

# TReR: A Lightweight Transformer Re-Ranking Approach for 3D LiDAR Place Recognition

Tiago Barros, Luís Garrote, Martin Aleksandrov, Cristiano Premebida, Urbano J. Nunes

**Abstract**—Autonomous driving systems often require reliable loop closure detection to guarantee reduced localization drift. Recently, 3D LiDAR-based localization methods have used retrieval-based place recognition to find revisited places efficiently. However, when deployed in challenging real-world scenarios, the place recognition models become more complex, which comes at the cost of high computational demand. This work tackles this problem from an information-retrieval perspective, adopting a first-retrieve-then-re-ranking paradigm, where an initial loop candidate ranking, generated from a 3D place recognition model, is re-ordered by a proposed lightweight transformer-based re-ranking approach (TReR). The proposed approach relies on global descriptors only, being agnostic to the place recognition model. The experimental evaluation, conducted on the KITTI Odometry dataset, where we compared TReR with s.o.t.a. re-ranking approaches such as  $\alpha QE$  and SGV, indicate the robustness and efficiency when compared to  $\alpha QE$  while offering a good trade-off between robustness and efficiency when compared to SGV.

## I. INTRODUCTION

Place recognition has become a popular approach to perform global localization and loop closure detection [1]. With the development of efficient deep learning (DL) techniques able to model geometrical information, 3D LiDAR information has gained increasing relevance in place recognition due to being more robust compared to visual data in appearance-changing environments: a feature of key importance in the long-term operation of autonomous vehicles [2]. In 3D LiDAR-based place recognition, DL networks are used to map the geometrical information (*i.e.* point clouds) into an adequate descriptor space, where descriptors from nearby physical places are close. In contrast, descriptors from farther physical places are apart. Nevertheless, as the DL-based models are exposed to more realistic scenarios, the more complex they become, which usually means more parameters to train and, consequently, more computational power is required for training. In other fields, such as information retrieval (IR), a common approach to this problem is the first-retrieve-then-re-ranking (RTrR) paradigm [3], [4].

In an RTrR architecture, two sequential stages take place: an initial coarse retrieval stage, where pre-trained models are used, and a re-ranking stage, where smaller re-ranking models are used, trained on a target domain dataset [5]. In the first stage, the goal is to learn global descriptors from the input point clouds. In the second stage, the goal is to learn

T. Barros, L. Garrote, C. Premebida, and U.J. Nunes are with the University of Coimbra, Institute of Systems and Robotics, Department of Electrical and Computer Engineering, Portugal. E-mails: {tiagobarros, garrote, cpremebida, urbano}@isr.uc.pt. Martin Aleksandrov is with Dahlem Center for Machine Learning and Robotics, Freie Universität Berlin, Berlin. E-mail: martin.aleksandrov@fu-berlin.de.

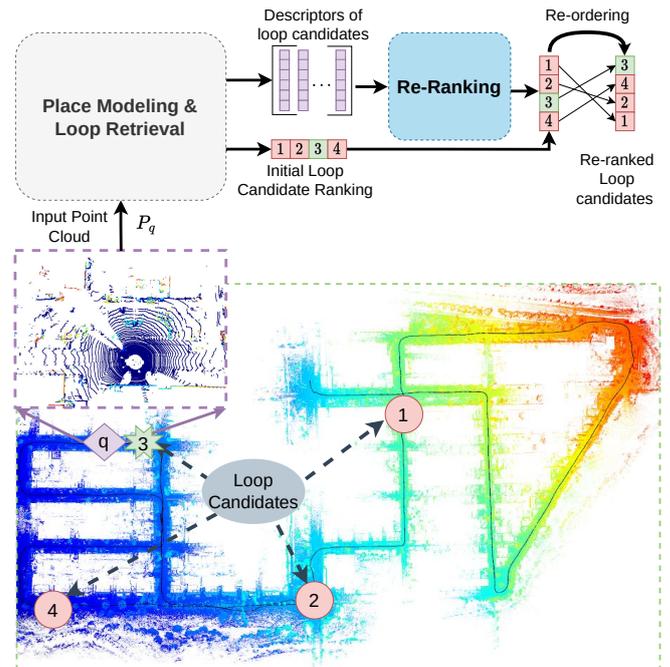


Fig. 1: A two stage framework with 3D LiDAR-based place recognition and re-ranking. For a given query  $q$  location, a pre-trained place recognition model first maps the point cloud  $P_q$  to a descriptor. The descriptor is used to query a database to return the  $k$  nearest neighbors, which represent an initial loop candidate ranking. The ranking is ordered from the most to the least similar. This process is illustrated with the following example: for the query point cloud  $P_q$  captured from the location  $q$ , four loop candidates are returned, which are identified in the map by a number and a shape. The number represents the initial rank, and the shape represents true loops (green star) or false loops (red circle).

interactions [6], associating a higher rank to more relevant candidates and a lower rank to less relevant candidates, such as false positives, as a way to refine the initial estimate. In traditional loop closure detection, false positives are identified using geometrical verification approaches such as RANSAC [7] and Iterative Closest Point [8], which are computationally demanding and, for this reason, often not being well-suited for large-scale applications.

In this work, we study the viability of the first-retrieve-then-re-ranking paradigm in place recognition, as shown in Figure 1. Inspired by the success of re-ranking-based applications in image retrieval [3], [4], our approach, called

TReR, is a lightweight, transformer-based model trained in a supervised regime. Unlike common transformer-based approaches for vision applications (e.g. [3], [4]), which rely on local features and return similarity scores for pairs of images, TReR relies on global descriptors only, making it lightweight and agnostic to place recognition models. Furthermore, TReR computes similarity scores for  $k$  loop candidates simultaneously instead of pairs.

