

IMAGE TAGGING AND ROAD OBJECT DETECTION

Guide : Prof. Anoop M. Namboodiri

Mentor : Sangeeth Reddy

- Richard Deepak
- Sridhar Deshpande
- Rahul Juluru
- Kshama Pandey

5

Agenda

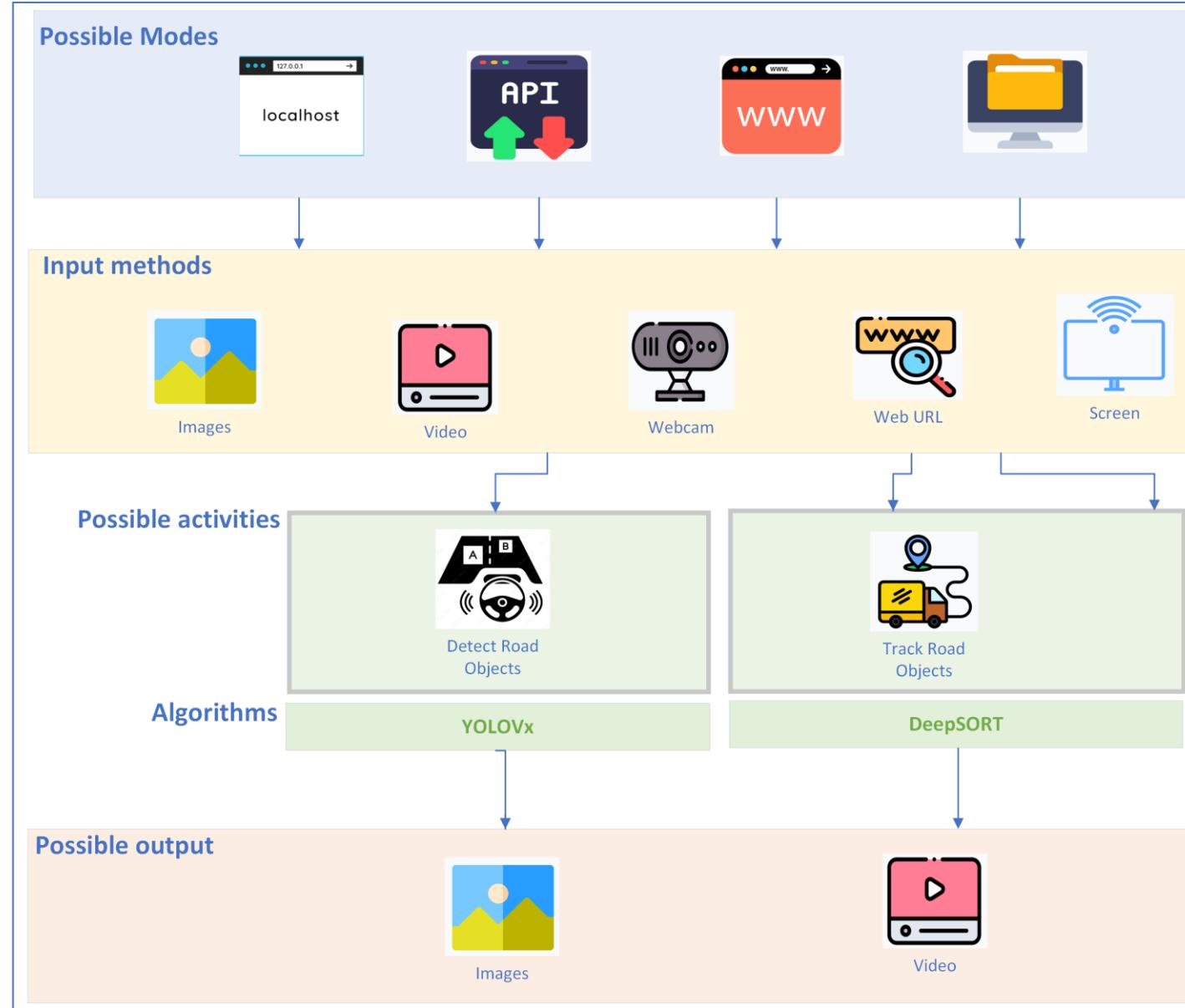
- Objective .
- Methodology and design considerations.
- Outcome in stages.
- Challenges.
- Application Demo.
- Going Beyond.
- Summary.



Objective

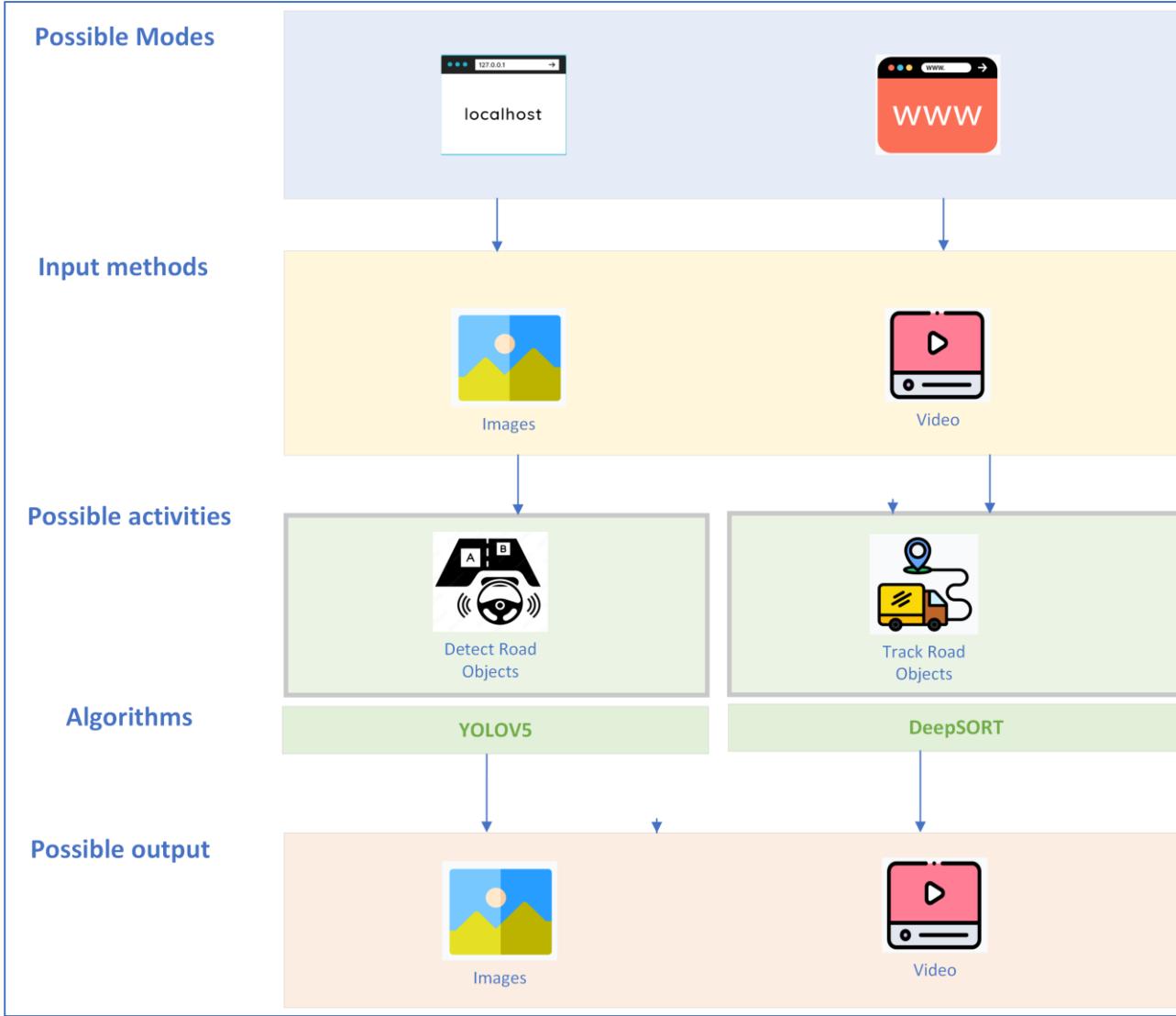
Build an application to detect multiple objects,
tag and tracking in a video.

Methodology and design considerations

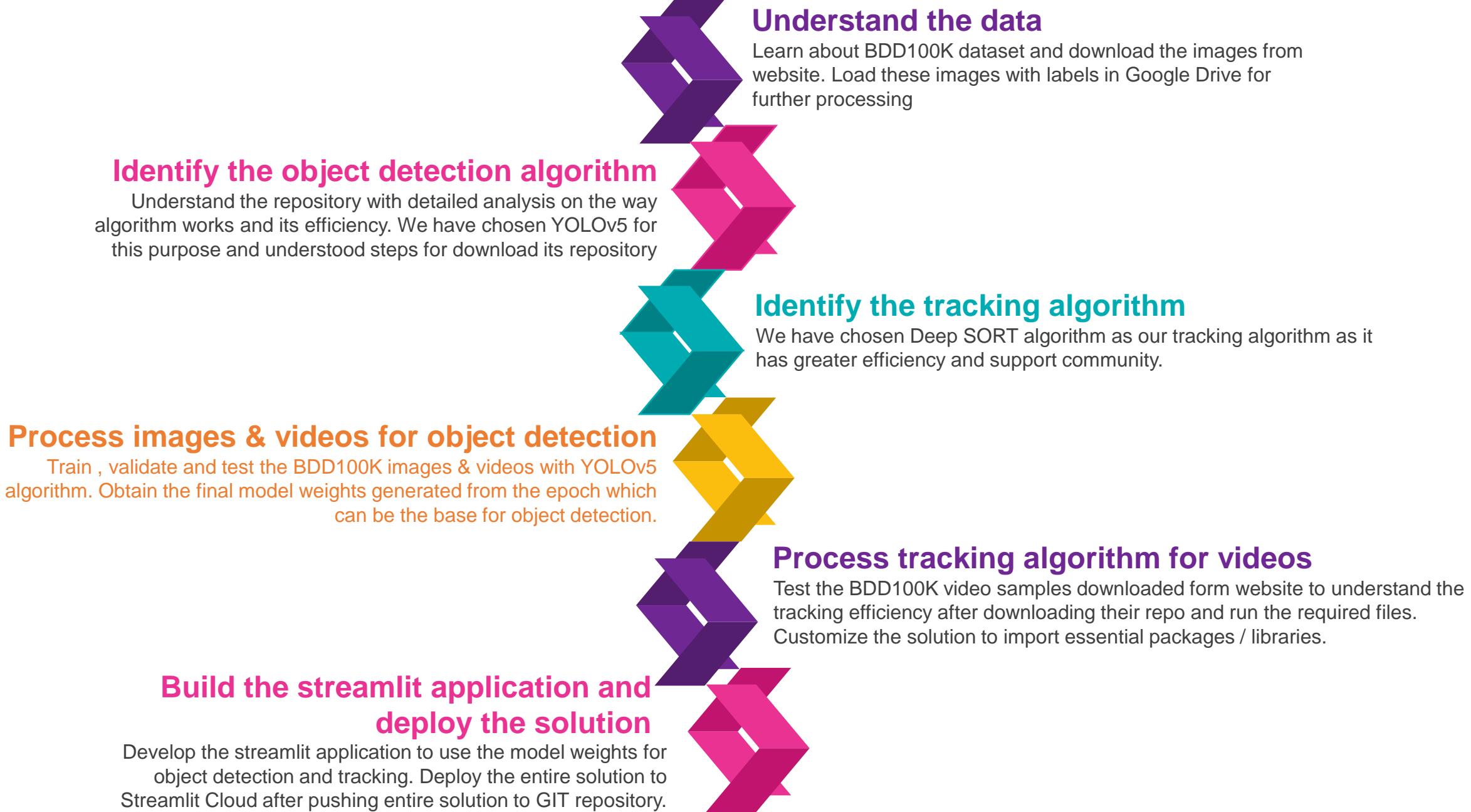


Methodology and design considerations

Project Scope



Outcome in stages



Outcome in stages

Stage 1 – Understand the data

Consists of 100,000 videos.
Each video is about 40 seconds long, 720p, and 30 fps

BDD100K
Dataset

Geographically diverse

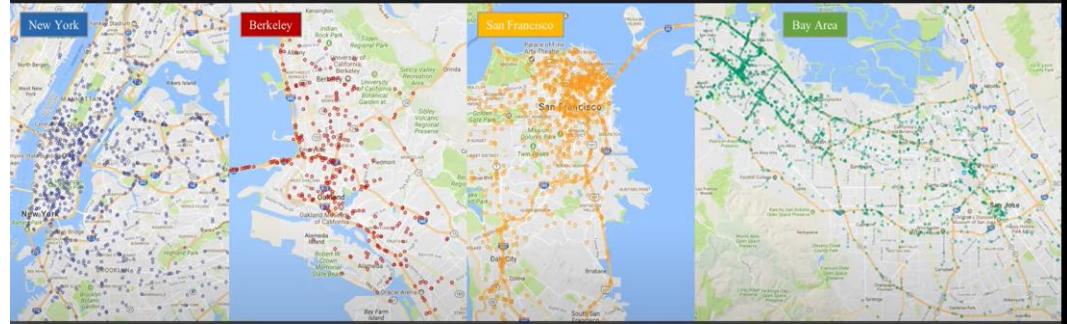
Scene Diversity : Scene = City , Tunnel , Highway , Parking , Residential

Time of the day = Dusk , Daytime , night

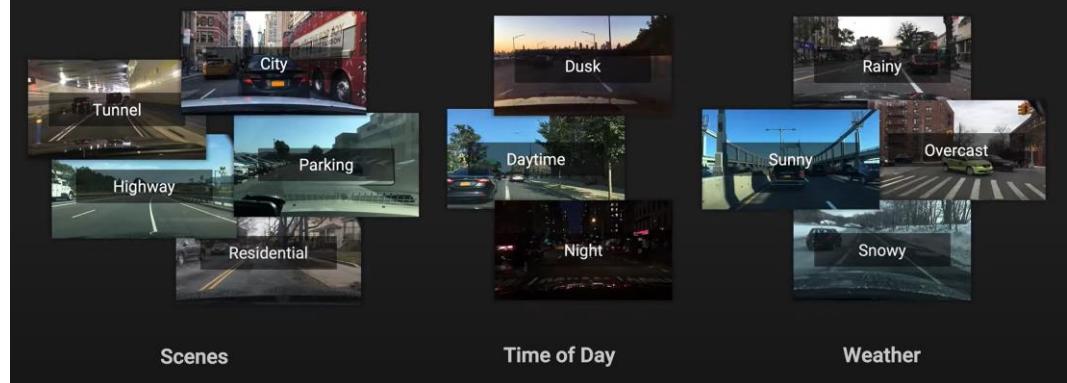
Weather = Rainy , Sunny , Overcast , Snowy

