# PointDifformer: Robust Point Cloud Registration with Neural Diffusion and Transformer

Rui She*, Qiyu Kang*, Sijie Wang*, Wee Peng Tay, *Senior Member, IEEE,* Kai Zhao, Yang Song, Tianyu Geng, Yi Xu, Diego Navarro Navarro and Andreas Hartmannsgruber

*Abstract*—**Point cloud registration is a fundamental technique in 3D computer vision with applications in graphics, autonomous driving, and robotics. However, registration tasks under challenging conditions, under which noise or perturbations are prevalent, can be difficult. We propose a robust point cloud registration approach that leverages graph neural Partial Differential Equations (PDEs) and heat kernel signatures. Our method first uses graph neural PDE modules to extract high-dimensional features from point clouds by aggregating information from the 3D point neighborhood, thereby enhancing the robustness of the feature representations. Then, we incorporate heat kernel signatures into an attention mechanism to efficiently obtain corresponding keypoints. Finally, a Singular Value Decomposition module with learnable weights is used to predict the transformation between two point clouds. Empirical experiments on a 3D point cloud dataset demonstrate that our approach not only does achieve state-of-the-art performance for point cloud registration but also exhibits better robustness to additive noise or 3D shape perturbations.**

*Index Terms*—**Point cloud registration, neural diffusion, graph neural network, heat kernel signature.**

## I. INTRODUCTION

IN the era of intelligent and smart perception, 3D computer vision techniques are increasingly being used in various fields, such as autonomous driving, robotics, and scene modeling [1]–[3]. Point cloud registration is a crucial task in 3D computer vision and has become an important tool in many applications, including object detection, odometry estimation, and SLAM [4]–[8], owing to its robustness against seasonal changes and illumination variations. Point cloud registration aims to estimate the transformation or relative pose between two given 3D point cloud frames [9].

Iterative algorithms are widely used for point cloud registration [10]–[12]. The Iterative Closest Point (ICP) algorithm is a well-known iterative method for point cloud registration that matches the closest points between two point clouds, iteratively updating the transformation matrix until convergence [10]. ICP has been successfully used in numerous fields, including robotic perception and autonomous driving.

Despite their usefulness, iterative algorithms face challenges that limit their effectiveness in certain scenarios. The non-convexity of the optimization problem presents a significant

challenge, making it difficult to obtain the global optimum [9]. As a result, iterative algorithms may converge to sub-optimal solutions, especially in complex and non-rigid scenes. Additionally, the performance of iterative algorithms heavily relies on the initialization of the algorithm, which can be time-consuming and computationally expensive. Sparse and non-uniform point clouds present another significant challenge for iterative algorithms in finding corresponding point pairs between two point clouds. Traditional approaches, such as nearest-neighbor search, may fail to find matching pairs in such cases, leading to errors in the registration result [9], [13].

To address these challenges, deep learning-based methods [9], [13]–[15] have been developed for predicting transformation matrices or relative poses. These methods are designed for various scenarios, including indoor and outdoor environments [16]–[20]. However, robust point cloud registration remains a challenging problem due to factors such as LiDAR scan distortion, dynamic objects, and environmental noise [21]–[25]. Efficiently estimating the transformation under scenarios with additive noise and perturbations remains an open problem.

In this paper, we propose a model for point cloud registration that utilizes a robust feature descriptor based on graph neural diffusion. We also present an end-to-end transformation estimation method by introducing the heat kernel signature into the attention module, without any prior prediction information. Our approach attempts to address the challenges faced by iterative algorithms, leading to robust and efficient point cloud registration. Our proposed approach is motivated by the following:

- Graph neural PDE learning has demonstrated robustness for representing graph-structured data, as highlighted in [26]. Our aim is to leverage this module for effective point cloud representation by constructing a neighborhood graph in the feature space.
- We believe that the shape isometry-invariance of the heat kernel signature, as described in [27], makes it beneficial to incorporate into attention mechanisms for improved robustness from a shape-preserving perspective.

Our main contributions are as follows:

- We design a 3D point cloud representation module based on graph neural PDE learning.
- We propose a robust 3D point cloud registration method using the graph neural diffusion modules and the attention mechanism with a heat kernel signature.
- We empirically demonstrate that our point cloud registration method outperforms other baselines not only under

*R. She, Q. Kang and S. Wang, contributed equally to this work.

R. She, Q. Kang, S. Wang, W. P. Tay, K. Zhao, Y. Song, T. Geng, and Y. Xu are with Nanyang Technological University, Singapore. {rui.she@; qiyu.kang@; wang1679@e.; kai.zhao@; songy@; tianyu.geng@; yi.xu@; wptay@}ntu.edu.sg

D. N. Navarro and A. Hartmannsgruber are with Continental Automotive Singapore Pte. Ltd., Singapore. {diego.navarro.navarro; andreas.hartmannsgruber}@continental.com

normal scenarios but also when noise and perturbations are present.

The rest of this paper is organized as follows. In Section II, we discuss the related works. In Section III, we describe our proposed model in detail. In Section IV, we present the experimental results to evaluate our model and compare it with several baselines. Finally, we conclude the paper in Section V.

## II. RELATED WORK

In this section, we summarize relevant literature in the areas of point cloud registration, point cloud feature representation, and neural diffusion, including works utilizing the heat kernel signature in point cloud feature descriptors.

### A. Point Cloud Registration Methods

**Iteration-based Methods.** Iteration-based methods, such as the Iterative Closest Point (ICP) [10] and the RANdom SAmple Consensus (RANSAC) [28], are classical approaches commonly used for point cloud registration. However, due to its slow convergence rate, RANSAC requires high computing resources and has a high running time complexity. The performance of ICP heavily relies on the accuracy of the initial value estimation, making it prone to suboptimal solutions. To address these challenges, several refinement methods for ICP have been proposed, such as the Branch-and-Bound (BnB) method [11], convex relaxation improvement [29], and mixed-integer programming [30]. However, these methods may be computationally expensive and do not ensure global optimality. Alternatively, updated ICP methods such as Voxelized ICP [12] and Generalized-ICP [31] have been developed to improve both acceleration and accuracy.

**Correspondence-based estimators.** Correspondence-based estimators are commonly used for point cloud registration, which involves estimating the transformation between two frames [17], [32]–[34]. This approach obtains correspondences between two point clouds and then uses pose estimators such as RANSAC [18], [28], [35], Singular Value Decomposition (SVD) [9], [10], [13], [36] and Maximal Cliques (MAC) [37] to predict the transformation. There are generally two types of correspondence-based estimators. One involves repeatable keypoint detection [18], [35], [38], [39], followed by using learned or handcrafted keypoint descriptors for correspondence acquisition [17], [19], [40] or similarity measures to obtain the correspondences [2], [36]. For example, DeepVCP [39] uses PointNet++ [41] to extract features for the point clouds and learns keypoint correspondences based on matching probabilities among candidates. D3Feat [18] employs 3D fully convolutional networks to output detection scores and descriptive features for 3D points. PREDATOR [35] uses an overlap-attention block for cross-information between two point clouds and makes good use of their overlap region to achieve registration. The other [9], [20] involves correspondence retrieval for all possible matching point pairs without keypoint detection. For instance, Deep Closest Point (DCP) [9] aligns features based on the interaction of the point clouds. CoFiNet [20] achieves hierarchical correspondences with coarse and finer scales, without keypoint detection. In both types of

estimators, point cloud descriptors play significant roles, mainly contributing to the robustness and accuracy of the entire pipeline.

**Learning-based estimators.** In order to achieve more robust non-handcrafted estimators, learning-based methods are introduced into the transformation prediction [42]. Since conventional estimators like RANSAC have drawbacks in terms of convergence speed and are unstable in the presence of numerous outliers, learning-based estimators [16], [43]–[49], such as StickyPillars [50], PointDSC [51], EDFNet [52], GeoTransformer [42], Lepard [53], RoITr [54], BUFFER [55] and RoReg [56], have attracted much interest. Moreover, auxiliary modules or prior information can be incorporated into learning-based estimators, such as Prior-embedded Explicit Attention Learning (PEAL) [57] and VBReg [58]. Some of these classification neural networks can filter out extreme outliers and some estimation neural networks are designed to output the transformation. From the perspective of accuracy and running efficiency, they perform better than those conventional robust estimators. While, they need extra neural network training, which holds more time and space complexity. In contrast, our model achieves robust and accurate registration without the need for training the estimation networks to compute the final transformation in the output.

### B. Point Cloud Feature Representation

To extract more efficient features for point clouds, methods using different neural networks are studied. In general, we can classify point cloud feature representation methods into three categories as follows.

The first category performs voxel alignment on the points and then obtains the corresponding features based on a 3D Convolutional Neural Network (CNN) [59]–[62]. In this regard, the full information in the point cloud is used to learn the representation. However, it takes more computational resources to deal with a sparse and irregular point cloud when using closely spaced 3D voxels for more precise quantization.

The second category reduces a 3D point cloud to a 2D map and then exploits the classical 2D CNN to extract features [63]. The commonly used 2D maps are the bird's-eye view map, cylindrical map, spherical map, and camera-plane map, for which computational cost is incurred during the preprocessing stage. Due to quantization errors, this approach can also introduce unexpected noise.

The third category is to extract features from the raw point clouds directly using specific neural networks. PointNet [64] and PointNet++ [41] extract local point features independently and obtain global features through max-pooling. To incorporate local neighborhood information, Dynamic Graph Convolutional Neural Networks (DGCNN) [65] uses a dynamic graph network, and LPDNet [66] jointly exploits the geometry space and feature space. KPConv [67] uses kernel points to achieve more flexible convolutions compared with fixed grid convolutions. PointGLR [68] considers not only local features but also global patterns in point clouds. DIP [69] and GeDi [46] extract local point cloud patches based on rotation-invariant compact descriptors, which can be used in different data domains. Point

Cloud Transformer (PCT) [70] utilizes the Transformer to generate permutation-invariant descriptors for point clouds. Moreover, there are also other learning-based point cloud representation methods, such as PointMLP [71], and PointNeXt [72].

### C. Neural Diffusion

Neural diffusion methods [73], [74] combine neural networks with ordinary and partial differential equations. For a neural Ordinary Differential Equation (ODE) layer [73] with input $\mathbf{Z}(0)$ and output $\mathbf{Z}(T)$, the relationship between $\mathbf{Z}(0)$ and $\mathbf{Z}(T)$ is given by

$$\frac{\mathrm{d}\mathbf{Z}(t)}{\mathrm{d}t} = f_{\mathrm{ODE}}(\mathbf{Z}(t), t), \tag{1}$$

where $f_{\mathrm{ODE}} : \mathbb{R}^n \times [0, T) \to \mathbb{R}^n$ is a learnable layer and $\mathbf{Z} : [0, T) \to \mathbb{R}^n$ denotes the state of the neural ODE. At the terminal time $T \in [0, \infty)$, the output $\mathbf{Z}(T)$ is given by

$$\mathbf{Z}(T) = \mathbf{Z}(0) + \int_0^T f_{\mathrm{ODE}}(\mathbf{Z}(t), t)\mathrm{d}t. \tag{2}$$

In this paper, we consider only the time-invariant (autonomous) case, i.e., $f_{\mathrm{ODE}}(\mathbf{Z}(t), t) = f_{\mathrm{ODE}}(\mathbf{Z}(t))$.

For graph-structured data, graph neural PDEs are designed based on continuous flows, which represent the graph features more concisely and stably [26], [74]–[80]. Neural ODEs/PDEs are more robust in defending perturbations and even attacks, compared with other deep neural networks without neural diffusion (cf. [26], [81]). Compared with conventional graph neural networks (GNNs), including the Graph ATtention network (GAT) or the Graph Convolutional Network (GCN), graph neural PDEs have superior performance in some applications such as the node classification for graph-structured data. To approximately solve the graph neural PDEs [74], the neural ODE solvers proposed in [73] can be exploited.

### D. Heat Kernel Signature

Based on the heat diffusion process, the heat kernel signature [27] is presented as an intrinsic feature and is given by

$$h(x, t) = \sum_{i=0}^{\infty} \exp\left(-\lambda_i t\right) \phi_i^2(x), \tag{3}$$

where $x$ is a 3D point in a point cloud, $\lambda_i$s denote eigenvalues and $\phi_i$s are the corresponding eigenfunctions of the Laplace-Beltrami operator. The feature $h(x, t)$ is a robust local geometric descriptor containing large-scale information [82], [83]. From the physics perspective, this feature descriptor represents the temperature evolution of a point at which a heat source is placed and removed immediately. The heat diffuses to the neighborhood of the point [27], [84]. This evolution is based on the temperature diffusion speed, which essentially depends on the geometry of the objects projected by the point clouds.

