

COMPUTER VISION

CS-6350

Prof. Sukhendu Das
Deptt. of Computer Science and Engg.,
IIT Madras, Chennai – 600036.

Email: sdas@cse.iitm.ac.in

URL: [//www.cse.iitm.ac.in/~sdas](http://www.cse.iitm.ac.in/~sdas)

[//www.cse.iitm.ac.in/~vplab/computer_vision.html](http://www.cse.iitm.ac.in/~vplab/computer_vision.html)

JULY – 2023.

INTRODUCTION

Contents to be covered

- 1 Introduction
- 2 Neighborhood and Connectivity of pixels
- 3 DFT, Filtering/Enhancement in spatial and spectral domains
- 4 3D transformations, projection and stereo reconstruction
- 5 Histogram based image processing & DHS
- 6 Concepts in Edge Detection
- 7 Hough Transform
- 8 Image segmentation
- 9 Texture analysis using Gabor filters
- 10 Pattern Recognition
- 11 Motion Analysis
- 12 Shape from Shading
- 13 Scale-Space - Image Pyramids
- 14 Feature extraction (recent trends) – detectors and descriptors
- 15 Bag of Words and Prob. Graphical Models
- 16 Object Recognition
- 17 Wavelet transform
- 18 Registration and Matching

20 Solid Modelling;
22. Hardware;

21. Color
23. Morphology

Use slides as brief :
Points, concepts, links

*These are not substitute
for materials in books*

References

- 1. "Digital Image Processing"; R. C. Gonzalez and R. E. Woods; Addison Wesley; 1992+.**
- 2. "Computer Vision: Algorithms and Applications"; by Richard Szeliski; Springer-Verlag London Limited 2011.**
- 3. "Multiple View geometry"; R. Hartley and A. Zisserman, 2002 Cambridge university Press.**
- 4. "Pattern Recognition and Machine Learning"; Christopher M. Bishop; Springer, 2006.**
- 5. "Digital Image Processing and Computer Vision"; Robert J. Schalkoff; John Wiley and Sons; 1989+.**
- 6. "Pattern Recognition: Statistical. Structural and Neural Approaches"; Robert J. Schalkoff; John Wiley and Sons; 1992+.**
- 7. "3-D Computer Vision"; Y. Shirai; Springer-Verlag, 1984**
- 8. "Computer Vision: A Modern Approach"; D. A. Forsyth and J. Ponce; Pearson Education; 2003+.**

References (Contd..)

Journals:

- IEEE-T-PAMI (Transactions on Pattern Analysis and Machine Intelligence)
- IEEE-T-IP (Transactions on Image processing)
- PR (Pattern Recognition)
- PRL (Pattern Recognition Letters)
- CVIU (Computer Vision, Image Understanding)
- IJCV (International Journal of Computer Vision)

Online links

1. CV online: <http://homepages.inf.ed.ac.uk/rbf/CVonline>
2. Computer Vision Homepage:
<http://www-2.cs.cmu.edu/afs/cs/project/cil/ftp/html/vision.html>

Typical Distribution of marks for Evaluation/grading

Quiz (50 mins.) - 15 - 20

End Sem exam (120-150 mins.) - 35 – 40

TPA - 30 - 35

TUTs - 10 - 15

Total 100

**+/- 05 marks variation at any part;
To be finalized well before End Sem Exam.**

*Pre-Req: - Linear Algebra; Geometry; Stat&Prob basics; Calculus basics;
DSP, Programming, Data Structure basics*

July-Nov '23

Days	8.00 - 8.50	9.00 - 9.50	10.00 - 10.50	11.00 - 11.50	12.00 – 12.50	14.00 - 15.15	15.30 - 16.45	17.00 - 17.50
Mon	A	B	C	D	G	P		J/J3
						H/H1	M/M2	
Tue	B	C	D	E	A	Q		F
						M/M1	H/H2	
Wed	C	D	E	F	B	R		G
						J/J1	K/K2	
Thu	E	F	G	A	D	S		H/H3
						L/L1	J/J2	
Fri	F	G	A	B	C	T		E
						K/K1	L/L2	

L
U
N
C
H



- TUTs – Altn. weeks; Mid-sem etc.

*May be held Online
Occasionally*

What is CVPR ?

<http://cvpr2022.thecvf.com/>

<https://openaccess.thecvf.com/menu>

<https://openaccess.thecvf.com/CVPR2022>

Also, check **ICCV (26)**, **ECCV (21)**, **NIPS**

	Publication	<u>h5-index</u>	<u>h5-median</u>
1.	Nature	<u>467</u>	707
2.	The New England Journal of Medicine	<u>439</u>	876
3.	Science	<u>424</u>	665
4.	IEEE/CVF Conference on Computer Vision and Pattern Recognition	<u>422</u>	681
5.	The Lancet	<u>368</u>	688
6.	Nature Communications	<u>349</u>	456
7.	Advanced Materials	<u>326</u>	415
8.	Cell	<u>316</u>	503
9.	Neural Information Processing Systems	<u>309</u>	503
10.	International Conference on Learning Representations	<u>303</u>	563
11.	JAMA	<u>286</u>	476
12.	Science of The Total Environment	<u>273</u>	375
13.	Nature Medicine	<u>268</u>	459
14.	Proceedings of the National Academy of Sciences	<u>268</u>	394
15.	Angewandte Chemie International Edition	<u>266</u>	362
16.	Chemical Reviews	<u>264</u>	459
17.	International Conference on Machine Learning	<u>254</u>	463

- **3D computer vision**
- **Action and behavior recognition**
- Adversarial learning, adversarial attack and defense methods
- **Biometrics, face**, gesture, body pose
- Computational photography, **image and video synthesis**
- Datasets and evaluation
- Efficient training and inference methods for networks
- Explainable AI, fairness, accountability, privacy, transparency ethics in vision
- **Image retrieval**
- **Low-level** and physics-based vision
- Machine learning architectures and formulations
- Medical, biological and **cell microscopy**
- **Motion and tracking**
- Neural generative models, **auto encoders, GANs**
- **Optimization** and learning methods
- **Recognition (object detection, categorization)**
- Representation learning, **deep learning**
- **Scene analysis** and understanding
- **Segmentation, grouping and shape**
- **Transfer**, low-shot, **semi- and un- supervised learning**
- **Video analysis** and understanding
- Vision + language, vision + other modalities
- Vision applications & systems, **vision for robotics & autonomous vehicles**
- **Visual reasoning** and logical representation

3D from multi-view & sensors	Photogrammetry and remote sensing
3D from single images	Physics-based vision and shape-from-X
Action/event recognition	Pose estimation and tracking
Adversarial attacks & defense	Privacy & federated learning
Behavior analysis	Recognition: detection, categorization, retrieval
Biometrics	Representation learning
Computational photography	RGBD sensors and analytics
Computer vision theory	Robot vision
Computer vision for social good	Scene analysis and understanding
Datasets and evaluation	Segmentation, grouping & shape analysis
Deep learning architectures & techniques	Statistical methods
Document analysis and understanding	Transfer/low-shot/long-tail learning
Efficient learning and inference	Transparency, fairness, accountability, privacy & ethics in vision
Explainable computer vision	Self-/semi-/meta-/unsupervised learning
Face and gesture	Video analysis and understanding
Image and video synthesis and generation	Vision + graphics
Low-level vision	Vision + language
Machine learning	Vision + X
Medical, biological and cell microscopy	Vision applications and systems
Motion and tracking	Visual reasoning
Navigation and autonomous driving	Others
Optimization methods	

3D from single images
Adversarial attack and defense
Autonomous driving
Biometrics
Computational imaging
Computer vision for social good
Computer vision theory
Datasets and evaluation
Deep learning architectures and techniques
Document analysis and understanding
Efficient and scalable vision
Embodied vision: Active agents, simulation
Explainable computer vision
Humans: Face, body, pose, gesture, movement
Image and video synthesis and generation
Low-level vision
Machine learning (other than deep learning)
Medical and biological vision, cell microscopy

Multi-modal learning
Optimization methods (other than deep

Photogrammetry and remote sensing
Physics-based vision and shape-from-X
Recognition: Categorization, detection, retrieval

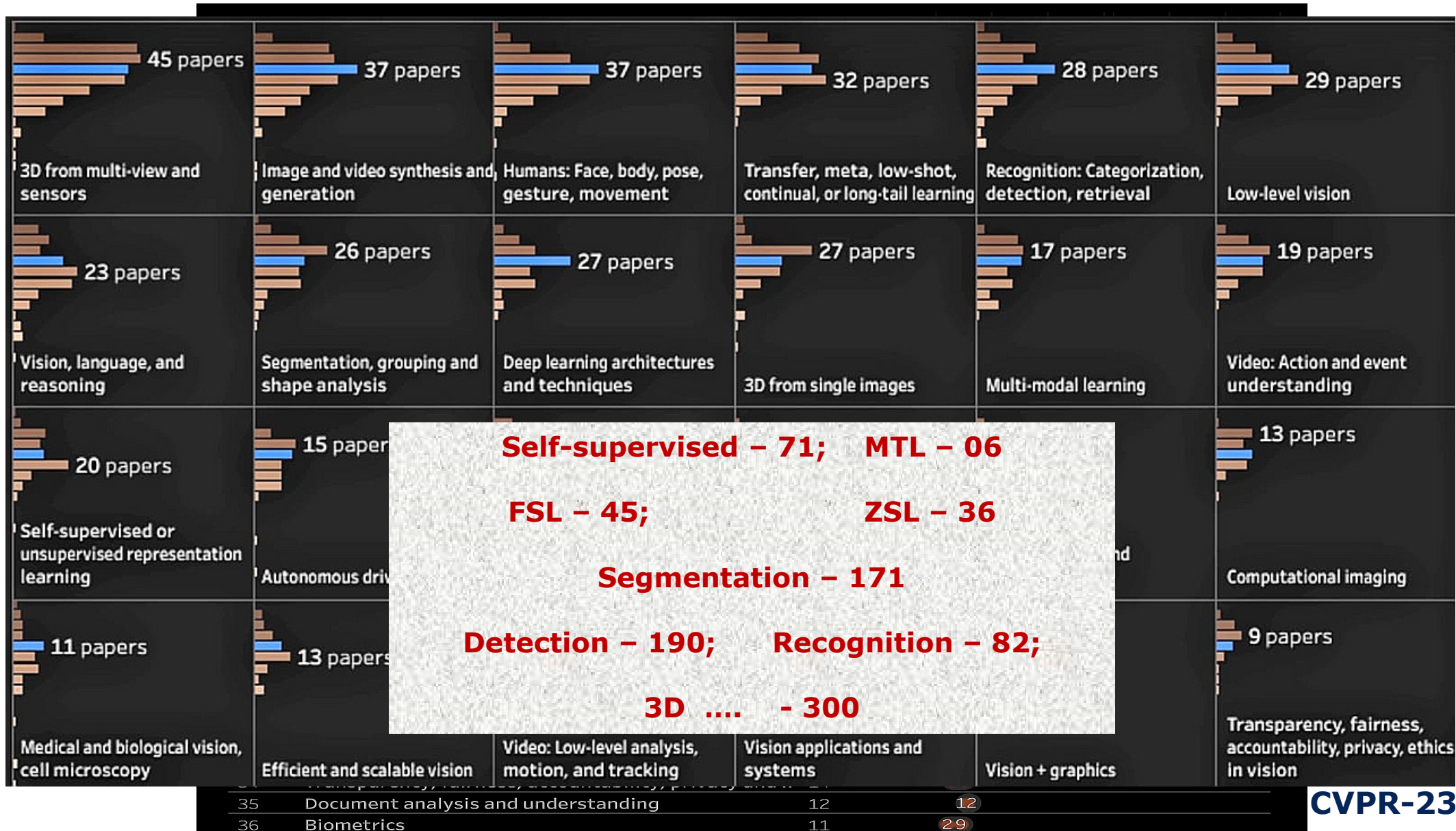
Robotics
Scene analysis and understanding
Segmentation, grouping and shape analysis
Self-supervised or unsupervised representation learning

