

Joint Monocular 3D Vehicle Detection and Tracking

Hou-Ning Hu^{1*}, Qi-Zhi Cai^{2*}, Dequan Wang³, Ji Lin^{4*},
 Min Sun¹, Philipp Krähenbühl⁵, Trevor Darrell³, Fisher Yu³
¹National Tsing Hua University ²Sinovation Ventures AI Institute
³UC Berkeley ⁴MIT ⁵UT Austin

Abstract

Vehicle 3D extents and trajectories are critical cues for predicting the future location of vehicles and planning future agent ego-motion based on those predictions. In this paper, we propose a novel online framework for 3D vehicle detection and tracking from monocular videos. The framework can not only associate detections of vehicles in motion over time, but also estimate their complete 3D bounding box information from a sequence of 2D images captured on a moving platform. Our method leverages 3D box depth-ordering matching for robust instance association and utilizes 3D trajectory prediction for re-identification of occluded vehicles. We also design a motion learning module based on an LSTM for more accurate long-term motion extrapolation. Our experiments on simulation, KITTI, and Argoverse datasets show that our 3D tracking pipeline offers robust data association and tracking. On Argoverse, our image-based method is significantly better for tracking 3D vehicles within 30 meters than the LiDAR-centric baseline methods.

1. Introduction

Autonomous driving motivates much of contemporary visual deep learning research. However, many commercially successful approaches to autonomous driving control rely on a wide array of views and sensors, reconstructing 3D point clouds of the surroundings before inferring object trajectories in 3D. In contrast, human observers have no difficulty in perceiving the 3D world in both space and time from simple sequences of 2D images rather than 3D point clouds, even though human stereo vision only reaches several meters. Recent progress in monocular object detection and scene segmentation offers the promise to make low-cost mobility widely available. In this paper, we explore architectures and datasets for developing similar capabilities using deep neural networks.

*Work was done while Hou-Ning Hu, Qi-Zhi Cai and Ji Lin were at the Berkeley DeepDrive Center

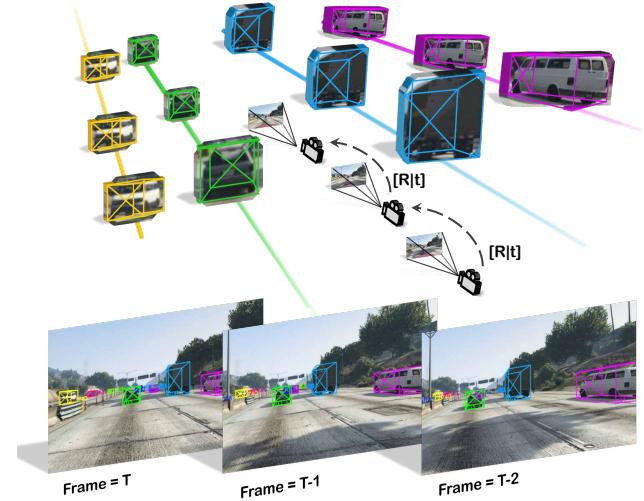


Figure 1: Joint online detection and tracking in 3D. Our dynamic 3D tracking pipeline predicts 3D bounding box association of observed vehicles in image sequences captured by a monocular camera with an ego-motion sensor.

Monocular 3D detection and tracking are inherently ill-posed. In the absence of depth measurements or strong priors, a single view does not provide enough information to estimate 3D layout of a scene accurately. Without a good layout estimate, tracking becomes increasingly difficult, especially in the presence of large ego-motion (e.g., a turning car). The two problems are inherently intertwined. Robust tracking helps 3D detection, as information along consecutive frames is integrated. Accurate 3D detection helps to track, as ego-motion can be factored out.

In this paper, we propose an online network architecture to jointly track and detect vehicles in 3D from a series of monocular images. Figure 1 provides an overview of our 3D tracking and detection task. After detecting 2D bounding boxes of objects, we utilize both world coordinates and re-projected camera coordinates to associate instances across frames. Notably, we leverage novel occlusion-aware associ-

ation and depth-ordering matching algorithms to overcome the occlusion and reappearance problems in tracking. Finally, we capture the movement of instances in a world coordinate system and update their 3D poses using LSTM motion estimation along a trajectory, integrating single-frame observations associated with the instance over time.

Like any deep network, our model is data hungry. The more data we feed it, the better it performs. However, existing datasets are either limited to static scenes [51], lack the required ground truth trajectories [35], or are too small to train contemporary deep models [16]. To bridge this gap, we resort to realistic video games. We use a new pipeline to collect large-scale 3D trajectories, from a realistic synthetic driving environment, augmented with dynamic meta-data associated with each observed scene and object.

To the best of our knowledge, we are the first to tackle the estimation of complete 3D vehicle bounding box tracking information from a monocular camera. We jointly track the vehicles across frames based on deep features and estimate the full 3D information of the tracks including position, orientation, dimensions, and projected 3D box centers of each object. The depth ordering of the tracked vehicles constructs an important perceptual cue to reduce the mismatch rate. Our occlusion-aware data association provides a strong prior for occluded objects to alleviate the identity switch problem. Our experiments show that 3D information improves predicted association in new frames compared to traditional 2D tracking, and that estimating 3D positions with a sequence of frames is more accurate than single-frame estimation.

2. Related Works

Object tracking has been explored extensively in the last decade [54, 46, 49]. Early methods [5, 15, 27] track objects based on correlation filters. Recent ConvNet-based methods typically build on pre-trained object recognition networks. Some generic object trackers are trained entirely online, starting from the first frame of a given video [20, 1, 25]. A typical tracker will sample patches near the target object which are considered as foreground and some farther patches as background. These patches are then used to train a foreground-background classifier. However, these online training methods cannot fully utilize a large amount of video data. Held *et al.* [22] proposed a regression-based method for offline training of neural networks, tracking novel objects at test-time at 100 fps. Siamese networks also found in use, including tracking by object verification [50], tracking by correlation [3], tracking by detection [13]. Yu *et al.* [53] enhance tracking by modeling a track-let into different states and explicitly learns an Markov Decision Process (MDP) for state transition. Due to the absence of 3D information, it just uses 2D location to decide whether a track-let is occluded.

All those methods only take 2D visual features into consideration, where the search space is restricted near the orig-

inal position of the object. This works well for a static observer, but fails in a dynamic 3D environment. Here, we further leverage 3D information to narrow down the search space, and stabilize the trajectory of target objects.

Sharma *et al.* [48] uses 3D cues for 2D vehicle tracking. Scheidegger *et al.* [47] also adds 3D kalman filter on the 3D positions to get more consistent 3D localization results. Because the goals are for 2D tracking, 3D box dimensions and orientation are not considered. Osep *et al.* [37] and Li *et al.* [31] studies 3D bounding box tracking with stereo cameras. Because the 3D depth can be perceived directly, the task is much easier, but in many cases such as ADAS, large-baseline stereo vision is not possible.

Object detection reaped many of the benefits from the success of convolutional representation. There are two mainstream deep detection frameworks: 1) two-step detectors: R-CNN [18], Fast R-CNN [17], and Faster R-CNN [41]. 2) one-step detectors: YOLO [38], SSD [33], and YOLO9000 [39].

We apply Faster R-CNN, one of the most popular object detectors, as our object detection input. The above algorithms all rely on scores of labeled images to train on. In 3D tracking, this is no different. The more training data we have, the better our 3D tracker performs. Unfortunately, getting a large amount of 3D tracking supervision is hard.

Driving datasets have attracted a lot of attention in recent years. KITTI [16], Cityscapes [10], Oxford RobotCar [34], BDD100K [57], NuScenes [6], and Argoverse [7] provide well annotated ground truth for visual odometry, stereo reconstruction, optical flow, scene flow, object detection and tracking. However, their provided 3D annotation is very limited compared to virtual datasets. Accurate 3D annotations are challenging to obtain from humans and expensive to measure with 3D sensors like LiDAR. Therefore these real-world datasets are typically small in scale or poorly annotated.

To overcome this difficulty, there has been significant work on virtual driving datasets: virtual KITTI [15], SYNTHIA [44], GTAV [43], VIPER [42], CARLA [11], and Free Supervision from Video Games (FSV) [26]. The closest dataset to ours is VIPER [42], which provides a suite of videos and annotations for various computer vision problems while we focus on object tracking. We extend FSV [26] to include object tracking in both 2D and 3D, as well as fine-grained object attributes, control signals from driver actions.

In the next section, we describe how to generate 3D object trajectories from 2D dash-cam videos. Considering the practical requirement of autonomous driving, we primarily focus on online tracking systems, where only the past and current frames are accessible to a tracker.

