# Semantic Geometric Fusion Multi-object Tracking and Lidar Odometry in Dynamic Environment

Tingchen Ma and Yongsheng Ou*

*Abstract*—The SLAM system based on static scene assumption will introduce huge estimation errors when moving objects appear in the field of view. This paper proposes a novel multi-object dynamic lidar odometry (MLO) based on semantic object detection technology to solve this problem. The MLO system can provide reliable localization of robot and semantic objects and build long-term static maps in complex dynamic scenes. For ego-motion estimation, we use the environment features that take semantic and geometric consistency constraints into account in the extraction process. The filtering features are robust to semantic movable and unknown dynamic objects. At the same time, a least square estimator using the semantic bounding box and object point cloud is proposed to achieve accurate and stable multi-object tracking between frames. In the mapping module, we further realize dynamic semantic object detection based on the absolute trajectory tracking list (ATTL). Then, static semantic objects and environmental features can be used to eliminate accumulated localization errors and build pure static maps. Experiments on public KITTI data sets show that the proposed system can achieve more accurate and robust tracking of the object and better real-time localization accuracy in complex scenes compared with existing technologies.

*Index Terms*—Lidar SLAM, semantic mapping, multi-object tracking, dynamic scene

## I. INTRODUCTION

Simultaneous localization and mapping (SLAM) [31], [41] technology using 3D lidar is widely used in industry due to its good scene robustness (lighting, low texture). Meanwhile, the calculation result provided by the SLAM system is also essential for robot control and task planning. Currently, most SLAM systems [2], [26] are built based on rigid assumptions, which will introduce huge estimation errors when there are many moving objects in the scene. Therefore, it is necessary to take reasonable measures to improve the performance of robot navigation systems in real scenes.

In recent years, some studies have considered detecting dynamic objects through spatial geometric constraints [8], [25], [27], then the SLAM system can avoid the interference of abnormal data association on ego localization. On the other hand, the learning-based semantic detection method is also an excellent way to process dynamic objects. For example, [30] and [28] directly identify moving 3d points in a single scan through the trained network model. [4] and [5] use semantic segmentation network [23] to obtain point-level semantic labels in the lidar scan, and realize 3D map construction with scene semantics (roads, buildings).

In fact, both geometry-based and semantic-based detection and elimination schemes can achieve reliable robot localization in the dynamic scene. However, for some practical applications such as planning tasks in autonomous driving and human-computer interaction in AR/VR, robots further need to stably track and represent semantic objects (cars and people) in the map coordinate. We propose a least square estimator for the above problems that uses semantic and geometric information for reliable multi-object tracking (SGF-MOT). Combined with the estimation results of ego odometry, the absolute localization for each object can be calculated. In the mapping module, the absolute trajectory tracking list of objects is used for dynamic semantic object detection. Then, static object and environment features can be used to eliminate the accumulated error of ego and object odometry. Finally, we realize the construction of a 4D scene map.

The main contributions of this paper are as follows:

- A complete multi-object lidar odometry system (MLO) can provide the absolute localization of both robot and semantic objects and build a 4D scene map (robot, semantic object localization and long-term static map).

- A least-square estimator considering the bounding box plane feature and geometric point cloud distribution is used for object state updating in the multi-object tracking module. The semantic bounding boxes can ensure the correct optimization direction. The relative motion model eliminates the point cloud distortion caused by robot and object motion. Meanwhile, the direct measurement point cloud can improve the estimation accuracy.

- A dynamic object detection method based on the absolute trajectory tracking list, which can identify slow-moving semantic objects.

The rest of this paper is structured as follows. First, in section 2, we discuss related works. Then, the proposed MLO system are described in section 3. The experimental setup,
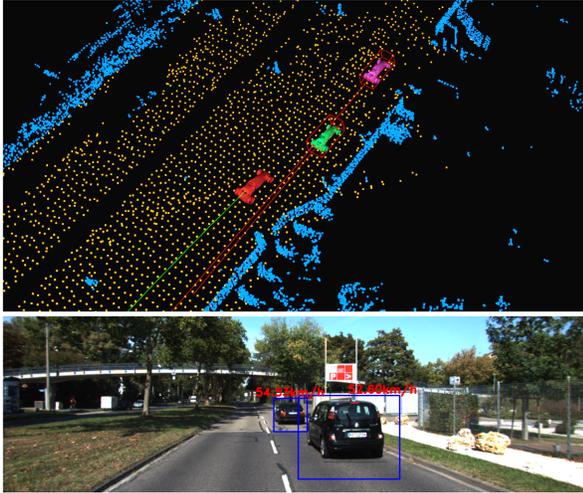
**Fig. 1. Qualitative results for sequence 0926-0013 in the KITTI-raw dataset.** The top half of the picture shows the results in rviz. The red car is a robot equipped with lidar, and green line is the ego motion path. The red lines are the tracked object motion paths. The bottom half (for visualization only) shows the projected 2D bounding box (blue) and corresponding object velocity (red) in the image with the synchronized timestamp.

results, and discussion will be presented in section 4. Finally, section 5 is the conclusion.

## II. RELATED WORK

### A. Object detection and tracking

For the multi-object tracking (MOT) method based on point cloud data, object detection and tracking are the two main steps of most solutions. In the early days, traditional methods focused on extracting objects from the spatial distribution information of point clouds. Then, object tracking is realized by constructing geometric error constraints and corresponding estimator. Moosmann et al. [25] implemented a general object segmentation method based on the local convexity criterion. The point-to-plane ICP [6] algorithm is used for MOT. Dewan et al. [9] used RANSAC [11] to estimate the motion model sequentially. The point clouds were then associated with the corresponding model by a bayesian approach. By assuming local geometric constancy and regularization of smooth motion field, Ferri et al. [10] estimated object motion through rigid scene flow, and their method can be suitable for non-rigid motion, such as pedestrians. Sualeh et al. [20] fitted the 3d bounding box from the cluster point cloud and updated the object state based on the Kalman filter [17].

Recently, the learning-based recognition method on point clouds has brought new solutions to the MOT problem. They can bring more stable object recognition and semantic information (such as object bounding boxes). Weng et al. [37] use Point-RCNN [32] to detect semantic objects in point clouds and track 3D bounding boxes (3D-BB) through the Kalman filter. Kim et al. [18] used 2D/3D-BB detected in image and point clouds for filter tracking, achieving better tracking accuracy. Wang et al. [36] and Huang et al. [16] designed an end-to-end MOT framework using different deep-learning models to complete object detection and data association tasks.

In summary, a reliable initial guess is indispensable for the optimization-based tracking method using point cloud geometric constraints. When the initial value of the estimator is poor, many point-level data associations may provide the wrong direction at the early stage of the calculation. On the other hand, extracting bounding boxes from the point cloud or learning method for filter-based tracking is relatively stable. However, it abandons the high-precision measured point cloud. In our SGF-MOT module, semantic bounding box plane and geometric point cloud distribution are used simultaneously in a least-square estimator, which achieves a better balance between tracking accuracy and robustness.

### B. Dynamic aware 3d lidar SLAM

At present, 3d lidar SLAM systems are mainly divided into two categories: feature-based and matching-based. For the feature-based method, Zhang et al. [41] and Shan et al. [31] extracted stable line and surface features from the point cloud by scan line smoothness. Then the extracted features are used for localization and mapping. Liu et al. [21] introduced the bundle adjustment (BA) [34] method, which is more commonly used in visual SLAM, to improve the system mapping accuracy. In the matching-based method, Behley et al. [2] used surfel to model the environment and designed a projection data association method for ego localization. Dellenbach et al. [7] used the implicit moving least squares surface to represent the map and realized high-precision mapping by considering the continuous time constraints between frames. The above systems are built on the assumption of the rigid scene and cannot work stably in a highly dynamic scene.

For dynamic-aware 3D lidar SLAM technology, some studies use geometric methods such as background model [27], local convexity [25], and point cloud clustering [8] to detect and eliminate information that may violate static assumptions. On the other hand, with the development of deep learning technology, some semantic-based lidar SLAM can also process dynamic objects in the field of view. Pfreundschuh et al. [28] designed a method for the automatic labelling of dynamic objects and detected the dynamic 3d points in the point cloud before the implementation of LOAM [41] on the 3D-MinNet network. Sualeh et al. [20] improved the accuracy of dynamic object detection by considering the space-time constraints of sequence frames input during model training. Both [4] and [5] using semantic segmentation networks realize semantic scene mapping and detect dynamic objects through point-by-point checks and feature extraction from the object point cloud.

