
Incremental Multiview Point Cloud Registration with Two-stage Candidate Retrieval

Shiqi Li Jihua Zhu Yifan Xie
Xi'an Jiaotong University

Mingchen Zhu
UC Davis

Abstract

Multiview point cloud registration serves as a cornerstone of various computer vision tasks. Previous approaches typically adhere to a global paradigm, where a pose graph is initially constructed followed by motion synchronization to determine the absolute pose. However, this separated approach may not fully leverage the characteristics of multiview registration and might struggle with low-overlap scenarios. In this paper, we propose an incremental multiview point cloud registration method that progressively registers all scans to a growing meta-shape. To determine the incremental ordering, we employ a two-stage coarse-to-fine strategy for point cloud candidate retrieval. The first stage involves the coarse selection of scans based on neighbor fusion-enhanced global aggregation features, while the second stage further reranks candidates through geometric-based matching. Additionally, we apply a transformation averaging technique to mitigate accumulated errors during the registration process. Finally, we utilize a Reservoir sampling-based technique to address density variance issues while reducing computational load. Comprehensive experimental results across various benchmarks validate the effectiveness and generalization of our approach.

1 Introduction

Point cloud registration constitutes a pivotal challenge within the domains of computer vision and robotics. Predominantly, pairwise registration emerges as the prevalent mode in point cloud registration, having witnessed significant strides in recent years [1; 2; 3; 4]. Nonetheless, a single pair of scans typically fails to encapsulate expansive scenes adequately. Consequently, the necessity arises for multiview point cloud registration, aimed at offering a comprehensive representation of the scenario by integrating data from multiple scans. The resultant complete scene reconstruction facilitates various applications including simultaneous localization and mapping (SLAM), autonomous driving, and Embodied AI. Despite the notable progress in pairwise registration, multiview point cloud registration garners comparatively less attention and presents persistent challenges in achieving satisfactory performance.

The majority of multiview point cloud registration methodologies adhere to a global approach, initially constructing a pose graph through pairwise registration and subsequently employing motion synchronization to ascertain the pose of individual point cloud frames. Recent progress in global multiview point cloud registration encompasses the integration of more robust pairwise registration techniques [5; 6], the formulation of diverse strategies for pose graph construction [7], and the innovation of optimization-based [8] or learning-based [9] motion synchronization methodologies.

Although these studies deliver some promising results, the global multiview registration paradigm presents inevitable drawbacks. The objective of motion synchronization is to determine a set of absolute poses \mathbf{T}_V that minimize consistency errors with the relative transformation $\tilde{\mathbf{T}}$ in pose graph

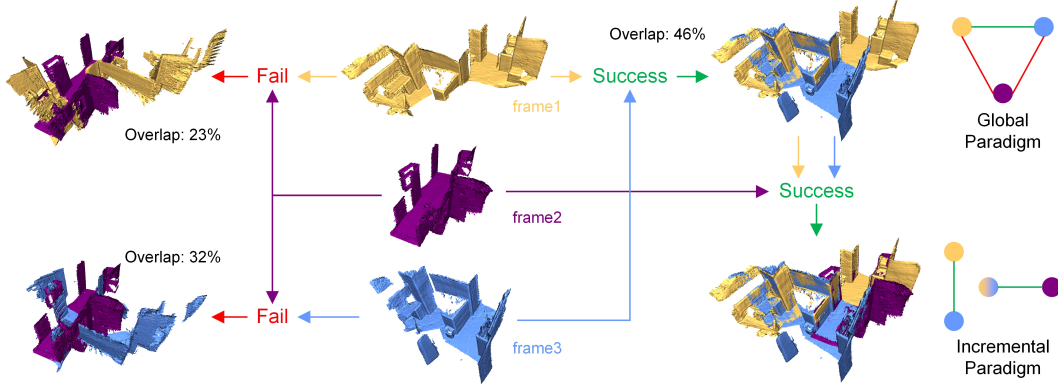


Figure 1: An example comprising **frame1**, **frame2**, and **frame3** within an indoor aisle. While frame1 and frame3 share some overlapping areas with frame2, these common areas predominantly consist of flat floors, which offer limited cues for registration. However, frame1 and frame3 share sufficient structurally significant areas that enable successful alignment. The motion synchronization mechanism cannot handle the pose graph, as depicted in the top right corner, due to the absence of a correct connection to the frame2. Nevertheless, our incremental paradigm first merges frame1 and frame3 to expand the meaningful overlapping areas with frame2, ultimately achieving a complete scene, as illustrated in the bottom right corner.

$$\mathcal{G} = \{\mathcal{V}, \mathcal{E}\},$$

$$\mathbf{T}_{\mathcal{V}} = \arg \min_{\{\mathbf{T}_1, \mathbf{T}_2, \dots, \mathbf{T}_{|\mathcal{V}|}\}} \sum_{e_{ij} \in \mathcal{E}} d(\tilde{\mathbf{T}}_{ij} - \mathbf{T}_i \mathbf{T}_j^{\top}), \quad (1)$$

where $d(\cdot)$ is a distance measure between two transformations. However, addressing Eq. 1 may not consistently yield satisfactory outcomes. The accuracy of the estimated absolute pose depends on both the noise level of pairwise measurements and the Laplacian of the pose graph structure [10]. While employing diverse graph construction and motion synchronization strategies can approximate certain lower error bounds, these bounds are significantly impacted by the accuracy and robustness of the pairwise method. For instance, as illustrated in Fig. 1, if a frame within the unorganized set lacks sufficient overlapping area with all neighbors, global multiview methods may struggle to handle it.

The incremental method represents an alternative paradigm for multiview point cloud registration. As its name suggests, incremental methods iteratively perform pairwise registration and expand the point cloud to acquire the absolute pose of each frame. In contrast to the global paradigm, the incremental method is relatively less efficient but holds the potential to yield more accurate results [11]. It is intuitively expected that as the point cloud iteratively grows, the overlap ratio between two frames will not be lower than that achieved through pairwise registration in the pose graph construction of the global paradigm. This characteristic may alleviate the challenges posed by low overlap issues, which are problematic for pairwise registration methods. Despite its significant potential, incremental approaches have not been fully explored; existing works primarily rely on handcrafted growth strategies [12] or additional prior assumptions [13], limiting their applicability across a wide range of scenarios.

