



# I2P Registration by Learning the Underlying Alignment Feature Space from Pixel-to-Point Similarities

YUNDA SUN and LIN ZHANG, School of Software Engineering, Tongji University, Shanghai, China

ZHONG WANG, Department of Automation, Shanghai Jiao Tong University, Shanghai, China

YANG CHEN and SHENGJIE ZHAO, School of Software Engineering, Tongji University, Shanghai, China

YICONG ZHOU, Department of Computer and Information Science, University of Macau, Macau, China

Estimating the relative pose between a camera and a LiDAR holds paramount importance in facilitating complex task execution within multi-agent systems. Nonetheless, current methodologies encounter two primary limitations. First, amid the cross-modal feature extraction, they typically employ separate modal branches to extract cross-modal features from images and point clouds. This approach results in the feature spaces of images and point clouds being misaligned, thereby reducing the robustness of establishing correspondences. Second, due to the scale differences between images and point clouds, one-to-many pixel-point correspondences are inevitably encountered, which will mislead the pose optimization. To address these challenges, we propose a framework named Image-to-Point cloud registration by learning the underlying alignment feature space from Pixel-to-Point SIMilarities ( $I2P_{ppsim}$ ). Central to  $I2P_{ppsim}$  is a Shared Feature Alignment Module (SFAM). It is designed under a coarse-to-fine architecture and uses a weight-sharing network to construct an alignment feature space. Benefiting from SFAM,  $I2P_{ppsim}$  can effectively identify the co-view regions between images and point clouds and establish high-reliability 2D-3D correspondences. Moreover, to mitigate the one-to-many correspondence issue, we introduce a similarity maximization strategy termed point-max. This strategy effectively filters out outliers, thereby establishing accurate 2D-3D correspondences. To evaluate the efficacy of our framework, we conduct extensive experiments on KITTI Odometry and Oxford Robotcar. The results corroborate the effectiveness of our framework in improving image-to-point cloud registration. To make our results reproducible, the source codes have been released at <https://cslinzhang.github.io/I2P>.

CCS Concepts: • Computing methodologies → Vision for robotics;

Additional Key Words and Phrases: image-to-point cloud registration, Data association, Cross-modal learning

---

This work was supported in part by the National Natural Science Foundation of China under Grant 62272343; in part by the Shuguang Program of Shanghai Education Development Foundation and Shanghai Municipal Education Commission under Grant 21SG23; and in part by the Fundamental Research Funds for the Central Universities.

Authors' Contact Information: Yunda Sun, School of Software Engineering, Tongji University, Shanghai, China; e-mail: 2110850@tongji.edu.cn; Lin Zhang (corresponding author), School of Software Engineering, Tongji University, Shanghai, China; e-mail: cslinzhang@tongji.edu.cn; Zhong Wang, Department of Automation, Shanghai Jiao Tong University, Shanghai, China; e-mail: cszhongwang@sjtu.edu.cn; Yang Chen, School of Software Engineering, Tongji University, Shanghai, China; e-mail: 2011439@tongji.edu.cn; Shengjie Zhao, School of Software Engineering, Tongji University, Shanghai, China; e-mail: shengjiezhaotongji.edu.cn; Yicong Zhou, Department of Computer and Information Science, University of Macau, Macau, China; e-mail: yicongzhou@um.edu.mo.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM 1551-6865/2024/11-ART388

<https://doi.org/10.1145/3697839>

**ACM Reference format:**

Yunda Sun, Lin Zhang, Zhong Wang, Yang Chen, Shengjie Zhao, and Yicong Zhou. 2024. I2P Registration by Learning the Underlying Alignment Feature Space from Pixel-to-Point Similarities. *ACM Trans. Multimedia Comput. Commun. Appl.* 20, 12, Article 388 (November 2024), 21 pages.

<https://doi.org/10.1145/3697839>

## 1 Introduction

**Image-to-Point cloud (I2P)** registration refers to estimating the relative pose between a **Light Detection And Ranging (LiDAR)** and a camera via their measurements (point clouds and images), where the image and the point cloud are captured from the same scene. This task is widely used in many robotics and computer vision applications, such as Simultaneous Localization and Mapping, robot navigation, and scene understanding [16, 18, 31, 42, 52].

The key to I2P registration is feature matching between images and point clouds. Unlike the widely studied homologous registration (image-to-image registration [4, 5, 47, 51]), **point cloud registration (PCR)** [7, 46, 50]), I2P registration is sporadically explored due to the challenging modality differences between images and point clouds.

As shown in Figure 1, the pipeline of I2P involves feature extraction, feature matching, correspondence establishment, and pose estimation. Previous studies relied on complex cross-modal manual feature designs or time-consuming optimization algorithms [10, 23], overlooking the differences between images and point clouds in feature space, perceptual range, and scale. Therefore, the performance of these studies is unsatisfactory. Specifically, to improve the performance of I2P registration, there are still three challenges to be faced with:

- (1) *Misaligned feature space.* Existing methods utilize separate modal branches to extract cross-modal features from images and point clouds [10], which poses a challenge in feature matching. Specifically, as different modalities of data, images and point clouds have significant differences in data structure and scene information captured. Due to the use of different modal branches, current approaches cannot effectively alleviate such modal differences, but instead lead to feature space misalignment. This misalignment in the feature space reduces the performance of feature matching.
- (2) *Insufficient feature fusion.* The current single-stage feature fusion scheme fails to meet the different requirements of feature receptive fields for image-point cloud feature matching [23, 30]. Generally, cross-modal features with global receptive fields are suitable for detecting co-view regions between images and point clouds, while those with local receptive fields are suitable for predicting pixel-point matches. Therefore, it is of necessity to extract cross-modal features with different receptive fields for co-view region detection and matched pixel-point prediction. Unfortunately, current methods lack this capability.
- (3) *Matching ambiguity.* The one-to-many correspondence problem caused by the scale difference between images and point clouds misleads pose optimization. For example, taking the camera center as the origin, as the perception distance gets farther, for one pixel, there are usually multiple points within the frustum. Moreover, the Euclidean distance among those points may vary greatly. These seemingly “correct” one-to-many correspondences can mislead the pose optimization.

To deal with the aforementioned challenges, we propose a novel I2P registration framework, called **Image-to-Point cloud registration by learning the underlying alignment feature space from Pixel-to-Point SIMilarities (I2P<sub>ppsim</sub>)**. I2P<sub>ppsim</sub> learns the underlying feature alignment space between images and point clouds via a **Shared Feature Alignment Module (SFAM)**, and designs

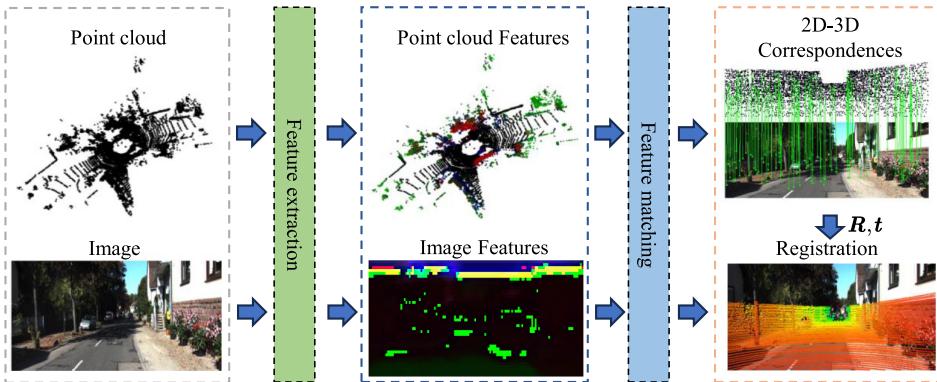


Fig. 1. The general pipeline of image-to-point cloud (I2P) registration. I2P registration first extracts the common features and then conducts the feature matching to build the 2D-3D correspondences. After that, a perspective-n-point (PnP) solver is used to estimate the relative pose.

a matching constraint called point-max based on feature similarity to alleviate the one-to-many correspondence dilemma. The characteristics of I2P<sub>ppsim</sub> and our contributions are as follows:

- (1) The first feature space alignment-based I2P registration framework is proposed, named I2P<sub>ppsim</sub>. It mines pixel-point similarities by learning aligned cross-modal feature spaces. Based on the cross-modal representation, co-view regions are detected and pixel-point correspondences are directly predicted. Extensive experiments demonstrate that our I2P<sub>ppsim</sub> achieves **State-of-the-Art (SOTA)**.
- (2) A novel SFAM is designed. Benefiting from the coarse-to-fine architecture, SFAM can extract coarse-grained features focusing on global expression and fine-grained features focusing on local expression. The former is helpful for co-view region detection and the latter is suitable for pixel-point matching estimation. In addition, SFAM uses a weight-sharing network to construct an aligned cross-modal feature space, which effectively alleviates the modality difference.
- (3) A plug-and-play matching strategy named *point-max* is introduced to solve the one-to-many correspondences. It does not rely on the feature learning ability of the network and aims to identify the best matching point for each pixel in the co-view region. By using point-max, a significant improvement in registration accuracy is achieved. Moreover, point-max can be seamlessly integrated as a plug-and-play module for other I2P registration methods, thereby improving their performance.

## 2 Related Work

### 2.1 Image Registration

Image registration is usually treated as a pre-processing step for applications such as Structure from Motion and image stitching [13, 20, 29]. The key to registration is to establish an accurate image matching in the  $\mathbb{R}^2$  space. Current studies on image matching can be divided into two categories: feature-based ones and matching-based ones. The typical pipeline of the former is to extract the feature descriptors of the image [24], then calculate the distance between those descriptors, and determine the matching feature pairs. Recent approaches expect to obtain better visual feature descriptions through Convolutional Neural Networks [4, 5, 8], and further improve the correct rate and number of matching pairs.

