

# HybridFusion: LiDAR and Vision Cross-Source Point Cloud Fusion

Yu Wang<sup>1</sup>, Shuhui Bu<sup>1\*</sup>, Lin Chen<sup>1</sup>, Yifei Dong<sup>1</sup>, Kun Li<sup>1</sup>, Xuefeng Cao<sup>2</sup>, Ke Li<sup>2</sup>

**Abstract**—Recently, cross-source point cloud registration from different sensors has become a significant research focus. However, traditional methods confront challenges due to the varying density and structure of cross-source point clouds. In order to solve these problems, we propose a cross-source point cloud fusion algorithm called HybridFusion. It can register cross-source dense point clouds from different viewing angle in outdoor large scenes. The entire registration process is a coarse-to-fine procedure. First, the point cloud is divided into small patches, and a matching patch set is selected based on global descriptors and spatial distribution, which constitutes the coarse matching process. To achieve fine matching, 2D registration is performed by extracting 2D boundary points from patches, followed by 3D adjustment. Finally, the results of multiple patch pose estimates are clustered and fused to determine the final pose. The proposed approach is evaluated comprehensively through qualitative and quantitative experiments. In order to compare the robustness of cross-source point cloud registration, the proposed method and generalized iterative closest point method are compared. Furthermore, a metric for describing the degree of point cloud filling is proposed. The experimental results demonstrate that our approach achieves state-of-the-art performance in cross-source point cloud registration.

## I. INTRODUCTION

With the rapid development of 3D data acquisition technology, numerous methods have emerged that can efficiently and accurately reconstruct 3D data, such as Structure from Motion (SfM), real-time dense 3D reconstruction based on monocular or depth camera vision, and LiDAR mapping. Although a single sensor may use sequential acquisition to create a comprehensive point cloud, the limited acquisition viewing angle may not cover all target surfaces, resulting in inevitable loss of some content.

For example, DenseFusion [1] uses the high-altitude UAV viewing angle to conduct real-time 3D dense reconstruction, but only the surface above the objects can be reconstructed, which is difficult to completely describe 3D information of all object surface. Moreover, the use of multi-view scanning with a single sensor not only incurs significant time cost but also present a challenge to ensure the stability of the mapping algorithm due to potential variations in scene structures. Therefore, it is especially important to use multi-view and multi-sensor to collect data and then conduct cross-source point cloud fusion. However, cross-source point cloud fusion is a challenging task due to various difficulties, including differences in density distribution and missing data.

There are many attempts to solve cross-source point cloud registration. Cross-source graph matching (CSGM) [2]

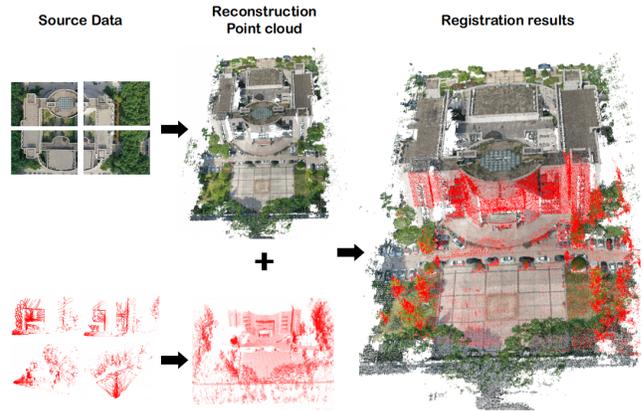


Fig. 1. Our method can register cross-source point clouds from different viewing angle. Left part is the dense point cloud generated by vision and LiDAR, and right part is the result of registration.

transforms the registration problem into a graph matching problem, subsequently, local Iterative Closest Point (ICP) refinement is conducted, leveraging the global graph matching results. Mellado *et al.* [3] employ a combination of global alignment and local refinement techniques, namely Random Sample Consensus (RANSAC) and Iterative Closest Point (ICP) algorithms, to register point clouds from different sources. Huang *et al.* [4] propose a method that consists of two components, weak area affinity and pixel thinning, to maintain the global and local information of 3D point cloud, but this is only support small scenes with high confidence. Point cloud registration often assumes strong structural consistency, but actual scenes may not have such ideal similarity. Most cross-source point clouds have different densities, noise models and part missing, making it difficult to obtain satisfactory results through conventional point cloud registration methods.

As shown in Fig. 1, it can be seen that cross-source point clouds often contain missing points, which can present a significant challenge for direct registration. Variations in point density between point clouds can complicate registration, especially when the point clouds are incomplete. To overcome these challenges, the fusion of heterogeneous point clouds can be divided into two stages, resulting in increased robustness and precision. In the first stage, the point clouds are divided into small point cloud patches in advance to extract effective information, and then the matching candidate sets are obtained by using the global descriptor and spatial distribution of the point cloud patches. The second stage involves two steps of registration using

\*Corresponding author (bushuhui@nwpu.edu.cn)

<sup>1</sup>School of Aeronautics, Northwestern Polytechnical University

<sup>2</sup>PLA Strategic Support Force Information Engineering University

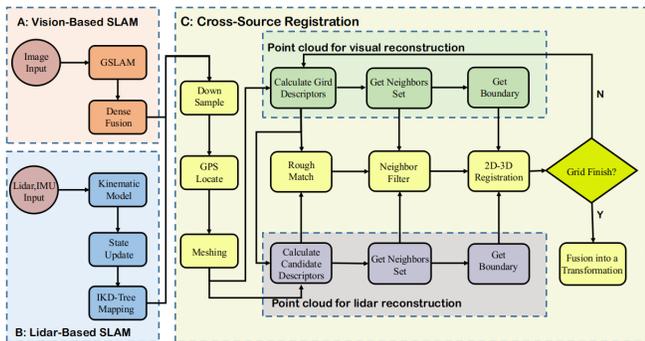


Fig. 2. System overview of Hybrid Fusion

the Normal Distribution Transform (NDT) algorithm. Firstly, the boundary points in patches projected onto the 2D plane are used for registration, and then 3D fine registration is carried out on this basis, the final pose is obtained by fusing the multiple patch pose. This approach is particularly useful when dealing with missing points in different regions, which cause ICP-based registration to fail. In summary, contributions of this work are summarized as follows:

- A multi-view cross-source point cloud fusion algorithm is proposed, which can quickly build a complete 3D point cloud.
- In order to overcome the difference in the spatial distribution of point clouds from different viewings, a point cloud registration method based on the joint optimization of global descriptors and 2D boundary is proposed. This method can effectively fuse cross-source point clouds.
- A cross-view fusion dataset<sup>1</sup>, including multiple simulation and real-world datasets, is established, which is composed of over-view and street-view data.