Unlike the above scheme, our system (MLO) integrates the results of robust ego localization and SGF-MOT module, thus achieving accurate and reliable localization of instance objects in the map coordinate. At the same time, the object absolute trajectory tracking list is used to detect dynamic objects in the scene. So we create a 4D scene map that contains both robot, semantic objects position and static environment.

## III. METHOD

The whole MLO system flowchart is shown in Fig.2, which includes three main parts: Multi-task Fusion Perception,
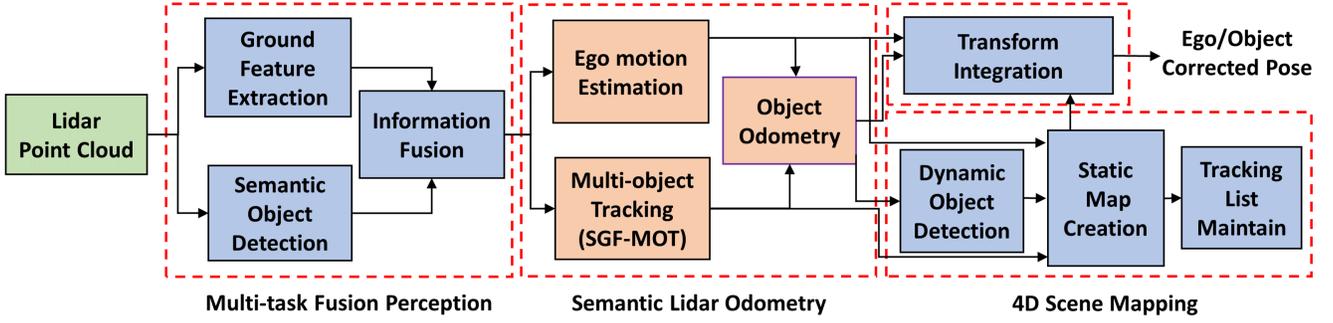
Fig. 2. **The proposed MLO system framework.** The orange box is the key part of this paper, it includes the semantic-geometric fusion multi-object tracking method (SGF-MOT) which will be discussed in detail. A 4D scene mapping module is also introduced for maintaining the long-term static map and movable object map separately, which can be used to further navigation applications.

Semantic Lidar Odometry, and 4D Scene Mapping. Firstly, the ground and objects are detected and refined by information fusion. Then, the proposed SGF-MOT tracker is integrated into the odometry module to achieve robust and accurate tracking of objects. After dynamic object detection, the mapping module uses both static object and environment features [41] to eliminate accumulated errors.

### A. Notation and coordinate frames

In this paper, we define the notation as follows: The lidar (mechanical) has pose $\mathbf{T}_{ol}^i \in \mathrm{SE}(3)$ at time $i$ in odom coordinate $o$. It can observe movable objects with pose $\mathbf{T}_{ot}^i \in \mathrm{SE}(3)$ in odom coordinate $o$. Meanwhile, the lidar pose in map coordinate $m$ is $\mathbf{T}_{ml}^i \in \mathrm{SE}(3)$, and the correction matrix which transforms the object and lidar from odom coordinate $o$ to map coordinate $m$ is $\mathbf{T}_{mo}^i \in \mathrm{SE}(3)$.

### B. Multi-task fusion perception

After getting the ROI point clouds $\mathbf{\Phi}^i$ of the current frame, we perform ground feature extraction and semantic object detection in parallel to obtain preliminary object point clouds and geometric ground features. Then, the information fusion module can receive more accurate object points segmentation and environment features through a mutual correction step.

*1) Ground feature extraction:* First, we fit ground parameters in the point cloud using the iterative PCA algorithm proposed in [40] and mark the corresponding points near the ground in point cloud $\mathbf{\Phi}_g^i$.

The distortion correction time (taking the value range [0,1]) is calculated by the relative offset of each 3D point angle from the 3D point angle at the initial time of the frame. The scan line is uniformly segmented. Through the pointwise inspection of ground labels (**0:** background point, **1:** ground point) calculated above, it is possible to record whether each segment contains the ground point or background point labels. The segment head and tail point id pairs containing the ground or background points are stored in the ground container $\mathbf{V}_g$ or background container $\mathbf{V}_b$.

Based on the smoothness (readers can refer to [41] for the calculation formula) calculated by segment, ground features $\chi_f^i$ can be extracted from the point cloud contained in $\mathbf{V}_g$. The

smoothness prior $S_a$ used for subsequent background surface feature check is obtained as follows:

$$S_a = \sum_{o=1}^{N} c_o \bigg/ N_c \tag{1}$$

where $c_o$ is the smoothness of $o$-th ground feature, and $N_c$ is the total number of ground features in current frame. After downsampling all ground points, the candidate ground point cloud $\chi_c^i$ is generated.

For point clouds contained in $\mathbf{V}_b$, only smoothness is calculated and sorted here. After that, the information fusion module will finely select the background features.

*2) Semantic object detection:* The semantic object detection module is to acquire 3D-BB and the foreground object point clouds they enclose. In addition, the 3D-BB will be used to provide optimization constraints for the motion estimation of movable objects. We accomplish the above tasks using 3DSSD [39], a single-stage, point-based lightweight object detection model that balances detection accuracy and efficiency. Assume that the $k$-th 3D-BB in current detection result is:

$$\tau^{k,i} = \left[ x^{k,i}, y^{k,i}, z^{k,i}, l^{k,i}, w^{k,i}, h^{k,i}, yaw^{k,i} \right]^T \tag{2}$$

where $(x, y, z)$ represents the object position, $(l, w, h)$ represents the 3D-BB size, and $yaw$ is the object direction angle. Each group of object point cloud in 3D-BB is stored as $\gamma^{k,i}$. Each 3D point is marked with an independent object id in $\mathbf{\Phi}_s^i$. Object semantic confidences are denoted as $\mu^{k,i}$.

*3) Information fusion:* Considering that 3D-BB obtained by semantic detection module usually encloses both object and ground points. Meanwhile, features obtained by ground feature extraction module also contain a small number of points belonging to the object. Therefore, it is necessary to perform mutual correction first to ensure their motion consistency.

According to the object id marked in $\mathbf{\Phi}_s^i$, object points in ground feature sets $\chi_f^i$ and $\chi_c^i$ are eliminated. Similarly, the ground feature in each object point cloud $\gamma^{k,i}$ can be rejected by labels in $\mathbf{\Phi}_g^i$. The ambiguous intersection information will not participate in the subsequent SLAM and MOT tasks.

For point clouds contained in $\mathbf{V}_b$, we perform two feature extractions based on the computed smoothness: The candidate object edge feature $\gamma_c^{k,i}$ is extracted from the point clouds belonging to objects, which will be used in mapping module to
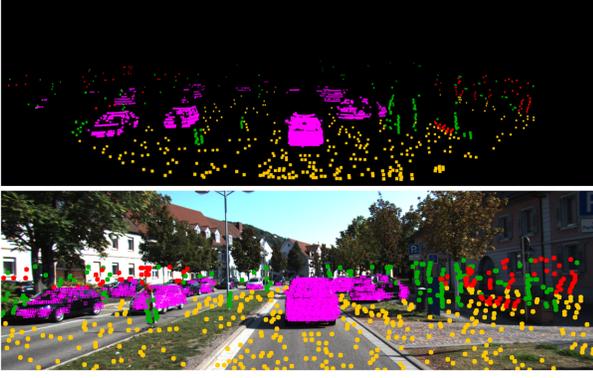
Fig. 3. **The results of multi-task fusion perception module for sequence 0926-0056 in the KITTI-raw dataset.** Pink point clouds represent moving cars, and yellow point clouds represent ground features. Green and red point clouds are background edge and surface features. The image projection results with the same timestamp are used for visualization only.

eliminate accumulated errors. Then, background edge features $\psi_f^i$, surface features $\zeta_f^i$, and the corresponding candidate sets $\psi_c^i$, $\zeta_c^i$ are extracted from the point clouds that do not belong to semantic objects and ground features, and they will participate in both localization and mapping modules.