Most feature-based methods determine the matching pairs through the neighbor search among features [9], while the matching-based ones no longer focus on the extraction of image features, but model the image matching as a learning problem. For example, SuperGlue [32] uses attention to aggregate the global and local features, formulates the image matching problem as a graph matching problem, and determines the matching pairs by approximately linear distribution. Furthermore, LoFTR [34] completely abandons the learning of image features and directly predicts dense pixel-by-pixel matches in an end-to-end manner.

## 2.2 PCR

PCR methods can be broadly categorized into two groups. The first category is characterized by its emphasis on the extraction of point cloud features, aiming to establish correspondences mainly based on feature matching [1, 15, 28, 33, 36, 37, 38, 39, 40, 44, 49]. These methods primarily rely on the feature extraction capabilities of neural networks. Early studies mainly use PointNet [3] to extract point clouds' global descriptors and optimize the network by minimizing distances between global descriptors [1, 33]. Recently, some approaches replaced PointNet with transformer [35]. Benefiting from the expanded receptive fields and enhanced contextual association capabilities brought by the transformer, these methods achieved impressive performance [28, 36].

The above-mentioned PCR methods are sensitive to noise, and another category of methods introduces additional geometric or optimization constraints to enhance the robustness of PCR [2, 6, 11, 41, 45, 48]. PointDSC [2] introduces spatial consistency to eliminate inaccurate matching pairs. RGM [11] uses deep map matching to implement PCR. MAC [48] searches for the maximum clique subsets among the matching pairs and selects the optimal transformation guided by the reprojection error.

## 2.3 I2P Registration

Compared with image registration and PCR, there are few studies on I2P registration. According to the ways of correspondence establishment, these I2P methods are mainly categorized into two classes: keypoint-based methods and keypoint-free ones. The core idea of the keypoint-based methods is to measure the distance and establish correspondences based on the keypoint descriptors extracted from images and point clouds [10].

In order to avoid complex cross-modal keypoint design, keypoint-free methods aim to learn point/pixel-wise features with strong repetitiveness [17, 19, 23, 30, 43]. DeepI2P [23] utilizes cross-attention to fuse the features of images and point clouds and proposes inverse camera projection for relative pose estimation. Building upon DeepI2P, CorrI2P establishes 2D-3D correspondences based on feature similarity metrics [30]. Similar to CorrI2P, EP2P-Loc achieves visual localization using images and point cloud submaps as inputs [19].

Although keypoint-free methods improve registration performance, they overlook the one-to-many correspondence dilemma caused by scale ambiguity. Furthermore, these methods extract features from images and point clouds through different modal branches, indicating that the feature spaces of images and point clouds are not aligned. The misalignment of feature spaces further reduces the repeatability of cross-modal features.

## 3 Methodology

### 3.1 Problem Definition and Framework Overview

Given an image  $I \in \mathbb{R}^{3 \times W \times H}$  ( $W$  and  $H$  represent the width and height of the image) and a point cloud  $\mathcal{P} = \{P_1, P_2, \dots, P_N \in \mathbb{R}^3\}$  ( $N$  is the number of points), the task of I2P registration is to estimate the relative rigid transformation  $T = [R|t] \in SE(3)$  ( $R \in SO(3)$ ,  $t \in \mathbb{R}^3$ ) from the LiDAR frame to

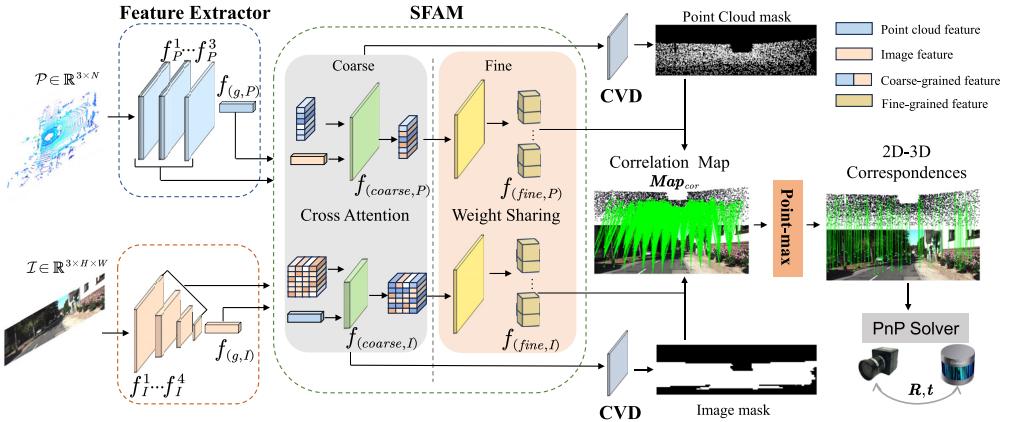


Fig. 2. Overview of I2P\_ppsim. I2P\_ppsim is composed of five modules: Feature Extractor, SFAM, CVD, point-max, and PnP solver. Firstly, the high-dimensional features of images and point clouds are separately extracted by the two branches of Feature Extractor. Such features are then fed into SFAM to obtain cross-modal features in a coarse-to-fine architecture. With the coarse-grained features from SFAM, the co-view region of the image and point cloud is determined by CVD. Further, based on the fine-grained features from SFAM and the outputs of CVD, a correlation map is predicted, where 2D-3D correspondences are selected by point-max. At last, the relative pose of the LiDAR to the camera is estimated via the PnP solver. CVD, co-view detector; PnP, perspective-n-point.

the camera frame. Generally, a standard registration problem is modeled as a **Perspective-n-Point (PnP)** or Iterative Closest Point problem. However, the point cloud collected by LiDAR has little geometric and appearance similarity with the RGB image. Also, to establish the correspondences among pixels and points is non-trivial. We expect to represent the two kinds of data in a higher-dimensional feature space through information fusion, build 2D-3D correspondences, and then regard I2P registration as a PnP problem. To this end, I2P\_ppsim is designed to comprise two parts: a correlation map estimation module and a pose estimator (as shown in Figure 2). In I2P\_ppsim, Feature Extractor, SFAM, and **Co-View Detector (CVD)** are used to estimate the correlation map. When performing inference, given a pair of  $\mathcal{I}$  and  $\mathcal{P}$ , the image and point cloud are first mapped to the high-dimensional space by Feature Extractor, and then SFAM is used to perform feature fusion and feature space alignment. Subsequently, the pixel-point correlation map is calculated based on the aligned cross-modal features in the co-view region and is further fed into the pose estimator. Based on our point-max strategy, we select matching candidates obtained from the correlation map. In this way, the 2D-3D correspondences can be established. Finally, **Efficient Perspective-n-Point (EPnP)** [21] and RANSAC are employed to iteratively optimize the pose.

### 3.2 Feature Extractor

In view of the inherent dissimilarities in the properties of images and point clouds, employing the same feature extraction network to process both of them is impractical. Inspired by DeepI2P, we resort to ResNet [14] to encode the image features, while a modified PointNet++ [22, 27] serves as the feature encoder for the point cloud. Through these feature encoders, the multi-scale features of images and point clouds can be obtained, expressed as  $F_I^i \in \mathbb{R}^{c_i \times W_i \times H_i}$ ,  $i \in \{1, 2, 3, 4\}$ , and  $F_P^j \in \mathbb{R}^{c_j \times N_j}$ ,  $j \in \{1, 2, 3\}$ , respectively, where  $c_i/c_j$  denotes the  $i$ th/ $j$ th features. Then, the global descriptions of the scene from images and point clouds can be obtained by performing max pooling, which are denoted by  $F_{(g,I)} \in \mathbb{R}^{c_I}$  and  $F_{(g,P)} \in \mathbb{R}^{c_P}$ , where  $c_I/c_P$  is the dimension of image/point

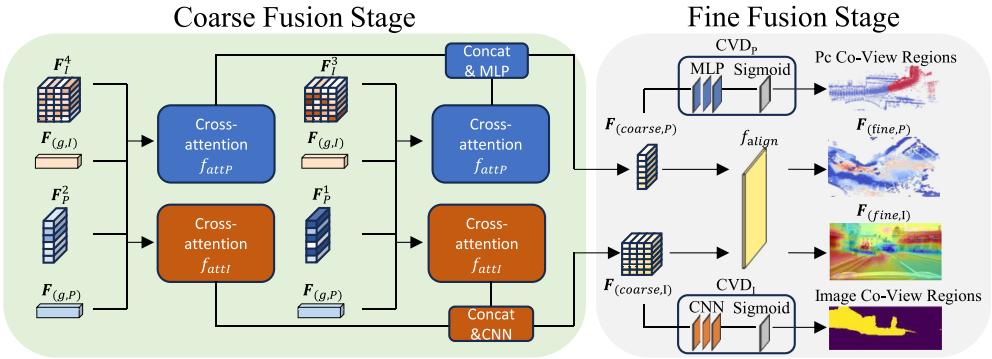


Fig. 3. The overall architecture of SFAM. SFAM fuses image and point cloud features in a coarse-to-fine manner. In the coarse fusion stage, symmetric cross-attention is leveraged to fuse multi-scale image and point cloud features. In the fine-grained fusion stage, fine-grained features  $F_{(fine,I)}$  and  $F_{(fine,P)}$  are extracted by a weight-sharing network  $f_{align}$ . In addition, coarse-grained features  $F_{(coarse,I)}$  and  $F_{(coarse,P)}$  are used as inputs of the CVD, where  $CVD_I$  is the image CVD and  $CVD_P$  is the point cloud CVD.

cloud global feature vector. We posit that taking into account both local and global features is more conducive to enhancing the ability of feature representation. Consequently, images and point clouds of various scales are employed for subsequent feature fusion.

### 3.3 SFAM

SFAM utilizes multi-scale features  $F_I^i$  and  $F_P^j$  as inputs and generates per-pixel and per-point cross-modal features  $F_{(fine,I)}$  and  $F_{(fine,P)}$ . Considering that forcibly aligning the visual-laser feature spaces is rude and meaningless, SFAM is designed as a coarse-to-fine architecture, employing a two-stage approach to extract common features, as illustrated in Figure 3.

