

# Density-invariant Features for Distant Point Cloud Registration

Quan Liu    Hongzi Zhu\*    Yunsong Zhou  
 Shanghai Jiao Tong University  
 {liuquan2017, hongzi, zhouyunsong}@sjtu.edu.cn

Hongyang Li  
 Shanghai AI Lab  
 hy@opendrivelab.com

Shan Chang  
 Donghua University  
 changshan@dhu.edu.cn

Minyi Guo  
 Shanghai Jiao Tong University  
 guo-my@cs.sjtu.edu.cn

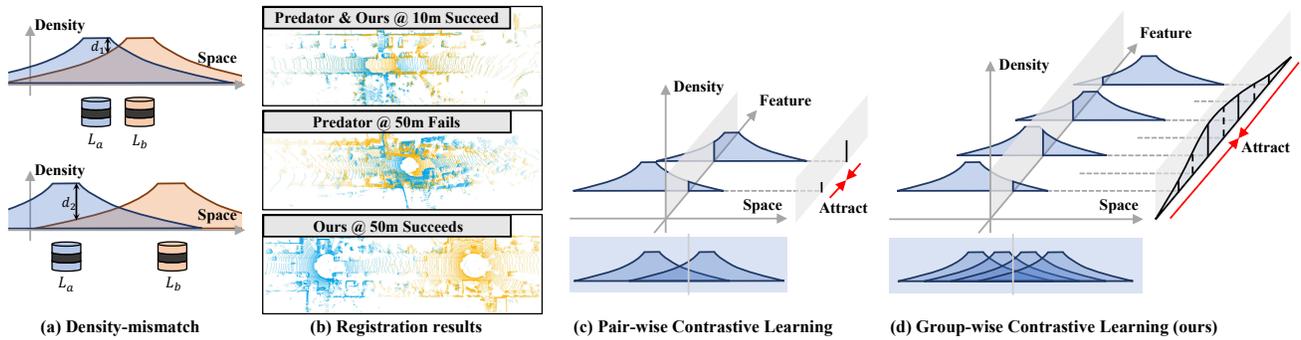


Figure 1: **Motivation.** (a) A schematic diagram of the density  $d$  of point clouds produced by two LiDARs  $L_a$  and  $L_b$ . Their densities diverge when two LiDARs drift from short-range to long-range scenario. (b) State-of-the-art method Predator [20] fails to register distant point clouds (middle) while our method succeeds (bottom). (c) Traditional Pair-wise Contrastive Learning takes a pair of positive samples (marked with the gray plane) from two point clouds, and pull positive features together, which suffer from negative density correlation given a fixed pair of distant point clouds. (d) Our Group-wise Contrastive Learning takes multiple positive samples (marked with the gray plane) from a point cloud series, and contract the feature distance between all positive examples.

## Abstract

Registration of distant outdoor LiDAR point clouds is crucial to extending the 3D vision of collaborative autonomous vehicles, and yet is challenging due to small overlapping area and a huge disparity between observed point densities. In this paper, we propose Group-wise Contrastive Learning (GCL) scheme to extract density-invariant geometric features to register distant outdoor LiDAR point clouds. We mark through theoretical analysis and experiments that, contrastive positives should be independent and identically distributed (i.i.d.), in order to train density-invariant feature extractors. We propose upon the conclusion a simple yet effective training scheme to force the feature of multiple point clouds in the same spatial location (referred to as positive groups) to be similar, which naturally avoids the sampling bias introduced by a pair of point

clouds to conform with the i.i.d. principle. The resulting fully-convolutional feature extractor is more powerful and density-invariant than state-of-the-art methods, improving the registration recall of distant scenarios on KITTI and nuScenes benchmarks by 40.9% and 26.9%, respectively. Code is available at <https://github.com/liuQuan98/GCL>.

## 1. Introduction

Point cloud registration is the cornerstone technique for various computer vision tasks such as SLAM [29, 30], scene flow estimation [26, 27], and early/late fusion in 3D scene understanding [54, 52, 22, 55]. Due to the complex scenarios and the large scale of outdoor point clouds, point cloud registration in driving scenarios is a more challenging and rewarding task than indoor scenes, which can help expand the perceptual field of collaborative vehicles for enhancing the driving safety. In such scenarios, point cloud registration should be accurate enough even when dealing with ex-

\* Corresponding author

remely low-overlap point clouds (e.g., two LiDARs of interest may be over 50 meters apart) to secure downstream tasks such as object detection [55, 22], segmentation [48], and tracking [44, 37].

As depicted in Figure 1(a), as a pair of overlapping point clouds drift apart, they scan the same location with increasingly different densities because point cloud local density reduces quadratically with the distance from the LiDAR. This is referred to as the *density-mismatch* problem exclusively found on distant outdoor point clouds. As generally discovered by literature [46, 39, 49, 20, 34], deep learning features are sensitive to point density so that density-mismatch results in low feature similarity and harms registration performance.

In recent years, there has been a boom in learning-based outdoor point cloud registration methods [10, 15, 32, 1, 49, 9, 5, 20, 50], all of which train feature extractors based on a pair of point clouds, referred to as the pair-wise contrastive learning (PCL) technique [10]. As stated by Arora *et al.* [3], a prerequisite for PCL network convergence is that data samples in a positive pair, which represents scans of a location made by two LiDARs, have to be independently and identically distributed (*i.i.d.*). However, severe density-mismatch leads to violation of the *i.i.d.* principle. More specifically, given a pair of distant point clouds  $S, T$ , it is highly likely that a high-density location in point cloud  $S$  corresponds to a low-density location in point cloud  $T$  and vice versa, *i.e.*, their densities are correlated, as depicted in Figure 1(c). Despite adopting several density-related techniques such as voxelization [15, 28, 51, 34, 49, 9, 20, 50] or density-adaptive calculation [5, 46], existing methods still fail when handling distant outdoor point clouds, as depicted in Figure 1 (b), where state-of-the-art method Predator [20] fails to register two point clouds of 50 meters apart.

Based on the analysis above, we remark that *i.i.d.* positive features *w.r.t.* density are essential for training *density-invariant* feature extractors to solve the distant point cloud registration problem.

In this paper, we propose *group-wise contrastive learning* (GCL) scheme for *density-invariant* feature extraction in order to register distant point clouds. The core idea of GCL is to utilize groups of highly-overlapped point clouds, continuously collected by moving vehicles of public datasets, to break the density correlation between positive examples. As illustrated in Figure 1(d), GCL aligns such point clouds, and collects all point correspondences and their features at one location, referred to as a *positive group*. A large enough positive group can better approximate the underlying feature distribution. Furthermore, given a positive group and a specific positive sample in the group, the density of its possible correspondence is unknown since we do not know which point cloud the other sample belongs to. As a result, GCL positive samples are approximately

*i.i.d.* and can be used for designing compelling group-wise loss to train the density-invariant feature extractor.

How to engage the extractor to derive consistent features over inconsistent densities in positive groups is non-trivial. One straightforward solution is to minimize the variance of all features, which is not sufficient to set the optimal convergence target. In contrast, we additionally ask the mean of a positive group to be close to its *finest feature*, which is defined as the feature derived from the point with the highest point density in the group. By doing this, we force features extracted with low point densities to be similar with the most descriptive one in a positive group, facilitating features in a positive group to converge towards a better consensus.

We implement GCL on both sparse voxel convolution and KPConv, and conduct extensive experiments on KITTI [14] and nuScenes [6]. Results demonstrate that GCL can achieve above 40.9% and 26.9% registration recall gains over SOTA point cloud registration methods when handling far point cloud pairs in KITTI and nuScenes, respectively, without performance loss on near point cloud registration benchmarks. In addition, GCL is lightweight, making it preferable for online distant point cloud registration on smart vehicles.

We highlight our main contributions as follows:

- We theoretically analyze the difficulty of registering distant point clouds, and mark that constructing more *i.i.d.* positive samples is the key to train a density-invariant feature extractor for this challenging task.
- We propose an effective density-invariant feature extraction scheme based on group-wise contrastive learning, where *i.i.d.* positive groups are neatly constructed and new contrastive learning loss are particularly designed.
- We conduct extensive trace-driven experiments on KITTI and nuScenes. Results demonstrate superior density-invariance along with +40.9% and +26.9% registration recall improvements on distant scenarios in KITTI and nuScenes, respectively.

## 2. Related Work

### 2.1. Deep Point Cloud Registration

Recent registration pipelines often adopt a learning-based 3D backbone with contrastive loss to extract local geometry for feature matching. This is the main focus in this work.

