



ABV-Indian Institute of Information Technology and Management  
Gwalior

---

# Course on Computer Vision (ITIT-9507)

---

Instructor – Dr. Sunil Kumar

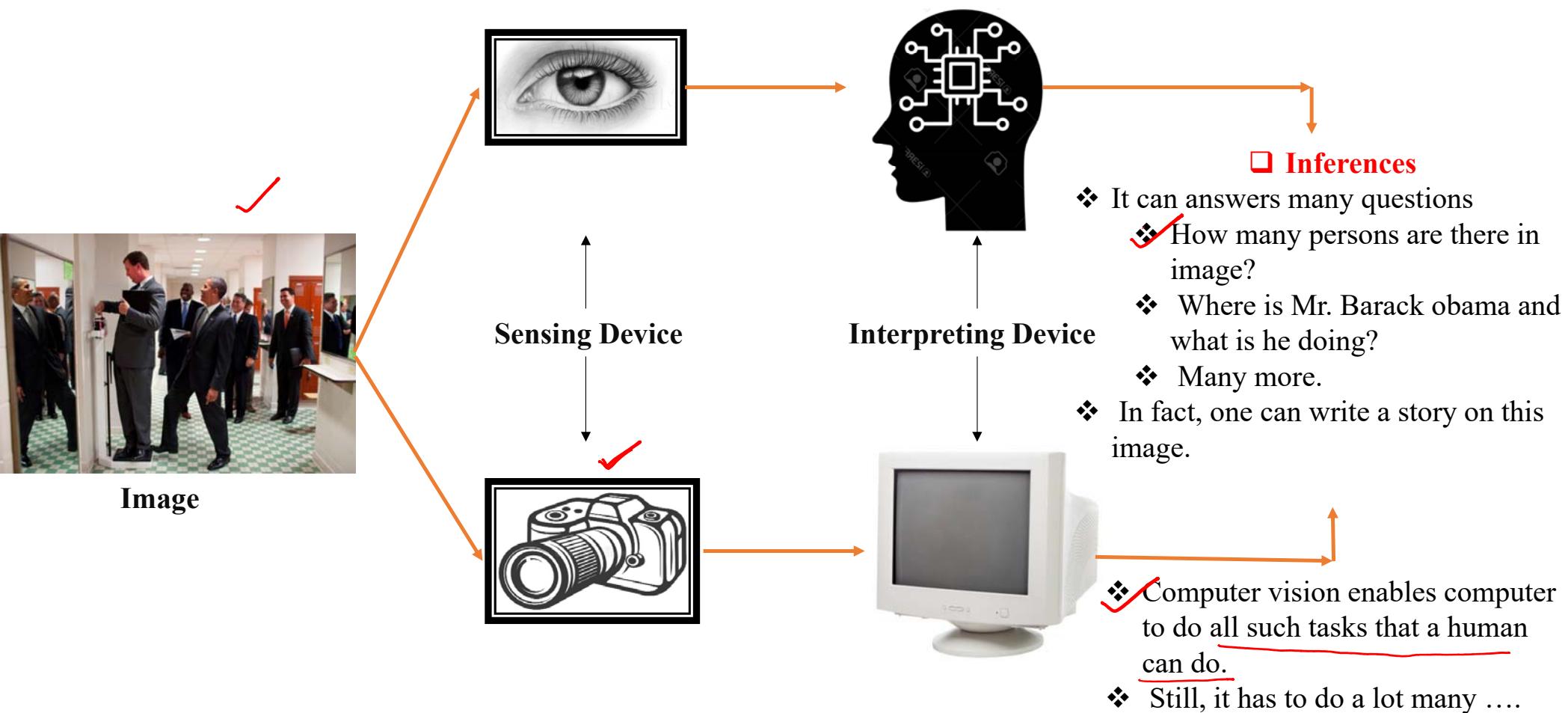
Office – 206, F-Block (V), Tel No – 0751-2449710 (O), Email - [snk@iiitm.ac.in](mailto:snk@iiitm.ac.in)

Mob - 8472842090

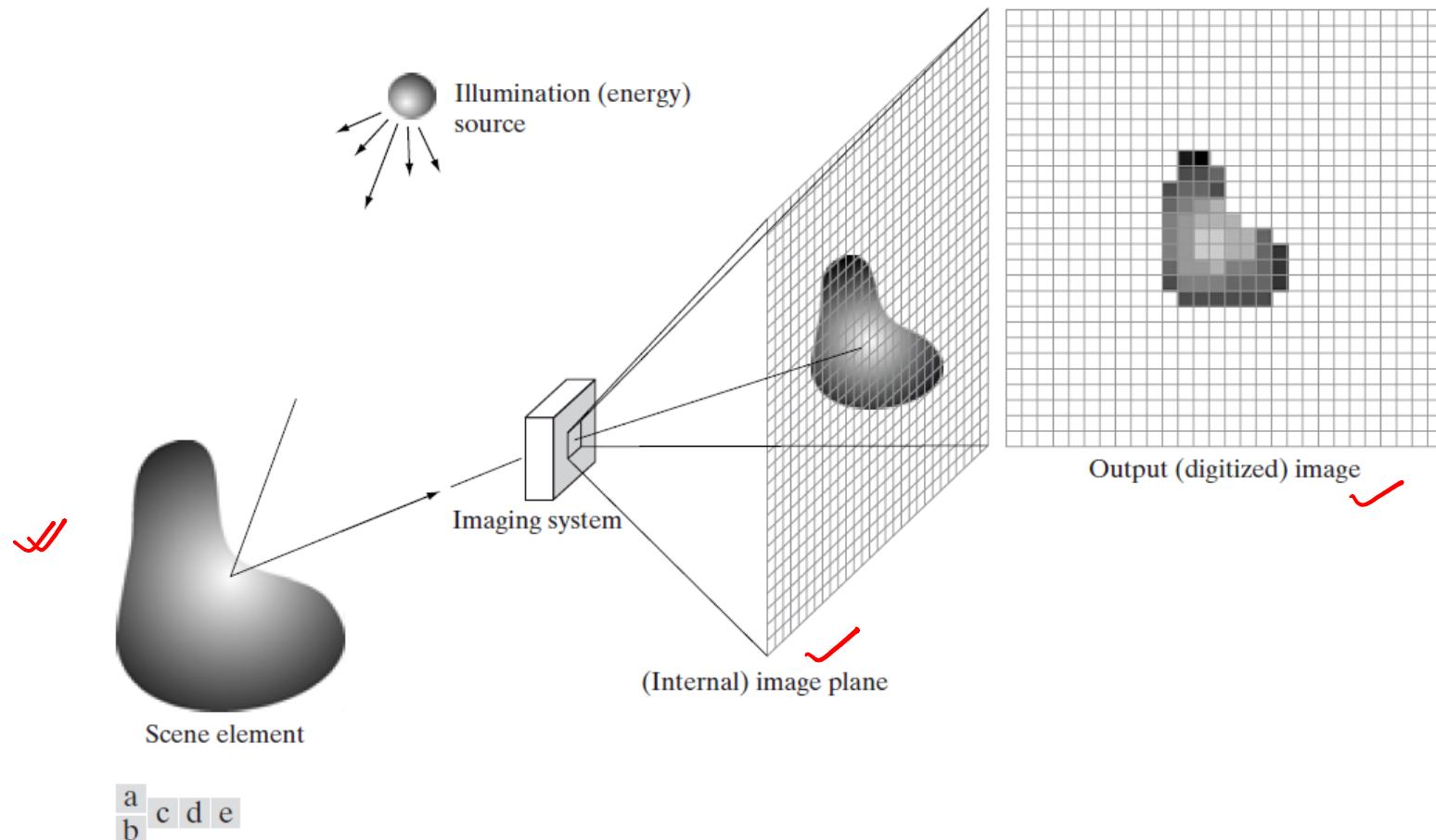
# Computer Vision

- The ability of Computer to see images and/or videos in a sense that
  - it can understand
  - Helps to take action based on what it infers
- ❖ Analogous to : Human vision system + Human brain

# Vision System Vs. Computer Vision



# What is an Image (Digital Image)?



**FIGURE 2.15** An example of the digital image acquisition process. (a) Energy (“illumination”) source. (b) An element of a scene. (c) Imaging system. (d) Projection of the scene onto the image plane. (e) Digitized image.

# Digital Image Representation



	0	1	2	3	4	5	6	7	...	N-1	
0	5	20	3	0	88	47	255	156	20	90	y-axis
1	50	210	30	10	128	40	27	61	17	25	
2	7	23	43	137	18	51	43	29	89	86	
3	42	42	63	96	128	48	127	161	127	203	
4	50	210	30	10	128	40	27	61	17	26	
5	56	0	30	77	256	52	224	58	66	83	
M-1	15	21	14	53	128	40	27	61	17	222	

Figure- 1 : ABV-IIITM Gwalior sports

# Digital Image Representation



5	20	3	0	88	47	255	156	20	90
50	210	30	10	128	40	27	61	17	25
7	23	43	137	18	51	43	29	89	86
42	42	63	96	128	48	127	161	127	203
50	210	30	10	128	40	27	61	17	26
56	0	30	77	256	52	224	58	66	83
15	21	14	53	128	40	27	61	17	222
15	21	14	53	128	40	27	61	17	222

Figure- 1 : ABV-IIITM Gwalior sports

Resolution : 356 x 512 = No. of rows \* No. of columns

~~8 bits → 256 levels  
2 bits → 4 levels~~

## Types of a Digital Image

Grayscale image : 0 to 255 for 8-bits



✓ Color image : [R-plane, G-plane, B-plane]



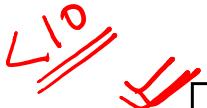
- Grayscale image with intensity value 0 (black) and 1 (white) Called “black and white image”.

# ✓ Types of a Digital Image

RGB Color image : [R-plane, G-plane, B-plane]



# Video



- Defined by a sequence of images (also called frames).
- Frame rate : 30 frames/second

Class: dribble



Class: kick\_ball



.....