Finally, we define the edge and surface (including ground) features used for state estimation in this paper. The distance from the point to the edge or surface features in lidar coordinate $l$ are as follows:

$$\mathbf{e}_{e_k} = \frac{\left| (^l\mathbf{p}_{i,k}^e - {}^l\bar{\mathbf{p}}_{i-1,m}^e) \times (^l\mathbf{p}_{i,k}^e - {}^l\bar{\mathbf{p}}_{i-1,n}^e) \right|}{\left| {}^l\bar{\mathbf{p}}_{i-1,m}^e - {}^l\bar{\mathbf{p}}_{i-1,n}^e \right|} \quad (3)$$

$$\mathbf{e}_{s_k} = \left| \frac{\left( {}^l\mathbf{p}_{i,k}^s - {}^l\bar{\mathbf{p}}_{i-1,m}^s \right)}{\left( {}^l\bar{\mathbf{p}}_{i-1,m}^s - {}^l\bar{\mathbf{p}}_{i-1,n}^s \right) \times \left( {}^l\bar{\mathbf{p}}_{i-1,m}^s - {}^l\bar{\mathbf{p}}_{i-1,o}^s \right)}{\left| \left( {}^l\bar{\mathbf{p}}_{i-1,m}^s - {}^l\bar{\mathbf{p}}_{i-1,n}^s \right) \times \left( {}^l\bar{\mathbf{p}}_{i-1,m}^s - {}^l\bar{\mathbf{p}}_{i-1,o}^s \right) \right|} \right| \quad (4)$$

where $k$ represents the edge or surface feature indices in current frame. $m, n, o$ represent the matching edge line or planar patch indices in last frame. For edge features $^l\mathbf{p}_{i,k}^e$ in set $\psi_f^i$, $^l\bar{\mathbf{p}}_{i-1,m}^e$ and $^l\bar{\mathbf{p}}_{i-1,n}^e$ are the matching points in set $\bar{\psi}_c^{i-1}$ that make up the edge line in last frame without motion-distortion. Similarly, $^l\mathbf{p}_{i,k}^s$ is the surface feature in set $\zeta_f^i$. $^l\bar{\mathbf{p}}_{i-1,m}^s$, $^l\bar{\mathbf{p}}_{i-1,n}^s$, and $^l\bar{\mathbf{p}}_{i-1,o}^s$ are the points that consist of the corresponding planar patch in set $\bar{\chi}_c^{i-1}$ or $\bar{\zeta}_c^{i-1}$.

### C. Semantic lidar odometry

Based on the perception results in Section 3.2, the odometry module firstly estimates the absolute increment for ego motion and relative motion increment of each semantic object in parallel. Then, each object pose is calculated in body-fixed coordinate [15].

*1) Robust ego odometry:* Before entering the ego-motion estimation process, **Algorithm 1** is used to achieve a geometric consistency check on the ground and background features, eliminating unknown dynamic influences that the object detection module cannot identify. The feature sets that passed the geometric consistency check are denoted as $(\tilde{\chi}_f^i, \tilde{\psi}_f^i, \tilde{\zeta}_f^i)$.

---

**Algorithm 1:** Geometric consistency feature check

**Input:** Current feature sets: $(\chi_f^i, \psi_f^i, \zeta_f^i)$, candidate sets without distortion: $(\bar{\chi}_c^{i-1}, \bar{\psi}_c^{i-1}, \bar{\zeta}_c^{i-1})$, increment $_l^{i-2}\delta^{i-1}$ from last loop, association update interval $\lambda$, thresholds: $(th_r, th_g, th_b)$

**Output:** Consistent feature sets: $(\tilde{\chi}_f^i, \tilde{\psi}_f^i, \tilde{\zeta}_f^i)$

The number of feature sets $(\chi_f^i, \psi_f^i, \zeta_f^i)$ is recorded as $N_f$. Let $_l^{i-2}\delta_{best}^{i-1} = {}_l^{i-2}\delta^{i-1}$;

**for** $i = 1 : N_{iter}$ **do**

  **if** $i \% \lambda = 0$ **then**

    Based on $_l^{i-2}\delta_{best}^{i-1}$ and $(\chi_f^i, \psi_f^i, \zeta_f^i)$ including distortion correction time, feature position at the initial time of current frame is computed; Search their closest points in kd-trees formed by $(\bar{\chi}_c^{i-1}, \bar{\psi}_c^{i-1}, \bar{\zeta}_c^{i-1})$. Then correct $(\chi_f^i, \psi_f^i, \zeta_f^i)$ to the end time of current frame, denoted as $(\mathbf{m}_{\chi_f}, \mathbf{m}_{\psi_f}, \mathbf{m}_{\zeta_f})$;

  **end**

  Random select matching pairs from $\mathbf{m}_{\chi_f}$, $\mathbf{m}_{\psi_f}$ and their candidate sets;

  Update $_l^{i-2}\delta^{i-1}$ with matching pairs and ICP model;

  Transform $(\chi_f^i, \psi_f^i, \zeta_f^i)$ to the initial time of current frame with $_l^{i-2}\delta^{i-1}$. Find closest points in kd-trees to perform Euclidean distance check;

  Record the number of points passing the check as $(\eta_{\chi_f}, \eta_{\psi_f}, \eta_{\zeta_f})$. Let $\eta_n = \eta_{\chi_f} + \eta_{\psi_f} + \eta_{\zeta_f}$;

  **if** $(\eta_n/N_f \geq th_r) \wedge (\eta_{\chi_f} \geq th_g) \wedge (\eta_{\psi_f} + \eta_{\zeta_f} \geq th_b) \wedge (\eta_n > \eta_{best})$ **then**

    Record feature sets $(\tilde{\chi}_f^i, \tilde{\psi}_f^i, \tilde{\zeta}_f^i)$;

    Let $\eta_{best} = \eta_n$, $_l^{i-2}\delta_{best}^{i-1} = {}_l^{i-2}\delta^{i-1}$;

  **end**

**end**

---

To improve the computational efficiency of ego odometry, we use the two-step algorithm proposed by [31] for pose estimation. First, frame increment $[t_z, \theta_{roll}, \theta_{pitch}]$ is computed based on matched points in the checked ground feature set $\tilde{\chi}_f^i$, and the distortion-corrected candidate feature set $\bar{\chi}_c^{i-1}$. Then, increment $[t_z, \theta_{roll}, \theta_{pitch}]$ is fixed, and another 3-DOF increment $[t_x, t_y, \theta_{yaw}]$ is estimated based on the edge feature $\tilde{\psi}_f^i$, surface feature in $\tilde{\zeta}_f^i$ with more "flat" smoothness than $S_a$ and their corresponding candidate sets. Finally, we obtain ego pose $\mathbf{T}_{ol}^i$ by accumulating the 6-DOF increment $_l^{i-1}\delta^i = [t_x, t_y, t_z, \theta_{roll}, \theta_{pitch}, \theta_{yaw}]$.

*2) Semantic-geometric fusion multi-object tracking:* In order to track semantic objects robustly and accurately, a least-squares multi-object tracking module fusing semantic 3D-BB and geometric point clouds is designed in this paper. The module maintains a tracking list containing object information (object point cloud, sample points on 3D-BB, and poses in lidar coordinate) and tracking states (unique id, optimization quality evaluation and tracking continuity results) for object data association and relative motion increment estimation.

First, we use the constant motion model to predict object

position, which, together with the object detection position, forms the position error term $\mathbf{e}_p$. Meanwhile, the direction error term $\mathbf{e}_d$ is formed by the relative object motion direction, and the 3D-BB size is used to calculate the scale error term $\mathbf{e}_s$. Then the association matrix is solved by Kuhn-Munkres algorithm to obtain the object matching relationship. The calculation method of each error term can refer to [1].

For objects that are tracked in last frame but no matching relationship is currently found, we use the constant motion model to infer their positions in the subsequent two frames. If still no successful matching is found, they will be removed.

Then, we use the voxelized G-ICP algorithm proposed in [19] to construct the cost function of object point cloud $\gamma^i$ (object id $k$ is omitted). Meanwhile, considering the point cloud distortion caused by the relative motion of object and robot between frames, we introduce a linear interpolation relative motion model to deal with this problem.

