# Fast and Accurate Deep Loop Closing and Relocalization for Reliable LiDAR SLAM

Chenghao Shi*, Xieyuanli Chen*, Junhao Xiao†, Bin Dai, Huimin Lu†

*Abstract*—Loop closing and relocalization are crucial techniques to establish reliable and robust long-term SLAM by addressing pose estimation drift and degeneration. This article begins by formulating loop closing and relocalization within a unified framework. Then, we propose a novel multi-head network LCR-Net to tackle both tasks effectively. It exploits novel feature extraction and pose-aware attention mechanism to precisely estimate similarities and 6-DoF poses between pairs of LiDAR scans. In the end, we integrate our LCR-Net into a SLAM system and achieve robust and accurate online LiDAR SLAM in outdoor driving environments. We thoroughly evaluate our LCR-Net through three setups derived from loop closing and relocalization, including candidate retrieval, closed-loop point cloud registration, and continuous relocalization using multiple datasets. The results demonstrate that LCR-Net excels in all three tasks, surpassing the state-of-the-art methods and exhibiting a remarkable generalization ability. Notably, our LCR-Net outperforms baseline methods without using a time-consuming robust pose estimator, rendering it suitable for online SLAM applications. To our best knowledge, the integration of LCR-Net yields the first LiDAR SLAM with the capability of deep loop closing and relocalization. The implementation of our methods will be made open-source.

*Index Terms*—Autonomous Driving, 3D Registration, Deep Learning, Loop Closing, Relocalization

## I. INTRODUCTION

SIMULTANEOUS localization and mapping, also known as SLAM, plays a fundamental role across domains such as autonomous driving, robotics, and computer vision. Ensuring the reliability and stability of a SLAM system is crucial for practical applications. External sensors like GPS and IMU are commonly employed to enhance SLAM in real-world scenarios, but challenges arise when they are inaccessible or unreliable. Hence, attaining enduring reliability in LiDAR-only SLAM becomes important yet challenging. Relocalization and loop closing are pivotal techniques within this context. Relocalization refers to recovering the global pose in local tracking failure, whereas loop closing involves identifying previously visited locations to correct the drift in pose estimation. Despite different objectives, both techniques share similar underlying concepts. They both first coarsely

C. Shi, X. Chen, J. Xiao and H. Lu are with College of Intelligence Science and Technology, National University of Defense Technology, China. B. Dai is with National Innovation Institution of Defense Technology, China.

* Joint first authors with equal contribution.

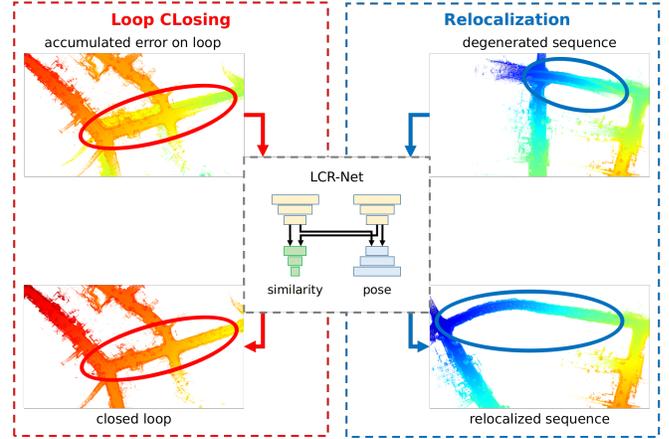† Joint Corresponding authors: {junhao.xiao, lhmnew}.nudt.edu.cn

Fig. 1: Our proposed LCR-Net for loop closing and relocalization. LCR-Net solves both tasks by first retrieving the most similar frame from the map and then estimating the 6-DoF pose.

locate the most similar candidate in the map relative to the current scan, followed by a fine relative pose estimation.

More specifically, current LiDAR loop closing usually first project point clouds into image-like formats to generate global descriptors and estimate 1 degrees of freedom (DoF) [1] or 3-DoF poses [2]–[4]. To achieve full 6-DoF LiDAR loop closing, multiple existing methods [5]–[7] using bag-of-words (BoW) [8]. Such methods extract local features to construct BoW models and utilize such models for loop closure detection and feature matching. The final poses are then obtained by registering the matched features. Since deep neural networks have shown great advances in perception tasks, recent studies have also employed deep learning approaches [9]–[11] for LiDAR loop closing. However, they still follow a similar structure as BoW, albeit using learning-based approaches instead of hand-crafted local features extraction and BoW model. These techniques face a dilemma: obtaining a comprehensive representation of environmental features often requires a deeper encoder. However, a deeper encoder diminishes the count of local features, potentially hindering accurate localization. Additionally, improving registration performance often involves integrating more complex designs into local feature extraction, thereby substantially reducing global description efficiency.

Therefore, despite loop closing and relocalization sharing similar underlying techniques, few works have specifically been proposed for LiDAR relocalization. When local pose tracking degenerates, some point-based SLAM [12] incorporates additional sensors and switches to camera [13] or IMU [14] odometry mode, while some surfel-based methods [15], [16] fall back from frame-to-map to frame-to-frame

pose estimation to avoid errors caused by distorted maps. To our best knowledge, no prior research has addressed both LiDAR loop closing and relocalization simultaneously. In this article, we aim to tackle LiDAR loop closing and relocalization using a joint framework with a novel multi-head network, named LCR-Net, offering four main contributions as follows:

Firstly, we revisit the challenges of addressing LiDAR loop closing and relocalization separately. We identify limitations within the current paradigm and propose **a new framework for solving loop closing and relocalization simultaneously** (Sec. III). Our framework leverages the shared techniques underlying these two tasks, integrating them within a coherent coarse-to-fine framework. This framework concurrently addresses both tasks, beginning with generating global descriptors for an initial coarse global candidate search. Subsequently, the framework generates dense local features to facilitate precise 6-DoF pose estimation. Through such a framework, we circumvent dilemmas arising from the requirements of these two tasks, providing a solid foundation for improved candidate retrieval and registration performance.

Following our framework, we then introduce **a novel multi-head network LCR-Net** (Sec. IV). It employs a shared encoder backbone to downsample and encode the point cloud into three types of sparse features. These features are then processed separately in two distinct heads. One head generates lightweight global descriptors for each scan, enabling fast candidate retrieving. The other head establishes initial sparse matches and then extends them to dense ones for accurate registration. Unlike existing methods using only sparse features, our approach exploits fast dense feature matching based on the neighborhood consistency, achieving accurate and fast pose estimation for online applications without requiring costly robust estimator or iterative pose refinements. To effectively train the multi-head network, we also introduce novel losses with a specific training strategy, which leads to state-of-the-art (SOTA) performance compared to all existing baselines.

To enhance the performance on both tasks, we present the third contribution as the **novel keypoint detection module (Sec. IV-A) comprising two novel sub-modules, 3D-RoFormer++ enabling geometric and contextual information aggregation, and VoteEncoder for reliable keypoint detection and encoding**. Keypoint detection module offers three types of features, enabling fast global descriptor generation and dense reliable match establishing. Initially, it rapidly samples and encodes uniform features across the point cloud for overall representation of the environment. Subsequently, to enhance the features and make them salient and discriminative for improving registration, we introduce 3D-RoFormer++ and VoteEncoder. 3D-RoFormer++ enhances the representation capability of local features with contextual and structural information by enabling the information exchange between the two point clouds. The other module, VoteEncoder effectively downsamples the points while identifying keypoints lying on geometrically significant regions and subsequently aggregating the features from their neighbors. The VoteEncdoer can significantly enhance registration robustness and accuracy by improving the coverage of matching points over the overlapping area of point clouds. The superior improvement brought by VoteEncoder highlights the importance of the match distribution OVER inlier ratio, offering valuable insights (Sec. VI-H) for future research on point cloud registration.

The fourth contribution is **the first LiDAR SLAM system with the capability of deep learning-based loop closing and relocalization** (Sec. V). We build a full LiDAR SLAM system based on our proposed deep loop closing and relocalization method. The system effectively tackles the local pose tracking, loop closing, and relocalization in parallel. Enormous test of our SLAM system in diverse environments and situations showcase the effectiveness and robustness after integrating our proposed LCR-Net.

We extensively evaluate our approach on three setups derived from loop closing and relocalization: candidate retrieval, closed-loop point cloud registration, and continuous relocalization. The results demonstrate that i) our approach outperforms respective baselines and dominates the SOTA in all three tasks, in particular, ii) our approach achieves the best candidate retrieval performance with a simple architecture benefiting from the feature representation capability of the backbone, iii) our approach boosts the baseline by a large margin in registration tasks, even outperforms the baseline method refined by ICP [17], iv) the SOTA registration performance can be achieved efficiently without requiring for a costly RANSAC estimator. We also conduct tests on multiple sequences to assess the performance of our SLAM system. The results depict that integrating our proposed LCR-Net, our SLAM system is capable of addressing the relocalization and loop closing challenges. In the loop closing task, our approach outperforms the most commonly employed loop closing approach, Scan Context [18] combined with ICP. In addition, we provide detailed ablation studies to demonstrate the effectiveness of our design.

## II. RELATED WORK

While various studies have been conducted for the foundational techniques underlying loop closing and relocalization, including candidate retrieval and point cloud registration, few work can simultaneously address both tasks. Therefore, we first introduce the underlying techniques and then explore the recent advancements in loop closing and relocalization.

**Candidate retrieval**, also known as place recognition or loop closure detection, compares the current sensor observation with pre-built maps to determine the approximate location of the robot within the map. LiDAR-based loop closure detection approach can be categorized into global descriptor-based methods and local descriptor-based methods.

Global descriptor-based methods, such as Scan Context [18], represent point clouds as overhead views and encode different segmented spaces to construct global descriptors. Wang et al. [19] extract descriptors using LoG-Gabor threshold filtering and measure similarity using Hamming distance. These methods require additional functions to evaluate the similarity of places, which significantly reduces computational efficiency when the map size increases. Recent works exploited advanced deep learning techniques and generated global descriptors for more robust loop closure detection.

For example, PointNetVLAD [20] first extracts local features and then aggregates them into global descriptors. Different from methods that directly operate on raw point clouds, Zhou et al. [21] transform point clouds into NDT and combine them with Transformer [22] to generate descriptors. OverlapNet [1] introduces a deep learning-based method, which estimates the overlap and relative yaw angles of a set of point clouds for place recognition and initial pose estimation. Ma et al. [23]–[25] combine OverlapNet with Transformer to propose a rotation-invariant global descriptor. While the global descriptor-based methods can identify loop closures, they lack the ability to accurately estimate the 6-DoF pose between the current scan and the loop candidate.

On the other hand, local feature-based methods typically extract local sparse features from point clouds [26], [27] and organize them using a bag-of-words model for place recognition [6], [7]. These methods are capable for 6-DoF pose estimation based on local feature correspondences. However, extracting stable and reliable local features from 3D LiDAR scan is a challenging task, which limits the performance of such methods. Recently, LCDNet [9] employs a deep learning-based method PVRCNN [28] for robust feature extraction and then utilizes NetVLAD [29] for global descriptor generation. FinderNet [4] circumvents the challenge of feature extraction from point clouds by converting them into Digital Elevation Maps (DEMs) and leveraging a CNN network to extract local features and global descriptors. However, converting point clouds to DEM makes it challenging to estimate 6-DoF poses. In contrast, our method operates directly on point clouds, bypassing the challenge of estimating pose on sparse correspondences while possessing both loop closure detection and accurate 6-DoF pose estimation capability.

**Point cloud registration** refers to the process of determining the relative spatial transformation that aligns two point clouds. Extracting accurate correspondence is the most challenging aspect. Once correspondences are established, the transformation can be solved using either a direct solver or a robust estimator [30]. The iterative closest point (ICP) algorithm [17] and its various variants [31], [32] are known and applied methods. These methods establish correspondences iteratively using nearest neighbor search or other heuristics. However, the common drawback among ICP-like methods is that they heavily rely on good initial estimates for the transformation.

To release the requirement of initial estimates, other methods opt to establish correspondences on local features [27], [33] to achieve global registration. Due to the powerful feature representation capabilities exhibited by deep learning, massive learning-based methods for feature extraction have been proposed. Deng et al. [34], [35] propose PPFNet and PPF-FoldNet, which combine point pair features (PPF) with PointNet [36] to generate local patch representations for matching. In contrast to PPFNet and PPF-FoldNet, which establish correspondences on uniformly sampled points, keypoint-based techniques sample points based on pre-defined [27] or learned saliency [37]–[40] to achieve better repeatability. Due to the inherent errors introduced by individual matches, more matches usually lead to higher registration accuracy.

However, the aforementioned methods establish matches on sparse keypoints generated through uniform sampling [34], [35] or keypoint detection [37]–[40], which limits their registration accuracy. Recently, some studies [41]–[43] employ a coarse-to-fine mechanism that initially seeks correspondences on sparse keypoints and then extends them to dense ones, showing potential in registration. To enhance the reliability of sparse keypoint correspondences, CoFiNet [41] exploits transformer for contextual information aggregation and GeoTransformer [42] introduces geometric transformer to incorporate relative geometric information. RDMNet [43] introduces 3D-RoFormer for fast and lightweight relative geometric information encoding and the voting scheme for keypoint detection. HRegNet [44] extracts multi-level features and refines the transformation hierarchically. Despite the rapid advancements in global registration methods, these studies have remained disconnected from the challenges of loop closing and relocalization, lacking the ability of similarity evaluation.

**Loop Closing and Relocalization** both need to first find a coarse location and then estimate the fine 6-DoF pose. Though the individual techniques are widely explored, as discussed in the previous review, few work can address both tasks simultaneously. OverlapNet [3] and FinderNet [4] are capable of estimating 1-DoF or 3-DoF poses while implementing loop closure detection. These methods have also shown success in achieving 6-DoF loop correction when combined with other local registration techniques. However, in more challenging scenarios involving large pose drift or relocalization, non-6-DoF pose estimation is insufficient. BoW3D [7] and LCD-Net [9] achieve loop closure detection and 6-DoF pose estimation by extracting local features. Nevertheless, the accuracy of sparse local feature-based registration is constrained, requiring additional refinement through local registration techniques such as ICP. Therefore, these methods are evidently unsuitable for relocalization tasks where local registration have already failed and rapid online pose recovery are required.

To address system degeneration, most approaches integrate additional sensors such as camera [13], IMU [14], or ultra-wideband (UWB) [45], and switch between different tracking modalities. However, to the best of our knowledge, no piror LiDAR-only method has been proposed to achieve relocalization handling system degeneracy. This can be attributed to the challenge in achieving accurate LiDAR-based global registration when local pose tracking has already failed.

In the field of visual SLAM, however, due to the success in visual feature extraction, loop closing and relocalization have been well-explored. ORB-SLAM [46] extracts ORB features to form a BoW model for loop closing. When local tracking fails, ORB-SLAM switches from frame-to-frame alignment to frame-to-map alignment to find more potential landmark matches for relocalization. OA-SLAM [47] achieves more robust relocalization by relocalizing with reconstructed objects instead of local landmarks. Our approach, however, seeks to generate global descriptors to achieve rapid similarity evaluation for candidate detection in loop closing and relocalization. Additionally, a robust and accurate enough global registration method is employed to achieve 6-DoF pose estimation.
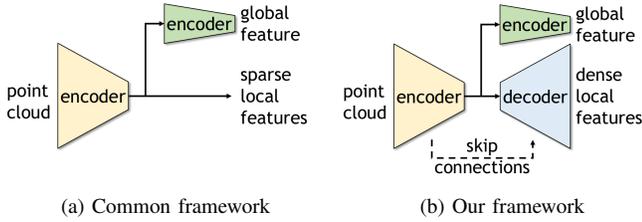
Fig. 2: Comparison of common framework and our framework. Our framework introduces decoder with skip connections for denser local feature generation.

## III. PROBLEM DEFINITION

We aim to address the challenges of loop closing and relocalization for LiDAR-based SLAM in outdoor driving environments. The underlying techniques of relocalization and loop closing are similar: both tasks involve the identification of the most similar candidate scan from the existing map and subsequently determining the 6-DoF pose. This commonality provides a foundation for addressing both tasks within a unified framework. However, the technical focus of the two tasks is quite different. Firstly, in the case of loop closing, the main challenge lies in rapidly and accurately identifying loop closures within a large database. In most cases, loop closures are identified when there is a substantial overlap between the current and candidate scans, simplifying the subsequent registration process once the loop closure has been correctly identified. Conversely, selecting the candidate scan for relocalization is relatively straightforward. In many autonomous driving situations, simply opting for the most recent scans can be sufficient. Instead, the primary challenge of relocalization lies in achieving precise and rapid global registration, as local pose tracking even with the aid of prior information has failed in such situations. This typically occurs in challenging scenarios that involve low overlap, extensive occlusions, or degraded scene features, posing significant challenges for point cloud registration. Secondly, loop closing can be executed at a relatively low frequency as a few correctly closed loops are sufficient to eliminate accumulated error. However, relocalization needs to be fast as it directly affects online localization. A longer relocalization process results in reduced overlap between the current scan and the map, decreasing the success rate of relocalization. While numerous works have focused on loop closing, the demanding requirements of robustness, accuracy, and speed in registration could be the primary reason for the limited attention to relocalization. To address this issue, we aim to initially study the framework to support the requirements of both tasks.

