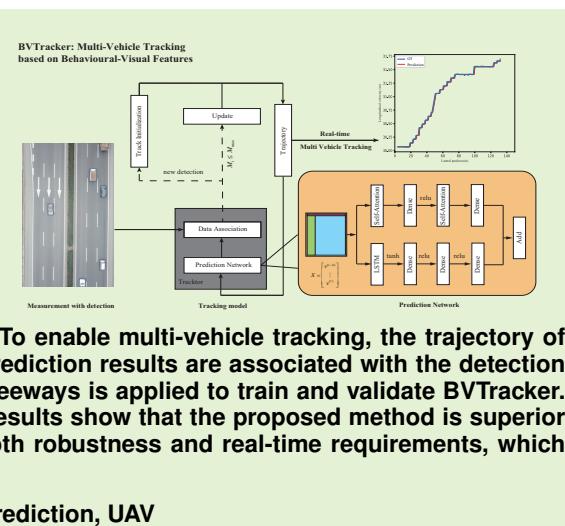


BVTracker: Multi-Vehicle Tracking based on Behavioural-Visual Features

Mingxu Li, Donghai Zhai, *Member, IEEE*, Da Yang, Lihua Xu

Abstract— In recent years, Unmanned Aerial Vehicle (UAV) is widely used in vehicle tracking. However, the objects in drone pictures always consist of fewer pixels due to their flying height. The lacking of visual information results in unreliable tracking results. Meanwhile, higher flying heights capture more targets, which makes it difficult to perform real-time inference with existing methods. To solve the above problems, we present a lightweight tracking model which employed behavioural information in the multi-vehicle tracking problem. Besides, BVTracker utilizes visual information as well as behavioural information independently and consists of two branches. One is a trajectory prediction model based on Long Short-Term Memory (LSTM) network and self-attention network. The other one is an association model base on the Hungarian Algorithm. To enable multi-vehicle tracking, the trajectory of each vehicle is predicted by the trajectory prediction model and the prediction results are associated with the detection findings. Moreover, a new dataset collected by UAVs for vehicles on freeways is applied to train and validate BVTracker. Compared with the state-of-art tracking algorithms, the experimental results show that the proposed method is superior to the existing algorithms at different frame rates and can achieve both robustness and real-time requirements, which profits the fast and effective traffic data analysis.

Index Terms— multi-vehicle tracking, trajectory extraction, trajectory prediction, UAV



I. INTRODUCTION

MULTI-VEHICLE TRACKING (MVT), which aims at identifying and tracking vehicles in image sequences, provides a reliable data source like velocity, attitude and direction for applications such as traffic flow analysis and security monitoring.

At present, many studies on the tracking problem have been conducted [1]–[7]. The idea behind the existing tracking algorithms is to find out and utilize more effective features to improve the detection and re-ID processes. However, due to the fact that from a bird's-eye view, the targets are similar to each other and lack sufficient visual information, the visual-based method needs a more complex model to deal with the confusing information. Besides, since the features between road and vehicles are totally different while the target is quite similar to other vehicles, sharing the features between the object detection and re-ID tasks may not perform well. To

meet the need of accuracy and low computation cost from the blooming traffic industry, a feasible idea is to introduce other information and realise tracking by detection. For a tracking-by-detection paradigm, the MVT problem is broken down into three parts: 1) localize the targets in the current frame; 2) predict the position based on the past information; 3) combine the predictions and detections together. The detection and re-ID tasks then can utilize the information independently.

Therefore, the object tracking model proposed in this paper is a hybrid model structure in which one is a trajectory prediction branch based on the LSTM network and self-attention network, while the other one is an association model base on the Hungarian Algorithm. The detection findings and the trajectory predictions model are associated by the Hungarian algorithm. To train and evaluate the proposed model, we collected a new multi-vehicle tracking dataset on freeways from a bird's-eye view. Compared with the state-of-the-art MVT models, our method significantly reduces the number of ID switches. Furthermore, BVTracker updates with an average frame per second (FPS) rate of 157.7, which is much higher than other trackers.

