

Multiview Point Cloud Registration via Optimization in an Autoencoder Latent Space

Luc Vedrenne, Sylvain Faisan, Denis Fortun

Abstract—Point cloud rigid registration is a fundamental problem in 3D computer vision. In the multiview case, we aim to find a set of 6D poses to align a set of objects. Methods based on pairwise registration rely on a subsequent synchronization algorithm, which makes them poorly scalable with the number of views. Generative approaches overcome this limitation, but are based on Gaussian Mixture Models and use an Expectation-Maximization algorithm. Hence, they are not well suited to handle large transformations. Moreover, most existing methods cannot handle high levels of degradations. In this paper, we introduce POLAR (POint cloud LAtent Registration), a multiview registration method able to efficiently deal with a large number of views, while being robust to a high level of degradations and large initial angles. To achieve this, we transpose the registration problem into the latent space of a pretrained autoencoder, design a loss taking degradations into account, and develop an efficient multistart optimization strategy. Our proposed method significantly outperforms state-of-the-art approaches on synthetic and real data. POLAR is available at github.com/pypolar/polar or as a standalone package which can be installed with `pip install polaregistration`.

Index Terms—Multiview point cloud registration, point cloud reconstruction and restoration, latent space

I. INTRODUCTION

DOWNSTREAM tasks in 3D computer vision often involve a rigid registration step [1]–[6], which consists of determining 6D rigid transformations to align objects. In the multiview context, the objective is to align a set of point clouds, rather than just a pair. This paper focuses on *object-level* registration, where each view represents the same object and the objective is to reconstruct an accurate 3D model of this reference [1], [7]. In particular, our primary goal is to address the problem of data acquired in a microscopy modality called SMLM (single molecule localization microscopy). The challenge of registration with this data is that the views have undergone a very high level of anisotropic noise and outliers, and a moderate level of occlusions. This is in contrast with *scene-level* applications, where the point clouds represent fragments of a large-scale scene obtained with LIDAR or RGB-D camera, with very low noise and outliers, but high

This work of the Interdisciplinary Thematic Institute HealthTech, as part of the ITI 2021-2028 program of the University of Strasbourg, CNRS and Inserm, was partially supported by IdEx Unistra (ANR-10-IDEX-0002) and SFRI (STRAT'US project, ANR-20-SFRI-0012) under the framework of the French Investments for the Future Program. It was also supported by the French National Research Agency (ANR) through the SP-Fluo project (ANR-20-CE45-0007).

Luc Vedrenne (corresponding author: vedrenne@unistra.fr), Sylvain Faisan, and Denis Fortun are with the ICube Laboratory, IMAGeS team, UMR 7357, CNRS, University of Strasbourg, France.

occlusion levels. In what follows, we use the term *degradation* to refer to the alteration of the shape of a point cloud due to noise, occlusions or outliers.

Multiview registration methods can be classified into two main categories. The first one, largely predominant, estimates all the pairwise relative motions, and then runs a subsequent synchronization algorithm to retrieve absolute poses from all the relative ones [8]–[16]. This approach has three main limitations: (i) it poorly scales with the number of views, as it requires $\mathcal{O}(N^2)$ registrations to retrieve N absolute poses, in addition to the cost of the synchronization; (ii) all failed pairwise registrations negatively impact the overall result; (iii) each pairwise registration is performed independently, without leveraging informations from other views. The second family of methods is the generative approach: a template of the reference object from which the views are observed is estimated, and all the views are registered onto this template [7], [17]–[24]. This is typically achieved by modeling the template as a probability distribution using a Gaussian Mixture Model (GMM), jointly estimated with the absolute poses using an Expectation-Maximization (EM) algorithm. These generative approaches have the primary benefit of simultaneously registering all point clouds, which mitigates the limitations of the synchronization approach. However, they are only able to converge to local minima, which limits their applicability to refining poses of already coarsely aligned objects. Moreover, they cannot benefit from the robustness and flexibility of point cloud descriptors learned by neural networks. Hence, almost all state-of-the-art multiview registration methods are currently based on synchronization.

We propose to revitalize the generative approach by formulating the multiview registration problem entirely in the latent space of an autoencoder, which has been pretrained to reconstruct clean views from degraded ones. The interest is fourfold: (i) as a generative approach, it enables the simultaneous registration of numerous views, (ii) the template is modeled by a global latent vector learned by a deep neural network, which arguably makes it more expressive than the usual mixture of Gaussians, (iii) leveraging a global descriptor enables correspondence-free registration, yielding more robustness to noise and occlusion than conventional methods based on feature matching, (iv) the latent space enables faster and more scalable optimization than its ambient counterpart. Moreover, we design a loss that takes degradations into account, such that the estimated template is progressively restored during the optimization. Finally, we propose an optimization scheme

able to retrieve the global optimum from potentially many local ones, which allows our method to handle arbitrarily large transformations between the views to register.

We begin with a brief overview of existing registration methods related to our work and their limitations in Sec. II. In Sec. III, we describe our method, POLAR. In Sec. IV, we provide a theoretical justification of the principle of POLAR. Implementation details are provided in Sec. V. Extensive experiments across many challenging scenarios on both synthetic and real-world data are conducted in Sec. VI, empirically demonstrating the benefits of our approach.

II. RELATED WORKS

We first introduce popular multiview registration methods, either based on synchronization or on a generative framework (Sec. II-A). We then focus on registration approaches that remain pairwise but nonetheless address specific challenges related to our work (Sec. II-B).

A. Multiview registration

Synchronization The main challenge of motion synchronization is to limit the impact of pairwise registration failures. Hence, many synchronization algorithms seek to achieve robust synchronization, by relying on convex relaxation [10], [12]–[15], Iterative Reweighted Least Squares [8], [16], or by operating on the quaternion or $\text{se}(3)$ Lie algebras [9], [11], [25]. Recently, some works have proposed to learn this synchronization [26], [27] + [28]. In [29], the whole process of pairwise registration and synchronization is learned end-to-end.

Generative approach Existing generative methods all rely on an EM algorithm, which is sensitive to initialization. The primary differences among these methods lie in how they model the reconstructed template onto which the objects are registered. Most generative methods rely on GMMs to represent the probability distribution of the reference model onto which the objects are registered [17]–[23]. Variants have been proposed to enhance robustness, by considering normal vector information with Hybrid mixture models [7], by handling outliers with the Student’s t-mixture model [24] or Laplacian model [30], or by performing fuzzy clustering [31], [32].

B. Pairwise registration

Correspondence-based registration The problem of registering a pair of objects is conventionally addressed through a three-step process, possibly iterated until convergence is achieved: (1) feature extraction, (2) feature matching, and (3) transformation estimation. In its simplest form, features correspond to point coordinates, matching is accomplished through a nearest neighbor projection, and the rigid transformation is estimated from correspondences by SVD. This is the ICP (Iterative Closest Point) algorithm [33], which has given rise to numerous variants [34]. Each of these three steps can be enhanced by deep learning algorithms. Numerous methods have sought to learn descriptive features, with an emphasis on rotation-invariant features [35]–[39] or additional

geometric properties [40]–[44]. The most recent ones leverage a geometric Transformer network [37], [38], [44], [45], often in a multi-scale coarse-to-fine approach [39], [44], [46]. Some further learn the matching procedure [39], [44], [45], [47]–[54], typically through a differentiable Sinkhorn algorithm. Estimating the transformation from the established correspondences is typically done through the RANSAC algorithm [55] and its variants [56]–[61].

Correspondence-free registration A distinct category of methods obviates the need for explicit feature matching by characterizing a point cloud with a single global descriptor, typically learned by a neural network [62]–[66]. This global descriptor is anticipated to be more robust to noise, local shape variations or repetitive patterns. However, methods based on global descriptors have been restricted to local convergence and can only be applied after a first step of coarse alignment. **Global convergence** To overcome the sensitivity to local minima, several works have developed global optimization strategies, such as Branch-and-Bound [67] and graduated non-convexity [68], [69], or designed richer descriptors to widen the basins of convergence [70]–[72].

Registration in a latent space Methods that register objects in a dedicated latent space can be traced back to pairwise correlation-based methods that operate in the GMM latent space: the source and target point clouds are both modeled as GMMs and the optimization is performed by minimizing a divergence (typically Kullback-Leibler) [17]–[19]. DeepGMR [73] extends this idea by letting neural networks learn the whole process. Some works proposed to use an autoencoder to learn relevant features in an unsupervised manner [36], [74] or even to register [64], [66], [75], [76], but remain non-generative, and therefore pairwise.

III. METHOD

A. Notations

In what follows, $\rho \mathbf{X} = [\rho \mathbf{x}_1, \dots, \rho \mathbf{x}_k] \in \mathbb{R}^{k \times 3}$ denotes the element-wise application of a rigid motion $\rho \in \text{SE}(3)$ to a point cloud $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_k] \in \mathbb{R}^{k \times 3}$. The use of bold mathematical font indicates a stacking of previously defined variables, whereas calligraphic font is strictly reserved for sets whose elements are of varying sizes. For instance, \mathbf{A} denotes a stacking $[\mathbf{A}_1 \in \mathbb{R}^k, \dots, \mathbf{A}_n \in \mathbb{R}^k] \in \mathbb{R}^{n \times k}$ whereas \mathcal{A} denotes a set $\{\mathbf{A}_1 \in \mathbb{R}^{k_1}, \dots, \mathbf{A}_n \in \mathbb{R}^{k_n}\}$.

B. Problem formulation

Let us consider an unknown reference point cloud \mathbf{X}^* from which N views $\mathcal{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_N\}$ are observed, oriented by rigid motions $\rho^* = [\rho_1^*, \dots, \rho_N^*]$ and corrupted by degradations φ_i :

$$\mathbf{X}_i = \varphi_i(\rho_i^* \mathbf{X}^*), \quad (1)$$

for $i = 1, \dots, N$. Generative multiview registration aims to jointly estimate the N rigid motions ρ^* and the reference point cloud \mathbf{X}^* . To make the problem feasible, the goal is to estimate a parametric model (typically a GMM in previous works) from which the reference \mathbf{X}^* can be generated. Note that there exist several solutions to this problem, corresponding to all possible poses of \mathbf{X}^* .

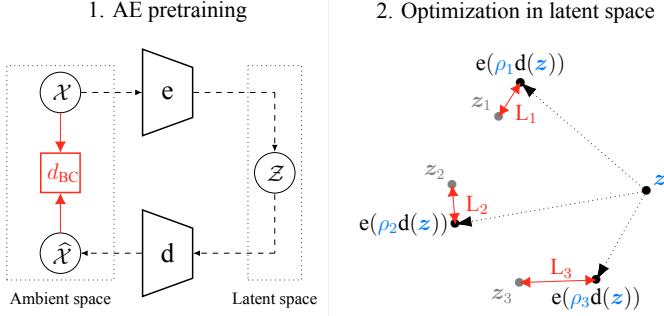


Fig. 1. Overview of the proposed method. 1. Once and for all, before any registration, an autoencoder is trained to reconstruct point clouds (Eq. (2)). 2. To register a set of views, an optimization problem is iteratively solved within the learnt latent space (Eq. (11)).

C. Overview

We propose a two-phase method. First, an autoencoder is trained to reconstruct point clouds, in order to learn a robust global descriptor (Fig. 1.1). This step is performed once and for all before registration, and the autoencoder does not need to be trained again for new point cloud data. The training is described in Sec. III-D. The actual registration is performed in the second phase, by optimizing a cost function defined in the latent space of the frozen autoencoder. We optimize this loss not only with respect to the pose parameters of the views, but also with respect to a latent vector that represents a reconstructed clean point cloud object, on which the views are registered. We describe our latent criterion in Sec. III-E. The optimization scheme to minimize this criterion is described in Sec. III-F.