# Motivation of this Course

- Why do you need to study this course?
- Billions of images/videos captured per day



flickr



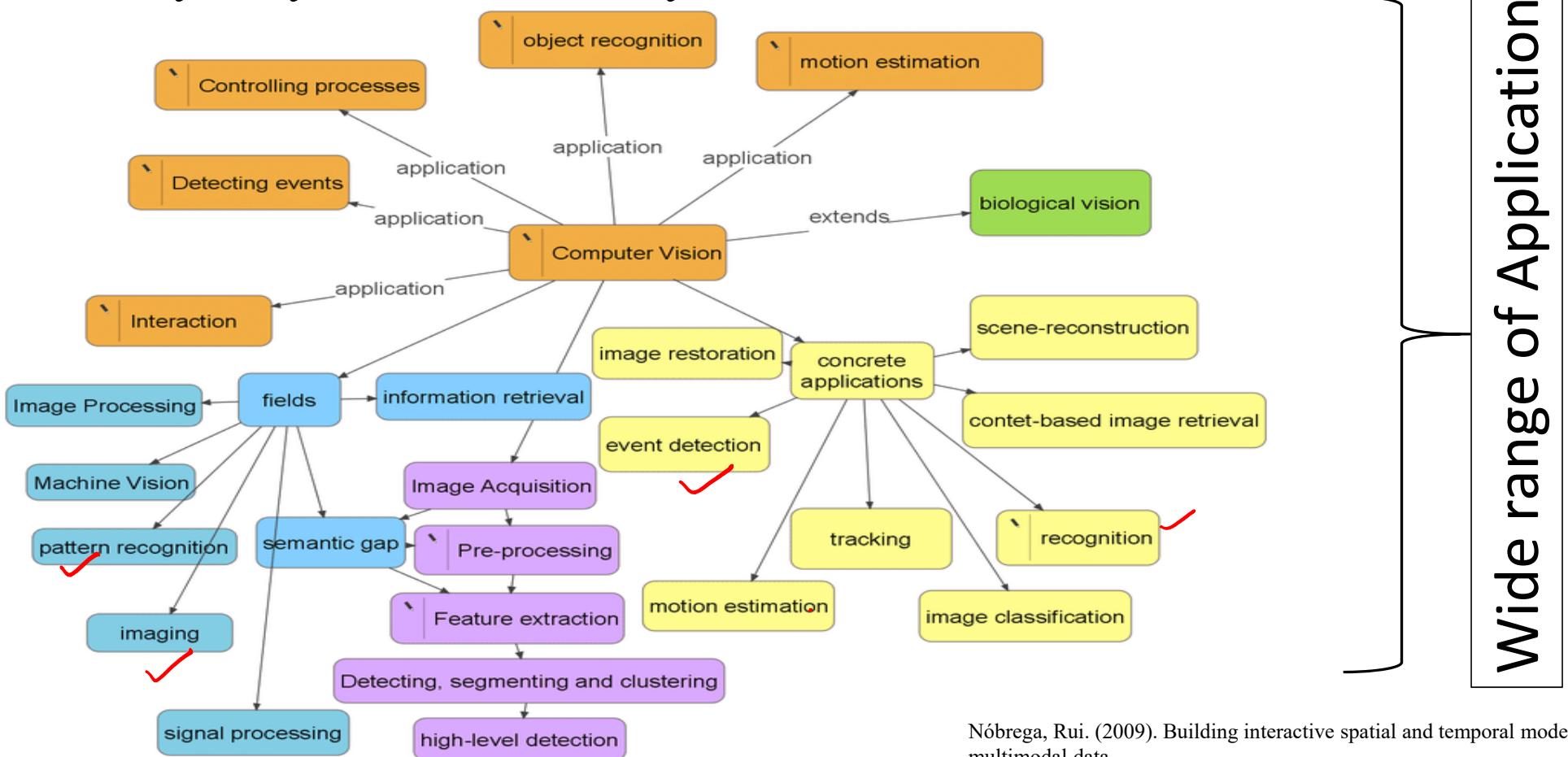
You Tube  
Broadcast Yourself™



- Huge number of useful applications

# Motivation of this Course

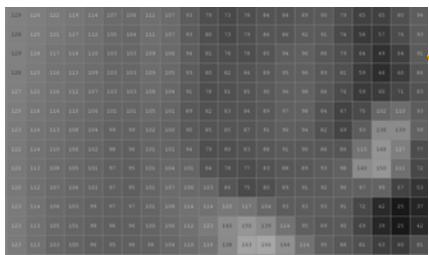
- Why do you need to study this course?



Nóbrega, Rui. (2009). Building interactive spatial and temporal models using multimodal data.

# Goal of Computer Vision

- To bridge the gap between “image pixels” and its “semantic meaning”



145	152	156	155	161	169	170	164	163	168	146	145	142	134	124	119	121	127	113	115
146	170	168	173	171	160	155	162	170	174	155	153	146	136	127	123	121	121	116	118
142	170	187	194	189	173	167	177	184	183	186	182	170	155	143	135	125	118	121	121
137	156	166	162	156	154	159	168	174	175	170	172	167	155	146	141	134	126	120	122
140	157	172	152	142	154	169	177	183	190	181	185	180	164	147	136	125	116	113	121
131	147	154	189	178	174	183	188	198	178	190	191	186	173	155	141	135	135	129	126
157	164	182	189	166	172	188	187	195	188	184	185	186	184	176	160	141	128	122	122
152	147	177	165	147	173	200	195	199	205	194	183	170	162	157	147	133	121	118	121
133	122	149	145	148	181	206	196	189	197	206	195	178	163	152	144	139	135	120	125
139	135	147	158	175	188	195	180	161	165	179	186	190	183	165	146	133	127	124	126

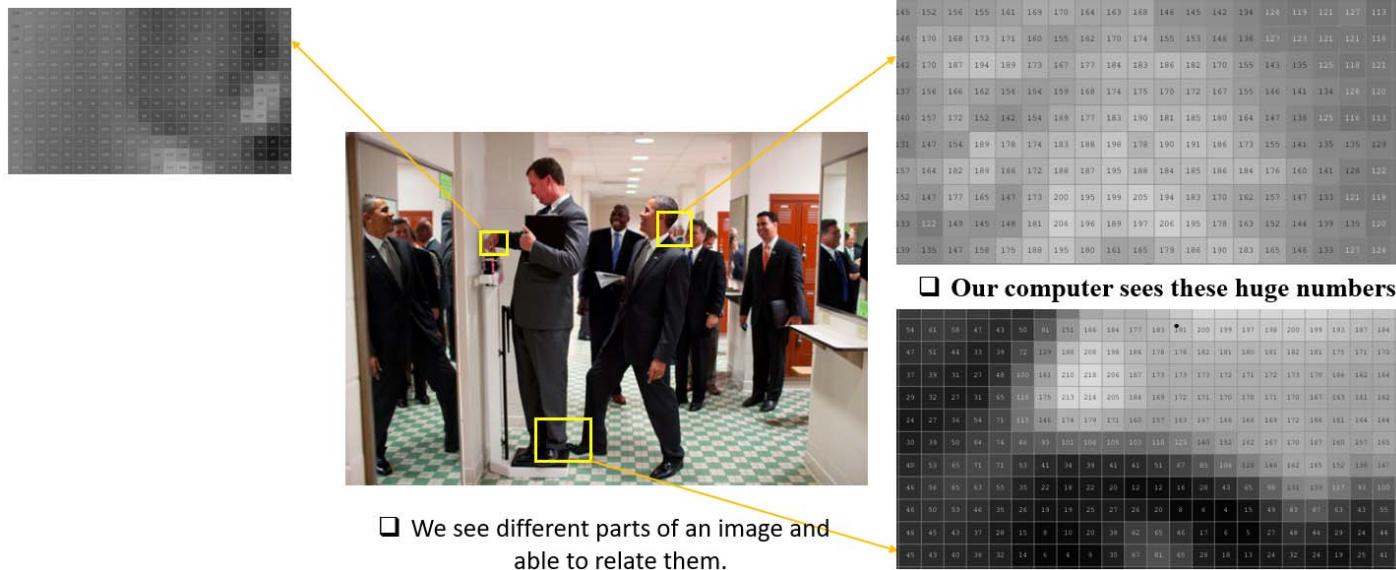
- Our computer sees these huge numbers.

54	61	58	47	43	50	91	151	186	184	177	183	•91	200	199	197	198	200	199	193	187	184
47	51	44	33	39	72	129	188	208	198	184	178	178	182	181	180	181	182	181	175	171	170
37	39	31	27	48	100	161	210	218	206	187	173	173	173	172	171	172	173	170	166	162	164
29	32	27	31	65	118	175	213	214	205	184	169	172	171	170	170	171	170	167	163	161	162
24	27	36	54	71	115	146	174	179	171	160	157	163	167	166	166	169	172	166	161	164	164
30	39	50	64	74	86	93	101	106	105	103	110	123	140	152	162	167	170	167	160	157	165
40	53	65	71	71	53	41	34	39	41	41	51	67	85	106	128	146	162	165	152	138	147
46	56	65	63	55	35	22	18	22	20	12	12	16	28	43	65	98	131	139	117	93	100
46	50	53	46	35	26	19	19	25	27	26	20	8	6	4	15	49	83	87	63	43	55
46	45	43	37	28	15	8	10	20	38	62	65	46	17	6	5	27	48	44	29	24	44
45	43	40	36	32	14	6	4	9	30	67	81	65	28	18	13	24	32	24	19	25	41

- We see different parts of an image and able to relate them.