Transfer, meta, low-shot, continual, or long-tail learning

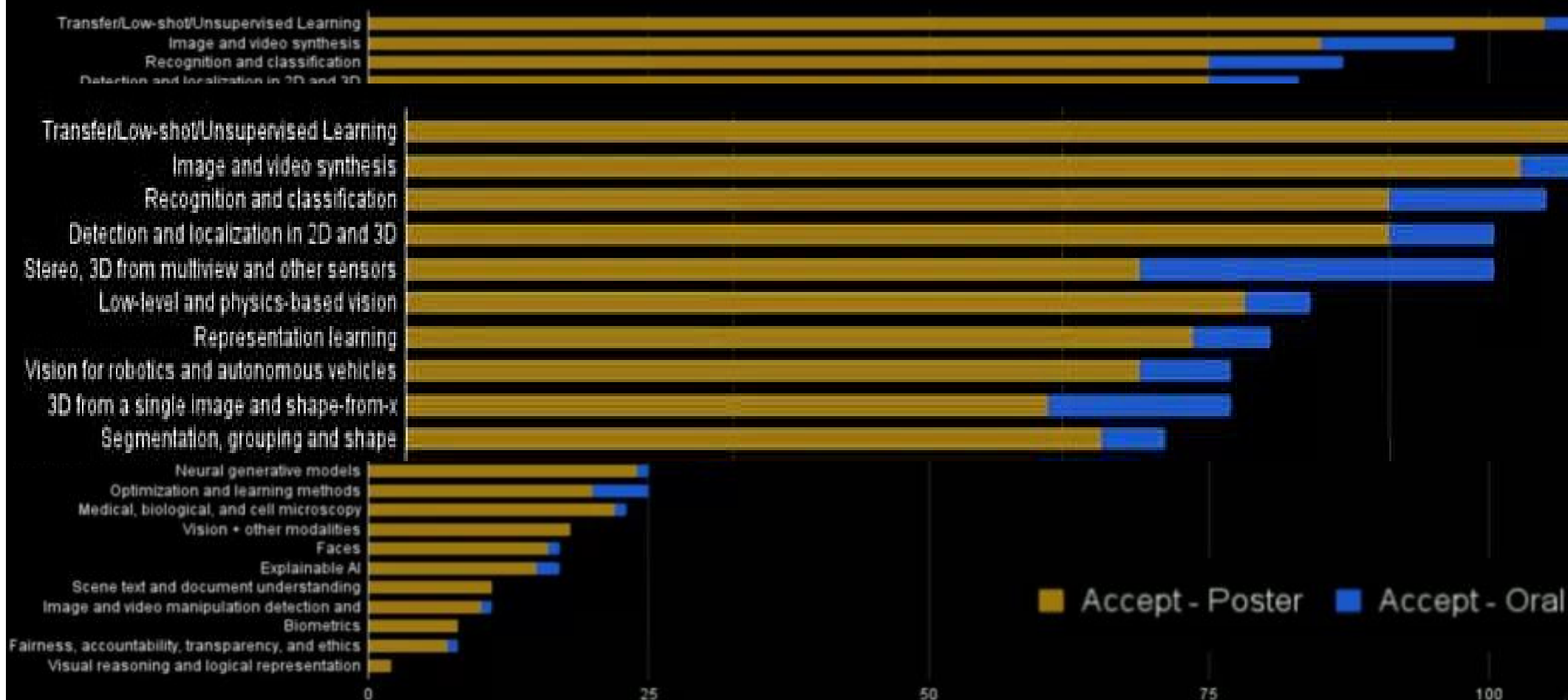
Transparency, fairness, accountability, privacy, ethics in vision

Video: Action and event understanding
Video: Low-level analysis, motion, and tracking
Vision + graphics
Vision, language, and reasoning
Vision applications and systems

