

CoFiI2P: Image-to-Point Cloud Registration with Coarse-to-Fine Correspondences for Intelligent Driving

Shuhao Kang^{*}, Youqi Liao^{*}, Jianping Li[†] *Member, IEEE*, Fuxun Liang, Yuhao Li, Xianghong Zou, Fangning Li, Zhen Dong *Member, IEEE*, Bisheng Yang[†]

Abstract—Image-to-point cloud (I2P) registration is a fundamental task in the field of autonomous vehicles and transportation systems for cross-modality data fusion and localization. Existing I2P registration methods estimate correspondences at the point/pixel level, often overlooking global alignment. However, I2P matching can easily converge to a local optimum when performed without high-level guidance from global constraints. To address this issue, this paper introduces CoFiI2P, a novel I2P registration network that extracts correspondences in a coarse-to-fine manner to achieve the globally optimal solution. First, the image and point cloud data are processed through a Siamese encoder-decoder network for hierarchical feature extraction. Second, a coarse-to-fine matching module is designed to leverage these features and establish robust feature correspondences. Specifically, in the coarse matching phase, a novel I2P transformer module is employed to capture both homogeneous and heterogeneous global information from the image and point cloud data. This enables the estimation of coarse super-point/super-pixel matching pairs with discriminative descriptors. In the fine matching module, point/pixel pairs are established with the guidance of super-point/super-pixel correspondences. Finally, based on matching pairs, the transform matrix is estimated with the EPnP-RANSAC algorithm. Extensive experiments conducted on the KITTI dataset demonstrate that CoFiI2P achieves impressive results, with a relative rotation error (RRE) of 1.14 degrees and a relative translation error (RTE) of 0.29 meters. These results represent a significant improvement of 84% in RRE and 89% in RTE compared to the current state-of-the-art (SOTA) method. Qualitative results are available at <https://youtu.be/ovbedasXuZE>. The source code will be publicly released at <https://github.com/kang-1-2-3/CoFiI2P>.

This study was jointly supported by the National Natural Science Foundation Project (No. 42130105, No. 42201477), China Postdoctoral Science Foundation (2022M712441, 2022TQ0234), and Open Fund of Key Laboratory of Urban Spatial Information, Ministry of Natural Resources (Grant No. 2023ZD001). (Shuhao Kang and Youqi Liao are co-first authors and contribute equally to the paper.) (Corresponding authors: Jianping Li and Bisheng Yang)

Shuhao Kang is with LIESMARS, Wuhan University, Wuhan 430079, China, and the school of Engineering and Design, Technical University of Munich, 80333 Munich, Germany (e-mail: shuhao.kang@tum.de).

Youqi Liao and Zhen Dong are with the Hubei LuoJia Laboratory, Wuhan University, Wuhan 430079, China and also with the LIESMARS, Wuhan University, Wuhan 430079, China (e-mail: martin_liao@whu.edu.cn, dongzhenwhu@whu.edu.cn).

Jianping Li is with the School of Electrical and Electronic Engineering, Nanyang Technological University, 50 Nanyang Avenue 639798, Singapore. (e-mail: jianping.li@ntu.edu.sg).

Fangning Li is with Beijing Urban Construction Exploration and Surveying Design Research Institute Co. Ltd, Beijing 100101, China. (e-mail: 64877558@qq.com)

Fuxun Liang, Yuhao Li, Xianghong Zou and Bisheng Yang is with the LIESMARS, Wuhan University, Wuhan 430079, China. (e-mail: liangfuxun@whu.edu.cn, yhaoli@whu.edu.cn, ericxhzou@whu.edu.cn, bshyang@whu.edu.cn).

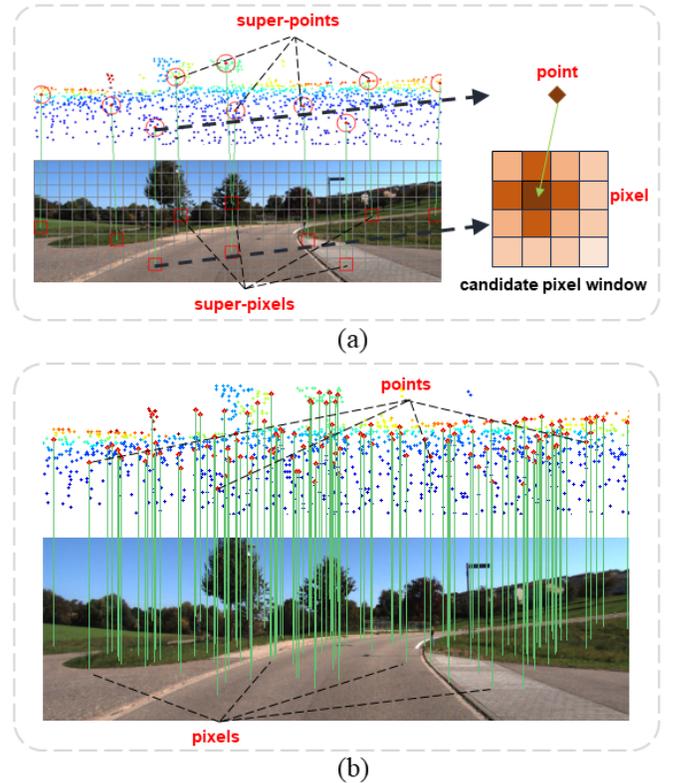


Fig. 1. Comparison of the proposed coarse-to-fine I2P registration scheme and the existing one-stage I2P registration pipeline. (a) shows the coarse-to-fine two-stage registration pipeline. Darker color represents higher similarity between the node point and the candidate pixel. (b) shows the one-stage registration pipeline.

Index Terms—Image-to-point cloud registration, coarse-to-fine correspondence, transformer

I. INTRODUCTION

Estimating six degrees of freedom (6-DoF) pose of a monocular image with respect to a pre-build point cloud map is a fundamental requirement for autonomous vehicles (AV) [1, 2, 3], object detection [4, 5, 6] and re-localization [7, 8]. Specifically, low-cost unmanned ground vehicles (UGVs) and autonomous vehicles (AVs) equipped with only a monocular camera frequently face challenges related to scale ambiguity in absolute localization and depth sensing. The establishment of a precise transformation relationship between

the image coordinate system and the pre-built point cloud coordinate system is of paramount importance, as it not only accurately localizes the camera but also effectively reduces data uncertainty inherent in monocular data. [9].

However, cross-modality registration has its inherent challenges. Some existing methods use hand-crafted detectors and descriptors for I2P registration [10, 11]. These approaches rely on structured features like edge or line features [10, 12], which are limited by specific environmental conditions. With the rapid development of deep learning (DL), learning-based I2P registration approaches [13, 14, 15] have been proposed to extract representative keypoints and descriptors. 2D3D-Matchnet [13] is the first learning-based I2P registration approach, which extracts SIFT [16] and ISS [17] keypoints from image and point cloud as registration primitives, respectively. Then it learns the descriptors with convolution neural networks (CNN). However, the manually designed detectors for keypoints extraction lead to poor correspondence accuracy. To alleviate the difficulty of correspondence construction and improve the registration success rate, DeepI2P [14] converts the registration problem to a classification problem. A novel binary classification network is designed to distinguish whether the projected points are within or beyond the camera frustum. The classification results are passed into an inverse camera projection solver to estimate the transformation between the camera and laser scanners. As a large number of points on the boundary are misclassified, the accuracy of the camera pose is still limited. CorrI2P [15] proposed an overlap region detector for both image and point cloud, then pixels and points in the overlap region are matched to obtain I2P correspondences. Feature fusion module is exploited to fuse the point cloud and image information. Although CorrI2P [15] has significant improvement over DeepI2P [14], matching merely on one stage, namely, the pixel-point level without global alignment guidance, can lead to local minimal and instability.

To sum up, there are two challenges of existing I2P registration methods: (1) Global context information is ignored and correspondences are established at the point/pixel level directly in the existing one-stage matching method. Fig. 1 illustrates the difference between the existing one-stage matching scheme and the coarse-to-fine matching scheme. (2) Existing methods incorporate global feature vectors into other modalities through feature fusion modules, but they tend to overlook the spatial relationships within the same-modality data and the cross-modality data. Consequently, each point or pixel retains limited knowledge of its global localization and lacks an understanding of its spatial affinity within both homogeneous and heterogeneous data. As a result, these existing I2P methods tend to be trapped in locally optimal solutions.

