

Subject Code:	21AI639
Subject Name:	Computer Vision
Roll No.:	Cb.en.p2cse21012
Name:	Raahul Varman
Branch:	M.Tech CSE

**Course Instructor : Dr. SENTHILKUMAR .T** , Associate Professor, Dept. of CSE  
*in partial fulfillment of the requirements for the COURSE – 21AI639- Computer Vision*

**MASTER OF TECHNOLOGY**

**IN**

**COMPUTER SCIENCE AND ENGINEERING**

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

**AMRITA SCHOOL OF ENGINEERING , AMRITA VISHWA VIDYAPEETHAM**

**COIMBATORE – 641112**

**MAY 2022**

# INDEX:

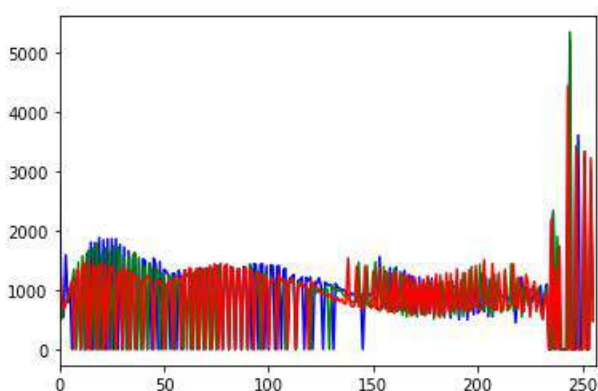
S.No	Lab Program	Page No.
1.	Basic Image Processing.	3
2.	Feature Detection & Matching.	9
3.	Data population.	20
4.	Optical flow Algorithm and Camshift & Meanshift.	22
4 b.	FACE DETECTION.	35
3 b.	Data Population – Customized.	36
5.	Data Analysis and Modeling.	38
6.	Deep Learning Optical Flow – RAFT.	41
7.	Action recognition with key frame extraction.	42
8.	Computer vision based alignment and tracking .	43



**Inference:**

Contrast refers to the brightness difference between various objects or parts of an image, whereas brightness refers to the overall lightness or darkness of an image.

Adding a positive constant to all of the image pixel values makes the image brighter.

**Question 3: Histogram Equalization****Output:****Inference:**

In image processing, a histogram is a crucial tool. It's a graphical representation of how data is distributed. An image histogram is a graphical depiction of a digital image's pixel intensity distribution. The x-axis represents the variable's possible range of values.



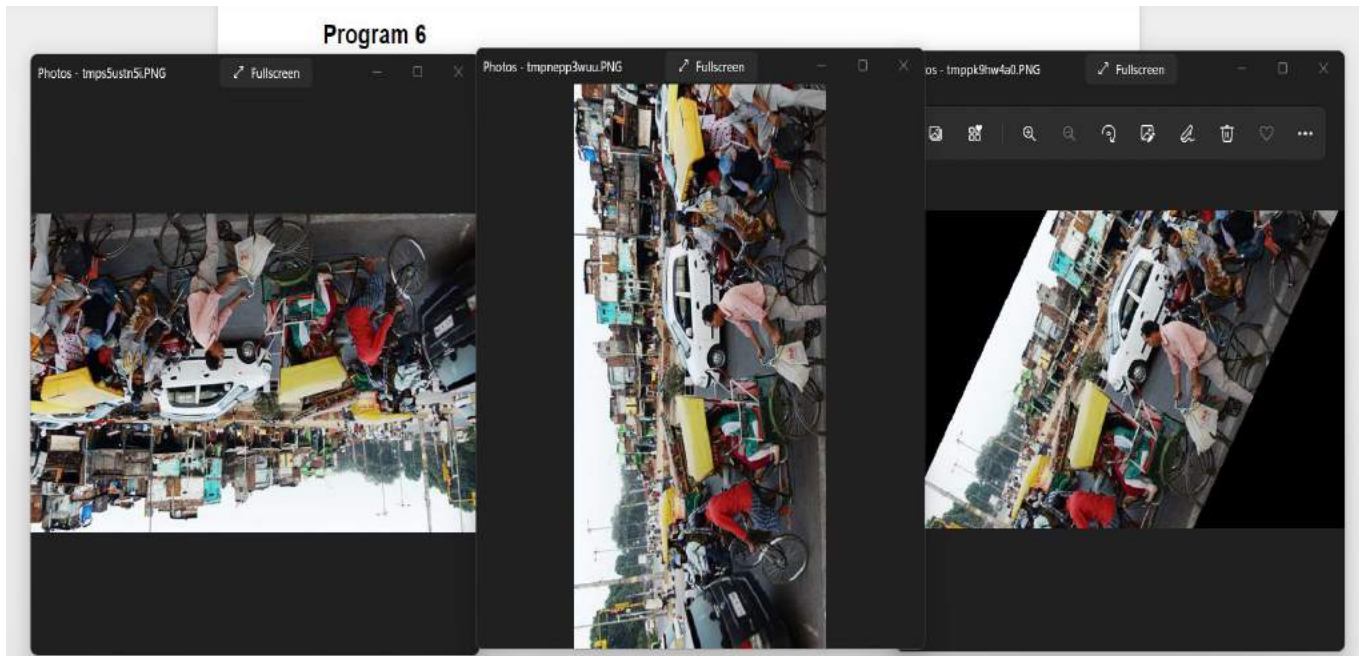
**Question 4: Averaging filter****Output:****Inference:**

Average filtering is a technique for smoothing photographs by lowering the intensity fluctuation between adjacent pixels. The average filter replaces each value with the average value of neighbouring pixels, including itself, as it moves through the image pixel by pixel.

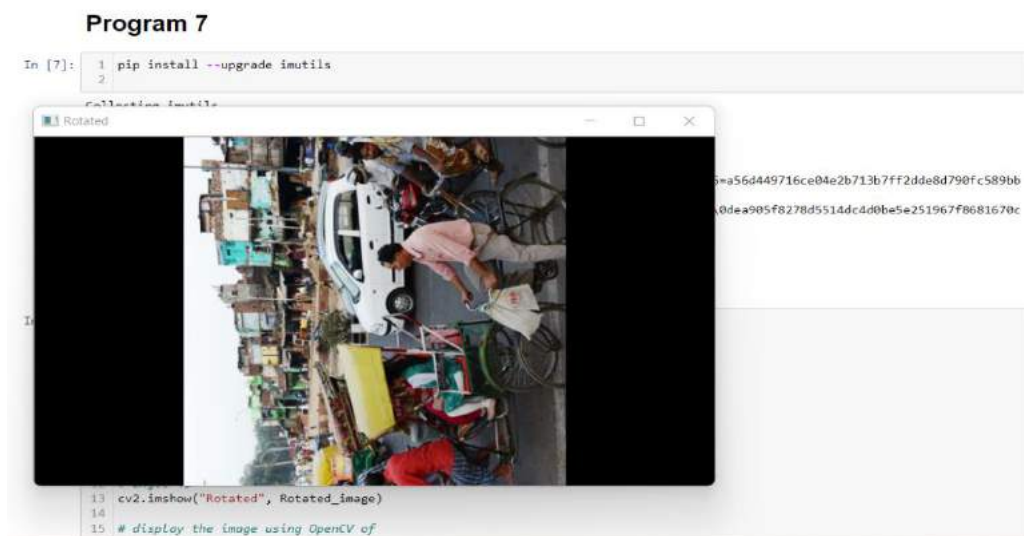
**Question 5: Median Filter****Output:**

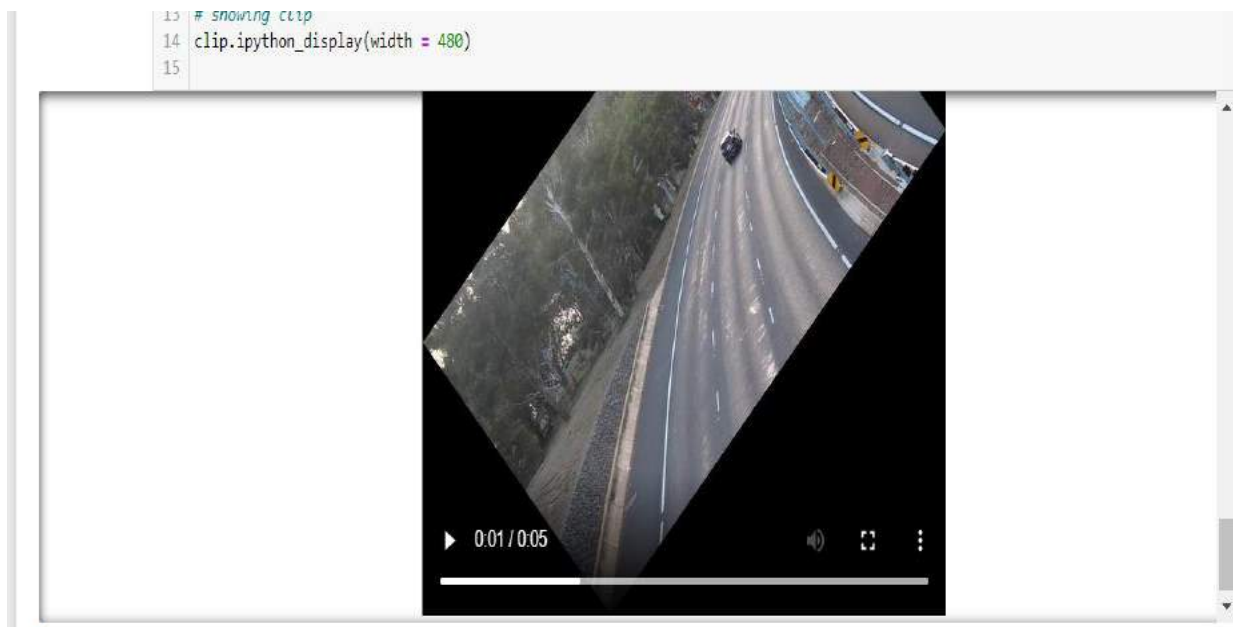
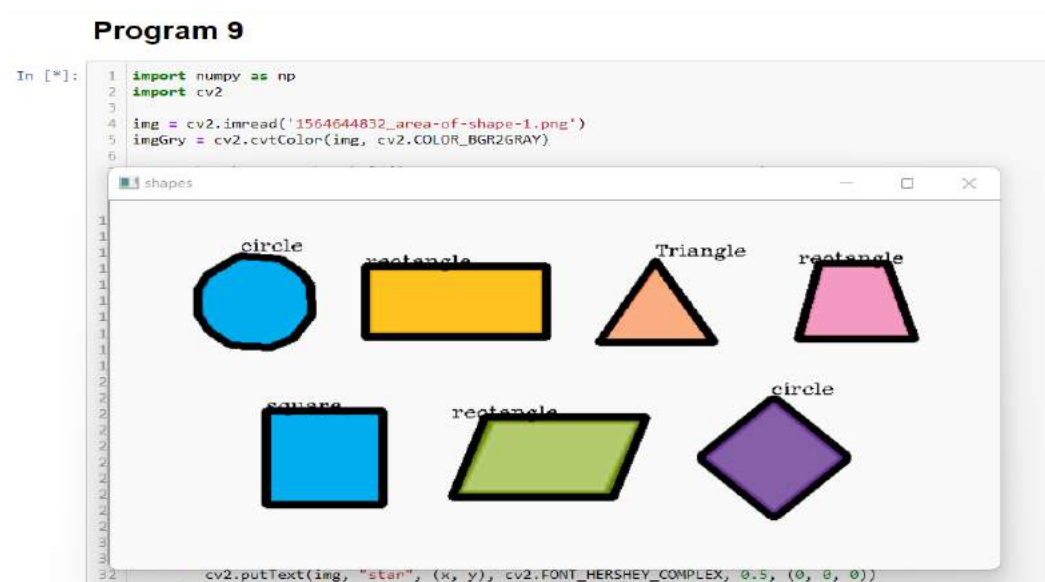
**Inference:**

The median filter is a non-linear digital filter that is commonly used to reduce noise from images.