A commonly employed framework for simultaneously similarity evaluation and registration is shown in Fig. 2a. For an incoming LiDAR scan $\mathcal{P} = \{\boldsymbol{p}_i \in \mathbb{R}^3\}_{i=1}^N$, it initially downsamples and encodes the point cloud into local features $[\hat{\mathcal{P}}|\hat{\mathbf{F}}] = f_{\text{Encoder}}(\mathcal{P})$ and then generate a global descriptor $V = f_{\text{Encoder2}}(\hat{\mathcal{P}}, \hat{\mathbf{F}})$ based on these local features. The global descriptor $V$ is exploited to exhibit similarity for candidate frame retrieving and the local features $[\hat{\mathcal{P}}|\hat{\mathbf{F}}]$ are matched for pose estimation. Such a framework encounters a dilemma: A deeper encoder is often required to obtain reliable global feature representations. Nevertheless, a deeper encoder

tends to yield a reduced number of local features, potentially undermine registration performance. Conversely, enhancing registration performance often entails the incorporation of more complex designs into the encoder, consequently reducing global description efficiency. This is a crucial consideration in candidate retrieval tasks. Based on this insight, we propose the framework shown in Fig. 2b. We leave the workflow of global descriptor generation untouched but incorporate a decoder with skip connections for denser local feature generation $\mathbf{F} = f_{\text{Decoder}}(\hat{\mathcal{P}}|\hat{\mathbf{F}})$. Though simple, this resolves the conflict between the requirements of the two tasks. The incorporation of dense local features has the potential to enhance registration performance by establishing more correct matches, while ensuring the efficient generation of global descriptors. However, maintaining the match quality in an increased search space can be difficult and time-consuming. We thereby have implemented a sparse-to-dense matching approach to improve the matching process for reliable and fast registration. By exploiting different types of features and a multi-head network, we concurrently address the disparities between loop closing and relocalization tasks while leveraging their shared characteristics to unify them within a single framework. More detailed description of the proposed network following our framework are introduced in the next section.

## IV. LCR-NET: LOOP CLOSING AND RELOCALIZATION NETWORK

To realize our proposed framework, we design a novel multi-head network, named LCR-Net. As shown in Fig. 3, it consists of a keypoint detection module (Sec. IV-A) to extract keypoints from the raw point cloud, a global description head (Sec. IV-B) for global descriptor generation and a dense point matching head (Sec. IV-C) for local feature generation and matching. The devised loss function and the training strategy of our approach are detailed in Sec. IV-D and Sec. IV-E respectively.

### A. Keypoint detection module

The keypoint detection module aims to downsample the point cloud into sparse keypoints for further processing in two heads. In this work, we utilize KPEncoder [48] as the starting point for extracting the features. KPEncoder comprises a series of downsampling and kernal point-based convolution (KP-Conv) blocks, enabling hierarchical encoding of the point cloud into the uniformly distributed keypoints with descriptors $[\hat{\mathcal{P}}|\hat{\mathbf{F}}]$. These features provide sufficient information about the overall structure of the point cloud and are well-suited for input into the global description head. However, the uniformly sampled keypoints can not satisfy the demand for accurate registration due to their limited repeatability and saliency. They suffer from a lack of information exchange between two scans. To address these limitations and enhance feature matching, we introduce the *3D-RoFormer++* to reason about contextual information in both point clouds. In addition, we propose a new *VoteEncoder* that shifts the keypoints to nearby significant regions based on enhanced features and generates the final
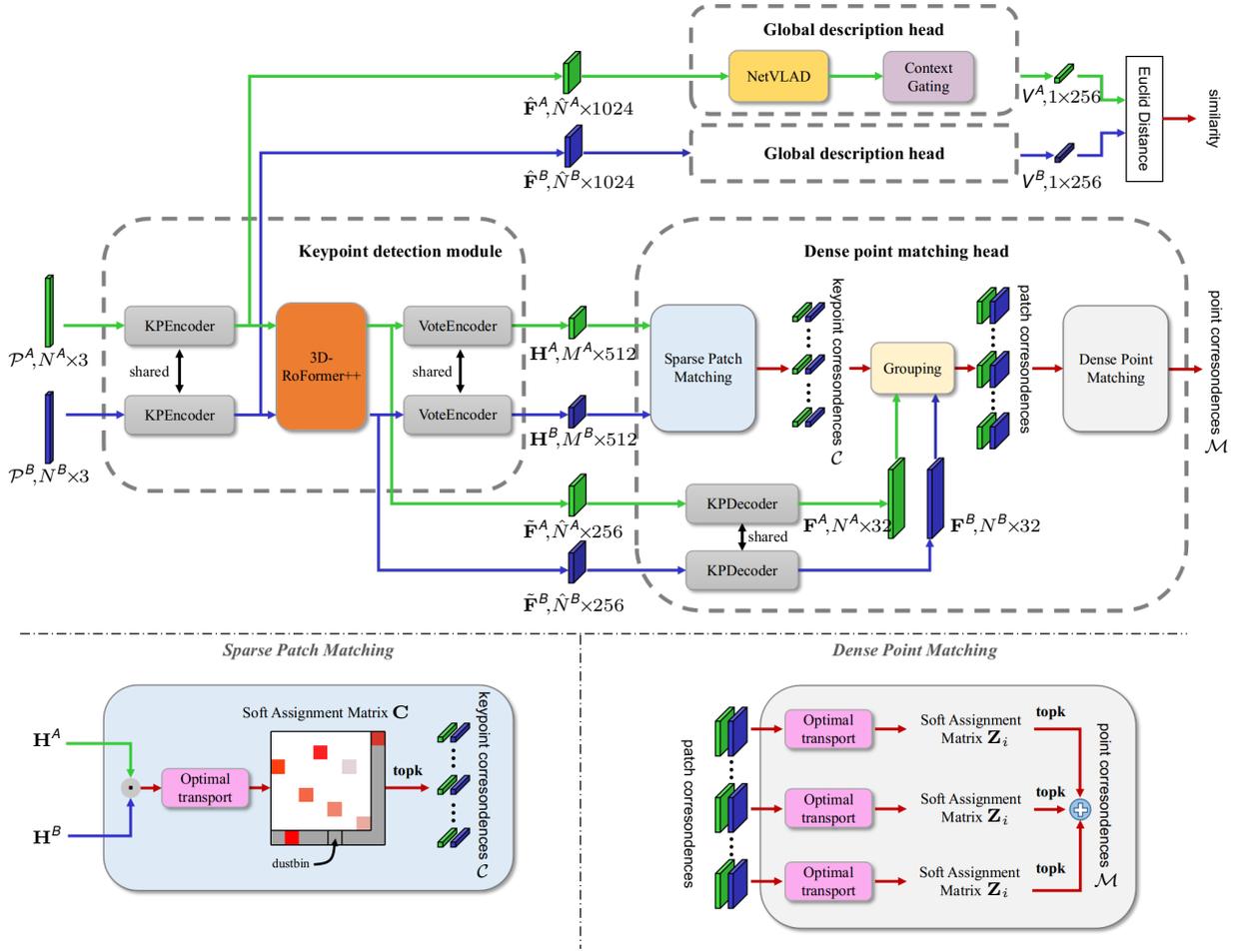
Fig. 3: **Pipeline overview.** LCR-Net consists of three main components: a keypoint detection module, a global description head, and a dense point matching head. The keypoint detection module extracts three types of features for further processing in two heads. The global description head generates a global descriptor for fast candidate retrieval. The dense point matching head exploits a sparse-to-fine approach to establish dense point matching for 6-DoF pose estimation.

keypoints by predicting the center points. We provide detailed descriptions of each component below.

**3D-RoFormer++.** In our previous work, the 3D-RoFormer [43] is introduced for lightweight relative pose-aware contextual aggregation. In this article, we have brought 3D-RoFormer to maturity and present the 3D-RoFormer++ by providing valuable translational invariance and enhanced feature representation performance. The 3D-RoFormer is built upon the vanilla transformer [22]. For a point $p_i^Q$ with its feature $h_i^Q$ in the query point cloud $Q$ and all the points in the source point cloud $S$, the transformer computes the query $q_i$, key $k_j$, and value $v_j$ feature maps with linear projections:

$$
\begin{aligned}
q_i &= W_1 f_i^Q + b_1, \\
k_j &= W_2 f_j^S + b_2, \\
v_j &= W_3 f_j^S + b_3.
\end{aligned} \tag{1}
$$

If $Q, S$ represent the same point cloud, Eq. (1) generates the feature maps for self-attention operation, otherwise cross-attention. In addition to the contextual features, 3D-RoFormer encodes the position $\hat{p}_i \in \mathbb{R}^3$ into the rotary embedding

$\Theta_i = [\theta_1, \theta_2, \cdots, \theta_{d/2}] \in \mathbb{R}^{\frac{d}{2}}$:

$$
\begin{aligned}
\Theta_i &= f_{\text{rot}}(\hat{p}_i)) \tag{2} \\
&= 2\pi \cdot \text{sigmoid}(\text{MLP}(\hat{p}_i)). \tag{3}
\end{aligned}
$$

By treating each element in $\Theta_i$ as a rotation in a 2D plane, it can be converted to a rotation matrix formulation $R_{\Theta_i} \in \mathbb{R}^{d \times d}$:

$$
R_{\Theta_i} = \begin{bmatrix}
\cos\theta_1 & -\sin\theta_1 & \cdots & 0 & 0 \\
\sin\theta_1 & \cos\theta_1 & \cdots & 0 & 0 \\
\vdots & \vdots & \ddots & \vdots & \vdots \\
0 & 0 & \cdots & \cos\theta_{\frac{d}{2}} & -\sin\theta_{\frac{d}{2}} \\
0 & 0 & \cdots & \sin\theta_{\frac{d}{2}} & \cos\theta_{\frac{d}{2}}
\end{bmatrix}. \tag{4}
$$

Applying $R_{\Theta_i}$ and $R_{\Theta_j}$ to query $q_i$ and key $k_j$ respectively in self-attention operation, the rotary self-attention in 3D-RoFormer can be written as:

$$
\alpha_{ij}'' = \text{softmax}_j((R_{\Theta_i} q_i)^\top R_{\Theta_j} k_j), \tag{5}
$$

$$
\tilde{f}_i = \sum_{j=1}^{|\hat{\mathcal{P}}|} \alpha_{ij}'' v_j. \tag{6}
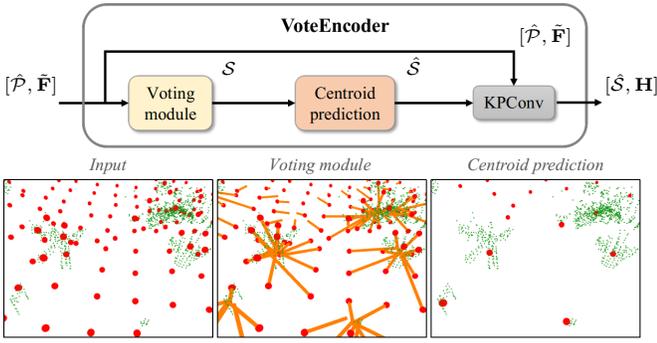$$

Fig. 4: The VoteEncoder takes sparse features as input, and generates offset from each keypoint to its nearest significant region. The final keypoints are estimated using centroid prediction algorithm, and their features are aggregated by KPConv. The ground points are removed from the visualization for clarity.

The Eq. (5) can be further written as:

$$\alpha''_{ij} = \text{softmax}_j(\boldsymbol{q}_i^\top R_{\Theta_i}^\top R_{\Theta_j} \boldsymbol{k}_j),$$
$$= \text{softmax}_j(\boldsymbol{q}_i^\top R_{\Theta_j - \Theta_i} \boldsymbol{k}_j). \tag{7}$$

The important advantage of 3D-RoFormer is that it explicitly encodes the relative geometric information neatly without requiring extra-large storage memory for relative position embedding. As in Eq. (7), relative "rotation" $\Theta_j - \Theta_i$ is naturally incorporated into the calculation and then fused with the output feature $\tilde{\boldsymbol{f}}_i$ in Eq. (6). Furthermore, if the mapping function $f_{\text{rot}}$ is linear, we can further derive:

$$\Theta_j - \Theta_i = f_{\text{rot}}(\hat{\boldsymbol{p}}_j - \hat{\boldsymbol{p}}_i). \tag{8}$$

This leads to a very important property for keypoint detection which is translation-invariance. However, designing a $f_{\text{rot}}$ that provides good rotary feature representation while maintaining linearity is challenging. To ensure the ability of rotary representation, the rotary embedding in the original 3D-RoFormer, as shown in Eq. (3), sacrifices linearity for rotary representation, leading to reduced generalization performance.

Based on this insight, we improve our 3D-RoFormer by adopting a learning-based linear mapping function:

$$\Theta_i = \text{Linear}(\hat{\boldsymbol{p}}_i), \tag{9}$$

with a boundary penalty loss (Sec. IV-D) as an auxiliary loss to supervise the network actively learning effective rotary representations. With this modification, our 3D-RoFormer++ significantly enhances the final output features $\tilde{\mathbf{F}}^A$ and $\tilde{\mathbf{F}}^B$ for point matching by interleaving the rotary self-attention and cross-attention for $l$ times.

The enhanced features $\tilde{\mathbf{F}}$ possess geometric and contextual information between two point clouds, which is then extended to dense features for further processing in dense point matching head, as detailed in Sec. IV-C. Nevertheless, the uniform sampling nature makes these features less salient and discriminative. We therefore propose the VoteEncoder.

**VoteEnoder**. To steer the evenly sampled features $[\hat{\mathcal{P}}|\tilde{\mathbf{F}}]$ towards nearby salient areas and obtain more meaningful features conducive to registration tasks, we introduce the VoteEncoder for additional feature encoding.

We use a voting module [43], [49] to estimate the geometric offset from the uniformly sampled keypoints to the proposal

---

**Algorithm 1** Centroid prediction

**Input:** proposal set $\mathcal{S}$, nearest neighbour search range $d$
**Output:** center point set $\hat{\mathcal{S}}$
1: **for** all $\boldsymbol{s}_i \in \mathcal{S}$ **do**
2:      **if** ISLABELED($\boldsymbol{s}_i$) is False **then**
3:          $\mathcal{N}_i \leftarrow$ NEARESTNEIGHBOUR($\boldsymbol{s}_i, \mathcal{S}, d$)
4:          $\hat{\boldsymbol{s}}_i \leftarrow$ MEAN($\mathcal{N}_i$)     ▷ Get centroid point
5:          $\hat{\mathcal{S}}$.APPEND($\hat{\boldsymbol{s}}_i$)
6:          **for** all $\boldsymbol{s}_j \in \mathcal{N}_i$ **do**
7:             ISLABELED($\boldsymbol{s}_j$) $\leftarrow$ True
8: **return** $\hat{\mathcal{S}}$

---

keypoints $\mathcal{S}$, i.e., $\Delta\mathbf{P} = \text{Vote}(\tilde{\mathbf{F}})$, $\mathcal{S} = \hat{\mathcal{P}} + \Delta\mathbf{P}$. The voting module comprises a collection of Multi-Layer Perceptrons (MLPs). Despite its simplicity, this module produces meaningful offsets (see Fig. 4), utilizing the features from our 3D-RoFormer++. These generated proposals subsequently forecast multiple central points, serving as the final keypoints. The process of predicting centroids is straightforward yet efficient, without the need for additional sampling strategies. It clusters all the proposals into various patches and predicts the centers, detailed in Alg. 1. To aggregate the descriptors $\mathbf{H}$ for each center point $\hat{\mathcal{S}}$, we employ a KPConv module that performs kernel-based convolution after finding nearest neighbors of $\hat{\mathcal{S}}$ in $[\hat{\mathcal{P}}|\tilde{\mathbf{F}}]$. We use a larger search range than that used in Alg. 1 to incorporate more related context near the keypoints.