Inspired by recent coarse-to-fine matching schedules and transformers in image-to-image (I2I) registration [18, 19] and point cloud-to-point cloud (P2P) registration approaches [20, 21], this paper proposes **Coarse-to-Fine Image-to-Point cloud (CoFiI2P)** network to tackle the challenges of existing one-stage I2P registration. I2P transformer with self-attention modules and cross-attention modules is embedded into the network for global alignment. CoFiI2P mainly contains four modules: feature extraction (FE), coarse matching (CM), fine matching

(FM), and pose estimation (PE). FE module embeds input image and point cloud into shared high-dimensional spaces, and then CM and FM modules establish super-point/super-pixel correspondences and point/pixel correspondences. In the PE module, the relative transformation between the camera and laser scanners is estimated with EPnP-RANSAC algorithm [22, 23]. Overall, the main contributions of this paper are:

- 1) A novel coarse-to-fine I2P registration network is proposed to align image and point cloud in a progressive way. The coarse matching step provides robust but rough super-point/super-pixel correspondences for the fine matching step, which filters out most mismatched pairs and reduces the computation burden. The fine matching step achieves accurate and reliable point/pixel correspondences within the global guidance.
- 2) A novel I2P transformer that incorporates both self-attention and cross-attention modules is proposed to enhance its global-aware capabilities in homogeneous and heterogeneous data. The self-attention module enables the capture of spatial context within the same modality data, while the cross-attention module facilitates the extraction of hybrid features from both the image and point cloud data.

The rest of the paper is organized as follows. Section II reviews the existing I2I, P2P, and I2P registration works. Section III elaborates on the network design of CoFiI2P. Extensive experiments are conducted in Section IV. Section V discusses the crucial modules and factors in the experiments. Conclusion and future work are drawn in Section VI.

II. RELATED WORK

A. Same-modality Registration

1) *I2I Registration*: Before the age of deep-learning, hand-crafted detectors and descriptors (i.e., SIFT [16] and ORB [24]) are widely used to extract correspondences. Compared to traditional methods, learning-based methods improve the robustness and accuracy of image matching with large view-point differences and illumination changes. MagicPoint [25] proposed a point-tracking system powered by two CNNs. The first network in MagicPoint detects salient points and the second network estimates the transformation relationship. Based on MagicPoint, SuperPoint [26] proposed a self-learning training method through homography adaptation. SuperGlue [27] proposed an attention-based graph neural network (GNN) for feature matching. Patch2Pix [28] is the first work to obtain patch matches and regress pixel-wise matches in a coarse-to-fine manner. Cotr [19] proposed a hybrid transformer and CNN network that applies to both coarse and dense matching problems. LoFTR [18] operates detector-free matching to get dense correspondences and overcomes small receive field of CNN with transformer. However, I2I registration methods can not be directly transferred to the I2P registration task, which should extract heterogeneous features from cross-modality data.

2) *P2P Registration*: Point cloud registration aims to estimate the optimal rigid transformation of two point clouds.

Correspondence-based methods [29, 30, 31] estimate the correspondence first and recover the transformation with robust estimation methods [23]. RoReg [21] embeds orientation information of point cloud to estimate local orientation and refine coarse rotation through residual regression to achieve fine registration. In the transformer era, point-based transformers [20, 32, 33] have emerged and shown great performance. CoFiNet [34] followed LoFTR's [18] design and proposed coarse-to-fine correspondences for registration, which computed coarse matches with descriptors strengthened by the transformer and refined coarse matches through density adaptive matching module. REGTR [35] predicts overlap point sets and directly outputs transformed points to estimate transformation. PREDATOR [36] first predicts the overlap region between two point clouds and finds corresponding points with overlap scores and matching ability scores. GeoTransformer [20] proposed a geometric Transformer to make full use of the 3D properties of the point cloud. BUFFER [33] proposed an efficient, accurate and generalized network that utilizes both point-wise and patch-wise information for registration. As the point cloud and image are cross-modality data, which are characterized by totally different information, above same-modality registration methods cannot be directly applied.

B. Cross-modality Registration

To address the cross-modality registration problem, a variety of I2P registration methods have been proposed, which could be roughly divided into two categories: I2P fine registration (initial transformation dependent) and I2P coarse registration (initial transformation free). The I2P fine registration methods [37, 38, 39, 40] rely on initial transform parameters, and are widely applied in sensors calibration. Although this paper focuses on the second class, namely the coarse I2P registration without any initial transformation knowledge, we give an exhaustive review of both categories registration methods here.

1) *Fine registration methods*: Fine registration methods have been thoroughly researched for several decades. Early-stage studies [41, 42, 43] utilize various artificial targets as calibration constraints. In recent years, some approaches have argued that structured features shared in images and point cloud could be used for target-less I2P registration, e.g., edge feature [44] and line feature [10, 11]. Stamos et al. [44] computes orientation with vanishing points and then calculates camera position by matching 2D and 3D line features. Pandey et al. [45] uses mutual information (MI) as the registration criterion to automatically align a 3D laser scanner and camera. Zhang et al. [10] proposed a line-based method to register image and point cloud. Levinson and Thrun [37] employs edge information and optimizes transformation according to the response value. Recently, research on 2D and 3D semantic segmentation motivates semantic feature-based I2P registration. Li et al. [38] extracts vehicles from panoramic image and point cloud and maximizes the overlapping area through Particle Swarm Optimization (PSO) [46] to achieve accurate transformation estimation. Liao et al. [40] further studied semantic edges in I2P registration to employ more common semantic information instead of a certain type of object.

Overall, fine registration methods achieve high registration accuracy but rely on initial transform knowledge.

2) *Coarse registration methods*: 2D3D-Matchnet [13] is the pioneering method to regress the relative transform parameters with CNN. It extracts SIFT [16] and ISS [17] keypoints from image and point cloud and learns descriptors with a Siamese network. However, the hand-crafted detectors from different modalities match poorly. DeepI2P [14] proposed a feature fusion module to merge image and point cloud information and classify points in/beyond the camera frustum. CorrI2P [15] predicts pixels and points in overlapping areas and matches with dense per-pixel/per-point features directly to get I2P correspondences. Although overlapping region detectors significantly reduce the number of false candidates, I2P registration only using the low-level feature without global guidance leads to too many mismatches. Inspired by the coarse-to-fine strategies in CoFiNet [34], we propose the coarse-to-fine I2P registration network CoFiI2P, which injects the high-level response information into low-level matching for rejecting mismatches.

III. METHODOLOGY

For convenient description, a pair of partially overlapped image and point cloud are defined as $I \in \mathbb{R}^{W \times H \times 3}$ and $P \in \mathbb{R}^{N \times 3}$, where W and H are the width and height, and N is the number of points. The purpose of I2P registration is to estimate the relative transformation between the image I and point cloud P , including rotation matrix $\mathbf{R} \in \mathbb{R}^{3 \times 3}$ and translation vector $\mathbf{t} \in \mathbb{R}^3$.

Our method adopts the coarse-to-fine manner to find the correct correspondences set $\Omega(I, P)$. The CoFiI2P mainly consists of four modules: feature extraction (FE), coarse matching (CM), fine matching (FM) and pose estimation (PE). FE is an encoder-decoder structure network, that encodes raw inputs from different modalities into high-dimension feature spaces and finds in-frustum super-points. CM and FM are cascaded two-stage matching modules. CM constructs coarse matching at super-pixel/super-point level, and then FM constructs fine matching at pixel/point level sequentially with the guidance of super-pixel/super-point correspondences. Lastly, the PE module exploits point-pixel matching pairs to regress the relative transform with the EPnP-RANSAC [22, 23] algorithm. Workflow of the proposed method is shown in Fig. 2.

A. Multi-scale Feature extraction

We utilize ResNet-34 [47] and KPConv-FPN [48] as the backbones for image and point cloud to extract multi-level features. The encoder progressively embeds raw inputs into high-dimensional features, and the decoder propagates high-level information to low-level details with skip-connection for dense per-pixel/point feature generation. Then, we extract features from multiple resolutions for coarse-to-fine matching. Specifically, super-points set \hat{P} and super-pixels set \hat{I} at the coarsest resolution are chosen as candidates for coarse matching, and local points group $G_{\hat{p}}$ and pixels patches $G_{\hat{i}}$ are generated from coarse matching pairs for fine matching.

