

# Point-cloud-based place recognition using CNN feature extraction

Ting Sun<sup>1</sup>, Ming Liu<sup>1</sup>, Haoyang Ye<sup>1</sup>, Dit-Yan Yeung<sup>2</sup>

**Abstract**—This paper proposes a novel point-cloud-based place recognition system that adopts a deep learning approach for feature extraction. By using a convolutional neural network pre-trained on color images to extract features from a range image without fine-tuning on extra range images, significant improvement has been observed when compared to using hand-crafted features. The resulting system is illumination invariant, rotation invariant and robust against moving objects that are unrelated to the place identity. Apart from the system itself, we also bring to the community a new place recognition dataset containing both point cloud and grayscale images covering a full 360° environmental view. In addition, the dataset is organized in such a way that it facilitates experimental validation with respect to rotation invariance or robustness against unrelated moving objects separately.

**Index Terms**—place recognition, point cloud, CNN

## I. INTRODUCTION

In autonomous driving, place recognition is to recognize a previously memorized place when the vehicle revisits it. With a stored map, place recognition can be used for localization when a high quality global positioning system (GPS) signal is unavailable. In simultaneous localization and mapping (SLAM), place recognition performs loop-closure detection, which is crucial for drifting error correction. In contrast to localization, place recognition only concerns the location of the robot, regardless of its orientation. A well-performing place recognition system is expected to correctly identify a previously visited place with a high probability in real-time.

Most of the place recognition methods are based on images [1]–[9], and few have reported using LiDAR [10]. The limitation of using images is that they are variant to illumination change, as cameras are passive photoreceptive sensors. LiDARs are active illumination invariant sensors. In this work we adopt Velodyne<sup>1</sup> for input, since its generated point-cloud covers 360° in the horizontal direction isotropically, and this enables us to design a rotation-invariance system, which is necessary for place recognition. As similar sensors become increasingly popular in autonomous driving [11]–[14] and the cost has been greatly reduced, it is promising and practical to conduct place recognition using similar sensors.

The traditional place recognition methods first detect key points, where some hand crafted features like SIFT [15] are

extracted. These features are commonly with high computational cost, and are further encoded by bag-of-words, then recognition is conducted by matching the encoded indexes [2], [16]. Since deep learning approaches have achieved state-of-the-art performance in many challenging computer vision tasks [17]–[21], there is a natural momentum to attempt place recognition using deep learning. In particular, convolutional neural networks (CNNs) [22] provide a powerful end-to-end framework for image-based vision tasks. Such a structure would also greatly inspire place recognition using LiDAR. The advantages of CNNs include the following:

- It directly takes raw images as input;
- It extracts hierarchical features automatically;
- Both the convolutional layer and pooling layer make the features of the higher layer shift invariant to some extent;
- Due to common characteristics of natural images, such as analogy, the CNN models trained on one dataset can be transferred to other datasets.

However, it is highly nontrivial to leverage a CNN to extract features from an unstructured point cloud. This is because images are typically a “dense representation”, for which each pixel in the image configuration space has a defined intensity value. Conversely, point-cloud is a “sparse representation”, for which not *all* (abstractly) points in the configuration space are defined. Only the locations with point observations are informative. Point clouds are much more expensive to obtain than color images, but training a deep CNN is data demanding. The features extracted by a deep CNN in different layers have different levels of abstraction, and its worth studying how the performance varies with the choice of layers. In brief, the solutions to the following issues are critical to use CNN with point-clouds:

- A mapping to convert an unstructured point cloud to an image is to be defined;
- The selection of a proper CNN model should be advised;
- The representation of the extracted deep feature should be associated with a layer of the network;
- The features directly obtained from a pre-trained CNN are redundant and noisy. Efficient and sufficient post-processing is necessary for more compact feature description.

In this paper, we propose a point-cloud-based place recognition method using a CNN for feature extraction. A point-cloud is first aligned with its principal directions, then converted to a range image. A CNN pre-trained on abundant RGB images is used to extract features from the range images. Principal component analysis (PCA) is used for further dimension reduction

<sup>1</sup> Department of Electronic & Computer Engineering, Hong Kong University of Science and Technology, Hong Kong, China. (tsun, eelium)@ust.hk

<sup>2</sup> Department of Computer Science, Hong Kong University of Science and Technology, Hong Kong, China. dyyeung@cse.ust.hk

<sup>1</sup><http://velodynelidar.com/>

following our previous guideline [23]. A scoring mechanism that concerns both cosine similarity and the discrimination of the best match among the top matches is used to make a decision.

A good place recognition system is expected to be invariant to illumination change and rotation, as well as to moving objects that are irrelevant to the place [24]. However, to the limit of our knowledge, there is no place recognition dataset that deliberately separates its content according to these three aspects in this context. In this paper, we introduce a place recognition dataset containing grayscale images (the images are not directly used in this paper) and point-clouds. The images are taken under severe illumination change and both types of data cover a full 360° environmental view. The content of our dataset is organized to especially facilitate tests of yaw rotation, which is the primary rotation case in autonomous driving, invariance, and robustness to moving objects separately. Details of the dataset are given in Subsec. III-D.

To summarize, we stress the following contributions in this paper:

- 1) We propose a novel end-to-end point-cloud-based place recognition system using CNN feature extraction, which is:
  - illumination invariant,
  - rotation invariant, and
  - robust to moving objects.
- 2) We analyze the properties of the system by testing the effectiveness of each module;
- 3) We introduce a new dataset for place recognition.

The remainder of this paper is organized as follows. Sec. II reviews some previous work on place recognition and the related CNN study. Our proposed method is presented in Sec. III, which is then followed by experiment results and discussion in Sec. IV. Sec. V concludes the paper.

## II. RELATED WORK

### A. Place recognition

Place recognition contains two main steps: feature extraction and feature retrieval. We separately introduce the related work in the two domains as follows.

1) *Feature extraction and description*: Descriptors are crucial, and improvement of illumination invariance, viewpoint invariance and calculation efficiency are the three main research directions of feature extraction and description.

In the traditional visual place recognition methods, descriptors roughly fall into two categories [25]: local descriptors and global descriptors (or holistic descriptors [26]). Local descriptors are extracted around detected keypoints like corners, while global descriptors describe the whole image. Widely adopted local descriptors include scale-invariant feature transforms (SIFT) [15], speeded-up robust features (SURF) [27], binary robust independent elementary features (BRIEF) [28], binary robust invariant scalable keypoints (BRISK) [29], oriented FAST and rotated BRIEF (ORB) [30], local difference binary (LDB) [31], [32] etc. Since the number of detected keypoints in each frame varies and directly matching the features can be inefficient [25], the bag-of-words model is used to further

encode the local features to facilitate frame-wise comparison. Global descriptors of a frame can be obtained by integrating the local descriptors [26], [33], or by directly extracting them from the whole image [34], [35]. A trade-off between local and global descriptors is explored in [36], [37], where features are extracted from the object region proposed by edge boxes [38], and [39] proposes a lightweight adaptive line descriptor based on color features and geometric information.

As a variety of range sensors become popular, place recognition starts to benefit from this new type of input. In [40], range value is used to obtain scale information assisting the detection of salient regions. [41] manually extracts 41 rotation invariant features from each frame of a 3D point cloud, and adopts AdaBoost to train a binary classifier to distinguish positive and negative laser pairs. [42] designs a local feature extracted from a point cloud, called a neighbor-binary landmark density descriptor (NBLD), and extracts the NBLD from detected keypoints to recognize places through a voting framework. However, point-cloud-based feature extraction methods are far from mature compared with image-based methods. It is more difficult to identify keypoints, lines, and objects in an unstructured point cloud than in an image, and the increased dimension also intensifies the computational cost. Some extraordinary frameworks like CNN are not designed for point clouds and it is definitely worth trying to design a point-cloud-based system so that the power of deep learning approaches can be leveraged.

2) *Feature Retrieval*: With the extracted descriptor of a new frame, retrieval is used to tell whether the current place matches a previously visited place, and if the answer is yes, where the place is. In most cases, the retrieval method is designed independently of descriptors, and is less time consuming to compute than feature extraction for most of the datasets, unless the stored map is of a very large scale.

Based on different concerns and assumptions, well-designed retrieval methods help to increase the precision-vs-recall performance. *E.g.* by exploring the structure of the bag-of-words data [2], [43], [44] or covisibility of landmarks [8], [9] to reduce perceptual aliasing; using the assumption that sequential frame queries are of adjacent or the same place as prior to narrow down the search range [2]; based on the assumption that the vehicle will repeat the same path rather than just revisit a single place, sequence matching is adopted [1], [3]–[7]; [4] propose a method to correct the wrong recognition when new information comes; the time consumed by many retrieval methods is proportional to the number of places stored, and this time can be reduced by narrowing down the search range as mentioned before, or using a kd-tree [45], hashing [46] etc.

Our focus in this work is point-cloud-based feature extraction using a CNN. In order to test the feature quality, we show precision-vs-recall performance using single scan feature matching without special prior. However, if the assumptions hold, any retrieval approaches can be adopted to further improve the performance.

## B. Convolutional neural networks

Deep learning [47], [48] approaches focus on learning features directly from raw data in a hierarchical manner. In particular, CNNs [22] provide a powerful end-to-end framework that achieves state-of-the-art performance in many challenging computer vision tasks [17]–[21], [49]. The rich features learned by a deep CNN ranging from low-level to high-level representations in the hidden layers have also aroused extensive research interest in investigating how to take advantage of them [49]–[52]. Moreover, the public availability of efficient CNN implementations [53], powerful pre-trained CNN models [17], [20], [21], [54], [55] and their ability to transfer to work on other tasks even without fine-tuning has further popularized the pervasive use of CNNs for various applications. Not surprisingly, visual place recognition also turns to CNNs for feature extraction [7], [25], [35]–[37], [56]. Our method is point-cloud-based, and in the next section we demonstrate how we leverage the feature extraction power of a CNN.

## III. PROPOSED METHOD

The overview of our proposed place recognition system is shown in Figure 1. A point cloud is first aligned with its principal directions, then is projected onto a cylinder image plane. After hole-closing, a CNN is used for feature extraction, followed by PCA dimension reduction. Retrieval is based on a score concerning both similarity and discrimination. A threshold applied to the score is used to trade off between precision and recall. The motivation and design details of each module are illustrated in the following. The last subsection will introduce our new data set for place recognition.