## II. RELATED WORKS

Our approach relies on three research fields: monocular camera dense reconstruction, LiDAR mapping and point cloud registration. Some related works are summarized in follows.

### A. Monocular Dense Reconstruction

Monocular dense reconstruction from image sequences has been well-studied research field in recent years. The dense matching algorithm used by SfM also has extensive updates. The earliest algorithms, such as the least square matching algorithm [5], [6], have evolved into more advanced methods like patch match stereo (MVS) method. Recent advancements include the semi global matching (SGM) [7] and patch based MVS [8]. To solve the depth map reconstruction stereo matching problem, typically, matching cost calculation and aggregation, as well as disparity calculation and refinement, are performed using local, global, or semi-global energy minimization techniques.

<sup>1</sup>github.com/npupilab/cvf-dataset

Another important research interest is real-time dense reconstruction. Bloesch *et al.* proposed CodeSLAM [9], which can obtain a compact implicit representation of a dense depth map using only a few parameters. Chen *et al.* proposed DenseFusion [1], which completed 3D reconstruction by building virtual stereo pairs from selected key frames to generate dense 3D point clouds after trimming.

### B. Lidar-based SLAM

The most important method of LiDAR-based simultaneous localization and mapping (SLAM) system is rely on sparse features [10], [11], such as line features and surface features. For the sparse features, most LiDAR-base SLAM developed from LiDAR odometry and mapping (LOAM) [10]. It is composed of three main modules: feature extraction, odometry, and mapping. To reduce the computation load, local smoothness is used to extract features from each frame of LiDAR scanning. The odometry module matches the extracted features from two consecutive frames to obtain a rough but real-time LiDAR odometry. Finally, the 3D points from the LiDAR scan are registered to the global map. Subsequent LiDAR-based SLAM works keep a framework similar to LOAM. For example, Lego-LOAM [12] introduces ground point segmentation to reduce the computation and use a closed-loop module to reduce scene drift. LIOM [13] proposes a tightly-coupled LiDAR inertial fusion method where the IMU pre-integrations are added into the odometry. FAST-LIO2 [14] uses tightly coupled iterative Kalman filtering to directly register the original LiDAR points to the map without extracting features. To maintain the map, it employs an incremental k-d tree structure called the Incremental Kalman-filtered k-d tree (IKD-tree).

### C. Point Cloud Registration

Our method mainly involves two types of point cloud registration: feature-based and NDT-based methods. Feature-based methods extract the feature from point cloud, and transforms the point cloud registration 3D place into feature space. Feature-based registration has two interests, the former extract features from the fields near a feature point of the point cloud, such as binary shape context (BSC) [15], Fast Point Feature Histograms (FPFH) [16], and rotational projection statistics (RoPS) [17]. However, these methods often encounter difficulties in cross-source point cloud registration due to differences in point cloud density and scanning range across different sources. The latter is a feature extraction method applied to the entire point cloud, describing its the overall characteristics, such as Ensemble of Shape Functions (ESF) [18], Globally Aligned Spatial Distribution (GASD) [19], and so on. The global descriptor is better than those feature-based points for cross-source point cloud registration, which pays more attention to the features of the whole point cloud rather than the local region.

The Normal Distribution Transform (NDT) method compares the normal distribution of points in a voxel grid between two point clouds to find the best match. It uses standard optimization techniques to minimize the difference

between the distributions. Das *et al.* [20] proposed a multi-scale k-means normal distribution transformation (MSKM-NDT), which addresses the issue of cost discontinuity in point cloud registration by dividing the data into clusters using k-means clustering and optimizing them at various scales. Lu *et al.* [21] proposed a nondestructive testing algorithm of variable size voxels, which improved the accuracy of the algorithm. NDT registration is to register point clouds by comparing the similarity of grids and nonlinear optimization, so it has certain robustness to cross-source point clouds. In the subsequent algorithm of this paper, the NDT algorithm is used for registration.

### III. METHODOLOGY

The overall processing pipeline is shown in the Fig. 2. It takes a sequence of images, LiDAR scans and IMU data as input, while high quality fused point clouds are generated.

The *module-A* takes charge of visual dense reconstruction. In this module, UAV is used to take aerial photos, whose position is estimated with the SLAM plug-in using the GSLAM framework. At the same time, in order to fix the scale of the visual reconstruction point cloud, the GNSS information is integrated into the optimization process. Finally, the DenseFusion method [1] is adopted to densify the point cloud.

The *module-B* adopts FAST-LIO2 algorithm [14] to generate dense point cloud. It uses an efficient tightly-coupled iterated Kalman filter, and directly registers the raw points to the map without extracting features. Then, IKD-tree is used to dynamically maintain the map, which has better search and access speed than traditional methods.

The *module-C* performs cross-source point cloud registration. Firstly, mapping LiDAR-based clouds into vision-based ones using GNSS position information. In addition, the approximate range obtained by the visual SLAM can also be used to replace GNSS information. Due to the presence of large GNSS error, only a coarse registration range can be determined. Point clouds are divided into small patches, and candidate patches are selected based on similar information. However, this alone cannot solve the problem of cross-source point cloud registration. To overcome limitations, a novel cross-view point cloud registration method is proposed. This method uses the boundary points of 2D projection for rough registration, and then performs 3D fine registration. The registered position of all patches is fused into a final position, resulting in precise cross-source point cloud registration with different viewing angles.

#### A. Dense Reconstruction

Real-time monocular dense reconstruction is the initial phase. It consists of two parts: the pose estimation algorithm and dense reconstruction.