Assuming that in lidar coordinate $l$, the $j$-th point on current processing object is $^l\mathbf{p}_{i,j}^g$, and its corresponding timestamp is $t_j$. Let $_l^{i-1}\mathbf{H}_j^i \in$ SE(3) denote pose transformation between interval $[t^{i-1}, t_j]$. Accordingly, $_l^{i-1}\mathbf{H}_j^i$ can be calculated by linear interpolation of the relative motion increment $_l^{i-1}\mathbf{H}^i \in$ SE(3):

$$_l^{i-1}\mathbf{H}_j^i = \frac{t_j - t^i}{t^{i-1} - t^i} {}_l^{i-1}\mathbf{H}^i \tag{5}$$

We use $_l^{i-1}\mathbf{R}^i$ and $_l^{i-1}\mathbf{t}^i$ to represent the rotation and translation parts of the increment $_l^{i-1}\mathbf{H}^i$. The geometric error term $\mathbf{e}_{g_j}$ containing motion distortion correction is expressed as:

$$\mathbf{e}_{g_j} = \frac{\sum {}^l\mathbf{p}_{i-1,m}^g}{N_j} - {}_l^{i-1}\mathbf{H}_j^i {}^l\mathbf{p}_{i,j}^g \tag{6}$$

where $^l\mathbf{p}_{i-1,m}^g$ is the $m$-th point in matched voxel $^l\mathbf{v}_{i-1}^g$ at the end time of the last frame, $N_j$ is the number of 3D points contained in the matched voxel. The covariance matrix $\mathbf{\Omega}_j$ corresponding to the error term $\mathbf{e}_{g_j}$ is as follows:

$$\mathbf{\Omega}_j = \frac{\sum C_{i-1,m}^g}{N_j} - {}_l^{i-1}\mathbf{H}_j^i C_{i,j}^g {}_l^{i-1}\mathbf{H}_j^{i\,T} \tag{7}$$

$C_{i-1,m}^g$ and $C_{i,j}^g$ denote the covariance matrix. They describe the surrounding shape distribution of the point in the matched voxel and point $^l\mathbf{p}_{i,j}^g$.

Note that objects detected in a single frame always have less constrained information. Meanwhile, the inaccurate initial value of the estimator may also lead to insufficient matchings at the beginning of the optimization. Eq. (6) sometimes fails to constrain the 6-DOF motion estimation. Therefore, while using geometric point cloud constraints, this paper further introduces semantic constraints to increment estimation for providing the right converged direction.

Specifically, we model a 3D-BB $\tau^i$ as a cuboid containing six planes and use $\pi = [\mathbf{n}, d]$ to represent a single plane. Where, $\mathbf{n}$ is the plane normal vector with $\|\mathbf{n}\|_2 = 1$. $d$ is the distance from the origin to the plane. Unlike the geometric point cloud error with iterative matching for estimation, the matching relationship for 3D-BB plane features is directly

available. Then, using the closest point $\mathbf{\Pi}$ [13] to parameterize $\pi$, i. e. $\mathbf{\Pi} = \mathbf{n}d$, we can get the plane-to-plane semantic error term $\mathbf{e}_{b_r}$ as follows:

$$\mathbf{e}_{b_r} = {}^l\mathbf{\Pi}_{i-1,r}^b - {}^l\hat{\mathbf{\Pi}}_{i-1,r}^b \tag{8}$$

$^l\mathbf{\Pi}_{i-1,r}^b$ and $^l\hat{\mathbf{\Pi}}_{i-1,r}^b$ are the r-th measured and estimated planes on 3D-BB at time $i$ in lidar coordinate $l$. As described in [14], the motion transformation of a plane can be expressed as:

$$\begin{bmatrix} {}^l\mathbf{n}_{i-1,r}^b \\ {}^l d_{i-1,r}^b \end{bmatrix} = \begin{bmatrix} {}_l^{i-1}\mathbf{R}^i & \mathbf{0} \\ -({}_l^{i-1}\mathbf{R}^{i\,T} {}_l^{i-1}\mathbf{t}^i)^T & 1 \end{bmatrix} \begin{bmatrix} {}^l\mathbf{n}_{i,r}^b \\ {}^l d_{i,r}^b \end{bmatrix} \tag{9}$$

Taking Eq.(9) into $^l\hat{\mathbf{\Pi}}_{i-1,r}^b$, we can further obtain:

$$^l\hat{\mathbf{\Pi}}_{i-1,r}^b = \left( {}_l^{i-1}\mathbf{R}_j^i {}^l\mathbf{n}_{i,r}^b \right) \left( -{}_l^{i-1}\mathbf{t}_j^{i\,T} {}_l^{i-1}\mathbf{R}_j^i {}^l\mathbf{n}_{i,r}^b + d \right) \tag{10}$$

When using the least squares method to estimate $_l^{i-1}\mathbf{H}^i$, it is necessary to obtain the derivative of $\mathbf{e}_{b_r}$ with respect to $_l^{i-1}\mathbf{H}^i$:

$$\frac{\partial \mathbf{e}_{b_r}}{\partial \delta_{{}_l^{i-1}\theta_j^i}} = -\left[ \mathbf{F}{}_l^{i-1}\mathbf{t}_j^{i\,T}\mathbf{F} \right]_\times - \mathbf{F}_{\mathbf{l}}^{\mathbf{i-1}}\mathbf{t}_{\mathbf{j}}^{\mathbf{i\,T}}[\mathbf{F}]_\times + [\mathbf{F}d]_\times \tag{11}$$

$$\frac{\partial \mathbf{e}_{b_r}}{\partial \delta_{{}_l^{i-1}\mathbf{t}_j^i}} = \mathbf{F}\mathbf{F}^T \tag{12}$$

where $\mathbf{F} = {}_l^{i-1}\mathbf{R}_j^i {}^l\mathbf{n}_{i,r}^b$. $\frac{\partial \mathbf{e}_{b_r}}{\partial \delta_{{}_l^{i-1}\theta_j^i}}$ and $\frac{\partial \mathbf{e}_{b_r}}{\partial \delta_{{}_l^{i-1}\mathbf{t}_j^i}}$ represent the jacobians of the error term $\mathbf{e}_{b_r}$ w.r.t. the rotation and translation component. The detailed derivation of Eq. (11) and (12) can be found in APPENDIX.

Finally, Eq. (13) is the comprehensive cost function for the estimation of $_l^{i-1}\mathbf{H}^i$ considering the semantic confidence $\mu^i$:

$$\underset{_l^{i-1}\mathbf{H}^i}{\arg\min} \frac{2 - \mu^i}{2N_g} \sum_{j=1}^{N_g} \left( N_j\, \mathbf{e}_{g_j}^T \mathbf{\Omega}_j^{-1} \mathbf{e}_{g_j} \right) + \frac{\mu^i}{2N_b} \sum_{r=1}^{N_b} \left( \mathbf{e}_{b_r}^T \mathbf{e}_{b_r} \right) \tag{13}$$

where $N_g$ is the number of object point cloud. $N_b$ is the number of plane features on 3D-BB.

In this paper, we use LM [22] algorithm to solve Eq. (9) and the estimated relative increment to get motion-distorted object point cloud $\bar{\gamma}^i$ and object edge feature $\bar{\gamma}_c^i$, which will be used for object tracking in next frame and accumulated error elimination in mapping module.

