

Image Analysis and Computer Vision

Professor Vincenzo Caglioti

Object Tracking and Vehicle Counting

(Full Project)

2nd Semester – 2022-2023 A.Y.

Team members

Milad Goudarzi

10765491

Bitra Rahmat Zadeh

10758773

Kiarash Rezaei

10809307



POLITECNICO
MILANO 1863

Contents

1. Introduction	3
1. 1. What is the project all about?	3
1. 2. State-of-the-art algorithms and our choice.....	3
2. Our approach	4
2. 1. Cascade classification.....	4
2. 2. YOLO.....	5
2. 3. SORT and DeepSORT.....	6
2. 4. ByteTrack.....	7
2. 5. Adaptive Histogram Equalization.....	9
3. Description of the code implementation	10
4. Conclusion	14
5. Appendix	15
5. 1. Kalman Filter.....	15
5. 2. ByteTrack: In great details.....	17

1. Introduction

1.1. What is the project all about?

Our project is about counting vehicles in a video. In first glance, it looks like detecting vehicles and counting them would lead us to what we want. However, this is true only for images and not videos. In videos, the whole process, i.e., object detection, will be repeated on each frame and therefore if our video has say 30 frames per second, each vehicle will be counted 30 times and it is obviously wrong. So, in addition to object detection, we need to have a technique to determine if we have already counted a vehicle and this could be achieved by tracking vehicles.

Object tracking is following and monitoring the movement of an object in a video or series of images like camera streams. Tracking objects could be quite a complex problem in some situations, e.g., when occlusion happens, or when the object is motion blurred, or even when the object leaves the frame completely and then comes into the picture after a few frames. Many techniques were developed to address this problem and we will have a look at some of them throughout this report.

1.2. State-of-the-art algorithms and our choice

It is always best to study what are state-of-the-art methods and get inspired by them. As of now, the best available tracker is SMILEtrack¹ according to MOT17 and MOT20 benchmarks (figure 1). SMILEtrack is proposed on Nov. 2022, and it proposes a Similarity Learning Module (SLM) motivated from the Siamese network to extract important object appearance features and a procedure to combine object motion and appearance features effectively. However, there is a really slight difference (less than 1%) between the top 5 trackers. SMILEtrack authors say a large portion of their work is borrowed from ByteTrack and YOLOv7 and it is of no surprise that ByteTrack is among top 5 available trackers and since we do not want to naively copy and paste SMILEtrack developers job, and also, since ByteTrack helped SMILEtrack to become the best tracker, it is quite rational to decide to jump into ByteTrack and YOLO details and try to combine them on our own because a newer version of YOLO was made available after the publication of SMILEtrack paper.

¹ SMILEtrack: SIMilarity LEarning for Multiple Object Tracking (<https://arxiv.org/pdf/2211.08824v2.pdf>)

Rank	Model	MOTA↑	IDF1	HOTA	Rank	Model	HOTA	MOTA↑	IDF1
1	SMILEtrack	78.2	77.5	63.4	1	SMILEtrack	65.24	81.06	80.5
2	SparseTrack	78.2	77.3	63.4	2	SparseTrack	65.1	81.0	80.1
3	BoT-SORT	77.8	77.5	63.3	3	BoT-SORT	65.0	80.5	80.2
4	ByteTrack	77.8	75.2	61.3	4	ByteTrack	63.1	80.3	77.3
5	STGT	77.5	75.2		5	StrongSORT	64.4	79.6	79.5

Figure 1. left: MOT20 ranking of trackers (<https://paperswithcode.com/sota/multi-object-tracking-on-mot20-1>). Right: MOT17 ranking of trackers (<https://paperswithcode.com/sota/multi-object-tracking-on-mot17>). Highlighted rows are the two discussed trackers.

2. Our approach

2.1. Cascade classification

Before trying to mess with ByteTrack, we tried a few other possibilities to see their shortcomings. Firstly, we tried to use a cascade classifier trained on vehicle images. The cascade classifier is composed of multiple stages, where each stage consists of several weak classifiers. The term "weak" refers to the fact that these classifiers alone may not be very accurate in detecting objects, but they are computationally inexpensive. It was not the worst thing ever, but it had lots of false negatives even on a single image as it is also seen from other experiments, let alone complicated conditions like occlusion.

