

# CORAL: Colored structural representation for bi-modal place recognition

Yiyuan Pan, Xuecheng Xu, Weijie Li, Yunxiang Cui, Yue Wang, Rong Xiong

**Abstract**—Place recognition is indispensable for a drift-free localization system. Due to the variations of the environment, place recognition using single-modality has limitations. In this paper, we propose a bi-modal place recognition method, which can extract a compound global descriptor from the two modalities, vision and LiDAR. Specifically, we first build the elevation image generated from 3D points as a structural representation. Then, we derive the correspondences between 3D points and image pixels that are further used in merging the pixel-wise visual features into the elevation map grids. In this way, we fuse the structural features and visual features in the consistent bird-eye view frame, yielding a semantic representation, namely CORAL. And the whole network is called CORAL-VLAD. Comparisons on the Oxford RobotCar show that CORAL-VLAD has superior performance against other state-of-the-art methods. We also demonstrate that our network can be generalized to other scenes and sensor configurations on cross-city datasets.

## I. INTRODUCTION

Loop closure is essential for the localization system because of the drift reduction, especially in large-scale outdoor scenes. A popular pipeline for loop closure usually employs place recognition as the first step, since it is able to find a place from the large map database that is close to the current place, based on purely sensor data similarity. Therefore, place recognition is widely applied in various autonomous robot systems for navigation.

The camera is the most popular sensor as it observes the texture of the environment at a low cost. Therefore, visual place recognition draws research attention for years. A traditional pipeline is to build a global feature descriptor by aggregating handcrafted sparse local features for each image, e.g. SIFT [14] and SURF [3]. Then the efficient searching method is used to look for the nearest global descriptor on the database as the most similar match. With the development of the deep neural network, convolution neural networks (CNNs) demonstrate promising retrieval performance [2, 18, 28] by extracting more reliable image descriptors. However, due to the susceptibility of the visual image under strong variations of seasons, weather, illumination, and viewpoints, building robust and discriminative image features remains a challenge.

To relieve the problem, LiDAR has been an alternative that provides accurate and relatively stable 3D structural information. Following the visual place recognition pipeline,

Yiyuan Pan, Xuecheng Xu, Weijie Li, Yunxiang Cui, Yue Wang, and Rong Xiong are with the State Key Laboratory of Industrial Control Technology and Institute of Cyber-Systems and Control, Zhejiang University, Zhejiang, China. Yue Wang is the corresponding author wangyue@iipc.zju.edu.cn.

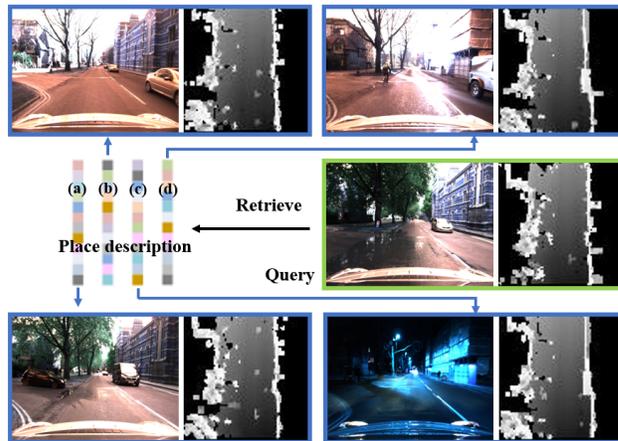


Fig. 1: The proposed bi-modal coupling of representation fusing visual image and elevation image retrieves the most similar sample in the reference maps with different environmental conditions, like (a) overcast, (b) sun, (c) night, and (d) dawn.

LiDAR-based place recognition methods employ deep neural networks to extract structural features from LiDAR scans [13, 1], demonstrating superior performance than the visual place recognition, especially in changing outdoor environment. However, LiDAR still has its weakness when the environmental structure has fewer features.

Due to the inherent shortages of both the camera and LiDAR, it is difficult to extract appropriate features using the single-modality sensor to describe numerous complex scenes. Thus, multi-sensor data fusion is becoming a feasible solution for place recognition. A common sensor for fusion image and geometry information is the RGB-D camera but it's not credible in the outdoor scenes. Using the camera and LiDAR is a more robust way to deal with multi-sensor fusion. Whereas, inconsistent viewpoints, diverse observation ranges and data structures of the different sensor data become major issues to generate a compound global descriptor effectively. In existing fusion-based place recognition methods, the visual and structural information are processed independently to generate their own descriptors and then concatenated directly as a compound descriptor without consistency in geometry.