*3) Object odometry:* Based on the calculation results of ego odometry and multi-object tracking, object pose in odom coordinate can be computed. First, object motion increment $_l^i\mathbf{X}^{i-1}$ in lidar coordinate is denoted as:

$$_l^i\mathbf{X}^{i-1} = {}_l^{i-1}\delta^i {}_l^i\mathbf{H}^{i-1} \tag{14}$$

Then, combined with the body-fixed coordinate proposed in [15] (as shown in Fig.4), and assuming that the transform matrix from object to lidar at time $i - 1$ is $^{i-1}\mathbf{L}_{lt}$, we can obtain absolute object motion increment $_t^{i-1}\mathbf{X}^i$ in body-fixed coordinate:
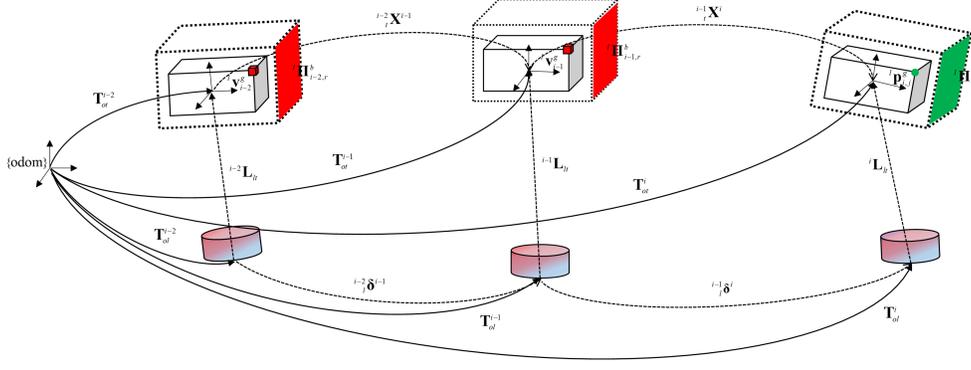
Fig. 4. **Notation and coordinate frames.** Coloured cylinders denote the lidar. White cuboids represent a moving object. Dashed cuboids are the result inferred by semantic detection module. Solid lines represent lidar and object pose in odom coordinate. Dashed lines indicate their motion increment in local coordinate. Small red cubes and parallelograms represent voxels whose motion distortion has been eliminated and 3D planes obtained by 3D-BB. Green dot and plane represent object points in current point cloud and 3D-BB plane feature without motion distortion.

$$_{t}^{i-1}\mathbf{X}^i = {}^{i-1}\mathbf{L}_{lt}^{-1}\ {}_{l}^{i}\mathbf{X}^{i-1}\ {}^{i-1}\mathbf{L}_{lt} \tag{15}$$

Odometry pose $\mathbf{T}_{ot}^i$ for each object at time $i$ can be computed by accumulating increment $_{t}^{i-1}\mathbf{X}^i$. This paper uses ego odometry $\mathbf{T}_{ol}$ and transform matrix $\mathbf{L}_{lt}$ extracted from semantic detection module to set object odometry $\mathbf{T}_{ot}$ at the first detected frame.

### D. 4D scene mapping

To reuse the created maps, we suggest maintaining the long-term static map and movable object maps separately. First, we use the object absolute trajectory tracking list (ATTL) to detect dynamic objects. Static map creation aims to update the correction matrix $\mathbf{T}_{mo}^i$ for accumulated error correction with static objects and environment features. Finally, we correct the object pose in map coordinate and update the object map based on the object motion state.

*1) Dynamic object detection:* As shown in Fig.5, the ATTL maintained by MLO system in the mapping module is used to analyze the motion state of each instance semantic object. First, the current frame objects are matched with the ATTL through the unique id obtained from SGF-MOT module. Then, we use the current error correction matrix $\mathbf{T}_{mo}^{i-1}$ to predict the object pose $\mathbf{T}_{mt}^{i'}$ in map coordinate:

$$\mathbf{T}_{mt}^{i'} = \mathbf{T}_{mo}^{i-1}\mathbf{T}_{ot}^i \tag{16}$$

By calculating the motion increment of the predicted object pose and the oldest static pose of the matching object in ATTL, we can determine whether each object is in motion. Note that the detected object in the first frame is set to be static by default.

*2) Static map creation:* By matching the static object features $\bar{\gamma}_c^i$ in the current frame with the temporary object map, ground feature $\bar{\chi}_c^i$, background edge feature $\bar{\psi}_c^i$ and surface feature $\bar{\zeta}_c^i$ with their surrounding point cloud map, ego pose $\mathbf{T}_{ol}^i$ can be corrected to map coordinate. Then, a long-term static map can be built by aligning the non-semantic features into map coordinate. Readers can refer to [41] for more map matching and optimization details.
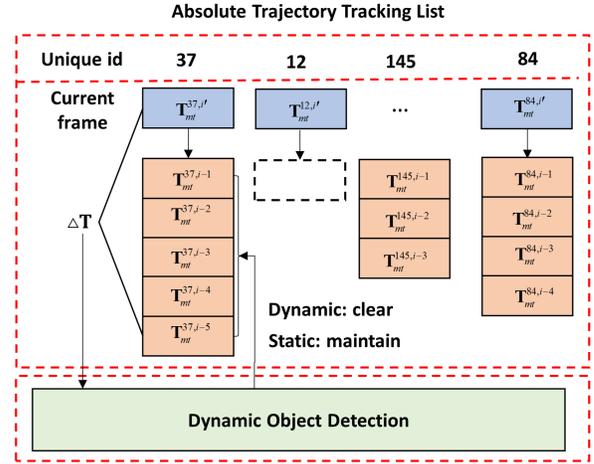


Fig. 5. **Dynamic object detection based on ATTL.** The blue boxes represent the object predicted pose in the current frame, and the orange boxes represent the object tracking list that has corrected the accumulated error. Objects with ids 37 and 84 are in normal tracking state. The object with id 12 will be initialized as a static object, and the object with id 145 will be deleted because the current frame does not observe it. $\Delta\mathbf{T}$ is used to determine the current motion state of each object.

According to the ego pose $\mathbf{T}_{ol}^i$ in odom coordinate and the corrected pose $\mathbf{T}_{ml}^i$ in map coordinate, the updated accumulated error correction matrix $\mathbf{T}_{mo}^i$ can be obtained.

*3) Tracking list maintain:* The ATTL can be updated based on the motion state of each object. If the object is moving, the map should be cleared first. Otherwise, the map will remain unchanged. Then, we use $\mathbf{T}_{mo}^i$ and $\mathbf{T}_{ml}^i$ to convert the object odometry pose $\mathbf{T}_{ot}^i$ and extracted object features $\bar{\gamma}_c^{k,i}$ into map coordinate, which will participate in subsequent mapping steps.

## IV. EXPERIMENT

The experimental environment of the proposed framework is AMD® Ryzen 7 5800h (8 cores @3.2 GHz), 16GB RAM, ROS Melodic, and NVIDIA GeForce RTX 3070. We use absolute trajectory error (ATE) [33], relative rotation error (RRE) and relative translation error (RTE) [12] to evaluate the accuracy and odometry drift of ego and object localization.
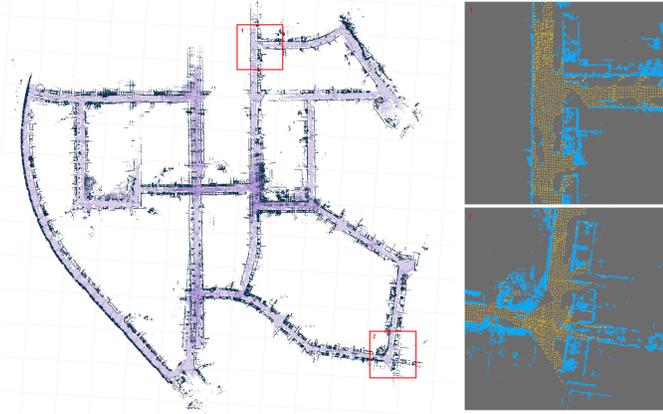
Fig. 6. **Point cloud map created by sequence 00 in the KITTI-odometry dataset.** Blue point cloud on the left represents the background map, and purple point cloud represents the ground. As shown in the enlarged section, the resulting map does not contain semantic objects that may change over time by maintaining semantic objects and static map separately.

| | Absolute Motion Trajectory RMSE [m] | | | | |
|------|---------|-----------|--------|---------|-------|
| Seq  | A-LOAM  | Lego-LOAM | F-LOAM | MLO-ddo | MLO   |
| 00   | 12.12   | 41.99     | 14.71  | 8.23    | **7.52** |
| 01   | 15.81   | **fail**  | 17.61  | **12.11** | 12.42 |
| 02   | **fail** | **fail** | **fail** | 18.38 | **17.65** |
| 03   | 0.57    | 14.19     | 0.90   | 0.52    | **0.50** |
| 04   | 0.45    | **fail**  | 0.36   | 0.29    | **0.27** |
| 05   | 5.72    | 8.12      | 8.06   | 2.74    | **1.92** |
| 06   | **0.44** | **fail** | 1.04   | 0.45    | 0.46  |
| 07   | 2.73    | 3.75      | 2.49   | 2.33    | **1.48** |
| 08   | **4.42** | **fail** | 6.02   | 5.18    | 4.86  |
| 09   | **3.48** | **fail** | 21.53  | 4.66    | 4.72  |
| 10   | **1.18** | 10.57    | 1.54   | 1.80    | 1.90  |
| Mean | 4.69    | **fail**  | 7.43   | 3.88    | **3.66** |

The CLEAR MOT metric [3] is used to evaluate estimated accuracy and tracking object configurations consistently over time for multi-object tracking module.