**Question 6: Image rotation****Output:****Inference:**

We rotate the image.

**Question 7: image rotation 2****Output:**

**Question 8: Video Rotation****Output:****Inference:** We rotate the video.**Question 9: Shape Detection****Output:****Inference:**

To detect shapes in an image, use OpenCV's `findContours()` and `approxPolyDP()` functions. The OpenCV functions `findContours()` and `approxPolyDP()` can be used to discover forms in an image. We can recognize shapes based on how many corners they have.



**Question 10: Erosion and Dilation****Output:****Inference:**

In an image, dilation adds pixels to object boundaries, whereas erosion removes pixels from object boundaries. The size and shape of the structuring element used to process the image determines the number of pixels added or removed from the objects in the image.

**Question 11: Inverting an image****Output:****Inference:**

Light areas are mapped to dark, and dark areas are mapped to light in this image processing technique.



- **Inference:**

Image processing is a method for performing certain operations on an image, in order to get an improved image or extract useful information.

This is a type of signal processing where the input is an image and the output can be an image or features/characteristics related to that image.

## ➤ LAB 2: Feature Extraction and Matching:

### **Aim:**

The aim of this evaluation is to learn the different feature extraction algorithms and feature matching algorithms and suggest which suitable approach can be applied for the data-set.

### **Data-set Description:**

#### **Image Sample 1:**



#### **Image Sample 2:**



**Identify the following:** [Features required to be detected & Features that need to be matched]

1. Corners.
2. Scaled Image.
3. Rotational Image.
4. Affine transformation.

- **Feature Extraction Algorithm Name:**

1. **Harris Corner Detection:**

**Working principle:** Harris Corner Detector is a corner detection operator that is commonly used in computer vision algorithms to extract corners and infer features of an image.

- **Adv of Harris corner detection algorithm:**

Commonly, Harris corner detector algorithm can be divided into five steps.

- Color to grayscale.
- Spatial derivative calculation.
- Structure tensor setup.
- Harris response calculation.
- Non-maximum suppression.

- **Limitation:**

The big drawback of Harris corner detector is the need to set a different threshold for each image in order to detect the most important interesting points

**Output:**

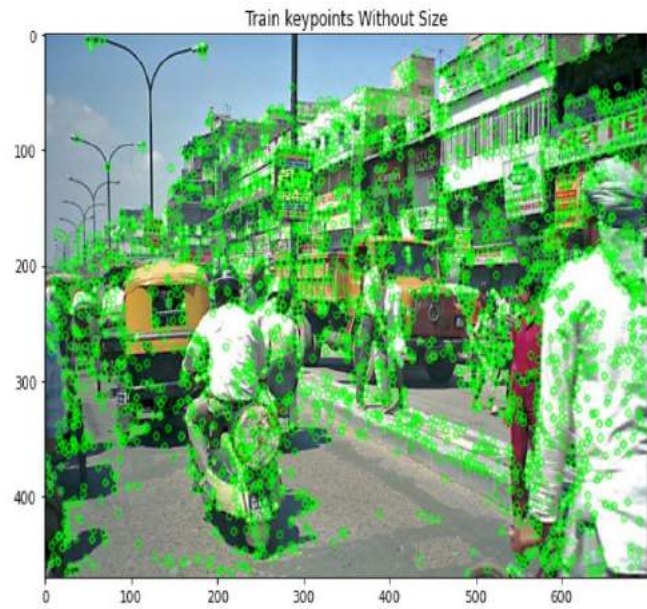
```
Out[3]: <matplotlib.image.AxesImage at 0x1de792f8610>
```



## 2. Scale-Invariant Feature Transform (SIFT):

### Output:

Number of Keypoints Detected In The Training Image: 3865  
 Number of Keypoints Detected In The Query Image: 389



## 3. Oriented FAST and Rotated BRIEF (ORB):

**Working principle:** ORB performs as well as SIFT on the task of feature detection (and is better than SURF) while being almost two orders of magnitude faster. ORB builds on the well-known FAST key-point detector and the BRIEF descriptor. Both these techniques are attractive because of their good performance and low cost.

### ▪ **ORB's Advantages:**

- The addition of a fast and accurate orientation component to FAST.
- The efficient computation of oriented BRIEF features.
- Analysis of variance and correlation of oriented BRIEF features.
- A learning method for decor-relating BRIEF features under rotational in variance, leading to better performance in nearest-neighbor applications.



## Output:

Number of Keypoints Detected In The Training Image: 500

Number of Keypoints Detected In The Query Image: 234



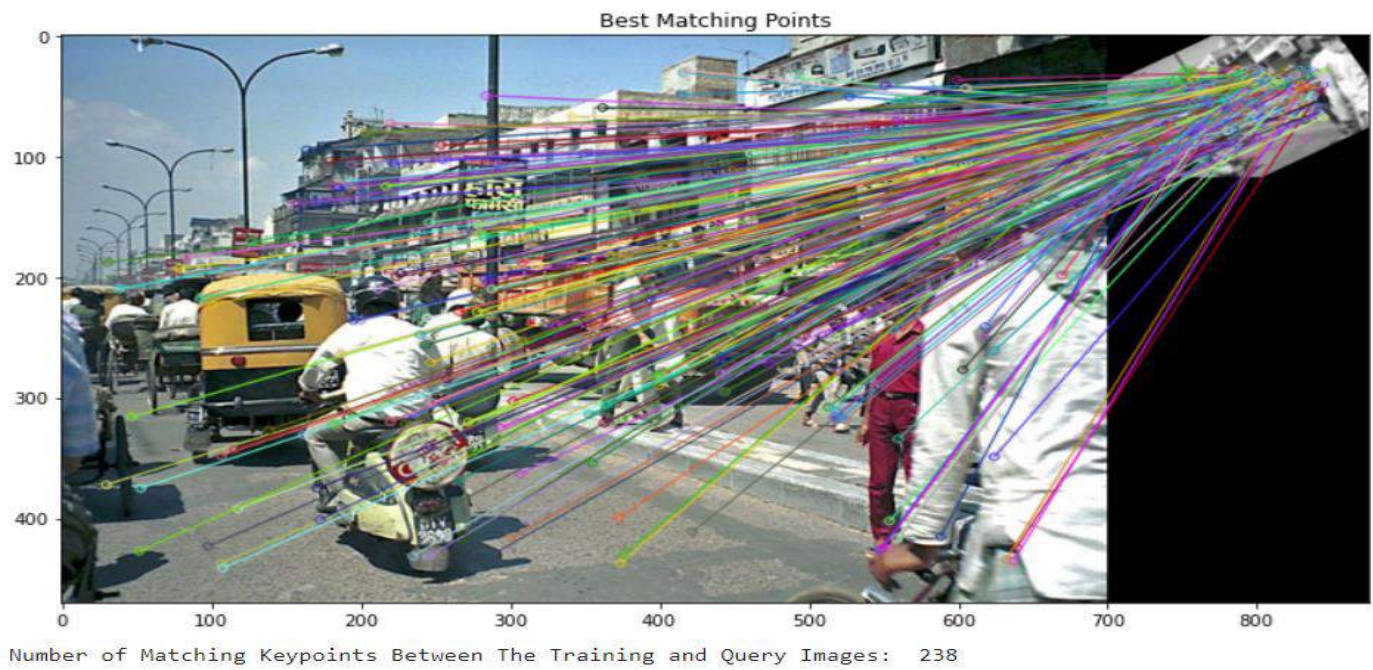
- **Feature Descriptor Algorithm Name:**

1. **BRIEF(Binary Robust Independent Elementary Features):**

**Working Principle:** Brief takes all key-points found by the fast algorithm and convert it into a binary feature vector so that together they can represent an object. Binary features vector also know as binary feature descriptor is a feature vector that only contains 1 and 0. In brief, each key-point is described by a feature vector which is 128–512 bits string.

**Advantages of BRIEF:**

Brief relies on a relatively small number of intensity difference tests to represent an image patch as a binary string. Not only is construction and matching for this descriptor much faster than for other state-of-the-art ones, but it also tends to yield higher recognition rates, as long as invariance to large in-plane rotations is not a requirement.



## 2. Oriented FAST and Rotated BRIEF (ORB):



- **Metrics used for Feature Matching:**

- 1. Brute Force Matcher:**

we identify image features, distinctive points in our images. Second, we associate a descriptor for each feature from its neighborhood. Finally, we use descriptors to match features across two or more images. Afterwards, we can use the matched features for a variety of applications including state estimation, visual odometry, and object detection.

The simplest solution to the matching problem is referred to as “brute force feature matching”, and is described as the following. First, define a distance function  $d$  that compares the descriptors of two features  $f_i$  and  $f_j$ , and defines the distance between them.

The more similar the two descriptors are to each other, the smaller the distance between them. Second, for every feature  $f_i$  in image one, we apply the distance function  $d$  to compute the distance with every feature  $f_j$  in image two. Finally, we will return the feature which we'll call  $f_c$  from image two with the minimum distance to the feature  $f_i$  in image one as our match.

This feature is known as the nearest neighbor, and it is the closest feature to the original one in the descriptor space. The most common “distance function” used to compare descriptors is the “sum of squared distances or SSD”.

- **Sum of Squared Differences (SSD):**

$$d(f_i, f_j) = \sum_{k=1}^D (f_{i,k} - f_{j,k})^2$$

**Limitation:**

Brute force matching is suitable when the number of features we want to match is reasonable, but has quadratic computational complexity making it ill-suited as the number of features increases.



## ■ Inference:

Edges and interest points in images convey a wealth of information about the image's content. They correspond to picture local regions and are crucial in many image analysis applications such as recognition, matching, reconstruction, and so on.

Feature extraction aids in the reduction of unnecessary data in a data set. Finally, reducing the data makes it easier to build the model with less machine effort, as well as speeding up the learning and generalization processes in the machine learning process.

Using a search distance method, feature matching finds equivalent features from two similar photos. The feature matching technique is used to either find or deduce and transfer properties from the source to the target image, with one image serving as the source and the other as the target.

## ● INTERACTIVE UI Feature Detection and Matching:

**Note:** I did this entire setup and execution in kali Linux.

1. There are two python files: one python files is for creating/ adding feature algorithms and one more for user interface setup for using featuring algorithms.