1) *Pose Estimation Algorithm*: For each frame of image, SiftGPU [22] is used to extract feature points and corresponding descriptors. In order to acquire the georeferenced position, the drone's pose estimation integrates both visual and GNSS-based positions. The Bundle Adjustment (BA)

process and optimize two loss functions: re-projection error  $e_r$  and GNSS error  $e_g$ . The overall error function is defined as follows:

$$\min f(x) = \frac{1}{2} \left[ \sum_{i=0}^n e_{r_i}^\top W_r e_{r_i} + e_g^\top W_g e_g \right], \quad (1)$$

where  $e_{r_i}$  is re-projection error for 3D points,  $e_g$  represents the scale error of both visual and GNSS,  $W_r$ ,  $W_g$  are the information matrix of re-projection error and GNSS error.  $e_r$  is defined as follows:

$$e_r = x'_i - \frac{1}{s_i} \mathbf{K} \exp(\xi^\wedge) \mathbf{P}_i, \quad (2)$$

where  $x'_i$  stands for the position of corresponding key point in current frame.  $\mathbf{K} \in \mathbb{R}^{3 \times 3}$  is the camera internal parameter matrix,  $\xi \in \mathfrak{se}(3)$  represents the camera position and pose represented by Lie algebra, and  $\mathbf{P}_i \in \mathbb{R}^3$  represents the map point position.

GNSS Error  $e_g$  is defined as follows:

$$e_g = t_{GNSS} - t_{SLAM}, \quad (3)$$

where  $t_{GNSS}$  represents the displacement vector of GNSS between two frames,  $t_{SLAM}$  represents the displacement vector of SLAM.

2) *Densification*: To ensure point cloud accuracy, dense reconstruction uses larger baseline images derived from adjacent time frames as stereo image pairs. Bouguet's rectification aligns stereo images with parallel optical axes, and CUDA-accelerated SGM [7] performs efficient stereo matching.

To minimize depth inaccuracies, we utilized a combination of depth consistency verification and photometric consistency verification.

#### B. Lidar-based SLAM

In order to obtain high quality LiDAR point clouds, FAST-LIO2 [14] is adopted, which primarily consists of two modules: the motion model of tightly coupled Kalman filter and the mapping module.

1) *Motion Model*: The Kalman filter-based motion model is employed for efficient processing, which performs forward propagation of IMU measurements and iterative updating of each LiDAR scan. The forward propagation estimates the current LiDAR position by integrating the IMU measurements when they arrive. Specifically, it transfers the current LiDAR state and covariance when the process noise is small. Based on the IMU forward propagation estimation, the scanned point cloud is projected into the world coordinate system. The point cloud is then aligned with the map points using nonlinear optimization to minimize errors, resulting in a more accurate pose estimation.

2) *Mapping*: The novel data structure IKD-tree is used to organize the map points for large scale range mapping. To control the size of the map, the system saves a specific length of area around the current location. When the LiDAR scanning area touches the map boundary, the map moves in the direction of the boundary and deletes map points in the



Fig. 3. GNSS Registration, red one is LiDAR-based point cloud, others are the vision-based point cloud. There is an obvious mismatch in the cylindrical tower.

opposite direction. All of these operations are based on the IKD-tree data structure.

### C. Cross-Source Registration

After generated two point clouds from dense reconstruction and LiDAR-based SLAM, a cross-source point cloud registration is required to generated a high quality point cloud. Our algorithm can register point clouds from different viewing angles. The detailed processing flow is described in following sections.

1) *Pre-processing*: Appropriate pre-processing, which is separated into initial adjustment and point clouds meshing, is required for dense point clouds. The initial adjustment of point cloud mainly includes: point cloud down sampling and coarse alignment. The first step is the point cloud down sampling. The VoxelGrid filter [23] is used for down sampling the point cloud, which reduces its density without losing details. The second step is coarse alignment which use GNSS constraints. However, the accuracy of GNSS is often limited, leading to inaccurate alignment, as shown in Fig. 3. Typically, only the approximate positioning range can be determined.

Point cloud gridding: For cross-source point clouds, only a few sub-regions are similar due to significant differences in regional feature distributions. To address this issue, the point cloud is segmented into smaller patches for individual registration. Then, the overall registration pose of the point cloud can be estimated based on the registered patches. To partition the point cloud into patches, the point cloud area is divided into a grid using a specified step length. The step length is set to one tenth of the scene size based on experience. Points are then assigned to different grids based on their positions.

We define the visual dense reconstructed point cloud as  $\mathbf{G} = \{g_1, g_2, g_3 \cdots g_n\}$  and the LiDAR reconstructed point cloud as  $\mathbf{L} = \{l_1, l_2, l_3 \cdots l_m\}$ . The point cloud patches in  $\mathbf{G}$  are defined as  $\mathbf{P}_{\mathbf{G}} = \{p_{g_1}, p_{g_2}, p_{g_3} \cdots p_{g_n}\}$ . In the same manner, point cloud  $\mathbf{L}$  is also divided into patches as  $\mathbf{P}_{\mathbf{L}} = \{p_{l_1}, p_{l_2}, p_{l_3} \cdots p_{l_w}\}$ .

2) *Coarse Match*: Even though point clouds  $\mathbf{G}$  and  $\mathbf{L}$  lack complete surface data, there still exist some overlapping regions that can be utilized to approximately match their

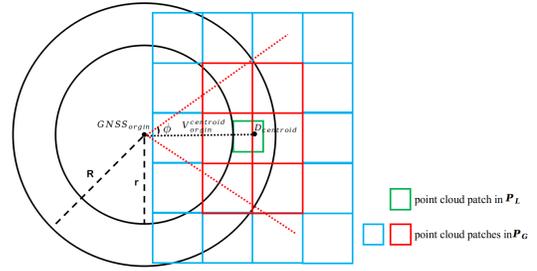


Fig. 4. 2D demonstration of point cloud selection. Green patch is the point cloud patch in  $\mathbf{P}_{\mathbf{L}}$ , blue red patches are in  $\mathbf{P}_{\mathbf{G}}$ , red patches mean the candidate patches we select, while blue indicates it is not a candidate.  $\phi$  is the limited angle for selecting the candidates.

patches. The proposed method divides this into two steps: selecting point cloud candidates and filtering them.

a) *Candidate Selecting*: The candidate range of the patches has a great impact on the subsequent matching accuracy, therefore it is necessary to select a good candidate range of patches.

