
Unsupervised Non-Rigid Point Cloud Matching through Large Vision Models

Anonymous Author(s)

Affiliation

Address

email

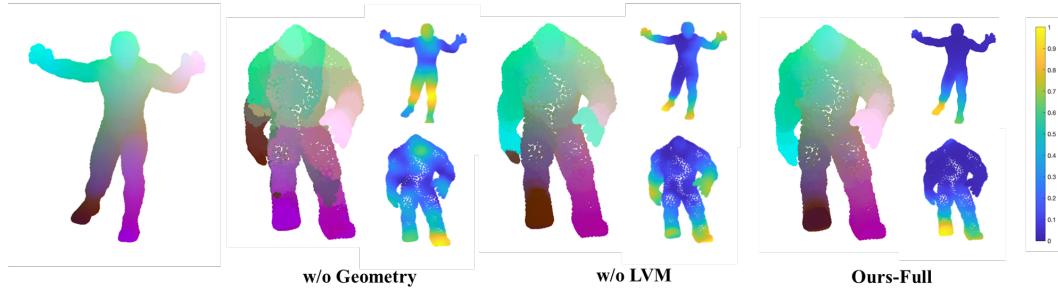


Figure 1: We match a challenging point cloud pair alien to the training set and visualize both maps and feature alignment. Our result surpasses that from purely LVM-based (w/o Geometry) and from purely geometry-based (w/o LVM) baselines in both ends. See text for more details.

Abstract

1 In this paper, we propose a novel learning-based framework for non-rigid point
2 cloud matching, which can be trained *purely* on point clouds without any correspon-
3 dence annotation but also be extended naturally to partial-to-full matching. Our key
4 insight is to incorporate semantic features derived from large vision models (LVMs)
5 to geometry-based shape feature learning. Our framework effectively leverages the
6 structural information contained in the semantic features to address ambiguities
7 arise from self-similarities among local geometries. Furthermore, our framework
8 also enjoys the strong generalizability and robustness regarding partial observations
9 of LVMs, leading to improvements in the regarding point cloud matching tasks.
10 In order to achieve the above, we propose a pixel-to-point feature aggregation
11 module, a local and global attention network as well as a geometrical similarity
12 loss function. Experimental results show that our method achieves state-of-the-art
13 results in matching non-rigid point clouds in both near-isometric and heterogeneous
14 shape collection as well as more realistic partial and noisy data.

15

1 Introduction

16 Estimating dense correspondences between non-rigid 3D shapes is a fundamental task in computer
17 vision and graphics, which plays a pivotal role in an array of applications including, including
18 3D reconstruction [63], 3D pose estimation [58], and animation [51] among others. In contrast
19 to the significant progress [10, 56] on matching well-structured shapes (*i.e.*, triangular meshes),
20 the advancement on *unstructured point clouds* is relatively lagged. Meanwhile, motivated by the
21 prevalence and ease of point cloud data scanning in practice, we propose a novel *unsupervised*
22 non-rigid shape matching framework, which is trained *purely* on point cloud base. Moreover, our

23 framework can be naturally and effectively extended to partial point cloud matching, which fits
24 squarely with the partial nature of point cloud scanning.

25 Before diving into our framework, we briefly overview the prior arts. Early approaches on unsu-
26 pervised non-rigid point cloud matching [23, 34, 64, 17] leverage point cloud reconstruction as a
27 proxy task to learn embeddings without correspondence labels. However, the reconstruction-based
28 approaches attain training ease at the cost of intrinsic geometric understanding – the widely used
29 Chamfer loss is purely *extrinsic*, falling short of efficiently capturing the geometric details of the
30 underlying surface. To this end, a recent trend [42, 30, 11, 31] is to transplant the success of
31 matching triangular meshes into the domain of point clouds. In essence, such approaches follow
32 a self-supervised scheme built on the bijective mapping between a mesh and its vertex set. While
33 achieving more competing performance than the above, these methods all require triangular meshes
34 as training data, either for full or partial shape matching [11]. This constraint undoubtedly limits their
35 practical utility, especially in the tasks where high-quality meshes are non-trivial to achieve (*e.g.*, 3D
36 medical data).

37 Facing the aforementioned challenges, our key insight is to go *beyond* geometry and leverage a
38 powerful tool from the world of a different dimension – large vision models (LVMs). First of all,
39 LVMs are typically trained on datasets not only orders of magnitudes larger, but also significantly
40 more versatile than the best possible in 3D domain. Based on such, LVMs have shown remarkable
41 generalizability on understanding ever-changing objects in the wild, which would be highly desired
42 in enhancing generalizability for our model. Secondly, LVMs are trained on images, which inherently
43 are partial (in the sense of 3D world). The induced capacity of robustly encoding partiality is beyond
44 valuable in dealing with partial point cloud matching.

45 In fact, there has been a recent trend on utilizing LVMs in 3D vision tasks, such as 3D model
46 pre-training [65], shape segmentation [20], keypoint detection [60], and, most relevantly, 3D shape
47 matching [1, 46]. Despite the simplicity and effectiveness, these methods essentially leverages LVMs
48 as virtual annotator for generating auxiliary cues to assist 3D vision tasks. For instance, [46] uses
49 DinoV2 to generate sparse landmarks as input to a second-stage neural surface mapping. Since
50 DinoV2 (among many LVMs) is not designed for dense, fine-grained matching, such strategy can
51 be sensitive to noisy LVM features. Moreover, most of the above are built on mesh render, which is
52 non-trivial to extend to point clouds, especially those of partiality.

53 In light of the above, we propose for the first time a fine-grained end-to-end feature learning framework
54 for unsupervised non-rigid point cloud matching, which make full use of both 3D geometric learning
55 and pre-trained LVMs. Firstly, we propose an efficient module for aggregating pixel-to-point features,
56 which adeptly assigns pixel-wise 2D representations to point-wise 3D point clouds. Secondly,
57 we introduce a novel local and global attention network that refines the integration of visual and
58 geometric features, transforming them into a more fine-grained canonical space. This attention
59 network enhances the model’s ability to capture details and complex relationships within the data.
60 Finally, in addition to the conventional reconstruction loss, we propose a new geometry loss designed
61 to further improve the model’s performance by encouraging the preservation of geometric integrity.

62 Fig. 1 demonstrate a challenging task, in which we *directly* infer a pair of alien shapes from DT4D-H
63 dataset, which are distinctive from SCAPE [5], the training set. LVM-dominated feature (w/o Geome-
64 try) leads to patch-wise correspondences, echoing the fact that it generally captures coarse semantics.
65 On the other hand, geometry-dominated feature (w/o LVM) delivers smoother correspondences but
66 fails to encode global structure, leading to severe mismatches around the hands. Finally, taking the
67 best of both worlds, our learned feature leads to smooth and precise maps, surpassing the former two
68 by a large margin. We also visualize the feature of the first channel on both shapes next to the target.
69 It is obvious that our feature is by far more localized and cleaner than the counterparts.

70 We conduct a rich set of experiments to verify the effectiveness of our pipeline, highlighting that
71 it achieves state-of-the-art results in matching non-rigid point clouds in both near-isometric and
72 heterogeneous shape collections. Remarkably, it generalizes well despite the distinctiveness between
73 the training set and test set. Moreover, our framework outperforms the competing methods in partial
74 point clouds and noisy real scans.

75 **2 Related Works**

76 **2.1 Non-rigid Shape Matching**

77 Non-rigid shape matching is a long-standing problem in computer vision and graphics. Unlike the
78 rigid counterpart, non-rigidly aligning shapes is more challenging owing to the complexity inherent
79 in deformation models.

80 Originating from the foundational work on functional maps [50], along with a series of follow-
81 ups [47, 27, 52, 44, 28, 39, 54, 10, 36, 19, 6, 56], spectral methods have made significant progress in
82 addressing the non-rigid shape matching problem, yielding state-of-the-art performance. However,
83 because of the heavy dependence of Laplace-Beltrami operators, DFM can suffer notable performance
84 drop when applied to point clouds without adaptation [11]. In fact, inspired by the success of DFM,
85 several approaches [30, 11, 31] have been proposed to leverage intrinsic geometry information carried
86 by meshes in the training of feature extractors tailored for non-structural point clouds. When it comes
87 to pure point cloud matching, there is a line of works[64, 34, 17] leverage point cloud reconstruction
88 as the proxy task to learn embeddings without correspondence labels. Since intrinsic information is
89 not explicitly formulated in these methods, they can suffer from significant intrinsic deformations
90 and often generalize poorly to unseen shapes.

91 **2.2 Large Vision Model for 3D Shape Analysis**

92 Recently, Large Vision Models have become increasingly popular in due to their remarkable ability
93 to understand data distributions from extensive image datasets. In the fields of shape analysis,
94 [65] proposes an alternative to obtain superior 3D representations from 2D pre-trained models via
95 Image-to-Point Masked Auto-encoders. [1] introduces a fully multi-stage method that exploits
96 the exceptional reasoning capabilities of recent foundation models in language [48] and vision[35]
97 to tackle difficult shape correspondence problems. In [46], before surface matching, the authors
98 propose to use the features extracted from DINOv2 [49] of multi-view images of the shapes to
99 perform co-alignment. In contrast to these approaches, which primarily utilize coarse patch features
100 for sparse landmarks or semantic matching, our approach introduces an end-to-end method that
101 aggregates pixel-level 2D features into point-wise 3D features.

102 **2.3 Non-rigid Partial Shape Matching**

103 While significant advancements have been made in full shape matching, there remains considerable
104 room for improvement in estimating dense correspondences between shapes with partiality. Functional
105 maps representation [53, 6, 11] has already been applied to partial shapes. However, both axiomatic
106 and learning-based lines of work typically assume the input to be a *connected mesh*, with the exception
107 of [11], which relies on graph Laplacian construction [55] in its preprocessing. For partial point
108 cloud matching, axiomatic registration approaches [4, 62, 38] assume the deformation of interest
109 can be approximated by local, small-to-moderate, rigid deformations, therefore suffer from large
110 intrinsic deformations. Simultaneously, there's a growing trend towards integrating deep learning
111 techniques [9, 8, 26, 37]. However, these methods often focus on addressing the partial sequence
112 point cloud registration problem.

113 **3 Methodology**

114 Fig. 2 shows the overall pipeline. We first introduce our pixel-to-point feature aggregation method in
115 Sec.3.1. Then our global and local attention network will be discussed in Sec.3.2. The training losses
116 are described in Sec.3.3.