**Patch-based features.** Patch based features [53, 10, 15, 32, 1] generally follow the pioneering 3DMatch [53] to extract deep features on pre-selected local patches, and apply contrastive loss on patch embeddings. PPF-Net [10]

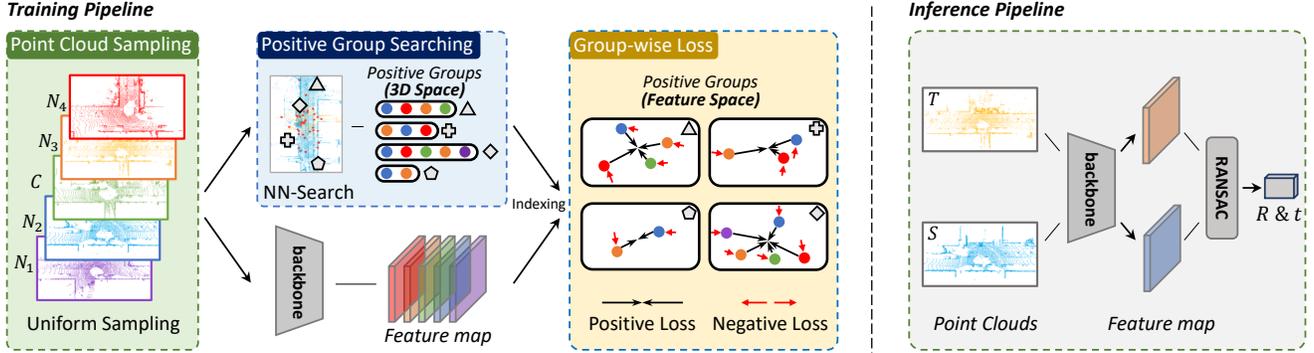


Figure 2: **Overview** for the proposed Group-Wise Contrastive Learning (GCL). The training pipeline of GCL composes of three stages: (1) Point Cloud Sampling, where multiple neighboring point clouds are uniformly sampled around a central point cloud; (2) Positive Group Searching, which collects all correspondences on the same spatial location to form a positive group, then collects their features from the feature map; (3) Group-wise Loss, where contrastive loss is applied on positive groups to forge density-invariant features. During inference, features are extracted for a pair of distant point clouds and then fed into RANSAC [12] to calculate the relative transformation.

further introduces the PointNet backbone [33], while PerfectMatch [15] improves feature robustness with smoothed density value. DIP [32] incorporates a patch reconstruction step. SpinNet [1] uses spherical convolution to achieve  $SO(2)$  equivalence. However, patch-based methods are usually slow due to repeated computation even when patches largely overlap each other, hindering their application in real-time scenarios such as self-driving.

**Fully convolutional features.** Following FCGF [9], fully convolutional methods [9, 5, 20] extract dense features for the whole point cloud in one forward pass and apply contrastive loss to points instead of patches. These methods achieve both state-of-the-art performance and low inference time. D3Feat [5] grants the KPConv [40] extractor the ability of a key-point detector. Predator [20] further improves the low-overlap scenario with an overlap attention module in the bottleneck. Despite their promising performance in indoor or close LiDAR point clouds, their negligence of density variance leads to degraded performance on distant outdoor point clouds, which this paper fixates on improving.

## 2.2. Methods against Density Variation

The conventional techniques aimed at addressing density variation have been extensively used in point cloud processing. However, they are insufficient to solve the distant point cloud registration problem. Voxelization slightly eases density variation and has been adopted by SOTA methods [15, 28, 51, 34, 49, 9, 20, 50], but has limited effect according to Fig. 1(b). Sampling methods such as Farthest Point Sampling [11, 33, 2, 35] achieve uniform density by aggressively dropping points, undermining feature descriptiveness. Density estimation methods, *e.g.*, distance-based [24], KDE [41, 46] or SDV [15], allow density-adaptive fea-

ture calculation, but also requires sufficient input sampling to work properly, which is parallel with our work.

## 2.3. Contrastive Learning

Contrastive Learning, also known as Deep Metric Learning, was first introduced to effectively extract dense visual representations [17, 36, 23, 19, 47, 7, 18, 42], then extended to process audio [45] or texts [13, 21]. A noticeable trend is that sampling more and harder negative samples will improve feature quality due to elevated stability and informativeness [31, 19, 47, 10, 18, 42], but improvements on positive samples are scarce. Our work is the first to examine density variation of 3-dimensional point clouds and apply loss on positives with multiple different densities, which a topic that has not been well studied yet.

## 3. Method

The GCL method overview is depicted in Figure 2, which adopts a feature-based architecture [9, 5, 20]. During training, GCL composes of three stages: (1) Point Cloud Sampling from neighboring views that are uniformly distributed on the road; (2) Positive Group Searching which finds positive groups containing observations of the same spot from various perspectives and distances, and collects their corresponding features; (3) Group-wise Loss that applies advanced constraints on positive groups instead of the positive pair loss adopted in PCL.

During inference, features are extracted on two point clouds  $S, T$  before being fed into a robust estimator such as RANSAC [12] to recover the relative transformation.

### 3.1. Preliminary

**Point cloud registration.** Given two distant point clouds  $S = \{p_S^i \in \mathbb{R}^3 | i = 1, 2, \dots, n\}, T = \{p_T^j \in \mathbb{R}^3 | j =$

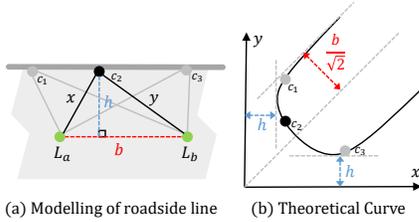


Figure 3: **Our hypothesis.** (a) Looking from bird’s eye view, PCL correspondences  $c$  cluster on roadside lines, and (b) they form a ‘U’ shaped curve when plotting the distance from a correspondence to both LiDARs  $L_a, L_b$  as the  $x$ - $y$  coordinate. For example, three spatially collinear correspondences  $c_1, c_2, c_3$  are bent with the ‘U’ curve when using  $x$  and  $y$  as their coordinates.

$1, 2, \dots, m\}$  with partial overlap, the point cloud registration problem is to uncover the optimal transformation  $R \in SO(3), t \in \mathbb{R}^3$  that resembles the spatial displacements between both LiDARs.

**Pair-wise Contrastive Learning in registration.** Contrastive learning is a pairwise optimization based training paradigm, where data in the same class (positives) are encouraged to have similar features, and data in different classes (negatives) are encouraged to have distinct features. In the context of point cloud registration, positives refer to points in the same spatial location (*i.e.*, correspondences) and negatives refer to points in different locations. In practice, PCL methods [10, 15, 32, 1, 49, 9, 5, 20, 50] usually find point correspondences on two point clouds, and apply contrastive loss on the corresponding features for supervision. They usually have a RANSAC-based inference pipeline as the one depicted in Figure 2.

### 3.2. Analysis

***i.i.d.* positive features.** It is a general basic assumption that positives should be sampled in an *i.i.d.* manner [3], so that the positive features of PCL can converge together to achieve density invariance. Given a specific location in world coordinates, we denote its feature distribution  $D$  as containing all possible features for this specific location observed from all densities and angles. We denote a positive pair as  $(p_S^i, p_T^j) \in C$ , and their features as  $f_S^i, f_T^j \sim D$  from point cloud  $S$  and  $T$ , respectively. Note that  $S$  and  $T$  are also variables.  $C$  denotes all pairs of positive correspondences in this spatial location. The positive loss  $L_{pos}$  on this location is formulated in Equation 1, where parameter  $r \in [1, +\infty)$  represents vector norm, and  $M \in \mathbb{R}^+$  is a tolerance margin.

$$L_{pos} = \frac{1}{|C|} \sum_{(p_S^i, p_T^j) \in C} \max(\|f_S^i - f_T^j\|_r - M, 0) \quad (1)$$

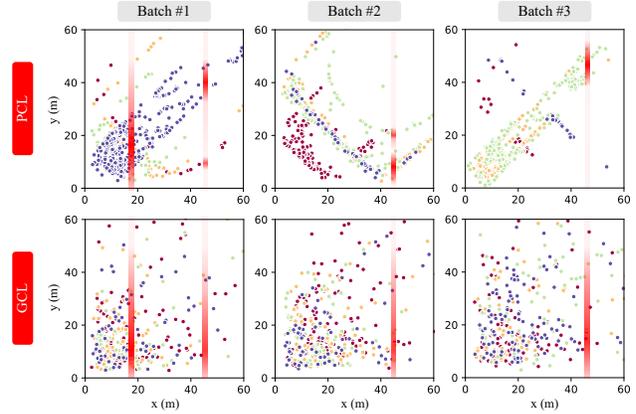


Figure 4: **Visual validation** of our hypothesis in Figure 3. We plot the curves of real correspondences in PCL (top row) and GCL (bottom row) with  $batch\_size = 4$ , where the  $x$  and  $y$  coordinates of a dot denotes the distance from a correspondence to both LiDARs. Points are colored according to in-batch indexes. The red stripes denote the conditional distribution of  $y$  on a fixed  $x$ . PCL positive examples form ‘U’ shapes and are highly correlated, their correlation constantly varying across batches. In contrast, GCL positive examples are more independent and obey a consistent distribution across batches.

**Lemma 1** *If  $f_S^i, f_T^j \sim D$  are i.i.d., then  $\exists \hat{f}$ , so that minimizing  $L_{pos}$  converges in probability to encouraging all features  $f \sim D$  to converge towards the same location  $\hat{f}$  in feature space; Otherwise, non-i.i.d. sampling of  $f_S^i, f_T^j$  will encourage all features to converge towards different locations in feature space.*

Detailed proof is provided in Appendix A. Lemma 1 is the fundamental reason why PCL cannot handle distant point clouds. Besides, density-invariance can be achieved through features of all densities converging to the same location. This hints that positive features must be independent in order to train a density-invariant feature extractor.