From a geometric perspective, the heat kernel signature is isometry-invariant, meaning that two isometric shapes have equivalent heat kernel signatures. If the heat kernel signatures of two shapes are equal, the corresponding shapes or parts of the shapes are similar under isometric transformations [27]. Therefore, this feature is somewhat robust, making it a desirable method for point description.

### III. POINT CLOUD REGISTRATION METHOD BASED ON GRAPH NEURAL PDE

In this section, we present our registration model that aims to predict the transformation between two 3D point clouds. However, point clouds may contain noise or perturbations that can compromise the robustness of the transformation prediction. Therefore, our goal is to develop a more robust method for the registration task.

Our model is called the *Point Cloud Diffusion Transformer (PointDifformer)*. First, we provide an overview of the PointDifformer framework, which is illustrated in Fig. 1. Then, we present the details of the modules and the loss function used in PointDifformer.

### A. Overview of PointDifformer

Before introducing in detail the modules of PointDifformer, we provide an overview of its pipeline as follows.

1) Within a 3D point cloud frame, the neighborhood graph of each point consisting of its $K$ nearest neighbors is constructed. Points are regarded as the vertices in the graph. The $\mathcal{L}_2$ distance between point features is used for neighbor acquisition. The initial features of the points are taken to be their 3D coordinates. The neighborhood graph for each point is an undirected complete graph. Then, graph neural PDE layers are applied to the neighborhood graph of each point to obtain a robust representation of the point.

2) Based on the robust feature representations, a self-cross attention module is applied to obtain an embedding containing point-level information interaction within a point cloud frame and between two point cloud frames. The heat kernel signature, as a robust feature descriptor for point clouds, is introduced into the attention module as the weights.

3) Using the above embeddings, an attention module is established to learn weights for points from different frames that indicate their correspondences.

4) Through the correspondence among points in the two point-cloud frames, the optimal transformation (including the translation and the rotation) can be estimated using optimization solution methods like the weighted SVD.

### B. Model Details

*1) Point Cloud Representation with Neural Diffusion:* To represent point cloud features efficiently and robustly, we design a neural diffusion network for point cloud representation, called *Point-Diffusion Net*. This module consists of GNN modules and graph neural PDE modules with different rewiring and is shown in Fig. 2. Its details are described as follows.

**3D Points Refinement.** To pre-process the point cloud, which has potential outliers, we first use a graph neural PDE module as a learning-based filter. Consider a point cloud denoted by $\mathbf{X} \in \mathbb{R}^{N \times 3}$. The 3-dimensional coordinates of each point are regarded as its feature map and $N$ is the number of points in the point cloud. We construct the neighborhood graphs for the points by means of the $K$-Nearest Neighbors
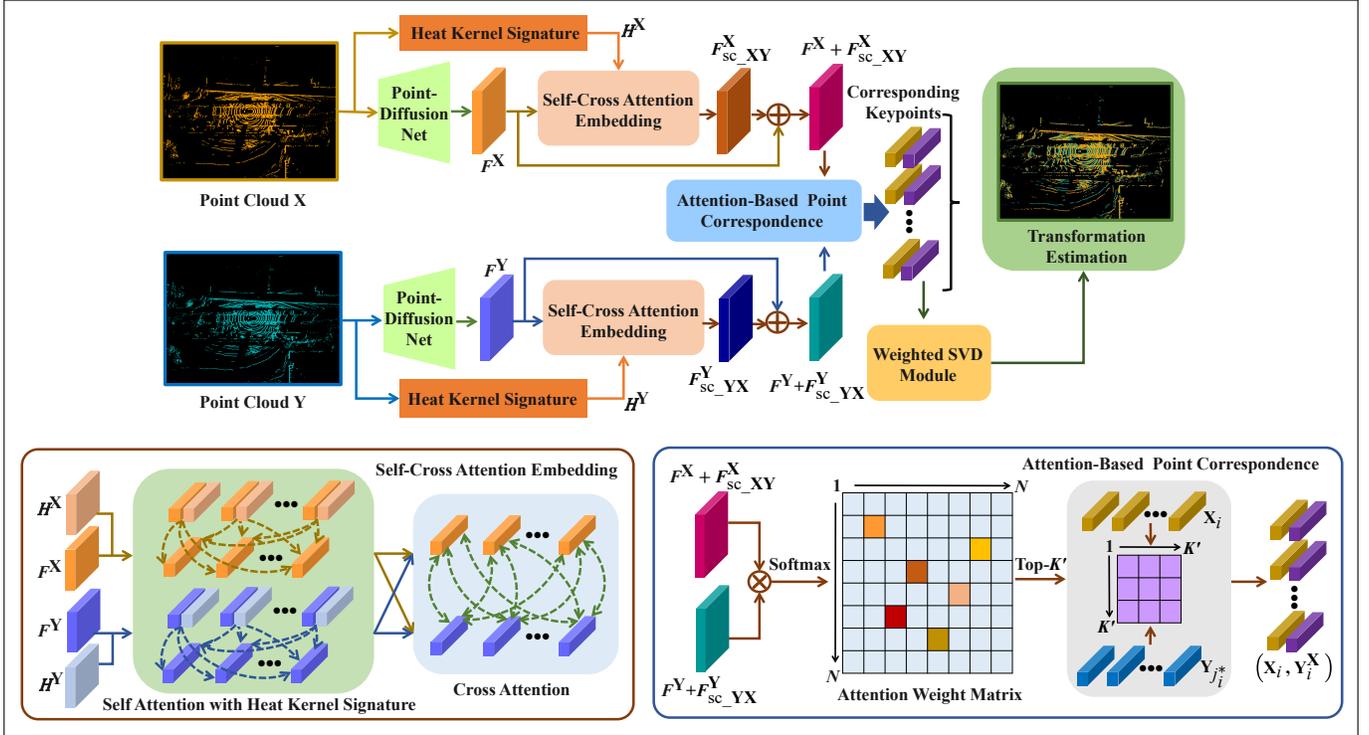
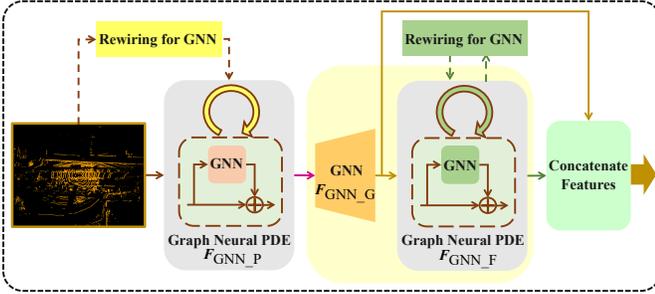Fig. 1. PointDifformer for point cloud registration. The details of the modules are provided in Section III-B.



Fig. 2. The Point-Diffusion Net based on the graph neural PDEs for the point cloud representation. The details are provided in Section III-B1.

($K$-NN) method using the Euclidean distance. Furthermore, we exploit a graph neural PDE module for feature updating, given by

$$\frac{d\mathbf{Z}_p(t)}{dt} = f_{\text{GNN\_P}}(\mathbf{Z}_p(t)), \quad (4)$$

where $f_{\text{GNN\_P}}(\cdot)$ denotes a graph learning module and $\mathbf{Z}_p(t)$ is the state at time $t$. The initial state is given by $\mathbf{Z}_p(0) = \mathbf{X}$. By integrating $f_{\text{GNN\_P}}(\cdot)$ from $t = 0$ to $t = T_p$ (using differential equation solvers [85]), we obtain the solution of (4) at time $T_p$, given by

$$\mathbf{Z}_p(T_p) = F_{\text{GNN\_P}}(\mathbf{Z}_p(0)) = F_{\text{GNN\_P}}(\mathbf{X}) \in \mathbb{R}^{N \times 3}, \quad (5)$$

where $F_{\text{GNN\_P}}(\cdot)$ can be regarded as the embedding function for the input $\mathbf{Z}_p(0)$. In addition, the output of this graph neural module also includes 3-dimensional features, which can also be viewed as "generated" points.

**High Dimensional Feature Extraction with Graph Neural PDE.** We extract high dimensional features for the preprocessed 3D points through a GNN module, e.g., DGCNN [65]. In this regard, 3-dimensional coordinates of points are extended into $d$-dimensional features. We construct a neighborhood graph for each point using the $K$-NN method. The output from the GNN module is denoted by

$$F_{\text{G}}(\mathbf{X}) = F_{\text{GNN\_G}} \circ F_{\text{GNN\_P}}(\mathbf{X})$$
$$= F_{\text{GNN\_G}}(\mathbf{Z}_p(T_p)) \in \mathbb{R}^{N \times d}, \quad (6)$$

where $\circ$ denotes function composition and $F_{\text{GNN\_G}}(\cdot)$ denotes the mapping of the GNN module. Then, we apply another graph neural PDE module to update the feature $F_{\text{G}}(\mathbf{X})$, which is described as

$$\frac{d\mathbf{Z}_f(t)}{dt} = f_{\text{GNN\_F}}(\mathbf{Z}_f(t)), \quad (7)$$

where $f_{\text{GNN\_F}}(\cdot)$ is also a graph learning module that deals with the neighborhood graph of input features. The initial state is given by $\mathbf{Z}_f(0) = F_{\text{G}}(\mathbf{X})$. The equation (7) is solved in the same way as that for (4). The output at time $T_f$ is given by

$$\mathbf{Z}_f(T_f) = F_{\text{GNN\_F}}(\mathbf{Z}_f(0)) = F_{\text{GNN\_F}} \circ F_{\text{G}}(\mathbf{X}), \quad (8)$$

where $\mathbf{Z}_f(T_f) \in \mathbb{R}^{N \times d}$, $F_{\text{GNN\_F}}(\cdot)$ can be regarded as the embedding function for the input $\mathbf{Z}_f(0)$.

Finally, we concatenate the output from the GNN module and the graph neural diffusion module as the eventual output for the point cloud representation module, which is denoted by

$$F^{\mathbf{X}} = F_{\text{G}}(\mathbf{X}) \,\big\|\, F_{\text{GNN\_F}} \circ F_{\text{G}}(\mathbf{X}), \quad (9)$$

where $F^{\mathbf{X}} \in \mathbb{R}^{N \times 2d}$ and $\|$ denotes the concatenation operation. Here, we follow the same approach as [13], [65] in concatenating the feature $\mathbf{Z}_f(T_f) = F_{\text{GNN\_F}} \circ F_{\text{G}}(\mathbf{X})$ with its corresponding hidden feature $F_{\text{G}}(\mathbf{X})$ to retain more information. Our numerical experiments indicate that this is a better approach than using only $\mathbf{Z}_f(T_f)$.

*2) Self-Cross Attention Embedding Based on Heat Kernel Signature:* Based on the previous high-dimensional feature, a self-cross attention mechanism is introduced to reinforce the static structure information in each point cloud and the interactive corresponding information between a pair of point clouds, respectively. To improve the robustness of the feature extraction, we also introduce the heat kernel signature [27] into the attention mechanism as additive weights.

For a pair of point clouds $(\mathbf{X}, \mathbf{Y})$, the corresponding input feature pair for the self-cross attention module is represented as $(F^{\mathbf{X}}, F^{\mathbf{Y}})$ using the embedding from (9). The embedding from the self-cross attention module with respect to (w.r.t.) $\mathbf{X}$ is given by

$$F^{\mathbf{X}}_{\text{sc\_att}}(F^{\mathbf{X}}, F^{\mathbf{Y}}) = F_{\text{s\_att}}(F^{\mathbf{X}}) + F^{\mathbf{X}}_{\text{c\_att}}(F^{\mathbf{X}}, F^{\mathbf{Y}}), \quad (10)$$

where $F_{\text{s\_att}}(F^{\mathbf{X}})$ and $F^{\mathbf{X}}_{\text{c\_att}}(F^{\mathbf{X}}, F^{\mathbf{Y}})$ are the features based on the self-attention and the cross-attention, respectively. The details of the features $F_{\text{s\_att}}(F^{\mathbf{X}})$ and $F^{\mathbf{X}}_{\text{c\_att}}(F^{\mathbf{X}}, F^{\mathbf{Y}})$ are described as follows.