In the coarse fusion stage, we concatenate the global features of the image ( $F_{(g,I)}$ ) and point cloud ( $F_{(g,P)}$ ) with multi-scale point cloud features ( $\{F_P^j\}_{j=1,2}$ ) and image features ( $\{F_I^i\}_{i=3,4}$ ), respectively. Next, fused features at different scales are obtained via symmetrical cross-attention ( $f_{attI} : \mathbb{R}^{c \times W \times H} \rightarrow \mathbb{R}^{c \times W \times H}$  and  $f_{attP} : \mathbb{R}^{c \times N} \rightarrow \mathbb{R}^{c \times N}$ ),

$$\begin{aligned} f'_{(coarse,I)} &= f_{attI}(F_{(g,P)}, F_I^4, F_P^2, F_{(g,I)}), \\ f'_{(coarse,P)} &= f_{attP}(F_{(g,P)}, F_I^4, F_P^2, F_{(g,I)}), \end{aligned} \quad (1)$$

where  $f'_{(coarse,I)}$  and  $f'_{(coarse,P)}$  are the fused features of the image and the point cloud, respectively. Similarly, by replacing  $F_I^4$  and  $F_P^2$  in Equation (2) with  $F_I^3$  and  $F_P^1$ , fused features of another scale can be obtained, which are indicated as  $f''_{(coarse,I)}$  and  $f''_{(coarse,P)}$ . Then the fused features from different scales are concatenated, and the feature encoding functions are used to further extract coarse-grained features  $f_{(coarse,I)}$  and  $f_{(coarse,P)}$ ,

$$\begin{aligned} f_{(coarse,I)} &= \text{CNN}(f'_{(coarse,I)}, f''_{(coarse,I)}), \\ f_{(coarse,P)} &= \text{MLP}(f'_{(coarse,P)}, f''_{(coarse,P)}). \end{aligned} \quad (2)$$

Coarse-grained features focus more on global information and provide a broader perspective of the correspondence between the two modalities. In I2P<sub>ppsim</sub>, they are regarded as the input of CVDs.

In the fine fusion stage, an aligned feature space is constructed to mine the consistent features of both images and point clouds. Considering that CNN is difficult to handle unordered point cloud features, while unordered network structures still have the ability to handle ordered features, we

resort to a variant of pointnet [23] to construct a weight-sharing network to model the feature space. With  $f_{(coarse,I)}$  and  $f_{(coarse,P)}$  as the network inputs, the fine-grained features  $f_{(fine,I)}$  and  $f_{(fine,P)}$  are obtained by

$$\begin{aligned} f_{(fine,P)} &= f_{align}(f_{(coarse,P)}), \\ f_{(fine,I)} &= f_{align}(f_{(coarse,I)}), \end{aligned} \quad (3)$$

where  $f_{align} : \mathbb{R}^{c \times (W \times H)} \rightarrow \mathbb{R}^{c \times (W \times H)}$  for images and  $\mathbb{R}^{c \times N} \rightarrow \mathbb{R}^{c \times N}$  for point clouds. These features focus on local similarity, enabling SFAM to capture delicate correspondences among pixels and points. After acquiring the fine-grained features, we expect to measure the similarity among them to build correspondences. With  $f_{(fine,I)}$ ,  $f_{(fine,P)}$ , we can calculate the correlation map  $\mathbf{Map}_{cor} \in \mathbb{R}^{(W \times H) \times N}$  between the image and the point cloud by

$$\mathbf{Map}_{cor} = f_{(fine,I)}^T f_{(fine,P)}. \quad (4)$$

$\mathbf{Map}_{cor}$  reflects the similarity of visual-laser data, which enables I2P<sub>ppsim</sub> to learn the feature matching process. We optimize SFAM by minimizing the similarity loss of  $\mathbf{Map}_{cor}$ . The specific loss function design and analysis will be presented in Section 3.6.

*The Motivation behind the Weight-Sharing Network.* Existing feature fusion scheme for I2P registration extracts cross-modal features of images and point clouds through different modal branches. This results in these cross-modal descriptions being in different feature spaces, which in turn hinders the prediction of 2D-3D correspondences. To address this issue, we design a weight-sharing network in SFAM. By utilizing this network, features from different modal branches are mapped to the same aligned feature space, which assists feature matching.

*The Advantages of SFAM.* Compared to other methods, our SFAM takes into account both global and local information, which enables I2P<sub>ppsim</sub> to better focus on semantic objects in the scene (such as cars, houses). Moreover, while other methods use separate modal branches to extract cross-modal features from images and point clouds, SFAM employs a weight-sharing network to directly construct an aligned feature space. As a result, the cross-modal features of images and point clouds are mapped to the same feature space, which enhances the reliability of feature similarity. Figure 4 illustrates the similarities of cross-modal features extracted by different methods. It can be seen that our I2P<sub>ppsim</sub> outperforms the others.

### 3.4 CVD

The number of pixels or points in an image or point cloud typically ranges from thousands to tens of thousands. In the case where the number of pixels is  $M$  and the number of points is  $N$ , the size of the correlation map would be  $M \times N$ . The computation and storage requirements for such a large-scale correlation map are substantial. However, in I2P registration, the associated data are typically concentrated in a fan-shaped area, occupying only a small portion of the image and point cloud. Leveraging this characteristic, CVDs for images and point clouds are designed to determine whether a pixel or point belongs to the co-view region, they are denoted by  $CVD_I$  and  $CVD_P$ , respectively.

During the calculation of the correlation map, only the pixels and points in the co-view region are considered. In this way, the data scale of the correlation map is significantly reduced, leading to accelerated network inference. To accomplish co-view region detection, we treat it as a binary classification problem. We employ two classification heads to analyze the coarse-grained features  $f_{(coarse,I)}$  and  $f_{(coarse,P)}$  separately. The outputs of the CVD correspond to the co-view scores for each pixel or point, indicating their likelihood of belonging to the co-view region.

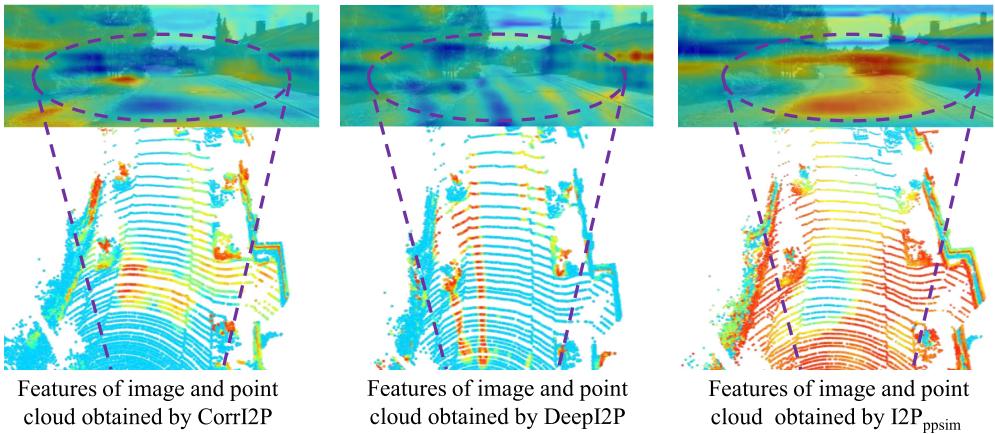


Fig. 4. Comparison of cross-modal features extracted by different feature fusion schemes. The co-view regions are marked, the warmer the feature, the higher the probability of a match.

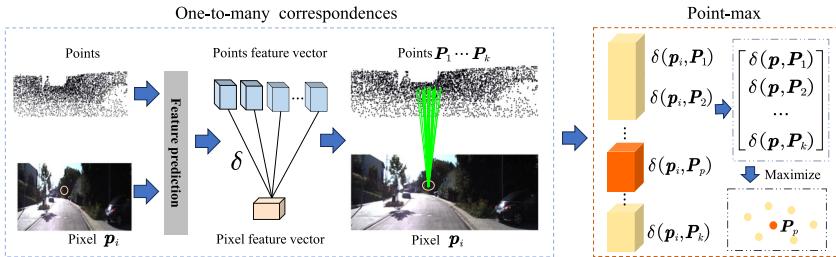


Fig. 5. The phenomenon of one-to-many correspondences between a pixel  $p_i$  and multiple points  $P_1, \dots, P_k$ . Point-max is designed to eliminate those outliers.

### 3.5 Pose Estimation

Following the network prediction, notable dissimilarities are expected to exist between the matched and unmatched pixel-point pairs within the correlation map  $\text{Map}_{cor}$ . With this in mind, one intuitive approach is to employ threshold screening to determine the correspondences between pixels and points. Alternatively, some methods establish 2D-3D associations by employing a classification-based approach. These methods aim to achieve the closest possible one-to-one correspondence between points and pixels. However, in practice, due to the different sensor measurement manners, the matching between point and pixel is not strictly one-to-one correspondence. As shown in Figure 5, in the reprojection of the point cloud, each pixel in the image corresponds to a frustum in the real world, and multiple points are distributed in the frustum. So, one pixel in the image often has high similarity with multiple points. Furthermore, as the image downscaling and visual depth increase, the adverse impact caused by one-to-many correspondences will become more serious. Those seemingly “correct” correspondences indeed do harm to the calculation of reprojection errors in pose optimization. For example, four adjacent pixels  $[u_i, v_i]^T, i \in \{1, 2, 3, 4\}$  are downsampled to one pixel  $[u_s, v_s]^T$ , which corresponds to  $k$  points  $P_1, \dots, P_k$ , where  $P_j = [x_j, y_j, z_j]^T, j = 1 \sim k$ . Before downscaling, for each pixel  $[u_i, v_i]^T$  and its matching points  $\{P_c = [x_c, y_c, z_c]^T : c = 1 \sim k_i\}$

(where  $k_i$  is the number of matching points for  $[u_i, v_i]^T$ ), the reprojection error is

$$e = \sum_i e_i = \sum_i \sum_c \left\| \begin{bmatrix} u_i \\ v_i \end{bmatrix} - \pi \left( \frac{1}{z_c} K T \bar{P}_c \right) \right\|_2, \quad (5)$$

where  $K$  refers to the intrinsic matrix of the camera and  $\bar{P}_c$  is the homogeneous coordinate of  $P_c$ ,  $\pi(\cdot)$  is an operator that takes the first two dimensions of a vector. After scaling, the reprojection error is reduced to

$$e_s = \sum_{j=1}^k \left\| \begin{bmatrix} u_s \\ v_s \end{bmatrix} - \frac{1}{4} \pi \left( \frac{1}{z_j} K T \bar{P}_j \right) \right\|_2. \quad (6)$$

Such a loss will undoubtedly lead to performance degradation of the network which regards the reprojection error as its supervision. Unfortunately, it is challenging for classification or threshold screening to eliminate those seemingly “correct” points. Besides, these points are indeed observed by the corresponding pixels. Consequently, relying solely on cross-modal learning is also insufficient for removing such correspondences, and it is necessary to design a matching constraint to further refine the matching based on the correlation map.

