

“© 2021 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.”

Globally optimal consensus maximization for robust visual inertial localization in point and line map

Yanmei Jiao¹, Yue Wang¹, Bo Fu¹, Qimeng Tan², Lei Chen², Minhong Wang³,
Shoudong Huang⁴ and Rong Xiong¹

Abstract—Map based visual inertial localization is a crucial step to reduce the drift in state estimation of mobile robots. The underlying problem for localization is to estimate the pose from a set of 3D-2D feature correspondences, of which the main challenge is the presence of outliers, especially in changing environment. In this paper, we propose a robust solution based on efficient global optimization of the consensus maximization problem, which is insensitive to high percentage of outliers. We first introduce *translation invariant measurements* (TIMs) for both points and lines to decouple the consensus maximization problem into rotation and translation subproblems, allowing for a two-stage solver with reduced search space. Then we show that (i) the rotation can be estimated by minimizing TIMs using only *1-dimensional branch-and-bound* (BnB), (ii) the translation can be estimated by running *1-dimensional search* for each of the three axes with *prioritized progressive voting*. Compared with the popular randomized solver, our solver achieves deterministic global convergence without requiring an initial value. Furthermore, ours is exponentially faster compared with existing BnB based methods. Finally, our experiments on both simulation and real-world datasets demonstrate that the proposed method gives accurate pose estimation even in the presence of 90% outliers (only 2 inliers).

I. INTRODUCTION

Visual inertial navigation system is popular for state estimation of mobile robots, autonomous vehicles and augmented reality applications. Many efforts have been paid to build accurate, consistent and efficient visual inertial odometry [1]. However, its inherent drift is unacceptable in long-term operation, requiring absolute pose estimation for correction. Map based visual inertial localization is therefore an important component in a complete navigation system, of which the goal is to estimate the absolute pose from a set of corresponding 2D image feature points and global 3D map points. In this problem, one main challenge is the robustness of the solver against outliers, i.e. incorrect feature correspondences. When high percentage of correspondences is outlier, the performance of the general pose estimator may be severely degenerated.

¹Yanmei Jiao, Yue Wang, Bo Fu, Rong Xiong are with the State Key Laboratory of Industrial Control and Technology, Zhejiang University, Hangzhou, P.R. China. ²Qimeng Tan and Lei Chen are with the Beijing Key Laboratory of Intelligent Space Robotic System Technology and Applications, Beijing Institute of Spacecraft System Engineering, Beijing, P.R. China. ³Minhang Wang is with the Application Lab, Huawei Incorporated Company, P.R. China. ⁴Shoudong Huang is with the Center for Autonomous Systems (CAS), the University of Technology, Sydney, Australia. Yue Wang is the corresponding author wangyue@iipc.zju.edu.cn. This work was supported in part by the National Key R&D Program of China (2018YFB1309300), and in part by the National Nature Science Foundation of China (61903332).



Fig. 1: The projected map points on the map image (left column) and the detected image key points on the query image (right column), with inlier correspondences in red and outliers in blue. The initial correspondences found by feature descriptor matching (top), and the consensus set correspondences searched by RANSAC (middle) and proposed consensus maximization algorithm (bottom).

Pose estimation with outliers is often stated as a consensus maximization problem. One popular solution is random sample consensus (RANSAC), which has lots of variants [2] [3] and has been employed in many visual localization methods [4] [5]. The advantage of RANSAC is the simplicity for implementation, and the effectiveness in many scenarios with moderate percentage of outliers. However, RANSAC cannot tolerate extreme percentage of outliers, say 90%. In addition, it cannot guarantee the deterministic global optimality due to the probabilistic convergence.

In contrast to RANSAC, another solution to consensus maximization is global optimization based methods. It gives globally optimal solution without relying on an initial value [6] [7], while it cannot perform in real-time due to the considerable computation time. Most existing global optimization methods aim at general pose estimation problems. They employ branch-and-bound (BnB) as the basic framework to reduce the search space [8], or mixed integer programming for further acceleration [9]. But the computational cost is still unsatisfactory as the pose space $SE(3)$ is coupled. Even inertial measurement is provided, it cannot be easily substituted into the problem for decoupling.

In this paper, we propose a deterministic visual inertial localization solution to achieve global convergence with much higher efficiency. The key idea is to divide $SE(3)$ search space into multiple 1-D search spaces. Specifically, inspired by the decoupling idea in [10], we build intermediate cost function for both point and line features, *translation invariant*

measurements (TIMs), to decouple consensus maximization into two cascaded subproblems only related to rotation $SO(3)$ and translation \mathbb{R}^3 respectively. Based on TIMs, the globally optimal rotation is then searched by *1-dimensional BnB* in $[-\pi, \pi]$ with the aid of inertial measurements. For the translation, \mathbb{R}^3 search is replaced with three 1-dimensional \mathbb{R} search for each axis using *prioritized progressive voting*. To the best of our knowledge, this is the first solver for visual inertial localization with deterministic global optimality. In summary, our contributions include

- TIMs based formulation of visual inertial localization that decouples the problem and enables 1D BnB based global optimization of the rotation.
- Prioritized progressive voting method that replaces \mathbb{R}^3 space search with three \mathbb{R} search for global optimization of the translation.
- Experiments on simulation and real-world cross-session datasets that validate the effectiveness and efficiency of the proposed method against comparative methods.

The remainder of the paper is organized as follows: Section II reviews the related literatures. Section III presents the decoupling of the consensus maximization problem. Section IV introduces the solutions of the subproblems. Section V demonstrates the experimental settings and results, followed by Section VI concluding the paper.

II. RELATED WORKS

A. Random sample consensus

Visual localization and navigation for mobile robots has been studied extensively in the robotics and computer vision communities in the recent decade [11] [12]. For robust localization given the feature correspondences containing outliers, RANSAC is the most popular solution employed in many visual navigation system. In [13] [14], point feature correspondences based RANSAC are studied. In [15], RANSAC is extended to line features. When inertial measurements are provided, the DoF of the problem is reduced, which is utilized by RANSAC to improve the robustness in [16], and extended to both point and line correspondences in [5]. As RANSAC is developed on randomized sampling theory, it is simple to implement and has good performance on scenarios with moderate outliers. But its disadvantage is also obvious, including low tolerance against extreme outliers, local convergence and no guarantee of the optimality [17].

