

Phase-SLAM: Phase Based Simultaneous Localization and Mapping for Mobile Structured Light Illumination Systems

Xi Zheng¹, Rui Ma¹, Rui Gao¹, Qi Hao^{*1,2}

Abstract—Structured Light Illumination (SLI) systems have been used for reliable indoor dense 3D scanning via phase triangulation. However, mobile SLI systems for 360 degree 3D reconstruction demand 3D point cloud registration, involving high computational complexity. In this paper, we propose a phase based Simultaneous Localization and Mapping (Phase-SLAM) framework for fast and accurate SLI sensor pose estimation and 3D object reconstruction. The novelty of this work is threefold: (1) developing a reprojection model from 3D points to 2D phase data towards phase registration with low computational complexity; (2) developing a local optimizer to achieve SLI sensor pose estimation (odometry) using the derived Jacobian matrix for the 6 DoF variables; (3) developing a compressive phase comparison method to achieve high-efficiency loop closure detection. The whole Phase-SLAM pipeline is then exploited using existing global pose graph optimization techniques. We build datasets from both the unreal simulation platform and a robotic arm based SLI system in real-world to verify the proposed approach. The experiment results demonstrate that the proposed Phase-SLAM outperforms other state-of-the-art methods in terms of the efficiency and accuracy of pose estimation and 3D reconstruction. The open-source code is available at <https://github.com/ZHENGXi-git/Phase-SLAM>.

I. INTRODUCTION

The SLI technology has been widely used for high-precision 3D scanning for many industrial applications with the camera-projector pair. There are usually two approaches for SLI systems to achieve 360 degree 3D reconstruction: controlled motion based and free motion based [1]. The former uses a servo motor to rotate the object along a pre-defined trajectory for multiple view scanning; the latter estimates sensor motions through local and global point cloud registration, such as Iterative Closest Point (ICP) and associated variants [2], [3]. The free-motion approach is advantageous in its flexibility but incurs high computational complexity and demands a high storage capacity.

Meanwhile, as the 2D phase data produced by SLI systems contain 3D information [4], it is appealing to utilize the phase to achieve high-efficiency pose estimation and loop closure detection. However, to develop a fully functional Phase-SLAM system has to cope with the following technological

This work is partially supported by the National Natural Science Foundation of China (No: 61773197); and the Shenzhen Nanshan District Science and Technology Innovation Bureau (No: LHTD20170007); and the Intel ICRI-IACV Research Fund (No: CG#52514373).

¹ Xi Zheng, Rui Ma, Rui Gao, and Qi Hao are with the Department of Computer Science and Engineering, Southern University of Science and Technology, Shenzhen, China. zhengx3@mail.sustech.edu.cn, mar@sustech.edu.cn, 12032493@mail.sustech.edu.cn

² Qi Hao is with the Research Institute of Trustworthy Autonomous Systems, Southern University of Science and Technology, Shenzhen, China.

* Corresponding author: Qi Hao (haoq@sustech.edu.cn)

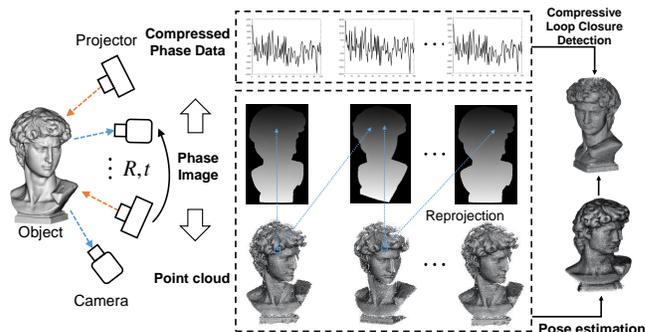


Fig. 1. A diagram of the proposed Phase-SLAM framework based on the camera-projector pair, which utilizes a 3D-to-2D reprojection model to predict the phase data for an assumed sensor pose, a local optimizer to achieve pose estimation, and a compressive method to enable fast loop closure detection.

challenges: (1) how to build the intrinsic relationship between the phase and the transformation of SLI; (2) how to develop a local optimization procedure for estimating 6 DOF motions of the SLI sensor (odometry); (3) how to achieve sparse representation and fast matching of phase data for loop closure detection. Our previous work [5] proposes a geometric reference plane to model the relationship between phases and motions of 6 DoF separately, which is complicated and inconvenient. Besides, if the loop closure detection is based on whole phase images, the memory footprint will also grow quickly as the scanning view increases. This paper presents an upgraded Phase-SLAM framework, which utilizes a 3D point to 2D phase reprojection method to build the model, a gradient based local optimizer to achieve odometry functionality and a compression method to enable efficient loop detection (Fig. 1). The main contributions of this work include,

- 1) proposing a reprojection model from 3D point to 2D phase data, which can be used to get phase estimations and measurements;
- 2) constructing a local pose optimizer with the reprojection model and the analytical expression of Jacobian matrix is derived for pose estimation;
- 3) developing a complete pipeline of Phase-SLAM framework with a compressive loop closure detection scheme and the pose graph optimization;
- 4) building simulation and real-world datasets and providing the open-source code for further development.

This paper is organized as follows. Section II introduces the related work. Section III gives an overview of the Phase-SLAM system pipeline. Section IV describes the proposed

phase-based pose estimation and compressive loop detection methods. Section V provides experiment results and discussions. Section VI concludes the paper and outline the future works. Appendix supplements the details of the Jacobian matrix in use.

II. RELATED WORK

Most visual SLAM systems are based on either direct or indirect schemes. Direct approaches [6], [7] sample pixels from image regions and minimize the photometric error. Indirect approaches [8], [9] require extra computational resources for detecting and describing features. In contrast, the proposed Phase-SLAM system is based on pixel-level phase data, which contain 3D depth information and can be extracted directly by selecting a region of interest (ROI).

A. Point Cloud Registration

SLI systems often use point cloud registration methods to achieve large fields of view scans, either local or global. Classical local registration methods, such as Point-to-Point ICP [2], minimizes the sum of distances between points and their nearest neighbours. Point-to-Plane ICP [3] assumes that each corresponding point is located on a plane, and introduces surface normals into the objective function to achieve more efficient data registration. Symmetrized objective function (SymICP) have been proposed to extend the planar convergence into a quadratic one at extra computational costs [10].

Local methods are limited by initial guesses, so structural features of point clouds are used to search for transformations globally. Point coordinates and surface normals have been used to compute the Fast Point Feature Histograms (FPFH) [11], and the coplanar 4-point sets have been chosen as features for registration (Super 4PCS) [12]. Besides, Go-ICP uses the branch-and-bound (BnB) scheme to avoid local optima [13]. Fast Global Registration (FGR) applies a Black-Rangarajan duality to achieve a more robust objective function [14]. BCPD++ formulates coherent point drift in a Bayesian setting to supervise the convergence of algorithm [15]. Compared with above 3D point cloud registration methods, our approach converts 3D point cloud registration into 2D phase data registration, resulting in much reduced computational complexity and memory footprint.