In this paper, we introduce a novel incremental multiview point cloud registration pipeline that achieves more accurate and robust estimation. Drawing inspiration from modern image retrieval [14] and visual place recognition systems [15], we propose a two-stage coarse-to-fine point cloud retrieval module to address the critical ordering problem inherent in the incremental process. The first coarse stage is instantiated with deep learning-based global semantic features aggregation to identify potential candidate frames, while the subsequent fine stage employs geometric matching to rerank the candidates and select the best frame for point cloud expansion. Given the issue of accumulated error during the growing process, we mitigate this challenge by employing single transformation averaging. The refinement of the final pose for each point cloud is achieved by leveraging all neighbor frames that exhibit sufficient overlapping. Finally, considering the computation cost and density variance caused by the point cloud growth, we propose a Reservoir sampling-based [16] strategy to maintain the keypoints and descriptors of the meta point cloud. We evaluate the effectiveness of our approach on 3D(Lo)Match [17; 1] and ScanNet [18] benchmarks. Experimental results demonstrate

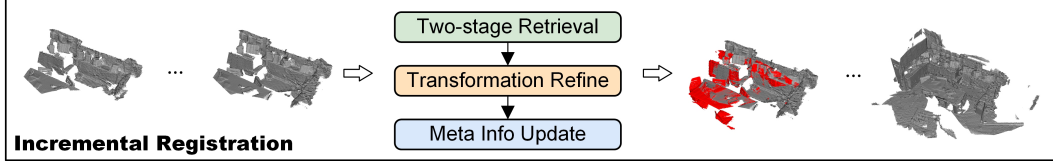


Figure 2: Overview of our pipeline: Point cloud frames are progressively incorporated into the same coordinate system, referred to as the meta-shape in this study. For each frame, the process entails two-stage frame retrieval, transformation refinement, and meta-information update.

the superiority of our design, with our method significantly outperforming state-of-the-art methods, achieving registration recalls of 97.1% and 87.9% on the 3DMatch/3DLoMatch datasets, respectively.

2 Related Work

Pairwise Registration Pairwise point cloud registration, which forms the basis of multiview registration, can be broadly classified into two types. The first category encompasses correspondence-based approaches [19; 2; 1; 3], which initially establish a series of tentative correspondences and then employ a robust estimator [20; 21; 22] or geometric hashing [23] to determine the transformation. Notably, the Iterative Closest Point (ICP) [24] series, as representative methods, establish correspondences based on various distances in Euclidean space, albeit they are susceptible to initialization and outliers. To enhance robustness, handcrafted features have been proposed to characterize the geometric nature of point clouds in a local area [25], and correspondences are established through feature matching. Recent advancements in deep learning tend to replace handcrafted features with learnable descriptors [19; 26]. Metric learning-guided neural networks automatically generate similar descriptors for consistent areas and distinct descriptors for non-overlapping areas, with cross-attention mechanisms typically applied to facilitate explicit information exchange [1; 3; 5]. The second category comprises correspondence-free methods [27; 28], which aim to directly regress the transformation from the point cloud pair instead of establishing correspondences. Although these methods have shown promising results on synthetic object-level shapes, they struggle with large-scale real-world scenes. In this study, we utilize off-the-shelf learnable correspondence-based pairwise registration methods as a plug-and-play module in our multiview registration pipeline.

Multiview Registration Multiview point cloud registration endeavors to reconstruct a holistic scene from a set of unorganized, partially overlapping point cloud scans. Analogous to the taxonomy of structure-from-motion (SfM), multiview registration can be classified into global and incremental methods. Global methods initially collect relative transformations to construct a pose graph, followed by utilizing global cycle consistency to optimize absolute poses [29; 30]. Recognizing that fully connected graphs may encompass numerous outliers, a pruning strategy is introduced to eliminate low-quality registration pairs [31]. Additionally, a more efficient approach in [7] directly builds a sparse graph. Motion synchronization aims to recover global poses with minimal consistency error, as mentioned in Eq. 1. Iteratively reweighted least-squares (IRLS) based schemes are prevalent in motion averaging [32; 8; 31], gradually reducing the weights of outlier edges and assigning higher weights to inlier edges by optimizing a robust loss function [33]. However, IRLS may be susceptible to local minima when undesired robust loss or reweighting schemes are employed. Beyond IRLS, recent approaches address the synchronization problem in a data-driven manner [9; 34], often projecting transformations into a high-dimensional latent space and utilizing graph neural networks (GNN) to facilitate information interaction among the pose graph. In contrast, incremental methods iteratively merge a new point cloud frame into the growing meta-point cloud until all scans are incorporated into the final scene point cloud [12]. The incremental paradigm finds more popularity in the fields of geoscience and remote sensing [35; 13], particularly given the extensive overlap typically present in terrestrial laser scans between successive frames.

3 Method

Given a set of unorganized partial overlap point cloud scans $\mathcal{P} = \{P_i | i = 1, 2, \dots, N\}$. The multiview registration aims to seek an absolute pose $\mathbf{T}_i = (\mathbf{R}_i, \mathbf{t}_i) \in SE(3)$ for each point cloud

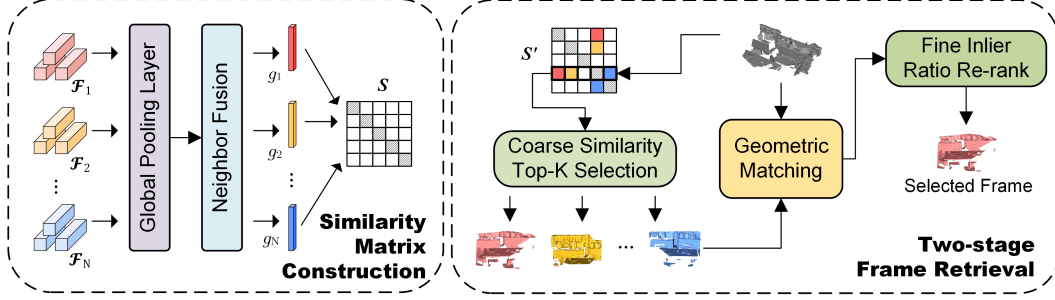


Figure 3: Left: Illustration depicting the construction of the similarity matrix. Right: Workflow of the proposed two-stage frame retrieval process.

frame P_i that correctly recovers the complete scene in a unified coordinate system. In this study, we employ local descriptor-based methods [19; 26] to align point cloud pairs. Consequently, for each point cloud P_i , we can extract a set of keypoints along with their associated local descriptors \mathcal{F}_i . The transformation is estimated based on the correspondences established through descriptor matching.