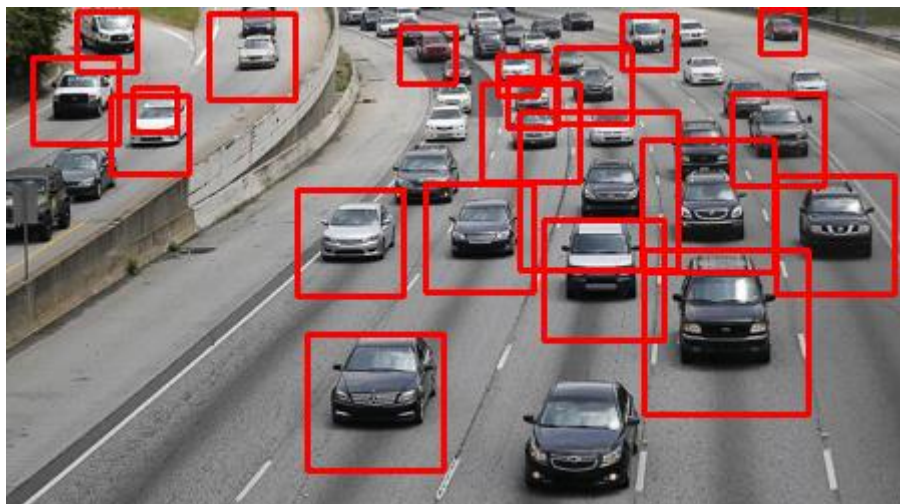


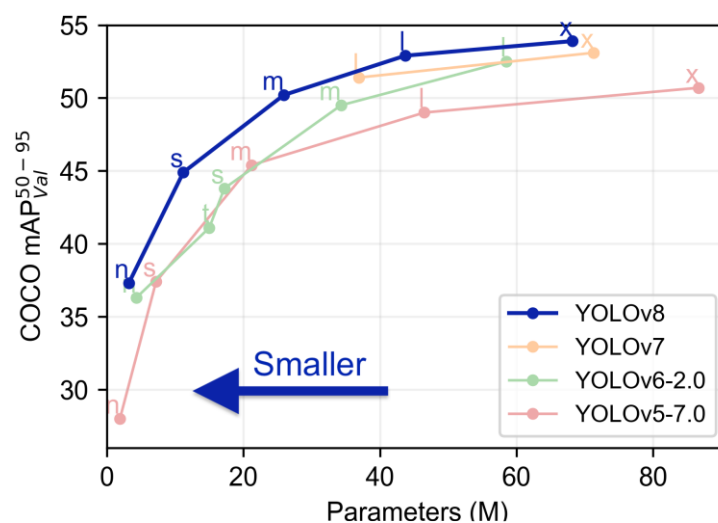
Figure 2. Cascade classifier on a single image. A lot of false negatives are present in the results.

2.2. YOLO

After doing some more research we found that YOLO² object detection model is a popular and widely used model to this end. YOLO is an acronym for “You Only Look Once” and we will see why it is called this way. YOLO is fast, has a simple architecture, and it does not use region-proposal methods for finding objects, as previous methods like R-CNN³, Faster R-CNN⁴. These methods first generate a set of region proposals or candidate object regions in the image using techniques like selective search or region proposal networks. These proposals are then classified and refined to obtain the final object detections. Region-proposal methods typically involve additional processing steps, which can increase the overall computation time. YOLO, in contrast, predicts bounding boxes (regions) and class probabilities in a single shot and hence the name “You Only Look Once”. Furthermore, it is common for object detection networks to generate multiple bounding boxes for a single object. YOLO employs Non-Maxima Suppression (NMS), which leverages the intersection over union metric, to eliminate redundant boxes. Each predicted bounding box has a confidence score and NMS removes redundant boxes by following these steps:

1. Overlapping Regions: It determines the overlap between bounding boxes.
2. Sorting: It sorts bounding boxes by confidence score.
3. Highest Confidence Box: It selects the box with the highest confidence as the reference.
4. Intersection over Union (IoU): It measures overlap using IoU.
5. Thresholding: It sets a minimum IoU threshold (e.g., 0.5).
6. Suppression: It discards overlapping boxes that exceed the threshold.
7. Iterative Process: It repeats steps 3-6 with the next highest confidence box.
8. Final Bounding Boxes: The remaining non-overlapping boxes are the final predictions.

In addition, YOLO has a few different versions and each of them has a different number of subversions. These subversions benefit from different number of parameters and hence, different accuracy and speed.



² You Only Look Once: Unified, Real-Time Object Detection (<https://arxiv.org/abs/1506.02640>)

³ Rich feature hierarchies for accurate object detection and semantic segmentation (<https://arxiv.org/abs/1311.2524v5>)

⁴ Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks (<https://arxiv.org/abs/1506.01497>)

The various versions of YOLO models are compared to one another using the mAP metric, utilizing IoU thresholds ranging from 50 to 95 with a step size of 5. mAP is the mean Average Precision which is calculated according to the following formula:

$$\text{mAP} = \frac{1}{n} \sum_{k=1}^n AP_k$$

AP_k = Average Precision of class k

n = the number of classes

Within this equation, we begin by determining distinct precisions for each class by applying various IoU thresholds. This involves categorizing a detection as a true positive when the IoU metric exceeds the specified threshold between the predicted box and the actual box. Once we have computed these precisions for a particular class, we calculate their average and repeat the process for other classes. Ultimately, we derive the mean of these Average Precisions across all classes, yielding the mAP value for the respective model.

Due to the higher mAP of YOLOv8 version x compared to other versions, we selected it as our detector. This means that we prioritized accuracy over speed, as we were not required to process videos in real-time.

2.3. SORT and DeepSORT

As our research in finding popular methods for our purpose (counting vehicles) kept going, we found out SORT⁵ algorithm did a good job when it was introduced. After SORT, an extension of it, i.e., DeepSORT⁶, was published. Both algorithms follow the "tracking by detection" paradigm, where they rely on bounding box detections provided by object detection algorithms as input, which in our case was coming from YOLOv8x. Both SORT and DeepSORT integrate Kalman filter to predict motion of objects and estimate their future positions based on past observations. The Kalman filter helps maintain smooth and consistent tracks. However, DeepSORT enhances SORT by incorporating deep appearance descriptor networks, typically based on convolutional neural networks (CNNs), to extract discriminative appearance features. These appearance features encode detailed visual information such as color, texture, and shape. According to the DeepSORT paper, it achieves better tracking accuracy and robustness in scenarios with occlusions and crowded scenes. In our case, when we implemented DeepSORT, we faced some issues with the tracking. It sometimes generated two different tracks for the same object, i.e., two tracks of the same class, e.g., "car", "truck", "motorbike", etc. This issue was mainly coming from YOLO but since the tracker benefits from appearance features, it is expected that it considers both as a single object and not two. After playing around with different parameters of YOLO and not succeeding in resolving this issue, finally, we came up with the idea of calculating the Euclidean distance between the two bounding boxes.