3. Joint 3D Detection and Tracking

Our goal is to track objects and infer their precise 3D location, orientation, and dimension from a single monocular video stream and a GPS sensor. Figure 2 shows an overview

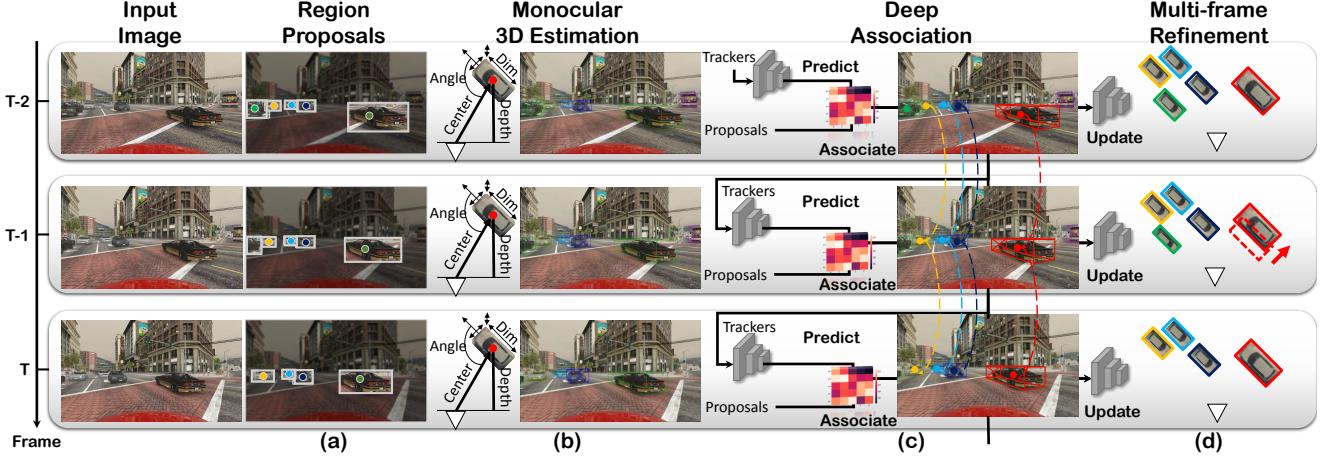


Figure 2: Overview of our monocular 3D tracking framework. Our online approach processes monocular frames to estimate and track region of interests (RoIs) in 3D (a). For each ROI, we learn 3D layout (*i.e.*, depth, orientation, dimension, a projection of 3D center) estimation (b). With 3D layout, our LSTM tracker produces robust linking across frames leveraging occlusion-aware association and depth-ordering matching (c). With the help of 3D tracking, the model further refines the ability of 3D estimation by fusing object motion features of the previous frames (d).

of our system. Images are first passed through a detector network trained to generate object proposals and centers. These proposals are then fed into a layer-aggregating network which infers 3D information. Using 3D re-projection to generate similarity metric between all trajectories and detected proposals, we leverage estimated 3D information of current trajectories to track them through time. Our method also solves the occlusion problem in tracking with the help of occlusion-aware data association and depth-ordering matching. Finally, we re-estimate the 3D location of objects using the LSTM through the newly matched trajectory.

3.1. Problem Formulation

We phrase the 3D tracking problem as a supervised learning problem. We aim to find N trajectories $\{\tau^1, \dots, \tau^N\}$, one for each object in a video. Each trajectory τ^i links a sequence of detected object states $\{s_a^{(i)}, s_{a+1}^{(i)}, \dots, s_b^{(i)}\}$ starting at the first visible frame a and ending at the last visible frame b . The state of an object at frame a is given by $s_a = (P, O, D, F, \Delta P)$, where P defines the 3D world location (x, y, z) of the object, and ΔP stands for its velocity $(\dot{x}, \dot{y}, \dot{z})$. O, D, F denotes for object orientation θ , dimension (l, w, h) and appearance feature f_{app} , respectively. In addition, we reconstruct a 3D bounding box X for each object, with estimated P, O, D and the projection $c = (x_c, y_c)$ of 3D box's center in the image. The bounding boxes enable the use of our depth-ordering matching and occlusion-aware association. Each bounding box X also forms a projected 2D box $M(X) = \{x_{\min}, y_{\min}, x_{\max}, y_{\max}\}$ projected onto a 2D image plane using camera parameters $M = K[R|t]$.

The intrinsic parameter K can be obtained from camera

calibration. The extrinsic parameter $[R|t]$ can be calculated from the commonly equipped GPS or IMU sensor. The whole system is powered by a convolutional network pipeline trained on a considerable amount of ground truth supervision. Next, we discuss each component in more detail.

3.2. Candidate Box Detection

In the paper, we employ Faster R-CNN [41] trained on our dataset to provide object proposals in the form of bounding boxes. Each object proposal (Figure 2(a)) corresponds to a 2D bounding box $\mathbf{d} = \{x_{\min}, y_{\min}, x_{\max}, y_{\max}\}$ as well as an estimated projection of the 3D box's center c . The detection results are used to locate the candidate vehicles and extract their appearance features. However, the centers of objects' 3D bounding boxes usually do not project directly to the center of their 2D bounding boxes. As a result, we have to provide an estimation of the 3D box center for better accuracy. More details about the estimation of the 3D center can be found in the supplementary material¹.

Projection of 3D box center. To estimate the 3D layout from single images more accurately, we extend the regression process to predict a projected 2D point of the 3D bounding box's center from an ROI-pooled feature F using L1 loss. Estimating a projection of 3D center is crucial since a small gap in the image coordinate will cause a gigantic shift in 3D. It is worth noting that our pipeline can be used with any off-the-shelf detector and our 3D box estimation module is extendable to estimate projected 2D points even if the detector is replaced. With the extended ROI head, the model

¹Supplementary material of Joint Monocular 3D Vehicle Detection and Tracking can be found at <https://eborboihuc.github.io/Mono-3DT/>

regresses both a bounding box \mathbf{d} and the projection of 3D box’s center c from an anchor point. ROIalign [21] is used instead of ROIpool to get the regional representation given the detected regions of interest (ROIs). This reduces the misalignment of two-step quantization.

3.3. 3D Box Estimation

We estimate complete 3D box information (Figure 2(b)) from an ROI in the image via a feature representation of the pixels in the 2D bounding box. The ROI feature vector F is extracted from a 34-layer DLA-up [56] using ROIalign. Each of the 3D information is estimated by passing the ROI features through a 3-layer 3x3 convolution sub-network, which extends the stacked Linear layers design of Mousavian *et al.* [36]. We focus on 3D location estimation consisting of object center, orientation, dimension and depth, whereas [36] focus on object orientation and dimension from 2D boxes. Besides, our approach integrates with 2D detection and has the potential to jointly training, while [36] crops the input image with pre-computed boxes. This network is trained using ground truth depth, 3D bounding box center projection, dimension, and orientation values. A convolutional network is used to preserve spatial information. In the case that the detector is replaced with another architecture, the center c can be obtained from this sub-network. More details of c can be found at Appendix A.

3D World Location. Contrasting with previous approaches, we also infer 3D location P from monocular images. The network regresses an inverse depth value $1/d$, but is trained to minimize the L1 loss of the depth value d and the projected 3D location P . A projected 3D location P is calculated using an estimated 2D projection of the 3D object center c as well as the depth d and camera transformation M .

Vehicle Orientation. Given the coordinate distance $\hat{x} = x_c - \frac{w}{2}$ to the horizontal center of an image and the focal length f , we can restore the global rotation θ in the camera coordinate from θ_l with simple geometry, $\theta = (\theta_l + \arctan \frac{\hat{x}}{f}) \bmod 2\pi$. Following [36] for θ_l estimation, we first classify the angle into two bins and then regress the residual relative to the bin center using Smooth L1 loss.

Vehicle Dimension. In driving scenarios, the high variance of the distribution of the dimensions of different categories of vehicles (e.g., car, bus) results in difficulty classifying various vehicles using unimodal object proposals. Therefore, we regress a dimension D to the ground truth dimension over the object feature representation using L1 loss.

The estimation of an object’s 3D properties provides us with an observation for its location P with orientation θ , dimension D and 2D projection of its 3D center c . For any new tracklet, the network is trained to predict monocular object state s of the object by leveraging ROI features. For any previously tracked object, the following association network is able to learn a mixture of a multi-view monocular

3D estimates by merging the object state from last visible frames and the current frame. First, we need to generate such a 3D trajectory for each tracked object in world coordinates.

3.4. Data Association and Tracking

Given a set of tracks $\{\tau^J, \dots, \tau^K\}$ at frame a where $1 \leq J \leq K \leq M$ from M trajectories, our goal is to associate each track with a candidate detection, spawn new tracks, or end a track (Figure 2(c)) in an online fashion.

We solve the data association problem by using a weighted bipartite matching algorithm. Affinities between tracks and new detections are calculated from two criteria: overlap between projections of current trajectories forward in time and bounding boxes candidates; and the similarity of the deep representation of the appearances of new and existing object detections. Each trajectory is projected forward in time using the estimated velocity of an object and camera ego-motion. Here, we assume that ego-motion is given by a sensor, like GPS, an accelerometer, gyro and/or IMU.

We define an affinity matrix $\mathbf{A}(\tau_a, \mathbf{s}_a)$ between the information of an existing track τ_a and a new candidate \mathbf{s}_a as a joint probability of appearance and location correlation.

$$\mathbf{A}_{\text{deep}}(\tau_a, \mathbf{s}_a) = \exp(-\|F_{\tau_a}, F_{\mathbf{s}_a}\|_1) \quad (1)$$

$$\mathbf{A}_{\text{2D}}(\tau_a, \mathbf{s}_a) = \frac{\mathbf{d}_{\tau_a} \cap \mathbf{d}_{\mathbf{s}_a}}{\mathbf{d}_{\tau_a} \cup \mathbf{d}_a} \quad (2)$$

$$\mathbf{A}_{\text{3D}}(\tau_a, \mathbf{s}_a) = \frac{\mathbf{M}(X_{\tau_a}) \cap \mathbf{M}(X_{\mathbf{s}_a})}{\mathbf{M}(X_{\tau_a}) \cup \mathbf{M}(X_{\mathbf{s}_a})}, \quad (3)$$

where $F_{\tau_a}, F_{\mathbf{s}_a}$ are the concatenation of appearance feature f_{app} , dimension D , center c , orientation θ and depth d . X_{τ_a} and $X_{\mathbf{s}_a}$ are the tracked and predicted 3D bounding boxes, \mathbf{M} is the projection matrix casting the bounding box to image coordinates, \mathbf{A}_{2D} and \mathbf{A}_{3D} is the Intersection of Union (IoU).

$$\begin{aligned} \mathbf{A}(\tau_a, \mathbf{s}_a) = & w_{\text{deep}} \mathbf{A}_{\text{app}}(\tau_a, \mathbf{s}_a) + w_{\text{2D}} \mathbf{A}_{\text{2D}}(\tau_a, \mathbf{s}_a) \\ & + w_{\text{3D}} \mathbf{A}_{\text{3D}}(\tau_a, \mathbf{s}_a) \end{aligned} \quad (4)$$

$w_{\text{deep}}, w_{\text{2D}}, w_{\text{3D}}$ are the weights of appearance, 2D overlap, and 3D overlap. We utilize a mixture of those factors as the affinity across frames, similar to the design of POI [55].

Comparing to 2D tracking, 3D-oriented tracking is more robust to ego-motion, visual occlusion, overlapping, and re-appearances. When a target is temporally occluded, the corresponding 3D motion estimator can roll-out for a period of time and relocate 2D location at each new point in time via the camera coordinate transformation.