There are >10 objects per image with pixel annotation and tracking

Geography Diversity



Scene Diversity



Outcome in stages

Stage 1 – Understand the data



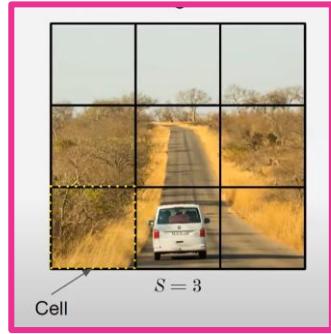
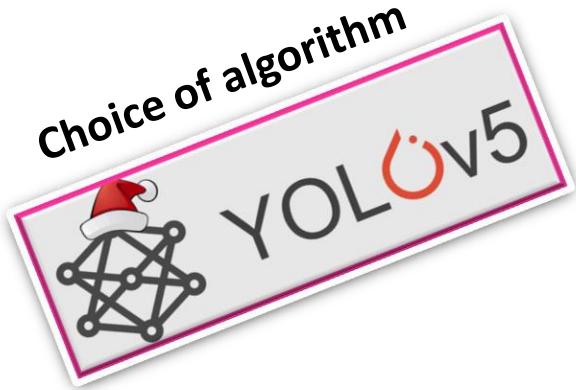
Original Image

BDD100k dataset bounding box JSON format

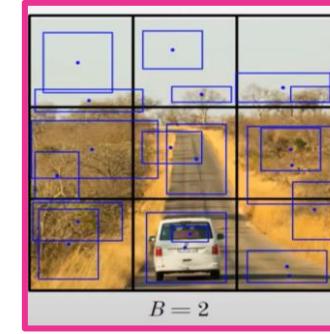
```
1  [
2  {
3      "name": "000d4f89-3bcbe37a.jpg",
4          "attributes": {
5              "weather": "overcast",
6              "scene": "city street",
7              "timeofday": "daytime"
8          },
9          "timestamp": 10000,
10         "labels": [
11             {
12                 "category": "traffic sign",
13                     "attributes": {
14                         "occluded": false,
15                         "truncated": false,
16                         "trafficLightColor": "none"
17                     },
18                     "manualShape": true,
19                     "manualAttributes": true,
20                     "box2d": {
21                         "x1": 1000.698742,
22                         "y1": 281.992415,
23                         "x2": 1040.626872,
24                         "y2": 326.91156
25                     },
26                     "id": 0
27             },
28         {
29             "category": "traffic sign",
30             "attributes": {
```

Outcome in stages

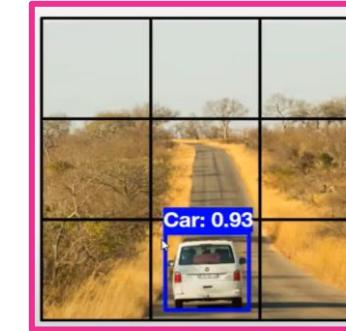
Stage 2 – Identify the object detection algorithm



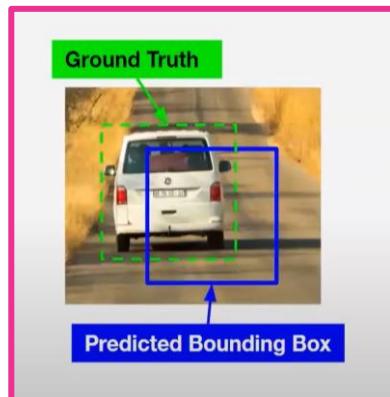
Divide the image to cells



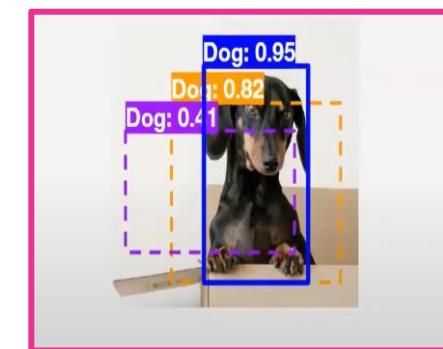
Predicting bounding boxes within each cell



Selection of bounding box which is above confidence threshold



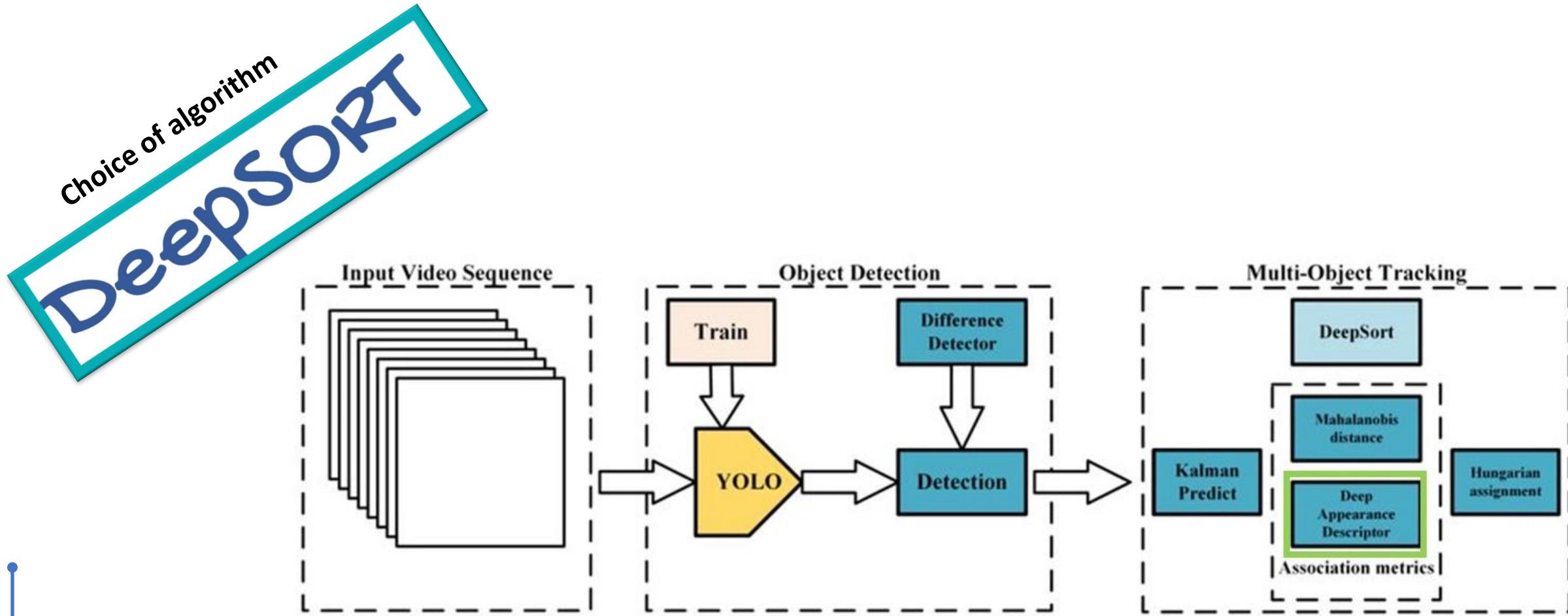
Intersection over Union



Non max suppression

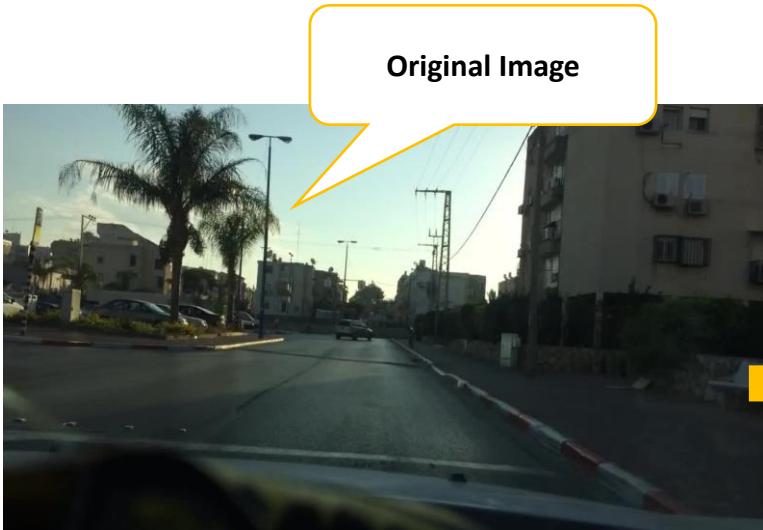
Outcome in stages

Stage 3 – Identify the tracking algorithm



Outcome in stages

Stage 4 – Process images and videos for object detection



Original Image

```
dd100k_label...ges_val.json* □ X
1 [
2   {
3     "name": "000d4f89-3bcbe37a.jpg",
4       "attributes": {
5         "weather": "overcast",
6         "scene": "city street",
7         "timeofday": "daytime"
8       },
9     "timestamp": 10000,
10    "labels": [
11      {
12        "category": "traffic sign",
13        "attributes": {
14          "occluded": false,
15          "truncated": false,
16          "trafficLightColor": "none"
17        },
18        "manualShape": true,
19        "manualAttributes": true,
20        "box2d": {
21          "x1": 1000.698742,
22          "y1": 281.992415,
23          "x2": 1040.626872,
24          "y2": 326.91156
25        },
26        "id": 0
27      },
28    {
29      "category": "traffic sign",
30      "attributes": {
```

BDD100k dataset bounding box JSON format

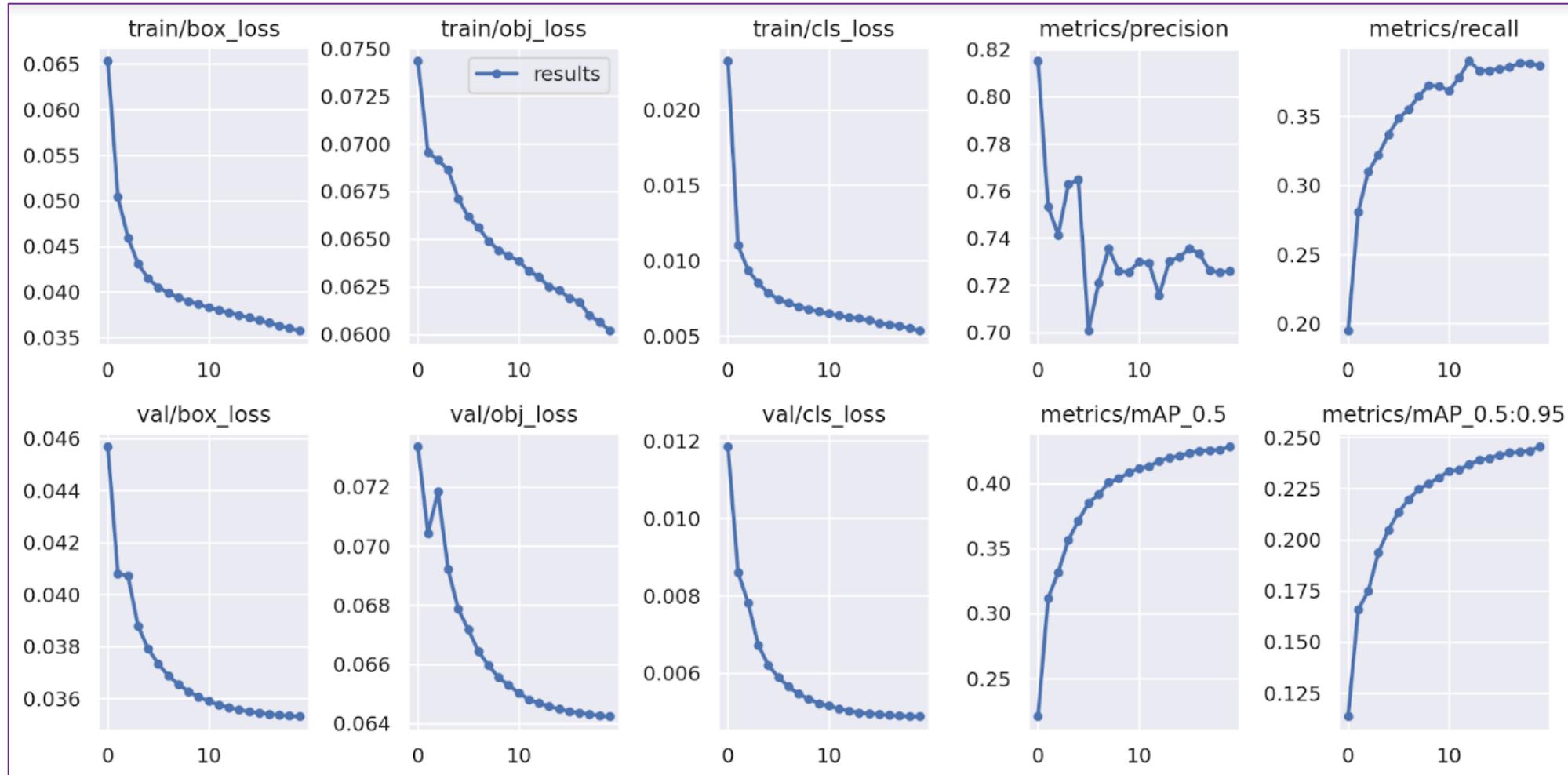
	bounding box x1,y1,x2,y2 normalized					
2	0.32104250507812504	0.4995917312499995	0.03926019609374998	0.0440816236111111		
2	0.2588552611822023	0.4901953158935303	0.062384705760595335	0.05467654567594956		
2	0.23528997304687502	0.4940815263888886	0.06612243515625002	0.06244896666666667		
2	0.1738167703125	0.48760738138567267	0.1456759875	0.07154148499356758		
2	0.058102510156249995	0.4720407166666667	0.03202805625	0.02571428055555608		
2	0.015361844921875001	0.4674488784722222	0.030723689843750002	0.06061223194444435		
2	0.46025830133564527	0.5259084653235868	0.05424690826620946	0.05630693620272922		