The main contributions of this work are as follows:

- We propose a new model for MVT problem that efficiently tracks vehicles from the bird's eye view of a stationary UAV. Our model minimizes ID switches and increases algorithm computation speed and accuracy.
- We utilize both behavioral information and environmental constraints to optimize tracking performance. By incor-

(Corresponding authors: Donghai Zhai)

Mingxu Li is with School of Computing and Artificial Intelligence, Southwest Jiaotong University, Chengdu China (e-mail: Li.Mingxu@my.swjtu.edu.cn).

Donghai Zhai is with School of Computing and Artificial Intelligence, Southwest Jiaotong University, Chengdu, Sichuan, 610031, China (e-mail: dhzhai@swjtu.edu.cn).

Da Yang is with School of Transportation and Logistics, Southwest Jiaotong University, Chengdu, Sichuan, 610031, China (e-mail: yang8@swjtu.edu.cn).

Lihua Xu is with School of Physical Science and Technology, Southwest Jiaotong University, Chengdu, Sichuan, 610031, China (e-mail: xulihua@home.swjtu.edu.cn)

porating these factors into our model, we enhance the accuracy and efficiency of our MVT system.

- We have collected a comprehensive dataset for MVT problem that covers various shooting altitudes and illumination conditions. This dataset is formatted in MOT16, facilitating compatibility with a range of tracking algorithms.

The paper is structured as follows: Section II reviews the related work on tracking methods. Section III presents our tracking model. In Section IV, we describe a new dataset and compare the outputs with different models. Finally, Section V makes conclusions.

II. LITERATURE REVIEW

With the emergence of excellent object detection algorithms [8]–[11], tracking-by-detection methods have been a feasible solution to MVT problems. Alex Bewley et al. proposed a simple online and real-time tracking (SORT) model [12]. In this model, the predictive coordinates were obtained by Kalman filtering, which was based on the target's historical motion state. Based on the SORT framework, Simple online and real-time tracking with a deep association metric (DEEP-SORT) [4] was proposed. This algorithm introduces visual information into the SORT framework, and uses visual features as measurement indicators for data association. Such work also included [2], [3], [6], [13]. For this series of methods, first, the position of the tracked vehicle is extracted by the detector based on Convolutional Neural Networks (CNN) and the prediction is calculated by the predictor. Then the detection results and prediction results are fed into the association method. Finally, tracks are formed by associating detections and predictions across adjacent frames.

Unlike traditional surveillance equipment that captures road sections discontinuously [14]–[19], UAV surveillance is able to cover continuous roads and provides more comprehensive data. Therefore, analysis of traffic conditions from a bird's-eye view has attracted the attention of many researchers recently. Shi et al. [20] proposed a novel vehicle trajectory extraction method that can extract high-granularity vehicle trajectories from aerial videos. Besides, they put forward a vehicle trajectory dataset named High-Granularity Highway Simulation (HIGH-SIM). Compared with the Next Generation Simulation dataset (NGSIM) [21], the HIGH-SIM dataset is more reasonable in speed and acceleration distributions. Barmpounakis et al. [22] created an urban traffic dataset to study congestion. In this dataset, traffic streams in a congested area of 1.3 km^2 are recorded by 10 drones, which helps researchers to develop and test their own models. Fukun et al. [23] proposed a context-based tracking method for remote sensing target tracking of UAV aerial videos to improve the tracking performance. The Response-Adaptive Context-Aware Correlation Filter (RA-CACF) model is also introduced into the network architecture to improve tracking performance. Ke et al. [24] presented an optical flow-based framework as a detector and an ensemble classifier as a tracker to match identified vehicles in adjacent frames. The model fills a gap in traffic flow analysis and has produced encouraging results,

but the tracking result is very sensitive to the detection result. Once a vehicle is lost in a few frames, wrong trajectories may be generated. Though UAVs are widely used in traffic analysis, as shown in [25]–[28], due to the challenges of tracking and detection from UAV view, vision-based tracking algorithms tend to underperform. Therefore, we propose a lightweight and accurate multi-vehicle tracking model to meet the growing need of the traffic industry.

III. METHODOLOGY

In this section, we introduce BVTracker including the trajectory prediction network, the association algorithm and the collected dataset.