Unlike the object detection [49], [50], where the object center can serve as a well-defined reference for supervising point shifts, our case lacks a readily available ground truth center for the significant areas, primarily due to the challenge in precisely defining the significance. Therefore, we train the matched keypoints to move closer to each other instead, which indirectly accomplishes our objective. In practice, the offset $\Delta\mathbf{P}$ is limited to a certain range to maintain an even distribution of keypoints throughout the point cloud. This prevents the keypoints from being only concentrated in significant areas while also avoiding potential degeneracy.

In sum, our keypoint detection module offers various options for the sparse features, including uniformly sampled features $\hat{\mathbf{F}}$ from KPEncoder [48], enhanced features $\tilde{\mathbf{F}}$ from 3D-RoFormer++ and voted features $\mathbf{H}$ from VoteEncoder. Uniformly sampled features $\hat{\mathbf{F}}$ are distributed evenly throughout the entire point cloud, enabling fast extraction and the capacity to represent the entire point cloud comprehensively. These features are utilized for the global description head, as detailed in Sec. IV-B. On the other hand, enhanced features $\tilde{\mathbf{F}}$, built upon the uniformly sampled features $\hat{\mathbf{F}}$, incorporate geometric information and the correlation between two scans. These features will be decoded into dense features, propagating the aforementioned advantages to them. Finally, voted features $\mathbf{H}$ exhibit sensitivity and expressive capability for local salient regions while maintaining good coverage. These features will be employed for initial sparse matching, which will then be extended to dense matching for accurate registration. Both enhanced features and voted features will be employed for the dense point matching head detailed in Sec. IV-C.

## B. Global Description Head

The objective of the global description head is to condense sparse features into a single global feature for fast candidate retrieval. We adopt the features derived from the KPEncoder as the input to our global description head.

When creating the global description head, we choose a widely applied simple method NetVLAD [29] to compress the features, and a context gating module [9] to enhance these features for retrieval. NetVLAD uses k-means clustering and defines $K$ learnable cluster centers $\{\boldsymbol{c}_1, \cdots, \boldsymbol{c}_K\}, \boldsymbol{c}_k \in \mathbb{R}^{\hat{d}}$, along with learnable weights $\boldsymbol{w}_k$ and offsets $b_k$. By weighting each feature to each cluster center:

$$a_k(\hat{\mathbf{F}}_i) = \frac{e^{\boldsymbol{w}_k^\top \hat{\mathbf{F}}_i + b_k}}{\sum_{k'=1}^{K} e^{\boldsymbol{w}_{k'}^\top \hat{\mathbf{F}}_i + b_{k'}}}, \quad (10)$$

the new descriptors for $K$ cluster centers are obtained:

$$\mathbf{FR} = [\mathbf{FR}_1, \cdots, \mathbf{FR}_K] \in \mathbb{R}^{K \times \hat{d}}, \quad (11)$$

$$\mathbf{FR}_K = \sum_{i=1}^{|\hat{\mathcal{P}}|} a_k(\hat{\mathbf{F}}_i)(\hat{\mathbf{F}}_i - \boldsymbol{c}_k). \quad (12)$$

Finally, a simple MLP compresses $\mathbf{FR}$ into a single descriptor $X \in \mathbb{R}^G$. Based on NetVLAD, the context gating module re-evaluates the weights of each channel of feature $X$ based on the self-attention mechanism and further enhances it to obtain the final global descriptor $V \in \mathbb{R}^G$:

$$V = \mathrm{CG}(X) = \sigma(WX + \boldsymbol{b}) \otimes X, \quad (13)$$

where $\sigma$ is the sigmoid activation function, $\otimes$ is element-wise multiplication, and $W$ and $\boldsymbol{b}$ are learnable weights and offsets.

The design of the global description head is straightforward yet remarkably effective, surpassing all baseline methods in our experiments. Its simplicity is also particularly important because LiDAR SLAM requires quick and accurate retrieval for real-time candidate retrieval, which narrows down the computational scope for subsequent fine 6-DoF pose estimation.

## C. Dense Point Matching Head

Once identifying the candidates, we leverage the dense point matching head to establish correspondences and subsequently recover precise 6-DoF pose estimation. The features obtained from our keypoint detection module $[\hat{\mathcal{S}}|\mathbf{H}]$ are sufficient for ensuring dependable point cloud registration. However, there are two factors that impact the accuracy of the final 6-DoF pose estimation. Firstly, despite VoteEncoder improving keypoint locations, there might still be noticeable distances between matched sparse features. These gaps can lead to errors that restrict the overall accuracy. Secondly, due to the sparse characteristics of these features, there might not be adequate feature matches to fully rectify errors arising from mismatches. Considering these limitations, we employ a two-step matching approach [41], [42]. Initially, we identify sparse yet dependable keypoint matches, and then we extend these point-to-point matches to patch-to-patch matches. By utilizing neighbor consistency, we enhance these patch matches into dense point matches, ensuring more precise and reliable registration.

**Sparse keypoint matching.** We firstly conduct sparse keypoint matching between $[\hat{\mathcal{S}}^A|\mathbf{H}^A]$ and $[\hat{\mathcal{S}}^B|\mathbf{H}^B]$. We compute a matching score matrix $\mathbf{C} \in \mathbb{R}^{|\hat{\mathcal{S}}^A| \times |\hat{\mathcal{S}}^B|}$ between $\mathbf{H}^A$ and $\mathbf{H}^B$:

$$\mathbf{C} = \mathbf{H}^A(\mathbf{H}^B)^\top / \sqrt{d_c}, \quad (14)$$

where $d_c$ refers to the feature dimension of $\mathbf{H}$. To handle non-matched points, we append a "dustbin" row and column for $\mathbf{C}$ filled with a learnable parameter $\alpha \in \mathbb{R}$. The Sinkhorn algorithm [51] is then used to solve the soft assignment matrix. It iteratively performs normalization along rows and columns. At the $t$ iteration, the score matrix is updated by:

$$^{(t)}\mathbf{C}'_{ij} = {}^{(t)}\mathbf{C}_{ij} - \log \sum_j e^{(t)\mathbf{C}_{ij}}, \quad (15)$$

$$^{(t+1)}\mathbf{C}_{ij} = {}^{(t)}\mathbf{C}'_{ij} - \log \sum_i e^{(t)\mathbf{C}'_{ij}}. \quad (16)$$

After $T$ iterations, we use the solution as the soft assignment matrix: $\hat{\mathbf{C}} = {}^{(T)}\mathbf{C}$. We choose the largest $N_c$ entries as the keypoint correspondences:

$$\mathcal{C} = \{(\hat{\boldsymbol{s}}_{x_i}^A, \hat{\boldsymbol{s}}_{y_i}^B)|(x_i, y_i) \in \mathrm{Top\text{-}k}_{x,y}(\hat{\mathbf{C}})\}. \quad (17)$$

**Patch grouping.** To achieve dense point matches from sparse keypoint matches, we expand correspondences between keypoints to encompass overlaps between their respective neighborhood patches and subsequently leverage these patches to identify more point matches.

For each keypoint $\hat{\boldsymbol{s}}_i$, we construct a local patch $\mathcal{G}_i$ using a point-to-node strategy [38], where each point is assigned to its nearest keypoint. Based on the grouped point patch, we can now extend each keypoint match $(\hat{\boldsymbol{s}}_{x_i}^A, \hat{\boldsymbol{s}}_{y_i}^B)$ to its corresponding patch match $(\mathcal{G}_{x_i}^A, \mathcal{G}_{y_i}^B)$.

**Dense point matching.** We then generate more point matches from the sparse patch matches. We leverage the KPDecoder [48] to recover point-level descriptors $\mathbf{F}$ from enhanced keypoint features $[\hat{\mathcal{P}}|\tilde{\mathbf{F}}]$. For each keypoint correspondence $(\hat{\boldsymbol{s}}_{x_i}^A, \hat{\boldsymbol{s}}_{y_i}^B)$, we compute a match score matrix $\mathbf{O}_i \in \mathbb{R}^{|\mathcal{G}_{x_i}^A| \times |\mathcal{G}_{y_i}^B|}$ of their corresponding patches $\mathcal{G}_{x_i}^A$ and $\mathcal{G}_{y_i}^B$:

$$\mathbf{O}_i = \mathbf{F}_{x_i}^A(\mathbf{F}_{y_i}^B)^\top / \sqrt{d_f}, \quad (18)$$

where $d_f$ refers to the feature dimension of $\mathbf{F}$. Same with our sparse keypoint matching module, we append a learnable "dustbin" row and column for $\mathbf{O}_i$ to handle non-matched points and use the sinkhorn algorithm to solve the soft assignment matrix $\mathbf{Z}_i \in \mathbb{R}^{(|\mathcal{G}_{x_i}^A|+1) \times (|\mathcal{G}_{y_i}^B|+1)}$. Unlike works [41], [42] that drops the dustbin and recovers the assignment by comparing the soft assignment score with a hand-tuned threshold, we directly find max entry both row-wise and column-wise on $\mathbf{Z}_i$ which is then recovered to assignment $\mathcal{M}_i$:

$$\mathcal{M}_i = \{(\mathcal{G}_{x_i}^A(m), \mathcal{G}_{y_i}^B(n))|(m, n) \in \mathrm{toprow}_{m,n}(\mathbf{Z}_{1:M_i, 1:(N_i+1)}^i)\} \cup$$
$$\{(\mathcal{G}_{x_i}^A(m), \mathcal{G}_{y_i}^B(n))|(m, n) \in \mathrm{topcolumn}_{m,n}(\mathbf{Z}_{1:(M_i+1), 1:N_i}^i)\}. \quad (19)$$

A point is either assigned to points in the matched patch or to the dustbin. By this, we do not need manual tuning but require a discriminative assignment matrix, which can be obtained by using our proposed loss function as detailed in

Sec. IV-D. Note that a point is not strictly assigned to a single point in our approach, as the strict one-to-one point correspondences do not hold in practice due to the sparsity nature of the LiDAR scans. Instead, we trust and keep the assignment results from both sides, i.e., matches from query to source and vice versa. This results in extensively more point matches while maintaining a high inlier ratio, which benefits the transformation estimation. The final correspondences are the combination of points matches from all patches:

$$\mathcal{M} = \bigcup_{i=1}^{N_c} \mathcal{M}_i. \tag{20}$$

**Local-to-global registration.** We use local-to-global registration (LGR) proposed in [42] for fast pose estimation. It is a hypothesize-and-verify approach specifically proposed for matching methods following a sparse-to-dense manner. For each matched patch, LGR solves a transformation $\{R_i, t_i\}$ based on its dense point matches using weighted SVD [17]:

$$R_i, t_i = \min_{R,t} \sum_{(p_{x_j}^A, p_{y_j}^B) \in \mathcal{M}_i} \omega_j^i \| R \cdot p_{x_j}^A + t - p_{y_j}^B \|_2^2, \tag{21}$$

where the soft assignment value in $\mathbf{Z}^i$ serves as the weight $\omega_j^i$. After obtaining the transformations for all matched patches, LGR selects the transformation that has the most inliers among all dense point matches:

$$R, t = \max_{R_i, t_i} \sum_{(p_{x_j}^A, p_{y_j}^B) \in \mathcal{M}} [\![ \| R_i \cdot p_{x_j}^A + t_i - p_{y_j}^B \|_2^2 < \tau_a ]\!], \tag{22}$$

where $[\![ \bullet ]\!]$ is an indicator function for which the statement is true. Finally, it solves the final transformation $R, t$ by solving Eq. (21) on surviving inliers for $N_r$ times.

LGR significantly reduces the number of iterations compared to RANSAC [30], achieving a substantial speed advantage with about 30 times faster in our experiments. However, the performance of LGR, particularly its robustness, can be heavily influenced by the quality of sparse patch matching. We significantly improve the matching quality of sparse patches through the powerful feature aggregation module 3D-RoFormer++ and the feature detection module VoteEncoder, achieving performance comparable to or surpassing RANSAC's accuracy and robustness.

### D. Loss function

To effectively guide our network in accomplishing various tasks, we construct our loss function with five components: the keypoint detection loss $L_s$, the boundary penalty loss for keypoint detection module, the triplet loss $L_t$ for global description head, and the sparse match loss $L_c$ and the dense match loss $L_f$ for dense point matching head.

**Keypoint detection loss.** The keypoint detection loss consists of two parts $L_s = L_{s1} + L_{s2}$. The first part $L_{s1}$ is designed to guide the corresponding keypoints from two point clouds close to each other lying within the significant region:

$$L_{s1} = \sum_{i=1}^{|\mathcal{S}^A|} \min_{s_j^B \in \mathcal{S}^B} \| s_i^A - s_j^B \|_2^2 + \sum_{i=1}^{|\mathcal{S}^B|} \min_{s_j^A \in \mathcal{S}^A} \| s_i^B - s_j^A \|_2^2. \tag{23}$$

Supervised by $L_{s1}$, we find that the keypoints tend to move to their nearest significant regions to indirectly minimize the distance between keypoint pairs.

The second part $L_{s2}$ is designed to make the keypoints close to the real measurement points. It minimizes the distance between the keypoint with its closest point:

$$L_{s2} = \sum_{i=1}^{|\mathcal{S}^A|} \min_{p_j^A \in \mathcal{P}^A} \| s_i^A - p_j^A \|_2^2 + \sum_{i=1}^{|\mathcal{S}^B|} \min_{p_j^B \in \mathcal{P}^B} \| s_i^B - p_j^B \|_2^2. \tag{24}$$

**Boundary penalty loss.** To guide the 3D-RoFormer++ learning a general rotary embedding representation, we add a boundary penalty loss to force the value of rotary embedding $\Theta$ lies between $[-\pi, \pi]$:

$$L_p^i = \frac{1}{M_i} \sum_{m=1}^{M_i} [\mathrm{abs}(\Theta) - \pi]_+. \tag{25}$$

**Triplet loss.** We use the triplet loss to train the global description head. For each scan, we define the scans with an overlap greater than 30% as positive, otherwise negative. For each triplet, we use one query scan, $N_p$ positive scans and $N_n$ negative scans. The triplet loss is calculated as:

$$L_t(V_q, \{V_p\}, \{V_n\}) = \tag{26}$$
$$N_p(\alpha + \max_p (d(V_q, V_p)) - \frac{1}{N_n} \sum_{N_n} (d(V_q, V_n))).$$

**Sparse match loss.** We utilize a gap loss [40] to learn a discriminative soft assignment matrix $\mathbf{C}$ for sparse keypoint matching. The ground truth of assignment matrix $\mathbf{P} \in \{0, 1\}^{(M+1) \times (N+1)}$ is generated based on the overlap ratio between the patches, where $M = |\hat{\mathcal{S}}^A|$ and $N = |\hat{\mathcal{S}}^B|$ are the keypoints number. Two patches are matched when they share at least 10% overlap. A patch is assigned to the dustbin when it has no match pair. We also generate a negative assignment matrix $\bar{\mathbf{P}} \in \{-\inf, 1\}^{(M+1) \times (N+1)}$, where 1 represents two patches are not overlapped and $\inf$ represents the value will not involve in the calculation of loss. Then the gap loss is calculated as:

$$L_c = \frac{1}{M} \sum_{m=1}^{M} \log\left(\sum_{n=1}^{N+1} (-r_m + \mathbf{C}_{m,n} \bar{\mathbf{P}}_{m,n} + \eta)_+ + 1\right)$$
$$+ \frac{1}{N} \sum_{n=1}^{N} \log\left(\sum_{m=1}^{M+1} (-c_n + \mathbf{C}_{m,n} \bar{\mathbf{P}}_{m,n} + \eta)_+ + 1\right), \tag{27}$$

where

$$r_m = \max_m(\mathbf{C}_{m,n} \mathbf{P}_{m,n}), \quad c_n = \max_n(\mathbf{C}_{m,n} \mathbf{P}_{m,n}), \tag{28}$$

refer to the soft assignment value for the hardest true match of $m$-th keypoint in $\hat{\mathcal{S}}^A$ and $n$-th keypoint in $\hat{\mathcal{S}}^B$ respectively, and $(\bullet)_+ = \max(\bullet, 0)$.