### A. Preprocessing

The preprocessing module is highlighted by a light blue rectangle in Figure 1. The output of this module is a range (i.e. grayscale) image that can be put into a CNN. The input we consider is a 3D point cloud that covers a full 360° environmental view. This kind of point cloud is commonly created by a Velodyne LiDAR. In order to align the point cloud, PCA is adopted to find the orthogonal directions sorted by the variance. However, the resultant bases are not unique, there are 8 possible combinations of signs. Considering the typical autonomous driving environment, we narrow down the cases to two. Let's denote the original basis of the point cloud as  $\mathbf{B} = [\mathbf{e}_x, \mathbf{e}_y, \mathbf{e}_z]$ , the basis obtained by PCA as  $\mathbf{B}' = [\mathbf{e}'_x, \mathbf{e}'_y, \mathbf{e}'_z]$  and  $\mathbf{B}' = \mathbf{T}\mathbf{B}$ . The two cases that we keep satisfy one of the two following constraints:

$$t_{11}, t_{22}, t_{33} \geq 0 \quad (1)$$

or

$$t_{11}, t_{22} < 0; t_{33} \geq 0, \quad (2)$$

where  $t_{ij}$  is the element of  $\mathbf{T}$  in  $i^{\text{th}}$  row and  $j^{\text{th}}$  column. Following the setup of the KITTI dataset [13], i.e.  $\mathbf{e}_x$  points to the front of the car and  $\mathbf{e}_z$  points up. For a common street

view point cloud, it is reasonable to assume in the PCA basis  $\mathbf{B}'$ ,  $\mathbf{e}'_x$  and  $\mathbf{e}'_y$  span a plane which is roughly parallel to the ground since they are the directions that capture most of the variance. In order to satisfy either (1) or (2),  $\mathbf{e}'_z$  will always point up, and  $[\mathbf{e}'_x, \mathbf{e}'_y]$  will be one of the cases in Figure 2. As will be shown in Sec. IV, for unidirectional loop closure, using one alignment is enough, considering both cases is mainly to handle bidirectional loop closure.

We then create the range image from the aligned point cloud using the Point Cloud Library (PCL) [57] implementation. One parameter we would like to mention is angular resolution<sup>2</sup>, which decides how fine-grained the grid in the cylinder image is when the range image is generated. If the angular resolution is set too small, there will be a lot of empty grids; and if it is too big, the grid will be too rough to preserve the details. We suggest setting the angular resolution similar to that of the LiDAR sensor used, then filling the small holes using morphological closing, i.e. dilation followed by erosion.

### B. CNN feature extraction

As mentioned in Sec. I, leveraging a CNN to extract features from an unstructured point cloud is highly nontrivial. It is laborious and expensive to collect a large amount of point cloud data to train a model from scratch, and we want to take advantage of the abundant pre-trained ones. Currently most of these CNN models are trained on RGB images of a square shape e.g. 224:224, but the range images we generate from point clouds are grayscale and of a very long rectangle shape about 4000:100. The range image is simply repeated 3 times to fit the color channels; but for the spatial resolution, naive resizing will severely distort the ratio, and brutal down sampling will cause great loss of details. We propose to preserve the original horizontal resolution of the image and may increase its vertical resolution to ensure that it will not reduce to zero through the pooling layers of a CNN.

A deep CNN trained on a large dataset can extract hierarchical generic features for other tasks [18]. From the lower layer to the higher layer, the precision of the hidden activation decreases while the abstraction and invariance increase [50]–[52]. We try to find the proper precision-abstraction trade-off by testing multiple layers of the networks [54], as presented in Sec. IV. Our range image is generated by projecting the point cloud on a cylinder plane then unfolding it, so the rotation of the point cloud around the axis of the cylinder becomes a shift in the resultant range image. The feature extracted by a CNN is spatially invariant to some extent, which contributes to the rotation invariance property of our system.

In general, a deeper and wider network is potentially more powerful [55]. We test our proposed method on three CNNs i.e. AlexNet [17], VGG-CNN-S [54] and Places-CNDS-8 [55]. VGG-CNN-S [54] is modified from AlexNet [17] by increasing the channels of the last three convolutional layers, and Places-CNDS-8 [55] has a deeper network structure. Because of the length of the input image, the size of the hidden layer

<sup>2</sup>The term ‘resolution’ here means the angle range corresponding to one pixel in the range image. In the next subsection, ‘resolution’ means the number of pixels of an image.

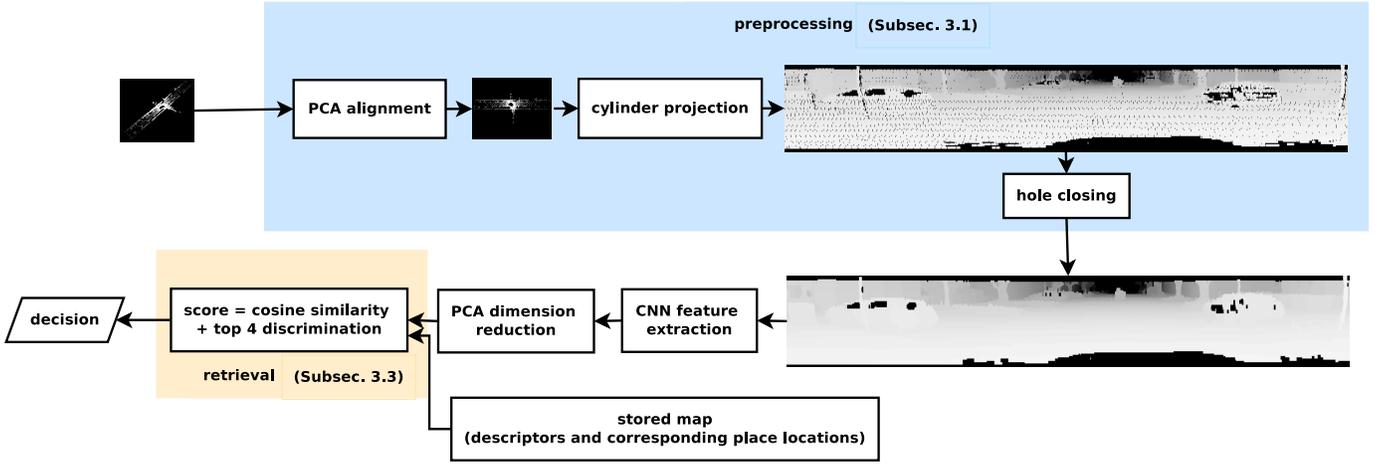


Fig. 1: System overview

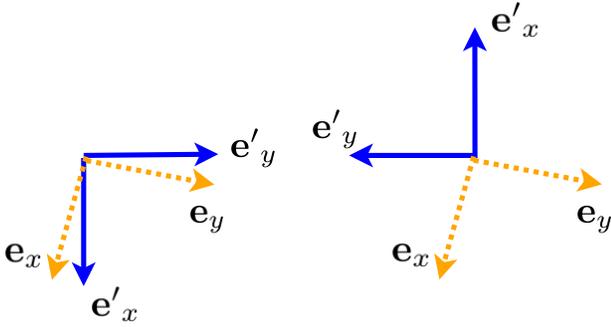


Fig. 2: Two cases of PCA alignment.

maps is large. Notice that although the linear transformations conducted by fully connected layers perform reshaping and dimension reduction, we cannot use the features processed by them; since the shape of our range image is different from that of the original CNN model input, the pre-trained weights of fully connected layers can not be loaded due to parameter number mismatch.

The features directly obtained from the CNN are highly redundant. The redundancy is introduced by:

- 1) the repetition of the 3 color channel input;
- 2) the uninformative region of the range image, like the boundary and the floor; and
- 3) the domain knowledge that only applies to the data used to train the CNN, but not to our application scenario.

In our system, features are extracted from convolutional layers or pooling layers. One more pooling layer may be manually inserted on top of the chosen hidden layer for preliminary dimension reduction. PCA is adopted as a postprocessing step to obtain more compact features. This is due to the nature of the place recognition task and the global descriptor we use. Since one place is supposed to have one descriptor, the variance of each dimension indicates its discrimination ability.

### C. Retrieval

We normalize the descriptor vector of each point cloud and use cosine distance as the similarity metric. Most image-based

retrieval methods only consider the best match in the memory. On the contrary, our method jointly considers the top  $k$  best match and shows that this design leads to a better result. The proposed retrieval method is described in Algorithm 1, where  $\mathbf{f}^i$  is the feature of place  $i$ ,  $\mathcal{S}_f^{stored}$  is the set of stored features of previously visited places,  $\mathbf{f}^{query}$  is the feature vector of the current scan,  $\mathcal{C}(\mathbf{f}^i, \mathbf{f}^j)$  returns the cosine similarity between  $\mathbf{f}^i$  and  $\mathbf{f}^j$ ,  $s$  is a score we define, and  $threshold$  is used to trade off between precision and recall.