B. Loop Closure Detection

Loop closure detection can effectively eliminate the accumulating error. A plain method is randomly sampling a number of keyframes to find loop closures [16]. Odometry based approaches judge whether there is a loop closure at the current position according to the calculated map [17]. Appearance based approaches determine the loop relationship based on the similarity of two scenes [8], [9]. Bag-of-Words (BoW) based the approach [18] uses descriptors (words) for loop closure detection instead of whole images. In this paper, our loop closure detection is based on compressed phase data to reduce both computational complexity and storage space without losing much detection performance.

III. SYSTEM SETUP AND PROBLEM STATEMENT

A. System Setup

The proposed Phase-SLAM pipeline is shown in Fig. 2. Based on the SLI sensor data, corresponding 3D point clouds and the reprojection model are used to obtain phase data (Section IV-B). Then, the local pose optimization module (Section IV-C) is used to estimate sensor poses by minimizing errors between predictions and measurements of phase data. Local pose graphs are updated until the compressive loop closure detection (Section IV-E) is triggered. The pose graph optimizer then performs global optimization to eliminate the cumulative errors and revise sensor poses. Finally, poses are used to align multi-view point clouds and achieve the overall 3D object reconstruction.

We define the notations used in this paper. The initial position of the projector is chosen as the origin of the world coordinate system. $(\cdot)^w$ is the world frame, $(\cdot)^c$ is the camera frame, $(\cdot)^p$ is the projector frame, and $(\cdot)_k$ means the k -th sensor pose. The Φ and ϕ stand for the phase image and phase value at each pixel location, respectively. $\hat{(\cdot)}$ denotes the estimated value. $\mathbf{P}(x, y, z)$ is the 3D coordinate of a point. The transformation between two sensor poses is represented by vector $\Delta\mathbf{X} = [\delta x, \delta y, \delta z, \delta\alpha, \delta\beta, \delta\gamma]$, Matrix \mathbf{R} and vector \mathbf{t} represent rotation and translation from $pose_k$ to $pose_{k+1}$. \mathbf{R} and \mathbf{t} can be obtained for a given $\Delta\mathbf{X}$.

B. Problem Statement

This work aims at developing a complete SLAM system that can estimate the SLI sensor pose transformation $\Delta\mathbf{X}$ through phase data registration and achieve global 360 degree dense 3D reconstruction through pose graph optimization. At each step, 3D points are projected into the sensor imaging plane with initialized pose rotation and translation $\mathbf{R}_k, \mathbf{t}_k$ by using

$$\mathbf{u}_{k+1}(\mu, \nu) = \pi_{\mathcal{M}}(\mathbf{R}_k \mathbf{P}_k^w + \mathbf{t}_k), \quad (1)$$

where the $\pi_{\mathcal{M}}$ is the perspective transformation with the projection matrix \mathcal{M} , means $\mathbb{R}^3 \rightarrow \mathbb{R}^2$ that projects a 3D point onto the imaging plane. \mathbf{u}_{k+1} is the pixel position, which is used for obtaining phase data estimations $\hat{\phi}$ and measurements ϕ . Obtained the $\hat{\phi}$ and ϕ , the sensor pose transformation $\Delta\mathbf{X}$ is estimated by

$$\Delta\mathbf{X}^* = \arg \min_{\Delta\mathbf{X}} \mathbf{F}(\Delta\mathbf{X}), \quad (2)$$

where

$$\mathbf{F}(\Delta\mathbf{X}) = \frac{1}{2} \sum \left\| \hat{\phi} - \phi \right\|^2. \quad (3)$$

Such a local optimization procedure requires computing the Jacobian matrix iteratively until it converges. The loop closure detection and pose graph optimization will be also needed to reduce estimation errors.

C. SLI Scanning

In the camera-projector based SLI system, the Phase Measuring Profilometry (PMP) method is used to calculate the phase image, as shown in Fig. 3. The camera captures

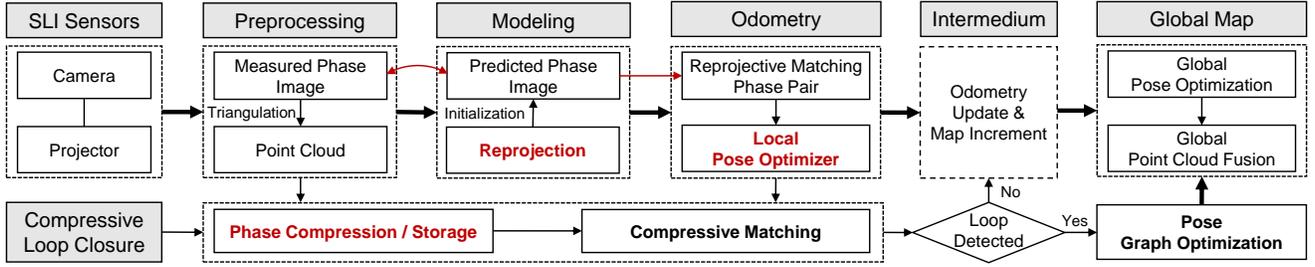


Fig. 2. The system diagram of the proposed Phase-SLAM. Based on the SLI phase image, 3D point clouds for each new sensor pose can be computed, which are used for phase data prediction (Section IV-B) and local pose estimation (Section IV-C). The compressive loop closure detection is performed to trigger the global pose graph optimizer (Section IV-E). Finally, the refined sensor poses are used to achieve 360 degree 3D point clouds of the object under scanning.

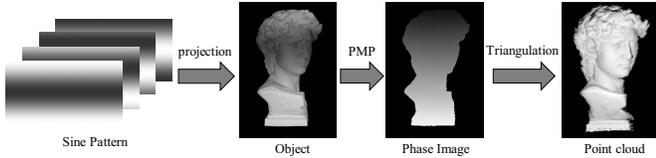


Fig. 3. An illustration of SLI imaging system. PMP uses the images of projection patterns to compute phase images; the 3D point clouds are then obtained by triangulation with the calibrated camera-projector parameters.

the raw images of sine patterns deformed by the scanned surface, given by

$$I_n^c(\mu, \nu) = A + B \cos\left(\Phi(\mu, \nu) - \frac{2\pi n}{N}\right), \quad (4)$$

where $n = 1, 2, \dots, N$ (the number of patterns), A and B are the background brightness and intensity modulation, respectively. The phase image $\Phi(\mu, \nu)$ can then be calculated by [4]

$$\Phi = \arctan\left[\frac{\sum_{n=1}^N I_n^c \sin(2\pi n/N)}{\sum_{n=1}^N I_n^c \cos(2\pi n/N)}\right]. \quad (5)$$

IV. PROPOSED METHODS

This section investigates the geometric model among 3D point, phase data and sensor pose. Comparing with our previous work [5], this work develops a more intuitive and simpler model based on reprojective transformation method. After phase data pairing, the sensor pose motion can be estimated through least-square optimization between phase predictions and measurements. A compressed sensing scheme is adopted to achieve fast loop closure detection for global pose graph optimization.