# Goal of Computer Vision

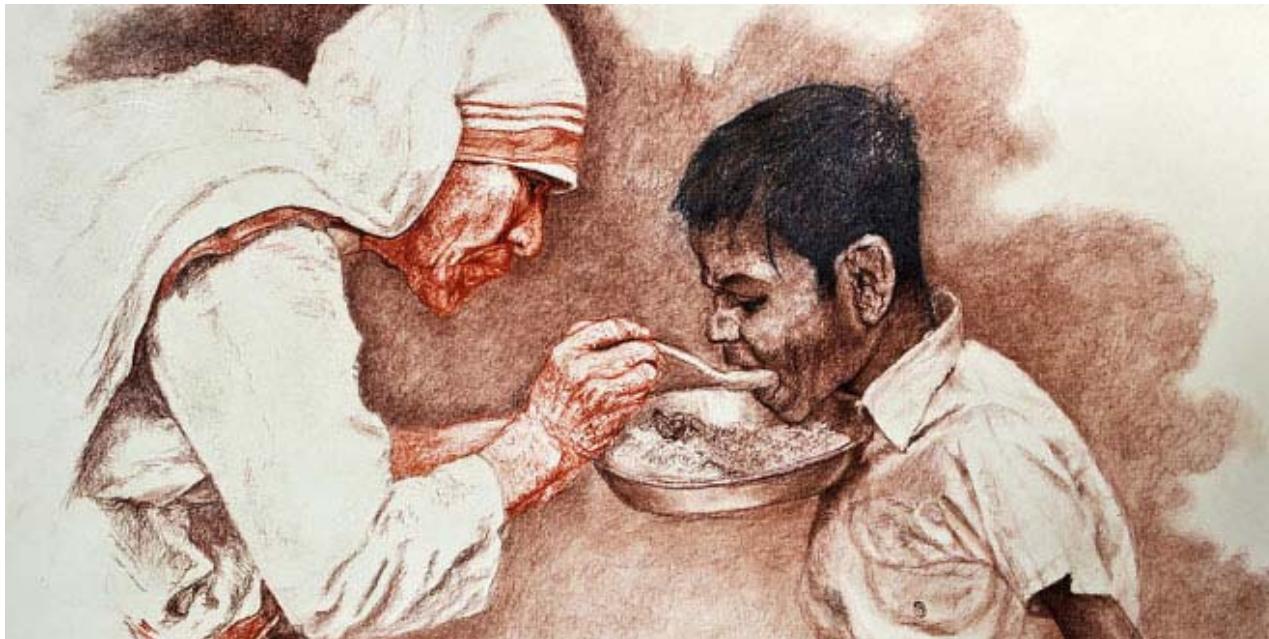
- To bridge the gap between “**image pixels**” and its “**semantic meaning**”



- In other words: The goal of Computer Vision is “to enable computers to extract meaning of every pixel in a image (frames of a given video) , relate them with every other pixels in order to interprets the image or a video.”
- So, Computer Vision mimics the human perception.

# Goal of Computer Vision

- ❑ Every image tells a story



A picture is worth a thousand words

- ❑ Can computer do so? ... It is still harder... and things are in progress

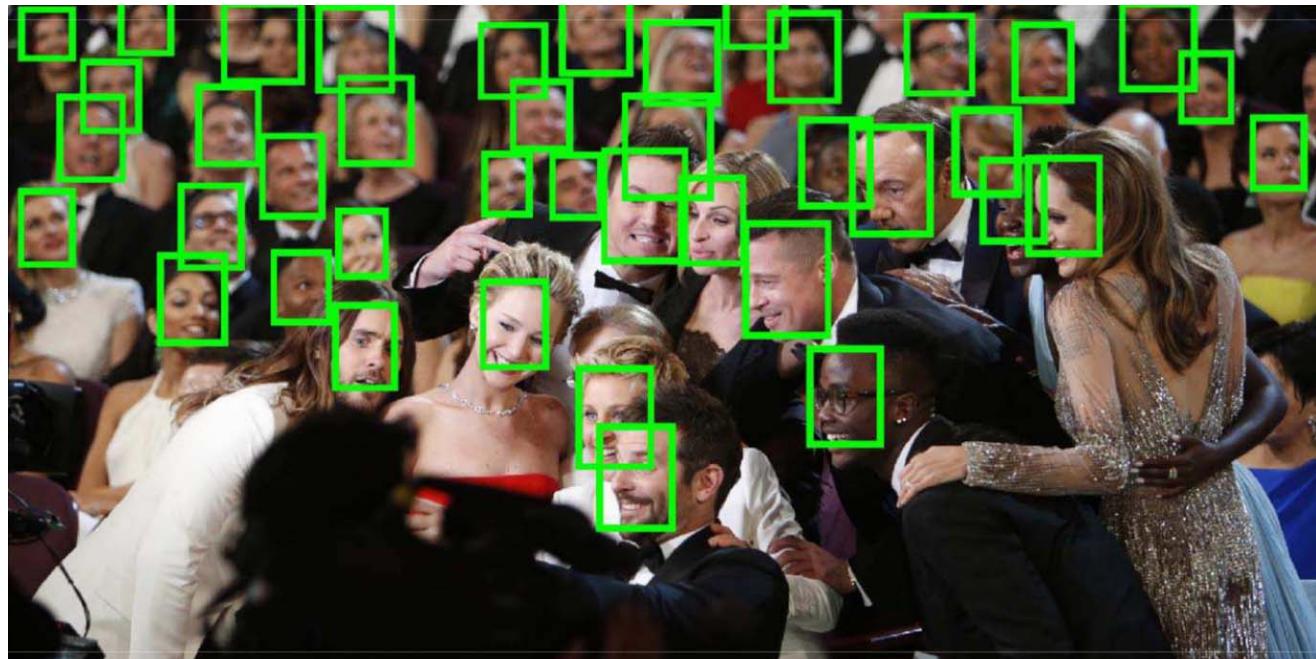
# Applications of Computer Vision

✓ There are many applications ...

- Face detection
- Human identification
- Facial expression recognition
- Object recognition and localization
- Segmentation : Medical imaging
- Robotics
- Self-driving car
- Virtual and augmented reality
- Lip reading
- Disease recognition (Covid-19)
- Person identification in crowd
- Multi-object tracking
- Action recognition
- Video surveillance
- Image and video synthesis
- Vision-based biometrics
- Many more ...

# CV Tasks : Face Detection

☛ Problem -1 : Does an image or a video contain one or more faces? If yes, localize them.

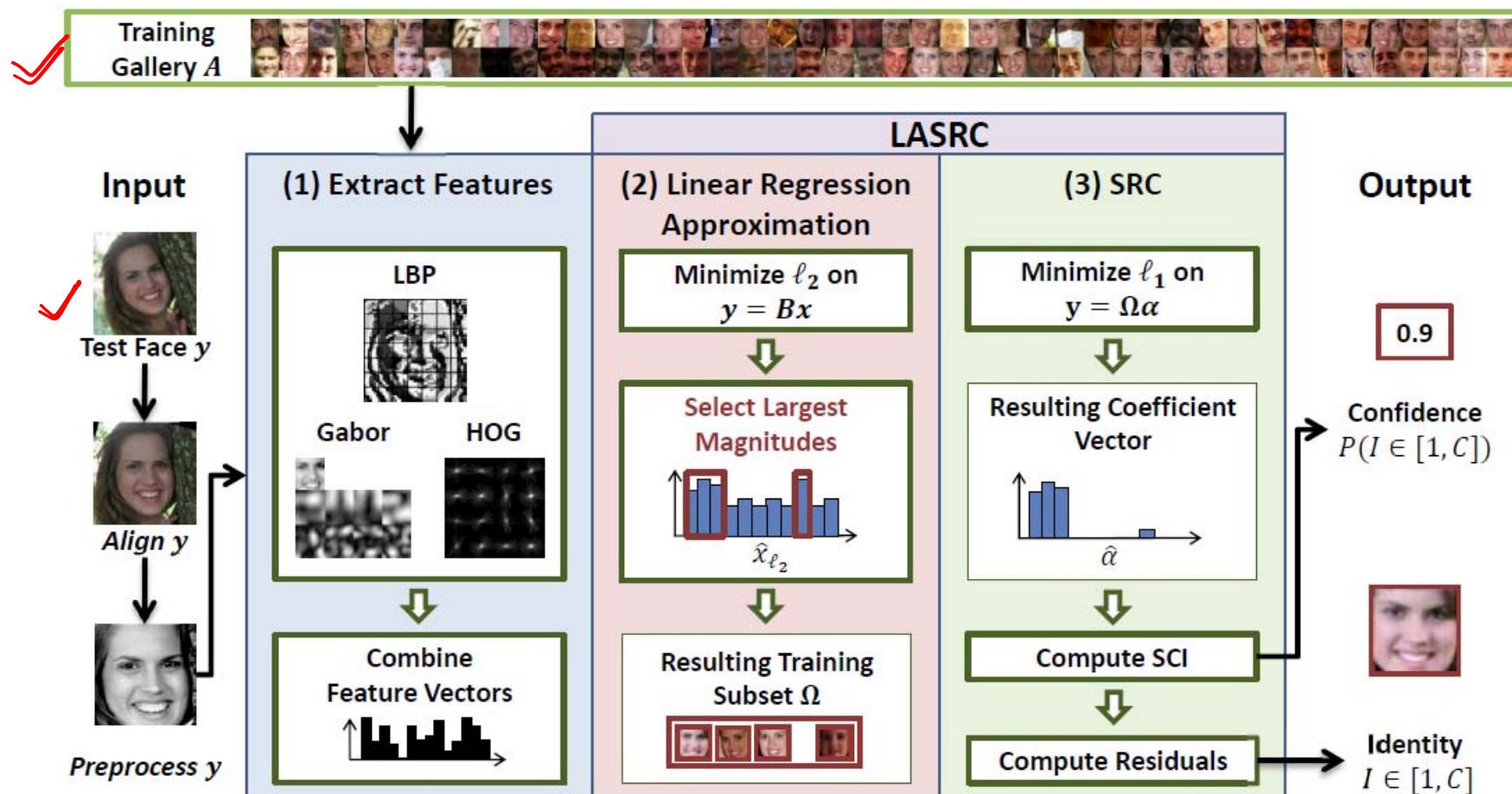


<https://modelcards.withgoogle.com/face-detection>

✓ Yang, Shuo, et al. "From facial parts responses to face detection: A deep learning approach." *Proceedings of the IEEE international conference on computer vision (ICCV)*. 2015.