D. Autoencoder pretraining

The fundamental component of POLAR is a pretrained autoencoder that provides a global descriptor of a point cloud in its latent space. Various architectures have been designed to represent a point cloud through a global descriptor [48], [74], [77]. As our approach is agnostic to the choice of the autoencoder, we consider in this section an abstract definition, and will present our specific implementation in Sec. V-A. The general structure of an autoencoder is the composition of two differentiable functions, an encoder $e : \mathbb{R}^{k \times 3} \mapsto \mathbb{R}^l$ and a decoder $d : \mathbb{R}^l \mapsto \mathbb{R}^{k \times 3}$. We train these functions upstream to reconstruct and restore a point cloud \mathbf{X} degraded by a function φ , *i.e.* to minimize $d_{BC}(d \circ e \circ \varphi(\mathbf{X}), \mathbf{X})$ where d_{BC} denotes the standard bidirectional Chamfer distance. For $\mathbf{P} = [\mathbf{p}_1, \dots, \mathbf{p}_m] \in \mathbb{R}^{m \times 3}$ and $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_n] \in \mathbb{R}^{n \times 3}$, d_{BC} is defined as

$$d_{BC}(\mathbf{P}, \mathbf{Q}) = \frac{1}{m} \sum_{i=1}^m d_{NN}(\mathbf{p}_i, \mathbf{Q}) + \frac{1}{n} \sum_{j=1}^n d_{NN}(\mathbf{q}_j, \mathbf{P}), \quad (2)$$

where

$$d_{NN}(\mathbf{p}, \mathbf{Q}) = \min_{\mathbf{q} \in \mathbf{Q}} \|\mathbf{p} - \mathbf{q}\|_2 \quad (3)$$

is the distance between \mathbf{p} and its nearest neighbor in \mathbf{Q} . Once trained, this autoencoder is frozen, and the registration will be performed in its latent space.

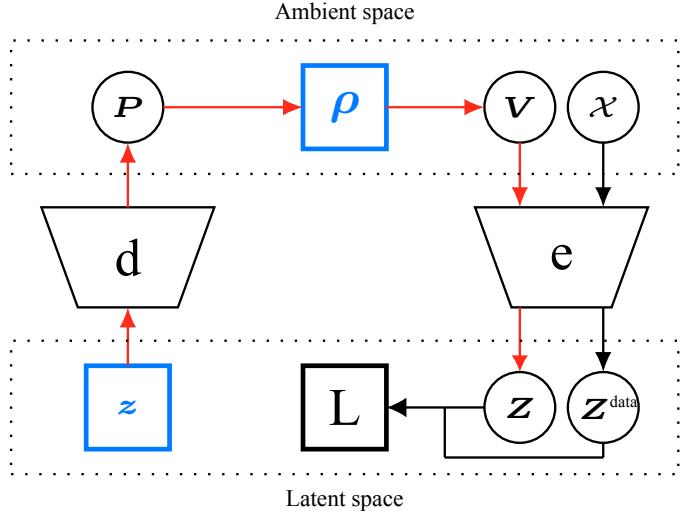


Fig. 2. Graphical illustration of the loss computation in POLAR. $\mathcal{X} = \{\mathbf{X}_1 \in \mathbb{R}^{k_1}, \dots, \mathbf{X}_N \in \mathbb{R}^{k_N}\}$ denotes the views to register. $\mathbf{Z}_{\text{data}} = [e(\mathbf{X}_1), \dots, e(\mathbf{X}_N)] \in \mathbb{R}^{N \times l}$ is the matrix of their encodings. $\mathbf{P} = d(\mathbf{z}) \in \mathbb{R}^{k \times 3}$ is the estimated template. $\mathbf{V} = [\rho_1 \mathbf{P}, \dots, \rho_N \mathbf{P}] \in \mathbb{R}^{N \times k \times 3}$ denotes the views obtained by applying the estimated motions ρ to the estimated template \mathbf{P} . Finally, $\mathbf{Z} = [e(\rho_1 \mathbf{P}), \dots, e(\rho_N \mathbf{P})] \in \mathbb{R}^{N \times l}$ is the matrix of the latent vectors obtained by encoding these estimated views.

E. Registration loss in the latent space

a) *Clean data:* We first consider the case of data without degradation. It follows from Eq. (1) that the registration task reduces to finding \mathbf{X}^* and ρ^* such that $\mathbf{X}_i = \rho_i^* \mathbf{X}^*$. Instead of directly comparing the point clouds \mathbf{X}_i and $\rho_i^* \mathbf{X}^*$, we compare their latent representation. We seek \mathbf{X}^* and ρ^* such that $e(\mathbf{X}_i) = e(\rho_i^* \mathbf{X}^*)$. Furthermore, the template itself is expressed through its latent representation: we estimate a latent vector \mathbf{z} such that after decoding, $d(\mathbf{z})$ represents \mathbf{X}^* . Thus, the encoding of a view to register $e(\mathbf{X}_i)$ is compared with the encoding of the reconstructed template $d(\mathbf{z})$, rigidly transformed with the estimated pose ρ_i , which writes $e(\rho_i d(\mathbf{z}))$. Hence, the optimization problem to solve is

$$\hat{\mathbf{z}}, \hat{\rho} = \underset{\mathbf{z}, \rho}{\operatorname{argmin}} L(\mathbf{z}, \rho) \quad (4)$$

with

$$L(\mathbf{z}, \rho) = \sum_i^N \|e(\rho_i d(\mathbf{z})) - e(\mathbf{X}_i)\|_2^2. \quad (5)$$

A graphical model of the computation of this loss is presented in Fig. 2.

POLAR can be viewed as the combination of generative and deep correspondence-free approaches. It is generative because it estimates a parameterized representation of the object X^* , which brings the benefits of simultaneous registration of all the point clouds. Since this parameterization is a global descriptor, we also avoid the need for local correspondences. Furthermore, this global descriptor is obtained from an autoencoder and thus exhibits an increased robustness and modeling expressiveness, in particular compared to the standard GMM representation.

The loss of Eq. (5) is designed for registration of non-degraded point clouds. In the following paragraphs, we extend this loss to handle anisotropic noise, partial visibility, and outliers.

b) Anisotropic Noise: We first consider the case of noisy data. We denote ν_i the function that randomly noises each point of a cloud, under the i.i.d. assumption, such that our degradation model is $\varphi_i(\mathbf{X}) = \nu_i(\mathbf{X})$. We make no hypothesis regarding the nature of this noise. In particular, it can be anisotropic and different for each view. Using the loss (5) in this case would lead to a noisy reconstruction $d(\mathbf{z})$ that would fit the noisy data \mathbf{X}_i . Therefore, in order to enforce the denoising of $d(\mathbf{z})$, we follow the generative model in Eq. (1) and account for this noise by applying ν_i to the reconstructed point cloud to create a noisy reconstruction. Thus, the loss (5) becomes

$$L(\mathbf{z}, \rho) = \sum_i^N \|e(\nu_i(\rho_i d(\mathbf{z}))) - e(\mathbf{X}_i)\|_2^2. \quad (6)$$

Note that this loss is not deterministic. However, since one realisation of noise is added to each point, when a cloud is composed of a fairly high amount of points, the variability at the shape scale remains low enough for the optimization to converge.

c) Partial visibility: We now consider the case where the views are partially occluded: some points of each view are masked out. An illustration of this degradation can be seen in Fig. 4a, where red points are occluded. Since the template is estimated from many views, each occluded in a unique way, it should be possible to reconstruct a complete object. Therefore, some regions of the complete reconstructed object $\rho_i d(\mathbf{z})$ will not have correspondences in the view \mathbf{X}_i , which introduces undesirable discrepancies in their encoding $e(\rho_i d(\mathbf{z}))$ and $e(\mathbf{X}_i)$ and affects the computation of the loss (5). To prevent this, we estimate the parts of $\rho_i d(\mathbf{z})$ missing in \mathbf{X}_i , in order to mask them in $\rho_i d(\mathbf{z})$. We denote this masked template as $\mathbf{M}_{\mathbf{X}_i}^v(\rho_i d(\mathbf{z}))$, where $v \in [0, 1]$ is a parameter representing the proportion of occluded points in the views. The loss (5) thus becomes

$$L(\mathbf{z}, \rho) = \sum_i^N \|e(\mathbf{M}_{\mathbf{X}_i}^v(\rho_i d(\mathbf{z}))) - e(\mathbf{X}_i)\|_2^2. \quad (7)$$

To compute the mask $\mathbf{M}_{\mathbf{X}_i}^v(\rho_i d(\mathbf{z}))$, we leverage the vector of nearest neighbor distances

$$\mathbf{D}_{P \rightarrow Q} = [d_{NN}(\mathbf{p}_1, \mathbf{Q}), \dots, d_{NN}(\mathbf{p}_m, \mathbf{Q})] \quad (8)$$

where $d_{NN}(\mathbf{p}_k, \mathbf{Q})$ is defined in Eq. (3). The elements of the nearest neighbor distances $\mathbf{D}_{\rho_i d(\mathbf{z}) \rightarrow \mathbf{X}_i}$ from the complete reconstructed template $\rho_i d(\mathbf{z})$ to the observed occluded point cloud \mathbf{X}_i are illustrated in Fig. 3.b. In the absence of noise, and assuming a perfect estimation of the pose, $\mathbf{D}_{\rho_i d(\mathbf{z}) \rightarrow \mathbf{X}_i}$ is zero for visible points of \mathbf{X}_i and non zero for occluded ones. Hence, $\mathbf{D}_{\rho_i d(\mathbf{z}) \rightarrow \mathbf{X}_i}$ can be used to select the points to mask: for a view \mathbf{X}_i and assuming that the proportion of occluded points v is known, the proportion v of elements of $\mathbf{D}_{\rho_i d(\mathbf{z}) \rightarrow \mathbf{X}_i}$ with the highest values are considered as occluded points and are masked.

d) Outliers: Finally, we now describe how we handle the presence of outliers in the views. In Fig. 4b we show an example of outliers on a view. In that case, we follow the same reasoning as for occluded data. As we seek a template

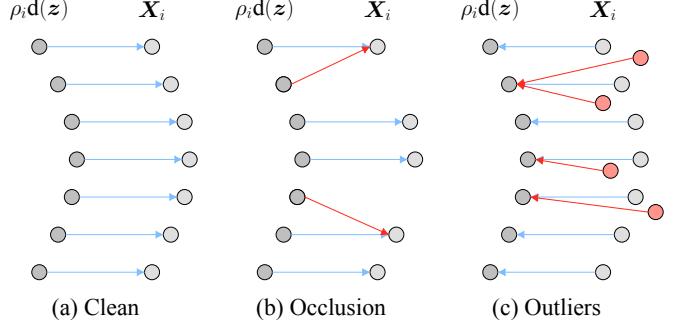


Fig. 3. Schematic representation of nearest neighbor distances in three scenarios. The estimated template in the i -th pose $\rho_i d(\mathbf{z})$ is translated away from the corresponding view \mathbf{X}_i for visualization purpose but should be seen as superimposed, such that blue arrows denote null distances. (a) When the two objects are identical, all distances are null. (b) In case of occlusion, the distance from a point in the template to its nearest neighbor in the view is non-zero if and only if this point is occluded in the view (red arrows). (c) Similarly, the distance from a point in the view to its nearest neighbor in the template is non-zero if and only if this point is not part of the template, i.e. if and only if it is an outlier (red arrows).