The overview of our proposed method is illustrated in Fig. 2.

3.1 Coarse-to-Fine Frame Selection

Determining the order of registration sequences is crucial for our incremental method. We adopt a two-stage coarse-to-fine approach to accurately select the newly added point cloud frame. In the first stage, point clouds similar to the meta-shape are coarsely retrieved using global features aggregated by a deep learning method. Then, the selected candidates undergo further reranking through a second stage based on geometric matching. Subsequently, the best matching frame is registered into the meta-shape, as shown in Fig. 3.

Our first stage is instantiated with a global pooling layer (e.g. NetVLAD [36], GeM [37]). This layer is responsible for aggregating a global feature for each point cloud frame from the associated local descriptor set \mathcal{F}_i ,

$$\tilde{g}_i = \text{Pooling}(\mathcal{F}_i). \quad (2)$$

To enhance the effectiveness of the vanilla global pooling features, we further introduce a fusion operation [14]. Formally, each global feature is refined using its neighboring features,

$$g_i = \frac{\tilde{g}_i + \sum_{j \in \mathcal{M}_i} \beta(\tilde{g}_j \cdot \tilde{g}_i) \tilde{g}_j}{1 + \sum_{j \in \mathcal{M}_i} \beta(\tilde{g}_j \cdot \tilde{g}_i)}, \quad (3)$$

where \mathcal{M}_i represents indexes of the nearest neighbors of \tilde{g}_i . This fusion operation facilitates the generation of similar global features for overlapping point clouds.

We construct the similarity matrix $S \in \mathbb{R}^{N \times N}$ by mapping the distances between any two global features to $[0, 1]$,

$$s_{ij} = \frac{2 - \|g_i - g_j\|_2}{2}. \quad (4)$$

With the matrix S , we initialize the meta-shape as the frame that is most similar to the others,

$$x = \arg \max_{1 \leq i \leq N} \sum_{j=1}^N s_{ij}, \quad (5)$$

The keypoints and descriptors of point cloud x are utilized to initialize the meta-shape, while the pose of point cloud x is set to the identity matrix \mathbf{I} .

After determining the initial meta-shape, we select the top- k similar point cloud frames to the meta-shape based on their scores in S as candidate frames. The procedure outlined here serves as an example using the initial meta-shape; however, the same approach is employed for arbitrary meta-shapes during the incremental process. Additionally, instead of employing the pooling layer to

aggregate a new global descriptor at each step, the global similarity matrix S is dynamically updated throughout the incremental registration process. Details of the global similarity matrix update will be provided in Sec. 3.3.

Recognizing that the global feature-based retrieval in the first stage may inadvertently lose some subtle patterns during the aggregation process, we further introduce a second stage based on geometric matching to rerank the candidates. Specifically, for each candidate frame, we establish a coarse correspondence set \mathcal{C} with the meta-shape using descriptor matching. Subsequently, a robust estimator (e.g. RANSAC) is employed to estimate the transformation \mathbf{T} between the meta-shape and the candidate frame. Finally, we compute the inlier count (IC) for each candidate frame, and the frame with the highest IC is selected to be merged into the meta-shape. IC measures the number of correspondences that confirm the estimated transformation, where a high IC value supports reliable registration. The IC can be calculated as:

$$\text{IC} = \sum_{(p,q) \in \mathcal{C}} \llbracket \|\mathbf{T}(p) - q\|_2 < \tau \rrbracket, \quad (6)$$

where $\llbracket \cdot \rrbracket$ is Iverson bracket, τ is a distance threshold.

3.2 Transformation Refinement

The selection process identifies the most suitable frame and provides an estimated transformation \mathbf{T} . However, we do not solely rely on the transformation from the robust estimator; instead, we further refine it with single translation averaging. For each frame in the meta-shape, we calculate the overlap ratio with the selected frame using:

$$o_{ij} = \frac{\sum_{p \in P_i} \llbracket \|\mathbf{T}(p) - \text{NN}(\mathbf{T}(p), P_j)\|_2 < \tau \rrbracket + \sum_{q \in P_j} \llbracket \|q - \text{NN}(q, \mathbf{T}(P_i))\|_2 < \tau \rrbracket}{|P_i| + |P_j|}, \quad (7)$$

where $\text{NN}(\cdot)$ denotes the spatial nearest neighbor.

For frames that contain more than 30% overlap with the selected frame, we solve the Orthogonal Procrustes problem to compute a new transformation. If the selected frame yields more than one new transformation, we apply a single transformation averaging method to average them. Otherwise, the transformation obtained from the robust estimator is directly employed.

The transformation averaging process comprises rotation averaging and translation averaging. For rotation averaging, we utilize a simplified single rotation averaging method from [38] to compute the averaged rotation matrix. At the beginning, the rotation \mathbf{R} from the robust estimator is set as the initial averaging estimation. Subsequently, both \mathbf{R} and the rotations $\hat{\mathbf{R}}$ from the newly observed transformations are converted into vectors. Next, we gather the residual vectors v between \mathbf{R} and $\hat{\mathbf{R}}$ in vector form and calculate associated norms. The averaging estimation is then updated to a weighted sum of residual vectors, where the weights correspond to the overlapping ratio in Eq. 7. This procedure is iterated until the averaging estimation converges or reaches the maximum iteration rounds. Finally, the averaging estimation vector is converted into matrix form and projected onto the $SO(3)$ space using SVD to obtain the final $\bar{\mathbf{R}}$. A detailed algorithm can be found in the Appendix A.