In image matching, the matching relationship satisfies the following constraint: one pixel in an image can at most match with one pixel in another image [32]. Inspired by that, we appropriately relax this constraint criterion, apply it to I2P registration, and propose the point-max matching constraints: a pixel matches at most one point. Considering the data structure of the correlation map, point-max is performed via maximizing similarity in practice. Specifically, given a pixel  $p$ , its receptive field contains  $k$  candidate matching points  $P_1, \dots, P_k$ . Define the vector correlation operator as  $\delta$ . The fine-grained feature vectors of  $p$  and  $P$  are  $f_{(fine,p)}$  and  $f_{(fine,P)}$ , respectively. Then the correlation between pixel  $p$  and point  $P$  is

$$\delta(p, P) = f_{(fine,p)} \cdot f_{(fine,P)} \quad (7)$$

Thus, the point  $P_p$  matching with the pixel  $p$  is

$$P_p = \arg \max_{P_j} \{ (\delta(p, P_1), \dots, \delta(p, P_k)), j = 1, \dots, k \}. \quad (8)$$

Equation (8) is similar to performing max pooling among the matching candidates and selecting the candidate point with the highest correlation for the pixel. By applying Equation (8) to the correlation map  $Map_{cor} \in \mathbb{R}^{M \times N}$ , we can get the 2D-3D correspondences  $Match_I$ , where map  $Map_{cor} \in \mathbb{R}^{M \times N}$  is expressed as

$$Map_{cor} = \begin{bmatrix} \delta(p_1, P_1) & \cdots & \delta(p_1, P_N) \\ \vdots & \ddots & \vdots \\ \delta(p_M, P_1) & \cdots & \delta(p_M, P_N) \end{bmatrix}, \quad (9)$$

and the final correspondence  $Match_I$  is

$$Match_I = \begin{bmatrix} P_{p_1} \\ \vdots \\ P_{p_i} \\ \vdots \\ P_{p_M} \end{bmatrix}. \quad (10)$$

$Match_I$  is the 2D-3D correspondence derived from  $Map_{cor} \in \mathbb{R}^{M \times N}$  by point-max selection, where each element embeds a pair of corresponding 2D pixel and the 3D point. Through point-max,

the most similar match is selected in the one-to-many pixel-point correspondences, which alleviates the adverse impact of the outliers. With  $\text{Match}_I$ , a set of equations that relate the observed 2D image coordinates to their corresponding 3D point coordinates can be established. In this way, the I2P registration problem becomes a PnP problem, and we iteratively optimize the pose through EPnP under RANSAC.

### 3.6 Loss Function

In I2P<sub>ppsim</sub>, the performance of SFAM and CVD is particularly important. In order to obtain better cross-modal descriptions and more accurate co-view region detection results, we propose a joint loss function, which consists of the correlation loss and the co-view loss. For the correlation loss, the matched pixel-point pairs are expected to have cross-modal features with higher similarity, and vice versa for unmatched pairs. It implies that the correlations between matched pairs and unmatched pairs should have significant differences. For the co-view loss, based on the output scores of CVD, we expect that the points and pixels in the co-view region have higher scores, and the scores for outliers should be lower.

*Correlation Loss.* For the input image-point cloud pair  $(\mathcal{I}, \mathcal{P})$ , with the ground truth of relative pose  $T \in SE(3)$  and intrinsic matrix  $K$  of the camera, the reprojection error  $e_{pro}$  between the pixel  $\mathbf{p}_i$  and point  $\mathbf{P}_j = [x_j, y_j, z_j]^T$  can be calculated by

$$e_{pro}(\mathbf{p}_i, \mathbf{P}_j) = \left\| \mathbf{p}_i - \pi \left( \frac{1}{z_j} K T \bar{\mathbf{P}}_j \right) \right\|_2. \quad (11)$$

When  $e_{pro}$  is less than a safety threshold  $e_t$ ,  $\mathbf{p}_i$  and  $\mathbf{P}_j$  can be considered as a matched pair, otherwise an unmatched pair. The matching ground truth is denoted by  $G(\mathbf{p}_i, \mathbf{P}_j)$

$$G(\mathbf{p}_i, \mathbf{P}_j) = \begin{cases} 1, & \text{if } e_{pro}(\mathbf{p}_i, \mathbf{P}_j) < e_t \\ 0, & \text{otherwise.} \end{cases}. \quad (12)$$

With  $G(\mathbf{p}_i, \mathbf{P}_j)$ , the correlation loss for a predicted pixel-point pair can be calculated according to the correlation between the pixel-point feature vectors. As mentioned in Equation (7), the correlation between the feature vectors is denoted by  $\delta(\mathbf{p}_i, \mathbf{P}_j)$ , where  $\mathbf{p}_i$  and  $\mathbf{P}_j$  stand for the pixel and point, respectively. Defining the logits function as  $\xi$ , the loss of each pixel-point pair can be given as

$$\begin{aligned} \mathcal{L}(\mathbf{p}_i, \mathbf{P}_j) = & -w[G(\mathbf{p}_i, \mathbf{P}_j) \cdot \log \xi(\delta(\mathbf{p}_i, \mathbf{P}_j)) \\ & + (1 - G(\mathbf{p}_i, \mathbf{P}_j)) \cdot \log(1 - \xi(\delta(\mathbf{p}_i, \mathbf{P}_j)))] \end{aligned} \quad (13)$$

where  $w$  is the weight parameter.

There are  $M \times N$  pixel-point pairs in  $\text{Map}_{cor} \in \mathbb{R}^{M \times N}$ , and most of them are unmatched pairs. In order to speed up the optimization and balance the sample distribution, during training, we randomly select  $n$  pixels and  $n$  points in the co-view region for loss calculation and construct a correlation map with the size of  $n \times n$ . Finally, based on  $\mathcal{L}(\mathbf{p}_i, \mathbf{P}_j)$ , the correlation loss is defined as

$$\mathcal{L}_c = \frac{1}{n^2} \sum_{i,j=0}^n \mathcal{L}(\mathbf{p}_i, \mathbf{P}_j). \quad (14)$$

*Co-View Loss.* Similar to the correlation loss, we sample  $n$  pixels  $\mathbf{I}_{pos}$  and  $n$  points  $\mathbf{P}_{pos}$  in the co-view region, and  $n$  pixels  $\mathbf{I}_{neg}$  and  $n$  points  $\mathbf{P}_{neg}$  out of the co-view region when calculating the

co-view loss. Instead of focusing on the correlation between pixels and points, the co-view loss concerns whether the pixel or point belongs to the co-view region, which is a binary classification problem. The classification scores of each pixel or point can be obtained by CVD, denoted by  $S_{I, pos}$ ,  $S_{P, pos}$ ,  $S_{I, neg}$ , and  $S_{P, neg}$ . We expect that CVD can make the pixels and points in the co-view region have higher scores and vice versa. So the co-view loss is defined as

$$\mathcal{L}_{cv} = \frac{1}{n} \sum (S_{I, neg} + S_{P, neg} - S_{I, pos} - S_{P, pos}). \quad (15)$$

Combining Equations (14) and (15) yields the final joint loss function

$$\mathcal{L}_{i2p} = \mathcal{L}_c + \mathcal{L}_{cv}. \quad (16)$$

## 4 Experiments

### 4.1 Setup

**4.1.1 Dataset.** Our I2P<sub>ppsim</sub> was evaluated on KITTI Odometry [12] and Oxford Robotcar [25]. **KITTI Odometry.** In KITTI Odometry, the images and point clouds were acquired from an RGB camera and a 3D LiDAR. The camera and LiDAR had fixed extrinsics  $T_{cam}^{pc} \in SE(3)$ . This fixed relative pose in training and testing would be prone to cause network overfitting. On this account, it was necessary to perform data augmentation. Therefore, we followed the design of CorrI2P, using a random pose  $T_r$  to transform  $T_{cam}^{pc}$  and the point cloud. After taking augmentation, the relative pose of image-point cloud pair became  $T_{gt} = T_{cam}^{pc} T_r^{-1}$ . Besides, the relative translation between the image and point cloud was guaranteed less than 10 m. We followed the settings of Ren et al. [30] and Li et al. [23] to use the 0th–8th sequences for training, and the 9th–10th ones for testing. During training and testing, the size of the image was set to  $160 \times 512$ , and the number of points was 20,480. In total, there were 40,818 image-point cloud pairs used for training, and 5,584 pairs for testing.