**Self-attention Feature.** To improve the robustness, the heat kernel signature is introduced into the self-attention module. For a point cloud pair $(\mathbf{X}, \mathbf{Y})$, the corresponding heat kernel signature pair is denoted by $(H^{\mathbf{X}}, H^{\mathbf{Y}})$. Using the normalized $F^{\mathbf{X}}$ and $H^{\mathbf{X}}$ as the inputs for the self-attention module, we have

$$F_{\text{s\_att}}(F^{\mathbf{X}})$$
$$= \mathbf{W}^{\text{S}} \bigg\|_{i=1}^{S_{\text{head}}} \bigg\{ F_{\text{softmax}}\Big( \frac{(\mathbf{W}^{\text{SQ}}_i F^{\mathbf{X}})(\mathbf{W}^{\text{SK}}_i F^{\mathbf{X}})^{\intercal}}{\sqrt{d^{\text{S}}_i}}$$
$$+ \frac{(\mathbf{W}^{\text{HQ}}_i H^{\mathbf{X}})(\mathbf{W}^{\text{HK}}_i H^{\mathbf{X}})^{\intercal}}{\sqrt{d^{\text{H}}_i}} \Big)(\mathbf{W}^{\text{SV}}_i F^{\mathbf{X}}) \bigg\} + F^{\mathbf{X}}, \quad (11)$$

where $(\cdot)^{\intercal}$ denotes the transpose operation, $S_{\text{head}}$ is the number of multi-heads for the attention, $\mathbf{W}^{\text{SQ}}_i$, $\mathbf{W}^{\text{SK}}_i$, $\mathbf{W}^{\text{SV}}_i$, $\mathbf{W}^{\text{HQ}}_i$, $\mathbf{W}^{\text{HK}}_i$, and $\mathbf{W}^{\text{S}}$ are learnable layers, and $d^{\text{S}}_i$ and $d^{\text{H}}_i$ are the dimensions for the point cloud features and heat kernel signatures in $i$-th attention head, respectively. The function $F_{\text{softmax}}(\cdot)$ denotes row-wise softmax.

The heat kernel signature is implemented as follows.

- *Heat kernel signature acquisition.* We compute this feature based on the point cloud using the formula (3). Since the function $h(x, t)$ in (3) is a robust local geometric descriptor containing large-scale information [82], we use it to robustly describe the repeatable features for point clouds.
- *Embedding.* We process the heat kernel signatures using the graph neural PDE module and the Fully Connected (FC) layer to obtain the embedding. Similar to (4), the graph neural PDE is used as a filter for the heat kernel signatures. The graph construction in the graph neural PDE is based on the $K$-NN for the heat kernel signatures.

- *Self-attention weights.* After embedding the heat kernel signatures, they are input into the self-attention module as the additive weights. By introducing extra information from the heat kernel signature, the robustness of the point cloud representation is enhanced in the self-attention module.

**Cross-attention Feature.** Based on the self-attention feature $F_{\text{s\_att}}(F^{\mathbf{Y}})$ for the point cloud $\mathbf{Y}$, we acquire the cross-attention features for $F^{\mathbf{X}}$. $F_{\text{s\_att}}(F^{\mathbf{Y}})$ is input into a Feed Forward Network (FFN) [86] to obtain the feature

$$F_{\text{s\_att\_n}}(F^{\mathbf{Y}}) = F_{\text{FNN}}(F_{\text{s\_att}}(F^{\mathbf{Y}})) + F_{\text{s\_att}}(F^{\mathbf{Y}}), \quad (12)$$

where $F_{\text{FNN}}$ denotes the FFN consisting of two linear layers with normalization operation and the Rectified linear activation function (ReLU). Inputting normalized $F_{\text{s\_att\_n}}(F^{\mathbf{Y}})$ and $F_{\text{s\_att}}(F^{\mathbf{x}})$ into the cross-attention module, we have

$$F^{\mathbf{X}}_{\text{c\_att}}(F^{\mathbf{X}}, F^{\mathbf{Y}}) = \mathbf{W}^{\text{C}} \times$$
$$\bigg\|_{i=1}^{C_{\text{head}}} \bigg\{ F_{\text{softmax}}\Big( \frac{(\mathbf{W}^{\text{CQ}}_i F_{\text{s\_att}}(F^{\mathbf{X}}))(\mathbf{W}^{\text{CK}}_i F_{\text{s\_att\_n}}(F^{\mathbf{Y}}))^{\intercal}}{\sqrt{d^{\text{C}}_i}} \Big)$$
$$\times (\mathbf{W}^{\text{CV}}_i F_{\text{s\_att\_n}}(F^{\mathbf{Y}})) \bigg\}, \quad (13)$$

where $\mathbf{W}^{\text{CQ}}_i$, $\mathbf{W}^{\text{CK}}_i$, $\mathbf{W}^{\text{CV}}_i$, and $\mathbf{W}^{\text{C}}$ are learnable layers, $d^{\text{C}}_i$ is the feature dimension for the $i$-th attention head. The remaining notations are similar to those in (11).

The joint feature $F^{\mathbf{X}}_{\text{sc\_att}}(F^{\mathbf{X}}, F^{\mathbf{Y}})$ based on the self-cross attention module is obtained as mentioned in (10). Inputting $F^{\mathbf{X}}_{\text{sc\_att}}(F^{\mathbf{X}}, F^{\mathbf{Y}})$ into the FNN, we have the embedding from the self-cross attention module as

$$F^{\mathbf{X}}_{\text{sc\_XY}} = F_{\text{FNN}}(F^{\mathbf{X}}_{\text{sc\_att}}(F^{\mathbf{X}}, F^{\mathbf{Y}})) + F^{\mathbf{X}}_{\text{sc\_att}}(F^{\mathbf{X}}, F^{\mathbf{Y}}), \quad (14)$$

whose normalization is regarded as the final output of this module.

Similarly, the above self-cross attention module is also available for the point cloud $\mathbf{Y}$ to obtain its output $F^{\mathbf{Y}}_{\text{sc\_YX}}$. The architectures of the self-attention and the cross-attention are shown in Fig. 3.

*3) Attention-Based Keypoint Correspondence:* Using the attention mechanism, the information of point cloud $\mathbf{X}$ can be involved in the embedding for the point cloud $\mathbf{Y}$. By resorting to the self-cross attention embeddings, we can obtain the weighted $\mathbf{Y}$ denoted by $\mathbf{Y}^{\mathbf{X}}$ which is regarded as the transformed point cloud corresponding to the $\mathbf{X}$. The details are given as follows.

*i)* We compute the attention weight matrix

$$\mathbf{W}^{\text{att}} = F_{\text{softmax}}\left( \frac{(F^{\mathbf{X}} + F^{\mathbf{X}}_{\text{sc\_XY}})(F^{\mathbf{Y}} + F^{\mathbf{Y}}_{\text{sc\_YX}})^{\intercal}}{\sqrt{d^{\text{att}}}} \right), \quad (15)$$

where the feature dimension $d^{\text{att}} = 2d$ and $F_{\text{softmax}}$ denotes the row-wise softmax function.

*ii)* For each point $\mathbf{x}_i$ in the point cloud $\mathbf{X}$, we select its corresponding point $\mathbf{y}_{j^*_i}$ in the point cloud $\mathbf{Y}$, which has the highest similarity with the $\mathbf{x}_i$. Furthermore, based on the Top-$K'$ similarity scores, we select the corresponding point pairs
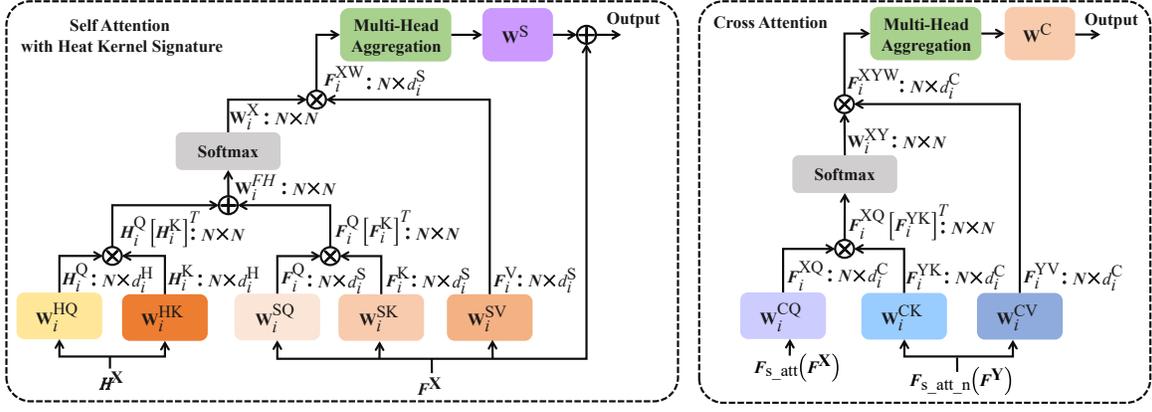
Fig. 3. The modules of the self-attention with heat kernel signature and the cross-attention.

$\left\{(\mathbf{x}_i, \mathbf{y}_{j_i^*}) : i = 1, 2, ..., K'\right\}$. Specifically, for the $i$-th row of $\mathbf{W}^{\text{att}}$, we have

$$j_i^* = \max_j \mathbf{w}_{i,j}^{\text{att}}, \tag{16}$$

where $\mathbf{w}_{i,j}^{\text{att}}$ ($i, j \in \{1, 2, ..., N\}$) denotes the element in the $i$-th row and $j$-th column among the $\mathbf{W}^{\text{att}}$. Then, based on the $\mathbf{w}_{i,j_i^*}^{\text{att}}$ ($i \in \{1, 2, ..., N\}$), we select the Top-$K'$ point pairs $\left\{(\mathbf{x}_i, \mathbf{y}_{j_i^*}) : i = 1, 2, ..., K'\right\}$ to obtain the updated point cloud pairs $(\mathbf{X}, \mathbf{Y})$, where abusing notations $\mathbf{X}$ and $\mathbf{Y}$ are used.

*iii)* Based on the updated point clouds $\mathbf{X}$ and $\mathbf{Y}$ consisting of the Top-$K'$ points, we have the corresponding updated attention weight matrix similar to (15), which is denoted by an abusing notation $\mathbf{W}^{\text{att}}$. Then, we have the weighted $\mathbf{Y}$ as

$$\mathbf{Y}^{\mathbf{X}} = \mathbf{W}^{\text{att}}\mathbf{Y}, \tag{17}$$

whose the point number is also $K'$ the same as that in the $\mathbf{X}$.

In general, the points $\mathbf{x}_i$ and $\mathbf{y}_i^{\mathbf{x}}$ ($i = 1, 2, ..., K'$) from the point clouds $\mathbf{X}$ and $\mathbf{Y}^{\mathbf{X}}$ are treated as the correspondence points.

*4) Transformation Prediction:* By resorting to the correspondence of points, we can predict the transformation or relative pose between two point clouds. Consider the Mean Squared Error (MSE) given by

$$\ell_{\text{MSE}}(\hat{\mathbf{R}}, \hat{\mathbf{t}}) = \frac{1}{K'} \sum_{i=1}^{K'} \|\hat{\mathbf{R}}\mathbf{x}_i + \hat{\mathbf{t}} - \mathbf{y}_i^{\mathbf{x}}\|_2, \tag{18}$$

in which $\|\cdot\|_2$ denotes the $\mathcal{L}_2$ norm, $\mathbf{x}_i = [x_i^{(1)}, x_i^{(2)}, x_i^{(3)}]^{\mathsf{T}}$ and $\mathbf{y}_i^{\mathbf{x}} = [y_i^{x(1)}, y_i^{x(2)}, y_i^{x(3)}]^{\mathsf{T}}$ where $x_i^{(l)}$ and $y_i^{x(l)}$ ($l \in \{1, 2, 3\}$) are elements from $\mathbf{x}_i$ and $\mathbf{y}_i^{\mathbf{x}}$, respectively. The $\hat{\mathbf{R}} \in \mathbb{R}^{3\times3}$ and $\hat{\mathbf{t}} \in \mathbb{R}^{3\times1}$ are the predicted results w.r.t. the ground-truth rotation $\mathbf{R} \in \mathbb{R}^{3\times3}$ and translation $\mathbf{t} \in \mathbb{R}^{3\times1}$. The optimal results of $\hat{\mathbf{R}}$ and $\hat{\mathbf{t}}$ are given by

$$\hat{\mathbf{R}}^*, \hat{\mathbf{t}}^* = \arg\min_{\hat{\mathbf{R}}, \hat{\mathbf{t}}} \ell_{\text{MSE}}(\hat{\mathbf{R}}, \hat{\mathbf{t}}). \tag{19}$$

Then, we use the weighted SVD [10], [42] to solve the optimization problem in (19). Specifically, the weighted mean

of the $\{\mathbf{x}_i : i = 1, 2, ..., K'\}$ and the $\{\mathbf{y}_i^{\mathbf{x}} : i = 1, 2, ..., K'\}$ are first computed as

$$\mathbf{x}_{\mathbf{w}} = \frac{1}{K'} \sum_{i=1}^{K'} \mathbf{w}_i^{\mathbf{x}} \mathbf{x}_i, \tag{20}$$

$$\mathbf{y}_{\mathbf{w}}^{\mathbf{x}} = \frac{1}{K'} \sum_{i=1}^{K'} \mathbf{w}_i^{\mathbf{y}} \mathbf{y}_i^{\mathbf{x}}, \tag{21}$$

where $\mathbf{w}_i^{\mathbf{x}} \in \mathbb{R}^{3\times1}$ and $\mathbf{w}_i^{\mathbf{y}} \in \mathbb{R}^{3\times1}$ are trainable weights.