Based on the rotation averaging result, the translation averaging can be formulated as follows:

$$\bar{\mathbf{t}} = \arg \min_{\mathbf{t} \in \mathbb{R}^{3 \times 1}} \sum_{i=1}^M w_i \|\hat{\mathbf{t}}_i - \hat{\mathbf{R}}_i \bar{\mathbf{R}}^\top \mathbf{t}\|_2, \quad (8)$$

where M the number of newly observed transformation, and w_i is the associated overlapping ratio. This problem can be solved using the least squares method. We construct three zero block matrices: $\mathbf{A} \in \mathbb{R}^{3M \times 3}$, $\mathbf{B} \in \mathbb{R}^{3M \times 1}$, and $\mathbf{W} \in \mathbb{R}^{3M \times 3M}$, where \mathbf{A} and \mathbf{W} consist of 3×3 blocks, and \mathbf{B} consists of 3×1 blocks. For each newly observed transformation, the i -th block on the diagonal of \mathbf{W} is set to $w_i \mathbf{I}$, while the i -th blocks in \mathbf{A} and \mathbf{B} are set to $\hat{\mathbf{R}}_i \bar{\mathbf{R}}^\top$ and $\hat{\mathbf{t}}_i$, respectively. The final averaged translation can then be calculated as:

$$\bar{\mathbf{t}} = (\mathbf{A}^\top \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{W} \mathbf{B}. \quad (9)$$

The result of the transformation averaging $\bar{\mathbf{T}} = (\bar{\mathbf{R}}, \bar{\mathbf{t}})$ is utilized to determine the pose of the selected frame. It is worth noting that to be consistent with the coordinate system definition commonly used in the registration field the final absolute pose is the inverse of the calculated transformation.

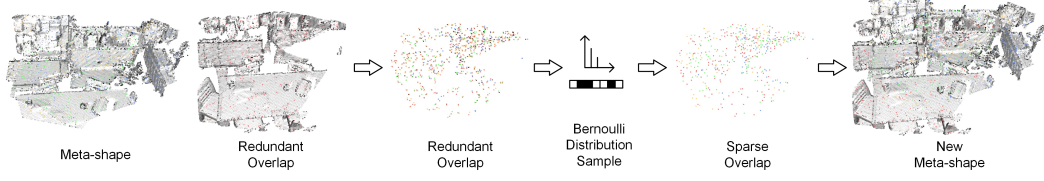


Figure 4: An illustration of the meta information update process.

3.3 Meta Information Update

After obtaining the pose of the selected frame, we utilize it to update the meta-shape. The simplest method involves directly adding the keypoints and descriptors of the selected frame to the meta-shape. However, this straightforward merging operation introduces numerous redundant points and descriptors in the overlapping areas, resulting in two drawbacks. Firstly, it increases memory usage and computational costs. Secondly, this approach causes density variance in the meta-shape, which adversely affects feature matching and registration processes.

To mitigate redundancy, we employ a sampling strategy to update the keypoints and descriptors, as illustrated in Fig. 4. Initially, we use the estimated pose to align the selected point cloud frame with the coordinate system of the meta-shape. Subsequently, we identify mutual nearest point pairs in Euclidean space between the meta-shape and the selected frame, considering pairs with distances below the threshold τ as redundant. Next, we randomly sample p_{k_1, k_2, \dots, k_n} from a multivariate Bernoulli distribution $P\{X_1 = k_1, X_2 = k_2, \dots, X_n = k_n\}$ to determine the source of the keypoints and descriptors in the overlapping area. Here, the number of variables n equals the count of redundant pairs, with each variable X_i following an independent Bernoulli distribution whose probability mass function can be expressed as (derivation in Appendix B):

$$P(X_i = k) = \frac{r_i^{1-k}}{r_i + 1}, k \in \{0, 1\}, \quad (10)$$

where r_i represents the redundancy count for point in pair i . This Reservoir sampling-based [16] design ensures that overlapping points covered by multiple point cloud frames are uniformly sampled, thus yielding a high-quality meta-shape for multiview registration. For each redundant pair, if the corresponding variable equals 0, the keypoint and descriptor from the meta-shape are used; otherwise, those from the selected frame are utilized. The points and descriptors from the sampled redundant pairs are then combined with other non-pairing ones to form the new meta-shape.

In addition to updating the keypoints and descriptors, the global similarity matrix S is also updated. We utilize the row corresponding to the selected frame to update the row of the meta-shape. The similarity is set to the maximum value between the two rows, except the frame already present in the meta-shape, which is set to 0.

Ultimately, the updated meta-shape and matrix S revert to the two-stage selection process outlined in Sec. 3.1 to incorporate a new frame until all frames are included in the meta-shape.

3.4 Loss

To train the coarse stage in our selection module, we employ the smooth L1 loss [39] between the predicted similarity score s_{ij} and the ground truth overlapping ratio.

4 Experiments

4.1 Dataset

The 3DMatch [17] is an indoor dataset commonly utilized for registration tasks. It comprises 46 training scenes, 8 validation scenes, and 8 testing scenes. Our evaluation adheres to the standard practice of assessing both the 3DMatch and 3DLoMatch [1] test splits. Notably, the overlap of point cloud pairs in 3DMatch exceeds 30%, whereas in 3DLoMatch, it ranges between 10% and 30%.

The ScanNet [18] dataset comprises RGB-D data collected from over 1500 indoor scenes. Following the operation outlined in [31; 7], we evaluate a subset of 32 scenes, consisting of 960 scans and a total of 13920 pairs. To generate point clouds, we randomly select 30 RGB-D images from each scene, spaced 20 frames apart, and convert them into point clouds. Notably, the temporal sequence of frames is disregarded during this process. The substantial temporal gap between frames, coupled with the omission of temporal sequence, renders the test setting exceptionally challenging.

4.2 Baseline and Protocol

We compare our method with EIGSE3 [32], L1-IRLS [33], RotAvg [33], LMVR [31], LITS [9], HARA [40], and SGHR [7]. EIGSE3 employs spectral decomposition for motion synchronization and integrates IRLS to enhance robustness. L1-IRLS and RotAvg are classical rotation averaging methods utilizing IRLS with ℓ_1 and $\ell_{1/2}$ loss functions, respectively. LMVR is an end-to-end learning approach that combines pairwise registration and motion synchronization within a single network. LITS is a motion synchronization method leveraging GNN for learning. HARA utilizes a cycle consistency-based spanning tree to reject outliers in rotation averaging. SGHR is a multiview registration method based on sparse pose graphs and IRLS transformation synchronization. We evaluate the multiview performance of all methods using the FCGF [19] and YOHO [26] as pairwise registration methods, with the exception of LMVR, which includes its own pairwise module.