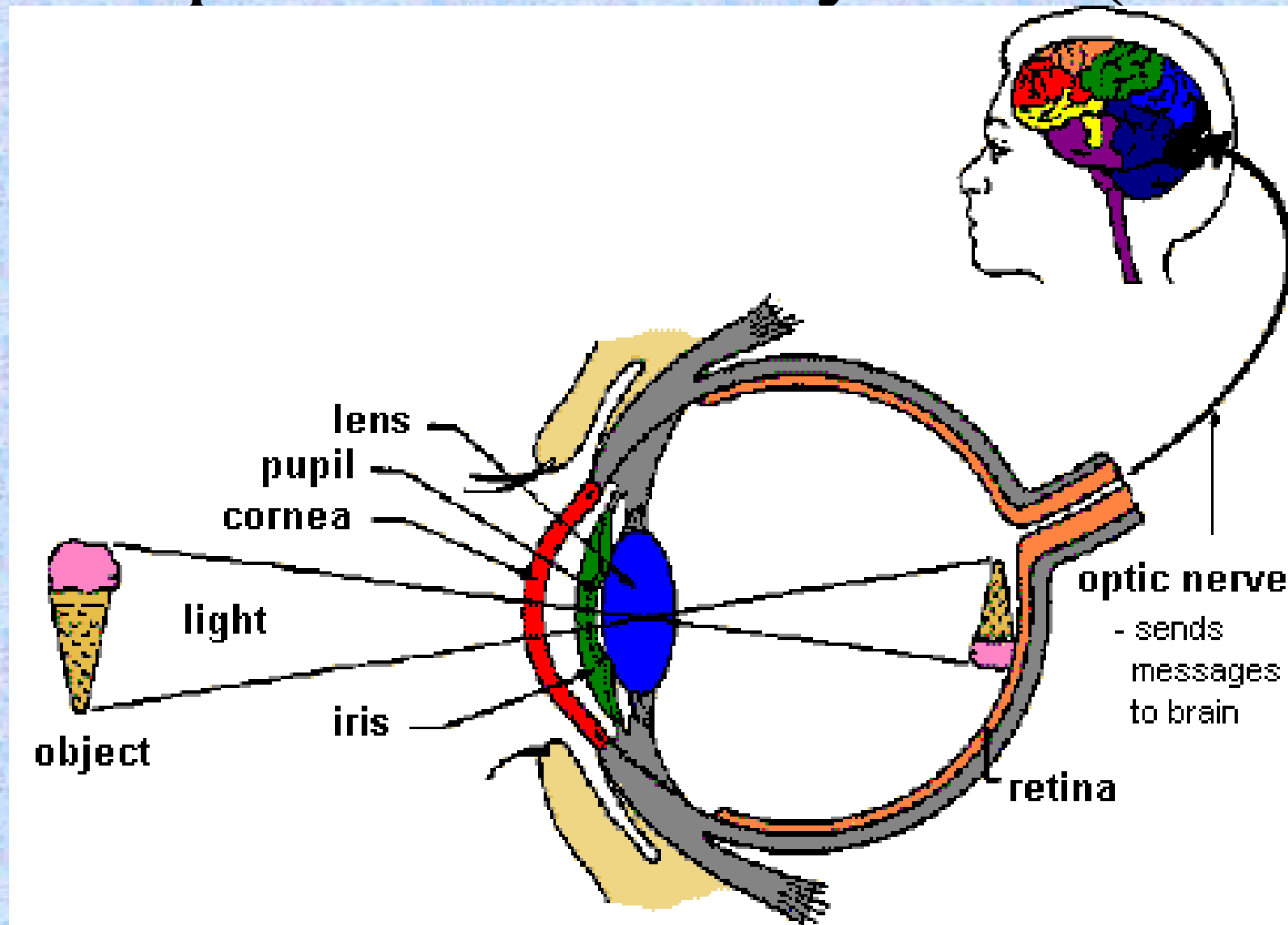
CVPR – 2022-3



Subject Areas of Accepted Papers

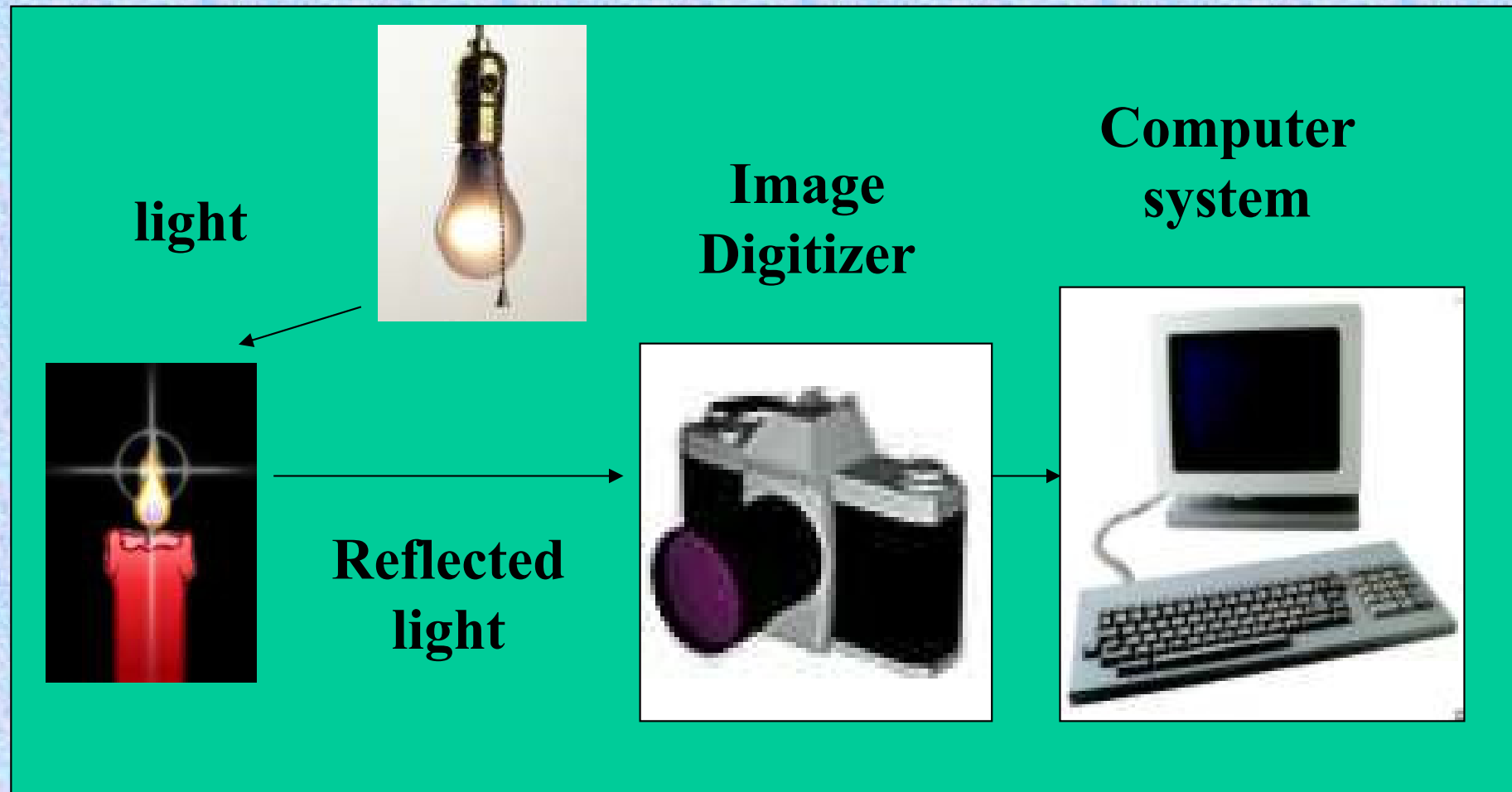


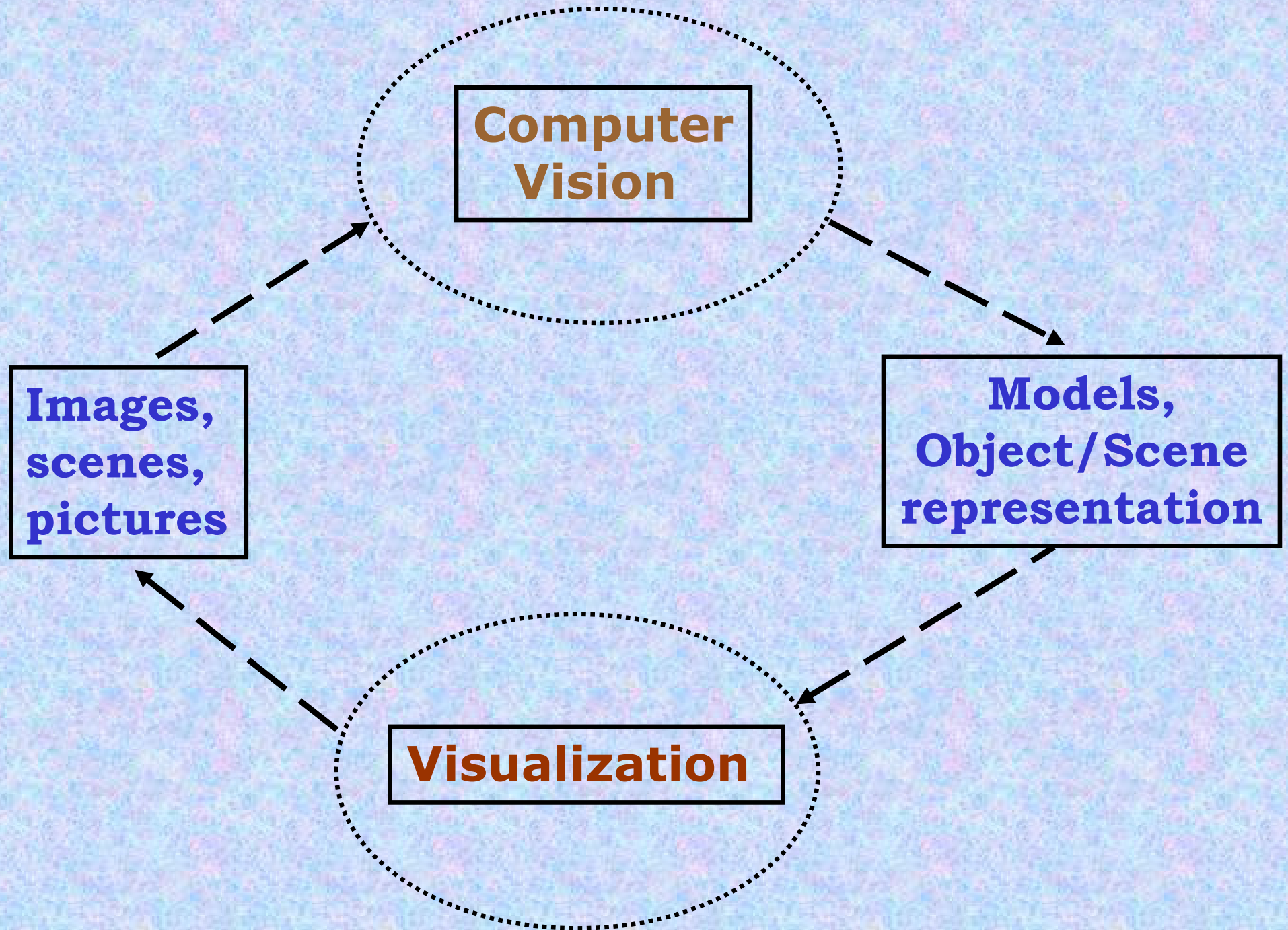
Human Vision System (HVS) Vs. Computer Vision System (CVS)



The Optics of the eye

A computer Vision System (CVS)



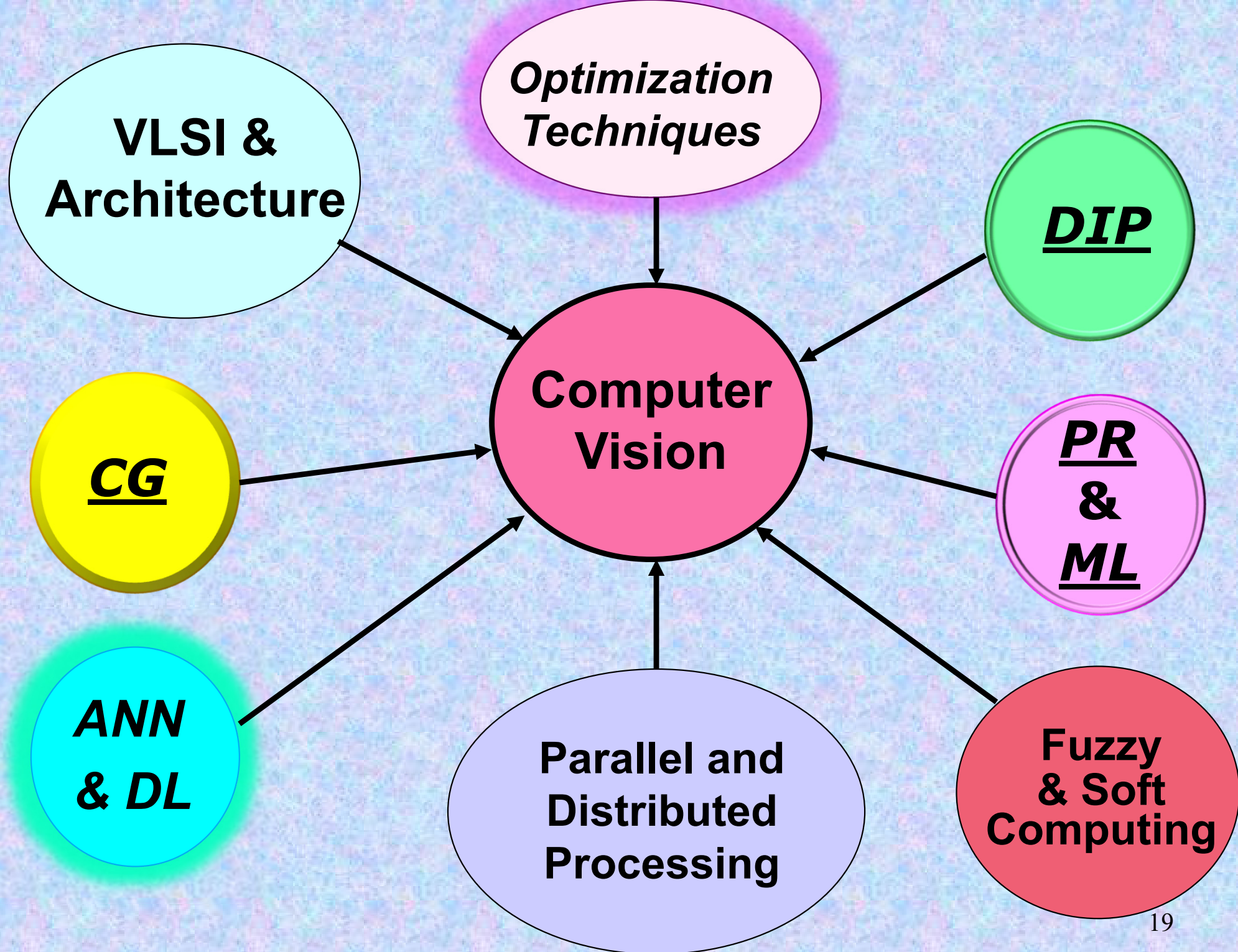


Computer Vision is an area of work, which is a combination of concepts, techniques and ideas from Digital Image Processing, Pattern Recognition, Artificial Intelligence and Computer Graphics.

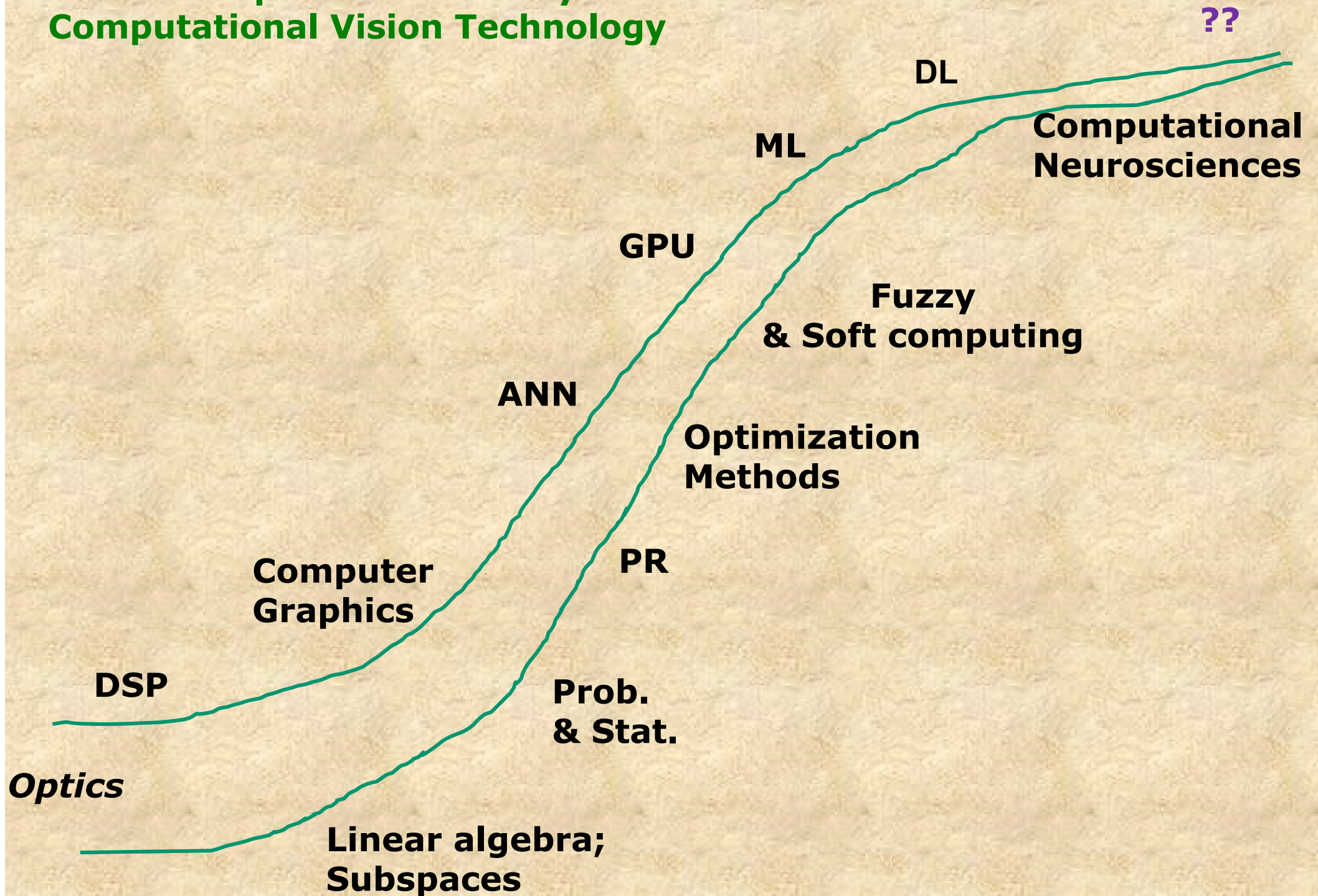
Majority of the tasks in the fields of Digital Image Processing or Computer Vision deals with the process of **understanding** or deriving the scene information or description, from the input scene (digital image/s). The methods used to solve a problem in digital image processing depends on the **application domain** and nature of data being analyzed.