---

#### Algorithm 1 Retrieval

---

**Input:**  $\mathbf{f}^{query}$ ,  $\mathcal{S}_f^{stored}$

**Output:** return place identity if match found, otherwise report not found

- 1:  $\mathcal{C}(\mathbf{f}^{query}, \mathbf{f}^{i_1}) > \mathcal{C}(\mathbf{f}^{query}, \mathbf{f}^{i_2}) \dots > \mathcal{C}(\mathbf{f}^{query}, \mathbf{f}^{i_k}) > \dots$   
 $\triangleright$  find top  $k$  match
  - 2:  $s = \mathcal{C}(\mathbf{f}^{query}, \mathbf{f}^{i_1}) \times 2 - \mathcal{C}(\mathbf{f}^{query}, \mathbf{f}^{i_k})$
  - 3: **if**  $s > threshold$  **then**
  - 4:     return  $i_1$
  - 5: **else**
  - 6:     report not found
- 

Notice that the score  $s = \mathcal{C}(\mathbf{f}^{query}, \mathbf{f}^{i_1}) \times 2 - \mathcal{C}(\mathbf{f}^{query}, \mathbf{f}^{i_k})$  is the sum of two terms: the cosine similarity of the best match  $\mathcal{C}(\mathbf{f}^{query}, \mathbf{f}^{i_1})$  and the gap between the best match and the  $k$ th best match  $\mathcal{C}(\mathbf{f}^{query}, \mathbf{f}^{i_1}) - \mathcal{C}(\mathbf{f}^{query}, \mathbf{f}^{i_k})$ . The intuition of this design is illustrated in Figure 3 where two typical retrieval examples in the KITTI dataset [13] are shown. We manually find the revisited trails and use one round as the stored map, colored blue, and leave the rest as queries colored green. A pair of scans is treated as a true match if all the differences of their 3 coordinates are smaller than 3 meters. In each example, the query and its top 5 matches, as well as their corresponding cosine similarity, are marked. The top example shows a case when the required place has a match in the memory, while the bottom example shows a case when there is no match. After comparison, our key observation is that when a true match exists, it has two properties: 1) it has high similarity with the query and 2) its cosine similarity distinguishes it from the rest by a large margin, which we call the ‘discrimination gap’ in

the following text. In our system, we empirically set  $k = 4$ . We do not use the gap between the 1st match and the 2nd, considering that it is rare for the vehicle to revisit exactly the same location. It is more likely to be somewhere between the stored places so there could be a few true matches instead of a unique one. The effectiveness of this design is shown in Figure 4. The left subfigure shows sequence 05 from the KITTI dataset [13]. The middle subfigure plots the histogram of the cosine similarity of all the matched scan pairs in blue, and those that do not match in orange. The right subfigure plots the histogram of our proposed score. Comparing with using cosine similarity only, it can be seen that in the histogram of our proposed score, the blue bars are slightly pushed to the right, while the orange bars have no visible change. This result is consistent with our observation that the discrimination gap is a property of the true match. Thus, modifying cosine similarity with the discrimination gap makes the matched scan pairs more distinguishable from the unmatched ones. This effect is also shown in Subsec. IV-D.

#### D. HKUST dataset

We capture our dataset using an omni stereo camera and a VLP-16 Velodyne LiDAR tied together and placed on a tripod, as shown in Figure 5. Each pair of shots contain a grayscale image and a point cloud. Both types of data cover a full  $360^\circ$  environmental view and they are synchronized within 1 millisecond. At each location (i.e. within a 1 meter shift), we collect two sets of data: one for rotation invariance testing, and one for robustness to unrelated (moving) objects testing. When collecting the first set of data, we turn the tripod about  $40^\circ$  for each pair of shots, and repeat this process about 10 times at one location, completing a  $360^\circ$  circle; while for the collection of the second set of data, the equipment is fixed, and we take a pair of shots when there are people or cars passing by. This is also repeated about 10 times at each location. We collect data in 7 different locations on the Hong Kong University of Science and Technology (HKUST) campus under different lightening conditions, gaining 171 pairs of shots in total. Among them, 77 pairs belong to the rotation invariance testing set and 94 pairs are for unrelated object testing. Examples from the 7 locations in our HKUST dataset are shown in Figure 6. The top images show the point cloud scans, the middle images show the grayscale images taken at the same location and at the same time, and the bottom images show the same location when unrelated objects are passing by.

## IV. EXPERIMENTS

We thoroughly test our proposed system in this section. The performance of features extracted from different CNNs and different layers are shown first, then the effectiveness of each module in our system is tested using the best performing layers. We test our system mainly on the KITTI dataset [13], which offers the ground truth locations for 11 sequences numbered from 00 to 10. Among the 11 sequences, sequences 01, 03, 04, 07, 09 and 10 do not contain a significant (i.e. less than 100 scans) revisited path; sequences 00, 05 and 06 only contain paths revisited from the same direction i.e.

unidirectional loop closures; sequence 08 only contains paths revisited from the opposite directions i.e. bidirectional loop closures; and sequence 02 has revisited paths from both directions. In order to clearly demonstrate the performance of our system, we use sequences 00, 05, 06 and 08. We classify our place recognition tasks into 3 classes according to their difficulty. In the easiest task, the stored scans are uniformly sampled from the sequences, so most of the true matches are only different from the queries by a small location shift and rotation. The tasks with median difficulty use one trail of the manually-selected unidirectional loop closure paths as stored scans, so most of the queries do not have a true match and the existing matches differ from the queries by a large location shift and rotation. The difficult tasks are to recognize the revisited places in bidirectional loop closures. Many works are tested on the easy tasks, the median difficulty ones are more practical, and few works can handle the difficult tasks. (Notice that in order to recognize revisited places under a bidirectional loop closure situation, the input sensor must cover a full  $360^\circ$  degree environmental view.) We found that for unidirectional loop closures, only using the alignment that satisfies (1) is sufficient to achieve good performance. We then show in Subsec. IV-F that bidirectional loop closures require considering both (1) and (2), then we pick the one with higher matching score. The revisited paths in each sequence are manually selected and are shown in figure 7, and the stored scans' indexes are listed in Table I

sequence	loop closure type	stored scans' index
00	unidirectional	1-200, 3280-3840
05	unidirectional	13-150, 530-885
06	unidirectional	1-276
08	bidirectional	1420-1503, 1620-1840

TABLE I: Manually selected stored scans' indexes.

We run our C++ code on a desktop with 4 cores @ 3.20 GHz without multi-threading for CPU implementation, and adopt Caffe [53] for CNN feature extraction. We also use GPU version of Caffe [53] on a desktop with 4 cores @ 3.07 GHz with a NVIDIA GTX 980.

#### A. Performance of features from different layers of different CNNs

As mentioned in Subsec III-B, we test our proposed method on AlexNet [17], VGG-CNN-S [54] and Places-CNDS-8 [55]. VGG-CNN-S [54] is modified from AlexNet [17] by increasing the channels of the last three convolutional layers, and Places-CNDS-8 [55] has a deeper network structure. For AlexNet [17], we test the feature from *conv3*, *conv4* and *pool5*. When extracting the features of *conv3* and *conv4*, we insert a pooling layer after them that has the same parameters as *pool5* for preliminary dimension reduction, and name the resultant features *conv3+pool5* and *conv4+pool5* respectively. Similarly, we use the features *conv3+pool5*, *conv4+pool5* and *pool5* from VGG-CNN-S [54]; and *pool3+pool5*, *pool4+pool5* and *pool5* from Places-CNDS-8 [55]. We also show the performance of directly using the vectorized range image as

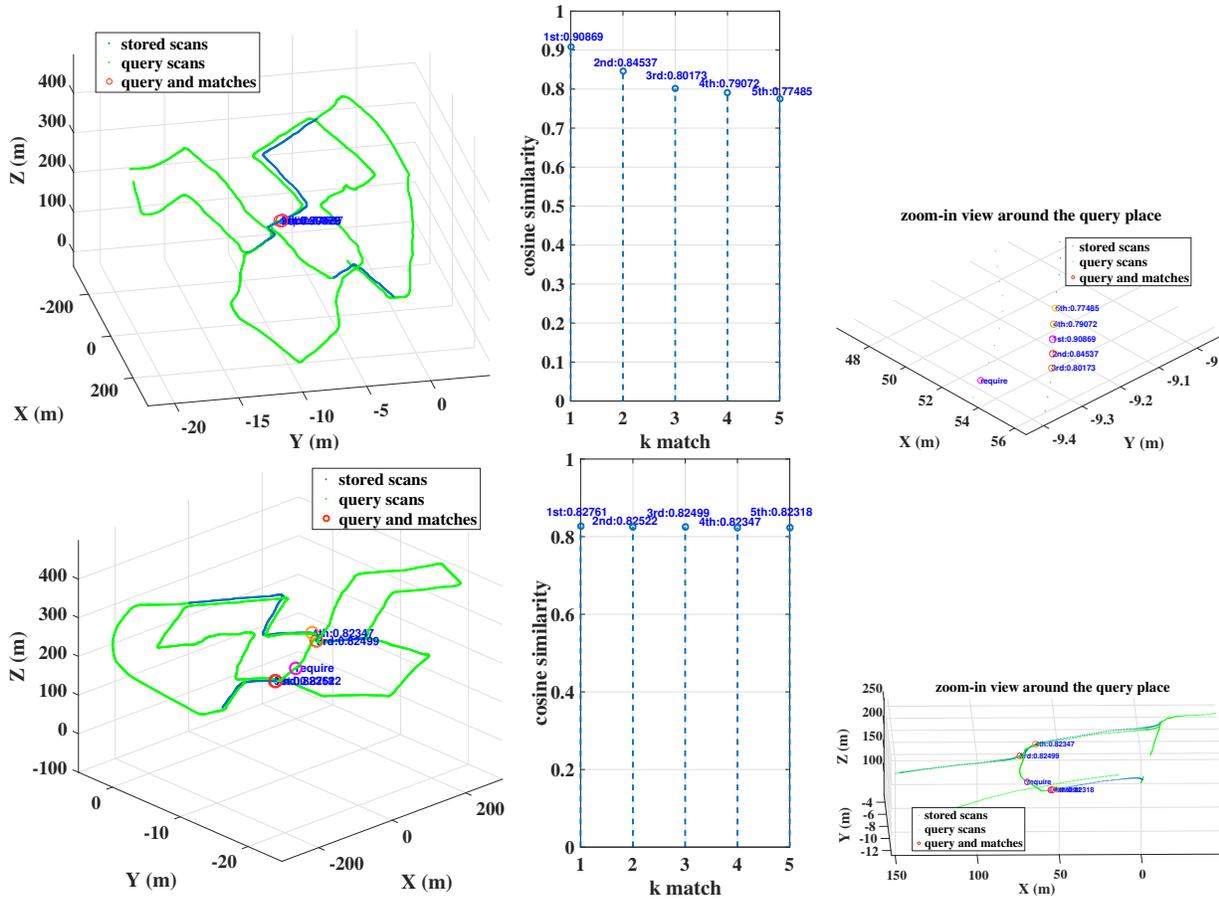


Fig. 3: Two typical retrieval examples in the sequence 00 of the KITTI dataset [13] using our proposed method, to illustrate the motivation of the design of score  $s$  in Algorithm 1. The blue dots are the manually selected stored scans and the green dots are queries. The top one shows a case when the required place feature has a match in the memory, while the bottom one shows a case when there is no match. The left subfigure shows the location of the query and that of its top 5 matches. The middle subfigure shows the cosine similarity between the query feature and its top 5 matches in the memory. The right subfigure shows the zoom-in view around the query place. A more detailed description can be found in the text.