117 **3.1 Pixel to Point Feature Aggregation**

118 **Depth aware projection:** As shown in Fig2(b), given a point cloud P consisting of N points,
119 we denote the i -th point by $p_i = (x_i, y_i, z_i)$. Following I2P-MAE [65], we project P on
120 $xy-$, $yz-$, $xz-$ plane to obtain three images. Without loss of generality, we consider the xy -plane
121 and let (u_i, v_i) be the projected pixel of p_i , which is computed as follows:

$$u_i = \lfloor \frac{x_i - x_{\min}}{\Delta_{xy}} \times H \rfloor, v_i = \lfloor \frac{y_i - y_{\min}}{\Delta_{xy}} \times W \rfloor, \quad (1)$$

122 where $\Delta = \max\{x_{\max} - x_{\min}, y_{\max} - y_{\min}\}$ and x_{\min}, x_{\max} are the minimum and maximum of
123 the x -coordinates of P , and H, W are the pre-determined image resolution. On the other hand,

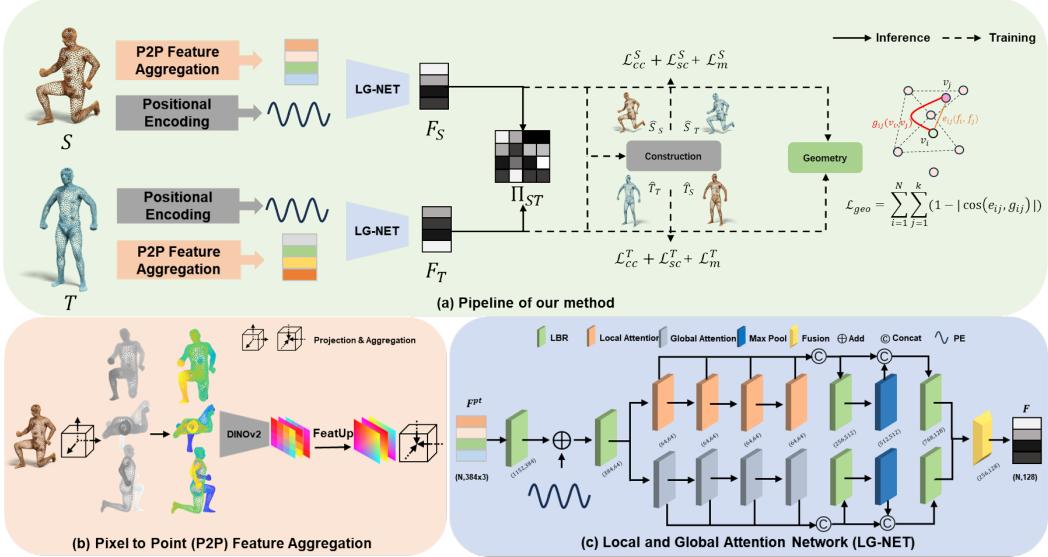


Figure 2: The schematic illustration of our pipeline.

124 since z -coordinates are eliminated, I2P-MAE proposes to assign $\text{sigmod}(z_i)$ as the intensity of
125 $(u_i, v_i), i = 1, 2, \dots, N$, and 0 for all the unprojected pixels.

126 Despite of the simplicity, the above scheme eventually gray images with holes (as only projected
127 pixels carry non-zero intensity), which are distinctive from the realistic training images used in LVMs.
128 To alleviate the discrepancy, we propose to 1) apply a 3×3 mean filter on the gray images and 2)
129 assign pseudo color on the pixel values with the PiYG colormap in MATLAB. We now denote by
130 $I_{\hat{z}}, I_{\hat{x}}, I_{\hat{y}}$ to resulting images, where \hat{x} indicates projection onto xy -plane (and similarly for \hat{x}, \hat{y}).

131 **2D to 3D feature aggregation:** In I2PMAE [65], the obtained three projected images are fed into
132 DINOv2 [49], resulting in a $D \times D \times C$ feature. Note that D is typically around 16, which is much
133 smaller than H (or W). Thanks to the recent advance on *super-resolution* features from LVMs –
134 FeatUp [21], we get rid of the resolution degradation and obtain

$$F_{\hat{z}}^{img} = \Theta(I_{\hat{z}}) \in \mathbb{R}^{H \times W \times C}, \quad (2)$$

135 where Θ is the per-pixel encoder of DinoV2-FeatUp [21] and C is the number of channels for each
136 pixel. According to Eqn. 1, we obtain the point-wise feature of p_i via a simple pull-back:

$$f_{\hat{z}}^i = F_{\hat{z}}(u_i, v_i, :) \in \mathbb{R}^C. \quad (3)$$

137 We then have $F_{\hat{z}}^{pt} \in \mathbb{R}^{N \times C}$ by stacking $f_{\hat{z}}^i$ in order. We compute $F_{\hat{x}}^{pt}, F_{\hat{y}}^{pt}$ in the same manner. We
138 emphasize that these computations are independent, as the pixel-point maps (Eqn. 1) vary. In the end,
139 we arrive at

$$F^{pt}(P) = [F_{\hat{z}}^{pt}, F_{\hat{x}}^{pt}, F_{\hat{y}}^{pt}] \in \mathbb{R}^{N \times 3C}. \quad (4)$$

140 The above procedure returns a set of per-point features for the input P , which essentially carry the
141 semantic information extracted by the LVM.

142 3.2 Local and Global Attention Network

143 In this part, we describe our Local and Global attention Network, which are depicted in Fig. 2(c).

144 **Input feature:** Given a point cloud P , we have computed the semantic per-point features based on
145 LVMs (Eqn. 4). In order to exploit both semantic (image-based) and geometric (point-based) features,
146 we propose to perform early fusion at the input stage as follows:

$$F^{in}(P) = \text{LBR}(F^{pt}(P)) + \gamma(P). \quad (5)$$

147 where $\gamma(P) \in \mathbb{R}^{N \times 384}$ is the positional encoding [45] and LBR is a module proposed in PCT [24]
148 for non-linearly converting $F^{pt}(P)$ into the same dimension of $\gamma(P)$.

149 **Architecture design:** In non-rigid shape matching, both local and global features provide critical
150 information for finding the precise correspondence. Intuitively, local features guide fine-detailed

151 matching while suffering from global structure ambiguity due to self-similarities across the shape.
 152 Global features, on the other hand, provide structural descriptions for making full use of the local
 153 ones.
 154 Motivated by the above, we propose a dual-pathway architecture in parallel, comprising global
 155 attention [61] and local attention [24] blocks. For each point, the global attention blocks systematically
 156 survey the features of the remaining points to achieve comprehensive global perceptual awareness.
 157 On the other hand, the receptive field of the local attention blocks is constrained to the local neighbor-
 158 hood of a point. In particular, we highlight the key difference between the usage of local attention
 159 blocks in PCT [24] and ours: The former used a fixed neighborhood computed w.r.t the input spatial
 160 distribution, while we employ KNN search to connect with the nearest k features in the *latent space*.
 161 This design is inspired by DGCNN [59], and motivated by a typical challenge in non-rigid point cloud
 162 matching. Namely, points with small Euclidean distance are not necessarily close on the surface (*e.g.*,
 163 when a human put hand close to head).
 164 The above dual-path design enables the extraction of profound information through a combination
 165 of layer-wise progression and cross-layer interactions. In the end, we leverage the fusion module,
 166 consisting of LBR and a three-layer stacked N2P [61] attention, to merge features from both global
 167 and local paths, resulting in our output feature. We refer readers to the appendix for more details.

168 3.3 Training Objectives and Matching Inference

169 In the following, we introduce the training losses, which consist of both reconstruction-based losses
 170 and our novel geometrical similarity loss. As shown in Fig. 2(a), our main model is a Siamese network.
 171 Given a pair of shapes (point clouds) $\mathcal{S}, \mathcal{T}, F_{\mathcal{S}}, F_{\mathcal{T}}$ are obtained by passing through the pixel-to-point
 172 feature aggregation (Sec. 3.1) and LG-Net (Sec. 3.2). We then estimate dense correspondences, $\Pi_{\mathcal{ST}}$,
 173 based on the cosine similarity between features of each pair of points $x_i \in \mathcal{S}, y_j \in \mathcal{T}$, which is used
 174 to define the following losses.

175 **Cross- and self-reconstruction losses** are proposed in DPC [34], which aims to cross-construct the
 176 shape using latent proximity between source and target points and the coordinates. Specifically, the
 177 cross-construction process is computed as follows:

$$\hat{y}_{x_i} = \sum_{j \in \mathcal{N}_{\mathcal{T}}(x_i)} \frac{e^{s_{ij}}}{\sum_{l \in \mathcal{N}_{\mathcal{T}}(x_i)} e^{s_{il}}} y_j, \quad (6)$$

178 where $x_i \in \mathcal{S}, y_j \in \mathcal{T}$, and s_{ij} is the cosine similarity between the latent feature of them. $\mathcal{N}_{\mathcal{T}}$
 179 represents the latent k-nearest neighbors of x_i in the target $F_{\mathcal{T}}$. The cross-construction of \mathcal{T} by the
 180 source point cloud \mathcal{S} is denoted $\hat{\mathcal{T}}_{\mathcal{S}} \in \mathbb{R}^{N \times 3}$, where $\hat{\mathcal{T}}_{\mathcal{S}}^i = \hat{y}_{x_i}$. The cross-construction is then be
 181 defined as:

$$\mathcal{L}_{cc} = \mathbf{CD}(\mathcal{S}, \hat{\mathcal{S}}_{\mathcal{T}}) + \mathbf{CD}(\mathcal{T}, \hat{\mathcal{T}}_{\mathcal{S}}), \quad (7)$$

182 where **CD** denotes Chamfer distance. In addition to the cross-construction Loss, we further employ
 183 a loss to enhance the smoothness within neighborhoods. This loss is equivalent to a special case of
 184 cross-construction, namely, the self-construction loss:

$$\mathcal{L}_{sc} = \mathbf{CD}(\mathcal{S}, \hat{\mathcal{S}}_{\mathcal{S}}) + \mathbf{CD}(\mathcal{T}, \hat{\mathcal{T}}_{\mathcal{T}}). \quad (8)$$

185 **Mapping loss** is also proposed in DPC [34], which enforces the mapped points of neighboring points
 186 in \mathcal{S} to be close to each other as well in \mathcal{T} . Specifically, it is defined as

$$\mathcal{L}_m = \mathcal{L}_m^{(\mathcal{S}, \hat{\mathcal{S}}_{\mathcal{T}})} + \mathcal{L}_m^{(\mathcal{T}, \hat{\mathcal{T}}_{\mathcal{S}})}, \quad (9)$$

187 where $\mathcal{L}_m^{(\mathcal{S}, \hat{\mathcal{S}}_{\mathcal{T}})}$ denotes the mapping loss from \mathcal{S} to \mathcal{T}

$$\mathcal{L}_m^{(\mathcal{S}, \hat{\mathcal{S}}_{\mathcal{T}})} = \frac{1}{Nk_m} \sum_i \sum_{l \in \mathcal{N}_{\mathcal{S}}(x_i)} e^{-\|x_i - x_l\|_2^2/\alpha} \|\hat{y}_{x_i} - \hat{y}_{x_l}\|_2^2, \quad (10)$$

188 where k_m, α are fixed constants. We then define similarly on the opposite direction $\mathcal{L}_m^{(\mathcal{T}, \hat{\mathcal{T}}_{\mathcal{S}})}$.