B. Outlier resistant estimator

Another branch to reject outliers is to refer other forms of cost functions instead of the squared error [18]. In [19], Geman-McClure cost function is utilized for 3D-3D registration, which is insensitive to outliers. In [20], M-estimators in several typical robotics problems are presented. Switchable cost function is employed to solve pose graph optimization with outlier loop closures [21]. A more compact solver for such cost function is dynamic covariance scaling which is introduced in [22]. More recently, in [23], several forms of robust cost functions are unified and solved using graduated non-convexity without an initial guess, which

demonstrates good performance in 3D-3D registration, pose graph optimization, and is extended to non-minimal solver for shape reconstruction from an image in [24]. Alternatively, in [25], the outlier rejection is solved by adaptively removing the measurements with large errors, which is simple but show superior performance than RANSAC. These methods achieve deterministic convergence, while some of them offer certifiable optimality (or sub-optimality guarantees).

C. Global optimization method

Global optimization methods are proposed to achieve the global optimality and deterministic convergence. In this branch of literatures, Branch-and-Bound (BnB) is mostly used, which gradually prunes the solution space by coarse-to-fine division. In [6], BnB is used to solve the 2D-2D registration problems. In [8], a general framework for point, line and plane features is proposed to solve 3D-3D registration via BnB. Integrated with mixed integer programming, the BnB optimization can converge faster [9]. In [17], the linear matrix inequality constraints are introduced to mixed integer programming, resulting in a general-purpose faster BnB for all 2D-2D, 2D-3D and 3D-3D geometric vision problems. In the works mentioned above, the rotation is modeled as a rotation matrix with matrix level constraints. Thus it is unclear about the incorporation of inertial measurements. In addition, there are also specialized globally optimal algorithms focusing on one class of problem. In [26] [27], pairs of features are used to decouple the 3D-3D registration. In [10], TEASER is proposed to decoupled scaled 3D-3D registration, achieving a fast three-stage optimization. These works show that it is possible to have superior performance with *specialized* algorithms rather than only the *general-purpose* framework, even also accelerated.

In this paper, we follow the idea of specialized solver to bridge the gap of globally optimal deterministic solution for visual inertial localization, which is a robust 3D-2D pose estimation problem with inertial measurements. To the best of our knowledge, this is the first work to study this problem in the context of global optimality. And the solution is accurate and efficient which is demonstrated in later experiments.

III. DECOUPLING TRANSLATION AND ROTATION

The underlying problem of visual inertial localization is the pose estimation from 3D-2D correspondences with outliers. Formally, given a set \mathfrak{P} consisting of correspondences between 3D global points $p_i \in \mathbb{R}^3$ and 2D visual points $u_i \in \mathbb{R}^2$, they satisfy

$$u_i = \pi(Rp_i + t, K) + o_i + e_i \quad (1)$$

where $R \in SO(3)$ and $t \in \mathbb{R}^3$ is the camera pose to be estimated, π is the camera projection function with known intrinsic parameters K , $|e_i| < n_i$ is assumed to be bounded random measurement noise, o_i is zero for inlier while an arbitrary number for outlier. To deal with outliers, the robust pose estimation generally begins with consensus maximization problem as

$$\begin{aligned} & \max_{R,t,\{z_i\}} \sum z_i & (2) \\ \text{s.t. } & z_i |u_i - \pi(Rp_i + t, K)| \leq n_i, \quad i \in \mathfrak{P} & (3) \end{aligned}$$

where z_i is binary, indicating whether o_i is zero. To solve the problem in global, general BnB algorithms search in $SE(3)$, which is a coupled space of $SO(3)$ and \mathbb{R}^3 . But this probably leads to exponential computational complexity in bad cases. For local techniques like RANSAC, inliers may be estimated conservatively, i.e. inliers regarded as outliers, especially when the noise is unavoidable.

A. Translation invariant measurements

1) *Point-TIM*: Inspired by the minimal solution in RANSAC, we develop an intermediate measurement which is invariant to the translation of the pose. Mathematically, given an image key point u_i , we have an un-normalized direction vector from the camera center as

$$\tilde{u}_i \triangleq \begin{pmatrix} \tilde{u}_{i,x} \\ \tilde{u}_{i,y} \\ 1 \end{pmatrix} = K^{-1} \begin{pmatrix} u_i \\ 1 \end{pmatrix} \quad (4)$$

Then the corresponding world point p_i is transformed to the camera coordinates and satisfies

$$\frac{R_1 p_i + t_x}{\tilde{u}_{i,x}} = \frac{R_2 p_i + t_y}{\tilde{u}_{i,y}} = R_3 p_i + t_z \quad (5)$$

where $R \triangleq (R_1^T, R_2^T, R_3^T)^T$ and $t \triangleq (t_x, t_y, t_z)^T$. Based on (5), we have two constraints from a correspondence. Naturally, given another correspondence u_j and p_j , we can have two more constraints as

$$\frac{R_1 p_j + t_x}{\tilde{u}_{j,x}} = \frac{R_2 p_j + t_y}{\tilde{u}_{j,y}} = R_3 p_j + t_z \quad (6)$$

According to (5) and (6), we have linear constraints of the translation t . With proper variable substitutions among the constraints, and the globally observable pitch and roll angles from inertial measurements, we can eliminate t , reduce $SO(3)$ to $[-\pi, \pi]$, and derive TIM as

$$d_p(\alpha) = d_{p,1} \sin \alpha + d_{p,2} \cos \alpha + d_{p,3} \quad (7)$$

where α is the unknown yaw angle, $d_{p,1}$, $d_{p,2}$, $d_{p,3}$ and the derivation details are presented in the Appendix I-A¹. Now we substitute the constraints which are related to both R and t in (2) with the TIM, leading to

$$\max_{R(\alpha), \{z_{ij}\}} \sum z_{ij} \quad (8)$$

$$\text{s.t. } z_{ij} |d_{p,ij}(\alpha)| \leq n_{ij}, \quad i, j \in \mathfrak{P} \quad (9)$$

where $n_{ij} = \min(n_i, n_j)$, $z_{ij} = 1$ indicates the i -th and j -th correspondence derived the constraint are inliers.

2) *Line-TIM*: Similar to a pair of point correspondences, given a set of line correspondences \mathfrak{L} , it is also possible to develop TIM. Given the end points of the image line segment u_{k1} and u_{k2} , we have two un-normalized directions as (4), denoted as \tilde{u}_{k1} and \tilde{u}_{k2} .