**General impact of density-mismatch on PCL.** Ideally, positive features should be independently sampled from all over the space. However, this is not true for PCL due to the sampling bias introduced by distant point clouds. PCL samples positives in the limited overlap of a pair of point clouds, trying to capture everything in the overlap while ignoring others out of the overlap. We will show that the spatially restricted sampling strategy in turn breaks the *i.i.d.* assumption of positive features.

**Theoretical correlation in PCL positive distribution.** Based on our observation, PCL positive correspondences are clustered on roadside lines with a strong density correlation. The equation describing a line  $h$  meters away from a pair of vehicles that are  $b$  meters apart follows equation

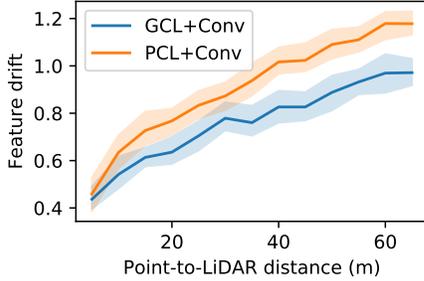


Figure 5: **Feature drift (as defined in Section 3.2) grows with increasing distance from LiDAR.** The quality of GCL features deteriorates slower with decreasing density compared to PCL, indicating the superior density-invariance of GCL.

2, where  $x, y$  denotes the distance from a correspondence  $c \in C$  to both LiDARs  $L_a, L_b$ , respectively. When using  $x, y$  as the coordinates, the roadside line forms a ‘U’ shaped pattern as depicted in Figure 3. Our observation is also supported by Figure 6.

$$\left| \sqrt{x^2 - h^2} \pm \sqrt{y^2 - h^2} \right| = b \quad (2)$$

**Empirical correlation in PCL positive distribution.** We plot the real correspondences for both PCL (top row) and GCL (bottom row) in three different batches with  $batch\_size = 4$  in Figure 4. Note that PCL (top row) exhibits several ‘U’ shapes identical to Figure 3(b). We conclude that the distribution of scan distance  $y$  and  $x$  are strongly correlated for PCL, violating the *i.i.d.* assumption.

**Feature drift and its correlation with density.** Though long-desired in 3D vision, complete density invariance is unattainable due to its ill-posed nature, and deep features are generally affected by input point density. To quantitatively measure the degree of sensitivity to density-variation, we define a metric called *feature drift* which represents the distance from a feature to the feature extracted from the densest observation of this spot. According to Figure 5, there is a roughly linear correlation between the scan distance (quadratic-reciprocal with density) and the feature drift for convolutional backbones. The strong correlation between feature and density is the last piece of the puzzle.

**Summary of the analysis.** Due to the density-mismatch in a pair of distant point clouds, PCL malfunctions under outdoor low-overlap scenario. Specifically, PCL samples correspondences in the limited overlap, resulting in a strong correlation between the point density in a correspondence. This undermines feature independence due to the linear relationship of density and feature. According to Lemma 1, dependent positive features converge towards different locations so that they cannot be matched through Nearest-Neighbor Search, damaging registration performance.

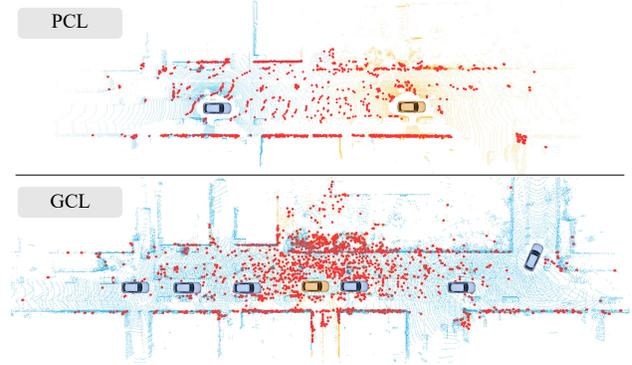


Figure 6: **Spatial distribution of 2000 random sampled positives denoted as red dots for PCL (top) and GCL (bottom).**  $T$  and  $C$  are tainted orange while  $S$  and  $\{N_k\}$  are tainted blue. PCL finds correspondences in the limited overlap, which severely biases towards roadside lines and ignores important object details both near and away from both LiDARs. In contrast, GCL positives are broadly scattered, obeying the same distribution as the central frame  $C$ .

Additionally, we believe that density correlation should also exist in other scenarios than self-driving. For example, a small indoor object will appear as clustered dots for PCL on Figure 4, which is another kind of density correlation. Consequently, our analysis should still hold on other scenarios. Please refer to Appendix C for more discussion.

### 3.3. Overview of GCL

As depicted in Figure 2, during Point Cloud Sampling, GCL divides the region of  $[-60m, 60m]$  around a central frame  $C$  into  $\phi$  segments and select one point cloud from each of the segments with uniform distribution, forming the neighborhood point clouds  $\{N_k | k = 1, \dots, \phi\}$ . Then, Positive Group Searching is carried out through nearest-neighbor search on the aligned point clouds to find multiple correspondences in the same spatial location, referred to as a positive group. As the point clouds complement each other to cover the whole space, positive groups could be found everywhere without special bias towards roadside lines, as depicted in Figure 6. The point clouds are then passed through a feature backbone to extract features for every point, and the corresponding features for positive groups are collected. Finally, group-wise losses are applied on positive groups to train the feature extraction network to be density-invariant.

### 3.4. Positive Group Searching

Intuitively, the whole central frame is the most consistent yet non-biased spatial distribution of points, which is ideal for positive sampling. We follow this idea to launch nearest-neighbor search from every single point in the central frame  $C$  to all neighboring point clouds  $\{N_k | k = 1, \dots, \phi\}$ , and

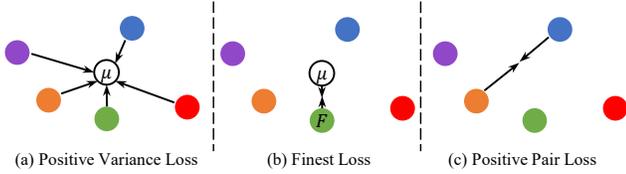


Figure 7: **Loss designs in positive groups.** Traditional Positive Pair Loss optimizes a pair of features, while GCL optimizes both the variance and the mean of a positive group through Positive Variance Loss and Finest Loss, respectively.  $\mu$  denotes group mean while  $F$  denotes finest feature (see Section 3.5 for definition).

gather all matches to form a positive group. Often having more than 2 matched points, the positive group can better approximate the underlying continuous high-dimensional feature distribution. It is still possible that some points in  $C$  have no nearest neighbor and we simply discard them. Empirically,  $87\% \pm 5.4\%$  of the points in the central frame are able to form positive groups on KITTI. Consequently, the density of GCL positives are far more independent and consistent than PCL, as depicted in Figure 4.

### 3.5. Group-wise Loss

In principle, the positive group loss should cater to positive groups instead of positive pairs to fully exploit in-group information. With a slight overload of subscript, we denote all positive groups as  $G = \{g_x\}$  which are formulated as sets of  $d$ -dimensional features  $g_x = \{f_x^i \in \mathbb{R}^d\}$ . Note that positive groups can be formulated using either points or features, as every point has its corresponding feature vector. All loss formulations are visualized in Figure 7. We formulate the positive pair loss  $L_{PP}$  as Equation 3, where  $f_x^i, f_x^j$  are a single pair of correspondence sampled from  $g_x$ . For all loss formulations, parameter  $r \in [1, +\infty)$  represents vector norm and  $m_1, m_2, m_3, m_4 \in \mathbb{R}^+$  are 4 tolerance margins.

$$L_{PP} = \frac{1}{|G|} \sum_{g_x \in G} \max(\|f_x^i - f_x^j\|_r - m_1, 0) \quad (3)$$

In contrast, we explore two new loss functions, namely Positive Variance Loss  $L_{PV}$  and Finest Loss  $L_F$ . Positive Variance Loss minimizes the in-group feature variance to reduce sampling instability, as shown in Equation 4, where  $\mu(\cdot)$  is the average operator. Finest Loss in turn improves the subtle in-group structure, asking the mean feature to be close to the *finest feature* as formulated in Equation 5. Finest feature is defined as the feature extracted on the point cloud with highest density in this group, and the function  $\mathcal{F}(\cdot)$  returns the finest feature of a group.

$$L_{PV} = \sum_{g_x \in G} \sum_{f_x^i \in g_x} \frac{\max(\|f_x^i - \mu(g_x)\|_r - m_2, 0)}{|G| \times |g|} \quad (4)$$

$$L_F = \sum_{g_x \in G} \frac{\max(\|\mathcal{F}(g_x) - \mu(g_x)\|_r - m_3, 0)}{|G|} \quad (5)$$

We adopt the hardest-negative loss [9], which is generalized to group-wise form as Equation 6, where  $\mathcal{H}(\cdot, \cdot, \cdot)$  defined in Equation 7 returns the distance from a feature to its hardest negative, the latter defined as the nearest neighbor of the feature among all non-correspondence features. The total loss is defined as  $L = \lambda_1 L_{PV} + \lambda_2 L_F + \lambda_3 L_{HN}$ , where  $\lambda_1, \lambda_2, \lambda_3$  controls the ratio between loss terms.

$$L_{HN} = \sum_{g_x \in G} \sum_{f_x^i \in g_x} \frac{\max(m_4 - \mathcal{H}(f_x^i, g, G), 0)}{|G| \times |g|} \quad (6)$$

$$\mathcal{H}(f_x^i, g^x, G) = \min_{g_y \neq g_x \in G, f_y^j \in g_y} \|f_x^i - f_y^j\|_r \quad (7)$$

## 4. Results

We test GCL design on two outdoor vehicle-mounted LiDAR registration datasets KITTI [14] and nuScenes [6] with comparison and ablation studies. We then compare GCL to a mass sampling baseline and showcase the density invariance of GCL.