$\rho_i d(\mathbf{z})$ without outliers, the outliers of the views \mathbf{X}_i have no correspondence in the template, and this will be reflected in their encodings. Thus, we estimate parts of \mathbf{X}_i missing in $\rho_i d(\mathbf{z})$ and mask them before computing the loss. Let $\mathbf{M}_{\rho_i d(\mathbf{z})}^o(\mathbf{X}_i)$ denotes the masked view, where $o \in [0, 1]$ is a parameter representing the proportion of outliers in each view. The loss with outliers handling writes

$$L(\mathbf{z}, \rho) = \sum_i^N \|e(\rho_i d(\mathbf{z})) - e(\mathbf{M}_{\rho_i d(\mathbf{z})}^o(\mathbf{X}_i))\|_2^2. \quad (9)$$

As for occluded data, we compute the mask $\mathbf{M}_{\rho_i d(\mathbf{z})}^o(\mathbf{X}_i)$ using the vector of nearest neighbor distances. But here we compute the distance from the view to the template $\mathbf{D}_{\mathbf{X}_i \rightarrow \rho_i d(\mathbf{z})}$. Fig. 3.c illustrates the values of this vector without noise and assuming perfect pose estimation: $\mathbf{D}_{\mathbf{X}_i \rightarrow \rho_i d(\mathbf{z})}$ is zero for inliers and non zero for outliers. Hence, assuming that the outliers ratio o is known, a point \mathbf{x} in \mathbf{X}_i is masked if its distance $d_{NN}(\mathbf{x}, \rho_i d(\mathbf{z}))$ is above the $1 - o$ percentile of $\mathbf{D}_{\mathbf{X}_i \rightarrow \rho_i d(\mathbf{z})}$.

e) Combining degradations: In the previous sections, we have presented how the loss (5) can be extended to take into account three types of degradation separately. However, in real scenarios, data to register can present any combination of these degradations. If the views \mathbf{X}_i are noised with a function ν_i , a ratio $1 - v$ of each view is occluded, and there is a ratio o of outliers points in each view, we use the loss

$$L(\mathbf{z}, \rho) = \sum_i^N L_i(\mathbf{z}, \rho) \quad (10)$$

with

$$L_i(\mathbf{z}, \rho) = \|e(\mathbf{M}_{\mathbf{X}_i}^v(\rho_i d(\mathbf{z}))) - e(\mathbf{M}_{\nu_i(\rho_i d(\mathbf{z}))}^o(\mathbf{X}_i))\|_2^2. \quad (11)$$

This approach allows us to incorporate prior knowledge of the degradation model into the criterion.

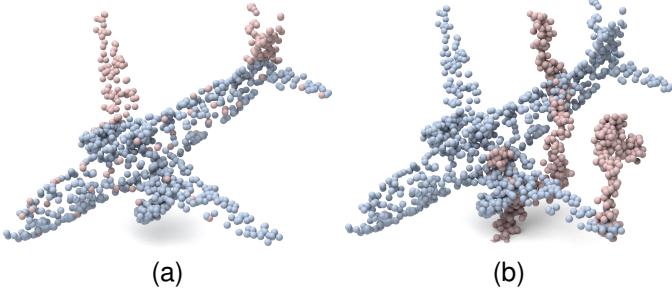


Fig. 4. Visual results of the masking operations (Sec. III-E.c, Sec. III-E.d) in case of occlusion and outliers. Red points in the template (a) and in the view (b) will be discarded for the loss computation.

f) Regularization: We empirically found that when dealing with highly degraded data, POLAR converges towards a template in which the distribution of points across the surface is not uniform: there are regions of high point density and regions of low point density. To penalize this behaviour, we add a regularization term that corresponds to the standard deviation of local point density. Precisely, given a point cloud $\mathbf{X} \in \mathbb{R}^{k \times 3}$, let $b_r(\mathbf{x}, \mathbf{X})$ be the number of points in \mathbf{X} lying inside the ball of center \mathbf{x} and radius r (*i.e.* the number of neighbors of \mathbf{x}). The mean point density of the point cloud \mathbf{X} is

$$\bar{b}_r(\mathbf{X}) = \frac{1}{k} \sum_{\mathbf{x} \in \mathbf{X}} b_r(\mathbf{x}, \mathbf{X}) \quad (12)$$

and its variance is

$$\sigma_{b_r}^2(\mathbf{X}) = \frac{1}{k-1} \sum_{\mathbf{x} \in \mathbf{X}} (b_r(\mathbf{x}, \mathbf{X}) - \bar{b}_r(\mathbf{X}))^2. \quad (13)$$

We define the regularization term $R(\mathbf{z})$ as

$$R(\mathbf{z}) = \sigma_{b_r}(\mathbf{d}(\mathbf{z})) \quad (14)$$

and the final loss optimized in POLAR is

$$\mathbf{L}^f(\mathbf{z}, \boldsymbol{\rho}) = \mathbf{L}(\mathbf{z}, \boldsymbol{\rho}) + \lambda R(\mathbf{z}) \quad (15)$$

where λ is a weighting hyperparameter.

F. Optimization procedure

We now describe how we minimize the criterion \mathbf{L}^f in Eq. (15). This optimization procedure is a hard task due to the highly non-convex nature of the problem. This is a fundamental issue in registration algorithms: as real-world objects often have symmetries or near-symmetries, discrepancy measures exhibit several local minima. This is illustrated in Fig. 5 for an airplane, where 2 minima exist, corresponding to a 180° rotation along its fuselage axis (Fig. 5). As gradient-based optimization may be trapped in local minima, we design an optimization procedure able to retrieve the global minimum from potentially many local ones. Our approach employs a multistart strategy, detailed in the following subsections. In summary, we perform multiple gradient descents in parallel, starting from various plausible initializations, aiming to span all basins of attraction. These plausible starts are determined by a coarse exhaustive search described in Sec. III-F.a. The whole procedure is repeated until convergence.

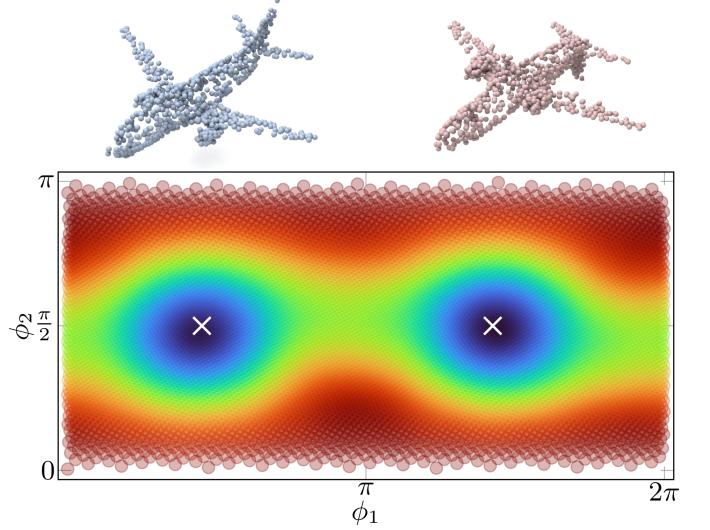


Fig. 5. Landscape of the loss in Eq. (5) with respect to the two first Euler angles (the loss is summed over the third Euler angle). The white crosses are the results of the SO(3) local minima search (Sec. III-F). In the case of an airplane, two minima coexist, corresponding to a 180° rotation along the fuselage axis.

a) Finding Local Minima over SO(3) (FLAMES):

Before describing our full optimization pipeline, we present a novel necessary tool for finding local minima of the loss with respect to the rotations only, that we name FLAMES. These local minima are found by an exhaustive search over a uniform sampling of L rotations in SO(3) denoted $\text{SO}(3)_L$. Note that obtaining such a sampling is a hard problem, which we tackle using the Super Fibonacci spirals [78]. A k -neighborhood graph is then defined in $\text{SO}(3)_L$, using the relative angle, a geodesic distance on the manifold of rotations:

$$\angle : \text{SO}(3) \times \text{SO}(3) \mapsto [0, \pi]$$

$$\mathbf{R}_1, \mathbf{R}_2 \mapsto \arccos \left(\frac{\text{tr}(\mathbf{R}_1 \mathbf{R}_2^\top) - 1}{2} \right). \quad (16)$$

This k -nn graph only depends on L and k and as such, can be pre-computed once and for all. We fix the current estimation of the template \mathbf{z} . Thus, the terms $\mathbf{L}_i(\mathbf{z}, \boldsymbol{\rho})$ Eq. (11) of the criterion (10) become independent. For each view, we fix the translation t_i as well, and compute the criterion $\mathbf{L}_i(\mathbf{z}, \boldsymbol{\rho})$ (11) for each sampled rotation. If a rotation has the lowest cost in its k -neighborhood, it is a local minimum for this view. An example of loss landscape and local minima obtained with this procedure in the case of an airplane is illustrated in Fig. 5.

b) Initialization: Given a matrix $\mathbf{Z}^{\text{data}} = [\mathbf{z}_1, \dots, \mathbf{z}_N]$ of N latent vectors of views to register, the parameters to estimate \mathbf{z} and $\boldsymbol{\rho}$ are initialized as follows. The template \mathbf{z}_{init} is set to the medoid of the latent vectors \mathbf{Z}^{data} , *i.e.* the latent vector whose sum of dissimilarities to all the others is minimal:

$$\mathbf{z}_{\text{init}} = \underset{\mathbf{z} \in \mathbf{Z}^{\text{data}}}{\operatorname{argmin}} \sum_{i=1}^N \|\mathbf{z} - \mathbf{z}_i\|_2 \quad (17)$$

The translations \mathbf{t} are set to zeros since the views are centered beforehand (*cf.* Sec. V-B). The rotations are initialized using

the FLAMES procedure described above, taking the local minimum of minimal cost for each view.

c) Joint gradient descent: The parameters z and ρ are jointly updated to minimize the criterion (15) by gradient descent, using the Adam algorithm [79] with decoupled weight decay (AdamW) [80]. The implementation is based on PyTorch’s automatic differentiation engine [81].

d) Parallel multistart: Using the current estimation of z and t , the FLAMES procedure is used to find the M best local minima over the rotation sampling $\text{SO}(3)_L$ for each of the N views. This results in $N \cdot M$ initializations, from which as many gradient descents are executed in parallel. In this step, we only optimize the rigid motions ρ and keep the current estimated template z fixed. Note that this amounts to no more than $N \cdot M \cdot 6$ parameters. Even when registering a large number of views, this remains a relatively low-dimensional problem, which allows us to run all the gradient descents in parallel. After convergence, M new poses and losses are obtained for each of the N views. First, we select the new pose of minimal loss for each view. We note $\rho^F = [\rho_1^F, \dots, \rho_N^F]$ the obtained new poses and $L^F = [L_1^F, \dots, L_N^F]$ their corresponding losses. Given the current poses $\rho = [\rho_1, \dots, \rho_N]$ and their corresponding losses $L = [L_1, \dots, L_N]$, a new pose replaces the current one if it has a lower loss. Formally, the update function u writes

$$u(\rho_i^F, \rho_i) = \begin{cases} \rho_i^F & \text{if } L_i^F \leq L_i \\ \rho_i & \text{if } L_i^F > L_i \end{cases} \quad (18)$$

e) Full procedure: After the first FLAMES and parameters initialization, the sequence (i) joint gradient descent, (ii) FLAMES, (iii) parallel multistart is repeated until a stopping criterion is met. To design this criterion, we detect if we have escaped from a local minimum after the multistart step, right before the update (18). This is achieved by verifying if a new rotation is more than T_R away from the current one and if its loss is better than the current one. Let $R^F = [R_1^F, \dots, R_N^F]$ and $R = [R_1, \dots, R_N]$ denote the rotation part of the new poses ρ^F and of the current ones ρ respectively. The function s detecting an escape from a local minimum is defined as

$$s(\rho_i^F, \rho_i) = \begin{cases} 1 & \text{if } L_i^F \leq L_i \text{ and } \angle(R_i^F, R_i) \geq T_R \\ 0 & \text{otherwise.} \end{cases} \quad (19)$$

The full optimization procedure is executed until no escape is detected, *i.e.* until

$$\sum_{i=1}^N s(\rho_i^F, \rho_i) = 0. \quad (20)$$

The POLAR optimization scheme is summarized in Alg. 1.