In this paper, we set to combine the two modalities, vision and LiDAR, for place recognition, as shown in Fig. 1. The main novelty is the construction of the bi-modal fused representation, namely colored structural representation (CORAL). As shown in previous works on LiDAR

place recognition, various 3D representations, including point cloud [1], histogram [32], and polar image [11, 31], are designed to represent LiDAR scans, which highly impact the efficiency and effectiveness. We propose to build a local dense elevation map to describe the environmental structure. The map also derives the correspondences between 3D points and image pixels. The pixel-wise visual features are then inserted into the elevation map grids to semantically ‘colorize’ the structural features. With such tightly bi-modal coupling, CORAL encodes both visual and structural features in the same consistent bird-eye view (BEV) frame. The contributions can be summarized as:

- A local dense elevation map representation is utilized for place recognition, which is irrelevant to LiDAR hardware configurations to encode the structural information.
- A semantic representation with corresponding geometry named CORAL is proposed for place recognition which is robust towards various environmental changes.
- A validation on experiments is conducted to evaluate the performance of the proposed method, which shows superior performance in testing and generalization datasets. The code is also released <sup>1</sup>.

## II. RELATED WORK

**Vision-based place recognition** Vision-based place recognition is typically regarded as the problem of image retrieval which is solved by searching the most similar match on the reference database. Traditionally, some salient image parts are encoded as handcrafted local features, such as SURF [3], SIFT [14], or ORB [16, 17], and then these local features can be aggregated to a global descriptor leveraging feature aggregated methods, such as Bag-of-visual-words [6], VLAD [10] and Fisher Vectors [9]. Obtaining the global descriptor of a query, efficient searching methods like KD-tree search can be used to look for the closest global descriptor as the most similar match on the reference database. In recent years, handcrafted local features have been increasingly replaced by learnable features using deep neural networks that have significant improvement in extracting descriptive descriptors. Several mature networks for extracting local features, like VGG-Net [27] and ResNet [8], achieve an amazing performance on place recognition. As for aggregating local features, inspired by the traditional method - VLAD, NetVLAD [2] is formulated as a learnable function that is better in local features clustering. Generalized-Mean Pooling [24] is also an efficient and differentiable aggregated method, allows the network to capture a compact global descriptor in an end-to-end fashion.

**Structural-based place recognition** Considering structural features are more robust in changing environments, structural-based methods become an alternative for place recognition. Handcrafted structural descriptors, like PFH [25] and SHOT [26], usually have poor generalizability, that they can only be used in specific tasks. To relieve this problem,

convolution neural networks are used to extract structural descriptors. Owing to orderless of point clouds, several works convert the raw point clouds into a 3D volume representation, such as VoxelNet [34] and volumetric CNNs [23]. However, the conversion process introduces high quantization loss and requires lots of computation time. PointNet [22] is a pioneering work that is able to capture structural features from raw point clouds directly. Combining PointNet and NetVLAD, PointNetVLAD [1] is the first approach to achieve large-scale long-term place recognition with the input of raw point data. However, PointNet operates each point independently thus ignoring the local structure relationship of points. In response to this problem, LPD-Net [13] proposes the adaptive local feature extraction module and the graph-based neighborhood aggregation module. PCAN [33] introduces an attention map to predict the significance of local regions. Both methods boost the performance of place recognition.

**Fusion-based place recognition** Multi-sensor fusion solution by incorporating visual and structural features is a more appropriate way to achieve place recognition under strong environment variations. In recent years, a few fusion-based works have emerged for place recognition. The naive strategy [30, 19] of generating a compound descriptor from different modality inputs is directly combining two global descriptors extracted from two independent network streams, ignoring the inner relationship of visual context and geometry structure in the same region. Referring to some 3D object detection methods, works of [29, 12] fuse the visual and structural information by operating the intermediate structural feature map and visual feature map in the same frame, which inspires us to construct a more discriminative and robust compound descriptor.

## III. METHODOLOGY

The architecture of the proposed fusion network is shown in Fig. 2. The network inputs consist of two streams: raw front-view images from the camera and filtered elevation images generated from the LiDAR scans. The two-stream network tightly couples the visual and structural features which enforce the final representation to encode the semantics into the geometric structure. Specifically, the dense elevation image representation encodes the structural information, and it also provides the point-to-pixel correspondences, which is leveraged to insert visual features into the consistent BEV frame. Therefore, aggregating the bi-modal features is geometrically sensible, yielding the final global descriptor for place recognition.

### A. Elevation map generation

We introduce the elevation map representation defined on a grid map. Each 2D grid indexes an elevation to describe the environment structure. Originated in the 2D occupancy grid map, the elevation map replaces the occupancy information with the elevation, which is capable of representing the 2.5D ground surface.

First, the 6 degree-of-freedom (DOF) pose of the sensor denoted as  $T$  is calculated by LiDAR inertial odometry

<sup>1</sup>[https://github.com/Panyiyuan96/CORAL\\_Pytorch.git](https://github.com/Panyiyuan96/CORAL_Pytorch.git)