A. Phase Values under Epipolar Constraint

In a SLI system, we regard the projector as another camera, which has similar projection parameters and perspective principles with it. As shown in the left of Fig. 4, according to the epipolar constraint, a phase value obtained from the phase image can correspond to a pixel location in the “camera” imaging plane (phase pattern), like stereo-vision [19]. In PMP method, the phase pattern is actively projected by the projector, so pixel locations and phase values on pattern plane have a fixed and known relevance. As shown in the right of Fig. 4, the phase value of each column in phase

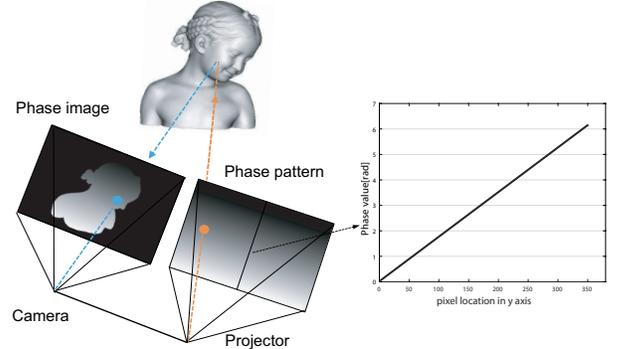


Fig. 4. An illustration of SLI imaging principle. The projector is regarded as a camera with the phase pattern. The phase pair should be performed under the epipolar constraint. The phase values of each column of phase pattern are linearly increased from 0 to 2π .

pattern is linearly increased from 0 to 2π and each row in the pattern is the same. This means when we get the phase value of a 3D point \mathbf{P} from the phase image, we can know its ordinate under the projector’s pattern coordinate. Vice versa, if we knew the projective coordinate (μ^p, ν^p) (just only ν^p) of a point \mathbf{P} in the pattern, we could get its corresponding phase value on the phase image by

$$\phi = 2\pi\nu^p/H^p, \quad (6)$$

where H^p is the row height of the projector’s imaging plane.

B. Phase Pairing Based on Reprojective Transformation

As shown in the Fig. 5, a 3D point \mathbf{P} is measured by the SLI sensor at the $pose_k$ with the coordinate $\mathbf{P}_k = [x, y, z]^\top$. Assuming the transformation: rotation matrix $\mathbf{R}(\delta\alpha, \delta\beta, \delta\gamma) \in SO(3)$ and translation vector $\mathbf{t} = [\delta x, \delta y, \delta z]^\top$, the SLI move to $pose_{k+1}$ by it and the point \mathbf{P} will have a new coordinate $\mathbf{P}_{k+1} = [x', y', z']^\top$, given by

$$\mathbf{P}_{k+1} = \mathbf{R}\mathbf{P}_k + \mathbf{t}. \quad (7)$$

Based on the new coordinate, the point \mathbf{P} is reprojected into camera and projector imaging plane in $pose_{k+1}$ to get two pixel locations on them: $\mathbf{u}_{k+1}^c(\mu_{k+1}^c, \nu_{k+1}^c)$ and $\mathbf{u}_{k+1}^p(\mu_{k+1}^p, \nu_{k+1}^p)$ by the transformation $\pi_{\mathcal{M}}$ (Eq. (1)), respectively. On the projector imaging plane (the phase pattern in Fig. 4), when the reprojection ordinate of point \mathbf{P} : ν_{k+1}^p is known, the phase value prediction $\hat{\phi}_{k+1}$ can be obtained by

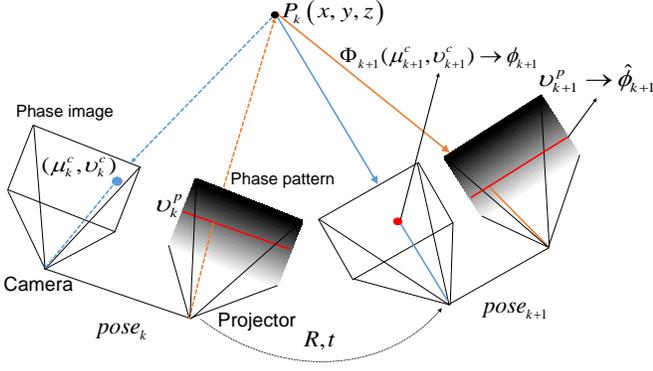


Fig. 5. An illustration of reprojection from 3D points to 2D phase data. A 3D point P obtained by $pose_k$ is reprojected into the imaging plane of camera and projector in $pose_{k+1}$ with a assumed rotation \mathbf{R} and translation \mathbf{t} . Then the errors between the predicted and measured phase data ($\hat{\phi}_{k+1}$ and ϕ_{k+1}) are minimized with respect to \mathbf{R} and \mathbf{t} .

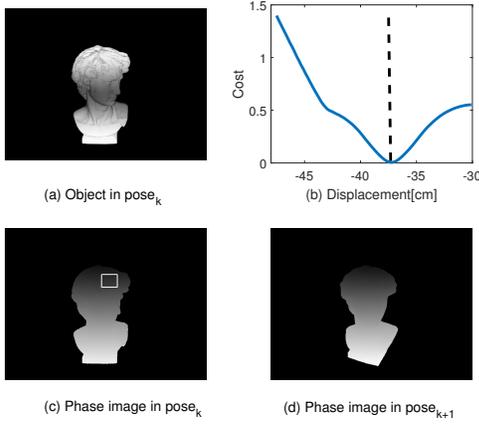


Fig. 6. An illustration of a simple local pose optimization process. (a) an object raw image acquired at $pose_k$; (c) the phase image acquired at $pose_k$ with a white ROI; (d) the phase images acquired at a new sensor pose. (b) shows the plot of errors between two sets of phase data (within the ROI) with respect to the displacements of the SLI sensor. It can be seen that such an optimization process can be converged to the local minimum [21].