**Dense match loss.** The dense match loss is calculated over all the matched patches. For each matched patch pair $\{\mathcal{G}_{x_i}^A, \mathcal{G}_{y_i}^B\}$, we generate its ground truth positive correspondences matrix $\mathbf{M}^i \in \{0, 1\}^{(M_i+1) \times (N_i+1)}$ and negative matrix $\bar{\mathbf{M}}^i \in \{10^{12}, 1\}^{(M_i+1) \times (N_i+1)}$ with a distance threshold $\tau$, where $M_i = |\mathcal{G}_{x_i}^A|$, $N_i = |\mathcal{G}_{x_i}^B|$. A point pair is positive when

the distance is below $\tau$ and is negative when it exceeds $2\tau$. To learn a discriminative soft assignment matrix, we also calculate a gap loss $L_f^i$ for patch correspondence's soft assignment matrix $\mathbf{Z}^i$. The final fine match loss is the average over all the matched patch pairs: $L_f = \frac{1}{2|\mathcal{M}|} \sum_{i=1}^{|\mathcal{M}|} L_f^i$.

*E. Training Strategy*

We seek a training strategy to stimulate the potential of each head with limited computing resources. As a result, a two-stage training strategy is employed. We first train the keypoint detection module and the dense point matching head for registration. Then, we exclusively train the global description head for the candidate retrieval for the following reasons:

Firstly, the input for the training of two heads differs. The training of the global description head requires at least three scans: an anchor, a positive, and a negative. Conversely, training the dense point matching head only requires two overlapped scans. Including additional input does not benefit the training of the dense point-matching head but consumes a significant amount of memory. Secondly, for tasks of candidate retrieval, a higher batch size typically results in better performance. By freezing the keypoint detection module and dense point matching head, we can use the pre-extracted features for input, thereby preserving substantial memory for expanding the batch size.

However, using pre-extracted features prevents the implementation of data augmentation techniques, thereby limiting performance. To address this problem, we utilize a training strategy, which we refer to as semi-online. In this approach, we utilize offline pre-extracted features for both positive and negative samples while generating features for the anchor online. This allows for applying data augmentation on the anchor. Since the anchor participates in loss calculations with all positive and negative samples, this can be the most efficient way to implement data augmentation.

In sum, our two-stage training strategy first trains the network in the registration task and then exclusively trains the global description head semi-online for the candidate retrieval task. This strategy offers several advantages. Firstly, it allows us to utilize larger batch sizes and sample quantities during training for the candidate retrieval task, leading to improved results. Secondly, the semi-online approach enables the application of data augmentation techniques. Lastly, this strategy provides great convenience in training as we only need to select the pre-trained network that performs best in registration for the second-stage training, followed by selecting the network that performs best in the candidate retrieval task. These advantages contribute to this training strategy's enhanced performance compared to fine-tuning on a pre-trained model and end-to-end training, as demonstrated in our experiments.

We implement and train our LCR-Net on 4 NVIDIA RTX 3090 GPUs. The network is trained using the Adam optimizer [52] with an initial learning rate as $10^{-4}$, which undergoes exponential decay by $0.05$ every $4$ epochs. When training the dense point matching head, we use a batch size of $1$. On the other hand, when training the global description
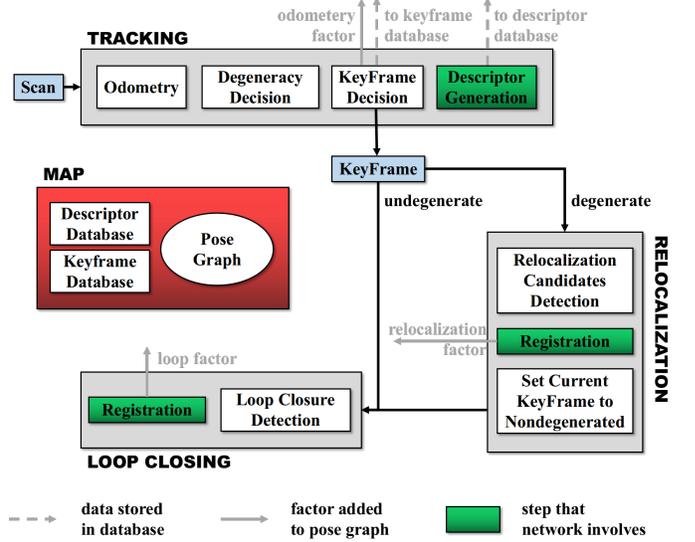


Fig. 5: Overview of the laser SLAM system integrated with LCR-Net, showing key steps executed by the tracking, relocalization and loop closing threads. LCR-Net serves for relocalization and loop closing.

head, we set the batch size to 6, utilizing $N_a = 1$ anchor scan, $N_p = 6$ positive scans, and $N_n = 6$ negative scans. Additionally, we apply the same data augmentation techniques as in [9].

## V. Robust Loop Closing and Relocalization Based LiDAR SLAM System

We integrate our LCR-Net into a SLAM system for relocalization and loop closing. We use an Incremental Smoothing and Mapping (iSAM2) [53] based pose-graph optimization (PGO) [1] framework to manage and optimize the global pose graph. The system comprises three threads that run in parallel: tracking, relocalization, and loop closing.

### A. Tracking

The tracking thread is responsible for tracking the LiDAR pose of every scan, i.e., odometry. It also decides whether the current odometry estimation is degenerated and when to insert a new keyframe. The process of odometry can be formulated as a state estimation problem:

$$\arg \min_{\boldsymbol{x}} f^2(\boldsymbol{x}), \qquad (29)$$

where $\boldsymbol{x} = [\mathsf{R}, \boldsymbol{t}]$ is the state vector. Given an initial guess of $\boldsymbol{x}$, most nonlinear optimization methods solve the function by computing the Jacobian matrix of $f$ w.r.t $\boldsymbol{x}$:

$$\mathsf{J} = \delta f(\boldsymbol{x})/\delta \boldsymbol{x}, \qquad (30)$$

and iteratively adjust $\boldsymbol{x}$ by utilizing $\mathsf{J}$ until convergence. The tracking thread can be implemented through LiDAR odometry, such as LOAM [12]. For degeneracy evaluation, we adopt the criterion proposed in [13], which considers the problem degenerate if the smallest eigenvalue of matrix $\mathsf{J}^{\mathsf{T}}\mathsf{J}$ falls below a certain threshold. In practice, the threshold is conservatively set to ensure that all degeneracy can be

[1] https://github.com/gisbi-kim/SC-A-LOAM.

detected. Though this leads to more frequent relocalization and increased computational load, it is necessary since even a single instance of severe degeneracy could inflict significant damage on the system.

LiDAR scans are keyframes if the robot has moved beyond a predefined threshold from the last saved keyframe. However, when odometry performance deteriorates, relying solely on the previous travel distance to establish keyframes becomes unreliable. To address this, we introduce additional scans as keyframes based on specific conditions as follows: i) The first degenerated scan follows nondegenerated scans. ii) The first nondegenerated scan follows a degenerated scan, provided the degenerated scan is not chosen as a keyframe. iii) Every $t$-second scan within a continuous sequence of degenerated scans. The keyframes selected according to these conditions are categorized as "degenerated", as poses estimated between these keyframes may include degenerated ones. Finally, for each keyframe, we use LCR-Net to generate a global descriptor. The similarity between our generated global descriptors is evaluated based on the Euclidean distance in the feature space. For fast retrieving, we utilize the FAISS library [54] for descriptor database management and search.

### B. Relocalization

The relocalization processes every "degenerated" keyframe, generating more accurate pose estimation to replace the unreliable odometry output. We select the most recent "nondegenerated" keyframe for each degenerated keyframe as the candidate and match it with the current keyframe using LCR-Net. For real-time efficiency, we utilize the rapid estimator LGR for pose estimation. Our experimental results show that employing LCR-Net with LGR produces comparable accuracy and robustness in registration compared to LCR-Net with RANSAC while offering a significantly faster processing speed of nearly 30 times. To assess whether the relocalization is successful, we calculate a current pose estimation reliability score by averaging the assignment score of all the found correspondences. The relocalization succeeds if the reliability score surpasses a threshold $\rho_r$. Once the calculated pose is included in the graph optimization, the label assigned to the current keyframe is updated to "nondegenerated". In most cases, this allows for recovery of the pose tracking. Otherwise, we retrieve the most similar candidate with descriptor distance below a threshold $\rho_s$ from the database and attempt relocalization.

### C. Loop closing

The loop closing thread functions by searching the loop candidate and subsequently estimating pose as a new node added to the pose graph. Specifically, we retrieve the candidate keyframe from the database that has the most similar descriptor for each new keyframe while excluding the 100 most recent keyframes. If the distance is below a threshold $\rho_s$, we set the retrieved keyframe as the loop and estimate the relative 6-DoF pose using LCR-Net with the LGR solver. Finally, the iSAM2-based pose graph optimization is performed to achieve global consistency.

## VI. EXPERIMENTAL EVALUATION

We conduct experiments to demonstrate the efficacy of our proposed LCR-Net in addressing loop closing and relocalization for online LiDAR SLAM. We derive three setups from these challenges: candidate retrieval, closed-loop point cloud registration, and continuous relocalization. In candidate retrieval experiments (Sec. VI-B), we examine the capability of LCR-Net in accurately retrieving the appropriate candidates from the previous map. In closed-loop point cloud registration (Sec. VI-C) and continuous relocalization experiments (Sec. VI-D), we assess the ability of LCR-Net in successfully registering point clouds during loop situations as well as continuous scenarios with low overlap. We also evaluate the runtime of our approach for candidate retrieving and point cloud registration (Sec. VI-E). We then evaluate our LCR-Net enhanced SLAM in multiple real-world scenes. Finally, we conduct ablation studies on the network design (Sec. VI-G) and provide valuable insights (Sec. VI-H).

### A. Experimental Setup

We evaluate LCR-Net and compare it with the SOTA methods on multiple publicly available datasets, including KITTI odometry [55], KITTI-360 [56], Apollo-SouthBay [57], Ford Campus [58] and Mulran [59] datasets. These datasets provide LiDAR scans collected in various environments in multiple countries with the corresponding ground truth poses.

In particular, we divide the KITTI odometry dataset into the following: sequences 01 and 03-07 for training, sequence 02 for validation, and sequences 00, 08-10 for testing. The other datasets are all used to test the models' generalization capabilities. We use the two-stage training strategy described in Sec. IV-E: pre-train the network with dense point matching head in the registration task and then linear probe the global description head for training in the candidate retrieval task. For the registration task training, we select point cloud pairs with a mix of continuous point cloud pairs with the distance varying from 0-10 m and close-loop point cloud pairs with an overlap ratio exceeding 0.3 [1]. In the candidate retrieval task training, we use all point cloud pairs with an overlap ratio exceeding 0.3. We denote the model trained following the above setup as *LCR-Net* and evaluate its performance on all subsequent tasks. This model is also utilized for integration with our SLAM system for evaluation.

In addition, to eliminate the impact of different training sets and ensure equitable comparisons with existing baselines, we also train two other models, denoted as *LCR-Net*[†] and *LCR-Net*[◇], using different dataset splittings that follow the baselines for specific tasks: *LCR-Net*[†] is pre-trained on close-loop point cloud pairs, while *LCR-Net*[◇] is pre-trained on continuous point cloud pairs. Both models are then trained in the candidate retrieval task. The subsequent experiments will delve into the specific training details for these two models.

### B. Candidate Retrieval Performance

To validate the candidate retrieval performance, we follow Chen et al. [1], [23] and test our approach on the KITTI odometry and Ford Campus dataset. For a fair comparison, we

TABLE I: Candidate Retrieval Results on KITTI and Ford Campus.

| Dataset | Method | AUC | F1max | Recall @1 | Recall @1% |
|---|---|---|---|---|---|
| KITTI | Histogram [60] | 0.826 | 0.825 | 0.738 | 0.871 |
| | Scan Context [18] | 0.836 | 0.835 | 0.820 | 0.869 |
| | LiDAR-Iris [19] | 0.843 | 0.848 | 0.835 | 0.877 |
| | PointNetVLAD [20] | 0.856 | 0.846 | 0.776 | 0.845 |
| | OverlapNet [1] | 0.867 | 0.865 | 0.816 | 0.908 |
| | NDT-Transformer-P [21] | 0.855 | 0.856 | 0.802 | 0.869 |
| | MinkLoc3D [61] | 0.894 | 0.869 | 0.876 | 0.920 |
| | OverlapTransformer [23] | 0.907 | 0.877 | 0.906 | 0.964 |
| | LCDNet | <u>0.933</u> | <u>0.883</u> | <u>0.915</u> | <u>0.974</u> |
| | LCR-Net† | **0.945** | **0.907** | **0.926** | **0.980** |
| | LCR-Net | 0.958 | 0.922 | 0.937 | 0.993 |
| Ford Campus | Histogram [60] | 0.841 | 0.800 | 0.812 | 0.897 |
| | Scan Context [18] | 0.903 | 0.842 | 0.878 | 0.958 |
| | LiDAR-Iris [19] | 0.907 | 0.842 | 0.849 | 0.937 |
| | PointNetVLAD [20] | 0.872 | 0.830 | 0.862 | 0.938 |
| | OverlapNet [1] | 0.854 | 0.843 | 0.857 | 0.932 |
| | NDT-Transformer-P [21] | 0.835 | 0.850 | 0.900 | 0.927 |
| | MinkLoc3D [61] | 0.871 | 0.851 | 0.878 | 0.942 |
| | OverlapTransformer [23] | 0.923 | 0.856 | 0.914 | 0.954 |
| | LCDNet | <u>0.961</u> | <u>0.908</u> | <u>0.949</u> | <u>0.984</u> |
| | LCR-Net† | **0.974** | **0.929** | **0.951** | **0.985** |
| | LCR-Net | 0.972 | 0.920 | 0.932 | 0.987 |

The best results are highlighted in bold, and the second best in underlines.

follow the setup of Chen et al. [1], [23] and train *LCR-Net†* on KITTI odomtery sequences 03-10 and validate it on KITTI odomtery sequence 02. Two scans are chosen as a candidate if their overlap value is larger than 0.3.

**Metrics** : We use four metrics to evaluate the performance of candidate retrieval: i) the area enclosed by the Receiver Operating Characteristic (ROC) curve and the coordinate axes; ii) Maximum F1 score (F1max), which is the highest F1 score at different threshold values. ; iii) Recall@1, which measures the recall when only the most similar candidate frame is selected; iv) Recall@1%, which measures the recall when the top 1% of the most similar candidate frames are selected.

**Results** : The baseline methods used in this experiment are SOTA place recognition methods. For Histogram [60], Scan Context [18], LiDAR-Iris [19], OverlapNet [1], Point-NetVLAD [20], NDT-Transformer-P [21], MinkLoc3D [61], and OverlapTranformer [23], we use the results reported in [23]. For LCDNet [9], we utilize its official implementation. As shown in Tab. I, LCR-Net and LCDNet achieve cutting-edge performance compared to current advanced methods. However, our method further boosts the baseline by a significant margin for the F1max metric while maintaining a leading position in Recall@1, Recall@1%, and AUC. It is worth noting that our global description head does not employ complex designs. The exceptional performance can be attributed to our backbone's outstanding feature representation capabilities and our training strategy, which permits larger batch sizes and more input samples.

### C. Closed-Loop Point Cloud Registration Performance

We validate the registration performance of our method for closing the loop. We follow Cattaneo et al. [9] and test our method on sequences 00 and 08 of the KITTI odometry dataset and sequences 02 and 09 of the KITTI-360 dataset. For a fair comparison, we follow Cattaneo et al. to train the *LCR-Net†*

on the KITTI sequence 05-07 and 09, validating it on KITTI sequence 02. The point cloud pairs with ground truth pose distances less than $4\,\mathrm{m}$ and time intervals greater than $50\,\mathrm{s}$ are chosen as loop closure samples.

**Metrics** : In line with [9], we employ three metrics to evaluate the registration performance at loop closure: i) Relative Translation Error (RTE), which measures the Euclidean distance between estimated and ground truth translation vectors, ii) Relative Yaw Error (RYE), which is the average difference between estimated and ground truth yaw angle, and iii) Registration Recall (RR), representing the fraction of scan pairs with RYE and RTE below certain thresholds, e.g., $5°$ and $2\,\mathrm{m}$.

**Results** : The baseline methods in this experiment include advanced traditional registration methods such as ICP [17] and RANSAC [30] with FPFH features [33]. Besides traditional methods, SOTA deep learning approaches are also included, such as RPMNet [65], FCGF [62], DGR [64], Predator [63], CofiNet [41], Geotransformer [42], RDMNet [43] and LCDNet [9]. Furthermore, we also include results from advanced place recognition methods that can output yaw angles, including Scan Context [18], LiDAR-Iris [19], and OverlapNet [1]. For the hand-crafted method, we use the results reported in [9]. For all DNN-based approaches, we use its official implementation along with open-source models trained on KITTI odometry datasets. We also report the results of Geotransformer and RDMNet using LGR. LCDNet also provides a fast version using weighted SVD. We denote it as LCDNet (fast) and report the results.