bounding box x1,y1,x2,y2 normalized

Converted to TXT file with required
class and Bbox with mapping data and
labels

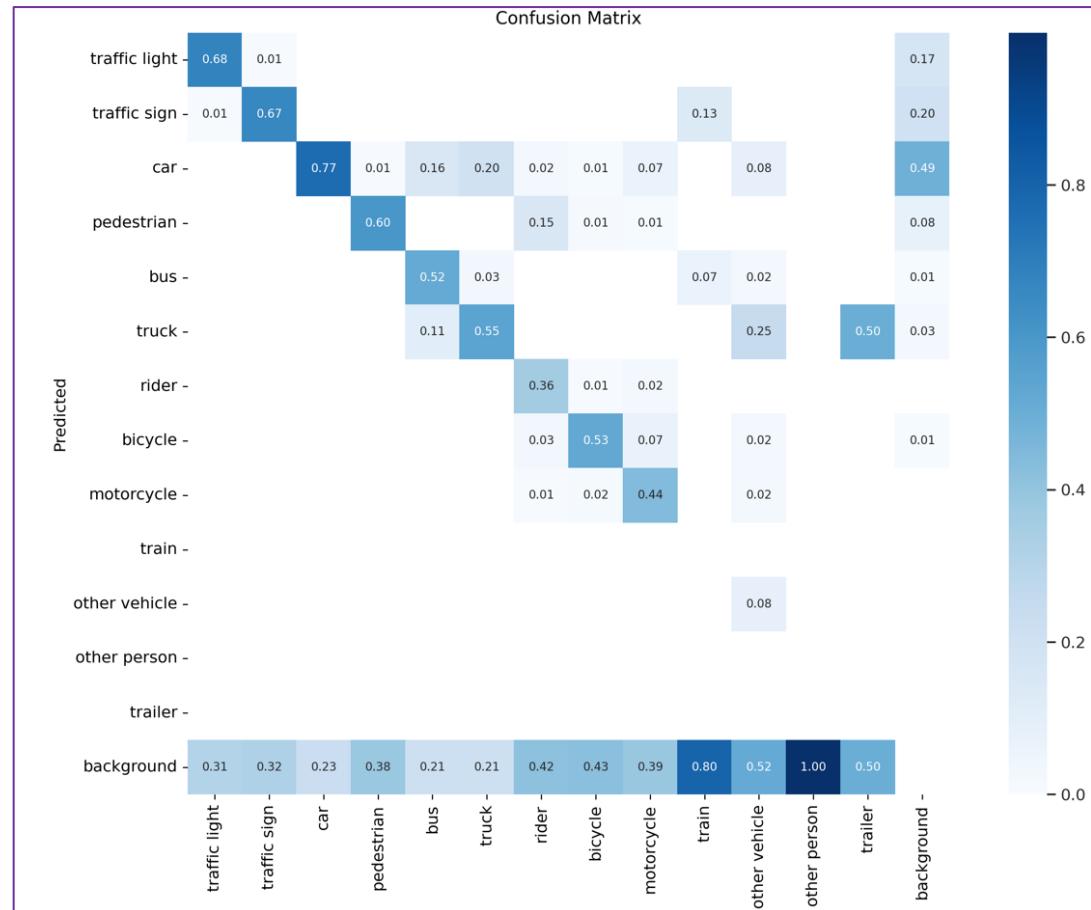
YOLO-v5

Training Metrics



YOLO-v5

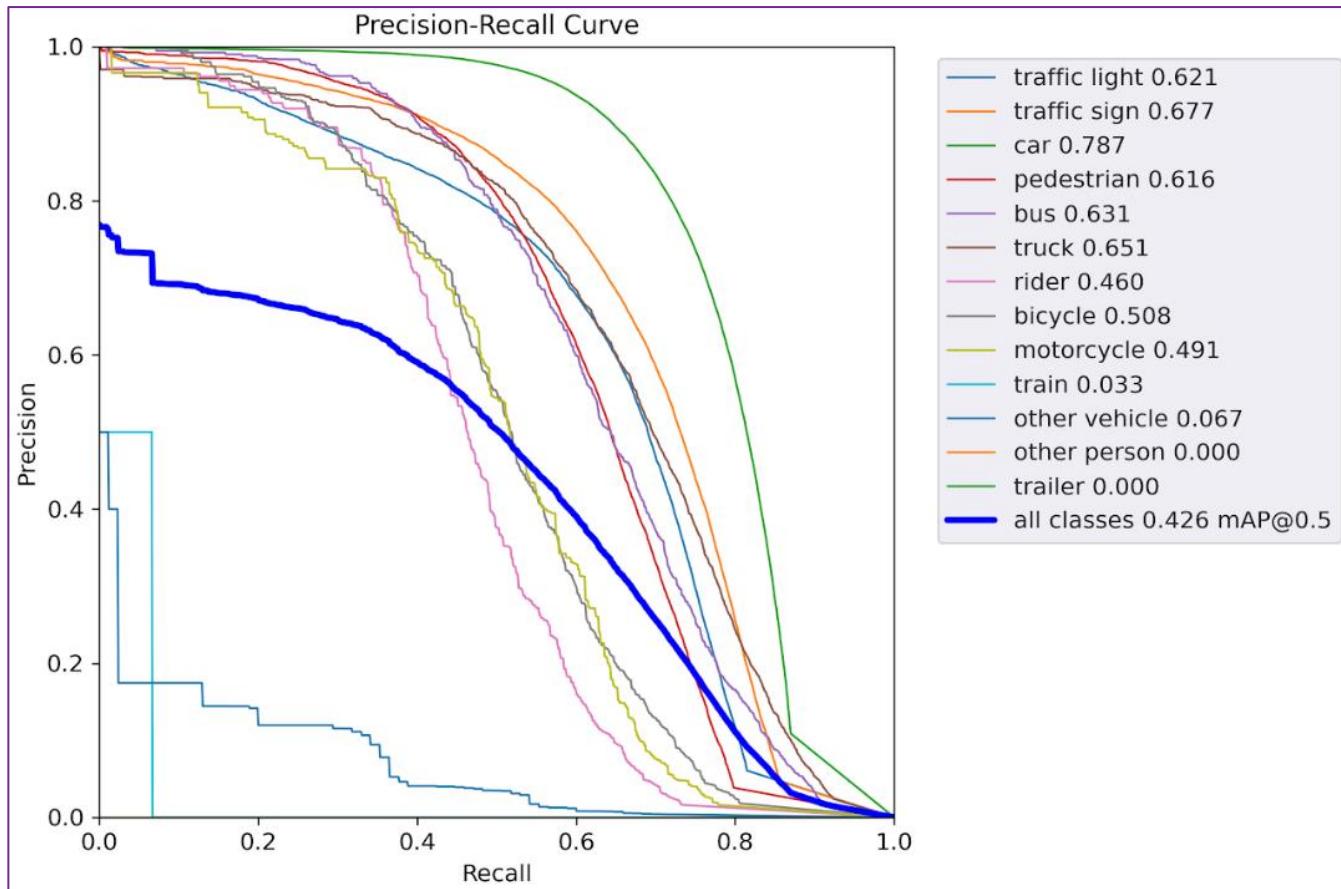
Confusion Matrix & validation results



Class	Images	Instances	P	R	mAP50	mAP50-95:
all	10000	186033	0.727	0.387	0.426	0.244
traffic light	10000	26884	0.676	0.601	0.621	0.234
traffic sign	10000	34724	0.728	0.623	0.677	0.36
car	10000	102837	0.796	0.722	0.787	0.505
pedestrian	10000	13425	0.719	0.553	0.616	0.307
bus	10000	1660	0.698	0.551	0.631	0.488
truck	10000	4243	0.691	0.594	0.651	0.477
rider	10000	658	0.733	0.388	0.46	0.231
bicycle	10000	1039	0.601	0.477	0.508	0.247
motorcycle	10000	460	0.667	0.446	0.491	0.244
train	10000	15	1	0	0.0326	0.0261
other vehicle	10000	85	0.143	0.0706	0.0672	0.0507
other person	10000	1	1	0	0	0
trailer	10000	2	1	0	0	0

YOLO-v5

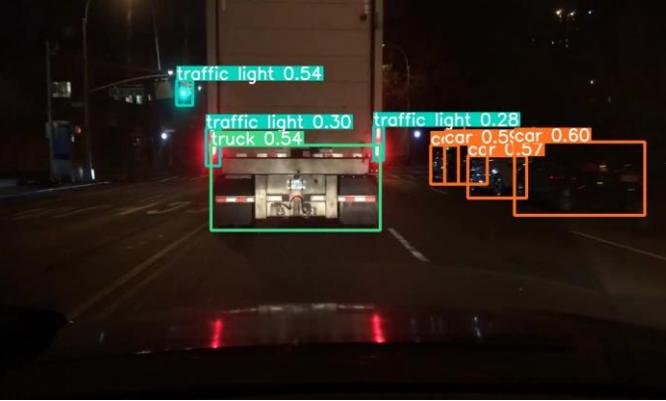
PR Curves



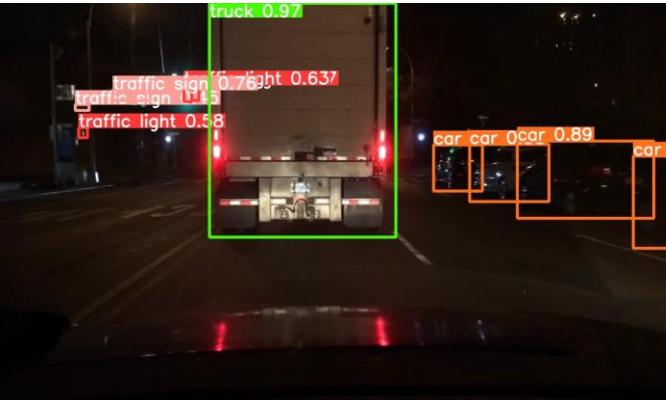
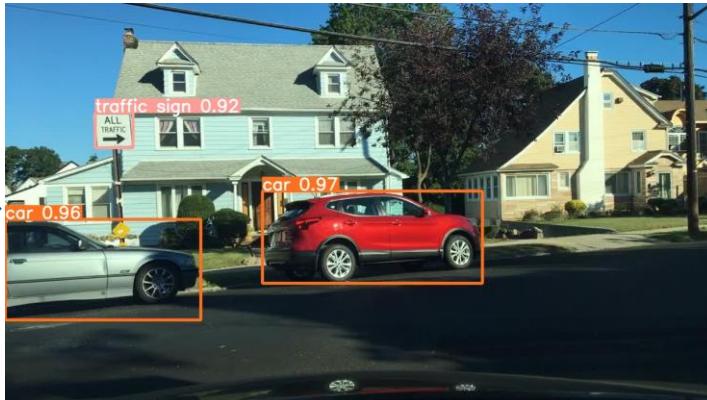
YOLO-v5

Trained model vs Pretrained model

Pre-trained
(COCO dataset)

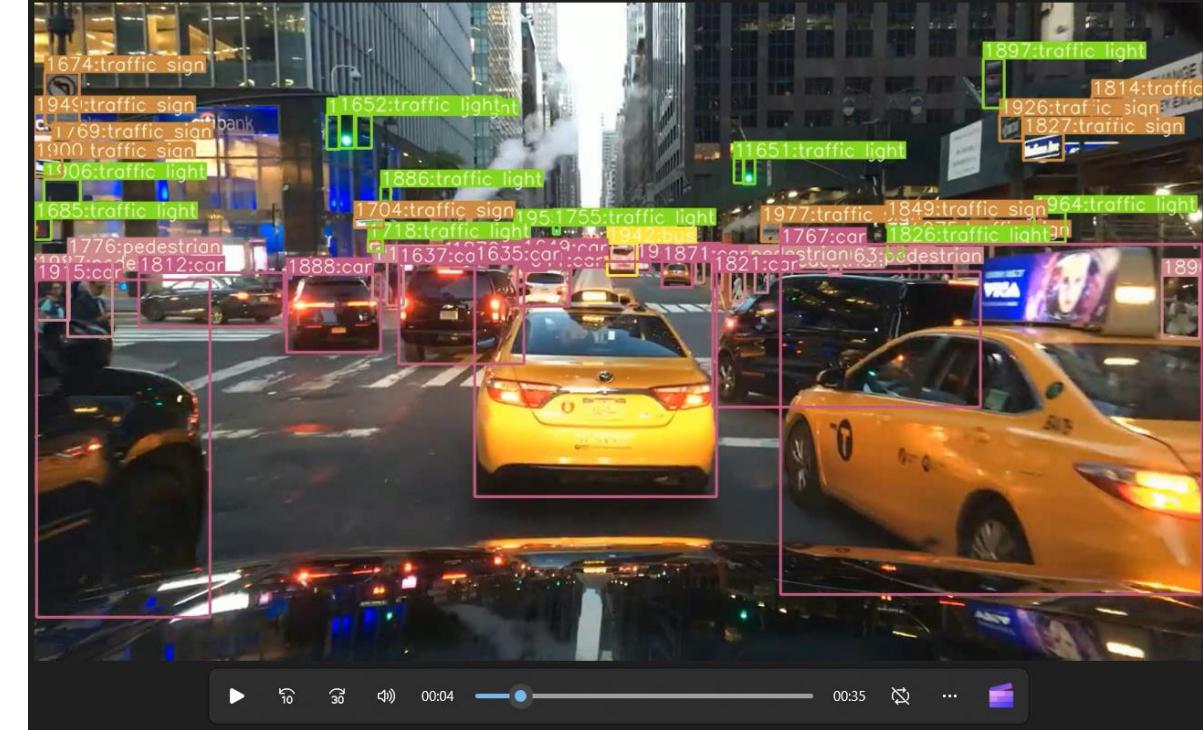
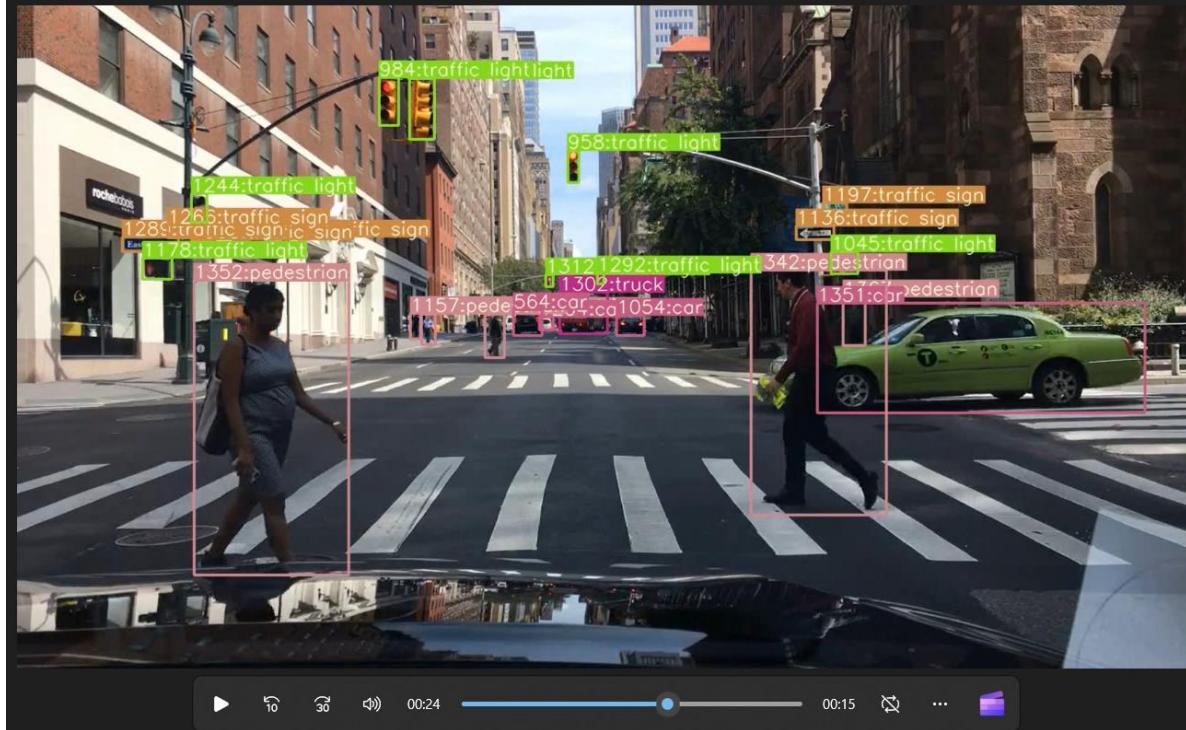


Trained model
BDD100K dataset)