Then following the fact that the point p_k on the world line lies on the plane spanned by the rays from camera center

along direction \tilde{u}_{k1} and \tilde{u}_{k2} , we have

$$(\tilde{u}_{k1} \times \tilde{u}_{k2})^T (Rp_k + t) = 0 \quad (10)$$

which is a constraint for both rotation and translation. Since arbitrary number of points can be sampled from a line, we sample another point on the same world line to formulate the constraint as (10). Then only one line correspondence can lead to line-TIM after proper substitution as

$$d_l(\alpha) = d_{l,1} \sin \alpha + d_{l,2} \cos \alpha + d_{l,3} \quad (11)$$

where the line-TIM has the same form as point-TIM in (7), but the coefficients are different. The derivation details are presented in the Appendix I-B¹.

TIMs based rotation only problem. Note that either (7) or (11) is only related to the yaw angle. By combining them together, we have a general consensus maximization problem with TIM constraints only related to rotation compatible to the map having both point and line features as

$$\max_{R(\alpha), \{z_*\}} \sum z_* \quad (12)$$

$$\text{s.t. } z_{ij} |d_{p,ij}(\alpha)| \leq n_{ij}, \quad i, j \in \mathfrak{P} \quad (13)$$

$$z_k |d_{l,k}(\alpha)| \leq n_k, \quad k \in \mathfrak{L} \quad (14)$$

B. Two-stage consensus maximization solver

With TIMs for both point and line correspondences, we decouple the original consensus maximization problem into rotation only problem, and translation only problem when the rotation is fixed. Accordingly, the proposed solver has two stages in cascade:

- We estimate the rotation \hat{R} by $R(\hat{\alpha})$ based on the TIMs in (12). This estimator solves a 1D optimization problem and is described in Section IV-A.
- We estimate the translation \hat{t} based on the original consensus maximization in (2) where the rotation is assigned with \hat{R} . This estimator solves a \mathbb{R}^3 optimization problem and is described in Section IV-B.

IV. ESTIMATORS OF ROTATION AND TRANSLATION

A. BnB based optimization for rotation

We employ BnB strategy to solve problem (12). The cost function in (12) relates to α and z_* . But it is obvious that when α is determined, $\{z_*\}$ is simply derived by evaluating the constraints. So we denote the cost function as $E(\alpha)$ that is explained as the number of inliers given a yaw angle α .

Upper bound of cost function. We then derive the upper bound of $E(\alpha)$ on the subset \mathbb{A} , denoted as $\bar{E}(\mathbb{A})$, where $\alpha \in \mathbb{A} \subseteq [-\pi, \pi]$. Recall (7) and (11), as the forms of point-TIM and line-TIM are the same, we denote them as $d(\alpha)$. The lower bound of $|d(\alpha)|$ on \mathbb{A} , denoted as $\underline{d}(\mathbb{A})$, is

$$\underline{d}(\mathbb{A}) = \min |a_1 \sin(\alpha + a_2) + d_3| \quad (15)$$

where the derivation of the coefficients are in Appendix I-C¹. Note that $\underline{d}(\mathbb{A})$ can be solved analytically without iterations. Now we formulate a consensus maximization problem as

$$\max_{R(\alpha), \{z_*\}, \alpha \in \mathbb{A}} \sum z_* \quad (16)$$

$$\text{s.t. } z_{ij} \underline{d}_{p,ij}(\mathbb{A}) \leq n_{ij}, \quad i, j \in \mathfrak{P} \quad (17)$$

$$z_k \underline{d}_{l,k}(\mathbb{A}) \leq n_k, \quad k \in \mathfrak{L} \quad (18)$$

¹The appendix is available on <https://arxiv.org/abs/2002.11905>

where the problem is defined on \mathbb{A} , and the TIMs constraints are replaced with tight lower bounds, relaxing the constraints and yielding an optimistic estimation of \hat{z}_* . We then have

$$E(\alpha) \leq \bar{E}(\mathbb{A}) = \sum \hat{z}_*, \quad \alpha \in \mathbb{A} \quad (19)$$

as a tight upper bound. The equality exists when all constraints give the same α with $d_{p,ij}(\alpha) = \underline{d}_{p,ij}(\mathbb{A})$ and $d_{l,k}(\alpha) = \underline{d}_{l,k}(\mathbb{A})$, which is only possible in noise-free condition.

Accelerate BnB optimization. With (12-19), we have the BnB search for globally optimal rotation, of which the pseudo code is listed in Algorithm 1. Note that the main idea of BnB is to prune the solution space \mathbb{A} when its upper bound $\bar{E}(\mathbb{A})$ is smaller than the current best estimates E^* . Therefore, if we have a fast solution to initialize a good E^* , most solution spaces can be pruned at early stage, significantly improving the search efficiency. To implement this idea, we use RANSAC [5] to generate a rough initial E^* . In addition, we introduce a heuristics to balance the global optimality and the efficiency. The best M estimated α during RANSAC is utilized to initialize M subsets among $[-\pi, \pi]$. Each subset centers at each estimated α with a width w . When w is large, global optimality is emphasized and vice versa. Another implementation trick is to store the respective inliers when evaluating (16) on each subset \mathbb{A} . When \mathbb{A} is further divided into smaller subsets, only the stored inliers within \mathbb{A} are evaluated, instead of all constraints, saving lots of computational cost. These techniques are all shown to accelerate the search in the experimental ablation study without drop of accuracy.

Algorithm 1: Globally Optimal Rotation Search

Input: 3D-2D feature correspondences $\mathfrak{P}, \mathfrak{L}$

Output: Optimal α^*

- 1 Initialize partition of $[-\pi, \pi]$ into subsets $\{\mathbb{A}_i\}$.
 - 2 Initialize best estimation E^*, α^* .
 - 3 Insert $\{\mathbb{A}_i\}$ into queue q .
 - 4 **while** q is not empty **do**
 - 5 Pop the first subset of q as \mathbb{A} .
 - 6 Compute $\bar{E}(\mathbb{A})$ as (16).
 - 7 **if** $\bar{E}(\mathbb{A}) > E^*$ **then**
 - 8 Assign center of \mathbb{A} as α_c .
 - 9 Compute $E(\alpha_c)$ as (12).
 - 10 **if** $E(\alpha_c) > E^*$ **then**
 - 11 Update $E^* \leftarrow E(\alpha_c), \alpha^* \leftarrow \alpha_c$.
 - 12 Subdivide \mathbb{A} into subsets and insert into q .
-

B. Prioritized progressive voting for translation

As the translation has 3DoF, it will take exponentially higher time if we still apply BnB-based method. Inspired by the polynomial-time algorithm *adaptive voting* in [10], we present the novel prioritized progressive voting for translation in the following.