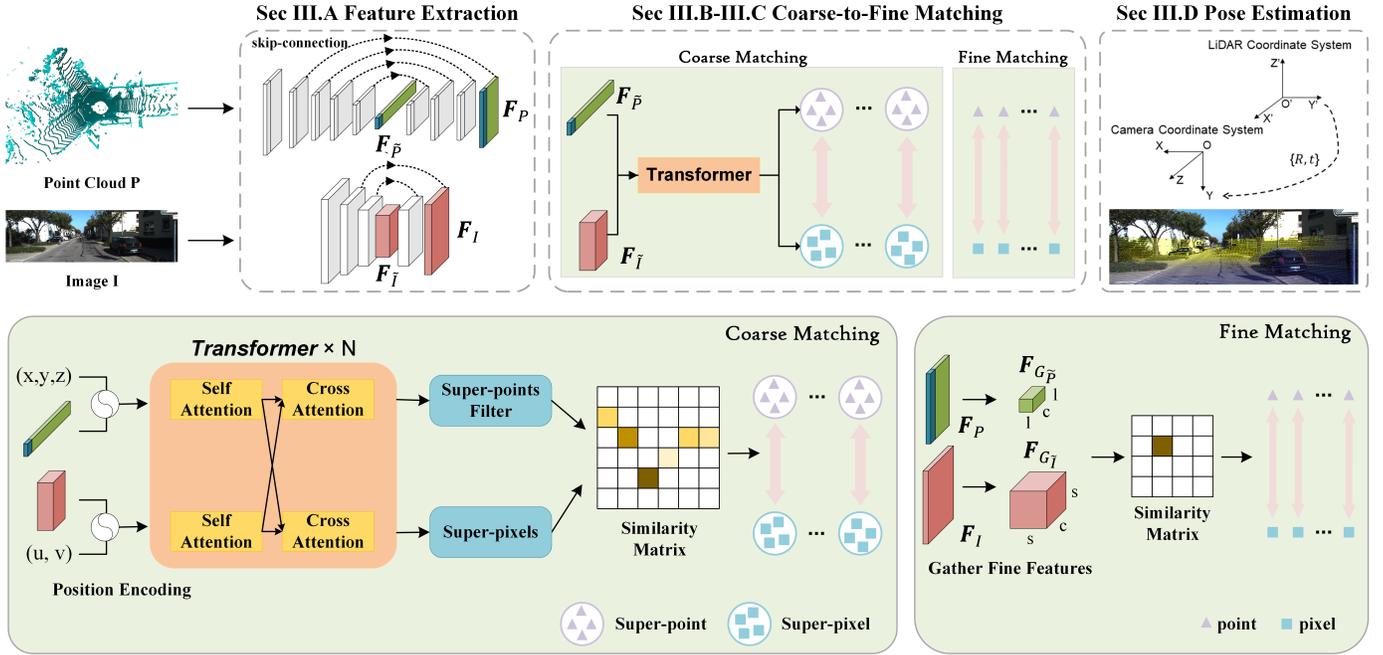


Fig. 2. Workflow of CoFil2P. The proposed method consists of feature extraction, coarse matching, fine matching and pose estimation modules. Image and point cloud are sent to the feature extraction module to obtain hierarchical deep features, respectively. The coarse-level features are strengthened by I2P transformer module and then matched with the cosine similarity rule. Fine features are gathered from the last layer of the decoder. In each super-point/super-pixel pair, the node point is set as the candidate and the corresponding pixel is selected from the super-pixel area, a $s \times s$ window. The generated dense matching pairs are utilized to regress the pose with the EPnP-RANSAC [22, 23] algorithm.

The local points group $G_{\tilde{p}_x}$ are constructed with the point-to-node strategy in the geometric space:

$$G_{\tilde{p}_x} = \{p \in P \mid \|p - \tilde{p}_x\| < r_g\}, \quad (1)$$

where r_g is the chosen radius. As the pixels array on the image with a rigid order, the local pixels patch $G_{\tilde{i}}$ construct simply with the pyramid matching strategy.

B. I2P Coarse Matching

In CM module, I2P transformer is utilized to capture the geometric and spatial consistency between image and point cloud. Each stage of the I2P transformer consists of a self-attention block for inter-modality long-range context and a cross-attention module for intra-modality feature exchange. The self-attention and cross-attention modules are repeated N times to extract well-mixed features for super-point/super-pixel correspondence matching.

1) *I2P Transformer*: [49, 50] have shown that Vision Transformer (ViT) outperforms traditional CNN-based methods with a large margin in classification, detection, segmentation and other tasks. Furthermore, recent approaches [18, 19, 32] have introduced transformer modules for I2I and P2P registration tasks. Therefore, we introduce the I2P transformer module customized for cross-modality registration task to enhance the representability and robustness of descriptors. Different from ViT used in the same-modality registration tasks, our I2P transformer contains both self-attention modules for space context capturing in homogeneous data and cross-attention for hybrid feature extraction among heterogeneous data.

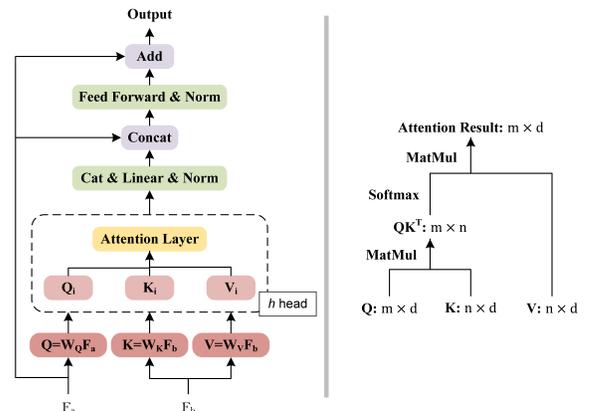


Fig. 3. Illustration of I2P Transformer module (left) and attention module (right).

For the self-attention module, given a coarse-level feature map $\mathbf{F} \in \mathbb{R}^{n \times C}$ of image or point cloud, the query, key and value vectors $\mathbf{q}, \mathbf{k}, \mathbf{v}$ are generated as:

$$\mathbf{q} = \mathbf{W}_q \mathbf{F}, \mathbf{k} = \mathbf{W}_k \mathbf{F}, \mathbf{v} = \mathbf{W}_v \mathbf{F}, \quad (2)$$

where $\mathbf{W}_q, \mathbf{W}_k \in \mathbb{R}^{C_{qk}^{sa} \times C}$, $\mathbf{W}_v \in \mathbb{R}^{C_v^{sa} \times C}$ are learnable weight matrixs, as shown in Fig. 3. Then, the global attention enhanced feature map \mathbf{F}^{sa} is calculated as:

$$\mathbf{F}^{sa} = \text{Softmax}\left(\frac{\mathbf{q}\mathbf{k}^\top}{\sqrt{C}}\right)\mathbf{v}. \quad (3)$$

Extracted global-aware features \mathbf{F}^{sa} are fed into the feed-forward network (FFN) to fuse the spatial relation information

in channel dimension. Given a feature map \mathbf{F} , the relative positions are encoded with multi-layer perception (MLP) [51].

Cross-attention is designed for fusing image and point cloud features in I2P registration task. Given the feature map $\mathbf{F}_{\tilde{P}}$, $\mathbf{F}_{\tilde{I}}$ for super-points set \tilde{P} and super-pixels set \tilde{I} , cross-attention enhanced feature maps $\mathbf{F}_{\tilde{P}}^{ca}$ of point cloud and $\mathbf{F}_{\tilde{I}}^{ca}$ of image are calculated as:

$$\begin{aligned}\mathbf{F}_{\tilde{P}}^{ca} &= \text{Softmax}\left(\frac{\mathbf{q}_{\tilde{P}}\mathbf{k}_{\tilde{I}}^\top}{\sqrt{C}}\right)\mathbf{v}_{\tilde{P}}, \\ \mathbf{F}_{\tilde{I}}^{ca} &= \text{Softmax}\left(\frac{\mathbf{q}_{\tilde{I}}\mathbf{k}_{\tilde{P}}^\top}{\sqrt{C}}\right)\mathbf{v}_{\tilde{I}},\end{aligned}\quad (4)$$

where $\mathbf{q}_{\tilde{P}}$, $\mathbf{k}_{\tilde{P}}$, $\mathbf{v}_{\tilde{P}}$ are query, key and value vectors of point cloud feature $\mathbf{F}_{\tilde{P}}$, and $\mathbf{q}_{\tilde{I}}$, $\mathbf{k}_{\tilde{I}}$, $\mathbf{v}_{\tilde{I}}$ are query, key and value vectors of image feature $\mathbf{F}_{\tilde{I}}$.

Remark. While self-attention module encodes the spatial and geometric features for each super-pixel and super-point, the cross-attention module injects the geometric structure information and texture information across image and point cloud respectively. Outputs of the I2P transformer carry powerful cross-modality information for matching.

2) *Super-point/Super-pixel Matching:* For the monocular camera, the field of view (FoV) is obviously smaller than the laser scans of 3D Lidar (e.g., Velodyne-H64), which usually sweeps 360 degrees in the horizontal direction. Therefore, only a small number of super-points are in the camera frustum. To filter out beyond-frustum super-points, we add a simple binary classification head to predict super-points in or beyond the frustum of the camera. After the beyond-frustum super-points are removed, the coarse level correspondences are estimated between in-frustum super-points set \tilde{P} and super-pixels set \tilde{I} by sorting the nearest super-pixel \tilde{i} in feature space:

$$\begin{aligned}\tilde{\Omega}_{\mathcal{M}} &= \{(\tilde{p}_x \in \tilde{P}, \tilde{i}_y \in \tilde{I}) \\ y &= \underset{|\tilde{I}|}{\text{argmin}} \|\mathbf{F}_{\tilde{P}}(\tilde{p}_x) - \mathbf{F}_{\tilde{I}}(\tilde{i}_y)\|\},\end{aligned}\quad (5)$$

where $\mathbf{F}_{\tilde{P}}$ and $\mathbf{F}_{\tilde{I}}$ are corresponding feature maps.