> Features.py

> FeaturesUIdisplay.py

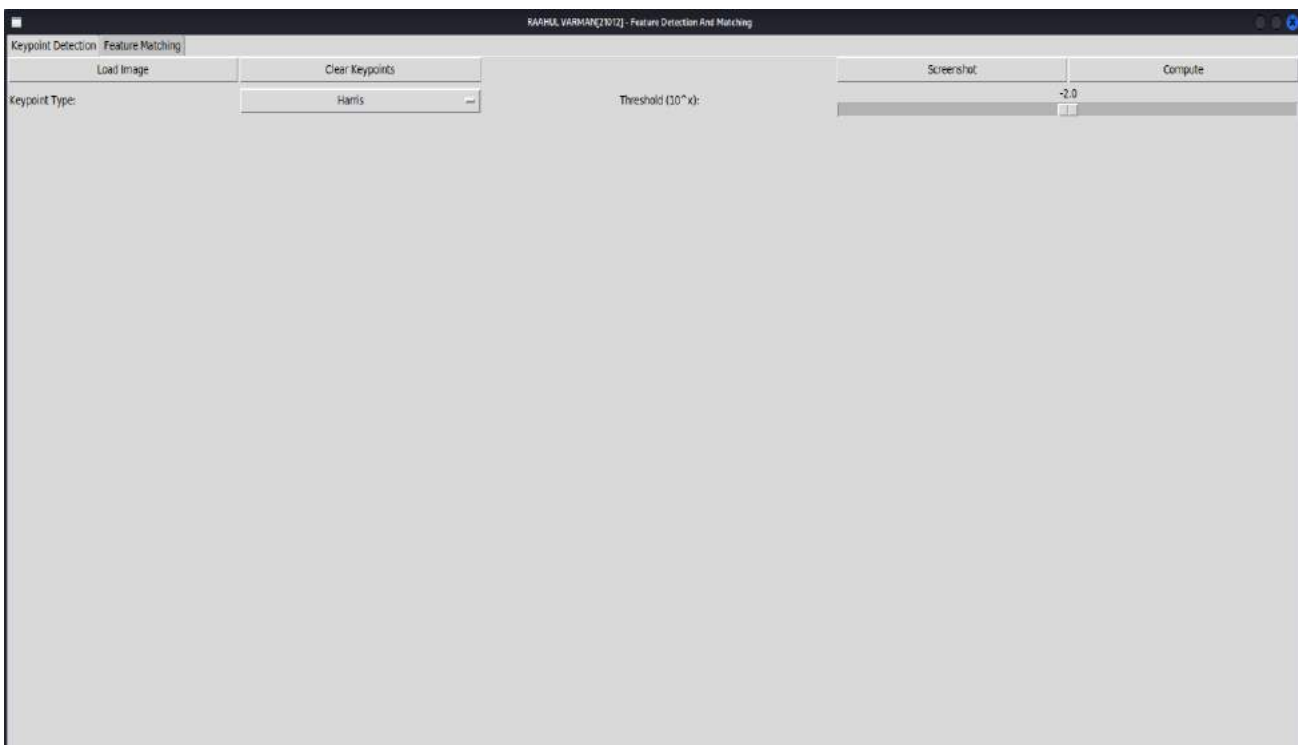
2. Using command prompt we set the path to python file location to call the python file for execution.
3. We use -> python <filename.py>.
4. In our case we call -> python FeaturesUIdisplay.py and press Enter.



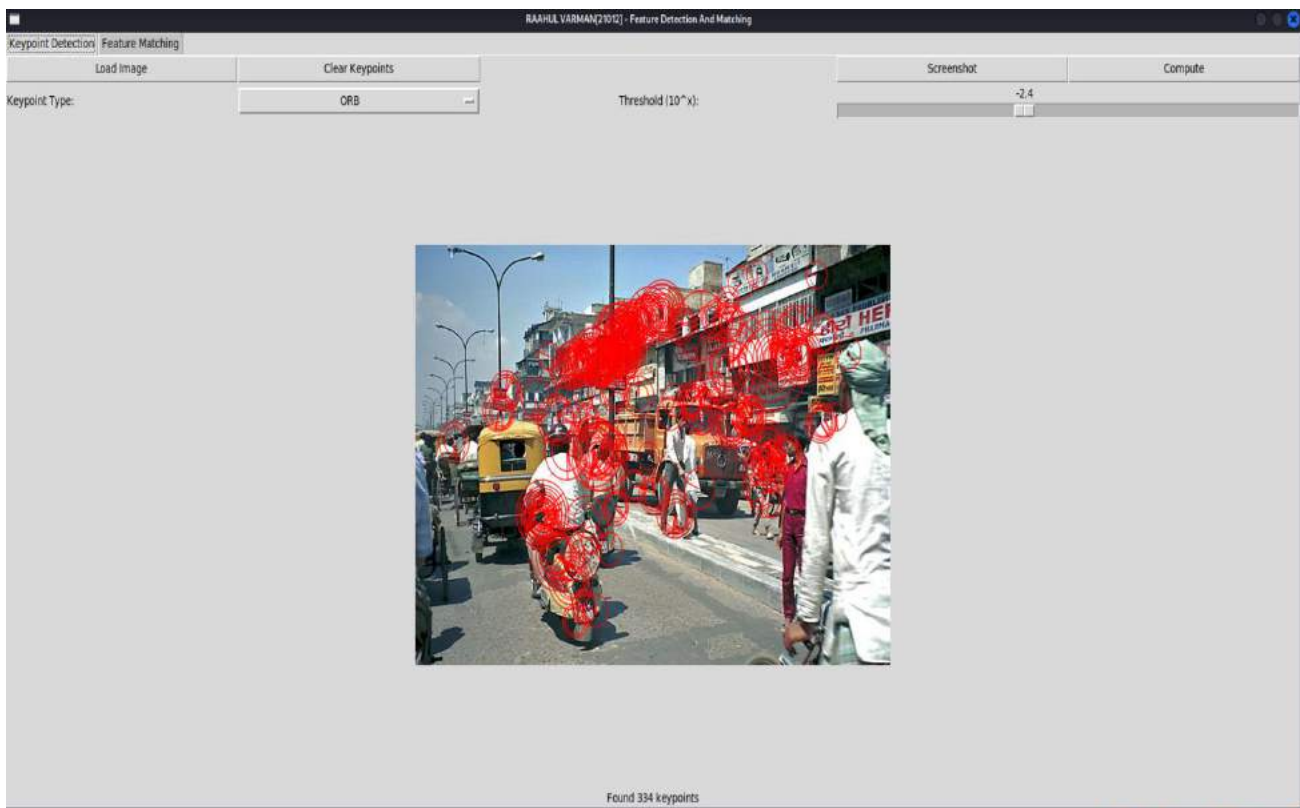
```
kali@kali:~/Desktop/Features_Detection_and_Matching-master/src
File Actions Edit View Help

(kali@kali)~/Desktop/Features_Detection_and_Matching-master/src
$ python featuresUIDisplay.py
/home/kali/Desktop/Features_Detection_and_Matching-master/src/featuresUIDisplay.py:3: DeprecationWarning: the imp module is deprecated in favour of importlib; see the module's documentation for alternative uses
from imp import reload
```

## 5. Output will be loaded.



## Key Feature Detection from my sample image:

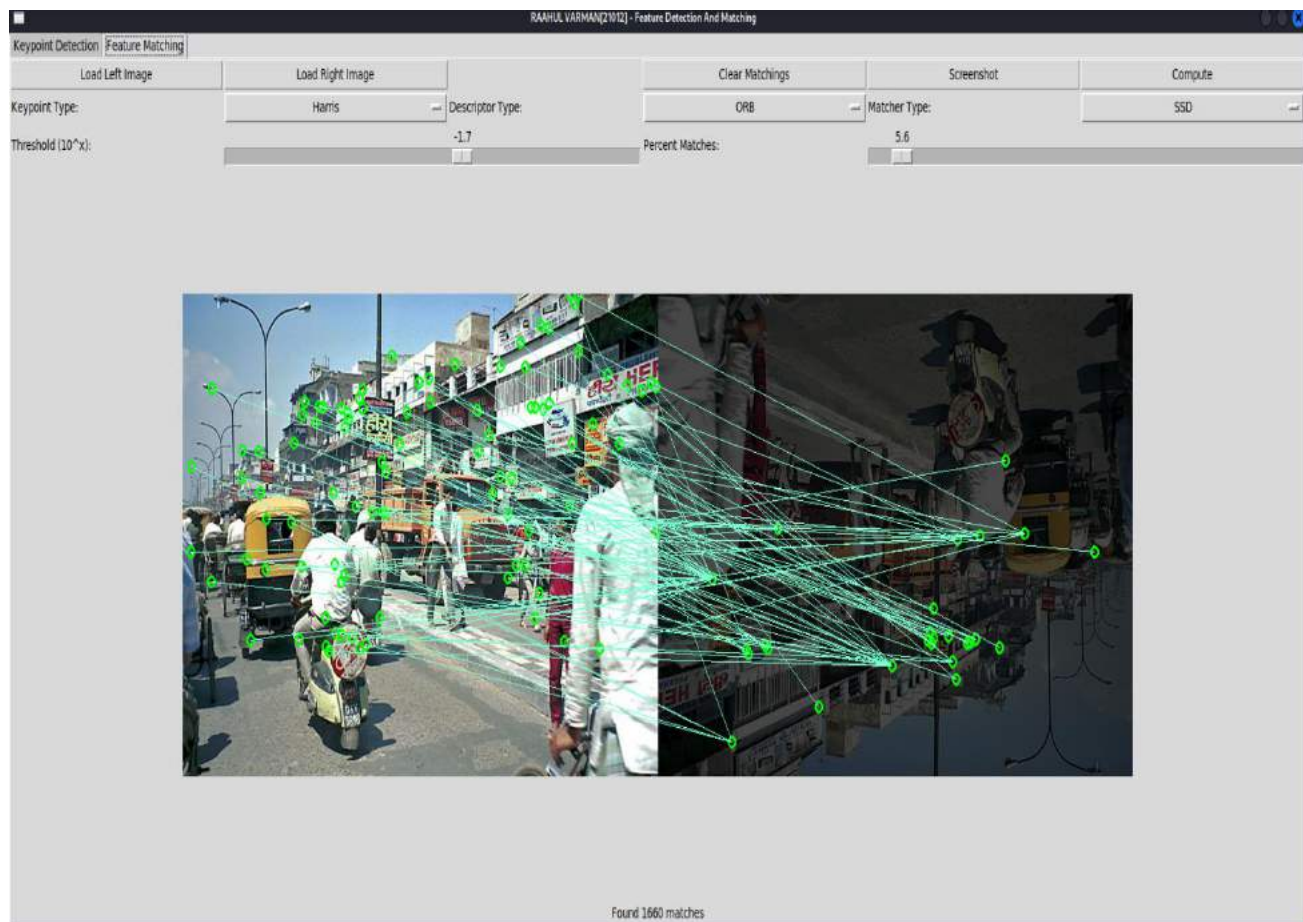


## Feature Matching:

**Sample image:** One image is original and second image is rotated and added some filter (noise).



## Output:



## ➤ Lab 03 : DATA POPULATION:

- **Title:** “Model For Assisting Blind People By Detecting The Surrounding Objects/People/Activities With Voice Description.”

- **Abstract:**

There have been several systems designed to support visually-impaired people and to improve the quality of their lives. Unfortunately, most of these systems are limited in their capabilities. We decided to build a model which detects all possible surrounding things and guide them through voice description. Our main objective is to detect the obstacles in the person's pathway and warn them using voice description.

