# Large-Scale LiDAR Consistent Mapping using Hierarchical LiDAR Bundle Adjustment

Xiyuan Liu, Zheng Liu, Fanze Kong, and Fu Zhang

*Abstract*—Reconstructing an accurate and consistent large-scale LiDAR point cloud map is crucial for robotics applications. The existing solution, pose graph optimization, though it is time-efficient, does not directly optimize the mapping consistency. LiDAR bundle adjustment (BA) has been recently proposed to resolve this issue; however, it is too time-consuming on large-scale maps. To mitigate this problem, this paper presents a globally consistent and efficient mapping method suitable for large-scale maps. Our proposed work consists of a bottom-up hierarchical BA and a top-down pose graph optimization, which combines the advantages of both methods. With the hierarchical design, we solve multiple BA problems with a much smaller Hessian matrix size than the original BA; with the pose graph optimization, we smoothly and efficiently update the LiDAR poses. The effectiveness and robustness of our proposed approach have been validated on multiple spatially and timely large-scale public spinning LiDAR datasets, i.e., KITTI, MulRan and Newer College, and self-collected solid-state LiDAR datasets under structured and unstructured scenes. With proper setups, we demonstrate our work could generate a globally consistent map with around 12% of the sequence time.

*Index Terms*—LiDAR Bundle Adjustment, Pose Graph Optimization, High-Resolution Mapping.

## I. INTRODUCTION

Reconstructing the three-dimensional (3D) high-resolution map is of great significance in the fields of robotics, environmental and civil engineering. This 3D map could be used as a prior for autonomous service robots and as an information model for buildings and geographical measurements. Compared with the traditional 3D laser scanner, the light detection and ranging (LiDAR) sensor extraordinarily fits into this purpose due to its fast scanning rate. Moreover, it is more lightweight, cost-effective and flexible to be carried on multiple platforms, e.g., ground or aerial vehicles and hand-held devices. In this paper, we focus on developing an accurate and consistent LiDAR mapping method for large-scale maps.

Rich research results have been presented on LiDAR-based mapping algorithms [1]–[3], which generate both point cloud maps and LiDAR odometry. Due to the accumulation of scan-to-map registration errors, odometry drift usually appears and further leads to divergence in the point cloud map. The most well-known method to refine the mapping quality (closing the gap) is pose graph optimization, which minimizes the relative pose errors between two LiDAR frames. In pose graph optimization, the relative pose estimation is assumed

to follow Gaussian distribution. However, this approximation might be overestimated in reality [4]. Moreover, the pose graph optimization does not directly optimize the consistency of the point cloud that the divergence within the point cloud map might only be narrowed but not fully eliminated (or not even aware of). This phenomenon is more obvious when the wrong loops are detected, or incorrect relative transformation estimations happen.

LiDAR bundle adjustment (BA) approach [5, 6] directly optimizes the mapping consistency by minimizing the overall point-to-plane distance, which often leads to high mapping quality that is necessary for mapping applications. In [5], the plane parameters are analytically solved first such that the final optimization problem is only related to the LiDAR pose. In [6], the plane parameters are eliminated in each iteration of the optimization by a Schur complement trick as in visual BA [7, 8]. Either way, the resultant optimization is (at least) the dimension of the LiDAR pose number $N$, requiring $O(N^3)$ time to solve [9]. The cubic growth of the computation time has prohibited the bundle adjustment for large-scale maps with large pose numbers.

To address the above issues, we propose a hierarchical LiDAR BA method to globally optimize the mapping consistency while maintaining time efficiency. The method constructs a pyramid structure of frame poses (see Fig. 1) and conduct a bottom-up hierarchical bundle adjustment and a top-down pose graph optimization (see Fig. 2). The bottom-up process conducts a hierarchical bundle adjustment within local windows from the bottom layer (local BA) to the top layer frames (global BA). Such design benefits the computation time since the process of local BA in each layer is suitable for parallel processing and the time complexity of each local BA is relatively low due to the small number of poses involved. One issue in the bottom-up process is that it neglects features co-visible across different local windows, which could lower the accuracy. To mitigate this issue, the top-down process constructs a pose graph from the top to the bottom layers and distributes the errors by pose graph optimization. The two process iterates until convergence.

With the hierarchical bundle adjustment design, we could both directly optimize the consistency of the planar surfaces within the point cloud and avoid solving a cost function with large dimensions. With pose graph optimization, we properly update the entire LiDAR poses towards convergence in a fast and reliable manner. To retain the smoothness between every two adjacent keyframes, we keep an overlap area between them by setting the stride size smaller than the window size. To

Xiyuan Liu, Zheng Liu, Fanze Kong and Fu Zhang are with the Department of Mechanical Engineering, The University of Hong Kong, Hong Kong. (email: {xliuaa,u3007335,kongfz}@connect.hku.hk, fuzhang@hku.hk) (*Corresponding author: Fu Zhang*).

further boost the optimization speed, we have applied a filter to remove outlier points and implemented CPU-based parallel processing when constructing the pyramid. In summary, our contributions are as follows:

- We propose a hierarchical bundle adjustment method to globally optimize the LiDAR mapping consistency and odometry accuracy. Our proposed approach improves the mapping quality given a good initial pose trajectory (e.g., from a pose graph optimization) and even closes the gap when the initial pose trajectory has large drifts.
- The effectiveness of our proposed work has been validated on multiple public mechanical spinning LiDAR datasets and our self-collected solid-state LiDAR dataset in both structured and unstructured scenes.

## II. RELATED WORKS

Multiple approaches have been discussed in the literature on improving the mapping quality, which could be mainly divided into two categories: pose graph optimization and plane (bundle) adjustment-based methods. In pose graph optimization, the relative transformation (pose constraint) between two frames is estimated by ICP [10] or its variants [11, 12]. This relative pose error is then weighted by the information matrix, which is usually the inverse of the corresponding Hessian [13] or simply a constant matrix [14]. The pose graph is optimized when the summed relative pose errors are minimized. Though computationally efficient, one important issue of pose graph optimization is that it does not directly optimize the consistency of the point cloud. Due to incorrect estimation or imprecise modelling of the relative pose constraint, the pose graph might converge to a local minimum that considerable divergence within the point cloud could still exist [3, 15].

The plane adjustment method directly optimizes the consistency of the point cloud by minimizing the summed point-to-plane distance. In plane adjustment, each plane feature is represented by two parameters, i.e., the distance from its center to the viewpoint and the estimated plane normal vector [6]. In [16], authors concurrently optimize both the LiDAR poses and the geometric plane features. This method needs to maintain and update the parameters of all features during the optimization, whereas the total number of features will rapidly grow when the scale of the map enlarges, leaving a huge dimension of the cost function to solve. Though with the Schur complement technique, the optimization variables could be reduced to LiDAR poses only, this method is prone to generate glitches in pose estimation during optimization in real-world practices.