Euclidean distance of two vectors in 2d = $\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$

Since the two bounding boxes were close to each other, we could decide for a threshold, which rationally must be a small number because otherwise we could combine two different bounding boxes from two different vehicles into a single one and define a condition like "if the distance

⁵ Simple Online and Realtime Tracking (<https://arxiv.org/abs/1602.00763>)

⁶ Simple Online and Realtime Tracking with a Deep Association Metric (<https://arxiv.org/abs/1703.07402>)

was less than λ , keep one of them”. But this approach had a really small, yet important bug. The bug becomes prominent when two distinct vehicles come close to each other, which is usually the case with cluttered environments, and the described approach would lead to a false decision for removing one of the bounding boxes. So, we decided to change our criteria and go for a dynamic threshold and put it equal to the minimum Euclidean distance of all bounding boxes from one another. This approach helped a lot but suddenly another problem arose. In really specific scenarios, i.e., when the vehicle was leaving the frame and was partially present in the frame, YOLO would detect it to be from another class of vehicles with respect to what it was before and since the tracker keeps the history of its tracks in memory and knows that this vehicle was for instance a “car” and not a “truck”, when YOLO predicts an object to be a “truck”, the tracker would generate a new track for it, even if the bounding boxes are super close because they have different classes.



Figure 3. DeepSORT tracker generates two different tracks for the same object with even two super close bounding boxes.

At this point we decided to leave DeepSORT as its appearance descriptor did not satisfy us and we moved on to ByteTrack.

2.4. ByteTrack⁷

As mentioned earlier, ByteTrack is among the top 5 available trackers until now and had a great impact on SMILEtrack which is the best tracker available. ByteTrack is the combination of the so-called BYTE algorithm and YOLOX.

BYTE demonstrates an effective association between detections and tracks. One notable advantage of BYTE is its ability to handle low confidence detection boxes by treating them as potentially occluded or motion-blurred, or in any case hard-to-recognise objects.

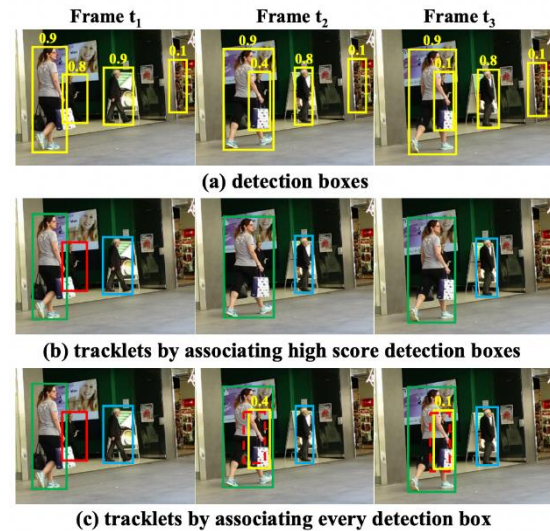
The figure on the right, which is borrowed directly from ByteTrack paper, demonstrates:

(a) In all frames, the detection model successfully identifies people in the scene but also generates false positives by detecting a portion of the background as an object.

⁷ ByteTrack: Multi-Object Tracking by Associating Every Detection Box (<https://arxiv.org/abs/2110.06864>)

(b) In the first phase, BYTE attempts to associate high-confidence detection boxes with previously initialized tracks. This association appears reasonably effective, but in frame t_2 and t_3 , it disregards the red bounding box due to its low confidence.

(c) In the second phase, represented in section c, BYTE endeavors to associate low-confidence bounding boxes with the remaining unmatched tracks. In this phase, we observe the recovery of the red object, which was previously ignored, and the elimination of the falsely detected background object.



This feature enables BYTE to track occluded vehicles in complex environments with clutter. Furthermore, BYTE has the capability to recognize the same object even if it temporarily leaves the frame and reappears later. ByteTrack incorporates the use of a Kalman filter to monitor the movement of objects across frames. This filter enables ByteTrack to predict the future location of an object based on its previous motion patterns, similar to anticipating the landing spot of a ball based on its current trajectory. As a result, even if an object is not detected in a particular frame, ByteTrack can utilize its predictions to maintain tracking of the object.

We have also observed it in our experiments. As you can see in the picture below, in the left image a car is detected and ByteTrack has assigned id number 5 to it. In the very next frame, i.e., the image in the middle, that car is occluded completely by the blue car. Again, in the next frame, which is the image on the right, that car is detected once more and is assigned the same id as before.



Figure 4. Three consecutive frames of a video being processed by YOLOv8x (for detection) and ByteTrack.

To briefly explain the Kalman filter, let's consider a simple example of tracking the position of a moving object. Suppose we have a sensor that measures the position of a ball as it moves horizontally along a straight line. However, the sensor measurements are noisy and may contain errors. The Kalman filter helps us estimate the true position of the ball by combining the noisy measurements with predictions based on the ball's previous position and velocity. Here's how it works:

1. Initialization: We start by initializing the filter with an initial estimate of the ball's position and velocity. We also define the uncertainty or noise associated with these initial estimates.

2. **Prediction:** Based on the previous estimated position and velocity, we predict the expected position of the ball in the next time step. This prediction takes into account the motion model of the ball, assuming it moves at a constant velocity.
3. **Measurement Update:** When a new measurement from the sensor arrives, we compare it with our predicted position. We calculate the difference or "residual" between the predicted and measured positions, accounting for the measurement noise.
4. **Kalman Gain:** The Kalman filter calculates a parameter called the Kalman gain, which determines the weight or influence of the measurement in updating our estimate. It takes into account the uncertainties of both the predicted state and the sensor measurements.
5. **State Update:** Using the Kalman gain, we update our estimate of the ball's position and velocity based on the measurement. This update combines the predicted state with the measurement, giving more weight to the more reliable information.
6. **Iteration:** The process of prediction, measurement update, and state update is repeated for each new measurement, iteratively improving the estimation of the ball's position as new data becomes available.