**Oxford Robotcar.** Different from the acquisition method of point cloud in KITTI Odometry, the point clouds in Oxford Robotcar were captured by 2D scanning using a 2D LiDAR. To make the point clouds more dense, following DeepI2P, we spliced the adjacent point clouds at an interval of 2 m, and finally merged all the point clouds in an area within the radius of 50 m. About 34 sequences were used for training and 4 sequences for testing. During training and testing, the image size was set to  $384 \times 640$ , and the number of points was the same as that in KITTI Odometry. Finally, 109,398 image-point cloud pairs were used for training and 13,545 pairs for testing.

**4.1.2 Implementation Details.** We conducted all experiments on a workstation equipped with an AMD Ryzen9 5900X processor and an NVIDIA GeForce RTX 3090 GPU. I2P<sub>ppsim</sub> was implemented by Pytorch [26]. The Adam optimizer was used for network training. We trained our network over 25 epochs on each dataset. The batch size for training was 16, and 8 for testing. The learning rate of the optimizer was initialized as  $10^{-3}$ , and decayed by 75% every 5 epochs. During training, we set the safe threshold of the reprojection error ( $e_t$ ) to 1 pixel.

Some important hyperparameters in I2P<sub>ppsim</sub> are reported in Table 1. The Feature Extractor of image is ResNet34. It outputs four different scale feature maps ( $f_I^1 - f_I^4$ ) of SFAM as mentioned in Section 3.3.

In pose estimation, we experimentally set the co-view threshold of CVD as 0.9. The relative pose was estimated by EPnP under the RANSAC framework. The number of iterations was 500, and the reprojection error threshold was set to 1 pixel.

Table 1. Network Hyperparameters of I2P<sub>ppsim</sub>

Module	Layer Type	K	Channel Dimensions
Point cloud Feature Extractor	Layer 1	-	[32, 128, 256]
	Layer 2	-	[64, 64]
	KNN Layer	32	[256, 256], [512, 256]
	Layer 3	-	[256, 512]
SFAM	Coarse layer Pa	-	[256, $H_4/\text{scale} \times W_4/\text{scale}$ ]
	Coarse fusion layer Pa	-	[1,024, 512, 512]
	Coarse layer Pb	-	[256, $H_4/\text{scale} \times W_4/\text{scale}$ ]
	Coarse fusion layer Pb	-	[512, 128, 128]
	Fine Layer P	-	[128, 256, 128, 64]
	Fine Layer I	-	[128, 256, 128, 64]
CVD	Detector Head P	-	[128, 128, 64, 1]
	Detector Head I	-	[64, 64, 64, 1]

$H_4/W_4$  denotes a quarter of image height/width and  $\text{scale}$  represents the reduction factor during the network inference.

**4.1.3 Evaluation Metrics.** The **relative rotation error (RRE)** and the **relative translation error (RTE)** are adopted to evaluate the performance of the registration, which are formulated as

$$\begin{aligned} \text{RRE} &= \sum_{i=1}^3 |\theta(i)| \\ \text{RTE} &= \|\mathbf{t}_{\text{pred}} - \mathbf{t}_g\|_2, \end{aligned} \quad (17)$$

where  $\theta(\cdot)$  denotes the Euler angle of  $R_{\text{pred}}^{-1} R_g$ , and  $\theta(1)$ ,  $\theta(2)$ , and  $\theta(3)$  are roll, pitch, and yaw, respectively. We denote the ground truth of the rotation matrix and relative translation vector by  $R_g$ ,  $t_g$ , and the predicted ones by  $R_{\text{pred}}$  and  $t_{\text{pred}}$ .

**4.1.4 Compared Methods.** We compared our I2P<sub>ppsim</sub> with four other approaches, which were Grid Cls.+PnP [23], DeepI2P (2D) [23], DeepI2P (3D) [23], CorrI2P [30], EFGHNet [17], and EP2P-Loc [19]:

- *Grid Cls.+EPnP*. This method was proposed in DeepI2P [23]. It divides the image into grids with the same size. For example, in the evaluation on KITTI Odometry, an image is divided into 80 grids, with a size of  $5 \times 16$ . Grid Cls.+EPnP predicts which grid the point cloud belongs to through a classification network. The 2D-3D correspondences are built from the classification results. With such correspondences, a PnP solver is used to estimate the relative pose.
- *DeepI2P*. Based on the idea of frustum binary classification, DeepI2P trains a frustum classifier to judge whether the point cloud is within the field of view of the camera. With the classification results, it proposes 2D/3D inverse camera projection to estimate the relative pose, called DeepI2P (2D)/DeepI2P (3D), respectively.
- *CorrI2P*. This method builds 2D-3D correspondences using the outputs of multi-modal branches. Based on those correspondences, the relative pose is estimated accordingly. We used the same network settings as those in the paper to reproduce the work. To our knowledge, CorrI2P is the SOTA approach for I2P registration.
- *EFGHNet*. This method adopts a divide-and-conquer strategy to decouple feature alignment and feature matching, and estimates the pose based on the feature matching results.

Table 2. Registration Accuracy on KITTI Odometry and Oxford Robotcar

	KITTI Odometry		Oxford Robotcar	
	RTE (m)	RRE (°)	RTE (m)	RRE (°)
Grid Cls.+EPnP	$3.22 \pm 3.58$	$10 \pm 13.74$	$1.91 \pm 1.56$	$2.94 \pm 10.72$
DeepI2P (3D)	$3.17 \pm 3.22$	$15.52 \pm 12.73$	$2.27 \pm 2.19$	$15.00 \pm 13.64$
DeepI2P (2D)	$3.28 \pm 3.09$	$7.56 \pm 7.63$	<b><math>1.65 \pm 1.36</math></b>	$4.14 \pm 4.90$
CorrI2P	$2.32 \pm 9.74$	$4.66 \pm 6.79$	$3.20 \pm 3.14$	$2.49 \pm 8.51$
EFGHNet	$4.83 \pm 2.92$	$4.58 \pm 8.67$	$3.78 \pm 3.48$	$4.76 \pm 5.69$
EP2P-Loc	$1.32 \pm 1.13$	$4.11 \pm 5.46$	$3.56 \pm 3.79$	$8.65 \pm 9.81$
I2P <sub>ppsim</sub>	<b><math>1.18 \pm 1.48</math></b>	<b><math>4.08 \pm 4.46</math></b>	$2.95 \pm 2.66$	<b><math>2.26 \pm 5.12</math></b>

The best results are highlighted in bold.

— *EP2P-Loc.* EP2P-Loc establishes 2D-3D correspondences through multi-scale matching and uses a differentiable PnP layer to directly estimate the relative pose. Due to the lack of open source code, we reproduced this method as described in the paper.

## 4.2 Quantitative Experiments

Different from CorrI2P which uses the error threshold to eliminate data with large RTE and RRE, we believe that using all test data to evaluate the accuracy of the algorithm can better reflect the robustness of the registration method. Therefore, we followed DeepI2P and compared I2P<sub>ppsim</sub> with competing methods on all test data of KITTI Odometry and Oxford Robotcar. The results achieved are reported in Table 2. It can be seen that I2P<sub>ppsim</sub> achieves significantly better performance over other counterparts on KITTI Odometry. As for Oxford Robotcar, although I2P<sub>ppsim</sub> achieves the best performance on RRE, it did not perform as well as other SOTA approaches on RTE. The main reason for this phenomenon is the preprocessing method of the point cloud. The point cloud in Oxford Robotcar is formed by accumulating the 2D LiDAR scanning results from nearby areas. As a result, ghosting and blurring of many dynamic objects may emerge, which makes it difficult for I2P<sub>ppsim</sub> to predict correct correspondences among pixels and points. However, Grid Cls.+EPnP, DeepI2P (2D), and DeepI2P (3D) only need to predict rough grid classification or point cloud visibility results, without establishing strict pixel-point correspondences, so their performance on translation estimation is relatively better.

To compare the registration performance in more detail, we drew the registration recall curve under different RRE and RTE thresholds on the two datasets and calculated the area under the curve in Figure 6. It can be seen that the performance of each approach is basically the same as Table 2 shows. However, we find that the leading edge of I2P<sub>ppsim</sub> over CorrI2P is narrow, which is different from the significant advantage of I2P<sub>ppsim</sub> shown in Table 2.

To dive deep into this phenomenon, we plotted the error distributions of the two approaches with a higher error threshold on KITTI Odometry in Figure 7. It can be found that many large errors emerge in the results of CorrI2P. In contrast, I2P<sub>ppsim</sub> performs more stably, and there are few cases where the errors are extremely large.

*Error Distribution.* I2P<sub>ppsim</sub>'s error distributions of RRE (°) and RTE (m) on KITTI Odometry and Oxford Robotcar are shown in Figure 8. Obviously, the translation estimation ability of I2P<sub>ppsim</sub> on

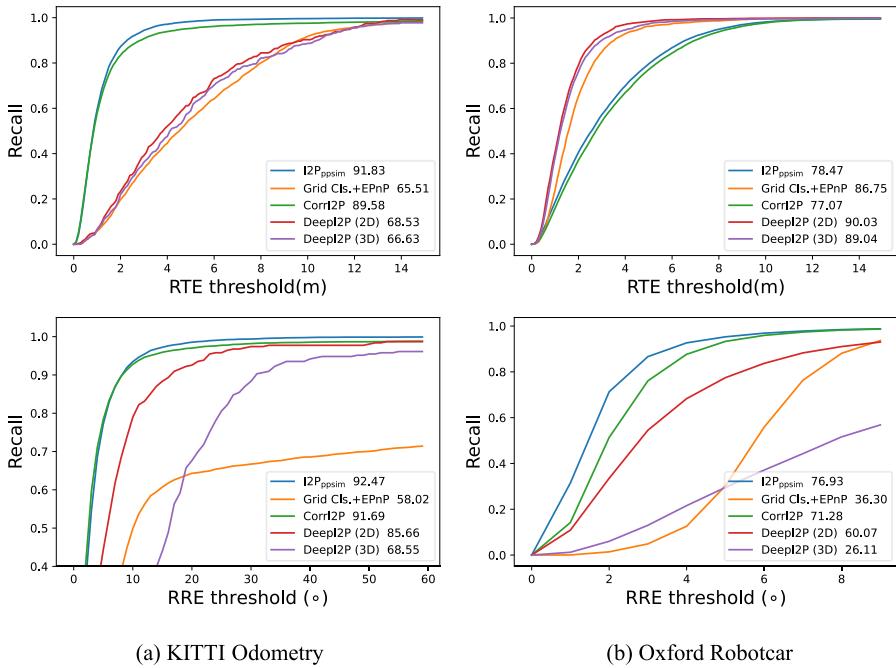


Fig. 6. Comparison of the registration recall of different methods with various RTE (m) and RRE (°) thresholds on KITTI Odometry and Oxford Robotcar. The area under each curve is presented behind the corresponding method's name.