# CV Tasks : Open-Universe Face Identification

Problem -2(a) : STILL-IMAGE, OPEN-UNIVERSE FACE IDENTIFICATION



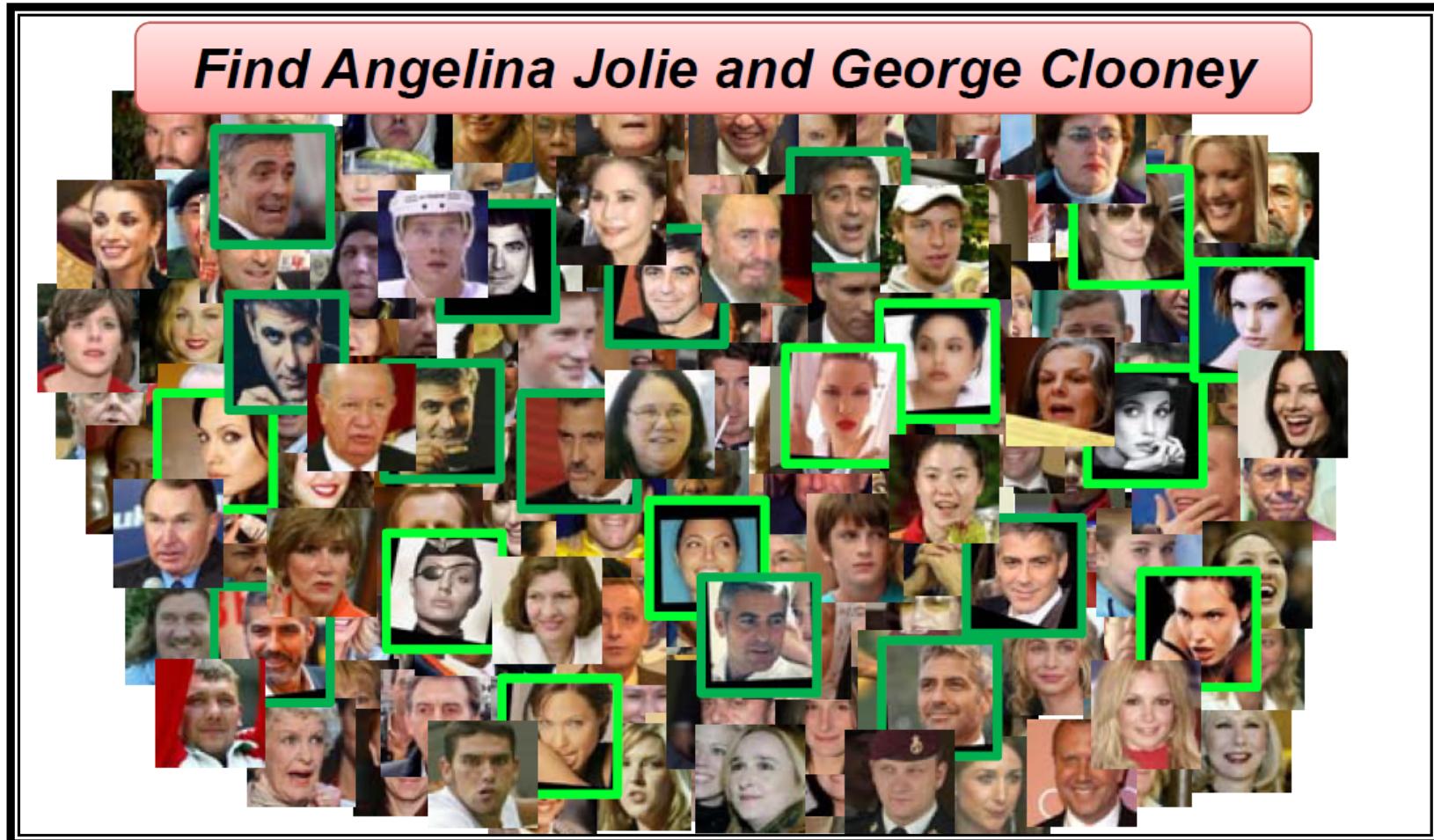
# CV Tasks : Open-Universe Face Identification

Problem -2(a) : STILL-IMAGE, OPEN-UNIVERSE FACE IDENTIFICATION



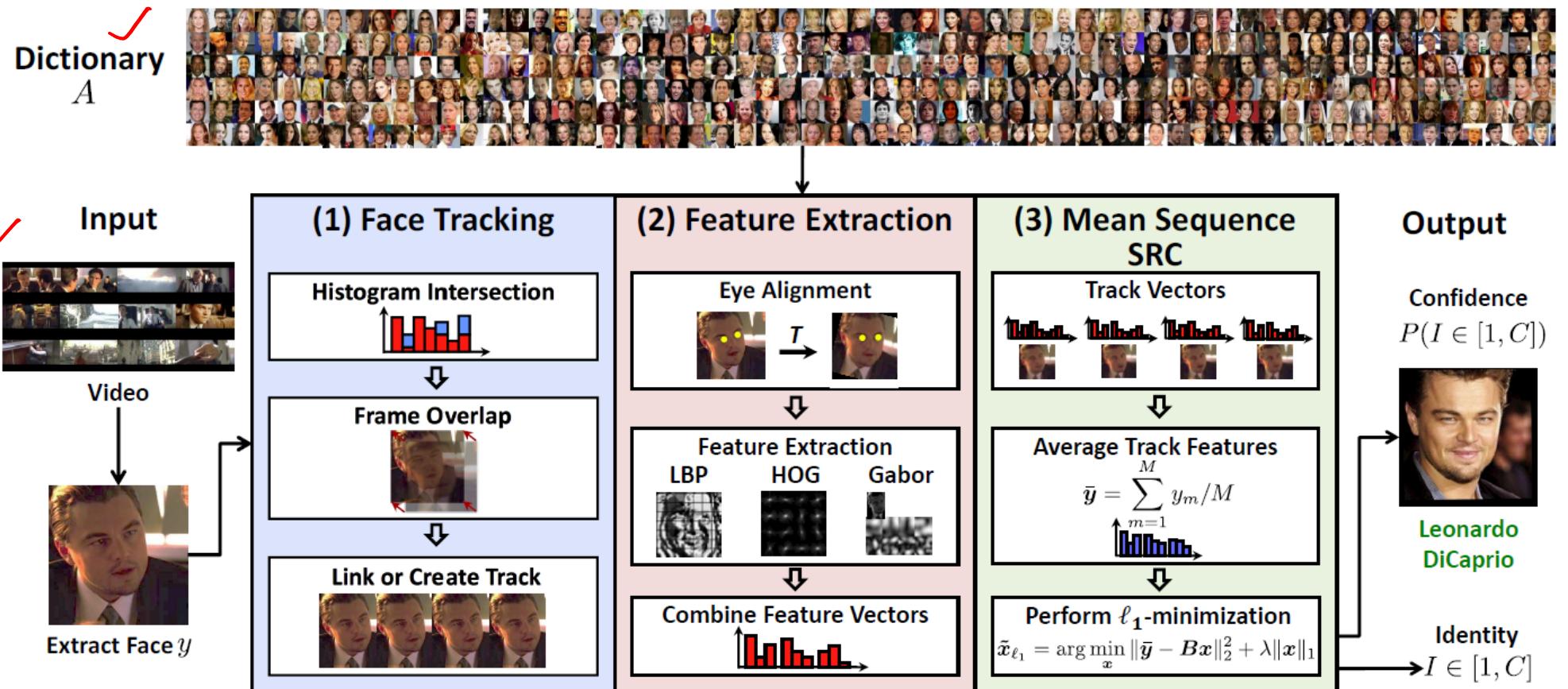
# CV Tasks : Open-Universe Face Identification

Problem -2(a) : STILL-IMAGE, OPEN-UNIVERSE FACE IDENTIFICATION



# CV Tasks : Open-Universe Face Identification

✓ Problem -2 (b) : Video, OPEN-UNIVERSE FACE IDENTIFICATION



Thesis : - Taming Wild Faces: Web-Scale, Open-Universe Face Identification in Still and Video Imagery

# CV Tasks : Object Recognition



**Problem -3** : Given an image I, does it contain image of a person?



✓  
Y  
E  
S



Andy Barron / Reno Gazette-Journal

# CV Tasks : Object Recognition

Problem -3 : Given an image I, does it contain image of a person?



Monroe County Sheriff's Department / Newsmakers



Mark Garkinkel / The Boston Herald



NATIONAL GEOGRAPHIC.COM  
© 2003 National Geographic Society. All rights reserved.



Jeff J. Mitchell / Reuters



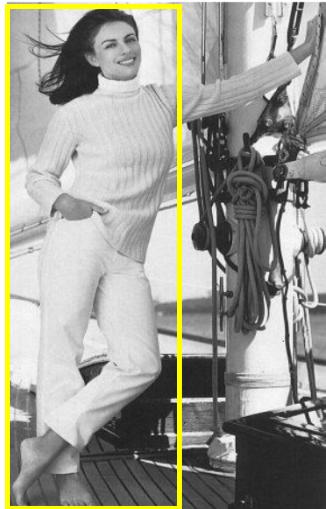
Uno Andersson / AP



✓  
N  
O

# CV Tasks : Object Localization

Problem -4 : Where is the object (say “person”) in a given image I?