Furthermore, the weighted cross-covariance matrix $\mathbf{M}$ is given by

$$\mathbf{M} = \sum_{i=1}^{K'} (\mathbf{x}_i - \mathbf{x}_{\mathbf{w}})(\mathbf{w}_i^{\mathbf{M}}(\mathbf{y}_i^{\mathbf{x}} - \mathbf{y}_{\mathbf{w}}^{\mathbf{x}}))^{\mathsf{T}}, \tag{22}$$

where $(\cdot)^{\mathsf{T}}$ denotes the transpose operation, $\mathbf{w}_i^{\mathbf{M}} \in \mathbb{R}^{3\times1}$ is a trainable weight.

Similar to the procedure of SVD mentioned in [9], [13], the matrix $\mathbf{M}$ can be decomposed as

$$\mathbf{M} = \mathbf{U}\Lambda\mathbf{V}^{\mathsf{T}}, \tag{23}$$

where $\mathbf{U}$ and $\mathbf{V}$ are unitary matrices and $\Lambda$ is a rectangular diagonal matrix with non-negative real diagonal elements. Furthermore, the transformation prediction (including the predicted rotation $\hat{\mathbf{R}}^*$ and the translation $\hat{\mathbf{t}}^*$) can be obtained as

$$\hat{\mathbf{R}}^* = \mathbf{V}\mathbf{U}^{\mathsf{T}}, \tag{24}$$

$$\hat{\mathbf{t}}^* = -\hat{\mathbf{R}}^*\mathbf{x}_{\mathbf{w}} + \mathbf{y}_{\mathbf{w}}^{\mathbf{x}}. \tag{25}$$

*5) Loss function:* As the point $\mathbf{y}_i^{\mathbf{x}}$ corresponds to the point $\mathbf{x}_i$, we use the corresponding point loss given by

$$\mathcal{L}_{\text{point}} = \frac{1}{K'} \sum_{i=1}^{K'} \|\hat{\mathbf{R}}^*\mathbf{x}_i + \hat{\mathbf{t}}^* - \mathbf{y}_i^{\mathbf{x}}\|_2. \tag{26}$$

To quantify the deviation between the ground truth and the predicted results w.r.t. rotation and translation, we use the loss given by

$$\mathcal{L}_{\text{rt}} = \exp(-\gamma_{\text{t}})\|\hat{\mathbf{t}}^* - \mathbf{t}\|_2 + \gamma_{\text{t}} + \exp(-\gamma_{\text{r}})\|\mathbf{R}^{\mathsf{T}}\hat{\mathbf{R}}^* - \mathbf{I}\|_2 + \gamma_{\text{r}}, \tag{27}$$

where $\mathbf{I}$ denotes the identity matrix, and $\gamma_{\mathrm{t}}$ and $\gamma_{\mathrm{r}}$ are learnable parameters. The learnable weights based on $\gamma_{\mathrm{t}}$ and $\gamma_{\mathrm{r}}$ are inpsired by the loss in [77], [87]. The total loss combining $\mathcal{L}_{\mathrm{point}}$ and $\mathcal{L}_{\mathrm{rt}}$ is given by

$$\mathcal{L}_{\mathrm{total}} = \exp\left(-\eta_{\mathrm{p}}\right)\mathcal{L}_{\mathrm{point}} + \eta_{\mathrm{p}} + \exp\left(-\eta_{\mathrm{rt}}\right)\mathcal{L}_{\mathrm{rt}} + \eta_{\mathrm{rt}}, \quad (28)$$

where $\eta_{\mathrm{p}}$ and $\eta_{\mathrm{rt}}$ are learnable parameters.

## IV. EXPERIMENTS

### A. Dataset preparation

**vReLoc Dataset.** The vReLoc Dataset is a publicly available indoor dataset,[1] containing LiDAR point clouds and camera images. In this paper, we randomly generate a transformation matrix for each point cloud to obtain a pair of point cloud frames. The transformation matrix is based on translation along the $x$, $y$, and $z$ axes, as well as rotation along the roll, pitch, and yaw axes. The generated transformation matrix is regarded as the ground truth. The generated translation values are uniformly sampled from the intervals $[-1, 1]$, $[-2, 2]$, and $[-0.5, 0.5]$ along the $x$-, $y$-, and $z$-axes, respectively. The generated rotation values are uniformly sampled from the intervals $[0°, 5°]$, $[0°, 5°]$, and $[0°, 30°]$ around the roll, pitch, and yaw axes, respectively.

**Boreas Dataset.** The Boreas dataset is a publicly available outdoor dataset[2] that comprises multi-sensor data, including LiDAR and camera data. It presents various environmental scenarios, such as sunny, night, and rainy conditions, as it was collected over the course of one year by repeatedly driving a specific route. Furthermore, the dataset provides post-processed ground-truth poses with centimeter-level accuracy, which offers the transformation matrix required for two consecutive LiDAR point clouds. The Boreas datasets undergo preprocessing involving distortion correction of LiDAR point clouds, as detailed in [88]. However, these preprocessing techniques do not completely eliminate all distortions in LiDAR point clouds, such as the tailing phenomenon [89]. Additionally, noise can still be present in environments with adverse weather conditions, dynamic objects or vehicles, and pedestrians, which can affect the accuracy of the LiDAR measurements.

**KITTI Dataset.** The KITTI dataset is a publicly available outdoor dataset[3] that provides multi-sensor data for autonomous driving. It includes LiDAR point clouds of street scenes captured using the Velodyne Laserscanner in Karlsruhe, Germany, with tens of thousands of LiDAR points in each frame. The dataset consists of 11 sequences (from sequence "0" to "10") depicting different street scenes, and global ground-truth poses are available for each sequence. Similar to the Boreas Dataset, we can use the ground-truth poses to obtain the transformation matrix between each pair of adjacent LiDAR point clouds in the KITTI dataset. While the KITTI dataset incorporates preprocessing for point cloud calibration [90] through auxiliary sensors, such as Global Positioning System (GPS) and Inertial Measurement Unit (IMU), perturbations similar to those in the Boreas datasets persist.

[1] https://github.com/loveoxford/vReLoc
[2] https://www.boreas.utias.utoronto.ca/
[3] http://www.cvlibs.net/datasets/kitti/

### B. Experimental Details

**Model Setting.** We set $d = 256$ in (4). To deal with the neighborhood graph of the $K$ nearest neighbors ($K = 20$), we set the GNN layer $f_{\mathrm{GNN\_P}}$ in the graph neural PDE (4) to be the union EdgeConv layers [65] which is also regarded as a kind of DGCNN. There are 5 EdgeConv layers used in the DGCNN block, whose hidden input and output dimensions are given by $[6, 16]$, $[16, 16]$, $[16, 32]$, $[32, 64]$ and $[128, 3]$, respectively. We set the graph learning module for the $F_{\mathrm{GNN\_G}}$ to another DGCNN, in which 5 EdgeConv layers are used with the hidden input and output dimensions $[6, 64]$, $[64, 64]$, $[64, 128]$, $[128, 256]$ and $[512, 256]$, respectively. When using the DGCNN-based graph neural PDE for $F_{\mathrm{GNN\_F}}$, there are 2 EdgeConv layers used in the DGCNN block with the hidden input and output dimensions $[256, 256]$ and $[768, 256]$. As for the self-cross attention module, there are 4 attention heads with 128 hidden features for each attention head, which implies 512 hidden features in total. We adopt the Adam optimizer [91] in the training, where the learning rate is set as 0.0001. We set the number of training epochs as 50.

**Baseline Methods.** To demonstrate the superior performance of PointDifformer, we compared it against several baseline methods, including ICP [31], DCP [9], HGNN [92], VCR-Net [13], PCT [70], and GeoTransformer [42]. ICP is an iterative optimization method that does not require neural networks for feature learning, meaning that it does not need a training process. On the other hand, the other methods utilize learned point cloud features to determine point correspondence such as DCP, VCR-Net and GeoTransformere. We further enhanced HGNN and PCT with attention modules for point correspondence registration, which we refer to as HGNN++ and PCT++, respectively.

### C. Point Cloud Registration Performance

*1) Evaluation on Indoor vReLoc Dataset:* To compare our method with other baselines, we evaluate them on their ability to predict transformations between two nearby point cloud frames from the vReLoc dataset. For training, we use sequences "3", "6", and "9", while for testing, we use sequences "14" and "16". We evaluate the performance of these methods using statistics such as Mean Absolute Error (MAE) and Root Mean Square Error (RMSE) for the predicted relative translation and rotation results. We also utilize Registration Recall (RR) as a metric, which is defined in [37], [42], [54]. RR measures the percentage of point cloud frames that achieve a certain threshold of registration accuracy. In our evaluation, we use a fine-tuned threshold for performance comparison.

From Table I, we observe that PointDifformer outperforms the other baselines without neural diffusions in terms of relative translation and rotation prediction. This suggests that the graph neural PDE modules play a positive role in point cloud registration. Further analysis of PointDifformer in Fig. 5 shows that the translation and rotation errors lie in a small region close to zero. In Fig. 5, the empirical probability or the relative frequency of an error value, is the ratio of the number of errors within a small bin around the error value to the total number of trials. The relative translation error and rotation error in Fig. 5
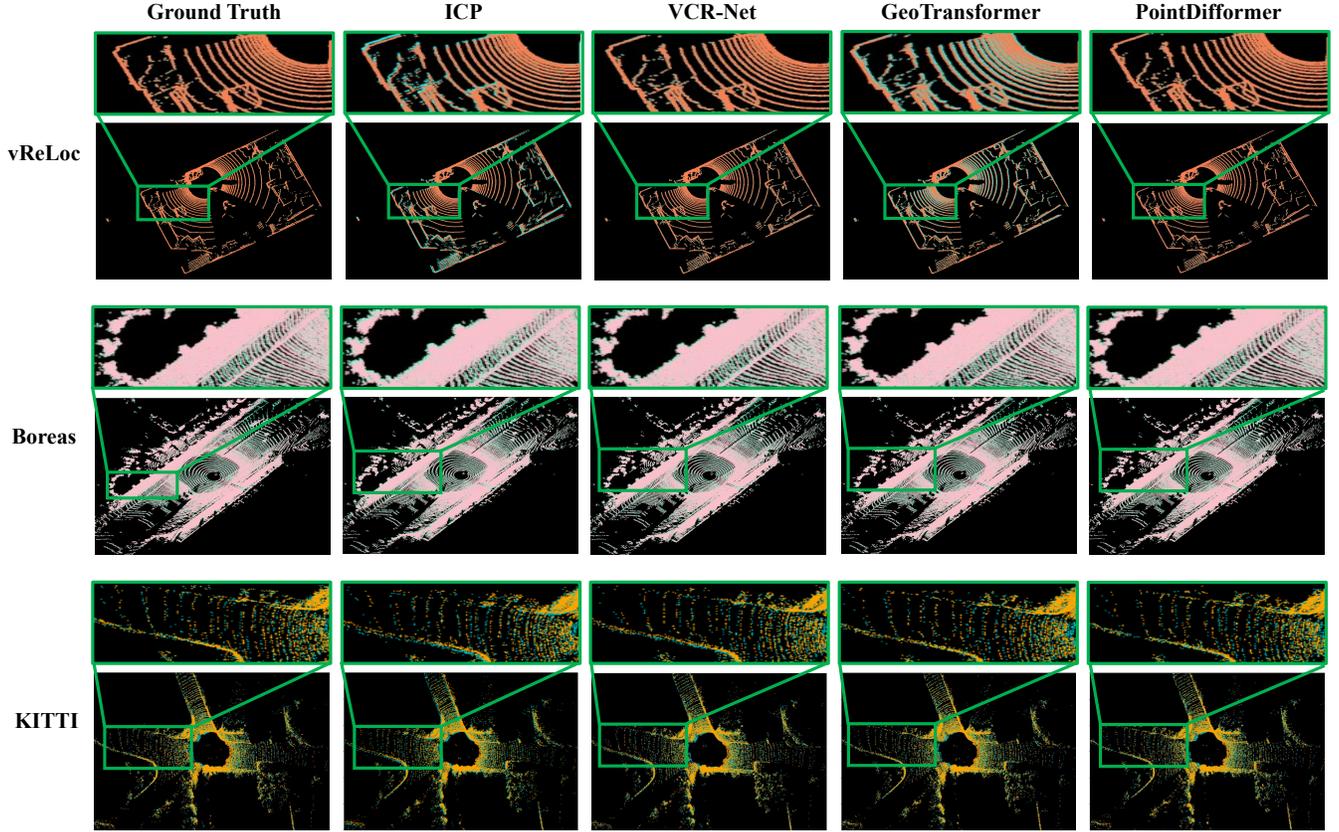
Fig. 4. Examples of point cloud frame pairs transformed using different prediction methods to align the second frame with the coordinate system of the first frame.

are calculated as the absolute values of the difference between the corresponding predictions and the ground truth. Using the predicted transformation between two point cloud frames, we can transform the second frame into the coordinate system of the first frame to achieve alignment of the point clouds. We show several examples of point cloud alignment based on the predicted transformation in Fig. 4, where the degree of overlap between the two frames increases with the accuracy of the predicted transformation.