Note that our multi-object tracking module is performed in lidar coordinate. Therefore, the relative motion tracking ability of the object is evaluated. On the other hand, object trajectories in map coordinate evaluate the overall accuracy of ego localization and multi-object tracking module.

Furthermore, the point cloud in the KITTI-raw benchmark does not use known GPS/IMU information to correct motion distortion nor provides ground-truth trajectories that can be directly used for localization evaluation. For obtaining a smooth and accurate 6-DOF pose ground truth, we use the extended Kalman filter [24] to process GPS/IMU data in the KITTI raw dataset. Combined with the high-precision artificial object annotation in lidar coordinate, the ground truth of object pose in map coordinate can be obtained.

### A. Ego odometry localization accuracy experiment

The KITTI-odometry benchmark contains rich static objects with long sequences. Meanwhile, the KITTI-raw dataset contains more challenging scenes, such as the city sequence with unknown dynamic semantic objects (trains). We evaluate the MLO system in these two datasets by comparing it with other state-of-the-art 3D lidar SLAM. Since the KITTI dataset labels semantic objects within the camera's field of view, only the ROI-cutting lidar data is input into the fusion perception module.

For the experimental lidar SLAM system, A-LOAM [29] is a simplified version of LOAM [41] with the removal of the IMU used. It also extracts edge and surface features from point cloud to perform frame-to-frame matching for odometry estimation and frame-to-map matching to complete mapping. Lego-LOAM [31] takes into account ground constraints. Use ground features and background edge features, respectively, to estimate the increment $[t_z, \theta_{roll}, \theta_{pitch}]$ and $[t_x, t_y, \theta_{yaw}]$ for improving efficiency. F-LOAM [35] proposes a two-step

motion distortion removal algorithm, and then localization and mapping are done directly through frame-to-map matching.

As shown in Table I, our method successfully runs all sequences under the KITTI-odometry benchmark and achieves the best results in most cases. When estimating $[t_x, t_y, \theta_{yaw}]$ based on background edge features, the Lego-LOAM [31] system shows severe drift on many sequences due to only using point clouds from the camera's view. Our system also uses the two-step estimation algorithm to ensure computational efficiency. However, using background surface features based on smoothness prior $S_a$ and geometric consistency check algorithm improve the MLO system robustness.

Further, detecting dynamic objects by ATTL and adding static object constraints for static map creation can better eliminate accumulated error compared to MLO-ddo system, which does not detect and use static semantic objects in the mapping module. In sequences 00, 02, 03, 05 and 07 containing rich static objects, the MLO system achieved better localization accuracy.
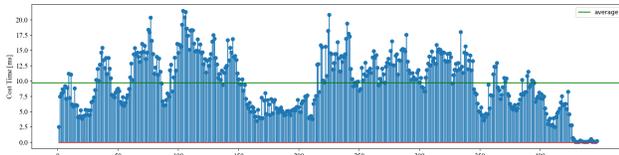
Table II shows that in complex scenes, the MLO system can still achieve better localization accuracy in most cases. Meanwhile, the MLO system outperforms the other two frameworks in terms of average RRE and RTE. For sequences 0926-0056, there is a train that the object detection module cannot recognize. Since the A-LOAM [29] system does not distinguish the ground and background features and makes use of all the information for localization and mapping, it can improve the system robustness to a certain extent. For the MLO system that also uses the two-step estimation algorithm, especially for the estimation of $[t_x, t_y, \theta_{yaw}]$, a geometric consistency check can reduce the proportion of abnormal associations and avoid the trajectory drift problem of the MLO-gc system (ego-motion estimation module uses the features extracted after removing the object point cloud, but does not consider the geometric consistency check algorithm).

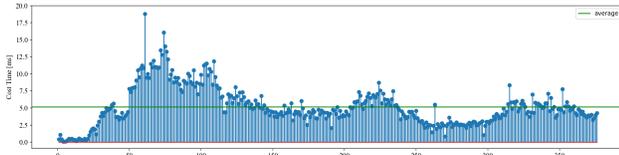### B. Object tracking accuracy and robustness experiment

First, we compare the SGF-MOT module proposed in this paper with AB3DMOT [37] and PC3T [38] trackers under the KITTI-tracking benchmark. The tracking results are evaluated in 2D image, which can be computed by the KITTI calibration matrix. AB3DMOT [37] only takes point cloud as input and

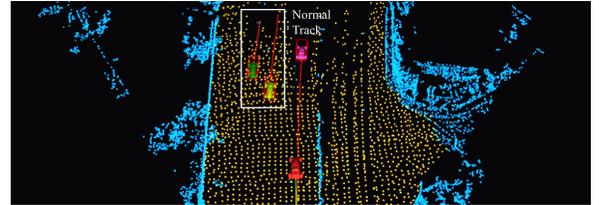| Sequence | A-LOAM | | | MLO-gc | | | MLO | | |
|---|---|---|---|---|---|---|---|---|---|
| | ATE[m] | RTE[m/f] | RRE[deg/f] | ATE[m] | RTE[m/f] | RRE[deg/f] | ATE[m] | RTE[m/f] | RRE[deg/f] |
| 0926-0011 | **0.36** | 0.09 | 0.11 | 0.39 | 0.06 | 0.08 | 0.39 | 0.06 | 0.07 |
| 0926-0014 | 0.42 | 0.11 | 0.21 | **0.26** | 0.11 | 0.14 | 0.27 | 0.11 | 0.13 |
| 0926-0056 | 0.37 | 0.16 | 0.21 | **fail** | 0.51 | 0.15 | **0.32** | 0.12 | 0.13 |
| 0926-0059 | **0.24** | 0.06 | 0.14 | 0.28 | 0.07 | 0.11 | 0.28 | 0.06 | 0.10 |
| 0929-0071 | **1.62** | 0.05 | 0.15 | 1.63 | 0.07 | 0.14 | 1.63 | 0.06 | 0.12 |
| 0926-0019 | 5.56 | 0.32 | 0.16 | **2.46** | 0.25 | 0.13 | 2.57 | 0.25 | 0.13 |
| 0926-0022 | 1.83 | 0.08 | 0.21 | **1.75** | 0.11 | 0.20 | 1.81 | 0.08 | 0.16 |
| 0926-0039 | 0.53 | 0.07 | 0.28 | 0.43 | 0.12 | 0.23 | **0.43** | 0.06 | 0.18 |
| 0926-0064 | 1.76 | 0.07 | 0.24 | 1.62 | 0.07 | 0.18 | **1.60** | 0.07 | 0.18 |
| Mean | 1.41 | 0.11 | 0.19 | 2.99 | 0.15 | 0.15 | **1.03** | **0.10** | **0.13** |



(a) Sequence 0001



(b) Sequence 0015

Fig. 7. Time-consuming results of the semantic-geometric fusion tracking module in sequences filled with objects under the KITTI-tracking dataset.

uses the Kalman filter for 3D-BB tracking in lidar coordinate. PC3T [38] uses the ground truth of ego localization provided by KITTI dataset to transform the detected object into map coordinate. Then, it performs absolute motion tracking based on the object kinetic model.
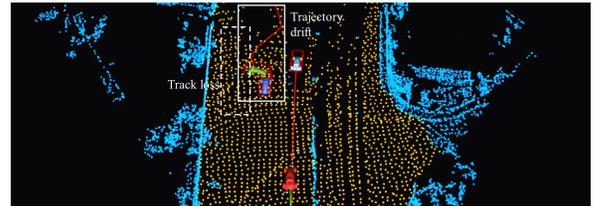
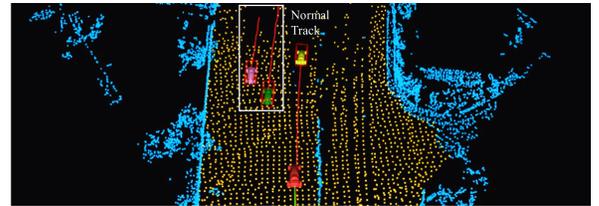| | $\text{MOTA}_{2d}$ | $\text{MOTP}_{2d}$ | Time [ms] |
|---|---|---|---|
| **AB3DMOT** | 68.83 | 87.23 | 4.82 |
| **PC3T** | **85.33** | 87.18 | **4.5** |
| **SGF-MOT** | 77.55 | **87.27** | 5.01 |

As shown in Table III, thanks to the accurate point cloud, three methods have little difference in the MOTP metric, which illustrates the estimation accuracy of the tracker. Since the object is transformed into map coordinate using the prior ego ground truth, PC3T [38] tracker only needs to consider the uncertainty caused by absolute object motion between frames. Moreover, it achieves the best results on the MOTA metric. The proposed SGF-MOT tracker does not use ego localization



(a) Semantic only MOT method



(b) Geometry only MOT method



(c) Fusion MOT Method

Fig. 8. Comparison of different tracking methods in sequence 0926-0056 under the KITTI-tracking dataset. Both semantic only and fusion methods can track objects robustly, but geometric only method have the problems of object tracking loss and trajectory drift.

information but also improves the tracking robustness by introducing a least-squares estimator fusing geometric and semantic constraints. So we achieve better MOTA results than the AB3DMOT [37] tracker.