The Kalman filter considers both the predicted state and the measurement to estimate the true state of the system. It handles noisy measurements and provides a robust estimation that balances the influence of the prediction and the measurement. In our use case, the Kalman filter's job is to predict movement of vehicles, instead of a ball. The overall procedure is the same, but for the mathematical details of it, you can check the appendix.

2.5. Adaptive Histogram Equalization (AHE)

What is more, in order to help YOLO and ByteTrack, we decided to increase the contrast of frames in our video. Adaptive histogram equalization is an image enhancement technique that improves the contrast and visibility of details in an image. It achieves this by redistributing the pixel intensities in different regions of the image.

To understand adaptive histogram equalization, consider our experiment. The image contains different regions with varying levels of brightness or contrast. Some regions might appear too dark, while others might be too bright, making it difficult to see the details clearly.

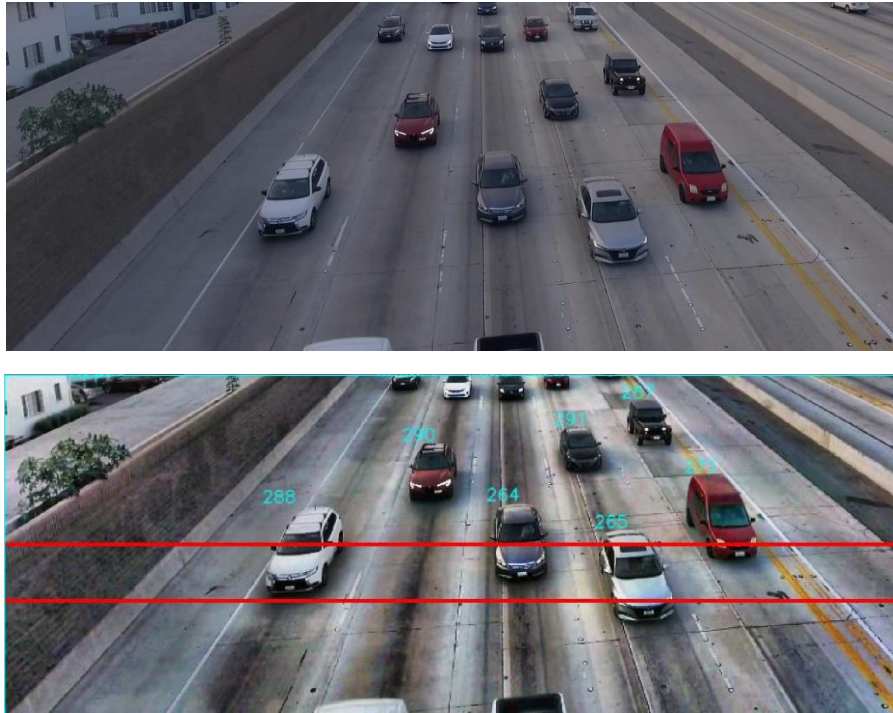


Figure 5. Before and after applying AHE. In the bottom image, the whole tracking process is carried out, and the two red lines and the number above each vehicle is their track id.

In adaptive histogram equalization, the image is divided into smaller sub-regions called "tiles" or "blocks." For each tile, a histogram of the pixel intensities is calculated, which shows the distribution of pixel values in that particular region.

Next, the histogram is equalized for each tile individually. Equalization means redistributing the pixel intensities to spread them out more evenly across the intensity range. This process helps to stretch the contrast and bring out the details in each tile.

However, one important aspect of adaptive histogram equalization is that it takes into account the local characteristics of each tile. Instead of applying a global equalization to the entire image, it adapts the equalization process based on the content of each tile.

3. Description of the code implementation

Let us explain how the code works step by step.

```

1 import sys
2 sys.path.append(".\ByteTrack")
3 from ultralytics import YOLO
4 import numpy as np
5 import cv2
6 from collections import deque
7 from math import sqrt
8 from deep_sort_pytorch_2nd.utils.parser import get_config
9 from deep_sort_pytorch_2nd.deep_sort import DeepSort
10 import torch
11 from yolox.tracker.byte_tracker import BYTETracker
12
13 # YOLO with Less parameters
14 # model = YOLO(model='yolov8n.pt')
15
16 # YOLO with Large number of parameters
17 model = YOLO('YOLOv8x.pt')
18

```

Figure 6. the code

First, we import the required libraries and add our specific system path because we are using a local version of ByteTrack, since the original one could not be run right out of the box.

```
1 class BYTETrackerArgs:
2     def __init__(self, track_thresh, track_buffer, mot20, match_thresh, \
3                 aspect_ratio_thresh, min_box_area):
4         self.track_thresh = track_thresh
5         self.track_buffer = track_buffer
6         self.mot20 = mot20
7         self.match_thresh = match_thresh
8         self.aspect_ratio_thresh = aspect_ratio_thresh
9         self.min_box_area = min_box_area
```

Figure 7. the code

Then we define a class for ByteTrack arguments, as it is designed in a way to receive arguments as an object only. These arguments are described as follows:

- `track_thresh`: Represents the threshold value for tracking. It is a value used to determine if a track should be initiated.
- `track_buffer`: Represents the buffer size for tracks. It determines how many past frames of track history should be stored and considered during the association process.
- `mot20`: A variable related to the MOT20 dataset or functionality specific to it. We do not really need it because we are not working with this dataset. However, we need to define it.
- `match_thresh`: Represents the threshold value for matching. It is a value used to determine the similarity required for matching a detection to an existing track.
- `aspect_ratio_thresh`: Represents the threshold value for the aspect ratio of a bounding box. It is used to filter out bounding boxes that have aspect ratios above a certain threshold.
- `min_box_area`: Represents the minimum area threshold for a bounding box. It is used to filter out bounding boxes that have an area smaller than the specified threshold.

```
12 def countVehicles(video_path, output_file_name, vertical, roi_xxyy=(0,0,0,0)):
13
14     assert type(video_path) == str, "video_path argument should be string"
15     assert type(output_file_name) == str, "output_file_name argument should be string"
16     assert type(vertical) == bool, "vertical argument should be boolean"
17
18     args = BYTETrackerArgs(track_thresh = 0.25,
19                           track_buffer = 30,
20                           mot20 = False,
21                           match_thresh = 0.8,
22                           aspect_ratio_thresh = 3.0,
23                           min_box_area = 1.0)
24
25     obj_tracker = BYTETracker(args)
26
```

Figure 8. the code

We have defined a pipeline for our project in order for everything to happen in the right order.

- `video_path`: It is the path to the video which is going to be processed.
- `out_put_file_name`: This is the name of the video file we are going to save our results in.
- `vertical`: It is a Boolean, indicating the approximate direction of vehicles movement. We consider different criteria for vertical and horizontal movements.

- `roi_xxyy`: roi stands for “Region Of Interest” and this argument defines the coordinates of the top left and bottom right points of the roi. It should be in xxyy format, meaning that the first two numbers are the two x-coordinates and the last two numbers are the two y-coordinates. We consider the default value of (0, 0, 0, 0) because if no special roi is passed to the pipeline, we want to use the default values.

Then, we define the arguments of the ByteTrack, which are shown to be working best in similar projects.

```

76     if not vertical:
77         areaLine1 = x_starting_point + int((x_ending_point - x_starting_point)/2) - 15
78         areaLine2 = x_starting_point + int((x_ending_point - x_starting_point)/2) + 15
79     else:
80         areaLine1 = y_ending_point - 150
81         areaLine2 = y_ending_point - 100

```

Figure 9. the code

`areaLine1` and `areaLine2` are the two red lines in different scene orientations. If vehicles are moving horizontally in the scene, we want the area lines to be in the middle with a gap of $15+15=30$ pixels. If the vehicles move vertically, we want the area lines to be placed close to the bottom of the frame with a 50 pixels gap between them.



Figure 10. The two area lines in vertical and horizontal orientations.

```

83     # apply adaptive histogram equalization (AHE) in order to increase the contrast in our region of interest.
84     clahe = cv2.createCLAHE(clipLimit=4.0, tileGridSize=(8,8))
85
86     R, G, B = cv2.split(frame[y_starting_point:y_ending_point, x_starting_point:x_ending_point]) # we don't need to
87                                                         # apply AHE
88                                                         # to the whole
89                                                         # frame
90
91     c11 = clahe.apply(R)
92     c12 = clahe.apply(G)
93     c13 = clahe.apply(B)
94
95     orig_frame = frame.copy() # we take a copy of our original frame before loosing it.
96     frame = cv2.merge((c11, c12, c13))
97     frame_h, frame_w = frame.shape[:2]
98     frame_size = np.array([frame_h, frame_w])
99     orig_frame[y_starting_point:y_ending_point, x_starting_point:x_ending_point] = frame # we replace the region of
100                                                         # interest with
101                                                         # the enhanced version of
102                                                         # it.

```

Figure 11. the code

Here, we applied the AHE, using OpenCV library for python. The point here is that since we do not need to process the whole frame, then we also do not need to apply AHE to the whole frame. That is why we defined roi. Therefore, we apply AHE on the roi only.

```

110     for result in res:
111         for box, r in zip(result.bboxes, result.bboxes.data):
112             x, y, w, h = box.xywh[0]
113             # we add x_starting_point and y_starting_point to x and y coordinate because we shrunk the frame earlier
114             x1, y1, x2, y2 = int(x) + x_starting_point - int(w/2), int(y) + y_starting_point - int(h/2), \
115                             int(x) + x_starting_point + int(w/2), int(y) + y_starting_point + int(h/2)
116
117             # if class of the detected object is not vehicle then discard it
118             if r[-1] > 0 and r[-1] < 8:
119                 xyxys.append([x1, y1, x2, y2, r[-2]]) # xyxy and score
120                 confss.append(r[-2])
121                 oids.append(r[-1]) # class of the detected object
122
123
124             if len(xyxys) > 0:
125                 tracks = obj_tracker.update(np.array(xyxys), frame_size, frame_size)
126             else:
127                 tracks = np.array([])

```

Figure 12. the code

Here, we looped through the bounding boxes and separated confidence scores and object ids and if there were any, we pass them to the tracker.

```

139     for track in tracks:
140
141         cv2.putText(orig_frame, str(track.track_id), (int(track.tlbr[0]), int(track.tlbr[1])), cv2.FONT_HERSHEY_SIMPL
142         if not vertical:
143             conditions = ((track.tlbr[0] > areaLine1 and track.tlbr[0] < areaLine2) and # upper left corner of the bb
144                         track.track_id not in ids and
145                         track.score > 0.6)
146         else:
147             conditions = ((track.tlbr[1] > areaLine1 and track.tlbr[1] < areaLine2) and # upper left corner of the bb
148                         track.track_id not in ids and
149                         track.score > 0.6)
150
151         if (conditions):
152
153             cv2.putText(orig_frame, str(track.track_id), (int(track.tlbr[0]), int(track.tlbr[1])), cv2.FONT_HERSHEY_S
154             ids.append(track.track_id)
155