Tab. II shows the results on sequences 00 and 08 of the KITTI odometry dataset. As can be seen, our LCR-Net achieves the best performance across both sequences. LCR-Net is superior in pose estimation, benefiting from dense, reliable matching. We highlight that its pose estimation accuracy exceeds baselines by a large margin and is comparable to LCDNet refined by ICP. As illustrated in Tab. II, the registration recall (RR) for several baselines has reached saturation, reaching $100\%$ in the dataset of the closed loop distance below $4\,\mathrm{m}$. To better showcase the superiority of our approach, we have constructed a new test set comprising point cloud pairs with an overlap ratio exceeding $0.3$. The overlap ratio is computed following [1]. These test sets are more challenging, including point cloud pairs with a distance of up to $15\,\mathrm{m}$. We present the results of several advanced baselines in former tests. As can be seen, our LCR-Net further amplifies its advantages over other competing approaches. Only LCR-Net maintains $100\%$ RR and a similar pose estimation accuracy, while others have all declined.

We highlight that the pose estimation accuracy of LCR-Net surpasses even LCDNet refined by ICP. Especially regarding translation estimation, LCR-Net reduces the error by $64\%$ on seq. 00 and by $16\%$ on seq. 08. This result is encouraging, as current global registration methods typically exhibit inferior registration accuracy compared to geometry-based local registration approaches based on a fine initial guess. As a result, they are commonly employed as initial estimations for methods like the ICP algorithm in practical applications. In this experiment, LCR-Net demonstrates a noteworthy ad-

TABLE II: Point Cloud Registration Results on Closed Loop of KITTI Odometry Dataset.

| | | *Closed loop with distance below 4 m* | | | | | |
| | | | Seq. 00 | | | Seq. 08 | |
| | | RR(%) | RTE(m)[succ./all] | RYE(°)[succ./all] | RR | RTE(m)[succ./all] | RYE(°)[succ./all] |
|---|---|---|---|---|---|---|---|
| RANSAC-based | RANSAC [30] | 33.95 | 0.98/2.75 | 1.37/12.01 | 15.61 | 1.33/4.57 | 1.79/37.31 |
| | FCGF [62] | 47.31 | 1.05/2.07 | 0.60/1.88 | 27.80 | 1.28/2.42 | 2.38/159.48 |
| | Predator [63] | <u>98.93</u> | 0.07/0.13 | 0.11/0.45 | 99.70 | 0.10/0.11 | 0.34/0.53 |
| | CofiNet [41] | **100** | <u>0.07/0.07</u> | <u>0.10/0.10</u> | **100** | <u>0.10/0.10</u> | **0.33/0.33** |
| | Geotransformer [42] | 98.11 | 0.09/0.28 | 0.11/0.13 | 99.12 | 0.11/0.18 | 0.35/1.76 |
| | RDMNet [43] | 98.36 | 0.07/0.27 | 0.11/0.28 | <u>99.80</u> | 0.10/0.14 | 0.34/0.49 |
| | LCDNet [9] | **100** | 0.11/0.11 | 0.12/0.12 | **100** | 0.15/0.15 | <u>0.34/0.34</u> |
| | LCR-Net[†] | **100** | **0.04/0.04** | **0.09/0.09** | **100** | **0.08/0.08** | **0.33/0.33** |
| | LCR-Net | 100 | 0.04/0.04 | 0.09/0.09 | 100 | 0.08/0.08 | 0.33/0.33 |
| RANSAC-free | Scan Context* [18] | 97.66 | -/- | 1.34/1.92 | 98.21 | -/- | 1.71/3.11 |
| | LiDAR-Iris* [19] | <u>98.83</u> | -/- | 0.65/1.69 | 99.29 | -/- | 0.93/1.84 |
| | OverlapNet* [1] | 83.86 | -/- | 1.28/3.89 | 0.10 | -/- | 2.03/65.45 |
| | ICP [17] | 35.57 | 0.97/2.08 | 1.36/8.98 | 0 | -/2.43 | -/160.46 |
| | DGR [64] | 95.11 | 0.57/0.74 | 0.41/2.70 | 2.42 | 1.13/7.72 | 3.13/135.53 |
| | RPMNet [65] | 47.31 | 1.05/2.07 | 0.60/1.88 | 27.80 | 1.28/2.42 | 1.77/13.13 |
| | LCDNet (fast) [9] | 93.03 | 0.65/0.77 | 0.86/1.07 | 60.71 | 1.02/1.62 | 1.65/3.13 |
| | Geotransformer (LGR) [42] | 96.72 | 0.10/0.41 | 0.12/0.99 | 97.06 | 0.16/0.35 | 0.46/2.85 |
| | RDMNet (LGR) [43] | 97.57 | <u>0.06/0.37</u> | 0.12/0.64 | <u>99.36</u> | <u>0.10/0.22</u> | 0.34/0.59 |
| | LCR-Net[†] (LGR) | **100** | **0.04/0.04** | **0.09/0.09** | **100** | **0.08/0.08** | **0.33/0.33** |
| | LCR-Net (LGR) | 100 | 0.04/0.04 | 0.09/0.09 | 100 | 0.07/0.07 | 0.33/0.33 |
| | LCDNet+ICP [9] | 100 | 0.04/0.04 | 0.09/0.09 | 100 | 0.09/0.09 | 0.33/0.33 |
| | | *Closed loop with overlap beyond 0.3* | | | | | |
| | | | Seq. 00 | | | Seq. 08 | |
| | | RR(%) | RTE(m)[succ./all] | RYE(°)[succ./all] | RR | RTE(m)[succ./all] | RYE(°)[succ./all] |
| RANSAC-based | Predator [63] | 80.52 | 0.10/2.13 | 0.14/13.29 | 81.98 | 0.15/1.93 | 0.37/23.05 |
| | CofiNet [41] | <u>98.21</u> | <u>0.09/0.29</u> | <u>0.13/1.13</u> | <u>99.05</u> | <u>0.13/0.24</u> | <u>0.35/1.49</u> |
| | Geotransformer [42] | 82.07 | 0.10/3.04 | 0.13/4.53 | 92.69 | 0.15/0.97 | 0.41/10.29 |
| | RDMNet [43] | 71.36 | 0.13/4.53 | 0.23/7.81 | 92.94 | 0.16/1.37 | 0.46/6.09 |
| | LCDNet [9] | 94.86 | 0.23/0.89 | 0.24/3.93 | 96.87 | 0.31/0.62 | 0.51/1.74 |
| | LCR-Net[†] | **100** | **0.05/0.05** | **0.10/0.10** | **100** | **0.08/0.08** | **0.34/0.34** |
| | LCR-Net | 100 | 0.05/0.05 | 0.10/0.10 | 100 | 0.08/0.08 | 0.34/0.34 |
| RANSAC-free | LCDNet (fast) [9] | 66.82 | 0.56/2.02 | 0.68/6.12 | 47.29 | 0.88/3.00 | 1.06/5.28 |
| | Geotransformer (LGR) [42] | <u>80.43</u> | <u>0.11/3.33</u> | <u>0.16/7.91</u> | 87.85 | 0.20/1.47 | 0.56/12.78 |
| | RDMNet (LGR) [43] | 70.19 | 0.13/5.06 | 0.23/11.14 | <u>91.62</u> | <u>0.15/1.98</u> | <u>0.48/7.16</u> |
| | LCR-Net[†] (LGR) | **100** | **0.04/0.04** | **0.10/0.10** | **100** | **0.08/0.08** | **0.33/0.33** |
| | LCR-Net (LGR) | 100 | 0.04/0.04 | 0.10/0.10 | 100 | 0.08/0.08 | 0.34/0.34 |
| | LCDNet+ICP [9] | 94.98 | 0.11/0.77 | 0.14/3.84 | 96.98 | 0.18/0.48 | 0.43/1.68 |

Superscript * means the approaches only estimate yaw angle. The best results are highlighted in bold, and the second best in underlines.

TABLE III: Point Cloud Registration Results on Closed Loop of KITTI360.

| | | | Seq. 02 | | | Seq. 09 | |
| | | RR(%) | RTE(m)[succ./all] | RYE(°)[succ./all] | RR | RTE(m)[succ./all] | RYE(°)[succ./all] |
|---|---|---|---|---|---|---|---|
| RANSAC-based | RANSAC [30] | 24.78 | 1.24/3.67 | 1.83/32.22 | 29.69 | 1.12/3.14 | 1.48/23.42 |
| | FCGF [62] | 12.43 | 0.56/11.32 | 0.62/146.77 | 62.40 | 0.47/5.06 | 0.44/60.64 |
| | Predator [63] | 97.56 | <u>0.22/0.34</u> | <u>0.28/1.19</u> | <u>98.44</u> | <u>0.14/0.22</u> | 0.18/0.66 |
| | CofiNet [41] | 98.56 | 0.25/0.30 | 0.30/0.35 | **100** | 0.15/0.15 | <u>0.19/0.19</u> |
| | Geotransformer [42] | 92.32 | 0.26/0.84 | 0.49/9.71 | 96.24 | 0.16/0.50 | 0.25/1.84 |
| | RDMNet [43] | 94.98 | 0.23/0.66 | 0.38/0.97 | 83.27 | 0.19/2.24 | 0.31/3.04 |
| | LCDNet [9] | <u>98.62</u> | 0.28/0.32 | 0.32/0.35 | **100** | 0.18/0.18 | 0.20/0.20 |
| | LCR-Net[†] | **98.63** | **0.21/0.26** | **0.27/0.30** | **100** | **0.11/0.11** | **0.16/0.16** |
| | LCR-Net | 98.64 | 0.21/0.26 | 0.27/0.30 | 100 | 0.11/0.11 | 0.16/0.16 |
| RANSAC-free | Scan Context* [18] | 92.31 | -/- | 1.60/5.49 | <u>95.29</u> | -/- | 1.40/6.80 |
| | LiDAR-Iris* [19] | <u>96.54</u> | -/- | 1.71/3.44 | 86.26 | -/- | 1.51/7.08 |
| | OverlapNet* [1] | 11.42 | -/- | 1.79/76.74 | 54.33 | -/- | 1.38/33.62 |
| | ICP [17] | 4.19 | 1.10/2.26 | 1.74/149.76 | 21.24 | 1.06/2.22 | 1.34/66.34 |
| | DGR [64] | 14.26 | 0.66/6.68 | 0.72/126.81 | 63.51 | 0.50/3.14 | 0.38/53.83 |
| | RPMNet [65] | 37.99 | 1.18/2.26 | 1.30/5.97 | 41.42 | 1.13/2.21 | 1.02/3.95 |
| | LCDNet (fast) [9] | 82.92 | 0.84/1.10 | 1.28/1.67 | 89.49 | 0.76/0.94 | 0.99/1.19 |
| | Geotransformer (LGR) [42] | 91.24 | 0.25/1.01 | 0.61/7.90 | 94.23 | <u>0.17/0.70</u> | <u>0.29/2.70</u> |
| | RDMNet (LGR) [43] | 94.44 | <u>0.19/0.72</u> | <u>0.36/1.43</u> | 80.36 | 0.18/2.82 | 0.31/5.64 |
| | LCR-Net[†] (LGR) | **98.63** | **0.16/0.21** | **0.24/0.27** | **100** | **0.10/0.10** | **0.14/0.14** |
| | LCR-Net (LGR) | 98.59 | 0.16/0.21 | 0.23/0.26 | 100 | 0.09/0.09 | 0.14/0.14 |
| | LCDNet+ICP [9] | 98.51 | 0.20/0.25 | 0.24/0.27 | 100 | 0.10/0.10 | 0.15/0.15 |

Superscript * means the approaches only estimate yaw angle. The best results are highlighted in bold, and the second best in underlines.

vancement in registration accuracy compared to the SOTA global registration methods refined by ICP. This significant outcome profoundly underscores the practical advantage of our proposed method in real-world applications.

Tab. III presents the results on KITTI-360 datasets. As can be seen, the advantages of our method over other baselines are still maintained on KITTI-360 datasets. Two baselines, i.e., CofiNet and LCDNet, demonstrate comparable registration recall to LCR-Net. However, both methods fall short compared to LCR-Net regarding registration accuracy. Additionally, both LCDNet and CofiNet require RANSAC for pose estimation, and LCDNet requires additional ground point filtering.

From the tests, we have observed that RANSAC-based methods generally outperform RANSAC-free methods. However, RANSAC is computationally intensive and can account for more than half of the entire registration process, as evident from our runtime experiments as in Sec. VI-E. Nevertheless, our proposed method attains the best performance without relying on RANSAC. By employing a fast solver, LGR, LCR-Net attains comparable registration performance to that of RANSAC and even surpasses it in terms of translation estimation accuracy. This characteristic offers a significant advantage for integrating our approach within real-time systems.

In sum, our LCR-Net achieves SOTA performance in closed-loop point cloud registration. It exhibits significant advantages over the baseline methods, offering: i) The highest level of registration robustness; ii) Exceptional registration accuracy compared to baseline methods, surpassing even the results obtained by baseline with ICP refinement; iii) The ability to achieve best performance without relying on RANSAC.

### D. Continuous Relocalization Performance

To evaluate the relocalization performance, we use LiDAR pairs at most 10 m apart as samples, which may cause odometry degeneration. We test our approach on sequences 08-10 of KITTI odometry, KITTI-360, Apollo-SouthBay, Ford Campus, and Mulran datasets. For a fair comparison, we follow the setup of prior work [42], [63], [66] and train the model *LCR-Net$^\diamond$* on KITTI odometry sequence 00-05 and validate it on KITTI odometry sequence 06-07. Notably, the sensors, environments, and platform setups differ between the KITTI odometry dataset and other datasets, thoroughly testing the approaches' generalization abilities.

**Metrics** : We use three metrics to evaluate the registration performance: i) Relative Translation Error (RTE), which measures the Euclidean distance between estimated and ground truth translation vectors, ii) Relative Rotation Error (RRE), which measures the geodesic distance between estimated and ground truth rotation matrices, and iii) Registration Recall (RR), which represents the fraction of scan pairs with RRE and RTE below certain thresholds, e.g., $5°$ and 2 m.

**Results** : We compare the results of our method with the recent RANSAC-based SOTA methods: Predator [63], CofiNet [41], NgeNet [66], Geotransformer [42], and RDMNet [43]. We also compare our method using Local-to-Global Registration (LGR) [42] with SOTA RANSAC-free methods: HReg-Net [44], Geotransformer [42] and RDMNet [43]. The results are shown in Tab. IV.