Outcome in stages

Stage 5 – Processed videos for object tracking



YOLO-v5

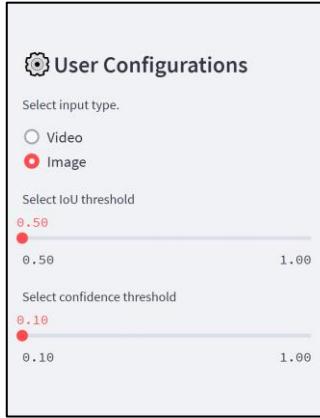
Detection + Deep SORT tracker output



Tracked GIF via YOLOv5 + DeepSORT

Outcome in stages

Stage 6 – Develop streamlit application and deploy onto Streamlit Cloud

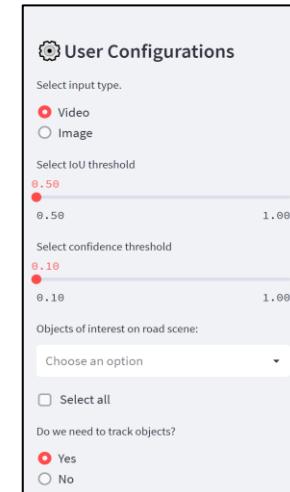
**User Configurations**
Select input type.
 Video
 Image
Select IoU threshold

0.50 1.00
Select confidence threshold

0.10 1.00

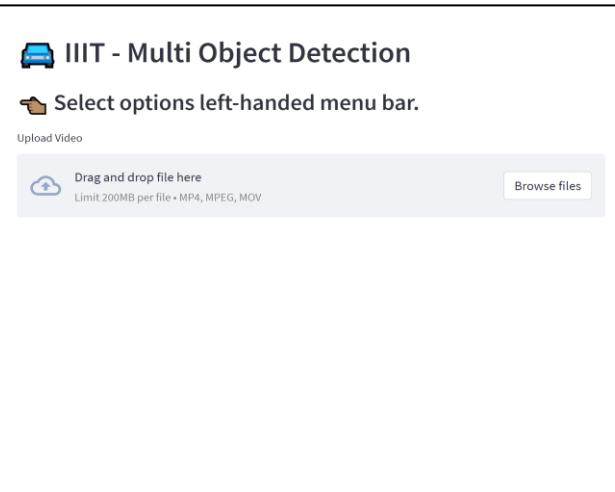
**IIIT - Multi Object Detection**
Select options left-handed menu bar.
Upload An Image
Drag and drop file here
Limit 200MB per file • PNG, JPEG, JPG
Browse files

Application for Image detection

**User Configurations**
Select input type.
 Video
 Image
Select IoU threshold

0.50 1.00
Select confidence threshold

0.10 1.00
Objects of interest on road scene:
Choose an option
 Select all
Do we need to track objects?
 Yes
 No

**IIIT - Multi Object Detection**
Select options left-handed menu bar.
Upload Video
Drag and drop file here
Limit 200MB per file • MP4, MPEG, MOV
Browse files

Application for Video detection and tracking

Application Demo

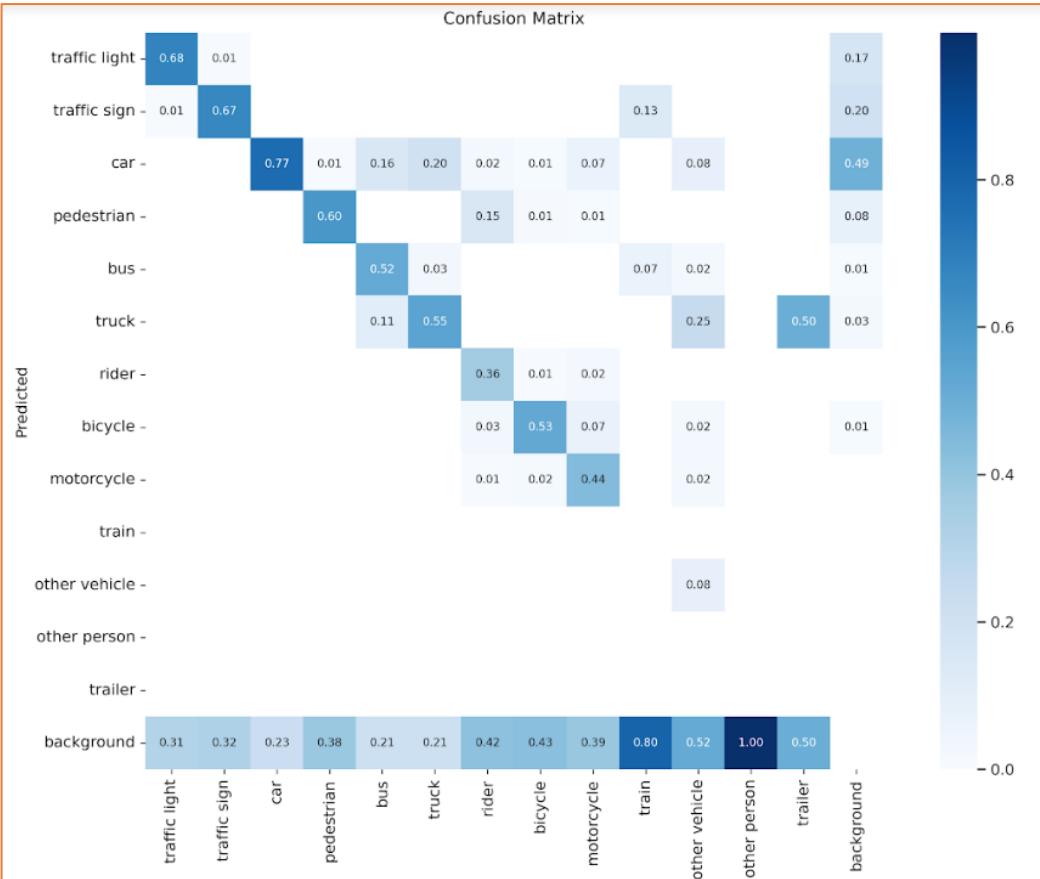
Challenges

1. Missing labels for some of the images of BDD100k data.
2. Downloading of BDD100k videos were time consuming.
3. Resource constraint – GPU causing disconnection of runtime in colab.
4. Understanding cloud constraints
5. Understanding installation of packages in streamlit cloud.
6. How to store the models so as to be loaded in the cloud environment.
7. Understanding Codec formats compatible with browsers.
8. Resource constraints while testing the web-UI locally.

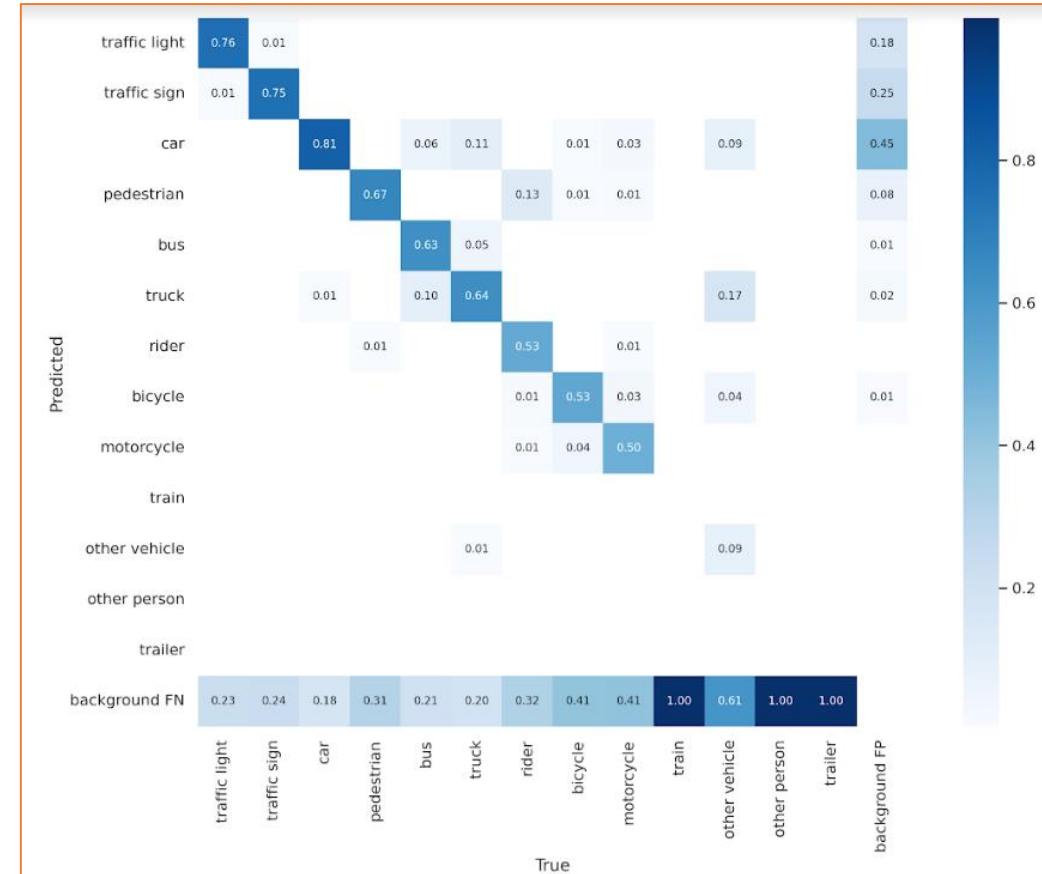
Going Beyond

1.YOLO-V7

YOLOv5 vs YOLOv7



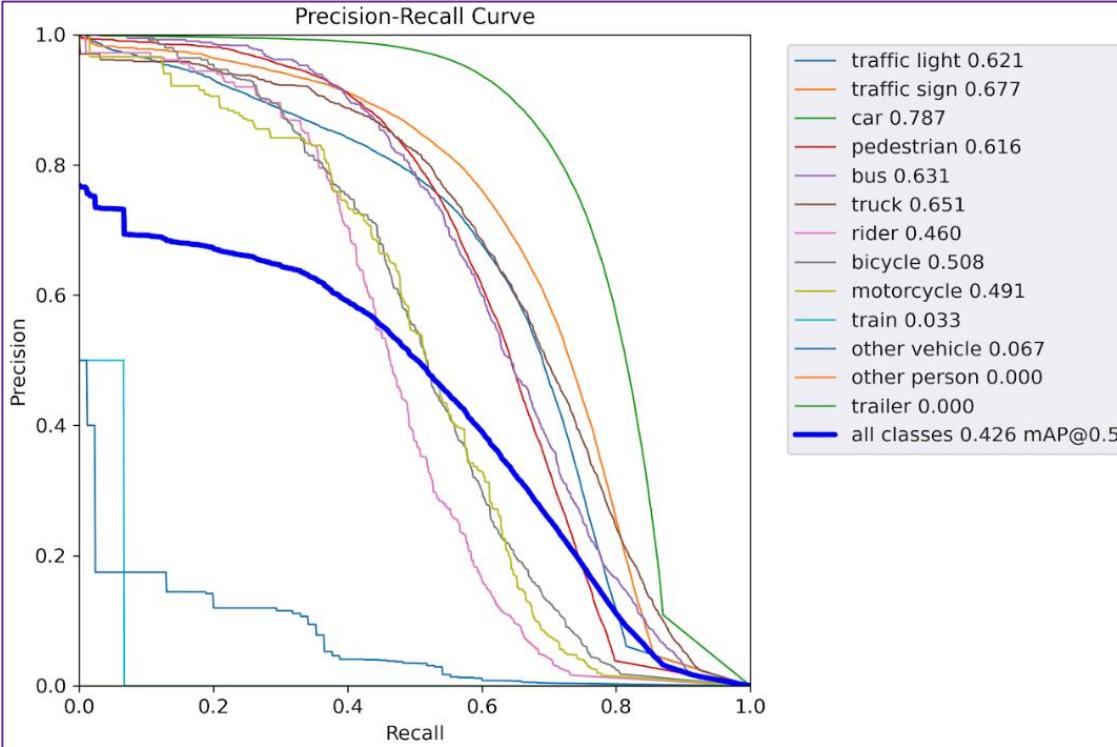
V5



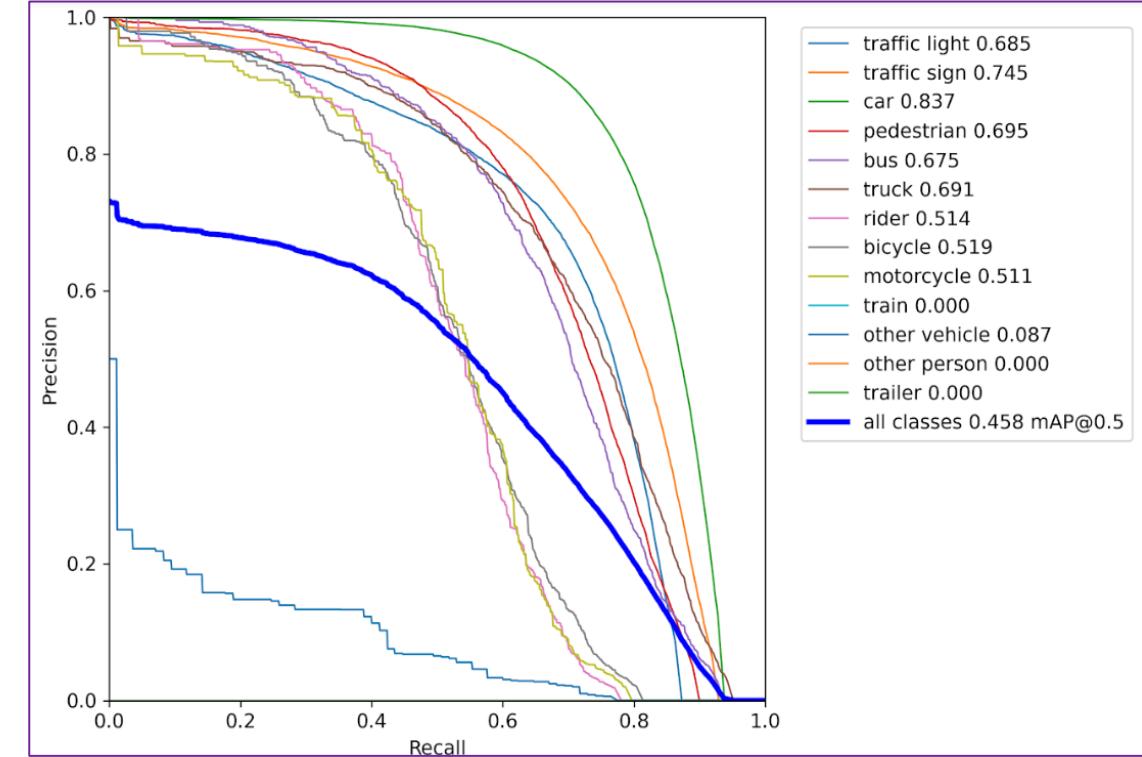
V7

1.YOLO-V7

YOLOv5 vs YOLOv7



V5



V7

1.YOLO-V7

Class	Images	Instances	P	R	mAP50	mAP50-95:
all	10000	186033	0.727	0.387	0.426	0.244
traffic light	10000	26884	0.676	0.601	0.621	0.234
traffic sign	10000	34724	0.728	0.623	0.677	0.36
car	10000	102837	0.796	0.722	0.787	0.505
pedestrian	10000	13425	0.719	0.553	0.616	0.307
bus	10000	1660	0.698	0.551	0.631	0.488
truck	10000	4243	0.691	0.594	0.651	0.477
rider	10000	658	0.733	0.388	0.46	0.231
bicycle	10000	1039	0.601	0.477	0.508	0.247
motorcycle	10000	460	0.667	0.446	0.491	0.244
train	10000	15	1	0	0.0326	0.0261
other vehicle	10000	85	0.143	0.0706	0.0672	0.0507
other person	10000	1	1	0	0	0
trailer	10000	2	1	0	0	0