The first step is to identify significant features in the LiDAR point cloud patches. The selected patches, denoted as  $\mathbf{P}_{\mathbf{L}}^{slt}$ , are chosen from  $\mathbf{P}_{\mathbf{L}}$  based on the criteria of having a point count in a patch greater than  $\eta$ . This threshold is set as insufficient patches do not possess adequate discriminate.

The patches in  $\mathbf{G}$  are selected based on the position of  $\mathbf{P}_{\mathbf{L}}^{slt}$ , as shown in Fig. 4. The origin of  $\mathbf{L}$  is mapped to  $\mathbf{G}$  through GNSS, resulting in  $\mathbf{GNSS}_{origin}$ . For each  $p_{l_i}^{slt}$  in  $\mathbf{P}_{\mathbf{L}}^{slt}$ , its centroid is defined as  $\mathbf{D}_{centroid}$  and the Euclidean distance  $E$  between  $\mathbf{GNSS}_{origin}$  and  $\mathbf{D}_{centroid}$  in the  $X-Y$  plane is calculated. To determine the candidate point cloud patches range, a circular ring on the  $X-Y$  plane is defined with a center at  $\mathbf{GNSS}_{origin}$ . The radii of the ring are defined as:

$$R = (1 + \lambda)E, \quad (4)$$

$$r = (1 - \lambda)E, \quad (5)$$

where  $\lambda$  is the scaling factor for determining the radius range.  $C_r^R$  denotes the circle ring with inner radius of  $r$  and outer radius of  $R$ . To find candidate patches, point cloud patches are selected if any point falls in the  $C_r^R$  and then point cloud patches with  $\theta$  satisfying  $-\phi \leq \theta \leq \phi$  are reserved, where  $\theta$  is the angle determined by the line connecting the patch centroid and  $\mathbf{GNSS}_{origin}$ , as well as the line  $\mathbf{V}_{origin}^{centroid}$ , as shown in Fig. 4. The candidate point cloud patches are stored into the set  $\mathbf{P}_{\mathbf{G} \rightarrow l_i}^{cdt}$ .

b) *Patch Filtering*: The use of descriptors helps to compare each selected LiDAR point cloud patch ( $p_{l_i}^{slt}$ ) with its corresponding candidate patches ( $\mathbf{P}_{\mathbf{G} \rightarrow l_i}^{cdt}$ ). The candidate patches with low similarity are removed to improve the precision of future registration. The ESF640 descriptor is used, but due to its high dimension, ordinary descriptor similarity calculation, such as Euclid distance is computation consuming. Hence, the Pearson correlation coefficient is used

instead.

$$\rho(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}. \quad (6)$$

Here,  $x$  and  $y$  represent descriptors.  $\bar{x}$  and  $\bar{y}$  respectively represent the mean values of  $x$  and  $y$ . It can be seen that the calculation can be optimized by pre-calculating the  $(x_i - \bar{x})$  and  $(y_i - \bar{y})$ , so that the calculation will be greatly accelerated. If a patch  $p_{g_j \rightarrow l_i}^{cdt}$  in  $\mathbf{P}_{\mathbf{G} \rightarrow l_i}^{cdt}$  and its corresponding patch  $p_{l_i}^{slt}$  do not satisfy the condition  $\rho(p_{l_i}^{slt}, p_{g_j \rightarrow l_i}^{cdt}) > 0.6$ , the two descriptors are regarded as irrelevant and the  $p_{g_j \rightarrow l_i}^{cdt}$  will be removed from  $\mathbf{P}_{\mathbf{G} \rightarrow l_i}^{cdt}$ .

3) *Neighbor Filter*: According to the above preliminary filtering, the LiDAR point cloud patches and its similar visual candidate patches have been obtained, but there are still some patches with poor similarity that need to be eliminated. This requires comparing the similarity among the neighbors of the point cloud patches. First, we need to get the neighbors of each patch in  $\mathbf{P}_{\mathbf{L}}^{slt}$  and its related  $\mathbf{P}_{\mathbf{G} \rightarrow l_i}^{cdt}$  sets. If the number of points in its neighbor point cloud patch is greater than a threshold, it will be selected into corresponding neighbor patch set whose descriptor should be calculated.

For each patch  $p_{l_i}^{slt}$  in  $\mathbf{P}_{\mathbf{L}}^{slt}$ , and its corresponding candidate set  $\mathbf{P}_{\mathbf{G} \rightarrow l_i}^{cdt}$ , compare their descriptors within their neighborhoods. Then record its average descriptor similarity value. Next, if the value is lower than a certain threshold, the patch  $p_{g_j \rightarrow l_i}^{cdt}$  will be removed from the  $\mathbf{P}_{\mathbf{G} \rightarrow l_i}^{cdt}$ .

Remember that matching in the same direction between neighbor sets is not mandatory (e.g. a patch's upper neighbor does not have to match with another patch's upper neighbor). This is because the point cloud segmentation process uses different criteria, which may result in poor matches if forced to match in the same direction. Therefore the matching in the same direction is not mandatory required.