Depth-Ordering Matching. We introduce instance depth ordering for assigning a detection to neighbor tracklets, which models the strong prior of relative depth ordering found in human perception. For each detection of interest (DOI), we consider potential associated tracklets in order of their depths. From the view of each DOI, we obtain the IOU

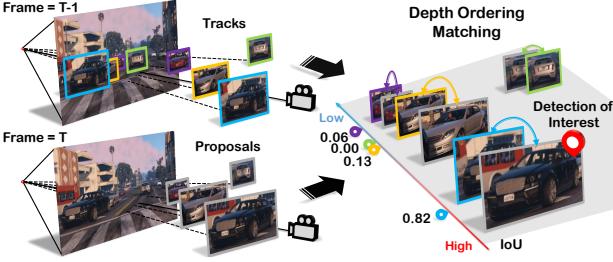


Figure 3: Illustration of depth-ordering matching. Given the tracklets and detections, we sort them into a list by depth order. For each detection of interest (DOI), we calculate the IOU between DOI and non-occluded regions of each tracklet. The depth order naturally provides higher probabilities to tracklets near the DOI.

of two non-occluded overlapping map from both ascending and descending ordering. To cancel out the ordering ambiguity of a distant tracklet, we filter out those tracklets with a larger distance to a DOI than a possible matching length. So Equation 3 becomes

$$\mathbf{A}_{3D}(\tau_a, \mathbf{s}_a) = \mathbb{1} \times \frac{\phi(\mathbf{M}(X_{\tau_a})) \cap \mathbf{M}(X_{\mathbf{s}_a})}{\phi(\mathbf{M}(X_{\tau_a})) \cup \mathbf{M}(X_{\mathbf{s}_a})}, \quad (5)$$

where $\mathbb{1}$ denotes if the tracklets is kept after depth filtering, and the overlapping function

$$\phi(\cdot) = \arg \min_x \{x | \text{ord}(x) < \text{ord}(x_0) \forall x_0 \in \mathbf{M}(X_{\tau_a})\}$$

captures pixels of non-occluded tracklets region with the nearest depth order. It naturally provides higher probabilities of linking neighbor tracklets than those layers away. In this way, we obtain the data association problem of moving objects with the help of 3D trajectories in world coordinates. Figure 3 depicts the pipeline of depth ordering. Finally, we solve data association using the Kuhn-Munkres algorithm.

Occlusion-aware Data Association. Similar to previous state-of-the-art methods [52, 53, 45], we model the lifespan of a tracker into four major subspaces in MDP state space: {birth, tracked, lost, death}. For each new set of detections, the tracker is updated using pairs with the highest affinities score (Equation 4). Each unmatched detection spawns a new tracklet; however, an unmatched tracklet is not immediately terminated, as tracklets can naturally disappear in occluded region and reappear later. We address the dynamic object inter-occlusion problem by separating a new state called “occluded” from a lost state. An object is considered occluded when covered by another object in the front with over 70% overlap. An occluded tracklet will not update its lifespan or its feature representation until it is clear from occlusion, but we still predict its 3D location using the estimated motion. Figure 4 illustrates how the

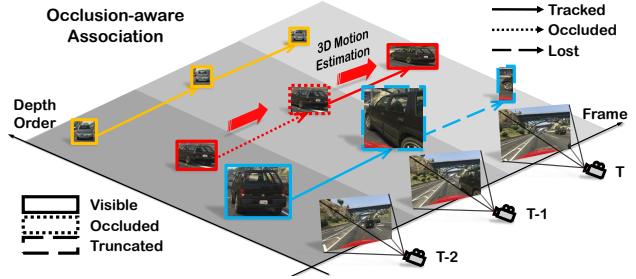


Figure 4: Illustration of Occlusion-aware association. A tracked tracklet (yellow) is visible all the time, while a tracklet (red) is occluded by another (blue) at frame $T - 1$. During occlusion, the tracklet does not update state but keep inference motion until reappearance. For a truncated or disappear tracklet (blue at frame T), we left it as lost.

occlusion-aware association works. More details of data association can be found at Appendix B.

In the next subsection, we show how to estimate that distance leveraging the associated tracklet and bounding box using a deep network.

3.5. Motion Model

Deep Motion Estimation and Update. To exploit the temporal consistency of certain vehicles, we associate the information across frames by using two LSTMs. We embed a 3D location P to a 64-dim location feature and use 128-dim hidden state LSTMs to keep track of a 3D location from the 64-dim output feature.

Prediction LSTM (P-LSTM) models dynamic object location in 3D coordinates by predicting object velocity from previously updated velocities $\hat{P}_{T-n:T-1}$ and the previous location \bar{P}_{T-1} . We use previous $n = 5$ frames of vehicle velocity to model object motion and acceleration from the trajectory. Given the current expected location of the object from 3D estimation module, the **Updating LSTM (U-LSTM)** considers both current \hat{P}_T and previously predicted location \tilde{P}_{T-1} to update the location and velocity (Figure 2(c)).

Modeling motion in 3D world coordinates naturally cancels out adverse effects of ego-motion, allowing our model to handle missed and occluded objects. The LSTMs continue to update the object state

$$\mathbf{s}_a^{(i)} = \mathbf{s}_{a-1}^{(i)} + \alpha(\mathbf{s}_a^* - \mathbf{s}_{a-1}^{(i)})$$

using the observation of matched detection state \mathbf{s}_a^* with an updating ratio $\alpha = 1 - \mathbf{A}_{deep}(\tau_a^i, \mathbf{s}_a^*)$, while assuming a linear velocity model if there is no matched bounding box. Therefore, we model 3D motion (Figure 2(d)) in world coordinates allowing occluded tracklet to move along motion plausible paths while managing the birth and death of moving objects. More details can be found at Appendix C.

Concretely, our pipeline consists of a single-frame monocular 3D object detection model for object-level pose inference and recurrent neural networks for inter-frame object association and matching. We extend the region processing to include 3D estimation by employing multi-head modules for each object instance. We introduced occlusion-aware association to solve inter-object occlusion problem. For tracklet matching, depth ordering lowers mismatch rate by filtering out distant candidates from a target. The LSTM motion estimator updates the velocity and states of each object independent of camera movement or interactions with other objects. The final pipeline produces accurate and dense object trajectories in 3D world coordinate system.

4. 3D Vehicle Tracking Simulation Dataset

It is laborious and expensive to annotate a large-scale 3D bounding box image dataset even in the presence of LiDAR data, although it is much easier to label 2D bounding boxes on tens of thousands of videos [57]. Therefore, no such dataset collected from real sensors is available to the research community. To resolve the data problem, we turn to driving simulation to obtain accurate 3D bounding box annotations at no cost of human efforts. Our data collection and annotation pipeline extend the previous works like VIPER [42] and FSV [26], especially in terms of linking identities across frames. Details on the thorough comparison to prior data collection efforts, and dataset statistics can be found in the Appendix D.

Our simulation is based on *Grand Theft Auto V*, a modern game that simulates a functioning city and its surroundings in a photo-realistic three-dimensional world. To associate object instances across frames, we utilize in-game API to capture global instance id and corresponding 3D annotations directly. In contrast, VIPER leverages a weighted matching algorithm based on a heuristic distance function, which can lead to inconsistencies. It should be noted that our pipeline is real-time, providing the potential of large-scale data collection, while VIPER requires expensive off-line processing.

5. Experiments

We evaluate our 3D detection and tracking pipeline on Argoverse Tracking benchmark [7], KITTI MOT benchmark [16] and our large-scale dataset, featuring real-world driving scenes and a wide variety of road conditions in a diverse virtual environment, respectively.

5.1. Training and Evaluation

Dataset. Our GTA raw data is recorded at 12 FPS, which is helpful for temporal aggregation. With the goal of autonomous driving in mind, we focus on vehicles closer than 150m, and also filtered out the bounding boxes whose areas are smaller than 256 pixels. The dataset is then split into

train, validation and test set with ratio 10 : 1 : 4. The KITTI Tracking benchmark provides real-world driving scenario. Our 3D tracking pipeline train on the whole training set and evaluate the performance on the public testing benchmark. The Argoverse Tracking benchmark offers novel 360 degree driving dataset. We train on the training set and evaluate the performance on the validation benchmark since the evaluation server is not available upon the time of submission.

Training Procedure. We train our 3D estimation network (Section 3.3) on each training set, separately. 3D estimation network produces feature maps as the input of ROIAlign [21]. The LSTM motion module (Section 3.5) is trained on the same set with a sequence of 10 images per batch. For GTA, all the parameters are searched using validation set with detection bounding boxes from Faster R-CNN. The training is conducted for 100 epochs using Adam optimizer with an initial learning rate 10^{-3} , momentum 0.9, and weight decay 10^{-4} . Each GPU has 5 images and each image with no more than 300 candidate objects before NMS. More training details can be found in Appendix E.

3D Estimation. We adapt depth evaluation metrics [12] from image-level to object-level, leveraging both error and accuracy metrics. Error metrics include absolute relative difference (Abs Rel), squared relative difference (Sq Rel), root mean square error (RMSE) and RMSE log. Accuracy metrics are percents of y_i that $\max(\frac{y_i}{y_i^*}, \frac{y_i^*}{y_i}) < \delta$ where $\delta = 1.25, 1.25^2, 1.25^3$. Following the setting of KITTI [16], we use orientation score (OS) for orientation evaluation.

We propose two metrics for evaluating estimated object dimension and 3D projected center position. A Dimension Score (*DS*) measures how close an object volume estimation to a ground truth. *DS* is defined as $DS = \min(\frac{V_{\text{pred}}}{V_{\text{gt}}}, \frac{V_{\text{gt}}}{V_{\text{pred}}})$ with an upper bound 1, where V is the volume of a 3D box by multiplying its dimension $l * w * h$. A Center Score (*CS*) measures distance of a projected 3D center and a ground truth. *CS* is calculated by $CS = (1 + \cos(a_{\text{gt}} - a_{\text{pd}}))/2$, with a upper bound 1, where a depicts an angular distance $((x_{\text{gt}} - x_{\text{pd}})/w_{\text{pd}}, (y_{\text{gt}} - y_{\text{pd}})/h_{\text{pd}})$, weighted by corresponding box width and height in the image coordinates.