IV. MOTIVATION

In this section, we interpret and motivate our model in (5) from a differential geometry perspective.

The set of rigid motion $\text{SE}(3)$ is endowed with a structure of Lie group. This property involves that the orbit of a point cloud X , denoted as $\mathcal{O}_X = \{\rho X, \rho \in \text{SE}(3)\}$, is a compact and smooth manifold of dimension at most 6 [82]. In the absence

```

Data:  $Z^{\text{data}} = e(\mathcal{X})$ ,  $\text{SO}(3)_L$ 
Result:  $\hat{z}, \hat{R}, \hat{t}$ 
 $z \leftarrow \text{medoid}(Z^{\text{data}})$ 
 $t \leftarrow 0$ 
 $R \leftarrow \text{FLAMES}_{\text{top } 1}(z, t, Z^{\text{data}}, \text{SO}(3)_L)$ 
while not converged do
     $z, R, t, L \leftarrow \text{joint gradient descent}(z, R, t, Z^{\text{data}})$ 
     $R^F \leftarrow \text{FLAMES}_{\text{top } M}(z, t, Z^{\text{data}}, \text{SO}(3)_L)$ 
     $\rho^F, L^F \leftarrow \text{multistart}(z, R^F, t, Z^{\text{data}}, L)$ 
     $n \leftarrow \sum_{i=1}^N s(\rho_i^F, \rho_i)$  (19)
     $\rho \leftarrow u(\rho^F, \rho)$  (18)
    if  $n = 0$  then
        | converged
    end
end
 $\hat{z}, \hat{R}, \hat{t} \leftarrow \text{joint gradient descent}(z, R, t, Z^{\text{data}})$ 

```

Algorithm 1: POLAR optimization algorithm. The notation $\text{FLAMES}_{\text{top } k}$ means that the k best local minima for each view are selected following the FLAMES procedure of Sec. III-F.

of degradation, the views to register would be a sampling of the unknown template’s orbit \mathcal{O}_{X^*} . In practice, data corruption moves away the observed views X_i from the orbit \mathcal{O}_{X^*} . Our loss (5) is defined in the autoencoder’s latent space and it aims at estimating $e(\mathcal{O}_{X^*})$, where $e(\mathcal{O}_X) = \{e(\rho X), \rho \in \text{SE}(3)\}$ is the encoding of the orbit in the latent space. One interest of performing the registration in the latent space is that the encoder can be guided to produce a latent representation that is robust to data corruption. Of course, the encoder is not strictly invariant to degradation. Nevertheless, it is still robust to a certain extent, such that $e(\varphi(X))$ should be closer to $e(X)$ than $\varphi(X)$ is to X .

In order to enable efficient gradient-based optimization in the latent space, it is crucial that $e(\mathcal{O}_X)$ preserves the manifold structure of \mathcal{O}_X . Since the image of a manifold under an embedding remains a manifold, it is sufficient for e to be an embedding (an injective function whose differential is injective). However, its smoothness is enough to ensure that $e(\mathcal{O}_X)$ is “almost” a smooth manifold thanks to the following result [82], [83]:

Corollary of the Whitney embedding theorem Suppose Γ is a compact smooth n -manifold with or without boundary. If $l \geq 2n + 1$, then every smooth map from Γ to \mathbb{R}^l can be uniformly approximated by embeddings.

With $\Gamma = \mathcal{O}_X$, this corollary asserts that the encoder e can be uniformly approximated by an embedding (for any $\varepsilon > 0$, there exists an embedding g such that $\|g - e\|_\infty < \varepsilon$), provided that e is a smooth function and the latent space has dimension $l > 13$. Thus, the set $e(\mathcal{O}_X)$ is arbitrarily close to a manifold.

V. IMPLEMENTATION DETAILS

A. Autoencoder architecture and training

a) Architecture: We define our encoder as a simple PointNet architecture [77] where we removed the T-net module. Our decoder is an MLP that takes the global feature of

PointNet as input and outputs a point cloud. Batch norm is used for all layers with ReLU. DropOut layers are used for the decoder. Note that while PointNet is theoretically able to process point clouds of varying sizes, it is not trivial to come up with an implementation that actually allows it. A naive way would be to duplicate some points to pad each point cloud so that the autoencoder receives inputs of a fixed size. Unfortunately, the invariance to point duplication only holds true in the absence of Batch Normalization layers. To fully allow varying sizes without computational overhead while maintaining the use of Batch Normalization layers, we employ the message passing scheme from the `torch geometric` library [84] for our implementation. Note that any other global autoencoder could be integrated to POLAR.

b) Training: We train a single autoencoder once and for all, and use it in all subsequent experiments, whether they involve simulated or real data, and regardless of the degradations. Our autoencoder is trained on the full ModelNet40 [85] training set. The training is performed for 200 epochs with the AdamW [80] optimizer. The initial learning rate is $1e^{-3}$. The learning rate is divided by a factor of 2 when the loss does not improve for 10 consecutive epochs. We then process the data as follows. First, basic pre-processing steps are applied: 1024 points are randomly sub-sampled (ModelNet is originally composed of dense point clouds of 5000 points), centered, and normalized to lie exactly in the unit sphere. Then, a random pose is applied. To obtain an autoencoder able to reconstruct objects in arbitrary poses, it is crucial to uniformly sample $\text{SO}(3)$, without relying on a discretization. To achieve this, we do not sample uniformly each Euler's angle, as it does not result in a uniform sampling over $\text{SO}(3)$. Instead, we leverage the exponential map from the Lie algebra to the underlying group $\text{SO}(3)$. We sample an axis $\mathbf{k} = [k_x, k_y, k_z]$ on the unit sphere and an angle $\theta \sim \mathcal{U}(0, \pi)$. The corresponding element \mathbf{K} of the Lie algebra is

$$\mathbf{K} = \begin{bmatrix} 0 & -k_z & k_y \\ k_z & 0 & -k_x \\ -k_y & k_x & 0 \end{bmatrix} \quad (21)$$

and its associated rotation \mathbf{R} in $\text{SO}(3)$ is

$$\mathbf{R} = \mathbf{I}_3 + \sin \theta \mathbf{K} + (1 - \cos \theta) \mathbf{K}^2. \quad (22)$$

Finally, degradations are applied with the following policy:

- 1) **Jit:** Add a centered multivariate Gaussian.
- 2) **Plane Cut:** Sample a plane normal, and retain points close enough to this plane, so that a visibility ratio v of the object is retained, with $v \sim \mathcal{U}(0.7, 1)$. [39], [86], [87]
- 3) **Center & Normalize:** Center and scale to lie exactly within the unit sphere.

On the standard ModelNet40 training set, the training takes about 20mn on a single NVIDIA GeForce RTX 3090.

B. Data normalization for registration

While it is common to normalize point clouds by centering and rescaling them to fit within a unit sphere, it is crucial for rigid registration to apply the same scaling factor for each view

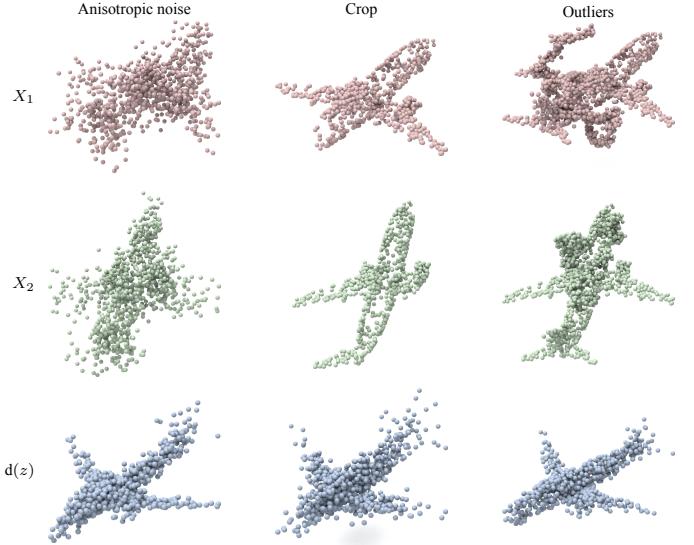


Fig. 6. **Visual results.** First two rows: Degraded views to register. Last row: Template obtained after the optimization of the criterion L^f (Eq. (15)) with known visibility ratio $r = 0.2$, outliers ratio $o = 0.3$, and covariance matrix $\Sigma = \text{diag}(0.03, 0.03, 0.15)$. Each optimization is done for $N = 100$ views (2 of which are shown). For each degradation, the estimated template $d(z)$ has compensated for the degradation: it is complete, without outliers, and denoised.

to be registered in order to ensure that point clouds maintain their relative sizes. This is achieved by independently calculating the scaling factor for each view and then normalizing them using the smallest previously calculated scaling factor.

C. Hyperparameter setting

For the $k\text{-nn}$ graph over $\text{SO}(3)$, we sample $L = 5e^4$ rotations and use $k = 256$ neighbors. The multistart step considers the top $M = 8$ minima per view, and T_R is set to 15° to detect escapes from local minima. The decoder outputs clouds of 1024 points, and the latent space dimension is $l = 1024$. The regularization is done with balls of radius 0.1 and the loss weighting λ is set to $1e^{-2}$. Each gradient descent is performed until the loss does not improve for 100 steps. The learning rate is $1e^{-2}$ at start, and is divided by 10 every time the loss doesn't improve for 10 consecutive steps.

VI. EXPERIMENTS

A. Experimental protocol

a) Baselines: We compare POLAR with different classes of methods: (i) non-learning-based methods, which can be either pairwise, like FGR [69] and more recently MAC [88], or based on a generative multiview paradigm, like JRMPG [20] and EMPMR [23]; (ii) deep learning-based methods, which are pairwise, namely PointNetLK [62], DCP [87], RPM-Net [86], DeepGMR [73], or more recently GeoT [39] and RoITr [38], with the exception of SGHR [16], a recent deep learning based multiview method. For pairwise methods, we use a synchronization algorithm to retrieve multiview absolute poses. We selected the method described in [89], which we empirically observed to provide the best results. It is based on

an iterative reweighted least squares method coupled with a message passing algorithm that estimates corruption levels.

b) Data: We use ModelNet40 [85], a comprehensive CAD model repository of 4602 objects spanning 40 categories as the synthetic dataset of reference. We follow the official training and validation split from [85]. Learning-based methods, including POLAR’s autoencoder are trained on this training set. Following the usual processing [39], [86], symmetric classes are removed. All subsequent experiments are conducted on the obtained validation set, considering 100 views for each registration task. Each view is obtained following the preprocessing of Sec. V-A followed by the application of a degradation. In all experiments involving degradations, we suppose that their parameters are known (namely the noise covariance matrix Σ , the visibility ratio v and the outliers ratio σ). The impact of fixing wrong values to these parameters is moderate and will be studied in Sec. VI-C.c. Since POLAR is tailored for object-level registration, motivated by the SMLM application (registration of many severely noised versions of an object), we do not consider scene-level datasets. POLAR is not optimally suited for scene registration for two reasons. First, in POLAR, the views are registered on an estimated template, hence requiring a global latent vector to encode the entire scene. Our autoencoder is able to accurately represent single objects, but its representational capacity is not sufficient to capture the details of large-scale scenes. Second, in POLAR, the template is initialized with one of the views. This can be problematic in large-scale scenes, where a single fragment represents only a small portion of the complete scene. As a result, the optimization is more likely to fall into a local minimum. In contrast, at the object level, each view is an occluded observation of the same complete object, leading to a higher overlap with the template. In Sec. VI-D, we provide results on challenging object-level real data, namely the Faust partial dataset [90] and SMLM data [91].

c) Metrics: Regardless of the approach employed, a rigid transformation is obtained for each point cloud, which enables their alignment in a shared reference frame. Since the pose in this frame can be arbitrary, we independently examine the $N(N - 1)/2$ relative poses between the pairs of views. The one related to X_i and X_j is computed from the estimated and ground truth rotations, respectively $\hat{\mathbf{R}}$ and \mathbf{R}^* as $\theta_{ij} = \angle(\hat{\mathbf{R}}_i \mathbf{R}_i^*, \hat{\mathbf{R}}_j \mathbf{R}_j^*)$. This is usually called Relative Rotation Error (RRE). In our experiments, we report the registration success rate, often called Registration Recall (RR), defined as the proportion of angles θ_{ij} that are below a threshold. We also show the cumulative distribution function (CDF) of θ_{ij} , which provides a complete view of the registration performance, independent from the choice of a threshold (the registration recall for one threshold t is a point of the CDF). In particular, the CDF allows us to distinguish the notions of *accuracy*, which corresponds to low values of θ_{ij} and reflects the quality of the alignment when the registration is successful, and *robustness*, which corresponds to higher values of θ_{ij} and denotes the ability to obtain coarse successful registration in challenging scenarios. In the context of global registration, we are more interested in *robustness*, since *accuracy* can be

TABLE I
REGISTRATION RECALL (RR, $t = 10^\circ$) EVOLUTION WITH INITIAL ANGLE.
GLOBAL METHODS ARE SHOWN IN BOLD.