189 **Geometrical similarity:** Previous methods [34, 17, 25] mainly rely on the above losses. It is worth
 190 noting, though, the involved cosine similarity emphasizes more on the *angular* difference between

191 features, which falls short of constraining features from the perspective of *magnitude*. Inspired by
 192 such, we propose geometrical similarity loss for taking magnitude into consideration.

193 In particular, we notice NIE [30] leverages geodesic supervision to enforce the Euclidean distance
 194 between learned feature to ensemble geodesic. However, estimating accurate geodesics on point
 195 clouds is a non-trivial but also heavy task. We therefore adopt the heat method [16] for point clouds
 196 to compute \mathbf{M}_S as the approximated geodesic distance matrix of S .

197 On the other hand, for each learned feature of $x_i \in S$, F_S^i , we consider $\text{NN}(i) = \{j_1, j_2, \dots, j_k\}$ be
 198 the set of ordered indices of the nearest neighborhood in the latent space and compute the distance
 199 vector $d_S^i = [\|F_S^i - F_S^{j_1}\|, \|F_S^i - F_S^{j_2}\|, \dots, \|F_S^i - F_S^{j_k}\|]$. Similarly, we can construct another
 200 vector given \mathbf{M}_S and $\text{NN}(i)$, i.e., $m_S^i = [\mathbf{M}_S(i, j_1), \mathbf{M}_S(i, j_2), \dots, \mathbf{M}_S(i, j_k)]$.

201 While it seems natural to minimize the residual between the above two vectors for each point, we opt
 202 for the following loss given the potential noise in estimating geodesics on unstructured point clouds:

$$\mathcal{L}_{geo}^S = \frac{1}{N} \sum_{i=1}^N \left(1 - \frac{d_S^i \cdot m_S^i}{\|d_S^i\| \|m_S^i\|} \right). \quad (11)$$

203 Similarly, we define the geometrical similarity loss for \mathcal{T} : $\mathcal{L}_{geo} = \mathcal{L}_{geo}^S + \mathcal{L}_{geo}^T$.

204 The overall objective function of our point correspondence learning scheme is:

$$\mathcal{L}_{total} = \lambda_{cc} \mathcal{L}_{cc} + \lambda_{sc} \mathcal{L}_{sc} + \lambda_m \mathcal{L}_m + \lambda_{geo} \mathcal{L}_{geo}, \quad (12)$$

205 where $\lambda_{cc}, \lambda_{sc}, \lambda_m, \lambda_{geo}$ are hyper-parameters, balancing the contribution of the different loss terms.

206 **Partial matching loss:** In the above, we entail the losses for training full-to-full non-rigid point
 207 cloud matching. Remarkably, our formulation can be easily extended to the challenging scenario of
 208 partial-to-full matching. In fact, we simply modify \mathcal{L}_{cc} to a unilateral loss, i.e. only $\mathbf{CD}(\mathcal{S}, \hat{\mathcal{S}}_{\mathcal{T}})$ is
 209 considered, and set $\lambda_m = 0$.

210 **Inference:** At inference time, we choose the nearest latent cross-neighborhood of $x_i \in S$ to be its
 211 corresponding point by KNN [15], thus get the shape matching result between point cloud S and \mathcal{T} .

212 4 Experiments

213 **Dataset:** We evaluate our method with several state-of-the-art techniques for estimating correspondences on a set of benchmarks as follows. **SCAPE_r:** The remeshed version of the SCAPE dataset[5]
 214 comprises 71 human shapes. We split the first 51 shapes for training and the rest 20 shapes for testing;
 215 **FAUST_r:** The remeshed version of FAUST dataset [7] comprises 100 human shapes. We split
 216 the first 80 shapes for training and the rest 20 for testing. **SHREC'19_r:** The remehsed version of
 217 SHREC19 dataset[43] comprises 44 shapes. We pair them into 430 annotated examples provided by
 218 [43] for testing. **DT4D-H:** A dataset from [41] comprises 10 categories of heterogeneous humanoid
 219 shapes. Following [31], we use it solely in testing, and evaluating the inter-class maps split in [41].
 220 **SHREC'07-H:** A subset of SHREC'07 dataset [22] comprises 20 heterogeneous human shapes.
 221 We use it solely in testing. **TOSCA:** Dataset from [66] comprises 41 different shapes of various
 222 animal species. Following [34, 17], we pair these shapes to create both for training and evaluation,
 223 respectively. **SHREC'16:** Partial shape dataset SHREC'16 [14] includes two subsets, namely CUTS
 224 with 120 pairs and HOLES with 80 pairs. Following [6, 11], we train our method for each subset
 225 individually and evaluate it on the corresponding unseen test set (200 shapes for each subset). More-
 226 over, we further conduct some practical experiments on partial real scan dataset processed from [33]
 227 and medical dataset from [3].

228 **Baseline:** We compare our method with a set of competitive baselines, including methods that can
 229 both train and test on point cloud; methods required mesh for geometry-based training but inference
 230 with point cloud. Methods are labelled [U], [S], [W] as unsupervised, supervised, weakly supervised.

231 **Evaluation metric:** Though we focus on the matching of point clouds, we primarily employ the
 232 widely-accepted geodesic error normalized by the square root of the total area of the mesh, to evaluate
 233 the performance of all methods.

234 **Hyper-parameters:** In Equation 12, $\lambda_{cc}, \lambda_{sc}, \lambda_m, \lambda_{geo}$ are normally set to 1, 10, 1, and 0.5 respec-
 235 tively. Model training utilizes the AdamW [40] optimizer with $\beta = (0.9, 0.99)$, learning rate of 2e-3,
 236 and batch size of 6. We provide more details of hyper-parameters in Tab. 9 in the appendix.

Table 1: Quantitative results on SCAPE_r, FAUST_r and SHREC'19_r in terms of mean geodesic errors ($\times 100$). The **best** results from the pure point cloud methods in each column are highlighted.

Method	Train Test	SCAPE_r			FAUST_r		
		SCAPE_r	FAUST_r	SHREC'19_r	FAUST_r	SCAPE_r	SHREC'19_r
3D-CODED[S] [23]	Trained on Mesh	31.0	33.0	\	2.5	31.0	\
TransMatch[S] [57]		18.6	18.3	38.8	2.7	33.6	21.0
DiffFMaps[S] [42]		12.0	12.0	17.6	3.6	19.0	16.4
NIE[W] [30]		11.0	8.7	15.6	5.5	15.0	15.1
SSMSM[W] [11]		4.1	8.5	7.3	2.4	11.0	9.0
CorrNet3D[U] [64]	Trained on PCD	58.0	63.0	\	63.0	58.0	\
SyNoRiM[S] [26]		9.5	24.6	\	7.9	21.9	\
DPC[U] [34]		17.3	11.2	28.7	11.1	17.5	31.0
SE-ORNet[U] [17]		24.6	22.8	23.6	20.3	18.9	23.0
Ours-w/o LVM[U]		13.8	10.4	14.7	8.5	16.1	15.8
Ours[U]		7.6	7.3	9.5	4.6	12.4	11.5

Table 2: Quantitative results on DT4D-H and SHREC'07-H in terms of mean geodesic errors ($\times 100$). The **best** results from the pure point cloud methods in each column are highlighted.

Method	Train Test	SCAPE_r		FAUST_r	
		DT4D-H	SHREC'07-H	DT4D-H	SHREC'07-H
TransMatch[S] [57]	Trained on Mesh	25.3	31.2	26.7	25.3
DiffFMaps[S] [42]		15.9	15.4	18.5	16.8
NIE[W] [30]		12.1	13.4	13.3	15.3
SSMSM[W] [11]		8.0	37.7	11.8	42.2
DPC[U] [34]		21.7	17.1	13.8	18.1
SE-ORNet[U] [17]	Trained on PCD	15.5	27.7	12.2	20.9
Ours-w/o LVM[U]		10.3	14.3	11.2	16.4
Ours[U]		7.7	8.9	9.0	8.9

4.1 Experimental Results

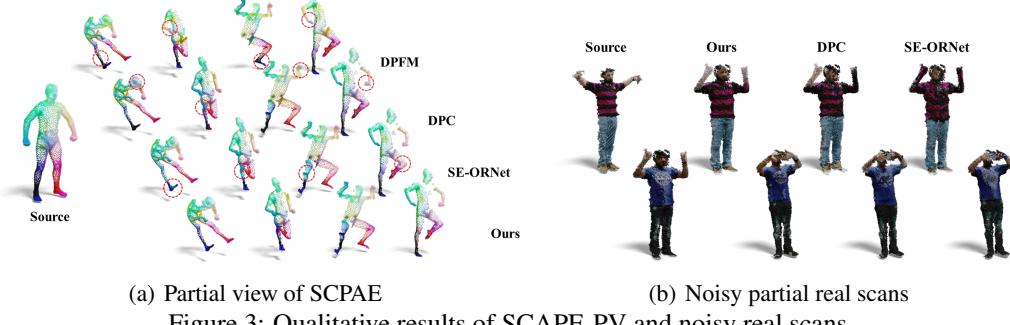
Near-isometric benchmarks: As illustrated in Tab. 1, our method consistently outperforms other pure point cloud methods in all settings. Especially, our method achieves a promising performance improvement of over **50%** compared to the previous SOTA approaches (**9.5** vs. 23.6; **11.5** vs. 23.0) in SCAPE_r/SHREC19'_r case and FAUST_r/SHREC19'_r case, i.e., many previous methods performs well in the standard seen datasets but generalizes poorly to unseen shapes. Remarkably, our method also indeed outperforms all of the baselines even the mesh-required method like SSMSM [11] (**7.3** vs. 11.2) in SCAPE_r/FAUST_r case.

Non-isometric benchmarks: We stress test our method on challenging non-isometric datasets including SHREC'07-H and DT4D-H. Our method achieves the SOTA performance for all kinds of methods as shown in Tab. 2, which indicates the excellent generalization ability to some unseen challenging cases. 1) Regarding DT4D-H, We follow the test setting of AttentiveFMaps [36], which only considers the more challenging inter-class mapping for testing of DT4D-H. The test on this non-isometric benchmark further confirms the robustness of our approach. 2) SHREC'07-H dataset comprises 20 heterogeneous human shapes with vertex numbers ranging from 3000 to 15000 and includes topological noise. Our method achieves a performance improvement of over **30%** compared to the previous SOTA approaches (**8.9** vs. 13.4; **8.9** vs. 15.3).

We attribute the above success to that the per-point features aggregated from LVMs carry rich semantic information, help to identify correspondence at the coarse level. In addition to that, our final features are further boosted by the geometric losses, leading to strong performance.

