# Model-free Vehicle Tracking and State Estimation in Point Cloud Sequences

Ziqi Pang, Zhichao Li, Naiyan Wang
TuSimple
{ziqipang.z, leeisabug, winsty}@gmail.com

*Abstract*— **Estimating the states of surrounding traffic participants stays at the core of autonomous driving. In this paper, we study a novel setting of this problem: model-free single-object tracking (SOT), which takes the object state in the first frame as input, and jointly solves state estimation and tracking in subsequent frames. The main purpose for this new setting is to break the strong limitation of the popular "detection and tracking" scheme in multi-object tracking. Moreover, we notice that shape completion by overlaying the point clouds, which is a by-product of our proposed task, not only improves the performance of state estimation but also has numerous applications. As no benchmark for this task is available so far, we construct a new dataset *LiDAR-SOT* and corresponding evaluation protocols based on the Waymo Open dataset [29]. We then propose an optimization-based algorithm called *SOTracker* involving point cloud registration, vehicle shapes, correspondence, and motion priors. Our quantitative and qualitative results prove the effectiveness of our SOTracker and reveal the challenging cases for SOT in point clouds, including the sparsity of LiDAR data, abrupt motion variation, etc. Finally, we also explore how the proposed task and algorithm may benefit other autonomous driving applications, including simulating LiDAR scans, generating motion data, and annotating optical flow. The code and protocols for our benchmark and algorithm are available at `https://github.com/TuSimple/LiDAR_SOT/`. A video demonstration is at `https://www.youtube.com/watch?v=BpHixKs91i8`.**

## I. INTRODUCTION

Autonomous driving calls for accurate state estimation of surrounding objects, including their positions, orientations, etc. Among all the sensors available for state estimation, LiDAR is unique in providing accurate 3D sensing. Nevertheless, the complexity of road scenarios and sparsity of point clouds pose great difficulties for state estimation. Therefore, estimating the states of objects using LiDAR data has long been valuable yet challenging for autonomous driving research.

The mainstream method for object tracking is to utilize the so-called "detection and tracking" paradigm. It mostly focuses on the association of bounding boxes across frames. The accuracy of tracking then heavily relies on the performance of detection, since most tracking algorithms directly accept the bounding boxes from detection with minor modifications. Current state-of-the-art object detection algorithms are mostly *model-based*: the algorithms carry out the detection task with their recognition power gained from learning on large scale training sets. The benchmarks nowadays [29][4][10][2] also mainly focus on this *model-based* MOT setting. However, we propose to pay attention to the *model-free* Single-Object Tracking (SOT) setting, which targets on some overlooked
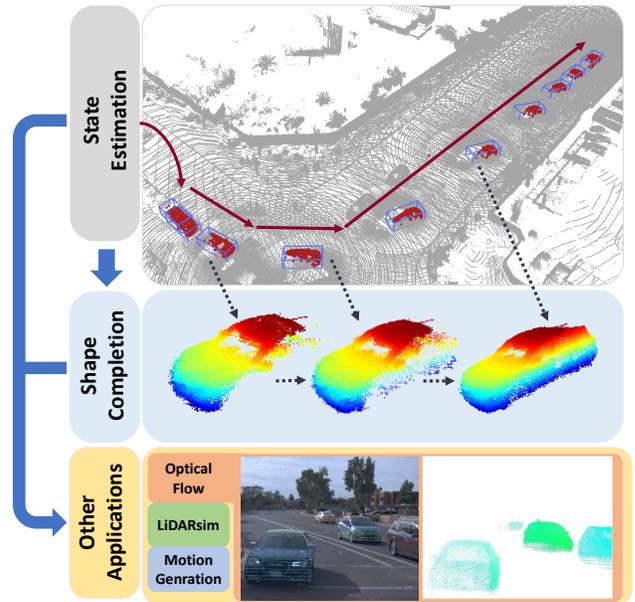


Fig. 1: Our vehicle tracking system and its applications. **Top**. Tracking results visualized at the interval of 2 seconds (20 frames). **Middle**. Shape completion by overlaying point clouds. **Bottom**. Applications using the estimated states and shapes, such as optical flow annotation. More usages are described in Sec. VI.

applications and challenges in MOT.

First, model-based detection is prone to fail on sparse point clouds and unfamiliar objects (see Fig. 2). By being model-free, the algorithms can generalize to such cases, since it is not assumed that similar objects have been met before. The task of *model-free* SOT is then useful for both online and offline applications from two aspects: 1) For offline settings, model-free SOT can augment the training set for model-based methods by providing cheap annotations or simulated data [21]. 2) For online settings, SOT serves as a back-up system for model-based methods on their failures (see Fig. 2). This feature is mandatory for a high-level autonomous system which requires no single failure point.

Second, model-free SOT algorithms rely on richer information for tracking, such as raw point clouds, compared to MOT methods that mostly use bounding boxes only. Therefore, the SOT algorithms specialize at providing more subtle state information with even centimeter level accuracy. This property is important for many applications in autonomous driving, such as generating first order motion information. It
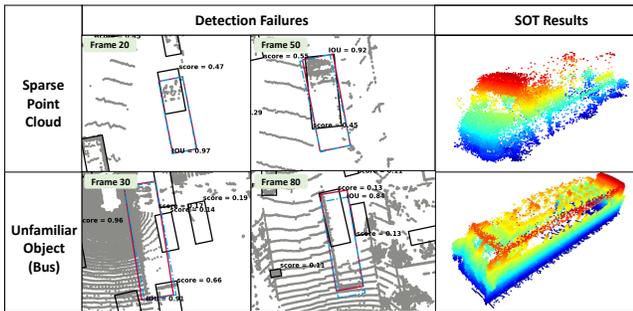
Fig. 2: **Left:** The detection algorithm (PointPillars [19]) fails on certain cases, including false negatives and wrong size. But our tracker still has high quality after a period of tracking, with only requiring an initial bounding box as input. We sample two frames from a continuous period of failed detection for demonstration. *Gray*: Point Cloud; *Red*: The Ground Truth Box; *Dashed Blue*: SOT Box; *Black*: Detection Boxes. **Right:** We visualize the quality of SOT with the aggregated point clouds, following [15][32]. Our algorithm can provide accurate tracking and shape aggregation results.

is discussed in Sec. VI-A.