We assess the adequacy of the approach within a 3D LiDAR place recognition framework for urban contexts. TReR is evaluated on five different place recognition models [9], [10], [11], which serve as baselines. Additionally, we also compare TReR to another global descriptor-based re-ranking approach ( $\alpha QE$  [12]) and a geometrical verification (GV) method (SGV[13]). The results are obtained on the KITTI dataset and indicate that TReR achieves a greater performance than the baselines (with no re-ranking) and  $\alpha QE$  while offering an acceptable trade-off between performance and efficiency in contrast to SGV[13]. In brief, the main contributions of this work are the following:

- An lightweight Transformer-based re-ranking approach (TReR) for 3D LiDAR place recognition that relies only on global descriptors and predicts the relative similarity scores of  $k$  candidates in one forward pass;
- TReR improves retrieval performances compared to other global-based methods and offers a good compromise between performance and efficiency.

## II. RELATED WORK

Place recognition has been an active field of research over the last decade, used in increasingly complex scenarios such as long-term operation in ever-changing environments [2]. The current endeavor to deal with such scenarios has been in the direction of proposing more complex models, requiring higher computational resources, which are not easily available. In other fields, such as IR, a common approach is to use an RTrT architecture, where, in the first stage, pre-trained models are used to generate coarse ranking estimates, and in the second stage, these initial estimates are refined using re-ranking models [5].

We next outline and discuss works on 3D LiDAR-based place recognition, which relate to the first stage, and works on re-ranking, which relate to the second stage.

### A. 3D-LiDAR-based Place Recognition

In 3D LiDAR-based place recognition, works such as PointNetVLAD [9] and LPD-Net [14] pioneered the idea of modeling physical places from point clouds directly by using end-to-end frameworks. Such frameworks, similar to those for vision-based place recognition, extract local features from point clouds and then aggregate the local features into a global descriptor. Over the years, the research interests have focused on different aspects of place recognition frameworks.

Some works focus on improving the feature extraction from point clouds e.g., addressing the rotation invariance problem in neural networks [15], proposing more suitable

feature extraction operators for point clouds, such as sparse 3D convolutions [10] and hierarchical graphs [16].

Other works have focused on aggregating local features into global descriptors. Most aggregators were initially proposed for images, such as the popular NetVLAD [17], which is a trainable generalization of VLAD [18]. Other approaches, for instance, resort the attention to re-adjusting the weights of local features [19], or compute various pooling operations on the local features, such as sum pooling (SPoC) [20] or generalized-mean pooling (GeM) [12]. In addition, the combination of the aforementioned aggregators into one global aggregator to leverage the strength of each aggregator has been explored in [11].

In this work, we use PointNetVLAD [9], LoGG3D-Net [10], and ORCHNet [11] for generating global descriptors and, additionally, we have adapted SPoC [20] and GeM [12] for working with point clouds.

### B. Re-Ranking

A re-ranking approach takes an initial ranking estimate and re-orders it based on a pre-defined criterion. In information retrieval, such a criterion is to rank higher content that is more relevant for a given query. In retrieval-based place recognition, the aim is to promote loop candidates that are true loops while ranking lower false positives.

These re-ranking methods can be categorized into two families: those that use only global descriptors and those that use local and global features combined. As for the global descriptors only, a widely used approach is Query Expansion (QE). In particular,  $\alpha QE$  has been used both in image retrieval [12] and retrieval-based place recognition [13]. On the other hand, the category of methods that use both local and global information is more common in applications with geometrical data. For instance, RANSAC [7] and Iterative Closest Point [8] are widely adopted in localization-related tasks, where they are employed to assess how well two distinct point clouds match. These methods, known as geometrical verification (GV), work well when seeking false positives but are usually computationally demanding when working with large point clouds. More recently, a more efficient approach was proposed in [13] called spectral geometrical verification (SGV), which relies on correspondence compatibility graphs [21] to assess the spatial consistency of two point clouds.

## III. PROPOSED APPROACH

In this section, we describe the proposed Transformer-based **Re-Ranking** framework (TReR) outlined in Figure 2. Firstly, we formulate the place recognition and re-ranking tasks in generic terms. Secondly, we present the proposed approach in more detail.

### A. Problem Formulation

The proposed framework is depicted in Figure 2. It can be split into two sub-tasks: 1) loop retrieval, using a place recognition approach for generating  $k$  candidates; 2) loop re-ranking, which re-orders the initial loop ranking, giving