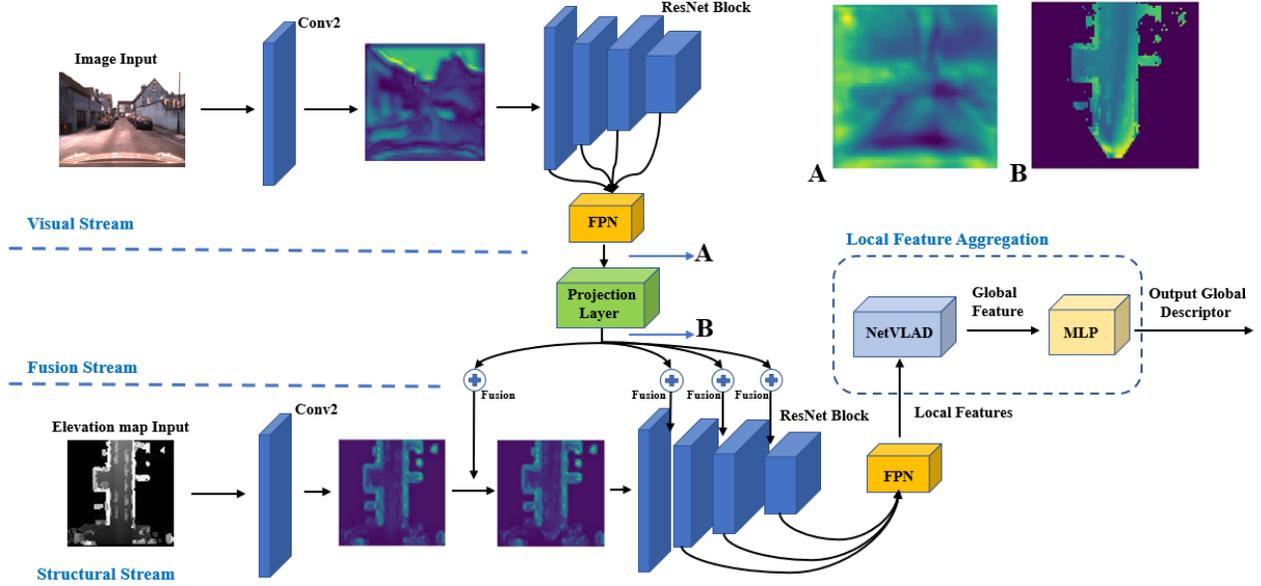


Fig. 2: The network architecture of CORAL-VLAD. The network inputs an elevation image generated by LiDAR scans and a visual image, and then outputs a global descriptor to predict the most similar scene on the reference database. *A* shows multi-scale visual feature map from the output of FPN layer, *B* shows the BEV visual feature map from the output of projection layer.

algorithm, with respect to the global frame  $G$ . Given a 3D point  $p_i$  of a new measurement in the sensor frame, we transform it into the global frame by  $Tp_i$  and calculate the elevation:

$$e_p = PTp_i \quad (1)$$

the projection matrix  $P = [0, 0, 1]$  retrieves 3rd entry of the transformed point as the elevation  $e_p$ , while the first two entries are rounded and then converted to the corresponding index of the grid map. For processing multiple observations of the same grid, the variance  $\sigma_p$  of elevation is introduced to describe the uncertainty of elevation measurement according to range sensor models [5]. In this way, each grid data  $(e_g, \sigma_g^2)$  is updated according to the corresponding measurement data  $(e_p, \sigma_p^2)$ . Note that only the measurement within the Mahalanobis distance threshold of the grid is fused by variance weighted strategy to obtain an updated grid data as follows:

$$e_g = \frac{\sigma_p^2 e_g + \sigma_g^2 e_p}{\sigma_p^2 + \sigma_g^2} \quad \sigma_g^2 = \frac{\sigma_p^2 \sigma_g^2}{\sigma_p^2 + \sigma_g^2} \quad (2)$$

Furthermore, to clear the dynamic objects due to the use of the accumulation mechanism, ray tracing is utilized to check whether the grid is crossed by a ray, and if the grid is occupied by a dynamic object, the elevation and variance of this grid are initialized. The detailed generation process can be found in [20].

**Elevation image** To utilize 2D convolution neural networks directly, the elevation map is converted into a one-channel grayscale image with the same size as the grid map, namely elevation image. As shown in Fig. 3, we scale the

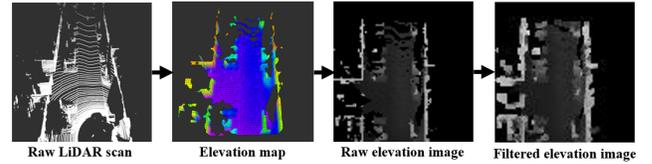


Fig. 3: Elevation image generation. The elevation map accumulated from raw LiDAR scans is projected on the image plane as a raw elevation image. The filtered elevation image is generated using a mean filter.

elevation value from 0 to 255 as the grayscale value of the elevation image. Additionally, a mean filter is applied to fill up some invalid pixels which have not been observed on the elevation image using a  $3 \times 3$  filter template. Compared with the single scan based BEV representation, the elevation image is much denser to describe the structure, also leading to more accurate point-to-pixel correspondences.

### B. Network architecture

The architecture of the proposed compound network is shown in Fig. 2. The whole network can be divided into three main components: the feature extraction module, the fusion layer, and the local feature aggregation module.

**Feature extraction** Our multi-sensor network has two streams for feature extraction. The visual stream employs the lightweight ResNet18 [8] as the backbone to efficiently extract the visual features. Due to the utilize of several convolutional layers with stride 2, the output feature map suffers a dramatic reduction in size which is not suitable for fusion. So we use feature pyramid network (FPN) to combine

four feature maps from the residual block group to recover the final feature map with the size same as the first residual block and exploit multi-scale feature information.