Eq. (6). So, combining Eq. (1, 6, 7), the phase value of the a 3D point in the phase pattern can be estimated by using

$$\hat{\phi}_{k+1} = \frac{2\pi}{Hp} \left(\frac{f_y^p (R_{21}x + R_{22}y + R_{23}z + \delta y)}{R_{31}x + R_{32}y + R_{33}z + \delta z} + C_y^p \right), \quad (8)$$

where R_{ij} is the ij^{th} element of \mathbf{R} , f_y^p , C_y^p is calibration parameter (projector's focal length and principal point along the row of the projector imaging plane, respectively).

On the camera imaging plane, the phase value measurement ϕ_{k+1} can be obtained from the phase image Φ_{k+1} at the pixel location $(\mu_{k+1}^c, \nu_{k+1}^c)$, which can be computed by

$$\begin{aligned} \mu_{k+1}^c &= \frac{m_{11}x' + m_{12}y' + m_{13}z'}{m_{31}x' + m_{32}y' + m_{33}z'} \\ \nu_{k+1}^c &= \frac{m_{21}x' + m_{22}y' + m_{23}z'}{m_{31}x' + m_{32}y' + m_{33}z'}, \end{aligned} \quad (9)$$

where m_{ij} is the ij^{th} element of projection matrix \mathcal{M} . When μ_{k+1}^c and ν_{k+1}^c are not integers, bilinear interpolation on Φ_{k+1} can be used to calculate phase data at integer indices.

C. Local Pose Optimizer

In the local optimizer, the state variables are defined as $\Delta\mathbf{X}(\delta x, \delta y, \delta z, \delta\alpha, \delta\beta, \delta\gamma)$, which is equivalent to \mathbf{R} and \mathbf{t} . The error \mathbf{e} between $\hat{\phi}_{k+1}$ and ϕ_{k+1} are given by

$$\mathbf{e} = \hat{\phi}_{k+1}(\Delta\mathbf{X}) - \Phi_{k+1}(\mu_{k+1}(\Delta\mathbf{X}), \nu_{k+1}(\Delta\mathbf{X})). \quad (10)$$

The objective function is shown in

$$\mathbf{F}(\Delta\mathbf{X}) = \frac{1}{2} \sum_{\langle \mu, \nu \rangle \in \mathbb{R}} \|\mathbf{e}\|^2, \quad (11)$$

the \mathbb{R} is a ROI in phase images, $(\cdot)^i$ is the i -th point in ROI. $\mathbf{e} = [\mathbf{e}^1, \mathbf{e}^2, \dots, \mathbf{e}^n]^\top$.

The proposed function Eq. (11) can be solved by iterative gradient-based methods [20]. Given the initial value $\Delta\tilde{\mathbf{X}}$, the cost function can be approximated by Taylor expand about $\Delta\tilde{\mathbf{X}}$, and $\mathbf{F}(\Delta\tilde{\mathbf{X}} + \Delta) \approx \mathbf{F}(\Delta\tilde{\mathbf{X}}) + \nabla\mathbf{F}\Delta$, where

$$\begin{aligned} \nabla\mathbf{F} &= \mathbf{J}^\top \mathbf{e} \\ \mathbf{J} &= \partial\mathbf{e}/\partial\Delta\mathbf{X}. \end{aligned} \quad (12)$$

\mathbf{J} is the Jacobian matrix, the optimization increment Δ is computed by $\lambda\Delta = -\mathbf{J}^\top \mathbf{e}$, which is the negative gradient direction of \mathbf{F} , λ controls the size of steps. The solution is updated by $\Delta\mathbf{X}_{i+1} = \Delta\mathbf{X}_i + \Delta_i$, where i is the iterative index [19], [20].

Fig. 6 shows a simple example the optimization process. Fig. 6 (a) shows an object image. (c) shows the corresponding phase image with a ROI. (d) shows the phase images acquired at a new sensor pose. (b) shows the plot of errors between two sets of phase data (within the ROI) with respect to the displacements of the SLI sensor. It can be seen that such an optimization process can be converged to the local minimum [21].

D. The Jacobian Matrix

According to Eq. (10, 12), the Jacobian matrix of \mathbf{e}^i ($i = 1, 2, \dots, n$) is given by

$$\mathbf{J}^i = \frac{\partial\mathbf{e}^i}{\partial\Delta\mathbf{X}} = \frac{\partial\hat{\phi}_{k+1}^i}{\partial\Delta\mathbf{X}} - \left(\frac{\partial\Phi_{k+1}}{\partial\mu_{k+1}^i} \frac{\partial\mu_{k+1}^i}{\partial\Delta\mathbf{X}} + \frac{\partial\Phi_{k+1}}{\partial\nu_{k+1}^i} \frac{\partial\nu_{k+1}^i}{\partial\Delta\mathbf{X}} \right), \quad (13)$$

where $\partial\Phi_{k+1}/\partial\mu_{k+1}$ and $\partial\Phi_{k+1}/\partial\nu_{k+1}$ are vertical and horizontal gradients of Φ_{k+1} , computed by pixel difference. More details in Eq. (13) are provided in the Appendix, and the other term in the Eq. (12) is substituted by

$$\mathbf{J}^\top \mathbf{e} = \sum_{i=1}^n \mathbf{e}^i \frac{\partial\mathbf{e}^i}{\partial\Delta\mathbf{X}}. \quad (14)$$

E. Loop Closure Detection

The proposed Phase-SLAM utilizes the Compressive Sensing (CS) technique to reduce computational complexity and data storage space for loop closure detection. The compressibility of an image is determined by its sparsity. More sparse images will lose less information after compression and the sparse image contains less high-frequency information [22]. A haar wavelet bases and L_1 norm are used to illustrate the

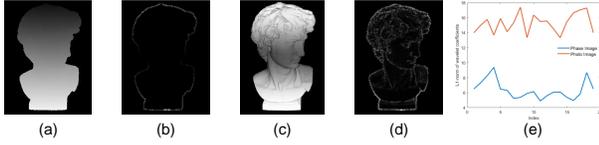


Fig. 7. An illustration of the sparsity of phase and photo images. (a,c) A phase image and a photo image; (b,d) the corresponding wavelet coefficients of two images; (e) the wavelet coefficient L_1 norms of two types of images within one SLAM loop.

degree of sparsity of phase images. As shown in Fig. 7, a phase image and a photo image are projected upon wavelet bases first. Then the L_1 norms of the wavelet coefficients of two types of images within one SLAM loop are compared in Fig. 7 (e). It can be seen that the L_1 norms of the wavelet coefficients of phase images are much smaller than photo images, indicating the degree of sparsity of phase images is much smaller than photo images.