TABLE IV: Continuous Registration Results on Multiple Datasets. All the Models Are Only Trained on the KITTI Dataset.

| | | KITTI | KITTI-360 | Apollo | Ford | Mulran |
|---|---|---|---|---|---|---|
| *Registration Recall (%)* | | | | | | |
| RANSAC-based | Predator [63] | **99.82** | 99.50 | 99.27 | 99.32 | 53.02 |
| | CofiNet [41] | **99.82** | 99.62 | **100** | 99.32 | 80.79 |
| | NgeNet [66] | **99.82** | **99.94** | **100** | **100** | 82.96 |
| | Geotransformer [42] | **99.82** | 99.86 | **100** | **100** | 75.68 |
| | RDMNet [43] | **99.82** | 99.89 | **100** | **100** | 87.09 |
| | LCR-Net$^\diamond$ | **99.82** | **99.94** | **100** | **100** | **97.89** |
| | LCR-Net | 99.82 | 99.94 | 100 | 100 | 98.22 |
| RANSAC-free | HRegNet [44] | 96.76 | 20.39 | 9.39 | 90.48 | - |
| | Geotransformer (LGR) [42] | **99.82** | 99.86 | **100** | 98.64 | 72.91 |
| | RDMNet (LGR) [43] | **99.82** | 99.90 | **100** | **100** | 83.68 |
| | LCR-Net$^\diamond$ (LGR) | **99.82** | **99.94** | **100** | **100** | **96.76** |
| | LCR-Net (LGR) | 99.82 | 99.94 | 100 | 100 | 97.53 |
| *Relative Rotation Error (°)* | | | | | | |
| RANSAC-based | Predator [63] | 0.25 | 0.29 | 0.21 | 0.38 | 1.03 |
| | CofiNet [41] | 0.37 | 0.44 | <u>0.18</u> | 0.37 | 0.52 |
| | NgeNet [66] | 0.26 | 0.30 | <u>0.18</u> | 0.26 | 0.35 |
| | Geotransformer [42] | 0.22 | 0.28 | 0.28 | 0.28 | <u>0.30</u> |
| | RDMNet [43] | <u>0.18</u> | <u>0.25</u> | **0.10** | <u>0.21</u> | 0.45 |
| | LCR-Net$^\diamond$ | **0.17** | **0.22** | **0.10** | **0.18** | **0.17** |
| | LCR-Net | 0.19 | 0.24 | 0.09 | 0.16 | 0.17 |
| RANSAC-free | HRegNet [44] | 1.04 | 2.18 | 2.14 | 0.91 | - |
| | Geotransformer (LGR) [42] | 0.31 | 0.36 | 0.29 | 0.50 | 0.49 |
| | RDMNet (LGR) [43] | **0.27** | **0.35** | **0.29** | **0.39** | <u>0.48</u> |
| | LCR-Net$^\diamond$ (LGR) | <u>0.35</u> | <u>0.36</u> | <u>0.30</u> | <u>0.42</u> | **0.45** |
| | LCR-Net (LGR) | 0.35 | 0.36 | 0.29 | 0.34 | 0.43 |
| *Relative Translation Error (cm)* | | | | | | |
| RANSAC-based | Predator [63] | 5.8 | 7.2 | 7.8 | 11.7 | 30.04 |
| | CofiNet [41] | 8.2 | 10.1 | 6.7 | 11.6 | 17.3 |
| | NgeNet [66] | 6.1 | 7.5 | 5.9 | 7.9 | <u>9.2</u> |
| | Geotransformer [42] | 6.7 | 8.1 | 6.1 | 10.6 | 12.0 |
| | RDMNet [43] | <u>5.3</u> | <u>7.0</u> | <u>4.6</u> | <u>7.6</u> | 14.4 |
| | LCR-Net$^\diamond$ | **3.9** | **5.7** | **3.6** | **7.1** | **7.5** |
| | LCR-Net | 3.9 | 5.8 | 3.4 | 6.6 | 7.4 |
| RANSAC-free | HRegNet [44] | 6.7 | 112.8 | 116.7 | 16.5 | - |
| | Geotransformer (LGR) [42] | 5.5 | 6.9 | 5.1 | <u>12.2</u> | 9.7 |
| | RDMNet (LGR) [43] | <u>3.9</u> | <u>5.9</u> | <u>3.5</u> | 6.1 | <u>9.3</u> |
| | LCR-Net$^\diamond$ (LGR) | **2.7** | **5.0** | **2.5** | **6.1** | **6.2** |
| | LCR-Net (LGR) | 2.8 | 5.4 | 2.5 | 6.5 | 6.2 |

The best results are highlighted in bold, and the second best in underlines.

As can be seen, when using RANSAC, our LCR-Net outperforms all the baselines on all the datasets. Especially on the Mulran dataset, RDMNet shows remarkable performance by boosting baselines with a large margin for all three metrics, i.e., increasing registration recall by 10%, reducing rotation error by 43%, and reducing translation error by 20%. The Mulran dataset poses significant challenges for generalization, as it loses approximately $70°$ FOV. The outstanding performance of our proposed method on the Mulran dataset demonstrates its exceptional generalization capability. When using LGR, our method attains remarkable results for translation estimation and achieves an average reduction of 35% in translation error compared to the best results attained through RANSAC-based methods across all datasets.

We conduct a qualitative comparison of registration using LCR-Net and LCDNet shown in Fig. 6. We also present ICP alignment by using the LCDNet prediction as an initial guess to demonstrate the accuracy of LCR-Net. The first two columns depict the successful cases of both methods. Our LCR-Net achieves better alignment than LCDNet and LCDNet
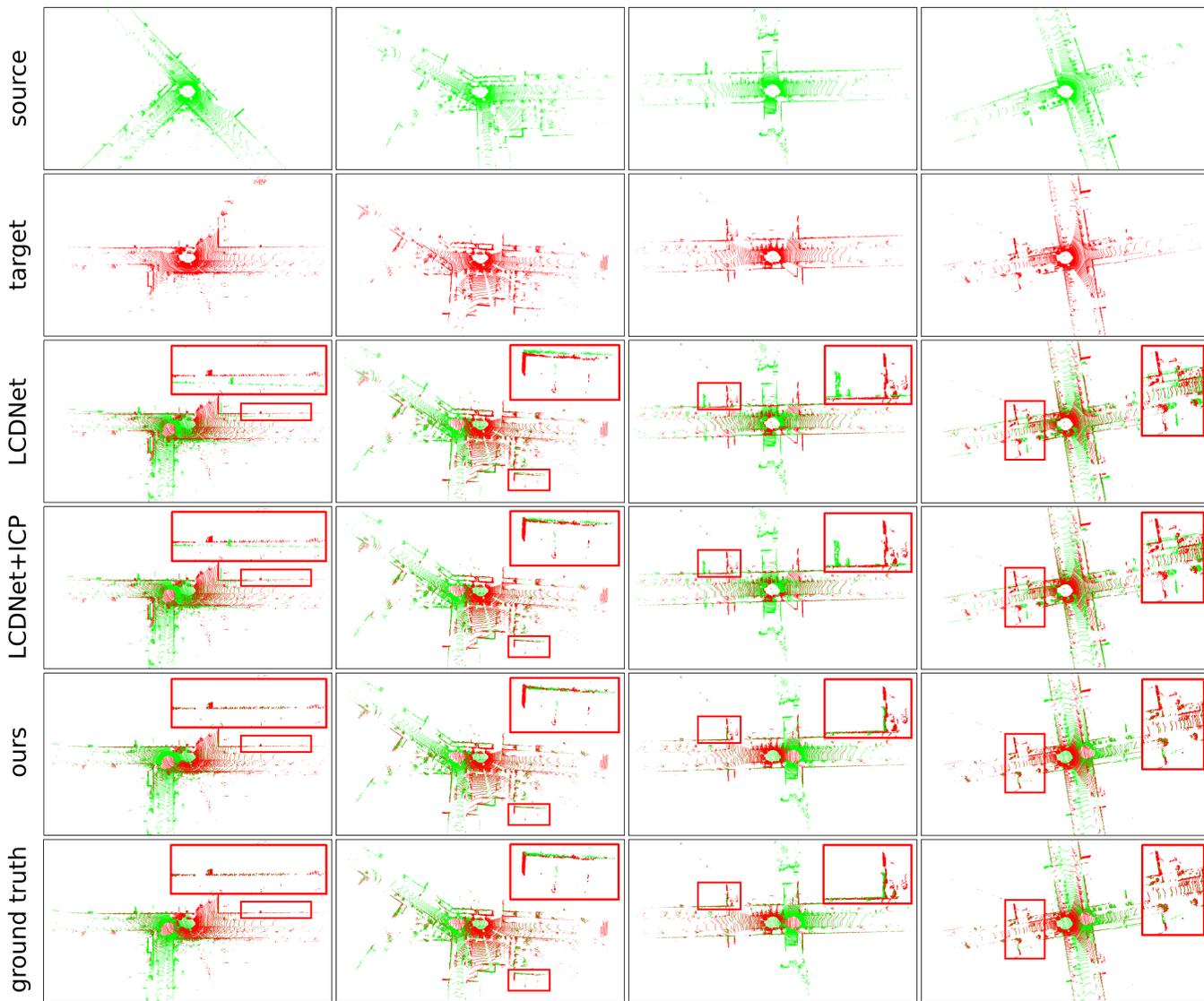
Fig. 6: Qualitative comparison of registration using LCDNet and ICP with using LCDNet as initial guess and LCR-Net. The first two columns depict the successful cases of both methods, where LCR-Net achieves better alignment. The last two columns show the failure cases of LCDNet, yet LCR-Net still perfectly aligns the two point clouds.
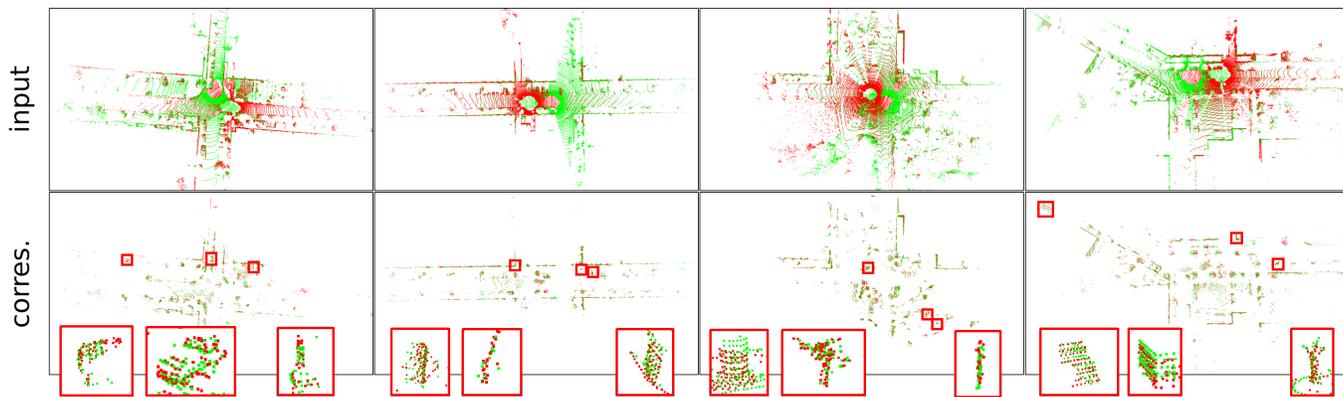


Fig. 7: We visualize the correspondent points founded by LCR-Net. For better visualization, we align the input point clouds and the founded correspondent points using LCR-Net based on LGR. For each cases, we zoom in on three representative regions for closer examination. Our LCR-Net effectively finds correspondences covering overall overlapping region and focuses on geometrically significant region.

TABLE V: Comparison of Inference Time.
(a) Inference time for point cloud registration.

| | Model | Descriptor Extraction (ms) | Pairwise Reg. (ms) | Total (ms) |
|---|---|---|---|---|
| Hand | RANSAC | 135 | 156 | 291 |
| | FGR | 135 | 246 | 381 |
| DNN-based | Predator | 173 | 244 | 417 |
| | CofiNet | 377 | 257 | 634 |
| | LCDNet | 235 | 431 | 666 |
| | LCDNet+ICP | 235 | 566 | 801 |
| | LCR-Net | 293 | 390 | 683 |
| | LCR-Net (LGR) | 293 | 13 | 306 |

(b) Inference time for candidate retrieval.

| | Model | Descriptor Extraction (ms) | Pairwise Comp. (ms) | Map querying (ms) |
|---|---|---|---|---|
| Hand | Scan Context | 1 | 0.09 | 6 |
| | LiDAR-Iris | 10 | 9 | 42037 |
| DNN | OverlapNet | 15 | 6 | 26634 |
| | LCDNet | 182 | 0.01 | 5 |
| | LCR-Net | 50 | 0.01 | 5 |

refined by ICP. The last two columns show the failure cases of LCDNet, while LCR-Net keeps achieving good alignment.

To provide more insights into the proposed LCR-Net, we visualize the correspondent points founded by LCR-Net as in Fig. 7. The points are aligned using LCR-Net based on LGR. We observe that there are four characteristics of the regions that the LCR-Net focuses on: i) It effectively captures overlapping regions of two point clouds. ii) It finds matches over all overlapping regions of two point clouds rather than being limited to a relatively small range, allowing for a wider baseline important for accurate registration. iii) It neglects the majority of ground points. iv) It focuses more on isolated landmarks like tree trunks, signage, vehicles, etc. v) It prioritizes important geometric structures like corners and sloping surfaces.

### E. Study on Runtime

We report the runtime of LCR-Net compared to existing SOTA baselines on a system with an Intel i9-10920X CPU and an NVIDIA GTX 3090 GPU. All the methods are evaluated on KITTI sequence 00 with the official implementation. We present the runtime of point cloud registration in Tab. Va in which the descriptor extraction time also includes the preprocessing required by the respective method. Pairwise Reg. refers to pair registration and uses RANSAC as default. As depicted, our LCR-Net using LGR is the fastest method among DNN-based approaches. More commendably, it is also the most robust and accurate one, as shown in Sec. VI-D and Sec. VI-C.

In Tab. Vb, the runtime of loop detection is evaluated. We emphasize that LCR-Net exhibits significantly higher efficiency in descriptor extraction when inferring loop detection compared to descriptor extraction during point cloud registration. This efficiency is achieved due to our framework, which only activates the KPEncoder and the global description head during inference, eliminating the need for complex point-level descriptor generation. In contrast, methods that rely on point-wise descriptors, like LCDNet, have no such substantial runtime reductions for global descriptor extraction. Regarding

map querying, we report the runtime for searching a descriptor in all its previous scans. For the pairwise comparison, both LiDAR-Iris and OverlapNet employ an ad-hoc comparison function that is hard to optimize in runtime, while LCR-Net and LCDNet use the Euclidean distance. Therefore, we can use the FAISS library [54] for fast searching in LCR-Net and LCDNet. We report the querying time of Scan Context using the ring key descriptor for fast searching. As depicted, our LCR-Net demonstrates exceptional efficiency in loop retrieval with the support of Faiss.

In summary, LCR-Net exhibits high efficiency in both registration and candidate retrieval tasks, rendering it well-suited for online SLAM applications.

### F. Evaluation of Complete SLAM System

In the previous experiments, we have demonstrated the superior performance and high efficiency of the proposed LCR-Net in specific subtasks of loop closing and relocalization. We have also verified that the model trained on our dataset splittings, *LCR-Net*, can achieve comparable SOTA performance with the models specifically trained for each task, *LCR-Net*[†] and *LCR-Net*[◇]. In this section, we will evaluate the performance of our network in solving loop closing and relocalization challenges when integrated with a SLAM system. We use the SLAM system with A-LOAM as the front-end odometry. We demonstrate the gradual improvement of SLAM results by integrating our *LCR-Net* using the LGR estimator for relocalization and loop closing, as depicted in Fig. 8. The relocalization operates at a frequency of 2.5 Hz, while loop closing is configured to operate at a frequency of 1 Hz.

The top row of Fig. 8a shows the original results of A-LOAM. As seen, A-LOAM suffers severe degeneration and fails to provide satisfactory results. We determine a keyframe's ground truth degeneration flag by comparing the relative pose estimation error to a pre-defined threshold. When the pose error exceeds 1 m, the current keyframe is considered degenerated. The bottom row of Fig. 8a highlights the degenerated keyframes on the ground truth trajectory in red.

The bottom row of Fig. 8b shows the degenerated keyframes detected by our approach. Our approach effectively detects all the degenerations. While this introduces false positives, it avoids missing degeneration detections, which is paramount in practical applications. The SLAM results are significantly improved by enabling the relocalization module, as in the top row of Fig. 8b. The trajectory is colored with the distance error between our estimations and ground truth poses.

The bottom row of Fig. 8c shows the loop detected and closed by our approach. As in the top row of Fig. 8c, by enabling both relocalization and loop closing module, the results are further refined compared to Fig. 8b.