The structural stream comprises a group of convolutional layers to capture structural features, and four groups of residual blocks to extract fusion features. Except for the first group, all groups start with the convolution layer with stride 2 and all other convolutions are with stride 1. The number of  $3 \times 3$  kernel convolutions in each group is 2, 4, 4, 6, 6, and each group outputs the feature vector with the corresponding dimensions of 64, 64, 128, 192, and 256 respectively. The outputs of the last three residual blocks also exploit multi-scale information as the visual stream does.

**Fusion Layer** The fusion layer comprises two components: the projection layer and the fusion module. The goal of the projection layer is to convert the front-view visual feature map into the BEV visual feature map through sparse matrix multiplication. The grid index of the elevation image and the corresponding elevation value can be converted into a 3D position in the global frame. So each grid can be denoted as a 3D point and transformed into the LiDAR frame according to the estimated pose. With the intrinsic parameters of the camera and extrinsic parameters between the LiDAR and the camera, each 3D point can be projected onto the 2D camera image plane which helps retrieve corresponding visual image features. Following this idea, we implement a parameterless projection layer to find the correspondences. Considering that the coordinates of 2D projected points are often non-integers, we combine the visual features from adjacent discrete pixels by bilinear interpolating. After that visual feature map is generated in the BEV frame consistent with the elevation image, thus achieving appropriate features corresponding to later fusion.

Then, we arrive at the generation of CORAL representation in the fusion module. Denote the  $i$ th layer of structural feature map as  $S_i$ , and the corresponding BEV visual feature map as  $V_i$ . To maintain the same shape of  $S_i$  and  $V_i$ , pooling operation  $Pool(\cdot)$  is used to adjust the output size of the raw projection layer  $V$  and  $1 \times 1$  kernel convolution  $Conv(\cdot)$  is applied to keep the same channel size. We try two different methods for aggregating features. The first one uses element-wise concatenation to obtain CORAL  $F_i$  defined by

$$F_i = [Conv(Pool(V)), S_i] \quad (3)$$

In the second method, feature maps are combined by element-wise summation given by

$$F_i = Conv(Pool(V)) + S_i \quad (4)$$

We also propose two fusion strategies that combine the structural feature map and the visual feature map from the first residual layer, or all four blocks, which are evaluated in Section IV.

**Local feature aggregation** The local feature aggregation module learns to further extract a global descriptor. Considering the capacity of the NetVLAD on aggregating features, we feed CORAL to a NetVLAD layer and generate a global descriptor. Furthermore, we utilize a multi-layer

perception (MLP) to process the raw output of NetVLAD for learning a dimension reduction mapping to decrease the size of the global descriptor which accelerates the nearest neighbor search.

### C. Training the fusion descriptor

We train our compound network in an end-to-end fashion to yield a bi-modal fused global descriptor. A margin-based loss is adopted to train pairs of samples labeled with positive or negative based on corresponding GPS position.

**Loss function** The training data is constructed as sets of tuples  $\mathcal{T} = (P_a, \{P_{pos}\}, \{P_{neg}\})$ .  $P_a$  denotes an anchor with a pair of visual image and elevation image,  $\{P_{pos}\}, \{P_{neg}\}$  represent the set of anchor’s positive matches and negative matches determined by their relative distances. We apply the squared Euclidean distance  $\delta$  to evaluate the similarity of two global descriptors. Margin-based loss is used to minimize the distance  $\delta_{a,pos}$  between the global descriptors of matching samples  $(P_a, P_{pos})$  while pushing apart the dissimilar matches  $(P_a, P_{neg})$ . To achieve faster convergence and better discrimination, we only use the closest/hardest negative  $P_{neg}^-$  in  $\{P_{neg}\}$  and the most dissimilar positive  $P_{pos}^+$  in  $\{P_{pos}\}$  during back propagation. Comparing with various retrieval loss functions, we utilize the lazy quadruplet [4] as the training loss:

$$\mathcal{L}(\mathcal{T}, P_{neg*}) = \max([\alpha + \delta_{a,pos} - \delta_{a,neg}]_+) + \max([\beta + \delta_{a,pos} - \delta_{a,neg*}]_+) \quad (5)$$

where  $[\cdot]_+$  is the hinge loss,  $\alpha, \beta$  are constant margin parameters, and  $P_{neg*}$  is randomly sampled from training data which is dissimilar to all observations of  $\mathcal{T}$ .

**Data sampling strategy** The original hard-negative training strategy usually offers faster convergence but it may lead to a collapsed model. To alleviate this problem, we divide the training process into two stages. In the first stage, negative matches  $\{P_{neg}\}$  are sampled randomly from all negative samples. This is followed by a second training stage, negative matches  $\{P_{neg}\}$  are generated by looking for the hardest pairs from all global descriptors of negatives calculated by the latest model. This hard-negative mining strategy ensures negative samples become progressively harder, which avoids prolonged convergence and boosts the performance of the converged model.