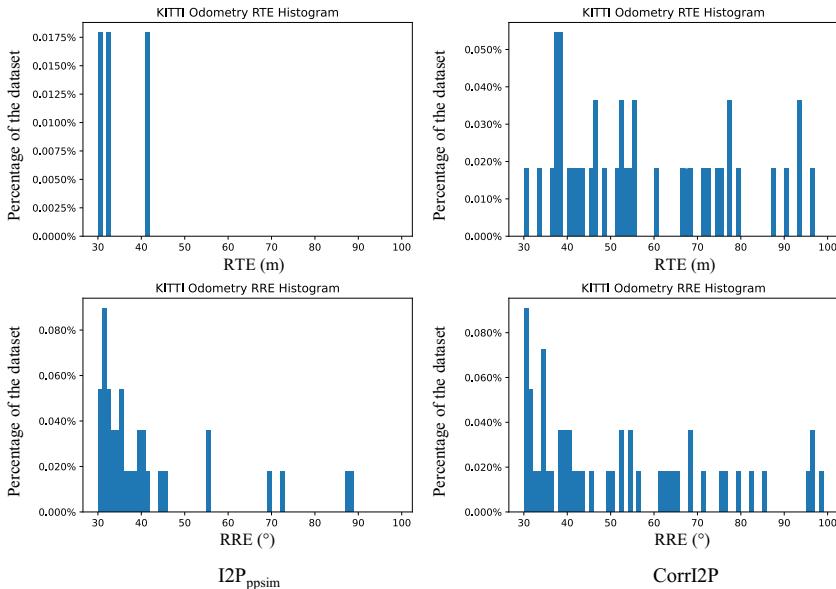


Fig. 7. Histograms of error distributions of I2P<sub>ppsim</sub> and CorrI2P at high RTE and RRE thresholds.

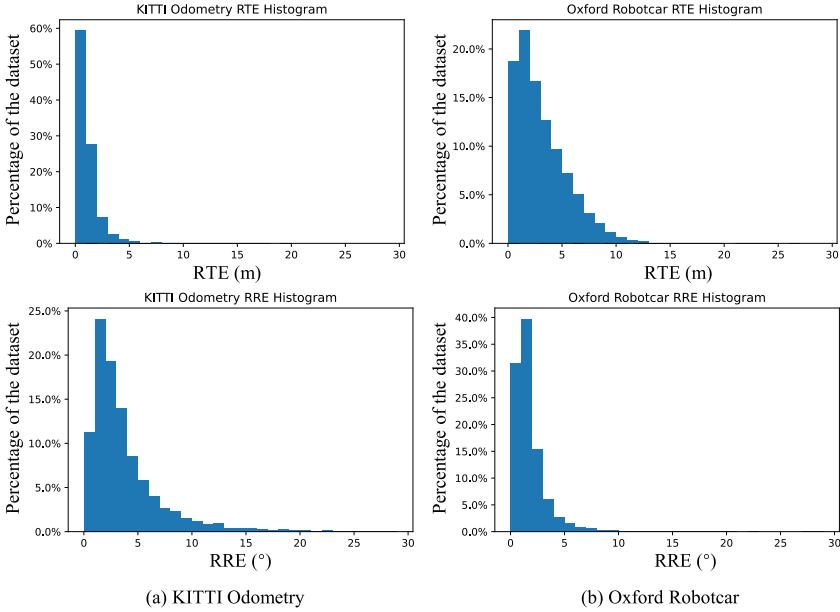


Fig. 8. Histograms of RTE (m) and RRE ( $^{\circ}$ ) on KITTI Odometry and Oxford Robotcar obtained by I2P<sub>ppsim</sub>.

KITTI Odometry is stronger than that on Oxford Robotcar, while the rotation estimation ability is weaker than that on Oxford Robotcar. The mode of RTE/RRE is  $\sim 0.5$  m/ $1^{\circ}$  on KITTI Odometry and  $\sim 2$  m/ $2^{\circ}$  on Oxford Robotcar.

#### 4.3 Qualitative Experiments

To compare the performance of different methods more intuitively, we demonstrated the 2D-3D correspondences they generated. When generating the correspondences, the reprojection error threshold was set to 1 pixel. We show the 2D-3D correspondences in Figure 9, where the wrong matches are in red and the correct ones in green.

Obviously, the “slopes” of the correspondence line segments are positively correlated with the reprojection errors. As mentioned in Section 3.5, as the visual depth grows, one pixel often corresponds to multiple points, which is particularly evident in Grid Cls.+EPnP and CorrI2P. For Grid Cls.+EPnP, amid the network’s inference, the final feature map’s size is downsampled at least 64 times compared with the raw input. After image downscaling, in KITTI Odometry, there are only  $5 \times 16$  image grids and 20,480 points classified into these grids. As shown in the column “Grid Cls.+EPnP” of Figure 9, one pixel may correspond to several or even dozens of points. Under the framework of grid classification of Grid Cls.+EPnP, even if the correspondences are correctly matched, their “slopes” can still be large. Such rough correspondences make it difficult to estimate the accurate pose. CorrI2P alleviates this phenomenon by reducing the times of image downscaling. However, there are still a lot of one-to-many correspondences, which mislead the pose estimation. Compared with the above methods, I2P<sub>ppsim</sub> handles the one-to-many correspondences through point-max. With such a strategy, the 2D-3D correspondences distribute more uniformly. Moreover, the wrong matches will be reduced from one cluster to one. These advantages are reflected in the notable improvement of the registration accuracy.

In addition, benefiting from SFAM proposed in Section 3.3, the 2D-3D correspondences generated by I2P<sub>ppsim</sub> are also more accurate than the competing methods. In Figure 9, to make the visualization

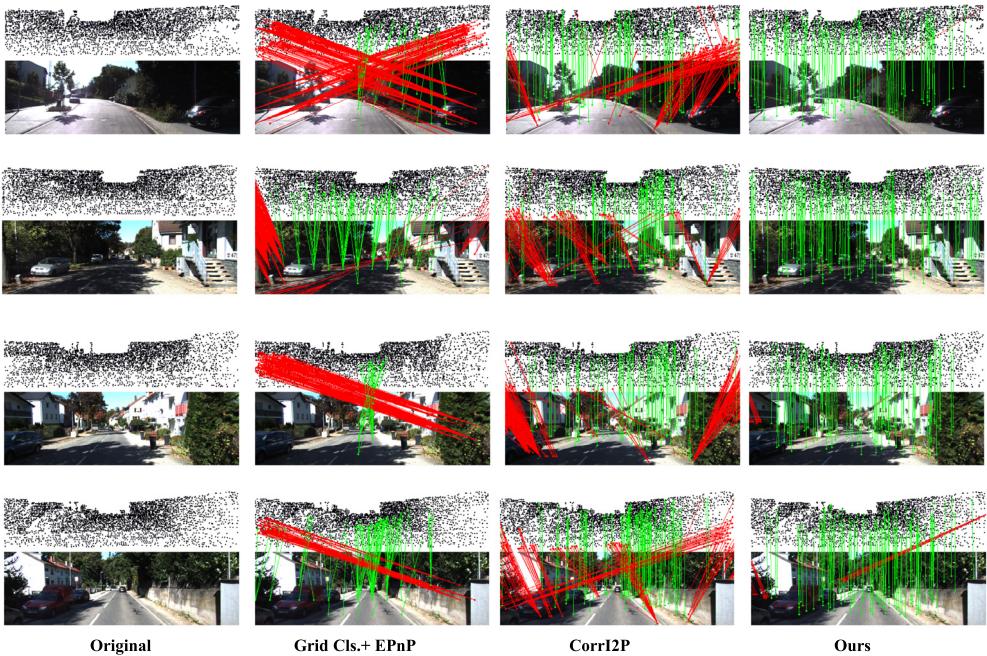


Fig. 9. 2D-3D correspondences established by different methods. The line segments represent 2D-3D correspondences, where red indicates false ones and green indicates correct ones.

clearer, we downsampled the correspondences and filtered out the wrong ones with reprojection errors of less than 15 pixels. It is notable that the correlation prediction of  $I2P_{ppsim}$  is more accurate, and there is basically no significant reprojection errors. We believe that the alignment feature space in SFAM plays an active role. Amid the coarse-to-fine feature alignment of  $I2P_{ppsim}$ , the fine-grained features are utilized to predict pixel-point correspondences, while the counterparts only conduct feature fusion in a one-stage pattern. Also, the effective guidance of the correlation loss enables  $I2P_{ppsim}$  to more directly predict the pixel-point correspondences, thus further improving the ability of correspondence prediction.