```

Figure 13. the code

After having the tracks, we write the id of each bounding box above it and then we define some conditions. The first condition determines if the bounding box is inside the area lines. The second condition determines if we have counted a track before. Since this process is repeated in each frame, each track would be counted more than once. So, by this condition, we ensure they are counted only once. The last condition enforces to count the tracks with a confidence score more than a threshold. If all these conditions get satisfied, we add the track id to the ids list.

```

161     n = 20 # starting row for displaying the counted vehicle image in the original frame
162     m = 250 # starting column for displaying the counted vehicle image in the original frame
163     for track in tracks:
164         if (track.track_id in ids):
165             cv2.rectangle(orig_frame, (int(track.tlbr[0]),int(track.tlbr[1])), (int(track.tlbr[2]),int(track.tlbr[3]
166             cv2.rectangle(orig_frame, (m, n), (m+int(track.tlwh[2]), n+int(track.tlwh[3])), (0,255,0), 2)
167             try:
168                 orig_frame[n:n+int(track.tlwh[3]), m:m+int(track.tlwh[2])] = \
169                 orig_frame[int(track.tlwh[1]):int(track.tlwh[1])+int(track.tlwh[3]), int(track.tlwh[0]):int(trac
170                 m += int(track.tlwh[2])+5
171             except:
172                 print("Error!")
173
174         # drawing our RoI (Region of Interest)
175         cv2.rectangle(orig_frame, (x_starting_point, y_starting_point), (x_ending_point,y_ending_point), (255,255,0), 1)
176
177         writer.write(orig_frame)
178         cv2.imshow('frame', orig_frame)
179
180         # press esc for quitting the video
181         if cv2.waitKey(1) & 0xFF == 27:
182             break

```

Figure 14. the code

Then, we draw a green rectangle around the counted vehicles and then show the vehicle being counted above the screen, just to make it easier to see which vehicles are counted and which ones are not. In the end, we save the video file and then make the user to stop the process any time by pressing escape from keyboard.

```
1 # options for videos:
2 # vehicles moving vertically (vertical: True): los_angeles.mp4, highway.mp4
3 # vehicles moving horizontally (vertical: False): driving1.mp4
4 countVehicles('highway.mp4', 'test3.mp4', True)
```

Figure 15. the code

We can run the whole pipeline just by calling the function `countVehicles`. We have tested our implementation on three different videos, and they are all commented to let you know what the file names are, just in case someone wants to try it out. Just remember that each file should be run with the correct “vertical” parameter.

4. Conclusion

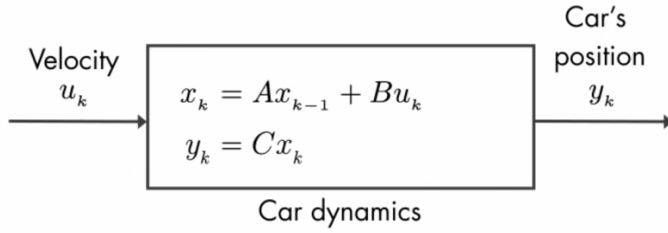
Overall, counting vehicles in a video or live stream camera could be a pretty challenging task to do due to complicated situations that might happen, such as occlusions, motion blurs, low quality camera, etc. We tried to tackle these challenges by utilizing different approaches, which in the end lead us to choose YOLOv8 and ByteTrack as our perfect combination. Although this combination works well in most situations, it has some shortcomings which we hope to be able to solve in future. For instance, ByteTrack has a hard time keeping track of small objects, meaning that it loses track when small objects are occluded even for a short time. It also experiences difficulties when there are two different cars from the same brand but with the same color and the same moving direction (moving direction is important because of Kalman Filter). In these settings, ByteTrack gets confused and there is still work to do in the mentioned situations.

5. Appendix

5.1. Kalman Filter

We can understand it best with an example. Suppose we want to design a self-driving car, which uses GPS sensor to measure its position. The car is supposed to drive 1 km on 100 different terrains.

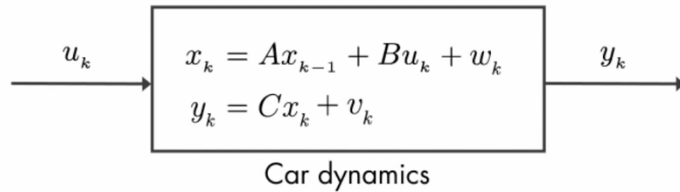
Each time, it must stop as close as possible to the finish line. Then, the average final position is computed; the goal is to have the smallest error variance and an average final position closest to 1 km. Here, we can estimate a car's position using a Kalman Filter.



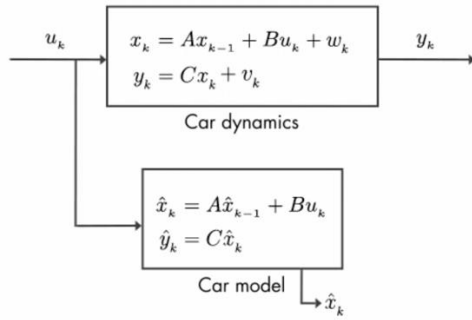
$$x_k = [position]$$

$$C = 1$$

Suppose we have the velocity of the car as an input, and we are interested to compute the position of the car. Since the GPS outputs are noisy, we won't have an accurate y . so we will show this measurement noise with V , which is a random variable. Similarly, there's process noise ω , which is also random and can represent the effect of wind or changes in the car's velocity. \mathcal{V} and ω are assumed to be from Gaussian distribution $\mathcal{V} \sim N(0, R)$ and $\omega \sim N(0, Q)$.