## IV. EXPERIMENTS

In this section, we discuss the datasets and settings for training and evaluation. Quantitative results are provided to demonstrate the performance of our method on diverse scene conditions and cross-city generalization experiments. Furthermore, the loss parameters are set as  $\alpha = 0.5, \beta = 0.2$ , the number of positive matches  $\{P_{pos}\}$  is 2 and negative matches  $\{P_{neg}\}$  is 18.

TABLE I: Comparison results with the average recall@1 and recall@1% of different networks on the Oxford dataset.

	Ave recall@1	Ave recall@1%
PN-VLAD	67.94	81.01
LPD-Net	86.28	94.42
Img-VLAD	64.47	85.24
Aug-Net	79.47	91.24
Vis-VLAD(ours)	57.62	83.05
Ele-VLAD(ours)	82.49	93.61
Sum-First(ours)	82.44	92.71
Con-First(ours)	84.82	93.62
Sum-Four(ours)	86.23	94.43
Con-Four(ours)	<b>88.93</b>	<b>96.13</b>

TABLE II: Average timing for computing a single instance using NVIDIA 2080 Ti.

Network	Avg. computational cost per descriptor
PN-VLAD	8.44ms
LPD-Net	15.28ms
Img-NetVLAD	62.28ms
Aug-Net	18.80ms
CORAL-VLAD(ours)	11.19ms

#### A. Datasets and settings

We choose the Oxford Robotcar dataset [15] for training and testing. The car equipped with LiDAR sensors and cameras repeatedly drives through the same regions in one year and collects data in large-scale and long-term conditions across weather, season and illumination. To create training tuples, we define the similar observations with a relative distance less than 10m and heading difference less than 30 degrees as positive pairs, and those at least 50m apart as negative pairs. Referring to the data splitting rules of [1], we obtain 21636 training tuples from 44 sequences of the original dataset and 3011 testing samples from 23 sequences. To keep large overlapping regions between the camera image and the elevation image, the elevation map is set with the size of  $80 \times 80$  and the resolution of 0.5m. The input of the visual image is downscaled into  $112 \times 112$  for extracting visual features efficiently. Furthermore, as scans are accumulated by 2D LiDAR scans at fixed intervals, there are some incomplete samples on the dataset. We discard them when generating the elevation image.

**Generalization settings** We choose KITTI dataset [7] to evaluate the generalization performance across city and sensor configurations. Only Sequence 00 which visits the same places repeatedly is used. The first 170s of Sequence 00 are used to construct the reference map and the rest are used as localization queries. The second dataset for generalization evaluation is collected by ourselves in different seasons, namely YQ. Data in overcast scene and snowy scene are used for mapping and query respectively. For a fair comparison, the same criteria on Oxford Robotcar dataset are applied.

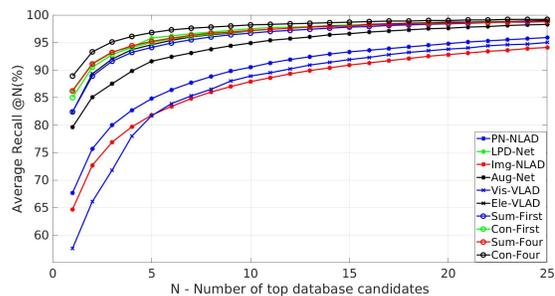


Fig. 4: Average recall@N(%) of different networks on the Oxford dataset.

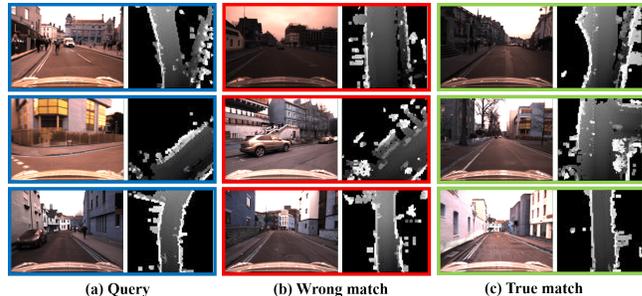


Fig. 5: Networks limitations. These are incorrect matches retrieved by our network, where (a) is the query, (b) is the wrong match, and (c) shows the true match.

#### B. Place recognition results on the Oxford dataset

We present qualitative results to demonstrate the feasibility of our CORAL-VLAD on the Oxford dataset under changing scene conditions. The common assessment indicators of place recognition - average recall@1 and the average recall@1% are used to evaluate the network performance. For fairness, the final global descriptor dimensions of all networks are set to 256.

**Ablation study on fusion strategy** We test our compound network with four feature fusion strategies operating on the intermediate visual feature map and the structural feature map, including (a) element-wise summarization in the first residual block (Sum-First), (b) element-wise concatenation in the first residual block (Con-First), (c) element-wise summarization in four residual blocks (Sum-Four), (d) element-wise concatenation in four residual blocks (Con-Four). The results are shown in Fig. 4 and Tab. I, the multi-scale fusion strategies outperform the single-scale due to the adequate information. And the better performance of concatenation than summation contributes to the intact information. In the following experiments, we use our network with the Con-Four fusion strategy as the final version, denoted as CORAL-VLAD.