A. Model Structure

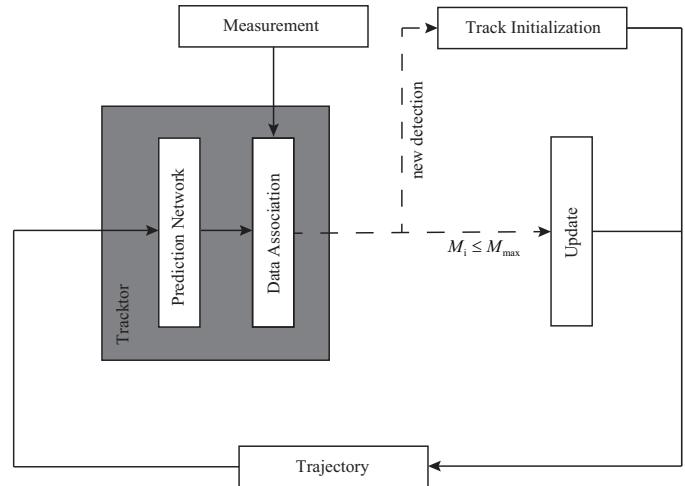


Fig. 1. Overview of BVTracker. The input image is first fed to a detection model based on CNN to extract the position of each vehicle. Then we put the measurement into the tracking model and associate the measurement with existing trajectories. Finally, a new trajectory will be generated or the existing trajectories will be updated.

The structure of BVTracker is displayed in Fig. 1. The bounding boxes of the vehicles shot from a drone are obtained by the detector. Then, the data association method is carried out to match the detected vehicles and the predicted trajectory. For each trajectory, we count the frame number since the last prediction result is successfully matched with the detection result, which is presented as M_i , where i is the trajectory's id. Meanwhile, to distinguish whether the vehicle has left the screen or is shortly lost, we define a threshold M_{\max} . Once M_i exceeds a predefined value M_{\max} , we stop to predict its position in the next frame. A new trajectory will be generated when a new detection is found. When a successful matching happens or a new trajectory is created, the detection result will be updated to its historical trajectory. After that, the historical trajectory is fed into our prediction network. The outputs of the prediction network are associated with the detection results in the next frame and realizing the online multi-vehicle tracking.

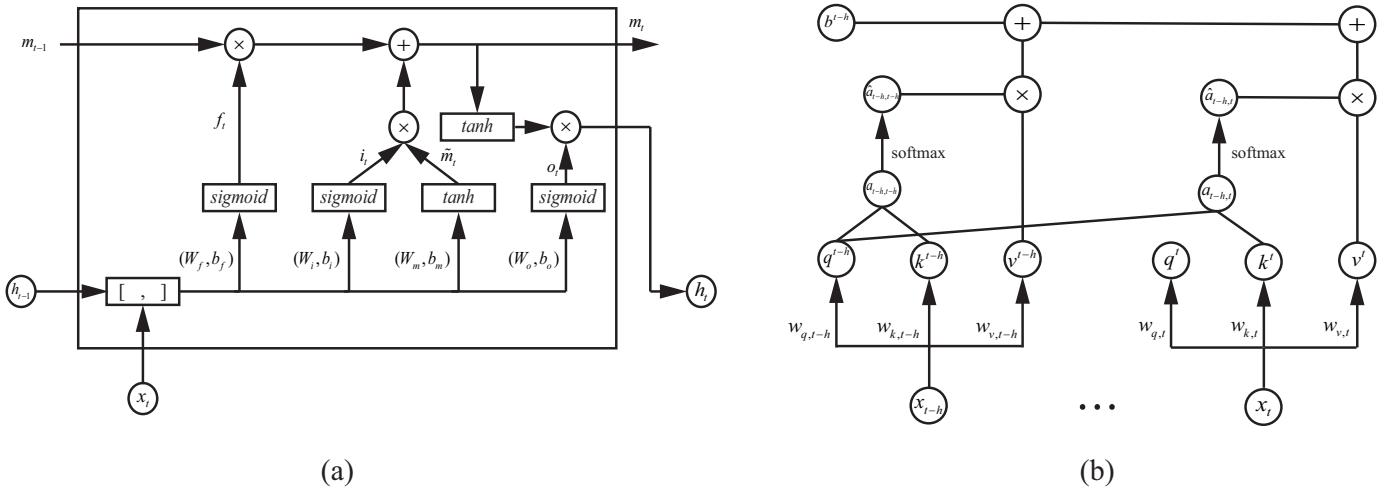


Fig. 3. (a). The detail of LSTM cell. *sigmoid* and *tanh* represent a sigmoid and tanh activation function respectively and $[\cdot, \cdot]$ indicates concatenate x_t and h_t together. (b). An illustration of the work process of the self-attention layer.