According to the CS theory [22], two signals A_1 and A_2 are distinguishable after compression if the matrix C satisfies

$$2(1 - \delta_{2s}) \leq \|CA_2 - CA_1\|_2^2 \leq 2(1 + \delta_{2s}), \quad (15)$$

where δ_{2s} is a constant, and C is Gaussian matrix. The compressed signals $y_{n \times 1} = C_{n \times N} A_{N \times 1}$ $n < N$, has quite smaller size than the original signals. For a 2D phase image Φ , we first reshape it into a 1D vector Φ' . The reshaped phase data Φ' can be recovered by the compressed signal $y = C\Phi'$, and the error between two compressive phase vector is shown in

$$dy = \|C\Phi'_2 - C\Phi'_1\|_2^2. \quad (16)$$

When dy is smaller than a threshold, the loop-closure is detected.

F. The Pipeline of Phase-SLAM

After successful loop-closure detection, the pose graph optimization technique [23] will be used to eliminates cumulative error and refine poses. The pose sequences in our system usually have a large interval during the scanning process, so every estimated pose is a vertex in the pose graph optimizer.

V. EXPERIMENT RESULTS AND DISCUSSIONS

The proposed Phase-SLAM system was evaluated with both the Unreal Engine 4 (UE4) simulator and real-world experiments. All experiments were implemented on a PC with an Intel Core i7-9800K CPU @ 3.6GHz.

A. Simulation Experiments

The simulation dataset was collected with the Airsim plugin in UE4. Different 3D models were used as targets, and the virtual SLI device moved around the target along a radius of 120 cm and with a rotation interval of 20 degrees. The simulated dataset is based on three models namely David, Elephant and Dancing girl, which contains calibration parameters, phase images and ground-truth poses. The baseline methods include four state-of-the-art local methods, namely Point-to-Point ICP [2], Point-to-Plant ICP [3], SymICP [10]

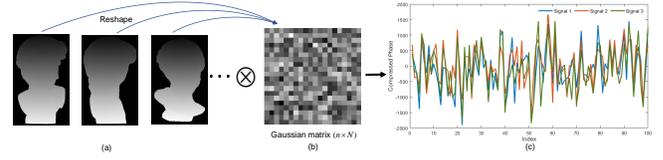


Fig. 8. An illustration of the proposed compressive loop closure detection. (a) Three phase images; (b) the Gaussian pseudo-random matrix used for compressive projection; (c) three sets of compressed signals corresponding to three phase images used for loop closure detection.

and FPFH [11], and two SOTA global methods, namely FGR [14] and BCPD++ [15]. Local methods were conducted based on Point Cloud Library (PCL) implementation [24]. Global methods were based on open-source code. The numbers of iterations of Point-to-Point ICP, Point-to-Plane ICP, and SymICP were chosen as 30; FPFH was 10000. The implementation of FGR and BCPD++ used the recommended parameters.

The compression of phase images is shown in Fig. 8. In simulation experiments, the resolution of phase images is 640×480 (Fig. 8 (a)), the size of Gaussian compressive random matrix is chosen as 100×307000 (Fig. 8 (b)), that is, the compression ratio is 3070:1 and the size of the compressed phase signal is 100×1 . Fig. 8 (c) shows the compressed signals corresponding to three phase images like Fig. 8 (a). It can be seen that the three sets of signals are distinguishable in terms of the peaks and valleys for loop closure detection. Furthermore, experiment results show that the time consumption of the back-end optimization using CS technique can be reduced by 20% than using original phase images.

Fig. 9 (a)-(c) shows the 3D reconstruction results (top) and ground truth (bottom) for the three simulation targets (David, Elephant and Dancing girl) using the proposed Phase-SLAM with loop closure. More quantified reconstruction errors are shown in Fig. 15. Fig. 10 shows the relative pose errors (RPE) [25] in rotation (top) and translation (bottom), respectively by using five methods with David dataset. It can be seen that the proposed Phase-SLAM method with loop closure (PhaseS-Loop) outperforms other four methods. The median RPE of Phase-SLAM is 0.81 degree and 0.94cm in rotation and translation, respectively. Table I shows the root mean squared error (RMSE) of absolute trajectory error (ATE) [25] and the computation time for three different datasets. The average RMSE of our approach is 1.06cm, which is better than other methods. Actually, Phase-SLAM with loop closure outperforms PhaseS by 38.5%. In simulations, the average number of 3D points corresponding to the image is around 50000. BCPD++ has the highest computation speed among the 6 existing methods. Our approach is still almost two times faster than BCPD++. And the average running time of the back-end optimization is 40.7 ms.

B. Real-World Experiments

Fig. 11 shows the experiment setup, where the SLI sensor, consisting of a projector (DLP3000 DMD from TI) and an industrial camera (1280×1024 resolution from HIKVISION),

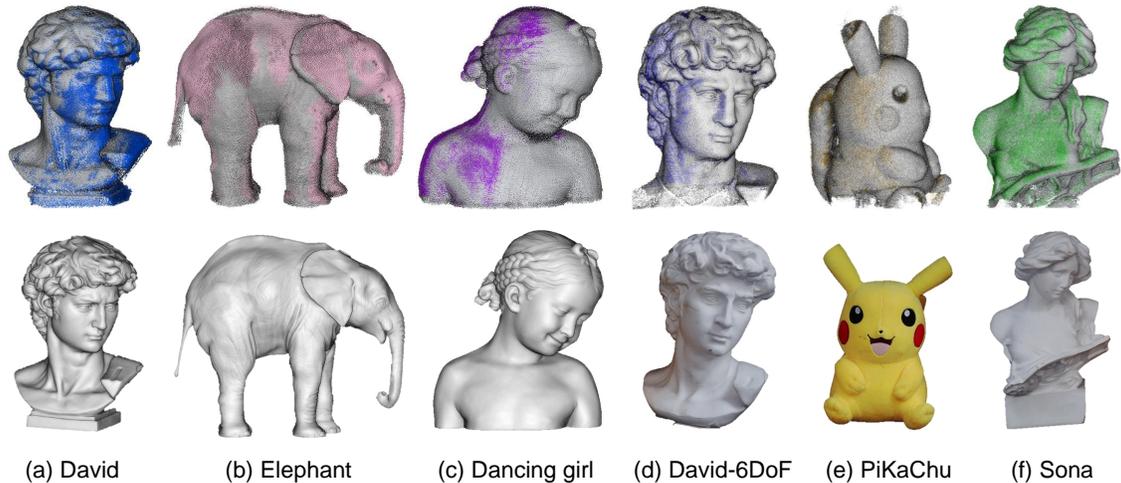


Fig. 9. A comparison of global 3D point cloud registration results along with ground-truth. (a-c) Simulation datasets named David, Elephant and Dancing girl; (e-f) real-world datasets named David-6DoF, PiKaChu and Sona. (Top Row) The ground truth (in gray) and reconstruction results (in other colors) by using the proposed Phase-SLAM; (Bottom Row) the 3D objects.