Object Tracking. We follow the metrics of CLEAR [2], including multiple object tracking accuracy (MOTA), multiple object tracking precision (MOTP), miss-match (MM), false positive (FP), and false negative (FN), etc.

Overall Evaluation. We evaluated the 3D IoU mAP of 3D layout estimation with refined depth estimation of different tracking methods. The metric reflects the conjunction of all 3D components, dimension, rotation, and depth.

5.2. Results

3D for tracking. The ablation study of tracking performance could be found in Table 1. Adding deep feature distinguishes two near-overlapping objects, our false negative (FN) rate drops with an observable margin. With depth-order

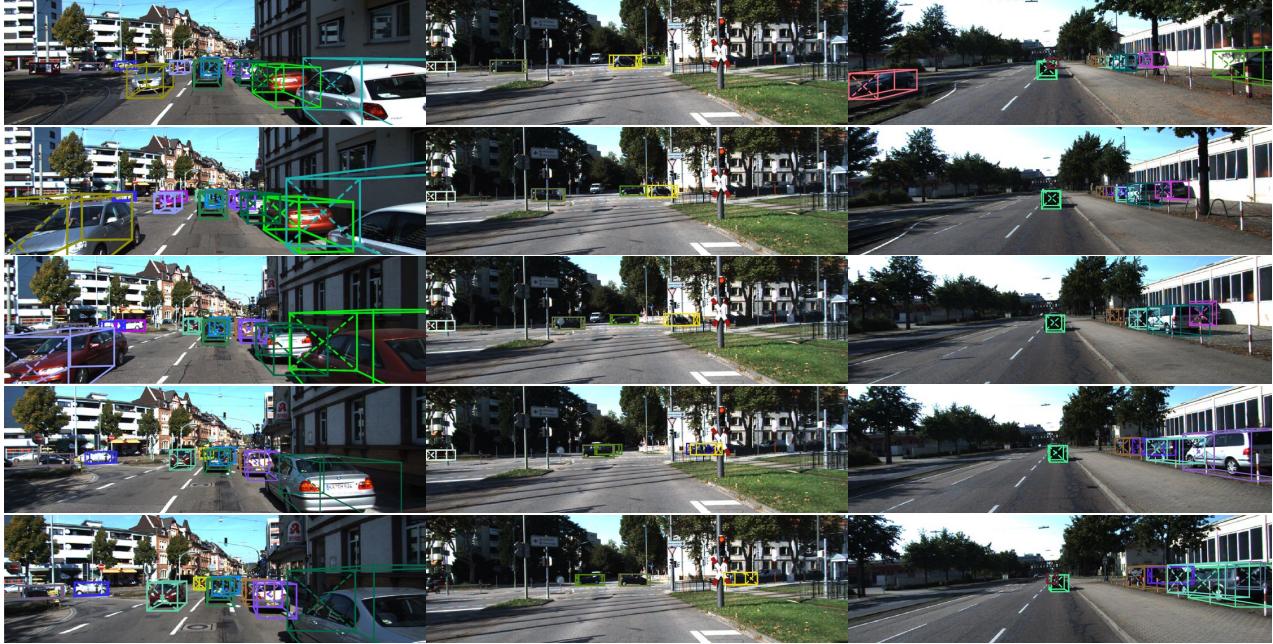


Figure 5: Qualitative results of 3D Estimation on KITTI testing set. We show predicted 3D layout colored with tracking IDs.

Motion	Deep	Order	MOTA↑	FN↓	MM↓	Ratio (%)↑
KF3D	-	-	69.616	21.298	17480	0
KF3D	-	✓	70.061	21.319	16042	8.222
KF3D	✓	-	69.843	18.547	27639	0
KF3D	✓	✓	70.126	18.177	25860	6.434
LSTM	✓	✓	70.439	18.094	23863	13.662

Table 1: Ablation study of tracking performance with different methods in GTA dataset. Motion column shows which motion model being used. KF stands for Kalman Filter. Ratio is the relative improvement in reducing the number of mismatch pairs. Using deep feature in correlation can reduce the false negative (FN) rate. Using depth-order matching and occlusion-aware association filter out relatively 6 – 8% ($\frac{(17480 - 16042)}{16042} \times 100\%$) mismatching trajectories. LSTM modeling dynamic motion benefits 3D IoU AP in Table 2.

matching and occlusion-aware association, our model filters out 6 – 8% possible mismatching trajectories. For a full ablation study, please refer to the Appendix F

Motion Modeling. We propose to use LSTM to model the vehicle motion. To analyze its effectiveness, we compare our LSTM model with traditional 3D Kalman filter (KF3D) and single frame 3D estimation using 3D IoU mAP. Table 2 shows that KF3D provides a small improvement via trajectory smoothing within prediction and observation. On the other hand, our LSTM module gives a learned estimation based on past n velocity predictions and current frame observation, which may compensate for the observation error. Our LSTM module achieves the highest accuracy among the

Method	Easy			Medium			Hard		
	AP_{3d}^{70}	AP_{3d}^{50}	AP_{3d}^{25}	AP_{3d}^{70}	AP_{3d}^{50}	AP_{3d}^{25}	AP_{3d}^{70}	AP_{3d}^{50}	AP_{3d}^{25}
None	6.13	35.12	69.52	4.93	24.25	53.26	4.05	17.26	41.33
KF3D	6.14	35.15	69.56	4.94	24.27	53.29	4.06	17.27	41.42
LSTM	7.89	36.37	73.39	5.25	26.18	53.61	4.46	17.62	41.96

Table 2: Comparison of tracking performance on 3D IoU AP 25, 50, 70 in GTA dataset. Noted that KF3D slightly improves the AP performance compare to single-frame estimation (None), while LSTM grants a clear margin. The difference comes from how we model object motion in the 3D coordinate. Instead of using Kalman filter smoothing between prediction and observation, we directly model vehicle movement using LSTMs.

	IDS ↓	FRAG ↓	FP ↓	FN ↓	MOTA ↑	MOTP ↑
2D Center	315	497	401	1435	91.06	90.36
3D Center	190	374	401	1435	91.58	90.36

Table 3: Importance of using projection of 3D bounding box center estimation on KITTI training set. We evaluate our proposed model using different center inputs c to reveal the importance of estimating projection of a 3D center. The reduction of ID Switch (IDS), track fragmentation (FRAG), and the increase of MOTA suggest that the projection of a 3D center benefits our tracking pipeline over the 2D center.

other methods with all the IoU thresholds.

3D Center Projection Estimation. We estimate the 3D location of a bounding box through predicting the projection of its center and depth, while Mousavian *et al.* [36] uses the center of detected 2D boxes directly. Table 3 shows the

Dataset	Amount	Depth Error Metric				Depth Accuracy Metric			Orientation OS ↑	Dimension DS ↑	Center CS ↑
		Abs Rel ↓	Sq Rel ↓	RMSE ↓	RMSElog ↓	$\delta < 1.25^1 \uparrow$	$\delta < 1.25^2 \uparrow$	$\delta < 1.25^3 \uparrow$			
GTA	1%	0.258	4.861	10.168	0.232	0.735	0.893	0.934	0.811	0.807	0.982
	10%	0.153	3.404	6.283	0.142	0.880	0.941	0.955	0.907	0.869	0.982
	30%	0.134	3.783	5.165	0.117	0.908	0.948	0.957	0.932	0.891	0.982
	100%	0.112	3.361	4.484	0.102	0.923	0.953	0.960	0.953	0.905	0.982
KITTI	100%	0.074	0.449	2.847	0.126	0.954	0.980	0.987	0.962	0.918	0.974

Table 4: Performance of 3D box estimation. The evaluation demonstrates the effectiveness of our model from each separate metrics. The models are trained and tested on the same type of dataset, either GTA or KITTI. With different amounts of training data in our GTA dataset, the results suggest that large data capacity benefits the performance of a data-hungry network.

Range	Method	MOTA ↑	MM ↓	#FP ↓	#FN ↓
30m	LiDAR [7]	73.02	19.75	92.80	350.50
	Ours	77.93	5.29	104.29	395.33
50m	LiDAR [7]	52.74	31.60	99.70	1308.25
	Ours	53.48	12.25	194.67	857.08
100m	LiDAR [7]	37.98	32.55	105.40	2455.30
	Ours	15.59	19.83	338.54	1572.33

Table 5: Tracking performance on the validation set of Argoverse tracking benchmark [7]. Note that the LiDAR-centric baseline[7] uses LiDAR sweeps, HD maps for evaluation.

Methods	MOTA ↑	MOTP ↑	MT ↑	ML ↓	FP ↓	FN ↓
Ours	84.52	85.64	73.38	2.77	705	4242
BeyondPixels [†] [48]	84.24	85.73	73.23	2.77	705	4247
PMBM [†] [47]	80.39	81.26	62.77	6.15	1007	5616
MDP [53]	76.59	82.10	52.15	13.38	606	7315

Table 6: Tracking performance on the testing set of KITTI tracking benchmark. Only published methods are reported. † sign means 3D information is used.

comparison of these two methods on KITTI dataset. The result indicates the correct 3D projections provides higher tracking capacity for motion module to associate candidates and reduces the ID switch (IDS) significantly.

Amount of Data Matters. We train the depth estimation module with 1%, 10%, 30% and 100% training data. The results show how we can benefit from more data in Table 4, where there is a consistent trend of performance improvement as the amount of data increases. The trend of our results with a different amount of training data indicates that large-scale 3D annotation is helpful, especially with accurate ground truth of far and small objects.

Real-world Evaluation. Besides evaluating on synthetic data, we resort to Argoverse [7] and KITTI [16] tracking benchmarks to compare our model abilities. For Argoverse, we use Faster R-CNN detection results of mmdetection [8] implementation pre-trained on COCO [32] dataset. Major results are listed in Table 5 and Table 6. For a full evaluation explanation, please refer to the Appendix F. Our monocular 3D tracking method outperforms all the published methods

on KITTI and beats LiDAR-centric baseline methods on 50m and 30m ranges of the Argoverse tracking validation set upon the time of submission.