### Output:

		sharpness	contrast	mean_sat	std_sat	min_sat	max_sat	mean_bright	std_bright	min_bright	max_bright	colorfulness	hue_hist_1	hue_hist_2	hue_hist_3	shannon_entropy
2	0	1477.053	0.263872	102.68141	114.6114	255	0	108.9957952	102.864254	255	0	3.064395899	14	147	46.7053067	8.888408288
3	1	4543.997	0.259352	76.52696	56.62634	255	0	121.5755413	73.9336103	255	0	34.51540471	18	84	46.4239517	13.54138463
4	2	1740.611	0.331498	63.754197	72.50882	255	0	88.46784044	70.0220229	255	0	27.15711482	20	97	59.3380661	12.52376452
5	3	2769.083	0.225151	148.447594	75.33106	255	0	88.77501302	50.7272271	255	0	51.92119559	16	217	40.3020655	13.15744944
6	4	3121.815	0.170011	83.5869466	55.83494	255	0	109.3293522	58.807363	255	0	39.18421205	14	135	30.4319237	12.30113589
7	5	1124.724	0.219065	71.7033189	64.18686	255	0	148.6929084	51.0884041	255	0	51.2284675	14	84	39.2125942	12.18934376
8	6	189.3224	0.301611	36.784528	33.12133	228	0	119.732819	73.2426575	255	22	21.7024955	15	88	53.9883826	10.75647046
9	7	703.414	0.309182	103.23351	73.71014	255	0	96.65379041	64.0930649	255	0	71.31867585	20	10	55.343617	12.31874489
10	8	3127.229	0.280199	46.5957454	53.59137	255	0	120.1293099	82.1876149	255	0	29.60098717	12	545	50.1556087	12.30173162
11	9	3395.363	0.148805	95.2081711	45.46012	255	0	97.271381	49.6309649	255	0	37.78645975	17	48	26.6360416	13.3826206
12	10	5264.446	0.28279	72.0500928	60.18805	255	0	100.4784204	64.934584	255	0	37.94815978	20	179	50.6193212	13.91473623
13	11	2289.165	0.326508	77.2764225	61.2455	255	0	93.06748698	69.8960605	255	0	33.93625138	20	300	58.4449094	13.20427646
14	12	386.9361	0.437889	206.779769	69.73481	255	0	28.06677409	59.6336499	255	0	81.20001489	11	458	78.3820949	6.48821717
15	13	4327.952	0.216743	173.421423	62.11964	255	0	85.24505706	62.5925097	255	0	53.96807458	20	30	38.7969305	13.421823
16	14	2424.411	0.272996	71.3469596	39.32545	255	0	139.3976544	50.176872	255	0	43.94417881	7	137	48.8662689	12.46852461
17	15	2208.279	0.319879	41.5486719	33.90807	255	0	154.4046224	68.7067761	255	2	23.41485499	20	450	57.2583126	13.06359323
18	16	3520.001	0.196402	88.9670302	65.37251	255	0	162.0028835	65.203602	255	0	61.72770794	20	54	35.1560414	13.09650071
19	17	10275.09	0.173935	105.891663	54.41864	255	0	133.4038574	72.8623815	255	0	48.92667297	20	165	31.1344259	14.51299042
20	18	4774.396	0.166907	188.781817	73.71669	255	0	42.16796196	64.6732957	255	0	46.11209813	20	43	29.8762723	11.07037375
21	19	7404.018	0.306592	62.35179	62.28873	255	0	89.66992467	59.218121	255	0	33.68074231	20	108	54.8799577	13.67333008
22	20	2763.642	0.270204	83.5134433	60.85452	255	0	106.4424825	66.3760119	255	0	59.76832494	20	52	48.3664732	13.05296298
23	21	1231.514	0.321021	31.9909819	50.78078	255	0	173.5219036	66.4348485	255	0	25.93276655	12	350	57.4627434	10.16870035
24	22	3712.637	0.117291	112.472061	52.75428	255	0	97.17611979	54.9159267	255	0	57.14263084	18	59	20.9951626	12.40858599
25	23	5228.616	0.249572	46.9773387	43.05727	255	0	129.5107627	62.8979166	255	0	25.24885236	6	302	44.673394	12.95778131
26	24	2063.501	0.204231	19.805168	24.36841	212	0	170.2842187	63.4532253	255	22	16.50098474	10	421	36.5573314	12.13066847
27	25	5822.647	0.290135	58.6510968	48.84845	255	0	128.5952315	60.2586314	255	0	41.7120145	19	79	51.9341521	13.66920677
28	26	16637.04	0.256034	87.2983735	72.47033	255	0	62.39764458	75.973646	255	0	22.42663286	20	41	45.8301732	12.72732377
29	27	952.1554	0.23965	73.0964213	51.40065	255	0	103.9513867	54.2206188	255	0	37.07146248	18	125	42.8974141	11.93527667
30	28	4354.452	0.286668	44.9596029	53.44046	255	0	114.8688542	69.1182519	255	0	42.59907131	20	412	51.3135862	12.29873017
31	29	3273.591	0.229014	56.6622546	45.38545	255	0	104.1372137	61.2149407	255	0	21.40022073	3	541	40.993588	12.43377402

Feature Name	Purpose	Data Type of the feature
Sharpness	Sharpness is a combination of two factors: resolution and acutance. Resolution is straightforward and not subjective. It's just the size, in pixels, of the image file. All other factors equal, the higher the resolution of the image—the more pixels it has—the sharper it can be. Acutance is a little more complicated. It's a subjective measure of the contrast at an edge. There's no unit for acutance—you either think an edge has contrast or think it doesn't. Edges that have more contrast appear to have a more defined edge to the human visual system.	Float64
Contrast	The term contrast refers to the amount of color or grayscale differentiation that exists between various image features in both analog and digital images. Images having a higher contrast level generally display a greater degree of color or grayscale variation than those of lower contrast.  Contrast enhancement processes adjust the relative brightness and darkness of objects in the scene to improve their visibility.	Float64



Saturation	The term saturation describes the depth or intensity of colour present within an image.	Float64
Brightness	The term saturation describes the depth or intensity of colour present within an image.	Float64
Colorfulness	Colorfulness is the "attribute of a visual perception according to which the perceived color of an area appears to be more or less chromatic. The colorfulness evoked by an object depends not only on its spectral <u>reflectance</u> but also on the strength of the illumination, and increases with the latter unless the <u>brightness</u> is very high	Float64
Hue_Histogram	An image histogram is a type of histogram that acts as a graphical representation of the tonal distribution in a digital image. It plots the number of pixels for each tonal value. By looking at the histogram for a specific image a viewer will be able to judge the entire tonal distribution at a glance.	Float64/int64
Shannon entropy	In Image, Entropy is defined as corresponding states of intensity level which individual pixels can adapt. It is used in the quantitative analysis and evaluation image details, the entropy value is used as it provides better comparison of the image details.  Shannon entropy as a measure of image information is extensively used in image processing applications.	Float64

## Lab 04: DEEP LEARNING BASED ON OPTICAL FLOW

### OBJECTIVE:

Our aim is to estimate the motion of the image intensities over time in a video using optical-flow algorithms like Lucas-Kanade, Horn-schunck and Dense optical flow. i.e. to calculate velocity of points within the images, and to provide the estimation of where the points could be in the next image sequence in our shopping video dataset.

### PERFORMANCE METRICS:

Metrics Name	Purpose	Formula	Expected Value
Point Rotational Error (PRE)	To prevent deviation (error) of a objects angular orientation from an expected, nominal, or commanded angle or displacement.	$PRE = \begin{cases} \cos^{-1} \left( \frac{u u_{GT} + v v_{GT}}{\sqrt{u^2 + v^2} \sqrt{u_{GT}^2 + v_{GT}^2}} \right), & \text{if } (u^2 + v^2) \neq 0 \wedge (u_{GT}^2 + v_{GT}^2) \neq 0 \\ \pi, & \text{if } (u^2 + v^2) \oplus (u_{GT}^2 + v_{GT}^2) = 0 \\ 0, & \text{if } (u^2 + v^2) = 0 \wedge (u_{GT}^2 + v_{GT}^2) = 0 \end{cases}$	Between 0.01 and 0.05
Angular Error (AE)	AE is very sensitive to small displacements, To determine the errors in angles.	$AE = \cos^{-1} \left( \frac{u u_{GT} + v v_{GT} + 1}{\sqrt{u^2 + v^2 + 1} \sqrt{u_{GT}^2 + v_{GT}^2 + 1}} \right)$	Between 1.00 – 2.00
End point Error (EPE)	It measures the distance between the endpoints of two optical flow vectors ( $u_0, v_0$ ) and ( $u_1, v_1$ )	$\text{sqrt}((u_0 - u_1)^2 + (v_0 - v_1)^2)$	Should have small/min value
Normalized Euclidean Error (NEE)	IT gives the squared distance between two vectors where there lengths have been scaled to have unit norm. This is helpful when the direction of the vector is meaningful but the magnitude is not.	$NEE = \begin{cases} \frac{\sqrt{(u - u_{GT})^2 + (v - v_{GT})^2}}{\min((u^2 + v^2), (u_{GT}^2 + v_{GT}^2))}, & \text{if } \min((u^2 + v^2), (u_{GT}^2 + v_{GT}^2)) > \epsilon \\ \frac{\sqrt{(u - u_{GT})^2 + (v - v_{GT})^2}}{\epsilon}, & \text{if } \min((u^2 + v^2), (u_{GT}^2 + v_{GT}^2)) = 0 \end{cases}$	Value should be between 0 to 1
Enhanced Normalized Euclidean Error (ENEE)	Another way to get over EPE drawbacks, is to calculate the relative distance between ~E and G~T vectors and to use different normalization methods.	$ENEE = \begin{cases} \frac{\sqrt{(\ P_{GT}\ )^2 + \tau(\ N_{GT}\ )^2}}{\min((u^2 + v^2), (u_{GT}^2 + v_{GT}^2))}, & \text{if } \min((u^2 + v^2), (u_{GT}^2 + v_{GT}^2)) > \epsilon \\ \frac{\sqrt{(\ P_{GT}\ )^2 + \tau(\ N_{GT}\ )^2}}{\epsilon}, & \text{if } \min((u^2 + v^2), (u_{GT}^2 + v_{GT}^2)) = 0 \end{cases}$	Value should be between 0 to 1
Linear Projection Error (LPE)	This metric is a kind of mixture between AE and EPE, where the magnitude difference between ( $u;v$ ) is added to the perpendicular distance between them.  We can the angular distance between the two non-null vectors based on the perpendicular distance between both vectors.	$LPE = \begin{cases} \ \vec{GT} - \vec{E}\  + \max(\ proj_{\vec{GT}} \vec{E}\ , \ proj_{\vec{E}} \vec{GT}\ ), & \text{if } \ \vec{GT} \cdot \vec{E}\  \neq 0 \\ \ \vec{GT} - \vec{E}\  + \max(\ \vec{GT}\ , \ \vec{E}\ ), & \text{if } \ \vec{GT} \cdot \vec{E}\  = 0 \end{cases}$	Should have small/min value

## Part-A :

1. Analyze the video for slow, fast and medium frame rate based videos.

Dataset Drive Link : <https://drive.google.com/drive/folders/1-1uNvgEhFRPvqYipvp6wQQhIjHLi-KFP?usp=sharing>

Dataset: Shopping Mall Video.


- Original Video Length is 13 seconds.
  - Fast Video Length is 3 seconds.
  - Medium video length is 7 seconds.
  - Slow video length is 27 seconds.
2. Try out Horn and Schunck, Lucas Kanade, Dense Optical Flow algorithm and present the results in a Table form.