Partial matching benchmarks: As shown in Sec. 3.3, our framework can be easily adapted for unsupervised partial-to-full shape matching. We evaluate our method in two types of partial shape matching, including the challenging SHREC'16[14] Cuts and Holes benchmark and two partial-view benchmarks built on SCAPE_r and FAUST_r datasets by ourselves, where we employ raycasting from the center of each face of a regular dodecahedron to observe the shapes, resulting 12 partial view point clouds. In all cases, we match a partial point cloud with a given null (complete) point cloud.

The challenge of partial view matching arises from the presence of numerous disconnected components in the partial shapes, and the sampling of partial point clouds differs from that of complete shapes. We split the training and testing set consistent with those of SCAPE_r, FAUST_r, respectively. As illustrated in Tab. 3, our method outperforms the recent unsupervised method SSMSM [11] in 3 out of 4 test cases, which requires meshes for training. Fig. 3(a) further shows qualitative that our framework outperforms the competing methods including DPFM [6], which is based on mesh input as well.



(a) Partial view of SCAPEAE

(b) Noisy partial real scans

Figure 3: Qualitative results of SCAPE-PV and noisy real scans.

Table 3: Quantitative results on partial cases including SCAPE-PV, FAUSR-PV and SHREC’16 in terms of mean geodesic errors ($\times 100$). * indicates its original checkpoint using SURREAL190K. The **best** results from the pure point cloud methods in each column are highlighted.

Method	Train/Test	SCAPE-PV		SHREC’16-CUTS		SHREC’16-HOLES	
		SCAPE-PV	FAUST-PV	CUTS	/	HOLEs	
ConsistFMaps unsup[U] [10]	Trained on Mesh	/	/	26.6	27.0		
DPFM unsup[U] [6]		11.5	15.2	20.9	22.8		
HCLV2S*[S] [29]		8.7	5.3	/	/		
SSMSM[W] [11]		8.8	8.0	12.2	16.7		
DPC[U] [34]	Trained on PCD	13.6	14.5	32.9	32.5		
SE-ORNet[U] [17]		15.4	13.9	40.5	27.6		
Ours-w/o LVM[U]		10.1	9.2	36.4	29.3		
Ours[U]		8.4	7.6	21.2	15.3		

Table 4: Generalization performance of the checkpoint trained on sampled point cloud with fixed 1024 points of SHREC’19. We test this checkpoint on the more dense original point cloud. The **best** is highlighted.

Method	DPC[U] [34]	SE-ORNet[U] [17]	Ours[U]
SHREC’19 (1024)	5.6	5.1	4.72
SHREC’19 (Ori.)	6.1 (+8.93%)	5.9 (+15.69%)	4.73 (+0.21%)

272 Regarding purely point cloud-based baselines, we modify the shared loss of [34, 17] to adapt partiality
 273 in the same way. In the end, we achieve the SOTA compared to them, exhibiting a significant over
 274 **45%** superiority in partial view matching (**8.4** vs. 13.6; **7.6** vs. 13.9) and over **35%** superiority in
 275 cuts/holes setting (**21.2** vs. 32.9; **15.3** vs. 27.6).

276 **Generalization & robustness analysis:** Reconstruction-based methods [34, 17, 25] typically perform
 277 down-sampling to $n = 1024$ points for *both* training and testing. On the one hand, over down-
 278 sampling leads to loss of geometric details (*e.g.*, human fingers). On the other hand, point clouds
 279 scanned in reality typically consist of tens of thousands of points, which is much denser. Such
 280 resolution gap can pose great challenge for purely geometric methods. To see that, we perform
 281 generalization test on the SHREC’19 benchmark. More specifically, we use the checkpoint trained on
 282 down-sampled data released by the regarding authors to evaluate performance in both down-sampled
 283 test data (1024 points) and original test data (~ 5000 points). As shown in Tab. 4, DPC and SE-
 284 ORNet [34, 17] both experience a degradation more than **8%**. On the other hand, our method only
 285 yields a 0.21% drop and achieves the best performance in both cases. Beyond the quantitative results,
 286 we also report qualitatively the generalization performance on TOSCA benchmark following the
 287 same setting as above, see Fig. 4 for more details.

288 We attribute the above robustness to our introduction of LVMs in point feature learning, which
 289 effectively compensates for the discrepancy of low-resolution geometry. Last but not least, we
 290 also compare our method following the same training scheme and evaluation protocol as [17]. The
 291 quantitative results are reported in Tab. 7 in the appendix, which again conforms our superiority over
 292 the baselines.

293 4.2 Realworld Applications

294 In this part, we showcase the utility of our framework in two real-world applications: **Matching real**
 295 **scans:** The Panoptic dataset [32] consists of partial point clouds derived from multi-view RGB-D
 296 images. We randomly select a subset of these views to recover partial point clouds. As shown in
 297 Fig. 3(b), we transfer texture from the source shape (left-most of each row) to target via maps from
 298 ours, DPC [34], SE-ORNet [17]. Our method demonstrates smoother texture transfer compared to

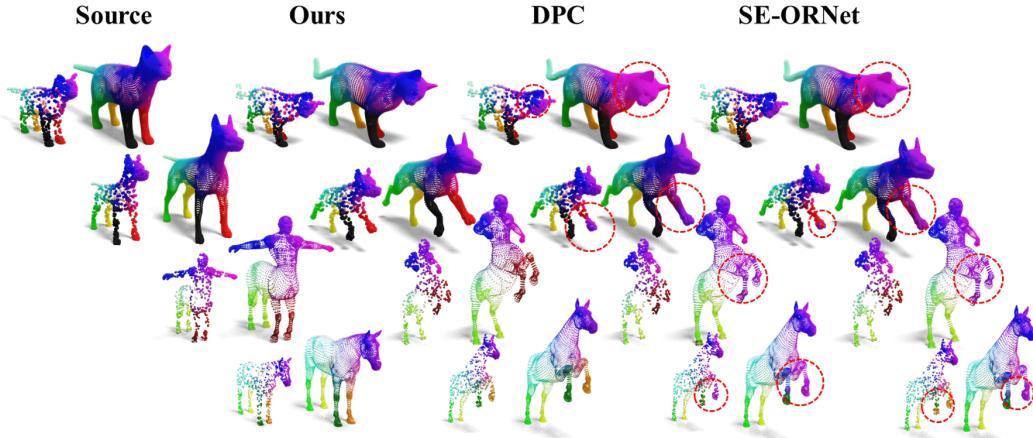


Figure 4: Qualitative results of TOSCA. Our method demonstrates enhanced generalization capabilities when transitioning from sparse point clouds in training to dense point clouds in testing.

Table 5: Statistical shape analysis on spleen medical dataset in terms of chamfer distance. The **best** is highlighted.

Model	PN-AE [2]	DG-AE [59]	CPAE [13]	ISR [12]	DPC [34]	Point2SSM [3]	Ours
CD (mm)	43.7	43.5	61.3	17.6	10.6	3.4	2.9

Table 6: Mean geodesic errors ($\times 100$) on different ablated settings, the models are all trained on SCAPE_r and test on SCAPE_r.

w/o pointwise proj	w/o LG-NET	w/o Geo. Loss	w/o Featup	w/o PE	w/o LA-NET	w/o GA-NET	w/o Fusion	Full
33.1	19.6	9.1	9.0	8.6	9.8	9.8	9.7	7.6

299 baselines (see particularly facial details and strips in the T-shirt); **Statistical shape models(SSM) for**
300 **medical data:** Following Point2SSM [3], we test our method on the anatomical SSM tasks. We stick
301 to the regarding experimental setting and report our score in the spleen subset. As shown in Tab. 5,
302 our method outperforms the second best by a 14% relative error reduction.

303 4.3 Ablation Study

304 We first justify our overall design in Tab. 6, where we sequentially remove each building block
305 from our pipeline and train/test model on SCAPE_r. The performance gaps well support our claims.
306 Beyond that, in Tab. 1, Tab. 2 and Tab. 3, we report experimental results [w/o LVM] to demonstrate
307 that the coarse semantic representations extracted by LVMs play a crucial role throughout the pipeline,
308 whether it is in the full or partial settings. Finally, we highlight that we have also performed robustness
309 evaluation regarding noisy data and rotation perturbations in Sec. B.3.

310 5 Conclusion and Limitation

311 In this paper, we address the challenge of unsupervised non-rigid point cloud matching. In conclusion,
312 we proposed a novel learning-based framework for non-rigid point cloud matching that can be trained
313 purely on point clouds without correspondence annotations and extends naturally to partial-to-full
314 matching. By incorporating semantic features from large vision models (LVMs) into geometry-based
315 shape feature learning, our framework resolves ambiguities from self-similarities and demonstrates
316 strong generalizability and robustness. Our method achieves state-of-the-art results in matching
317 non-rigid point clouds, even in challenging scenarios with partial and noisy data.

318 **Limitation & Future Work** The primary limitation of our method is its assumption of roughly
319 aligned input point clouds. In the future, we plan to further explore the ability of LVM to address this
320 limitation.