As the code of LCDNet with SLAM is currently unavailable, we compare the results with the most commonly employed loop closing approach, which exploits Scan Context for loop detection and ICP for registration, as in Fig. 8d. The bottom row of Fig. 8d shows the loop detected and closed by Scan Context. As depicted, it introduces several false positives, whereas our approach successfully avoids such occurrences.
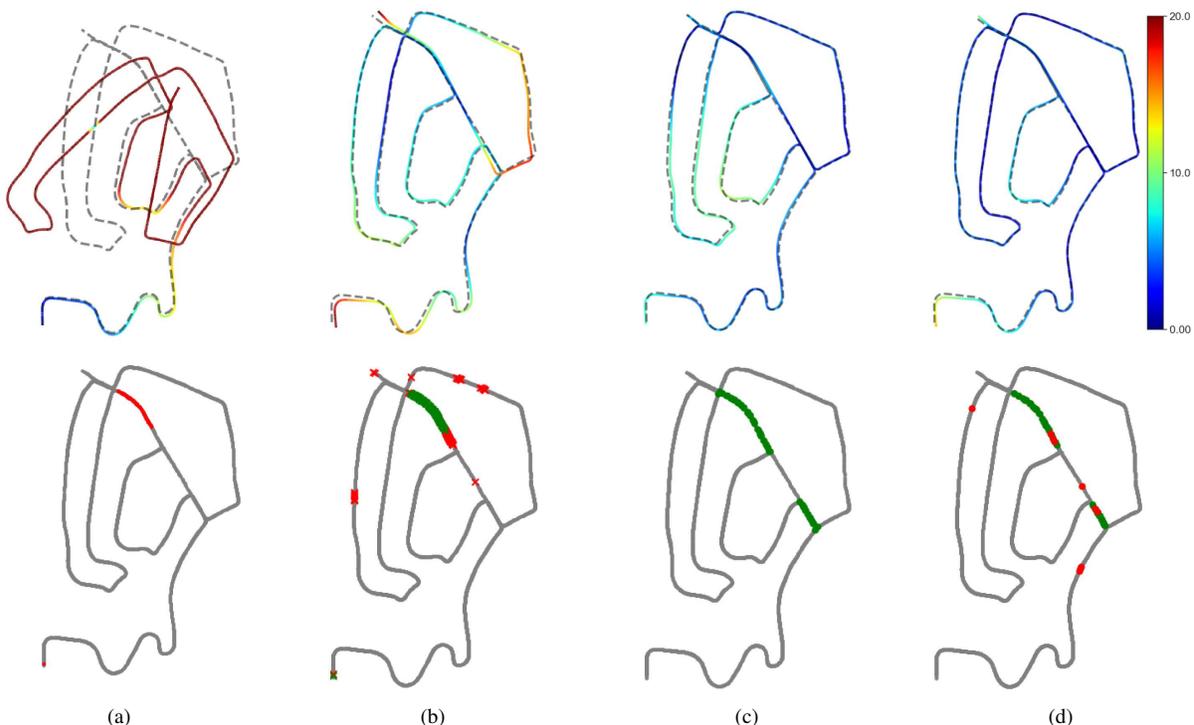
Fig. 8: Performance of A-LOAM with integrating our approach compared to integrating Scan Context on KITTI sequence 02. We present the results of (a) original A-LOAM (top) and degenerated keyframes colored in red (bottom), (b) A-LOAM with enabling relocalization using our approach (top) and the detected degenerated keyframes (bottom, green crosses $\times$ are true positive and red crosses $\times$ are false positive detections), (c) A-LOAM with enabling relocalization and loop closing using our approach (top) and the detected loop closures (bottom, green dots ● are true positive detections), (d) A-LOAM with enabling relocalization using our approach and loop closing using Scan Context and ICP (top) and the detected loop closures (bottom, green dots ● are true positive detections and red dots ● are false positive).
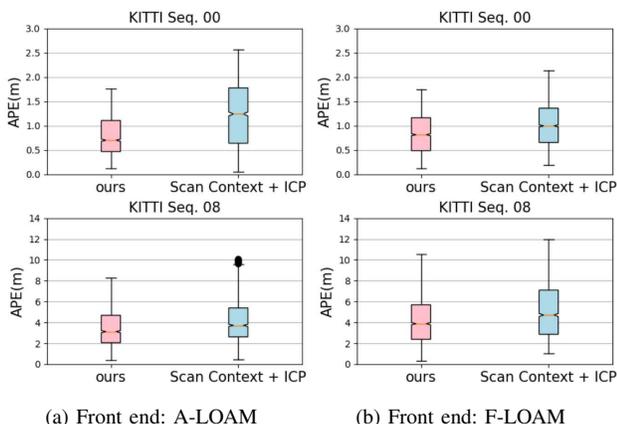


Fig. 9: Comparison of SLAM systems with different front-end odometry approaches on different sequences.

We assess different systems' Absolute Position Error (APE) on KITTI odometry sequences 00 and 08. The results of our system and the system employing Scan Context and ICP as the loop closing pipeline are presented in Fig. 9. By leveraging the flexible structure of the PGO framework, we also substitute the front-end odometry with F-LOAM [67], thereby demonstrating the capability of our method to integrate diverse odometry systems. As depicted in Fig. 9, regardless of the chosen odometry front-end, our approach consistently yields more precise localization results than systems utilizing Scan Context and ICP. Moreover, our approach exhibits smaller minimum and maximum errors and enhanced stability.

### G. Ablation Study

We conduct ablation studies to better understand the effectiveness of each module in our proposed LCR-Net, as in Tab. VI. Tab. VIa and Tab. VIb study specific designs of 3D-RoFormer++ and VoteEncoder in continuous registration tasks using Mulran dataset and closed-loop registration task using KITTI dataset (loop with distance below 4 m). Tab. VIc and Tab. VId study overall structure and training strategy in registration task and candidate retrieval task using KITTI and Ford Campus datasets. We use the model trained on our dataset splittings, *LCR-Net*, as the default model. The registration performance is evaluated using the RANSAC estimator.

**Rotary embedding.** Tab. VIa studies the influence of different rotary position embedding used in 3D-RoFormer++, where i) nonlinear refers to the embedding using Eq. (3), ii) linear w/o penalty refers to the embedding using Eq. (9) but trained without penalty loss. Using a linear rotary position embedding improves registration robustness and accuracy. The reason could be attributed to the network's ability to consistently represent relative positions through a linear embedding, thus enhancing the registration performance. Using the penalty loss can further help the position embedding learning.

**VoteEncoder.** Tab. VIb studies the designs of VoteEncoder. Leveraging the central prediction algorithm, VoteEncoder acquires a more robust estimation for the center of the interested region, thus enhancing registration accuracy. The incorporation of feature aggregation capability via KPConv improves both robustness and accuracy.

TABLE VI: Ablation Studies of Individual Modules.

(a) **Rotary embedding in 3D-RoFormer++**. Linear embedding is effective.

| case | Mulran RR(%) | RRE(°) | RTE(cm) | KITTI-loop RR(%) | RRE(°) | RTE(cm) |
|---|---|---|---|---|---|---|
| nonlinear | 94.97 | 0.18 | 8.1 | 99.86 | 0.29 | 5.9 |
| linear | **98.22** | **0.17** | **7.4** | **100** | **0.21** | **4.3** |
| linear w/o penalty | 94.29 | 0.18 | 7.5 | **100** | 0.23 | **4.3** |

(b) **VoteEncoder**. The full VoteEncoder performs the best.

| case | Mulran RR(%) | RRE(°) | RTE(cm) | KITTI-loop RR(%) | RRE(°) | RTE(cm) |
|---|---|---|---|---|---|---|
| full | **98.22** | **0.17** | **7.4** | **100** | **0.21** | **4.3** |
| KPConv w/o | 93.40 | 0.18 | 8.0 | 99.97 | 0.26 | 5.0 |
| central prediction w/o | 92.80 | 0.19 | 8.0 | 99.99 | 0.24 | 4.5 |

(c) **Overall structure**. Both VoteEncoder and 3D-RoFormer++ contributes to enhance the registration robustness and accuracy.

| case | Registration Mulran RR(%) | RRE(°) | RTE(cm) | KITTI-loop RR(%) | RRE(°) | RTE (cm) | Candidate retrieval KITTI AUC | F1max | Recall@1 | Recall@1% | Campus AUC | F1max | Recall@1 | Recall@1% |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| full | **98.22** | **0.17** | **7.4** | **100** | **0.21** | **4.3** | 0.958 | 0.922 | 0.937 | **0.993** | 0.972 | **0.920** | 0.932 | 0.987 |
| VoteEncoder w/o | 85.46 | 0.41 | 20.1 | 99.99 | 0.24 | 6.6 | **0.975** | **0.940** | **0.948** | 0.990 | **0.973** | 0.917 | 0.927 | 0.981 |
| 3D-RoFormer++ w/o | 93.67 | 0.22 | 10.1 | **100** | 0.26 | 4.8 | 0.950 | 0.919 | 0.920 | 0.962 | 0.951 | 0.871 | **0.939** | **0.994** |

(d) **Training strategy**. Linear probe ensures the best registration performance while allowing for enhanced candidate retrieval performance.

| case | b p n | Registration Mulran RR(%) | RRE(°) | RTE(cm) | KITTI-loop RR(%) | RRE(°) | RTE (cm) | Candidate retrieval KITTI AUC | F1max | Recall@1 | Recall@1% | Campus AUC | F1max | Recall@1 | Recall@1% |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| linear probe | 6 6 6 | **98.22** | **0.17** | **7.4** | **100** | **0.21** | **4.3** | **0.958** | **0.922** | **0.937** | **0.993** | **0.972** | **0.920** | **0.932** | **0.987** |
| fine-tune | 1 1 1 | 90.13 | 0.18 | 8.1 | **100** | **0.21** | 4.5 | 0.900 | 0.834 | 0.911 | 0.975 | 0.883 | 0.818 | 0.877 | 0.971 |
| end-to-end | 1 1 1 | 85.70 | 0.21 | 9.6 | **100** | 0.22 | 4.7 | 0.868 | 0.808 | 0.872 | 0.966 | 0.834 | 0.782 | 0.871 | 0.979 |

Default structure are marked in gray and the best results are highlighted in bold.

**Overall structure.** Tab. VIc studies the impact of our VoteEncoder and 3D-RoFormer++. As can be seen, VoteEncoder significantly enhances the registration performance. The underlying factors contributing to this noteworthy improvement are examined in Sec. VI-H. We study the 3D-RoFormer++ by replacing it with a vanilla Transformer using an absolute position encoder [40]. As depicted, 3D-RoFormer++ exhibits superior advantages in registration compared to the vanilla Transformer.

Though neither 3D-RoFormer++ nor VoteEncoder directly contribute to global descriptor generation, they still impact performance. As illustrated, 3D-RoFormer++ is beneficial in enhancing the capacity of the global descriptor, which suggests that 3D-RoFormer++, compared to vanilla Transformer, is more effective in improving the feature learning of the backbone. On the other hand, using the VoteEncoder does not have a significant impact. The performance trade-off is justified by the substantial improvement in registration brought about by VoteEncoder.

**Training strategy.** In Tab. VId, we study the influence of training strategy. We compare the training strategy used in this article, linear probe, with two alternative strategies. The first strategy is fine-tuning, where we also pre-train the model in the registration task and subsequently fine-tune the model along with the global descriptor head in the candidate retrieval task. The second strategy is end-to-end, which starts training the whole network from scratch. We are constrained to smaller batch sizes and fewer positive/negative scans than the linear probe strategy in these two strategies. With the principle of maximizing computational resource utilization, we have set the batch size (b) to 1, with $N_p = 1$ positive scan (p) and $N_n = 1$ negative scan (n) for both strategies. As can be seen, the batch size, positive scans, and negative scans used in these two strategies are significantly less than that used in linear probes, which largely constrains the performance in candidate retrieval. Note that the linear probe also uses the same batch

size as 1 when trained on registration, yet it still achieves the best registration performance. This can be attributed to the linear probe strategy not facing the trade-off between global descriptor capacity and registration capability during training.

### H. Insights on How VoteEncoder Contributes to Registration

In our ablation study, we observed significant performance enhancement brought about by our VoteEncoder in registration tasks. To provide a comprehensive understanding of the benefits of VoteEncoder for registration, we further conducted insight experiments on six additional metrics: i) Inlier Ratio (IR), the ratio of correct correspondences with residuals below a certain threshold, e.g., 0.6 m, after applying the ground truth transformation. ii) Match Recall (MR), the ratio of true point matches found among all the true matches. iii) Hit Ratio (HR), the fraction of true points that have matches among all points that have matches. iv) Patch Inlier Ratio (PIR), v) Patch Match Recall (PMR), and vi) Patch Hit Ratio (PHR) are with the same meaning as IR, MR, and HR, respectively, but at sparse matching levels. The true patch match represents that patch pairs exhibit actual overlap. MR and PMR evaluate the coverage ratio of found matches among all true matches, while HR and PHR better assess the match distribution over the overlapping region since a point/patch may possess multiple matches.

We present the result of LCR-Net and LCR-Net without VoteEncoder in Tab. VII. To evaluate the efficacy of the VoteEncoder in contrast to an alternative downsampling and aggregation approach, we also present a variant, denoted as *5KPEncoder*, that replaces VoteEncoder with another KPEncoder, comprising 5 layers of KPEncoder. As can be seen, using VoteEncoder has shown a decrease in IR and PIR but has notably increased PMR and HR. The rationale behind this can be explained as follows: the VoteEncoder can be seen as a downsampler that filters redundant points and aggregates them

TABLE VII: A Comprehensive Evaluation of VoteEncoder on Continuous Registration Task Using the KITTI Dataset.

| Model | Direct metrics | | | Indirect metrics | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | RR(%) | RRE(°) | RTE(cm) | IR | MR | HR | PIR | PMR | PHR |
| full | **99.82** | **0.17** | **3.9** | 78.2 | **9.2** | **69.2** | 80.9 | **62.8** | **86.3** |
| VoteEncoder w/o | **99.82** | 0.22 | 6.4 | **84.5** | 3.7 | 35.1 | **97.0** | 15.7 | 43.2 |
| 5KPEncoder | **99.82** | 0.20 | 4.6 | 81.2 | 6.3 | 51.0 | 88.6 | 39.4 | 77.5 |

The best results are highlighted in bold.



(a) Input point clouds.



(b) Matching points found by LCR-Net ablating VoteEncoder.



(c) Matching points found by LCR-Net.

Fig. 10: We visualize the matching points identified by LCR-Net to demonstrate the focus of our LCR-Net. The matching points obtained by LCR-Net after ablating VoteEncoder is also visualized for comparison. For better visualization, we align the points in two scans using the pose estimated by LCR-Net based on LGR. In each case, we zoom in on three representative regions for closer examination.

into new keypoints. After filtering massive points with similar neighbors, it becomes more difficult to find true matches in sparser candidates (as observed by the decrease in PIR). However, this also allows LCR-Net to focus on searching different geometrically salient regions over the point cloud rather than repeatedly sampling the same salient region (as observed by the increase in PMR and HR). Moreover, this also leads to a more reasonable point patch grouping, making it easier to match an object as a whole and simplifying the subsequent dense matching process (as observed by a decrease in PIR by 17.9%, while IR only decreases by 6.3%). As depicted, the significant improvement in PMR and HR brought by VoteEncoder, allowing for wider baselines and more dimensional constraints in registration, greatly reduces registration errors (as observed by a decrease in RRE by 23% and RTE by 39%). Comparison between *5KPEncoder* and full LCR-Net demonstrates that the VoteEncoder achieves more effective downsampling and better registration than a KPEncoder.

We provide qualitative comparisons of matching points as in Fig. 10. As depicted, matching points obtained by full LCR-Net offers enhanced coverage of overlapping regions within the point clouds. We zoom in and examine four specific local areas: car, tree, signage, and shrub within both scenes. It is evident that the full LCR-Net effectively groups and matches points of local areas.

**Insights**. The above results lead us to rethink which indirect metrics are more important in point cloud registration. Apart from the final performance evaluated by direct metrics of Registration Recall (RR), Relative Rotation Error (RRE), and Relative Translation Error (RTE), most prior work [34], [35], [37]–[39], [41]–[43], [62], [63], [66] focuses only on the keypoint salience evaluated by indirect metric Inlier Ratio (IR), also referred to as precision. However, higher IR does not guarantee higher robustness or accuracy, as observed in Tab. VII. Conversely, the coverage ratio of found matches among all true matches, as evaluated by metrics MR and HR, should be a crucial factor to consider. These metrics provide insights into the distribution of matching points and can guide the design of new keypoint detection methods by considering this distribution. It is important to note that the coverage ratio is just one minor aspect within the broader research context of distribution, and there are still numerous other factors, such as divergence and distribution across dimensions, waiting to be explored in this field.

Indeed, numerous existing keypoint detection and matching methods [38], [39] have proposed to prevent keypoints from being excessively concentrated in local areas. Earlier work for LiDAR odometry, LOAM [12] and its variants [67], [68] propose to extract keypoints uniformly from all LiDAR scan lines. However, they have not provided a clear explanation for the underlying reasons. We are the first to specifically propose this insight and validate it through experimental cases by providing quantitative metrics.

## VII. CONCLUSION

This article analyzes the challenges and commonalities in loop closing and relocalization tasks and proposes a unified framework to address both tasks. Based on this framework, we present LCR-Net, a novel multi-head network for candidate retrieving and point cloud registration. LCR-Net incorporates a novel keypoint detection module, offering three types of features to meet different requirements for distinct post-processing. Two novel sub-modules, 3D-RoFormer++ and VoteEncoder, are proposed for feature generation. 3D-RoFormer++ learns contextual information between two point clouds and leverages a novel transformation-invariant rotary position embedding for geometric information aggregation. VoteEncoder, on the other hand, learns to detect keypoints from distinct objects in the scene while maintaining a uniform, non-repeating distribution. Both 3D-RoFormer++ and VoteEncoder significantly contribute to our network's loop closing and relocalization performance. Given the impressive improvement in registration performance brought by VoteEncoder, we conduct experiments and evaluation metrics to analyze the underlying reasons and obtain an important insight that

point matching should consider distribution rather than solely focusing on precision. We demonstrate the SOTA performance of our LCR-Net in three sub-tasks related to loop closing and relocalization. It is shown that the best performance of LCR-Net does not rely on a costly robust estimator and outperforms the baseline method refined by ICP. Finally, we integrate our network as the loop closing and relocalization module in a complete SLAM system and verify its capability to address loop closing and relocalization tasks.