For a target object, the input to SOT algorithms is its initial bounding box in the first frame, which may come from either human annotations or highly confident detections. The algorithms then estimate the object's states in every subsequent frame, as is in the top of Fig. 1. During solving state estimation, we notice the value of its by-product: aggregating point clouds along the tracklets completes denser shapes of objects (shown in the middle of Fig. 1). In our experiments, leveraging these shapes improves the performance of tracking. This utility of shape aggregation also benefits a series of other applications in autonomous driving. For example, with the generated shapes and motion information, we can create vehicle bank for creating virtual LiDAR scenes [21], optical flow annotations [22] (our result visualized in Fig. 1, bottom row), etc. These tasks so far mostly require manual labeling or artificial CAD models, so automating the process of directly using raw point cloud data could greatly improve their efficiency and scalability. The details of such downstream applications are discussed in Sec. VI.

To summarize, our contributions are as follows:

- We propose a novel LiDAR-based object tracking task[1]. To facilitate further study, we create a new benchmark called LiDAR-SOT based on the WOD dataset [29].
- We propose a model-free SOT algorithm called SO-Tracker. It involves optimizing for point cloud registration, object shapes and motions.
- We evaluate our SOTracker on LiDAR-SOT and discuss the potential usages in autonomous driving, such as motion data generation, LiDAR scan simulation [21], and optical flow annotation [22].

---

[1]We focus on vehicle so far, and are working on extending to other types of objects.

## II. RELATED WORK

### A. Datasets and Benchmarks

Many existing autonomous driving benchmarks, including Waymo Open Dataset (WOD) [29], Argoverse [4], KITTI [10], and nuScenes [2], have rich point clouds and tracklets for object tracking. However, they mainly adopt the MOT setting while our task demands the SOT setting, and not all the tracklets for MOT evaluation suits the requirements for SOT. So we design a test bed specially for SOT evaluation. Among the mentioned benchmarks, WOD [29] is superior in scale and comprehensiveness, so we build the benchmark LiDAR-SOT based on the data of WOD.

To evaluate the SOT algorithms, we inherit the idea from VOT [18] for the tracking performance on 2D images, concerning both the accuracy and robustness. As for shape completion, we consider it for evaluation because of its usages on downstream applications. Specifically, we refer to ShapeNet [3] and use Chamfer-Distance(CD) as the metric.

### B. State Estimation and Tracking in Point Cloud

State estimation in raw point clouds is challenging because it involves associating the LiDAR points with the target object (a.k.a. tracking) and then accurately estimate the state.

Many methods adopt a model-based approach to discover the object in LiDAR points. [25][8][35][17] fit the vehicles with boxes and [30][6] utilizes occupancy grid for detection. In the era of deep learning, the learned neural networks replace the manually designed models. P2B [28] generates and verifies bounding boxes guided by a learned PointNet++ [27], Zou *et al.* [41] manages to fuse RGB information using frustums, Zarzar *et al.* [39] refines the detection results by association, while [37][20] apply convolutional neural networks(CNN) on the Bird-Eye-View images of LiDAR, and Giancola *et al.* [11] learns to complete the partial point cloud, and then searches for it.

These model-based approaches heavily rely on detection or shape completion in every single frame. So point cloud sparsity and unfamiliar objects could easily lead to failures. However, our model-free setting seeks to tackle this since it does not need to fit a model and can consider multi-frame information from LiDAR sequences.

Different from the aforementioned works involving searching for the target object in raw LiDAR data, some assume the point cloud could be well-segmented or the segmentation information is available. Hence their focus is state estimation in relatively clean point cloud data. Some representative methods are as follows, Held *et al.* [15] minimizes an energy function on the inconsistency between observed point cloud and object states. Groß *et al.* [13] registers the motion between point cloud snapshots. Goforth *et al.* [12] learns to complete the shapes of vehicles and jointly decode the states. The problem of these methods lies in the assumption of well-segmented point cloud, which is rare in practice. Therefore, by only providing the first frame annotation, our setting enforces the algorithms to face the challenge of solving tracking and state estimation simultaneously.

Many model-free works [23][36][9][26][1][31][33] mainly focus on the problem of MOT. They first operate object detection leveraging either the motion information or assuming that LiDAR points draw the boundaries of objects. Then they associate the detected bounding boxes according to the estimated motion information or motion models. Therefore, these methods differ from our objectives for using SOT.

The only work we found falls into our setting is Ushani *et al.* [32]. Dating back to that period, there is no large-scale dataset fits the evaluation needs, thus [32] only tests on a private dataset, which makes the comparisons infeasible. This is just the reason to propose our new LiDAR-SOT benchmark, which aims to promote the studies in this area.

### C. Point Cloud Completion

Following the paradigm of PCN [38], researchers can complete point clouds with neural networks trained from artificial CAD models in graphics datasets like ShapeNet [3]. [12][11] adopt the same idea for vehicles. However, these methods could hardly scale up as CAD models are expensive and limited to some categories. Recently, Gu *et al.* [14] discards the need for CAD models by using multi-view sensors. But multi-view sensor information are also not universal for real-world scenarios.

Considering the mentioned drawbacks of supervised methods, completing point cloud by tracking interests us most. [7][5][16] jointly use LiDAR and camera for this task. [15][32] are closest to our objective in using LiDAR only: Held *et al.* [15] aggregates the observed point clouds, while Ushani *et al.* [32] directly optimizes a point cloud model.

## III. LiDAR-SOT Benchmark

### A. Problem Formulation

Suppose the dataset has $N$ tracklets $\mathcal{T}^{1:N} = \{\mathcal{T}^1, \mathcal{T}^2, ... \mathcal{T}^N\}$ with lengths $L^{1:N} = \{L^1, L^2, ..., L^N\}$. We refer to the $i$-th frame in the $j$-th sequence $\mathcal{T}^j$ as $\mathcal{T}_i^j$, which contains the point cloud $P_i^j$ and a bounding box $\widetilde{B}_i^j$. The box $\widetilde{B}_i^j \triangleq \{\widetilde{C}_i^j, \widetilde{S}_i^j\}$ describes the size and state of an object. In our benchmark, $\widetilde{C}_i^j$ are the length, width and height of a 3D cuboid, and the state $\widetilde{S}_i^j \triangleq [\widetilde{x}_i^j, \widetilde{y}_i^j, \widetilde{z}_i^j, \widetilde{\theta}_i^j]$ represents the 3D coordinates of the object's center and the heading. We only consider 4DOF instead of 6DOF because roll and pitch deviations are usually aligned to the road for autonomous driving. For a tracklet $\mathcal{T}^j$, it also contains the shape ground-truth of target object $\widetilde{M}^j$.