### ➤ Lucas Kanade Algorithm:

#### ➤ Frame 1: original video



### Output:

Original video		Time taken for original video is 3mins 41secs
----------------	--	---

Fast video		Time taken for fast video is 1mins 13secs
Medium video		Time taken for medium video is 2mins 9secs
Slow video		Time taken for slow video is 3mins 23secs
Inference	<p>Lucas Kanade algorithm don't works well in fast video footage, compare to orginal,medium and slow video footage.</p> <p>Algorithm was able to detect the corners and was able to estimate the where the flow could be in next image sequence.</p>	

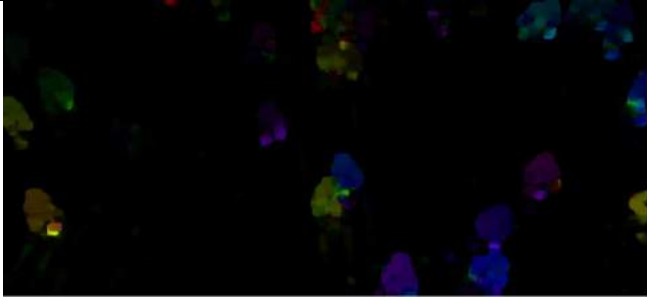
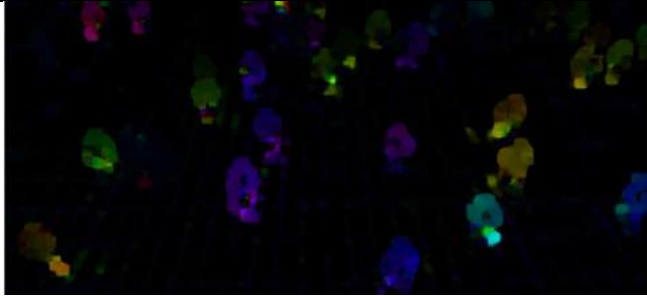
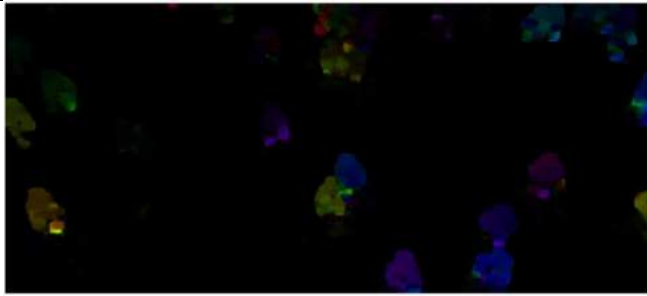
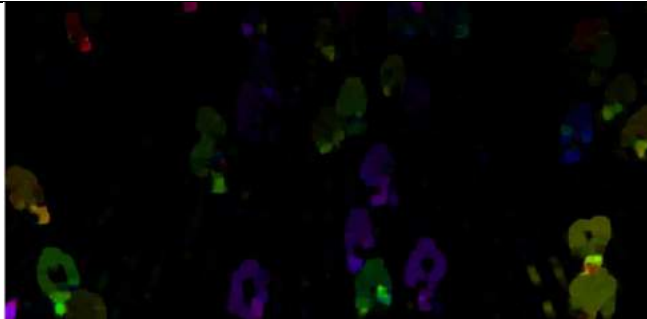
## ➤ Dense Optical Flow algorithm:

Frame 1: Original video





**Output:**




<b>Original video</b>		Time taken for original video is 3mins 53secs
<b>Fast video</b>		Time taken for fast video is 1mins 21secs
<b>Medium video</b>		Time taken for medium video is 2mins 12secs
<b>Slow video</b>		Time taken for slow video is 3mins 47secs
<b>Inference</b>	<p>Dense Optical Flow algorithm works well in fast video footage and Slow video footage, compare to original and medium video footage.</p> <p>Dense optical flow able to compute optical flow vector for each and every pixel of each frame.</p>	


## ➤ Horn and Schunck:

In horn and schunck we use Image(Each video frame) as the input to estimate the where the points could be in the next image sequence.

Frame 1: Original video



<b>Original video (Frame 137)</b>	
<b>Fast video (Frame 35)</b>	
<b>Medium video (Frame 71)</b>	

<p><b>Slow video (Frame 50)</b></p>	
<p><b>Inference</b></p>	<p>As the people in the video don't have constant velocity movement the estimation in all the footage is little poor, which result in error at the boundaries.</p>

### 3. Write a Discussion section and present the findings in a paragraph.

In Horn and Schunck, as we know it can only handle object/person/things which has constant velocity. In our dataset we don't have any constant movement kind of objects/person/things, which resulted in poor estimation.

In Lucas-Kanade, this algorithm works fine with slow motion footages/recordings. And it don't work for fast motion footages/recordings. We can clearly see in our output, where we got better estimation of flow in next image sequence in slow motion footage then fast motion footage. In fast motion footage, the estimation was mismatched and uneven.

In Dense Optical Flow, as it compute optical flow vector for each and every pixel of each frame. We where able to estimate the flow in all the footage clearly. Even though the computational time is much higher then other two algorithm. Its give accurate result and denser result which is suitable for application like motion and video segmentation.

From comparison, Dense optical flow was able to give better estimation to describe image motion.

## Part-B:

1. Introduce the noise in the video and complete the following:

### a. Apply Filtering in the video:

Dataset Drive Link:

[https://drive.google.com/drive/folders/1IL6MeTv1drSc\\_cgQ7WPesSsvg4EvJB2N?usp=sharing](https://drive.google.com/drive/folders/1IL6MeTv1drSc_cgQ7WPesSsvg4EvJB2N?usp=sharing)

Dataset: Shopping Mall Video.

- Original Noise Video Length is 13 seconds.
- Fast Noise Video Length is 3 seconds.
- Medium Noise video length is 7 seconds.
- Slow Noise video length is 27 seconds.

b. Try out Horn and Schunck, LucasKanade, Dense Optical Flow algorithm:



#### ➤ Lucas Kanade:

Frame 1:

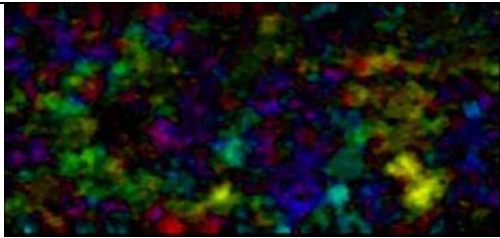
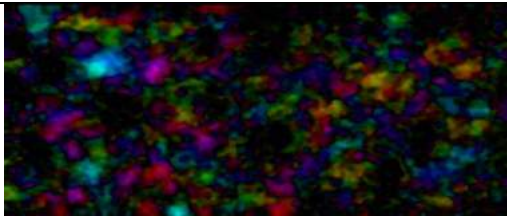


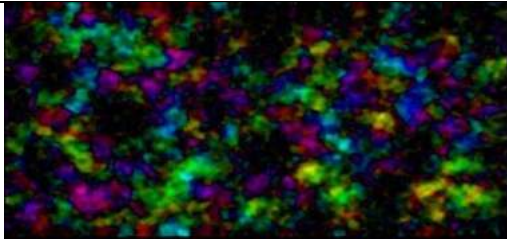
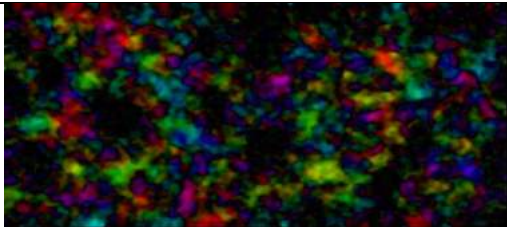
Original video		Time taken for original video is 2mins 49secs
Fast video		Time taken for original video is 3mins 14secs



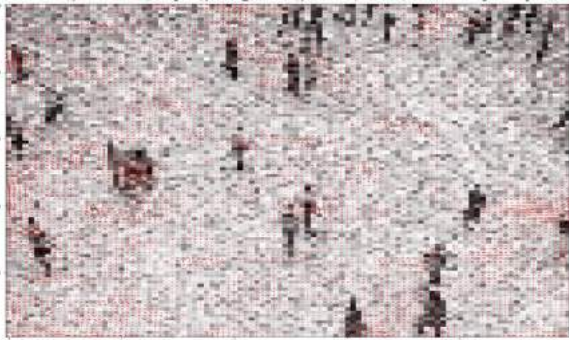
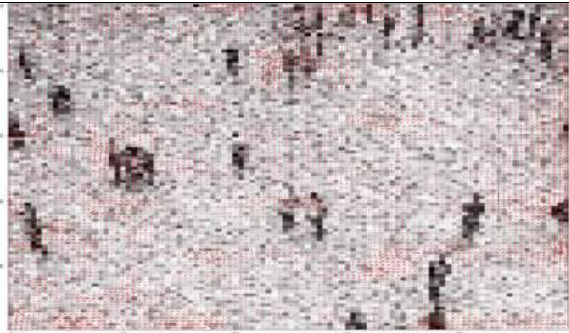
<b>Medium video</b>		<b>Time taken for original video is 2mins 40secs</b>
<b>Slow video</b>		<b>Time taken for original video is 2mins 43secs</b>
<b>Inference</b>	<p>Lucas Kanade algorithm didn't work well in any of the noise footage properly.</p> <p>Algorithm was not able to detect the corners and was not able to estimate the where the flow could in next image sequence</p>	



➤ **Dense Optical Flow:**

<b>Original video</b>		<b>Time taken for original video is 10mins 9secs</b>
<b>Fast video</b>		<b>Time taken for original video is 2mins 43secs</b>

Medium video		Time taken for original video is 5mins 20secs
Slow video		Time taken for original video is 20mins 51secs
Inference	<p>Dense Optical Flow algorithm didn't work well in any of the noise video footage.</p> <p>Dense optical flow able to compute optical flow vector for each and every pixel of each frame.</p>	

➤ Horn and Schunck:

Original video	
Fast video	

<b>Medium video</b>	
<b>Slow video</b>	
<b>Inference</b>	<b>Horn and Schunck Algorithm didn't work for noise video footage. The estimation is very poor.</b>

### C. Write a Discussion section and present the findings in a paragraph:

In Horn and Schunck, as we know it can only handle object/person/things which has constant velocity. In our dataset we have noise also, which made the estimation very poor or not estimated properly the flow of motion in the image sequence.

In Lucas-Kanade, this algorithm didn't work well for the noise dataset. The estimation was very poor. It was unable to detect the flow properly. Which gave an poor estimation of where the flow could be in next image sequence. in

In Dense Optical Flow, as it compute optical flow vector for each and every pixel of each frame. We where not able to estimate the flow in all the footage clearly. And the computational time is much higher then other two algorithm. It didn't give accurate result for noise dataset.

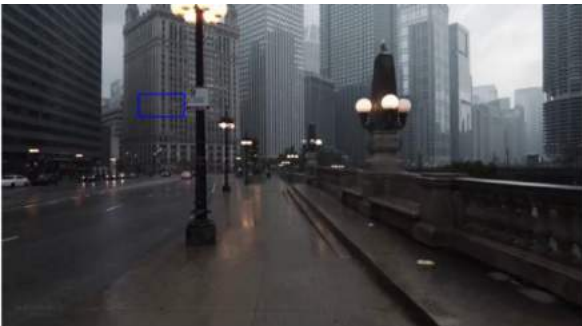
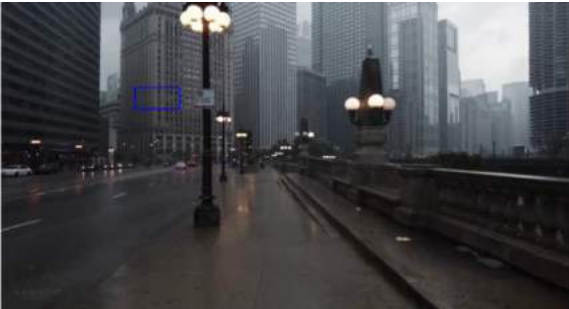
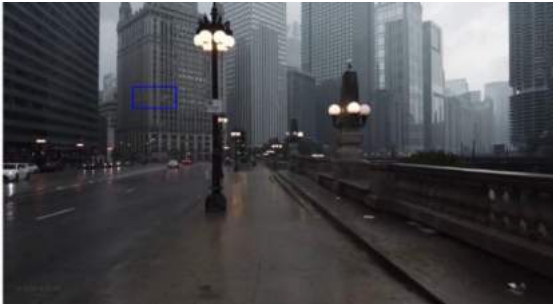
From comparison, None of the algorithm worked for the dataset with noise.