METHOD	$\theta \leq \frac{\pi}{4}$	$\theta \in [\frac{\pi}{4}, \frac{\pi}{2}]$	$\theta \in [\frac{\pi}{2}, \frac{3\pi}{4}]$	$\theta \in [\frac{3\pi}{4}, \pi]$
FGR	99.46	97.9	96.67	95.9
MAC	99.46	97.29	93.19	92
FMR	96.05	80.18	51.88	32
JRMPG	96.05	80.18	51.88	32
EMPMR	69.21	10.14	0.28	0
PNLK	76.57	31.56	3.68	0
DCP	71.53	24.73	8.16	3
RPM-NET	99.73	99.81	99.37	99
DEEPGMR	99.59	98.18	97.29	98
GEOTR	99.80	99.42	99.00	98.66
RoITr	100	100	100	100
SGHR	100	100	100	100
POLAR	100	100	100	100

achieved by subsequent refinement. Translation parameters are not considered for evaluation because we observed that a proper estimation of the rotation consistently leads to a proper estimation of the translation.

B. Comparison with other methods

a) Ability to handle large transformations: We first study the ability of selected methods to handle arbitrarily large transformations. We use the ModelNet40 validation set, pre-processed following the policy of Sec. V-A (which include a uniformly random pose over the whole SO(3) group), without degradation. Once N absolute poses have been estimated, the RRE and initial angle are computed for each pair. Finally, the RR is computed for four ranges of initial angles $\angle(\mathbf{R}_i^*, \mathbf{R}_j^*)$: $\leq 45^\circ, \in [45^\circ, 90^\circ], \in [90^\circ, 135^\circ], \in [135^\circ, 180^\circ]$. The results are presented in Tab. I. Global methods are shown in bold. As advertised, methods based on the EM algorithm (JMPRC and EMPMR) can only converge locally. As for deep learning based methods, PointNetLK is also local due to the use of the Lucas Kanade algorithm. DCP relies on an SVD algorithm to retrieve a transformation from correspondences estimated through an attention map between pose variant features, hence it remains local. In subsequent experiments, these local methods are not considered. On the other hand, FGR and MAC are almost global, even if for large angles they may be fooled by near-symmetries, while RPM-Net, DeepGMR, GeoT, RoITr, SGHR, and POLAR maintain near perfect performances regardless of the initial angle.

b) Isotropic Noise: We then study the robustness of these global methods to isotropic additive Gaussian noise. We consider increasing levels of noise following $\mathcal{N}(0, \sigma)$ with σ ranging from 0 to 0.15. For one given noise level, each view is uniquely degraded by adding noise. The Registration recall is then computed for each level of noise (Fig. 7). Methods purely based on local correspondences (FGR, MAC) struggle to handle large noise levels and as such, we do not consider them in further experiments. Even with their more robust correspondences estimation mechanisms, RoITr and GeoT fail due to the large local dissimilarity induced by strong noise.

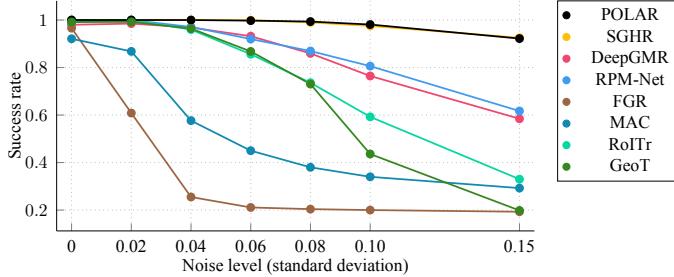


Fig. 7. Registration Recall evolution with noise standard deviation.

It should be noted that POLAR and SGHR exhibit a strong robustness to isotropic noise.

c) *Anisotropic Noise*: A far harder degradation is the case of anisotropic noise, as it deforms each view in a unique way depending on its orientation. Two examples of such views are shown in Fig. 6. This often occurs in microscopic acquisition, in which the resolution in the microscope’s axis is far lower than in its orthogonal plane. The greater the anisotropy of resolution, the greater the deformation undergone by objects. We study the impact of this anisotropy factor on registration performances, measured by the full Registration Recall CDF, reported in Fig. 8. Under such a degradation, views are no longer locally similar. Therefore, methods purely based on local correspondences such as DeepGMR and RPM-Net cannot handle such cases. When the anisotropy factor and the noise level are relatively low, RoITr, GeoT, SGHR, and POLAR obtain good performances (Fig. 8.a). With the same noise level but a stronger anisotropy factor, performances decrease and only POLAR maintains a high success rate (Fig. 8.b). Similarly, with a relatively low anisotropy factor but stronger base noise level, POLAR maintains a high success rate when other methods performances drop (Fig. 8.c). An example of template obtained by POLAR is shown in Fig. 6.b. Compared to the input views, the reconstructed template is denoised and the anisotropy is corrected. To conclude, SGHR, RoITr and GeoT keep relatively good performances, as they partially leverage global descriptors: in SGHR, YOHO [92] features are processed by a NetVLAD [93] module, while RoITr and GeoT gradually transforms local correspondences in more global ones using an attention mechanism. However, the POLAR approach obtains the best performances.

d) *Varying noise level*: To further study the ability of POLAR to be robust to local shape variations, we first consider the case of noisy data, with a different noise level for each view. Specifically, each view is noised with an isotropic Gaussian noise of variance $\sigma \sim \mathcal{U}(0.01, 0.2)$. The results are shown in Fig. 9. Methods based on local correspondences (DeepGMR, RPM-Net) are inherently limited, whereas POLAR maintains robust performance. SGHR and RoITr are partially based on global descriptors, hence they keep decent performances in this scenario. This experiment also highlights the benefit brought by the rotation-invariant cross-frame position awareness of RoITr over GeoT.

e) *Point density*: Similarly, we study the case where the views have different point densities. Specifically, the number of point of a view is randomly sampled in $\mathcal{U}(205, 1024)$. The

results are shown in Fig. 10. POLAR is almost invariant to such degradation, and as for the varying noise level experiment, it obtains the best performances, on par with SGHR, closely followed by RoITr and GeoT, since they both partially leverage global descriptors.

f) *Partial visibility*: Point cloud data often suffer from partial occlusion. We study the case of data where a fixed ratio of points has been cropped out for each view (two examples are shown in Fig. 6). The results are shown in Fig. 11. For high overlap, POLAR exhibits the best performance with SGHR, thanks to our occlusion handling masking in the loss (see Eq. (7)). The coarse-to-fine learnable Sinkhorn from attention approach of RoITr and GeoT is the best algorithm for low overlaps up to a certain extent, at the price of a far greater computational complexity (studied in Sec. VI.h). For even lower overlaps, SGHR is the best method, thanks to its NetVLAD module, specifically designed to handle such scenarios. These behaviors highlight specific design choices: the sparse pose graph initialized from local correspondences in SGHR is able to deal with very low overlaps as long as views remain locally similar, whereas the global descriptor of POLAR is suited to handle strong local dissimilarities, at the price of impaired performances for low overlaps. Indeed, the main limiting point of POLAR for low overlaps is the initialization of the template with the medoid of the encoded views (Eq. (17)), which gets increasingly incomplete with the occlusion. An example of template reconstructed with POLAR is shown in Fig. 6. Despite the missing parts of the cropped view, the template is a complete plane.

g) *Outliers*: Finally, the last degradation considered is the presence of outliers. We simulate outliers by sampling points along a random curve starting from the object’s surface, as can be seen on Fig. 6. This simulates the kind of outliers observed in SMLM data (Sec. VI-D) and presents a more challenging case compared to usual simulations such as uniform sampling in the unit sphere, in which case outliers could be segmented out by a dedicated method easily. POLAR exhibits a strong robustness to outliers, on par with that of RoITr and GeoT (Fig. 12) which were specifically crafted to handle wrong correspondences. SGHR performances decrease gradually with the amount of outliers. DeepGMR and RPM-Net are not robust to outliers at all. An example of template obtained by POLAR is shown in Fig. 6. The template is free from outliers.

h) *Time efficiency*: We study the scalability of registration methods with the number of views for the three best performing methods. As pairwise methods leveraging big Transformer networks and approximated optimal transport, the computation cost of RoITr and GeoT increases very rapidly with the number of views. SGHR is also pairwise, but the NetVLAD module estimates the highest overlaps, and only the resulting subset of pairwise registrations is computed, thereby enhancing scalability. As a generative approach, POLAR scales linearly with the number of views, and yields the lowest computation cost.

C. Study of POLAR

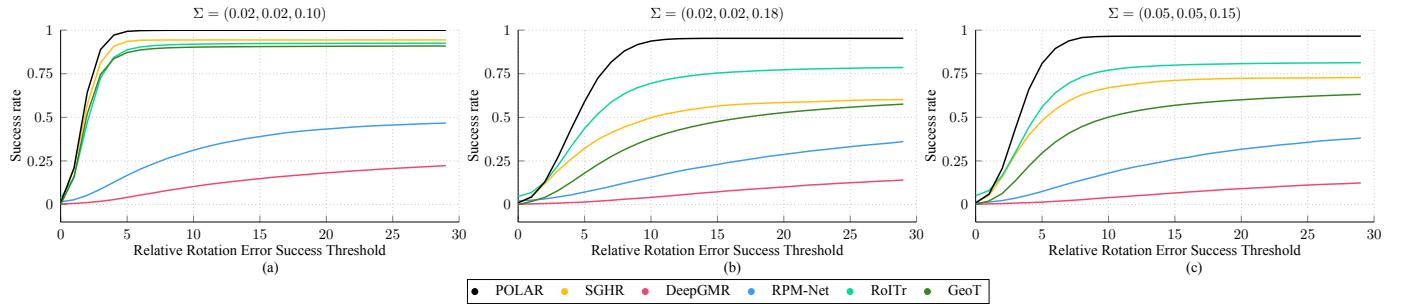


Fig. 8. Cumulative distribution function of the registration recall for three anisotropic noises. (a) Low noise, low anisotropy. (b) Low noise, high anisotropy. (c) High noise, low anisotropy.