C. I2P Fine Matching

The first-stage matching at the coarse level constructs coarse super-point/super-pixel pairs $\tilde{\Omega}_{\mathcal{M}}$ but leads to poor registration accuracy. In order to obtain high-quality I2P correspondences, we generate fine correspondence based on the coarse matching results. During the decoder process, in each super-point/super-pixel correspondence $(\tilde{p}_x, \tilde{i}_y)$, the super-point \tilde{p}_x reverse to local points group $G_{\tilde{p}_x}$ and the super-pixel \tilde{i}_y reverse to local pixels patch $G_{\tilde{i}_y}$. Considering the uneven distribution of point cloud and computational efficiency, only the node points p_n in local patch group $G_{\tilde{p}_x}$ are selected to establish correspondences. For each node point p_n , we select the pixel $i_k \in G_{\tilde{i}_y}$ that lies nearest in the feature space. Point-pixel pairs in each super-point-super-pixel pair are staked together as the dense corresponding pairs $\Omega_{\mathcal{M}}$. With point feature map $\mathbf{F}_{G_{\tilde{p}_x}}$ and pixel feature map $\mathbf{F}_{G_{\tilde{i}_y}}$ of local patch, the fine matching process is defined as:

$$\begin{aligned}\Omega_{\mathcal{M}} &= \{(p_n \in G_{\tilde{p}_x}, i_k \in G_{\tilde{i}_y}) \\ k &= \underset{|G_{\tilde{i}_y}|}{\text{argmin}} \|\mathbf{F}_{G_{\tilde{p}_x}}(p_n) - \mathbf{F}_{G_{\tilde{i}_y}}(i_k)\|\}.\end{aligned}\quad (6)$$

D. EPnP-RANSAC based Pose Estimation

With the precise point-pixel pairs $\Omega_{\mathcal{M}}$, the relative transformation can be solved with EPnP [22] algorithm. As mentioned in previous approaches, wrong matching may infiltrate into the point-pixel pairs and decrease the registration accuracy. In the CoFiI2P, EPnP-RANSAC [22, 23] algorithm is used for robust estimating camera relative pose.

E. Loss Function

In order to supervise the network simultaneously learning coarse level descriptors, fine level descriptors and in/beyond-frustum super-points classification, we introduce a joint loss function consisting of coarse level descriptor loss \mathcal{L}_{coarse} , fine level descriptor loss \mathcal{L}_{fine} and classification loss $\mathcal{L}_{classify}$. The classification loss encourages the network to correctly label each super-point, and two descriptor losses pull positive matching pairs closer and push negative matching pairs farther in feature space.

The cosine similarity rule $s(p_x, i_y)$ of point cloud feature vector $\mathbf{F}(p_x)$ and image feature vector $\mathbf{F}(i_y)$ is defined as :

$$s(p_x, i_y) = \frac{\langle \mathbf{F}(p_x), \mathbf{F}(i_y) \rangle}{\|\mathbf{F}(p_x)\| \|\mathbf{F}(i_y)\|}, \quad (7)$$

and the distance $d(p_x, i_y)$ is defined as:

$$d(p_x, i_y) = 1 - s(p_x, i_y). \quad (8)$$

On the coarse level, the positive anchor \tilde{i}_{pos} for each in-frustum super-point \tilde{p}_x is sampled from the ground-truth pairs set $\tilde{\Omega}_{\mathcal{M}^*}$:

$$\tilde{\Omega}_{\mathcal{M}^*} = \{(\tilde{p}_x \in \tilde{P}, \tilde{i}_{pos} \in \tilde{I}) | \tilde{i}_{pos} = \Gamma(\mathbf{T}\tilde{p}_x)\}, \quad (9)$$

where \mathbf{T} is the ground-truth transform matrix from point cloud coordinate system to image frustum coordinate system, and Γ represents the mapping function that converts points from camera frustum to image plane coordinate system. The negative anchor \tilde{i}_{neg} is defined as the super-pixel with the smallest distance to the \tilde{p}_x in the feature space:

$$\tilde{i}_{neg} = \underset{\tilde{i} \in \tilde{I}}{\text{argmin}} \|d(\tilde{p}_x, \tilde{i}_{neg})\| \quad s.t. \quad \|\tilde{i}_{neg} - \tilde{i}_{pos}\| > r, \quad (10)$$

where r is the safe radius to remove adversarial samples. Finally, with positive margin Δ_{pos} and negative margin Δ_{neg} , coarse level descriptor loss is defined in a triplet way as Eq. (11) :

$$\begin{aligned}\mathcal{L}_{coarse} &= \sum_{(\tilde{p}_x, \tilde{i}_{pos}, \tilde{i}_{neg})} \left(\max(0, d(\tilde{p}_x, \tilde{i}_{pos}) - \Delta_{pos}) + \right. \\ &\quad \left. \max(0, \Delta_{neg} - d(\tilde{p}_x, \tilde{i}_{neg})) \right).\end{aligned}\quad (11)$$

Fine level descriptor loss is defined as modified circle loss [52]. We randomly select n in-frustum point and their positive

anchor pixels and negative anchors, as Eq. (10), then the descriptor loss is defined as:

$$\mathcal{L}_{fine} = \log \left(1 + \exp \left(-\gamma \alpha_{pos} (s(p_x, i_{pos}) - \Delta_{pos}) \right) \sum_{(p_x, i_{pos}, i_{neg})} \exp \left(\gamma \alpha_{neg} (p_x, s(i_{neg}) - \Delta_{neg}) \right) \right), \quad (12)$$

where α_{neg} and α_{pos} are the dynamic optimizing rates towards negative and positive pairs, and γ is the scale factor. As in [52], the α_{neg} and α_{pos} are defined as:

$$\begin{aligned} \alpha_{neg} &= \max(0, s(p_x, i_{neg}) + \Delta_{neg}), \\ \alpha_{pos} &= \max(0, 1 + \Delta_{pos} - s(p_x, i_{pos})). \end{aligned} \quad (13)$$

Super-points classification loss is a binary cross-entropy loss:

$$\mathcal{L}_{classify} = - \sum_{\tilde{p}_x \in \tilde{P}} (\tilde{p}_x \log(\tilde{p}_x^*) + (1 - \tilde{p}_x) \log(1 - \tilde{p}_x^*)). \quad (14)$$

Overall, our loss function is

$$\mathcal{L} = \lambda_1 \mathcal{L}_{coarse} + \lambda_2 \mathcal{L}_{fine} + \lambda_3 \mathcal{L}_{classify}, \quad (15)$$

where $\lambda_1, \lambda_2, \lambda_3$ are hyperparameters that control the weights between losses.

IV. EXPERIMENTS

A. Experiment setup

1) *Dataset*: we evaluate our method on KITTI Odometry dataset [53], a benchmark dataset extensively employed in the field of I2P registration. KITTI Odometry dataset consists of 22 sequences of images and point cloud data collected simultaneously in urban environments, encompassing a variety of road scenes and environmental conditions, and 11 sequences provide ground-truth calibration files. The camera intrinsic function Γ extracted from calibration files is thought unbiased during experiments. For a fair comparison with existing approaches [14, 15], sequences 0-8 are used for training and 9-10 for testing.

2) *Baseline methods*: we compare proposed CoFiI2P with two open-sourced I2P methods:

- DeepI2P¹ [14]: proposed the frustum classification and inverse camera projection to estimate the camera pose. Due to the time-consuming data pre-processing and training, we reference the registration results of the paper.
- CorrI2P² [15]: is SOTA comparable I2P method. CorrI2P [15] predicts the overlapping area and establishes correspondences densely for pose estimation. We use officially released codes to reimplement this method. The only difference is the number of points changed from 40960 to 20480 for fair comparison.

3) *Evaluation metrics*: we report the relative rotation error (RRE), relative translation error (RTE) and registration recall (RR) to evaluate the registration results. Inlier ratio (IR) and root mean square error (RMSE) used in I2P and P2P approaches [34, 54] are introduced to evaluate the quality of correspondences. RRE and RTE are defined as :

$$\begin{aligned} \text{RRE} &= \sum_{i=1}^3 |\mathbf{r}(i)|, \\ \text{RTE} &= \|\mathbf{t}_{gt} - \mathbf{t}_e\|, \end{aligned} \quad (16)$$

where \mathbf{r} is the Euler angle vector of $\mathbf{R}_{gt}^{-1} \mathbf{R}_e$, \mathbf{R}_{gt} and \mathbf{t}_{gt} are the ground-truth rotation and translation matrix, \mathbf{R}_e and \mathbf{t}_e are the estimated rotation and translation matrix.