## MEANSHIFT AND CAMSHIFT:

**Dataset:** <https://drive.google.com/drive/folders/1GA449sOZ4AzdH80C3b1f3A83Vla0MCSq?usp=sharing>

### MEANSHIFT-

**Output:**

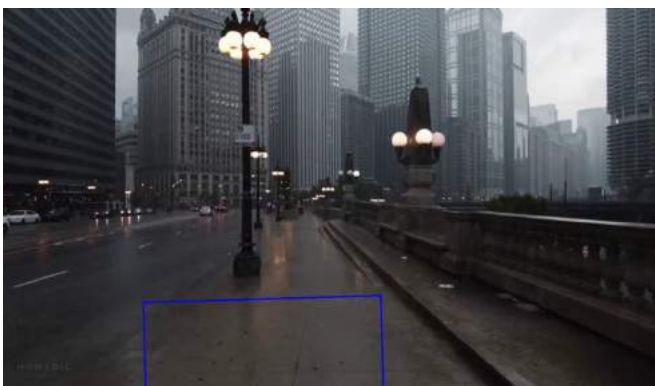
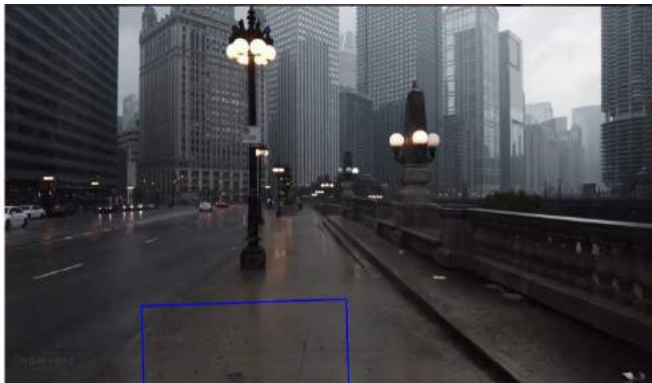


- The mean shift method is a statistical concept that has to do with clustering. The mean shift technique, like most other clustering algorithms, looks for regions in the data set with a large concentration of data points, or clusters.
- The Mean Shift clustering algorithm is an unsupervised clustering algorithm that organizes data without using labeled data to train it. The Mean Shift clustering algorithm is hierarchical in nature, which means it is based on a cluster hierarchy.
- The Mean Shift Algorithm is a clustering algorithm that uses the highest density points or mode value as the major parameter for machine learning development. It's a supervised machine learning algorithm that's not supervised. Kernel Density Estimation (KDE) is the basis for the algorithm. The mode searching algorithm is another name for it. The Kernel is involved in mathematical calculations including data point weighting. The flat kernel and the Gaussian Kernel are two of the most common kernel functions linked with the mean Shift Algorithm. Computer vision and picture segmentation are the most common applications for this approach.



## CAMSHIFT:

### Output:



- The CAMShift algorithm is based on the mean shift algorithm, which is responsible for locating the centre of the object's probability distribution. The primary distinction is that CAMShift adapts to the size of the search window, which is useful when object sizes change as they move closer or farther away from the camera.

## LAB 4b: FACE DETECTION:

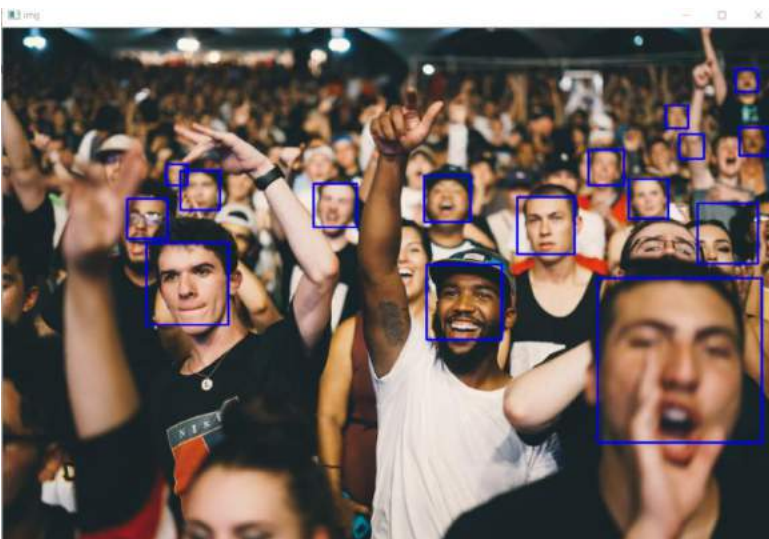
**Aim:** The primary aim of face detection algorithms is to determine whether there is any face in an image or not.

Face identification with **Haar cascades** is a machine learning strategy that involves training a cascade function with a collection of input data. We'll use the face classifier from OpenCV, which already has several pre-trained classifiers for faces, eyes, smiles, and so forth.

### Input:



### Output:



Face identification with Haar cascades is a machine learning strategy that involves training a cascade function with a collection of input data. We'll use the face classifier from OpenCV, which already has several pre-trained classifiers for faces, eyes, smiles, and so forth.

- Accuracy is obtained for the Haar cascade is **96.24%**.

## LAB 3B: (Data population-customized)

### Data Description:

- Customized Dataset Drive Link:

<https://drive.google.com/drive/folders/1oE2vpCSvuDcnnDGpulBO-m3LlpYI9iD?usp=sharing>

- Training Sample Size : 735.
- Testing Sample Size : 183.

Sample Images:



- No. of Positive image: 558
- No. of Negative Image: 336

The dataset contain common scenario scene images where it can be used for object/human detection and tracking.

Challenges in the image dataset:

1. Some of the image has noise and unclear.
2. shadows are there.
3. Image perspective is different, which make difficult yo algorithm for training.
4. Luminous issues in many.
5. Occlusion problem.

## Feature Extraction from the customized dataset:

	sharpness	contrast	mean_sat	std_sat	min_sat	max_sat	mean_bright	std_bright	min_bright	max_bright	colorfulness	hue_hist_1	hue_hist_2	hue_hist_3	shannon_entropy
0	4339.562485	0.09054085506	189.7331085	45.8707028	255	0	119.2013514	56.72508793	255	0	52.11395331	7	44	11.80789252	13.24471705
1	1264.54331	0.2523955587	67.73518107	50.95556063	255	0	88.3499254	47.26059575	255	0	28.83025425	20	268	44.80344147	13.73329963
2	10104.18587	0.2041440555	62.78518926	73.21670064	255	0	139.3614217	53.81622987	255	0	49.87820242	20	34	36.54178655	14.25544132
3	210.8174342	0.3333985518	75.6708705	45.40454651	255	0	81.05027938	45.35023288	252	1	31.30190651	15	77	57.87045568	10.84089054
4	30863.12014	0.2705379285	118.5050255	81.80320625	255	0	101.924505	83.85481448	255	0	44.9155785	20	25	49.33918521	14.67334359
5	1737.510781	0.1223710578	60.33580145	40.81655439	255	0	136.4059045	34.24219207	254	0	22.81212575	14	112	21.30442113	11.02849201
6	4482.415493	0.290606738	67.15709483	53.00101288	255	0	109.8017385	66.6632988	255	0	23.8016552	15	500	48.23961855	13.8895455
7	3652.491726	0.2726632584	83.05838193	61.13382637	255	0	86.15143229	57.85118749	253	0	30.88982689	20	226	48.80673059	13.0321734
8	1253.534837	0.1334883987	37.80309029	27.18032348	255	0	119.8717221	87.28503335	243	0	21.438422	9	18	23.59442301	10.896118
9	3936.885528	0.3203811188	47.72773608	35.61012718	255	0	138.0521918	82.53482925	250	0	25.95934738	20	118	57.44888882	12.87476824
10	6463.191967	0.2001770363	167.952011	82.37404691	255	0	94.39414705	51.48530848	250	0	64.81939326	18	108	36.30968448	13.78474328
11	1794.462913	0.1026865644	18.15059207	21.98257207	255	0	176.5591954	34.82079447	253	0	15.25216673	10	246	15.40815574	10.54361264
12	735.8834151	0.2202855483	163.0568287	49.84034695	255	0	60.88378133	41.806865	251	0	47.83765025	14	138	40.44125354	11.83447577
13	1326.193871	0.301218558	37.55564352	25.55530611	255	0	132.3368945	54.72354751	255	0	29.75871278	10	27	35.01813521	12.47879505
14	1626.304049	0.3073428408	176.2647353	94.62431705	255	0	108.5449479	51.15796617	255	0	85.04494585	17	132	36.68195882	11.33537835
15	1540.187494	0.261685785	39.72861054	16.78463447	255	0	183.7468504	23.99601398	255	0	20.34233075	8	111	47.74474555	10.21089518
16	2884.875185	0.2895980059	63.71385758	58.53047018	255	0	123.7405482	64.39751897	253	0	58.2179485	19	180	51.21943306	12.67219044
17	340.8223981	0.2959896338	28.78777659	31.94578173	255	0	125.806525	32.87951723	254	0	20.84594325	20	17	45.82422515	11.37128825
18	389.2339991	0.1303353813	43.66548267	21.41077405	255	0	62.08117333	48.76816653	253	0	9.803087811	10	48	23.36883325	10.6268848
19	4543.997416	0.2503518552	76.52955	56.6254123	255	0	121.3750413	73.93381826	253	0	34.53494711	18	84	46.42380165	13.54136483
20	1740.513038	0.331487578	63.75419859	72.30881852	255	0	88.45794044	70.82502285	253	0	27.19711482	20	97	59.3380581	12.53376452
21	3273.593783	0.2290144378	58.65225458	45.35548223	255	0	104.1372157	51.21484659	253	0	21.40222073	3	541	49.93587785	12.43377402
22	2769.083494	0.2251612041	148.4475544	75.33105108	255	0	88.77101302	50.72272709	255	0	51.82119959	16	217	49.3026554	13.15744944
23	3121.814996	0.1709107487	83.5895461	55.83454147	255	0	109.3295522	58.89736301	255	0	59.18421205	14	135	59.43192371	12.30119189
24	1134.723621	0.2106647718	71.70331084	64.18585732	255	0	148.5692984	51.08565497	253	0	51.2284575	14	84	39.2126616	12.19334376
25	189.3223723	0.3036110781	36.74452789	33.12132813	228	0	119.732918	73.24265748	255	22	31.702495	15	88	53.9838352	10.75847048
26	1477.05281	0.2638717899	102.6814096	114.8113723	255	0	108.5987952	105.8542544	255	0	3.064398969	14	147	46.70930274	8.888452688
27	793.4140196	0.3051822179	183.2336102	73.71014059	255	0	86.45179541	54.89356492	253	0	71.31987585	20	148	55.343517	12.31874489
28	4354.452187	0.2956860793	44.90560286	53.44545857	255	0	114.5688542	59.11825192	250	0	42.39897191	20	412	51.3158516	12.39873019
29	452.105426	0.2396853579	73.98482193	51.40354710	255	0	103.4513807	54.23251875	255	0	37.57148248	18	105	42.99741406	11.83257787
30	16887.73955	0.2059349872	87.20837549	72.47933205	255	0	82.33784458	75.87384959	255	0	22.45863385	20	41	45.8387732	12.7182377
31	4774.205445	0.1806865493	186.7818173	73.71683454	255	0	42.18796106	54.8722672	255	0	46.11099513	20	43	29.57827233	11.0737375
32	7304.919332	0.30091915424	62.35178959	62.28972905	255	0	89.05962407	59.21812101	253	0	39.8874221	20	108	54.87965789	13.67330305
33	1331.514119	0.3210305106	31.89398183	50.78078437	255	0	173.2314036	56.43548853	253	0	35.32796555	12	350	57.46274335	10.73690338
34	2353.601399	0.3043305016	19.805168	24.3584147	212	0	170.2842187	63.4532353	255	22	16.50296474	10	421	38.5573135	12.13068847
35	8328.615779	0.24029720333	46.9733887	43.05728581	255	0	129.5107827	62.89701839	255	0	35.24888336	6	302	44.87330329	12.5778131
36	2783.543547	0.2708032612	85.51344332	80.88430304	255	0	106.4424828	66.37801192	259	0	59.78932664	20	92	48.38847328	13.0398208
37	3712.638803	0.1123816134	112.4738058	82.75458134	255	0	97.12611978	54.91582865	253	0	57.14263084	18	59	39.39518284	12.40886593