Analysis of two-dimensional pictures are generally not applicable of processing three-dimensional scenes, and vice-versa. The choice of processing, techniques and methods and '**features**' to be used for a particular application is made after some amount of trial and error, and hence experience in handling images is crucial in most of these cases.

For example, analysis of remote sensed or satellite imagery involves techniques based on classification or analysis of texture imagery. These techniques are not useful for analyzing optical images of indoor or outdoor scenes.



The Developmental Pathway of Computational Vision Technology



Digital Image processing is in many cases concerned with taking one array of pixels as input and producing another array of pixels as output which in some way represents an improvement to the original array.

Purpose:

1. Improvement of Pictorial Information

- improve the contrast of the image,
- remove noise,
- remove blurring caused by movement of the camera during image acquisition,
- it may correct for geometrical distortions caused by the lens.

2. Automatic Machine perception (termed Computer Vision, Pattern Recognition or Visual Perception) for intelligent interpretation of scenes or pictures.

Elements of a Digital Image Processing System

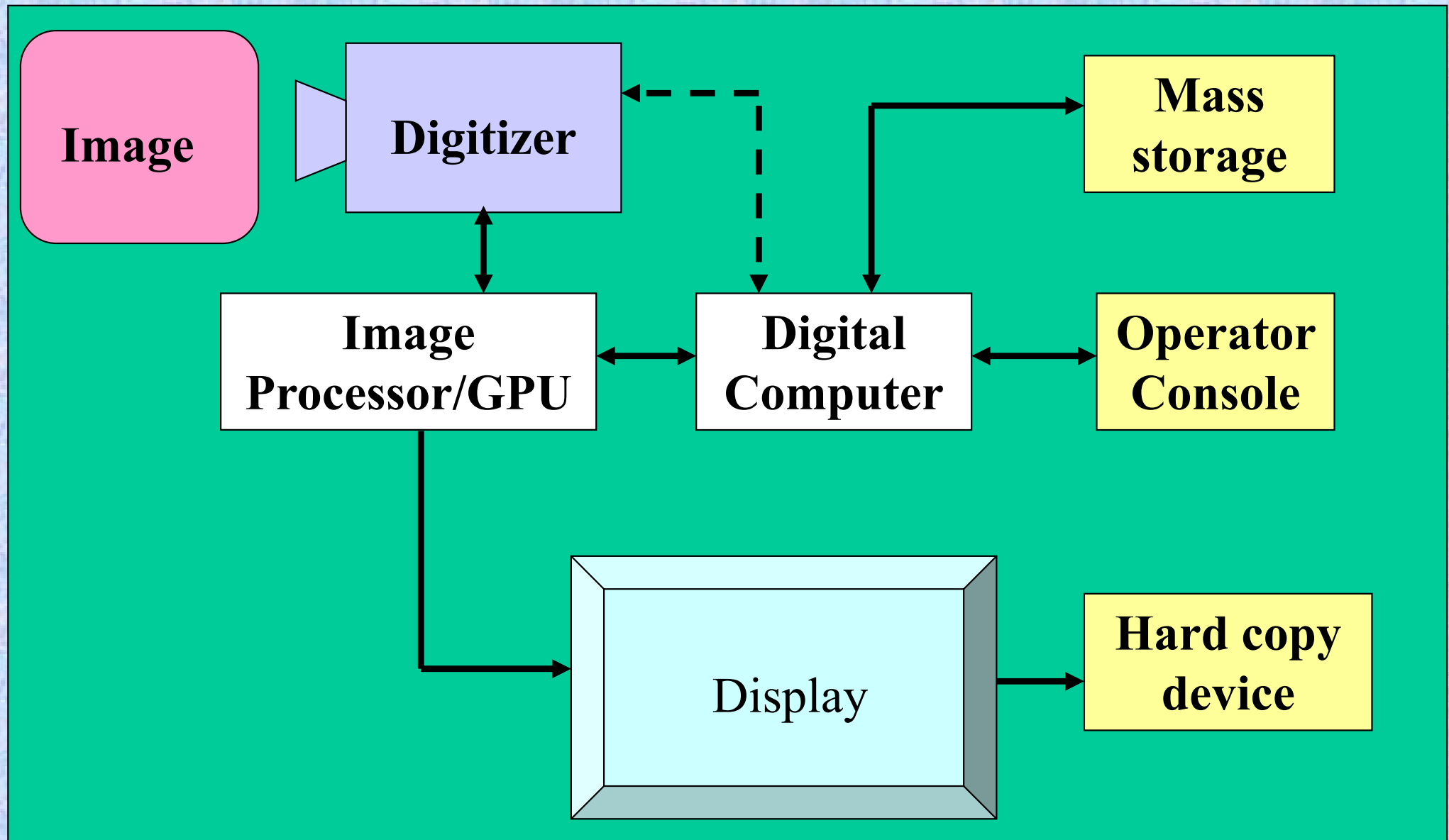
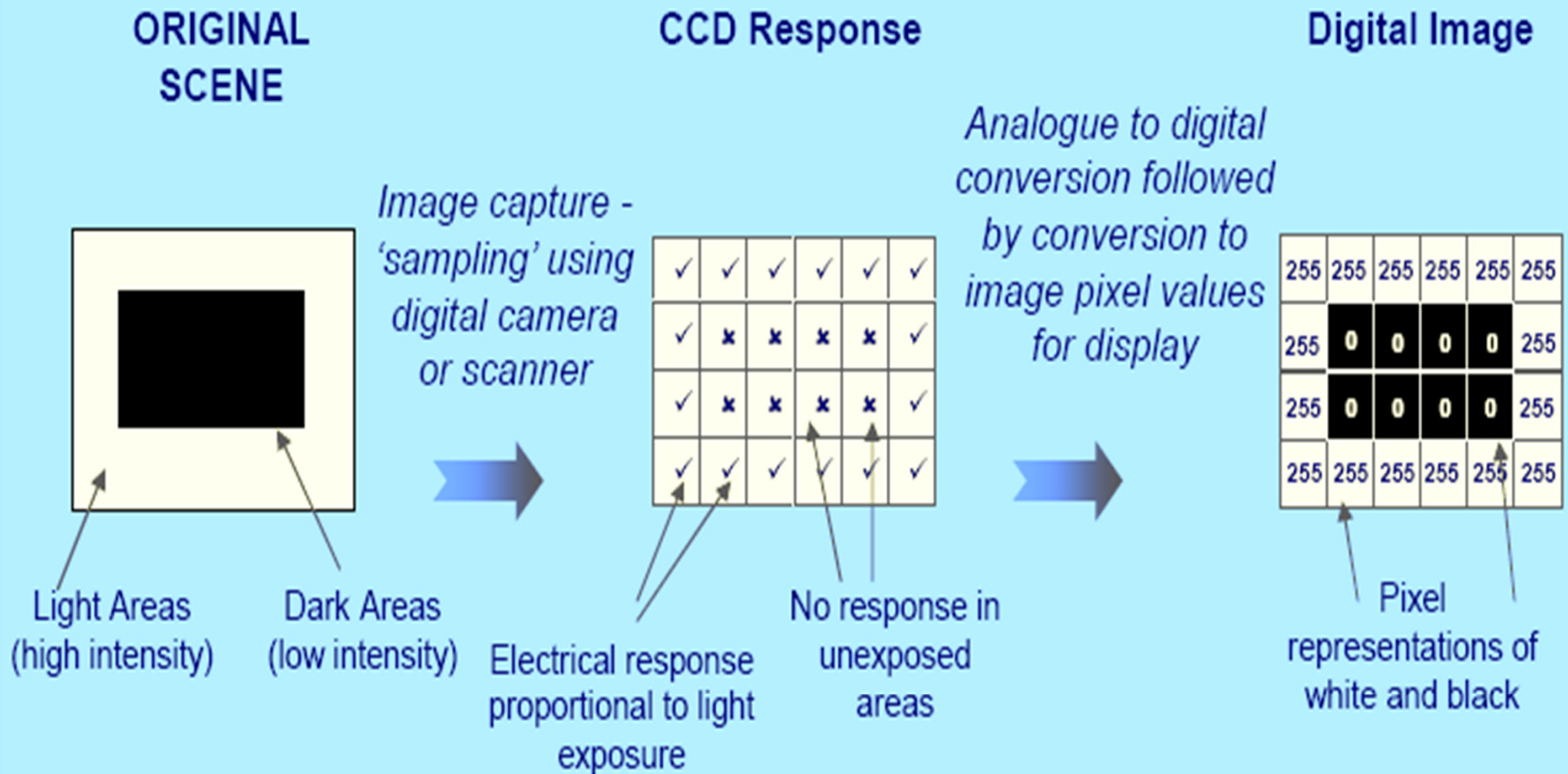


Image processors: Consists of set of hardware modules that perform 4 basic functions:

- Image acquisition: frame grabber**
 - Storage: frame buffer**
 - Low-level processing: specialized hardware device designed to perform Arithmetic Logic operations on pixels in parallel**
 - Display: read from image memory (frame buffer) and convert to analog video signal**
- Digitizers: Converts image into numerical representation suitable for input to a digital computer**
 - Digital Computers: Interfaced with the image processor to provide versatility and ease of programming.**
 - Storage Devices: For bulk storage. e.g:- Magnetic disks, magnetic tapes, optical disks**
 - Display and Recording devices : Monochrome and Color Television monitors, CRT, Laser printers, heat-sensitive paper devices, and ink spray systems.**

Image acquisition using a CCD camera



Resolution standards: HDMI - 1024*768;

UHD -

A digital Image

Image is an array of integers: $f(x,y) \in \{0,1,\dots,I_{\max}-1\}$,
where, $x,y \in \{0,1,\dots,N-1\}$

- N is the resolution of the image and I_{\max} is the level of discretized brightness value
- Larger the value of N , more is the clarity of the picture (larger resolution), but more data to be analyzed in the image
- If the image is a gray-level (8-bit per pixel - termed raw, gray) image, then it requires N^2 Bytes for storage
- If the image is color - RGB, each pixel requires 3 Bytes of storage space.

Image Size (resolution)	Storage space required	
	Raw - Gray	Color (RGB)
64*64	4K	12K
256*256	64K	192K
512*512	256K	768K

$2048 \times 1536 =$  megapixels \rightarrow  MB for RGB

A **digital image** is a two-dimensional (3-D image is called range data) array of intensity values, $f(x, y)$, which represents 2-D intensity function discretized both in spatial coordinates (**spatial sampling**) and brightness (**quantization**) values.

The elements of such an array are called **pixels** (picture elements).

The storage requirement for an image depends on the **spatial resolution** and number of bits necessary for **pixel quantization**.