V5

Class	Images	Labels	P	R	mAP@.5	mAP@.5:.95:
all	10000	186033	0.742	0.422	0.458	0.258
traffic light	10000	26884	0.708	0.668	0.685	0.258
traffic sign	10000	34724	0.727	0.701	0.745	0.396
car	10000	102837	0.833	0.76	0.837	0.522
pedestrian	10000	13425	0.756	0.617	0.695	0.344
bus	10000	1660	0.752	0.586	0.675	0.517
truck	10000	4243	0.704	0.63	0.691	0.505
rider	10000	658	0.611	0.491	0.514	0.257
bicycle	10000	1039	0.611	0.497	0.519	0.249
motorcycle	10000	460	0.73	0.457	0.511	0.251
train	10000	15	1	0	0	0
other vehicle	10000	85	0.218	0.0824	0.0873	0.0519
other person	10000	1	1	0	0	0
trailer	10000	2	1	0	0	0

V7

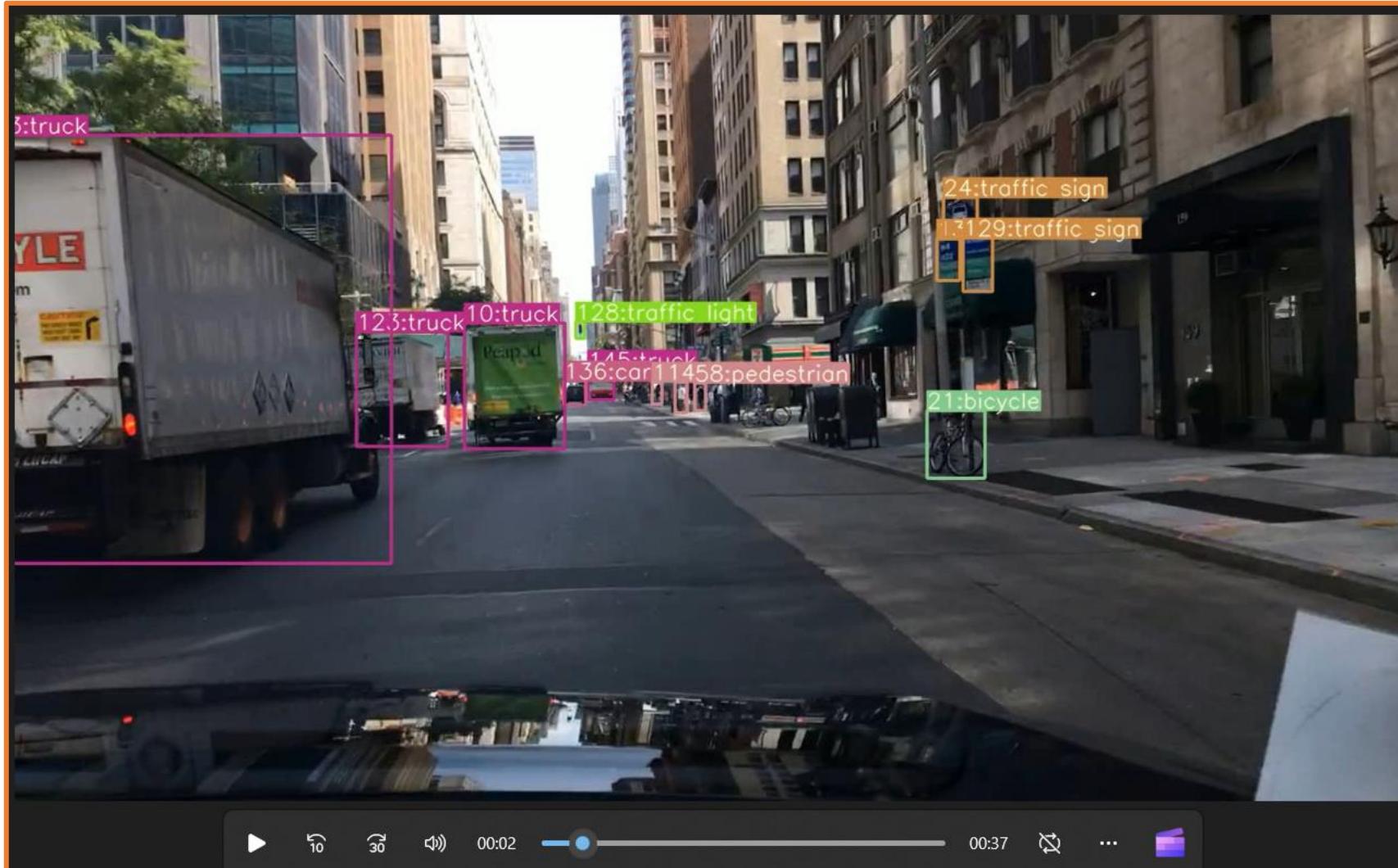
1.YOLO-V7

YOLOv5 vs YOLOv7



2.YOLO-V7 – Detection + Tracker

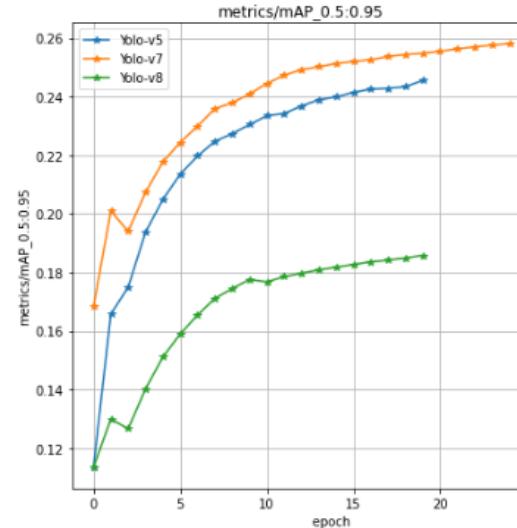
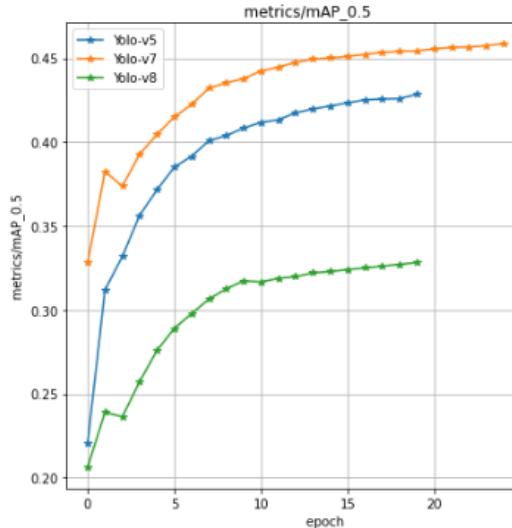
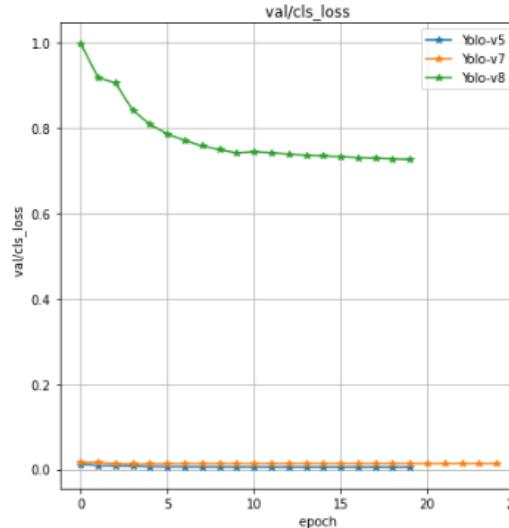
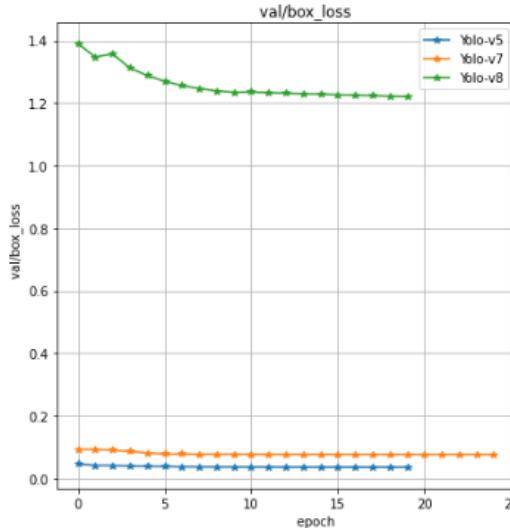
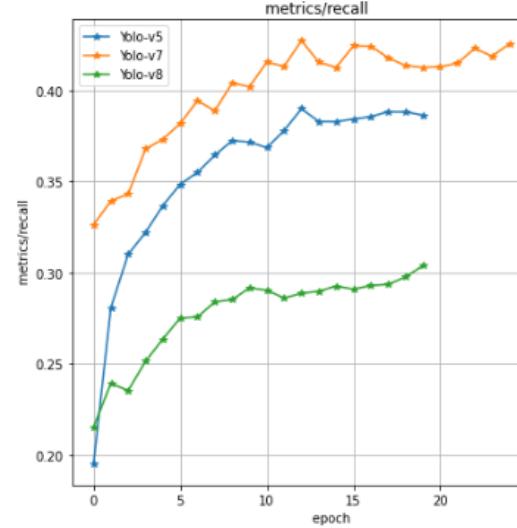
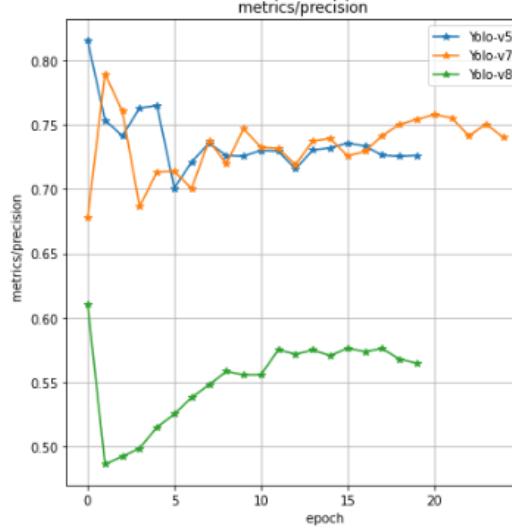
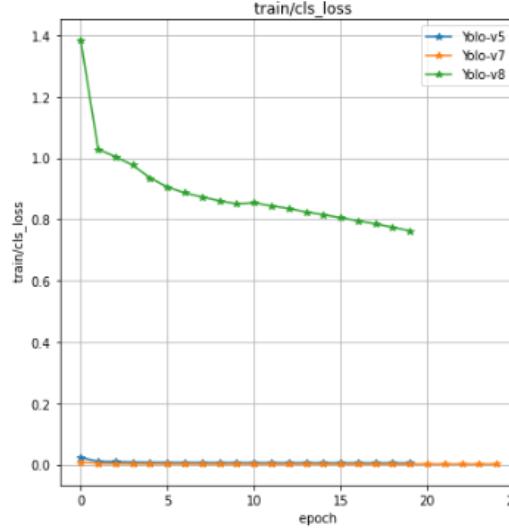
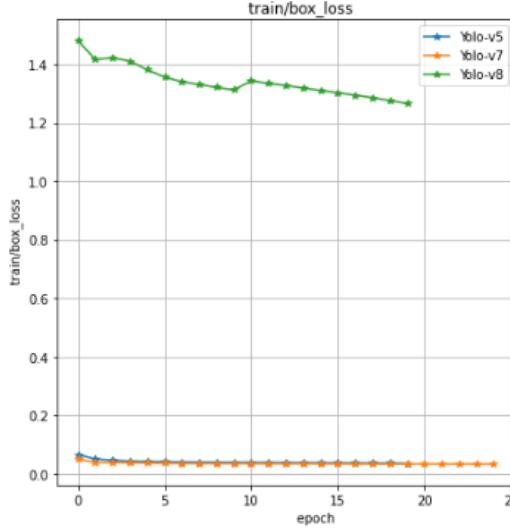
Results:



3.YOLO-V8

YOLOv5 vs YOLOv7 vs YOLOv8

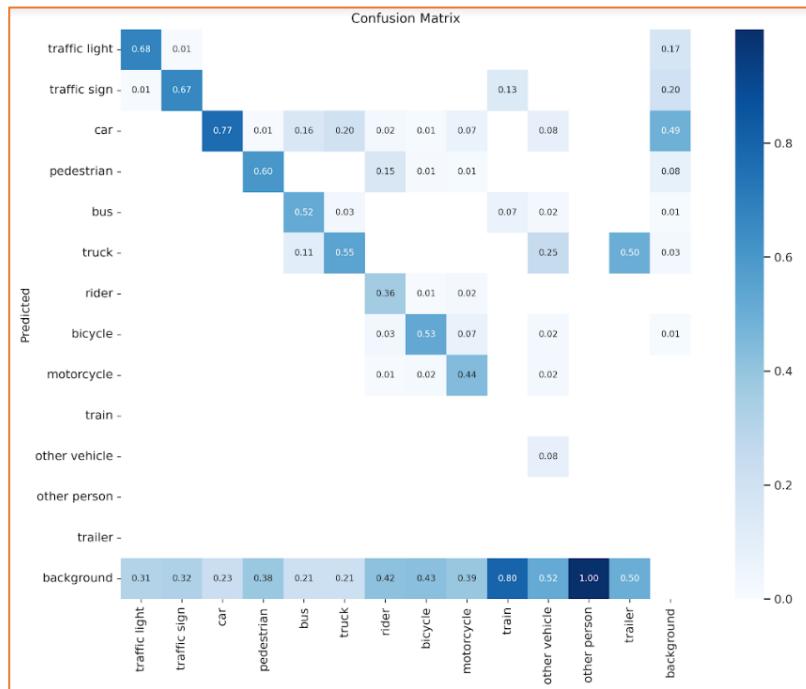
(Trained v8 nano model architecture with all BDD100k data for 20 epochs)