Like the other two methods, our method runs in real-time on the CPU. Fig.7 shows the time-consuming results when the fusion tracking module executes sequence 0001 and 0015. It can be concluded that the time-consuming of our method is always less than 25ms, even during peak hours.

Then, we selected 8 sequences containing moving objects

from KITTI-raw City and Road sequences. The semantic-geometric fusion tracking method (Fusion Method) is compared with methods using only semantic boxes (Semantic Only) and geometric point clouds (Geometry Only) by evaluating the object trajectory accuracy in map coordinate.

The object localization accuracy of geometry-only method exhibit a distinct "bipolar" distribution. Furthermore, we found that it can only achieve stable and accurate tracking for objects perceived within 15m. The main reason for this phenomenon is that the lidar points successfully hit on the object will decrease rapidly with the increase of measurement range. When the constraints are not sufficient, and the initial value of optimization is not precise enough, the optimization results diverge. However, if the object is in the near field of view, the estimation method based on the measured point cloud will be more accurate than that of 3D-BB inferred by the semantic detection module. We can draw this conclusion by comparing the semantic and converged geometric tracking results. In addition, the object motion trajectories of different tracking methods under sequence 0926-0056 in Fig.8 can also qualitatively verify our conclusion.

We hope to use the directly measured object point cloud to improve state estimation accuracy. Meanwhile, the tracker's robustness should be guaranteed when the object moves away from the robot. Therefore, using a fusion estimation method with semantic and geometric information seems natural. As shown in Table IV, our tracking method does not exhibit significant object trajectory drift and achieves better localization accuracy on the test sequences.

## V. CONCLUSIONS AND FUTURE WORK

In this paper, we propose a novel multi-object lidar odometry (MLO) system that aims to solve the problem of simultaneous localization, mapping, and multi-object tracking using only a lidar sensor. Specifically, we propose a fused least squares estimator using the semantic bounding box and object point cloud for object state update of multi-object modules (SGF-MOT). In the mapping module, dynamic semantic objects are detected based on the maintained absolute trajectory tracking list, and the system can achieve reliable mapping in highly dynamic scenarios. Experiments on open datasets show that, compared with the state of art MOT works, the SGF-MOT method can achieve a better balance between the accuracy and robustness of object tracking than the semantic-only and geometry-only methods. Meanwhile, compared with the representative feature-based 3D lidar SLAM methods, the MLO system with the SGF-MOT tracker can provide more stable and accurate localization in complex scenes. In the future, based on the existing MLO system, we will do an in-depth analysis and make use of variable object information to achieve a more robust environment awareness and mapping system for mobile robots.

## APPENDIX

The appendix will explain the Jacobian of Eq. (11) and (12) in detail. In the process of derivation, we omit the superscript and subscript of variables without ambiguity. First, we use the perturbation model $\delta_\theta$ for the rotating part to obtain the new closest point ${}^l\mathbf{\Pi}_{i-1,r}^{b}{}^{*}$:

$$
\begin{aligned}
{}^l\mathbf{\Pi}_{i-1,r}^{b}{}^{*} &= -\left(\mathbf{I} + [\delta_\theta]_\times\right)\mathbf{Rnt}^T\left(\mathbf{I} + [\delta_\theta]_\times\right)\mathbf{Rn} \\
&\qquad + \left(\mathbf{I} + [\delta_\theta]_\times\right)\mathbf{Rn}d \\
&= -\left(\mathbf{Rnt}^T\mathbf{Rn} + [\delta_\theta]_\times\mathbf{Rnt}^T\mathbf{Rn} + \mathbf{Rnt}^T[\delta_\theta]_\times\mathbf{Rn}\right) \\
&\qquad + \mathbf{Rn}d + [\delta_\theta]_\times\mathbf{Rn}d \\
&= -\mathbf{Rnt}^T\mathbf{Rn} + \mathbf{Rn}d - [\delta_\theta]_\times\mathbf{Rnt}^T\mathbf{Rn} \\
&\qquad - \mathbf{Rnt}^T[\delta_\theta]_\times\mathbf{Rn} + [\delta_\theta]_\times\mathbf{Rn}d \\
&= {}^l\hat{\mathbf{\Pi}}_{i-1,r}^{b} + \left[\mathbf{Rnt}^T\mathbf{Rn}\right]_\times\delta_\theta + \mathbf{Rnt}^T[\mathbf{Rn}]_\times\delta_\theta \\
&\qquad - \mathbf{Rn}d\delta_\theta
\end{aligned}
\tag{17}
$$

Then, the Jacobian of $\mathbf{e}_{b_r}$ with respect to the axis angle increment $\delta_\theta$ is:

$$
\frac{\partial \mathbf{e}_{b_r}}{\partial \delta_\theta} = -\left[\mathbf{Rnt}^T\mathbf{Rn}\right]_\times - \mathbf{Rnt}^T[\mathbf{Rn}]_\times + \mathbf{Rn}d \tag{18}
$$

Similarly, we can get the new closest point ${}^l\mathbf{\Pi}_{i-1,r}^{b}{}^{*}$ by using the perturbation model $\delta_\mathbf{t}$ for the translation part:

$$
\begin{aligned}
{}^l\mathbf{\Pi}_{i-1,r}^{b}{}^{*} &= -\mathbf{Rn}(\mathbf{t} + \delta_\mathbf{t})^T\mathbf{Rn} + \mathbf{Rn}d \\
&= -\mathbf{Rnt}^T\mathbf{Rn} + \mathbf{Rn}d - \mathbf{Rn}\delta_\mathbf{t}^T\mathbf{Rn} \\
&= {}^l\hat{\mathbf{\Pi}}_{i-1,r}^{b} - \mathbf{Rn}(\mathbf{Rn})^T\delta_\mathbf{t}
\end{aligned}
\tag{19}
$$

The Jacobian of $\mathbf{e}_{b_r}$ with respect to the translation $\delta_\mathbf{t}$ is:

$$
\frac{\partial \mathbf{e}_{b_r}}{\partial \delta_\mathbf{t}} = \mathbf{Rn}(\mathbf{Rn})^T \tag{20}
$$

## REFERENCES

[1] Baidu, "Apolloauto," https://github.com/ApolloAuto/apollo, 2022.

[2] J. Behley and C. Stachniss, "Efficient surfel-based slam using 3d laser range data in urban environments." in *Robotics: Science and Systems*, vol. 2018, 2018, p. 59.

[3] K. Bernardin and R. Stiefelhagen, "Evaluating multiple object tracking performance: the clear mot metrics," *EURASIP Journal on Image and Video Processing*, vol. 2008, pp. 1–10, 2008.

[4] G. Chen, B. Wang, X. Wang, H. Deng, B. Wang, and S. Zhang, "Psf-lo: Parameterized semantic features based lidar odometry," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 5056–5062.

[5] X. Chen, A. Milioto, E. Palazzolo, P. Giguere, J. Behley, and C. Stachniss, "Suma++: Efficient lidar-based semantic slam," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2019, pp. 4530–4537.

[6] Y. Chen and G. Medioni, "Object modelling by registration of multiple range images," *Image and vision computing*, vol. 10, no. 3, pp. 145–155, 1992.

[7] P. Dellenbach, J.-E. Deschaud, B. Jacquet, and F. Goulette, "Ct-icp: Real-time elastic lidar odometry with loop closure," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 5580–5586.

