

Multiple Object Tracking in Urban Traffic Scenes with a Multiclass Object Detector

Hui-Lee Ooi, Guillaume-Alexandre Bilodeau, Nicolas Saunier, and
David-Alexandre Beaupré

Polytechnique Montréal, Canada
<hui-lee.ooi, gabilodeau, nicolas.saunier,
david-alexandre.beaupre>@polymtl.ca

Abstract. Multiple object tracking (MOT) in urban traffic aims to produce the trajectories of the different road users that move across the field of view with different directions and speeds and that can have varying appearances and sizes. Occlusions and interactions among the different objects are expected and common due to the nature of urban road traffic. In this work, a tracking framework employing classification label information from a deep learning detection approach is used for associating the different objects, in addition to object position and appearances. We want to investigate the performance of a modern multiclass object detector for the MOT task in traffic scenes. Results show that the object labels improve tracking performance, but that the output of object detectors are not always reliable.

Keywords: Multiple object tracking · road user detection · urban traffic.

1 Introduction

The objective of multiple object tracking (MOT) is extracting the trajectories of the different objects of interest in the scene (camera field of view). It is a common computer vision problem that is still open in complex applications. This paper deals with one of these complex applications, urban traffic, that involves different kinds of road users such as drivers of motorized and non-motorized vehicles, and pedestrians (see Figure 1). The various road users exhibit different properties of moving speeds and directions in the urban environment. Their size vary because of perspective. Besides, road users are frequently interacting and occluding each other, which makes it even more challenging.

In this work, we want to investigate the performance of a modern multiclass object detector [2] for the MOT task in traffic scenes. We are interested in testing MOT in urban traffic settings with road users of varying sizes using an object detector while most previous works in such applications employ background subtraction or optical flow to extract the objects of interest regardless of their size. Our contributions in this work is an assessment of a typical model object detector for tracking in urban traffic scenes, and the introduction of label

information for describing the objects in the scenes. Due to the variability of objects found in urban scenes, the label information should be a useful indicator to distinguish and associate the objects of interests across frames, thereby producing a more accurate trajectory. In this paper, the improvements obtained thanks to classification labels are evaluated with respect to a baseline tracker that uses a Kalman filter, bounding box positions and color information.

The results show that using classification labels from a detector improves significantly tracking performances on an urban traffic dataset. Therefore, multiple object trackers should capitalize on this information when it is available. However, they also show that the outputs of a multiclass object detector are not always reliable and not always easy to interpret.



Fig. 1: A frame from the urban traffic dataset that shows several road users in an intersection.

2 Related Works

MOT in urban traffic scenes was previously studied in [3], where the use of background subtraction is proposed for detecting the objects of interest followed by updating the object model with a state machine that uses feature points and spatial information. In fact, most previous work in MOT uses background subtraction or optical flow to detect the objects. The reason is that historically, methods based on pre-trained bounding box detectors are difficult to apply to road user tracking scenarios because it is difficult to design a detector that can detect and classify every possible type of road user from various viewpoints. However, recent progress in deep learning [2,21] make this avenue now possible and worth investigating.

When using background subtraction, the detection results give blobs that can correspond to parts of objects, one object, or many objects grouped together. The task is then to distinguish between merging, fragmentation, and splitting of objects. This is the main drawback of this method, since under congested traffic

conditions, road users may partially occlude each other and therefore be merged into a single blob. Examples of trackers based on background subtraction include the work of Fuentes and Velastin [18], Torabi et al. [16], Jun et al. [13], Kim et al. [12], Mendes et al. [10], and Jodoin et al. [20]. For data association, they typically use the overlap of foreground blobs between two frames or a graph-based model for data association using appearance information, such as textures, color or keypoints. These approaches track objects in a merge-split manner as objects are tracked as groups during occlusion. The Hungarian algorithm is a classical graph-based choice for solving object assignment problems. To compensate for the missing detections, the Kalman filter is a popular option for estimating the location of the object of interest. A basic implementation of multiple object tracking is proposed in [4] using this approach.