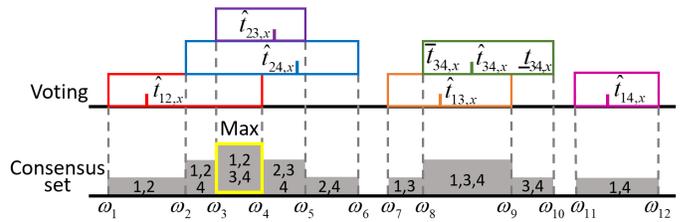


Fig. 2: The voting illustration of \hat{t}_x . Each $\hat{t}_{ij,x}$ derived by i -th and j -th correspondence votes for the interval if $[\omega_i, \omega_{i+1}] \subseteq [\underline{t}_{ij,x}, \bar{t}_{ij,x}]$, which means the corresponding consensus set contains i and j .

When $R(\hat{\alpha})$ is estimated, the co-linear and co-planar constraints (5) and (10) are all linear constraints for t . Thus we can transform the consensus maximization problem with point and line constraints as

$$\max_{t, \{z_i\}} \sum z_i \quad (20)$$

$$s.t. \quad z_i |A_i t + b_i| \leq n_i, \quad i \in \mathfrak{P} \cup \mathfrak{L} \quad (21)$$

where $A_i \in \mathbb{R}^{1 \times 3}$ and $b_i \in \mathbb{R}$ are the coefficients for linear constraints derived from (5) or (10) with estimated $R(\hat{\alpha})$. However, this problem still has coupled constraints for t so that \mathbb{R}^3 search is indispensable.

Decoupled linear constraints. Note that for a point correspondence constraint (5), we have two linear equations, while for a line correspondence constraint (10), we have one. Therefore, given a pair of correspondences including at least one point correspondence, say the i -th point correspondence and the j -th point or line correspondence, it is sufficient to solve \hat{t}_{ij} for this small linear system (see Appendix II¹ for details), then we have

$$\max_{t, \{z_{ij}\}} \sum z_{ij} \quad (22)$$

$$s.t. \quad z_{ij} |\hat{t}_{ij} - t| \leq n_{ij}, \quad i \in \mathfrak{P}, j \in \mathfrak{P} \cup \mathfrak{L} \quad (23)$$

Now we find that the constraints are decoupled for each dimension of t . Set the x -dimension as example, we have

$$\max_{t_x, \{z_{ij}\}} \sum z_{ij} \quad (24)$$

$$s.t. \quad z_{ij} |\hat{t}_{ij,x} - t_x| \leq n_{ij,x}, \quad i \in \mathfrak{P}, j \in \mathfrak{P} \cup \mathfrak{L} \quad (25)$$

arriving at the resultant three dimension-wise linear constrained consensus maximization problems.

Dimension-wise voting algorithm. We use a voting algorithm to solve the problem. We first specify the noise bound $n_{ij,x}$ in (24). Given the noise bound n_i in (20), we have the noise bound for t following the techniques in [28] as

$$\underline{t}_{ij} \leq \hat{t}_{ij} \leq \bar{t}_{ij} \quad (26)$$

The details can be found in Appendix II¹.

Still taking x -dimension as example, each estimated $\hat{t}_{ij,x}$ defines an interval $[\underline{t}_{ij,x}, \bar{t}_{ij,x}]$. If the real t_x lies in this interval, then the real inlier set contains the two correspondences deriving $\hat{t}_{ij,x}$. According to [10], the insight is that the inlier set only changes its membership when real t_x enters a new interval. Besides, given K estimations, the *maximum number of possible consensus sets*, i.e. the cardinality of the solution space, is $2K-1$, where K is in quadratic w.r.t the number of

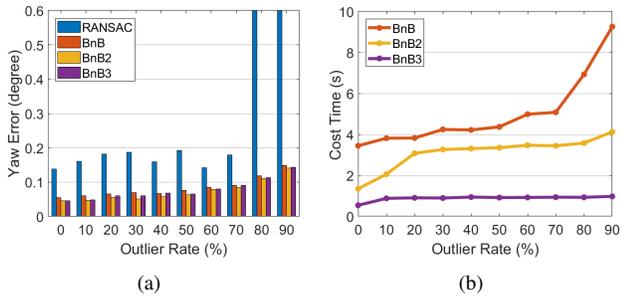


Fig. 3: The rotation accuracy and computation time over the increasing outlier rate. *BnB2* denotes the BnB with RANSAC initialization. *BnB3* denotes the BnB with both RANSAC initialization and the implementation trick.

correspondences. This complexity enables a voting algorithm for all $2K - 1$ sets. By counting the unique correspondences of the votes in each set, we get the corresponding consensus set. Then the maximal consensus set can lead to an estimation of \hat{t}_x . An illustrative case is shown in Fig. 2 and the pseudo code is listed in Algorithm 2 with x -dimension as example. For simplicity, we replace $\hat{t}_{ij,x}$ with $\hat{t}_{k,x}$ in the pseudo code. Following the similar idea in [10], by repeating the voting algorithm for three times, \hat{t} is estimated as $[\hat{t}_x, \hat{t}_y, \hat{t}_z]^T$.

Algorithm 2: Voting

Input: $\{\hat{t}_{k,x}\}, \{\underline{t}_{k,x}\}, \{\bar{t}_{k,x}\}, k = 1..K$
Output: Consensus sets S

- 1 Initialize key-value map S .
- 2 $\omega = \text{sort}([\underline{t}_{1,x}, \bar{t}_{1,x}, \underline{t}_{2,x}, \bar{t}_{2,x}, \dots, \underline{t}_{K,x}, \bar{t}_{K,x}])$.
- 3 **for** $i = 1..2K - 1$ **do**
- 4 $S([\omega_i, \omega_{i+1}]) = \emptyset$.
- 5 **for** $k = 1..K$ **do**
- 6 **if** $[\omega_i, \omega_{i+1}] \subseteq [\underline{t}_{k,x}, \bar{t}_{k,x}]$ **then**
- 7 $S([\omega_i, \omega_{i+1}]) = S([\omega_i, \omega_{i+1}]) \cup k$.