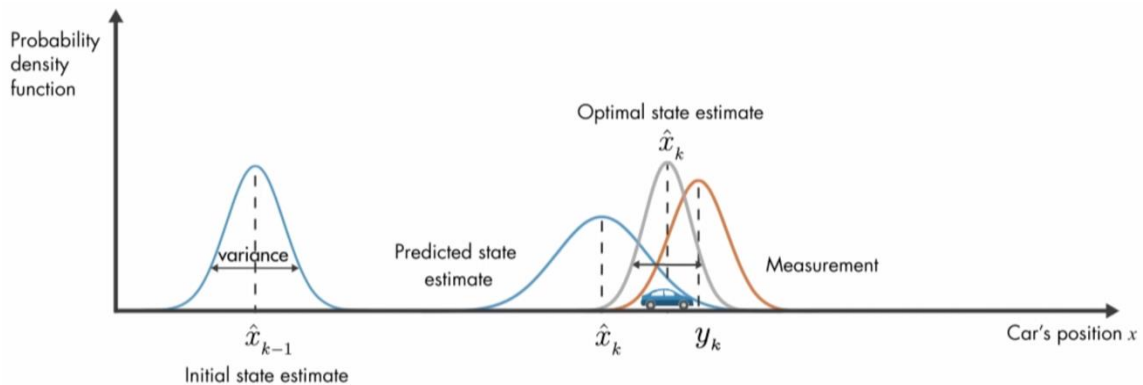


So, if we know the car evolution model, we can run the input through it to estimate the position. But this estimate also won't be perfect. Because now we're estimating x , which is uncertain due to the process noise. This is where the Kalman Filter comes into play; it combines these two pieces of information to come up with the best estimate of the car's position in the presence of process and measurement noise.



Now, let's use probability density function to discuss the working principle of the Kalman filter. At the initial time step $k-1$, the actual car position can be anywhere around the estimate \hat{x}_{k-1} . At the next step, the uncertainty in the estimate has increased, which is shown with larger variance. This is because between time step $k-1$ and k , the car might have run over a pothole, or maybe the wheels may have slipped. Therefore, it may have traveled a different distance than what we've predicted by model. Another source of information on a car's position comes from the measurement. Here, the variance is represented in the noise measurement. The optimal way to estimate the car's position is by combining these two pieces of information and this is by combining these two probability functions together.

The resulting product is also a Gaussian function, and the mean of this probability density function gives us the optimal estimate of the car's position. This is the basic idea of Kalman filters.



The Kalman filter equation being:

$$\hat{x}_k = A\hat{x}_{k-1} + Bu_k + K_k(y_k - C(A\hat{x}_{k-1} + Bu_k))$$

The part $A\hat{x}_{k-1} + Bu_k$ predicts the current state by using state estimate from previous time step and the current input, also called A Priori Estimate (\hat{x}_k^-). (predict)

Part x uses the measurement and incorporates it into the prediction to update the a-priori estimate. And we'll call the result the A Posteriori Estimate. (update)

$$\hat{x}_k = \hat{x}_k^- + K_k(y_k - C\hat{x}_k^-)$$

Prediction:

$$\hat{x}_k^- = A\hat{x}_{k-1} + Bu_k$$

$$P_k^- = AP_{k-1}A^T + Q$$

P_k^- : error covariance

Update:

$$K_k = \frac{P_k^- C^T}{CP_k^- C^T + R}$$

$$\hat{x}_k = \hat{x}_k^- + K_k(y_k - C\hat{x}_k^-)$$

$$P_k = (I - K_k C)P_k^-$$

K_k : Kalman gain

Once we've calculated the update equations, in the next time step the A Posteriori estimates are used to predict the new A priori estimates. And the algorithm repeats itself.

5.2. ByteTrack: In great details

BYTETrack is an advanced object tracking algorithm that combines state-of-the-art techniques from computer vision and deep learning to achieve robust and accurate tracking performance. It is specifically designed for multi-object tracking. It utilizes a two-stage tracking framework that incorporates both detection and association stages to handle complex tracking scenarios. It leverages a high-quality object detector (YOLOX) to detect objects in each frame as well as employing sophisticated data association methods called BYTE to link detections across frames and maintain track identities. One of the key strengths of BYTETrack is its integration with the Kalman filter for object motion modeling and prediction. These predictions are then used to guide the association process and improve the tracking accuracy. The main difference between BYTETrack and other trackers is considering the fact that low confidence detection boxes sometimes indicate the existence of objects. Thus, unlike other trackers, they will be used in the association phase.

Useful terminologies:

1- Detection boxes: (spatial extends of an object)

Represents the rectangle enclosing the object of interest. In other words, the spatial extent of an object as detected in a single frame of a video or image. It is typically obtained through an object detection algorithm or technique and provides information about the location and size of the object within the frame.

2- Track: (long temporal sequence of detection boxes)

a track refers to the temporal sequence of detection boxes across multiple frames. It represents the trajectory or path followed by the object over time.

3- Tracklet: (short temporal sequence of detection boxes)

It refers to a shorter or more localized segment of a track that could have been corrupted due to different reasons (e.g. occlusion)

4- Data Association: It is the process of computing the similarity between tracklets and detection boxes in order to determine which detection in the current frame corresponds to which object or track from previous frames. There are different similarity metrics that can be location, motion and appearance based that are context dependent.

5- Matching strategy: it is the process of ID assignment to the objects to match them to the tracklets. This can be done with different approaches such as Hungarian Algorithm.

6- Hungarian Algorithm: Is an algorithm aims to associate detection boxes with existing tracks in an optimized way.