With optical flow, objects are detected by studying the motion of tracked points in a video. Feature points that are moving together belongs to the same object. Several methods accomplish this process using the Kanade-Lucas-Tomasi (KLT) tracker [15]. The following researchers have proposed such trackers, often called feature-based: Beymer et al. [17], Coifman et al. [14], Saunier et al. [11] and Aslani and Mahdavi-Nasab [19]. For example, the algorithm proposed by Saunier et al. [11], named Traffic Intelligence, tracks road users at urban intersections by continuously detecting new features. The main issue is to select the right parameters to segment objects moving at similar speeds, while at the same time not oversegmenting smaller non-rigid objects such as pedestrians. Because objects are identified only by their motion, nearby road users moving at similar speed are often merged together. The exact bounding box occupied by the road user is unknown because it depends on the position of sparse feature points. Furthermore, when an object stops, its features flow becomes zero and feature trajectories are interrupted, which leads to fragmented object trajectories. Using a deep learning-based detector on road users is expected to provide objects that are less fragmented and that can be tracked whether they are moving or not.

3 Method

The proposed method consists of two main components: object detection and data association. It is illustrated in Algorithm 1. Object detection involves the extraction of objects of interest from the frames for further processing. Data association determines the tracking architecture to ensure the formation of the trajectories of each object in the scene. In order to match the objects correctly, an assignment cost based on a measure of similarity is computed for all the potential matches.

3.1 Object Detection

The road users from each frame are detected by using a deep-learning object detection model from the Region-based Fully Convolutional Network (RFCN) [2] framework due to its efficiency and accuracy. This detector was selected because

it was the best performing approach on the MIO-TCD localization challenge [1]. The pre-trained model is further refined by using the MIO-TCD dataset [1] to provide the labels of the different road users found in traffic scenes, belonging to one of the eleven categories or labels: articulated truck, bicycle, bus, car, motorcycle, motorized vehicle, non-motorized vehicle, pedestrian, pickup truck, single unit truck and work van.

A non-maximal suppression (NMS) method [7,8] is applied to reduce the redundant detections of the same road users in each frame.

3.2 Data Association

The object assignment or data association is essentially performed on a set of detected objects from the current frame and a list of actively tracked objects that are accumulated from previous frames.

For the matched pairings, the latest position of the corresponding object in the track list is updated from the detected object. In the case of new detection, a new object will be initialized and added to the track list. In the case of objects in the track list without a matched candidate from the detection list, i.e. a missing detection, a Kalman filter [9] is applied to predict its subsequent location in the scene and the track information is updated using the prediction.

For the matching of objects across frames, if the total cost of assigning object pairs is higher than a set threshold T_{match} , the paired object would be reassigned to unmatched detection and unmatched track respectively due to the high probability of them not being a good match.

Actively tracked objects that are not assigned a corresponding object from the new detections after $N_{timeout}$ frames are removed from the list, under the assumption that the object has left the scene or the object was an anomaly from the detection module.

Object Assignment Cost Once objects are detected, the subsequent step is to link the correct objects by using sufficient information about the objects to compute the cost of matching the objects. The Hungarian algorithm [5] is applied to match the list of active objects with the list of new detections in the current frame so that the matchings are exclusive and unique. The bipartite matching indicates that each active object can only be paired with one other candidate object (the detection) from the current frame. The algorithm can make use of different costs of assignment, with higher costs given to objects that are likely to be different road users.

Label Cost In order to describe the properties of the detected objects, the labels and corresponding confidence score from the detections are taken into account. Setting the range of scores between 0 and 1, object pairs across frames that are more similar will be given a lower cost. Using the classification labels, object pairs with different labels are less likely to be the correct matchings, therefore they will be given cost of 1. Meanwhile, when the pairing labels are the same,

Algorithm 1 MOT algorithm

```

1: procedure MOT
2:   for  $i^{th}$  frame do
3:     Extract detections with multiclass object detector
4:     if  $i == 1$  then
5:       Assign all detections as tracks
6:     else
7:       for each detection do
8:         Compute cost of detection with respect to each track
9:         Run Hungarian algorithm for assigning pairing of detection and track
10:      for each matched detection do
11:        if  $Cost > T_{match}$  then
12:          Reassign as unmatched detection and unmatched track
13:        else
14:          Update the track information from the detection
15:      for each unmatched detection do
16:        Initialize as new track
17:      for each unmatched track do
18:        if  $N > N_{timeout}$  then
19:          Remove track
20:        else
21:          Update track information using prediction from Kalman filter

```

the average of the confidence score of each detection are being taken as the label cost. The label cost is defined as

$$C_{label} = \begin{cases} 1 - 0.5 \times (Conf_i + Conf_j) & \text{if } L_i = L_j \\ 1 & \text{if } L_i \neq L_j \end{cases} \quad (1)$$

where L_n denotes the label of detection n and $Conf_n$ denotes the confidence of the corresponding label of the n^{th} detection.