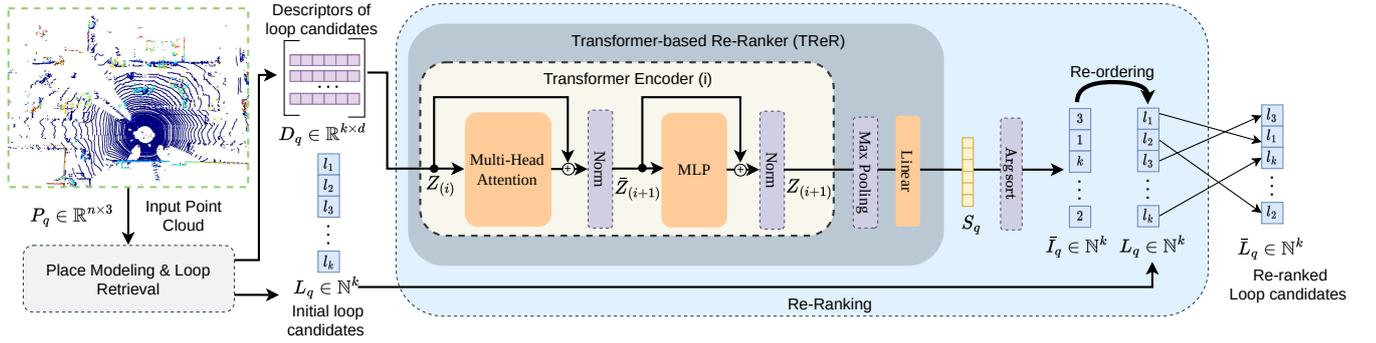


Fig. 2: The place recognition framework with two sequential tasks: (1) place modeling (Loop retrieval) and (2) re-ranking. The first model (Loop retrieval) receives a query point cloud  $P_q$  and then returns a loop ranking  $L_q$  and the corresponding descriptors  $D_q$ . The proposed re-ranking approach TRER receives the descriptors as input and returns a new ranking  $\bar{L}_q$  as output.

greater ranks to the true loops and lower ranks to the false loops.

1) *Loop Retrieval*: A point cloud-based loop retrieval task can be formulated as follows: given a query 3D point cloud  $P_q \in \mathbb{R}^{n \times 3}$  with  $n$  points, and a database with descriptors of already visited places, the task goal is to map  $P_q$  to a descriptor vector  $d_q \in \mathbb{R}^d$  of  $d$  dimensions, and use  $d_q$  to query the database to retrieve the  $k$  nearest neighbors based on a similarity score computed using the Euclidean distance. The retriever returns a set of database indices  $L_q = \{l_i\}_{i=1}^k \in \mathbb{N}^k$ , where  $l_i$  is the index (*i.e.* position in the database) of the  $i$ -th retrieved loop candidate. The candidates are ordered from the most similar to the least similar.

In this work, a loop candidate  $l_i \in L_q$  is a true loop, when the loop candidate is within a given  $\zeta$  Euclidean distance from the query: *i.e.* both the query and the respective candidate originated from the same place. In generic terms, this can be formulated by the following function:

$$\Gamma : \mathbb{R}^{n \times 3} \rightarrow \mathbb{N}^k, \quad (1)$$

where  $\Gamma$  can be learned. Since this work focuses on re-ranking, we refer the reader to [22] for more details on the training procedure for this learning task.

In the main framework, shown in Figure 2, the loop retrieval task is performed by the ‘‘Place Modeling & Loop Retrieval’’ module, which is detailed in Figure 3. This module returns the loop candidates in  $L_q$  and their descriptor vectors in  $D_q = \{d_i\}_{i=1}^k \in \mathbb{R}^{k \times d}$ , where  $k$  is the number of retrieved candidates, each of dimension size  $d$ .

2) *Re-Ranking*: The re-ranking task works on the premise that the initial loop ranking over candidates from  $L_q$ , contains samples that are not loops (*i.e.* False Positives), in which case we need an approach that is specialized in re-ordering the candidates, ranking higher loop candidates that correspond to real loops and lower loop candidates that correspond to false positives. Mathematically, the re-ranking process can be formulated as follows: given an initial loop ranking over candidates from  $L_q$  for a given query point cloud  $P_q$ , a function  $\Theta$  such that:

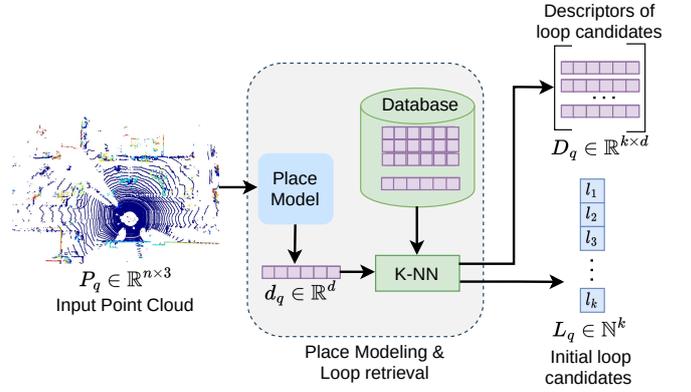


Fig. 3: The loop retrieval framework, where the input point cloud  $P_q$  is mapped to the descriptor vector  $d_q$ . The descriptor vector  $d_q$  is then used to query the database, after which the  $k$  nearest neighbors are returned, serving as loop candidates  $L_q$ . Both  $L_q$  and the respective descriptors  $D_q$  are returned.

$$\Theta : \mathbb{N}^k \rightarrow \mathbb{N}^k, \quad (2)$$

maps the initial ranking to a new ranking over the same  $k$  candidates, whose corresponding set is  $\bar{L}_q \in \mathbb{N}^k$ , where the true positives are promoted higher, and the false positives are promoted lower.