---

#### Algorithm 1 2D-3D Registration

---

**Input:**  $p_{l_i}^{slt}, p_{g_j \rightarrow l_i}^{cdt}$   
**Output:**  $\mathbf{K}$

- 1:  $\mathbf{P}_{\mathbf{L}}^{slt+nbr} = \text{splicing\_neighbors}(p_{l_i}^{slt})$
- 2:  $\mathbf{P}_{\mathbf{G} \rightarrow l_i}^{cdt+nbr} = \text{splicing\_neighbors}(p_{g_j \rightarrow l_i}^{cdt})$
- 3: **for**  $p_{l_i}^{slt+nbr}$  in  $\mathbf{P}_{\mathbf{L}}^{slt+nbr}$  **do**
- 4:   **for**  $p_{g_j \rightarrow l_i}^{cdt+nbr}$  in  $\mathbf{P}_{\mathbf{G} \rightarrow l_i}^{cdt+nbr}$  **do**
- 5:      $p_{l_i}^{bdy} = \text{project2D\_get\_boundary\_points}(p_{l_i}^{slt+nbr})$
- 6:      $p_{g_j \rightarrow l_i}^{bdy} = \text{project2D\_get\_boundary\_points}(p_{g_j \rightarrow l_i}^{cdt+nbr})$
- 7:      $\text{score\_2d} = \text{find\_min\_2D\_score}(p_{l_i}^{bdy}, p_{g_j \rightarrow l_i}^{bdy})$
- 8:     **if**  $\text{score\_2d} \leq \text{min\_score}$  **then**
- 9:        $p_{g_{min} \rightarrow l_i}^{cdt+nbr} = \text{find\_min\_2d\_score\_patch}(\text{score\_2d})$
- 10:        $\mathbf{T}_{2d} = \text{get\_2D\_trans}(\text{score\_2d})$
- 11:        $\text{min\_score} = \text{score\_2d}$
- 12:      $\text{score\_3d} = \text{find\_min\_3D\_score}(p_{l_i}^{slt+nbr}, p_{g_{min} \rightarrow l_i}^{cdt+nbr})$
- 13:      $\mathbf{T}_{3d} = \text{get\_3D\_trans}(\text{score\_3d})$
- 14:     **if**  $\text{NDT\_3D\_convergence}(\text{score\_3d})$  **then**
- 15:        $\mathbf{T}_{patch} = \mathbf{T}_{3d} * \mathbf{T}_{2d}$
- 16:        $\mathbf{K}.\text{push\_back}(\mathbf{T}_{patch})$
- 17: **return**  $\mathbf{K}$

---

4) *2D-3D Registration*: The matched set for each element of  $\mathbf{P}_{\mathbf{L}}^{slt}$  has been roughly found. Point cloud registration will be performed next. The surface of point clouds may be incomplete due to different viewing angles, and ESF descriptor

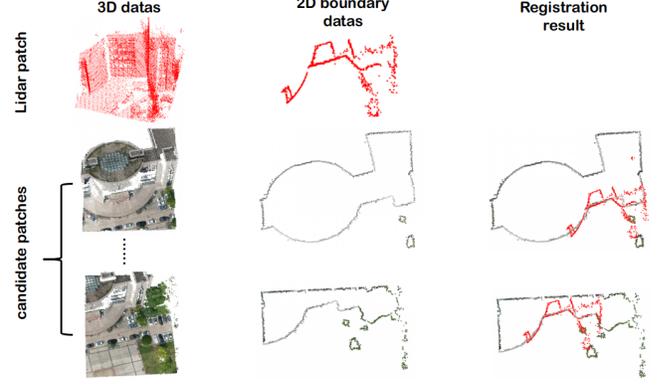


Fig. 5. 2D boundary points registration. LiDAR patch and its candidate patch set are projected onto the  $z$ -plane and boundary points are extracted, then registered on this basis, which can overcome the problems of point cloud missing and inconsistent structure to a great extent.

can only help to some extent, but it is not suitable for high accurate registration. However, for different reconstruction viewing angle, the boundary points of the object are mostly complete, so it can be used to register.

Two large point cloud patch sets for registering are denoted as  $\mathbf{P}_{\mathbf{L}}^{slt+nbr}$  and  $\mathbf{P}_{\mathbf{G} \rightarrow l_i}^{cdt+nbr}$ , which are generated by combining original  $\mathbf{P}_{\mathbf{L}}^{slt}$  and  $\mathbf{P}_{\mathbf{G} \rightarrow l_i}^{cdt}$  with their neighbor sets.

In the following processing, each point cloud patches  $p_{l_i}^{slt+nbr}$  in  $\mathbf{P}_{\mathbf{L}}^{slt+nbr}$  are used for estimating their own transformation. Boundary points are extracted for each  $p_{l_i}^{slt+nbr}$  and its corresponding set  $\mathbf{P}_{\mathbf{G} \rightarrow l_i}^{cdt+nbr}$ . Points with  $d > h$  are extracted from each patch to filter out ground noise, where  $d$  is the distance from the point to the ground plane and  $h$  is the height threshold value. Then, the collected points are projected to the  $X - Y$  plane, and the boundary points are calculated, through the boundary estimation algorithm [23]. Next, the two-dimensional NDT algorithm is used,  $p_{l_i}^{slt+nbr}$  is registered with each element in its corresponding set  $\mathbf{P}_{\mathbf{G} \rightarrow l_i}^{cdt+nbr}$ . An example of 2D registration is shown in the Fig. 5.

Finally, through the cost of the minimum average distance of points the  $p_{l_i}^{slt+nbr}$  with it best matching point cloud patch can be found. In addition, their optimal transformation  $\mathbf{T}_{2d}$  is also estimated.

Based on 2D registration results, a 3D registration is conducted through 3D NDT, which adjusts pose in the  $z$ -axis direction. The registration will be adopted if the matching process can converge, otherwise it will be abandoned.

After 3D registration process, the transformation set  $\mathbf{T}_{3d}$  is obtained, the registration pose of this  $p_{l_i}^{slt+nbr}$  is defined as  $\mathbf{T}_{patch} = \mathbf{T}_{3d} \cdot \mathbf{T}_{2d}$ . If one pair point cloud patches can estimate the transformation  $\mathbf{T}_{patch}$ , it is append to the final results set  $\mathbf{K} = \mathbf{K} \cup \mathbf{T}_{patch}$ . The detail of 2D-3D registration is depicted in Algorithm 1.