For global paradigm methods, we offer three types of pose graphs: full, pruned, and sparse. "Full" denotes the exhaustive registration of all scan pairs. "Pruned" entails retaining only those scan pairs where the median point distance in the registered overlapping region is less than 0.05m from the fully connected graph [31]. "Sparse" involves selecting pairs based solely on the technique outlined in [7].

4.3 Metrics

Following the protocol in [31; 7], we assess multiview registration by utilizing relative transformations calculated from the recovered absolute poses.

For the 3D(Lo)Match datasets, we utilize the Registration Recall (RR) metric. RR measures the fraction of point cloud pairs with a transformation error smaller than a specified threshold (0.2m).

For ScanNet, we present the mean, median, and empirical cumulative distribution functions (eCDF) of the rotation error re and translation error te . The re and te are defined as:

$$re = \arccos\left(\frac{\text{tr}(\mathbf{R}^{pred\top}\mathbf{R}^{gt}) - 1}{2}\right), te = \|\mathbf{t}^{pred} - \mathbf{t}^{gt}\|_2. \quad (11)$$

4.4 Results

The quantitative results for 3DMatch, 3DLoMatch, and ScanNet are presented in Tab. 1, Tab. 2, and Tab. 3, respectively. Our method consistently demonstrates outstanding performance across all three benchmarks, irrespective of the pairwise registration methods used. For instance, our method achieves 97.1% and 87.9% on 3DMatch and 3DLoMatch when integrated with YOHO, surpassing the previous best by 0.9 and 5.6 percentage points, respectively. When adapted to ScanNet, our method continues to exhibit strong performance, further underscoring the effectiveness and generalizability of our design. Fig. 5 provides some qualitative results.

4.5 Ablation Studies

To verify the effectiveness of the proposed components, we conducted ablation studies on the 3D(Lo)Match dataset. Tab. 4 presents a summary of the ablation results. In Exp. (a), we tested the influence of coarse semantic selection by using only geometric clues to determine the growth ordering. In Exp. (b), we removed geometric reranking and relied solely on global features to determine the incremental sequence. Additionally, in Exp. (c), we used an overlap-based criterion from [12]. These experiments demonstrate the necessity of our two-stage design. In Exp. (d), we verified the effectiveness of our transformation averaging design. Finally, we conducted two experiments to validate the effectiveness of our meta-information update process. In Exp. (e), we simply merged newly registered frames by concatenation, while in Exp. (f), we calculated the mean of paired keypoints and descriptors instead of sampling.

Table 1: Registration recall on 3DMatch dataset. \heartsuit Full graph, \clubsuit Pruned graph, \diamond Sparse graph.

#Samples	FCGF					YOHO				
	5000	2500	1000	500	250	5000	2500	1000	500	250
EIGSE3 \heartsuit	13.2	16.1	13.9	8.8	7.6	16.6	12.6	12.4	9.8	7.6
EIGSE3 \clubsuit	60.3	62.5	66.5	68.5	73.9	43.6	45.7	50.3	49.9	53.7
EIGSE3 \diamond	45.9	46.8	40.1	37.7	35.3	58.9	54.5	39.2	36.6	25.4
IRLS \heartsuit	47.0	44.4	47.8	44.3	35.5	41.9	41.3	45.8	45.5	33.6
IRLS \clubsuit	77.6	77.9	81.5	83.1	81.0	67.6	68.9	72.3	72.8	80.6
IRLS \diamond	81.2	80.5	80.4	76.7	77.9	83.1	80.1	77.9	80.1	70.8
RotAvg \heartsuit	62.0	65.3	69.8	59.7	61.0	66.6	65.6	68.3	60.8	56.5
RotAvg \clubsuit	86.5	84.8	84.6	86.9	82.5	77.3	78.0	79.0	79.2	80.9
RotAvg \diamond	83.5	83.5	84.8	80.4	82.7	88.5	87.3	82.2	81.1	76.8
LITS \heartsuit	72.8	77.7	76.8	73.0	67.5	78.4	68.1	77.3	71.2	67.2
LITS \clubsuit	82.3	83.2	83.8	85.2	78.7	82.5	82.5	86.5	83.6	83.4
LITS \diamond	73.3	70.4	67.8	66.5	70.4	74.2	70.3	81.3	75.1	70.1
HARA \heartsuit	83.5	83.7	85.2	79.8	77.8	78.8	79.1	80.8	79.6	76.7
HARA \clubsuit	88.2	91.7	89.7	88.9	89.1	83.6	84.7	89.7	88.9	85.7
HARA \diamond	87.7	88.0	86.9	85.9	81.9	88.0	88.3	91.5	88.4	80.0
SGHR \heartsuit	91.5	90.5	90.0	89.8	86.6	93.2	92.0	89.6	<u>92.8</u>	80.0
SGHR \clubsuit	<u>93.9</u>	<u>94.2</u>	<u>91.8</u>	<u>91.7</u>	<u>91.8</u>	95.2	93.8	92.7	89.9	<u>90.8</u>
SGHR \diamond	92.6	90.5	91.3	90.4	87.0	<u>96.2</u>	<u>95.7</u>	<u>95.3</u>	92.2	90.7
Ours	95.4	94.3	93.8	93.6	92.8	97.1	95.4	95.6	94.7	93.4

Table 2: Registration recall on 3DLoMatch dataset. \heartsuit Full graph, \clubsuit Pruned graph, \diamond Sparse graph.