321 **References**

- 322 [1] Ahmed Abdelreheem, Abdelrahman Eldekokey, Maks Ovsjanikov, and Peter Wonka. Zero-shot 3d shape
323 correspondence. In *SIGGRAPH Asia 2023 Conference Papers*, pages 1–11, 2023.
- 324 [2] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas J Guibas. Learning representations
325 and generative models for 3d point clouds. *Proceedings of the 35th International Conference on Machine
326 Learning*, 2018.
- 327 [3] Jadie Adams and Shireen Elhabian. Point2ssm: Learning morphological variations of anatomies from
328 point cloud. *arXiv preprint arXiv:2305.14486*, 2023.
- 329 [4] Brian Amberg, Sami Romdhani, and Thomas Vetter. Optimal step nonrigid icp algorithms for surface
330 registration. 2007.
- 331 [5] Dragomir Anguelov, Praveen Srinivasan, Daphne Koller, Sebastian Thrun, Jim Rodgers, and James Davis.
332 SCAPE: Shape Completion and Animation of People. 2005.
- 333 [6] Souhaib Attaiki, Gautam Pai, and Maks Ovsjanikov. Dpfm: Deep partial functional maps. In *2021
334 International Conference on 3D Vision (3DV)*, pages 175–185. IEEE, 2021.
- 335 [7] Federica Bogo, Javier Romero, Matthew Loper, and Michael J Black. Faust: Dataset and evaluation for 3d
336 mesh registration. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
337 pages 3794–3801, 2014.
- 338 [8] Aljaz Bozic, Pablo Palafox, Michael Zollhöfer, Angela Dai, Justus Thies, and Matthias Nießner. Neural
339 non-rigid tracking. In *NeurIPS*, volume 33, pages 18727–18737, 2020.
- 340 [9] Aljaz Bozic, Michael Zollhofer, Christian Theobalt, and Matthias Nießner. Deepdeform: Learning
341 non-rigid rgb-d reconstruction with semi-supervised data. In *CVPR*, pages 7002–7012, 2020.
- 342 [10] Dongliang Cao and Florian Bernard. Unsupervised deep multi-shape matching. In *ECCV*, 2022.
- 343 [11] Dongliang Cao and Florian Bernard. Self-supervised learning for multimodal non-rigid shape matching.
344 In *CVPR*, 2023.
- 345 [12] Nenglun Chen, Lingjie Liu, Zhiming Cui, Runnan Chen, Duygu Ceylan, Changhe Tu, and Wenping
346 Wang. Unsupervised learning of intrinsic structural representation points. In *Proceedings of the IEEE/CVF
347 conference on computer vision and pattern recognition*, pages 9121–9130, 2020.
- 348 [13] An-Chieh Cheng, Xuetong Li, Min Sun, Ming-Hsuan Yang, and Sifei Liu. Learning 3d dense correspon-
349 dence via canonical point autoencoder. *Advances in Neural Information Processing Systems*, 34:6608–6620,
350 2021.
- 351 [14] Luca Cosmo, Emanuele Rodola, Michael M Bronstein, Andrea Torsello, Daniel Cremers, Y Sahillioğlu,
352 et al. Shrec’16: Partial matching of deformable shapes. In *Eurographics Workshop on 3D Object Retrieval,
353 EG 3DOR*, pages 61–67. Eurographics Association, 2016.
- 354 [15] Thomas Cover and Peter Hart. Nearest neighbor pattern classification. *IEEE transactions on information
355 theory*, 13(1):21–27, 1967.
- 356 [16] Keenan Crane, Clémence Weischedel, and Max Wardetzky. The heat method for distance computation.
357 *Commun. ACM*, 60(11):90–99, Oct. 2017.
- 358 [17] Jiacheng Deng, Chuxin Wang, Jiahao Lu, Jianfeng He, Tianzhu Zhang, Jiyang Yu, and Zhe Zhang. Se-
359 ornet: Self-ensembling orientation-aware network for unsupervised point cloud shape correspondence. In
360 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5364–5373,
361 2023.
- 362 [18] Theo Deprelle, Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry.
363 Learning elementary structures for 3d shape generation and matching. *arXiv preprint arXiv:1908.04725*,
364 2019.
- 365 [19] Nicolas Donati, Etienne Cormann, and Maks Ovsjanikov. Deep orientation-aware functional maps: Tackling
366 symmetry issues in shape matching. In *CVPR*, pages 742–751, 2022.
- 367 [20] Niladri Shekhar Dutt, Sanjeev Muralikrishnan, and Niloy J Mitra. Diffusion 3d features (diff3f): Decorating
368 untextured shapes with distilled semantic features. *arXiv preprint arXiv:2311.17024*, 2023.
- 369 [21] Stephanie Fu, Mark Hamilton, Laura E. Brandt, Axel Feldmann, Zhoutong Zhang, and William T. Freeman.
370 Featup: A model-agnostic framework for features at any resolution. In *The Twelfth International Conference
371 on Learning Representations*, 2024.
- 372 [22] Daniela Giorgi, Silvia Biasotti, and Laura Paraboschi. Shape retrieval contest 2007: Watertight models
373 track. *SHREC competition*, 8(7):7, 2007.
- 374 [23] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. 3d-coded: 3d
375 correspondences by deep deformation. In *ECCV*, 2018.
- 376 [24] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min Hu. Pct:
377 Point cloud transformer. *Computational Visual Media*, 7:187–199, 2021.
- 378 [25] Jianfeng He, Jiacheng Deng, Tianzhu Zhang, Zhe Zhang, and Yongdong Zhang. Hierarchical shape-
379 consistent transformer for unsupervised point cloud shape correspondence. *IEEE Transactions on Image
380 Processing*, 2023.

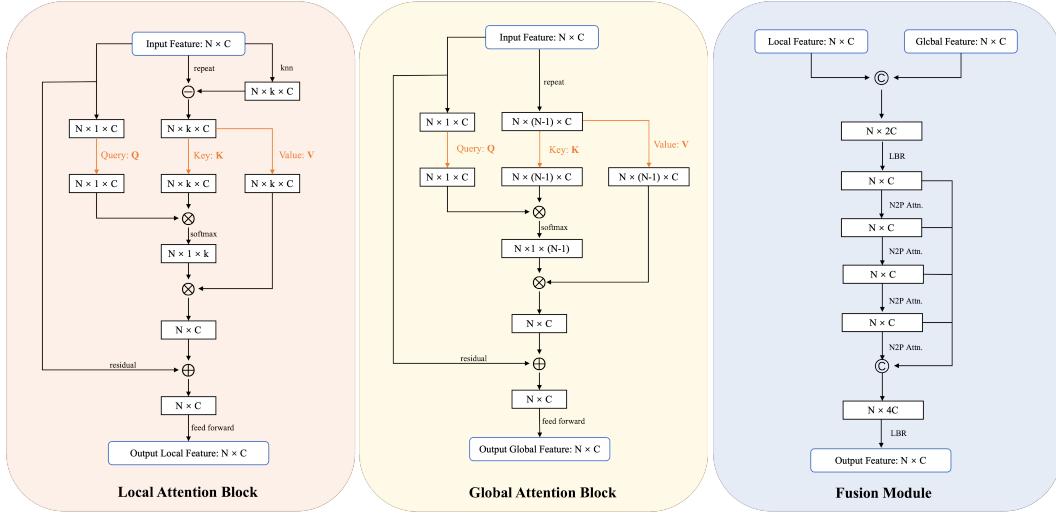
- 381 [26] Jiahui Huang, Tolga Birdal, Zan Gojcic, Leonidas J. Guibas, and Shi-Min Hu. Multiway Non-rigid Point
 382 Cloud Registration via Learned Functional Map Synchronization. *IEEE Transactions on Pattern Analysis*
 383 and *Machine Intelligence*, pages 1–1, 2022.
- 384 [27] Ruqi Huang and Maks Ovsjanikov. Adjoint map representation for shape analysis and matching. In
 385 *Computer Graphics Forum*, volume 36, pages 151–163. Wiley Online Library, 2017.
- 386 [28] Ruqi Huang, Jing Ren, Peter Wonka, and Maks Ovsjanikov. Consistent zoomout: Efficient spectral map
 387 synchronization. In *Computer Graphics Forum*, volume 39, pages 265–278. Wiley Online Library, 2020.
- 388 [29] Xiangru Huang, Haitao Yang, Etienne Vouga, and Qixing Huang. Dense correspondences between human
 389 bodies via learning transformation synchronization on graphs. In *NeurIPS*, 2020.
- 390 [30] Puhua Jiang, Mingze Sun, and Ruqi Huang. Neural intrinsic embedding for non-rigid point matching. In
 391 *CVPR*, 2023.
- 392 [31] Puhua Jiang, Mingze Sun, and Ruqi Huang. Non-rigid shape registration via deep functional maps prior.
 393 In *NeurIPS*, 2023.
- 394 [32] Hanbyul Joo, Hao Liu, Lei Tan, Lin Gui, Bart Nabbe, Iain Matthews, Takeo Kanade, Shohei Nobuhara, and
 395 Yaser Sheikh. Panoptic studio: A massively multiview system for social motion capture. In *Proceedings of*
 396 *the IEEE International Conference on Computer Vision*, pages 3334–3342, 2015.
- 397 [33] Hanbyul Joo, Tomas Simon, Xulong Li, Hao Liu, Lei Tan, Lin Gui, Sean Banerjee, Timothy Scott Godisart,
 398 Bart Nabbe, Iain Matthews, Takeo Kanade, Shohei Nobuhara, and Yaser Sheikh. Panoptic studio: A
 399 massively multiview system for social interaction capture. *IEEE Transactions on Pattern Analysis and*
 400 *Machine Intelligence*, 2017.
- 401 [34] Itai Lang, Dvir Ginzburg, Shai Avidan, and Dan Raviv. Dpc: Unsupervised deep point correspondence
 402 via cross and self construction. In *2021 International Conference on 3D Vision (3DV)*, pages 1442–1451.
 403 IEEE, 2021.
- 404 [35] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training
 405 with frozen image encoders and large language models. In *International conference on machine learning*,
 406 pages 19730–19742. PMLR, 2023.
- 407 [36] Lei Li, Nicolas Donati, and Maks Ovsjanikov. Learning multi-resolution functional maps with spectral
 408 attention for robust shape matching. In *NeurIPS*, 2022.
- 409 [37] Yang Li and Tatsuya Harada. Lepard: Learning partial point cloud matching in rigid and deformable
 410 scenes. In *CVPR*, 2022.
- 411 [38] Yang Li and Tatsuya Harada. Non-rigid point cloud registration with neural deformation pyramid. *Advances*
 412 in *Neural Information Processing Systems*, 35:27757–27768, 2022.
- 413 [39] Or Litany, Tal Remez, Emanuele Rodolà, Alexander M. Bronstein, and Michael M. Bronstein. Deep
 414 functional maps: Structured prediction for dense shape correspondence. In *ICCV*, 2017.
- 415 [40] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
 416 *arXiv:1711.05101*, 2017.
- 417 [41] Robin Magnet, Jing Ren, Olga Sorkine-Hornung, and Maks Ovsjanikov. Smooth non-rigid shape matching
 418 via effective dirichlet energy optimization. In *2022 International Conference on 3D Vision (3DV)*, pages
 419 495–504. IEEE, 2022.
- 420 [42] Riccardo Marin, Marie-Julie Rakotosaona, Simone Melzi, and Maks Ovsjanikov. Correspondence learning
 421 via linearly-invariant embedding. *Advances in Neural Information Processing Systems*, 33:1608–1620,
 422 2020.
- 423 [43] Simone Melzi, Riccardo Marin, Emanuele Rodolà, Umberto Castellani, Jing Ren, Adrien Poulenard, et al.
 424 Shrec’19: matching humans with different connectivity. In *Eurographics Workshop on 3D Object Retrieval*.
 425 The Eurographics Association, 2019.
- 426 [44] Simone Melzi, Jing Ren, Emanuele Rodolà, Peter Wonka, and Maks Ovsjanikov. Zoomout: Spectral
 427 upsampling for efficient shape correspondence. *Proc. SIGGRAPH Asia*, 2019.
- 428 [45] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren
 429 Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*,
 430 65(1):99–106, 2021.
- 431 [46] Luca Morreale, Noam Aigerman, Vladimir G. Kim, and Niloy J. Mitra. Semantic neural surface maps. In
 432 *Eurographics*, 2024.
- 433 [47] Dorian Nogneng and Maks Ovsjanikov. Informative descriptor preservation via commutativity for shape
 434 matching. *Computer Graphics Forum*, 36(2):259–267, 2017.
- 435 [48] OpenAI. Gpt-3.5 language model. <https://www.openai.com/research/gpt-3>, 2021. Accessed:
 436 May 21, 2023.
- 437 [49] Maxime Oquab, Timothée Darzet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre
 438 Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu
 439 Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel
 440 Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski.
 441 Dinov2: Learning robust visual features without supervision, 2023.