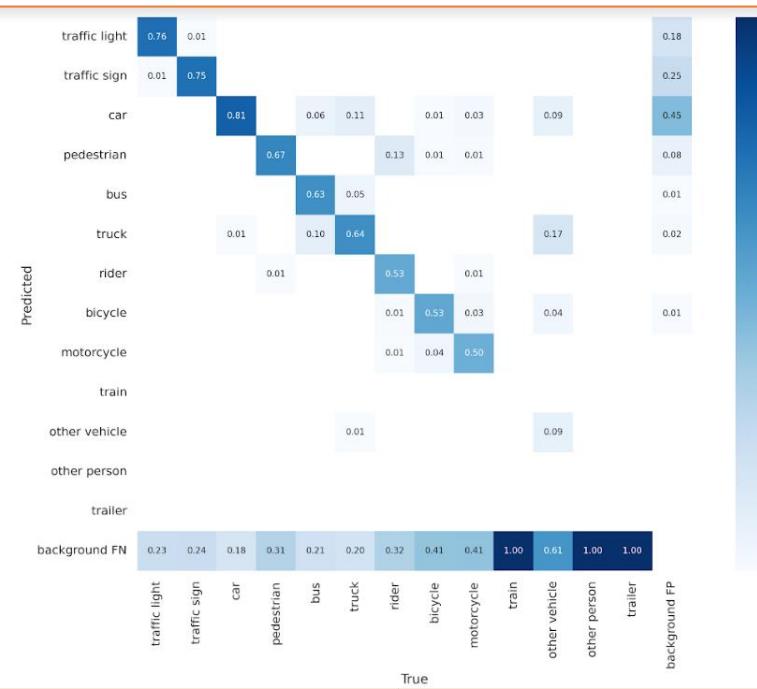
3.YOLO-V8

YOLOv5 vs YOLOv7 vs YOLOv8

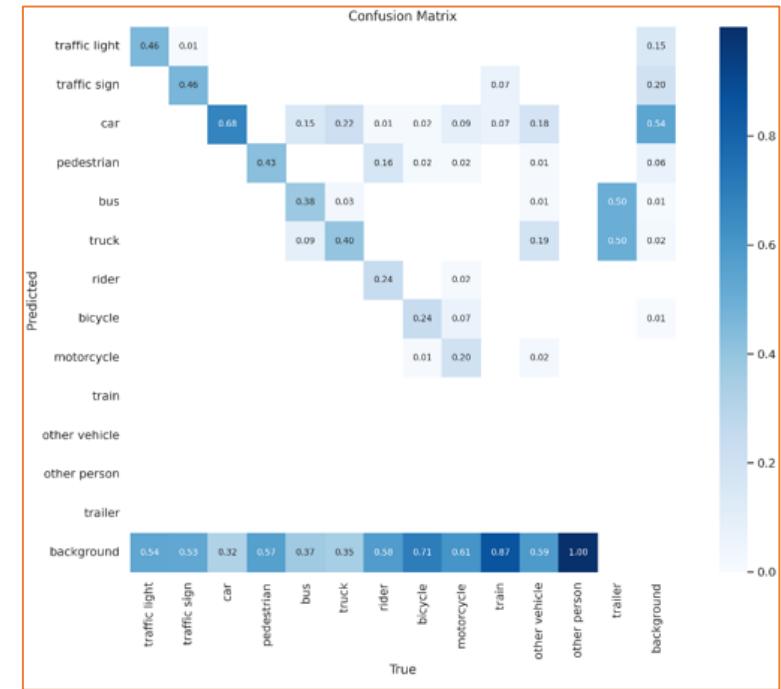
Confusion Matrix



V5



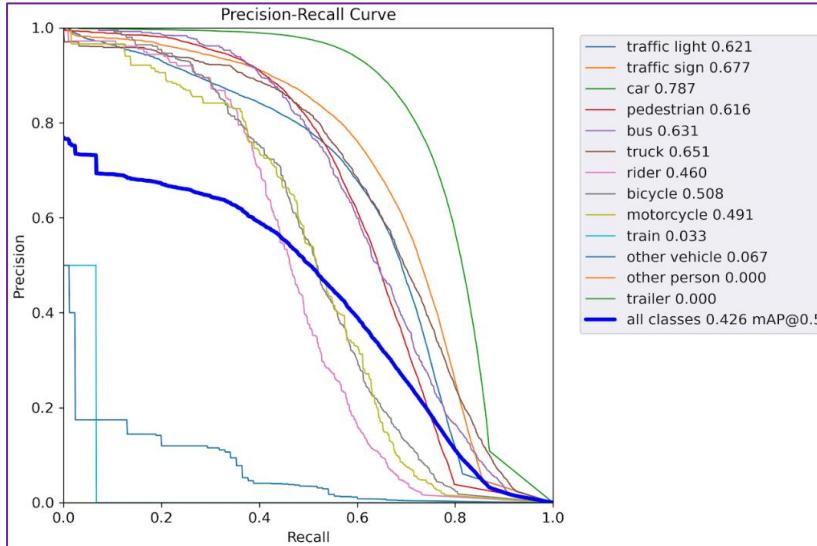
V7



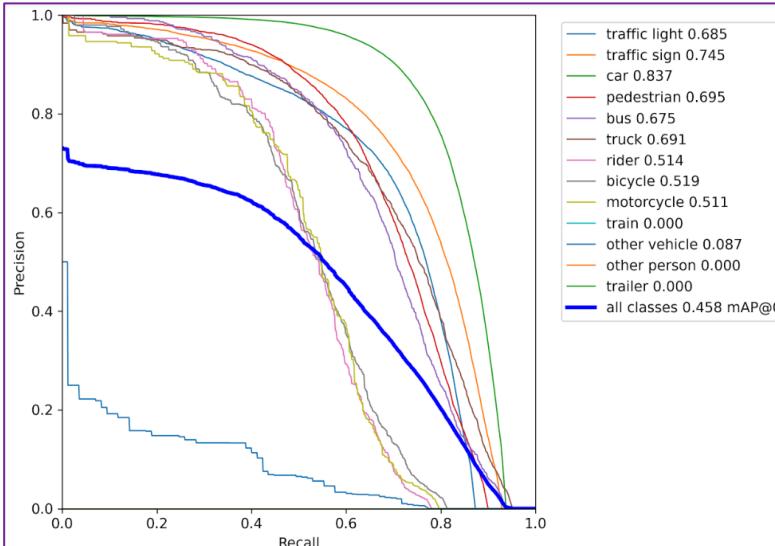
V8

3.YOLO-V8

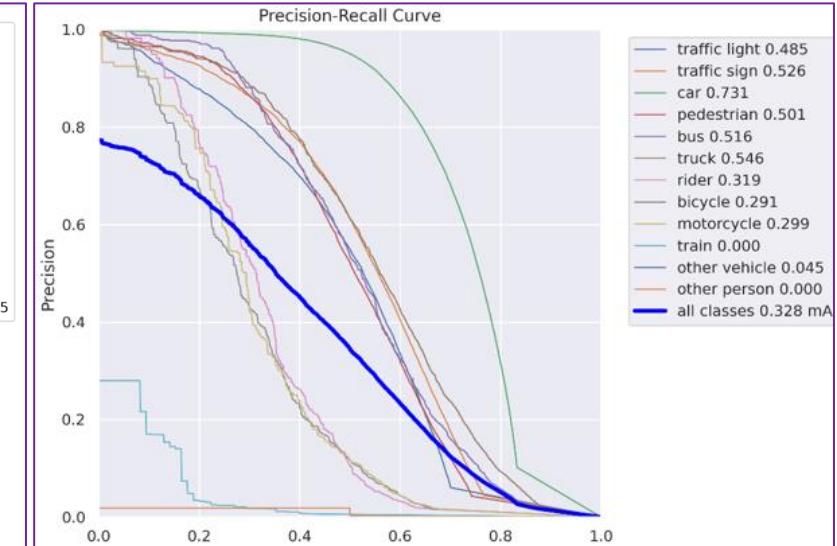
YOLOv5 vs YOLOv7 vs YOLOv8



V5



V7



V8

3.YOLO-V8

YOLOv5 vs YOLOv7 vs YOLOv8

Class	Images	Instances	P	R	mAP50	mAP50-95:
all	10000	186033	0.727	0.387	0.426	0.244
traffic light	10000	26884	0.676	0.601	0.621	0.234
traffic sign	10000	34724	0.728	0.623	0.677	0.36
car	10000	102837	0.796	0.722	0.787	0.505
pedestrian	10000	13425	0.719	0.553	0.616	0.307
bus	10000	1660	0.698	0.551	0.631	0.488
truck	10000	4243	0.691	0.594	0.651	0.477
rider	10000	658	0.733	0.388	0.46	0.231
bicycle	10000	1039	0.601	0.477	0.508	0.247
motorcycle	10000	460	0.667	0.446	0.491	0.244
train	10000	15	1	0	0.0326	0.0261
other vehicle	10000	85	0.143	0.0706	0.0672	0.0507
other person	10000	1	1	0	0	0
trailer	10000	2	1	0	0	0

V5

Class	Images	Labels	P	R	mAP@.5	mAP@.5:.95:
all	10000	186033	0.742	0.422	0.458	0.258
traffic light	10000	26884	0.708	0.668	0.685	0.258
traffic sign	10000	34724	0.727	0.701	0.745	0.396
car	10000	102837	0.833	0.76	0.837	0.522
pedestrian	10000	13425	0.756	0.617	0.695	0.344
bus	10000	1660	0.752	0.586	0.675	0.517
truck	10000	4243	0.704	0.63	0.691	0.505
rider	10000	658	0.611	0.491	0.514	0.257
bicycle	10000	1039	0.611	0.497	0.519	0.249
motorcycle	10000	460	0.73	0.457	0.511	0.251
train	10000	15	1	0	0	0
other vehicle	10000	85	0.218	0.0824	0.0873	0.0519
other person	10000	1	1	0	0	0
trailer	10000	2	1	0	0	0

V7

Class	Images	Instances	Box(P)	R	mAP50	mAP50-95:
all	10000	186033	0.565	0.303	0.328	0.186
traffic light	10000	26884	0.597	0.473	0.486	0.178
traffic sign	10000	34724	0.633	0.49	0.526	0.274
car	10000	102837	0.696	0.694	0.732	0.466
pedestrian	10000	13425	0.598	0.462	0.502	0.239
bus	10000	1660	0.579	0.48	0.516	0.398
truck	10000	4243	0.621	0.496	0.546	0.395
rider	10000	658	0.52	0.299	0.317	0.155
bicycle	10000	1039	0.48	0.281	0.292	0.143
motorcycle	10000	460	0.545	0.257	0.299	0.139
train	10000	15	0	0	0	0
other vehicle	10000	85	0.0815	0.0118	0.0447	0.0264
other person	10000	1	1	0	0	0
trailer	10000	2	1	0	0.00988	0.00593