RR metric [34, 54] estimates the percentage of correctly matched pairs, indicating the descriptor learning ability of networks. RR is defined as:

$$\text{RR}_{\mathcal{M}} = \frac{1}{|\Omega_{\mathcal{M}}|} \sum_{i=1}^{|\Omega_{\mathcal{M}}|} \mathbb{1} \left(\sqrt{\frac{1}{|\Omega_{\mathcal{M}}|} \sum_{(p_x, i_y) \in \Omega_{\mathcal{M}}} \|\Gamma(\mathbf{T}p_x) - i_y\|} < \tau \right), \quad (17)$$

τ is the threshold to remove false registration results, i.e. $10^\circ/5m$. IR denotes the inlier ratio of matching pairs, which measures the accuracy of correspondences. The IR for point/pixel correspondences set $\Omega_{\mathcal{M}}$ is defined as :

$$\text{IR}_{\mathcal{M}} = \frac{1}{|\Omega_{\mathcal{M}}|} \sum_{(p_x, i_y) \in \Omega_{\mathcal{M}}} \mathbb{1}(\|\Gamma(\mathbf{T}p_x) - i_y\| < \tau), \quad (18)$$

in which τ is used to control the reprojection error tolerance. RMSE denotes the average reprojection error over the correspondences set $\Omega_{\mathcal{M}}$:

$$\text{RMSE}_{\mathcal{M}} = \sqrt{\frac{1}{|\Omega_{\mathcal{M}}|} \sum_{(p_x, i_y) \in \Omega_{\mathcal{M}}} \|\Gamma(\mathbf{T}p_x) - i_y\|}. \quad (19)$$

B. Implementation Details

Raw images and point clouds are preprocessed to filter out noise and reduce the computational burden. The top 50 rows of images are cropped out as most of these pixels see sky. Then images are resized to 160×512 for training and testing. The points are downsampled with $0.1m \times 0.1m \times 0.1m$ voxel and then randomly sampled 20480 points as input. The resolution for coarse matching is 1/8 of original resolution of image and 1/16 of original resolution of point cloud, while 1/2 for both image and point cloud in fine matching.

During the training process, with correct transform parameters provided by the calibration files, the ground truth correspondences $\Omega_{\mathcal{M}^*}$ are established to supervise the network. We trained the whole network 25 epochs with batch size of 1. We use the Adam [55] to optimize the network, and the initial learning rate is 0.001 and multiple by 0.25 after every 5 epochs. For our joint loss, we set $\lambda_1 = \lambda_2 = \lambda_3 = 1$. The safe radius r , positive margin Δ_{pos} and negative margin Δ_{neg} in loss function are set to 1, 0.2 and 1.8 respectively. Scale factor γ is set to 10. During the coarse matching, super-points are projected to the frustum of the image, and beyond-frustum super-points are removed. Then we uniformly sample

¹<https://github.com/lijx10/DeepI2P>

²<https://github.com/rsy6318/CorrI2P>

TABLE I
REGISTRATION ACCURACY ON KITTI DATASET

	Threshold($^{\circ}/m$)	RRE($^{\circ}$)	RTE(m)	RR(%)
DeepI2P (2D)	10/5	7.56 ± 7.63	3.28 ± 3.09	-
DeepI2P (3D)	10/5	15.52 ± 12.73	3.17 ± 3.22	-
	none/none	6.93 ± 29.03	2.68 ± 10.51	-
CorrI2P	45/10	3.41 ± 3.64	1.48 ± 1.35	97.07
	10/5	2.70 ± 1.97	1.24 ± 0.87	90.66
	none/none	$1.14 \pm 0.78(-84\%)$	$0.29 \pm 0.19(-89\%)$	-
CoFiI2P	45/10	$1.14 \pm 0.78(-67\%)$	$0.29 \pm 0.19(-80\%)$	100.00(+2.93)
	10/5	$1.14 \pm 0.78(-58\%)$	$0.29 \pm 0.19(-77\%)$	100.00(+9.34)

64 super-points from in-frustum ones. The corresponding super-pixels are found with Eq. 4. During the testing process, we set the confidence threshold as 0.9 to obtain reliable in-frustum super-points in coarse matching. All the experiments are conducted on a single RTX 4090 GPU. The training of CoFiI2P takes about 49 hours.

C. Registration Accuracy

Different from [13, 14, 15] which use a specific set of thresholds (i.e. $10^{\circ}/5m$) to reject false registration frames, we report the RRE and RTE of our CoFiI2P under three different settings as listed in Table I. When we use *none/none* as the threshold, it implies that no specific thresholds are applied to filter out false registration frames. Conversely, when we set the thresholds to $10^{\circ}/5m$ and $45^{\circ}/10m$, it means that any frames with registration errors exceeding these specified thresholds are excluded during the evaluation process. By employing these diverse settings, we can conduct a more equitable analysis of the performance of the I2P registration methods.

The proposed method outperforms all baseline methods by a large margin, which indicates that the CoFiI2P is both robust and accurate in urban environments. Specifically, with different registration error threshold *none/none*, $45^{\circ}/10m$, $10^{\circ}/5m$, the RRE reduces 84%, 67%, 58% and the RTE reduces 89%, 80%, 77% compared to CorrI2P [15] respectively. Meanwhile, the RR improves by 2.93%, 9.34% under two pairs of thresholds. Notably, our method significantly improves registration accuracy and reduces inferior results. We indicate that the CorrI2P [15] takes poor performance under large view and illumination variations, but our CoFiI2P operates stably and robustly. Therefore, the global RRE and RTE are significantly better than the CorrI2P [15]. We project the points to image plane with transform parameters provided by ground-truth files, CorrI2P [15] and CoFiI2P respectively and render with depth and deploy the qualitative registration results in Fig. 5. It can be seen that the CoFiI2P remains stable and provides more accurate results, which proves our indication. More qualitative results are released in the video demo³.

The distribution of relative rotation error (RRE) and relative translation error (RTE) on the KITTI dataset are shown in Fig. 6. Neither thresholds nor other constraints are used to remove

false registration results. Fig. 6 shows that the predictions of our CoFiI2P are clearly concentrated in the interval with much smaller errors and serious misregistration cases are completely eliminated.

The precision, recall, accuracy and F1-score of frustum classification are shown in Table II. The results show that our CoFiI2P gets higher recall, accuracy and F1-score than CorrI2P [15], where the recall improves about 10.69% over the CorrI2P [15], but the precision drops about 5.96%. We indicate that the coarse-to-fine matching schedule excavates more overlapping areas and gets higher recall value, but classification on coarse level leads to lower precision. Overall, we have achieved a better balance between precision and recall.

TABLE II
FRUSTUM CLASSIFICATION RESULT

	Precision	Recall	Accuracy	F1-score
DeepI2P	-	-	94.00	-
CorrI2P	97.14	89.07	97.25	92.81
CoFiI2P	91.18	99.76	98.26	95.26

In order to evaluate the quality of correspondences directly, we introduce the inlier ratio (IR) metric as in P2P registration approach [34]. The quantitative results are shown in Fig. 7 and qualitative results are shown in Fig. 4. From Fig. 7, CoFiI2P gets obviously higher IR scores than CorrI2P [15] under various thresholds. As we relax the constraint, the advantage of CoFiI2P becomes more apparent. When the tolerance threshold is set to 5 pixels, IR score of our CoFiI2P is 80.75%, 49.91% higher than CorrI2P [15]; when set to 10 pixels, IR score of CoFiI2P rises to 97.49%, 43.35% higher than CorrI2P [15]. The fact that our method obtains higher IR scores than CorrI2P [15] with a small threshold indicates that the CoFiI2P corrects a large partition of severe misregistration correspondences. The RMSE distribution is shown in Fig. 8 and demonstrates that our CoFiI2P owns higher correspondence quality. Specifically, CoFiI2P reduces the global average RMSE value from 25.78 pixels of CorrI2P [15] to 3.58 pixels, about 86.11% improvement. From the pairs matching aspect, our method obviously improves the quality and quantity of correspondences in a global view and benefits downstream registration tasks.