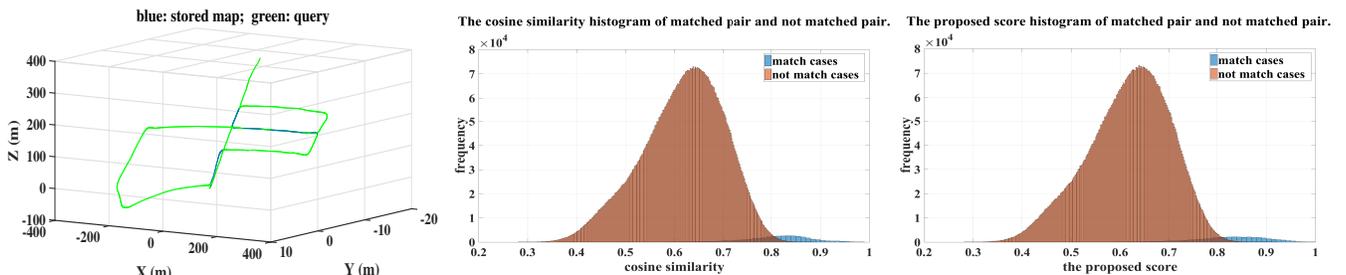


Fig. 4: The effectiveness of proposed score is shown in this figure. The left subfigure shows sequence 05 from the KITTI dataset [13], where the blue dots are the manually selected stored scans and the green dots are the queries. The middle subfigure plots the histogram of the cosine similarity of all the matched scan pairs in blue, and scans that do not match in orange. The right subfigure plots the histogram of our proposed score.



Fig. 5: The equipment used to collect our dataset. We use an omni stereo camera and a VLP-16 Velodyne LiDAR tied together and placed on a tripod.

the global descriptor for comparison. As shown in Figure 8, all the CNN features outperform the vectorized range images. Both the wider VGG-CNN-S [54] and deeper Places-CNDS-8 [55] have a better performance than AlexNet [17]. The best performing layers differ from network to network. In AlexNet [17], *conv3+pool5* offers the best features, in VGG-CNN-S [54] it is *conv4+pool5* and in Places-CNDS-8 [55] it is *pool5*.

### B. Easy task vs median difficulty task

In the easy task, we uniformly sample every 3rd scene in a sequence as a stored map, and leave the rest as queries; while the median difficulty task stores one path in the manually selected unidirectional loop closures. The results of two types of tasks are plotted together in Figure 9. It can be seen that for sequence 00 and sequence 05, the easy task achieves much better performance than the median difficulty task, especially in the high recall region. Sequence 06 roughly shapes like a rectangle (see Figure 7), and the car passes two of the edges at high speed and the other two at relatively low speed. The unidirectional loop closure is completely contained in the ‘low-speed edge’, where the scans densely cover the road. Using uniform sampling, we actually create a lot of true matches along the ‘high-speed edge’, where even the consecutive scans have a large position change, and we also create a lot of unmatched but close pairs at the same time. Thus, the performance of ‘easy task’ in sequence 06 drops.

### C. Effectiveness of PCA alignment

The precision-vs-recall curves with and without the PCA alignment mentioned in Subsec. III-A are plotted together in Figure 10, where each row contains the performance of the same feature of different test sequences. It can be seen that PCA alignment in our preprocessing module significantly helps to promote the performance in all cases. Another observation is that which regions on the curves show more

increase depends on the sequences. Except sequence 00, the performance of both sequence 05 and 06 is promoted obviously in the high and middle part of the curves, but tends to merge in the low precision region. This is a common situation when the similarity between true match pairs is decreased by the direction misalignment, and our PCA alignment module eliminates this effect. Sequence 00 contains the most sharp turnings, and it can be seen that without PCA alignment, some of the true match pairs have such a low score that they fail to stand out among unmatched pairs even when the threshold is lowered, leaving the gap between two curves in the low precision region.

### D. Retrieval using cosine similarity only vs using proposed score

As discussed in Subsec. III-C, during place retrieval, we jointly consider the top 4 best matches according to cosine similarity of the features, and threshold a score modified from cosine similarity by adding a discrimination gap. The effectiveness of this design is shown in Figure 11. It can be seen that using the proposed score leads to slightly better performance in most of the cases, especially in the high precision region.

### E. PCA dimension reduction

In place recognition, each place is expected to have one unique and discriminative descriptor, so the variance of the feature indicates its discrimination ability. Thus, in our system, PCA is adopted as a postprocessing step to obtain more compact features. The precision-vs-recall performance with different numbers of remaining dimensions is shown in Figure 12 (better viewed in color). It can be seen that in most of the cases, the curves are quite close to each other, and the performance remains approximately the same even when the dimension is reduced from around 100,000 to a few hundred. This result is consistent with our analysis, i.e. 1) the feature extracted by a pre-trained CNN has redundancy; and 2) after PCA dimension reduction, the main useful information of the features remains. Notice that the lower dimensional features can outperform the higher dimensional ones. This means that the PCA actually denoises the raw CNN features, e.g. there may be some cars on the road in a few scans, but in most of the scans the road region is flat and the variance corresponding to these elements is small, so that they are pruned during dimension reduction, eliminating the disturbance of the cars at the same time.

### F. Place recognition in bidirectional loop closure

In this subsection, we test our system on sequence 08, which mainly contains bidirectional loop closures. We tried using condition (1) or (2) alone, and considering both then choosing the match with the higher score. The results are shown in Figure 13. It can be seen that only using one alignment case fails in this sequence, but combining both solves the problem. We suggest that if the prior knowledge of only the unidirectional loop closure existing is available,

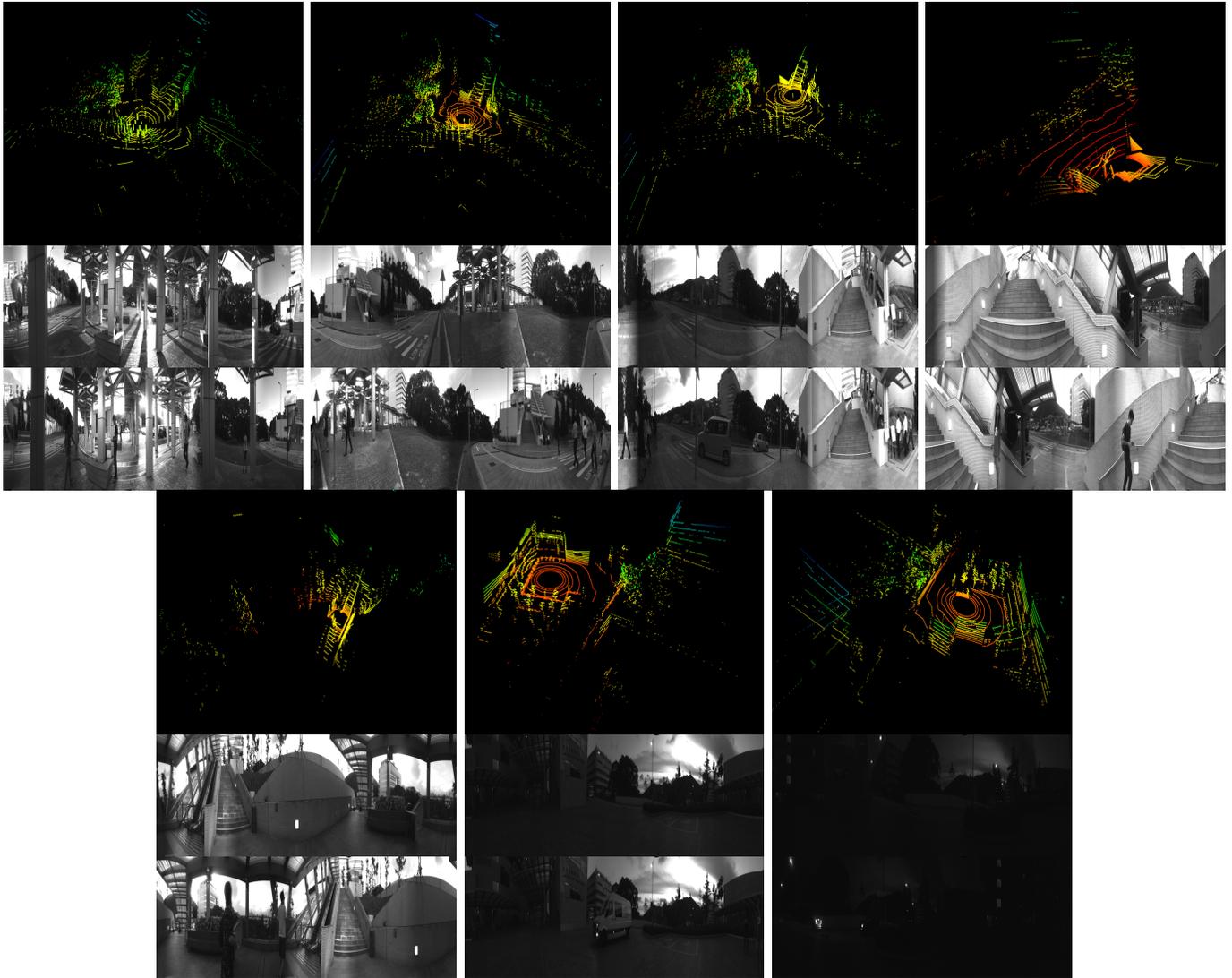


Fig. 6: Examples of the 7 locations in our HKUST dataset. The top images show the point cloud scans, the middle images show the grayscale images taken at the same location and at the same time, and the bottom images show the same location with unrelated objects passing by.

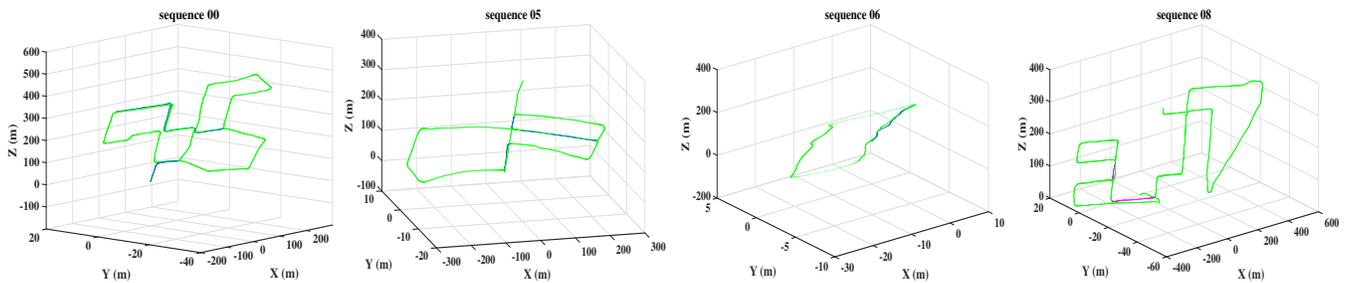


Fig. 7: The manually selected revisited paths in each sequence are shown in this figure. The unidirectional loop closure trails are blue and the bidirectional loop closure trails are pink. The precise blue/pink scans' indexes are listed in Table I. All the green scans are queries.