It is worth noting that the baseline methods on Argoverse tracking benchmark leveraging HD maps, locations, and LiDAR points for 3D detection, in addition to images. Our monocular 3D tracking approach can reach competitive results with image stream only. It is interesting to see that the 3D tracking accuracy based on images drops much faster with increasing range threshold than LiDAR-based method. This is probably due to different error characteristics of the two measurement types. The farther objects occupy smaller number of pixels, leading to bigger measurement variance on the images. However, the distance measurement errors of LiDAR hardly change for faraway objects. At the same time, the image-based method can estimate the 3D positions and associations of nearby vehicles accurately. The comparison reveals the potential research directions to combine imaging and LiDAR signals.

6. Conclusion

In this paper, we learn 3D vehicle dynamics from monocular videos. We propose a novel framework, combining spatial visual feature learning and global 3D state estimation, to track moving vehicles in a 3D world. Our pipeline consists of a single-frame monocular 3D object inference model and motion LSTM for inter-frame object association and updating. In data association, we introduced occlusion-aware association to solve inter-object occlusion problem. In tracklet matching, depth ordering filters out distant candidates from a target. The LSTM motion estimator updates the velocity of each object independent of camera movement. Both qualitative and quantitative results indicate that our model takes advantage of 3D estimation leveraging our collection of dynamic 3D trajectories.

7. Acknowledgements

The authors gratefully acknowledge the support of Berkeley AI Research, Berkeley DeepDrive and MOST-107 2634-F-007-007, MOST Joint Research Center for AI Technology and All Vista Healthcare.

References

- [1] Boris Babenko, Ming-Hsuan Yang, and Serge Belongie. Visual tracking with online multiple instance learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009. 2
- [2] Keni Bernardin and Rainer Stiefelhagen. Evaluating multiple object tracking performance: the clear mot metrics. *EURASIP Journal on Image and Video Processing (JIVP)*, 2008. 6
- [3] Luca Bertinetto, Jack Valmadre, Joao F Henriques, Andrea Vedaldi, and Philip HS Torr. Fully-convolutional siamese networks for object tracking. In *European Conference on Computer Vision (ECCV)*, 2016. 2
- [4] Alex Bewley, Zongyuan Ge, Lionel Ott, Fabio Ramos, and Ben Upcroft. Simple online and realtime tracking. In *IEEE International Conference on Image Processing (ICIP)*, 2016. 15, 16
- [5] David S Bolme, J Ross Beveridge, Bruce A Draper, and Yui Man Lui. Visual object tracking using adaptive correlation filters. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010. 2
- [6] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Lioung, Qiang Xu, Anush Krishnan, Yu Pan, Giacarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. *ArXiv:1903.11027*, 2019. 2
- [7] Ming-Fang Chang, John Lambert, Patsorn Sangkloy, Jagjeet Singh, Slawomir Bak, Andrew Hartnett, De Wang, Peter Carr, Simon Lucey, Deva Ramanan, and James Hays. Argoverse: 3d tracking and forecasting with rich maps. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2, 6, 8, 14
- [8] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tian-heng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark. *ArXiv:1906.07155*, 2019. 8, 13
- [9] Wongun Choi. Near-online multi-target tracking with aggregated local flow descriptor. In *IEEE International Conference on Computer Vision (ICCV)*, pages 3029–3037, 2015. 15
- [10] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 13
- [11] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla: An open urban driving simulator. In *Conference on Robot Learning (CoRL)*, 2017. 2
- [12] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in Neural Information Processing Systems (NIPS)*, 2014. 6
- [13] Christoph Feichtenhofer, Axel Pinz, and Andrew Zisserman. Detect to track and track to detect. In *IEEE International Conference on Computer Vision (ICCV)*, 2017. 2
- [14] Davi Frossard and Raquel Urtasun. End-to-end learning of multi-sensor 3d tracking by detection. In *IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018. 15
- [15] Adrien Gaidon, Qiao Wang, Yohann Cabon, and Eleonora Vig. Virtual worlds as proxy for multi-object tracking analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 12, 13
- [16] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012. 2, 6, 8, 12, 13, 14
- [17] Ross Girshick. Fast r-cnn. In *IEEE International Conference on Computer Vision (ICCV)*, 2015. 2
- [18] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014. 2
- [19] Gultekin Gunduz and Tankut Acarman. A lightweight online multiple object vehicle tracking method. In *IEEE Intelligent Vehicles Symposium (IV)*, pages 427–432. IEEE, 2018. 15
- [20] Sam Hare, Stuart Golodetz, Amir Saffari, Vibhav Vineet, Ming-Ming Cheng, Stephen L Hicks, and Philip HS Torr. Struck: Structured output tracking with kernels. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2016. 2
- [21] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *IEEE International Conference on Computer Vision (ICCV)*, 2017. 4, 6
- [22] David Held, Sebastian Thrun, and Silvio Savarese. Learning to track at 100 fps with deep regression networks. In *European Conference on Computer Vision (ECCV)*, 2016. 2
- [23] Ju Hong Yoon, Chang-Ryeol Lee, Ming-Hsuan Yang, and Kuk-Jin Yoon. Online multi-object tracking via structural constraint event aggregation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1392–1400, 2016. 15
- [24] Soonmin Hwang, Jaesik Park, Namil Kim, Yukyung Choi, and In So Kweon. Multispectral pedestrian detection: Benchmark dataset and baseline. *Integrated Computer-Aided Engineering (ICAE)*, 2013. 12, 13
- [25] Zdenek Kalal, Krystian Mikolajczyk, and Jiri Matas. Tracking-learning-detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2012. 2
- [26] Philipp Krähenbühl. Free supervision from video games. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 6
- [27] Matej Kristan, Jiri Matas, Ales Leonardis, Michael Felsberg, Luka Cehovin, Gustavo Fernández, Tomas Vojir, Gustav Hager, Georg Nebehay, and Roman Pflugfelder. The visual object tracking vot2015 challenge results. In *IEEE International Conference on Computer Vision Workshops (ICCV Workshops)*, 2015. 2
- [28] Laura Leal-Taixé, Anton Milan, Ian Reid, Stefan Roth, and Konrad Schindler. Motchallenge 2015: Towards a benchmark for multi-target tracking. *ArXiv:1504.01942*, 2015. 12, 13

- [29] Byungjae Lee, Enkhbayar Erdenee, SongGuo Jin, Mi Young Nam, Young Giu Jung, and Phill-Kyu Rhee. Multi-class multi-object tracking using changing point detection. In *European Conference on Computer Vision Workshops (ECCV Workshops)*, 2016. 15
- [30] Philip Lenz, Andreas Geiger, and Raquel Urtasun. Followme: Efficient online min-cost flow tracking with bounded memory and computation. In *IEEE International Conference on Computer Vision (ICCV)*, pages 4364–4372, 2015. 15
- [31] Peiliang Li, Tong Qin, and andShaojie Shen. Stereo vision-based semantic 3d object and ego-motion tracking for autonomous driving. In *European Conference on Computer Vision (ECCV)*, 2018. 2
- [32] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision (ECCV)*, pages 740–755. Springer, 2014. 8, 13
- [33] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European Conference on Computer Vision (ECCV)*, 2016. 2
- [34] Will Maddern, Geoffrey Pascoe, Chris Linegar, and Paul Newman. 1 year, 1000 km: The oxford robotcar dataset. *International Journal of Robotics Research (IJRR)*, 2017. 2
- [35] Anton Milan, Laura Leal-Taixé, Ian Reid, Stefan Roth, and Konrad Schindler. Mot16: A benchmark for multi-object tracking. *ArXiv:1603.00831*, 2016. 2, 12, 13
- [36] Arsalan Mousavian, Dragomir Anguelov, John Flynn, and Jana Košecká. 3d bounding box estimation using deep learning and geometry. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 4, 7
- [37] Aljosa Osep, Wolfgang Mehner, Markus Mathias, and Bastian Leibe. Combined image- and world-space tracking in traffic scenes. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2017. 2, 15
- [38] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [39] Joseph Redmon and Ali Farhadi. Yolo9000: Better, faster, stronger. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [40] Jimmy Ren, Xiaohao Chen, Jianbo Liu, Wenxiu Sun, Jia-hao Pang, Qiong Yan, Yu-Wing Tai, and Li Xu. Accurate single stage detector using recurrent rolling convolution. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 13
- [41] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in Neural Information Processing Systems (NIPS)*, 2015. 2, 3
- [42] Stephan R Richter, Zeeshan Hayder, and Vladlen Koltun. Playing for benchmarks. In *IEEE International Conference on Computer Vision (ICCV)*, 2017. 2, 6, 12, 13
- [43] Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In *European Conference on Computer Vision (ECCV)*, 2016. 2
- [44] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M Lopez. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [45] Amir Sadeghian, Alexandre Alahi, and Silvio Savarese. Tracking the untrackable: Learning to track multiple cues with long-term dependencies. In *IEEE International Conference on Computer Vision (ICCV)*, pages 300–311, 2017. 5
- [46] Samuele Salti, Andrea Cavallaro, and Luigi Di Stefano. Adaptive appearance modeling for video tracking: Survey and evaluation. *IEEE Transactions on Image Processing (TIP)*, 2012. 2
- [47] Samuel Scheidegger, Joachim Benjaminsson, Emil Rosenberg, Amrit Krishnan, and Karl Granstrom. Mono-camera 3d multi-object tracking using deep learning detections and pmbm filtering. In *IEEE Intelligent Vehicles Symposium (IV)*, 2018. 2, 8, 15
- [48] Sarthak Sharma, Junaid Ahmed Ansari, J. Krishna Murthy, and K. Madhava Krishna. Beyond pixels: Leveraging geometry and shape cues for online multi-object tracking. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2018. 2, 8, 13, 15
- [49] Arnold WM Smeulders, Dung M Chu, Rita Cucchiara, Simone Calderara, Afshin Dehghan, and Mubarak Shah. Visual tracking: An experimental survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2014. 2
- [50] Ran Tao, Efstratios Gavves, and Arnold WM Smeulders. Siamese instance search for tracking. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [51] Longyin Wen, Dawei Du, Zhaowei Cai, Zhen Lei, Ming-Ching Chang, Honggang Qi, Jongwoo Lim, Ming-Hsuan Yang, and Siwei Lyu. Ua-detrac: A new benchmark and protocol for multi-object detection and tracking. *ArXiv:1511.04136*, 2015. 2, 12, 13
- [52] Nicolai Wojke, Alex Bewley, and Dietrich Paulus. Simple online and realtime tracking with a deep association metric. In *IEEE International Conference on Image Processing (ICIP)*, 2017. 5, 15, 16
- [53] Yu Xiang, Alexandre Alahi, and Silvio Savarese. Learning to track: Online multi-object tracking by decision making. In *IEEE International Conference on Computer Vision (ICCV)*, pages 4705–4713, 2015. 2, 5, 8, 15
- [54] Alper Yilmaz, Omar Javed, and Mubarak Shah. Object tracking: A survey. *ACM computing surveys (CSUR)*, 2006. 2
- [55] Fengwei Yu, Wenbo Li, Quanquan Li, Yu Liu, Xiaohua Shi, and Junjie Yan. Poi: multiple object tracking with high performance detection and appearance feature. In *European Conference on Computer Vision (ECCV)*, 2016. 4
- [56] Fisher Yu, Dequan Wang, Evan Shelhamer, and Trevor Darrell. Deep layer aggregation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 4
- [57] Fisher Yu, Wenqi Xian, Yingying Chen, Fangchen Liu, Mike Liao, Vashisht Madhavan, and Trevor Darrell. Bdd100k: A diverse driving video database with scalable annotation tooling. *ArXiv:1805.04687*, 2018. 2, 6