**Ablation study on sensor modal** Furthermore, to investigate the contribution of the fusion step, we separate our fusion architecture into an independent visual stream and structural stream with a shared NetVLAD module to generate the global descriptor, called Vision NetVLAD(Vis-VLAD) and Elevation image NetVLAD(Ele-VLAD). The results are

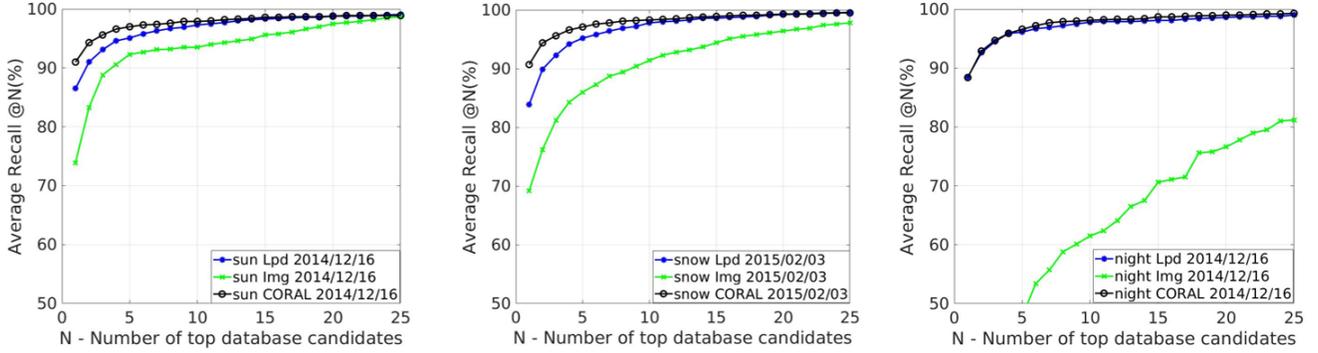


Fig. 6: Average recall@N(%) with LPD-Net(Lpd), Img-VLAD(Img) and CORAL-VLAD(CORAL) under different scene conditions on the Oxford dataset.

also shown in Fig. 4 and Tab. I. Except for the fusion strategy of Sum-First, results have been significantly improved using composite features. Furthermore, the results for single modal place recognition are shown, validating the correct design and implementation of our method.

**Comparison with the-state-of-art methods** We compare our method with the state-of-the-art single modal place recognition methods, PointNetVLAD (PN-VLAD) [1], LPD-Net [13] and Img-VLAD [2], as well as bi-modal place recognition method, compound network (Aug-Net)[19].

Using only structural data input, the performance of our Ele-VLAD is better than PN-VLAD, yet slightly worse than LPD-Net. Inputs of LPD-Net and PN-VLAD are 4096 filtered points with detailed 3D structural information of the environment while the elevation image only has  $40 \times 40$  grids with one-channel elevation. It also proves effectiveness of the elevation image for representing 3D geometric data. Although providing a rich appearance context, Vis-VLAD and Img-VLAD cannot outperform Ele-VLAD, which proves that the elevation image representation is more robust under environment variations. Results are shown in Fig. 1. Meanwhile, there are many over-exposed images on the Oxford dataset, causing loss of visual features.

Using composite descriptors, the performance of CORAL-VLAD exceeds all other methods, including Aug-Net. Note that there is a difference in the Aug-Net from [19], since we employ the PN-VLAD splitting rules clarified in [1] to ensure consistency among all methods. All comparison network training sessions last no longer than 24h. We suppose that the main reason for the slightly worse results of Aug-Net is the incompleteness of the 3D volume based representation.

**Comparison under different conditions** We compare the performance of LPD-Net, Img-VLAD, CORAL-VLAD, which use inputs in different representation modalities under changing scene conditions. Fig. 6 shows the results when queries are taken from three different scene conditions against the testing database on the Oxford database (except query run). It can be seen that the huge variants of environmental conditions have little impact on retrieval performance using structural cues. As almost all corresponding visual features are lost at night scene, the visual-based approach

TABLE III: Comparison generalization results of the average recall@1(%) on the KITTI and YQ dataset

Network	KITTI_laser	KITTI_stereo	YQ
PN-VLAD	72.43	65.43	40.36
LPD-Net	74.58	65.82	62.35
Aug-Net	75.60	70.56	66.91
CORAL-VLAD(ours)	<b>76.43</b>	<b>70.77</b>	<b>73.82</b>

obtains lots of incorrect matches. Accordingly, CORAL-VLAD has minimal improvement over LPD-Net in such condition, showing the limitation of additional visual modal.

**Cases study** Fig. 1 shows some of the successfully matched results in changing environments. We can observe that our network has learned robust features and alleviated the negative effect brought by dynamic obstacles and variations in illumination. Fig. 5 shows three wrong cases. We can see that the network is confused in the same scenes with opposite viewpoints (top row) and a few overlapping areas (middle row and bottom row).