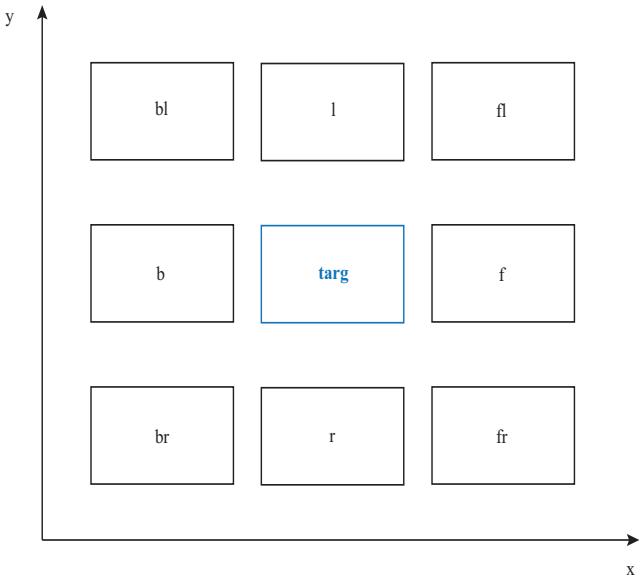


Fig. 4. The relative position between vehicles. The left and right vehicles of the target are denoted as l ($y_l > y_{targ}$), r ($y_r < y_{targ}$). The vehicles preceding l , $targ$, and r are denoted as fl , f , and fr . Similarly, the vehicles at the back of the target are denoted as bl , b , and br .

vector into a sequence format like $[seq_len, feature_len]$. The sequence includes information on the current situation and the environmental constraint.

C. Data Association

Based on the detection and trajectory prediction results, the problem of multi-vehicle tracking is transformed into a frame matching problem which is solved optimally using the Hungarian algorithm [4], [12]. The objective function of the Hungarian algorithm is as follows,

$$\min z = \sum_q \sum_p C_{qp} x_{qp} \quad (11)$$

where x_{qp} denotes whether the q_{th} detection frame is assigned to match the p_{th} prediction and equals 1 or 0, and C_{qp} denotes

the cost between q_{th} detection and the p_{th} prediction.

In BVTracker, we use the Intersection over Union (IoU) [30] to measure the cost between prediction which is denoted as,

$$C_{ij} = -\frac{\text{Area}(prediction) \cap \text{Area}(detection)}{\text{Area}(prediction) \cup \text{Area}(detection)} \quad (12)$$

Additionally, we found that IoU distance is able to handle the occlusion problem implicitly. Once an occlusion happens, only the occluder is captured by the detector. Since the target length-width ratio is stable from a bird's-eye view, and the target tracking box always keeps a similar scale, the occluder bounding box will not be matched with the target vehicle.

IV. EXPERIMENTS

To evaluate the effectiveness of BVTracker, we compare the proposed model with several state-of-the-art algorithms. Both designed experiments and the dataset are described in this section. All experiments were conducted on Ubuntu 18 platform. The computer has an Intel(R) Core(TM) i7-10700 CPU and NVIDIA GeForce RTX 3060 GPU with a memory of 32 GB.

A. Dataset

Although new challenges are brought to tracking tasks by UAVs, limited datasets have been published for vehicle tracking. Due to the movement of the camera, the published datasets don't provide the trajectory information. To validate BVTracker, we collect a new dataset for the MVT problem from a bird's-eye view. Based on the different situations on the freeway, the dataset is captured at a fixed position over a freeway and includes scenes varying from day, evening, 50 m and 100 m. Specifically, this dataset covers different scenes like the interference from street lights and shadows, which poses a challenge for MVT tasks. The video is taken at 30 fps with a resolution of 1920*1080. About 47000 frames are extracted from the videos and finely annotated