For an arbitrary tracklet $\mathcal{T}^j$, the algorithms are fed with point cloud $\{P_{1:L^j}^j\}$ and the ground-truth bounding box $\widetilde{B}_1^j$ in the first frame. The outputs are the estimated states $\{S_{2:L^j}^j\}$ and completed object shape $M^j$. Note that we assume the size of an object remains unchanged, thus it keeps the value as $\widetilde{C}_1^j$ in subsequent frames. This assumption is generally held for rigid objects such as vehicles.

The algorithms can run under either online or offline setting. The online setting restricts the algorithms to access point cloud data $\{P_{1:L^j}^L\}$ by order, and output the estimation results before getting the next frame input. The offline setting, in comparison, poses no constraints. In this paper,

we concentrate on the online setting because it fits both on-the-road autonomous driving and offline applications.

### B. Discussion with Closely Related Tasks

In this section, we highlight the differences between our proposed task and several relevant tasks.
*Multiple-object Tracking:* As mentioned above, MOT relies on detections and mainly focuses on associating them. Moreover, MOT usually does not modify or generate new bounding boxes, while our SOT task creates bounding boxes according to the point cloud and motion information.
*Shape Completion:* Generally speaking, shape completion complements the partial point cloud in a single frame under the supervised setting. It needs paired training data to learn generative models, while ours could utilize temporal LiDAR data to complete the shapes in an unsupervised manner.
*Point Cloud Registration:* Registration focuses on matching the point clouds between two frames, while more previous frames are accessible in our task. Consequently, rich priors could be incorporated to improve the performance, as demonstrated later in the experiments, the performance of our state estimation is significantly higher than using ICP only.

### C. Dataset Construction

To facilitate the study on the aforementioned task, we propose a new dataset called LiDAR-SOT. It is largely based on the most comprehensive and largest point cloud dataset: Waymo Open Dataset (WOD) [29]. Considering that the vehicles are the most abundant rigid objects for autonomous driving, we only concern the type of vehicles for now. The WOD dataset contains full annotations of 3D bounding boxes and their associations across frames for the evaluation of MOT, thus, enabling our evaluation for state estimation in point cloud sequences. As for the evaluation for shape completion, since WOD contains no ground truth, we create pseudo-ground-truth shapes (PGTs) for each tracklet by selecting and aggregating the LiDAR points from every frame guided by the ground-truth bounding boxes.

Our algorithm in the paper requires no additional training sets. However, we still include a training set for LiDAR-SOT by inheriting all the sequences from WOD training set and leaving users with freedom of training/validation split. As for the test set, we select the suitable tracklets from WOD validation set to ensure meaningful evaluation. In total, 1121 vehicle tracklets are included based on the following criteria:
**Tracklet length.** Short sequences may result in unstable evaluation and they are also not challenging enough for SOT tasks. We therefore filter out the tracklets less than 100 frames.
**Initial number of points.** The numbers of LiDAR points in the beginning frames are crucial for the meaningful initialization, where VOT [18] sets up a warm-up period of 10 frames. Enough number of LiDAR points is also necessary for human annotators to provide a initial bounding box with high quality for offline applications. We therefore take both aspects into consideration, and only keep the tracklets whose initial 10 frames all have more than 20 points.
**Vehicle mobility.** We exclude the tracklets in which the vehicle stays static all the time, as these cases have trivial
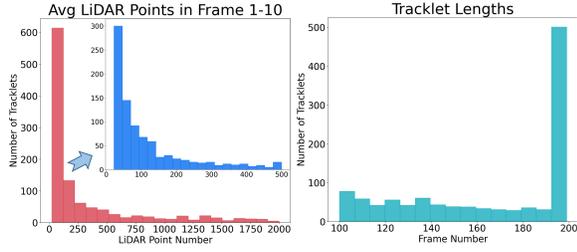
Fig. 3: Statistics of LiDAR-SOT. **Left**. Average per-frame LiDAR point number in the first 10 frames. **Right**. Distribution of tracklet lengths.

solutions and bring negative biases to the evaluation. Note that we still keep the tracklets where the vehicles start or end with static states.

As visualized in the left of Fig. 3, although we guarantee the number of LiDAR points in the beginning, the long-tail distribution still indicates the challenges in point cloud sparsity. Moreover, the right of Fig. 3 illustrates the difficulty of LiDAR-SOT in long tracklets (note that the original WOD [29] has a maximum of 200 frames).

For future analyses, we follow the paradigm of KITTI [10] and divide the test set into easy, medium and hard subsets. Intending on the difficulty of point cloud sparsity, we compute the average point number of each tracklet's first 10 frames, then equally split the tracklets into three sets according to the average point number. The specific thresholds are 38.2 and 808.3 points per frame.

### D. Evaluation Metrics

Regarding the output on tracklet $\mathcal{T}^j$: states $\{S^j_{2:L^j}\}$ and shape $M^j$, we use multiple metrics for thorough evaluation. Our focus is on the quality of estimated states in SOT and shape completion.

**State Estimation:** For state estimation, we are interested in how long the algorithm could track objects and its quality. Consequently, inspired by Visual Object Tracking (VOT) [18] work, we propose two metrics for single pass evaluation.

**Accuracy:** We denote $\text{IOU}(B^j_i, \widetilde{B}^j_i)$ as the IOU of the estimated and ground-truth bounding boxes. We then propose a metric reflecting the state estimation accuracy by averaging the IOU on each frame as:

$$\text{Acc} = \frac{1}{\sum^N_{k=1} L^k} \sum^N_{k=1} \sum^{L^k}_{j=1} \text{IOU}(B^k_j, \widetilde{B}^k_j) \qquad (1)$$

**Robustness:** Given a pre-defined failure threshold $t$, we let the tracking length $l^j(t)$ be the first frame number before IOU drops lower than $t$ in tracklet $\mathcal{T}^j$. This can be served as measurement for robustness. However, as a single threshold $t$ may not be comprehensive, we compute the area under $t$ - $\text{Rob}(t)$ curve by Eq. 2. In implementation, we approximate the integral by sampling $t$ from 0 to 1 with intervals of 0.05.

$$\text{Rob}(t) = \frac{\sum^N_{k=1} l^k(t)}{\sum^N_{k=1} L^k}, \quad \text{Rob} = \int^1_0 \text{Rob}(t) \mathrm{d}t \qquad (2)$$

**Shape Completion** As the PGTs comes from overlaying all the point cloud from tracklets, their densities are non-uniform.