**Comparison on efficiency** We further evaluate the running time of the network. Due to the use of lightweight network backbones and compact structural representation, our network implementation takes about 11ms which is faster than all methods except PN-VLAD shown in Tab. II. For generating the elevation map, we have a GPU-based implementation at almost 30Hz as shown in [21], achieving real-time performance for robotics applications.

### C. Generalization evaluation

To analyze the generalization of our network, we evaluate our network on the YQ and KITTI datasets using the model trained on the Oxford Robotcar dataset. These cross-city datasets consist of unobserved conditions with different sensor configurations, including sensor types and extrinsic parameters. Specifically, KITTI\_laser uses the point cloud collected by LiDAR while KITTI\_stereo uses the point cloud generated from stereo images.

In the case of KITTI\_laser dataset, the elevation image is generated in real-time using Velodyne 64 HDL LiDAR

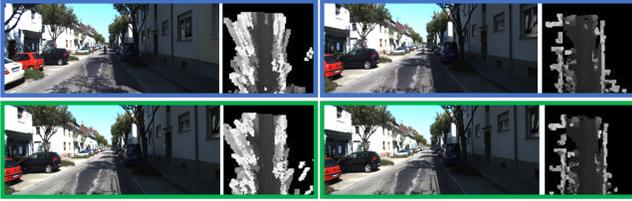


Fig. 7: Matching examples on the KITTI dataset. The left column shows the correct match on the KITTI\_stereo and the right column shows the correct match on the KITTI\_laser.



Fig. 8: The collection of the YQ dataset. The left image shows the data collection platform. The yellow line of the right image shows the trajectory of the data collection on our campus.

shown in Fig. 7 and our network outperforms other methods in Tab. III. Although there is noise involved during the calculation of the disparity map on the KITTI\_stereo dataset, generated elevation maps become blurred, which loses some structural information and introduces inaccurate matches between pixels and 3D points. CORAL-VLAD still achieves the best results, demonstrating the validity of our network for pure vision-based applications.

Additionally, we conduct experiments under different weather conditions on our campus. Unlike evaluating with similar conditions at a short period of time on the KITTI dataset, the unmanned ground vehicle equipped with a LiDAR and an RGB camera collects sensor data on sunny days in spring and snowy days in winter denoted as YQ dataset. The data collection platform is shown in Fig. 8. In addition, DGPS provides the ground-truth position to check the correctness of the matching results. The summer-overcast scene of the YQ dataset is utilized as the reference database, and the winter-snow scene as the query dataset. Matching examples are shown in Fig. 9, formulating a generalization evaluation on both conditions and sensor configurations. CORAL-VLAD still achieves the best performance in this scenario as shown in Tab. III, because of the elevation image based structural representation, transforming the texture into a geometric sensible representation.

## V. CONCLUSION

In this paper, we introduce the elevation map as the structural information and propose the bi-modal environment representation CORAL to fuse the structural and visual

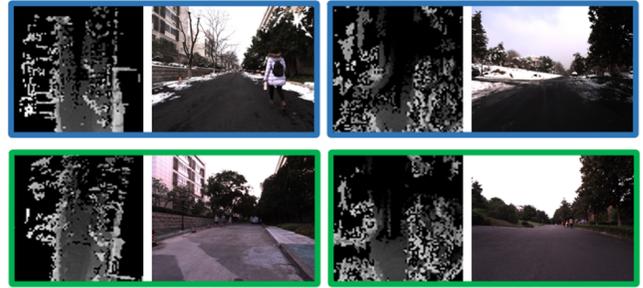


Fig. 9: Matching examples on the YQ dataset. The left column shows the correct matching example at point A in Fig. 8. The right column shows the matching result at point B in Fig. 8.

features in the same consistent BEV frame, which can handle various environmental variances like viewpoint changes, illuminations, and structure losses. We show that the method performs best on the Oxford Robotcar dataset, as well as generalization test on conditions and sensor configurations using the cross-city datasets of the KITTI dataset and YQ dataset.

## VI. ACKNOWLEDGMENT

This work was supported in part by the National Nature Science Foundation of China under Grant 61903332, and in part by the Natural Science Foundation of Zhejiang Province under grant number LGG21F030012.

## REFERENCES

- [1] Mikaela Angelina Uy and Gim Hee 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.
- [2] Relja Arandjelovic et al. “NetVLAD: CNN architecture for weakly supervised place recognition”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 5297–5307.
- [3] Herbert Bay et al. “Speeded-up robust features (SURF)”. In: *Computer vision and image understanding* 110.3 (2008), pp. 346–359.
- [4] Weihua Chen et al. “Beyond triplet loss: a deep quadruplet network for person re-identification”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2017, pp. 403–412.
- [5] Péter Fankhauser et al. “Kinect v2 for mobile robot navigation: Evaluation and modeling”. In: *2015 International Conference on Advanced Robotics (ICAR)*. IEEE. 2015, pp. 388–394.
- [6] Dorian Gálvez-López and Juan D Tardos. “Bags of binary words for fast place recognition in image sequences”. In: *IEEE Transactions on Robotics* 28.5 (2012), pp. 1188–1197.