Feature Name	Purpose	Data Type of the feature
Sharpness	Sharpness is a combination of two factors: resolution and acutance. Resolution is straightforward and not subjective. It's just the size, in pixels, of the image file. All other factors equal, the higher the resolution of the image—the more pixels it has—the sharper it can be. Acutance is a little more complicated. It's a subjective measure of the contrast at an edge. There's no unit for acutance—you either think an edge has contrast or think it doesn't. Edges that have more contrast appear to have a more defined edge to the human visual system.	Float64
Contrast	<p>The term contrast refers to the amount of color or grayscale differentiation that exists between various image features in both analog and digital images. Images having a higher contrast level generally display a greater degree of color or grayscale variation than those of lower contrast.</p> <p>Contrast enhancement processes adjust the relative brightness and darkness of objects in the scene to improve their visibility.</p>	Float64
Saturation	The term saturation describes the depth or intensity of colour present within an image.	Float64
Brightness	The term saturation describes the depth or intensity of colour present within an image.	Float64
Colorfulness	Colorfulness is the "attribute of a visual perception according to which the perceived color of an area appears to be more or less chromatic. The colorfulness evoked by an object depends not only on its spectral reflectance but also on the strength of the illumination, and increases with the latter unless the brightness is very high	Float64
Hue_Histogram	An image histogram is a type of histogram that acts as a graphical representation of the tonal distribution in a digital image. It plots the number of pixels for each tonal value. By looking at the histogram for a specific image a viewer will be able to judge the entire tonal distribution at a glance.	Float64/int64
Shannon entropy	<p>In Image, Entropy is defined as corresponding states of intensity level which individual pixels can adapt. It is used in the quantitative analysis and evaluation image details, the entropy value is used as it provides better comparison of the image details.</p> <p>Shannon entropy as a measure of image information is extensively used in image processing applications.</p>	Float64



## LAB 5: Data Analysis and Modeling

- Dataset Used for the project use case is: "COCO DATASET". [[Link](#)]

The "Common Objects In Context" (COCO) dataset is a collection of challenging, quality datasets for computer vision, mostly using cutting-edge neural networks.

COCO is a large-scale object detection, segmentation, and captioning dataset.

- COCO has several features:
  - Object segmentation.
  - Recognition in context.
  - Super pixel stuff segmentation.
  - 330K images (>200K labeled).
  - 1.5 million object instances.
  - 80 object categories.

### **Type of Challenge addressed in the dataset:**

When developing Object Detection and Image Segmentation models, common data issues include:

1. Dimensions and aspect ratios of images (especially dealing with extreme values).
2. Identifies the composition's imbalances, bounding box dimensions, and aspect ratios (for instance a lot of small objects)
3. Not the right data preparation for your dataset.
4. Modeling strategy is not consistent with the data.

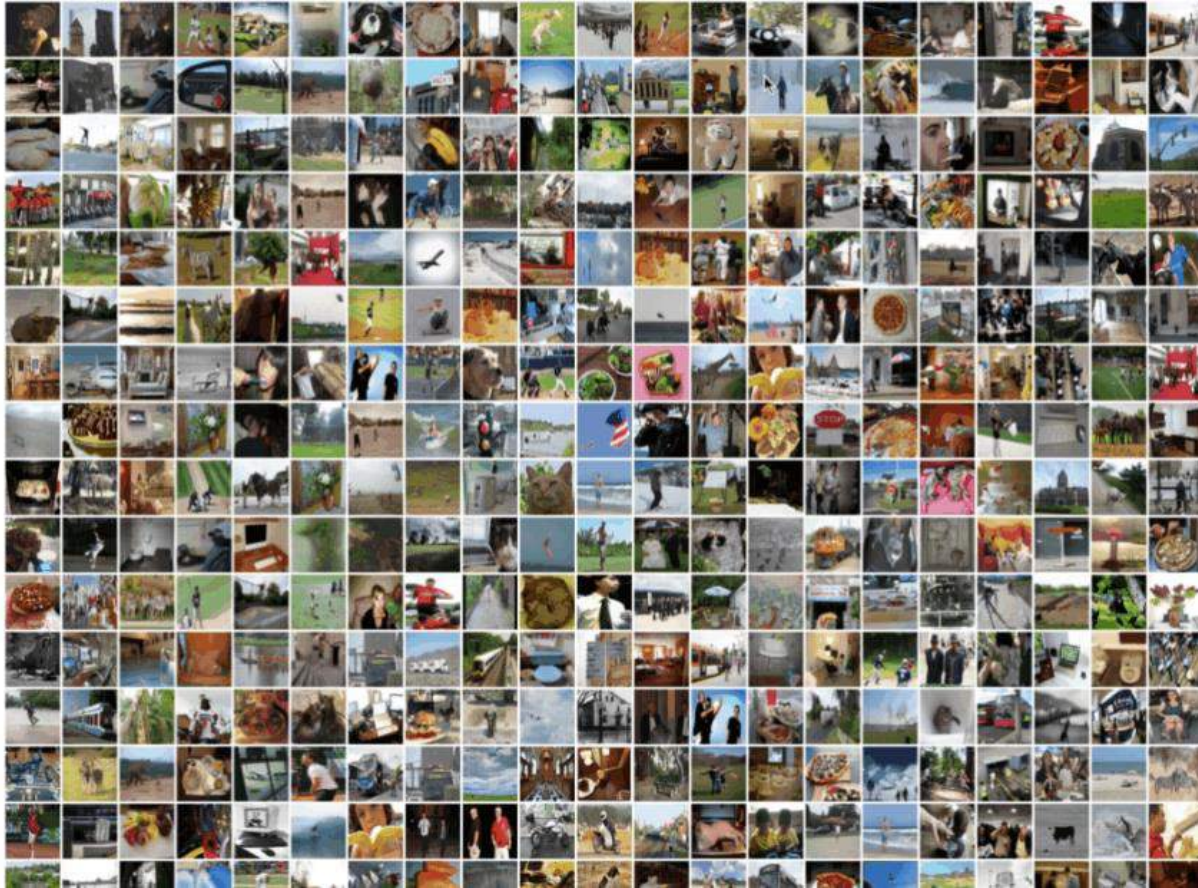
### **General data quality:**

1. Acquire a basic sense of a dataset and visually examine it.
2. Ensure sure it isn't damaged and doesn't have any glaring distortions.

3. Verify that all of the files can be read; you don't want to discover this during training.

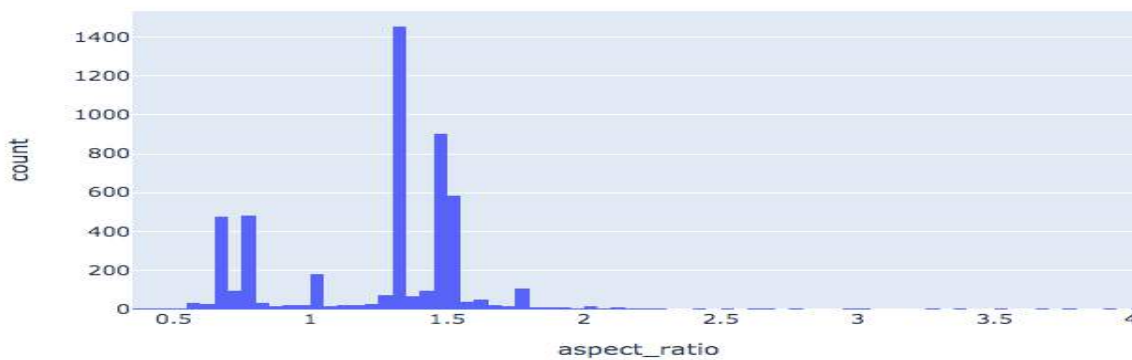
Here is to visualize as many pictures as possible:

- Plot them in a jupyter notebook using matplotlib:



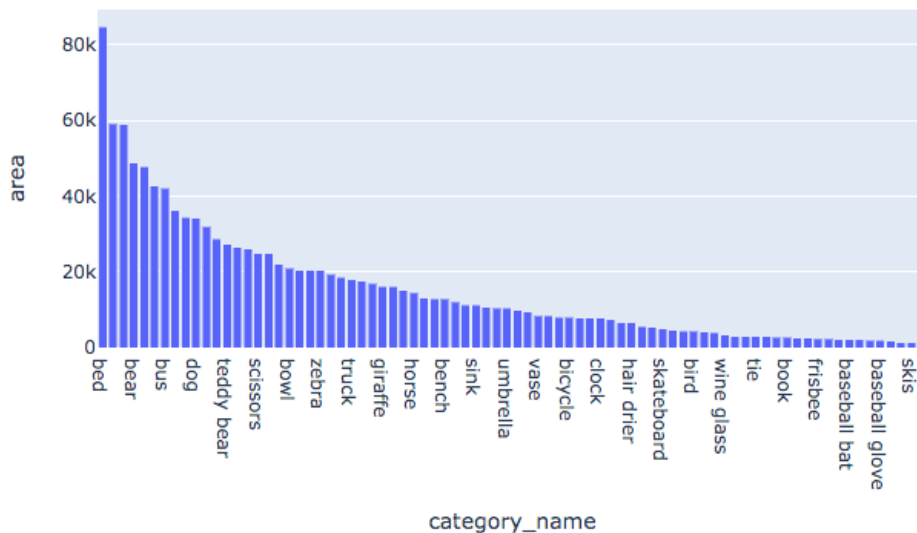
### Image sizes and aspect ratios:

- Training your model on image patches (randomly selected during training or extracted before training)
- Resizing the entire dataset to avoid doing this every time you load your data.



### Label (objects) sizes and dimensions :

- For instance, it might be able to greatly reduce the model if your dataset solely contains really large items.
- On the other hand, if you have photographs with small objects (10x10px, for example), it's possible that you won't be able to train the model with this configuration.
- When it comes to box or mask dimensions, the following factors should be taken into account:
  1. Aspect ratios.
  2. Size (Area).



### Understanding preprocessing sequences:

Any computer vision system must include preprocessing.

It is more difficult to apply it to object recognition and segmentation problems than to simple picture classification since various changes (like rotation or crop) must be applied to both the source and destination images (masks or bounding boxes). The following frequent transformations need a target transform:

1. Affine transformations,
2. Cropping,
3. Distortions,
4. Scaling,
5. Rotations
6. and many more.

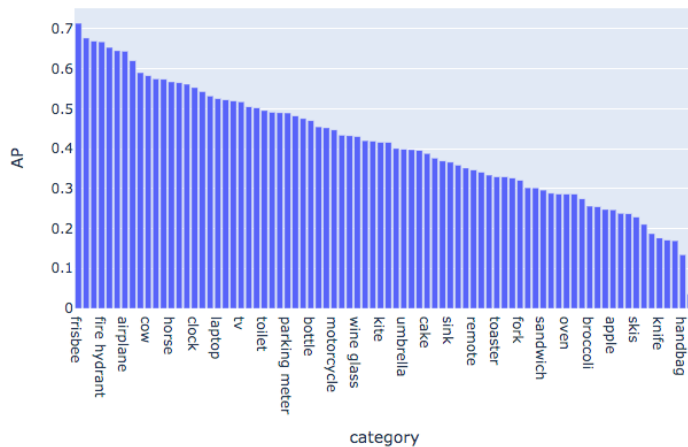
## Results understanding:

Let's take a look at the formal Coco Challenge and how that system of assessment functions