The processing of an image depends on the application domain and the methodology used to solve a problem. There exists four broad categories of tasks in digital image processing:

- | | |
|------------------------------|---------------------------|
| (i) Compression, | (ii) Segmentation, |
| (iii) Recognition and | (iv) motion. |

Segmentation deals with the process of fragmenting the image into homogeneous meaningful parts, regions or sub-images. Segmentation is generally based on the analysis of the histogram of images using gray level values as features. Other features used are edges or lines, colors and textures.

Recognition deals with identification or classification of objects in an image for the purpose of interpretation or identification. Recognition is based on models, which represent an object. A system is trained (using HMM, GMM, ANN etc.) to learn or store the models, based on training samples. The test data is then matched with all such models to identify the object with a certain measure of confidence.

Compression involves methodologies for efficient storage and retrieval of image data, which occupies large disk space. Typical methods are, JPEG-based, Wavelet based, Huffman Coding, Run length coding etc. for still images and MPEG-I, II, IV & VII for digital video or sequence of frames.

Motion analysis (or dynamic scene analysis) involves techniques for the purpose of tracking and estimation of the path of movement of object/s from a sequence of frames (digital video). Methods for dynamic scene analysis are based on (i) tracking, (ii) obtaining correspondence between frames and then (iii) estimating the motion parameters and (iv) structure of moving objects. Typical methods for analysis are based on optical flow, iterative Kalman filter and Newton/Euler's equations of dynamics.

There are generally three main categories of tasks involved in a complete computer vision system. They are:

- ***Low level processing:*** Involves image processing tasks in which the quality of the image is improved for the benefit of human observers and higher level routines to perform better.
- ***Intermediate level processing:*** Involves the processes of feature extraction and pattern detection tasks. The algorithms used here are chosen and tuned in a manner as may be required to assist the final tasks of high level vision.
- ***High level vision:*** Involves autonomous interpretation of scenes for pattern classification, recognition and identification of objects in the scenes as well as any other information required for human understanding.

A **top down approach**, rather than a bottom-up approach is used in the **design** of these systems in many applications. The **methods** used to solve a problem in digital image processing depends on the **application domain** and **nature of data** being analyzed.

Different fields of applications include:

- **Character Recognition,**
- **Document processing,**
- **Commercial (signature & seal verification) application,**
- **Biometry and Forensic (authentication: recognition and verification of persons using face, palm & fingerprint),**
- **Pose and gesture identification,**
- **Automatic inspection of industrial products,**
- **Industrial process monitoring,**
- **Biomedical Engg. (Diagnosis and surgery),**
- **Military surveillance and target identification,**
- **Navigation and mobility (for robots and unmanned vehicles - land, air and underwater),**
- **Remote sensing (using satellite imagery),**
- **GIS**
- **Safety and security (night vision),**
- **Traffic monitoring,**
- **Sports (training and incident analysis)**
- **VLDB (organization and retrieval)**
- **Entertainment and virtual reality.**

TARGETED INDUSTRIAL APPLICATIONS

Intelligent Traffic Control

Anti-forging Stamps

Card Counting Systems

Drive Quality Test

Camera Flame Detection

CCTV Fog Penetration

Key Image Search/Index

Security Monitoring

Robust Shadow Detection

Vehicle Segmentation

Visual Tracking Systems

Illegal content (adult) Filter

Scratch Detection

Smart Traffic Monitoring

Vehicle Categorization

Vehicle Wheel alignment

Number Plate Identification

Referrals for Line calls

Different categories of work being done in CV, to solve problems:

**2-D image analysis –
segmentation, target detection,
matching, CBIR;**

**Pattern Recognition
for Objects, scenes;**

**Feature extraction:
Canny, GHT, Snakes,
DWT, Corners,
SIFT, GLOH, LESH;**

**Image and Video-based
Rendering;**

**3-D multi-camera calibration;
Correspondence and stereo;
Reconstruction of
3-D Objects and surfaces;**

**Video and motion analysis;
Video analytics; CBVR;
Compression;**

**Multi-sensor data,
Decision and feature fusion;**

**Steganography and
Watermarking;**

The various sub-categories of technology in these related fields are:

- *image enhancement,*
- *image restoration and filtering,*
- *representation and description,*
- *feature extraction,*
- *image segmentation,*
- *image matching,*
- *color image processing,*
- *image synthesis,*
- *image representation,*
- *image reconstruction*
- *range data processing,*
- *stereo image processing*
- *computational geometry,*
- *image morphology,*
- *artificial neural networks,*
- *Neuro-fuzzy techniques,*
- *computational geometry,*
- *parallel architectures & algorithms.*

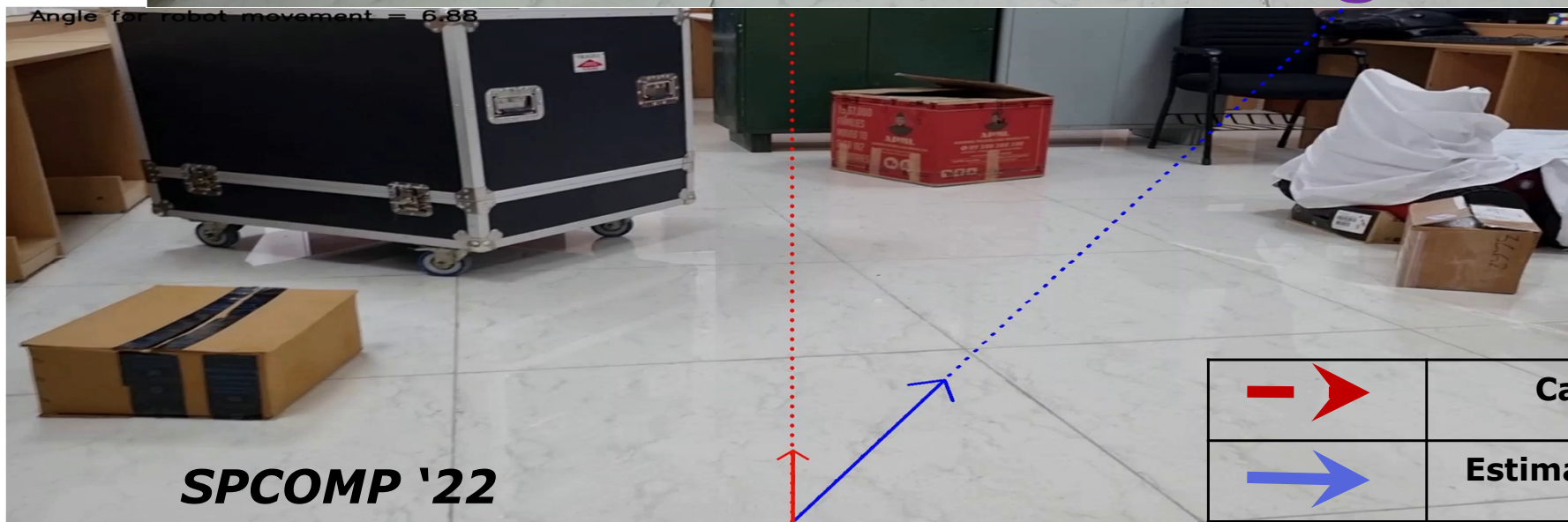
Few DEMOS and ILLUSTRATIONS

Courtesy: TA/students of VPLAB - CSE-IITM

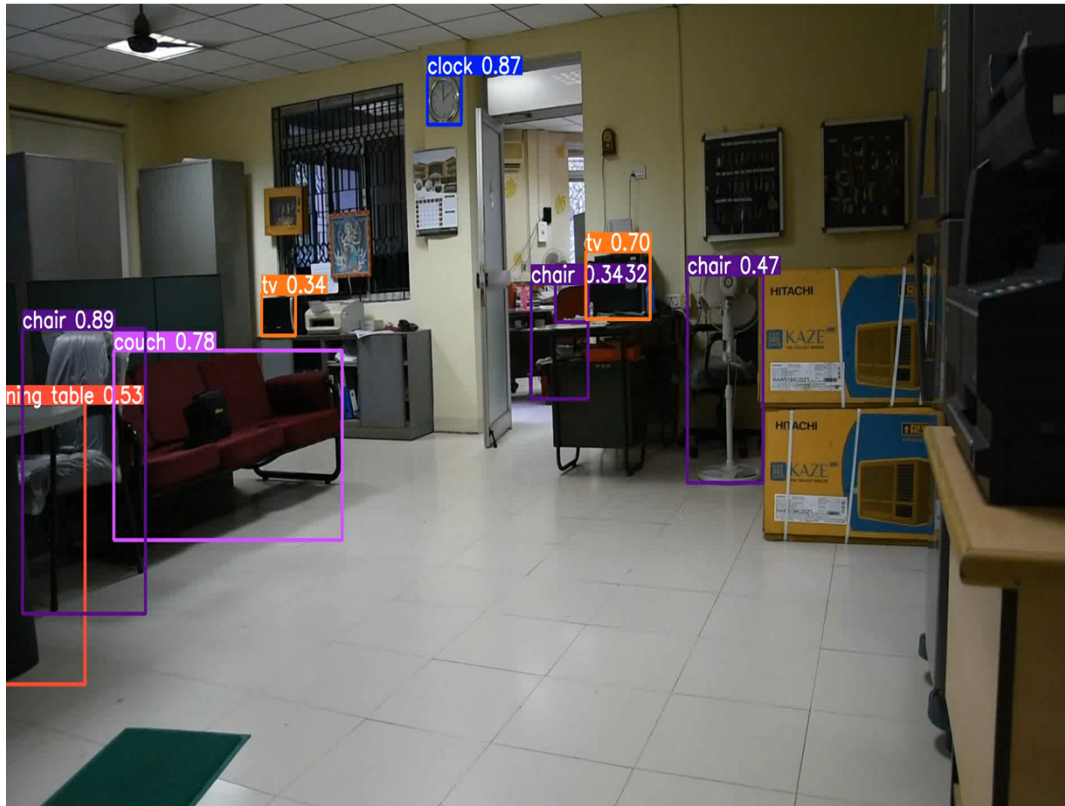
Video Object Segmentation



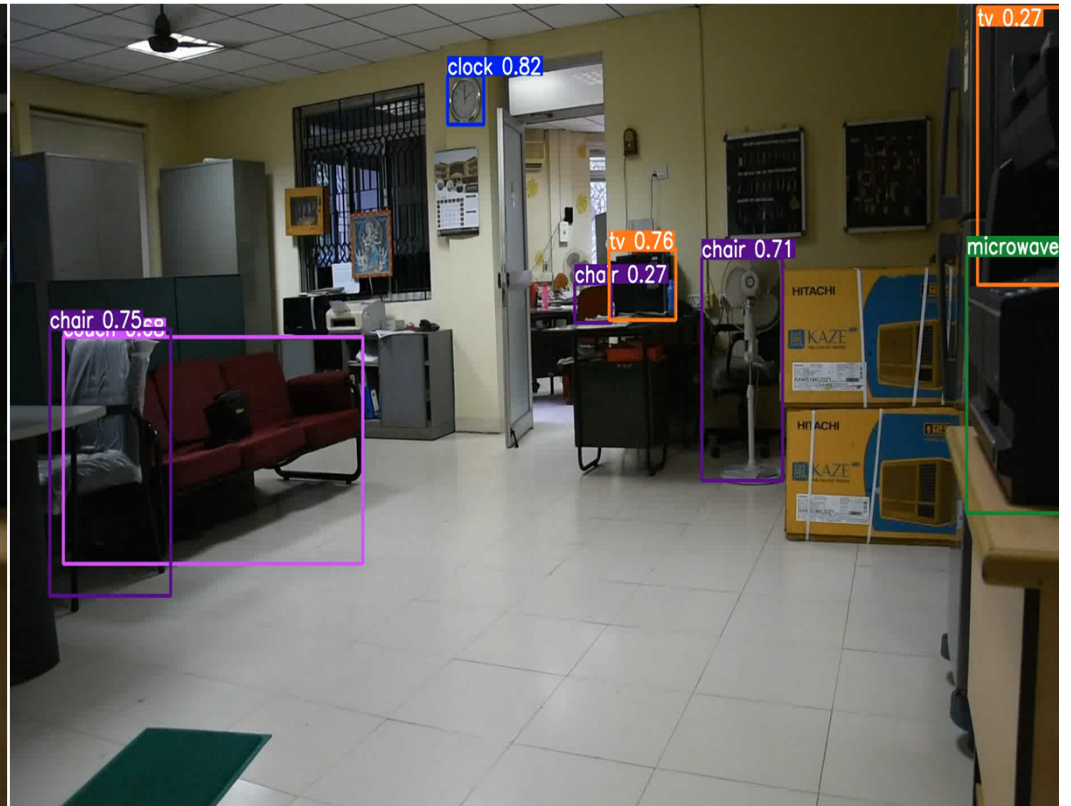
Best Student Paper Award - "Motion-based Occlusion-aware Pixel Graph Network for Video Object Segmentation", Saptakatha Adak and Sukhendu Das; In 26th International Conference on Neural Information Processing (**ICONIP, Rank A**), Sydney, Australia, December 12-15, 2019.