V8

3.YOLO-V8

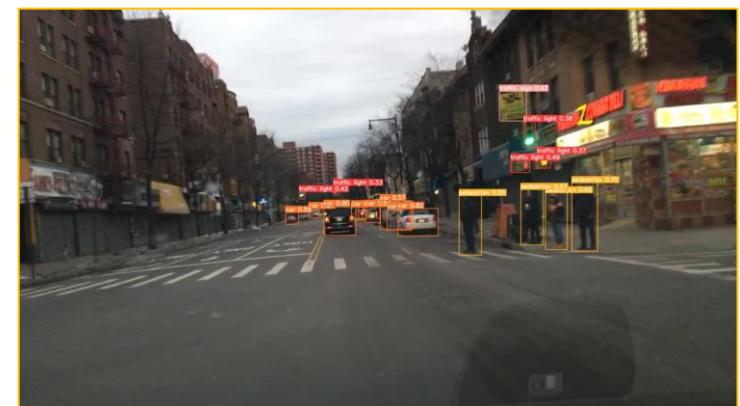
YOLOv5 vs YOLOv7 vs YOLOv8



V5



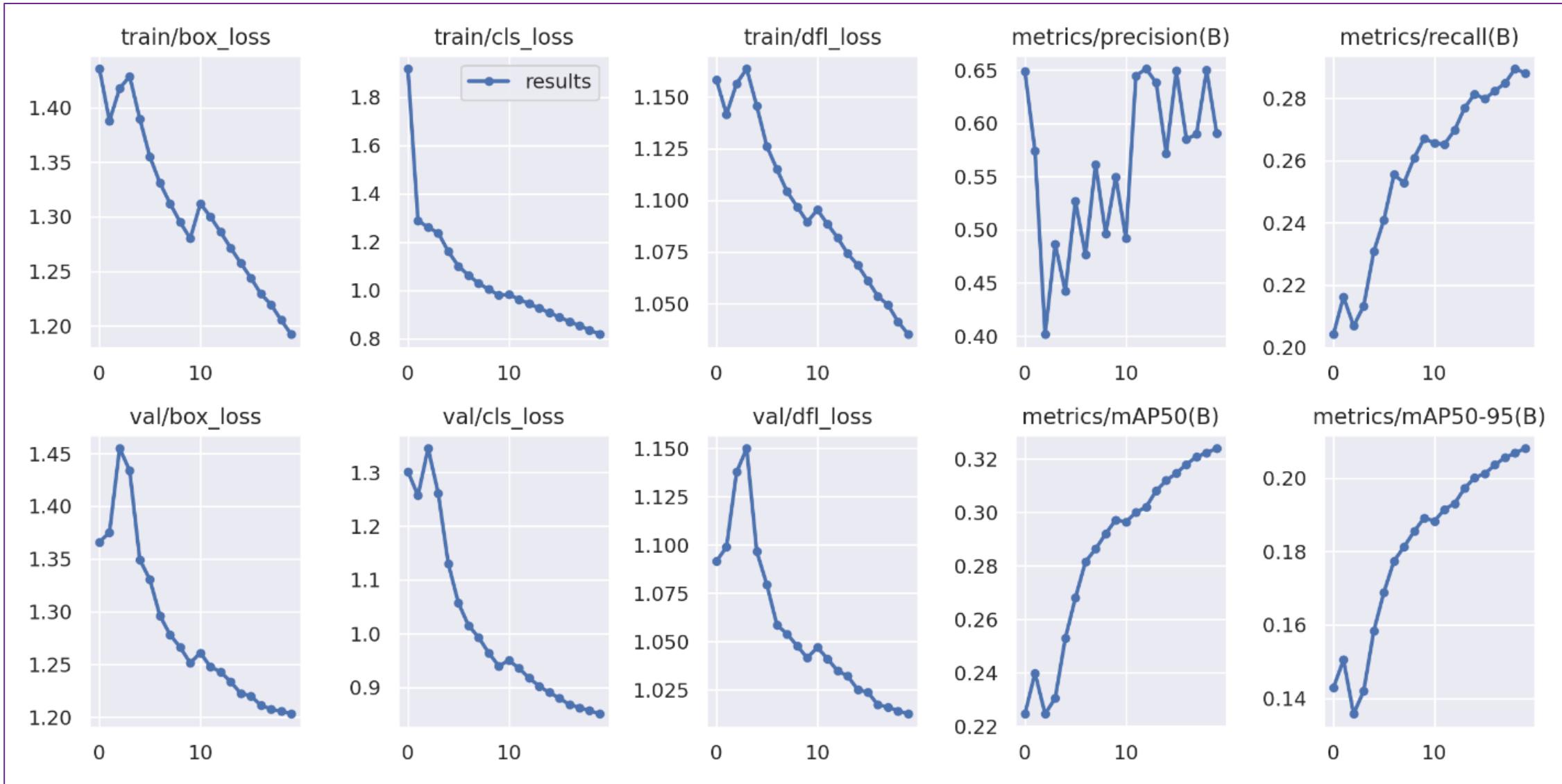
V7



V8

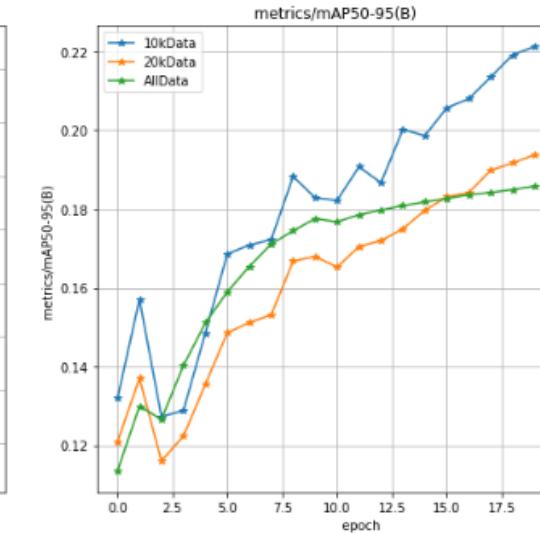
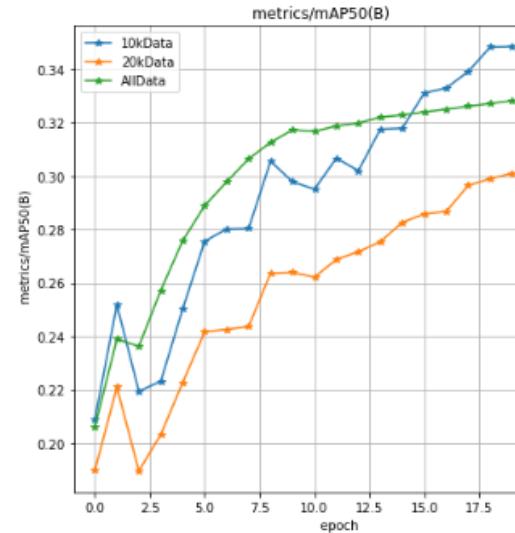
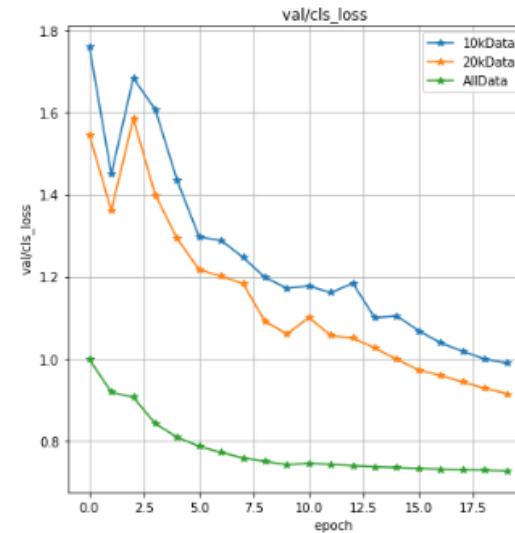
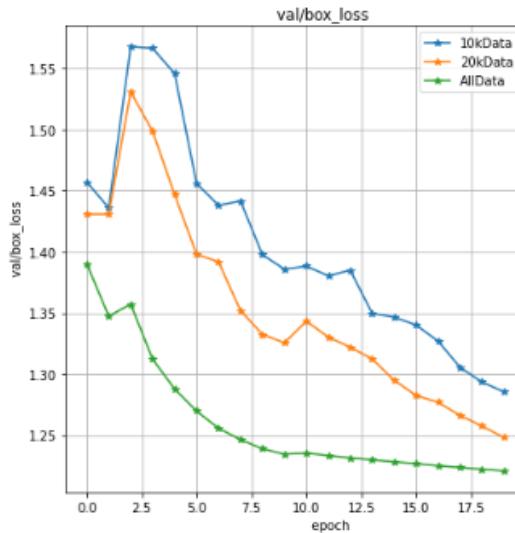
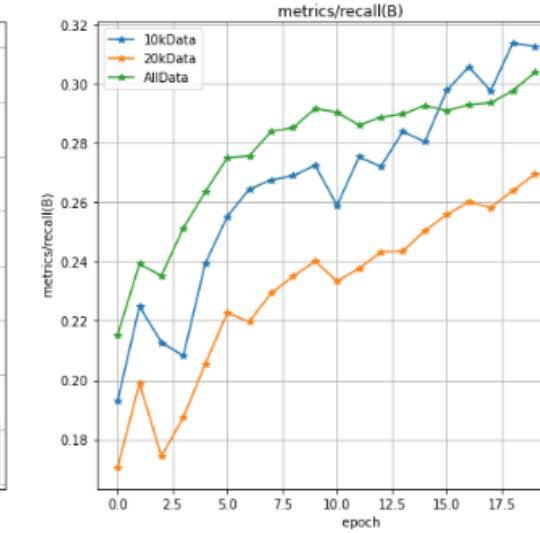
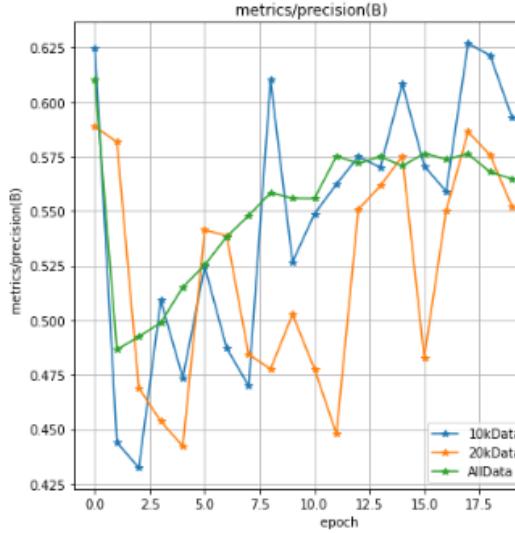
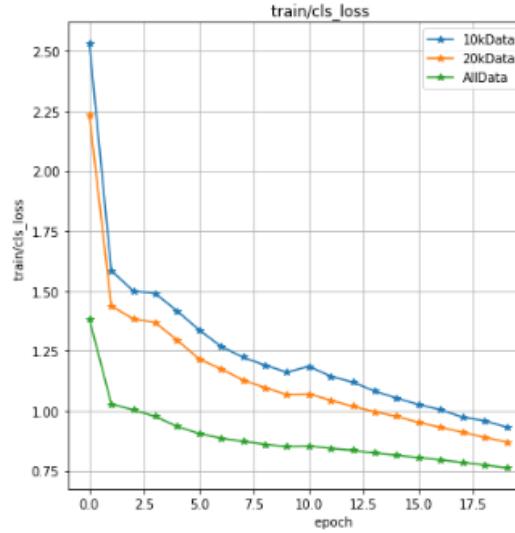
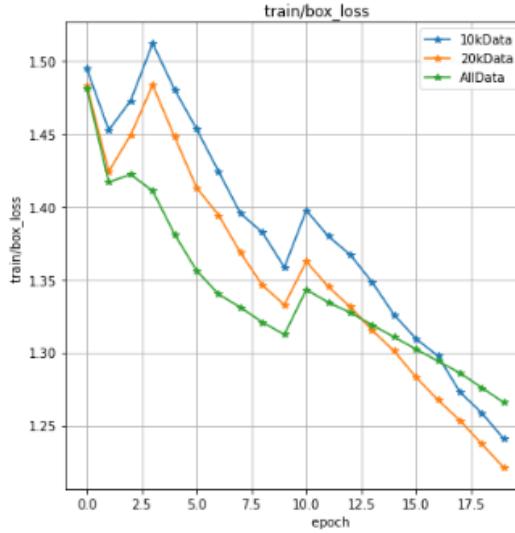
4.YOLO-V8 : IDD Dataset

Train Metrics



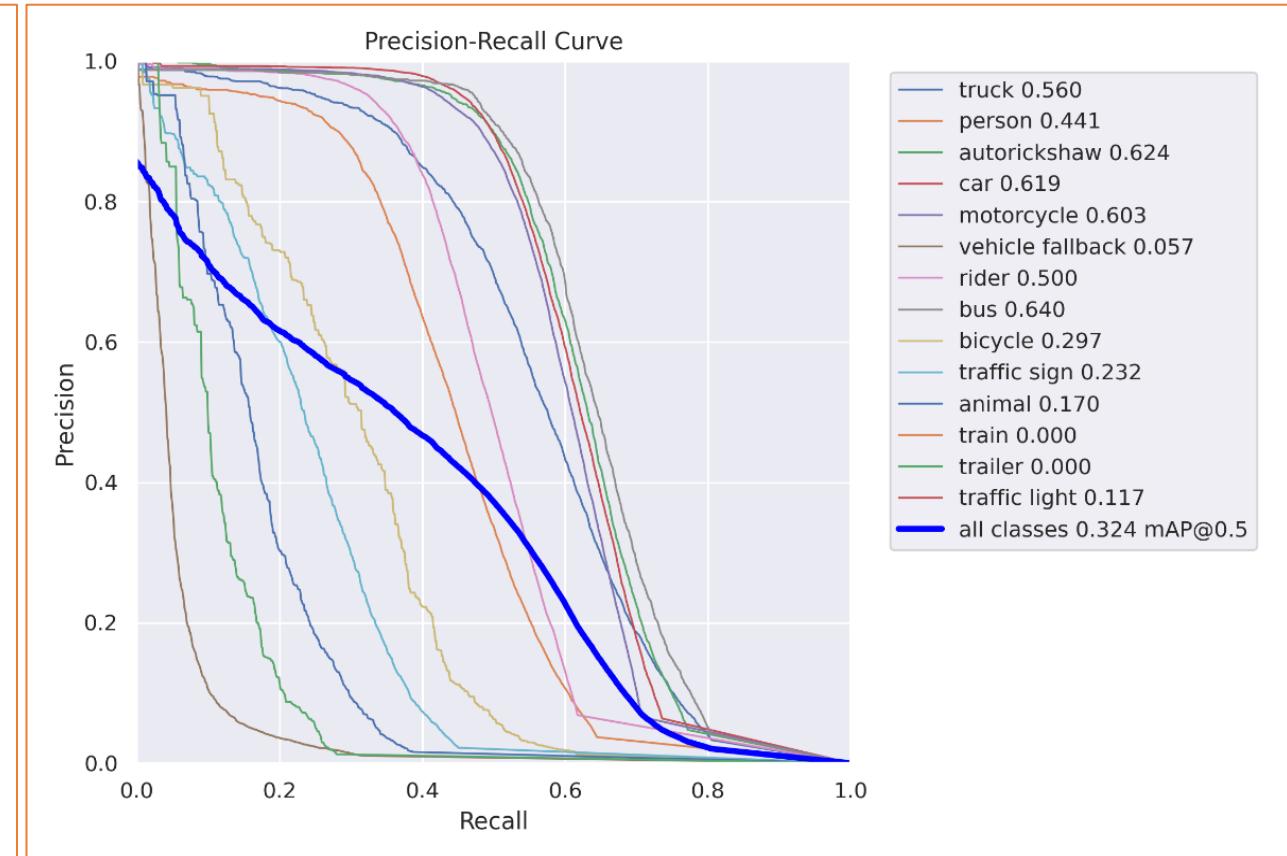
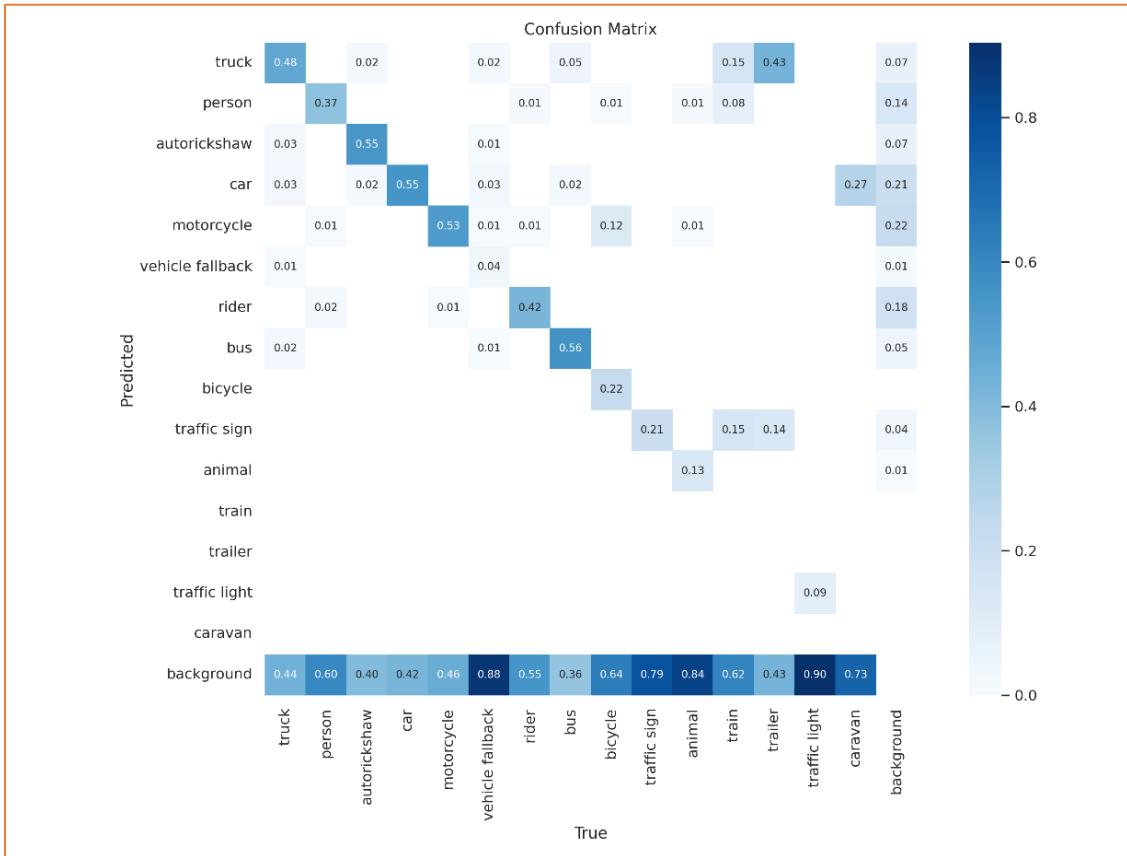
4.YOLO-V8 : IDD Dataset

Train Metrics comparison against different Data sizes



4.YOLO-V8 : IDD Dataset

YOLOv8 – Trained on IDD dataset
 [Confusion Matrix & P-R Curves]



4.YOLO-V8 : IDD Dataset

[Validation Metrics]