Prioritized progressive voting algorithm. When the number of inliers is high, independent voting along three dimensions is possible. But when the number of inliers is low and outlier rate is high, independent dimension-wise voting may lead to failure. The reason is that, though it is almost impossible that there are more outliers than inliers having the similar t , it is possible that there are more outliers than inliers having the similar t_x . In such scenario, search along x -dimension leads to incorrect \hat{t}_x , which cannot be corrected in the successive voting along y or z -dimension.

To deal with such scenario while keeping a low computational complexity, we propose a prioritized progressive voting for translation in Algorithm 3. The main idea is that we progressively vote on the three dimensions, but there is a priority, i.e. number of votes, for early termination. The experimental results show that the computational complexity of prioritized progressive voting is almost similar to the dimension-wise voting. Otherwise, it is also possible to use 3D BnB translation search for better accuracy, but it is slower

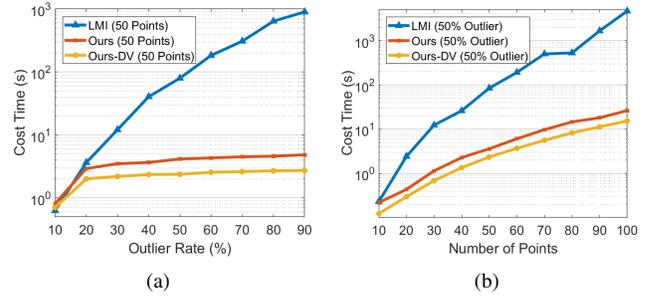


Fig. 4: Computation time comparison over increasing (a) outlier rate (b) number of points. *Ours* denotes the proposed method with prioritized progressive voting, while *Ours-DV* denotes the dimension-wise voting.

because of the coupled multi-dimensional solution space. Finally, we apply nonlinear refinement to achieve the best accuracy when the maximum consensus set is found.

Algorithm 3: Prioritized Progressive Voting

Input: $\{\hat{t}_k\}, \{\underline{t}_k\}, \{\bar{t}_k\}, k = 1..K$
Output: Maximum consensus set \hat{t}

- 1 Initialize best estimation $E^* = 0$.
- 2 $S_x = \text{Voting}(\{\hat{t}_{k,x}\}, \{\underline{t}_{k,x}\}, \{\bar{t}_{k,x}\})$.
- 3 Sort S_x in decreasing cardinality.
- 4 **for each key** $[i]$ **in** S_x **do**
- 5 **if** $|S_x([i])| < E^*$ **then**
- 6 **break;**
- 7 $S_y = \text{Voting}(\{\hat{t}_{k,y}\}, \{\underline{t}_{k,y}\}, \{\bar{t}_{k,y}\}, k \in S_x([i]))$.
- 8 **for each key** $[j]$ **in** S_y **do**
- 9 **if** $|S_y([j])| < E^*$ **then**
- 10 **break;**
- 11 $S_z = \text{Voting}(\{\hat{t}_{k,z}\}, \{\underline{t}_{k,z}\}, \{\bar{t}_{k,z}\}, k \in S_y([j]))$.
- 12 **if** $\max_{S_z([m])} |S_z([m])| > E^*$ **then**
- 13 Update $E^* \leftarrow \max_{S_z([m])} |S_z([m])|$.
- 14 Update $S^* \leftarrow \arg \max_{S_z([m])} |S_z([m])|$.

V. EXPERIMENTAL RESULTS

In the experiments, we evaluate the proposed consensus maximization solver on (i) the feasibility and effectiveness of the subproblem solvers, (ii) the accuracy and robustness compared with existing methods, and (iii) the performance in real world visual inertial localization applications. We implement the proposed solver in MATLAB on a desktop with CPU Intel i7-7700 3.60GHz and 8G RAM.

A. Ablation study

We build the synthetic world consisting of 3D points and lines in the cube $[-1, 1]^3$. The 2D image projections are generated with randomly sampled camera poses in $[-2, 2]^3 \times [-\pi, \pi]^3$, as well as their inlier correspondences. All the projected 2D image points are added with bounded random noise e_i with the bound $n_i = 2$. Each outlier

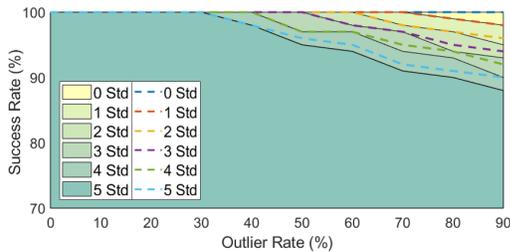


Fig. 5: The sensitivity experiment result using proposed algorithm with dimension-wise voting (solid) and prioritized progressive voting (dash).

correspondence is generated from other randomly sampled camera pose different to ground truth pose. The total number of correspondences is fixed as 50. Specifically, there are 50 point correspondences when evaluating point only methods, while 25 point and 25 line correspondences for the point and line methods. We vary the outlier percentage from 10% to 90% with a step of 10%. Statistic performance indicators are evaluated with an average of 100 Monte Carlo runs. Denoting the ground truth pose as $[R_{gt}|t_{gt}]$, we compute the translation error as $\Delta T = |\hat{t} - t_{gt}|$ in meter and the rotation error as the angle of $\Delta R = \hat{R}R_{gt}^T$ in degree.

BnB heuristics. We first evaluate the heuristics introduced in Section IV-A from the aspect of accuracy and efficiency. As shown in Fig. 3, with the heuristics, the efficiency is improved while the accuracy stays similar. Since the final pose is refined by nonlinear optimization, slight rotation error after BnB can be ignored. As a baseline, we also show the error of estimated rotation giving the most inliers in RANSAC, of which the performance is much worse, indicating inconsistency between the identified inliers and the real inliers. In following experiments, heuristics are applied with BnB as default setting.

Translation voting. We then compare the voting strategies introduced in Section IV-B. Now we can evaluate the final accuracy after nonlinear refinement. In addition to efficiency and accuracy, we also evaluate the consistency between the estimated consensus set and the real inlier set (CCI) using precision and recall. As shown in Fig. 4, the computation of the prioritized progressive voting is slightly higher than the dimension-wise voting. More importantly, the increased time keeps almost consistent w.r.t outlier rate and correspondences number, which might be explained as no complexity growth for prioritized progressive voting. The CCI and accuracy are shown in the right columns in Tab. I. We see that all variants achieve perfect CCI, naturally leading to high accuracy.