3) *Re-ranking training*: In this work, we resort to the global descriptor-based learning-to-rank approach depicted in Figure 2. In a learning-to-rank setting, the function  $\Theta$  is trained. To define the training set, let us define a set of queries  $Q \in \mathbb{N}^M$ , where  $M$  is the number of queries, for each query  $q \in Q$  exists a set of loop candidates  $L_q \in \mathbb{N}^k$ , their respective descriptors  $D_q \in \mathbb{R}^{k \times d}$  and the ground-truth relevance scores  $Y_q = \{y_i = \|p_q - p_{l_i}\|_2 \mid l_i \in L_q\} \in \mathbb{R}^k$ , where  $p_q \in \mathbb{R}^3$  and  $p_{l_i} \in \mathbb{R}^3$  are the 3D coordinates of the query  $q$  and the  $l_i$ -th loop candidate, respectively. The training set can thus be defined as follows:

$$\Upsilon = \{(D_q, Y_q) \mid q \in Q\}. \quad (3)$$

Given a training set  $\Upsilon$ , the goal of training is to find  $\hat{\Theta}$ , given by:

$$\hat{\Theta} = \arg \min(\mathcal{L}(\Theta)), \quad (4)$$

that minimizes the overall loss, given by:

$$\mathcal{L}(\Theta) = \frac{1}{|\Upsilon|} \sum_{(D_q, Y_q) \in \Upsilon} l(Y_q, \Theta(D_q)), \quad (5)$$

with  $l$  being a ranking loss function of a single query. This loss is computed between the ground-truth relevance scores  $Y_q$ , and the predicted relevance scores  $S_q = \Theta(D_q)$ .

### B. Transformer-based Re-Ranking

This section describes the proposed transformer-based re-ranking approach, TReR, depicted in Figure 2. Instead of following traditional point cloud-based loop re-ranking approaches, which rely predominantly on geometrical cues present in the scene [13], we adopt a more general information re-ranking approach, mainly used in image-related ranking tasks [3]. We are furthermore inspired by other works in context ranking, such as [23]. However, our approach differs in two ways: (a) the re-ranking model in [23] relies on local features, whereas the TReR model relies only on the global descriptors, which allows TReR to be used in applications where it is difficult to extract the local features; (b) the re-ranking model in [23] and the TReR model also have different output layers. More specifically, we add a max-polling layer followed by a linear transformation to compute the relative relevance scores used for re-ranking loop candidates.

1) *Architecture*: The proposed transformer-based re-ranker TReR is based on the multi-layer transformer encoders as proposed in [24]. We let  $C$  denote the number of transformer encoders. Pick  $i \in \{1, \dots, C\}$ . The  $i$ -th transformer encoder is outlined in Figure 2. We denote the input of the  $i$ -th transformer encoder as  $Z_{(i)}$ . For example, the input of the first encoder (*i.e.*  $Z_{(1)}$ ) is the set of global descriptor vectors  $D_q \in \mathbb{R}^{k \times d}$ . The multi-layer transformer encoder can thus be defined by the following equations:

$$\bar{Z}_{(i+1)} = l_N(Z_{(i)} + h_A(Z_{(i)})) \quad \text{and} \quad (6)$$

$$Z_{(i+1)} = l_N(f(\bar{Z}_{(i+1)})), \quad (7)$$

where  $l_N(\cdot)$  is a normalization layer as proposed in [25] and  $f(\cdot)$  is a multi-layer perceptron network with two linear transformations and a ReLU activation in between, given by:

$$f(Z) = \max(0, ZW_1 + b_1)W_2 + b_2, \quad (8)$$

where  $W_1 \in \mathbb{R}^{d \times d_h}$ ,  $W_2 \in \mathbb{R}^{d_h \times d}$ ,  $b_1 \in \mathbb{R}^{k \times d_h}$ , and  $b_2 \in \mathbb{R}^{k \times d}$  are weights and biases of linear projections, with  $d$  being the input/output dimensionality and  $d_h$  the hidden dimensionality. The function  $h_A(\cdot)$  is a multi-head attention module with  $n$  heads, defined as follows:

$$h_A(Z) = W_{AO} [h_1, \dots, h_j, \dots, h_n]^T, \quad (9)$$

where  $W_{AO} \in \mathbb{R}^{d \times dn}$  and  $h_j$  is the  $j$ -th attention head given by the following equation:

$$h_j = \sigma\left(\frac{QK^T}{\sqrt{d}}\right)V, \quad (10)$$

with  $\sigma(\cdot)$  being the *softmax* function and  $d$  is the feature dimensionality, while  $Q$ ,  $K$ , and  $V$  are, respectively, the queries, keys, and values. In this work, the attention mechanism is transformed to self-attention, which results in  $Q = Z_{(i)}W_j^Q$ ,  $K = Z_{(i)}W_j^K$ , and  $V = Z_{(i)}W_j^V$ , where  $W_j^Q$ ,  $W_j^K$ , and  $W_j^V \in \mathbb{R}^{d \times d}$  are linear projection matrices of the  $j$ -th head.

The  $i$ -th transformer encoder outputs re-weighted descriptors  $Z_{(i+1)} \in \mathbb{R}^{k \times d}$ , to which a max-pooling operation is applied alongside the feature dimensions followed by a linear projection:

$$S_q = \max(Z_{(i+1)})W_o + b_o, \quad (11)$$

with  $S_q \in \mathbb{R}^k$ , whereas  $W_o \in \mathbb{R}^{k \times k}$  and  $b_o \in \mathbb{R}^k$  encode the weights and biases, respectively. To obtain the final re-ranking, TReR sorts  $S_q$  in a descending direction as follows:

$$\bar{I}_q = \arg \text{sort}(S_q), \quad (12)$$

with  $\bar{I}_q \in \mathbb{N}^k$  being a vector of  $k$  dimensions, containing the reordering indices, which can be used for transforming the initial ranking of candidates from  $L_q$  into the final ranking of candidates from  $\bar{L}_q$ , as illustrated in Figure 2.