## REFERENCES

[1] X. Chen, T. Läbe, A. Milioto, T. Röhling, O. Vysotska, A. Haag, J. Behley, and C. Stachniss, "Overlapnet: Loop closing for lidar-based SLAM," in *Proc. of Robotics: Science and Systems (RSS)*, 2020.

[2] L. Li, X. Kong, X. Zhao, T. Huang, and Y. Liu, "Ssc: Semantic scan context for large-scale place recognition," in *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, 2021, pp. 2092–2099.

[3] X. Chen, T. Läbe, A. Milioto, T. Röhling, J. Behley, and C. Stachniss, "OverlapNet: A Siamese Network for Computing LiDAR Scan Similarity with Applications to Loop Closing and Localization," *Autonomous Robots*, vol. 46, pp. 61–81, 2021.

[4] S. S. Harithas, G. Singh, A. Chavan, S. Sharma, S. Patni, C. Arora, and K. M. Krishna, "Findernet: A data augmentation free canonicalization aided loop detection and closure technique for point clouds in 6-dof separation," *arXiv preprint*, vol. abs/2304.01074, 2023. [Online]. Available: https://api.semanticscholar.org/CorpusID:257912690

[5] B. Steder, M. Ruhnke, S. Grzonka, and W. Burgard, "Place Recognition in 3D Scans Using a Combination of Bag of Words and Point Feature Based Relative Pose Estimation," in *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems*, 2011. [Online]. Available: proceedings:steder2011iros.pdf

[6] M. Bosse and R. Zlot, "Place recognition using keypoint voting in large 3d lidar datasets," in *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*, 2013, pp. 2677–2684.

[7] Y. Cui, X. Chen, Y. Zhang, J. Dong, Q. Wu, and F. Zhu, "Bow3d: Bag of words for real-time loop closing in 3d lidar slam," *IEEE Robotics and Automation Letters (RA-L)*, vol. 8, pp. 2828–2835, 2022.

[8] J. Sivic and A. Zisserman, "Video Google: A Text Retrieval Approach to Object Matching in Videos," in *Proc. of the IEEE Intl. Conf. on Computer Vision*, Oct 2003, pp. 1470–1477 vol.2. [Online]. Available: http://www.robots.ox.ac.uk/~vgg/publications/papers/sivic03.pdf

[9] D. Cattaneo, M. Vaghi, and A. Valada, "Lcdnet: Deep loop closure detection and point cloud registration for lidar slam," *IEEE Trans. on Robotics (TRO)*, vol. 38, pp. 2074–2093, 2021.

[10] J. Arce, N. Vodisch, D. Cattaneo, W. Burgard, and A. Valada, "Padloc: Lidar-based deep loop closure detection and registration using panoptic attention," *IEEE Robotics and Automation Letters (RA-L)*, vol. 8, pp. 1319–1326, 2022. [Online]. Available: https://api.semanticscholar.org/CorpusID:252383609

[11] K. Vidanapathirana, P. Moghadam, S. Sridharan, and C. Fookes, "Spectral geometric verification: Re-ranking point cloud retrieval for metric localization," *IEEE Robotics and Automation Letters (RA-L)*, vol. 8, pp. 2494–2501, 2022. [Online]. Available: https://api.semanticscholar.org/CorpusID:252780185

[12] J. Zhang and S. Singh, "LOAM: Lidar Odometry and Mapping in Real-time," in *Proc. of Robotics: Science and Systems*, 2014. [Online]. Available: proceedings:zhang2014rss.pdf

[13] J. Zhang, M. Kaess, and S. Singh, "On Degeneracy of Optimization-Based State Estimation Problems," in *Proc. of the IEEE Intl. Conf. on Robotics & Automation*, 2016.

[14] W. Xu and F. Zhang, "Fast-lio: A fast, robust lidar-inertial odometry package by tightly-coupled iterated kalman filter," *IEEE Robotics and Automation Letters (RA-L)*, vol. 6, pp. 3317–3324, 2020. [Online]. Available: https://api.semanticscholar.org/CorpusID:223953478

[15] J. Behley and C. Stachniss, "Efficient Surfel-Based SLAM using 3D Laser Range Data in Urban Environments," in *Proc. of Robotics: Science and Systems*, 2018. [Online]. Available: http://www.roboticsproceedings.org/rss14/p16.pdf

[16] X. Chen, A. Milioto, E. Palazzolo, P. Giguère, J. Behley, and C. Stachniss, "SuMa++: Efficient LiDAR-based Semantic SLAM," in *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots*

[17] P. J. Besl and N. D. McKay, "A Method for Registration of 3D Shapes," *IEEE Trans. on Pattern Analalysis and Machine Intelligence*, vol. 14, no. 2, pp. 239–256, 1992. [Online]. Available: http://www-evasion.inrialpes.fr/people/Franck.Hetroy/Teaching/ProjetsImage/2007/Bib/besl_mckay-pami1992.pdf

[18] G. Kim and A. Kim, "Scan context: Egocentric spatial descriptor for place recognition within 3d point cloud map," in *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, 2018, pp. 4802–4809.

[19] Y. Wang, Z. Sun, J. Yang, and H. Kong, "Lidar iris for loop-closure detection," in *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, 2019, pp. 5769–5775.

[20] A. Uy and G. Lee, "PointNetVLAD: Deep point cloud based retrieval for large-scale place recognition," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2018, pp. 4470–4479. [Online]. Available: proceedings:angelina2018cvpr.pdf

[21] Z. Zhou, C. Zhao, D. Adolfsson, S. Su, Y. Gao, T. Duckett, and L. Sun, "Ndt-transformer: Large-scale 3d point cloud localisation using the normal distribution transform representation," in *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*, 2021, pp. 5654–5660.

[22] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. of the Advances in Neural Information Processing Systems (NIPS)*, 2017, pp. 5998–6008.

[23] J. Ma, J. Zhang, J. Xu, R. Ai, W. Gu, and X. Chen, "Overlaptransformer: An efficient and yaw-angle-invariant transformer network for lidar-based place recognition," *IEEE Robotics and Automation Letters (RA-L)*, vol. 7, pp. 6958–6965, 2022.

[24] J. Ma, X. Chen, J. Xu, and G. Xiong, "Seqot: A spatial–temporal transformer network for place recognition using sequential lidar data," *IEEE Trans. on Industrial Electronics*, vol. 70, no. 8, pp. 8225–8234, 2023.

[25] J. Ma, G. Xiong, J. Xu, and X. Chen, "CVTNet: A Cross-View Transformer Network for LiDAR-Based Place Recognition in Autonomous Driving Environments," *IEEE Trans. on Industrial Informatics (TII)*, 2023.

[26] J. Knopp, J. Sivic, and T. Pajdla, "Avoiding confusing features in place recognition," in *Proc. of the Europ. Conf. on Computer Vision*. Springer, 2010, pp. 748–761. [Online]. Available: https://www.di.ens.fr/~josef/publications/knopp10.pdf

[27] Y. Zhong, "Intrinsic shape signatures: A shape descriptor for 3d object recognition," in *Proc. of the Int. Conf. on Computer Vision (ICCV) Workshops*, 2009, pp. 689–696.

[28] S. Shi, C. Guo, L. Jiang, Z. Wang, J. Shi, X. Wang, and H. Li, "Pv-rcnn: Point-voxel feature set abstraction for 3d object detection," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 10 526–10 535.

[29] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "NetVLAD: CNN Architecture for Weakly Supervised Place Recognition," *IEEE Trans. on Pattern Analalysis and Machine Intelligence (TPAMI)*, vol. 40, pp. 1437–1451, 2015.

[30] M. Fischler and R. Bolles, "Random Sample Consensus: A Paradigm for Model Fitting with Applications to Image Analysis and Automated Cartography," *Commun. ACM*, vol. 24, no. 6, pp. 381–395, 1981.

[31] Z. Zhang, "Iterative point matching for registration of free-form curves and surfaces," *Intl. Journal of Computer Vision*, vol. 13, no. 2, pp. 119–152, 1994.

[32] A. Segal, D. Haehnel, and S. Thrun, "Generalized-ICP," in *Proc. of Robotics: Science and Systems*, 2009.

[33] R. B. Rusu, N. Blodow, and M. Beetz, "Fast point feature histograms (fpfh) for 3d registration," in *Proc. of the IEEE Intl. Conf. on Robotics & Automation*, 2009, pp. 3212–3217. [Online]. Available: proceedings:rusu2009icra.pdf

[34] H. Deng, T. Birdal, and S. Ilic, "Ppfnet: Global context aware local features for robust 3d point matching," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 195–205.

[35] ——, "Ppf-foldnet: Unsupervised learning of rotation invariant 3d local descriptors," in *Proc. of the Europ. Conf. on Computer Vision (ECCV)*, vol. 11209. Springer, 2018, pp. 620–638.

[36] C. Qi, H. Su, K. Mo, and L. Guibas, "PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2017. [Online]. Available: https://arxiv.org/pdf/1612.00593.pdf

[37] Z. J. Yew and G. H. Lee, "3dfeat-net: Weakly supervised local 3d features for point cloud registration," in *Proc. of the Europ. Conf. on*

*Computer Vision (ECCV)*, ser. Lecture Notes in Computer Science, vol. 11219.  Springer, 2018, pp. 630–646.

[38] J. Li and G. H. Lee, "Usip: Unsupervised stable interest point detection from 3d point clouds," in *Proc. of the IEEE Intl. Conf. on Computer Vision (ICCV)*, 2019, pp. 361–370.

[39] 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. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 6358–6366.

[40] C. Shi, X. Chen, K. Huang, J. Xiao, H. Lu, and C. Stachniss, "Keypoint matching for point cloud registration using multiplex dynamic graph attention networks," *IEEE Robotics and Automation Letters (RA-L)*, vol. 6, pp. 8221–8228, 2021.

[41] H. Yu, F. Li, M. Saleh, B. Busam, and S. Ilic, "Cofinet: Reliable coarse-to-fine correspondences for robust pointcloud registration," in *Proc. of the Advances in Neural Information Processing Systems (NIPS)*, 2021, pp. 23 872–23 884.

[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. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 11 133–11 142.

[43] C. Shi, X. Chen, H. Lu, W. Deng, J. Xiao, and B. Dai, "Rdmnet: Reliable dense matching based point cloud registration for autonomous driving," *IEEE Trans. on Intelligent Transportation Systems (ITS)*, 2023.

[44] 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. of the IEEE Intl. Conf. on Computer Vision (ICCV)*. IEEE, 2021, pp. 15 994–16 003.

[45] H. Zhou, Z. Yao, and M. Lu, "Lidar/uwb fusion based slam with anti-degeneration capability," *IEEE Trans. on Vehicular Technology*, vol. 70, pp. 820–830, 2021. [Online]. Available: https://api.semanticscholar.org/CorpusID:231919869

[46] C. Campos, R. Elvira, J. J. G. Rodr'iguez, J. M. M. Montiel, and J. D. Tardós, "Orb-slam3: An accurate open-source library for visual, visual–inertial, and multimap slam," *IEEE Trans. on Robotics (TRO)*, vol. 37, pp. 1874–1890, 2020. [Online]. Available: https://api.semanticscholar.org/CorpusID:220713377

[47] M. Zins, G. Simon, and M.-O. Berger, "Oa-slam: Leveraging objects for camera relocalization in visual slam," *Proc. of the Intl. Symposium on Mixed and Augmented Reality (ISMAR)*, pp. 720–728, 2022. [Online]. Available: https://api.semanticscholar.org/CorpusID:252367374

[48] H. Thomas, C. Qi, J. Deschaud, B. Marcotegui, F. Goulette, and L. Guibas, "KPConv: Flexible and Deformable Convolution for Point Clouds," in *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision*, 2019. [Online]. Available: proceedings:thomas2019iccv.pdf

[49] C. R. Qi, O. Litany, K. He, and L. J. Guibas, "Deep Hough Voting for 3D Object Detection in Point Clouds," in *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision*, 2019. [Online]. Available: https://openaccess.thecvf.com/content_ICCV_2019/papers/Qi_Deep_Hough_Voting_for_3D_Object_Detection_in_Point_Clouds_ICCV_2019_paper.pdf

[50] Y. Zhang, Q. Hu, G. Xu, Y. Ma, J.-H. Wan, and Y. Guo, "Not all points are equal: Learning highly efficient point-based detectors for 3d lidar point clouds," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 18 931–18 940. [Online]. Available: https://api.semanticscholar.org/CorpusID:247594094

[51] R. Sinkhorn, "A relationship between arbitrary positive matrices and doubly stochastic matrices," *Annals of Mathematical Statistics*, vol. 35, pp. 876–879, 1964. [Online]. Available: https://api.semanticscholar.org/CorpusID:120846714

[52] D. Kingma and J.Ba, "Adam: A method for stochastic optimization," *arXiv preprint*, vol. abs/1412.6980, 2014. [Online]. Available: http://arxiv.org/pdf/1412.6980

[53] M. Kaess, H. Johannsson, R. Roberts, V. Ila, J. J. Leonard, and F. Dellaert, "isam2: Incremental smoothing and mapping using the bayes tree," *Intl. Journal of Robotics Research (IJRR)*, vol. 31, pp. 216 – 235, 2012.

[54] J. Johnson, M. Douze, and H. Jégou, "Billion-scale similarity search with gpus," *IEEE Trans. on Big Data*, vol. 7, pp. 535–547, 2017. [Online]. Available: https://api.semanticscholar.org/CorpusID:926364

[55] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2012. [Online]. Available: http://www.cvlibs.net/publications/Geiger2012CVPR.pdf

[56] Y. Liao, J. Xie, and A. Geiger, "Kitti-360: A novel dataset and benchmarks for urban scene understanding in 2d and 3d," *IEEE Trans. on Pattern Analalysis and Machine Intelligence (TPAMI)*, vol. 45, pp. 3292–3310, 2021.

[57] W. Lu, Y. Zhou, G. Wan, S. Hou, and S. Song, "L3-Net: Towards Learning Based LiDAR Localization for Autonomous Driving," in *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2019. [Online]. Available: proceedings:lu2019cvpr-ltlb.pdf

[58] D. Barnes, M. Gadd, P. Murcutt, P. Newman, and I. Posner, "The oxford radar robotcar dataset: A radar extension to the oxford robotcar dataset," *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*, pp. 6433–6438, 2019. [Online]. Available: https://api.semanticscholar.org/CorpusID:202122594

[59] X. Zhang, L. Wang, and Y. Su, "Visual place recognition: A survey from deep learning perspective," *Pattern Recognition*, vol. 113, p. 107760, 2021.

[60] T. Röhling, J. Mack, and D. Schulz, "A Fast Histogram-Based Similarity Measure for Detecting Loop Closures in 3-D LIDAR Data," in *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems*, 2015, pp. 736–741. [Online]. Available: proceedings:rhling2015iros.pdf

[61] J. Komorowski, "Minkloc3d: Point cloud based large-scale place recognition," in *Proc. of the IEEE Winter Conf. on Applications of Computer Vision (WACV)*, 2020, pp. 1789–1798.

[62] C. B. Choy, J. Park, and V. Koltun, "Fully convolutional geometric features," in *Proc. of the IEEE Intl. Conf. on Computer Vision (ICCV)*, 2019, pp. 8957–8965.

[63] S. Huang, Z. Gojcic, M. Usvyatsov, A. Wieser, and K. Schindler, "Predator: Registration of 3d point clouds with low overlap," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. Computer Vision Foundation / IEEE, 2021, pp. 4267–4276.

[64] C. B. Choy, W. Dong, and V. Koltun, "Deep global registration," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 2511–2520.

[65] Z. J. Yew and G. H. Lee, "Rpm-net: Robust point matching using learned features," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 11 821–11 830.

[66] L. Zhu, H. Guan, C. Lin, and R. Han, "Neighborhood-aware geometric encoding network for point cloud registration," *arXiv preprint*, vol. abs/2201.12094, 2022.

[67] H. Wang, C. Wang, C.-L. Chen, and L. Xie, "F-loam : Fast lidar odometry and mapping," *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, pp. 4390–4396, 2021. [Online]. Available: https://api.semanticscholar.org/CorpusID:235727638

[68] T. Shan and B. Englot, "Lego-loam: Lightweight and ground-optimized lidar odometry and mapping on variable terrain," in *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, 2018, pp. 4758–4765. [Online]. Available: https://api.semanticscholar.org/CorpusID:53706363