#Samples	FCGF					YOHO				
	5000	2500	1000	500	250	5000	2500	1000	500	250
EIGSE3 \heartsuit	5.7	6.3	4.9	1.7	1.0	6.4	6.0	5.4	1.7	1.0
EIGSE3 \clubsuit	43.0	45.5	47.6	51.8	51.5	33.8	35.0	35.7	34.4	37.2
EIGSE3 \diamond	29.8	30.3	25.5	21.7	20.9	39.8	36.9	26.1	24.5	13.0
IRLS \heartsuit	35.7	30.7	32.9	25.4	20.1	25.9	27.2	28.0	27.0	17.4
IRLS \clubsuit	60.9	64.3	61.8	64.6	60.8	52.3	47.4	53.7	57.0	59.5
IRLS \diamond	62.0	57.7	54.6	54.9	47.2	56.8	56.5	54.0	54.3	44.0
RotAvg \heartsuit	47.0	49.8	54.1	43.0	42.4	45.0	46.0	46.1	45.8	37.2
RotAvg \clubsuit	71.4	72.8	71.8	73.4	64.8	58.6	59.2	62.7	64.9	64.2
RotAvg \diamond	63.2	66.7	64.9	60.0	56.4	63.2	63.3	58.6	62.6	53.4
LITS \heartsuit	58.9	61.0	59.1	52.5	46.2	62.8	51.3	59.9	50.7	44.8
LITS \clubsuit	67.3	68.8	68.3	66.5	64.2	66.0	68.0	68.3	66.2	62.7
LITS \diamond	44.0	39.2	40.2	37.4	37.3	45.0	39.9	51.3	46.3	35.6
HARA \heartsuit	64.5	67.3	65.4	66.7	59.8	60.5	60.7	65.7	62.0	60.6
HARA \clubsuit	75.9	80.4	76.8	74.3	<u>74.7</u>	66.6	67.7	74.2	76.2	69.5
HARA \diamond	68.8	71.3	70.6	66.5	55.9	68.9	68.8	74.1	68.0	56.6
SGHR \heartsuit	80.6	79.5	75.5	77.0	69.1	76.8	77.5	74.4	73.0	63.6
SGHR \clubsuit	82.8	<u>81.1</u>	<u>80.2</u>	<u>79.5</u>	73.1	<u>82.3</u>	76.6	79.7	<u>76.1</u>	<u>74.1</u>
SGHR \diamond	79.5	78.8	77.6	77.4	70.0	81.6	<u>80.5</u>	<u>80.9</u>	76.0	70.6
Ours	<u>81.8</u>	81.6	84.6	81.5	76.7	87.9	87.1	86.3	82.9	77.6

Table 3: Registration results on Scannet dataset with YOHO #Samples=5000 setting. \heartsuit Full graph, \clubsuit Pruned graph, \diamond Sparse graph.

Method	Rotation Error					Translation Error(m)					
	3°	5°	10°	30°	45°	Mean/Med	0.05	0.1	0.25	0.5	0.75
LMVR \heartsuit	48.3	53.6	58.9	63.2	64.0	48.1°/33.7°	34.5	49.1	58.5	61.6	63.9
EIGSE3 \heartsuit	19.7	24.4	32.3	49.3	56.9	53.6°/48.0°	11.2	19.7	30.5	45.7	56.7
EIGSE3 \clubsuit	40.8	46.3	51.9	61.2	65.7	40.6°/37.1°	23.9	38.5	51.0	59.3	66.1
L1-IRLS \heartsuit	38.1	44.2	48.8	55.7	56.5	53.9°/47.1°	18.5	30.4	40.7	47.8	54.4
L1-IRLS \clubsuit	46.3	54.2	61.6	64.3	66.8	41.8°/34.0°	24.1	38.5	48.3	55.6	60.9
RotAvg \heartsuit	44.1	49.8	52.8	56.5	57.3	53.1°/44.0°	28.2	40.8	48.6	51.9	56.1
RotAvg \clubsuit	50.2	60.1	65.3	66.8	68.8	38.5°/31.6°	31.8	49.0	58.8	63.3	65.6
LITS \heartsuit	52.8	67.1	74.9	77.9	79.5	26.8°/27.9°	29.4	51.1	68.9	75.0	77.0
LITS \clubsuit	54.3	69.4	75.6	78.5	80.3	24.9°/19.9°	31.4	54.4	72.3	76.7	79.6
HARA \heartsuit	54.9	64.3	71.3	74.1	74.2	32.1°/29.2°	35.8	54.4	66.3	69.7	72.9
HARA \clubsuit	55.7	63.7	69.0	70.8	72.1	34.7°/31.3°	35.2	53.6	65.4	68.6	71.7
SGHR \heartsuit	57.2	68.5	75.1	78.1	78.8	26.4°/19.5°	39.4	61.5	72.0	75.2	77.6
SGHR \clubsuit	59.4	71.9	80.0	82.1	82.6	21.7°/19.1°	39.9	63.0	74.3	77.6	80.2
SGHR \diamond	59.1	73.1	80.8	82.5	83.0	21.7°/19.0°	39.9	64.1	76.7	79.0	81.9
Ours	59.9	74.0	85.1	87.9	88.5	17.8°/13.7°	41.8	63.1	79.3	83.1	85.5

Table 4: Ablation experiments. SS: Semantic stage. GS: Geometric stage. OR: Overlap-based retrieval. TR: Transformation refinement. RS: Reservoir sampling. MU: Mean update. Tested with YOHO #Samples=1000 setting.

	(a) w/o SS	(b) w/o GS	(c) w/ OR	(d) w/o TR	(e) w/o RS	(f) w/ MU	Full
3DM	94.0	91.3	87.9	94.4	94.7	95.6	95.6
3DLo	82.3	68.9	64.6	81.2	82.1	83.0	86.3

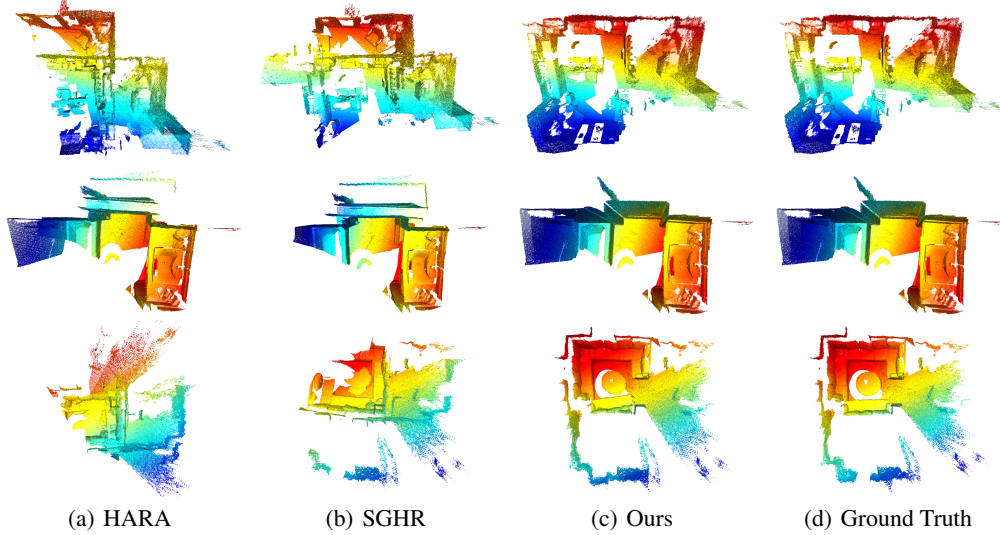


Figure 5: Qualitative comparison results.