2) *Training*: The proposed re-ranking approach is trained in a supervised regime, where the objective is to maximize the correctness of the relative descriptor order instead of the absolute relevance. As in other ranking works such as [26] and [27], we adopt the logistic loss, which is a pairwise loss based on the cross-entropy and is given by:

$$l(Y_q, S_q) = \sum_{y_i > y_j} \log_2(1 + e^{-\sigma(s_i - s_j)}), \quad (13)$$

where  $\sigma$  is an adjustable parameter, and  $s_i, s_j \in S_q$  and  $y_i, y_j \in Y_q$ . The advantage of this loss is that it penalizes each out-of-order pair  $(i, j)$  that has  $y_i > y_j$  but  $s_i < s_j$ .

## IV. EXPERIMENTAL EVALUATION

In this section, we describe the experimental and the implementation parts. We also present and discuss the obtained empirical results.

TABLE I: The number of frames and loops in the dataset.

Sequence:	00	02	05	06	08
Length [frames]:	4051	4661	2761	1101	4071
Queries [loops]:	1097	450	692	284	521

TABLE II: The loop retrieval results obtained with the baseline models and the respective re-ranking models TReR and  $\alpha$ QE on the dataset. All models were evaluated using a cross-validation approach, tested on a fixed sequence, and trained on the remaining sequences. The presented results are the mean recall scores for the test sequence,  $Recall@N$  for  $N \in \{1, 5, 10\}$ , which are denoted as R1, R5, and R10, respectively, computed by using the  $N$  most relevant loop candidates.

	00			02			05			06			08		
	R1	R5	R10												
LoGG3D-Net	<b>0.69</b>	0.75	<b>0.79</b>	<b>0.55</b>	0.58	0.60	<b>0.68</b>	0.76	0.81	<b>0.91</b>	<b>0.94</b>	0.94	0.09	0.25	0.38
LoGG3D-Net + $\alpha$ QE [12]	0.58	0.71	0.78	0.48	0.58	0.60	0.59	0.72	0.79	0.75	0.91	0.95	0.08	0.25	0.35
LoGG3D-Net + TReR	<b>0.69</b>	<b>0.77</b>	<b>0.79</b>	<b>0.55</b>	<b>0.59</b>	<b>0.61</b>	<b>0.68</b>	<b>0.77</b>	<b>0.82</b>	<b>0.91</b>	<b>0.94</b>	<b>0.95</b>	<b>0.11</b>	<b>0.29</b>	<b>0.42</b>
ORCHNet	<b>0.72</b>	0.77	0.79	<b>0.53</b>	<b>0.57</b>	<b>0.58</b>	<b>0.65</b>	0.71	0.73	<b>0.89</b>	<b>0.94</b>	<b>0.95</b>	<b>0.52</b>	0.70	0.76
ORCHNet + $\alpha$ QE [12]	0.60	0.72	0.77	0.40	0.49	0.53	0.52	0.68	0.75	0.74	0.83	0.91	0.52	0.68	0.75
ORCHNet + TReR	<b>0.72</b>	<b>0.78</b>	<b>0.81</b>	<b>0.53</b>	<b>0.57</b>	<b>0.58</b>	<b>0.65</b>	<b>0.72</b>	<b>0.74</b>	<b>0.89</b>	<b>0.94</b>	<b>0.95</b>	<b>0.52</b>	<b>0.73</b>	<b>0.78</b>
PointNetVLAD	<b>0.75</b>	0.79	0.81	<b>0.51</b>	0.60	0.65	<b>0.68</b>	<b>0.74</b>	0.77	<b>0.92</b>	<b>0.95</b>	0.96	<b>0.61</b>	0.72	0.77
PointNetVLAD + $\alpha$ QE [12]	0.58	0.73	0.78	0.42	0.62	0.66	0.64	0.75	0.79	0.79	0.93	0.94	0.57	0.70	0.78
PointNetVLAD + TReR	<b>0.75</b>	<b>0.81</b>	<b>0.82</b>	<b>0.51</b>	<b>0.63</b>	<b>0.66</b>	<b>0.68</b>	<b>0.74</b>	<b>0.78</b>	<b>0.92</b>	<b>0.95</b>	<b>0.97</b>	<b>0.61</b>	<b>0.75</b>	<b>0.80</b>
PoinNetSPoC	<b>0.68</b>	<b>0.76</b>	0.78	<b>0.47</b>	0.56	0.59	<b>0.63</b>	0.70	0.73	<b>0.90</b>	<b>0.93</b>	0.94	<b>0.60</b>	<b>0.74</b>	0.77
PoinNetSPoC + $\alpha$ QE [12]	0.29	0.56	0.72	0.28	0.52	0.61	0.36	0.67	0.77	0.53	0.83	0.92	0.40	0.65	0.74
PoinNetSPoC + TReR	<b>0.68</b>	<b>0.76</b>	<b>0.79</b>	<b>0.47</b>	<b>0.58</b>	<b>0.62</b>	<b>0.63</b>	<b>0.71</b>	<b>0.76</b>	<b>0.90</b>	<b>0.93</b>	<b>0.95</b>	<b>0.60</b>	<b>0.74</b>	<b>0.82</b>
PointNetGeM	<b>0.68</b>	0.77	0.80	<b>0.49</b>	0.60	0.64	<b>0.63</b>	0.71	0.75	<b>0.91</b>	<b>0.95</b>	<b>0.97</b>	<b>0.63</b>	0.75	0.82
PointNetGeM + $\alpha$ QE [12]	0.23	0.51	0.68	0.29	0.54	0.63	0.32	0.59	0.73	0.59	0.79	0.90	0.55	0.72	0.79
PointNetGeM + TReR	<b>0.68</b>	<b>0.78</b>	<b>0.81</b>	<b>0.49</b>	<b>0.61</b>	<b>0.66</b>	<b>0.63</b>	<b>0.74</b>	<b>0.78</b>	<b>0.91</b>	<b>0.97</b>	<b>0.97</b>	<b>0.63</b>	<b>0.80</b>	<b>0.84</b>