³<https://youtu.be/ovbedasXuZE>

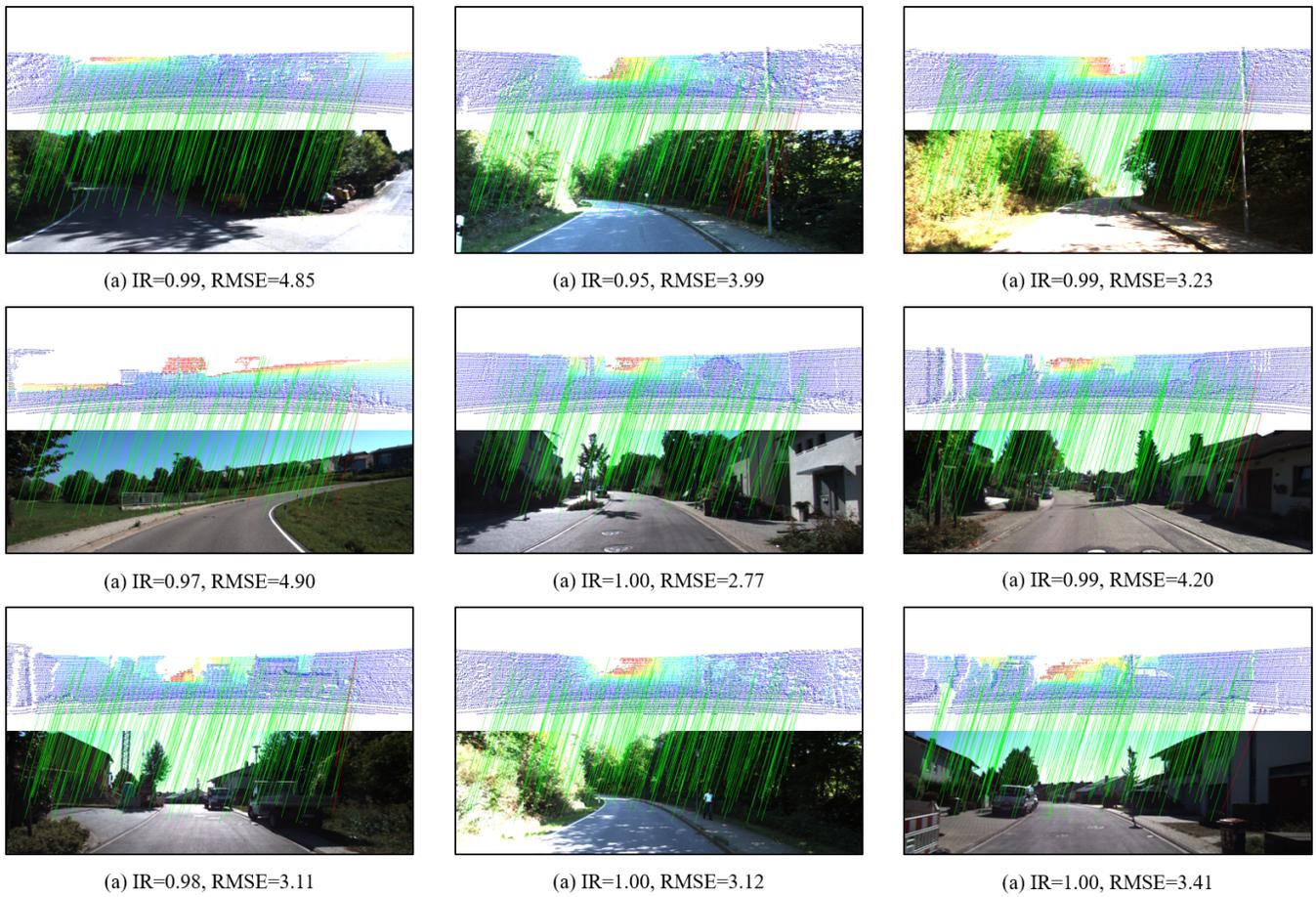


Fig. 4. Qualitative results of correspondences estimated by CoFiI2P.

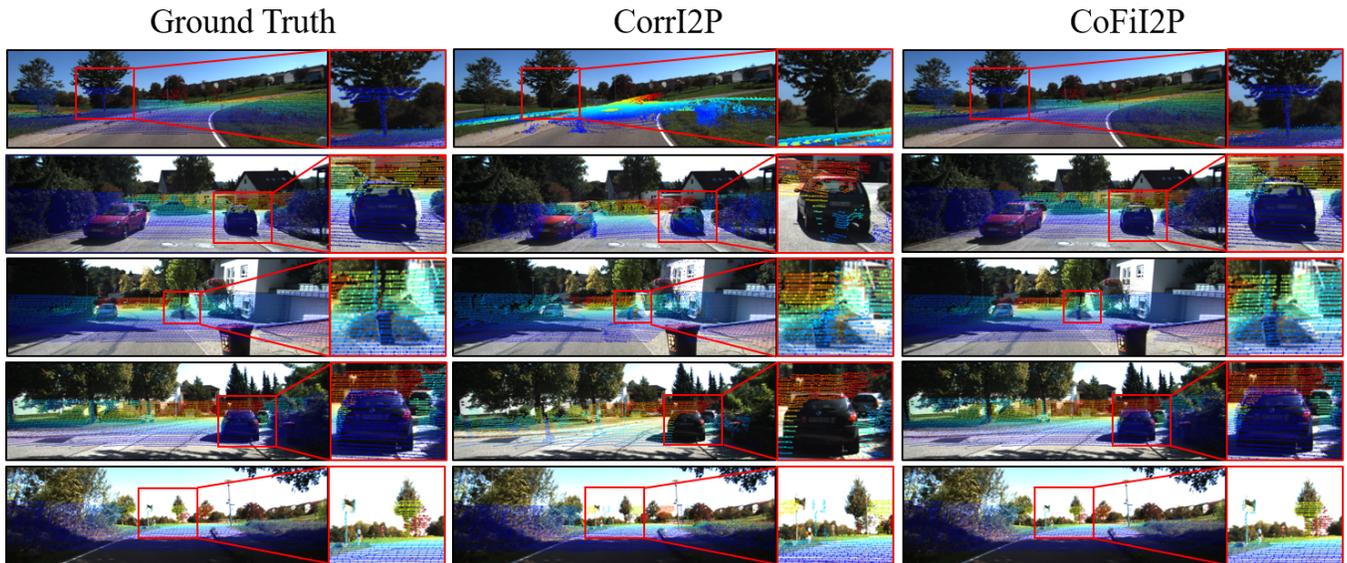


Fig. 5. Quantitative registration results.

V. DISCUSSION

In this part, we analyze three crucial factors in our CoFiI2P: I2P transformer module, coarse-to-fine matching scheme and point cloud density. We conduct ablation studies on them

to prove the effectiveness of each module and review the influence of point cloud density. In order to reduce the computational burden, we train the variant networks 10 epochs in the discussion part, instead of 25 epochs in the experiments

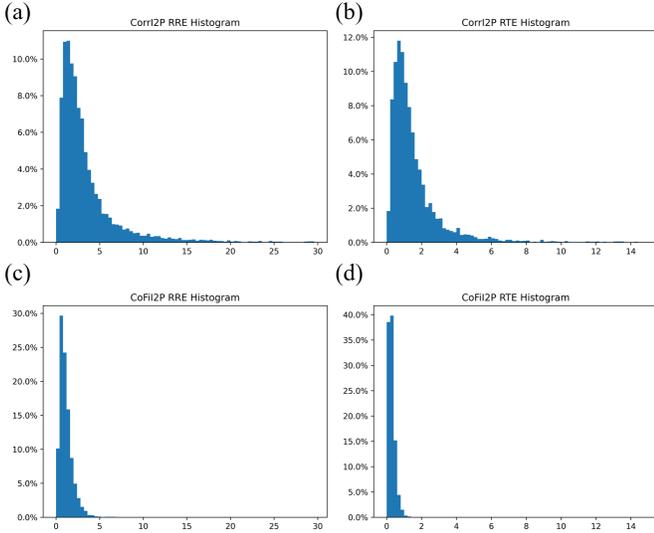


Fig. 6. Error distribution of RRE and RTE. (a) and (b) show the RRE and RTE distribution of Corri2P [15], and (c) and (d) show the RRE and RTE distribution of CoFiI2P respectively.