- 442 [50] Maks Ovsjanikov, Mirela Ben-Chen, Justin Solomon, Adrian Butscher, and Leonidas Guibas. Functional
 443 Maps: A Flexible Representation of Maps Between Shapes. *ACM Transactions on Graphics (TOG)*,
 444 31(4):30, 2012.
- 445 [51] Gianluca Paravati, Fabrizio Lamberti, Valentina Gatteschi, Claudio Demartini, and Paolo Montuschi. Point
 446 cloud-based automatic assessment of 3d computer animation courseworks. *IEEE Transactions on Learning
 447 Technologies*, 10(4):532–543, 2016.
- 448 [52] Jing Ren, Adrien Poulenard, Peter Wonka, and Maks Ovsjanikov. Continuous and orientation-preserving
 449 correspondences via functional maps. *ACM Trans. Graph.*, 37(6):248:1–248:16, Dec. 2018.
- 450 [53] Emanuele Rodolà, Luca Cosmo, Michael M Bronstein, Andrea Torsello, and Daniel Cremers. Partial
 451 Functional Correspondence. In *Computer Graphics Forum*, 2016.
- 452 [54] Nicholas Sharp, Souhaib Attaiki, Keenan Crane, and Maks Ovsjanikov. Diffusionnet: Discretization
 453 agnostic learning on surfaces. *ACM Transactions on Graphics*, 2022.
- 454 [55] Nicholas Sharp and Keenan Crane. A laplacian for nonmanifold triangle meshes. *Computer Graphics
 455 Forum*, 2020.
- 456 [56] Mingze Sun, Shiwei Mao, Puhua Jiang, Maks Ovsjanikov, and Ruqi Huang. Spatially and spectrally
 457 consistent deep functional maps. In *ICCV*, 2023.
- 458 [57] Giovanni Trappolini, Luca Cosmo, Luca Moschella, Riccardo Marin, Simone Melzi, and Emanuele Rodolà.
 459 Shape registration in the time of transformers. *Advances in Neural Information Processing Systems*,
 460 34:5731–5744, 2021.
- 461 [58] Jinbao Wang, Shujie Tan, Xiantong Zhen, Shuo Xu, Feng Zheng, Zhenyu He, and Ling Shao. Deep 3d
 462 human pose estimation: A review. *Computer Vision and Image Understanding*, 2021.
- 463 [59] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E. Sarma, Michael M. Bronstein, and Justin M. Solomon.
 464 Dynamic graph cnn for learning on point clouds. *ACM Transactions on Graphics (TOG)*, 2019.
- 465 [60] Thomas Wimmer, Peter Wonka, and Maks Ovsjanikov. Back to 3d: Few-shot 3d keypoint detection with
 466 back-projected 2d features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
 467 Recognition*, 2024.
- 468 [61] Chengzhi Wu, Junwei Zheng, Julius Pfrommer, and Jürgen Beyerer. Attention-based point cloud edge
 469 sampling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
 470 pages 5333–5343, 2023.
- 471 [62] Yuxin Yao, Bailin Deng, Weiwei Xu, and Juyong Zhang. Fast and robust non-rigid registration using
 472 accelerated majorization-minimization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
 473 2023.
- 474 [63] Tao Yu, Zerong Zheng, Kaiwen Guo, Jianhui Zhao, Qionghai Dai, Hao Li, Gerard Pons-Moll, and Yebin
 475 Liu. Doublefusion: Real-time capture of human performances with inner body shapes from a single depth
 476 sensor. In *CVPR*, 2018.
- 477 [64] Yiming Zeng, Yue Qian, Zhiyu Zhu, Junhui Hou, Hui Yuan, and Ying He. Corinet3d: Unsupervised
 478 end-to-end learning of dense correspondence for 3d point clouds. In *CVPR*, pages 6052–6061, 2021.
- 479 [65] Renrui Zhang, Liuhui Wang, Yu Qiao, Peng Gao, and Hongsheng Li. Learning 3d representations from 2d
 480 pre-trained models via image-to-point masked autoencoders. In *Proceedings of the IEEE/CVF Conference
 481 on Computer Vision and Pattern Recognition*, pages 21769–21780, 2023.
- 482 [66] Silvia Zuffi, Angjoo Kanazawa, David W Jacobs, and Michael J Black. 3d menagerie: Modeling the
 483 3d shape and pose of animals. In *Proceedings of the IEEE conference on computer vision and pattern
 484 recognition*, pages 6365–6373, 2017.

485 In this appendix, we provide more technical details and experimental results, including 1) A detailed
 486 description of the building blocks of our LG-Net in Sec. A; 2) Further qualitative results on matching
 487 heterogeneous shapes from SHREC'07 and DT4D-H in Sec. B.1; 3) Quantitative results following
 488 the setting from [34, 17, 25], where train/test with the sparse point clouds of fixed 1024 points in
 489 Sec. B.2; 4) Robustness evaluation of our method with respect to several perturbations in Sec. B.3; 5)
 490 More high-dimensional feature visualization and matching results of medical dataset in Sec. B.4; 6)
 491 Run-time analysis, hyper-parameter instruction in Sec. B.5 and Sec. B.6 respectively. Finally, the
 492 broader impacts are discussed in Sec. C.

493 A Technical Details

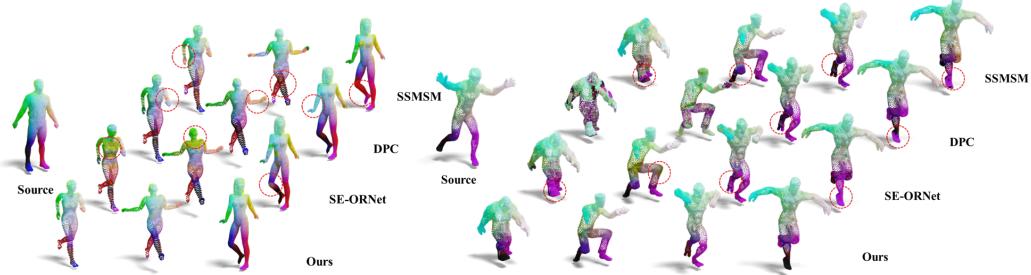
494 Fig. 5 depicts from left to right the architecture diagrams of our local attention, global attention, and
 495 fusion module.



495 Figure 5: The schematic illustration of the main blocks of our network.

496 B Experiments

497 B.1 Further Qualitative Results



498 Figure 6: We estimate correspondences between heterogeneous shapes from SHREC'07 and DT4D-H
 499 with DPC, SE-ORNET and one SSMSM, all trained on the SCAPE_r dataset. Our method outperforms
 500 the competing methods by a large margin.

498 In Fig. 6, we qualitatively visualize maps obtained by different methods tested in the SHREC'07
 499 and DT4D-H benchmark. It is obvious that our results outperform all the competing the competing
 500 methods, showing superior generalization performance.

Table 7: Quantitative results on human and animals datasets. Acc signifies correspondence accuracy at 0.01 error tolerance, and err denotes average correspondence error. The **best** results in each column are highlighted.

Train Test	SHREC'19		SURREAL		TOSCA		SMAL	
	SHREC'19		SURREAL		TOSCA		SMAL	
	acc ↑	err ↓						
3D-CODED[S] [23]	/	/	2.1%	8.1	/	/	0.5%	19.2
Elementary[S] [18]	/	/	0.5%	13.7	/	/	2.3%	7.6
CorrNet3D[U] [64]	0.4%	33.8	6.0%	6.9	0.3%	32.7	5.3%	9.8
DPC[U] [34]	15.3%	5.6	17.7%	6.1	34.7%	2.8	33.2%	5.8
SE-ORNet[U] [17]	17.5%	5.1	21.5%	4.6	38.3%	2.7	36.4%	3.9
HSTR[U] [25]	19.3%	4.9	19.4%	5.6	52.3%	1.2	33.9%	5.6
Ours [U]	20.4%	4.7	23.4%	4.4	43.7%	1.0	37.9%	3.6

Table 8: Mean geodesic errors ($\times 100$) on under different perturbations. Noisy PC means the input point clouds are perturbed by Gaussian noise. Rotated PC means the input point clouds are randomly rotated within ± 30 degrees. The standard deviation value is shown in parentheses.

Method	Unperturbed	Noisy PC	Rotated PC
DiffFMaps[S] [42]	Trained on Mesh	12.0	14.9(2.57)
NIE[W] [30]		11.0	11.5(0.32)
SSMSM[W] [11]		4.1	5.4(0.11)
DPC[U] [34]	Trained on PCD	17.3	18.2(0.80)
SE-ORNet[U] [17]		24.6	24.7(0.15)
Ours [U]		7.6	7.8(0.10)
			8.7(0.60)

501 B.2 Further Quantitative Results

502 **Sparse Humans/Animals Benchmarks:** Following the prior works [34, 17, 25], we conduct the
503 experiments with a consistent sampling point number of $n = 1024$. In addition to the datasets
504 mentioned in the Sec. 4, two more datasets are included for training purposes including SURREAL
505 and SMAL. **SURREAL** is the large-scale dataset from [23] comprises 230,000 training shapes, from
506 which we select the first 2,000 shapes and use them solely for training. **SMAL** is from [66], which
507 includes parameterized animal models for generating shapes. We employ the model to generate 2000
508 instances of diverse poses for each animal category, resulting in a training dataset comprising 10000
509 shapes.

510 Specifically, we train on the SURREAL and SHREC'19 dataset respectively, and then test on the
511 SHREC'19 dataset. Similarly, we train respectively on SMAL and TOSCA dataset, and then test on
512 the TOSCA dataset. As shown in Tab. 7, unlike HSTR[25], which achieves the best performance
513 on its intra-dataset but lags behind SE-ORNet[17] on cross-dataset generalization, our approach
514 excels in both intra-dataset and cross-dataset tests, surpassing all existing methods by over **4%**. This
515 also complements Tab. 1 and Tab. 2, demonstrating that our method yields robust results whether
516 trained/tested on dense or sparse point clouds.

517 B.3 Robustness

518 Moreover, we evaluate the robustness of our model with respect to noise and rotation perturbation
519 and report in Tab. 8. More specifically, we perturb the point clouds by: 1) Adding per-point Gaussian
520 noise with i.i.d $\mathcal{N}(0, 0.02)$ along the normal direction on each point; 2) Randomly rotating ± 30
521 degree along some randomly sampled direction. We perform 3 rounds of test, and report both mean
522 error and the standard deviation in parentheses. Note that SE-ORNET[17] is designed for rotational
523 robustness, which enjoys better rotation performance. Apart from that, our pipeline delivers the most
524 robust performance among the baselines.

525 B.4 More Visualizations

526 **High-dimensional feature visualization:** To further validate the characteristics of the representations
527 learned by our method, we present a set of more comprehensive visualizations of the features. As
528 shown in Fig. 7, our feature distribution is more clean and localized. However, upon losing geometric
529 or semantic information, the features across different dimensions become divergent, resulting in the
530 loss of regular fine-grained representation at various levels.