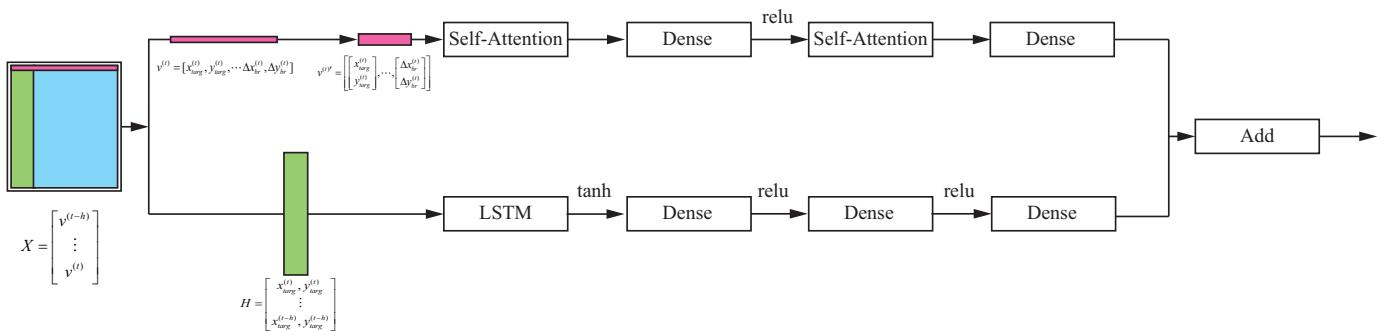


Fig. 5. Overview of our prediction model. The matrix Eq. 8 is split into two slices. One matrix is presented as H and denotes the motion of the target vehicle. The other slice is named $v^{(t)}$ and denotes the current state. The two slices are fed into different branches of our prediction model. Finally, the two branches are added together and generate the final output.



Fig. 6. Sample of annotated frame in our dataset.

with bounding boxes by the vatic tool [31]. Finally, about 2.3 million bounding boxes are labelled, as well as 2740 different vehicles are marked. We select about 30% of the dataset collected on both day and night as the evaluation data, and the rest of the dataset as the training data. The whole dataset is saved in the MOT16 [32] format. Fig. 6 shows some sample frames in our dataset.

B. Implementation Details

We use the standard tracking metrics [33] to evaluate models:

- Multi-object tracking accuracy (MOTA): a measure that combines overall tracking accuracy including identity switches, false positives as well as false negatives which is defined by

$$MOTA = 1 - \frac{\sum_t (fn_t + fp_t + mme_t)}{\sum_t gt} \quad (13)$$

where fn_t is the number of false negative, t is the frame number, fp_t is the number of false positive, mme_t is the number of ID switches and gt is the number of ground-truth.

- Multi-object tracking precision (MOTP): a measure that combines overall tracking precision including bounding

boxes overlap between ground-truth and reported positions, which is defined by

$$MOTP = \frac{\sum_{k,t} d_t^k}{\sum_t c_t} \quad (14)$$

where k is the identity of each match, c_t is the number of successful matches between detections and predictions at frame t , d_t^k is the distance between each detection-prediction pair.

- Mostly tracked objects (MT): a ratio of ground-truth trajectories which have the same label with their hypothesis at least 80% in their life span.
- Mostly lost objects (ML): a ratio of ground-truth trajectories which have the same label with their hypothesis at most 20% in their life span.
- Identity switches (IDs): the number of identity switches.
- IDF1 Score (IDF1): a ratio of correctly identified targets over the ground truth and the average number of calculated targets. IDF1 is given by:

$$IDF1 = \frac{2 \times IDTP}{2 \times IDTP + IDFP + IDFN} \quad (15)$$

where $IDTP$ is the number of true positive identities and $IDFP$, $IDFN$ are the number of false positive, the number of false negative respectively.

In our evaluation, to avoid the influence of different units, the pixel distance extracted from videos is converted into the Système International d'Unités (SI) units. In addition, YOLO v4 proposed by Bochkovskiy et al. [8] is used as the detector and trained by the collected dataset. The outputs of the detector with confidence larger than 0.6 are selected. The comparison algorithms are trained in the same experimental setting as ours.