Sensitivity to noisy inertial measurements. As inertial measurements are noisy, it's necessary to evaluate the sensitivity. We add Gaussian noise with zero mean and increasing standard deviation up to 5 degree on both pitch and roll angle. The threshold to judge a successful localization is 0.1m for translation error and 0.5 degree for rotation error as in [5]. The result is shown in Fig. 5, indicating the algorithm can achieve over 90% success rate when the noise increases to 5 degree. This level of noise is far more than the pitch and roll estimations in practice [29]. In addition, the performance is

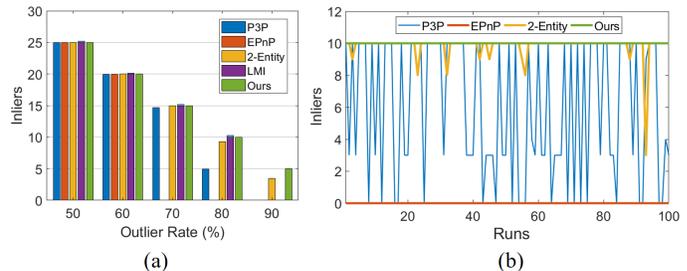


Fig. 6: (a) The number of inliers in the estimated maximal consensus set w.r.t increasing outliers of successful estimation. (b) The number of inliers in the estimated maximal consensus set for 100 runs when the outlier rate is 80%.

better when employing prioritized progressive search.

B. Comparison on synthetic datasets

The comparative methods include the RANSAC-based methods EPnP [14], P3P [13], 2-Entity [5] and globally optimal method LMI [17]. We use the OpenCV [30] implementation of EPnP and P3P. For LMI, we modify their open source code in MATLAB following the paper, since only code for 3D-3D registration is released. The \mathcal{M} and ϵ are set to 10^4 and 10^{-2} , respectively. In addition, we control the evaluation data having rotation angle less than 60° and add it as the constraint of LMI, as suggested in [17]. The 2-Entity RANSAC is implemented in MATLAB and we select the mixed sampling strategy which utilize both points and lines for pose estimation. For all RANSAC-based method, 100 iterations are performed and the threshold of reprojection error for counting inliers is 8 pixel [5]. All methods are followed by nonlinear refinement on the identified consensus set. We still use the synthetic dataset as in the ablation study.

Efficiency of globally optimal methods. We first compare the efficiency between the proposed method and the LMI. We evaluate the computational cost with respect to the number of feature correspondences and the percentage of outliers. The result is shown in Fig. 4, the computational cost of LMI is significantly higher than the proposed methods both for increasing number of correspondences, and the percentage of outliers. The growing gap may also indicate that the complexity of LMI is higher than ours.

Deterministic convergence. The vital difference between RANSAC and globally optimal method is the convergence. We compare the number of inliers in the estimated maximal consensus set with respect to increasing outliers when the final pose estimation is successful. The result is shown in Fig. 6, which indicates that the proposed solution achieves deterministic perfect CCI, while RANSAC gives conservative estimations with less inliers and LMI finds optimistic estimations by incorrectly regarding outliers as inliers. In addition, both RANSAC and LMI fail when the outlier rate is 90%. The results for all 100 runs when the outlier rate is 80% are also shown in Fig. 6. We can see that the proposed algorithm deterministically finds the globally optimal consensus, while RANSAC achieves global optimality probabilistically.

Note that the BnB-based methods are usually be used at the worst cases where RANSAC-based methods cannot

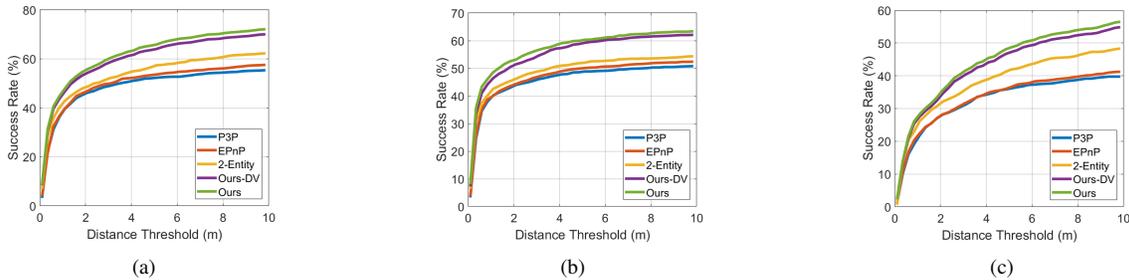


Fig. 7: Success rate with respect to threshold on the whole three sessions 0827 (left), 0828 (center) and 0129 (right).

TABLE I: Accuracy and CCI comparison.

Outlier	Method	P3P	EPnP	2-Entity	LMI	Ours-DV	Ours
60%	ΔT (m)	0.0010	0.0009	0.0008	0.0128	0.0005	0.0006
	ΔR ($^\circ$)	0.0196	0.0170	0.0059	0.0083	0.0019	0.0020
	Precision	1.00	1.00	1.00	0.96	1.00	1.00
	Recall	0.99	0.99	1.00	0.98	1.00	1.00
	Success%	100	100	100	65	100	100
70%	ΔT (m)	0.0013	-	0.0011	0.0209	0.0005	0.0006
	ΔR ($^\circ$)	0.0213	-	0.0211	0.1059	0.0017	0.0028
	Precision	1.00	0	1.00	0.93	1.00	1.00
	Recall	0.98	0	0.99	0.93	1.00	1.00
	Success%	100	0	100	54	100	100
80%	ΔT (m)	0.0017	-	0.0017	0.0246	0.0007	0.0006
	ΔR ($^\circ$)	0.0267	-	0.0257	0.4778	0.0050	0.0032
	Precision	1.00	0	1.00	0.46	1.00	1.00
	Recall	0.49	0	0.93	0.58	1.00	1.00
	Success%	52	0	96	37	100	100
90%	ΔT (m)	-	-	0.0027	-	0.0007	0.0007
	ΔR ($^\circ$)	-	-	0.0411	-	0.0073	0.0043
	Precision	0	0	1.00	0.27	1.00	1.00
	Recall	0	0	0.70	0.35	1.00	1.00
	Success%	0	0	86	0	100	100

¹ The accuracy is evaluated for successful trials, the precision and recall of CCI are for all test trails.