Jaccard Distance-based Position Cost The bounding box coordinates of the detected objects are a useful indicator for matching the objects across frames as well. To judge the similarity of two bounding boxes in terms of proximity and size, the Jaccard distance is computed from the coordinates of the paired object, where the ratio of intersection over union of the bounding boxes is computed. This is calculated using

$$C_{position} = 1 - \frac{|Box_i \cap Box_j|}{|Box_i \cup Box_j|} \quad (2)$$

where Box_n denotes the set of pixels of the bounding box of the detected object n .

Color Cost The visual appearance of the objects is characterized by their color histograms that are used to compute the color cost. In this work, the Bhat-

tacharyya distance is applied to compute the distance of the color histogram of detections across frames with

$$C_{color} = \sqrt{1 - \frac{1}{\sqrt{H_i H_j} N^2} \sum \sqrt{H_i H_j}} \quad (3)$$

where H_i denotes the color histogram of detection i , H_j denotes the color histogram of detection j and N is the total number of histogram bins.

4 Results and Discussion

To test the RFCN multiclass object detector in MOT and to assess the usefulness of the classification labels, we used the Urban Tracker dataset [3] since it contains a variety of road users in an urban environment. Fig 2 shows some sample frames from the Urban Tracker dataset with RFCN detections. The MOT performance is evaluated by using the CLEAR MOT metrics [6]:

- multiple object tracking accuracy (MOTA) that evaluates the quality of the tracking, if all road users are correctly detected in the frames they are visible and if there are no false alarms;
- and multiple object tracking precision (MOTP) that evaluates the quality of the localization of the matches.

To test the contribution of using labels in MOT, the proposed baseline method is applied with and without object classification labels in the cost computation for data association. The following parameters are used in the experiments: T_{match} is set at 1.5 and the value of $N_{timeout}$ is set at 5.

Table 1 summarizes the results obtained with the baseline tracker. First of all, we were not able to obtain interesting results on the René-Lévesque video. From the evaluation, it is observed that the size of the objects greatly influences the performance of the proposed method because of the shortcomings of RFCN. When the size of the road users is too small, there are not enough details for the detector to distinguish the different types of objects reliably. Mis-detections are common in such cases, as observed in video René-Lévesque, for example in Figure 3. Since the frames are captured at a higher altitude than the other urban scenes, the object detector has difficulties in detecting and classifying the objects clearly due to the lack of details. On the other hand, larger objects such as buildings have the tendency of being detected as they share similarities with the features learned by the detector.

Secondly, it can be noticed from Table 1 that the MOTA results are negative and disappointing. This comes from the difficulty of interpreting the detections of RFCN. The same object is sometimes detected as several instances from the object detection module, as shown in Figure 4. This often causes confusion and unnecessarily increases the number of detected objects and degrades the reported tracking performance. When there are no consecutive redundant detections, these redundant instances of the same object will usually be removed after a few frames since the object assignments are exclusive.

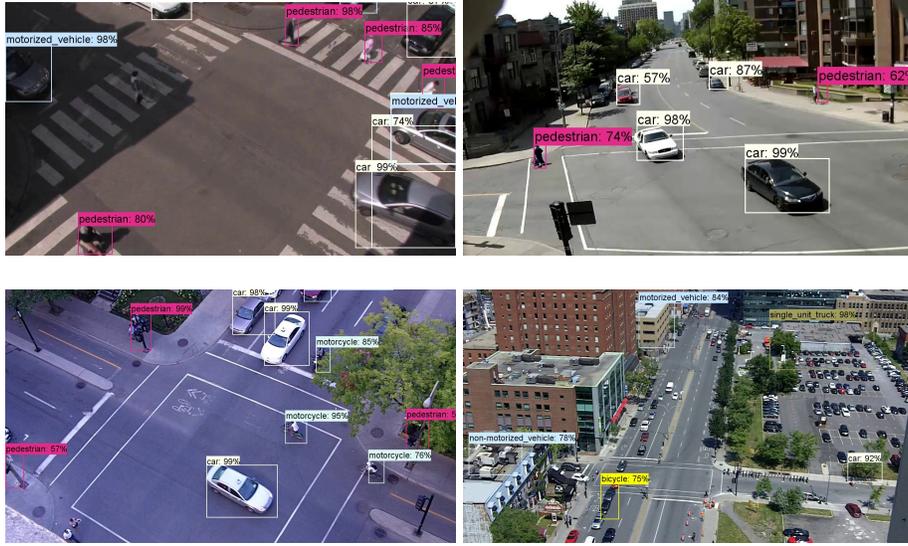


Fig. 2: Samples frames with detections from the Urban Tracker dataset



Fig. 3: Typical detections obtained from the René-Lévesque video.