Foreword

This appendix provides technical details about our monocular 3D detection and tracking network and our virtual GTA dataset, and more qualitative and quantitative results. Section A shows the importance of learning a 3D center projection. Section B provides more details about occlusion aware association. Section C will talk about details of two LSTMs predicting and updating in 3D coordinates. Section D offers frame- and object-based statistical summaries. Section E describes our training procedure and network setting of each dataset. Section F illustrates various comparisons, including ablation, inference time and network settings, of our method on GTA, Argoverse and KITTI data.

A. Projection of 3D Center

We estimate the vehicle 3D location by first estimating the depth of the detected 2D bounding boxes. Then the 3D position is located based on the observer's pose and camera calibration. We find that accurately estimating the vehicle 3D center and its projection is critical for accurate 3D box localization. However, the 2D bounding box center can be far away from the projection of 3D center for several reasons. First, there is always a shift from the 2D box center if the 3D bounding box is not axis aligned in the observer's local coordinate system. Second, the 2D bounding box is only detecting the visible area of the objects after occlusion and truncation, while the 3D bounding box is defined on the object's full physical dimensions. The projected 3D center can be even out of the detected 2D boxes. For instance, the 3D bounding box of a truncated object can be out of the camera view. These situations are illustrated in Figure 6.

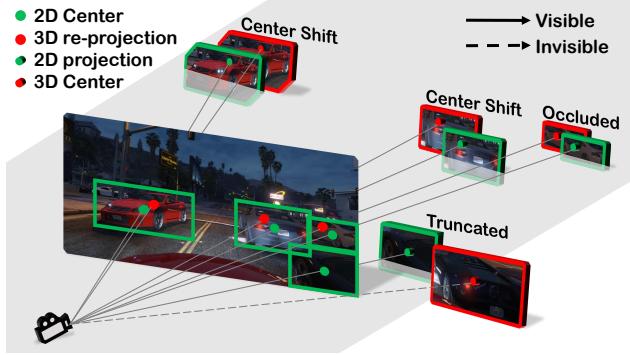


Figure 6: The illustration of the difference in 2D bounding box and 3D projected centers. A visible 2D center projection point may wrongly locate above or under the ground plane in the 3D coordinates. Other states, occluded or truncated, would inevitably suffer from the center shift problem. The center shift causes misalignment for GT and predicted tracks and harms 3D IOU AP performance.

B. Data Association Details

Due to the length limitation in the main paper, we supplement more details of data association here.

Occlusion-Aware Association. We continue to predict the 3D location of unmatched tracklets until they disappear from our tracking range (e.g. 0.15m to 100m) or die out after 20 time-steps. The sweeping scheme benefits tracker from consuming a huge amount of memory to store all the occluded tracks.

Depth-Ordering Matching. The equation mentioned in the main paper

$$A_{3D}(\tau_a, s_a) = \mathbb{1} \times \frac{\phi(M(X_{\tau_a})) \cap M(X_{s_a})}{\phi(M(X_{\tau_a})) \cup M(X_{s_a})}, \quad (6)$$

where

$$\mathbb{1} = (d_{\tau_a} - d_{s_a}) < l_{\tau_a} + w_{\tau_a} + l_{s_a} + w_{s_a}$$

denotes if the tracklets is kept after depth filtering, which can be interpreted as a loose bound

$$\sqrt{l_{\tau_a}^2 + w_{\tau_a}^2} + \sqrt{l_{s_a}^2 + w_{s_a}^2} < l_{\tau_a} + w_{\tau_a} + l_{s_a} + w_{s_a}$$

of two car within a reachable distance range, *i.e.* the diagonal length of two car. The overlapping function

$$\phi(\cdot) = \arg \min_x \{x | \text{ord}(x) < \text{ord}(x_0) \forall x_0 \in M(X_{\tau_a})\}$$

captures pixels of non-occluded tracklets region with the nearest depth order $\text{ord}(\cdot)$. Note that the $\text{ord}(\cdot)$ returns the same order if the distance is within 1 meter. Tracklets closer to the DOI occlude the overlapped area of farther ones. So each DOI matches to the tracklet with highest IOU while implicitly modeling a strong prior of depth sensation called "Occultation".

C. Motion Model

Kalman Filter. In our experiments, we compare a variant of approaches in estimating locations. 2D Kalman filter (KF2D) models

$$\{x, y, s, a, \Delta x, \Delta y, \Delta a\}$$

in 2D coordinates, where s stands for the width/high ratio and a for the area of the object. KF2D often misses tracking due to substantial dynamic scene changes from sudden ego-motion. We also propose a 3D Kalman filter (KF3D) which models

$$\{x, y, z, \Delta x, \Delta y, \Delta z\}$$

in 3D coordinates. However, Kalman filter tries to find a position in between observation and prediction location that minimize a mean of square error, while LSTM does not be bounded in such a restriction.

Deep Motion. The proposed final motion model consists of two LSTMs, one for predicting and the other for updating. The latter is in charge of refining 3D locations.

Predicting LSTM (P-LSTM) learns a velocity in 3D $\{\Delta x, \Delta y, \Delta z\}$ using hidden state based on the previously n updated velocities $\hat{P}_{T-n:T-1}$ and location \bar{P}_{T-1} and predicts a smoothed step toward an estimated location \tilde{P}_T . We embed velocity for the past $n = 5$ consecutive frames with an initial value of all 0 to model object motion and acceleration from the trajectory. As the next time step comes, the velocity array updates a new velocity into the array, and discards the oldest one.

For Updating LSTM (U-LSTM) module, we encode both predicted location \hat{P}_{T-1} and current single-frame observation \hat{P}_T into two location features. By concatenating them together, we obtained a 128-dim feature as the input of U-LSTM. The LSTM learns a smooth estimation of location difference from the observation \hat{P}_T to the predicted location \hat{P}_{T-1} and updates the 3D location \bar{P}_T and velocity $\dot{P}_{T-n+1:T}$ from hidden state leveraging previous observations.

Both LSTM modules are trained with ground truth location projected using the depth value, a projection of 3D center, camera intrinsic and extrinsic using 2 losses. The L1 loss $L_1 = |\dot{P}_T - \dot{P}_{T-1}|$ reduces the distance of estimated and ground truth location. The linear motion loss $L_{linear} = (1 - \cos(\dot{P}_T, \dot{P}_{T-1})) + L1(\dot{P}_T, \dot{P}_{T-1})$ focuses on the smooth transition of location estimation.

D. Dataset Statistics

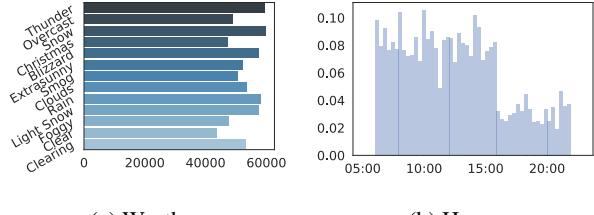
To help understand our dataset and its difference, we show more statistics.

Dataset Comparison. Table 7 demonstrates a detailed comparison with related datasets, including detection, tracking, and driving benchmarks. KITTI-D [16] and KAIST [24] are mainly detection datasets, while KITTI-T [16], MOT15 [28], MOT16 [35], and UA-DETRAC [51] are primarily 2D tracking benchmarks. The common drawback could be the limited scale, which cannot meet the growing demand for training data. Compared to related synthetic dataset, Virtual KITTI [15] and VIPER [42], we additionally provide fine-grained attributes of object instances, such as color, model, maker attributes of vehicle, motion and control signals, which leaves the space for imitation learning system.

Weather and Time of Day. Figure 7 shows the distribution of weather, hours of our dataset. It features a full weather cycle and time of a day in a diverse virtual environment. By collecting various weather cycles (Figure 7a), our model learns to track with a higher understanding of environments. With different times of a day (Figure 7b), the network handles changeable perceptual variation.