The bundle adjustment method improves the technique of plane adjustment by eliminating the feature parameter prior to the optimization using a closed-form solution [5]. In [5], authors segment the point cloud into multiple voxels, each containing a plane feature. The original point-to-plane minimization problem is transformed into the minimization of the eigenvalue of points covariance in each voxel. Such a method needs to iterate through every point within each feature to derive the Hessian matrix, whose time complexity is the square of the number of points, causing a great computation demand.

This problem is addressed in the following-up work [9], which aggregates all the points of a feature observed by the same pose and fundamentally removes the dependence of time complexity on the point number. Despite this, in all the above-mentioned plane and bundle adjustment methods, the computation complexity is still cubic to the pose number, which is not practical for large-scale maps. Moreover, when the divergence in the map is larger than or approximates the maximum voxel size, these methods might have a slow convergence rate.

Our proposed hierarchical bundle adjustment approach takes advantage of both BA and pose graph optimization. We use BA to minimize the point-to-plane distance directly and pose graph optimization to smoothly and efficiently update the LiDAR poses to avoid glitches in pose estimation. With the hierarchical design, we could parallelly solve multiple BA problems with a much smaller Hessian matrix size compared with the original problem using [5]. Moreover, we could flexibly set BA parameters from the bottom layers to the top layers in accordance with the quality of the initial pose trajectory.

## III. METHODOLOGY

### A. Overview

The system workflow of our proposed method is illustrated in Fig. 2. The inputs are raw or deskewed points from each LiDAR scan and an initial estimation of their corresponding poses in the global frame, which could be obtained from the general LiDAR odometry or simultaneous locomotion and mapping (SLAM) algorithms. The method consists of two processes, bottom-up (see Sec. III-B) and top-down (see Sec. III-C), which iterate until convergence. In the bottom-up process, a local BA is performed on LiDAR frames within smaller local windows to construct keyframes from the first layer to the second layer (see Fig. 1). This process is hierarchically performed until the optimal layer number is met and a global BA is performed on the top layer keyframes. Then a pose graph is constructed using the factors contributed from each optimized layer and between adjacent layers (see Fig. 1).

As shown in Fig. 1, the term *first layer*, also referred to as *bottom layer* in the context above, describes the collection of initial LiDAR frames and poses. Similarly, the term *second layer* represents the collection of LiDAR keyframes and poses created from the first layer using the local BA. The term *top layer* means the collection of the last remaining LiDAR keyframes (in Fig. 1, the top layer refers to the third layer). The process of hierarchically creating LiDAR keyframes from the bottom layer to the top layer is called the bottom-up process. The process of updating bottom layer LiDAR poses by pose graph optimization is called the top-down process.

### B. Bottom-Up Hierarchical BA

We denote $\mathbb{F}_j^i$ the $j$-th LiDAR frame of the $i$-th layer and $\mathbf{x}_j^i \triangleq \mathbf{T}_j^i = (\mathbf{R}, \mathbf{t}) \in SE(3)$ its corresponding pose. We denote $\mathbf{T}_{j,k}^i$ the relative pose between $\mathbf{T}_j^i$ and $\mathbf{T}_k^i$, i.e., $\mathbf{T}_{j,k}^i = \left(\mathbf{T}_j^i\right)^{-1} \cdot \mathbf{T}_k^i$. It is noted that points in $\mathbb{F}_j^i$ is represented in the LiDAR local frame, and $\mathbf{T}_j^i$ is in the global frame. We denote $w$ as the local window size and $s$ as the stride size during the bottom-up construction of the LiDAR keyframes, as illustrated
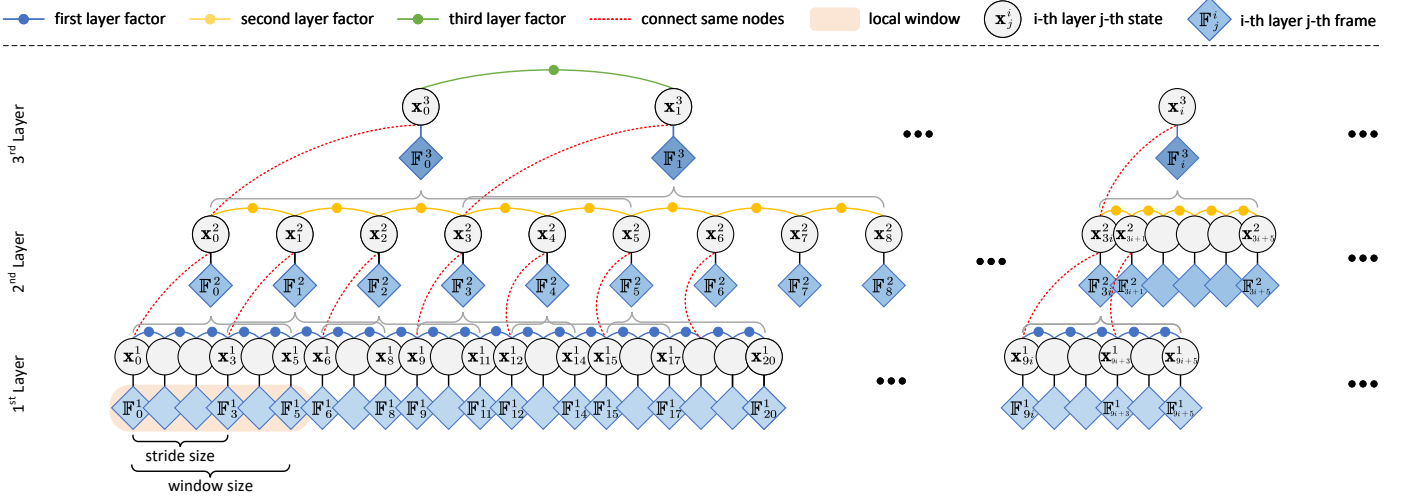
Fig. 1: Pyramid structure of the proposed hierarchical bundle adjustment with *layer number l=3*, *stride size s=3* and *window size w=6*. The first, second and third layer factors connect every two adjacent nodes within the same layer. The red dashed line connects the nodes ought to be the same, e.g., $\mathbf{x}_{9i}^1$, the first node of the local window from the lower layer to the node constructed from this local window from the upper layer, e.g., $\mathbf{x}_{3i}^2$ and $\mathbf{x}_i^3$.
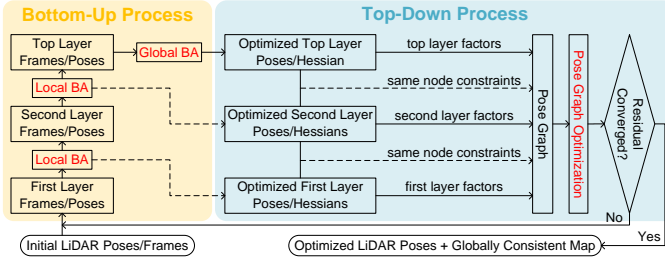