TABLE I
PERFORMANCE OF POINT CLOUD REGISTRATION PREDICTION ON THE
vReLoc DATASET. THE BEST AND SECOND-BEST RESULTS UNDER
DIFFERENT METRICS ARE HIGHLIGHTED IN **BOLD** AND <u>UNDERLINED</u>,
RESPECTIVELY.

| Method | Relative Translation Error (centimeter [cm]) | | Relative Rotation Error (degree [°]) | | RR (%) |
|---|---|---|---|---|---|
| | *MAE* | *RMSE* | *MAE* | *RMSE* | |
| ICP | 2.20 | 12.69 | 0.20 | 1.35 | 96.2 |
| DCP | 1.35 | 3.26 | 0.45 | 1.31 | 85.4 |
| HGNN++ | 3.33 | 11.29 | 0.36 | 1.57 | 83.5 |
| VCR-Net | <u>0.25</u> | 0.45 | <u>0.04</u> | <u>0.11</u> | **99.9** |
| PCT++ | <u>0.25</u> | <u>0.44</u> | 0.05 | 0.12 | **99.9** |
| GeoTransformer | 0.66 | 1.10 | 0.07 | 0.16 | **99.9** |
| PointDifformer | **0.14** | **0.40** | **0.03** | **0.10** | **99.9** |

*2) Evaluation on Outdoor Boreas Dataset:* We compare PointDifformer with other baselines on the outdoor Boreas dataset, where the training dataset is collected under sunny
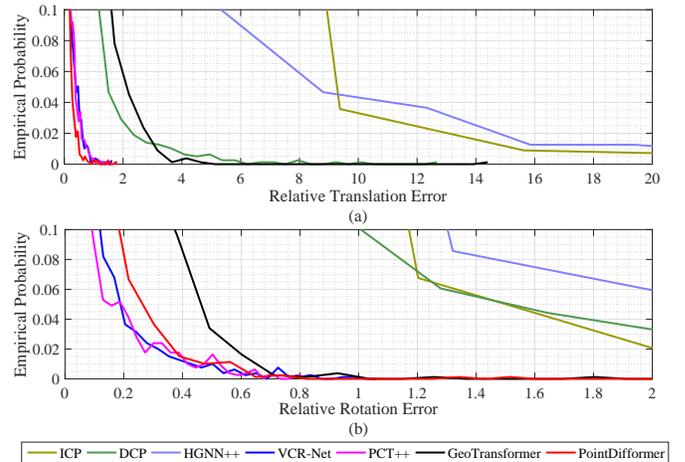


Fig. 5. The empirical probability of relative translation (centimeter [cm]) and rotation (degree [°]) errors on the vReLoc dataset.

weather, and the test dataset is collected under night weather. As shown in Table II, PointDifformer outperforms the other baselines under all criteria, except for the relative translation MAE, for which GeoTransformer is slightly better. However, Fig. 6 shows that GeoTransformer has longer probability tails in the translation and rotation errors than PointDifformer, indicating that our method is more robust. We also present

several examples of point cloud alignment using the predicted transformations in Fig. 4.
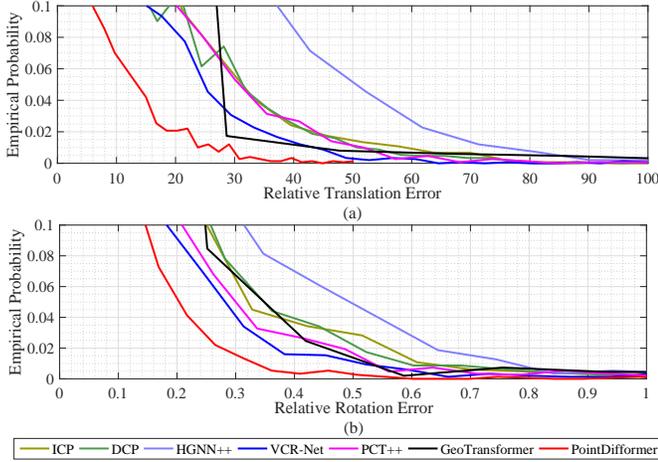


Fig. 6. The empirical probability of relative translation (centimeter [cm]) and rotation (degree [°]) errors on the Boreas dataset.



Fig. 7. The empirical probability of relative translation (centimeter [cm]) and rotation (degree [°]) errors on the KITTI dataset.

TABLE II
POINT CLOUD REGISTRATION PERFORMANCE ON THE BOREAS DATASET.

| Method | Relative Translation Error (centimeter [cm]) | | Relative Rotation Error (degree [°]) | | RR (%) |
|---|---|---|---|---|---|
| | *MAE* | *RMSE* | *MAE* | *RMSE* | |
| ICP | 10.83 | 18.28 | 0.11 | 0.21 | 75.4 |
| DCP | 11.63 | 17.36 | 0.12 | 0.21 | 70.5 |
| HGNN++ | 14.41 | 23.16 | 0.14 | 0.25 | 56.1 |
| VCR-Net | 8.71 | 13.56 | 0.10 | 0.17 | 84.7 |
| PCT++ | 9.81 | 15.77 | 0.10 | 0.19 | 79.6 |
| GeoTransformer | **4.58** | 15.78 | **0.08** | 0.22 | 94.9 |
| PointDifformer | 6.12 | **8.84** | **0.07** | **0.12** | **96.1** |

TABLE III
PERFORMANCE OF POINT CLOUD REGISTRATION PREDICTION ON THE KITTI DATASET.

| Method | Relative Translation Error (centimeter [cm]) | | Relative Rotation Error (degree [°]) | | RR (%) |
|---|---|---|---|---|---|
| | *MAE* | *RMSE* | *MAE* | *RMSE* | |
| ICP | 9.86 | 19.48 | 0.17 | 0.27 | 87.9 |
| DCP | 9.28 | 15.34 | 0.26 | 0.49 | 95.0 |
| HGNN++ | 8.86 | 17.20 | 0.20 | 0.31 | 89.9 |
| VCR-Net | 5.31 | 11.07 | 0.16 | 0.24 | 97.3 |
| PCT++ | 6.16 | 13.96 | 0.18 | 0.28 | 95.4 |
| GeoTransformer | **3.93** | 13.50 | 0.18 | 0.50 | **97.8** |
| PointDifformer | 4.14 | **8.86** | **0.14** | **0.23** | 97.7 |

TABLE IV
PERFORMANCE OF POINT CLOUD REGISTRATION PREDICTION ON THE KITTI DATASET WITH SEQUENCE "0" TO "8" FOR TRAINING AND SEQUENCE "9" TO "10" FOR THE TEST.

| Method | Relative Translation Error (centimeter [cm]) | | Relative Rotation Error (degree [°]) | | RR (%) |
|---|---|---|---|---|---|
| | *MAE* | *RMSE* | *MAE* | *RMSE* | |
| VCR-Net | 7.17 | 10.88 | 0.21 | 0.36 | 97.2 |
| GeoTransformer | 3.45 | 12.91 | 0.14 | 0.76 | **99.0** |
| PointDifformer | **3.10** | **5.93** | **0.11** | **0.17** | **99.0** |

*3) Evaluation on Outdoor KITTI Dataset:* We conduct point cloud registration methods on the KITTI dataset, selecting around 1600 and 1200 point cloud pairs for training and testing, respectively. From Table III and Fig. 7, we observe that PointDifformer demonstrates superior performance compared to other baselines on the KITTI dataset, with shorter tails of relative rotation and translation error probabilities. This is similar to its performance on the Boreas dataset. We also present several examples of our results in Fig. 4. Furthermore, we train on sequences "0" to "8" and test on sequences "9" to "10". We observe that PointDifformer surpasses other baselines when the size of the training data is larger, as shown in Table IV. The LiDAR point clouds in the KITTI dataset have practical noise due to dynamic objects and complex environments. However, the graph neural PDE modules in PointDifformer exhibit robustness to input perturbations, as demonstrated in [81]. This may be the reason why PointDifformer achieves more accurate predicted results when there is more practical noise in larger-sized data.

*4) Generalization from KITTI Dataset to Boreas dataset:* To cross-validate, we conduct point cloud registration by training the models on the KITTI dataset and evaluating them on the Boreas dataset. Specifically, we train the models on sequence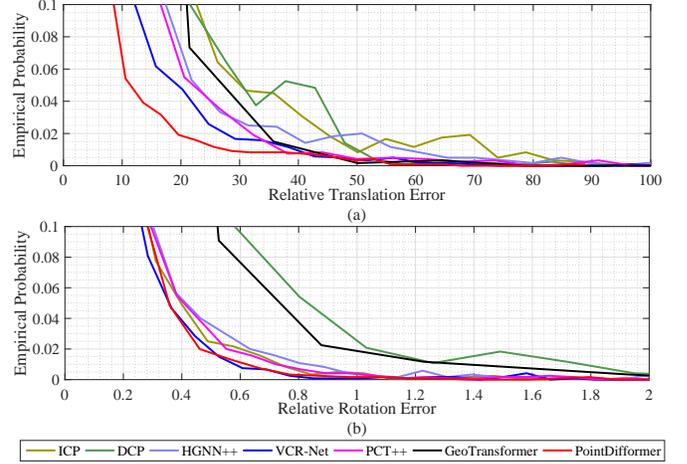 "9" of the KITTI dataset and test them on the sequence "night" of the Boreas dataset. From Table V, we observe that PointDifformer has competitive performance compared with the current state-of-the-art.

TABLE V
THE PERFORMANCE OF POINT CLOUD REGISTRATION ON THE BOREAS DATASET FOR TESTING (USING THE MODEL PRE-TRAINED ON THE KITTI DATASET).

| Method | Relative Translation Error (centimeter [cm]) | | Relative Rotation Error (degree [°]) | | RR (%) |
|---|---|---|---|---|---|
| | *MAE* | *RMSE* | *MAE* | *RMSE* | |
| HGNN++ | 16.06 | 25.86 | 0.15 | 0.27 | 49.8 |
| VCR-Net | 11.97 | 19.78 | 0.11 | 0.19 | 68.6 |
| PCT++ | 11.61 | 19.57 | 0.13 | 0.31 | 72.4 |
| GeoTransformer | **5.97** | 27.90 | 0.09 | 0.33 | 93.1 |
| PointDifformer | 6.63 | **10.07** | **0.08** | **0.14** | **93.3** |

## D. Robutness Evaluation

*1) Robustness against synthetic noise on the KITTI Dataset:*
We investigate the robustness of PointDifformer under various types of synthetic noise, as discussed below.

TABLE VI
POINT CLOUD REGISTRATION PERFORMANCE ON THE KITTI DATASET
WITH THE GAUSSIAN NOISE THAT FOLLOWS $\mathcal{N}(0, \sigma = 0.25)$.

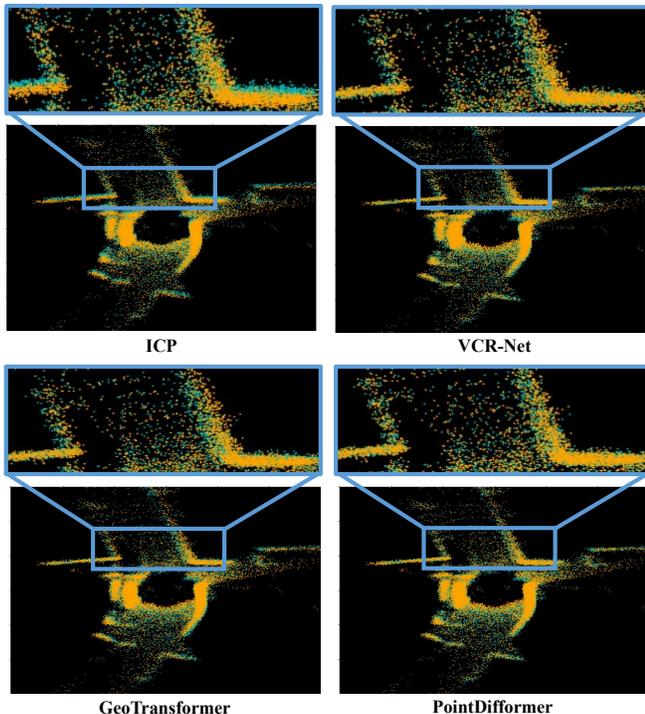| Method | Relative Translation Error (centimeter [cm]) | | Relative Rotation Error (degree [°]) | | RR (%) |
|---|---|---|---|---|---|
| | *MAE* | *RMSE* | *MAE* | *RMSE* | |
| ICP | 14.97 | 26.09 | 0.20 | 0.32 | 69.3 |
| DCP | 9.97 | 15.84 | 0.29 | 0.52 | 94.3 |
| HGNN++ | 10.62 | 18.76 | 0.22 | 0.34 | 88.9 |
| VCR-Net | 6.40 | 12.40 | 0.18 | 0.27 | 96.3 |
| PCT++ | 6.85 | 14.03 | 0.20 | 0.30 | 95.3 |
| GeoTransformer | 5.37 | 14.43 | 0.25 | 0.50 | 97.5 |
| PointDifformer | **5.23** | **9.00** | **0.17** | **0.25** | **97.7** |



Fig. 8. Examples for the point cloud frame pairs from the KITTI dataset with Gaussian noise using different transformation prediction methods for alignment.