Furthermore, contrarily to background subtraction or optical flow-based methods, RFCN gives detection outputs also for cars that are parked or for a car on a advertising billboard. Therefore, the data association process is distracted by many irrelevant objects. In such cases, standard NMS is not very useful in a traffic scene. Although NMS is used, it is insufficient to eliminate all the redundancies.

Since the proposed method is intrinsically dependent on the results from the detection module, the mis-detections propagate and deteriorate the overall MOT performance. In this case, the existence of redundant tracks severely affects the MOTA score such that it falls into the negative range, as shown in Table 1. The

Table 1: Comparison of MOTA and MOTP scores for three videos of the Urban Tracker dataset with the inclusion and exclusion of label cost in the data association (the best results are in boldface).

	No. objects	MOTP		MOTA	
		with labels	without labels	with labels	without labels
Rouen	16	0.6870	0.6893	-0.1877	-0.4176
Sherbrooke	20	0.7488	0.7324	0.0266	-0.0023
St-Marc	28	0.7234	0.7136	-0.3657	-0.2749

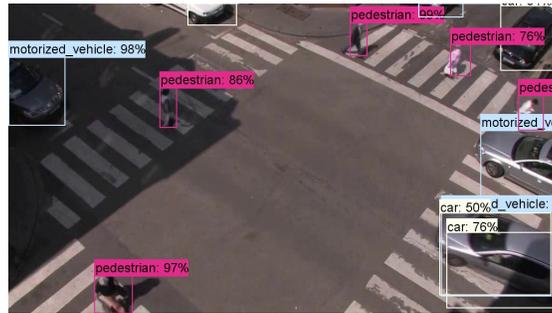


Fig. 4: An example of the redundant detection output for the same object.

MOTA takes into account the number of misses, false positives and mismatches from the produced trajectories.

However, it can be noted that MOT with inclusion of classification label generally gives higher MOTA. Among the different classes of labels from the detection module, the non-motorized vehicle label is currently excluded in the tracking framework since the occurrence of non-motorized vehicles is very rare in this experiment while parts of the background are sometimes mistakenly identified as objects from this class. MOTP is sometimes slightly better without labels as there are cases where tracking of an object fails in successive frames due to the switch of labels from the detection results. This is because with the labels, some matches are penalized and rejected because they are higher than the cost threshold. Therefore, the total number of matches is different, leading to slightly different values for MOTP. This occurrence is common among classes that share similarities such as pedestrians, bicycles and motorcycles, resulting in redundant tracks or fragmented tracks for the same object and thus lowering the overall tracking performance.

5 Conclusion

In this paper, the use of a modern multiclass object detector was investigated for the MOT task in traffic scenes. It was integrated in a baseline multiple object tracker. Results show that classification labels can be beneficial in MOT.

However, the outputs of the multiclass object detector are hardly usable because they include a large number of false detections, or detections of objects that are not of interest in the current application (e.g. parked cars). Small objects are also difficult to detect. As a result, to use such a detector, its output needs to be combined with another detector that can focus more precisely on objects of interest such as background subtraction or optical flow.

6 Acknowledgement

This research is partly funded by Fonds de Recherche du Quebec -Nature et Technologies(FRQ-NT) with team grant No. 2016-PR-189250 and Polytechnique Montréal PhD Fellowship. We gratefully acknowledge the support of NVIDIA Corporation with the donation of the Titan X GPU used for this work.