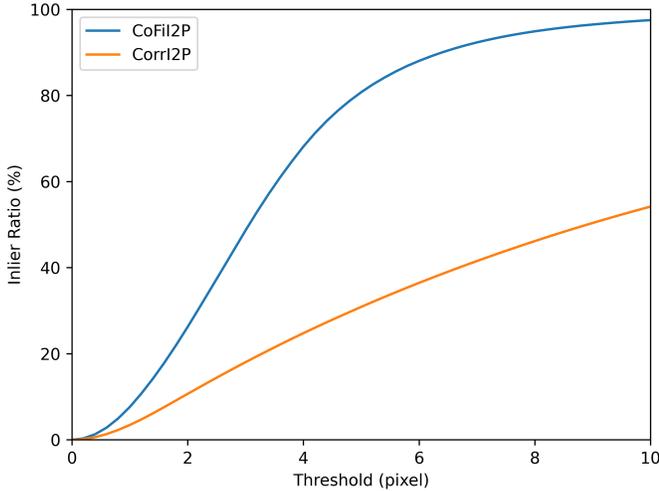


Fig. 7. Comparison of inlier ratio with different thresholds

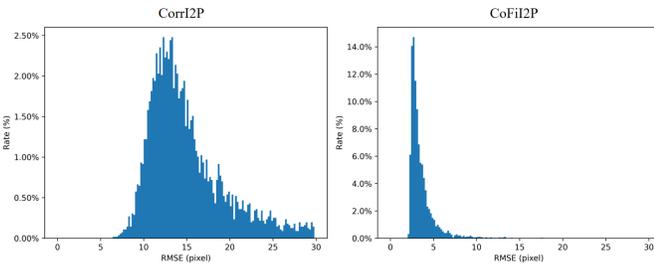


Fig. 8. RMSE distribution of Corri2P [15] and CoFiI2P

part. We report the global RRE and RTE as evaluation metrics and no thresholds are used to reject false registration scenes.

A. Analyze of the I2P transformer

I2P transformer [56] with self-attention module and cross-attention module is crucial to image-to-point cloud alignment

at the global level. In this part, we conduct ablation studies to assess the effectiveness of the I2P transformer. We train the CoFiI2P without any attention module as baseline. Then, the self-attention modules and cross-attention modules are added on the coarse level respectively. Table III shows registration results. As shown in Table III, the baseline method performance drops significantly without any attention module. Besides, the self-attention module reduces the RRE about 1.06° and the RTE about $0.47m$, and the cross-attention module reduces the RRE about 0.92° and the RTE about $0.38m$ respectively. With both the self-attention and cross-attention modules, the CoFiI2P achieves the smallest registration error. Moreover, with both self-attention and cross-attention modules, the variance reduces by a large margin, which indicates that the I2P transformer block enhances both the accuracy and robustness of network.

TABLE III
ABLATION STUDY ON I2P TRANSFORMER BLOCKS. SA DENOTES SELF-ATTENTION MODULE AND CA DENOTES CROSS-ATTENTION MODULE.

Baseline	SA	CA	RRE($^\circ$)	RTE(m)
✓			2.88 ± 5.32	0.90 ± 2.44
✓	✓		1.82 ± 2.71	0.43 ± 0.75
✓		✓	1.96 ± 1.95	0.52 ± 0.59
✓	✓	✓	1.28 ± 0.94	0.33 ± 0.23

B. Comparison of coarse-to-fine matching scheme and one-stage matching scheme

Our CoFiI2P proposes to estimate coarse correspondences at super-point/super-pixel level first and then generate dense correspondences at point-pixel level sequentially. To demonstrate that the progressive two-stage registration operates better than the one-stage registration used in previous I2P approaches, we conduct ablation experiments on the coarse-to-fine matching scheme. This ablation study employs the backbone with full I2P transformer blocks as baseline and evaluates the registration accuracy with only the coarse matching scheme or fine matching scheme. For coarse matching only, the matching pairs are established on the coarse level and remapped to the original resolution for pose estimation. By contrast, for fine matching only, matching pairs are established on the fine level directly, without guidance of coarse level correspondences. Experiment results in Table IV show that removing either coarse matching stage or fine matching stage leads to higher registration error and variance. We indicate that coarse-level registration provides robust correspondences and fine-level registration provides accurate matching pairs. Combining the coarse-level and fine-level registration sequentially makes it easier to access the global optimal solution in I2P registration task.

C. Analysis of point cloud density

Given the significant impact of point cloud density on the representation learning process, we conducted ablation studies to examine this influence. The results in Table V present

TABLE IV
ABLATION STUDY ON I2P TRANSFORMER BLOCKS. **CM** DENOTES COARSE MATCHING AND **FM** FINE MATCHING.

Baseline	CM	FM	RRE(°)	RTE(m)
✓	✓		1.58 ± 1.64	0.41 ± 0.39
✓		✓	1.62 ± 2.89	0.44 ± 0.82
✓	✓	✓	1.28 ± 0.94	0.33 ± 0.23

registration accuracy and computational complexity for various point cloud densities. All experimental settings and hyperparameters remained consistent with those outlined in Section IV-A. The findings in Table V reveal that as point cloud density decreases, the qualitative metrics deteriorate, while higher point cloud densities correspond to a sharp increase in computational complexity. It's an intuitive observation that low-density point clouds lose local structured information, and high-density point clouds place a heavy computational burden. As a result, we opted for a compromise and selected 20,480 points, striking a balance between efficiency and accuracy.

TABLE V
ABLATION STUDY ON POINT CLOUD DENSITY.

#Points	RRE(°)	RTE(m)	FLOPs
5120	10.68 ± 31.69	3.42 ± 22.85	37.42G
10240	1.67 ± 2.07	0.42 ± 0.41	56.00G
20480	1.28 ± 0.94	0.33 ± 0.23	93.80G
40960	1.23 ± 1.30	0.32 ± 0.28	171.91G

VI. CONCLUSION

In this paper, we propose CoFiI2P, a novel network for I2P registration. The core is the coarse-to-fine matching strategy, which establishes robust correspondence on the global level first and learns the precise correspondences on the local level progressively. Furthermore, the I2P transformer with self-attention and cross-attention module is introduced to enhance the global-aware ability in homogeneous and heterogeneous data. Compared with one-stage dense prediction and matching approaches, CoFiI2P filters out a large number of false correspondences and holds a safe lead in all the metrics. Extensive experiments on the KITTI dataset have demonstrated that CoFiI2P owns accuracy, robustness, and efficiency in various environments. We hope that the released code of CoFiI2P could motivate the related societies.

In the near future, we intend to design lightweight I2P networks for compact robots and auto-driving vehicles. We believe that achieving real-time regression of transformation parameters is of paramount importance for seamlessly fusing multiple modalities, enhancing navigation capabilities, and enabling downstream perception tasks.

REFERENCES

- [1] J. S. Berrio, M. Shan, S. Worrall, and E. Nebot, "Camera-lidar integration: Probabilistic sensor fusion for semantic mapping," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 7, pp. 7637–7652, 2021.
- [2] L. Nunes, L. Wiesmann, R. Marcuzzi, X. Chen, J. Behley, and C. Stachniss, "Temporal consistent 3d lidar representation learning for semantic perception in autonomous driving," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 5217–5228.
- [3] L. Wang, X. Zhang, W. Qin, X. Li, J. Gao, L. Yang, Z. Li, J. Li, L. Zhu, H. Wang *et al.*, "Camo-mot: Combined appearance-motion optimization for 3d multi-object tracking with camera-lidar fusion," *IEEE Transactions on Intelligent Transportation Systems*, 2023.
- [4] K. Zhao, L. Ma, Y. Meng, L. Liu, J. Wang, J. M. Junior, W. N. Gonçalves, and J. Li, "3d vehicle detection using multi-level fusion from point clouds and images," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 9, pp. 15 146–15 154, 2022.
- [5] C. Lin, D. Tian, X. Duan, J. Zhou, D. Zhao, and D. Cao, "Cl3d: Camera-lidar 3d object detection with point feature enhancement and point-guided fusion," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 10, pp. 18 040–18 050, 2022.
- [6] Y. Peng, Y. Qin, X. Tang, Z. Zhang, and L. Deng, "Survey on image and point-cloud fusion-based object detection in autonomous vehicles," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 12, pp. 22 772–22 789, 2022.
- [7] K. Chen, H. Yu, W. Yang, L. Yu, S. Scherer, and G.-S. Xia, "I2d-loc: Camera localization via image to lidar depth flow," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 194, pp. 209–221, 2022.
- [8] C. Chen, B. Wang, C. X. Lu, N. Trigoni, and A. Markham, "Deep learning for visual localization and mapping: A survey," *arXiv preprint arXiv:2308.14039*, 2023.
- [9] X. Yan, J. Gao, C. Zheng, C. Zheng, R. Zhang, S. Cui, and Z. Li, "2dpass: 2d priors assisted semantic segmentation on lidar point clouds," in *European Conference on Computer Vision*. Springer, 2022, pp. 677–695.
- [10] X. Zhang, S. Zhu, S. Guo, J. Li, and H. Liu, "Line-based automatic extrinsic calibration of lidar and camera," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 9347–9353.
- [11] S. Wang, X. Zhang, G. Zhang, Y. Xiong, G. Tian, S. Guo, J. Li, P. Lu, J. Wei, and L. Tian, "Temporal and spatial online integrated calibration for camera and lidar," in *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2022, pp. 3016–3022.
- [12] C. Yuan, X. Liu, X. Hong, and F. Zhang, "Pixel-level extrinsic self calibration of high resolution lidar and camera in targetless environments," *IEEE Robotics and Automation Letters*, vol. 6, no. 4, pp. 7517–7524, 2021.
- [13] M. Feng, S. Hu, M. H. Ang, and G. H. Lee, "2d3d-matchnet: Learning to match keypoints across 2d image and 3d point cloud," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 4790–4796.
- [14] J. Li and G. H. Lee, "Deepi2p: Image-to-point cloud