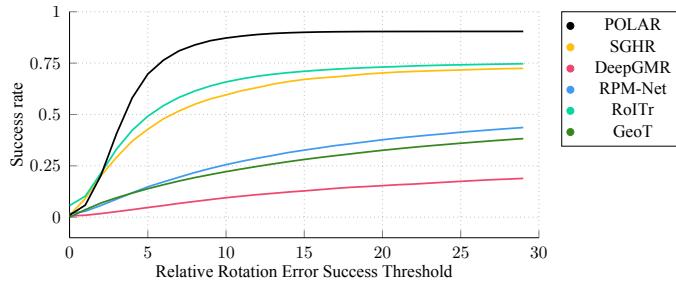


Fig. 9. Cumulative distribution function of registration recall when registering views degraded by a varying noise level ($\sigma \sim \mathcal{U}(0.01, 0.2)$).

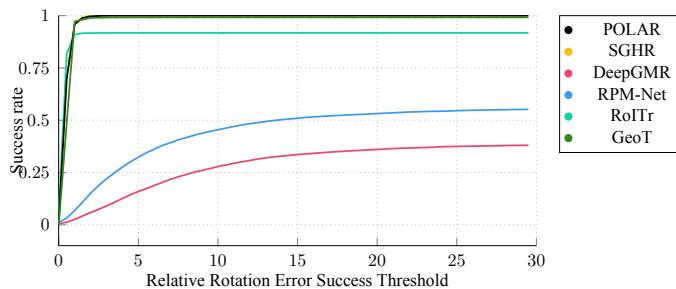


Fig. 10. Cumulative distribution function of registration recall when registering views with varying point density ($n \sim \mathcal{U}(205, 1024)$). The curve of SGHR is not visible because it is superimposed on the POLAR curve.

a) Number of views: We study the impact of the number of views on the registration performances. We generate 100 views combining three degradations (20% of outliers, 20% of occlusion, isotropic noise of standard deviation 0.02). POLAR is optimized on the first n views, with n in $\{5, 10, 25, 50, 50\}$. The results are shown on Fig. 14. Performances gradually

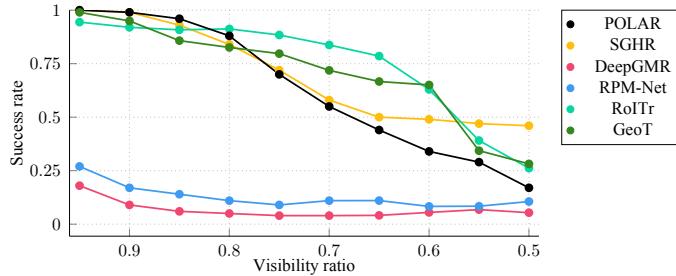


Fig. 11. Evolution of the registration recall with the visibility ratio.

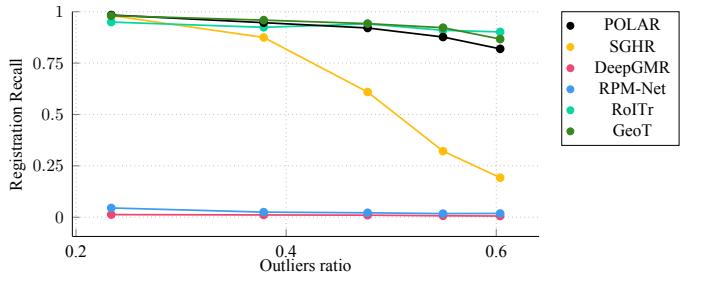


Fig. 12. Evolution of the registration recall with the outliers ratio.

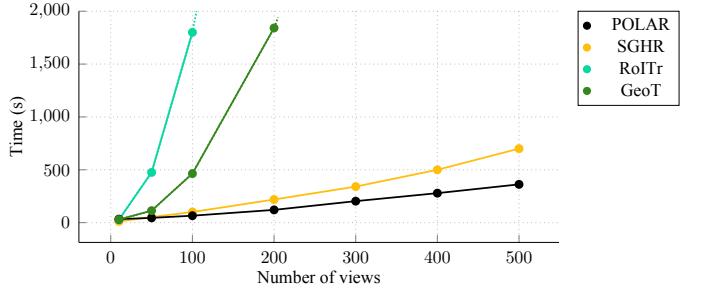


Fig. 13. Evolution of the computation time with the number of views.

increase with the number of views. POLAR benefits from an increasing number of views since it allows the template to be more accurately estimated. This is illustrated in Fig. 15 where the estimated template for $n = 5$ (Fig. 15a) and $n = 100$ (Fig. 15b) are displayed. In the latter case, the reconstruction is complete and more precise.

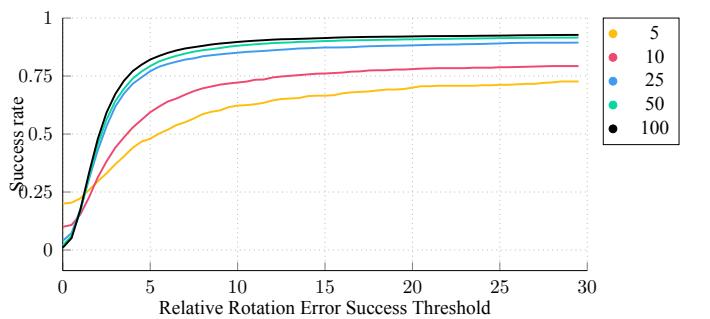


Fig. 14. Evolution of the registration recall of POLAR with the number of views to register.

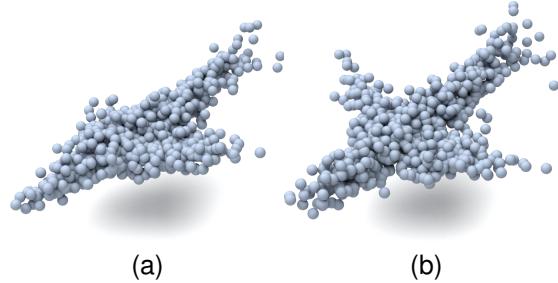


Fig. 15. Visual comparison of template reconstructed by POLAR from data combining the three degradations with (a) 5 and (b) 100 views.

b) Ablation studies: We now study the impact of two components of the criterion (15): the regularization, and the degradation model. The full version of the criterion, *i.e.* the standard POLAR algorithm using the loss (15), is named *degraded + regularized*. In order to examine the influence of the regularization term, we set $\lambda = 0$ in (15), obtaining the *degraded* version. We study a third version of POLAR that we name *regularized*, by keeping the regularization term but replacing in (15) the degraded criterion $L(z, \rho)$ by its basic version (5). These three variations of the optimization criterion are compared on data corrupted by the three types of degradation considered: anisotropic noise (Fig. 16.a), occlusion (Fig. 16.b), and outliers (Fig. 16.c). We observe that the regularization has no impact from a pure registration performance perspective. It should be noted though, that obtained templates with regularization are visually more appealing to a human eye. For each of the addressed degradation, while the basic latent approach already brings decent performance, the degraded version of the criterion greatly enhances the algorithm's robustness.

c) Robustness to visibility and outlier ratios: For the anisotropic case (Fig. 16.a), the noise properties must be known. These properties can often be estimated separately, which is notably the case on SMLM data studied in Sec. VI-D. However, estimating the ratio of outliers and occlusion is more difficult and limiting in practice. To evaluate the sensitivity of POLAR to these ratios, we study the impact of wrong ratio values on the registration performances.

Specifically, we consider two cases: data with 20% cropped out (Fig. 17.a), and data with $\sim 50\%$ of outliers (Fig. 17.b). We then run several registrations using varying values for the visibility ratio v and the outliers ratio o . In both cases, a rough estimate of the ground truth ratio does not severely impact performances, and slightly underestimating the level of degradation improves results. This proves that POLAR is usable even when the visibility and outliers ratios are unknown, making it suitable for application where the degradation model is unknown. Moreover, an extension of the present method could consist in estimating this degradation model alongside the registration task.

D. Real Data

a) Realistic occlusion with FAUST-partial: We now study the registration performances on data from the FAUST-

partial dataset [90]. FAUST-partial is comprised of 100 human body scans. Realistic occlusions are generated by applying the hidden point removal algorithm on icosahedron points around each scan. On these data, each view has a specific visibility ratio. Thus, we select a subset of views from which this ratio is in $[0.7, 1]$ and run POLAR with a fixed ratio of 0.9, which in most cases slightly underestimates the degradation, following the results from Fig. 17.b. POLAR is the best method in that case (Fig. 18). Note that this also highlights the generalization capability of POLAR. Indeed, the shapes here are unknown to the autoencoder. Hence, reconstructed templates are less precise. Nonetheless, the estimated poses are still correct.

b) Highly degraded SMLM data: We use data obtained with direct optical reconstruction microscopy (dStorm) [91], combined with expansion microscopy. The data has been acquired in the group of Markus Sauer in University of Würzburg. It is composed of nine identical particles called centrioles, observed from different views. This is the most challenging dataset, as each view is heavily corrupted with strong anisotropic noise, outliers, missing parts. The SMLM acquisition process provides point clouds with a known 3D uncertainty Σ . The resolution along the microscope axis (z in Fig. 19a) is much worse than in the orthogonal plane, which induces a large deformation of the object in the XZ and YZ planes. The shape of a centriole can be coarsely approximated by a cylinder. Hence, when it is aligned with the microscope axis, the cylinder shape is clearly visible (blue particle in Fig. 19b) whereas other orientations result in a loss of this characteristic shape due to the anisotropy of resolution. The centriole is often attached to tubular structure called microtubules, which can be considered as outliers (see the yellow particle in Fig. 19b). Finally, SMLM data are usually corrupted by a high level of partial visibility, since each point corresponds to the fluorescent emission of a fluorophore, and the fluorophores do not cover uniformly the surface of the particle. Hence, this modality combines a high level of each of the three degradations addressed by POLAR. In Fig. 19a, we show the whole set of nine centrioles, alongside POLAR results. Then in Fig. 19b, for the sake of clarity, we select a subset of three particles to visually compare all methods. Among all the tested methods, POLAR is the only one able to correctly register centrioles.

VII. CONCLUSION

We introduced Point Cloud Latent Registration (POLAR), an algorithm designed to simultaneously estimate a set of rigid motions that align numerous severely degraded views. POLAR integrates a global descriptor derived from a pretrained autoencoder, a global optimization framework, and an informed criterion taking degradations into account. By combining these elements, POLAR demonstrates state-of-the-art performance on both synthetic and real-world datasets that are significantly affected by anisotropic noise, partial visibility, and outliers.

REFERENCES

- [1] H. Heydarian, M. Joosten, A. Przybylski, F. Schueder, R. Jungmann, B. v. Werkhoven, J. Keller-Findeisen, J. Ries, S. Stallinga, M. Bates *et al.*, “3d particle averaging and detection of macromolecular symmetry

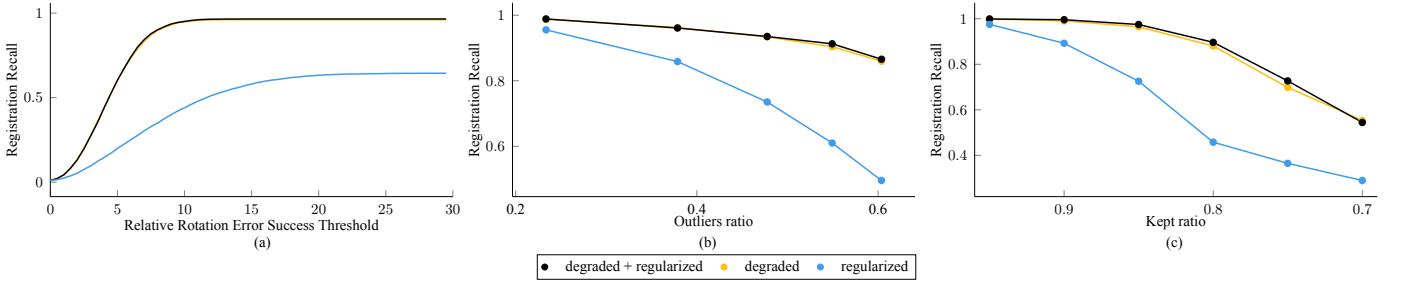


Fig. 16. **Ablation studies on the three degradations.** (a) Cumulative distribution function of registration recall on data corrupted by an anisotropic noise of covariance $\Sigma = \text{diag}(0.02, 0.02, 0.18)$. (b) Registration recall evolution on data gradually corrupted by outliers. (c) Registration recall evolution on data gradually cropped out. *regularized* is the criterion in Eq. (5), without the degradation model, but with the regularization of Sec. III-E-f. *degraded* is the degraded criterion of Eq. (10), but without the regularization. *degraded + regularized* is the criterion of Eq. (15) that combines both the degradation model and the regularization.