#### 4.4 Ablation Study

To verify the effectiveness of each module in our approach, we conducted ablation studies on our  $I2P_{ppsim}$  using KITTI Odometry. The baselines involved in the ablation study are elaborated as follows:

- *PC-CVD*: CVD for point clouds;
- *IMG-CVD*: CVD for images;
- *w/o CVD*:  $I2P_{ppsim}$  without CVD;
- *Direct Regression*:  $I2P_{ppsim}$  without SFAM and CVD;
- *CorrI2P (point-max)*: CorrI2P with point-max;
- $I2P_{ppsim}$  (*w/o point-max*):  $I2P_{ppsim}$  without point-max;
- $I2P_{ppsim}$  (*mutual check*):  $I2P_{ppsim}$  with mutual check.
- $I2P_{ppsim}$  (*w/o fine*):  $I2P_{ppsim}$  without fine stage in SFAM.
- $I2P_{ppsim}$  (*w/o coarse*):  $I2P_{ppsim}$  without coarse stage in SFAM.

Table 3. Accuracy of CVD on KITTI Odometry

	Recall	Precision	F2-Score	Accuracy
PC-CVD	0.962	0.901	0.949	0.973
IMG-CVD	0.782	0.732	0.771	0.791

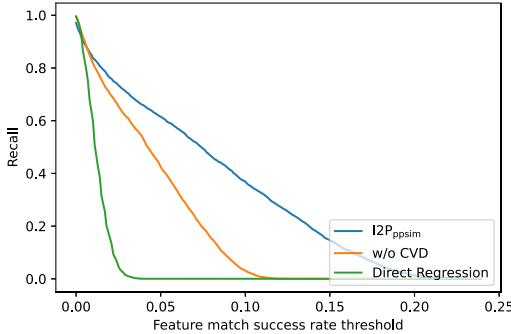


Fig. 10. Ablation study of SFAM on KITTI Odometry.

*Analysis on CVD.* I2P<sub>ppsim</sub> detects the co-view region between the camera and the LiDAR through CVD and establishes 2D-3D correspondences based on the points and pixels in that region. In this way, we can effectively reduce the time cost and storage space of  $\text{Map}_{cor}$ . Before CVD is used, the storage space of  $\text{Map}_{cor}$  is 12.5 MB. After CVD selecting,  $\text{Map}_{cor}$  is notably compressed to 1.15 MB. Also, the time cost of the correlation map computing is largely reduced from 33.3 seconds to 0.51 seconds. Furthermore, we evaluated CVD’s recall, precision, F2-score, and accuracy, and the results are reported in Table 3. It can be seen that CVD has advanced co-view region detection performance, achieving accuracies of 97.3% and 79.1% for images and point clouds, respectively.

*Analysis on SFAM.* SFAM can fully integrate the multi-scale features of images and point clouds and mine the correlation between them. It enables I2P<sub>ppsim</sub> to perform co-view region detection, improves the ability to extract cross-modal features, and enhances the accuracy of 2D-3D matching. To verify such a claim, we carried out ablation experiments to analyze the importance of SFAM. Specifically, we designed a network without SFAM and CVD, called Direct Regression. Direct Regression only extracts the respective features of the image and the point cloud and predicts the matching relationship between pixel and point via Direct Regression. Feature matching recall was utilized to analyze the performance of Direct Regression and the results are shown in Figure 10. It can be observed that Direct Regression without feature fusion has a poor ability to predict correct feature matching. This is mainly due to the huge modal difference between the image and the point cloud, which makes the network unable to extract effective common features. Consequently, the poor matching prediction ability of Direct Regression further justifies the necessity of feature fusion.

In addition, to explore the impact of the coarse-to-fine architecture on feature matching, we constructed a network without CVD (w/o CVD) which only uses fine-grained features to generate  $\text{Map}_{cor}$ . Also, the feature matching recall was utilized as the metric to compare the impact of using coarse-grained features for co-view region judgment on the establishment of 2D-3D correspondences. The corresponding I2P<sub>ppsim</sub> uses coarse-grained features for co-view region screening and then generates 2D-3D correspondences by fined-grained features. The relevant experimental results

Table 4. Ablation Results of Point-Max and SFAM

Method	RTE (m)	RRE (°)
CorrI2P	$2.32 \pm 9.74$	$4.66 \pm 6.79$
CorrI2P (point-max)	$1.45 \pm 1.62$	$4.32 \pm 5.03$
I2P <sub>ppsim</sub> (w/o point-max)	$1.95 \pm 2.97$	$4.31 \pm 6.41$
I2P <sub>ppsim</sub> (mutual check)	$1.44 \pm 2.19$	$4.26 \pm 6.13$
I2P <sub>ppsim</sub> (w/o fine)	$2.12 \pm 2.67$	$4.41 \pm 6.26$
I2P <sub>ppsim</sub> (w/o coarse)	$2.46 \pm 2.80$	$4.82 \pm 6.84$
I2P <sub>ppsim</sub>	<b><math>1.18 \pm 1.48</math></b>	<b><math>4.08 \pm 4.46</math></b>

The best results are highlighted in bold.

are shown in Figure 10. Compared with “w/o CVD,” it is evident that the introduction of coarse-grained features has a significant positive effect on the establishment of 2D-3D correspondences.

To analyze the impact of each stage of SFAM, we conducted ablation experiments on SFAM, and the results are presented in Table 4. We removed the fine stage of SFAM (I2P<sub>ppsim</sub> (w/o fine)) to analyze the impact of the weight-sharing network. Similarly, we removed the coarse stage of SFAM (I2P<sub>ppsim</sub> (w/o coarse)), which means that SFAM only adopts the weight-sharing network to fuse the features of images and point clouds. From Table 4, it can be seen that both the coarse stage and fine stage of SFAM have a positive effect on registration performance.

*Analysis on Point-Max.* Based on  $\text{Map}_{cor}$ , point-max is employed to find the matching point for each pixel. In this way, I2P<sub>ppsim</sub> filters out a large number of outliers that are difficult to eliminate only via the feature distances, therefore improving the registration accuracy.

To demonstrate the effectiveness of point-max, we compared the registration performance of I2P<sub>ppsim</sub>, I2P<sub>ppsim</sub> without point-max, and I2P<sub>ppsim</sub> with mutual check. Among them, I2P<sub>ppsim</sub> with mutual check was achieved by applying point-max to both pixels and point clouds and conducting consistency check on the point-max results. In addition, since point-max is a feature-independent matching constraint strategy that can be considered as a plug-and-play module for other approaches, we also applied it to CorrI2P and compared the registration performance of point-max before and after use. All the above experimental results are reported in Table 4. It is notable that the introduction of point-max significantly improves the performance of the I2P registration approaches, especially in terms of translation estimation. Such a result also confirms the negative impact of the one-to-many correspondences on registration. Specifically, the introduction of mutual check leads to another serious problem, i.e., the number of available 2D-3D correspondences will sharply decrease. The average number of matching pairs is only 307 for I2P<sub>ppsim</sub> with mutual check, while the count for the original I2P<sub>ppsim</sub> is 2,264. This sparsity of matching pairs has a detrimental effect on the performance of subsequent pose estimation which relies on nonlinear optimization techniques.

## 5 Conclusion

In order to fulfill the task of I2P registration, this article presents a novel framework based on alignment feature space learning, namely I2P<sub>ppsim</sub>. I2P<sub>ppsim</sub> leverages SFAM to enhance the network’s ability to extract cross-modal features. Moreover, a matching strategy called point-max is

proposed to address the one-to-many correspondences caused by scale ambiguity. Extensive experiments were conducted on benchmark datasets, namely KITTI Odometry and Oxford Robotcar, to demonstrate the outstanding performance of I2P<sub>ppsim</sub>. Furthermore, we conducted ablation studies to validate the efficacy of each module within the framework. The promising results obtained from these experiments suggest that I2P<sub>ppsim</sub> holds potential for utilization in other tasks that require cross-modal fusion, such as multi-sensor calibration. In future work, we will devote our efforts to further improve the scalability of our framework, e.g., to make it leverage the pose as supervisory information and support multi-agent systems.

## References

- [1] Yasuhiro Aoki, Hunter Goforth, Rangaprasad Arun Srivatsan, and Simon Lucey. 2019. PointNetLK: Robust & efficient point cloud registration using PointNet. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 7156–7165.
- [2] Xuyang Bai, Zixin Luo, Lei Zhou, Hongkai Chen, Lei Li, Zeyu Hu, Hongbo Fu, and Chiew-Lan Tai. 2021. PointDSC: Robust point cloud registration using deep spatial consistency. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 15854–15864.
- [3] R. Qi Charles, Hao Su, Mo Kaichun, and Leonidas J. Guibas. 2017. PointNet: Deep learning on point sets for 3D classification and segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 77–85.
- [4] Hongkai Chen, Zixin Luo, Jiahui Zhang, Lei Zhou, Xuyang Bai, Zeyu Hu, Chiew-Lan Tai, and Long Quan. 2021. Learning to match features with seeded graph matching network. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 6281–6290.
- [5] Hongkai Chen, Zixin Luo, Lei Zhou, Yurun Tian, Mingmin Zhen, Tian Fang, David McKinnon, Yanghai Tsin, and Long Quan. 2022. ASpanFormer: Detector-free image matching with adaptive span transformer. In *Proceedings of the European Conference on Computer Vision*, 20–36.
- [6] Zhi Chen, Kun Sun, Fan Yang, and Wenbing Tao. 2022. SC2-PCR: A second order spatial compatibility for efficient and robust point cloud registration. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 13211–13221.
- [7] Christopher Choy, Jaesik Park, and Vladlen Koltun. 2019. Fully convolutional geometric features. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 8957–8965.
- [8] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. 2018. SuperPoint: Self-Supervised Interest Point Detection and Description. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshop*, 337–349.
- [9] Patrick Ebel, Eduard Trulls, Kwang Moo Yi, Pascal Fua, and Anastasiia Mishchuk. 2019. Beyond cartesian representations for local descriptors. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 253–262.
- [10] Mengdan Feng, Sixing Hu, Marcelo H. Ang, and Gim Hee Lee. 2019. 2D3D-matchnet: Learning to match keypoints across 2D image and 3D point cloud. In *Proceedings of the International Conference on Robotics and Automation*, 4790–4796.
- [11] Kexue Fu, Shaolei Liu, Xiaoyuan Luo, and Manning Wang. 2021. Robust point cloud registration framework based on deep graph matching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 8889–8898.
- [12] Andreas Geiger, Philip Lenz, and Raquel Urtasun. 2012. Are we ready for autonomous driving? The KITTI vision benchmark suite. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3354–3361.
- [13] Md. Nazmul Haque, Moyuresh Biswas, Mark R. Pickering, and Michael R. Frater. 2012. A low-complexity image registration algorithm for global motion estimation. *IEEE Trans. Circuits Syst. Video Tech.* 22, 3 (2012), 426–433.
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 770–778.
- [15] Shengyu Huang, Zan Gojcic, Mikhail Usvyatsov, Andreas Wieser, and Konrad Schindler. 2021. PREDATOR: Registration of 3D point clouds with low overlap. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 4265–4274.
- [16] Youngseok Jang, Changsuk Oh, Yunwoo Lee, and H. Jin Kim. 2021. Multirobot collaborative monocular SLAM utilizing rendezvous. *IEEE Trans. Robot.* 37, 5 (2021), 1469–1486.
- [17] Yurim Jeon and Seung-Woo Seo. 2022. EFGHNet: A versatile image-to-point cloud registration network for extreme outdoor environment. *IEEE Robot. Autom. Letters* 7, 3 (2022), 7511–7517.
- [18] Marco Karrer, Patrik Schmuck, and Margarita Chli. 2018. CVI-SLAM-collaborative visual-inertial SLAM. *IEEE Robot. Autom. Letters* 3, 4 (2018), 2762–2769.