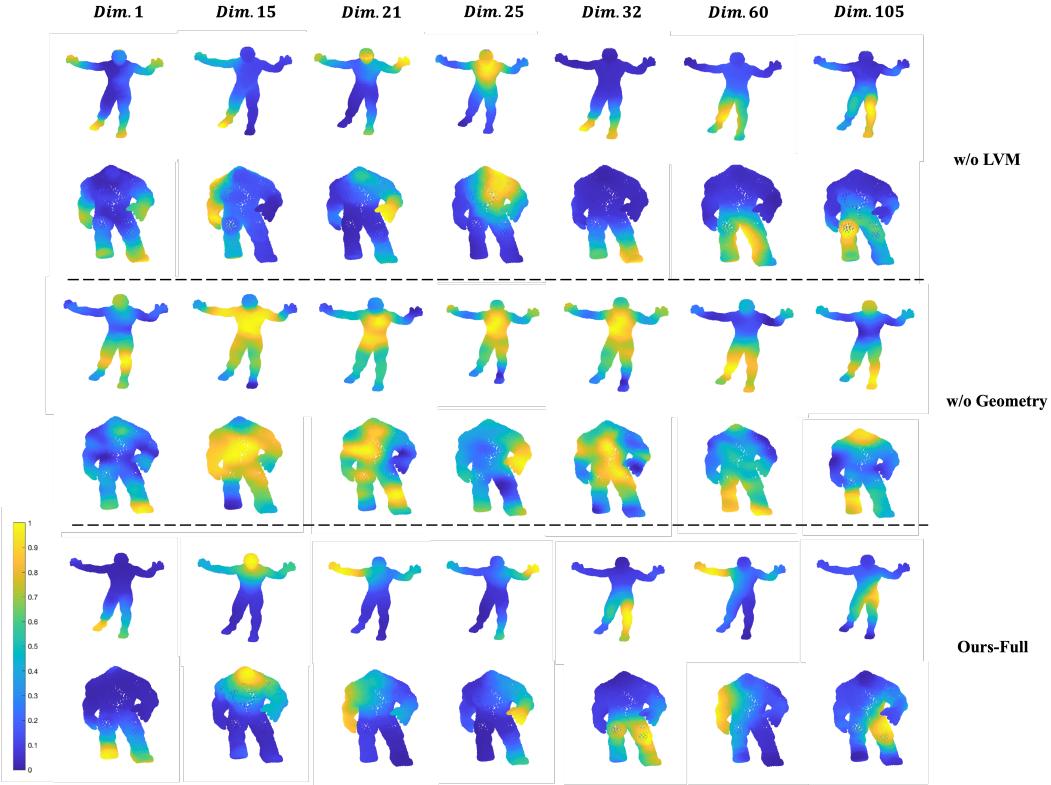
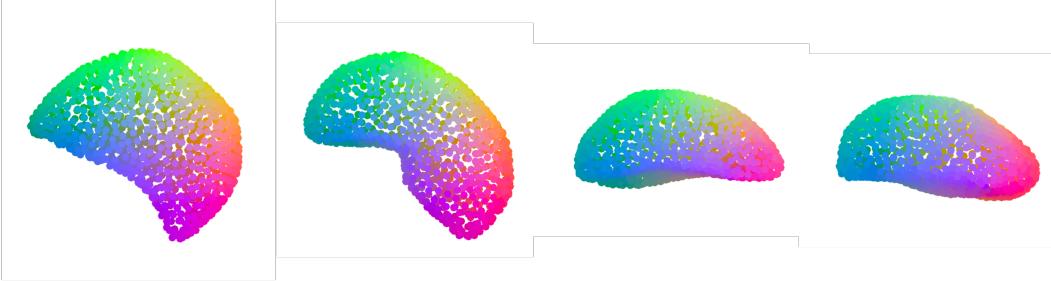


Figure 7: Visualization of different feature dimensions. **Dim.i** denotes the features of the $i - th$ dimension, where $i \leq 128$.

531 **Matching results of medical dataset:** To supplement Tab.5, we further visualize the matching results
 532 on the spleen organ in Fig.8, where excellent mapping is achieved regardless of whether the spleen
 exhibits various shapes or is positioned at different angles.



533 Figure 8: Matching result of the spleen dataset from [3].

534 **More qualitative results:** We further visualize the results of TOSCA, DT4D, and SCAPE-PV, which
 535 respectively serve as qualitative validation supplements for learning sparse point clouds in Tab. 7, the
 536 generalization capability in Tab. 2, and the adaptability to partial shapes in Tab. 3. The training and
 537 testing procedures align with the methods described in the aforementioned table, with quantitative
 538 supplements presented respectively in Fig. 4, Fig. 10 and Fig. 11, respectively.

539 **B.5 Running Time**

540 We perform all the experiments on a machine with NVIDIA A100-SMX4 80GB and Intel(R) Xeon(R)
 541 CPU E5-2680 v4 @ 2.40GHz using the PyTorch 2.2.0 framework. Benchmarking on SCAPE_r
 542 dataset, our method achieves an average processing time of approximately **0.21** seconds per pair on
 543 testing, and approximately **2.35** seconds per iteration (6 batches) on training. In addition, it is also

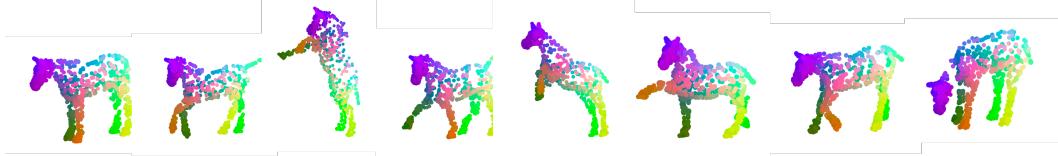


Figure 9: More qualitative results of TOSCA. All horse shapes from the dataset have been showcased.



Figure 10: More qualitative results of DT4D. Our method demonstrates a notable improvement over other baselines.

Table 9: Hyper-parameters. The tables details the hyperparameter values that we used for the training of our method.

symbol	Description	Value
k	The nearest number for computing geometrically similarity loss	500
k_{attn}	The number for searching latent nearest features in local attention	40
k_c	The number for self-construction/cross-construction neighborhood size	40
k_m	Mapping loss neighborhood size	40
k'_{attn}	The number for local attention when training point cloud with fixed $n = 1024$ points	10
k'_m	The mapping loss neighborhood size when training point cloud with fixed $n = 1024$ points	10
k'_c	The number for construction neighborhood size when training point cloud with fixed $n = 1024$ points	40
α	Mapping loss neighbor sensitivity	8
TEs	Training epochs	50
H, W	The size of our projected image	224,224

544 feasible to train on a single NVIDIA GeForce RTX 3090 24GB, only necessitating a reduced batch
 545 size.

546 B.6 Additional Hyper-parameter Details

547 For a comprehensive understanding of the specific hyper-parameter configurations, please refer to
 548 Tab. 9.



Figure 11: Qualitative results of SCPAE-PV.

549 C Broader Impacts

550 We fail to see any immediate ethical issue with the proposed method. On the other hand, since
 551 our method is extensively evaluated in matching human shapes and achieves excellent results, one
 552 potential misuse can be surveillance, which may pose negative societal impact.

553 **NeurIPS Paper Checklist**

554 **1. Claims**

555 Question: Do the main claims made in the abstract and introduction accurately reflect the
556 paper's contributions and scope?

557 Answer: [Yes]

558 Justification: We have made the main claims in the abstract and introduction.

559 Guidelines:

- 560 • The answer NA means that the abstract and introduction do not include the claims
561 made in the paper.
- 562 • The abstract and/or introduction should clearly state the claims made, including the
563 contributions made in the paper and important assumptions and limitations. A No or
564 NA answer to this question will not be perceived well by the reviewers.
- 565 • The claims made should match theoretical and experimental results, and reflect how
566 much the results can be expected to generalize to other settings.
- 567 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
568 are not attained by the paper.

569 **2. Limitations**

570 Question: Does the paper discuss the limitations of the work performed by the authors?

571 Answer: [Yes]

572 Justification: We have discussed the limitations of the work in the end of paper.

573 Guidelines:

- 574 • The answer NA means that the paper has no limitation while the answer No means that
575 the paper has limitations, but those are not discussed in the paper.
- 576 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 577 • The paper should point out any strong assumptions and how robust the results are to
578 violations of these assumptions (e.g., independence assumptions, noiseless settings,
579 model well-specification, asymptotic approximations only holding locally). The authors
580 should reflect on how these assumptions might be violated in practice and what the
581 implications would be.
- 582 • The authors should reflect on the scope of the claims made, e.g., if the approach was
583 only tested on a few datasets or with a few runs. In general, empirical results often
584 depend on implicit assumptions, which should be articulated.
- 585 • The authors should reflect on the factors that influence the performance of the approach.
586 For example, a facial recognition algorithm may perform poorly when image resolution
587 is low or images are taken in low lighting. Or a speech-to-text system might not be
588 used reliably to provide closed captions for online lectures because it fails to handle
589 technical jargon.
- 590 • The authors should discuss the computational efficiency of the proposed algorithms
591 and how they scale with dataset size.
- 592 • If applicable, the authors should discuss possible limitations of their approach to
593 address problems of privacy and fairness.
- 594 • While the authors might fear that complete honesty about limitations might be used by
595 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
596 limitations that aren't acknowledged in the paper. The authors should use their best
597 judgment and recognize that individual actions in favor of transparency play an impor-
598 tant role in developing norms that preserve the integrity of the community. Reviewers
599 will be specifically instructed to not penalize honesty concerning limitations.

600 **3. Theory Assumptions and Proofs**

601 Question: For each theoretical result, does the paper provide the full set of assumptions and
602 a complete (and correct) proof?

603 Answer: [NA]

604 Justification: Our paper mainly focuses on methods and their applications which does not
605 include theoretical results.

606 Guidelines:

- 607 • The answer NA means that the paper does not include theoretical results.
- 608 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
609 referenced.
- 610 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 611 • The proofs can either appear in the main paper or the supplemental material, but if
612 they appear in the supplemental material, the authors are encouraged to provide a short
613 proof sketch to provide intuition.
- 614 • Inversely, any informal proof provided in the core of the paper should be complemented
615 by formal proofs provided in appendix or supplemental material.
- 616 • Theorems and Lemmas that the proof relies upon should be properly referenced.

617 4. Experimental Result Reproducibility

618 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
619 perimental results of the paper to the extent that it affects the main claims and/or conclusions
620 of the paper (regardless of whether the code and data are provided or not)?