Fig. 2: System overview. The light yellow region depicts the bottom-up process and the light blue region depicts the top-down process.

in Fig. 1. Assume we have $N_i$ total number of LiDAR frames from the $i$-th layer. In the bottom-up process, a local BA is performed in each local window using the provided initial pose trajectory, and the relative pose between each frame and the first frame in this window is optimized. The derived Hessian matrix $\mathbf{H}$ from the BA in each local window is also recorded and used as the information matrix in the latter top-down pose graph construction. Given a local window of $w$ LiDAR frames $\{\mathbb{F}_{sj+k}^i \mid j = 0, \cdots, \lfloor \frac{N_i-w}{s} \rfloor; k = 0, \cdots, w-1\}$ in $i$-th layer and its optimized relative poses $\mathbf{T}_{j,k}^{i*}$ in the window, we aggregate these frames into a keyframe for the $(i+1)$-th layer. Points in this keyframe are all transformed into the first frame of the local window, and the pose of the keyframe, denoted by $\mathbf{T}_j^{i+1}$, is set as the pose of the first frame optimized in the preceding local window, i.e.,

$$\mathbb{F}_j^{i+1} \triangleq \bigcup_{k=0}^{w-1} \left( \mathbf{T}_{s\cdot j, s\cdot j+k}^{i*} \cdot \mathbb{F}_{s\cdot j+k}^i \right)$$
$$\mathbf{T}_j^{i+1} = \mathbf{T}_{s\cdot j}^{i*} = \prod_{k=1}^{j} \mathbf{T}_{s\cdot k-s, s\cdot k}^{i*} \tag{1}$$

This process could be repeatedly performed from the lower layer to the upper layer until the optimal layer number $l$ is reached. It is noticed that the construction of new keyframes (local BA) does not rely on the frames outside the local window, making it suitable to use multiple local windows for

parallel processing in the same layer. Suppose we have $N$ total number of LiDAR frames, i.e., $N_1 = N$, and each time we choose to aggregate $w$ frames from the lower layer into one frame to the upper layer with stride frame size $s$. Let $n$ be the number of threads that could be used for parallel processing. Since the computation time of BA is $O\left(M^3\right)$ with $M$ is the number of involved poses, we could derive the overall time consumption $O\left(T_l\right)$ for a $l$-th layer pyramid.

The total time consumption of $l$-th layer pyramid includes that consumed by the local BA in each layer and that from the global BA in the top layer. For a $l$-th layer pyramid, the number of local windows in $i$-th layer ($i < l$) is $\frac{N}{s^i}$ and each local window consumes $O(w^3)$ of time. With $n$ number of parallel threads, the total time consumption of the local BA equals to the sum of the local BA in each layer which is $w^3 \cdot \left( \sum_{i=1}^{l-1} \frac{N}{s^i} \cdot \frac{1}{n} \right)$ and the global BA in the $l$-th layer takes $O\left( \left( \frac{N}{s^{l-1}} \right)^3 \right)$ of time. In summary, $O\left(T_l\right)$ is expressed as

$$T_l = \begin{cases} N^3 & (l=1) \\ \dfrac{w^3}{n} \cdot \displaystyle\sum_{i=1}^{l-1} \dfrac{N}{s^i} + \left( \dfrac{N}{s^{l-1}} \right)^3 & (l>1) \end{cases} \tag{2}$$

We view $T_l$ as a function of $l$ and calculate the optimal $l^*$ by letting the derivative of $T_l$ equal to zero, which leads to

$$l^* = \left\lfloor \frac{1}{2} \log_s \left( \frac{3N^2 \left(s^3 - s\right) n}{w^3} \right) \right\rfloor \tag{3}$$

Fig. 3 shows an example of the computation time $T_l$ versus the layer number $l$ at different frame number $N$. As can be seen, as the layer number increases from $l=1$ (original BA) to $l^*$, the computation time is greatly reduced, suggesting the effectiveness of the proposed hierarchical BA. When $l > l^*$, the computation time does not increase much and keeps almost constant, suggesting that any layer number greater than $l^*$ will work equally well.
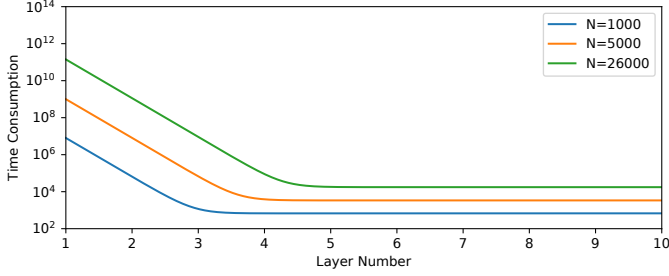
Fig. 3: Example plot of time consumption $T_l$ w.r.t. the layer number $l$ and the pose number $N$ using $w=10$, $s=5$ and $n=8$. It is seen the overall time consumption is significantly reduced from the original BA ($l=1$) to our proposed hierarchical BA ($l=3$ or $l=4$).

For feature extraction and association in each layer of the bottom-up hierarchical BA, we use an adaptive voxelization method proposed in [5], which extracts plane features of different sizes suitable for environments of different structures. To extract these various size plane features, the entire point cloud, after being transformed into the same global frame using the initial pose trajectory, is split into multiple voxels of size $V$, each goes through a plane test by checking the minimum and maximum eigenvalue ratio of the contained points is smaller than a threshold, i.e., $\frac{\lambda_1}{\lambda_3} < \theta$. If the plane test passes, points in the voxel will be viewed as lying on the same plane and used in the BA. Otherwise, the voxel will be recursively split until the contained points form a plane.

The above adaptive voxelization process is time-consuming when the number of points is extremely large. To mitigate this issue, we notice that points that are not considered as plane features in lower layers will not form a plane in upper layers either. We thus use only the plane feature points from each voxel in local BA to construct the keyframe for the upper layer in the bottom-up process. This procedure further saves time for next-layer adaptive-voxel map construction and improves the computation accuracy in local BA.