### A. Dataset

The proposed approach is evaluated on the KITTI Odometry benchmark [28], a widely used dataset in various localization and place recognition works. This dataset comprises 22 sequences with 3D LiDAR, RGB, and RTK-GPS data from urban, highway, and rural environments, representing thus a variety of scenarios that an autonomous vehicle may encounter in real scenarios. Since this work is partly focused on loop retrieval, from within the 22 original sequences, we only used the sequences that contain revisited segments, which are 00, 02, 05, 06, and 08. Sequence 08 is particularly important because it is the only sequence that has revisits from the opposing direction. As a result, sequence 08 promotes an additional challenge for place recognition methods. In this work, we assume that a loop exists in a given sequence whenever two point clouds were captured from places within a range of 25 m as proposed in [22]. Table I contains the number of loops in each sequence under this assumption and the number of sequence frames.

### B. Experimental Setup

All experiments were conducted in Python 3.8 using the PyTorch framework with CUDA 11.6. All models were evaluated using a cross-validation approach: *e.g.* the descriptors of sequence 08 were obtained using a model trained on sequences 00, 02, 05, and 06. As for the training, all models were trained by using the AdamW optimizer with a learning rate (Lr) of 0.0001 and a weight decay (Wd) of 0.0005.

1) *Baselines*: TReR was evaluated on global descriptors with 256 dimensions, which were generated by five different models: PointNetVLAD [22], ORCHNet [11], LoGG3D-Net [10], PoinNetSPoC, and PoinNetGeM. PoinNetVLAD and ORCHNet were originally proposed with PointNet as backbone [29]. SPoC [20] and GeM [30] were adapted to

be implemented on a PointNet-based framework. For clouds with less than 10 thousand points, preliminary experiments conducted with PointNetVLAD, ORCHNet, PoinNetSPoC, and PoinNetGeM indicated that the performance of these models increased with the number of points in the input point clouds. For point clouds with more than 10 thousand points, the performance stagnated. Given these findings, we concluded that a point cloud with 10 thousand points was acceptable for preserving performance and computational demand, which led us to down-sample all point clouds to 10k points using a random sampling approach. The weights of PointNetVLAD [22], ORCHNet [11], PoinNetSPoC, and PoinNetGeM were trained based on the training protocol described in [9]. The weights of LoGG3D-Net [10] were pre-trained.

2) *Re-Ranking*: We performed experiments on a TReR architecture with one self-attention head (*i.e.*  $n = 1$ ) and one transformer encoder (*i.e.*  $C = 1$ ). As for the input, the descriptors had a size of 256 (*i.e.*  $d = 256$ ), the number of maximum candidates was 25 (*i.e.*  $k = 25$ ), and  $d_h = 512$ . With this architecture, TReR had in total 396,033 learnable parameters, which is 4% of LoGG3D-Net’s parameters (8,834,583) and 2% of PoinNetVLAD’s parameters (19,779,145).

3) *Evaluation*: The recall metric is a popular metric for measuring the performance of methods not only in place recognition but also in machine learning. The recall is defined as follows:

$$Recall = \frac{TP}{TP + FN}, \quad (14)$$

*i.e.* the ratio between true positive (TP) and the sum of the true positive and false negative (FN) loops in the ranking.

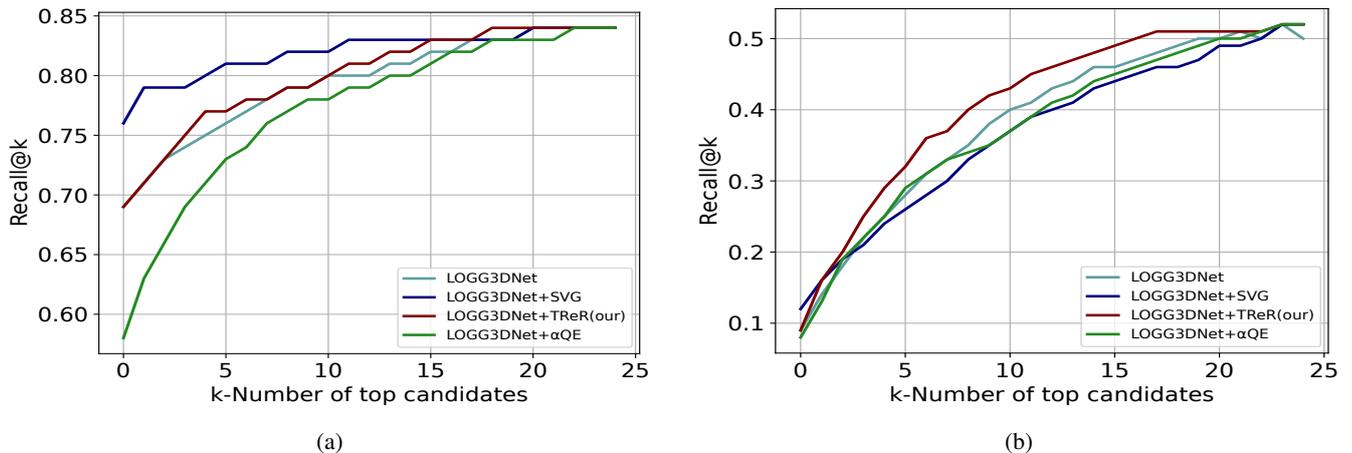


Fig. 4: Comparison between TReR, Spectral Geometric Verification (SGV), and  $\alpha$ -Query Expansion ( $\alpha$ QE). LoGG3D-Net was used as a baseline. Key: a) the results for  $Recall@k$  on sequence 00; b) the results for  $Recall@k$  on sequence 08.