To avoid unstable evaluation, we downsample both the PGTs $\{\widetilde{M}^{1:N}\}$ and our completion $\{M^{1:N}\}$ to the resolution of 5cm during evaluation. For shape $M^i$, we use Chamfer Distance (CD) to measure its distance to $\widetilde{M}^i$ as Eq. 3. Then the Shape metric is the average of CD on all tracklets.

$$\text{CD}^i = \frac{1}{|M^i|} \sum_{a \in M^i} \min_{b \in \widetilde{M}^i} ||a-b||_2 + \frac{1}{|\widetilde{M}^i|} \sum_{b \in \widetilde{M}^i} \min_{a \in M^i} ||a-b||_2 \qquad (3)$$

### IV. SOTRACKER

In this section, we elaborate our proposed SOTracker. The SOT task is essentially a cross frame registration task, but a series of priors could be explored for enhancement. Briefly speaking, our method estimates the relative motions between frames using registration and various motion priors, and aggregate them. Along the process of tracking, we also combine the shapes on the fly for further improvement.

### A. Pipeline

We take an arbitrary tracklet $\mathcal{T}$ for example to explain the algorithm. The notation $\Delta S_k = S_k - S_{k-1}$ is the motion between $S_{k-1}$ and $S_k$. Its entries $\Delta S_k = [\Delta x_k, \Delta y_k, \Delta z_k, \Delta \theta_k]$ are the 3D translation and rotation of heading. In the following sections, we use $O_k$ to denote the estimated LiDAR points on the surface of the object in frame $\mathcal{T}_k$, in comparison with the raw point cloud $P_k$.

The pipeline of SOTracker is in Algo. 1, where it successively estimates the state in each frame with function `FrameEstimate` in Algo. 2. In frame $\mathcal{T}_k$, SOTracker first solves $S_k$ with `FrameEstimate`; then forms shape $M_k$ by combining previous shape $M_{k-1}$ and object point cloud $O_k$ together; it eventually updates the motion information $\Delta S_{prior}$, which is used for initializing the next frame's state estimation. In `FrameEstimate`, the function initializes the state $\widehat{S}_k$ from adding the $\Delta S_{prior}$ onto previous state $S_{k-1}$. Then it iteratively selects the relevant LiDAR points $\widehat{O}_k$ and optimizes the state $S_k$ towards a loss function $\mathcal{L}$.

Some other details are clarified as follows:

*1) :* We implement $\Delta S_{prior}$ as the moving average of motions and use $\alpha = 0.5$ for the update in Algo. 1.

*2) :* `PointInBox`$(P, S, C, \gamma)$ selects LiDAR points from $P$ that locates in the cuboid with center $S$ and size $\gamma C$. We set different $\gamma$ values for different usages. $\gamma_{in} = 1.5$ is for the selection inside an estimation loop, while $\gamma_{aft} = 1.1$ is for updating shapes.

*3) :* The shape update happens in `OverlayShape`, where we initialize it using all the points in $\widetilde{S}_1$. As taking LiDAR points from every frame brings heavy burden for computation, we only combine new point clouds at key frames (that is every 5 frames in our implementation).

*4) :* The second frame is special for `FrameEstimate` as previous motion information $\Delta S_{prior}$ is not available. During estimating $S_2$, we let $\Delta S_{prior} = 0$ and use a larger search region for `PointInBox` by setting $\gamma_{in} = 3.0$. After this, we initialize the motion prior by $\Delta S_{prior} \leftarrow \Delta S_2$.

**Algorithm 1** Algorithm for tracklet $\mathcal{T}$

---

**Input:** Point cloud: $P_{1:L}$, bounding box $\widetilde{B}_1$, Size Scaling: $\gamma$, Motion Prior Update Factor: $\alpha$
**Output:** vehicle states: $S_{2:L}$, shape $M_L$
1: Frame number $k \leftarrow 2$
2: $\Delta S_{prior} \leftarrow 0$
3: $O_1 \leftarrow \texttt{PointInBox}(P_1, \widetilde{S}_1, \widetilde{C}_1, 1.0)$
4: $M_1 \leftarrow O_1$
5: **for** $k <= L$ **do**
6:    $S_k \leftarrow \texttt{FrameEstimate}(P_k, O_{k-1}, \Delta S_{prior}, S_{k-1},$
             $\widetilde{C}_1, M_1, M_{k-1}, \gamma_{in})$
7:    $O_k \leftarrow \texttt{PointInBox}(P_k, S_K, \widetilde{C}_1, \gamma_{aft})$
8:    $M_k \leftarrow \texttt{OverlayShape}(O_k, S_k, M_k, k)$
9:    $\Delta S_{prior} \leftarrow \alpha \Delta S_{prior} + (1 - \alpha) \Delta S_k$
10:   $k \leftarrow k + 1$
11: **end for**

---

**Algorithm 2** `FrameEstimate` on frame $\mathcal{T}_k$

---

**Input:** Point Cloud $P_k$ and $O_{k-1}$, Motion Prior: $\Delta S_{prior}$, State: $S_{k-1}$, Size $\widetilde{C}_1$, Shape $M_1, M_{k-1}$, parameters $\gamma_{in}$
**Output:** State Estimation $S_k$
1: $\widehat{S}_k \leftarrow S_{k-1} + \Delta S_{prior}$
2: **for** Iteration Number $<$ MaxIter **do**
3:    $\widehat{O}_k \leftarrow \texttt{PointInBox}(P_k, \widehat{S}_k, \widetilde{C}_1, \gamma_{in})$
4:    $S_k^* \leftarrow \underset{S_k}{\arg\min} \mathcal{L}(S_k, \widehat{O}_k, O_{k-1}, S_{k-1}, M_{k-1})$
5:    $\widehat{S}_k \leftarrow S_k^*$
6: **end for**
7: $S_k \leftarrow \widehat{S}_k$,

---

### B. Objective Function

The objective function is in Eq. 4. It is the weighted sum of the four terms: ICP Term, Shape Term, Motion Consistency Term, and Motion Prior Term. We set their weights by $[1, 1, 0.1, 0.1]$. Eventually, we use the BFGS solver from scipy [34] to optimize $\mathcal{L}$.

$$\mathcal{L} = w_I \mathcal{L}_{ICP} + w_S \mathcal{L}_{Shape} + w_{MC} \mathcal{L}_{MC} + w_{MP} \mathcal{L}_{MP} \quad (4)$$

In the sequel, we use the notation $\Delta S \otimes (P, S)$ to denote apply transformation $\Delta S = [\Delta x, \Delta y, \Delta z, \Delta \theta]$ to the point cloud $P$, relative to the coordinate system centered at $S = [x, y, z, \theta]$. $\Delta S \otimes (P, S)$ equals to the first 3 rows of Eq. 5.

$$\begin{bmatrix} \cos\Delta\theta & -\sin\Delta\theta & 0 & x + \Delta x - x\cos\Delta\theta + y\sin\Delta\theta \\ \sin\Delta\theta & \cos\Delta\theta & 0 & y + \Delta y - x\sin\Delta\theta - y\cos\Delta\theta \\ 0 & 0 & 1 & \Delta z \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} P \\ 1 \end{bmatrix}$$
$$(5)$$

We explain each term in the following sections.