### 4.1. Experiment Setup

**Datasets.** We validate GCL design on two commonly used outdoor registration datasets KITTI [14] and nuScenes [6]. We sub-divide the PCL datasets with different registration difficulty, measured by the distance between two LiDARs, using  $[b_1, b_2]$  to denote that the distance is uniformly sampled between  $b_1$  and  $b_2$  in meters. These PCL datasets are used both during training of previous methods and testing of all methods. As SpinNet does not open-source its dataset preparation code, we report the test results on its pretrained models. The GCL datasets are only used during training, and are created differently as discussed in Section 3.3 but consists of similar size to PCL datasets. We also distill low-overlap datasets with  $\leq 30\%$  overlap, denoted as LoKITTI and LoNuScenes following the methodology of Predator [20]. We follow the general protocols [9, 6] to divide KITTI and nuScenes into train-val-test splits. Please refer to Appendix B for details.

**Training.** As previous methods often struggle to converge on distant LiDAR point cloud pairs, we start by training a baseline model on [5, 20] and fine-tune 4 additional models on [5, 30], [5, 40], [5, 50], [5, 60], then report the best performance of all models. However, there is no pair-wise LiDAR distance difference for GCL, so only one GCL model is trained and tested on arbitrary distances.

Dataset	mRR	[5,10]	[10,20]	[20,30]	[30,40]	[40,50]
FCGF [9]	55.2	97.0	85.4	54.1	25.0	14.3
Predator [20]	74.7	<u>99.3</u>	<b>96.8</b>	<u>90.2</u>	60.6	26.7
SpinNet [1]	35.6	97.6	73.1	7.3	0.0	0.0
D3Feat [5]	52.5	98.7	86.8	52.7	20.0	4.5
CoFiNet [51]	68.6	<b>99.6</b>	94.2	80.0	44.8	24.3
GeoTransformer [34]	39.0	97.9	88.3	8.3	0.7	0.0
GCL+KPCnv (ours)	<u>83.5</u>	99.1	<u>96.5</u>	89.3	<u>78.6</u>	<u>54.1</u>
GCL+Conv (ours)	<b>88.8</b>	98.4	96.1	<b>94.1</b>	<b>87.6</b>	<b>67.6</b>

Table 1: Comparison of RR (%) between SOTA methods and GCL on five *KITTI* [ $b_1, b_2$ ] datasets, with increasing LiDAR distance and registration difficulty. The mean RR is displayed in the first column.

Dataset	mRR	[5,10]	[10,20]	[20,30]	[30,40]	[40,50]
FCGF [9]	37.0	78.4	46.6	27.6	22.0	10.2
Predator [20]	39.5	96.6	50.9	32.8	9.8	7.4
GCL-Conv (ours)	<u>70.2</u>	<u>97.6</u>	<u>88.0</u>	<u>71.7</u>	<u>56.7</u>	<b>37.1</b>
GCL-KPCnv (ours)	<b>71.5</b>	<b>99.0</b>	<b>91.0</b>	<b>77.7</b>	<b>57.3</b>	<u>32.5</u>

Table 2: Comparison of RR (%) between SOTA methods and GCL on five *nuScenes* [ $b_1, b_2$ ] datasets, with increasing LiDAR distance and registration difficulty. The mean RR is displayed in the first column.

## 4.2. Performance Comparison

We compare GCL against SOTA methods under five [ $b_1, b_2$ ] datasets on both *KITTI* and *nuScenes* with increasing registration difficulty, under a rigid registration criterion of  $RTE \leq 0.6m$ ,  $RRE \leq 1.5^\circ$ . The RR of all methods on *KITTI* and *nuScenes* are listed in Table 1 and Table 2, respectively. The mean RR is shown in the first column.

On *KITTI* dataset, unfortunately, the training of D3Feat [5] and GeoTransformer [34] do not converge, resulting in huge performance drop on harder datasets. Observing from the mRR, both GCL+Conv and GCL+KPCnv are the best across all methods, achieving 88.8% (+14.1%) and 83.5% (+8.8%) mRR improvement over the closest competitor Predator. Despite performing on par with SOTA methods on the saturated easy datasets including *KITTI* [5,10] and *KITTI* [10,20], GCL receives drastic improvements on scenarios with larger registration difficulty, achieving 67.6% (+40.9%) and 54.1% (+27.4%) RR on *KITTI* [40,50] dataset compared to Predator, which proves the outstanding discriminative power and density invariance of GCL.

On *nuScenes* dataset, the improvements of GCL is much more broadly aware than those on *KITTI*, achieving 70.2% (+30.7%) and 71.5% (+32.0%) mRR with GCL+Conv and GCL+KPCnv, respectively. Dramatic improvements are visible even in close-range datasets including *nuScenes* [10,20], where 88.0% (+37.1%) and 91.0% (+40.1%) RR improvements are made by GCL+Conv and GCL+KPCnv, respectively. While GCL+KPCnv can deal with slight density variation on *nuScenes* [5,10] to [30,40], GCL+Conv is better under extreme density variation as it performs the

Loss	<i>LoKITTI</i>			<i>KITTI</i> [10,10]		
	RR	RTE	RRE	RR	RTE	RRE
C	70.1	27.0	<u>1.07</u>	<u>99.0</u>	7.25	0.28
F	53.2	37.1	1.43	<b>99.2</b>	6.86	<u>0.26</u>
PP	70.3	<u>25.9</u>	1.06	98.6	6.90	0.28
PV	64.4	<b>24.2</b>	<u>1.07</u>	<b>99.2</b>	<b>6.38</b>	<b>0.24</b>
BF+PP	50.1	34.7	1.43	<b>99.2</b>	7.40	0.31
F+PP	<u>72.1</u>	27.3	0.97	<b>99.2</b>	7.11	<u>0.26</u>
<b>F+PV</b>	<b>72.3</b>	25.9	<b>1.03</b>	98.6	<u>6.62</u>	<u>0.26</u>

Table 3: Ablation of loss designs for GCL on *KITTI* [10,10] and *LoKITTI*, measured by RR (%), RTE (cm), and RRE ( $^\circ$ ). F+PV is selected due to having the best performance on *LoKITTI*. The gray column is the main metric.

$\phi$	RR (%) on <i>LoKITTI</i>	RR (%) on <i>KITTI</i> [10,10]	Training time (h)	Feature extraction time (ms)
2	46.3	98.6	32.1	44.2
4	68.8	<b>98.8</b>	38.6	44.2
<b>6</b>	<b>72.3</b>	98.6	58.6	44.2
8	71.4	98.6	71.0	44.2
10	69.7	98.4	78.1	44.2

Table 4: Ablation on number of point clouds  $\phi$  from 2 to 10 for GCL on *KITTI* [10,10] and *LoKITTI*.  $\phi = 6$  is selected to achieve the highest RR on *LoKITTI*. The gray column is the main metric.

best on [40,50] dataset both on *KITTI* and *nuScenes*. We conclude that the representative capability of GCL is even more significant on harder datasets such as *nuScenes*. Example registration results are placed in Figure 8.

## 4.3. Ablation Study

**Loss ablation.** We ablate various GCL loss components and display the registration performance of GCL+Conv on both *KITTI* [10,10] and *LoKITTI* in Table 3. *C* denotes Circle Loss [38] reimplemented on positive groups; *PP*, *F*, *PV* are Positive Pair Loss, Finest Loss, and Positive Variance Loss, respectively; *BF* is a variant of Finest Loss where the gradient of the finest feature is blocked, under the assumption that the finest feature should not be moved. Generally, most configurations perform similarly on *KITTI* [10,10]. While *F* and *PV* are inferior to *PP* on their own, they perform better when combined. *F+PP* and *F+PV* performs the best on *LoKITTI*, achieving 72.1% and 72.3% RR, respectively. We keep *F+PV* as the optimal loss configuration during other experiments.