References

1. Luo, Zhiming and Frederic, B and Lemaire, Carl and Konrad, Janusz and Li, Shaozi and Mishra, Akshaya and Achkar, Andrew and Eichel, Justin and Jodoin, Pierre-Marc and others: MIO-TCD: A new benchmark dataset for vehicle classification and localization. *IEEE Transactions on Image Processing* (2018)
2. Dai, Jifeng and Li, Yi and He, Kaiming and Sun, Jian: R-fcn: Object detection via region-based fully convolutional networks. In: *Advances in neural information processing systems*, pp. 379–387. (2016)
3. Jodoin, Jean-Philippe and Bilodeau, Guillaume-Alexandre and Saunier, Nicolas: Tracking all road users at multimodal urban traffic intersections. *IEEE Transactions on Intelligent Transportation Systems* **17**(11), 99–110 (2016)
4. Bewley, Alex and Ge, Zongyuan and Ott, Lionel and Ramos, Fabio and Upcroft, Ben: Simple online and realtime tracking. In: *2016 IEEE International Conference on Image Processing (ICIP) on Proceedings*, pp. 3464-3468. (2016)
5. Kuhn, Harold W: The Hungarian method for the assignment problem. *Naval research logistics quarterly* **2**(5), 83–97 (1955)
6. Bernardin, Keni and Stiefelhagen, Rainer: Evaluating multiple object tracking performance: the CLEAR MOT metrics. *Journal on Image and Video Processing* (2008)
7. Malisiewicz, Tomasz and Gupta, Abhinav and Efron, Alexei A: Ensemble of exemplar-svms for object detection and beyond. In: *9Computer Vision (ICCV), 2011 IEEE International Conference*, pp. 89–96. Publisher, Location (2011)
8. Felzenszwalb, Pedro F and Girshick, Ross B and McAllester, David and Ramanan, Deva: Object detection with discriminatively trained part-based models. In: *IEEE transactions on pattern analysis and machine intelligence*, pp. 1627–1645. (2010)
9. Kalman, Rudolph Emil: A new approach to linear filtering and prediction problems. *Journal of basic Engineering* **82**(1), 35–45 (1960)
10. Mendes, Jean Carlo and Bianchi, Andrea Gomes Campos and Júnior, Alvaro R Pereira: Vehicle Tracking and Origin-Destination Counting System for Urban Environment. In: *Proceedings of the International Conference on Computer Vision Theory and Applications*, (2015)

11. Saunier, N. and Sayed, T.: A feature-based tracking algorithm for vehicles in intersections. In: Saunier, N. and Sayed, T.: A feature-based tracking algorithm for vehicles in intersections. In: The 3rd Canadian Conference on Computer and Robot Vision, pp. 59–59. (2006)
12. Kim, ZuWhan: Real time object tracking based on dynamic feature grouping with background subtraction, In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1-8. (2008)
13. Jun, Goo and Aggarwal, J. K. and Gokmen, Muhittin: Tracking and Segmentation of Highway Vehicles in Cluttered and Crowded Scenes. In: Proceedings of the 2008 IEEE Workshop on Applications of Computer Vision, pp. 1-6. (2008)
14. Benjamin Coifman and David Beymer and Philip McLauchlan and Jitendra Malik: A real-time computer vision system for vehicle tracking and traffic surveillance. In: Transportation Research Part C: Emerging Technologies, pp. 271–288. (1998)
15. Shi, J. and Tomasi, C.: Good features to track. In: Computer Vision and Pattern Recognition. pp. 593–600. (1994)
16. Torabi, A. and Bilodeau, G. -A: A Multiple Hypothesis Tracking Method with Fragmentation Handling. In: Canadian Conference on Computer and Robot Vision. pp. 8–15. (2009)
17. Beymer, D. and McLauchlan, P. and Coifman, B. and Malik, J.: A real-time computer vision system for measuring traffic parameters. In: IEEE Computer Society Conference on Computer Vision and Pattern Recognition. pp. 495–501. (1997)
18. Luis M. Fuentes and Sergio A. Velastin: People tracking in surveillance applications. In: Image and Vision Computing. pp. 1165–1171. (2006)
19. Sepehr Aslani and Homayoun Mahdavi-Nasab: Optical Flow Based Moving Object Detection and Tracking for Traffic Surveillance. In: International Journal of Electrical, Robotics, Electronics and Communications Engineering. pp. 773–777. (2013)
20. J. P. Jodoin and G. A. Bilodeau and N. Saunier: In: IEEE Winter conference on Applications of Computer Vision. pp. 885-892. (2016)
21. Redmon, Joseph and Divvala, Santosh and Girshick, Ross and Farhadi, Ali: You only look once: Unified, real-time object detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 779–788. (2016)