TABLE VII
THE MAE OF THE PREDICTED RELATIVE TRANSLATION/ROTATION
(CENTIMETER [CM] / DEGREE [°]) ON DIFFERENT NOISE POWERS.

| Metric | Methods | $\mathcal{N}(0, \sigma = 0.5)$ | $\mathcal{N}(0, \sigma = 0.75)$ | $\mathcal{N}(0, \sigma = 1.0)$ |
|---|---|---|---|---|
| | VCR-Net | 7.97 / **0.23** | 9.70 / 0.29 | 11.62 / 0.36 |
| MAE | GeoTransformer | 7.95 / 0.34 | 10.04 / 0.44 | 13.27 / 0.54 |
| | PointDifformer | **7.18** / **0.23** | **8.83** / **0.27** | **10.38** / **0.31** |
| | VCR-Net | 14.08 / 0.35 | 15.92 / 0.44 | 17.95 / 0.55 |
| RMSE | GeoTransformer | 14.40 / 0.68 | **15.00** / 0.78 | 20.15 / 0.90 |
| | PointDifformer | **12.90** / **0.34** | 15.08 / **0.40** | **17.31** / **0.45** |

**Performance on the KITTI dataset with Gaussian noise.**
To evaluate the robustness of PointDifformer, we add white

TABLE VIII
POINT CLOUD REGISTRATION PERFORMANCE ON THE KITTI DATASET
WITH 3D SHAPE PERTURBATIONS.

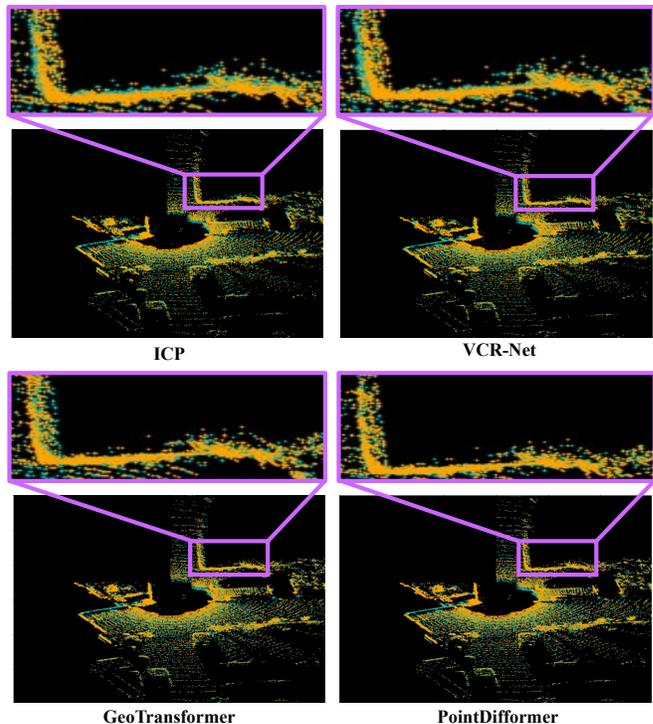| Method | Relative Translation Error (centimeter [cm]) | | Relative Rotation Error (degree [°]) | | RR (%) |
|---|---|---|---|---|---|
| | *MAE* | *RMSE* | *MAE* | *RMSE* | |
| ICP | 12.60 | 24.92 | 0.18 | 0.32 | 80.4 |
| DCP | 12.63 | 23.04 | 0.28 | 0.48 | 85.8 |
| HGNN++ | 11.39 | 22.27 | 0.21 | 0.36 | 85.1 |
| VCR-Net | 6.39 | 13.69 | 0.18 | 0.36 | 95.4 |
| PCT++ | 7.23 | 16.23 | 0.21 | 0.36 | 93.6 |
| GeoTransformer | **3.89** | 13.08 | 0.18 | 0.44 | **97.8** |
| PointDifformer | 4.14 | **8.99** | **0.14** | **0.22** | **97.8** |



Fig. 9. Examples for noisy point cloud frames with 3D shape perturbations using different transformation prediction methods for alignment.

Gaussian noise $\mathcal{N}(0, \sigma)$ to the original KITTI dataset, similar to the experiments for Table III. Based on the results presented in Table VI, we observe that PointDifformer surpasses the other benchmark methods in terms of relative translation and rotation errors. Additionally, we present examples of point cloud alignment using the predicted transformation in Fig. 8. We also evaluate the robustness of our method and other baselines to different noise powers. As shown in Table VII, PointDifformer demonstrates superior robustness w.r.t. additive Gaussian noise compared to the other baselines across different noise powers.

**Performance on the noisy KITTI dataset with 3D shape perturbations.** We next introduce 3D shape perturbations to the original KITTI dataset. This is achieved by removing certain parts of the original point clouds, resulting in an imperfect 3D shape. Specifically, we remove a 25 m ×15 m region in the lower left corner of each point cloud frame measuring 60 m ×30 m. The results presented in Table VIII show that

PointDifformer still outperforms the baselines in terms of both criteria for relative rotation error and the RMSE for relative translation prediction. Some examples are presented in Fig. 9.

*2) Robustness against natural noise on the Boreas Dataset:* We conduct experiments on the Boreas dataset under rainy conditions, which is considered natural noise on point clouds. From Table IX, we observe that PointDifformer outperforms the baseline methods in terms of relative translation and rotation RMSEs, while having comparable performance in terms of MAEs. This indicates that our method produces fewer outliers in the predicted results, demonstrating its superior robustness compared to the baselines. We also provide several examples of point cloud alignment in Fig. 10.

TABLE IX
POINT CLOUD REGISTRATION PERFORMANCE ON THE BOREAS DATASET UNDER THE RAINING ENVIRONMENT.

| Method | Relative Translation Error (centimeter [cm]) | | Relative Rotation Error (degree [°]) | | RR (%) |
|---|---|---|---|---|---|
| | *MAE* | *RMSE* | *MAE* | *RMSE* | |
| ICP | 11.90 | 20.57 | 0.15 | 0.27 | 70.8 |
| DCP | 10.60 | 16.00 | 0.14 | 0.22 | 75.7 |
| HGNN++ | 15.02 | 25.63 | 0.18 | 0.32 | 52.5 |
| VCR-Net | 8.81 | 14.09 | 0.13 | 0.20 | 84.1 |
| PCT++ | 10.39 | 16.86 | 0.14 | 0.24 | 77.4 |
| GeoTransformer | **4.96** | 16.75 | **0.10** | 0.25 | 94.9 |
| PointDifformer | 5.91 | **8.45** | **0.10** | **0.14** | **96.9** |



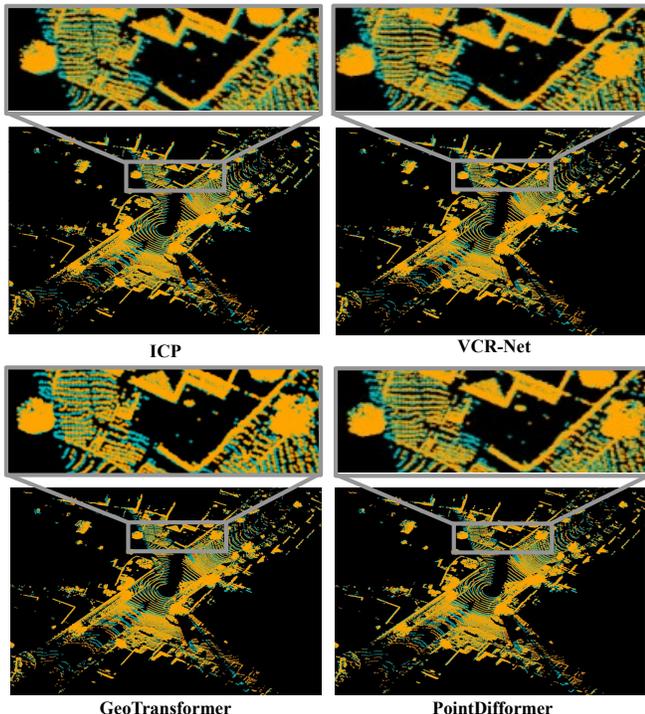**ICP**      **VCR-Net**

**GeoTransformer**      **PointDifformer**

Fig. 10. Examples for the point cloud frame pairs from the Boreas dataset under a raining environment using different transformation prediction methods for alignment.

### E. Evaluation on Datasets with Lower Overlaps

We evaluate the performance of our method on datasets with lower overlaps, namely, 3DMatch and 3DLoMatch, as shown in Table X. The standard benchmark metric RR is used, and the experimental configuration follows that of [42], [54]. The baselines include FCGF [17], D3Feat [18], SpinNet [19], Predator [35], YOHO [44], CoFiNet [20], GeoTransformer [42], RoITr [54], and RoReg [56]. To perform the registration task on the 3DMatch and 3DLoMatch datasets, similar to several current state-of-the-art models like CoFiNet, GeoTransformer, RoITr and RoReg, we first employ a module to identify the areas of higher overlap in the point clouds. Specifically, we use superpoint matching [42]. Subsequently, keypoint matching under the PointDifformer model is employed to achieve more precise correspondences for pairs of point clouds. Table X demonstrates that PointDifformer achieves state-of-the-art performance, affirming its feasibility on datasets with lower overlaps.

TABLE X
PERFORMANCE ON THE 3DMATCH AND 3DLOMATCH DATASETS, USING THE SAME EXPERIMENTAL CONFIGURATION AS THAT IN [42], [54]. THE RESULTS OF BASELINES ARE BORROWED FROM [17]–[20], [35], [42], [44], [54], [56].

| Method | 3DMatch RR (%) | 3DLoMatch RR (%) |
|---|---|---|
| FCGF | 85.1 | 40.1 |
| D3Feat | 81.6 | 37.2 |
| SpinNet | 88.6 | 59.8 |
| Predator | 89.0 | 59.8 |
| YOHO | 90.8 | 65.2 |
| CoFiNet | 89.3 | 67.5 |
| GeoTransformer | 92.0 | 75.0 |
| RoITr | 91.9 | 74.8 |
| RoReg | 92.9 | 70.3 |
| PointDifformer | **93.0** | **75.2** |

### F. Computational Complexity

In Table XI, we present the average inference time and graphics processing unit (GPU) memory required for registering each point cloud pair based on the KITTI dataset. We test the methods on an NVIDIA RTX A5000 GPU. The average inference time and GPU memory are measured in seconds (s) and gigabytes (GB), respectively. From Table XI, we observe that PointDifformer requires higher GPU memory and incurs longer inference time compared to other baselines due to its higher complexity. On average, the inference time is still acceptable for real-time applications. A possible future work is to optimize and prune [93]–[95] the PointDifformer model to reduce its memory and inference time footprints.

TABLE XI
THE AVERAGE INFERENCE TIME AND GPU MEMORY FOR EACH POINT CLOUD PAIR ON THE KITTI DATASET.

| Method | VCR-Net | PCT++ | GeoTransformer | PointDifformer |
|---|---|---|---|---|
| Inference time | 0.047s | 0.648s | 0.061s | 0.072s |
| GPU memory | 2.29GB | 2.38GB | 1.51GB | 2.44GB |

### G. Ablation Study

We perform an ablation study using the KITTI dataset under the same experimental settings as described in Section IV-C3

for Table III. We first evaluate the efficiency of the self-attention module with heat kernel signature by comparing it with the vanilla self-attention module and the module without self-attention. As shown in Table XII, the introduction of the heat kernel signature as weights into the self-attention module improves the transformation prediction accuracy. Furthermore, we observe that the vanilla self-attention module also contributes to the point cloud registration performance.