**Number of point clouds  $\phi$ .** The number of point clouds is the major difference between GCL and PCL, as GCL with  $\phi = 1$  degrades to PCL. We ablate  $\phi$  from 2 to 10 and list GCL performance and training time on *KITTI* [10,10] and *LoKITTI* in Table 4. While GCL with larger  $\phi$  receives extended training time, the feature extraction time during inference stay constant. Considering the major metric (*i.e.*, RR

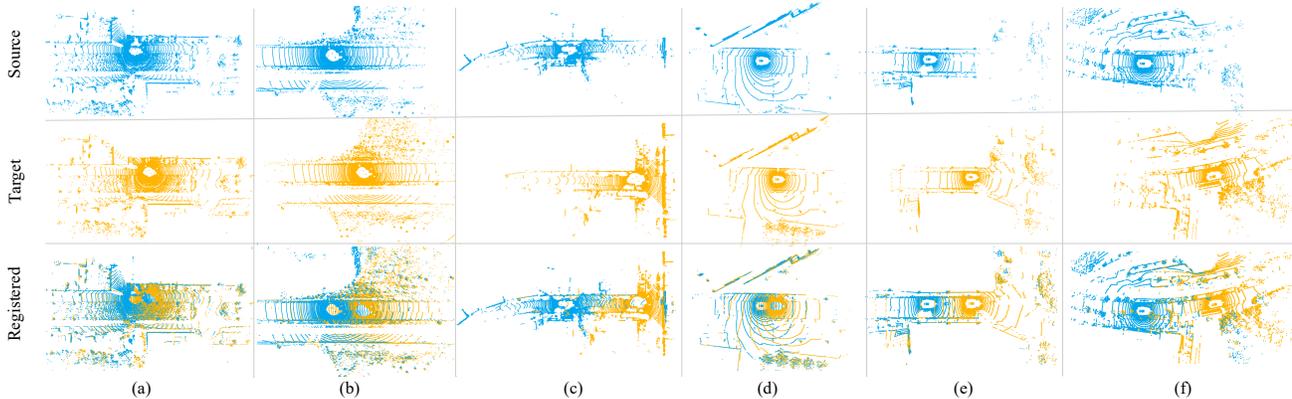


Figure 8: **Example registration results of GCL on KITTI (a-c) and nuScenes (d-f).** The point clouds are sampled from (a) KITTI [5,10], (b) KITTI [20,30], (c) KITTI [40,50], (d) nuScenes [5,10], (e) nuScenes [20,30], and (f) nuScenes [40,50] datasets, respectively. Distant point clouds are both significantly harder to register and more rewarding to downstream tasks.

Data upscale	1×	2×	3×	5×	10×	20×	GCL
RR (%) on <i>LoKITTI</i>	24.0	24.5	25.2	25.8	27.9	28.8	72.3
Time per epoch (min)	6	12	19	34	59	127	58

Table 5: **Comparison to a mass sampling baseline.** Scaling up the number of training point clouds of FCGF by up to 20× the original amount results in marginal improvements. Performance of GCL+Conv is displayed in the last column.

on *LoKITTI*), GCL performance quickly degrades with low  $\phi$ , while peaks at  $\phi = 6$ . With  $\phi > 6$ , the gain from GCL marginalizes, but the motion blur worsens with multiple frames, causing marginal degradation. We pick  $\phi = 6$  to achieve the best performance.

**Comparison to a mass sampling baseline.** A common misconception is that GCL succeeds simply through sampling several times more point clouds than PCL. In this study, we craft a baseline by increasing the number of sampled point cloud pairs for FCGF from 1× to 20×. In particular, PCL with a 3.5× upscale uses an equivalent number of point clouds as GCL with  $\phi = 6$ . The performance of FCGF on *LoKITTI* is displayed in Table 5, which indicates that mass sampling will only contribute marginally to network performance.

**Density Invariance.** We visualize the distribution of feature drift for PCL and GCL with two types of backbones in Figure 5 and Figure 9. Feature drift reveals the backbone’s sensitivity to density change and should ideally stay low for complete density invariance. GCL methods have consistently lower feature drift compared to the baselines, with both a lower average and smaller variance. This indicates that GCL-trained FCNs are more density-invariant than PCL-trained models.

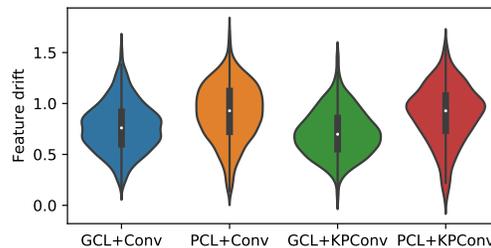


Figure 9: **Distribution of feature drift** (see Section 3.2 for definition) of GCL and PCL with sparse convolution and KPConv. GCL has constantly lower feature drift on two different convolutional architectures, which indicates the superior density-invariance and wide applicability of GCL.

## 5. Conclusion

We have proposed GCL, a density-invariant feature extraction scheme for distant outdoor point cloud registration. Through the joint sampling and optimization on a group of point clouds, GCL naturally avoids being affected by the low-overlap and density-mismatch problem on a pair of distant point clouds. Fully convolutional networks trained with GCL are more representative and density-invariant than SOTA methods trained with PCL, thus being preferable for distant point cloud registration. GCL performs on par with SOTA methods on close point cloud pairs, but exhibits a drastic performance increase on distant point cloud pairs, resulting in a much higher overall performance. Specifically, GCL surpasses SOTA methods with +40.9% RR on KITTI and +26.9% RR on nuScenes in distant scenarios.

**Acknowledgement.** This work was supported in part by the National Natural Science Foundation of China (Grant No. 61972081), the Natural Science Foundation of Shanghai (Grant No. 22ZR1400200), and the Fundamental Research Funds for the Central Universities (No. 2232023Y-01). We would also like to thank reviewers and Li Chen from OpenDriveLab for fruitful discussions.

## References

- [1] Sheng Ao, Qingyong Hu, Bo Yang, Andrew Markham, and Yulan Guo. Spinnet: Learning a general surface descriptor for 3d point cloud registration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11753–11762, 2021. [2](#), [3](#), [4](#), [7](#), [13](#)
- [2] Yasuhiro Aoki, Hunter Goforth, Rangaprasad Arun Srivatsan, and Simon Lucey. Pointnetlk: Robust & efficient point cloud registration using pointnet. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7163–7172, 2019. [3](#)
- [3] Sanjeev Arora, Hrishikesh Khandeparkar, Mikhail Khodak, Orestis Plevrakis, and Nikunj Saunshi. A theoretical analysis of contrastive unsupervised representation learning. *arXiv preprint arXiv:1902.09229*, 2019. [2](#), [4](#)
- [4] Xuyang Bai, Zixin Luo, Lei Zhou, Hongkai Chen, Lei Li, Zeyu Hu, Hongbo Fu, and Chiew-Lan Tai. Pointdsc: Robust point cloud registration using deep spatial consistency. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15859–15869, 2021. [13](#)
- [5] Xuyang Bai, Zixin Luo, Lei Zhou, Hongbo Fu, Long Quan, and Chiew-Lan Tai. D3feat: Joint learning of dense detection and description of 3d local features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6359–6367, 2020. [2](#), [3](#), [4](#), [7](#), [12](#), [13](#)
- [6] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multi-modal dataset for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, June 2020. [2](#), [6](#)
- [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *Proceedings of the International Conference on Machine Learning*, pages 1597–1607, 2020. [3](#)
- [8] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3075–3084, 2019. [13](#)
- [9] Christopher Choy, Jaesik Park, and Vladlen Koltun. Fully convolutional geometric features. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019. [2](#), [3](#), [4](#), [6](#), [7](#), [12](#), [13](#)
- [10] Haowen Deng, Tolga Birdal, and Slobodan Ilic. Ppfnet: Global context aware local features for robust 3d point matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 195–205, 2018. [2](#), [3](#), [4](#), [13](#)
- [11] Yuval Eldar, Michael Lindenbaum, Moshe Porat, and Yehoshua Y Zeevi. The farthest point strategy for progressive image sampling. *IEEE Transactions on Image Processing*, 6(9):1305–1315, 1997. [3](#)
- [12] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981. [3](#)
- [13] Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821*, 2021. [3](#)
- [14] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2012. [2](#), [6](#)
- [15] Zan Gojcic, Caifa Zhou, Jan D Wegner, and Andreas Wieser. The perfect match: 3d point cloud matching with smoothed densities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5545–5554, 2019. [2](#), [3](#), [4](#), [13](#)
- [16] Zan Gojcic, Caifa Zhou, and Andreas Wieser. Learned compact local feature descriptor for tls-based geodetic monitoring of natural outdoor scenes. *International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, 4:113–120, 2018. [12](#)
- [17] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 1735–1742. IEEE, 2006. [3](#)
- [18] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020. [3](#)
- [19] Alexander Hermans, Lucas Beyer, and Bastian Leibe. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737*, 2017. [3](#)
- [20] Shengyu Huang, Zan Gojcic, Mikhail Usvyatsov, Andreas Wieser, and Konrad Schindler. Predator: Registration of 3d point clouds with low overlap. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4267–4276, 2021. [1](#), [2](#), [3](#), [4](#), [6](#), [7](#), [12](#), [13](#), [14](#), [15](#)
- [21] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *Proceedings of the International Conference on Machine Learning*, pages 4904–4916, 2021. [3](#)
- [22] Jason Ku, Melissa Mozifian, Jungwook Lee, Ali Harakeh, and Steven L Waslander. Joint 3d proposal generation and object detection from view aggregation. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1–8. IEEE, 2018. [1](#), [2](#)
- [23] Vijay Kumar BG, Gustavo Carneiro, and Ian Reid. Learning local image descriptors with deep siamese and triplet convolutional networks by minimising global loss functions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5385–5394, 2016. [3](#)
- [24] Felix Järemo Lawin, Martin Danelljan, Fahad Shahbaz Khan, Per-Erik Forssén, and Michael Felsberg. Density adaptive point set registration. In *Proceedings of the*

- IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3829–3837, 2018. [3](#)
- [25] Yang Li and Tatsuya Harada. Leopard: Learning partial point cloud matching in rigid and deformable scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5554–5564, 2022. [15](#)
- [26] Haisong Liu, Tao Lu, Yihui Xu, Jia Liu, Wenjie Li, and Lijun Chen. Camliflow: bidirectional camera-lidar fusion for joint optical flow and scene flow estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5791–5801, 2022. [1](#)
- [27] Xingyu Liu, Charles R Qi, and Leonidas J Guibas. FlowNet3D: Learning scene flow in 3D point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 529–537, 2019. [1](#)
- [28] Fan Lu, Guang Chen, Yinlong Liu, Lijun Zhang, Sanqing Qu, Shu Liu, and Rongqi Gu. HregNet: A hierarchical network for large-scale outdoor lidar point cloud registration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16014–16023, 2021. [2](#), [3](#)
- [29] Michael Montemerlo, Sebastian Thrun, Daphne Koller, Ben Wegbreit, et al. FastSLAM: A factored solution to the simultaneous localization and mapping problem. *Aaai/iaai*, 593598, 2002. [1](#)
- [30] Raul Mur-Artal and Juan D Tardós. Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras. *IEEE transactions on robotics*, 33(5):1255–1262, 2017. [1](#)
- [31] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4004–4012, 2016. [3](#)
- [32] Fabio Poiesi and Davide Boscaini. Distinctive 3D local deep descriptors. In *Proceedings of the International Conference on Pattern Recognition*, pages 5720–5727. IEEE, 2021. [2](#), [3](#), [4](#), [13](#)
- [33] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. PointNet: Deep learning on point sets for 3D classification and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 652–660, 2017. [3](#)
- [34] Zheng Qin, Hao Yu, Changjian Wang, Yulan Guo, Yuxing Peng, and Kai Xu. Geometric transformer for fast and robust point cloud registration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11143–11152, 2022. [2](#), [3](#), [7](#), [13](#), [15](#)
- [35] Vinit Sarode, Xueqian Li, Hunter Goforth, Yasuhiro Aoki, Rangaprasad Arun Srivatsan, Simon Lucey, and Howie Choset. PcnNet: Point cloud registration network using pointnet encoding. *arXiv preprint arXiv:1908.07906*, 2019. [3](#)
- [36] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 815–823, 2015. [3](#)
- [37] Martin Simon, Karl Amende, Andrea Kraus, Jens Honer, Timo Samann, Hauke Kaulbersch, Stefan Milz, and Horst Michael Gross. Complexer-yolo: Real-time 3D object detection and tracking on semantic point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019. [2](#)
- [38] Yifan Sun, Changmao Cheng, Yuhang Zhang, Chi Zhang, Liang Zheng, Zhongdao Wang, and Yichen Wei. Circle loss: A unified perspective of pair similarity optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6398–6407, 2020. [7](#)
- [39] Gusi Te, Wei Hu, Amin Zheng, and Zongming Guo. Rgcnn: Regularized graph CNN for point cloud segmentation. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 746–754, 2018. [2](#)
- [40] Hugues Thomas, Charles R Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, François Goulette, and Leonidas J Guibas. Kpconv: Flexible and deformable convolution for point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6411–6420, 2019. [3](#), [13](#), [14](#)
- [41] Berwin A Turlach. Bandwidth selection in kernel density estimation: A review. In *CORE and Institut de Statistique*. Citeseer, 1993. [3](#)
- [42] Xun Wang, Haozhi Zhang, Weilin Huang, and Matthew R Scott. Cross-batch memory for embedding learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6388–6397, 2020. [3](#)
- [43] Bingli Wu, Jie Ma, Gaojie Chen, and Pei An. Feature interactive representation for point cloud registration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5530–5539, 2021. [15](#)
- [44] Hai Wu, Wenkai Han, Chenglu Wen, Xin Li, and Cheng Wang. 3D multi-object tracking in point clouds based on prediction confidence-guided data association. *IEEE Transactions on Intelligent Transportation Systems*, 23(6):5668–5677, 2021. [2](#)
- [45] Ho-Hsiang Wu, Prem Seetharaman, Kundan Kumar, and Juan Pablo Bello. Wav2clip: Learning robust audio representations from clip. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 4563–4567, 2022. [3](#)
- [46] Wenxuan Wu, Zhongang Qi, and Li Fuxin. PointConv: Deep convolutional networks on 3D point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9621–9630, 2019. [2](#), [3](#)
- [47] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3733–3742, 2018. [3](#)
- [48] Chenfeng Xu, Bichen Wu, Zining Wang, Wei Zhan, Peter Vajda, Kurt Keutzer, and Masayoshi Tomizuka. SqueezeSegV3: Spatially-adaptive convolution for efficient point-cloud segmentation. In *Proceedings of the European Conference on Computer Vision*, pages 1–19. Springer, 2020. [2](#)
- [49] Zi Jian Yew and Gim Hee Lee. 3Dfeat-net: Weakly supervised local 3D features for point cloud registration. In *Proceedings of the European Conference on Computer Vision*, pages 607–623, 2018. [2](#), [3](#), [4](#), [13](#)

- [50] Zi Jian Yew and Gim Hee Lee. Regtr: End-to-end point cloud correspondences with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6677–6686, 2022. [2](#), [3](#), [4](#), [13](#)
- [51] Hao Yu, Fu Li, Mahdi Saleh, Benjamin Busam, and Slobodan Ilic. Cofinet: Reliable coarse-to-fine correspondences for robust pointcloud registration. *Advances in Neural Information Processing Systems*, 34:23872–23884, 2021. [2](#), [3](#), [7](#), [13](#), [15](#)
- [52] Haibao Yu, Yizhen Luo, Mao Shu, Yiyi Huo, Zebang Yang, Yifeng Shi, Zhenglong Guo, Hanyu Li, Xing Hu, Jirui Yuan, et al. Dair-v2x: A large-scale dataset for vehicle-infrastructure cooperative 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21361–21370, 2022. [1](#)
- [53] Andy Zeng, Shuran Song, Matthias Nießner, Matthew Fisher, Jianxiong Xiao, and Thomas Funkhouser. 3dmatch: Learning local geometric descriptors from rgb-d reconstructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1802–1811, 2017. [2](#)
- [54] Xumiao Zhang, Anlan Zhang, Jiachen Sun, Xiao Zhu, Y Ethan Guo, Feng Qian, and Z Morley Mao. Emp: Edge-assisted multi-vehicle perception. In *Proceedings of the 27th Annual International Conference on Mobile Computing and Networking*, pages 545–558, 2021. [1](#)
- [55] Hanqi Zhu, Jiajun Deng, Yu Zhang, Jianmin Ji, Qiuyu Mao, Houqiang Li, and Yanyong Zhang. Vpfnet: Improving 3d object detection with virtual point based lidar and stereo data fusion. *IEEE Transactions on Multimedia*, 2022. [1](#), [2](#)

## A. Proof of Lemma 1

Lemma 1 serves as the basis of our analysis, indicating the fundamental incompetence for PCL, while hinting an *i.i.d.* solution towards density-invariance.

*Proof.* Since we need to investigate the effect of loss on a single feature  $f_i^S$ , we need to marginalize the effect of  $f_j^T$ . We start by selecting a specific first feature  $\hat{f}_i^S$ , regarding it as a constant, and take out all correspondences for the specific feature, which is  $\hat{C} = \{(\hat{f}_i^S, f_j^T) \in C\}$ . We focus on a part of the loss involving this specific  $\hat{f}_i^S$ , referred to as a function  $\hat{L}_{pos}(\hat{f}_i^S)$ , in equation 8.

$$\hat{L}_{pos}(\hat{f}_i^S) = \frac{1}{|\hat{C}|} \sum_{(\hat{f}_i^S, f_j^T) \in \hat{C}} \max(\|\hat{f}_i^S - f_j^T\|_p - m, 0) \quad (8)$$

Next, we marginalize the effect of  $f_j^T$  through sampling infinitely many  $f_j^T$ . Assuming  $\hat{f}_i^S, f_j^T$  are *i.i.d.*, then countless  $f_j^T$  approximates the distribution  $D$ . We can write out the limitation of  $\hat{L}_{pos}$  when  $|\hat{C}| \rightarrow \infty$  as Equation 9.

$$\lim_{|\hat{C}| \rightarrow \infty} \hat{L}_{pos}(\hat{f}_i^S) = \mathbb{E}_{f_j^T \sim D} \max(\|\hat{f}_i^S, f_j^T\|_p - m, 0) \quad (9)$$

Equation 9 is convex and has a single global minimum at  $\hat{f}$  which solely depends on  $D$  (cases where a minimal plateau exists is impossible in real setup). The effect of minimizing  $L_{pos}$  converges in probability to all features  $f \sim D$  heading towards the same location  $\hat{f}$  in feature space.

$$\lim_{|\hat{C}| \rightarrow \infty} \hat{L}_{pos}(\hat{f}_i^S) = \mathbb{E}_{(\hat{f}_i^S, f_j^T) \in \hat{C}} \max(\|\hat{f}_i^S, f_j^T\|_p - m, 0) \quad (10)$$

Otherwise, if  $\hat{f}_i^S, f_j^T$  are non-*i.i.d.*, it is impossible to marginalize  $f_j^T$ , and the loss in Equation 10 is the expectation on a subset of correspondences  $\hat{C}$  that is correlated with  $\hat{f}_i^S$ . All likely features have different loss formulation with different global minimums. This means that different features will converge towards different locations.