# CV Tasks : Image Segmentation

- Image segmentation is the process of assigning a label to every pixel in an image in such way that pixels with the same label share certain characteristics.
- In image segmentation :
  - Input – image
  - Output – regions, structures

Example-1



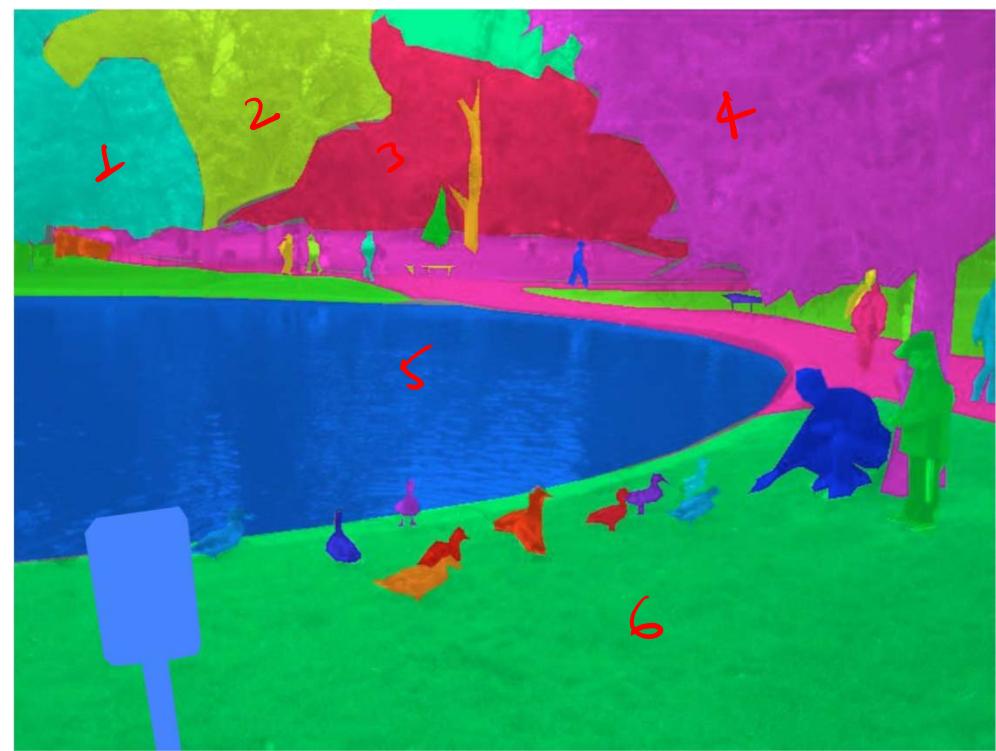
# CV Tasks : Image Segmentation

Problem -5 :

- In image segmentation :
  - Input – image
  - Output – regions, structures



(a) : Input image



Example-2

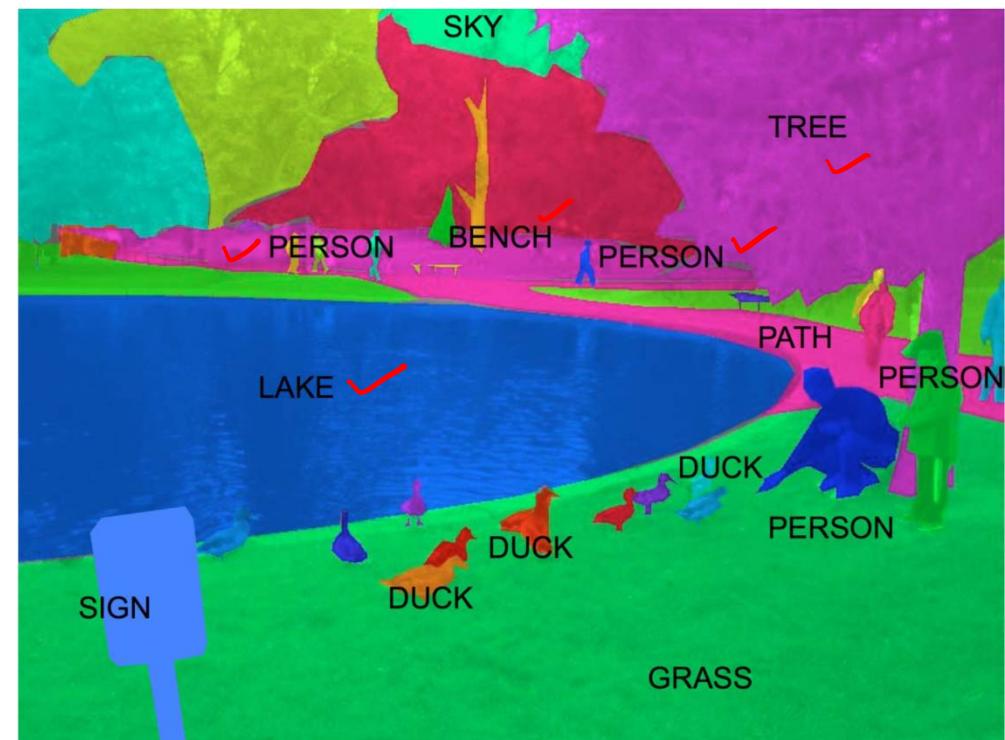
(b) Output labelled image

# CV Tasks : Semantic Segmentation

- ❑ Problem -6 : In semantic image segmentation,
- ❑ Idea : recognizing, understanding what is in the image at pixel level.
- ❑ It is a challenging task due to many reasons.



(a) : Input image



(b) Output : recognized different objects with their labels

# CV Tasks : Semantic Segmentation

- ❑ Problem -7 : In semantic image segmentation,
- ❑ Idea : recognizing, understanding what is in the image at pixel level.
- ❑ Going one more deeper level and try to understand what is happening in the image



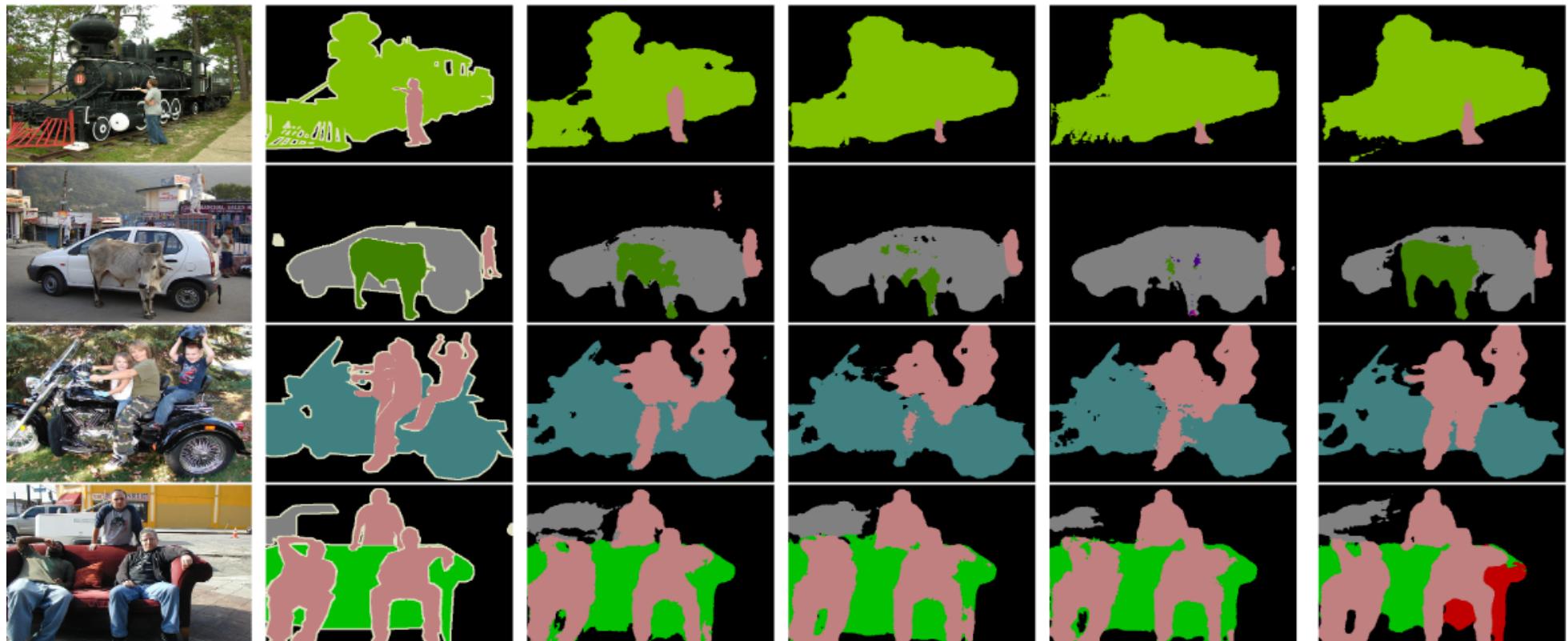
(a) : Input image



(b) Output Understanding from the image

# CV Tasks : Semantic Segmentation

- ❑ Problem -7 : In semantic image segmentation,
- ❑ Idea : recognizing, understanding what is in the image at pixel level.