### C. Top-Down Pose Graph Optimization

The top-down pose graph optimization process aims to reduce the pose estimation errors in the bottom-up hierarchical BA process, which considers only features co-visible in the same local windows but ignores those observed across different local windows. As shown in Fig. 1, the pose graph is constructed in a top-down manner in the pyramid structure. In each layer of the pyramid, the factors are relative poses between adjacent frames. Specifically, the cost item in the $i$-th layer factor between node $\mathbf{x}_j^i$ and $\mathbf{x}_{j+1}^i$ is defined as

$$\mathbf{e}_{j,j+1}^i = \mathrm{Log}\left(\left(\mathbf{T}_{j,j+1}^{i*}\right)^{-1}\left(\mathbf{T}_j^i\right)^{-1}\mathbf{T}_{j+1}^i\right)^{\vee}$$
$$c\left(\mathbf{x}_j^i, \mathbf{x}_{j+1}^i\right) = \left(\mathbf{e}_{j,j+1}^i\right)^T \left(\mathbf{\Omega}_{j,j+1}^i\right)^{-1} \mathbf{e}_{j,j+1}^i \quad (4)$$

where $\mathbf{\Omega}_{j,j+1}^i$ is correlation matrix between pose $\mathbf{x}_j^i$ and $\mathbf{x}_{j+1}^i$, $i \in \mathcal{L}$ and $j \in \mathcal{F}^i$, and is computed by inverting the Hessian matrix $\mathbf{H}$ obtained from the bottom-up hierarchical BA process. The $\mathcal{L} = \{1, \cdots, l\}$ represents the set of $l$ layers and $\mathcal{F}^i = \{0, \cdots, N_i - 1\}$ represents the set of $N_i$ numbers. Since the node $\mathbf{x}_j^{i+1}$ and $\mathbf{x}_{s \cdot j}^i$ are essentially the same, we have

$$\mathbf{T}_j^i = \mathbf{T}_{s \cdot j}^{i-1} = \cdots = \mathbf{T}_{s^{i-1} \cdot j}^1, \ \forall i \in \mathcal{L}, j \in \mathcal{F}^i$$
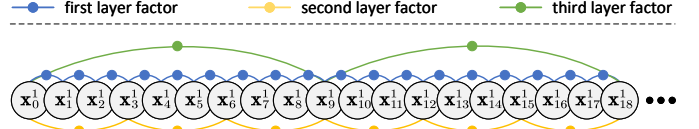


Fig. 4: Final factor graph of our proposed approach with *layer number l=3, stride size s=3* and *window size w=6.*

Therefore, the cost item in (4) reduces to

$$\mathbf{e}_{j,j+1}^i = \mathrm{Log}\left(\left(\mathbf{T}_{j,j+1}^{i*}\right)^{-1}\left(\mathbf{T}_{s^{i-1} \cdot j}^1\right)^{-1}\mathbf{T}_{s^{i-1} \cdot (j+1)}^1\right)^{\vee}$$
$$c\left(\mathbf{x}_{s^{i-1} \cdot j}^1, \mathbf{x}_{s^{i-1} \cdot (j+1)}^1\right) = \left(\mathbf{e}_{j,j+1}^i\right)^T \left(\mathbf{\Omega}_{j,j+1}^i\right)^{-1} \mathbf{e}_{j,j+1}^i \quad (5)$$

where $i \in \mathcal{L}$ and $j \in \mathcal{F}^i$, and the original pose graph in Fig. 1 reduces to Fig. 4. It is noted that when $s < w$, the frames appearing within the overlap area of consecutive local windows could contribute multiple cost items (each local window contributes one). The objective function to be minimized is thus

$$f(\mathbb{F}, \mathcal{X}) = \sum_{i \in \mathcal{L}} \sum_{j \in \mathcal{F}^i} c\left(\mathbf{x}_{s^{i-1} \cdot j}^1, \mathbf{x}_{s^{i-1} \cdot (j+1)}^1\right) \quad (6)$$

where $\mathbb{F} = \{\mathbb{F}_i^1 \mid i \in \mathcal{F}^1\}$ is the collection of all first layer frames. This factor graph is then solved by the Levenberg-Marquardt method using GTSAM[1].

### IV. EXPERIMENTS

#### A. Accuracy Analysis

*1) Initial Odometry with Loop Closure:* In this section, we take the odometry results from the state-of-the-art (SOTA) LiDAR SLAM algorithms [1]–[3] as the input and further optimize them with our proposed work. We demonstrate that our work could both improve the mapping quality (consistency) and the pose estimation accuracy even when the initial poses have already been pose-graph optimized. The BA and pyramid parameters used in all our following experiments, without further specification, are listed in Table I. The optimal $l^*$ is obtained by calculation using (3) based on the actual pose number $N$ in each sequence. For data length $N < 5 \times 10^3$ (KITTI [17], MulRan [18] and Newer College [19]), we have $l^* = 3$ and for larger sequence, i.e., $N \geq 5 \times 10^3$ (New College [20]), we have $l^* = 4$.

TABLE I: Hierarchical Bundle Adjustment Parameter Setting

| Parameter | Value | Description |
|---|---|---|
| $w$ | 10 | window size |
| $s$ | 5 | stride size |
| $n$ | 8 | threads number for parallel processing |
| $\theta_{local}$ | 0.05 | eigenvalue ratio threshold for local BA |
| $V_{local}$ | 4 | initial voxel size for local BA |
| $\theta_{global}$ | 0.1 | eigenvalue ratio threshold for global BA |
| $V_{global}$ | 4 | initial voxel size for global BA |

We first test on the KITTI dataset [17] and use the pose estimation results from MULLS with loop closure [3] as our initial input. The RMSE of the absolute rotation (degree) and translation (meter) errors are summarized in Table II. We

---

[1]https://github.com/borglab/gtsam

TABLE II: RMSE of the ATE ($°/m$) on KITTI Dataset with Loop Closure

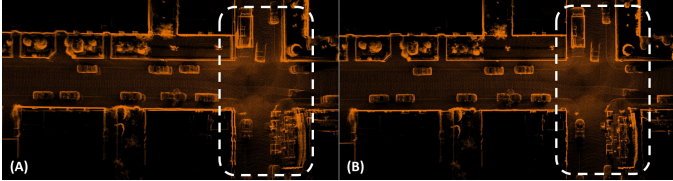| Method | Seq. 00* | Seq. 01 | Seq. 02* | Seq. 03 | Seq. 04 | Seq. 05* | Seq. 06* | Seq. 07* | Seq. 08* | Seq. 09* | Seq. 10 | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Proposed** | **0.7/0.8** | **0.9/1.9** | **1.2**/5.1 | **0.7/0.6** | **0.1**/0.8 | **0.5/0.4** | **0.4/0.2** | **0.4/0.3** | 1.3/2.7 | **0.8**/1.3 | **0.6**/1.1 | **0.7/1.4** |
| MULLS [3] | **0.7**/1.1 | **0.9/1.9** | **1.2**/5.4 | **0.7**/0.7 | 0.2/0.9 | 0.6/1.0 | **0.4**/0.3 | **0.4**/0.4 | 1.3/2.9 | 1.0/2.1 | **0.6**/1.1 | **0.7**/1.6 |
| LiTAMIN2 [21] | 0.8/1.3 | 3.5/15.9 | 1.3/**3.2** | 2.6/0.8 | 2.3/0.7 | 0.7/0.6 | 0.8/0.8 | 0.6/0.5 | **0.9/2.1** | 1.7/2.1 | 1.2/**1.0** | 1.2/2.4 |
| SuMa [22] | **0.7**/1.0 | 3.2/13.8 | 1.7/7.1 | 1.5/0.9 | 1.8/**0.4** | 0.5/0.6 | 0.7/0.6 | 1.1/1.0 | 1.2/3.4 | 0.8/**1.1** | 0.9/1.3 | 1.1/3.2 |

\* Sequence with loops.