5) *Post Processing*: After the transformation set is estimated, an overall optimized transformation is fused from the transformation set  $\mathbf{K}$ . The transformation ( $\mathbf{T}$ ) in  $\mathbf{K}$  is split into  $\mathbf{Q}$  (a quaternion representing rotation) and  $\mathbf{t}$

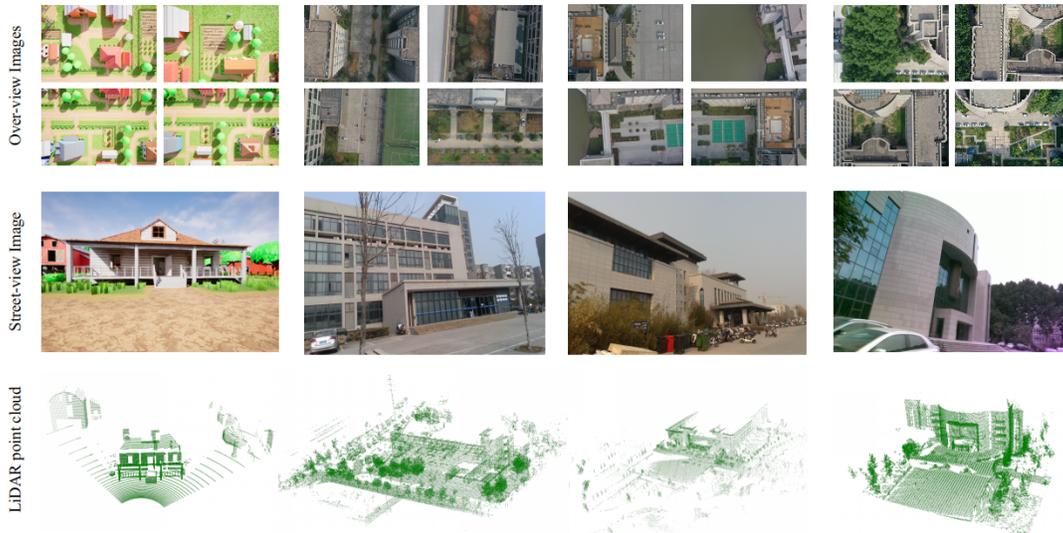


Fig. 6. Cross-source dataset. The first row is the over-view image set taken by UAV. The second row is a street-view image example taken by a UGV. The last row is the point cloud generated by the LiDAR.

(representing translation), which represented in the four-dimensional space  $(x, y, z, \theta)$ , where  $x, y, z$  represent  $t$  and  $\theta$  represents the angle of two quaternions. Because each  $p_{l_i}^{slt+nbr}$  have its own transformation, cluster according to the following requirements:

$$\sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2} < \epsilon, \quad (7)$$

$$fabs\{R[\mathbf{Q}_i, \mathbf{Q}_j]\} < \omega,$$

where,  $R[\mathbf{Q}_i, \mathbf{Q}_j]$  represents the angle between  $\mathbf{Q}_i$  and  $\mathbf{Q}_j$ . If both translation and angle between two different transformation are less than the threshold values  $\epsilon$  and  $\omega$ , they can be grouped into one cluster. After all elements in the  $\mathbf{K}$  are processed, the cluster with the largest number of transformations is extracted and its average value is calculated as the optimal transformation, where the average value of  $\mathbf{Q}$  is calculated by spherical linear interpolation. The registration can be completed by transforming  $\mathbf{L}$  to  $\mathbf{G}$ . Later, minor adjustments can be adjusted through 3D NDT once again.

#### IV. EXPERIMENTS

In order to evaluate the performance of the proposed method, we first establish multiple simulated and real-world datasets, demonstrating the registration results on these datasets. Then, evaluation metrics are proposed to assess the completeness and alignment accuracy of the point clouds. The results show that our system achieves highly robust and accurate cross-source point cloud registration.

The algorithm is implemented based on C++, and all experiments are tested on a desktop PC equipped with Intel i7-8700 CPU and 16 GB RAM. We run our approach in 64-bit Linux system.

#### A. Evaluation Metrics

For evaluating performance of point cloud registration, we propose a novel evaluation metric. First, the result of the point cloud after registration is defined as  $\mathbf{O}$ . Then,  $\mathbf{O}$  and point cloud  $\mathbf{G}$  reconstructed from monocular are described by octree. The evaluation metric of point cloud supplement degree is defined as follows:

$$S_{pt} = \frac{V_{\mathbf{O}} - V_{\mathbf{G}}}{V_{\mathbf{G}}}, \quad (8)$$

where  $S_{pt}$  represents the supplement degree,  $V_{\mathbf{O}}$  and  $V_{\mathbf{G}}$  represent the leaves volume of  $\mathbf{O}$  and  $\mathbf{G}$  after being wrapped by octree.

For the accuracy evaluation of registration results, due to the missing of point clouds, the traditional methods cannot be used. However, the point cloud boundary is usually complete, the accuracy is compared by detecting the nearest distance between the boundary points of the point clouds. In the evaluation, the average value of nearest distance is obtained as the registration accuracy.

#### B. Datasets

To evaluate our entire pipeline, a cross-view fusion dataset is created which contains simulation data and real-world data. These datasets consist of aerial images captured by UAVs, LiDAR data, and ground-based video recordings. In Fig. 6, a qualitative example of a dataset is shown.

Synthetic datasets are composed of over-view image set taken by UAV as well as street-view video and LiDAR data taken by UGV (Unmanned Ground Vehicle). The environment is created using Unreal Engine 4, and Airsim simulation is used to control the UAV, UGV, and various sensors. For the over-view image set, the GNSS information is stored as well as the color image. For street-view scene, the camera and LiDAR are loaded in front of the UGV to capture