✓ Ibrahim, Mostafa S., et al. "Semi-supervised semantic image segmentation with self-correcting networks." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2020. ✓

# CV Tasks : Counting in Extremely Dense Crowd Images

□ Problem -8 : Count the no. of people in a given extremely dense crowed image



Ground truth/proposed [1] =  $673/\underline{\underline{653}}$

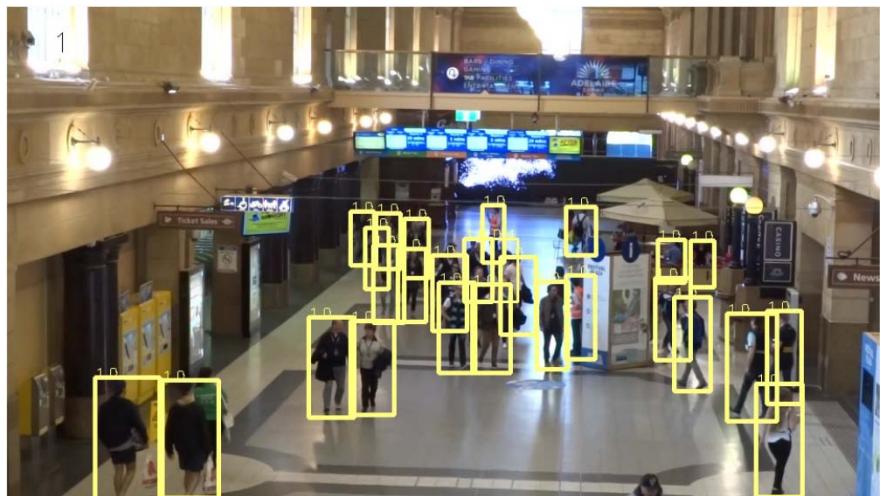
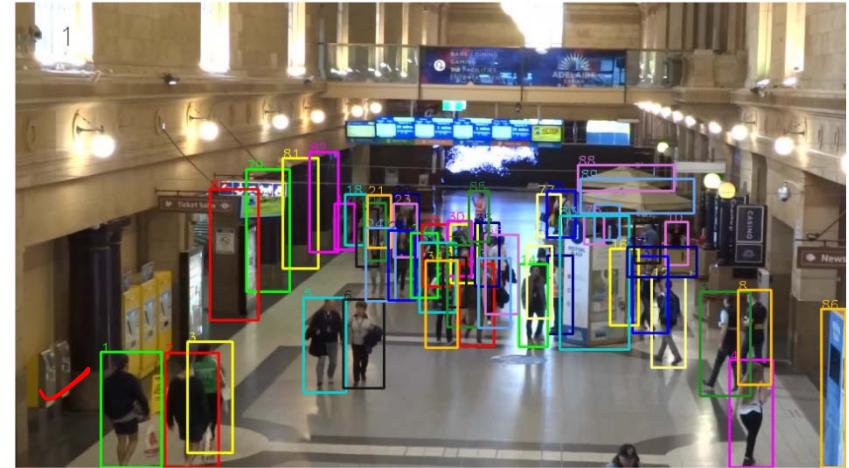


Ground truth/proposed [1] =  $2322/\underline{\underline{2203}}$

✓ [1] Idrees, Haroon, et al. "Multi-source multi-scale counting in extremely dense crowd images." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2013.

# CV Tasks : MOT20 Dataset ✓

✓ Problem -9 : A benchmark for multi object tracking in crowded scenes



- (a) Image frame of a video
- (b) Ground truth
- (c) Faster-RCNN

# CV Tasks : Multiple-Object Tracking

Problem -9 : Multiple Object Tracking with Correlation Learning

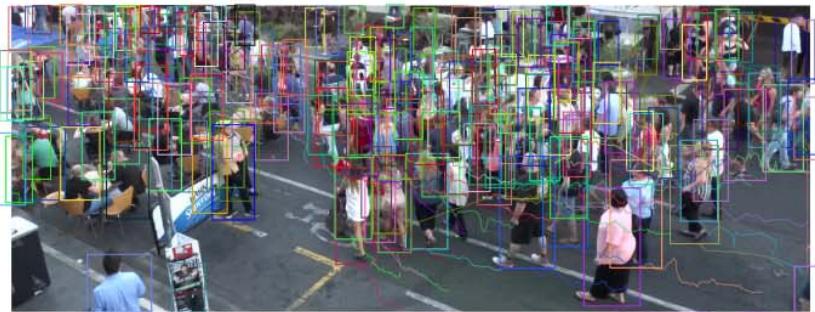
MOT17-01



MOT17-12



MOT20-06



MOT17-07



MOT17-14



MOT20-08



✓ Wang, Qiang, et al. "Multiple Object Tracking with Correlation Learning." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021. ✓

# CV Tasks : Multiple-Object Tracking

**Problem -9 :** Discriminative Appearance Modeling with Multi-track Pooling for Real-time Multi-object Tracking

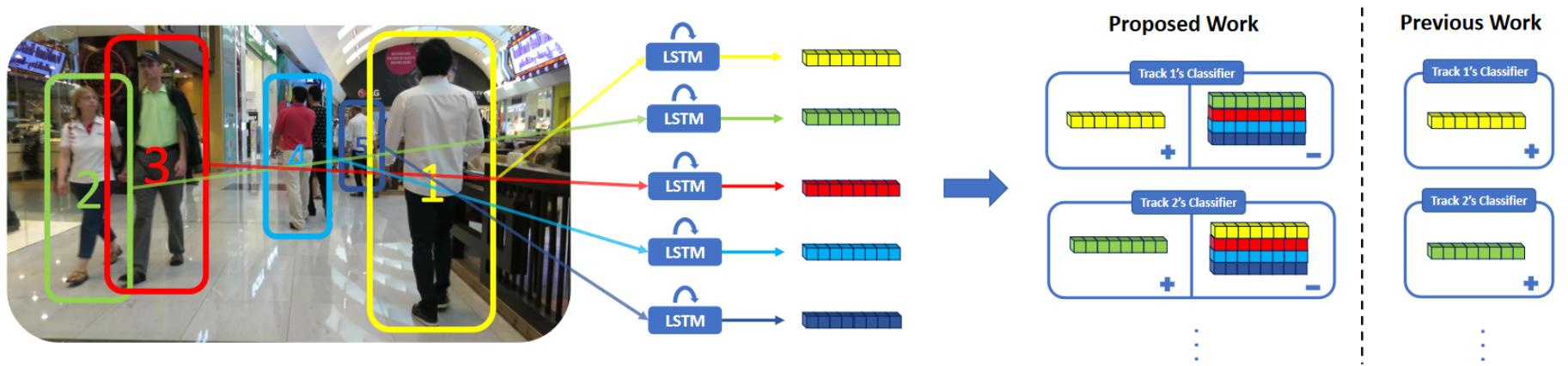


Figure 1: Existing recurrent neural network-based track classifiers used only matched detections for updating its appearance memory during tracking. This does not consider other objects in the scene (i.e. negative examples), which may have similar appearances. We propose to improve the predicted likelihood of such a classifier by augmenting its memory with appearance information about other tracks in the scene with multi-track pooling, leveraging the appearance information from the full set of tracks in the scene. The resulting classifier learns to adapt its prediction based on the information from other tracks in the scene.

Kim, Chanho, et al. "Discriminative Appearance Modeling with Multi-track Pooling for Real-time Multi-object Tracking." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021.

The code and trained models are available at <https://github.com/chkim403/blstm-mtp>.

# CV Tasks : Multiple-Object Tracking

Problem -10 : Detection, Tracking, and Counting Meets Drones in Crowds: A Benchmark

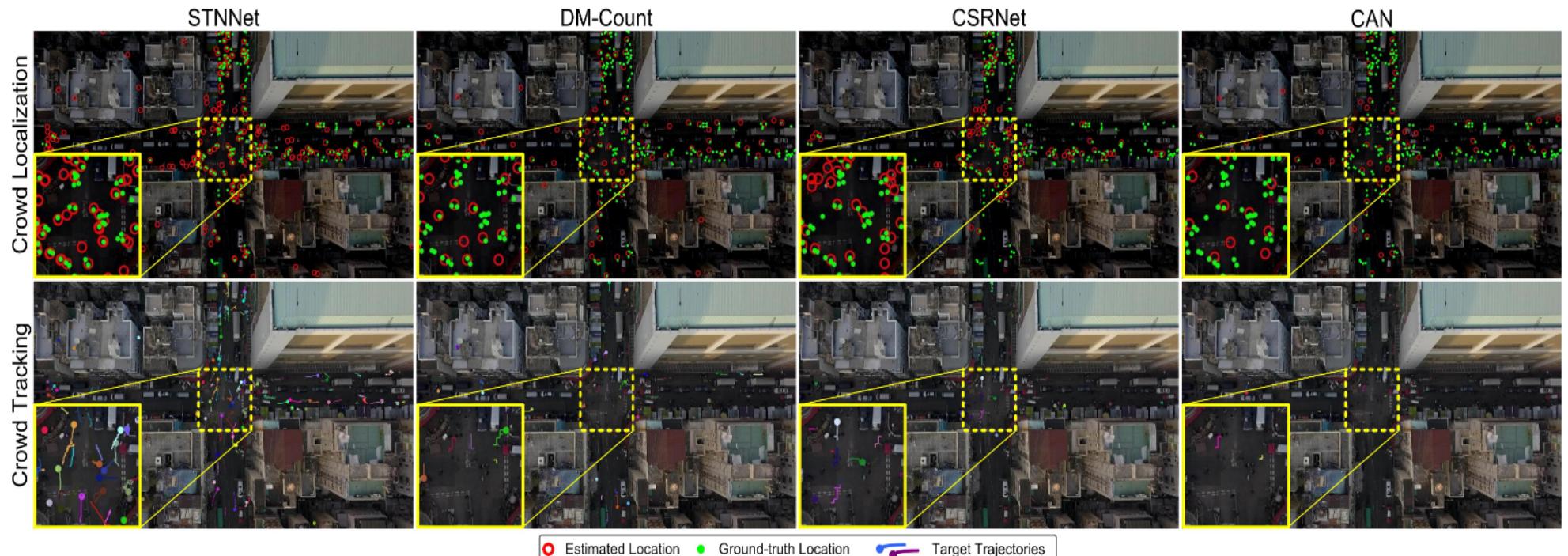
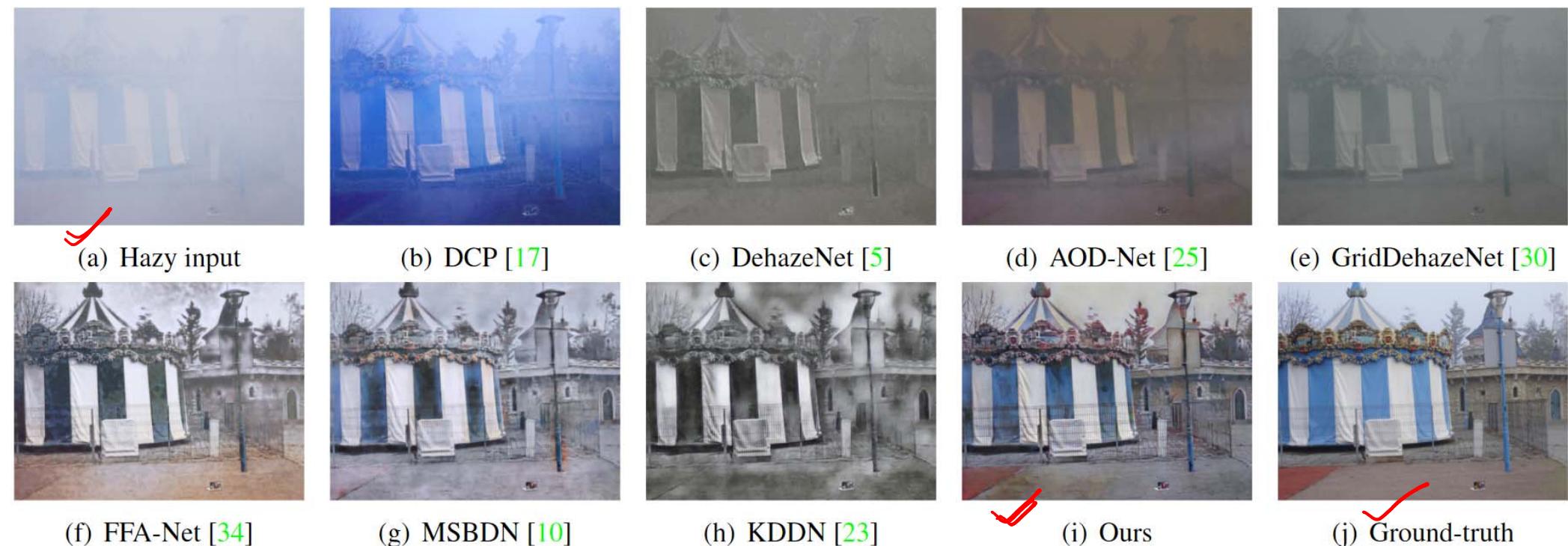