- [7] Andreas Geiger, Philip Lenz, and Raquel Urtasun. "Are we ready for autonomous driving? the kitti vision benchmark suite". In: *2012 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE. 2012, pp. 3354–3361.
- [8] Kaiming He et al. "Deep residual learning for image recognition". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 770–778.
- [9] Herve Jegou et al. "Aggregating local image descriptors into compact codes". In: *IEEE transactions on pattern analysis and machine intelligence* 34.9 (2011), pp. 1704–1716.
- [10] Hervé Jégou et al. "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.
- [11] G. Kim, B. Park, and A. Kim. "1-Day Learning, 1-Year Localization: Long-Term LiDAR Localization Using Scan Context Image". In: *IEEE Robotics and Automation Letters* 4.2 (Apr. 2019), pp. 1948–1955.
- [12] Ming Liang et al. "Deep continuous fusion for multi-sensor 3d object detection". In: *Proceedings of the European Conference on Computer Vision (ECCV)*. 2018, pp. 641–656.
- [13] Zhe Liu et al. "Lpd-net: 3d point cloud learning for large-scale place recognition and environment analysis". In: *Proceedings of the IEEE International Conference on Computer Vision*. 2019, pp. 2831–2840.
- [14] David G Lowe. "Distinctive image features from scale-invariant keypoints". In: *International journal of computer vision* 60.2 (2004), pp. 91–110.
- [15] Will Maddern et al. "1 year, 1000 km: The Oxford RobotCar dataset". In: *The International Journal of Robotics Research* 36.1 (2017), pp. 3–15.
- [16] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. "ORB-SLAM: a versatile and accurate monocular SLAM system". In: *IEEE transactions on robotics* 31.5 (2015), pp. 1147–1163.
- [17] Raul Mur-Artal and Juan D Tardós. "Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras". In: *IEEE Transactions on Robotics* 33.5 (2017), pp. 1255–1262.
- [18] Hyeonwoo Noh et al. "Large-scale image retrieval with attentive deep local features". In: *Proceedings of the IEEE international conference on computer vision*. 2017, pp. 3456–3465.
- [19] Amadeus Oertel, Titus Cieslewski, and Davide Scaramuzza. "Augmenting Visual Place Recognition with Structural Cues". In: *arXiv preprint arXiv:2003.00278* (2020).
- [20] Yiyuan Pan et al. "GEM: online globally consistent dense elevation mapping for unstructured terrain". In: *IEEE Transactions on Instrumentation and Measurement* (2020).
- [21] Yiyuan Pan et al. "GPU accelerated real-time traversability mapping". In: *2019 IEEE International Conference on Robotics and Biomimetics (ROBIO)*. IEEE. 2019, pp. 734–740.
- [22] Charles R Qi et al. "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.
- [23] Charles R Qi et al. "Volumetric and multi-view cnns for object classification on 3d data". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 5648–5656.
- [24] Filip Radenović, Giorgos Tolias, and Ondřej Chum. "Fine-tuning CNN image retrieval with no human annotation". In: *IEEE transactions on pattern analysis and machine intelligence* 41.7 (2018), pp. 1655–1668.
- [25] Radu Bogdan Rusu et al. "Aligning point cloud views using persistent feature histograms". In: *2008 IEEE/RSJ international conference on intelligent robots and systems*. IEEE. 2008, pp. 3384–3391.
- [26] Samuele Salti, Federico Tombari, and Luigi Di Stefano. "SHOT: Unique signatures of histograms for surface and texture description". In: *Computer Vision and Image Understanding* 125 (2014), pp. 251–264.
- [27] Karen Simonyan and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition". In: *arXiv preprint arXiv:1409.1556* (2014).
- [28] Li Tang et al. "Adversarial feature disentanglement for place recognition across changing appearance". In: *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE. 2020, pp. 1301–1307.
- [29] Zining Wang, Wei Zhan, and Masayoshi Tomizuka. "Fusing bird's eye view lidar point cloud and front view camera image for 3d object detection". In: *2018 IEEE Intelligent Vehicles Symposium (IV)*. IEEE. 2018, pp. 1–6.
- [30] Shaorong Xie et al. "Large-Scale Place Recognition Based on Camera-LiDAR Fused Descriptor". In: *Sensors* 20.10 (2020), p. 2870.
- [31] Xuecheng Xu et al. "DiSCO: Differentiable Scan Context with Orientation". In: *arXiv preprint arXiv:2010.10949* (2020).
- [32] Huan Yin et al. "LocNet: Global localization in 3D point clouds for mobile vehicles". In: *2018 IEEE Intelligent Vehicles Symposium (IV)*. IEEE. 2018, pp. 728–733.
- [33] Wenxiao Zhang and Chunxia Xiao. "PCAN: 3D attention map learning using contextual information for point cloud based retrieval". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2019, pp. 12436–12445.
- [34] Yin Zhou and Oncel Tuzel. "Voxelnet: End-to-end learning for point cloud based 3d object detection". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2018, pp. 4490–4499.