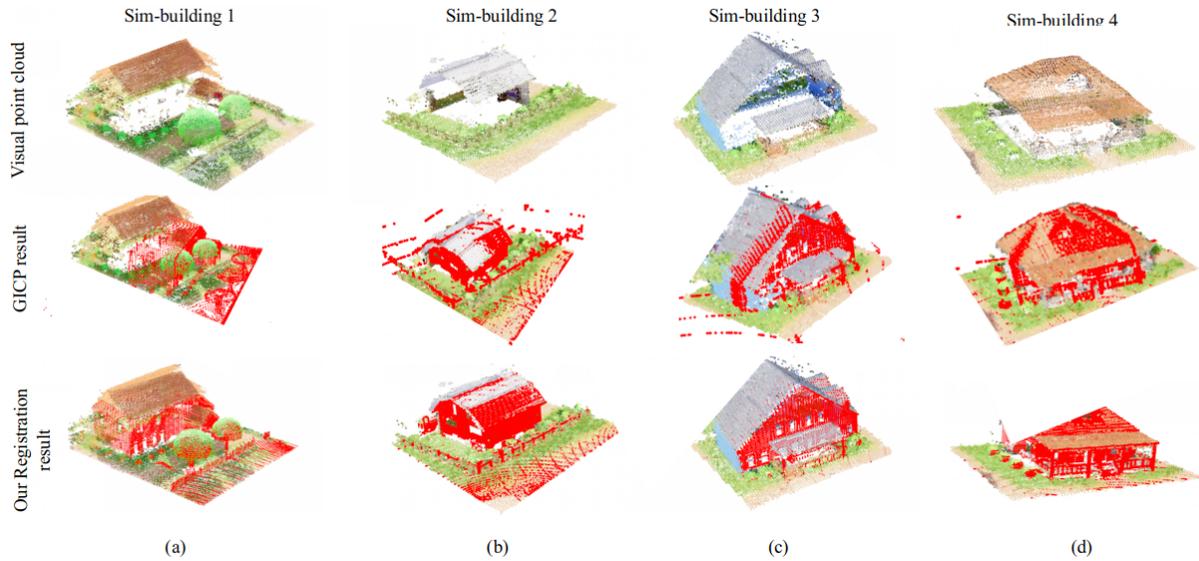


Fig. 7. Simulation experiment results. The first row represents the dense point cloud reconstructed from an overhead view, the second row represents the LiDAR fusion result obtained using the GICP method, and the third row represents the LiDAR fusion result obtained using our method.

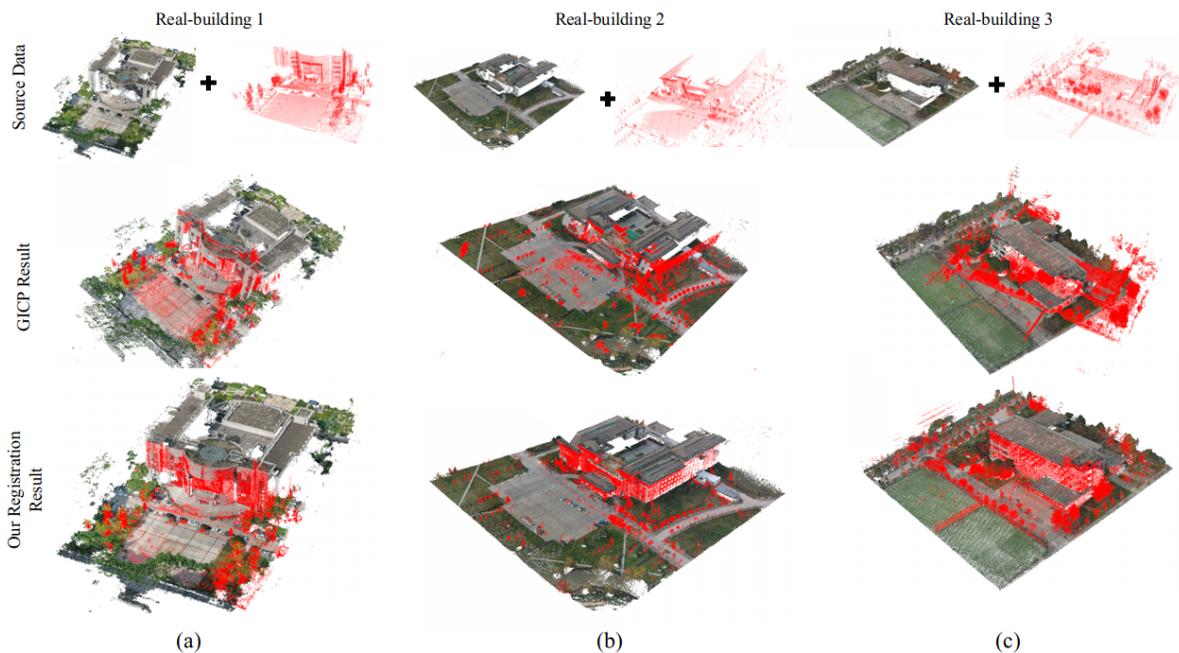


Fig. 8. Real world experiment results. The first row represents the dense point cloud reconstructed from an overhead view and the LiDAR point cloud. The second row represents GICP method result. The third row represents our method result.

TABLE I  
POINT CLOUD RESULT EVALUATION METRICS

Scene	Visual point cloud volume ( $m^3$ )	Result volume ( $m^3$ )	Volume growth ratio	GICP Accuracy ( $m$ )	Our Accuracy ( $m$ )
<i>sim-building1</i>	217.49	287.57	32.21%	2.54	0.41
<i>sim-building2</i>	76.05	92.49	21.60%	1.56	0.29
<i>sim-building3</i>	103.67	140.02	35.05%	1.62	0.41
<i>sim-building4</i>	84.05	107.74	28.18%	1.81	0.37
<i>real-building1</i>	2667.33	3025.60	13.43%	2.06	0.66
<i>real-building2</i>	4432.07	4876.19	10.02%	2.31	0.85
<i>real-building3</i>	2495.45	2896.52	16.08%	4.70	0.72

the environmental information in real time and record the information through the recording function of ROS.

In order to verify the results of the algorithm, we enter the dataset in the actual world. The DJI Mavic 2 Pro UAV is used to cruise the urban area at a certain image repetition rate, take photos at regular intervals, record the GNSS position information when taking photos, and generate an over-view image set. The flight altitude is about 65 meters. For the collection of street-view environment on the ground side, we use RealSense D435 camera, Livox Avia LiDAR and GNSS receiver, which are connected to Jetson Xavier NX for data collection.