- a) **Cost matrix:** a cost matrix is constructed that represents the pair-wise dissimilarities between all prediction boxes and tracks. It is a N-by-M matrix, where N is the number of prediction boxes and M is the number of tracks.
- b) **Assignment matrix:** Initializing the assignment matrix as a zero matrix with the same dimension to the cost matrix, indicating no assignments.
- c) **Row and Column reduction:** is simplifying the cost matrix by adjusting its values. involve subtracting the minimum value from each row and column of the cost matrix to create zeros and simplify the assignment problem. This reduction process prepares the cost matrix for the subsequent steps of the algorithm.
- d) **Zero Assignment Check:** The algorithm checks if there is a zero assignment in the cost matrix. If there is no zero assignment, the algorithm proceeds to the previous step. Otherwise, it moves to the assignment step.
- e) **Assignment:** The algorithm attempts to find an optimal assignment by iteratively selecting assignments and updating the assignment matrix. It starts by selecting a zero assignment from the cost matrix and marking it as an assigned task. Then, it examines the other assignments that involve the same

detection or track and updates the assignment matrix accordingly. This process continues until all the assignments are made.

- f) **Augmentation:** If the number of assignments made is less than the total number of detections or tracks, the algorithm enters the augmentation phase. It searches for an alternative assignment path by exploring the assigned and unassigned elements of the assignment matrix. This phase aims to improve the assignment and find the optimal solution.
 - g) **Final Assignment:** Once the augmentation phase is complete, the algorithm generates the final assignment based on the updated assignment matrix. The assignments in the matrix represent the optimal associations between the detections and tracks, minimizing the total cost.
- 7- **IOU (Intersection over Union):** In the context of object tracking is a similarity metric that is calculated by dividing the area of intersection between two prediction boxes by the area of their union. It provides a value between 0 and 1, where 0 indicates no overlap, and 1 indicates complete overlap or identical regions.
- 8- **Re-ID (Re-Identification):** In the context of object tracking, it refers to the metric used to measure the visual similarity between objects or individuals for the purpose of re-identifying them across frames or scenes. It involves comparing visual features or embeddings to determine the degree of similarity.
- 9- **Track rebirth:** the process of creating new tracks for objects that appear in a video sequence for the first time or after being absent for a certain duration.

STEPS OF BYTE ALGORITHM:

- 1- **Predict the detection boxes:** The bounding box corresponding to each object is predicted as well as their score using an object detector model. Then we separate them into two groups (D_high and D_low) based on comparing each detection box score and the detection score threshold.
- 2- **Kalman Filter:** We estimate and predict the state of tracked objects which is represented as a vector. Generally, using the location of the object in the previous frame (previous state of the object) we predict the new location (current state) in the current frame for each track based on the predefined motion model (evolution model).
- 3- **First association and matching:** it is performed between the high score detection boxes(D_high) and all the tracks (T). This association can be done based on either computing IoU or Re-ID distances between the detection boxes(D_high) and the predicted box of tracks(T). Then the distance matrix will be used as the cost matrix in

Hungarian algorithm in order to finish the matching. It should be noted that unmatched detections from D_{high} and unmatched tracks from T will be stored (in D_{remain} and T_{remain} , respectively.)

- 4- **Second association:** it is performed between low score detection boxes(D_{low}) and the remaining tracks from the previous association (T_{remain}). After association, the unmatched low score detection boxes will be deleted but the remaining tracks from this step will be stored in $T_{re-remain}$. As the similarity metric for this association, it is suggested to use IoU. Since the appearance features of low score detections are likely to be unreliable (severe occlusion, motion blur and, etc) Re-ID cannot perform well. Regarding the unmatched tracks ($T_{re-remain}$) we put them in T_{lost} . If each track in T_{lost} exists for more than a certain number of frames (e.g 30) we remove it from the list of tracks (T).
- 5- **New track initialization:** For the unmatched high score detection boxes (D_{high}) new tracks are initialized.

Algorithm 1: Pseudo-code of BYTE.

```

Input: A video sequence  $V$ ; object detector  $Det$ ; detection score threshold  $\tau$ 
Output: Tracks  $\mathcal{T}$  of the video
1 Initialization:  $\mathcal{T} \leftarrow \emptyset$ 
2 for frame  $f_k$  in  $V$  do
  /* Figure 2(a) */
  /* predict detection boxes & scores */
3    $\mathcal{D}_k \leftarrow Det(f_k)$ 
4    $\mathcal{D}_{high} \leftarrow \emptyset$ 
5    $\mathcal{D}_{low} \leftarrow \emptyset$ 
6   for  $d$  in  $\mathcal{D}_k$  do
7     if  $d.score > \tau$  then
8        $\mathcal{D}_{high} \leftarrow \mathcal{D}_{high} \cup \{d\}$ 
9     end
10    else
11       $\mathcal{D}_{low} \leftarrow \mathcal{D}_{low} \cup \{d\}$ 
12    end
13  end
  /* predict new locations of tracks */
14  for  $t$  in  $\mathcal{T}$  do
15     $t \leftarrow KalmanFilter(t)$ 
16  end

  /* Figure 2(b) */
  /* first association */
17  Associate  $\mathcal{T}$  and  $\mathcal{D}_{high}$  using Similarity#1
18   $\mathcal{D}_{remain} \leftarrow$  remaining object boxes from  $\mathcal{D}_{high}$ 
19   $\mathcal{T}_{remain} \leftarrow$  remaining tracks from  $\mathcal{T}$ 

  /* Figure 2(c) */
  /* second association */
20  Associate  $\mathcal{T}_{remain}$  and  $\mathcal{D}_{low}$  using similarity#2
21   $\mathcal{T}_{re-remain} \leftarrow$  remaining tracks from  $\mathcal{T}_{remain}$ 

  /* delete unmatched tracks */
22   $\mathcal{T} \leftarrow \mathcal{T} \setminus \mathcal{T}_{re-remain}$ 

  /* initialize new tracks */
23  for  $d$  in  $\mathcal{D}_{remain}$  do
24     $\mathcal{T} \leftarrow \mathcal{T} \cup \{d\}$ 
25  end
26 end
27 Return:  $\mathcal{T}$ 

```

Figure 16. BYTE algorithm