² Ours-DV denotes the proposed method with dimension-wise voting.

handle, due to the high computation time (tens of seconds to minutes). These cases are those with high outlier rate. While in [17], the absolute pose estimation algorithm is only tested under conditions with 50%-60% outlier rate where RANSAC can also perform well. The cases with higher outlier rate (say higher than 80%) are not presented. According to the algorithmic property, the computational complexity of BnB grows exponentially on the outlier rate. Therefore, the high time consuming of LMI is reasonable.

Robustness and accuracy. We finally show the performance of all methods on the synthetic data, including accuracy, precision and recall to measure the CCI, with respect to percentage of outliers ranging from 60% to 90%. Note that we only evaluate the accuracy for successful trials, since result on incorrectly identified consensus set can lead to very large error, disturbing the accuracy. The result in Tab. I first confirms that CCI is highly related to the accuracy, validating the feasibility of maximizing consensus set. RANSAC gives consistent conservative estimations, as the precision remains at a higher level compared with the recall. For LMI, the estimation is prone to regard the outliers as inliers, thus the recall is higher compared with precision. Considering that LMI, P3P and EPnP are designed for general visual localization, the better performance achieved by 2-Entity and the proposed method, designed for visual inertial localization, is reasonable. But we can still summarize that superior result can be found by specialized globally optimal method.

TABLE II: Performance on selected cases in real world.

Method	ExpID	$ \zeta_P /N_P^1$	$ \zeta_L /N_L^1$	ExpID	$ \zeta_P /N_P$	$ \zeta_L /N_L$
	01	9/18	0/0	02	15/39	0/0
	$\Delta T/\Delta R$	Inliers ²	Time	$\Delta T/\Delta R$	Inliers ²	Time
	(m/ $^\circ$)	$ \zeta^* / \zeta $	(s)	(m/ $^\circ$)	$ \zeta^* / \zeta $	(s)
EPnP	0.99/0.80	7/12	0.15	0.90/1.33	11/21	0.15
P3P	0.82/0.63	7/11	0.08	1.98/0.60	10/20	0.09
2-Entity	0.67/0.44	8/10	0.10	0.57/0.34	12/21	0.10
LMI	0.16/0.20	9/13	8.45	0.28/0.22	14/19	344.9
Ours-DV	0.12/0.13	9/9	1.73	0.18/0.16	14/14	22.89
Ours	0.12/0.13	9/9	2.28	0.17/0.13	15/15	53.37
	ExpID	$ \zeta_P /N_P$	$ \zeta_L /N_L$	ExpID	$ \zeta_P /N_P$	$ \zeta_L /N_L$
	03	21/65	0/2	04	23/48	7/15
EPnP	0.45/0.97	10/29	0.13	0.55/0.78	19/28	0.11
P3P	0.32/0.88	13/27	0.11	0.37/0.41	19/27	0.09
2-Entity	0.31/0.46	15/27	0.11	0.14/0.21	27/33	0.12
LMI	0.30/0.38	19/44	1724.78	0.28/0.17	22/28	1163.82
Ours-DV	0.14/0.17	21/23	83.14	0.03/0.16	28/29	67.43
Ours	0.13/0.17	21/23	135.21	0.03/0.15	30/30	116.68
	ExpID	$ \zeta_P /N_P$	$ \zeta_L /N_L$	ExpID	$ \zeta_P /N_P$	$ \zeta_L /N_L$
	05	21/38	8/13	06	96/134	3/4
EPnP	1.09/0.81	13/25	0.14	0.27/0.52	93/112	0.13
P3P	1.09/0.81	13/25	0.13	0.17/0.52	90/98	0.10
2-Entity	0.17/0.27	27/29	0.11	0.12/0.46	95/108	0.10
LMI	0.76/0.64	16/28	712.93	0.09/0.28	96/102	2454.3
Ours-DV	0.16/0.11	29/29	27.06	0.08/0.27	99/99	128.46
Ours	0.16/0.11	29/29	62.69	0.08/0.27	99/99	202.79

¹ N_P denotes the total number of points in the case, $|\zeta_P|$ denotes the number of point inliers, while N_L and $|\zeta_L|$ are numbers for lines.

² $|\zeta|$ denotes the number of identified inliers, while $|\zeta^*|$ the true inliers.

C. Comparison on visual inertial localization

Finally, we evaluate all the methods on a real world cross-session visual inertial localization task. The dataset employed is YQ-dataset [31]. In the dataset, there are three sessions collected in summer 2017, denoted as 2017-0823, 2017-0827 and 2017-0828, and one session in winter 2018 after snow denoted as 2018-0129. The 3D map is built with 2017-0823 session and the other three sessions are used to evaluate the localization performance, indicating the changing environment. The details to obtain the 3D-2D point and line correspondences can be found in Appendix III¹. For evaluation, we compute the ground truth relative pose between the query camera and the map by aligning the synchronized LiDAR scans. For the pitch and roll angle, we use the estimation of visual inertial odometry [32].

Selected cases performance. We first select several typical examples as in [17] and the results are shown in Tab. II. The Exp01, Exp02 and Exp03 are cases with pure point features where Exp03 has lines as disturbance and the outlier rate in these three cases are all more than 50%. One thing to note is that in real world dataset, dimension-wise voting brings slight performance drop, but still achieves superior performance against comparative methods. Also note that in Exp03, the proposed method gives optimistic results by

regarding 2 outliers as inliers, which may be caused by unknown noise bound thus inappropriate threshold in real world data. In Exp04, Exp05 and Exp06, the utilization of good line features promotes the performance of point line methods obviously (2-Entity and ours). Overall, the results still confirm the conclusions in simulation.

Full dataset performance. Finally, we arrive at the success rate on the whole three sessions as shown in Fig. 7. As LMI is too slow to finish all the dataset, here we only show the result of ours and RANSAC methods. We first see that the proposed globally optimal methods consistently outperform the RANSAC methods on all three sessions. The other fact is that progressive prioritized voting brings the best accuracy over the one with dimension-wise voting, because of the consideration on extremely low number of inliers.