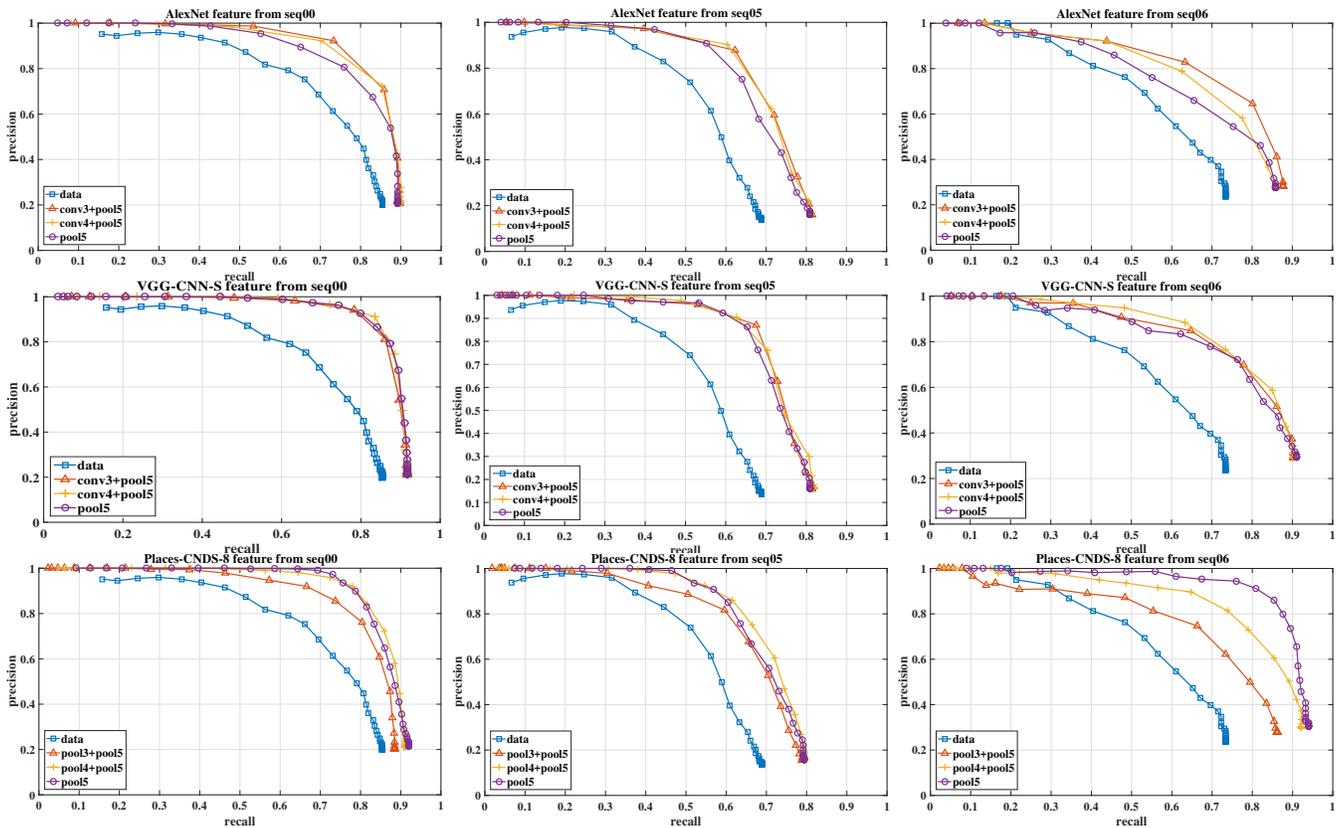


Fig. 8: Precision-vs-recall performance of features extracted from different CNNs and different layers. ‘data’ means directly using the vectorized range image as a feature, ‘conv3+pool5’ means extracting the feature from the ‘conv3’ layer, and using a pooling layer with the same structure as the ‘pool5’ layer for preliminary dimension reduction. Similarly for ‘conv4+pool5’, ‘pool3+pool5’ and ‘pool4+pool5’.

we can use either (1) or (2) alone to generate one descriptor for each frame, otherwise consider using both, which doubles the computation.

### G. Comparison with other methods

Our proposed method is compared with the famous FAB-MAP [2], and 3 hand-crafted global features of point clouds: Viewpoint Feature Histogram (VFH) [58], Ensemble of Shape Functions (ESF) [59] and Global Radius-based Surface Descriptor (GRSD) [60]. We use the openFABMAP [61] implementation in opencv library [62]. For all the hand-crafted features of point clouds we adopt the implementation in the Point Cloud Library (PCL) [57] and normalize the features, then use cosine similarity for place retrieval. The experiment is conducted on sequence 00, 05 and 06 of the KITTI dataset [13] under ‘median difficulty’ task settings. The properties of each method are summarized in Table II, and the precision-vs-recall performance of all the methods is shown in Figure 14. It can be seen from Figure 14 that our proposed CNN features significantly outperform the hand-crafted point clouds features and FAB-MAP [2] with comparable speed. In Table II, notice that both FAB-MAP [2] and our proposed method need training data. However, FAB-MAP [2] needs the testing data to be similar to the training data, which is expensive to collect for robotic application, while our method leverages the transfer

learning ability of CNN, i.e. using the models pre-trained on abundant RGB images to extract features from a point cloud.

### H. Test on the HKUST dataset

As mentioned in Subsec. III-D, we collect two sets of data to test rotation invariance and robustness to small objects that are unrelated to the place identity. We conduct place recognition on the two sets of data separately. In each experiment, we sample 1/3, 1/6 and 1/10 of the scans as a stored map, and leave the rest as queries. Since every one among the 7 locations has at least one scan stored, true matches exist for every query. More stored scans means the query can find more similar matches, resulting in an easier task. We collect about 10 scans at each location, so storing every 10th scan means that on average, each place has one scan stored, and all the remaining scans should be able to match it. The precision-vs-recall performance of rotation invariant testing and robustness to small unrelated objects are shown in Figure 15 and Figure 16, respectively. It can be seen from Figure 15 that our system achieves very good performance under rotation (notice that the precision/vertical axis starts from 0.8), which is due to the PCA alignment and the spacial-invariant property of the CNN features. Figure 16 shows that our system can always recognize the correct place (notice that the precision/vertical axis starts form 0.9), which indicates that small unrelated

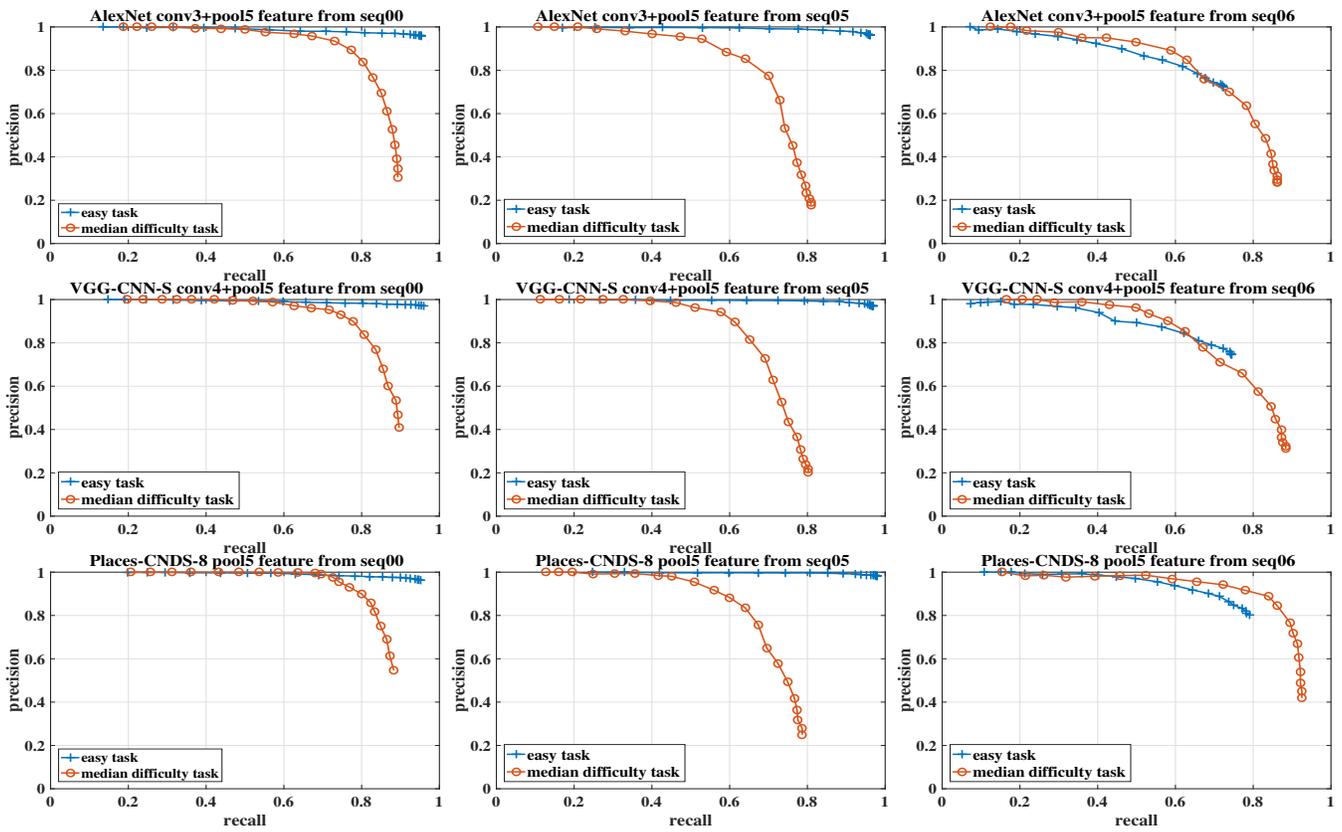


Fig. 9: Precision-vs-recall curves of easy task and median difficulty task are plotted together in this figure. A detailed description can be found in the text.