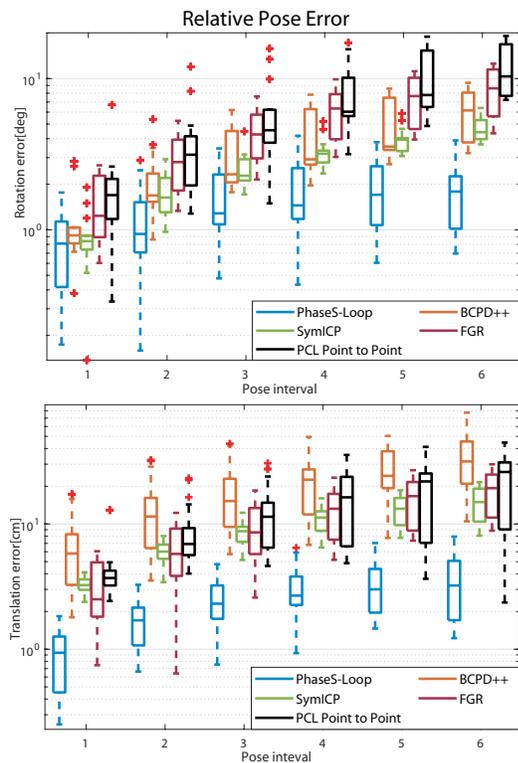


Fig. 10. The plot of Relative Pose Errors [25] of (top) rotation and (bottom) translation in simulation. PhaseS-Loop denotes Phase-SLAM with pose graph optimization.

is mounted on a UR5 robotic arm. Fig. 9 (d-f) show two plaster statues (David, Sona) and a plush toy (PiKaChu) used to build real-world datasets, namely David-6DoF, David-3DoF, Sona-3DoF and PiKaChu-3DoF, where 6DoF and 3DoF stand for six and three degrees of freedom motions, respectively. The David-6DoF dataset includes 31 random poses; Each 3DoF dataset has 37 poses at equal intervals of 10 degrees and a radius of 60cm.

Fig. 12 shows the RPE results of five different methods

TABLE I
RMSE OF ATE (CM) / COMPUTATION TIME (S)

Method	David	Elephant	Dancing Girl
PhaseS-Loop[ours]	1.39 / 1.52	0.72 / 1.58	1.07 / 0.82
PhaseS[ours]	2.32 / 1.49	2.40 / 1.23	2.05 / 0.76
BCPD++[15]	25.53 / 2.87	87.09 / 3.46	15.57 / 2.68
SymICP[10]	6.35 / 109.10	7.06 / 146.71	8.17 / 104.22
Point to Plane[3]	6.76 / 55.15	13.07 / 65.73	10.57 / 45.89
Point to Point[2]	17.17 / 32.56	19.35 / 44.68	40.26 / 30.75
FGR[14]	11.12 / 25.56	15.67 / 34.68	38.74 / 26.90
FPFH[11]	8.99 / 70.59	25.10 / 79.09	32.70 / 59.68

RMSE of ATE: The root mean squared error of absolute trajectory error.

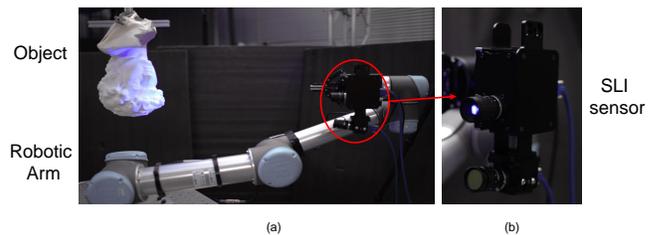


Fig. 11. (a)The real-world experiment setup where an object is fixed on a bracket and the SLI sensor is installed on a UR5 robotic arm. (b)The SLI sensor consists of a DLP3000 projector and a HIKVISION camera.

using the David-6DoF dataset. It is clear that the proposed method (PhaseS-Loop) has a better performance than other four methods. Table II is the RMSE of ATE and the computation time for four different datasets using eight different methods. It can be seen that the proposed method outperforms other six methods in both terms of accuracy and computation time. The 3D object reconstruction results using the proposed method under real-world datasets are shown in Fig. 9 (d-f). Fig. 13 illustrates the estimated SLI sensor trajectory and the ground truth under David-6DoF dataset, where the total trajectory length is 3.967m. Fig. 14 shows the pose estimation errors by using local methods (SymICP, Point-to-Plant and Point-to-Point ICP)

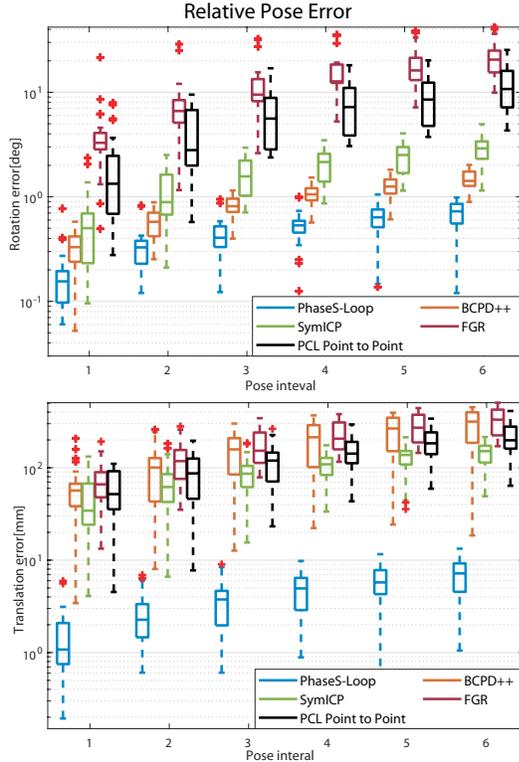


Fig. 12. The plot of Relative Pose Errors [25] of (top) rotation and (bottom) translation in real-world experiments.

TABLE II
RMSE OF ATE (MM) / COMPUTATION TIME (S)

Method	David-6DoF	David-3DoF	PiKaChu	Sona
PhaseS-Loop	4.69/4.20	4.71/3.19	2.09/3.30	1.83/3.18
PhaseS	6.12/4.17	5.74/3.17	3.27/3.27	2.29/3.15
BCPD++	244.41/2.79	53.01/2.93	21.72/3.72	22.39/3.81
SymICP	99.28/374.26	28.78/345.25	35.68/268.69	30.88/242.12
Point to Plane	101.5/168.53	33.66/152.13	33.97/119.23	36.65/109.15
Point to Point	170.2/118.42	89.28/101.54	70.36/81.34	84.44/77.9
FGR	282.3/238.25	224.21/213.15	231.17/302.45	149.92/191.7
FPFH	109.85/386.8	95.51/153.66	254.34/202.14	90.81/217.9

RMSE of ATE: The root mean squared error of absolute trajectory error.

and our method without global optimization under different initial values. We can see that the proposed method is least sensitive to initial values. Fig. 15 shows a radar chart of seven methods without global optimization for all datasets using five performance metrics (Hausdorff distance, computation time, translation/rotation errors, and storage space) for a comprehensive evaluation. The Hausdorff distance is used to describe the dissimilarity between reconstructed point clouds and the ground-truth [26]. It is obvious that the proposed method has the superior performance in all those metrics.