5 Conclusion

In this paper, we propose a novel incremental multiview point cloud registration method. Our approach utilizes a two-stage coarse-to-fine point cloud retrieval process that combines the advantages of both semantic and geometric clues. Additionally, a single transformation averaging scheme and a Reservoir sampling-based update strategy further enhance robustness while reducing computational costs. Experiments on various datasets demonstrate the effectiveness of our proposed method.

References

- [1] Shengyu Huang, Zan Gojcic, Mikhail Usvyatsov, Andreas Wieser, and Konrad Schindler. Predator: Registration of 3d point clouds with low overlap. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, pages 4267–4276, 2021.
- [2] Hao Yu, Fu Li, Mahdi Saleh, Benjamin Busam, and Slobodan Ilic. Cofinet: Reliable coarse-to-fine correspondences for robust pointcloud registration. *Advances in Neural Information Processing Systems*, 34:23872–23884, 2021.
- [3] Zheng Qin, Hao Yu, Changjian Wang, Yulan Guo, Yuxing Peng, and Kai Xu. Geometric transformer for fast and robust point cloud registration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11143–11152, 2022.
- [4] Fan Yang, Lin Guo, Zhi Chen, and Wenbing Tao. One-inlier is first: Towards efficient position encoding for point cloud registration. *Advances in Neural Information Processing Systems*, 35:6982–6995, 2022.
- [5] Junle Yu, Luwei Ren, Yu Zhang, Wenhui Zhou, Lili Lin, and Guojun Dai. Peal: Prior-embedded explicit attention learning for low-overlap point cloud registration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17702–17711, 2023.
- [6] Shengze Jin, Iro Armeni, Marc Pollefeys, and Daniel Barath. Multiway point cloud mosaicking with diffusion and global optimization. *arXiv preprint arXiv:2404.00429*, 2024.
- [7] Haiping Wang, Yuan Liu, Zhen Dong, Yulan Guo, Yu-Shen Liu, Wenping Wang, and Bisheng Yang. Robust multiview point cloud registration with reliable pose graph initialization and history reweighting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9506–9515, 2023.
- [8] Xiangru Huang, Zhenxiao Liang, Xiaowei Zhou, Yao Xie, Leonidas J Guibas, and Qixing Huang. Learning transformation synchronization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8082–8091, 2019.
- [9] Zi Jian Yew and Gim Hee Lee. Learning iterative robust transformation synchronization. In *2021 International Conference on 3D Vision (3DV)*, pages 1206–1215. IEEE, 2021.
- [10] Nicolas Boumal, Amit Singer, P-A Absil, and Vincent D Blondel. Cramér-rao bounds for synchronization of rotations. *Information and Inference: A Journal of the IMA*, 3(1):1–39, 2014.
- [11] Xiang Gao, Lingjie Zhu, Zexiao Xie, Hongmin Liu, and Shuhan Shen. Incremental rotation averaging. *International Journal of Computer Vision*, 129:1202–1216, 2021.
- [12] Yulan Guo, Ferdous Sohel, Mohammed Bennamoun, Jianwei Wan, and Min Lu. An accurate and robust range image registration algorithm for 3d object modeling. *IEEE transactions on multimedia*, 16(5):1377–1390, 2014.
- [13] Hao Wu, Li Yan, Hong Xie, Pengcheng Wei, and Jicheng Dai. A hierarchical multiview registration framework of tfs point clouds based on loop constraint. *ISPRS Journal of Photogrammetry and Remote Sensing*, 195:65–76, 2023.
- [14] Shihao Shao, Kaifeng Chen, Arjun Karpur, Qinghua Cui, André Araujo, and Bingyi Cao. Global features are all you need for image retrieval and reranking. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11036–11046, 2023.
- [15] Paul-Edouard Sarlin, Cesar Cadena, Roland Siegwart, and Marcin Dymczyk. From coarse to fine: Robust hierarchical localization at large scale. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12716–12725, 2019.
- [16] Jeffrey S Vitter. Random sampling with a reservoir. *ACM Transactions on Mathematical Software (TOMS)*, 11(1):37–57, 1985.

- [17] Andy Zeng, Shuran Song, Matthias Nießner, Matthew Fisher, Jianxiong Xiao, and Thomas Funkhouser. 3dmatch: Learning local geometric descriptors from rgb-d reconstructions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1802–1811, 2017.
- [18] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5828–5839, 2017.
- [19] Christopher Choy, Jaesik Park, and Vladlen Koltun. Fully convolutional geometric features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8958–8966, 2019.
- [20] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981.
- [21] Haobo Jiang, Zheng Dang, Zhen Wei, Jin Xie, Jian Yang, and Mathieu Salzmann. Robust outlier rejection for 3d registration with variational bayes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1148–1157, 2023.
- [22] Xiyu Zhang, Jiaqi Yang, Shikun Zhang, and Yanning Zhang. 3d registration with maximal cliques. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17745–17754, 2023.
- [23] Stefan Hinterstoisser, Vincent Lepetit, Naresh Rajkumar, and Kurt Konolige. Going further with point pair features. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part III 14*, pages 834–848. Springer, 2016.
- [24] Paul J Besl and Neil D McKay. Method for registration of 3-d shapes. In *Sensor fusion IV: control paradigms and data structures*, volume 1611, pages 586–606. Spie, 1992.
- [25] Radu Bogdan Rusu, Nico Blodow, and Michael Beetz. Fast point feature histograms (fpfh) for 3d registration. In *2009 IEEE international conference on robotics and automation*, pages 3212–3217. IEEE, 2009.
- [26] Haiping Wang, Yuan Liu, Zhen Dong, and Wenping Wang. You only hypothesize once: Point cloud registration with rotation-equivariant descriptors. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 1630–1641, 2022.
- [27] Hao Xu, Nianjin Ye, Guanghui Liu, Bing Zeng, and Shuaicheng Liu. Finet: Dual branches feature interaction for partial-to-partial point cloud registration. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 2848–2856, 2022.
- [28] Haobo Jiang, Mathieu Salzmann, Zheng Dang, Jin Xie, and Jian Yang. Se (3) diffusion model-based point cloud registration for robust 6d object pose estimation. *Advances in Neural Information Processing Systems*, 36, 2024.
- [29] Federica Arrigoni, Beatrice Rossi, and Andrea Fusiello. Global registration of 3d point sets via lrs decomposition. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 489–504. Springer, 2016.
- [30] Federica Arrigoni, Beatrice Rossi, Pasqualina Fragneto, and Andrea Fusiello. Robust synchronization in so (3) and se (3) via low-rank and sparse matrix decomposition. *Computer Vision and Image Understanding*, 174:95–113, 2018.
- [31] Zan Gojcic, Caifa Zhou, Jan D Wegner, Leonidas J Guibas, and Tolga Birdal. Learning multiview 3d point cloud registration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1759–1769, 2020.
- [32] Federica Arrigoni, Beatrice Rossi, and Andrea Fusiello. Spectral synchronization of multiple views in se (3). *SIAM Journal on Imaging Sciences*, 9(4):1963–1990, 2016.