Fig. 5: Mapping result from (A) MULLS [3] and (B) our proposed method on KITTI Seq. 07. The white dashed rectangle emphasizes the divergence.
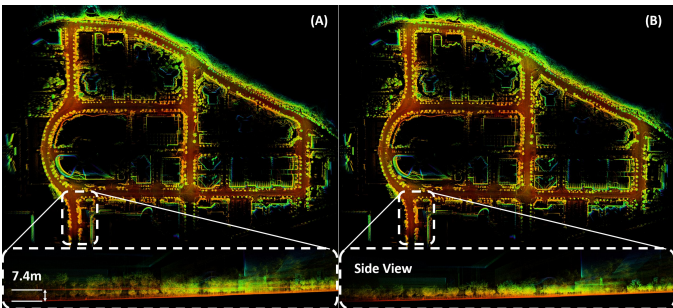


Fig. 6: Point cloud mapping in sequence DCC03 of MulRan. (A): The mapping result from LIO-SAM [2] with loop closure. (B): The mapping result from our proposed method. The start (red) and ending (blue) positions have been emphasized by the white dashed rectangle and enlarged at the bottom (zoomed view is recommended).

choose the absolute trajectory error (ATE) as the evaluation criterion since each LiDAR frame has one unique ground truth. As can be seen, our proposed work could further improve the pose estimation accuracy, especially in translation, even when they are pose-graph optimized. Though our approach does not achieve the best result in every sequence, our work produces the optimal ATE results ($0.7°/1.4m$) on average compared with other SOTA methods. Moreover, since the pose graph optimization does not directly optimize the consistency of the point cloud, the divergence in the map is not fully eliminated in [3] (see Fig. 5). This divergence could be iteratively solved with our proposed hierarchical BA and pose graph optimization, which in return, reduces the ATE on pose estimations.

We then test our proposed method on another spatially large-scale spinning LiDAR dataset, MulRan [18]. The average lengths of the contained sequences DCC, KAIST and RIVER-SIDE, are 4.9km, 6.1km and 6.8km, respectively. Unlike the KITTI dataset, which collects most of the data within the urban area, MulRan dataset includes more challenging scenes from the viaduct, river and woods. We choose to use the pose-graph optimized pose trajectory results from LIO-SAM [2] as our input. The RMSE of the absolute translation error has been summarized in Table III. It is seen our proposed work could still improve the pose estimation accuracy regardless of these challenging unstructured woods scenes. However, these scenes

TABLE III: RMSE of the ATE ($m$) on MulRan Dataset with Loop Closure

| Sequence | | **Proposed** | LIO-SAM [2] | LEGO-LOAM [23] |
|---|---|---|---|---|
| DCC | 01 | **5.20** | 5.67 | 6.95 |
| | 02 | **3.22** | 3.48 | 5.49 |
| | 03 | **2.54** | 2.84 | 6.29 |
| KAIST | 01 | **3.36** | 3.55 | 5.45 |
| | 02 | **3.75** | 3.81 | 5.49 |
| | 03 | **3.53** | 3.60 | 5.70 |
| RIVER-SIDE | 01 | **8.92** | 9.25 | 19.05 |
| | 02 | **7.94** | 8.04 | 16.04 |
| | 03 | **10.26** | 10.37 | 30.91 |
| Avg. | | **5.41** | 5.62 | 11.62 |

make the LIO-SAM generate poor relative pose and covariance estimations, which further leads to partial failure in the loop closure, causing a large divergence in altitude (see Fig. 6). With our proposed hierarchical BA and pose graph optimization mechanism, this divergence could be fully eliminated, and the pose trajectory accuracy is thus further improved.

Lastly, we test our work on the sequence *long_experiment* which is the longest sequence ($N$=26557) in the New College dataset [20]. We choose the FAST-LIO2 [1] to provide the pose trajectory as our initial input. It is noted that these initial poses are generated without loop closure. Table IV shows the absolute translation error of our proposed and other SOTA methods. It is seen our proposed work outperforms other loop closure enabled SOTA methods and achieves the optimal accuracy on this large time-scale sequence.

TABLE IV: RMSE of the ATE ($m$) on Newer College Dataset with Loop Closure

| Method | long_experiment |
|---|---|
| **Proposed** | **0.26** |
| GICP Matching Factor [4] | 0.28 |
| FAST-LIO2 [1] | 0.35 |
| LIO-SAM [2] | 0.53 |
| CT-ICP [24] | 0.58 |

*2) Initial Odometry without Loop Closure:* In this section, we demonstrate that our proposed work can converge well even when the loop is not closed in the initial pose trajectory. We first test on the KITTI dataset [17] using the initial pose trajectory estimated from MULLS [3] without loop closure. The RMSE of the absolute rotation and translation errors are summarized in Table V. As can be seen, other SOTA methods get much worse ATEs due to the lack of loop closure function, whereas our proposed work could still produce reliable pose estimations, with the absolute rotation and translation errors being largely reduced (e.g., Seq. 00 and Seq. 08) and achieving the optimal ATE on average ($0.8°/1.9m$). This is due to the

TABLE V: RMSE of the ATE ($°/m$) on KITTI Dataset without Loop Closure

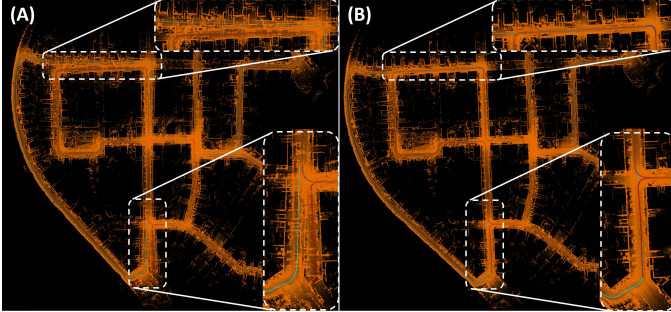| Method | Seq. 00* | Seq. 01 | Seq. 02* | Seq. 03 | Seq. 04 | Seq. 05* | Seq. 06* | Seq. 07* | Seq. 08* | Seq. 09* | Seq. 10 | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Proposed** | 1.0/**1.2** | **1.0/2.4** | 2.3/9.0 | **0.7/0.6** | **0.2**/0.9 | **0.5/0.7** | **0.3/0.2** | **0.4/0.3** | 1.3/2.5 | **0.7/1.5** | **0.5**/1.1 | **0.8/1.9** |
| MULLS [3] | 1.7/6.1 | **1.0/2.4** | 2.4/10.7 | **0.7**/0.7 | **0.2**/0.9 | 1.0/2.4 | 0.4/0.6 | 0.5/0.6 | 1.9/4.3 | 1.0/3.1 | **0.5**/1.1 | 1.0/3.0 |
| Voxel Map [25] | **0.9**/2.8 | 1.9/7.8 | **1.7/6.1** | 1.2/0.7 | 0.6/**0.3** | 0.8/1.2 | 0.4/0.4 | 0.7/0.7 | **1.1/2.3** | 1.0/1.9 | 1.0/1.1 | 1.2/2.9 |
| SuMa [22] | 1.0/2.9 | 3.2/13.8 | 2.2/8.4 | 1.5/0.9 | 1.8/0.4 | 0.7/1.2 | 0.4/0.4 | 0.7/0.5 | 1.5/2.8 | 1.1/2.9 | 0.8/1.3 | 1.4/3.9 |
| LiTAMIN2 [21] | 1.6/5.8 | 3.5/15.9 | 2.7/10.7 | 2.6/0.8 | 2.3/0.7 | 1.1/2.4 | 1.1/0.9 | 1.0/0.6 | 1.3/2.5 | 1.7/2.1 | 1.2/**1.0** | 1.8/5.1 |