Class	Images	Instances	Box(P)	R	mAP50	mAP50-95):
all	10224	126004	0.59	0.288	0.324	0.208
truck	10224	7075	0.673	0.508	0.561	0.4
person	10224	18070	0.692	0.382	0.442	0.244
autorickshaw	10224	7781	0.745	0.568	0.624	0.445
car	10224	24831	0.751	0.557	0.619	0.432
motorcycle	10224	25484	0.765	0.544	0.602	0.357
vehicle fallback	10224	6078	0.457	0.0439	0.0576	0.0313
rider	10224	24510	0.736	0.434	0.5	0.286
bus	10224	4910	0.728	0.589	0.64	0.491
bicycle	10224	569	0.616	0.254	0.297	0.175
traffic sign	10224	4287	0.57	0.213	0.233	0.121
animal	10224	1460	0.588	0.138	0.171	0.0831
train	10224	13	1	0	0	0
trailer	10224	7	0	0	0	0
traffic light	10224	918	0.533	0.0937	0.116	0.0579
caravan	10224	11	0	0	0	0

4.YOLO-V8 : IDD Dataset

[Object Detection – Images]



BDD100K
dataset
training

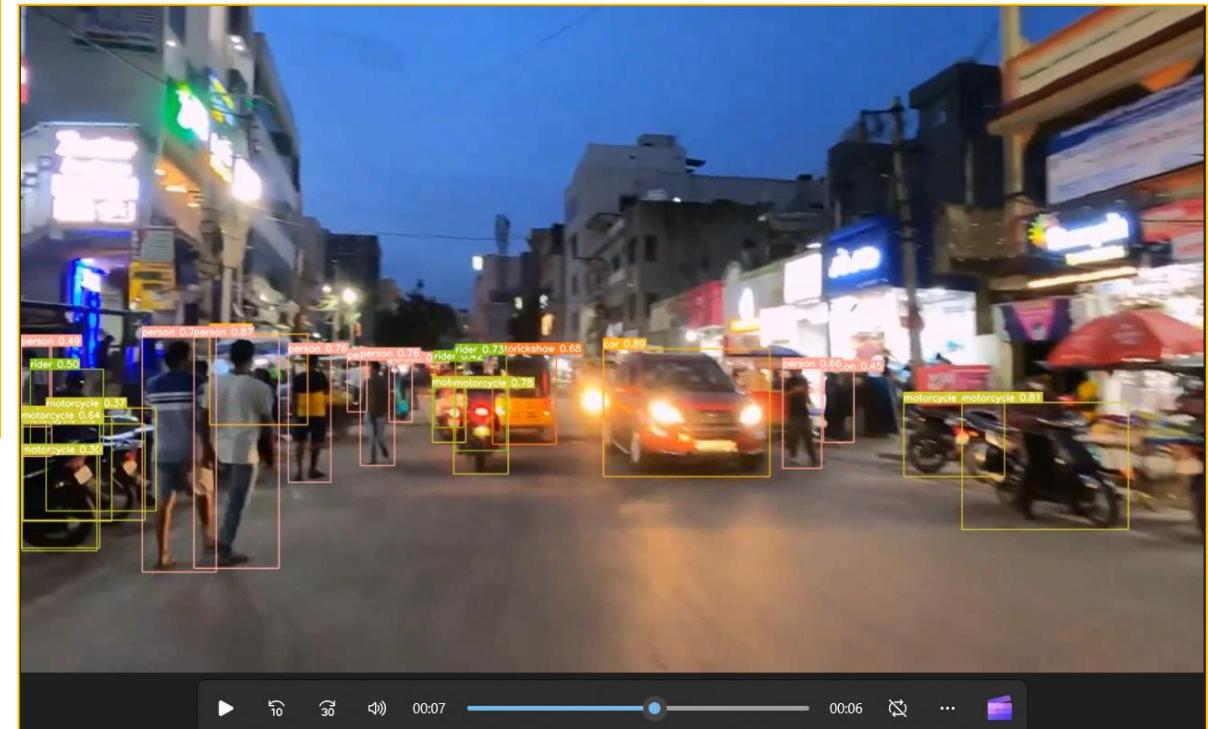


4.YOLO-V8 : IDD Dataset

[Object Detection – Videos]



training



5.YOLO-V8 : IDD Dataset – Semantic Segmentation

Preparing Labels

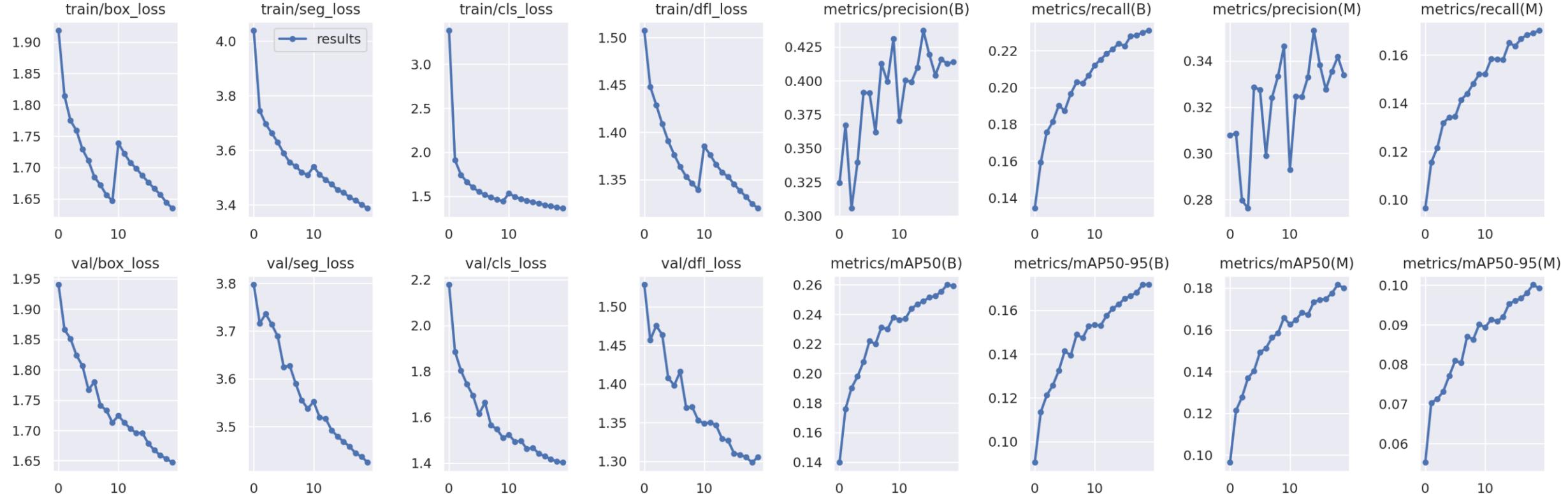


```
{  
    "imgHeight": 1080,  
    "imgWidth": 1920,  
    "objects": [  
        {  
            "date": "14-Jun-2019 17:26:53",  
            "deleted": 0,  
            "draw": true,  
            "id": 0,  
            "label": "road", |  
            "polygon": [  
                [ 1094.3204419889503,  
                  572.817679558011  
                ],  
                [ 1372.375690607735,  
                  553.7237569060774  
                ],  
                [ 1919.0,  
                  699.2307692307692  
                ],  
                [ 1919.0,  
                  708.4615384615385  
                ],  
                [ 1094.3204419889503,  
                  572.817679558011  
                ]  
            ]  
        }  
    ]  
}
```

```
0 0.5699585635359116 0.5303867403314917 0.714779005524862 0.5127071823  
1 0.6084944751381216 0.5414364640883979 0.5991712707182321 0.550276243  
1 0.0 0.7585470085470085 0.055288461538461536 0.7318376068376068 0.108  
1 0.6015625 0.5373931623931624 0.5835336538461537 0.5405982905982906 0  
1 0.7160220994475138 0.5138121546961326 0.9994791666666667 0.593370165  
2 0.2535911602209945 0.5558011049723757 0.2548342541436464 0.560220994  
3 0.0 0.24640883977900555 0.0 0.0 0.9988259668508288 0.001104972375690  
1 0.0037292817679558015 0.7624309392265194 0.04848066298342542 0.73149  
4 0.6899171270718232 0.5314917127071824 0.6843232044198896 0.531491712  
2 0.6116022099447513 0.5259668508287293 0.6140883977900553 0.529281767  
2 0.15290055248618784 0.5878453038674034 0.1522790055248619 0.58563535  
5 0.6681629834254145 0.4375690607734807 0.6681629834254145 0.432044198  
6 0.619682320441989 0.4751381215469614 0.6221685082872929 0.4751381215  
6 0.6631906077348066 0.5116022099447515 0.6607044198895029 0.511602209  
7 0.5065607734806631 0.4243093922651934 0.5065607734806631 0.365745856  
8 0.5363950276243094 0.36243093922651937 0.5525552486187846 0.36022099  
6 0.5792817679558012 0.5226519337016575 0.5780386740331491 0.519337016
```

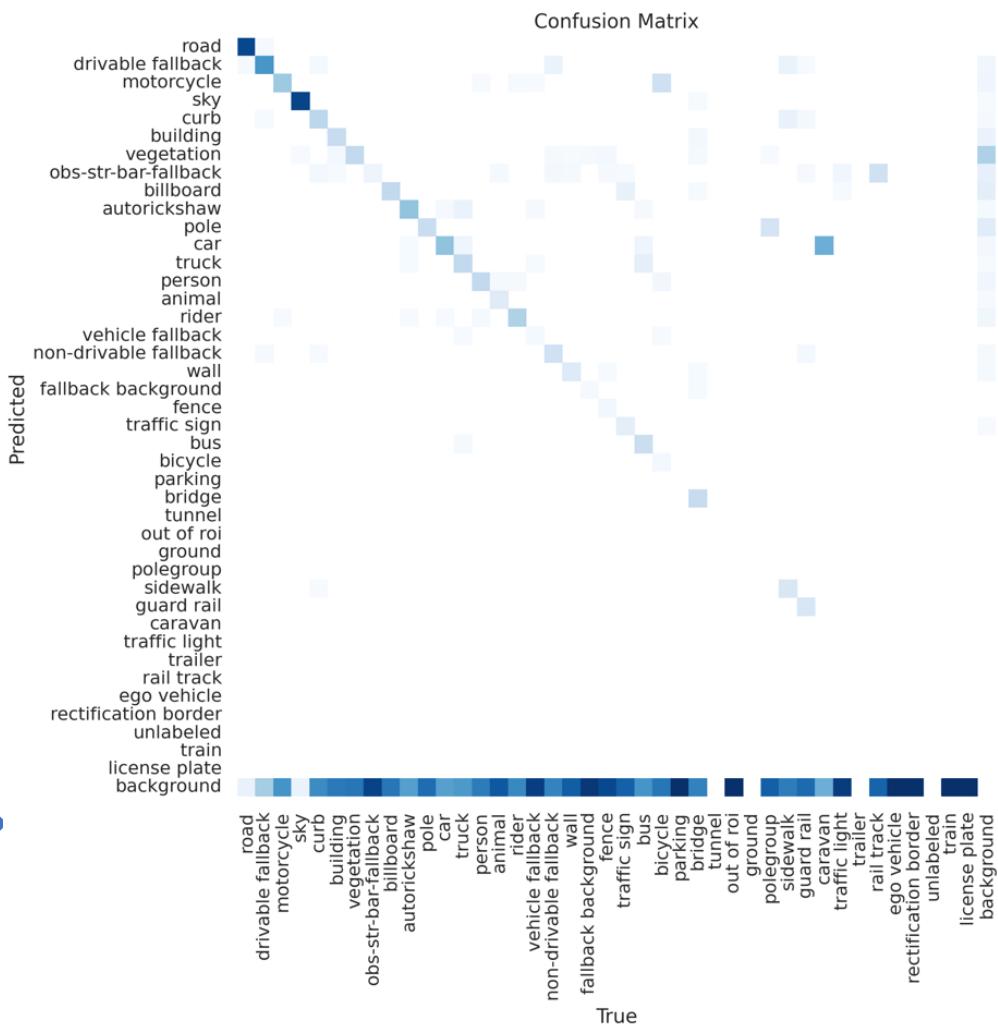
5.YOLO-V8 : IDD Dataset – Semantic Segmentation

[Train Metrics]



5.YOLO-V8 : IDD Dataset – Semantic Segmentation

[Confusion Matrix & P-R Curves]

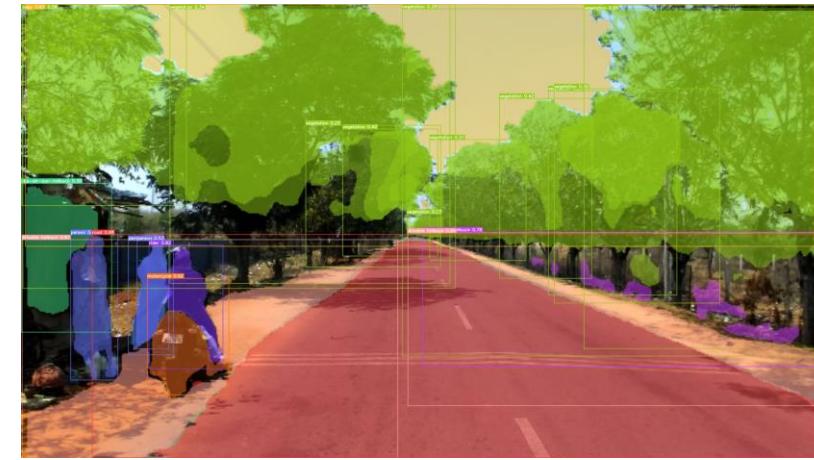


5.YOLO-V8 : IDD Dataset – Semantic Segmentation

[Validation Metrics]

5.YOLO-V8 : IDD Dataset – Semantic Segmentation

[Semantic Segmentation – Images]



5.YOLO-V8 : IDD Dataset – Semantic Segmentation

[Semantic Segmentation – Videos]

Pre-trained
COCO
dataset

After
BDD100K
dataset
training



Summary:

	YOLO-V5 [BDD100k]	YOLO-V7 [BDD100k]	YOLO-V8 [BDD100k]	YOLO-V8 [IDD Detection]	YOLO-V8 [IDD Seg]
Pretrained model	Yolov5m	yolov7	yolov8n	Yolov8n	yolov8n-seg
# Images for training (Train:Test:Val)	~70% of 100k (70:20:10)	~70% of 100k (70:20:10)	~70% of 100k (70:20:10)	~68% of 46588 (70:20:10)	~70% of 10k* (70:20:10) * -10k PNG images were not used (~20GB) due to resource constraint
# Classes	13	13	13	15	41
# Batch Size	64	32	32	32	32
Image Size	640	640	640	640	640

Summary:

	YOLO-V5 [BDD100k]	YOLO-V7 [BDD100k]	YOLO-V8 [BDD100k]	YOLO-V8 [IDD Detection]	YOLO-V8 [IDD Seg]
# epochs	20	25	20	20	20
Pre-trained dataset					
Train time per epoch	~00:27:40	~00:25:20	~00:08:00	~00:04:00	~00:12:40
<u>map@0.5</u> [box]	Train: 0.426 Valid: 0.426	Train: 0.459 Valid: 0.458	Train: 0.328 Valid: 0.328	Train: 0.324 Valid: 0.324	Train: 0.260 Valid: 0.260
<u>map@0.5</u> [mask]	NA	NA	NA	NA	Train: 0.182 Valid: 0.180

GPU used for training - NVIDIA A100-SXM4-40GB, 40536MiB)



5