Note that the loss we investigate is a partial representation of the complete loss function, as negative losses are not considered. However, the result is highly likely true even with negative loss added. That is because positive loss controls the sub-structure inside a specific positive cluster, while negative loss controls the large-scale relative structure between different positive clusters, and the negative loss should not disturb positive structures too much when the feature representation stabilizes after the first few epochs.

## B. Detailed Experiment Setup

### B.1. Dataset Preparation

Two kinds of datasets are used in this paper, *i.e.*, pairwise contrastive learning (PCL) datasets and group-wise contrastive learning (GCL) datasets. The PCL datasets contain point cloud pairs that are sampled with a random distance interval  $b$  denoting the distance between two LiDARs. The distance  $b$  is randomly picked for every point cloud pair, and we refer to a sub-divided dataset where  $b_1 \leq b \leq b_2$  as  $[b_1, b_2]$ . Both during training and testing, we always reset the random seed to 0 before finding the required point clouds to produce the exact same point cloud pairs for repeatable results. To create the GCL datasets, we sample central point clouds  $C$  at a fixed interval of 11 frames, then randomly sample neighboring point clouds around each central point cloud according to the process described in Section 3.3. The GCL datasets are never used during testing.

Following Huang *et al.* [20], we define *overlap*  $O$  between a pair of point clouds  $S \in \mathbb{R}^{N \times 3}, T \in \mathbb{R}^{M \times 3}$  as a subset of  $S$  according to Equation 11.

$$O = \left\{ p_S^i \in S \mid \min_{p_T^j \in T} \|p_S^i - p_T^j\|_2 \leq \delta \right\} \quad (11)$$

The overlap denotes the part of  $S$  where at least a corresponding point in  $T$  could be found through nearest-neighbor search of radius  $\delta = 0.45m$ .  $S$  and  $T$  are down-sampled using a voxel size of 0.3m before the search. Overlap ratio is then defined as  $\frac{|O|}{|S|}$ . All point cloud pairs with  $\leq 30\%$  overlap ratio in  $[5, 20]$ ,  $[20, 30]$ ,  $[30, 40]$ ,  $[40, 50]$  datasets are collected on *KITTI* and *nuScenes*, referred to as *LoKITTI* and *LoNuScenes*, respectively. They represent the hardest cases for the distant point cloud registration task.

We follow previous literature [5] to divide *OdometryKITTI* with sequences 0-5 for training, 6-7 for validation, and 8-10 for testing. *NuScenes* is divided sequentially with the first 700 sequences for training, the next 150 sequences for validation and the last 150 sequences for testing.

### B.2. Metrics

Both traditional and new metrics are used during evaluation. Following previous work [16, 20, 5, 9], we report 3 metrics including Registration Recall (RR) defined as percentage of pairs successfully registered, Relative Rotation Error (RRE) defined as the geodesic distance between estimated rotation and ground-truth rotation, and Relative Translation Error (RTE) defined as the euclidean distance between estimated translation and the ground-truth translation. We forge a new metric as the average of RR on  $[5, 10]$ ,  $[10, 20]$ ,  $[20, 30]$ ,  $[30, 40]$ ,  $[40, 50]$  datasets, referred

	Dataloader	Inference	RANSAC	Total
FCGF	6.2	45.4	576.1	627.7
GCL+Conv (ours)	5.4	44.2	523.0	572.6
Predator	635.1	78.7	66.3	780.1
GCL+KPCConv (ours)	637.5	64.4	76.4	778.3

Table 6: **Inference time (ms) analysis on LoKITTI.** GCL is always more lightweight than their existing counterparts with the same backbone (FCGF: Conv; Predator: KPCConv) in terms of inference time.

to as mean Registration Recall (mRR), which measures the overall registration performance.

### B.3. Network Structure

We adopt the popular Res-UNet network structure [9], and implement it on both sparse voxel convolution [8] and KPCConv [40], referred to as GCL+Conv and GCL+KPCConv, respectively. As depicted in Figure 10, both GCL+Conv and GCL+KPCConv adopt three layers of skip connections with a roughly symmetric encoder-decoder design. Features are all normalized onto a unit sphere after the final layer.

### B.4. Loss Configuration

There are several parameters that need specifying for network convergence. The distance margins are set to  $m_1 = 0.1, m_2 = 0.1, m_3 = 0.2, m_4 = 1.4$ . The loss terms are reweighed differently on two datasets, where we set  $\lambda_1 = \lambda_2 = \lambda_3 = 1$  on KITTI and  $\lambda_1 = \lambda_2 = 0.7, \lambda_3 = 1$  on nuScenes.

## C. Additional Experiments

**Inference time.** We list the inference time breakdown for FCGF, Predator, GCL+Conv and GCL+KPCConv in Table 6. The inference time of GCL is always lower than counterparts with the same backbone. While GCL+KPCConv performs faster registration, it requires extended data loading time due to underlying KPCConv architecture conducting repeated nearest neighbor calculation. In contrast, GCL+Conv runs faster during data loading and inference, and the extended RANSAC registration time can be reduced given recent progress on fast registration pipelines [4]. The focus of GCL is to propose a contrastive learning based training method which can be plugged into any existing registration pipelines that incorporate feature matching in it [10, 15, 32, 1, 49, 9, 5, 20, 50], and GCL is the general solution to the distant registration problem on all these methods since they are all based on either Voxel Convolution [8] or KPCConv [40]. We conclude that GCL is a universal lightweight feature extraction method.

Loss	LoKITTI			KITTI [10,10]		
	RR	RTE	RRE	RR	RTE	RRE
C	<u>53.8</u>	32.5	1.41	<u>99.0</u>	7.8	0.27
F	18.3	38.9	1.92	98.8	<u>7.6</u>	<b>0.25</b>
PP	<u>53.8</u>	<b>27.2</b>	<b>1.28</b>	<b>99.2</b>	<u>7.6</u>	<u>0.26</u>
PV	45.0	29.1	1.39	98.6	<u>7.6</u>	<b>0.25</b>
BF+PP	45.7	31.1	1.40	98.6	<u>7.6</u>	<u>0.26</u>
F+PV	50.5	28.4	<u>1.30</u>	<b>99.2</b>	<b>7.5</b>	<b>0.25</b>
<b>F+PP</b>	<b>55.4</b>	<u>27.8</u>	<b>1.28</b>	<b>99.2</b>	7.9	<u>0.26</u>

Table 7: **Ablation of loss designs** for GCL+KPCConv on KITTI [10,10] and LoKITTI, measured by RR (%), RTE (cm), and RRE ( $^{\circ}$ ). F+PP is selected according to performance on LoKITTI. The gray column is the main metric.

Dataset	mRR	[5,10]	[10,20]	[20,30]	[30,40]	[40,50]
FCGF [9]	77.4	98.4	95.3	86.8	69.7	36.9
Predator [20]	87.9	<b>100.0</b>	98.6	<b>97.1</b>	80.6	63.1
SpinNet [1]	39.1	99.1	82.5	13.7	0.0	0.0
D3Feat [5]	66.4	99.8	98.2	90.7	38.6	4.5
CoFiNet [51]	82.1	<u>99.9</u>	<b>99.1</b>	94.1	78.6	38.7
GeoTransformer [34]	42.2	<b>100.0</b>	93.9	16.6	0.7	0.0
GCL+KPCConv (ours)	<u>89.6</u>	<b>100.0</b>	98.2	93.2	<u>88.3</u>	<u>68.5</u>
GCL+Conv (ours)	<b>93.5</b>	99.0	<u>98.8</u>	<u>96.1</u>	<b>91.7</b>	<b>82.0</b>

Table 8: **Comparison of RR (%) between SOTA methods and GCL on five KITTI [ $b_1, b_2$ ] datasets**, with increasing LiDAR distance and registration difficulty. Registration metrics are loosened to  $5^{\circ}$ , 2m compared to Table 1. The mean RR is displayed in the first column.

Dataset	mRR	[5,10]	[10,20]	[20,30]	[30,40]	[40,50]
FCGF [9]	39.5	87.9	63.9	23.6	11.8	10.2
Predator [20]	51.0	<u>99.7</u>	72.2	52.8	16.2	14.3
GCL-Conv (ours)	<u>85.5</u>	99.3	<u>97.7</u>	<u>91.8</u>	<u>77.8</u>	<u>60.7</u>
GCL-KPCConv (ours)	<b>90.3</b>	<b>99.9</b>	<b>98.5</b>	<b>96.1</b>	<b>85.4</b>	<b>71.6</b>

Table 9: **Comparison of RR (%) between SOTA methods and GCL on five nuScenes [ $b_1, b_2$ ] datasets**, with increasing LiDAR distance and registration difficulty. Registration metrics are loosened to  $5^{\circ}$ , 2m compared to Table 2. The mean RR is displayed in the first column.