- [33] Avishek Chatterjee and Venu Madhav Govindu. Robust relative rotation averaging. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):958–972, 2017.
- [34] Heng Li, Zhaopeng Cui, Shuaicheng Liu, and Ping Tan. Rago: Recurrent graph optimizer for multiple rotation averaging. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15787–15796, 2022.
- [35] Zhen Dong, Bisheng Yang, Fuxun Liang, Ronggang Huang, and Sebastian Scherer. Hierarchical registration of unordered tls point clouds based on binary shape context descriptor. *ISPRS Journal of Photogrammetry and Remote Sensing*, 144:61–79, 2018.
- [36] Relja Arandjelovic, Petr Gronat, Akihiko Torii, Tomas Pajdla, and Josef Sivic. Netvlad: Cnn architecture for weakly supervised place recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5297–5307, 2016.
- [37] Filip Radenović, Giorgos Tolias, and Ondřej Chum. Fine-tuning cnn image retrieval with no human annotation. *IEEE transactions on pattern analysis and machine intelligence*, 41(7):1655–1668, 2018.
- [38] Seong Hun Lee and Javier Civera. Robust single rotation averaging. *arXiv preprint arXiv:2004.00732*, 2020.
- [39] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015.
- [40] Seong Hun Lee and Javier Civera. Hara: A hierarchical approach for robust rotation averaging. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15777–15786, 2022.

Appendix

A Details about Single Rotation Averaging

Denote the vectorization of an $n \times m$ matrix by $\text{vec}(\cdot)$ and its inverse by $\text{vec}_{n \times m}^{-1}(\cdot)$, the procedure of single rotation averaging is shown in Alg. 1.

Algorithm 1 Single rotation averaging

Input: Rotation \mathbf{R} from robust estimator and list of observed rotation with weight $\{\hat{\mathbf{R}}_i, w_i\}_{i=1}^M$.

Output: Averaged rotation $\bar{\mathbf{R}}$.

```

1:  $\mathbf{s} = \text{vec}(\mathbf{R})$ 
2: for  $it = 1, 2, \dots, 10$  do
3:   for  $i = 1, 2, \dots, M$  do
4:      $\mathbf{v}_i = \text{vec}(\hat{\mathbf{R}}_i) - \mathbf{s}$ 
5:      $d_i = \|\mathbf{v}_i\|$ 
6:    $\mathbf{s}_{pre} = \mathbf{s}$ 
7:    $\mathbf{s} = \frac{\sum_{i=1}^M w_i \mathbf{v}_i / d_i}{\sum_{i=1}^M w_i / d_i}$ 
8:   if  $\|\mathbf{s} - \mathbf{s}_{pre}\| < 0.001$  then
9:     break
10:  $\bar{\mathbf{R}} = \text{proj}_{SO(3)}(\text{vec}_{3 \times 3}^{-1}(\mathbf{s}))$ 
11: return  $\bar{\mathbf{R}}$ 

```

B Probability Mass Function of the Bernoulli Distribution

During the meta-information update, we use a multivariate Bernoulli distribution to address the redundancy issue, with the probability mass function of each variable X_i specified. Here, we provide

a more detailed derivation.

$$\begin{aligned}
P(X_i = k) &= \left(\frac{1}{r_i + 1}\right)^k \left(1 - \frac{1}{r_i + 1}\right)^{1-k} \\
&= \left(\frac{1}{r_i + 1}\right)^k \left(\frac{r_i}{r_i + 1}\right)^{1-k} \\
&= \frac{r_i^{1-k}}{(r_i + 1)^k (r_i + 1)^{1-k}} \\
&= \frac{r_i^{1-k}}{r_i + 1}, k \in \{0, 1\},
\end{aligned} \tag{12}$$

where r_i represents the redundancy count for point in pair i . As a result, for pair i , the possibility of selecting the new keypoint and descriptor is $\frac{1}{r_i + 1}$, while the probability of keeping the existing ones is $\frac{r_i}{r_i + 1}$.

C Implementation Details

We conducted experiments using PyTorch. For the global pooling layer, we utilized a NetVLAD with 64 clusters. Following the practice in [7], we train the model on an augmented 3DMatch training split. For each scene, we sampled alpha frames and for each frame, we randomly selected one keypoint, then picked beta nearest neighbor keypoints to yield varying overlap. We trained the model for 300 epochs using the Adam optimizer with an initial learning rate of 10^{-3} and a weight decay of 10^{-4} , employing a cosine learning rate schedule. For the fusion operation in Eq. 3, the cardinality of all \mathcal{M} was set to 3. The value of k in the first retrieval stage is set to 10, and the distance threshold τ is 0.07.

D Limitation

The primary limitation of our method is its dependence on a target-agnostic descriptor-based pairwise registration method. Some recent pairwise methods use cross-attention to achieve information interaction, causing the extracted source features to vary with different target point clouds. This variability hampers our meta-shape construction and dynamic updates. Therefore, incorporating target-aware pairwise methods into our pipeline could be a promising direction for future research.