621 Answer: [Yes]

622 Justification: We have provided the Implementation details in Section 4.

623 Guidelines:

- 624 • The answer NA means that the paper does not include experiments.
- 625 • If the paper includes experiments, a No answer to this question will not be perceived
626 well by the reviewers: Making the paper reproducible is important, regardless of
627 whether the code and data are provided or not.
- 628 • If the contribution is a dataset and/or model, the authors should describe the steps taken
629 to make their results reproducible or verifiable.
- 630 • Depending on the contribution, reproducibility can be accomplished in various ways.
631 For example, if the contribution is a novel architecture, describing the architecture fully
632 might suffice, or if the contribution is a specific model and empirical evaluation, it may
633 be necessary to either make it possible for others to replicate the model with the same
634 dataset, or provide access to the model. In general, releasing code and data is often
635 one good way to accomplish this, but reproducibility can also be provided via detailed
636 instructions for how to replicate the results, access to a hosted model (e.g., in the case
637 of a large language model), releasing of a model checkpoint, or other means that are
638 appropriate to the research performed.
- 639 • While NeurIPS does not require releasing code, the conference does require all submis-
640 sions to provide some reasonable avenue for reproducibility, which may depend on the
641 nature of the contribution. For example
 - 642 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
643 to reproduce that algorithm.
 - 644 (b) If the contribution is primarily a new model architecture, the paper should describe
645 the architecture clearly and fully.
 - 646 (c) If the contribution is a new model (e.g., a large language model), then there should
647 either be a way to access this model for reproducing the results or a way to reproduce
648 the model (e.g., with an open-source dataset or instructions for how to construct
649 the dataset).
 - 650 (d) We recognize that reproducibility may be tricky in some cases, in which case
651 authors are welcome to describe the particular way they provide for reproducibility.
652 In the case of closed-source models, it may be that access to the model is limited in
653 some way (e.g., to registered users), but it should be possible for other researchers
654 to have some path to reproducing or verifying the results.

655 5. Open access to data and code

656 Question: Does the paper provide open access to the data and code, with sufficient instruc-
657 tions to faithfully reproduce the main experimental results, as described in supplemental
658 material?

659 Answer: [Yes]

660 Justification: We will release the code after acceptance.

661 Guidelines:

- 662 • The answer NA means that paper does not include experiments requiring code.
- 663 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 664 • While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- 665 • The instructions should contain the exact command and environment needed to run to 666 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 667 • The authors should provide instructions on data access and preparation, including how 668 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 669 • The authors should provide scripts to reproduce all experimental results for the new 670 proposed method and baselines. If only a subset of experiments are reproducible, they 671 should state which ones are omitted from the script and why.
- 672 • At submission time, to preserve anonymity, the authors should release anonymized 673 versions (if applicable).
- 674 • Providing as much information as possible in supplemental material (appended to the 675 paper) is recommended, but including URLs to data and code is permitted.
- 676

677 6. Experimental Setting/Details

678 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
679 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
680 results?

681 Answer: [Yes]

682 Justification: We have provided the Implementation details in Section 4.

683 Guidelines:

- 684 • The answer NA means that the paper does not include experiments.
- 685 • The experimental setting should be presented in the core of the paper to a level of detail
686 that is necessary to appreciate the results and make sense of them.
- 687 • The full details can be provided either with the code, in appendix, or as supplemental
688 material.

689 7. Experiment Statistical Significance

690 Question: Does the paper report error bars suitably and correctly defined or other appropriate
691 information about the statistical significance of the experiments?

692 Answer: [Yes]

693 Justification: We have reported the robustness of our method regarding rotation perturbation
694 in Tab. 8.

695 Guidelines:

- 696 • The answer NA means that the paper does not include experiments.
- 697 • The authors should answer “Yes” if the results are accompanied by error bars, confi-
698 dence intervals, or statistical significance tests, at least for the experiments that support
699 the main claims of the paper.
- 700 • The factors of variability that the error bars are capturing should be clearly stated (for
701 example, train/test split, initialization, random drawing of some parameter, or overall
702 run with given experimental conditions).
- 703 • The method for calculating the error bars should be explained (closed form formula,
704 call to a library function, bootstrap, etc.)
- 705 • The assumptions made should be given (e.g., Normally distributed errors).
- 706

- 710 • It should be clear whether the error bar is the standard deviation or the standard error
711 of the mean.
712 • It is OK to report 1-sigma error bars, but one should state it. The authors should
713 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
714 of Normality of errors is not verified.
715 • For asymmetric distributions, the authors should be careful not to show in tables or
716 figures symmetric error bars that would yield results that are out of range (e.g. negative
717 error rates).
718 • If error bars are reported in tables or plots, The authors should explain in the text how
719 they were calculated and reference the corresponding figures or tables in the text.

720 **8. Experiments Compute Resources**

721 Question: For each experiment, does the paper provide sufficient information on the com-
722 puter resources (type of compute workers, memory, time of execution) needed to reproduce
723 the experiments?

724 Answer: [Yes]

725 Justification: We have reported computational resources in Sec. B.5.

726 Guidelines:

- 727 • The answer NA means that the paper does not include experiments.
728 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
729 or cloud provider, including relevant memory and storage.
730 • The paper should provide the amount of compute required for each of the individual
731 experimental runs as well as estimate the total compute.
732 • The paper should disclose whether the full research project required more compute
733 than the experiments reported in the paper (e.g., preliminary or failed experiments that
734 didn't make it into the paper).

735 **9. Code Of Ethics**

736 Question: Does the research conducted in the paper conform, in every respect, with the
737 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

738 Answer: [Yes]

739 Justification: Of course.

740 Guidelines:

- 741 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
742 • If the authors answer No, they should explain the special circumstances that require a
743 deviation from the Code of Ethics.
744 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
745 eration due to laws or regulations in their jurisdiction).

746 **10. Broader Impacts**

747 Question: Does the paper discuss both potential positive societal impacts and negative
748 societal impacts of the work performed?

749 Answer: [Yes]

750 Justification: We have discussed broader impacts in Sec. C

751 Guidelines:

- 752 • The answer NA means that there is no societal impact of the work performed.
753 • If the authors answer NA or No, they should explain why their work has no societal
754 impact or why the paper does not address societal impact.
755 • Examples of negative societal impacts include potential malicious or unintended uses
756 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
757 (e.g., deployment of technologies that could make decisions that unfairly impact specific
758 groups), privacy considerations, and security considerations.

- 759 • The conference expects that many papers will be foundational research and not tied
 760 to particular applications, let alone deployments. However, if there is a direct path to
 761 any negative applications, the authors should point it out. For example, it is legitimate
 762 to point out that an improvement in the quality of generative models could be used to
 763 generate deepfakes for disinformation. On the other hand, it is not needed to point out
 764 that a generic algorithm for optimizing neural networks could enable people to train
 765 models that generate Deepfakes faster.
- 766 • The authors should consider possible harms that could arise when the technology is
 767 being used as intended and functioning correctly, harms that could arise when the
 768 technology is being used as intended but gives incorrect results, and harms following
 769 from (intentional or unintentional) misuse of the technology.
- 770 • If there are negative societal impacts, the authors could also discuss possible mitigation
 771 strategies (e.g., gated release of models, providing defenses in addition to attacks,
 772 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
 773 feedback over time, improving the efficiency and accessibility of ML).

774 11. Safeguards

775 Question: Does the paper describe safeguards that have been put in place for responsible
 776 release of data or models that have a high risk for misuse (e.g., pretrained language models,
 777 image generators, or scraped datasets)?

778 Answer: [NA]

779 Justification: We do not see such danger.

780 Guidelines:

- 781 • The answer NA means that the paper poses no such risks.
- 782 • Released models that have a high risk for misuse or dual-use should be released with
 783 necessary safeguards to allow for controlled use of the model, for example by requiring
 784 that users adhere to usage guidelines or restrictions to access the model or implementing
 785 safety filters.
- 786 • Datasets that have been scraped from the Internet could pose safety risks. The authors
 787 should describe how they avoided releasing unsafe images.
- 788 • We recognize that providing effective safeguards is challenging, and many papers do
 789 not require this, but we encourage authors to take this into account and make a best
 790 faith effort.

791 12. Licenses for existing assets

792 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
 793 the paper, properly credited and are the license and terms of use explicitly mentioned and
 794 properly respected?

795 Answer: [Yes]

796 Justification: We have carefully checked and confirmed this.

797 Guidelines:

- 798 • The answer NA means that the paper does not use existing assets.
- 799 • The authors should cite the original paper that produced the code package or dataset.
- 800 • The authors should state which version of the asset is used and, if possible, include a
 801 URL.
- 802 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 803 • For scraped data from a particular source (e.g., website), the copyright and terms of
 804 service of that source should be provided.
- 805 • If assets are released, the license, copyright information, and terms of use in the
 806 package should be provided. For popular datasets, paperswithcode.com/datasets
 807 has curated licenses for some datasets. Their licensing guide can help determine the
 808 license of a dataset.
- 809 • For existing datasets that are re-packaged, both the original license and the license of
 810 the derived asset (if it has changed) should be provided.

- 811 • If this information is not available online, the authors are encouraged to reach out to
812 the asset's creators.

813 **13. New Assets**

814 Question: Are new assets introduced in the paper well documented and is the documentation
815 provided alongside the assets?

816 Answer: [NA]

817 Justification: There is no new assets provided.

818 Guidelines:

- 819 • The answer NA means that the paper does not release new assets.
820 • Researchers should communicate the details of the dataset/code/model as part of their
821 submissions via structured templates. This includes details about training, license,
822 limitations, etc.
823 • The paper should discuss whether and how consent was obtained from people whose
824 asset is used.
825 • At submission time, remember to anonymize your assets (if applicable). You can either
826 create an anonymized URL or include an anonymized zip file.

827 **14. Crowdsourcing and Research with Human Subjects**

828 Question: For crowdsourcing experiments and research with human subjects, does the paper
829 include the full text of instructions given to participants and screenshots, if applicable, as
830 well as details about compensation (if any)?

831 Answer: [NA]

832 Justification: Our work is irrelavent to crowd-sourcing or human subjects.

833 Guidelines:

- 834 • The answer NA means that the paper does not involve crowdsourcing nor research with
835 human subjects.
836 • Including this information in the supplemental material is fine, but if the main contribu-
837 tion of the paper involves human subjects, then as much detail as possible should be
838 included in the main paper.
839 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
840 or other labor should be paid at least the minimum wage in the country of the data
841 collector.

842 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
843 Subjects**

844 Question: Does the paper describe potential risks incurred by study participants, whether
845 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
846 approvals (or an equivalent approval/review based on the requirements of your country or
847 institution) were obtained?

848 Answer: [NA]

849 Justification: Our work is irrelavant to this front.

850 Guidelines:

- 851 • The answer NA means that the paper does not involve crowdsourcing nor research with
852 human subjects.
853 • Depending on the country in which research is conducted, IRB approval (or equivalent)
854 may be required for any human subjects research. If you obtained IRB approval, you
855 should clearly state this in the paper.
856 • We recognize that the procedures for this may vary significantly between institutions
857 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
858 guidelines for their institution.
859 • For initial submissions, do not include any information that would break anonymity (if
860 applicable), such as the institution conducting the review.