VI. CONCLUSION

This paper presents a phase based Simultaneous Localization and Mapping (Phase-SLAM) pipeline for fast and accurate SLI sensor pose estimation and 3D object reconstruction. The proposed reprojection model and local pose optimizer can achieve the odometry functionality with high efficiency,

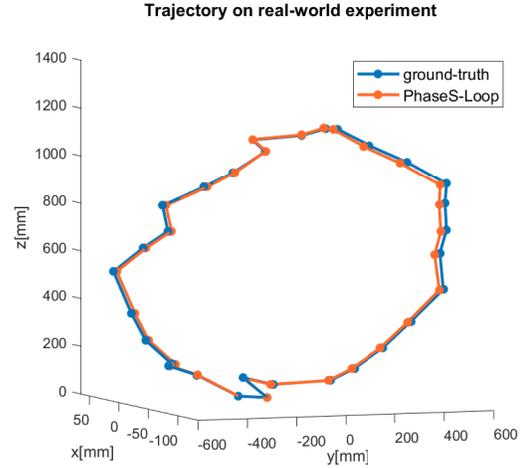


Fig. 13. The plot of the estimated sensor trajectory by using the full pipeline of the proposed Phase-SLAM on David-6DoF. The ground-truth is obtained via the UR5 robotic arm.

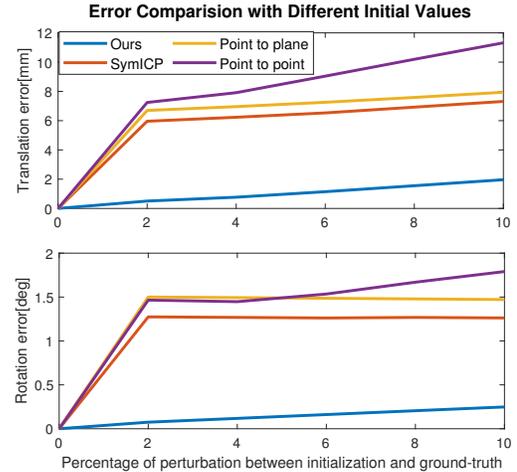


Fig. 14. The plot of registration errors of (top) translation and (bottom) rotation with respect to different sensor pose initializations. The x-axis is the percentage of perturbation for pose initialization with respect to the ground-truth.

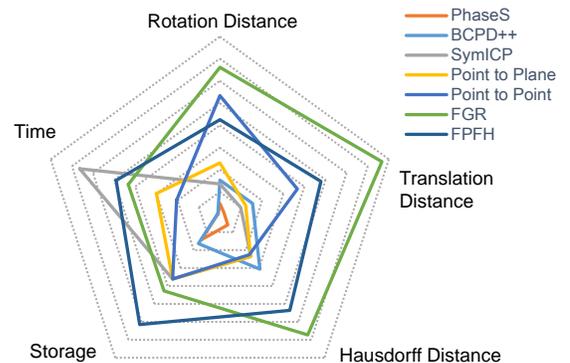


Fig. 15. The radar chart of 5 performance metrics for 7 different algorithms. The rotation and translation errors are measured via the Euler distances; the Hausdorff distance is used to measure the dissimilarity between two point clouds.

accuracy and low sensitivity to initial pose knowledge. The proposed compressive loop closure detection technique can reduce both the loop closure computational time and data storage space. Even without global optimization, the proposed local data registration method outperforms six other existing 3D point cloud based methods in terms of sensor pose estimation accuracy, storage space, computation time and 3D reconstruction errors. The code of our framework and the dataset in use are available online.

APPENDIX

The analytic expression of the Jacobian of \mathbf{e} with respect to $\delta x, \delta y, \delta z, \delta \alpha, \delta \beta, \delta \gamma$ is provided in this section. The intermediate terms are given by

$$\begin{aligned}
g_x &= \partial \Phi_{k+1} / \partial \mu_{k+1}, g_y = \partial \Phi_{k+1} / \partial \nu_{k+1} \\
K &= H_p / (2\pi), s_{k+1} = m_{31}x' + m_{32}y' + m_{33}z' \\
\mu_1 &= m_{11} - m_{31}\mu_{k+1}^c, \nu_1 = m_{21} - m_{31}\nu_{k+1}^c \\
\mu_2 &= m_{12} - m_{32}\mu_{k+1}^c, \nu_2 = m_{22} - m_{32}\nu_{k+1}^c \\
\mu_3 &= m_{13} - m_{33}\mu_{k+1}^c, \nu_3 = m_{23} - m_{33}\nu_{k+1}^c \\
J_{x\alpha} &= R_{13}y - R_{12}z, J_{y\alpha} = R_{23}y - R_{22}z \\
J_{z\alpha} &= R_{33}y - R_{32}z \\
J_{\mu\alpha} &= \mu_1 J_{x\alpha} + \mu_2 J_{y\alpha} + \mu_3 J_{z\alpha} / s_{k+1} \\
J_{\nu\alpha} &= \nu_1 J_{x\alpha} + \nu_2 J_{y\alpha} + \nu_3 J_{z\alpha} / s_{k+1} \\
J_{x\beta} &= -x \sin \delta \beta \cos \delta \gamma - y \sin \delta \alpha \cos \delta \beta \cos \delta \gamma \\
&\quad - z \cos \delta \alpha \cos \delta \beta \cos \delta \gamma \\
J_{y\beta} &= -x \sin \delta \beta \sin \delta \gamma - y \sin \delta \alpha \cos \delta \beta \sin \delta \gamma \\
&\quad - z \cos \delta \alpha \cos \delta \beta \sin \delta \gamma \\
J_{z\beta} &= -x \cos \delta \beta - y \sin \delta \alpha \sin \delta \beta - z \cos \delta \alpha \sin \delta \beta \\
J_{\mu\beta} &= \mu_1 J_{x\beta} + \mu_2 J_{y\beta} + \mu_3 J_{z\beta} / s_{k+1} \\
J_{\nu\beta} &= \nu_1 J_{x\beta} + \nu_2 J_{y\beta} + \nu_3 J_{z\beta} / s_{k+1} \\
J_{x\gamma} &= \delta y - y', J_{y\gamma} = x' - \delta x, J_{y\gamma} = 0 \\
J_{\mu\gamma} &= \mu_1 J_{x\gamma} + \mu_2 J_{y\gamma} + \mu_3 J_{z\gamma} / s_{k+1} \\
J_{\nu\gamma} &= \nu_1 J_{x\gamma} + \nu_2 J_{y\gamma} + \nu_3 J_{z\gamma} / s_{k+1}
\end{aligned} \tag{17}$$