Heavy Version of YOLOV5



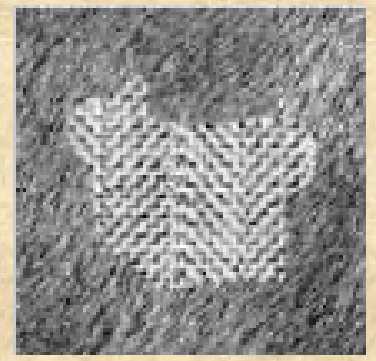
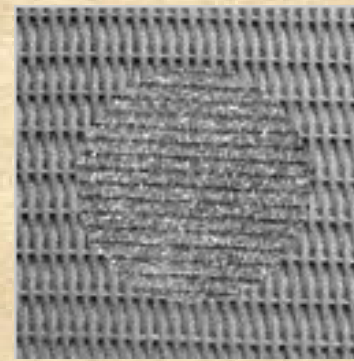
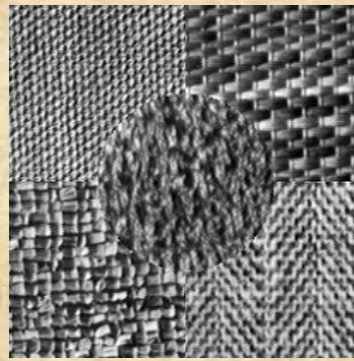
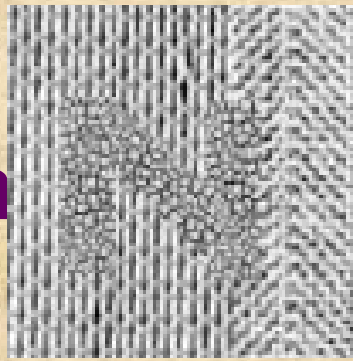
Light Version of YOLOV5



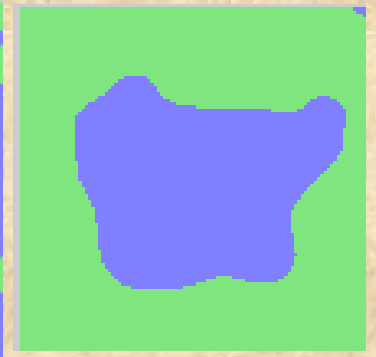
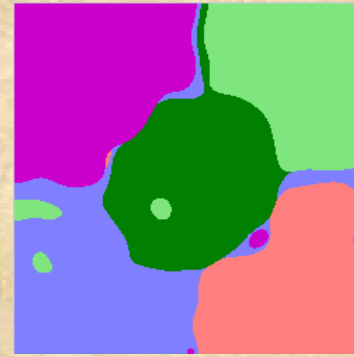
SI	Process	GPU - NVIDIA GeForce RTX 2080	CPU CORE i7 8th Generation
1	Yolov5 - Heavy	40 fps	-
2	Yolov5 – Light	149 fps	18 fps
3	Yolov7 (2022)	23 fps	-

Results of Segmentation

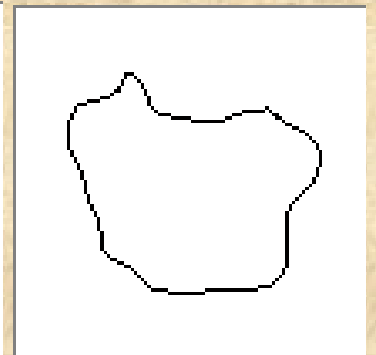
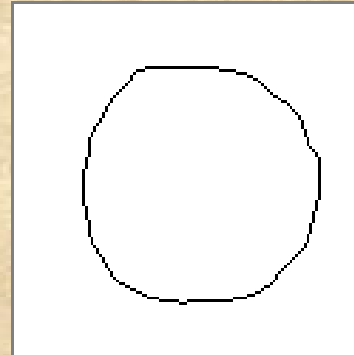
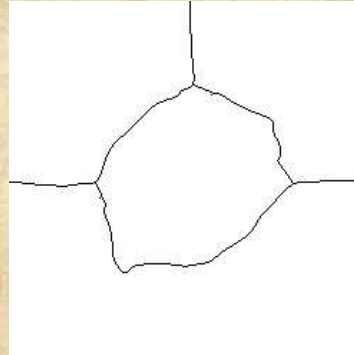
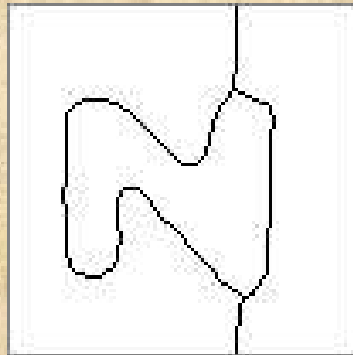
Input Image



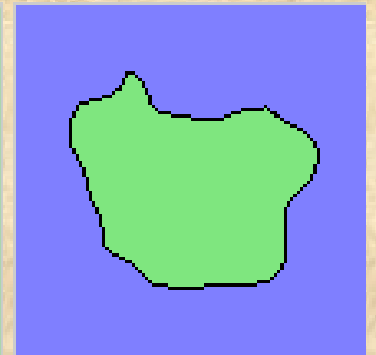
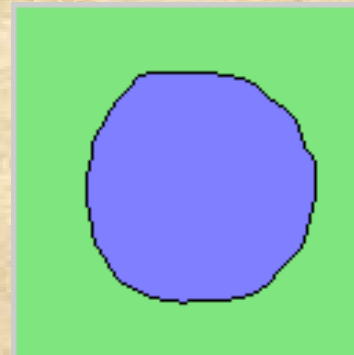
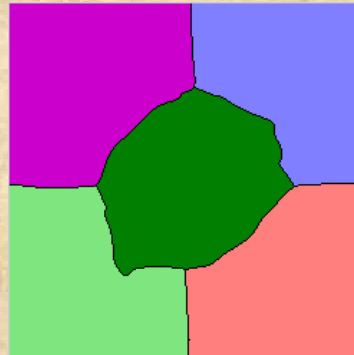
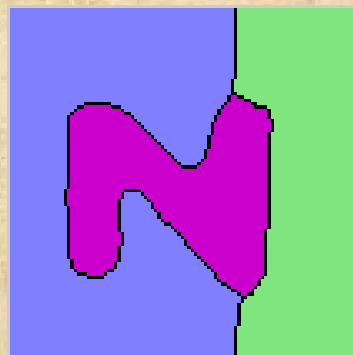
Segmented map
before integration



Edge map before
integration

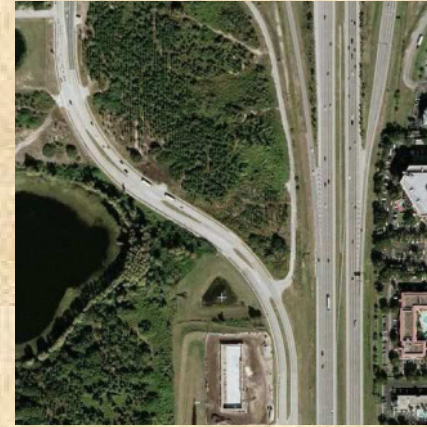


Segmented map
and Edge map
after integration

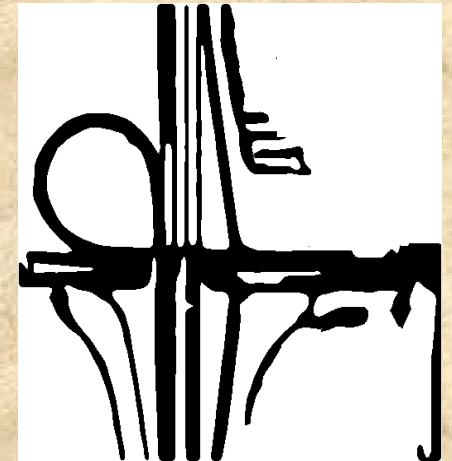
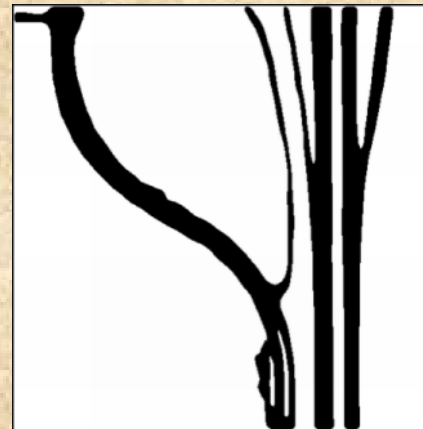
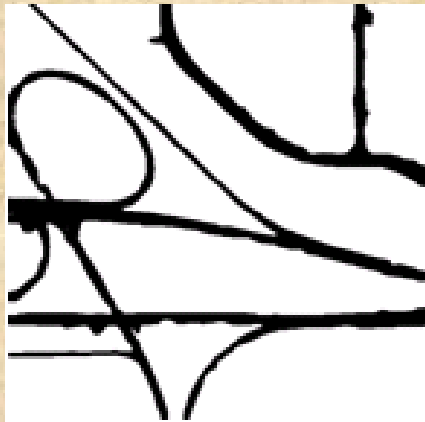