TABLE XII
ABLATION STUDY FOR THE SELF-ATTENTION MODULE WITH HEAT KERNEL
SIGNATURE.

| Module | Relative Translation Error (centimeter [cm]) | | Relative Rotation Error (degree [°]) | | RR (%) |
|---|---|---|---|---|---|
| | MAE | RMSE | MAE | RMSE | |
| Self Attention (No) | 5.75 | 11.73 | 0.16 | 0.25 | 92.1 |
| Self Attention (Vanilla) | 4.56 | 10.03 | 0.15 | 0.26 | 94.5 |
| Self Attention (Heat Kernel) | 4.14 | 8.86 | 0.14 | 0.23 | 97.7 |

We investigate the influence of our proposed Point-Diffusion Net module on point cloud representation by comparing it with other graph learning methods. The results in Table XIII show that the point cloud registration model with our Point-Diffusion Net outperforms those with other graph learning methods including HGNN [92] and DGCNN [65]. This indicates that the Point-Diffusion Net can achieve a more robust representation of the point cloud.

TABLE XIII
ABLATION STUDY FOR THE EFFECTIVENESS OF POINT-DIFFUSION NET.

| Method | Relative Translation Error (centimeter [cm]) | | Relative Rotation Error (degree [°]) | | RR (%) |
|---|---|---|---|---|---|
| | MAE | RMSE | MAE | RMSE | |
| HGNN | 8.22 | 16.40 | 0.19 | 0.31 | 90.1 |
| DGCNN | 4.46 | 9.55 | 0.14 | 0.23 | 97.3 |
| Point-Diffusion Net | 4.14 | 8.86 | 0.14 | 0.23 | 97.7 |

We investigate the impact of the number of selected top $K'$ corresponding keypoints from the attention-based correspondence module. Specifically, we select 25%, 50%, 75%, and 100% corresponding keypoints from the entire set of keypoints. The results in Table XIV show that selecting 50% and 75% corresponding keypoints achieves better performance.

TABLE XIV
ABLATION STUDY FOR THE NUMBER OF SELECTED CORRESPONDING
KEYPOINTS IN THE ATTENTION-BASED CORRESPONDENCE.

| Proportion | Relative Translation Error (centimeter [cm]) | | Relative Rotation Error (degree [°]) | | RR (%) |
|---|---|---|---|---|---|
| | MAE | RMSE | MAE | RMSE | |
| 25% | 8.88 | 15.84 | 0.20 | 0.32 | 88.1 |
| 50% | 4.06 | 7.95 | 0.15 | 0.25 | 97.3 |
| 75% | 4.14 | 8.86 | 0.14 | 0.23 | 97.7 |
| 100% | 7.20 | 16.14 | 0.19 | 0.34 | 92.2 |

To evaluate the effectiveness of our total loss $\mathcal{L}_{\text{total}}$ in (28), we compare it with the corresponding point loss $\mathcal{L}_{\text{point}}$ in (26) and the ground-truth-to-prediction loss $\mathcal{L}_{\text{rt}}$ in (27). Based on the results presented in Table XV, we observe that the total loss $\mathcal{L}_{\text{total}}$ outperforms the others, suggesting that incorporating more information in the loss function has a positive effect on training. Additionally, we find that $\mathcal{L}_{\text{point}}$ surpasses $\mathcal{L}_{\text{rt}}$, which implies that the point loss plays a more important role in the $\mathcal{L}_{\text{total}}$.

TABLE XV
ABLATION STUDY FOR THE LOSS FUNCTION.

| Loss | Relative Translation Error (centimeter [cm]) | | Relative Rotation Error (degree [°]) | | RR (%) |
|---|---|---|---|---|---|
| | MAE | RMSE | MAE | RMSE | |
| $\mathcal{L}_{\text{rt}}$ | 6.32 | 11.55 | 0.21 | 0.39 | 94.2 |
| $\mathcal{L}_{\text{point}}$ | 5.41 | 11.44 | 0.15 | 0.24 | 96.3 |
| $\mathcal{L}_{\text{total}}$ | 4.14 | 8.86 | 0.14 | 0.23 | 97.7 |

## V. CONCLUSIONS

In order to develop a robust 3D point cloud registration approach that is able to handle noise or perturbations, we have utilized graph neural PDE modules to learn point cloud feature representations. We have also designed attention modules with heat kernel signatures to establish correspondence between points from two point clouds. Our approach has been extensively evaluated through experiments, which demonstrate that it generally outperforms baselines not only on raw point clouds but also on point clouds with additive noise and 3D shape perturbations. These results suggest that graph neural PDEs are beneficial for the task of point cloud registration.

In this paper, we have conducted a robustness study limited only to perturbations through Gaussian noise, rain, and partial removal of a frame. A future work of interest is to further investigate the robustness of our method under diverse perturbations, including *adversarial attacks*. Furthermore, we aim to enhance our model's adaptability to noisy datasets under different environmental factors.

## ACKNOWLEDGMENT

## REFERENCES

[1] L. Zhang, J. Guo, Z. Cheng, J. Xiao, and X. Zhang, "Efficient pairwise 3-D registration of urban scenes via hybrid structural descriptors," *IEEE Trans. Geosci. Remote. Sens.*, vol. 60, pp. 1–17, 2022.

[2] S. Quan and J. Yang, "Compatibility-guided sampling consensus for 3-D point cloud registration," *IEEE Trans. Geosci. Remote. Sens.*, vol. 58, no. 10, pp. 7380–7392, 2020.

[3] Q. Kang, R. She, S. Wang, W. Tay, D. Navarro, and A. Hartmannsgruber, "Location learning for AVs: LiDAR and image landmarks fusion localization with graph neural networks," in *Proc. IEEE Int. Conf. Intell. Transp. Syst.*, 2022, pp. 3032–3037.

[4] L. Sun, Z. Zhang, R. Zhong, D. Chen, L. Zhang, L. Zhu, Q. Wang, G. Wang, J. Zou, and Y. Wang, "A weakly supervised graph deep learning framework for point cloud registration," *IEEE Trans. Geosci. Remote. Sens.*, vol. 60, pp. 1–12, 2022.

[5] P. Zhou, X. Guo, X. Pei, and C. Chen, "T-LOAM: Truncated least squares LiDAR-only odometry and mapping in real time," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–13, 2021.

[6] W. Shi and R. Rajkumar, "Point-GNN: Graph neural network for 3D object detection in a point cloud," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2020, pp. 1711–1719.

[7] T. Shan, B. Englot, D. Meyers, W. Wang, C. Ratti, and D. Rus, "LIO-SAM: Tightly-coupled LiDAR inertial odometry via smoothing and mapping," in *Proc. IEEE Int. Conf. Intell. Robot. Syst.*, 2020, pp. 5135–5142.

[8] J. Zhang and S. Singh, "LOAM: LiDAR odometry and mapping in real-time," in *Proc. Robot. Sci. Syst.*, 2014, pp. 1–9.

[9] Y. Wang and J. M. Solomon, "Deep closest point: Learning representations for point cloud registration," in *Proc. IEEE Int. Conf. Comput. Vision*, 2019, pp. 3523–3532.

[10] P. J. Besl and N. D. McKay, "Method for registration of 3-D shapes," in *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 14, no. 2, 1992, pp. 239–256.

[11] J. Yang, H. Li, and Y. Jia, "Go-ICP: Solving 3D registration efficiently and globally optimally," in *Proc. IEEE Int. Conf. Comput. Vision*, 2013, pp. 1457–1464.

[12] K. Koide, M. Yokozuka, S. Oishi, and A. Banno, "Voxelized GICP for fast and accurate 3D point cloud registration," in *Proc. IEEE Int. Conf. Robot. Autom.*, 2021, pp. 11 054–11 059.

[13] H. Wei, Z. Qiao, Z. Liu, C. Suo, P. Yin, Y. Shen, H. Li, and H. Wang, "End-to-end 3D point cloud learning for registration task using virtual correspondences," in *Proc. IEEE Int. Conf. Intell. Robot. Syst.*, 2020, pp. 2678–2683.

[14] J. Xu, Y. Huang, Z. Wan, and J. Wei, "GLORN: Strong generalization fully convolutional network for low-overlap point cloud registration," *IEEE Trans. Geosci. Remote. Sens.*, vol. 60, pp. 1–14, 2022.

[15] Y. Zhang, J. Xu, Y. Zou, P. X. Liu, and J. Liu, "PS-Net: Point shift network for 3-D point cloud completion," *IEEE Trans. Geosci. Remote. Sens.*, vol. 60, pp. 1–13, 2022.

[16] F. Lu, G. Chen, Y. Liu, L. Zhang, S. Qu, S. Liu, and R. Gu, "HRegNet: A hierarchical network for large-scale outdoor LiDAR point cloud registration," in *Proc. IEEE Int. Conf. Comput. Vision*, 2021, pp. 16 014–16 023.

[17] C. Choy, J. Park, and V. Koltun, "Fully convolutional geometric features," in *Proc. IEEE Int. Conf. Comput. Vision*, 2019, pp. 8958–8966.

[18] 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," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2020, pp. 6359–6367.

[19] S. Ao, Q. Hu, B. Yang, A. Markham, and Y. Guo, "SpinNet: Learning a general surface descriptor for 3D point cloud registration," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2021, pp. 11 753–11 762.

[20] H. Yu, F. Li, M. Saleh, B. Busam, and S. Ilic, "CoFiNet: Reliable coarse-to-fine correspondences for robust pointcloud registration," *Adv. Neural Inform. Process. Syst.*, vol. 34, pp. 23 872–23 884, 2021.

[21] J. Li, Q. Hu, and M. Ai, "Point cloud registration based on one-point ransac and scale-annealing biweight estimation," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 11, pp. 9716–9729, 2021.

[22] F. Wang, H. Hu, X. Ge, B. Xu, R. Zhong, Y. Ding, X. Xie, and Q. Zhu, "Multientity registration of point clouds for dynamic objects on complex floating platform using object silhouettes," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 1, pp. 769–783, 2020.

[23] S. Chen, L. Nan, R. Xia, J. Zhao, and P. Wonka, "PLADE: A plane-based descriptor for point cloud registration with small overlap," *IEEE Trans. Geosci. Remote. Sens.*, vol. 58, no. 4, pp. 2530–2540, 2019.

[24] J. Yu, Y. Lin, B. Wang, Q. Ye, and J. Cai, "An advanced outlier detected total least-squares algorithm for 3-D point clouds registration," *IEEE Trans. Geosci. Remote. Sens.*, vol. 57, no. 7, pp. 4789–4798, 2019.

[25] Y. Zhao, Y. Li, H. Zhang, V. Monga, and Y. C. Eldar, "A convergent neural network for non-blind image deblurring," in *Proc. IEEE Int. Conf. Image Process.*, 2023, pp. 1505–1509.

[26] Y. Song, Q. Kang, S. Wang, Z. Kai, and W. P. Tay, "On the robustness of graph neural diffusion to topology perturbations," in *Adv. Neural Inform. Process. Syst.*, vol. 35, 2022, pp. 6384–6396.

[27] J. Sun, M. Ovsjanikov, and L. Guibas, "A concise and provably informative multi-scale signature based on heat diffusion," *Comput. Graphics Forum*, vol. 28, no. 5, pp. 1383–1392, 2009.

[28] M. A. Fischler and R. C. Bolles, "Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography," *ACM Commun.*, vol. 24, no. 6, pp. 381–395, 1981.

[29] H. Maron, N. Dym, I. Kezurer, S. Kovalsky, and Y. Lipman, "Point registration via efficient convex relaxation," *ACM Trans. Graphics*, vol. 35, no. 4, pp. 1–12, 2016.

[30] T. Guérout, Y. Gaoua, C. Artigues, G. Da Costa, P. Lopez, and T. Monteil, "Mixed integer linear programming for quality of service optimization in clouds," *Future Gener. Comput. Syst.*, vol. 71, pp. 1–17, 2017.

[31] A. Segal, D. Haehnel, and S. Thrun, "Generalized-ICP," in *Proc. Robot. Sci. Syst.*, vol. 2, no. 4, 2009, p. 435.

[32] H. Deng, T. Birdal, and S. Ilic, "PPF-FoldNet: Unsupervised learning of rotation invariant 3D local descriptors," in *Proc. Eur. Conf. Comput. Vision*, 2018, pp. 602–618.

[33] H. Deng, T. Birdal, and et al., "PPFNet: Global context aware local features for robust 3D point matching," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2018, pp. 195–205.

[34] Z. Gojcic, C. Zhou, J. D. Wegner, and A. Wieser, "The perfect match: 3D point cloud matching with smoothed densities," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2019, pp. 5545–5554.

[35] S. Huang, Z. Gojcic, M. Usvyatsov, A. Wieser, and K. Schindler, "Predator: Registration of 3D point clouds with low overlap," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2021, pp. 4267–4276.

[36] Z. Chen, K. Sun, F. Yang, and W. Tao, "SC$^2$-PCR: A second order spatial compatibility for efficient and robust point cloud registration," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2022, pp. 13 221–13 231.