The analytic expression of Jacobian is then given by

$$\begin{aligned}
\partial \mathbf{e} / \partial \delta x &= -(g_x \mu_1 + g_y \nu_1) / s_{k+1} \\
\partial \mathbf{e} / \partial \delta y &= f_p / (K z') - (g_x \mu_2 + g_y \nu_2) / s_{k+1} \\
\partial \mathbf{e} / \partial \delta z &= f_p y' / (K z'^2) - (g_x \mu_3 + g_y \nu_3) / s_{k+1} \\
\partial \mathbf{e} / \partial \delta \alpha &= f_p (J_{y\alpha} z' - J_{z\alpha} y') / (K z'^2) - (g_x J_{\mu\alpha} + g_y J_{\nu\alpha}) \\
\partial \mathbf{e} / \partial \delta \beta &= f_p (J_{y\beta} z' - J_{z\beta} y') / (K z'^2) - (g_x J_{\mu\beta} + g_y J_{\nu\beta}) \\
\partial \mathbf{e} / \partial \delta \gamma &= f_p (J_{y\gamma} z' - J_{z\gamma} y') / (K z'^2) - (g_x J_{\mu\gamma} + g_y J_{\nu\gamma})
\end{aligned} \tag{18}$$

REFERENCES

- [1] J. Salvi, J. Pages, and J. Batlle, "Pattern codification strategies in structured light systems," *Pattern recognition*, vol. 37, no. 4, pp. 827–849, 2004.
- [2] P. J. Besl and N. D. McKay, "Method for registration of 3-d shapes," in *Sensor fusion IV: control paradigms and data structures*, vol. 1611. International Society for Optics and Photonics, 1992, pp. 586–606.
- [3] K.-L. Low, "Linear least-squares optimization for point-to-plane icp surface registration," *Chapel Hill, University of North Carolina*, vol. 4, no. 10, pp. 1–3, 2004.

- [4] Y. Wang, K. Liu, Q. Hao, X. Wang, D. L. Lau, and L. G. Hassebrook, "Robust active stereo vision using kullback-leibler divergence," *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 3, pp. 548–563, 2012.
- [5] X. Zheng, R. Ma, R. Gao, and Q. Hao, "Phase-slam: Mobile structured light illumination for full body 3d scanning," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, pp. 1617–1624.
- [6] J. Engel, T. Schöps, and D. Cremers, "Lsd-slam: Large-scale direct monocular slam," in *Computer Vision – ECCV 2014*, D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, Eds. Cham: Springer International Publishing, 2014, pp. 834–849.
- [7] J. Engel, V. Koltun, and D. Cremers, "Direct sparse odometry," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 3, pp. 611–625, 2017.
- [8] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos, "Orb-slam: a versatile and accurate monocular slam system," *IEEE transactions on robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.
- [9] T. Qin, P. Li, and S. Shen, "Vins-mono: A robust and versatile monocular visual-inertial state estimator," *IEEE Transactions on Robotics*, vol. 34, no. 4, pp. 1004–1020, 2018.
- [10] S. Rusinkiewicz, "A symmetric objective function for icp," *ACM Transactions on Graphics (TOG)*, vol. 38, no. 4, pp. 1–7, 2019.
- [11] R. B. Rusu, N. Blodow, and M. Beetz, "Fast point feature histograms (fpfh) for 3d registration," in *2009 IEEE international conference on robotics and automation*. IEEE, 2009, pp. 3212–3217.
- [12] N. Mellado, D. Aiger, and N. J. Mitra, "Super 4pcs fast global pointcloud registration via smart indexing," in *Computer graphics forum*, vol. 33, no. 5. Wiley Online Library, 2014, pp. 205–215.
- [13] 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.
- [14] Q.-Y. Zhou, J. Park, and V. Koltun, "Fast global registration," in *European conference on computer vision*. Springer, 2016, pp. 766–782.
- [15] O. Hirose, "Acceleration of non-rigid point set registration with downsampling and gaussian process regression," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [16] F. Endres, J. Hess, J. Sturm, D. Cremers, and W. Burgard, "3d mapping with an rgb-d camera," *IEEE transactions on robotics*, vol. 30, no. 1, pp. 177–187, 2013.
- [17] D. Hahnel, W. Burgard, D. Fox, and S. Thrun, "An efficient fastslam algorithm for generating maps of large-scale cyclic environments from raw laser range measurements," in *Proceedings 2003 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2003)(Cat. No. 03CH37453)*, vol. 1. IEEE, 2003, pp. 206–211.
- [18] D. Gálvez-López and J. D. Tardos, "Bags of binary words for fast place recognition in image sequences," *IEEE Transactions on Robotics*, vol. 28, no. 5, pp. 1188–1197, 2012.
- [19] A. M. Andrew, "Multiple view geometry in computer vision," *Kybernetes*, 2001.
- [20] J. Nocedal and S. Wright, *Numerical optimization*. Springer Science & Business Media, 2006.
- [21] Y. Zhou, G. Gallego, H. Rebecq, L. Kneip, H. Li, and D. Scaramuzza, "Semi-dense 3d reconstruction with a stereo event camera," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 235–251.
- [22] E. J. Candes, "The restricted isometry property and its implications for compressed sensing," *Comptes rendus mathématique*, vol. 346, no. 9-10, pp. 589–592, 2008.
- [23] G. Grisetti, R. Kümmerle, C. Stachniss, and W. Burgard, "A tutorial on graph-based slam," *IEEE Intelligent Transportation Systems Magazine*, vol. 2, no. 4, pp. 31–43, 2010.
- [24] R. B. Rusu and S. Cousins, "3d is here: Point cloud library (pcl)," in *2011 IEEE international conference on robotics and automation*. IEEE, 2011, pp. 1–4.
- [25] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of rgb-d slam systems," in *2012 IEEE/RSJ international conference on intelligent robots and systems*. IEEE, 2012, pp. 573–580.
- [26] A. A. Taha and A. Hanbury, "An efficient algorithm for calculating the exact hausdorff distance," *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 11, pp. 2153–2163, 2015.