Figure 6. Qualitative results of DM-Count [34], CSRNet [15], CAN [17], and our STNNet on DroneCrowd. Best view in color version.

STNNet : Wen, Longyin, et al. "Detection, Tracking, and Counting Meets Drones in Crowds: A Benchmark." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021.

# CV Tasks : Image Dehazing

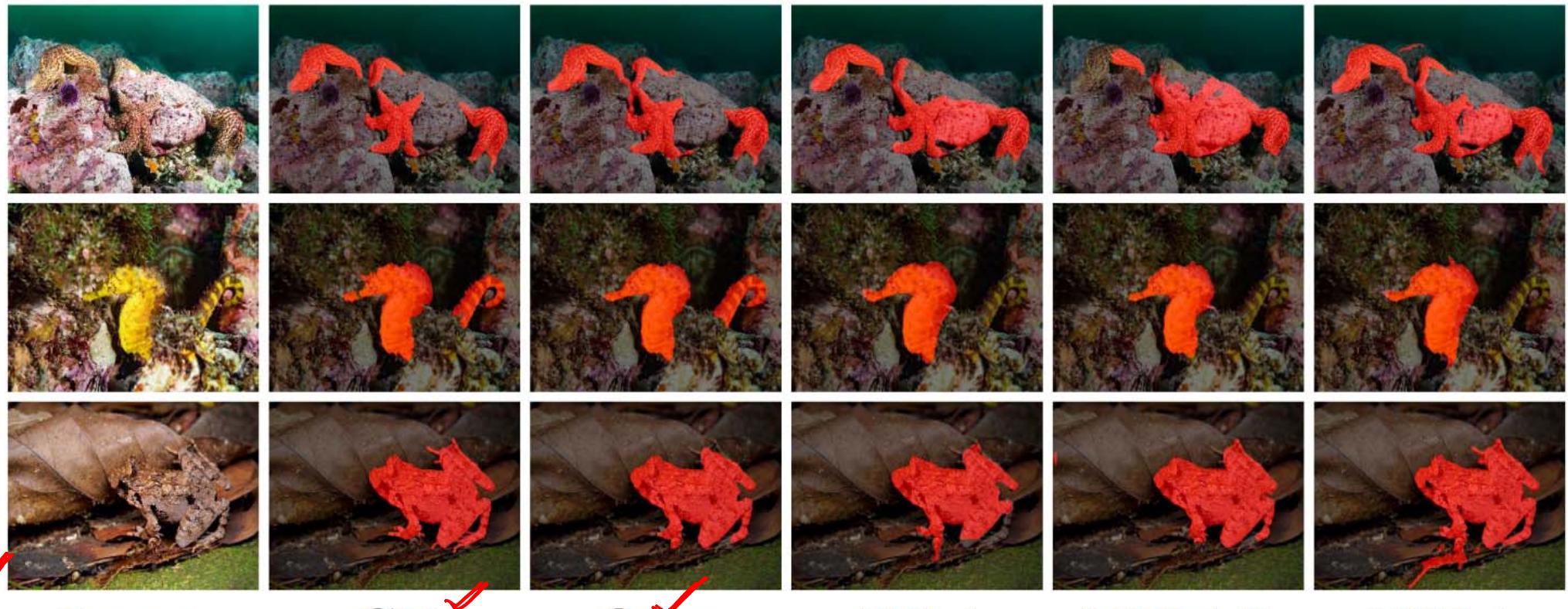
✓**Problem -11 :** Given a haze image I, How can we remove haze in order to get clear sharpened image.



✓  
Wu, Haiyan, et al. "Contrastive Learning for Compact Single Image Dehazing." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021. (Code available)

# CV Tasks : Camouflaged Object Segmentation

**Problem -12 :** What if the object color is similar to the background? How do we come up with methods which can solve segmentation problem in underline environment.



**Image**

**GT**

**Ours**

SINet

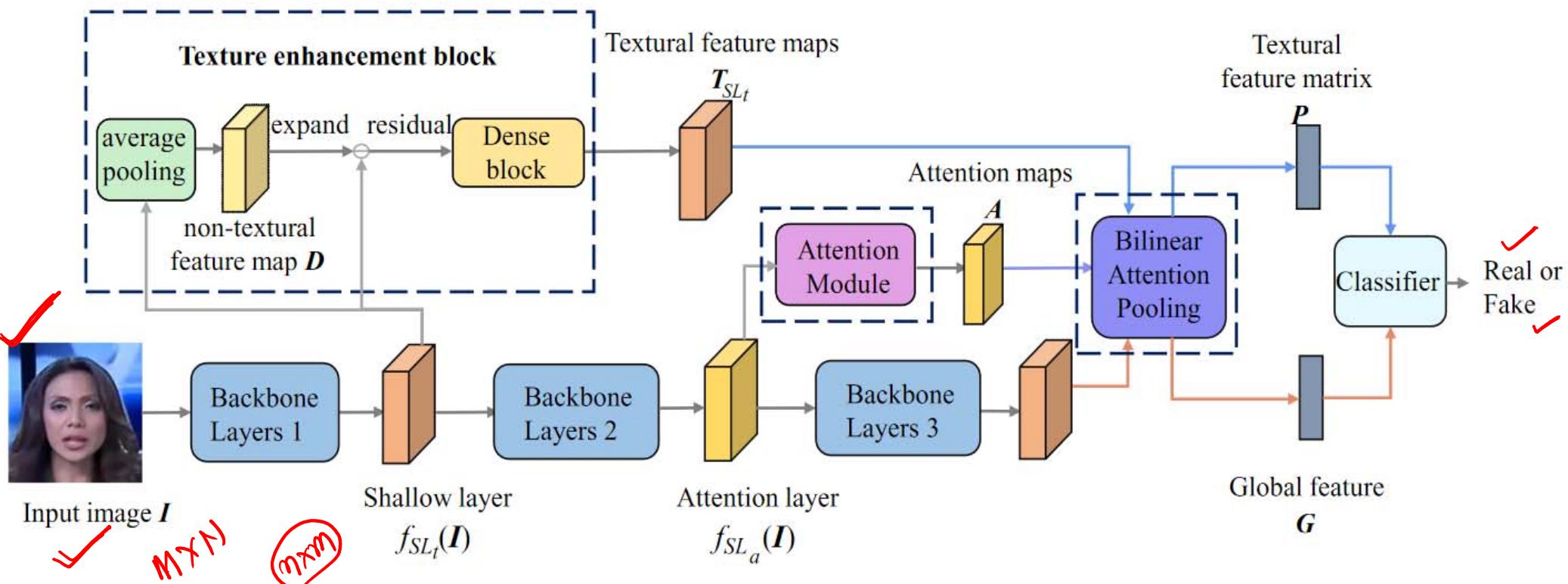
MINet-R

F3Net

Mei, Haiyang, et al. "Camouflaged object segmentation with distraction mining." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021. (Code available)

# CV Tasks : Image Forgery

Problem -13 : How to detect forgery (global and/or local) encountered in a given image.