*1) ICP Term:* Registering the point cloud $\widehat{O}_k$ and $O_{k-1}$ indicates the motion between these two frames. Therefore, we propose a registration term $\mathcal{L}_{ICP}$ as in Eq. 6. Note that $p_{S_{k-1}}$ indicates applying transformation to $p$ relative to $S_{k-1}$, as Eq. 5. The details for finding associated point pairs $\mathcal{A}_I^{O_{k-1}, \widehat{O}_k}$ are in Sec. IV-C.3.

$$\mathcal{L}_{ICP} = \frac{1}{|\mathcal{A}_I^{O_{k-1}, \widehat{O}_k}|} \sum_{(p,q) \in \mathcal{A}_I^{O_{k-1}, \widehat{O}_k}} ||\Delta S_k \otimes (p, S_{k-1}) - q||_2^2$$
$$(6)$$

*2) Shape Term:* Registering the point cloud in neighboring frames could fail in long term tracking due to drifting or shifting in point cloud distribution. We remedy this by adding regularization between the shapes of objects and current frame point clouds. Specifically, we penalize the inconsistency between the current frame point cloud $\widehat{O}_k$ and the latest shape $M_{k-1}$ by registering them. In the implementation, we first use the first frame point cloud as initialization, since the initial bounding box has high quality and the related point clouds are accurate. Then we aggregate new LiDAR points along the process of tracking, and use this updated shape to compute the loss function. Formally, the Shape Term is computed as Eq. 7. $\mathcal{A}_S^{M_{k-1}, \widehat{O}_k}$ is the set of associate point pairs, explained in Sec. IV-C.3.

$$\mathcal{L}_{Shape} = \frac{1}{|\mathcal{A}_S^{M_{k-1}, \widehat{O}_k}|} \sum_{(p,q) \in \mathcal{A}_S^{M_{k-1}, \widehat{O}_k}} ||(S_k - \widetilde{S}_1) \otimes (p, \widetilde{S}_1) - q||_2^2$$
$$(7)$$

*3) Motion Consistency Term:* The objects cannot be treated as point mass, so their movements have to follow certain constraints. For simplicity, we enforce the most general constraint: their directions should be consistent with their own headings. The corresponding term is computed as Eq. 8, where $v = [\Delta x_k, \Delta y_k]$ represents the velocity.

$$\mathcal{L}_{MC} = \left|\left| ||v||_2 \cos(\frac{\theta_{k-1} + \theta_k}{2}) - \Delta x_k \right|\right|_2^2 \quad (8)$$

*4) Motion Prior Term:* To combine the motion model into the optimization framework, we design the "Motion Prior" term. It penalizes the inconsistency between the optimization results $\Delta S_k$ and the prediction from motion models $\Delta S_{prior}$. This term encourages the consensus between the latest optimization and trajectory histories. The computation of Motion Prior Term is as Eq. 9.

$$\mathcal{L}_{MP} = ||\Delta S - \Delta S_{prior}||_2^2 \quad (9)$$

### C. Design Choices

We identify several key implementation factors that are helpful for the performance of state estimation.

*1) Ground Removal:* The ground points are irrelevant with the vehicle's states and shapes, and they usually confuse the algorithm. So we remove the ground points using [40] as a pre-processing operation.

*2) Subshape:* Due to the sparsity of LiDAR point clouds, it is challenging to only use single frame for registration. In the ICP term, we can enhance the performance by overlaying the LiDAR points in nearby frames. In the implementation, we overlay additional 2 frames $O_{k-3:k-2}$ onto $\widehat{O}_{k-1}$ in Eq. 6.

*3) LiDAR Point Association:* As a prerequisite of computing the ICP Term and Shape Term in Sec. IV-B, the LiDAR points have to be associated according to the rules of nearest
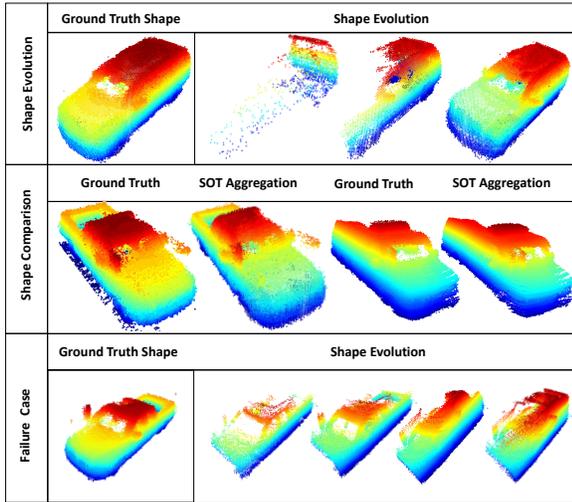
Fig. 4: Visualization SOT results. **1**. The evolution of shapes during the tracking process. **2**. Comparison between the results from SOT and ground truth. **3**. Visualization of failure cases.

neighbors. As for the ICP term described in Sec. IV-B.1, we follow the common approach and construct the associated point set $\mathcal{A}_I^{O_{k-1},\widehat{O}_k}$ as Eq. 10. In the equation, the function $\mathrm{NN}(p, \mathcal{B})$ means finding the nearest neighbor for point $p$ in the point set $\mathcal{B}$.

$$\left\{ (x, y) \mid y = \mathrm{NN}(\Delta S_k \otimes (x, S_{k-1}), \widehat{O}_k), \forall x \in O_{k-1} \right\} \quad (10)$$

However, as for the Shape Term in Sec. IV-B.2, the point clouds $M_{k-1}$ and $\widehat{O}_k$ may differ significantly during the process of tracking. The outliers of association lead to deviated solutions and degrade the state estimation. Therefore, we use RANSAC for Shape Term to exclude the outlier pairs. The result $\mathcal{A}_S^{M',\widehat{O}_k}$ is computed as Eq. 11.

$$\mathrm{RANSAC}(\left\{ (x, y) \mid y = \mathrm{NN}((S_k - \widetilde{S}_1) \otimes (x, \widetilde{S}_1), \widehat{O}_k), \forall x \in M_{k-1} \right\}) \quad (11)$$

### D. Limitations and Future Work

For simplicity, our LiDAR-SOT is currently limited to the rigid type of objects: vehicles. Extending it to more general types of objects, such as cyclists and pedestrians, requires the modeling for the periodically and continuously changing shapes. Using multiple templates or implicit functions [24] for their shapes could be the solution.