In Subsection III-A, M_{\max} , the maximum number that is allowed to skip for unmatched tracks is set to 2. To more realistically model the decision-making during driving, we train the network using inputs with a time window of 30 frames, representing the observed data within the past 1 s. We utilize the Adam optimizer [34] with a starting learning rate of 10^{-4} to update network parameters. The proposed model is implemented by PyTorch. Furthermore, the batch size is set to 128 and Root Mean Square Error (RMSE) is used as the loss function.

$$RMSE = \sqrt{\frac{1}{N} \sum_{t=1}^N \|P_{pre} - P_{gt}\|^2} \quad (16)$$

where P_{pre} is the predicted value (x_{pre}, y_{pre}) , P_{gt} is the ground truth (x_{gt}, y_{gt}) and N is the number of predictions in each batch.

C. Ablation Studies

To demonstrate the impact of different factors of the model, we have experimented with the following variations of the model:

- 1) BVTracker.
- 2) Only uses target car information.
- 3) Only LSTM network.
- 4) Adding a third self-attention layer with 256 neurons.
- 5) Only one self-attention block.
- 6) Using Kalman Filtering to replace deep learning model.

The results are summarized in Table I. By comparing whether we use information from environmental targets, it can be found that the experimental results are consistent with our expectations. In experiments, compared to using only the target vehicle or only the LSTM branch, the proposed model reduced ID switches by 94% (1835) and 99% (12856), respectively. BVTracker had better performance on the main indicators but had only slight differences on other indicators, as its model with or without self-attention layers was similar to the proposed model. On the other hand, the Kalman filter-based

method had similar ID switch results to the proposed method, but due to its inability to handle environmental constraint information, it had a significant difference in MOTA compared to the proposed method.

For scenarios such as lane changing and acceleration, the model considering environmental factors can track the target more effectively. Fig. 7 demonstrates a comparison between whether the model uses environmental constraints. As shown in the upper sequence, since the model uses only previously observed data, the predicted response delay when the vehicle state changed, which causes an ID switch. Fig. 8 shows the performance of the two models mentioned above when a vehicle accelerates. As seen in Fig. 8, the incorrect tracking of two vehicles happens at the same time in the upper sequence. The missing environmental constraints make it hard for the model to distinguish them. While BVTracker keeps the target's ID stable.

D. Comparing with SOTA methods

To verify the effectiveness of the proposed method, we compare BVTracker with several state-of-the-art algorithms.

- DEEPSORT [4]: DEEPSORT is based on the two-step MOT method, which takes visual information into consideration and a cascade matching mechanism to improve the accuracy of inter-frame association.
- FairMOT [35]: FairMOT combines the detection task and the tracking task together. This model implements tracking based on Re-ID features.
- MPNTracker [36]: This framework uses the message-passing mechanism in graph convolutional networks directly. And tracking by using features based on Re-ID as well as position data. MPNTracker has been ranked at the top of various MOT Challenge lists [37].
- UniTrack [36]: UniTrack is a unified framework to handle tracking problems. This model contains a task-agnostic appearance model and proposes a novel reconstruction-based similarity metric for association tasks.

The results are given in Table II. MPNTracker and Uni-Tracker are powerful trackers which perform well on FN, MT and ML indicators but bring much more complexity and lead to a drop in FPS. Besides, as can be seen from the comparison between the results of UniTrack(resnet18) and UniTrack(resnet50), a deeper appearance model helps little on IDs and FP indicators but lead to a significantly slower updates rate. FairMOT and DEEPSORT get remarkable

TABLE I
COMPARISON OF THE DIFFERENT MODELS

| Model | MOTA \uparrow | MOTP \uparrow | FP \downarrow | FN \downarrow | MT \uparrow | ML \downarrow | IDs \downarrow |
|---|-----------------|-----------------|-----------------|-----------------|---------------|-----------------|------------------|
| BVTracker | 76.0 | 76.0 | 48700 | 39552 | 1321 | 41 | 108 |
| Only use target car information | 75.0 | 75.9 | 48392 | 39559 | 1321 | 41 | 1943 |
| Only LSTM network | 59.5 | 75.5 | 40099 | 92500 | 685 | 202 | 12964 |
| Adding a third self-attention layer with 256 neurons | 75.4 | 75.9 | 48836 | 39427 | 1322 | 40 | 110 |
| Only one self-attention block | 75.4 | 75.9 | 48903 | 39417 | 1322 | 60 | 111 |
| Using Kalman Filtering to replace deep learning model | 75.3 | 75.9 | 48372 | 40231 | 1315 | 42 | 120 |