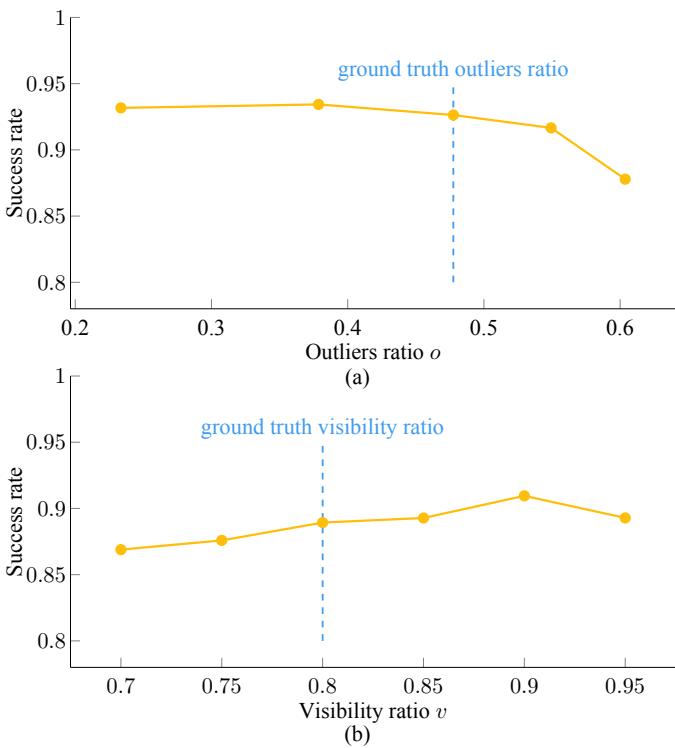


Fig. 17. Registration recall evolution when using varying masking ratios to register the same data.

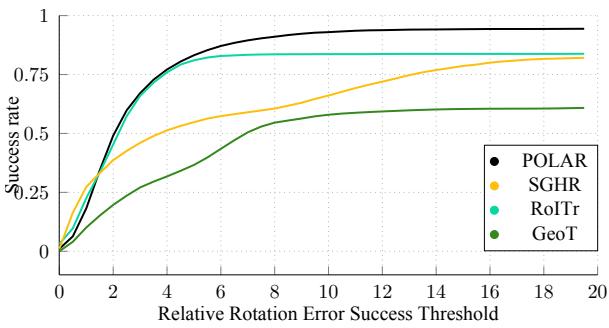


Fig. 18. Cumulative distribution function of the registration recall on realistically occluded data from the Faust partial dataset [90].

- in localization microscopy,” *Nature communications*, vol. 12, no. 1, p. 2847, 2021. 1
- [2] G. Blais and M. D. Levine, “Registering multiview range data to create 3d computer objects,” *IEEE TPAMI*, 1995. 1
 - [3] A. N. üchter, K. Lingemann, J. Hertzberg, and H. Surmann, “6d slam–3d mapping outdoor environments,” *J. Field Robotics*, 2007. 1
 - [4] R. Newcombe, D. Fox, and S. Seitz, “Dynamicfusion: Reconstruction and tracking of non-rigid scenes in real-time,” *CVPR*, 2015. 1
 - [5] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon, “Kinectfusion: Real-time dense surface mapping and tracking,” *IEEE International Symposium on Mixed and Augmented Reality*, 2011. 1
 - [6] J. Tang, D. Xu, K. Jia, and L. Zhang, “Learning parallel dense correspondence from spatio-temporal descriptors for efficient and robust 4d reconstruction,” *CVPR*, 2021. 1
 - [7] Z. Min, J. Wang, and M. Q. H. Meng, “Robust generalized point cloud registration using hybrid mixture model,” *International Conference on Robotics and Automation (ICRA)*, 2018. 1, 2
 - [8] A. Chatterjee and V. Govindu, “Robust relative rotation averagings,” *IEEE TPAMI*, 2018. 1, 2
 - [9] A. Torsello, E. Rodola, and A. Albarelli, “Multiview registration via graph diffusion of dual quaternions,” *CVPR*, 2011. 1, 2
 - [10] E. Maset, F. Arrigoni, and A. Fusiello, “Practical and efficient multi-view matching,” *ICCV*, 2017. 1, 2
 - [11] T. Birdal, U. Simsekli, M. O. Eken, and S. Ilic, “Bayesian pose graph optimization via bingham distributions and tempered geodesic mcmc,” *NeurIPS*, 2018. 1, 2
 - [12] F. Bernard, J. Thunberg, P. Gemmar, F. Hertel, A. Husch, and J. Goncalves, “A solution for multi-alignment by transformation synchronisation,” *CVPR*, 2015. 1, 2
 - [13] F. Arrigoni, B. Rossi, and A. Fusiello, “Spectral synchronization of multiple views in $se(3)$,” *Journal on Imaging Sciences*, 2016. 1, 2
 - [14] F. Arrigoni, L. Magri, B. Rossi, P. Fragneto, and A. Fusiello, “Robust absolute rotation estimation via low-rank and sparse matrix decomposition,” *IEEE International Conference on 3D Vision (3DV)*, 2014. 1, 2
 - [15] M. Arie-Nachimson, S. Z. Kovalsky, I. Kemelmacher-Shlizerman, A. Singer, and R. Basri, “Global motion estimation from point matches,” *International Conference on 3D Imaging, Modeling, Processing, Visualization and Transmission*, 2012. 1, 2
 - [16] H. Wang, Y. Liu, Z. Dong, Y. Guo, Y.-S. Liu, W. Wang, and B. Yang, “Robust multiview point cloud registration with reliable pose graph initialization and history reweighting,” *CVPR*, 2023. 1, 2, 7
 - [17] B. Jian and B. C. Vemuri, “Robust point set registration using Gaussian mixture models,” *IEEE TPAMI*, 2011. 1, 2
 - [18] B. Eckart, K. Kim, A. Troccoli, A. Kelly, and J. Kautz, “Mlmd: Maximum likelihood mixture decoupling for fast and accurate point cloud registration,” *International Conference on 3D Vision*, 2015. 1, 2
 - [19] H. Chui and A. Rangarajan, “A feature registration framework using mixture models,” *IEEE Workshop on Mathematical Methods in Biomedical Image Analysis*, 2000. 1, 2
 - [20] G. D. Evangelidis and R. Horaud, “Joint alignment of multiple point sets with batch and incremental expectation-maximization,” *IEEE TPAMI*, 2017. 1, 2, 7

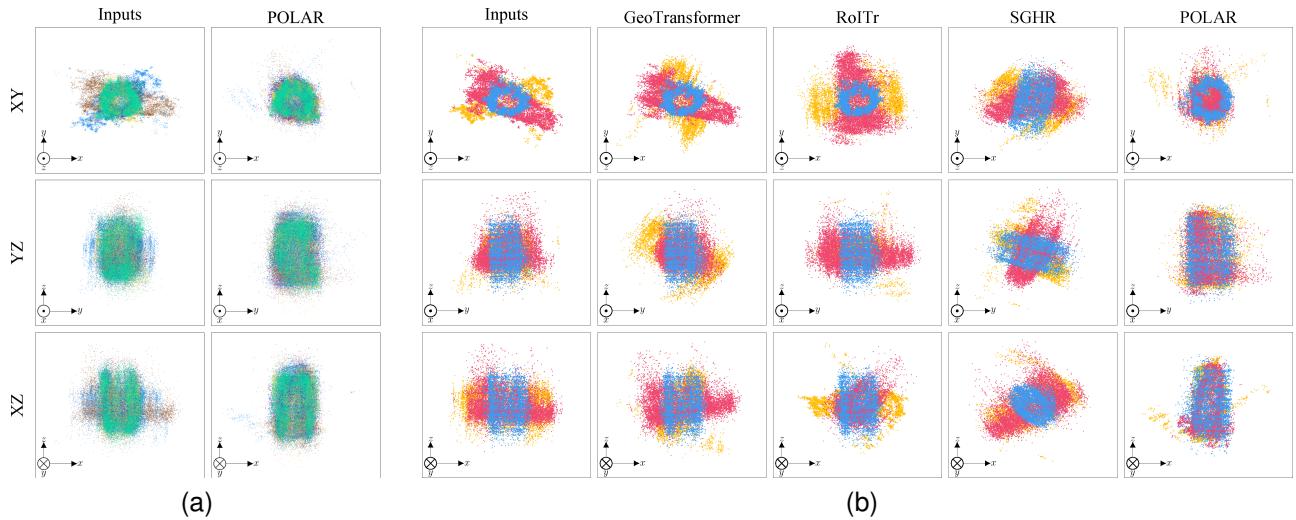


Fig. 19. Visual comparison of registration results on SMLM data. (a) Whole dataset of 9 centrioles, in initial poses (first column) and after registration with POLAR (second column). (b) Selected subset of three centrioles. In each line, the point clouds are observed in a different plane (XY, YZ and XZ). The resolution anisotropy can be seen in the inputs point clouds (columns 1 and 3): the resolution in the plane XY is much higher than in the planes YZ and XZ.