\* Sequence with loops.



Fig. 7: Closure of the gap on KITTI dataset Seq. 00 with our proposed method. The mapping result (A) is provided by MULLS [3] without loop closure, and (B) our proposed method. The odometry is colored by the moving distance from the start (red) to the end (blue). The main gaps are detailed by white dashed rectangles. The full experiment video is available on https://youtu.be/CuLnTnXVujw.
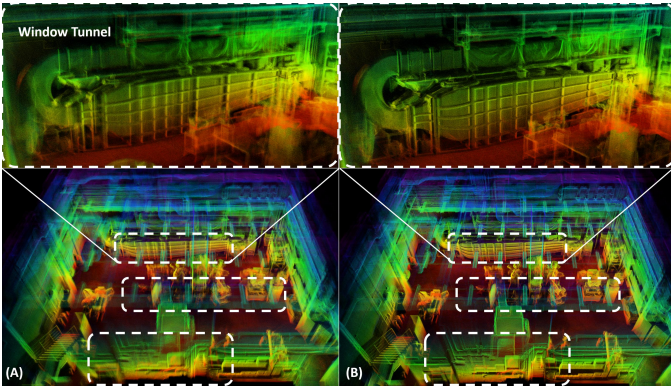


Fig. 9: Reconstructed point cloud map of scene-2 using (A) FAST-LIO2 [1] (MME=-2.60) and (B) our proposed work (MME=-2.69). The main differences are emphasized by the white dashed rectangle. The second row depicts the side view of divergence in height. The third row shows the divergence from the bird-eye's view.
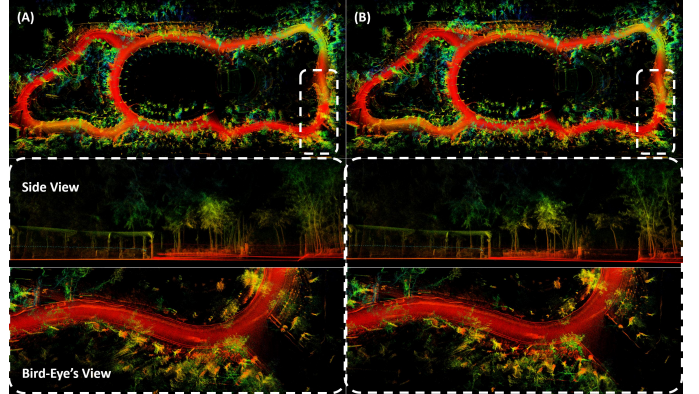


Fig. 8: Reconstructed point cloud map of scene-1 using (A) FAST-LIO2 [1] (MME=-2.99) and (B) our proposed work (MME=-3.06). The main differences are emphasized by the white dashed rectangles, including the walls (bottom), ceiling lights (middle) and wind tunnel (top). Please view our experiment video for more details.

7339 frames. The pyramid parameters are set the same as in Table I except that the initial voxel size is decreased due to the smaller size of the indoor scene ($V_{local}=V_{global}=1m$). The second test scene is an unstructured outdoor park with bush, grassland and woods around (see Fig. 9). The size of this scene is around 95×195m, and the length of this sequence is 3407 frames. We thus adjust the initial voxel size ($V_{local}=V_{global}=2m$) without changing the other pyramid parameters from Table I. For both test sequences, we use the pose estimations from FAST-LIO2 [1] without loop closure as our input. Since the measurement of the ground truth is not available, we instead use the mean map entropy (MME [27]) to evaluate the mapping quality (see Table VI). Since the MME calculates the natural logarithm of the determinant of the covariance matrix, the smaller MME is, the more consistent the point cloud is. It is seen our proposed work could further improve the mapping consistency and close the gap in both structured and unstructured scenes regardless of the LiDAR types.

reason that, in local BA, the strict parameters (see Table I) ensure the frames within each local window do not diverge, and the loose parameters of the top layer global BA could implicitly identify potential divergences, generating correct relative pose constraints (see Fig. 7). Though false feature correspondence matching in global BA might happen that incorrect top layer factor is added to the pose graph, the dense bottom layer factors ensure this incorrect factor will not drag the poses away from the correct direction.

We further validate the versatility of our proposed work on our self-collected dataset using solid-state LiDAR [26] in both structured and unstructured scenes. The first test scene is a structured indoor factory with several irregular-shaped pipelines and machines (see Fig. 8). The size of this scene is around 14×16×8m, and the length of this sequence is

TABLE VI: MME on Self Collected Dataset

| Method | scene-1 | scene-2 |
|---|---|---|
| FAST-LIO2 [1] | -2.99 | -2.60 |
| **Proposed** | **-3.06** | **-2.69** |

### B. Ablation Study

*1) Pose Graph Optimization versus Direct Assign:* In this section, we demonstrate our proposed top-down design is non-trivial. To update the bottom layer poses, a trivial method is to directly assign the optimized upper layer poses to the lower layer poses, e.g., an optimized pose of a keyframe from the