### C. Experiment Results and Analysis

As shown in the Figs. 7 and 8, experiments are conducted on both simulation and real datasets. It compares the fusion data obtained through the GICP [24] method and our method. The GICP algorithm is used with default parameters. The figure shows that ICP-based methods are not well-suited for supplementing incomplete point clouds, since they rely on an iterative point-to-point optimization process. Gaps in incomplete point clouds can cause the ICP loss function to rapidly increase, which forces the algorithm to search for matches in the gaps and can result in deviations in global registration. In contrast, our algorithm focuses on matching point cloud patches and prioritizes patches with prominent features for matching. Through an iterative coarse-to-fine registration process, the final matching error is reduced to tens of centimeters. Quantitative results are listed in Table I. It can be seen that the registration accuracy and filling effect are better than traditional methods.

## V. CONCLUSION

In this paper, we propose a framework that can register cross-source point clouds and fill in missing points from various viewing angles. This framework has the following advantages: 1) High accuracy, our method can correctly fit cross-source point clouds, even if there is a large range of point missing. 2) Point cloud filling and fusion: 3D cross-source point clouds reconstructed from a single viewing angle usually contain missing points. Through cross-source point cloud filling from multiple viewing angle, the gap of point clouds can be filled.

Although the proposed method achieves good performance, there are still some limitations. In noisy scenes, the registration effect may not be good enough. At this time, parameters are needed to adjust then get a better registration result.

## REFERENCES

- [1] L. Chen, Y. Zhao, S. Xu, S. Bu, P. Han, and G. Wan, "Densefusion: Large-scale online dense pointcloud and dsm mapping for uavs," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 4766–4773.
- [2] X. Huang, J. Zhang, L. Fan, Q. Wu, and C. Yuan, "A systematic approach for cross-source point cloud registration by preserving macro and micro structures," *IEEE Transactions on Image Processing*, vol. 26, no. 7, pp. 3261–3276, 2017.
- [3] N. Mellado, M. Dellepiane, and R. Scopigno, "Relative scale estimation and 3d registration of multi-modal geometry using growing least squares," *IEEE transactions on visualization and computer graphics*, vol. 22, no. 9, pp. 2160–2173, 2015.
- [4] X. Huang, L. Fan, Q. Wu, J. Zhang, and C. Yuan, "Fast registration for cross-source point clouds by using weak regional affinity and pixel-wise refinement," in *2019 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2019, pp. 1552–1557.
- [5] F. Ackermann, "Digital image correlation: performance and potential application in photogrammetry," *The Photogrammetric Record*, vol. 11, no. 64, pp. 429–439, 1984.
- [6] A. Gruen, "Adaptive least squares correlation: a powerful image matching technique," *South African Journal of Photogrammetry, Remote Sensing and Cartography*, vol. 14, no. 3, pp. 175–187, 1985.
- [7] H. Hirschmuller, "Stereo processing by semiglobal matching and mutual information," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 30, no. 2, pp. 328–341, 2007.
- [8] Y. Furukawa, "Patch-based multi-view stereo software," <http://grail.cs.washington.edu/software/pmvs>, 2010.
- [9] M. Bloesch, J. Czarowski, R. Clark, S. Leutenegger, and A. J. Davison, "Codeslam—learning a compact, optimisable representation for dense visual slam," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2560–2568.
- [10] J. Zhang and S. Singh, "Loam: Lidar odometry and mapping in real-time," in *Robotics: Science and Systems*, vol. 2, no. 9. Berkeley, CA, 2014, pp. 1–9.
- [11] F. Moosmann and C. Stiller, "Velodyne slam," in *2011 IEEE intelligent vehicles symposium (iv)*. IEEE, 2011, pp. 393–398.
- [12] T. Shan and B. Englot, "Lego-loam: Lightweight and ground-optimized lidar odometry and mapping on variable terrain," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 4758–4765.
- [13] H. Ye, Y. Chen, and M. Liu, "Tightly coupled 3d lidar inertial odometry and mapping," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 3144–3150.
- [14] W. Xu, Y. Cai, D. He, J. Lin, and F. Zhang, "Fast-direct lidar-inertial odometry," *IEEE Transactions on Robotics*, 2022.
- [15] Z. Dong, B. Yang, F. Liang, R. Huang, and S. Scherer, "Hierarchical registration of unordered tfs point clouds based on binary shape context descriptor," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 144, pp. 61–79, 2018.
- [16] 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.
- [17] Y. Guo, F. Sohel, M. Bennamoun, M. Lu, and J. Wan, "Rotational projection statistics for 3d local surface description and object recognition," *International journal of computer vision*, vol. 105, no. 1, pp. 63–86, 2013.
- [18] W. Wohlkinger and M. Vincze, "Ensemble of shape functions for 3d object classification," in *2011 IEEE international conference on robotics and biomimetics*. IEEE, 2011, pp. 2987–2992.
- [19] J. P. S. do Monte Lima and V. Teichrieb, "An efficient global point cloud descriptor for object recognition and pose estimation," in *2016 29th SIBGRAPI conference on graphics, patterns and images (SIBGRAPI)*. IEEE, 2016, pp. 56–63.
- [20] A. Das and S. L. Waslander, "Scan registration with multi-scale k-means normal distributions transform," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2012, pp. 2705–2710.
- [21] L. Jun, L. Wei, D. Donglai, and S. Qiang, "Point cloud registration algorithm based on ndt with variable size voxel," in *2015 34th Chinese Control Conference (CCC)*. IEEE, 2015, pp. 3707–3712.
- [22] C. Wu, "Siftgpu: A gpu implementation of scale invariant feature transform (sift)(2007)," [URL http://cs.unc.edu/~ccwu/siftgpu](http://cs.unc.edu/~ccwu/siftgpu), 2011.
- [23] R. B. Rusu and S. Cousins, "3D is here: Point Cloud Library (PCL)," in *IEEE International Conference on Robotics and Automation (ICRA)*, Shanghai, China, May 9-13 2011.
- [24] A. Segal, D. Haehnel, and S. Thrun, "Generalized-icp," in *Robotics: science and systems*, vol. 2, no. 4. Seattle, WA, 2009, p. 435.