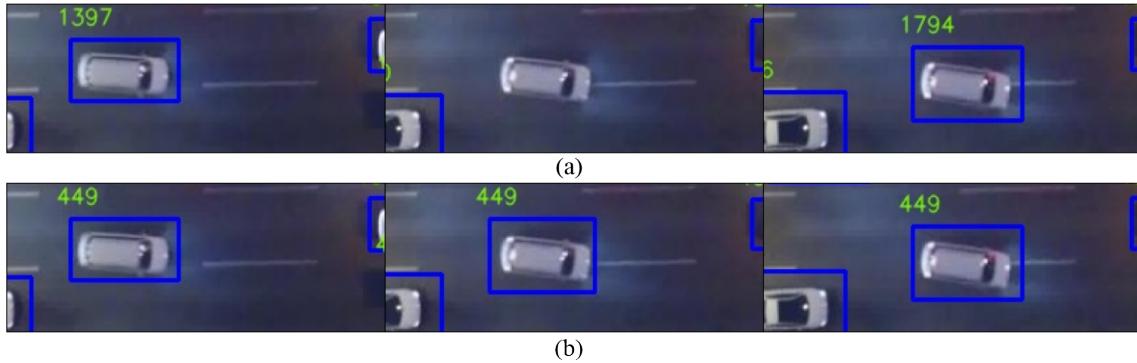


Fig. 7. Comparison of BVTracker in the lane change scenario with a model that does not use environmental vehicle information. **(a)** Model that without environmental information. **(b)** Our model.

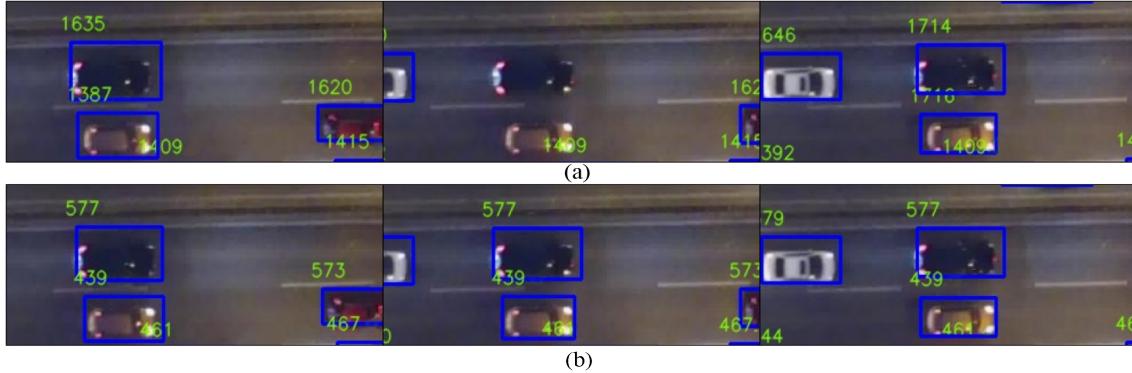


Fig. 8. Acceleration and vehicles run parallel scenarios in which BVTracker is compared with a model that does not apply information about the surrounding vehicles. **(a)** Model that without environmental information. **(b)** Our model.

TABLE II
TRACKING RESULTS ON THE PROPOSED DATASET

| model | MOTA \uparrow | MOTP \uparrow | FP \downarrow | FN \downarrow | MT \uparrow | ML \downarrow | IDs \downarrow | IDF1 \uparrow | FPS \uparrow |
|--------------------|-----------------|-----------------|-----------------|-----------------|---------------------|------------------|------------------|-----------------|----------------|
| BVTracker | 76.0 | 76.0 | 48700 | 39552 | 1321 (92.3%) | 41 (2.9%) | 108 | 86.4 | 157.7 |
| FairMOT | 52.5 | 71.8 | 89079 | 81113 | 1153 (80.6%) | 96 (6.7%) | 501 | 74.8 | 15.4 |
| MPNTracker | 74.3 | 76.0 | 53546 | 38681 | 1332 (93.1%) | 43 (3.0%) | 321 | 83.5 | 38.7 |
| DEEPSORT | 72.9 | 75.5 | 54138 | 43156 | 1313 (91.8%) | 40 (2.8%) | 184 | 84.5 | 39.6 |
| UniTrack(resnet18) | 74.1 | 75.5 | 51818 | 40764 | 1321 (92.4%) | 38 (2.7%) | 689 | 85.9 | 6.4 |
| UniTrack(resnet50) | 74.1 | 75.5 | 51784 | 40784 | 1319 (92.2%) | 38 (2.7%) | 676 | 85.9 | 2.9 |