For evaluating the methods, we use  $Recall@k$ , where  $k$  is the number of candidates retrieved from the database.

### C. Re-Ranking Evaluation

We next present and discuss the empirical results. Firstly, we give the results for the baseline methods without re-ranking and with re-ranking, using TReR (our approach) and  $\alpha$ QE [12]. Secondly, we compare TReR to the spectral geometrical verification (SGV) method [13] and  $\alpha$ QE, using as place recognition model LoGGNet3D-Net.

1) *Global Descriptors Only*: We compared the proposed approach TReR with the baselines and  $\alpha$ QE [12], which takes top-ranking candidates of a query and generates an updated query. The results are presented in Table II. Overall, they indicate that adding TReR on top of the baseline methods improves the overall performance in terms of  $Recall@k$  with  $k = [1, 5, 10]$ . TReR tended to have more effect in the middle of the ranking, lifting lower-ranked candidates while leaving unchanged upper-ranked candidates. The results were consistent throughout all models and sequences, indicating the robustness of TReR. On the one hand, TReR had a low impact on ranking estimates that required substantially fewer re-ranking moves, which was the case for sequence 06. On the other hand, TReR had a high impact on ranking estimates that required substantially more re-ranking moves, which was the case for sequence 08. This might be because, as we mentioned, sequence 08 is the most challenging, due to having revisits from the opposing direction. Finally,  $\alpha$ QE underperformed in all experiments. For this reason, we concluded that this approach might not be adequate for this application.

2) *Comparison to SGV*: Furthermore, we compared TReR to the SGV method [13], which relies on local features and key points to find local feature correspondence. GV-based approaches are usually known for their robustness, whereas approaches that use only global descriptors are usually known for their efficiency. We also observed such performance in our experiments. In Figure 4, we report the

results on sequences 00 and 08 for the baseline method LoGGNet3D-Net without re-ranking and LoGGNet3D-Net with re-ranking obtained with TReR, SGV, and  $\alpha$ QE. Sequences 00 and 08 represented two interesting edge cases for the methods. On sequence 00, SGV outperformed TReR and  $\alpha$ QE, confirming the robustness of GV-based approaches. By comparison, on sequence 08, TReR outperformed SGV and  $\alpha$ QE, confirming the efficiency of global-descriptors-only-based approaches. The fact that SGV performed worse than TReR on sequence 08 might be explained by the inability of SGV to find consistencies in pairs of point clouds that are significantly misaligned. In contrast, TReR does not use such local features and key points and, for this reason, performed better than SGV on sequence 08.

### D. Runtime Analysis

For each query, TReR re-ranked 25 candidates in a single forward pass, which took approximately 0.3 s of processing time on average on an AMD Ryzen 9 5900X 12-Core CPU with 64 GB RAM, and 0.001 s of processing time on average on a GeForce RTX 3090 24GB GPU. These and other runtime results are outlined in Table III. The approach that achieved the lowest performance was  $\alpha$ QE. Regarding SGV, we were not able to run it on our GPU due to memory constraints. In more detail, SGV required over 27GB of memory when computing key points and local features extracted from the clouds of 10 thousand points, whereas the GeForce RTX3090 had only 24GB for storing such points. Nevertheless, even when accounting just for the CPU time, we could observe that TReR is substantially more efficient than SGV.

## V. CONCLUSION

We proposed a 3D LiDAR-based place recognition re-ranking approach called TReR. TReR uses a lightweight transformer-based architecture with global descriptors to re-rank, in a first-retrieve-then-re-ranking paradigm, an initial ranking estimate from a 3D place recognition model.

TABLE III: The average processing time per query.

Methods	CPU(s)	GPU(s)
TReR	$0.30 \pm 0.12$	$0.001 \pm 0.0005$
SGV	$18 \pm 2.1$	-
$\alpha$ QE	$0.00016 \pm 0.000002$	$0.00016 \pm 0.000002$

The experimental evaluation was conducted on the KITTI Odometry dataset and highlighted the importance of re-ranking predictions for improving recall over several baseline methods that use no re-ranking and methods such as  $\alpha$ QE that use re-ranking. The improvement in recall was consistent over the baselines, showing the robustness of TReR. Similar results were observed when considering the method SGV. TReR improved place recognition model performance, specifically when considering challenging sequences such as sequence 08. Finally, TReR showed a good trade-off between performance and efficiency compared to the SGV method.

#### ACKNOWLEDGMENTS

This work has been partially supported by the project “GreenBotics” (ref. PTDC/EEI-ROB/2459/2021), funded by Fundação para a Ciência e a Tecnologia (FCT), Portugal. It was also supported by FCT through grant UIDB/00048/2020 and under the Ph.D. grant with reference 2021.06492.BD. Martin Damyanov Aleksandrov was supported by the DFG Individual Research Grant on “Fairness and Efficiency in Emerging Vehicle Routing Problems” (497791398).