Road extraction from Satellite Images

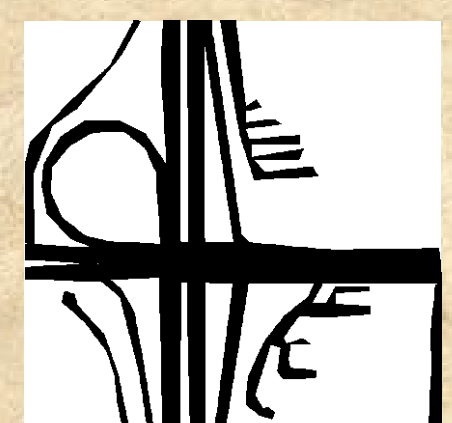
SAT
Images



Results

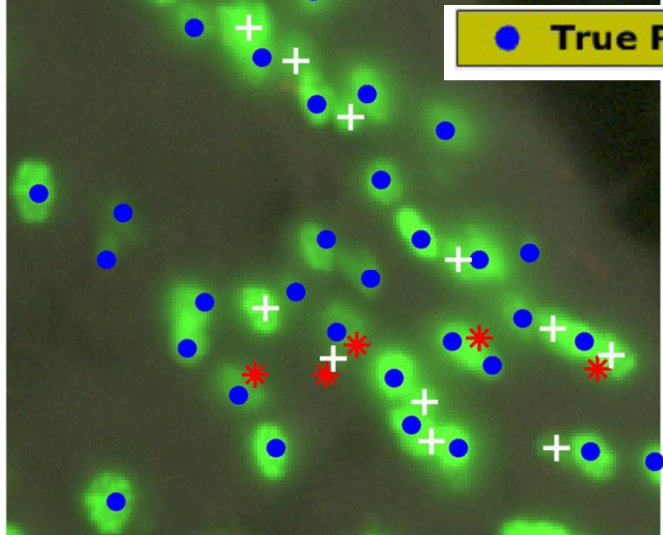
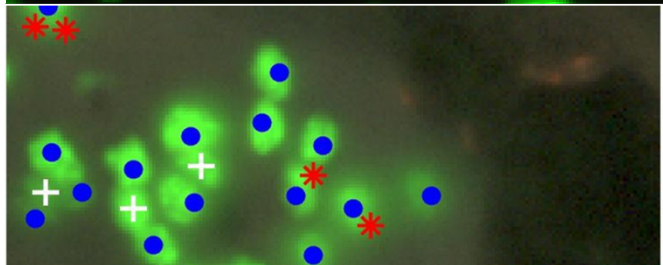
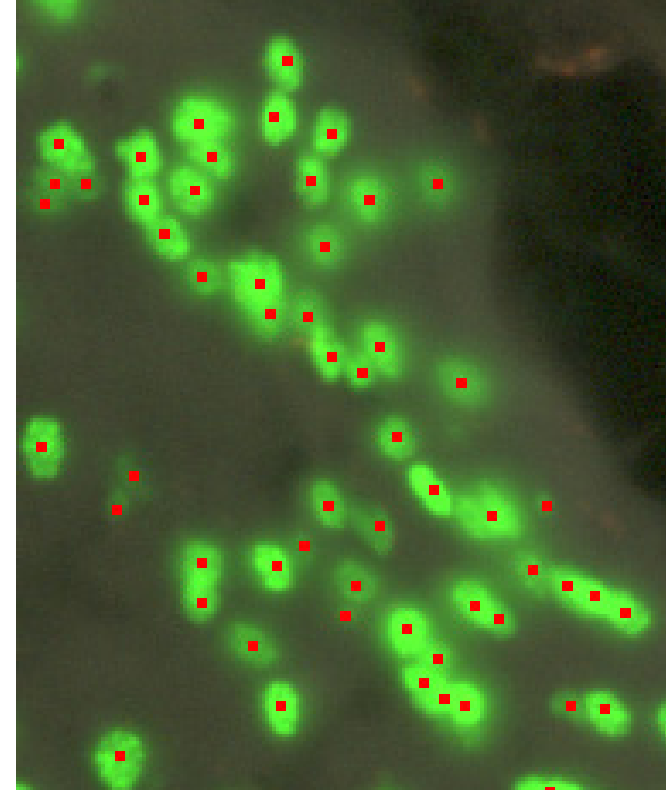
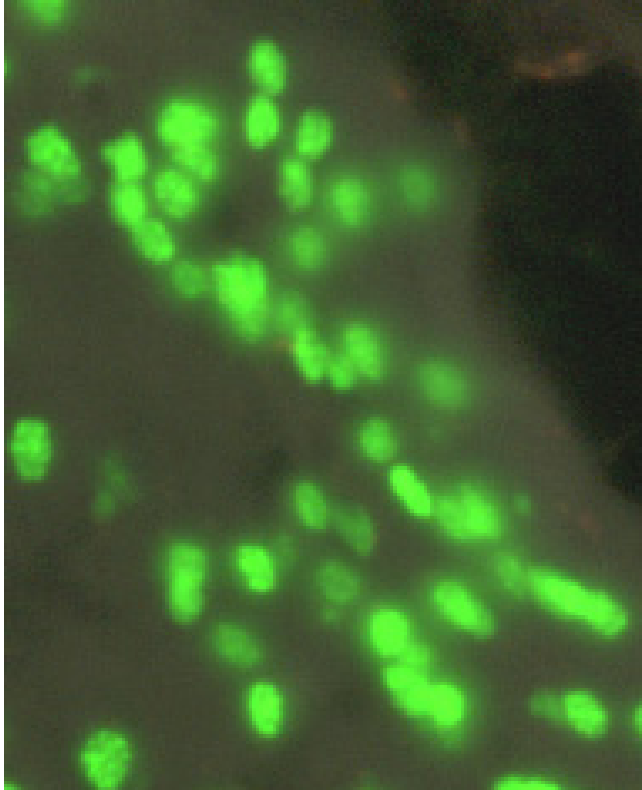
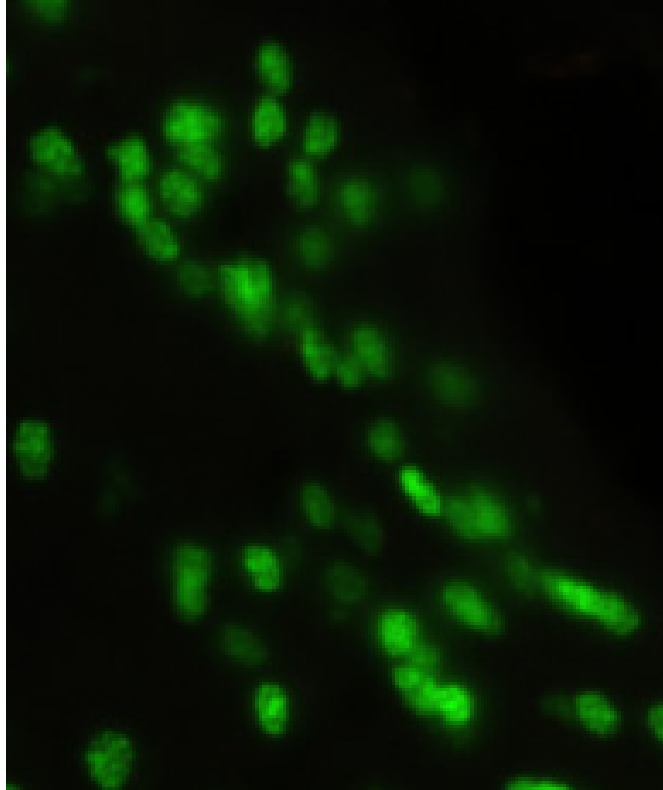


Hand-
drawn



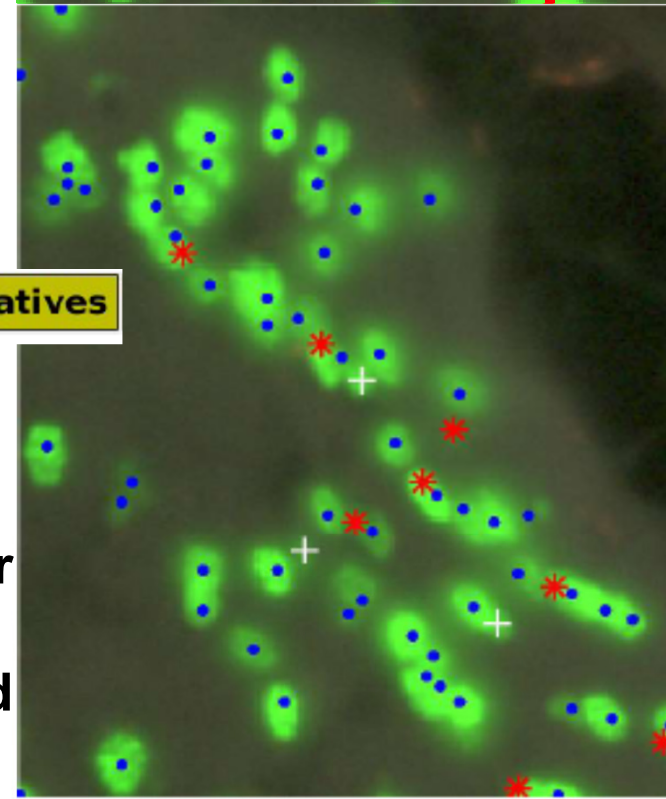
Object Extraction From an Image



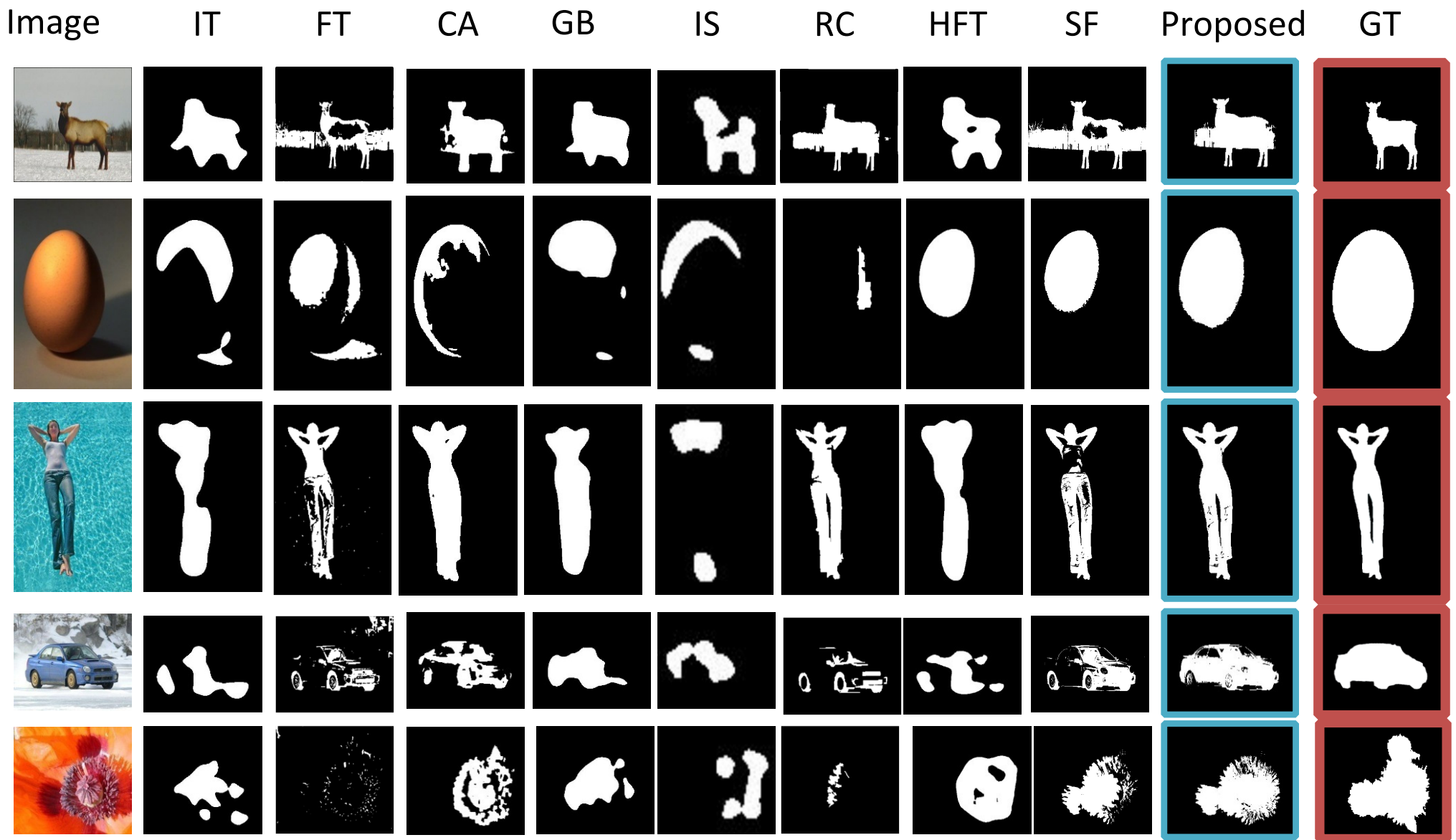


RCNN

Our
Unsupervised
method





Unsupervised Saliency

Images from MSRA B 5000 image Dataset

<http://research.microsoft.com/en-us/um/people/jiansun/SalientObject/salient object.htm>