- [19] Minjung Kim, Junseo Koo, and Gunhee Kim. 2023. EP2P-Loc: End-to-end 3D point to 2D pixel localization for large-scale visual localization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 21470–21480.
- [20] Ik-Hyun Lee and Tae-Sun Choi. 2013. Accurate registration using adaptive block processing for multispectral images. *IEEE Trans. Circuits Syst. Video Tech.* 23, 9 (2013), 1491–1501.
- [21] Vincent Lepetit, Francesc Moreno-Noguer, and Pascal Fua. 2009. EPnP: An accurate O(n) solution to the PnP problem. *Int. J. Comput. Vis.* 81 (2009), 155–166.
- [22] Jiaxin Li, Ben M. Chen, and Gim Hee Lee. 2018. SO-net: Self-organizing network for point cloud analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 9397–9406.
- [23] Jiaxin Li and Gim Hee Lee. 2021. DeepI2P: Image-to-point cloud registration via deep classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 15955–15964.
- [24] David G. Lowe. 2004. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vis.* 60, 2 (2004), 91–110.
- [25] Will Maddern, Geoff Pascoe, Chris Linegar, and Paul Newman. 2017. 1 Year, 1000km: The Oxford RobotCar dataset. *Int. J. Robot. Res.* 36, 1 (2017), 3–15.
- [26] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An imperative style, high-performance deep learning library. In *Proceedings of the Advances in Neural Information Processing Systems*, 8024–8035.
- [27] Charles R. Qi, Li Yi, Hao Su, and Leonidas J. Guibas. 2017. PointNet++: Deep hierarchical feature learning on point sets in a metric space. In *Proceedings of the Advances in Neural Information Processing Systems*, 5105–5114.
- [28] Zheng Qin, Hao Yu, Changjian Wang, Yulan Guo, Yuxing Peng, Slobodan Ilic, Dewen Hu, and Kai Xu. 2023. Geo-Transformer: Fast and robust point cloud registration with geometric transformer. *IEEE Trans. Circuits Syst. Video Tech.* 33, 8 (2023), 9806–9821.
- [29] B. S. Reddy and B. N. Chatterji. 1996. An FFT-based technique for translation, rotation, and scale-invariant image registration. *IEEE Trans. Imag. Process.* 5, 8 (1996), 1266–1271.
- [30] Siyu Ren, Yiming Zeng, Junhui Hou, and Xiaodong Chen. 2023. CorrI2P: Deep image-to-point cloud registration via dense correspondence. *IEEE Trans. Circuits Syst. Video Tech.* 33, 3 (2023), 1198–1208.
- [31] Sajad Saeedi, Michael Trentini, Mae Seto, and Howard Li. 2016. Multiple-robot simultaneous localization and mapping: A review. *J. Field Robot.* 33, 1 (2016), 3–46.
- [32] PaulEdouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. 2020. SuperGlue: Learning feature matching with graph neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 4937–4946.
- [33] Vinit Sarode, Xueqian Li, Hunter Goforth, Yasuhiro Aoki, Rangaprasad Arun Srivatsan, Simon Lucey, and Howie Choset. 2019. PCRNNet: Point cloud registration network using pointnet encoding. arXiv:1908.07906. Retrieved from <https://arxiv.org/abs/1908.07906>
- [34] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. 2021. LoFTR: Detector-free local feature matching with transformers. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 8918–8927.
- [35] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Proceedings of the Advances in Neural Information Processing Systems*, 6000–6010.
- [36] Yue Wang and Justin Solomon. 2019. Deep closest point: Learning representations for point cloud registration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 3522–3531.
- [37] Yue Wu, Xidao Hu, Yue Zhang, Maoguo Gong, Wenping Ma, and Qiguang Miao. 2023. SACF-Net: Skip-attention based correspondence filtering network for point cloud registration. *IEEE Trans. Circuits Syst. Video Tech.* 33, 8 (2023), 3585–3595.
- [38] Yue Wu, Jiaming Liu, Maoguo Gong, Peiran Gong, Xiaolong Fan, A. K. Qin, Qiguang Miao, and Wenping Ma. 2024. Self-supervised intra-modal and cross-modal contrastive learning for point cloud understanding. *IEEE Trans. Multimedia* 26 (2024), 1626–1638.
- [39] Yue Wu, Jiaming Liu, Yongzhe Yuan, Xidao Hu, Xiaolong Fan, Kunkun Tu, Maoguo Gong, Qiguang Miao, and Wenping Ma. 2023. Correspondence-free point cloud registration via feature interaction and dual branch. *IEEE Comput. Intell. Mag.* 18, 4 (2023), 66–79.
- [40] Yue Wu, Yue Zhang, Xiaolong Fan, Maoguo Gong, Qiguang Miao, and Wenping Ma. 2023. INENet: Inliers estimation network with similarity learning for partial overlapping registration. *IEEE Trans. Circuits Syst. Video Tech.* 33, 3 (2023), 1413–1426.

- [41] Yue Wu, Yue Zhang, Wenping Ma, Maoguo Gong, Xiaolong Fan, Mingyang Zhang, A. K. Qin, and Qiguang Miao. 2023. RORNet: Partial-to-partial registration network with reliable overlapping representations. *IEEE Trans. Neural Netw. Learn. Syst.* (2023), 1–14.
- [42] Yuting Xie, Yuchen Zhang, Long Chen, Hui Cheng, Wei Tu, Dongpu Cao, and Qingquan Li. 2022. RDC-SLAM: A real-time distributed cooperative SLAM system based on 3D LiDAR. *IEEE Trans. Intell. Transp. Syst.* 23, 9 (2022), 14721–14730.
- [43] Shen Yan, Maojun Zhang, Yang Peng, Yu Liu, and Hanlin Tan. 2022. AgentI2P: Optimizing image-to-point cloud registration via behaviour cloning and reinforcement learning. *Remote Sens.* 14, 24 (2022), 6301.
- [44] Jiaqi Yang, Zhiqiang Huang, Siwen Quan, Qian Zhang, Yanning Zhang, and Zhiguo Cao. 2022. Toward efficient and robust metrics for RANSAC hypotheses and 3D rigid registration. *IEEE Trans. Circuits Syst. Video Tech.* 32, 2 (2022), 893–906.
- [45] Yongzhe Yuan, Yue Wu, Xiaolong Fan, Maoguo Gong, Wenping Ma, and Qiguang Miao. 2023. EGST: Enhanced geometric structure transformer for point cloud registration. *IEEE Trans. Visualization Comput. Graph.* 30 (2023), 6222–6234.
- [46] Ji Zhang and Sanjiv Singh. 2014. LOAM: Lidar odometry and mapping in real-time. In *Proceedings of the Robotics: Science and Systems*.
- [47] Tianjun Zhang, Hao Deng, Lin Zhang, Shengjie Zhao, Xiao Liu, and Yicong Zhou. 2022. Online correction of camera poses for the surround-view system: A sparse direct approach. *ACM Trans. Multimedia Comput. Commun. Appl.* 18, 4 (2022), Article 106, 24 pages.
- [48] Xiyu Zhang, Jiaqi Yang, Shikun Zhang, and Yanning Zhang. 2023. 3D registration with maximal cliques. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 17745–17754.
- [49] Qianyi Zhou, Jaesik Park, and Vladlen Koltun. 2016. Fast global registration. In *Proceedings of the European Conference on Computer Vision*, 766–782.
- [50] Guanyu Zhu, Yong Zhou, Rui Yao, Hancheng Zhu, and Jiaqi Zhao. 2023. Cyclic self-attention for point cloud recognition. *ACM Trans. Multimedia Comput. Commun. Appl.* 19, 1s (2023), Article 49, 19 pages.
- [51] Jianke Zhu, Steven C. H. Hoi, Michael R. Lyu, and Shuicheng Yan. 2011. Near-duplicate keyframe retrieval by semi-supervised learning and nonrigid image matching. *ACM Trans. Multimedia Comput. Commun. Appl.* 7, 1 (2011), Article 4, 24 pages.
- [52] Danping Zou and Ping Tan. 2013. CoSLAM: Collaborative visual SLAM in dynamic environments. *IEEE Trans. Patt. Anal. Mach. Intell.* 35, 2 (2013), 354–366.

Received 19 April 2024; revised 15 August 2024; accepted 20 September 2024