#### REFERENCES

- [1] D. Cattaneo, M. Vaghi, and A. Valada, “LCDNet: Deep loop closure detection and point cloud registration for LiDAR SLAM,” *IEEE Transactions on Robotics*, vol. 38, no. 4, pp. 2074–2093, 2022.
- [2] T. Barros, R. Pereira, L. Garrote, C. Premebida, and U. J. Nunes, “Place recognition survey: An update on deep learning approaches,” *arXiv preprint arXiv:2106.10458*, 2021.
- [3] F. Tan, J. Yuan, and V. Ordonez, “Instance-level image retrieval using reranking transformers,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 12 105–12 115.
- [4] H. Zhang, X. Chen, H. Jing, Y. Zheng, Y. Wu, and C. Jin, “Etr: An efficient transformer for re-ranking in visual place recognition,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 5665–5674.
- [5] W. Chen, Y. Liu, W. Wang, E. M. Bakker, T. Georgiou, P. Fieguth, L. Liu, and M. S. Lew, “Deep learning for instance retrieval: A survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [6] Y. Fan, X. Xie, Y. Cai, J. Chen, X. Ma, X. Li, R. Zhang, and J. Guo, “Pre-training Methods in Information Retrieval,” *Foundations and Trends® in Information Retrieval*, vol. 16, no. 3, pp. 178–317, 2022.
- [7] M. A. Fischler and R. C. Bolles, “Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography,” *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981.
- [8] Z. Zhang, “Iterative point matching for registration of free-form curves and surfaces,” *International journal of computer vision*, vol. 13, no. 2, pp. 119–152, 1994.
- [9] M. A. Uy and G. H. Lee, “PointNetVLAD: Deep point cloud based retrieval for large-scale place recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4470–4479.
- [10] K. Vidanapathirana, M. Ramezani, P. Moghadam, S. Sridharan, and C. Fookes, “LoGG3D-Net: Locally guided global descriptor learning for 3d place recognition,” in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 2215–2221.
- [11] T. Barros, L. Garrote, P. Conde, M. Coombes, C. Liu, C. Premebida, and U. Nunes, “ORCHNet: A robust global feature aggregation approach for 3D LiDAR-based place recognition in orchards,” *arXiv preprint arXiv:2303.00477*, 2023.
- [12] F. Radenović, G. Toliás, and O. Chum, “Fine-tuning CNN image retrieval with no human annotation,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 7, pp. 1655–1668, 2018.
- [13] 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*, vol. 8, no. 5, pp. 2494–2501, 2023.
- [14] Z. Liu, S. Zhou, C. Suo, P. Yin, W. Chen, H. Wang, H. Li, and Y.-H. Liu, “LPD-Net: 3D point cloud learning for large-scale place recognition and environment analysis,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2831–2840.
- [15] L. Li, X. Kong, X. Zhao, T. Huang, W. Li, F. Wen, H. Zhang, and Y. Liu, “RINet: Efficient 3D LiDAR-based place recognition using rotation invariant neural network,” *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 4321–4328, 2022.
- [16] D. W. Shu and J. Kwon, “Hierarchical bidirected graph convolutions for large-scale 3-D point cloud place recognition,” *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [17] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, “NetVLAD: CNN architecture for weakly supervised place recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5297–5307.
- [18] H. Jégou, M. Douze, C. Schmid, and P. Pérez, “Aggregating local descriptors into a compact image representation,” in *2010 IEEE computer society conference on computer vision and pattern recognition*. IEEE, 2010, pp. 3304–3311.
- [19] T. Barros, L. Garrote, R. Pereira, C. Premebida, and U. J. Nunes, “AttDLNet: Attention-based deep network for 3D LiDAR place recognition,” in *ROBOT2022: Fifth Iberian Robotics Conference: Advances in Robotics, Volume 1*. Springer, 2022, pp. 309–320.
- [20] A. Babenko and V. Lempitsky, “Aggregating local deep features for image retrieval,” in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1269–1277.
- [21] M. Leordeanu and M. Hebert, “A spectral technique for correspondence problems using pairwise constraints,” in *Tenth IEEE International Conference on Computer Vision (ICCV’05) Volume 1*, vol. 2. IEEE, 2005, pp. 1482–1489.
- [22] M. A. Uy and G. H. Lee, “PointNetVLAD: Deep Point Cloud Based Retrieval for Large-Scale Place Recognition,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4470–4479.
- [23] P. Pobrotyn, T. Bartczak, M. Synowiec, R. Białobrzeski, and J. Bojar, “Context-aware learning to rank with self-attention,” *arXiv preprint arXiv:2005.10084*, 2020.
- [24] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [25] J. L. Ba, J. R. Kiros, and G. E. Hinton, “Layer normalization,” *arXiv preprint arXiv:1607.06450*, 2016.
- [26] C. Burges, T. Shaked, E. Renshaw, A. Lazier, M. Deeds, N. Hamilton, and G. Hullender, “Learning to rank using gradient descent,” in *Proceedings of the 22nd international conference on Machine learning*, 2005, pp. 89–96.
- [27] X. Wang, C. Li, N. Golbandi, M. Bendersky, and M. Najork, “The lambdaloss framework for ranking metric optimization,” in *Proceedings of the 27th ACM international conference on information and knowledge management*, 2018, pp. 1313–1322.
- [28] A. Geiger, P. Lenz, and R. Urtasun, “Are we ready for autonomous driving? The KITTI vision benchmark suite,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [29] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, “PointNet: Deep learning on point sets for 3D classification and segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.
- [30] J. Kumorowski, “MinkLoc3D: Point cloud based large-scale place recognition,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 1790–1799.