Oct 24, 2014



Visual Results on PASCAL

Image

SF

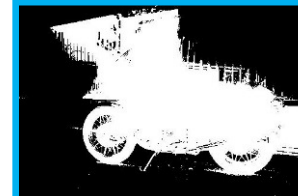
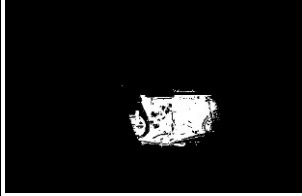
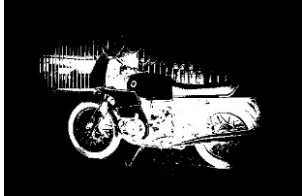
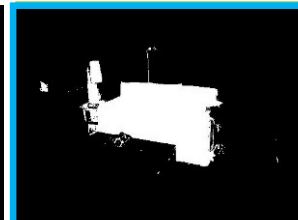
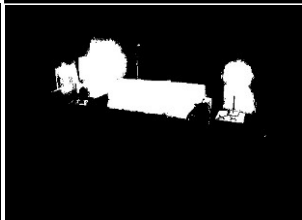
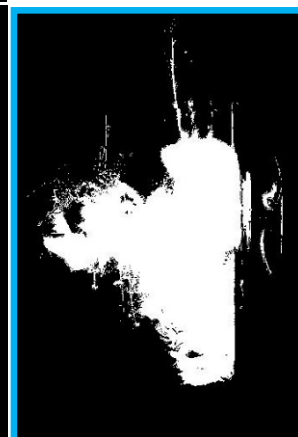
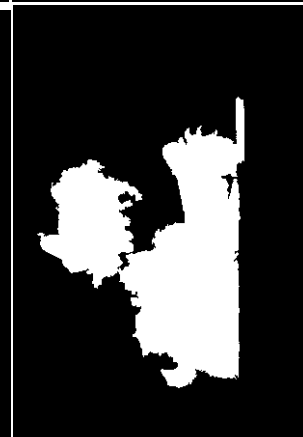
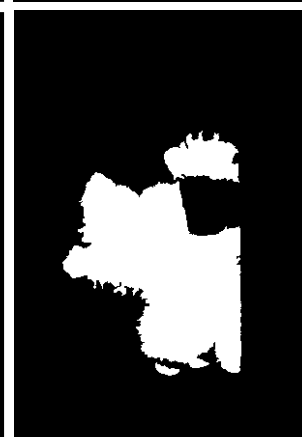
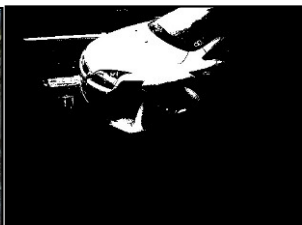
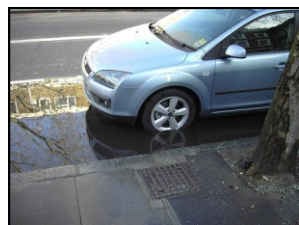
PARAM

MR

wCrt

Proposed

GT



Snake
Output

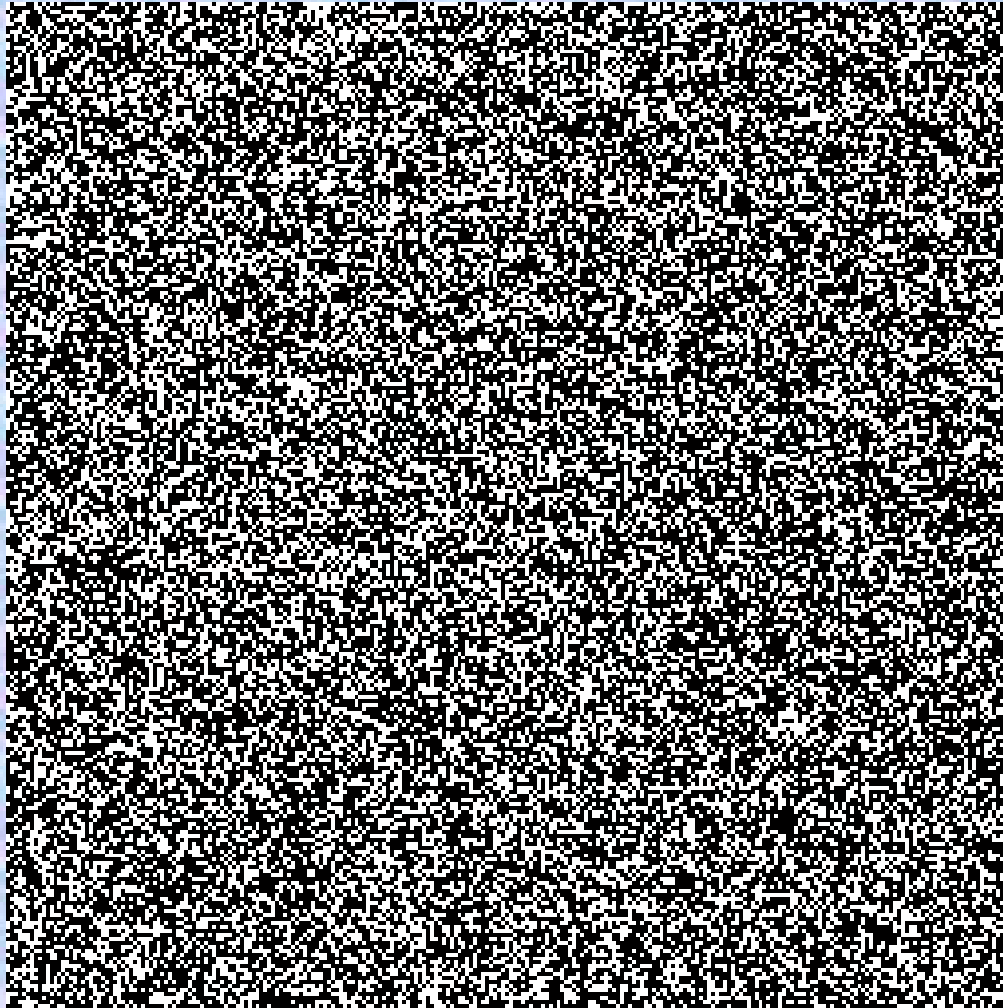


GrabCut Output

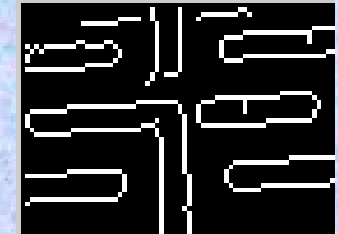


SnakeCut
Output

The Problem Definition



IMRN



IMT

Given a bitmap template (IMT) and a noisy bitmap image IMRN which contains IMT (believe me):

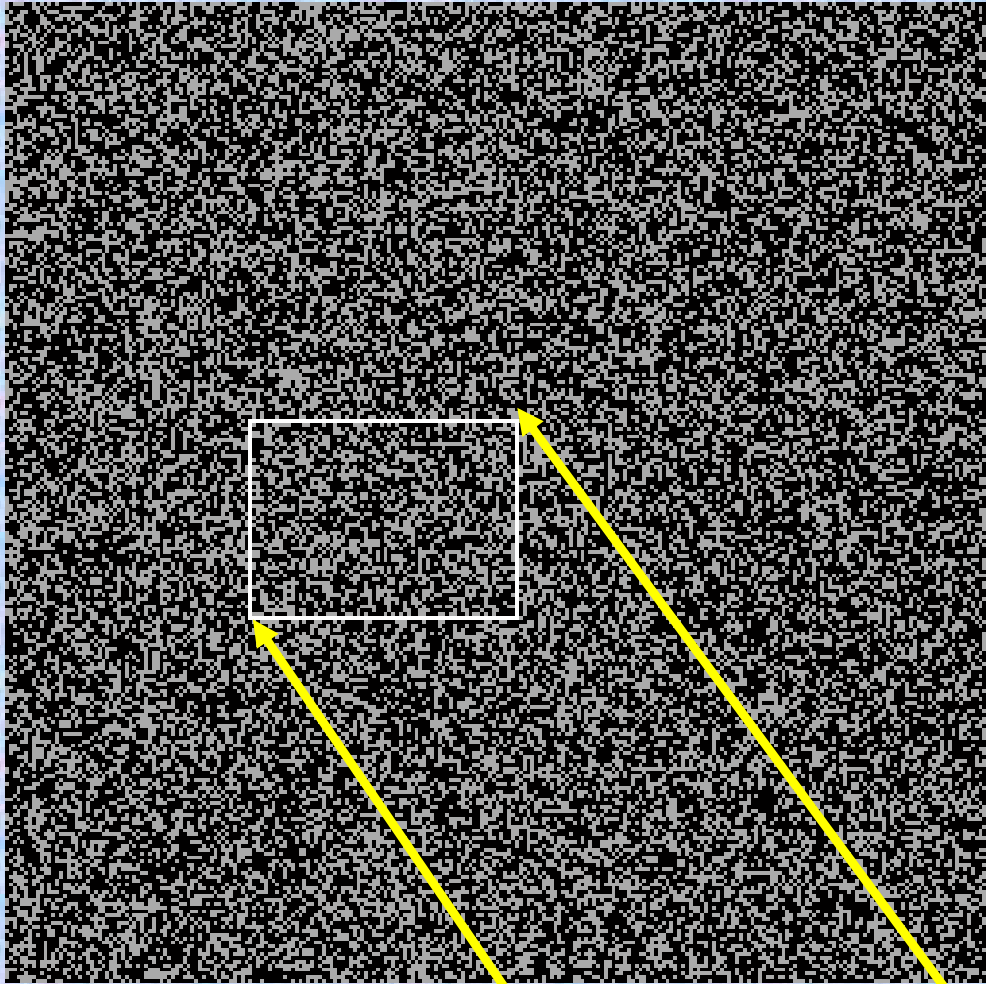
FIND OUT the location of IMT in IMRN !

Problem explanation for pessimists.



- IMRN (in previous page) is obtained by adding a large level of “Salt and Pepper” noise onto IMR bitmap image.
- IMT is also obtained from IMR as shown above.

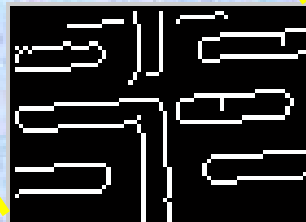
The RESULT beats the human EYE



IMRN



IMR



IMT

Published almost 3 decades ago;
Without GPU and DL

Thank you