Combining SOT methods with the multiple detections in different frames also poses great challenges, as the quality of detection bounding boxes may vary across different frames. Therefore, we have to model the quality of them and combine our optimization-based algorithm into the MOT framework. Also, the disappearance or occlusion of an object need to be considered in the setting. We leave these two issues, and will try to solve them in future work.

## V. RESULTS AND ANALYSES

### A. Quantitative Results

Table I evaluates the contributions of each individual factor. We start from the baseline of using ICP and motion model, then add the other terms one by one. When ablating

the optimization factors from Sec. IV-B, we remove them separately from the strongest SOTracker to show their benefits.

Compared to the baseline of using ICP and motion model only (ICP+MP), adding the Shape-Term (S(1), using the first frame point cloud as shape in the Shape Term) significantly improves the performance. During this process, aggregating shapes along the process of tracking (S(All), using aggregated shape in the Shape Term) assist the state estimation even further. On top of the point cloud related factors, the regularization on motion consistency (MC) is helpful especially on medium and hard cases. As for the design choices in Sec. IV-C, removing any of them greatly affects both the accuracy and robustness, where Ground Removal and RANSAC are the most critical ones.

### B. Comparison with Model-based Methods

*1) Compare with Model-based Baseline:* To demonstrate the effectiveness of our model-free approach, we compare our method with one common model-based baseline: utilizing Kalman filters to associate the detections, and then directly output the associated bounding boxes. Specifically, we use *Point Pillars* [19] for detection. Then we change the detection bounding boxes' sizes to ground-truth following LiDARSim [21] for fair comparison with SOT methods, since other SOT methods have ground-truth size in the first frame.

Contrasting Row 4 and "Det" in Tab. I, we observe that our SOTracker can approach the performance of "Det" on the easy set without external training data, even slightly better on robustness, while has large disadvantages on the medium and hard sets. The results are reasonable because model-based methods could benefit from the large-scale training set, while model-free methods simply don't use it, but should be more robust on the rare and unseen cases.

*2) Combine Model-free and Model-based Methods:* Furthermore, we try to incorporate the detections into our framework by *Detection Term*. Briefly speaking, it selects the detection with the largest IOU with the motion model prediction, while enforcing a minimum score of 0.5 and IOU of 0.1. If no bounding box fulfills the requirement, this term is discarded. The term is illustrated in Eq. 12, where $S_{\mathrm{D}}$ is the state for selected detection. This hybrid method is the "Det+SOT" in Tab. I. Compared to the "Det" row, "Det+SOT" consistently outperforms, proving our point.

$$\mathcal{L}_{Det} = \|\Delta S_k - (S_{\mathrm{D}} - S_{k-1})\|_2^2 \quad (12)$$

### C. Analyses

*1) **Vehicle Shape Visualization:*** Following the approach in [15], [32], we visualize the aggregated point clouds in Fig. 4 to indicate the quality of tracking. Moreover, the vehicle banks constructed during this process are also beneficial to a series of tasks discussed in Sec. VI. In Fig. 4 and Fig. 1, the evolution of point clouds and the quality of shapes prove that our algorithm can produce high quality state information and 3D shapes. This further poses the necessity and effectiveness of using LiDAR sequences for creating 3D shapes.

| Type | Method | All | | | Easy | | | Medium | | | Hard | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Acc↑ | Rob↑ | Shape↓ | Acc↑ | Rob↑ | Shape↓ | Acc↑ | Rob↑ | Shape↓ | Acc↑ | Rob↑ | Shape↓ |
| Model-free | ICP+MP | 0.3927 | 0.3734 | 0.1836 | 0.5156 | 0.4997 | 0.1067 | 0.3578 | 0.3340 | 0.1792 | 0.2828 | 0.2638 | 0.2652 |
| | ICP+MP+S(1) | 0.5931 | 0.5113 | 0.1226 | 0.7386 | 0.6801 | 0.0612 | 0.5380 | 0.4522 | 0.1218 | 0.4779 | 0.3721 | 0.1851 |
| | ICP+MP+S(All) | 0.6027 | 0.5348 | 0.1239 | 0.7479 | 0.6992 | 0.0591 | 0.5613 | 0.4819 | 0.1261 | 0.4789 | 0.3946 | 0.1867 |
| | **ICP+MP+S(All)+MC** | **0.6146** | **0.5467** | **0.1164** | **0.7496** | **0.7002** | **0.0589** | **0.5743** | **0.4910** | **0.1173** | **0.4959** | **0.4224** | **0.1727** |
| | w/o GM | 0.5715 | 0.4994 | 0.1284 | 0.6924 | 0.6450 | 0.0676 | 0.5320 | 0.4512 | 0.1279 | 0.4662 | 0.3783 | 0.1889 |
| | w/o RANSAC | 0.5707 | 0.5130 | 0.1276 | 0.7016 | 0.6615 | 0.0705 | 0.5185 | 0.4480 | 0.1239 | 0.4728 | 0.4049 | 0.1883 |
| | w/o SUB | 0.6025 | 0.5367 | 0.1221 | 0.7436 | 0.6937 | 0.0595 | 0.5574 | 0.4813 | 0.1199 | 0.4814 | 0.4079 | 0.1871 |
| Model-based | Det | 0.7574 | 0.6184 | 0.0714 | 0.7850 | 0.6918 | 0.0589 | 0.7409 | 0.5993 | 0.0780 | 0.7423 | 0.5508 | 0.0774 |
| | Det+SOT | 0.7690 | 0.6329 | 0.0644 | 0.8017 | 0.7092 | 0.0544 | 0.7597 | 0.6169 | 0.0637 | 0.7398 | 0.5587 | 0.0753 |

TABLE I: Ablation studies. ↑(↓) means the performance is better with larger(smaller) values. Optimization terms: **ICP** – ICP Term; **S** – Shape Term, **S(1)** – only using the first frame point cloud, **S(all)** – using the updated shape; **MC** – Motion Consistency Term, **MP** – Motion Prior Term. Implementation details: **GM** – ground removal, **RANSAC** – using RANSAC, **SUB** – subshape. **Det** and **Det+SOT** are covered in Sec. V-B.

