | B:19 |
| : 4 | R: 6 | R:14 | R: 2 | R:13 |
| :11 | G: 6 | G: 6 | G: 4 | G: 7 |
| :17 | B: 8 | B: 8 | B:19 | B:17 |
| : 8 | R:15 | R: 5 | R: 7 | R:13 |
| :13 | G: 9 | G: 8 | G:11 | G:12 |

Pixel info: (X, Y) [R G B]

✓ Color images can be represented using three 2D arrays of same size, one for each color channel: **red (R), green (G), and blue (B) .**

✓ Each array element contains an 8-bit value, indicating the amount of red, green, or blue at that point in a [0, 255] scale .

✓ The combination of the three 8-bit values into a 24-bit number allows **$2^{24}$** (16,777,216, usually referred to as 16 million or **16 M**) color combinations.

# Image Transformation

- An **image processing** operation typically defines a new image $g$ in terms of an existing image $f$.

- We can transform either the range of $f$.
$$g(x, y) = t(f(x, y))$$

- Or the domain of $f$:
$$g(x, y) = f\big(t_x(x, y), t_y(x, y)\big)$$

- What kinds of operations can each perform?