- [21] A. Myronenko and X. Song, “Point set registration : Coherent point drift,” *IEEE TPAMI*, 2010. [1](#) [2](#)
- [22] W. Gao and R. Tedrake, “Filterreg : Robust and efficient probabilistic point-set registration using gaussian filter and twist parameterization,” *CVPR*, 2019. [1](#) [2](#)
- [23] J. Zhu, R. Guo, Z. Li, J. Zhang, and S. Pang, “Registration of multi-view point sets under the perspective of expectation-maximization,” *IEEE TIP*, 2020. [1](#) [2](#) [7](#)
- [24] Y. Ma, J. Zhu, Z. Tian, and Z. Li, “Effective multiview registration of point clouds based on student’s-t mixture model,” *Information Sciences*, 2022. [1](#) [2](#)
- [25] V. Govindu, “Lie-algebraic averaging for globally consistent motion estimation,” *CVPR*, 2004. [2](#)
- [26] X. Huang, Z. Liang, X. Zhou, Y. Xie, L. Guibas, and Q. Huang, “Learning transformation synchronization,” *CVPR*, 2019. [2](#)
- [27] L. Ding and C. Feng, “Deepmapping: Unsupervised map estimation from multiple point clouds,” *CVPR*, 2019. [2](#)
- [28] Z. J. Yew and G. H. Lee, “Learning iterative robust transformation synchronization,” *2021 International Conference on 3D Vision (3DV)*, 2021. [2](#)
- [29] Z. Gojcic, C. Zhou, J. D. Wegner, L. J. Guibas, and T. Birdal, “Learning multiview 3d point cloud registration,” *CVPR*, 2020. [2](#)
- [30] J. Zhang, M. Zhao, X. Jiang, and D. Yan, “Robust multi-view registration of point sets with laplacian mixture model,” *ACCV*, 2021. [2](#)
- [31] Q. Liao, D. wei Sun, S. Zhang, A. Loutfi, and H. Andreasson, “Fuzzy cluster-based group-wise point set registration with quality assessment,” *IEEE TIP*, 2022. [2](#)
- [32] W. Wang and K. Lin, “Information granule-based multi-view point sets registration using fuzzy c-means clustering,” *Multimedia Tools and Applications*, 2022. [2](#)
- [33] P. J. Besl and N. McKay, “A method for registration of 3-d shapes,” *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 1992. [2](#)
- [34] S. Rusinkiewicz and M. Levoy, “Efficient variants of the icp algorithm,” *Third International Conference on 3-D Digital Imaging and Modeling*, 2001. [2](#)
- [35] S. Ao, Q. Hu, B. Yang, A. Markham, and Y. Guo, “Spinnet: Learning a general surface descriptor for 3d point cloud registration,” *CVPR*, 2021. [2](#)
- [36] H. Deng, T. Birdal, , and S. Ilic, “Ppf-foldnet: Unsupervised learning of rotation invariant 3d local descriptors,” *ECCV*, 2018. [2](#)
- [37] H. Yu, J. Hou, Z. Qin, M. Saleh, I. S. Shugurov, K. Wang, B. Busam, and S. Ilic, “Riga: Rotation-invariant and globally-aware descriptors for point cloud registration,” *IEEE TPAMI*, 2022. [2](#)
- [38] H. Yu, Z. Qin, J. Hou, M. Saleh, D. Li, B. Busam, and S. Ilic, “Rotation-invariant transformer for point cloud matching,” *CVPR*, 2023. [2](#) [7](#)
- [39] Z. Qin, H. Yu, C. Wang, Y. Guo, Y. Peng, and K. Xu, “Geometric transformer for fast and robust point cloud registration,” *CVPR*, 2022. [2](#) [7](#) [8](#)
- [40] X. Bai, Z. Luo, L. Zhou, H. Fu, L. Quan, and C.-L. Tai, “D3feat: Joint learning of dense detection and description of 3d local features,” *CVPR*, 2020. [2](#)
- [41] M. Saleh, S. Dehghani, B. Busam, N. Navab, and F. Tombari, “Graphite: Graph-induced feature extraction for point cloud registration,” *International Conference on 3D Vision (3DV)*, 2020. [2](#)
- [42] Z. Gojcic, C. Zhou, J. D. Wegner, and A. Wieser, “The perfect match: 3d point cloud matching with smoothed densities,” *CVPR*, 2019. [2](#)
- [43] Y. Wu, J. Sheng, H. Ding, P. Gong, H. Li, M. Gong, W. Ma, and Q. Miao, “Evolutionary multitasking descriptor optimization for point cloud registration,” *IEEE Transactions on Evolutionary Computation*, 2024. [2](#)
- [44] Y. Yuan, Y. Wu, X. Fan, M. Gong, W. Ma, and Q. Miao, “Egst: Enhanced geometric structure transformer for point cloud registration,” *IEEE Transactions on Visualization and Computer Graphics*, 2024. [2](#)
- [45] Z. J. Yew and G. H. Lee, “Regtr: End-to-end point cloud correspondences with transformers,” *CVPR*, 2022. [2](#)
- [46] Y. Wu, J. Liu, M. Gong, Z. Liu, Q. Miao, and W. Ma, “Mpct: Multiscale point cloud transformer with a residual network,” *IEEE Transactions on Multimedia*, 2024. [2](#)
- [47] C. Choy, J. Park, and V. Koltun, “Fully convolutional geometric features,” *ICCV*, 2019. [2](#)
- [48] H. Deng, T. Birdal, and S. Ilic, “Ppfnet: Global context aware local features for robust 3d point matching,” *CVPR*, 2018. [2](#) [3](#)
- [49] S. Huang, Z. Gojcic, M. M. Usyntsov, A. Wieser, and K. Schindler, “Predator: Registration of 3d point clouds with low overlap,” *CVPR*, 2021. [2](#)
- [50] Y. Li and T. Harada, “Lepard: Learning partial point cloud matching in rigid and deformable scenes,” *CVPR*, 2021. [2](#)
- [51] M. Saleh, S. cheng Wu, L. D. Cosmo, N. Navab, B. Busam, and F. Tombari, “Bending graphs: Hierarchical shape matching using gated optimal transport,” *CVPR*, 2022. [2](#)
- [52] H. Yu, F. Li, M. Saleh, B. Busam, and S. Ilic, “Cofinet: Reliable coarse-to-fine correspondences for robust point cloud registration,” *NeurIPS*, 2021. [2](#)
- [53] A. Zeng, S. Song, M. Nießner, M. Fisher, J. Xiao, and T. Funkhouser, “3dmatch : Learning local geometric descriptors from rgbd reconstructions,” *CVPR*, 2017. [2](#)
- [54] Y. Zhang, J. Yu, X. Huang, W. Zhou, and J. Hou, “Pcr-cg: Point cloud registration via deep explicit color and geometry,” in *ECCV*, 2023. [2](#)
- [55] M. A. Fischler and R. C. Bolles, “Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography,” *COMMUN ACM*, 1981. [2](#)
- [56] D. Baráth and J. Matas, “Graph-cut ransac,” *CVPR*, 2017. [2](#)

- [57] S. Quan and J. Yang, “Compatibility-guided sampling consensus for 3-d point cloud registration,” *IEEE Transactions on Geoscience and Remote Sensing*, 2020. [2](#)
- [58] R. B. Rusu, N. Blodow, and M. Beetz, “Fast point feature histograms (fpfh) for 3d registration,” *IEEE International Conference on Robotics and Automation (ICRA)*, 2009. [2](#)
- [59] J. Yang, J. Chen, S. Quan, W. Wang, and Y. Zhang, “Correspondence selection with loose-tight geometric voting for 3-d point cloud registration,” *IEEE Transactions on Geoscience and Remote Sensing*, 2022. [2](#)
- [60] J. Yang, Z. Huang, S. Quan, Z. Qi, and Y. Zhang, “Sac-cot: Sample consensus by sampling compatibility triangles in graphs for 3-d point cloud registration,” *IEEE Transactions on Geoscience and Remote Sensing*, 2021. [2](#)
- [61] J. Yang, Z. Huang, S. Quan, Q. Zhang, Y. Zhang, and Z. Cao, “Toward efficient and robust metrics for ransac hypotheses and 3d rigid registration,” *IEEE Transactions on Circuits and Systems for Video Technology*, 2021. [2](#)
- [62] Y. Aoki, H. Goforth, R. A. Srivatsan, and S. Lucey, “Pointnetlk: Robust and efficient point cloud registration using pointnet,” *CVPR*, 2019. [2](#), [7](#)
- [63] V. Sarode, X. Li, H. Goforth, Y. Aoki, R. A. Srivatsan, S. Lucey, and H. Choset, “Pcnnet: Point cloud registration network using pointnet encoding,” in *ArXiv e-prints*, 2019. [2](#)
- [64] X. Huang, G. Mei, and J. Zhang, “Feature-metric registration: A fast semi-supervised approach for robust point cloud registration without correspondences,” *CVPR*, 2020. [2](#)
- [65] H. Xu, S. Liu, G. Wang, G. Liu, and B. Zeng, “Omnnet: Learning overlapping mask for partial-to-partial point cloud registration,” *ICCV*, 2021. [2](#)
- [66] M. Zhu, M. Ghaffari, and H. Peng, “Correspondence-free point cloud registration with so(3)-equivariant implicit shape representations,” *Conference on Robot Learning*, 2022. [2](#)
- [67] J. Yang, H. Li, D. Campbell, and Y. Jia, “Go-icp: A globally optimal solution to 3d icp point-set registration,” *IEEE TPAMI*, 2016. [2](#)
- [68] H. Yang, P. Antonante, V. Tzoumas, and L. Carlone, “Graduated non-convexity for robust spatial perception: From non-minimal solvers to global outlier rejection,” *International Conference on Robotics and Automation (ICRA)*, 2019. [2](#)
- [69] Q.-Y. Zhou, J. Park, and V. Koltun, “Fast global registration,” *ECCV*, 2016. [2](#), [7](#)
- [70] R. I. Hartley and F. Kahl, “Global optimization through searching rotation space and optimal estimation of the essential matrix,” *ICCV*, 2007. [2](#)
- [71] O. Enqvist and F. Kahl, “Robust optimal pose estimation,” *ECCV*, 2008. [2](#)
- [72] C. Olsson, F. Kahl, and M. Oskarsson, “Branch-and-bound methods for euclidean registration problems,” *IEEE TPAMI*, 2009. [2](#)
- [73] W. Yuan, B. Eckart, K. Kim, V. Jampani, D. Fox, and J. Kautz, “Deepgmr : Learning latent gaussian mixture models for registration,” *ECCV*, 2020. [2](#), [7](#)
- [74] S. Xie, J. Gu, D. Guo, C. R. Qi, L. Guibas, and O. Litany, “Pointcontrast: Unsupervised pretraining for 3d point cloud understanding,” *ECCV*, 2020. [2](#), [3](#)
- [75] Y. Shen, L. Hui, H. Jiang, J. Xie, and J. Yang, “Reliable inlier evaluation for unsupervised point cloud registration,” *AAAI*, 2022. [2](#)
- [76] G. Elbaz, T. Avraham, and A. Fischer, “3d point cloud registration for localization using a deep neural network auto-encoder,” *CVPR*, 2017. [2](#)
- [77] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, “Pointnet: Deep learning on point sets for 3d classification and segmentation,” *CVPR*, 2017. [3](#), [6](#)
- [78] M. Alexa, “Super-fibonacci spirals: Fast, low-discrepancy sampling of so(3),” *CVPR*, 2022. [5](#)
- [79] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *ICLR*, 2014. [6](#)
- [80] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” *ICLR*, 2019. [6](#), [7](#)
- [81] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style, high-performance deep learning library,” *NeurIPS*, 2019. [6](#)
- [82] J. M. Lee, *Introduction to Smooth Manifolds*. Springer New York, NY, 2012. [6](#)
- [83] M. W. Hirsch, *Differential Topology*. Springer New York, NY, 1976. [6](#)
- [84] M. Fey and J. E. Lenssen, “Fast graph representation learning with PyTorch Geometric,” *ICLR*, 2019. [7](#)
- [85] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, “3d shapenets: A deep representation for volumetric shapes,” *CVPR*, 2015. [7](#), [8](#)
- [86] Z. J. Yew and G. H. Lee, “Rpm-net: Robust point matching using learned features,” *CVPR*, 2020. [7](#), [8](#)
- [87] Y. Wang and J. M. Solomon, “Deep closest point : Learning representations for point cloud registration,” *ICCV*, 2019. [7](#)
- [88] X. Zhang, J. Yang, S. Zhang, and Y. Zhang, “3d registration with maximal cliques,” *CVPR*, 2023. [7](#)
- [89] Y. Shi and G. Lerman, “Message passing least squares framework and its application to rotation synchronization,” *ICML*, 2020. [7](#)
- [90] D. Bojanic, K. Bartol, J. Forest, T. Petkovic, and T. Pribanic, “Addressing the generalization of 3d registration methods with a featureless baseline and an unbiased benchmark,” *Machine Vision and Applications*, 2024. [8](#), [11](#), [12](#)
- [91] M. Heilemann, S. V. D. Linde, M. SchuNtpelz, R. Kasper, B. Seefeldt, A. Mukherjee, P. Tinnefeld, and M. Sauer, “Subdiffraction-resolution fluorescence imaging with conventional fluorescent probes,” *Angewandte Chemie International Edition*, 2008. [8](#), [11](#)
- [92] H. Wang, Y. Liu, Z. Dong, and W. Wang, “You only hypothesize once: Point cloud registration with rotation-equivariant descriptors,” *ACM MM*, 2022. [9](#)
- [93] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, “Netvlad: Cnn architecture for weakly supervised place recognition,” *IEEE TPAMI*, 2018. [9](#)