*2) Shape Completion ≠ State Estimation:* Although in Table. I, better shapes generally come along with better state estimation performance, we argue that these two are not equivalent tasks. The LiDAR points are aggregated until the last frame in the experiments of Table I, but this is not the optimal approach and leads to worse shapes on the failures of state estimation, just as the last row in Fig. 4. Therefore, ideal shape completion needs to reason whether and how to aggregate the point clouds, thus differs from mere state estimation. We will investigate this direction in future study.

### D. Challenging Scenarios

*1) Point Cloud Sparsity:* In Table I, the performance gaps among the three difficulty levels reflect the challenge of sparse point clouds. The hard set has less LiDAR points, which causes inaccurate registration at the beginning and leads to lower performance. In long-term tracking, the ubiquitous cases of occlusion or distant objects indicate the needs for handling this challenge.

*2) Drifting:* Accumulated error is well known in SLAM and commonly exists in odometry based approaches. The errors can accumulate and cause drifting over time if without absolute observation or extra constraints like loop closure. We alleviate this by adding the shape terms, which treats the shape as "Map" and enforce its consistency with the observed point cloud. Although the terms are effective in many cases as in Table I, drifting still exists as in Fig. 1 and requires future research on it.

*3) Abrupt Motion Change:* Accurate motion initialization is crucial since it provides the initial point set correspondence. Although current motion models are able to adapt to most cases, they still struggle where the objects change the motions greatly in short periods, as the incorrectly associated point clouds can quickly lead to failures. We will explore more advanced motion prediction methods, and close the loop of state estimation and motion prediction in the future.

## VI. EXPERIMENTS ON DOWNSTREAM APPLICATIONS

### A. Motion Ground Truth Generation

SOTracker can estimate the motion of a designated vehicle in wild point cloud sequences. Therefore, it is capable of providing motion data for prediction and planning systems at a low cost. As in Fig. 5, the deviation of generated motion
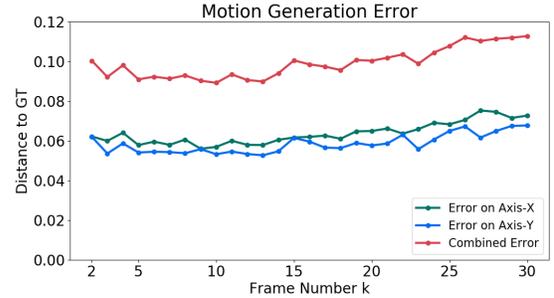


Fig. 5: Motion generation results. The average distance (meters) between the motion from ground-truth and our algorithm at frame k.

is less than 10cm within 30 frames, which is even smaller than the size of a finger. With Argoverse [4] providing 30 data points for the motion prediction task, we believe that our system meets the standard.

### B. Optical Flow Annotation

With the estimated shapes $M$ and states $S_{k:k+1}$, we can annotate optical flows for the tracked objects in corresponding images $I_{k:k+1}$. For each point $m \in M$, we compute its location $(qx_k, qy_k)$ on $I_k$ by moving the shape according to $S_k$, then project it onto the image plane. In the same way, we compute $(qx_{k+1}, qy_{k+1})$ on image $I_{k+1}$, then the optical flow for pixel $(qx_k, qy_k)$ is $(qx_{k+1} - qx_k, qy_{k+1} - qy_k)$. In Fig. 1, our annotated optical flow can act as high quality pseudo-ground-truth. Such an application is infeasible for the single frame point cloud (the left image), as it contains neither enough points nor the motion information for annotating the optical flow. In KITTI [22], researchers require all three of LiDAR data, CAD models and ground-truth 3D bounding boxes to infer the optical flows of vehicles, while our algorithm only requires an initial bounding box and raw LiDAR scenes. Therefore, our method can reduce the cost of optical flow annotation and scale it up.

### C. Simulating LiDAR Scans Using Vehicle Bank

In LiDARsim [21] researchers propose to simulate single LiDAR scans for autonomous driving by creating object shapes with manually labeled 3D locations. As mentioned in LiDARsim [21], augmenting the training data with the simulation can improve the performance of perception systems. In

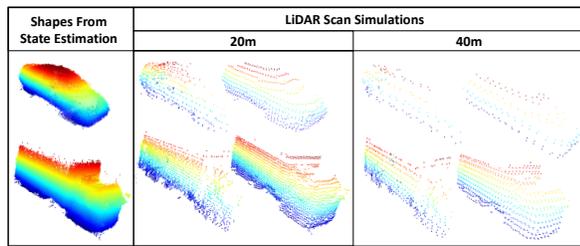| Shapes From State Estimation | LiDAR Scan Simulations | |
|---|---|---|
| | 20m | 40m |

Fig. 6: LiDAR scans created using our completed shapes and physical simulation. The visualization includes two types of vehicles: car and truck. The LiDAR simulation fits different angles and different distances (20m and 40m).

Fig. 6 we visualize the results of our LiDAR scan simulation on different types of vehicles, angles, and distances. For both types of vehicles, we can build up high quality shapes and simulated scans, which proves the effectiveness of our model-free SOT approach. Moreover, unlike LiDARsim [21], which requires exhausted annotation of objects, we only require an initial bounding box and raw point clouds. Thus, we make the process cheaper and more scalable, and produce unlimited LiDAR data with minimum human annotation.

## VII. CONCLUSION

In this paper, we calls for the attention on a new setting of object tracking in point cloud sequences: model-free single-object tracking. This new setting targets on the drawbacks of model-based detection and MOT methods, and can benefit a series of applications. To facilitate the research, we construct the LiDAR-SOT from WOD and propose a strong algorithm called SOTracker. According to the qualitative and quantitative results, our method not only achieves strong baseline on state estimation, but also generalizes its estimated states and 3D shapes to the applications of simulating LiDAR scans, optical flow annotation, and motion data generation.

## REFERENCES

[1] A. Azim and O. Aycard. Detection, classification and tracking of moving objects in a 3d environment. In *IV*, 2012.

[2] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom. nuScenes: A multimodal dataset for autonomous driving. In *CVPR*, 2020.

[3] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, et al. ShapeNet: An information-rich 3D model repository. *arXiv:1512.03012*, 2015.

[4] M. Chang, J. Lambert, P. Sangkloy, J. Singh, S. Bak, A. Hartnett, D. Wang, P. Carr, S. Lucey, D. Ramanan, and J. Hays. Argoverse: 3D tracking and forecasting with rich maps. In *CVPR*, 2019.