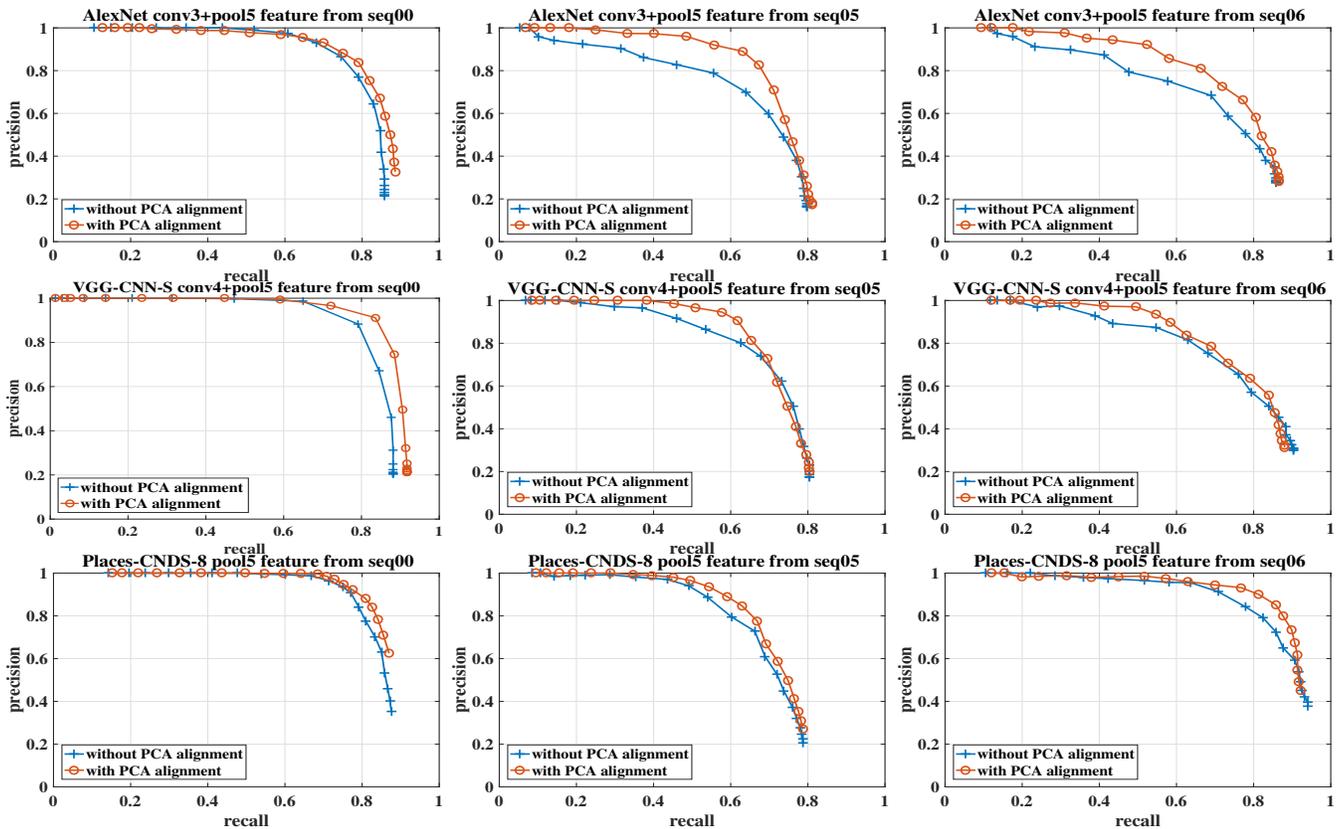


Fig. 10: The precision-vs-recall curves with and without PCA alignment mentioned in Subsec. III-A are plotted together in this figure.

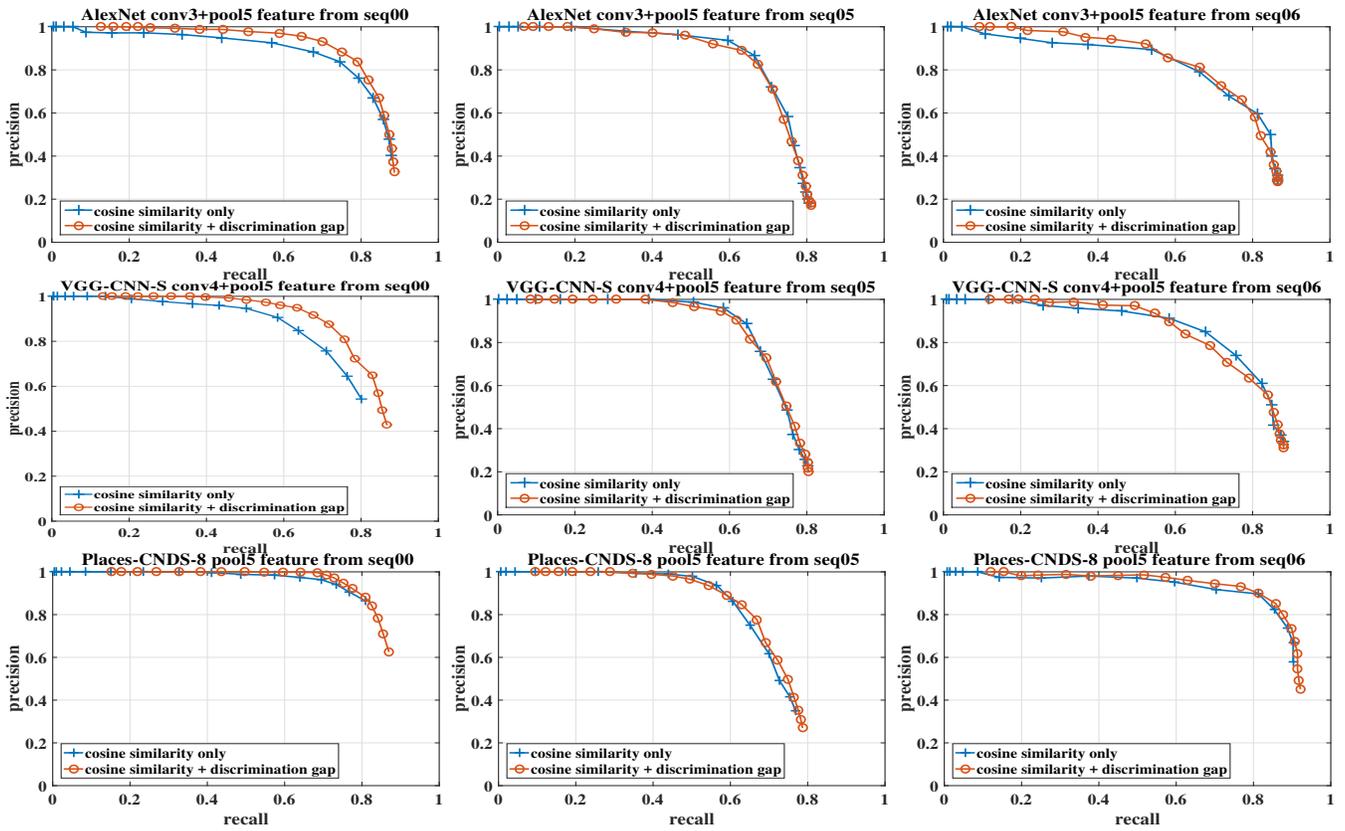


Fig. 11: Retrieval using cosine similarity only vs using our proposed score = cosine similarity + discrimination gap.

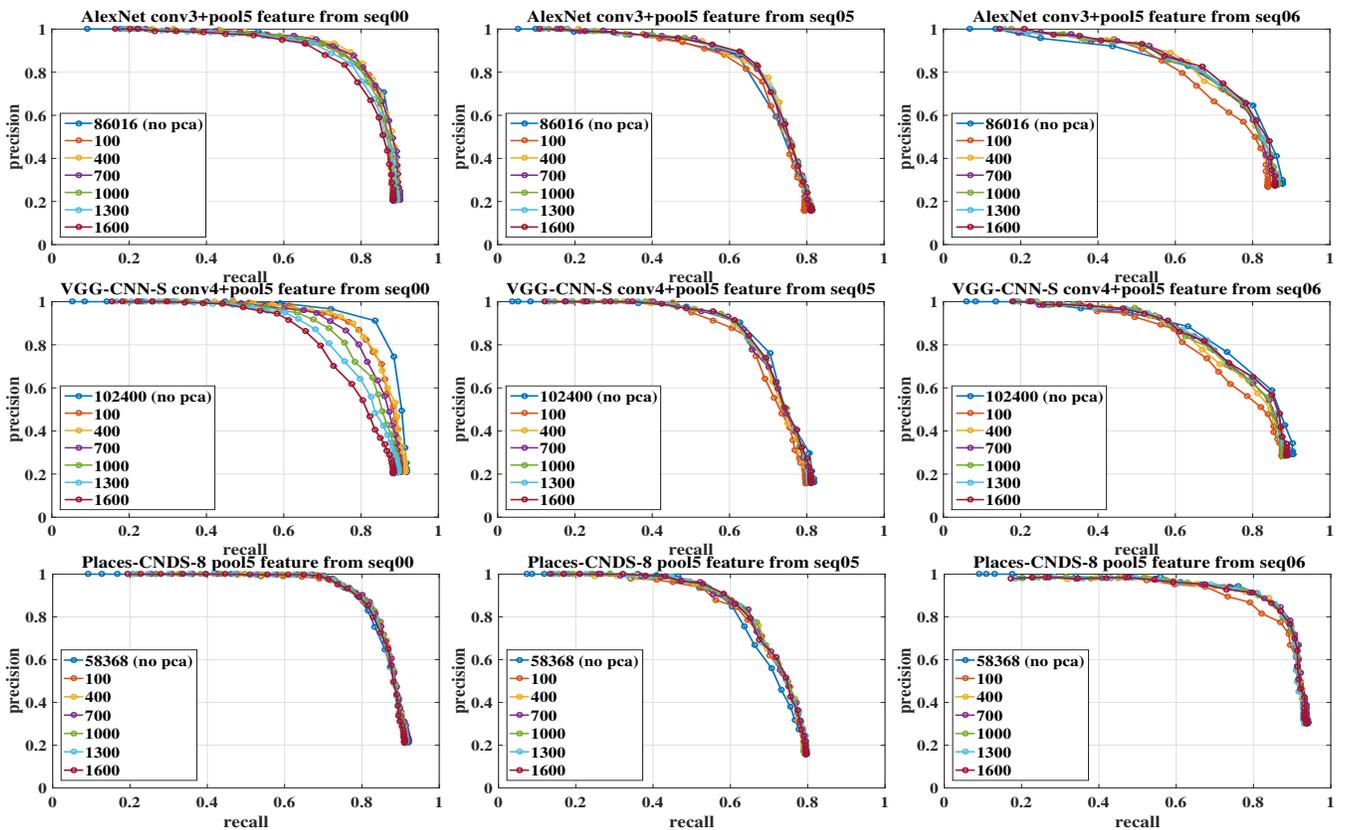


Fig. 12: The precision-vs-recall performance with different remaining dimensions after PCA. The original feature dimension and the remaining dimensions after PCA are listed in the legend.