- registration via deep classification,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 15 960–15 969.
- [15] S. Ren, Y. Zeng, J. Hou, and X. Chen, “Corri2p: Deep image-to-point cloud registration via dense correspondence,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 33, no. 3, pp. 1198–1208, 2022.
- [16] D. G. Lowe, “Distinctive image features from scale-invariant keypoints,” *International journal of computer vision*, vol. 60, pp. 91–110, 2004.
- [17] Y. Zhong, “Intrinsic shape signatures: A shape descriptor for 3d object recognition,” in *2009 IEEE 12th international conference on computer vision workshops, ICCV Workshops*. IEEE, 2009, pp. 689–696.
- [18] J. Sun, Z. Shen, Y. Wang, H. Bao, and X. Zhou, “Loftr: Detector-free local feature matching with transformers,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 8922–8931.
- [19] W. Jiang, E. Trulls, J. Hosang, A. Tagliasacchi, and K. M. Yi, “Cotr: Correspondence transformer for matching across images,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 6207–6217.
- [20] Z. Qin, H. Yu, C. Wang, Y. Guo, Y. Peng, and K. Xu, “Geometric transformer for fast and robust point cloud registration,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 11 143–11 152.
- [21] H. Wang, Y. Liu, Q. Hu, B. Wang, J. Chen, Z. Dong, Y. Guo, W. Wang, and B. Yang, “Roreg: Pairwise point cloud registration with oriented descriptors and local rotations,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [22] V. Lepetit, F. Moreno-Noguer, and P. Fua, “Ep n p: An accurate o (n) solution to the p n p problem,” *International journal of computer vision*, vol. 81, pp. 155–166, 2009.
- [23] M. A. Fischler and R. C. Bolles, “Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography,” *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981.
- [24] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, “Orb: An efficient alternative to sift or surf,” in *2011 International conference on computer vision*. Ieee, 2011, pp. 2564–2571.
- [25] D. DeTone, T. Malisiewicz, and A. Rabinovich, “Toward geometric deep slam,” *arXiv preprint arXiv:1707.07410*, 2017.
- [26] DeTone, Daniel and Malisiewicz, Tomasz and Rabinovich, Andrew, “Superpoint: Self-supervised interest point detection and description,” in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018, pp. 224–236.
- [27] P.-E. Sarlin, D. DeTone, T. Malisiewicz, and A. Rabinovich, “Superglue: Learning feature matching with graph neural networks,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 4938–4947.
- [28] Q. Zhou, T. Sattler, and L. Leal-Taixe, “Patch2pix: Epipolar-guided pixel-level correspondences,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 4669–4678.
- [29] H. Deng, T. Birdal, and S. Ilic, “Ppf-foldnet: Unsupervised learning of rotation invariant 3d local descriptors,” in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 602–618.
- [30] H. Wang, Y. Liu, Z. Dong, and W. Wang, “You only hypothesize once: Point cloud registration with rotation-equivariant descriptors,” in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 1630–1641.
- [31] Z. Gojcic, C. Zhou, J. D. Wegner, and A. Wieser, “The perfect match: 3d point cloud matching with smoothed densities,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 5545–5554.
- [32] H. Yu, Z. Qin, J. Hou, M. Saleh, D. Li, B. Busam, and S. Ilic, “Rotation-invariant transformer for point cloud matching,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 5384–5393.
- [33] S. Ao, Q. Hu, H. Wang, K. Xu, and Y. Guo, “Buffer: Balancing accuracy, efficiency, and generalizability in point cloud registration,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 1255–1264.
- [34] H. Yu, F. Li, M. Saleh, B. Busam, and S. Ilic, “Cofinet: Reliable coarse-to-fine correspondences for robust point-cloud registration,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 23 872–23 884, 2021.
- [35] Z. J. Yew and G. H. Lee, “Regtr: End-to-end point cloud correspondences with transformers,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 6677–6686.
- [36] S. Huang, Z. Gojcic, M. Usvyatsov, A. Wieser, and K. Schindler, “Predator: Registration of 3d point clouds with low overlap,” in *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, 2021, pp. 4267–4276.
- [37] J. Levinson and S. Thrun, “Automatic online calibration of cameras and lasers,” in *Robotics: science and systems*, vol. 2, no. 7. Citeseer, 2013.
- [38] J. Li, B. Yang, C. Chen, R. Huang, Z. Dong, and W. Xiao, “Automatic registration of panoramic image sequence and mobile laser scanning data using semantic features,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 136, pp. 41–57, 2018.
- [39] Y. Wang, Y. Li, Y. Chen, M. Peng, H. Li, B. Yang, C. Chen, and Z. Dong, “Automatic registration of point cloud and panoramic images in urban scenes based on pole matching,” *International Journal of Applied Earth Observation and Geoinformation*, vol. 115, p. 103083, 2022.
- [40] Y. Liao, J. Li, S. Kang, Q. Li, G. Zhu, S. Yuan, Z. Dong, and B. Yang, “Se-calib: Semantic edges based lidar-camera boresight online calibration in urban scenes,”

- IEEE Transactions on Geoscience and Remote Sensing*, 2023.
- [41] Q. Zhang and R. Pless, “Extrinsic calibration of a camera and laser range finder (improves camera calibration),” in *2004 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)(IEEE Cat. No. 04CH37566)*, vol. 3. IEEE, 2004, pp. 2301–2306.
- [42] R. Unnikrishnan and M. Hebert, “Fast extrinsic calibration of a laser rangefinder to a camera,” *Robotics Institute, Pittsburgh, PA, Tech. Rep. CMU-RI-TR-05-09*, 2005.
- [43] D. Scaramuzza, A. Harati, and R. Siegwart, “Extrinsic self calibration of a camera and a 3d laser range finder from natural scenes,” in *2007 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2007, pp. 4164–4169.
- [44] I. Stamos, L. Liu, C. Chen, G. Wolberg, G. Yu, and S. Zokai, “Integrating automated range registration with multiview geometry for the photorealistic modeling of large-scale scenes,” *International Journal of Computer Vision*, vol. 78, pp. 237–260, 2008.
- [45] G. Pandey, J. McBride, S. Savarese, and R. Eustice, “Automatic targetless extrinsic calibration of a 3d lidar and camera by maximizing mutual information,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 26, no. 1, 2012, pp. 2053–2059.
- [46] M. Clerc and J. Kennedy, “The particle swarm-explosion, stability, and convergence in a multidimensional complex space,” *IEEE transactions on Evolutionary Computation*, vol. 6, no. 1, pp. 58–73, 2002.
- [47] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [48] H. Thomas, C. R. Qi, J.-E. Deschaud, B. Marcotegui, F. Goulette, and L. J. Guibas, “Kpconv: Flexible and deformable convolution for point clouds,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 6411–6420.
- [49] W. Zhang, Z. Huang, G. Luo, T. Chen, X. Wang, W. Liu, G. Yu, and C. Shen, “Topformer: Token pyramid transformer for mobile semantic segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 083–12 093.
- [50] Q. Wan, Z. Huang, J. Lu, G. Yu, and L. Zhang, “Seaformer: Squeeze-enhanced axial transformer for mobile semantic segmentation,” *arXiv preprint arXiv:2301.13156*, 2023.
- [51] K. Wu, H. Peng, M. Chen, J. Fu, and H. Chao, “Re-thinking and improving relative position encoding for vision transformer,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10033–10041.
- [52] Y. Sun, C. Cheng, Y. Zhang, C. Zhang, L. Zheng, Z. Wang, and Y. Wei, “Circle loss: A unified perspective of pair similarity optimization,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 6398–6407.
- [53] A. Geiger, P. Lenz, and R. Urtasun, “Are we ready for autonomous driving? the kitti vision benchmark suite,” in *2012 IEEE conference on computer vision and pattern recognition*. IEEE, 2012, pp. 3354–3361.
- [54] B. Wang, C. Chen, Z. Cui, J. Qin, C. X. Lu, Z. Yu, P. Zhao, Z. Dong, F. Zhu, N. Trigoni *et al.*, “P2-net: Joint description and detection of local features for pixel and point matching,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 16 004–16 013.
- [55] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [56] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.