[5] D. Christie, C. Jiang, D. Paudel, and C. Demonceaux. 3D reconstruction of dynamic vehicles using sparse 3D-laser-scanner and 2D image fusion. In *ICIC*, 2016.

[6] R. Danescu, F. Oniga, and S. Nedevschi. Modeling and tracking the driving environment with a particle-based occupancy grid. *T-ITS*, 12(4):1331–1342, 2011.

[7] M. H. Daraei, A. Vu, and R. Manduchi. Velocity and shape from tightly-coupled LiDAR and camera. In *IV*, 2017.

[8] M. Darms, P. Rybski, and C. Urmson. Classification and tracking of dynamic objects with multiple sensors for autonomous driving in urban environments. In *IV*, 2008.

[9] A. Dewan, T. Caselitz, G. D. Tipaldi, and W. Burgard. Motion-based detection and tracking in 3d lidar scans. In *ICRA*, 2016.

[10] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the KITTI vision benchmark suite. In *CVPR*, 2012.

[11] S. Giancola, J. Zarzar, and B. Ghanem. Leveraging shape completion for 3D siamese tracking. In *CVPR*, 2019.

[12] H. Goforth, X. Hu, M. Happold, and S. Lucey. Joint pose and shape estimation of vehicles from LiDAR data. *arXiv:2009.03964*, 2020.

[13] J. Groß, A. Osep, and B. Leibe. Alignnet-3D: Fast point cloud registration of partially observed objects. In *3DV*, 2019.

[14] J. Gu, W.-C. Ma, S. Manivasagam, W. Zeng, Z. Wang, Y. Xiong, H. Su, and R. Urtasun. Weakly-supervised 3D shape completion in the wild. *arXiv:2008.09110*, 2020.

[15] D. Held, J. Levinson, S. Thrun, and S. Savarese. Combining 3D shape, color, and motion for robust anytime tracking. In *RSS*, 2014.

[16] C. Jiang, D. Christie, D. P. Paudel, and C. Demonceaux. High quality reconstruction of dynamic objects using 2D-3D camera fusion. In *ICIP*, 2017.

[17] R. Kaestner, J. Maye, Y. Pilat, and R. Siegwart. Generative object detection and tracking in 3D range data. In *ICRA*, 2012.

[18] M. Kristan, J. Matas, A. Leonardis, T. Vojír, R. P. Pflugfelder, G. Fernández, G. Nebehay, F. Porikli, and L. Cehovin. A novel performance evaluation methodology for single-target trackers. *TPAMI*, 38(11):2137–2155, 2016.

[19] A. H. Lang, S. Vora, H. Caesar, L. Zhou, J. Yang, and O. Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *CVPR*, 2019.

[20] W. Luo, B. Yang, and R. Urtasun. Fast and Furious: Real time end-to-end 3D detection, tracking and motion forecasting with a single convolutional net. In *CVPR*, 2018.

[21] S. Manivasagam, S. Wang, K. Wong, W. Zeng, M. Sazanovich, S. Tan, B. Yang, W. Ma, and R. Urtasun. LiDARsim: Realistic LiDAR simulation by leveraging the real world. In *CVPR*, 2020.

[22] M. Menze and A. Geiger. Object scene flow for autonomous vehicles. In *CVPR*, 2015.

[23] F. Moosmann and C. Stiller. Joint self-localization and tracking of generic objects in 3d range data. In *ICRA*, 2013.

[24] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *CVPR*, 2019.

[25] A. Petrovskaya and S. Thrun. Model based vehicle tracking for autonomous driving in urban environments. In *RSS*, 2008.

[26] F. Pomerleau, P. Krüsi, F. Colas, P. Furgale, and R. Siegwart. Long-term 3d map maintenance in dynamic environments. In *ICRA*, 2014.

[27] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. PointNet++: Deep hierarchical feature learning on point sets in a metric space. In *NeurIPS*, 2017.

[28] H. Qi, C. Feng, Z. Cao, F. Zhao, and Y. Xiao. P2B: Point-to-box network for 3D object tracking in point clouds. In *CVPR*, 2020.

[29] P. Sun, H. Kretzschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine, V. Vasudevan, W. Han, J. Ngiam, H. Zhao, A. Timofeev, S. Ettinger, M. Krivokon, A. Gao, A. Joshi, Y. Zhang, J. Shlens, Z. Chen, and D. Anguelov. Scalability in perception for autonomous driving: Waymo Open Dataset. *arXiv:1912.04838*, 2019.

[30] G. Tanzmeister, J. Thomas, D. Wollherr, and M. Buss. Grid-based mapping and tracking in dynamic environments using a uniform evidential environment representation. In *ICRA*, 2014.

[31] G. D. Tipaldi and F. Ramos. Motion clustering and estimation with conditional random fields. In *IROS*, 2009.

[32] A. K. Ushani, N. Carlevaris-Bianco, A. G. Cunningham, E. Galceran, and R. M. Eustice. Continuous-time estimation for dynamic obstacle tracking. In *IROS*, 2015.

[33] J. Van De Ven, F. Ramos, and G. D. Tipaldi. An integrated probabilistic model for scan-matching, moving object detection and motion estimation. In *ICRA*, 2010.

[34] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, et al. Scipy 1.0: fundamental algorithms for scientific computing in python. *Nature methods*, 17(3):261–272, 2020.

[35] T. Vu and O. Aycard. Laser-based detection and tracking moving objects using data-driven markov chain monte carlo. In *ICRA*, 2009.

[36] D. Z. Wang, I. Posner, and P. Newman. Model-free detection and tracking of dynamic objects with 2d lidar. *IJRR*, 2015.

[37] B. Yang, W. Luo, and R. Urtasun. PIXOR: Real-time 3D object detection from point clouds. In *CVPR*, 2018.

[38] W. Yuan, T. Khot, D. Held, C. Mertz, and M. Hebert. PCN: Point completion network. In *3DV*, 2018.

[39] J. Zarzar, S. Giancola, and B. Ghanem. PointRGCN: Graph convolution networks for 3D vehicles detection refinement. *arXiv:1911.12236*, 2019.

[40] D. Zermas, I. Izzat, and N. Papanikolopoulos. Fast segmentation of 3D point clouds: A paradigm on LiDar data for autonomous vehicle applications. In *ICRA*, 2017.

[41] H. Zou, J. Cui, X. Kong, C. Zhang, Y. Liu, F. Wen, and W. Li. F-siamese tracker: A frustum-based double siamese network for 3d single object tracking. *arXiv:2010.11510*, 2020.