TABLE VII: RMSE of the ATE ($°/m$) on KITTI Dataset

| Method | Seq. 00* | Seq. 01 | Seq. 02* | Seq. 03 | Seq. 04 | Seq. 05* | Seq. 06* | Seq. 07* | Seq. 08* | Seq. 09* | Seq. 10 | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Direct Assign★ | 0.68/0.81 | **0.87/1.88** | 1.44/5.25 | 0.72/**0.62** | 0.19/**0.81** | **0.51**/0.52 | **0.35**/0.29 | **0.40/0.26** | 1.28/2.76 | 0.86/1.59 | **0.55**/1.13 | 0.71/1.47 |
| **Proposed★** | **0.67/0.79** | **0.87**/1.91 | **1.19/5.08** | **0.71**/0.63 | **0.14**/0.82 | 0.53/**0.39** | 0.36/**0.22** | **0.40**/0.30 | **1.25/2.68** | **0.83/1.26** | 0.57/**1.12** | **0.68/1.38** |
| Direct Assign | 1.30/1.32 | **0.97/2.37** | 2.18/9.47 | 0.67/**0.58** | 0.28/**0.87** | 0.66/0.83 | 0.35/0.31 | 0.47/0.32 | 1.71/3.42 | 1.02/2.12 | 0.54/1.08 | 0.92/2.06 |
| **Proposed** | **1.00/1.17** | 0.98/2.41 | 2.26/8.95 | **0.65/0.59** | **0.21**/0.90 | **0.48/0.74** | **0.33/0.22** | **0.43/0.29** | **1.29/2.53** | **0.70/1.47** | **0.53/1.07** | **0.81/1.85** |

\* Sequence with loops.
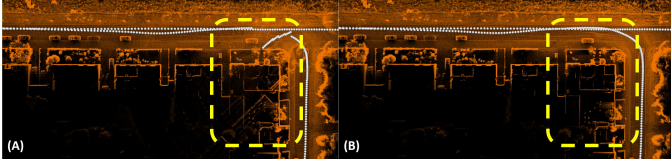★ The initial poses trajectory is generated with loop closure.



Fig. 10: Mapping results (A) without and (B) with our proposed pose graph optimization in KITTI Seq. 08. The main difference is emphasized by the yellow dashed rectangle. The white dots represent the trajectory. Inconsistency occurs at the higher layer and is iteratively assigned to the bottom layer poses if no pose graph optimization is applied.

upper layer is used to update the first $s$ poses from the lower layer within the corresponding local window. For example, in Fig. 1, the poses $\mathbf{x}_0^2, \mathbf{x}_1^2, \mathbf{x}_2^2$ are updated by $\mathbf{x}_0^{3*}$ and the poses $\mathbf{x}_3^2, \mathbf{x}_4^2, \mathbf{x}_5^2$ are updated by $\mathbf{x}_1^{3*}$, respectively. We test this trivial method ("Direct Assign") with our proposed pose graph optimization mechanism on KITTI dataset [17] both using the poses generated before and after loop closure from MULLS [3] as the input. The RMSE of the absolute rotation and translation errors are summarized in Table VII. As can be seen, our proposed pose graph optimization outperforms the direct assigning one both in pose estimation accuracy and mapping consistency (see Fig. 10). Though the above trivial strategy could still improve the pose estimation accuracy, such a method neglects the relative pose constraints from each local window, e.g., the pose $\mathbf{x}_3^2$ is involved in two local windows, whereas it is updated by $\mathbf{x}_1^{3*}$ only without considering the relative constraint from $\mathbf{x}_2^2$. This will lead to a mapping inconsistency between frames $\mathbb{F}_2^2$ and $\mathbb{F}_3^2$ and further between frames $\mathbb{F}_8^1$ and $\mathbb{F}_9^1$. In our proposed pose graph optimization approach, the first layer factor ensures the consistency between every adjacent frame while the second and above layer factors ensure the gap is converged towards the correct direction.

*2) Hierarchical BA versus Reduced BA:* In this section, we demonstrate that our proposed bottom-up design is non-trivial. To accelerate the Hessian matrix solving process, a trivial way is to keep only the block diagonal elements (of size $s$ of the stride length) of the original Hessian matrix and solve this reduced matrix without considering the relative pose constraints among different local windows as in our method. We verify this reduced BA with the original BA and our approach on the DCC sequence of the MulRan dataset [18]. The RMSE of the absolute translation error and the total optimization time of all methods are summarized in Table VIII. It is seen our proposed work achieves the optimal precision as the original BA method while drastically reducing the computation time. This is due to the reason that the original BA needs to construct

an adaptive voxel map using all points (adaptive-voxel map), which quickly escalates as the involved pose number increases. Our method conducts BA in local windows, so we only have to construct an adaptive voxel map using a very small amount of points in the local window, and different local windows can be paralleled. The reduced BA actually takes longer time than the original BA due to two reasons. First, it needs to construct an adaptive-voxel map similar to the original BA. Second, the reduced BA zeros the off-diagonal block elements, which leads to inaccurate Hessian estimation and significantly slows down the convergence speed. In our experiments, the reduced BA may even fail to converge when the maximum iteration number is reached (*max_iter=10*) while the original BA converges within a few steps. Moreover, since we use relatively strict parameters on the local BA, factors from these layers ensure that glitches will not appear between every adjacent frame. For the simplified and the original BA, only the strict parameter could be adopted, otherwise, false feature matching will frequently happen (points within the voxel do not form a plane feature).

TABLE VIII: RMSE of ATE ($m$) and Optimization Time on DCC Sequence of MulRan Dataset

| Method | DCC01 | | DCC02 | | DCC03 | |
|---|---|---|---|---|---|---|
| | RMSE | Time | RMSE | Time | RMSE | Time |
| Initial | 5.67m | - | 3.48m | - | 2.84m | - |
| Original BA | **5.20m** | 2897.54s | **3.20m** | 2853.76s | **2.54m** | 7219.38s |
| Reduced BA | 5.66m | 4972.10s | 3.46m | 4923.90s | 2.83m | 8176.46s |
| **Proposed** | **5.20m** | **206.04s** | 3.22m | **328.59s** | **2.54m** | **256.27s** |

*C. Computation Cost*

In this section, we demonstrate that our proposed approach is computationally efficient, especially on the large-scale dataset. We test our proposed method with multiple setups of used layers on New College [20] and Newer College [19] datasets whose data length varies from $10^3$ to $10^4$ frames. The total computation time (including the adaptive-voxel map construction and BA time) and the maximum RAM memory consumption recorded for each setup at all sequences are illustrated in Fig. 11. Due to the huge time and RAM consumption of the original BA method, we do not test it on the last two sequences since the plot could already depict the trend.