method	feature dimension	need training data	parameters (if any)	time (frames/sec)
VFH [58]	308	No	normal estimation search radius = 0.03	4.8
ESF [59]	640	No	–	18.0
GRSD [60]	21	No	normal estimation search radius = 0.05, GRSD search radius = 0.1	1.6
FAB-MAP [2]	400	Yes	feature: SIFT, vocabulary size = 400	8.7
proposed	400	Yes	PCA remaining dimension = 400	2.0 (CPU), 5.4 (GPU)

TABLE II: Summary of the the properties of the following methods: Viewpoint Feature Histogram (VFH) [58], Ensemble of Shape Functions (ESF) [59], Global Radius-based Surface Descriptor (GRSD) [60], FAB-MAP [2], and our proposed method. The GPU time of our proposed method is obtained by using the GPU version of Caffe [53] for CNN feature extraction.

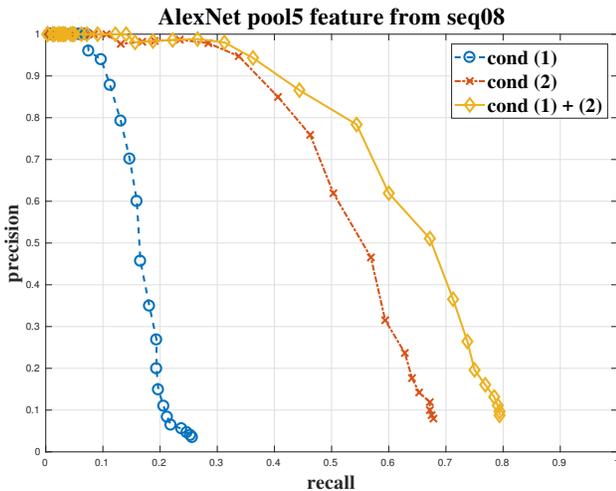


Fig. 13: The precision-vs-recall performance with different PCA alignments.

objects basically do not affect our place recognition system. For the sensor that covers a full 360° environmental view, the proportion occupied by the unrelated objects like pedestrian and cars is very small, in other words, these sensors are unlikely to be severely blocked, which results in a robust system.

## V. CONCLUSION

In this paper, we proposed a novel point-cloud-based place recognition system leveraging CNN feature extraction. Our method bridges the gap between powerful image-based deep learning approaches and point-cloud recognition. Without using any range images for training, the CNN features obtained in our system significantly outperform hand-crafted features and the resultant system is illumination invariant, rotation invariant and robust to unrelated small moving objects. We also introduce a new place recognition dataset containing both point cloud and grayscale images. Both types of data cover a full 360° environmental view, and the content of the dataset is organized to especially facilitate tests of rotation invariance and robustness to unrelated (moving) objects separately.

## REFERENCES

- [1] K. L. Ho and P. Newman, “Detecting loop closure with scene sequences,” *International Journal of Computer Vision*, vol. 74, no. 3, pp. 261–286, 2007.
- [2] M. Cummins and P. Newman, “Fab-map: Probabilistic localization and mapping in the space of appearance,” *The International Journal of Robotics Research*, vol. 27, no. 6, pp. 647–665, 2008.
- [3] M. J. Milford and G. F. Wyeth, “Seqslam: Visual route-based navigation for sunny summer days and stormy winter nights,” in *proceedings of IEEE International Conference on Robotics and Automation*. IEEE, 2012, pp. 1643–1649.
- [4] Y. Latif, C. Cadena, and J. Neira, “Robust loop closing over time for pose graph slam,” *The International Journal of Robotics Research*, vol. 32, no. 14, pp. 1611–1626, 2013.
- [5] P. Hansen and B. Browning, “Visual place recognition using hmm sequence matching,” in *proceedings of IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2014, pp. 4549–4555.
- [6] R. Arroyo, P. F. Alcantarilla, L. M. Bergasa, and E. Romera, “Towards life-long visual localization using an efficient matching of binary sequences from images,” in *proceedings of IEEE International Conference on Robotics and Automation*. IEEE, 2015, pp. 6328–6335.
- [7] T. Naseer, M. Ruhnke, C. Stachniss, L. Spinello, and W. Burgard, “Robust visual slam across seasons,” in *proceedings of IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2015, pp. 2529–2535.
- [8] E. S. Stumm, C. Mei, and S. Lacroix, “Building location models for visual place recognition,” *The International Journal of Robotics Research*, vol. 35, no. 4, pp. 334–356, 2016.
- [9] S. Cascianelli, G. Costante, E. Bellocchio, P. Valigi, M. L. Fravolini, and T. A. Ciarfuglia, “Robust visual semi-semantic loop closure detection by a covisibility graph and cnn features,” *Robotics and Autonomous Systems*, vol. 92, pp. 53–65, 2017.
- [10] M. Bosse and R. Zlot, “Place recognition using keypoint voting in large 3d lidar datasets,” in *proceedings of IEEE International Conference on Robotics and Automation*. IEEE, 2013, pp. 2677–2684.
- [11] W. Zhang, “Lidar-based road and road-edge detection,” in *Intelligent Vehicles Symposium (IV), 2010 IEEE*. IEEE, 2010, pp. 845–848.
- [12] Q. Zhu, L. Chen, Q. Li, M. Li, A. Nüchter, and J. Wang, “3d lidar point cloud based intersection recognition for autonomous driving,” in *Intelligent Vehicles Symposium (IV), 2012 IEEE*. IEEE, 2012, pp. 456–461.
- [13] A. Geiger, P. Lenz, and R. Urtasun, “Are we ready for autonomous driving? the kitti vision benchmark suite,” in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. IEEE, 2012, pp. 3354–3361.
- [14] A. Ibsch, S. Stümper, H. Altinger, M. Neuhausen, M. Tschentscher, M. Schlipfing, J. Salinen, and A. Knoll, “Towards autonomous driving in a parking garage: Vehicle localization and tracking using environment-embedded lidar sensors,” in *Intelligent Vehicles Symposium (IV), 2013 IEEE*. IEEE, 2013, pp. 829–834.
- [15] D. G. Lowe, “Object recognition from local scale-invariant features,” in *Proceedings of the International Conference on Computer Vision*, vol. 2. Ieee, 1999, pp. 1150–1157.
- [16] 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 *proceedings of IEEE/RSJ International*

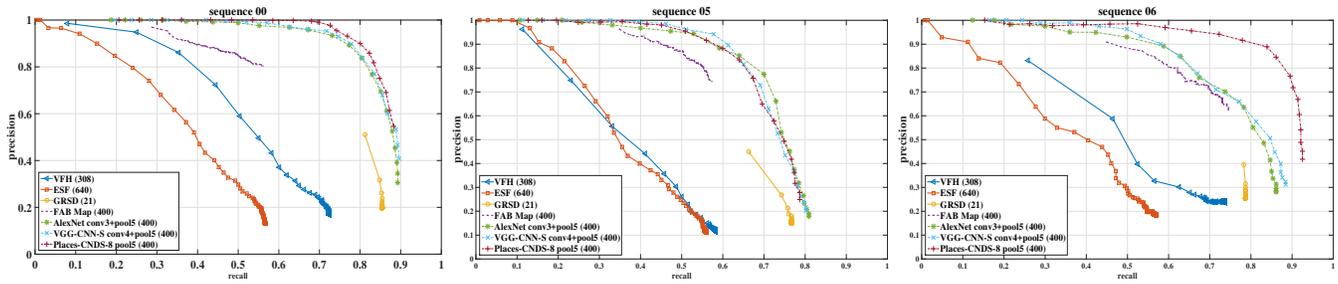


Fig. 14: The precision-vs-recall performance with different methods. The feature dimension of each method is annotated in the brackets.

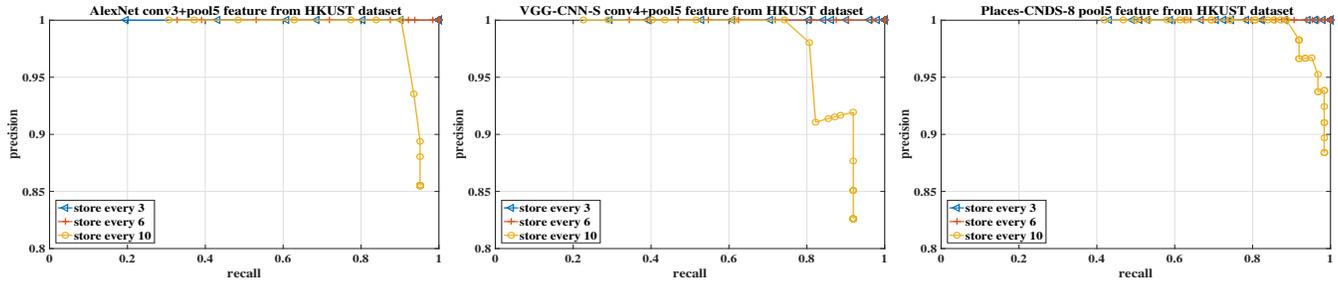


Fig. 15: The precision-vs-recall performance on the rotation invariance testing dataset of HKUST.