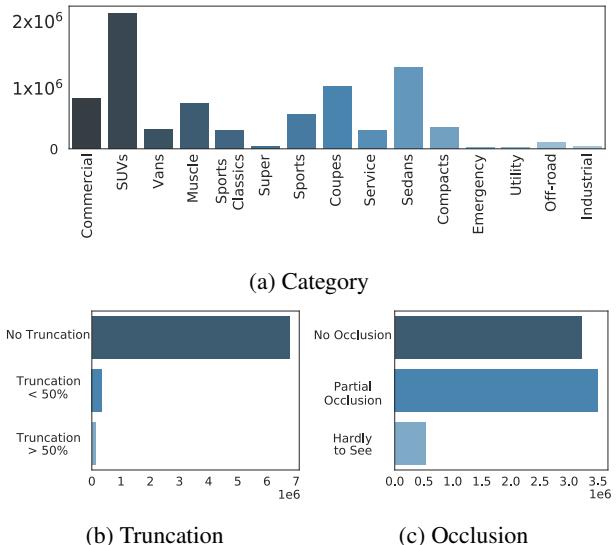
Number of Instances in Each Category. The vehicle diversity is also very large in the GTA world, featuring 15

fine-grained subcategories. We analyzed the distribution of the 15 subcategories in Figure 8a. Besides instance categories, we also show the distribution of occluded (Figure 8c) and truncated (Figure 8b) instances to support the problem of partially visible in the 3D coordinates.



(a) Weather (b) Hours

Figure 7: The statistics of scene in our dataset.



(a) Category (b) Truncation (c) Occlusion

Figure 8: The statistics of object in our dataset.

Dataset Statistics. Compared to the others, our dataset has more diversity regarding instance scales (Figure 9a) and closer instance distribution to real scenes (Figure 9b).

Examples of Our Dataset. Figure 10 shows some visual examples in different time, weather and location.

E. Training Details

We apply different training strategies to our model to optimize the performance on different datasets.

Training Procedure. We train on 4 GPUs (effective mini-batch size is 20). To keep the original image resolution can be divided by 32, we use 1920×1080 resolution for GTA, 1920×1216 for Argoverse, and 1248×384 for KITTI. The anchors of RPN in Faster R-CNN span 4 scales and 3 ratios, in order to detect the small objects at a distance. For GTA,

Dataset	Task	Object	Frame	Track	Boxes	3D	Weather	Occlusion	Ego-Motion
KITTI-D [16]	D	C,P	7k	-	41k	✓	-	✓	-
KAIST [24]	D	P	50k	-	42k	-	✓	✓	-
KITTI-T [16]	T	C	8k	938	65k	✓	-	✓	✓
MOT15-3D [28]	T	P	974	29	5632	✓	✓	-	-
MOT15-2D [28]	T	P	6k	500	40k	-	✓	-	-
MOT16 [35]	T	C,P	5k	467	110k	-	✓	✓	-
UA-DETRAC [51]	D,T	C	84k	6k	578k	-	✓	✓	-
Virtual KITTI [15]	D,T	C	21k	2610	181k	✓	✓	✓	✓
VIPER [42]	D,T	C,P	184k	8272	5203k	✓	✓	✓	✓
Ours	D,T	C,P	688k	325k	10068k	✓	✓	✓	✓

Table 7: Comparison to related dataset for detection and tracking (Upper half: real-world, Lower half: synthetic). We only count the size and types of annotations for training and validation (D=detection, T=tracking, C=car, P=pedestrian). To our knowledge, our dataset is the largest 3D tracking benchmark for dynamic scene understanding, with control signals of driving, sub-categories of object.

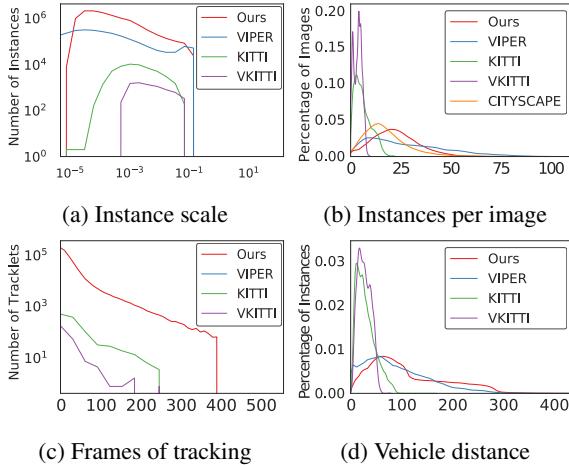


Figure 9: Statistical summary of our dataset in comparison of KITTI [16], VKITTI [15], VIPER [42], and Cityscapes[10]

we estimate both 2D detection and 3D center on the Faster RCNN pr-trained on Pascal VOC dataset. For KITTI, we use RRC [40] detection results following a state-of-the-art method, BeyondPixels [48]. For Argoverse, we use Faster RCNN detection results of mmdetection [8] implementation pre-trained on COCO [32] dataset. Since both the results do not come with projected 3D centers, we estimate the center c from the network in Section 3.3 of the main paper.

F. Experiments

Visual Odometry vs Ground Truth. To analyze the impact of ego-motion on tracking, we compare tracking results based on poses from GT camera and Mono-VO approach in Table 8. The proposed method maintains similar MOTA using estimated camera pose. Note that 3D tracking is compensated with GT ego-motion in the text.

DATA	POSE	MOTA↑	MOTP↑	MM↓
KITTI	VO	71.079	88.233	0.784
	GT	71.215	88.210	0.642

Table 8: Ablation study of tracking performance with estimated (VO) or ground truth (GT) camera pose in KITTI training set. We find that VO can lead to similar results as GT.

Inference Time. The total inference time is 203 ms on a single P100 GPU (see Table 9 for details). Note that the inference time on KITTI benchmark focus only on the non-detection part (92 ms).

	Detection	3D Estimation	Tracking
Run Time	0.111	0.062	0.030

Table 9: Inference Time (second) of our proposed framework on KITTI tracking benchmark. Note that elapsed time for object detection is not included in the specified runtime of KITTI benchmark.

Tracking policy. In the association, we keep all tracklets until they disappear from our tracking range (e.g. 10m to 100m) or die out after 20 time-steps. Table 10 shows the ablation comparison of using location and max-age. Since the max-age affects the length of each tracklet, we prefer to keep the tracklet a longer period (*i.e.* > 15) to re-identify possible detections using 3D location information.

Data Association Weights. During our experiment, we use different weights of appearance, 2D overlap, and 3D overlap for corresponding methods. We select weight ratios based on the results of our validation set.

For None and 2D, we use $w_{app} = 0.0, w_{2D} = 1.0, w_{3D} =$



Figure 10: Examples of our dataset.

0.0. For 3D related methods, we give $w_{3D} = 1.0$ and $w_{2D} = 0.0$. For Deep related methods, we times w_{2D} or w_{3D} with 0.7 and 0.3 for w_{app} .

Tracking Performance on KITTI dataset. As mentioned in the main paper, we resort to Argoverse [7] and KITTI [16] tracking benchmarks to compare our model abilities in the

Max-Age	Location	IDS↓	FRAG↓	FP↓	FN↓	MOTA↑	MOTP↑
30	-	513	710	341	1571	89.93	90.40
20	-	373	571	341	1571	90.51	90.40
-	✓	316	513	341	1571	90.74	90.40
30	✓	290	488	341	1571	90.85	90.40
20	✓	261	459	341	1571	90.97	90.40

Table 10: Comparison of performance using different location and maxage in our tracking pipeline. The hyphen - represents not using the information in tracking.

real-world scenario. Our monocular 3D tracking method outperforms all the published methods upon the time of submission. Results are listed in Table 11.

Benchmark	MOTA ↑	MOTP ↑	MT ↑	ML ↓	FP ↓	FN ↓
Ours	84.52	85.64	73.38	2.77	705	4242
BeyondPixels [†] [48]	84.24	85.73	73.23	2.77	705	4247
PMBM [†] [47]	80.39	81.26	62.77	6.15	1007	5616
extraCK [19]	79.99	82.46	62.15	5.54	642	5896
MCMOT-CPD [29]	78.90	82.13	52.31	11.69	316	6713
NOMT* [9]	78.15	79.46	57.23	13.23	1061	6421
MDP [53]	76.59	82.10	52.15	13.38	606	7315
DSM [14]	76.15	83.42	60.00	8.31	578	7328
SCEA [23]	75.58	79.39	53.08	11.54	1306	6989
CIWT [37]	75.39	79.25	49.85	10.31	954	7345
NOMT-HM* [9]	75.20	80.02	50.00	13.54	1143	7280
mbodSSP [30]	72.69	78.75	48.77	8.77	1918	7360

Table 11: Tracking performance on the testing set of KITTI tracking benchmark. Only published methods are reported. † sign means 3D information is used.

Ablation Study. The ablation study of tracking performance can be found in Table 12. From the ablation study, we observe the huge performance improvement, especially in track recall (TR) for a 8% and Accuracy (MOTA) for a 18% relative improvement, from 2D to 3D, which indicates 3D information helps to recover unmatched pairs of objects. We implement our 2D baselines deriving from Sort [4], Deep Sort [52] tracker for inter-frame ROIs matching in 2D. Modeling motion in 3D coordinate is more robust in linking candidates than 2D motion. Moreover, adding deep feature distinguishes two near-overlapping objects, our false negative (FN) rate drops with an observable margin. With depth-order matching and occlusion-aware association, our model filters out 6 – 8% possible mismatching trajectories.

Qualitative results. We show our evaluation results in Figure 11 on our GTA test set and in Figure 12 on Argoverse validation set. The method comparison on GTA shows that our proposed 3D Tracking outperforms the strong baseline which using 3D Kalman Filter. The light dashed rectangular represents the ground truth vehicle while the solid rectangular stands for the predicted vehicle in bird’s eye view. The qualitative results on Argoverse demonstrate our monocular 3D tracking ability under day and night road conditions without the help of LiDAR or HD maps. The figures are best

viewed in color.

Evaluation Video. We have uploaded a demo video which summaries our main concepts and demonstrates video inputs with estimated bounding boxes and tracked trajectories in bird’s eye view. We show the comparison of our method with baselines in both our GTA dataset and real-world videos.