It is seen for all test scenes the more layers are used, the less computation time the pyramid takes. When the pose number $N < 5 \times 10^3$, the time and RAM consumption of 3-layer and 4-layer pyramids are similar, and when $N > 5 \times 10^3$, the 4-layer pyramid becomes optimal. All these phenomenons
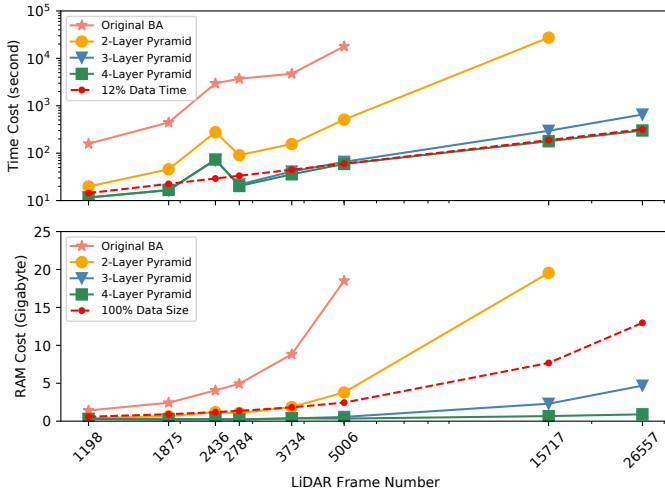
Fig. 11: Comparison of computation time and RAM costs with multiple setups of layers used in the pyramid under Newer College Dataset [19, 20]. Different setups are represented using different colors. The red dashed lines represent 12% of the total data time and total point cloud size of the corresponding sequence.

are in accordance with our theoretical analysis shown in Fig. 3 and (3). Since the third test scene ($N$=2436) contains a more complex environment (thus more adaptive-voxel map construction time), the total time consumption is actually larger than the latter scene. Despite this, by choosing the optimal layer setup, our work could converge within around 12% of the whole data time and consumes much smaller RAM during operation, which is suitable for practical usage.

## V. CONCLUSION

In this paper, we propose a hierarchical BA and pose graph optimization-based work to optimize the pose estimation accuracy and mapping consistency globally for the large-scale LiDAR point cloud. With the bottom-up hierarchical BA, we parallelly solve multiple BA problems with a much smaller Hessian matrix size than the original BA method. With the top-down pose graph optimization, we smoothly and efficiently update the LiDAR poses without generating glitches. We validate the effectiveness of our work on spatially and timely large-scale LiDAR datasets with structured and unstructured scenes, given a good initial pose trajectory or with large drifts. We demonstrate our proposed work outperforms other SOTA methods in pose estimation accuracy and mapping consistency on multiple public spinning LiDAR and our self-collected solid-state LiDAR datasets. In our future work, we could combine the IMU pre-integration and LiDAR measurement noise model into our hierarchical BA work.

## REFERENCES

[1] W. Xu, Y. Cai, D. He, J. Lin, and F. Zhang. Fast-lio2: Fast direct lidar-inertial odometry. *IEEE Transactions on Robotics*, 38(4):2053–2073, 2022.

[2] T. Shan, B. Englot, D. Meyers, W. Wang, C. Ratti, and Rus D. Lio-sam: Tightly-coupled lidar inertial odometry via smoothing and mapping. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5135–5142. IEEE, 2020.

[3] Y. Pan, P. Xiao, Y. He, Z. Shao, and Z. Li. Mulls: Versatile lidar slam via multi-metric linear least square. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11633–11640, 2021.

[4] K. Koide, M. Yokozuka, S. Oishi, and A. Banno. Globally consistent and tightly coupled 3d lidar inertial mapping. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 5622–5628, 2022.

[5] Z. Liu and F. Zhang. Balm: Bundle adjustment for lidar mapping. *IEEE Robotics and Automation Letters*, 6(2):3184–3191, 2021.

[6] L. Zhou, D. Koppel, and M. Kaess. Lidar slam with plane adjustment for indoor environment. *IEEE Robotics and Automation Letters*, 6(4):7073–7080, 2021.

[7] H. Zhang, Z. Hou, N. Li, and S. Song. A graph-based hierarchical SLAM framework for large-scale mapping. In *Intelligent Robotics and Applications*, pages 439–448. Springer Berlin Heidelberg, 2012.

[8] Y. Tian, Y. Wang, M. Ouyang, and X. Shi. Hierarchical segment-based optimization for slam. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 6573–6580, 2021.

[9] Z. Liu, X. Liu, and F. Zhang. Efficient and consistent bundle adjustment on lidar point clouds, 2022.

[10] P. J. Besl and N. D. McKay. A method for registration of 3-d shapes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 14(2):239–256, 1992.

[11] A. Segal, D. Haehnel, and S. Thrun. Generalized-ICP. In *Robotics: Science and Systems V*. Robotics: Science and Systems Foundation, June 2009.

[12] H. Yang, J. Shi, and L. Carlone. Teaser: Fast and certifiable point cloud registration. *IEEE Transactions on Robotics*, 37(2):314–333, 2021.

[13] W. Hess, D. Kohler, H. Rapp, and D. Andor. Real-time loop closure in 2d lidar slam. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1271–1278, 2016.

[14] J. Behley and C. Stachniss. Efficient surfel-based SLAM using 3d laser range data in urban environments. In *Robotics: Science and Systems XIV*. Robotics: Science and Systems Foundation, June 2018.

[15] S. Liang, Z. Cao, C. Wang, and J. Yu. A novel 3d lidar slam based on directed geometry point and sparse frame. *IEEE Robotics and Automation Letters*, 6(2):374–381, 2021.

[16] K. Ćwian, M. R. Nowicki, J. Wietrzykowski, and P. Skrzypczyński. Large-scale lidar slam with factor graph optimization on high-level geometric features. *Sensors*, 21(10):3445, May 2021.

[17] 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*, pages 3354–3361, 2012.

[18] G. Kim, Y. S. Park, Y. Cho, J. Jeong, and A. Kim. Mulran: Multimodal range dataset for urban place recognition. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6246–6253, 2020.

[19] L. Zhang, M. Camurri, D. Wisth, and M. Fallon. Multi-camera lidar inertial extension to the newer college dataset, 2021.

[20] M. Ramezani, Y. Wang, M. Camurri, D. Wisth, M. Mattamala, and M. Fallon. The newer college dataset: Handheld lidar, inertial and vision with ground truth. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4353–4360, 2020.

[21] M. Yokozuka, K. Koide, S. Oishi, and A. Banno. Litamin2: Ultra light lidar-based slam using geometric approximation applied with kl-divergence. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11619–11625, 2021.

[22] J. Behley and C. Stachniss. Efficient surfel-based slam using 3d laser range data in urban environments. In *Proc. of Robotics: Science and Systems (RSS)*, 2018.

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

[24] P. Dellenbach, J. 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)*, pages 5580–5586, 2022.

[25] C. Yuan, W. Xu, X. Liu, X. Hong, and F. Zhang. Efficient and probabilistic adaptive voxel mapping for accurate online lidar odometry. *IEEE Robotics and Automation Letters*, 7(3):8518–8525, 2022.

[26] Z. Liu, F. Zhang, and X. Hong. Low-cost retina-like robotic lidars based on incommensurable scanning. *IEEE/ASME Transactions on Mechatronics*, 27(1):58–68, 2022.

[27] J. Razlaw, D. Droeschel, D. Holz, and S. Behnke. Evaluation of registration methods for sparse 3d laser scans. In *2015 European Conference on Mobile Robots (ECMR)*, pages 1–7, 2015.