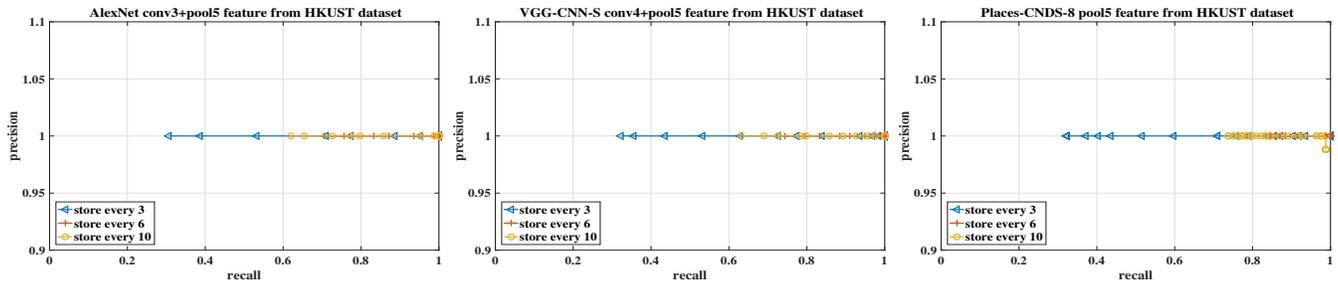


Fig. 16: The precision-vs-recall performance on the robustness to unrelated small objects testing dataset of HKUST.

*Conference on Intelligent Robots and Systems*. IEEE, 2011, pp. 1249–1255.

[17] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[18] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell, “Decaf: A deep convolutional activation feature for generic visual recognition,” in *International conference on machine learning*, 2014, pp. 647–655.

[19] A. Sharif Razavian, H. Azizpour, J. Sullivan, and S. Carlsson, “Cnn features off-the-shelf: an astounding baseline for recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2014, pp. 806–813.

[20] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.

[21] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1–9.

[22] Y. LeCun, B. E. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. E. Hubbard, and L. D. Jackel, “Handwritten digit recognition with a back-propagation network,” in *Proceedings of Conference on Neural Information Processing Systems (NIPS)*, 1990, pp. 396–404.

[23] L. Tai, Q. Ye, and M. Liu, “Pca-aided fully convolutional networks for semantic segmentation of multi-channel fmri,” *arXiv preprint arXiv:1610.01732*, 2016.

[24] Y. Sun, M. Liu, and M. Q.-H. Meng, “Improving rgb-d slam in dynamic environments: A motion removal approach,” *Robotics and Autonomous Systems*, vol. 89, pp. 110–122, 2017.

[25] S. Lowry, N. Sünderhauf, P. Newman, J. J. Leonard, D. Cox, P. Corke, and M. J. Milford, “Visual place recognition: A survey,” *IEEE Transactions on Robotics*, vol. 32, no. 1, pp. 1–19, 2016.

[26] N. Sünderhauf and P. Protzel, “Brief-gist-closing the loop by simple means,” in *proceedings of IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2011, pp. 1234–1241.

[27] H. Bay, T. Tuytelaars, and L. Van Gool, “Surf: Speeded up robust features,” in *Proceedings of the European Conference on Computer Vision*. Springer, 2006, pp. 404–417.

[28] M. Calonder, V. Lepetit, P. Fua, K. Konolige, J. Bowman, and P. Michelich, “Compact signatures for high-speed interest point description and matching,” in *Computer Vision, 2009 IEEE 12th International Conference on*. IEEE, 2009, pp. 357–364.

[29] S. Leutenegger, M. Chli, and R. Y. Siegwart, “Brisk: Binary robust invariant scalable keypoints,” in *Proceedings of the International Conference on Computer Vision*. IEEE, 2011, pp. 2548–2555.

[30] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, “Orb: An efficient alternative to sift or surf,” in *Proceedings of the International Conference on Computer Vision*. IEEE, 2011, pp. 2564–2571.

[31] X. Yang and K.-T. Cheng, “Ldb: An ultra-fast feature for scalable augmented reality on mobile devices,” in *International Symposium on Mixed and Augmented Reality (ISMAR)*. IEEE, 2012, pp. 49–57.

[32] R. Arroyo, P. F. Alcantarilla, L. M. Bergasa, J. J. Yebe, and S. Bronte, “Fast and effective visual place recognition using binary codes and disparity information,” in *proceedings of IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2014, pp. 3089–3094.

[33] H. Badino, D. Huber, and T. Kanade, “Real-time topometric localization,” in *proceedings of IEEE International Conference on Robotics and Automation*. IEEE, 2012, pp. 1635–1642.

- [34] I. Ulrich and I. Nourbakhsh, "Appearance-based place recognition for topological localization," in *proceedings of IEEE International Conference on Robotics and Automation*, vol. 2. IEEE, 2000, pp. 1023–1029.
- [35] N. Sünderhauf, S. Shirazi, F. Dayoub, B. Upcroft, and M. Milford, "On the performance of convnet features for place recognition," in *proceedings of IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2015, pp. 4297–4304.
- [36] N. Sünderhauf, S. Shirazi, A. Jacobson, F. Dayoub, E. Pepperell, B. Upcroft, and M. Milford, "Place recognition with convnet landmarks: Viewpoint-robust, condition-robust, training-free," *Proceedings of Robotics: Science and Systems XII*, 2015.
- [37] S. Cascianelli, G. Costante, E. Bellocchio, P. Valigi, M. L. Fravolini, and T. A. Ciarfuglia, "Robust visual semi-semantic loop closure detection by a covisibility graph and cnn features," *Robotics and Autonomous Systems*, volume=92, pages=53–65, year=2017, publisher=Elsevier.
- [38] C. L. Zitnick and P. Dollár, "Edge boxes: Locating object proposals from edges," in *Proceedings of the European Conference on Computer Vision*. Springer, 2014, pp. 391–405.
- [39] M. Liu and R. Siegwart, "Topological mapping and scene recognition with lightweight color descriptors for an omnidirectional camera," *IEEE Transactions on Robotics*, vol. 30, no. 2, pp. 310–324, 2014.
- [40] P. Newman and K. Ho, "Slam-loop closing with visually salient features," in *proceedings of IEEE International Conference on Robotics and Automation*. IEEE, 2005, pp. 635–642.
- [41] K. Granström and T. B. Schön, "Learning to close the loop from 3d point clouds," in *proceedings of IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2010, pp. 2089–2095.
- [42] T. Cieslewski, E. Stumm, A. Gawel, M. Bosse, S. Lynen, and R. Siegwart, "Point cloud descriptors for place recognition using sparse visual information," in *proceedings of IEEE International Conference on Robotics and Automation*. IEEE, 2016, pp. 4830–4836.
- [43] M. Liu and R. Siegwart, "Dp-fact: Towards topological mapping and scene recognition with color for omnidirectional camera," in *proceedings of IEEE International Conference on Robotics and Automation*. IEEE, 2012, pp. 3503–3508.
- [44] E. Johns and G.-Z. Yang, "Generative methods for long-term place recognition in dynamic scenes," *International Journal of Computer Vision*, vol. 106, no. 3, pp. 297–314, 2014.
- [45] R. Zlot and M. Bosse, "Place recognition using keypoint similarities in 2d lidar maps," in *Experimental Robotics*. Springer, 2009, pp. 363–372.
- [46] O. Chum, M. Perdoch, and J. Matas, "Geometric min-hashing: Finding a (thick) needle in a haystack," in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. IEEE, 2009, pp. 17–24.
- [47] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [48] J. Schmidhuber, "Deep learning in neural networks: An overview," *Neural networks*, vol. 61, pp. 85–117, 2015.
- [49] T. Sun, L. Sun, and D.-Y. Yeung, "Fine-grained categorization via cnn-based automatic extraction and integration of object-level and part-level features," *Image and Vision Computing*, 2017.
- [50] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Proceedings of the European Conference on Computer Vision*. Springer, 2014, pp. 818–833.
- [51] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2921–2929.
- [52] S. Guo, W. Huang, L. Wang, and Y. Qiao, "Locally supervised deep hybrid model for scene recognition," *IEEE Transactions on Image Processing*, vol. 26, no. 2, pp. 808–820, 2017.
- [53] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *Proceedings of the 22nd ACM international conference on Multimedia*. ACM, 2014, pp. 675–678.
- [54] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman, "Return of the devil in the details: Delving deep into convolutional nets," *arXiv preprint arXiv:1405.3531*, 2014.
- [55] L. Wang, C.-Y. Lee, Z. Tu, and S. Lazebnik, "Training deeper convolutional networks with deep supervision," *arXiv preprint arXiv:1505.02496*, 2015.
- [56] R. Gomez-Ojeda, M. Lopez-Antequera, N. Petkov, and J. Gonzalez-Jimenez, "Training a convolutional neural network for appearance-invariant place recognition," *arXiv preprint arXiv:1505.07428*, 2015.
- [57] R. B. Rusu and S. Cousins, "3d is here: Point cloud library (pcl)," in *proceedings of IEEE International Conference on Robotics and Automation*. IEEE, 2011, pp. 1–4.
- [58] R. B. Rusu, G. Bradski, R. Thibaux, and J. Hsu, "Fast 3d recognition and pose using the viewpoint feature histogram," in *proceedings of IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2010, pp. 2155–2162.
- [59] W. Wohlkinger and M. Vincze, "Ensemble of shape functions for 3d object classification," in *IEEE International Conference on Robotics and Biomimetics (ROBIO)*. IEEE, 2011, pp. 2987–2992.
- [60] Z.-C. Marton, D. Pangercic, R. B. Rusu, A. Holzbach, and M. Beetz, "Hierarchical object geometric categorization and appearance classification for mobile manipulation," in *The 10th IEEE-RAS International Conference on Humanoid Robots (Humanoids)*. IEEE, 2010, pp. 365–370.
- [61] A. Glover, W. Maddern, M. Warren, S. Reid, M. Milford, and G. Wyeth, "Openfabmap: An open source toolbox for appearance-based loop closure detection," in *proceedings of IEEE International Conference on Robotics and Automation*. IEEE, 2012, pp. 4730–4735.
- [62] G. Bradski, "The opencv library," *Dr. Dobb's Journal: Software Tools for the Professional Programmer*, vol. 25, no. 11, pp. 120–123, 2000.