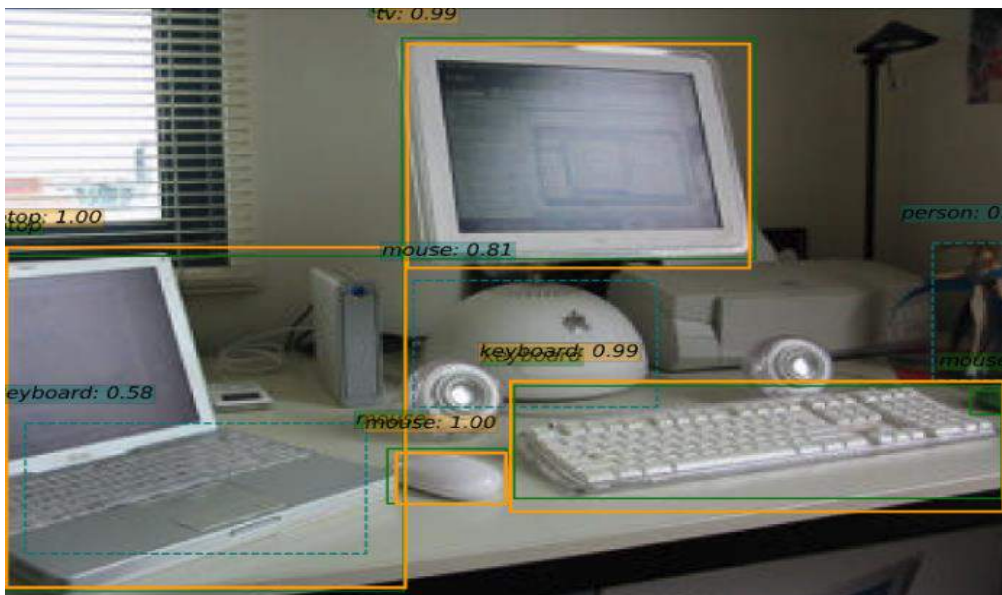
```

Average Precision (AP) @[ IoU=0.50:0.95 | area= all | maxDets=100 ] = 0.373
Average Precision (AP) @[ IoU=0.50 | area= all | maxDets=100 ] = 0.590
Average Precision (AP) @[ IoU=0.75 | area= all | maxDets=100 ] = 0.402
Average Precision (AP) @[ IoU=0.50:0.95 | area= small | maxDets=100 ] = 0.219
Average Precision (AP) @[ IoU=0.50:0.95 | area=medium | maxDets=100 ] = 0.409
Average Precision (AP) @[ IoU=0.50:0.95 | area= large | maxDets=100 ] = 0.481
Average Recall (AR) @[ IoU=0.50:0.95 | area= all | maxDets= 1 ] = 0.310
Average Recall (AR) @[ IoU=0.50:0.95 | area= all | maxDets= 10 ] = 0.495
Average Recall (AR) @[ IoU=0.50:0.95 | area= all | maxDets=100 ] = 0.521
Average Recall (AR) @[ IoU=0.50:0.95 | area= small | maxDets=100 ] = 0.323
Average Recall (AR) @[ IoU=0.50:0.95 | area=medium | maxDets=100 ] = 0.560
Average Recall (AR) @[ IoU=0.50:0.95 | area= large | maxDets=100 ] = 0.667

```



## Visualizing results:





## LAB 6 : Optical Flow based Deep Learning Tracking [RAFT].

**Aim:** Our aim is to estimate the motion of the image intensities over time in a video using optical-flow algorithms like RAFT. i.e. to calculate Pixel points within the images, and to provide the estimation of where the points could be in the next image sequence in our shopping video dataset.

➤ **RAFT :**

RAFT can be divided into three stages:

1. **Feature extractors:** The network input consists of two consecutive frames. To extract features from these two images, we use two CNNs with shared weights.
2. **Visual Similarity:** Visual similarity is calculated as the inner product of all pairs of feature maps.
3. **Iterative update:** An iterative update is a sequence of Gated Recurrent Unit (GRU) cells that combine all data we have calculated before.

**Input:**



**Output:**



**Inference:**

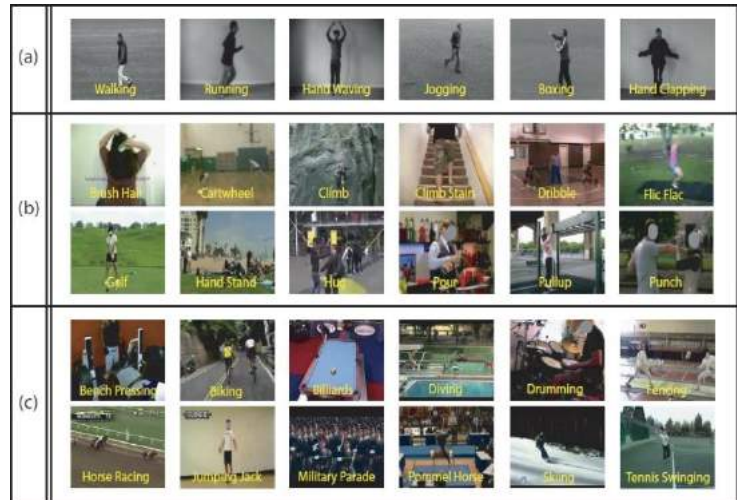
- Here the Optical flow is a per pixel prediction and the main idea is that it assumes a brightness constancy, meaning it tries to estimate how the pixels brightness moves across the screen over time.

## LAB 7: Activity recognition: (CNN+LSTM,CNN+RNN)



**Aim:** Need to identify one or more agents' activities and objectives from a sequence of observations on those actions and the surrounding circumstances.

### ➤ UCF50 - Action Recognition Data Set :

UCF50 data set's 50 action categories collected from youtube are: Baseball Pitch, Basketball Shooting, Bench Press, Biking, bicycling, Billiards Shot, Breaststroke, Clean and Jerk, Diving, Drumming, Fencing, Golf Swing, Playing Guitar, High Jump, Horse Race, Horse Riding, Hula Hoop, Javelin Throw, Juggling Balls, Jump Rope, Jumping Jack, Kayaking, Lunges, Military Parade, Mixing Batter, Nun chucks, Playing Piano, Pizza Tossing, Pole Vault, Pommel Horse, Pull Ups, Punch, Push Ups, Rock Climbing Indoor, Rope Climbing, Rowing, Salsa Spins, Skate Boarding, Skiing, Skijet, Soccer Juggling, Swing, Playing Tabla, TaiChi, Tennis Swing, Trampoline Jumping, Playing Violin, Volleyball Spiking, Walking with a dog, and Yo Yo.



### Outputs:

ACTION	Sample image	Description	Output
Riding a bicycle		The video of a kid riding a bicycle on a road	[bicycling]
Riding motorcycle		The video obtained from the camera fixed o the helmet of a person riding a motorcycle.	[Biking]

Doing  
YOGA



This video obtained  
from the camera, a  
person doing yoga.

[yoga]

## LAB 8 : Computer Vision based library [solvePnp() library]

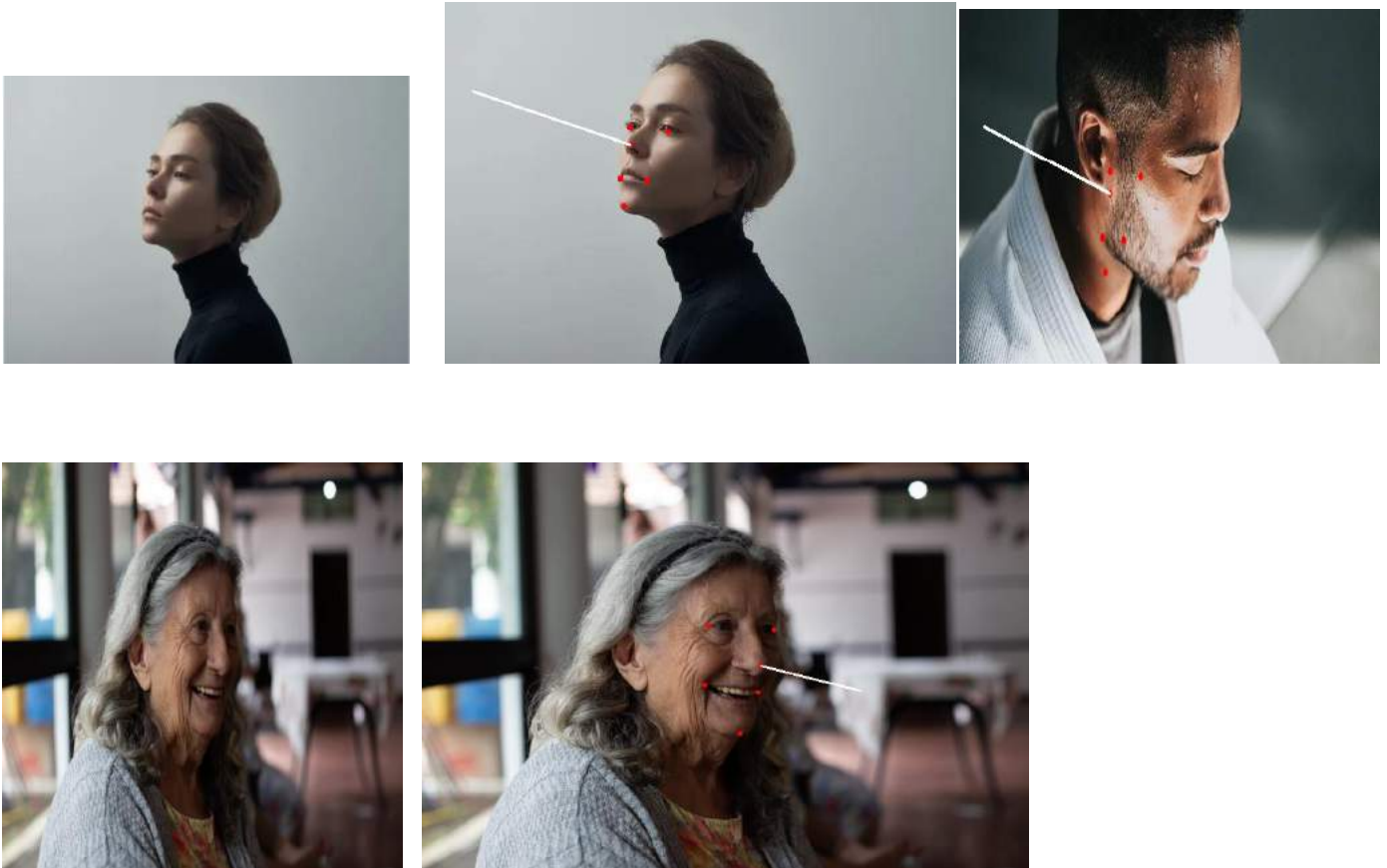
### ➤ Understand the PnP Problem:

- The PnP problem is very common in Computer Vision and stands for the Perspective n-Points problem. In this issue, we cannot determine the pose of an object with respect to the camera after being provided with the 2D and 3D coordinates.
- This can be understood with the example of face tracking during an online exam. The pose of an object with respect can change with the change in direction.
- The following two types of motions facilitate this change:
  - The first type of motion is translational motion, which can happen along any of the three axes. The object moves in a uniform motion in any particular direction, thereby changing its coordinates.
  - The second type of motion is the rotational motion, in which the object can revolve around any of the three axes.

### ➤ Use the opencv.solvepnp() Function to Solve the PnP Problem:

- The solvepnp() function from the OpenCV library is used for the pose estimation of a given object with respect to the camera, thus solving the PnP problem. It returns rotational and translational vectors.
- It uses the 2D and 3D coordinates of the object with the camera matrix. The coordinates supplied are of the different features of the face.
- These features are the nose, corners of the mouth, chin, and both of the eyes.

### Input and Output:





The camera matrix and the 2D and 3D points of the image's face features serve as the major parameters. It returns the vectors that identify the 3D points of the pose using these values.

Using the `projectPoints()` function, we can project these points in 2D with respect to the camera. We next use these points to plot a line that will serve as the determined pose's representation in the image.

**Inference:**

Finally, we plot a line to represent the determined pose in the image using these points.