[37] X. Zhang, J. Yang, S. Zhang, and Y. Zhang, "3D registration with maximal cliques," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2023, pp. 17 745–17 754.

[38] W. Lu, Y. Zhou, G. Wan, S. Hou, and S. Song, "L$^3$-Net: Towards learning based LiDAR localization for autonomous driving," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2019, pp. 6389–6398.

[39] W. Lu, G. Wan, Y. Zhou, X. Fu, P. Yuan, and S. Song, "DeepVCP: An end-to-end deep neural network for point cloud registration," in *Proc. IEEE Int. Conf. Comput. Vision*, 2019, pp. 12–21.

[40] Y. Aoki, H. Goforth, R. A. Srivatsan, and S. Lucey, "PointNetLK: Robust & efficient point cloud registration using pointnet," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2019, pp. 7163–7172.

[41] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "PointNet++: Deep hierarchical feature learning on point sets in a metric space," *Adv. Neural Inform. Process. Syst.*, vol. 30, pp. 1–10, 2017.

[42] Z. Qin, H. Yu, C. Wang, Y. Guo, Y. Peng, and K. Xu, "Geometric transformer for fast and robust point cloud registration," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2022, pp. 11 143–11 152.

[43] C. Choy, W. Dong, and V. Koltun, "Deep global registration," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2020, pp. 2514–2523.

[44] H. Wang, Y. Liu, Z. Dong, and W. Wang, "You only hypothesize once: Point cloud registration with rotation-equivariant descriptors," in *Proc. ACM Int. Conf. Multimedias*, 2022, pp. 1630–1641.

[45] M. Zhao, L. Ma, X. Jia, D.-M. Yan, and T. Huang, "GraphReg: Dynamical point cloud registration with geometry-aware graph signal processing," *IEEE Trans. Image Process.*, vol. 31, pp. 7449–7464, 2022.

[46] F. Poiesi and D. Boscaini, "Learning general and distinctive 3D local deep descriptors for point cloud registration," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 45, no. 3, pp. 3979–3985, 2023.

[47] W. Tang and D. Zou, "Multi-instance point cloud registration by efficient correspondence clustering," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2022, pp. 6667–6676.

[48] J. Lee, S. Kim, M. Cho, and J. Park, "Deep hough voting for robust global registration," in *Proc. IEEE Int. Conf. Comput. Vision*, 2021, pp. 15 994–16 003.

[49] G. D. Pais, S. Ramalingam, V. M. Govindu, J. C. Nascimento, R. Chellappa, and P. Miraldo, "3DRegNet: A deep neural network for 3D point registration," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2020, pp. 7193–7203.

[50] K. Fischer, M. Simon, F. Olsner, S. Milz, H.-M. Gross, and P. Mader, "StickyPillars: Robust and efficient feature matching on point clouds using graph neural networks," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2021, pp. 313–323.

[51] X. Bai, Z. Luo, L. Zhou, H. Chen, L. Li, Z. Hu, H. Fu, and C.-L. Tai, "PointDSC: Robust point cloud registration using deep spatial consistency," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2021, pp. 15 859–15 869.

[52] Z. Zhang, Y. Dai, B. Fan, J. Sun, and M. He, "Learning a task-specific descriptor for robust matching of 3D point clouds," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 32, no. 12, pp. 8462–8475, 2022.

[53] Y. Li and T. Harada, "Lepard: Learning partial point cloud matching in rigid and deformable scenes," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2022, pp. 5554–5564.

[54] H. Yu, Z. Qin, J. Hou, M. Saleh, D. Li, B. Busam, and S. Ilic, "Rotation-invariant transformer for point cloud matching," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2023, pp. 5384–5393.

[55] S. Ao, Q. Hu, H. Wang, K. Xu, and Y. Guo, "BUFFER: Balancing accuracy, efficiency, and generalizability in point cloud registration," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2023, pp. 1255–1264.

[56] H. Wang, Y. Liu, Q. Hu, B. Wang, J. Chen, Z. Dong, Y. Guo, W. Wang, and B. Yang, "RoReg: Pairwise point cloud registration with oriented descriptors and local rotations," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 45, no. 8, pp. 10 376–10 393, 2023.

[57] J. Yu, L. Ren, Y. Zhang, W. Zhou, L. Lin, and G. Dai, "PEAL: Prior-embedded explicit attention learning for low-overlap point cloud registration," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2023, pp. 17 702–17 711.

[58] H. Jiang, Z. Dang, Z. Wei, J. Xie, J. Yang, and M. Salzmann, "Robust outlier rejection for 3D registration with variational bayes," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2023, pp. 1148–1157.

[59] Y. Zhou and O. Tuzel, "VoxelNet: End-to-end learning for point cloud based 3D object detection," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2018, pp. 4490–4499.

[60] V. A. Sindagi, Y. Zhou, and O. Tuzel, "MVX-Net: Multimodal voxelnet for 3D object detection," in *Proc. IEEE Int. Conf. Robot. Autom.*, 2019, pp. 7276–7282.

[61] O. Kopuklu, N. Kose, A. Gunduz, and G. Rigoll, "Resource efficient 3D convolutional neural networks," in *Proc. IEEE Int. Conf. Comput. Vision Workshops*, 2019, pp. 1–10.

[62] S. Kumawat and S. Raman, "LP-3DCNN: Unveiling local phase in 3D convolutional neural networks," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2019, pp. 4903–4912.

[63] H. Su, S. Maji, E. Kalogerakis, and E. Learned-Miller, "Multi-view convolutional neural networks for 3D shape recognition," in *Proc. IEEE Int. Conf. Comput. Vision*, 2015, pp. 945–953.

[64] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "PointNet: Deep learning on point sets for 3D classification and segmentation," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2017, pp. 652–660.

[65] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic graph CNN for learning on point clouds," *ACM Trans. Graphics*, vol. 38, no. 5, pp. 1–12, 2019.

[66] Z. Liu, S. Zhou, C. Suo, P. Yin, W. Chen, H. Wang, H. Li, and Y.-H. Liu, "LPD-Net: 3D point cloud learning for large-scale place recognition and environment analysis," in *Proc. IEEE Int. Conf. Comput. Vision*, 2019, pp. 2831–2840.

[67] H. Thomas, C. R. Qi, J.-E. Deschaud, B. Marcotegui, F. Goulette, and L. J. Guibas, "KPConv: Flexible and deformable convolution for point clouds," in *Proc. IEEE Int. Conf. Comput. Vision*, 2019, pp. 6411–6420.

[68] Y. Rao, J. Lu, and J. Zhou, "Global-local bidirectional reasoning for unsupervised representation learning of 3D point clouds," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2020, pp. 5376–5385.

[69] F. Poiesi and D. Boscaini, "Distinctive 3D local deep descriptors," in *Proc. IEEE Int. Conf. Pattern Recognit.*, 2021, pp. 5720–5727.

[70] M.-H. Guo, J.-X. Cai, Z.-N. Liu, T.-J. Mu, R. R. Martin, and S.-M. Hu, "PCT: Point cloud transformer," *Comput. Vis. Media*, vol. 7, pp. 187–199, 2021.

[71] X. Ma, C. Qin, H. You, H. Ran, and Y. Fu, "Rethinking network design and local geometry in point cloud: A simple residual MLP framework," *arXiv preprint arXiv:2202.07123*, 2022.

[72] G. Qian, Y. Li, H. Peng, J. Mai, H. Hammoud, M. Elhoseiny, and B. Ghanem, "PointNeXt: Revisiting PointNet++ with improved training and scaling strategies," *Adv. Neural Inform. Process. Syst.*, vol. 35, pp. 23 192–23 204, 2022.

[73] R. T. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud, "Neural ordinary differential equations," *Adv. Neural Inform. Process. Syst.*, vol. 31, pp. 1–13, 2018.

[74] B. P. Chamberlain, J. Rowbottom, M. Goronova, S. Webb, E. Rossi, and M. M. Bronstein, "GRAND: Graph neural diffusion," in *Proc. Int. Conf. Mach. Learn.*, 2021, pp. 1407–1418.

[75] Z.-Q. Chen, Z. Qian, Y. Hu, and W. Zheng, "Stability and approximations of symmetric diffusion semigroups and kernels," *J. functional anal.*, vol. 152, no. 1, pp. 255–280, 1998.

[76] B. P. Chamberlain, J. Rowbottom, D. Eynard, F. Di Giovanni, D. Xiaowen, and M. M. Bronstein, "Beltrami flow and neural diffusion on graphs," *Adv. Neural Inform. Process. Syst.*, pp. 1–16, 2021.

[77] S. Wang, Q. Kang, R. She, W. P. Tay, A. Hartmannsgruber, and D. N. Navarro, "RobustLoc: Robust camera pose regression in challenging driving environments," in *Proc. AAAI Conf. Artificial Intell.*, 2022, pp. 1–8.

[78] Q. Kang, Y. Song, Q. Ding, and W. P. Tay, "Stable neural ODE with Lyapunov-stable equilibrium points for defending against adversarial attacks," *Adv. Neural Inform. Process. Syst.*, vol. 34, pp. 14 925–14 937, 2021.

[79] Y. Zhao, F. Dai, and J. Shi, "A dredge traffic algorithm for maintaining network stability," in *Proc. Int. Conf. Commun. Signal Process. Syst.*, 2020, pp. 1100–1108.

[80] K. Zhao, Q. Kang, Y. Song, R. She, S. Wang, and W. P. Tay, "Adversarial robustness in graph neural networks: A Hamiltonian approach," *arXiv preprint arXiv:2310.06396*, 2023.

[81] R. She, Q. Kang, S. Wang, Y.-R. Yang, K. Zhao, Y. Song, and W. P. Tay, "RobustMat: Neural diffusion for street landmark patch matching under challenging environments," *IEEE Trans. Image Process.*, vol. 32, pp. 5550–5563, 2023.

[82] Z. Chen, W. Zeng, Z. Yang, L. Yu, C.-W. Fu, and H. Qu, "LassoNet: Deep lasso-selection of 3D point clouds," *IEEE Trans. Vis. Comput. Graphics*, vol. 26, no. 1, pp. 195–204, 2019.

[83] M. M. Bronstein and I. Kokkinos, "Scale-invariant heat kernel signatures for non-rigid shape recognition," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2010, pp. 1704–1711.

[84] G. K. Tam, Z.-Q. Cheng, Y.-K. Lai, F. C. Langbein, Y. Liu, D. Marshall, R. R. Martin, X.-F. Sun, and P. L. Rosin, "Registration of 3D point clouds and meshes: A survey from rigid to nonrigid," *IEEE Trans. Vis. Comput. Graphics*, vol. 19, no. 7, pp. 1199–1217, 2012.

[85] R. T. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud, "Neural ordinary differential equations," *Adv. Neural Inform. Process. Syst.*, vol. 31, pp. 1–13, 2018.

[86] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Adv. Neural Inform. Process. Syst.*, vol. 30, pp. 1–11, 2017.

[87] A. Kendall and R. Cipolla, "Geometric loss functions for camera pose regression with deep learning," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2017, pp. 5974–5983.

[88] K. Burnett, D. J. Yoon, Y. Wu, A. Z. Li, H. Zhang, S. Lu, J. Qian, W.-K. Tseng, A. Lambert, K. Y. Leung *et al.*, "Boreas: A multi-season autonomous driving dataset," *Int. J. Rob. Res.*, vol. 42, no. 1-2, pp. 33–42, 2023.

[89] W. Song, D. Li, S. Sun, X. Xu, and G. Zu, "Registration for 3D LiDAR datasets using Pyramid Reference Object," *IEEE Trans. Instrum. Meas.*, vol. 72, pp. 1–9, 2023.

[90] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? The KITTI vision benchmark suite," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.* IEEE, 2012, pp. 3354–3361.

[91] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proc. Int. Conf. Learn. Representations*, 2015, pp. 1–15.

[92] Y. Feng, H. You, Z. Zhang, R. Ji, and Y. Gao, "Hypergraph neural networks," in *Proc. AAAI Conf. Artificial Intell.*, 2019, pp. 3558–3565.

[93] Z. Liu, M. Sun, T. Zhou, G. Huang, and T. Darrell, "Rethinking the value of network pruning," in *Proc. Int. Conf. Learn. Representations*, 2019, pp. 1–13.

[94] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, "Optuna: A next-generation hyperparameter optimization framework," in *Proc. ACM SIGKDD Int. Conf. Knowl. Discov. & Data Mining*, 2019, pp. 2623–2631.

[95] G. Fang, X. Ma, and X. Wang, "Structural pruning for diffusion models," *arXiv preprint arXiv:2305.10924*, 2023.

[96] OpenAI, "ChatGPT [Large language model]," 2023. [Online]. Available: https://chat.openai.com/chat