Zhao, Hanqing, et al. "Multi-attentional deepfake detection." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021. Code available

# CV Tasks : Image Forgery

**Problem -13 :** How to detect forgery (global and/or local) encountered in a given image.



DeepFakes

Face2Face

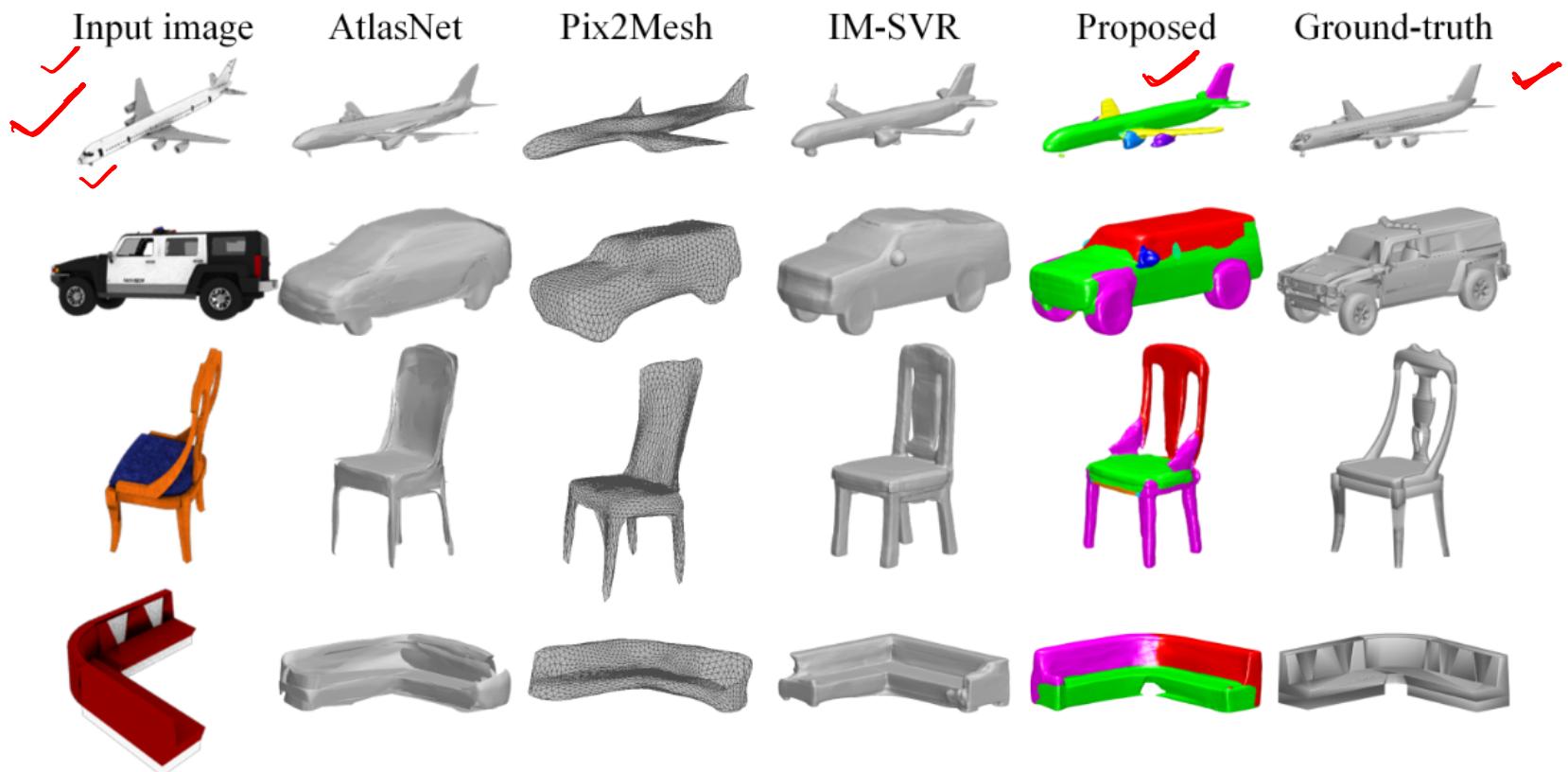
StyleGAN

PGGAN

Wang, Chengrui, and Weihong Deng. "Representative Forgery Mining for Fake Face Detection." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021. (<https://github.com/crywang/RFM.>)

# CV Tasks : 3D-Reconstruction from 2D-Images

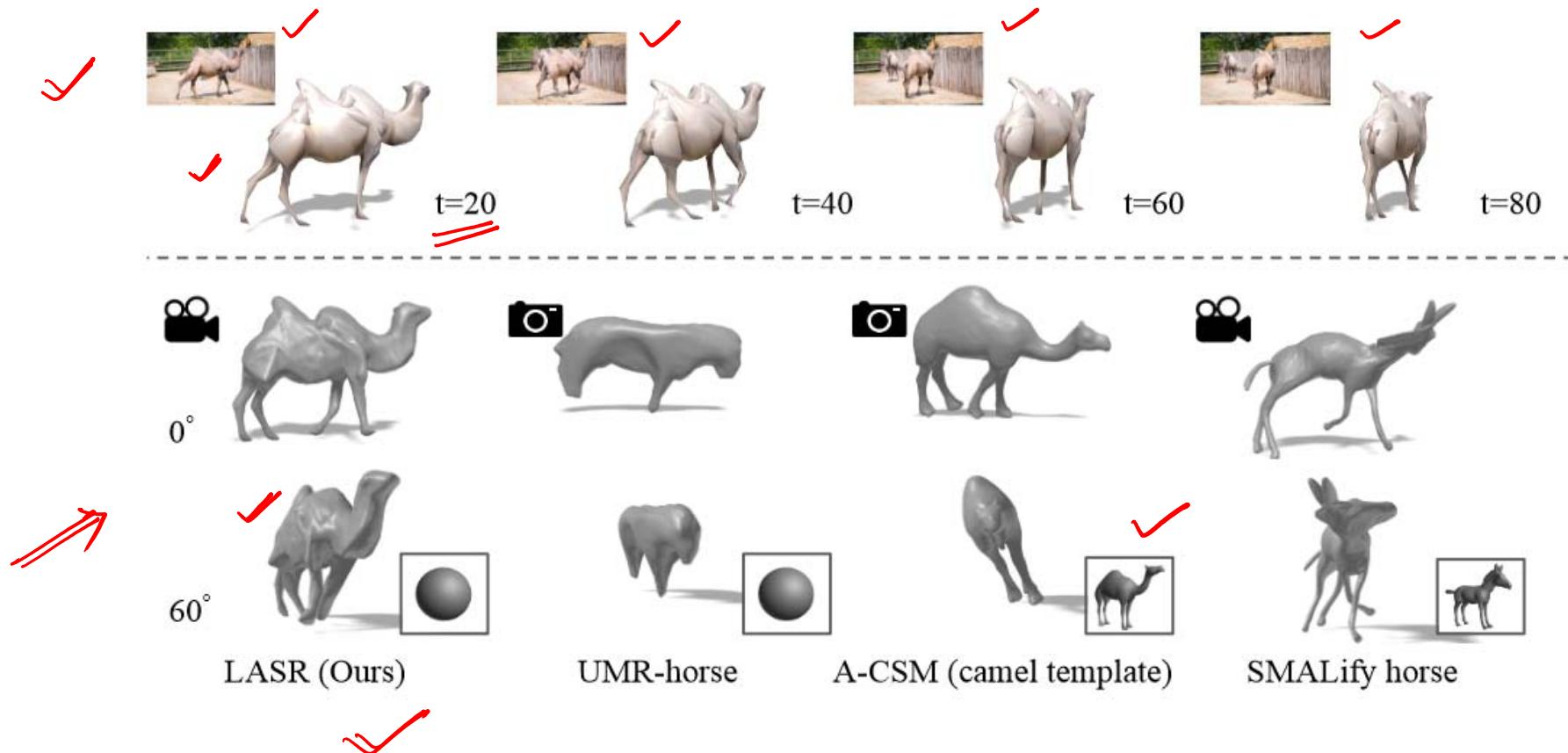
**✓ Problem -14 :** 3D structure of objects observed from a single view is a fundamental computer vision problem



**✓** Liu, Feng, Luan Tran, and Xiaoming Liu. "Fully Understanding Generic Objects: Modeling, Segmentation, and Reconstruction." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021. Code is available at <http://cvlab.cse.msu.edu/project-fully3dobject.html>.

## CV Tasks : Articulated 3D-shape Reconstruction

**Problem -15** : Articulated 3D-shape reconstruction from video as input.



Yang, Gengshan, et al. "LASR: Learning Articulated Shape Reconstruction from a Monocular Video." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021. Code is available at [lasr-google.github.io](https://lasr-google.github.io)

# CV Tasks : Self-Driving Car



**Problem -16** : How does the Computer Vision algorithms help in self-driving cars?



Waymo

[Tesla's Model S Autopilot is Amazing! - YouTube](#)

# Some of the Major Challenges in CV Tasks

1. **Viewpoint variation** : same object appears differently



View-1



View-2



View-3



2. **Illumination** – poor illumination is one of the major challenges in CV



# Some of the Major Challenges in CV Tasks

## 3. Intra-class variability:



Car



Car



Car

## 4. Background clutter and Occlusion



CV Tasks

# Conferences and Journals in Computer Vision

## Conferences:

- Computer Vision and Pattern Recognition (CVPR)
- International Conference on Computer Vision (ICCV)
- European Conference on Computer Vision (ECCV)
- International Conference on Pattern Recognition (ICPR)
- Asian Conference on Computer Vision (ACCV)
- International Conference on Image Processing (ICIP)

## Journals :

- IEEE Tran. on Pattern Analysis and Machine Intelligence (PAMI)
- IEEE Transactions on Image Processing
- International Journal of Computer Vision ←
- Computer Vision and Image Understanding ←
- ...so on

# Evaluation Pattern

✓ Major Exam ----- 25 ✓

✗ Assignments / Quiz ----- 25

✓ Project work ----- 50

✓ Work done – 20/25

✓ Report ----- 10

76  
≤ 6

✓ Paper (Journal/Conference) – 20/15

✓ Group formation / problem selection : 5<sup>th</sup> Oct. 2021

✓ Final submission date : 20<sup>th</sup> Nov. 2021

(5 students)

01 - 05 → 61  
06 - 10 → 62

Total = 100

# Pre-requisite of this course

- Fundamental of image processing
- Machine learning
- Basic probability and linear algebra
- Python

# Reference

- ✗ Richard Szeliski, Computer Vision: Algorithms and Applications,  
Springer, 2010 (online draft),
- ✗ Mubarak Shah, “Fundamentals of Computer Vision” (Online available)
- ✓ Ian Goodfellow, Yoshua Bengio and Aaron Courville, “Deep Learning”  
(Online available)

# Acknowledgement!



Sources for this lecture include materials from works by Szeliski, Abhijit Mahalanobis, Sedat Ozer, Ulas Bagci, Mubarak Shah, Antonio Torralba, and others.  
References are given for the source image contents.

# Queries!

Prof. Fei Fei Li