**Loss ablation with GCL+KPCConv.** We ablate various loss components for GCL+KPCConv and display the registration performance of on both KITTI [10,10] and LoKITTI in Table 7. Similar to results with GCL+Conv, Finest Loss in combination with a positive loss performs the best among all methods, as F+PP achieves both the best RR of 55.4% on LoKITTI and 99.2% KITTI [10,10]. All methods perform roughly the same on the close point cloud dataset KITTI [10,10]. With the KPCConv backbone, however, Finest loss alone does not lead to a decent performance on LoKITTI as it does with the voxel convolution backbone. We select F+PP as the optimal configuration for GCL+KPCConv during all other experiments.

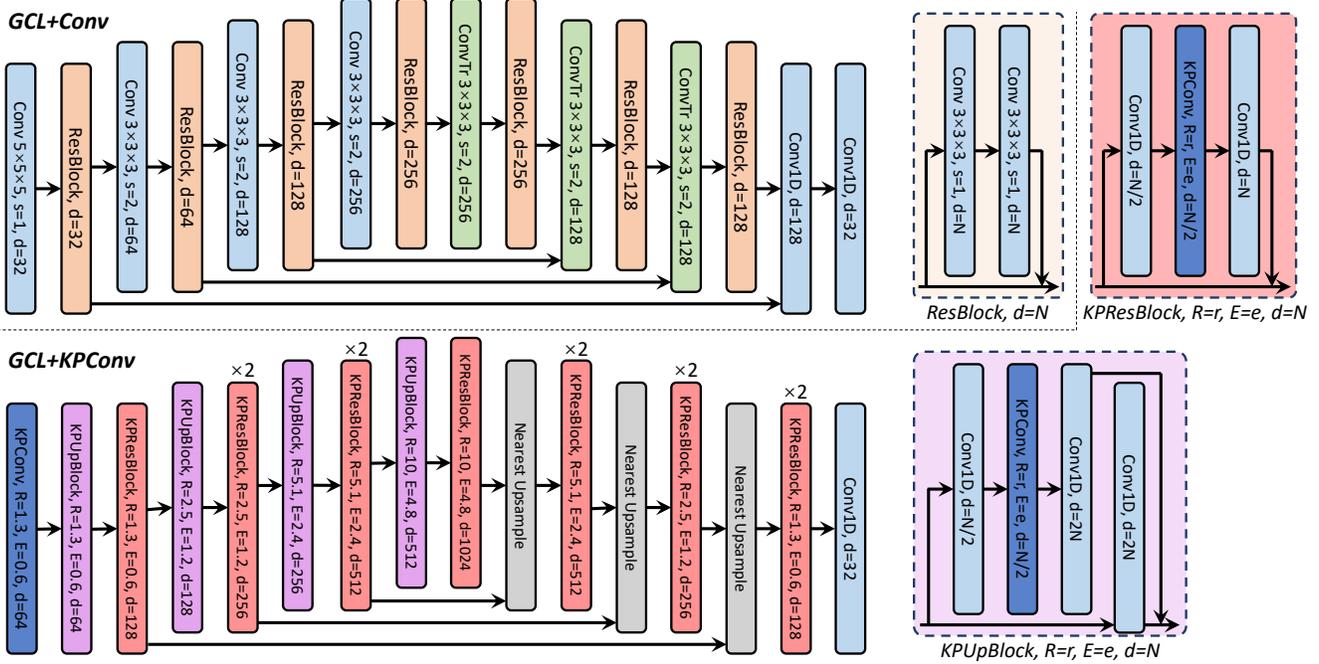


Figure 10: **Network structures for GCL.** Batch Normalization and ReLU activation are used after all Conv blocks except for the last layer, while batch normalization and leaky ReLU are used in KPConv blocks with a 0.1 slope. Voxel Convolution is parameterized by the kernel size, stride  $s$ , and output dimension  $d$ . The kernel size and stride are both omitted for Conv1D. Non-deformable KPConv [40] is parameterized by kernel point offset radius  $R$ , kernel point influence extent  $E$ , and feature dimension  $d$ .

**Performance comparison under loose registration criterion.** We additionally provide the comparison between GCL and SOTA methods on both *KITTI* and *nuScenes* under a loose registration criterion of  $RTE \leq 2m$ ,  $RRE \leq 5^\circ$ , where the registration recalls are generally elevated due to the loosen criterion. The mean RR is shown in the first column. As listed in Table 8, GCL+Conv and GCL+KPCnv achieve the highest overall performance on *KITTI* with 89.6% (+1.7%) and 93.5% (+4.6%) mRR over Predator [20], respectively. Furthermore, GCL methods receive greater improvements on distant scenarios including [30,40] and [40,50] on *KITTI*. On the other hand, GCL methods beat SOTA methods by a larger margin on *nuScenes* than on *KITTI*, achieving 85.5% (+34.5%) and 90.3% (+39.3%) mRR for GCL+Conv and GCL+KPCnv compared to Predator [20], respectively on *nuScenes* according to Table 9. We mark that GCL methods beat SOTAs on every sub-divided dataset on *nuScenes*, and that GCL+KPCnv always performs the best. We conclude that, under a loose registration criterion, GCL still achieves giant improvements comparable to the scenario under a stricter criterion, setting a new SOTA for the distant point cloud registration problem.

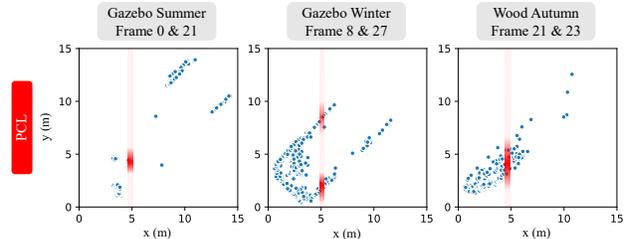


Figure 11: **Distribution of correspondences in ETH dataset** with PCL and  $batch\_size = 1$ , where the  $x$  and  $y$  coordinates of a dot denotes the distance from a correspondence to both LiDARs. The red stripes denote the conditional distribution of  $y$  on a fixed  $x$ . Density correlation still exists in ETH, where linear structures are less dominant.

**Generalization to ETH.** We demonstrate the generalization results from *KITTI* to *ETH* in Table 10, by shrinking the voxel sizes from 0.3m to 0.05m during testing without any finetuning. *ETH* is an outdoor dataset featuring a majority of vegetation over linear structures. However, density correlation still exists in *ETH*, as depicted in Figure 11, which confirms the wide applicability of our analysis in Section 3.2. It can be seen that GCL effectively improves the generalization capability of two baseline back-

	Gazebo		Wood		Avg.
	Summer	Winter	Autumn	Summer	
Predator	21.2	20.8	23.5	30.4	24.0
FCGF	40.2	26.0	54.8	67.2	47.0
GCL+KPCConv (ours)	46.2	28.4	56.5	72.0	50.8
GCL+Conv (ours)	<b>46.7</b>	<b>30.8</b>	<b>61.7</b>	<b>73.6</b>	<b>53.2</b>

Table 10: **Generalization test from KITTI to ETH**, by shrinking voxel size from 0.3m to 0.05m during testing. The FMR scores at  $\tau_1 = 10\text{cm}$ ,  $\tau_2 = 5\%$  are compared.

bones, where GCL+Conv achieves the best overall FMR of 53.2% (+6.2%). We conclude that GCL can generalize to other scenarios other than autonomous driving.

## D. Discussion and Limitation

**More explanation on non-*i.i.d.* PCL positives.** A pair of close-range point clouds also have non-*i.i.d.* positives, as their positives have roughly the same density, *i.e.*, their densities are positively correlated. This may sound weird, as close-range LiDAR point cloud registration has already been well-solved [20, 51]. Actually, non-*i.i.d.* positives will not hinder close-range registration problems because the problem is so simple that even a density-variant feature extractor will solve the problem nicely. Now consider a hand-crafted density-variant feature that upon the input coordinate  $(x, y, z)$ , outputs the vector length of the coordinate  $\sqrt{x^2 + y^2 + z^2}$ . Intuitively, this density-variant feature combined with RANSAC will likely produce a decent guess for two concentric (*i.e.*, extremely close) point clouds. However, this special solution will not work for distant scenarios with severe density mismatch, which means that a more powerful solution like GCL is needed to solve the distant point cloud registration problem.

**Training time.** As listed in Table 4, GCL has a linearly growing training time consumption *w.r.t.*  $\phi$ . This is mainly caused by increased data loading time where repeated nearest neighbor searches are carried out from the central point cloud to all neighborhood point clouds. However, the heavy time consumption is a necessary cost for building the positive groups. Luckily, only training time is affected for GCL and the testing time remain unchanged when registering two point clouds.

**Information exchange.** Information exchange serves as a key source of improvement for SOTA registration methods [20, 51, 43, 25, 34]. It is carried out between a pair of point clouds, which calls for a non-trivial extension of GCL that contains  $\phi + 1$  point clouds. Note that features after the exchange will vary according to different companion point clouds. Consequently, a naive traversal of  $C_{\phi+1}^2$  pairs for

GCL will not only suffer from  $O(\phi^2)$  complexity but also have to deal with  $\phi$  different features for a single point. We hope to extend the information exchange module (mainly composed of cross-attention) to a group-wise version in future work.