Motion	Deep	Order	MOTA \uparrow	MOTP \uparrow	TP \uparrow	TR \uparrow	MM \downarrow	NM \downarrow	RM \downarrow	FP \downarrow	FN \downarrow
-	-	-	5.056	67.726	79.445	71.102	6.31	1.636	5.839	25.997	62.637
2D [4]	-	-	57.042	82.794	86.905	77.058	3.739	1.064	3.418	6.085	33.134
2D* [52]	✓	-	65.892	81.86	90.607	87.186	4.796	4.096	3.898	10.099	19.213
3D	-	-	69.616	84.776	85.947	84.202	1.435	0.511	1.283	7.651	21.298
3D	✓	-	69.843	84.523	90.242	87.113	2.269	1.798	1.672	9.341	18.547
3D	-	✓	70.061	84.845	85.952	84.377	1.317	0.403	1.184	7.303	21.319
3D	✓	✓	70.126	84.494	90.155	87.432	2.123	1.717	1.512	9.574	18.177
LSTM	✓	✓	70.439	84.488	90.651	87.374	1.959	1.547	1.37	9.508	18.094
-	-	-	4.609	71.064	74.941	83.594	14.535	3.735	13.834	31.326	49.531
2D [4]	-	-	42.748	81.878	70.003	84.292	9.172	2.077	8.755	16.731	31.350
2D* [52]	✓	-	48.518	81.479	66.381	88.222	7.270	2.683	6.738	21.223	22.989
3D	-	-	54.324	83.914	64.986	90.869	3.032	0.799	2.820	23.574	19.070
3D	✓	-	54.855	83.812	68.235	91.687	2.087	0.776	1.864	25.223	17.835
3D	-	✓	54.738	84.045	64.001	90.623	3.242	0.533	3.102	22.120	19.899
3D	✓	✓	55.218	83.751	68.628	92.132	1.902	0.723	1.688	25.578	17.302
LSTM	✓	✓	55.150	83.780	69.860	92.040	2.150	0.800	1.920	25.220	17.470

Table 12: Ablation study of tracking performance with different methods and features in GTA dataset. The upper half of the table shows evaluation of targets within 100 meters while the lower within 150 meters. Modeling motion in 3D coordinate with the help of our depth estimation is more robust in linking candidate tracklets than 2D motion. Using deep feature in correlation can reduce the false negative (FN) rate. Using depth-order matching and occlusion-aware association filter out relatively 6 – 8% mismatching trajectories. Replacing 3D kalman filter with LSTM module to re-estimate the location using temporal information.

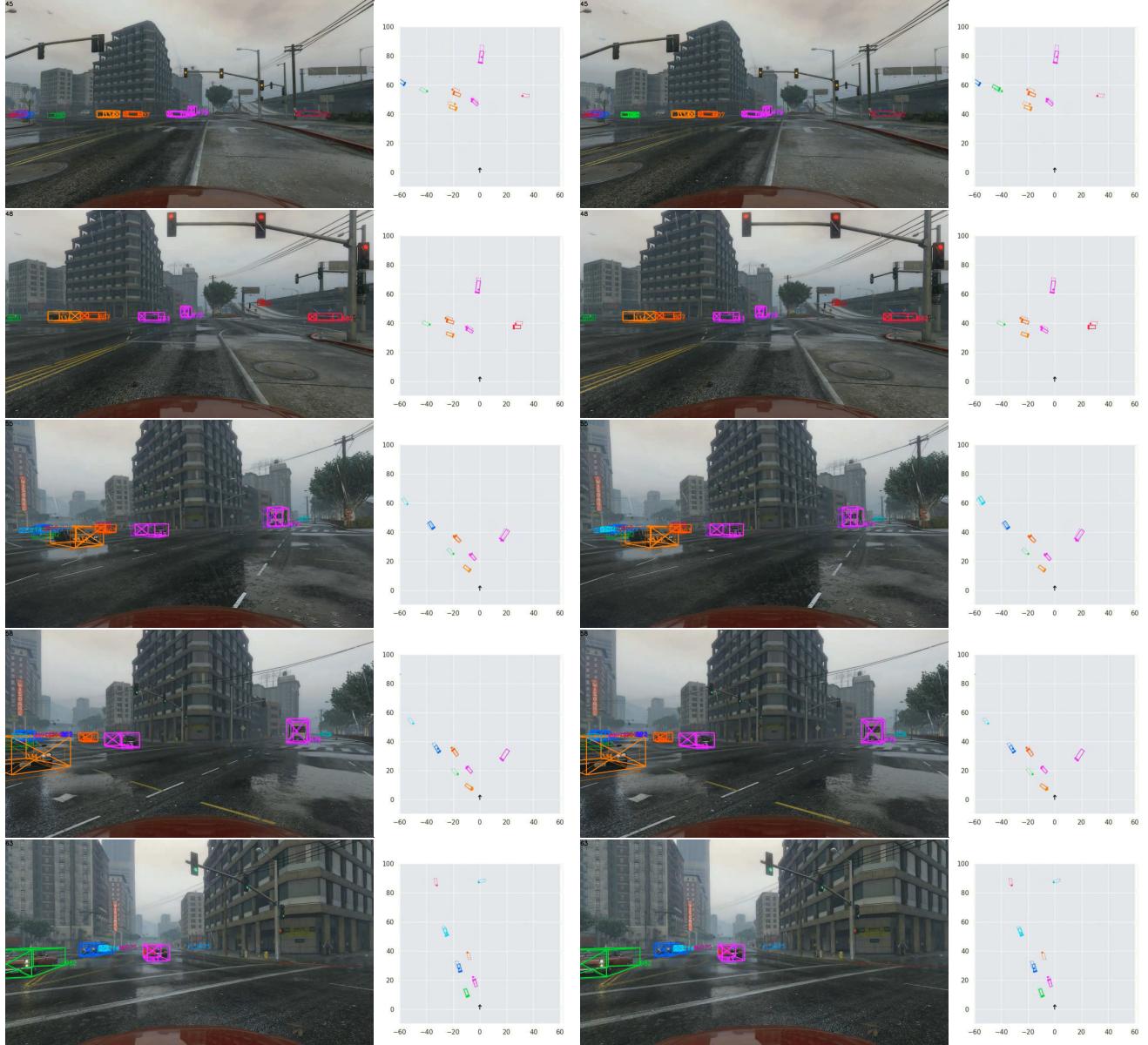


Figure 11: Qualitative comparison results of 3D Kalman and our method.

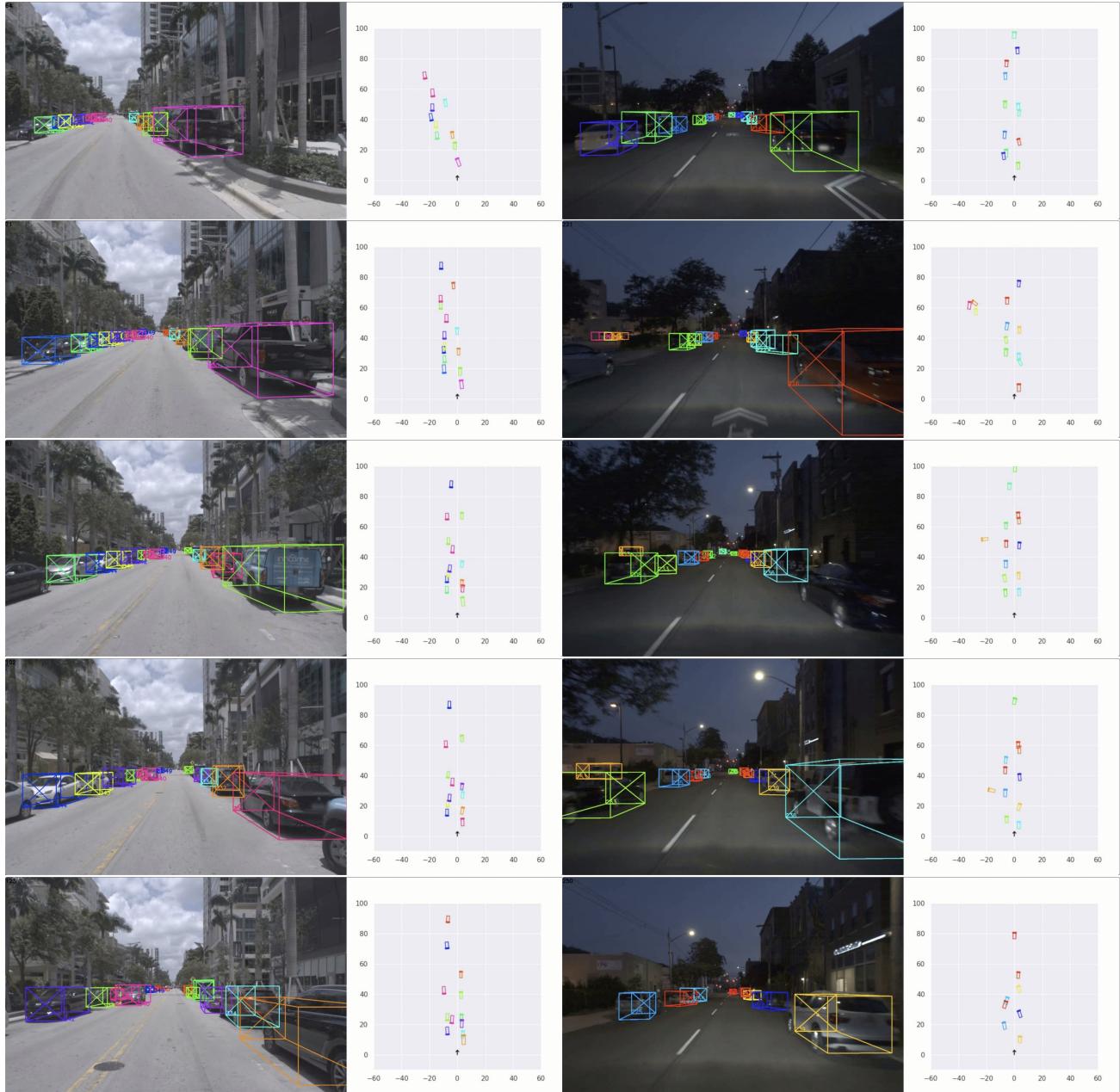


Figure 12: Qualitative results of our method on Argoverse.