VI. CONCLUSIONS

In this paper, we propose a robust solver designed for visual inertial localization, achieving global optimization of the consensus maximization problem, even when the percentage of outliers is very high, say 90%. The key step in our solver is the derivation of *translation invariant measurements* for both points and lines, thus decoupling the problem into two smaller subproblems. Then we propose 1D BnB and prioritized progressive voting to find globally optimal rotation and translation respectively, accelerating the search efficiency. The effectiveness of the proposed method is validated on both synthetic and real world dataset.

REFERENCES

- [1] M. Li and A. I. Mourikis, "High-precision, consistent ekf-based visual-inertial odometry," *The International Journal of Robotics Research*, vol. 32, no. 6, pp. 690–711, 2013.
- [2] 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.
- [3] S. Choi, T. Kim, and W. Yu, "Performance evaluation of ransac family," *Journal of Computer Vision*, vol. 24, no. 3, pp. 271–300, 1997.
- [4] R. Hartley and A. Zisserman, *Multiple view geometry in computer vision*. Cambridge university press, 2003.
- [5] Y. Jiao, Y. Wang, X. Ding, B. Fu, S. Huang, and R. Xiong, "2-entity ransac for robust visual localization: Framework, methods and verifications," *IEEE Transactions on Industrial Electronics*, 2020.
- [6] T. M. Breuel, "Implementation techniques for geometric branch-and-bound matching methods," *Computer Vision and Image Understanding*, vol. 90, no. 3, pp. 258–294, 2003.
- [7] D. Campbell, L. Petersson, L. Kneip, and H. Li, "Globally-optimal inlier set maximisation for simultaneous camera pose and feature correspondence," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1–10, 2017.
- [8] C. Olsson, F. Kahl, and M. Oskarsson, "Branch-and-bound methods for euclidean registration problems," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 5, pp. 783–794, 2008.
- [9] H. Li, "Consensus set maximization with guaranteed global optimality for robust geometry estimation," in *2009 IEEE 12th International Conference on Computer Vision*, pp. 1074–1080, IEEE, 2009.
- [10] H. Yang and L. Carlone, "A polynomial-time solution for robust registration with extreme outlier rates," in *Robotics: Science and Systems (RSS)*, 2019.
- [11] D. Nistér, O. Naroditsky, and J. Bergen, "Visual odometry," in *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004.*, vol. 1, pp. I–I, Ieee, 2004.
- [12] L. Tang, Y. Wang, X. Ding, H. Yin, R. Xiong, and S. Huang, "Topological local-metric framework for mobile robots navigation: a long term perspective," *Autonomous Robots*, pp. 1–15, 2018.
- [13] X.-S. Gao, X.-R. Hou, J. Tang, and H.-F. Cheng, "Complete solution classification for the perspective-three-point problem," *IEEE transactions on pattern analysis and machine intelligence*, vol. 25, no. 8, pp. 930–943, 2003.
- [14] V. Lepetit, F. Moreno-Noguer, and P. Fua, "Epnnp: An accurate o (n) solution to the pnp problem," *International journal of computer vision*, vol. 81, no. 2, p. 155, 2009.
- [15] S. Ramalingam, S. Bouaziz, and P. Sturm, "Pose estimation using both points and lines for geo-localization," in *ICRA 2011-IEEE International Conference on Robotics and Automation*, pp. 4716–4723, IEEE Computer Society, 2011.
- [16] L. Kneip, M. Chli, and R. Y. Siegwart, "Robust real-time visual odometry with a single camera and an imu," in *Proceedings of the British Machine Vision Conference 2011*, British Machine Vision Association, 2011.
- [17] P. Speciale, D. Pani Paudel, M. R. Oswald, T. Kroeger, L. Van Gool, and M. Pollefeys, "Consensus maximization with linear matrix inequality constraints," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4941–4949, 2017.
- [18] K. MacTavish and T. D. Barfoot, "At all costs: A comparison of robust cost functions for camera correspondence outliers," in *2015 12th Conference on Computer and Robot Vision*, pp. 62–69, IEEE, 2015.
- [19] Q.-Y. Zhou, J. Park, and V. Koltun, "Fast global registration," in *European Conference on Computer Vision*, pp. 766–782, Springer, 2016.
- [20] M. Bosse, G. Agamennoni, I. Gilitschenski, *et al.*, "Robust estimation and applications in robotics," *Foundations and Trends® in Robotics*, vol. 4, no. 4, pp. 225–269, 2016.
- [21] G. H. Lee, F. Fraundorfer, and M. Pollefeys, "Robust pose-graph loop-closures with expectation-maximization," in *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 556–563, IEEE, 2013.
- [22] P. Agarwal, G. D. Tipaldi, L. Spinello, C. Stachniss, and W. Burgard, "Robust map optimization using dynamic covariance scaling," in *2013 IEEE International Conference on Robotics and Automation*, pp. 62–69, Ieee, 2013.
- [23] H. Yang, P. Antonante, V. Tzoumas, and L. Carlone, "Graduated non-convexity for robust spatial perception: From non-minimal solvers to global outlier rejection," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 1127–1134, 2020.
- [24] H. Yang and L. Carlone, "In perfect shape: Certifiably optimal 3d shape reconstruction from 2d landmarks," in *Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [25] V. Tzoumas, P. Antonante, and L. Carlone, "Outlier-robust spatial perception: Hardness, general-purpose algorithms, and guarantees," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019.
- [26] J. Yang, H. Li, D. Campbell, and Y. Jia, "Go-icp: A globally optimal solution to 3d icp point-set registration," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 11, pp. 2241–2254, 2015.
- [27] Y. Liu, C. Wang, Z. Song, and M. Wang, "Efficient global point cloud registration by matching rotation invariant features through translation search," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 448–463, 2018.
- [28] H. D. Sherali and A. Alameddine, "A new reformulation-linearization technique for bilinear programming problems," *Journal of Global optimization*, vol. 2, no. 4, pp. 379–410, 1992.
- [29] M. Bloesch, S. Omari, M. Hutter, and R. Siegwart, "Robust visual inertial odometry using a direct ekf-based approach," in *2015 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pp. 298–304, IEEE.
- [30] <https://opencv.org/>.
- [31] X. Ding, Y. Wang, D. Li, L. Tang, H. Yin, and R. Xiong, "Laser map aided visual inertial localization in changing environment," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 4794–4801, IEEE, 2018.
- [32] R. Mur-Artal and J. D. Tardós, "Visual-inertial monocular slam with map reuse," *IEEE Robotics and Automation Letters*, vol. 2, no. 2, pp. 796–803, 2017.