[8] J.-E. Deschaud, "Imls-slam: Scan-to-model matching based on 3d data," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 2480–2485.

[9] A. Dewan, T. Caselitz, G. D. Tipaldi, and W. Burgard, "Motion-based detection and tracking in 3d lidar scans," in *2016 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2016, pp. 4508–4513.

[10] ——, "Rigid scene flow for 3d lidar scans," in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2016, pp. 1765–1770.

[11] M. A. Fischler and R. C. Bolles, "Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography," *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981.

TABLE IV
EVALUATION FOR ABSOLUTE OBJECT TRAJECTORY ACCURACY UNDER KITTI-RAW CITY AND ROAD SEQUENCES.

| Sequence | id | Geometry Only | | | Semantic Only | | | Fusion Method | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | ATE[m] | RTE[m/m] | RRE[deg/m] | ATE[m] | RTE[m/m] | RRE[deg/m] | ATE[m] | RTE[m/m] | RRE[deg/m] |
| 0926-0009 | 87 | **0.11** | 0.11 | 0.35 | 0.25 | 0.33 | 2.10 | 0.17 | 0.13 | 0.36 |
| | 89 | **0.06** | 0.06 | 0.32 | 0.16 | 0.16 | 1.55 | 0.14 | 0.11 | 0.86 |
| 0926-0013 | 1 | **0.21** | 0.15 | 0.93 | 0.39 | 0.29 | 3.85 | 0.40 | 0.26 | 2.76 |
| | 2 | <span style="color:red">2.54</span> | 1.03 | 9.80 | **0.22** | 0.18 | 1.50 | 0.23 | 0.16 | 1.30 |
| 0926-0018 | 6 | 0.11 | 0.12 | 1.77 | 0.12 | 0.13 | 2.53 | **0.08** | 0.12 | 1.09 |
| | 11 | **0.08** | 0.10 | 1.39 | 0.15 | 0.16 | 1.90 | **0.11** | 0.12 | 1.02 |
| 0926-0051 | 21 | <span style="color:red">2.28</span> | 1.86 | 2.83 | 0.20 | 0.18 | 1.73 | **0.18** | 0.13 | 1.22 |
| | 26 | **0.07** | 0.12 | 0.27 | 0.14 | 1.54 | 1.21 | 0.09 | 0.09 | 1.13 |
| 0926-0084 | 14 | <span style="color:red">4.81</span> | 1.11 | 7.11 | 0.27 | 0.38 | 6.25 | **0.23** | 0.35 | 6.14 |
| 0926-0015 | 32 | **0.14** | 0.25 | 1.67 | 0.19 | 0.21 | 2.56 | 0.16 | 0.16 | 1.96 |
| | 35 | <span style="color:red">3.29</span> | 1.07 | 2.28 | 0.30 | 0.28 | 5.44 | **0.27** | 0.27 | 7.15 |
| 0926-0028 | 1 | **fail** | **fail** | **fail** | 0.31 | 0.22 | 2.22 | **0.31** | 0.21 | 2.06 |
| 0926-0032 | 15 | 0.18 | 0.16 | 0.43 | 0.20 | 0.21 | 1.06 | **0.17** | 0.18 | 1.31 |
| Mean | | 1.16 | 0.52 | 2.43 | 0.22 | 0.34 | 2.64 | **0.19** | **0.17** | **2.19** |

[12] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *2012 IEEE conference on computer vision and pattern recognition*. IEEE, 2012, pp. 3354–3361.

[13] P. Geneva, K. Eckenhoff, Y. Yang, and G. Huang, "Lips: Lidar-inertial 3d plane slam," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 123–130.

[14] R. Hartley and A. Zisserman, *Multiple view geometry in computer vision*. Cambridge university press, 2003.

[15] M. Henein, J. Zhang, R. Mahony, and V. Ila, "Dynamic slam: the need for speed," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 2123–2129.

[16] K. Huang and Q. Hao, "Joint multi-object detection and tracking with camera-lidar fusion for autonomous driving," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 6983–6989.

[17] R. E. Kalman, "A new approach to linear filtering and prediction problems," 1960.

[18] A. Kim, A. Ošep, and L. Leal-Taixé, "Eagermot: 3d multi-object tracking via sensor fusion," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 11 315–11 321.

[19] K. Koide, M. Yokozuka, S. Oishi, and A. Banno, "Voxelized gicp for fast and accurate 3d point cloud registration," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 11 054–11 059.

[20] W. Liu, W. Sun, and Y. Liu, "Dloam: Real-time and robust lidar slam system based on cnn in dynamic urban environments," *IEEE Open Journal of Intelligent Transportation Systems*, 2021.

[21] Z. Liu and F. Zhang, "Balm: Bundle adjustment for lidar mapping," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 3184–3191, 2021.

[22] K. Madsen, H. B. Nielsen, and O. Tingleff, "Methods for non-linear least squares problems," 2004.

[23] A. Milioto, I. Vizzo, J. Behley, and C. Stachniss, "Rangenet++: Fast and accurate lidar semantic segmentation," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2019, pp. 4213–4220.

[24] T. Moore and D. Stouch, "A generalized extended kalman filter implementation for the robot operating system," in *Intelligent autonomous systems 13*. Springer, 2016, pp. 335–348.

[25] F. Moosmann and T. Fraichard, "Motion estimation from range images in dynamic outdoor scenes," in *2010 IEEE International Conference on Robotics and Automation*. IEEE, 2010, pp. 142–147.

[26] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos, "Orb-slam: a versatile and accurate monocular slam system," *IEEE transactions on robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.

[27] J. Park, Y. Cho, and Y.-S. Shin, "Nonparametric background model-based lidar slam in highly dynamic urban environments," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 12, pp. 24 190–24 205, 2022.

[28] P. Pfreundschuh, H. F. Hendrikx, V. Reijgwart, R. Dubé, R. Siegwart, and A. Cramariuc, "Dynamic object aware lidar slam based on automatic generation of training data," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 11 641–11 647.

[29] T. Qin and S. Cao, "A-loam: A lidar odometry and mapping," https://github.com/HKUST-Aerial-Robotics/A-LOAM, 2019.

[30] P. Ruchti and W. Burgard, "Mapping with dynamic-object probabilities calculated from single 3d range scans," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 6331–6336.

[31] T. Shan and B. Englot, "Lego-loam: Lightweight and ground-optimized lidar odometry and mapping on variable terrain," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 4758–4765.

[32] S. Shi, X. Wang, and H. Li, "Pointrcnn: 3d object proposal generation and detection from point cloud," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 770–779.

[33] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of rgb-d slam systems," in *2012 IEEE/RSJ international conference on intelligent robots and systems*. IEEE, 2012, pp. 573–580.

[34] B. Triggs, P. F. McLauchlan, R. I. Hartley, and A. W. Fitzgibbon, "Bundle adjustment—a modern synthesis," in *International workshop on vision algorithms*. Springer, 1999, pp. 298–372.

[35] H. Wang, C. Wang, C.-L. Chen, and L. Xie, "F-loam: Fast lidar odometry and mapping," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 4390–4396.

[36] S. Wang, P. Cai, L. Wang, and M. Liu, "Ditnet: End-to-end 3d object detection and track id assignment in spatio-temporal world," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 3397–3404, 2021.

[37] X. Weng, J. Wang, D. Held, and K. Kitani, "3d multi-object tracking: A baseline and new evaluation metrics," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 10 359–10 366.

[38] H. Wu, W. Han, C. Wen, X. Li, and C. Wang, "3d multi-object tracking in point clouds based on prediction confidence-guided data association," *IEEE Transactions on Intelligent Transportation Systems*, 2021.

[39] Z. Yang, Y. Sun, S. Liu, and J. Jia, "3dssd: Point-based 3d single stage object detector," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 11 040–11 048.

[40] D. Zermas, I. Izzat, and N. Papanikolopoulos, "Fast segmentation of 3d point clouds: A paradigm on lidar data for autonomous vehicle applications," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 5067–5073.

[41] J. Zhang and S. Singh, "Loam: Lidar odometry and mapping in real-time." in *Robotics: Science and Systems*, vol. 2, no. 9. Berkeley, CA, 2014, pp. 1–9.