results on public datasets but get poor performances on the MVT problem due to their results depending more on visual information. Since BVTracker utilizes visual and behavioural information independently, the proposed model gets the best performance on the leading indicators. It further decreases the number of ID switches reduces from 501, 312, 184, 689 and 676 to 108 and runs at the highest frame rate.

Besides, as demonstrated in Fig. 9, with the decrement of the frame rate, the advantages of BVTracker become more obvious. Compared with the state-of-the-art algorithms, our method achieves much better results in ID-related indicators, while making a slight sacrifice in some metrics like MOTA when the frame rate comes to 6. This is mainly due to the fact that BVTracker takes the influence of surrounding

vehicles into account. It is noteworthy that the DEEPSORT algorithm also shows a significant drop in the IDs metric when the frame rate is lowered. However, from the MOTA figure, the comprehensive index of the DEEPSORT algorithm has a more dramatic drop. This means that the algorithm is only tracking a very small number of vehicles at low frame rates, hence the drop in the IDs metric. The UniTrack(resnet18) and UniTrack(resnet50) perform almost the same on different indexes. It means that a more complex appearance model is not helpful to the robustness and accuracy in the MVT problem. Notably, our paper uses the YOLOv4 algorithm, which has been officially released and is widely used in multi-object tracking algorithms. Thus the impact from the detection phase is stable.

- [33] B. Kene and S. Rainer, "Evaluating multiple object tracking performance: The clear mot metrics," *EURASIP Journal on Image and Video Processing*, vol. 2008, no. 1, 2008.
- [34] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [35] Y. Zhang, C. Wang, X. Wang, W. Zeng, and W. Liu, "Fairmot: On the fairness of detection and re-identification in multiple object tracking," *International Journal of Computer Vision*, vol. 129, no. 11, pp. 3069–3087, 2021.
- [36] G. Brasó and L. Leal-Taixé, "Learning a neural solver for multiple object tracking," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [37] L. Leal-Taixé, A. Milan, I. Reid, S. Roth, and K. Schindler, "Motchallenge 2015: Towards a benchmark for multi-target tracking," *arXiv preprint arXiv:1504.01942*, 2015.



Lihua Xu received the M.S. degree in Communication and Information System and the Ph.D. degree in Optics from Sichuan University, Chengdu, China, in 1999 and 2002 respectively.

He is currently an Associate Professor with the School of Physical Science and Technology, Southwest Jiaotong University. His research interests include automated vehicles and deep learning.



Mingxu Li received the M.S. degree in software engineering from Southwest Jiaotong University, Chengdu, China, in 2021. He is currently pursuing his Ph.D. degree in School of Computing and Artificial Intelligence from Southwest Jiaotong University. His research interests include automatic driving and computer vision.



Donghai Zhai (IEEE Member) received his Ph.D. degree in traffic information engineering and control from Southwest Jiaotong University, China, in 2003.

From 2003 to 2005 he was employed at IBM China Research Laboratory. Since 2006 he has been associated with Southwest Jiaotong University in School of Computing and Artificial Intelligence. He has been a visiting scholar at Louisiana State University, Baton Rouge in 2016. His research interests include autonomous driving, digital image processing, computer vision, and pattern recognition.

autonomous driving, digital image processing, computer vision, and pattern recognition.



Da Yang received the B.S. and M.S. degrees in logistics engineering and the Ph.D. degree in transportation engineering from Southwest Jiaotong University, Chengdu, China, in 2007, 2009, and 2013, respectively. He studied in the University of Wisconsin-Madison for two years between August 2010 and August 2012.

He is currently an Associate Professor with the School of Transportation and Logistics, Southwest Jiaotong University. He is also a Post-Doctoral Research Associate with the Traffic Management Research Institute of the Ministry of Public Security. His research interests include automated vehicles, traffic flow, and vehicular networks.